

9/4/25 ALGORITHM: (mechanized)

1: correctness on every possible input

2: efficiency (metrics: time, memory, )

asymptotic,  $O(n)$ , hidden constants

3: optimality

ADDING INTEGERS: two ints  $x$  &  $y$ , compute  $x + y$   
 $y$  times

preschool:  $x + \underbrace{1+1+\dots+1}_{y \text{ times}}$



Ex:  $n=3 \rightarrow 123$  ( $n$  is # digits is exponential)

$O(10^n)$

grade-school:  $\begin{array}{r} 123 \\ + 456 \\ \hline 579 \end{array}$  adding 3 digits  $n$  times

1 → need  $2n$  digits      2:  $O(n)$   
( $2 n$ -digit #'s)

2 →  $n+1$  digit # in output

MULTIPLICATION:  $x * y$

$\begin{array}{r} 123 \\ \times 456 \\ \hline 546 \end{array}$  single dig. mult.  $x_i * y_j$  for  $\forall i, j$  (multiplies each pair)

$\Theta(n^2)$  multiples



is it optimal??  $\rightarrow O(n)$  is possible (read / write)

DIVIDE & CONQUER (assume  $n$ =even) Ex:  $x=1234 \dots$  same for  $y$

multiply ( $x, y$ )      (no floors/ceil)

$x_{hi}=12$

$x_{lo}=34$

base case:  $n=1$

lower order (R)

$$x * y = (x_{hi} \cdot 10^{n/2} + x_{lo})(y_{hi} \cdot 10^{n/2} + y_{lo})$$

$$\text{recursion: } x_{lo} = 10^0 x_0 + \dots + 10^{n/2-1} x_{n/2-1} = (x_{hi} \cdot x_{lo}) 10^n + (x_{hi} y_{lo} + x_{lo} y_{hi}) 10^{n/2} + x_{lo} y_{lo}$$

higher order (left)  
 $n/2$  digits

$$A = \text{mul}(x_{lo}, y_{lo})$$

$$B = \text{mul}(x_{lo}, y_{hi})$$

$$C = \text{mul}(x_{hi}, y_{lo})$$

$$D = \text{mul}(x_{hi}, y_{hi})$$

]

$$\rightarrow \text{out: } A + 10^{n/2} (B+C) + 10^n D$$

BASIC OPS:  
+, -, \*, shift

$T(n)$ : # basic ops for multiply ( $x, y$ )  $\sim n$  digits

$$T(n) = 4T(n/2) + C, n$$



RECURRENCE SOLN: masters thm  $\rightarrow \Theta(n^2)$

KARATSUBA: 1st math algo. better than  $\Theta(n^2)$   $k < 2$

**Master theorem:** Let  $T$  be a recurrence of the form  $T(n) = aT(n/b) + g(n)$ .

Case 1: If  $g(n) = O(n^{\log_b(a)-\epsilon})$  for some constant  $\epsilon > 0$ , then  $T(n) = \Theta(n^{\log_b(a)})$ .

Case 2: If  $g(n) = \Theta(n^{\log_b(a)} \log^k(n))$  for some constant  $k \geq 0$ , then  $T(n) = \Theta(n^{\log_b(a)} \log^{k+1}(n))$

Case 3: If  $g(n) = \Omega(n^{\log_b(a)+\epsilon})$  for some constant  $\epsilon > 0$  and  $ag(n/b) < cg(n)$  for some constant  $c < 1$  and sufficiently large  $n$ , then  $T(n) = \Theta(g(n))$ .

### KARATSUBA:

$$A = \text{multiply}(x_{10}, y_{10})$$

$$B = (x_{10}, y_{hi})$$

$$E = \text{multiply}(x_{10} + x_{hi}, y_{10} + y_{hi})$$

$$B+C = E-A-D$$

$$T(n) = 3T(n/2) + C_2 n \leftarrow \text{apply master thm}$$

$$\Theta(n^{\log_2 3}) \quad \log_2 3 \approx 1.58$$

3 multiplications, each on size  $n/2$

### HIGH PRECISION DIVISION:

- requires many bits
- need large integer multiplication

**CORRECTNESS:** right answer on all inputs

**EFFICIENCY:** if it makes least/small use of resources (time, space, ...)

### WORST CASE ANALYSIS:

- $T(n)$  = number of "operations" algorithm takes on worst input
- Asymptotics to summarize  $T(n) \in \Theta(f(n))$

### DEF: BIG-O

$O(f(n))$  algorithm A has (worst-case) runtime  $O(f(n))$  if  $\exists$  constants  $c, n_0 > 0$  s.t. on all inputs of length  $n > n_0$  it holds worst case runtime of A  $\leq c \cdot f(n)$



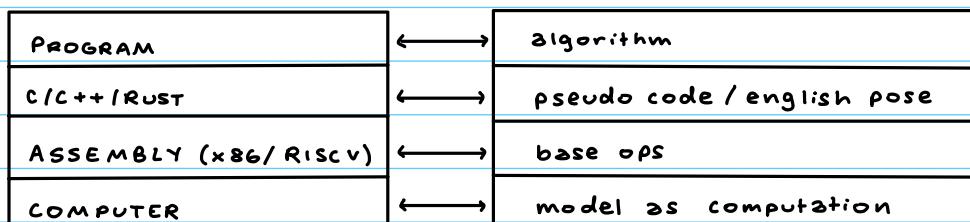
$\Omega(f(n))$  algorithm A has (worst-case) runtime  $\Omega(f(n))$  if  $\exists$   $c, n_0 > 0$  s.t.

an input for all  $n > n_0$  s.t. A takes time  $\geq c \cdot f(n)$  on that input



CAVEAT: alg A  $T(n) \leq 2^{1000} \cdot n \in O(n)$  asymptotically:  $O(n)$  beats  $O(n^2)$

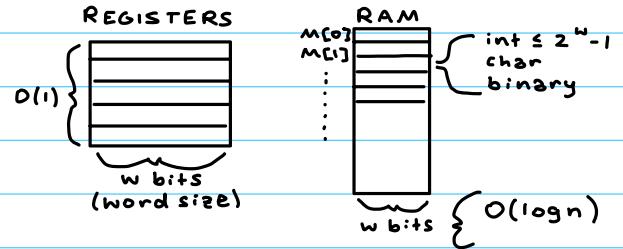
alg B  $T(n) \leq n^2 \in O(n^2)$  practically: constants matter, for  $n \leq 2500$ ,  $O(n^2)$  alg. might be faster.



### WORD-RAM MODEL:

- mem
- registers
- processor

int math: #, +, //, %  
 logic: and, or...  
 bitwise: &, |, ~  
 load:  $R[j] \leftarrow M[R[i]]$   
 store:  $M[R[i]] \leftarrow R[j]$   
 branch  
 half



### EXECUTION ON WORD-RAM:

1) Programs input starts out in mem locations  $M[0], \dots, M[n-1]$

2) processor steps through alg until HALT

3) output must be in first 17 mem. locations

$2^w$  addresses  $\geq n$  input size

$$w \geq \log_2 n \quad w = \text{word}$$

$$w \geq 2 \log_2 n \rightarrow n^2 \text{ mem. locs.}$$

\*use Pythonic model instead & under the hood is word RAM

$$\begin{aligned} \text{input len } n, \\ \text{word len } w = O(\log n) \end{aligned}$$

Correct answer

Question 4 1/1 pts

Select all correct answers. Assume the word-RAM model of computation with word size  $w$ .

Finding the minimum of a list of positive integers of length  $n$  takes:

$O(n)$  time, assuming all numbers fit in one word.  
We can scan the list and keep track of the current minimum.

$O(n)$  time, assuming all numbers are strictly less than  $2^w$ .  
If all numbers are smaller than  $2^w$ , then they fit into a single word. We can scan the list and keep track of the current minimum, which also fits into a single word.

$O(1)$  time.

$O(1)$  space.

### COMMON Python objects & running times:

INTEGER     $x+y$               Schoolbook     $O(|x|+|y|)$

$x*y$               karatsuba     $O((|x|+|y|)^{\log 3}) \leftarrow T(m) = 3T\left(\frac{m}{2}\right) + O(m)$

LIST               $L.append(x)$               ?               $\approx O(m^{\log 2})$

len(L)

### PEAK FINDING

INPUT: a list "A" of  $n$  integers      index is "peak" of A

OUTPUT: a "peak" of list A      if  $A[i]$  is at least as large as neighbor

[1, 3, 4, 2]



[1, 1, 1, 1, 1]

### ALG FOR PEAK FINDING:

• find-peak(A):

```

n=len(A) // O(1)           runs in worst-case time linearly
for i=0,...,n-1: // O(n)   O(n) worst case
  if i is peak return i // O(1)
  
```

### STEPS TO SUCCESS in algs:

1) write alg. in English prose

↳ pseudocode if necessary

2) prove alg. is correct on all inputs

3) analyze running time

## DIVIDE & CONQUER ALG (more efficient peak finding)

A 



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

$$\leq O(\log n)$$

`find-peak(A):`

```
n = len(A)
if n == 1: return 0
j = n/2
if A[j] peak: return j
if A[j] < A[j-1]: return find-peak(A[:j])
else: return find-peak(A[j+1:])
```

9/11/25 SEQUENCES + AMORTIZATION

### DATA STRUCTURES

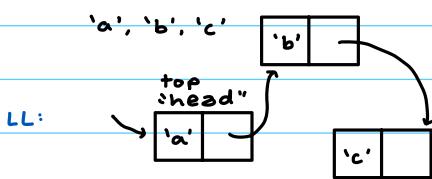
#### 2 COMPONENTS:

- interface (api, addt)

- implementation

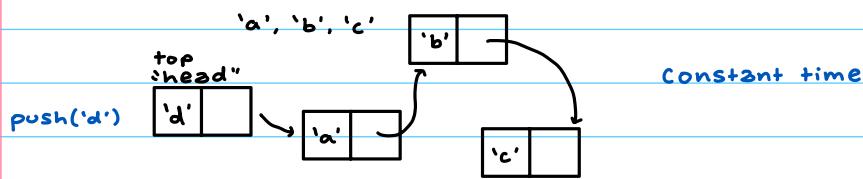
#### STACK: LIFO access pattern

- interface {
- `push(x):` stores x on "top" of stack
  - `pop:` remove & output "top" of stack
  - `peek:` return "top" of stack w/o removing
  - `is empty:` check if empty
- } constant time!!

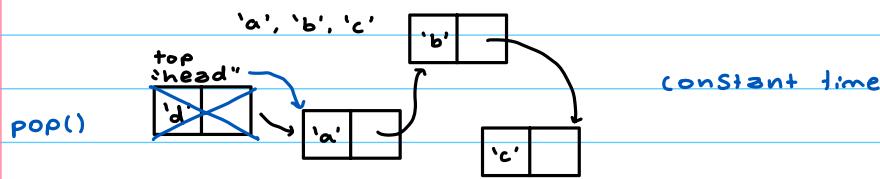


#### LINKED LIST:

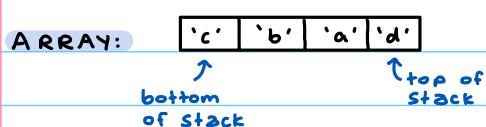
- pointer-based data structure



Constant time



constant time



\*problem: fixed-len arrays. we need to reallocate a new array & will take  $O(n)$  time

HOWEVER, if we copy into much larger array, can do many operations



## AMORTIZATION:

- charge expensive operations to the cheaper ones

Ex: teaholic buy 9 get 1 free

- averaging cost over sequence of operations & bounding

- amortized cost  $T$  if any sequence of  $m$  ops has total cost  $\leq m \cdot T$
- $\uparrow$  function of state of data structure  
 $\overbrace{\phantom{m \cdot T}}$  start from empty data structure

Ex: what is amortized cost of array if we make additional  $n$  slots?  $\rightarrow$  constant time

cheap: total cost of  $m$  pushes (starting from empty data struct):  $m \cdot \Theta(1)$

expensive: total cost of copy operations:  $\sum_{i=1}^{\lfloor \log_2 m \rfloor} \Theta(2^i)$   $\leftarrow$  happens whenever you hit a power of 2  
Sum:  $m \cdot \Theta(1) + \sum_{i=1}^{\lfloor \log_2 m \rfloor} \Theta(2^i) = \Theta(m)$

amortized cost:  $\Theta(m) / m \rightarrow \Theta(1)$  amortized (constant)

\* PROBLEM: push many values then pop them all, we end up with a massive empty array (lots of memory)



**MEMORY OVERHEAD:** additional space used beyond what's actually used to store the data.

- dynamic arrays: linear  $\Theta(n)$

- linked list: linear

allowed to use several diff. ops, each of which we can use diff.  $T$ ?

