

CS 240

Input for sorting algs

- list of ints
- no particular order.

e.g. 2, 5, 7, 23, 14, 8

what size is this instance?

$$n = 5 \text{ ints} \quad \text{or} \quad n = 5 \times 4 = 20 \text{ bytes}$$

e.g.: Matrix Mult

$$\begin{bmatrix} C_{11} & \dots & C_{1n} \\ \vdots & \ddots & \vdots \\ C_{n1} & \dots & C_{nn} \end{bmatrix} = \begin{bmatrix} A_{11} & \dots & A_{1n} \\ \vdots & \ddots & \vdots \\ A_{n1} & \dots & A_{nn} \end{bmatrix} \begin{bmatrix} B_{11} & \dots & B_{1n} \\ \vdots & \ddots & \vdots \\ B_{n1} & \dots & B_{nn} \end{bmatrix}$$

C A B

Size of input?

o # of ints in both matrices

o # of bytes

o dimension of matrices

e.g. matrices that are 10x10

$$\text{Dimension: } n = 10$$

200 ints, 800 bytes.
 n^2

$$C_{11} = a_{11} \cdot b_{11} + \dots + a_{1n} \cdot b_{1n} \quad \text{n mults}$$

$\frac{n-1 \text{ adds}}{2n-1 \text{ ops}}$



扫描全能王 创建

n^2 positions in C.

$$\text{Total } n^2(2n-1) = 3n^3 - n^2 \in O(n^3)$$

2017.9.12

Big-O

- upper bound on growth rate

- "grows no faster than"

Relation Basic idea functions = E

$f \leq g$

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

$f \geq g$

Ω

$f = g$

Θ

$f < g$

$\in O$

$f > g$

ω

e.g. $O \leq \frac{2n^2 + 3n + 1}{f(n)} \leq C \cdot \frac{n^3}{g(n)}$ $\forall n \geq n_0$

Find an n_0 and C .

* You must ~~choose~~ C and derive

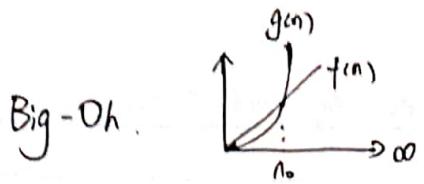
* Do not simply guess $C = 1/n^2$, $g(n)$

Start $\frac{2n^2 + 3n + 1}{f(n)} \leq \frac{2n^2 + 3n^2 + 1}{n^2} \leq n^2$ $\forall n \geq 1$

$\therefore (n \geq 1) \quad \text{justify}$



扫描全能王 创建



Typically we think $g(n)$ must dominate $f(n)$.

Eg. Prove $\frac{(n+1)^5}{f(n)} \in O(\frac{n^5}{g(n)})$

: $g(n)$ will never dominate $f(n)$
 $\Rightarrow C > 1$

Note, not all values of c will work.

$$(n+1) \leq 2 \cdot n \quad \forall n \geq 1$$

$$(n+1)^5 \leq (2 \cdot n)^5 = 2^5 \cdot n^5 = 32n^5 \quad \forall n \geq 1$$

eg $\underbrace{n^2 + n \log n + n}_{f(n)} \in O(\underbrace{n^2}_{g(n)})$

$$n^2 + n \log n + n \leq n^2 + n^2 + n^2 \quad \forall n \geq 1$$

since $n^2 > n \log n \quad \forall n \geq 1$
 $n^2 \geq n \quad \forall n > 0$

$$n^2 + n^2 + n^2 = \underbrace{3n^2}_{c g(n)}$$

Don't get lazy in proofs!

- take small steps. justify your statements.



扫描全能王 创建

From "First Principles" we explicitly find c, n_0 .

Properties: Assume $f(n) \geq 0, g(n) \geq 0, h(n) \geq 0$

1. $a f(n) \in O(f(n))$ for any constant $a > 0$

$$\text{pf: } c = a, n_0 = 1.$$

2. If $f(n) \in O(g(n))$ and $g(n) \in O(h(n))$

then $f(n) \in O(h(n))$

From

$$\begin{aligned} \textcircled{1} \quad f(n) &\leq C_1 \cdot g(n) \quad \forall n \geq n_1, \\ \textcircled{2} \quad g(n) &\leq C_2 \cdot h(n) \quad \forall n \geq n_2 \end{aligned}$$

Substitute \textcircled{2} into \textcircled{1}

Show \textcircled{2} $f(n) \leq C \cdot h(n) \quad \forall n \geq n_0$

$$C = C_1 \cdot C_2$$

$$n_0 = \max(n_1, n_2)$$

3. $\max(f(n), g(n)) \in O(\max(f(n), g(n)))$

$$\text{pf: } C = 1, n_0 = 1$$

Exercise: Show $f(n) + g(n) \in O(\max(f(n), g(n)))$

4. $a_0 + a_1 x + a_2 x^2 + \dots + a_n x^n \in O(x^n)$

5. $n^x \in O(a^n), x > 0, a > 1$

e.g. $n^{100000000} \in O(1.000001^n)$

6. $(\log n)^x \in O(n^y) \quad x, y > 0$



扫描全能王 创建

Draft:

Show $f(n) + g(n) \in O(\max(f(n), g(n)))$.

3. $f(n) \leq f(n) + g(n)$

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

$$\max(f(n), g(n)) \in O(f(n) + g(n))$$

$$f(n) + g(n) \leq 2 \max(f(n), g(n)).$$

Thus $f(n) + g(n) \in O(\max(f(n), g(n))) \quad \forall c > 1, n_0 = 1$

al. 1.

a) $n > \log n \quad \forall n > 0$

$$206n > 206 \quad \forall n > 1$$

$$27n^7 + 17n^3 \log n + 206 \leq 27n^7 + 17n^7 + 206n^7 = 2060n^7$$

$$= \frac{2060}{n^2} n^9$$

$$\leq cn^9.$$

Since $c \geq \frac{2060}{n^2}$
 $n^9 \geq \sqrt[7]{\frac{2060}{c}}$

For any c we choose $n \geq \sqrt[7]{\frac{2060}{c}}$, thus $n^9 \leq cn^9$, which means ...

b. Since $(\log n)^{1000} \geq 1 \quad \forall n > 1$

$$n^2 (\log n)^{1000} \geq n^2 \quad \text{when } c=1, \forall n \geq 2$$

c.



扫描全能王 创建

$$a \frac{n^2}{n+\log n} = \frac{n}{\frac{1+\log n}{n}} / (n \geq 1) \leq \log n \leq n \leq \frac{\log n}{n} \leq 1.$$

$$1 \leq \frac{\log n}{n} + 1 \leq 2$$

$$\frac{n}{2} < \frac{n}{\frac{\log n}{n} + 1} \leq n$$

$$d. n^{20} \leq 20n^{20}$$

$$\log n \quad n \leq 20.$$

$$\leq \frac{20}{n^{20-n}} n^n$$

$$n^n \geq \frac{1}{n^{n-20}} n^{20}.$$

$$> cn^{20} \geq \frac{1}{n^{n-20}} n^{20}.$$

$$c \geq \frac{1}{n^{n-20}} n^{n-20}$$

$$\log c \geq n-20$$

$$n^n \geq cn^{20}$$

$$\log n^n \geq \log cn^{20}.$$

$$n \geq \log c + 20.$$

$$n \log n \geq \log + 20 \log n$$

h.

$$(n-20) \log n \geq \log c.$$

$$n^{20} \leq n \cdot n^{20} \quad n \geq 1$$

$$20 \log n \leq \log n \leq c \log$$

$$\log n + \frac{n-20}{n}$$

$$\log n + 1 - \frac{20}{n}$$



扫描全能王 创建

eg. $(\log n)^{100000000} \in O(n^{0.000000001})$

• Big-Omega "grows no slower than"
lower bound.

eg. $n^3 \log n \in \Omega(n^3)$

$$\frac{1 \cdot n^3}{c \cdot g(n)} \leq n^3 / \log n \quad \forall n \geq 2$$

Since $\log n \geq 1 \quad \forall n \geq 2$

Θ -Big-Theta $f(n)$ grows at same rate as $g(n)$.

Show $f(n) \in \Theta(g(n)) \Rightarrow$ show $f(n) \in O(g(n))$
and $f(n) \in \Omega(g(n))$

$$0 \leq c_1 \cdot n^3 \leq 2n^3 - n^2 \leq 2n^3 \quad \forall n \geq n_0$$

O -Little-O $f(n)$ grows slower than $g(n)$

key difference "Hc"

eg $n \in o(n^2)$

Given any c , we can always find n_0 , $c = \frac{1}{10}, \frac{1}{100}, \frac{1}{100000}$

ω Little-omega $f(n)$ grows faster than $g(n)$

eg $n^2 \in \omega(n)$ no matter how large c we choose, $\exists n_0$



扫描全能王 创建

$$2010n^2 + 1388n \in \Theta(n^3)$$

$$\forall c > 0 \exists n_0 > 0 \text{ s.t. } 2010n^2 + 1388n < cn^3$$

Proof: Let $c > 0$ be given

Note $0 < c < 1$ is possible

$$2010n^2 + 1388n < 5000n^3 \quad \forall n \geq 1$$
$$\leq \frac{5000}{n} n^3 \quad \forall n \geq 1$$

$$\leq c \cdot n^3$$

$$\forall n \geq n_0 \text{ Note } \frac{5000}{n} \leq c$$

$$\forall n \geq \frac{5000}{c}$$

* ~~the~~ no ~~m~~ may depend on c .
can go as far as required.

For any c we choose $n_0 = \frac{5000}{c}$

2017.9.14

Insertion sort:

Instance: Given an ~~an~~ array of length n .

Best case: array ~~is~~ already sorted $\Rightarrow \Theta(n)$

Worst case: array reverse sorted $\Rightarrow \Theta(n^2)$

Partially sorted - included in average case

Average case:

- must consider all inputs of a certain size n
- all possible permutations of array.



扫描全能王 创建

- Sometimes algs have bad worst case but perform well in practice
- we often focus on worst case.

- insertion sort
 - worst case $\Theta(n^2)$
 - average case $\Theta(n^2)$

$$f(n) = \begin{cases} 3^n + \left\lfloor \frac{n^2}{2} \right\rfloor & \text{for } n \leq 10 \\ 15 \cdot 3n^3 + 2n^2 + 12 & \text{for } n \geq 11 \end{cases}$$

$$\begin{aligned} f(n) &\in O(n^3) & f(n) &\in O(n^3) \\ f(n) &\in O(3^n) \end{aligned}$$

n	$n \log n$	n^2	n^3	3^n
4	8	16	64	16
16	64	256	4096	95,376
64	384	4096	262,144	1.84×10^{19}

Matrix Mult

A₁: Standard: $2n^3 - n^2 \in \Theta(n^3)$

A₂: Strassen: $42n^{2.8074...} \in \Theta(n^{\log 2})$

$$n=10 \quad A_1 = 1900$$

$$A_2 = 24955$$

Let $i > 0$

$$\lim_{n \rightarrow \infty} \frac{\log n}{n^i} = \lim_{n \rightarrow \infty} \frac{\frac{1}{n}}{\ln n^{i-1}} = \lim_{n \rightarrow \infty} \frac{1}{i n^{i-1}} = 0$$

So $\log(n) \in o(n^i)$ for any $i > 0$.

$$n \leq n(2 + \sin \frac{n\pi}{2}) \leq 3n$$



扫描全能王 创建

~~Arithmetic~~ Arithmetic Sequence - a, d, r , are constants.

$$\text{eg. } \sum_{i=1}^n i = \underbrace{1+2+\cdots+n}_{\text{arithmetic}} = (n+1) \left(\frac{n}{2}\right) = \frac{n^2+n}{2}$$

$$\text{eg. } S_1 = \sum_{i=1}^{\infty} \frac{1}{2^i} = 1 + \frac{1}{2} + \frac{1}{4} + \dots$$

Geometric series: $a=1, r=\frac{1}{2}, 0 < r < 1$ case 3.

$$\Rightarrow a \frac{1-r^n}{1-r} = \frac{1-(\frac{1}{2})^n}{1-\frac{1}{2}} = 2(1-(\frac{1}{2})^n)$$

$$\text{"closed form"} \lim_{n \rightarrow \infty} \sum_{i=0}^{n-1} \left(\frac{1}{2}\right)^i = \lim_{n \rightarrow \infty} 2(1-(\frac{1}{2})^n) = 2$$

$$S_2 = \sum_{i=0}^{\infty} \frac{1}{2^i} = \frac{0}{1} + \frac{1}{2} + \frac{2}{2^2} + \frac{3}{2^3} + \dots$$

How to solve this?

$$\begin{aligned} 2 \cdot S_2 - S_2 &= 1 + \frac{2}{2} + \frac{3}{2^2} + \frac{4}{2^3} + \dots \\ &\quad - \frac{1}{2} - \frac{2}{2^2} - \frac{3}{2^3} - \dots \\ &= 1 + \frac{1}{2} + \frac{1}{4} + \frac{1}{8} + \dots \end{aligned}$$

$$\Rightarrow 2 \cdot S_2 - S_2 = S_1$$

$$\Rightarrow S_2 = 2$$

$$H_n = \sum_{i=1}^n \frac{1}{i}$$

harmonic series

diverges very slowly

$$\sum_{i=1}^{1000} \frac{1}{i} \approx 230.84$$

$$\lg 2^n + n^2 2^{3/g} n = \cancel{A \lg 2} + \cancel{\frac{1}{2} g n^3} n \lg 2 + n^2 2^{3/g} n^3$$

$$= n \cdot 1 + n^2 \cdot n^3 \lg 2$$

$$= n + n^5$$



扫描全能王 创建

- each iteration of inner loop $\Theta(1)$
- How many times do the loops iterate
- work from innermost to outermost.

