

COS 528

Disjoint Sets and Compressed Trees

© Robert E. Tarjan 2013

Disjoint set union

Devise a data structure for an intermixed sequence of the following kinds of operations:

make-set(x) (x in no set): create a set $\{x\}$, with *root* x .

find(x): (x in a set): return the root of the set containing x .

link(x, y) ($x \neq y$): combine the sets whose roots are x and y into a single set; choose x or y as the root of the new set.

Each element is in at most one set (sets are *disjoint*).

The root of a set serves to identify it, can store information about the set (size, name, etc.)

Applications

Global greedy MST algorithm

FORTRAN compilers: COMMON and
EQUIVALENCE statements

Incremental connected components

Percolation

Additional operations

unite(x, y) ($find(x) \neq find(y)$): *link*($find(x), find(y)$)

contingent-unite(x, y):

if $find(x) = find(y)$ **then return** false

else {*link*($find(x), find(y)$); **return** true}

make-set, *contingent-unite* suffice to implement
global greedy MST algorithm

Variant: named sets

make-set(x, g): create a set $\{x\}$, named g , with root x

find-name(x): return the name of the set containing x

unite(x, y, g) ($find(x) \neq find(y)$): combine the sets containing x and y ; name the new set g

Disjoint set implementation

Represent each set by a rooted tree, whose nodes are the elements of the set, with the set root the tree root, and each node x having a pointer to its parent $a(x)$. Store information about set (such as name) in root.

The shape of the tree is *arbitrary*.

$n = \text{\#elements}$, $m = \text{\#finds}$, $n > 1$, $n = O(m)$

Set operations

make-set(x): make x the root of a new one-node

tree: $a(x) \leftarrow \text{null}$

find(x): follow parent pointers from x to the tree

root: **if** $a(x) = \text{null}$ **then return** x

else return $\text{find}(a(x))$

link(x, y): make y the parent of x (or x the parent

of y): $a(x) \leftarrow y$ (or $a(y) \leftarrow x$)

A bad sequence of links can create a tree that is a path of n nodes, on which each *find* can take $\Omega(n)$ time, totaling $\Omega(mn)$ time for m *finds*

Goal: reduce the amortized time per *find*:
reduce node depths

Improve links: linking by *size* or by *rank*

Improve finds: *compress* the trees

Linking by size: maintain the number of nodes in each tree (store in root). Link root of smaller tree to larger. Break a tie arbitrarily.

make-set(x): $\{a(x) \leftarrow x; s(x) \leftarrow 1\}$

link(x, y):

if $s(x) < s(y)$ **then** $\{a(x) \leftarrow y; s(y) \leftarrow s(y) + s(x)\}$

else $\{a(y) \leftarrow x; s(x) \leftarrow s(x) + s(y)\}$

Linking by rank: Maintain an integer *rank* for each root, initially 0. Link root of smaller rank to root of larger rank. If tie, increase rank of new root by 1.

make-set(*x*): {*a*(*x*) \leftarrow *x*; *r*(*x*) \leftarrow 0}

link(*x*, *y*): {**if** *r*(*x*) = *r*(*y*) **then** *r*(*y*) \leftarrow *r*(*y*) + 1;
if *r*(*x*) < *r*(*y*) **then** *a*(*x*) \leftarrow *y* **else** *a*(*y*) \leftarrow *x*}

r(*x*) = *h*(*x*), the height of *x*

Linking by size and linking by rank have similar efficiency. Linking by rank needs fewer bits ($\lg \lg n$ for rank vs. $\lg n$ for size) and less time: use linking by rank

For any x , $r(a(x)) > r(x)$

Proof: Immediate.

#nodes of rank $\geq k \leq n/2^k$

Proof: Only roots increase in rank. Production of one root of rank $k + 1$ consumes two roots of rank k .

