# Distributed Computing on Core-Periphery Networks: Axiom-based Design

Chen Avin[a,1,2], Michael Borokhovich[a,2], Zvi Lotker[a,2], David Peleg[b,2]

[a]*Ben-Gurion University of the Negev, Israel.*
[b]*The Weizmann Institute, Israel.*

## Abstract

Inspired by social networks and complex systems, we propose a *core-periphery* network architecture that supports fast computation for many distributed algorithms and is robust and efficient in number of links. Rather than providing a concrete network model, we take an axiom-based design approach. We provide three intuitive and independent algorithmic axioms and prove that any network that satisfies all axioms enjoys an efficient algorithm for a range of tasks (such as MST, sparse matrix multiplication, and more). We also show the *minimality* of our axiom set: for networks that satisfy any subset of the axioms, the same efficiency cannot be guaranteed for *any* deterministic algorithm.

## 1. Introduction

A fundamental goal in distributed computing concerns finding network architectures that allow fast running times for various distributed algorithms, but at the same time are cost-efficient, in terms of minimizing the number of communication links between the machines and the amount of memory used by each processor.

For illustration, let's consider three basic network topologies: a star, a clique and a constant degree expander. The *star graph* has only a linear number of links, and can compute every computable function in one round of communication. But clearly, such an architecture has two major disadvantages: the memory requirements of the central node do not scale, and the network is not robust (in the sense that a failure of the central node is enough to disable the network). The *complete graph*, on the other hand, is very robust, and can support extremely high performance for tasks such as information dissemination,

---

distributed sorting and minimum spanning tree, to name a few [26, 24, 23]. Also, in a complete graph, the amount of memory used by a single processor is minimal. The main drawback of that architecture is the high number of links it uses. *Constant degree expanders* are a family of graphs that support efficient computation for many tasks. They also have linear number of links, and can effectively balance the workload between many machines. But the diameter of these graphs is lower bounded by $\Omega(\log n)$, which implies a similar lower bound on the time required for most of the interesting tasks one can consider.

Therefore, a natural question is whether there are other candidate topologies with guaranteed good performance. We are interested in the best compromise solution: a network on which distributed algorithms have low running times, memory requirements at each node are limited, the architecture is robust to link and node failures and the total number of links is minimized (preferably linear in the number of nodes).

To try to answer this question, we adopt in this paper an *axiomatic* approach to the design of efficient networks. In contrast to the direct approach to network design, which is based on providing a *concrete* type of networks (by deterministic or randomized construction) and showing its efficiency, the axiomatic approach attempts to abstract away the algorithmic requirements that are imposed on the concrete model. This allows one to isolate and identify the basic requirements that a network needs for a certain type of tasks. While there are obvious similarities between the traditional direct modeling approach and our axiom-based one, there are also some marked differences. Perhaps the main difference is that whereas the direct modeling approach focuses on rules governing the concrete *structural* properties of a network, our axiomatic approach relies on defining necessary *operational* or *algorithmic* properties required from the network, without committing to any a specific topology. This approach may allow us to abstract away many of the less significant details of the structure, and thus enable us to derive a simpler analysis of the essential properties of the structure under consideration, developing efficient algorithms for various purposes, relying solely on those abstract features.

A closely related distinction is that while the performance of distributed algorithms is usually expressed by specific *structural* network parameters (e.g., diameter, degree, etc.), the axioms proposed in this work are expressed in terms of desired *algorithmic* properties that the network should have.

Our axiomatic approach is also influenced by the recent concept of *Software Defined Networks* (SDN) [14]. The concept of SDN implies layering and abstractions in the networking control plane, which allows easy configuration of overlay architectures with a specific *behavior* (as opposed to specific *structure*). By only defining the desired behavior, we may allow many network implementations, as long as the desired behavior is maintained. Using our axioms, we are able to define the exact abstractions required for devising efficient distributed algorithms,thus decoupling them from specific topology details.

