# Algorithm Design and Analysis

## Union-Find (More Amortized Analysis!)

# Reminder

- **Midterm One is next Tuesday at 7:00pm**

- If you 100% can not make this for a legitimate reason, email us or post on Piazza **by the end of today**.

# Roadmap for today

- Design the *Union-Find* data structure for the *disjoint sets problem*

- Practice *potential functions* by analyzing Union-Find

# Motivation: Kruskal's Algorithm

**Review (Minimum Spanning Tree)**: A spanning tree of an undirected graph with the least total (edge) cost of all possible spanning trees

**Review (Kruskal's Algorithm)**: For each edge $(u, v)$ in sorted order by cost, add the edge to the spanning tree if $u$ and $v$ are not connected.

**How do we do that part??**

# The disjoint-sets problem

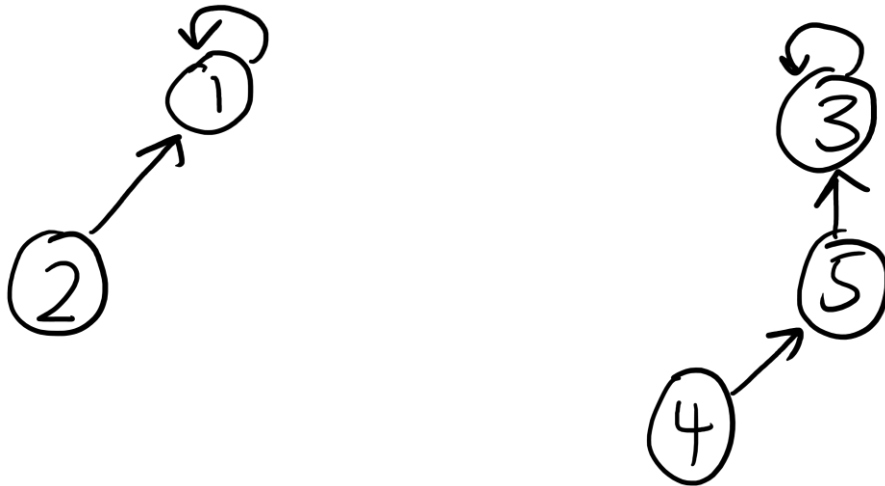*Problem (Disjoint Sets)*: We want to support the following API:

- **MakeSet**$(x)$: Create a set consisting of the single element $\{x\}$
- **Find**$(x)$: Return the *representative element* of the set containing $x$
- **Union**$(x, y)$: Merge the two sets $S_x \ni x$ and $S_y \ni y$ into a single set.

**Simple but inefficient #1:** Maintain a representative for each element. Union loops over every element and updates the ID of the representative: **costs $O(n)$**. Find **costs $O(1)$**.

**Simple but inefficient #2:** Maintain a graph with an adjacency list. Union just adds a new edge: **costs $O(1)$**. But find must search the entire connected component: **costs $O(n)$**.

# The disjoint-set forest data structure

- ***Key idea***: Represent the sets as **trees**.  Use the roots of the trees as the representative element.

- Representation:  Store a parent pointer for each node. Roots have no parent (by convention, $p(x) = x$ for roots).

# Implementation (basic version)

- **MakeSet**$(x)$:
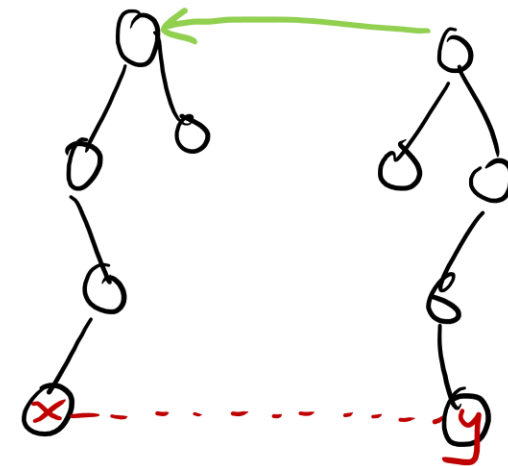
    Create $x$

    Set $p(x) \leftarrow x$

- **Find**$(x)$:

    Walk up parent chain
    until root.