$\rightarrow r(x) \leq \lg n$, $find(x)$ takes $O(\lg n)$ time

Compression

compress(x) (a(a(x)) ≠ null): a(x) ← a(a(x))

Reduces depth of x , reducing find time for x and increasing no find time; preserves sets

Compression preserves $r(a(x)) > r(x)$

With compression, $r(x) \geq h(x)$, but not necessarily equal

Collapsing (fast find)

After a link that makes x a child of y , compress each child of x (make the grandchildren of y children of y). Each tree is *flat*: each node is a root or a child of a root \rightarrow *find* takes $O(1)$ time. To implement: for each tree, maintain a circular linked list of its nodes; during a link, catenate lists

Collapsing: total time is $O(n^2 + m)$: each node changes parent $\leq n$ times

Collapsing with union by rank: total time is $O(n \lg n + m)$: each node changes parent $\leq \lg n$ times

But collapsing uses one extra pointer per node, dominated by path compression (next) no matter how links are done

Collapsing is too eager: better to do compressions only on *find* paths

Path compression

During each find, make the root the parent of each node on the find path, by doing compression top-down along the find path

```
find(x): if a(x) = null then return x  
           else {if a(a(x)) ≠ null then a(x) ← find(a(x));  
                  return a(x)}
```

Alternative implementations

Since parent of root is null, can use parent field of root to store rank (or size), but violates type, bad programming practice

Another use of parent of root: set equal to root instead of null. Saves one test in *find* loop, but does an unneeded assignment

make-set(x): $\{a(x) \leftarrow x; r(x) \leftarrow 0\}$

find(x): **if** $a(a(x)) \neq a(x)$ **then** $a(x) \leftarrow a(a(x));$
return $a(x)$

Collapsing vs. path compression

For any sequence of operations and any linking rule, path compression changes no more parents than collapsing: path compression dominates

Proof: Compare three scenarios: (i) no compression, (ii) collapsing, (iii) path compression. If in (iii) w becomes a parent of v , then in (i) w becomes a proper ancestor of v , and in (ii) w becomes a parent of v .

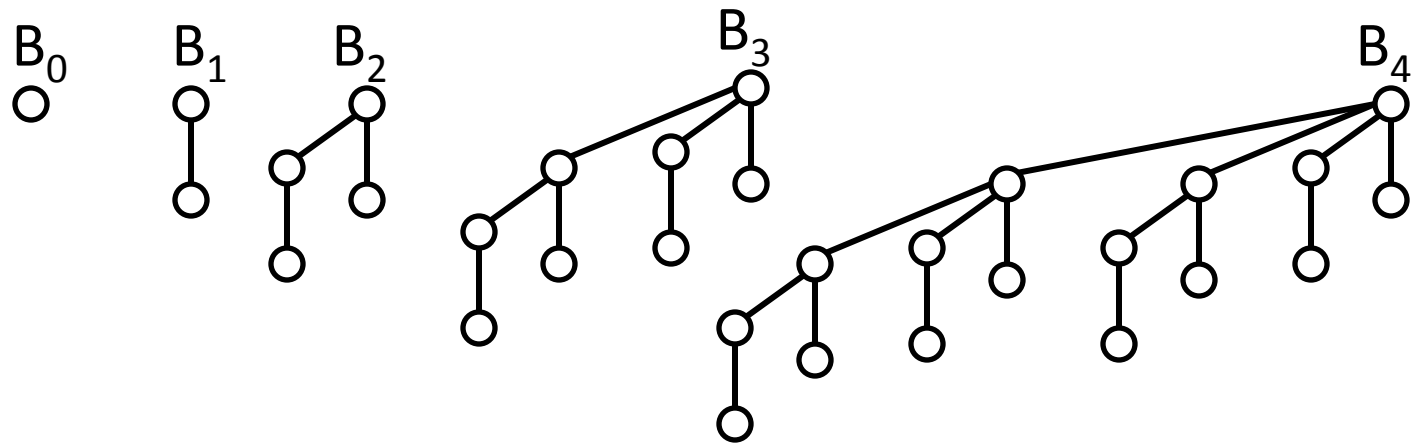
Path compression with naïve linking

Bad example for path compression?