Direct method

$$\# \text{ iterations} = \sum_{i=1}^n (n-i+1) = \sum_{i=1}^n i = (n+1)\left(\frac{n}{2}\right) = \frac{n^2+n}{2} \in \Theta(n^2)$$

\Rightarrow complexity: $\Theta(n^2)$

Sloppy Method to get O -bound.

$$\# \text{ iter} = \sum_{i=1}^n (n-i+1) \leq \sum_{i=1}^n n = n^2$$

$\Rightarrow O(n^2)$

Caveat: still need to prove $\omega(n^2)$

$$\sum_{i=1}^n (n-i+1) \geq \sum_{i=\frac{n}{2}}^{\frac{n}{2}} \frac{n}{2} = \frac{n}{2} \cdot \frac{n}{2} = \frac{n^2}{4}$$

start half way

use smallest value.

Page 42 each iteration has cost $\Theta(1)$

Determine # iterations

Direct Method:

$$\# \text{ iter} \sum_{i=1}^n \sum_{j=i}^n (j-i+1)$$

separate out part of sum that depends on the index.

$$= \sum_{i=1}^n \left[\sum_{j=i}^n j - \sum_{j=i}^n (i-1) \right] \boxed{(n-i+1)(i-1)}$$

$$= \sum_{i=1}^n \left[\sum_{j=1}^{i-1} j - \sum_{j=1}^{i-1} j - (n-i+1)(i-1) \right]$$



扫描全能王 创建

$$= \sum_{i=1}^n \left[\frac{n(n+1)}{2} - \frac{(i-1)i}{2} - [(n-i+1)(i-1)] \right]$$

separate terms, expand, simplify.

$$= \sum_{i=1}^n \left[\frac{1}{2}i^2 + \left(-n - \frac{3}{2}\right)i + \frac{n(n+1)}{2} + n + 1 \right]$$

$$= \frac{1}{2} \sum_{i=1}^n i^2 + \left(-n - \frac{3}{2}\right) \sum_{i=1}^n i + n \left[\frac{n(n+1)}{2} + n + 1 \right]$$

$\underbrace{(n(n+1)(2n+1)/4)}$ $\underbrace{\frac{n(n+1)}{2}}$

$$= \frac{1}{4}n^3 + \frac{1}{2}n^2 + \frac{1}{3}n \Rightarrow \# \text{ iteration is } \Theta(n^3)$$

Easy to get O-bound.

$$\sum_{i=1}^n \sum_{j=1}^n (j-i+1) \leq \sum_{i=1}^n \sum_{j=1}^n n = n^3 \Rightarrow O(n^3)$$

Page 43.

Floors and ceilings are difficult to work with!

Fact: $\lfloor j/2 \rfloor \leq j/2 \quad \forall j \geq 1$

Create a modified alg.

o replace $(j < \lfloor j/2 \rfloor)$ with $j \leftarrow j/2$

\Rightarrow O-bound.

For a given i , find # iter of inner loop.

Solve for k : $i + \left(\frac{i}{2}\right)^k = 1$ \nwarrow while condition.

$$(gi + klg(\frac{i}{2})) = 0 \Rightarrow \begin{cases} gi = k \\ i <= gi \end{cases}$$

Is k # iterations?

No, but it close

- real number $lg^7 \approx 2.8$

- out by one $lg 8 = 3$ but loop 4 times.



扫描全能王 创建

iters $(k_j + 1) \leq k + 1$

$$\forall i \geq 1 \quad \sum_{j=1}^n c(\lfloor jg_i \rfloor + 1) \leq \sum_{j=1}^n (lg_i + 1) \leq \sum_{j=1}^n (lg_{n+1})$$

$$= n \lg n + n \in O(n \lg n)$$

2017.9.19

Divide and conquer

- Always recourse on a smaller problem

Method 1: Expand "Plug & Chug"

$$T(n) = 2T\left(\frac{n}{2}\right) + cn$$

$$T(k) = 2T\left(\frac{k}{2}\right) + cn$$

$$= 2[2T\left(\frac{n}{4}\right) + c \cdot \frac{n}{2}] + cn$$

$$= 4T\left(\frac{n}{4}\right) + C\left(2\left(\frac{n}{2}\right) + n\right)$$

$$= 4 \left[2T\left(\frac{n}{8}\right) + C \cdot \frac{n}{4} \right] + C \left[2\left(\frac{n}{2}\right) + n \right]$$

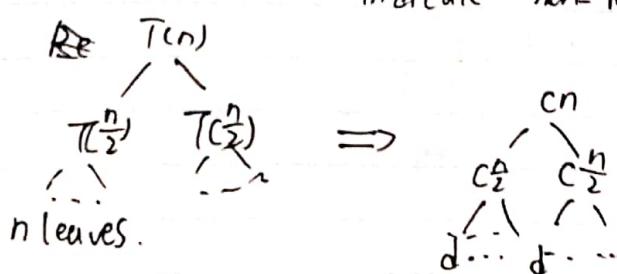
$$= 8T\left(\frac{n}{8}\right) + C\left(4\left(\frac{n}{4}\right) + 2\left(\frac{n}{2}\right) + n\right)$$

$$= \cancel{nT(1)} + C\left(n + 2\left(\frac{n}{2}\right) + 4\left(\frac{n}{4}\right) + \dots + \left(\frac{n}{2}\right)\left(\frac{n}{n}\right)\right)$$

$$= n \cdot d + cn(\lg n - 1) \in \Theta(n \lg n)$$

Method 2: Recursion tree with costs

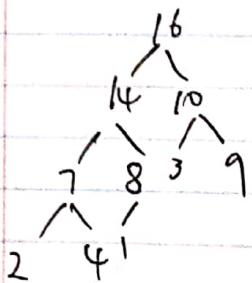
- indicate non-recursive cost at each node.



扫描全能王 创建

Heap: - Not a BST

shape \triangleleft ← last row may not full.



level 0	1	node	2^0
Level 1	2	nodes	2^1
level 2	4	nodes	2^2

Lemma: Height with n nodes is $\Theta(\log n)$

Notation: height of node : # of edges in longest simple path down to leaf.

depth of node: # of edges in longest path up to the root

e.g. node ④ height : 2
depth : 1

pf: Suppose height of tree is h .

$$\begin{aligned}n &\geq 2^0 + 2^1 + 2^2 + \dots + 2^{h-1} + 1 = 2^h \\n &\geq 2^h \Rightarrow \lg n \geq h \lg 2 \\h &\leq \lg n \text{ so } h \in \Theta(\log n)\end{aligned}$$

Lower bound is smaller

$$\begin{aligned}n &\leq 2^0 + 2^1 + 2^2 + \dots + 2^h = 2^{h+1} - 1 \\&\Rightarrow h \geq \lg(n+1) - 1 \\&\text{so } h \in \Omega(\log n) \\&\Rightarrow h \in \Theta(\log n)\end{aligned}$$



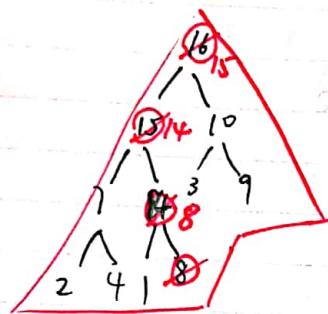
扫描全能王 创建



Insert is

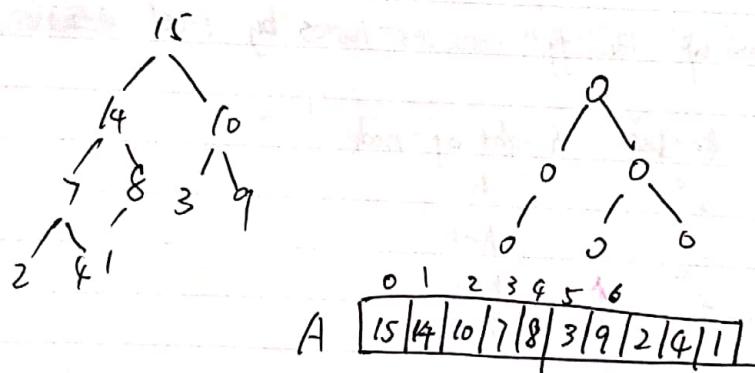
- maintain structure
- maintain order

Bubble-up $\in O(\log n)$



DeleteMax

Bubble-down $O(\log n)$

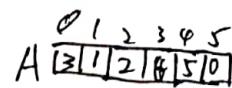


root A[0]

left

Bottom up Heapify

for $i=n-1$ down to 0
bubble down at position i



$n=6$

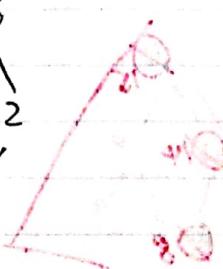
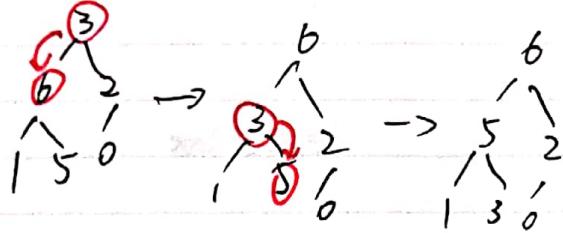
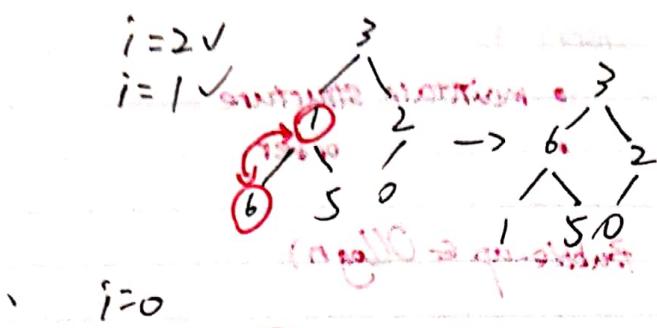
$i=5$

$=4$

$=3$ ↗ We will never bubble down on a leaf node (or node at level h)
 \Rightarrow use $\lfloor \frac{n}{2} \rfloor$ in loop



扫描全能王 创建



2017.9.21.

Runtime: Bottom up Heapify: consider notes by level ~~nodes~~.

level	#nodes	height of node
0	2^0	h
1	2^1	$h-1$
2	2^2	$h-2$
\vdots	\vdots	\vdots
h	2^h	0

nodes with height $i \leq 2^{h-i}$

swaps for bubble down from height i is $\leq i$

$$\text{Total # Swaps} : \sum_{i=0}^h i \cdot 2^{h-i} \leq \sum_{i=0}^h n \cdot \frac{i}{2^i} \text{ and } < n \sum_{i=0}^h \frac{i}{2^i} = 2n \in O(n)$$



扫描全能王 创建

Given a number, can we find its index?

9 3 2 8 7 4 11 12 22 1
v v v v x x x v

$\Rightarrow 9$ belongs at index 6 if sorted

Given a index, what # belongs at that index?

Index q: more difficult or at least not as straight forward.

$\frac{k}{\leq A(p)} \rightarrow A(p)$ input $A[* \dots * A(p) * \dots *] = x$
 $A(p)$

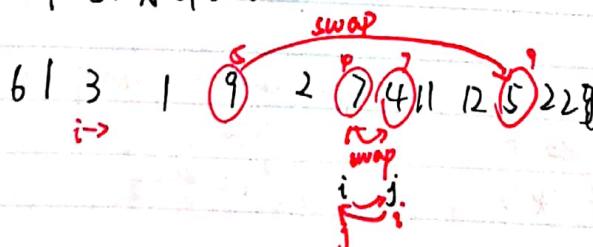
after $A[* \dots * \overset{q}{x} * \dots *] = x$

' P may not be the same as q.

participating is done "in place"

$A[7 3 1 9 2 6 4 11 12 5 22]$

$p=5, A[p]=6$



finally swap $A[i]$ with $A[j]$

$A[6 \underbrace{3 1}_{i} \underbrace{5 2}_{j} \underbrace{6}_{R} 7 \underbrace{11 12}_{i} \underbrace{9 22}_{j}]$

Quick sort selected . choose pivot

given k

partition

Figure out how to proceed.



