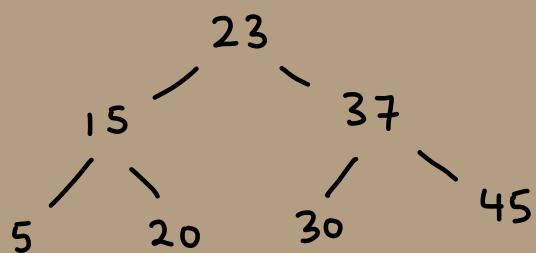


# CS 240E

## 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: $O(f(n))$

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

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

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

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  $O \leq 2n^2 + 3n^2 + 11 \leq cn^2 \quad \forall n \geq n_0$ .

Sol<sup>2</sup>. Consider  $n_0=1$ . Then

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

#### $\Omega$ -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

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

#### $\Theta$ -NOTATION (BIG THETA): $f(n) \in \Theta(g(n))$

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

$$f(n) \in \Theta(g(n)) \Leftrightarrow f(n) \in O(g(n)) \quad \& \quad f(n) \in \Omega(g(n)).$$

#### $\mathcal{O}$ -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  $g(n)$ .

#### $\omega$ -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  $A[i]$ , we need  $i-1$  swaps.

So

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

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))

$\Leftrightarrow$  (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  $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = L$ .  
 $\Rightarrow \forall \epsilon > 0: \exists n_0$  s.t.  $|\frac{f(n)}{g(n)} - L| < \epsilon \quad \forall n \geq n_0$ .

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.

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

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

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

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

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))

$\Leftrightarrow$  (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))

$\Leftrightarrow$  (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))

$\Leftrightarrow$  (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  $f(n) \in O(g(n)), g(n) \in \Theta(h(n)) \forall n \geq N \Rightarrow f(n) \in \Theta(h(n))$
- ⑥  $f(n) \in \Omega(g(n)), g(n) \in h(n) \forall n \geq N \Rightarrow f(n) \in \Omega(h(n))$
- 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  $h(n) \in O(f(n) + g(n)) \Rightarrow h(n) \in O(\max(f(n), g(n)))$
- ⑩  $h(n) \in \Omega(f(n) + g(n)) \Rightarrow h(n) \in \Omega(\max(f(n), g(n)))$

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

$\Leftrightarrow$  (let  $f(n) \in P_d(\mathbb{R})$ , 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  $\gg$

$\Leftrightarrow$  (let  $b > 1$ . 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" = "log<sub>2</sub>".

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

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

Proof follows from the limit rule.  $\square$

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

$\Leftrightarrow$  Suppose  $f(n) \in o(g(n))$ . Then  $f(n) \in O(g(n))$ .

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

$\Leftrightarrow$  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

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

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

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

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

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

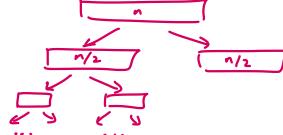
Suppose  $f(n) \in \omega(g(n))$ . Then  $f(n) \in \Omega(g(n))$ .

**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  $\pi$ , we define its inverse  $\pi^{-1}$  to be the array which entries have exactly the same "relative order" as in  $A$ .

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

We denote " $\Pi_n$ " 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  $\Pi_n$  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  $\pi^{-1}$ .

# 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... $\lfloor \frac{n}{2} \rfloor - 1$ ,  $\lfloor \frac{n}{2} \rfloor$ ) // 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

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

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

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

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

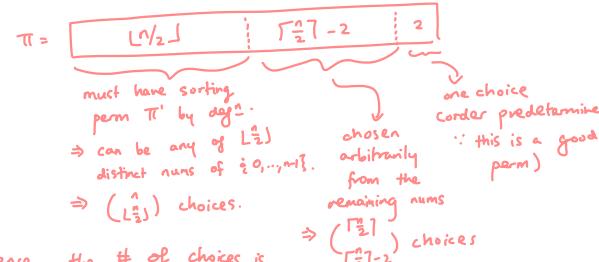
Note that  $\pi_{\text{half}} \in \Pi_{\lfloor \frac{n}{2} \rfloor}$ . Then see that

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

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  $\pi \in \Pi_n^{\text{good}}(\pi')$ . See that



Hence, the # of choices is  $\binom{C(n, 2)}{C(n, 2) - 2}$

$$\left( \binom{C(n, 2)}{C(n, 2) - 2} \right) \cdot \left( \binom{C(n, 2) - 2}{C(n, 2) - 3} \right) \cdots = \frac{n!}{C(n, 2)! \cdot 2},$$

as needed.  $\blacksquare$

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

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

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

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

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

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: $T_0^{\text{amort}}$

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]: $\Phi$

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  $\Phi_i$  is the value of  $\Phi$  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  $\Phi$ , 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  $\Phi_{\text{before}}$  &  $\Phi_{\text{after}}$  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

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

## 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  $\Phi$  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  $\Phi$  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  $\leq 1$  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

Q<sub>1</sub>: A "priority queue" stores items that have a priority, or key.

Q<sub>2</sub>: Operations:

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

Q<sub>3</sub>: We can use a PQ to sort:

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

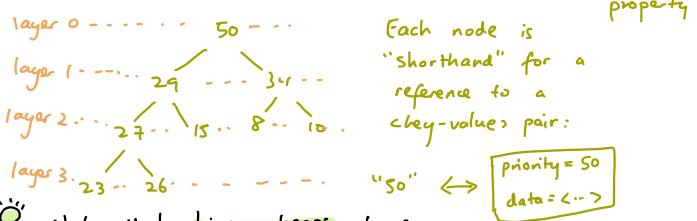
```
init PQ to an empty PQ  
for i=0 to n-1 do  
    PQ.insert(A[i])  
for i=n-1 down to 0 do  
    A[i] ← PQ.deleteMax()
```

Q<sub>4</sub>: The run-time of the above algorithm is O(initialization + n·insert + n·deleteMax).

## BINARY HEAPS

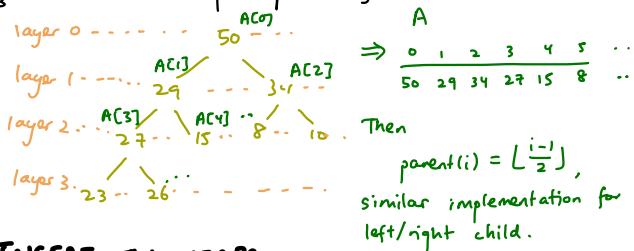
Q<sub>1</sub>: "Binary heaps" are binary trees with the following properties:

- 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"



Q<sub>2</sub>: Note that binary heaps have height  $O(\log n)$ .

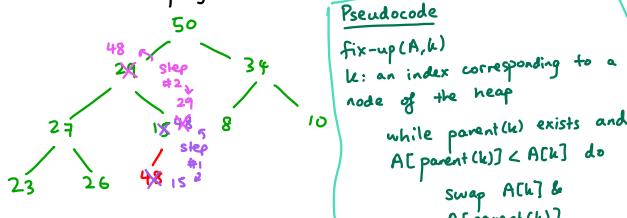
Q<sub>3</sub>: We store binary heaps using an array:



## INSERT IN HEAPS

Q<sub>1</sub>: To insert into a heap, we just place the new key at the first free leaf.

Q<sub>2</sub>: We also employ "fix-up":



### Pseudocode

fix-up(A, k)

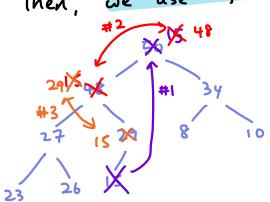
k: an index corresponding to a node of the heap  
while parent(k) exists and  $A[\text{parent}(k)] < A[k]$  do  
 swap  $A[k]$  &  $A[\text{parent}(k)]$   
 $k \leftarrow \text{parent}(k)$

## 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  $k_{\text{new}}$ , and decreases the key of  $i$  to  $k_{\text{new}}$  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).

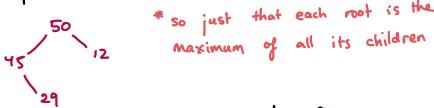
- for bin heap: increase value of key at  $i$  to  $\infty$  (or  $\text{findMax}().\text{key}() + 1$ );

- 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.



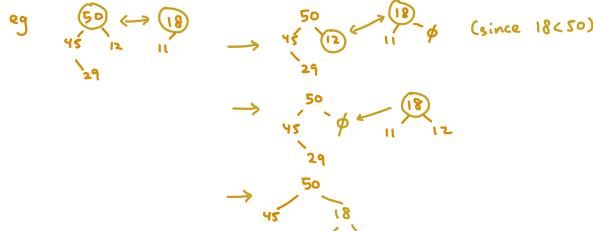
Operations for meldable heaps:

- 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

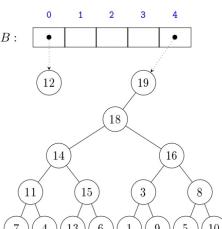
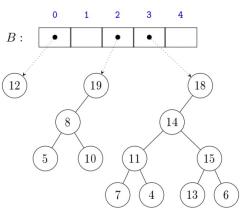
To make a BH proper, we do the following:  
if the BH has two flagged trees  $T, T'$  of the same height  $h$ , then they both have  $2^h$  nodes.

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 B[h]
        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( $k, v$ )

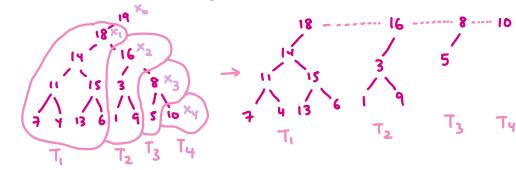
- like w/ meldable heaps; create a single-node binomial heap, then call merge
- takes time  $O(\log n)$

### 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

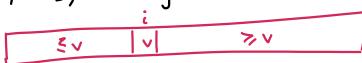
$[0, 10, 10, 30, 40, 50, 60, 70, 80, 90]$

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

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

\* 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:

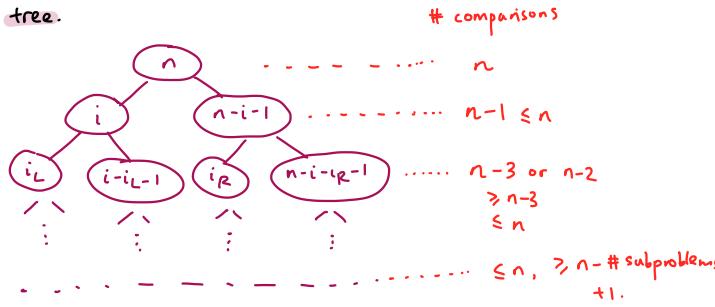
```
quick-sort(A)
  // A: array of size n
  1. if n ≤ 1 then return
  2. p ← choose-pivot(A)
  3. i ← partition(A, p)
  4. quick-sort(A[0, ..., i-1])
  5. quick-sort(A[i+1, ..., n-1])
```

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

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

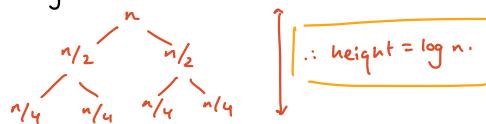
So what can we say about the recursion tree's height?

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  $\leq 10$
- 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,

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

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.

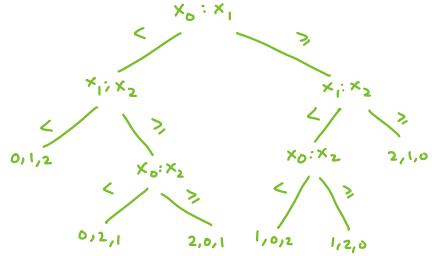
$$\begin{aligned} \text{Step: } T(n) &\leq n + \frac{2}{n} \sum_{i=2}^{n-1} 2i \ln(i) \\ &= n + \frac{4}{n} \sum_{i=2}^{n-1} i \ln(i) \\ &\leq n + \frac{4}{n} \int_2^n x \ln x \, dx \quad (\text{since } x \ln x \text{ is increasing}) \\ &\leq n + \frac{4}{n} \left( \frac{1}{2} n^2 \ln(n) - \frac{1}{4}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

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  $\leq 3$  key-comparisons.

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

Executing  $A$  on  $I_{\pi}$  leads to a leaf that stores  $\pi$ .

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  $\pi$ .

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

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

Therefore  $2^n > n!$ , or

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

Finally, consider the sorting perm  $\pi$  that has its leaf on layer  $h$ .

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

# SORTING INTEGERS

## BUCKET-SORT

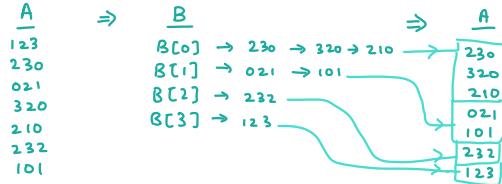
Bucket-sort can be used to sort a collection of integers by a specific "position" of digit.

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**

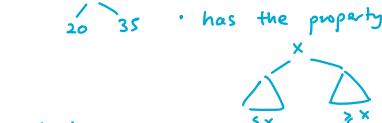
- fast insert, slow search, slow delete

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.

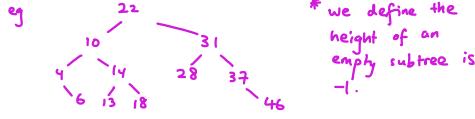
# AVL-TREES

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.      $z \leftarrow z.\text{parent}$

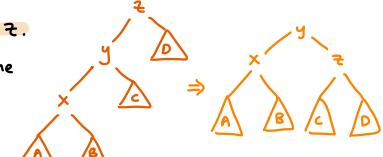
## ROTATIONS [OF BSTS]

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;  $x < y < z$ .

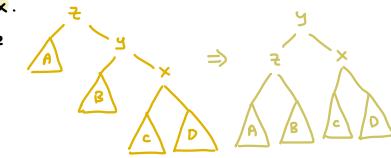
We want  $y$  to become the root.  
pseudocode similar to left rotation; see below



② Left rotation;  $z < y < x$ .

We want  $y$  to be the root.

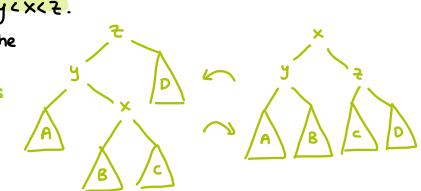
- rotate-left( $z$ )
1.  $y \leftarrow z.\text{right}$
  2.  $z.\text{right} \leftarrow y.\text{left}$
  3.  $y.\text{left} \leftarrow z$
  4. setHeightFromSubtrees( $z$ )
  5. setHeightFromSubtrees( $y$ )
  6. return  $y$



③ Double-right rotation;  $y < x < z$ .

We want  $x$  to be the root.

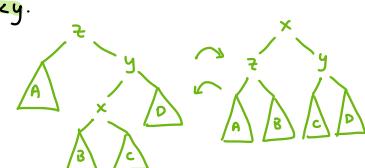
# can also be written as single-left rotation at  $y$  and then a single-right rotation at  $z$  (double-rotation)



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  $y = z.\text{right}$  &  $x = y.\text{left}$  //  $x-y-z$
  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?  

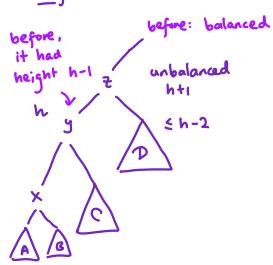
- height of tree =  $O(\log n)$
- rotations take  $O(1)$  time
- so delete takes  $O(\log n)$ .
- tight if we insert/delete at lowest level & never rotate.

## 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.

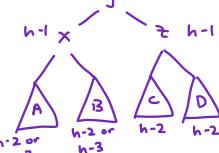
Proof.



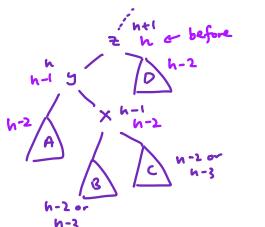
Case 1: "x" is the left node  

- x has height  $h-2$  before
- C has height  $h-2$

After rotation:



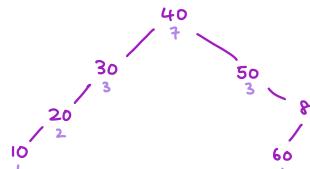
Case 2: x is the right child.



