

A New Self-stabilizing Minimum Spanning Tree Construction with Loop-Free Property

Lélia Blin¹, Maria Potop-Butucaru¹, Stephane Rovedakis², and Sébastien Tixeuil¹

¹ Univ. Pierre & Marie Curie - Paris 6,
LIP6-CNRS UMR 7606, France

{lelia.blin, maria.gradinariu, sebastien.tixeuil}@lip6.fr

² Université d'Evry, IBISC, CNRS FRE 3190, France
stephane.rovedakis@ibisc.univ-evry.fr

Abstract. The minimum spanning tree (MST) construction is a classical problem in Distributed Computing for creating a globally minimized structure distributedly. Self-stabilization is versatile technique for forward recovery that permits to handle any kind of transient faults in a unified manner. The loop-free property provides interesting safety assurance in dynamic networks where edge-cost changes during operation of the protocol.

We present a new self-stabilizing MST protocol that improves on previous known approaches in several ways. First, it makes fewer system hypotheses as the size of the network (or an upper bound on the size) need *not* be known to the participants. Second, it is loop-free in the sense that it guarantees that a spanning tree structure is always preserved while edge costs change dynamically and the protocol adjusts to a new MST. Finally, time complexity matches the best known results, while space complexity results show that this protocol is the most efficient to date.

1 Introduction

Since its introduction in a centralized context [25,22], the minimum spanning tree (or MST) construction problem gained a benchmark status in distributed computing thanks to the influential seminal work of [13]. Given an edge-weighted graph $G = (V, E, w)$, where w denotes the edge-weight function, the MST problem consist in computing a tree T spanning V , such that T has minimum weight among all spanning trees of G .

One of the most versatile techniques to ensure forward recovery of distributed systems is that of *self-stabilization* [6,7]. A distributed algorithm is self-stabilizing if after faults and attacks hit the system and place it in some arbitrary global state, the system recovers from this catastrophic situation without external (*e.g.* human) intervention in finite time. A recent trend in self-stabilizing research is to complement the self-stabilizing abilities of a distributed algorithm with some additional *safety* properties that are guaranteed when the permanent and intermittent failures that hit the system satisfy some conditions. In addition to being self-stabilizing, a protocol could thus also tolerate a limited number of topology changes [9], crash faults [15,2], nap faults [10,23], Byzantine faults [11,3], and sustained edge cost changes [4,20].

This last property is especially relevant when building spanning trees in dynamic networks, since the cost of a particular edge is likely to evolve through time. If an MST

protocol is *only* self-stabilizing, it may adjust to the new costs in such a way that a previously constructed MST evolves into a disconnected or a looping structure (of course, in the absence of new edge cost changes, the self-stabilization property guarantees that *eventually* a new MST is constructed). Of course, if edge costs change unexpectedly and continuously, an MST cannot be maintained at all times. Now, a packet routing algorithm is *loop free* [14,12] if at any point in time the routing tables are free of loops, despite possible modification of the edge-weights in the graph (*i.e.*, for any two nodes u and v , the actual routing tables determines a simple path from u to v , at any time). The *loop-free* property [4,20] in self-stabilization guarantees that, a spanning tree being constructed (not necessarily an MST), then the self-stabilizing convergence to a “minimal” (for some metric) spanning tree maintains a spanning tree at all times (obviously, this spanning tree is not “minimal” at all times). The consequence of this safety property in addition to that of self-stabilization is that the spanning tree structure can still be used (*e.g.* for routing) while the protocol is adjusting, and makes it suitable for networks that undergo such very frequent dynamic changes.

Related works. Gupta and Srimani [18] have presented the first self-stabilizing algorithm for the MST problem. It applies on graphs whose nodes have unique identifiers, whose edges have integer edge weights, and a weight can appear at most once in the whole network. To construct the (unique) MST, every node performs the same algorithm. The MST construction is based on the computation of all the shortest paths (for a certain cost function) between all the pairs of nodes. While executing the algorithm, every node stores the cost of all paths from it to all the other nodes. To implement this algorithm, the authors assume that every node knows the number n of nodes in the network, and that the identifiers of the nodes are in $\{1, \dots, n\}$. Every node u stores the weight of the edge $e_{u,v}$ placed in the MST for each node $v \neq u$. Therefore the algorithm requires $\Omega(\sum_{v \neq u} \log w(e_{u,v}))$ bits of memory at node u . Since all the weights are distinct integers, the memory requirement at each node is $\Omega(n \log n)$ bits.

Higham and Lyan [19] have proposed another self-stabilizing algorithm for the MST problem. As in [18], their work applies to undirected connected graphs with unique integer edge weights and unique node identifiers, where every node has an upper bound on the number of nodes in the system. The algorithm performs roughly as follows: every edge aims at deciding whether it eventually belongs to the MST or not. For this purpose, every non tree-edge e floods the network to find a potential cycle, and when e receives its own message back along a cycle, it uses information collected by this message (*i.e.*, the maximum edge weight of the traversed cycle) to decide whether e could potentially be in the MST or not. If the edge e has not received its message back after the time-out interval, it decides to become tree edge. The core memory of each node holds only $O(\log n)$ bits, but the information exchanged between neighboring nodes is of size $O(n \log n)$ bits, thus only slightly improving that of [18].

To our knowledge, *none* of the self-stabilizing MST construction protocols is loop-free. Since the aforementioned two protocols also make use of the knowledge of the global number of nodes in the system, and assume that no two edge costs can be equal, these extra hypotheses make them suitable for static networks only.

Relatively few works investigate merging self-stabilization and loop free routing, with the notable exception of [4,20]. While [4] still requires that a upper bound on the

Table 1. Distributed Self-Stabilizing algorithms for the MST and loop-free SP problems

	metric	size known	unique weights	memory usage	loop-free
[18]	MST	yes	yes	$O(n \log n)$	no
[19]	MST	upper bound	yes	$O(n \log n)$	no
[4]	SP	upper bound	no	$\Theta(\log n)$	yes
[20]	SP	no	no	$\Theta(\log n)$	yes
This paper	MST	no	no	$O(\log n)$	yes

network diameter is known to every participant, no such assumption is made in [20]. Also, both protocols use only a reasonable amount of memory ($O(\log n)$ bits per node). However, the metrics that are considered in [4,20] are derivative of the shortest path (*a.k.a.* SP) metric, that is considered a much easier task in the distributed setting than that of the MST, since the associated metric is *locally optimizable* [17], allowing essentially locally greedy approaches to perform well. By contrast, some sort of *global optimization* is needed for MST, which often drives higher complexity costs and thus less flexibility in dynamic networks.