扫描全能王 创建

- ① if $i=k$ found element at k , done.
- ② if $i>k$, recurse, find k in L .
- ③ if $i < k$, recurse, find $k-i-1$ in R .

worse case: one subproblem of size $n-1$ to recurse on.

$$\begin{aligned}
 T(n) &= T(n-1) + cn \\
 &= [T(n-2) + c(n-1)] + cn \\
 &= T(n-3) + c(n-2) + c(n-1) + cn \\
 &\stackrel{\text{ODD}}{=} d + c \cdot 2 + c \cdot 3 + \dots + c(n-1) + c(n) \\
 &= d + c[(n+2)(\frac{n-1}{2})] \in \Theta(n^2)
 \end{aligned}$$

Base case: we're lucky - find the item in the first pass.

still must partition once $\in \Theta(n)$

Average case?:

calls to printf on average?

```

    foo (k)
    for i=1 to k
        printf('*')
    
```

$1 \leq k \leq 6$

$$\frac{1}{6} \sum_{k=1}^6 k = \frac{21}{6} = \frac{7}{2}$$

(a) A. Avg [b] Consider all inputs
 all inputs are "equally likely."
 take average.

Assumption:

each input is equally likely - "uniform distribution"

Need to consider

- all inputs of a certain size, then take average.
- sum up runtime.
- divide by # inputs.



扫描全能王 创建

Assumption

- no items are repeated.
- behavior of alg depends on relative ordering of keys
- not actual values.

e.g. [2 4 6 8] & [11 12 13 14] both worse case.

\Rightarrow may assume keys are 1, 2, ..., n.

\Rightarrow need to consider all orderings: $n!$

We'll count # comparisons.

Remember: n possible positions for pivot

• each one has $(n-1)!$ permutation of non-pivot elements.

• each pivot location is "equally-likely"
 \Rightarrow divide by n.

Case: ① $i=k$ done

② $i>k$ Search L = $A[0 \dots i-1]$ for item k

③ $i>k$ search R = $A[i+1 \dots n-1]$ for item $k-i-1$

$$T(n, k) = (n+1) \left[\sum_{i=0}^k T(n-i-1, k-i-1) + \sum_{i=k+1}^{n-1} T(i, k) + 0 \right]$$

Let's be lazy: Define $T(n) = \max_{0 \leq k \leq n} T(n, k)$

$A[\underline{\text{xxx}}][\underline{\text{xxx}}][\underline{\text{xxx}}]$ $\frac{1}{2}$ time pivot is here
• subproblem size at most n.
 xxx $\frac{1}{2}$ time pivot is here

• subproblem size at most $(\frac{3}{4}n)$

$$T(n) \leq \begin{cases} cn + \frac{1}{2}(T(n) + T(\lfloor \frac{3n}{4} \rfloor)) & n \geq 2 \\ 1 & n=1 \end{cases}$$



扫描全能王 创建

$$2T(n) \leq 2Cn + T(n) + T\left(\frac{n}{4}\right)$$

$$Th) \leq 2Cn + T\left(\frac{n}{4}\right)$$

$$\leq 2Cn + 2C\left(\frac{1}{4}n\right) + 2C\left(\frac{1}{16}n\right) + \dots + d$$

$$\leq d + 2n \sum_{i=0}^{\infty} \left(\frac{1}{4}\right)^i \stackrel{>4}{\cancel{}} \in O(n), \text{ Average case}$$

Best case: $\Theta(n)$

$\Rightarrow \Theta(n)$, Average case

worse case: $\Theta(n^2)$

Module 3 Page 10

foo1(k) // $1 \leq k \leq 6$

for i=1 to k

printf("*")

$$\text{average case: } \frac{1}{6} \sum_{i=1}^6 i = \frac{7}{2}$$

foo2

k ← ~~rand~~ roll a 6 sided die

foo1(k)

random

"expected" number of calls to print.

$$\sum_{i=1}^6 \left(\frac{1}{6}\right) \cdot i = \frac{7}{2}$$

probability cost of case.
uniform distribution

requires $A[0..n-1]$ is a
permutation of $\underbrace{1, 1, 1, \dots, 1}_{n-1 \text{ ones}}, n$

foo3(A, n)

k=0

for i=1 to A[k]

printf("*")

prints in

best case?

worst case?

foo4(A, n)

k ← random int in 0..n-1

for i=1 to A[k]

printf("*")

let I_1 be fixed $A[1..n]$

$$T^{(exp)}(I_1) = \sum_{k=0..3} T(I_1, k) \cdot \frac{1}{4}$$



扫描全能王 创建

$$T^{(\text{exp})}(I_1) = \sum_{k=0 \dots 3} T(I_1, k) * \frac{1}{4}$$

$$= 1 \cdot \frac{1}{4} + 1 \cdot \frac{1}{4} + 4 \cdot \frac{1}{4} + 1 \cdot \frac{1}{4} = 1 \frac{3}{4}$$

let I_n be input $A[n, 1, 1, \dots, 1]$

$$T^{(\text{exp})}(I_n) = \sum_{k=0 \dots n-1} T(I_n, k) * \frac{1}{n}$$

$$= 1 \cdot \frac{n-1}{n} + h \cdot \frac{1}{n}$$

$$= \frac{2n-1}{n} < 2$$

- input is fixed
- t is random probability of each case uniform dist

Worse case Expected runtime

$$T^{(\text{exp})}(n) = \max_{\text{size}(I)=n} T^{(\text{exp})}(I)$$

- max of all inputs of a certain size

- each input has $\binom{n-1}{\frac{n-1}{2}}$ ones

~~BFR~~ BFPRT

- partition the array into groups of 5 items
- add number, few left over

- take median of each groups of 5
 - ~ 2 items smaller
 - ~ 2 items bigger

- take median of medians



$$T(n) \leq T(\lceil \frac{n}{5} \rceil) + T(\frac{4}{5}n) + Q_n$$

median of medians reorder partition



扫描全能王 创建

Worst case : array sorted.

choose index 0 as pivot.

after partition A [0 1 2 ... n]

0 subproblem n pivot subproblem 2

Recursion Tree:

Worst Case Technique

badly partitioned

0 T(n-1)

0 T(n-2)

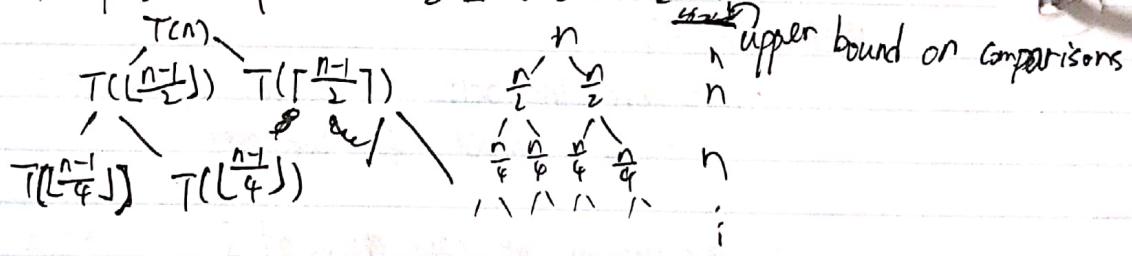
n
0 n-1
0 n-2

$$T^{\text{worse}}(n) = T^{\text{worse}}(n-1) + \Theta(n)$$
$$\Rightarrow \Theta(n^2)$$

to worst case min. /
best max.

Ideally want Balance Partitioning

subproblems of size $\lfloor \frac{n-1}{2} \rfloor$ and $\lceil \frac{n-1}{2} \rceil$



Best Case $\Rightarrow \Theta(n \log n)$

Average Case : - not worst or best.

• maybe $\frac{1}{4}n$ and $\frac{3}{4}n$

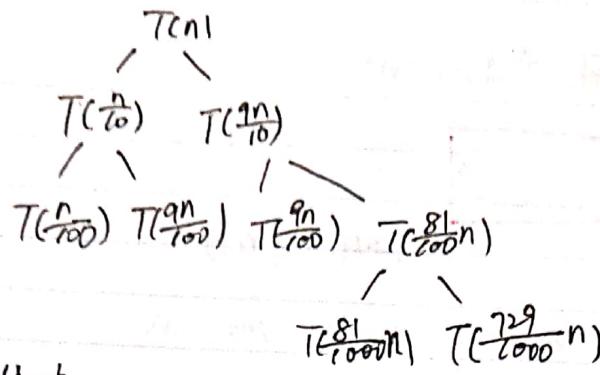
Suppose partitioning always split $\frac{1}{10}n$ and $\frac{9}{10}n$.

$$\# \text{comparisons } T(n) = \begin{cases} T\left(\frac{9n}{10}\right) + T\left(\frac{n}{10}\right) + n & \text{if } n > 1 \\ 0 & \text{if } n \leq 1 \end{cases}$$

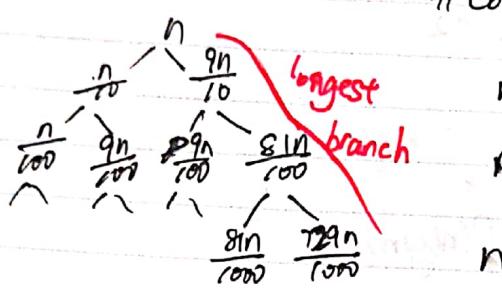


扫描全能王 创建

Recursion Tree



Work



Comparisons

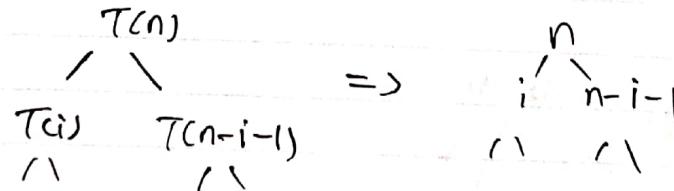
Deepest leaf

Height of tree is minimum ^{al} integer h such that

$$\left(\frac{9}{10}\right)^h \cdot n \leq 1 \Rightarrow h \lceil \log_{\frac{10}{9}} n \rceil \in O(\log n)$$

Total # comparisons $\in O(n \log n)$

Average cost



$$T(n) = \frac{1}{n} \sum_{i=0}^{n-1} [T(i) + T(n-i-1)] + \theta(n)$$

Let $H(n)$ be the average height of the recursion tree.



扫描全能王 创建

$$f(n) \leq 1 + \frac{1}{2}H\left(\lfloor \frac{3}{4}n \rfloor\right) + \frac{1}{2}f(n)$$

$$2f(n) \leq 2 + H\left(\lfloor \frac{3}{4}n \rfloor\right) + f(n)$$

$$H(n) \leq 2 + f\left(\lfloor \frac{3}{4}n \rfloor\right)$$

$$\leq 2 + 2 + H\left(\lfloor \frac{9}{16}n \rfloor\right)$$

By

$\leq 2 \cdot h$ where h is minimal s.t.

$$\left(\frac{3}{4}\right)^h \cdot n < 2$$

$$\Rightarrow f(n) \in O(\log n) \quad : \text{height}$$

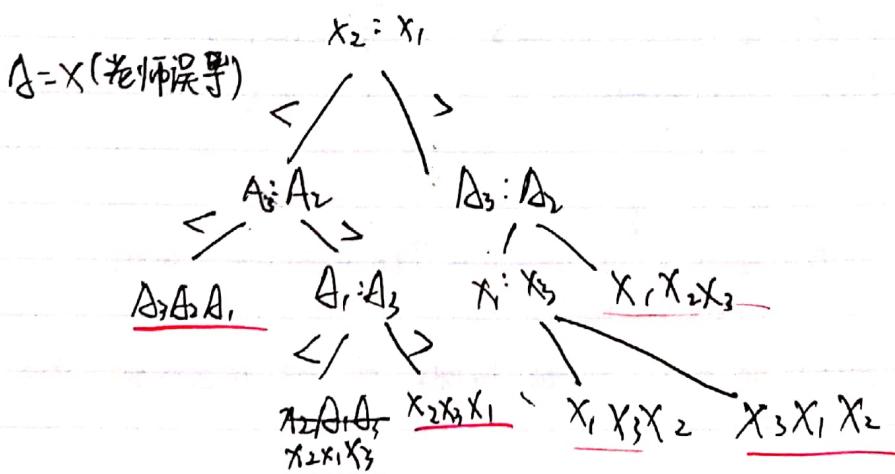
\Rightarrow Average case: $O(n \log n)$

Best case: $\Theta(n \log n) \Rightarrow$ Average case: $\Theta(n \log n)$

Comparisons: $<$, $>$, $=$

Decision tree:

eg. $n=3$ $\{x_1, x_2, x_3\}$ $\left\{ \begin{array}{l} [3|1|2] [2|1|3] [1|2|3] \\ [3|2|1] [2|3|1] [1|3|2] \end{array} \right\}$



