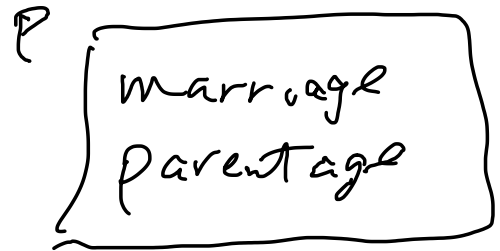
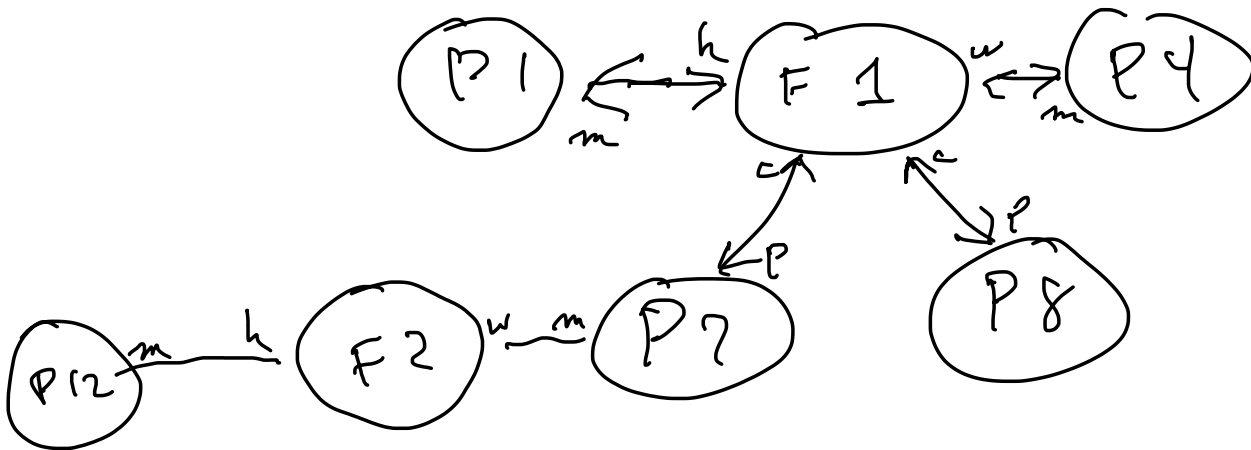
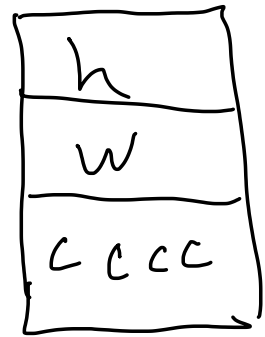


Data Structures #5

F:



Bipartite graph:

families only point to persons.
persons only point to families.

edges always connect a person + a family.

Pattern matching, — Text search problem.

find a match for a pattern (string) p in a text t

$t = \overset{0}{h} \overset{1}{e} \overset{2}{l} \overset{3}{l} \overset{4}{o} \overset{5}{ } \overset{6}{w} \overset{7}{o} \overset{8}{r} \overset{9}{l} \quad \text{result} = 7$

