

Disjoint Set Union

Problem: Set union is slow because eliminating duplicates requires us to find the intersection.

Solution: Disallow duplicates

$$S_i \cap S_j = \emptyset$$

Example: Assassin Game

Each player issued a target card with name of another player.

← MakeSet

If you "kill" your target, you take over his cards.

← Union

Last player has all of the cards.

Question: Who has my card?

← Find

Union-Find operations

Make-Set (x): create set $\{x\}$

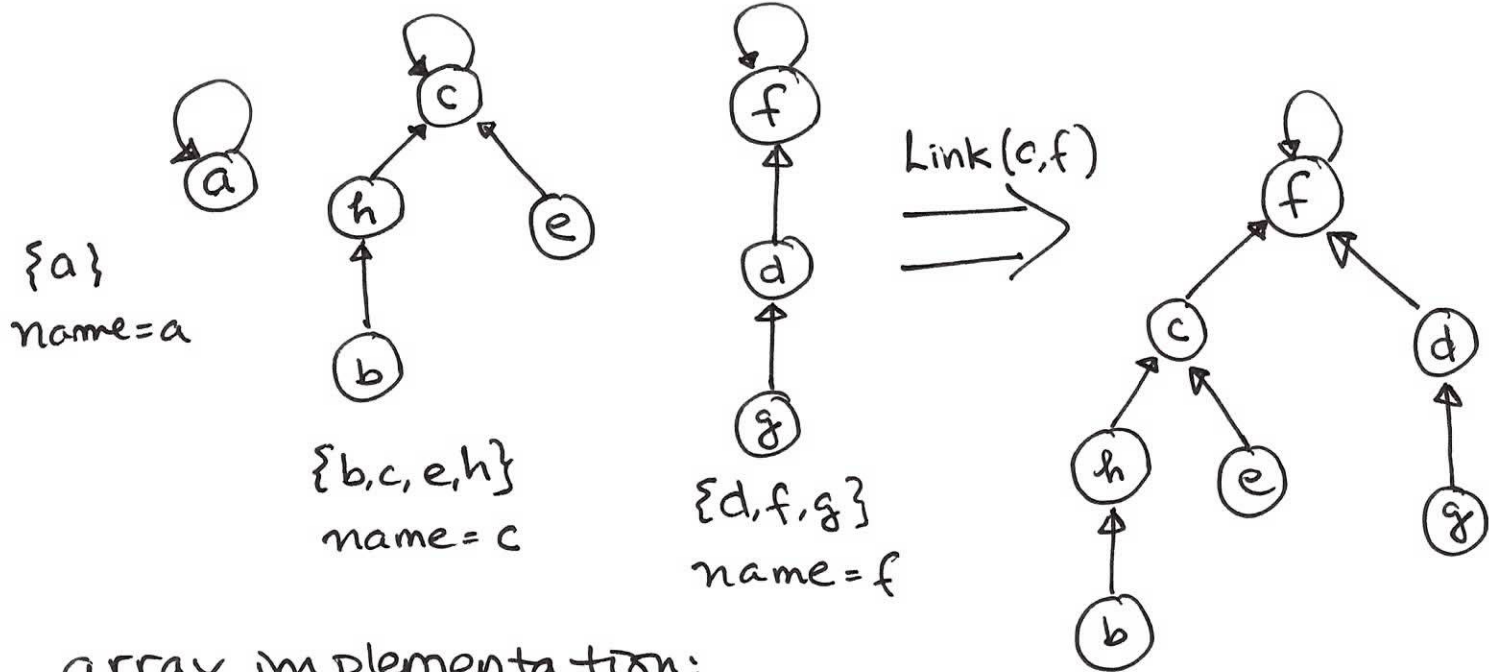
Find-Set (x): return "name" of set
that contains x .

changes over time

Link ($name_1, name_2$): union of sets named
 $name_1$ and $name_2$

Union (x, y): join sets containing x & y .

