

# CS 373: Combinatorial Algorithms

University of Illinois, Urbana-Champaign

Instructor: Jeff Erickson

## Teaching Assistants:

- **Spring 1999:** Mitch Harris and Shripad Thite
- **Summer 1999 (IMCS):** Mitch Harris
- **Summer 2000 (IMCS):** Mitch Harris
- **Fall 2000:** Chris Neihengen, Ekta Manaktala, and Nick Hurlburt
- **Spring 2001:** Brian Ensink, Chris Neihengen, and Nick Hurlburt
- **Summer 2001 (I2CS):** Asha Seetharam and Dan Bullok
- **Fall 2002:** Erin Wolf, Gio Kao, Kevin Small, Michael Bond, Rishi Talreja, Rob McCann, and Yasutaka Furakawa

© Copyright 1999, 2000, 2001, 2002, 2003 Jeff Erickson.

This work may be freely copied and distributed.

It may not be sold for more than the actual cost of reproduction.

This work is distributed under a Creative Commons license; see <http://creativecommons.org/licenses/by-nc-sa/1.0/>.

For the most recent edition, see <http://www.uiuc.edu/~jeffe/teaching/373/>.



*For junior faculty, it may be a choice between a book and tenure.*

— George A. Bekey, “The Assistant Professor’s Guide to the Galaxy” (1993)

*I’m writing a book. I’ve got the page numbers done.*

— Stephen Wright

## About These Notes

This course packet includes lecture notes, homework questions, and exam questions from the course ‘CS 373: Combinatorial Algorithms’, which I taught at the University of Illinois in Spring 1999, Fall 2000, Spring 2001, and Fall 2002. Lecture notes and videotapes lectures were also used during Summer 1999, Summer 2000, Summer 2001, and Fall 2002 as part of the UIUC computer science department’s Illinois Internet Computer Science (I2CS) program.

The recurrences handout is based on samizdat, probably written by Ari Trachtenberg, based on a paper by George Lueker, from an earlier semester taught by Ed Reingold. I wrote most of the lecture notes in Spring 1999; I revised them and added a few new notes in each following semester. Except for the infamous Homework Zero, which is entirely my doing, homework and exam problems and their solutions were written mostly by the teaching assistants: Asha Seetharam, Brian Ensink, Chris Neihengen, Dan Bullok, Ekta Manaktala, Erin Wolf, Gio Kao, Kevin Small, Michael Bond, Mitch Harris, Nick Hurlburt, Rishi Talreja, Rob McCann, Shripad Thite, and Yasu Furakawa. Lecture notes were posted to the course web site a few days (on average) after each lecture. Homeworks, exams, and solutions were also distributed over the web. I have deliberately excluded solutions from this course packet.

The lecture notes, homeworks, and exams draw heavily on the following sources, all of which I can recommend as good references.

- Alfred V. Aho, John E. Hopcroft, and Jeffrey D. Ullman. *The Design and Analysis of Computer Algorithms*. Addison-Wesley, 1974. (This was the textbook for the algorithms classes I took as an undergrad at Rice and as a masters student at UC Irvine.)
- Sara Baase and Allen Van Gelder. *Computer Algorithms: Introduction to Design and Analysis*. Addison-Wesley, 2000.
- Mark de Berg, Marc van Kreveld, Mark Overmars, and Otfried Schwarzkopf. *Computational Geometry: Algorithms and Applications*. Springer-Verlag, 1997. (This is the required textbook in my computational geometry course.)
- Thomas Cormen, Charles Leiserson, Ron Rivest, and Cliff Stein. *Introduction to Algorithms*, second edition. MIT Press/McGraw-Hill, 2000. (This is the required textbook for CS 373, although I never actually use it in class. Students use it as a educated second opinion. I used the first edition of this book as a teaching assistant at Berkeley.)
- Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman, 1979.
- Michael T. Goodrich and Roberto Tamassia. *Algorithm Design: Foundations, Analysis, and Internet Examples*. John Wiley & Sons, 2002.
- Dan Gusfield. *Algorithms on Strings, Trees, and Sequences: Computer Science and Molecular Biology*. Cambridge University Press, 1997.
- Udi Manber. *Introduction to Algorithms: A Creative Approach*. Addison-Wesley, 1989. (I used this textbook as a teaching assistant at Berkeley.)

- Rajeev Motwani and Prabhakar Raghavan. *Randomized Algorithms*. Cambridge University Press, 1995.
- Ian Parberry. *Problems on Algorithms*. Prentice-Hall, 1995. (This was a recommended textbook for early versions of CS 373, primarily for students who needed to strengthen their prerequisite knowledge. This book is out of print, but can be downloaded as karmaware from <http://hercule.csci.unt.edu/~ian/books/poa.html> .)
- Robert Sedgewick. *Algorithms*. Addison-Wesley, 1988. (This book and its sequels have by far the best algorithm *illustrations* anywhere.)
- Robert Endre Tarjan. *Data Structures and Network Algorithms*. SIAM, 1983.
- Class notes from my own algorithms classes at Berkeley, especially those taught by Dick Karp and Raimund Seidel.
- Various journal and conference papers (cited in the notes).
- Google.

Naturally, everything here owes a great debt to the people who taught me this algorithm stuff in the first place: Abhiram Ranade, Bob Bixby, David Eppstein, Dan Hirshberg, Dick Karp, Ed Reingold, George Lueker, Manuel Blum, Mike Luby, Michael Perlman, and Raimund Seidel. I've also been helped immensely by many discussions with colleagues at UIUC—Ed Reingold, Edgar Raoms, Herbert Edelsbrunner, Jason Zych, Lenny Pitt, Mahesh Viswanathan, Shang-Hua Teng, Steve LaValle, and especially Sarel Har-Peled—as well as voluminous feedback from the students and teaching assistants. I stole the overall course structure (and the idea to write up my own lecture notes) from Herbert Edelsbrunner.

We did the best we could, but I'm sure there are still plenty of mistakes, errors, bugs, gaffes, omissions, snafus, kludges, typos, mathos, grammaros, thinkos, brain farts, nonsense, garbage, cruft, junk, and outright lies, all of which are entirely Steve Skiena's fault. I revise and update these notes every time I teach the course, so please let me know if you find a bug. (Steve is unlikely to care.)

When I'm teaching CS 373, I award extra credit points to the first student to post an explanation and correction of any error in the lecture notes to the course newsgroup ([uiuc.class.cs373](mailto:uiuc.class.cs373)). Obviously, the number of extra credit points depends on the severity of the error and the quality of the correction. If I'm not teaching the course, encourage your instructor to set up a similar extra-credit scheme, and forward the bug reports to ~~Steve~~ me!

Of course, any other feedback is also welcome!

Enjoy!

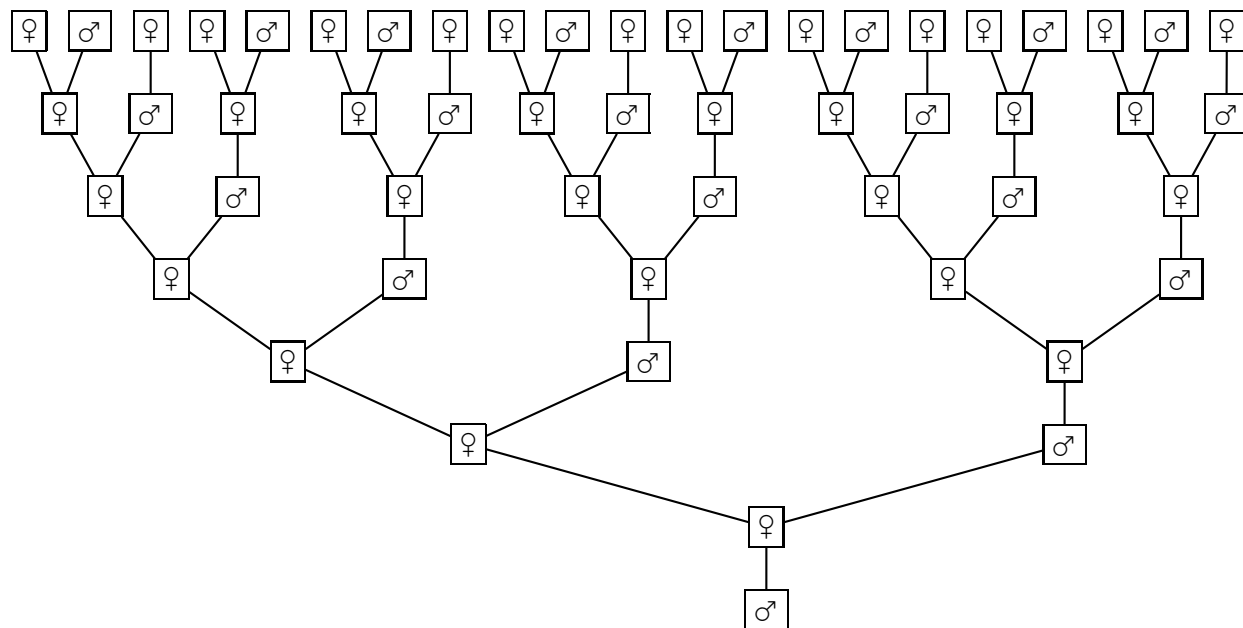
— Jeff

*It is traditional for the author to magnanimously accept the blame for whatever deficiencies remain. I don't. Any errors, deficiencies, or problems in this book are somebody else's fault, but I would appreciate knowing about them so as to determine who is to blame.*

— Steven S. Skiena, *The Algorithm Design Manual*, Springer-Verlag, 1997, page x.

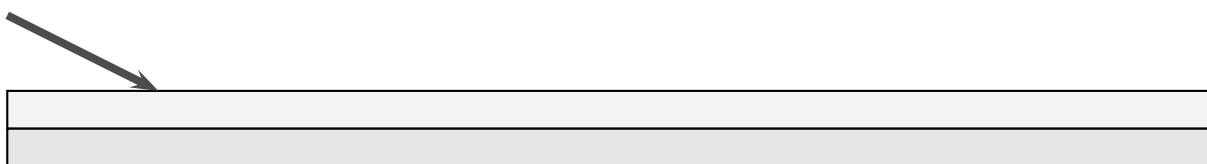
# 1 Fun with Fibonacci numbers

Consider the reproductive cycle of bees. Each male bee has a mother but no father; each female bee has both a mother and a father. If we examine the generations we see the following family tree:

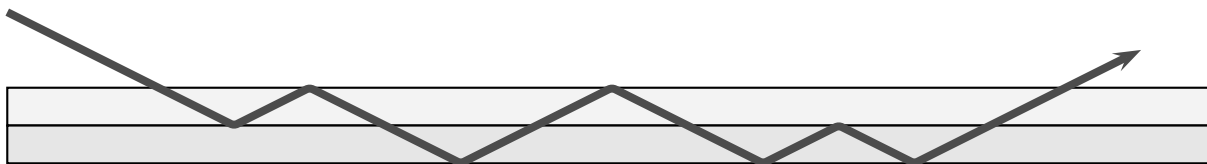


We easily see that the number of ancestors in each generation is the sum of the two numbers before it. For example, our male bee has three great-grandparents, two grandparents, and one parent, and  $3 = 2 + 1$ . The number of ancestors a bee has in generation  $n$  is defined by the Fibonacci sequence; we can also see this by applying the rule of sum.

As a second example, consider light entering two adjacent planes of glass:

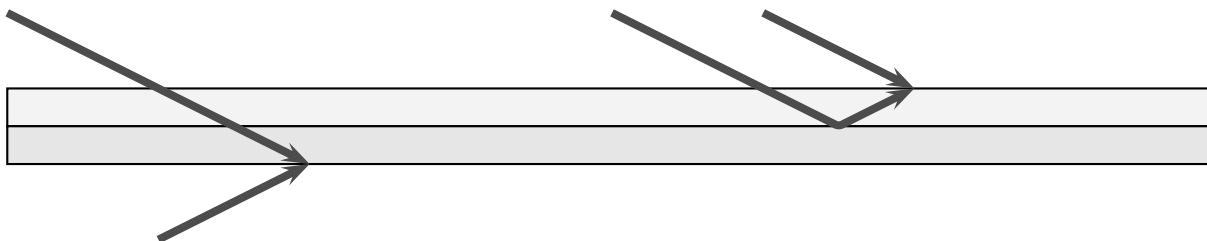


At any meeting surface (between the two panes of glass, or between the glass and air), the light may either reflect or continue straight through (refract). For example, here is the light bouncing seven times before it leaves the glass.



In general, how many different paths can the light take if we are told that it bounces  $n$  times before leaving the glass?

The answer to the question (in case you haven't guessed) rests with the Fibonacci sequence. We can apply the rule of sum to the event  $E$  constituting all paths through the glass in  $n$  bounces. We generate two separate sub-events,  $E_1$  and  $E_2$ , illustrated in the following picture.



- **Sub-event  $E_1$ :** Let  $E_1$  be the event that the first bounce is *not* at the boundary between the two panes. In this case, the light must first bounce off the bottom pane, or else we are dealing with the case of having zero bounces (there is only one way to have zero bounces). However, the number of remaining paths after bouncing off the bottom pane is the same as the number of paths entering through the bottom pane and bouncing  $n - 1$  bounces more. Entering through the bottom pane is the same as entering through the top pane (but flipped over), so  $E_1$  is the number of paths of light bouncing  $n - 1$  times.
- **Sub-event  $E_2$ :** Let  $E_2$  be the event that the first bounce *is* on the boundary between the two panes. In this case, we consider the two options for the light after the first bounce: it can either leave the glass (in which case we are dealing with the case of having one bounce, and there is only one way for the light to bounce once) or it can bounce yet again on the top of the upper pane, in which case it is equivalent to the light entering from the top with  $n - 2$  bounces to take along its path.

By the rule of sum, we thus get the following recurrence relation for  $F_n$ , the number of paths in which the light can travel with exactly  $n$  bounces. There is exactly one way for the light to travel with no bounces—straight through—and exactly two ways for the light to travel with only one bounce—off the bottom and off the middle. For any  $n > 1$ , there are  $F_{n-1}$  paths where the light bounces off the bottom of the glass, and  $F_{n-2}$  paths where the light bounces off the middle and then off the top.

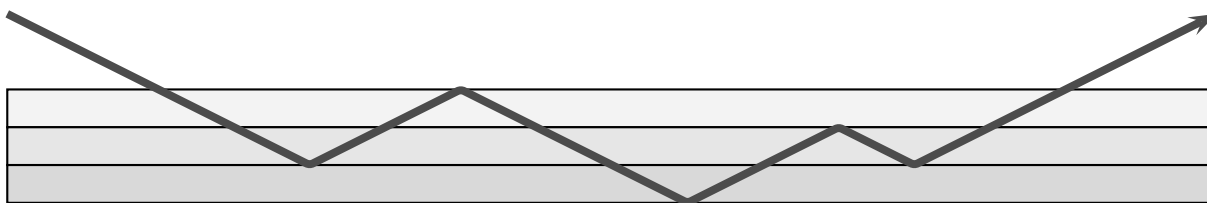
$$F_0 = 1$$

$$F_1 = 2$$

$$F_n = F_{n-1} + F_{n-2}$$

### Stump a professor

What is the recurrence relation for *three* panes of glass? This question once stumped an anonymous professor<sup>1</sup> in a science discipline, but now you should be able to solve it with a bit of effort. Aren't you proud of your knowledge?



<sup>1</sup>Not me! —Jeff

## 2 Sequences, sequence operators, and annihilators

We have shown that several different problems can be expressed in terms of Fibonacci sequences, but we don't yet know how to explicitly compute the  $n$ th Fibonacci number, or even (and more importantly) roughly how big it is. We can easily write a program to compute the  $n$ th Fibonacci number, but that doesn't help us much here. What we really want is a *closed form solution* for the Fibonacci recurrence—an explicit algebraic formula without conditionals, loops, or recursion.

In order to solve recurrences like the Fibonacci recurrence, we first need to understand *operations* on infinite sequences of numbers. Although these sequences are formally defined as *functions* of the form  $A : \mathbb{N} \rightarrow \mathbb{R}$ , we will write them either as  $A = \langle a_0, a_1, a_2, a_3, a_4, \dots \rangle$  when we want to emphasize the entire sequence<sup>2</sup>, or as  $A = \langle a_i \rangle$  when we want to emphasize a generic element. For example, the Fibonacci sequence is  $\langle 0, 1, 1, 2, 3, 5, 8, 13, 21, \dots \rangle$ .

We can naturally define several sequence operators:

- We can add or subtract any two sequences:

$$\begin{aligned}\langle a_i \rangle + \langle b_i \rangle &= \langle a_0, a_1, a_2, \dots \rangle + \langle b_0, b_1, b_2, \dots \rangle = \langle a_0 + b_0, a_1 + b_1, a_2 + b_2, \dots \rangle = \langle a_i + b_i \rangle \\ \langle a_i \rangle - \langle b_i \rangle &= \langle a_0, a_1, a_2, \dots \rangle - \langle b_0, b_1, b_2, \dots \rangle = \langle a_0 - b_0, a_1 - b_1, a_2 - b_2, \dots \rangle = \langle a_i - b_i \rangle\end{aligned}$$

- We can multiply any sequence by a constant:

$$c \cdot \langle a_i \rangle = c \cdot \langle a_0, a_1, a_2, \dots \rangle = \langle c \cdot a_0, c \cdot a_1, c \cdot a_2, \dots \rangle = \langle c \cdot a_i \rangle$$

- We can shift any sequence to the left by removing its initial element:

$$\mathbf{E}\langle a_i \rangle = \mathbf{E}\langle a_0, a_1, a_2, a_3, \dots \rangle = \langle a_1, a_2, a_3, a_4, \dots \rangle = \langle a_{i+1} \rangle$$

**Example:** We can understand these operators better by looking at some specific examples, using the sequence  $T$  of powers of two.

$$\begin{aligned}T &= \langle 2^0, 2^1, 2^2, 2^3, \dots \rangle = \langle 2^i \rangle \\ \mathbf{E}T &= \langle 2^1, 2^2, 2^3, 2^4, \dots \rangle = \langle 2^{i+1} \rangle \\ 2T &= \langle 2 \cdot 2^0, 2 \cdot 2^1, 2 \cdot 2^2, 2 \cdot 2^3, \dots \rangle = \langle 2^1, 2^2, 2^3, 2^4, \dots \rangle = \langle 2^{i+1} \rangle \\ 2T - \mathbf{E}T &= \langle 2^1 - 2^1, 2^2 - 2^2, 2^3 - 2^3, 2^4 - 2^4, \dots \rangle = \langle 0, 0, 0, 0, \dots \rangle = \langle 0 \rangle\end{aligned}$$

### 2.1 Properties of operators

It turns out that the distributive property holds for these operators, so we can rewrite  $\mathbf{E}T - 2T$  as  $(\mathbf{E} - 2)T$ . Since  $(\mathbf{E} - 2)T = \langle 0, 0, 0, 0, \dots \rangle$ , we say that the operator  $(\mathbf{E} - 2)$  *annihilates*  $T$ , and we call  $(\mathbf{E} - 2)$  an *annihilator* of  $T$ . Obviously, we can trivially annihilate any sequence by multiplying it by zero, so as a technical matter, we do not consider multiplication by 0 to be an annihilator.

What happens when we apply the operator  $(\mathbf{E} - 3)$  to our sequence  $T$ ?

$$(\mathbf{E} - 3)T = \mathbf{E}T - 3T = \langle 2^{i+1} \rangle - 3\langle 2^i \rangle = \langle 2^{i+1} - 3 \cdot 2^i \rangle = \langle -2^i \rangle = -T$$

The operator  $(\mathbf{E} - 3)$  did very little to our sequence  $T$ ; it just flipped the sign of each number in the sequence. In fact, we will soon see that *only*  $(\mathbf{E} - 2)$  will annihilate  $T$ , and all other simple

<sup>2</sup>It really doesn't matter whether we start a sequence with  $a_0$  or  $a_1$  or  $a_5$  or even  $a_{-17}$ . Zero is often a convenient starting point for many recursively defined sequences, so we'll usually start there.

operators will affect  $T$  in very minor ways. Thus, if we know how to annihilate the sequence, we know what the sequence must look like.

In general,  $(\mathbf{E} - c)$  annihilates any geometric sequence  $A = \langle a_0, a_0c, a_0c^2, a_0c^3, \dots \rangle = \langle a_0c^i \rangle$ :

$$(\mathbf{E} - c)\langle a_0c^i \rangle = \mathbf{E}\langle a_0c^i \rangle - c\langle a_0c^i \rangle = \langle a_0c^{i+1} \rangle - \langle c \cdot a_0c^i \rangle = \langle a_0c^{i+1} - a_0c^{i+1} \rangle = \langle 0 \rangle$$

To see that this is the only operator of this form that annihilates  $A$ , let's see the effect of operator  $(\mathbf{E} - d)$  for some  $d \neq c$ :

$$(\mathbf{E} - d)\langle a_0c^i \rangle = \mathbf{E}\langle a_0c^i \rangle - d\langle a_0c^i \rangle = \langle a_0c^{i+1} - da_0c^i \rangle = \langle (c - d)a_0c^i \rangle = (c - d)\langle a_0c^i \rangle$$

So we have a more rigorous confirmation that an annihilator annihilates exactly one type of sequence, but multiplies other similar sequences by a constant.

We can use this fact about annihilators of geometric sequences to solve certain recurrences. For example, consider the sequence  $R = \langle r_0, r_1, r_2, \dots \rangle$  defined recursively as follows:

$$\begin{aligned} r_0 &= 3 \\ r_{i+1} &= 5r_i \end{aligned}$$

We can easily prove that the operator  $(\mathbf{E} - 5)$  annihilates  $R$ :

$$(\mathbf{E} - 5)\langle r_i \rangle = \mathbf{E}\langle r_i \rangle - 5\langle r_i \rangle = \langle r_{i+1} \rangle - \langle 5r_i \rangle = \langle r_{i+1} - 5r_i \rangle = \langle 0 \rangle$$

Since  $(\mathbf{E} - 5)$  is an annihilator for  $R$ , we must have the closed form solution  $r_i = r_05^i = 3 \cdot 5^i$ . We can easily verify this by induction, as follows:

$$\begin{aligned} r_0 &= 3 \cdot 5^0 = 3 && \text{[definition]} \\ r_i &= 5r_{i-1} && \text{[definition]} \\ &= 5 \cdot (3 \cdot 5^{i-1}) && \text{[induction hypothesis]} \\ &= 5^i \cdot 3 && \text{[algebra]} \end{aligned}$$

## 2.2 Multiple operators

An operator is a function that transforms one sequence into another. Like any other function, we can apply operators one after another to the same sequence. For example, we can multiply a sequence  $\langle a_i \rangle$  by a constant  $d$  and then by a constant  $c$ , resulting in the sequence  $c(d\langle a_i \rangle) = \langle c \cdot d \cdot a_i \rangle = (cd)\langle a_i \rangle$ . Alternatively, we may multiply the sequence by a constant  $c$  and then shift it to the left to get  $\mathbf{E}(c\langle a_i \rangle) = \mathbf{E}\langle c \cdot a_i \rangle = \langle c \cdot a_{i+1} \rangle$ . This is exactly the same as applying the operators in the reverse order:  $c(\mathbf{E}\langle a_i \rangle) = c\langle a_{i+1} \rangle = \langle c \cdot a_{i+1} \rangle$ . We can also shift the sequence twice to the left:  $\mathbf{E}(\mathbf{E}\langle a_i \rangle) = \mathbf{E}\langle a_{i+1} \rangle = \langle a_{i+2} \rangle$ . We will write this in shorthand as  $\mathbf{E}^2\langle a_i \rangle$ . More generally, the operator  $\mathbf{E}^k$  shifts a sequence  $k$  steps to the left:  $\mathbf{E}^k\langle a_i \rangle = \langle a_{i+k} \rangle$ .

We now have the tools to solve a whole host of recurrence problems. For example, what annihilates  $C = \langle 2^i + 3^i \rangle$ ? Well, we know that  $(\mathbf{E} - 2)$  annihilates  $\langle 2^i \rangle$  while leaving  $\langle 3^i \rangle$  essentially unscathed. Similarly,  $(\mathbf{E} - 3)$  annihilates  $\langle 3^i \rangle$  while leaving  $\langle 2^i \rangle$  essentially unscathed. Thus, if we apply both operators one after the other, we see that  $(\mathbf{E} - 2)(\mathbf{E} - 3)$  annihilates our sequence  $C$ .

In general, for any integers  $a \neq b$ , the operator  $(\mathbf{E} - a)(\mathbf{E} - b)$  annihilates any sequence of the form  $\langle c_1a^i + c_2b^i \rangle$  but nothing else. We will often 'multiply out' the operators into the shorthand notation  $\mathbf{E}^2 - (a + b)\mathbf{E} + ab$ . It is left as an exhilarating exercise to the student to verify that this



shorthand actually makes sense—the operators  $(\mathbf{E} - a)(\mathbf{E} - b)$  and  $\mathbf{E}^2 - (a + b)\mathbf{E} + ab$  have the same effect on every sequence.

We now know finally enough to solve the recurrence for Fibonacci numbers. Specifically, notice that the recurrence  $F_i = F_{i-1} + F_{i-2}$  is annihilated by  $\mathbf{E}^2 - \mathbf{E} - 1$ :

$$\begin{aligned} (\mathbf{E}^2 - \mathbf{E} - 1)\langle F_i \rangle &= \mathbf{E}^2\langle F_i \rangle - \mathbf{E}\langle F_i \rangle - \langle F_i \rangle \\ &= \langle F_{i+2} \rangle - \langle F_{i+1} \rangle - \langle F_i \rangle \\ &= \langle F_{i-2} - F_{i-1} - F_i \rangle \\ &= \langle 0 \rangle \end{aligned}$$

Factoring  $\mathbf{E}^2 - \mathbf{E} - 1$  using the quadratic formula, we obtain

$$\mathbf{E}^2 - \mathbf{E} - 1 = (\mathbf{E} - \phi)(\mathbf{E} - \hat{\phi})$$

where  $\phi = (1 + \sqrt{5})/2 \approx 1.618034$  is the golden ratio and  $\hat{\phi} = (1 - \sqrt{5})/2 = 1 - \phi = -1/\phi$ . Thus, the operator  $(\mathbf{E} - \phi)(\mathbf{E} - \hat{\phi})$  annihilates the Fibonacci sequence, so  $F_i$  must have the form

$$F_i = c\phi^i + \hat{c}\hat{\phi}^i$$

for some constants  $c$  and  $\hat{c}$ . We call this the *generic solution* to the recurrence, since it doesn't depend at all on the base cases. To compute the constants  $c$  and  $\hat{c}$ , we use the base cases  $F_0 = 0$  and  $F_1 = 1$  to obtain a pair of linear equations:

$$\begin{aligned} F_0 = 0 &= c + \hat{c} \\ F_1 = 1 &= c\phi + \hat{c}\hat{\phi} \end{aligned}$$

Solving this system of equations gives us  $c = 1/(2\phi - 1) = 1/\sqrt{5}$  and  $\hat{c} = -1/\sqrt{5}$ .

We now have a closed-form expression for the  $i$ th Fibonacci number:

$$F_i = \frac{\phi^i - \hat{\phi}^i}{\sqrt{5}} = \frac{1}{\sqrt{5}} \left( \frac{1 + \sqrt{5}}{2} \right)^i - \frac{1}{\sqrt{5}} \left( \frac{1 - \sqrt{5}}{2} \right)^i$$

With all the square roots in this formula, it's quite amazing that Fibonacci numbers are integers. However, if we do all the math correctly, all the square roots cancel out when  $i$  is an integer. (In fact, this is pretty easy to prove using the binomial theorem.)

### 2.3 Degenerate cases

We can't quite solve *every* recurrence yet. In our above formulation of  $(\mathbf{E} - a)(\mathbf{E} - b)$ , we assumed that  $a \neq b$ . What about the operator  $(\mathbf{E} - a)(\mathbf{E} - a) = (\mathbf{E} - a)^2$ ? It turns out that this operator annihilates sequences such as  $\langle ia^i \rangle$ :

$$\begin{aligned} (\mathbf{E} - a)\langle ia^i \rangle &= \langle (i + 1)a^{i+1} - (a)ia^i \rangle \\ &= \langle (i + 1)a^{i+1} - ia^{i+1} \rangle \\ &= \langle a^{i+1} \rangle \end{aligned}$$

$$(\mathbf{E} - a)^2\langle ia^i \rangle = (\mathbf{E} - a)\langle a^{i+1} \rangle = \langle 0 \rangle$$

More generally, the operator  $(\mathbf{E} - a)^k$  annihilates any sequence  $\langle p(i) \cdot a^i \rangle$ , where  $p(i)$  is any polynomial in  $i$  of degree  $k - 1$ . As an example,  $(\mathbf{E} - 1)^3$  annihilates the sequence  $\langle i^2 \cdot 1^i \rangle = \langle i^2 \rangle = \langle 1, 4, 9, 16, 25, \dots \rangle$ , since  $p(i) = i^2$  is a polynomial of degree  $n - 1 = 2$ .

As a review, try to explain the following statements:

- $(\mathbf{E} - 1)$  annihilates any constant sequence  $\langle \alpha \rangle$ .
- $(\mathbf{E} - 1)^2$  annihilates any arithmetic sequence  $\langle \alpha + \beta i \rangle$ .
- $(\mathbf{E} - 1)^3$  annihilates any quadratic sequence  $\langle \alpha + \beta i + \gamma i^2 \rangle$ .
- $(\mathbf{E} - 3)(\mathbf{E} - 2)(\mathbf{E} - 1)$  annihilates any sequence  $\langle \alpha + \beta 2^i + \gamma 3^i \rangle$ .
- $(\mathbf{E} - 3)^2(\mathbf{E} - 2)(\mathbf{E} - 1)$  annihilates any sequence  $\langle \alpha + \beta 2^i + \gamma 3^i + \delta i 3^i \rangle$ .

## 2.4 Summary

In summary, we have learned several operators that act on sequences, as well as a few ways of combining operators.

Operator	Definition
Addition	$\langle a_i \rangle + \langle b_i \rangle = \langle a_i + b_i \rangle$
Subtraction	$\langle a_i \rangle - \langle b_i \rangle = \langle a_i - b_i \rangle$
Scalar multiplication	$c\langle a_i \rangle = \langle ca_i \rangle$
Shift	$\mathbf{E}\langle a_i \rangle = \langle a_{i+1} \rangle$
Composition of operators	$(\mathbf{X} + \mathbf{Y})\langle a_i \rangle = \mathbf{X}\langle a_i \rangle + \mathbf{Y}\langle a_i \rangle$ $(\mathbf{X} - \mathbf{Y})\langle a_i \rangle = \mathbf{X}\langle a_i \rangle - \mathbf{Y}\langle a_i \rangle$ $\mathbf{X}\mathbf{Y}\langle a_i \rangle = \mathbf{X}(\mathbf{Y}\langle a_i \rangle) = \mathbf{Y}(\mathbf{X}\langle a_i \rangle)$
$k$ -fold shift	$\mathbf{E}^k\langle a_i \rangle = \langle a_{i+k} \rangle$

Notice that we have not defined a multiplication operator for two sequences. This is usually accomplished by *convolution*:

$$\langle a_i \rangle * \langle b_i \rangle = \left\langle \sum_{j=0}^i a_j b_{i-j} \right\rangle.$$

Fortunately, convolution is unnecessary for solving the recurrences we will see in this course.

We have also learned some things about annihilators, which can be summarized as follows:

Sequence	Annihilator
$\langle \alpha \rangle$	$\mathbf{E} - 1$
$\langle \alpha a^i \rangle$	$\mathbf{E} - a$
$\langle \alpha a^i + \beta b^i \rangle$	$(\mathbf{E} - a)(\mathbf{E} - b)$
$\langle \alpha_0 a_0^i + \alpha_1 a_1^i + \cdots + \alpha_n a_n^i \rangle$	$(\mathbf{E} - a_0)(\mathbf{E} - a_1) \cdots (\mathbf{E} - a_n)$
$\langle \alpha i + \beta \rangle$	$(\mathbf{E} - 1)^2$
$\langle (\alpha i + \beta) a^i \rangle$	$(\mathbf{E} - a)^2$
$\langle (\alpha i + \beta) a^i + \gamma b^i \rangle$	$(\mathbf{E} - a)^2(\mathbf{E} - b)$
$\langle (\alpha_0 + \alpha_1 i + \cdots + \alpha_{n-1} i^{n-1}) a^i \rangle$	$(\mathbf{E} - a)^n$
If $\mathbf{X}$ annihilates $\langle a_i \rangle$ , then $\mathbf{X}$ also annihilates $c\langle a_i \rangle$ for any constant $c$ .	
If $\mathbf{X}$ annihilates $\langle a_i \rangle$ and $\mathbf{Y}$ annihilates $\langle b_i \rangle$ , then $\mathbf{X}\mathbf{Y}$ annihilates $\langle a_i \rangle \pm \langle b_i \rangle$ .	

## 3 Solving Linear Recurrences

### 3.1 Homogeneous Recurrences

The general expressions in the annihilator box above are really the most important things to remember about annihilators because they help you to solve any recurrence for which you can write down an annihilator. The general method is:

1. Write down the annihilator for the recurrence
2. Factor the annihilator
3. Determine the sequence annihilated by each factor
4. Add these sequences together to form the generic solution
5. Solve for constants of the solution by using initial conditions

**Example:** Let's show the steps required to solve the following recurrence:

$$\begin{aligned}r_0 &= 1 \\r_1 &= 5 \\r_2 &= 17 \\r_i &= 7r_{i-1} - 16r_{i-2} + 12r_{i-3}\end{aligned}$$

1. *Write down the annihilator.* Since  $r_{i+3} - 7r_{i+2} + 16r_{i+1} - 12r_i = 0$ , the annihilator is  $\mathbf{E}^3 - 7\mathbf{E}^2 + 16\mathbf{E} - 12$ .
2. *Factor the annihilator.*  $\mathbf{E}^3 - 7\mathbf{E}^2 + 16\mathbf{E} - 12 = (\mathbf{E} - 2)^2(\mathbf{E} - 3)$ .
3. *Determine sequences annihilated by each factor.*  $(\mathbf{E} - 2)^2$  annihilates  $\langle(\alpha i + \beta)2^i\rangle$  for any constants  $\alpha$  and  $\beta$ , and  $(\mathbf{E} - 3)$  annihilates  $\langle\gamma 3^i\rangle$  for any constant  $\gamma$ .
4. *Combine the sequences.*  $(\mathbf{E} - 2)^2(\mathbf{E} - 3)$  annihilates  $\langle(\alpha i + \beta)2^i + \gamma 3^i\rangle$  for any constants  $\alpha, \beta, \gamma$ .
5. *Solve for the constants.* The base cases give us three equations in the three unknowns  $\alpha, \beta, \gamma$ :

$$\begin{aligned}r_0 = 1 &= (\alpha \cdot 0 + \beta)2^0 + \gamma \cdot 3^0 = \beta + \gamma \\r_1 = 5 &= (\alpha \cdot 1 + \beta)2^1 + \gamma \cdot 3^1 = 2\alpha + 2\beta + 3\gamma \\r_2 = 17 &= (\alpha \cdot 2 + \beta)2^2 + \gamma \cdot 3^2 = 8\alpha + 4\beta + 9\gamma\end{aligned}$$

We can solve these equations to get  $\alpha = 1$ ,  $\beta = 0$ ,  $\gamma = 1$ . Thus, our final solution is  $r_i = i2^i + 3^i$ , which we can verify by induction.

### 3.2 Non-homogeneous Recurrences

A *height balanced tree* is a binary tree, where the heights of the two subtrees of the root differ by at most one, and both subtrees are also height balanced. To ground the recursive definition, the empty set is considered a height balanced tree of height  $-1$ , and a single node is a height balanced tree of height  $0$ .

Let  $T_n$  be the smallest height-balanced tree of height  $n$ —how many nodes does  $T_n$  have? Well, one of the subtrees of  $T_n$  has height  $n - 1$  (since  $T_n$  has height  $n$ ) and the other has height either  $n - 1$  or  $n - 2$  (since  $T_n$  is height-balanced and as small as possible). Since both subtrees are themselves height-balanced, the two subtrees must be  $T_{n-1}$  and  $T_{n-2}$ .

We have just derived the following recurrence for  $t_n$ , the number of nodes in the tree  $T_n$ :

$$\begin{aligned}t_{-1} &= 0 && \text{[the empty set]} \\t_0 &= 1 && \text{[a single node]} \\t_n &= t_{n-1} + t_{n-2} + 1\end{aligned}$$

The final '+1' is for the root of  $T_n$ .

We refer to the terms in the equation involving  $t_i$ 's as the *homogeneous* terms and the rest as the *non-homogeneous* terms. (If there were no non-homogeneous terms, we would say that the recurrence itself is homogeneous.) We know that  $\mathbf{E}^2 - \mathbf{E} - 1$  annihilates the homogeneous part  $t_n = t_{n-1} + t_{n-2}$ . Let us try applying this annihilator to the entire equation:

$$\begin{aligned} (\mathbf{E}^2 - \mathbf{E} - 1)\langle t_i \rangle &= \mathbf{E}^2\langle t_i \rangle - \mathbf{E}\langle a_i \rangle - 1\langle a_i \rangle \\ &= \langle t_{i+2} \rangle - \langle t_{i+1} \rangle - \langle t_i \rangle \\ &= \langle t_{i+2} - t_{i+1} - t_i \rangle \\ &= \langle 1 \rangle \end{aligned}$$

The leftover sequence  $\langle 1, 1, 1, \dots \rangle$  is called the *residue*. To obtain the annihilator for the entire recurrence, we compose the annihilator for its homogeneous part with the annihilator of its residue. Since  $\mathbf{E} - 1$  annihilates  $\langle 1 \rangle$ , it follows that  $(\mathbf{E}^2 - \mathbf{E} - 1)(\mathbf{E} - 1)$  annihilates  $\langle t_n \rangle$ . We can factor the annihilator into

$$(\mathbf{E} - \phi)(\mathbf{E} - \hat{\phi})(\mathbf{E} - 1),$$

so our annihilator rules tell us that

$$t_n = \alpha\phi^n + \beta\hat{\phi}^n + \gamma$$

for some constants  $\alpha, \beta, \gamma$ . We call this the *generic solution* to the recurrence. Different recurrences can have the same generic solution.

To solve for the unknown constants, we need three equations in three unknowns. Our base cases give us two equations, and we can get a third by examining the next nontrivial case  $t_1 = 2$ :

$$\begin{aligned} t_{-1} = 0 &= \alpha\phi^{-1} + \beta\hat{\phi}^{-1} + \gamma = \alpha/\phi + \beta/\hat{\phi} + \gamma \\ t_0 = 1 &= \alpha\phi^0 + \beta\hat{\phi}^0 + \gamma = \alpha + \beta + \gamma \\ t_1 = 2 &= \alpha\phi^1 + \beta\hat{\phi}^1 + \gamma = \alpha\phi + \beta\hat{\phi} + \gamma \end{aligned}$$

Solving these equations, we find that  $\alpha = \frac{\sqrt{5}+2}{\sqrt{5}}$ ,  $\beta = \frac{\sqrt{5}-2}{\sqrt{5}}$ , and  $\gamma = -1$ . Thus,

$$t_n = \frac{\sqrt{5}+2}{\sqrt{5}} \left( \frac{1+\sqrt{5}}{2} \right)^n + \frac{\sqrt{5}-2}{\sqrt{5}} \left( \frac{1-\sqrt{5}}{2} \right)^n - 1$$

Here is the general method for non-homogeneous recurrences:

1. Write down the homogeneous annihilator, directly from the recurrence  $1\frac{1}{2}$ . 'Multiply' by the annihilator for the residue
2. Factor the annihilator
3. Determine what sequence each factor annihilates
4. Add these sequences together to form the generic solution
5. Solve for constants of the solution by using initial conditions

### 3.3 Some more examples

In each example below, we use the base cases  $a_0 = 0$  and  $a_1 = 1$ .

- $a_n = a_{n-1} + a_{n-2} + 2$

- The homogeneous annihilator is  $\mathbf{E}^2 - \mathbf{E} - 1$ .
- The residue is the constant sequence  $\langle 2, 2, 2, \dots \rangle$ , which is annihilated by  $\mathbf{E} - 1$ .
- Thus, the annihilator is  $(\mathbf{E}^2 - \mathbf{E} - 1)(\mathbf{E} - 1)$ .
- The annihilator factors into  $(\mathbf{E} - \phi)(\mathbf{E} - \hat{\phi})(\mathbf{E} - 1)$ .
- Thus, the generic solution is  $a_n = \alpha\phi^n + \beta\hat{\phi}^n + \gamma$ .
- The constants  $\alpha, \beta, \gamma$  satisfy the equations

$$\begin{aligned} a_0 = 0 &= \alpha + \beta + \gamma \\ a_1 = 1 &= \alpha\phi + \beta\hat{\phi} + \gamma \\ a_2 = 3 &= \alpha\phi^2 + \beta\hat{\phi}^2 + \gamma \end{aligned}$$

- Solving the equations gives us  $\alpha = \frac{\sqrt{5}+2}{\sqrt{5}}$ ,  $\beta = \frac{\sqrt{5}-2}{\sqrt{5}}$ , and  $\gamma = -2$

- So the final solution is 
$$a_n = \frac{\sqrt{5}+2}{\sqrt{5}} \left( \frac{1+\sqrt{5}}{2} \right)^n + \frac{\sqrt{5}-2}{\sqrt{5}} \left( \frac{1-\sqrt{5}}{2} \right)^n - 2$$

(In the remaining examples, I won't explicitly enumerate the steps like this.)

- $a_n = a_{n-1} + a_{n-2} + 3$

The homogeneous annihilator  $(\mathbf{E}^2 - \mathbf{E} - 1)$  leaves a constant residue  $\langle 3, 3, 3, \dots \rangle$ , so the annihilator is  $(\mathbf{E}^2 - \mathbf{E} - 1)(\mathbf{E} - 1)$ , and the generic solution is  $a_n = \alpha\phi^n + \beta\hat{\phi}^n + \gamma$ . Solving the equations

$$\begin{aligned} a_0 = 0 &= \alpha + \beta + \gamma \\ a_1 = 1 &= \alpha\phi + \beta\hat{\phi} + \gamma \\ a_2 = 4 &= \alpha\phi^2 + \beta\hat{\phi}^2 + \gamma \end{aligned}$$

gives us the final solution 
$$a_n = \frac{\sqrt{5}+3}{\sqrt{5}} \left( \frac{1+\sqrt{5}}{2} \right)^n + \frac{\sqrt{5}-3}{\sqrt{5}} \left( \frac{1-\sqrt{5}}{2} \right)^n - 3$$

- $a_n = a_{n-1} + a_{n-2} + 2^n$

The homogeneous annihilator  $(\mathbf{E}^2 - \mathbf{E} - 1)$  leaves an exponential residue  $\langle 4, 8, 16, 32, \dots \rangle = \langle 2^{i+2} \rangle$ , which is annihilated by  $\mathbf{E} - 2$ . Thus, the annihilator is  $(\mathbf{E}^2 - \mathbf{E} - 1)(\mathbf{E} - 2)$ , and the generic solution is  $a_n = \alpha\phi^n + \beta\hat{\phi}^n + \gamma 2^n$ . The constants  $\alpha, \beta, \gamma$  satisfy the following equations:

$$\begin{aligned} a_0 = 0 &= \alpha + \beta + \gamma \\ a_1 = 1 &= \alpha\phi + \beta\hat{\phi} + 2\gamma \\ a_2 = 5 &= \alpha\phi^2 + \beta\hat{\phi}^2 + 4\gamma \end{aligned}$$

- $a_n = a_{n-1} + a_{n-2} + n$

The homogeneous annihilator  $(\mathbf{E}^2 - \mathbf{E} - 1)$  leaves a linear residue  $\langle 2, 3, 4, 5 \dots \rangle = \langle i + 2 \rangle$ , which is annihilated by  $(\mathbf{E} - 1)^2$ . Thus, the annihilator is  $(\mathbf{E}^2 - \mathbf{E} - 1)(\mathbf{E} - 1)^2$ , and the generic solution is  $a_n = \alpha\phi^n + \beta\hat{\phi}^n + \gamma + \delta n$ . The constants  $\alpha, \beta, \gamma, \delta$  satisfy the following equations:

$$\begin{aligned} a_0 = 0 &= \alpha + \beta + \gamma \\ a_1 = 1 &= \alpha\phi + \beta\hat{\phi} + \gamma + \delta \\ a_2 = 3 &= \alpha\phi^2 + \beta\hat{\phi}^2 + \gamma + 2\delta \\ a_3 = 7 &= \alpha\phi^3 + \beta\hat{\phi}^3 + \gamma + 3\delta \end{aligned}$$

- $a_n = a_{n-1} + a_{n-2} + n^2$

The homogeneous annihilator  $(\mathbf{E}^2 - \mathbf{E} - 1)$  leaves a quadratic residue  $\langle 4, 9, 16, 25 \dots \rangle = \langle (i + 2)^2 \rangle$ , which is annihilated by  $(\mathbf{E} - 1)^3$ . Thus, the annihilator is  $(\mathbf{E}^2 - \mathbf{E} - 1)(\mathbf{E} - 1)^3$ , and the generic solution is  $a_n = \alpha\phi^n + \beta\hat{\phi}^n + \gamma + \delta n + \varepsilon n^2$ . The constants  $\alpha, \beta, \gamma, \delta, \varepsilon$  satisfy the following equations:

$$\begin{aligned} a_0 = 0 &= \alpha + \beta + \gamma \\ a_1 = 1 &= \alpha\phi + \beta\hat{\phi} + \gamma + \delta + \varepsilon \\ a_2 = 5 &= \alpha\phi^2 + \beta\hat{\phi}^2 + \gamma + 2\delta + 4\varepsilon \\ a_3 = 15 &= \alpha\phi^3 + \beta\hat{\phi}^3 + \gamma + 3\delta + 9\varepsilon \\ a_4 = 36 &= \alpha\phi^4 + \beta\hat{\phi}^4 + \gamma + 4\delta + 16\varepsilon \end{aligned}$$

- $a_n = a_{n-1} + a_{n-2} + n^2 - 2^n$

The homogeneous annihilator  $(\mathbf{E}^2 - \mathbf{E} - 1)$  leaves the residue  $\langle (i + 2)^2 - 2^{i-2} \rangle$ . The quadratic part of the residue is annihilated by  $(\mathbf{E} - 1)^3$ , and the exponential part is annihilated by  $(\mathbf{E} - 2)$ . Thus, the annihilator for the whole recurrence is  $(\mathbf{E}^2 - \mathbf{E} - 1)(\mathbf{E} - 1)^3(\mathbf{E} - 2)$ , and so the generic solution is  $a_n = \alpha\phi^n + \beta\hat{\phi}^n + \gamma + \delta n + \varepsilon n^2 + \eta 2^i$ . The constants  $\alpha, \beta, \gamma, \delta, \varepsilon, \eta$  satisfy a system of six equations in six unknowns determined by  $a_0, a_1, \dots, a_5$ .

- $a_n = a_{n-1} + a_{n-2} + \phi^n$

The annihilator is  $(\mathbf{E}^2 - \mathbf{E} - 1)(\mathbf{E} - \phi) = (\mathbf{E} - \phi)^2(\mathbf{E} - \hat{\phi})$ , so the generic solution is  $a_n = \alpha\phi^n + \beta n\phi^n + \gamma\hat{\phi}^n$ . (Other recurrence solving methods will have a “interference” problem with this equation, while the operator method does not.)

Our method does not work on recurrences like  $a_n = a_{n-1} + \frac{1}{n}$  or  $a_n = a_{n-1} + \lg n$ , because the functions  $\frac{1}{n}$  and  $\lg n$  do not have annihilators. Our tool, as it stands, is limited to linear recurrences.

## 4 Divide and Conquer Recurrences and the Master Theorem

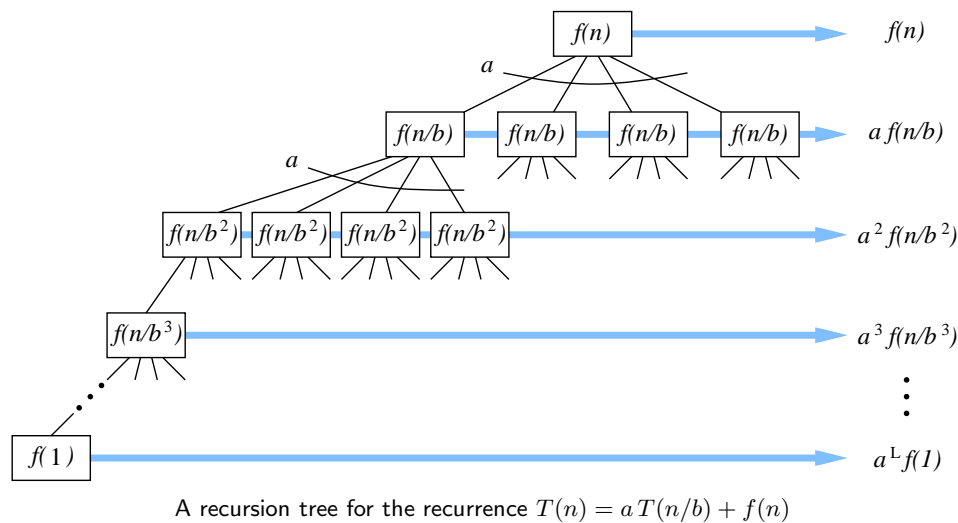
Divide and conquer algorithms often give us running-time recurrences of the form

$$T(n) = aT(n/b) + f(n) \tag{1}$$

where  $a$  and  $b$  are constants and  $f(n)$  is some other function. The so-called ‘Master Theorem’ gives us a general method for solving such recurrences  $f(n)$  is a simple polynomial.

Unfortunately, the Master Theorem doesn't work for all functions  $f(n)$ , and many useful recurrences don't look like (1) at all. Fortunately, there's a general technique to solve most divide-and-conquer recurrences, even if they don't have this form. This technique is used to *prove* the Master Theorem, so if you remember this technique, you can forget the Master Theorem entirely (which is what I did). Throw off your chains!

I'll illustrate the technique using the generic recurrence (1). We start by drawing a *recursion tree*. The root of the recursion tree is a box containing the value  $f(n)$ , it has  $a$  children, each of which is the root of a recursion tree for  $T(n/b)$ . Equivalently, a recursion tree is a complete  $a$ -ary tree where each node at depth  $i$  contains the value  $a^i f(n/b^i)$ . The recursion stops when we get to the base case(s) of the recurrence. Since we're looking for asymptotic bounds, it turns out not to matter much what we use for the base case; for purposes of illustration, I'll assume that  $T(1) = f(1)$ .



Now  $T(n)$  is just the sum of all values stored in the tree. Assuming that each level of the tree is full, we have

$$T(n) = f(n) + a f(n/b) + a^2 f(n/b^2) + \dots + a^i f(n/b^i) + \dots + a^L f(n/b^L)$$

where  $L$  is the depth of the recursion tree. We easily see that  $L = \log_b n$ , since  $n/b^L = 1$ . Since  $f(1) = \Theta(1)$ , the last non-zero term in the summation is  $\Theta(a^L) = \Theta(a^{\log_b n}) = \Theta(n^{\log_b a})$ .

Now we can easily state and prove the Master Theorem, in a slightly different form than it's usually stated.

**The Master Theorem.** The recurrence  $T(n) = aT(n/b) + f(n)$  can be solved as follows.

- If  $a f(n/b) = \kappa f(n)$  for some constant  $\kappa < 1$ , then  $T(n) = \Theta(f(n))$ .
- If  $a f(n/b) = K f(n)$  for some constant  $K > 1$ , then  $T(n) = \Theta(n^{\log_b a})$ .
- If  $a f(n/b) = f(n)$ , then  $T(n) = \Theta(f(n) \log_b n)$ .

**Proof:** If  $f(n)$  is a *constant factor larger* than  $a f(b/n)$ , then by induction, the sum is a descending geometric series. The sum of any geometric series is a constant times its largest term. In this case, the largest term is the first term  $f(n)$ .

If  $f(n)$  is a *constant factor smaller* than  $a f(b/n)$ , then by induction, the sum is an ascending geometric series. The sum of any geometric series is a constant times its largest term. In this case, this is the last term, which by our earlier argument is  $\Theta(n^{\log_b a})$ .

Finally, if  $a f(b/n) = f(n)$ , then by induction, each of the  $L + 1$  terms in the summation is equal to  $f(n)$ .  $\square$

Here are a few canonical examples of the Master Theorem in action:

- **Randomized selection:**  $T(n) = T(3n/4) + n$

Here  $a f(n/b) = 3n/4$  is smaller than  $f(n) = n$  by a factor of  $4/3$ , so  $T(n) = \Theta(n)$

- **Karatsuba's multiplication algorithm:**  $T(n) = 3T(n/2) + n$

Here  $a f(n/b) = 3n/2$  is bigger than  $f(n) = n$  by a factor of  $3/2$ , so  $T(n) = \Theta(n^{\log_2 3})$

- **Mergesort:**  $T(n) = 2T(n/2) + n$

Here  $a f(n/b) = f(n)$ , so  $T(n) = \Theta(n \log n)$

- $T(n) = 4T(n/2) + n \lg n$

In this case, we have  $a f(n/b) = 2n \lg n - 2n$ , which is not quite twice  $f(n) = n \lg n$ . However, for sufficiently large  $n$  (which is all we care about with asymptotic bounds) we have  $2f(n) > a f(n/b) > 1.9f(n)$ . Since the level sums are bounded both above and below by ascending geometric series, the solution is  $T(n) = \Theta(n^{\log_2 4}) = \Theta(n^2)$ . (This trick will *not* work in the second or third cases of the Master Theorem!)

Using the same recursion-tree technique, we can also solve recurrences where the Master Theorem doesn't apply.

- $T(n) = 2T(n/2) + n/\lg n$

We can't apply the Master Theorem here, because  $a f(n/b) = n/(\lg n - 1)$  isn't equal to  $f(n) = n/\lg n$ , but the difference isn't a constant factor. So we need to compute each of the level sums and compute their total in some other way. It's not hard to see that the sum of all the nodes in the  $i$ th level is  $n/(\lg n - i)$ . In particular, this means the depth of the tree is at most  $\lg n - 1$ .

$$T(n) = \sum_{i=0}^{\lg n - 1} \frac{n}{\lg n - i} = \sum_{j=1}^{\lg n} \frac{n}{j} = n H_{\lg n} = \Theta(n \lg \lg n)$$

- **Randomized quicksort:**  $T(n) = T(3n/4) + T(n/4) + n$

In this case, nodes in the same level of the recursion tree have different values. This makes the tree lopsided; different leaves are at different levels. However, it's not hard to see that the nodes in any *complete* level (*i.e.*, above any of the leaves) sum to  $n$ , so this is like the last case of the Master Theorem, and that every leaf has depth between  $\log_4 n$  and  $\log_{4/3} n$ . To derive an upper bound, we overestimate  $T(n)$  by ignoring the base cases and extending the tree downward to the level of the *deepest* leaf. Similarly, to derive a lower bound, we overestimate  $T(n)$  by counting only nodes in the tree up to the level of the *shallowest* leaf. These observations give us the upper and lower bounds  $n \log_4 n \leq T(n) \leq n \log_{4/3} n$ . Since these bound differ by only a constant factor, we have  $T(n) = \Theta(n \log n)$ .



- **Deterministic selection:**  $T(n) = T(n/5) + T(7n/10) + n$

Again, we have a lopsided recursion tree. If we look only at complete levels of the tree, we find that the level sums form a descending geometric series  $T(n) = n + 9n/10 + 81n/100 + \dots$ , so this is like the first case of the master theorem. We can get an upper bound by ignoring the base cases entirely and growing the tree out to infinity, and we can get a lower bound by only counting nodes in complete levels. Either way, the geometric series is dominated by its largest term, so  $T(n) = \Theta(n)$ .

- $T(n) = \sqrt{n} \cdot T(\sqrt{n}) + n$

In this case, we have a complete recursion tree, but the *degree* of the nodes is no longer constant, so we have to be a bit more careful. It's not hard to see that the nodes in any level sum to  $n$ , so this is like the third Master case. The depth  $L$  satisfies the identity  $n^{2^{-L}} = 2$  (we can't get all the way down to 1 by taking square roots), so  $L = \lg \lg n$  and  $T(n) = \Theta(n \lg \lg n)$ .

- $T(n) = 2\sqrt{n} \cdot T(\sqrt{n}) + n$

We still have at most  $\lg \lg n$  levels, but now the nodes in level  $i$  sum to  $2^i n$ . We have an increasing geometric series of level sums, like the second Master case, so  $T(n)$  is dominated by the sum over the deepest level:  $T(n) = \Theta(2^{\lg \lg n} n) = \Theta(n \log n)$

- $T(n) = 4\sqrt{n} \cdot T(\sqrt{n}) + n$

Now the nodes in level  $i$  sum to  $4^i n$ . Again, we have an increasing geometric series, like the second Master case, so we only care about the leaves:  $T(n) = \Theta(4^{\lg \lg n} n) = \Theta(n \log^2 n)$  Ick!

## 5 Transforming Recurrences

### 5.1 An analysis of mergesort: domain transformation

Previously we gave the recurrence for mergesort as  $T(n) = 2T(n/2) + n$ , and obtained the solution  $T(n) = \Theta(n \log n)$  using the Master Theorem (or the recursion tree method if you, like me, can't remember the Master Theorem). This is fine if  $n$  is a power of two, but for other values of  $n$ , this recurrence is incorrect. When  $n$  is odd, then the recurrence calls for us to sort a fractional number of elements! Worse yet, if  $n$  is not a power of two, we will *never* reach the base case  $T(1) = 0$ .

To get a recurrence that's valid for *all* integers  $n$ , we need to carefully add ceilings and floors:

$$T(n) = T(\lceil n/2 \rceil) + T(\lfloor n/2 \rfloor) + n.$$

We have almost no hope of getting an exact solution here; the floors and ceilings will eventually kill us. So instead, let's just try to get a tight asymptotic upper bound for  $T(n)$  using a technique called *domain transformation*. A domain transformation rewrites a function  $T(n)$  with a difficult recurrence as a nested function  $S(f(n))$ , where  $f(n)$  is a simple function and  $S()$  has an easier recurrence.

First we overestimate the time bound, once by pretending that the two subproblem sizes are equal, and again to eliminate the ceiling:

$$T(n) \leq 2T(\lceil n/2 \rceil) + n \leq 2T(n/2 + 1) + n.$$

Now we define a new function  $S(n) = T(n + \alpha)$ , where  $\alpha$  is a unknown constant, chosen so that  $S(n)$  satisfies the Master-ready recurrence  $S(n) \leq 2S(n/2) + O(n)$ . To figure out the correct value of  $\alpha$ , we compare two versions of the recurrence for the function  $T(n + \alpha)$ :

$$\begin{aligned} S(n) \leq 2S(n/2) + O(n) &\implies T(n + \alpha) \leq 2T(n/2 + \alpha) + O(n) \\ T(n) \leq 2T(n/2 + 1) + n &\implies T(n + \alpha) \leq 2T((n + \alpha)/2 + 1) + n + \alpha \end{aligned}$$

For these two recurrences to be equal, we need  $n/2 + \alpha = (n + \alpha)/2 + 1$ , which implies that  $\alpha = 2$ . The Master Theorem now tells us that  $S(n) = O(n \log n)$ , so

$$T(n) = S(n - 2) = O((n - 2) \log(n - 2)) = O(n \log n).$$

A similar argument gives a matching lower bound  $T(n) = \Omega(n \log n)$ . So  $\boxed{T(n) = \Theta(n \log n)}$  after all, just as though we had ignored the floors and ceilings from the beginning!

Domain transformations are useful for removing floors, ceilings, and lower order terms from the arguments of any recurrence that otherwise looks like it ought to fit either the Master Theorem or the recursion tree method. But now that we know this, we don't need to bother grinding through the actual gory details!

## 5.2 A less trivial example

There is a data structure in computational geometry called *ham-sandwich trees*, where the cost of doing a certain search operation obeys the recurrence  $\mathbf{T}(n) = \mathbf{T}(n/2) + \mathbf{T}(n/4) + \mathbf{1}$ . This doesn't fit the Master theorem, because the two subproblems have different sizes, and using the recursion tree method only gives us the loose bounds  $\sqrt{n} \ll T(n) \ll n$ .

Domain transformations save the day. If we define the new function  $t(k) = T(2^k)$ , we have a new recurrence

$$t(k) = t(k - 1) + t(k - 2) + 1$$

which should immediately remind you of Fibonacci numbers. Sure enough, after a bit of work, the annihilator method gives us the solution  $t(k) = \Theta(\phi^k)$ , where  $\phi = (1 + \sqrt{5})/2$  is the golden ratio. This implies that

$$T(n) = t(\lg n) = \Theta(\phi^{\lg n}) = \boxed{\Theta(n^{\lg \phi})} \approx \Theta(n^{0.69424}).$$

It's possible to solve this recurrence without domain transformations and annihilators—in fact, the inventors of ham-sandwich trees did so—but it's much more difficult.

## 5.3 Secondary recurrences

Consider the recurrence  $\mathbf{T}(n) = \mathbf{2T}(\frac{n}{3} - \mathbf{1}) + \mathbf{n}$  with the base case  $T(1) = 1$ . We already know how to use domain transformations to get the tight asymptotic bound  $T(n) = \Theta(n)$ , but how would we obtain an *exact* solution?

First we need to figure out how the parameter  $n$  changes as we get deeper and deeper into the recurrence. For this we use a *secondary recurrence*. We define a sequence  $n_i$  so that

$$T(n_i) = 2T(n_{i-1}) + n_i,$$

So  $n_i$  is the argument of  $T()$  when we are  $i$  recursion steps away from the base case  $n_0 = 1$ . The original recurrence gives us the following secondary recurrence for  $n_i$ :

$$n_{i-1} = \frac{n_i}{3} - 1 \implies n_i = 3n_{i-1} + 3.$$

The annihilator for this recurrence is  $(\mathbf{E} - 1)(\mathbf{E} - 3)$ , so the generic solution is  $n_i = \alpha 3^i + \beta$ . Plugging in the base cases  $n_0 = 1$  and  $n_1 = 6$ , we get the exact solution

$$n_i = \frac{5}{2} \cdot 3^i - \frac{3}{2}.$$

Notice that our original function  $T(n)$  is only well-defined if  $n = n_i$  for some integer  $i \geq 0$ .

Now to solve the original recurrence, we do a range transformation. If we set  $t_i = T(n_i)$ , we have the recurrence  $t_i = 2t_{i-1} + \frac{5}{2} \cdot 3^i - \frac{3}{2}$ , which by now we can solve using the annihilator method. The annihilator of the recurrence is  $(\mathbf{E} - 2)(\mathbf{E} - 3)(\mathbf{E} - 1)$ , so the generic solution is  $\alpha' 3^i + \beta' 2^i + \gamma'$ . Plugging in the base cases  $t_0 = 1$ ,  $t_1 = 8$ ,  $t_2 = 37$ , we get the exact solution

$$t_i = \frac{15}{2} \cdot 3^i - 8 \cdot 2^i + \frac{3}{2}$$

Finally, we need to substitute to get a solution for the original recurrence in terms of  $n$ , by inverting the solution of the secondary recurrence. If  $n = n_i = \frac{5}{2} \cdot 3^i - \frac{3}{2}$ , then (after a little algebra) we have

$$i = \log_3 \left( \frac{2}{5}n + \frac{3}{5} \right).$$

Substituting this into the expression for  $t_i$ , we get our exact, closed-form solution.

$$\begin{aligned} T(n) &= \frac{15}{2} \cdot 3^i - 8 \cdot 2^i + \frac{3}{2} \\ &= \frac{15}{2} \cdot 3^{\log_3 \left( \frac{2}{5}n + \frac{3}{5} \right)} - 8 \cdot 2^{\log_3 \left( \frac{2}{5}n + \frac{3}{5} \right)} + \frac{3}{2} \\ &= \frac{15}{2} \left( \frac{2}{5}n + \frac{3}{5} \right) - 8 \cdot \left( \frac{2}{5}n + \frac{3}{5} \right)^{\log_3 2} + \frac{3}{2} \\ &= 3n - 8 \cdot \left( \frac{2}{5}n + \frac{3}{5} \right)^{\log_3 2} + 6 \end{aligned}$$

Isn't that special? Now you know why we stick to asymptotic bounds for most recurrences.

## 6 References

Methods for solving recurrences by annihilators, domain transformations, and secondary recurrences are nicely outlined in G. Lueker, Some Techniques for Solving Recurrences, *ACM Computing Surveys* 12(4):419-436, 1980. The master theorem is presented in sections 4.3 and 4.4 of CLR.

Sections 1–3 and 5 of this handout were written by Ed Reingold and Ari Trachtenberg and substantially revised by Jeff Erickson. Section 4 is entirely Jeff's fault.

Partly because of his computational skills, Gerbert, in his later years, was made Pope by Otto the Great, Holy Roman Emperor, and took the name Sylvester II. By this time, his gift in the art of calculating contributed to the belief, commonly held throughout Europe, that he had sold his soul to the devil.

— Dominic Olivastro, *Ancient Puzzles*, 1993

## 0 Introduction (August 29)

### 0.1 What is an algorithm?

This is a course about algorithms, specifically combinatorial algorithms. An algorithm is a set of simple, unambiguous, step-by-step instructions for accomplishing a specific task. Note that the word “computer” doesn’t appear anywhere in this definition; algorithms don’t necessarily have anything to do with computers! For example, here is an algorithm for singing that annoying song ‘99 Bottles of Beer on the Wall’ for arbitrary values of 99:

```

BOTTLESOFBEER( $n$ ):
  For  $i \leftarrow n$  down to 1
    Sing “ $i$  bottles of beer on the wall,  $i$  bottles of beer,”
    Sing “Take one down, pass it around,  $i - 1$  bottles of beer on the wall.”
  Sing “No bottles of beer on the wall, no bottles of beer,”
  Sing “Go to the store, buy some more,  $x$  bottles of beer on the wall.”

```

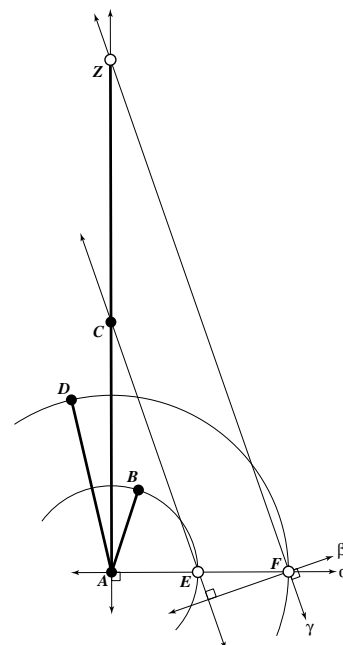
Algorithms have been with us since the dawn of civilization. Here is an algorithm, popularized (but almost certainly not discovered) by Euclid about 2500 years ago, for multiplying or dividing numbers using a ruler and compass. The Greek geometers represented numbers using line segments of the right length. In the pseudo-code below,  $\text{CIRCLE}(p, q)$  represents the circle centered at a point  $p$  and passing through another point  $q$ ; hopefully the other instructions are obvious.

```

«Construct the line perpendicular to  $\ell$  and passing through  $P$ .»
RIGHTANGLE( $\ell, P$ ):
  Choose a point  $A \in \ell$ 
   $A, B \leftarrow \text{INTERSECT}(\text{CIRCLE}(P, A), \ell)$ 
   $C, D \leftarrow \text{INTERSECT}(\text{CIRCLE}(A, B), \text{CIRCLE}(B, A))$ 
  return LINE( $C, D$ )

«Construct a point  $Z$  such that  $|AZ| = |AC| \cdot |AD| / |AB|$ .»
MULTIPLYORDIVIDE( $A, B, C, D$ ):
   $\alpha \leftarrow \text{RIGHTANGLE}(\text{LINE}(A, C), A)$ 
   $E \leftarrow \text{INTERSECT}(\text{CIRCLE}(A, B), \alpha)$ 
   $F \leftarrow \text{INTERSECT}(\text{CIRCLE}(A, D), \alpha)$ 
   $\beta \leftarrow \text{RIGHTANGLE}(\text{LINE}(E, C), F)$ 
   $\gamma \leftarrow \text{RIGHTANGLE}(\beta, F)$ 
  return INTERSECT( $\gamma, \text{LINE}(A, C)$ )

```



Multiplying or dividing using a compass and straight-edge.

This algorithm breaks down the difficult task of multiplication into simple primitive steps: drawing a line between two points, drawing a circle with a given center and boundary point, and so on. The primitive steps need not be quite this primitive, but each primitive step must be something that the person or machine executing the algorithm already knows how to do. Notice in this example that we have made constructing a right angle a primitive operation in the MULTIPLYORDIVIDE algorithm by writing a subroutine.

As a bad example, consider “Martin’s algorithm”:<sup>1</sup>

```
BECOMEAMILLIONAIREANDNEVERPAYTAXES:
```

```
  Get a million dollars.
```

```
  Don't pay taxes.
```

```
  If you get caught,
```

```
    Say “I forgot.”
```

Pretty simple, except for that first step; it’s a doozy. A group of billionaire CEOs would consider this an algorithm, since for them the first step is both unambiguous and trivial. But for the rest of us poor slobs who don’t have a million dollars handy, Martin’s procedure is too vague to be considered an algorithm. [On the other hand, this is a perfect example of a *reduction*—it *reduces* the problem of being a millionaire and never paying taxes to the ‘easier’ problem of acquiring a million dollars. We’ll see reductions over and over again in this class. As hundreds of businessmen and politicians have demonstrated, if you know how to solve the easier problem, a reduction tells you how to solve the harder one.]

Although most of the previous examples are algorithms, they’re not the kind of algorithms that computer scientists are used to thinking about. In this class, we’ll focus (almost!) exclusively on algorithms that can be reasonably implemented on a computer. In other words, each step in the algorithm must be something that either is directly supported by your favorite programming language (arithmetic, assignments, loops, recursion, etc.) or is something that you’ve already learned how to do in an earlier class (sorting, binary search, depth first search, etc.).

For example, here’s the algorithm that’s actually used to determine the number of congressional representatives assigned to each state.<sup>2</sup> The input array  $P[1..n]$  stores the populations of the  $n$  states, and  $R$  is the total number of representatives. (Currently,  $n = 50$  and  $R = 435$ .)

```
APPORTIONCONGRESS( $P[1..n], R$ ):
```

```
   $H \leftarrow \text{NEWMAXHEAP}$ 
```

```
  for  $i \leftarrow 1$  to  $n$ 
```

```
     $r[i] \leftarrow 1$ 
```

```
    INSERT( $H, i, P[i]/\sqrt{2}$ )
```

```
   $R \leftarrow R - n$ 
```

```
  while  $R > 0$ 
```

```
     $s \leftarrow \text{EXTRACTMAX}(H)$ 
```

```
     $r[s] \leftarrow r[s] + 1$ 
```

```
    INSERT( $H, i, P[i]/\sqrt{r[i](r[i] + 1)}$ )
```

```
     $R \leftarrow R - 1$ 
```

```
  return  $r[1..n]$ 
```

<sup>1</sup>S. Martin, “You Can Be A Millionaire”, Saturday Night Live, January 21, 1978. Reprinted in *Comedy Is Not Pretty*, Warner Bros. Records, 1979.

<sup>2</sup>The congressional apportionment algorithm is described in detail at <http://www.census.gov/population/www/censusdata/apportionment/computing.html>, and some earlier algorithms are described at <http://www.census.gov/population/www/censusdata/apportionment/history.html>.

Note that this description assumes that you know how to implement a max-heap and its basic operations NEWMAXHEAP, INSERT, and EXTRACTMAX. Moreover, the correctness of the algorithm doesn't depend at all on how these operations are implemented.<sup>3</sup> The Census Bureau implements the max-heap as an unsorted array, probably inside an Excel spreadsheet. (You should have learned a more efficient solution in CS 225.)

So what's a *combinatorial* algorithm? The distinction is fairly artificial, but basically, this means something distinct from a numerical algorithm. Numerical algorithms are used to approximate computation with ideal real numbers on finite precision computers. For example, here's a numerical algorithm to compute the square root of a number to a given precision. (This algorithm works remarkably quickly—every iteration doubles the number of correct digits.)

<p>SQUAREROOT(<math>x, \varepsilon</math>):</p> <p><math>s \leftarrow 1</math></p> <p>while <math> s - x/s  &gt; \varepsilon</math></p> <p style="padding-left: 2em;"><math>s \leftarrow (s + x/s)/2</math></p> <p>return <math>s</math></p>
--

The output of a numerical algorithm is necessarily an approximation to some ideal mathematical object. Any number that's close enough to the ideal answer is a correct answer. Combinatorial algorithms, on the other hand, manipulate discrete objects like arrays, lists, trees, and graphs that can be represented *exactly* on a digital computer.

## 0.2 Writing down algorithms

Algorithms are *not* programs; they should never be specified in a particular programming language. The whole *point* of this course is to develop computational techniques that can be used in *any* programming language.<sup>4</sup> The idiosyncratic syntactic details of C, Java, Visual Basic, ML, Smalltalk, Intercal, or Brainfunk<sup>5</sup> are of absolutely no importance in algorithm design, and focusing on them will only distract you from what's really going on. What we really want is closer to what you'd write in the *comments* of a real program than the code itself.

On the other hand, a plain English prose description is usually not a good idea either. Algorithms have a lot of structure—especially conditionals, loops, and recursion—that are far too easily hidden by unstructured prose. Like any language spoken by humans, English is full of ambiguities, subtleties, and shades of meaning, but algorithms must be described as accurately as possible. Finally and more seriously, many people have a tendency to describe loops informally: “Do this first, then do this second, and so on.” As anyone who has taken one of those ‘what comes next in this sequence?’ tests already knows, specifying what happens in the first couple of iterations of a loop doesn't say much about what happens later on.<sup>6</sup> Phrases like ‘and so on’ or ‘do X over and over’ or ‘et cetera’ are a good indication that the algorithm *should* have been described in terms of

<sup>3</sup>However, as several students pointed out, the algorithm's correctness does depend on the populations of the states being not too different. If Congress had only 384 representatives, then the 2000 Census data would have given Rhode Island only one representative! If state populations continue to change at their current rates, The United States have have a serious constitutional crisis right after the 2030 Census. I can't wait!

<sup>4</sup>See <http://www.ionet.net/~timtroyr/funhouse/beer.html> for implementations of the BOTTLESOFBEER algorithm in over 200 different programming languages.

<sup>5</sup>Pardon my thinly bowdlerized Anglo-Saxon. Brainfork is the well-deserved name of a programming language invented by Urban Mueller in 1993. Brainflick programs are written entirely using the punctuation characters `<>+-, . []`, each representing a different operation (roughly: shift left, shift right, increment, decrement, input, output, begin loop, end loop). See <http://www.catseye.mb.ca/esoteric/bf/> for a complete definition of Brainfoeey, sample Brainfudge programs, a Brainfiretruck interpreter (written in just 230 characters of C), and related shit.

<sup>6</sup>See <http://www.research.att.com/~njas/sequences/>.

loops or recursion, and the description should have specified what happens in a *generic* iteration of the loop.<sup>7</sup>

The best way to write down an algorithm is using pseudocode. Pseudocode uses the structure of formal programming languages and mathematics to break the algorithm into one-sentence steps, but those sentences can be written using mathematics, pure English, or some mixture of the two. Exactly how to structure the pseudocode is a personal choice, but here are the basic rules I follow:

- Use standard imperative programming keywords (if/then/else, while, for, repeat/until, case, return) and notation (array[index], pointer→field, function(args), etc.)
- The block structure should be visible from across the room. Indent everything carefully and consistently. Don't use syntactic sugar (like C/C++/Java braces or Pascal/Algol begin/end tags) unless the pseudocode is absolutely unreadable without it.
- *Don't* typeset keywords in a different **font** or **style**. Changing type style emphasizes the keywords, making the reader think the syntactic sugar is actually important—it isn't!
- Each statement should fit on one line, each line should contain only one statement. (The only exception is extremely short and similar statements like  $i \leftarrow i + 1$ ;  $j \leftarrow j - 1$ ;  $k \leftarrow 0$ .)
- Put each structuring statement (for, while, if) on its own line. The order of nested loops matters a great deal; make it absolutely obvious.
- Use short but mnemonic algorithm and variable names. Absolutely *never* use pronouns!