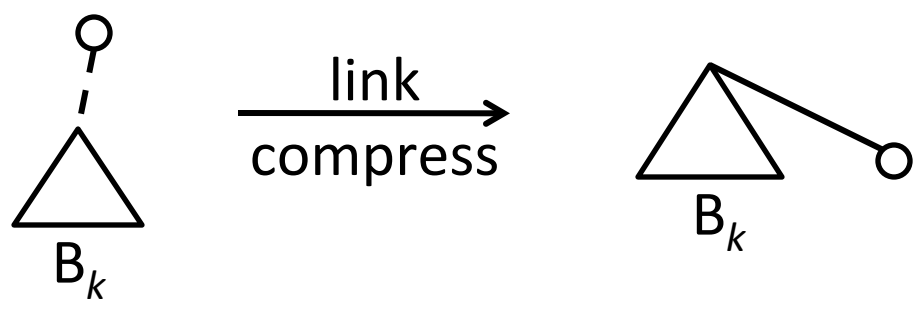
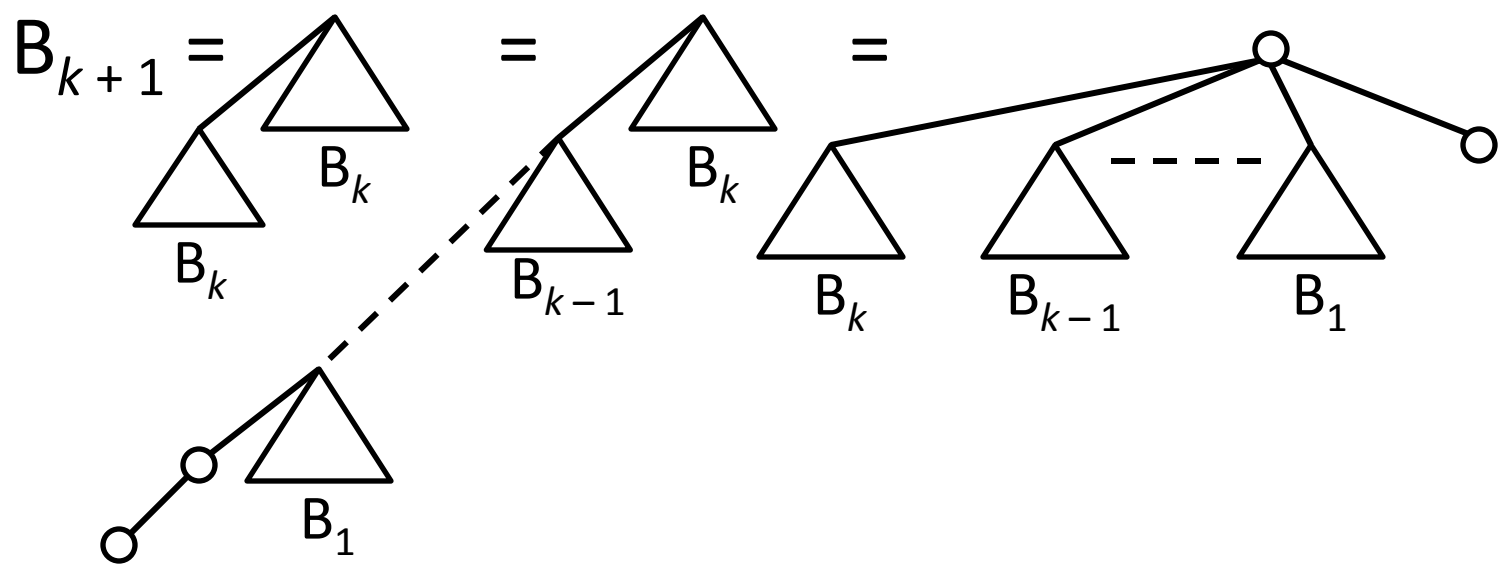
A path of n nodes can result in one find path of n nodes, but compression flattens the tree:
not repeatable