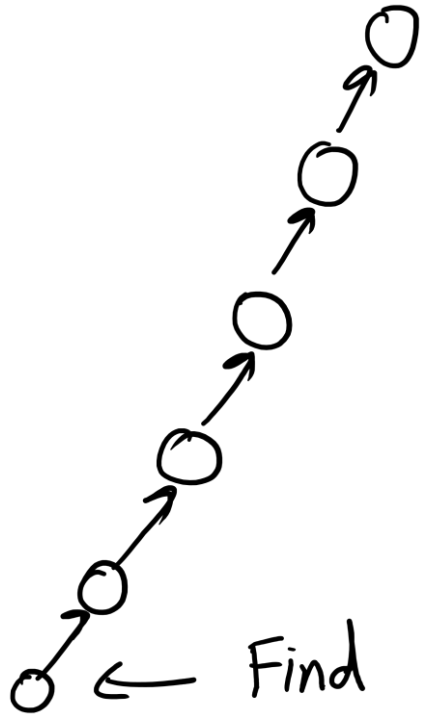- **Union**$(x, y)$:

    Link ( Find (x), Find (y))

- **Link**$(x, y)$

    $P(y) \leftarrow x$

# Performance

Theorem: Let $n$ be the current number of elements in the sets (i.e., the number of MakeSet operations performed so far). There exists inputs for which every find costs $\Theta(n)$.

# Making Union better?

- The bad performance was caused by long chains of nodes…
- Can we just… not do that?

> *Idea (Union-by-size)*: When performing a Union, make the smaller tree a child of the larger tree. If they're the same size, then pick arbitrarily.

- We should store an extra field $s(x)$ that knows the size of the trees
- $s(x)$ is the size of the tree rooted at $x$ (we don't care about non-roots)

# Union-by-size implementation

- **Link$(x, y)$:**

    If $S(x) < S(y)$ then swap$(x, y)$

    $P(y) \leftarrow x$

    $S(x) \leftarrow S(x) + S(y)$

# Performance of union-by-size

**Theorem**: Let $n$ be the current number of elements in the sets. Using *union-by-size*, every Link operation costs $O(1)$ and every Find operation costs $O(\log n)$ in the **worst-case**.
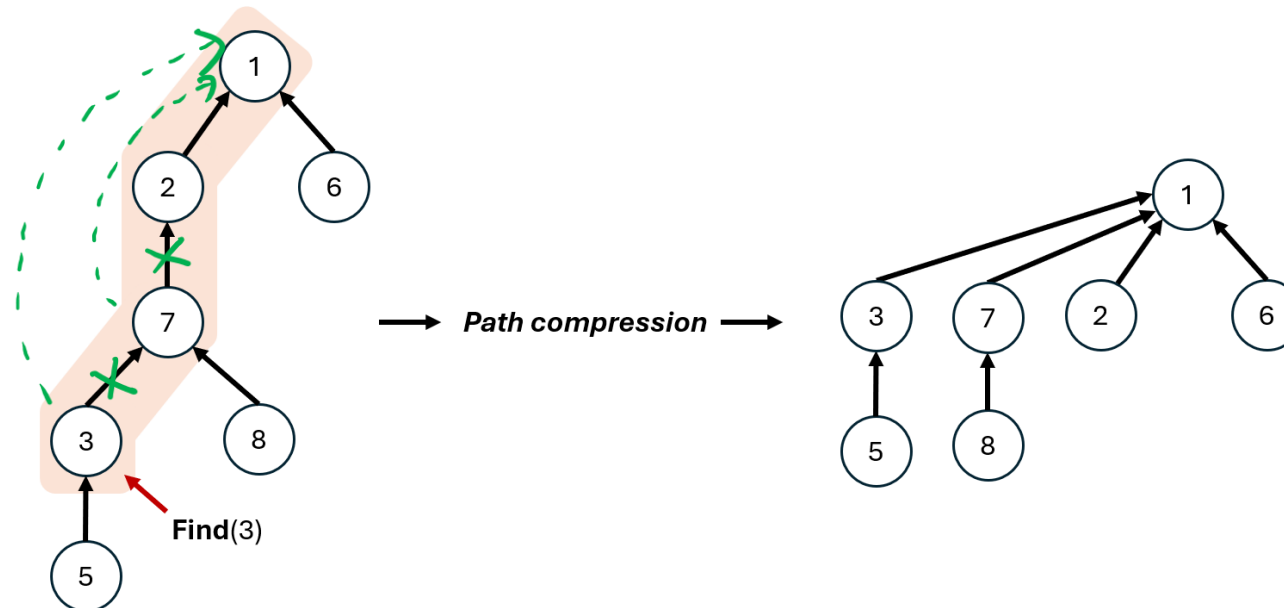


