

Recognizing Union-Find trees built up using union-by-rank strategy is NP-complete^{*}

Kitti Gelle, Szabolcs Iván

Department of Computer Science, University of Szeged, Hungary
 {kgelle,szabivan}@inf.u-szeged.hu

Abstract. Disjoint-Set forests, consisting of Union-Find trees, are data structures having a widespread practical application due to their efficiency. Despite them being well-known, no exact structural characterization of these trees is known (such a characterization exists for Union trees which are constructed without using path compression) for the case assuming union-by-rank strategy for merging. In this paper we provide such a characterization by means of a simple PUSH operation and show that the decision problem whether a given tree (along with the rank info of its nodes) is a Union-Find tree is **NP**-complete, complementing our earlier similar result for the union-by-size strategy.

1 Introduction

Disjoint-Set forests, introduced in [10], are fundamental data structures in many practical algorithms where one has to maintain a partition of some set, which supports three operations: *creating* a partition consisting of singletons, *querying* whether two given elements are in the same class of the partition (or equivalently: *finding* a representative of a class, given an element of it) and *merging* two classes. Practical examples include e.g. building a minimum-cost spanning tree of a weighted graph [4], unification algorithms [18] etc.

To support these operations, even a linked list representation suffices but to achieve an almost-constant amortized time cost per operation, Disjoint-Set forests are used in practice. In this data structure, sets are represented as directed trees with the edges directed towards the root; the CREATE operation creates n trees having one node each (here n stands for the number of the elements in the universe), the FIND operation takes a node and returns the root of the tree in which the node is present (thus the SAME-CLASS(x, y) operation is implemented as $\text{FIND}(x) == \text{FIND}(y)$), and the MERGE(x, y) operation is implemented by merging the trees containing x and y , i.e. making one of the root nodes to be a child of the other root node (if the two nodes are in different classes).

In order to achieve near-constant efficiency, one has to keep the (average) height of the trees small. There are two “orthogonal” methods to do that: first, during the merge operation it is advisable to attach the “smaller” tree below the “larger” one. If the “size” of a tree is the number of its nodes, we say the trees

^{*} Research was supported by the NKFI grant no. 108448.

are built up according to the *union-by-size* strategy, if it's the depth of a tree, then we talk about the *union-by-rank* strategy. Second, during a FIND operation invoked on some node x of a tree, one can apply the *path compression* method, which reattaches each ancestor of x directly to the root of the tree in which they are present. If one applies both the path compression method and either one of the union-by-size or union-by-rank strategies, then any sequence of m operations on a universe of n elements has worst-case time cost $O(m\alpha(n))$ where α is the inverse of the extremely fast growing (not primitive recursive) Ackermann function for which $\alpha(n) \leq 5$ for each practical value of n (say, below 2^{65535}), hence it has an amortized almost-constant time cost [23]. Since it's proven [9] that *any* data structure maintaining a partition has worst-case time cost $\Omega(m\alpha(n))$, the Disjoint-Set forests equipped with a strategy and path compression offer a theoretically optimal data structure which performs exceptionally well also in practice. For more details see standard textbooks on data structures, e.g. [4].

Due to these facts, it is certainly interesting both from the theoretical as well as the practical point of view to characterize those trees that can arise from a forest of singletons after a number of merge and find operations, which we call Union-Find trees in this paper. One could e.g. test Disjoint-Set implementations since if at any given point of execution a tree of a Disjoint-Set forest is not a valid Union-Find tree, then it is certain that there is a bug in the implementation of the data structure (though we note at this point that this data structure is sometimes regarded as one of the “primitive” data structures, in the sense that it is possible to implement a correct version of them that needs not be certifying [21]). Nevertheless, only the characterization of Union trees is known up till now [2], i.e. which correspond to the case when one uses one of the union-by- strategies but *not* path compression. Since in that case the data structure offers only a theoretic bound of $\Theta(\log n)$ on the amortized time cost, in practice all implementations imbue path compression as well, so for a characterization to be really useful, it has to cover this case as well.

In this paper we show that the recognition problem of Union-Find trees is **NP**-complete when the union-by-rank strategy is used, complementing our earlier results [13] where we proved **NP**-completeness for the union-by-size strategy. The proof method applied here resembles to that one, but the low-level details for the reduction (here we use the PARTITION problem, there we used the more restricted version 3-PARTITION as this is a very canonical strongly **NP**-complete problem) differ greatly. This result also confirms the statement from [2] that the problem “seems to be much harder” than recognizing Union trees. As (up to our knowledge) in most of the actual software libraries having this data structure implemented the union-by-rank strategy is used (apart from the cases when one quickly has to query the size of the sets as well), for software testing purposes the current result is more relevant than the one applying union-by-size strategy.

Related work. There is an increasing interest in determining the complexity of the recognition problem of various data structures. The problem was considered for suffix trees [17,22], (parametrized) border arrays [15,20,8,14,16], suffix

arrays [1,7,19], KMP tables [6,12], prefix tables [3], cover arrays [5], and directed acyclic word- and subsequence graphs [1].