Representation: parent pointers

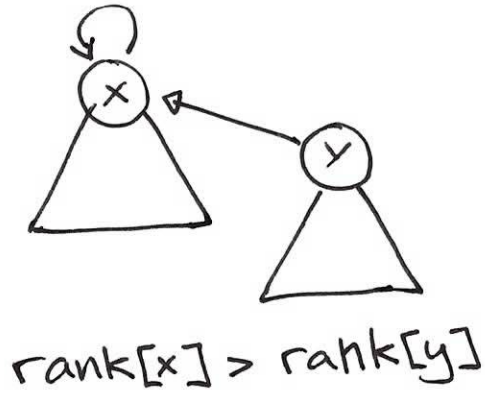


array implementation:

a	b	c	d	e	f	g	h
a	h	c	f	c	f	d	c

Problem: Link(x,y) can create long chains

Solution: Union-by-Rank = higher ranked node is root.



```
Link(x,y) {  
    if rank[x] > rank[y]  
        p[y] = x  
    else {  
        p[x] = y  
        if rank[x] == rank[y]  
            rank[y]++  
    }  
}
```

initial rank = 0

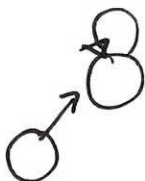
rank increases by 1 when
2 equal ranked nodes are linked.

Smallest trees with rank r .

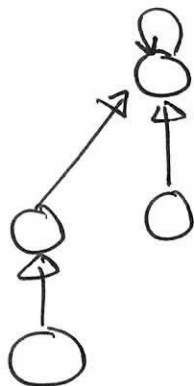
$r=0$



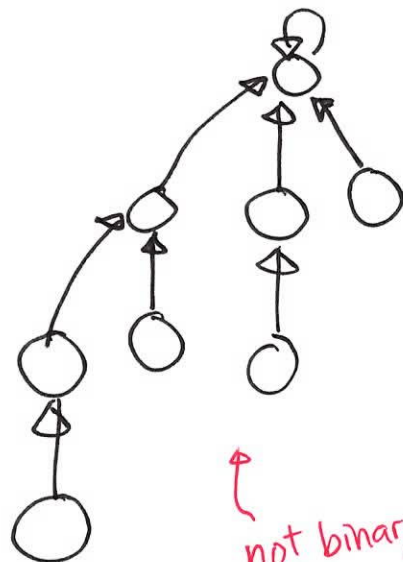
$r=1$



$r=2$



$r=3$



not binary trees

Claim: A node with rank r must have
 $\geq 2^r$ descendants (incl. itself).

Pf: (induction on r)

true for $r=0$.

The only way to get a node with rank $r+1$ is to link 2 nodes with rank r . By induction, these have $\geq 2^r$ descendants each. New root has $2 \cdot 2^r = 2^{r+1}$ descendants.



Claim: A node with rank r has height $\leq r$.

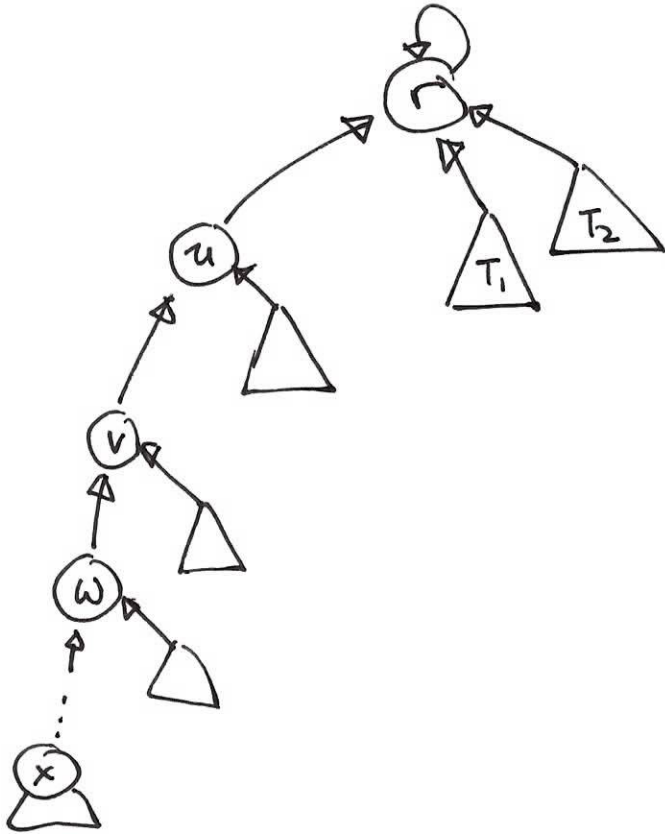
Pf: Easy induction.