Our contributions. We describe a new self-stabilizing algorithm for the MST problem. Contrary to previous self-stabilizing MST protocols, our algorithm does not make any assumption about the network size (including upper bounds) or the unicity of the edge weights. Moreover, our solution improves on the memory space usage since each participant needs only $O(\log n)$ bits¹, and node identifiers are not needed.

In addition to improving over system hypotheses and complexity, our algorithm provides additional safety properties to self-stabilization, as it is loop-free. Compared to previous protocols that are both self-stabilizing and loop-free, our protocol is the first to consider non-monotonous tree metrics.

The key techniques that are used in our scheme include fast construction of a spanning tree, that is continuously improved by means of a pre-order construction over the nodes. The cycles that are considered over time are precisely those obtained by adding one edge to the evolving spanning tree. Considering solely that type of cycles reduces the memory requirement at each node compared to [18,19] because the latter consider all possible paths connecting pairs of nodes. Moreover, constructing and using a pre-order on the nodes allows our algorithm to proceed in a completely asynchronous manner, and without any information about the size of the network, as opposed to [18,19]. The main characteristics of our solution are presented in Table 1, where a boldface denotes the most useful (or efficient) feature for a particular criterium.

2 Model and Notations

We consider an undirected weighted connected network $G = (V, E, w)$ where V is the set of nodes, E is the set of edges and $w : E \rightarrow \mathbb{R}^+$ is a positive cost function. Nodes

¹ Note that one may use the techniques proposed in [1] in order to construct a self-stabilizing MST starting from non-stabilizing solutions. This technique would increase the memory complexity.

represent processors and edges represent bidirectional communication links. Additionally, we consider that $G = (V, E, w)$ is a network in which the weight of the communication links may change value. We consider anonymous networks (i.e., the processors have no IDs), with one distinguished node, called the *root*². Throughout the paper, the root is denoted r . We denote by $\deg(v)$ the number of v 's neighbors in G . The $\deg(v)$ edges incident to any node v are labeled from 1 to $\deg(v)$, so that a processor can distinguish the different edges incident to a node.

The processors asynchronously execute their programs consisting of a set of variables and a finite set of rules. The variables are part of the shared register which is used to communicate with the neighbors. A processor can read and write its own registers and can read the shared registers of its neighbors. Each processor executes a program consisting of a sequence of guarded rules. Each *rule* contains a *guard* (boolean expression over the variables of a node and its neighborhood) and an *action* (update of the node variables only). Any rule whose guard is *true* is said to be *enabled*. A node with one or more enabled rules is said to be *privileged* and may make a *move* executing the action corresponding to the chosen enabled rule.

A *local state* of a node is the value of the local variables of the node and the state of its program counter. A *configuration* of the system $G = (V, E)$ is the cross product of the local states of all nodes in the system. The transition from a configuration to the next one is produced by the execution of an action at a node. A *computation* of the system is defined as a *weakly fair, maximal* sequence of configurations, $e = (c_0, c_1, \dots, c_i, \dots)$, where each configuration c_{i+1} follows from c_i by the execution of a single action of at least one node. During an execution step, one or more processors execute an action and a processor may take at most one action. *Weak fairness* of the sequence means that if any action in G is continuously enabled along the sequence, it is eventually chosen for execution. *Maximality* means that the sequence is either infinite, or it is finite and no action of G is enabled in the final global state.

In the sequel we consider the system can start in any configuration. That is, the local state of a node can be corrupted. Note that we don't make any assumption on the bound of corrupted nodes. In the worst case all the nodes in the system may start in a corrupted configuration. In order to tackle these faults we use self-stabilization techniques.

Definition 1 (self-stabilization). Let \mathcal{L}_A be a non-empty legitimacy predicate³ of an algorithm A with respect to a specification predicate $Spec$ such that every configuration satisfying \mathcal{L}_A satisfies $Spec$. Algorithm A is self-stabilizing with respect to $Spec$ iff the following two conditions hold:

² Observe that the two self-stabilizing MST algorithms mentioned in the Previous Work section assume that the nodes have distinct IDs with no distinguished nodes. Nevertheless, if the nodes have distinct IDs then it is possible to elect one node as a leader in a self-stabilizing manner. Conversely, if there exists one distinguished node in an anonymous network, then it is possible to assign distinct IDs to the nodes in a self-stabilizing manner [8]. Note that it is not possible to compute deterministically an MST in a fully anonymous network (i.e., without any distinguished node), as proved in [18].

³ A legitimacy predicate is defined over the configurations of a system and is an indicator of its correct behavior.

(i) Every computation of \mathcal{A} starting from a configuration satisfying $\mathcal{L}_{\mathcal{A}}$ preserves $\mathcal{L}_{\mathcal{A}}$ (closure).

(ii) Every computation of \mathcal{A} starting from an arbitrary configuration contains a configuration that satisfies $\mathcal{L}_{\mathcal{A}}$ (convergence).

We define below a *loop-free* configuration of a system as a configuration which contains paths with no cycle between any couple of nodes in the system. Given two nodes $u, v \in V$, we note $P(u, v)$ the path between u and v .

Definition 2 (Loop-Free Configuration). Let $Cycle(u, v)$ be the following predicate defined for two nodes $u, v \in V$ on configuration C : $Cycle(u, v) \equiv \exists P(u, v), P(v, u) : P(u, v) \cap P(v, u) = \emptyset$.

A *loop-free configuration* is a configuration of the system which satisfies: $\forall u, v \in V, Cycle(u, v) = false$.

We use the definition of a loop-free configuration to define a *loop-free stabilizing* system.

Definition 3 (Loop-Free Stabilization). A distributed system is called *loop-free stabilizing* if and only if it is self-stabilizing and there exists a non-empty set of configurations such that the following conditions hold: (i) Every computation starting from a loop-free configuration reaches a loop-free configuration (closure). (ii) Every computation starting from an arbitrary configuration contains a loop-free configuration (convergence).