2 Notation

A (*ranked*) *tree* is a tuple $t = (V_t, \text{ROOT}_t, \text{RANK}_t, \text{PARENT}_t)$ with V_t being the finite set of its *nodes*, $\text{ROOT}_t \in V_t$ its *root*, $\text{RANK}_t : V_t \rightarrow \mathbb{N}_0$ mapping a nonnegative integer to each node, and $\text{PARENT}_t : (V_t - \{\text{ROOT}_t\}) \rightarrow V_t$ mapping each nonroot node to its parent (so that the graph of PARENT_t is a directed acyclic graph, with edges being directed towards the root). We require $\text{RANK}_t(x) < \text{RANK}_t(\text{PARENT}_t(x))$ for each nonroot node x , i.e. the rank strictly decreases towards the leaves.

For a tree t and a node $x \in V_t$, let $\text{CHILDREN}(t, x)$ stand for the set $\{y \in V_t : \text{PARENT}_t(y) = x\}$ of its children and $\text{CHILDREN}(t)$ stand as a shorthand for $\text{CHILDREN}(t, \text{ROOT}_t)$, the set of depth-one nodes of t . Also, let $x \preceq_t y$ denote that x is a (non-strict) *ancestor* of y in t , i.e. $x = \text{PARENT}_t^k(y)$ for some $k \geq 0$. For $x \in V_t$, let $t|_x$ stand for the *subtree* $(V_x = \{y \in V : x \preceq_t y\}, x, \text{RANK}_t|_{V_x}, \text{PARENT}_t|_{V_x})$ of t rooted at x . As shorthand, let $\text{RANK}(t)$ stand for $\text{RANK}_t(\text{ROOT}_t)$, the rank of the root of t .

Two operations on trees are that of *merging* and *collapsing*. Given two trees $t = (V_t, \text{ROOT}_t, \text{RANK}_t, \text{PARENT}_t)$ and $s = (V_s, \text{ROOT}_s, \text{RANK}_s, \text{PARENT}_s)$ with V_t and V_s being disjoint and $\text{RANK}(t) \geq \text{RANK}(s)$, then their *merge* $\text{MERGE}(t, s)$ (in this order) is the tree $(V_t \cup V_s, \text{ROOT}_t, \text{RANK}, \text{PARENT})$ with $\text{PARENT}(x) = \text{PARENT}_t(x)$ for $x \in V_t$, $\text{PARENT}(\text{ROOT}_s) = \text{ROOT}_t$ and $\text{PARENT}(y) = \text{PARENT}_s(y)$ for each nonroot node $y \in V_s$ of s , and

$$\text{RANK}(\text{ROOT}_t) = \begin{cases} \text{RANK}(t) & \text{if } \text{RANK}(s) < \text{RANK}(t), \\ \text{RANK}(t) + 1 & \text{otherwise,} \end{cases}$$

and $\text{RANK}(x) = \text{RANK}_t(x)$, $\text{RANK}_s(x)$ resp. for each $x \in V_t - \{\text{ROOT}_t\}$, $x \in V_s$ resp.

Given a tree $t = (V, \text{ROOT}, \text{RANK}, \text{PARENT})$ and a node $x \in V$, then $\text{COLLAPSE}(t, x)$ is the tree $(V, \text{ROOT}, \text{RANK}, \text{PARENT}')$ with $\text{PARENT}'(y) = \text{ROOT}$ if y is a nonroot ancestor of x in t and $\text{PARENT}'(y) = \text{PARENT}(y)$ otherwise. For examples, see Figure 1.

Observe that both operations indeed construct a ranked tree (e.g. the rank remains strictly decreasing towards the leaves).

We say that a tree is a *singleton* tree if it has exactly one node, and this node has rank 0.

The class of Union trees is the least class of trees satisfying the following two conditions: every singleton tree is a Union tree, and if t and s are Union trees with $\text{RANK}(t) \geq \text{RANK}(s)$, then $\text{MERGE}(t, s)$ is a Union tree as well.

Analogously, the class of Union-Find trees is the least class of trees satisfying the following three conditions: every singleton tree is a Union-Find tree, if t and