Running time of Union-by-Rank:

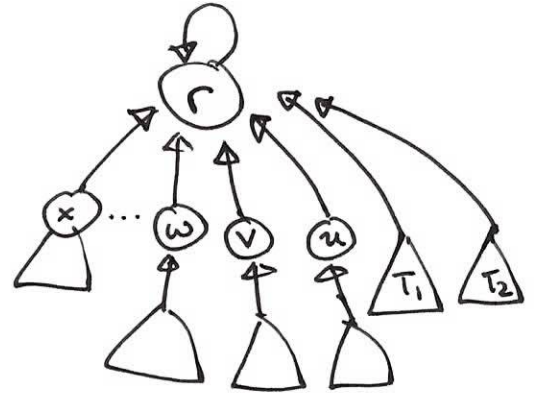
Make-Set	$O(1)$
Find-Set	$O(\log n)$
link	$O(1)$
Union	$O(\log n)$

↖ We can do even better
(Good enough for Kruskal's)

Path Compression



Find-Set(x)
⇒



Find-Set (x) {

$r \leftarrow x$

while ($r \neq p[r]$)

$r \leftarrow p[r]$

/ r is root */*

while ($x \neq r$) {

$temp \leftarrow p[x]$

$p[x] \leftarrow r$

$x \leftarrow temp$

}

return (r)

}

parent
of

Path compression:

- doubles time for Find-Set
- future finds much faster

everybody
point to
root r

Amortized Analysis:

Suppose we have m operations including n Make-Sets.

If we use Union-Find in another algorithm (e.g. Kruskal's), we don't care how long each operation takes. We just want to know the total running time for all m operations.

↙
divide by m gives
per operation
"amortized" running time

With Path Compression and Union-by-Rank, m disjoint set union operations (including n Make-Set operations) takes time $O(m \lg^* n)$.

very slowly growing function.

$$\lg^* n = \min \{ i \geq 0 : \underbrace{\log_2(\log_2(\dots(\log_2 n)\dots))}_{i \text{ times}} \leq 1 \}$$

$$\lg^* 2^{2^{\dots^2}} \}^n = n + 1$$

$$\lg^* 2 = 1$$

$$\lg^* 4 = 2$$

$$\lg^* 16 = 3$$

$$\lg^* 64k = 4$$

$$\lg^* 2^{64k} = 5$$

$$\lg^* 10^{80} < 5$$

of atoms in observable universe

Ackerman's function

$$A(1, j) = 2^j, \quad j \geq 1$$

$$A(i, 1) = A(i-1, 2) \quad i \geq 2$$

$$A(i, j) = A(i-1, A(i, j-1)) \quad i, j \geq 2$$

A	$j=1$	2	3	4
$i=1$	2	2^2	2^3	2^4
$i=2$	2^2	2^{2^2}	$2^{2^{2^2}}$	$2^{2^{2^{2^2}}}$
$i=3$	2^{2^2}	$2^{2^{\dots^2}} \} 16$	$2^{2^{\dots^2}} \} 2^{2^{\dots^2}} \} 16$	$2^{2^{\dots^2}} \} 16$

← stack of j 2's

Ackerman's function (cont'd)

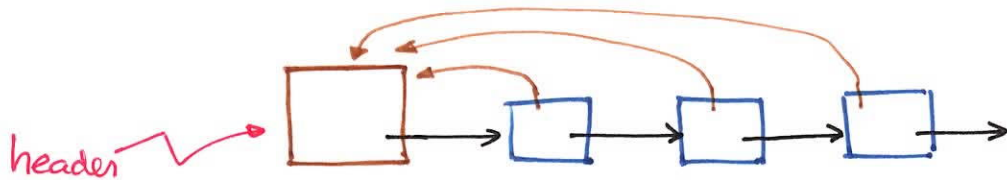
$$\alpha(m, n) = \min \{ i \geq 1 : A(i, \lfloor m/n \rfloor) > \lg n \}$$

$$A(4, 1) > 10^{80} = \# \text{ of atoms in universe}$$

$$\alpha(m, n) < 4 \text{ for all practical purposes}$$

Example of Amortized Analysis of Union-Find.