$$\text{Size}(p) > 2 \, \text{size}(u)$$

# Another improvement

- We just made Union better. Can we instead/also make Find better?

> *Idea (Path compression)*: When performing a Find, point every node along the path at its current root/representative element.



Find(3) → Path compression →

# Path compression implementation

- **Find**($x$):

$$\text{if } p(x) \neq x:$$
$$\quad \text{set } p(x) \leftarrow \text{Find}(p(x))$$
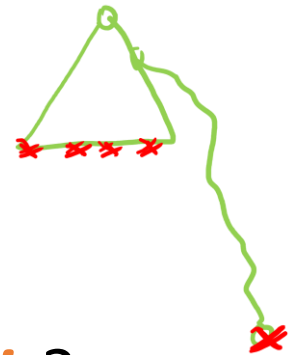$$\text{return } p(x)$$

# Cost model for amortization

- To avoid arbitrary constants in the analysis, we will once again work in a simplified cost model.  All our analyses will be asymptotically valid in the word RAM up to constant factors.

  - **MakeSet** costs 1
  - **Link** costs 1
  - **Find** costs **number of nodes touched**

- **Goal**: *Amortized costs* of $O(\log n)$ for each operation

# Performance of path compression

**Theorem**: Let $n$ be the current number of elements in the sets. Using *path compression* (but not union-by-size*)*, in our cost model, the **amortized cost** of MakeSet is 1, Link is $(1 + \log n)$, and Find is $(2 + \log n)$.

- **Observation**: *Balance* is what matters

- Balanced trees are always fast, imbalanced trees are slow

- How do we measure how balanced a tree is at a *per-node basis*?
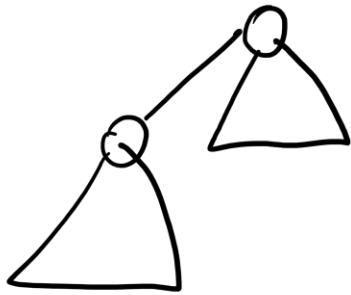
15

# Balanced or imbalanced?

**Definition (heavy/light)**: Given a node $u$ and its parent $p$, call a node:

- **Heavy** if $\text{size}(u) > \frac{1}{2}\,\text{size}(p)$, i.e., $u$ contains a majority of $p$'s descendants

- **Light** if $\text{size}(u) \leq \frac{1}{2}\,\text{size}(p)$, i.e., $u$ contains at most half of $p$'s descendants

- Root is neither heavy nor light (it has no parent)

- In a perfectly balanced tree, every node is light (except root)

- In a chain (the worst-balanced tree), every node is heavy (except root)

*2 leaf.*

# Balanced or imbalanced?

**Lemma (light lemma)**: On any root-to-leaf path in any tree of $n$ nodes, there are at most $\log n$ light nodes.

**Definition (Light):**
$$\text{size}(u) \leq \frac{1}{2}\text{size}(p),$$



Same proof as before!

# Reaching your potential

- Find costs #nodes touched = $(1 + \text{#heavy} + \text{#light})$.