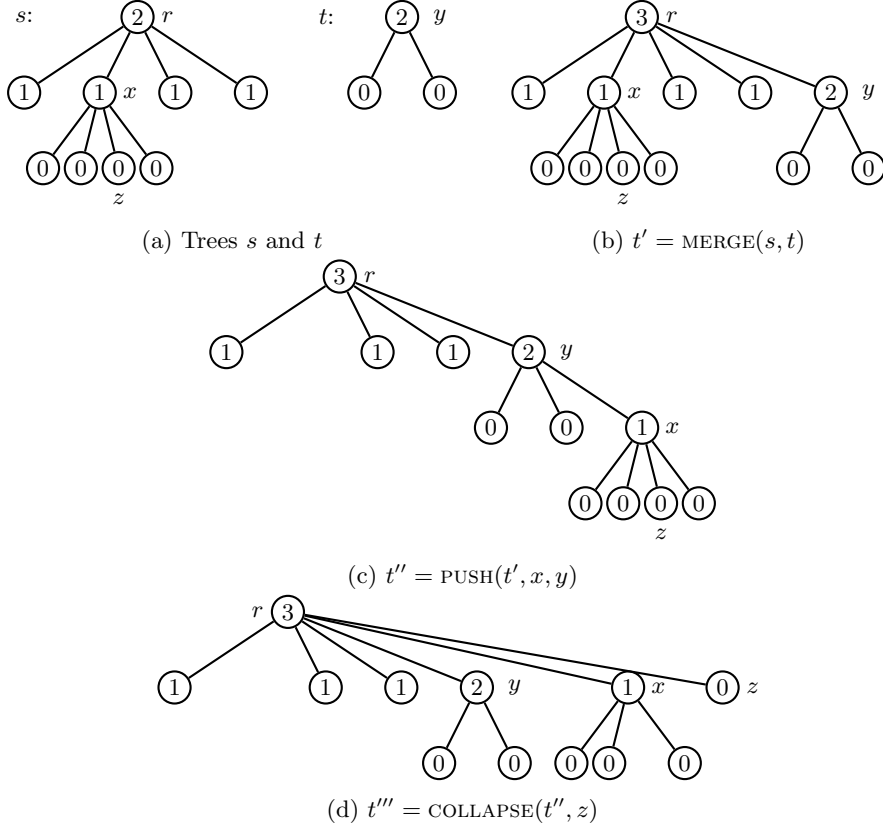


Fig. 1: Merge, collapse and push.

s are Union-Find trees with $\text{RANK}(t) \geq \text{RANK}(s)$, then $\text{MERGE}(t, s)$ is a Union-Find tree as well, and if t is a Union-Find tree and $x \in V_t$ is a node of t , then $\text{COLLAPSE}(t, x)$ is also a Union-Find tree.

We say that a node x of a tree t *satisfies the Union condition* if

$$\{\text{RANK}_t(y) : y \in \text{CHILDREN}(t, x)\} = \{0, 1, \dots, \text{RANK}_t(x) - 1\}.$$

Then, the characterization of Union trees from [2] can be formulated in our terms as follows:

Theorem 1. *A tree t is a Union tree if and only if each node of t satisfies the Union condition.*

Note that the rank of a Union tree always coincides by its *height*. (And, any subtree of a Union tree is also a Union tree.) In particular, the leaves are exactly those nodes of rank 0.

3 Structural characterization of Union-Find trees

Suppose s and t are trees on the same set V of nodes, with the same root ROOT and the same rank function RANK . We write $s \preceq t$ if $x \preceq_s y$ implies $x \preceq_t y$ for each $x, y \in V$.

Clearly, \preceq is a partial order on any set of trees (i.e. is a reflexive, transitive and antisymmetric relation). It is also clear that $s \preceq t$ if and only if $\text{PARENT}_s(x) \preceq_t x$ holds for each $x \in V - \{\text{ROOT}\}$ which is further equivalent to requiring $\text{PARENT}_s(x) \preceq_t \text{PARENT}_t(x)$ since $\text{PARENT}_s(x)$ cannot be x .

Another notion we define is the (partial) operation PUSH on trees as follows: when t is a tree and $x \neq y \in V_t$ are *siblings* in t , i.e. have the same parent, and $\text{RANK}_t(x) < \text{RANK}_t(y)$, then $\text{PUSH}(t, x, y)$ is defined as the tree $(V_t, \text{ROOT}_t, \text{RANK}_t, \text{PARENT}')$ with

$$\text{PARENT}'(z) = \begin{cases} y & \text{if } z = x, \\ \text{PARENT}_t(z) & \text{otherwise,} \end{cases}$$

that is, we “push” the node x one level deeper in the tree just below its former sibling y . (See Figure 1.)

We write $t \vdash t'$ when $t' = \text{PUSH}(t, x, y)$ for some x and y , and as usual, \vdash^* denotes the reflexive-transitive closure of \vdash .

Proposition 1. *For any pair s and t of trees, the following conditions are equivalent:*

- (i) $s \preceq t$,
- (ii) *there exists a sequence $t_0 = s, t_1, t_2, \dots, t_n$ of trees such that for each $i = 1, \dots, n$ we have $t_i = \text{PUSH}(t_{i-1}, x, y)$ for some depth-one node $x \in \text{CHILDREN}(t_{i-1})$, moreover, $\text{CHILDREN}(t_n) = \text{CHILDREN}(t)$ and $t_n|_x \preceq t|_x$ for each $x \in \text{CHILDREN}(t)$,*
- (iii) $s \vdash^* t$.