Linked list with root pointer.



Make-Set $\Theta(1)$

Find-Set $\Theta(1)$

Link $\Theta(n)$

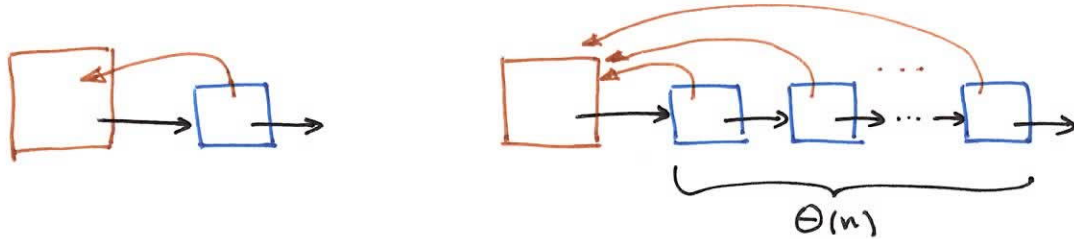
Union $\Theta(n)$

} worst case running time

Link & Union worst case: update root pointers of $\Theta(n)$ nodes

Each Union might take $\Theta(n)$ time.

If we had m unions, how much time does that take?



Idea: make nodes in shorter list point to root of longer list.

- ① Fewer pointer updates
- ② When a node's root pointer is updated, the size of its list at least doubles
- ③ Each node's pointer can be updated at most $\log n$ times. *Size of list cannot exceed n .*
- ④ Worst case Union cannot happen very often.

Amortized Analysis

- Estimate the running time of a sequence of instructions (rather than a single one).
- Often "object-oriented":
Estimate the time spent on each piece of data rather than time spent on each line of code
- Credits for time.

Back to linked lists with root pointers:

- Time for Union is proportional to number of pointer updates.
- m operations including n Make-Sets:
total time for all Unions is $O(n \log n)$.

There are n nodes, each root pointer can be updated at most $\log n$ times.

- Amortized time for Union: $\log n$.

$$\frac{O(n \log n)}{m} \leq \log n$$

Since $m \geq n$

Credits: ↗ Alternate method for amortized analysis

- ① Think prepaid cell phones: use "units" to pay for running time.
- ② Use 1 credit for each pointer update.
- ③ Prepay for $\log n$ credits at each node.
- ④ Can "charge" the credits to other operations. ↗ think parents & cell phones
- ⑤ Charge credits to Make-Set operations:

Make-Set	$\Theta(\log n)$	Link	$\Theta(1)$	} pay for pointer updates using prepaid credits
Find set	$\Theta(1)$	Union	$\Theta(1)$	

Getting $\log \log n$ amortized time...

- Use disjoint set forests, union-by-rank, path compression
- Link is $\Theta(1)$
- Must pay for path compression
 - + Charge some credits to Find-Set
 - + Charge some credits to nodes
 - + Don't have to charge nodes equally

} depends
on rank of
node &
rank of parent

Rank(x) - not exactly what you think it is.

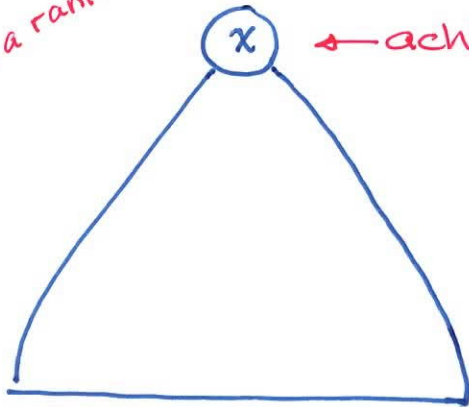
- Rank(x) used to be height of x, but path compression messes that up.
- Once x becomes non-root, rank(x) does not change
- Let $\text{size}(x) = \#$ of items in subtree rooted at x including x.
- Claim: $\text{size}(x) \geq 2^{\text{rank}(x)}$ if x is a root.
Easy induction. Rank increases only during union-by-rank. Size must at least double.
↳ when equal ranked roots are linked.