Need a class of trees preserved by path compression: *binomial trees*

Binomial trees



$$B_0 = \circ,$$



Given $n = 2^k$, build a B_{k-1} . Then repeat $n/2$ times: link with a singleton, do a find on deepest element: $\Theta(\lg n)$ time per find

Let the *density of finds* $d = \lceil m/n \rceil$. As d increases, the amortized time per find decreases: $\Theta(\log_{d+1} n)$

Lower bound: class of trees preserved by a link with a singleton followed by d finds (generalized binomial trees: exercise)

Upper bound: charging argument

Count changes of parent once a node becomes a non-root: undercounts number of nodes on each find path by 2

For purposes of the analysis we give each node a *rank*: when $\text{make-set}(x)$ occurs, $r(x) \leftarrow 0$; when $a(y) \leftarrow x$ in a link, $r(x) \leftarrow \max\{r(x), r(y) + 1\}$

Without compression, $r(x) = h(x)$; with compression, $r(x) \geq h(x)$. With or without compression, $r(x) < r(a(x))$; $r(a(x))$ never decreases

Let x be a non-root. We charge a change in $a(x)$ during a find to the corresponding increase in $r(a(x))$. We define $k(x)$ and $j(x)$, the *level of x* and the *index of x* , as follows:

$$k(x) = \max\{k \mid (d + 1)^k \leq r(a(x)) - r(x)\}$$

$$j(x) = \max\{j \mid j(d + 1)^{k(x)} \leq r(a(x)) - r(x)\}$$

$$0 \leq k(x) \leq \lfloor \log_{d+1} n \rfloor, \quad 1 \leq j(x) \leq d$$

When a find occurs, if x is a node whose parent changes, we charge x unless it is the last node in its level along the find path, in which case we charge the find.

Along each find path, there is at most one node per level that is last in the level, totaling $1 + \lfloor \log_{d+1} n \rfloor$ charges per find. The remaining nodes whose parents change are charged for the change.

How many node charges?

Let x be a node. Then $k(x)$ never decreases.

Suppose x gets charged. Then there is a node y after x on the find path such that $k(x) = k(y)$.

Let a, a' , respectively, be the parent functions before and after the compression. Then

$$\begin{aligned} r(a'(x)) - r(x) &\geq r(a(y)) - r(x) \\ &\geq r(a(y)) - r(y) + r(a(x)) - r(x) \\ &\geq (d + 1)^{k(x)} + j(x)(d + 1)^{k(x)} \\ &\geq (j(x) + 1)(d + 1)^{k(x)} \end{aligned}$$

Thus when x is charged, its index or its level increases. Since $j(x)$ can only increase $d - 1$ times before $k(x)$ increases, x can incur at most $d \lfloor \log_{d+1} n \rfloor$ charges. Summing over all nodes, $\text{\#charges} = O(m \log_{d+1} n)$

→ amortized time per find = $O(\log_{d+1} n)$

Path compression with linking by rank

History of bounds (amortized time per find)

1971 $O(1)$ (false)

1972 $O(\lg \lg n)$ M. Fisher

1973 $O(\lg^* n)$ Hopcroft & Ullman

1975 $\Theta(\alpha(n, d))$ Tarjan

later $\Omega(\lg \lg n)$ (false)

2005 top-down analysis Seidel & Sharir

Ackermann's function (Péter & Robinson)

$$A(k, j) = j + 1 \text{ if } k = 0$$

$$= A(k - 1, 1) \text{ if } k > 0, j = 0$$

$$= A(k - 1, A(k, j - 1)) \text{ if } k > 0, j > 0$$

$A(1, j) = j + 2$, $A(2, j) = 2j + 3$, $A(3, j) > 2^j$, $A(4, j) >$
tower of j 2's, $A(4, 2)$ has 19,729 decimal digits

$A(k, j)$ is strictly increasing in both arguments

$$\alpha(n, d) = \min\{k > 0 \mid A(k, d) > n\}$$

Upper bound: charging argument

Count changes of parent once a node becomes a non-root: undercounts number of nodes on each find path by 2. Charge a change in $a(x)$ to the corresponding increase in $r(a(x))$.