Proof. i) \Rightarrow ii). It is clear that \preceq is equality on singleton trees, thus \preceq implies \vdash^* for trees of rank 0. Assume $s \preceq t$ for the trees $s = (V, \text{ROOT}, \text{RANK}, \text{PARENT})$ and $t = (V, \text{ROOT}, \text{RANK}, \text{PARENT}')$ and let X stand for the set $\text{CHILDREN}(s)$ of the depth-one nodes of s and Y stand for $\text{CHILDREN}(t)$. Clearly, $Y \subseteq X$ since by $s \preceq t$, any node x of s having depth at least two has to satisfy $\text{PARENT}(x) \preceq_t \text{PARENT}'(x)$ and since $\text{PARENT}(x) \neq \text{ROOT}$ for such nodes, x has to have depth at least two in t as well. Now there are two cases: either $\text{ROOT} = \text{PARENT}(x) = \text{PARENT}'(x)$ for each $x \in X$, or $\text{PARENT}(x) \prec_t \text{PARENT}'(x)$ for some $x \in X$.

If $\text{PARENT}'(x) = \text{ROOT}$ for each $x \in X$, then $X = Y$ and we only have to show that $s|_x \preceq t|_x$ for each $x \in X$. For this, let $u, v \in V(s|_x)$ with $u \preceq_{s|_x} v$. Since $s|_x$ is a subtree of s , this holds if and only if $x \preceq_s u \preceq_s v$. From $s \preceq t$ this implies $x \preceq_t u \preceq_t v$, that is, $u \preceq_{t|_x} v$, hence $s|_x \preceq t|_x$.

Now assume $\text{PARENT}(x) \prec_t \text{PARENT}'(x)$ for some $x \in X$. Then $\text{PARENT}'(x) \neq \text{ROOT}$, thus there exists some $y \in Y$ with $y \preceq_t \text{PARENT}'(x)$. By $Y \subseteq X$, this y is a member of X as well, and $\text{RANK}_s(y) = \text{RANK}_t(y) > \text{RANK}_t(x) = \text{RANK}_s(x)$, thus $s' = \text{PUSH}(s, x, y)$ is well-defined. Moreover, $s' \preceq t$ since $\text{PARENT}_{s'}(z) \preceq_t z$ for

each $z \in V$: either $z \neq x$ in which case $\text{PARENT}_{s'}(z) = \text{PARENT}(z) \preceq_t z$ by $s \preceq t$, or $z = x$ and then $\text{PARENT}_{s'}(z) = y \preceq_t \text{PARENT}'(x) \preceq_t x = z$ also holds. Thus, there exists a tree $s' = \text{PUSH}(s, x, y)$ for some $x \in \text{CHILDREN}(s)$ with $s' \preceq t$; since $\text{CHILDREN}(s') = X - \{x\}$, by repeating this construction we eventually arrive to a tree t_n with $|\text{CHILDREN}(t_n)| = |Y|$, implying $\text{CHILDREN}(t_n) = Y$ by $Y \subseteq \text{CHILDREN}(t_n)$.

ii) \Rightarrow iii). We apply induction on $\text{RANK}(s) = \text{RANK}(t)$. When $\text{RANK}(s) = 0$, then s is a singleton tree and the condition in ii) ensures that t is a singleton tree as well. Thus, $s = t$ and clearly $s \vdash^* t$.

Now let assume the claim holds for each pair of trees of rank less than $\text{RANK}(s)$ and let t_0, \dots, t_n be trees satisfying the condition. Then, by construction, $s \vdash^* t_n$. Since $\text{RANK}(t_n|_x) < \text{RANK}(t_n) = \text{RANK}(s)$ for each node $x \in \text{CHILDREN}(t_n)$, by $t_n|_x \preceq t|_x$ we get applying the induction hypothesis that $t_n|_x \vdash^* t|_x$ for each depth-one node x of t_n , thus $t_n \vdash^* t$, hence $s \vdash^* t$ as well.

iii) \Rightarrow i). For \vdash^* implying \preceq it suffices to show that \vdash implies \preceq since the latter is reflexive and transitive. So let $s = (V, r, \text{RANK}, \text{PARENT})$ and $x \neq y \in V$ be siblings in s with the common parent z , $\text{RANK}(x) < \text{RANK}(y)$ and let $t = \text{PUSH}(s, x, y)$. Then, since $\text{PARENT}_s(x) = z = \text{PARENT}_t(y) = \text{PARENT}_t(\text{PARENT}_t(x))$, we get $\text{PARENT}_s(x) \preceq_t x$, and by $\text{PARENT}_s(w) = \text{PARENT}_t(w)$ for each node $w \neq x$, we have $s \preceq t$.

□

The relations \preceq and \vdash^* are introduced due to their intimate relation to Union-Find and Union trees (similarly to the case of the union-by-size strategy [13], but there the PUSH operation itself was slightly different):

Theorem 2. *A tree t is a Union-Find tree if and only if $t \vdash^* s$ for some Union tree s .*

Proof. Let t be a Union-Find tree. We show the claim by structural induction. For singleton trees the claim holds since any singleton tree is a Union tree as well. Suppose $t = \text{MERGE}(t_1, t_2)$. Then by the induction hypothesis, $t_1 \vdash^* s_1$ and $t_2 \vdash^* s_2$ for the Union trees s_1 and s_2 . Then, for the tree $s = \text{MERGE}(s_1, s_2)$ we get that $t \vdash^* s$. Finally, assume $t = \text{COLLAPSE}(t', x)$ for some node x . Let $x = x_1 \succ x_2 \succ \dots \succ x_k = \text{ROOT}_{t'}$ be the ancestral sequence of x in t' . Then, defining $t_0 = t$, $t_i = \text{PUSH}(t_{i-1}, x_i, x_{i+1})$ we get that $t \vdash^* t_{k-2} = t'$ and $t' \vdash^* s$ for some Union tree s applying the induction hypothesis, thus $t \vdash^* s$ also holds.