The axioms proposed in the current work are motivated and inspired by the *core-periphery* structure exhibited by many social networks and complex systems. A core-periphery network is a network structured of two distinct groups

of nodes, namely, a large, sparse and weakly connected group of nodes identified as the *periphery*, which is loosely organized around a small, cohesive and densely connected group identified as the *core*. Such a dichotomic structure appears in many domains of our life, and has been observed in many social organizations, including modern social networks [2]. It can also be found in urban and even global systems (e.g., in global economy, the wealthiest countries constitute the core, which is highly connected by trade and transportation routes) [15, 21, 19]. An analysis conducted in [29] for many software systems, revealed that $75-80\%$ of the systems examined possess a core-periphery structure. There are also peer-to-peer networks that use a similar hierarchical structure, e.g., FastTrack [25] and Skype [4], in which the supernodes can be viewed as the core, while the regular users constitute the periphery. Various client-server systems can also be thought of in these terms.

The vast presence of the core-periphery structure in our life suggests that it is a natural and effective design pattern. Consequently, we argue that a distributed network architecture based on it may support distributed algorithms that are easy to devise and reason about, given that we are used to thinking in "core-periphery terms" (for instance, in employing information processes based on collecting inputs from various peripheral sources to the core nodes, processing the data centrally, and then sending the results back to the periphery).

The main technical contribution of this paper is in proposing a minimal set of simple core-periphery-oriented axioms, and demonstrating that networks satisfying these axioms achieve efficient running time for various distributed computing tasks, while being able to maintain a linear number of edges and limited memory use. We identify three basic, abstract and conceptually simple (parameterized) properties, which turn out to be highly relevant to the effective interplay between core and periphery. For each of these three properties, we propose a corresponding axiom, which in our opinion captures some intuitive aspect of the desired behavior expected of a network based on a core-periphery structure. Let us briefly describe our three properties, along with their "real life" interpretation, technical formulation and associated axioms.

The three properties are: (i) balanced boundary between the core and periphery, (ii) clique-like structure of the core and (iii) fast convergecast from the periphery to the core. The first property (i) concerns the boundary between the core and the periphery. Drawing an analogy from the world of social networks, the core can be thought of as a highly influential group, that exerts its influence (in the form of instructions, opinions, or other means) on the periphery. The mechanism through which this is done is based on certain core nodes, which are each connected to many nodes in the periphery, and act as "*ambassadors*" of the core. Ambassadors serve as bidirectional channels, through which information flows into the core and influence flows from the core to the periphery. However, to be effective as an ambassador, the core node must maintain a balance between its interactions with the external periphery, and its interactions with the other core members, serving as its natural "support"; a core node that is significantly more connected to the periphery than to the core, becomes ineffective as a channel of influence. In distributed computing terms, a core node that has many

3

| Task | Running time on $\mathcal{CP}$ networks | Lower bounds | |
|---|---|---|---|
| | | All Axioms | Any 2 Axioms |
| MST * | $O(\log^2 n)$ | $\Omega(1)$ | $\tilde{\Omega}(\sqrt[4]{n})$ |
| Matrix transposition | $O(k)$ | $\Omega(k)$ | $\Omega(n)$ |
| Vector by matrix multiplication | $O(k)$ | $\Omega(k/\log n)$ | $\Omega(n/\log n)$ |
| Matrix multiplication | $O(k^2)$ | $\Omega(k^2)$ | $\Omega(n/\log n)$ |
| Rank finding | $O(1)$ | $\Omega(1)$ | $\Omega(n)$ |
| Median finding | $O(1)$ | $\Omega(1)$ | $\Omega(\log n)$ |
| Mode finding | $O(1)$ | $\Omega(1)$ | $\Omega(n/\log n)$ |
| Number of distinct values | $O(1)$ | $\Omega(1)$ | $\Omega(n/\log n)$ |
| Top $r$ ranked by areas | $O(r)$ | $\Omega(r)$ | $\Omega(r\sqrt{n})$ |

$k$ - maximum number of nonzero entries in a row or column. * - randomized algorithm

Table 1: Summary of algorithms for core-periphery networks.

connections to the periphery has to be able to distribute all the information it collected from them, to other core nodes. Hence the relevant property is having a *balanced boundary*: a set $S$ of nodes is said to have an $\alpha$-*balanced* boundary if for each of its nodes, the ratio between the sizes of its neighborhoods outside and inside $S$ is at most $O(\alpha)$. The corresponding Axiom $\mathcal{A}_B$ states that the core must have a $\Theta(1)$-balanced boundary.

