

# UNION BY SIZE ANALYSIS

ARUN KONA GURTHU  
FIT3155

Using union-by-size, any root node's height has the following relationship with the size (number of elements) in/of the tree:

$$\begin{array}{c} \text{root node} \leftarrow \text{height}[r] \\ \text{Size}[r] \geq 2 \end{array}$$

This can be proved inductively:

Base case:  $\text{Size}[r] = 1, \text{height}[r] = 0;$   
 $\Rightarrow \text{height}[r] = 0 \Rightarrow \text{Size}[r] = 1$

Singleton tree (base case)

$\text{Size} = 1$   
 $\text{height} = 0$

Inductive case:

Assume statement is true for/after 'k' union operations.

$\text{Size}[r] \geq 2^{\text{height}[r]}$  (assumption)

Let another tree rooted at 's' is unioned/merged with the current tree rooted at 'r'.

Case 1:  $\text{height}[r] > \text{height}[s]$



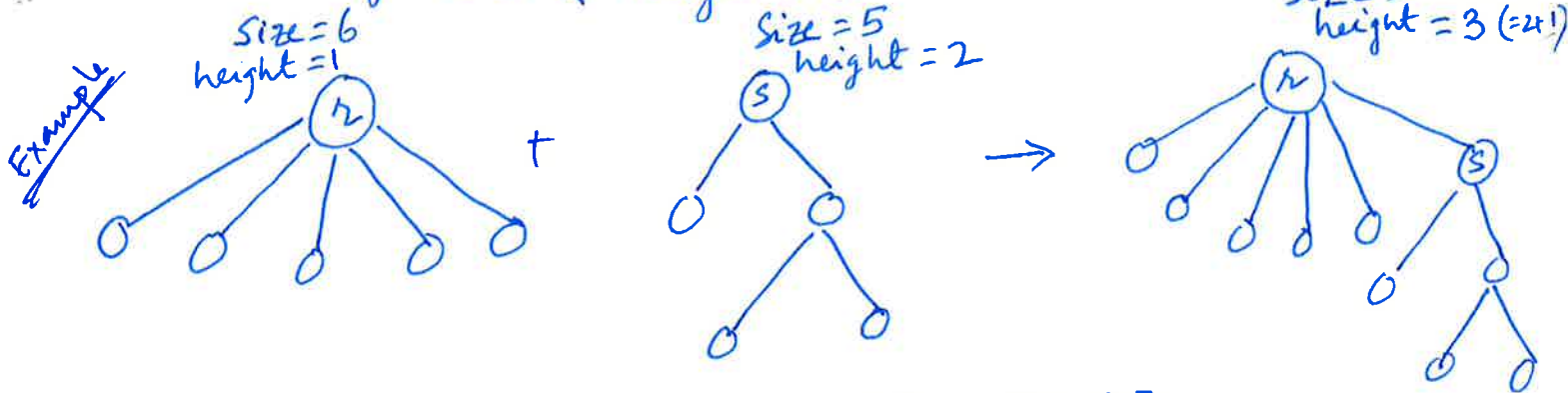
$$\begin{aligned} (\text{new}) \text{Size}[r] &\geq (\text{old}) \text{Size}[r] \\ &\geq 2^{(\text{old}) \text{height}[r]} \\ &\geq 2 \end{aligned}$$

(Inductive step)

But  $(\text{new}) \text{height}[r] = (\text{old}) \text{height}[r]$

$\Rightarrow (\text{new}) \text{Size}[r] \geq 2^{(\text{new}) \text{height}[r]}$

Case 2:  $\text{height}[r] \leq \text{height}[s]$



$$(\text{new}) \text{size}[r] = (\text{old}) \text{size}[r] + \text{size}[s]$$

$$\geq 2 \times \text{size}[s]$$

$$\geq 2 \times 2^{\text{height}[s]} \quad (\text{Inductive step})$$

$$\geq 2^{\text{height}[s]+1} = 2^{(\text{new}) \text{height}[r]}$$

$$\Rightarrow (\text{new}) \text{size}[r] \geq 2$$



So, this proves that for any tree under Union-by-Size

$\text{size}[r] \geq 2^{\text{height}[r]}$   
 in worst-case  
 But  $\text{size}[r]$  is bounded by 'n' (Total number of elems)  
 $n \geq 2^{\text{height}[r]}$  (worst case)