Now assume $t \vdash^* s$ (equivalently, $t \preceq s$) for some Union tree s . We show the claim by induction on the height of t . For singleton trees the claim holds since any singleton tree is a Union-Find tree.

Now assume $t = (V, \text{ROOT}, \text{RANK}, \text{PARENT})$ is a tree and $t \vdash^* s$ for some Union tree s . Then by Proposition 1, there is a set $X = \text{CHILDREN}(s) \subseteq \text{CHILDREN}(t)$ of depth-one nodes of t and a function $f : Y \rightarrow X$ with $Y = \{y_1, \dots, y_\ell\} = \text{CHILDREN}(t) - X$ such that for the sequence $t_0 = t$, $t_i = \text{PUSH}(t_{i-1}, y_i, f(y_i))$ we have that $t_\ell|_x \preceq s|_x$ for each $x \in X$. As each $s|_x$ is a Union tree (since so is s), we have by the induction hypothesis that each $t_\ell|_x$ is a Union-Find tree. Now let $X = \{x_1, \dots, x_k\}$ be ordered nondecreasingly by rank; then, as s is a Union tree

and $X = \text{CHILDREN}(s)$, we get that $\{\text{RANK}(x_i)\} = \{0, 1, \dots, \text{RANK}(\text{ROOT}) - 1\}$ by Theorem 1. Hence for the sequence t'_i defined as t'_0 being a singleton tree with root ROOT and for each $i \in \{1, \dots, k\}$, $t'_i = \text{MERGE}(t'_{i-1}, t_{\ell|_{x_i}})$, we get that $t_\ell = t'_k$ is a Union-Find tree. Finally, we get t from t_ℓ by applying successively one COLLAPSE operation on each node in Y , thus t is a Union-Find tree as well. \square

4 Complexity

In order to show **NP**-completeness of the recognition problem, we first make a useful observation.

Proposition 2. *In any Union-Find tree t there are at least as many rank-0 nodes as nodes of positive rank.*

Proof. We apply induction on the structure of t . The claim holds for singleton trees (having one single node of rank 0). Let $t = \text{MERGE}(t_1, t_2)$ and suppose the claim holds for t_1 and t_2 . There are two cases.

- Assume $\text{RANK}(t_1) = 0$. Then, since $\text{RANK}(t_1) \geq \text{RANK}(t_2)$ we have that $\text{RANK}(t_2)$ is 0 as well, i.e. both t_1 and t_2 are singleton trees (of rank 0). In this case t has one node of rank 1 and one node of rank 0.
- If $\text{RANK}(t_1) > 0$, then (since ROOT_{t_1} is the only node in $V_t = V_{t_1} \cup V_{t_2}$ whose rank can change at all, in which case it increases) neither the total number of rank-0 nodes nor the total number of nodes with positive rank changes, thus the claim holds.

Let $t = \text{COLLAPSE}(s, x)$ and assume the claim holds for s . Then, since the COLLAPSE operation does not change the rank of any of the nodes, the claim holds for t as well. \square

In order to define a reduction from the strongly **NP**-complete problem PARTITION we introduce several notions on trees:

An *apple* of weight a for an integer $a > 0$ is a tree consisting of a root node of rank 2, a depth-one node of rank 0 and a depth-one nodes of rank 1.

A *basket* of size H for an integer $H > 0$ is a tree consisting of $H + 4$ nodes: the root node having rank 3, $H + 1$ depth-one children of rank 0 and one depth-one child of rank 1, which in turn has a child of rank 0.

A *flat tree* is a tree t of the following form: the root of t has rank 4. The immediate subtrees of t are:

- a node of rank 0, having no children;
- a node of rank 1, having a single child of rank 0;
- a node of rank 2, having two children: a single node of rank 0 and a node of rank 1, having a single child of rank 0;
- an arbitrary number of apples,
- and an arbitrary number of baskets for some fixed size H .

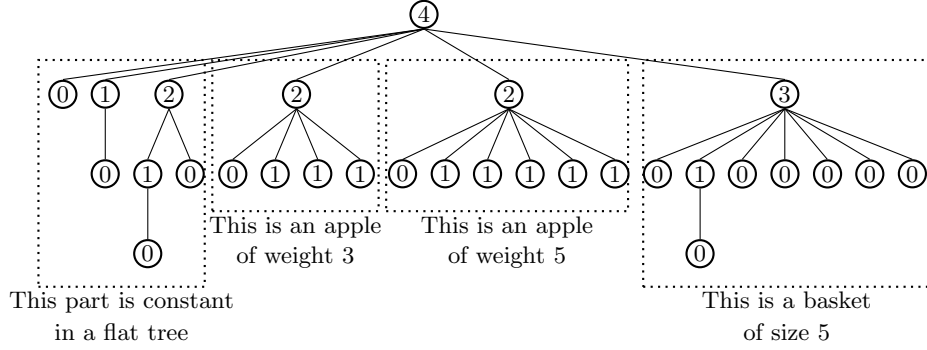


Fig. 2: A flat tree.