A good description of an algorithm reveals the internal structure, hides irrelevant details, and can be implemented easily by any competent programmer in any programming language, even if they don't understand why the algorithm works. Good pseudocode, like good code, makes the algorithm much easier to understand and analyze; it also makes mistakes much easier to spot. The algorithm descriptions in the textbooks and lecture notes are good examples of what we want to see on your homeworks and exams..

### 0.3 Analyzing algorithms

It's not enough just to write down an algorithm and say 'Behold!' We also need to convince ourselves (and our graders) that the algorithm does what it's supposed to do, and that it does it quickly.

**Correctness:** In the real world, it is often acceptable for programs to behave correctly most of the time, on all 'reasonable' inputs. Not in this class; our standards are higher<sup>8</sup>. We need to *prove* that our algorithms are correct on *all possible* inputs. Sometimes this is fairly obvious, especially for algorithms you've seen in earlier courses. But many of the algorithms we will discuss in this course will require some extra work to prove. Correctness proofs almost always involve induction. We *like* induction. Induction is our *friend*.<sup>9</sup>

Before we can formally prove that our algorithm does what we want it to, we have to formally state what we want the algorithm to do! Usually problems are given to us in real-world terms,

---

<sup>7</sup>Similarly, the appearance of the phrase 'and so on' in a proof is a good indication that the proof *should* have been done by induction!

<sup>8</sup>or at least different

<sup>9</sup>If induction is *not* your friend, you will have a hard time in this course.

not with formal mathematical descriptions. It's up to us, the algorithm designer, to restate the problem in terms of mathematical objects that we can prove things about: numbers, arrays, lists, graphs, trees, and so on. We also need to determine if the problem statement makes any hidden assumptions, and state those assumptions explicitly. (For example, in the song “ $n$  Bottles of Beer on the Wall”,  $n$  is always a positive integer.) Restating the problem formally is not only required for proofs; it is also one of the best ways to really understand what the problem is asking for. The hardest part of solving a problem is figuring out the right way to ask the question!

An important distinction to keep in mind is the distinction between a problem and an algorithm. A problem is a task to perform, like “Compute the square root of  $x$ ” or “Sort these  $n$  numbers” or “Keep  $n$  algorithms students awake for  $t$  minutes”. An algorithm is a set of instructions that you follow if you want to execute this task. The same problem may have hundreds of different algorithms.

**Running time:** The usual way of distinguishing between different algorithms for the same problem is by how fast they run. Ideally, we want the fastest possible algorithm for our problem. In the real world, it is often acceptable for programs to run efficiently most of the time, on all ‘reasonable’ inputs. Not in this class; our standards are different. We require algorithms that *always* run efficiently, even in the worst case.

But how do we measure running time? As a specific example, how long does it take to sing the song BOTTLESOFBEER( $n$ )? This is obviously a function of the input value  $n$ , but it also depends on how quickly you can sing. Some singers might take ten seconds to sing a verse; others might take twenty. Technology widens the possibilities even further. Dictating the song over a telegraph using Morse code might take a full minute per verse. Ripping an mp3 over the Web might take a tenth of a second per verse. Duplicating the mp3 in a computer’s main memory might take only a few microseconds per verse.

Nevertheless, what’s important here is how the singing time changes as  $n$  grows. Singing BOTTLESOFBEER( $2n$ ) takes about twice as long as singing BOTTLESOFBEER( $n$ ), no matter what technology is being used. This is reflected in the asymptotic singing time  $\Theta(n)$ . We can measure time by counting how many times the algorithm executes a certain instruction or reaches a certain milestone in the ‘code’. For example, we might notice that the word ‘beer’ is sung three times in every verse of BOTTLESOFBEER, so the number of times you sing ‘beer’ is a good indication of the total singing time. For this question, we can give an exact answer: BOTTLESOFBEER( $n$ ) uses exactly  $3n + 3$  beers.

There are plenty of other songs that have non-trivial singing time. This one is probably familiar to most English-speakers:

```

NDAYSOFCHRISTMAS(gifts[2..n]):
  for i ← 1 to n
    Sing “On the ith day of Christmas, my true love gave to me”
    for j ← i down to 2
      Sing “j gifts[j”
    if i > 1
      Sing “and”
    Sing “a partridge in a pear tree.”

```

The input to NDAYSOFCHRISTMAS is a list of  $n - 1$  gifts. It’s quite easy to show that the singing time is  $\Theta(n^2)$ ; in particular, the singer mentions the name of a gift  $\sum_{i=1}^n i = n(n + 1)/2$  times (counting the partridge in the pear tree). It’s also easy to see that during the first  $n$  days of Christmas, my true love gave to me exactly  $\sum_{i=1}^n \sum_{j=1}^i j = n(n + 1)(n + 2)/6 = \Theta(n^3)$  gifts.



Other songs that take quadratic time to sing are “Old MacDonald”, “There Was an Old Lady Who Swallowed a Fly”, “Green Grow the Rushes O”, “The Barley Mow” (which we’ll see in Homework 1), “Echad Mi Yode’a” (“Who knows one?”), “Allouette”, “Ist das nicht ein Schnitzelbank?”<sup>10</sup> etc. For details, consult your nearest pre-schooler.

For a slightly more complicated example, consider the algorithm APPORTIONCONGRESS. Here the running time obviously depends on the implementation of the max-heap operations, but we can certainly bound the running time as  $O(N + RI + (R - n)E)$ , where  $N$  is the time for a NEWMAXHEAP,  $I$  is the time for an INSERT, and  $E$  is the time for an EXTRACTMAX. Under the reasonable assumption that  $R > 2n$  (on average, each state gets at least two representatives), this simplifies to  $O(N + R(I + E))$ . The Census Bureau uses an unsorted array of size  $n$ , for which  $N = I = \Theta(1)$  (since we know a priori how big the array is), and  $E = \Theta(n)$ , so the overall running time is  $\Theta(Rn)$ . This is fine for the federal government, but if we want to be more efficient, we can implement the heap as a perfectly balanced  $n$ -node binary tree (or a heap-ordered array). In this case, we have  $N = \Theta(1)$  and  $I = R = O(\log n)$ , so the overall running time is  $\Theta(R \log n)$ .

Incidentally, there is a faster algorithm for apportioning Congress. I’ll give extra credit to the first student who can find the faster algorithm, analyze its running time, and prove that it always gives exactly the same results as APPORTIONCONGRESS.

Sometimes we are also interested in other computational resources: space, disk swaps, concurrency, and so forth. We use the same techniques to analyze those resources as we use for running time.

#### 0.4 Why are we here, anyway?

We will try to teach you two things in CS373: how to *think* about algorithms and how to *talk* about algorithms. Along the way, you’ll pick up a bunch of algorithmic facts—mergesort runs in  $\Theta(n \log n)$  time; the amortized time to search in a splay tree is  $O(\log n)$ ; the traveling salesman problem is NP-hard—but these aren’t the point of the course. You can always look up facts in a textbook, provided you have the intuition to know what to look for. That’s why we let you bring cheat sheets to the exams; we don’t want you wasting your study time trying to memorize all the facts you’ve seen. You’ll also practice a lot of algorithm design and analysis skills—finding useful (counter)examples, developing induction proofs, solving recurrences, using big-Oh notation, using probability, and so on. These skills are useful, but they aren’t really the point of the course either. At this point in your educational career, you should be able to pick up those skills on your own, once you know what you’re trying to do.

The main goal of this course is to help you develop algorithmic *intuition*. How do various algorithms really work? When you see a problem for the first time, how should you attack it? How do you tell which techniques will work at all, and which ones will work best? How do you judge whether one algorithm is better than another? How do you tell if you have the best possible solution?

The flip side of this goal is developing algorithmic *language*. It’s not enough just to understand how to solve a problem; you also have to be able to explain your solution to somebody else. I don’t mean just how to turn your algorithms into code—despite what many students (and inexperienced programmers) think, ‘somebody else’ is *not* just a computer. Nobody programs alone. Perhaps more importantly in the short term, explaining something to somebody else is one of the best ways of clarifying your own understanding.

---

<sup>10</sup>Wakko: Ist das nicht Otto von Schnitzelpusskrankengescheitmeyer?

Yakko and Dot: Ja, das ist Otto von Schnitzelpusskrankengescheitmeyer!!

You'll also get a chance to develop brand new algorithms and algorithmic techniques on your own. Unfortunately, this is not the sort of thing that we can really teach you. All we can really do is lay out the tools, encourage you to practice with them, and give you feedback.

Good algorithms are extremely useful, but they can also be elegant, surprising, deep, even beautiful. But most importantly, algorithms are *fun*!! Hopefully this class will inspire at least a few of you to come play!

*The control of a large force is the same principle as the control of a few men:  
it is merely a question of dividing up their numbers.*

— Sun Zi, *The Art of War* (c. 400 C.E.), translated by Lionel Giles (1910)

## 1 Divide and Conquer (September 3)

### 1.1 MergeSort

Mergesort is one of the earliest algorithms proposed for sorting. According to Knuth, it was suggested by John von Neumann as early as 1945.

1. Divide the array  $A[1..n]$  into two subarrays  $A[1..m]$  and  $A[m+1..n]$ , where  $m = \lfloor n/2 \rfloor$ .
2. Recursively mergesort the subarrays  $A[1..m]$  and  $A[m+1..n]$ .
3. Merge the newly-sorted subarrays  $A[1..m]$  and  $A[m+1..n]$  into a single sorted list.

Input:	S	O	R	T	I	N	G	E	X	A	M	P	L	
Divide:	S	O	R	T	I	N		G	E	X	A	M	P	L
Recurse:	I	N	O	S	R	T		A	E	G	L	M	P	X
Merge:	A	E	G	I	L	M	N	O	P	S	R	T	X	

A Mergesort example.

The first step is completely trivial; we only need to compute the median index  $m$ . The second step is also trivial, thanks to our friend the recursion fairy. All the real work is done in the final step; the two sorted subarrays  $A[1..m]$  and  $A[m+1..n]$  can be merged using a simple linear-time algorithm. Here's a complete specification of the Mergesort algorithm; for simplicity, we separate out the merge step as a subroutine.

```

MERGESORT( $A[1..n]$ ):
  if ( $n > 1$ )
     $m \leftarrow \lfloor n/2 \rfloor$ 
    MERGESORT( $A[1..m]$ )
    MERGESORT( $A[m+1..n]$ )
    MERGE( $A[1..n], m$ )

```

```

MERGE( $A[1..n], m$ ):
   $i \leftarrow 1$ ;  $j \leftarrow m + 1$ 
  for  $k \leftarrow 1$  to  $n$ 
    if  $j > n$ 
       $B[k] \leftarrow A[i]$ ;  $i \leftarrow i + 1$ 
    else if  $i > m$ 
       $B[k] \leftarrow A[j]$ ;  $j \leftarrow j + 1$ 
    else if  $A[i] < A[j]$ 
       $B[k] \leftarrow A[i]$ ;  $i \leftarrow i + 1$ 
    else
       $B[k] \leftarrow A[j]$ ;  $j \leftarrow j + 1$ 
  for  $k \leftarrow 1$  to  $n$ 
     $A[k] \leftarrow B[k]$ 

```

To prove that the algorithm is correct, we use our old friend induction. We can prove that MERGE is correct using induction on the total size of the two subarrays  $A[i..m]$  and  $A[j..n]$  left to be merged into  $B[k..n]$ . The base case, where at least one subarray is empty, is straightforward; the algorithm just copies it into  $B$ . Otherwise, the smallest remaining element is either  $A[i]$  or  $A[j]$ , since both subarrays are sorted, so  $B[k]$  is assigned correctly. The remaining subarrays—either

$A[i + 1 .. m]$  and  $A[j .. n]$ , or  $A[i .. m]$  and  $A[j + 1 .. n]$ —are merged correctly into  $B[k + 1 .. n]$  by the inductive hypothesis.<sup>1</sup> This completes the proof.

Now we can prove MERGESORT correct by another round of straightforward induction.<sup>2</sup> The base cases  $n \leq 1$  are trivial. Otherwise, by the inductive hypothesis, the two smaller subarrays  $A[1 .. m]$  and  $A[m + 1 .. n]$  are sorted correctly, and by our earlier argument, merged into the correct sorted output.

What’s the running time? Since we have a recursive algorithm, we’re going to get a recurrence of some sort. MERGE clearly takes linear time, since it’s a simple for-loop with constant work per iteration. We get the following recurrence for MERGESORT:

$$T(1) = O(1), \quad T(n) = T(\lceil n/2 \rceil) + T(\lfloor n/2 \rfloor) + O(n).$$

## 1.2 Aside: Domain Transformations

Except for the floor and ceiling, this recurrence falls into case (b) of the Master Theorem [CLR, §4.3]. If we simply ignore the floor and ceiling, the Master Theorem suggests the solution  $T(n) = O(n \log n)$ . We can easily check that this answer is correct using induction, but there is a simple method for solving recurrences like this directly, called *domain transformation*.

First we overestimate the time bound, once by pretending that the two subproblem sizes are equal, and again to eliminate the ceiling:

$$T(n) \leq 2T(\lceil n/2 \rceil) + O(n) \leq 2T(n/2 + 1) + O(n).$$

Now we define a new function  $S(n) = T(n + \alpha)$ , where  $\alpha$  is a constant chosen so that  $S(n)$  satisfies the Master-ready recurrence  $S(n) \leq 2S(n/2) + O(n)$ . To figure out the appropriate value for  $\alpha$ , we compare two versions of the recurrence for  $T(n + \alpha)$ :

$$\begin{aligned} S(n) \leq 2S(n/2) + O(n) &\implies T(n + \alpha) \leq 2T(n/2 + \alpha) + O(n) \\ T(n) \leq 2T(n/2 + 1) + O(n) &\implies T(n + \alpha) \leq 2T((n + \alpha)/2 + 1) + O(n + \alpha) \end{aligned}$$

For these two recurrences to be equal, we need  $n/2 + \alpha = (n + \alpha)/2 + 1$ , which implies that  $\alpha = 2$ . The Master Theorem tells us that  $S(n) = O(n \log n)$ , so

$$T(n) = S(n - 2) = O((n - 2) \log(n - 2)) = O(n \log n).$$

We can use domain transformations to remove floors, ceilings, and lower order terms from any recurrence. But now that we know this, we won’t bother actually grinding through the details!

## 1.3 QuickSort

Quicksort was discovered by Tony Hoare in 1962. In this algorithm, the hard work is splitting the array into subsets so that merging the final result is trivial.

1. Choose a *pivot* element from the array.

<sup>1</sup>“The inductive hypothesis” is just a technical nickname for our friend the recursion fairy.

<sup>2</sup>Many textbooks draw an artificial distinction between several different flavors of induction: standard/weak (‘the principle of mathematical induction’), strong (‘the second principle of mathematical induction’), complex, structural, transfinite, decaffeinated, etc. Those textbooks would call this proof “strong” induction. I don’t. *All* induction proofs have precisely the same structure: Pick an arbitrary object, make one or more simpler objects from it, apply the inductive hypothesis to the simpler object(s), infer the required property for the original object, and check the base cases. Induction is just recursion for proofs.

2. Split the array into three subarrays containing the items less than the pivot, the pivot itself, and the items bigger than the pivot.
3. Recursively quicksort the first and last subarray.

Input:	S	O	R	T	I	N	G	E	X	A	M	P	L
Choose a pivot:	S	O	R	T	I	N	G	E	X	A	M	P	L
Partition:	M	A	E	G	I	L	N	R	X	O	S	P	T
Recurse:	A	E	G	I	L	M	N	O	P	S	R	T	X

A Quicksort example.

Here's a more formal specification of the Quicksort algorithm. The separate PARTITION subroutine takes the original position of the pivot element as input and returns the post-partition pivot position as output.

```

QUICKSORT( $A[1..n]$ ):
  if ( $n > 1$ )
    Choose a pivot element  $A[p]$ 
     $k \leftarrow$  PARTITION( $A, p$ )
    QUICKSORT( $A[1..k-1]$ )
    QUICKSORT( $A[k+1..n]$ )

```

```

PARTITION( $A[1..n], p$ ):
  if ( $p \neq n$ )
    swap  $A[p] \leftrightarrow A[n]$ 
   $i \leftarrow 0$ ;  $j \leftarrow n$ 
  while ( $i < j$ )
    repeat  $i \leftarrow i + 1$  until ( $i = j$  or  $A[i] \geq A[n]$ )
    repeat  $j \leftarrow j - 1$  until ( $i = j$  or  $A[j] \leq A[n]$ )
    if ( $i < j$ )
      swap  $A[i] \leftrightarrow A[j]$ 
  if ( $i \neq n$ )
    swap  $A[i] \leftrightarrow A[n]$ 
  return  $i$ 

```

Just as we did for mergesort, we need two induction proofs to show that QUICKSORT is correct—weak induction to prove that PARTITION correctly partitions the array, and then straightforward strong induction to prove that QUICKSORT correctly sorts assuming PARTITION is correct. I'll leave the gory details as an exercise for the reader.

The analysis is also similar to mergesort. PARTITION runs in  $O(n)$  time:  $j - i = n$  at the beginning,  $j - i = 0$  at the end, and we do a constant amount of work each time we increment  $i$  or decrement  $j$ . For QUICKSORT, we get a recurrence that depends on  $k$ , the rank of the chosen pivot:

$$T(n) = T(k-1) + T(n-k) + O(n)$$

If we could choose the pivot to be the median element of the array  $A$ , we would have  $k = \lceil n/2 \rceil$ , the two subproblems would be as close to the same size as possible, the recurrence would become

$$T(n) = 2T(\lceil n/2 \rceil - 1) + T(\lfloor n/2 \rfloor) + O(n) \leq 2T(n/2) + O(n),$$

and we'd have  $T(n) = O(n \log n)$  by the Master Theorem.

Unfortunately, although it is *theoretically* possible to locate the median of an unsorted array in linear time, the algorithm is incredibly complicated, and the hidden constant in the  $O()$  notation is quite large. So in practice, programmers settle for something simple, like choosing the first or last element of the array. In this case,  $k$  can be anything from 1 to  $n$ , so we have

$$T(n) = \max_{1 \leq k \leq n} (T(k-1) + T(n-k) + O(n))$$

In the worst case, the two subproblems are completely unbalanced—either  $k = 1$  or  $k = n$ —and the recurrence becomes  $T(n) \leq T(n - 1) + O(n)$ . The solution is  $T(n) = O(n^2)$ . Another common heuristic is ‘median of three’—choose three elements (usually at the beginning, middle, and end of the array), and take the middle one as the pivot. Although this is better in practice than just choosing one element, we can still have  $k = 2$  or  $k = n - 1$  in the worst case. With the median-of-three heuristic, the recurrence becomes  $T(n) \leq T(1) + T(n - 2) + O(n)$ , whose solution is still  $T(n) = O(n^2)$ .

Intuitively, the pivot element will ‘usually’ fall somewhere in the middle of the array, say between  $n/10$  and  $9n/10$ . This suggests that the *average-case* running time is  $O(n \log n)$ . Although this intuition is correct, we are still far from a *proof* that quicksort is usually efficient. I’ll formalize this intuition about average cases in a later lecture.

## 1.4 The Pattern

Both mergesort and quicksort follow the same general three-step pattern of all divide and conquer algorithms:

1. **Split** the problem into several *smaller independent* subproblems.
2. **Recurse** to get a subsolution for each subproblem.
3. **Merge** the subsolutions together into the final solution.

If the size of any subproblem falls below some constant threshold, the recursion bottoms out. Hopefully, at that point, the problem is trivial, but if not, we switch to a different algorithm instead.

Proving a divide-and-conquer algorithm correct usually involves strong induction. Analyzing the running time requires setting up and solving a recurrence, which often (but unfortunately not always!) can be solved using the Master Theorem, perhaps after a simple domain transformation.

## 1.5 Multiplication

Adding two  $n$ -digit numbers takes  $O(n)$  time by the standard iterative ‘ripple-carry’ algorithm, using a lookup table for each one-digit addition. Similarly, multiplying an  $n$ -digit number by a one-digit number takes  $O(n)$  time, using essentially the same algorithm.

What about multiplying two  $n$ -digit numbers? At least in the United States, every grade school student (supposedly) learns to multiply by breaking the problem into  $n$  one-digit multiplications and  $n$  additions:

$$\begin{array}{r}
 31415962 \\
 \times 27182818 \\
 \hline
 251327696 \\
 31415962 \\
 251327696 \\
 62831924 \\
 251327696 \\
 31415962 \\
 219911734 \\
 \hline
 62831924 \\
 \hline
 853974377340916
 \end{array}$$

We could easily formalize this algorithm as a pair of nested for-loops. The algorithm runs in  $O(n^2)$  time—altogether, there are  $O(n^2)$  digits in the partial products, and for each digit, we spend constant time.

We can do better by exploiting the following algebraic formula:

$$(10^m a + b)(10^m c + d) = 10^{2m} ac + 10^m (bc + ad) + bd$$

Here is a divide-and-conquer algorithm that computes the product of two  $n$ -digit numbers  $x$  and  $y$ , based on this formula. Each of the four sub-products  $e, f, g, h$  is computed recursively. The last line does not involve any multiplications, however; to multiply by a power of ten, we just shift the digits and fill in the right number of zeros.

<pre> MULTIPLY(<math>x, y, n</math>):   if <math>n = 1</math>     return <math>x \cdot y</math>   else     <math>m \leftarrow \lceil n/2 \rceil</math>     <math>a \leftarrow \lfloor x/10^m \rfloor</math>; <math>b \leftarrow x \bmod 10^m</math>     <math>d \leftarrow \lfloor y/10^m \rfloor</math>; <math>c \leftarrow y \bmod 10^m</math>     <math>e \leftarrow \text{MULTIPLY}(a, c, m)</math>     <math>f \leftarrow \text{MULTIPLY}(b, d, m)</math>     <math>g \leftarrow \text{MULTIPLY}(b, c, m)</math>     <math>h \leftarrow \text{MULTIPLY}(a, d, m)</math>     return <math>10^{2m}e + 10^m(g + h) + f</math> </pre>
--

You can easily prove by induction that this algorithm is correct. The running time for this algorithm is given by the recurrence

$$T(n) = 4T(\lceil n/2 \rceil) + O(n), \quad T(1) = 1,$$

which solves to  $T(n) = O(n^2)$  by the Master Theorem (after a simple domain transformation). Hmm... I guess this didn't help after all.

But there's a trick, first suggested by Anatoliĭ Karatsuba in 1962. We can compute the middle coefficient  $bc + ad$  using only *one* recursive multiplication, by exploiting yet another bit of algebra:

$$ac + bd - (a - b)(c - d) = bc + ad$$

This trick lets us replace the last three lines in the previous algorithm as follows:

<pre> FASTMULTIPLY(<math>x, y, n</math>):   if <math>n = 1</math>     return <math>x \cdot y</math>   else     <math>m \leftarrow \lceil n/2 \rceil</math>     <math>a \leftarrow \lfloor x/10^m \rfloor</math>; <math>b \leftarrow x \bmod 10^m</math>     <math>d \leftarrow \lfloor y/10^m \rfloor</math>; <math>c \leftarrow y \bmod 10^m</math>     <math>e \leftarrow \text{FASTMULTIPLY}(a, c, m)</math>     <math>f \leftarrow \text{FASTMULTIPLY}(b, d, m)</math>     <math>g \leftarrow \text{FASTMULTIPLY}(a - b, c - d, m)</math>     return <math>10^{2m}e + 10^m(e + f - g) + f</math> </pre>
---

The running time of Karatsuba's FASTMULTIPLY algorithm is given by the recurrence

$$T(n) \leq 3T(\lceil n/2 \rceil) + O(n), \quad T(1) = 1.$$

After a domain transformation, we can plug this into the Master Theorem to get the solution  $T(n) = O(n^{\lg 3}) = O(n^{1.585})$ , a significant improvement over our earlier quadratic-time algorithm.<sup>3</sup>

Of course, in practice, all this is done in binary instead of decimal.

We can take this idea even further, splitting the numbers into more pieces and combining them in more complicated ways, to get even faster multiplication algorithms. Ultimately, this idea leads to the development of the *Fast Fourier transform*, a complicated divide-and-conquer algorithm that can be used to multiply two  $n$ -digit numbers in  $O(n \log n)$  time.<sup>4</sup> We'll talk about Fast Fourier transforms later in the semester.

## 1.6 Exponentiation

Given a number  $a$  and a positive integer  $n$ , suppose we want to compute  $a^n$ . The standard naïve method is a simple for-loop that does  $n - 1$  multiplications by  $a$ :

<p style="margin: 0;">SLOWPOWER(<math>a, n</math>):</p> <p style="margin: 0; padding-left: 20px;"><math>x \leftarrow a</math></p> <p style="margin: 0; padding-left: 20px;">for <math>i \leftarrow 2</math> to <math>n</math></p> <p style="margin: 0; padding-left: 40px;"><math>x \leftarrow x \cdot a</math></p> <p style="margin: 0; padding-left: 20px;">return <math>x</math></p>
---

This iterative algorithm requires  $n$  multiplications.

Notice that the input  $a$  could be an integer, or a rational, or a floating point number. In fact, it doesn't need to be a number at all, as long as it's something that we know how to multiply. For example, the same algorithm can be used to compute powers modulo some finite number (an operation commonly used in cryptography algorithms) or to compute powers of matrices (an operation used to evaluate recurrences and to compute shortest paths in graphs). All that's required is that  $a$  belong to a multiplicative group.<sup>5</sup> Since we don't know what kind of things we're multiplying, we can't know how long a multiplication takes, so we're forced to analyze the running time in terms of the number of multiplications.

There is a much faster divide-and-conquer method, using the simple formula  $a^n = a^{\lfloor n/2 \rfloor} \cdot a^{\lceil n/2 \rceil}$ . What makes this approach more efficient is that once we compute the first factor  $a^{\lfloor n/2 \rfloor}$ , we can compute the second factor  $a^{\lceil n/2 \rceil}$  using at most one more multiplication.

<sup>3</sup>Karatsuba actually proposed an algorithm based on the formula  $(a+c)(b+d) - ac - bd = bc + ad$ . This algorithm also runs in  $O(n^{\lg 3})$  time, but the actual recurrence is a bit messier:  $a - b$  and  $c - d$  are still  $m$ -digit numbers, but  $a + b$  and  $c + d$  might have  $m + 1$  digits. The simplification presented here is due to Donald Knuth.

<sup>4</sup>This fast algorithm for multiplying integers using FFTs was discovered by Arnold Schönhage and Volker Strassen in 1971.

<sup>5</sup>A *multiplicative group*  $(G, \otimes)$  is a set  $G$  and a function  $\otimes : G \times G \rightarrow G$ , satisfying three axioms:

1. There is a unit element  $1 \in G$  such that  $1 \otimes g = g \otimes 1$  for any element  $g \in G$ .
2. Any element  $g \in G$  has a *inverse* element  $g^{-1} \in G$  such that  $g \otimes g^{-1} = g^{-1} \otimes g = 1$ .
3. The function is associative: for any elements  $f, g, h \in G$ , we have  $f \otimes (g \otimes h) = (f \otimes g) \otimes h$ .



```
FASTPOWER( $a, n$ ):  
  if  $n = 1$   
    return  $a$   
  else  
     $x \leftarrow$  FASTPOWER( $a, \lfloor n/2 \rfloor$ )  
    if  $n$  is even  
      return  $x \cdot x$   
    else  
      return  $x \cdot x \cdot a$ 
```

The total number of multiplications is given by the recurrence  $T(n) \leq T(\lfloor n/2 \rfloor) + 2$ , with the base case  $T(1) = 0$ . After a domain transformation, the Master Theorem gives us the solution  $T(n) = O(\log n)$ .

Incidentally, this algorithm is asymptotically optimal—any algorithm for computing  $a^n$  must perform  $\Omega(\log n)$  multiplications. In fact, when  $n$  is a power of two, this algorithm is *exactly* optimal. However, there are slightly faster methods for other values of  $n$ . For example, our divide-and-conquer algorithm computes  $a^{15}$  in six multiplications ( $a^{15} = a^7 \cdot a^7 \cdot a$ ;  $a^7 = a^3 \cdot a^3 \cdot a$ ;  $a^3 = a \cdot a \cdot a$ ), but only five multiplications are necessary ( $a \rightarrow a^2 \rightarrow a^3 \rightarrow a^5 \rightarrow a^{10} \rightarrow a^{15}$ ). Nobody knows of an algorithm that always uses the minimum possible number of multiplications.

*Those who cannot remember the past are doomed to repeat it.*

— George Santayana, *The Life of Reason, Book I: Introduction and Reason in Common Sense* (1905)

*The “Dynamic-Tension<sup>®</sup>” bodybuilding program takes only 15 minutes a day in the privacy of your room.*

— Charles Atlas

## 2 Dynamic Programming (September 5 and 10)

### 2.1 Exponentiation (Again)

Last time we saw a “divide and conquer” algorithm for computing the expression  $a^n$ , given two integers  $a$  and  $n$  as input: first compute  $a^{\lfloor n/2 \rfloor}$ , then  $a^{\lceil n/2 \rceil}$ , then multiply. If we computed both factors  $a^{\lfloor n/2 \rfloor}$  and  $a^{\lceil n/2 \rceil}$  recursively, the number of multiplications would be given by the recurrence

$$T(1) = 0, \quad T(n) = T(\lfloor n/2 \rfloor) + T(\lceil n/2 \rceil) + 1.$$

The solution is  $T(n) = n - 1$ , so the naïve recursive algorithm uses *exactly* the same number of multiplications as the naïve iterative algorithm.<sup>1</sup>

In this case, it’s obvious how to speed up the algorithm. Once we’ve computed  $a^{\lfloor n/2 \rfloor}$ , we don’t need to start over from scratch to compute  $a^{\lceil n/2 \rceil}$ ; we can do it in at most one more multiplication. This same simple idea—**don’t solve the same subproblem more than once**—can be applied to lots of recursive algorithms to speed them up, often (as in this case) by an exponential amount. The technical name for this technique is *dynamic programming*.

### 2.2 Fibonacci Numbers

The Fibonacci numbers  $F_n$ , named after Leonardo Fibonacci Pisano<sup>2</sup>, are defined as follows:  $F_0 = 0$ ,  $F_1 = 1$ , and  $F_n = F_{n-1} + F_{n-2}$  for all  $n \geq 2$ . The recursive definition of Fibonacci numbers immediately gives us a recursive algorithm for computing them:

<pre> RECFIBO(n):   if (n &lt; 2)     return n   else     return RECFIBO(n - 1) + RECFIBO(n - 2) </pre>
---

How long does this algorithm take? Except for the recursive calls, the entire algorithm requires only a constant number of steps: one comparison and possibly one addition. If  $T(n)$  represents the number of recursive calls to RECFIBO, we have the recurrence

$$T(0) = 1, \quad T(1) = 1, \quad T(n) = T(n - 1) + T(n - 2) + 1.$$

This looks an awful lot like the recurrence for Fibonacci numbers! In fact, it’s fairly easy to show by induction that  $T(n) = 2F_{n+1} - 1$ . In other words, computing  $F_n$  using this algorithm takes more than twice as many steps as just counting to  $F_n$ !

<sup>1</sup>But less time. If we assume that multiplying two  $n$ -digit numbers takes  $O(n \log n)$  time, then the iterative algorithm takes  $O(n^2 \log n)$  time, but this recursive algorithm takes only  $O(n \log^2 n)$  time.

<sup>2</sup>Literally, “Leonardo, son of Bonacci, of Pisa”.

Another way to see this is that the RECFIBO is building a big binary tree of additions, with nothing but zeros and ones at the leaves. Since the eventual output is  $F_n$ , our algorithm must call RECFIBO(1) (which returns 1) exactly  $F_n$  times. A quick inductive argument implies that RECFIBO(0) is called exactly  $F_{n-1}$  times. Thus, the recursion tree has  $F_n + F_{n-1} = F_{n+1}$  leaves, and therefore, because it's a full binary tree, it must have  $2F_{n+1} - 1$  nodes. (See Homework Zero!)

### 2.3 Aside: The Annihilator Method

Just how slow is that? We can get a good asymptotic estimate for  $T(n)$  by applying the annihilator method, described in the 'solving recurrences' handout:

$$\begin{aligned} \langle T(n+2) \rangle &= \langle T(n+1) \rangle + \langle T(n) \rangle + \langle 1 \rangle \\ \langle T(n+2) - T(n+1) - T(n) \rangle &= \langle 1 \rangle \\ (E^2 - E - 1) \langle T(n) \rangle &= \langle 1 \rangle \\ (E^2 - E - 1)(E - 1) \langle T(n) \rangle &= \langle 0 \rangle \end{aligned}$$

The characteristic polynomial of this recurrence is  $(r^2 - r - 1)(r - 1)$ , which has three roots:  $\phi = \frac{1+\sqrt{5}}{2} \approx 1.618$ ,  $\hat{\phi} = \frac{1-\sqrt{5}}{2} \approx -0.618$ , and 1. Thus, the generic solution is

$$T(n) = \alpha\phi^n + \beta\hat{\phi}^n + \gamma.$$

Now we plug in a few base cases:

$$\begin{aligned} T(0) &= 1 = \alpha + \beta + \gamma \\ T(1) &= 1 = \alpha\phi + \beta\hat{\phi} + \gamma \\ T(2) &= 3 = \alpha\phi^2 + \beta\hat{\phi}^2 + \gamma \end{aligned}$$

Solving this system of linear equations gives us

$$\alpha = 1 + \frac{1}{\sqrt{5}}, \quad \beta = 1 - \frac{1}{\sqrt{5}}, \quad \gamma = -1,$$

so our final solution is

$$T(n) = \left(1 + \frac{1}{\sqrt{5}}\right) \phi^n + \left(1 - \frac{1}{\sqrt{5}}\right) \hat{\phi}^n - 1 = \Theta(\phi^n).$$

Actually, if we only want an asymptotic bound, we only need to show that  $\alpha \neq 0$ , which is much easier than solving the whole system of equations. Since  $\phi$  is the largest characteristic root with non-zero coefficient, we immediately have  $T(n) = \Theta(\phi^n)$ .

### 2.4 Memo(r)ization and Dynamic Programming

The obvious reason for the recursive algorithm's lack of speed is that it computes the same Fibonacci numbers over and over and over. A single call to RECURSIVEFIBO( $n$ ) results in one recursive call to RECURSIVEFIBO( $n-1$ ), two recursive calls to RECURSIVEFIBO( $n-2$ ), three recursive calls to RECURSIVEFIBO( $n-3$ ), five recursive calls to RECURSIVEFIBO( $n-4$ ), and in general,  $F_{k-1}$  recursive calls to RECURSIVEFIBO( $n-k$ ), for any  $0 \leq k < n$ . For each call, we're recomputing some Fibonacci number from scratch.

We can speed up the algorithm considerably just by writing down the results of our recursive calls and looking them up again if we need them later. This process is called *memoization*.<sup>3</sup>

<sup>3</sup>"My name is Elmer J. Fudd, millionaire. I own a mansion and a yacht."

```

MEMFIBO(n):
  if (n < 2)
    return n
  else
    if F[n] is undefined
      F[n] ← MEMFIBO(n - 1) + MEMFIBO(n - 2)
    return F[n]

```

If we actually trace through the recursive calls made by MEMFIBO, we find that the array  $F[]$  gets filled from the bottom up: first  $F[2]$ , then  $F[3]$ , and so on, up to  $F[n]$ . Once we realize this, we can replace the recursion with a simple for-loop that just fills up the array in that order, instead of relying on the complicated recursion to do it for us. This gives us our first explicit *dynamic programming* algorithm.

```

ITERFIBO(n):
  F[0] ← 0
  F[1] ← 1
  for i ← 2 to n
    F[i] ← F[i - 1] + F[i - 2]
  return F[n]

```

ITERFIBO clearly takes only  $O(n)$  time and  $O(n)$  space to compute  $F_n$ , an exponential speedup over our original recursive algorithm. We can reduce the space to  $O(1)$  by noticing that we never need more than the last two elements of the array:

```

ITERFIBO2(n):
  prev ← 1
  curr ← 0
  for i ← 1 to n
    next ← curr + prev
    prev ← curr
    curr ← next
  return curr

```

(This algorithm uses the non-standard but perfectly consistent base case  $F_{-1} = 1$ .)

But even this isn't the fastest algorithm for computing Fibonacci numbers. There's a faster algorithm defined in terms of matrix multiplication, using the following wonderful fact:

$$\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix} \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} y \\ x + y \end{bmatrix}$$

In other words, multiplying a two-dimensional vector by the matrix  $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$  does exactly the same thing as one iteration of the inner loop of ITERFIBO2. This might lead us to believe that multiplying by the matrix  $n$  times is the same as iterating the loop  $n$  times:

$$\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}^n \begin{bmatrix} 1 \\ 0 \end{bmatrix} = \begin{bmatrix} F_{n-1} \\ F_n \end{bmatrix}.$$

A quick inductive argument proves this. So if we want to compute the  $n$ th Fibonacci number, all we have to do is compute the  $n$ th power of the matrix  $\begin{bmatrix} 0 & 1 \\ 1 & 1 \end{bmatrix}$ .

We saw in the previous lecture, and the beginning of this lecture, that computing the  $n$ th power of something requires only  $O(\log n)$  multiplications. In this case, that means  $O(\log n)$   $2 \times 2$  matrix multiplications, but one matrix multiplication can be done with only a constant number of integer multiplications and additions. By applying our earlier dynamic programming algorithm for computing exponents, we can compute  $F_n$  in only  $O(\log n)$  steps.

This is an exponential speedup over the standard iterative algorithm, which was already an exponential speedup over our original recursive algorithm. Right?

## 2.5 Uh... wait a minute.

Well, not exactly. Fibonacci numbers grow exponentially fast. The  $n$ th Fibonacci number is approximately  $n \log_{10} \phi \approx n/5$  decimal digits long, or  $n \log_2 \phi \approx 2n/3$  bits. So we can't possibly compute  $F_n$  in logarithmic time — we need  $\Omega(n)$  time just to write down the answer!

I've been cheating by assuming we can do arbitrary-precision arithmetic in constant time. As we discussed last time, multiplying two  $n$ -digit numbers takes  $O(n \log n)$  time. That means that the matrix-based algorithm's actual running time is given by the recurrence

$$T(n) = T(\lfloor n/2 \rfloor) + O(n \log n),$$

which solves to  $T(n) = O(n \log n)$  by the Master Theorem.

Is this slower than our "linear-time" iterative algorithm? No! Addition isn't free, either. Adding two  $n$ -digit numbers takes  $O(n)$  time, so the running time of the iterative algorithm is  $O(n^2)$ . (Do you see why?) So our matrix algorithm really is faster than our iterative algorithm, but not exponentially faster.

Incidentally, in the recursive algorithm, the extra cost of arbitrary-precision arithmetic is overwhelmed by the huge number of recursive calls. The correct recurrence is

$$T(n) = T(n-1) + T(n-2) + O(n),$$

which still has the solution  $O(\phi^n)$  by the annihilator method.

## 2.6 The Pattern

Dynamic programming is essentially *recursion without repetition*. Developing a dynamic programming algorithm generally involves two separate steps.

1. **Formulate the problem recursively.** Write down a formula for the whole problem as a simple combination of the answers to smaller subproblems.
2. **Build solutions to your recurrence from the bottom up.** Write an algorithm that starts with the base cases of your recurrence and works its way up to the final solution by considering the intermediate subproblems in the correct order.

Of course, you have to prove that each of these steps is correct. If your recurrence is wrong, or if you try to build up answers in the wrong order, your algorithm won't work!

Dynamic programming algorithms need to store the results of intermediate subproblems. This is often *but not always* done with some kind of table.

## 2.7 Edit Distance

The *edit distance* between two words is the minimum number of letter insertions, letter deletions, and letter substitutions required to transform one word into another. For example, the edit distance between FOOD and MONEY is at most four:

$$\underline{\text{FOOD}} \rightarrow \text{MOOD} \rightarrow \text{MON}\underline{\text{D}} \rightarrow \text{MONED}\underline{\text{ }} \rightarrow \text{MONEY}$$

A better way to display this editing process is to place the words one above the other, with a gap in the first word for every insertion, and a gap in the second word for every deletion. Columns with two *different* characters correspond to substitutions. Thus, the number of editing steps is just the number of columns that don't contain the same character twice.

F	O	O	D	
M	O	N	E	Y

It's fairly obvious that you can't get from FOOD to MONEY in three steps, so their edit distance is exactly four. Unfortunately, this is not so easy in general. Here's a longer example, showing that the distance between ALGORITHM and ALTRUISTIC is at most six. Is this optimal?

A	L	G	O	R		I		T	H	M
A	L		T	R	U	I	S	T	I	C

To develop a dynamic programming algorithm to compute the edit distance between two strings, we first need to develop a recursive definition. Let's say we have an  $m$ -character string  $A$  and an  $n$ -character string  $B$ . Then define  $E(i, j)$  to be the edit distance between the first  $i$  characters of  $A$  and the first  $j$  characters of  $B$ . The edit distance between the entire strings  $A$  and  $B$  is  $E(m, n)$ .

This gap representation for edit sequences has a crucial "optimal substructure" property. Suppose we have the gap representation for the shortest edit sequence for two strings. **If we remove the last column, the remaining columns must represent the shortest edit sequence for the remaining substrings.** We can easily prove this by contradiction. If the substrings had a shorter edit sequence, we could just glue the last column back on and get a shorter edit sequence for the original strings.

There are a couple of obvious base cases. The only way to convert the empty string into a string of  $j$  characters is by doing  $j$  insertions, and the only way to convert a string of  $i$  characters into the empty string is with  $i$  deletions:

$$E(i, 0) = i, \quad E(0, j) = j.$$

In general, there are three possibilities for the last column in the shortest possible edit sequence:

- **Insertion:** The last entry in the bottom row is empty. In this case,  $E(i, j) = E(i - 1, j) + 1$ .
- **Deletion:** The last entry in the top row is empty. In this case,  $E(i, j) = E(i, j - 1) + 1$ .
- **Substitution:** Both rows have characters in the last column. If the characters are the same, we don't actually have to pay for the substitution, so  $E(i, j) = E(i - 1, j - 1)$ . If the characters are different, then  $E(i, j) = E(i - 1, j - 1) + 1$ .

To summarize, the edit distance  $E(i, j)$  is the smallest of these three possibilities:

$$E(i, j) = \min \left\{ \begin{array}{l} E(i - 1, j) + 1 \\ E(i, j - 1) + 1 \\ E(i - 1, j - 1) + [A[i] \neq B[j]] \end{array} \right\}$$

[The bracket notation  $[P]$  denotes the *indicator variable* for the logical proposition  $P$ . Its value is 1 if  $P$  is true and 0 if  $P$  is false. This is the same as the C/C++/Java expression  $P ? 1 : 0$ .]

If we turned this recurrence directly into a recursive algorithm, we would have the following horrible double recurrence for the running time:

$$T(m, 0) = T(0, n) = O(1), \quad T(m, n) = T(m, n - 1) + T(m - 1, n) + T(n - 1, m - 1) + O(1).$$

Yuck!! The solution for this recurrence is an exponential mess! I don't know a general closed form, but  $T(n, n) = \Theta((1 + \sqrt{2})^n)$ . Obviously a recursive algorithm is not the way to go here.

Instead, let's build an  $m \times n$  table of all possible values of  $E(i, j)$ . We can start by filling in the base cases, the entries in the 0th row and 0th column, each in constant time. To fill in any other entry, we need to know the values directly above it, directly to the left, and both above and to the left. If we fill in our table in the standard way—row by row from top down, each row from left to right—then whenever we reach an entry in the matrix, the entries it depends on are already available.

```

EDITDISTANCE( $A[1..m], B[1..n]$ ):
  for  $i \leftarrow 1$  to  $m$ 
    Edit[ $i, 0$ ]  $\leftarrow i$ 
  for  $j \leftarrow 1$  to  $n$ 
    Edit[ $0, j$ ]  $\leftarrow j$ 
  for  $i \leftarrow 1$  to  $m$ 
    for  $j \leftarrow 1$  to  $n$ 
      if  $A[i] = B[j]$ 
        Edit[ $i, j$ ]  $\leftarrow \min \{$  Edit[ $i - 1, j$ ] + 1,
                                Edit[ $i, j - 1$ ] + 1,
                                Edit[ $i - 1, j - 1$ ]  $\}$ 
      else
        Edit[ $i, j$ ]  $\leftarrow \min \{$  Edit[ $i - 1, j$ ] + 1,
                                Edit[ $i, j - 1$ ] + 1,
                                Edit[ $i - 1, j - 1$ ] + 1  $\}$ 
  return Edit[ $m, n$ ]

```

Since there are  $\Theta(n^2)$  entries in the table, and each entry takes  $\Theta(1)$  time once we know its predecessors, the total running time is  $\Theta(n^2)$ .

Here's the resulting table for **ALGORITHM**  $\rightarrow$  **ALTRUISTIC**. Bold numbers indicate places where characters in the two strings are equal. The arrows represent the predecessor(s) that actually define each entry. Each direction of arrow corresponds to a different edit operation: horizontal=deletion, vertical=insertion, and diagonal=substitution. Bold diagonal arrows indicate "free" substitutions of a letter for itself. A path of arrows from the top left corner to the bottom right corner of this table represents an optimal edit sequence between the two strings. There can be many such paths.

		A	L	G	O	R	I	T	H	M										
		0	→	1	→	2	→	3	→	4	→	5	→	6	→	7	→	8	→	9
A	↓	1	↘	0	→	1	→	2	→	3	→	4	→	5	→	6	→	7	→	8
L	↓	2	↓	1	↘	0	→	1	→	2	→	3	→	4	→	5	→	6	→	7
T	↓	3	↓	2	↓	1	↘	1	→	2	→	3	→	4	→	4	→	5	→	6
R	↓	4	↓	3	↓	2	↓	2	↘	2	→	2	↘	2	↘	2	↘	2	↘	2
U	↓	5	↓	4	↓	3	↓	3	↓	3	↓	3	↓	3	↘	3	↘	3	↘	3
I	↓	6	↓	5	↓	4	↓	4	↓	4	↓	4	↓	4	↘	4	↘	4	↘	4
S	↓	7	↓	6	↓	5	↓	5	↓	5	↓	5	↓	5	↘	5	↘	5	↘	5
T	↓	8	↓	7	↓	6	↓	6	↓	6	↓	6	↓	6	↘	6	↘	6	↘	6
I	↓	9	↓	8	↓	7	↓	7	↓	7	↓	7	↓	7	↘	7	↘	7	↘	7
C	↓	10	↓	9	↓	8	↓	8	↓	8	↓	8	↓	8	↘	8	↘	8	↘	8

The edit distance between ALGORITHM and ALTRUISTIC is indeed six. There are three paths through this table from the top left to the bottom right, so there are three optimal edit sequences:

```

A L G O R I T H M
A L T R U I S T I C

A L G O R I T H M
A L T R U I S T I C

A L G O R I T H M
A L T R U I S T I C
    
```

### 2.8 Danger! Greed kills!

If we're very very very lucky, we can bypass all the recurrences and tables and so forth, and solve the problem using a *greedy* algorithm. The general greedy strategy is look for the best first step, take it, and then continue. For example, a greedy algorithm for the edit distance problem might look for the longest common substring of the two strings, match up those substrings (since those substitutions dont cost anything), and then recursively look for the edit distances between the left halves and right halves of the strings. If there is no common substring—that is, if the two strings have no characters in common—the edit distance is clearly the length of the larger string.

If this sounds like a hack to you, pat yourself on the back. It isn't even close to the correct solution. Nevertheless, for many problems involving dynamic programming, many student's first intuition is to apply a greedy strategy. This almost never works; problems that can be solved correctly by a greedy algorithm are *very* rare. Everyone should tattoo the following sentence on the back of their hands, right under all the rules about logarithms and big-Oh notation:

**Greedy algorithms never work!**  
**Use dynamic programming instead!**

Well...hardly ever. Later in the semester we'll see correct greedy algorithms for minimum spanning trees and shortest paths.



## 2.9 Optimal Binary Search Trees

You all remember that the cost of a successful search in a binary search tree is proportional to the depth of the target node plus one. As a result, the worst-case search time is proportional to the height of the tree. To minimize the worst-case search time, the height of the tree should be as small as possible; ideally, the tree is perfectly balanced.

In many applications of binary search trees, it is more important to minimize the total cost of several searches than to minimize the worst-case cost of a single search. If  $x$  is a more ‘popular’ search target than  $y$ , we can save time by building a tree where the depth of  $x$  is smaller than the depth of  $y$ , even if that means increasing the overall height of the tree. A perfectly balanced tree is *not* the best choice if some items are significantly more popular than others. In fact, a totally unbalanced tree of depth  $\Omega(n)$  might actually be the best choice!

This situation suggests the following problem. Suppose we are given a sorted array of  $n$  keys  $A[1..n]$  and an array of corresponding *access frequencies*  $f[1..n]$ . Over the lifetime of the search tree, we will search for the key  $A[i]$  exactly  $f[i]$  times. Our task is to build the binary search tree that minimizes the *total* search time.

Before we think about how to solve this problem, we should first come up with the right way to describe the function we are trying to optimize! Suppose we have a binary search tree  $T$ . Let  $\text{depth}(T, i)$  denote the depth of the node in  $T$  that stores the key  $A[i]$ . Up to constant factors, the total search time  $S(T)$  is given by the following expression:

$$S(T) = \sum_{i=1}^n (\text{depth}(T, i) + 1) \cdot f[i]$$

This expression is called the *weighted external path length* of  $T$ . We can trivially split this expression into two components:

$$S(T) = \sum_{i=1}^n f[i] + \sum_{i=1}^n \text{depth}(T, i) \cdot f[i].$$

The first term is the total number of searches, which doesn’t depend on our choice of tree at all. The second term is called the *weighted internal path length* of  $T$ .

We can express  $S(T)$  in terms of the recursive structure of  $T$  as follows. Suppose the root of  $T$  contains the key  $A[r]$ , so the left subtree stores the keys  $A[1..r-1]$ , and the right subtree stores the keys  $A[r+1..n]$ . We can actually define  $\text{depth}(T, i)$  recursively as follows:

$$\text{depth}(T, i) = \begin{cases} \text{depth}(\text{left}(T), i) + 1 & \text{if } i < r \\ 0 & \text{if } i = r \\ \text{depth}(\text{right}(T), i) + 1 & \text{if } i > r \end{cases}$$

If we plug this recursive definition into our earlier expression for  $S(T)$ , we get the following:

$$S(T) = \sum_{i=1}^n f[i] + \sum_{i=1}^{r-1} (\text{depth}(\text{left}(T), i) + 1) \cdot f[i] + \sum_{i=1}^{r-1} (\text{depth}(\text{right}(T), i) + 1) \cdot f[i]$$

This looks complicated, until we realize that the second and third look exactly like our initial expression for  $S(T)$ !

$$S(T) = \sum_{i=1}^n f[i] + S(\text{left}(T)) + S(\text{right}(T))$$

Now our task is to compute the tree  $T_{\text{opt}}$  that minimizes the total search time  $S(T)$ . Suppose the root of  $T_{\text{opt}}$  stores key  $A[r]$ . The recursive definition of  $S(T)$  immediately implies that the left subtree  $\text{left}(T_{\text{opt}})$  must also be the optimal search tree for the keys  $A[1..r-1]$  and access frequencies  $f[1..r-1]$ . Similarly, the right subtree  $\text{right}(T_{\text{opt}})$  must also be the optimal search tree for the keys  $A[r+1..n]$  and access frequencies  $f[r+1..n]$ . Thus, once we choose the correct key to store at the root, the recursion fairy will automatically construct the rest of the optimal tree for us!

More formally, let  $S(i, j)$  denote the total search time for the *optimal* search tree containing the subarray  $A[1..j]$ ; our task is to compute  $S(1, n)$ . To simplify notation a bit, let  $F(i, j)$  denote the total frequency counts for all the keys in the subarray  $A[i..j]$ :

$$F(i, j) = \sum_{k=i}^j f[k]$$

We now have the following recurrence:

$$S(i, j) = \begin{cases} 0 & \text{if } j = i - 1 \\ F(i, j) + \min_{1 \leq r \leq n} (S(1, r - 1) + S(r + 1, n)) & \text{otherwise} \end{cases}$$

The base case might look a little weird, but all it means is that the total cost for searching an empty set of keys is zero. We could use the base cases  $S(i, i) = f[i]$  instead, but this would lead to extra special cases when  $r = 1$  or  $r = n$ . Also, the  $F(i, j)$  term is outside the max because it doesn't depend on the root index  $r$ .

Now, if we try to evaluate this recurrence directly using a recursive algorithm, the running time will have the following evil-looking recurrence:

$$T(n) = \Theta(n) + \sum_{k=1}^n (T(k - 1) + T(n - k))$$

The  $\Theta(n)$  term comes from computing  $F(1, n)$ , which is the total number of searches. A few minutes of pain and suffering by a professional algorithm analyst gives us the solution  $T(n) = \Theta(3^n)$ . Once again, top-down recursion is not the way to go.

In fact, we're not even computing the access counts  $F(i, j)$  as efficiently as we could. Even if we memoize the answers in an array  $F[1..n][1..n]$ , computing each value  $F(i, j)$  using a separate for-loop requires a total of  $O(n^3)$  time. A better approach is to turn the recurrence

$$F(i, j) = \begin{cases} f[i] & \text{if } i = j \\ F(i, j - 1) + f[j] & \text{otherwise} \end{cases}$$

into the following  $O(n^2)$ -time dynamic programming algorithm:

```

INITF( $f[1..n]$ ):
  for  $i \leftarrow 1$  to  $n$ 
     $F[i, i - 1] \leftarrow 0$ 
    for  $j \leftarrow i$  to  $n$ 
       $F[i, j] \leftarrow F[i, j - 1] + f[j]$ 

```

This will be used as an initialization subroutine in our final algorithm.

So now let's compute the optimal search tree cost  $S(1, n)$  from the bottom up. We can store all intermediate results in a table  $S[1..n, 0..n]$ . Only the entries  $S[i, j]$  with  $j \geq i - 1$  will actually be used. The base case of the recursion tells us that any entry of the form  $S[i, i - 1]$  can immediately be set to 0. For any other entry  $S[i, j]$ , we can use the following algorithm fragment, which comes directly from the recurrence:

```

COMPUTES( $i, j$ ):
   $S[i, j] \leftarrow \infty$ 
  for  $r \leftarrow i$  to  $j$ 
     $tmp \leftarrow S[i, r - 1] + S[r + 1, j]$ 
    if  $S[i, j] > tmp$ 
       $S[i, j] \leftarrow tmp$ 
   $S[i, j] \leftarrow S[i, j] + F[i, j]$ 

```

The only question left is what order to fill in the table.

Each entry  $S[i, j]$  depends on all entries  $S[i, r - 1]$  and  $S[r + 1, j]$  with  $i \leq k \leq j$ . In other words, every entry in the table depends on all the entries directly to the left or directly below. In order to fill the table efficiently, we must choose an order that computes all those entries before  $S[i, j]$ . There are at least three different orders that satisfy this constraint. The one that occurs to most people first is to scan through the table one diagonal at a time, starting with the trivial base cases  $S[i, i - 1]$ . The complete algorithm looks like this:

```

OPTIMALSEARCHTREE( $f[1..n]$ ):
  INITF( $f[1..n]$ )
  for  $i \leftarrow 1$  to  $n$ 
     $S[i, i - 1] \leftarrow 0$ 
  for  $d \leftarrow 0$  to  $n - 1$ 
    for  $i \leftarrow 1$  to  $n - d$ 
      COMPUTES( $i, i + d$ )
  return  $S[1, n]$ 

```

We could also traverse the array row by row from the bottom up, traversing each row from left to right, or column by column from left to right, traversing each columns from the bottom up. These two orders give us the following algorithms:

```

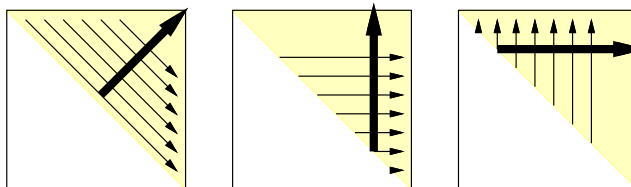
OPTIMALSEARCHTREE2( $f[1..n]$ ):
  INITF( $f[1..n]$ )
  for  $i \leftarrow n$  downto 1
     $S[i, i - 1] \leftarrow 0$ 
    for  $j \leftarrow i$  to  $n$ 
      COMPUTES( $i, j$ )
  return  $S[1, n]$ 

```

```

OPTIMALSEARCHTREE3( $f[1..n]$ ):
  INITF( $f[1..n]$ )
  for  $j \leftarrow 0$  to  $n$ 
     $S[j + 1, j] \leftarrow 0$ 
    for  $i \leftarrow j$  downto 1
      COMPUTES( $i, j$ )
  return  $S[1, n]$ 

```



Three different orders to fill in the table  $S[i, j]$ .

No matter which of these three orders we actually use, the resulting algorithm runs in  $\Theta(n^3)$  time and uses  $\Theta(n^2)$  space.

We could have predicted this from the original recursive formulation.

$$S(i, j) = \begin{cases} 0 & \text{if } j = i - i \\ F(i, j) + \min_{i \leq r \leq j} (S(i, r - 1) + S(r + 1, j)) & \text{otherwise} \end{cases}$$

First, the function has two arguments, each of which can take on any value between 1 and  $n$ , so we probably need a table of size  $O(n^2)$ . Next, there are *three* variables in the recurrence ( $i$ ,  $j$ , and  $r$ ), each of which can take any value between 1 and  $n$ , so it should take us  $O(n^3)$  time to fill the table.

In general, you can get an easy estimate of the time and space bounds for any dynamic programming algorithm by looking at the recurrence. The time bound is determined by how many values *all* the variables can have, and the space bound is determined by how many values the parameters of the function can have. For example, the (completely made up) recurrence

$$F(i, j, k, l, m) = \min_{0 \leq p \leq i} \max_{0 \leq q \leq j} \sum_{r=1}^{k-m} F(i-p, j-q, r, l-1, m-r)$$

should immediately suggest a dynamic programming algorithm that uses  $O(n^8)$  time and  $O(n^5)$  space. This rule of thumb immediately usually gives us the right time bound to shoot for.

But not always! In fact, the algorithm I've described is *not* the most efficient algorithm for computing optimal binary search trees. Let  $R[i, j]$  denote the root of the optimal search tree for  $A[i..j]$ . Donald Knuth proved the following nice monotonicity property for optimal subtrees: if we move either end of the subarray, the optimal root moves in the same direction or not at all, or more formally:

$$\boxed{R[i, j - 1] \leq R[i, j] \leq R[i + 1, j] \text{ for all } i \text{ and } j.}$$

This (nontrivial!) observation leads to the following more efficient algorithm:

<pre> <b>FASTEROPTIMALSEARCHTREE</b>(<math>f[1..n]</math>):   INITF(<math>f[1..n]</math>)   for <math>i \leftarrow n</math> downto 1     <math>S[i, i - 1] \leftarrow 0</math>     <math>R[i, i - 1] \leftarrow i</math>   for <math>j \leftarrow i</math> to <math>n</math>     COMPUTESANDR(<math>i, j</math>)   return <math>S[1, n]</math> </pre>	<pre> <b>COMPUTESANDR</b>(<math>f[1..n]</math>):   <math>S[i, j] \leftarrow \infty</math>   for <math>r \leftarrow R[i, j - 1]</math> to <math>j</math>     <math>tmp \leftarrow S[i, r - 1] + S[r + 1, j]</math>     if <math>S[i, j] &gt; tmp</math>       <math>S[i, j] \leftarrow tmp</math>       <math>R[i, j] \leftarrow r</math>   <math>S[i, j] \leftarrow S[i, j] + F[i, j]</math> </pre>
---	---

It's not hard to see the  $r$  increases monotonically from  $i$  to  $n$  during each iteration of the *outermost* for loop. Consequently, the innermost for loop iterates at most  $n$  times during a single iteration of the outermost loop, so the total running time of the algorithm is  $O(n^2)$ .

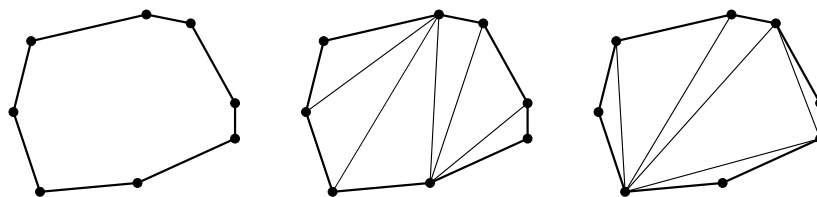
If we formulate the problem slightly differently, this algorithm can be improved even further. Suppose we require the optimum *external* binary tree, where the keys  $A[1..n]$  are all stored at the leaves, and intermediate pivot values are stored at the internal nodes. An algorithm due to Te Ching Hu and Alan Tucker<sup>4</sup> computes the optimal binary search tree in this setting in only

<sup>4</sup>T. C. Hu and A. C. Tucker, Optimal computer search trees and variable length alphabetic codes, *SIAM J. Applied Math.* 21:514-532, 1971. For a slightly simpler algorithm with the same running time, see A. M. Garsia and M. L. Wachs, A new algorithms for minimal binary search trees, *SIAM J. Comput.* 6:622-642, 1977. The original correctness proofs for both algorithms are rather intricate; for simpler proofs, see Marek Karpinski, Lawrence L. Larmore, and Wojciech Rytter, Correctness of constructing optimal alphabetic trees revisited, *Theoretical Computer Science*, 180:309-324, 1997.

$O(n \log n)$  time!

## 2.10 Optimal Triangulations of Convex Polygons

A *convex polygon* is a circular chain of line segments, arranged so none of the corners point inwards—imagine a rubber band stretched around a bunch of nails. (This is technically not the best definition, but it'll do for now.) A *diagonal* is a line segment that cuts across the interior of the polygon from one corner to another. A simple induction argument (hint, hint) implies that any  $n$ -sided convex polygon can be split into  $n-2$  triangles by cutting along  $n-3$  different diagonals. This collection of triangles is called a *triangulation* of the polygon. Triangulations are incredibly useful in computer graphics—most graphics hardware is built to draw triangles incredibly quickly, but to draw anything more complicated, you usually have to break it into triangles first.



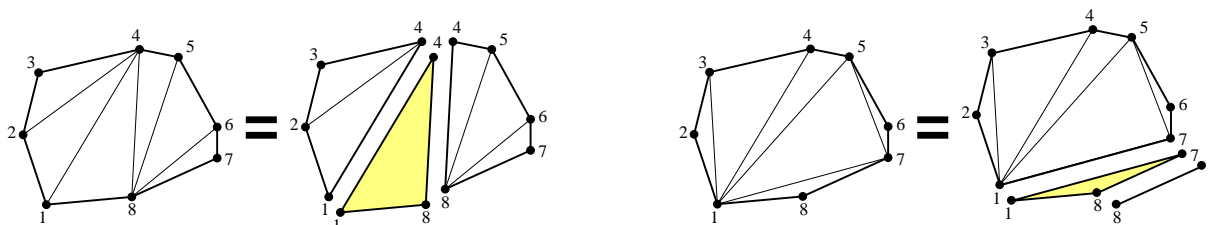
A convex polygon and two of its many possible triangulations.

There are several different ways to triangulate any convex polygon. Suppose we want to find the triangulation that requires the least amount of ink to draw, or in other words, the triangulation where the total perimeter of the triangles is as small as possible. To make things concrete, let's label the corners of the polygon from 1 to  $n$ , starting at the bottom of the polygon and going clockwise. We'll need the following subroutines to compute the perimeter of a triangle joining three corners using their  $x$ - and  $y$ -coordinates:

$$\Delta(i, j, k) : \\ \text{return DIST}(i, j) + \text{DIST}(j, k) + \text{DIST}(i, k)$$

$$\text{DIST}(i, j) : \\ \text{return } \sqrt{(x[i] - x[j])^2 + (y[i] - y[j])^2}$$

In order to get a dynamic programming algorithm, we first need a recursive formulation of the minimum-length triangulation. To do that, we really need some kind of recursive definition of a *triangulation*! Notice that in any triangulation, exactly one triangle uses both the first corner and the last corner of the polygon. If we remove that triangle, what's left over is two smaller triangulations. The base case of this recursive definition is a 'polygon' with just two corners. Notice that at any point in the recursion, we have a polygon joining a contiguous subset of the original corners.



Two examples of the recursive definition of a triangulation.

Building on this recursive definition, we can now recursively define the total length of the minimum-length triangulation. In the best triangulation, if we remove the 'base' triangle, what

remains must be the optimal triangulation of the two smaller polygons. So we just have choose the best triangle to attach to the first and last corners, and let the recursion fairy take care of the rest:

$$M(i, j) = \begin{cases} 0 & \text{if } j = i + 1 \\ \min_{i < k < j} (\Delta(i, j, k) + M(i, k) + M(k, j)) & \text{otherwise} \end{cases}$$

What we're looking for is  $M(1, n)$ .

If you think this looks similar to the recurrence for  $S(i, j)$ , the cost of an optimal binary search tree, you're absolutely right. We can build up intermediate results in a two-dimensional table, starting with the base cases  $M[i, i + 1] = 0$  and working our way up. We can use the following algorithm fragment to compute a generic entry  $M[i, j]$ :

```

COMPUTEM(i, j):
  M[i, j] ← ∞
  for k ← i + 1 to j - 1
    tmp ← Δ(i, j, k) + M[i, k] + M[k, j]
    if M[i, j] > tmp
      M[i, j] ← tmp

```

As in the optimal search tree problem, each table entry  $M[i, j]$  depends on all the entries directly to the left or directly below, so we can use any of the orders described earlier to fill the table.

```

MINTRIANGULATION:
  for i ← 1 to n - 1
    M[i, i + 1] ← 0
  for d ← 2 to n - 1
    for i ← 1 to n - d
      COMPUTEM(i, i + d)
  return M[1, n]

```

```

MINTRIANGULATION2:
  for i ← n downto 1
    M[i, i + 1] ← 0
  for j ← i + 2 to n
    COMPUTEM(i, j)
  return M[1, n]

```

```

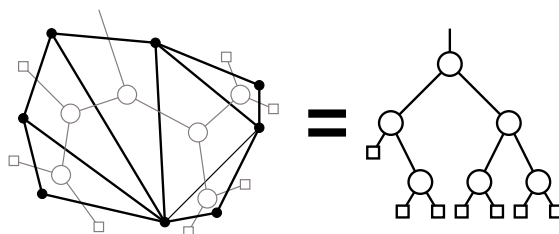
MINTRIANGULATION3:
  for j ← 2 to n
    M[j - 1, j] ← 0
  for i ← j - 1 downto 1
    COMPUTEM(i, j)
  return M[1, n]

```

In all three cases, the algorithm runs in  $\Theta(n^3)$  time and uses  $\Theta(n^2)$  space, just as we should have guessed from the recurrence.

## 2.11 It's the same problem!

Actually, the last two problems are both special cases of the same meta-problem: computing optimal *Catalan* structures. There is a straightforward one-to-one correspondence between the set of triangulations of a convex  $n$ -gon and the set of binary trees with  $n - 2$  nodes. In effect, these two problems differ only in the cost function for a single node/triangle.



A polygon triangulation and the corresponding binary tree. (Squares represent null pointers.)

A third problem that fits into the same mold is the infamous matrix chain multiplication problem. Using the standard algorithm, we can multiply a  $p \times q$  matrix by a  $q \times r$  matrix using

$O(pqr)$  arithmetic operations; the result is a  $p \times r$  matrix. If we have three matrices to multiply, the cost depends on which pair we multiply first. For example, suppose  $A$  and  $C$  are  $1000 \times 2$  matrices and  $B$  is a  $2 \times 1000$  matrix. There are two different ways to compute the threefold product  $ABC$ :

- **$(AB)C$** : Computing  $AB$  takes  $1000 \cdot 2 \cdot 1000 = 2\,000\,000$  operations and produces a  $1000 \times 1000$  matrix. Multiplying this matrix by  $C$  takes  $1000 \cdot 1000 \cdot 2 = 2\,000\,000$  additional operations. So the total cost of  $(AB)C$  is  $4\,000\,000$  operations.
- **$A(BC)$** : Computing  $BC$  takes  $2 \cdot 1000 \cdot 2 = 4\,000$  operations and produces a  $2 \times 2$  matrix. Multiplying  $A$  by this matrix takes  $1000 \cdot 2 \cdot 2 = 4\,000$  additional operations. So the total cost of  $A(BC)$  is only  $8\,000$  operations.

Now suppose we are given an array  $D[0..n]$  as input, indicating that each matrix  $M_i$  has  $D[i-1]$  rows and  $D[i]$  columns. We have an exponential number of possible ways to compute the  $n$ -fold product  $\prod_{i=1}^n M_i$ . The following dynamic programming algorithm computes the number of arithmetic operations for the best possible parenthesization:

<pre> <u>MATRIXCHAINMULT:</u> for <math>i \leftarrow n</math> downto 1   <math>M[i, i+1] \leftarrow 0</math>   for <math>j \leftarrow i+2</math> to <math>n</math>     COMPUTEM(<math>i, j</math>) return <math>M[1, n]</math> </pre>	<pre> <u>COMPUTEM(<math>i, j</math>):</u> <math>M[i, j] \leftarrow \infty</math> for <math>k \leftarrow i+1</math> to <math>j-1</math>   <math>tmp \leftarrow (D[i] \cdot D[j] \cdot D[k]) + M[i, k] + M[k, j]</math>   if <math>M[i, j] &gt; tmp</math>     <math>M[i, j] \leftarrow tmp</math> </pre>
---	---

The derivation of this algorithm is left as a simple exercise.

The first nuts and bolts appeared in the middle 1400's. The bolts were just screws with straight sides and a blunt end. The nuts were hand-made, and very crude. When a match was found between a nut and a bolt, they were kept together until they were finally assembled.

In the Industrial Revolution, it soon became obvious that threaded fasteners made it easier to assemble products, and they also meant more reliable products. But the next big step came in 1801, with Eli Whitney, the inventor of the cotton gin. The lathe had been recently improved. Batches of bolts could now be cut on different lathes, and they would all fit the same nut.

Whitney set up a demonstration for President Adams, and Vice-President Jefferson. He had piles of musket parts on a table. There were 10 similar parts in each pile. He went from pile to pile, picking up a part at random. Using these completely random parts, he quickly put together a working musket.

— Karl S. Kruszelnicki ('Dr. Karl'), *Karl Trek*, December 1997

### 3 Randomized Algorithms (September 12)

#### 3.1 Nuts and Bolts

Suppose we are given  $n$  nuts and  $n$  bolts of different sizes. Each nut matches exactly one bolt and vice versa. The nuts and bolts are all almost exactly the same size, so we can't tell if one bolt is bigger than the other, or if one nut is bigger than the other. If we try to match a nut with a bolt, however, the nut will be either too big, too small, or just right for the bolt.

Our task is to match each nut to its corresponding bolt. But before we do this, let's try to solve some simpler problems, just to get a feel for what we can and can't do.

Suppose we want to find the nut that matches a particular bolt. The obvious algorithm — test every nut until we find a match — requires exactly  $n - 1$  tests in the worst case. We might have to check every bolt except one; if we get down the the last bolt without finding a match, we know that the last nut is the one we're looking for.<sup>1</sup>

Intuitively, in the 'average' case, this algorithm will look at approximately  $n/2$  nuts. But what exactly does 'average case' mean?

#### 3.2 Deterministic vs. Randomized Algorithms

Normally, when we talk about the running time of an algorithm, we mean the *worst-case* running time. This is the maximum, over all problems of a certain size, of the running time of that algorithm on that input:

$$T_{\text{worst-case}}(n) = \max_{|X|=n} T(X).$$

On extremely rare occasions, we will also be interested in the *best-case* running time:

$$T_{\text{best-case}}(n) = \min_{|X|=n} T(X).$$

The *average-case* running time is best defined by the *expected value*, over all inputs  $X$  of a certain size, of the algorithm's running time for  $X$ :<sup>2</sup>

$$T_{\text{average-case}}(n) = \mathbf{E}_{|X|=n} [T(X)] = \sum_{|X|=n} T(x) \cdot \Pr[X].$$

<sup>1</sup>"Whenever you lose something, it's always in the last place you look. So why not just look there first?"

<sup>2</sup>The notation  $\mathbf{E}[\ ]$  for expectation has nothing to do with the shift operator  $E$  used in the annihilator method for solving recurrences!



The problem with this definition is that we rarely, if ever, know what the probability of getting any particular input  $X$  is. We could compute average-case running times by assuming a particular probability distribution—for example, every possible input is equally likely—but this assumption doesn't describe reality very well. Most real-life data is decidedly non-random.

Instead of considering this rather questionable notion of average case running time, we will make a distinction between two kinds of algorithms: *deterministic* and *randomized*. A deterministic algorithm is one that always behaves the same way given the same input; the input completely *determines* the sequence of computations performed by the algorithm. Randomized algorithms, on the other hand, base their behavior not only on the input but also on several *random* choices. The same randomized algorithm, given the same input multiple times, may perform different computations in each invocation.

This means, among other things, that the running time of a randomized algorithm on a given input is no longer fixed, but is itself a random variable. When we analyze randomized algorithms, we are typically interested in the *worst-case expected* running time. That is, we look at the average running time for each input, and then choose the maximum over all inputs of a certain size:

$$T_{\text{worst-case expected}}(n) = \max_{|X|=n} \mathbb{E}[T(X)].$$

It's important to note here that we are making *no* assumptions about the probability distribution of possible inputs. All the randomness is inside the algorithm, where we can control it!

### 3.3 Back to Nuts and Bolts

Let's go back to the problem of finding the nut that matches a given bolt. Suppose we use the same algorithm as before, but at each step we choose a nut *uniformly at random* from the untested nuts. 'Uniformly' is a technical term meaning that each nut has exactly the same probability of being chosen.<sup>3</sup> So if there are  $k$  nuts left to test, each one will be chosen with probability  $1/k$ . Now what's the expected number of comparisons we have to perform? Intuitively, it should be about  $n/2$ , but let's formalize our intuition.

Let  $T(n)$  denote the number of comparisons our algorithm uses to find a match for a single bolt out of  $n$  nuts.<sup>4</sup> We still have some simple base cases  $T(1) = 0$  and  $T(2) = 1$ , but when  $n > 2$ ,  $T(n)$  is a random variable.  $T(n)$  is always between 1 and  $n - 1$ ; its actual value depends on our algorithm's random choices. We are interested in the *expected value* or *expectation* of  $T(n)$ , which is defined as follows:

$$\mathbb{E}[T(n)] = \sum_{k=1}^{n-1} k \cdot \Pr[T(n) = k]$$

If the target nut is the  $k$ th nut tested, our algorithm performs  $\min\{k, n - 1\}$  comparisons. In particular, if the target nut is the last nut chosen, we don't actually test it. Because we choose the next nut to test uniformly at random, the target nut is equally likely—with probability exactly  $1/n$ —to be the first, second, third, or  $k$ th bolt tested, for any  $k$ . Thus:

$$\Pr[T(n) = k] = \begin{cases} 1/n & \text{if } k < n - 1, \\ 2/n & \text{if } k = n - 1. \end{cases}$$

<sup>3</sup>This is what most people think 'random' means, but they're wrong.

<sup>4</sup>Note that for this algorithm, the input is completely specified by the number  $n$ . Since we're choosing the nuts to test at random, even the order in which the nuts and bolts are presented doesn't matter. That's why I'm using the simpler notation  $T(n)$  instead of  $T(X)$ .