**POTENTIAL:** if  $\exists \phi \dots$

function  
data  
struct

$\phi(A) \geq 0$   
operation that takes  $A_0 \rightarrow A_1$ , takes time  $t$ :  $t + \phi(A_1) - \phi(A_0) \in O(f(|A_0|))$

$\Rightarrow O(f)$  amortized cost

array potential func:

$$\phi(A) = c |n - m/2| \quad n = \# \text{elements in stack}$$

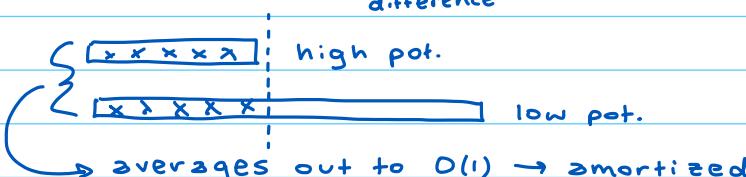
m = capacity of array

- potential decrease by linear amount

Ex: push w/o copy  $\Theta(1) + \Theta(1) = \Theta(1)$  cost

Ex: push with copy  $\Theta(n) - \Theta(n) = \Theta(1)$

cost      potential  
difference



**SEQUENCE:** maintains sequence  $x_0, x_1, \dots, x_{n-1}$

**ITER-SEQ():** iterates over sequence

**BUILD(A):** build sequence out of A  $\leftarrow$  container op

**GET-AT(:):** output  $x_i$

**SET-AT( $i, x$ ):** set  $x_i = x$

} static op

**INDEX-OF( $x$ ):** output  $i$  s.t.  $x_i = x$

**INSERT ( $at, first, last$ )**

**DELETE ( $at, first, last$ )**

SEQUENCE	Container	Operations					
		Static		Dynamic			
		build	get_at	index_of	insert_first	insert_last	insert_at
Linked List	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
Doubly-Linked List	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$	$\Theta(1)$	$\Theta(1)$	$\Theta(n)$
Dynamic Array	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$	$\Theta(n)$	$\Theta(1)$ (Am.)	$\Theta(1)$ (Am.)	$\Theta(n)$
2-End Dynamic Array	$\Theta(n)$	$\Theta(1)$	$\Theta(n)$	$\Theta(1)$ (Am.)	$\Theta(1)$ (Am.)	$\Theta(1)$ (Am.)	$\Theta(n)$

(Am.) = amortized running time

## 9/16 SETS (key-value maps)

Ex: {38, Henry}       $\{\text{key, value}\}$  pairs  
(63, Srinivas)      ↓  
(35, Brynmor)      keys are comparable, unique

OPERATIONS:

$\Theta(n)$  • build(A)

• length()

$\Theta(n)$  • find(key)  $\rightarrow$  returns value associated w/ key

$\Theta(n)$  • insert(key, value) might be able to do in amortized  $\Theta(1)$

$\Theta(n)$  • delete(key)

• find\_min()  $\rightarrow$  returns value w/ lowest key

• find\_max()

• find-next(key)  $\rightarrow$  returns value w/ smallest key larger than key

• find-prev(key)  $\rightarrow$  returns value w/ largest key smaller than key

$\Theta(n^2)$  • iter\_ord()  $\rightarrow$  each find takes n, so all finds will be  $n^2$  time

\* if we sort array, time could be  $\Theta(\log n)$

## ARRAY:

a[0], a[1], ...

## SORTING:

input: list or array, A, of comparable elements

output: list or array, B, permutation of A that is sorted.

$$\forall i < |A|, B[i] \leq B[i+1]$$

PERMUTATION SORT:  $\Omega(n!)$   $\gg 2^n$

- generates permutations of elements until sorted

INSERTION SORT \*stable & in place

Ex: 8 6 7 5 3 0 9  
  ~

1: 6 8 7 5 3 0 9  
  ~

2: 6 7 8 5 3 0 9

⋮

insert-last(A, i)

WORST CASE: every iteration

while  $i \geq 0$  and  $A[i] > A[i+1]$ :  $O(n)$   $\Omega(n) \leftarrow @ith position, A[i] smaller than anything in prefix$

swap  $A[i]$  with  $A[i+1]$   $O(1)$

$i--$   $O(1)$

insertion-sort(A):  $\Theta(n^2)$

for  $i$  in  $1 \dots |A|$ :  $O(n)$

insert-last(A, i)  $O(n)$

in-place:  $\Theta(1)$  memory overhead (modifying input directly - no extra mem.)

STABLE SORT: preserves order of equal items

in: (2, 3), (1, 1), (2, 4)

out 1: (1, 1), (2, 3), (2, 4)  $\leftarrow$  valid & stable

out 2: (1, 1), (2, 4), (2, 3)  $\leftarrow$  another valid but unstable out

```

1 def merge_sort(A, a = 0, b = None):                                # T(b - a = n)
2   """Sort A[a:b]"""
3   if b is None: b = len(A)
4   if 1 < b - a:
5     c = (a + b + 1) // 2
6     merge_sort(A, a, c)
7     merge_sort(A, c, b)
8     L, R = A[a:c], A[c:b]
9     merge(L, R, A, len(L), len(R), a, b)                            # S(n)
10
11 def merge(L, R, A, i, j, a, b):                                     # S(b - a = n)
12   """Merge sorted L[:i] and R[:j] into A[a:b]"""
13   if a < b:
14     if (j <= 0) or (i > 0 and L[i - 1] > R[j - 1]): # O(1)
15       A[b - 1] = L[i - 1]                                # O(1)
16       i = i - 1                                         # O(1)
17     else:
18       A[b - 1] = R[j - 1]                                # O(1)
19       j = j - 1                                         # O(1)
20   merge(L, R, A, i, j, a, b - 1)                                    # S(n - 1)

```

MERGE SORT: divide + conquer

Ex: 8 6 7 5 3 0 9 • has many edge cases

6 7 8 0 3 5 9 ) merge  $\Theta(n)$

0 3 5 6 7 8 9

\* not in place

\* can make stable by preferring left array

merge-sort(A): pseudocode

size of A

if  $|A| < 2$ : return A (already sorted)

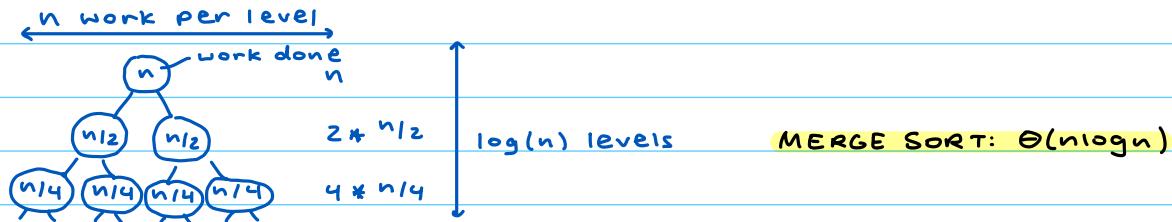
mid =  $|A|/2$

L = merge-sort(A[:mid])

R = merge-sort(A[mid:])

return merge(L, R)

$T(n) = \Theta(n) + 2T(n/2)$  each level is n work! equal!!

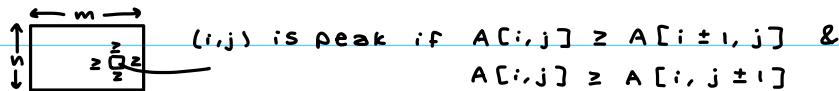


2D PEAK FINDING:



\* looking for ANY peak, not max peak

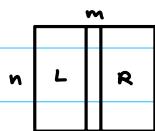
· Find-peak-1D(A) =  $\Theta(\log n)$



GOAL: algorithm that returns  $(i,j)$  a peak in A

BRUTE FORCE:  $\Theta(n \times m)$

IDEA: DIVIDE + CONQUER  $\leftarrow$  THIS IS A BUGGY SOLUTION



$$j_{mid} = \lfloor \frac{m}{2} \rfloor$$

col-peak = find-peak-1D

col-peak returns some value  $i$  s.t.  $(i, j)$  is a peak in column  $j$ .

check if  $(i, j)$  is peak in  $A$ .

CASE 1: it is return  $(i, j)$

CASE 2: it isn't. didn't fail N or S b/c we used col-peak(). must have failed E or W. recurse in uphill direction (L or R).

PROOF: induction on  $m$ .

INDUCTIVE HYPOTHESIS (preposterous): assume the alg works for arrays up to  $m-1$ . INDUCTION: show it works for  $m$  columns

BASE CASE:  $m=1$ : 1D peak finding, which we know to be correct & returns a peak in the column, which is a peak in  $A$  b/c 1D.

acquire  $j_{mid}$  & call col-peak( $A[i, j_{mid}]$ )  $\rightarrow i^*$

CASE 1:  $(i^*, j_{mid})$  is peak in  $A$ . ✓

CASE 2:  $A[i^*, j_{mid}]$  has greater value in L or R

WLOG (without loss of generality), higher value is in L. recursively call on L.

$(i, j)$ .  $A[i, j]$  is peak in L. b/c of preposterous hypothesis

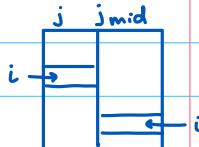
CASE 2a:  $j < j_{mid}-1$ . by inductive hypothesis  $(i, j)$  is peak in  $A$ .

CASE 2b:  $j = j_{mid}-1$ . show  $(i, j)$  is peak in  $A$ .

BUG: we don't know relationship b/t  $i$  &  $i^*$ . need to establish a relationship

& can't use 1D peak finder - must find max peak of columns

NEW INDUCTIVE HYPOTHESIS: assume alg. works up to  $m-1$ . ALSO assume it returns a max in col  $j$ .



## 9/18 FASTER SORTING

so far...

- insertion  $\Theta(n^2)$
- selection  $\Theta(n^2)$
- merge  $\Theta(n \log n)$

### COMPARISON MODEL:

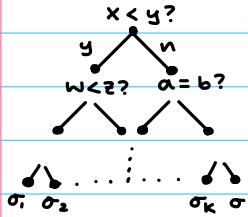
**Comparable:** black box that supports binary comparison operations

$<$ ,  $=$ ,  $>$ ,  $\leq$ ,  $\geq$ ,  $\neq$

**INPUT:** list  $x_1, x_2, \dots, x_n$ . key is comparable

**OUT:**  $\sigma \in S_n$  s.t.  $x_i \sigma$  key are increasing with  $i$ .

### DECISION TREE:



\*need at least 1 leaf for each permutation

every permutation must have  $\geq 1$  leaf

POSSIBLE  $\sigma$ 's  $\geq n!$  leafs

ROOT-TO-LEAF PATH  $\leftrightarrow$  runtime

**HEIGHT:** worst case runtime  $\geq \log(n!) \geq n \log n$

\*branches always binary (2 per level)

$$\geq \log\left(\left(\frac{n}{2}\right)^{n/2}\right)$$

$$= \frac{n}{2} \log \frac{n}{2} \in \Omega(n \log n)$$

\*need superconstant branching to beat  $n \log n$

**REFERENCING / DEREFERENCING:**  $A[B[i]]$  [5, 7, 2, 0, 4]

↓

$$i = 7+1 = 8$$

(direct access array)

DAA-SORT: assume all keys are IN  $\in [0, u]$  & DISTINCT

$$A[x.key] = x$$

$$\begin{bmatrix} 0 & 2 & 4 & 5 & 7 \\ 0 & 1 & 2 & 3 & 4 & 5 & 6 & 7 \end{bmatrix} \text{ array of size } 8$$

typically,  
just use  
this instead

- Stable by default b/c nothing can appear twice
- (no duplicate keys to preserve man...)

[0 2 4 5 7] copy array back

what if keys not distinct??

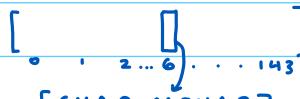


COUNTING SORT

→ allocate a sequence within list

Sorting by #s	
CHAR	0 0 6
PIKA	0 2 5
SNDR	1 4 3
MCHAR	0 0 6

key collision for diff. values



- maintain stability: must ensure sequence memorizes insertion order (queue, etc)

- store "chain" or queue of all elts. w/ same key
- concat all chains together to create output list
- stable, but not in place (uses extra mem. to sort)

assume keys are tuples of equal length.

sort by Col A, Col B lexicographically

**TUPLE LENGTH:** sort by key components individually in order of significance

Ex: sort cards by suit (clubs, diamonds, hearts, spades) then #

· counting sort: put cards into suit order, sort within suit,

put all back together OR

sort by #s, then by suit  $\Rightarrow$  same output!

TLDR: doesn't matter whether we sort by less/more significant key; if sort by less sig. first, need to make sure we preserve order (stability)

**RADIX SORT:** what if u is larger? e.g.  $u=n^2$

$3n + 5 \rightarrow (3, 5)$

**TUPLE SORT** using **COUNTING SORT** as our auxillary algorithm  
keys in base n  
MUST SAY: "use tuple sort with auxillary counting sort alg"

Ex:  $[17, 3, 24, 22, 12] \rightarrow [(3, 2), (0, 3), (4, 4), (4, 2), (2, 2)]$

$\Theta(n + n \log n)$        $(0, 3) \quad (2, 2) \quad (3, 2) \quad (4, 2) \quad (4, 4) \leftarrow 4 \cdot 5 + 4 = 24$

when  $\log u = O(1)$  if u is small:  $[3, 12, 17, 22, 24]$   
 $\rightarrow$  sorting takes linear time

$n^4 v, n \log n v,$

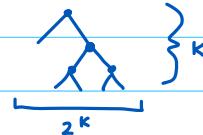
$2^n x$

SORTING ALGS:	RUNTIME:	STABLE?	IN PLACE?
insertion	$\Theta(n^2)$	✓	✓
selection	$\Theta(n^2)$	✗	✓
merge	$\Theta(n \log n)$	✓	✗
DAA	$\Theta(n + u)$	✓	✗
counting	$\Theta(n + u)$	✓	✗
radix	$\Theta(n + u)$	✓	✗

	BUILD	FIND	INSERT/DELETE	FIND-PREV / FIND-NEXT
ARRAY	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$	$\Theta(n)$
SORTED ARRAY		$\Theta(n \log n)$	$\Theta(n)$	$\Theta(\log n)$
DAA	$\Theta(u)$	$\Theta(1)$	$\Theta(1)$	$\Theta(u)$
HASH TABLE	$\Theta(n)$	$\Theta(1)$ exp.	$\Theta(1)$ exp./amor.	$\Theta(n)$

bottom line: know when to use which

- SORTED ARRAY: range queries
- HASH TABLES: lookups by key
- ARRAYS: almost never



`find(key) → item`

$$2^k \geq n \Rightarrow k = \log_2 n$$

#dict is implemented w/ hash table

WARM-UP:  $O(1)$ -time find w/ huge memory overhead

use array of  $u$  elements: if all keys are in  $\{0, 1, \dots, u-1\} = [u]$



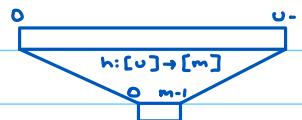
insert(3, item)

find(s)

Ex: student id's  $n < 10^9$  id #s  $\rightarrow$   $n$  students,  $u = 10^9 \rightarrow$  lots of wasted space

### HASHING:

- USE DAA OF SIZE  $m \ll u$



- collision when  $k \neq k' \in [u]$  s.t.  $h(k) = h(k')$

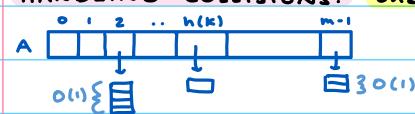
### HANDLING COLLISIONS: open addressing

- A  $h(k) \in [m]$
- store  $h(k')$  & item in next available empty bucket
- don't need another data struct
- however, insert/delete makes alg. complicated

### SIZING:

$n$  items  
 $m$  table size  
 $m \geq n$

### HANDLING COLLISIONS: chaining



- each "bucket" (spot) now points to another DAA
- hash table shouldn't be too small b/c we want DAA's in each index to be small

`find(k):  $\Theta(1 + L)$`

$$\Theta(1)$$

$$C = A[h(k)]$$

$$\Theta(1)$$

$$\text{return } C.\text{find}(k)$$

$$\Theta(L)$$

`insert(item):`

$$C = A[h(item.key)]$$

$$C = \text{insert}(item)$$

Ex: find a collision

$$h: \mathbb{Z} \rightarrow [10] \quad h_1(x) = x \bmod 10$$
$$h_2(x) = x \bmod 10000 \quad * \text{seems like it's doomed}$$
$$h_2(1) = h_2(10001)$$
$$h(x) = \text{Lcos}(\tan^{-1}(x+15))$$

\* pick hash function in secret after given keys.