The second property (ii) deals with the flow of information within the core. It is guided by the key observation that to be influential, the core must be able to accomplish fast information dissemination internally among its members. The extreme example of a dissemination-efficient network is the complete graph, so the core's efficiency in information flow should naturally be measured against that benchmark. Formally, a set of nodes is said to be a $\beta$-*clique emulator* if it can accomplish full communication (namely, message exchange between every pair of its members) in $\beta$ time (communication rounds). The corresponding Axiom $\mathcal{A}_E$ postulates that the core must be a $\Theta(1)$-clique emulator. Note that this requirement is stronger than just requiring the core to possess a dense interconnection subgraph, since the latter permits the existence of "bottlenecks", whereas the requirement dictated by the axiom disallows such bottlenecks.

The third and the last property (iii) focuses on the flow of information from the periphery to the core and measures its efficiency. The core-periphery structure of the network is said to be a $\gamma$-*convergecaster* if this data collection operation can be performed in time $\gamma$. The corresponding Axiom $\mathcal{A}_C$ postulates that information can flow from the periphery nodes to the core efficiently (i.e., in constant time), namely, the core and periphery must form a $\Theta(1)$-convergecaster. Note that one implication of this requirement is that the presence of periphery nodes that are far away from the core, or bottleneck edges that bridge between many periphery nodes and the core, is forbidden.

One may raise the "semi-philosophical" question whether the properties we defined should be referred to as "axioms". Our answer is that adopting the ax-

iomatic view yields the added benefit that it immediately raises the fundamental issues of minimality, independence and necessity, thus allowing us to carefully verify the role and usefulness of each property / axiom. Indeed, we partially address these issues. First, we establish the independence of our axioms, by showing that neither of them is implied by the other two. Second, for each task, we establish the necessity of the axioms. Specifically, we show that if at least one of the axioms required by the algorithm is omitted, then there exists a network that satisfies the other axioms, but for which the running time (of *any* distributed algorithm) is larger by at least a factor of $\log n$ and at most a factor $n$, see Table 1.

To support and justify our selection of axioms, we examine their usefulness for effective distributed computations on core-periphery networks. We consider a collection of different types of tasks, and show that they can be efficiently solved on core-periphery networks, by providing a distributed algorithm for each task and bounding its running time.

Table 1 provides a summary of the main tasks we studied, along with the upper and lower bounds on the running time when the network satisfies our axioms, and a worst case lower bound on the time required when at least one of the axioms is not satisfied. For each task we provide an algorithm, and prove formally its running time and the necessity of the axioms. As it turns out, some of the necessity proofs make use of an interesting connection to known communication complexity results.

The most *technically* challenging part of the paper is the distributed construction of a *minimum-weight spanning tree* (MST), a significant task in both the distributed systems world, cf. [28, 34], and the social networks world [1, 6, 9]. Thus, the main *algorithmic* result of the current paper is proving that MST can be computed efficiently (in $O(\log^2 n)$ rounds) on core-periphery networks, namely, networks that comply with our axioms (interestingly, our algorithm is randomized, which is a rarity in the distributed MST literature). To position this result in context, let us briefly review the state of the art on the problem of distributed MST construction. The problem was first studied in [16, 3], where the main focus was on low communication costs, and the run-time was at least linear in $n$. The study of sublinear-time distributed MST construction was initiated in [17]. Currently, the best upper bound for general graphs is $O\left(\sqrt{n}\log^* n + D\right)$ [22], where $D$ denotes the network diameter. Conversely, it was shown in [35] that there is a graph for which any deterministic algorithm for MST requires $\Omega(\sqrt{n}/\log n)$ time. More efficient algorithms exist for specific families of graphs. For the complete graph $G = K_n$, an MST, can be constructed in a distributed manner in $O(\log \log n)$ time [26]. For the wider class of graphs, of diameter at most 2, this task can still be performed in time $O(\log n)$. In contrast, taking the next step, and considering graphs of diameter 3, drastically changes the picture, as there are examples of such graphs for which any distributed MST construction requires $\Omega\left(\sqrt[4]{n}\right)$ time [27].

