

CS 240E

Personal Notes

Marcus Chan

Taught by Therese Biedl

UW CS '25



Chapter 1: Algorithm Analysis

HOW TO "SOLVE" A PROBLEM

When solving a problem, we should
① write down exactly what the problem is.

e.g. Sorting Problem
→ given n numbers in an array,
put them in sorted order

② Describe the idea;

e.g. Insertion Sort



Idea:
repeatedly move
one item into the
correct space of
the sorted part.

③ Give a detailed description; usually pseudocode.

e.g. Insertion Sort:

```
for i=1,..,n-1
    j=i
    while j>0 and A[j-1] > A[j]
        swap A[j] and A[j-1]
        j-=1
```

④ Argue the correctness of the algorithm.
→ In particular, try to point out loop invariants & variants.

⑤ Argue the run-time of the program.

→ We want a theoretical bound.
(Using asymptotic notation).

To do this, we count the # of primitive operations.

PRIMITIVE OPERATIONS

In our computer model,
① our computer has memory cells
② all cells are equal
③ all cells are big enough to store our numbers

Then, "primitive operations" are $+, -, \times, \div$, load & following references.

We also assume each primitive operation takes the same amount of time to run.

ASYMPTOTIC NOTATION

BIG-O NOTATION: $f(n) \in O(g(n))$

We say that " $f(n) \in O(g(n))$ " if there exist $c > 0, n_0 > 0$ s.t.

$$|f(n)| \leq c|g(n)| \quad \forall n \geq n_0.$$

e.g. $f(n) = 75n + 500$ & $g(n) = 5n^2$,
 $c=1$ & $n_0=20$

Usually, " n " represents input size.

SHOW $2n^2 + 3n^2 + 11 \in O(n^2)$

To show the above, we need to find c, n_0 such that $O \leq 2n^2 + 3n^2 + 11 \leq cn^2 \quad \forall n \geq n_0$.

Sol². Consider $n_0=1$. Then

$$\begin{aligned} 1 \leq n &\Rightarrow 1 \leq 1n^2 \\ 1 \leq n &\Rightarrow n \leq n^2 \Rightarrow 3n \leq 3n^2 \\ &\xrightarrow{\text{(+)}} 2n^2 \leq 2n^2 \\ &\Rightarrow 2n^2 + 3n^2 \leq 11n^2 + 2n^2 + 3n^2 = 16n^2 \\ \text{Hence } c=16 \text{ & } n_0=1, \text{ so } 2n^2 + 3n^2 + 11 \in O(n^2). \end{aligned}$$

Ω -NOTATION (BIG OMEGA): $f(n) \in \Omega(g(n))$

We say " $f(n) \in \Omega(g(n))$ " if there exist $c > 0, n_0 > 0$ such that

$$c|g(n)| \leq |f(n)| \quad \forall n \geq n_0.$$

Θ -NOTATION (BIG OMEGA): $f(n) \in \Theta(g(n))$

We say " $f(n) \in \Theta(g(n))$ " if there exist $c_1, c_2 > 0, n_0 > 0$ such that

$$c_1|g(n)| \leq |f(n)| \leq c_2|g(n)|.$$

Note that

$$f(n) \in \Theta(g(n)) \Leftrightarrow f(n) \in O(g(n)) \text{ & } f(n) \in \Omega(g(n)).$$

\mathcal{O} -NOTATION (SMALL O): $f(n) \in \mathcal{O}(g(n))$

We say " $f(n) \in \mathcal{O}(g(n))$ " if for any $c > 0$, there exists some $n_0 > 0$ such that

$$|f(n)| < c|g(n)| \quad \forall n \geq n_0.$$

If $f(n) \in \mathcal{O}(g(n))$, we say $f(n)$ is "asymptotically strictly smaller" than $g(n)$.

ω -NOTATION (SMALL OMEGA): $f(n) \in \omega(g(n))$

We say $f(n) \in \omega(g(n))$ if for all $c > 0$, there exists some $n_0 > 0$ such that

$$0 \leq c|g(n)| < |f(n)| \quad \forall n \geq n_0.$$

FINDING RUNTIME OF A PROGRAM

- To evaluate the run-time of a program, given its pseudocode, we do the following:
- Annotate any primitive operations with just " $\Theta(1)$ ";
 - For any loops, find the worst-case bound for how many times it will execute;
 - Calculate the big-O run time of the program;
 - Argue this bound is tight (ie show program is also in $\Omega(g(n))$, so runtime $\in \Theta(g(n))$.)

e.g. insertion sort

```
for i=1, ..., n-1
    j=i  $\Theta(1)$ 
    while j>0 and A[j-1] > A[j]
        swap A[j] and A[j-1]  $\Theta(1)$ 
        j-1  $\Theta(1)$ 
```

Then, let c be a const s.t. the upper bounds all the times needed to execute one line.
 $\text{So runtime} \leq n \cdot n \cdot c = c \cdot n^2 \in O(n^2)$.

Next, consider the worst pos. case of insertion sort.

For each $A[i:j]$, we need $i-1$ swaps.

So

$$\begin{aligned}\text{runtime} &\geq \sum_{i=1}^{n-1} (i-1) \\ &= \sum_{i=1}^{n-1} i = \frac{(n-1)n}{2} \\ &\in \Omega(n^2),\end{aligned}$$

and so
 $\text{runtime} \in \Theta(n^2)$.

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = L < \infty \Rightarrow f(n) \in O(g(n))$$

\ll LIMIT RULE I \gg (LI.1(2))

(\exists (let $f(n), g(n)$ be such that $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = L < \infty$).

Then necessarily $f(n) \in O(g(n))$.

Proof. We know $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = L$.
 $\Rightarrow \forall \epsilon > 0: \exists n_0 \text{ s.t. } \left| \frac{f(n)}{g(n)} - L \right| < \epsilon \forall n \geq n_0$.

We want to show
 $\exists C > 0, \exists n_0 \Rightarrow \forall n \geq n_0, f(n) \in O(g(n))$.

Choose $\epsilon = 1$. Then there exists n_1 , s.t.

$$\forall n \geq n_1, \left| \frac{f(n)}{g(n)} - L \right| \leq 1.$$

$$\Leftrightarrow \frac{f(n)}{g(n)} - L \leq \left| \frac{f(n)}{g(n)} - L \right| \leq 1.$$

$$\Leftrightarrow \frac{f(n)}{g(n)} \leq L + 1$$

$$\Leftrightarrow f(n) \in O(g(n)) + g(n) \quad (\text{since } g(n) > 0)$$

Choose $C = L+1$. Note $f(n), g(n) > 0$, so $L+1 > 0$, and since $L < \infty$, thus $C < \infty$.
Now, for all $n \geq n_1$, $f(n) \leq C \cdot g(n)$, and so $f(n) \in O(g(n))$. \square

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0 \Leftrightarrow f(n) \in o(g(n))$$

\ll LIMIT RULE II \gg (LI.1(1))

(\exists (let $f(n), g(n)$). Then $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$

iff $f(n) \in o(g(n))$).

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = L > 0 \Rightarrow f(n) \in \Omega(g(n))$$

\ll LIMIT RULE III \gg (LI.1(3))

(\exists (let $f(n), g(n)$ such that $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = L > 0$).

Then necessarily $f(n) \in \Omega(g(n))$.

$$\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \infty \Rightarrow f(n) \in \omega(g(n))$$

\ll LIMIT RULE IV \gg (LI.1(4))

(\exists (let $f(n), g(n)$ such that $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = \infty$).

Then necessarily $f(n) \in \omega(g(n))$.

OTHER LIMIT RULES

The following are corollaries of the limit rules:

- ① $f(n) \in \Theta(f(n))$ } (Identity)
- ② $K \cdot f(n) \in \Theta(f(n)) \forall K \in \mathbb{R}$ } (Constant multiplication)
- ③ $f(n) \in O(g(n)), g(n) \in O(h(n)) \Rightarrow f(n) \in O(h(n))$ } (Transitivity)
- ④ $f(n) \in \Omega(g(n)), g(n) \in \Omega(h(n)) \Rightarrow f(n) \in \Omega(h(n))$
- ⑤ $f(n) \in O(g(n)), g(n) \in \Theta(h(n)) \forall n \geq N \Rightarrow f(n) \in O(h(n))$
- ⑥ $f(n) \in \Omega(g(n)), g(n) \geq h(n) \forall n \geq N \Rightarrow f(n) \in \Omega(h(n))$
- ⑦ $f_1(n) \in O(g_1(n)), f_2(n) \in O(g_2(n)) \Rightarrow f_1(n) + f_2(n) \in O(g_1(n) + g_2(n))$
- ⑧ $f_1(n) \in \Omega(g_1(n)), f_2(n) \in \Omega(g_2(n)) \Rightarrow f_1(n) + f_2(n) \in \Omega(g_1(n) + g_2(n))$
- ⑨ $h(n) \in O(f(n) + g(n)) \Rightarrow h(n) \in O(\max(f(n), g(n)))$ } (Maximum-rule for O)
- ⑩ $h(n) \in \Omega(f(n) + g(n)) \Rightarrow h(n) \in \Omega(\max(f(n), g(n)))$ } (Maximum-rule for Ω)

$f(n) \in P_d(\mathbb{R}) \Rightarrow f(n) \in \Theta(n^d)$ \ll POLYNOMIAL RULE \gg

(\exists (at $f(n) \in P_d(\mathbb{R})$, ie of the form $f(n) = c_0 + c_1 n + \dots + c_d n^d$).

Then necessarily $f(n) \in \Theta(n^d)$.

$b > 1; \log_b(n) \in \Theta(\log n)$ \ll LOG RULE I \gg

(\exists (at $b > 1$. Then necessarily $\log_b(n) \in \Theta(\log n)$).

Proof. Note

$$\lim_{n \rightarrow \infty} \frac{\log_b(n)}{\log n} = \lim_{n \rightarrow \infty} \frac{\left(\frac{\log(n)}{\log(b)} \right)}{\log(n)} = \lim_{n \rightarrow \infty} \frac{1}{\log(b)} > 0,$$

so by Limit Rules 1 & 3, $\log_b(n) \in \Theta(\log n)$. \square

*convention:
"log" = "log₂".

$c, d > 0; \log^c n \in O(n^d)$ \ll LOG RULE II \gg

(\exists (at $c, d > 0$. Then necessarily $\log^c n \in O(n^d)$),

where $\log^c n \equiv (\log n)^c$.

Proof. See that

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{\ln^k n}{n} &= \lim_{n \rightarrow \infty} \frac{k \ln^{k-1} n \cdot \frac{1}{n}}{1} \\ &= \dots \\ &= \lim_{n \rightarrow \infty} \frac{k!}{n} = 0,\end{aligned}$$

so $\ln^k n \in o(n)$.

Fix $c, d > 0$. Then

$$\begin{aligned}\lim_{n \rightarrow \infty} \frac{\ln^c n}{n^d} &= \left(\lim_{n \rightarrow \infty} \frac{\ln^c n}{\ln^d n} \right)^d \\ &\leq \left(\lim_{n \rightarrow \infty} \frac{\ln^c n}{n} \right)^d \\ &= 0^d = 0.\end{aligned}$$

As $\log^c n = \left(\frac{1}{\ln 2}\right)^c \ln^c n$, thus $\lim_{n \rightarrow \infty} \frac{\log^c n}{n^d} = \left(\frac{1}{\ln 2}\right)^c \cdot 0 = 0$.

Proof follows from the limit rule. \square

$f(n) \in o(g(n)) \Rightarrow f(n) \in O(g(n))$

(\exists (suppose $f(n) \in o(g(n))$. Then $f(n) \in O(g(n))$).

$f(n) \in O(g(n)) \Rightarrow f(n) \in \Omega(g(n))$

(\exists (suppose $f(n) \in O(g(n))$. Then $f(n) \in \Omega(g(n))$).

Proof. Prove by contrapositive:

$$f(n) \in \Omega(g(n)) \Rightarrow f(n) \notin o(g(n)).$$

Consider cases for $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)}$.

Case 1 It DNE.

$$\Rightarrow f(n) \notin o(g(n)).$$

Case 2 Limit exists.

Then by $f(n) \in \Omega(g(n))$, thus

$$f(n) \geq c \cdot g(n) \text{ for some } c > 0 \text{ & } n \geq n_0.$$

$$\Rightarrow \lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} \geq c > 0$$

$$\Rightarrow \text{limit} \neq 0$$

$$\Rightarrow f(n) \notin o(g(n)). \quad \square$$

$f(n) \in \omega(g(n)) \Rightarrow f(n) \in \Omega(g(n))$

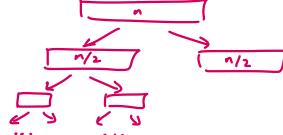
Suppose $f(n) \in \omega(g(n))$. Then $f(n) \in \Omega(g(n))$.

WORST-CASE RUNTIME: $T_A^{\text{worst}}(n)$

The "worst-case runtime" for an algorithm, denoted $T_A^{\text{worst}}(n)$, is the max run-time among all instances of size n .

ANALYZING RECURSIVE ALGORITHMS

Consider merge sort:



Analysis of MergeSort:

```
MergeSort(A, n, l=0, r=n-1, S=NULL)
A: array of size n, 0 ≤ l ≤ r ≤ n-1
if S is NIL init it as array S[0...n-1]
if (r < l) then
    return
else
    m = (l+r)/2
    MergeSort(A, n, l, m, S)
    MergeSort(A, n, m+1, r, S)
    Merge(A, l, m, r, S)
```

Merge(A, l, m, r, S)

$A[0, \dots, n-1]$ is an array. $A[l, \dots, m]$ is sorted.
 $A[m+1, \dots, r]$ is sorted. $S[0, \dots, n-1]$ is an array.

```
copy A[l...r] into S[l...r]
int i_L=l; int i_R=m+1
for (k=L; k < R; k++) do
    if (i_L > m) A[k] ← S[i_R++]
    else if (i_R > r) A[k] ← S[i_L++]
    else if (S[i_L] ≤ S[i_R]) A[k] ← S[i_L++]
    else A[k] ← S[i_R++]
```

Arguing run-time:

Let $T(n)$ = run-time of merge sort with n items.

$$\therefore T(n) \leq \begin{cases} c, & n \leq 1 \\ 2 \cdot T\left(\frac{n}{2}\right) + c \cdot n, & \text{otherwise} \end{cases}$$

$\in \Theta(n \log n)$ (see below)

SOME RECURRENCE RELATIONS

Note:

Recursion...	... resolves to ...	Example
$T(n) = T(n/2) + \Theta(1)$	$T(n) \in \Theta(\log n)$	Binary search
$T(n) = 2T(n/2) + \Theta(n)$	$T(n) \in \Theta(n \log n)$	Mergesort
$T(n) = 2T(n/2) + \Theta(\log n)$	$T(n) \in \Theta(n)$	Heapify
$T(n) = T(cn) + \Theta(n), 0 < c < 1$	$T(n) \in \Theta(n)$	Selection
$T(n) = 2T(n/4) + \Theta(1)$	$T(n) \in \Theta(N^{\sqrt{2}})$	Range Search
$T(n) = T(N^{\sqrt{n}}) + \Theta(1)$	$T(n) \in \Theta(\log \log n)$	Interpolation Search

SORTING PERMUTATION [OF AN ARRAY]

A "sorting permutation" of an array A is the permutation $\pi: \{0, \dots, n-1\} \rightarrow \{0, \dots, n-1\}$ such that

$$A[\pi(0)] \leq A[\pi(1)] \leq \dots \leq A[\pi(n-1)].$$

* A off would be "sorted".

eg if $A = [14, 3, 2, 6, 1, 11, 7]$, then

$$\pi = [4, 2, 1, 3, 6, 5, 0].$$

For a sorting permutation π , we define its inverse π^{-1} to be the array which entries have exactly the same "relative order" as in A .

eg $\pi^{-1} = [6, 2, 1, 3, 0, 5, 4]$ (if A is as above)

We denote " Π_n " to be the set of all sorting permutations of $\{0, \dots, n-1\}$.

AVERAGE-CASE RUNTIME: $T_A^{\text{avg}}(n)$

The "average-case run-time" of an algorithm is defined to be

$$T_A^{\text{avg}}(n) := \underset{I \in \mathcal{X}_n}{\text{avg}} T_A(I) = \frac{1}{|\mathcal{X}_n|} \sum_{I \in \mathcal{X}_n} T_A(I),$$

where \mathcal{X}_n is the set of instances of size n .

In particular, if we can map each instance I to a permutation $\pi \in \Pi_n$, where Π_n is the set of all permutations of $\{0, \dots, n-1\}$, then we may alternatively state that

$$T^{\text{avg}}(n) = \frac{1}{n!} \sum_{\pi \in \Pi_n} T(\pi)$$

where $T(\pi)$ is the number of comparisons on instance π^{-1} .