$$\log_2 n \geq \text{height}[r]$$

$$\Rightarrow \text{Find (any element)} \text{ is } O(\text{height}[r]) = O(\log_2 n)$$

$$\text{Union (two elements)} \text{ involves } 2 \text{ Find operations} + O(1) \text{ effort} \\ = O(\log_2 n)$$

# Union by rank (height) Analysis

(WITHOUT PATH COMPRESSION)

## OBSERVATION 1:

For any element 'x' not = root node,  
 $\text{rank}[x] < \text{rank}[\text{parent}[x]]$

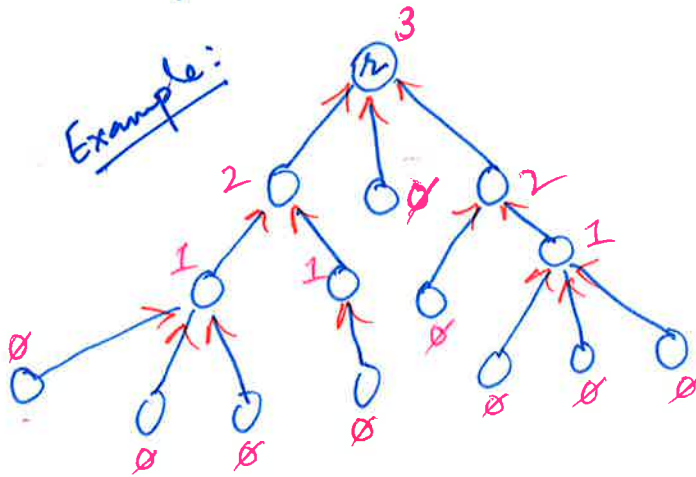
## OBSERVATION 2:

For any element 'x' not = root,  $\text{rank}[x]$  does NOT change under further union operations.

## OBSERVATION 3:

- For <sup>any</sup> path from 'x' to root, the ranks of nodes along that path is a strictly increasing sequence.

Example:



Lemma 1: Any node with rank 'k' has  $\geq 2^k$  nodes in its (sub)tree.

Inductive Proof:

Base case: when  $k=0$  (tree with singleton (root) node)  
 $2^0 = 1$  (TRUE).

Inductive case: Assume true for some rank  $k-1$ .

A node of rank  $k$  is created only when merging 2 trees with rank  $k-1$  each.

But size of each rank  $= k-1$  tree has  $\geq 2^{k-1}$  nodes (inductive step)

$\Rightarrow$  size of (rank  $= k$ ) tree  $\geq 2 \times 2^{k-1} = 2^k$



COROLLARY 1: From OBSERVATION 1 and LEMMA 1, the highest rank for any tree (under union-by-rank) operations is  $\leq \log_2 n$  (where  $n$  is the total # of elems)

Lemma 2 : (AKA "rank lemma")