SIMPLIFYING ASSUMPTION: have data structure pick truly random hash function

↳ hash is picked independently of the keys  $h: [n] \rightarrow [m]$

- SUHA

'EXPECTED' TIME for find / insert:

· imagine we've inserted items w/ distinct keys  $k_1, k_2, \dots, k_n \in [n]$

· insert  $k_{\text{new}} \notin \{k_1, k_2, \dots, k_n\}$

· look @ prob. that  $k$  hashes to same bucket as  $k_{\text{new}}$

$$\Pr[h(k_j) = h(k_{\text{new}})] = \sum_{i=0}^{m+1} \Pr[h(k_j) = i \wedge h(k_{\text{new}}) = i]$$
$$= \sum_{i=0}^{m+1} \Pr[h(k_j) = i] \cdot \Pr[h(k_{\text{new}}) = i]$$
$$= m \cdot \left(\frac{1}{m}\right) \cdot \left(\frac{1}{m}\right) = \frac{1}{m}$$



LOAD FACTOR =  $\alpha$

$$\text{the expected chain length } \mathbb{E}[L] = \mathbb{E}[L] = 1 + \sum_{j=1}^m \Pr[h(k_j) = h(k_{\text{new}})] = 1 + \frac{n}{m} = 1 + \alpha$$

EXPECTED:

$$\text{time to insert: } 1 + \alpha = 1 + \frac{n}{m} = 1 + \frac{n}{3m} \in O(1)$$

$m = \text{table size}$

$n = \# \text{ total items inserted}$

$L = \text{size of the bucket}$

STORAGE FOR TABLE:  $m+n$      $m > n$ ,  $m = \Theta(n)$

· takes lots of time to double hash table size like dynamic array.  
 $h: [n] \rightarrow [2m]$



Ex: python: `hash('cow')` → always get same # until you quit python.

· starting Python over → get a different #

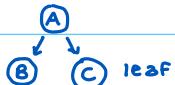
· good for preventing hackers

9/25 BINARY TREES & BINARY SEARCH TREES:

DATA STRUCTURE:	BUILD	FIND	INSERT/DELETE	FIND-MIN/FIND-MAX	FIND-PREV/FIND-NEXT	....
SORTED ARRAY:	$n \log n$	$\log n$	$n$	1	$\log n$	
HASH TABLE:	$n(e)$	$1(e)$	$1(a)(e)$	$n$	$n$	
GOAL:	$n \log n$	$\log n$	$\log n$	$\log n$	$\log n$	

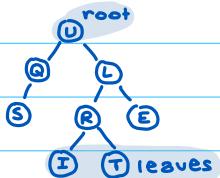
EX: key prime  
X o/p P

### BINARY TREE:



B = left child of A  
C = right child of A

A = parent of B & C



DEPTH of node:  $d(x) = \# \text{ of ancestors}$   
 $\rightarrow \text{root} = 0, \text{leaf} = 3$

HEIGHT of node:  $h(x) = \# \text{ edges of longest path from node to a leaf}$   
 $\rightarrow \text{root} = 3, \text{leaf} = 0$   
 $\rightarrow h = O(\log n)$

### TRAVERSALS: $O(n)$ , $n = \# \text{ nodes}$

IN-ORDER: visits first, recursively the left child, root, right child.

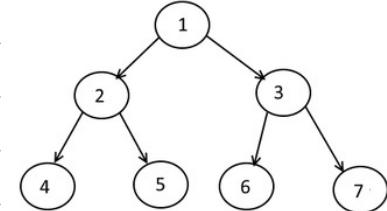
$\rightarrow S, Q, U, I, R, T, L, E$

PRE-ORDER: visits recursively the root, then left, then right nodes

$\rightarrow U, Q, S, L, R, I, T, E$

POST-ORDER: visits recursively the left, then right, then node.

$\rightarrow S Q I T R E L U$



Inorder Traversal: 4 2 5 1 6 3 7

Preorder Traversal: 1 2 4 5 3 6 7

Postorder Traversal: 7 6 3 5 4 2 1

Breadth-First Search: 1 2 3 4 5 6 7

Depth-First Search: 1 2 4 5 3 6 7

### BINARY SEARCH TREES: # distinct keys

for every node q, all the keys stored in the left subtree of q

are smaller than q's key & all the keys stored in the nodes

of q's right subtree are greater than q's key.



11 ← NOT BST b/c everything on left of root should be strictly smaller, but 11 < 10

$\text{find}(T, k)$ : find node in T whose item has # key k if it exists, else None

if T is None: return None

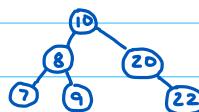
$O(h)$        $h = \text{height}$

if  $k < T.item.key$ : return  $\text{find}(T.left, key)$

if  $k > T.item.key$ : return  $\text{find}(T.right, key)$

return T

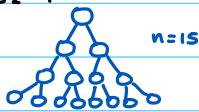
valid BST:



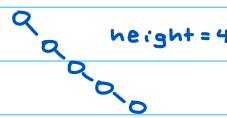
$O(h)$

$h = O(\log_2 n)$

Ex: height=3



$n=15$



height=4

**find\_min( $T$ ):** keep going left until can't anymore

**find\_max( $T$ ):** keep going right until can't anymore



**find-next( $T, x$ ):** next node after  $x$  in  $T$ 's inorder traversal

**CASE 1:** if  $x$  has right child, successor of  $x$  is smallest element in the right subtree of  $x$ .  $\text{find-min}(T.\text{right}, x)$

**CASE 2:** successor of  $x$  is 1st ancestor which we reach through left child.

```

1 def find_next(T):
2     # Finds next node after T in the in-order traversal
3     if T is None:
4         return None
5     if T.right:
6         return find_min(T.right)
7     child = T
8     while child.parent and child.parent.right == child:
9         child = child.parent
10    return child.parent

```

**delete( $T, x$ ):**  $\Theta(h)$

**CASE 1:** if node is a leaf, just delete



**CASE 2:** node only has one child, replace with child



**CASE 3:** node is internal. Swap with successor ( $\text{find-next}$ )

successor will not have a left child, or else it



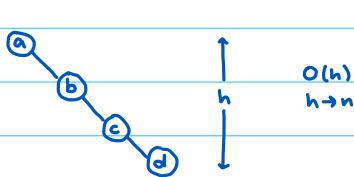
would not be the successor.

predecessor will not have right child.

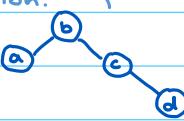
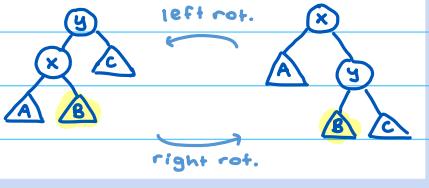


9/30

```
def in-order(T):
    if T is None: return
    in-order(T.left)
    print(T.item)
    in-order(T.right)
```



left rotation:

 $O(h)$  $h \rightarrow n$  $\Delta, \times, \triangle, \circlearrowleft, \circlearrowright$ 

```
def right_rotate(T):
    return Tree(T.left.left, T.left.item,
               left      node
               Tree(T.left.right, T.item, T.right))
```

left + node right

```
def left_rotate(T):
    return Tree(Tree(T.left, T.item, T.right.left),
               left
               T.right.item, ← node
               T.right.right) ← right
```

**PERFECTLY BALANCED TREE:** if  $T$  is full  $\forall$  every level!  $h = O(\log n)$



**SKEW(n):**  $height(n.right) - height(n.left)$

**AVL PROPERTY:** if  $\forall n \in T$ ,  $|skew(n)| \leq 1$

- as long as deepest leaf on one side is  $\leq 1$  deepest leaf on the other

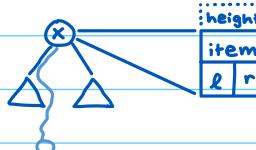
- a tree with AVL property has  $h = O(\log n)$



$n_h = n_{h-2} + n_{h-1} + 1 \leftarrow$  Fibonacci! exponential growth

$$n_h > 2 * n_{h-2}$$

$$n_h > 2^{\frac{h}{2}} \Rightarrow 2 \log n > h \Rightarrow h = O(\log n)$$



**T.height:** insert  $\mathcal{O}$  leaves

new node gets height = 0

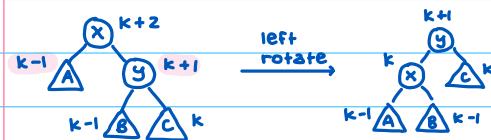
\*insert & delete are logarithmic!

### Deletions:

1) leaf: throw it away; update the ancestors

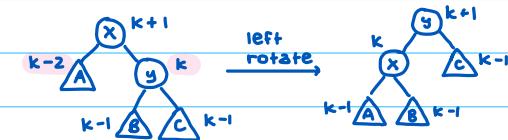
2) internal nodes have  $\leq 1$  child: replace node w/ child; update ancestors

#### CASE 1: skew=2



skew=0

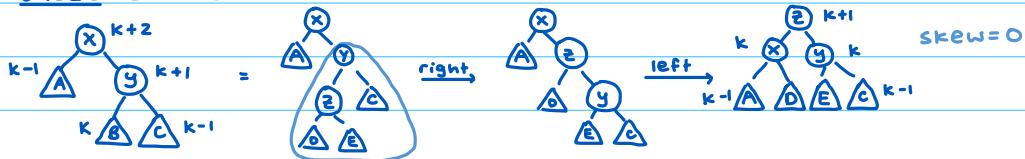
#### CASE 2: skew=2



skew=1

\*can't happen during insertion

#### CASE 3: skew=2



skew=0

**AVL SORT:** use AVL tree to represent our set.  $\Theta(n \log n)$

`def avl-sort(A):`

```
for a in A: tree.insert(a) ← each insert = log n, total insert [Θ(n log n)]
for x in tree: iter_order(): B.append(x) ← Θ(n)
return B
```

10/1

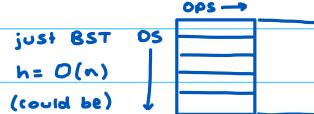
### AUGMENTATION OF TREES

- dynamic order statistics ← size field  
 $O(\log n)$  for dynamic sets for O.S.

- maintainery augmentations

- interval trees

last time:  $BST \rightarrow AVL$   
 $O(h)$   $O(\log n)$



### ORDER-STATISTIC TREE

- get  $i$ th order statistic: element in set with  $i$ th smallest key  
give  $i \rightarrow$  return element
- rank: give element  $\rightarrow$  return position  $i$  in ascending order
- $x.size = x.left.size + x.right.size + 1$

### RETRIEVING ELEMENT w/ GIVEN RANK:

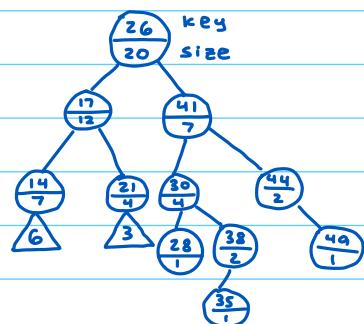
`OS-select(x, i):`

```
r = x.left.size + 1
if i == r: return x
else if i < r: return OS-select(x.left, i)
else: return OS-select(x.right, i - r) shifted rank
```

`OS-select(17, 17-13)`

4th

4 - 2 → 2nd

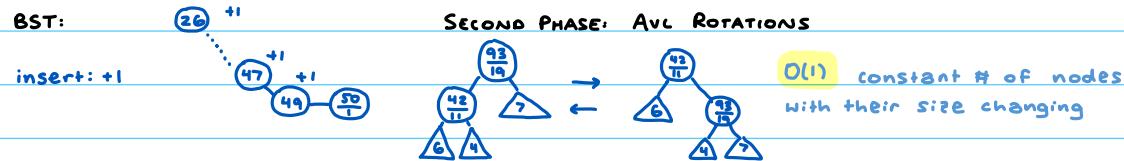


FIND RANK OF ELEMENT: given pointer to node  $x$

OS-Rank ( $T, x$ ):

```
r = x.left.size + 1  
y = x  
while y ≠ T:  
    if y == y.parent.right: r = r + y.parent.left.size + 1  
    y = y.parent  
return r
```

DYNAMIC SETS:

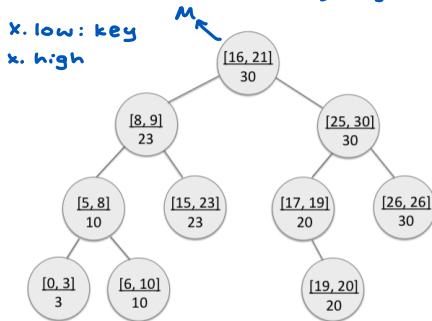


INTERVAL TREES: interval-search ( $T, i$ ) → find a node who's interval intersects  $i$  if it exists

low( $i$ ) high( $i$ )