扫描全能王 创建

leaves : $n!$

$$\# \text{leaves} = \# \text{internal nodes} + 1$$
$$n! = n! - 1 + 1$$

height $\geq \lg(n!)$

$\geq \Omega(n \lg n)$

Page 23

Input $A[4 5 7 0 5 5 10]$ $R=11$

1. $C[0 0 0 0 0 0 0 0 0 0]$ $n=7$

2. $C[0 1 2 3 4 5 6 7 8 9 10]$

3. $B[4 5 7 0 5 5 10]$ $\# \quad \text{start index}$

4	1
5	2
7	5
10	6

$A[0 4 5 7 10]$

$$A[1] = 4$$

"stable sort"

$$A[2] = 5$$

$$A[5] = 7$$

~~$A[6] =$~~

$$A[0] = 0$$

$$A[3] = 5$$

$\Theta(n+R)$

what if $n=R \Rightarrow \Theta(n)$

{3, 14, 159, 9265432}

8 comparisons
 $n=4 \quad R=9265432$



扫描全能王 创建

MSD Radix sort (Most Significant Digit)

$A[829 \ 331 \ 457 \ 330 \ 436 \ 720 \ 355]$

$n=7, R=10, 3\text{-digit number } (m=3)$

MSD-Radix sort (A)

order by 1st digit.

min. 10. 6
231] 331] 330
330] 330] 331
355] 355] 355

MSD

leftmost to right
sort by digit.
"stable"

457] 434] 434
434] 457] 457

many recursions

720] 720] 720 works on variable-length

829] 829] 829

recursively, using 2nd digit
sort each digit

LSD-Radix Sort:

$A[829 \ 331 \ 457 \ 330 \ 436 \ 720 \ 355]$
 $n=7, R=10, m=3$

330 720 331 355 436 457 829

720 829 330 331 436 355 457

330 331 355 436 457 720 829



扫描全能王 创建

to LSD to MSD

right to left

$\Theta(m(n+r))$

\uparrow

count sort

m digits.

calls to count sort



2017. 10. 3. Module 04

What are we counting for height?

- vertices or edges?
 \Rightarrow edges.

Worst case:

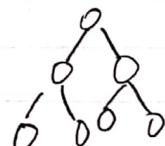
n nodes



$n-1$ edges

Best case:

a balanced

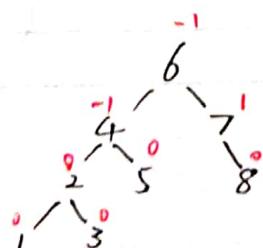


average case: $h\Theta(\log(n))$

- over all $n!$ tree
- similar to quicksort analysis

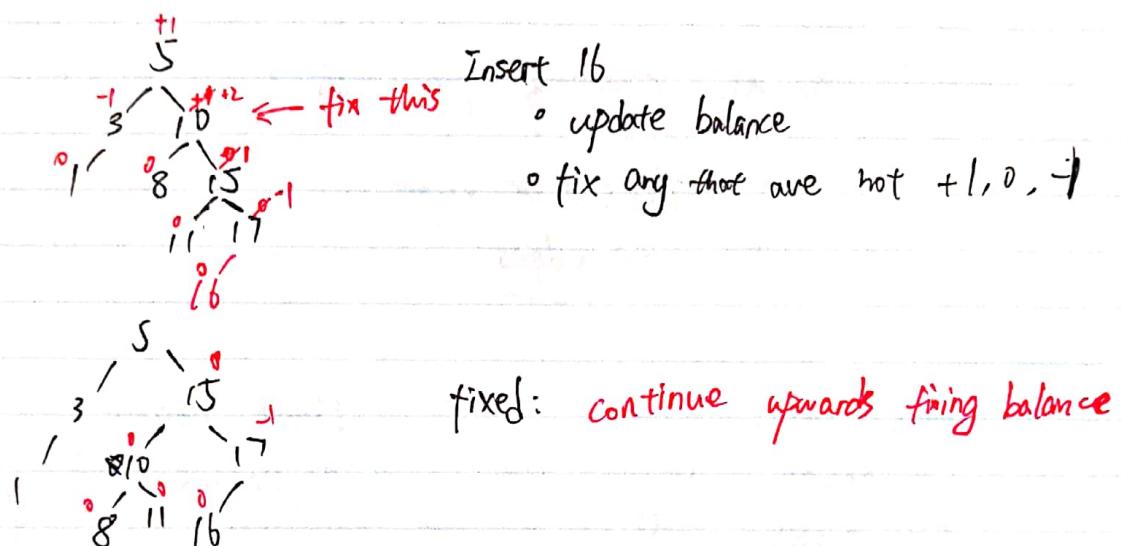
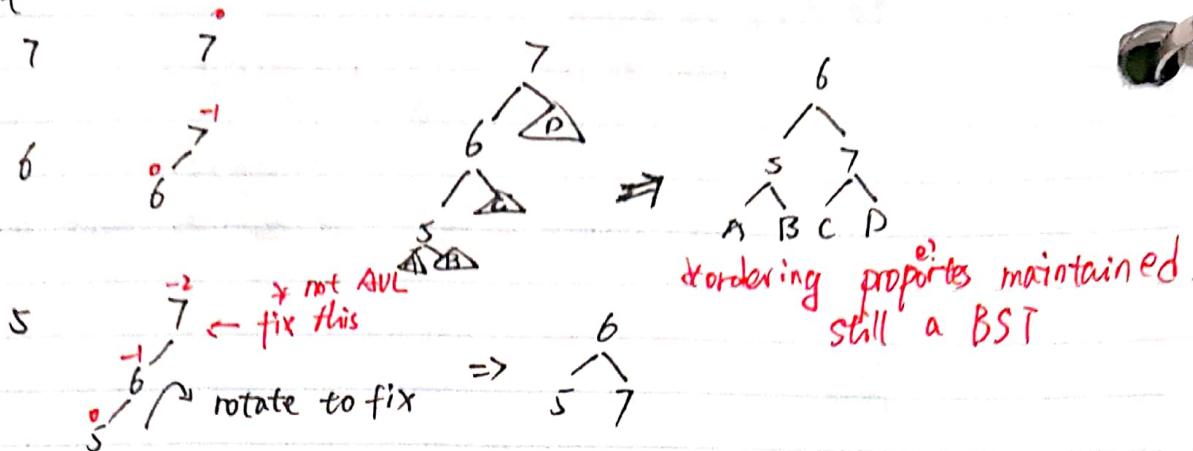
— AVL

- "balanced" - slight imperfection
- every node is "balanced"
- store ~~balance~~ balance factor in the node ??

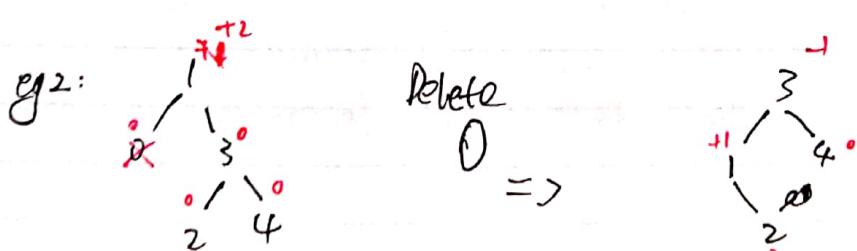
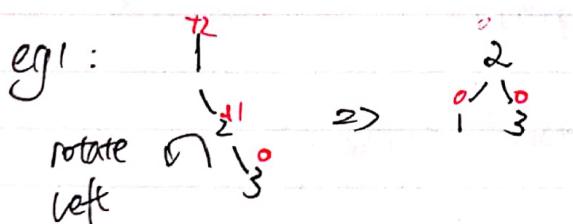


扫描全能王 创建

Insert



* Note: functions return pointer to the new tree.



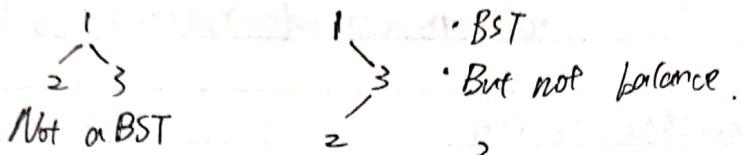
Note: Balance factors are different for case (+2,0) then for (+1,1)



扫描全能王 创建

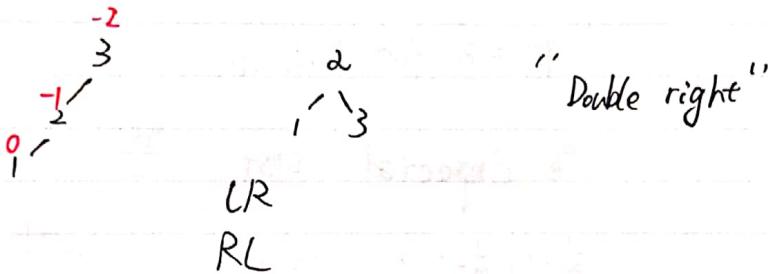


Can 1 be the root?



\Rightarrow ① Cannot be the root.

Idea: do left rotation at ① first then right rotation



2017.10.5.

Self Organizing Search.

- Think of a messg desk with papers all over it.
- The things that are frequency used or recently used are conveniently located on top.
- Could sort by frequency of access
 - sort items by their prop probability of access



扫描全能王 创建

Model.

- Search items 1, 2, 3.
- only count comparisons
 - all unsuccessful searches, cost n.
- successful search.
 - Worst case: n comparisons.
- P_i - probability to search for item x_i
what if uniform distribution?

$$P_i = \frac{1}{n} \text{ for } i=1\dots n.$$

Expected cost: $\frac{n+1}{2}$

$$\frac{1}{n} \cdot 1 + \frac{1}{n} \cdot 2 + \dots + \frac{1}{n} \cdot n = \dots = \frac{n+1}{2}$$

Non-uniform Distribution.

order $A_1 \dots A_n$ where $p_1 \geq \dots \geq p_n$.

Uniform Distribution: Expect $O(n)$

What if $P_i \geq \frac{1}{2^i}$ for $i=1\dots n-1$

$$\text{and } P_n = 2^{-n+1} = \frac{1}{2^{n-1}}$$

$$\text{Expected} = \sum_{i=1}^{n-1} \frac{1}{2^i} + n \cdot 2^{-n} = 2 - 2^{-n} < 2$$



扫描全能王 创建

In practice, Distribution is usually unknown

In practice, 8% of searches are on 20% of keys.

Optimal static Ordering.

proper order $\frac{1}{2} \cdot 1 + \frac{1}{4} \cdot 1$

Move-to-front

Given: 1 2 3 4 5 6

Search 3 3 1 2 4 5 6
4 4 3 1 2 5 6
6 6 4 3 1 2 5

- quick ^{to} steady state
- bad in an array
- affected by a rare lookup
- lots of work to shuffle it back.

Can show $C_{MTF} \leq C_{OPT}$

C_{MTF} cost of MTF
 C_{OPT} cost of optimal ordering

Transpose:

Given 1 2 3 4 5 6
Search 3 1 3 2 4 5 6
4 1 3 4 2 5 6
6 1 3 4 2 6 5



扫描全能王 创建

- slow to steady state
- unaffected by a rare lookup.
- good in an array.

Given 0 1 2 3 4 5 6 7
 Accesses 5 3 5 6 4 6 5 0 3 5 6 4

MTF : 4 6 5 3 0 1 2 7

- 54 comparisons.

Transpose: 0 1 5 3 6 4 2 7

- 62 comp.

2017.10.12

Skip Lists - Binary search or linked list.

• Sorted arrays: binary search.
 • must shift items - insert & delete

• Linked Lists:
 • easy to move items around.
 if we know where they are.
 • no binary search

AVL:

Discovered by Bill Pugh (1989)

Competes with balanced tree structures.
 • often after wins.



扫描全能王 创建

Randomized - random element as part of alg.

Idea: Given singly linked list - search O(n) so

add an express lane S_1

Search(23)

Implementation: Linked Nodes

- different # pointers per node so vs S_1, S_2, S_3, \dots

- pointers will change with delete and insert.

Doubly Linked List?

- not necessary lots of nodes
- lots of waste space.

Ideally, want shortcut (towers)

Search S_2

- get stack S from search

- Random element: height of tower for new node.

coin toss: H T $\Rightarrow i=1$ (2 levels)

- add s_2 to S_0 and S_1

- add s_2 after s_4 on S_0

- after s_7 on S_1

Insert 100.

coin toss H/H/T $\Rightarrow i=3$ (4 levels)

- insert new level

- height increases.

* top level is always $\boxed{-\infty} \rightarrow \boxed{+\infty}$ by defn.