Plugging this into the definition of expectation gives us our answer.

$$\begin{aligned} \mathbb{E}[T(n)] &= \sum_{k=1}^{n-2} \frac{k}{n} + \frac{2(n-1)}{n} \\ &= \sum_{k=1}^{n-1} \frac{k}{n} + \frac{n-1}{n} \\ &= \frac{n(n-1)}{2n} + 1 - \frac{1}{n} \\ &= \frac{n+1}{2} - \frac{1}{n} \end{aligned}$$

We can get exactly the same answer by thinking of this algorithm recursively. We always have to perform at least one test. With probability  $1/n$ , we successfully find the matching nut and halt. With the remaining probability  $1 - 1/n$ , we recursively solve the same problem but with one fewer nut. We get the following recurrence for the expected number of tests:

$$T(1) = 0, \quad \mathbb{E}[T(n)] = 1 + \frac{n-1}{n} \mathbb{E}[T(n-1)]$$

To get the solution, we define a new function  $t(n) = n \mathbb{E}[T(n)]$  and rewrite:

$$t(1) = 0, \quad t(n) = n + t(n-1)$$

This recurrence translates into a simple summation, which we can easily solve.

$$\begin{aligned} t(n) &= \sum_{k=2}^n k = \frac{n(n+1)}{2} - 1 \\ \implies \mathbb{E}[T(n)] &= \frac{t(n)}{n} = \frac{n+1}{2} - \frac{1}{n} \end{aligned}$$

### 3.4 Finding All Matches

Not let's go back to the problem introduced at the beginning of the lecture: finding the matching nut for every bolt. The simplest algorithm simply compares every nut with every bolt, for a total of  $n^2$  comparisons. The next thing we might try is repeatedly finding an arbitrary matched pair, using our very first nuts and bolts algorithm. This requires

$$\sum_{i=1}^n (i-1) = \frac{n^2 - n}{2}$$

comparisons in the worst case. So we save roughly a factor of two over the really stupid algorithm. Not very exciting.

Here's another possibility. Choose a *pivot* bolt, and test it against every nut. Then test the matching pivot nut against every other bolt. After these  $2n - 1$  tests, we have one matched pair, and the remaining nuts and bolts are partitioned into two subsets: those smaller than the pivot pair and those larger than the pivot pair. Finally, recursively match up the two subsets. The worst-case number of tests made by this algorithm is given by the recurrence

$$\begin{aligned} T(n) &= 2n - 1 + \max_{1 \leq k \leq n} \{T(k-1) + T(n-k)\} \\ &= 2n - 1 + T(n-1) \end{aligned}$$

Along with the trivial base case  $T(0) = 0$ , this recurrence solves to

$$T(n) = \sum_{i=1}^n (2i - 1) = n^2.$$

In the worst case, this algorithm tests *every* nut-bolt pair! We could have been a little more clever—for example, if the pivot bolt is the smallest bolt, we only need  $n - 1$  tests to partition everything, not  $2n - 1$ —but cleverness doesn't actually help that much. We still end up with about  $n^2/2$  tests in the worst case.

However, since this recursive algorithm looks almost exactly like quicksort, and everybody 'knows' that the 'average-case' running time of quicksort is  $\Theta(n \log n)$ , it seems reasonable to guess that the average number of nut-bolt comparisons is also  $\Theta(n \log n)$ . As we shall see shortly, if the pivot bolt is always chosen *uniformly at random*, this intuition is exactly right.

### 3.5 Reductions to and from Sorting

The second algorithm for matching up the nuts and bolts looks exactly like quicksort. The algorithm not only matches up the nuts and bolts, but also sorts them by size.

In fact, the problems of sorting and matching nuts and bolts are equivalent, in the following sense. If the bolts were sorted, we could match the nuts and bolts in  $O(n \log n)$  time by performing a binary search with each nut. Thus, if we had an algorithm to sort the bolts in  $O(n \log n)$  time, we would immediately have an algorithm to match the nuts and bolts, starting from scratch, in  $O(n \log n)$  time. This process of *assuming* a solution to one problem and using it to solve another is called *reduction*—we can *reduce* the matching problem to the sorting problem in  $O(n \log n)$  time.

There is a reduction in the other direction, too. If the nuts and bolts were matched, we could sort them in  $O(n \log n)$  time using, for example, merge sort. Thus, if we have an  $O(n \log n)$  time algorithm for either sorting or matching nuts and bolts, we automatically have an  $O(n \log n)$  time algorithm for the other problem.

Unfortunately, since we aren't allowed to directly compare two bolts or two nuts, we can't use heapsort or mergesort to sort the nuts and bolts in  $O(n \log n)$  worst case time. In fact, the problem of sorting nuts and bolts *deterministically* in  $O(n \log n)$  time was only 'solved' in 1995<sup>5</sup>, but both the algorithms and their analysis are incredibly technical, the constant hidden in the  $O(\cdot)$  notation is extremely large, and worst of all, the solutions are nonconstructive—We know the algorithms exist, but we don't know what they look like!

Reductions will come up again later in the course when we start talking about lower bounds and NP-completeness.

### 3.6 Recursive Analysis

Intuitively, we can argue that our quicksort-like algorithm will usually choose a bolt of approximately median size, and so the average numbers of tests should be  $O(n \log n)$ . We can now finally formalize this intuition. To simplify the notation slightly, I'll write  $\bar{T}(n)$  in place of  $E[T(n)]$  everywhere.

Our randomized matching/sorting algorithm chooses its pivot bolt *uniformly at random* from the set of unmatched bolts. Since the pivot bolt is equally likely to be the smallest, second smallest,

<sup>5</sup>János Komlós, Yuan Ma, and Endre Szemerédi, Sorting nuts and bolts in  $O(n \log n)$  time, *SIAM J. Discrete Math* 11(3):347–372, 1998. See also Phillip G. Bradford, Matching nuts and bolts optimally, Technical Report MPI-I-95-1-025, Max-Planck-Institut für Informatik, September 1995. Bradford's algorithm is *slightly* simpler.

or  $k$ th smallest for any  $k$ , the expected number of tests performed by our algorithm is given by the following recurrence:

$$\begin{aligned}\bar{T}(n) &= 2n - 1 + \mathbb{E}_k[\bar{T}(k - 1) + \bar{T}(n - k)] \\ &= \boxed{2n - 1 + \frac{1}{n} \sum_{k=0}^{n-1} (\bar{T}(k - 1) + \bar{T}(n - k))}\end{aligned}$$

The base case is  $T(0) = 0$ . (We can save a few tests by setting  $T(1) = 1$ , but the analysis will be easier if we're a little stupid.)

Yuck. At this point, we could simply *guess* the solution, based on the incessant rumors that quicksort runs in  $O(n \log n)$  time in the average case, and prove our guess correct by induction. A similar inductive proof appears in [CLR, pp. 166–167], but it was removed from the new edition [CLRS]. That's okay; nobody ever understood that proof anyway.

However, if we're only interested in asymptotic bounds, we can afford to be a little conservative. What we'd *really* like is for the pivot bolt to be the median bolt, so that half the bolts are bigger and half the bolts are smaller. This isn't very likely, but there is a good chance that the pivot bolt is close to the median bolt. Let's say that a pivot bolt is *good* if it's in the middle half of the final sorted set of bolts, that is, bigger than at least  $n/4$  bolts and smaller than at least  $n/4$  bolts. If the pivot bolt is good, then the *worst* split we can have is into one set of  $3n/4$  pairs and one set of  $n/4$  pairs. If the pivot bolt is bad, then our algorithm is still better than starting over from scratch. Finally, a randomly chosen pivot bolt is good with probability  $1/2$ .

These simple observations give us the following simple recursive *upper bound* for the expected running time of our algorithm:

$$\bar{T}(n) \leq 2n - 1 + \frac{1}{2} \left( \bar{T}\left(\frac{3n}{4}\right) + \bar{T}\left(\frac{n}{4}\right) \right) + \frac{1}{2} \cdot \bar{T}(n)$$

A little algebra simplifies this even further:

$$\bar{T}(n) \leq 4n - 2 + \bar{T}\left(\frac{3n}{4}\right) + \bar{T}\left(\frac{n}{4}\right)$$

We can easily solve this recurrence using the recursion tree method, giving us the completely unsurprising upper bound  $\boxed{\bar{T}(n) = O(n \log n)}$ . Similar observations give us the matching lower bound  $\bar{T}(n) = \Omega(n \log n)$ .

### 3.7 Iterative Analysis

By making a simple change to our algorithm, which will have no effect on the number of tests, we can analyze it much more directly and exactly, without needing to solve a recurrence.

The recursive subproblems solved by quicksort can be laid out in a binary tree, where each node corresponds to a subset of the nuts and bolts. In the usual recursive formulation, the algorithm partitions the nuts and bolts at the root, then the left child of the root, then the leftmost grandchild, and so forth, recursively sorting everything on the left before starting on the right subproblem.

But we don't have to solve the subproblems in this order. In fact, we can visit the nodes in the recursion tree in any order we like, as long as the root is visited first, and any other node is visited after its parent. Thus, we can recast quicksort in the following iterative form. Choose a pivot bolt, find its match, and partition the remaining nuts and bolts into two subsets. Then pick a second pivot bolt and partition whichever of the two subsets contains it. At this point, we have

two matched pairs and three subsets of nuts and bolts. Continue choosing new pivot bolts and partitioning subsets, each time finding one match and increasing the number of subsets by one, until every bolt has been chosen as the pivot. At the end, every bolt has been matched, and the nuts and bolts are sorted.

Suppose we always choose the next pivot bolt *uniformly at random* from the bolts that haven't been pivots yet. Then no matter which subset contains this bolt, the pivot bolt is equally likely to be any bolt *in that subset*. That means our randomized iterative algorithm performs *exactly* the same set of tests as our randomized recursive algorithm, just in a different order.

Now let  $B_i$  denote the  $i$ th smallest bolt, and  $N_j$  denote the  $j$ th smallest nut. For each  $i$  and  $j$ , define an indicator variable  $X_{ij}$  that equals 1 if our algorithm compares  $B_i$  with  $N_j$  and zero otherwise. Then the total number of nut/bolt comparisons is exactly

$$T(n) = \sum_{i=1}^n \sum_{j=1}^n X_{ij}.$$

We are interested in the expected value of this double summation:

$$E[T(n)] = E \left[ \sum_{i=1}^n \sum_{j=1}^n X_{ij} \right] = \sum_{i=1}^n \sum_{j=1}^n E[X_{ij}].$$

This equation uses a crucial property of random variables called *linearity of expectation*: for any random variables  $X$  and  $Y$ , the sum of their expectations is equal to the expectation of their sum:  $E[X + Y] = E[X] + E[Y]$ . To analyze our algorithm, we only need to compute the expected value of each  $X_{ij}$ . By definition of expectation,

$$E[X_{ij}] = 0 \cdot \Pr[X_{ij} = 0] + 1 \cdot \Pr[X_{ij} = 1] = \Pr[X_{ij} = 1],$$

so we just need to calculate  $\Pr[X_{ij} = 1]$  for all  $i$  and  $j$ .

First let's assume that  $i < j$ . The only comparisons our algorithm performs are between some pivot bolt (or its partner) and a nut (or bolt) in the same subset. The only thing that can prevent us from comparing  $B_i$  and  $N_j$  is if some intermediate bolt  $B_k$ , with  $i < k < j$ , is chosen as a pivot before  $B_i$  or  $B_j$ . In other words:

**Our algorithm compares  $B_i$  and  $N_j$  if and only if the first pivot chosen from the set  $\{B_i, B_{i+1}, \dots, B_j\}$  is either  $B_i$  or  $B_j$ .**

Since the set  $\{B_i, B_{i+1}, \dots, B_j\}$  contains  $j - i + 1$  bolts, each of which is equally likely to be chosen first, we immediately have

$$E[X_{ij}] = \frac{2}{j - i + 1} \quad \text{for all } i < j.$$

Symmetric arguments give us  $E[X_{ij}] = \frac{2}{i - j + 1}$  for all  $i > j$ . Since our algorithm is a little stupid, every bolt is compared with its partner, so  $X_{ii} = 1$  for all  $i$ . (In fact, if a pivot bolt is the only bolt in its subset, we don't need to compare it against its partner, but this improvement complicates the analysis.)

Putting everything together, we get the following summation.

$$\begin{aligned}
 \mathbb{E}[T(n)] &= \sum_{i=1}^n \sum_{j=1}^n \mathbb{E}[X_{ij}] \\
 &= \sum_{i=1}^n \mathbb{E}[X_{ii}] + 2 \sum_{i=1}^n \sum_{j=i+1}^n \mathbb{E}[X_{ij}] \\
 &= \boxed{n + 4 \sum_{i=1}^n \sum_{j=i+1}^n \frac{1}{j-i+1}}
 \end{aligned}$$

This is quite a bit simpler than the recurrence we got before. In fact, with just a few lines of algebra, we can turn it into an exact, closed-form expression for the expected number of comparisons.

$$\begin{aligned}
 \mathbb{E}[T(n)] &= n + 4 \sum_{i=1}^n \sum_{j=2}^{n-i+1} \frac{1}{k} && \text{[substitute } k = j - i + 1\text{]} \\
 &= n + 4 \sum_{k=2}^n \sum_{i=1}^{n-k+1} \frac{1}{k} && \text{[reorder summations]} \\
 &= n + 4 \sum_{k=2}^n \frac{n-k+1}{k} \\
 &= n + 4 \left( (n-1) \sum_{k=2}^n \frac{1}{k} - \sum_{k=2}^n 1 \right) \\
 &= n + 4((n+1)(H_n - 1) - (n-1)) \\
 &= \boxed{4nH_n - 7n + 4H_n}
 \end{aligned}$$

Sure enough, it's  $\Theta(n \log n)$ .

### \*3.9 Masochistic Analysis

If we're feeling particularly masochistic, it is possible to solve the recurrence directly, all the way to an exact closed-form solution. [I'm including this only to show you it can be done; this won't be on the test.] First we simplify the recurrence slightly by combining symmetric terms.

$$\begin{aligned}
 T(n) &= 2n - 1 + \frac{1}{n} \sum_{k=0}^{n-1} (T(k+1) + T(n-k)) \\
 &= 2n - 1 + \frac{2}{n} \sum_{k=0}^{n-1} T(k)
 \end{aligned}$$

We then convert this 'full history' recurrence into a 'limited history' recurrence by shifting, and subtracting away common terms. (I call this "Magic step #1".) To make this slightly easier, we

first multiply both sides of the recurrence by  $n$  to get rid of the fractions.

$$\begin{aligned} nT(n) &= 2n^2 - n + 2 \sum_{k=0}^{n-1} T(k) \\ (n-1)T(n-1) &= \underbrace{2(n-1)^2 - (n-1)}_{2n^2 - 5n + 3} + 2 \sum_{k=0}^{n-2} T(k) \\ nT(n) - (n-1)T(n-1) &= 4n - 3 + 2T(n-1) \\ T(n) &= 4 - \frac{3}{n} + \frac{n+1}{n}T(n-1) \end{aligned}$$

To solve this limited-history recurrence, we define a new function  $t(n) = T(n)/(n+1)$ . (I call this “Magic step #2”.) This gives us an even simpler recurrence for  $t(n)$  in terms of  $t(n-1)$ :

$$\begin{aligned} t(n) &= \frac{T(n)}{n+1} \\ &= \frac{1}{n+1} \left( 4 - \frac{3}{n} + (n+1) \frac{T(n-1)}{n} \right) \\ &= \frac{4}{n+1} - \frac{3}{n(n+1)} + t(n-1) \\ &= \frac{7}{n+1} - \frac{3}{n} + t(n-1) \end{aligned}$$

I used the technique of partial fractions (remember calculus?) to replace  $\frac{1}{n(n+1)}$  with  $\frac{1}{n} - \frac{1}{n+1}$  in the last step. The base case for this recurrence is  $t(0) = 0$ . Once again, we have a recurrence that translates directly into a summation, which we can solve with just a few lines of algebra.

$$\begin{aligned} t(n) &= \sum_{i=1}^n \left( \frac{7}{i+1} - \frac{3}{i} \right) \\ &= 7 \sum_{i=1}^n \frac{1}{i+1} - 3 \sum_{i=1}^n \frac{1}{i} \\ &= 7(H_{n+1} - 1) - 3H_n \\ &= 4H_n - 7 + \frac{7}{n+1} \end{aligned}$$

The last step uses the recursive definition of the harmonic numbers:  $H_{n+1} = H_n + \frac{1}{n+1}$ . Finally, substituting  $T(n) = (n+1)t(n)$  and simplifying gives us the exact solution to the original recurrence.

$$T(n) = 4(n+1)H_n - 7(n+1) + 7 = \boxed{4nH_n - 7n + 4H_n}$$

Surprise, surprise, we get exactly the same solution!

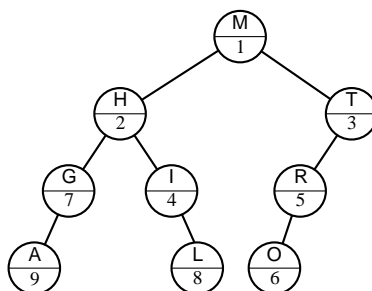
*I thought the following four [rules] would be enough, provided that I made a firm and constant resolution not to fail even once in the observance of them. The first was never to accept anything as true if I had not evident knowledge of its being so. . . . The second, to divide each problem I examined into as many parts as was feasible, and as was requisite for its better solution. The third, to direct my thoughts in an orderly way. . . . establishing an order in thought even when the objects had no natural priority one to another. And the last, to make throughout such complete enumerations and such general surveys that I might be sure of leaving nothing out.*

— René Descartes, *Discours de la Méthode* (1637)

## 4 Randomized Treaps (September 17)

### 4.1 Treaps

In this lecture, we will consider binary trees where every internal node has both a *search key* and a *priority*. In our examples, we will use letters for the search keys and numbers for the priorities. A *treap* is a binary tree where the inorder sequence of search keys is sorted and each node's priority is smaller than the priorities of its children.<sup>1</sup> In other words, a treap is simultaneously a binary search tree for the search keys and a (min-)heap for the priorities.



A treap. The top half of each node shows its search key and the bottom half shows its priority.

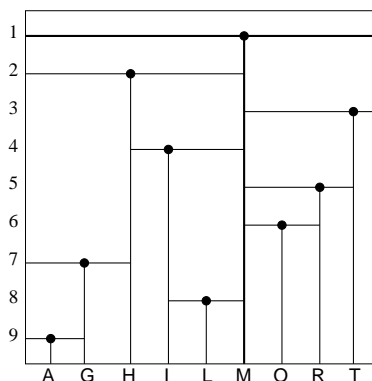
I'll assume from now on that all the keys and priorities are distinct. Under this assumption, we can easily prove by induction that the structure of a treap is completely determined by the search keys and priorities of its nodes. Since it's a heap, the node  $v$  with highest priority must be the root. Since it's also a binary search tree, any node  $u$  with  $key(u) < key(v)$  must be in the left subtree, and any node  $w$  with  $key(w) > key(v)$  must be in the right subtree. Finally, since the subtrees are treaps, by induction, their structures are completely determined. The base case is the trivial empty treap.

Another way to describe the structure is that a treap is exactly the binary tree that results by inserting the nodes one at a time into an initially empty tree, in order of increasing priority, using the usual insertion algorithm. This is also easy to prove by induction.

A third way interprets the keys and priorities as the coordinates of a set of points in the plane. The root corresponds to a T whose joint lies on the topmost point. The T splits the plane into three parts. The top part is (by definition) empty; the left and right parts are split recursively. This interpretation has some interesting applications in computational geometry, which (unfortunately) we probably won't have time to talk about.

<sup>1</sup>Sometimes I hate English. Normally, 'higher priority' means 'more important', but 'first priority' is also more important than 'second priority'. Maybe 'posteriority' would be better; one student suggested 'unimportance'.





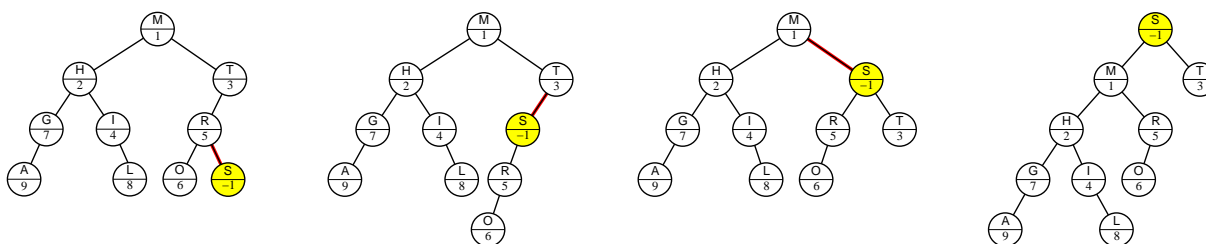
A geometric interpretation of the same treap.

Treaps were first discovered by Jean Vuillemin in 1980, but he called them *Cartesian trees*.<sup>2</sup> The word ‘treap’ was first used by Edward McCreight around 1980 to describe a slightly different data structure, but he later switched to the more prosaic name *priority search trees*.<sup>3</sup> Treaps were rediscovered and used to build randomized search trees by Cecilia Aragon and Raimund Seidel in 1989.<sup>4</sup> A different kind of randomized binary search tree, which uses random rebalancing instead of random priorities, was later discovered and analyzed by Conrado Martínez and Salvador Roura in 1996.<sup>5</sup>

## 4.2 Binary Search Tree Operations

The search algorithm is the usual one for binary search trees. The time for a successful search is proportional to the depth of the node. The time for an unsuccessful search is proportional to the depth of either its successor or its predecessor.

To insert a new node  $z$ , we start by using the standard binary search tree insertion algorithm to insert it at the bottom of the tree. At the point, the search keys still form a search tree, but the priorities may no longer form a heap. To fix the heap property, as long as  $z$  has smaller priority than its parent, perform a rotation at  $z$ . The running time is proportional to the depth of  $z$  before the rotations—we have to walk down the treap to insert  $z$ , and then walk back up the treap doing rotations. Another way to say this is that the time to insert  $z$  is roughly twice the time to perform an unsuccessful search for  $key(z)$ .



Left to right: After inserting  $(S, 10)$ , rotate it up to fix the heap property.

Right to left: Before deleting  $(S, 10)$ , rotate it down to make it a leaf.

<sup>2</sup>J. Vuillemin, A unifying look at data structures. *Commun. ACM* 23:229–239, 1980.

<sup>3</sup>E. M. McCreight. Priority search trees. *SIAM J. Comput.* 14(2):257–276, 1985.

<sup>4</sup>R. Seidel and C. R. Aragon. Randomized search trees. *Algorithmica* 16:464–497, 1996.

<sup>5</sup>C. Martínez and S. Roura. Randomized binary search trees. *J. ACM* 45(2):288–323, 1998. The results in this paper are virtually identical (including the constant factors!) to the corresponding results for treaps, although the analysis techniques are quite different.

Deleting a node is *exactly* like inserting a node, but in reverse order. Suppose we want to delete node  $z$ . As long as  $z$  is not a leaf, perform a rotation at the child of  $z$  with smaller priority. This moves  $z$  down a level and its smaller-priority child up a level. The choice of which child to rotate preserves the heap property everywhere except at  $z$ . When  $z$  becomes a leaf, chop it off.

We sometimes also want to *split* a treap  $T$  into two treaps  $T_<$  and  $T_>$  along some pivot key  $\pi$ , so that all the nodes in  $T_<$  have keys less than  $\pi$  and all the nodes in  $T_>$  have keys bigger than  $\pi$ . A simple way to do this is to insert a new node  $z$  with  $key(z) = \pi$  and  $priority(z) = -\infty$ . After the insertion, the new node is the root of the treap. If we delete the root, the left and right sub-treaps are exactly the trees we want. The time to split at  $\pi$  is roughly twice the time to (unsuccessfully) search for  $\pi$ .

Similarly, we may want to *merge* two treaps  $T_<$  and  $T_>$ , where every node in  $T_<$  has a smaller search key than any node in  $T_>$ , into one super-treap. Merging is just splitting in reverse—create a dummy root whose left sub-treap is  $T_<$  and whose right sub-treap is  $T_>$ , rotate the dummy node down to a leaf, and then cut it off.

### 4.3 Analysis

The cost of each of these operations is proportional to the depth  $d(v)$  of some node  $v$  in the treap.

- **Search:** A successful search for key  $k$  takes  $O(d(v))$  time, where  $v$  is the node with  $key(v) = k$ . For an unsuccessful search, let  $v^-$  be the inorder *predecessor* of  $k$  (the node whose key is just barely smaller than  $k$ ), and let  $v^+$  be the inorder *successor* of  $k$  (the node whose key is just barely larger than  $k$ ). Since the last node examined by the binary search is either  $v^-$  or  $v^+$ , the time for an unsuccessful search is either  $O(d(v^+))$  or  $O(d(v^-))$ .
- **Insert/Delete:** Inserting a new node with key  $k$  takes either  $O(d(v^+))$  time or  $O(d(v^-))$  time, where  $v^+$  and  $v^-$  are the predecessor and successor of the new node. Deletion is just insertion in reverse.
- **Split/Merge:** Splitting a treap at pivot value  $k$  takes either  $O(d(v^+))$  time or  $O(d(v^-))$  time, since it costs the same as inserting a new dummy root with search key  $k$  and priority  $-\infty$ . Merging is just splitting in reverse.

Since the depth of a node in a treap is  $\Theta(n)$  in the worst case, each of these operations has a worst-case running time of  $\Theta(n)$ .

### 4.4 Random Priorities

A *randomized binary search tree* is a treap in which the priorities are *independently and uniformly distributed continuous random variables*. That means that whenever we insert a new search key into the treap, we generate a random real number between (say) 0 and 1 and use that number as the priority of the new node. The only reason we're using real numbers is so that the probability of two nodes having the same priority is zero, since equal priorities make the analysis messy. In practice, we could just choose random integers from a large range, like 0 to  $2^{31} - 1$ , or random floating point numbers. Also, since the priorities are independent, each node is equally likely to have the smallest priority.

The cost of all the operations we discussed—search, insert, delete, split, join—is proportional to the depth of some node in the tree. Here we'll see that the *expected* depth of *any* node is  $O(\log n)$ , which implies that the expected running time for any of those operations is also  $O(\log n)$ .

Let  $x_k$  denote the node with the  $k$ th smallest search key. To analyze the expected depth, we define an indicator variable

$$A_k^i = [x_i \text{ is a proper ancestor of } x_k].$$

(The superscript doesn't mean power in this case; it just a reminder of which node is supposed to be further up in the tree.) Since the depth  $d(v)$  of  $v$  is just the number of proper ancestors of  $v$ , we have the following identity:

$$d(x_k) = \sum_{i=1}^n A_k^i.$$

Now we can express the *expected* depth of a node in terms of these indicator variables as follows.

$$E[d(x_k)] = \sum_{i=1}^n \Pr[A_k^i = 1]$$

(Just as in our analysis of matching nuts and bolts in Lecture 3, we're using linearity of expectation and the fact that  $E[X] = \Pr[X = 1]$  for any indicator variable  $X$ .) So to compute the expected depth of a node, we just have to compute the probability that some node is a proper ancestor of some other node.

Fortunately, we can do this easily once we prove a simple structural lemma. Let  $X(i, k)$  denote either the subset of treap nodes  $\{x_i, x_{i+1}, \dots, x_k\}$  or the subset  $\{x_k, x_{k+1}, \dots, x_i\}$ , depending on whether  $i < k$  or  $i > k$ .  $X(i, k)$  and  $X(k, i)$  always denote precisely the same subset, and this subset contains  $|k - i| + 1$  nodes.  $X(1, n) = X(n, 1)$  contains all  $n$  nodes in the treap.

**Lemma 1.** *For all  $i \neq k$ ,  $x_i$  is a proper ancestor of  $x_k$  if and only if  $x_i$  has the smallest priority among all nodes in  $X(i, k)$ .*

**Proof:** If  $x_i$  is the root, then it is an ancestor of  $x_k$ , and by definition, it has the smallest priority of *any* node in the treap, so it must have the smallest priority in  $X(i, k)$ .

On the other hand, if  $x_k$  is the root, then  $x_i$  is not an ancestor of  $x_k$ , and indeed  $x_i$  does not have the smallest priority in  $X(i, k)$  —  $x_k$  does.

On the gripping hand<sup>6</sup>, suppose some other node  $x_j$  is the root. If  $x_i$  and  $x_k$  are in different subtrees, then either  $i < j < k$  or  $i > j > k$ , so  $x_j \in X(i, k)$ . In this case,  $x_i$  is not an ancestor of  $x_k$ , and indeed  $x_i$  does not have the smallest priority in  $X(i, k)$  —  $x_j$  does.

Finally, if  $x_i$  and  $x_k$  are in the same subtree, the lemma follows inductively (or, if you prefer, recursively), since the subtree is a smaller treap. The empty treap is the trivial base case.  $\square$

Since each node in  $X(i, k)$  is equally likely to have smallest priority, we immediately have the probability we wanted:

$$\Pr[A_k^i = 1] = \frac{[i \neq k]}{|k - i| + 1} = \begin{cases} \frac{1}{k - i + 1} & \text{if } i < k \\ 0 & \text{if } i = k \\ \frac{1}{i - k + 1} & \text{if } i > k \end{cases}$$

---

<sup>6</sup>See Larry Niven and Jerry Pournelle, *The Gripping Hand*, Pocket Books, 1994.

To compute the expected depth of a node  $x_k$ , we just plug this probability into our formula and grind through the algebra.

$$\begin{aligned}
 E[d(x_k)] &= \sum_{i=1}^n \Pr[A_k^i = 1] \\
 &= \sum_{i=1}^{k-1} \frac{1}{k-i+1} + \sum_{i=k+1}^n \frac{1}{i-k+1} \\
 &= \sum_{j=2}^k \frac{1}{j} + \sum_{i=2}^{n-k} \frac{1}{i} \\
 &= H_k - 1 + H_{n-k} - 1 \\
 &< \ln k + \ln(n-k) - 2 \\
 &< 2 \ln n - 2.
 \end{aligned}$$

In conclusion, every search, insertion, deletion, split, and merge operation in an  $n$ -node randomized binary search tree takes  $O(\log n)$  expected time.

Since a treap is exactly the binary tree that results when you insert the keys in order of increasing priority, a randomized treap is the result of inserting the keys in *random* order. So our analysis also automatically gives us the expected depth of any node in a binary tree built by random insertions (without using priorities).

#### 4.5 Randomized Quicksort (Again?!)

We've already seen two completely different ways of describing randomized quicksort. The first is the familiar recursive one: choose a random pivot, partition, and recurse. The second is a less familiar iterative version: repeatedly choose a new random pivot, partition whatever subset contains it, and continue. But there's a third way to describe randomized quicksort, this time in terms of binary search trees.

RANDOMIZEDQUICKSORT:  
 $T \leftarrow$  an empty binary search tree  
 insert the keys into  $T$  in *random order*  
 output the inorder sequence of keys in  $T$

Our treap analysis tells us that this algorithm will run in  $O(n \log n)$  expected time, since each key is inserted in  $O(\log n)$  expected time.

Why is this quicksort? Just like last time, all we've done is rearrange the order of the comparisons. Intuitively, the binary tree is just the recursion tree created by the normal version of quicksort. In the recursive formulation, we compare the initial pivot against everything else and then recurse. In the binary tree formulation, the first "pivot" becomes the root of the tree without any comparisons, but then later as each other key is inserted into the tree, it is compared against the root. Either way, the first pivot chosen is compared with everything else. The partition splits the remaining items into a left subarray and a right subarray; in the binary tree version, these are exactly the items that go into the left subtree and the right subtree. Since both algorithms define the same two subproblems, by induction, both algorithms perform the same comparisons.

We even saw the probability  $\frac{1}{|k-i|+1}$  before, when we were talking about sorting nuts and bolts with a variant of randomized quicksort. In the more familiar setting of sorting an array of numbers,

the probability that randomized quicksort compares the  $i$ th largest and  $k$ th largest elements is exactly  $\frac{2}{|k-i|+1}$ . The binary tree version compares  $x_i$  and  $x_k$  if and only if  $x_i$  is an ancestor of  $x_k$  or vice versa, so the probabilities are exactly the same.

**Jaques:** *But, for the seventh cause; how did you find the quarrel on the seventh cause?*

**Touchstone:** *Upon a lie seven times removed:—bear your body more seeming, Audrey:—as thus, sir. I did dislike the cut of a certain courtier's beard: he sent me word, if I said his beard was not cut well, he was in the mind it was: this is called the Retort Courteous. If I sent him word again 'it was not well cut,' he would send me word, he cut it to please himself: this is called the Quip Modest. If again 'it was not well cut,' he disabled my judgment: this is called the Reply Churlish. If again 'it was not well cut,' he would answer, I spake not true: this is called the Reproof Valiant. If again 'it was not well cut,' he would say I lied: this is called the Counter-cheque Quarrelsome: and so to the Lie Circumstantial and the Lie Direct.*

**Jaques:** *And how oft did you say his beard was not well cut?*

**Touchstone:** *I durst go no further than the Lie Circumstantial, nor he durst not give me the Lie Direct; and so we measured swords and parted.*

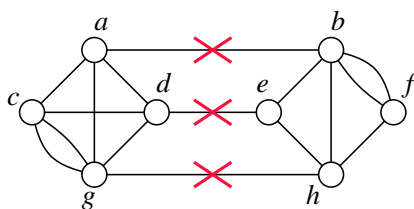
— William Shakespeare, *As You Like It* Act V, Scene 4 (1600)

## 5 Randomized Minimum Cut (September 19)

### 5.1 Setting Up the Problem

This lecture considers a problem that arises in robust network design. Suppose we have a connected multigraph<sup>1</sup>  $G$  representing a communications network like the UIUC telephone system, the internet, or Al-Qaeda. In order to disrupt the network, an enemy agent plans to remove some of the edges in this multigraph (by cutting wires, placing police at strategic drop-off points, or paying street urchins to ‘lose’ messages) to separate it into multiple components. Since his country is currently having an economic crisis, the agent wants to remove as few edges as possible to accomplish this task.

More formally, a *cut* partitions the nodes of  $G$  into two nonempty subsets. The *size* of the cut is the number of *crossing edges*, which have one endpoint in each subset. Finally, a *minimum* cut in  $G$  is a cut with the smallest number of crossing edges. The same graph may have several minimum cuts.



A multigraph whose minimum cut has three edges.

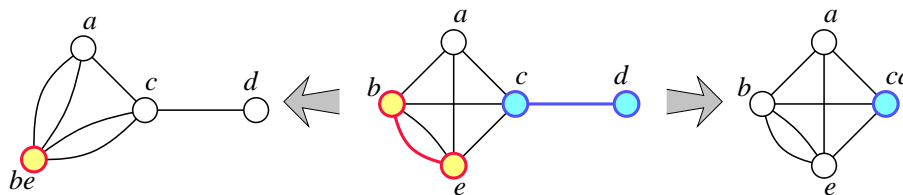
This problem has a long history. The classical deterministic algorithms for this problem rely on *network flow* techniques, which are discussed in Chapter 26 of CLRS. The fastest such algorithms run in  $O(n^3)$  time and are quite complex and difficult to understand (unless you’re already familiar with network flows). Here I’ll describe a relatively simple randomized algorithm published by David Karger<sup>2</sup>, how was a Ph.D. student at Stanford at the time.

Karger’s algorithm uses a primitive operation called *collapsing an edge*. Suppose  $u$  and  $v$  are vertices that are connected by an edge in some multigraph  $G$ . To collapse the edge  $\{u, v\}$ , we create a new node called  $uv$ , replace any edge of the form  $u, w$  or  $v, w$  with a new edge  $uv, w$ , and

<sup>1</sup>A multigraph allows multiple edges between the same pair of nodes. Everything in this lecture could be rephrased in terms of simple graphs where every edge has a non-negative weight, but this would make the algorithms and analysis slightly more complicated.

<sup>2</sup>D. R. Karger\*. Random sampling in cut, flow, and network design problems. Proc. 25th STOC, 648–657, 1994.

then delete the original vertices  $u$  and  $v$ . Equivalently, collapsing the edge shrinks the edge down to nothing, pulling the two endpoints together. The new collapsed graph is denoted  $G/\{u, v\}$ . We don't allow self-loops in our multigraphs; if there are multiple edges between  $u$  and  $v$ , collapsing any one of them deletes them all.



A graph  $G$  and two collapsed graphs  $G/\{b, e\}$  and  $G/\{c, d\}$ .

I won't describe how to actually implement collapsing an edge—it will be a homework exercise later in the course—but it can be done in  $O(n)$  time. Let's just accept collapsing as a black box subroutine for now.

The correctness of our algorithms will eventually boil down to the following simple observation: For any cut in  $G/\{u, v\}$ , there is a cut in  $G$  with exactly the same number of crossing edges. In fact, in some sense, the 'same' edges form the cut in both graphs. The converse is not necessarily true, however. For example, in the picture above, the original graph  $G$  has a cut of size 1, but the collapsed graph  $G/\{c, d\}$  does not.

This simple observation has two immediate but important consequences. First, collapsing an edge cannot decrease the minimum cut size. More importantly, collapsing an edge increases the minimum cut size if and only if that edge is part of *every* minimum cut.

## 5.2 Blindly Guessing

Let's start with an algorithm that tries to *guess* the minimum cut by randomly collapsing edges until the graph has only two vertices left.

```

GUESSMINCUT( $G$ ):
  for  $i \leftarrow n$  downto 2
    pick a random edge  $e$  in  $G$ 
     $G \leftarrow G/e$ 
  return the only cut in  $G$ 

```

Since each collapse takes  $O(n)$  time, this algorithm runs in  $O(n^2)$  time. Our earlier observations imply that as long as we never collapse an edge that lies in every minimum cut, our algorithm will actually guess correctly. But how likely is that?

Suppose  $G$  has only one minimum cut—if it actually has more than one, just pick your favorite—and this cut has size  $k$ . Every vertex of  $G$  must lie on at least  $k$  edges; otherwise, we could separate that vertex from the rest of the graph with an even smaller cut. Thus, the number of incident vertex-edge pairs is at least  $kn$ . Since every edge is incident to exactly two vertices,  $G$  must have at least  $kn/2$  edges. That implies that if we pick an edge in  $G$  uniformly at random, the probability of picking an edge in the minimum cut is at most  $2/n$ . In other words, the probability that we don't screw up on the very first step is at least  $1 - 2/n$ .

Once we've collapsed the first random edge, the rest of the algorithm proceeds recursively (with independent random choices) on the remaining  $(n - 1)$ -node graph. So the overall probability that

GUESSMINCUT returns the true minimum cut is given by the following recurrence:

$$P(n) \geq \frac{n-2}{n} \cdot P(n-1).$$

The base case for this recurrence is  $P(2) = 1$ . We can immediately expand this recurrence into a product, most of whose factors cancel out immediately.

$$P(n) \geq \prod_{i=3}^n \frac{i-2}{i} = \frac{\prod_{i=3}^n (i-2)}{\prod_{i=3}^n i} = \frac{\prod_{i=1}^{n-2} i}{\prod_{i=3}^n i} = \boxed{\frac{2}{n(n-1)}}$$

### 5.3 Blindly Guessing Over and Over

That's not very good. Fortunately, there's a simple method for increasing our chances of finding the minimum cut: run the guessing algorithm many times and return the smallest guess. Randomized algorithms folks like to call this idea *amplification*.

```

KARGERMINCUT(G):
  mink ← ∞
  for i ← 1 to N
    X ← GUESSMINCUT(G)
    if |X| < mink
      mink ← |X|
      minX ← X
  return minX

```

Both the running time and the probability of success will depend on the number of iterations  $N$ , which we haven't specified yet.

First let's figure out the probability that KARGERMINCUT returns the actual minimum cut. The only way for the algorithm to return the wrong answer is if GUESSMINCUT fails  $N$  times in a row. Since each guess is independent, our probability of success is at least

$$1 - \left(1 - \frac{2}{n(n-1)}\right)^N.$$

We can simplify this using one of the most important (and easy) inequalities known to mankind:

$$1 - x \leq e^{-x}$$

So our success probability is at least

$$1 - e^{-2N/n(n-1)}.$$

By making  $N$  larger, we can make this probability arbitrarily close to 1, but never equal to 1. In particular, if we set  $N = c \binom{n}{2} \ln n$  for some constant  $c$ , then KARGERMINCUT is correct with probability at least

$$1 - e^{-c \ln n} = 1 - \frac{1}{n^c}.$$



When the failure probability is a polynomial fraction, we say that the algorithm is correct *with high probability*. Thus, KARGERMINCUT computes the minimum cut of any  $n$ -node graph in  $O(n^4 \log n)$  time.

If we make the number of iterations even larger, say  $N = n^2(n - 1)/2$ , the success probability becomes  $1 - e^{-n}$ . When the failure probability is exponentially small like this, we say that the algorithm is correct with *very high probability*. In practice, very high probability is usually overkill; high probability is enough. (Remember, there is a small but non-zero probability that your computer will transform itself into a kitten before your program is finished.)

## 5.4 Not-So-Blindly Guessing

The  $O(n^4 \log n)$  running time is actually comparable to some of the simpler flow-based algorithms, but it's nothing to get excited about. But we can improve our guessing algorithm, and thus decrease the number of iterations in the outer loop, by observing that *as the graph shrinks, the probability of collapsing an edge in the minimum cut increases*. At first the probability is quite small, only  $2/n$ , but near the end of execution, when the graph has only three vertices, we have a  $2/3$  chance of screwing up!

A simple technique for working around this increasing probability of error was developed by David Karger and Cliff Stein.<sup>3</sup> Their idea is to group the first several random collapses a 'safe' phase, so that the cumulative probability of screwing up is small—less than  $1/2$ , say—and a 'dangerous' phase, which is much more likely to screw up.

The safe phase shrinks the graph from  $n$  nodes to  $n/\sqrt{2}$  nodes,<sup>4</sup> using a sequence of  $n - n/\sqrt{2}$  random collapses. Following our earlier analysis, the probability that any of these safe collapses touches the minimum cut is at most

$$\prod_{i=n/\sqrt{2}+1}^n \frac{i-2}{i} = \frac{(n/\sqrt{2})((n/\sqrt{2})-1)}{n(n-1)} < 1/2.$$

Now, to get around the danger of the dangerous phase, we use amplification. Instead of running through the dangerous phase one, we run it twice and keep the best of the two answers. And of course, we treat the dangerous phase recursively, so we actually obtain a larger binary recursion tree, which gets wider and wider as we get closer to the base case, instead of a single path. More formally, the algorithm looks like this:

$\overline{\text{CONTRACT}(G, m):}$ for $i \leftarrow n$ downto $m$ pick a random edge $e$ in $G$ $G \leftarrow G/e$ return $G$	$\overline{\text{BETTERGUESS}(G):}$ $X_1 \leftarrow \text{BETTERGUESS}(\text{CONTRACT}(G, n/\sqrt{2}))$ $X_2 \leftarrow \text{BETTERGUESS}(\text{CONTRACT}(G, n/\sqrt{2}))$ return $\min\{X_1, X_2\}$
---	--

This might look like we're just doing the same thing twice, but remember that CONTRACT (and thus BETTERGUESS) is randomized. Each call to CONTRACT contracts a different random set of edges;  $X_1$  and  $X_2$  are almost always different cuts.

<sup>3</sup>D. R. Karger\* and C. Stein. An  $\tilde{O}(n^2)$  algorithm for minimum cuts. Proc. 25th STOC, 757–765, 1993. Yes, *that* Cliff Stein.

<sup>4</sup>Strictly speaking, I should say  $\lfloor n/\sqrt{2} \rfloor$ , because otherwise, we'd have an irrational number of nodes! But this only adds some floors and ceilings to our recurrences, which we know that we can remove with domain transformations, so I'll just take them out now.

BETTERGUESS correctly returns the minimum cut unless *both* of the first two lines gives the wrong result. A single call to BETTERGUESS(CONTRACT( $G, n/\sqrt{2}$ )) gives the correct answer if none of the min cut edges are CONTRACTED and if the recursive BETTERGUESS is correct. So we have the following recurrence for the probability of success for an  $n$ -node graph:

$$P(n) \geq 1 - \left(1 - \frac{1}{2}P\left(\frac{n}{\sqrt{2}}\right)\right)^2 = P\left(\frac{n}{\sqrt{2}}\right) - \frac{1}{4}P\left(\frac{n}{\sqrt{2}}\right)^2$$

I don't know how to derive it from scratch, but we can easily prove by induction that  $P(n) = 1/\lg n$  is a solution for this recurrence. The base case is  $P(2) = 1 = 1/\lg 2$ .

$$\begin{aligned} P(n) &\geq P\left(\frac{n}{\sqrt{2}}\right) - \frac{1}{4}P\left(\frac{n}{\sqrt{2}}\right)^2 \\ &= \frac{1}{\lg \frac{n}{\sqrt{2}}} - \frac{1}{4\lg^2 \frac{n}{\sqrt{2}}} \\ &= \frac{1}{\lg n - \frac{1}{2}} - \frac{1}{4(\lg n - \frac{1}{2})^2} \\ &= \frac{4\lg n - 3}{4\lg^2 n - 4\lg n + 1} \\ &= \frac{1}{\lg n} + \frac{1 - \frac{1}{\lg n}}{4\lg^2 n - 4\lg n + 1} \\ &> \frac{1}{\lg n} \checkmark \end{aligned}$$

The last step used the fact that  $n > 2$ ; otherwise, that ugly fraction we threw out might be zero or negative.

For the running time, we get a simple recurrence that is easily solved using the Master theorem.

$$T(n) = O(n^2) + 2T\left(\frac{n}{\sqrt{2}}\right) = \boxed{O(n^2 \log n)}$$

So all this splitting and recursing has slowed down the guessing algorithm slightly, but the probability of failure is *exponentially* smaller!

Now if we call BETTERGUESS  $N = c(\lg n)(\ln n)$  times, for some constant  $c$ , the overall probability of success is

$$1 - \left(1 - \frac{1}{\lg n}\right)^{c(\lg n)(\ln n)} \geq 1 - e^{-c \ln n} = 1 - \frac{1}{n^c}.$$

In other words, we now have an algorithm that computes the minimum cut with high probability in only  $\boxed{O(n^2 \log^3 n)}$  time!

<i>Aitch</i>	<i>Ex</i>
<i>Are</i>	<i>Eye</i>
<i>Ay</i>	<i>Gee</i>
<i>Bee</i>	<i>Jay</i>
<i>Cue</i>	<i>Kay</i>
<i>Dee</i>	<i>Oh</i>
<i>Double U</i>	<i>Pea</i>
<i>Ee</i>	<i>See</i>
<i>Ef</i>	<i>Tee</i>
<i>El</i>	<i>Vee</i>
<i>Em</i>	<i>Wy</i>
<i>En</i>	<i>Yu</i>
<i>Ess</i>	<i>Zee</i>

— Sidney Harris, “The Alphabet in Alphabetical Order”

## 6 Hash Tables (September 24)

### 6.1 Introduction

A *hash table* is a data structure for storing a set of items, so that we can quickly determine whether an item is or is not in the set. The basic idea is to pick a *hash function*  $h$  that maps every possible item  $x$  to a small integer  $h(x)$ . Then we store  $x$  in slot  $h(x)$  in an array. The array is the hash table.

Let’s be a little more specific. We want to store a set of  $n$  items. Each item is an element of some finite<sup>1</sup> set  $\mathcal{U}$  called the *universe*; we use  $u$  to denote the size of the universe, which is just the number of items in  $\mathcal{U}$ . A hash table is an array  $T[1..m]$ , where  $m$  is another positive integer, which we call the *table size*. Typically,  $m$  is much smaller than  $u$ . A *hash function* is a function

$$h: \mathcal{U} \rightarrow \{0, 1, \dots, m - 1\}$$

that maps each possible item in  $\mathcal{U}$  to a slot in the hash table. We say that an item  $x$  *hashes* to the slot  $T[h(x)]$ .

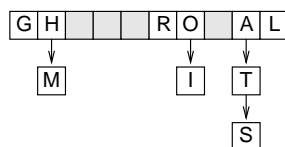
Of course, if  $u = m$ , then we can always just use the trivial hash function  $h(x) = x$ . In other words, use the item itself as the index into the table. This is called a *direct access table* (or more commonly, an *array*). In most applications, though, the universe of possible keys is orders of magnitude too large for this approach to be practical. Even when it is possible to allocate enough memory, we usually need to store only a small fraction of the universe. Rather than wasting lots of space, we should make  $m$  roughly equal to  $n$ , the number of items in the set we want to maintain.

What we’d like is for every item in our set to hash to a different position in the array. Unfortunately, unless  $m = u$ , this is too much to hope for, so we have to deal with *collisions*. We say that two items  $x$  and  $y$  *collide* if they have the same hash value:  $h(x) = h(y)$ . Since we obviously can’t store two items in the same slot of an array, we need to describe some methods for *resolving collisions*. The two most common methods are called *chaining* and *open addressing*.

<sup>1</sup>This finiteness assumption is necessary for several of the technical details to work out, but can be ignored in practice. To hash elements from an infinite universe (for example, the positive integers), pretend that the universe is actually finite but very very large. In fact, in *real* practice, the universe actually *is* finite but very very large. For example, on most modern computers, there are only  $2^{64}$  integers (unless you use a big integer package like GMP, in which case the number of integers is closer to  $2^{2^{32}}$ .)

## 6.2 Chaining

In a *chained* hash table, each entry  $T[i]$  is not just a single item, but rather (a pointer to) a linked list of all the items that hash to  $T[i]$ . Let  $\ell(x)$  denote the length of the list  $T[h(x)]$ . To see if an item  $x$  is in the hash table, we scan the entire list  $T[h(x)]$ . The worst-case time required to search for  $x$  is  $O(1)$  to compute  $h(x)$  plus  $O(1)$  for every element in  $T[h(x)]$ , or  $O(1 + \ell(x))$  overall. Inserting and deleting  $x$  also take  $O(1 + \ell(x))$  time.



A chained hash table with load factor 1.

In the worst case, every item would be hashed to the same value, so we'd get just one long list of  $n$  items. In principle, for any deterministic hashing scheme, a malicious adversary can always present a set of items with exactly this property. In order to defeat such malicious behavior, we'd like to use a hash function that is as random as possible. Choosing a truly random hash function is completely impractical, but since there are several heuristics for producing hash functions that behave close to randomly (on real data), we will analyze the performance *as though our hash function were completely random*. More formally, we make the following assumption.

**Simple uniform hashing assumption:** If  $x \neq y$  then  $\Pr[h(x) = h(y)] = 1/m$ .

Let's compute the expected value of  $\ell(x)$  under the simple uniform hashing assumption; this will immediately imply a bound on the expected time to search for an item  $x$ . To be concrete, let's suppose that  $x$  is not already stored in the hash table. For all items  $x$  and  $y$ , we define the indicator variable

$$C_{x,y} = [h(x) = h(y)].$$

(In case you've forgotten the bracket notation,  $C_{x,y} = 1$  if  $h(x) = h(y)$  and  $C_{x,y} = 0$  if  $h(x) \neq h(y)$ .) Since the length of  $T[h(x)]$  is precisely equal to the number of items that collide with  $x$ , we have

$$\ell(x) = \sum_{y \in T} C_{x,y}.$$

We can rewrite the simple uniform hashing assumption as follows:

$$x \neq y \implies \mathbb{E}[C_{x,y}] = \frac{1}{m}.$$

Now we just have to grind through the definitions.

$$\mathbb{E}[\ell(x)] = \sum_{y \in T} \mathbb{E}[C_{x,y}] = \sum_{y \in T} \frac{1}{m} = \frac{n}{m}$$

We call this fraction  $n/m$  the *load factor* of the hash table. Since the load factor shows up everywhere, we will give it its own symbol  $\alpha$ .

$$\alpha = \frac{n}{m}$$

Our analysis implies that the expected time for an unsuccessful search in a chained hash table is  $\Theta(1+\alpha)$ . As long as the number of items  $n$  is only a constant factor bigger than the table size  $m$ , the search time is a constant. A similar analysis gives the same expected time bound<sup>2</sup> for a successful search.

Obviously, linked lists are not the only data structure we could use to store the chains; any data structure that can store a set of items will work. For example, if the universe  $\mathcal{U}$  has a total ordering, we can store each chain in a balanced binary search tree. This reduces the worst-case time for a search to  $O(1 + \log \ell(x))$ , and under the simple uniform hashing assumption, the expected time for a search is  $O(1 + \log \alpha)$ .

Another possibility is to keep the overflow lists in hash tables! Specifically, for each  $T[i]$ , we maintain a hash table  $T_i$  containing all the items with hash value  $i$ . To keep things efficient, we make sure the load factor of each secondary hash table is always a constant less than 1; this can be done with only constant *amortized* overhead.<sup>3</sup> Since the load factor is constant, a search in any secondary table always takes  $O(1)$  expected time, so the total expected time to search in the top-level hash table is also  $O(1)$ .

### 6.3 Open Addressing

Another method we can use to resolve collisions is called *open addressing*. Here, rather than building secondary data structures, we resolve collisions by looking elsewhere in the table. Specifically, we have a sequence of hash functions  $\langle h_0, h_1, h_2, \dots, h_{m-1} \rangle$ , such that for any item  $x$ , the *probe sequence*  $\langle h_0(x), h_1(x), \dots, h_{m-1}(x) \rangle$  is a permutation of  $\langle 0, 1, 2, \dots, m-1 \rangle$ . In other words, different hash functions in the sequence always map  $x$  to different locations in the hash table.

We search for  $x$  using the following algorithm, which returns the array index  $i$  if  $T[i] = x$ , ‘absent’ if  $x$  is not in the table but there is an empty slot, and ‘full’ if  $x$  is not in the table and there are no empty slots.

<pre> OPENADDRESSSEARCH(x):   for i ← 0 to m - 1     if T[h<sub>i</sub>(x)] = x       return h<sub>i</sub>(x)     else if T[h<sub>i</sub>(x)] = ∅       return ‘absent’   return ‘full’ </pre>
--

The algorithm for inserting a new item into the table is similar; only the second-to-last line is changed to  $T[h_i(x)] \leftarrow x$ . Notice that for an open-addressed hash table, the load factor is never bigger than 1.

Just as with chaining, we’d like the sequence of hash values to be random, and for purposes of analysis, there is a stronger uniform hashing assumption that gives us constant expected search and insertion time.

**Strong uniform hashing assumption:** For any item  $x$ , the probe sequence  $\langle h_0(x), h_1(x), \dots, h_{m-1}(x) \rangle$  is equally likely to be any permutation of the set  $\{0, 1, 2, \dots, m-1\}$ .

<sup>2</sup>but with smaller constants hidden in the  $O()$ —see p.225 of CLR for details.

<sup>3</sup>This means that a single insertion or deletion may take more than constant time, but the total time to handle any sequence of  $k$  insertions or deletions, for any  $k$ , is  $O(k)$  time. We’ll discuss amortized running times after the first midterm. This particular result will be an easy homework problem.

Let's compute the expected time for an unsuccessful search using this stronger assumption. Suppose there are currently  $n$  elements in the hash table. Our strong uniform hashing assumption has two important consequences:

- The initial hash value  $h_0(x)$  is equally likely to be any integer in the set  $\{0, 1, 2, \dots, m - 1\}$ .
- If we ignore the first probe, the remaining probe sequence  $\langle h_1(x), h_2(x), \dots, h_{m-1}(x) \rangle$  is equally likely to be any permutation of the smaller set  $\{0, 1, 2, \dots, m - 1\} \setminus \{h_0(x)\}$ .

The first sentence implies that the probability that  $T[h_0(x)]$  is occupied is exactly  $n/m$ . The second sentence implies that if  $T[h_0(x)]$  is occupied, *our search algorithm recursively searches the rest of the hash table!* Since the algorithm will never again probe  $T[h_0(x)]$ , for purposes of analysis, we might as well pretend that slot in the table no longer exists. Thus, we get the following recurrence for the expected number of probes, as a function of  $m$  and  $n$ :

$$E[T(m, n)] = 1 + \frac{n}{m} E[T(m - 1, n - 1)].$$

The trivial base case is  $T(m, 0) = 1$ ; if there's nothing in the hash table, the first probe always hits an empty slot. We can now easily prove by induction that  $E[T(m, n)] \leq m/(m - n)$ :

$$\begin{aligned} E[T(m, n)] &= 1 + \frac{n}{m} E[T(m - 1, n - 1)] \\ &\leq 1 + \frac{n}{m} \cdot \frac{m - 1}{m - n} && \text{[induction hypothesis]} \\ &< 1 + \frac{n}{m} \cdot \frac{m}{m - n} && [m - 1 < m] \\ &= \frac{m}{m - n} \checkmark && \text{[algebra]} \end{aligned}$$

Rewriting this in terms of the load factor  $\alpha = n/m$ , we get  $E[T(m, n)] \leq 1/(1 - \alpha)$ . In other words, the expected time for an unsuccessful search is  $O(1)$ , unless the hash table is almost completely full.

In practice, however, we can't generate truly random probe sequences, so we use one of the following heuristics:

- **Linear probing:** We use a single hash function  $h(x)$ , and define  $h_i(x) = (h(x) + i) \bmod m$ . This is nice and simple, but collisions tend to make items in the table clump together badly, so this is not really a good idea.
- **Quadratic probing:** We use a single hash function  $h(x)$ , and define  $h_i(x) = (h(x) + i^2) \bmod m$ . Unfortunately, for certain values of  $m$ , the sequence of hash values  $\langle h_i(x) \rangle$  does not hit every possible slot in the table; we can avoid this problem by making  $m$  a prime number. (That's often a good idea anyway.) Although quadratic probing does not suffer from the same clumping problems as linear probing, it does have a weaker clustering problem: If two items have the same initial hash value, their entire probe sequences will be the same.
- **Double hashing:** We use two hash functions  $h(x)$  and  $h'(x)$ , and define  $h_i$  as follows:

$$h_i(x) = (h(x) + i \cdot h'(x)) \bmod m$$

To guarantee that this can hit every slot in the table, the *stride* function  $h'(x)$  and the table size  $m$  must be relatively prime. We can guarantee this by making  $m$  prime, but

a simpler solution is to make  $m$  a power of 2 and choose a stride function that is always odd. Double hashing avoids the clustering problems of linear and quadratic probing. In fact, the actual performance of double hashing is almost the same as predicted by the uniform hashing assumption, at least when  $m$  is large and the component hash functions  $h$  and  $h'$  are sufficiently random. This is the method of choice!<sup>4</sup>

## 6.4 Deleting from an Open-Addressed Hash Table

Deleting an item  $x$  from an open-addressed hash table is a bit more difficult than in a chained hash table. We can't simply clear out the slot in the table, because we may need to know that  $T[h(x)]$  is occupied in order to find some other item!

Instead, we should delete more or less the way we did with scapegoat trees. When we delete an item, we mark the slot that used to contain it as a *wasted* slot. A sufficiently long sequence of insertions and deletions could eventually fill the table with marks, leaving little room for any real data and causing searches to take linear time.

However, we can still get good *amortized* performance by using two rebuilding rules. First, if the number of items in the hash table exceeds  $m/4$ , double the size of the table ( $m \leftarrow 2m$ ) and rehash everything. Second, if the number of wasted slots exceeds  $m/2$ , clear all the marks and rehash everything in the table. Rehashing everything takes  $m$  steps to create the new hash table and  $O(n)$  expected steps to hash each of the  $n$  items. By charging a \$4 tax for each insertion and a \$2 tax for each deletion, we expect to have enough money to pay for any rebuilding.

In conclusion, the *expected amortized* cost of any insertion or deletion is  $O(1)$ , under the uniform hashing assumption. Notice that we're doing two very different kinds of averaging here. On the one hand, we are averaging the possible costs of *each individual search* over all possible probe sequences ('expected'). On the other hand, we are also averaging the costs of *the entire sequence* of operations to 'smooth out' the cost of rebuilding ('amortized'). Both randomization and amortization are necessary to get this constant time bound.

## 6.5 Universal Hashing

Now I'll describe how to generate hash functions that (at least in expectation) satisfy the uniform hashing assumption. We say that a set  $\mathcal{H}$  of hash function is *universal* if it satisfies the following property: For any items  $x \neq y$ , if a hash function  $h$  is chosen *uniformly at random* from the set  $\mathcal{H}$ , then  $\Pr[h(x) = h(y)] = 1/m$ . Note that this probability holds for *any* items  $x$  and  $y$ ; the randomness is entirely in choosing a hash function from the set  $\mathcal{H}$ .

To simplify the following discussion, I'll assume that the universe  $\mathcal{U}$  contains exactly  $m^2$  items, each represented as a pair  $(x, x')$  of integers between 0 and  $m - 1$ . (Think of the items as two-digit numbers in base  $m$ .) I will also assume that  $m$  is a prime number.

For any integers  $0 \leq a, b \leq m - 1$ , define the function  $h_{a,b}: \mathcal{U} \rightarrow \{0, 1, \dots, m - 1\}$  as follows:

$$h_{a,b}(x, x') = (ax + bx') \bmod m.$$

Then the set

$$\mathcal{H} = \{h_{a,b} \mid 0 \leq a, b \leq m - 1\}$$

of all such functions is universal. To prove this, we need to show that for any pair of distinct items  $(x, x') \neq (y, y')$ , exactly  $m$  of the  $m^2$  functions in  $\mathcal{H}$  cause a collision.

---

<sup>4</sup>...unless your hash tables are really huge, in which case linear probing has *far* better cache behavior, especially when the load factor is small.

Choose two items  $(x, x') \neq (y, y')$ , and assume without loss of generality<sup>5</sup> that  $x \neq y$ . A function  $h_{a,b} \in \mathcal{H}$  causes a collision between  $(x, x')$  and  $(y, y')$  if and only if

$$\begin{aligned} h_{a,b}(x, x') &= h_{a,b}(y, y') \\ (ax + bx') \bmod m &= (ay + by') \bmod m \\ ax + bx' &\equiv ay + by' \pmod{m} \\ a(x - y) &\equiv b(y' - x') \pmod{m} \\ a &\equiv \frac{b(y' - x')}{x - y} \pmod{m}. \end{aligned}$$

In the last step, we are using the fact that  $m$  is prime and  $x - y \neq 0$ , which implies that  $x - y$  has a unique multiplicative inverse modulo  $m$ .<sup>6</sup> Now notice for each possible value of  $b$ , the last identity defines a *unique* value of  $a$  such that  $h_{a,b}$  causes a collision. Since there are  $m$  possible values for  $b$ , there are exactly  $m$  hash functions  $h_{a,b}$  that cause a collision, which is exactly what we needed to prove.

Thus, if we want to achieve the constant expected time bounds described earlier, we should choose a random element of  $\mathcal{H}$  as our hash function, by generating two numbers  $a$  and  $b$  uniformly at random between 0 and  $m - 1$ . (Notice that this is *exactly* the same as choosing a element of  $\mathcal{U}$  uniformly at random.)

One perhaps undesirable ‘feature’ of this construction is that we have a small chance of choosing the trivial hash function  $h_{0,0}$ , which maps everything to 0. So in practice, if we happen to pick  $a = b = 0$ , we should reject that choice and pick new random numbers. By taking  $h_{0,0}$  out of consideration, we reduce the probability of a collision from  $1/m$  to  $(m - 1)/(m^2 - 1) = 1/(m + 1)$ . In other words, the set  $\mathcal{H} \setminus \{h_{0,0}\}$  is slightly *better* than universal.

This construction can be generalized easily to larger universes. Suppose  $u = m^r$  for some constant  $r$ , so that each element  $x \in \mathcal{U}$  can be represented by a vector  $(x_0, x_1, \dots, x_{r-1})$  of integers between 0 and  $m - 1$ . (Think of  $x$  as an  $r$ -digit number written in base  $m$ .) Then for each vector  $a = (a_0, a_1, \dots, a_{r-1})$ , define the corresponding hash function  $h_a$  as follows:

$$h_a(x) = (a_0x_0 + a_1x_1 + \dots + a_{r-1}x_{r-1}) \bmod m.$$

Then the set of all  $m^r$  such functions is universal.

<sup>5</sup>‘Without loss of generality’ is a phrase that appears (perhaps too) often in combinatorial proofs. What it means is that we are considering one of many possible cases, but once we see the proof for one case, the proofs for all the other cases are obvious thanks to some inherent symmetry. For this proof, we are not explicitly considering what happens when  $x = y$  and  $x' \neq y'$ .

<sup>6</sup>For example, the multiplicative inverse of 12 modulo 17 is 10, since  $12 \cdot 10 = 120 \equiv 1 \pmod{17}$ .



*For example, creating shortcuts by sprinkling a few diversely connected individuals throughout a large organization could dramatically speed up information flow between departments. On the other hand, because only a few random shortcuts are necessary to make the world small, subtle changes to networks have alarming consequences for the rapid spread of computer viruses, pernicious rumors, and infectious diseases.*

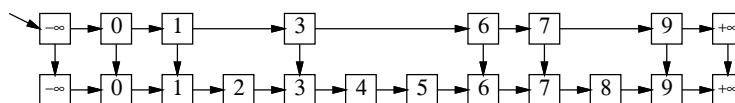
— Ivars Peterson, *Science News*, August 22, 1998.

## A Skip Lists

This lecture is about a probabilistic data structure called *skip lists*, first discovered by Bill Pugh in the late 1980's.<sup>1</sup> Skip lists have many of the desirable properties of balanced binary search trees, but their structure is completely different.

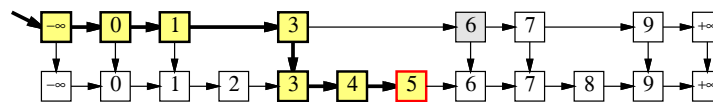
### A.1 Shortcuts

At a high level, a skip list is just a sorted, singly linked list with some shortcuts. To do a search in a normal singly-linked list of length  $n$ , we obviously need to look at  $n$  items in the worst case. To speed up this process, we can make a second-level list that contains roughly half the items from the original list. Specifically, for each item in the original list, we duplicate it with probability  $1/2$ . We then string together all the duplicates into a second sorted linked list, and add a pointer from each duplicate back to its original. Just to be safe, we also add sentinel nodes at the beginning and end of both lists.



A linked list with some randomly-chosen shortcuts.

Now we can find a value  $x$  in this augmented structure using a two-stage algorithm. First, we scan for  $x$  in the shortcut list, starting at the  $-\infty$  sentinel node. If we find  $x$ , we're done. Otherwise, we reach some value bigger than  $x$  and we know that  $x$  is not in the shortcut list. Let  $w$  be the largest item less than  $x$  in the shortcut list. In the second phase, we scan for  $x$  in the original list, starting from  $w$ . Again, if we reach a value bigger than  $x$ , we know that  $x$  is not in the data structure.



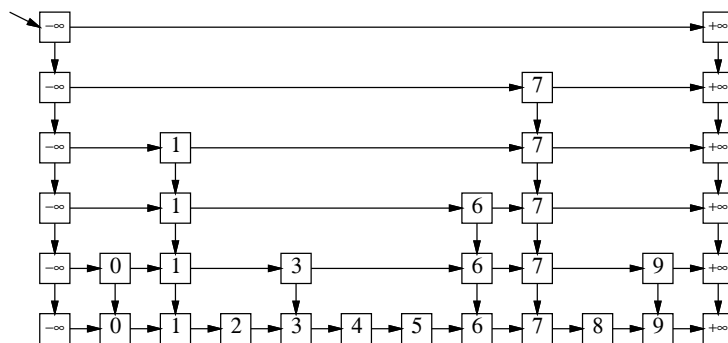
Searching for 5 in a list with shortcuts.

Since each node appears in the shortcut list with probability  $1/2$ , the expected number of nodes examined in the first phase is at most  $n/2$ . Only one of the nodes examined in the second phase has a duplicate. The probability that any node is followed by  $k$  nodes without duplicates is  $2^{-k}$ , so the expected number of nodes examined in the second phase is at most  $1 + \sum_{k \geq 0} 2^{-k} = 2$ . Thus, by adding these random shortcuts, we've reduced the cost of a search from  $n$  to  $n/2 + 2$ , roughly a factor of two in savings.

<sup>1</sup>William Pugh. Skip lists: A probabilistic alternative to balanced trees. *Commun. ACM* 33(6):668–676, 1990.

## A.2 Skip lists

Now there's an obvious improvement—add shortcuts to the shortcuts, and repeat recursively. That's exactly how skip lists are constructed. For each node in the original list, we flip a coin over and over until we get tails. For each heads, we make a duplicate of the node. The duplicates are stacked up in levels, and the nodes on each level are strung together into sorted linked lists. Each node  $v$  stores a search key ( $\text{key}(v)$ ), a pointer to its next lower copy ( $\text{down}(v)$ ), and a pointer to the next node in its level ( $\text{right}(v)$ ).



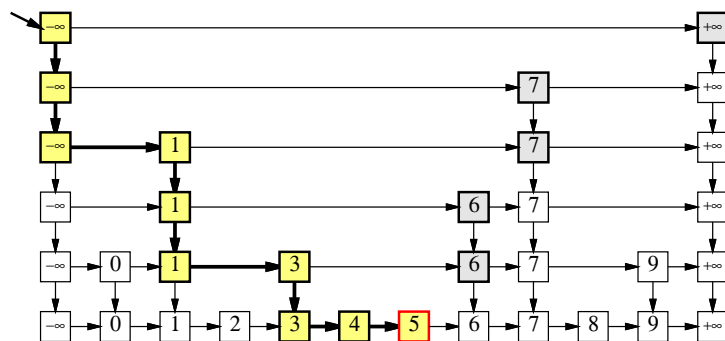
A skip list is a linked list with recursive random shortcuts.

The search algorithm for skip lists is very simple. Starting at the leftmost node  $L$  in the highest level, we scan through each level as far as we can without passing the target value  $x$ , and then proceed down to the next level. The search ends when we either reach a node with search key  $x$  or fail to find  $x$  on the lowest level.

```

SKIPLISTFIND( $x, L$ ):
   $v \leftarrow L$ 
  while ( $v \neq \text{NULL}$  and  $\text{key}(v) \neq x$ )
    if  $\text{key}(\text{right}(v)) > x$ 
       $v \leftarrow \text{down}(v)$ 
    else
       $v \leftarrow \text{right}(v)$ 
  return  $v$ 

```



Searching for 5 in a skip list.

Intuitively, Since each level of the skip lists has about half the number of nodes as the previous level, the total number of levels should be about  $O(\log n)$ . Similarly, each time we add another level of random shortcuts to the skip list, we cut the search time in half except for a constant overhead. So after  $O(\log n)$  levels, we should have a search time of  $O(\log n)$ . Let's formalize each of these two intuitive observations.

### A.3 Number of Levels

The actual values of the search keys don't affect the skip list analysis, so let's assume the keys are the integers 1 through  $n$ . Let  $L(x)$  be the number of levels of the skip list that contain some search key  $x$ , not counting the bottom level. Each new copy of  $x$  is created with probability  $1/2$  from the previous level, essentially by flipping a coin. We can compute the expected value of  $L(x)$  recursively—with probability  $1/2$ , we flip tails and  $L(x) = 0$ ; and with probability  $1/2$ , we flip heads, increase  $L(x)$  by one, and recurse:

$$E[L(x)] = \frac{1}{2} \cdot 0 + \frac{1}{2}(1 + E[L(x)])$$

Solving this equation gives us  $E[L(x)] = 1$ .

In order to analyze the expected cost of a search, however, we need a bound on the number of levels  $L = \max_x L(x)$ . Unfortunately, we can't compute the average of a maximum the way we would compute the average of a sum. Instead, we will derive a stronger result, showing that the depth is  $O(\log n)$  with high probability. 'High probability' is a technical term that means the probability is at least  $1 - 1/n^c$  for some constant  $c \geq 1$ .

In order for a search key  $x$  to appear on the  $k$ th level, we must have flipped  $k$  heads in a row, so  $\Pr[L(x) \geq k] = 2^{-k}$ . In particular,

$$\Pr[L(x) \geq 2 \lg n] = \frac{1}{n^2}.$$

(There's nothing special about the number 2 here.) The skip list has at least  $2 \lg n$  levels if and only if  $L(x) \geq 2 \lg n$  for at least one of the  $n$  search keys.

$$\Pr[L \geq 2 \lg n] = \Pr[(L(1) \geq 2 \lg n) \vee (L(2) \geq 2 \lg n) \vee \dots \vee (L(n) \geq 2 \lg n)]$$

Since  $\Pr[A \vee B] \leq \Pr[A] + \Pr[B]$  for any random events  $A$  and  $B$ , we can simplify this as follows:

$$\Pr[L \geq 2 \lg n] \leq \sum_{x=1}^n \Pr[L(x) \geq 2 \lg n] = \sum_{x=1}^n \frac{1}{n^2} = \frac{1}{n}.$$

So with high probability, a skip list has  $O(\log n)$  levels.

### A.4 Logarithmic Search Time

It's a little easier to analyze the cost of a search if we imagine running the algorithm backwards. `UPLEFTSEARCH` takes the output from `SKIPLISTFIND` as input and traces back through the data structure to the upper left corner. Skip lists don't really have up and left pointers, but we'll pretend that they do so we don't have to write '`(v)up`' or '`(v)left`'.<sup>2</sup>

$\overline{\text{UPLEFTSEARCH}(v):}$ while ( $v \neq L$ ) if <code>up</code> ( $v$ ) exists $v \leftarrow \text{up}(v)$ else $v \leftarrow \text{left}(v)$
---

<sup>2</sup> Leonardo da Vinci used to write everything this way, but not because he wanted to keep his discoveries secret. He just had really bad arthritis in his hand!

Now for *every* node  $v$  in the skip list,  $\text{up}(v)$  exists with probability  $1/2$ . So for purposes of analysis,  $\text{FIND}$  is equivalent to the following algorithm:

<pre>FLIPWALK(<math>v</math>):   while (<math>v \neq L</math>)     if COINFLIP = HEADS       <math>v \leftarrow \text{up}(v)</math>     else       <math>v \leftarrow \text{left}(v)</math></pre>
---

Obviously, the expected number of heads is exactly the same as the expected number of tails. Thus, the expected running time of this algorithm is twice the expected number of upward jumps. Since we already know that the number of upward jumps is  $O(\log n)$  with high probability, we can conclude that the expected search time is  $O(\log n)$ .

*The goode workes that men don whil they ben in good lif al amortised by synne folwyng.*

— Geoffrey Chaucer, “The Persones [Parson’s] Tale” (c.1400)

*I will gladly pay you Tuesday for a hamburger today.*

— J. Wellington Wimpy, “Thimble Theatre” (1931)

*I want my two dollars!*

— Johnny Gasparini [Demian Slade], “Better Off Dead” (1985)

## 7 Amortized Analysis (October 3)

### 7.1 Incrementing a Binary Counter

One of the questions in Homework Zero asked you to prove that any number could be written in binary. Although some of you (correctly) proved this using strong induction—pulling off either the least significant bit or the most significant bit and letting the recursion fairy convert the remainder—the most common proof used weak induction as follows:

**Proof:** *Base case:*  $1 = 2^0$ .

*Inductive step:* Suppose we have a set of distinct powers of two whose sum is  $n$ . If we add  $2^0$  to this set, we get a ‘set’ of powers of two whose sum is  $n + 1$ , but there might be two copies of  $2^0$ . To fix this, as long as there are two copies of any  $2^i$ , delete them both and add  $2^{i+1}$ . The value of the sum is unchanged by this process, since  $2^{i+1} = 2^i + 2^i$ . Since each iteration decreases the number of powers of two in our ‘set’, this process must eventually terminate. At the end of this process, we have a set of distinct powers of two whose sum is  $n + 1$ .  $\square$

Here’s a more formal (and shorter!) description of the algorithm to add one to a binary numeral. The input  $B$  is an array of bits, where  $B[i] = 1$  if and only if  $2^i$  appears in the sum.

<pre style="margin: 0;">INCREMENT(<math>B</math>):   <math>i \leftarrow 0</math>   while <math>B[i] = 1</math>     <math>B[i] \leftarrow 0</math>     <math>i \leftarrow i + 1</math>   <math>B[i] \leftarrow 1</math></pre>
--

We’ve already argued that INCREMENT must terminate, but how quickly? Obviously, the running time depends on the array of bits passed as input. If the first  $k$  bits are all 1s, then INCREMENT takes  $\Theta(k)$  time. Thus, if the number represented by  $B$  is between 0 and  $n$ , INCREMENT takes  $\Theta(\log n)$  time in the worst case, since the binary representation for  $n$  is exactly  $\lfloor \lg n \rfloor + 1$  bits long.