(See Figure 2.)

At this point we recall that the following problem PARTITION is **NP**-complete in the strong sense [11]: given a list a_1, \dots, a_m of positive integers and a value $k > 0$ such that the value $B = \frac{\sum_{i=1}^m a_i}{k}$ is an integer, does there exist a partition $\mathcal{B} = \{B_1, \dots, B_k\}$ of the set $\{1, \dots, m\}$ satisfying $\sum_{i \in B_j} a_i = B$ for each $1 \leq j \leq k$?

(Here “in the strong sense” means that the problem remains **NP**-complete even if the numbers are encoded in unary.)

Proposition 3. Assume t is a flat tree having k basket children, each having the size H , and m apple children of weights a_1, \dots, a_m respectively, satisfying $H \cdot k = \sum_{1 \leq i \leq m} a_i$.

Then t is a Union-Find tree if and only if the instance (a_1, \dots, a_m, k) is a positive instance of the PARTITION problem.

Proof. (For an example, the reader is referred to Figure 3.)

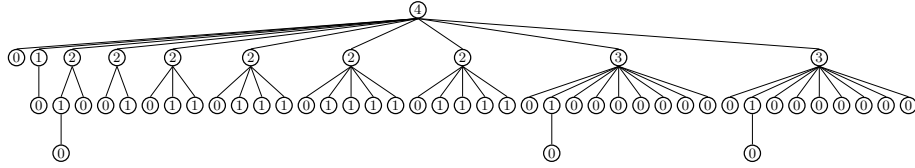


Fig. 3: The initial flat tree t corresponding to the PARTITION instance $(1, 2, 3, 4, 4, k = 2)$. The size of each basket is $(1 + 2 + 3 + 4 + 4)/k = 7$.

Suppose $\mathcal{I} = (a_1, \dots, a_m, k)$ is a positive instance of the PARTITION problem. Let H stand for the target sum $\frac{\sum a_i}{k}$. Let $\mathcal{B} = \{B_1, \dots, B_k\}$ be a solution of \mathcal{I} , i.e., $\sum_{i \in B_j} a_i = H$ for each $j = 1, \dots, k$. Let $x_1, \dots, x_k \in \text{CHILDREN}(t)$ be the

nodes corresponding to the baskets of t and let $y_1, \dots, y_m \in \text{CHILDREN}(t)$ be the nodes corresponding to the apples of t .

We define the following sequence t_0, t_1, \dots, t_m of trees: $t_0 = t$ and for each $i = 1, \dots, m$, let $t_i = \text{PUSH}(t_{i-1}, y_i, x_j)$ with $1 \leq j \leq k$ being the unique index with $i \in \mathcal{B}_j$. Then, $\text{CHILDREN}(t_m)$ consists of x_1, \dots, x_k and the three additional nodes having rank 0, 1 and 2. Note that the subtrees rooted at the latter three nodes are Union trees. Thus, if each of the trees $t_m|_{x_j}$ is a Union-Find tree, then so is t .

Consider a subtree $t' = t_m|_{x_j}$. By construction, t' is a tree whose root has rank 3 and has

- $H + 1$ children of rank 0,
- a single child of rank 1, having a child of rank 0,
- and several (say, ℓ) apple children with total weight H .

We give a method to transform t' into a Union tree. First, we push a_i rank-0 nodes to each apple child of weight a_i . After this stage t' has one child of rank 0, one child of rank 1 and ℓ “filled” apple children, having a root of rank 2, thus the root of the transformed t' satisfies the Union condition. We only have to show that each of these “filled” apples is a Union-Find tree.

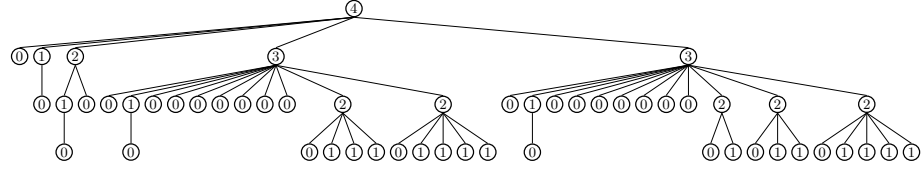
Such a subtree has a root node of rank 2, a_i depth-one nodes of rank 1 and $a_i + 1$ depth-one nodes of rank 0. Then, one can push into each node of rank 1 a node of rank 0 and arrive to a tree with one depth-one node of rank 0, and a_i depth-one nodes of rank 1, each having a single child of rank 0, which is indeed a Union tree, showing the claim by Theorem 2.

For an illustration of the construction the reader is referred to Figure 4.