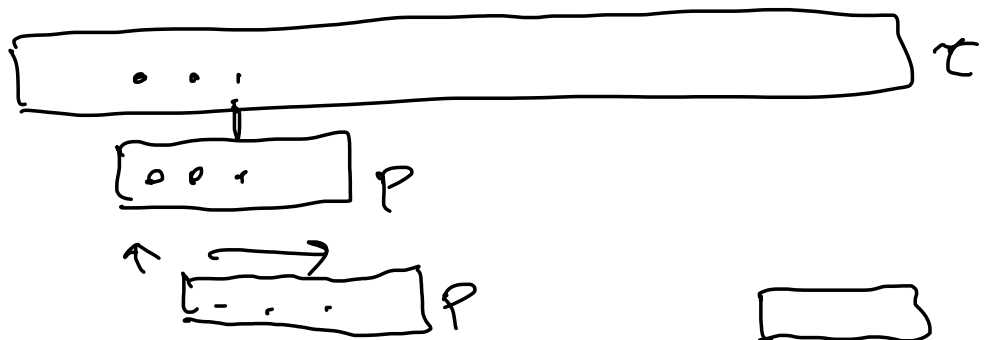
$p = \text{orl}$

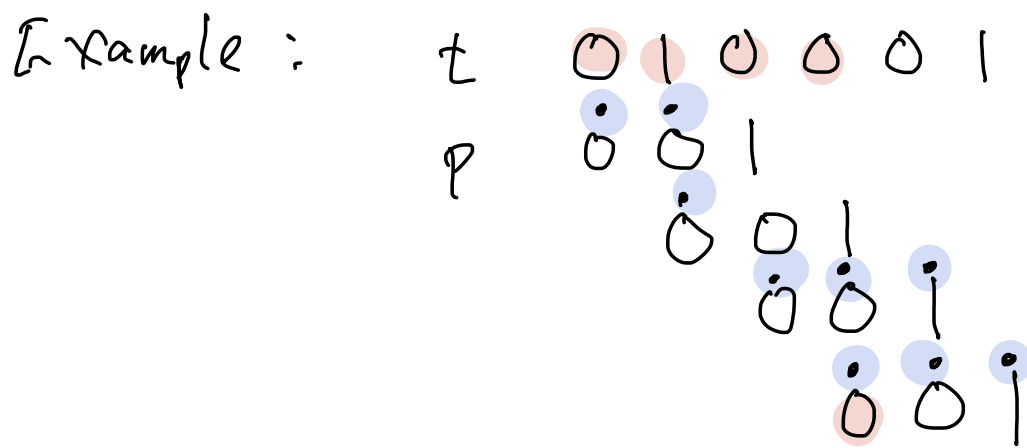
return value: -1 if not found.

index of start of pattern if found.

$|p| = m \quad |t| = n$

① Brute force.



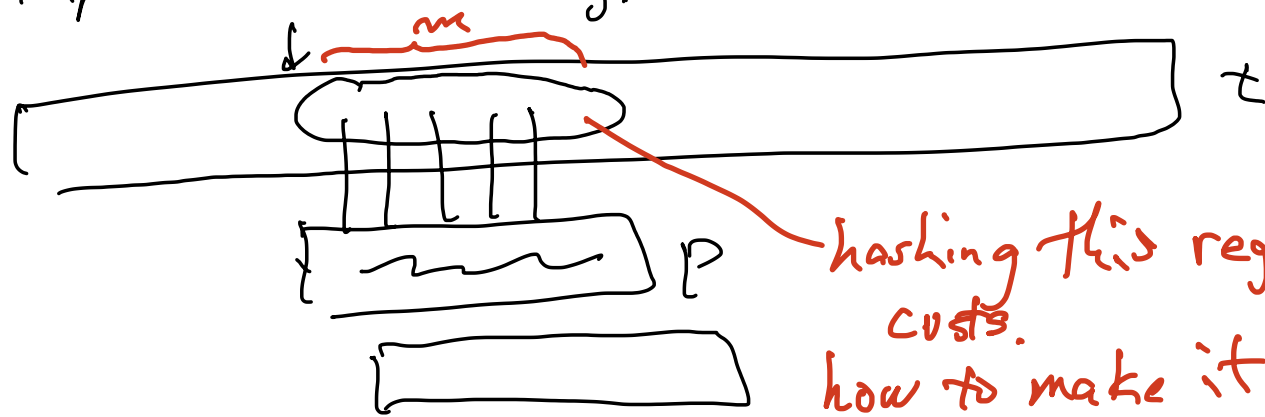


Complexity: (comparisons)
 best case: $n - m$
 worst case: \uparrow
 $(n - m)m$ \uparrow $\Theta(n)$
 \uparrow
 $\Theta(n \cdot m)$

In practice: $\ll n$ comparisons.

② Rabin-Karp 1987

Idea: do a quick hash-based check every time you slide the pattern over to avoid impossible matching.



Stage 1 of the development.

"fingerprinting"

based on parity: (# of 1 bits in a string) $\text{mod } 2$.

assume t, p are just 0,1 strings.

$$\text{parity of } p = \left(\sum_{0 \leq j < m} p[j] \right) \text{ mod } 2.$$

window
 \downarrow

\uparrow
 implement by
 & 01

t 0 0 1 0 1 1 0 1 0 1 0 0 1 0 1 0 0 1 1
 sliding parity of t \cdot 1 1 0 1 1 0 1 0 1 ...
 p 0 1 0 1 1 1 parity = 0
 ↑ ↑ ↑

sliding parity of $t = t \text{Parity}_j$

$t \text{Parity}_0 =$ direct computation

$$t \text{Parity}_{j+1} = (t \text{Parity}_j + t[j] + t[j+m]) \pmod 2$$

Stage 2: use XOR (\wedge) instead of parity.

Now we expect full comparison of characters $1/28$ th of the position of the window.

Cost: $O(n)$ parity updates, each constant time.

Stage 3: Reduce work not to $\frac{1}{2}$ (parity) or $\frac{1}{28}$ (XOR) but $1/g$ for a very high g .

Hash function for window starting at j

$$\left(\sum_{0 \leq i < m} 2^{m-1-i} t[j+i] \right) \pmod g.$$

Experience: g should be prime $> m$.