In the sequel we study the loop-free self-stabilizing **LoopFreeMST** problem. The legitimacy predicate $\mathcal{L}_{\mathcal{A}}$ for the **LoopFreeMST** problem is the conjunction of the following two predicates: (i) a tree T spanning the network is constructed. (ii) T is a minimum spanning tree of G (i.e., $\forall T', W(T) \leq W(T')$, with T' be a spanning tree of G and $W(S) = \sum_{e \in S} w(e)$ be the cost of the subgraph S).

3 The Algorithm LoopFreeMST

In this section, we describe our self-stabilizing algorithm for the MST problem. We call this algorithm **LoopFreeMST**. Let us begin by an informal description of **LoopFreeMST** aiming at underlining its main features.

3.1 High Level Description

LoopFreeMST is based on the red rule. That is, for constructing an MST, the algorithm successively deletes the edges of maximum weight within every cycle. For this purpose, a spanning tree is maintained, together with a pre-order labeling of its nodes. Given the current spanning tree T maintained by our algorithm, every edge e of the graph that is not in the spanning tree creates a unique cycle in the graph when added to T . This cycle is called *fundamental cycle*, and is denoted by C_e . (Formally, this cycle depends on T ; Nevertheless no confusion should arise from omitting T in the notation of C_e). If $w(e)$ is not the maximum weight of all the edges in C_e , then, according to the red rule, our algorithm swaps e with the edge f of C_e with maximum weight. This swapping

procedure is called an *improvement*. A straightforward consequence of the red rule is that if no improvements are possible then the current spanning tree is a minimum one.

Algorithm **LoopFreeMST** can be decomposed in three procedures: (i) Tree construction, (ii) Token label circulation, (iii) Cycle improvement.

The latter procedure (Cycle improvement) is in fact the core of our contribution. Indeed, the two first procedures are simple modifications of existing self-stabilizing algorithms, one for building a spanning tree, and the other for labelling its nodes. We will show how to compose the original procedure “Cycle improvement” with these two existing procedures. Note that “Cycle improvement” differs from the previous self-stabilizing implementation of the improvement swapping in [19] by the fact that it does not require any a priori knowledge of the network, and it is loop-free.

LoopFreeMST starts by constructing a spanning tree of the graph, using the self-stabilizing loop-free algorithm “Tree construction” described in [21]. The two other procedures are performed concurrently. A token circulates along the edges of the current spanning tree, in a self-stabilizing manner. This token circulation uses algorithms proposed in [5,24] as follows. A non-tree-edge can belong to at most one fundamental cycle, but a tree-edge can belong to several fundamental cycles. Therefore, to avoid simultaneous possibly conflicting improvements, our algorithm considers the cycles in order. For this purpose, the token labels the nodes of the current tree in a DFS order (pre-order). This labeling is then used to find the unique path between two nodes in the spanning tree in a distributed manner, and enables computing the fundamental cycle resulting from adding one edge to the current spanning tree.

We now sketch the description of the procedure “Cycle improvement” (see Figure 1). When the token arrives at a node u in a state **Done**, it checks whether u has some incident edges not in the current spanning tree T connecting u with some other node v with smaller label. If it is the case, then enters state **Verify**. Let $e = \{u, v\}$. Node u then initiates a traversal of the fundamental cycle C_e for finding the edge f with maximum weight in this cycle. If $w(f) = w(e)$ then no improvement is performed. Else an improvement is possible, and u enters State **Improve**. Exchanging e and f in T results in a new tree T' . When the improvement is terminated u enters in State **end**. The key issue here is to perform this exchange in a loop-free manner. Indeed, one cannot be sure that

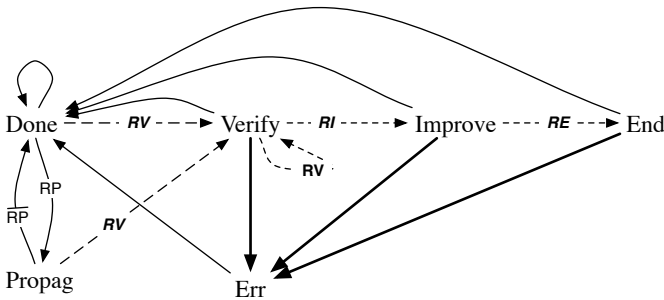


Fig. 1. Evolution of the node’s state in cycle improvement module. Rule R_D is depicted in plain. Rule R_{Err} is depicted in bold.

two modifications of the current tree (i.e., removing f from T , and adding e to T) that are applied at two distant nodes will occur simultaneously. And if they do not occur simultaneously, then there will a time interval during which the nodes will not be connected by a spanning tree. Our solution for preserving loop-freeness relies on a sequence of successive local and atomic changes, involving a single variable. This variable is a pointer to the current parent of a node in the current spanning tree. To get the flavor of our method, let us consider the example depicted on Figure 2. In this example, our algorithm has to exchange the edge $e = \{10, 12\}$ of weight 9, with the edge $f = \{7, 8\}$ of weight 10 (Figure 2(a)). Currently, the token is at node 12. The improvement is performed in two steps, by a sequence of two local changes. First, node 10 switches its parent from 8 to 12 (Figure 2(b)). Next, node 8 switches its parent from 7 to 10 (Figure 2(c)). A spanning tree is preserved at any time during the execution of these changes.

Note that any modification of the spanning tree makes the current labeling globally inaccurate, i.e., it is not necessarily a pre-order anymore. However, the labeling remains a pre-order in the portion of the tree involved in the exchange. For instance, consider again the example depicted on Figure 2(c). When the token will eventually reach node A , it will label it by some label $\ell > 12$. The exchange of $e = \{10, 12\}$ and $f = \{7, 8\}$ has not changed the pre-order for the fundamental cycle including edge $\{A, 12\}$. However, when the token will eventually reach node B and label it $\ell' > \ell$, the exchange of $e = \{10, 12\}$ and $f = \{7, 8\}$ has changed the pre-order for the fundamental cycle including edge $\{B, 9\}$: the parent of node labeled 10 is labeled 12 whereas it should have a label smaller than 10 in a pre-order. When the pre-order is modified by an exchange, the inaccurately labeled node changes its state to `Err`, and stops the traversal of the fundamental cycle. The token is then informed that it can discard this cycle, and carry on the traversal of the tree.

3.2 Detailed Level Description

We now enter into the details of Algorithm `LoopFreeMST`. First, let us state all variables used by the algorithm. Later on, we will describe its predicates and its rules.