EXAMPLE: AVG CASE DEMO

Consider the algorithm

```

avgCaseDemo(A, n)
// array A stores n distinct numbers
1. if n ≤ 2 return
2. if A[n-2] ≤ A[n-1] then
3.   avgCaseDemo(A[0...n-2], n-2) // good case
4. else
5.   avgCaseDemo(A[0...n-3], n-2) // bad case
    
```

We claim $T^{\text{avg}}(n) \in O(\log n)$.

Proof. To avoid constants, let $T(\cdot) := \# \text{ of recursions}$; the run-time is proportional to this.
As all numbers are distinct, we may associate each array with a sorting permutation.
So for $\pi \in \Pi_n$, let $T(\pi) = \# \text{ of recursions done if the input array has sorting permutation } \pi$.
Note we have two kinds of permutations:
① "Good" permutations — if $A[n-2] < A[n-1]$; or
② "Bad" permutations — if $A[n-2] > A[n-1]$.

Denote
 $\Pi_n^{\text{good}} = \# \text{ of good permutations of size } n$
& $\Pi_n^{\text{bad}} = \# \text{ of bad permutations of size } n$.

Then, we claim that

$$\sum_{\substack{\pi \in \Pi_n \\ \text{good}}} T(\pi) \leq |\Pi_n^{\text{good}}| (1 + T^{\text{avg}}(\lfloor \frac{n}{2} \rfloor))$$

$$\& \sum_{\substack{\pi \in \Pi_n \\ \text{bad}}} T(\pi) \leq |\Pi_n^{\text{bad}}| (1 + T^{\text{avg}}(n-2)).$$

Proof. We only prove this for good permutations; the other claim is similar.

Fix $\pi \in \Pi_n^{\text{good}}$, and let π_{half} be the permutation of the recursion;
ie $\pi_{\text{half}} = \text{the sorting perm of } \pi[0, \dots, \lfloor \frac{n}{2} \rfloor]$ & $T(\pi) = 1 + T(\pi_{\text{half}})$.

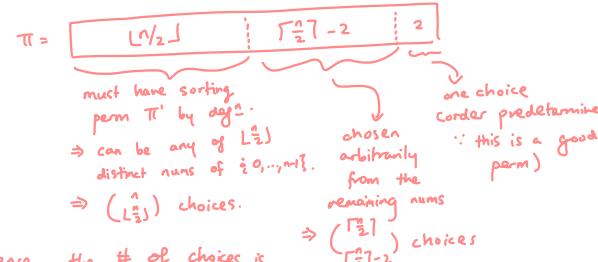
Note that $\pi_{\text{half}} \in \Pi_{\lfloor \frac{n}{2} \rfloor}$. Then see that

$$\begin{aligned} \sum_{\pi \in \Pi_n^{\text{good}}} T(\pi) &= \sum_{\pi \in \Pi_n^{\text{good}}} (1 + T(\pi_{\text{half}})) \\ &= |\Pi_n^{\text{good}}| + \sum_{\pi \in \Pi_n^{\text{good}}} T(\pi_{\text{half}}) \\ &= |\Pi_n^{\text{good}}| + \sum_{\substack{\pi' \in \Pi_{\lfloor \frac{n}{2} \rfloor} \\ \exists \pi \in \Pi_n^{\text{good}} \text{ for which } \pi_{\text{half}} = \pi[0, \dots, \lfloor \frac{n}{2} \rfloor]} | \cdot T(\pi'). \end{aligned}$$

We next prove the following claim:

Claim If $n \geq 3$, then $|\Pi_n^{\text{good}}(\pi')| = \frac{n!}{2!(\lfloor \frac{n}{2} \rfloor)!} \quad \forall \pi' \in \Pi_{\lfloor \frac{n}{2} \rfloor}$.

Proof. Fix $\pi' \in \Pi_{\lfloor \frac{n}{2} \rfloor}$. See that



Hence, the # of choices is $\binom{\lceil \frac{n}{2} \rceil}{\lfloor \frac{n}{2} \rfloor} \cdot (\lceil \frac{n}{2} \rceil - 2)!$

$$\left(\binom{\lceil \frac{n}{2} \rceil}{\lfloor \frac{n}{2} \rfloor} \cdot \binom{\lceil \frac{n}{2} \rceil}{\lceil \frac{n}{2} \rceil - 2} \cdot (\lceil \frac{n}{2} \rceil - 2)! \right) = \dots = \frac{n!}{\lfloor \frac{n}{2} \rfloor! \cdot 2},$$

as needed. \blacksquare

Then, since $|\Pi_n^{\text{good}}| = \frac{1}{2} |\Pi_n| = \frac{n!}{2}$, it follows that

$$|\Pi_n^{\text{good}}| / |\Pi_{\lfloor \frac{n}{2} \rfloor}|, \text{ and so}$$

$$\begin{aligned} \sum_{\pi \in \Pi_n^{\text{good}}} T(\pi) &= |\Pi_n^{\text{good}}| + \sum_{\pi' \in \Pi_{\lfloor \frac{n}{2} \rfloor}} |\Pi_n^{\text{good}}(\pi')| \cdot T(\pi') \\ &= |\Pi_n^{\text{good}}| + \sum_{\pi' \in \Pi_{\lfloor \frac{n}{2} \rfloor}} \frac{|\Pi_n^{\text{good}}|}{|\Pi_{\lfloor \frac{n}{2} \rfloor}|} \cdot T(\pi') \\ &= |\Pi_n^{\text{good}}| \left(1 + \frac{1}{|\Pi_{\lfloor \frac{n}{2} \rfloor}|} \sum_{\pi' \in \Pi_{\lfloor \frac{n}{2} \rfloor}} T(\pi') \right) \\ &= |\Pi_n^{\text{good}}| (1 + T^{\text{avg}}(\lfloor \frac{n}{2} \rfloor)) \end{aligned}$$

as needed. \blacksquare

Next, see that $|\Pi_n^{\text{good}}| = |\Pi_n^{\text{bad}}|$, since we can map any bad perm to a good one (and v.v.) by swapping $A[n-2]$ & $A[n-1]$.

Thus $|\Pi_n^{\text{good}}| = |\Pi_n^{\text{bad}}| = \frac{n!}{2}$, and so

$$\begin{aligned} T^{\text{avg}}(n) &\leq \frac{1}{|\Pi_n|} \sum_{\pi \in \Pi_n} T(\pi) \leq \frac{1}{|\Pi_n|} \left(\sum_{\pi \in \Pi_n^{\text{good}}} T(\pi) + \sum_{\pi \in \Pi_n^{\text{bad}}} T(\pi) \right) \\ &\leq \frac{1}{|\Pi_n|} (|\Pi_n^{\text{good}}| (1 + T^{\text{avg}}(\lfloor \frac{n}{2} \rfloor)) + |\Pi_n^{\text{bad}}| (1 + T^{\text{avg}}(n-2))) \\ &= 1 + \frac{1}{2} T^{\text{avg}}(\lfloor \frac{n}{2} \rfloor) + \frac{1}{2} T^{\text{avg}}(n-2). \end{aligned}$$

Finally, we show $T^{\text{avg}}(n) \leq 2 \log n$ by induction.

Clearly, this holds for $n \leq 2$.

So, assume $n \geq 3$. Assume the inductive hypothesis. Then see that

$$\begin{aligned} T^{\text{avg}}(n) &\leq 1 + \frac{1}{2} T^{\text{avg}}(\lfloor \frac{n}{2} \rfloor) + \frac{1}{2} T^{\text{avg}}(n-2) \\ &\leq 1 + \frac{1}{2} (2 \log(\lfloor \frac{n}{2} \rfloor)) + \frac{1}{2} (2 \log(n-2)) \\ &\leq 1 + \log(n-1) + \log(n) \\ &= 2 \log(n), \end{aligned}$$

which suffices to prove the claim. \blacksquare

EXPECTED-CASE RUNTIME: $T_A^{\text{exp}}(n)$

The "expected-case runtime" of an algorithm is defined to be

$$T_A^{\text{exp}}(I) = E[T_A(I, R)] = \sum_{\text{all } R} T_A(I, R) P(R)$$

where R is a sequence of random outcomes, and $P(R)$ denotes the probability the random variables in the algorithm A have outcomes R .

AMORTIZED RUN-TIME: T_0^{amort}

Let \mathcal{O} be an operation, and let $T^{\text{actual}}(\mathcal{O})$ be the actual run-time of \mathcal{O} .

Then, we say \mathcal{O} has "amortized run-time $T^{\text{amort}}(\mathcal{O})$ " if for any sequence of operations $\mathcal{O}_1, \dots, \mathcal{O}_k$ that could occur, we have

$$\sum_{i=1}^k T^{\text{actual}}(\mathcal{O}_i) \leq \sum_{i=1}^k T^{\text{amort}}(\mathcal{O}_i).$$

POTENTIAL FUNCTION [OF A DATA STRUCTURE]: Φ

A "potential function" is a function $\Phi(\cdot)$ that depends on the status of the data structure.

In particular, for any sequence $\mathcal{O}_1, \dots, \mathcal{O}_k$ of operations:

- ① $\Phi(i) \geq 0 \quad \forall i \geq 0$, where Φ_i is the value of Φ after $\mathcal{O}_1, \dots, \mathcal{O}_i$ have been executed; &
- ② $\Phi(0) = 0$.

$$T^{\text{amort}}(\mathcal{O}) := T^{\text{actual}}(\mathcal{O}) + \Phi_{\text{after}} - \Phi_{\text{before}}$$

UPPER BOUNDS ACTUAL RUN-TIME

For any potential function Φ , the function

$$T^{\text{amort}}(\mathcal{O}) := T^{\text{actual}}(\mathcal{O}) + \Phi_{\text{after}} - \Phi_{\text{before}}$$

upper-bounds the actual run-time,

where Φ_{before} & Φ_{after} denote the state of the potential function before & after \mathcal{O} .

Proof. Fix a sequence of operations $\mathcal{O}_1, \dots, \mathcal{O}_k$.

Summing up the amortized times and using a telescoping sum, we get

$$\begin{aligned} \sum_{i=1}^k T^{\text{amort}}(\mathcal{O}_i) &= \sum_{i=1}^k (T^{\text{actual}}(\mathcal{O}_i) + \Phi(i) - \Phi(i-1)) \\ &= \sum_{i=1}^k T^{\text{actual}}(\mathcal{O}_i) + \sum_{i=1}^k (\Phi(i) - \Phi(i-1)) \\ &= \sum_{i=1}^k T^{\text{actual}}(\mathcal{O}_i) + \Phi(k) - \Phi(0) \\ &\geq \sum_{i=1}^k T^{\text{actual}}(\mathcal{O}_i). \quad \blacksquare \end{aligned}$$

STEPS TO PERFORM AMORTIZED ANALYSIS USING THE POTENTIAL FUNCTION METHOD

To do amortized analysis using potential functions, we do the following:

- ① Define a "time unit", so that an operation with run-time $\Theta(k)$ takes at most k time units.
- ② Define a potential function Φ and verify $\Phi(0) = 0$ & $\Phi(i) \geq 0 \quad \forall i \geq 0$.
- ③ For each operation \mathcal{O} , compute

$$T^{\text{amort}}(\mathcal{O}) := T^{\text{actual}}(\mathcal{O}) + \Phi_{\text{after}} - \Phi_{\text{before}},$$

and find an asymptotic upper bound for it.

EXAMPLE: DYNAMIC ARRAYS

Consider dynamic arrays with two operations:

- ① insert; &
- ② rebuild, where we "lengthen" the array by a factor of 2.

Hence, insert takes $\Theta(1)$ time, & rebuild takes $\Theta(n)$ time (where n is the size of the array).

Define a "time unit" such that insert takes "one unit of time" & rebuild takes " n units of time".

We claim $T^{\text{amort}}(\text{insert}) = 3$ & $T^{\text{amort}}(\text{rebuild}) = 0$.

Proof. Let the potential function Φ be defined by $\Phi(i) = \max\{0, 2 \cdot \text{size} - \text{capacity}\}$.

Clearly $\Phi(i) \geq 0$. Also initially size=0 & cap=0, so $\Phi(0) = 0$ as desired.

Now, the amortized run-time for insert is $T^{\text{amort}}(\text{insert}) = T^{\text{actual}}(\text{insert}) + \Phi_{\text{after}} - \Phi_{\text{before}} \leq 3$,

as the actual time is ≤ 1 unit, the size increases by 1 & the capacity does not change.

Similarly,
 $T^{\text{amort}}(\text{rebuild}) = T^{\text{actual}}(\text{rebuild}) + \Phi_{\text{after}} - \Phi_{\text{before}}$
 $\leq n + (0-n)$
 $= 0. \quad \blacksquare$

Chapter 2: Priority Queues and Heaps

ADT PRIORITY QUEUES

Q₁: A "priority queue" stores items that have a priority, or key.

Q₂: Operations:

- ① insert (a given item & priority as a k-v pair)
- ② deleteMax (return item with largest priority)
- ③ size, isEmpty

Q₃: We can use a PQ to sort:

PQ-Sort (A[0, ..., n-1])

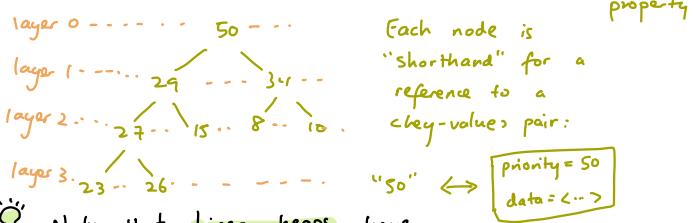
```
init PQ to an empty PQ  
for i=0 to n-1 do  
    PQ.insert(A[i])  
for i=n-1 down to 0 do  
    A[i] ← PQ.deleteMax()
```

Q₄: The run-time of the above algorithm is O(initialization + n·insert + n·deleteMax).

BINARY HEAPS

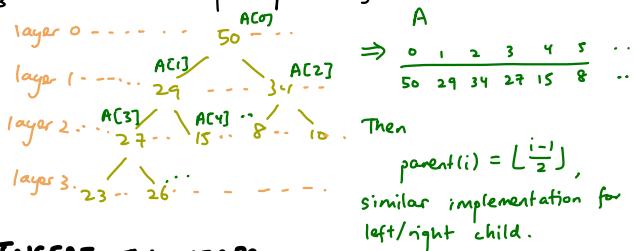
Q₁: "Binary heaps" are binary trees with the following properties:

- ① Each level is filled except the last, which is filled from the left; } "structural" property
- ② key(i) ≤ key(parent(i)) } "heap-order property"



Q₂: Note that binary heaps have height $O(\log n)$.

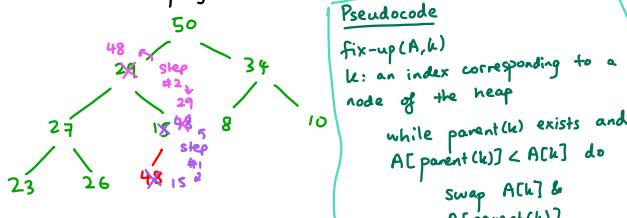
Q₃: We store binary heaps using an array:



INSERT IN HEAPS

Q₁: To insert into a heap, we just place the new key at the first free leaf.

Q₂: We also employ "fix-up":



Pseudocode

fix-up(A, k)

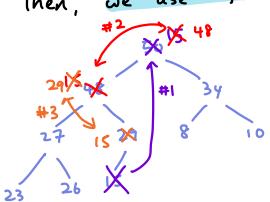
 k: an index corresponding to a node of the heap
 while parent(k) exists and $A[\text{parent}(k)] < A[k]$ do
 swap $A[k]$ & $A[\text{parent}(k)]$
 k ← parent(k)

DELETMAX IN HEAPS

The maximum item of a heap is just the root node.

We replace the root by the last leaf, which is taken out.

Then, we use "fix-down":



* we "swap down" from the root-node until the heap-ordering is satisfied.

Run-time: $O(\text{height}) = O(\log n)$

HEAPSORT

If we use a heap as a PQ and sort using it, the run-time of said algorithm is

$$\begin{aligned} T(n) &\in O(\text{initialization} + n \cdot \text{insert} + n \cdot \text{deleteMax}) \\ &= O(\Theta(1) + n \cdot O(\log n) + n \cdot O(\log n)) \\ &= O(n \log n). \end{aligned}$$

Pseudocode:

HeapSort(A, n)

```
// heapify
n ← A.size()
for i ← parent(last()) down to 0 do
    fix-down(A, i, n)
// repeatedly find maximum
while n > 1
    // 'delete' maximum by moving to end and
    // decreasing n
    swap item at A[root()] and A[last()]
    n --
    fix-down(A, root(), n)
```

* Heapify:

Given: all items that should be on the heap

It builds the heap all at once.

How? → by fixing-down incrementally.

Heapify has run-time $\Theta(n)$.