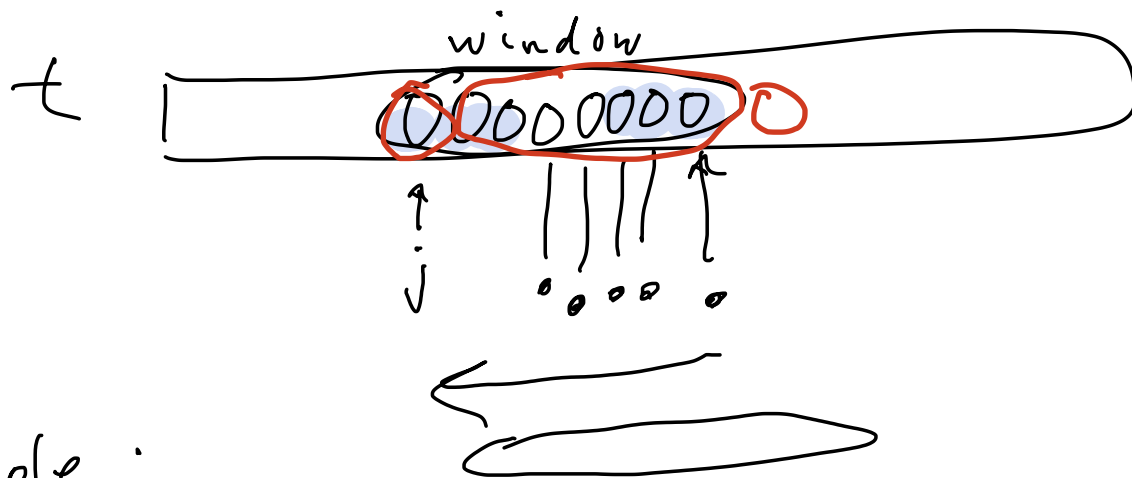
$$\text{Update } t \text{Parity}_{j+1} = t[j+m] +$$

$$2 \left(t \text{Parity}_j - 2^{m-1} t[j] \right)$$

$\pmod g$ ← still expensive

Use shifting to multiply by powers of 2.

Chances of getting a mismatch if hash agrees. is very low.
 \Rightarrow skip the inner loop, just assume success.



Example:

$$m = 10$$

$$n = 1000$$

choose g : prime, close to $m n^2$

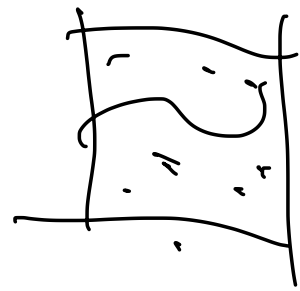
$$g = 9999991$$

probability of "mistake" = "false positive"

$$= \text{"hash collision"} = \frac{1}{g}$$

Monte-carlo: omit the point-by-point verification
 not worth the benefit.

Las Vegas version: check anyway.



Experience: $O(n)$ comparisons.

work: $7n$

not worth it.

③ Knuth-Morris-Pratt 1970-1977.

t: Tweedledee and Tweedledum

p: Tweedledum

p | Tweedledum
 1 2 3 4 5 6 7 8 9 10

slide window a large amount.

t: p a p p a p p a p p a r

p: p a p p a r

p a p p a r

p | p a p p a r
 1 | 1 2 3 3 6

p a p p a r

Precompute a "shift table" based on p.

t: p a p p a p p a s s s s s s s s

p: p a p p a r

p a p p a r

p a p p a r

p a p p a r

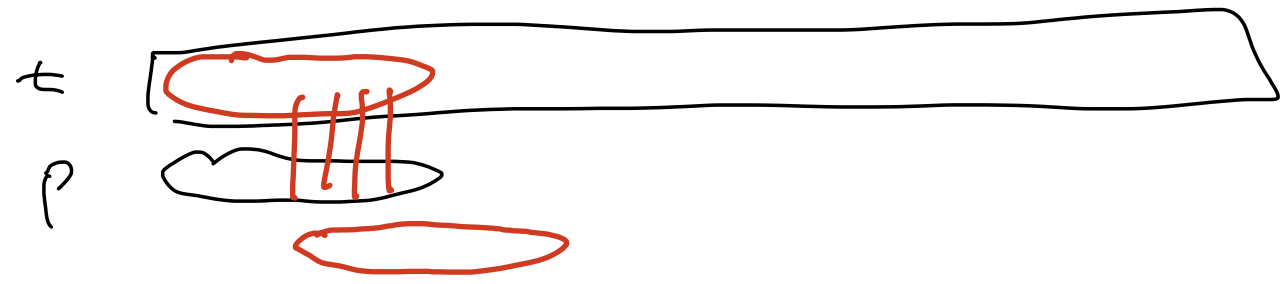
p a p p a r

How to build the shift table:

omitted. $\Theta(m)$

④ Boyer-Moore 1977

step 1: always search from end of pattern



Step 2: Occurrence heuristic

↑ method that is likely to improve performance.