If $r(x) \geq d$, we define $k(x)$ and $j(x)$, the *level of x*, and the *index of x*, as follows:

$$k(x) = \max\{k \mid A(k, r(x)) \leq r(a(x))\}$$

$$j(x) = \max\{j \mid A(k(x) + 1, j) \leq r(a(x))\}$$

$$A(0, r(x)) = r(x) + 1 \rightarrow k(x) \geq 0$$

$$A(\alpha(n, d), d) > n \rightarrow k(x) < \alpha(n, d) \quad (r(x) \geq d)$$

$$A(k(x) + 1, 0) = A(k(x), 1) \rightarrow j(x) \geq 0 \quad (r(x) \geq 1)$$

$$A(k(x) + 1, r(x)) > r(a(x)) \rightarrow j(x) < r(x)$$

$\rightarrow 0 \leq k(x) < \alpha(n, d), 0 \leq j(x) < r(x)$ if $r(x) \geq d$

We charge for nodes whose parent changes as a result of a find. If x is such a node, we charge x if $r(x) < d$ and $r(a(x)) < d$, or if $r(x) \geq d$ and x is not last in its level on the find path. Every other node on the find path is either last in its level, at most $\alpha(n, d)$ nodes per find, or its rank is $< d$ but the rank of its parent is $\geq d$, at most one node per find, for a charge per find of at most $\alpha(n, d) + 1$.

Each charged node x of rank $< d$ can accumulate a charge of at most $d - 1$ before its parent has rank $\geq d$ and it is never charged again. Thus the total charge accrued by such nodes is at most $n(d - 1) = O(m)$.

It remains to bound the charges accrued by nodes of rank $\geq d$. Suppose such a node x accrues a charge. Let y be a node after x on the find path with $k(y) = k(x)$. Let a, a' be the parent functions before and after the find, respectively.

$$\begin{aligned}
r(a'(x)) &\geq r(a(y)) \geq A(k(x), r(y)) \geq A(k(x), r(a(x))) \\
&\geq A(k(x), A(k(x) + 1, j(x))) \\
&= A(k(x) + 1, j(x) + 1)
\end{aligned}$$

→ $j(x)$ or $k(x)$ increases as a result of the find

→ #charges accrued by nodes of high rank
 $\leq \alpha(n, d)r$ per node of rank $r \geq d$

The sum of node ranks is at most n (Why?)

→ #charges per find = $O(\alpha(n, d))$

= amortized time per find

This argument can be tightened: the function A can grow even faster; the inverse function α can grow even more slowly, e. g. $\alpha(n, d) = \min\{k > 0 \mid A(k, d) > \lg n\}$, but this only improves the bound by an additive constant.

Seidel and Sharir have shown that for any feasible problem size, the number of parent changes during compressions is at most $m + 2n$

Can extend the bound to the case $m = o(n)$; can tighten the bound to $\alpha(n', d)$, where n' is the number of elements in the set on which the find is done.

The bound holds for some one-pass variants of path compression: *path halving*, *path splitting*

Is the bound tight?

We use **double** induction to build examples that change k pointers per find, for any k . The number of nodes is $B(k, j)$, defined as follows:

$$B(k, j) = 1 \text{ if } k = 0$$

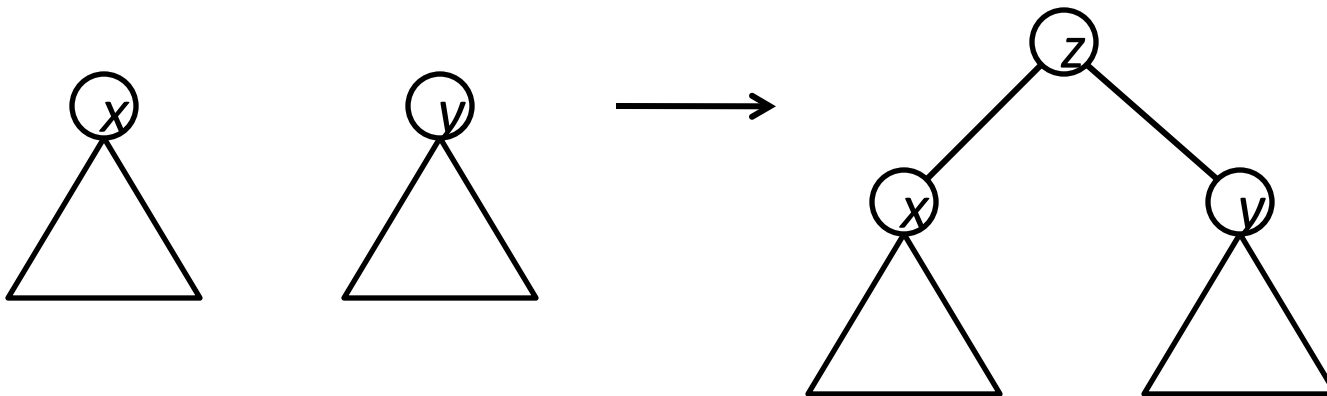
$$= 2B(k - 1, 2) \text{ if } k > 0, j = 1$$

$$= B(k, j - 1) \times B(k - 1, B(k, j - 1)) \text{ if } k > 0, j > 1$$

$B(k, j)$ grows even faster than $A(k, j)$, but the inverses are within an additive constant.

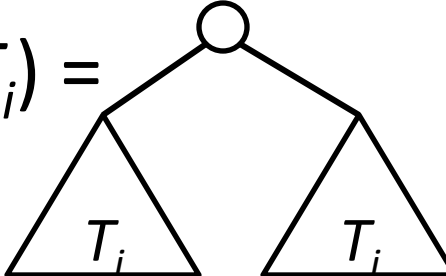
To simplify the argument, we change the way linking is done:

link(x, y): let z be a new root; $a(x) \leftarrow z$; $a(y) \leftarrow z$



We shall only link identical trees. Given a tree T ,

$$T_0 = T,$$

$$T_{i+1} = \text{link}(T_i, T_i) =$$


The diagram shows a tree structure. At the top is a small circle representing a root node. Two lines extend downwards from this root node to the top vertices of two identical triangles. Each triangle is labeled T_i in its center. The two triangles are positioned side-by-side, with the root node centered above the space between them. This represents the operation of linking two identical trees T_i under a new root node to form T_{i+1} .

By such links we can build perfect binary trees.

Theorem: Let T be a tree with j leaves other than the root. If $s = B(k, j)$, then starting with s copies of T , there is a sequence of intermixed links and finds such that each find changes at least k pointers, each link combines two T_i trees to form a T_{i+1} tree, and there are js finds, one on each leaf of an original copy of T .

Proof: By double induction on k and j .

Let $k = 0$. Then can do finds on all leaves of $T = T_0$. Each find changes no pointers; no links are needed.

Suppose true for $k - 1$, any j . Let T have one leaf other than the root. Link two copies of T to form T_1 . Let T' be T_1 with its two leaves deleted. Then T' has two leaves, the parents of the deleted leaves. By the induction hypothesis there is a sequence of intermixed links and finds on $B(k - 1, 2)$ copies of T' that does finds on all the leaves of the copies, each of find changing $k - 1$ pointers.

The corresponding sequence of intermixed links of copies of T_1 and finds on leaves of the copies changes k pointers per find: each find on a parent is replaced by a find on its child; one more pointer is changed since the find path is one edge longer. Thus the theorem holds for $k, j = 1$: $B(k, 1) = 2B(k - 1, 2)$.