### 7.2 Counting from 0 to $n$ : The Aggregate Method

Now suppose we want to use INCREMENT to count from 0 to  $n$ . If we only use the worst-case running time for each call, we get an upper bound of  $O(n \log n)$  on the total running time. Although this bound is correct, it isn’t the best we can do. The easiest way to get a tighter bound is to observe that we don’t need to flip  $\Theta(\log n)$  bits *every* time INCREMENT is called. The least significant bit  $B[0]$  does flip every time, but  $B[1]$  only flips every other time,  $B[2]$  flips every 4th time, and in general,  $B[i]$  flips every  $2^i$ th time. If we start from an array full of zeros, a sequence of  $n$

INCREMENTs flips each bit  $B[i]$  exactly  $\lfloor n/2^i \rfloor$  times. Thus, the total number of bit-flips for the entire sequence is

$$\sum_{i=0}^{\lfloor \lg n \rfloor} \left\lfloor \frac{n}{2^i} \right\rfloor < \sum_{i=0}^{\infty} \frac{n}{2^i} = 2n.$$

Thus, *on average*, each call to INCREMENT flips only two bits, and so runs in constant time.

This ‘on average’ is quite different from the averaging we did in the previous lecture. There is no probability involved; we are averaging over a sequence of operations, not the possible running times of a single operation. This averaging idea is called *amortization*—the *amortized* cost of each INCREMENT is  $O(1)$ . Amortization is a ~~slazy~~ clever trick used by accountants to average large one-time costs over long periods of time; the most common example is calculating uniform payments for a loan, even though the borrower is paying interest on less and less capital over time.

There are several different methods for deriving amortized bounds for a sequence of operations. CLR calls the technique we just used the *aggregate* method, which is just a fancy way of saying sum up the total cost of the sequence and divide by the number of operations.

**The Aggregate Method.** *Find the worst case running time  $T(n)$  for a sequence of  $n$  operations. The amortized cost of each operation is  $T(n)/n$ .*

### 7.3 The Taxation (Accounting) Method

The second method we can use to derive amortized bounds is called the *accounting* method in CLR, but a better name for it might be the *taxation* method. Suppose it costs us a dollar to toggle a bit, so we can measure the running time of our algorithm in dollars. Time is money!

Instead of paying for each bit flip when it happens, the Increment Revenue Service charges a two-dollar *increment tax* whenever we want to set a bit from zero to one. One of those dollars is spent changing the bit from zero to one; the other is stored away as *credit* until we need to reset the same bit to zero. The key point here is that we always have enough credit saved up to pay for the next INCREMENT. The amortized cost of an INCREMENT is the total tax it incurs, which is exactly 2 dollars, since each INCREMENT changes just one bit from 0 to 1.

It is often useful to assign various parts of the tax income to specific pieces of the data structure. For example, for each INCREMENT, we could store one of the two dollars on the single bit that is set for 0 to 1, so that *that* bit can pay to reset itself back to zero later on.

**Taxation Method 1.** *Certain steps in the algorithm charge you taxes, so that the total money it spends is never more than the total taxes you pay. The amortized cost of an operation is the overall tax charged to you during that operation.*

Perhaps a more optimistic way of looking at the taxation method is to have the bits in the array pay *us* a tax for the privilege of being updated at the proper time. Regardless of whether we change the bit or not, we charge each bit  $B[i]$  a tax of  $1/2^i$  dollars for each INCREMENT. The total tax we collect is  $\sum_{i \geq 0} 2^{-i} = 2$  dollars. Every time  $B[i]$  actually needs to be flipped, it has paid us a total of \$1 since the last change, which is just enough for us to pay for the flip.

**Taxation Method 2.** *Charge taxes to certain items in the data structure at each operation, so that the total money you spend is never more than the total taxes you collect. The amortized cost of an operation is the overall tax you collect during that operation.*

In both of the taxation methods, our task as algorithm analysts is to come up with an appropriate ‘tax schedule’. Different ‘schedules’ can result in different amortized time bounds. The tightest bounds are obtained from tax schedules that *just barely* stay in the black.

## 7.4 The Potential Method

The most powerful method (and the hardest to use) builds on a physics metaphor of ‘potential energy’. Instead of associating costs or taxes with particular operations or pieces of the data structure, we represent prepaid work as *potential* that can be spent on later operations. The potential is a function of the entire data structure.

Let  $D_i$  denote our data structure after  $i$  operations, and let  $\Phi_i$  denote its potential. Let  $c_i$  denote the actual cost of the  $i$ th operation (which changes  $D_{i-1}$  into  $D_i$ ). Then the *amortized* cost of the  $i$ th operation, denoted  $a_i$ , is defined to be the actual cost plus the change in potential:

$$a_i = c_i + \Phi_i - \Phi_{i-1}$$

So the *total* amortized cost of  $n$  operations is the actual total cost plus the total change in potential:

$$\sum_{i=1}^n a_i = \sum_{i=1}^n (c_i + \Phi_i - \Phi_{i-1}) = \sum_{i=1}^n c_i + \Phi_n - \Phi_0.$$

Our task is to define a potential function so that  $\Phi_0 = 0$  and  $\Phi_i \geq 0$  for all  $i$ . Once we do this, the total *actual* cost of any sequence of operations will be less than the total *amortized* cost:

$$\sum_{i=1}^n c_i = \sum_{i=1}^n a_i - \Phi_n \leq \sum_{i=1}^n a_i.$$

For our binary counter example, we can define the potential  $\Phi_i$  after the  $i$ th INCREMENT to be the number of bits with value 1. Initially, all bits are equal to zero, so  $\Phi_0 = 0$ , and clearly  $\Phi_i > 0$  for all  $i > 0$ , so this is a legal potential function. We can describe both the actual cost of an INCREMENT and the change in potential in terms of the number of bits set to 1 and reset to 0.

$$\begin{aligned} c_i &= \text{\#bits changed from 0 to 1} + \text{\#bits changed from 1 to 0} \\ \Phi_i - \Phi_{i-1} &= \text{\#bits changed from 0 to 1} - \text{\#bits changed from 1 to 0} \end{aligned}$$

Thus, the amortized cost of the  $i$ th INCREMENT is

$$a_i = c_i + \Phi_i - \Phi_{i-1} = 2 \times \text{\#bits changed from 0 to 1}$$

Since INCREMENT changes only *one* bit from 0 to 1, the amortized cost INCREMENT is 2.

**The Potential Method.** Define a potential function for the data structure that is initially equal to zero and is always nonnegative. The amortized cost of an operation is its actual cost plus the change in potential.

For this particular example, the potential is exactly equal to the total unspent taxes paid using the taxation method, so not too surprisingly, we have exactly the same amortized cost. In general, however, there may be no way of interpreting the change in potential as ‘taxes’.

Different potential functions will lead to different amortized time bounds. The trick to using the potential method is to come up with the best possible potential function. A good potential function goes up a little during any cheap/fast operation, and goes down a lot during any expensive/slow operation. Unfortunately, there is no general technique for doing this other than playing around with the data structure and trying lots of different possibilities.

## 7.5 Incrementing and Decrementing

Now suppose we wanted a binary counter that we could both increment and decrement efficiently. A standard binary counter won't work, even in an amortized sense, since alternating between  $2^k$  and  $2^k - 1$  costs  $\Theta(k)$  time per operation.

A nice alternative is represent a number as a pair of bit strings  $(P, N)$ , where for any bit position  $i$ , at most one of the bits  $P[i]$  and  $N[i]$  is equal to 1. The actual value of the counter is  $P - N$ . Here are algorithms to increment and decrement our double binary counter.

INCREMENT( $P, N$ ):	DECREMENT( $P, N$ ):
$i \leftarrow 0$	$i \leftarrow 0$
while $P[i] = 1$	while $N[i] = 1$
$P[i] \leftarrow 0$	$N[i] \leftarrow 0$
$i \leftarrow i + 1$	$i \leftarrow i + 1$
if $N[i] = 1$	if $P[i] = 1$
$N[i] \leftarrow 0$	$P[i] \leftarrow 0$
else	else
$P[i] \leftarrow 1$	$N[i] \leftarrow 1$

Here's an example of these algorithms in action. Notice that any number other than zero can be represented in multiple (in fact, infinitely many) ways.

$$\begin{array}{ccccccc}
 P = 10001 & P = 10010 & P = 10011 & P = 10000 & P = 10000 & P = 10000 & P = 10001 \\
 N = 01100 & \overset{++}{\rightarrow} N = 01100 & \overset{++}{\rightarrow} N = 01100 & \overset{++}{\rightarrow} N = 01000 & \overset{--}{\rightarrow} N = 01001 & \overset{--}{\rightarrow} N = 01010 & \overset{++}{\rightarrow} N = 01010 \\
 P - N = 5 & P - N = 6 & P - N = 7 & P - N = 8 & P - N = 7 & P - N = 6 & P - N = 7
 \end{array}$$

Incrementing and decrementing a double-binary counter.

Now suppose we start from  $(0, 0)$  and apply a sequence of  $n$  INCREMENTS and DECREMENTS. In the worst case, operation takes  $\Theta(\log n)$  time, but what is the amortized cost? We can't use the aggregate method here, since we don't know what the sequence of operations looks like.

What about the taxation method? It's not hard to prove (by induction, of course) that after either  $P[i]$  or  $N[i]$  is set to 1, there must be at least  $2^i$  operations, either INCREMENTS or DECREMENTS, before that bit is reset to 0. So if each bit  $P[i]$  and  $N[i]$  pays a tax of  $2^{-i}$  at each operation, we will always have enough money to pay for the next operation. Thus, the amortized cost of each operation is at most  $\sum_{i \geq 0} 2(\cdot 2^{-i}) = 4$ .

We can get even better bounds using the potential method. Define the potential  $\Phi_i$  to be the number of 1-bits in both  $P$  and  $N$  after  $i$  operations. Just as before, we have

$$\begin{aligned}
 c_i &= \text{\#bits changed from 0 to 1} + \text{\#bits changed from 1 to 0} \\
 \Phi_i - \Phi_{i-1} &= \text{\#bits changed from 0 to 1} - \text{\#bits changed from 1 to 0} \\
 \implies a_i &= 2 \times \text{\#bits changed from 0 to 1}
 \end{aligned}$$

Since each operation changes *at most* one bit to 1, the  $i$ th operation has amortized cost  $a_i \leq 2$ .

**Exercise:** Modify the binary double-counter to support a new operation SIGN, which determines whether the number being stored is positive, negative, or zero, in constant time. The amortized time to increment or decrement the counter should still be a constant. [Hint: If  $P$  has  $p$  significant bits, and  $N$  has  $n$  significant bits, then  $p - n$  always has the same sign as  $P - N$ . For example, if  $P = 17 = 10001_2$  and  $N = 0$ , then  $p = 5$  and  $n = 0$ . But how do you store  $p$  and  $n$ ??]



**Exercise:** Suppose instead of powers of two, we represent integers as the sum of Fibonacci numbers. In other words, instead of an array of bits, we keep an array of *fits*, where the  $i$ th least significant fit indicates whether the sum includes the  $i$ th Fibonacci number  $F_i$ . For example, the fitstring 101110<sub>F</sub> represents the number  $F_6 + F_4 + F_3 + F_2 = 8 + 3 + 2 + 1 = 14$ . Describe algorithms to increment and decrement a single fitstring in constant amortized time. [Hint: Most numbers can be represented by more than one fitstring!]

## 7.6 Aside: Gray Codes

An attractive alternate solution to the increment/decrement problem was independently suggested by several students. *Gray codes* (named after Frank Gray, who discovered them in the 1950s) are methods for representing numbers as bit strings so that successive numbers differ by only one bit. For example, here is the four-bit *binary reflected* Gray code for the integers 0 through 15:

0000, 0001, 0011, 0010, 0110, 0111, 0101, 0100, 1100, 1101, 1111, 1110, 1010, 1011, 1001, 1000

The general rule for incrementing a binary reflected Gray code is to invert the bit that would be set from 0 to 1 by a normal binary counter. In terms of bit-flips, this is the perfect solution; each increment or decrement *by definition* changes only one bit. Unfortunately, it appears that *finding* the single bit to flip still requires  $\Theta(\log n)$  time in the worst case, so the total cost of maintaining a Gray code is actually the same as that of maintaining a normal binary counter.

Actually, this is only true of the naïve algorithm. The following algorithm, discovered by Gideon Ehrlich<sup>1</sup> in 1973, maintains a Gray code counter in constant *worst-case* time per increment! The algorithm uses a separate ‘focus’ array  $F[0..n]$  in addition to a Gray-code bit array  $G[0..n-1]$ .

<pre> EHRlichGRAYINIT(n):   for i ← 0 to n - 1     G[i] ← 0   for i ← 0 to n     F[i] ← i </pre>
--

<pre> EHRlichGRAYINCREMENT(n):   j ← F[0]   F[0] ← 0   if j = n     G[n - 1] ← 1 - G[n - 1]   else     G[j] = 1 - G[j]     F[j] ← F[j + 1]     F[j + 1] ← j + 1 </pre>
--

The EHRlichGRAYINCREMENT algorithm obviously runs in  $O(1)$  time, even in the worst case. Here’s the algorithm in action with  $n = 4$ . The first line is the Gray bit-vector  $G$ , and the second line shows the focus vector  $F$ , both in reverse order:

$G$  : 0000, 0001, 0011, 0010, 0110, 0111, 0101, 0100, 1100, 1101, 1111, 1110, 1010, 1011, 1001, 1000  
 $F$  : 3210, 3211, 3220, 3212, 3310, 3311, 3230, 3213, 4210, 4211, 4220, 4212, 3410, 3411, 3240, 3214

Voodoo! I won’t explain in detail how Ehrlich’s algorithm works, except to point out the following invariant. Let  $B[i]$  denote the  $i$ th bit in the *standard* binary representation of the current number. **If  $B[j] = 0$  and  $B[j - 1] = 1$ , then  $F[j]$  is the smallest integer  $k > j$  such that  $B[k] = 1$ ; otherwise,  $F[j] = j$ .** Got that?

But wait — this algorithm only handles increments; what if we also want to decrement? Sorry, I don’t have a clue. Extra credit, anyone?

<sup>1</sup>G. Ehrlich. Loopless algorithms for generating permutations, combinations, and other combinatorial configurations. *J. Assoc. Comput. Mach.* 20:500–513, 1973.

*Everything was balanced before the computers went off line. Try and adjust something, and you unbalance something else. Try and adjust that, you unbalance two more and before you know what's happened, the ship is out of control.*

— Blake, *Blake's 7*, “Breakdown” (March 6, 1978)

*A good scapegoat is nearly as welcome as a solution to the problem.*

— Anonymous

*Let's play.*

— El Mariachi [Antonio Banderas], *Desperado* (1992)

## 8 Dynamic Binary Search Trees (February 8)

### 8.1 Definitions

I'll assume that everyone is already familiar with the standard terminology for binary search trees—node, search key, edge, root, internal node, leaf, right child, left child, parent, descendant, sibling, ancestor, subtree, preorder, postorder, inorder, etc.—as well as the standard algorithms for searching for a node, inserting a node, or deleting a node. Otherwise, see Chapter 12 of CLRS.

For this lecture, we will consider only *full* binary trees—where every internal node has *exactly* two children—where only the *internal* nodes actually store search keys. In practice, we can represent the leaves with null pointers.

Recall that the *depth*  $d(v)$  of a node  $v$  is its distance from the root, and its *height*  $h(v)$  is the distance to the farthest leaf in its subtree. The height (or depth) of the tree is just the height of the root. The *size*  $|v|$  of  $v$  is the number of nodes in its subtree. The size of the whole tree is just the total number of nodes, which I'll usually denote by  $n$ .

A tree with height  $h$  has at most  $2^h$  leaves, so the minimum height of an  $n$ -leaf binary tree is  $\lceil \lg n \rceil$ . In the worst case, the time required for a search, insertion, or deletion to the height of the tree, so in general we would like keep the height as close to  $\lg n$  as possible. The best we can possibly do is to have a *perfectly balanced* tree, in which each subtree has as close to half the leaves as possible, and both subtrees are perfectly balanced. The height of a perfectly balanced tree is  $\lceil \lg n \rceil$ , so the worst-case search time is  $O(\log n)$ . However, even if we started with a perfectly balanced tree, a malicious sequence of insertions and/or deletions could make the tree arbitrarily unbalanced, driving the search time up to  $\Theta(n)$ .

To avoid this problem, we need to periodically modify the tree to maintain ‘balance’. There are several methods for doing this, and depending on the method we use, the search tree is given a different name. Examples include AVL trees, red-black trees, height-balanced trees, weight-balanced trees, bounded-balance trees, path-balanced trees,  $B$ -trees, treaps, randomized binary search trees, skip lists,<sup>1</sup> and jumplists.<sup>2</sup> Some of these trees support searches, insertions, and deletions, in  $O(\log n)$  *worst-case* time, others in  $O(\log n)$  *amortized* time, still others in  $O(\log n)$  *expected* time.

In this lecture, I'll discuss two binary search tree data structures with good *amortized* performance. The first is the *scapegoat tree*, discovered by Arne Andersson in 1989 and independently by

<sup>1</sup>Yeah, yeah. Skip lists aren't really binary search trees. Whatever you say, Mr. Picky.

<sup>2</sup>These are essentially randomized variants of the Phobian binary search trees you saw in the first midterm! [H. Brönnimann, F. Cazals, and M. Durand. Randomized jumplists: A jump-and-walk dictionary data structure. Manuscript, 2002. <http://photon.poly.edu/~hbr/publi/jumplist.html>.] So now you know who to blame.

Igal Galperin and Ron Rivest in 1993.<sup>3</sup> The second is the *splay tree*, discovered by Danny Sleator and Bob Tarjan in 1985.<sup>4</sup>

## 8.2 Lazy Deletions: Global Rebuilding

First let's consider the simple case where we start with a perfectly-balanced tree, and we only want to perform searches and deletions. To get good search and delete times, we will use a technique called *global rebuilding*. When we get a delete request, we locate and mark the node to be deleted, *but we don't actually delete it*. This requires a simple modification to our search algorithm—we still use marked nodes to guide searches, but if we search for a marked node, the search routine says it isn't there. This keeps the tree more or less balanced, but now the search time is no longer a function of the amount of data currently stored in the tree. To remedy this, we also keep track of how many nodes have been marked, and then apply the following rule:

**Global Rebuilding Rule.** *As soon as half the nodes in the tree have been marked, rebuild a new perfectly balanced tree containing only the unmarked nodes.*<sup>5</sup>

With this rule in place, a search takes  $O(\log n)$  time in the worst case, where  $n$  is the number of unmarked nodes. Specifically, since the tree has at most  $n$  marked nodes, or  $2n$  nodes altogether, we need to examine at most  $\lg n + 1$  keys. There are several methods for rebuilding the tree in  $O(n)$  time, where  $n$  is the size of the new tree. (Homework!) So a single deletion can cost  $\Theta(n)$  time in the worst case, but only if we have to rebuild; most deletions take only  $O(\log n)$  time.

We spend  $O(n)$  time rebuilding, but only after  $\Omega(n)$  deletions, so the *amortized* cost of rebuilding the tree is  $O(1)$  per deletion. (Here I'm using a simple version of the 'taxation method'. For each deletion, we charge a \$1 tax; after  $n$  deletions, we've collected \$ $n$ , which is just enough to pay for rebalancing the tree containing the remaining  $n$  nodes.) Since we also have to find and mark the node being 'deleted', the total amortized time for a deletion is  $O(\log n)$ .

## 8.3 Insertions: Partial Rebuilding

Now suppose we only want to support searches and insertions. We can't 'not really insert' new nodes into the tree, since that would make them unavailable to the search algorithm.<sup>6</sup> So instead, we'll use another method called *partial rebuilding*. We will insert new nodes normally, but whenever a *subtree* becomes unbalanced enough, we rebuild it. The definition of 'unbalanced enough' depends on an arbitrary constant  $\alpha > 1$ .

Each node  $v$  will now also store its height  $h(v)$  and the size of its subtree  $|v|$ . We now modify our insertion algorithm with the following rule:

<sup>3</sup>A. Andersson. General balanced trees. *J. Algorithms* 30:1-28, 1999. I. Galperin and R. L. Rivest. Scapegoat trees. *Proc. SODA 1993*, pp. 165–174, 1993. The claim of independence is Andersson's; the conference version of his paper appeared in 1989. The two papers actually describe very slightly different rebalancing algorithms. The algorithm I'm using here is closer to Andersson's, but my analysis is closer to Galperin and Rivest's.

<sup>4</sup>D. D. Sleator and R. Tarjan. Self-adjusting binary search trees. *J. ACM* 32:652–686, 1985.

<sup>5</sup>Alternately: When the number of unmarked nodes is one less than an exact power of two, rebuild the tree. This rule ensures that the tree is always *exactly* balanced.

<sup>6</sup>Actually, there is another dynamic data structure technique, first described by Jon Bentley and James Saxe in 1980, that *doesn't* really insert the new node! Instead, their algorithm puts the new node into a brand new data structure all by itself. Then as long as there are two trees of exactly the same size, those two trees are merged into a new tree. So this is exactly like a binary counter—instead of one  $n$ -node tree, you have a collection of  $2^i$ -nodes trees of distinct sizes. The amortized cost of inserting a new element is  $O(\log n)$ . Unfortunately, we have to look through up to  $\lg n$  trees to find anything, so the search time goes up to  $O(\log^2 n)$ . [J. L. Bentley and J. B. Saxe. Decomposable searching problems I: Static-to-dynamic transformation. *J. Algorithms* 1:301-358, 1980.] See also Homework 3.

**Partial Rebuilding Rule.** *After we insert a node, walk back up the tree updating the heights and sizes of the nodes on the search path. If we encounter a node  $v$  where  $h(v) > \alpha \cdot \lg|v|$ , rebuild its subtree into a perfectly balanced tree (in  $O(|v|)$  time).*

If we always follow this rule, then after an insertion, the height of the tree is at most  $\alpha \cdot \lg n$ . Thus, since  $\alpha$  is a constant, the worst-case search time is  $O(\log n)$ . In the worst case, insertions require  $\Theta(n)$  time—we might have to rebuild the entire tree. However, the *amortized* time for each insertion is again only  $O(\log n)$ . Not surprisingly, the proof is a little bit more complicated than for deletions.

Define the *imbalance*  $I(v)$  of a node  $v$  to be one less than the absolute difference between the sizes of its two subtrees, or zero, whichever is larger:

$$I(v) = \max \{ ||\text{left}(v)| - |\text{right}(v)|| - 1, 0 \}$$

A simple induction proof implies that  $I(v) = 0$  for every node  $v$  in a perfectly balanced tree. So immediately after we rebuild the subtree of  $v$ , we have  $I(v) = 0$ . On the other hand, each insertion into the subtree of  $v$  increments either  $|\text{left}(v)|$  or  $|\text{right}(v)|$ , so  $I(v)$  changes by at most 1.

The whole analysis boils down to the following lemma.

**Lemma 1.** *Just before we rebuild  $v$ 's subtree,  $I(v) = \Omega(|v|)$ .*

Before we prove this, let's first look at what it implies. If  $I(v) = \Omega(|v|)$ , then  $\Omega(|v|)$  keys have been inserted in the  $v$ 's subtree since the last time it was rebuilt from scratch. On the other hand, rebuilding the subtree requires  $O(|v|)$  time. Thus, if we amortize the rebuilding cost across all the insertions since the last rebuilding,  $v$  is charged *constant* time for each insertion into its subtree. Since each new key is inserted into at most  $\alpha \cdot \lg n = O(\log n)$  subtrees, the total amortized cost of an insertion is  $O(\log n)$ .

**Proof:** Since we're about to rebuild the subtree at  $v$ , we must have  $h(v) > \alpha \cdot \lg |v|$ . Without loss of generality, suppose that the node we just inserted went into  $v$ 's left subtree. Either we just rebuilt this subtree or we didn't have to, so we also have  $h(\text{left}(v)) \leq \alpha \cdot \lg |\text{left}(v)|$ . Combining these two inequalities with the recursive definition of height, we get

$$\alpha \cdot \lg |v| < h(v) \leq h(\text{left}(v)) + 1 \leq \alpha \cdot \lg |\text{left}(v)| + 1.$$

After some algebra, this simplifies to  $|\text{left}(v)| > |v|/2^{1/\alpha}$ . Combining this with the identity  $|v| = |\text{left}(v)| + |\text{right}(v)| + 1$  and doing some more algebra gives us the inequality

$$|\text{right}(v)| < (1 - 1/2^{1/\alpha}) |v| - 1.$$

Finally, we combine these two inequalities using the recursive definition of imbalance.

$$I(v) \geq |\text{left}(v)| - |\text{right}(v)| - 1 > (2/2^{1/\alpha} - 1)|v|$$

Since  $\alpha$  is a constant bigger than 1, the factor in parentheses is a positive constant. □

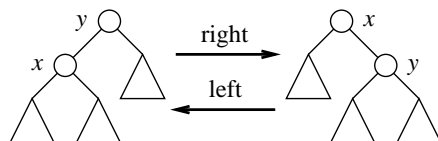
## 8.4 Scapegoat Trees

Finally, to handle both insertions and deletions efficiently, *scapegoat trees* use both of the previous techniques. We use partial rebuilding to re-balance the tree after insertions, and global rebuilding to re-balance the tree after deletions. Each search takes  $O(\log n)$  time in the worst case, and the amortized time for any insertion or deletion is also  $O(\log n)$ . There are a few small technical details left (which I won't describe), but no new ideas are required.

Once we've done the analysis, we can actually simplify the data structure. It's not hard to prove that at most one subtree (the *scapegoat*) is rebuilt during any insertion. Less obviously, we can even get the same amortized time bounds (except for a small constant factor) if we only maintain the three integers in addition to the actual tree: the size of the entire tree, the height of the entire tree, and the number of marked nodes. Whenever an insertion causes the tree to become unbalanced, we can compute the sizes of all the subtrees on the search path, starting at the new leaf and stopping at the scapegoat, in time proportional to the size of the scapegoat subtree. Since we need that much time to re-balance the scapegoat subtree, this computation increases the running time by only a small constant factor! Thus, unlike almost every other kind of balanced trees, scapegoat trees require only  $O(1)$  extra space.

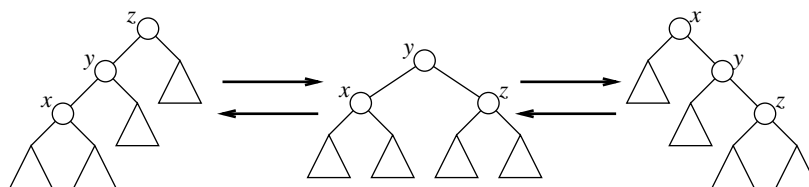
## 8.5 Rotations, Double Rotations, and Splaying

Another method for maintaining balance in binary search trees is by adjusting the shape of the tree locally, using an operation called a *rotation*. A rotation at a node  $x$  decreases its depth by one and increases its parent's depth by one. Rotations can be performed in constant time, since they only involve simple pointer manipulation.

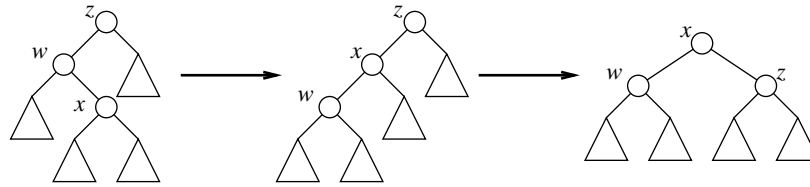


A right rotation at  $x$  and a left rotation at  $y$  are inverses.

For technical reasons, we will need to use rotations two at a time. There are two types of double rotations, which might be called *roller-coaster* and *zig-zag*. A roller-coaster at a node  $x$  consists of a rotation at  $x$ 's parent followed by a rotation at  $x$ , both in the same direction. A zig-zag at  $x$  consists of two rotations at  $x$ , in opposite directions. Each double rotation decreases the depth of  $x$  by two, leaves the depth of its parent unchanged, and increases the depth of its grandparent by either one or two, depending on the type of double rotation. Either type of double rotation can be performed in constant time.

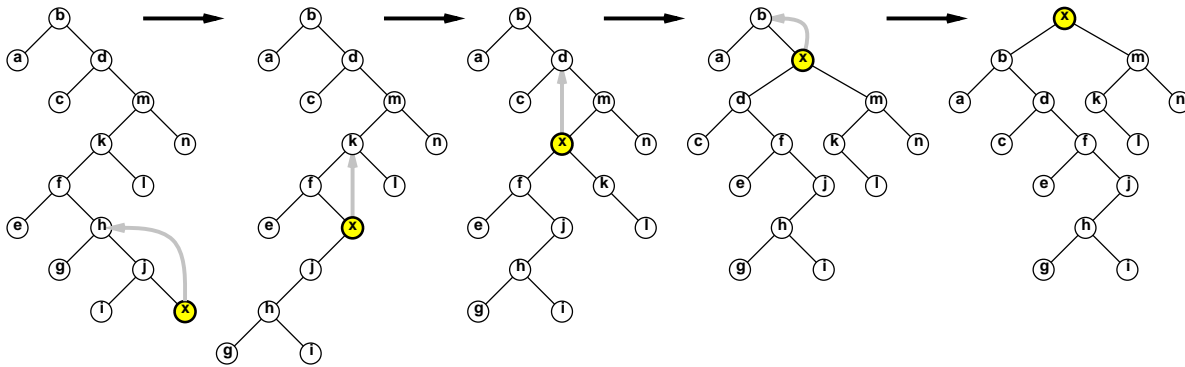


A right roller-coaster at  $x$  and a left roller-coaster at  $z$ .



A zig-zag at  $x$ . The symmetric case is not shown.

Finally, a *splay* operation moves an arbitrary node in the tree up to the root through a series of double rotations, possibly with one single rotation at the end. Splaying a node  $v$  requires time proportional to  $d(v)$ . (Obviously, this means the depth *before* splaying, since after splaying  $v$  is the root and thus has depth zero!)

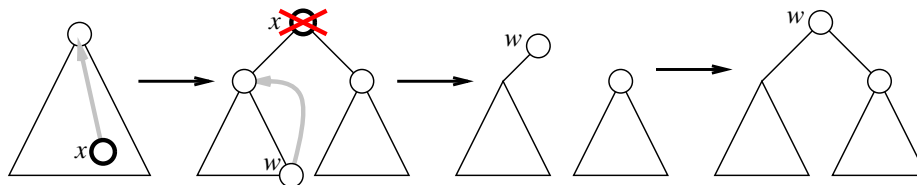


Splaying a node. Irrelevant subtrees are omitted for clarity.

### 8.6 Splay Trees

A *splay tree* is a binary search tree that is kept more or less balanced by splaying. Intuitively, after we access any node, we move it to the root with a splay operation. In more detail:

- **Search:** Find the node containing the key using the usual algorithm, or its predecessor or successor if the key is not present. Splay whichever node was found.
- **Insert:** Insert a new node using the usual algorithm, then splay that node.
- **Delete:** Find the node  $x$  to be deleted, splay it, and then delete it. This splits the tree into two subtrees, one with keys less than  $x$ , the other with keys bigger than  $x$ . Find the node  $w$  in the left subtree with the largest key (*i.e.*, the inorder predecessor of  $x$  in the original tree), splay it, and finally join it to the right subtree.



Deleting a node in a splay tree.

Each search, insertion, or deletion consists of a constant number of operations of the form *walk down to a node, and then splay it up to the root*. Since the walk down is clearly cheaper than the splay up, all we need to get good amortized bounds for splay trees is to derive good amortized bounds for a single splay.

Believe it or not, the easiest way to do this uses the potential method. The *rank* of any node  $v$  is defined as  $r(v) = \lfloor \lg |v| \rfloor$ . In particular, if  $v$  is the root of the tree, then  $r(v) = \lfloor \lg n \rfloor$ . We define the *potential* of a splay tree to be the sum of the ranks of all its nodes:

$$\Phi = \sum_v r(v) = \sum_v \lfloor \lg |v| \rfloor$$

The amortized analysis of splay trees boils down to the following lemma. Here,  $r(v)$  denotes the rank of  $v$  before a (single or double) rotation, and  $r'(v)$  denotes its rank afterwards.

**Lemma 2.** *The amortized cost of a single rotation at  $v$  is at most  $1 + 3r'(v) - 3r(v)$ , and the amortized cost of a double rotation at  $v$  is at most  $3r'(v) - 3r(v)$ .*

Proving this lemma requires an easy but boring case analysis of the different types of rotations, which takes up almost a page in Sleator and Tarjan's original paper. (They call it the 'Access Lemma'.) I won't repeat it here.

By adding up the amortized costs of all the rotations, we find that the total amortized cost of splaying a node  $v$  is at most  $1 + 3r'(v) - 3r(v)$ , where  $r'(v)$  is the rank of  $v$  after the splay. (The intermediate ranks cancel out in a nice telescoping sum.) But after the splay,  $v$  is the root, so  $r'(v) = \lfloor \lg n \rfloor$ , which means that the amortized cost of a splay is at most  $3 \lfloor \lg n \rfloor - 1 = O(\log n)$ . Thus, every insertion, deletion, or search in a splay tree takes  $O(\log n)$  amortized time, which is optimal.

Actually, splay trees are optimal in a much stronger sense. If  $p(v)$  denotes the *probability* of searching for a node  $v$ , then the amortized search time for  $v$  is  $O(\log(1/p(v)))$ . If every node is equally likely, we have  $p(v) = 1/n$  so  $O(\log(1/p(v))) = O(\log n)$ , as before. Even if the nodes aren't equally likely, though, the *optimal static* binary search tree for this set of *weighted* nodes has a search time of  $\Theta(\log(1/p(v)))$ . Splay trees give us this optimal amortized weighted search time with *no* change to the data structure or the splaying algorithm.

*E pluribus unum* (Out of many, one)

— Official motto of the United States of America

**John:** *Who's your daddy? C'mon, you know who your daddy is! Who's your daddy? D'Argo, tell him who his daddy is!*

**D'Argo:** *I'm your daddy.*

— *Farscape*, "Thanks for Sharing" (June 15, 2001)

## 9 Data Structures for Disjoint Sets (October 10 and 15)

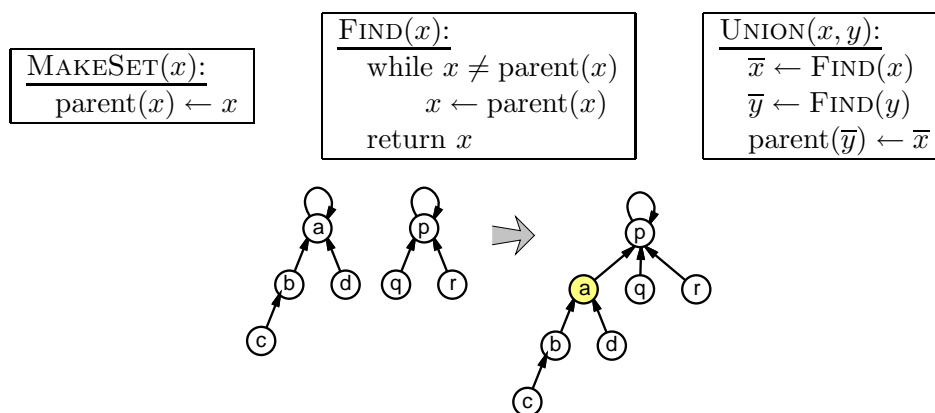
In this lecture, we describe some methods for maintaining a collection of disjoint sets. Each set is represented as a pointer-based data structure, with one node per element. Each set has a 'leader' element, which uniquely identifies the set. (Since the sets are always disjoint, the same object cannot be the leader of more than one set.) We want to support the following operations.

- **MAKESET( $x$ ):** Create a new set  $\{x\}$  containing the single element  $x$ . The element  $x$  must not appear in any other set in our collection. The leader of the new set is obviously  $x$ .
- **FIND( $x$ ):** Find (the leader of) the set containing  $x$ .
- **UNION( $A, B$ ):** Replace two sets  $A$  and  $B$  in our collection with their union  $A \cup B$ . For example, **UNION( $A, \text{MAKESET}(x)$ )** adds a new element  $x$  to an existing set  $A$ . The sets  $A$  and  $B$  are specified by arbitrary elements, so **UNION( $x, y$ )** has exactly the same behavior as **UNION(FIND( $x$ ), FIND( $y$ ))**.

Disjoint set data structures have lots of applications. For instance, Kruskal's minimum spanning tree algorithm relies on such a data structure to maintain the components of the intermediate spanning forest. Another application might be maintaining the connected components of a graph as new vertices and edges are added. In both these applications, we can use a disjoint-set data structure, where we keep a set for each connected component, containing that component's vertices.

### 9.1 Reversed Trees

One of the easiest ways to store sets is using trees. Each object points to another object, called its *parent*, except for the leader of each set, which points to itself and thus is the root of the tree. **MAKESET** is trivial. **FIND** traverses the parent pointers up to the leader. **UNION** just redirects the parent pointer of one leader to the other. Notice that unlike most tree data structures, objects do *not* have pointers down to their children.



Merging two sets stored as trees. Arrows point to parents. The shaded node has a new parent.



MAKE-SET clearly takes  $\Theta(1)$  time, and UNION requires only  $O(1)$  time in addition to the two FINDs. The running time of  $\text{FIND}(x)$  is proportional to the depth of  $x$  in the tree. It is not hard to come up with a sequence of operations that results in a tree that is a long chain of nodes, so that  $\text{FIND}$  takes  $\Theta(n)$  time in the worst case.

However, there is an easy change we can make to our UNION algorithm, called *union by depth*, so that the trees always have logarithmic depth. Whenever we need to merge two trees, we always make the root of the *shallower* tree a child of the *deeper* one. This requires us to also maintain the depth of each tree, but this is quite easy.

```

MAKESET(x):
  parent(x) ← x
  depth(x) ← 0

```

```

FIND(x):
  while x ≠ parent(x)
    x ← parent(x)
  return x

```

```

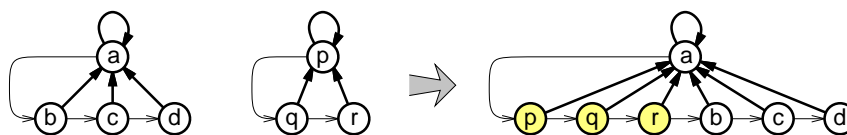
UNION(x, y)
  x̄ ← FIND(x)
  ȳ ← FIND(y)
  if depth(x̄) > depth(ȳ)
    parent(ȳ) ← x̄
  else
    parent(x̄) ← ȳ
    if depth(x̄) = depth(ȳ)
      depth(ȳ) ← depth(ȳ) + 1

```

With this simple change,  $\text{FIND}$  and  $\text{UNION}$  both run in  $\Theta(\log n)$  time in the worst case.

## 9.2 Shallow Threaded Trees

Alternately, we could just have every object keep a pointer to the leader of its set. Thus, each set is represented by a shallow tree, where the leader is the root and all the other elements are its children. With this representation,  $\text{MAKESET}$  and  $\text{FIND}$  are completely trivial. Both operations clearly run in constant time.  $\text{UNION}$  is a little more difficult, but not much. Our algorithm sets all the leader pointers in one set to point to the leader of the other set. To do this, we need a method to visit every element in a set; we will ‘thread’ a linked list through each set, starting at the set’s leader. The two threads are merged in the  $\text{UNION}$  algorithm in constant time.



Merging two sets stored as threaded trees.

Bold arrows point to leaders; lighter arrows form the threads. Shaded nodes have a new leader.

```

MAKESET(x):
  leader(x) ← x
  next(x) ← x

```

```

FIND(x):
  return leader(x)

```

```

UNION(x, y):
  x̄ ← FIND(x)
  ȳ ← FIND(y)
  y ← ȳ
  leader(y) ← x̄
  while (next(y) ≠ NULL)
    y ← next(y)
    leader(y) ← x̄
  next(y) ← next(x̄)
  next(x̄) ← ȳ

```

The worst-case running time of UNION is a constant times the size of the *larger* set. Thus, if we merge a one-element set with another  $n$ -element set, the running time can be  $\Theta(n)$ . Generalizing this idea, it is quite easy to come up with a sequence of  $n$  MAKESET and  $n - 1$  UNION operations that requires  $\Theta(n^2)$  time to create the set  $\{1, 2, \dots, n\}$  from scratch.

<p style="text-align: center;"><u>WORSTCASESEQUENCE(<math>n</math>):</u>          MAKESET(1)          for <math>i \leftarrow 2</math> to <math>n</math>              MAKESET(<math>i</math>)              UNION(1, <math>i</math>)</p>
--

We are being stupid in two different ways here. One is the order of operations in WORSTCASESEQUENCE. Obviously, it would be more efficient to merge the sets in the other order, or to use some sort of divide and conquer approach. Unfortunately, we can't fix this; we don't get to decide how our data structures are used! The other is that we always update the leader pointers in the larger set. To fix this, we add a comparison inside the UNION algorithm to determine which set is smaller. This requires us to maintain the size of each set, but that's easy.

<p style="text-align: center;"><u>MAKEWEIGHTEDSET(<math>x</math>):</u>          leader(<math>x</math>) <math>\leftarrow x</math>          next(<math>x</math>) <math>\leftarrow x</math>          size(<math>x</math>) <math>\leftarrow 1</math></p>
--

<p style="text-align: center;"><u>WEIGHTEDUNION(<math>x, y</math>)</u>  <math>\bar{x} \leftarrow \text{FIND}(x)</math>  <math>\bar{y} \leftarrow \text{FIND}(y)</math>          if size(<math>\bar{x}</math>) &gt; size(<math>\bar{y}</math>)              UNION(<math>\bar{x}, \bar{y}</math>)              size(<math>\bar{x}</math>) <math>\leftarrow</math> size(<math>\bar{x}</math>) + size(<math>\bar{y}</math>)          else              UNION(<math>\bar{y}, \bar{x}</math>)              size(<math>\bar{x}</math>) <math>\leftarrow</math> size(<math>\bar{x}</math>) + size(<math>\bar{y}</math>)</p>
---

The new WEIGHTEDUNION algorithm still takes  $\Theta(n)$  time to merge two  $n$ -element sets. However, in an amortized sense, this algorithm is much more efficient. Intuitively, before we can merge two large sets, we have to perform a large number of MAKEWEIGHTEDSET operations.

**Theorem 1.** *A sequence of  $m$  MAKEWEIGHTEDSET operations and  $n$  WEIGHTEDUNION operations takes  $O(m + n \log n)$  time in the worst case.*

**Proof:** Whenever the leader of an object  $x$  is changed by a WEIGHTEDUNION, the size of the set containing  $x$  increases by at least a factor of two. By induction, if the leader of  $x$  has changed  $k$  times, the set containing  $x$  has at least  $2^k$  members. After the sequence ends, the largest set contains at most  $n$  members. (Why?) Thus, the leader of any object  $x$  has changed at most  $\lceil \lg n \rceil$  times.

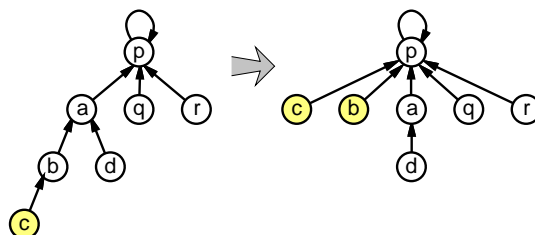
Since each WEIGHTEDUNION reduces the number of sets by one, there are  $m - n$  sets at the end of the sequence, and at most  $n$  objects are *not* in singleton sets. Since each of the non-singleton objects had  $O(\log n)$  leader changes, the total amount of work done in updating the leader pointers is  $O(n \log n)$ .  $\square$

The aggregate method now implies that each WEIGHTEDUNION has amortized cost  $O(\log n)$ .

### 9.3 Path Compression

Using unthreaded trees, FIND takes logarithmic time and everything else is constant; using threaded trees, UNION takes logarithmic amortized time and everything else is constant. A third method allows us to get both of these operations to have *almost* constant running time.

We start with the original unthreaded tree representation, where every object points to a parent. The key observation is that in any FIND operation, once we determine the leader of an object  $x$ , we can speed up future FINDs by redirecting  $x$ 's parent pointer directly to that leader. In fact, we can change the parent pointers of all the ancestors of  $x$  all the way up to the root; this is easiest if we use recursion for the initial traversal up the tree. This modification to FIND is called *path compression*.



Path compression during FIND( $c$ ). Shaded nodes have a new parent.

$\text{FIND}(x)$ if $x \neq \text{parent}(x)$ $\text{parent}(x) \leftarrow \text{FIND}(\text{parent}(x))$ return $\text{parent}(x)$
---

If we use path compression, the ‘depth’ field we used earlier to keep the trees shallow is no longer correct, and correcting it would take way too long. But this information still ensures that FIND runs in  $\Theta(\log n)$  time in the worst case, so we’ll just give it another name: *rank*.

$\text{MAKESET}(x):$ $\text{parent}(x) \leftarrow x$ $\text{rank}(x) \leftarrow 0$
---

$\text{UNION}(x, y)$ $\bar{x} \leftarrow \text{FIND}(x)$ $\bar{y} \leftarrow \text{FIND}(y)$ if $\text{rank}(\bar{x}) > \text{rank}(\bar{y})$ $\text{parent}(\bar{y}) \leftarrow \bar{x}$ else $\text{parent}(\bar{x}) \leftarrow \bar{y}$ if $\text{rank}(\bar{x}) = \text{rank}(\bar{y})$ $\text{rank}(\bar{y}) \leftarrow \text{rank}(\bar{y}) + 1$
---

Ranks have several useful properties that can be verified easily by examining the UNION and FIND algorithms. For example:

- If an object  $x$  is not a set leader, then the rank of  $x$  is strictly less than the rank of its parent.
- Whenever  $\text{parent}(x)$  changes, the new parent has larger rank than the old parent.
- The size of any set is exponential in the rank of its leader:  $\text{size}(\bar{x}) \geq 2^{\text{rank}(\bar{x})}$ . (This is easy to prove by induction hint hint.)
- In particular, since there are only  $n$  objects, the highest possible rank is  $\lfloor \lg n \rfloor$ .

We can also derive a bound on the number of nodes with a given rank  $r$ . Only set leaders can change their rank. When the rank of a set leader  $\bar{x}$  changes from  $r - 1$  to  $r$ , mark all the nodes in that set. At least  $2^r$  nodes are marked. The next time these nodes get a new leader  $\bar{y}$ , the rank of  $\bar{y}$  will be at least  $r + 1$ . Thus, any node is marked at most once. There are  $n$  nodes altogether, and any object with rank  $r$  marks  $2^r$  of them. Thus, there can be at most  $n/2^r$  objects of rank  $r$ .

Purely as an accounting tool, we will also partition the objects into several numbered *blocks*. Specifically, each object  $x$  is assigned to block number  $\lg^*(\text{rank}(x))$ . In other words,  $x$  is in block  $b$  if and only if

$$2 \uparrow\uparrow (b - 1) < \text{rank}(x) \leq 2 \uparrow\uparrow b,$$

where  $2 \uparrow\uparrow b$  is the *tower function*<sup>1</sup>

$$2 \uparrow\uparrow b = 2^{2^{2^{\dots^2}}} \Big\}^b = \begin{cases} 1 & \text{if } b = 0 \\ 2^{2 \uparrow\uparrow (b-1)} & \text{if } b > 0 \end{cases}$$

Since there are at most  $n/2^r$  objects with any rank  $r$ , the total number of objects in block  $b$  is at most

$$\sum_{r=2 \uparrow\uparrow (b-1)+1}^{2 \uparrow\uparrow b} \frac{n}{2^r} < \sum_{r=2 \uparrow\uparrow (b-1)+1}^{\infty} \frac{n}{2^r} = \frac{n}{2^{2 \uparrow\uparrow (b-1)}} = \frac{n}{2 \uparrow\uparrow b}.$$

Every object has a rank between 0 and  $\lfloor \lg n \rfloor$ , so there are  $\lg^* n$  blocks, numbered from 0 to  $\lg^* \lfloor \lg n \rfloor = \lg^* n - 1$ .

**Theorem 2.** *If we use both union-by-rank and path compression, the worst-case running time of a sequence of  $m$  operations,  $n$  of which are MAKESET operations, is  $O(m \log^* n)$ .*

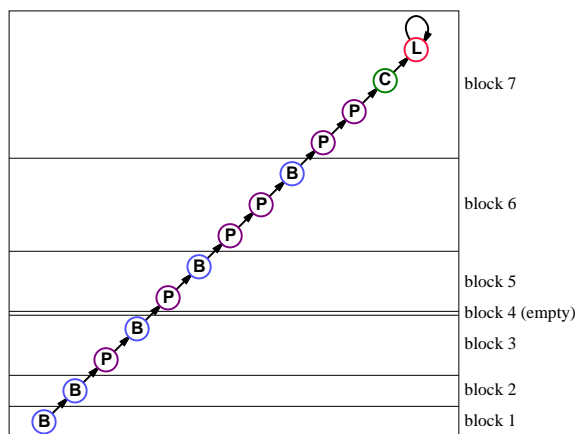
**Proof:** Since each MAKESET and UNION operation takes constant time, it suffices to show that any sequence of  $m$  FIND operations requires  $O(m \log^* n)$  time in the worst case.

The cost of  $\text{FIND}(x_0)$  is proportional to the number of nodes on the *find path* from  $x_0$  up to its leader (before path compression). To count up the total cost of all FINDs, we use an accounting method—each object  $x_0, x_1, x_2, \dots, x_l$  on the find path pays a \$1 tax into one of several different bank accounts. After all the FIND operations are done, the total amount of money in these accounts will tell us the total running time.

- The leader  $x_l$  pays into the *leader* account.
- The child of the leader  $x_{l-1}$  pays into the *child* account.
- Any other object  $x_i$  in a different block from its parent  $x_{i+1}$  pays into the *block* account.
- Any other object  $x_i$  in the same block as its parent  $x_{i+1}$  pays into the *path* account.

During any FIND operation, one dollar is paid into the leader account, at most one dollar is paid into the child account, and at most one dollar is paid into the block account for each of the  $\lg^* n$  blocks. Thus, when the sequence of  $m$  operations ends, those three accounts share a total of at most  $2m + m \lg^* n$  dollars. The only remaining difficulty is the path account.

<sup>1</sup>The arrow notation  $a \uparrow\uparrow b$  was introduced by Don Knuth in 1976.



Different nodes on the find path pay into different accounts: B=block, P=path, C=child, L=leader. Horizontal lines are boundaries between blocks. Only the nodes on the find path are shown.

So consider an object  $x_i$  in block  $b$  that pays into the path account. This object is not a set leader, so its rank can never change. The parent of  $x_i$  is also not a set leader, so after path compression,  $x_i$  acquires a new parent—namely  $x_l$ —whose rank is strictly larger than its old parent  $x_{i+1}$ . Since  $\text{rank}(\text{parent}(x))$  is always increasing, the parent of  $x_i$  must eventually lie in a different block than  $x_i$ , after which  $x_i$  will never pay into the path account. Thus,  $x_i$  can pay into the path account at most once for every rank in block  $b$ , or less than  $2 \uparrow \uparrow b$  times overall.

Since block  $b$  contains less than  $n/(2 \uparrow \uparrow b)$  objects, these objects contribute less than  $n$  dollars to the path account. There are  $\lg^* n$  blocks, so the path account receives less than  $n \lg^* n$  dollars altogether.

Thus, the total amount of money in all four accounts is less than  $2m + m \lg^* n + n \lg^* n = O(m \lg^* n)$ , and this bounds the total running time of the  $m$  FIND operations.  $\square$

The aggregate method now implies that each FIND has  $\text{amortized cost } O(\log^* n)$ , which is significantly better than its  $\text{worst-case cost } \Theta(\log n)$ .

### 9.4 Ackermann's Function and Its Inverse

But this amortized time bound can be improved even more! Just to *state* the correct time bound, I need to introduce a certain function defined by Wilhelm Ackermann in 1928. The function can be<sup>2</sup> defined by the following two-parameter recurrence.

$$A_i(n) = \begin{cases} 2 & \text{if } n = 1 \\ 2j & \text{if } i = 1 \text{ and } n > 1 \\ A_{i-1}(A_i(n-1)) & \text{otherwise} \end{cases}$$

Clearly, each  $A_i(n)$  is a monotonically increasing function of  $n$ , and these functions grow faster and faster as the index  $i$  increases— $A_2(n)$  is the power function  $2^n$ ,  $A_3(n)$  is the tower function  $2 \uparrow \uparrow n$ ,  $A_4(n)$  is the *wower* function  $2 \uparrow \uparrow \uparrow n = \underbrace{2 \uparrow \uparrow 2 \uparrow \uparrow \dots \uparrow \uparrow 2}_n$  (so named by John Conway), *et cetera ad infinitum*.

<sup>2</sup>Ackermann didn't define his function this way—I'm actually describing a different function defined 35 years later by R. Creighton Buck—but the exact details of the definition are surprisingly irrelevant!

$i$	$A_i(n)$	$n = 1$	$n = 2$	$n = 3$	$n = 4$	$n = 5$
$i = 1$	$2n$	2	4	6	8	10
$i = 2$	$2 \uparrow n$	2	4	8	16	32
$i = 3$	$2 \uparrow \uparrow n$	2	4	16	65536	$2^{65536}$
$i = 4$	$2 \uparrow \uparrow \uparrow n$	2	4	65536	$2^{2^{2^{\dots^2}}} \Big\}^{65536}$	$2^{2^{2^{\dots^2}}} \Big\}^{2^{2^{2^{\dots^2}}} \Big\}^{65536}}$
$i = 5$	$2 \uparrow \uparrow \uparrow \uparrow n$	2	4	$2^{2^{2^{\dots^2}}} \Big\}^{65536}$	$2^{2^{\dots^2}} \Big\}^{2^{\dots^2}} \Big\}^{65536}$	$2^{2^{2^{\dots^2}}} \Big\}^{2^{2^{\dots^2}}} \Big\}^{65536}$ «Yeah, right.»

Small(!!) values of Ackermann's function.

The *functional inverse* of Ackermann's function is defined as follows:

$$\alpha(m, n) = \min \{i \mid A_i(\lfloor m/n \rfloor) > \lg n\}$$

For all practical values of  $n$  and  $m$ , we have  $\alpha(m, n) \leq 4$ ; nevertheless, if we increase  $m$  and keep  $n$  fixed,  $\alpha(m, n)$  is eventually bigger than any fixed constant.

Bob Tarjan proved the following surprising theorem. The proof of the upper bound<sup>3</sup> is very similar to the proof of Theorem 2, except that it uses a more complicated 'block' structure. The proof of the matching lower bound<sup>4</sup> is, unfortunately, way beyond the scope of this class.<sup>5</sup>

**Theorem 3.** *Using both union by rank and path compression, the worst-case running time of a sequence of  $m$  operations,  $n$  of which are MAKESETS, is  $\Theta(m\alpha(m, n))$ . Thus, each operation has amortized cost  $\Theta(\alpha(m, n))$ . This time bound is optimal: any pointer-based data structure needs  $\Omega(m\alpha(m, n))$  time to perform these operations.*

<sup>3</sup>R. E. Tarjan. Efficiency of a good but not linear set union algorithm. *J. Assoc. Comput. Mach.* 22:215–225, 1975.

<sup>4</sup>R. E. Tarjan. A class of algorithms which require non-linear time to maintain disjoint sets. *J. Comput. Syst. Sci.* 19:110–127, 1979.

<sup>5</sup>But if you like this sort of thing, google for "Davenport-Schinzel sequences".

## B Fibonacci Heaps

### B.1 Mergeable Heaps

A *mergeable heap* is a data structure that stores a collection of *keys*<sup>1</sup> and supports the following operations.

- **Insert:** Insert a new key into a heap. This operation can also be used to create a new heap containing just one key.
- **FindMin:** Return the smallest key in a heap.
- **DeleteMin:** Remove the smallest key from a heap.
- **Merge:** Merge two heaps into one. The new heap contains all the keys that used to be in the old heaps, and the old heaps are (possibly) destroyed.

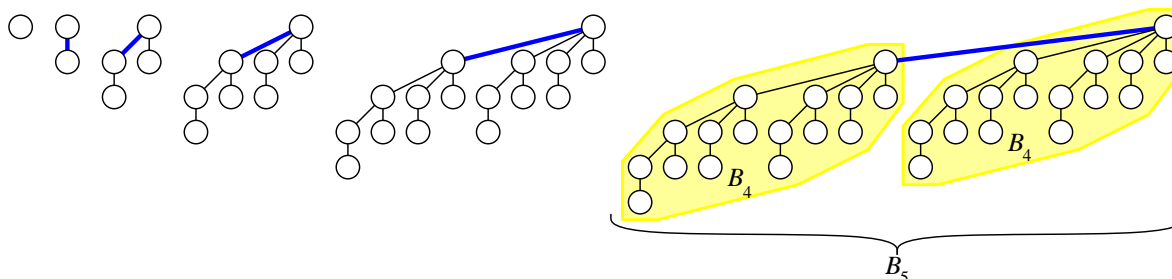
If we never had to use DELETEMIN, mergeable heaps would be completely trivial. Each “heap” just stores to maintain the single record (if any) with the smallest key. INSERTS and MERGES require only one comparison to decide which record to keep, so they take constant time. FINDMIN obviously takes constant time as well.

If we need DELETEMIN, but we don’t care how long it takes, we can still implement mergeable heaps so that INSERTS, MERGES, and FINDMINS take constant time. We store the records in a circular doubly-linked list, and keep a pointer to the minimum key. Now deleting the minimum key takes  $\Theta(n)$  time, since we have to scan the linked list to find the new smallest key.

In this lecture, I’ll describe a data structure called a *Fibonacci heap* that supports INSERTS, MERGES, and FINDMINS in constant time, even in the worst case, and also handles DELETEMIN in  $O(\log n)$  amortized time. That means that any sequence of  $n$  INSERTS,  $m$  MERGES,  $f$  FINDMINS, and  $d$  DELETEMINS takes  $O(n + m + f + d \log n)$  time.

### B.2 Binomial Trees and Fibonacci Heaps

A *Fibonacci heap* is a circular doubly linked list, with a pointer to the minimum key, but the elements of the list are not single keys. Instead, we collect keys together into structures called *binomial heaps*. Binomial heaps are trees<sup>2</sup> that satisfy the *heap property* — every node has a smaller key than its children — and have the following special structure.



Binomial trees of order 0 through 5.

<sup>1</sup>In the earlier lecture on treaps, I called these keys *priorities* to distinguish them from search keys.

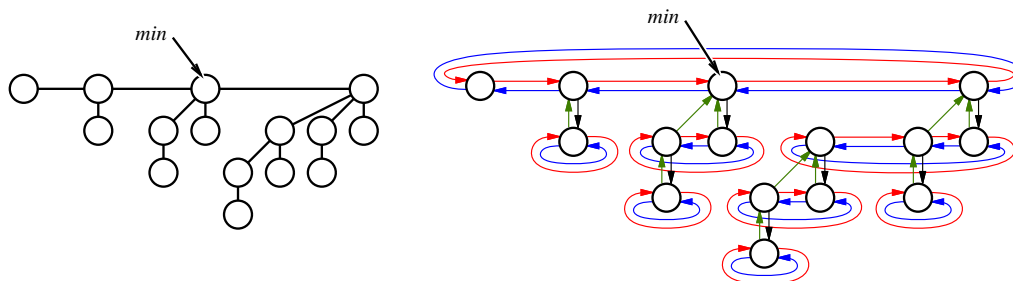
<sup>2</sup>CLR uses the name ‘binomial heap’ to describe a more complicated data structure consisting of a set of heap-ordered binomial trees, with at most one binomial tree of each order.

A  $k$ th order binomial tree, which I'll abbreviate  $B_k$ , is defined recursively.  $B_0$  is a single node. For all  $k > 0$ ,  $B_k$  consists of two copies of  $B_{k-1}$  that have been *linked* together, meaning that the root of one  $B_{k-1}$  has become a new child of the other root.

Binomial trees have several useful properties, which are easy to prove by induction (hint, hint).

- The root of  $B_k$  has degree  $k$ .
- The children of the root of  $B_k$  are the roots of  $B_0, B_1, \dots, B_{k-1}$ .
- $B_k$  has height  $k$ .
- $B_k$  has  $2^k$  nodes.
- $B_k$  can be obtained from  $B_{k-1}$  by adding a new child to every node.
- $B_k$  has  $\binom{k}{d}$  nodes at depth  $d$ , for all  $0 \leq d \leq k$ .
- $B_k$  has  $2^{k-h-1}$  nodes with height  $h$ , for all  $0 \leq h < k$ , and one node (the root) with height  $k$ .

Although we normally don't care in this class about the low-level details of data structures, we need to be specific about how Fibonacci heaps are actually implemented, so that we can be sure that certain operations can be performed quickly. Every node in a Fibonacci heap points to four other nodes: its parent, its 'next' sibling, its 'previous' sibling, and one of its children. The sibling pointers are used to join the roots together into a circular doubly-linked *root list*. In each binomial tree, the children of each node are also joined into a circular doubly-linked list using the sibling pointers.



A high-level view and a detailed view of the same Fibonacci heap. Null pointers are omitted for clarity.

With this representation, we can add or remove nodes from the root list, merge two root lists together, link one binomial tree to another, or merge a node's list of children with the root list, in constant time, and we can visit every node in the root list in constant time per node. Having established that these primitive operations can be performed quickly, we never again need to think about the low-level representation details.

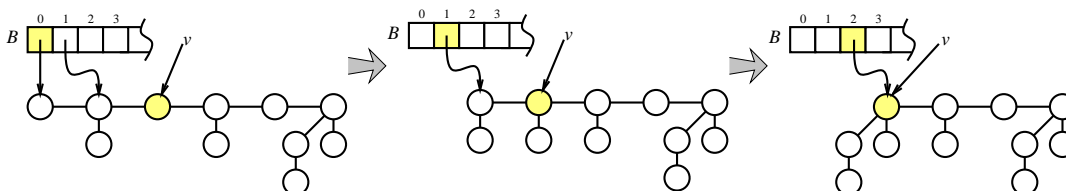
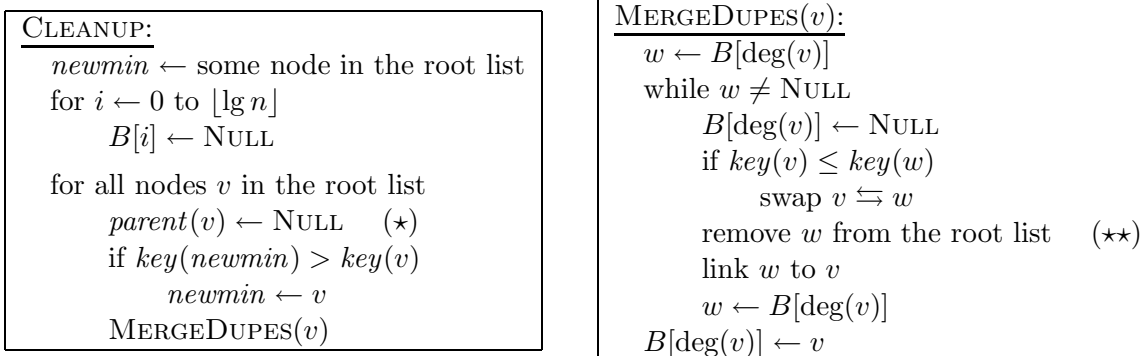
### B.3 Operations on Fibonacci Heaps

The INSERT, MERGE, and FINDMIN algorithms for Fibonacci heaps are exactly like the corresponding algorithms for linked lists. Since we maintain a pointer to the minimum key, FINDMIN is trivial. To insert a new key, we add a single node (which we should think of as a  $B_0$ ) to the root list and (if necessary) update the pointer to the minimum key. To merge two Fibonacci heaps, we just merge the two root lists and keep the pointer to the smaller of the two minimum keys. Clearly, all three operations take  $O(1)$  time.



Deleting the minimum key is a little more complicated. First, we remove the minimum key from the root list and splice its children into the root list. Except for updating the parent pointers, this takes  $O(1)$  time. Then we scan through the root list to find the new smallest key and update the parent pointers of the new roots. This scan could take  $\Theta(n)$  time in the worst case. To bring down the *amortized* deletion time, we apply a CLEANUP algorithm, which links pairs of equal-size binomial heaps until there is only one binomial heap of any particular size.

Let me describe the CLEANUP algorithm in more detail, so we can analyze its running time. The following algorithm maintains a global array  $B[1.. \lceil \lg n \rceil]$ , where  $B[i]$  is a pointer to some previously-visited binomial heap of order  $i$ , or NULL if there is no such binomial heap. Notice that CLEANUP simultaneously resets the parent pointers of all the new roots and updates the pointer to the minimum key. I've split off the part of the algorithm that merges binomial heaps of the same order into a separate subroutine MERGEDUPES.



MERGEDUPES( $v$ ), ensuring that no earlier root has the same degree as  $v$ .

Notices that MERGEDUPES is careful to merge heaps so that the heap property is maintained—the heap whose root has the larger key becomes a new child of the heap whose root has the smaller key. This is handled by swapping  $v$  and  $w$  if their keys are in the wrong order.

The running time of CLEANUP is  $O(r')$ , where  $r'$  is the length of the root list just before CLEANUP is called. The easiest way to see this is to count the number of times the two starred lines can be executed: line (★) is executed once for every node  $v$  on the root list, and line (★★) is executed *at most* once for every node  $w$  on the root list. Since DELETEMIN does only a constant amount of work before calling CLEANUP, the running time of DELETEMIN is  $O(r') = O(r + \text{deg}(\text{min}))$  where  $r$  is the number of roots before DELETEMIN begins, and min is the node deleted.

Although  $\text{deg}(\text{min})$  is at most  $\lg n$ , we can still have  $r = \Theta(n)$  (for example, if nothing has been deleted yet), so the worst-case time for a DELETEMIN is  $\Theta(n)$ . After a DELETEMIN, the root list has length  $O(\log n)$ , since all the binomial heaps have unique orders and the largest has order at most  $\lceil \lg n \rceil$ .

## B.4 Amortized Analysis of DeleteMin

To bound the amortized cost, observe that each insertion increments  $r$ . If we charge a constant ‘cleanup tax’ for each insertion, and use the collected tax to pay for the CLEANUP algorithm, the

unpaid cost of a DELETEMIN is only  $O(\deg(\min)) = O(\log n)$ .

More formally, define the *potential* of the Fibonacci heap to be the number of roots. Recall that the amortized time of an operation can be defined as its actual running time plus the increase in potential, provided the potential is initially zero (it is) and we never have negative potential (we never do). Let  $r$  be the number of roots before a DELETEMIN, and let  $r''$  denote the number of roots afterwards. The actual cost of DELETEMIN is  $r + \deg(\min)$ , and the number of roots increases by  $r'' - r$ , so the amortized cost is  $r'' + \deg(\min)$ . Since  $r'' = O(\log n)$  and the degree of any node is  $O(\log n)$ , the amortized cost of DELETEMIN is  $O(\log n)$ .

Each INSERT adds only one root, so its amortized cost is still constant. A MERGE actually doesn't change the number of roots, since the new Fibonacci heap has all the roots from its constituents and no others, so its amortized cost is also constant.

## B.5 Decreasing Keys

In some applications of heaps, we also need the ability to delete an arbitrary node. The usual way to do this is to decrease the node's key to  $-\infty$ , and then use DELETEMIN. Here I'll describe how to decrease the key of a node in a Fibonacci heap; the algorithm will take  $O(\log n)$  time in the worst case, but the amortized time will be only  $O(1)$ .

Our algorithm for decreasing the key at a node  $v$  follows two simple rules.

1. Promote  $v$  up to the root list. (This moves the whole subtree rooted at  $v$ .)
2. As soon as two children of any node  $w$  have been promoted, immediately promote  $w$ .

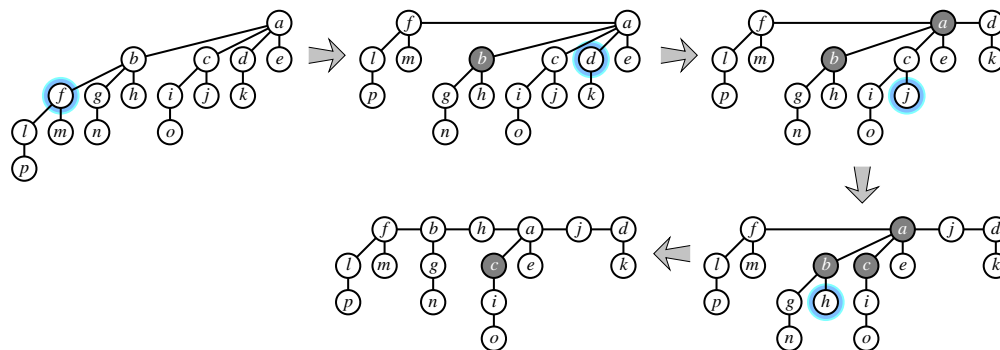
In order to enforce the second rule, we now *mark* certain nodes in the Fibonacci heap. Specifically, a node is marked if exactly one of its children has been promoted. If some child of a marked node is promoted, we promote (and unmark) that node as well. Whenever we promote a marked node, we unmark it; this is the *only* way to unmark a node. (Specifically, splicing nodes into the root list during a DELETEMIN is not considered a promotion.)

Here's a more formal description of the algorithm. The input is a pointer to a node  $v$  and the new value  $k$  for its key.

DECREASEKEY( $v, k$ ):  
 $key(v) \leftarrow k$   
 update the pointer to the smallest key  
 PROMOTE( $v$ )

PROMOTE( $v$ ):  
 unmark  $v$   
 if  $parent(v) \neq \text{NULL}$   
   remove  $v$  from  $parent(v)$ 's list of children  
   insert  $v$  into the root list  
   if  $parent(v)$  is marked  
     PROMOTE( $parent(v)$ )  
 else  
   mark  $parent(v)$

The PROMOTE algorithm calls itself recursively, resulting in a 'cascading promotion'. Each consecutive marked ancestor of  $v$  is promoted to the root list and unmarked, otherwise unchanged. The lowest unmarked ancestor is then marked, since one of its children has been promoted.



Decreasing the keys of four nodes: first  $f$ , then  $d$ , then  $j$ , and finally  $h$ . Dark nodes are marked.  $\text{DECREASEKEY}(h)$  causes nodes  $b$  and  $a$  to be recursively promoted.

The time to decrease the key of a node  $v$  is  $O(1 + \#\text{consecutive marked ancestors of } v)$ . Binomial heaps have logarithmic depth, so if we still had only full binomial heaps, the running time would be  $O(\log n)$ . Unfortunately, promoting nodes destroys the nice binomial tree structure; our trees no longer have logarithmic depth! In fact,  $\text{DECREASEKEY}$  runs in  $\Theta(n)$  time in the worst case.

To compute the amortized cost of  $\text{DECREASEKEY}$ , we'll use the potential method, just as we did for  $\text{DELETEMIN}$ . We need to find a potential function  $\Phi$  that goes up a little whenever we do a little work, and goes down a lot whenever we do a lot of work.  $\text{DECREASEKEY}$  unmarks several marked ancestors and possibly also marks one node. So *the number of marked nodes* might be an appropriate potential function here. Whenever we do a little bit of work, the number of marks goes up by at most one; whenever we do a lot of work, the number of marks goes down a lot.

More precisely, let  $m$  and  $m'$  be the number of marked nodes before and after a  $\text{DECREASEKEY}$  operation. The actual time (ignoring constant factors) is

$$t = 1 + \#\text{consecutive marked ancestors of } v$$

and if we set  $\Phi = m$ , the increase in potential is

$$m' - m \leq 1 - \#\text{consecutive marked ancestors of } v.$$

Since  $t + \Delta\Phi \leq 2$ , the amortized cost of  $\text{DECREASEKEY}$  is  $O(1)$ .

## B.6 Bounding the Degree

But now we have a problem with our earlier analysis of  $\text{DELETEMIN}$ . The amortized time for a  $\text{DELETEMIN}$  is still  $O(r + \text{deg}(\min))$ . To show that this equaled  $O(\log n)$ , we used the fact that the maximum degree of any node is  $O(\log n)$ , which implies that after a  $\text{CLEANUP}$  the number of roots is  $O(\log n)$ . But now that we don't have complete binomial heaps, this 'fact' is no longer obvious!

So let's prove it. For any node  $v$ , let  $|v|$  denote the number of nodes in the subtree of  $v$ , including  $v$  itself. Our proof uses the following lemma, which *finally* tells us why these things are called Fibonacci heaps.

**Lemma 1.** *For any node  $v$  in a Fibonacci heap,  $|v| \geq F_{\text{deg}(v)+2}$ .*

**Proof:** Label the children of  $v$  in the chronological order in which they were linked to  $v$ . Consider the situation just before the  $i$ th oldest child  $w_i$  was linked to  $v$ . At that time,  $v$  had at least  $i - 1$  children (possibly more). Since  $\text{CLEANUP}$  only links trees with the same degree, we had  $\text{deg}(w_i) =$

$\deg(v) \geq i - 1$ . Since that time, at most one child of  $w_i$  has been promoted away; otherwise,  $w_i$  would have been promoted to the root list by now. So currently we have  $\deg(w_i) \geq i - 2$ .

We also quickly observe that  $\deg(w_i) \geq 0$ . (Duh.)

Let  $s_d$  be the minimum possible size of a tree with degree  $d$  in any Fibonacci heap. Clearly  $s_0 = 1$ ; for notational convenience, let  $s_{-1} = 1$  also. By our earlier argument, the  $i$ th oldest child of the root has degree at least  $\max\{0, i - 2\}$ , and thus has size at least  $\max\{1, s_{i-2}\} = s_{i-2}$ . Thus, we have the following recurrence:

$$s_d \geq 1 + \sum_{i=1}^d s_{i-2}$$

If we assume inductively that  $s_i \geq F_{i+2}$  for all  $-1 \leq i < d$  (with the easy base cases  $s_{-1} = F_1$  and  $s_0 = F_2$ ), we have

$$s_d \geq 1 + \sum_{i=1}^d F_i = F_{d+2}.$$

(The last step was a practice problem in Homework 0.) By definition,  $|v| \geq s_{\deg(v)}$ . □

You can easily show (using either induction or the annihilator method) that  $F_{k+2} > \phi^k$  where  $\phi = \frac{1+\sqrt{5}}{2} \approx 1.618$  is the golden ratio. Thus, Lemma 1 implies that

$$\deg(v) \leq \log_{\phi}|v| = O(\log|v|).$$

Thus, since the size of any subtree in an  $n$ -node Fibonacci heap is obviously at most  $n$ , the degree of any node is  $O(\log n)$ , which is exactly what we wanted. Our earlier analysis is still good.

## B.7 Analyzing Everything Together

Unfortunately, our analyses of DELETETMIN and DECREASEKEY used two different potential functions. Unless we can find a *single* potential function that works for *both* operations, we can't claim both amortized time bounds simultaneously. So we need to find a potential function  $\Phi$  that goes up a little during a cheap DELETETMIN or a cheap DECREASEKEY, and goes down a lot during an expensive DELETETMIN or an expensive DECREASEKEY.

Let's look a little more carefully at the cost of each Fibonacci heap operation, and its effect on both the number of roots and the number of marked nodes, the things we used as our earlier potential functions. Let  $r$  and  $m$  be the numbers of roots and marks before each operation, and let  $r'$  and  $m'$  be the numbers of roots and marks after the operation.

operation	actual cost	$r' - r$	$m' - m$
INSERT	1	1	0
MERGE	1	0	0
DELETETMIN	$r + \deg(\min)$	$r' - r$	0
DECREASEKEY	$1 + m - m'$	$1 + m - m'$	$m' - m$

In particular, notice that promoting a node in DECREASEKEY requires constant time and increases the number of roots by one, and that we promote (at most) one unmarked node.