Why? → analyze a recursive version:

```
heapify(node i)
    if i has left child
        heapify(left child i)
    if i has right child
        heapify(right child i)
    fix-down(i)
```

Now, let

$T(n)$:= runtime of heapify on n items.

(Assume n divisible as needed.)

Then:

$$\text{size(left child)} = \text{size(right child)} = \frac{n-1}{2}.$$

$$\therefore T(n) = \begin{cases} \Theta(1), & n \leq 1 \\ T\left(\frac{n-1}{2}\right) + T\left(\frac{n-1}{2}\right) + \Theta(\log n), & \text{otherwise} \end{cases}$$

$$\Rightarrow T(n) \leq 2T\left(\frac{n}{2}\right) + \Theta(\log n) \in O(n).$$

Note this uses $O(1)$ auxiliary space, since we use the same input-array A for storing the heap.



OTHER PQ OPERATIONS

We can also support the following operations:

① "findMax" — finds the max element without removing it;

- in a bin heap this takes $\Theta(1)$

- since it is just the root node

② "decreaseKey" — takes in a ref i to the location of one item of the heap and a key k_{new} , and decreases the key of i to k_{new} if $k_{\text{new}} < \text{key}(i)$

- does nothing if $k_{\text{new}} \geq \text{key}(i)$

- easy to do in a bin heap; just need to change the key & then call fix-down on i to its children

③ "increaseKey" — "opposite" of decreaseKey.

- easy in bin heap, but just call fix-up instead of fix-down

④ "delete" — delete the item i (which we have a ref to).

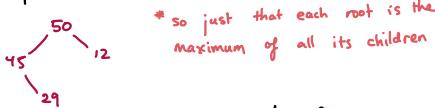
- for bin heap: increase value of key at i to ∞ (or $\text{findMax}().\text{key}() + 1$);

- then call deleteMax

- takes $O(\log n)$

MELDABLE HEAP

A "meldable heap" is the same as a "binary heap", except it drops the structural property.



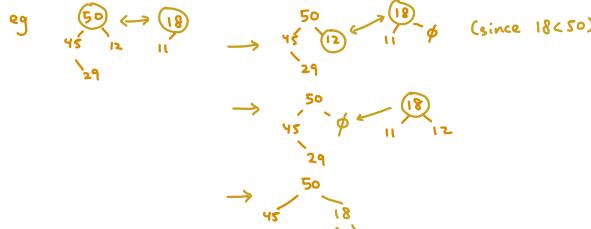
Operations for meldable heaps:

- ① insert
 - create 1-node meldable heap P' w/ new k-v pair we want to add
 - call $\text{merge}(P, P')$ w/ existing & newly-created meldable heap

- ② deleteMax
 - max is root, so just remove that
 - then return $\text{merge}(P_L, P_R)$, where P_L & P_R are the "subheaps" of the original heap

③ merge

```
meldableHeap::merge(r1, r2)
input: r1, r2 (roots of two meldable heaps)
      - ri ≠ NIL
output: returns root of merged heap
if r2 is NIL then return r1
if r1.key < r2.key then swap r1, r2
randomly pick one child c of r1
replace c by result of merge(r2, c)
return r1
```



$$\text{AVG. RUN-TIME (MERGE)} = O(\log n_1 + \log n_2)$$

(L2.3)

Note that

avg run-time (merge) = $O(\log n_1 + \log n_2)$
where n_1, n_2 are the sizes of the heaps to be merged.

BINOMIAL HEAPS

FLAGGED TREES

- 1 A "flagged tree" is one where every level is full, but the root node only has a left child.
- 2 Note that a flagged tree of height h has 2^h nodes.



BINOMIAL HEAP (C02.3)

A "binomial heap" is a list L of binary trees such that

- any tree in L is a flagged tree (structural property); &
- for any node v , all keys in the left subtree of v are no bigger than $v.key$. (order property).



PROPER BINOMIAL HEAP

- We say a BH is "proper" if no two flagged trees in L have the same height.



A PBH OF SIZE n CONTAINS $\leq \log(n)+1$

FLAGGED TREES (C02.2)

Let a PBH have size n .

Then it contains at most $\log(n)+1$ flagged trees.

Proof: Let T be the tree w/ max height, say h , in the list L of flagged trees.

Then T has 2^h nodes, so $2^h \leq n$, ie $h \leq \log(n)$.

Since the trees in L have distinct heights, we have at most one tree for each height $0, \dots, h$, and so at most $h+1 \leq \log(n)+1$ trees. \square

MAKING A BH PROPER

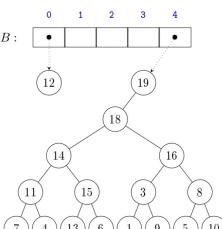
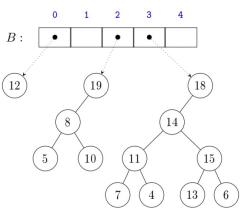
- To make a BH proper, we do the following: if the BH has two flagged trees T, T' of the same height h , then they both have 2^h nodes.

We want to combine them into one flagged tree of height $h+1$.

To do so, we use the following algorithm:

(A2.1)

```
binomialHeap::makeProper()
n <-- size of the binomial heap
for (l=0; n>1; n <-> [n/2]) do l++ // compute [log n]
B <- array of size l+1, initialized at NIL
L <- list of flagged trees
while L is non-empty do
    T <- L.pop(), h <- T.height
    while T' <- B[h] is not NIL do
        if T.root.key < T'.root.key then swap
        T & T'
        T'.right <- T.left, T.left <- T', T.height <- h+1 // merge T with T'
        B[h] <- NIL, h++
    B[h] <- T
for (h=0; h<l; h++) do
    if B[h] ≠ NIL then L.append(B[h]) // copy B back to the list
```



PQ OPERATIONS FOR PBHS

Each of the PQ operations can be performed in $O(\log n)$ time with makeProper.

① merge(P_1, P_2)

- concat lists of P_1, P_2 into one
- then call makeProper
- takes time $O(\log n_1 + \log n_2) \leq O(\log n)$

② insert(k, v)

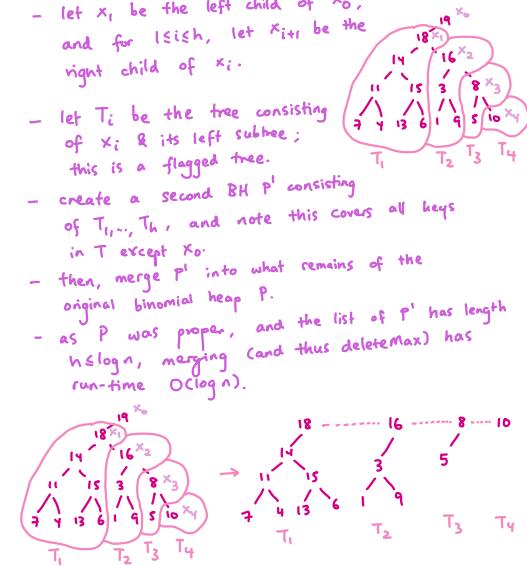
- like w/ meldable heaps: create a single-node binomial heap, then call merge
- takes time $O(\log n)$

③ findMax()

- Scan through L and compare keys of the roots
- largest is just maximum
- takes time $O(|L|) \leq O(\log n)$

④ deleteMax()

- assume we found the max at root x_0 of flagged tree T , say w/ height h (this takes $O(\log n)$ time if heights not already known)
- remove T from L and split it as follows:
 - let x_i be the left child of x_0 ,
 - and for $i \neq h$, let x_{i+1} be the right child of x_i .
- let T_i be the tree consisting of x_i & its left subtree; this is a flagged tree.
- create a second BH P' consisting of T_1, \dots, T_h , and note this covers all keys in T except x_0 .
- then, merge P' into what remains of the original binomial heap P .
- as P was proper, and the list of P' has length $h \leq \log n$, merging (and thus deleteMax) has run-time $O(\log n)$.



Chapter 3:

Sorting

THE SELECTION PROBLEM

The "selection problem" is:

"Given an array $A[0, \dots, n-1]$ and an index $0 \leq k \leq n$, $\text{select}(A, k)$ should return the element in A that would be at index k if we sorted A ".

e.g. if $A[0, \dots, 9] = [30, 60, 10, 0, 50, 80, 90, 10, 40, 70]$
then the sorted array would be

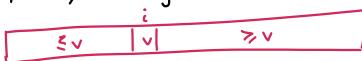
$[0, 10, 10, 30, 40, 50, 60, 70, 80, 90]$

so

$\text{select}(3) = 30$.

PARTITION [THE FUNCTION]

The "partition" function does the following:
given the pivot-value $v = A[p]$ and an array $A[0, \dots, n-1]$, rearrange A s.t.



and return the pivot-index i .

Partition algorithm ("efficient in-place partition —

Hoare"):

```
partition(A, p)
// A: array of size n
// p: integer s.t.  $0 \leq p \leq n$ 
1. swap(A[n-1], A[n])
2. i ← 1, j ← n-1, v ← A[n-1]
3. loop
4.   do i ← i+1 while  $A[i] < v$ 
5.   do j ← j-1 while  $j \geq 1$  &  $A[j] > v$ 
6.   if i ≥ j then break (goto 9)
7.   else swap(A[i], A[j])
8. end loop
9. swap(A[n-1], A[i])
10. return i
```

* we keep swapping
the outer-most
wrongly-positioned
pairs

QUICK-SELECT

The "quick-select" algorithm:

```
quick-select(A, k)
// A: array of size n
// k: integer s.t.  $0 \leq k \leq n$ 
1. p ← choose-pivot(A) // for now, p=n-1
2. i ← partition(A, p)
3. if i=k then
4.   return A[i]
5. else if i < k then
6.   return quick-select(A[0, ..., i-1], k)
7. else if i > k then
8.   return quick-select(A[i+1, ..., n-1], k-i-1)
```

* intuition:

Have: $\begin{array}{c|c} \leq v & | v | \\ \hline \end{array}$

Want: $\begin{array}{c|c} \leq m & | m | \\ \hline \end{array}$

* Run-time analysis of quick-select:

We analyze the # of key-comparisons.
↳ so we don't mess with constants.

In particular, partition uses n key comparisons.

Then, the run-time on an instance I is

$$T(I) = \underbrace{n + T(I')}_\text{partition subarray}$$

How big is I' ?

best case: don't recurse at all $\rightarrow O(n)$

worst case: $|I'| = n-1$

$$\begin{aligned} \therefore T^{\text{worst}}(n) &= \max_I T(I) \\ &= n + T^{\text{worst}}(n-1) \\ &= n + (n-1) + T^{\text{worst}}(n-2) \\ &= \dots \\ &= \frac{n(n+1)}{2} \in O(n^2). \end{aligned}$$

Average-case:

$$\begin{aligned} T^{\text{avg}}(n) &= \frac{1}{|\mathcal{X}_n|} \sum_{I \in \mathcal{X}_n} T(I) \\ &= \frac{1}{n} \sum_{k=0}^{n-1} \frac{1}{n!} \sum_{\pi} T(\pi, k). \end{aligned}$$

We will do

- ① Creating a random quick select;
- ② Analyze exp. runtime; then
- ③ Argue that this implies bound on avg-case of quick-select.

RANDOMIZED QUICK-SELECT

Consider the "randomized quick-select" algorithm:

```
quick-select(A, k)
  // A: array of size n
  // k: integer s.t. 0 <= k < n
  1. p ← random(n) // key step
  2. i ← partition(A, p)
  3. if i = k then
  4.   return A[i]
  5. else if i > k then
  6.   return quick-select(A[0, ..., i-1], k)
  7. else if i < k then
  8.   return quick-select(A[i+1, i+2, ..., n-1], k-i+1)
```

Then, what is $P(\text{pivot-index} = i)$?

\Rightarrow pivot-value is equally likely to be any of $A[0], \dots, A[n-1]$

$\Rightarrow \therefore$ pivot-index is equally likely to be any of $0, \dots, n-1$

$$\Rightarrow P(\text{pivot index} = i) = \frac{1}{n}$$

We claim $T^{\text{exp}}(n) \in O(n)$.

Proof. Recall that

$$T^{\text{exp}}(n) = \max_I \sum_{\substack{\text{random} \\ \text{outcomes } R}} P(R) \cdot T(I, R)$$

In particular, note that

$$T(I, R) = \begin{cases} n + T(A[0, \dots, i-1], k, R') & i > k \\ n + T(\text{right subarray}, k-i-1, R'), & i < k \\ n & i = k \end{cases}$$

(we count # of comparisons)

Then,

$$\begin{aligned} \sum_R P(R) T(I, R) &= \sum_R P(R) T(A, k, R) \\ &= \sum_{(i, R')} \underbrace{P(i)}_{\frac{1}{n}} \underbrace{P(R')}_{\sum_{R'} P(R')} \underbrace{T(A, k, R')}_{\begin{cases} n + \dots & i > k \\ \dots & i < k \\ 0 & i = k \end{cases}} \\ &= n + \frac{1}{n} \sum_{i=0}^{n-1} \left\{ \begin{array}{l} \sum_{R'} P(R') T(A_2, k, R') \\ \sum_{R'} P(R') T(A_r, k-i-1, R') \\ 0 \end{array} \right\}_{\begin{array}{l} i > k \\ i < k \\ i = k \end{array}} \\ &\leq n + \frac{1}{n} \sum_{i=0}^{n-1} \left\{ \begin{array}{l} \max_{k'} \max_{A_2} \sum_{R'} P(R') T(A_2, k', R') \\ \max_{k'} \max_{A_r} \sum_{R'} P(R') T(A_r, k', R') \\ 0 \end{array} \right\}_{\begin{array}{l} i > k \\ i < k \\ i = k \end{array}} \\ &= n + \frac{1}{n} \sum_{i=0}^{n-1} \left\{ \begin{array}{l} T^{\text{exp}}(i) \\ T^{\text{exp}}(n-i-1) \\ 0 \end{array} \right\}_{\begin{array}{l} i > k \\ i < k \\ i = k \end{array}} \\ &\leq n + \frac{1}{n} \sum_{i=0}^{n-1} \max_i \{ T^{\text{exp}}(i), T^{\text{exp}}(n-i-1) \} \end{aligned}$$

Then, we show $T^{\text{exp}}(n) \in 8n$.

Proof. By induction.

$$n=1: \text{Comp} = 0 < 8(1) = 8.$$

Step:

$$T(n) \leq n + \frac{1}{n} \sum_{i=0}^{n-1} \max_i \{ 8i, 8(n-i-1) \}$$

bad	good	bad
$n/4$	$3n/4$	

$$\begin{aligned} &\leq n + \frac{1}{n} \sum_{i \text{ good}} 8\left(\frac{3}{4}n\right) + \frac{1}{n} \sum_{i \text{ bad}} 8n \\ &= n + \frac{1}{n} \cdot \frac{n}{2}(6n) + \frac{1}{n} \cdot \frac{n}{2}(8n) \\ &= n + 3n + 4n = 8n. \quad \blacksquare \end{aligned}$$

$\therefore T^{\text{exp}}(n) \in O(n)$. \blacksquare

QUICK-SORT

Pseudocode:

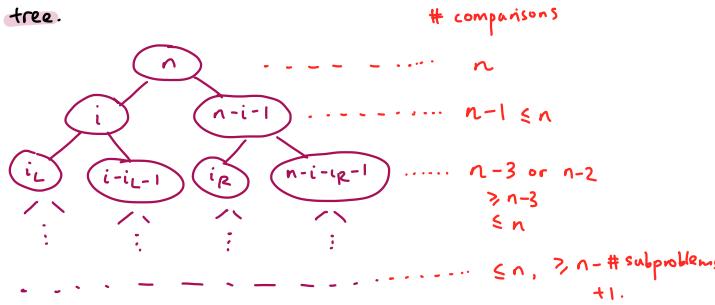
```
quick-sort(A)
  // A: array of size n
  1. if n ≤ 1 then return
  2. p ← choose-pivot(A)
  3. i ← partition(A, p)
  4. quick-sort(A[0, ..., i-1])
  5. quick-sort(A[i+1, ..., n-1])
```

There are better implementations of this idea.

Runtime:

$$T(n) = n + T(i) + T(n-i-1).$$

But there's a simpler method: the recursion tree.



Thus

$$\# \text{ comparisons} \leq n \cdot \# \text{ layers} = n \cdot \text{height of the recursion tree.}$$

So what can we say about the recursion tree's height?

① Worst case: height = $\Theta(n)$

- * tight if array is sorted.
- we need n recursions.

Thus run-time $\in \Theta(n^2)$.

② Best case: if the pivot-index is $\sim \frac{n}{2}$ always.



Thus run-time $\in \Theta(n \log n)$.

We can improve QuickSort's run-time by

① Not passing sub-arrays, but instead passing "boundaries";

② Stopping recursion early;

- stop recursing when array of subproblem is ≤ 10
- then use insertion-sort to sort the array
- which has $\Theta(n)$ best-case run time if the array is (almost) sorted

③ Avoid recursions;

④ Reduce auxiliary space;

- in the worst case (currently), the auxiliary space is $|S| \in \Theta(n)$.
- but if we put the bigger subproblem on the stack, the space can be reduced to $|S| \in O(\log n)$.