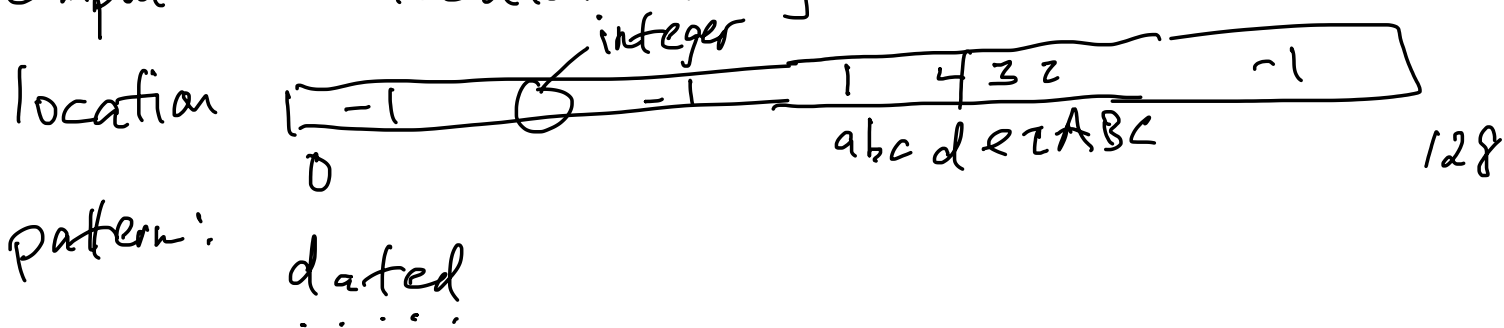
At a mismatch: Say see letter α in t ,
 shift p so that the rightmost occurrence of
 α in p lines up with that α in t .
 (don't slide left)
 (if α is absent from p , slide p to one
 place beyond α)

$t = \text{conundrum}$	}	$t = \text{conundrum}$
$p = \text{rum} \quad \text{rum} \quad \text{rum}$		$p = \text{drum} \quad \text{drum}$

$t = \text{conundrum}$	
$p = \text{natu} \quad $	
natu	

$t = \text{detective}$
$p = \text{date} \quad \text{date}$

precompute a "location" array



for (charVal = 0; charVal < 128; charVal += 1) {
 location[charVal] = -1;

}
for (pIndex = 0; pIndex < m; pIndex += 1) {
 location[p[pIndex]] = pIndex;

}

on mismatch, tIndex += max(1, pIndex - location[d])
(see d intent)

case: d not in p, pIndex = m-1

$$\text{location}[d] = -1$$

$$\max(1, m-1 - (-1)) = m$$

case: d not in p, pIndex = j

partial shift: j + 1

case d is in p.

shift enough so last instance of d in p
aligns with d in t.

Match heuristic: Use a shift table, like

Kaath-Morris-Pratt.

Nigel Horspool's heuristic:

at a mismatch, consider β , the character in t
where we started matching. Shift pattern

so β in t aligns with its rightmost
occurrence in p (not counting last place)

t d e t e c t i v e
p d a t e d a t e

Typical complexity:

n/m
 $O(nm)$

Other pattern-matching developments
approximate matching (Ed: Manber)

agrep
exact matching: grep ← print
 ↑ ↓
 global regular
 expression

example (in ed) g/date/p

grep date FILE ...

egrep (more complex expressions)

zgrep (for compressed files)

Advanced: regular expression (Perl regexp)

• exact strings: conandrum

• don't care symbol: con.ndr.um

↑ ↑
don't care

• character classes: c[ou1-5]nandrum

• predefined character classes:

\w	word-like	\W	not word-like
\s	space-like	\S	not space
\d	digit	\D	not digit

• Unicode character classes:

$\backslash p \{ASCII\}$

$\backslash p \{digit\}$

$\backslash p \{Final-Punctuation\}$

• pseudocharacters

\wedge start of text

$\$$ end of text

example: \wedge drum $\$$

• building expressions from others

alternation: drum | cone

repetition $um+$: 1 or more times

um^* : 0 or more times.

$\{uoth\}sthe$

$k\{2,4\}$

• capture groups

$(drum|cone)(foo)$

$\$1$

$\$2$

• zero-width assertions.

$(?=con)$ condition

90 man perlre

perlunicode

Edit distance : how hard is it to change one

string into another by: (1) delete char,

(2) insert char (3) replace char.

ghost \xrightarrow{D} host \xrightarrow{I} houst \xrightarrow{R} house

peseta → presto

		0	1	2	3	4	5	6
			p	e	s	e	t	a
1	p	0	-	1	-	2	-	3
2	r	1	0	-	1	-	2	-
3	e	2	1	0	-	1	-	2
4	s	3	2	1	0	-	1	-
5	t	4	3	2	1	0	-	1
6	a	5	4	3	2	1	0	-