How tall can a tower be? very tall.



扫描全能王 创建

Starts: probability of H (or T) = $\frac{1}{2}$

probability of tower of height 3.

$$\frac{1}{16} \Rightarrow P(HHT)$$

probability of no H's followed by one T:

$$(\frac{1}{2})^n = \frac{1}{2^{10}} \leq 0.05\%$$

controlling Height:

1) Set max height to be a constant
as $n \rightarrow \infty$, not enough levels.

2) have height $h \leq f(n)$, eg $f(n) = \lceil 3 \log n \rceil$ works well.

3) allow $h \rightarrow \infty$, not likely to get too tall.

Delete

- update links
- free nodes
- remove any extra top levels.
• may decrease height.

space: expect number of keys of level i :

• probability of getting i or move H in a row.

n node on S_0

Expected # nodes on $S_1 = \frac{n}{2}$

$$S_2 = \frac{n}{4}$$

$$S_3 = \frac{n}{8}$$

$$S_4 = \frac{n}{16}$$

Expected # levels: $\Theta(\log n)$

Expected # nodes at level $S_i = \frac{n}{2^i}$



扫描全能王 创建

Sum over all levels: $\sum_{i=0}^{\log n} \frac{n}{2^i} = 2n - \frac{n}{2^{\log n}} \approx 2n \in O(n)$

Time:

worst case for search/insert/delete: $O(n \log n)$

Expected height: $O(\log n)$

Expected number: $O(n \log n)$?

- Don't usually scan all ~~n~~ n elements

Refined analysis:

Idea: Backwards analysis: follow path. in reverse. always move up if you can.

How long is backwards path before rising k levels..

Probability Node will have a tower to go up?

$$C(n) = \frac{1}{2}$$

$C(n)$ = # steps. to go up n levels.

• 1 step then either.

• go up, prob $\frac{1}{2}$: $C(n-1)$

• go left, prob. $\frac{1}{2}$: $C(n)$

\Rightarrow Expected # steps to go up n levels is

$$C(n) = 1 + \frac{1}{2} C(n-1) + \frac{1}{2} C(n)$$

$$2C(n) = 2 + C(n-1) + C(n)$$

$$C(n) = 2 + C(n-1)$$

$$C(n) = 2 + 2C(n-2)$$

...

$$= 2n \quad \text{where } n \text{ is # of levels}$$



扫描全能王 创建

Expected

$O(\log n)$ levels $\Rightarrow O(\log n)$ steps

• constant time at each level
expected

2017.10.17.

Tries "try"

Binary Trees

- store collection of binary strings (key)
- runtime analysis based on length of string not number of string in tree.

Compressed Trie (Patricia Trie).

- eliminate unflagged nodes with only one child.
 \Rightarrow save space

Multway-Tries

• extension to arbitrary alphabet.

- keys can be variable length. could stop associated value
- keys are not explicitly stored in trie.
- edge labels are implicit: left: 0 } we show them for right: 1 } convenience

Module 6 Slides 3/27

Slides 18/27

Search (01001)
1 2 3 4

Check Bit

0 0 \Rightarrow go left

1 1 \Rightarrow go right

2 0 \Rightarrow go left



扫描全能王 创建

At leaf: compare stored key (in tree) with search key 0/001.

• match (still $\Theta(1 \times 1)$)

Search (101)

Check Bit

0 1 \Rightarrow go right

1 1 \Rightarrow go right.

At leaf: compare stored - search key

• x match

Search (1110)

Check Bit

0 1 \Rightarrow go right

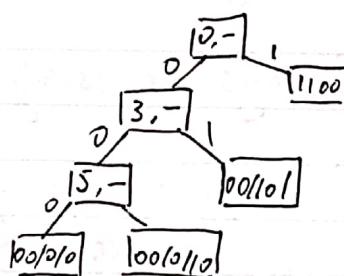
1 1 \Rightarrow go right

At leaf: compare search - stored key.

• x match (arrive at some leaf with search (101))

Slide 25/27.

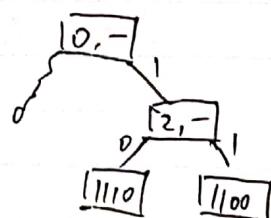
Bit	0	1	2	3	4	5	6
0	0	1	1	0	1		
0	0	1	0	1	0		
0	0	1	0	1	1	0	
1	1	0	0				



Try to find bits where strings disagree.

Insert a new key (1110)

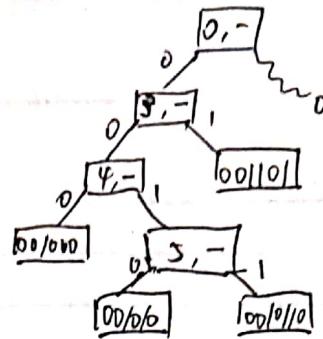
- search ends at leaf [1100]
- disagree at index 2.
- insert a new node



扫描全能王 创建

Insert (001000)

- search ends at leaf [001000]
- disagree at index 4
- insert a new node
- goes between nodes with index 3 and 5.

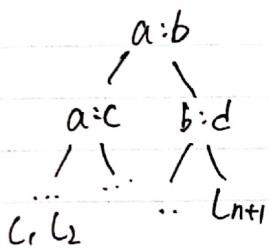


Module 07.

- n keys

Thm: Comparison model $\Omega(\log n)$ comparisons required.

- each comparison has two outcomes,
- any algorithm can be modeled by binary tree.



• Total # of outcomes is: $n+1$

- one of n keys.

- or "not found"

\Rightarrow need $n+1$ leaves

• binary tree \Rightarrow ~~n~~ $\lceil \log n \rceil$ nodes

\Rightarrow height $\approx \lceil \log n \rceil$

Dictionary as an array - problems?

① Keys are usually not simple eg strings.

② storage if $n \ll M$

often # keys used \ll # possible keys.

Direct Address is OK if we can allocate array with one position per key.

Idea: try to reduce amount of wasted space.



扫描全能王 创建

Let M be # buckets in hash table $T \rightarrow \text{array } T[0, \dots M-1]$
let $n = \# \text{ actual KVPs stored in table}$

Goal:

- reduce storage to $\Theta(n)$
- still get $O(1)$ time for all operations at least in average
not necessarily 1 but a constant #.

a key k "hashes" to index $h(k)$ in T
if $M < n$ then ... Problem? two or more keys hash to the same bucket.
 $\Rightarrow \text{Collision}$

more possible keys than indices

- if $n \leq M$ ~ may or may not happen
- if $n > M$ ~ definitely happen. **pigeon hole principle**

What is the likelihood of a collision?

Birthday Paradox: suppose n people with "random" birthdays
 $M = 365$.

What is the probability that 2 people have same birthday?
if $n > 365 \Rightarrow$ probability is 1 - 100%

Probability of a shared birthday is: $1 - \frac{\binom{365}{n}}{365^n}$

$$n=2 \Rightarrow 0.2\%$$

$$n=10 \Rightarrow 12\%$$

$$n=23 \Rightarrow 50\%$$

$$n=50 \Rightarrow 97\%$$

Conclusion: collisions are very likely even when $n \ll M$

must be prepared to handle them

Challenges:

- good hash function

- dealing with collisions $h(k_1) = h(k_2)$
 $k_1 \neq k_2$



扫描全能王 创建

"Good" Hash function

- be unrelated to patterns in the data.

- depend on all parts of the key

- be efficient to compute

Chaining Chaining analysis : n kVps in table T.

~~M buckets~~ in table T.

$$\text{Load factor: } \alpha = \frac{n}{M}$$

Search : cost depends on length of particular list that we hash to ~~unsuccessful~~ unsuccessful search.

compute

- assume ~~complete~~ hash value in O(n)

- assume uniform hashing - each bucket in T is equally likely

- average length of list $\frac{n}{M} = d$

$\Rightarrow d$ ~~comparisons~~ comparisons on average

$\Rightarrow \Theta(1 + \alpha)$ why not $\Theta(\alpha)$? It's possible that $\alpha < 1, = 1, > 1$

hash to a bucket

Worst case! $\Theta(n)$

Successful Search

- # comparisons to find k.

- compute hash value + average search time on average length list.

e.g. list of size 4:

$$\text{Average # searches : } (1+2+3+4) / 4 = 2.5$$

$$= 1 + \frac{1}{2} \left(\frac{n}{M} + 1 \right) = 1 + \frac{\alpha}{2} + \frac{1}{2}$$

$$= \Theta(1 + \alpha)$$

$(\frac{n}{M} = \alpha)$

$\frac{1}{2} \left(\frac{n}{M} + 1 \right)$, half list (average) search



扫描全能王 创建

Search (average case) : $\Theta(1/\alpha)$

Insert : $O(1)$

Delete : same as search

if $n \leq O(M)$

\Rightarrow average cost are all $O(1)$

Open Addressing - Linear Probe Probing.

Delete : don't simply delete item.

* must flag as deleted - lazy deletion

Insert : scan until empty or flagged (deleted) slot found.

Search : scan until empty (continue if flagged (deleted)) or found

Many deletions \Rightarrow many slots flagged

- Adv:
 - space efficient
 - easy to implement

Thm: The average # Probes required to search in hash table of size M with n keys is approx:

$$\text{successful search} : \frac{1}{\alpha} \left(1 + \frac{1}{1-\alpha} \right) \quad \left(\frac{1}{2} \left(1 + \frac{1}{1-\alpha} \right) \right)$$

$$\text{unsuccessful} : \frac{1}{\alpha} \left(1 + \frac{1}{(1-\alpha)^2} \right) \quad \left(\frac{1}{2} \left(1 + \frac{1}{(1-\alpha)^2} \right) \right)$$

for α not too close to 1.

Primary Clustering

- a sequence of slots gets full
- the bigger the cluster gets, the quicker it continues to grow.
- the larger area it spans
 - \Rightarrow larger chance that new random insertions will go in the cluster
 - \Rightarrow cluster expands



扫描全能王 创建

long search/insert	α	$\frac{1}{2}$	$\frac{2}{3}$	$\frac{3}{4}$	$\frac{9}{10}$
$\frac{1}{2}(1 + \frac{1}{1-\alpha})$ hit		1.5	2.0	2.5 3.0	5.5
$\frac{1}{2}(1 + \frac{1}{(1-\alpha)^2})$ miss		2.5 2.5	5.0	8.0 8.5	58.5

Double Hashing

Idea: Stop formation of clusters by using 2 independent hash functions.

- if $h_1(k_1) = h_2(k_2)$ collide
- then $h_1(k_1) \neq h_2(k_2)$ unlikely but possible.

eg. $h_1 = k \bmod 11$ some floating point constant.
 $h_2 = L \lceil A \cdot k \rceil - \lfloor A \cdot k \rfloor + 1$ $h_1: U \rightarrow \{0, 1, \dots, M-1\}$
 $h_2: U \rightarrow \{1, \dots, M-1\}$

h_2 is an offset from original hash location

Double hashing: $h(k, i) = [h_1(k) + i \cdot h_2(k)] \bmod M$

try $i = 0, 1, 2, \dots$ until find open slot.

$M=11$ eg $k=11$ probe sequence: 0, 8, 5, 2, ... linear probing
 $k=22$ 0, 6, 1, 7, ... both are 0, 1, 2 ...

eg $M=2^{10}=1024$ Suppose $h_1(k)=18$ and $h_2(k)=256$ for some k
probes! 10, 266, 522, 778, 10, 288, 522.
 h_2 should be relatively prime to Hash Table size.

\Rightarrow no common factors.

\Rightarrow choose prime # for M and $1 < h_2(k) < M$.