Let us now briefly review related work. Core periphery structures in social networks provided the inspiration to our approach. A number of excellent books provide a general review of social networks. A brief description of the

core-periphery structure, and the core's possible utilization as an interconnection mechanism among the network's users (albeit with no formal models or algorithmic procedures), can be found in Easley and Kleinberg's book [12]. The first explicit treatment of the core-periphery structure in social networks is given by Borgatti and Everett in [8], which provides a descriptive model for the core-periphery structure, and reviews some of its occurrences in a variety of social settings, lending support to our intuition regarding the centrality of the core-periphery structure in social networks. It does not, however, provide any systematic mathematical model for this structure. Similar phenomena of core-periphery structure were shown to exist in economics, e.g., the core–periphery model of Krugman [21] and other network formation models [18]. Several books present a formal rigorous analysis for the Preferential Attachment model. The interested reader is invited to look at, cf., Chapter 3 of [10], Chapter 4 of [7], or Chapter 14 of [32].

Turning to the distributed realm, in recent years there have been many studies focusing on the analysis of distributed algorithms on the complete graph $K_n$, cf. [33, 5, 11, 20, 13, 27, 35, 26, 24]. We believe that some of those algorithms can be extended to, and applied in, the context of core-periphery networks, which provide a much larger family of graphs.

The rest of the paper is organized as follows. Section 2 formally describes core-periphery networks, the axioms and their basic structural implications. Section 3 provides a description of the MST algorithm, and Section 4 presents the rest of the tasks we study.

## 2. Axiomatic design for core-periphery networks

### 2.1. Preliminaries

Let $G(V, E)$ denote our (simple undirected) network, where $V$ is the set of nodes, $|V| = n$, and $E$ is the set of edges, $|E| = m$. The network can be thought of as representing a distributed system. We assume the synchronous CONGEST model (cf. [34]), where communication proceeds in *rounds*, and in each round each node can send a message of at most $O(\log n)$ bits to each of its neighbors. Initially, each node has a unique ID of $O(\log n)$ bits.

For a node $v$, let $N(v)$ denote its set of neighbors and $d(v) = |N(v)|$ its degree. For a set $S \subset V$ and a node $v \in S$, let $N_{in}(v, S) = N(v) \cap S$ denote its set of neighbors within $S$, and denote the number of neighbors of $v$ in the set $S$ by $d_{in}(v, S) = |N_{in}(v, S)|$. Analogously, let $N_{out}(v, S) = N(v) \cap V \setminus S$ denote $v$'s set of neighbors outside the set $S$, and let $d_{out}(v) = |N_{out}(v, S)|$. For two subsets $S, T \subseteq V$, let $\partial(S, T)$ be the *edge boundary* (or cut) of $S$ and $T$, namely, the set of edges with exactly one endpoint in $S$, one in $T$ and $|\partial(S, T)| = \sum_{v \in S} |N_{out}(v, S) \cap T|$. Let $\partial(S)$ denote the boundary in the special case where $T = V \setminus S$.

### 2.2. Core-periphery networks

Given a network $G(V, E)$, a $\langle \mathcal{C}, \mathcal{P} \rangle$-partition is a partition of the nodes of $V$ into two sets, the *core* $\mathcal{C}$ and the *periphery* $\mathcal{P}$. Denote the sizes of the core and

the periphery by $n_c$ and $n_{\mathcal{P}}$, respectively. To represent the partition along with the network itself, we denote the *partitioned network* by $G(V, E, \mathcal{C}, \mathcal{P})$.

Intuitively, the core $\mathcal{C}$ consists of a relatively small group of strong and highly connected machines, designed to act as central servers, whereas the periphery $\mathcal{P}$ consists of the remaining nodes, typically acting as clients. The periphery machines are expected to be weaker and less well connected than the core machines, and they perform much of their communication via the dense interconnection network of the core. In particular, a central component in many of our algorithms, for various coordination and computational tasks, is based on assigning each node $v$ a *representative* core node $r(v)$, which is essentially a neighbor acting as a "channel" between $v$ and the core. The representative chosen for each periphery node is fixed.