If we guess that the correct potential function is a linear combination of our old potential functions  $r$  and  $m$  and play around with various possibilities for the coefficients, we will eventually stumble across the correct answer:

$$\Phi = r + 2m$$

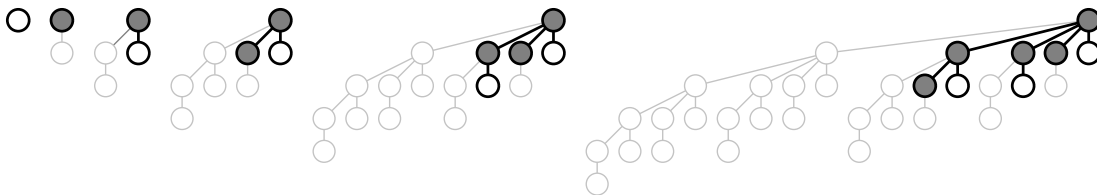
To see that this potential function gives us good amortized bounds for every Fibonacci heap operation, let's add two more columns to our table.

operation	actual cost	$r' - r$	$m' - m$	$\Phi' - \Phi$	amortized cost
INSERT	1	1	0	1	2
MERGE	1	0	0	0	1
DELETEMIN	$r + \text{deg}(\text{min})$	$r' - r$	0	$r' - r$	$r' + \text{deg}(\text{min})$
DECREASEKEY	$1 + m - m'$	$1 + m - m'$	$m' - m$	$1 + m' - m$	2

Since Lemma 1 implies that  $r' + \text{deg}(\text{min}) = O(\log n)$ , we're finally done! (Whew!)

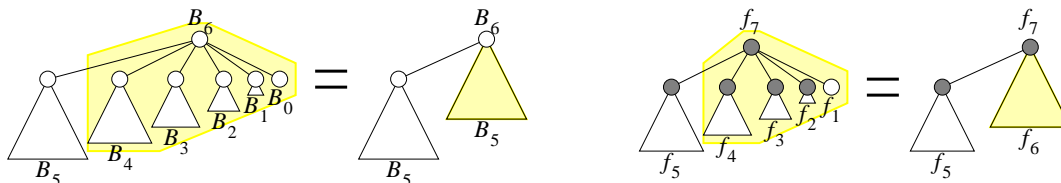
### B.8 Fibonacci Trees

To give you a little more intuition about how Fibonacci heaps behave, let's look at a worst-case construction for Lemma 1. Suppose we want to remove as many nodes as possible from a binomial heap of order  $k$ , by promoting various nodes to the root list, but without causing any cascading promotions. The most damage we can do is to promote the largest subtree of every node. Call the result a *Fibonacci tree* of order  $k + 1$ , and denote it  $f_{k+1}$ . As a base case, let  $f_1$  be the tree with one (unmarked) node, that is,  $f_1 = B_0$ . The reason for shifting the index should be obvious after a few seconds.



Fibonacci trees of order 1 through 6. Light nodes have been promoted away; dark nodes are marked.

Recall that the root of a binomial tree  $B_k$  has  $k$  children, which are roots of  $B_0, B_1, \dots, B_{k-1}$ . To convert  $B_k$  to  $f_{k+1}$ , we promote the root of  $B_{k-1}$ , and recursively convert each of the other subtrees  $B_i$  to  $f_{i+1}$ . The root of the resulting tree  $f_{k+1}$  has degree  $k - 1$ , and the children are the roots of smaller Fibonacci trees  $f_1, f_2, \dots, f_{k-1}$ . We can also consider  $B_k$  as two copies of  $B_{k-1}$  linked together. It's quite easy to show that an order- $k$  Fibonacci tree consists of an order  $k - 2$  Fibonacci tree linked to an order  $k - 1$  Fibonacci tree. (See the picture below.)



Comparing the recursive structures of  $B_6$  and  $f_7$ .

Since  $f_1$  and  $f_2$  both have exactly one node, the number of nodes in an order- $k$  Fibonacci tree is exactly the  $k$ th Fibonacci number! (That's why we changed in the index.) Like binomial trees, Fibonacci trees have lots of other nice properties that easy to prove by induction (hint, hint):

- The root of  $f_k$  has degree  $k - 2$ .
- $f_k$  can be obtained from  $f_{k-1}$  by adding a new unmarked child to every marked node and then marking all the old unmarked nodes.

- $f_k$  has height  $\lceil k/2 \rceil - 1$ .
- $f_k$  has  $F_{k-2}$  unmarked nodes,  $F_{k-1}$  marked nodes, and thus  $F_k$  nodes altogether.
- $f_k$  has  $\binom{k-d-2}{d-1}$  unmarked nodes,  $\binom{k-d-2}{d}$  marked nodes, and  $\binom{k-d-1}{d}$  total nodes at depth  $d$ , for all  $0 \leq d \leq \lceil k/2 \rceil - 1$ .
- $f_k$  has  $F_{k-2h-1}$  nodes with height  $h$ , for all  $0 \leq h \leq \lceil k/2 \rceil - 1$ , and one node (the root) with height  $\lceil k/2 \rceil - 1$ .

Obie looked at the seein' eye dog. Then at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one. . . and then he looked at the seein' eye dog. And then at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one and began to cry.

Because Obie came to the realization that it was a typical case of American blind justice, and there wasn't nothin' he could do about it, and the judge wasn't gonna look at the twenty-seven 8 by 10 color glossy pictures with the circles and arrows and a paragraph on the back of each one explainin' what each one was, to be used as evidence against us.

And we was fined fifty dollars and had to pick up the garbage. In the snow.

But that's not what I'm here to tell you about.

— Arlo Guthrie, “Alice’s Restaurant” (1966)

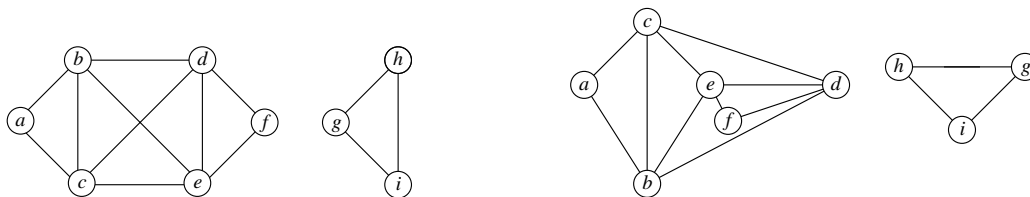
## 10 Basic Graph Properties (October 17)

### 10.1 Definitions

A *graph*  $G$  is a pair of sets  $(V, E)$ .  $V$  is a set of arbitrary objects which we call *vertices*<sup>1</sup> or *nodes*.  $E$  is a set of vertex pairs, which we call *edges* or occasionally *arcs*. In an *undirected* graph, the edges are unordered pairs, or just sets containing two vertices. In a *directed* graph, the edges are ordered pairs of vertices. We will only be concerned with *simple* graphs, where there is no edge from a vertex to itself and there is at most one edge from any vertex to any other.

Following standard (but admittedly confusing) practice, I’ll also use  $V$  to denote the *number* of vertices in a graph, and  $E$  to denote the *number* of edges. Thus, in an undirected graph, we have  $0 \leq E \leq \binom{V}{2}$ , and in a directed graph,  $0 \leq E \leq V(V - 1)$ .

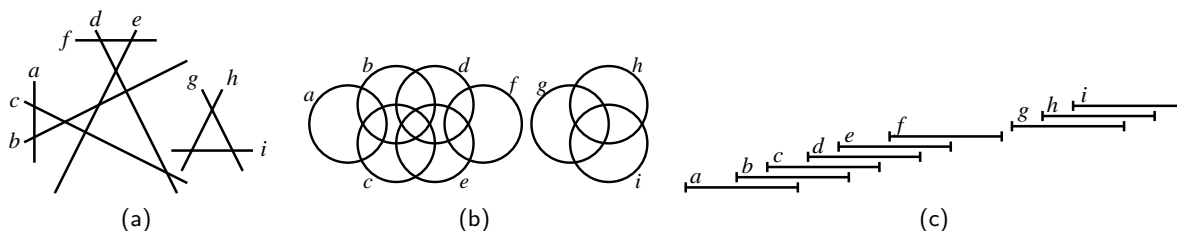
We usually visualize graphs by looking at an *embedding*. An embedding of a graph maps each vertex to a point in the plane and each edge to a curve or straight line segment between the two vertices. A graph is *planar* if it has an embedding where no two edges cross. The same graph can have many different embeddings, so it is important not to confuse a particular embedding with the graph itself. In particular, planar graphs can have non-planar embeddings!



A non-planar embedding of a planar graph with nine vertices, thirteen edges, and two connected components, and a planar embedding of the same graph.

There are other ways of visualizing and representing graphs that are sometimes also useful. For example, the *intersection graph* of a collection of objects has a node for every object and an edge for every intersecting pair. Whether a particular graph can be represented as an intersection graph depends on what kind of object you want to use for the vertices. Different types of objects—line segments, rectangles, circles, etc.—define different classes of graphs. One particularly useful type of intersection graph is an *interval graph*, whose vertices are intervals on the real line, with an edge between any two intervals that overlap.

<sup>1</sup>The singular of ‘vertices’ is **vertex**. The singular of ‘matrices’ is **matrix**. Unless you’re speaking Italian, there is no such thing as a vertice, a matrice, an indice, an appendice, a helice, an apice, a vortice, a radice, a simplice, a directrice, a dominatrice, a Unice, a Kleenice, or Jimi Hendrice!



The example graph is also the intersection graph of (a) a set of line segments, (b) a set of circles, or (c) a set of intervals on the real line (stacked for visibility).

If  $(u, v)$  is an edge in an undirected graph, then  $u$  is a *neighbor* of  $v$  and vice versa. The *degree* of a node is the number of neighbors. In directed graphs, we have two kinds of neighbors. If  $u \rightarrow v$  is a directed edge, then  $u$  is a *predecessor* of  $v$  and  $v$  is a *successor* of  $u$ . The *in-degree* of a node is the number of predecessors, which is the same as the number of edges going into the node. The *out-degree* is the number of successors, or the number of edges going out of the node.

A graph  $G' = (V', E')$  is a *subgraph* of  $G = (V, E)$  if  $V' \subseteq V$  and  $E' \subseteq E$ .

A *path* is a sequence of edges, where each successive pair of edges shares a vertex, and all other edges are disjoint. A graph is *connected* if there is a path from any vertex to any other vertex. A disconnected graph consists of several *connected components*, which are maximal connected subgraphs. Two vertices are in the same connected component if and only if there is a path between them.

A *cycle* is a path that starts and ends at the same vertex, and has at least one edge. A graph is *acyclic* if no subgraph is a cycle; acyclic graphs are also called *forests*. *Trees* are special graphs that can be defined in several different ways. You can easily prove by induction (hint, hint, hint) that the following definitions are equivalent.

- A tree is a connected acyclic graph.
- A tree is a connected component of a forest.
- A tree is a connected graph with *at most*  $V - 1$  edges.
- A tree is a minimal connected graph; removing any edge makes the graph disconnected.
- A tree is an acyclic graph with *at least*  $V - 1$  edges.
- A tree is a maximal acyclic graph; adding an edge between any two vertices creates a cycle.

A *spanning tree* of a graph  $G$  is a subgraph that is a tree and contains every vertex of  $G$ . Of course, a graph can only have a spanning tree if it's connected. A *spanning forest* of  $G$  is a collection of spanning trees, one for each connected component of  $G$ .

## 10.2 Explicit Representations of Graphs

There are two common data structures used to explicitly represent graphs: *adjacency matrices*<sup>2</sup> and *adjacency lists*.

The adjacency matrix of a graph  $G$  is a  $V \times V$  matrix of indicator variables. Each entry in the matrix indicates whether a particular edge is or is not in the graph:

$$A[i, j] = [(i, j) \in E].$$

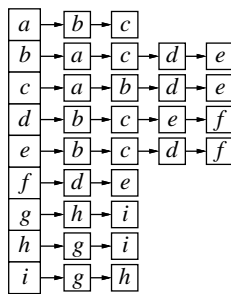
<sup>2</sup>See footnote 1.



For undirected graphs, the adjacency matrix is always *symmetric*:  $A[i, j] = A[j, i]$ . Since we don't allow edges from a vertex to itself, the diagonal elements  $A[i, i]$  are all zeros.

Given an adjacency matrix, we can decide in  $\Theta(1)$  time whether two vertices are connected by an edge just by looking in the appropriate slot in the matrix. We can also list all the neighbors of a vertex in  $\Theta(V)$  time by scanning the corresponding row (or column). This is optimal in the worst case, since a vertex can have up to  $V - 1$  neighbors; however, if a vertex has few neighbors, we may still have to examine every entry in the row to see them all. Similarly, adjacency matrices require  $\Theta(V^2)$  space, regardless of how many edges the graph actually has, so it is only space-efficient for very *dense* graphs.

	<i>a</i>	<i>b</i>	<i>c</i>	<i>d</i>	<i>e</i>	<i>f</i>	<i>g</i>	<i>h</i>	<i>i</i>
<i>a</i>	0	1	1	0	0	0	0	0	0
<i>b</i>	1	0	1	1	1	0	0	0	0
<i>c</i>	1	1	0	1	1	0	0	0	0
<i>d</i>	0	1	1	0	1	1	0	0	0
<i>e</i>	0	1	1	1	0	1	0	0	0
<i>f</i>	0	0	0	1	1	0	0	0	0
<i>g</i>	0	0	0	0	0	0	0	0	1
<i>h</i>	0	0	0	0	0	0	0	1	1
<i>i</i>	0	0	0	0	0	0	1	1	0



Adjacency matrix and adjacency list representations for the example graph.

For *sparse* graphs—graphs with relatively few edges—we're better off using adjacency lists. An adjacency list is an array of linked lists, one list per vertex. Each linked list stores the neighbors of the corresponding vertex.

For undirected graphs, each edge  $(u, v)$  is stored twice, once in  $u$ 's neighbor list and once in  $v$ 's neighbor list; for directed graphs, each edge is stored only once. Either way, the overall space required for an adjacency list is  $O(V + E)$ . Listing the neighbors of a node  $v$  takes  $O(1 + \deg(v))$  time; just scan the neighbor list. Similarly, we can determine whether  $(u, v)$  is an edge in  $O(1 + \deg(u))$  time by scanning the neighbor list of  $u$ . For undirected graphs, we can speed up the search by simultaneously scanning the neighbor lists of both  $u$  and  $v$ , stopping either we locate the edge or when we fall off the end of a list. This takes  $O(1 + \min\{\deg(u), \deg(v)\})$  time.

The adjacency list structure should immediately remind you of hash tables with chaining. Just as with hash tables, we can make adjacency list structure more efficient by using something besides a linked list to store the neighbors. For example, if we use a hash table with constant load factor, when we can detect edges in  $O(1)$  expected time, just as with an adjacency list. In practice, this will only be useful for vertices with large degree, since the constant overhead in both the space and search time is larger for hash tables than for simple linked lists.

You might at this point ask why anyone would ever use an adjacency matrix. After all, if you use hash tables to store the neighbors of each vertex, you can do everything as fast or faster with an adjacency list as with an adjacency matrix, only using less space. The answer is that many graphs are only represented *implicitly*. For example, intersection graphs are usually represented implicitly by simply storing the list of objects. As long as we can test whether two objects overlap in constant time, we can apply any graph algorithm to an intersection graph by *pretending* that it is stored explicitly as an adjacency matrix. On the other hand, any data structure built from records with pointers between them can be seen as a directed graph. Algorithms for searching graphs can be applied to these data structures by *pretending* that the graph is represented explicitly using an adjacency list.

To keep things simple, we'll consider only undirected graphs for the rest of this lecture, although the algorithms I'll describe also work for directed graphs.

### 10.3 Traversing connected graphs

Suppose we want to visit every node in a connected graph (represented either explicitly or implicitly). The simplest method to do this is an algorithm called *depth-first search*, which can be written either recursively or iteratively. It's exactly the same algorithm either way; the only difference is that we can actually see the 'recursion' stack in the non-recursive version. Both versions are initially passed a *source* vertex  $s$ .

<pre> RECURSIVEDFS(<math>v</math>):   if <math>v</math> is unmarked     mark <math>v</math>     for each edge <math>(v, w)</math>       RECURSIVEDFS(<math>w</math>) </pre>	<pre> ITERATEDFS(<math>s</math>):   PUSH(<math>s</math>)   while stack not empty     <math>v \leftarrow</math> POP     if <math>v</math> is unmarked       mark <math>v</math>       for each edge <math>(v, w)</math>         PUSH(<math>w</math>) </pre>
---	--

Depth-first search is one (perhaps the most common) instance of a general family of graph traversal algorithms. The generic graph traversal algorithm stores a set of candidate edges in some data structure that I'll call a 'bag'. The only important properties of a 'bag' are that we can put stuff into it and then later take stuff back out. (In C++ terms, think of the 'bag' as a template for a real data structure.) Here's the algorithm:

<pre> TRAVERSE(<math>s</math>):   put <math>(\emptyset, s)</math> in bag   while the bag is not empty     take <math>(p, v)</math> from the bag    (*)     if <math>v</math> is unmarked       mark <math>v</math>       parent(<math>v</math>) <math>\leftarrow</math> <math>p</math>       for each edge <math>(v, w)</math>    (†)         put <math>(v, w)</math> into the bag    (**) </pre>
---

Notice that we're keeping *edges* in the bag instead of *vertices*. This is because we want to remember, whenever we visit a vertex  $v$  for the first time, which previously-visited vertex  $p$  put  $v$  into the bag. The vertex  $p$  is called the *parent* of  $v$ .

**Lemma 1.** TRAVERSE( $s$ ) marks every vertex in any connected graph exactly once, and the set of edges  $(v, \text{parent}(v))$  with  $\text{parent}(v) \neq \emptyset$  form a spanning tree of the graph.

**Proof:** First, it should be obvious that no node is marked more than once.

Clearly, the algorithm marks  $s$ . Let  $v \neq s$  be a vertex, and let  $s \rightarrow \dots \rightarrow u \rightarrow v$  be the path from  $s$  to  $v$  with the minimum number of edges. Since the graph is connected, such a path always exists. (If  $s$  and  $v$  are neighbors, then  $u = s$ , and the path has just one edge.) If the algorithm marks  $u$ , then it must put  $(u, v)$  into the bag, so it must later take  $(u, v)$  out of the bag, at which point  $v$  must be marked (if it isn't already). Thus, by induction on the shortest-path distance from  $s$ , the algorithm marks every vertex in the graph.

Call an edge  $(v, \text{parent}(v))$  with  $\text{parent}(v) \neq \emptyset$  a *parent edge*. For any node  $v$ , the path of parent edges  $v \rightarrow \text{parent}(v) \rightarrow \text{parent}(\text{parent}(v)) \rightarrow \dots$  eventually leads back to  $s$ , so the set of

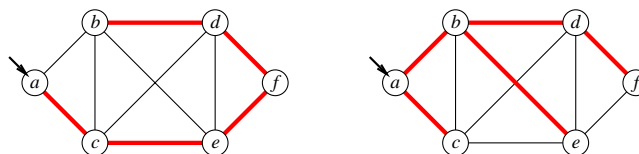
parent edges form a connected graph. Clearly, both endpoints of every parent edge are marked, and the number of parent edges is exactly one less than the number of vertices. Thus, the parent edges form a spanning tree.  $\square$

The exact running time of the traversal algorithm depends on how the graph is represented and what data structure is used as the ‘bag’, but we can make a few general observations. Since each vertex is visited at most once, the for loop ( $\dagger$ ) is executed at most  $V$  times. Each edge is put into the bag exactly twice; once as  $(u, v)$  and once as  $(v, u)$ , so line ( $\star\star$ ) is executed at most  $2E$  times. Finally, since we can’t take more things out of the bag than we put in, line ( $\star$ ) is executed at most  $2E + 1$  times.

## 10.4 Examples

Let’s first assume that the graph is represented by an adjacency list, so that the overhead of the for loop ( $\dagger$ ) is only a constant per edge.

- If we implement the ‘bag’ by using a *stack*, we have *depth-first search*. Each execution of ( $\star$ ) or ( $\star\star$ ) takes constant time, so the overall running time is  $O(V + E)$ . Since the graph is connected,  $V \leq E + 1$ , so we can simplify the running time to  $O(E)$ . The spanning tree formed by the parent edges is called a *depth-first spanning tree*. The exact shape of the tree depends on the order in which neighbor edges are pushed onto the stack, but in general, depth-first spanning trees are long and skinny.
- If we use a *queue* instead of a stack, we have *breadth-first search*. Again, each execution of ( $\star$ ) or ( $\star\star$ ) takes constant time, so the overall running time is still  $O(E)$ . In this case, the *breadth-first spanning tree* formed by the parent edges contains *shortest paths* from the start vertex  $s$  to every other vertex in its connected component. The exact shape of the shortest path tree depends on the order in which neighbor edges are pushed onto the queue, but in general, shortest path trees are short and bushy. We’ll see shortest paths again next week.



A depth-first spanning tree and a breadth-first spanning tree of one component of the example graph, with start vertex  $a$ .

- Suppose the edges of the graph are weighted. If we implement the ‘bag’ using a *priority queue*, always extracting the minimum-weight edge in line ( $\star$ ), then we have what might be called *shortest-first search*. In this case, each execution of ( $\star$ ) or ( $\star\star$ ) takes  $O(\log E)$  time, so the overall running time is  $O(V + E \log E)$ , which simplifies to  $O(E \log E)$  if the graph is connected. For this algorithm, the set of parent edges form the *minimum spanning tree* of the connected component of  $s$ . We’ll see minimum spanning trees again in the next lecture.

If the graph is represented using an adjacency matrix, the finding all the neighbors of each vertex in line ( $\dagger$ ) takes  $O(V)$  time. Thus, depth- and breadth-first search take  $O(V^2)$  time overall, and ‘shortest-first search’ takes  $O(V^2 + E \log E) = O(V^2 \log V)$  time overall.

## 10.5 Searching disconnected graphs

If the graph is disconnected, then  $\text{TRAVERSE}(s)$  only visits the nodes in the connected component of the start vertex  $s$ . If we want to visit all the nodes in every component, we can use the following ‘wrapper’ around our generic traversal algorithm. Since  $\text{TRAVERSE}$  computes a spanning tree of one component,  $\text{TRAVERSEALL}$  computes a spanning *forest* of the entire graph.

$\text{TRAVERSEALL}(s)$ :  
for all vertices  $v$   
if  $v$  is unmarked  
     $\text{TRAVERSE}(v)$

There is a rather unfortunate mistake on page 477 of CLR:

Unlike breadth-first search, whose predecessor subgraph forms a tree, the predecessor subgraph produced by depth-first search may be composed of several trees, because the search may be repeated from multiple sources.

This statement seems to imply that depth-first search is always called with the  $\text{TRAVERSEALL}$ , and breadth-first search never is, **but this is not true!** The choice of whether to use a stack or a queue is completely independent of the choice of whether to use  $\text{TRAVERSEALL}$  or not.

## 11 Shortest Paths (October 22)

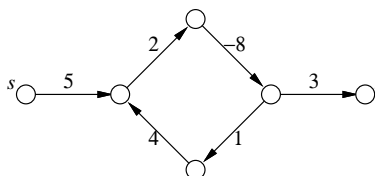
### 11.1 Introduction

Given a weighted *directed* graph  $G = (V, E, w)$  with two special vertices, a *source*  $s$  and a *target*  $t$ , we want to find the shortest directed path from  $s$  to  $t$ . In other words, we want to find the path  $p$  starting at  $s$  and ending at  $t$  minimizing the function

$$w(p) = \sum_{e \in p} w(e).$$

For example, if I want to answer the question ‘What’s the fastest way to drive from my old apartment in Champaign, Illinois to my wife’s old apartment in Columbus, Ohio?’, we might use a graph whose vertices are cities, edges are roads, weights are driving times,  $s$  is Champaign, and  $t$  is Columbus.<sup>1</sup> The graph is directed since the driving times along the same road might be different in different directions.<sup>2</sup>

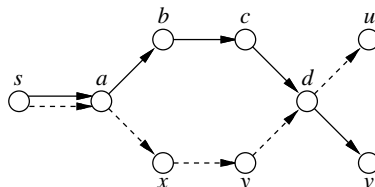
Perhaps counter to intuition, we will allow the weights on the edges to be negative. Negative edges make our lives complicated, since the presence of a negative cycle might mean that there is no shortest path. In general, a shortest path from  $s$  to  $t$  exists if and only if there is *at least one* path from  $s$  to  $t$ , but there is no path from  $s$  to  $t$  that touches a negative cycle. If there is a negative cycle between  $s$  and  $t$ , then we can always find a shorter path by going around the cycle one more time.



There is no shortest path from  $s$  to  $t$ .

Every algorithm known for solving this problem actually solves (large portions of) the following more general *single source shortest path* or *SSSP* problem: find the shortest path from the source vertex  $s$  to *every* other vertex in the graph. In fact, the problem is usually solved by finding a *shortest path tree* rooted at  $s$  that contains all the desired shortest paths.

It’s not hard to see that if the shortest paths are unique, then they form a tree. To prove this, it’s enough to observe that sub-paths of shortest paths are also shortest paths. If there are multiple shortest paths to the same vertices, we can always choose one path to each vertex so that the union of the paths is a tree. If there are shortest paths to two vertices  $u$  and  $v$  that diverge, then meet, then diverge again, we can modify one of the paths so that the two paths only diverge once.

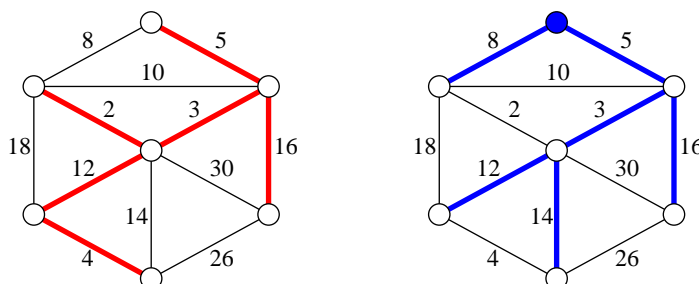


If  $s \rightarrow a \rightarrow b \rightarrow c \rightarrow d \rightarrow v$  and  $s \rightarrow a \rightarrow x \rightarrow y \rightarrow d \rightarrow u$  are both shortest paths, then  $s \rightarrow a \rightarrow b \rightarrow c \rightarrow d \rightarrow u$  is also a shortest path.

<sup>1</sup>West on Church, north on Prospect, east on I-74, south on I-465, east on Airport Expressway, north on I-65, east on I-70, north on Grandview, east on 5th, north on Olentangy River, east on Dodridge, north on High, west on Kelso, south on Neil. Depending on traffic. We both live in Urbana now.

<sup>2</sup>There is a speed trap on I-70 just inside the Ohio border, but only for eastbound traffic.

I should emphasize here that shortest path trees and minimum spanning trees are usually very different. For one thing, there is only one minimum spanning tree, but in general, there is a different shortest path tree for every source vertex.



A minimum spanning tree and a shortest path tree (rooted at the topmost vertex) of the same graph.

All of the algorithms I'm describing in this lecture also work for undirected graphs, with some slight modifications. Most importantly, we must specifically prohibit alternating back and forth across the same undirected negative-weight edge. Our unmodified algorithms would interpret any such edge as a negative cycle of length 2.

## 11.2 The Only SSSP Algorithm

Just like graph traversal and minimum spanning trees, there are several different SSSP algorithms, but they are all special cases of the a single generic algorithm. Each vertex  $v$  in the graph stores two values, which describe a *tentative* shortest path from  $s$  to  $v$ .

- $\text{dist}(v)$  is the length of the tentative shortest  $s \rightsquigarrow v$  path.
- $\text{pred}(v)$  is the predecessor of  $v$  in the tentative shortest  $s \rightsquigarrow v$  path.

Notice that the predecessor pointers automatically define a tentative shortest path tree. We already know that  $\text{dist}(s) = 0$  and  $\text{pred}(s) = \text{NULL}$ . For every vertex  $v \neq s$ , we initially set  $\text{dist}(v) = \infty$  and  $\text{pred}(v) = \text{NULL}$  to indicate that we do not know of *any* path from  $s$  to  $v$ .

We call an edge  $u \rightarrow v$  *tense* if  $\text{dist}(u) + w(u \rightarrow v) < \text{dist}(v)$ . If  $u \rightarrow v$  is tense, then the tentative shortest path  $s \rightsquigarrow v$  is incorrect, since the path  $s \rightsquigarrow u \rightarrow v$  is shorter. Our generic algorithm repeatedly finds a tense edge in the graph and *relaxes* it:

$\begin{array}{l} \text{RELAX}(u \rightarrow v): \\ \text{dist}(v) \leftarrow \text{dist}(u) + w(u \rightarrow v) \\ \text{pred}(v) \leftarrow u \end{array}$
---

If there are no tense edges, our algorithm is finished, and we have our desired shortest path tree.

The correctness of the relaxation algorithm follows directly from three simple claims:

1. If  $\text{dist}(v) \neq \infty$ , then  $\text{dist}(v)$  is the total weight of the predecessor chain ending at  $v$ :

$$s \rightarrow \dots \rightarrow \text{pred}(\text{pred}(v)) \rightarrow \text{pred}(v) \rightarrow v.$$

This is easy to prove by induction on the number of relaxation steps. (Hint, hint.)

2. If the algorithm halts, then  $\text{dist}(v) \leq w(s \rightsquigarrow v)$  for *any* path  $s \rightsquigarrow v$ . This is easy to prove by induction on the number of edges in the path  $s \rightsquigarrow v$ . (Hint, hint.)

3. The algorithm halts if and only if there is no negative cycle reachable from  $s$ . The ‘only if’ direction is easy—if there is a reachable negative cycle, then after the first edge in the cycle is relaxed, the cycle *always* has at least one tense edge. The ‘if’ direction follows from the fact that every relaxation step reduces either the number of vertices with  $\text{dist}(v) = \infty$  by 1 or reduces the sum of the finite shortest path lengths by some positive amount.

I haven’t said anything about how we detect which edges can be relaxed, or what order we relax them in. In order to make this easier, we can refine the relaxation algorithm slightly, into something closely resembling the generic graph traversal algorithm. We maintain a ‘bag’ of vertices, initially containing just the source vertex  $s$ . Whenever we take a vertex  $u$  out of the bag, we scan all of its outgoing edges, looking for something to relax. Whenever we successfully relax an edge  $u \rightarrow v$ , we put  $v$  in the bag.

<u>INITSSSP(<math>s</math>):</u>
$\text{dist}(s) \leftarrow 0$
$\text{pred}(s) \leftarrow \text{NULL}$
for all vertices $v \neq s$
$\text{dist}(v) \leftarrow \infty$
$\text{pred}(v) \leftarrow \text{NULL}$

<u>GENERICSSSP(<math>s</math>):</u>
INITSSSP( $s$ )
put $s$ in the bag
while the bag is not empty
take $u$ from the bag
for all edges $u \rightarrow v$
if $u \rightarrow v$ is tense
RELAX( $u \rightarrow v$ )
put $v$ in the bag

Just as with graph traversal, using different data structures for the ‘bag’ gives us different algorithms. There are three obvious choices to try: a stack, a queue, and a heap. Unfortunately, if we use a stack, we have to perform  $\Theta(2^V)$  relaxation steps in the worst case! (This is a problem in the current homework.) The other two possibilities are much more efficient.

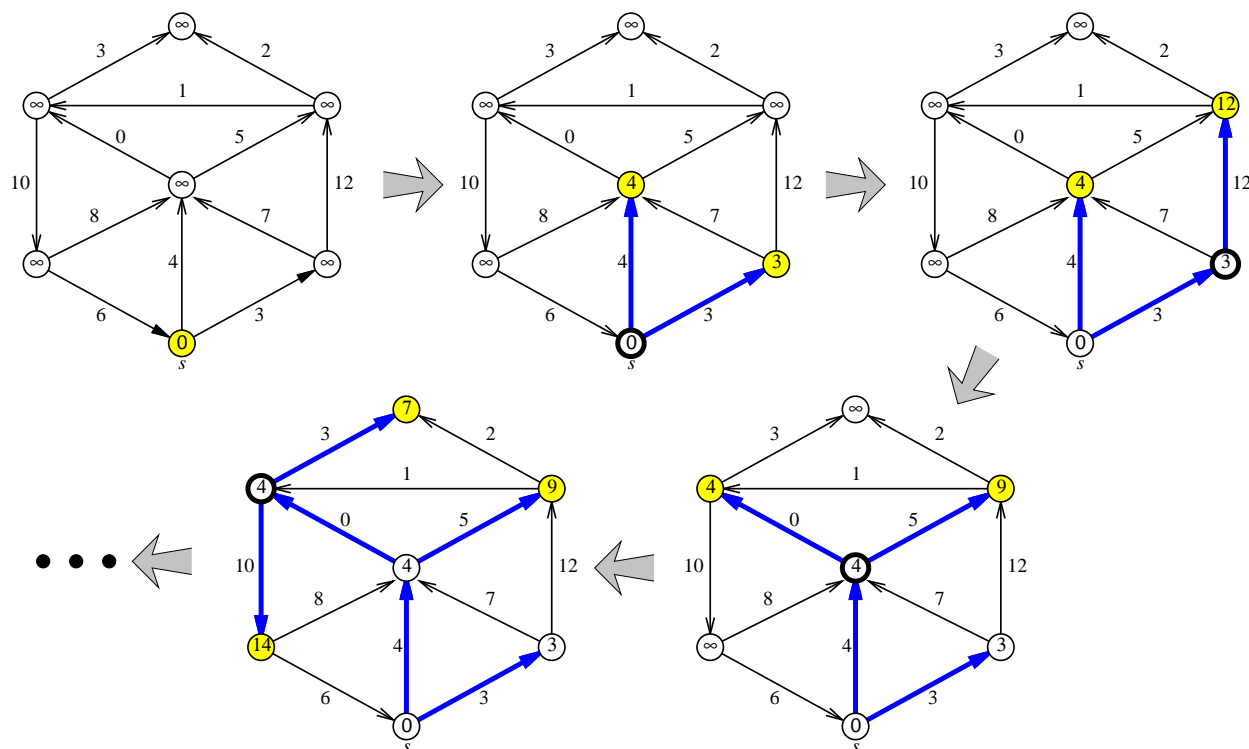
### 11.3 Dijkstra’s Algorithm

If we implement the bag as a heap, where the key of a vertex  $v$  is  $\text{dist}(v)$ , we obtain an algorithm first published by Edsger Dijkstra in 1959.

Dijkstra’s algorithm is particularly well-behaved if the graph has no negative-weight edges. In this case, it’s not hard to show (by induction, of course) that the vertices are scanned in increasing order of their shortest-path distance from  $s$ . It follows that each vertex is scanned at most once, and thus that each edge is relaxed at most once. Since the key of each vertex in the heap is its tentative distance from  $s$ , the algorithm performs a DECREASEKEY operation every time an edge is relaxed. Thus, the algorithm performs at most  $E$  DECREASEKEYS. Similarly, there are at most  $V$  INSERT and EXTRACTMIN operations. Thus, if we store the vertices in a Fibonacci heap, the total running time of Dijkstra’s algorithm is  $O(E + V \log V)$ .

This analysis assumes that no edge has negative weight. Dijkstra’s algorithm (in the form I’m presenting here) is still *correct* if there are negative edges<sup>3</sup>, but the worst-case running time could be exponential. (I’ll leave the proof of this unfortunate fact as an extra credit problem.)

<sup>3</sup>The version of Dijkstra’s algorithm presented in CLRS gives incorrect results for graphs with negative edges.



Four phases of Dijkstra's algorithm run on a graph with no negative edges.  
 At each phase, the shaded vertices are in the heap, and the bold vertex has just been scanned.  
 The bold edges describe the evolving shortest path tree.

## 11.4 The $A^*$ Heuristic

A slight generalization of Dijkstra's algorithm, commonly known as the  $A^*$  algorithm, is frequently used to find a shortest path from a single source node  $s$  to a single target node  $t$ .  $A^*$  uses a black-box function  $\text{GUESSDISTANCE}(v, t)$  that returns an estimate of the distance from  $v$  to  $t$ . The only difference between Dijkstra and  $A^*$  is that the key of a vertex  $v$  is  $\text{dist}(v) + \text{GUESSDISTANCE}(v, t)$ .

The function  $\text{GUESSDISTANCE}$  is called *admissible* if  $\text{GUESSDISTANCE}(v, t)$  never overestimates the actual shortest path distance from  $v$  to  $t$ . If  $\text{GUESSDISTANCE}$  is admissible and the actual edge weights are all non-negative, the  $A^*$  algorithm computes the actual shortest path from  $s$  to  $t$  at least as quickly as Dijkstra's algorithm. The closer  $\text{GUESSDISTANCE}(v, t)$  is to the real distance from  $v$  to  $t$ , the faster the algorithm. However, in the worst case, the running time is still  $O(E + V \log V)$ .

The heuristic is especially useful in situations where the actual graph is not known. For example,  $A^*$  can be used to solve puzzles (15-puzzle, Freecell, Shanghai, Minesweeper, Sokoban, Atomix, Rush Hour, Rubik's Cube, ...) and other path planning problems where the starting and goal configurations are given, but the graph of all possible configurations and their connections is not given explicitly.

## 11.5 Moore's Algorithm ('Bellman-Ford')

If we replace the heap in Dijkstra's algorithm with a queue, we get an algorithm that was first published by Moore in 1957, and then independently by Bellman in 1958. (Since Bellman used the idea of relaxing edges, which was first proposed by Ford in 1956, this algorithm is usually called 'Bellman-Ford'.) Moore's algorithm is efficient even if there are negative edges, and it can be used to quickly detect the presence of negative cycles. If there are no negative edges, however, Dijkstra's

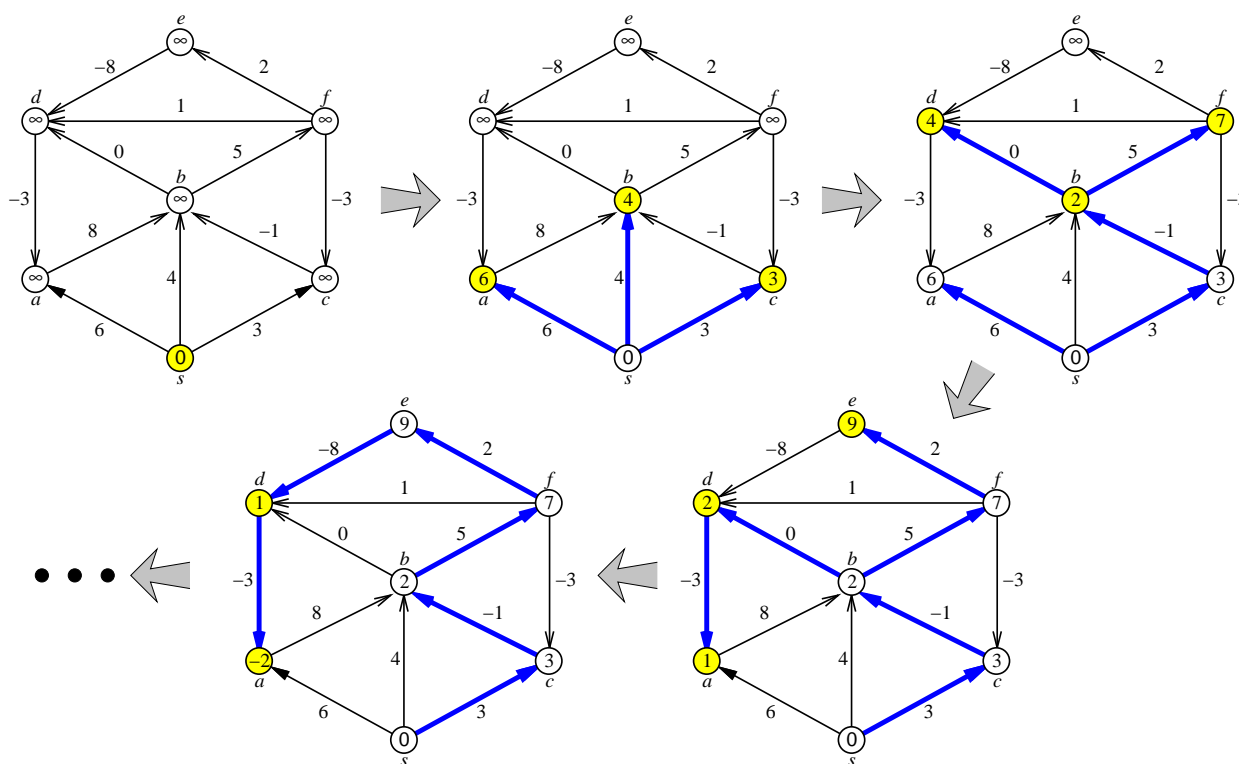


algorithm is faster. (In fact, in practice, Dijkstra's algorithm is often faster even for graphs with negative edges.)

The easiest way to analyze the algorithm is to break the execution into phases. Before we even begin, we insert a token into the queue. Whenever we take the token out of the queue, we begin a new phase by just reinserting the token into the queue. The 0th phase consists entirely of scanning the source vertex  $s$ . The algorithm ends when the queue contains *only* the token. A simple inductive argument (hint, hint) implies the following invariant:

At the end of the  $i$ th phase, for every vertex  $v$ ,  $\text{dist}(v)$  is less than or equal to the length of the shortest path  $s \rightsquigarrow v$  consisting of  $i$  or fewer edges.

Since a shortest path can only pass through each vertex once, either the algorithm halts before the  $V$ th phase, or the graph contains a negative cycle. In each phase, we scan each vertex at most once, so we relax each edge at most once, so the running time of a single phase is  $O(E)$ . Thus, the overall running time of Moore's algorithm is  $O(VE)$ .



Four phases of Moore's algorithm run on a directed graph with negative edges.

Nodes are taken from the queue in the order  $s \diamond a b c \diamond d f b \diamond a e d \diamond d a \diamond \diamond$ , where  $\diamond$  is the token.

Shaded vertices are in the queue at the end of each phase. The bold edges describe the evolving shortest path tree.

Now that we understand how the phases of Moore's algorithm behave, we can simplify the algorithm considerably. Instead of using a queue to perform a partial breadth-first search of the graph in each phase, let's just scan through the adjacency list directly and try to relax every edge in the graph.

```

MOORESSSP( $s$ )
  INITSSSP( $s$ )
  repeat  $V$  times:
    for every edge  $u \rightarrow v$ 
      if  $u \rightarrow v$  is tense
        RELAX( $u \rightarrow v$ )
  for every edge  $u \rightarrow v$ 
    if  $u \rightarrow v$  is tense
      return 'Negative cycle!'

```

This is closer to how CLRS presents the ‘Bellman-Ford’ algorithm. The  $O(VE)$  running time of this version of the algorithm should be obvious, but it may not be clear that the algorithm is still correct. To prove correctness, we just have to show that our earlier invariant holds; as before, this can be proved by induction on  $i$ .

## 11.6 Greedy Shortest Paths?

Here’s another algorithm that fits our generic framework, but which I’ve never seen analyzed.

Repeatedly relax the tensest edge.

Specifically, let’s define the ‘tenseness’ of an edge  $u \rightarrow v$  as follows:

$$\text{tenseness}(u \rightarrow v) = \max \{0, \text{dist}(v) - \text{dist}(u) - w(u \rightarrow v)\}$$

(This is defined even when  $\text{dist}(v) = \infty$  or  $\text{dist}(u) = \infty$ , as long as we treat  $\infty$  just like some indescribably large but finite number.) If an edge has zero tenseness, it’s not tense. If we relax an edge  $u \rightarrow v$ , then  $\text{dist}(v)$  decreases by the edge’s tenseness.

Intuitively, we can keep the edges of the graph in some sort of heap, where the key of each edge is its tenseness. Then we repeatedly pull out the tensest edge  $u \rightarrow v$  and relax it. Then we need to recompute the tenseness of other edges adjacent to  $v$ . Edges leaving  $v$  possibly become more tense, and edges coming into  $v$  possibly become less tense. So we need a heap that efficiently supports the operations INSERT, EXTRACTMAX, INCREASEKEY, and DECREASEKEY.

If there are no negative cycles, this algorithm eventually halts with a shortest path tree, but how quickly? Can the same edge be relaxed more than once, and if so, how many times? Is it faster if all the edge weights are positive? Hmm....

## 12 All-Pair Shortest Paths (October 24)

### 12.1 The Problem

In the last lecture, we saw algorithms to find the shortest path from a source vertex  $s$  to a target vertex  $t$  in a directed graph. As it turns out, the best algorithms for this problem actually find the shortest path from  $s$  to every possible target (or from every possible source to  $t$ ) by constructing a shortest path tree. The shortest path tree specifies two pieces of information for each node  $v$  in the graph

- $\text{dist}(v)$  is the length of the shortest path (if any) from  $s$  to  $v$ .
- $\text{pred}(v)$  is the second-to-last vertex (if any) the shortest path (if any) from  $s$  to  $v$ .

In this lecture, we want to generalize the shortest path problem even further. In the *all pairs shortest path* problem, we want to find the shortest path from *every* possible source to *every* possible destination. Specifically, for every pair of vertices  $u$  and  $v$ , we need to compute the following information:

- $\text{dist}(u, v)$  is the length of the shortest path (if any) from  $u$  to  $v$ .
- $\text{pred}(u, v)$  is the second-to-last vertex (if any) on the shortest path (if any) from  $u$  to  $v$ .

For example, for any vertex  $v$ , we have  $\text{dist}(v, v) = 0$  and  $\text{pred}(v, v) = \text{NULL}$ . If the shortest path from  $u$  to  $v$  is only one edge long, then  $\text{dist}(u, v) = w(u \rightarrow v)$  and  $\text{pred}(u, v) = u$ . If there is *no* shortest path from  $u$  to  $v$ —either because there’s no path at all, or because there’s a negative cycle—then  $\text{dist}(u, v) = \infty$  and  $\text{pred}(u, v) = \text{NULL}$ .

The output of our shortest path algorithms will be a pair of  $V \times V$  arrays encoding all  $V^2$  distances and predecessors. Many maps include a distance matrix—to find the distance from (say) Champaign to (say) Columbus, you would look in the row labeled ‘Champaign’ and the column labeled ‘Columbus’. In these notes, I’ll focus almost exclusively on computing the distance array. The predecessor array, from which you would compute the actual shortest paths, can be computed with only minor additions to the algorithms I’ll describe (hint, hint).

### 12.2 Lots of Single Sources

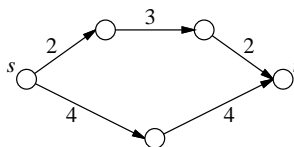
The most obvious solution to the all pairs shortest path problem is just to run a single-source shortest path algorithm  $V$  times, once for every possible source vertex! Specifically, to fill in the one-dimensional subarray  $\text{dist}[s][\ ]$ , we invoke either Dijkstra’s or Moore’s algorithm starting at the source vertex  $s$ .

OBVIOUSAPSP( $V, E, w$ ):  
 for every vertex  $s$   
 $\text{dist}[s][\ ] \leftarrow \text{SSSP}(V, E, w, s)$

The running time of this algorithm depends on which single source algorithm we use. If we use Moore’s algorithm, the overall running time is  $\Theta(V^2E) = O(V^4)$ . If all the edge weights are positive, we can use Dijkstra’s algorithm instead, which decreases the running time to  $\Theta(VE + V^2 \log V) = O(V^3)$ . For graphs with negative edge weights, Dijkstra’s algorithm can take exponential time, so we can’t get this improvement directly.

### 12.3 Reweighting

One idea that occurs to most people is increasing the weights of all the edges by the same amount so that all the weights become positive, and then applying Dijkstra's algorithm. Unfortunately, this simple idea doesn't work. Different paths change by different amounts, which means the shortest paths in the reweighted graph may not be the same as in the original graph.



Increasing all the edge weights by 2 changes the shortest path  $s$  to  $t$ .

However, there is a more complicated method for reweighting the edges in a graph. Suppose each vertex  $v$  has some associated *cost*  $c(v)$ , which might be positive, negative, or zero. We can define a new weight function  $w'$  as follows:

$$w'(u \rightarrow v) = c(u) + w(u \rightarrow v) - c(v)$$

To give some intuition, imagine that when we leave vertex  $u$ , we have to pay an exit tax of  $c(u)$ , and when we enter  $v$ , we get  $c(v)$  as an entrance gift.

Now it's not too hard to show that the shortest paths with the new weight function  $w'$  are exactly the same as the shortest paths with the original weight function  $w$ . In fact, for *any* path  $u \rightsquigarrow v$  from one vertex  $u$  to another vertex  $v$ , we have

$$w'(u \rightsquigarrow v) = c(u) + w(u \rightsquigarrow v) - c(v).$$

We pay  $c(u)$  in exit fees, plus the original weight of the path, minus the  $c(v)$  entrance gift. At every intermediate vertex  $x$  on the path, we get  $c(x)$  as an entrance gift, but then immediately pay it back as an exit tax!

### 12.4 Johnson's Algorithm

Johnson's all-pairs shortest path algorithm finds a cost  $c(v)$  for each vertex, so that when the graph is reweighted, every edge has non-negative weight.

Suppose the graph has a vertex  $s$  that has a path to every other vertex. Johnson's algorithm computes the shortest paths from  $s$  to every other vertex, using Moore's algorithm (which doesn't care if the edge weights are negative), and then sets

$$c(v) = \text{dist}(s, v),$$

so the new weight of every edge is

$$w'(u \rightarrow v) = \text{dist}(s, u) + w(u \rightarrow v) - \text{dist}(s, v).$$

Why are all these new weights non-negative? Because otherwise, Moore's algorithm wouldn't be finished! Recall that an edge  $u \rightarrow v$  is *tense* if  $\text{dist}(s, u) + w(u \rightarrow v) < \text{dist}(s, v)$ , and that single-source shortest path algorithms eliminate all tense edges. The only exception is if the graph has a negative cycle, but then shortest paths aren't defined, and Johnson's algorithm simply aborts.

But what if the graph *doesn't* have a vertex  $s$  that can reach everything? Then no matter where we start Moore's algorithm, some of those vertex costs will be infinite. Johnson's algorithm

avoids this problem by adding a new vertex  $s$  to the graph, with zero-weight edges going from  $s$  to every other vertex, but *no* edges going back into  $s$ . This addition doesn't change the shortest paths between any other pair of vertices, because there are no paths into  $s$ .

So here's Johnson's algorithm in all its glory.

```

JOHNSONAPSP( $V, E, w$ ) :
  create a new vertex  $s$ 
  for every vertex  $v \in V$ 
     $w(s \rightarrow v) \leftarrow 0$ ;  $w(v \rightarrow s) \leftarrow \infty$ 
   $\text{dist}[s][\ ] \leftarrow \text{MOORE}(V, E, w, s)$ 
  abort if MOORE found a negative cycle

  for every edge  $(u, v) \in E$ 
     $w'(u \rightarrow v) \leftarrow \text{dist}[s][u] + w(u \rightarrow v) - \text{dist}[v][s]$ 
  for every vertex  $v \in V$ 
     $\text{dist}[v][\ ] \leftarrow \text{DIJKSTRA}(V, E, w', v)$ 

```

The algorithm spends  $\Theta(V)$  time adding the artificial start vertex  $s$ ,  $\Theta(VE)$  time running MOORE,  $O(E)$  time reweighting the graph, and then  $\Theta(VE + V^2 \log V)$  running  $V$  passes of Dijkstra's algorithm. The overall running time is  $\Theta(VE + V^2 \log V)$ .

## 12.5 Dynamic Programming

There's a completely different solution to the all-pairs shortest path problem that uses dynamic programming instead of a single-source algorithm. For *dense* graphs where  $E = \Omega(V^2)$ , the dynamic programming approach gives the same running time as Johnson's algorithm, but with a much simpler algorithm. In particular, the new algorithm avoids Dijkstra's algorithm, which gets its efficiency from Fibonacci heaps, which are rather easy to screw up in the implementation.

To get a dynamic programming algorithm, we first need to come up with a recursive formulation of the problem. If we try to recursively define  $\text{dist}(u, v)$ , we might get something like this:

$$\text{dist}(u, v) = \begin{cases} 0 & \text{if } u = v \\ \min_x (\text{dist}(u, x) + w(x \rightarrow v)) & \text{otherwise} \end{cases}$$

In other words, to find the shortest path from  $u$  to  $v$ , try all possible predecessors  $x$ , compute the shortest path from  $u$  to  $x$ , and then add the last edge  $x \rightarrow v$ . **Unfortunately, this recurrence doesn't work!** In order to compute  $\text{dist}(u, v)$ , we first have to compute  $\text{dist}(u, x)$  for every other vertex  $x$ , but to compute any  $\text{dist}(u, x)$ , we first need to compute  $\text{dist}(u, v)$ . We're stuck in an infinite loop!

To avoid this circular dependency, we need some additional parameter that decreases at each recursion, eventually reaching zero at the base case. One possibility is to include the number of edges in the shortest path as this third magic parameter. So let's define  $\text{dist}(u, v, k)$  to be the length of the shortest path from  $u$  to  $v$  that uses *at most*  $k$  edges. Since we know that the shortest path between any two vertices has at most  $V - 1$  vertices, what we're really trying to compute is  $\text{dist}(u, v, V - 1)$ .

After a little thought, we get the following recurrence.

$$\text{dist}(u, v, k) = \begin{cases} 0 & \text{if } u = v \\ \infty & \text{if } k = 0 \text{ and } u \neq v \\ \min_x (\text{dist}(u, x, k - 1) + w(x \rightarrow v)) & \text{otherwise} \end{cases}$$

Just like last time, the recurrence tries all possible predecessors of  $v$  in the shortest path, but now the recursion actually bottoms out when  $k = 0$ .

Now it's not difficult to turn this recurrence into a dynamic programming algorithm. Even before we write down the algorithm, though, we can tell that its running time will be  $\Theta(V^4)$  simply because recurrence has four variables— $u$ ,  $v$ ,  $k$ , and  $x$ —each of which can take on  $V$  different values. Except for the base cases, the algorithm itself is just four nested for loops. To make the algorithm a little shorter, let's assume that  $w(v \rightarrow v) = 0$  for every vertex  $v$ .

```

DYNAMICPROGRAMMINGAPSP( $V, E, w$ ):
  for all vertices  $u \in V$ 
    for all vertices  $v \in V$ 
      if  $u = v$ 
         $\text{dist}[u][v][0] \leftarrow 0$ 
      else
         $\text{dist}[u][v][0] \leftarrow \infty$ 
    for  $k \leftarrow 1$  to  $V - 1$ 
      for all vertices  $u \in V$ 
        for all vertices  $v \in V$ 
           $\text{dist}[u][v][k] \leftarrow \infty$ 
          for all vertices  $x \in V$ 
            if  $\text{dist}[u][v][k] > \text{dist}[u][x][k - 1] + w(x \rightarrow v)$ 
               $\text{dist}[u][v][k] \leftarrow \text{dist}[u][x][k - 1] + w(x \rightarrow v)$ 

```

The last four lines actually evaluate the recurrence.

In fact, this algorithm is almost exactly the same as running Moore's algorithm once for every source vertex. The only difference is the innermost loop, which in Moore's algorithm would read "for all edges  $x \rightarrow v$ ". This simple change improves the running time to  $\Theta(V^2E)$ , assuming the graph is stored in an adjacency list.

## 12.6 Divide and Conquer

But we can make a more significant improvement. The recurrence we just used broke the shortest path into a slightly shorter path and a single edge, by considering all predecessors. Instead, let's break it into two shorter paths at the middle vertex on the path. This idea gives us a different recurrence for  $\text{dist}(u, v, k)$ . Once again, to simplify things, let's assume  $w(v \rightarrow v) = 0$ .

$$\text{dist}(u, v, k) = \begin{cases} w(u \rightarrow v) & \text{if } k = 1 \\ \min_x (\text{dist}(u, x, k/2) + \text{dist}(x, v, k/2)) & \text{otherwise} \end{cases}$$

This recurrence only works when  $k$  is a power of two, since otherwise we might try to find the shortest path with a fractional number of edges! But that's not really a problem, since  $\text{dist}(u, v, 2^{\lceil \lg V \rceil})$  gives us the overall shortest distance from  $u$  to  $v$ . Notice that we use the base case  $k = 1$  instead of  $k = 0$ , since we can't use half an edge.

Once again, a dynamic programming solution is straightforward. Even before we write down the algorithm, we can tell the running time is  $\Theta(V^3 \log V)$ —we consider  $V$  possible values of  $u$ ,  $v$ , and  $x$ , but only  $\lceil \lg V \rceil$  possible values of  $k$ .

```

FASTDYNAMICPROGRAMMINGAPSP( $V, E, w$ ):
  for all vertices  $u \in V$ 
    for all vertices  $v \in V$ 
       $\text{dist}[u][v][0] \leftarrow w(u \rightarrow v)$ 
  for  $i \leftarrow 1$  to  $\lceil \lg V \rceil$        $\langle\langle k = 2^i \rangle\rangle$ 
    for all vertices  $u \in V$ 
      for all vertices  $v \in V$ 
         $\text{dist}[u][v][i] \leftarrow \infty$ 
        for all vertices  $x \in V$ 
          if  $\text{dist}[u][v][i] > \text{dist}[u][x][i-1] + \text{dist}[x][v][i-1]$ 
             $\text{dist}[u][v][i] \leftarrow \text{dist}[u][x][i-1] + \text{dist}[x][v][i-1]$ 

```

### 12.7 Aside: ‘Funny’ Matrix Multiplication

There is a very close connection between computing shortest paths in a directed graph and computing powers of a square matrix. Compare the following algorithm for multiplying two  $n \times n$  matrices  $A$  and  $B$  with the inner loop of our first dynamic programming algorithm. (I’ve changed the variable names in the second algorithm slightly to make the similarity clearer.)

```

MATRIXMULTIPLY( $A, B$ ):
  for  $i \leftarrow 1$  to  $n$ 
    for  $j \leftarrow 1$  to  $n$ 
       $C[i][j] \leftarrow 0$ 
      for  $k \leftarrow 1$  to  $n$ 
         $C[i][j] \leftarrow C[i][j] + A[i][k] \cdot B[k][j]$ 

```

```

APSPINNERLOOP:
  for all vertices  $u$ 
    for all vertices  $v$ 
       $D'[u][v] \leftarrow \infty$ 
      for all vertices  $x$ 
         $D'[u][v] \leftarrow \min\{D'[u][v], D[u][x] + w[x][v]\}$ 

```

The *only* difference between these two algorithms is that we use addition instead of multiplication and minimization instead of addition. For this reason, the shortest path inner loop is often referred to as ‘funny’ matrix multiplication.

DYNAMICPROGRAMMINGAPSP is the standard iterative algorithm for computing the  $(V - 1)$ th ‘funny power’ of the weight matrix  $w$ . The first set of for loops sets up the ‘funny identity matrix’, with zeros on the main diagonal and infinity everywhere else. Then each iteration of the second main for loop computes the next ‘funny power’. FASTDYNAMICPROGRAMMINGAPSP replaces this iterative method for computing powers with repeated squaring, exactly like we saw at the beginning of the semester. The fast algorithm is simplified slightly by the fact that unless there are negative cycles, every ‘funny power’ after the  $V$ th is the same.

There are faster methods for multiplying matrices, similar to Karatsuba’s divide-and-conquer algorithm for multiplying integers. (See ‘Strassen’s algorithm’ in CLR.) Unfortunately, these algorithms use subtraction, and there’s no ‘funny’ equivalent of subtraction. (What’s the inverse operation for min?) So at least for general graphs, there seems to be no way to speed up the inner loop of our dynamic programming algorithms.

Fortunately, this isn't true. There is a beautiful randomized algorithm, due to Noga Alon, Zvi Galil, Oded Margalit\*, and Moni Naor,<sup>1</sup> that computes all-pairs shortest paths in undirected graphs in  $O(M(V) \log^2 V)$  expected time, where  $M(V)$  is the time to multiply two  $V \times V$  integer matrices. A simplified version of this algorithm for *unweighted* graphs, due to Raimund Seidel<sup>2</sup>, appears in the current homework.

## 12.8 Floyd and Warshall's Algorithm

Our fast dynamic programming algorithm is still a factor of  $O(\log V)$  slower than Johnson's algorithm. A different formulation due to Floyd and Warshall removes this logarithmic factor. Their insight was to use a different third parameter in the recurrence.

Number the vertices arbitrarily from 1 to  $V$ , and define  $\text{dist}(u, v, r)$  to be the length of the shortest path from  $u$  to  $v$ , where all the *intermediate* vertices (if any) are numbered  $r$  or less. If  $r = 0$ , we aren't allowed to use any intermediate vertices, so the shortest legal path from  $u$  to  $v$  is just the edge (if any) from  $u$  to  $v$ . If  $r > 0$ , then either the shortest legal path from  $u$  to  $v$  goes through vertex  $r$  or it doesn't. We get the following recurrence:

$$\text{dist}(u, v, r) = \begin{cases} w(u \rightarrow v) & \text{if } r = 0 \\ \min \{ \text{dist}(u, v, r-1), \text{dist}(u, r, r-1) + \text{dist}(r, v, r-1) \} & \text{otherwise} \end{cases}$$

We need to compute the shortest path distance from  $u$  to  $v$  with no restrictions, which is just  $\text{dist}(u, v, V)$ .

Once again, we should immediately see that a dynamic programming algorithm that implements this recurrence will run in  $\Theta(V^3)$  time: three variables appear in the recurrence ( $u$ ,  $v$ , and  $r$ ), each of which can take on  $V$  possible values. Here's one way to do it:

```

FLOYDWARSHALL( $V, E, w$ ):
  for  $u \leftarrow 1$  to  $V$ 
    for  $v \leftarrow 1$  to  $V$ 
       $\text{dist}[u][v][0] \leftarrow w(u \rightarrow v)$ 
  for  $r \leftarrow 1$  to  $V$ 
    for  $u \leftarrow 1$  to  $V$ 
      for  $v \leftarrow 1$  to  $V$ 
        if  $\text{dist}[u][v][r-1] < \text{dist}[u][r][r-1] + \text{dist}[r][v][r-1]$ 
           $\text{dist}[u][v][r] \leftarrow \text{dist}[u][v][r-1]$ 
        else
           $\text{dist}[u][v][r] \leftarrow \text{dist}[u][r][r-1] + \text{dist}[r][v][r-1]$ 

```

<sup>1</sup>N. Alon, Z. Galil, O. Margalit\*, and M. Naor. Witnesses for Boolean matrix multiplication and for shortest paths. *Proc. 33rd FOCS* 417-426, 1992. See also N. Alon, Z. Galil, O. Margalit\*. On the exponent of the all pairs shortest path problem. *Journal of Computer and System Sciences* 54(2):255-262, 1997.

<sup>2</sup>R. Seidel. On the all-pairs-shortest-path problem in unweighted undirected graphs. *Journal of Computer and System Sciences*, 51(3):400-403, 1995. This is one of the few algorithms papers where (in the conference version at least) the algorithm is completely described and analyzed *in the abstract* of the paper.



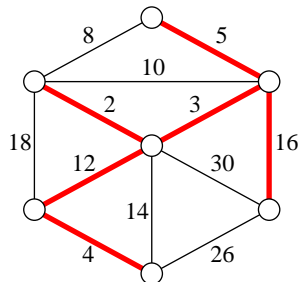
## 13 Minimum Spanning Trees (October 29)

### 13.1 Introduction

Suppose we are given a connected, undirected, *weighted* graph. This is a graph  $G = (V, E)$  together with a function  $w: E \rightarrow \mathbb{R}$  that assigns a *weight*  $w(e)$  to each edge  $e$ . For this lecture, we'll assume that the weights are real numbers. Our task is to find the *minimum spanning tree* of  $G$ , *i.e.*, the spanning tree  $T$  minimizing the function

$$w(T) = \sum_{e \in T} w(e).$$

To keep things simple, I'll assume that all the edge weights are distinct:  $w(e) \neq w(e')$  for any pair of edges  $e$  and  $e'$ . Distinct weights guarantee that the minimum spanning tree of the graph is unique. Without this condition, there may be several different minimum spanning trees. For example, if all the edges have weight 1, then *every* spanning tree is a minimum spanning tree with weight  $V - 1$ .



A weighted graph and its minimum spanning tree.

If we have an algorithm that assumes the edge weights are unique, we can still use it on graphs where multiple edges have the same weight, as long as we have a consistent method for breaking ties. One way to break ties consistently is to use the following algorithm in place of a simple comparison. `SHORTEREDGE` takes as input four integers  $i, j, k, l$ , and decides which of the two edges  $(i, j)$  and  $(k, l)$  has ‘smaller’ weight.

```

SHORTEREDGE( $i, j, k, l$ )
  if  $w(i, j) < w(k, l)$  return  $(i, j)$ 
  if  $w(i, j) > w(k, l)$  return  $(k, l)$ 
  if  $\min(i, j) < \min(k, l)$  return  $(i, j)$ 
  if  $\min(i, j) > \min(k, l)$  return  $(k, l)$ 
  if  $\max(i, j) < \max(k, l)$  return  $(i, j)$ 
   $\langle\langle$  if  $\max(i, j) < \max(k, l)$   $\rangle\rangle$  return  $(k, l)$ 

```

### 13.2 The Only MST Algorithm

There are several different ways to compute minimum spanning trees, but really they are all instances of the following generic algorithm. The situation is similar to the previous lecture, where we saw that depth-first search and breadth-first search were both instances of a single generic traversal algorithm.

The generic MST algorithm maintains an acyclic subgraph  $F$  of the input graph  $G$ , which we will call an *intermediate spanning forest*.  $F$  is a subgraph of the minimum spanning tree of  $G$ ,

and every component of  $F$  is a minimum spanning tree of its vertices. Initially,  $F$  consists of  $n$  one-node trees. The generic MST algorithm merges trees together by adding certain edges between them. When the algorithm halts,  $F$  consists of a single  $n$ -node tree, which must be the minimum spanning tree. Obviously, we have to be careful about *which* edges we add to the evolving forest, since not every edge is in the eventual minimum spanning tree.

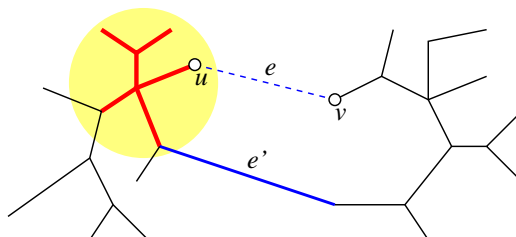
The intermediate spanning forest  $F$  induces two special types of edges. An edge is *useless* if it is not an edge of  $F$ , but both its endpoints are in the same component of  $F$ . For each component of  $F$ , we associate a *safe* edge—the minimum-weight edge with exactly one endpoint in that component.<sup>1</sup> Different components might or might not have different safe edges. Some edges are neither safe nor useless—we call these edges *undecided*.

All minimum spanning tree algorithms are based on two simple observations.

**Lemma 1.** *The minimum spanning tree contains every safe edge and no useless edges.*

**Proof:** Let  $T$  be the minimum spanning tree. Suppose  $F$  has a ‘bad’ component whose safe edge  $e = (u, v)$  is not in  $T$ . Since  $T$  is connected, it contains a unique path from  $u$  to  $v$ , and at least one edge  $e'$  on this path has exactly one endpoint in the bad component. Removing  $e'$  from the minimum spanning tree and adding  $e$  gives us a new spanning tree. Since  $e$  is the bad component’s safe edge, we have  $w(e') > w(e)$ , so the new spanning tree has smaller total weight than  $T$ . But this is impossible— $T$  is the *minimum* spanning tree. So  $T$  must contain every safe edge.

Adding any useless edge to  $F$  would introduce a cycle. □



Proving that every safe edge is in the minimum spanning tree. The ‘bad’ component of  $F$  is highlighted.

So our generic minimum spanning tree algorithm repeatedly adds one or more safe edges to the evolving forest  $F$ . Whenever we add new edges to  $F$ , some undecided edges become safe, and others become useless. To specify a particular algorithm, we must decide which safe edges to add, and how to identify new safe and new useless edges, at each iteration of our generic template.

### 13.3 Borůvka’s Algorithm

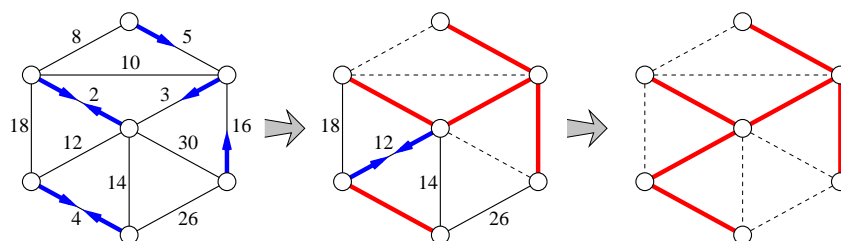
The oldest and possibly simplest minimum spanning tree algorithm was discovered by Borůvka in 1926, long before computers even existed, and practically before the invention of graph theory!<sup>2</sup> The algorithm was rediscovered by Choquet in 1938; again by Florek, Lukaziewicz, Perkal, Stienhaus, and Zubrzycki in 1951; and again by Sollin some time in the early 1960s.

The Borůvka/Choquet/Florek/Lukaziewicz/Perkal/Stienhaus/Zubrzycki/Sollin algorithm can be summarized in one line:

<sup>1</sup>This is actually a special case of a more general definition: For any partition of  $F$  into two subforests, the minimum-weight edge with one endpoint in each subforest is light. A few minimum spanning tree algorithms require this more general definition, but we won’t talk about them here.

<sup>2</sup>The first book on graph theory, written by D. König, was published in 1936. Leonard Euler published his famous theorem about the bridges of Königsburg (HW3, problem 2) in 1736. Königsburg was not named after *that* König.

BORŮVKA: Add all the safe edges and recurse.



Borůvka's algorithm run on the example graph. Thick edges are in  $F$ . Arrows point along each component's safe edge. Dashed edges are useless.

At the beginning of each phase of the Borůvka algorithm, each component elects an arbitrary 'leader' node. The simplest way to hold these elections is a depth-first search of  $F$ ; the first node we visit in any component is that component's leader. Once the leaders are elected, we find the safe edges for each component, essentially by brute force. Finally, we add these safe edges to  $F$ .

BORŮVKA( $V, E$ ):

$F = (V, \emptyset)$

while  $F$  has more than one component

    choose leaders using DFS

    FINDSAFEEDGES( $V, E$ )

    for each leader  $\bar{v}$

        add safe( $\bar{v}$ ) to  $F$

FINDSAFEEDGES( $V, E$ ):

for each leader  $\bar{v}$

    safe( $\bar{v}$ )  $\leftarrow \infty$

for each edge  $(u, v) \in E$

$\bar{u} \leftarrow \text{leader}(u)$

$\bar{v} \leftarrow \text{leader}(v)$

    if  $\bar{u} \neq \bar{v}$

        if  $w(u, v) < w(\text{safe}(\bar{u}))$

            safe( $\bar{u}$ )  $\leftarrow (u, v)$

        if  $w(u, v) < w(\text{safe}(\bar{v}))$

            safe( $\bar{v}$ )  $\leftarrow (u, v)$

Each call to FINDSAFEEDGES takes  $O(E)$  time, since it examines every edge. Since the graph is connected, it has at most  $E + 1$  vertices. Thus, each iteration of the while loop in BORŮVKA takes  $O(E)$  time, assuming the graph is represented by an adjacency list. Each iteration also reduces the number of components of  $F$  by at least a factor of two—the worst case occurs when the components coalesce in pairs. Since there are initially  $V$  components, the while loop iterates  $O(\log V)$  times. Thus, the overall running time of Borůvka's algorithm is  $O(E \log V)$ .

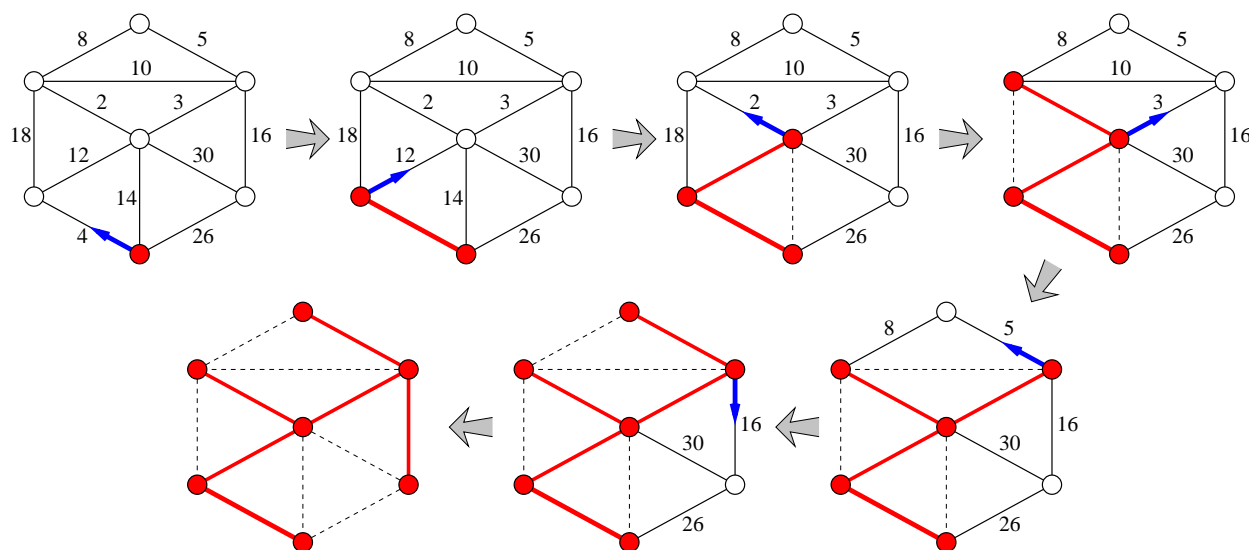
Despite its relatively obscure origin, early algorithms researchers were aware of Borůvka's algorithm, but dismissed it as being "too complicated"! As a result, despite its simplicity and efficiency, Borůvka's algorithm is rarely mentioned in algorithms and data structures textbooks.

### 13.4 Jarník's ('Prim's') Algorithm

The next oldest minimum spanning tree algorithm was discovered by the Vojtěch Jarník in 1930, but it is usually called Prim's algorithm. Prim independently rediscovered the algorithm in 1956 and gave a much more detailed description than Jarník. The algorithm was rediscovered again in 1958 by Dijkstra, but he already had an algorithm named after him. Such is fame in academia.

In Jarník's algorithm, the forest  $F$  contains only one nontrivial component  $T$ ; all the other components are isolated vertices. Initially,  $T$  consists of an arbitrary vertex of the graph. The algorithm repeats the following step until  $T$  spans the whole graph:

JARNÍK: Find  $T$ 's safe edge and add it to  $T$ .



Jarník's algorithm run on the example graph, starting with the bottom vertex.

At each stage, thick edges are in  $T$ , an arrow points along  $T$ 's safe edge, and dashed edges are useless.

To implement Jarník's algorithm, we keep all the edges adjacent to  $T$  in a heap. When we pull the minimum-weight edge off the heap, we first check whether both of its endpoints are in  $T$ . If not, we add the edge to  $T$  and then add the new neighboring edges to the heap. In other words, Jarník's algorithm is just another instance of the generic graph traversal algorithm we saw last time, using a heap as the 'bag'! If we implement the algorithm this way, its running time is  $O(E \log E) = O(E \log V)$ .

However, we can speed up the implementation by observing that the graph traversal algorithm visits each vertex only once. Rather than keeping edges in the heap, we can keep a heap of vertices, where the key of each vertex  $v$  is the length of the minimum-weight edge between  $v$  and  $T$  (or  $\infty$  if there is no such edge). Each time we add a new edge to  $T$ , we may need to decrease the key of some neighboring vertices.

To make the description easier, we break the algorithm into two parts. JARNÍKINIT initializes the vertex heap. JARNÍKLOOP is the main algorithm. The input consists of the vertices and edges of the graph, plus the start vertex  $s$ .