'no common factors'

Analysis: Average # probes:

successful: $\frac{1}{2} \ln \left(\frac{1}{1-\alpha} \right)$

unsuccessful: $\frac{1}{1-\alpha}$

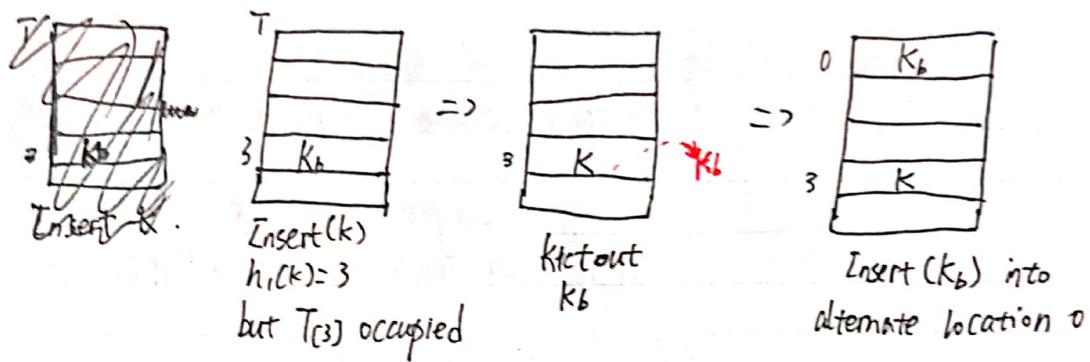
α	$\frac{1}{2}$	$\frac{2}{3}$	$\frac{3}{4}$	$\frac{9}{10}$
hit	1.4	1.4	1.8	2.6
miss	1.5	2.0	3.0	5.5



扫描全能王 创建

Cuckoo Hashing:

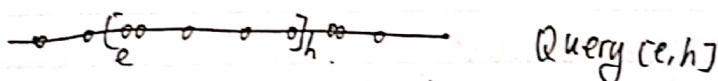
- * One KVP per slot
- * 2 independent hash functions
- * always insert k to $T[h_1(k)]$
 - may "kick out" another item
 - if so, insert "kick out" item at its alternate location



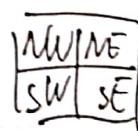
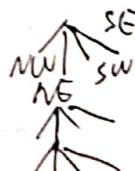
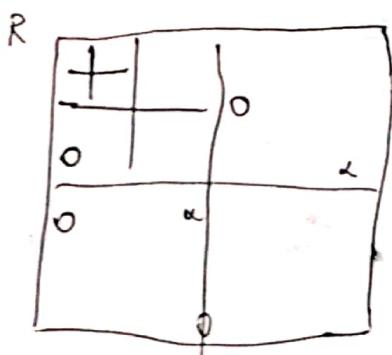
- Keys:
- letters
 - strings "apple" \Rightarrow ~~is~~ index in hash table
each letter has an ASCII value.

Module 8

Range Search 1-D



- binary search for l - $O(\log n)$
- go right until $> h$
- let k be size of output
 $\Rightarrow O(\log n + k)$ * ~~output sensitive~~



扫描全能王 创建

If point $p \in P$ on x mid, then p in NW v SW

If point $p \in P$ on y mid, then p in SW v SE

Lemma - Depth of Quadtree is at most $\log(s/c) + \frac{3}{2}$

where s ($= d_{\max}$) is side length between of initial square

c ($= d_{\min}$) is minimum distance between any two points in P .

Proof: at depth i , size of square: $\frac{s}{2^i}$

\Rightarrow max distance between 2 points in the same region:



$$\text{length diagonal } \sqrt{2(s/2^i)^2} = \frac{\sqrt{2}s}{2^i}$$

depth i of internal nodes satisfies $\sqrt{s/2^i} \geq c \Rightarrow \frac{\sqrt{s}}{c} \geq 2^i$

$$\text{Solve for } i: \log(\frac{s}{c}) + \frac{1}{2}$$

2011. 10. 31

Quadtree

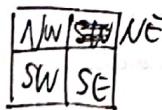
• Define a bounding square

\Rightarrow root node

• all points

• if region contains > 1 pt, \Rightarrow each square is a child
subdivide region into 4 squares

• if no points in region,
do not add child.



convention! (left & lower)

• order matter

• keep subdividing until only 1 pt per region

\Rightarrow all points are leaves
• not based on n .

Depth of Quadtree $\leq \log(\frac{d_{\max}}{d_{\min}}) + \frac{3}{2}$

Build tree: $O(nh)$

Range Search: $O(nh)$ ~ follow branches of tree that intersect query rectangle
Often performs well in practice.



扫描全能王 创建

- easy to implement
- easy to extend to higher dimensions



main problem: based on splitting region equally regardless of distribution of points.

- Some region will have many points \Rightarrow region many subdivisions.

\Rightarrow bad performance for non-uniformly distributed points.

kd-tree

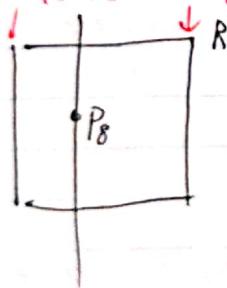
Idea: split points instead of by area

- \Rightarrow always \approx half points on each side of dividing line
- alternate between vertical & horizontal divisions.

Convention: "other" points on vertical split line \Rightarrow left
horizontal \Rightarrow lower

make sure to sort points by x and/or y coord
ONLY ONCE!

R.leftside R.rightside



- Search left and right



Search only left



search right only



扫描全能王 创建

Analysis: We want better than $\Theta(n)$

Runtime is bounded by # calls to

- ① reporting a point in R
- + ② recursing on C because C.region $\cap R \neq \emptyset$

① let k be # points reported $\Rightarrow \Theta(k)$

② $\leq \# \text{ nodes } C \text{ in } T$:

$S \cap C.\text{point} \in R$

+ U number of regions we check unsuccessfully

- regions that intersect R but not fully in R.

- points not in R but are close

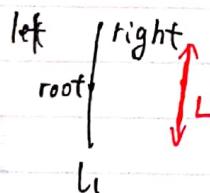
$$|S| \in \Theta(k)$$

Need to bound U

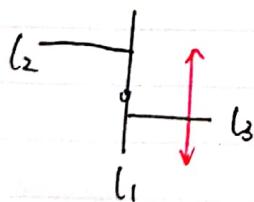
$|U| \leq \# \text{ nodes in } T \text{ such that } C.\text{region is intersected by } L$
an edge of R

Let C be a fixed vertical line (eg an edge of R)

Let $Q(n)$ be # of nodes in T s.t. C.region is intersected by L

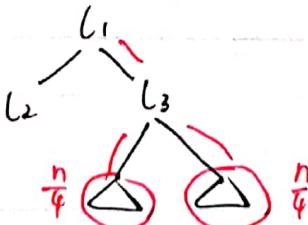


- L intersects exactly one region
- either left or right.



- L intersects both upper & lower region of l3
- but neither for l2

\Rightarrow



扫描全能王 创建

$$Q(n) = \begin{cases} Q(1) & \text{if } n=1 \\ \cancel{Q(1)} + \underline{2Q(\frac{n}{4})} & \text{if } n>1. \end{cases}$$

Solve:

$$Q(n) = 2Q\left(\frac{n}{4}\right) + \underline{\theta(1)} = 2Q\left(\frac{n}{4}\right) + c$$

a constant c

$$= 4Q\left(\frac{n}{16}\right) + 2c + c$$

$$= 8Q\left(\frac{n}{64}\right) + 4c + 3c + c$$

$$= 2^{\log_4 n} \cdot Q(1) + \sum_{i=0}^{\log_4 n} 2^i \cdot c = O(2^{\log_4 n})$$

$$2^{\log_4 n} = \cancel{2^{\log_2 n}} \cdot 2^{\frac{\log_2 n}{4}} = (2^{\log_2 n})^{\frac{1}{4}} = (n^{\log_2 2})^{\frac{1}{4}} = (n^{\log_2 2})^{\frac{1}{4}}$$

$$= n^{\frac{1}{2}} = \sqrt{n}$$

• Same holds for all 4 edges of R: $O(4\sqrt{n}) \Rightarrow O(\sqrt{n})$

Runtime: $O(k + \sqrt{n})$ where k is number points reported.

2017.11.2

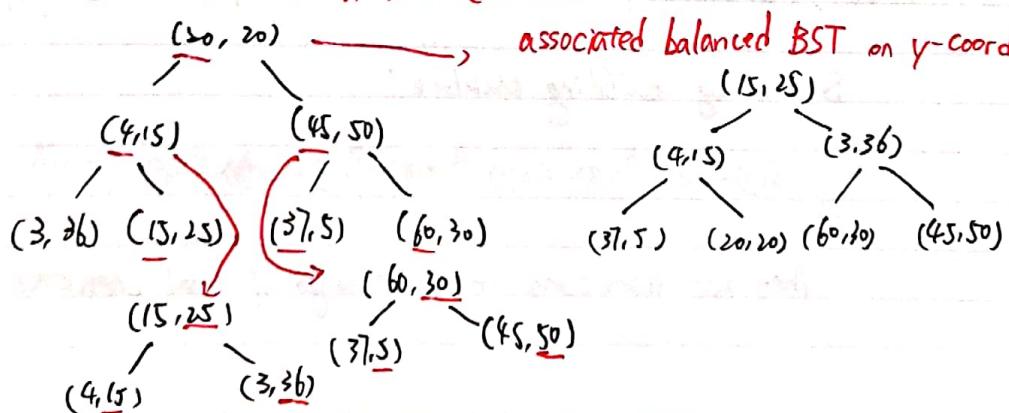
Range Trees 2-D

Let $P(v)$: points stored in subtree rooted at v

• for every initial v, point to an associated BST (balanced) on y-coordinates storing points $P(v)$

e.g. given points: (3, 36), (4, 15), (15, 25), (20, 20), (37, 5), (45, 50), (60, 30)

Balanced BST on x-coord



扫描全能王 创建

Inside: "top" inside node v

\Rightarrow perform range search on y-coords for $[y_1, y_2]$ in ~~Yasser (D)~~

- do not need to check x-coords.

~~Step~~

Space: each point is stored in $O(\log n)$ y-BSTs

- in each tree on path up to root

- n points

$$\Rightarrow O(n \log n)$$

Runtime: at most $2 \log n$ internal "top" nodes (subtrees)

- do range search on y-coords BSTs $\sim O(\log n + k)$ for each.

$$\Rightarrow O((\log n)^2 + k)$$

Note: Can be improved to $O(c \log n + k)$

Idea: We ~~search~~ search the same y range $[y_1, y_2]$ over & over

- only do it once & use the result many times

\Rightarrow fractional cascading.

Pattern matching

e.g. \Rightarrow check $\Theta(m)$ characters

$\Rightarrow \Theta(n-m+1)$ times

\Rightarrow total $\Theta(nm)$ usually $m \ll n$.

DFA \circ given a pattern, the DFA is a pattern matching checker

Start by building backbone:



Then add transitions on where to go if final character found.



扫描全能王 创建

- for a string s , $l(s) \in \{0, \dots, m\}$ is the length of the longest prefix of P that is also a suffix of s .

eg. $P = \underline{ababaca}$
 $s = \underline{ababaa}$

- a is the longest prefix of P that is a suffix of s .
 $\Rightarrow l(s) = 1 \quad \delta(s, a) = 1$
 from state s on symbol a goto state 1 .

eg. $P = ababaca$
 $s = ababab$
 $\Rightarrow l(s) = 4$

$$\delta(s, b) = 4$$

T has size n

P has size m

Goal: $O(n)$ runtime in worst case

KMP

- Compare pattern left to right
- when a mismatch occurs shift pattern "intelligently"
- similar to mismatch transitions in DFA

g. Guess: $i:j$ check: $i:j$
 $T: \dots b a c b \overset{i}{\cancel{a}} \overset{j}{\cancel{b}} a b a a b c$
 $P: \overset{a}{\cancel{b}} \overset{a}{\cancel{b}} b a a c$
 $0 \ 1 \ 2 \dots$
 j

Main loop invariant: if checking index j of P with i of T
 then $P[0:j-1] = T[i:j-1]$

- previous characters of P match T .



扫描全能王 创建

Case to consider when checking $T[i:j]$ with $P[j]$

- cases: $T[i:j] = P[j]$
- $T[i:j] \neq P[j]$

Boundary cases: $j=0, j=m-1$

While $i < n$ do

- case1: $T[i:j] = P[j]$ and $j < m-1$
 - continue matching $\Rightarrow i++, j++$.

- Case2: $T[i:j] = P[j]$ and $j = m-1$
 - whole pattern matched
 - return $i-j$, first occurrence of P in T

- case3: $T[i:j] \neq P[j]$ and $j = 0$

eg. $T: \dots cb ab ab \dots$ mismatch on 1st character
 $P: \cancel{a} b a b c a$ of pattern
 j
 $\Rightarrow i++, \text{keep } j=0$

- case4: $T[i:j] \neq P[j]$ and $j > 0$

$T: \dots b a c b a b ab a | a c b a \dots$
 $P: d b \cancel{a} b \cancel{a} d | c a$
 j

- mismatch not at the first character.

Brute Force: go back in text to $i-j+1$ and start checking pattern at $j=0$.

KMP - shift pattern more intelligently.

- keep i same
- decrement j by correct amount *intelligently*.