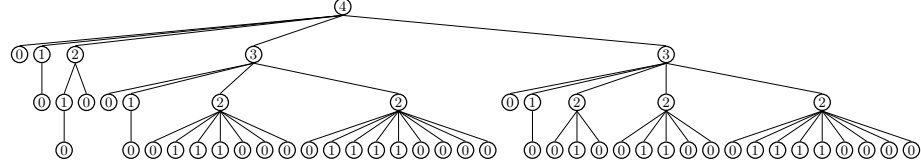
For the other direction, suppose t is a Union-Find tree. By Theorem 2 and Proposition 1, there is a subset $X \subseteq \text{CHILDREN}(t)$ and a mapping $f : Y \rightarrow X$ with $Y = \{y_1, \dots, y_\ell\} = \text{CHILDREN}(t) - X$ such that for the sequence $t_0 = t$, $t_i = \text{PUSH}(t_{i-1}, y_i, f(y_i))$ we have that each immediate subtree of t_ℓ is a Union-Find tree and moreover, the root of t_ℓ satisfies the Union condition.

The root of t has rank 4, t_ℓ has to have at least one child having rank 0, 1, 2 and 3 respectively. Since t has exactly one child with rank 0 and rank 1, these nodes has to be in X . This implies that no node gets pushed into the apples at this stage (because the apples have rank 2). Thus, since the apples are *not* Union-Find trees (as they have strictly less rank-0 nodes than positive-rank nodes, cf. Proposition 2), all the apples have to be in Y . Apart from the apples, t has exactly one depth-one node of rank 2 (which happens to be a root of a Union tree), thus this node has to stay in X as well. Moreover, we cannot push the baskets as they have the maximal rank 3, hence they cannot be pushed.

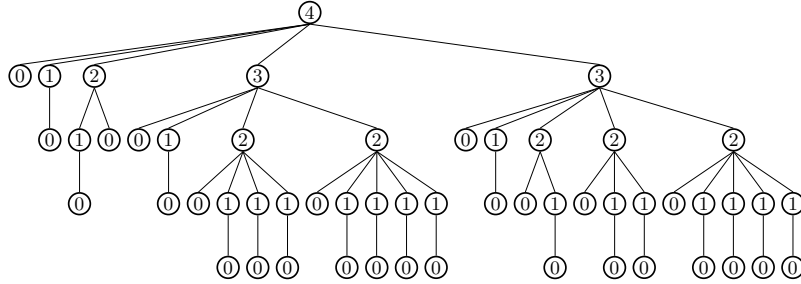
Thus, we have to push all the apples, and we can push apples only into baskets (as exactly the baskets have rank greater than 2). Let $x \in X$ be a basket node, let t' stand for $t_\ell|_x$ and let $\{y'_1, \dots, y'_j\} \subseteq Y$ be the set of those apples that get pushed into x during the operation. Then, the total number of nodes having rank 0 in t' is $H + 2 + j$ (j of them coming from the apples and the other ones coming from the basket) while the total number of nodes having a positive rank is $2 + j + A$ where A is the total weight of the apples in $\{y'_1, \dots, y'_j\}$. Applying



(a) Apples of size 3 and 4 are pushed into the first basket, apples of size 1, 2 and 4 are pushed into the second basket.



(b) The apples get filled from the baskets' surplus rank-0 leaves.



(c) The filling of the apples is pushed a level deeper and we have a Union tree.

Fig. 4: Pushing t of Figure 3 into a Union tree according to the solution $3 + 4 = 1 + 2 + 4$ of the PARTITION instance.

Proposition 2 we get that $A \leq H$ for each basket. Since the total weight of all apples is $H \cdot k$ and each apple gets pushed into exactly one basket, we get that $A = H$ actually holds for each basket. Thus, \mathcal{I} is a positive instance of the PARTITION problem. \square

Theorem 3. *The recognition problem of Union-Find trees is **NP**-complete.*

Proof. By Proposition 3 we get **NP**-hardness. For membership in **NP**, we make use of the characterization given in Theorem 2 and that the possible number of pushes is bounded above by n^2 : upon pushing x below y , the depth of x and its descendants increases, while the depth of the other nodes remains the same. Since the depth of any node is at most n , the sum of the depths of all the nodes is at most n^2 in any tree. Hence, it suffices to guess nondeterministically a sequence $t = t_0 \vdash t_1 \vdash \dots \vdash t_k$ for some $k \leq n^2$ with t_k being a Union tree (which also can be checked in polynomial time). \square

5 Conclusion, future directions

We have shown that unless $\mathbf{P} = \mathbf{NP}$, there is no efficient algorithm to check whether a given tree is a valid Union-Find tree, assuming union-by-rank strategy, since the problem is \mathbf{NP} -complete, complementing our earlier results assuming union-by-size strategy. A very natural question is the following: does there exist a merging strategy under which the time complexity remains amortized almost-constant, and at the same time allows an efficient recognition algorithm? Although this data structure is called “primitive” in the sense that it does not really need an automatic run-time certifying system, but we find the question to be also interesting from the mathematical point of view as well. It would be also an interesting question whether the recognition problem of Union-Find trees built up according to the union-by-rank strategy is still \mathbf{NP} -complete if the nodes of the tree are not tagged with the rank, that is, given a tree without rank info, does there exist a Union-Find tree with the same underlying tree?