$i \in [t_1, t_2]$

high( $i$ ) ≥ low( $i$ )



AVL BST

unsuccessful: [11, 14]

$x.M = \max(x.\text{high}, x.\text{left}.M, x.\text{right}.M)$

Interval-Search ( $T, i$ ):

$x = T$

while  $x \neq \text{None}$  &  $i$  doesn't overlap  $[x.\text{low}, x.\text{high}]$ :

if  $x.\text{left} \neq \text{None}$  and  $x.\text{left}.M \geq \text{low}(i)$ :

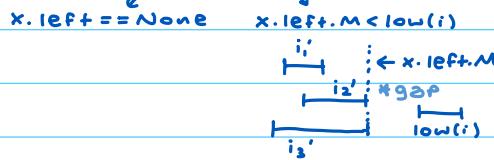
$x = x.\text{left}$

else:  $x = x.\text{right}$

return  $x$

### Interval-Search Proof:

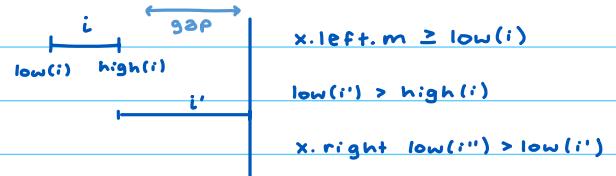
(a) suppose we go right



(b) Suppose we go left

need to show if no interval in  $x.\text{left}$  overlaps

then no interval in  $x.\text{right}$  will overlap



~ MIDTERM 1 ~

10/7 GRAPHS:  $G = (V, E)$  is a set of vertices  $V$  & a set of pairs of vertices  $E \subseteq V \times V$

#### UNDIRECTED:



- $E \subseteq \{u, v\}$  unordered edges
- size of  $E$ :  $|E| \leq \binom{|V|}{2} = O(|V|^2)$

#### DIRECTED:



- vertices are ordered pairs
- $E \subseteq (u, v)$  distinct edges
- $|E| \leq 2 \binom{|V|}{2} = O(|V|^2)$

#### SIMPLE GRAPHS:

- ✗ 1: no duplicate edges  $(u, v)$  for unique  $u, v$
- ✗ 2: no self looping  $(u, v)$  for  $u \neq v$

$\text{adj}^+(u) = \text{set of outgoing edges from } u$

$\text{adj}^-(u) = \text{set of incoming edges to } u$   $\text{adj}(u) = \text{defaults } \text{adj}^+(u)$  (outgoing)

$\deg^+(u) = |\text{adj}^+(u)|$

$\deg^-(u) = |\text{adj}^-(u)|$

how to represent in computers?  $\rightarrow$  adjacency matrix

	A	B	C	D	E
A	1	1	1	0	
B	1	0	1	0	
C					
D					
E					

$(u, v)$  index into adjacency matrix & find if there's an edge in  $\Theta(1)$

space: uses  $\Theta(|V|^2)$

\* matrix is useful when  $G$  is dense  $\rightarrow$  lots of edges with respect to  $|V|^2$   
not ideal when  $G$  is sparse  $\rightarrow$  not so many edges w/ respect to  $|V|^2$

IDEA: use a set with  $\Theta(1)$  lookup operations (DAA or Hash Table) to store adjacency lists

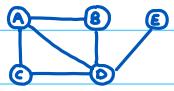
ADJACENCY LISTS: use DAA, hash table, linked list.

MATRIX:  $\Theta(|V|^2)$  space  
ADJ. SET:  $\Theta(|V| + |E|)$

$\sum_{v \in V} \deg^+(u) = |E| + \text{add all out degrees of vertices} \rightarrow \# \text{edges!}$

- set backed by direct array
- each elt has vertices it goes to

A	→ (B, C, D)
B	→ (A, D)
C	
D	
E	



**PATH:** sequence of vertices connected by edges

**EX:** path  $A \rightarrow B: (A, B), (A, C, D, A, B)$

fewest # edges  
needed to traverse  
from  $u$  to  $v$

$\cdot d(u, v)$ : distance from  $u$  to  $v$  following the shortest path

$\# d(u, v) = \infty$  by convention if there is no path from  $u$  to  $v$ .

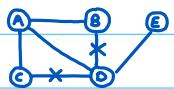
**single-pair-reachability( $G, s, t$ ):** true iff there's a path from  $s$  to  $t$  in  $G$ .

**single-pair-shortest-path( $G, s, t$ ):** return shortest path & its distance  $d(s, t)$  for  $s$  &  $t$  in  $G$ .

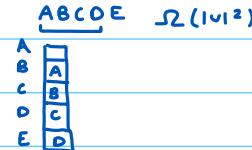
→ solving this problem solves reachability issue

**Single-Source-Shortest-Path( $G, s$ ):** return shortest path tree & set of distances for all vertices starting @  $s$ .

→ solving this solves single-pair-shortest-path

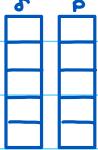


how to represent shortest path tree?



idea: store elt. before current

go in layers → **BFS:**



$L_0 = \{s\}$

```

for  $v \in L_i$ 
  for  $u \in \text{Adj}(v)$ 
    if  $d[u] == \infty$ :
       $L_{i+1} \text{ add } u$ 
       $d[u] = i+1$ 
       $P[u] = v$ 
  
```

$\sum_{u \in V} \deg^+(u) = |V|$

**COST OF BFS:**  $O(|V| + |E|)$  ← runtime

**PROOF OF CORRECTNESS:** induction  
(layers)



**DFS: single-source-reachability( $G, s$ ):**

```

def visit(P, v):
  for  $u \in \text{Adj}(v)$ 
    if  $u \notin P$ :
       $P[u] = v$ 
      visit(P, u)
  
```

$P$  = hash table

A	/
B	A
C	D
D	B
E	D

**RUNNING TIME:**  $O(|E|)$

\* improved R.T. @ sacrifice  
of shortest path

# can see if vertex is in graph  
in  $O(1)$ .

10/14

## LAST TIME:

- BFS time  $O(|V| + |E|)$ , linear in graph size
- DFS
- Single-source shortest path → solve w/ BFS
- Single-source reachability → solve w/ DFS

$$T(n) = O(n^2)$$

$$T'(n) = O(n^2)$$

**WRONG!**  
 $T(n) - T'(n) = O(n^2) - O(n^2) = 0 \times$   
**COUNTER:**  
 $T(n) = n^2$   
 $T'(n) = 7 \in O(n) \rightarrow T(n) - T'(n) = n^2 - 7 \neq 0$   
 don't subtract → usually add!

## (undirected graphs)

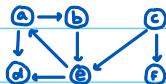
CONNECTED COMPONENTS: two vertices  $u, v$  are in the same CC if there is  $u-v$  path in graphhow to find connected components? input:  $G = (V, E)$ , output: CCs $S \leftarrow \{\}$  set holding CCs

TOTAL RUNTIME: linear (DFS/BFS only touch each node only once)

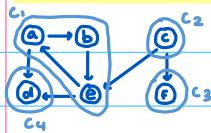
for each  $v \in V$  not already in a CC in  $S$ : sloppy:  $(O(|V| + |E|)) \cdot |V|$   
 $\nwarrow O(|V|)$ run  $C \leftarrow \text{SSReachability}(G, v)$  ← ex: a only returns b b/c only reachable node  
 $\nwarrow \text{DFS/BFS } O(|V| + |E|)$ add  $C$  to set  $S$ return  $S$ 

$$\begin{aligned} \text{time: } & O(|V|) + \sum_{\text{comps}} \left( O(\# \text{vertices in } C) + \# \text{edges in } C \right) \\ & = O(|V|) + O(|V|) + O(|E|) = O(|V| + |E|) \end{aligned}$$

## DIRECTED GRAPH:



## (for directed graphs)

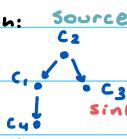
STRONGLY CONNECTED COMPONENTS (SCC): two vertices  $u, v$  are in the same SCC if∃  $u-v$  path AND  $v-u$  path in the (directed) graph

#AS soon as edges have direction, be careful of single vs. bidirectional

CONDENSATION GRAPH of directed graph:  $\xrightarrow{\text{Gcon}}$ 

$$V = \text{SCCS of } G = \{C_1, C_2, C_3, C_4\}$$

E = there is an edge  $(C_i, C_j) \in G_{\text{con}}$   
 if ∃ an edge from vertex in  $C_i$  to vertex in  $C_j$

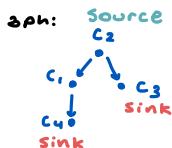


#draw edges going b/t CCs → careful of directions!

#nice way to summarize structure in directed graphs

(directed acyclic graph)

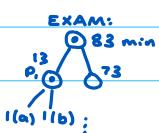
CLAIM: for every directed graph  $G$ , its condensation graph  $G_{\text{con}}$  is acyclic (no cycle) - DAGPF BY CONTRADICTION: assume ∃ cycle in  $G_{\text{con}}$ :  $C_1 \rightarrow C_2 \rightarrow \dots \rightarrow C_m$ all vertices in  $G_1, \dots, G_m$  are in the same SCC  $\Rightarrow C_1, C_2$  are in same SCC. contradiction.



IDEA:

- 3 Sink  $C_{sink}$  in  $G_{con}$
- BFS/DFS from vertex in  $C_{sink}$  finds entire  $C_{sink}$ .
- Q1: how to find vertex in  $C_{sink}$ ? Found  $C_{source}$ , but want  $C_{sink}$ . reverse arrows and see which finishing time is longer.
- Q2: what's next?

#Finishing time of root always > subprob.



$O(|V| + |E|)$

FULL DFS: run DFS from every unexplored vertex in  $G$  until all are explored

FINISHING TIME: of a vertex  $v$  in Full DFS run is time at which run has explore  $v$

and all its neighbors

KOSARAJU - SHAVIR: finding SCCs

CLAIM: if  $C$  &  $C'$  are SCCs and  $\exists$  edge  $C \rightarrow C'$

then if run full DFS from any vertex,

largest finishing time in  $C$  is bigger than largest finishing time in  $C'$

PROOF: two cases

1)  $C$  explored first



2)  $C'$  explored first

FINDING ALL VERTICES IN  $C_{source}$ :

FINDSOURCE( $G$ ):

- run full-DFS on  $G$ , recording finishing times
- return vertex w/ biggest finishing time.

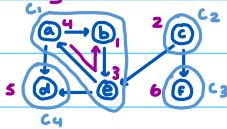
FIND ALL OF  $C_{sink}(G)$ :

- reverse all edges in  $G$  to get  $G_{rev}$
- $V \leftarrow$  FINDSOURCE( $G_{rev}$ )
- return SSReachability( $G, V$ )

FINDSCC( $G$ ):  $O(|V| + |E|)$

- let  $G_{rev}$  be  $G$  w/ all edges reversed
- run full DFS on  $G_{rev}$  recording finishing times  $f$
- run full DFS on  $G$  w/ vertices ordered by reverse finish time  $f$  (highest  $\rightarrow$  lowest)
- each time need a new tree, mark as new SCC

Finishing time:



10/16 TOPO SORT: (must be on acyclic, directed graphs)

TOPO ORDER on  $G = (V, E)$  is order  $\sigma$   $\forall$  edge  $(u, v)$   $\sigma(u) < \sigma(v)$

DFS, reverse finish order is topo order

C

$A \leq B$

REDUCTION: algorithm for transforming one problem to another

Ex: squaring mult.  $a^2 = a \cdot a$  (Turing & Coke Reduction)

Ex: mult. to squaring  $(a+b)^2 = a^2 + 2ab + b^2$

$$a \times b = \frac{(a+b)^2 - a^2 - b^2}{2} \quad (\text{Turing & Coke})$$

Ex: reachability to SSSP: see if  $d(s, v)$  is finite (Turing & Coke)

**TURING REDUCTION:**  $A \leq_T B$  iff  $\exists A$  that solves  $A$  while using  $B$  as atomic subroutine

- "oracle subroutine" bc soln to  $B$  is sometimes called oracle

**COK REDUCTION:** polynomial time Turing Reduction (polynomial upper bound)

THM:  $A \leq_T B$  in time  $f(n)$ ,  $B$  solveable in time  $g(n)$

$$f(n) \cdot g(n + f(n))$$

**MANY-ONE REDUCTION:**  $A \leq_m B$   
 $f_n : A \rightarrow B$

Ex: square to mult

$$a \mapsto (a, a)$$

$$\text{square}(a) = \text{mult}(a, a)$$

take input  $x$  for  $A$   
transform to  $y = f(x)$   
 $B(y)$  should be same  
answer as  $A(x)$

if reduction is poly time, called Karp Reduction

### REDUCTIONS & GRAPHS:



#### GRAPH DUPLICATION:

PROBLEM: given  $G = (V, E)$ , two vertices  $s, t \in V$

OUTPUT: shortest even-length  $s-t$  path

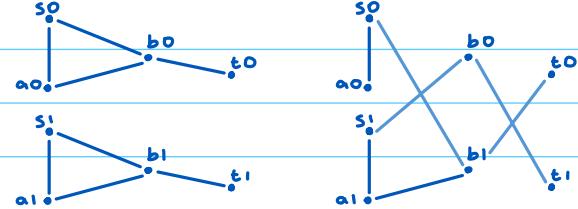
$$V' = V \times \{0, 1\} = \{(v, b) : v \in V, b = 0 \text{ or } 1\}$$

$$E' = \{(u, b), (v, 1-b) : (u, v) \in E, b \in \{0, 1\}\}$$

$$G' = (V', E')$$

$(G, s_t) \mapsto (G', s_0, t_0)$  reduces SPSP to SPSP

runtime:  $|V'| = 2|V|$  linear time to compute  $G'$   
 $|E'| = 2|E|$   $O(2|V| + 2|E|) = O(|V| + |E|)$  linear