Must maintain loop invariant: $P[0, \dots, j-1] = T[i:j, \dots, i-1]$



扫描全能王 创建

failureArray(P)

P = abacaba

$$F[0] = 0$$

$$i = 1$$

$$j = 0$$

$$5 < 7$$

$$P[5] = P[0]$$

$$F[5] = 1 + 1 = 2$$

$$i = 6$$

$$j = 2$$

$$1 < m$$

$$P[1] \neq P[0] \quad (a \neq b)$$

$$j = 0 \Rightarrow$$

$$F[0] = 0$$

$$i = 2$$

$$6 < 7$$

$$P[6] = P[2]$$

$$\Rightarrow F[6] = 3 + 1 = 3$$

$$i = 6$$

$$j = 3$$

$$2 < m$$

$$P[2] = P[0] \quad (a = a)$$

$$F[2] = 1$$

$$i = 3$$

$$j = 1$$

end

$$3 < m$$

$$P[3] \neq P[1] \quad (b \neq c)$$

$$j > 0 \Rightarrow$$

$$j = F[1 - 1] = F[0] = 0$$

$$3 < m$$

$$P[3] \neq P[0] \quad (a \neq c)$$

$$j = 0 \Rightarrow$$

$$F[3] = 0$$

$$i = 4$$

$$4 < 7$$

$$P[4] = P[0] \quad (a = a)$$

$$F[4] = 0 + 1 = 1$$

$$j = 1$$

$$i = 5$$



扫描全能王 创建

How do we update j ?

- without moving to backwards in T - currently at index i
- still need previous characters $T[i-1]$ $T[i-2]$... to match pattern

\Rightarrow we need to choose a maximal $j_{\text{new}} < j$ such that:

$$P[0 \dots j_{\text{new}}-1] = T[i-j_{\text{new}}, \dots i-1]$$

Define $F[j-1]$ as the index we will have to check in P after we bring the pattern to its next possible position.

i.e. $F[j-1] = j_{\text{new}}$

eg: $T: \dots b a c b a b a b a |^i a c b a$
 $P: \begin{array}{cccccc} | & | & | & | & | & \\ a & b & a & b & a & x \\ 0 & 1 & 2 & 3 & 4 & 5 & 6 \\ & & & & & | & \\ & & & & & & j \end{array}$

$j=5$ mismatch

$$F[4] = ?$$

Find maximal prefix of $P[0 \dots j-1]$: ababa

that matches suffix of $P[1 \dots j-1]$ ba ba
length 3

* proper suffix

$$F[4] = 3$$

P=abacaba

$j=5$

$P[0 \dots j]: abacab$

$P[1 \dots j]: bacab$

length 2

$$F[5]=2$$

Confusing part: Suppose there is a mismatch at $j=6$

look up $F[5]$, $j_{\text{new}} = 2$

shift pattern to right.



扫描全能王 创建

Case 4: $T[i] = P[j]$ and $j > 0$

- update $j = F[j-1]$

Last Occurrence

- mapping for all characters in alphabet Σ

e.g. $\Sigma = \{a, b, c, d, \dots\}$

C	a	b	c	d	...	Σ
Loc	4	5	3	-1	...	-1

Good suffix Arrays or Suffix skip arrays

- Similar to KMP failure array + extra condition
- Tougher to understand.

To compute: $S[i]$ • mismatch occurred at $P[i:i]$

- Shift pattern to right by positive amount so that
 - ① any pattern characters below already matched test characters must match up.

And ② $P[j:j]$ does not match $P[i:i]$

- return $j = i - \text{shift amount}$.

Find the largest j s.t. $P[i:i, \dots, m-1] = P[j:j+1, \dots, j+m-1-i]$ ①

and $P[j:j] \neq P[i:i]$ ②

ex: $P = \underset{01234567}{babobobo} \Rightarrow m=8$

$i=7 \quad P[8..7] = P[j:j+1, \dots, j] \Rightarrow j=7$ but $P[7]=P[7] \times$ try again
 $\Sigma - \text{no characters} \Rightarrow j=6$ so $P[6] \neq P[7]$

$i=6 \quad P[7..7] = P[j:j+1, \dots, j+1] \Rightarrow 'o' \text{ at index } 5$
~~'o'~~ $P[4]=P[6]$ try again



扫描全能王 创建

$\Rightarrow 'o'$ at index 3

$P[2] \neq P[6]$ $\leftarrow j=2$

$i=5 \quad P[6 \dots 7] = P[j+1 \dots j+2] \Rightarrow "bo"$ at index 4/5, $j=3$
 $P[3] = P[5]$ try again

$\Rightarrow "bo"$ at index 0/1, $j=1$
 $P[1] \checkmark$

$i=4 \quad P[5 \dots 7] = P[j+1 \dots j+3]$
 $"obo"$ $\Rightarrow "obo"$ at 3/4 $\Rightarrow j=2$
 $P[2] \neq P[4]$ \checkmark

$P = \underline{b} \underline{o} \underline{b} \underline{o} \underline{b} \underline{o} \underline{b} \underline{o}$

$i=4 \quad \text{match } "obo" \ \cancel{167}$

$\text{mismatch at } 'b' \ ^{\cancel{4}}$

try: $j=2 \times "obo"$ at 3/4, 5
 $j=0 \times "obo"$ at 1/2, 3

~~43-2-101234567~~
~~***bobobobo~~

$j=-2 \checkmark$

Rabin-Karp Fingerprint Alg

Idea: don't compare string
• compare hash values of string instead
finger print

Don't explicitly create the hash table
but need prime numbers (table size) for hash function.

Recall: Good Hash function properties



扫描全能王 创建

\Rightarrow If 2 hash values are the same, they are likely the same key.

Idea: Find hash value for Pattern $h(x) = x \bmod 97$

$P: 59265$

$h(P) = 95$

• only compute this once

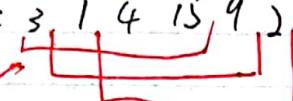
• then use $h(P)$ to find possible matches.

Search Text:

- computing hash values for strings of similar length to P

Graph①: $T: 3 1 4 15 9 2 6 53 \dots$

compute
hashvalue



"rolling hash"

• keeps on shifting right by 1 character

- compute hash for substring in text compare with $h(P)$
 - if match: "probably" the string we are looking for
 - use native approach to verify match $\Theta(m)$

Pattern has length $m \Rightarrow \Theta(m)$ to hash pattern.

If we must compute a hash for all possible entries of text,

• $n-m+1$ possible places to check for a match.

$\Rightarrow \Theta(nm)$ to search for a miss.

Robind Karp found a way to roll hash more efficiently.

2017.11.14

Update: how to compute the next hash value of text from previous.

Suppose Initial value is $41592 \bmod 97 = 76$

$$\frac{4x_{10}^0 + 1x_{10}^1 + 5x_{10}^2 + 9x_{10}^3 + 2x_{10}^4}{(1) \quad (2)}$$

\Rightarrow Graph①

If the next hash is: $15926 \bmod 97$

$$= \underline{\underline{4x_{10}^4}} - 0U$$

$$+ \underline{\underline{(1592)}} \text{ by } 10 - 0U?$$

$$+ \underline{\underline{6}} \text{ by } 10 - 0U$$

$$\bmod 97 - 0U$$

or to keep numbers smaller
work with mod values

$$10^4 \bmod 97 = 9$$

$$\text{old hash value: } (76 - (4 \times 9)) * 10 + 6 = 406 \bmod 97 \\ \Rightarrow 18$$



扫描全能王 创建

Work with binary keys - strings of {0,1}

$$H(T[i], \dots, i+m] = H(T[i, \dots, i+m-1]) \times 2 \\ - T[i] \times 2^m + T[i+m] \bmod M \quad / \text{bitwise operations}$$

Table size:

- prime number
- big enough so collisions are likely

When do we get collisions?

- $x = y \Rightarrow H(x) = H(y)$
- Also can have $H(x) = H(y)$ but $x \neq y$

iff M divides $|x-y| \Rightarrow$ naive check of keys = failure, no match

From Number Theory: # primes that divide $A < 2^n$ is ~~$\frac{n}{\ln n}$~~ $\frac{n}{\ln n}$

Pick $M = 2^{n^2}$ "large enough"

size of text

$$\text{then } \Pr[\text{failure}] \frac{n}{\ln n} \cdot \frac{\ln n^2}{n^2} = \frac{2}{n}$$

$$\text{if } n=4000, \Pr[\text{failure}] \leq 10^{-3}$$

Module 10

Reliability • very important for hard drives, signal transmission

=> error correcting codes

- store & send some extra info to check if packet arrived ok



扫描全能王 创建

eg. transmit 7 bit characters: 0110100, 1101100, 0101110
parity bit { 1 if odd # of 1's
0 if even # of 1's

=> 0110100 1

=> 1101100 0

=> 0101110 0

Consider the following code:

E=1010 ~~S~~

S=11

O=1011

Y=01

N=0110

• all codes are distinct

• ask the Oracle a question

It will give you an encoded answer

Will I pass CS 260?

Answer: 0 1101011
Y E S

0110 1011
N O

2017.11.16

Huffman coding

• Decoding Dictionary - a binary trie

=> left branch: 0

right branch: 1



扫描全能王 创建

encoding :

- From tree, build dictionary of letter-codes: $\Sigma \rightarrow \{0, 1\}^*$

L : 000

N : 001

A : 01

O : 100

E : 101

T : 11

ANANT: 010010000100111

* Output: is a stream of bits!

* No delimiters

need 2 bit per char

In a fixed length approach

Compression Ratio "LOSSLESS"

S = LOSSLESS

C = 0100110100111

$$\Sigma_s = \{E, O, L, S\} \quad |\Sigma_s| = 4 \quad \log_2 4 = 2$$

$$\Sigma_c = \{0, 1\} \quad |\Sigma_c| = 2 \quad \log_2 2 = 1$$

$$\frac{|C| \cdot \log |\Sigma_c|}{|S| \cdot \log |\Sigma_s|} = \frac{14 \cdot 1}{8 \cdot 2} = \frac{14}{16} = 88\% \quad * \text{English Text - reduced by up to } 50\%$$

RLE:

Binary length of k, $\text{len}(k) = \lfloor \lg k \rfloor + 1$

encode $\text{len}(k) - 1 \Rightarrow \lfloor \lg k \rfloor$

Number k

1

$\lfloor \lg k \rfloor$

0

Binary rep of k.

1

RLG 1

3

$\lfloor \lg 3 \rfloor = 1$

11

$\Rightarrow 011$

5

$\lfloor \lg 5 \rfloor = 2$

101

$\Rightarrow 00101$

23

$\lfloor \lg 23 \rfloor = 4$

10111

$\Rightarrow 000010111$



扫描全能王 创建

Indicate what first bit is : 0 or 1 , then if alternates runs

... 00 100 ... = 4

#100 Expansion: encoded version is bigger than the source run.

11.2/

• Compel - Ziv - Welch LZW.

- fixed length cube
- typically 12-bits $2^{12} = 4096$ possible patterns can be encoded.
- Dictionary will have 2^{12} entries

Dictionary

- Adaptive - non static
- Start with an initial dictionary D_0 such as all ASCII values ~fixed codes
- will add new entries as we encode $D_0 \Rightarrow D_1 \Rightarrow D_2 \Rightarrow \dots$
- As we find common patterns in the given input, add them to dictionary
- Remember, after we encode the input, it will also need to be decoded.
- Dictionary will be rebuilt by decoder

Advantages

- pass in add patterns to dictionary as we are encoding
 - streaming, we encode on the fly
 - for later "text" we will hopefully have better codes to use - get better compression.



扫描全能王 创建

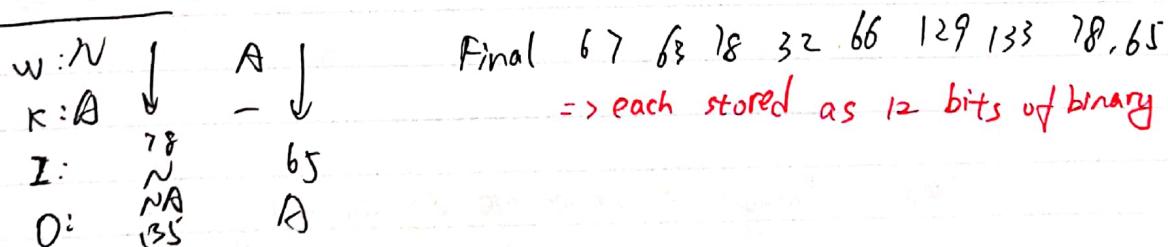
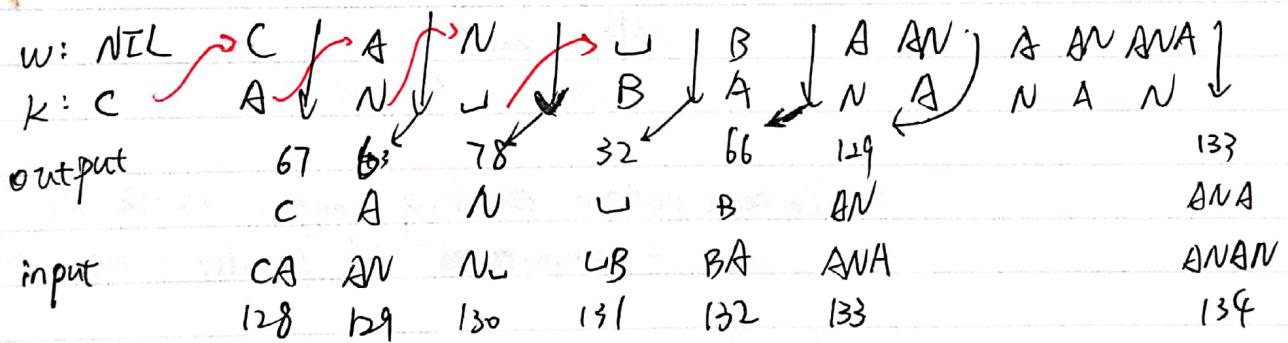
- for large files, space for encoder ~~and~~ decoder is small
 - typically better compression than Huffman

Disadvantage:

- early text of input will not have much compression
 - generally better than ASCII codes

eg. from 22: 1010101010101010
Later 1010101010101010 You
stored as 12-bits 79 89 33 still(12-bits=) 128 132
of binary

CAN- BANANA



Decorating:

- Decoder is "one step behind."
 - Encoder added w with look forward ahead char k but decoder doesn't know k until one step later



扫描全能王 创建

LZW

Decoding • any dictionary implementation will do

Encoding • need to find longest prefix of s in D

• trie insert & ~~top~~ lookup in time proportional to length of string read from input

What if Dictionary fills up?