JARNÍK( $V, E, s$ ):  
 JARNÍKINIT( $V, E, s$ )  
 JARNÍKLOOP( $V, E, s$ )

JARNÍKINIT( $V, E, s$ ):  
 for each vertex  $v \in V \setminus \{s\}$   
   if  $(v, s) \in E$   
     edge( $v$ )  $\leftarrow (v, s)$   
     key( $v$ )  $\leftarrow w(v, s)$   
   else  
     edge( $v$ )  $\leftarrow \text{NULL}$   
     key( $v$ )  $\leftarrow \infty$   
 INSERT( $v$ )

JARNÍKLOOP( $V, E, s$ ):  
 $T \leftarrow (\{s\}, \emptyset)$   
 for  $i \leftarrow 1$  to  $|V| - 1$   
    $v \leftarrow \text{EXTRACTMIN}$   
   add  $v$  and edge( $v$ ) to  $T$   
   for each edge  $(u, v) \in E$   
     if  $u \notin T$  and key( $u$ )  $> w(u, v)$   
       edge( $u$ )  $\leftarrow (u, v)$   
       DECREASEKEY( $u, w(u, v)$ )

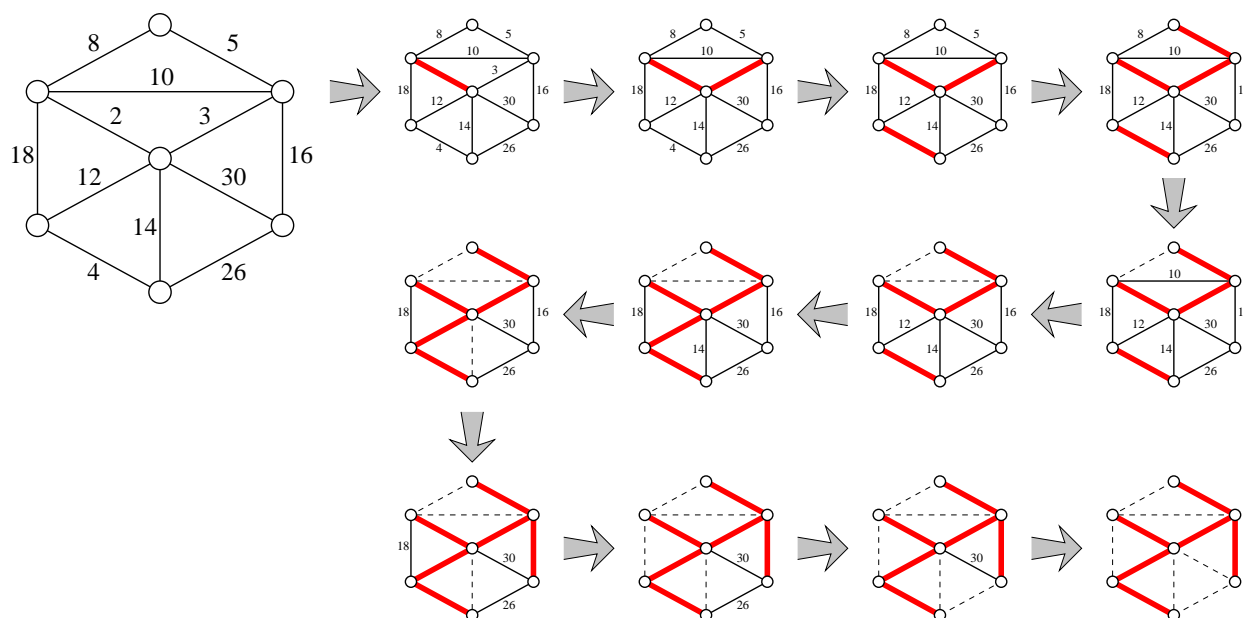
The running time of JARNÍK is dominated by the cost of the heap operations INSERT, EXTRACTMIN, and DECREASEKEY. INSERT and EXTRACTMIN are each called  $O(V)$  times once for each vertex except  $s$ , and DECREASEKEY is called  $O(E)$  times, at most twice for each edge. If we use a Fibonacci heap, the amortized costs of INSERT and DECREASEKEY is  $O(1)$ , and the amortized cost of EXTRACTMIN is  $O(\log n)$ . Thus, the overall running time of JARNÍK is  $O(E + V \log V)$ . This is faster than Borůvka's algorithm unless  $E = O(V)$ .

### 13.5 Kruskal's Algorithm

The last minimum spanning tree algorithm I'll discuss was discovered by Kruskal in 1956.

KRUSKAL: Scan all edges in increasing weight order; if an edge is safe, add it to  $F$ .

Since we examine the edges in order from lightest to heaviest, any edge we examine is safe if and only if its endpoints are in different components of the forest  $F$ . To prove this, suppose the edge  $e$  joins two components  $A$  and  $B$  but is not safe. Then there would be a lighter edge  $e'$  with exactly one endpoint in  $A$ . But this is impossible, because (inductively) any previously examined edge has both endpoints in the same component of  $F$ .



Kruskal's algorithm run on the example graph. Thick edges are in  $F$ . Dashed edges are useless.

Just as in Borůvka's algorithm, each component of  $F$  has a 'leader' node. An edge joins two components of  $F$  if and only if the two endpoints have different leaders. But unlike Borůvka's algorithm, we do not recompute leaders from scratch every time we add an edge. Instead, when two components are joined, the two leaders duke it out in a nationally-televised no-holds-barred steel-cage grudge match.<sup>3</sup> One of the two emerges victorious as the leader of the new larger component. More formally, we will use our earlier algorithms for the UNION-FIND problem, where the vertices are the elements and the components of  $F$  are the sets. Here's a more formal description of the algorithm:

<sup>3</sup>Live at the Assembly Hall! Only \$49.95 on Pay-Per-View!

```
KRUSKAL( $V, E$ ):
  sort  $E$  by weight
   $F \leftarrow \emptyset$ 
  for each vertex  $v \in V$ 
    MAKESET( $v$ )
  for  $i \leftarrow 1$  to  $|E|$ 
     $(u, v) \leftarrow$   $i$ th lightest edge in  $E$ 
    if FIND( $u$ )  $\neq$  FIND( $v$ )
      UNION( $u, v$ )
      add  $(u, v)$  to  $F$ 
  return  $F$ 
```

In our case, the sets are components of  $F$ , and  $n = V$ . Kruskal's algorithm performs  $O(E)$  FIND operations, two for each edge in the graph, and  $O(V)$  UNION operations, one for each edge in the minimum spanning tree. Using union-by-rank and path compression allows us to perform each UNION or FIND in  $O(\alpha(E, V))$  time, where  $\alpha$  is the not-quite-constant inverse-Ackerman function. So ignoring the cost of sorting the edges, the running time of this algorithm is  $O(E\alpha(E, V))$ .

We need  $O(E \log E) = O(E \log V)$  additional time just to sort the edges. Since this is bigger than the time for the UNION-FIND data structure, the overall running time of Kruskal's algorithm is  $O(E \log V)$ , exactly the same as Borůvka's algorithm, or Jarník's algorithm with a normal (non-Fibonacci) heap.

Blech! Ack! Oop! THPPFFT!

— Bill the Cat, “Bloom County” (1980)

## 14 Fast Fourier Transforms (November 7)

### 14.1 Polynomials

In this lecture we’ll talk about algorithms for manipulating *polynomials*: functions of one variable built from additions subtractions, and multiplications (but no divisions). The most common representation for a polynomial  $p(x)$  is as a sum of weighted powers of a variable  $x$ :

$$p(x) = \sum_{j=0}^n a_j x^j.$$

The numbers  $a_j$  are called *coefficients*. The *degree* of the polynomial is the largest power of  $x$ ; in the example above, the degree is  $n$ . Any polynomial of degree  $n$  can be specified by a sequence of  $n + 1$  coefficients. Some of these coefficients may be zero, but not the  $n$ th coefficient, because otherwise the degree would be less than  $n$ .

Here are three of the most common operations that are performed with polynomials:

- **Evaluate:** Give a polynomial  $p$  and a number  $x$ , compute the number  $p(x)$ .
- **Add:** Give two polynomials  $p$  and  $q$ , compute a polynomial  $r = p + q$ , so that  $r(x) = p(x) + q(x)$  for all  $x$ . If  $p$  and  $q$  both have degree  $n$ , then their sum  $p + q$  also has degree  $n$ .
- **Multiply:** Give two polynomials  $p$  and  $q$ , compute a polynomial  $r = p \cdot q$ , so that  $r(x) = p(x) \cdot q(x)$  for all  $x$ . If  $p$  and  $q$  both have degree  $n$ , then their product  $p \cdot q$  has degree  $2n$ .

Suppose we represent a polynomial of degree  $n$  as an array of  $n + 1$  coefficients  $P[0..n]$ , where  $P[j]$  is the coefficient of the  $x^j$  term. We learned simple algorithms for all three of these operations in high-school algebra:

<p>EVALUATE(<math>P[0..n], x</math>):</p> $X \leftarrow 1 \quad \langle\langle X = x^j \rangle\rangle$ $y \leftarrow 0$ for $j \leftarrow 0$ to $n$ $y \leftarrow y + P[j] \cdot X$ $X \leftarrow X \cdot x$ return $y$
--

<p>ADD(<math>P[0..n], Q[0..n]</math>):</p> for $j \leftarrow 0$ to $n$ $R[j] \leftarrow P[j] + Q[j]$ return $R[0..n]$
---

<p>MULTIPLY(<math>P[0..n], Q[0..m]</math>):</p> for $j \leftarrow 0$ to $n + m$ $R[j] \leftarrow 0$ for $j \leftarrow 0$ to $n$ for $k \leftarrow 0$ to $m$ $R[j + k] \leftarrow P[j] \cdot Q[k]$ return $R[0..n + m]$
---

EVALUATE uses  $O(n)$  arithmetic operations.<sup>1</sup> This is the best you can do in theory, but we can cut the number of multiplications in half using *Horner’s rule*:

$$p(x) = a_0 + x(a_1 + x(a_2 + \dots + xa_n)).$$

<sup>1</sup>I’m going to assume in this lecture that each arithmetic operation takes  $O(1)$  time. This may not be true in practice; in fact, one of the most powerful applications of FFTs is fast *integer* multiplication. One of the fastest integer multiplication algorithms, due to Schönhage and Strassen, multiplies two  $n$ -bit binary numbers in  $O(n \log n \log \log n \log \log \log n \log \log \log \log n \dots)$  time. The algorithm uses an  $n$ -element Fast Fourier Transform, which requires several  $O(\log n)$ -bit integer multiplications. These smaller multiplications are carried out recursively (of course!), which leads to the cascade of logs in the running time. Needless to say, this is a can of worms.

```

HORNER( $P[0..n], x$ ):
 $y \leftarrow P[n]$ 
for  $i \leftarrow n - 1$  downto 0
     $y \leftarrow x \cdot y + P[i]$ 
return  $y$ 

```

The addition algorithm also runs in  $O(n)$  time, and this is clearly the best we can do.

The multiplication algorithm, however, runs in  $O(n^2)$  time. In the very first lecture, we saw a divide and conquer algorithm (due to Karatsuba) for multiplying two  $n$ -bit integers in only  $O(n^{\lg 3})$  steps; precisely the same algorithm can be applied here. Even cleverer divide-and-conquer strategies lead to multiplication algorithms whose running times are arbitrarily close to linear— $O(n^{1+\epsilon})$  for your favorite value  $\epsilon > 0$ —but with great cleverness comes great confusion. These algorithms are difficult to understand, even more difficult to implement correctly, and not worth the trouble in practice thanks to large constant factors.

## 14.2 Alternate Representations: Roots and Samples

Part of what makes multiplication so much harder than the other two operations is our input representation. Coefficient vectors are the most common representation for polynomials, but there are at least two other useful representations.

The first exploits the fundamental theorem of algebra: Every polynomial  $p$  of degree  $n$  has  $n$  roots  $r_1, r_2, \dots, r_n$  such that  $p(r_j) = 0$  for all  $j$ . Some of these roots may be irrational; some of these roots may be complex; and some of these roots may be repeated. Despite these complications, we do get a unique representation of any polynomial of the form

$$p(x) = s \prod_{j=1}^n (x - r_j)$$

where the  $r_j$ 's are the roots and  $s$  is a scale factor. Once again, to represent a polynomial of degree  $n$ , we need a list of  $n + 1$  numbers: one scale factor and  $n$  roots.

Given a polynomial in root representation, we can clearly evaluate it in  $O(n)$  time. Given two polynomials in root representation, we can easily multiply them in  $O(n)$  time by multiplying their scale factors and just concatenating the two root sequences; in fact, if we don't care about keeping the old polynomials around, we can compute their product in  $O(1)$  time! Unfortunately if we want to *add* two polynomials in root representation, we're pretty much out of luck; there's essentially no correlation between the roots of  $p$ , the roots of  $q$ , and the roots of  $p + q$ . We could convert the polynomials to the more familiar coefficient representation first—this takes  $O(n^2)$  time using the high-school algorithms—but there's no easy way to convert the answer back. In fact, given a polynomial in coefficient form, it's usually *impossible* to compute its roots exactly.<sup>2</sup>

Our third representation for polynomials comes from the following consequence of the fundamental theorem of algebra. Given a list of  $n + 1$  pairs  $\{(x_0, y_0), (x_1, y_1), \dots, (x_n, y_n)\}$ , there is *exactly one* polynomial  $p$  of degree  $n$  such that  $p(x_j) = y_j$  for all  $j$ . This is just a generalization of the fact that any two points determine a unique line, since a line is (the graph of) a polynomial of degree 1. We say that the polynomial  $p$  *interpolates* the points  $(x_j, y_j)$ . As long as we agree on the sample locations  $x_j$  in advance, we once again need exactly  $n + 1$  numbers to represent a polynomial of degree  $n$ .

---

<sup>2</sup>This is where numerical analysis comes from.



Adding or multiplying two polynomials in this sample representation is easy, as long as they use the same sample locations  $x_j$ . To add the polynomials, just add their sample values. To multiply two polynomials, just multiply their sample values; however, if we're multiplying two polynomials of degree  $n$ , we need to *start* with  $2n + 1$  sample values for each polynomial, since that's how many we need to uniquely represent the product polynomial. Both algorithms run in  $O(n)$  time.

Unfortunately, evaluating a polynomial in this representation is no longer trivial. The following formula, due to Lagrange, allows us to compute the value of any polynomial of degree  $n$  at any point, given a set of  $n + 1$  samples.

$$p(x) = \sum_{j=0}^{n-1} \left( y_j \frac{\prod_{k \neq j} (x - x_k)}{\prod_{k \neq j} (x_j - x_k)} \right) = \sum_{j=0}^{n-1} \left( \frac{y_j}{\prod_{k \neq j} (x_j - x_k)} \prod_{k \neq j} (x - x_k) \right)$$

Hopefully it's clear that formula actually describes a polynomial, since each term in the rightmost sum is written as a scaled product of monomials. It's also not hard to check that  $p(x_j) = y_j$  for all  $j$ . As I mentioned earlier, the fact that this is *the only* polynomial that interpolates the points  $\{(x_j, y_j)\}$  is an easy consequence of the fundamental theorem of algebra. We can easily transform this formula into an  $O(n^2)$ -time algorithm.

We find ourselves in the following frustrating situation. We have three representations for polynomials and three basic operations. Each representation allows us to almost trivially perform a different pair of operations in linear time, but the third takes at least quadratic time, if it can be done at all!

	evaluate	add	multiply
coefficients	$O(n)$	$O(n)$	$O(n^2)$
roots + scale	$O(n)$	$\infty$	$O(n)$
samples	$O(n^2)$	$O(n)$	$O(n)$

### 14.3 Converting Between Representations?

What we need are fast algorithms to convert quickly from one representation to another. That way, when we need to perform an operation that's hard for our default representation, we can switch to a different representation that makes the operation easy, perform that operation, and then switch back. This strategy immediately rules out the root representation, since (as I mentioned earlier) finding roots of polynomials is impossible in general, at least if we're interested in exact results.

So how do we convert from coefficients to samples and back? Clearly, once we choose our sample positions  $x_j$ , we can compute each sample value  $y_j = p(x_j)$  in  $O(n)$  time from the coefficients. So we can convert a polynomial of degree  $n$  from coefficients to samples in  $O(n^2)$  time. The Lagrange formula gives us an explicit conversion algorithm from the sample representation back to the more familiar coefficient representation. If we use the naïve algorithms for adding and multiplying polynomials (in coefficient form), this conversion takes  $O(n^3)$  time.

This looks pretty bad, until we realize there's a degree of freedom we haven't exploited yet. ***Whenever we convert from coefficients to samples, we get to choose the sample points!*** Our slow algorithms may be slow only because we're trying to be too general. Perhaps, if we choose a set of sample points with just the right kind of recursive structure, we can do the conversion more quickly. In fact, there is a set of sample points that's perfect for the job.

### 14.4 The Discrete Fourier Transform

Given a polynomial of degree  $n - 1$ , we'd like to find  $n$  sample points that are somehow as symmetric as possible. The most natural choice for those  $n$  points are the  *$n$ th roots of unity*; these are the

roots of the polynomial  $x^n - 1 = 0$ . These  $n$  roots are spaced exactly evenly around the unit circle in the complex plane.<sup>3</sup> Every  $n$ th root of unity is a power of the *primitive root*

$$\omega_n = e^{2\pi i/n} = \cos \frac{2\pi}{n} + i \sin \frac{2\pi}{n}.$$

A typical  $n$ th root of unity has the form

$$\omega_n^j = e^{(2\pi i/n)j} = \cos \left( \frac{2\pi}{n} j \right) + i \sin \left( \frac{2\pi}{n} j \right).$$

These complex numbers have several useful properties for any integers  $n$  and  $k$ :

- There are only  $n$  different  $n$ th roots of unity:  $\omega_n^k = \omega_n^{k \bmod n}$ .
- If  $n$  is even, then  $\omega_n^{k+n/2} = -\omega_n^k$ ; in particular,  $\omega_n^{n/2} = -\omega_n^0 = -1$ .
- $1/\omega_n^k = \omega_n^{-k} = \overline{\omega_n^k} = (\overline{\omega_n})^k$ , where the bar represents complex conjugation:  $\overline{a + bi} = a - bi$
- $\omega_n = \omega_{kn}^k$ . Thus, every  $n$ th root of unity is also a  $(kn)$ th root of unity.

If we sample a polynomial of degree  $n - 1$  at the  $n$ th roots of unity, the resulting list of sample values is called the *discrete Fourier transform* of the polynomial (or more formally, of the coefficient vector). Thus, given an array  $P[0..n - 1]$  of coefficients, the discrete Fourier transform computes a new vector  $P^*[0..n - 1]$  where

$$P^*[j] = p(\omega_n^j) = \sum_{k=0}^{n-1} P[k] \cdot \omega_n^{jk}$$

We can obviously compute  $P^*$  in  $O(n^2)$  time, but the structure of the  $n$ th roots of unity lets us do better. But before we describe that faster algorithm, let's think about how we might invert this transformation.

It's not hard to see that the discrete Fourier transform—in fact, any conversion from a vector of coefficients to a vector of sample values—is a *linear* transformation. The DFT just multiplies the coefficient vector by a matrix  $V$  to obtain the sample vector. Each entry in  $V$  is an  $n$ th root of unity; specifically,  $v_{jk} = \omega_n^{jk}$  for all  $j, k$ .

$$V = \begin{bmatrix} 1 & 1 & 1 & 1 & \cdots & 1 \\ 1 & \omega_n & \omega_n^2 & \omega_n^3 & \cdots & \omega_n^{n-1} \\ 1 & \omega_n^2 & \omega_n^4 & \omega_n^6 & \cdots & \omega_n^{2(n-1)} \\ 1 & \omega_n^3 & \omega_n^6 & \omega_n^9 & \cdots & \omega_n^{3(n-1)} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & \omega_n^{n-1} & \omega_n^{2(n-1)} & \omega_n^{3(n-1)} & \cdots & \omega_n^{(n-1)^2} \end{bmatrix}$$

<sup>3</sup>In this lecture,  $i$  always represents the square root of  $-1$ . Most computer scientists are used to thinking of  $i$  as an integer index into a sequence, an array, or a for-loop, but we obviously can't do that here. The physicist's habit of using  $j = \sqrt{-1}$  just delays the problem (how do physicists write quaternions?), and typographical tricks like  $I$  or  $\mathbf{i}$  or Mathematica's  $\mathbf{i}$  are just stupid.

To invert the discrete Fourier transform, we just have to multiply  $P^*$  by the inverse matrix  $V^{-1}$ . But this is almost the same as multiplying by  $V$  itself, because of the following amazing fact:

$$\boxed{V^{-1} = \overline{V}/n}$$

In other words, if  $W = V^{-1}$  then  $w_{jk} = \overline{v_{jk}}/n = \overline{\omega_n^{jk}}/n = \omega_n^{-jk}/n$ . It's not hard to prove this fact with a little linear algebra.

**Proof:** We just have to show that  $M = VW$  is the identity matrix. We can compute a single entry in this matrix as follows:

$$m_{jk} = \sum_{l=0}^{n-1} v_{jl} \cdot w_{lk} = \sum_{l=0}^{n-1} \omega_n^{jl} \cdot \overline{\omega_n^{lk}}/n = \frac{1}{n} \sum_{l=0}^{n-1} \omega_n^{j-lk} = \frac{1}{n} \sum_{l=0}^{n-1} (\omega_n^{j-k})^l$$

If  $j = k$ , then  $\omega_n^{j-k} = 1$ , so

$$m_{jk} = \frac{1}{n} \sum_{l=0}^{n-1} 1 = \frac{n}{n} = 1,$$

and if  $j \neq k$ , we have a geometric series

$$m_{jk} = \sum_{l=0}^{n-1} (\omega_n^{j-k})^l = \frac{(\omega_n^{j-k})^n - 1}{\omega_n^{j-k} - 1} = \frac{(\omega_n^n)^{j-k} - 1}{\omega_n^{j-k} - 1} = \frac{1^{j-k} - 1}{\omega_n^{j-k} - 1} = 0.$$

That's it! □

What this means for us computer scientists is that any algorithm for computing the discrete Fourier transform can be easily modified to compute the inverse transform as well.

### 14.5 Divide and Conquer

The structure of the matrix  $V$  also allows us to compute the discrete Fourier transform efficiently using a divide and conquer strategy. The basic structure of the algorithm is almost the same as MergeSort, and the  $O(n \log n)$  running time will ultimately follow from the same recurrence. The *Fast Fourier Transform* algorithm, popularized by Cooley and Tukey in 1965<sup>4</sup>, assumes that  $n$  is a power of two; if necessary, we can just pad the coefficient vector with zeros.

To get an idea of how the divide-and-conquer strategy works, let's look at the DFT matrixes for  $n = 8$ . To simplify notation, let  $\omega = \omega_8 = \sqrt{2}/2 + i\sqrt{2}/2$ .

$$\begin{pmatrix} 1 & 1 & 1 & 1 & 1 & 1 & 1 & 1 \\ 1 & \omega & \omega^2 & \omega^3 & \omega^4 & \omega^5 & \omega^6 & \omega^7 \\ 1 & \omega^2 & \omega^4 & \omega^6 & \omega^8 & \omega^{10} & \omega^{12} & \omega^{14} \\ 1 & \omega^3 & \omega^6 & \omega^9 & \omega^{12} & \omega^{15} & \omega^{18} & \omega^{21} \\ 1 & \omega^4 & \omega^8 & \omega^{12} & \omega^{16} & \omega^{20} & \omega^{24} & \omega^{28} \\ 1 & \omega^5 & \omega^{10} & \omega^{15} & \omega^{20} & \omega^{25} & \omega^{30} & \omega^{35} \\ 1 & \omega^6 & \omega^{12} & \omega^{18} & \omega^{24} & \omega^{30} & \omega^{36} & \omega^{42} \\ 1 & \omega^7 & \omega^{14} & \omega^{21} & \omega^{28} & \omega^{35} & \omega^{42} & \omega^{49} \end{pmatrix} = \begin{pmatrix} \boxed{1} & 1 & \boxed{1} & 1 & \boxed{1} & 1 & \boxed{1} & 1 \\ \boxed{1} & \omega & \boxed{i} & \overline{\omega} & \boxed{-1} & -\omega & \boxed{-i} & -\overline{\omega} \\ \boxed{1} & i & \boxed{-1} & -i & \boxed{1} & i & \boxed{-1} & -i \\ \boxed{1} & \overline{\omega} & \boxed{i} & \omega & \boxed{-1} & -\overline{\omega} & \boxed{-i} & -\omega \\ 1 & -1 & 1 & -1 & 1 & -1 & 1 & -1 \\ 1 & -\omega & i & -\overline{\omega} & -1 & \omega & -i & \overline{\omega} \\ 1 & -i & -1 & i & 1 & -i & -1 & i \\ 1 & -\overline{\omega} & -i & -\omega & -1 & \overline{\omega} & i & \omega \end{pmatrix}$$

<sup>4</sup>Actually, the FFT algorithm was previously published by Runge and König in 1924, and again by Yates in 1932, and again by Stumpf in 1937, and again by Danielson and Lanczos in 1942. But it was first *used* by Gauss in the 1800s for calculating the paths of asteroids from a finite number of equally-spaced observations. By hand. Fourier always did it the hard way. Cooley and Tukey apparently developed their algorithm to help detect Soviet nuclear tests without actually visiting Soviet nuclear facilities, by interpolating off-shore seismic readings. Without their rediscovery of the FFT algorithm, the nuclear test ban treaty would never have been ratified, and we'd all be speaking Russian, or more likely, whatever language radioactive glass speaks.

The boxed entries actually form the DFT matrix for  $n = 4$ !

$$\begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & \omega_4 & \omega_4^2 & \omega_4^3 \\ 1 & \omega_4^2 & \omega_4^4 & \omega_4^6 \\ 1 & \omega_4^3 & \omega_4^6 & \omega_4^9 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & i & -1 & -i \\ 1 & -1 & 1 & -1 \\ 1 & -i & -1 & i \end{bmatrix}$$

The input to the FFT algorithm is an array  $P[0..n-1]$  of coefficients of a polynomial  $p(x)$  with degree  $n-1$ . We start by splitting  $p$  into two smaller polynomials  $u$  and  $v$ , each with degree  $n/2-1$ , by setting

$$U[k] = P[2k] \quad \text{and} \quad V[k] = P[2k-1].$$

In other words,  $u$  has all the even-degree coefficients of  $p$ , and  $v$  has all the odd-degree coefficients. For example, if  $p(x) = 3x^3 - 4x^2 + 7x + 5$ , then  $u(x) = -4x + 5$  and  $v(x) = 3x + 7$ . These three polynomials satisfy the equation

$$p(x) = u(x^2) + x \cdot v(x^2).$$

In particular, if  $x$  is an  $n$ th root of unity, we have

$$P^*[k] = p(\omega_n^k) = u(\omega_n^{2k}) + \omega_n^k \cdot v(\omega_n^{2k}).$$

Now we can exploit those roots of unity again. Since  $n$  is a power of two,  $n$  must be even, so we have  $\omega_n^{2k} = \omega_{n/2}^k = \omega_{n/2}^{k \bmod n/2}$ . In other words, the values of  $p$  at the  $n$ th roots of unity depend on the values of  $u$  and  $v$  at  $(n/2)$ th roots of unity. Those are just coefficients in the DFTs of  $u$  and  $v$ !

$$P^*[k] = U^*[k \bmod n/2] + \omega_n^k \cdot V^*[k \bmod n/2]$$

So once we recursively compute  $U^*$  and  $V^*$ , we can compute  $P^*$  in linear time. The base case for the recurrence is  $n = 1$ . The overall running time satisfies the recurrence  $T(n) = \Theta(n) + 2T(n/2)$ , which as we all know solves to  $T(n) = \Theta(n \log n)$ .

Here's the complete FFT algorithm, along with its inverse.

<pre> FFT(<math>P[0..n-1]</math>):   if <math>n = 1</math>     return <math>P</math>   for <math>j \leftarrow 0</math> to <math>n/2 - 1</math>     <math>U[j] \leftarrow P[2j]</math>     <math>V[j] \leftarrow P[2j + 1]</math>   <math>U^* \leftarrow \text{FFT}(U[0..n/2 - 1])</math>   <math>V^* \leftarrow \text{FFT}(V[0..n/2 - 1])</math>   <math>\omega_n \leftarrow \cos(\frac{2\pi}{n}) + i \sin(\frac{2\pi}{n})</math>   <math>\omega \leftarrow 1</math>   for <math>j \leftarrow 0</math> to <math>n/2 - 1</math>     <math>P^*[j] \leftarrow U^*[j] + \omega \cdot V^*[j]</math>     <math>P^*[j + n/2] \leftarrow U^*[j] - \omega \cdot V^*[j]</math>     <math>\omega \leftarrow \omega \cdot \omega_n</math>   return <math>P^*[0..n-1]</math> </pre>
--

<pre> INVERSEFFT(<math>P^*[0..n-1]</math>):   if <math>n = 1</math>     return <math>P</math>   for <math>j \leftarrow 0</math> to <math>n/2 - 1</math>     <math>U^*[j] \leftarrow P^*[2j]</math>     <math>V^*[j] \leftarrow P^*[2j + 1]</math>   <math>U \leftarrow \text{INVERSEFFT}(U^*[0..n/2 - 1])</math>   <math>V \leftarrow \text{INVERSEFFT}(V^*[0..n/2 - 1])</math>   <math>\bar{\omega}_n \leftarrow \cos(\frac{2\pi}{n}) - i \sin(\frac{2\pi}{n})</math>   <math>\omega \leftarrow 1</math>   for <math>j \leftarrow 0</math> to <math>n/2 - 1</math>     <math>P[j] \leftarrow 2(U[j] + \omega \cdot V[j])</math>     <math>P[j + n/2] \leftarrow 2(U[j] - \omega \cdot V[j])</math>     <math>\omega \leftarrow \omega \cdot \bar{\omega}_n</math>   return <math>P[0..n-1]</math> </pre>
---

Given two polynomials  $p$  and  $q$ , each represented by an array of coefficients, we can multiply them in  $\Theta(n \log n)$  arithmetic operations as follows. First, pad the coefficient vectors and with zeros until the size is a power of two greater than or equal to the sum of the degrees. Then compute the DFTs of each coefficient vector, multiply the sample values one by one, and compute the inverse DFT of the resulting sample vector.

```

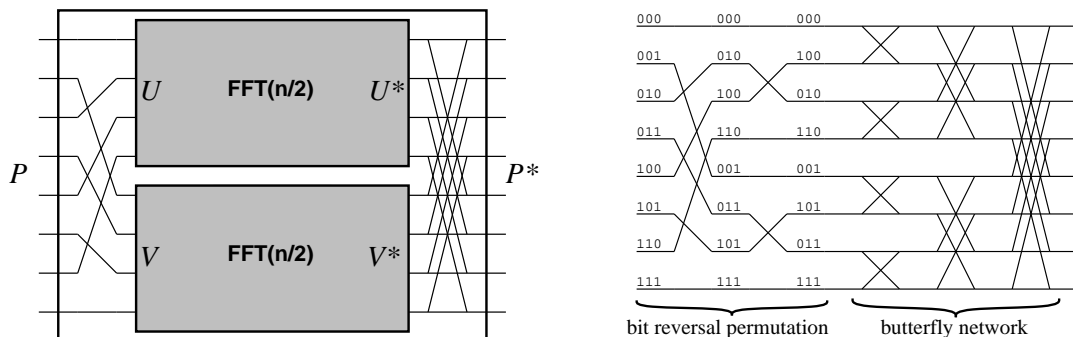
FFTMULTIPLY( $P[0..n-1], Q[0..m-1]$ ):
   $\ell \leftarrow \lceil \lg(n+m) \rceil$ 
  for  $j \leftarrow n$  to  $2^\ell - 1$ 
     $P[j] \leftarrow 0$ 
  for  $j \leftarrow m$  to  $2^\ell - 1$ 
     $Q[j] \leftarrow 0$ 

   $P^* \leftarrow FFT(P)$ 
   $Q^* \leftarrow FFT(Q)$ 
  for  $j \leftarrow 0$  to  $2^\ell - 1$ 
     $R^*[j] \leftarrow P^*[j] \cdot Q^*[j]$ 
  return INVERSEFFT( $R^*$ )

```

## 14.6 Inside the FFT

FFTs are often implemented in hardware as circuits. To see the recursive structure of the circuit, let's connect the top-level inputs and outputs to the inputs and outputs of the recursive calls. On the left we split the input  $P$  into two recursive inputs  $U$  and  $V$ . On the right, we combine the outputs  $U^*$  and  $V^*$  to obtain the final output  $P^*$ .



The recursive structure of the FFT algorithm.

If we expand this recursive structure completely, we see that the circuit splits naturally into two parts. The left half computes the *bit-reversal permutation* of the input. To find the position of  $P[k]$  in this permutation, write  $k$  in binary, and then read the bits backwards. For example, in an 8-element bit-reversal permutation,  $P[3] = P[011_2]$  ends up in position  $6 = 110_2$ . The right half of the FFT circuit is a *butterfly network*. Butterfly networks are often used to route between processors in massively-parallel computers, since they allow any processor to communicate with any other in only  $O(\log n)$  steps.

## 15 Number Theoretic Algorithms (November 12 and 14)

And it's one, two, three,  
 What are we fighting for?  
 Don't tell me, I don't give a damn,  
 Next stop is Vietnam; [or: This time we'll kill Saddam]  
 And it's five, six, seven,  
 Open up the pearly gates,  
 Well there ain't no time to wonder why,  
 Whoopee! We're all going to die.

— Country Joe and the Fish  
 “I-Feel-Like-I'm-Fixin'-to-Die Rag” (1967)

There are 0 kinds of mathematicians:  
 Those who can count modulo 2 and those who can't.

— anonymous

### 15.1 Greatest Common Divisors

Before we get to any actual algorithms, we need some definitions and preliminary results. **Unless specifically indicated otherwise, all variables in this lecture are integers.**

The symbol  $\mathbb{Z}$  (from the German word “Zahlen”, meaning ‘numbers’ or ‘to count’) to denote the set of integers. We say that one integer  $d$  *divides* another integer  $n$ , or that  $d$  is a *divisor* of  $n$ , if the quotient  $n/d$  is also an integer. Symbolically, we can write this definition as follows:

$$d \mid n \iff \left\lfloor \frac{n}{d} \right\rfloor = \frac{n}{d}$$

In particular, zero is not a divisor of any integer— $\infty$  is *not* an integer—but every other integer is a divisor of zero. If  $d$  and  $n$  are positive, then  $d \mid n$  immediately implies that  $d \leq n$ .

Any integer  $n$  can be written in the form  $n = qd + r$  for some non-negative integer  $0 \leq r < |d|$ . Moreover, the choices for the quotient  $q$  and remainder  $r$  are unique:

$$q = \left\lfloor \frac{n}{d} \right\rfloor \quad \text{and} \quad r = n \bmod d = n - d \left\lfloor \frac{n}{d} \right\rfloor.$$

Note that the remainder  $n \bmod d$  is *always* non-negative, even if  $n < 0$  or  $d < 0$  or both.<sup>1</sup>

If  $d$  divides two integers  $m$  and  $n$ , we say that  $d$  is a *common divisor* of  $m$  and  $n$ . It's trivial to prove (by definition crunching) that any common divisor of  $m$  and  $n$  also divides any integer linear combination of  $m$  and  $n$ :

$$(d \mid m) \text{ and } (d \mid n) \implies d \mid (am + bn)$$

The *greatest common divisor* of  $m$  and  $n$ , written  $\gcd(m, n)$ ,<sup>2</sup> is the largest integer that divides both  $m$  and  $n$ . Sometimes this is also called the greater common *denominator*. The greatest common divisor has another useful characterization as the *smallest* element of another set.

**Lemma 1.**  $\gcd(m, n)$  is the smallest positive integer of the form  $am + bn$ .

<sup>1</sup>The sign rules for the C/C++/Java % operator are just plain stupid. I can't count the number of times I've had to write  $x = (x+n)\%n$ ; instead of  $x \% = n$ ; Idiots.

<sup>2</sup>Do *not* use the notation  $(m, n)$  for greatest common divisor. *Ever*.

**Proof:** Let  $s$  be the smallest positive integer of the form  $am + bn$ . Any common divisor of  $m$  and  $n$  is also a divisor of  $s = am + bn$ . In particular,  $\gcd(m, n)$  is a divisor of  $s$ , which implies that  $\gcd(m, n) \leq s$ .

To prove the other inequality, let's show that  $s \mid m$  by calculating  $m \bmod s$ .

$$m \bmod s = m - s \left\lfloor \frac{m}{s} \right\rfloor = m - (am + bn) \left\lfloor \frac{m}{s} \right\rfloor = m \left( 1 - a \left\lfloor \frac{m}{s} \right\rfloor \right) + n \left( -b \left\lfloor \frac{m}{s} \right\rfloor \right)$$

We observe that  $m \bmod s$  is an integer linear combination of  $m$  and  $n$ . Since  $m \bmod s < s$ , and  $s$  is the smallest *positive* integer linear combination,  $m \bmod s$  cannot be positive. So it must be zero, which implies that  $s \mid m$ , as we claimed. By a symmetric argument,  $s \mid n$ . Thus,  $s$  is a common divisor of  $m$  and  $n$ . A common divisor can't be greater than the *greatest* common divisor, so  $s \leq \gcd(m, n)$ .

These two inequalities imply that  $s = \gcd(m, n)$ , completing the proof.  $\square$

## 15.2 Euclid's GCD Algorithm

The first part of this lecture is about computing the greatest common divisor of two integers. Our first algorithm for computing greatest common divisors follows immediately from two simple observations:

$$\gcd(m, n) = \gcd(m, n - m) \quad \text{and} \quad \gcd(n, 0) = n$$

The algorithm uses the first observation as a way to reduce the input and recurse; the second observation provides the base case.

```

SLOWGCD(m, n):
  m ← |m|; n ← |n|
  if m < n
    swap m ↔ n
  while n > 0
    m ← m - n
    if m < n
      swap m ↔ n
  return m
```

The first few lines just ensure that  $m \geq n \geq 0$ . Each iteration of the main loop decreases one of the numbers by at least 1, so the running time is  $O(m + n)$ . This bound is tight in the worst case; consider the case  $n = 1$ . Unfortunately, this is terrible. The input consists of just  $\log m + \log n$  bits; as a function of the input size, this algorithm runs in *exponential* time.

Let's think for a moment about what the main loop computes between swaps. We start with two numbers  $m$  and  $n$  and repeatedly subtract  $n$  from  $m$  until we can't any more. This is just a (slow) recipe for computing  $m \bmod n$ ! That means we can speed up the algorithm by using mod instead of subtraction.

```

EUCLIDGCD(m, n):
  m ← |m|; n ← |n|
  if m < n
    swap m ↔ n
  while n > 0
    m ← m mod n    (★)
    swap m ↔ n
  return m
```

This algorithm swaps  $m$  and  $n$  at *every* iteration, because  $m \bmod n$  is always less than  $n$ . This is usually called Euclid's algorithm, because the main idea is included in Euclid's *Elements*.<sup>3</sup>

The easiest way to analyze this algorithm is to work backward. First, let's consider the number of iterations of the main loop, or equivalently, the number times line ( $\star$ ) is executed. To keep things simple, let's assume that  $m > n > 0$ , so the first three lines are redundant, and the algorithm performs at least one iteration. Recall that the Fibonacci numbers(!) are defined as  $F_0 = 0$ ,  $F_1 = 1$ , and  $F_k = F_{k-1} + F_{k-2}$  for all  $k > 1$ .

**Lemma 2.** *If the algorithm performs  $k$  iterations, then  $m \geq F_{k+2}$  and  $n \geq F_{k+1}$ .*

**Proof (by induction on  $k$ ):** If  $k = 1$ , we have the trivial bounds  $n \geq 1 = F_2$  and  $m \geq 2 = F_3$ .

Suppose  $k > 1$ . The first iteration of the loop replaces  $(m, n)$  with  $(n, m \bmod n)$ . Since the algorithm performs  $k - 1$  more iterations, the inductive hypothesis implies that  $n \geq F_{k+1}$  and  $m \bmod n \geq F_k$ . We've assumed that  $m > n$ , so  $m \geq m + n(1 - \lfloor m/n \rfloor) = n + (m \bmod n)$ . We conclude that  $m \geq F_{k+1} + F_k = F_{k+2}$ .  $\square$

**Theorem 1.**  $\text{EUCLIDGCD}(m, n)$  runs in  $O(\log m)$  iterations.

**Proof:** Let  $k$  be the number of iterations. Lemma 2 implies that  $m \geq F_{k+2} \geq \phi^{k+2}/\sqrt{5} - 1$ , where  $\phi = (1 + \sqrt{5})/2$  (by the annihilator method). Thus,  $k \leq \log_\phi(\sqrt{5}(m + 1)) - 2 = O(\log m)$ .  $\square$

What about the actual running time? Every number used by the algorithm has  $O(\log m)$  bits. Computing the remainder of one  $b$ -bit integer by another using the grade-school long division algorithm requires  $O(b^2)$  time. So crudely, the running time is  $O(b^2 \log m) = O(\log^3 m)$ . More careful analysis reduces the time bound to  $O(\log^2 m)$ . We can make the algorithm even faster by using a fast integer division algorithm (based on FFTs, for example).

### 15.3 Modular Arithmetic and Algebraic Groups

Modular arithmetic is familiar to anyone who's ever wondered how many minutes are left in an exam that ends at 9:15 when the clock says 8:54.

When we do arithmetic 'modulo  $n$ ', what we're really doing is a funny kind of arithmetic on the elements of following set:

$$\mathbb{Z}_n = \{0, 1, 2, \dots, n-1\}$$

Modular addition and subtraction satisfies all the axioms that we expect implicitly:

- $\mathbb{Z}_n$  is *closed* under addition mod  $n$ : For any  $a, b \in \mathbb{Z}_n$ , their sum  $a + b \bmod n$  is also in  $\mathbb{Z}_n$

<sup>3</sup>However, Euclid's exposition was a little, erm, informal by current standards, primarily because the Greeks didn't know about induction. He basically said "Try one iteration. If that doesn't work, try three iterations." In modern language, Euclid's algorithm would be written as follows, assuming  $m \geq n > 0$ .

<pre> ACTUALEUCLIDGCD(m, n):   if n   m     return n   else     return n mod (m mod n) </pre>
---

This algorithm is *obviously* incorrect; consider the input  $m = 3$ ,  $n = 2$ . Nevertheless, mathematics and algorithms students have applied 'Euclidean induction' to a vast number of problems, only to scratch their heads in dismay when they don't get any credit.



- Addition is *associative*:  $(a + b \bmod n) + c \bmod n = a + (b + c \bmod n) \bmod n$ .
- Zero is an additive *identity* element:  $0 + a \bmod n = a + 0 \bmod n = a \bmod n$ .
- Every element  $a \in \mathbb{Z}_n$  has an *inverse*  $b \in \mathbb{Z}_n$  such that  $a + b \bmod n = 0$ . Specifically, if  $a = 0$ , then  $b = 0$ ; otherwise,  $b = n - a$ .

Any set with a binary operator that satisfies the closure, associativity, identity, and inverse axioms is called a *group*. Since  $\mathbb{Z}_n$  is a group under an ‘addition’ operator, we call it an *additive group*. Moreover, since addition is commutative ( $a + b \bmod n = b + a \bmod n$ ), we can call  $(\mathbb{Z}_n, + \bmod n)$  is an *abelian* additive group.

What about multiplication?  $\mathbb{Z}_n$  is closed under multiplication mod  $n$ , multiplication mod  $n$  is associative (and commutative), and 1 is a multiplicative identity, but some elements do not have multiplicative inverses. Formally, we say that  $\mathbb{Z}_n$  is a *ring* under addition and multiplication modulo  $n$ .

If  $n$  is composite, then the following theorem shows that we can factor the ring  $\mathbb{Z}_n$  into two smaller rings. The Chinese Remainder Theorem is named for Sun Tsu (or Sun Zi), the author of *the Art of War*, who proved a special case. (See the quotation for Lecture 1!)

**The Chinese Remainder Theorem.** *If  $p \perp q$ , then  $\mathbb{Z}_{pq} \cong \mathbb{Z}_p \times \mathbb{Z}_q$ .*

Okay, okay, before we prove this, let’s define all the notation. The product  $\mathbb{Z}_p \times \mathbb{Z}_q$  is the set of ordered pairs  $\{(a, b) \mid a \in \mathbb{Z}_p, b \in \mathbb{Z}_q\}$ , where addition, subtraction, and multiplication are defined as follows:

$$\begin{aligned}(a, b) + (c, d) &= (a + c \bmod p, b + d \bmod q) \\(a, b) - (c, d) &= (a - c \bmod p, b - d \bmod q) \\(a, b) \cdot (c, d) &= (ac \bmod p, bd \bmod q)\end{aligned}$$

It’s not hard to check that  $\mathbb{Z}_p \times \mathbb{Z}_q$  is a ring under these operations, where  $(0, 0)$  is the additive identity and  $(1, 1)$  is the multiplicative identity. The funky equal sign  $\cong$  means that these two rings are *isomorphic*: there is a bijection between the two sets that is consistent with the arithmetic operations.

As an example, the following table describes the bijection between  $\mathbb{Z}_{15}$  and  $\mathbb{Z}_3 \times \mathbb{Z}_5$ :

	0	1	2	3	4
0	0	6	12	3	9
1	10	1	7	13	4
2	5	11	2	8	14

For instance, we have  $8 = (2, 3)$  and  $13 = (1, 3)$ , and

$$\begin{aligned}(2, 3) + (1, 3) &= (2 + 1 \bmod 3, 3 + 3 \bmod 5) = (0, 1) = 6 = 21 \bmod 15 = (8 + 13) \bmod 15. \\(2, 3) \cdot (1, 3) &= (2 \cdot 1 \bmod 3, 3 \cdot 3 \bmod 5) = (2, 4) = 14 = 104 \bmod 15 = (8 \cdot 13) \bmod 15.\end{aligned}$$

**Proof:** The functions  $n \mapsto (n \bmod p, n \bmod q)$  and  $(a, b) \mapsto aq(q \bmod p) + bp(p \bmod q)$  are inverses of each other, and each map preserves the ring structure.  $\square$

We can extend the Chinese remainder theorem inductively as follows:

**The Real Chinese Remainder Theorem.** Suppose  $n = \prod_{i=1}^r p_i$ , where  $p_i \perp p_j$  for all  $i$  and  $j$ . Then  $\mathbb{Z}_n \cong \prod_{i=1}^r \mathbb{Z}_{p_i} = \mathbb{Z}_{p_1} \times \mathbb{Z}_{p_2} \times \cdots \times \mathbb{Z}_{p_r}$ .

If we want to perform modular arithmetic where the modulus  $n$  is very large, we can improve the performance of our algorithms by breaking  $n$  into several relatively prime factors, and performing modular arithmetic separately modulo each factor.

So we can do modular addition, subtraction, and multiplication; what about division? As I said earlier, not every element of  $\mathbb{Z}_n$  has a multiplicative inverse. The most obvious example is 0, but there can be others. For example, 3 has no multiplicative inverse in  $\mathbb{Z}_{15}$ ; there is no integer  $x$  such that  $3x \bmod 15 = 1$ . On the other hand, 0 is the only element of  $\mathbb{Z}_7$  without a multiplicative inverse:

$$1 \cdot 1 \equiv 2 \cdot 4 \equiv 3 \cdot 5 \equiv 6 \cdot 6 \equiv 1 \pmod{7}$$

These examples suggest (I hope) that  $x$  has a multiplicative inverse in  $\mathbb{Z}_n$  if and only if  $a$  and  $x$  are relatively prime. This is easy to prove as follows. If  $xy \bmod n = 1$ , then  $xy + kn = 1$  for some integer  $k$ . Thus, 1 is an integer linear combination of  $x$  and  $n$ , so Lemma 1 implies that  $\gcd(x, n) = 1$ . On the other hand, if  $x \perp n$ , then  $ax + bn = 1$  for some integers  $a$  and  $b$ , which implies that  $ax \bmod n = 1$ .

Let's define the set  $\mathbb{Z}_n^*$  to be the set of elements of  $\mathbb{Z}_n$  that have multiplicative inverses.

$$\mathbb{Z}_n^* = \{a \in \mathbb{Z}_n \mid a \perp n\}$$

It is a tedious exercise to show that  $\mathbb{Z}_n^*$  is an abelian group under multiplication modulo  $n$ . As long as we stick to elements of this group, we can reasonably talk about 'division mod  $n$ '.

We denote the number of elements in  $\mathbb{Z}_n^*$  by  $\phi(n)$ ; this is called Euler's *totient* function. This function is remarkably badly-behaved, but there is a relatively simple formula for  $\phi(n)$  (not surprisingly) involving prime numbers and division:

$$\phi(n) = n \prod_{p|n} \frac{p-1}{p}$$

I won't prove this formula, but the following intuition is helpful. If we start with  $\mathbb{Z}_n$  and throw out all  $n/2$  multiples of 2, all  $n/3$  multiples of 3, all  $n/5$  multiples of 5, and so on. Whenever we throw out multiples of  $p$ , we multiply the size of the set by  $(p-1)/p$ . At the end of this process, we're left with precisely the elements of  $\mathbb{Z}_n^*$ . *This is not a proof!* On the one hand, this argument throws out some numbers (like 6) more than once, so our estimate seems too low. On the other hand, there are actually  $\lfloor n/p \rfloor$  multiples of  $p$  in  $\mathbb{Z}_n$ , so our estimate seems too high. Surprisingly, these two errors exactly cancel each other out.

## 15.4 Toward Primality Testing

In this last section, we discuss algorithms for detecting whether a number is prime. Large prime numbers are used primarily (but not exclusively) in cryptography algorithms.

A positive integer is *prime* if it has exactly two positive divisors, and *composite* if it has more than two positive divisors. The integer 1 is neither prime nor composite. Equivalently, an integer  $n \geq 2$  is prime if  $n$  is relatively prime with every positive integer smaller than  $n$ . We can rephrase this definition yet again:  $n$  is prime if and only if  $\phi(n) = n - 1$ .

The obvious algorithm for testing whether a number is prime is *trial division*: simply try every possible nontrivial divisor between 2 and  $\sqrt{n}$ .

```

TRIALDIVPRIME( $n$ ):
  for  $d \leftarrow 1$  to  $\lfloor \sqrt{n} \rfloor$ 
    if  $n \bmod d = 0$ 
      return COMPOSITE
  return PRIME

```

Unfortunately, this algorithm is horribly slow. Even if we could do the remainder computation in constant time, the overall running time of this algorithm would be  $\Omega(\sqrt{n})$ , which is exponential in the number of input bits.

This might seem completely hopeless, but fortunately most composite numbers are quite easy to detect as composite. Consider, for example, the related problem of deciding whether a given integer  $n$ , whether  $n = m^e$  for any integers  $m > 1$  and  $e > 1$ . We can solve this problem in polynomial time with the following straightforward algorithm. The subroutine  $\text{ROOT}(n, i)$  computes  $\lfloor n^{1/i} \rfloor$  essentially by binary search. (I'll leave the analysis as a simple exercise.)

```

EXACTPOWER?( $n$ ):
  for  $i \leftarrow 2$  to  $\lg n$ 
    if  $(\text{ROOT}(n, i))^i = n$ 
      return TRUE
  return FALSE

```

```

ROOT( $n, i$ ):
   $r \leftarrow 0$ 
  for  $\ell \leftarrow \lceil (\lg n)/i \rceil$  down to 1
    if  $(r + 2^\ell)^i \leq n$ 
       $r \leftarrow r + 2^\ell$ 
  return  $r$ 

```

To distinguish between arbitrary prime and composite numbers, we need to exploit some results about  $\mathbb{Z}_n^*$  from group theory and number theory. First, we define the *order* of an element  $x \in \mathbb{Z}_n^*$  as the smallest positive integer  $k$  such that  $x^k \equiv 1 \pmod{n}$ . For example, in the group

$$\mathbb{Z}_{15}^* = \{1, 2, 4, 7, 8, 11, 13, 14\},$$

the number 2 has order 4, and the number 11 has order 2. For any  $x \in \mathbb{Z}_n^*$ , we can partition the elements of  $\mathbb{Z}_n^*$  into equivalence classes, by declaring  $a \sim_x b$  if  $a \equiv b \cdot x^k$  for some integer  $k$ . The size of every equivalence class is exactly the order of  $x$ . Since the equivalence classes must be disjoint, we can conclude that the order of any element divides the size of the group. We can express this observation more succinctly as follows:

**Euler's Theorem.**  $a^{\phi(n)} \equiv 1 \pmod{n}$ .<sup>4</sup>

The most interesting special case of this theorem is when  $n$  is prime.

**Fermat's Little Theorem.** If  $p$  is prime, then  $a^p \equiv a \pmod{p}$ .<sup>5</sup>

This theorem leads to the following efficient *pseudo*-primality test.

<sup>4</sup>This is not Euler's only theorem; he had thousands. It's not even his most famous theorem. His *second* most famous theorem is the formula  $v + e - f = 2$  relating the vertices, edges and faces of any planar map. His most famous theorem is the magic formula  $e^{\pi i} + 1 = 0$ . Naming something after a mathematician or physicist (as in 'Euler tour' or 'Gaussian elimination' or 'Avogadro's number') is considered a high compliment. Using a lower case letter ('abelian group') is even better; abbreviating ('volt', 'amp') is better still. The number  $e$  was named after Euler.

<sup>5</sup>This is not Fermat's only theorem; he had hundreds, most of them stated without proof. Fermat's Last Theorem wasn't the last one he published, but the last one proved. Amazingly, despite his dislike of writing proofs, Fermat was almost always right. In that respect, he was *very* different from you and me.

```

FERMATPSEUDOPRIME( $n$ ) :
  choose an integer  $a$  between 1 and  $n - 1$ 
  if  $a^n \bmod n \neq a$ 
    return COMPOSITE!
  else
    return PRIME?

```

In practice, this algorithm is both fast and effective. The (empirical) probability that a random 100-digit composite number will return PRIME? is roughly  $10^{-30}$ , even if we always choose  $a = 2$ . Unfortunately, there are composite numbers that always pass this test, no matter which value of  $a$  we use. A *Carmichael number* is a composite integer  $n$  such that  $a^n \equiv a \pmod{n}$  for every integer  $a$ . Thus, Fermat's Little Theorem can be used to distinguish between two types of numbers: (primes and Carmichael numbers) and everything else. Carmichael numbers are extremely rare; in fact, it was proved only a decade ago that there are an infinite number of them.

To deal with Carmichael numbers effectively, we need to look more closely at the structure of the group  $\mathbb{Z}_n^*$ . We say that  $\mathbb{Z}_n^*$  is *cyclic* if it contains an element of order  $\phi(n)$ ; such an element is called a *generator*. Successive powers of any generator *cycle* through every element of the group in some order. For example, the group  $\mathbb{Z}_9^* = \{1, 2, 4, 5, 7, 8\}$  is cyclic, with two generators: 2 and 5, but  $\mathbb{Z}_{15}^*$  is not cyclic. The following theorem completely characterizes which groups  $\mathbb{Z}_n^*$  are cyclic.

**The Cycle Theorem.**  $\mathbb{Z}_n^*$  is cyclic if and only if  $n = 2, 4, p^e$ , or  $2p^e$  for some odd prime  $p$  and positive integer  $e$ .

This theorem has two relatively simple corollaries.

**The Discrete Log Theorem.** Suppose  $\mathbb{Z}_n^*$  is cyclic and  $g$  is a generator. Then  $g^x \equiv g^y \pmod{n}$  if and only if  $x \equiv y \pmod{\phi(n)}$ .

**Proof:** Suppose  $g^x \equiv g^y \pmod{n}$ . By definition of 'generator', the sequence  $\langle 1, g, g^2, \dots \rangle$  has period  $\phi(n)$ . Thus,  $x \equiv y \pmod{\phi(n)}$ . On the other hand, if  $x \equiv y \pmod{\phi(n)}$ , then  $x = y + k\phi(n)$  for some integer  $k$ , so  $g^x = g^{y+k\phi(n)} = g^y \cdot (g^{\phi(n)})^k$ . Euler's Theorem now implies that  $(g^{\phi(n)})^k \equiv 1^k \equiv 1 \pmod{n}$ , so  $g^x \equiv g^y \pmod{n}$ .  $\square$

**The  $\sqrt{1}$  Theorem.** Suppose  $n = p^e$  for some odd prime  $p$  and positive integer  $e$ . The only elements  $x \in \mathbb{Z}_n^*$  that satisfy the equation  $x^2 \equiv 1 \pmod{n}$  are  $x = 1$  and  $x = n - 1$ .

**Proof:** Obviously  $1^2 \equiv 1 \pmod{n}$  and  $(n - 1)^2 = n^2 - 2n + 1 \equiv 1 \pmod{n}$ .

Suppose  $x^2 \equiv 1 \pmod{n}$  where  $n = p^e$ . By the Cycle Theorem,  $\mathbb{Z}_n^*$  is cyclic. Let  $g$  be a generator of  $\mathbb{Z}_n^*$ , and suppose  $x = g^k$ . Then we immediately have  $x^2 = g^{2k} \equiv g^0 = 1 \pmod{p^e}$ . The Discrete Log Theorem implies that  $2k \equiv 0 \pmod{\phi(p^e)}$ . Since  $p$  is an odd prime, we have  $\phi(p^e) = (p - 1)p^{e-1}$ , which is even. Thus, the equation  $2k \equiv 0 \pmod{\phi(p^e)}$  has just two solutions:  $k = 0$  and  $k = \phi(p^e)/2$ . By the Cycle Theorem, either  $x = 1$  or  $x = g^{\phi(n)/2}$ . Because  $x = n - 1$  is also a solution to the original equation, we must have  $g^{\phi(n)/2} \equiv n - 1 \pmod{n}$ .  $\square$

This theorem leads to a different *pseudo*-primality algorithm:

```

SQRT1PSEUDOPRIME( $n$ ):
  choose a number  $a$  between 2 and  $n - 2$ 
  if  $a^2 \bmod n = 1$ 
    return COMPOSITE!
  else
    return PRIME?

```

As with the previous pseudo-primality test, there are composite numbers that this algorithm cannot identify as composite: powers of primes, for instance. Fortunately, however, the set of composites that always pass the  $\sqrt{\cdot}$  test is disjoint from the set of numbers that always pass the Fermat test. In particular, Carmichael numbers *never* have the form  $p^e$ .

### 15.5 The Miller-Rabin Primality Test

The following randomized algorithm, adapted by Michael Rabin from an earlier deterministic algorithm of Gary Miller\*, combines the Fermat test and the  $\sqrt{\cdot}$  test. The algorithm repeats the same two tests  $s$  times, where  $s$  is some user-chosen parameter, each time with a random value of  $a$ .

MILLERRABIN( $n$ ):	
write $n - 1 = 2^t u$ where $u$ is odd	
for $i \leftarrow 1$ to $s$	
$a \leftarrow \text{RANDOM}(2, n - 2)$	
if $\text{EUCLIDGCD}(a, n) \neq 1$	
return COMPOSITE!	⟨⟨obviously!⟩⟩
$x_0 \leftarrow a^u \bmod n$	
for $j \leftarrow 1$ to $t$	
$x_j \leftarrow x_{j-1}^2 \bmod n$	
if $x_j = 1$ and $x_{j-1} \neq 1$ and $x_{j-1} \neq n - 1$	
return COMPOSITE!	⟨⟨by the $\sqrt{\cdot}$ Theorem⟩⟩
if $x_t \neq 1$	⟨⟨ $x_t = a^{n-1} \bmod n$ ⟩⟩
return COMPOSITE!	⟨⟨by Fermat's Little Theorem⟩⟩
return PRIME?	

First let's consider the running time; for simplicity, we assume that all integer arithmetic is done using the quadratic-time grade school algorithms. We can compute  $u$  and  $t$  in  $O(\log n)$  time by scanning the bits in the binary representation of  $n$ . Euclid's algorithm takes  $O(\log^2 n)$  time. Computing  $a^u \bmod n$  requires  $O(\log u) = O(\log n)$  multiplications, each of which takes  $O(\log^2 n)$  time. Squaring  $x_j$  takes  $O(\log^2 n)$  time. Overall, the running time for one iteration of the outer loop is  $O(\log^3 n + t \log^2 n) = O(\log^3 n)$ , since  $t \leq \lg n$ . Thus, the total running time of this algorithm is  $O(s \log^3 n)$ . If we set  $s = O(\log n)$ , this running time is polynomial in the size of the input.

Fine, so it's fast, but is it correct? Like the earlier pseudoprime testing algorithms, a prime input will always cause MILLERRABIN to return PRIME?. Composite numbers, however, may not always return COMPOSITE!; because we choose the number  $a$  at random, there is a small probability of error.<sup>6</sup> Fortunately, the error probability can be made ridiculously small—in practice, less than the probability that random quantum fluctuations will instantly transform your computer into a kitten—by setting  $s \approx 1000$ .

**Theorem 2.** *If  $n$  is composite, MILLERRABIN( $n$ ) returns COMPOSITE! with probability at least  $1 - 2^{-s}$ .*

<sup>6</sup>If instead, we try all possible values of  $a$ , we obtain an exact primality testing algorithm, but it runs in exponential time. Miller's original deterministic algorithm examined every value of  $a$  in a carefully-chosen subset of  $\mathbf{Z}_n^*$ . If the Extended Riemann Hypothesis holds, this subset has logarithmic size, and Miller's algorithm runs in polynomial time. The Riemann Hypothesis is a century-old open problem about the distribution of prime numbers. A solution would be at least as significant as proving Fermat's Last Theorem or  $P \neq NP$ .

**Proof:** First, suppose  $n$  is not a Carmichael number. Let  $F$  be the set of elements of  $\mathbb{Z}_n^*$  that pass the Fermat test:

$$F = \{a \in \mathbb{Z}_n^* \mid a^{n-1} \equiv 1 \pmod{n}\}.$$

Since  $n$  is not a Carmichael number,  $F$  is a *proper* subset of  $\mathbb{Z}_n^*$ . Given any two elements  $a, b \in F$ , their product  $a \cdot b \pmod{n}$  in  $\mathbb{Z}_n^*$  is also an element of  $F$ :

$$(a \cdot b)^{n-1} \equiv a^{n-1}b^{n-1} \equiv 1 \cdot 1 \equiv 1 \pmod{n}$$

We also easily observe that 1 is an element of  $F$ , and the multiplicative inverse  $\pmod{n}$  of any element of  $F$  is also in  $F$ . Thus,  $F$  is a proper *subgroup* of  $\mathbb{Z}_n^*$ , that is, a proper subset that is also a group under the same binary operation. A standard result in group theory states that if  $F$  is a subgroup of a finite group  $G$ , the number of elements of  $F$  divides the number of elements of  $G$ . (We used a special case of this result in our proof of Euler's Theorem.) In our setting, this means that  $|F|$  divides  $\phi(n)$ . Since we already know that  $|F| < \phi(n)$ , we must have  $|F| \leq \phi(n)/2$ . Thus, at most half the elements of  $\mathbb{Z}_n^*$  pass the Fermat test.

The case of Carmichael numbers is more complicated, but the main idea is the same: at most half the possible values of  $a$  pass the  $\sqrt{1}$  test. See CLRS for further details.  $\square$

## C String Matching

### C.1 Brute Force

The basic object that we're going to talk about for the next two lectures is a *string*, which is really just an array. The elements of the array come from a set  $\Sigma$  called the *alphabet*; the elements themselves are called *characters*. Common examples are ASCII text, where each character is an seven-bit integer<sup>1</sup>, strands of DNA, where the alphabet is the set of nucleotides  $\{A, C, G, T\}$ , or proteins, where the alphabet is the set of 22 amino acids.

The problem we want to solve is the following. Given two strings, a *text*  $T[1..n]$  and a *pattern*  $P[1..m]$ , find the first *substring* of the text that is the same as the pattern. (It would be easy to extend our algorithms to find *all* matching substrings, but we will resist.) A substring is just a contiguous subarray. For any *shift*  $s$ , let  $T_s$  denote the substring  $T[s..s+m-1]$ . So more formally, we want to find the smallest shift  $s$  such that  $T_s = P$ , or report that there is no match. For example, if the text is the string 'AMANAPLANACATACANALPANAMA'<sup>2</sup> and the pattern is 'CAN', then the output should be 15. If the pattern is 'SPAM', then the answer should be 'none'. In most cases the pattern is much smaller than the text; to make this concrete, I'll assume that  $m < n/2$ .

Here's the 'obvious' brute force algorithm, but with one immediate improvement. The inner while loop compares the substring  $T_s$  with  $P$ . If the two strings are not equal, this loop stops at the first character mismatch.

```

ALMOSTBRUTEFORCE( $T[1..n], P[1..m]$ ):
  for  $s \leftarrow 1$  to  $n - m + 1$ 
    equal  $\leftarrow$  true
     $i \leftarrow 1$ 
    while equal and  $i \leq m$ 
      if  $T[s + i - 1] \neq P[i]$ 
        equal  $\leftarrow$  false
      else
         $i \leftarrow i + 1$ 
    if equal
      return  $s$ 
  return 'none'

```

In the worst case, the running time of this algorithm is  $O((n-m)m) = O(nm)$ , and we can

<sup>1</sup>Yes, *seven*. Most computer systems use some sort of 8-bit character set, but there's no universally accepted standard. Java supposedly uses the Unicode character set, which has variable-length characters and therefore doesn't really fit into our framework. Just think, someday you'll be able to write ' $\P = \aleph[\infty+]/\mathbb{U}$ ;' in your Java code! Joy!

<sup>2</sup>Dan Hoey (or rather, his computer program) found the following 540-word palindrome in 1984. We have better online dictionaries now, so I'm sure you could do better.

A man, a plan, a caret, a ban, a myriad, a sum, a lac, a liar, a hoop, a pint, a catalpa, a gas, an oil, a bird, a yell, a vat, a caw, a pax, a wag, a tax, a nay, a ram, a cap, a yam, a gay, a tsar, a wall, a car, a luger, a ward, a bin, a woman, a vassal, a wolf, a tuna, a nit, a pall, a fret, a watt, a bay, a daub, a tan, a cab, a datum, a gall, a hat, a fag, a zap, a say, a jaw, a lay, a wet, a gallop, a tug, a trot, a trap, a tram, a torr, a caper, a top, a tonk, a toll, a ball, a fair, a sax, a minim, a tenor, a bass, a passer, a capital, a rut, an amen, a ted, a cabal, a tang, a sun, an ass, a maw, a sag, a jam, a dam, a sub, a salt, an axon, a sail, an ad, a wadi, a radian, a room, a rood, a rip, a tad, a pariah, a revel, a reel, a reed, a pool, a plug, a pin, a peek, a parabola, a dog, a pat, a cud, a nu, a fan, a pal, a rum, a nod, an eta, a lag, an eel, a batik, a mug, a mot, a nap, a maxim, a mood, a leek, a grub, a gob, a gel, a drab, a citadel, a total, a cedar, a tap, a gag, a rat, a manor, a bar, a gal, a cola, a pap, a yaw, a tab, a raj, a gab, a nag, a pagan, a bag, a jar, a bat, a way, a papa, a local, a gar, a baron, a mat, a rag, a gap, a tar, a decal, a tot, a led, a tic, a bard, a leg, a bog, a burg, a keel, a doom, a mix, a map, an atom, a gum, a kit, a baleen, a gala, a ten, a don, a mural, a pan, a faun, a ducat, a pagoda, a lob, a rap, a keep, a nip, a gulp, a loop, a deer, a leer, a lever, a hair, a pad, a tapir, a door, a moor, an aid, a raid, a wad, an alias, an ox, an atlas, a bus, a madam, a jag, a saw, a mass, an anus, a gnat, a lab, a cadet, an em, a natural, a tip, a caress, a pass, a baronet, a minimax, a sari, a fall, a ballot, a knot, a pot, a rep, a carrot, a mart, a part, a tort, a gut, a poll, a gateway, a law, a jay, a sap, a zag, a fat, a hall, a gamut, a dab, a can, a tabu, a day, a batt, a waterfall, a patina, a nut, a flow, a lass, a van, a nib, a draw, a regular, a call, a war, a stay, a gam, a yap, a cam, a ray, an ax, a tag, a wax, a paw, a cat, a valley, a drib, a lion, a saga, a plat, a catnip, a pooh, a rail, a calamus, a dairyman, a bater, a canal—Panama!

actually achieve this running time by searching for the pattern  $AAA \dots AAAB$  with  $m - 1$  A's, in a text consisting of  $n$  A's.

In practice, though, breaking out of the inner loop at the first mismatch makes this algorithm quite practical. We can wave our hands at this by assuming that the text and pattern are both random. Then on average, we perform a constant number of comparisons at each position  $i$ , so the total expected number of comparisons is  $O(n)$ . Of course, neither English nor DNA is really random, so this is only a heuristic argument.

## C.2 Strings as Numbers

For the rest of the lecture, let's assume that the alphabet consists of the numbers 0 through 9, so we can interpret any array of characters as either a string or a decimal number. In particular, let  $p$  be the numerical value of the pattern  $P$ , and for any shift  $s$ , let  $t_s$  be the numerical value of  $T_s$ :

$$p = \sum_{i=1}^m 10^{m-i} \cdot P[i] \quad t_s = \sum_{i=1}^m 10^{m-i} \cdot T[s+i-1]$$

For example, if  $T = 31415926535897932384626433832795028841971$  and  $m = 4$ , then  $t_{17} = 2384$ .

Clearly we can rephrase our problem as follows: Find the smallest  $s$ , if any, such that  $p = t_s$ . We can compute  $p$  in  $O(m)$  arithmetic operations, without having to explicitly compute powers of ten, using *Horner's rule*:

$$p = P[m] + 10(P[m-1] + 10(P[m-2] + \dots + 10(P[2] + 10 \cdot P[1]) \dots))$$

We could also compute any  $t_s$  in  $O(m)$  operations using Horner's rule, but this leads to essentially the same brute-force algorithm as before. But once we know  $t_s$ , we can actually compute  $t_{s+1}$  in constant time just by doing a little arithmetic — subtract off the most significant digit  $T[s] \cdot 10^{m-1}$ , shift everything up by one digit, and add the new least significant digit  $T[s+m]$ :

$$t_{s+1} = 10(t_s - 10^{m-1} \cdot T[s]) + T[s+m]$$

To make this fast, we need to precompute the constant  $10^{m-1}$ . (And we know how to do that quickly. Right?) So it seems that we can solve the string matching problem in  $O(n)$  worst-case time using the following algorithm:

```

NUMBERSEARCH( $T[1..n], P[1..m]$ ):
   $\sigma \leftarrow 10^{m-1}$ 
   $p \leftarrow 0$ 
   $t_1 \leftarrow 0$ 
  for  $i \leftarrow 1$  to  $m$ 
     $p \leftarrow 10 \cdot p + P[i]$ 
     $t_1 \leftarrow 10 \cdot t_1 + T[i]$ 
  for  $s \leftarrow 1$  to  $n - m + 1$ 
    if  $p = t_s$ 
      return  $s$ 
     $t_{s+1} \leftarrow 10 \cdot (t_s - \sigma \cdot T[s]) + T[s+m]$ 
  return 'none'

```

Unfortunately, the most we can say is that the number of *arithmetic operations* is  $O(n)$ . These operations act on numbers with up to  $m$  digits. Since we want to handle arbitrarily long patterns, we can't assume that each operation takes only constant time!



### C.3 Karp-Rabin Fingerprinting

To make this algorithm efficient, we will make one simple change, discovered by Richard Karp and Michael Rabin in 1981:

Perform all arithmetic modulo some prime number  $q$ .

We choose  $q$  so that the value  $10q$  fits into a standard integer variable, so that we don't need any fancy long-integer data types. The values  $(p \bmod q)$  and  $(t_s \bmod q)$  are called the *fingerprints* of  $P$  and  $T_s$ , respectively. We can now compute  $(p \bmod q)$  and  $(t_1 \bmod q)$  in  $O(m)$  time using Horner's rule 'mod  $q$ '

$$p \bmod q = P[m] + (\dots + (10 \cdot (P[2] + (10 \cdot P[1] \bmod q) \bmod q) \bmod q) \dots) \bmod q$$

and similarly, given  $(t_s \bmod q)$ , we can compute  $(t_{s+1} \bmod q)$  in constant time.

$$t_{s+1} \bmod q = (10 \cdot (t_s - ((10^{m-1} \bmod q) \cdot T[s] \bmod q) \bmod q) \bmod q) + T[s + m] \bmod q$$

Again, we have to precompute the value  $(10^{m-1} \bmod q)$  to make this fast.

If  $(p \bmod q) \neq (t_s \bmod q)$ , then certainly  $P \neq T_s$ . However, if  $(p \bmod q) = (t_s \bmod q)$ , we can't tell whether  $P = T_s$  or not. All we know for sure is that  $p$  and  $t_s$  differ by some integer multiple of  $q$ . If  $P \neq T_s$  in this case, we say there is a *false match* at shift  $s$ . To test for a false match, we simply do a brute-force string comparison. (In the algorithm below,  $\tilde{p} = p \bmod q$  and  $\tilde{t}_s = t_s \bmod q$ .)

KARPRABIN( $T[1..n], P[1..m]$ ):

```

choose a small prime  $q$ 
 $\sigma \leftarrow 10^{m-1} \bmod q$ 
 $\tilde{p} \leftarrow 0$ 
 $\tilde{t}_1 \leftarrow 0$ 
for  $i \leftarrow 1$  to  $m$ 
     $\tilde{p} \leftarrow (10 \cdot \tilde{p} \bmod q) + P[i] \bmod q$ 
     $\tilde{t}_1 \leftarrow (10 \cdot \tilde{t}_1 \bmod q) + T[i] \bmod q$ 
for  $s \leftarrow 1$  to  $n - m + 1$ 
    if  $\tilde{p} = \tilde{t}_s$ 
        if  $P = T_s$        $\langle\langle$ brute-force  $O(m)$ -time comparison $\rangle\rangle$ 
            return  $s$ 
     $\tilde{t}_{s+1} \leftarrow (10 \cdot (\tilde{t}_s - (\sigma \cdot T[s] \bmod q) \bmod q) \bmod q) + T[s + m] \bmod q$ 
return 'none'

```

The running time of this algorithm is  $O(n + Fm)$ , where  $F$  is the number of false matches.

Intuitively, we expect the fingerprints  $t_s$  to jump around between 0 and  $q - 1$  more or less at random, so the 'probability' of a false match 'ought' to be  $1/q$ . This intuition implies that  $F = n/q$  'on average', which gives us an 'expected' running time of  $O(n + nm/q)$ . If we always choose  $q \geq m$ , this simplifies to  $O(n)$ . But of course all this intuitive talk of probabilities is just frantic meaningless handwaving, since we haven't actually done anything random yet.

### C.4 Random Prime Number Facts

The real power of the Karp-Rabin algorithm is that by choosing the modulus  $q$  *randomly*, we can actually formalize this intuition! The first line of KARPRABIN should really read as follows:

Let  $q$  be a random prime number less than  $nm^2 \log(nm^2)$ .

For any positive integer  $u$ , let  $\pi(u)$  denote the number of prime numbers less than  $u$ . There are  $\pi(nm^2 \log nm^2)$  possible values for  $q$ , each with the same probability of being chosen.

Our analysis needs two results from number theory. I won't even try to prove the first one, but the second one is quite easy.

**Lemma 1 (The Prime Number Theorem).**  $\pi(u) = \Theta(u/\log u)$ .

**Lemma 2.** Any integer  $x$  has at most  $\lceil \lg x \rceil$  distinct prime divisors.

**Proof:** If  $x$  has  $k$  distinct prime divisors, then  $x \geq 2^k$ , since every prime number is bigger than 1. □

Let's assume that there are no true matches, so  $p \neq t_s$  for all  $s$ . (That's the worst case for the algorithm anyway.) Let's define a strange variable  $X$  as follows:

$$X = \prod_{s=1}^{n-m+1} |p - t_s|.$$

Notice that by our assumption,  $X$  can't be zero.

Now suppose we have false match at shift  $s$ . Then  $p \bmod q = t_s \bmod q$ , so  $p - t_s$  is an integer multiple of  $q$ , and this implies that  $X$  is also an integer multiple of  $q$ . In other words, if there is a false match, then  $q$  must one of the prime divisors of  $X$ .

Since  $p < 10^m$  and  $t_s < 10^m$ , we must have  $X < 10^{2m}$ . Thus, by the second lemma,  $X$  has  $O(2m)$  prime divisors. Since we chose  $q$  randomly from a set of  $\pi(nm^2 \log(nm^2)) = \Omega(nm^2)$  prime numbers, the probability that  $q$  divides  $X$  is at most

$$\frac{O(2m)}{\Omega(nm^2)} = O\left(\frac{1}{m}\right).$$

We have just proven the following amazing fact.

The probability of getting a false match is  $O(1/m)$ .

Recall that the running time of KARPRABIN is  $O(n + mF)$ , where  $F$  is the number of false matches. By using the *really* loose upper bound  $E[F] \leq \Pr[F > 0] \cdot n$ , we can conclude that the expected number of false matches is  $O(n/m)$ . Thus, the expected running time of the KARPRABIN algorithm is  $O(n)$ .