Suppose true for $k - 1, any j ; and for $k, j - 1$. Let T be a tree with j leaves. Starting with $B(k, j - 1)$ copies of T , can do links intermixed with finds on $j - 1$ leaves (all but one) in each copy of T .$

Each find changes k pointers, and the final tree T' is a compressed version of T_i for $i = \lg B(k, j - 1)$. Repeat this process $B(k - 1, B(k, j - 1))$ times, resulting in this many copies of T' . Let T'' be T' with all nodes deleted except for proper ancestors of the original leaves which finds have not yet been done. Then T'' has $B(k, j - 1)$ leaves. Starting with $B(k - 1, B(k, j - 1))$ copies of T'' , can do links intermixed with finds on all the leaves, each of which changes $k - 1$ pointers. Instead do the corresponding sequence of operations on the copies of T' , replacing each find by a find on its child that was an original leaf. Each of these finds changes k pointers.

Thus the theorem is true for k, j :

$$B(k, j) = B(k, j - 1) \times B(k - 1, B(k, j - 1))$$

Corollary: Starting with $B(k + 1, j)$ singletons, can do an intermixed sequence of links and finds such that there are j finds of each node and each find changes k pointers.

Corollary: Path compression with linking by rank takes $\Omega(n, d)$ amortized time per find.

Proof: Map original links to new links, add extra pointers (shortcuts) for free.

Proof: Do links by rank and corresponding new links concurrently. Label each new node with the root of the new tree in link by rank. Then each node in link by rank corresponds to a path of nodes in the tree built by new links. Add “shortcut” arcs from each node to all ancestors of the same label. Then if x is a node in an original tree that is k steps from the root, the leaf in the new tree corresponding to x is at most $2k + 1$ steps from the root: each step in the original tree corresponds to at most two steps in the new tree, one via a shortcut arc to the deepest node with the same label, one to a node labeled with the parent of the node.

In the new forest, do finds by generalized compression: if x is the node found and y is the root, add shortcut arcs from every node z on the original path from x to y to every proper ancestor of z on the original path. For the find, charge only the length of the shortest path from x to y before the new shortcuts are added. Add the extra shortcuts only when they are added by finds. The total cost of finds in the new forest is at most twice the cost of finds in the old forest, minus $O(1)$ per find, minus $O(1)$ per extra shortcut. Since there are only $O(n)$ extra shortcuts, a lower bound on finds in the new forest gives a lower bound in the old forest.

Can one make this argument work without introducing generalized compression?

Upper bound by top-down analysis (extra)

Bound the number of parent changes by a
divide-and-conquer recurrence

Solve the recurrence (or just plug it into itself
repeatedly)

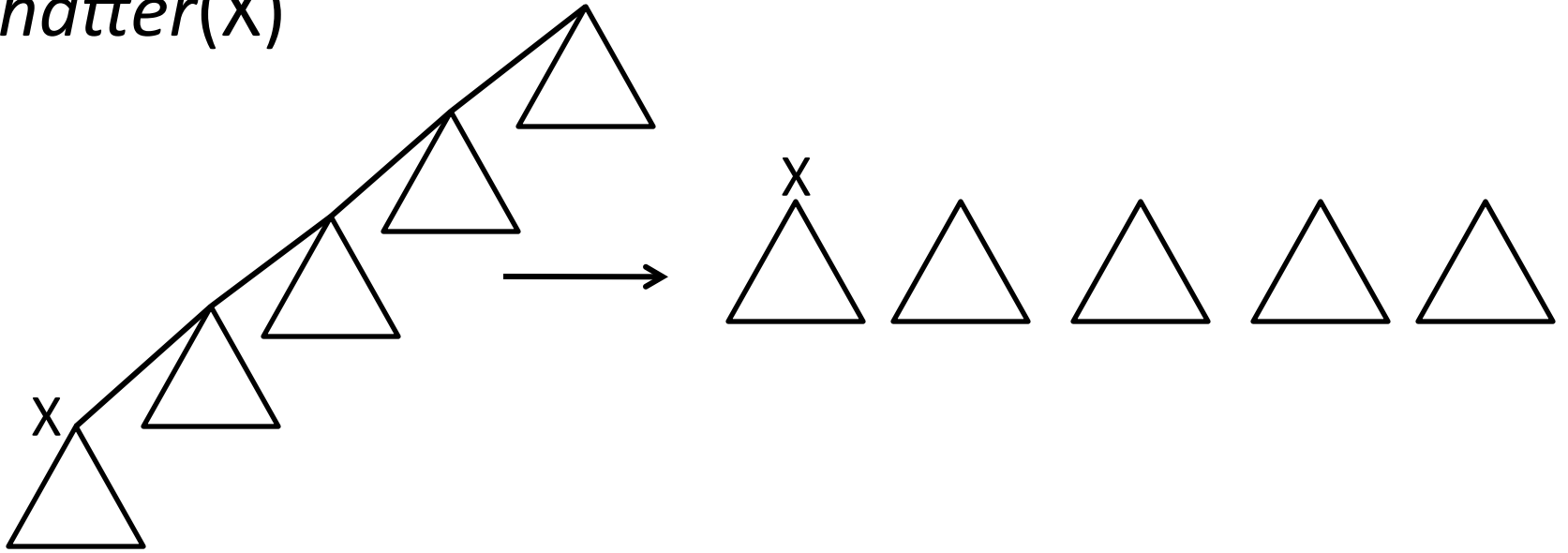
To obtain a closed recurrence, we need a
“funny” form of compression