#### SELF-REDUCTION

##### OUTPUT FORM:

v1: output  $d^*(s, -)$

v2: output P shortest path tree

$v1 \leq v2$  compute depth with pre-order traversal

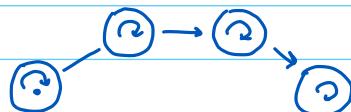
$v2 \leq v1$  for each vertex  $v$ ,  $P(v)$  is in-neighbor w/  $d^*(s, v) = 1 + d^*(s, P(v))$

**INPUT RESTRICTION:**  $G$  is semi-connected iff  $\forall u, v$  vertices, there is either  $u-v$  or  $v-u$  path

problem: given  $G$ , is  $G$  semiconnected?

problem 2:  $E \subseteq \subset$  see if all edges  $(v, v+1)$  exist ✓

problem 3:  $G$  acyclic reduce to P2 with DFS



10/21

## PRIORITY QUEUES

implementation:  
AVL

HEAP:

sequence:

build(A)  $O(n \log n)$ insert(x)  $O(\log n)$  $O(\log n)$ EITHER OR!  
not both.delete\_max()  $O(\log n)$  $O(\log n)$ delete\_min()  $O(\log n)$ 

if multiple, deletes one of them

 $O(n)$ def pq-sort(A):  $O(n \log n)$  sorting. NOT in place (have to build tree)B = []  $O(1)$ q = build(A)  $O(n \log n)$ 

for i = 0 ... |A|:

    B.append(q.delete\_max())  $O(\log n)$ 

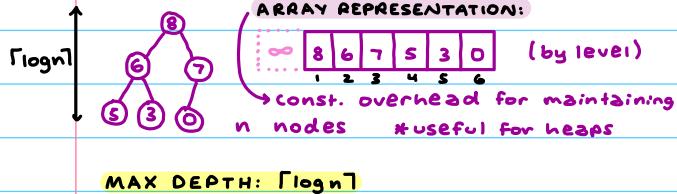
return B

AVL:

item
height
l r

binary  
COMPACT TREE: all nodes pushed as far left as can be

→ has 1:1 correspondence w/ array implementation



find children / parent:

$\text{LEFT}(x) = 2x \quad x = \text{parent}$

left child =  $2x$ 

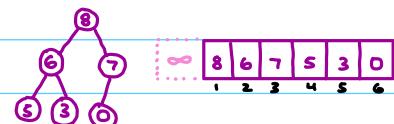
$\text{RIGHT}(x) = 2x + 1 \quad \text{right child} = 2x + 1$

$\text{PARENT}(x) = \lfloor \frac{x}{2} \rfloor$

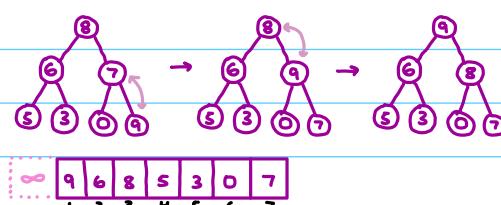
MAX HEAP PROPERTY: for  $x \in \text{Tree}$ ,  $x \geq \text{left}(x)$ ,  $x \geq \text{right}(x)$ 

· don't know relationship b/t L/R children, but know that parent is biggest

· MIN HEAP PROP: opposite (parent = min elt.)



insert(9):



def heapify-up(x):

$p = \frac{x}{2}$

if  $A[x] > A[p]$ :swap  $A[x]$  with  $A[p]$ 

heapify-up(p)

 $\Theta(\log n)$ 

delete\_max:

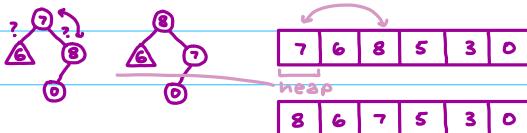
swap root w/ last node

heapify-down(1) ← start from root

return removed value



heapify-down: compare w/ children



case 1:  $A[x] \geq A[L(x)]$

$\geq A[R(x)]$

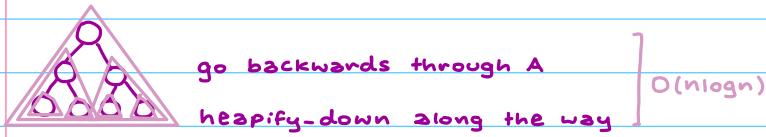
case 2:  $A[x] < A[L(x)]$

$\geq A[R(x)]$

swap w/  $L(x)$

case 3: both children  $> A[x]$

swap w/ larger child



at every level, # nodes =  $\frac{n}{2^{h+1}}$

$$\text{Work} = \sum_{h=0}^{\log n} \frac{1}{2^{h+1}} ch$$

converges on a constant

$$2\Theta(1)$$

$$= nc \sum_{h=0}^{\log n} \frac{h}{2^{h+1}} < nc \sum_{h=0}^{\infty} \frac{h}{2^h} < nc \sum_{h=0}^{\infty} \frac{\frac{1}{2}}{(1-\frac{1}{2})^2} = nc \frac{\frac{1}{2}}{(1-\frac{1}{2})^2} \rightarrow \text{Work} \in \Theta(n)$$

def heap-sort(A):

    build(A)      $\Theta(n)$

    for i = 1 ... |A|:

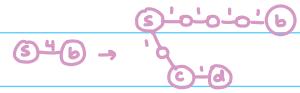
        decr. values → delete-max()      $\Theta(\log n)$

10/23

## SHORTEST PATHS: DIJKSTRA

Weighted graphs, dijkstra, analysis

SSSP: single source shortest paths



on undirected graph w/ equal weights: BFS  $\Theta(V+E)$

on undirected graph w/ non-equal weights (small positive ints): graph duplication  $\Theta(W(V+E))$

W: max weight of graph

if W is ugly #, have to go through & compare. ex: 1/2, 2.12, π, e

\*for now, positive non-zero #s.

$w(u,v)$ : weight of edge b/t u & v

$d(u,v)$ : shortest path weight b/t u & v

$d(s)=0$

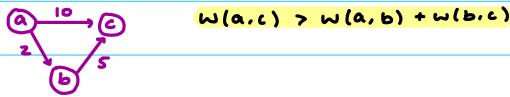
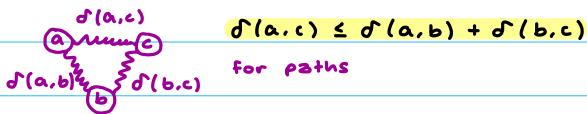
$d(u)=\infty \leftarrow u \notin S$

$d(s,v) = d(v)$

$d^*(u,v)$ : distance estimate of shortest path b/t u & v in V

$d(s,v) \equiv d(v)$

### TRIANGLE INEQUALITY: Num paths



### RELAXATION:



`relax(u, v):`

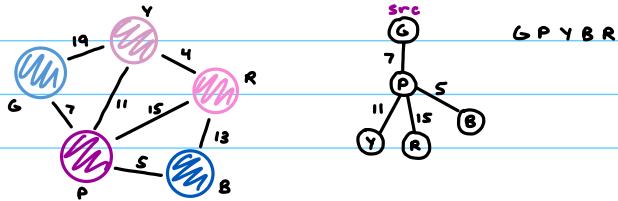
\*relaxation is safe

set  $d[v] = \min(d(v), d(u) + w(u,v))$

\*what order to relax in?

if  $d[v]$  changed, set `parent[v] = u`

order in terms of  $d$ : from G to others

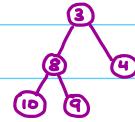


### DJIKSTRA:

- relax edges from each vertex in incrementing order of distance from source
- will use priority queue Q on vertices (unique IDs) S, U, V, etc. and keys correspond to  $d[v]$ 's

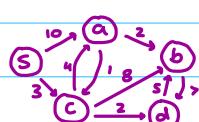
### ALGO:

- set  $d[s] = 0$ ,  $d[v] = \infty$  for  $v \in V$  in  $V$
- build priority queue  $Q$  for item  $(v, d(v))$ ,  $v \in V$   
while  $Q$  not empty
  - delete  $(v, d(v))$  from  $Q$ , with min  $d[u]$
  - for  $u \in \text{Adj}^+(v)$ : relax  $(u, v)$



decrease-key( $i$ , new-key):

- change val & heapify-up
- $\log(n)$  time heapify up



delete v from Q	S	A	B	C	D
S	0	$\infty$	$\infty$	$\infty$	$\infty$
C		10		3	
D		7	11		5
A		7	10		
B				9	

$$d(c) = 3$$

$$d(d) = 5$$

**CORRECTNESS:** @ end of Dijkstra,  $d(v) = d^*(v)$  for all  $v \in V$

- if relaxation sets  $d(v) = d^*(v)$ , then  $d(v) = d^*(v)$  on termination
- relaxation is safe & only decreases  $d(v)$

Show  $d(v) = d^*(v)$  when  $v$  is removed from  $Q$ .

**STRONG IND** of number  $k$  vertices removed from  $Q$ .

**BASE CASE:**  $s$   $d(s) = 0 = d^*(s)$

**IND. STEP:** we are removing  $v$  from  $Q$  and all vertices  $u$  removed from  $Q$ ,



assume  $d(u) = d^*(u)$

let  $\pi$  be a shortest path  $w(\pi) = d^*(v)$   
 $\xrightarrow{s} \xrightarrow{x} \xrightarrow{y} v$   $y$  is 1st vertex in  $\pi$  that is still in  $Q$ .  
 predecessor of  $y$   
 $x$  has been removed

I.H.  $d(x) = d^*(x)$

$(x, y)$  was relaxed when  $x$  was removed  $\Rightarrow d(y) \leq d^*(x) + w(x, y)$

subpaths of shortest path  $\pi$  are shortest paths  $\Rightarrow d(y) = d^*(y)$

$d^*(y) \leq d^*(v)$  non-negative edge weights

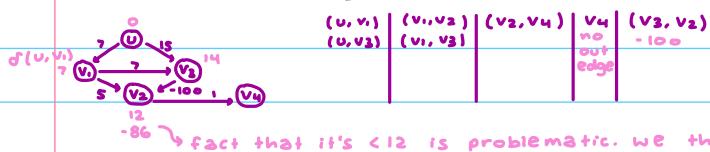
$\leq d(v)$  relaxation is safe

$\leq d(y)$   $v$  is vertex with min  $d(v)$  on  $Q$

$$d(v) = d^*(v) (= d(y))$$

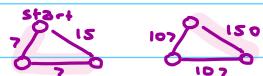
10/28 **BELLMAN FORD**

Dijkstra fails for negative edge weights



fact that it's < 12 is problematic. We thought we had the shortest path at that point.

→ if we run Dijkstra with +100 to every edge, discourages shortest path.

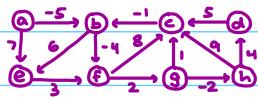


**NEGATIVE WEIGHT EDGES:**

$$d(a, b) = -\infty$$

$$d(a, f) = -\infty$$

path is under

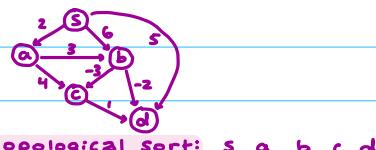


**NEGATIVE WEIGHT CYCLE:**  $b \rightarrow f \rightarrow g \rightarrow c \rightarrow b \rightarrow -2$  weight cycle

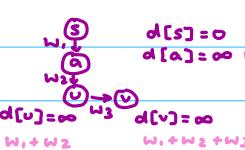


**DAG SHORTEST PATHS: negative-weight edges**

- no cycles. ∴ no negative cycles



topological sort:  $s \rightarrow a \rightarrow b \rightarrow c \rightarrow d$

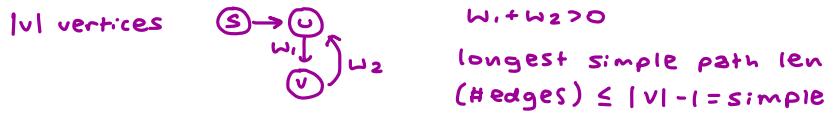


IF different order:

→ 1st: relax  $a \rightarrow u$

order of relaxation matters.

**CLAIM:** DAG topological sort computes shortest path even w/ negative weights



SIMPLE SHORTEST PATHS:

no neg. cycles

CLAIM: if  $d^*(s, v)$  is finite, there exists a shortest path to  $v$  that is simple.

could I have path w/ positive weight cycle that is shortest path?



K-EDGE DISTANCE  $d_k(s, v)$

is minimum weight path from  $s$  to  $v$  with  $\leq k$  edges

compute  $d^{|V|-1}(s, v)$ ,  $d^{|V|}(s, v)$

CLAIM 1: no negative cycles  $d^*(s, v) = d^{|V|-1}(s, v)$

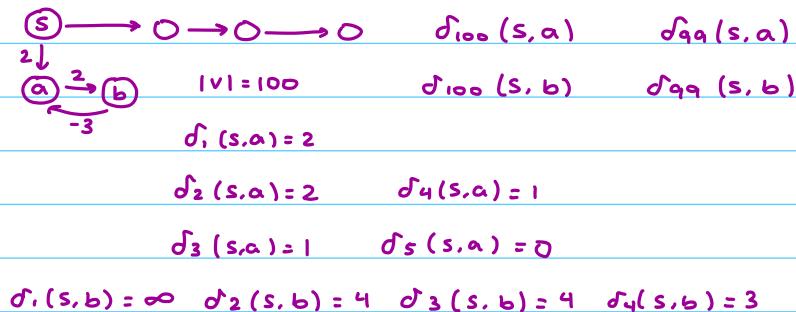
NEGATIVE CYCLE WITNESS:

if  $d^{|V|}(s, v) < d^{|V|-1}(s, v)$

shortest non-simple path to  $v$ , so  $d^*(s, v) = -\infty$

call  $v$  a negative cycle witness

CLAIM 2: every negative weight cycle  $C$  reachable from  $S$  contains a witness.