Variables. For any node $v \in V(G)$, we denote by $N(v)$ the set of all neighbors of v in G . Algorithm `LoopFreeMST` maintains the set $N(v)$ at every node v . We use the following notations:

- `parentv`: the parent of v in the current spanning tree;
- `labelv`: the integer label assigned to v ;
- `dv`: the distance (in hops) from v to the root in the current spanning tree;
- `statev`: the state of node v , with values in `{Done, Verify, Improve, End, Propag, Err}`⁴;
- `DefCyclev`: Let C_e the current fundamental cycle with $e = \{x, y\}$, `DefCyclev = (x, y)`.
- `VarCyclev`: a pair of variables: the first one is the maximum edge-weight in the current fundamental cycle; the second one is a (boolean) variable in `{Before, After}`⁵;
- `sucv`: the successor of v in the current fundamental cycle.

⁴ The state `Propag` is detailed in Consistency rules.

⁵ For details see paragraph 3.2 Cycle improvement rules.

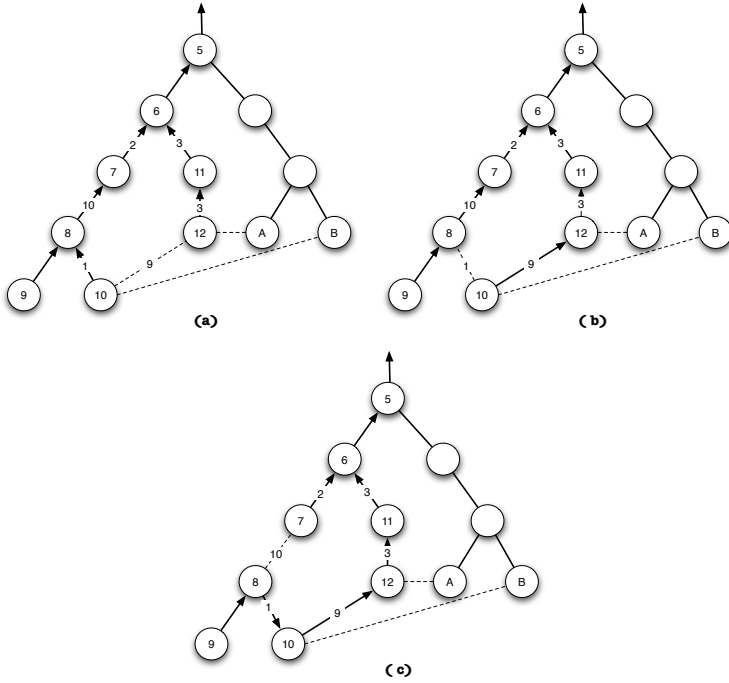


Fig. 2. Example of a loop-free improvement of the current spanning tree. The direction of the edges indicates the parent relation. Edges in the spanning tree are depicted as plain lines; Edges not in the spanning tree are denoted by dotted lines.

Consistency rules. The first task executed by LoopFreeMST is to check the consistency of the variables of each node see Figure 1. Done is the standard state of a node when this node does not have the token, or is not currently visited by the traversal of a fundamental cycle. When the variables of a node are detected to be not coherent, the state of the node becomes Err thanks to rule R_{Err} . There is one predicate in R_{Err} for each state, except for state Propag, to check whether the variables of the node are consistent (see Figure 3). The rule R_D allows the node to return to the standard state Done. More precisely, rule R_D resets the variables, and stops the participation of the node to any improvement.

R_{Err} : (Bad label)

If $CoherentCycle(v) \wedge Error(v) \wedge DefCycle[0]_v \neq label_v \wedge EndPropag(v)$
then $state_v := Err$;

R_D : (Improvement consistency)

If $\neg CoherentCycle(v) \wedge EndPropag(v)$
then $state_v := Done$; $DefCycle_v := (label_v, done)$; $VarCycle_v := (0, Before)$;
 $suc_v := \emptyset$;

Tree construction. LoopFreeMST starts by constructing a spanning tree of the graph, using the self-stabilizing loop-free algorithm “Tree construction” described in [21].

$\text{CoherentCycle}(v) \equiv \text{Coherent_Done}(v) \vee \text{Coherent_Verify}(v) \vee \text{Coherent_Improve}(v) \vee \text{Coherent_End}(v) \vee \text{Coherent_Error}(v)$ $\text{Coherent_Done}(v) \equiv \text{state}_v = \text{Done} \wedge \text{succ}_v = \emptyset \wedge \text{DefCycle}_v = (\text{label}_v, \text{done}) \wedge \text{VarCycle}_v = (0, \text{Before})$ $\text{Coherent_Verify}(v) \equiv \text{state}_v = \text{Verify} \wedge \text{succ}_v = \text{Succ}(v) \wedge [(\text{Init}(v) \wedge \text{VarCycle}_x = (0, \text{Before})) \vee \text{Nds_Verify}(v)]$ $\text{Coherent_Improve}(v) \equiv \text{state}_v = \text{Improve} \wedge \text{succ}_v = \text{Succ}(v) \wedge \text{DefCycle}_v = \text{DefCycle}_{\text{parent}_v} \wedge \text{VarCycle}_v = \text{VarCycle}_{\text{parent}_v}$ $\text{Coherent_End}(v) \equiv \text{state}_v = \text{End} \wedge \text{DefCycle}_v = \text{DefCycle}_{\text{parent}_v} \wedge (\text{NdDel}(v) \vee \text{Ask_EI}(v))$ $\text{Coherent_Error}(v) \equiv \text{state}_v = \text{Err} \wedge (\text{succ}_v = \text{Succ}(v) = \emptyset \vee \text{Ask_E}(v)) \wedge \text{DefCycle}_v = \text{DefCycle}_{\text{Pred}(v)}$ $\text{CoherentTree}(v)^a \equiv (v = r \wedge \text{d}_v = 0 \wedge \text{st}_v = N) \vee (v \neq r \wedge \text{d}_v = \text{d}_{\text{parent}_v} + 1 \wedge \text{st}_v = N \wedge \text{rw}_v = \text{d}_v) \vee \text{state}_{\text{parent}_v} = \text{Improve} \vee \text{state}_{\text{parent}_v} = \text{Propag}$ $\text{Ask_V}(v) \equiv \text{state}_{\text{Pred}(v)} = \text{Verify}$ $\text{Ask_I}(v) \equiv (\text{state}_{\text{Pred}(v)} = \text{Improve} \wedge \text{VarCycle}[1]_{\text{Pred}(v)} = \text{Before}) \vee (\text{state}_{\text{succ}_v} = \text{Improve} \wedge \text{VarCycle}[1]_{\text{succ}_v} = \text{After})$ $\text{Ask_EI}(v) \equiv (\exists u \in N(v), \text{parent}_u = v \wedge \text{state}_u = \text{End} \wedge \text{DefCycle}_u = \text{DefCycle}_v)$ $\text{Ask_E}(v) \equiv \text{succ}_v \neq \emptyset \wedge \text{state}_{\text{succ}_v} = \text{Err} \wedge \text{DefCycle}_v = \text{DefCycle}_{\text{succ}_v}$ $\text{Tree_Edge}(v, u) \equiv \text{parent}_v = u \vee \text{parent}_u = v$ $\text{C_Ancestor}(v) \equiv \text{parent}_v \neq \text{succ}_v \wedge \text{parent}_v \neq \text{Pred}(v)$ $\text{Init}(v) \equiv \text{DFS_F}(v) \wedge \text{DefCycle}[0]_v = \text{label}_v$ $\text{Nds_Verify}(v) \equiv [(\text{Ask_V}(v) \wedge \text{VarCycle}_v = (\text{Max_C}(v), \text{Way_C}(v))) \vee \text{Ask_I}(v)] \wedge \text{DefCycle}_v = \text{DefCycle}_{\text{Pred}(v)}$ $\text{NdDel}(v) \equiv \text{state}_{\text{parent}_v} \neq \text{Done} \wedge \text{state}_{\text{parent}_v} \neq \text{Propag} \wedge \neg \text{Improve}(v)$