↓ insert
 → delete
 ↘ no change
 ↙ replace

peseta
 presto

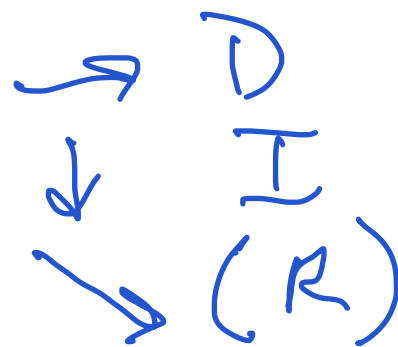
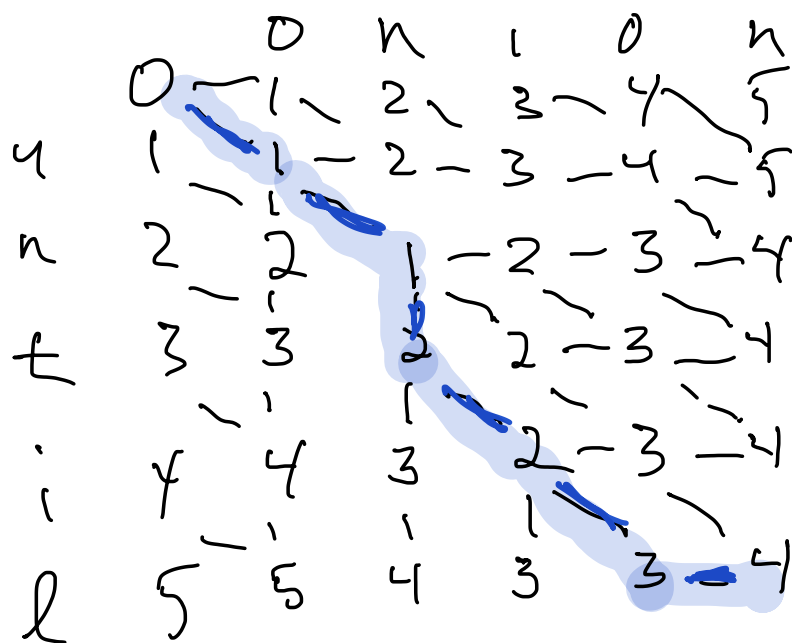
Complexity: lengths are n, m
 $O(nm)$

↘ (I)
 ↙ (R)

		0	1	2	3	4	5	6
			b	a	n	a	n	a
1	a	0	-	1	-	2	-	3
2	n	1	0	-	1	-	2	-
3	t	2	1	0	-	1	-	2
4	e	3	2	1	0	-	1	-
5	r	4	3	2	1	0	-	1
6	a	5	4	3	2	1	0	-

1
 b a n a n a
 a n a n a
 2
 a n t a n a
 3
 a n t e n a
 n

Onion \rightarrow until



- onion
- union
- untion
- untiln
- until

Category: dynamic programming
 solve all smaller problems first.

• Divide + conquer

if problem is trivial, solve it.

divide problem into a smaller problems.

of size n/b

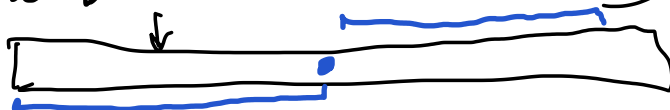
do those easier problems

combine the answers (at cost n^k)

In general $C_n = n^k + a C_{n/b}$

Examples

① j th smallest element of array

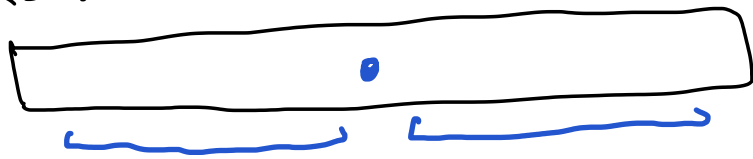


$a = 1$
 $b = 2$
 $k = 1$

$$a \quad b^k$$

$$1 < 2 \quad : \quad \Theta(n^k) = \Theta(n)$$

② Quicksort



$$a = 2$$

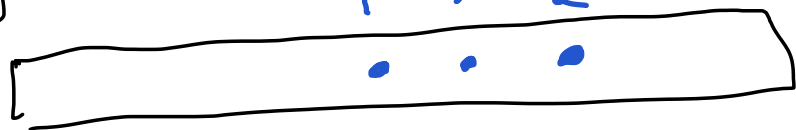
$$b = 2$$

$$k = 1$$

$$a \quad b^k$$

$$2 = 2 \quad : \quad \Theta(n^k \log n) = \Theta(n \log n)$$

③ Binary search (array)



$$a = 1$$

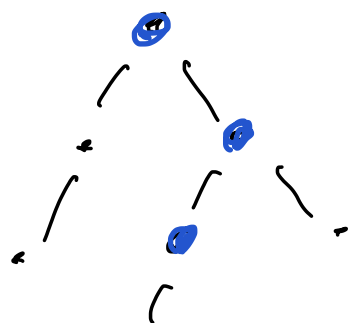
$$b = 2$$

$$k = 0$$

$$a \quad b^k$$

$$1 = 1 \quad : \quad \Theta(n^k \log n) = \Theta(\log n)$$

④ Binary search (tree)