⑤ Choose the pivot-index efficiently;

- don't let $p=n-1$. (run-time = $\Theta(n^2)$)
- use "median-of-3": use the median of $\{A[0], A[\frac{n}{2}], A[n-1]\}$ as the pivot-value

AVERAGE-CASE QUICK-SORT

① We accomplish this via randomization of the algorithm:

```
RandQS(A)
  // A: array of size n
  1. if n ≤ 1 then return
  2. p ← random(n)
  3. i ← partition(A, p)
  4. quick-sort(A[0, ..., i-1])
  5. quick-sort(A[i+1, ..., n-1])
```

② Then,

$$T^{\exp}(n) = \exp \# \text{ of comparisons of the algo.}$$

③ Note that

$$T(A, R) = n + T(\text{left subarray, } R') + T(\text{right subarray, } R').$$

↑ instance ↑ outcomes
size i size $n-i$

④ Hence

$$\begin{aligned} \sum_{R} P(R) T(A, R) &= n + \frac{1}{n} \sum_{i=0}^{n-1} T^{\exp}(i) + \frac{1}{n} \sum_{i=0}^{n-1} T^{\exp}(n-i-1) \\ &= n + \frac{2}{n} \sum_{i=2}^{n-1} T^{\exp}(i), \end{aligned}$$

and so

$$T(n) = \begin{cases} 0, & n \leq 1 \\ n + \frac{2}{n} \sum_{i=2}^{n-1} T(i), & \text{otherwise.} \end{cases}$$

⑤ We claim $T(n) \in O(n \log n)$, so the expected run-time of RandQS is in $O(n \log n)$, and so the average-case run time of QS is in $O(n \log n)$.

Proof. Specifically, we prove $T(n) \leq 2 \cdot n \cdot \ln(n)$.

By induction. Base ($n=1$) is trivial.

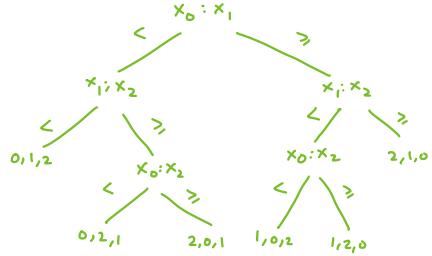
$$\begin{aligned} \text{Step: } T(n) &\leq n + \frac{2}{n} \sum_{i=2}^{n-1} 2i \ln(i) \\ &= n + \frac{4}{n} \sum_{i=2}^{n-1} i \ln(i) \\ &\leq n + \frac{4}{n} \int_2^n x \ln x \, dx \quad (\text{since } x \ln x \text{ is increasing}) \\ &\leq n + \frac{4}{n} \left(\frac{1}{2} n^2 \ln(n) - \frac{1}{4} n^2 \right) \\ &= n + 2n \ln(n) - n \\ &= 2n \ln(n). \quad \blacksquare \end{aligned}$$

ANY COMPARISON BASED SORTING ALGORITHM USES $\Omega(n \log n)$ KEY-COMPARISONS IN THE WORST-CASE

 As above.

Proof. Fix an arbitrary comparison-based sorting algorithm A , and consider how it sorts an instance of size n .

Since A uses only key-comparisons, we can express it as a decision tree T .



A decision tree for an algo to sort 3 elements x_0, x_1, x_2 with ≤ 3 key-comparisons.

For each sorting perm π , let I_π be an instance that has distinct items and sorting perm π (ie $I_\pi = \pi^{-1}$).

Executing A on I_{π} leads to a leaf that stores π .

Note that no two sorting perms can lead to the same leaf, as otherwise the output would be incorrect for one (as the items are distinct).

We then have $n!$ sorting perms, and hence at least $n!$ leaves that are reached for some π .

Let h be the largest layer-number of a leaf reached by some I_{π} .

Since γ is binary, for any $0 \leq l \leq h$ there are at most 2^l leaves in layers $0, \dots, l$.

Therefore $2^n > n!$, or

$$\begin{aligned}
 n > \log(n!) &= \log(n(n-1) \cdots 1) \\
 &= \log(n) + \log(n-1) + \cdots + \log(1) \\
 &> \log(n) + \log(n-1) + \cdots + \log(\lceil \frac{n}{2} \rceil) \\
 &> \log\left(\frac{n}{2}\right) + \log\left(\frac{n}{2}\right) + \cdots + \log\left(\frac{n}{2}\right) \\
 &\geq \frac{n}{2} \log\left(\frac{n}{2}\right) \\
 &= \frac{n}{2} \log(n) - \frac{n}{2} \in \mathcal{L}(n \log n).
 \end{aligned}$$

Finally, consider the sorting perm π that has its leaf on layer h .

Executing algorithm A on Π hence takes $h \in \Omega(n \log n)$ key-comparisons, so the worst-case bound holds. \square

SORTING INTEGERS

BUCKET-SORT

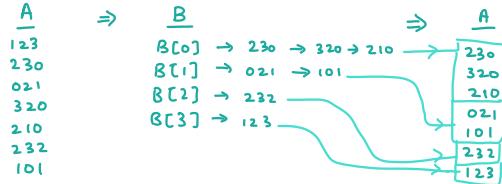
Bucket-sort can be used to sort a collection of integers by a specific "position" of digit.

e.g. the last digit 12345

Pseudocode:

```
bucket-sort(A, n < A.size, l=0, r=n-1, d)
// A: array of size  $\geq n$  with numbers with m digits in {0,...,R-1}
// d: 1st elem
// output: A[l...r] is sorted by "d" digit.
1. initialize an array B[0...R-1] of empty lists
2. for i < l to r do
3.   append A[i] at the end of B[dth digit of A[i]] // move array-items to buckets
4. i < l
5. for j < 0 to R-1 do
6.   while B[j] is not empty do // move bucket-items to array
7.     move first element of B[j] to A[i]
8.     i++
9. 
```

e.g.



See that

- ① Run-time = $\Theta(n+R)$; and
- ② Auxiliary space = $\Theta(n+R)$.

MSD-RADIX-SORT

MSD (Most Significant Digit) radix sort can be used to sort multi-digit numbers.

Pseudocode:

```
MSD-radix-sort(A, n < A.size, l=0, r=n-1, d=1)
// A: array of size  $\geq n$ , contains numbers with m digits in {0,...,R-1}, m, R are global variables
// l, r: range (ie A[l...r]) we wish to sort
// d: digit we wish to sort by
1. if l < r then
2.   bucket-sort(A, n, l, r, d)
3.   if d < m then
4.     // find sub-arrays that have the same dth digit and recurse
5.     int l' < l
6.     while l' < r do
7.       int r' < l'
8.       while r' < r & dth digit of A[r'+1] = dth digit of A[l'] do r'++
9.       MSD-radix-sort(A, n, l', r', d+1)
10.      l' <= r' + 1
11. 
```

Note that

- ① run-time = $\Theta(mRn)$; &
- ② auxiliary space = $\Theta(n+R+m)$.

LSD-RADIX-SORT

"LSD-Radix-Sort" is a better way of sorting multi-digit numbers (than MSD) because

- ① it has a faster runtime;
- ② it uses less auxiliary space; and
- ③ it uses no recursion.

Pseudocode:

```
LSD-radix-sort(A, n < A.size)
// A: array of size n, contains m-digit radix-R numbers, m, R are global
1. for d=m down to 1 do
2.   bucket-sort(A, d)
```

Clearly

- ① run-time = $\Theta(m(n+R))$; &
- ② auxiliary space = $\Theta(n+R)$.

Chapter 4: ADT Dictionaries

DICTIONARIES

💡 Dictionaries store key-value pairs, or "KVP".



In particular,

- ① `search(key)`: return the KVP for this key.
- ② `insert(key, value)`: add the KVP to the dictionary.
 - key is distinct from existing keys
- ③ `delete(key)`: remove KVP with this key.

💡 Assumptions:

- ① we assume all keys are distinct.
- ② we also assume keys can be compared.

💡 Implementations:

① **Unsorted list**

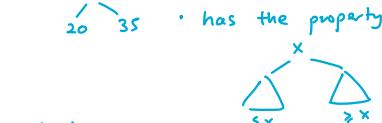
- fast insert, slow search, slow delete

② **Sorted array**

- fast search, slow insert, slow delete

③ **Binary search tree / BST**

- we only show the keys (implied each node is a KVP)
- has the property



In particular,

insert, delete, search $\in O(\text{height of tree})$.

* Height is $\Theta(n)$ is worst case, but typically much better ($O(\log n)$)

LAZY DELETION

💡 Consider a sorted array:

`delete(20)`:

10 20 30 40 60 70 80

- usually takes $\Theta(n)$ time to "backtrack" all other elements
- but we can avoid this if we instead just mark the box as "isDeleted".
- thus run-time (search) = $O(\log n)$.
- but we don't get any space back!
- however, we can occasionally "clean up":
 - create a new initialization; &
 - move all items into the new array if they were "real".

· clean-up takes $O(n \times \text{insert})$.

But we can frequently do better; in this example, we can perform clean-up in $O(n)$ time.

Then, the amortized time is $O(\text{insert} + \text{delete})$ for doing a deletion, and sometimes better.

💡 Why do we do **lazy deletion**?

- ① It is simpler.
- ② It might be faster.
 - for sorted arrays: $\Theta(n)$ worst-case for delete but $\Theta(\log n)$ amortized time for lazy deletion

③ It sometimes is required.

💡 Why not?

- ① It wastes space.
 - we have to allocate space to store the "isDeleted" flag.
- ② Occasionally deletion is very slow.

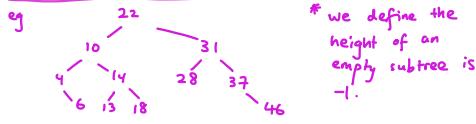
AVL-TREES

Goal: $O(\log n)$ worst case time.

Structural condition:

For any node z , we have

$$\text{balance}(z) = |\text{height}(z.\text{left}) - \text{height}(z.\text{right})| \leq 1.$$



- notice the structural property holds for every node.

We also store the height of the subtree at every node, or we can also store the balance.

(usually we will store the height).

$h \in O(\log n)$

Let h be the height of an AVL tree with n nodes.

Then necessarily $h \in O(\log n)$. # nodes

height 0: 1

height 1: or 2

height 2: to preserve the structural property! 4
has $h=1$

height 3: 7
 $h=2$ $h=1$

In general, let $N(h)$ = smallest # of nodes of

height h .

$N(h) = N(h-1) + N(h-2) + 1$.

$\therefore N(h) = N(h-1) + N(h-2) + 1$.

Thus

h	0	1	2	3	4	5	...
$N(h)$	1	2	4	7	12	20	...
$N(h)+1$	2	3	5	8	13	21	

We see $N(h)+1$ is the Fibonacci numbers!

(ie $F(0)=0$, $F(1)=1$, $F(i)=F(i-1)+F(i-2)$ $\forall i \geq 2$.)

By a easy induction proof, we can show

$$N(h)+1 = F(h+3).$$

In particular, we know

$$F(i) = \frac{1}{\sqrt{5}} \phi^i + \Theta(1), \quad \phi = \text{the golden ratio}.$$

Therefore

$$N(h) = \frac{1}{\sqrt{5}} \phi^{h+3} + \Theta(1)$$

Hence, for any AVL tree of height h ,

$$\# \text{nodes} = n \geq N(h) \approx \frac{1}{\sqrt{5}} \phi^{h+3},$$

and so

$$n \approx \log_{\phi} (\sqrt{5} n) - 3 \in O(\log n).$$

AVL OPERATIONS

INSERT, PART 1

First, we call `BST::insert`, and then rebalance the ancestors of z .

`AVL::insert(k, v)`

1. $z \leftarrow \text{BST}::\text{insert}(k, v)$
2. while z is not NIL do
3. if $|z.\text{left}.\text{height} - z.\text{right}.\text{height}| > 1$ then
4. let $y \leftarrow$ taller child of z
5. let $x \leftarrow$ taller child of y
6. $z \leftarrow \text{restructure}(x, y, z)$
7. break
8. $z.\text{height} \leftarrow 1 + \max\{z.\text{left}.\text{height}, z.\text{right}.\text{height}\}$ // we alias this as "setHeightFromSubtrees"
9. $z \leftarrow z.\text{parent}$

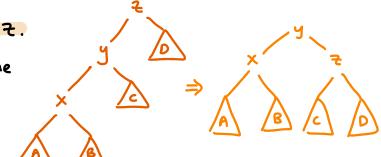
ROTATIONS [OF BSTS]

Let T be a BST with a node z that has a child y and a grandchild x . Then, a "rotation at z with respect to y & x " is a restructuring of T such that the result is again a BST, and sub-tree references have been changed only at x, y, z .

For restoring balances at AVL-trees, we want the four rotations that make the median of x, y, z the new root of the subtree:

① Right rotation; $x < y < z$.

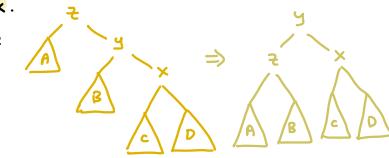
We want y to become the root.
pseudocode similar to left rotation; see below



② Left rotation; $z < y < x$.

We want y to be the root.

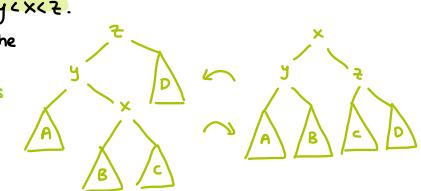
- rotate-left(z)
1. $y \leftarrow z.\text{right}$
 2. $z.\text{right} \leftarrow y.\text{left}$
 3. $y.\text{left} \leftarrow z$
 4. setHeightFromSubtrees(z)
 5. setHeightFromSubtrees(y)
 6. return y



③ Double-right rotation; $y < x < z$.

We want x to be the root.

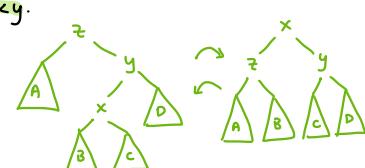
* can also be written as single-left rotation at y and then a single-right rotation at z (double-rotation)



④ Double-left rotation; $z < x < y$.

We want x to be the root.

* can also be written as single-right rotation at y and then a single-left rotation at z (double-rotation)



INSERT, PART 2

We now give the implementation of `restructure` in the `insert` function from the first part;

`restructure(x, y, z)`

- // node x has parent y & grandparent z
1. if $y = z.\text{left}$ & $x = y.\text{left}$ // $x-y-z$
 2. return rotate-right(z) // right rotation
 3. else if $y = z.\text{left}$ & $x = y.\text{right}$ // $y-z-x$
 4. $z.\text{left} \leftarrow$ rotate-left(y)
 5. return rotate-right(z) // double-right rotation
 6. else if $y = z.\text{right}$ & $x = y.\text{left}$ // $x-z-y$
 7. $z.\text{right} \leftarrow$ rotate-right(y)
 8. return rotate-left(z) // double-left rotation
 9. else // $x-y-z$
 10. return rotate-left(z) // left rotation

DELETE

Pseudocode:

```

AVL::delete(u)
1. z ∈ BST::delete(u) // z is the parent of the
   BST node that was
   removed
2. while z is not NIL do
   if |z.left.height - z.right.height| > 1 then
   let y = taller child of z
   let x = taller child of y
   // (break ties to prefer single rotation)
   6. z ← restructure(x,y,z)
   // do not break — continue up the path &
   rotate if needed.
7. setHeightFromSubtrees(z)
8. z ← z.parent

```

RUNTIME

Both insert & delete have run-time $O(\log n)$.

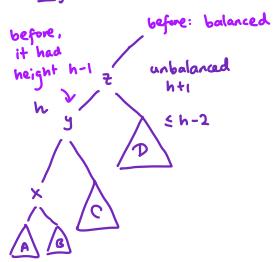
Why?
- height of tree = $O(\log n)$
- rotations take $O(1)$ time
- so delete takes $O(\log n)$.
- tight if we insert/delete at lowest level & never rotate.

CORRECTNESS [FOR INSERTION]

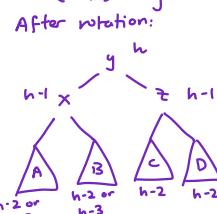
If we restructure at z during an insertion, then

- ① the subtree is balanced; and
- ② the subtree has now the height that it had before the insertion.

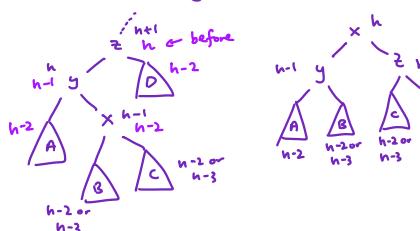
Proof.