For a partitioned network to be effective, the $\langle \mathcal{C}, \mathcal{P} \rangle$-partition must possess certain desirable properties. In particular, a partitioned network $G(V, E, \mathcal{C}, \mathcal{P})$ is called a *core-periphery network*, or $\mathcal{CP}$-*network* for short, if the $\langle \mathcal{C}, \mathcal{P} \rangle$-partition satisfies three properties, defined formally later on, in the form of three axioms.

### 2.3. Core-periphery properties and axioms

We define certain key parametrized properties of node groups, in networks that are of particular relevance to the relationships between core and periphery in our partitioned network architectures. We then state our axioms, which capture the expected behavior of those properties in core-periphery networks, and demonstrate their independence and necessity. Our three basic properties are the following.

**(i) $\alpha$-Balanced Boundary.** *A subset of nodes $S$ is said to have an $\alpha$-balanced boundary iff $\frac{d_{\mathrm{out}}(v, S)}{d_{\mathrm{in}}(v, S) + 1} = O(\alpha)$ for every node $v \in S$.*

**(ii) $\beta$-Clique Emulation.** *The task of* clique emulation *on an $n$-node graph $G$ involves delivering a distinct message $M_{v,w}$, from $v$ to $w$, for every pair of nodes $v, w$ in $V(G)$. An $n$-node graph $G$ is a $\beta$-clique-emulator, if it is possible to perform clique emulation on $G$ within $\beta$ rounds (in the CONGEST model).*

**(iii) $\gamma$-convergecast.** *For $S, T \subseteq V$, the task of $\langle S, T \rangle$-convergecast, on a graph $G$, involves delivering $|S|$ distinct messages $M_v$, originating at the nodes $v \in S$, to some nodes in $T$ (i.e., each message must reach at least one node in $T$). The sets $S, T \subset V$ form a $\gamma$-convergecaster if it is possible to perform $\langle S, T \rangle$-convergecast on $G$ in $\gamma$ rounds (in the CONGEST model).*

Consider a partitioned network $G(V, E, \mathcal{C}, \mathcal{P})$. We propose the following set of *axioms*, concerning the core $\mathcal{C}$ and periphery $\mathcal{P}$.

$\mathcal{A}_B$. **Core Boundary.** The core $\mathcal{C}$ has a $\Theta(1)$-balanced boundary.

$\mathcal{A}_E$. **Clique Emulation.** The core $\mathcal{C}$ is a $\Theta(1)$-clique emulator.

$\mathcal{A}_C$. **Periphery-Core Convergecast**. The periphery $\mathcal{P}$ and the core $\mathcal{C}$ form a $\Theta(1)$-convergecaster.
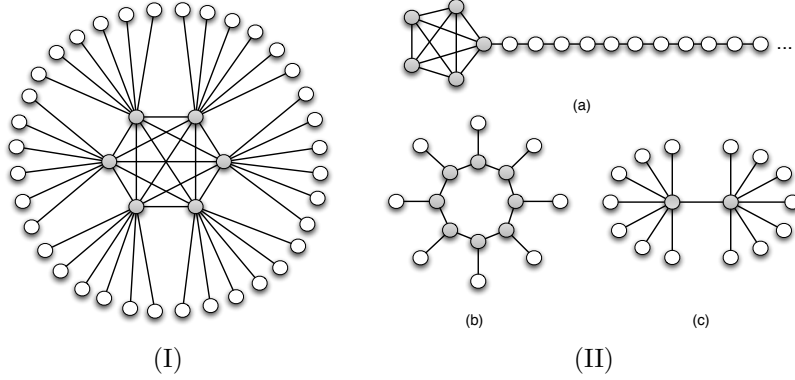
7

Figure 1: (I) An example for a 36-node $\mathcal{CP}$-network that satisfies all three axioms. The 6 core nodes (in gray) are connected in clique. In this example every core node is also an ambassador with equal number of edges to the core and outside the core. The core and periphery form a convergecaster since the periphery can send all its information to the core in one round. (II) Networks used in proofs: (a) The "lollipop partitioned" network $L_{25}$. (b) The "sun partitioned" network $S_{16}$. (c) The "dumbbell partitioned" network $D_{16}$ .

