

# Union-Find data structures

Maintain a collection of disjoint sets over  $[n]$ ,

supporting

$\{1, \dots, n\}$

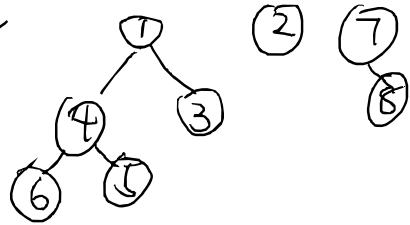
- $\text{union}(x, y)$ : merge sets containing  $x$  and  $y$
- $\text{find}(x)$ : return the representative of set containing  $x$ .

Initially, all sets are singletons  $\{i\}$ .

Recap: maintain a tree for each set, use the root as the representative

$\text{find}(x)$ :

keep going up,  
return the root.



$\text{union}(x, y)$

$r_x = \text{find}(x), r_y = \text{find}(y)$

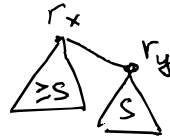
if  $r_x \neq r_y$

$\text{parent}[r_x] = r_y$

Union-by-size:

maintain the size of the set at root,

when union two trees, set parent [smaller tree] to root of bigger tree



Union-by-size guarantees  $\text{depth} \leq \log n$ .

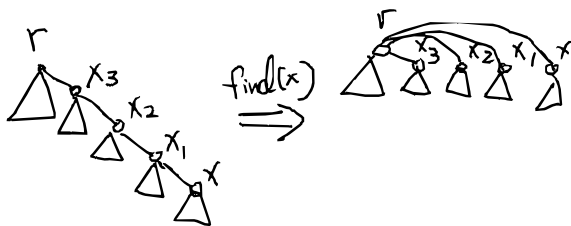
- going up one level (at least) doubles the size.

each operation takes  $O(\log n)$  time

Path compression:

$\text{find}(x)$ : keep going up, set the parent of every node visited on the way to the root.

return root.



amortized  $O(\log n)$

Union-by-size + path compression:

amortized  $O(\alpha(n))$  time

↑ inverse Ackermann function

Ackermann function

$$A(1, n) = n + 2 \quad (+2)$$

$$A(m, 1) = 2$$

$$A(m+1, n+1) = A(m, A(m+1, n))$$

$$A(2, n) = A(1, A(2, n-1)) = A(2, n-1) + 2 = 2 \cdot n \quad (\times 2)$$

$$A(3, n) = A(2, A(3, n-1)) = 2 \times A(3, n-1) = 2^n \quad (\text{power})$$

$$A(4, n) = 2^{2^{\dots^2}} \Big\}^n \quad A(5, n) = 2^{2^{2^{\dots^2}}} \Big\}^n$$

$A(m, n)$  is very fast growing even for  $m=4$ .

$$\alpha(n) = \min_{k \geq 1} \{k: A(k, k) \geq n\} \ll \log n$$

$$\ll \log^* n$$

only path compression, no union-by-size.

Prove: Amortized time of union-find is  $O(\log n)$ .

Consider a sequence of  $m$  union-find operations.

Connecting roots takes  $O(1)$  time, thus, focus on Find!

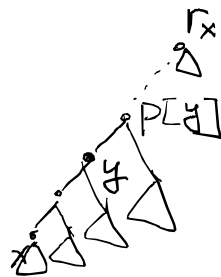
Let  $S(y)$  be the size of the subtree rooted at  $y$ .

Consider a find( $x$ ) operation

- visits all nodes on  $x \rightarrow r_x$  path.

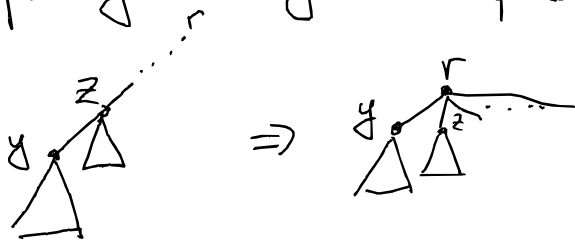
Call a node  $y$  on the path "light"

- if  $S(y) \leq \frac{1}{2} S(p[y])$ .



Claim: There are at most  $\log n$  light  $y$  on this path.

Suppose  $y$  is heavy and  $z = p[y]$  before path-compression.  
( $z$  is not the root)



In this case,  $S(z)$  drop by a factor of two!

Claim: If a node  $z$  is not a root, then  $z$  cannot become a root, and  $S(z)$  cannot increase.

Thus, every time we visit a heavy  $y$  s.t.

$z = P[y]$  is not the root,

$S(z)$  drops by a factor of two. This can happen to

$z$  at most  $\log n$  times.

During the  $O(m)$  Find ops,

- visit  $O(m \log n)$  light vertices

- visit  $O(m \log n)$  heavy  $y$ , s.t.  $P[y]$  is not the root.

- visit  $O(m)$  heavy  $y$ , s.t.  $P[y]$  is the root.

only  $O(m)$  nodes  $z$  can be part of a tree that is not a singleton.

Thus, all Find visits  $O(m \log n)$  nodes.

Amortized time is  $O(\log n)$ .

Path compression + union-by-size.

Prove: amortized time  $\leq O(\log^+ n)$  (weaker bound)

$$\begin{aligned} & \min_k \{k: \underbrace{\log \log \dots \log n}_k \leq 1\} \\ & = \min_k \{A(4, k) \geq n\}. \end{aligned}$$

Consider  $O(m)$  operations.

For a node  $y$ , define

$\bar{S}(y) :=$  max subtree size rooted at  $y$  during the process.

For a nonroot  $y$ ,  $\bar{S}(y) = S(y)$  when  $y$  becomes nonroot

$$2\bar{S}(y) \leq \bar{S}(p[y]) \text{ union-by-size}$$

Partition vertices into  $\log^+ n$  buckets based on  $\bar{S}(y)$ .

bucket  $i$  contains all  $y$  s.t.  $2^i \leq \bar{S}(y) < 2^{i+1}$

$\uparrow$   $B_i$   $\uparrow$   $B_{i+1} = 2^{B_i}$

Claim: #nodes in bucket  $i$  is

$$O\left(\frac{m}{B_i}\right).$$

Consider a node  $y$  visited during a find:

- $y$  or  $p[y]$  is root:  $O(1)$  per op
- $y, p[y]$  in different buckets:  $\log^* n$  per op.
- $y, p[y]$  in same bucket:

$$r \text{ be root, } \bar{S}(r) \geq 2\bar{S}(p[y]).$$

$p[y]$  is set to  $y$

→ thus,  $\bar{S}(p[y])$  increases by  $\times 2$ .  
can happen  $\leq \log_{\frac{B_i+1}{B_i}} \leq B_i$  times to  $y$ .

In total, this case happens

$$\sum_{\text{bucket } i} \sum_{y \in \text{bucket } i} B_i = \sum_{i \leq \log^* n} O\left(\frac{m}{B_i} \cdot B_i\right) = O(m \log^* n).$$