## C.5 Random Prime Number?

Actually choosing a random prime number is not particularly easy. The best method known is to repeatedly generate a random integer and test to see if it's prime. In practice, it's enough to choose a random *probable* prime. You can read about probable primes in the textbook *Randomized Algorithms* by Rajeev Motwani and Prabhakar Raghavan (Cambridge, 1995).

## D More String Matching

### D.1 Redundant Comparisons

Let's go back to the character-by-character method for string matching. Suppose we are looking for the pattern 'ABRACADABRA' in some longer text using the (almost) brute force algorithm described in the previous lecture. Suppose also that when  $s = 11$ , the substring comparison fails at the fifth position; the corresponding character in the text (just after the vertical line below) is not a C. At this point, our algorithm would increment  $s$  and start the substring comparison from scratch.

```

HOCUSPOCUSABRACADABRA...
                |
                C
ABRACADABRA
ABRACADABRA

```

If we look carefully at the text and the pattern, however, we should notice right away that there's no point in looking at  $s = 12$ . We already know that the next character is a B — after all, it matched  $P[2]$  during the previous comparison — so why bother even looking there? Likewise, we already know that the next two shifts  $s = 13$  and  $s = 14$  will also fail, so why bother looking there?

```

HOCUSPOCUSABRACADABRA...
                |
                C
ABRACADABRA
ABRACADABRA
ABRACADABRA
ABRACADABRA
ABRACADABRA

```

Finally, when we get to  $s = 15$ , we can't immediately rule out a match based on earlier comparisons. However, for precisely the same reason, we shouldn't start the substring comparison over from scratch — we already know that  $T[15] = P[4] = A$ . Instead, we should start the substring comparison at the *second* character of the pattern, since we don't yet know whether or not it matches the corresponding text character.

If you play with this idea long enough, you'll notice that the character comparisons should always advance through the text. **Once we've found a match for a text character, we never need to do another comparison with that character again.** In other words, we should be able to optimize the brute-force algorithm so that it always *advances* through the text.

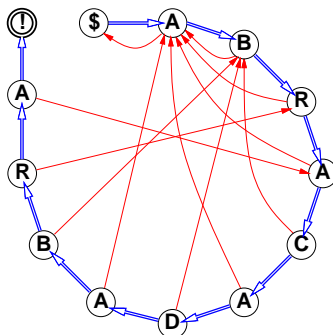
You'll also eventually notice a good rule for finding the next 'reasonable' shift  $s$ . A *prefix* of a string is a substring that includes the first character; a *suffix* is a substring that includes the last character. A prefix or suffix is *proper* if it is not the entire string. Suppose we have just discovered that  $T[i] \neq P[j]$ . **The next reasonable shift is the smallest value of  $s$  such that  $T[s..i-1]$ , which is a suffix of the previously-read text, is also a proper prefix of the pattern.**

In this lecture, we'll describe a string matching algorithm, published by Donald Knuth, James Morris, and Vaughn Pratt in 1977, that implements both of these ideas.

### D.2 Finite State Machines

If we have a string matching algorithm that follows our first observation (that we always advance through the text), we can interpret it as feeding the text through a special type of *finite-state machine*. A finite state machine is a directed graph. Each node in the graph, or *state*, is labeled with a character from the pattern, except for two special nodes labeled  $\$$  and  $\textcircled{1}$ . Each node has two outgoing edges, a *success* edge and a *failure* edge. The success edges define a path through the

characters of the pattern in order, starting at  $\textcircled{\$}$  and ending at  $\textcircled{!}$ . Failure edges always point to earlier characters in the pattern.



A finite state machine for the string 'ABRADACABRA'.  
Thick arrows are the success edges; thin arrows are the failure edges.

We use the finite state machine to search for the pattern as follows. At all times, we have a current text character  $T[i]$  and a current node in the graph, which is usually labeled by some pattern character  $P[j]$ . We iterate the following rules:

- If  $T[i] = P[j]$ , or if the current label is  $\textcircled{\$}$ , follow the success edge to the next node and increment  $i$ . (So there is no failure edge from the start node  $\textcircled{\$}$ .)
- If  $T[i] \neq P[j]$ , follow the failure edge back to an earlier node, but do not change  $i$ .

For the moment, let's simply assume that the failure edges are defined correctly—we'll come back to this later. If we ever reach the node labeled  $\textcircled{!}$ , then we've found an instance of the pattern in the text, and if we run out of text characters ( $i > n$ ) before we reach  $\textcircled{!}$ , then there is no match.

The finite state machine is really just a (very!) convenient metaphor. In a real implementation, we would not construct the entire graph. Since the success edges always go through the pattern characters in order, we only have to remember where the failure edges go. We can encode this *failure function* in an array  $fail[1..n]$ , so that for each  $j$  there is a failure edge from node  $j$  to node  $fail[j]$ . Following a failure edge back to an earlier state exactly corresponds, in our earlier formulation, to shifting the pattern forward. The failure function  $fail[j]$  tells us how far to shift after a character mismatch  $T[i] \neq P[j]$ .

Here's what the actual algorithm looks like:

```

KNUTHMORRISPRATT( $T[1..n], P[1..m]$ ):
   $j \leftarrow 1$ 
  for  $i \leftarrow 1$  to  $n$ 
    while  $j > 0$  and  $T[i] \neq P[j]$ 
       $j \leftarrow fail[j]$ 
    if  $j = m$      $\langle\langle Found\ it!\rangle\rangle$ 
      return  $i - m + 1$ 
     $j \leftarrow j + 1$ 
  return 'none'

```

Before we discuss computing the failure function, let's analyze the running time of KNUTHMORRISPRATT under the assumption that a correct failure function is already known. At each character comparison, either we increase  $i$  and  $j$  by one, or we decrease  $j$  and leave  $i$  alone.

We can increment  $i$  at most  $n - 1$  times before we run out of text, so there are at most  $n - 1$  successful comparisons. Similarly, there can be at most  $n - 1$  failed comparisons, since the number of times we decrease  $j$  cannot exceed the number of times we increment  $j$ . In other words, we can amortize character mismatches against earlier character matches. Thus, the total number of character comparisons performed by KNUTHMORRISPRATT in the worst case is  $O(n)$ .

### D.3 Computing the Failure Function

We can now rephrase our second intuitive rule about how to choose a reasonable shift after a character mismatch  $T[i] \neq P[j]$ :

$P[1..fail[j] - 1]$  is the longest proper prefix of  $P[1..j - 1]$  that is also a suffix of  $T[1..i - 1]$ .

Notice, however, that if we are comparing  $T[i]$  against  $P[j]$ , then we must have already matched the first  $j - 1$  characters of the pattern. In other words, we already know that  $P[1..j - 1]$  is a suffix of  $T[1..i - 1]$ . Thus, we can rephrase the prefix-suffix rule as follows:

$P[1..fail[j] - 1]$  is the longest proper prefix of  $P[1..j - 1]$  that is also a suffix of  $P[1..j - 1]$ .

This is the definition of the Knuth-Morris-Pratt failure function  $fail[j]$  for all  $j > 1$ .<sup>1</sup> By convention we set  $fail[1] = 0$ ; this tells the KMP algorithm that if the first pattern character doesn't match, it should just give up and try the next text character.

$P[i]$	A	B	R	A	C	A	D	A	B	R	A
$fail[i]$	0	1	1	1	2	1	2	1	2	3	4

Failure function for the string 'ABRACADABRA'

(Compare with the finite state machine on the previous page.)

We could easily compute the failure function in  $O(m^3)$  time by checking, for each  $j$ , whether every prefix of  $P[1..j - 1]$  is also a suffix of  $P[1..j - 1]$ , but this is not the fastest method. The following algorithm essentially uses the KMP search algorithm to look for the pattern inside itself!

```

COMPUTEFAILURE( $P[1..m]$ ):
   $j \leftarrow 0$ 
  for  $i \leftarrow 1$  to  $m$ 
     $fail[i] \leftarrow j$       (*)
    while  $j > 0$  and  $P[i] \neq P[j]$ 
       $j \leftarrow fail[j]$ 
     $j \leftarrow j + 1$ 

```

Here's an example of this algorithm in action. In each line, the current values of  $i$  and  $j$  are indicated by superscripts; \$ represents the beginning of the string. (You should imagine pointing at  $P[j]$  with your left hand and pointing at  $P[i]$  with your right hand, and moving your fingers according to the algorithm's directions.)

<sup>1</sup>CLR defines a similar *prefix function*, denoted  $\pi[j]$ , as follows:

$P[1..\pi[j]]$  is the longest proper prefix of  $P[1..j]$  that is also a suffix of  $P[1..j]$ .

These two functions are not the same, but they are related by the simple equation  $\pi[j] = fail[j + 1] - 1$ . The off-by-one difference between the two functions adds a few extra +1s to CLR's version of the algorithm.

$j \leftarrow 0, i \leftarrow 1$ $fail[i] \leftarrow j$	$\$^j$ A <sup>i</sup> B R A C A D A B R X ... 0
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$ $j \leftarrow fail[j]$	$\$$ A <sup>j</sup> B <sup>i</sup> R A C A D A B R X ... 0 1 $\$^j$ A B <sup>i</sup> R A C A D A B R X ...
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$ $j \leftarrow fail[j]$	$\$$ A <sup>j</sup> B R <sup>i</sup> A C A D A B R X ... 0 1 1 $\$^j$ A B R <sup>i</sup> A C A D A B R X ...
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$	$\$$ A <sup>j</sup> B R A <sup>i</sup> C A D A B R X ... 0 1 1 1
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$ $j \leftarrow fail[j]$ $j \leftarrow fail[j]$	$\$$ A B <sup>j</sup> R A C <sup>i</sup> A D A B R X ... 0 1 1 1 2 $\$$ A <sup>j</sup> B R A C <sup>i</sup> A D A B R X ... $\$^j$ A B R A C <sup>i</sup> A D A B R X ...
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$	$\$$ A <sup>j</sup> B R A C A <sup>i</sup> D A B R X ... 0 1 1 1 2 1
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$ $j \leftarrow fail[j]$ $j \leftarrow fail[j]$	$\$$ A B <sup>j</sup> R A C A D <sup>i</sup> A B R X ... 0 1 1 1 2 1 2 $\$$ A <sup>j</sup> B R A C A D <sup>i</sup> A B R X ... $\$^j$ A B R A C A D <sup>i</sup> A B R X ...
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$	$\$$ A <sup>j</sup> B R A C A D A <sup>i</sup> B R X ... 0 1 1 1 2 1 2 1
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$	$\$$ A B <sup>j</sup> R A C A D A B <sup>i</sup> R X ... 0 1 1 1 2 1 2 1 2
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$	$\$$ A B R <sup>j</sup> A C A D A B R <sup>i</sup> X ... 0 1 1 1 2 1 2 1 2 3
$j \leftarrow j + 1, i \leftarrow i + 1$ $fail[i] \leftarrow j$ $j \leftarrow fail[j]$ $j \leftarrow fail[j]$	$\$$ A B R A <sup>j</sup> C A D A B R X <sup>i</sup> ... 0 1 1 1 2 1 2 1 2 3 4 ... $\$$ A <sup>j</sup> B R A C A D A B R X <sup>i</sup> ... $\$^j$ A B R A C A D A B R X <sup>i</sup> ...

COMPUTEFAILURE in action. Do this yourself by hand.

Just as we did for KNUTHMORRISPRATT, we can analyze COMPUTEFAILURE by amortizing character mismatches against earlier character matches. Since there are at most  $m$  character matches, COMPUTEFAILURE runs in  $O(m)$  time.

Let's prove (by induction, of course) that COMPUTEFAILURE correctly computes the failure function. The base case  $fail[1] = 0$  is obvious. Assuming inductively that we correctly computed  $fail[1]$  through  $fail[i]$  in line (\*), we need to show that  $fail[i + 1]$  is also correct. Just after the  $i$ th iteration of line (\*), we have  $j = fail[i]$ , so  $P[1..j - 1]$  is the longest proper prefix of  $P[1..i - 1]$  that is also a suffix.

Let's define the iterated failure functions  $fail^c[j]$  inductively as follows:  $fail^0[j] = j$ , and

$$fail^c[j] = fail[fail^{c-1}[j]] = \overbrace{fail[fail[\dots[fail[j]]\dots]]}^c.$$

In particular, if  $fail^{c-1}[j] = 0$ , then  $fail^c[j]$  is undefined. We can easily show by induction (see [CLR, p.872]) that every string of the form  $P[1..fail^c[j] - 1]$  is both a proper prefix and a proper suffix of  $P[1..i - 1]$ , and in fact, these are the only examples. Thus, the longest proper prefix/suffix of  $P[1..i]$  must be the longest string of the form  $P[1..fail^c[j]]$  — i.e., the one with smallest  $c$  — such that  $P[fail^c[j]] = P[i]$ . This is exactly what the while loop in COMPUTEFAILURE computes; the  $(c + 1)$ th iteration compares  $P[fail^c[j]] = P[fail^{c+1}[i]]$  against  $P[i]$ . COMPUTEFAILURE is actually a *dynamic programming* implementation of the following recursive definition of  $fail[i]$ :

$$fail[i] = \begin{cases} 0 & \text{if } i = 0, \\ \max_{c \geq 1} \{ fail^c[i - 1] + 1 \mid P[i - 1] = P[fail^c[i - 1]] \} & \text{otherwise.} \end{cases}$$

## D.4 Optimizing the Failure Function

We can speed up KNUTHMORRISPRATT slightly by making one small change to the failure function. Recall that after comparing  $T[i]$  against  $P[j]$  and finding a mismatch, the algorithm compares  $T[i]$  against  $P[fail[j]]$ . With the current definition, however, it is possible that  $P[j]$  and  $P[fail[j]]$  are actually the same character, in which case the next character comparison will automatically fail. So why do the comparison at all?

We can optimize the failure function by ‘short-circuiting’ these redundant comparisons with some simple post-processing:

```

OPTIMIZEFAILURE( $P[1..m]$ ,  $fail[1..m]$ ):
  for  $i \leftarrow 2$  to  $m$ 
    if  $P[i] = P[fail[i]]$ 
       $fail[i] \leftarrow fail[fail[i]]$ 

```

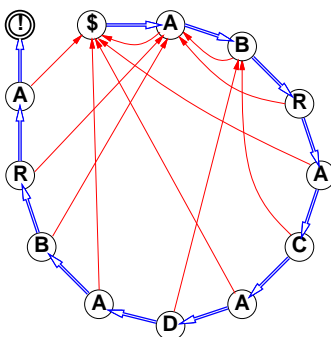
We can also compute the optimized failure function directly by adding three new lines (in bold) to the COMPUTEFAILURE function.

```

COMPUTEOPTFAILURE( $P[1..m]$ ):
   $j \leftarrow 0$ 
  for  $i \leftarrow 1$  to  $m$ 
    if  $P[i] = P[j]$ 
       $fail[i] \leftarrow fail[j]$ 
    else
       $fail[i] \leftarrow j$ 
    while  $j > 0$  and  $P[i] \neq P[j]$ 
       $j \leftarrow fail[j]$ 
     $j \leftarrow j + 1$ 

```

This optimization slows down the preprocessing slightly, but it may significantly decrease the number of comparisons at each text character. The worst-case running time is still  $O(n)$ ; however, the constant is about half as big as for the unoptimized version, so this could be a significant improvement in practice.



Optimized finite state machine for the string ‘ABRADACABRA’

$P[i]$	A	B	R	A	C	A	D	A	B	R	A
$fail[i]$	0	<b>1</b>	<b>1</b>	<b>0</b>	<b>2</b>	<b>0</b>	<b>2</b>	<b>0</b>	<b>1</b>	<b>1</b>	<b>0</b>

Optimized failure function for ‘ABRADACABRA’, with changes in bold.

Here are the unoptimized and optimized failure functions for a few more patterns:

$P[i]$	A	N	A	N	A	B	A	N	A	N	A	N	A
unoptimized $fail[i]$	0	1	1	2	3	4	1	2	3	4	5	6	5
optimized $fail[i]$	0	1	0	1	0	4	0	1	0	1	0	6	0

Failure functions for 'ANANABANANANA'.

$P[i]$	A	B	A	B	C	A	B	A	B	C	A	B	C
unoptimized $fail[i]$	0	1	1	2	3	1	2	3	4	5	6	7	8
optimized $fail[i]$	0	1	0	1	3	0	1	0	1	3	0	1	8

Failure functions for 'ABABCABABCABC'.

$P[i]$	A	B	B	A	B	B	A	B	A	B	B	A	B
unoptimized $fail[i]$	0	1	1	1	2	3	4	5	6	2	3	4	5
optimized $fail[i]$	0	1	1	0	1	1	0	1	6	1	1	0	1

Failure functions for 'ABBABBABABBAB'.

$P[i]$	A	A	A	A	A	A	A	A	A	A	A	A	B
unoptimized $fail[i]$	0	1	2	3	4	5	6	7	8	9	10	11	12
optimized $fail[i]$	0	0	0	0	0	0	0	0	0	0	0	0	12

Failure functions for 'AAAAAAAAAAAAAB'.



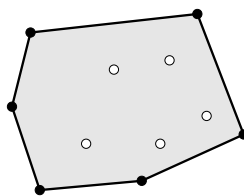
## E Convex Hulls

### E.1 Definitions

We are given a set  $P$  of  $n$  points in the plane. We want to compute something called the *convex hull* of  $P$ . Intuitively, the convex hull is what you get by driving a nail into the plane at each point and then wrapping a piece of string around the nails. More formally, the convex hull is the smallest convex polygon containing the points:

- **polygon:** A region of the plane bounded by a cycle of line segments, called *edges*, joined end-to-end in a cycle. Points where two successive edges meet are called *vertices*.
- **convex:** For any two points  $p, q$  inside the polygon, the line segment  $\overline{pq}$  is completely inside the polygon.
- **smallest:** Any convex proper subset of the convex hull excludes at least one point in  $P$ . This implies that every vertex of the convex hull is a point in  $P$ .

We can also define the convex hull as the *largest* convex polygon whose vertices are all points in  $P$ , or the *unique* convex polygon that contains  $P$  and whose vertices are all points in  $P$ . Notice that  $P$  might have *interior* points that are not vertices of the convex hull.



A set of points and its convex hull.  
Convex hull vertices are black; interior points are white.

Just to make things concrete, we will represent the points in  $P$  by their Cartesian coordinates, in two arrays  $X[1..n]$  and  $Y[1..n]$ . We will represent the convex hull as a circular linked list of vertices in counterclockwise order. If the  $i$ th point is a vertex of the convex hull,  $next[i]$  is index of the next vertex counterclockwise and  $pred[i]$  is the index of the next vertex clockwise; otherwise,  $next[i] = pred[i] = 0$ . It doesn't matter which vertex we choose as the 'head' of the list. The decision to list vertices counterclockwise instead of clockwise is arbitrary.

To simplify the presentation of the convex hull algorithms, I will assume that the points are in *general position*, meaning (in this context) that *no three points lie on a common line*. This is just like assuming that no two elements are equal when we talk about sorting algorithms. If we wanted to really implement these algorithms, we would have to handle colinear triples correctly, or at least consistently. This is fairly easy, but definitely not trivial.

### E.2 Simple Cases

Computing the convex hull of a single point is trivial; we just return that point. Computing the convex hull of two points is also trivial.

For three points, we have two different possibilities — either the points are listed in the array in clockwise order or counterclockwise order. Suppose our three points are  $(a, b)$ ,  $(c, d)$ , and  $(e, f)$ , given in that order, and for the moment, let's also suppose that the first point is furthest to the

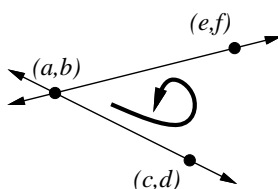
left, so  $a < c$  and  $a < f$ . Then the three points are in counterclockwise order if and only if the line  $\overleftrightarrow{(a,b)(c,d)}$  is less than the slope of the line  $\overleftrightarrow{(a,b)(e,f)}$ :

$$\text{counterclockwise} \iff \frac{d-b}{c-a} < \frac{f-b}{e-a}$$

Since both denominators are positive, we can rewrite this inequality as follows:

$$\text{counterclockwise} \iff (f-b)(c-a) > (d-b)(e-a)$$

This final inequality is correct even if  $(a,b)$  is not the leftmost point. If the inequality is reversed, then the points are in clockwise order. If the three points are colinear (remember, we're assuming that never happens), then the two expressions are equal.



Three points in counterclockwise order.

Another way of thinking about this counterclockwise test is that we're computing the *cross-product* of the two vectors  $(c,d) - (a,b)$  and  $(e,f) - (a,b)$ , which is defined as a  $2 \times 2$  determinant:

$$\text{counterclockwise} \iff \begin{vmatrix} c-a & d-b \\ e-a & f-b \end{vmatrix} > 0$$

We can also write it as a  $3 \times 3$  determinant:

$$\text{counterclockwise} \iff \begin{vmatrix} 1 & a & b \\ 1 & c & d \\ 1 & e & f \end{vmatrix} > 0$$

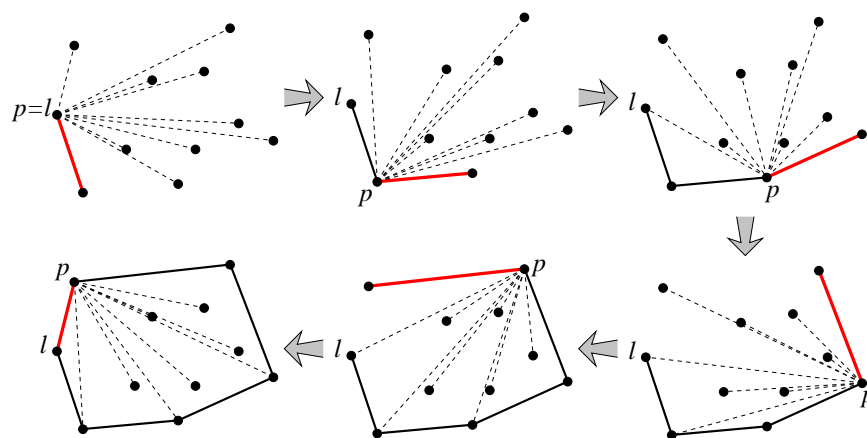
All three boxed expressions are algebraically identical.

This counterclockwise test plays *exactly* the same role in convex hull algorithms as comparisons play in sorting algorithms. Computing the convex hull of three points is analogous to sorting two numbers: either they're in the correct order or in the opposite order.

### E.3 Jarvis's Algorithm (Wrapping)

Perhaps the simplest algorithm for computing convex hulls simply simulates the process of wrapping a piece of string around the points. This algorithm is usually called *Jarvis's march*, but it is also referred to as the *gift-wrapping* algorithm.

Jarvis's march starts by computing the leftmost point  $\ell$ , *i.e.*, the point whose  $x$ -coordinate is smallest, since we know that the left most point must be a convex hull vertex. Finding  $\ell$  clearly takes linear time.



The execution of Jarvis's March.

Then the algorithm does a series of *pivoting* steps to find each successive convex hull vertex, starting with  $\ell$  and continuing until we reach  $\ell$  again. The vertex immediately following a point  $p$  is the point that appears to be furthest to the right to someone standing at  $p$  and looking at the other points. In other words, if  $q$  is the vertex following  $p$ , and  $r$  is any other input point, then the triple  $p, q, r$  is in counter-clockwise order. We can find each successive vertex in linear time by performing a series of  $O(n)$  counter-clockwise tests.

```

JARVISMARCH( $X[1..n], Y[1..n]$ ):
 $\ell \leftarrow 1$ 
for  $i \leftarrow 2$  to  $n$ 
    if  $X[i] < X[\ell]$ 
         $\ell \leftarrow i$ 

 $p \leftarrow \ell$ 
repeat
     $q \leftarrow p + 1$      $\langle\langle$ Make sure  $p \neq q$  $\rangle\rangle$ 
    for  $i \leftarrow 2$  to  $n$ 
        if CCW( $p, i, q$ )
             $q \leftarrow i$ 
     $next[p] \leftarrow q$ ;  $prev[q] \leftarrow p$ 
     $p \leftarrow q$ 
until  $p = \ell$ 

```

Since the algorithm spends  $O(n)$  time for each convex hull vertex, the worst-case running time is  $O(n^2)$ . However, this naïve analysis hides the fact that if the convex hull has very few vertices, Jarvis's march is extremely fast. A better way to write the running time is  $O(nh)$ , where  $h$  is the number of convex hull vertices. In the worst case,  $h = n$ , and we get our old  $O(n^2)$  time bound, but in the best case  $h = 3$ , and the algorithm only needs  $O(n)$  time. Computational geometers call this an *output-sensitive* algorithm; the smaller the output, the faster the algorithm.

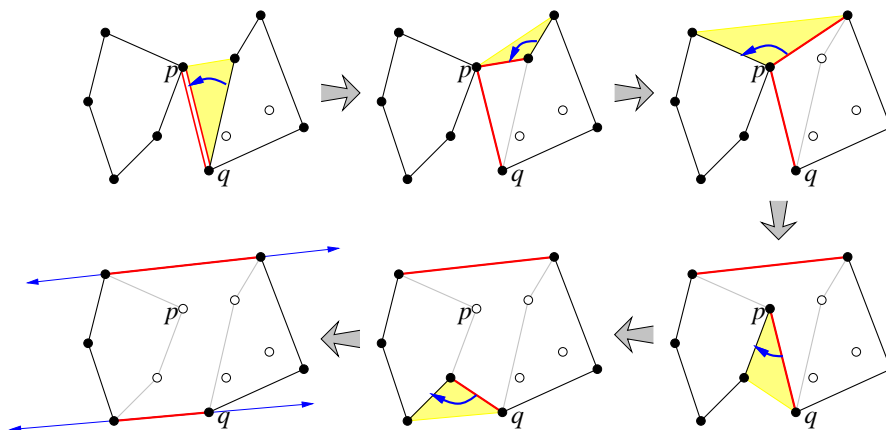
#### E.4 Divide and Conquer (Splitting)

The behavior of Jarvis's march is very much like selection sort: repeatedly find the item that goes in the next slot. In fact, most convex hull algorithms resemble some sorting algorithm.

For example, the following convex hull algorithm resembles quicksort. We start by choosing a *pivot* point  $p$ . Partitions the input points into two sets  $L$  and  $R$ , containing the points to the left

of  $p$ , including  $p$  itself, and the points to the right of  $p$ , by comparing  $x$ -coordinates. Recursively compute the convex hulls of  $L$  and  $R$ . Finally, merge the two convex hulls into the final output.

The merge step requires a little explanation. We start by connecting the two hulls with a line segment between the rightmost point of the hull of  $L$  with the leftmost point of the hull of  $R$ . Call these points  $p$  and  $q$ , respectively. (Yes, it's the same  $p$ .) Actually, let's add *two* copies of the segment  $\overline{pq}$  and call them *bridges*. Since  $p$  and  $q$  can 'see' each other, this creates a sort of dumbbell-shaped polygon, which is convex except possibly at the endpoints off the bridges.

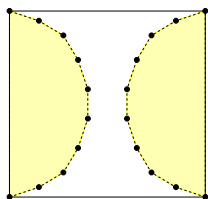


Merging the left and right subhulls.

We now expand this dumbbell into the correct convex hull as follows. As long as there is a clockwise turn at either endpoint of either bridge, we remove that point from the circular sequence of vertices and connect its two neighbors. As soon as the turns at both endpoints of both bridges are counter-clockwise, we can stop. At that point, the bridges lie on the *upper* and *lower common tangent* lines of the two subhulls. These are the two lines that touch both subhulls, such that both subhulls lie below the upper common tangent line and above the lower common tangent line.

Merging the two subhulls takes  $O(n)$  time in the worst case. Thus, the running time is given by the recurrence  $T(n) = O(n) + T(k) + T(n - k)$ , just like quicksort, where  $k$  the number of points in  $R$ . Just like quicksort, if we use a naïve deterministic algorithm to choose the pivot point  $p$ , the worst-case running time of this algorithm is  $O(n^2)$ . If we choose the pivot point randomly, the expected running time is  $O(n \log n)$ .

There are inputs where this algorithm is clearly wasteful (at least, clearly to *us*). If we're really unlucky, we'll spend a long time putting together the subhulls, only to throw almost everything away during the merge step. Thus, this divide-and-conquer algorithm is *not* output sensitive.

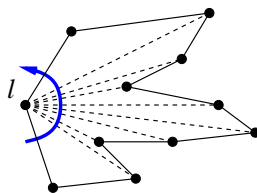


A set of points that shouldn't be divided and conquered.

## E.5 Graham's Algorithm (Scanning)

Our third convex hull algorithm, called *Graham's scan*, first explicitly sorts the points in  $O(n \log n)$  and then applies a linear-time scanning algorithm to finish building the hull.

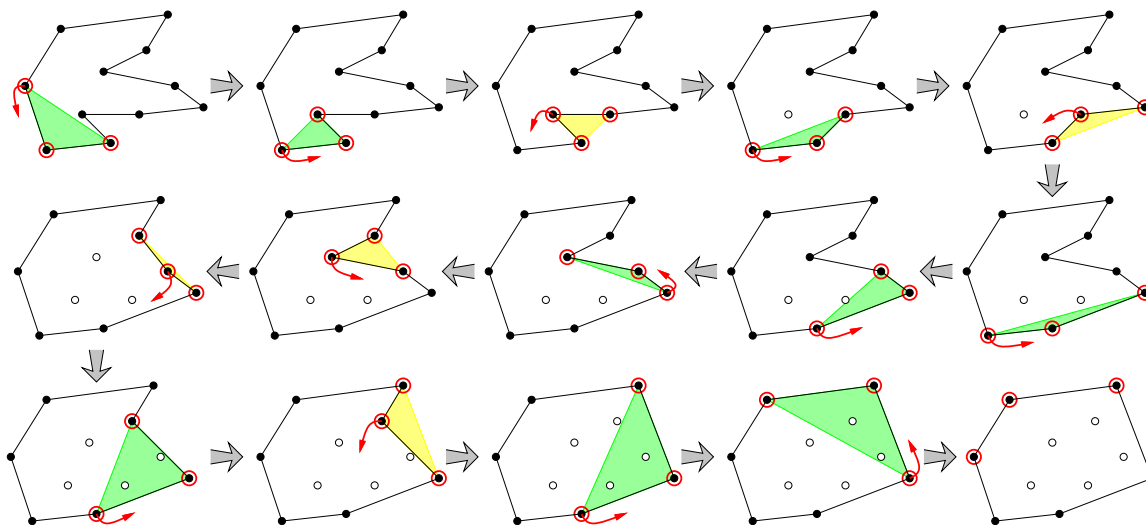
We start Graham's scan by finding the leftmost point  $\ell$ , just as in Jarvis's march. Then we sort the points in counterclockwise order around  $\ell$ . We can do this in  $O(n \log n)$  time with any comparison-based sorting algorithm (quicksort, mergesort, heapsort, etc.). To compare two points  $p$  and  $q$ , we check whether the triple  $\ell, p, q$  is oriented clockwise or counterclockwise. Once the points are sorted, we connected them in counterclockwise order, starting and ending at  $\ell$ . The result is a *simple* polygon with  $n$  vertices.



A simple polygon formed in the sorting phase of Graham's scan.

To change this polygon into the convex hull, we apply the following 'three-penny algorithm'. We have three pennies, which will sit on three consecutive vertices  $p, q, r$  of the polygon; initially, these are  $\ell$  and the two vertices after  $\ell$ . We now apply the following two rules over and over until a penny is moved forward onto  $\ell$ :

- If  $p, q, r$  are in counterclockwise order, move the back penny forward to the successor of  $r$ .
- If  $p, q, r$  are in clockwise order, remove  $q$  from the polygon, add the edge  $pr$ , and move the middle penny backward to the predecessor of  $p$ .



The 'three-penny' scanning phase of Graham's scan.

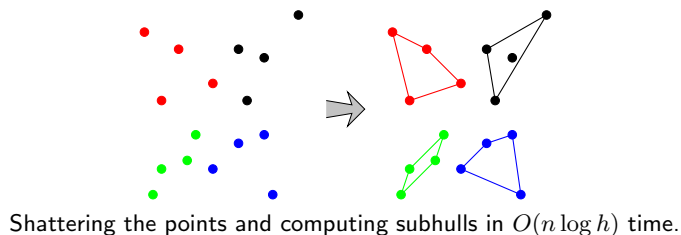
Whenever a penny moves forward, it moves onto a vertex that hasn't seen a penny before (except the last time), so the first rule is applied  $n - 2$  times. Whenever a penny moves backwards, a vertex is removed from the polygon, so the second rule is applied exactly  $n - h$  times, where  $h$  is as usual the number of convex hull vertices. Since each counterclockwise test takes constant time, the scanning phase takes  $O(n)$  time altogether.

## E.6 Chan's Algorithm (Shattering)

The last algorithm I'll describe is an output-sensitive algorithm that is never slower than either Jarvis's march or Graham's scan. The running time of this algorithm, which was discovered by

Timothy Chan in 1993, is  $O(n \log h)$ . Chan’s algorithm is a combination of divide-and-conquer and gift-wrapping.

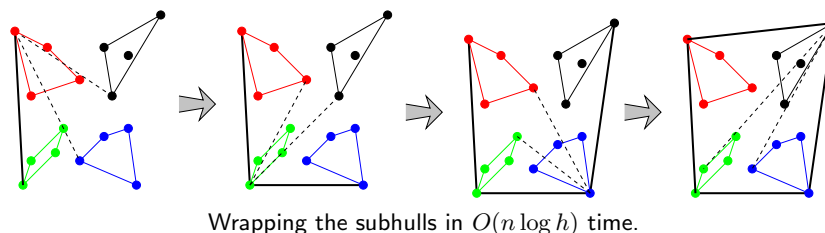
First suppose a ‘little birdie’ tells us the value of  $h$ ; we’ll worry about how to implement the little birdie in a moment. Chan’s algorithm starts by *shattering* the input points into  $n/h$  arbitrary<sup>1</sup> subsets, each of size  $h$ , and computing the convex hull of each subset using (say) Graham’s scan. This much of the algorithm requires  $O((n/h) \cdot h \log h) = O(n \log h)$  time.



Once we have the  $n/h$  subhulls, we follow the general outline of Jarvis’s march, ‘wrapping a string around’ the  $n/h$  subhulls. Starting with  $p = \ell$ , where  $\ell$  is the leftmost input point, we successively find the convex hull vertex that follows  $p$  and counterclockwise order until we return back to  $\ell$  again.

The vertex that follows  $p$  is the point that appears to be furthest to the right to someone standing at  $p$ . This means that the successor of  $p$  must lie on a *right tangent line* between  $p$  and one of the subhulls—a line from  $p$  through a vertex of the subhull, such that the subhull lies completely on the right side of the line from  $p$ ’s point of view. We can find the right tangent line between  $p$  and any subhull in  $O(\log h)$  time using a variant of binary search. (This is a practice problem in the homework!) Since there are  $n/h$  subhulls, finding the successor of  $p$  takes  $O((n/h) \log h)$  time altogether.

Since there are  $h$  convex hull edges, and we find each edge in  $O((n/h) \log h)$  time, the overall running time of the algorithm is  $O(n \log h)$ .



Unfortunately, this algorithm only takes  $O(n \log h)$  time if a little birdie has told us the value of  $h$  in advance. So how do we implement the ‘little birdie’? Chan’s trick is to *guess* the correct value of  $h$ ; let’s denote the guess by  $h^*$ . Then we shatter the points into  $n/h^*$  subsets of size  $h^*$ , compute their subhulls, and then find the first  $h^*$  edges of the global hull. If  $h < h^*$ , this algorithm computes the complete convex hull in  $O(n \log h^*)$  time. Otherwise, the hull doesn’t wrap all the way back around to  $\ell$ , so we know our guess  $h^*$  is too small.

Chan’s algorithm starts with the optimistic guess  $h^* = 3$ . If we finish an iteration of the algorithm and find that  $h^*$  is too small, we *square*  $h^*$  and try again. In the final iteration,  $h^* < h^2$ , so the last iteration takes  $O(n \log h^*) = O(n \log h^2) = O(n \log h)$  time. The total running time of Chan’s algorithm is given by the sum

$$O(n \log 3 + n \log 3^2 + n \log 3^4 + \cdots + n \log 3^{2^k}),$$

<sup>1</sup>In the figures, in order to keep things as clear as possible, I’ve chosen these subsets so that their convex hulls are disjoint. This is not true in general!

for some integer  $k$ . We can rewrite this as a geometric series:

$$O(n \log 3 + 2n \log 3 + 4n \log 3 + \cdots + 2^k n \log 3).$$

Since any geometric series adds up to a constant times its largest term, the total running time is a constant times the time taken by the last iteration, which is  $O(n \log h)$ . So Chan's algorithm runs in  $O(n \log h)$  time overall, even without the little birdie.

## F Line Segment Intersection

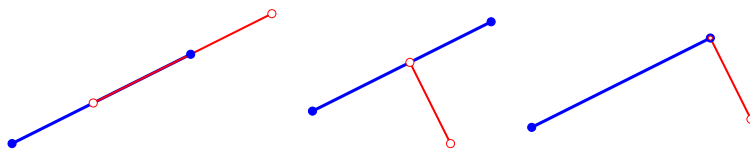
### F.1 Introduction

In this lecture, I'll talk about detecting line segment intersections. A line segment is the convex hull of two points, called the *endpoints* (or *vertices*) of the segment. We are given a set of  $n$  line segments, each specified by the  $x$ - and  $y$ -coordinates of its endpoints, for a total of  $4n$  real numbers, and we want to know whether any two segments intersect.

To keep things simple, just as in the previous lecture, I'll assume the segments are in *general position*.

- No three endpoints lie on a common line.
- No two endpoints have the same  $x$ -coordinate. In particular, no segment is vertical, no segment is just a point, and no two segments share an endpoint.

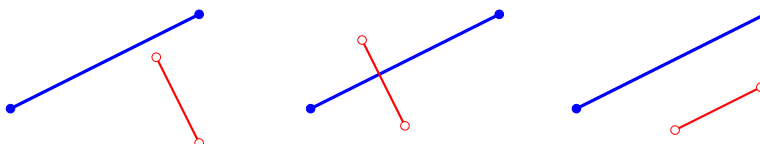
This general position assumption lets us avoid several annoying degenerate cases. Of course, in any real implementation of the algorithm I'm about to describe, you'd have to handle these cases. Real-world data is *full* of degeneracies!



Degenerate cases of intersecting segments that we'll pretend never happen: Overlapping colinear segments, endpoints inside segments, and shared endpoints.

### F.2 Two segments

The first case we have to consider is  $n = 2$ . ( $n \leq 1$  is obviously completely trivial!) How do we tell whether two line segments intersect? One possibility, suggested by a student in class, is to construct the convex hull of the segments. Two segments intersect if and only if the convex hull is a quadrilateral whose vertices alternate between the two segments. In the figure below, the first pair of segments has a triangular convex hull. The last pair's convex hull is a quadrilateral, but its vertices don't alternate.



Some pairs of segments.

Fortunately, we don't need (or want!) to use a full-fledged convex hull algorithm just to test two segments; there's a much simpler test.

Two segments  $\overline{ab}$  and  $\overline{cd}$  intersect if and only if

- the endpoints  $a$  and  $b$  are on opposite sides of the line  $\overleftrightarrow{cd}$ , and
- the endpoints  $c$  and  $d$  are on opposite sides of the line  $\overleftrightarrow{ab}$ .



To test whether two points are on opposite sides of a line through two other points, we use the same counterclockwise test that we used for building convex hulls. Specifically,  $a$  and  $b$  are on opposite sides of line  $\overleftrightarrow{cd}$  if and only if exactly one of the two triples  $a, c, d$  and  $b, c, d$  is in counterclockwise order. So we have the following simple algorithm.

```

INTERSECT( $a, b, c, d$ ):
  if  $CCW(a, c, d) = CCW(b, c, d)$ 
    return FALSE
  else if  $CCW(a, b, c) = CCW(a, b, d)$ 
    return FALSE
  else
    return TRUE

```

Or even simpler:

```

INTERSECT( $a, b, c, d$ ):
  return  $[CCW(a, c, d) \neq CCW(b, c, d)] \wedge [CCW(a, b, c) \neq CCW(a, b, d)]$ 

```

### F.3 A Sweep Line Algorithm

To detect whether there's an intersection in a set of more than just two segments, we use something called a *sweep line* algorithm. First let's give each segment a unique *label*. I'll use letters, but in a real implementation, you'd probably use pointers/references to records storing the endpoint coordinates.

Imagine sweeping a vertical line across the segments from left to right. At each position of the sweep line, look at the sequence of (labels of) segments that the line hits, sorted from top to bottom. The only times this sorted sequence can change is when the sweep line passes an endpoint or when the sweep line passes an intersection point. In the second case, the order changes because two adjacent labels swap places.<sup>1</sup> Our algorithm will simulate this sweep, looking for potential swaps between adjacent segments.

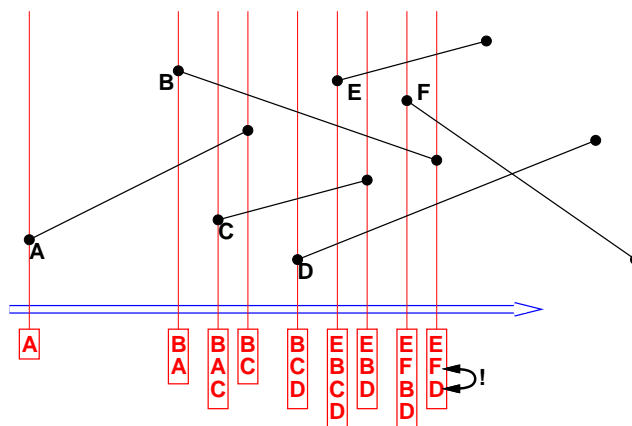
The sweep line algorithm begins by sorting the  $2n$  segment endpoints from left to right by comparing their  $x$ -coordinates, in  $O(n \log n)$  time. The algorithm then moves the sweep line from left to right, stopping at each endpoint.

We store the vertical label sequence in some sort of balanced binary tree that supports the following operations in  $O(\log n)$  time. Note that the tree does not store any explicit search keys, only segment labels.

- **Insert** a segment label.
- **Delete** a segment label.
- Find the **neighbors** of a segment label in the sorted sequence.

$O(\log n)$  amortized time is good enough, so we could use a scapegoat tree or a splay tree. If we're willing to settle for an expected time bound, we could use a treap or a skip list instead.

<sup>1</sup>Actually, if more than two segments intersect at the same point, there could be a larger reversal, but this won't have any effect on our algorithm.



The sweep line algorithm in action. The boxes show the label sequence stored in the binary tree. The intersection between F and D is detected in the last step.

Whenever the sweep line hits a left endpoint, we insert the corresponding label into the tree in  $O(\log n)$  time. In order to do this, we have to answer questions of the form ‘Does the new label X go above or below the old label Y?’ To answer this question, we test whether the new left endpoint of X is above segment Y, or equivalently, if the triple of endpoints  $\text{left}(Y)$ ,  $\text{right}(Y)$ ,  $\text{left}(X)$  is in counterclockwise order.

Once the new label is inserted, we test whether the new segment intersects either of its two neighbors in the label sequence. For example, in the figure above, when the sweep line hits the left endpoint of F, we test whether F intersects either B or E. These tests require  $O(1)$  time.

Whenever the sweep line hits a right endpoint, we delete the corresponding label from the tree in  $O(\log n)$  time, and then check whether its two neighbors intersect in  $O(1)$  time. For example, in the figure, when the sweep line hits the right endpoint of C, we test whether B and D intersect.

If at any time we discover a pair of segments that intersects, we stop the algorithm and report the intersection. For example, in the figure, when the sweep line reaches the right endpoint of B, we discover that F and D intersect, and we halt. Note that we may not discover the intersection until long after the two segments are inserted, and the intersection we discover may not be the one that the sweep line would hit first. It’s not hard to show by induction (hint, hint) that the algorithm is correct. Specifically, if the algorithm reaches the  $n$ th right endpoint without detecting an intersection, none of the segments intersect.

For each segment endpoint, we spend  $O(\log n)$  time updating the binary tree, plus  $O(1)$  time performing pairwise intersection tests—at most two at each left endpoint and at most one at each right endpoint. Thus, the entire sweep requires  $O(n \log n)$  time in the worst case. Since we also spent  $O(n \log n)$  time sorting the endpoints, the overall running time is  $O(n \log n)$ .

Here’s a slightly more formal description of the algorithm. The input  $S[1..n]$  is an array of line segments. The sorting phase in the first line produces two auxiliary arrays:

- $\text{label}[i]$  is the label of the  $i$ th leftmost endpoint. I’ll use indices into the input array  $S$  as the labels, so the  $i$ th vertex is an endpoint of  $S[\text{label}[i]]$ .
- $\text{isleft}[i]$  is TRUE if the  $i$ th leftmost endpoint is a left endpoint and FALSE if it’s a right endpoint.

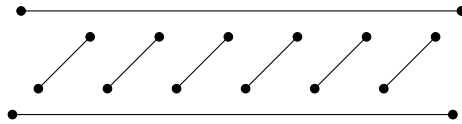
The functions INSERT, DELETE, PREDECESSOR, and SUCCESSOR modify or search through the sorted label sequence. Finally, INTERSECT tests whether two line segments intersect.

```

ANYINTERSECTIONS( $S[1..n]$ ):
  sort the endpoints of  $S$  from left to right
  create an empty label sequence
  for  $i \leftarrow 1$  to  $2n$ 
     $\ell \leftarrow \text{label}[i]$ 
    if isleft[ $i$ ]
      INSERT( $\ell$ )
      if INTERSECT( $S[\ell], S[\text{SUCCESSOR}(\ell)]$ )
        return TRUE
      if INTERSECT( $S[\ell], S[\text{PREDECESSOR}(\ell)]$ )
        return TRUE
    else
      if INTERSECT( $S[\text{SUCCESSOR}(\ell)], S[\text{PREDECESSOR}(\ell)]$ )
        return TRUE
      DELETE( $\text{label}[i]$ )
  return FALSE

```

Note that the algorithm doesn't try to avoid redundant pairwise tests. In the figure below, the top and bottom segments would be checked  $n - 1$  times, once at the top left endpoint, and once at the right endpoint of every short segment. But since we've already spent  $O(n \log n)$  time just sorting the inputs,  $O(n)$  redundant segment intersection tests make no difference in the overall running time.

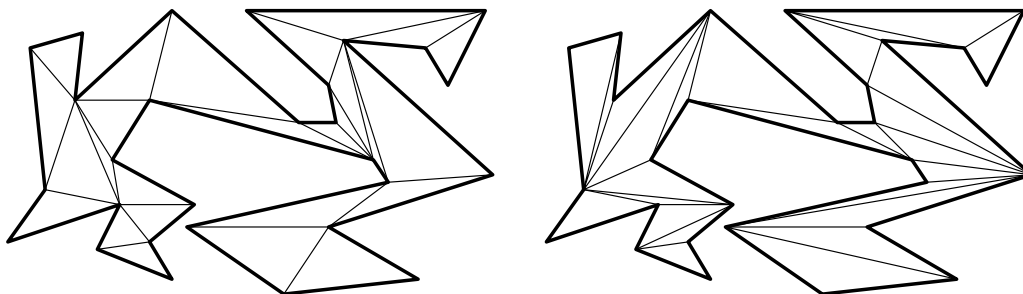


The same pair of segments might be tested  $n - 1$  times.

## G Polygon Triangulation

### G.1 Introduction

Recall from last time that a *polygon* is a region of the plane bounded by a cycle of straight edges joined end to end. Given a polygon, we want to decompose it into triangles by adding *diagonals*: new line segments between the vertices that don't cross the boundary of the polygon. Because we want to keep the number of triangles small, we don't allow the diagonals to cross. We call this decomposition a *triangulation* of the polygon. Most polygons can have more than one triangulation; we don't care which one we compute.

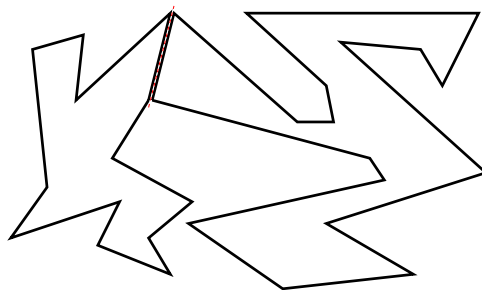


Two triangulations of the same polygon.

Before we go any further, I encourage you to play around with some examples. Draw a few polygons (making sure that the edges are straight and don't cross) and try to break them up into triangles.

### G.2 Existence and Complexity

If you play around with a few examples, you quickly discover that every triangulation of an  $n$ -sided polygon has  $n - 2$  triangles. You might even try to prove this observation by induction. The base case  $n = 3$  is trivial: there is only one triangulation of a triangle, and it obviously has only one triangle! To prove the general case, let  $P$  be a polygon with  $n$  edges. Draw a diagonal between two vertices. This splits  $P$  into two smaller polygons. One of these polygons has  $k$  edges of  $P$  plus the diagonal, for some integer  $k$  between 2 and  $n - 2$ , for a total of  $k + 1$  edges. So by the induction hypothesis, this polygon can be broken into  $k - 1$  triangles. The other polygon has  $n - k + 1$  edges, and so by the induction hypothesis, it can be broken into  $n - k - 1$  triangles. Putting the two pieces back together, we have a total of  $(k - 1) + (n - k - 1) = n - 2$  triangles.



Breaking a polygon into two smaller polygons with a diagonal.

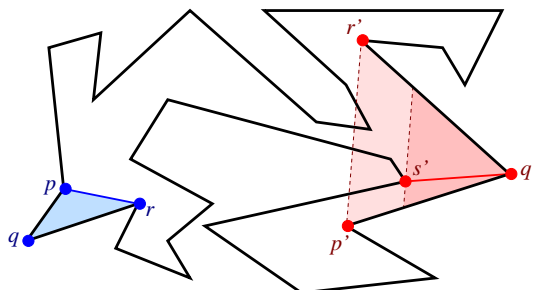
This is a fine induction proof, which any of you could have discovered on your own (right?), except for one small problem. How do we know that every polygon *has* a diagonal? This seems

patently obvious, but it's surprisingly hard to prove, and in fact many incorrect proofs were actually published as late as 1975. The following proof is due to Meisters in 1975.

**Lemma 1.** *Every polygon with more than three vertices has a diagonal.*

**Proof:** Let  $P$  be a polygon with more than three vertices. Every vertex of a  $P$  is either *convex* or *concave*, depending on whether it points into or out of  $P$ , respectively. Let  $q$  be a convex vertex, and let  $p$  and  $r$  be the vertices on either side of  $q$ . For example, let  $q$  be the leftmost vertex. (If there is more than one leftmost vertex, let  $q$  be the lowest one.) If  $\overline{pr}$  is a diagonal, we're done; in this case, we say that the triangle  $\triangle pqr$  is an *ear*.

If  $pr$  is not a diagonal, then  $\triangle pqr$  must contain another vertex of the polygon. Out of all the vertices inside  $\triangle pqr$ , let  $s$  be the vertex furthest away from the line  $\overleftrightarrow{pr}$ . In other words, if we take a line parallel to  $\overleftrightarrow{pr}$  through  $q$ , and translate it towards  $\overleftrightarrow{pr}$ , then  $s$  is the first vertex that the line hits. Then the line segment  $\overline{qs}$  is a diagonal.  $\square$



The leftmost vertex  $q$  is the tip of an ear, so  $pr$  is a diagonal.  
The rightmost vertex  $q'$  is not, since  $\triangle p'q'r'$  contains three other vertices. In this case,  $q's'$  is a diagonal.

### G.3 Existence and Complexity

Meister's existence proof immediately gives us an algorithm to compute a diagonal in linear time. The input to our algorithm is just an array of vertices in counterclockwise order around the polygon. First, we can find the (lowest) leftmost vertex  $q$  in  $O(n)$  time by comparing the  $x$ -coordinates of the vertices (using  $y$ -coordinates to break ties). Next, we can determine in  $O(n)$  time whether the triangle  $\triangle pqr$  contains any of the other  $n - 3$  vertices. Specifically, we can check whether one point lies inside a triangle by performing three counterclockwise tests. Finally, if the triangle is not empty, we can find the vertex  $s$  in  $O(n)$  time by comparing the areas of every triangle  $\triangle pqs$ ; we can compute this area using the counterclockwise determinant.

Here's the algorithm in excruciating detail. We need three support subroutines to compute the area of a polygon, to determine if three points are in counterclockwise order, and to determine if a point is inside a triangle.

«Return twice the signed area of  $\triangle P[i]P[j]P[k]$ »

AREA( $i, j, k$ ):  
return  $(P[k].y - P[i].y)(P[j].x - P[i].x) - (P[k].x - P[i].x)(P[j].y - P[i].y)$

«Are  $P[i], P[j], P[k]$  in counterclockwise order?»

CCW( $i, j, k$ ):  
return  $\text{AREA}(i, j, k) > 0$

«Is  $P[i]$  inside  $\triangle P[p]P[q]P[r]$ ?»

INSIDE( $i, p, q, r$ ):  
return  $\text{CCW}(i, p, q)$  and  $\text{CCW}(i, q, r)$  and  $\text{CCW}(i, r, p)$

```

FINDDIAGONAL( $P[1..n]$ ):
   $q \leftarrow 1$ 
  for  $i \leftarrow 2$  to  $n$ 
    if  $P[i].x < P[q].x$ 
       $q \leftarrow i$ 
   $p \leftarrow q - 1 \bmod n$ 
   $r \leftarrow q + 1 \bmod n$ 

   $ear \leftarrow \text{TRUE}$ 
   $s \leftarrow p$ 
  for  $i \leftarrow 1$  to  $n$ 
    if  $i \leq p$  and  $i \neq q$  and  $i \neq r$  and  $\text{INSIDE}(i, p, q, r)$ 
       $ear \leftarrow \text{FALSE}$ 
      if  $\text{AREA}(i, r, p) > \text{AREA}(s, r, p)$ 
         $s \leftarrow i$ 

  if  $ear = \text{TRUE}$ 
    return  $(p, r)$ 
  else
    return  $(q, s)$ 

```

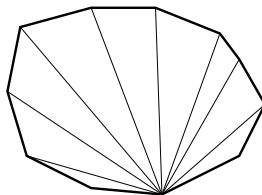
Once we have a diagonal, we can recursively triangulate the two pieces. The worst-case running time of this algorithm satisfies almost the same recurrence as quicksort:

$$T(n) \leq \max_{2 \leq k \leq n-2} T(k+1) + T(n-k+1) + O(n).$$

Just like quicksort, the solution is  $T(n) = O(n^2)$ . So we can now triangulate any polygon in quadratic time.

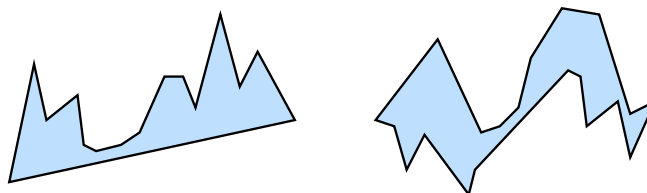
#### G.4 Faster Special Cases

For certain special cases of polygons, we can do much better than  $O(n^2)$  time. For example, we can easily triangulate any convex polygon by connecting any vertex to every other vertex. Since we're given the counterclockwise order of the vertices as input, this takes only  $O(n)$  time.



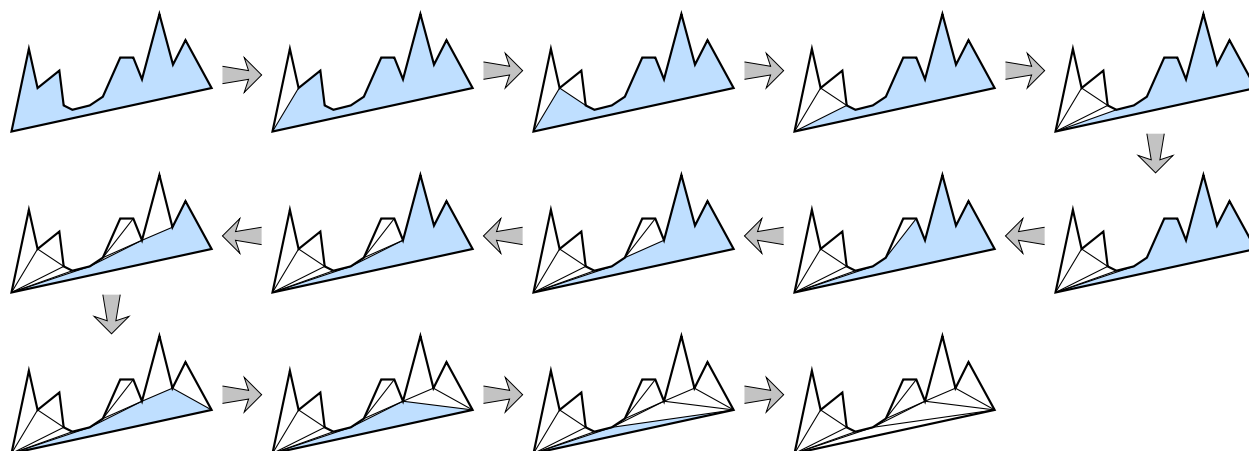
Triangulating a convex polygon is easy.

Another easy special case is *monotone mountains*. A polygon is *monotone* if any vertical line intersects the boundary in at most two points. A monotone polygon is a *mountain* if it contains an edge from the rightmost vertex to the leftmost vertex. Every monotone polygon consists of two chains of edges going left to right between the two extreme vertices; for mountains, one of these chains is a single edge.



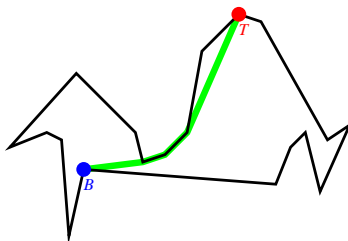
A monotone mountain and a monotone non-mountain.

Triangulating a monotone mountain is extremely easy, since every convex vertex is the tip of an ear, except possibly for the vertices on the far left and far right. Thus, all we have to do is scan through the intermediate vertices, and when we find a convex vertex, cut off the ear. The simplest method for doing this is probably the three-penny algorithm used in the “Graham’s scan” convex hull algorithm—instead of filling in the outside of a polygon with triangles, we’re filling in the inside, both otherwise it’s the same process. This takes  $O(n)$  time.



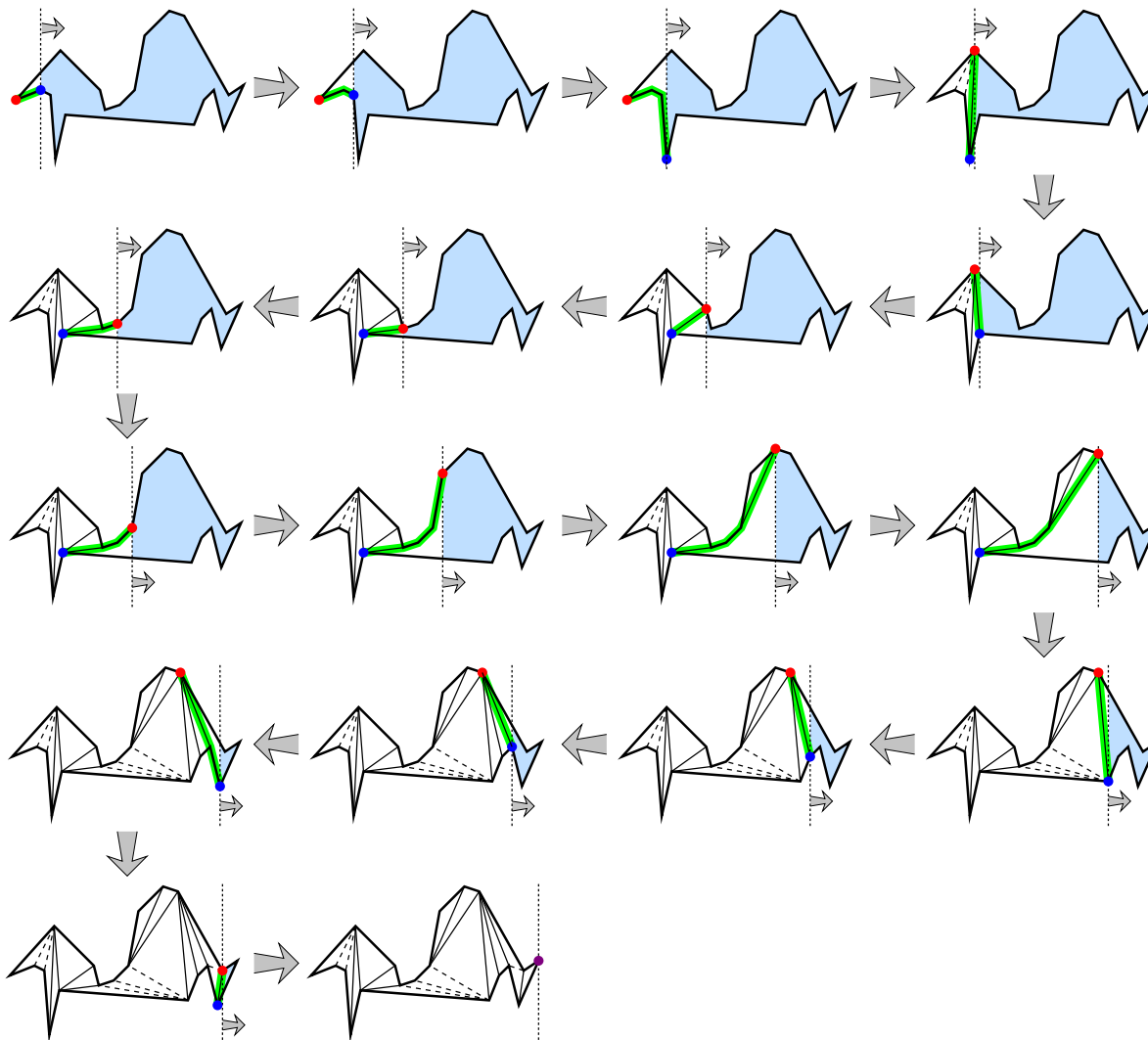
Triangulating a monotone mountain. (Some of the triangles are very thin.)

We can also triangulate general monotone polygons in linear time, but the process is more complicated. A good way to visualize the algorithm is to think of the polygon as a complicated room. Two people named Tom and Bob are walking along the top and bottom walls, both starting at the left end and going to the right. At all times, they have a rubber band stretched between them that can never leave the room.



A rubber band stretched between a vertex on the top and a vertex on the bottom of a monotone polygon.

Now we loop through *all* the vertices of the polygon in order from left to right. Whenever we see a new bottom vertex, Bob moves onto it, and whenever we see a new top vertex Tom moves onto it. After either person moves, we cut the polygon along the rubber band. (In fact, this will only cut the polygon along a single diagonal at any step.) When we’re done, the polygon is decomposed into triangles and *boomerangs*—nonconvex polygons consisting of two straight edges and a concave chain. A boomerang can only be triangulated in one way, by joining the *apex* to every vertex in the concave chain.



Triangulating a monotone polygon by walking a rubber band from left to right.

I don't want to go into too many implementation details, but a few observations should convince you that this algorithm can be implemented to run in  $O(n)$  time. Notice that at all times, the rubber band forms a concave chain. The leftmost edge in the rubber band joins a top vertex to a bottom vertex. If the rubber band has any other vertices, either they are all on top or all on bottom. If all the other vertices are on top, there are just three ways the rubber band can change:

1. The bottom vertex changes, the rubber band straightens out, and we get a new boomerang.
2. The top vertex changes and the rubber band gets a new concave vertex.
3. The top vertex changes, the rubber band loses some vertices, and we get a new boomerang.

Deciding between the first case and the other two requires a simple comparison between  $x$ -coordinates. Deciding between the last two requires a counterclockwise test.



Number Six: What do you want?  
 Number Two: Information!  
 Number Six: Whose side are you on?  
 Number Two: That would be telling. We want information!  
 Number Six: You won't get it!  
 Number Two: By hook or by crook, we will!  
 — Opening sequence of 'The Prisoner' (1967–68)

## H Lower Bounds

### H.1 What *Are* Lower Bounds?

So far in this class we've been developing algorithms and data structures for solving certain problems and analyzing their time and space complexity.

Let  $T_A(X)$  denote the running of algorithm  $A$  given input  $X$ . Then the worst-case running time of  $A$  for inputs of size  $n$  is defined as follows:

$$T_A(n) = \max_{|X|=n} (T_A(X)).$$

The worst-case complexity of a *problem*  $\Pi$  is the worst-case running time of the *fastest* algorithm for solving it:

$$T_{\Pi}(n) = \min_{A \text{ solves } \Pi} (T_A(n)) = \min_{A \text{ solves } \Pi} \left( \max_{|X|=n} (T_A(X)) \right).$$

Now suppose we've shown that the worst-case running time of an algorithm  $A$  is  $O(f(n))$ . Then we immediately have an *upper bound* for the complexity of  $\Pi$ :

$$T_{\Pi}(n) \leq T_A(n) = O(f(n)).$$

The faster our algorithm, the better our upper bound. In other words, when we give a running time for an algorithm, what we're really doing — and what most theoretical computer scientists devote their entire careers doing<sup>1</sup> — is bragging about how *easy* some problem is.

Starting with this lecture, we've turned the tables. Instead of bragging about how easy problems are, now we're arguing that certain problems are *hard* by proving *lower bounds* on their complexity. This is a little harder, because it's no longer enough to examine a single algorithm. To show that  $T_{\Pi}(n) = \Omega(f(n))$ , we have to prove that *every* algorithm that solves  $\Pi$  has a worst-case running time  $\Omega(f(n))$ , or equivalently, that *no* algorithm runs in  $o(f(n))$  time.

---

<sup>1</sup>This sometimes leads to long sequences of results that sound like an obscure version of "Name that Tune":

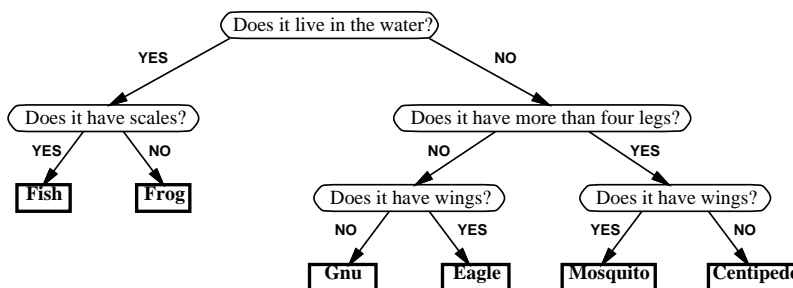
Lennes: "I can triangulate that polygon in  $O(n^2)$  time."  
 Shamos: "I can triangulate that polygon in  $O(n \log n)$  time."  
 Tarjan: "I can triangulate that polygon in  $O(n \log \log n)$  time."  
 Seidel: "I can triangulate that polygon in  $O(n \log^* n)$  time."  
 [Audience gasps.]  
 Chazelle: "I can triangulate that polygon in  $O(n)$  time."  
 [Audience gasps and applauds.]  
 "Triangulate that polygon!"

## H.2 Decision Trees

Unfortunately, there is no formal definition of the phrase ‘all algorithms’!<sup>2</sup> So when we derive lower bounds, we first have to specify, formally, what an algorithm is and how to measure its running time. This specification is called a *model of computation*.

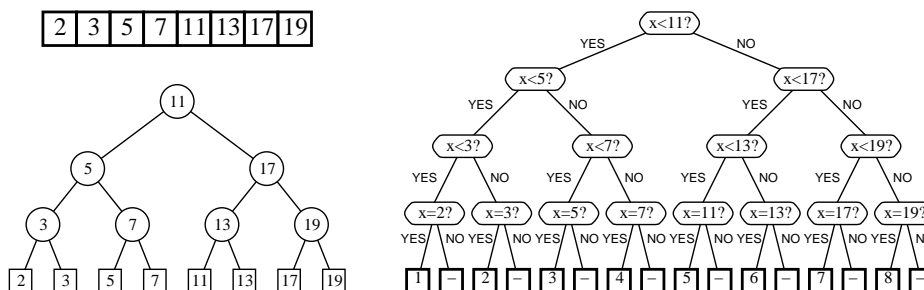
One rather powerful model of computation is *decision trees*. A decision tree is (as the name suggests) a tree. Each internal node in the tree is labeled by a *query*, which is just a question about the input. The edges out of a node correspond to the various answers to the query. Each leaf of the tree is labeled with an *output*. To compute with a decision tree, start at the root and follow a path down to a leaf. At each internal node, the answer to the query tells you which node to visit next. When you reach a leaf, output its label.

For example, the guessing game where one person thinks of an animal and the other person tries to figure it out with a series of yes/no questions can be modeled as a decision tree. Each internal node is labeled with a question and has two edges labeled ‘yes’ and ‘no’. Each leaf is labeled with an animal.



A decision tree to choose one of six animals.

Here’s another simple example, called the *dictionary problem*. Let  $A$  be a fixed array with  $n$  numbers. Suppose want to determine, given a number  $x$ , the position of  $x$  in the array  $A$ , if any. One solution to the dictionary problem is to sort  $A$  (remembering every element’s original position) and then use binary search. The (implicit) binary search tree can be used almost directly as a decision tree. Each internal node the the *search* tree stores a key  $k$ ; the corresponding node in the *decision* tree stores the question ‘Is  $x < k$ ?’. Each leaf in the *search* tree stores some value  $A[i]$ ; the corresponding node in the *decision* tree asks ‘Is  $x = A[i]$ ?’ and has two leaf children, one labeled ‘ $i$ ’ and the other ‘none’.



Left: A binary search tree for the first eight primes.

Right: The corresponding binary decision tree for the dictionary problem (– = ‘none’).

<sup>2</sup>Complexity-theory purists might argue that ‘all algorithms’ is just a synonym for ‘all Turing machines’. (If you want to know what a Turing machine is, take 375.) In my opinion, this is nonsense. Or it might not be nonsense, but it isn’t a particularly *useful* definition. Turing machines are just another model of computation.

We *define* the running time of a decision tree algorithm for a given input to be the number of queries in the path from the root to the leaf. For example, in the ‘Guess the animal’ tree above,  $T(\text{frog}) = 2$ . Thus, the worst-case running time of the algorithm is just the depth of the tree. This definition ignores other kinds of operations that the algorithm might perform that have nothing to do with the queries. (Even the most efficient binary search problem requires more than one machine instruction per comparison!) But the number of decisions is certainly a *lower bound* on the actual running time, which is good enough to prove a lower bound on the complexity of a problem.

Both of the examples describe *binary* decision trees, where every query has only two answers. We may sometimes want to consider decision trees with higher degree. For example, we might use queries like ‘Is  $x$  greater than, equal to, or less than  $y$ ?’ or ‘Are these three points in clockwise order, colinear, or in counterclockwise order?’ A  $k$ -ary decision tree is one where every query has (at most)  $k$  different answers. **From now on, I will only consider  $k$ -ary decision trees where  $k$  is a constant.**

### H.3 Information Theory

Most lower bounds for decision trees are based on the following simple observation: *the answers to the queries must give you enough information to specify any possible output*. If a problem has  $N$  different outputs, then obviously any decision tree must have at least  $N$  leaves. (It’s possible for several leaves to specify the same output.) Thus, if every query has at most  $k$  possible answers, then the depth of the decision tree must be at least  $\lceil \log_k N \rceil = \Omega(\log N)$ .

Let’s apply this to the dictionary problem for a set  $S$  of  $n$  numbers. Since there are  $n + 1$  possible outputs, any decision tree must have at least  $n + 1$  leaves, and thus any decision tree must have depth at least  $\lceil \log_k(n + 1) \rceil = \Omega(\log n)$ . So the complexity of the dictionary problem, in the decision-tree model of computation, is  $\Omega(\log n)$ . This matches the upper bound  $O(\log n)$  that comes from a perfectly-balanced binary search tree. That means that the standard binary search algorithm, which runs in  $O(\log n)$  time, is *optimal*—there is no faster algorithm in this model of computation.

### H.4 But wait a second...

We can solve the membership problem in  $O(1)$  expected time using hashing. Isn’t this inconsistent with the  $\Omega(\log n)$  lower bound?

No, it isn’t. The reason is that hashing involves a query with more than a constant number of outcomes, specifically ‘What is the hash value of  $x$ ?’ In fact, if we don’t restrict the degree of the decision tree, we can get constant running time even without hashing, by using the obviously unreasonable query ‘For which index  $i$  (if any) is  $A[i] = x$ ?’. No, I am *not* cheating — remember that the decision tree model allows us to ask *any* question about the input!

This example illustrates a common theme in proving lower bounds: *choosing the right model of computation is absolutely crucial*. If you choose a model that is too powerful, the problem you’re studying may have a completely trivial algorithm. On the other hand, if you consider more restrictive models, the problem may not be solvable at all, in which case any lower bound will be meaningless! (In this class, we’ll just tell you the right model of computation to use.)

### H.5 Sorting

Now let’s consider the *sorting* problem — Given an array of  $n$  numbers, arrange them in increasing order. Unfortunately, decision trees don’t have any way of describing moving data around, so we have to rephrase the question slightly:

Given a sequence  $\langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find the permutation  $\pi$  such that  $x_{\pi(1)} < x_{\pi(2)} < \dots < x_{\pi(n)}$ .

Now a  $k$ -ary decision-tree lower bound is immediate. Since there are  $n!$  possible permutations  $\pi$ , any decision tree for sorting must have at least  $n!$  leaves, and so must have depth  $\Omega(\log(n!))$ . To simplify the lower bound, we apply *Stirling's approximation*

$$n! = \left(\frac{n}{e}\right)^n \sqrt{2\pi n} \left(1 + \Theta\left(\frac{1}{n}\right)\right) > \left(\frac{n}{e}\right)^n.$$

This gives us the lower bound

$$\lceil \log_k(n!) \rceil > \left\lceil \log_k \left(\frac{n}{e}\right)^n \right\rceil = \lceil n \log_k n - n \log_k e \rceil = \Omega(n \log n).$$

This matches the  $O(n \log n)$  upper bound that we get from mergesort, heapsort, or quicksort, so those algorithms are optimal. The decision-tree complexity of sorting is  $\Theta(n \log n)$ .

Well... we're not quite done. In order to say that those algorithms are optimal, we have to demonstrate that they fit into our model of computation. A few minutes thought will convince you that they can be described as a special type of decision tree called a *comparison tree*, where every query is of the form 'Is  $x_i$  bigger or smaller than  $x_j$ ?' These algorithms treat any two input sequences exactly the same way as long as the same comparisons produce exactly the same results. This is a feature of any comparison tree. In other words, *the actual input values don't matter, only their order*. Comparison trees describe almost all sorting algorithms: bubble sort, selection sort, insertion sort, shell sort, quicksort, heapsort, mergesort, and so forth — but *not* radix sort or bucket sort.

## H.6 Finding the Maximum and Adversaries

Finally let's consider the *maximum* problem: Given an array  $X$  of  $n$  numbers, find its largest entry. Unfortunately, there's no hope of proving a lower bound in this formulation, since there are an infinite number of possible answers, so let's rephrase it slightly.

Given a sequence  $\langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find the index  $m$  such that  $x_m$  is the largest element in the sequence.

We can get an upper bound of  $n - 1$  comparisons in several different ways. The easiest is probably to start at one end of the sequence and do a linear scan, maintaining a current maximum. Intuitively, this seems like the best we can do, but the information-theoretic lower bound is only  $\lceil \log_2 n \rceil$ .

To prove that  $n - 1$  comparisons are actually necessary, we use something called an *adversary argument*. The idea is that an all-powerful malicious adversary *pretends* to choose an input for the algorithm. When the algorithm asks a question about the input, the adversary answers in whatever way will make the algorithm do the most work. If the algorithm does not ask enough queries before terminating, then there will be several different inputs, each consistent with the adversary's answers, the should result in different outputs. In this case, whatever the algorithm outputs, the adversary can 'reveal' an input that is consistent with its answers, but contradicts the algorithm's output, and then claim that that was the input that he was using all along.