$$a = 1$$

$$b = 2 \quad \Theta(\log n)$$

$$k = 0$$

⑤ Multiplication (Karatsuba)

$$x = a \cdot 10^{n/2} + b \quad a = 3$$

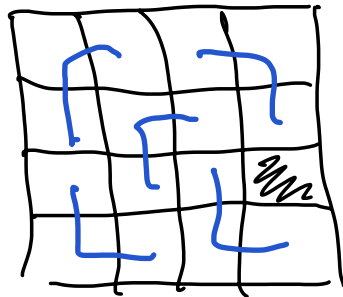
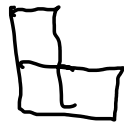
$$y = c \cdot 10^{n/2} + d \quad b = 2$$

$$k = 1$$

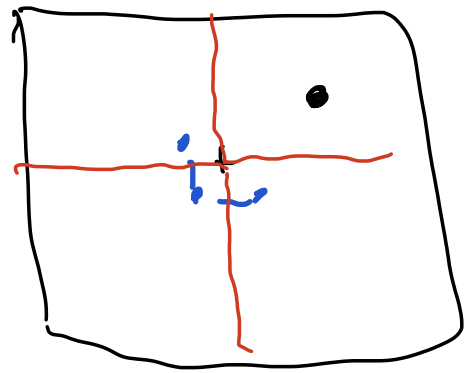
$$a \quad b^k$$

$$3 > 2 \quad \Theta(n^{\log_b a}) = \Theta(n^{\log_2 3}) \\ \approx \Theta(n^{1.5})$$

⑤ Tile an $n \times n$ board with triminos missing 1 cell. ($n = 2^p$)



n



$a \quad 4$
 $b \quad 2$
 $k \quad 0$

$a \quad b^k$

$4 > 1$

$$\Theta(n^{\log_b a}) = \Theta(n^{\log_2 4})$$

$$= \Theta(n^2)$$

(if n is the side of the board)

if $n = \#$ of cells:

$a \quad 4$
 $b \quad 4$
 $k \quad 0$

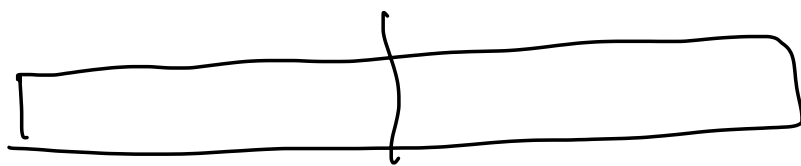
$a \quad b^k$
 $4 > 1$

$$\Theta(n^{\log_b a})$$

$$= \Theta(n^{\log_4 4})$$

$$= \Theta(n)$$

⑥ Mergesort



$a \quad 2$
 $b \quad 2$
 $k \quad 1$

$$\frac{a}{2} = \frac{b^k}{2} \quad \Theta(n^k \log n) = \Theta(n \log a)$$

• Greedy: enlarge a partial solution by best single step.

① Compute coins for amount
use biggest coins first

76¢: 25, 25, 25, 1

not optimal: coins 1, 6, 10,

to get 12¢: 10, 10, 1, 1

1, 5, 10 coins: greedy optimal

0z : T : ... : C : p : g : ...

weights on barbell:

② Kruskal's algorithm,
greedily add edges

③ Prim's algorithm
greedily add vertices

④ Dijkstra's algorithm for all shortest paths
expand the shortest path so far.

⑤ Huffman codes for data compression

frequencies					
L	60	S	6	1	1
A	22	T	4	A	0
O	16			R	1
R	13			S	10
				T	6

text with all those letters, at those frequencies \Rightarrow 121 chars
 $7 \text{ bits/char} = 847 \text{ bits}$

code table:

W	0
A	111
D	110
R	101
S	1001
T	1000

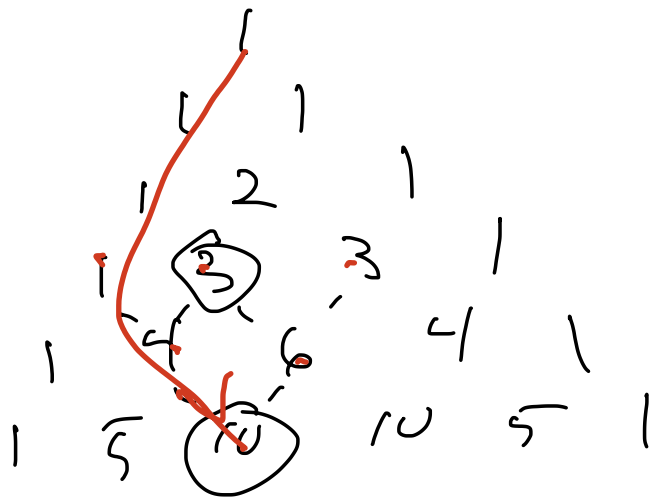
$$\frac{600}{W} + \frac{22 \cdot 3}{A} + \dots$$

$$= 253 \text{ bits}$$

⑥ Add a million floats, all in $[0 \dots 1]$ with best precision.

repeatedly: add two smallest numbers, put result back in the set.