- We know that $\text{#light} \leq \log n$           **(light lemma)**

- So, Find costs at most $(1 + \log n + \text{#heavy})$.

- We want to define a *potential function* that will save up and **pay for the cost of touching the heavy nodes**

$$\mathbf{ac} = \underbrace{1 + \log n + \text{#heavy}}_{\textit{Actual cost}} + \underbrace{\Delta \Phi}_{\textit{Change in potential}}$$

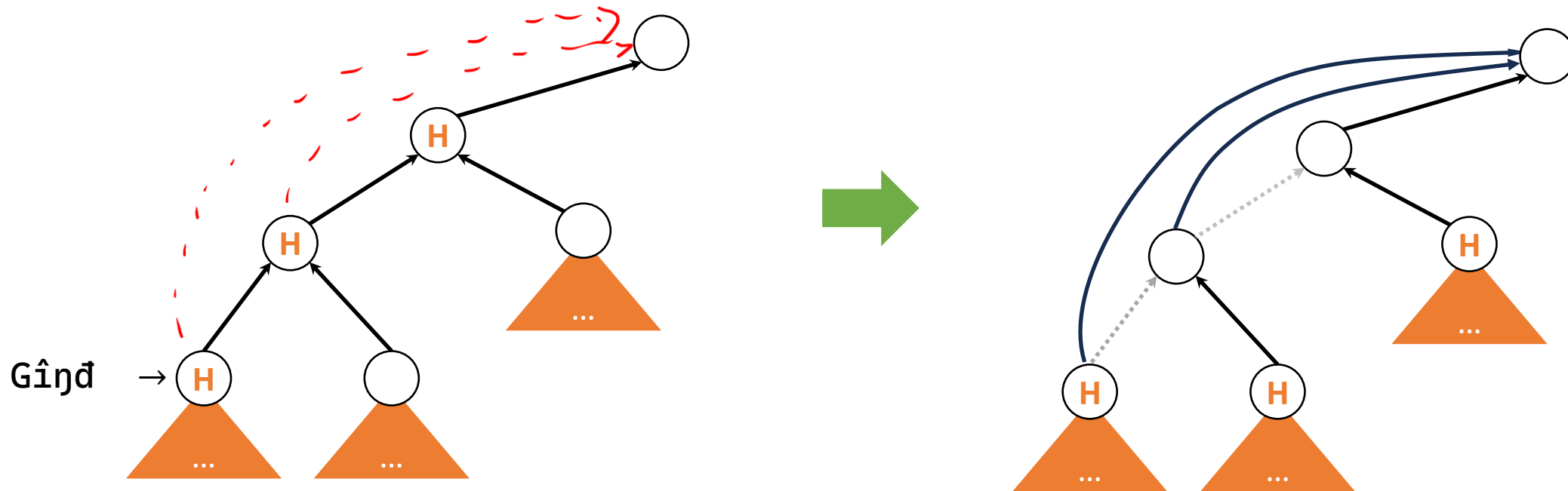*want this to be* $\approx -\text{#heavy}$

18

# A *potential* idea

- **Observation**: A node can have *at most one heavy child*

- *Potential* potential function idea

$$\Phi(F) = \text{\# heavy nodes}$$

- **Problem**: Can we prove that the number of heavy nodes decreases?

- **Exercise***: Draw a scenario where a path compression operation does not decrease the number of heavy nodes.

# Too many heavy nodes!

- *Exercise:* Draw a scenario where a path compression operation does not decrease the number of heavy nodes.

# Refining the potential

- ***Observation***: A node can have ***at most one heavy child* but** moving it does not necessarily reduce the number of heavy children (a sibling may become heavy in its place)

- However, is there a maximum number of times this can happen?

- Moving the heavy child *at least halves* the size of the subtree!

- Therefore, the subtree of node $x$ with size $s(x)$ can have its heavy child moved at most $\underline{\log_2(\text{size}(x))}$ times!