For the maximum problem, the adversary originally pretends that  $x_i = i$  for all  $i$ , and answers all comparison queries appropriately. Whenever the adversary reveals that  $x_i < x_j$ , he marks  $x_i$  as an item that the algorithm knows (or should know) is not the maximum element. At most

one element  $x_i$  is marked after each comparison. Note that  $x_n$  is never marked. If the algorithm does less than  $n - 1$  comparisons before it terminates, the adversary must have at least one other unmarked element  $x_k \neq x_n$ . In this case, the adversary can change the value of  $x_k$  from  $k$  to  $n + 1$ , making  $x_k$  the largest element, without being inconsistent with any of the comparisons that the algorithm has performed. In other words, the algorithm cannot tell that the adversary has cheated. However,  $x_n$  is the maximum element in the original input, and  $x_k$  is the largest element in the modified input, so the algorithm cannot possibly give the correct answer for both cases. Thus, in order to be correct, any algorithm must perform at least  $n - 1$  comparisons.

It is very important to notice that the adversary makes *no* assumptions about the order in which the algorithm does its comparisons. The adversary forces *any* algorithm (in this model of computation<sup>3</sup>) to either perform  $n - 1$  comparisons, or to give the wrong answer for at least one input sequence. Notice also that no algorithm can distinguish between a malicious adversary and an honest user who actually chooses an input in advance and answers all queries truthfully.

In the next lecture, we'll see several more complicated adversary arguments.

---

<sup>3</sup>Actually, the  $n - 1$  lower bound for finding the maximum holds in a much powerful model called *algebraic* decision trees, which are binary trees where every query is a comparison between two polynomial functions of the input values, such as 'Is  $x_1^2 - 3x_2x_3 + x_4^{17}$  bigger or smaller than  $5 + x_1x_3^5x_5^2 - 2x_7^{42}$ ?'

*It is possible that the operator could be hit by an asteroid and your \$20 could fall off his cardboard box and land on the ground, and while you were picking it up, \$5 could blow into your hand. You therefore could win \$5 by a simple twist of fate.*

— Penn Jillette, explaining how to win at Three-Card Monte (1999)

## I Adversary Arguments

### I.1 Three-Card Monte

Until Times Square was turned into TimesSquareLand by Mayor Guiliani, you could often find dealers stealing tourists' money using a game called 'Three Card Monte' or 'Spot the Lady'. The dealer has three cards, say the Queen of Hearts and the two and three of clubs. The dealer shuffles the cards face down on a table (usually slowly enough that you can follow the Queen), and then asks the tourist to bet on which card is the Queen. In principle, the tourist's odds of winning are at least one in three.

In practice, however, the tourist never wins, because the dealer cheats. There are actually *four* cards; before he even starts shuffling the cards, the dealer palms the queen or sticks it up his sleeve. No matter what card the tourist bets on, the dealer turns over a black card. If the tourist gives up, the dealer slides the queen under one of the cards and turns it over, showing the tourist 'where the queen was all along'. If the dealer is really good, the tourist won't see the dealer changing the cards and will think maybe the queen *was* there all along and he just wasn't smart enough to figure that out.<sup>1</sup> As long as the dealer doesn't reveal all the black cards at once, the tourist has no way to prove that the dealer cheated!

### I.2 $n$ -Card Monte

Now let's consider a similar game, but with an algorithm acting as the tourist and with bits instead of cards. Suppose we have an array of  $n$  bits and we want to determine if any of them is a 1. Obviously we can figure this out by just looking at every bit, but can we do better? Is there maybe some complicated tricky algorithm to answer the question "Any ones?" without looking at every bit? Well, of course not, but how do we prove it?

The simplest proof technique is called an *adversary* argument. The idea is that an all-powerful malicious adversary (the dealer) *pretends* to choose an input for the algorithm (the tourist). When the algorithm wants looks at a bit (a card), the adversary sets that bit to whatever value will make the algorithm do the most work. If the algorithm does not look at enough bits before terminating, then there will be several different inputs, each consistent with the bits already seen, the should result in different outputs. Whatever the algorithm outputs, the adversary can 'reveal' an input that is has all the examined bits but contradicts the algorithm's output, and then claim that that was the input that he was using all along. Since the only information the algorithm has is the set of bits it examined, the algorithm cannot distinguish between a malicious adversary and an honest user who actually chooses an input in advance and answers all queries truthfully.

For the  $n$ -card monte problem, the adversary originally pretends that the input array is all zeros—whenever the algorithm looks at a bit, it sees a 0. Now suppose the algorithms stops before looking at all three bits. If the algorithm says 'No, there's no 1,' the adversary changes one of the unexamined bits to a 1 and shows the algorithm that it's wrong. If the algorithm says 'Yes, there's a 1,' the adversary reveals the array of zeros and again proves the algorithm wrong. Either way, the algorithm cannot tell that the adversary has cheated.

---

<sup>1</sup>He's right about the second part!



In general, we say that a bit pattern is *evasive* if we have to look at every bit to decide if a string of  $n$  bits contains the pattern. So the pattern 1 is evasive for all  $n$ , and the pattern 01 is evasive if and only if  $n$  is even. It turns out that the *only* patterns that are evasive for *all* values of  $n$  are the one-bit patterns 0 and 1.

## I.4 Evasive Graph Properties

Another class of problems for which adversary arguments give good lower bounds is graph problems where the graph is represented by an adjacency matrix, rather than an adjacency list. Recall that the adjacency matrix of an undirected  $n$ -vertex graph  $G = (V, E)$  is an  $n \times n$  matrix  $A$ , where  $A[i, j] = [(i, j) \in E]$ . We are interested in deciding whether an undirected graph has or does not have a certain *property*. For example, is the input graph connected? Acyclic? Planar? Complete? A tree? We call a graph property *evasive* if we have to look at all  $\binom{n}{2}$  entries in the adjacency matrix to decide whether a graph has that property.

An obvious example of an evasive graph property is *emptiness*: Does the graph have any edges at all? We can show that emptiness is evasive using the following simple adversary strategy. The adversary maintains *two* graphs  $E$  and  $G$ .  $E$  is just the empty graph with  $n$  vertices. Initially  $G$  is the complete graph on  $n$  vertices. Whenever the algorithm asks about an edge, the adversary removes that edge from  $G$  (unless it's already gone) and answers 'no'. If the algorithm terminates without examining every edge, then  $G$  is not empty. Since both  $G$  and  $E$  are consistent with all the adversary's answers, the algorithm must give the wrong answer for one of the two graphs.

## I.5 Connectedness Is Evasive

Now let me give a more complicated example, *connectedness*. Once again, the adversary maintains two graphs,  $Y$  and  $M$  ('yes' and 'maybe').  $Y$  contains all the edges that the algorithm knows are definitely in the input graph.  $M$  contains all the edges that the algorithm thinks *might* be in the input graph, or in other words, all the edges of  $Y$  plus all the unexamined edges. Initially,  $Y$  is empty and  $M$  is complete.

Here's the strategy that adversary follows when the adversary asks whether the input graph contains the edge  $e$ . I'll assume that whenever an algorithm examines an edge, it's in  $M$  but not in  $Y$ ; in other words, algorithms never ask about the same edge more than once.

<pre> HIDECONNECTEDNESS(<math>e</math>):   if <math>M \setminus \{e\}</math> is connected     remove <math>(i, j)</math> from <math>M</math>     return 0   else     add <math>e</math> to <math>Y</math>     return 1 </pre>
---

Notice that  $Y$  and  $M$  are both consistent with the adversary's answers. The adversary strategy maintains a few other simple invariants.

- $Y$  is a subgraph of  $M$ .
- $M$  is connected.
- If  $M$  has a cycle, none of its edges are in  $Y$ . If  $M$  has a cycle, then deleting any edge in that cycle leaves  $M$  connected.



- **$Y$  is acyclic.** Obvious from the previous invariant.
- **If  $Y \neq M$ , then  $Y$  is disconnected.** The only connected acyclic graph is a tree. Suppose  $Y$  is a tree and some edge  $e$  is in  $M$  but not  $Y$ . Then there is a cycle in  $M$  that contains  $e$ , all of whose other edges are in  $Y$ . This is impossible.

Now, if an algorithm terminates before examining all  $\binom{n}{2}$  edges, then there is an edge in  $M$  that is not in  $Y$ . Since the algorithm cannot distinguish between  $M$  and  $Y$ , even though  $M$  is connected and  $Y$  is not, the algorithm cannot possibly give the correct output for both graphs. Thus, in order to be correct, any algorithm must examine every edge—*Connectedness is evasive!*

## I.6 An Evasive Conjecture

A graph property is *nontrivial* if there is at least one graph with the property and at least one graph without the property. (The only trivial properties are ‘Yes’ and ‘No’.) A graph property is *monotone* if it is closed under taking subgraphs — if  $G$  has the property, then any subgraph of  $G$  has the property. For example, emptiness, planarity, acyclicity, and *non*-connectedness are monotone. The properties of being a tree and of having a vertex of degree 3 are not monotone.

**Conjecture 1 (Aanderra, Karp, and Rosenberg).** *Every nontrivial monotone property of  $n$ -vertex graphs is evasive.*

The Aanderra-Karp-Rosenberg conjecture has been proven when  $n = p^e$  for some prime  $p$  and positive integer exponent  $e$ —the proof uses some interesting results from algebraic topology<sup>2</sup>—but it is still open for other values of  $n$ .

Incidentally, this was really the point of the ‘scorpion’ practice homework problem. ‘Scorpionhood’ is an example of a nontrivial graph property that is *not* evasive, since there is an algorithm to identify scorpions by examining only  $O(n)$  of their edges. However, this property is not monotone—a subgraph of a scorpion is not necessarily a scorpion.

## I.7 Finding the Minimum and Maximum

Last time, we saw an adversary argument that finding the largest element of an unsorted set of  $n$  numbers requires at least  $n - 1$  comparisons. Let’s consider the complexity of finding the largest *and* smallest elements. More formally:

Given a sequence  $X = \langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find indices  $i$  and  $j$  such that  $x_i = \min X$  and  $x_j = \max X$ .

How many comparisons do we need to solve this problem? An upper bound of  $2n - 3$  is easy: find the minimum in  $n - 1$  comparisons, and then find the maximum of everything else in  $n - 2$  comparisons. Similarly, a lower bound of  $n - 1$  is easy, since any algorithm that finds the min and the max certainly finds the max.

We can improve both the upper and the lower bound to  $\lceil 3n/2 \rceil - 2$ . The upper bound is established by the following algorithm. Compare all  $\lfloor n/2 \rfloor$  consecutive pairs of elements  $x_{2i-1}$  and  $x_{2i}$ , and put the smaller element into a set  $S$  and the larger element into a set  $L$ . If  $n$  is odd, put

<sup>2</sup>Let  $\Delta$  be a contractible simplicial complex whose automorphism group  $\text{Aut}(\Delta)$  is vertex-transitive, and let  $\Gamma$  be a vertex-transitive subgroup of  $\text{Aut}(\Delta)$ . Suppose there are normal subgroups  $\Gamma_1 \triangleleft \Gamma_2 \triangleleft \Gamma$  such that  $|\Gamma_1| = p^\alpha$  for some prime  $p$  and integer  $\alpha$ ,  $|\Gamma/\Gamma_2| = q^\beta$  for some prime  $q$  and integer  $\beta$ , and  $\Gamma_2/\Gamma_1$  is cyclic. Then  $\Delta$  is a simplex.

No, this will not be on the final exam.

$x_n$  into both  $L$  and  $S$ . Then find the smallest element of  $S$  and the largest element of  $L$ . The total number of comparisons is at most

$$\underbrace{\left\lfloor \frac{n}{2} \right\rfloor}_{\text{build } S \text{ and } L} + \underbrace{\left\lceil \frac{n}{2} \right\rceil - 1}_{\text{compute min } S} + \underbrace{\left\lceil \frac{n}{2} \right\rceil - 1}_{\text{compute max } L} = \left\lceil \frac{3n}{2} \right\rceil - 2.$$

For the lower bound, we use an adversary argument. The adversary marks each element  $+$  if it might be the maximum element, and  $-$  if it might be the minimum element. Initially, the adversary puts both marks  $+$  and  $-$  on every element. If the algorithm compares two double-marked elements, then the adversary declares one smaller, removes the  $+$  mark from the smaller element, and removes the  $-$  mark from the larger one. In every other case, the adversary can answer so that at most one mark needs to be removed. For example, if the algorithm compares a double marked element against one labeled  $-$ , the adversary says the one labeled  $-$  is smaller and removes the  $-$  mark from the other. If the algorithm compares to  $+$ 's, the adversary must unmark one of the two.

Initially, there are  $2n$  marks. At the end, in order to be correct, exactly one item must be marked  $+$  and exactly one other must be marked  $-$ , since the adversary can make any  $+$  the maximum and any  $-$  the minimum. Thus, the algorithm must force the adversary to remove  $2n - 2$  marks. At most  $\lfloor n/2 \rfloor$  comparisons remove two marks; every other comparison removes at most one mark. Thus, the adversary strategy forces any algorithm to perform at least  $2n - 2 - \lfloor n/2 \rfloor = \lceil 3n/2 \rceil - 2$  comparisons.

## I.8 Finding the Median

Finally, let's consider the *median* problem: Given an unsorted array  $X$  of  $n$  numbers, find its  $n/2$ th largest entry. (I'll assume that  $n$  is even to eliminate pesky floors and ceilings.) More formally:

Given a sequence  $\langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find the index  $m$  such that  $x_m$  is the  $n/2$ th largest element in the sequence.

To prove a lower bound for this problem, we can use a combination of information theory and two adversary arguments. We use one adversary argument to prove the following simple lemma:

**Lemma 1.** *Any comparison tree that correctly finds the median element also identifies the elements smaller than the median and larger than the median.*

**Proof:** Suppose we reach a leaf of a decision tree that chooses the median element  $x_m$ , and there is still some element  $x_i$  that isn't known to be larger or smaller than  $x_m$ . In other words, we cannot decide based on the comparisons that we've already performed whether  $x_i < x_m$  or  $x_i > x_m$ . Then in particular no element is known to lie between  $x_i$  and  $x_m$ . This means that there must be an input that is consistent with the comparisons we've performed, in which  $x_i$  and  $x_m$  are adjacent in sorted order. But then we can swap  $x_i$  and  $x_m$ , without changing the result of any comparison, and obtain a different consistent input in which  $x_i$  is the median, not  $x_m$ . Our decision tree gives the wrong answer for this 'swapped' input.  $\square$

This lemma lets us rephrase the median-finding problem yet again.

Given a sequence  $X = \langle x_1, x_2, \dots, x_n \rangle$  of  $n$  distinct numbers, find the indices of its  $n/2 - 1$  largest elements  $L$  and its  $n/2$ th largest element  $x_m$ .

Now suppose a ‘little birdie’ tells us the set  $L$  of elements larger than the median. This information fixes the outcomes of certain comparisons — any item in  $L$  is bigger than any element not in  $L$  — and so we can ‘prune’ those comparisons from the comparison tree. The pruned tree finds the largest element of  $X \setminus L$  (the median of  $X$ ), and thus must have depth at least  $n/2 - 1$ . In fact, an adversary argument similar to last lecture’s implies that *every* leaf in the pruned tree must have depth at least  $n/2 - 1$ , so the pruned tree has at least  $2^{n/2-1}$  leaves.

There are  $\binom{n}{n/2-1} \approx 2^n / \sqrt{n/2}$  possible choices for the set  $L$ . Every leaf in the original comparison tree is also a leaf in exactly one of the  $\binom{n}{n/2-1}$  pruned trees, so the original comparison tree must have at least  $\binom{n}{n/2-1} 2^{n/2-1} \approx 2^{3n/2} / \sqrt{n/2}$  leaves. Thus, any comparison tree that finds the median must have depth at least

$$\left\lceil \frac{n}{2} - 1 + \lg \binom{n}{n/2-1} \right\rceil = \frac{3n}{2} - O(\log n).$$

A more complicated adversary argument (also involving pruning the tree with little birdies) improves this lower bound to  $2n - o(n)$ .

A similar argument implies that at least  $n - k + \left\lceil \lg \binom{n}{k-1} \right\rceil$  comparisons are required to find the  $k$ th largest element in an  $n$ -element set.

## J Reductions

### J.1 Introduction

A common technique for deriving algorithms is *reduction*—instead of solving a problem directly, we use an algorithm for some other related problem as a subroutine or black box.

For example, when we talked about nuts and bolts, I argued that once the nuts and bolts are sorted, we can match each nut to its bolt in linear time. Thus, since we can sort nuts and bolts in  $O(n \log n)$  expected time, then we can also match them in  $O(n \log n)$  expected time:

$$T_{\text{match}}(n) \leq T_{\text{sort}}(n) + O(n) = O(n \log n) + O(n) = O(n \log n).$$

Let's consider (as we did in lecture 3) a decision tree model of computation, where every query is a comparison between a nut and a bolt—too big, too small, or just right? The output to the matching problem is a permutation  $\pi$ , where for all  $i$ , the  $i$ th nut matches the  $\pi(i)$ th bolt. Since there are  $n!$  permutations of  $n$  items, any nut/bolt comparison tree that matches  $n$  nuts and bolts has at least  $n!$  leaves, and thus has depth at least  $\lceil \log_3(n!) \rceil = \Omega(n \log n)$ .

Now the same reduction from matching to sorting can be used to prove a lower bound for *sorting* nuts and bolts, just by reversing the inequality:

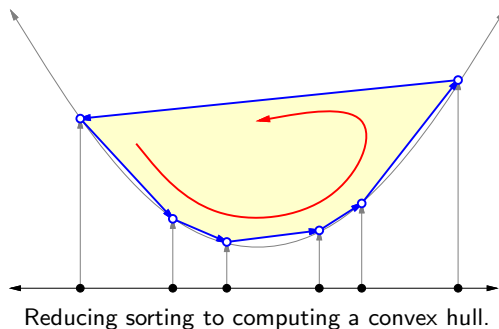
$$T_{\text{sort}}(n) \geq T_{\text{match}}(n) - O(n) = \Omega(n \log n) - O(n) = \Omega(n \log n).$$

Thus, any nut-bolt comparison tree that sorts  $n$  nuts and bolts has depth  $\Omega(n \log n)$ , and our randomized quicksort algorithm is optimal.<sup>1</sup>

### J.2 Sorting to Convex Hulls

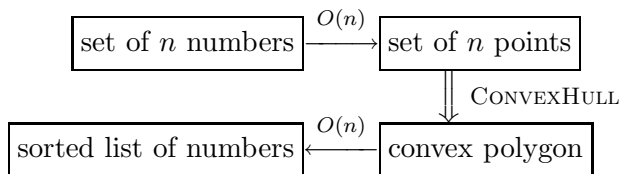
Here's a slightly less trivial example. Suppose we want to prove a lower bound for the problem of computing the convex hull of a set of  $n$  points in the plane. To do this, we demonstrate a reduction from sorting to convex hulls.

To sort a list of  $n$  numbers  $\{a, b, c, \dots\}$ , we first transform it into a set of  $n$  points  $\{(a, a^2), (b, b^2), (c, c^2), \dots\}$ . You can think of the original numbers as a set of points on a horizontal real number line, and the transformation as lifting those point up to the parabola  $y = x^2$ . Then we compute the convex hull of the parabola points. Finally, to get the final sorted list of numbers, we output the first coordinate of every convex vertex, starting from the leftmost vertex and going in counterclockwise order.



<sup>1</sup>We could have proved this lower bound directly. The output to the sorting problem is *two* permutations, so there are  $n!^2$  possible outputs, and we get a lower bound of  $\lceil \log_3(n!^2) \rceil = \Omega(n \log n)$ .

Transforming the numbers into points takes  $O(n)$  time. Since the convex hull is output as a circular doubly-linked list of vertices, reading off the final sorted list of numbers also takes  $O(n)$  time. Thus, given a black-box convex hull algorithm, we can sort in linear extra time. In this case, we say that *there is a linear time reduction from sorting to convex hulls*. We can visualize the reduction as follows:



(I *strongly* encourage you to draw a picture like this whenever you use a reduction argument, at least until you get used to them.) The reduction gives us the following inequality relating the complexities of the two problems:

$$T_{\text{sort}}(n) \leq T_{\text{convex hull}}(n) + O(n)$$

Since we can compute convex hulls in  $O(n \log n)$  time, our reduction implies that we can also sort in  $O(n \log n)$  time. More importantly, by reversing the inequality, we get a lower bound on the complexity of computing convex hulls.

$$T_{\text{convex hull}}(n) \geq T_{\text{sort}}(n) - O(n)$$

Since any binary decision tree requires  $\Omega(n \log n)$  time to sort  $n$  numbers, it follows that any binary decision tree requires  $\Omega(n \log n)$  time to compute the convex hull of  $n$  points.

### J.3 Watch the Model!

This result about the complexity of computing convex hulls is often misquoted as follows:

*Since we need  $\Omega(n \log n)$  comparisons to sort, we also need  $\Omega(n \log n)$  comparisons (between  $x$ -coordinates) to compute convex hulls.*

Although this statement is true, **it's completely trivial**, since it's impossible to compute convex hulls using *any* number of comparisons! In order to compute hulls, we *must* perform counterclockwise tests on triples of points.

The convex hull algorithms we've seen — Graham's scan, Jarvis's march, divide-and-conquer, Chan's shatter — can all be modeled as binary<sup>2</sup> decision trees, where every query is a counterclockwise test on three points. So our binary decision tree lower bound is meaningful, and several of those algorithms are optimal.

This is a subtle but important point about deriving lower bounds using reduction arguments. In order for any lower bound to be meaningful, it must hold in a model in which the problem can be solved! Often the problem we are reducing *from* is much simpler than the problem we are reducing *to*, and thus can be solved in a more restrictive model of computation.

<sup>2</sup>or ternary, if we allow colinear triples of points

## J.4 Element Uniqueness (A Bad Example)

The *element uniqueness* problem asks, given a list of  $n$  numbers  $x_1, x_2, \dots, x_n$ , whether any two of them are equal. There is an obvious and simple algorithm to solve this problem: sort the numbers, and then scan for adjacent duplicates. Since we can sort in  $O(n \log n)$  time, we can solve the element uniqueness problem in  $O(n \log n)$  time.

We also have an  $\Omega(n \log n)$  lower bound for sorting, but our reduction does *not* give us a lower bound for element uniqueness. The reduction goes the wrong way! Inscribe the following on the back of your hand<sup>3</sup>:

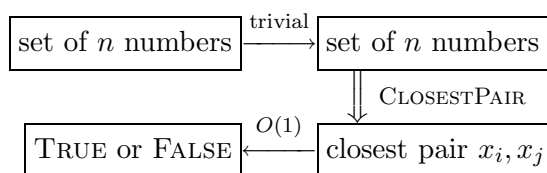
**To prove that problem A is harder than problem B, reduce B to A.**

There isn't (as far as I know) a reduction from sorting to the element uniqueness problem. However, using other techniques (which I won't talk about), it is possible to prove an  $\Omega(n \log n)$  lower bound for the element uniqueness problem. The lower bound applies to so-called *algebraic decision trees*. An algebraic decision tree is a ternary decision tree, where each query asks for the sign of a constant-degree polynomial in the variables  $x_1, x_2, \dots, x_n$ . A comparison tree is an example of an algebraic decision tree, using polynomials of the form  $x_i - x_j$ . The reduction from sorting to element uniqueness implies that any algebraic decision tree requires  $\Omega(n \log n)$  time to sort  $n$  numbers. But since algebraic decision trees are ternary decision trees, we already knew that.

## J.5 Closest Pair

The simplest version of the *closest pair* problem asks, given a list of  $n$  numbers  $x_1, x_2, \dots, x_n$ , to find the closest pair of elements, that is, the elements  $x_i$  and  $x_j$  that minimize  $|x_i - x_j|$ .

There is an obvious reduction from element uniqueness to closest pair, based on the observation that the elements of the input list are distinct if and only if the distance between the closest pair is bigger than zero. This reduction implies that the closest pair problem requires  $\Omega(n \log n)$  time in the algebraic decision tree model.



There are also higher-dimensional closest pair problems; for example, given a set of  $n$  points in the plane, find the two points that closest together. Since the one-dimensional problem is a special case of the 2d problem — just put all  $n$  point son the  $x$ -axis — the  $\Omega(n \log n)$  lower bound applies to the higher-dimensional problems as well.

## J.6 3SUM to Colinearity...

Unfortunately, lower bounds are relatively few and far between. There are thousands of computational problems for which we cannot prove any good lower bounds. We can still learn something useful about the complexity of such a problem by from reductions, namely, that it is harder than some other problem.

Here's an example. The problem 3SUM asks, given a sequence of  $n$  numbers  $x_1, \dots, x_n$ , whether any three of them sum to zero. There is a fairly simple algorithm to solve this problem in  $O(n^2)$

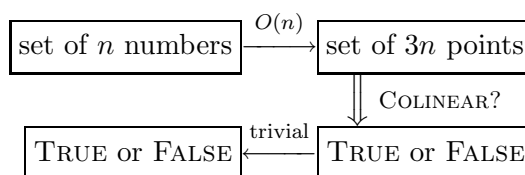
<sup>3</sup>right under all those rules about logarithms, geometric series, and recurrences

time (**hint, hint**). This is widely believed to be the fastest algorithm possible. There is an  $\Omega(n^2)$  lower bound for 3SUM, but only in a fairly weak model of computation.<sup>4</sup>

Now consider a second problem: given a set of  $n$  points in the plane, do any three of them lie on a common non-horizontal line? Again, there is an  $O(n^2)$ -time algorithm, and again, this is believed to be the best possible. The following reduction from 3SUM offers some support for this belief. Suppose we are given an array  $A$  of  $n$  numbers as input to 3SUM. Replace each element  $a \in A$  with three points  $(a, 0)$ ,  $(-a/2, 1)$ , and  $(a, 2)$ . Thus, we replace the  $n$  numbers with  $3n$  points on three horizontal lines  $y = 0$ ,  $y = 1$ , and  $y = 2$ .

If any three points in this set lie on a common non-horizontal line, they consist of one point on each of those three lines, say  $(a, 0)$ ,  $(-b/2, 1)$ , and  $(c, 2)$ . The slope of the common line is equal to both  $-b/2 - a$  and  $c + b/2$ ; since these two expressions are equal, we must have  $a + b + c = 0$ . Similarly, if any three elements  $a, b, c \in A$  sum to zero, then the resulting points  $(a, 0)$ ,  $(-b/2, 1)$ , and  $(c, 2)$  are colinear.

So we have a valid reduction from 3SUM to the colinear-points problem:



$$T_{3\text{SUM}}(n) \leq T_{\text{colinear}}(3n) + O(n) \quad \implies \quad T_{\text{colinear}}(n) \geq T_{3\text{SUM}}(n/3) - O(n).$$

Thus, if we could detect colinear points in  $o(n^2)$  time, we could also solve 3SUM in  $o(n^2)$  time, which seems unlikely. Conversely, if we could prove an  $\Omega(n^2)$  lower bound for 3SUM in a sufficiently powerful model of computation, it would imply an  $\Omega(n^2)$  lower bound for the colinear points problem as well.

The existing  $\Omega(n^2)$  lower bound for 3SUM does *not* imply a lower bound for finding colinear points, because the model of computation is too weak. It is possible to prove an  $\Omega(n^2)$  lower bound directly using an adversary argument, but only in a fairly weak decision-tree model of computation.

Note that in order to prove that the reduction is correct, we have to show that both yes answers and no answers are correct: the numbers sum to zero *if and only if* three points lie on a line. **Even though the reduction itself only goes one way, from the ‘easier’ problem to the ‘harder’ problem, the proof of correctness must go both ways.**

Anka Gajentaan and Mark Overmars<sup>5</sup> defined a whole class of computational geometry problems that are harder than 3SUM; they called these problems 3SUM-*hard*. A sub-quadratic algorithm for any 3SUM-hard problem would imply a subquadratic algorithm for 3SUM. I’ll finish the lecture with two more examples of 3SUM-hard problems.

## J.7 ... to Segment Splitting ...

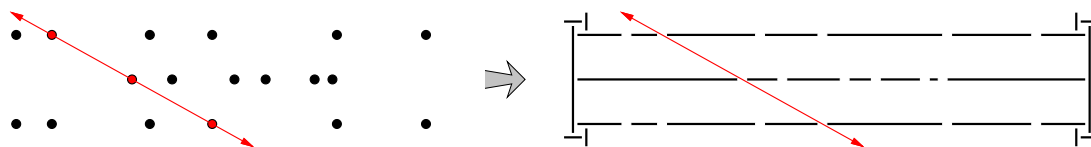
Consider the following *segment splitting problem*: Given a collection of line segments in the plane, is there a line that does not hit any segment and splits the segments into two non-empty subsets?

To show that this problem is 3SUM-hard, we start with the collection of points produced by our last reduction. Replace each point by a ‘hole’ between two horizontal line segments. To make sure

<sup>4</sup>The  $\Omega(n^2)$  lower bound holds in a decision tree model where every query asks for the sign of a linear combination of three of the input numbers. For example, ‘Is  $5x_1 + x_{42} - 17x_5$  positive, negative, or zero?’ See my paper ‘Lower bounds for linear satisfiability problems’ (<http://www.uiuc.edu/~jeffe/pubs/linsat.html>) for the gory(!) details.

<sup>5</sup>A. Gajentaan and M. Overmars, On a class of  $O(n^2)$  problems in computational geometry, *Comput. Geom. Theory Appl.* 5:165–185, 1995. <ftp://ftp.cs.ruu.nl/pub/RUU/CS/techreps/CS-1993/1993-15.ps.gz>

that the only way to split the segments is by passing through three colinear holes, we build two ‘gadgets’, each consisting of five segments, to cap off the left and right ends as shown in the figure below.

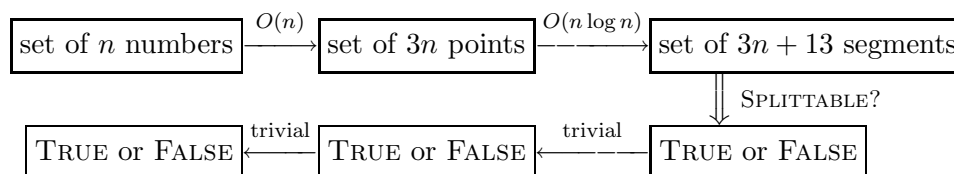


Top:  $3n$  points, three on a non-horizontal line.

Bottom:  $3n + 13$  segments separated by a line through three colinear holes.

This reduction could be performed in linear time if we could make the holes infinitely small, but computers can’t really deal with infinitesimal numbers. On the other hand, if we make the holes too big, we might be able to thread a line through three holes that don’t quite line up. I won’t go into details, but it is possible to compute a working hole size in  $O(n \log n)$  time by first computing the distance between the closest pair of points.

Thus, we have a valid reduction from 3SUM to segment splitting (by way of colinearity):



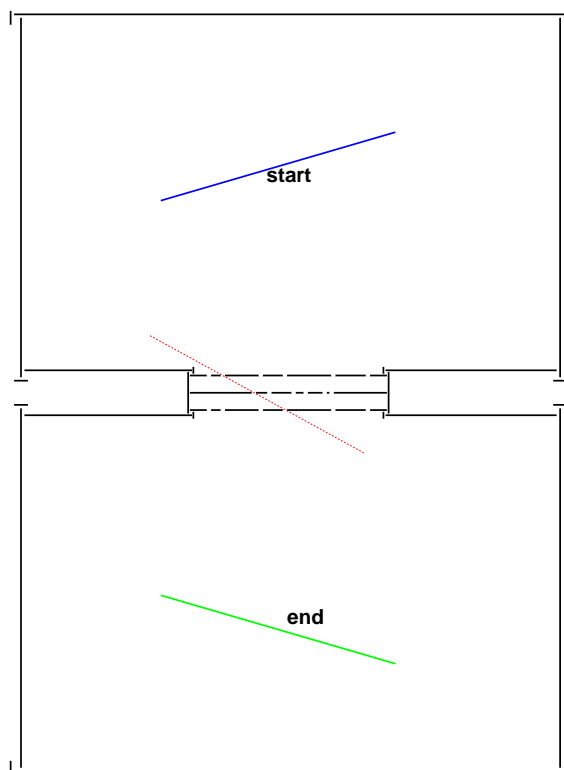
$$T_{3\text{SUM}}(n) \leq T_{\text{split}}(3n + 13) + O(n \log n) \quad \implies \quad T_{\text{split}}(n) \geq T_{3\text{SUM}}\left(\frac{n - 13}{3}\right) - O(n \log n).$$

## J.8 ... to Motion Planning

Finally, suppose we want to know whether a robot can move from one position and location to another. To make things simple, we’ll assume that the robot is just a line segment, and the environment in which the robot moves is also made up of non-intersecting line segments. Given an initial position and orientation and a final position and orientation, is there a sequence of translations and rotations that moves the robot from start to finish?

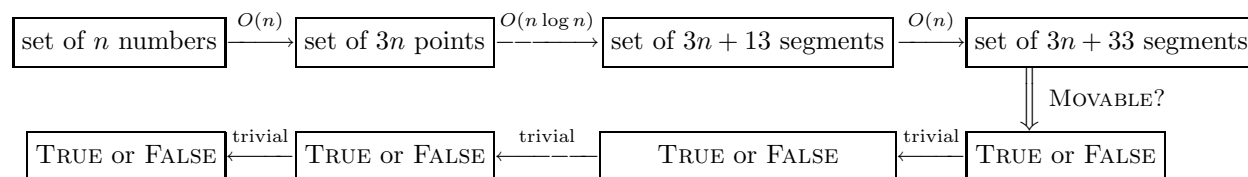
To show that this *motion planning* problem is 3SUM-hard, we do one more reduction, starting from the set of segments output by the previous reduction algorithm. Specifically, we use our earlier set of line segments as a ‘screen’ between two large rooms. The rooms are constructed so that the robot can enter or leave each room only by passing through the screen. We make the robot long enough that the robot can pass from one room to the other if and only if it can pass through three colinear holes in the screen. (If the robot isn’t long enough, it could get between the ‘layers’ of the screen.) See the figure below:





The robot can move from one room to the other if and only if the screen between the rooms has three colinear holes.

Once we have the screen segments, we only need linear time to compute how big the rooms should be, and then  $O(1)$  time to set up the 20 segments that make up the walls. So we have a fast reduction from 3SUM to motion planning (by way of colinearity and segment splitting):



$$T_{3\text{SUM}}(n) \leq T_{\text{motion}}(3n + 33) + O(n \log n) \implies T_{\text{motion}}(n) \geq T_{3\text{SUM}}\left(\frac{n - 33}{3}\right) - O(n \log n).$$

Thus, a sufficiently powerful  $\Omega(n^2)$  lower bound for 3SUM would imply an  $\Omega(n^2)$  lower bound for motion planning as well. The existing  $\Omega(n^2)$  lower bound for 3SUM does *not* imply a lower bound for this problem — the model of computation in which the lower bound holds is too weak to even solve the motion planning problem. In fact, the best lower bound anyone can prove for this motion planning problem is  $\Omega(n \log n)$ , using a (different) reduction from element uniqueness. But the reduction does give us *evidence* that motion planning ‘should’ require quadratic time.

*Math class is tough!*

— Teen Talk Barbie (1992)

*That's why I like it!*

— What she should have said next

*The Manhattan-based Barbie Liberation Organization claims to have performed corrective surgery on 300 Teen Talk Barbies and Talking Duke G.I. Joes—switching their sound chips, repackaging the toys, and returning them to store shelves. Consumers reported their amazement at hearing Barbie bellow, 'Eat lead, Cobra!' or 'Vengeance is mine!', while Joe chirped, 'Will we ever have enough clothes?' and 'Let's plan our dream wedding!'*

— Mark Dery, "Hacking Barbie's Voice Box: Vengeance is Mine!", *New Media* (May 1994)

## 16 NP-Hard Problems (December 3 and 5)

### 16.1 'Efficient' Problems

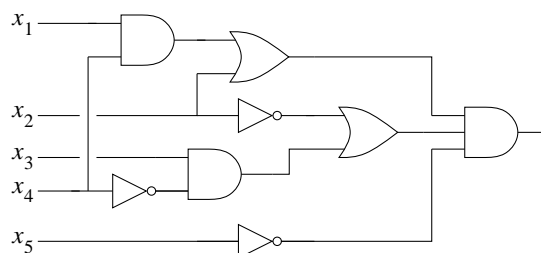
A long time ago<sup>1</sup>, theoretical computer scientists like Steve Cook and Dick Karp decided that a minimum requirement of any efficient algorithm is that it runs in polynomial time:  $O(n^c)$  for some constant  $c$ . People recognized early on that not all problems can be solved this quickly, but we had a hard time figuring out exactly which ones could and which ones couldn't. So Cook, Karp, and others, defined the class of *NP-hard* problems, which most people believe *cannot* be solved in polynomial time, even though nobody can prove a super-polynomial lower bound.

*Circuit satisfiability* is a good example of a problem that we don't know how to solve in polynomial time. In this problem, the input is a *boolean circuit*: a collection of and, or, and not gates connected by wires. We will assume that there are no loops in the circuit (so no delay lines or flip-flops). The input to the *circuit* is a set of  $m$  boolean (true/false) values  $x_1, \dots, x_m$ . The output is a single boolean value. Given specific input values, we can calculate the output in polynomial (actually, *linear*) time using depth-first-search and evaluating the output of each gate in constant time.

The circuit satisfiability problem asks, given a circuit, whether there is an input that makes the circuit output TRUE, or conversely, whether the circuit *always* outputs FALSE. Nobody knows how to solve this problem faster than just trying all  $2^m$  possible inputs to the circuit, but this requires exponential time. On the other hand, nobody has ever proved that this is the best we can do; maybe there's a clever algorithm that nobody has discovered yet!



An AND gate, an OR gate, and a NOT gate.



A boolean circuit. Inputs enter from the left, and the output leaves to the right.

<sup>1</sup>... in a galaxy far far away ...

## 16.2 P, NP, and co-NP

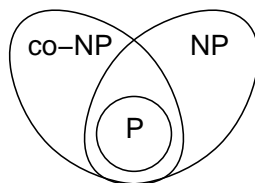
Let me define three classes of problems:

- **P** is the set of yes/no problems<sup>2</sup> that can be solved in polynomial time. Intuitively, P is the set of problems that can be solved quickly.
- **NP** is the set of yes/no problems with the following property: If the answer is yes, then there is a *proof* of this fact that can be checked in polynomial time. Intuitively, NP is the set of problems where we can verify a YES answer quickly if we have the solution in front of us. For example, the circuit satisfiability problem is in NP. If the answer is yes, then any set of  $m$  input values that produces TRUE output is a proof of this fact; we can check the proof by evaluating the circuit in polynomial time.
- **co-NP** is the exact opposite of NP. If the answer to a problem in co-NP is *no*, then there is a proof of this fact that can be checked in polynomial time.

If a problem is in P, then it is also in NP — to verify that the answer is yes in polynomial time, we can just throw away the proof and recompute the answer from scratch. Similarly, any problem in P is also in co-NP.

The (or at least, *a*) central question in theoretical computer science is whether or not  $P=NP$ . Nobody knows. Intuitively, it should be obvious that  $P \neq NP$ ; the homeworks and exams in this class have (I hope) convinced you that problems can be incredibly hard to solve, even when the solutions are obvious once you see them. But nobody can prove it.

Notice that the definition of NP (and co-NP) is not symmetric. Just because we can verify every yes answer quickly, we may not be able to check no answers quickly, and vice versa. For example, as far as we know, there is no short proof that a boolean circuit is *not* satisfiable. But again, we don't have a proof; everyone believes that  $NP \neq co-NP$ , but nobody really knows.



What we *think* the world looks like.

## 16.3 NP-hard, NP-easy, and NP-complete

A problem  $\Pi$  is **NP-hard** if a polynomial-time algorithm for  $\Pi$  would imply a polynomial-time algorithm for *every problem in NP*. In other words:

$\Pi$  is NP-hard  $\iff$  If  $\Pi$  can be solved in polynomial time, then  $P=NP$

Intuitively, this is like saying that if we could solve one particular NP-hard problem quickly, then we could quickly solve *any* problem whose solution is easy to understand, using the solution to that one special problem as a subroutine. NP-hard problems are at least as hard as any problem in NP.

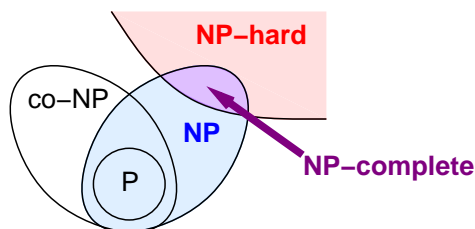
<sup>2</sup>Technically, I should be talking about *languages*, which are just sets of bit strings. The language associated with a yes/no problem is the set of bit strings for which the answer is yes. For example, if the problem is ‘Is the input graph connected?’, then the corresponding language is the set of connected graphs, where each graph is represented as a bit string (for example, its adjacency matrix). P is the set of languages that can be *recognized* in polynomial time by a single-tape Turing machine. Take 375 if you want to know more.

Saying that a problem is NP-hard is like saying ‘If I own a dog, then it can speak fluent English.’ You probably don’t know whether or not I own a dog, but you’re probably pretty sure that I don’t own a *talking* dog. Nobody has a mathematical *proof* that dogs can’t speak English—the fact that no one has ever heard a dog speak English is evidence, as are the hundreds of examinations of dogs that lacked the proper mouth shape and brainpower, but mere evidence is not a proof. Nevertheless, no sane person would believe me if I said I owned a dog that spoke fluent English. So the statement ‘If I own a dog, then it can speak fluent English’ has a natural corollary: No one in their right mind should believe that I own a dog! Likewise, if a problem is NP-hard, no one in their right mind should believe it can be solved in polynomial time.

The following theorem was proved by Steve Cook in 1971. I won’t even sketch the proof (since I’ve been deliberately vague about the definitions). If you want more details, take CS 375 next semester.

**Cook’s Theorem.** *Circuit satisfiability is NP-hard.*

Finally, a problem is **NP-complete** if it is both NP-hard and an element of NP (‘NP-easy’). NP-complete problems are the hardest problems in NP. If anyone finds a polynomial-time algorithm for even one NP-complete problem, then that would imply a polynomial-time algorithm for *every* NP-complete problem. Literally *thousands* of problems have been shown to be NP-complete, so a polynomial-time algorithm for one (*i.e.*, all) of them seems incredibly unlikely.



More of what we *think* the world looks like.

## 16.4 Reductions (again) and SAT

To prove that a problem is NP-hard, we use a reduction argument, *exactly* like we’re trying to prove a lower bound. So we can use a special case of that statement about reductions the you tattooed on the back of your hand last time.

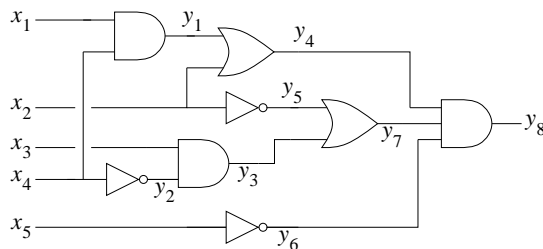
**To prove that problem  $A$  is NP-hard, reduce a known NP-hard problem to  $A$ .**

For example, consider the *formula satisfiability* problem, usually just called *SAT*. The input to SAT is a boolean *formula* like

$$(a \vee b \vee c \vee \bar{d}) \Leftrightarrow ((b \wedge \bar{c}) \vee \overline{(\bar{a} \Rightarrow d)}) \vee (c \neq a \wedge b),$$

and the question is whether it is possible to assign boolean values to the variables  $a, b, c, \dots$  so that the formula evaluates to TRUE.

To show that SAT is NP-hard, we need to give a reduction from a known NP-hard problem. The only problem we know is NP-hard so far is circuit satisfiability, so let’s start there. Given a boolean circuit, we can transform it into a boolean formula by creating new output variables for each gate, and then just writing down the list of gates separated by and. For example, we could transform the example circuit into a formula as follows:

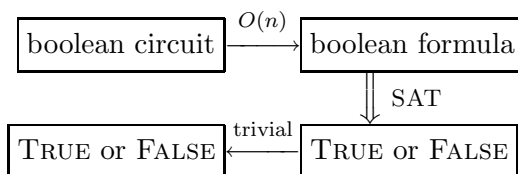


$$(y_1 = x_1 \wedge x_4) \wedge (y_2 = \overline{x_4}) \wedge (y_3 = x_3 \wedge y_2) \wedge (y_4 = y_1 \vee x_2) \wedge \\ (y_5 = \overline{x_2}) \wedge (y_6 = \overline{x_5}) \wedge (y_7 = y_3 \vee y_5) \wedge (y_8 = y_4 \wedge y_7 \wedge y_6) \wedge y_8$$

A boolean circuit with gate variables added, and an equivalent boolean formula.

Now the original circuit is satisfiable if and only if the resulting formula is satisfiable. Given a satisfying input to the circuit, we can get a satisfying assignment for the formula by computing the output of every gate. Given a satisfying assignment for the formula, we can get a satisfying input to the circuit by just ignoring the gate variables  $y_i$ .

We can transform any boolean circuit into a formula in linear time using depth-first search, and the size of the resulting formula is only a constant factor larger than the size of the circuit. Thus, we have a polynomial-time reduction from circuit satisfiability to SAT:



$$T_{\text{CSAT}}(n) \leq O(n) + T_{\text{SAT}}(O(n)) \quad \implies \quad T_{\text{SAT}}(n) \geq T_{\text{CSAT}}(\Omega(n)) - O(n)$$

The reduction implies that if we had a polynomial-time algorithm for SAT, then we'd have a polynomial-time algorithm for circuit satisfiability, which would imply that  $P=NP$ . So SAT is NP-hard.

To prove that a boolean formula is satisfiable, we only have to specify an assignment to the variables that makes the formula true. We can check the proof in linear time just by reading the formula from left to right, evaluating as we go. So SAT is also in NP, and thus is actually NP-complete.

## 16.5 3SAT

A special case of SAT that is incredibly useful in proving NP-hardness results is *3SAT* (or as [CLRS] insists on calling it, *3-CNF-SAT*).

A boolean formula is in *conjunctive normal form* (CNF) if it is a conjunction (AND) of several *clauses*, each of which is the disjunction (OR) of several *literals*, each of which is either a variable or its negation. For example:

$$\overbrace{(a \vee b \vee c \vee d)}^{\text{clause}} \wedge (b \vee \bar{c} \vee \bar{d}) \wedge (\bar{a} \vee c \vee d) \wedge (a \vee \bar{b})$$

A *3CNF* formula is a CNF formula with exactly three literals per clause; the previous example is not a 3CNF formula, since its first clause has four literals and its last clause has only two. 3SAT

is just SAT restricted to 3CNF formulas — given a 3CNF formula, is there an assignment to the variables that makes the formula evaluate to true?

We could prove that 3SAT is NP-hard by a reduction from the more general SAT problem, but it's easier just to start over from scratch, with a boolean circuit. We perform the reduction in several stages.

1. *Make sure every AND and OR gate has only two inputs.* If any gate has  $k > 2$  inputs, replace it with a binary tree of  $k - 1$  two-input gates.
2. *Write down the circuit as a formula, with one clause per gate.* This is just the previous reduction.
3. *Change every gate clause into a CNF formula.* There are only three types of clauses, one for each type of gate:

$$\begin{aligned} a = b \wedge c &\longmapsto (a \vee \bar{b} \vee \bar{c}) \wedge (\bar{a} \vee b) \wedge (\bar{a} \vee c) \\ a = b \vee c &\longmapsto (\bar{a} \vee b \vee c) \wedge (a \vee \bar{b}) \wedge (a \vee \bar{c}) \\ a = \bar{b} &\longmapsto (a \vee b) \wedge (\bar{a} \vee \bar{b}) \end{aligned}$$

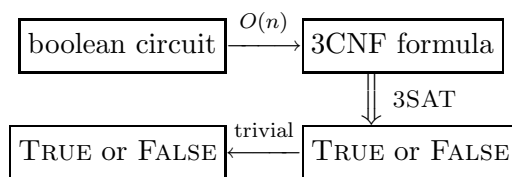
4. *Make sure every clause has exactly three literals.* Introduce new variables into each one- and two-literal clause, and expand it into two clauses as follows:

$$\begin{aligned} a &\longmapsto (a \vee x \vee y) \wedge (a \vee \bar{x} \vee y) \wedge (a \vee x \vee \bar{y}) \wedge (a \vee \bar{x} \vee \bar{y}) \\ a \vee b &\longmapsto (a \vee b \vee x) \wedge (a \vee b \vee \bar{x}) \end{aligned}$$

For example, if we start with the same example circuit we used earlier, we obtain the following 3CNF formula. Although this may look a lot more ugly and complicated than the original circuit at first glance, it's actually only a constant factor larger. Even if the formula were larger than the circuit by a *polynomial*, like  $n^{373}$ , we would have a valid reduction.

$$\begin{aligned} &(y_1 \vee \bar{x}_1 \vee \bar{x}_4) \wedge (\bar{y}_1 \vee x_1 \vee z_1) \wedge (\bar{y}_1 \vee x_1 \vee \bar{z}_1) \wedge (\bar{y}_1 \vee x_4 \vee z_2) \wedge (\bar{y}_1 \vee x_4 \vee \bar{z}_2) \\ &\wedge (y_2 \vee x_4 \vee z_3) \wedge (y_2 \vee x_4 \vee \bar{z}_3) \wedge (\bar{y}_2 \vee \bar{x}_4 \vee z_4) \wedge (\bar{y}_2 \vee \bar{x}_4 \vee \bar{z}_4) \\ &\wedge (y_3 \vee \bar{x}_3 \vee \bar{y}_2) \wedge (\bar{y}_3 \vee x_3 \vee z_5) \wedge (\bar{y}_3 \vee x_3 \vee \bar{z}_5) \wedge (\bar{y}_3 \vee y_2 \vee z_6) \wedge (\bar{y}_3 \vee y_2 \vee \bar{z}_6) \\ &\wedge (\bar{y}_4 \vee y_1 \vee x_2) \wedge (y_4 \vee \bar{x}_2 \vee z_7) \wedge (y_4 \vee \bar{x}_2 \vee \bar{z}_7) \wedge (y_4 \vee \bar{y}_1 \vee z_8) \wedge (y_4 \vee \bar{y}_1 \vee \bar{z}_8) \\ &\wedge (y_5 \vee x_2 \vee z_9) \wedge (y_5 \vee x_2 \vee \bar{z}_9) \wedge (\bar{y}_5 \vee \bar{x}_2 \vee z_{10}) \wedge (\bar{y}_5 \vee \bar{x}_2 \vee \bar{z}_{10}) \\ &\wedge (y_6 \vee x_5 \vee z_{11}) \wedge (y_6 \vee x_5 \vee \bar{z}_{11}) \wedge (\bar{y}_6 \vee \bar{x}_5 \vee z_{12}) \wedge (\bar{y}_6 \vee \bar{x}_5 \vee \bar{z}_{12}) \\ &\wedge (\bar{y}_7 \vee y_3 \vee y_5) \wedge (y_7 \vee \bar{y}_3 \vee z_{13}) \wedge (y_7 \vee \bar{y}_3 \vee \bar{z}_{13}) \wedge (y_7 \vee \bar{y}_5 \vee z_{14}) \wedge (y_7 \vee \bar{y}_5 \vee \bar{z}_{14}) \\ &\wedge (y_8 \vee \bar{y}_4 \vee \bar{y}_7) \wedge (\bar{y}_8 \vee y_4 \vee z_{15}) \wedge (\bar{y}_8 \vee y_4 \vee \bar{z}_{15}) \wedge (\bar{y}_8 \vee y_7 \vee z_{16}) \wedge (\bar{y}_8 \vee y_7 \vee \bar{z}_{16}) \\ &\wedge (y_9 \vee \bar{y}_8 \vee \bar{y}_6) \wedge (\bar{y}_9 \vee y_8 \vee z_{17}) \wedge (\bar{y}_9 \vee y_8 \vee \bar{z}_{17}) \wedge (\bar{y}_9 \vee y_6 \vee z_{18}) \wedge (\bar{y}_9 \vee y_6 \vee \bar{z}_{18}) \\ &\wedge (y_9 \vee z_{19} \vee z_{20}) \wedge (y_9 \vee \bar{z}_{19} \vee z_{20}) \wedge (y_9 \vee z_{19} \vee \bar{z}_{20}) \wedge (y_9 \vee \bar{z}_{19} \vee \bar{z}_{20}) \end{aligned}$$

At the end of this process, we've transformed the circuit into an equivalent 3CNF formula. The formula is satisfiable if and only if the original circuit is satisfiable. As with the more general SAT problem, the formula is only a constant factor larger than then any reasonable description of the original circuit, and the reduction can be carried out in polynomial time. Thus, we have a polynomial-time reduction from circuit satisfiability to 3SAT:



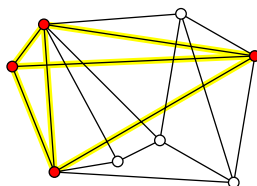
$$T_{\text{CSAT}}(n) \leq O(n) + T_{\text{3SAT}}(O(n)) \implies T_{\text{3SAT}}(n) \geq T_{\text{CSAT}}(\Omega(n)) - O(n)$$

So 3SAT is NP-hard.

Finally, since 3SAT is a special case of SAT, it is also in NP, so 3SAT is NP-complete.

## 16.6 Maximum Clique Size (from 3SAT)

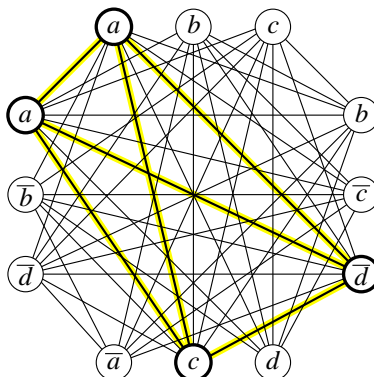
The last problem I'll consider in this lecture is a graph problem. A *clique* is another name for a complete graph. The *maximum clique size* problem, or simply MAXCLIQUE, is to compute, given a graph, the number of nodes in its largest complete subgraph.



A graph with maximum clique size 4.

I'll prove that MAXCLIQUE is NP-hard (but not NP-complete, since it isn't a yes/no problem) using a reduction from 3SAT. I'll describe a reduction algorithm that transforms a 3CNF formula into a graph that has a clique of a certain size if and only if the formula is satisfiable.

The graph has one node for each instance of each literal in the formula. Two nodes are connected by an edge if (1) they correspond to literals in different clauses and (2) those literals do not contradict each other. In particular, all the nodes that come from the same literal (in different clauses) are joined by edges. For example, the formula  $(a \vee b \vee c) \wedge (b \vee \bar{c} \vee \bar{d}) \wedge (\bar{a} \vee c \vee d) \wedge (a \vee \bar{b} \vee \bar{d})$  is transformed into the following graph. (Look for the edges that *aren't* in the graph.)

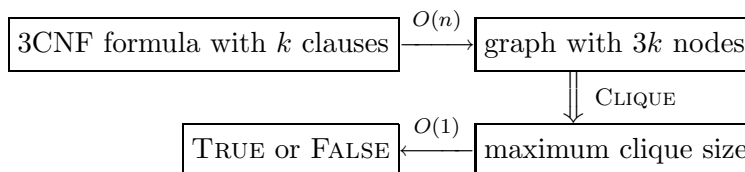


A graph derived from a 3CNF formula, and a clique of size 4.

Now suppose the original formula had  $k$  clauses. Then I claim that the formula is satisfiable if and only if the graph has a clique of size  $k$ .

1.  **$k$ -clique  $\implies$  satisfying assignment:** If the graph has a clique of  $k$  vertices, then each vertex must come from a different clause. To get the satisfying assignment, we declare that each literal in the clique is true. Since we only connect non-contradictory literals with edges, this declaration assigns a consistent value to several of the variables. There may be variables that have no literal in the clique; we can set these to any value we like.
2. **satisfying assignment  $\implies k$ -clique:** If we have a satisfying assignment, then we can choose one literal in each clause that is true. Those literals form a clique in the graph.

Thus, the reduction is correct. Since the reduction from 3CNF formula to graph can be done in polynomial time, so MAXCLIQUE is NP-hard. Here's a diagram of the reduction:

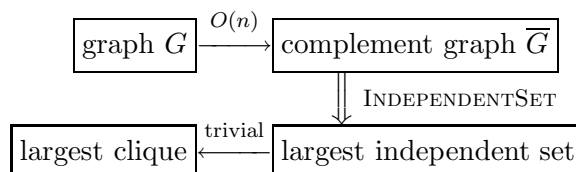


$$T_{3\text{SAT}}(n) \leq O(n) + T_{\text{MAXCLIQUE}}(O(n)) \implies T_{\text{MAXCLIQUE}}(n) \geq T_{3\text{SAT}}(\Omega(n)) - O(n)$$

## 16.7 Independent Set (from Clique)

An *independent set* is a collection of vertices in a graph with no edges between them. The INDEPENDENTSET problem is to find the largest independent set in a given graph.

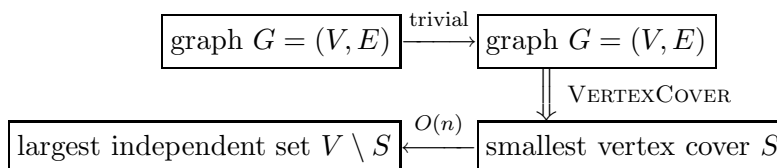
There is an easy proof that INDEPENDENTSET is NP-hard, using a reduction from CLIQUE. Any graph  $G$  has a *complement*  $\overline{G}$  with the same vertices, but with exactly the opposite set of edges— $(u, v)$  is an edge in  $\overline{G}$  if and only if it is *not* an edge in  $G$ . A set of vertices forms a clique in  $G$  if and only if the same vertices are an independent set in  $\overline{G}$ . Thus, we can compute the largest clique in a graph simply by computing the largest independent set in the complement of the graph.



## 16.8 Vertex Cover (from Independent Set)

A *vertex cover* of a graph is a set of vertices that touches every edge in the graph. The VERTEXCOVER problem is to find the smallest vertex cover in a given graph.

Again, the proof of NP-hardness is simple, and relies on just one fact: If  $I$  is an independent set in a graph  $G = (V, E)$ , then  $V \setminus I$  is a vertex cover. Thus, to find the *largest* independent set, we just need to find the vertices that aren't in the *smallest* vertex cover of the same graph.





### 16.9 Graph Coloring (from 3SAT)

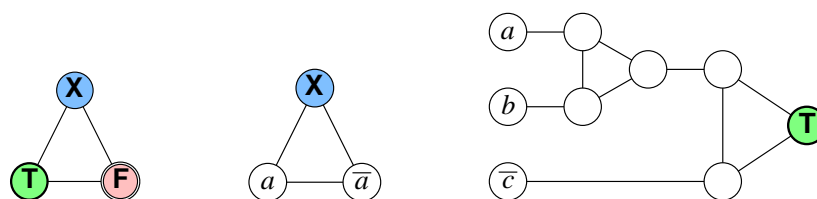
A  $c$ -coloring of a graph is a map  $C : V \rightarrow \{1, 2, \dots, c\}$  that assigns one of  $c$  ‘colors’ to each vertex, so that every edge has two different colors at its endpoints. The graph coloring problem is to find the smallest possible number of colors in a legal coloring. To show that this problem is NP-hard, it’s enough to consider the special case 3COLORABLE: Given a graph, does it have a 3-coloring?

To prove that 3COLORABLE is NP-hard, we use a reduction from 3SAT. Given a 3CNF formula, we produce a graph as follows. The graph consists of a *truth gadget*, one *variable gadget* for each variable in the formula, and one *clause gadget* for each clause in the formula.

The truth gadget is just a triangle with three vertices  $T$ ,  $F$ , and  $X$ , which intuitively stand for TRUE, FALSE, and OTHER. Since these vertices are all connected, they must have different colors in any 3-coloring. For the sake of convenience, we will name those colors TRUE, FALSE, and OTHER. Thus, when we say that a node is colored TRUE, all we mean is that it must be colored the same as the node  $T$ .

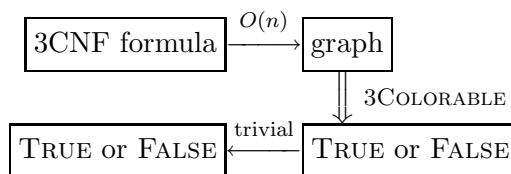
The variable gadget for a variable  $a$  is also a triangle joining two new nodes labeled  $a$  and  $\bar{a}$  to node  $X$  in the truth gadget. Node  $a$  must be colored either TRUE or FALSE, and so node  $\bar{a}$  must be colored either FALSE or TRUE, respectively.

Finally, each clause gadget joins three literal nodes to node  $T$  in the truth gadget using five new unlabeled nodes and ten edges; see the figure below. If all three literal nodes in the clause gadget are colored FALSE, then the rightmost vertex in the gadget cannot have one of the three colors. Since the variable gadgets force each literal node to be colored either TRUE or FALSE, in any valid 3-coloring, at least one of the three literal nodes is colored TRUE. I need to emphasize here that the final graph contains only *one* node  $T$ , only *one* node  $F$ , only *one* node  $\bar{a}$  for each variable  $a$ , and so on.

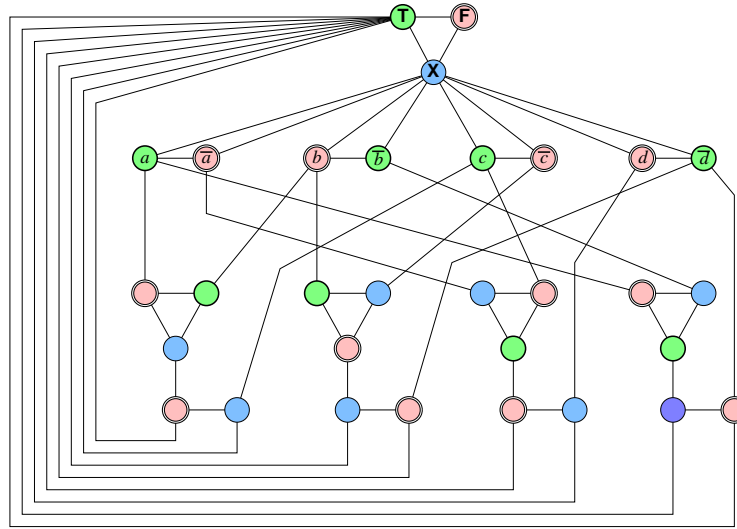


Gadgets for the reduction from 3SAT to 3-Colorability:  
The truth gadget, a variable gadget for  $a$ , and a clause gadget for  $(a \vee b \vee \bar{c})$ .

The proof of correctness is just brute force. If the graph is 3-colorable, then we can extract a satisfying assignment from any 3-coloring—at least one of the three literal nodes in every clause gadget is colored TRUE. Conversely, if the formula is satisfiable, then we can color the graph according to any satisfying assignment.



For example, the formula  $(a \vee b \vee c) \wedge (b \vee \bar{c} \vee \bar{d}) \wedge (\bar{a} \vee c \vee d) \wedge (a \vee \bar{b} \vee \bar{d})$  that I used to illustrate the MAXCLIQUE reduction would be transformed into the following graph. The 3-coloring is one of several that correspond to the satisfying assignment  $a = c = \text{TRUE}$ ,  $b = d = \text{FALSE}$ .



A 3-colorable graph derived from a satisfiable 3CNF formula.

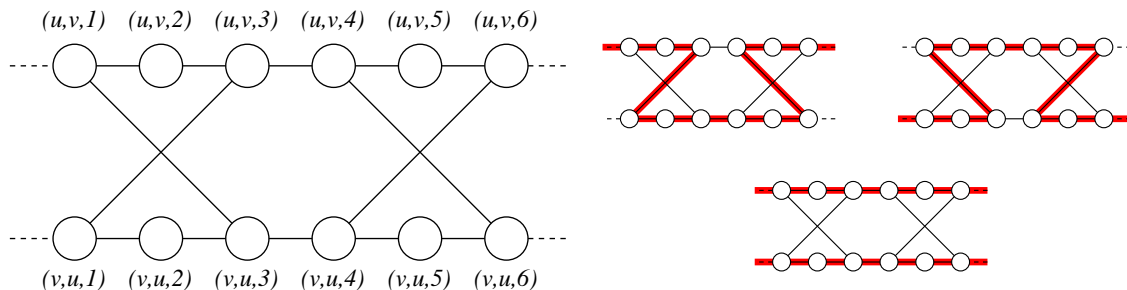
We can easily verify that a graph has been correctly 3-colored in linear time: just compare the endpoints of every edge. Thus, 3COLORING is in NP, and therefore NP-complete. Moreover, since 3COLORING is a special case of the more general graph coloring problem—What is the minimum number of colors?—the more problem is also NP-hard, but *not* NP-complete, because it's not a yes/no problem.

### 16.10 Hamiltonian Cycle (from Vertex Cover)

A *Hamiltonian cycle* in a graph is a cycle that visits every vertex exactly once. This is very different from an *Eulerian cycle*, which is actually a closed *walk* that traverses every *edge* exactly once. Eulerian cycles are easy to find and construct in linear time using a variant of depth-first search. Finding Hamiltonian cycles, on the other hand, is NP-hard.

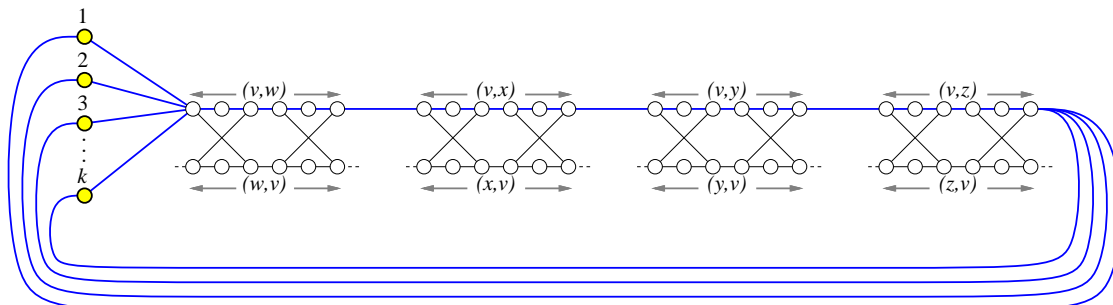
To prove this, we use a reduction from the vertex cover problem. Given a graph  $G$  and an integer  $k$ , we need to transform it into another graph  $G'$ , such that  $G'$  has a Hamiltonian cycle if and only if  $G$  has a vertex cover of size  $k$ . As usual, our transformation consists of putting together several gadgets.

- For each edge  $(u, v)$  in  $G$ , we have an *edge gadget* in  $G'$  consisting of twelve vertices and fourteen edges, as shown below. The four corner vertices  $(u, v, 1)$ ,  $(u, v, 6)$ ,  $(v, u, 1)$ , and  $(v, u, 6)$  each have an edge leaving the gadget. A Hamiltonian cycle can only pass through an edge gadget in one of three ways. Eventually, these will correspond to one or both of the vertices  $u$  and  $v$  being in the vertex cover.



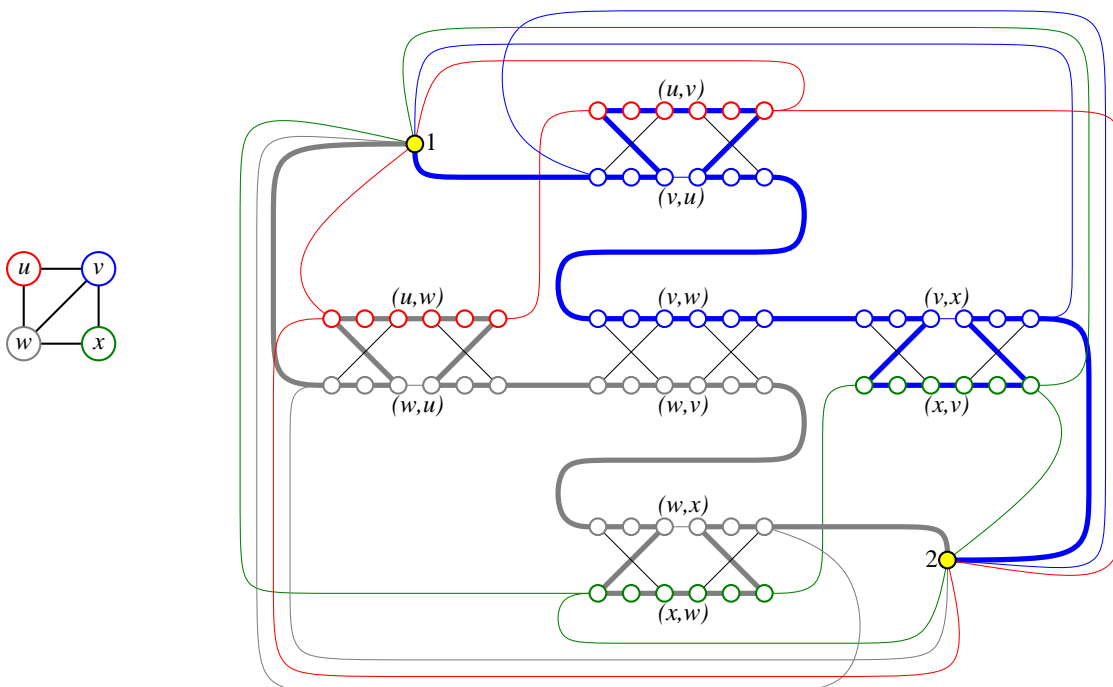
An edge gadget for  $(u, v)$  and the only possible Hamiltonian paths through it.

- $G'$  also contains  $k$  cover vertices, simply numbered 1 through  $k$ .
- Finally, for each vertex  $u$  in  $G$ , we string together all the edge gadgets for edges  $(u, v)$  into a single *vertex chain*, and then connect the ends of the chain to all the cover vertices. Specifically, suppose  $u$  has  $d$  neighbors  $v_1, v_2, \dots, v_d$ . Then  $G'$  has  $d - 1$  edges between  $(u, v_i, 6)$  and  $(u, v_{i+1}, 1)$ , plus  $k$  edges between the cover vertices and  $(u, v_1, 1)$ , and finally  $k$  edges between the cover vertices and  $(u, v_d, 6)$ .



The vertex chain for  $v$ : all edge gadgets involving  $v$  are strung together and joined to the  $k$  cover vertices.

It's not hard to prove that if  $\{v_1, v_2, \dots, v_k\}$  is a vertex cover of  $G$ , then  $G'$  has a Hamiltonian cycle—start at cover vertex 1, through traverse the vertex chain for  $v_1$ , then visit cover vertex 2, then traverse the vertex chain for  $v_2$ , and so forth, eventually returning to cover vertex 1. Conversely, any Hamiltonian cycle in  $G'$  alternates between cover vertices and vertex chains, and the vertex chains correspond to the  $k$  vertices in a vertex cover of  $G$ . (This is a little harder to prove.) Thus,  $G$  has a vertex cover of size  $k$  if and only if  $G'$  has a Hamiltonian cycle.



The original graph  $G$  with vertex cover  $\{v, w\}$ , and the transformed graph  $G'$  with a corresponding Hamiltonian cycle. Vertex chains are colored to match their corresponding vertices.

The transformation from  $G$  to  $G'$  takes at most  $O(n^2)$  time, so the Hamiltonian cycle problem is NP-hard. Moreover, since we can easily verify a Hamiltonian cycle in linear time, the Hamiltonian cycle problem is in NP, and therefore NP-complete.

A closely related problem to Hamiltonian cycles is the famous *traveling salesman problem*—Given a *weighted* graph  $G$ , find the shortest cycle that visits every vertex. Finding the shortest cycle is obviously harder than determining if a cycle exists at all, so the traveling salesman problem is also NP-hard.

### 16.11 Minesweeper (from Circuit SAT)

In 1999, Richard Kaye proved that the solitaire game Minesweeper is NP-complete, using a reduction from the original circuit satisfiability problem.<sup>3</sup> The reduction involves setting up gadgets for every possible feature of a boolean circuit: wires, AND gates, OR gates, NOT gates, wire crossings, and so forth. For all the gory details, see <http://www.mat.bham.ac.uk/R.W.Kaye/minesw/minesw.pdf>!

### 16.12 Other Useful NP-hard Problems

Literally thousands of different problems have been proved to be NP-hard. I want to close this note by listing a few NP-hard problems that are useful in deriving reductions. I won't describe the NP-hardness for these problems, but you can find them in Garey and Johnson's *Angry Black Book of NP-Completeness*.<sup>4</sup>

- **PLANARCIRCUITSAT**: Given a boolean circuit that can be embedded in the plane so that no two wires cross, is there an input that makes the circuit output TRUE? This can be proved NP-hard by reduction from the general circuit satisfiability problem, by replacing each crossing with a small series of gates. (This is an easy exercise.)
- **NOTALLEQUAL3SAT**: Given a 3CNF formula, is there an assignment of values to the variables so that every clause contains at least one true literal *and* at least one false literal? This can be proved NP-hard by reduction from the usual 3SAT.
- **PLANAR3SAT**: Given a 3CNF boolean formula, consider a bipartite graph whose vertices are the clauses and variables, where an edge indicates that a variable (or its negation) appears in a clause. If this graph is planar, the 3CNF formula is also called planar. The **PLANAR3SAT** problem asks, given a planar 3CNF formula, whether it has a satisfying assignment. This can be proven NP-hard by reduction from **PLANARSAT**.
- **PLANARNOTALLEQUAL3SAT**: You get the idea.
- **EXACT3DIMENSIONALMATCHING** or **X3M**: Given a set  $S$  and a collection of three-element subsets of  $S$ , called *triples*, is there a subcollection of disjoint triples that exactly cover  $S$ ? This can be proved NP-hard by a reduction from 3SAT.
- **PARTITION**: Given a set  $S$  of  $n$  integers, are there subsets  $A$  and  $B$  such that  $A \cup B = S$ ,  $A \cap B = \emptyset$ , and

$$\sum_{a \in A} a = \sum_{b \in B} b?$$

This can be proved NP-hard by a reduction from...

<sup>3</sup>Minesweeper is NP-complete. *Mathematical Intelligencer* 22(2):9–15, 2000.

<sup>4</sup>Michael Garey and David Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman and Co., 1979.

- SUBSETSUM: Given a set  $S$  of  $n$  integers and an integer  $T$ , is there a subset  $A \subseteq S$  such that

$$\sum_{a \in A} a = T?$$

This is a generalization of PARTITION, where  $T = (\sum S)/2$ . SUBSETSUM can be proved NP-hard by a (nontrivial!) reduction from either 3SAT or X3M.

- 3PARTITION: Given a set  $S$  with  $3n$  elements, can it be partitioned into  $n$  disjoint subsets, each with 3 elements, such that every subset has the same sum. Note that this is *very* different from the PARTITION problem; I didn't make up the names. This can be proved NP-hard by reduction from X3M. The similar problem of dividing a set of  $2n$  into  $n$  equal-weight two-element sets can be solved in  $O(n \log n)$  time.
- SETCOVER: Given a collection of sets  $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$ , find the smallest subcollection of  $S_i$ 's that contains all the elements of  $\bigcup_i S_i$ . This is a generalization of both VERTEXCOVER and X3M.
- HITTINGSET: Given a collection of sets  $\mathcal{S} = \{S_1, S_2, \dots, S_m\}$ , find the minimum number of elements of  $\bigcup_i S_i$  that hit every set in  $\mathcal{S}$ . This is also a generalization of VERTEXCOVER.
- LONGESTPATH: Given a non-negatively weighted graph  $G$  and two vertices  $u$  and  $v$ , what is the longest simple path from  $u$  to  $v$  in the graph? A path is *simple* if it visits each vertex at most once. This is a generalization of the HAMILTONIANPATH problem. Of course, the corresponding *shortest* path problem is in P.
- STEINERTREE: Given a weighted, undirected graph  $G$  with some of the vertices marked, what is the minimum-weight subtree of  $G$  that contains every marked vertex? If *every* vertex is marked, the minimum Steiner tree is just the minimum spanning tree; if exactly two vertices are marked, the minimum Steiner tree is just the shortest path between them. This can be proved NP-hard by reduction to HAMILTONIANPATH.
- TETRIS: Given a Tetris board and a finite sequence of future pieces, can you survive? This was recently proved NP-hard by reduction from 3PARTITION.