## 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

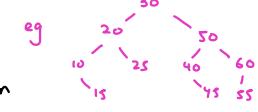
$$d \leq \log_{(\frac{1}{\alpha})}(n) \in O(\log n).$$

## 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

```

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

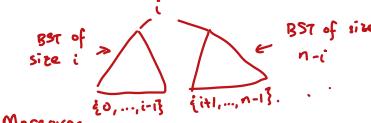
# RANDOMIZED BUILT BST HAS EXPECTED HEIGHT $\Theta(\log n)$

$\exists_1$  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  $\pi$  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  $\pi$  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:  $n=1 \Rightarrow \text{height}=0 \Rightarrow Y(1)=2^0=1 \leq 8$ .

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

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

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

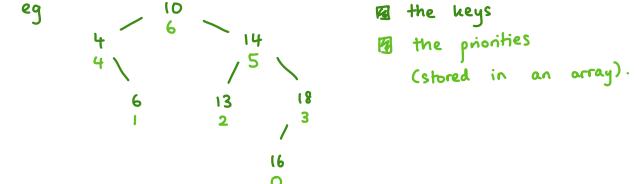
$\exists_1$  A "treap" is a randomized version of a BST, also called a "priority search tree".

$\exists_2$  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; &
- 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

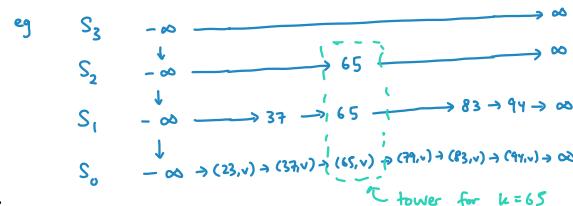
$\exists_1$  Note that treaps are

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  $-\infty$  &  $\infty$ , 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.



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)

💡 'getPredecessors(k)' is a helper routine for insert & delete.

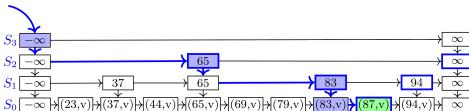
Algorithm 5.3: skipList::getPredecessors(k)

```

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)
2 p0 ← P.top()                                     // predecessor of k in S0
3 if p0.after.key = k then
4   | return p0.after
5 else
6   | return "not found, but would be after p0"

```

## INSERT

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

💡 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)

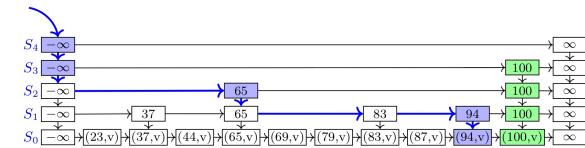
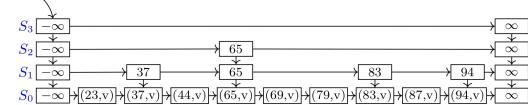
```

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 S0
10  while i > 0 do
11    p ← P.pop()                                              // insert k in S1, ..., 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

💡 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.

💡 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.

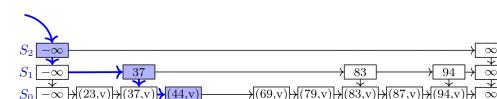
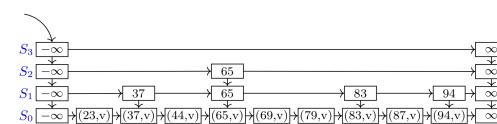
Algorithm 5.6: skipList::delete(k)

```

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.  $\square$

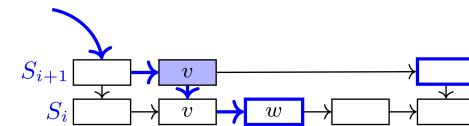
## getPredecessors: $E(\text{forward steps of } S_i) \leq 1$

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

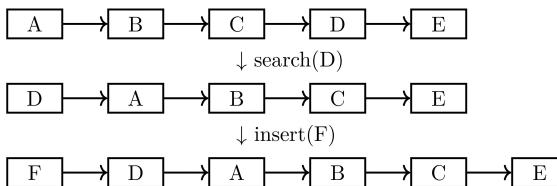
**B1** Here, we do not know how frequently a key is going to be accessed, so we cannot hope to build the best-possible data structure.

**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":

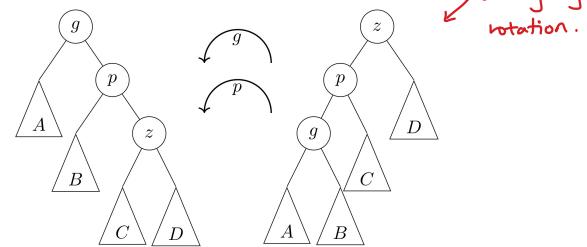
it takes at most twice as many comparisons that would have been taken using the optimal static ordering.

## 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.



a zig-zig rotation.

## INSERT

**B1** Pseudocode:

---

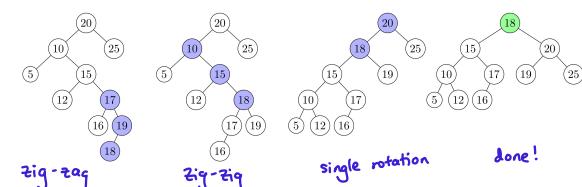
```

Algorithm 5.8: splayTree::insert(k, v)
1 z ← BST::insert(k, v)
2 while (z has a parent p and a grand-parent g) do
3     if (g is a left child of p) then                                // Zig-zig rotation rightward
4         | p ← rotate-right(g), z ← rotate-right(p)
5     else if (g is a right child of p) then                            // Zig-zig rotation rightward
6         | g.left ← rotate-left(p), z ← rotate-right(g)
7     else if (g is a left child of p) then                            // Zig-zig rotation leftward
8         | g.right ← rotate-right(p), z ← rotate-left(g)
9     else // (g is a right child of p)                                 // Zig-zig-rotation leftward
10    | p ← rotate-left(g), z ← rotate-left(p)
11 if (z has a parent p) then                                         // single rotation
12    if (z = p.left) then z ← rotate-right(p)
13    else z ← rotate-left(p)

```

---

**B2** Example:



done!

## RUN-TIME ANALYSIS

The amortized run-time of  $\text{insert}$  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  $\phi$  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  $\phi$  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  $\text{insert}$  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

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

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

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

$\Theta_1$ : we cannot do better than binary search.

- asymptotically, comparison-based

$\Theta_2$ : But, we can do a bit better;

- shave constant

$\Theta_3$ : But, we can do a lot better! & drop comparison based

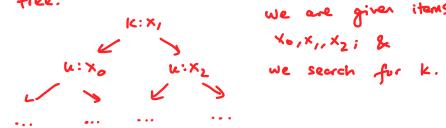
$\Theta_4$ : we can do even better!

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

$\Theta_1$ : we want to prove a lower bound for searching.

$\Theta_2$ : 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  $k \in (-\infty, x_0), k \in (x_0, x_1), \dots, k \in (x_n, +\infty)$ .

$\Theta_3$ : 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}⌋
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.

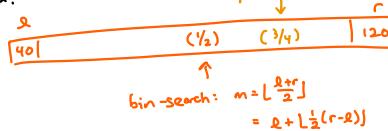
### INTERPOLATION SEARCH

### INTERPOLATION SEARCH

💡 "Interpolation search" assumes that

- ① we are given an array  $A[0 \dots n-1]$  of numbers in sorted order;
- ② we want to search for the number  $k$ .

💡 Idea: interpolation-search



Consider search(100).

⇒ bin-search searches naively in the middle.

⇒ interpolation search searches based on the "endpoints" of the keys, and interpolates where the key would be.

In particular, we notice that

$A[r] - A[l] = 80$  (the "distance" between nums).

Then  $k - A[r] = 60$ .

So, we should search for 100 roughly  $\frac{60}{80} = \frac{3}{4}$  in the range.

More generally, we search at the index,

$$l + \frac{A[r] - k}{A[r] - A[l]} (r - l)$$

start index      ratio      indices in range

and otherwise, the search works like bin-search.

### 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{k - 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 l ← m+1
8.   else r=m-1
// we always return from somewhere within the while loop
```

$T^{\text{avg}}$  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.

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

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

e.g. Consider

0 1 2 ... 8 9 11 "

and let's search for 10.

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

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

💡 In particular,

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

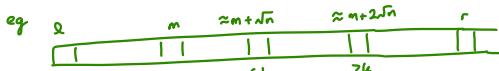
# OPTIMIZED INTERPOLATION-SEARCH

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.     l ← l + ⌈ k - A[l] ⌈ A[r] - A[l] ⌉ · (r-l-1) ⌉
6.     if (A[l] ≤ k)
7.         for h=1, 2, ...
8.             l ← m + (h-1) ⌈ n ⌉, r' ← min{r, m+h ⌈ 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  $A[m] \leq k$ .
- ③ and we keep going, hopping by  $\sqrt{n}$  each time until
  - $A[m+\sqrt{n}] > k$ ; or
  - $m+\sqrt{n}$  is out of bounds.

So: the idea is we use more probes to guarantee the subarray has size  $O(\sqrt{n})$ .

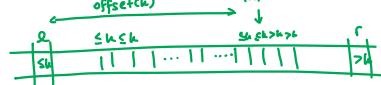
If  $T(n)$  = # of comparisons on  $n$  numbers, then

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

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

Moreover, avg # of probes  $\leq 2.5$ .

Proof.



Assume nums in  $A[r-l-1]$  have been chosen uniformly at random bw  $A[l]$  &  $A[r]$ . We want  $E[\# \text{probes}] \leq 2.5$ .

Then see that  $P(\# \text{probes} \geq 1) = 1$ ,  $P(\# \text{probes} \geq 2) \leq 1$ , &

$$\begin{aligned} P(\# \text{probes} \geq 3) &= P(A[\text{id of 2nd probe}] \leq k) \\ &= P(idx(k) \geq m + \sqrt{n}) \end{aligned}$$

(where  $idx(k)$  is a rv that is the largest  $i$  w/  $A[i] \leq k$ )

$$\begin{aligned} &\leq P(idx(k) - m \geq \sqrt{n}) \\ &= P(|idx(k) - E(idx(k))| \geq t) \\ &\leq \frac{V[idx(k)]}{t^2} \quad (V(X) = \text{Var}(X)). \\ &\text{Chebyshev} \end{aligned}$$

Then, see that  $offset(k) = idx(k) - l$ , &

$$P(offset(k) = i) = P(\text{exactly } i \text{ of the } r-l-1 \text{ randomly chosen numbers are } \leq k)$$

Then

$$\begin{aligned} P(\text{one number } x \text{ is } \leq k) &= P(\text{--- } k \text{ ---}) \\ &= \frac{k - A[l]}{A[r] - A[l]} \end{aligned}$$

and so

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

Thus

- $E[offset(k)] = Np$
- $E[idx(k)] = l + Np$
- $V(idx(k)) = V(offset(k)) = Np(1-p) \leq \frac{N}{4}$ .

$$So \quad P(\# \text{probes} \geq 3) \leq \frac{V[idx(k)]}{N} \leq \frac{1}{4},$$

and in general

$$P(\# \text{probes} \geq h) \leq \frac{1}{4^{(h-2)^2}}.$$

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

Therefore,

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

Proof. Specifically, we prove

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

for  $n \geq 4$ .

Let's consider  $L \in \mathbb{Z}^+$  s.t.

$$2^{L-1} < n \leq 2^L.$$

$$\Rightarrow 2^{L-1} < \log(n) \leq 2^L$$

$$\Rightarrow L-1 < \log \log(n) \leq L$$

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

$$\text{Observe that } \sqrt{n} \leq \sqrt{2^L} = \frac{1}{2} 2^L = 2^{L-1}$$

$$\Rightarrow \log(\sqrt{n}) \leq L-1 = \lceil \log \log(n) \rceil - 1.$$

Proof by induction, and we only consider the step:

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

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

$$\leq 2.5 \lceil \log \log(\sqrt{n}) \rceil + 2.5$$

$$\leq 2.5(\lceil \log \log n \rceil - 1) + 2.5 = 2.5 \lceil \log \log n \rceil,$$

as needed.  $\square$

# DICTIONARIES OF WORDS

## WORDS

$\exists_1$  "Words" are strings; ie an array of chars in some alphabet  $\Sigma$ .  
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.

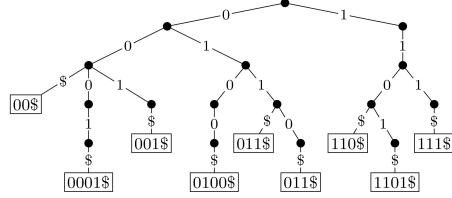
## TRIES

$\exists_1$  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

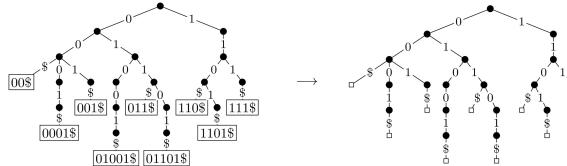
$\exists$  Note that in a trie, search, insert & delete all take  $O(|x|)$  time.

## NO-LEAF-LABEL TRIES

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

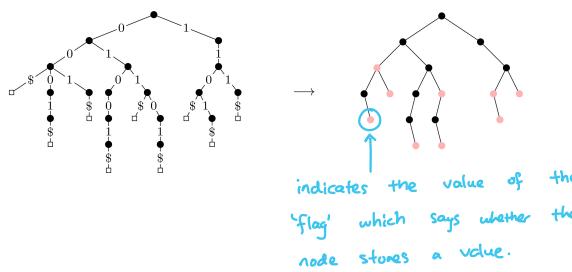
$\exists_2$  The key is stored implicitly through the characters along the path to the leaf.

$\exists_3$  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.

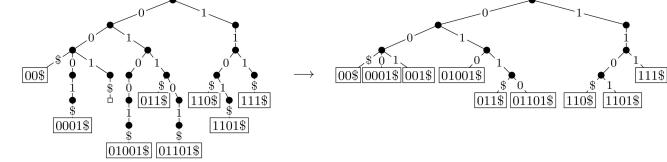


## PRUNED TRIES

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

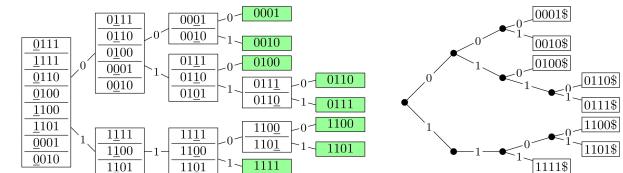
$\exists_2$  In particular,

- 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!)



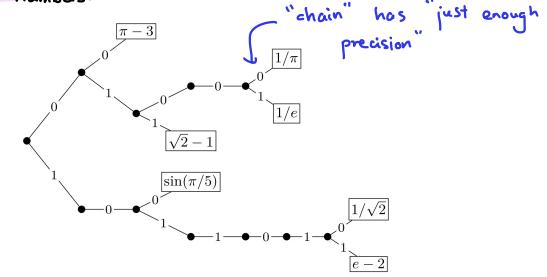
$\exists_3$  Note that we must store the full keys.

$\exists_4$  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:



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  $E[I_i] \leq 1$ , 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

Q<sub>1</sub> Here, we compress paths of nodes with only one child.

Q<sub>2</sub> Each node stores an index corresponding to the depth in the uncompressed trie.

Q<sub>3</sub> This gives the next bit to be tested during a search.

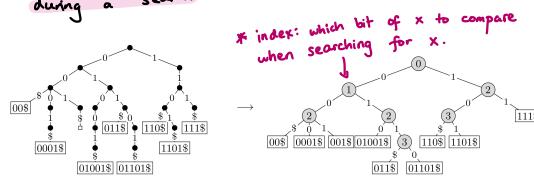


Figure 6.13: An example of a compressed trie.

Q<sub>4</sub> In particular, we know every node has  $\geq 2$  children.

Q<sub>5</sub> Thus, if the trie has  $n$  leaves, then it has  $\leq n-1$  internal nodes, and so

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

Q<sub>6</sub> 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.

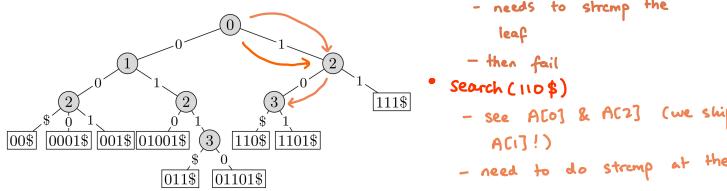
eg Consider



So space is much bigger in a pruned trie.