# The balance potential

- Define our potential function to be:

$$\Phi(F) = \sum_{u \in F} \log(\text{size}(u))$$

all nodes $\longrightarrow$

- **Nice properties**:
  - Initially zero (all trees start at size 1)
  - Always non-negative
  - Increases when we perform a Link
  - Decreases when we perform a Find

**no debt** $\qquad \Phi(S_m) \geqslant \Phi(S_0)$

**Links save up \$\$\$ to pay for Finds**

# Analysis of MakeSet

$$\overline{\Phi} = \sum \log(size(u))$$

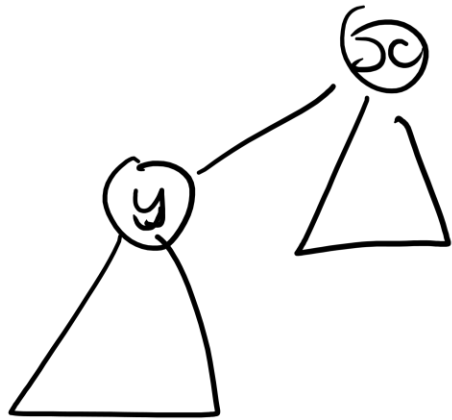Lemma (cost of MakeSet): MakeSet does not change $\Phi(F)$

$$\Delta \Phi = \log(1) = 0$$

Corollary: The MakeSet operation has an amortized cost of 1

# Analysis of Link

$$\overline{\Phi} = \sum log(size(u))$$

$$\Delta \overline{\Phi} = + log(size'(x)) - log(size(x))$$
$$\leq log(n) - log(size(x))$$
$$\leq log(n)$$

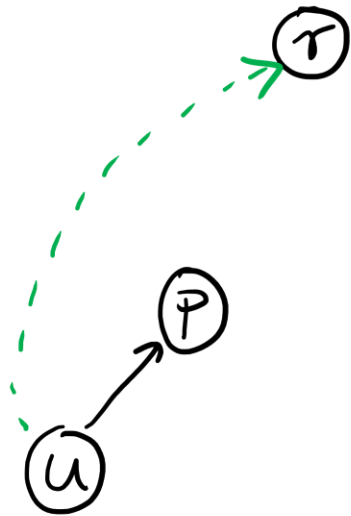$$aC = actual + \Delta \overline{\Phi}$$

Corollary: A link operation has amortized cost at most $1 + \log n$

# Analysis of Find

$$\overline{\Phi} = \sum \log_2(size(u))$$

*Lemma*: A Find operation decreases $\Phi(F)$ by at least #heavy nodes $-1$

- Consider *heavy nodes* $u$ with parent $p$ (other than $r$) on the **Find** path

$$size'(p) = size(p) - size(u)$$

$$< \frac{1}{2} size(p)$$

$$\triangle \overline{\Phi}_{(at\ node\ p)}$$

$$= \log(size'(p)) - \log(size(p))$$

$$< \log\left(\frac{1}{2} size(p)\right) - \log(size(p)) = -1$$

Add this up for all heavy $u$

25

# Analysis of Find

*Corollary (cost of Find)*: Find has amortized cost at most $(2 + \log n)$

$$\text{a.c. FIND} \leq \underbrace{1 + \#\cancel{heavy} + \#light}_{\text{actual cost}} \underbrace{- (\#\cancel{heavy} - 1)}_{\Delta \Phi}$$

$$\leq 2 + \#light$$

$$\leq 2 + \log n$$

26

# Summary of Union-Find

- *Union-Find with union by size*:
  - *Link*: $O(1)$
  - *Find*: $O(\log n)$ **worst-case**

- *Union-Find with path compression*:
  - *Link*: $O(\log n)$ **amortized**
  - *Find*: $O(\log n)$ **amortized**

- *Union-Find with both! (Not proven in this class)*:
  - *Link*: $O(\alpha(n))$ **amortized**
  - *Find*: $O(\alpha(n))$ **amortized**
  - $\Omega(\alpha(n))$ is also a lower bound so this is optimal!

$\alpha(n)$

( Text book, CLRS )