Dynamic: $\binom{10}{5}$



Fibonacci Numbers

1 1 2 3 5 8 13



Hashing

Basic idea: store data that has key k in an array at index $h(k)$

↑
hash function.

Handling collisions:

open addressing: (clustering)

find an alternative index

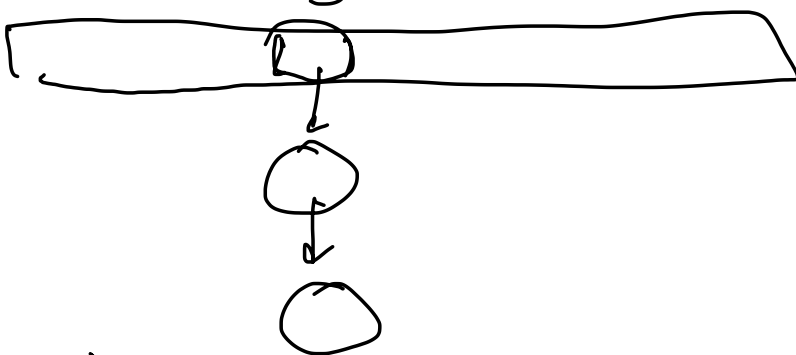
linear probing $p_i = h(k) + i$

quadratic probing $p_i = h(k) + i^2$

add-the-hash $p_i = h(k)(i+1)$

double hashing $p_i = h_1(k) + i h_2(k)$

external chaining



Good hash function

fast

uniform (equally likely to result in $0 \dots s-1$)

How big should the array be?

open addressing: $3 \times$ # of insertions.

external chaining: $1 \times$ # of insertions

If too small:
rehash into larger table.

Cryptographic hashes (digests)

one-way function from strings to bits
↑ ↑
arbitrary length example: 256

Purpose:

store passwords.

noticing plagiarism.

for authentication (signatures)

intrusion detection.

Graphs

vertex / vertices
edges

directed
undirected

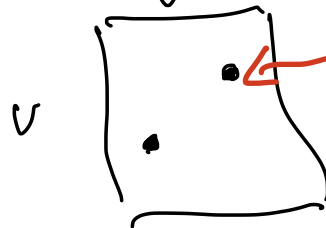
weighted
unweighted

dense
sparse

representations

special-purpose (genealogical)

adjacency matrix



1 if there is an edge
or weight

adjacency list

Algorithms

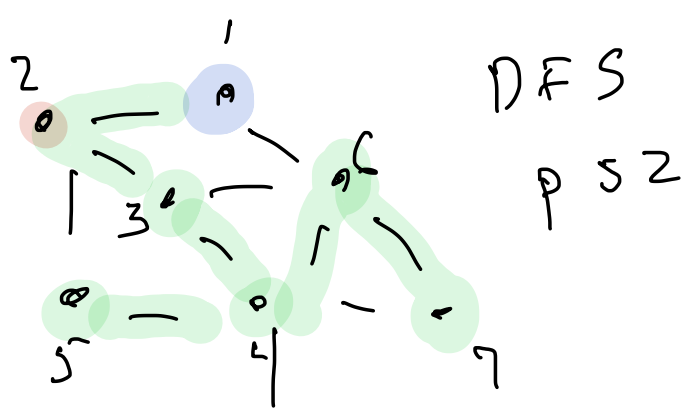
compute degree of each vertex

Depth-first search (DFS)

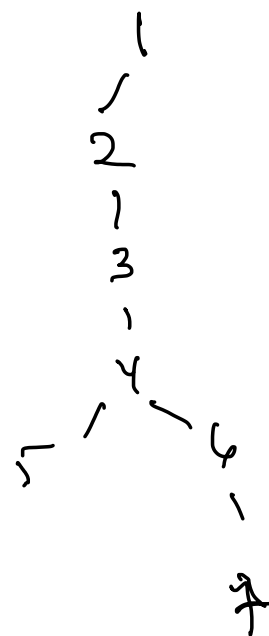
keeping track of what vertices have been visited