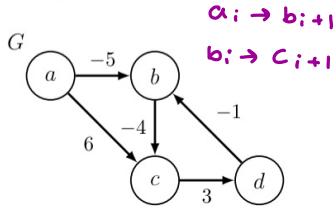


$\leq \sum_{v \in V} d_i(s, v) + \sum_{v \in V} w(v, v) < 0$

$\Rightarrow \sum_{v \in V} d_i(s, v) < \sum_{v \in V} d_{i+1}(s, v)$

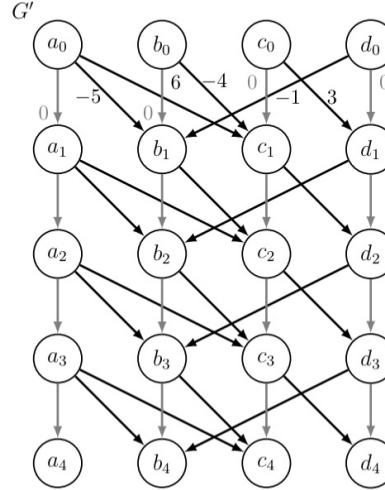
### BELLMAN-FORD:

#### Example



$\delta(a_0, v_k)$

$k \setminus v$	$a$	$b$	$c$	$d$
0	0	$\infty$	$\infty$	$\infty$
1	0	-5	6	$\infty$
2	0	-5	-9	9
3	0	-5	-9	-6
4	0	-7	-9	-6
$\delta(a, v)$	0	$-\infty$	$-\infty$	$-\infty$



Note that  $b$  is a witness, and  $c$  and  $d$  are not. The negative-weight cycle is  $bcd$ . Claim 2 simply says that there always exists a witness in the negative-weight cycle. Consider the case where there is an edge from  $a$  to a different vertex  $e$ , meaning that  $|V| = 5$ . In this case,  $\delta_5(a, b) = -7$ , so  $b$  is not a witness in this new graph, but  $\delta_5(a, c) = -11$ , and  $c$  is a witness.

produce a duplicated DAG level by level

Compute  $k$  edge dist. for each  $k$

$|V|$  passes of relaxation over all edges  $|E|$        $|V|-1$  passes → S.P. wts.

$O(|V||E|)$  alg.

one more pass → witness

10/30



ALL-PAIRS SHORTEST PATHS:

- Definition
- Johnson's Alg
- All-pairs reliability

### THE PROBLEM: APSP

INPUT: directed graph  $G = (V, E)$  weights  $w: E \rightarrow \mathbb{Z}$

OUTPUT:  $\delta: |V|^2 \rightarrow \mathbb{Z} \dots \delta(u, v)$  length of shortest  $u-v$  path  
 OR "FAIL" if  $\exists$  neg. weight cycle in  $G$ .

output len  $\Omega(|V|^2)$

GRAPH:	WEIGHTS:	ASAP ALG:	TIME O(·)
any	non-neg.	$ V  \cdot$ Dijkstra	$ V ^2 \log  V  +  V   E $
any	any	Johnson	$ V ^2 \log  V  +  V   E $
any	any	Floyd-Warshall	$ V ^3 / 2^{\Omega(\sqrt{\log  V })} =  V ^{3-o(1)}$

OPEN QUESTION:  $\exists$  alg for APSP that runs in time  $O(|V|^{2.99})$ ?

### why Study APSP?

- useful
- equivalent to many other problems
  - finding neg.-weight triangle
  - min-weight cycle
  - replacement paths problem

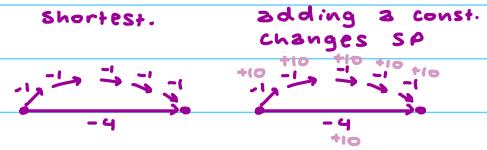
### 'MOST-BEAUTIFUL PATH PROBLEM'

$\{-10, \dots, 10\} \rightarrow$  use Johnsons.

### JOHNSON'S ALG:

PLAN:

- 1) find negative cycle, if exists. FAILS if so.
  - Bellman Ford,  $|V| \cdot |E|$
- 2) make all edge weights non-negative  $\leftarrow \star$ 
  - reweighting must preserve shortest paths
- 3) use Dijkstra from every vertex... clean up

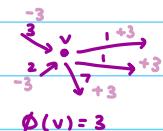


### 2 REWEIGHTING that preserves shortest paths:

let  $\phi: V \rightarrow \mathbb{Z}$

key idea: for every  $v \in V$ : add  $\phi(v)$  to all outgoing edges from  $v$

Subtract  $\phi(v)$  from all incoming edges from  $v$



CLAIM: this transformation changes weight of all  $uv$  paths by  $\phi(u) - \phi(v)$

$$\text{PF: weight of } \pi = (v_0, v_1, \dots, v_k) \text{ is } \sum_{i=1}^k w(v_{i-1}, v_i) = \sum_{i=1}^k [w(v_{i-1}, v_i) + \phi(v_{i-1}) - \phi(v_i)] \\ = (\sum w(v_{i-1}, v_i)) + \phi(v_0) - \phi(v_k)$$

### COMPUTING THE REWEIGHTING:

$\exists \phi(\cdot)$  s.t. every weight  $> 0$

let  $d_s(v)$  be weight of shortest  $s-v$  path

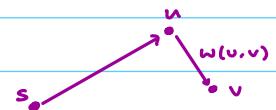
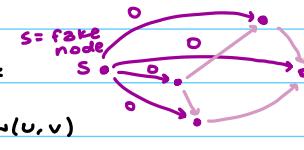
define  $\phi(v) = d_s(v)$

CLAIM: after reweighting by  $\phi(v) = d_s(v)$ , all edge weight  $\geq 0$

PF: by triangle inequality,  $d_s(v) \leq d_s(u) + w(u, v)$

$$\phi(v) \leq \phi(u) + w(u, v)$$

$$0 \leq \phi(u) - \phi(v) + w(u, v)$$



JOHNSON'S ALG:  $G = (V, E)$   $w: E \rightarrow \mathbb{Z}$

RT:

- add vertex  $s$ , connected to all  $v \in V$  w/ weight zero edge  $|V| \cdot |E|$
- use Bellman Ford to look for negative cycles. FAIL if found.
- compute  $d^*(v)$  from Bellman Ford. reweight  $(u, v) \mapsto w(u, v) + \phi(u) - \phi(v)$   $|E|$
- use Dijkstras ABSP to find  $d^*(u, v) \quad \forall u, v \in V$   $(|V| \log |V| + |E|) |V|$
- $\forall u, v \in V \quad d(u, v) = d^*(u, v) - \phi(u) + \phi(v)$  ← search every pair  $|V|^2$

### ALL PAIRS REACHABILITY:

→ could solve w/ BFS from every vertex

INPUT: directed graph  $G$  unweighted

$$|V| \cdot (|V| + |E|) = |V|^2 + |V||E| \sim |V|^3$$

OUTPUT:  $\forall u, v \in V$  "reachable" or not

$$\text{will show alg } |V|^W = |V|^{2.372}$$

multiply  $n$  by  $n$  matrices:

naive:  $O(n^3)$

clever:  $O(n^W) = O(n^{2.372..})$

### ALG (Munro)

- []
- let  $A$  be adjacency matrix of  $G$  w/ 1s on diagonal
  - use repeated squaring to compute  $A, A^2, A^4, A^8, \dots, A^{|V|} \log(|V|)$
  - $\exists (u, v)$  path in  $G$  iff  $A_{u,v}^{|V|} > 0$



$$A = \begin{pmatrix} a & b & c \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{pmatrix} \quad O(|V|^2)$$

Space & time

$$A^2 = \begin{pmatrix} 1 & 2 & 1 \\ 0 & 1 & 2 \\ 0 & 0 & 1 \end{pmatrix} \rightarrow \text{tells you how many paths there are from vertex to vertex}$$

RUNTIME:

squaring matrix mult

$$\log |V| + |V|^2 \rightarrow O(|V|^W - \log |V|) \quad W = \text{how fast we can do matrix mult.}$$

CLAIM: for all  $\ell \geq 1$ ,  $A_{u,v}^\ell$  counts # of paths from  $u \rightarrow v$  length  $\ell$ .

PF: by induction on  $\ell$ :

BC  $\ell=1$ : by construction

IND STEP: assume for  $\ell-1$ , the # of length  $\ell$   $u \rightarrow v$  paths is:

$$\sum_{w \in V} (\# \text{ of length } (\ell-1) \text{ paths } u \rightarrow w) \begin{pmatrix} 1 & \text{if } (w, v) \in E \\ 0 & \text{otherwise} \end{pmatrix}$$

$$= \sum_{w \in V} A_{u,w}^{\ell-1} \cdot A_{w,v} = A_{u,v}^\ell \quad \text{matrix mult.}$$

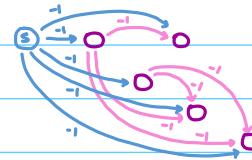
TEST 2



## 11/4 GREEDY ALGORITHMS

ACTIVITY SCHEDULING:  $A = \{[s, f]\}$ goal: create set  $S$  containing activities from  $A$  s.t. no overlapand  $S$  is maximal size. (most # of activities)input:  $\{[1, 8], [2, 4], [3, 6], [4, 7], [7, 8]\}$ output:  $\{[2, 4], [4, 7], [7, 8]\}$ let  $A' = A$  without  $S$ , or any other activities that overlap with it. $S \times \{S\}$  is optimal in  $A'$ 

• turn every activity into vertex:

• add an edge from every activity  $a$  to every activity that starts after  $a$  finishes

• add supernode connecting to all tasks.

cycles impossible!

• DAG SP on weights -1. (finding max path)

• find SP ("longest", most negative path)

- covers most edges

 $\Omega(n^2) \rightarrow$  we can do better!

## GREEDY CHOICE PROPERTY approach:

let  $a \in A$  be the activity with earliest finish time. $\exists S$ , optimal solution that contains  $a$ .PROOF: let  $S$  be optimal  $a \notin S$ .

$S = \{s_1, s_2, \dots, s_k\}$

swap  $a$  for  $s_i$  to get  $S'$ 

$S' = \{a, s_2, s_3, \dots, s_k\}$

 $|S'| = |S| \therefore$  must be an optimal solution.  $a = \text{greedy choice}$ 

GREEDY CHOICE PROPERTY: there's at least one optimal solution that contains greedy choice

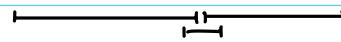
counterex:

interval w/ earliest start time:

$[1, 5], [2, 3], [3, 4]$



shortest interval counterex:



another greedy choice:

latest start time

### GREEDY ALGORITHM:

- 1) sort by finish time
- 2) make greedy choice
- 3) walk the list until non-overlapping activity and add it
- 4) go to 2

PF: claim: greedy alg produces optimal solution.

\* use induction

let  $S$  be an optimal solution.  $S$  starts with greedy choice  $g$ .

\* reduce greedy to a

let GREEDY be soln produced by Greedy algorithm.

Subproblem of itself

let  $A'$  be activities in  $A$  that don't overlap w/  $g$ .

\* can show that its a

GREEDY \ { $g$ } is optimal for  $A'$

Subproblem

$$|GREEDY| - 1 \geq |S| - 1$$

$$\therefore |GREEDY| \geq |S|$$

### ACTIVITY SCHEDULING

INPUT:  $A = \{(t, d)\}$

OUTPUT: schedule,  $S$ , that minimizes the max late time.

$l(a)$  = how late, after due time,  $a$  finishes.  $\min = 0$

$$l(A) = \max(l(a_i) \forall a_i \in A)$$

$$A = \{(1, 1), (2, 2), (3, 3)\}$$

$$S_1 = \overbrace{1}^{t=0} \ a_1 \ \overbrace{2}^3 \ a_2 \ \overbrace{3}^6 \ a_3$$

$a_1$  ends @ 1. supposed to end there.  $a_3$  supposed to end @ 3, but ends @ 6.

$$l(S_1) = 3 \leftarrow 1 + \text{est } (6-3)$$

$$S_2 = \overbrace{1}^{t=0} \ a_3 \ \overbrace{3}^1 \ a_2 \ \overbrace{5}^6 \ a_1$$

$$l(S_2) = 5$$

observation 1: there exists an optimal schedule with no idle time.

observation 2: let  $S$  be optimal for  $A$ . let  $a_i$  be scheduled first in  $S$ .

if remove  $a_i$ , then  $S \setminus \{a_i\}$  is optimal for  $A \setminus \{a_i\}$

observation 3 GCP:  $\exists$  optimal schedule where the 1st activity due is scheduled first.

PF 3: let  $S$  be optimal schedule where  $a$ , activity with earliest due date, is not first

let  $b$  directly precede  $a$  in  $S$ .  $S = [b|a]$  for  $a$ ,  $l(a)$  in  $S'$  improved vs  $l(a)$  in  $S$ .

let  $S' = S$  but with  $a, b$  swapped  $S' = [a|b]$   $l(b)$  might have gotten worse.

$l(b)$  in  $S' \leq l(a)$  in  $S$ .

### GREEDY ALG:

1) sort by due date  $O(n \log n)$

2) make greedy choice

3) go to 2

CLAIM: greedy alg produces optimal strategy, GREEDY.

PF: let  $S$  be optimal in  $A$ .  $S$  begins with  $g_1$ .

let  $t_1$  be duration of  $g_1$ .

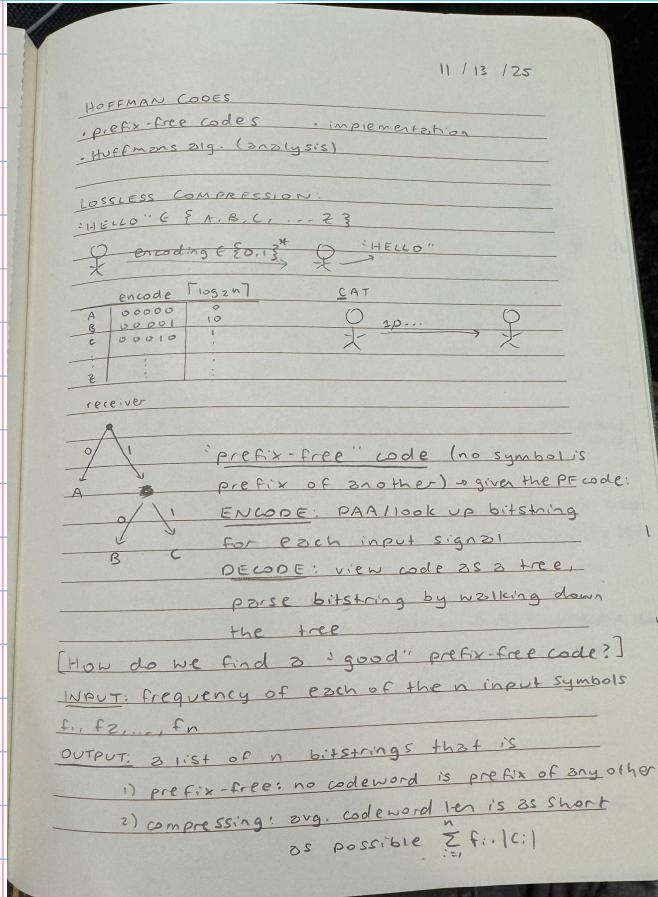
GREEDY \  $g_1$  is optimal in  $A \setminus g_1$ .