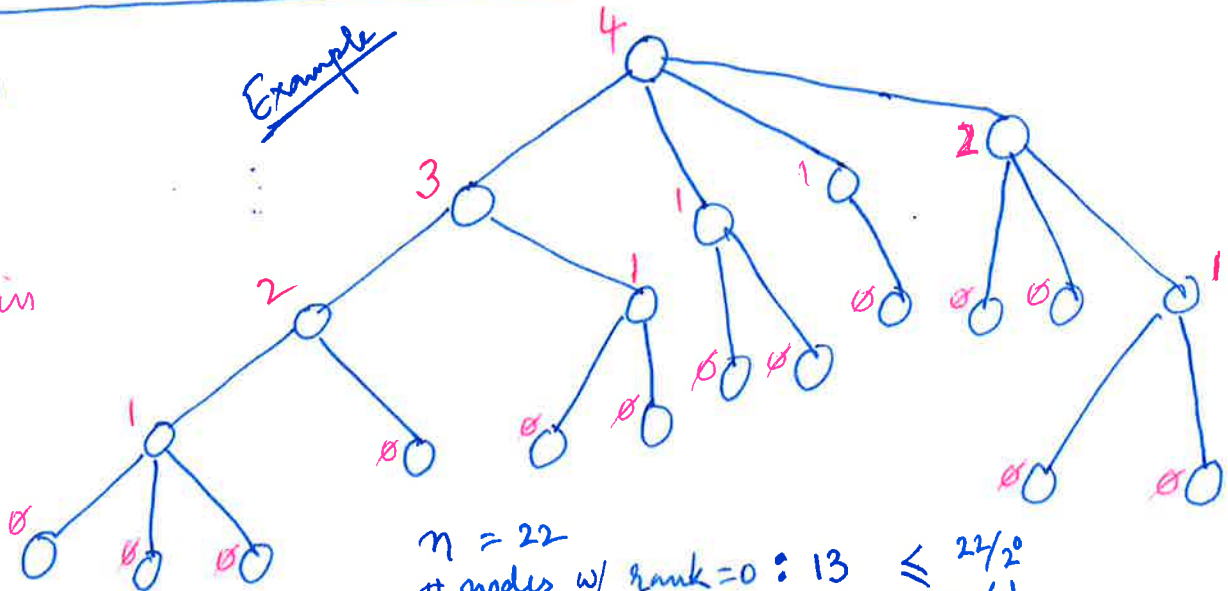
Lemma 2: For any rank  $k \geq 0$ , there are  $\leq \frac{n}{2^k}$  nodes in the tree with rank  $= k$ .

Proof:

① From Lemma 1, <sup>each</sup> ~~any~~ element <sub>n</sub> <sup>in the tree</sup> with rank = k has  $\geq 2^k$  elements in its (sub)tree.

② Different <sup>elements/nodes</sup> ~~nodes~~ with same rank cannot have common descendants (from observation 1)

Q. From the above 2 points, if we only have a total of 'n' elements, and each rank-k element has AT LEAST  $2^k$  elements (from ①), and they do not share common descendants (from ②), then it follows that there are AT MOST  $\frac{n}{2^k}$  such elements/nodes



$n = 22$

# nodes	w/	rank = 0 :	13	$\leq$	$22/2^0$
# nodes	w/	rank = 1 :	5	$\leq$	$22/2^1$
# nodes	w/	rank = 2 :	2	$\leq$	$22/2^2$
# nodes	w/	rank = 3 :	1	$\leq$	$22/2^3$
# nodes	w/	rank = 4 :	1	$\leq$	$22/2^4$

Theorem

## Hopcroft - Ullman's Theorem

Statement: Starting from a fully disjoint set of  $n$  elements and for any sequence of  $m \geq n$  Find/Union operations using ① Union by rank/height & ② path compression, the complexity for this sequence of  $m$  operations is  $O(m \log^* n)$ .

Before proving:

→ Iterated Logarithm function.

Definition of Iterated Log function ( $\log^*$ )

$$\log^* n = \begin{cases} 0 & \text{if } n \leq 1 \\ 1 + \log^*(\log n) & \text{otherwise.} \end{cases}$$

By this definition

$$n=1, \log^*(1) = 0$$

$$n=2, \log^*(2) = 1$$

$$n \in [3, 4], \log^*(n) = 2$$

$$n \in [5, 16], \log^*(n) = 3$$

$$n \in [17, 65536], \log^*(n) = 4$$

$$n \in [65537, 2^{65536}], \log^*(n) = 5$$

SIDE NOTE: # of atoms in the universe estimated  $< 2^{270}$  :)

BOTTOM LINE:  $\log^*(n)$  is an extremely slowly growing function.