Let us briefly explain the axioms. Axiom $\mathcal{A}_B$ talks about the boundary between the core and periphery. Core nodes with a high out-degree (i.e., with many links to the periphery) are thought of as ambassadors of the core to the periphery. Axiom $\mathcal{A}_B$ states that while not all nodes in the core must serve as ambassadors, if a node is indeed an ambassador, then it must also have many links *within* the core. Axiom $\mathcal{A}_E$ talks about the flow of information within the core, and postulates that the core must be dense, and in a sense, behave almost like a complete graph: "everyone must know almost everyone else". The clique-emulation requirement is actually stronger than just being a dense subgraph, since the latter permits the existence of "bottlenecks" nodes, which a clique-emulator must avoid. Axiom $\mathcal{A}_C$ also concerns the boundary between the core and periphery, but in addition it refers also to the structure of the periphery. It postulates that information can flow efficiently from the periphery to the core. For example, it forbids the presence of periphery nodes that are far away from the core, or bottleneck edges that bridge between many periphery nodes and the core. Fig. 1(I) provides an example for a $\mathcal{CP}$-network satisfying the three axioms.

We next show that the axioms are independent. Later, we prove the necessity of the axioms for the efficient performance of a variety of computational tasks.

**Theorem 1.** *Axioms $\mathcal{A}_B$, $\mathcal{A}_E$, $\mathcal{A}_C$ are independent, namely, assuming any two of them does not imply the third.*

*Proof.* We prove the theorem by considering three examples of partitioned networks, described next. Each of these networks satisfies two of the axioms, but

violates the third (hence, they are not $\mathcal{CP}$-networks), implying independence.

*The lollipop partitioned network $L_n$ (Fig. 1(II)(a))*. The lollipop graph consists of a $\lfloor\sqrt{n}\rfloor$-node clique and a $n - \lfloor\sqrt{n}\rfloor$-node line, attached to some node of the clique. The corresponding partitioned network is obtained by setting its core $\mathcal{C}$ to be the clique, and its periphery $\mathcal{P}$ to be the line. Observe that $L_n$ is not a $\mathcal{CP}$-network. Indeed, Axiom $\mathcal{A}_E$ holds on $L_n$, and Axiom $\mathcal{A}_B$ also holds since the outgoing degree of each node in the core is 0 or 1. However, $\mathcal{A}_C$ is not satisfied on $L_n$ since the periphery consists of a line of length $n - \lfloor\sqrt{n}\rfloor$, so it will take linear time for the periphery $\mathcal{P}$ to convergecast to the core $\mathcal{C}$.

*The sun partitioned network $S_n$ (Fig. 1(II)(b))*. The sun graph consists of an $\lceil n/2 \rceil$-node cycle, with an additional leaf node attached to each cycle node. The corresponding partitioned network is obtained by setting its core $\mathcal{C}$ to be the cycle, and its periphery $\mathcal{P}$ to contain all other $\lfloor n/2 \rfloor$ nodes. Clearly, $\mathcal{A}_C$ holds on $S_n$, since each node in $\mathcal{P}$ is only one hop away from some node in $\mathcal{C}$. Axiom $\mathcal{A}_B$ also holds, since the outgoing degree of each node in $\mathcal{C}$ is 1. Axiom $\mathcal{A}_E$, however, does not hold, since the distance between two diametrically opposing nodes in the cycle is at least $n/4$, so it is not possible to perform $\Theta(1)$ clique emulation.