^a In [21], variable st_v indicates if v propagates a new distance (state P) or not (state N), and rw_v is used to propagate the new distance in the tree.

Fig. 3. Corrections predicates used by the algorithm

This algorithm constructs a BFS, and uses two variables $parent$ and $distance$. During the execution of our algorithm, these two variables are subject to the same rules as in [21]. After each modification of the spanning tree, the new distance to the parent is propagated in sub-trees by Rules R_P and \bar{R}_P .

R_P : (Distance propagation)

If $\text{Coherent_Done}(v) \wedge \neg \text{Ask_V}(v) \wedge \text{succ}_v \neq \text{parent}_v \wedge \text{Pred}(v) \neq \text{parent}_v \wedge \text{d}_v \neq \text{d}_{\text{parent}_v} + 1 \wedge (\text{state}_{\text{parent}_v} = \text{Improve} \vee \text{state}_{\text{parent}_v} = \text{Propag})$
 then $\text{state}_v := \text{Propag}; \text{d}_v := \text{d}_{\text{parent}_v} + 1;$

\bar{R}_P : (End distance propagation)

If $\text{state}_v = \text{Propag} \wedge \text{EndPropag}(v)$
 then $\text{state}_v := \text{Done}; \text{DefCycle}_v := (\text{label}_v, \text{done}); \text{VarCycle}_v := (0, \text{Before});$
 $\text{succ}_v := \emptyset;$

Token circulation and pre-order labeling. LoopFreeMST uses the algorithm described in [5] to provide each node v with a label label_v . Each label is unique in the network traversed by the token. This labeling is used to find the unique path between two nodes in the spanning tree, in a distributed manner. For this purpose, we use the snap-stabilizing algorithm described in [24] for the circulation of a token in the spanning tree (a snap-stabilizing algorithm stabilizes in 0 steps thus algorithm in [24] allows to always have a correct token circulation). We have slightly modified this algorithm because LoopFreeMST stops the token circulation at a node during the “Cycle improvement” procedure. A node v knows if it has the token by applying predicate $\text{Init}(v)$ (Predicate $\text{DFS_F}(v)$ is true at node v if the token was forwarded by its parent). Rule R_{DFS} guides the circulation of the token. The token carries on its tree traversal if one of the following three conditions is satisfied: (i) there is no improvement which could be initiated by the node which holds the token, (ii) an improvement was performed in the current cycle, or (iii) inconsistent node labels were detected in the current cycle. The latter is under the control of Predicate $\text{ContinueDFS}(v)$.

R_{DFS} : (**Continue DFS token circulation**)

If $\text{CoherentCycle}(v) \wedge \text{Init}(v) \wedge \text{ContinueDFS}(v)$
then $\text{state}_v := \text{Done}; \text{DefCycle}[1]_v = \text{done};$

Cycle improvement rules. The procedure “Cycle improvement” is the core of LoopFreeMST. Its role is to avoid disconnection of the current spanning tree, while successively improving the tree until reaching an MST. The procedure can be decomposed in four tasks: (1) to check whether the fundamental cycle of the non-tree edge has an improvement or not, (2) perform the improvement if any, (3) update the distances, and (4) resume the token circulation.

Let us start by describing the first task. A node u in state `Done` changes its state to `Verify` if its variables are in consistent state, it has a token, and it has identified a candidate (i.e., an incident non-tree edge $e = \{u, v\}$ whose other extremity v has a smaller label than the one of u). The latter is under the control of Predicate $\text{InitVerify}(v)$, and the variable VarCycle_v contains the label of u and v . If the three conditions are satisfied, then the verification of the fundamental cycle C_e is initiated from node u , by applying rule R_V . The goal of this verification is twofold: first, to verify whether C_e exists or not, and, second, to save information about the maximum edge weight and the location of the edge of maximum weight in C_e . These information are stored in the variable $\text{Way_C}(v)$. In order to respect the orientation in the current spanning tree, the node u or v that initiates the improvement depends on the localization of the maximum weight edge f in C_e . More precisely, let r be the least common ancestor of nodes u and v in the current tree. If f occurs before r in T in the traversal of C_e from u starting by edge (u, v) , then the improvement starts from u , otherwise the improvement starts from v . To get the flavor of our method, let us consider the example depicted on Figure 2. In this example, f occurs after the least common ancestor (node 6). Therefore node 10 atomically swaps its parent to respect the orientation. However, if one replaces in the same example the weight of edge $\{11, 6\}$ by 11 instead of 3, then f would occur before r , and thus node 12 would have to atomically swap its parent. The relative places of f and r in the cycle is indicated by Predicate $\text{Way_C}(v)$ that