## SEARCH IN COMPRESSED TRIES

Q<sub>1</sub> 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.

Q<sub>2</sub> 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)
  
```

Q<sub>3</sub> Run-time:  $O(|x|)$ , where  $|x|$  is the length of the word to search.

## INSERT IN COMPRESSED TRIES

Q<sub>1</sub> Idea:

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

Q<sub>2</sub> Details are messy & omitted.

Q<sub>3</sub> Run-time:  $O(|x|)$ .

## DELETE IN COMPRESSED TRIES

Q<sub>1</sub> Run-time:  $O(|x|)$

## PREFIX-SEARCH

Q<sub>1</sub> Given  $x$  & a compressed trie, is  $x$  a prefix of a stored word?

Q<sub>2</sub> Can also be done in  $O(|x|)$  time.

## MULTIWAY TRIES

Q<sub>1</sub> These are used to represent larger alphabets.

Q<sub>2</sub> Main question: how do we store the children?
 

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



Q<sub>3</sub> Options:

- ① Array
- ② List
- ③ Dictionary

# Chapter 7:

## Hashing

### DIRECT ADDRESSING

**E<sub>1</sub>** Assume each key  $k$  is an integer with  $0 \leq k < M$ .

**E<sub>2</sub>** Then, we can store the  $k$ 's by using an array  $A$  of size  $M$  that stores  $(k, v)$  via  $A[k] \leftarrow v$ .



**E<sub>3</sub>** search, insert, delete have  $\Theta(1)$  run-time.

**E<sub>4</sub>** But the total space is  $\Theta(M)$ .

### HASHING

**E<sub>1</sub>** We can do direct addressing after modifying (ie "hashing") the key.

**E<sub>2</sub>** 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

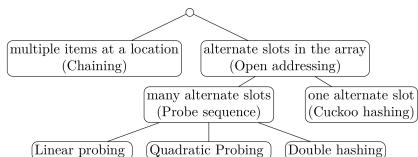
eg  $U = \mathbb{N}$ ,  $M = 11$ ,  $h(k) = k \bmod 11$

slot  $T[i]$

### COLLISIONS

**E<sub>1</sub>** Collisions occur if we want to insert  $(k, v)$  into the table, but  $T[h(k)]$  is already occupied.

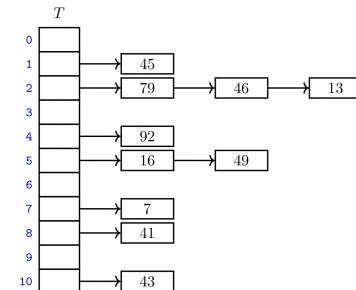
**E<sub>2</sub>** Solutions:



### HASHING WITH CHAINING

**E<sub>1</sub>** Idea: use lists to resolve collisions.

**E<sub>2</sub>** Each slot stores a list.



\* we use the MTF heuristic when inserting.

**E<sub>3</sub>** 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

**E<sub>1</sub>** we randomly pick the hash function among all possible hash functions uniformly.

**E<sub>2</sub>** This is called the "uniform hashing assumption".

**E<sub>3</sub>** Under the uniform hashing assumption,

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

list at  $T[i]$

Here, " $\alpha$ " 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  $O(\alpha)$ .

However, if  $k$  is in the dictionary, then

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

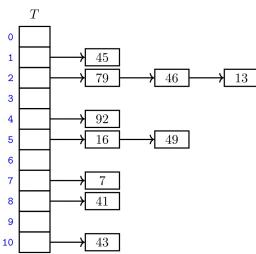
$\uparrow$   
 $k$  is in  
the bucket  
 $\leq \alpha$

Hence, under uniform hashing,

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

## REHASHING

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



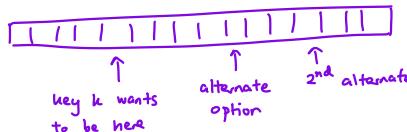
**💡** If with rehashing we keep  $O(1)$ , then all operations take  $O(1)$  time.

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

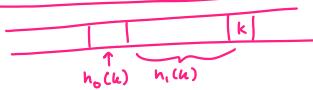
$$h(k,i) = (h(k) + c_1 i + c_2 i^2) \bmod M$$

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

## DOUBLE HASHING

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

$$h(k, i) = (h_0(k) + i h_1(k)) \bmod M$$



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

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

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

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)

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

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

## CHOOSING HASH FUNCTIONS

### MODULAR METHOD

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

Here, we let

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

$A \in [0, 1)$

$M \in [0, m)$

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

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

### HASHING MULTI-DIMENSIONAL DATA

Suppose we wanted to hash keys that are words (ie in  $\Sigma^*$ ).

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

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

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

ie  $(((((65R+80) \bmod M) R+80) \bmod M) R+76) \bmod M$

### RANDOMLY-CHOSEN HASH FUNCTIONS

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

But then we cannot compute the hash value quickly!

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

$$\begin{aligned} \Rightarrow (ak + b) \bmod p &= (ak' + b) \bmod p && (\because p \geq \text{mod } p) \\ \Rightarrow ak + b - ak' - b &\equiv_p 0 && (\equiv_p \Leftrightarrow \text{mod } p) \\ \Rightarrow a(k - k') &\equiv_p 0 \end{aligned}$$

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".

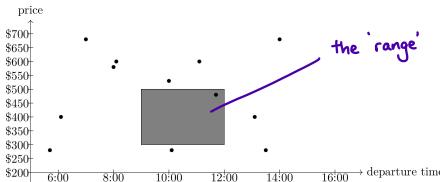
## MULTI-DIMENSIONAL DATA

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

$$A = [x, x'] \times [y, y']$$

- ② We want the points in A.

Alternative depiction:

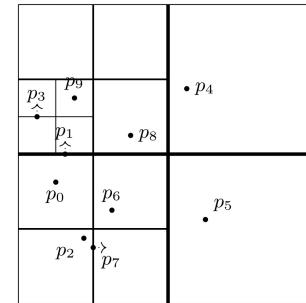
## QUAD-TREES

For "quad-trees":

- ① stores points in 2D;
- ② we assume the points are in a

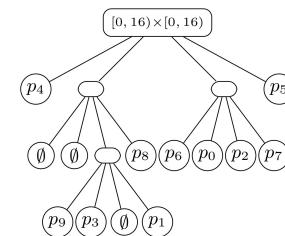
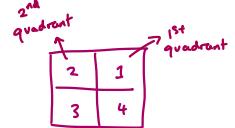
bounding box  $[0, 2^2] \times [0, 2^2]$

eg  $[0, 16] \times [0, 16]$



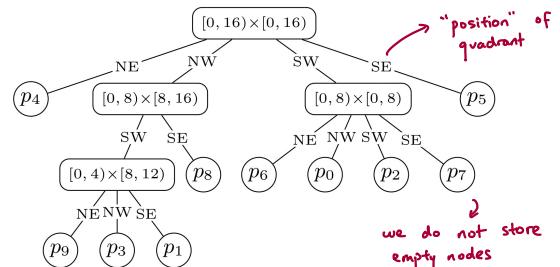
How to construct?

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



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

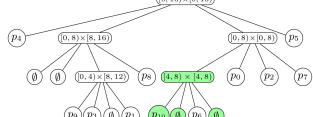
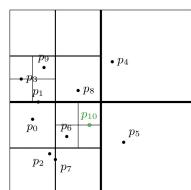
Alternative depiction:



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

# OPERATIONS OF QUAD-TREES

**Q1** Inset:



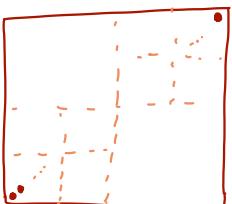
**Q2** Run-time:

- 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.

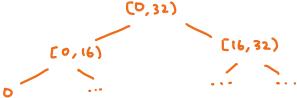
**Q2** However, we can bound the height if the coordinates are integers; specifically in the range  $\{0, \dots, 2^l - 1\}$ .

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

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

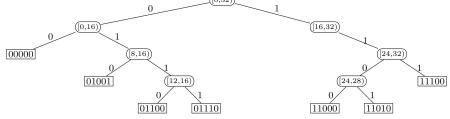


$\downarrow$



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 trie!

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  $l$

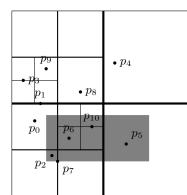
This argument generalizes to higher dimensions (e.g. 2D, i.e. quad trees)

**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:

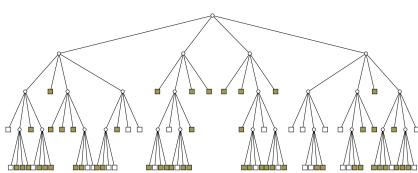
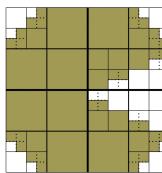
① We have no bounds in terms of  $n$  or  $s$ .

② Only thing we can say:

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

## OTHER USES OF QUAD-TREE

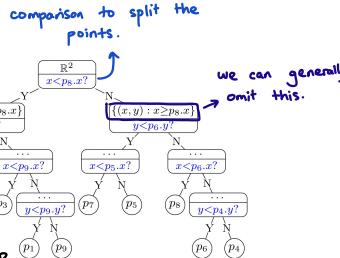
**Q1** We can use quad-trees to store pixelated images:



## kd-TREE

\* d = "dimensional"

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



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.

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

💡 Under this, each split leads to  
①  $\lceil \frac{n}{2} \rceil$  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; &  
 $\rightarrow O(n)$  expected
- ② Recurse in children  
 $\rightarrow$  size  $\approx \frac{n}{2}$  each.

## RANGE-SEARCH FOR kd-TREES

💡 Pseudocode:

Algorithm 8.2:  $kdTree::rangeSearch(r \leftarrow root, A)$

```

Input : Node r of a kd-tree. Query-rectangle A
1  $R \leftarrow$  region associated with node r // inside node
2 if  $R \subseteq A$  then
3   report all points in subtree at r and return // outside node
4 else if  $R \cap A$  is empty then
5   return // boundary node
6 else
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

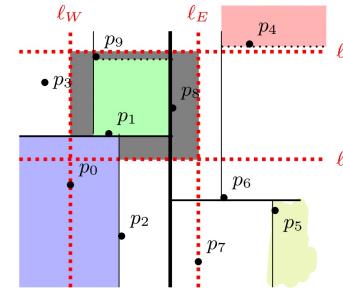
💡 Then, we see that

# of boundary nodes  $\in O(\sqrt{n})$ ,

so

run-time  $= O(\sqrt{n} + s)$ .

Proof:



How many nodes have an associated region that intersects one of these lines?

If the point does not intersect, then either  
- it is an outside region (eg  $p_5$ ); or  
- it is an inside region (eg  $p_1$ ):

ie  
# boundary nodes  $\leq$  nodes whose associated region intersects one of  $\ell_w, \ell_E, \ell_N, \ell_S$ .

Let  $Q(n, \ell) =$  # of boundary nodes.

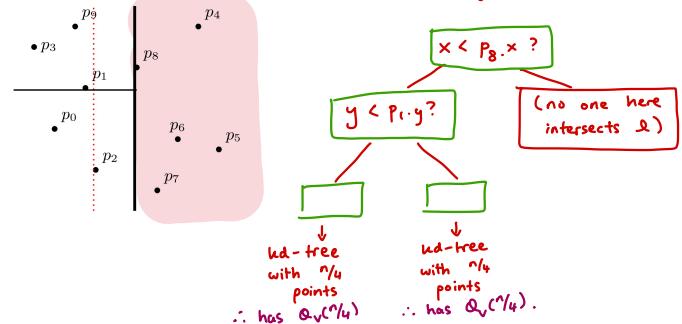
Let  $Q(n, \ell) = \max_{\text{kd-trees w/ } n \text{ pts}} \text{ that intersect a given line } \ell$ .  
This is independent of  $\ell$  (shift points), so only consider whether  $\ell$  is horizontal/vertical  $\rightsquigarrow Q_v(n), Q_h(n)$ .

Then  
 $Q(n) \leq Q(n, \ell_w) + Q(n, \ell_N) + Q(n, \ell_E) + Q(n, \ell_S)$   
 $\leq 2Q_v(n) + 2Q_h(n)$ .

So, if we prove  $Q_v(n) \in O(\sqrt{n})$  (and symmetrically  $Q_h(n) \in O(\sqrt{n})$ ), we are done.

See that  $Q_v(n) =$  # of nodes whose associated region intersects a vertical line  $\ell$ .

Consider making the kd-tree:

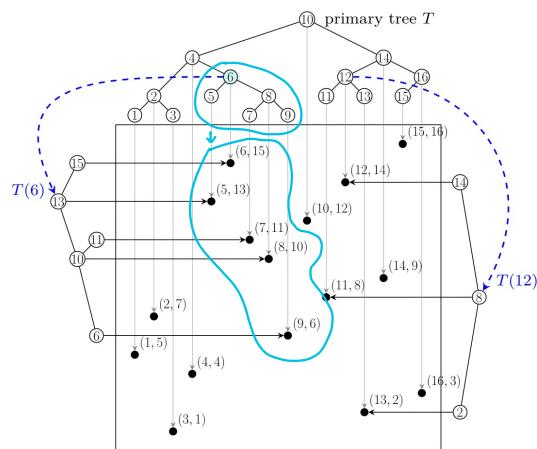


Therefore

$$Q_v(n) \leq 2Q_v(\frac{n}{4}) + 2.$$

This resolves to  $O(\sqrt{n})$ , as needed.  $\blacksquare$

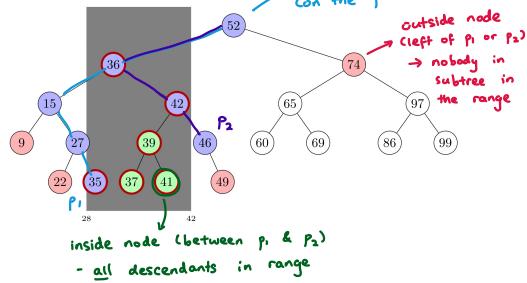
# RANGE TREE



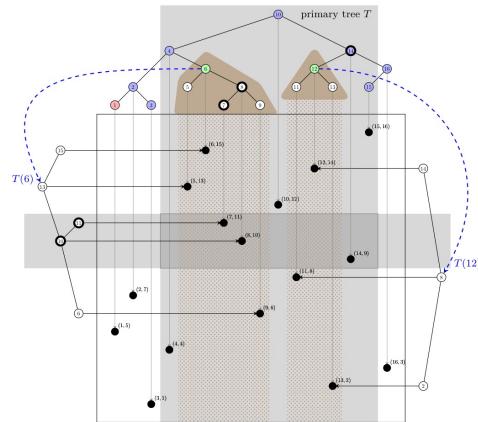
## RANGE-SEARCH

💡 Steps:

- ① A range-search in the primary tree (with a small twist) by the x-coord.
  - search for left boundary (gives a path  $P_1$ )
  - search for right boundary (gives a path  $P_2$ )



- ② Return all the boundary nodes & the top-most inside nodes.
- ③ For the boundary nodes, we explicitly check.
- ④ For the inside nodes,
  - all points are in range with x-coord.
  - so, we run a range-search on the associated tree by the y-coord.



💡 Run-time:  $O(\log^2 n + s)$ .

- we do one range-search in the primary tree ( $O(\log n)$ )
- we do range searches in  $O(\log n)$  associated trees (each of these takes  $O(\log n)$  time  $\rightarrow O(\log^2 n)$  time)

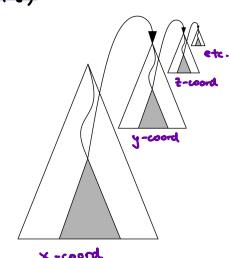
- report all points in range:

- $O(\log^2 n)$
- boundary nodes in primary tree
  - boundary nodes in same associated tree we searched in
  - only need → inside nodes in some associated tree to report

## HIGHER DIMENSIONS

💡 Run-time of range search:  $O(\log^d n + s)$

💡 Space:  $O(\log^d n \cdot n)$



# 3-SIDED RANGE SEARCH

Idea: return  $(x,y)$  with  $x_1 \leq x \leq x_2$  and  $y \geq y'$ .

Naive approach: Range-tree.

- runtime of range-search =  $\Theta(\log^2 n + s)$

- space =  $\Theta(n \log n)$ .