Case 1: "x" is the left node
- x has height $h-2$ before
- C has height $h-2$



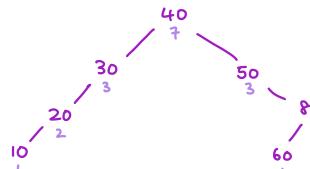
Case 2: x is the right child.



SCAPEGOAT TREES

"Scapegoat trees" are BSTs such that

- ① Each node v stores the size of its subtree; and
- ② $v.size \leq \alpha \cdot v.parent.size$ for all nodes that are not the root.



* Example, with $\alpha = \frac{2}{3}$. Verify that $p.size \geq \frac{3}{2}v.size$ for any parent p of any node v .

HEIGHT = $O(\log n)$

Any scapegoat tree has height $O(\log n)$.

Proof: At any leaf v (which has size 1), the parent has size $\geq \frac{1}{\alpha}$ by defn of a scapegoat tree.

Repeating this argument, the grand-parent has size $\geq (\frac{1}{\alpha})^2$, and so on.

As the root has size n , it follows that $(\frac{1}{\alpha})^d \leq n$,

where $d = \max \text{ depth of a leaf} = \text{height}$. Thus

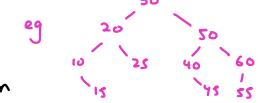
$$d \leq \log_{(\frac{1}{\alpha})}(n) \in O(\log n).$$

PERFECTLY BALANCED BST

A "perfectly balanced BST" is one where for any node v , we have

$$|v.left.size - v.right.size| \leq 1.$$

Given any n -node BST T , we can build a perfectly balanced BST with the same KVPs in $O(n)$ time.



INSERT

Pseudocode:

```

ScapegoatTree::insert(k,v)
1. z ∈ BST::insert(k,v)
2. S ← stack initialized w/ z
3. while p ≠ z.parent ≠ NIL do
4.   p.size++
5.   S.push(p)
6.   z = p
7. while S.size ≥ 2 do
8.   p = S.pop()
9.   if p.size <  $\frac{1}{\alpha} \max\{p.left.size, p.right.size\}$  then
10.    completely rebuild the subtree rooted
        at p as a perfectly balanced BST
11.    break

```

insert has worst case runtime $O(n)$, but amortized runtime $O(\log n)$.

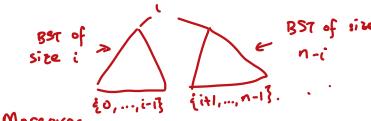
RANDOMIZED BUILT BST HAS EXPECTED HEIGHT $\Theta(\log n)$

\exists_1 A randomized built BST has expected height $\Theta(\log n)$.

\exists_2 Here, "randomly built" means we take a permutation (randomly) and insert items in the order of said permutation.

Proof. First, the first item of the perm π becomes the root.

Then, once we know the root i :



Moreover, $P[\text{first item of } \pi \text{ is } i] = \frac{1}{n}$.

Thus exp. height of the tree w/ perm π is

$$H(\pi) = 1 + \max\{H(\pi_L), H(\pi_R)\}.$$

Now, define $Y(\pi) = 2^{H(\pi)}$, & $Y(n) = E[Y(\pi)]$.

Then

$$\begin{aligned} E[H(\pi)] &= E[\log(Y(\pi))] \\ &\leq \log(E[Y(\pi)]). \quad (\because \log \text{ is concave}) \end{aligned}$$

We will show $E[Y(\pi)] \leq (n+1)^3$, which implies

$$H(n) \leq \log((n+1)^3) = 3\log(n+1).$$

So, let's do so. We note

$$\begin{aligned} Y(\pi) &= 2^{1 + \max\{H(\pi_L), H(\pi_R)\}} \\ &= 2 \cdot \max\left\{2^{H(\pi_L)}, 2^{H(\pi_R)}\right\} \\ &\leq 2 \left(2^{H(\pi_L)} + 2^{H(\pi_R)}\right) \\ &= 2(Y(\pi_L) + Y(\pi_R)), \end{aligned}$$

and so

$$\begin{aligned} Y(n) &= \frac{1}{n} \sum_{i=0}^{n-1} (2Y(i) + 2Y(n-i-1)) \\ &= \frac{4}{n} \sum_{i=0}^{n-1} Y(i). \end{aligned}$$

Now, we show $Y_n \leq (n+1)^3$.

Base: $n=1 \Rightarrow \text{height}=0 \Rightarrow Y(1)=2^0=1 \leq 8$.
 $n=0 \Rightarrow \text{height}=-1 \Rightarrow Y(0)=2^{-1}=\frac{1}{2} \leq 1$.

Step: $Y(n) = \frac{4}{n} \sum_{i=0}^{n-1} Y(i)^3$
 $= \frac{4}{n} \sum_{i=1}^{n-1} i^3 = \frac{4}{n} \cdot \frac{n^2(n+1)^2}{4}$
 $= n(n+1)^2 \leq (n+1)^3,$

as needed. \square

EXPECTED VS AVE

\exists_1 We know

$$E[\text{height of BST}] \in \Theta(\log n).$$

But

avg of height of BST is not $\Theta(\log n)$!

\exists_2 We can show the average height of a BST is in $\Theta(\sqrt{n})$.

\exists_3 Intuition on why:

$P[\text{randomly built BST has shape } \triangle] = \frac{1}{n}$

But

$$\begin{aligned} \frac{\# \text{BSTs w/ shape } \triangle}{\text{total } \# \text{ of BSTs}} &= \frac{\# \text{BST on } n-1 \text{ items}}{\# \text{BST on } n \text{ items}} \\ &= \frac{C(n-1)}{C(n)} \approx \frac{1}{4} \end{aligned}$$

where $C(n)$ are the Catalan numbers.

TREAPS

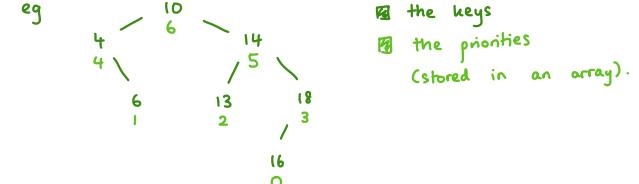
\exists_1 A "treap" is a randomized version of a BST, also called a "priority search tree".

\exists_2 A treap has the following properties:

- ① Each node has a KVP and a priority p .
- ② The treap acts like a BST wrt to the KVPs; &
* everyone in left \leq root \leq everyone in right
- ③ The treap acts like a heap wrt to the priority.

* each node stores the item with max priority among all nodes in its subtree

eg



TREAP INSERTION

\exists_1 When inserting into the treap, we just

- ① BST insert; and
- ② do "fix-up" to restore the heap-order property for the priorities.

treap::fix-up-with-rotations(z)

// z: node whose priority may have increased
1. while (y < z.parent is not NIL & z.priority > y.priority) do
2. if z is the left child of y then rotate-right(y)
3. else rotate-left(y)

treap::insert(k, v)

1. $n \in P[\text{size}]$
2. $z \leftarrow \text{BST}::\text{insert}(k, v)$; $n++$ // z is the leaf where k is now stored

3. $p \leftarrow \text{random}(n)$
4. if $p < n-1$ then // change priority of other node
5. $z' \leftarrow P[p]$, $z'.\text{priority} \leftarrow n-1$, $P[n-1] \leftarrow z'$
6. fix-up-with-rotations(z')
7. $z.\text{priority} \leftarrow p$, $P[p] \leftarrow z$
8. fix-up-with-rotations(z)

RUN-TIME / SPACE

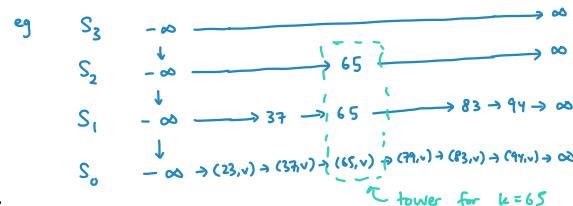
\exists_1 Note that treaps are

① BSTs with expected height $\Theta(\log n)$, and so the expected runtime of all operations is $\Theta(\log n)$;

② have a large space overhead, since we store the BST, and parent-references & priorities for each node.

SKIP LISTS

- A "skip list" is a hierarchy S of ordered linked lists (levels) S_0, \dots, S_h , such that
- ① Each list S_i contains the special keys $-\infty$ & ∞ , called "sentinels";
 - ② S_0 contains the KVP of S in non-decreasing order, and the other lists store only keys;
 - ③ Each list is a subsequence of the previous one, ie
 $S_0 \supseteq S_1 \supseteq \dots \supseteq S_h$; &
 - ④ S_h contains only the sentinels.



- \exists_2 Notes:
- ① A "node" in a skip list is any node in the linked lists of the hierarchy;
 - ② Each node has two references:
 - "after" — points to the next node in the LL
 - "below" — points to the copy of the node in S_{i-1}
 - ③ The "tower" of a key k is the set of all nodes that contain k ;
 - ④ The "height" of a tower is the maximum index i such that the tower includes a node in S_i ;
 - ⑤ The "height" of the skip list is the maximum height of a tower, which is equal to h .

getPredecessors(k)

- \exists_1 'getPredecessors(k)' is a helper routine for insert & delete.

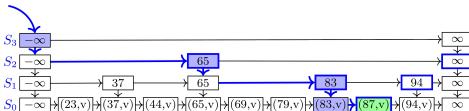
Algorithm 5.3: skipList::getPredecessors(k)

```

1 p ← root
2 P ← stack of nodes, initially containing p
3 while p.below ≠ NIL do
4   p ← p.below                                // drop down
5   while p.after.key < k do
6     p ← p.after                                // step forward
7   P.push(p)
8 return P

```

Example:



SEARCH

- 'search' is very simple; it just uses getPredecessors.

Algorithm 5.4: skipList::search(k)

```

1 P ← getPredecessors(k)                      // predecessor of k in  $S_0$ 
2 p0 ← P.top()
3 if p0.after.key = k then
4   | return p0.after
5 else
6   | return "not found, but would be after p0"

```

INSERT

- \exists_1 For 'insert', we first call 'getPredecessors', which tells us which node would precede the inserted KVP at each S_i .

- \exists_2 But, we determine whether k should be in S_i randomly; in particular,

$$P(\text{tower of key } k \text{ has height } \geq i) = \frac{1}{2^i}.$$

Algorithm 5.5: skipList::insert(k, v)

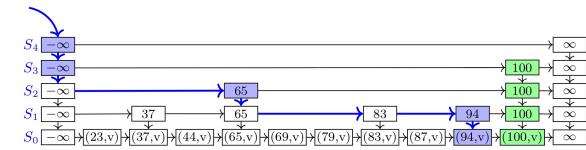
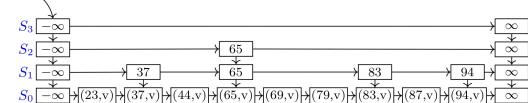
```

1 for (i ← 0; random(2) = 1;) do i++           // random tower height
2 for (h ← 0, p ← root.below; p ≠ NIL; p ← p.below) do h++    // compute height
3 while (i ≥ h) do
4   | create a new sentinel-only list and link it to the previous root-list appropriately
5   | root ← leftmost node of this new list
6   | h++
7   // Actual insertion
8   P ← getPredecessors(k)
9   z.below ← new node with (k, v), inserted after p           // insert (k, v) in  $S_0$ 
10  while i > 0 do
11    | p ← P.pop()                                         // insert k in  $S_1, \dots, S_i$ 
12    | z ← new node with k, inserted after p
13    | z.below ← z.below; z.below ← z
14    | i ← i - 1

```

Example:

Insertion of key = 100, with the determined tower height = 3.



DELETE

- \exists_1 To delete in a skip list, we find the key (which gives the stack of predecessors) and then remove it from all lists that it was in.

- \exists_2 We also "clean-up" the stack: if deleting a key results in multiple lists that store only sentinels, then delete all but one of them.

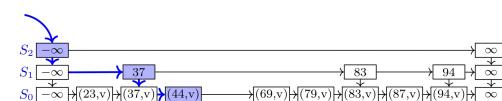
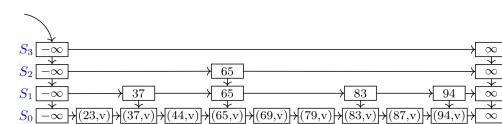
Algorithm 5.6: skipList::delete(k)

```

1 P ← getPredecessors(k)
2 while P is non-empty do
3   | p ← P.pop()                               // p could be predecessor of k in some list
4   | if p.after.key = k then remove p.after from the list
5   | else break                                 // no more copies of k
6   | p ← root                                  // clean up duplicate empty lists
7   while p.below ≠ NIL and p.below.after is the ∞-sentinel do
8   | remove the list that begins with p.below

```

Example:



* deleting the KVP w/ key=65.

$$E[\text{len}(S_i)] = \frac{n}{2^i}$$

In a skip list, the length of a list S_i is $\frac{n}{2^i}$.

Proof. Let X_k be the rv that denotes the height of the tower w/ key k .

Let $I_{i,k}$ be an indicator variable that is 1 if $X_k \geq i$ (ie list S_i contains key k) & 0 otherwise.

Then

$$|S_i| = \sum_{\text{key } k} I_{i,k}$$

and so

$$\begin{aligned} E[|S_i|] &= \sum_{\text{key } k} E[I_{i,k}] \\ &= \sum_{\text{key } k} P(X_k \geq i) \\ &= \sum_{\text{key } k} \frac{1}{2^i} \\ &= \frac{n}{2^i}. \quad \square \end{aligned}$$

$$E[\text{height of SL}] \leq \log n + O(1)$$

The expected height of a skip list is at most $\log(n) + O(1)$.

Proof. Let I_i be an ind var s.t. $I_i = 1$ if $|S_i| \geq 1$ & 0 otherwise.

Recall $\text{height}(SL) = h$, where the lists are S_0, \dots, S_h .

Since S_0, \dots, S_{h-1} all contain keys, thus

$$h = \sum_{i \geq 0} I_i.$$

Then, note that by $\cup_{i \geq 0} I_i \leq |S_i|$,

so

$$E[I_i] \leq \min\{\frac{1}{2}, \frac{n}{2^i}\}.$$

If $i \approx \log n$ then $\frac{n}{2^i} \approx \frac{1}{2}$, so we use this to "break up" the sum.

In particular, notice

$$\begin{aligned} E[h] &= \sum_{i \geq 0} E[I_i] \leq \sum_{i=0}^{\lceil \log n \rceil - 1} E[I_i] + \sum_{i \geq \lceil \log n \rceil} E[|S_i|] \\ &\leq \sum_{i=0}^{\lceil \log n \rceil - 1} (1) + \sum_{i \geq \lceil \log n \rceil} \frac{n}{2^i} \\ &\leq \lceil \log n \rceil + \sum_{j \geq 0} \frac{2}{2^{j+\lceil \log n \rceil}} \\ &\leq \log n + 1 + \sum_{j \geq 0} \frac{1}{2^j} \\ &= 3 + \log n, \end{aligned}$$

as needed. \square

EXPECTED SPACE = $\Theta(n)$

The expected space of a skip list is in $\Theta(n)$.

In particular, the expected number of nodes is $2n + o(n)$.

Proof. Each list S_i has $|S_i|$ nodes that store keys.

Hence expected # of nodes with keys is

$$E\left[\sum_{i \geq 0} |S_i|\right] = \sum_{i \geq 0} \frac{n}{2^i} = n \sum_{i \geq 0} \frac{1}{2^i} = 2n.$$

There are $2h+2$ nodes that do not store keys (the sentinels on each list S_0, \dots, S_h), but we have $E[h] \in \log(n) + O(1) = o(n)$,

& so the bound holds. \square

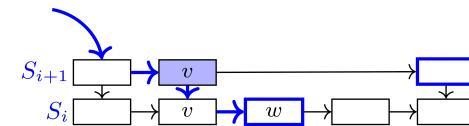
getPredecessors: $E(\text{forward steps of } S_i) \leq 1$

Proof. During 'getPredecessors', the expected number of forward-steps within list S_i is at most 1.

Proof. Let v be the leftmost node in list S_i that we visited during the search. If $i=h$ (the topmost list) then we do not step forward at all and are done, so assume $i < h$ & we reached v by dropping down from list S_{i+1} .

Let w be the item after v in S_i . Consider the prob. we step forward from v to w : if w also exists in list S_{i+1} , then we compare (before dropping down to v) the search key K with $w.\text{key}$.

So, we must have had $K \leq w.\text{key}$, else we would not have dropped down from v . Thus in list S_i we will immediately drop down.



Taking the contrapositive, if we step forward from v in list S_i , then the next node w in S_i did not exist in S_{i+1} .

In other words, the tower of w had height exactly i . The probability of this is $\frac{1}{2}$, because the decision to expand the tower of w into the list above was based on a "coin flip".

So we step forward from v w/ prob $< \frac{1}{2}$.