1. Stop adding
 - fast
 - bad if data changes structure
2. Clear and start over with only ASCII
 - temporary poor compression again
3. ~~Discard~~ Discard less frequently used patterns
 - may be expensive
 - ~~find min value~~
store a counter
4. Increase pattern encoding length $12=13$ bits
 - Complicated but possible - new codes are longer

Note: Decoder must do exactly the same thing as encoder.

MTF

• Note : Encoder & Decoder are very similar
⇒ they must ~~not~~ update the Dictionary in some way

n characters : $\lceil \log_2 n \rceil$ bits to represent n characters
if alphabet is ASCII , need 7-bit for each # in C
⇒ same as each letter for c.



扫描全能王 创建

prefix-free integer encoding

- good because small integers take fewer bits than large integers.

Huffman Coding

- good because we might expect small integers to appear more frequently.

Burrows-Wheeler Transform

- permutes input, but is reversible
- cyclic Shift

thequick...dog\$ ← end-of-word character.
hequick...dog\$t
euquick...dog\$th
uquick...dog\$the
:

* assume this has lowest ASCII value
So, '\$' < 'u'

lexicographic Order: compare characters, if tie, check next one.

eg. A = a a a b => tie, tie, tie, 'b' < 'c'
B = a a a c A comes before B.

C = asff\$fuelllaaata

eg: S = bubba\$

bubba\$

ubbab\$

bba\$bu

ba\$bub

a\$bubbb

① List all cyclic shifts

=>

\$ b u b b a
a \$ b u b b
b a \$ b u b ...
b b \$ b u
b u b b a \$
u b b a \$ b

② Sort lexicographically

C = abbu\$b.

③ extract last character of each row



扫描全能王 创建

Decode: what do we know? How to fill column 1

0	1	2	3	4	5	6	7	8	9
\$		a		a	\$		b		
a		b		b	a		\$		
b		b		b	b		a		
\$		u		u	b		b		
b		\$		\$	b		u		
u		b		b	u		b		

Column 2?

0	0	1	2	3	4	5	6	7	8	9
a	\$	b			\$	b	u			
b	\$	a	b		a	\$	b			
b	b	a			b	a	\$			
u	b	b			b	b	a			
\$	b	u			b	u	b			
b	u	b			u	b	\$			

0	1	2	3	4	5	6	7	8	9
\$	b	u							
a	\$	b							
b	a	\$							
b	b	a							
u	b	b							
\$	b	u							
b	u	b							

Decode Alg:

$j = \text{index of } \$ \text{ in } c$

$s = \text{empty string}$

$j = N[j]$

while $c[j] \neq \$$

append $c[j]$ to s

$j = N[j]$

$c = abbub\$b$

0 1 2 3 4 5

$A = (a, 0)$

$(b, 1)$

$(b, 2) \Rightarrow$

$(bu, 3) \text{ stable}$

$(\$, 4) \text{ sre}$

$(b, 5)$

$\$, 4$

$a, 0$

$b, 1$

$b, 2$

$b, 3$

$u, 4$

N

0 4

1 0

2 1

3 2

4 5

5 3

$j = 4 \bar{j} 3 2 1 0 4$

$s = \epsilon b a b b a \$$



扫描全能王 创建

$\text{Ex 2 } \& s = \text{banana\$}$

banana\\$ \\$banana\\$
 ana\\$ba ana\\$ba
 ana\\$ba ana\\$ba c = annb\\$aa
 ana\\$ba a\\$banan
 na\\$bana bana\\$na
 a\\$banan ana\\$ba
 \\$banana na\\$bana

Decide: $c = \text{annb\$aa}$

0 1 2 3 4 5 6

$A =$	a_0	$\$04$	$j = 43625104$
	n_1	a_{10}	
	n_2	a_{25}	$s = \$\text{banana\$}$
	b_3	a_{36}	
	$\$4$ stable	b_{43}	
	a_5 sort	n_{51}	
	a_6	n_{62}	

Decide $C = \text{nkoeopm\$}$

2017. 11. 28

photo in ~~my~~ phone.

2-3 trees

- each non-leaf node: 2 or 3 children (1 or 2 KVP)
- all leaves are at the same level

Insert

- find lowest internal node
- if overflow, promote middle

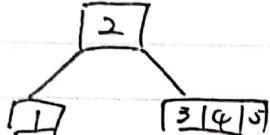
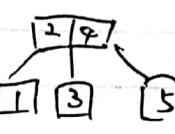
insert 1 2



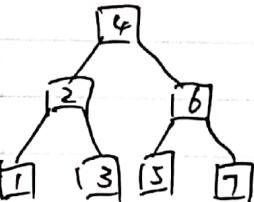
扫描全能王 创建

Insert 3:  \Rightarrow promote middle \Rightarrow 
 NIL NIL NIL NIL \leftarrow we don't continue
 to show this but remember
 the links are there.

Insert 4:  \Rightarrow height of tree grows
 iff root overflows.

5:  \Rightarrow promote middle \Rightarrow 

6. 

7.  \Rightarrow  \Rightarrow 

Delete

- Swap with inorder successor and delete
 - if no underflow, done
 - else first try transfer from immediate sibling
 - If transfer not possible, merge.

tree decrease in height iff root node underflow.



扫描全能王 创建

2017. 11. 30

B-tree of order M

- each node has at most $M-1$ KVP (M-children)
- each non-root node has at least $\lceil M/2 \rceil - 1$ KVPs ($\lceil M/2 \rceil$ children)



Suppose Hard drive sector has size 200 bytes

- KVP is 12 bytes
- child pointer (64-bit OS) 8 bytes

How many KVPs per sector?

~~m-1 < m~~

$$\frac{m}{(m-1)} \cdot 12 + (m+1) \cdot 8 \leq 200$$

find m

$$m = 9 \text{ KVPs}$$

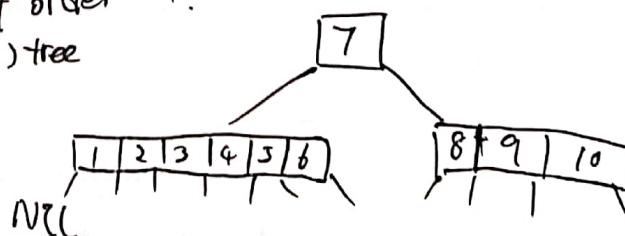
or 10 children

(5, 10) tree of order 10

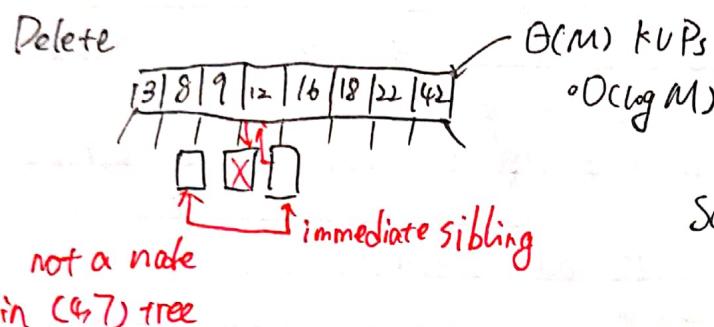
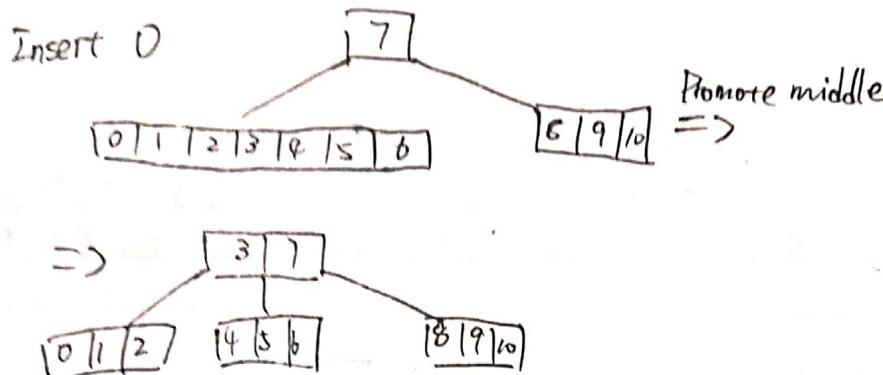
B^+
 B^*

- Search is easy
- insert may involve splitting \Rightarrow might increase height
- delete may involve transfer from immediate sibling or merge (possibly repeated all the way up to root)
 \Rightarrow may decrease the height

B-tree of order 7.
(4, 7) tree



扫描全能王 创建



$$\text{Search} = \left(\frac{\log n}{\log M} \cdot (\log M) \right) = O(\log n)$$

Extendible hashing

- hash table doesn't fit in RAM → store on hard drive
- Which hashing method might be Okay?
- cuckoo & double hashing jump all over the hash table
⇒ each one could be a hard drive access
 - none so far

Idea: want to fill up entire sector

- similar to B-tree
- similar to radix-sort - use bits to split up the data into tables.

Use a tree of height 1

⇒ directory stored in RAM

- root node

- contains hash table of size 2^d "order d"



扫描全能王 创建

\Rightarrow each entry of ~~directo~~ directory points to a block on HD

Hash function $h: \text{key} \rightarrow \text{hash value } \{0, 1, \dots, 2^L - 1\}$
for some $L > 0$.

- represent hash value in binary as a bit string
- use part of hash value to look up in directory of to determine the block on HD where key is stored

\Rightarrow "extendible": grows and shrinks as needed.

Block on HD

- each block contains at most S items • stored by hash value.

S - Size of block # items per block
 d - order of directory
 K_B - local depth of block B

- Directory stored in RAM
- Block stored in HD
- Hash values: $\{0, 1, \dots, 2^L - 1\}$
bit strings of length L .

Idea: Store keys: 00001, 01001, 01110.

Directory of order 2 K_B
 $d=2$ $0 \ 0 \longrightarrow [2 \ 00001]$ $L=5$ # bits in hash value
 $2^d=4$ $0 \ 1 \longrightarrow [2 \ 01001]$ $S=2$ # items per block
 $1 \ 0 \longrightarrow [\text{empty}]$ • block stores $\leq S$ items.
 $1 \ 1$

↑ pointers to blocks

• pointers may go to the same block

• directory has 2^d pointers to blocks

Local block depth

• "local depth" of each block B is $K_B \leq d$

• n_B is # items currently stored in block B • $n_B \leq S$



扫描全能王 创建

Search

- Hash (k)
 - use first d bits to look up in directory
 - follow pointer to corresponding block B
- Load block into RAM
- Search B for matching key (hash(k))
 - ↑ eg. binary search

Ex $L=5, S=2$ 2 entities per block allowed

Insert 01001

$$d=0$$

$$2^0 = 1 \quad [] - [\underset{k_B}{0} \quad 01001]$$

• Order 0

• One block

$$\cdot k_B = 0$$

Insert 00001

$$[] - [0 \underset{0}{\overset{k_B}{00001}}]$$

Insert

Insert 01110

• no room \Rightarrow

• But $k_B = d \Rightarrow$ double ~~hashing~~ directory size

$$d=1 \quad 0 - [\underset{1}{\overset{k_B}{00001}} \quad 01001]$$

$$2^1 = 2 \quad 1 - [\underset{1}{\overset{empty}{}}]$$

where does 01110 go?

• no room and $k_B=d$

\Rightarrow split block and double directory

$$d=2 \quad 0 \quad 0 - [2 \quad 00001]$$

$$0 \quad 1 - [2 \quad 01001]$$

$$1 \quad 0 \quad \nearrow [\underset{1}{\overset{empty}{}}]$$

$$1 \quad 1 \quad \nearrow [\underset{1}{\overset{empty}{}}]$$

Exercise: I



扫描全能王 创建