## ASSOCIATED HEAPS

Idea:

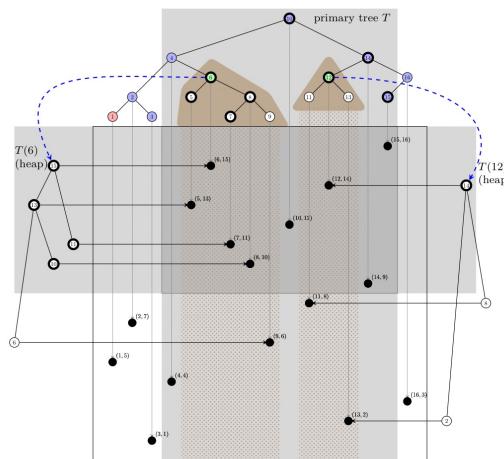
① Primary tree: balanced BST

② Associated tree: binary heap

③ Space:  $\Theta(n \log n)$

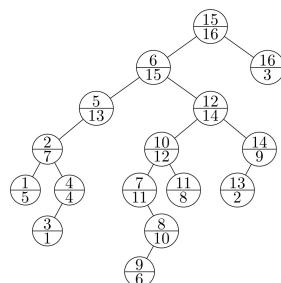
Range search:

- ① Search in primary tree as before; &
- ② In associated heaps: search by y-coordinate in  $O(1+s)$  time.
- ③ Run-time:  $O(\log n + s)$ .



## CARTESIAN TREES

Idea: Use a tree, but use x-coordinate as the key & the y-coordinate as the priority.



Space:  $\Theta(n)$

Range-search:

- ① Do BST::range-search( $x_1, x_2$ ) to get boundary and topmost inside nodes.
- ② For each inside node, do 1-sided search in heap at that node by y-coordinate.
- ③ Run-time:  $O(\text{height} + s)$

## PRIORITY SEARCH TREES

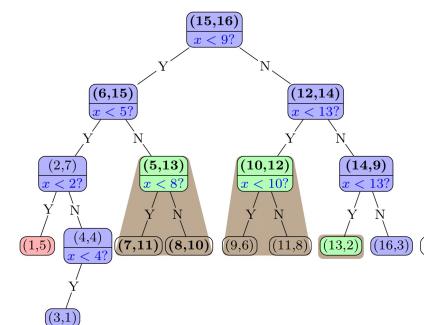
Idea:

① Store y-coordinates in "heap-order"

② Split in half by x-coordinate

- split-line coordinates: use median x-coordinate among the points in subtree

- so height =  $\lceil \log n \rceil$



Note: x-coordinate stored for splitting can be different from the x-coordinate of the stored point.

Range-search time:  $O(\log n + s)$ .

Space:  $\Theta(n)$ .

## SUMMARY OF 3-SIDED RANGE SEARCH ADTS

Summary:

	associated heaps	Cartesian tree	priority search tree
time for range-search	$O(\log n + s)$ (optimal)	$O(\text{height} + s)$ (height is bad)	$O(\log n + s)$ (good!)
Space	$\Theta(n \log n)$ (bad!)	$\Theta(n)$ (good!)	$\Theta(n)$ (good!)
insert/delete	$\Theta(\log^2 n)$ (amort.)	$\Theta(\text{height})$	$O(\log n)$ (feasible - worst case)
build from scratch	$\Theta(n \log n)$	sort by x-coord + $O(n)$ time	$O(n \log n)$ or sort + $O(n)$

# Chapter 9: Pattern Matching

Idea:

- ① Given text  $T[0 \dots n-1]$  of length  $n$  & a pattern  $P[0 \dots m-1]$  of length  $m$
- ② Want to know: is  $P$  a substring of  $T$ ?

eg  $T: ab\boxed{ba}baab$      $P: ab$   
 substring prefix    substring    substring suffix

Q In this course: report one occurrence  
 (usually the leftmost).  
 (CIRL: report all of them)

## BRUTE-FORCE

Idea: check every possible guess.

Algorithm 9.1: `bruteForce::patternMatching(T, P)`

```
Input : Text T of length n, Pattern P of length m
1 for i ← 0 to n - m do
2   if strcmp(T[i..i+m-1], P) = 0 then return found at guess i
3 return FAIL
```

Run-time:

- ① `strcmp` takes  $\Theta(m)$  time.
- ② Worst possible input:  $P = a^m, T = a^n$
- ③ So worst case runtime =  $\Theta((n-m+1)m)$
- ④ This is  $\Theta(mn)$  if  $m = \lceil n/2 \rceil$ .

## PRE-PROCESSING IDEA

Idea: break a problem into 2 parts:

- ① Build a data structure/info that will make later queries easy.  
 (this part can be slow)

} preprocessing

- ② Do the actual query.  
 (this part can be fast)

Q How to improve?

- ① Do extra preprocessing on the pattern  $P$ 
  - eliminate guesses based on completed matches & mismatches
  - Karp-Rabin, Boyer-Moore, DFA, KMP (length  $m$ )
  - used for web searches

- ② Do extra preprocessing on the text  $T$ 
  - we create a data structure to find matches easily
  - Suffix tree, Suffix array (length  $n$ )
  - bioinformatics

## KARP-RABIN FINGERPRINT ALGORITHM

Idea: Use hashing to eliminate guesses

- ① Compute hash function for each guess, compare with pattern hash

- ② If values are unequal, guess cannot be an occurrence

eg  $P = 5 \ 9 \ 2 \ 6 \ 5 \quad T = 3 \ 1 \ 4 \ 1 \ 5 \ 9 \ 2 \ 6 \ 5 \ 3 \ 5$

Use standard hash function: flattening + modular (radix  $R=10$ ):

$$h(x_0 \dots x_4) = (x_0 \dots x_4)_{10} \bmod 97$$

$$\Rightarrow h(P) = 59265 \bmod 97 = 95.$$

T:	3	1	4	1	5	9	2	6	5	3	5
	hash-value 84										
		hash-value 94									
			hash-value 76								
				hash-value 18							
					hash-value 95						

If  $P$  occurs at guess  $i \Rightarrow h(T[i \dots i+m-1]) = h(P)$ .

Q First attempt:

Algorithm 9.2: `KarpRabin::patternMatching-naive(T, P)`

```
1  $h_P \leftarrow \text{word-hash-value}(P, 10, M)$ 
2 for  $i \leftarrow 0$  to  $n - m$  do
3    $h_T \leftarrow \text{word-hash-value}(T[i..i+m-1], 10, M)$ 
4   if  $h_T = h_P$  then
5     if strcmp(T, P, i, i+m-1) = 0 then
6       return "found at guess i"
7 return FAIL
```

- never misses a match

-  $h(T[i \dots i+m-1])$  depends on  $m$  characters, so naive computation takes  $\Theta(m)$  time per guess

- running time =  $\Theta(mn)$  if  $P$  not in  $T$ .

Q Idea to be faster: our hash function

is a "rolling" hash function.

- ie given  $h(T[i \dots i+m-1])$ , compute  $h(T[i+1 \dots i+m])$  quickly.

eg Know  $h(41592) = 76$ , what is  $h(15926)$ ?

$$15926 \bmod 97 = [41592 \bmod 97 - 4(10000) \bmod 97] \bmod 97$$

- in general,

$$T[i \dots i+m] \bmod M = [(T[i \dots i+m-1] \bmod M - T[i] \cdot 10^m \bmod M) * 10 + T[i+m]] \bmod M$$

↑ next hash value      h(previous guess)      precompute

Q Second attempt:

Algorithm 9.3: `KarpRabin::patternMatching(T, P)`

```
Input :  $T$  and  $P$  are texts of length  $n$  and  $m$  over alphabet  $\Sigma = \{0, \dots, R-1\}$ 
1 bool needToReset ← TRUE
2 for  $i \leftarrow 0$  to  $n - m$  do
3   if needToReset then // get random prime and initialize fingerprints
4      $M \leftarrow$  prime number randomly chosen in  $\{1, \dots, mn^2 \log R\}$ 
5      $h_P \leftarrow \text{word-hash-value}(P, R, M)$ 
6      $h_T \leftarrow \text{word-hash-value}(T[i..i+m-1], R, M)$ 
7      $s \leftarrow \text{word-hash-value}(10..0, R, M)$  // passed word has  $m$  chars
8     needToReset ← FALSE;
9   else // get next fingerprint from previous
10     $h_T \leftarrow ((h_T - T[i-1] \cdot s) \cdot R + T[i+m-1]) \bmod M$ 
11   if  $h_T = h_P$  then
12     if strcmp(T, P, i, i+m-1) = 0 then
13       return "found at guess i"
14     else
15       needToReset ← TRUE;
16 return FAIL
```

- choose "table size"  $M$  to be a random prime in  $\{2, \dots, mn^2\}$

-  $T^{\text{exp}} = O(mn)$

- improvement: reset  $M$  if no match at  $h_T = h_P$ .

# KARP-RABIN ANALYSIS

If we never match hash-value:

①  $\Theta(m)$  preprocessing

②  $\Theta(1)$  per guess

$\therefore$  total runtime =  $\Theta(mn)$ .

If  $h_p = h_T$  and there was a match:

①  $\Theta(m)$  preprocessing

②  $\Theta(1)$  per guess

③  $\Theta(m)$  strcmp & break

$\therefore$  total run-time =  $\Theta(mn)$

Worst-case: If  $h_p = h_T$  but no match  
(false positive)

- if this happens  $\Theta(n)$  times, then run-time  
is  $\Theta(m(n-m))$ .

We can show

$$P(\geq 1 \text{ false positive}) \leq \frac{2c}{n}$$

for a constant  $c$ .

Thus

$$\begin{aligned} T^{\text{exp}} &\leq P(\text{false positive}) \cdot (n+m) + P(\text{false positive}) \cdot n \cdot m \\ &\leq n+m + 2c(m) \\ &\in O(n+m). \end{aligned}$$

Auxiliary space:  $O(1)$ .

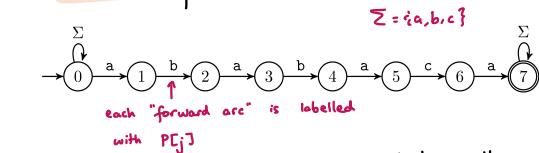
## DFA / NFA

Idea: Use a DFA to pattern match.

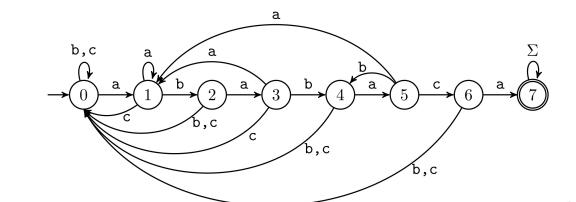
- we parse  $T$ , and reach state  $(M) \Leftrightarrow$

$P$  occurs in  $T$ .

We set our DFA up so state  $j$  means  
we have just now seen  $P[0 \dots j-1]$  in  
what was parsed.



We can show there is an equivalent small DFA.

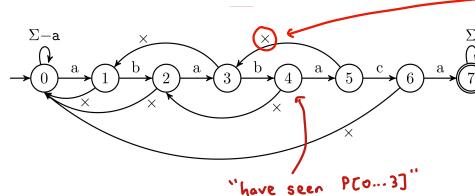


- then, it is easy to test whether  $P$  is in  $T$ .

# KNUTH-MORNS-PRATT / KMP

💡 DFA: would work, but complicated.

💡 Idea: use a new kind of finite automata.



- "failure":
- use this transition if no other fits
- does not consume a character
- every state has  $\leq 1$  failure arc

💡 With these rules, computations of the automaton are deterministic.

💡 We want to set up the failure arcs so that  $P$  is in  $T \Leftrightarrow$  we are at  $(M)$  at some point.

## KMP ALGORITHM

💡 Pseudocode:

preprocessing

```
Algorithm 9.4: KMP::patternMatching(T, P)
1  $F \leftarrow \text{failureArray}(P)$  // index of current character of  $T$  to parse
2  $i \leftarrow 0$  // index of current state
3  $j \leftarrow 0$  // index of current state
4 while  $i < n$  do
5   if  $P[j] = T[i]$  then → "can I use the forward arc"
6     if  $j = m - 1$  then
7       return "found at guess  $i - m + 1$ " // all characters were matched
8     else
9        $i \leftarrow i + 1$  // go along forward-arc
10       $j \leftarrow j + 1$  [use forward arc]
11    else
12      if  $j > 0$  then
13         $j \leftarrow F[j - 1]$  // go along failure-arc from state  $j$ 
14      else
15         $i \leftarrow i + 1$  // at leftmost state; go along loop
16 return FAIL
```

- follow failure arc
- $F[j-1]$ : stores where the failure arc from state  $j$  goes
- note the  $j-1$ !!
- we are at state 0, the initial state
- just consume (character is "useless")

💡 We can show

$$\# \text{ of executions of while loop} \leq 2i - j + \chi(\text{found match})$$

where  $\chi(\text{found match})$  is 1 if we found a match, & 0 otherwise.

Proof: Initially,  $i=j=0$  & no executions of while-loop.

Assume we had  $\leq 2i-j$  executions at some later time

(the "inductive hypothesis").

In the next execution, we had 4 cases: (see above)

① We found pattern  $P$ .

Then

$$\# \text{ of executions} \leq 2i - j + 1 = 2i - j + \chi(\text{found match}).$$

② We use the forward arc.

$$\Rightarrow i' = i+1, j' = j+1 \Rightarrow 2i' - j' = 2i - j + 1$$

③ We use the failure arc.  $\geq$  # of executions.

$\Rightarrow j$  decreases,  $i$  is the same

$$\text{So } 2i' - j' \geq 2i - j + 1$$

④ We consume  $T[i]$ .

$$\Rightarrow i \text{ increases, } j \text{ stays the same}$$

$$\text{So } 2i' - j' \geq 2i - j + 1.$$

Thus,

run-time without computation of failureArray  
is  $O(\max i) = O(n)$ .

## KMP FAILURE ARRAY

**B1** Assume we reach state  $j+1$  and now have a mismatch.



- we can eliminate "shift by 1" if  $P[1..j] \neq P[0..j-1]$ .

- in general, we can eliminate "shift by  $k$ " if  $P[1..j]$  does not end with  $P[0..j-k]$

**B2** So, we want the longest prefix  $P[0..l-1]$  that is a suffix of  $P[1..j]$ ;

**B3** Thus our failure-function is defined by

$$F[j] = \begin{cases} \text{head of failure arc from state } j+1 \\ = \text{length of longest prefix of } P \text{ that is a} \\ \text{suffix of } P[1..j]. \end{cases}$$

## COMPUTING THE FAILURE-ARRAY QUICKLY

**B1** Recall: we reach state  $l$

$\Leftrightarrow$  we have seen  $P[0..l-1]$

$\Leftrightarrow P[0..l-1]$  is a suffix of what was parsed

**B2** Let's instead consider parsing  $P[1..j]$ .

So

$$F[j] = \text{state reached when parsing } P[1..j] \text{ on} \\ \text{the KMP-automaton for } P.$$

**B3** Pseudocode:

```
Algorithm 9.5: KMP::failureArray(P)
1  $F \leftarrow$  array of length  $m$ ; initialize  $F[0] \leftarrow 0$ 
2  $j \leftarrow 1 \rightarrow$  start at index // index of current character of  $P[1..m]$  to parse
3  $\ell \leftarrow 0 \quad |$  // index of current state
4 while  $i < m$  do
5   if  $P[\ell] = P[j]$  then // does the next character
6      $\ell \leftarrow \ell + 1$  // go along forward-arc
7      $F[j] \leftarrow \ell$  // Have parsed  $P[1..j]$  and reached  $\ell$ 
8      $j \leftarrow j + 1$ 
9   else
10    if  $\ell > 0$  then // go along failure-arc
11       $\ell \leftarrow F[\ell - 1]$ 
12    else // at leftmost state
13       $F[j] \leftarrow 0$ 
14       $j \leftarrow j + 1$ 
```

set failure function

- in particular, we can use the KMP-automaton to build itself.