Repeating this argument, we step forward from w with prob $< \frac{1}{2}$, presuming we arrived at w in the first place, so the probability of this happening is at most $\frac{1}{4}$. Repeating, we see the prob of stepping forward i times is $< \frac{1}{2^i}$.

Thus

$$\begin{aligned} E[\# \text{ of forward steps}] &= \sum_{k \geq 1} P(\# \text{ of forward-steps is } \geq k) \\ &= \sum_{k \geq 1} \frac{1}{2^k} \leq 1. \quad \square \end{aligned}$$

EXPECTED RUN-TIME OF SEARCH / INSERT / DELETE IS $O(\log n)$

As above.

Proof. First, exp. run-time for getPredecessors is

$$O(E[h + \sum_{i \geq 0} F_i]),$$

where $F_i = \# \text{ of forward steps on level } S_i \in O(\log n)$.

By the previous results, the exp. time for

getPredecessors is in $O(\log n)$.

Once the predecessors are found, all other operations take $O(h)$ time.

This has exp $O(\log n)$. \square

BIASED SEARCH REQUESTS

STATIC SCENARIO

B1 In the "static scenario", we know beforehand how frequently a key is going to be accessed.

key	A	B	C	D	E
access-freq	2	8	1	10	5
access-prob	2/26	8/26	1/26	10/26	5/26

- B2** Terms:
- ① "Access frequency" — amount of times key is accessed
 - ② "Access probability" — proportion of accesses for a key

B3 In particular, we want to find the optimum assignment of keys to locations; ie the assignment that minimizes

$$\text{exp. access cost} = \sum_{\text{key } k} p(\text{want to access } k) \cdot (\text{cost of accessing } k)$$

DYNAMIC SCENARIO

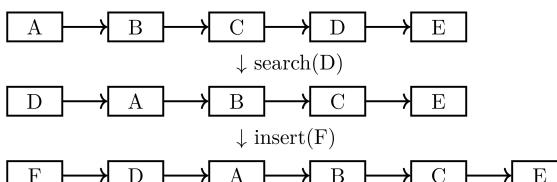
B1 Here, we do not know how frequently a key is going to be accessed, so we cannot hope to build the best-possible data structure.

B2 However, we can still change the data structure when we have accesses, to bring those items that were recently accessed to a place where the next access will be cheap.

- if we access an item, it is fairly likely we will access it again soon
- "temporal locality"

MOVE-TO-FRONT / MTF HEURISTIC IN A LIST

B1 The MTF heuristic involves moving the most recently accessed item in an unsorted list to the front.



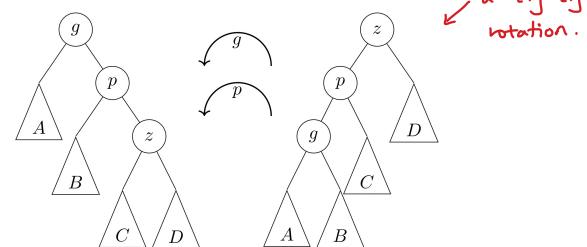
B2 The heuristic is "2-competitive": it takes at most twice as many comparisons that would have been taken using the optimal static ordering.

SPLAY TREES

B1 Splay trees are BSTs where after every operation, we apply zig-zag or zig-zig rotations (and perhaps one single rotation at the root) so the accessed item is at the root.

B2 Here,

- ① "zig-zag rotations" are just double rotations; we do these if the "z-p-g" path contains a left & right child; &
- ② "zig-zig rotations" are applied if the "z-p-g" path contains two left or two right children.



INSERT

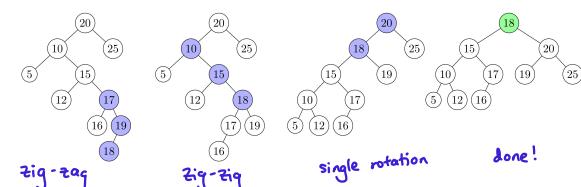
B1 Pseudocode:

```

Algorithm 5.8: splayTree::insert(k, v)
1  $z \leftarrow BST::insert(k, v)$ 
2 while ( $z$  has a parent  $p$  and a grand-parent  $g$ ) do
3   if  $\begin{cases} z \\ p \end{cases}$  then                                // Zig-zig rotation rightward
4     |  $p \leftarrow rotate-right(g)$ ,  $z \leftarrow rotate-right(p)$ 
5   else if  $\begin{cases} z \\ p \end{cases}$  then                         // Zig-zag rotation rightward
6     |  $g.left \leftarrow rotate-left(p)$ ,  $z \leftarrow rotate-right(g)$ 
7   else if  $\begin{cases} z \\ p \end{cases}$  then                           // Zig-zag rotation leftward
8     |  $g.right \leftarrow rotate-right(p)$ ,  $z \leftarrow rotate-left(g)$ 
9   else //  $\begin{cases} z \\ p \end{cases}$                                 // Zig-zig-rotation leftward
10  |  $p \leftarrow rotate-left(g)$ ,  $z \leftarrow rotate-left(p)$ 
11 if ( $z$  has a parent  $p$ ) then
12  if ( $z = p.left$ ) then  $z \leftarrow rotate-right(p)$            // single rotation
13  else  $z \leftarrow rotate-left(p)$ 

```

B2 Example:



RUN-TIME ANALYSIS

The amortized run-time of insert is $O(\log n)$, where n is the size of the tree.

Proof: For a splay tree S , let the put. func.

$$\phi(i) = \sum_{v \in S} \log(n_v^{(i)})$$

where n_v = size of the subtree rooted at v . Clearly this is a potential func ($\phi(0)=0$) & $\phi(i) \geq 0 \Leftrightarrow n_v^{(i)} \geq 1 \forall i, v$.

Insert has three phases:

- ① BST::insert;
- ② Bringing up the node with zig-zig & zig-zag rotation;
- ③ Doing the (last) single rotation.

Lemma #1: ① increases ϕ by at most $\log n$.

Proof: Let the nodes in the path from x to the root be x_1, x_2, \dots, x_d , where $x_d = \text{root}$.

After adding x , the size of the subtrees at all x_1, \dots, x_d increase by 1, whilst it is unchanged at all other nodes.

So, the contribution to ϕ only changes at x_1, \dots, x_d , and in particular

$$n_{x_u}^{(\text{after})} = n_{x_u}^{(\text{before})} + 1 \leq n_{x_{u+1}}^{(\text{before})} \quad \forall 1 \leq u \leq d.$$

Thus

$$\begin{aligned} \Delta\phi &= \sum_v \log(n_v^{(\text{after})}) - \log(n_v^{(\text{before})}) \\ &= \log(n_{x_d}^{(\text{after})}) + \sum_{u=1}^{d-1} (\log(n_{x_u}^{(\text{after})}) - \log(n_{x_u}^{(\text{before})})) \\ &\leq 0 + \sum_{u=1}^{d-1} (\log n_{x_{u+1}}^{(\text{before})} - \log n_{x_u}^{(\text{before})}) \\ &\quad + \log n_{x_d}^{(\text{after})} - \log n_{x_d}^{(\text{before})} \\ &= \log n_{x_d}^{(\text{before})} - \log n_{x_1}^{(\text{before})} + \log n_{x_d}^{(\text{after})} - \log n_{x_1}^{(\text{before})} \\ &\leq \log n \end{aligned}$$

as needed. \square

Lemma #2: Let D_i be a zig-zag / zig-zig rotation

that moves x two levels up.

$$\text{Then } \Phi_{\text{after}(D_i)} - \Phi_{\text{before}(D_i)} \leq 3\log(n_x^{(\text{after})}) - 3\log(n_x^{(\text{before})}) - 2.$$

Proof: See TB.

Main proof: we finally show the amortized run-time of insert is $O(\log n)$.

Note $T_{\text{actual}}(\text{insert}) = 1+td$, where $d = \text{depth from root to } x$.

Let the seq of operations in insert be

$$D_i, D_{i+1}, \dots, D_{i+R}.$$

$\underbrace{}_{\text{BST::insert}}$ $\underbrace{}_{\text{zig/zig}}$

Then

$$\begin{aligned} \Delta\phi(\text{insert}) &= \Phi(i+R) - \Phi(i-1) \\ &= \sum_{j=i}^{i+R} (\Phi(j) - \Phi(j-1)) \\ &= \sum_{j=i}^{i+R} \Delta\Phi(D_j) \\ &\leq \log n + \sum_{j=i+1}^{i+R-1} (\underbrace{3\log(n_x^{(j)}) - 3\log(n_x^{(j-1)}) - 2}_{\Delta\phi \text{ for zig-zig or zig-zag}}) + \underbrace{(3\log(n_x^{(i+R)}) - 3\log(n_x^{(i+R-1)}))}_{\Delta\phi(D_{i+R})} \\ &\leq \log n + 3\log(n_x^{(i+R)}) - 3\log(n_x^{(i)}) - 2(R-1) \\ &\leq \log n + 3\log n - 2R + 2, \end{aligned}$$

and so

$$\begin{aligned} T_{\text{actual}}(\text{insert}) + \Delta\phi(\text{insert}) &\leq (1+d) + 4\log n - 2R + 2 \\ &\leq 4\log n + 4 \\ &\in O(\log n) \end{aligned}$$

as needed. \square

BINARY SEARCH REVISITED

Θ_1 : we cannot do better than binary search.

- asymptotically, comparison-based

Θ_2 : But, we can do a bit better;

- shave constant

Θ_3 : But, we can do a lot better! & drop comparison based

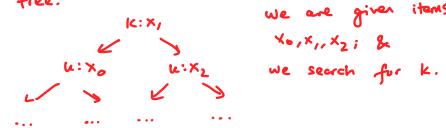
Θ_4 : we can do even better!

ANY COMPARISON-BASED SEARCH TALES $\Omega(\log n)$ TIME

Θ_1 : we want to prove a lower bound for searching.

Θ_2 : Any search that is comparison-based among n elements takes $\Omega(\log n)$ worst-case time, even if the elements are sorted.

Proof: Fix an algorithm, and look at its decision tree.



We want to show we have a lot of (accessible) leaves.

In particular, we have at least n leaves (one for each X_i).

Thus, the height is $\geq \log n \in \Omega(\log n)$. \square

We can also note

weight $\geq \lceil \log n \rceil$.

However, we also have $n+1$ "not found" leaves, corresponding to $k \in (-\infty, x_0), k \in (x_0, x_1), \dots, k \in (x_n, +\infty)$.

Θ_3 : In fact, we can show the height $\geq \lceil \log(2n+1) \rceil$.

Proof: In particular, we now show

of leaves $\geq 2n+1$

(and so height $\geq \lceil \log(2n+1) \rceil$).

To do so, we create $2n+1$ instances

$x_0 < x_1 < x_2 < \dots < x_{n-1}$.

Searching for x_i could result in the possibilities:

$x_0, x_1, x_2, \dots, x_{n-1}$

Searching between x_i :
 $(x_0, x_1), (x_1, x_2), \dots, (x_{n-2}, x_{n-1})$

Searching outside:
 $(-\infty, x_0), (x_{n-1}, +\infty)$.

Claim: no two reach the same leaf.

Proof by contradiction.

Assume $\exists k, k'$, $k \neq k'$, s.t. we reach the same leaf.

Then $k \neq k_i$ (as otherwise $k=k_i$), so it follows the leaf must be a 'not found'.

leaf reached w/
 k, k' .

Consider the above example.

In particular, we have $x_0 < k, k' < x_1$, and so on.

Since k, k' one 'not found', there must exist an x_i between a k & k' .

Thus $k < x_i < k'$

for some x_i .

Then, a search for x_i would reach the exact leaf, which is a contradiction since the leaf is a 'not found'.

Chapter 6:

Special Key Dictionaries

OPTIMIZED BIN SEARCH

💡 Normal binary search takes
 $T(n) = 2 + T(\frac{n}{2}) \approx 2\log(n)$.

But can we do better?

Pseudocode:

```
binary-search-optimized(A, n, k)
1. l=0, r=n-1, x←0
   // x is a bool var that tells us whether
   // we're in the left subarray
2. while (l < r)
3.   m ← ⌊\frac{l+r}{2}\rfloor
4.   if (A[m] < k) then r ← m+1
5.   else r=m, x←1
6. if (k < A[r]) then
7.   return "not found, bw A[r-1] and A[r]"
8. else if (x=1) or (k <= A[r]) then
9.   return "found at A[r]"
10. else return "not found, bw A[r] & A[r+1]
```

💡 We claim

- ① this algorithm terminates;
- ② this gives the correct answer;
- ③ this uses at most $\lceil \log n \rceil + 2$ comparisons (without x), which is about $\lceil \log(2n+1) \rceil + 1$ comparisons.
- ④ with the x, this uses $\leq \lceil \log(2n+1) \rceil$ comparisons, so this is optimal.

Proof. See TB.

💡 Pseudocode:

interpolation-search(A, n, k)

```
1. l=0, r=n-1
2. while (l < r)
3.   if (k < A[r] or k > A[l]) return "not found"
4.   if (k = A[r]) then return "found at A[r]"
5.   m ← l + ⌈ \frac{A[r]-A[l]}{A[r]-A[l]} · (r-l) ⌉
6.   if (A[m] == k) then return "found at A[m]"
7.   else if (A[m] < k) then r ← m+1
8.   else r=m-1
   // we always return from somewhere within the
   // while loop
```

T^{avg} OF INTERPOLATION SEARCH IS $O(\log \log n)$

💡 We can show under some assumptions,

$$T^{\text{avg}}(n) \in O(\log \log n).$$

→ see next page for proof under a realization.

T^{worst} OF INTERPOLATION SEARCH IS $\Theta(n)$

💡 We can show that T^{worst} of interpolation-search is bad!

eg Consider

0 1 2 ... 8 9 11 "

and let's search for 10.

Then $A[r] - A[l]$ is huge! \Rightarrow so $m \approx l+0 = 0$.

Then our next lower bound is 1, and so on, and we look at every element!

💡 In particular,

$$T^{\text{worst}}(n) \in \Theta(n).$$

OPTIMIZED INTERPOLATION-SEARCH

Pseudocode:

```

interpolation-search-modified(A, n, k)
1. if (k <= A[0]) or (k > A[n-1]) return "not found"
2. if (k == A[n-1]) return "found at index n-1"
3. r <= 0, r = n-1 // A[r] <= k <= A[r]
4. while (N <= (r-l-1) > 1)
5.     m = l + ⌈  $\frac{k - A[r]}{A[r] - A[l]}$  ⌉ · (r-l-1)
6.     if (A[m] <= k)
7.         for h=1,2,...,
8.             l = m + (h-1)  $\sqrt{n}$ , r' = min{r, m+h  $\sqrt{n}$ }
9.             if (l=r' or A[r'] > k) then r=r' and break
10.    else ...
11.    if (k == A[l]) return "found at index l"
12.    else return "not found"

```

- ① we compare here...
(a 'probe')

② but we also probe here, if $A[m] < k$.

③ and we keep going, hopping by \sqrt{m} each time until
 $- A[m + \sqrt{m}] > k$; or
 \rightarrow bounds.

So: the idea is we use more probes to guarantee the subarray has size $O(n^{\alpha})$. - $m \times n^{\alpha}$ is

Q2 If $T(n) = \#$ of comparisons on n numbers, then

$$T(n) = T(n') + \# \text{ of probes}$$

where $n' \leq \sqrt{n}$.

$$P(\text{offset}(k) = i) = \binom{N}{i} p^i (1-p)^{N-i} \quad (\text{ie } \text{offset}(k) \in \text{Bin}(N, p))$$

Thu

- $E[\text{offset}(u)] = N_p$
 - $E[\text{id}_{X(u)}] = \ell + N_p$
 - $V(\text{id}_{X(u)}) = V(\text{offset}(u)) = N_p(1-p) \leq \frac{N}{4}$

$$\text{So } P(\#\text{probes} \geq 3) \leq \frac{VC(\text{id} \times \{k\})}{N} \leq \frac{1}{4},$$

and in general

$$\begin{aligned}
 \text{So } E(\# \text{ probes}) &= \sum_{n=0}^{\infty} n \cdot P(\# \text{ probes} = n) \\
 &\leq \sum_{n=1}^{\infty} P(\# \text{ probes} \geq n) \\
 &\leq \underbrace{1}_{n=1} + \underbrace{\frac{1}{n=2}}_{n=2} + \sum_{k=3}^{\infty} \frac{1}{4(k-2)^2} \\
 &= 2 + \frac{1}{4} \sum_{i=1}^{\infty} \frac{1}{i^2} = 2 + \frac{1}{4} \frac{\pi^2}{6} \leq 2.5. \quad \text{PQ}
 \end{aligned}$$

Therefore

$$T^{\text{avg}}(n) \in O(\log \log n)$$

Proof. Specifically, we prov

$$T^{\text{avg}}(n) \leq 2.5 \sqrt{\log \log n}$$

for $n \geq 4$.