$$\begin{array}{l}
\text{Pred}(v) \equiv \arg \min \{ \text{label}_u : u \in N(v) \wedge \text{state}_u \neq \text{Done} \wedge \text{state}_u \neq \text{Propag} \wedge \text{succ}_u = v \} \\
\quad \text{if } u \text{ exists, } \emptyset \text{ otherwise} \\
\text{MaxLab}(v, x) \equiv \arg \max \{ \text{label}_s : s \in N(v) \wedge \text{label}_s < x \} \\
\text{Succ}(v) \equiv \begin{cases} \text{VarCycle}[0]_v & \text{if } \text{DefCycle}[1]_v = \text{label}_v \\ \text{parent}_v & \text{if } (\text{label}_v > \text{DefCycle}[1]_v \wedge \text{state}_v = \text{Verify}) \vee \\ & (\text{label}_v < \text{DefCycle}[1]_v \wedge \\ & (\text{state}_v = \text{Improve} \vee \text{state}_v = \text{End})) \\ \text{MaxLab}(v, \text{DefCycle}[1]_v) & \text{if } (\text{label}_v < \text{DefCycle}[1]_v \wedge \text{state}_v = \text{Verify}) \\ \text{MaxLab}(v, \text{label}_v) & \text{if } (\text{label}_v > \text{DefCycle}[1]_v \wedge \\ & (\text{state}_v = \text{Improve} \vee \text{state}_v = \text{End})) \end{cases} \\
\text{Max}_C(v) \equiv \max \{ \text{VarCycle}[0]_{\text{Pred}(v)}, w(v, \text{Pred}(v)) \} \\
\text{Way}_C(v) \equiv \begin{cases} \text{After} & \text{if } \text{VarCycle}[0]_v \neq \text{VarCycle}[0]_{\text{Pred}(v)} \wedge \text{label}_v > \text{label}_{\text{Pred}(v)} \\ \text{VarCycle}[1]_{\text{Pred}(v)} & \text{otherwise} \end{cases} \\
\text{LabCand}(v) \equiv \min \{ \text{label}_u : u \in N(v) \wedge \text{label}_u < \text{label}_v \wedge \neg \text{Tree_Edge}(v, u) \wedge \\ \quad \text{label}_u \succ \text{DefCycle}[1]_v \}^a \quad \text{if } u \text{ exists, end otherwise} \\
\hline
^a \succ \text{ order on neighbor labels for which 'end' is the biggest element and 'done' is the smallest} \\
\text{one.}
\end{array}$$

Fig. 4. Predicates used by the algorithm

returns two different values: **Before** or **After**. During the improvement of the tree, the fundamental cycle is modified. It is crucial to save information about this cycle during this modification. In particular, the successor of a node w in a cycle, stored in the variable succ_w , must be preserved. Its value is computed by Predicate **Succ**(v) which uses node labels to identify the current examined fundamental cycle. Each node is able to compute its predecessor in the fundamental cycle by applying Predicate **Pred**(v). The state of a node is compared with the ones of its successor and predecessor to detect potential inconsistent values. At the end of this task, the node u learns the maximum weight of the cycle C_e and can decide whether it is possible to make an improvement or not. If not, but there is another non-tree edge e' that is candidate for potential replacement, then u verifies $C_{e'}$. Otherwise the token carries on its traversal, and rule \bar{R}_P is applied.

R_V : (Verify rule)

```

If  $\text{CoherentCycle}(v) \wedge \neg \text{Error}(v) \wedge (\text{InitVerify}(v) \vee [\neg \text{Init}(v) \wedge (\text{Coherent\_Done}(v) \vee \text{state}_v = \text{Propag}) \wedge \text{Ask}_V(v)])$ 
then  $\text{state}_v := \text{Verify}$ ;
    If  $\text{DFS\_F}(v)$  then  $\text{DefCycle}[1]_v := \text{LabCand}(v)$ ;
    Else  $\text{DefCycle}_v := \text{DefCycle}_{\text{Pred}(v)}$ ;  $\text{VarCycle}_v := (\text{Max}_C(v), \text{Way}_C(v))$ ;
     $\text{succ}_v := \text{Succ}(v)$ ;

```

If C_e can yield an improvement, then rule R_I is executed. By this rule, a node enters in state **Improve**, and changes its parent to its predecessor if $\text{VarCycle}[1]_v = \text{Before}$ (respectively to its successor if $\text{VarCycle}[1]_v = \text{After}$). For this purpose, it uses the variable succ_v and the predicate **Pred**(v).

R_I: (Improve rule)

If $\text{CoherentCycle}(v) \wedge \neg \text{Error}(v) \wedge \text{Coherent_Verify}(v) \wedge \text{Improve}(v) \wedge$
 $\neg \text{C_Ancestor}(v) \wedge [(\text{DFS_F}(v) \wedge \text{Ask_V}(v)) \vee \text{Ask_I}(v)]$
then $\text{state}_v := \text{Improve};$
If $\text{DFS_F}(v) \vee \text{state}_{\text{Pred}(v)} = \text{Improve}$ **then** $\text{VarCycle}_v := \text{VarCycle}_{\text{Pred}(v)}$
If $(\text{DFS_F}(v) \wedge \text{VarCycle}[1]_v = \text{Before}) \vee \neg \text{DFS_F}(v)$ **then** $\text{parent}_v := \text{Pred}(v);$
If $\text{state}_{\text{Suc}_v} = \text{Improve}$ **then** $\text{VarCycle}_v := \text{VarCycle}_{\text{Suc}_v}; \text{parent}_v := \text{Suc}_v;$
If $w(v, \text{Suc}_v) \geq \text{VarCycle}[0]_v$ **then** $\text{Suc}_v = \text{Succ}(v)$
 $d_v := d_{\text{parent}_v} + 1;$

At the end of an improvement, it is necessary to inform the node holding the token that it has to carry on its traversal. This is the role of rule R_E. It is also necessary to inform all nodes impacted by the modification that they have to update their distances to the root (see Section 3.2).

R_E: (End of improvement rule)

If $\text{CoherentCycle}(v) \wedge \neg \text{Error}(v) \wedge \text{End_Improve}(v) \wedge \text{EndPropag}(v)$
then $\text{state}_v := \text{End};$

Module composition. All the different modules presented, except the tree construction parts of the correction module, need the presence of a spanning tree in G . Thus, we must execute the tree construction rules first if an incoherency in the spanning tree is detected. To this end, these rules are composed using the level composition defined in [16], i.e., if Predicate $\text{CoherentTree}(v)$ (see Fig. 3) is not verified then the tree construction rules are executed, otherwise the other modules can be executed. The token circulation algorithm and the naming algorithm are composed together using the conditional composition described in [5], i.e., the naming algorithm is executed when a logical expression (based on guards of token circulation algorithm) is true. Finally, we compose the token circulation algorithm and the cycle improvement module with a conditional composition using Predicate $\text{ContinueDFS}(v)$ (see Fig. 5). This allows to execute the token circulation algorithm only if the cycle improvement module does not need the token on a node. Figure 6 shows how the modules are composed together.