**B4** Run-time = time to parse  $P[1..m-1] = O(m)$  time.

**B5** Therefore,

$$\begin{aligned} \text{total worst-case time for KMP} &= O(m+n) \\ \text{aux. space} &= O(m) \quad (\text{failure array}) \end{aligned}$$

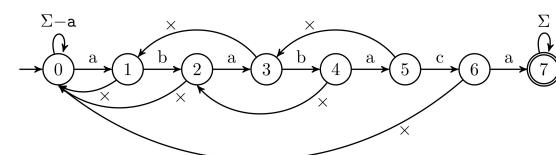
## USING KMP-AUTOMATON TO COMPUTE THE DFA IN $O(n)$ TIME

**B** Note: from the KMP-automaton, we could compute the DFA in  $O(|\Sigma| \cdot m)$  time.

## IMPROVING THE FAILURE FUNCTION

**B1** We can also improve the failure-function.

**B2** Consider the "bad case":



Suppose we reach state 4, and the next char is not a.

$\rightarrow$  failure arc goes to state 2

$\rightarrow$  forward arc at state 2 also has 'a'.

We will fail again!  $\Rightarrow$  go to  $F[4]$  right away.

**B3** Improved failure-function:

$$F^+[j] = \begin{cases} \text{length } l \text{ of the longest prefix of } P \text{ that} \\ \text{is a suffix of } P[1..j] \text{ and } P[l] \neq P[j+1]. \\ \text{or } 0 \text{ if no such } l \text{ exists} \end{cases}$$

**B4**  $F^+$  takes into account whether forward arcs mismatch.

**B5** To compute  $F^+$ :

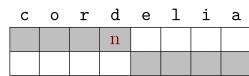
$$F^+[j] = \begin{cases} F[j], & P[j+1] \neq P[F[j]] \text{ or } F[j] = 0 \\ F^+[F[j]-1], & \text{otherwise} \end{cases}$$

# BOYER-MOORE

## REVERSE SEARCHING

💡 Idea of reverse searching:

- ① compare from right to left
- ② if we have a mismatch, shift so the new guess fits the character of T.



## BAD CHARACTER HEURISTIC

💡 We use the "bad character heuristic":

$P$ : p a p e r	
$T$ : f e e d a l l p o o r p a r r o o t s	

shift so [a] in paper is here.  
shift as little as possible (so we shift to the "rightmost" p)  
shift "past" the character if char not in P.

- we shift the guess so that  $T[i]$  matches the last occurrence in  $P$  (ie  $P[L[T[i]]]$ .)
- $L$  = last occurrence array.

## LAST OCCURRENCE ARRAY

💡 We need to precompute the "last occurrence array":

chart c	p	a	e	r	all others
$L[c]$	2	1	3	4	-1

if no match:  
- shift "past" the character.

- 💡 If we were comparing  $T[:]$  with  $P[j:]$ , and  $\ell = L[T[i]]$ , then we want to shift  $j-\ell$  units.

## SIMPLIFIED BOYER-MOORE

💡 Pseudocode:

```
Algorithm 9.6: BoyerMoore::patternMatching( $T, P$ )
1  $L \leftarrow$  last occurrence array computed from  $P$  - preprocessing
2 for  $i \leftarrow 0 \dots n-1$  do
3  $i \leftarrow m-1$  // currently inspected character of  $T$ 
4  $j \leftarrow m-1$  // currently inspected character of  $P$ 
5 while  $i < n$  and  $j \geq 0$  do
   // invariant: The current guess begins at  $T[i-j]$ 
6   if  $T[i] = P[j]$  then // match, go one character leftward
7      $i \leftarrow i-1$ 
8      $j \leftarrow j-1$ 
9   else
10     $i \leftarrow m-1+i-\min\{L[T[i]], j-1, S[j]\}$  // mismatch, find next guess to try
11     $j \leftarrow m-1$ 
12 if  $j = -1$  then
13   return "found at guess  $i+1$ "
14 else
15   return FAIL
```

### Algorithm 9.7: BoyerMoore::lastOccurrenceArray( $P$ )

```
1  $L \leftarrow$  array indexed by alphabet  $\Sigma$ 
2 Initialize  $L$  to be  $-1$  everywhere
3 for  $i \leftarrow 0 \dots m-1$  do  $L[P[i]] \leftarrow i$ 
4 return  $L$ 
```

💡 If we only do

- ① reverse-order searching; &
- ② bad-character jumps.

this is called the "simplified Boyer-Moore" algorithm.

💡 This works well in practice  
(skips ~25% of  $T$  in experiments)

Worst case run-time =  $O(mn)$ .

## GOOD SUFFIX ARRAY

💡 Idea: the matched suffix

$P = onobobo$	
o n o o b o b o n o i n b o b o l a n d	
(o) (b) (o) b (o) o	
↑	
(b) (o) b (o) o	
↑	
(o) (b) (o) b (o)	
↑	
(b) (o)	

we shift here to get the "last" occurrence of the suffix

- ① we use the same strategy as KMP:  
use what we matched to eliminate some guesses.

$P[j] \leftarrow Q = P[j+1 \dots m-1] \rightarrow$	
x ✓ ✓ ✓ ✓ ✓	

- ② if there is no matched string, move forward by 1 (the last-occurrence array might help).

- ③ if the matched string is not again in the pattern, try matching a suffix of the matched string to a prefix of  $P$ .

- ④ If only the empty suffix fits, shift past the word.

- ⑤ if  $Q = ""$ , use the last occurrence.

💡 Pseudocode:

```
Algorithm 9.6: BoyerMoore::patternMatching( $T, P$ )
1  $L \leftarrow$  last occurrence array computed from  $P$ 
2  $S \leftarrow$  good suffix array computed from  $P$ 
3  $i \leftarrow m-1$  // currently inspected character of  $T$ 
4  $j \leftarrow m-1$  // currently inspected character of  $P$ 
5 while  $i < n$  and  $j \geq 0$  do
   // invariant: The current guess begins at  $T[i-j]$ 
6   if  $T[i] = P[j]$  then // match, go one character leftward
7      $i \leftarrow i-1$ 
8      $j \leftarrow j-1$ 
9   else
10     $i \leftarrow m-1+i-\min\{L[T[i]], j-1, S[j]\}$  // shift according to L & S
11     $j \leftarrow m-1$ 
12 if  $j = -1$  then
13   return "found at guess  $i+1$ "
14 else
15   return FAIL
```

## WILDCARDS: $P^*$

💡 We have 3 main cases, depending on how much of  $Q$  exists again in  $P$ .

💡 To unify them, we add "wildcards" to  $P$ .

$$P^* = \underbrace{\_ \_ \dots \_ \_}_{m} P[0 \dots m-1]$$

where a wildcard matches any character.

💡 Then, the 3 cases become

- ①  $Q$  is a substring of  $P$ .

$P[j] \leftarrow Q = P[j+1 \dots m-1] \rightarrow$	
x ✓ ✓ ✓ ✓ ✓	

- ② A suffix of  $P$  is a prefix of  $Q$ .

$P[j] \leftarrow Q = P[j+1 \dots m-1] \rightarrow$	
x ✓ ✓ ✓ ✓ ✓	
* * *	← wildcards → $P[0]$

- ③ No character of  $Q$  is matched.

$P[j] \leftarrow Q = P[j+1 \dots m-1] \rightarrow$	
x ✓ ✓ ✓ ✓ ✓	
* * * * *	$P^*[0] \leftarrow$ add wildcards → $P[0]$

💡 Then,  $Q$  is a substring of  $P^*$  and take the rightmost occurrence except not at the very end.

$P^*:$

index	-m	-2	-1	0	m-1
*	...		* *	$P[0]$	...

let the "beginning" position to match be  $\ell+1$ .

Then  $Q$  is a prefix of  $P^*[\ell+1 \dots m-1]$ ;  
ie  $P[\ell+1 \dots m-1]$  is a prefix of  $P^*[\ell+1 \dots m-1]$ .

# COMPUTING THE GOOD SUFFIX ARRAY AS A HUMAN

j	$P[j+1..m-1]$	$P^*[-m..m-2]$										match starts at index	$S[j]$	
		*	*	*	*	*	*	*	b	o	o	b		
5	o												4	3
4	bo								b	o			3	2
3	obo								o	b	o		2	1
2	bobo				b	o	b	o					-2	-3
1	obobo			o	b	o	b	o					-3	-4
0	oobobo		o	o	b	o	b	o					-4	-5

- we use

$$S[j] = \max_l P[j+l..m-1] \text{ is a prefix of } P^*[l+1..m-2].$$

## FINDING $S[j]$ IN LINEAR TIME

Θ₁ We note

$$\begin{aligned} S[j] &= \max_l P[j+l..m-1] \text{ is a prefix of } P^*[l+1..m-2] \\ &= \max_l P[m-1..j+1] \text{ is a suffix of } P^*[m-2..l+1] \\ &\quad "P[j+1..m-1]^\text{rev}" \quad "P[l+1..m-2]^\text{rev}" \\ &\quad \text{the "reverse" string} \\ &= \max_l \left\{ \begin{array}{l} \text{when parsing } P^*[m-2..l+1] \text{ at the KMP-automaton} \\ \text{for } P[m-1..j+1], \text{ then we reach state } j+1. \end{array} \right\} \end{aligned}$$

Θ₂ This evaluates to

$$S[m-q-1] = m-2 - \min_k \left\{ \begin{array}{l} \text{parsing } R[1..k] \text{ on the KMP-automaton for} \\ R = P^\text{rev} * * \dots * * \text{ brings us to state } q \end{array} \right\}.$$

Θ₃ We can do this by parsing in  $O(m)$  time, to compute  $S$ .

## SUMMARY

Θ₁ With the good suffix heuristic, we can skip even more.

Θ₂ But, we can skip even more using  $S^+[j]$  vs.  $S[j]$ .

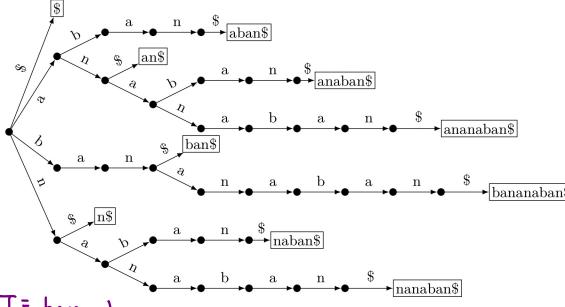
Θ₃ This algorithm can be made to run in  $O(m+n)$  time.

# SUFFIX TREE

- 💡 Idea: Is  $P$  a substring of text  $T$ ?  
 ↳ is  $P$  a prefix of a suffix of  $T$ ?  
 So, we can store all suffixes of  $T$ , and then do a prefix-search for  $P$ .

## SUFFIXES IN A TRIE

- 💡 We can store the suffixes in a trie:



$T = \text{bananaban}$

- 💡 To see if  $P$  is in  $T$ :  
 ① Traverse the trie using  $P$ .  
 ② If we are still "in" the trie after ①, then  $P$  is a substring of  $T$ .

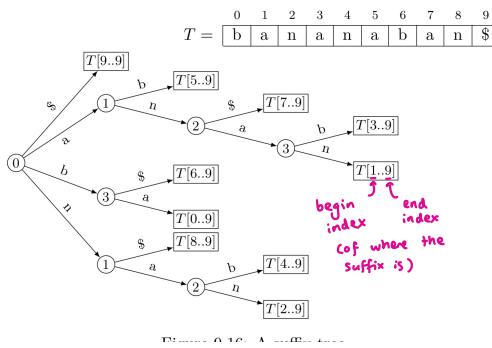
💡 Run-time:  $O(|P|) = O(m)$ .

💡 But the price to pay:

- $T$  has length  $n$
- so we have  $n+1$  suffixes of length  $0, 1, \dots, n$
- thus, the trie can have  $\Omega(n^2)$  nodes.

## SUFFIX TREE

- 💡 We can instead store the suffixes in a compressed trie:



- 💡 Prefix-search for  $P$  takes  $O(|P|) = O(m)$  time.  
 (details omitted)

💡 We can build the trie in  $O(n)$  time.

(no details - too complex!)

💡 Space:  $O(n)$ .

## SUFFIX ARRAY

- 💡 We can also store the suffixes in a "suffix array": the sorting permutations of the suffixes.

i	
0	bananaban\$
1	ananaban\$
2	nanaban\$
3	anaban\$
4	naban\$
5	aban\$
6	ban\$
7	an\$
8	n\$
9	\$

j	$A^s[j]$	↓ start index of suffix
0	9	\$
1	5	aban\$
2	7	an\$
3	3	anaban\$
4	1	ananaban\$
5	6	ban\$
6	0	bananaban\$
7	8	n\$
8	4	naban\$
9	2	nanaban\$

the sorting permutation (ie the "suffix array").

- 💡 Space:  $O(n)$  (and a very small constant).

## Building:

- ① Sort (using MSD-radix-sort)
  - run-time =  $O(n^2)$  worst case, but usually faster.
- ② A special algorithm
  - run-time =  $O(n \log n)$
  - covered in CS 482

## Pattern-matching:

- use binary search.

$\ell \rightarrow$	$j$	$A^s[j]$
1	0	9
2	1	5
3	2	7
4	3	3
$\nu \rightarrow$	4	1
5	5	6
6	6	0
7	7	8
8	8	4
$r \rightarrow$	9	2

$\ell \rightarrow$	$j$	$A^s[j]$
1	0	9
2	1	5
3	2	7
4	3	3
5	4	1
$\nu = \ell \rightarrow$	5	6
$r \rightarrow$	6	0
7	7	8
8	8	4
$r \rightarrow$	9	2

$\ell \rightarrow$	$j$	$A^s[j]$
0	0	9
1	1	5
2	2	7
3	3	3
4	4	1
5	5	6
6	6	0
7	7	8
8	8	4
9	9	2

- 💡 Run-time:  $O(\log n)$  comparisons.

- comparisons = compare  $P$  to a suffix:  $\Theta(m)$

So run-time =  $\Theta(m \log n)$ .

(slightly slower than suffix-trees.)

# Chapter 10: Text Compression

Idea:

- ① Input: source texts  $S$  (huge!)
  - ② Output: map text  $S$  into a new text  $C$  (smaller).
- Note  $S$  lives on alphabet  $\Sigma_S$ , and  $C$  lives on alphabet  $\Sigma_C$ .

Objectives:

- ① Minimize the "compression ratio"; where

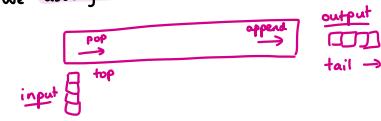
$$\text{comp. ratio} = \frac{|C|}{|S|} \log |\Sigma_C|$$

- ② Fast encoding & decoding.

All compressions we look at are lossless; we can get  $S$  back from  $C$  uniquely.

STREAMS

We usually store  $S$  and  $C$  as "streams".



Input stream: Read one char at a time via top/pop.  
- also supports 'isEmpty' and 'reset'.

Output stream: Write one character at a time (via append)

This is convenient for handling large texts (we start processing while loading)

CHAR-BY-CHAR ENCODING

Idea: assign to each  $c \in \Sigma_S$  a codeword  $E(c)$ .

char	A	B	C	D	...	(Caesar shift)
$E(c)$	B	C	D	E	...	

char	A	B	C	D	...	(ASCII)
$E(c)$	10000010	10000011		...		

The above are "fixed-length encodings": all codewords have the same length.

Better: "variable-length encoding" — more frequent characters get shorter codewords.

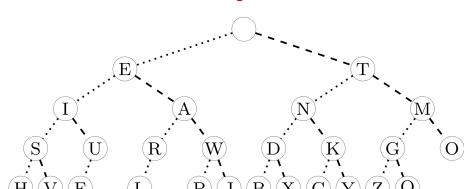
VARIABLE-LENGTH CODES

Overall goal: find a short encoding.

Idea: Some letters in  $\Sigma$  occur more often than others, so use shorter codes for more frequent characters.

Example: encoding the alphabet in Morse Code.

We build an "encoding trie":



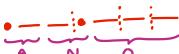
→ more frequent letters (eg E & T) go "higher" on the trie than less frequent ones (eg F)

MORSE CODE

Using the encoding trie in the previous section, we can convert the alphabet into Morse Code (eg if "left" = · & "right" = —)

Problem: Morse Code is not lossless!

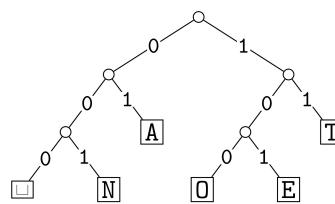
e.g. WATT



(unless we use an "end of char" pause)

PREFIX-FREE ENCODING/DECODING

- To achieve lossless encoding, we use an encoding trie where the codewords are prefix-free. (ie no codeword is a prefix of another codeword.) In particular, we make it so the encoded trie has codewords only at the leaves.
- Run-time to decode  $C$ :  $O(|C|)$   
Run-time to decode  $S$ :  $O(\sum_{c \in S} |E(c)|) = |C|$



ENCODING ALGORITHM

- Imagine we have some character encoding  $E: \Sigma_S \rightarrow \Sigma_C^*$ . Note:  $E$  is a dictionary with keys in  $\Sigma_S$ .