check for connectivity

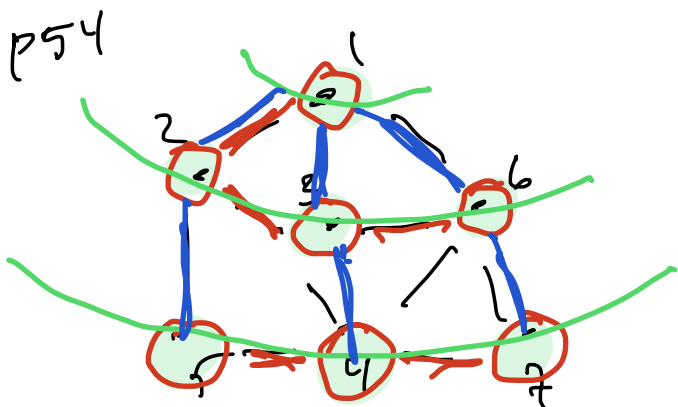
Breadth-first search (BFS)



navigation tree



if weighted,
queue \rightarrow heap.
p 55



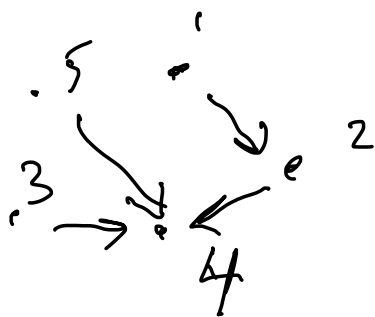
Dijkstra's algorithm: all shortest paths
from a starting vertex

Spanning trees (subset of edges connecting vertices)

Prim: add best vertex (closest to growing tree)

Kruskal: add best edge (shortest, without
introducing a cycle)

Union-Find



- u 1, 2
- u 4, 3
- u 5, 4
- u 1, 3

Numerical algorithms

Euclid GCD (greatest common divisor)

$$\begin{array}{r|l|l} a & 34 & 12 & 3 \\ b & 12 & 3 & \boxed{0} \end{array}$$

Fast exponentiation

$$243 \stackrel{745}{\leftarrow} \text{mod } 452$$

represent in binary (0011101001)

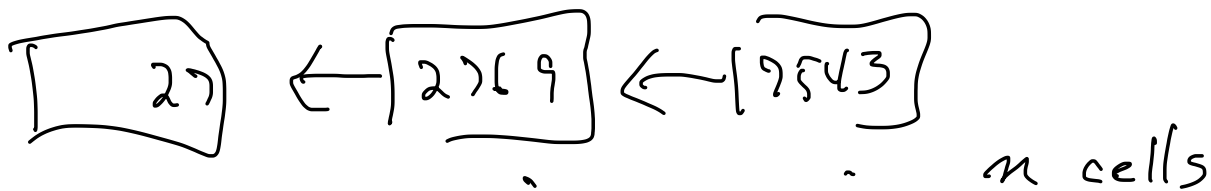
$$\left(\left(\left(\left(\left((1 \cdot a)^2 \right) a \right)^2 a \right)^2 a \right)^2 a \right)^2 a$$

14 multiplies
14 mods

Integer multiplication

Big nums: linked list of smaller chunks.

$$7543296927$$



Karatsuba's method.

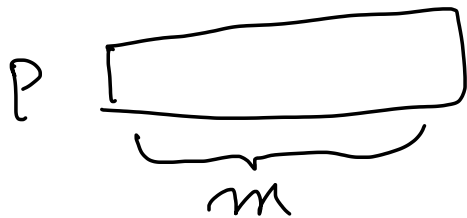
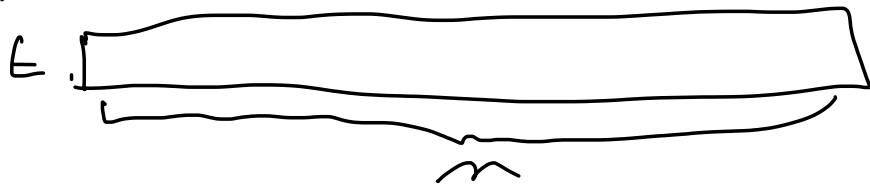
$$\text{if } n \geq 2, \quad \boxed{J} \in \boxed{J}$$

$$\boxed{J} \in \boxed{J}$$

4 multiplies
 \Rightarrow 3 multiplies

$$C_n = n + 3C_{n/2} \quad \Theta(n^{\log_2 3}) \approx \Theta(n^{1.6})$$

Strings, pattern matching



- 1) Brute force: try p at every location in t
- 2) Rabin-Karp: only try p if signature of p matches signature of the slice in t .
Constant-time update of slice's signature
- 3) Knuth-Morris-Pratt
on failure, slide p based on a shift table
- 4) Boyer-Moore
test p in reverse order, shift based on occurrence heuristic.
shift table

Horspools
Edit distance

Categories

- greedy: each step, do the best thing
- divide and conquer: subdivide into smaller problem: Recursion theorem
- dynamic programming: solve all smaller problems

Tractability $O(n)$ good $O(n^3)$ ok. $O(n^{20})$ tractable $O(2^n)$ bad.