shatter(x): make every ancestor of x a root, by setting its parent to null

Once $a(x)$ becomes null as a result of a shatter, x can no longer be linked; its tree is only subject to compressions and shatterings

The parent changes (to null) that occur during a shatter are counted outside the recursive subproblem in which the shatter occurs; within the subproblem, these changes are free

shatter(X)



To use the method if linking is naïve, we assign ranks to nodes as described previously

To bound parent changes, we partition nodes into *low* and *high* based on their final (maximum) ranks: if $r(x) < k$, x is low, otherwise x is high (k is a parameter)

Given the original problem, we form two subproblems, the *low* and *high* problems, on the low and high nodes, respectively

We do all the make-set operations before all the links and finds, and we give all nodes their final ranks (so that links do not change any ranks)

Each node in the low problem has the same rank as in the original problem; each node in the high problem has rank k less than its rank in the original problem.

Mapping of operations

Consider $link(x, y)$ in the original problem. Let y be the new root. (Proceed symmetrically if x is the new root.) If y is low, do the link in the low problem, if x is high, do the link in the high problem, and if x is low and y is high, do nothing in either subproblem.

Suppose $find(x)$ in the original problem returns node y . If y is low, do $find(x)$ in the low problem. If x is high, do $find(x)$ in the high problem. If x is low and y is high, let z be the first high node along the find path in the original problem; do $shatter(x)$ in the low problem and $find(z)$ in the high problem.

Each link in the original problem maps to a link in the low or the high problem or to nothing

Each find in the original problem maps to a find in the low or the high problem and, if to a find in the high problem, possibly to a shatter in the low problem.

If x is any non-root in either the low or the high problem, $r(x) < r(a(x))$.

If linking is naïve, in both the low and high problems there is at least one node per rank, from 0 up to the maximum rank.

If linking is by rank, the number of nodes of rank $\geq j$ in the both low and high problems is at most $1/2^j$ times the total number of nodes in the problem, and the number of nodes in the high problem is at most $1/2^k$ times the total number of nodes in the original problem.

The total cost of finds (number of parent changes) in the original problem is at most the total cost of finds in the high and low problems plus the number of nodes in the low problem plus the number of finds in the high problem:

Each find path that contains both low and high nodes contains one low node whose parent is already a high node, and zero or more nodes whose parent is low but whose new parent is high. Over all finds, this is at most one per find in the high problem plus at most one per node in the low problem.

The recurrence

$$C(n, m, r) \leq C(n', m', r') + C(n'', m'', r'') + n' + m''$$

n , m , and r are the number of nodes, the number of finds, and the maximum rank; single and double primes denote the low and high problems, respectively: $n = n' + n''$, $m = m' + m''$, $r = r' + r''$.

One can use this recurrence to bound the amortized time for finds with path compression, with naïve linking or linking by rank.