- For all x , $\text{rank}(x) \leq L \log n$.

Otherwise, $\text{size}(x) > n$.

- At most $n/2^r$ items ever achieve rank r .

Fix a rank r .



← achieves rank r

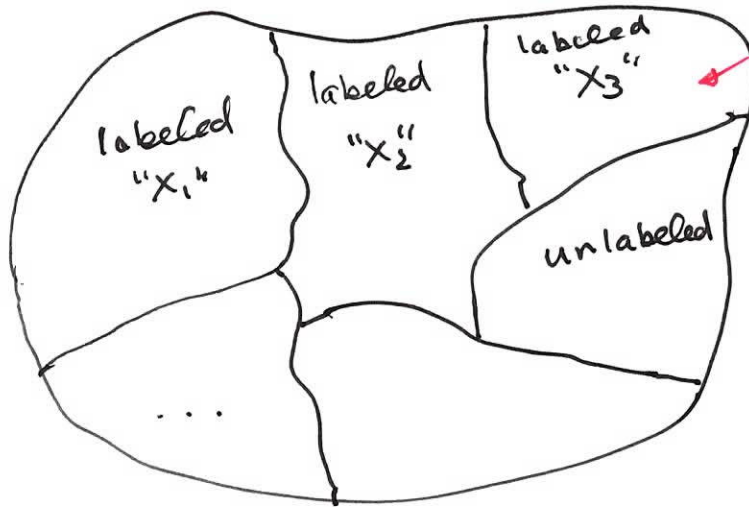
label all descendants "x"

Each item can only ever have 1 ancestor w/ rank r

At least 2^r nodes in subtree

n nodes either unlabeled or labeled by unique "x".

Can't have more than $n/2^r$ labels.



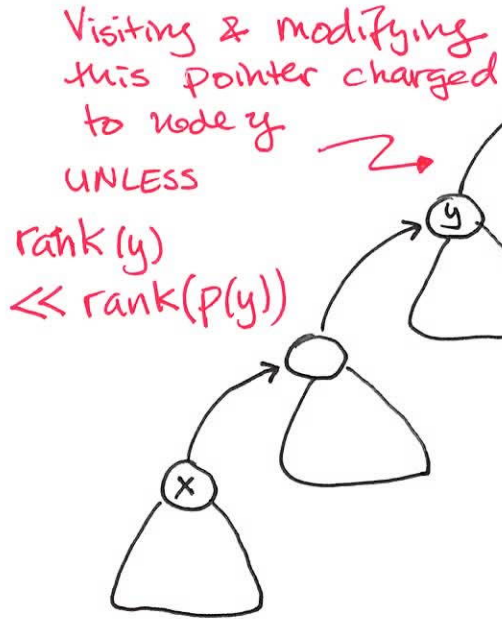
each partition
has at least 2^r items

Can't have more
than $n/2^r$
partitions.

Path Compression

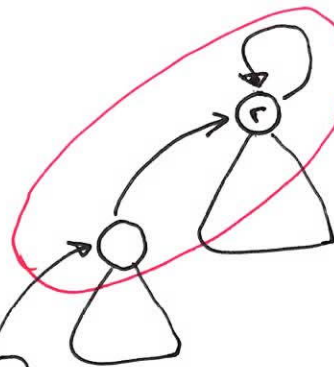
Find-Set(x)

must make pointers point to root r.



parent of y

$p(y)$



Visit to root & child of root charged to Find-Set

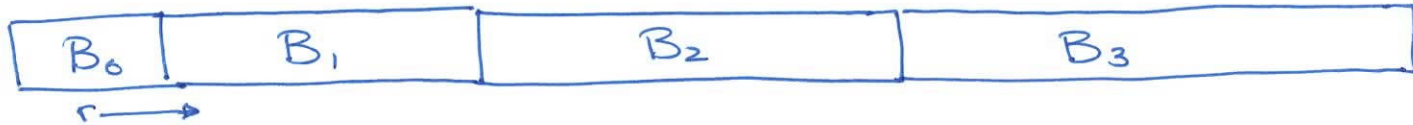
What does " \ll " mean?

note: ① $\text{rank}(y) < \text{rank}(p(y))$
is always true.