References

1. Hideo Bannai, Shunsuke Inenaga, Ayumi Shinohara, and Masayuki Takeda. Inferring strings from graphs and arrays. In Branislav Rovan and Peter Vojtáš, editors, *Mathematical Foundations of Computer Science 2003: 28th International Symposium, MFCS 2003*, pages 208–217. Springer Berlin Heidelberg, Berlin, Heidelberg, 2003.
2. Leizhen Cai. The recognition of Union trees. *Inf. Process. Lett.*, 45(6):279–283, 1993.
3. Julien Clément, Maxime Crochemore, and Giuseppina Rindone. Reverse engineering prefix tables. In Susanne Albers and Jean-Yves Marion, editors, *26th International Symposium on Theoretical Aspects of Computer Science, STACS 2009*, volume 3 of *LIPICs*, pages 289–300. Schloss Dagstuhl - Leibniz-Zentrum fuer Informatik, Germany, 2009.
4. Thomas H. Cormen, Clifford Stein, Ronald L. Rivest, and Charles E. Leiserson. *Introduction to Algorithms*. McGraw-Hill Higher Education, 2nd edition, 2001.
5. Maxime Crochemore, Costas S. Iliopoulos, Solon P. Pissis, and German Tischler. Cover array string reconstruction. In Amihoud Amir and Laxmi Parida, editors, *Combinatorial Pattern Matching: 21st Annual Symposium, CPM 2010*, pages 251–259, Berlin, Heidelberg, 2010. Springer Berlin Heidelberg.
6. J.-P. Duval, T. Lecroq, and A. Lefebvre. Efficient validation and construction of border arrays and validation of string matching automata. *RAIRO - Theoretical Informatics and Applications*, 43(2):281–297, 2009.
7. J.-P. Duval and A. Lefebvre. Words over an ordered alphabet and suffix permutations. *Theoretical Informatics and Applications*, 36(3):249–259, 2002.
8. Jean-Pierre Duval, Thierry Lecroq, and Arnaud Lefebvre. Border array on bounded alphabet. *J. Autom. Lang. Comb.*, 10(1):51–60, 2005.
9. M. Fredman and M. Saks. The cell probe complexity of dynamic data structures. In *Proceedings of the Twenty-first Annual ACM Symposium on Theory of Computing, STOC ’89*, pages 345–354, New York, NY, USA, 1989. ACM.
10. Bernard A. Galler and Michael J. Fisher. An improved equivalence algorithm. *Commun. ACM*, 7(5):301–303, May 1964.

11. Michael R. Garey and David S. Johnson. *Computers and Intractability: A Guide to the Theory of NP-Completeness*. W. H. Freeman & Co., New York, NY, USA, 1979.
12. P. Gawrychowski, A. Jez, and L. Jez. Validating the Knuth-Morris-Pratt failure function, fast and online. *Theory of Computing Systems*, 54(2):337–372, 2014.
13. Kitti Gelle and Szabolcs Iván. Recognizing Union-Find trees is NP-complete. *CoRR*, abs/1510.07462, 2015.
14. Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Counting parameterized border arrays for a binary alphabet. In *Language and automata theory and applications*, volume 5457 of *Lecture Notes in Comput. Sci.*, pages 422–433. Springer, Berlin, 2009.
15. Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Verifying and enumerating parameterized border arrays. *Theoretical Computer Science*, 412(50):6959 – 6981, 2011.
16. Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Verifying and enumerating parameterized border arrays. *Theoret. Comput. Sci.*, 412(50):6959–6981, 2011.
17. Tomohiro I, Shunsuke Inenaga, Hideo Bannai, and Masayuki Takeda. Inferring strings from suffix trees and links on a binary alphabet. *Discrete Applied Mathematics*, 163, Part 3:316 – 325, 2014. Stringology Algorithms.
18. Kevin Knight. Unification: A multidisciplinary survey. *ACM Comput. Surv.*, 21(1):93–124, 1989.
19. Gregory Kucherov, Lilla Tóthmérész, and Stéphane Vialette. On the combinatorics of suffix arrays. *Information Processing Letters*, 113(22–24):915 – 920, 2013.
20. Weilin Lu, P. J. Ryan, W. F. Smyth, Yu Sun, and Lu Yang. Verifying a border array in linear time. *J. Comb. Math. Comb. Comput.*, 42:223–236, 2000.
21. Ross M. McConnell, Kurt Mehlhorn, Stefan Näher, and Pascal Schweitzer. Certifying algorithms. *Computer Science Review*, 5(2):119–161, 2011.
22. Tatiana Starikovskaya and Hjalte Wedel Vildhøj. A suffix tree or not a suffix tree? *Journal of Discrete Algorithms*, 32:14 – 23, 2015. StringMasters 2012 2013 Special Issue (Volume 2).
23. Robert Endre Tarjan. Efficiency of a good but not linear set union algorithm. *J. ACM*, 22(2):215–225, 1975.