Let's consider $L \in \mathbb{Z}^+$

$$2^L < n \leq 2^{L+1}$$

$$\Rightarrow \lceil \log \log(n) \rceil = L$$

$\Rightarrow |a - b| \leq f(x) \leq 1 - 1 = 0$ for all $x \in [0, 1]$.

$\Rightarrow \log \log(nn) \leq L-1 = \lceil \log \log(n) \rceil - 1$

$$T^{\text{avg}}(n) \leq T^{\text{avg}}(n') + 2.5,$$

$$\leq 2.5 \lceil \log \log(n') \rceil + 2.5$$

$$\leq 2.5 \lceil \log \log(n\ln n) \rceil + 2.5$$

as needed.

DICTIONARIES OF WORDS

WORDS

\exists_1 "Words" are strings; ie an array of chars in some alphabet Σ .
eg $\Sigma = \{0,1\}$, $\Sigma = \text{ASCII}$, etc

\exists_2 Note words have arbitrary length.

\exists_3 We sort words lexicographically:

cat
fish
color
coloring

\exists_4 Comparing strings (we denote this by `'strcmp'`) takes $T_{\text{worst}} = \Theta(\min\{|w_1|, |w_2|\})$ run-time.

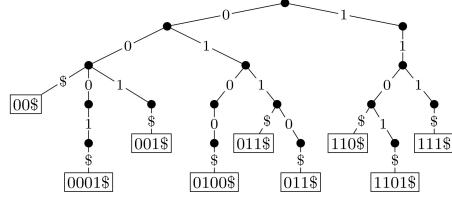
TRIES

\exists_1 A "trie" is a dictionary of words.

\exists_2 It is "prefix-free" — no string is a prefix of another.

\exists_3 This is satisfied if

- ① all strings have the same length; or
- ② all strings end with an "end-of-word" character '\$'.



*Note: items (keys) are only stored at the leaves.

OPERATIONS ON TRIE TAKE $O(|x|)$

TIME

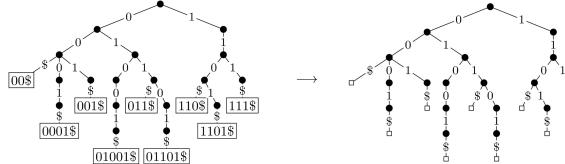
\exists Note that in a trie, search, insert & delete all take $O(|x|)$ time.

NO-LEAF-LABEL TRIES

\exists_1 Idea: we don't store the actual keys at the leaves.

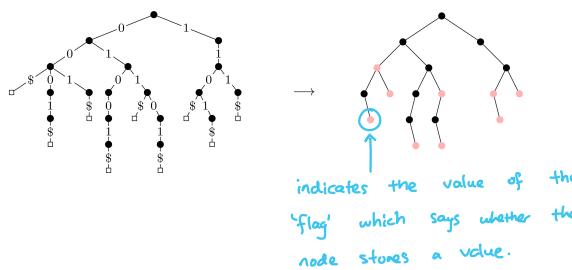
\exists_2 The key is stored implicitly through the characters along the path to the leaf.

\exists_3 Hence, this halves the amount of space needed.



ALLOW-PROPER-PREFIXES TRIES

\exists_1 Idea: we permit storing words at interior vertices, so we can accommodate words that are prefixes of others.

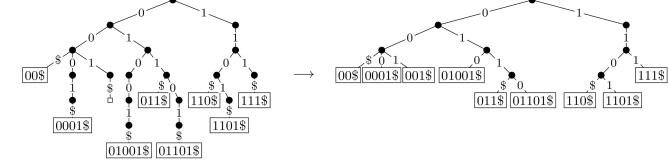


PRUNED TRIES

\exists_1 Idea: a pruned trie is such that we stop adding nodes as soon as the key is unique.

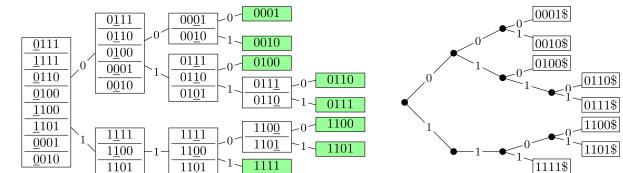
\exists_2 In particular,

- ① a node has a child only if it has two descendants;
- ② we save space if there are only a few long bitstrings; &
- ③ we can store infinite bitstrings (real numbers!)



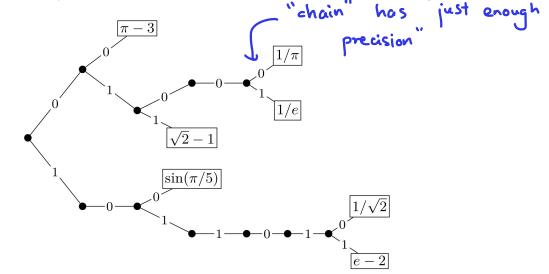
\exists_3 Note that we must store the full keys.

\exists_4 Also note that pruned tries are the recursion trees of MSI-radix-sort.



PRUNED TRIES FOR REAL NUMBERS

\exists we can use pruned tries to store infinite length numbers, eg real numbers:



"chain" has "just enough precision"

T IS A PRUNED TRIE THAT STORES n RANDOMLY CHOSEN INFINITE BITSTRINGS
 $\Rightarrow T^{\text{exp}}(\text{search}(B)) = O(\log n)$

As above. (B is an infinite bitstring).

Proof. Say we store B_1, \dots, B_n , and search for B.
 Then, each bit of each B_i was chosen in $\{0, 1\}$ uniformly.

Let the indicator variable

$$I_i = \begin{cases} 1, & \text{if we compared } B[i] \text{ to someone} \\ 0, & \text{otherwise.} \end{cases}$$

In particular,

$$\# \text{ comps for search}(B) = \sum_{i=0}^{\infty} I_i.$$

Then, let

$$I_{i,u} = \begin{cases} 1, & \text{we compared } B[i] \text{ w/ } B_u[i] \\ 0, & \text{otherwise.} \end{cases}$$

Then

$$\begin{aligned} P(I_{i,u}=1) &= P(B \text{ & } B_u \text{ agree on first } i \text{ bits}) \\ &= \left(\frac{1}{2}\right)^i \quad (\text{since bits of } B_i \text{ chosen randomly}). \end{aligned}$$

Hence

$$E[I_{i,u}] = P(I_{i,u}=1) = \left(\frac{1}{2}\right)^i,$$

and so

$$\begin{aligned} E[I_i] &\leq E\left[\sum_{u=1}^n I_{i,u}\right] \\ &\leq \sum_{u=1}^n \underbrace{E[I_{i,u}]}_{1/2^i} = \frac{n}{2^i}. \end{aligned}$$

But since $E[I_i] \in 1$, hence $E[I_i] \leq \min\{1, \frac{n}{2^i}\}$,

and so $E\left[\sum_{i=0}^{\infty} I_i\right] \leq O(1) + O(n)$.

The rest of the proof follows like the skip-list proof.

COMPRESSED TRIE

Q₁ Here, we compress paths of nodes with only one child.

Q₂ Each node stores an index corresponding to the depth in the uncompressed trie.

Q₃ This gives the next bit to be tested during a search.

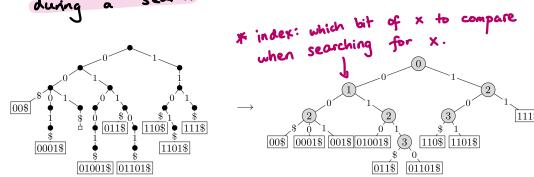


Figure 6.13: An example of a compressed trie.

Q₄ In particular, we know every node has ≥ 2 children.

Q₅ Thus, if the trie has n leaves, then it has $\leq n-1$ internal nodes, and so

$$\text{total \# of nodes} \leq 2n+1.$$

Q₆ Hence the trie takes $O(n)$ space, whereas other tries uses space that depends on the length of the words stored.

* if we only consider the space taken by the nodes.

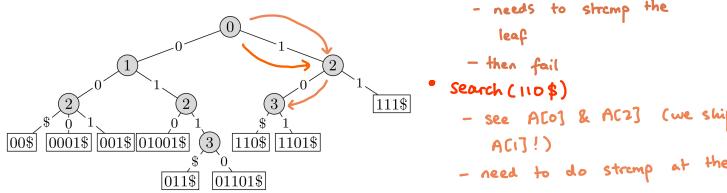
eg Consider



So space is much bigger in a pruned trie.

SEARCH IN COMPRESSED TRIES

Q₁ Consider:



- Search(001\$)
 - no such key
- search(0000\$)
 - needs to stomp the leaf
 - then fail
- search(110\$)
 - see A[0] & A[2] (we skip A[1]!)
 - need to do stomp at the leaf
- search(1\$)
 - needs to fail with "too short"
 - since A[2] DNE.

Q₂ Pseudocode:

```
Algorithm 6.8: compressedTrie::search(v ← root, x)
Input : Node v of the trie; word x to search for
1 if v is a leaf then
2   return strcmp(x, v.key)
3 else
4   d ← index stored at v
5   if x has at most d bits then
6     | return "not found"
7   else
8     let v' be the child of v for which the link from v is labeled with x[d]
9     if there is no such child then
10    | return "not found"
11   else
12    | compressedTrie::search(v', x)
```

Q₃ Run-time: $O(|x|)$, where $|x|$ is the length of the word to search.

INSERT IN COMPRESSED TRIES

Q₁ Idea:

- ① Find where it should be;
- ② Modify trie to put it there.

Q₂ Details are messy & omitted.

Q₃ Run-time: $O(|x|)$.

DELETE IN COMPRESSED TRIES

Q₁ Run-time: $O(|x|)$

PREFIX-SEARCH

Q₁ Given x & a compressed trie, is x a prefix of a stored word?

Q₂ Can also be done in $O(|x|)$ time.

MULTIWAY TRIES

Q₁ These are used to represent larger alphabets.

Q₂ Main question: how do we store the children?

- need to find (during search(x)) the child that stores $x[d]$.



Q₃ Options:

- ① Array
- ② List
- ③ Dictionary

Chapter 7:

Hashing

DIRECT ADDRESSING

E₁ Assume each key k is an integer with $0 \leq k < M$.

E₂ Then, we can store the k 's by using an array A of size M that stores (k, v) via $A[k] \leftarrow v$.



E₃ search, insert, delete have $\Theta(1)$ run-time.

E₄ But the total space is $\Theta(M)$.

HASHING

E₁ We can do direct addressing after modifying (ie "hashing") the key.

E₂ Assumptions:

① keys come from a universe U
- typically, U integers, $|U|$ finite

② Size of hash table, M , is pre-determined

③ We use a hash function $h: U \rightarrow \{0, \dots, M-1\}$ that maps indexes in U to an index in the array

0	1	2	3	4	5	6	7	8	9	10
	45	13		92	49		7			43

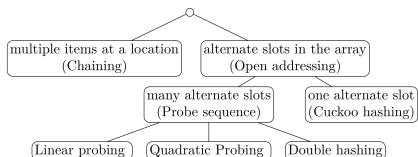
eg $U = \mathbb{N}$, $M = 11$, $h(k) = k \bmod 11$

slot $T[i]$

COLLISIONS

E₁ Collisions occur if we want to insert (k, v) into the table, but $T[h(k)]$ is already occupied.

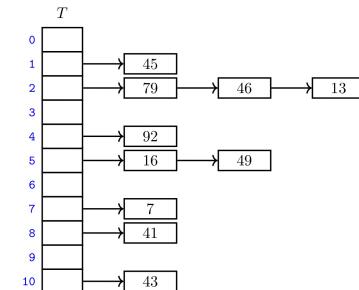
E₂ Solutions:



HASHING WITH CHAINING

E₁ Idea: use lists to resolve collisions.

E₂ Each slot stores a list.



* we use the MTF heuristic when inserting.

E₃ Run-time:

① Insert: $\Theta(1)$ + time to compute hash-value

- should choose $h(k)$ so time to compute $h(k)$ is $\Theta(1)$

② Search/delete: worst-case $\Theta(\text{length of } T[h(k)])$

↳ this is $\Theta(n)$ in worst-case

- we search in $T[h(k)]$, which is a list

RANDOMIZED VERSION OF HASHING

E₁ we randomly pick the hash function among all possible hash functions uniformly.

E₂ This is called the "uniform hashing assumption".

E₃ Under the uniform hashing assumption,

$$E[\text{length of bucket } h(k)] = \frac{n}{M} = \alpha.$$

list at $T[i]$

Here, " α " is called the "load factor".

Proof. Note $P(h(k)=i)$, for some key k & slot i , is $P(h(k)=i) = \frac{1}{M}$.

We have n keys $\rightarrow \frac{n}{M}$.

Thus, the exp. run-time for an unsuccessful search is $O(n)$.

However, if k is in the dictionary, then

$$E[\text{length of bucket } h(k)] = 1 + \frac{n-1}{M} \leq 1 + \alpha.$$

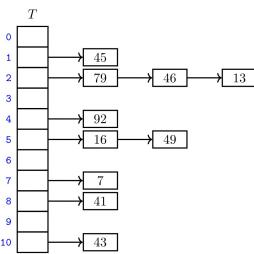
\uparrow
k is in
the bucket
 $\leq \alpha$

Hence, under uniform hashing,

search, delete take time $O(1+\alpha)$.

REHASHING

- 💡** We choose α by choosing M .
So, as n increases, we increase M so that α stays small.
💡 This is called "rehashing".



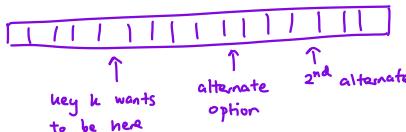
💡 If with rehashing we keep $O(1)$, then all operations take $O(1)$ time.

Moreover,

$$\text{space} \approx M + n = \frac{n}{\alpha} + n = O(n).$$

HASHING BY PROBING

- 💡** We want to avoid lists, since they have massive overhead.
💡 Idea: we allow keys to use multiple slots.



💡 For each key, we have a "probe sequence"

$$\langle h(k,0), h(k,1), \dots, h(k,M-1) \rangle$$

↑ index of option

LINEAR PROBING

💡 Here, we define

$$h(k,i) = (h(k) + i) \bmod M$$

0	1	2	3	4	5	6	7	8	9	10
45	13		92	49	37	7	41		43	

Consider $\text{insert}(37)$.

$$\Rightarrow h(k) = 4.$$

Probe seq is

$$\langle 4, 5, 6, 7, 8, 9, 10, 0, 1, 2, 3 \rangle.$$

💡 Linear probing builds big "clusters" of elements.

OPERATIONS OF PROBE HASHING

💡 In the array,

- ① To delete an key k , "mark" it as deleted.
- ② To search, we
 - follow the probe sequence;
 - ignore any "deleted" entries; &
 - continue until we find k or NIL.
- ③ To insert, we
 - follow the probe sequence; and
 - continue until we find a vacant spot (ie empty / deleted) or we reach the end of the probe sequence.

* this is called the "lazy deletion" technique.

💡 We also track how many elements are "deleted"; if this gets too large, we re-hash.

💡 Pseudocode:

Algorithm 7.1: $\text{probeSequenceHash::insert}(k, v)$

```

1 re-hash the table if the load factor is too big
2 for ( $j = 0; j < M; j++$ ) do
3   if  $T[h(k,j)]$  is NIL or 'deleted' then
4      $T[h(k,j)] \leftarrow (k, v)$ 
5     return "item inserted at index  $h(k,j)$ "
6   else
7     return "failure to insert" // need to re-hash

```

Algorithm 7.2: $\text{probeSequenceHash::search}(k)$

```

1 for ( $j = 0; j < M; j++$ ) do
2   if  $T[h(k,j)]$  is NIL then
3     return "item not found"
4   else if  $T[h(k,j)]$  has key  $k$  then
5     return "item found at index  $h(k,j)$ "
6   else // the current slot was 'deleted' or contains a different key
7     // ignore this and keep searching
7 return "item not found"

```

Algorithm 7.3: $\text{probeSequenceHash::delete}(k)$

```

1  $i \leftarrow$  index returned by  $\text{probeSequenceHash::search}(k)$ 
2  $T[i] \leftarrow$  'deleted'
3 re-hash the table if there are too many 'deleted' items

```

QUADRATIC PROBING

💡 Here, we define

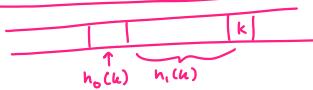
$$h(k,i) = (h(k) + c_1 i + c_2 i^2) \bmod M$$

for constants c_1, c_2 , which are picked so the probe sequence visits all slots in the hash table.