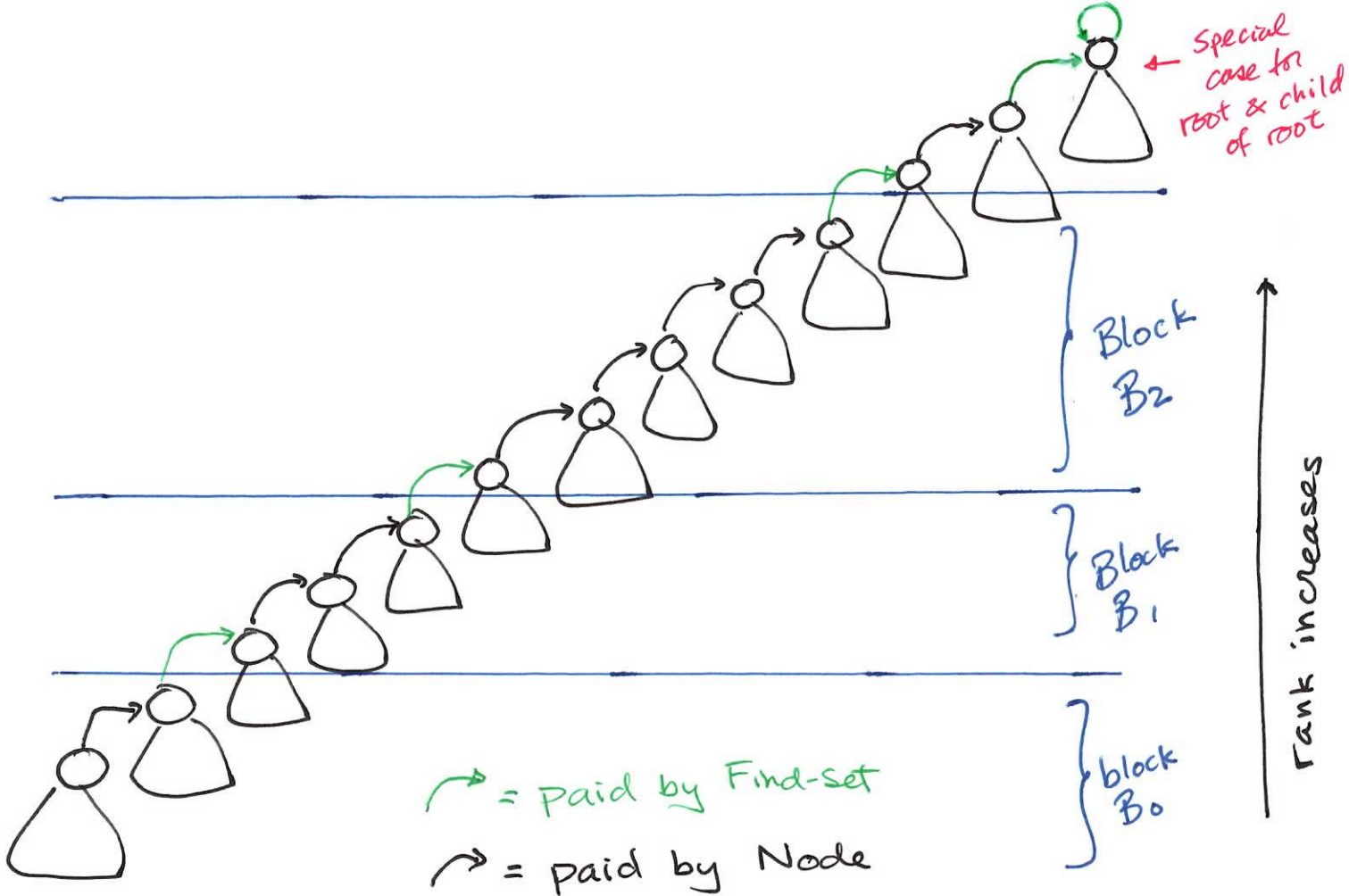
② $\text{rank}(p(y))$ can increase because $p(y)$ changes after path compression.