Algorithm 10.1:  $\text{charByChar::encoding}(S, C)$

```

Input : input-stream  $S$ , output-stream  $C$ 
// We assume the class stores an encoding dictionary  $E$ 
1 while  $S$  is not empty do
2    $c \leftarrow S.pop()$ 
3    $x \leftarrow E.search(c)$ 
4    $C.append(x)$ 

```

DECODING OF PREFIX-FREE CODES

Pseudocode:

Algorithm 10.2:  $\text{prefixFree::decoding}(C, S)$

```

Input : input-stream  $C$ , output-stream  $S$ 
// We assume the class stores a prefix-free encoding trie  $T$ 
1 while  $C$  is not empty do
2    $r \leftarrow T.root$ 
3   while  $r$  is not a leaf do
4     if  $C$  is empty or  $r$  has no child labeled with  $C.top()$  then
5       return "invalid encoding"
6     else
7        $r \leftarrow$  child of  $r$  that is labeled with  $C.pop()$ 
8    $S.append(\text{character stored at } r)$ 

```

Run-time:  $O(|C|)$ .

ENCODING FROM THE TRIE

We can also encode directly from the trie:

Algorithm 10.3:  $\text{prefixFree::encoding}(S, C)$

```

// We assume the class stores a prefix-free encoding trie  $T$ 
1 Initialize encoding-array  $E$  indexed by characters of  $\Sigma_S$ 
2 for all leaves  $l$  in  $T$  do  $E[\text{character at } l] \leftarrow l$  // setup
3 while  $S$  is non-empty do
4    $v \leftarrow E[S.pop()]$  // Find code-word  $w$  by going up from leaf  $v$ .
5    $w \leftarrow$  empty string
6   while  $v$  is not the root do
7      $w.prepend(\text{character on link from } v \text{ to its parent})$  ] run-time:  $O(|w|)$ 
8    $v \leftarrow$  parent of  $v$ 
9    $C.append(w)$ 

```

Run-time:  $O(|T| + |C|)$ .

(This is in  $O(|\Sigma_S| + |C|)$  if  $T$  has no node with child.)

## HUFFMAN - ENCODING

Idea: frequent chars should have short codewords.  
In particular,

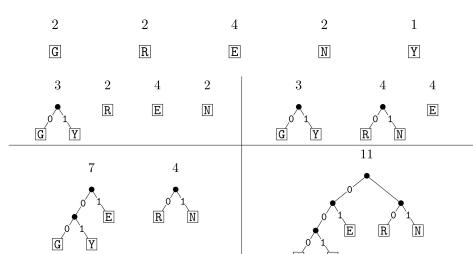
$$\begin{aligned} ICI &= \sum_{c \in S} I(c) \\ &= \sum_{c \in S} (\text{freq of } c \text{ in } S) \cdot I(c) \\ &= \sum_{c \in S} f(c) \cdot (\text{depth of } c \text{ in the } T) \\ &= \sum_{c \in S} f(c) \cdot d_T(c) \end{aligned}$$

We define

$$\text{cost}(T) = f(c) d_T(c).$$

Assumption:  $|S| \geq 2$ .

Idea:  
eg GREENENERGY,  $S = \{G, R, E, N, Y\}$   
char freq: G:2, R:2, E:4, N:2, Y:1



- ① Set up lots of 1-char tries;
- ② Find the 2 least frequent chars  $a, a'$  and make them siblings in the trie.



- ③ Repeat ②.

## PSEUDOCODE

```
Algorithm 10.5: Huffman::encoding(S, C)
Input : Input-stream S where S contains  $\geq 2$  distinct characters, output-stream C
1  $f \leftarrow$  array indexed by  $\Sigma_S$ , initially all-0 // compute frequencies
2 while  $S$  is non-empty do increase  $f[S.pop()]$  by 1
3  $T \leftarrow$  Huffman::buildTrie( $\Sigma_S, f$ ) // find encoding trie
4 Store  $T$  in the class
5  $C.append(trie.T encoded suitably)$ 
6 reset input-stream  $S$ 
7 prefixFree::encoding(S, C) // do actual encoding
```

## RUN-TIME OF HUFFMAN-ENCODING

- ① Priority queue has  $|S|$  items initially;
- ② We do delete in  $|S|$  time.  
- so, building  $T$  takes time  $O(|S| \log |S|)$ .
- ③ Rest of the encoding takes time  $O(|S| + ICI)$ .
- ④ But,  
① We need to send the tree along.  
② We need to go through  $S$  twice.  
→ but,  $\text{cost}(\text{Huffman tree})$  is small.

## $T_H$ HAS MIN COST AMONG ALL PREFIX-FREE BINARY ENCODING TRIES

The Huffman-tree  $T_H$  has minimum cost among all prefix-free binary encoding tries.

Proof: We will show for any prefix-free binary encoding tree  $T_0$ , we have

$$\text{cost}(T_H) \leq \text{cost}(T_0).$$

We do this by induction on  $|S|$ .

Base case:  $|S| = 2$ .



Step:

Huffman tree  $T_H$   
fix 2 chars  $a, a'$   
( $a$  &  $a'$  are siblings in  $T_H$ )

replace  
by

$T'_H$ : encoding trie for  
 $\Sigma_S \setminus \{a, a'\} \cup \{q\}$

some other  
trie  $T_0$

without increasing  
cost

encoding trie  $T_1$ : no node  
has exactly one child

w/o increasing cost

encoding trie  $T_2$ :  $a, a'$  are  
siblings

\* exchange  $aabb$ ,  
 $a'b'b$ , where  
 $b, b'$  are two  
siblings on the  
lowest level

$T'_1$ : encoding trie for  
 $\Sigma_S \setminus \{a, a'\} \cup \{q\}$

By induction,

$$\text{cost}(T'_H) \leq \text{cost}(T'_1)$$

$$\Rightarrow \text{cost}(T'_H) + f(a) + f(a') \leq \text{cost}(T'_1) + f(a) + f(a')$$

$$\Rightarrow \text{cost}(T_H) \leq \text{cost}(T_1) \quad (\leq \text{cost}(T_0)).$$

Proof follows.  $\square$

This is not necessarily true for other bases.

This also does not always make shorter.

## NO LOSSLESS COMPRESSION ALGORITHM HAS CMP RATIO $< 1$ FOR ALL INPUT STRINGS

No lossless compression algorithm can have compression ratio  $< 1$  for all input strings.

Proof: Assume  $\Sigma_S = \Sigma_C = \{0, 1\}$ . Assume we had such an algo.

Consider all the bitstrings of length  $n$ . There are  $2^n$  such strings.

The algorithm maps these to the bitstrings of length  $\leq n-1$ , of which there are

$$1 + 2^1 + \dots + 2^{n-1} = 2^n - 1 < 2^n.$$

This set is smaller; so, by Pigeonhole Principle,  $\exists$  bitstrings  $x_1, x_2$  of length  $n$  that both gets encoded as bitstring  $w$ . But this contradicts the losslessness of the algo!

Proof follows.  $\square$

# MULTI-CHARACTER ENCODING

**Idea:** Encode multiple characters (a substring of  $S$ ) using one code word.

## RUN-LENGTH ENCODING

**B1:** We assume  $\Sigma_S = \{0,1\}$

(the inputs are bitstrings).

**B2:** A "run" is a maximal substring that uses only one character.

eg  $S = \underbrace{1111}_{1} \underbrace{0001}_{3} \underbrace{111}_{4}$

**Idea:**

① We write down the lengths of the runs in  $S$ ;

& then

② we encode the string by these lengths.

③ We also need to specify the first bit we started with.

eg  $S = \underbrace{1}_{1} \underbrace{5}_{2} \underbrace{3}_{1} \underbrace{4}_{1}$   
first bit the run-lengths

**B4:** However, we want  $\sum_c$  to be finite, and preferably  $\{0,1\}$ .

**B5:** So, we need to map integers to bitstrings, so we do not need separations.

## ELIAS GAMMA CODING

**B1:** Idea: to encode  $k$ :

- ① Write  $\lfloor \log k \rfloor$  copies of 0; then
- ② The binary representation of  $k$  (always starts with 1). (has length  $\lfloor \log k \rfloor + 1$ .)

$k$	$\lfloor \log k \rfloor$	$k$ in binary	encoding $E(k)$
1	0	1	1
2	1	10	010
3	1	11	011
4	2	100	00100
5	2	101	00101
6	2	110	00110
⋮	⋮	⋮	⋮

**B2:** Pseudocode:

Algorithm 10.6: RLE::encoding( $S, C$ )

```

Input : Non-empty input-stream  $S$  of bits, output-stream  $C$ 
1  $b \leftarrow S.\text{top}()$  // bit-value for the current run
2  $C.append(b)$ 
3 while  $S$  is not empty do
4    $k \leftarrow 1$  // length of run
5   while  $S$  is not empty and  $S.pop() = b$  do  $k \leftarrow k + 1$ 
6    $K \leftarrow$  empty string // binary encoding of  $k$ 
7   while  $k > 1$  do
8      $C.append(0)$ 
9      $K.prepend(k \bmod 2)$ 
10     $k \leftarrow \lfloor k/2 \rfloor$ 
11    $K.prepend(1)$ 
12    $C.append(K)$ 
13    $b \leftarrow 1 - b$  // flip bit for next run

```

## DECODING

**B1:** Idea:

- ① Extract the leading bit;
- ② Compute the length of the 0-run,  $\ell$ , and extract  $\ell+1$  next bits.
- ③ Convert this into a number, and write this number of bits.
- ④ Repeat steps ②-③ until the string is empty.

eg  $C = \underbrace{000011}_{1} \underbrace{01001}_{2} \underbrace{0010}_{2}$   
 leading bit = 0  
 3 zeroes  
 $\ell = 13$  2 zeroes  
 $\ell = 4$  0 zeroes (so 1)  
 $\Rightarrow S = \underbrace{0000000000000}_{13} \underbrace{1111}_{4} \underbrace{011}_{2}$

**B2:** Pseudocode:

Algorithm 10.7: RLE::decoding( $C, S$ )

```

Input : Non-empty input-stream  $C$  of bits, output-stream  $S$ 
// pre:  $C$  is a valid run-length encoding
1  $b \leftarrow C.pop()$  // bit-value for the current run
2 while  $C$  is not empty do
3    $\ell \leftarrow 0$  // length of binary encoding, minus 1
4   while  $C.pop() = 0$  do  $\ell \leftarrow \ell + 1$  // The last pop() removed the leading bit of the binary encoding
5    $k \leftarrow 1$ 
6   for  $j \leftarrow 1$  to  $\ell$  do // get binary encoding and convert
7      $k \leftarrow k * 2 + C.pop()$ 
8   for  $j \leftarrow 1$  to  $k$  do  $S.append(b)$  // append the run
9    $b \leftarrow 1 - b$  // flip bit for next run

```

**B3:** Run-time:  $O(|S| + |C|)$ .

**B4:** We also note

① This works well for long runs.

eg  $0^n \rightarrow \underbrace{\lfloor \log n \rfloor + 1}_{\text{binary representation}} + \underbrace{\lfloor \log n \rfloor + 1}_{\text{the leading 0's}} + 1 = 2\lfloor \log n \rfloor + 2$ .

② But this works really badly for short runs.

eg A run of length 2.

....11.... → ....010.....

③ It is useful for transmitting black & white pictures (especially text).

Why? → long runs of white/black pixels.

# LEMPEL-ZIV-WELCH ENCODING

💡 This is used in compress & GIF.

💡 Idea: We instead automatically detect longer substrings that got codewords.

💡 To do this, we repeatedly

① Find the longest substring, in the part we want to encode, for which we have a codeword.

② Add this string + next char to our dictionary of strings with codewords.

eg  
D: 

Idea:  
D: dictionary that maps strings to codewords (usually a tree)

We keep adding to the trie in this way until we run out of input.

💡 We assume

- ① the text is in ASCII; &
- ② D initially stores ASCII.

💡 For now, the output is a list of integers.

We need to convert this into bitstrings.

- ① We could use Elias-Gamma codes - but this gets too long.
- ② Instead, we use "12-bit encoding" of the numbers.  
- but we stop adding codewords after code 4095.

000001000001 000001001110 000010000000 000001000001 000001010011 000010000000 000010000001  
65 78 128 65 83 128 129

💡 Pseudocode:

Algorithm 10.8: LZW::encoding( $S, C$ )

```

Input : Input-stream  $S$  of ASCII-characters, output-stream  $C$ 
1 Initialize dictionary  $D$  as a trie that maps ASCII to  $\{0, \dots, 127\}$ 
2  $idx \leftarrow 128$  // global counter for first free code-number
3 while  $S$  is not empty do
4    $v \leftarrow$  root of trie  $D$ 
5   while  $S$  is non-empty and  $v$  has a child  $c$  labeled with  $S.top()$  do
6      $v \leftarrow c$ 
7      $S.pop()$ 
8      $C.append(v \text{ code-number stored at } v)$ 
9   if  $S$  is non-empty then
10    create child of  $v$  labeled  $S.top()$  with code-number  $idx$ 
11     $idx++$ 
```

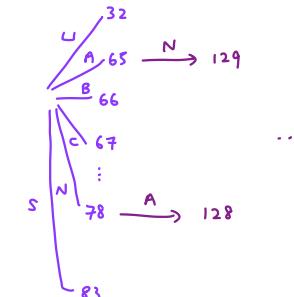
💡 Encoding-time: O(1SI)

## DECODING

💡 Idea:

- ① Initialize  $D$  with ASCII;
- ② Read the next code number;
- ③ Look up corresponding string in  $D$ ;
- ④ Expand  $D$  as the encoder would have done;  
(but we are one step behind)
- ⑤ Repeat ②-④.

eg 65 78 128 65 83 128 129



input	decodes to	Code #	String
67	C		
65	A	128	CA
78	N	129	AN
32	□	130	Ν□
66	B	131	□B
129	AN	132	BA
133	???	133	

💡 Note that at the end, we could get a code number that was about to be added.

💡 If we want to use the decoding of the codeword we are about to add:  
we can show the string is then the previous string + 1st char of the prev. string.

💡 Pseudocode:

Algorithm 10.9: LZW::decoding( $C, S$ )

```

Input : Input-stream  $C$  of integers, output-stream  $S$ 
1 Initialize  $D$  as a dictionary that maps  $\{0, \dots, 127\}$  to ASCII
2  $idx \leftarrow 128$  // global counter for first free code-number
3  $code \leftarrow C.pop()$  // first number creates no entry in  $D$ 
4  $s \leftarrow LZW::dictionaryLookup(D, code)$ 
5  $S.append(s)$ 
6 while  $C$  is not empty do
7    $s_{prev} \leftarrow s$  // previous string
8    $code \leftarrow C.pop()$ 
9   if  $code < idx$  then  $s \leftarrow LZW::dictionaryLookup(D, code)$ 
10  else if  $code = idx$  then  $s \leftarrow s_{prev} + s[0]$  // concatenation
11  else return "invalid encoding"
12   $S.append(s)$ 
13  insert  $s_{prev} + s[0]$  into  $D$  with code-word  $idx$ 
14   $idx++$ 
```

💡 Decoding time:

- each round of the while takes time proportional to the # of chars written to the output

Thus the run-time is O(1SI).

💡 Notes:

- ① This compresses well in practice; but
- ② This is very bad if no repeated substrings.