<p> $\text{Candidate}(v) \equiv \text{LabCand}(v) \neq \text{end}$ $\text{InitVerify}(v) \equiv \text{Init}(v) \wedge \text{Candidate}(v) \wedge (\text{Coherent_Done}(v) \vee [\text{Coherent_Verify}(v) \wedge$ $\neg \text{Improve}(v) \wedge \neg \text{C_Ancestor}(v) \wedge \text{Ask_V}(v)])$ $\text{ImproveF}(v, x) \equiv \neg \text{Tree_Edge}(v, x) \wedge \max(\text{VarCycle}[0]_v, \text{VarCycle}[0]_x) > w(v, x)$ $\text{Improve}(v) \equiv \text{ImproveF}(v, \text{Pred}(v)) \vee \text{ImproveF}(v, \text{Suc}_v)$ $\text{End_Improve}(v) \equiv \text{Coherent_Improve}(v) \wedge (\text{NdDel}(v) \vee \text{Ask_EI}(v))$ $\text{ContinueDFS}(v) \equiv (\text{Init}(v) \wedge [([\text{Coherent_Done}(v) \vee (\text{Coherent_Verify}(v) \wedge$ $\neg \text{ImproveF}(v, \text{Pred}(v)) \wedge \text{Ask_V}(v)]) \wedge \neg \text{Candidate}(v)] \vee$ $\text{Coherent_End}(v) \vee \text{Error}(v)]) \vee \neg \text{DFS_F}(v)$ $\text{Error}(v) \equiv \text{state}_v \neq \text{Done} \wedge \text{state}_v \neq \text{Err} \wedge (\text{Suc}_v = \text{Succ}(v) = \emptyset \vee \text{Ask_E}(v))$ $\text{EndPropag}(v) \equiv (\forall u \in N(v), \text{parent}_u = v \wedge \text{state}_u = \text{Done} \wedge d_u = d_v + 1)$ </p>

Fig. 5. Predicates used by the algorithm



Fig. 6. Composition of the presented modules

3.3 Complexity

Definition 4 (Red Rule). If C is a cycle in $G = (V, E)$ with no red edges then color in red the maximum edge weight in C .

Theorem 1 (Tarjan et al. [26]). Let G be a connected graph. If it is not possible to apply Red Rule then the set of not colored edges forms a minimum spanning tree of G .

Lemma 1. Starting from a configuration where an arbitrary spanning tree is constructed, in at most $O(mn)$ rounds the cycle improvement module produces a minimum spanning tree of G , with respectively m and n the number of edges and nodes of the network G .

Proof. In a given network $G = (V, E)$, if a spanning tree of G is constructed then there are exactly $m - (n - 1)$ fundamental cycles in G since there are $n - 1$ edges in any spanning tree of G . Thus, a tree edge can be contained in at most $m - n + 1$ fundamental cycles. Consider a configuration where a spanning tree T of G is constructed and a tree edge e_0 is contained in $m - n + 1$ fundamental cycles and all tree edges have a weight equal to 1, except e_0 of weight $w(e_0) > 1$. Suppose that T is not a minimum spanning tree of G such that $\forall e_i \in E, i = 1, \dots, m - n + 1, w(e_{i-1}) > w(e_i)$ with $e_0 \in T$ and $\forall i = 1, \dots, m - n + 1, e_i \notin T$ and $w(e_i) > 1$ (see the graph of Figure 7(a)). Consider the following sequence of improvements: $\forall i, i = 1, \dots, m - n + 1$, exchange the tree edge e_{i-1} by the not tree edge e_i (see a sequence of improvements in Figure 7). In this sequence, we have exactly $m - n + 1$ improvements and this is the maximum number of improvements to obtain a minimum spanning tree since there are $m - n + 1$ fundamental

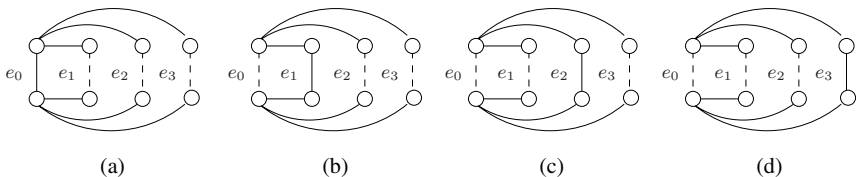


Fig. 7. (a) a spanning tree with plain lines in a graph with $m - n + 1$ improvements, (b) the spanning tree obtained after the first improvement, (c) the spanning tree obtained after the second improvement, (d) the minimum spanning tree of the graph obtained after the third improvement

cycles and for each one we apply the Red rule (see Definition 4 and Theorem 1). An improvement can be initiated in the cycle improvement module by a node with the DFS token. The DFS token performs a tree traversal in $O(n)$ rounds. Moreover, each improvement needs to cross a cycle a constant number of times and each cross requires $O(n)$ rounds. Since at most $m - n + 1$ improvements are needed to obtain a minimum spanning tree, at most $O(mn)$ rounds are needed to construct a minimum spanning tree.

Lemma 2. *Starting from a legitimate configuration, after a weight edge modification the system reaches a legitimate configuration in at most $O(mn)$ rounds.*

Proof. After a weight edge change the system is no more in a legitimate configuration in the following cases: (1) the weight of a not tree edge is less than the weight of the heaviest tree edge in its fundamental cycle, or (2) the weight of a tree edge is greater than the weight of a not tree edge in one of the fundamental cycles including the tree edges. In each case above, the algorithm must verify if improvements must be performed to reach again a legitimate configuration, otherwise the system is still in a legitimate configuration. Thus, in case (1) it is only sufficient to verify if an improvement must be performed in the fundamental cycle associated to the not tree edge (i.e. to apply the Red rule a single time). To this end, its fundamental cycle must be crossed at most three times: the first time to verify if an improvement is possible, a second time to perform the improvement and a last time to end the improvement, each one needs at most $O(n)$ rounds. Case (2) is more complicated, indeed the weight of a tree edge can change which leads to a configuration where at most $m - n + 1$ improvements must be performed to reach a legitimate configuration, since a tree edge can be contained in at most $m - n + 1$ fundamental cycles as described in proof of Lemma 1. Since each improvement phase needs $O(n)$ rounds (see case (1)) at most $O(mn)$ rounds are needed to reach a legitimate configuration. The complexity of case (2) dominates the complexity of the first case. Therefore, after a weight edge change at most $O(mn)$ rounds are needed to reach a legitimate configuration.