$l(\text{GREEDY} \setminus g_1) \leq l(S \setminus g_1)$  in  $A \setminus g_1$ .

adding  $g_1$  adds lateness of  $t_1$  to all activities evenly

$\therefore l(\text{GREEDY}) \leq l(S)$

11/13



A	00
B	010
C	10
D	1011

[TREE VIEW OF PREFIX - [FREE CODE.]

$$f_A = 5, f_B = 1, f_C = 2$$

• weights of nodes = sum of children

$$\text{define cost}(T) = \sum_{\text{nodes } x} f_x$$

$$\text{claim: cost}(T) = \sum_{i=1}^n f_i + c_i$$

• each value  $f_i$  in the tree

is counted  $l(i)$  times - once per level in tree until hit leaf  $i$ .

goal: zig for

INPUT: frequencies  $f_1, f_2, f_3, \dots, f_n$  | Huffman's  
Fits this

OUT: prefix-code tree of min cost

Observation 1: any optimal tree is a full bin. tree

Observation 2: in optimal tree,  $f_i \leq f_j$ , then  $l(j) \leq l(i)$

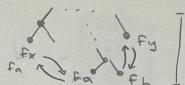
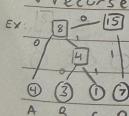
ALGORITHM: input:  $f_1, f_2, \dots, f_n$

• start by giving each frequency a leaf node  $\rightarrow n$  leaves

• find "smallest" nodes  $f_a, f_b$  & replace them w/ new

internal node  $w_l$  weight  $f_a + f_b$

• recurse.



ANALYSIS: claim: GCP: there is an optimal tree that has 2 smallest freqs. as siblings  
PF: start w/ optimal tree, transform it to one w/ GCP. let  $f_x, f_y$  be two smallest freqs  $f_x \leq f_y$   
let  $f_x$  be freq. deepest in tree,  $f_b$  be sibling by  $O_2$ . can swap  $f_x, f_b$   
know  $f_y \leq f_b$ ,  $f_y$  is 2nd smallest  $\Rightarrow$  By  $O_2$ , swap  $f_y \leftrightarrow f_b$

CLAIM: Huffman's zig. outputs min cost tree

PF: by induction on  $n$ , # of frequencies

BC:  $n=2 \rightarrow \frac{f_1 + f_2}{2} \downarrow$

IND STEP: let  $H$  = tree that HA outputs

• assume that Huffman's is optimal for  $n-1$  freq/symbols.

• let  $f_x, f_y$  be smallest freqs

• in  $H$ ,  $f_x$  &  $f_y$  are siblings

• by GCP,  $\exists$  an optimal tree  $T$  w/  $f_x, f_y$  siblings

• let  $T' \wedge H'$  are  $T, H$  but w/ leaves  $f_x, f_y$  removed

• imagine running HA on freq + table of  $n-1$  freqs

w/ new symbol  $z$  replacing  $x, y$  s.t.  $f_z = f_x + f_y$

• outputs exactly  $H'$ : let  $T'$  be any other tree

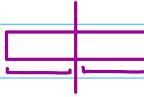
on  $n-1$  symbols, by Ind. hyp.  $\text{cost}(H') \leq \text{cost}(T')$

• by defn of cost of tree:  $\text{cost}(H) = \text{cost}(H') + f_x + f_y$

• by defn of cost of tree:  $\text{cost}(T) = \text{cost}(T') + f_x + f_y$

$\text{cost}(H') \leq \text{cost}(T')$ ,  $\text{cost}(H) \leq \text{cost}(T)$   $\square$

11/18 DYNAMIC PROGRAMMING



- generalizes recursion
  - divide & conquer, except overlapping problems

### COIN ROW PROBLEM:

ARRAY A of size n: represents a row of coins w/ positive values

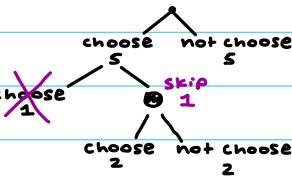
**GOAL:** pick max amount subject to: no two adjacent coins can be picked.

$$A = [c_0, c_1, \dots, c_{n-1}]$$

$$A = \begin{matrix} 5 & 1 & 2 & 10 & 6 & 2 \end{matrix} \rightarrow \max = 17$$

## RECURSIVE / EXHAUSTIVE SEARCH:

- 1) decide to not pick  $c_0$ .  
use recursion to find optimal solution  
 $[c_1, \dots, c_{n-1}]$
  - 2) decide to pick  $c_0$ . cannot choose  $c_0$ .  
use recursion on  $[c_2, \dots, c_{n-1}]$



## IMPLEMENTATION

Coins(A):

Base Case: if  $|A| < 2$ : output  $\sum A$

Else output  $\max(\text{coins}(A[1:]), A[0] + \text{coins}(A[2:]))$

#cells is exponential in n

$$T(n) = T(n-1) + T(n-2) + \Theta(1)$$

$T(n) = \sum (2^{n/2})$   $\approx$  slow - reduce the work! (there's many of the same branch that we are redoing)

## MEMOIZATION

`def Coins(A): linear time!`

`memo = [None] * len(A)`

\*only going to  
solve each suffix (end part)  
problem ONCE

```
def sum_values(i):
```

```
if i > len(A): return 0
```

if ~~more~~ [i] == None:

```
without i = coins-helper(i+1)
```

With  $i = \text{A}[i] + \text{coins\_helper}(i+2)$

`memo[i] = max (with i, without i)`

return memo[i];

```
return coin-helper[0]
```

### PURE DP SOLUTION O(n) runtime

`memo[i]` depends on `memo[i+1]` & `memo[i+2]`

```
def Coins(A):
```

```
    memo = [0] * (len(A) + 2)
```

```
    for i in range(len(A)-1, -1, -1):
```

```
        memo[i] = max(memo[i+1], A[i] + memo[i+2])
```

```
    return memo[0]
```

### SRT BOT:

S: subproblems



R: relationship b/t subproblems

T: topsort (must be acyclic) - dependence b/t subproblems = DAG

B: identify base case

D: show how output depends on entries of M.

T: analyze runtime.

### Coin Row - SRT BOT:

S:  $M(i) = \max$  value of coins we can pick from  $A[i:]$   $i: 0$  to  $n-1$

$\text{memo}[5] = 2$

$\text{memo}[4] = 6$

$\text{memo}[3] = 12$

:

$\text{memo}[8] = 17$

R:  $M[i] = \max \dots i=0, n-1$

T: each  $M[i]$  depends on  $M[i']$ ,  $i' > i$

B:  $M[i] = 0$ ,  $i \geq n$

D: output is  $M[0]$

T:  $n$  values  $M(i)$ , each takes  $O(1)$ ,  $\therefore O(n)$  runtime

`Sol = [] report coin choices for M[0]`

`while i < n:`

`if memo[i] == memo[i+1]:`

`i += 1 ← i skipped`

`else:`

`Sol.append(A[i])`

`i += 2`

`return Sol ← i chosen`

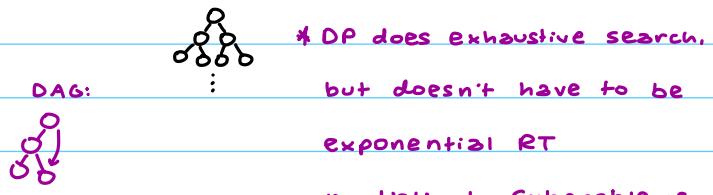
BOWLING PINS  $B = [1, -2, -4, \underline{-3}, 10]$  23

drop a pin, earn point value

drop 2 consecutive pins, product of values maximize score

1st step: choices (what they are)

### exhaustive search:



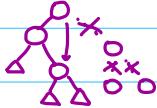
\* DP does exhaustive search.

but doesn't have to be

exponential RT

n distinct subproblems

11/20



### LONGEST COMMON SUBSEQUENCE:

input: A, B strings

A: GOOZILLA

out: m, len of LCS

B: GHIDORAH

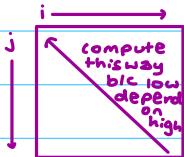
GDA ] out m=3  
GIA

if 1st char same → if  $A[0] == B[0]$ :  $1 + \text{LCS}(A[1:], B[1:])$

always  $\text{LCS}(A[1:], B)$

always  $\text{LCS}(A, B[1:])$

] return max of these 3 → this is exponential. we want to come up w/ memo so we don't recompute



if  $M(i, j)$  is None:  
 $M(i, j) = \max \left\{ \begin{array}{l} \text{help}(i+1, j+1) + 1 \text{ if } A[i] == B[j] \\ \text{help}(i+1, j) \\ \text{help}(i, j+1) \end{array} \right\}$   
 return  $M(i, j)$

RUNTIME:  $O(|A| \cdot |B|)$  quadratic!

S:  $M(i, j)$  is length of LCS of  $A[i:j], B[j:]$  where  $0 \leq i \leq |A|, 0 \leq j \leq |B|$   
 $(1 + M(i+1, j+1))$  if  $A[i] == B[j]$

R:  $M(i, j) = \max \sum M(i, j)$

T: increasing  $i+j$  (never read uninitialized indices)

B: 0, where  $j > |B|$  or  $i > |A|$

D:  $M(0, 0)$

\* can reconstruct LCS by walking diagonally down the table

T:  $O(|A| \cdot |B|) \leftarrow \text{size of table}$

### OSSP: correctness:

1:  $A[:], B[:]$  if we have optimal LCS, then  $1 + \text{LCS}(A[:], B[:])$  is optimal w/  $A[0] == B[0]$

2:  $A[:], B[:]$  if we have optimal LCS, then optimal among solns that didn't pick  $A[0]$

3:  $A, B[:]$ : if we have optimal LCS, then optimal among solns that didn't pick  $B[0]$

### LONGEST INCREASING SUBSEQUENCE:

input: array A of integers

EX:  $A = [8, \underline{6}, \underline{7}, 5, 3, 0, \underline{9}]$

\* for each val, choose or not

out: length LIS

out: 6, 7, 9 → 3 (len)

PROPAGATE w!

LIS(A, w=(-∞));

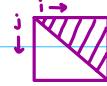
if A is empty: return 0  
 $\left\{ \begin{array}{l} 1 + \text{LIS}(A[:], A[0]) \text{ if } A[0] > w \\ \text{LIS}(A[:], w) \end{array} \right.$

indexify A → i

indexify w → j (define  $A[-1] = -\infty$ )

S:  $M(i, j)$ : length of LIS of  $A[i:j]$  w/ all numbers  $> A[j]$

R:  $M(i, j) = \max \begin{cases} 1 + M(i+1, j) \\ M(i+1, j) \end{cases}$



\* don't care about w itself.

T: increasing  $i$  (start from right end)

B:  $M(|A|, j) \forall j$   
 $\downarrow = 0$

O:  $M(0, -1)$

T:  $O(|A|^2)$  \* have space & time optimizations

S:  $M(i)$ , len of LIS of  $A[i:]$  that begins w/  $A[i]$ ,  $0 \leq i < |A|$

R:  $M(i) = 1 + \max \{ M(j) \text{ where } j > i \text{ and } A[j] > A[i] \}$



T: increasing  $i$  (only looking @ suffixes)

B: no base case! (there is no  $j > i$ )

O: max in M

set  $A[k] \rightarrow M[k]$

What if we used a range query (implement using AVL BST)  
augment w/ max

T:  $O(|A|^2)$

T: implement  $M$  using AVL BST augmented w/ max.