## bzip2

**Idea:** We transform the text into something that is not necessarily shorter, but has other desirable qualities.

Transformation name	Properties	Alphabet
Text $T_0$ : barbarabarbarbaren\$		ASCII
↓ Burrows-Wheeler transform	If $T_0$ has repeated longer substrings, then $T_1$ has long runs of characters.	
Text $T_1$ : nrbbbbb\$arreaaaa		ASCII
↓ Move-to-front transform	If $T_1$ has long runs of characters, then $T_2$ has long runs of zeros.	
Text $T_2$ : 110,114,100,0,0,0,1,6,100,2,0,0,103,2,0,0,0		{0,...,127}
↓ Modified RLE	If $T_2$ has long runs of zeroes, then $T_3$ has chars $A'$ and $B'$ very frequently	
Text $T_3$ : 110,114,100,A',B',1,6,100,2,B',103,2,A',B'		{1,...,127} $\cup \{A', B'\}$
↓ Huffman encoding	Compresses well since input-chars are unevenly distributed	
Text $T_4$ : 000111011110001000000111110101110010110001		{0,1}

## MODIFIED RLE

**Idea:** Encode only runs of 0, using binary bijection numeration.

$$(\Sigma = \{0, \dots, 127\})$$

$$\text{eg } S = 110, 114, 100, \underbrace{0, 0, 0}_{3}, 1, 6, 100, 2, \underbrace{0, 0}_{2}, 103$$

$$\Rightarrow C = 110, 114, 100, A', B', 1, 6, 100, 2, B', 103$$

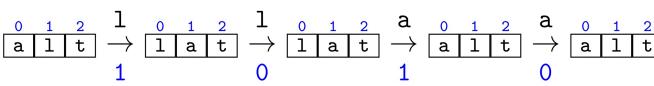
$$\begin{array}{ccccccccc} k & 0 & 1 & 2 & 3 & 4 & \dots & \leftarrow \text{maps runs of 0's of} \\ \text{Eg)} & A' & B' & A', A' & A', B' & \dots & \leftarrow \text{this length} & \\ & & & & & & \leftarrow \text{to this substring.} & \end{array}$$

$$\text{Output alphabet: } (\Sigma = \{1, \dots, 127\} \cup \{A', B'\})$$

## MOVE-TO-FRONT TRANSFORM

**Idea:**

- ① Initialize  $D$ : array of size  $|\Sigma_S|$  that stores  $\Sigma_S$  (typically ASCII)
- ② Get char  $c$  from the input.
- ③ Write  $D^{-1}(c)$  to the output.
- ④ Update  $D$  by bringing  $c$  to the "front" of  $D$  (the MTF-heuristic).



**I<sub>2</sub>:** If we have  $k$ -consecutive same characters, we get a run of  $k-1$  0's in the output.

**I<sub>3</sub>:** We should have lots of 0, 1, 2, ..., & very few of 125, 126, ...

**I<sub>4</sub>:** Decoding: Same except  $D^{-1}$  becomes  $D$ .

## BURROWS - WHEELER TRANSFORM

**Idea:**

$$\text{eg } S = \text{alf\_eats\_alfalfa\$}$$

→  $\begin{array}{l} \text{alf\_eats\_alfalfa\$} \\ \text{lf\_eats\_alfalfa\$a} \\ \text{f\_eats\_alfalfa\$al} \\ \text{eats\_alfalfa\$alf} \\ \text{ats\_alfalfa\$alf\_e} \\ \text{ts\_alfalfa\$alf\_ea} \\ \text{s\_alfalfa\$alf\_eat} \\ \text{alfalfa\$alf\_eats} \\ \text{alfalfa\$alf\_eats\_a} \\ \text{fa\$alf\_eats\_alfal} \\ \text{fa\$alf\_eats\_alfal} \\ \text{a\$alf\_eats\_alfal} \\ \text{a\$alf\_eats\_alfal} \end{array}$

C ↗

① Write all cyclic shifts; ie  $S[i:i+n-1] + S[0:i]$  Vi.

② Sort them lexicographically.

( $\$ < \text{a} < \text{c} < \text{e} < \dots$ )

③ C = rightmost column of the result "matrix".

**I<sub>2</sub>:** Observe that if  $S$  has repeated substrings, then C likely has long runs of chars.

eg in our example, "alf" shows up 3 times.

⇒ there exists 3 cyclic shifts that start with "alf".

⇒ there exist 3 cyclic shifts that start with "lf" and end at "a".

Likely, these shifts will end up consecutive.

If they are consecutive, this implies the 3 'a's are consecutive in C.

For the same reasoning, we also get 3 'l's that are likely to be consecutive in C.

## FAST BURROW-WHEELER TRANSFORM

**Idea:**

start-index of cyclic shift		$\pi(k)$	As the corresponding suffixes in the suffix array.
0	alf_eats_alfalfa\$	0	\$alf_eats_alfalfa
1	lf_eats_alfalfa\$a	1	alfalfa\$alf_eats
2	f_eats_alfalfa\$al	2	eats_alfalfa\$alf
3	eats_alfalfa\$alf	3	a\$alf_eats_alfalfa
4	ats_alfalfa\$alf_e	4	alf_eats_alfalfa\$
5	ts_alfalfa\$alf_ea	5	alfalfa\$alf_eats_alf
6	s_alfalfa\$alf_eat	6	alfalfa\$alf_eats_alf
7	alfalfa\$alf_eat	7	ats_alfalfa\$alf_eats
8	alfalfa\$alf_eats	8	eats_alfalfa\$alf
9	alfalfa\$alf_eats_a	9	f_alf_eats_alfalfa\$al
10	alfalfa\$alf_eats_a	10	fa\$alf_eats_alfalfa
11	fa\$alf_eats_alfal	11	falfa\$alf_eats_alfal
12	alfa\$alf_eats_alf	12	1f_eats_alfalfa\$al
13	ifa\$alf_eats_alfa	13	lf_eats_alfalfa\$al
14	fa\$alf_eats_alfal	14	lfalfa\$alf_eats_alf
15	as\$alf_eats_alfalf	15	s_alfalfa\$alf_eat
16	\$alf_eats_alfalfa	16	ts_alfalfa\$alf_eat

① We need the sorting permutation of the cyclic shifts.

② This is the same as the sorting permutation of the suffixes.

③ So, to compute the encoding, we

- compute the suffix array ( $O(n \log n)$ )

- output  $C[i] = S[(A_s[i]-1) \bmod n]$ .

## DRAWBACKS OF BWT

**I<sub>1</sub>:** Run-time:  $O(n \log n)$  (maybe  $O(n)$ )

**I<sub>2</sub>:** Space:  $O(n)$  (with suffix array S).

**I<sub>3</sub>:** Here, we need the entire text at once.

# DECODING IN BWT

Example 1:

eg  $C = annb\$aa$

Reconstruct the matrix of cyclic shifts:

$a_0$	$a_0$
$a_1$	$a_1$
$a_2$	$a_2$
$a_3$	$\$ b_3$
$b_4$	$y$
$n_5$	$a_5$
$n_6$	$a_6$

indexes by  
row #

- rightmost column is just  $C$ .
- leftmost column is the same as  $C$ , but sorted.
- First character of  $S$  is  $b$ .
- Second character of  $S$  is  $a$ .

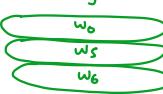
The char after  $b$  is the  $a$  in row 6.

Proof. Define  $w_i$  = the word in front of char  $x$  in row  $i$ .

$w_0$	$a_0$	$o$
$w_1$	$a_1$	$1$
$w_2$	$a_2$	$2$
$\vdots$		
$w_5$	$a_5$	$5$
$w_6$	$a_6$	$6$

We note  $w_0 \leq w_5 \leq w_6$   
(no equality  $\therefore$  end of char symbol  
is in different positions).

We have 3 cyclic shifts



The  $a$ 's do not change the lexicographic order,  
so

Restating:  $a_0 w_0 \leq a_5 w_5 \leq a_6 w_6$ .

The characters in the first column are in the same relative order as they were in the last column.

$a_0$	$a_0$	$o$
$a_5$	$a_5$	$1$
$a_6$	$a_6$	$2$
$b_7$	$\$ b_3$	$3$
$n_8$	$a_8$	$4$
$n_9$	$a_9$	$5$
$b_{10}$	$a_{10}$	$6$
$c_{11}$	$a_{11}$	$7$
$d_{12}$	$a_{12}$	$8$
$r_{13}$	$b_{13}$	$9$
$r_{14}$	$b_{14}$	$10$

indexes by  
row #

$S: b_a_n_a_n_a_o \$$

Final string:

abacadabra\$

(see course notes for  
more details.)

$g_3 \rightarrow a_0 \rightarrow r_1 \rightarrow b_{10} \rightarrow a_6 \rightarrow \dots$

(append backwards.)

Example 2:

eg  $C = ard\$recaaaaabb$

$\begin{array}{l} g_3 \\ a_0 \\ a_5 \\ a_6 \\ a_7 \\ a_8 \\ a_9 \\ b_{10} \\ b_{11} \\ c_5 \\ d_2 \\ r_1 \\ r_4 \end{array} \quad \begin{array}{l} \cdots \cdots \cdots \\ \cdots \cdots \cdots \end{array} \quad \begin{array}{l} a_0 \\ r_1 \\ d_2 \\ g_3 \\ a_7 \\ a_8 \\ r_4 \\ a_9 \\ b_{10} \\ a_5 \\ a_6 \\ a_3 \\ a_2 \\ a_1 \\ a_0 \end{array}$

Pseudocode:

Algorithm 10.14:  $BWT::decoding(C, S)$

Input : Coded text  $C$  as array (not stream), output-stream  $S$

```

1  $A \leftarrow$  array of size  $n \leftarrow C.size$ 
2 for  $i = 0$  to  $n - 1$  do  $A[i] \leftarrow (C[i], i)$ 
3 Stably sort  $A$  by first aspect // eg radix sort, merge sort
4 for  $j = 0$  to  $n - 1$  do if  $C[j] = \$$  then break           // find $-char
5 repeat                                         // actual decoding
6   |  $S.append(first\ entry\ of\ A[j])$ 
7   |  $j \leftarrow second\ entry\ of\ A[j]$ 
8 until last appended character was $
```

This is simple, & runs in  $O(n)$  time.

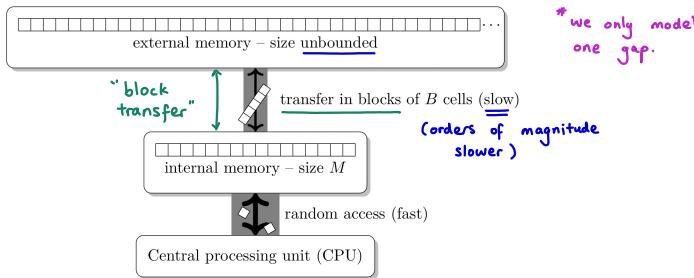
# Chapter 11:

## External Memory

- Q<sub>1</sub> So far, our computer model assumed "all memory cells are equal" and "we have infinitely many of them".
- Q<sub>2</sub> Now, we still assume we have infinite memory, but distinguish how easy it is access.

### THE EXTERNAL-MEMORY MODEL / EMM

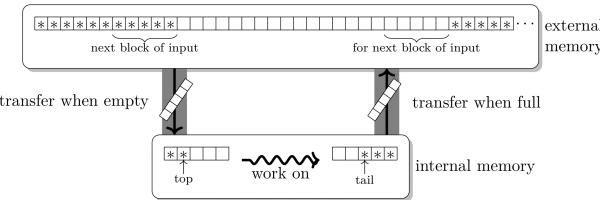
Idea:



- Q<sub>2</sub> In particular,
- $B$  is the "block size"; the amount transferred at once ( $\sim$  few MB)
  - $M$  is the size of the internal memory.  
(can load  $\leq \frac{M}{B}$  blocks at once)  
( $\sim$  few GB)
  - $n$  is the size of the input. ( $\sim$  few TB/PB)

### STREAMS

Idea:



- Q<sub>2</sub> Total # of block transfers:  $\Theta(\frac{n}{B})$   
(this is optimal).

### GOOD/BAD ALGOS FOR EM

Q<sub>1</sub> Good algorithms:

- Huffman
- LZW
- KMP
- Boyer-Moore

} use streams

Q<sub>2</sub> Bad algorithms:

- BWT
- Suffix array & tree

} require whole input

### MERGE

Q<sub>1</sub> Merge takes  $\Theta(\frac{n}{B})$  block transfers.

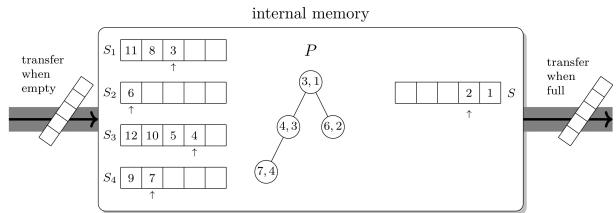
### MERGESORT

Q<sub>1</sub> Mergesort uses  $\lceil \log_2(n) \rceil$  rounds of "splitting" & "merging".

Q<sub>2</sub> So, we need  $\Theta(\frac{n}{B} \cdot \log n)$  block transfers.

### d-WAY MERGE

Idea:



Q<sub>2</sub> Pseudocode:

Algorithm 11.2:  $d$ -way merge( $S_1, \dots, S_d, S$ )

```

Input : Input-streams  $S_1, \dots, S_d$  that are in sorted order. Output-stream  $S$ .
1  $P \leftarrow$  empty priority queue that stores pairs and is min-oriented w.r.t. first aspect
2 for  $i \leftarrow 1$  to  $d$  do  $P.insert(\langle S_i.top(), i \rangle)$  // Populate priority queue
3 while  $P$  is not empty do // repeatedly extract minimum
4    $\langle x, i \rangle \leftarrow P.deleteMin()$  // Item  $x$  is top-item of stream  $S_i$ 
5    $S.append(S_i.pop())$ 
6   if  $S_i$  is not empty then  $P.insert(\langle S_i.top(), i \rangle)$ 

```

Q<sub>3</sub> If we split & merge  $d$  ways, we only need  $\lceil \log_d n \rceil$  rounds.

Q<sub>4</sub> So, # of block transfers =  $\Theta(\log_d n \cdot \frac{n}{B})$ .

Q<sub>5</sub> But we can improve further (details omitted)

$(\Theta(\log_{\frac{n}{B}} \frac{n}{M}) \cdot \frac{n}{B}) \rightarrow$  "optimal" bound for EMM.

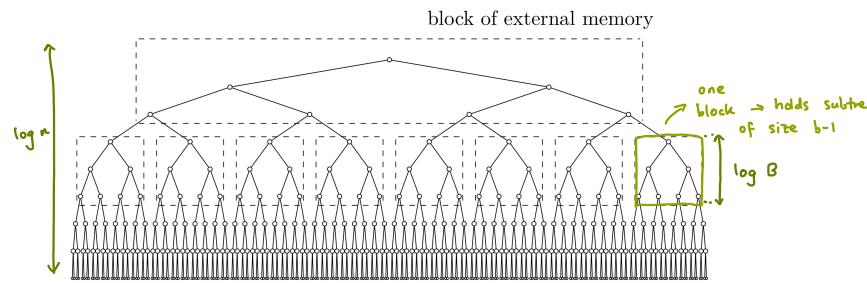
# COMPARISON-BASED DICTIONARIES

## BSTs

- $\Theta_1$  we know search/insert/delete can be done with  $O(\log n)$  block transfers (naively).
- $\Theta_2$  We can show we need at least  $\Omega(\log n)$  block transfers for search.
- $\Theta_3$  Can we achieve  $O(\log n)$  block transfers?

## IDEAL STRUCTURE FOR SEARCH

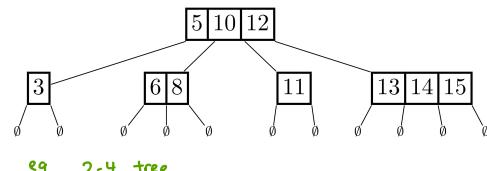
- $\Theta_1$  Idea:



- $\Theta_2$  Search path hits  $\frac{\Theta(\log n)}{\log b} \Rightarrow \Theta(\log_b n)$  block transfers, where  $b = \Theta(B)$ .
- $\Theta_3$  We can also view each block as one node in a multiway tree ( $b-1$  KVPs,  $b$  subtrees)
- $\Theta_4$

## a-b-TREE

- $\Theta_1$  Idea:

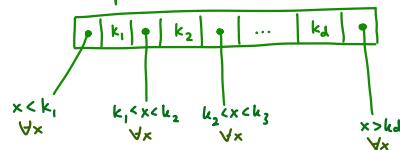


eg 2-4 tree

- $\Theta_1$   $a, b \in \mathbb{Z}^+, 2 \leq a \leq \lfloor \frac{b+1}{2} \rfloor$ . (usually we set  $a = \lfloor \frac{b+1}{2} \rfloor$ )
- $\Theta_2$  Every node has  $a \leq \# \text{subtrees} \leq b$ , except at the root, which has  $2 \leq \# \text{subtrees} \leq b$ .
- $\Theta_3$  All empty subtrees are on the same layer.
- $\Theta_4$  A node with  $d$  KVPs has exactly  $d+1$  subtrees.