# Hopcroft - Ullman's theorem proof (union-by-rank with path compression)

① Divide up ranks into the following "RANK BLOCKS"

Block 1: Rank  $k = \{1\}$

Block 2: Rank  $k = \{2\}$

Block 3: Ranks  $k = \{3, 4\}$

Block 4: Ranks  $k = \{5, 6, \dots, 16\}$

Block 5: Ranks  $k = \{17, 18, \dots, 2^{16}\}$

Block 6: Ranks  $k = \{2^{16}+1, 2^{16}+2, \dots, 2^{65536}\}$

$\vdots$

There are only  $\log^*(n)$  possible rank blocks for a given  $n$

② The same properties/observations that were discussed for union-by-rank (without path compression; see pages 3-4) also hold here (with path compression)

OBSERVATION 1: If node  $x \neq \text{root}$ ,  $\text{rank}[x] < \text{rank}[\text{parent}[x]]$

OBSERVATION 2: If node  $x \neq \text{root}$ ,  $\text{rank}[x]$  does NOT change.

OBSERVATION 3: Ranks along any path from  $x$  to root, follows a strictly increasing sequence

OBSERVATION 4: Any node with rank  $= k$  has  $\geq 2^k$  nodes in its (sub)tree.

OBSERVATION 5: Highest rank of a node is  $\leq \log_2 n$

OBSERVATION 6: For any integer  $k \geq 0$ , there are  $\leq \frac{n}{2^k}$  nodes with rank  $k$ .

③ Running Time of Find/Union is bounded by number of parent pointers traversed/followed. Two cases arise:

Case 1:  $\text{rank}(\text{parent}[x])$  is in a higher RANK BLOCK than  $\text{rank}(x)$

Case 2:  $\text{rank}(\text{parent}[x])$  is in the same RANK BLOCK as  $\text{rank}(x)$ .

Case 1: rank[parent(x)] is in a bigger rank block than rank[x]

Since there are at most  $\log^*(n)$  rank blocks, there can only be  $O(\log^*(n))$  nodes that satisfy this property. For  $m$  find/union operations, this becomes  $O(m \log^* n)$ .

Case 2: rank[parent(x)] is in the same rank block as rank[x]

↳ Fix some arbitrary rank block:  $\{2^k, 2^{k+1}, \dots, 2^{k+2}\}$

↳ The number of times a node in this rank block gets visited before its parent goes into a different/next rank block (over union operations) is  $\leq 2^k$

↳ But, using "rank lemma" (Page 4), the total # of elements/nodes in the tree whose rank lies in this rank block is bounded by:

$$\frac{n}{2^{k+1}} + \frac{n}{2^{k+2}} + \frac{n}{2^{k+3}} + \dots + \frac{n}{2^k} \leq \frac{n}{2^k}$$

↳  $\Rightarrow$  Total visits to nodes whose ranks lie in this rank block is  $\leq \frac{n}{2^k} \times 2^k \Rightarrow \leq n$

↳ But we only have at most  $O(\log^* n)$  rank blocks

↳ Total work for this case  $O(\underline{n} \log^* n)$

SUMMING UP

Total work covering cases 1 & 2 above is:

$$O(m \log^* n) + O(n \log^* n)$$

[note  $m \geq n$  in this theorem]

$$= O(m \log^* n)$$



NOT EXAMINABLE (but worth knowing)

Tarjan (another brilliant contributor to Computer Science) showed that, Hopcraft-Ullman's bound is in fact loose, and that it can be shown that any 'm' sequence of find/union operation grows as

~~$O(m \log n)$~~   $O(m \alpha(n))$

where  $\alpha(n)$  is the Inverse Ackerman function

Compare this with Hopcraft-Ullman's bound of

$$O(m \log^*(n))$$

Iterated log function

↳ It turns out that

$\alpha(n)$  grows excruciatingly more slowly than  $\log^*(n)$ ...

↳ ... and, in fact,  $O(m \alpha(n))$  is the optimal bound for disjoint-set data structures, amortised over 'm' find/union operations.

END OF ANALYSIS