Idea: Node y pays for modifying its parent pointer until $\text{rank}(p(y))$ gets too large. Then, Find-Set pays.

Divide ranks into "blocks"



if $\text{rank}(y) \notin \text{rank}(p(y))$ are in different blocks,
then we think $\text{rank}(y) \ll \text{rank}(p(y))$ and Find-Set pays.



We want Find-Set to run in $O(\log \log n)$ amortized time, so we can't have too many blocks.

- Limit # of blocks to $\log \log n$.
- Have exponentially increasing block sizes.

useful later

$$B_0 = \{0\}$$

$$B_1 = \{1\}$$

$$B_2 = \{2, 3\}$$

$$B_3 = \{4, 5, 6, 7\}$$

$$B_4 = \{8, 9, 10, \dots, 15\}$$

$$B_5 = \{16, 17, 18, \dots, 31\}$$

$$B_6 = \{32, 33, 34, \dots, 63\}$$

$$B_i = \{2^{i-1}, \dots, 2^i - 1\}$$

$$|B_i| = 2^{i-1} \quad \text{for } i > 0$$

$$\# \text{ of blocks} \leq \log \log n.$$

max rank $\leq \log n$

Items with higher rank are charged more often until rank of their parents move to next block.

(Each time y is involved in path compression, the rank of its parent must increase by at least 1.)

Example: $B_3 = \{4, 5, 6, 7\}$

if $\text{rank}(y) = 4$, then after 3 path compressions $\text{rank}(p(y))$ must be in B_4 or higher.

Example: $B_6 = \{32, 33, 34, \dots, 63\}$

if $\text{rank}(y) = 32$, then it might take 31 path compressions before $\text{rank}(p(y))$ moves to block B_7 .

THIS ^{IS} _^ OK because there are very few items with high rank.

Recap: n items, m operations w/ n MakeSets

Link $\Theta(1)$

Union = 2 Find Sets + 1 Link

Find Set $\Theta(\log \log n)$ amortized

Make Set $\Theta(1) + \underbrace{\text{node charges}} = \Theta(\log \log n)$
amortized

depends on rank, but
we don't know the rank, yet.

Solution: pool together all node
credits.

Total # of node charges:

$$\sum_{i=1}^{\log \log n} \sum_{r \in B_i} \left(\frac{n}{2^r} \right) 2^{i-1} = n \sum_{i=1}^{\log \log n} \sum_{r=2^{i-1}}^{2^i-1} \frac{2^{i-1}}{2^r}$$

Annotations for the first equation:
 - $\sum_{i=1}^{\log \log n}$: Sum over blocks
 - $\sum_{r \in B_i}$: each rank r in block
 - $\left(\frac{n}{2^r} \right) 2^{i-1}$: # of items w/ rank r (max # of node charges for block B_i)

$$\leq n \sum_{i=1}^{\log \log n} \left(2^{i-1} \cdot \frac{2^{i-1}}{2^{2^{i-1}}} \right)$$

Annotations for the second equation:
 - 2^{i-1} : # of r 's
 - $\frac{2^{i-1}}{2^{2^{i-1}}}$: smallest r in B_i

If we charge each of the n MakeSets $\log \log n$ credits, then there are enough credits for all node charges. Some nodes charged more than others.

$$= n \sum_{i=1}^{\log \log n} \frac{2^{2i-2}}{2^{2^{i-1}}}$$

$$< n \sum_{i=1}^{\log \log n} 1 = n \log \log n$$

Even better bounds

$m \log \log \log n$: Make block sizes double exponential

$$B_i = \{2^{2^{i-1}}, \dots, 2^{2^i} - 1\}$$

$m \log^* n$: block size = stack of 2's

$$B_i = \{b_i, \dots, b_{i+1} - 1\}$$

where $b_i = 2^{2^{\dots^2}}$ } $i-1$ times.

See notes online for calculations. ↙ easy