- $\Theta_1$  Order-property:

Close up on a node.



## SEARCH IN a-b-TREES

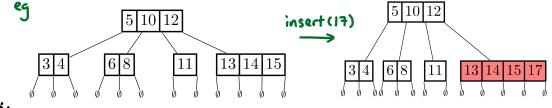
- $\Theta_1$  Similar to BSTs:
  - We load the root;
  - Find the best place for the search key in the node;
  - If  $k$  is not in the node, repeat in the appropriate subtree.
  - Otherwise, return  $k$  & the corresponding value.
- $\Theta_2$  Run-time (in RAM model):  $O(\text{height} \cdot \log b)$ .

## INSERT IN a-b-TREES

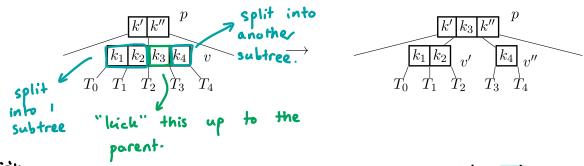
- $\Theta_1$  Idea:

- Call search, and add key and an empty subtree at the leaf.

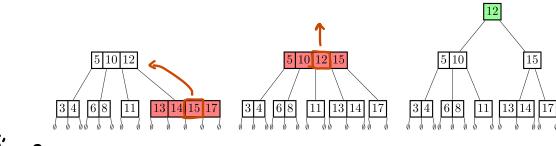
- $\Theta_2$  However, if we do this, we might get that a node is too big.



- $\Theta_3$  To solve this, we use "node splitting":



- $\Theta_4$  We have to repeat node-splitting until all nodes are "happy".



### Pseudocode:

```

Algorithm 11.4: 2dTree::insert(k)
1 v ← 2dTree::search(k) // leaf where k should be
2 Add k and an empty subtree at the appropriate place in the key-subtree-list of v
3 while v has 4 keys do // overflow leads to node split
4   Let  $\langle T_0, k_1, \dots, k_3, T_4 \rangle$  be key-subtree list at v
5   if v has no parent then // tree grows upwards
6     create a parent of v without keys
7   p ← parent of v
8   v' ← new node with keys  $k_1, k_2$  and subtrees  $T_0, T_1, T_2$ 
9   v'' ← new node with key  $k_4$  and subtrees  $T_3, T_4$ 
10 Replace  $\langle v \rangle$  by  $\langle v', k_3, v'' \rangle$  in key-subtree-list of p
11 v ← p
    
```

- $\Theta_5$

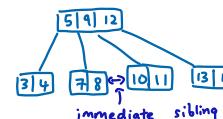
- Note:
  - overflow means  $b$  keys &  $b+1$  subtrees
  - after the node split, new nodes have  $\geq \lfloor \frac{b-1}{2} \rfloor$  keys
  - since we required  $a \leq \lfloor \frac{b+1}{2} \rfloor$ , this is  $\geq a-1$  keys as required.

- $\Theta_6$  Run-time: search ( $O(\text{height} \cdot \log b)$ )

+ update ( $O(\text{height} \cdot \log b)$ )

∴ run-time =  $O(\text{height} \cdot \log b)$ .

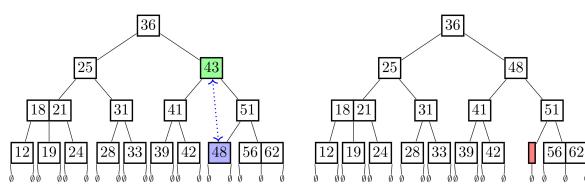
## IMMEDIATE SIBLING



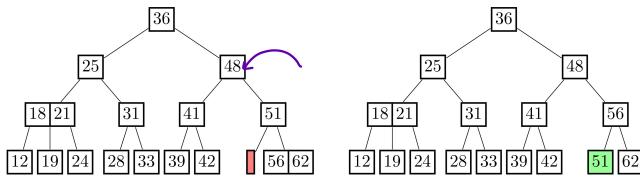
## DELETE IN a-b-TREES

Idea:

- ① Search for  $k$ , then trade with successor if KVP is not at a leaf.
- ② Remove the key & one empty subtree from leaf.
- ③ But this could lead to underflow: ie < a subtrees.



To resolve underflow, we can use "rotations", and "steal" from the immediate sibling if it has "extras".



Otherwise, we can perform a "node merge", by stealing from the parent.



- we repeat this until all nodes are "happy".

If this process continues until the root has only one subtree, it self-destructs.

Pseudocode:

Algorithm 11.5: 2dTree::delete( $k$ )

```

1  $v \leftarrow 2dTree::search(k)$                                 // node containing  $k$ 
2 if  $v$  is not a leaf then
3    $v \leftarrow$  leaf that contains successor  $k'$  of  $k$ , and exchange  $k$  and  $k'$ 
4 delete  $k$  and one empty subtree in  $v$                          // underflow
5 while  $v$  has 0 keys do
6   if  $v$  is the root then delete  $v$ , make its child the new root and break
7   if  $v$  has an immediate sibling  $u$  with 2 or more keys then
8     | rotate a key and subtree from  $u$  to  $v$  as in Figure 11.13 and break
9   else
10    | merge  $v$  with an immediate sibling as in Figure 11.15
11    |  $v \leftarrow$  parent of  $v$  // and repeat the while-loop

```

Run-time:  $O(\text{height} \cdot \log b)$ .

## HEIGHT OF AN a-b-TREE

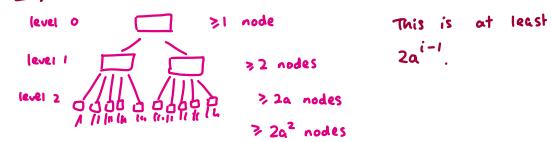
Claim:

$$\text{height} \in O(\log_a n)$$

where  $n = \# \text{ of KVPs}$ .

\* note: # of KVPs >> # of nodes.

Proof: Consider # of nodes on level  $i \geq 1$ .



Thus

$$\begin{aligned} \# \text{ of nodes} &\geq 1 + \sum_{i=1}^h 2a^{i-1} & (\text{if } h = \text{height}) \\ &= 1 + 2 \frac{a^h - 1}{a - 1}. \end{aligned}$$

Then

$$\begin{aligned} n = \# \text{ of KVPs} &\geq \underbrace{1}_{\text{root}} + \underbrace{(a-1) 2 \frac{a^h - 1}{a - 1}}_{\text{everything else}} \\ &= 1 + 2(a^h - 1) \\ &= 2a^h - 1. \end{aligned}$$

$$\Rightarrow h \leq \log_a \left( \frac{n+1}{2} \right) \in O(\log_a n)$$

as needed.  $\square$

With a similar proof, we can show

$$\text{height} \in \Omega(\log_b n).$$

We also choose  $a \in \Theta(b)$ , so the a-b-tree has height

$$\text{height} \in \Theta(\log_b n).$$

## RUN-TIMES OF a-b-TREE OPERATIONS

All the run-times (in the RAM model) were asymptotically  $\text{height} \cdot \log b$ .

Then

$$\text{height} \leq \log_a \left( \frac{n+1}{2} \right) \in O(\log_b n) = O\left(\frac{\log n}{\log b}\right).$$

Thus

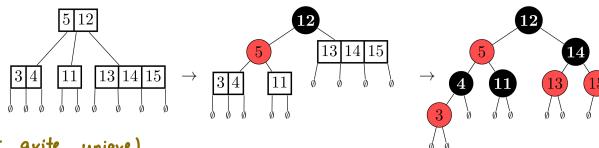
$$\text{run-time} \in O(\log n).$$

This is no faster than AVL trees.

## CONVERT 2-4-TREES TO BSTS ; RED-BLACK TREES

**B<sub>1</sub>** This has  $O(\log n)$  run-time, but faster than AVL-trees in practice.

**B<sub>2</sub>** A d-node becomes a black node with d-1 red children.



(not quite unique)

**B<sub>3</sub>** The resultant trees have height  $O(\log n)$ .

A "red-black tree" has the following properties:

- ① It is a BST;
- ② Every node is either red or black;
- ③ Every red node has a black parent;
- ④ Every empty subtree T has the same "black depth";  
ie the # of black nodes from the root to T.

### Insert / delete:

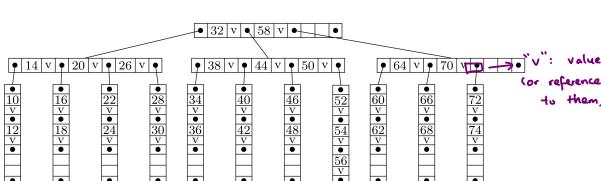
- ① Convert the red-black tree into a 2-4 tree;
- ② Perform the operation on the 2-4-tree; then
- ③ Convert the 2-4-tree back into a red-black tree.

## B-TREES

**B<sub>1</sub>** A "B-tree" is an a-b-tree tailored to Emm.

### Idea:

- ① Every node is one block of memory (of size B)
- ② b is chosen maximally s.t. a node with b-1 KVPs fit into a block; & (b is called the "order" of the tree.) \* be  $\Theta(B)$  typically.
- ③  $a = \lceil \frac{b}{2} \rceil$ .



**B<sub>3</sub>** Note that

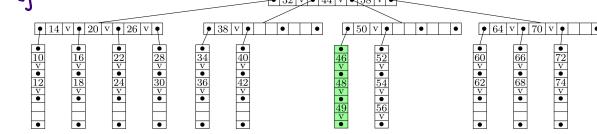
$$\begin{aligned} \# \text{ block transfers in B-tree} &= O(\text{height}) \\ &= \Theta(\log_b n) \\ &= \Theta(\log_B n) . \end{aligned}$$

**B<sub>4</sub>** This is asymptotically optimal.

## B-TREE IMPROVEMENT 1 — PRE-EMPTIVE SPLITTING / MERGING

### Motivation:

e.g. insert(49)



**B<sub>1</sub>** Standard insert does the following:

- load the blocks "on the way down"
- load them again "on the way up".

**B<sub>2</sub>** We want to avoid this.

### Idea:

- ① If node is not full, keep searching
- ② If node is full, immediately split.
- ③ Then keep searching in appropriate new node.

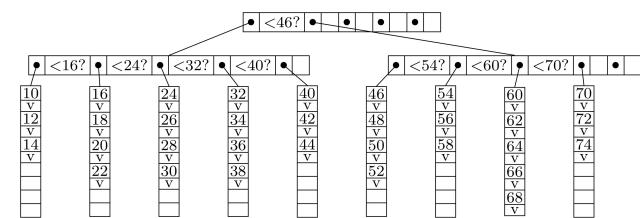
**B<sub>3</sub>** ie we split "just in case".

**B<sub>4</sub>** This can save half of the block transfers.

## B-TREE IMPROVEMENT 2 — B+ TREES

### Idea:

- ① Each node is one block of memory.
- \* ② All KVPs are stored at leaves, each of which are at least half full.
- ③ Interior nodes store only keys for comparison during search.
- ④ Interior non-root nodes have at least half of the possible subtrees.
- ⑤ Insert / delete use pre-emptive merging / splitting.



→ This can store  $5^3 = 125$  KVPs.

**B<sub>2</sub>** Advantages:

- ① Bigger order;
- ② Can store more KVPs for the same height.

## B-TREE IMPROVEMENT 3 — LSM-TREES

### Idea:

- ① Store dictionary in internal memory that logs all changes;
- ② To search: first search in C<sub>0</sub>, then (if needed) in C<sub>1</sub>
- ③ If internal memory full; do lots of updates in C<sub>1</sub> at once.

$$(32,v) \circ (56,v) \circ (-,-)$$

C<sub>0</sub> (log of changes)

C<sub>1</sub> (B-tree)



**B<sub>2</sub>** Note: many further improvements.

# HASHING IN EMM

💡 Each operation takes expected amortized  $\Theta(1)$  runtime.

💡 This means each operation also takes  $O(1)$  block-transfers.

💡 However, amortized bounds are a problem!  
→ rehashing is the issue.

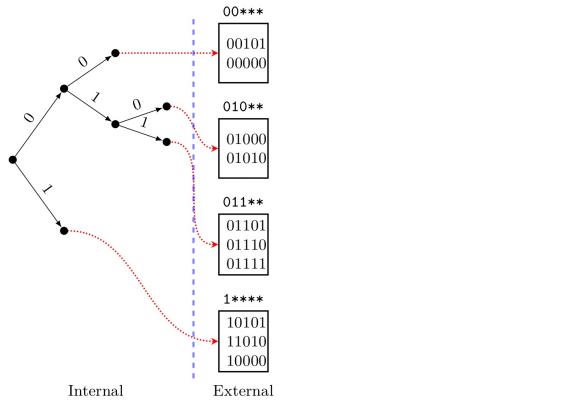
💡 So, how do we do hashing without re-hashing?

## TRIE OF BLOCKS

💡 Assumption: we store non-negative integers (bitstrings)

💡 Idea:

- ① Build pruned trie  $D$  (the "directory") of integers in internal memory.
- ② Stop splitting in trie when remaining items fit in one block.
- ③ Each leaf of  $D$  refers to block of external memory that stores the items.



💡 Search( $k$ ):

- ① get bitstring for  $k$ , say 00101
- ② Search in the trie
- ③ load block at leaf
- ④ search in it

→ This is 1 block transfer

(assuming trie fits into internal memory)

💡 Insert:

- ① Find the block that needs it, & insert
- ② If the block is full, split the block by the next bit.
- ③ Repeat ② until the blocks have space.

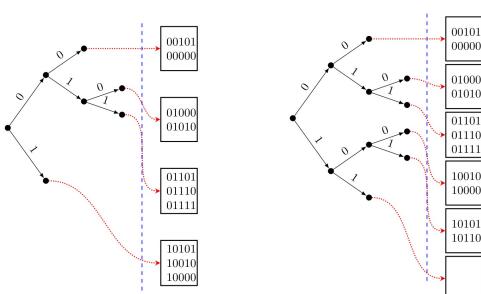


Figure 11.29: Inserting 10110 into a trie of blocks.

→ It is very likely we are done after 2-3 block transfers.

💡 Delete: similar.

💡 If no split works (ie duplicate bitstrings), we extend the hash function.

- ① Suppose  $h(k)$  yields the bitstring corresponding to  $k$ .
- ② Let  $h_i: k \rightarrow \text{bitstring}$  be another hash-function.
- ③ Then, use

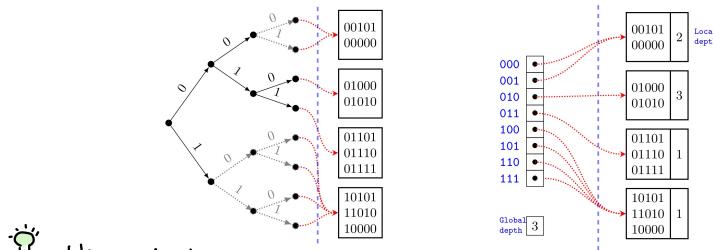
$$h + h_i : k \rightarrow \text{bitstring} = h(k) + h_i(k)$$

concatenation

## EXTENDIBLE HASHING IN TRIE OF BLOCKS

💡 We can also save links by "expanding" the trie and only storing the leaves.

💡 We store the  $D$  as an array, rather than a trie:

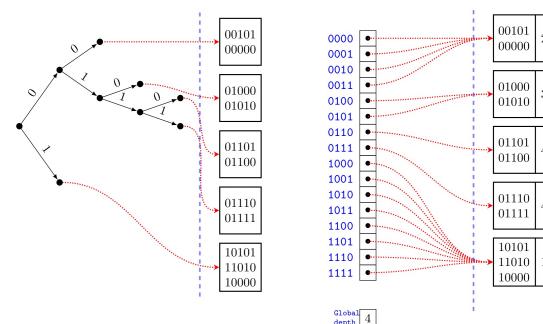


💡 We need to store

- ① the global depth  $d_D$ ; &
- ② the local depth.

💡 Insert: convert trie into table & v.v.

💡 If the trie's height increases, we need to update the array.



- but we don't need to load blocks to update the array.

💡 This uses no more block transfers than with tries.