DOUBLE HASHING

Θ_1 "Double hashing" uses two hash functions h_0, h_1 , and we define

$$h(k, i) = (h_0(k) + ih_1(k)) \bmod M$$



Θ_2 Requirements:

- ① $h_1(k) \neq 0$
- ② $h_1(k)$ is relatively prime with M
- ③ h_1 & h_0 should be independent.

eg $h_0(k) = k \bmod 11$, $h_1(k) = \lfloor 10(\varphi k - \lfloor \varphi k \rfloor) \rfloor + 1$,
 $\varphi = \frac{\sqrt{5}-1}{2}$

CUCKOO HASHING

Θ_1 In "cuckoo hashing", we use two hash tables & two hash functions, and promise key k is always at $T_0[h_0(k)]$ or $T_1[h_1(k)]$.

eg $M=11$, $h_0(k) = k \bmod 11$, $h_1(k) = \lfloor 11(\varphi k - \lfloor \varphi k \rfloor) \rfloor$

Θ_2 In particular, search & delete have $O(1)$ worst-case run-time.

Pseudocode for insert:

Algorithm 7.4: `cuckooHash::insert(k, v)`

```

1  $i \leftarrow 0$ 
2 repeat
3   if  $T_i[h_i(k)]$  is NIL then
4      $T_i[h_i(k)] \leftarrow (k, v)$ 
5     return "success"
6   else
7     swap((k, v),  $T_i[h_i(k)]$ )
8      $i \leftarrow 1 - i$ 
9 until we have tried  $2n$  times
// If we reach this point then the hash table is probably too full.
10 return "unsuccessful, should re-hash"
```

COMPLEXITY OF OPEN ADDRESSING STRATEGIES

Θ_1 Note:

	Number of key-comparisons is at most:		
	insert	search	delete (successful search)
Chaining	1 (worst-luck)	$1 + \alpha$	$1 + \frac{1}{2}\alpha$ (avg-inst.)
Linear Probing	$\frac{1}{2}\alpha \frac{2-\alpha}{(1-\alpha)^2}$	$\frac{1}{2}\alpha \frac{2-\alpha}{(1-\alpha)^2}$	$\frac{1}{2}\alpha \frac{2-\alpha}{1-\alpha}$ (avg-inst.)
Double Hashing	$\frac{1}{1-\alpha} + o(1)$	$\frac{1}{1-\alpha} + o(1)$	$\frac{1}{1-\alpha} + o(1)$
Uniform Probing	$\frac{1}{1-\alpha}$	$\frac{1}{1-\alpha}$	$\frac{1}{\alpha} \log \left(\frac{1}{1-\alpha} \right)$ (avg-inst.)
Cuckoo Hashing	$\frac{\alpha}{(1-2\alpha)^2}$	1 (worst-luck)	1 (worst-luck)

Θ_2 Thus, all operations have $O(1)$ expected run-time if hash function is chosen randomly & α is kept sufficiently small.

Θ_3 But for a fixed hash function, the worst-case run-time is $\Theta(n)$.

CHOOSING HASH FUNCTIONS

MODULAR METHOD

Θ_1 Here, we let

$$h(k) = k \bmod M,$$

where we usually pick M to be prime.

*don't pick $M=2^m$ or $M=10^m$!

MULTIPLICATION METHOD

Θ_1 Here, we let

$$h(k) = \lfloor M(A \cdot k - \lfloor A \cdot k \rfloor) \rfloor$$

$\underbrace{A}_{\in [0, 1)}$

$\underbrace{M}_{\in [0, m]}$

where $A \cdot k$ for some $k \geq 0$; &

- ① $M = 2^k$ for some $k \geq 0$; &
- ② $A \in (0, 1)$ (preferably an irrational).

HASHING MULTI-DIMENSIONAL DATA

Θ_1 Suppose we wanted to hash keys that are words (ie in Σ^*).

Θ_2 We can "flatten" the string w into an integer:

eg $A \cdot P \cdot P \cdot L \cdot E \rightarrow (65, 80, 80, 76, 69)$ (ASCII)

$$\rightarrow 65R^4 + 80R^3 + 80R^2 + 76R + 69, R=255$$

Θ_3 We can then combine this with a modular hash function (ie $h(w) = f(w) \bmod M$).

Θ_4 To compute this in $O(|w|)$ time without overflow, we use "Horner's rule" and apply mod early:

$$\text{ie } ((((((65R+80) \bmod M)R+80) \bmod M)R+76) \bmod M)R+69 \bmod M$$

RANDOMLY-CHOSEN HASH FUNCTIONS

Θ_1 There are $|U|^M$ many possible hash functions, & ideally we would choose randomly amongst them.

Θ_2 But then we cannot compute the hash value quickly!

Θ_3 Idea: Fix a family \mathcal{H} of hash-functions that are easy to compute, & choose uniformly among them.

UNIVERSAL HASH-FUNCTIONS

For analysis, we want uniform hash-values;
ie

$$P(h(k) = i) = \frac{1}{m}.$$

But we also need small probability of collisions
(ie "universal hashing");
in other words,

$$P(h(k) = h(k')) \leq \frac{1}{m} \quad \forall k \neq k'$$

CARTER-WEGMAN HASH FUNCTION

The "Carter-Wegman hash-function" is defined to be

$$\begin{aligned} h_{a,b}(k) &= f_{a,b}(k) \bmod M \\ &= (a \cdot k + b \bmod p) \bmod M \end{aligned}$$

where $k \in \mathbb{Z}_p$, p is prime, & $a \neq 0$, b are chosen randomly.

We claim $f_{a,b}$ defines a permutation of \mathbb{Z}_p ; ie

$$f_{a,b}(k) \neq f_{a,b}(k') \text{ for } k \neq k'$$

Proof. Assume $f_{a,b}(k) = f_{a,b}(k')$.

$$\begin{aligned} &\Rightarrow f_{a,b}(ak) \equiv f_{a,b}(ak') \pmod{p} \quad (\because p \in \mathbb{Z}_p) \\ &\Rightarrow (ak+b) \equiv (ak'+b) \pmod{p} \quad (\equiv_p \Leftrightarrow \pmod{p}) \\ &\Rightarrow ak+b - ak' - b \equiv 0 \pmod{p} \\ &\Rightarrow a(k-k') \equiv 0 \pmod{p} \end{aligned}$$

Since $a \in \{1, \dots, p-1\}$ & $k-k' \in \{-p+1, \dots, p-1\}$, thus

$$k-k'=0,$$

ie $k=k'$, and we're done. \blacksquare

CARTER-WEGMAN FUNCTIONS ARE UNIVERSAL

See that

$$P(h_{a,b}(k) = h_{a,b}(k')) \leq \frac{1}{m}.$$

Proof. Assume $h_{a,b}(k) = h_{a,b}(k')$ for $k \neq k' \in \mathbb{Z}_p$.

We know that

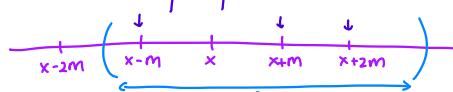
$$\underbrace{f_{a,b}(k)}_x \neq \underbrace{f_{a,b}(k')}_{x'},$$

but $x \bmod M = x' \bmod M$ (by defn of $h_{a,b}$).

Thus

$$x - x' \equiv_M 0.$$

How many such "bad" pairs (x, x') could there be in $\mathbb{Z}_p \times \mathbb{Z}_p$?



x' is one of these, but $x' \neq x$.

Hence, x' is among $\lceil \frac{p}{m} \rceil - 1$ numbers, and so

$$\# \text{ choices for } x' \leq \frac{p-1}{m}.$$

Fixing x , it follows that

$$\# \text{ bad pairs} \leq \underbrace{p}_{\text{choices for } x} \cdot \underbrace{\frac{p-1}{m}}_{\text{choices for } x'}.$$

Therefore

$$\begin{aligned} P(h_{a,b}(k) = h_{a,b}(k')) &= P((x, x') \text{ formed a "bad" pair}) \\ &= \sum_{\text{bad pairs } x, x'} P(f_{a,b}(k)=x, f_{a,b}(k')=x') \\ &= \sum_{\text{bad pairs } x, x'} P(a = (k-k')^{-1}(x-x') \bmod p, b = (x-ak) \bmod p) \\ &= \sum_{\text{bad pairs } x, x'} \frac{1}{p-1} \cdot \frac{1}{p} \\ &= \frac{1}{p(p-1)} \cdot \# \text{ bad pairs} \leq \frac{1}{m}. \quad \blacksquare \end{aligned}$$

Chapter 8:

Range Search

💡 Idea:

- ① We are given two keys $x_1 \leq x_2$; &
- ② We want to get all items between x_1 & x_2 .

💡 Sorted array:

5 10 12 17 23 32 45 62 71 ...

> rangeSearch(13, 36)

- ① Search for x_1
→ will tell us where x_1 would be "between"
- ② Search for x_2
→ likewise.
- ③ Repeat everyone in between.

* this approach works for tries, SIs, etc.

💡 Run-time: $O(\log n + s)$,
where "s" is the "output size".

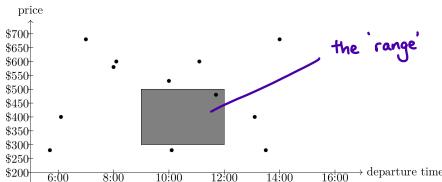
MULTI-DIMENSIONAL DATA

💡 Idea: "points" instead of single values.

ie (x_0, x_1, \dots, x_n)

* we will assume all data are integers.

💡 Range search then looks like:



💡 This is called an "orthogonal range-search":

- ① We are given a query rectangle

$$A = [x, x'] \times [y, y'] ; \quad \&$$

- ② We want the points in A.

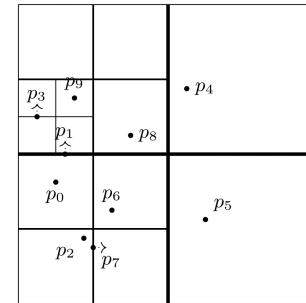
💡 Alternative depiction:

QUAD-TREES

💡 For "quad-trees":

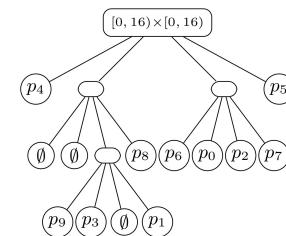
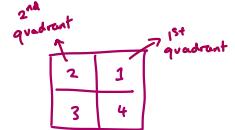
- ① stores points in 2D;
- ② we assume the points are in a bounding box $[0, 2^k] \times [0, 2^k]$
eg $[0, 16] \times [0, 16]$

💡 $\bullet p_4$



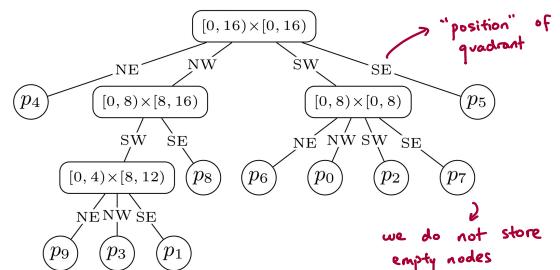
💡 How to construct?

- ① Split the bounding box into quadrants;
- ② Define the corresponding children;
- ③ Repeat at children until ≤ 1 point left in all regions.



💡 Note: points on the split line go above/right.

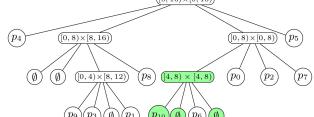
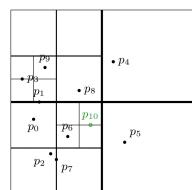
💡 Alternative depiction:



- ① We do not need to store the "sub-areas"
(we always cut in half)

OPERATIONS OF QUAD-TREES

Q1 Inset:



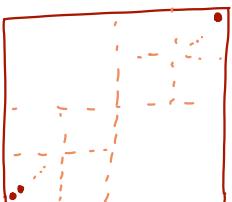
Q2 Run-time:

- ① No good bound depending on n ;
- ② We can only say it is $O(\text{tree height})$

HEIGHT OF A QUAD-TREE

Q1 But we have no bound on the height either!

eg

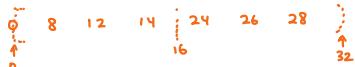


→ height arbitrary.

Q2 However, we can bound the height if the coordinates are integers; specifically in the range $\{0, \dots, 2^l - 1\}$.

Then the height of the quad-tree is $\leq l$.

Proof. Detour: let's look at a quad-tree in 1d.

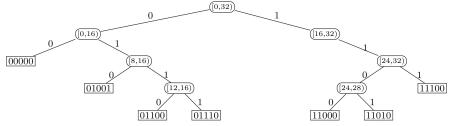


\downarrow



But we can also represent the points in base-2:

"Points:" 0 9 12 14 24 26 28
(in base-2) 00000 01001 01100 01110 11000 11010 11100



This is a pruned tree!

The height of a 1d quad-tree of integers in $\{0, \dots, 2^l - 1\}$ \leq the length of longest bitstring = l bitstrings of length l

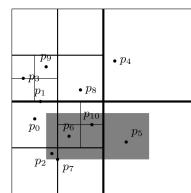
This argument generalizes to higher dimensions (e.g. 2D, i.e. quad trees)

Q3 In particular, a quad-tree is a pruned trie where we split by two keys in parallel.

Q4 The expected height of a quad-tree of randomly chosen points is $O(\log n)$.

RANGE-SEARCH IN QUAD-TREES

Q1 Idea:



"outside" node:
region disjoint from A
→ stop search

"boundary" node:
region overlaps with A
→ continue search

Q2 Pseudocode:

Algorithm 8.1: `quadTree::rangeSearch(r ← root, A)`

```

Input : Node r of a quad-tree. Query-rectangle A
1  $R \leftarrow$  square associated with node r
2 if  $R \subseteq A$  then
3   | report all points in subtree at r and return
4 else if  $R \cap A$  is empty then
5   | return
6 else
7   if r is a leaf then
8     | if r stores a point p and p ∈ A then
9       |   | return p
10    | else
11      |   | return
12  else
13    foreach child v of r do
14      |   | quadTree::rangeSearch(v, A)

```

① Checking " $R \cap A = \emptyset$ " is $O(1)$.

- since we use rectangles.

Q3 Run-time:

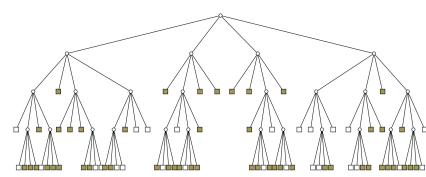
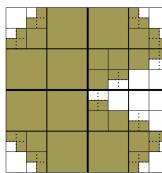
① We have no bounds in terms of n or s .

② Only thing we can say:

run-time $\in O(\text{size of quad tree})$.

OTHER USES OF QUAD-TREES

Q1 We can use quad-trees to store pixelated images:



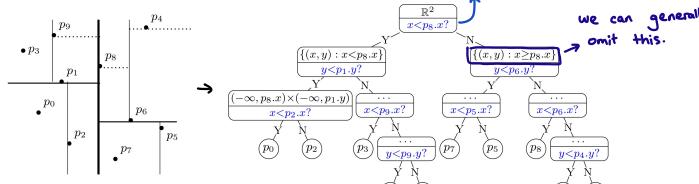
Q1 We can use quad-trees to store pixelated images:

kd-TREE

* d = "dimensional"

- 💡 Idea:
 - ① Similar to quad-trees; but
 - ② We split points in half (rather than the region);

comparison to split the points.



💡 ③ Repeat ② until 1 point left; &

💡 ④ Alternate split by x & y.

💡 Note: we always split at $\text{QuickSelect}(\lfloor \frac{n}{2} \rfloor)$.
ie the median if n is odd, & the upper median
if n is even.

💡 ⑤ we assume the points are in "general position":
no two x coordinates are the same, nor no two y coordinates are the same.

💡 ⑥ Under this, each split leads to

① $\lfloor \frac{n}{2} \rfloor$ points in the left child; &

② $\lceil \frac{n}{2} \rceil$ points in the right child.

💡 Therefore, the height is $O(\log n)$.

OPERATIONS OF kd-TREE

💡 Search: $O(\log n)$

💡 Insert/delete: basically impossible!

💡 However, we can build the whole thing, given all n points, in $O(n \log n)$ expected time:

① Find the median; &

→ $O(n)$ expected

② Recurse in children

→ size $\approx \frac{n}{2}$ each.

RANGE-SEARCH FOR kd-TREES

💡 Pseudocode:

Algorithm 8.2: kdTree::rangeSearch(r ← root, A)

Input : Node r of a kd-tree. Query-rectangle A

```

1 R ← region associated with node r
2 if  $R \subseteq A$  then                                // inside node
3   report all points in subtree at r and return
4 else if  $R \cap A$  is empty then                      // outside node
5   return
6 else                                                 // boundary node
7   if r is a leaf then
8     if the point p stored at r satisfies  $p \in A$  then
9       return p
10    else
11      return
12  else
13    foreach child v of r do
14      kdTree::rangeSearch(v, A)

```

* nearly identical to that in quad-trees.

💡 Run-time: $O(\# \text{ of boundary nodes} + s)$.

↑
output size

💡 Then,

of boun