S maps  $A[k] \rightarrow M[k]$  in  $O(\log |A|)$  time

walk back along input:  $|A|$  steps

$O(|A| \log |A|)$

## 11/25 DYNAMIC PROGRAMMING III:

Shortest paths using DP, DP using SP algs

Subpaths of SPs are SPs

"cut and paste" arg



## DAG Single-Source SP:

S:  $T(v) = d^*(s, v)$  for all  $v \in V$

$$\xrightarrow{w(u,v)} \text{if } w(u,v) < \infty$$

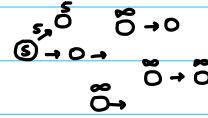
R:  $\min \{ T(u) + w(u,v) \text{ for } u \in \text{Adj}^-(v) \}$

$$\xrightarrow{w(u,v)}$$

T: topological sort/order

SSSP for General Graphs: Bellman Ford

→ need to handle cycles & negative cycles



compute shortest paths that are restricted

to a single edge



$\leq k$ -edge shortest paths

$|V|-1$  edges → longest simple path

$$K = |V| - 1$$

$k+1$  iterations

compare  $v$ th iteration (find witnesses)

\*remember how many cycles used so far as we build paths

S:  $T(v, k) = \text{weight of SP from } s \text{ to } v \text{ using at most } k \text{ edges}$  for  $v \in V$ ,

weight  $\leq$  most  $k$  edges for  $v \in V$ ,  $k \in \{0, \dots, \min(|V|, |E|+1)\}$

↳ implies cycle

R: for every  $v \neq k$ , guess the last edge used in path of length  $k$ , or that it's best to

use  $\leq k-1$  edges

$$T(v, k) = \min \left\{ \begin{array}{l} \min \{ T(u, k-1) + w(u, v) \mid u \in \text{Adj}^-(v) \} \\ T(v, k-1) \end{array} \right.$$



T: depend on entries w/ smaller  $k$

B:  $T(s, 0) = 0$ ,  $T(v, 0) = \infty$  for  $v \neq s$

O: Simple path can have  $\leq$  most  $|V|-1$  edges

if for some  $u$  we have  $T(u, |V|) < T(u, |V|-1)$ , we have witness

for a neg. cycle. otherwise,  $d^*(s, u) = T(u, |V|-1) = T(u, |V|)$

T:  $O(|V| \cdot |E|)$

## All Pairs Shortest Paths

dense graphs:  $E = O(|V|^2)$

Bellman Ford DP:  $O(|V|^2 |E|)$

BF:  $O(|V|^4)$

Johnson's:  $O(|V|^2 \log |V| + |V| |E|)$

GOAL: make dense graph  $O(|V|^3)$

### FLOYD-WARSHALL DP

→ restrict paths to be within a subset of  $V$ . ignore negative cycles



$v_i$

$v_j$

vertices are all numbered 0 through  $|V|-1$ . 0 intermediate

$$w(v, v) = 0 \text{ for all } v$$

$$w(u, v) = \infty \text{ if } (u, v) \notin E$$

$$\text{ex: } \{0, 1\} \leftarrow v_0, v_1$$



$$\text{ex: } \{0\} \rightarrow \{1\} \text{ no intermediate vertices } \{\}$$



$\{0\} \rightarrow \{1\} \rightarrow \{0\}$

$$S: T(k, u, v) = \min \text{ weight path from } u \text{ to } v \text{ through vertices in } \{0, 1, \dots, k\}$$

if  $k=0 \rightarrow$  no intermediate vertices allowed. single edge if it exists

$$B: T(0, u, v) = w(u, v) \text{ # or } \infty \text{ for all } u \neq v$$



$$T(0, u, v) = \min \{ T(0, \dots)$$

$$\begin{aligned} & \xrightarrow{0, 1, \dots, k-1} \xrightarrow{0, 1, \dots, k-1} \\ & \xrightarrow{\text{SP}} \xrightarrow{\text{SP}} \\ T(k, u, v) &= \min \{ T(k-1, u, k) + T(k-1, k, v) \\ & T(k-1, u, v) \} \\ (|V|-1)|V| \cdot |V| &= O(|V|^3) \end{aligned}$$

for no neg. cycles

### 12/2 DP IV. Subset sum

#### SANDWICH CUTTING:

- $V(L)$  maps to value  $\{1, 2, \dots, L\}$

- sandwich of len  $L$

- task: divide up so as to maximize total value

$$\text{ex: } 1, 10, 13, 18, 20, 31, 32$$

$$\begin{array}{ccccccc} L & 1 & 2 & 3 & 4 & 5 & 6 & 7 \\ & \underline{+} & & & & & & \\ & 12 & & & & & & \end{array}$$

$$10+10+13=33 \quad 31+1=32$$

what is 1st choice?

- pick where to make 1st cut      # length is changing → ∴ parameterize by length

- if we cut  $k$  inches,  $v(k) + sc(v, L-k)$       input:  $(L, v)$  size  $L+1$   
represented as  $v(1), v(2), \dots, v(L)$

$$S: sc(l) = \max \text{ value from cutting length } l \quad 0 \leq l \leq L$$

$$R: sc(l) = \max_{1 \leq k \leq l} v(k) + sc(l-k)$$

T: decreasing  $l$  (depends on  $l' < l$ )

$$B: sc(0) = 0$$

$$O: sc(L)$$

$$T: O(L^2) \quad L \text{ SPs, } O(L) \text{ each}$$

(quadratic in input size)

POLY TIME: bounded by polynomial in input size

PSEUDO-POLY: bounded by polynomial in input values

### SUBSET SUM:

positive

• input:  $A = [a_1, a_2, \dots, a_n]$ , len  $L$

• output: YES iff is a subset of  $A$  with sum  $L$

EX:  $A = [2, 5, 7, 8, 9]$ ,  $L=21$ ? YES ( $5, 7, 9$ )

$L=25$ ? NO

DECISION PROBLEM: output is binary (T/F, Y/N, etc)

• do we include  $a_n$ ?  $A' = [a_1, \dots, a_{n-1}]$   $L' = L - a_n$  o/w  $L' = L$



S:  $SS(i, L') = \text{YES}$  iff there is a subset of  $A[1, \dots, i]$  with sum  $L'$   $0 \leq i \leq n$ ,  $0 \leq L' \leq L$

R:  $SS(i, L') = \text{OR} \left\{ \begin{array}{l} SS(i-1, L') \\ SS(i-1, L'-a_i) \end{array} \right\}$  ← include only if valid

T: decreasing  $i$

B:  $SS(0, 0) = \text{yes}$

$SS(0, L') = \text{no}$  if  $L' \neq 0$

$SS(i, L') = \text{no}$  if  $L' < 0$

O:  $SS(n, L)$

T:  $O(nL)$  SPs,  $O(1)$  each,  $O(nL)$  total ( $2^n$  reachable SPs)

CAN WE DO BETTER?? nobody knows!!

if so, we'd get efficient algs for all kinds of problems, e.g.



PARTITION:

• input:  $[a_1, a_2, \dots, a_n]$

• output: YES if exists a bipartition of  $A$  with equal sums

REDUCE TO SS:  $A \mapsto (A, \frac{1}{2}\sum A)$  pseudopoly w/ SS algo.

$(A, L) \mapsto A + [L+1, \sum A - L+1]$  goes into diff halves of partition

NEGATIVE SUBSET SUM: allowed neg. vals in  $A$  as well

EX:  $A = [-1, 4, 7, 0, 0, 0]$   $L=3$  ← fill in remainder w/ 0's if needed

1, 6, 9, 2, 2, 2  $L'=9$

$IM-1, IM+4, \dots, IM$   $L=3M+3$

pseudopoly!  $M = \text{poly input.}$

12/4 last lecture:

- pseudopolynomial time
- equiv. subset sum, neg. subset sum, partition

today:

- decision problems, complexity, classes, P vs. NP

NP-hardness      Better than  $O(nL)$

what about longest simple path?

computational vs. relative complexity

compare APSP vs. longest path vs. subset sum

Sorting	$O(n!)$	$O(n \log n)$
SCC	$O(n!)$	$O(n+m)$
Subset Sum	$O(2^n)$	$O(nL)$

DECISION: given some input, does it have some property

all of these  
can be  
reduced  
to this

SEARCH: given input, which defines a search space, output a valid element

OPTIMIZATION: given input, constraint, cost function, output a soln which minimizes/maximizes cost

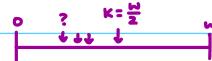
FUNCTION: given input, compute some output (math, counting)

MIN WEIGHT S-T PATH:

INPUT: graph G, nodes s-t

OUTPUT: weight of minimum s-t path

DECISION: Graph G, s-t, int k → does there exist an s-t path w/ weight ≤ k



exponential time  $O(2^{\text{poly}(n)})$

COMPLEXITY CLASSES: P, NP, EXPTIME

A: most natural have small exponents.  $O(n^2)$  APSP

P: easy & efficient: solvable in  $O(\text{poly}(n))$

B: Strong Church Theory Thesis: 2nd SP

$O(n)$  ✓  $O(n^2)$  ✓  $O(n^{6/210})$  ✓  $O(n \log \log n)$  X

SUBSET + LONGEST PATH

- easy to verify

2 VERIFIER  $v(x, c)$ , x is instance, c is a witness

- if x is YES, there exists c s.t.  $v(x, c) = T$
- if x is NO, for all c,  $v(x, c) = F$

V must run in poly time

$|c| = \text{poly}(n)$

EX: SUBSET SUM

C = list of numbers

VERIFY V = sum up the list. check  $v \leq L$

LONGEST PATH

SEARCH: does there exist a C s.t.  $V(x, c) = \text{True}$

C = a path

usually binary strings

V = is C a valid path? is  $|c| \geq k?$

$P \subseteq NP$        $NP \subseteq EXPTIME$

$C = \epsilon$       number of c's =  $O(2^{\text{poly}(n)})$

V = algo. loop over all c if  $V(x, c) = \text{True} \rightarrow \text{accept/return yes}$

$O(2^{\text{poly}(n)}) \cdot O(\text{poly}(n))$   
much bigger v

$P \subseteq NP \subseteq \text{Exp. Time}$

↑  
at least one is strict

We believe  $P \neq NP$

REDUCTIONS:

$A \leq B$  means A reduces to B,

$A \leq_p B$

- B is as hard as A,

- poly time algo. for B solves A in poly time

- algorithm for B solves A

- if  $B \leq P$ , then  $A \leq P$

NP-HARD: for all  $A \in NP$ ,  $A \leq_p B$ , B is NP-hard

key: B is as hard as everything in NP

NP-COMPLETE: if...

- in NP

- NP-Hard



CANONICAL PROBLEM: SAT

input: boolean func  $\phi$

output: does there exist Assignment X s.t.  $\phi(x) = T$

- if there exists a poly time alg. for any NP-C problem  $\Rightarrow P=NP$
- $P \neq NP \Rightarrow$  there does not exist poly time algo for any NP-C problem

## 2 WORLDS:

- 1: automated theorem proving
- $P=NP \Rightarrow$  2: public key encryption FAILS
- $P \neq NP \Rightarrow$  best  $2^{O(n)}$

- NP hard is worst case
- approx.
- SAT solvers

## 12/9 COMPUTABILITY: What is uncomputable?

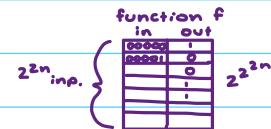
TURING (1936)

1. what is an algorithm? Turing machine
2. what is computable? what problems have algs?  
 linear  
 poly  
 expon.
3. what is uncomputable?

### SET UP:

- $\{0,1\}^*$  = set of all bitstrings
- program = bitstring (ex: Python in ASCII)
- input as bitstring in  $\{0,1\}^*$
- decision problem is just a function  $\{0,1\}^* \rightarrow \{0,1\}$   
 input YES/NO
- def: PROGRAM P computes f if  $\forall$  inputs  $x \in \{0,1\}^*$ ,  $P(x) = f(x)$   
 ↳ P always returns output ("Halts")

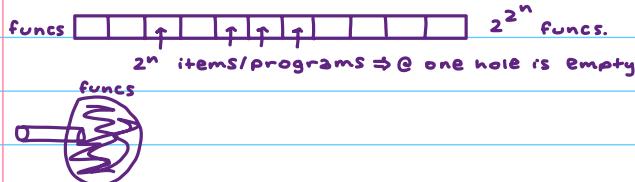
Life is unfair: there are more problems than solutions.



Fact: there exist functions  $f: \{0,1\}^{2^n} \rightarrow \{0,1\}$  that can't be computed by any length n program

Pf: there are  $2^{2^n}$  such functions f.e. if a program computes f,  
 $2^n$  bitstrings of length n.  
 ↳ it doesn't compute any other fn f.

⇒ by pigeonhole → ∃ some fn s.t. # program length n that computes it.



there are "natural" problems w/ no alg soln/no program that computes them.

## HALTING PROBLEM:

**input:** binary string interpreted as a pair  $(P, x)$   
→  $P$  as a program (Python prog. in ASCII)  
→  $x$  as input

**Output:** YES if  $P$  on input  $x$  HALTS. no otherwise.

Program takes  
a program as  
input

THM (TURING): no alg computes HALT (''Halt is undecidable'')

PF: by contradiction. Suppose you have alg. that computes HALT.

define program  $P(x)$ : if  $H(x, x) = \text{yes}$  → loop forever

otherwise → output yes

if we run  $P(P)$ :

if  $H(P,P) = \text{yes}$ , loops forever // if  $P(P)$  halts.  $P(P)$  doesn't halt.

otherwise, output yes

// if  $\rho(\rho)$  doesn't halt,  $\rho(\rho)$  halts.

## CONTRADICT!

## DIAGONALIZATION:

Programs P and Q

	0	1	2	3	4	5	6	7	8	9	...
inputs X	0	1	0	0	1	1	...				
P	0	N	Y	N	Y						
Q	1	N	Y	N	Y						
00	0	N	N	N	N						
01	0	1	1	1	1						

#diagonals are different

- there can exist programs that solve HALT on some/many inputs

\* run for 1 day, output "unsure" if still running

# if python prog. consist of only 1-9, +, -, \*, /



**TOTALITY:** takes program Q as input

Output: yes if Q halts on all inputs

no otherwise

THM: totality is undecidable

PF: towards contradiction, assume zig T compiles totally

then here is Halting alg:  $H(P,x)$ : define program  $Q :=$  "ignore input, runs  $P(x)$ , outputs YES"

return  $T(Q) \leftarrow$  totality big.

Claim: H computes Halt.  $\Rightarrow \Leftarrow$

1)  $P(x)$  halts  $\rightarrow Q$  halts on all inputs  $\rightarrow T(Q) = \text{YES}$  ✓

2)  $P(x)$  doesn't halt  $\rightarrow Q$  doesn't halt  $\rightarrow T(Q) = \text{No}$  ✓

### Equiv ("Program equivalence")

takes programs P, Q as input

outputs YES if  $\forall x \ P(x) = Q(x)$

NO otherwise

THM: Equiv. is undecidable

PF: by contradiction, assume E computes Equiv

consider H(P, x): defines Q = "ignore input, run P(x), output YES"

defines R = "ignore input, output YES"

return E(Q, R)

Claim: H' computes/solves HALT

$P(x)$  halts  $\Rightarrow Q$  outputs "yes"  $\Rightarrow \text{Equiv}(Q, R) = \text{"yes"}$

$P(x)$  doesn't halt  $\Rightarrow Q$  loops forever, R doesn't  $\Rightarrow \text{Equiv}(Q, R) = \text{"No"}$

### UNDECIDABILITY IN UNEXP. PLACES:

1) testing whether multivar. poly equation w/ int coeff has an int. soln.  $x^2 + 12y^3 + z^2 = 0$

2) whether a string is compressible given string x as input, is there length-l program that outputs x?