Note that the presented algorithm uses only a constant number of variables of size $O(\log n)$. Therefore, $O(\log n)$ bits of memory are needed at each node to execute the algorithm. Moreover, due to space constraints correctness proof are given in [27].

4 Concluding Remarks

We presented a new solution to the distributed MST construction that is both self-stabilizing and loop-free. It improves on memory usage from $O(n \log n)$ to $O(\log n)$, yet doesn't make strong system assumptions such as knowledge of network size or unicity of edge weights, making it particularly suited to dynamic networks. Two important open questions are raised:

1. For depth first search tree construction, self-stabilizing solutions that use only constant memory space do exist. It is unclear how the obvious constant space lower bound can be raised with respect to metrics that minimize a global criterium (such as MST).

2. Our protocol pioneers the design of self-stabilizing loop-free protocols for *non* locally optimizable tree metrics. We expect the techniques used in this paper to be useful to add loop-free property for other metrics that are only globally optimizable, yet designing a generic such approach is a difficult task.

References

1. Katz, S., Perry, K.J.: Self-stabilizing extensions for message-passing systems. *Distributed Computing* 7, 17–26 (1993)
2. Anagnostou, E., Hadzilacos, V.: Tolerating transient and permanent failures (extended abstract). In: Schiper, A. (ed.) *WDAG 1993*. LNCS, vol. 725, pp. 174–188. Springer, Heidelberg (1993)
3. Ben-Or, M., Dolev, D., Hoch, E.N.: Fast self-stabilizing byzantine tolerant digital clock synchronization. In: Bazzi, R.A., Patt-Shamir, B. (eds.) *PODC*, pp. 385–394. ACM Press, New York (2008)
4. Cobb, J.A., Gouda, M.G.: Stabilization of general loop-free routing. *J. Parallel Distrib. Comput.* 62(5), 922–944 (2002)
5. Datta, A.K., Gurumurthy, S., Petit, F., Villain, V.: Self-stabilizing network orientation algorithms in arbitrary rooted networks. *Stud. Inform. Univ.* 1(1), 1–22 (2001)
6. Dijkstra, E.W.: Self-stabilizing systems in spite of distributed control. *Commun. ACM* 17(11), 643–644 (1974)
7. Dolev, S.: *Self-stabilization*. MIT Press, Cambridge (2000)
8. Dolev, S.: *Self-Stabilization*. MIT Press, Cambridge (2000)
9. Dolev, S., Herman, T.: *Superstabilizing protocols for dynamic distributed systems*. Chicago J. Theor. Comput. Sci. (1997)
10. Dolev, S., Welch, J.L.: Wait-free clock synchronization. *Algorithmica* 18(4), 486–511 (1997)
11. Dolev, S., Welch, J.L.: Self-stabilizing clock synchronization in the presence of byzantine faults. *J. ACM* 51(5), 780–799 (2004)
12. Gafni, E.M., Bertsekas, P.: Distributed algorithms for generating loop-free routes in networks with frequently changing topology. *IEEE Transactions on Communications* 29, 11–18 (1981)
13. Gallager, R.G., Humblet, P.A., Spira, P.M.: A distributed algorithm for minimum-weight spanning trees. *ACM Trans. Program. Lang. Syst.* 5(1), 66–77 (1983)
14. Garcia-Luna-Aceves, J.J.: Loop-free routing using diffusing computations. *IEEE/ACM Trans. Netw.* 1(1), 130–141 (1993)
15. Gopal, A.S., Perry, K.J.: Unifying self-stabilization and fault-tolerance (preliminary version). In: *PODC*, pp. 195–206 (1993)
16. Gouda, M.G., Herman, T.: Adaptive programming. *IEEE Trans. Software Eng.* 17(9), 911–921 (1991)
17. Gouda, M.G., Schneider, M.: Stabilization of maximal metric trees. In: Arora, A. (ed.) *WSS*, pp. 10–17. IEEE Computer Society Press, Los Alamitos (1999)
18. Gupta, S.K.S., Srimani, P.K.: Self-stabilizing multicast protocols for ad hoc networks. *J. Parallel Distrib. Comput.* 63(1), 87–96 (2003)
19. Higham, L., Liang, Z.: Self-stabilizing minimum spanning tree construction on message-passing networks. In: Welch, J.L. (ed.) *DISC 2001*. LNCS, vol. 2180, pp. 194–208. Springer, Heidelberg (2001)
20. Johnen, C., Tixeuil, S.: Route Preserving Stabilization. In: Huang, S.-T., Herman, T. (eds.) *SSS 2003*. LNCS, vol. 2704, pp. 184–198. Springer, Heidelberg (2003)
21. Johnen, C., Tixeuil, S.: Route preserving stabilization. In: *Self-Stabilizing Systems*, pp. 184–198 (2003)

22. Kruskal, J.B.: On the shortest spanning subtree of a graph and the travelling salesman problem. *Proc. Amer. Math. Soc.* 7, 48–50 (1956)
23. Papatriantafyllou, M., Tsigas, P.: On self-stabilizing wait-free clock synchronization. *Parallel Processing Letters* 7(3), 321–328 (1997)
24. Petit, F., Villain, V.: Optimal snap-stabilizing depth-first token circulation in tree networks. *J. Parallel Distrib. Comput.* 67(1), 1–12 (2007)
25. Prim, R.C.: Shortest connection networks and some generalizations. *Bell System Tech. J.*, 1389–1401 (1957)
26. Sleator, D.D., Tarjan, R.E.: A data structure for dynamic trees. *J. Comput. Syst. Sci.* 26(3), 362–391 (1983)
27. Blin, L., Potop-Butucaru, M.G., Rovedakis, S., Tixeuil, S.: A new self-stabilizing minimum spanning tree construction with loop-free property. Research Report, inria-00384041, INRIA (2009)