*The dumbbell partitioned network $D_n$ (Fig. 1(II)(c))*. The dumbbell graph is composed of two stars, each consisting of a center node connected to $\lceil n/2 \rceil - 1$ leaves, whose centers are connected by an edge. The corresponding partitioned network is obtained by setting its core $\mathcal{C}$ to be the two centers, and the periphery $\mathcal{P}$ to consist of the $n - 2$ leaves of the two stars (each of degree 1). It is easy to see that Axioms $\mathcal{A}_E$ and $\mathcal{A}_C$ hold on $D_n$, while Axiom $\mathcal{A}_B$ does not. $\square$

*2.4. Structural implications of the axioms*

The axioms imply a number of simple properties of the network structure.

**Theorem 2.** *If the partitioned network $G(V, E, \mathcal{C}, \mathcal{P})$ is a core-periphery network (i.e., it satisfies Axioms $\mathcal{A}_B$, $\mathcal{A}_E$ and $\mathcal{A}_C$), then the following properties hold:*

1. *The core size satisfies $\Omega(\sqrt{n}) \leq n_c \leq O(\sqrt{m})$.*

2. *Every $v \in \mathcal{C}$ satisfies $d_{\mathrm{out}}(v, \mathcal{C}) = O(n_c)$ and $d_{\mathrm{in}}(v, \mathcal{C}) = \Omega(n_c)$.*

3. *The number of outgoing edges of the core is $|\partial(\mathcal{C})| = \Theta(n_c^2)$.*

4. *The core is dense, i.e., the number of edges in it is $\sum_{v \in \mathcal{C}} d_{\mathrm{in}}(v, \mathcal{C}) = \Theta(n_c^2)$.*

*Proof.* Axiom $\mathcal{A}_E$ necessitates that the inner degree of each node $v$ is $d_{\mathrm{in}}(v, \mathcal{C}) = \Omega(n_c)$ (or else it would not be possible to complete clique emulation in constant time), implying the second part of claim 2. It follows that the number of edges in the core is $\sum_{v \in \mathcal{C}} d_{\mathrm{in}}(v, \mathcal{C}) = \Theta(n_c^2)$, hence it is dense; claim 4 follows. Since

also $\sum_{v \in \mathcal{C}} d_{\text{in}}(v, \mathcal{C}) \leq 2m$, we must have the upper bound of claim 1, that is, $n_c = O(\sqrt{m})$. Axiom $\mathcal{A}_B$ yields that for every $v$, $d_{\text{out}}(v, \mathcal{C}) = O(n_c)$, so the first part of claim 2 follows. Note that $|\partial(\mathcal{C})| = \sum_{v \in \mathcal{C}} d_{\text{out}}(v, \mathcal{C}) = O(n_c^2)$, so the upper bound of claim 3 follows. To give a lower bound on $n_c$, note that by Axiom $\mathcal{A}_C$ we have $|\partial(\mathcal{C})| = \Omega(n - n_c)$ (otherwise the information from the $n - n_c$ nodes of $\mathcal{P}$ could not flow in $O(1)$ time to $\mathcal{C}$), so $n_c = \Omega(\sqrt{n})$ and the lower bounds of claims 1 and 3 follow. $\qquad \square$

An interesting case for efficient networks is where the number of edges is linear in the number of nodes. In this case, Theorem 2 implies the following.

**Corollary 1.** *In a core-periphery network $G(V, E, \mathcal{C}, \mathcal{P})$ where $m = O(n)$, the following properties hold:*

1. *The core size satisfies $n_c = \Theta(\sqrt{n})$.*

2. *The number of outgoing edges from the core is $|\partial(\mathcal{C})| = \Theta(n)$.*

3. *The number of edges in the core is $\sum_{v \in \mathcal{C}} d_{\text{in}}(v, \mathcal{C}) = \Theta(n)$.*

Now we show a key property relating our axioms to the network diameter.

**Claim 1.** *If the partitioned network $G(V, E, \mathcal{C}, \mathcal{P})$ satisfies Axioms $\mathcal{A}_E$ and $\mathcal{A}_C$, then its diameter is $\Theta(1)$.*

*Proof.* Suppose the partitioned network $G(V, E, \mathcal{C}, \mathcal{P})$ satisfies Axioms $\mathcal{A}_E$ and $\mathcal{A}_C$. Let $u, v$ be two nodes in $V$. There are three cases to consider: (1) Both $u, v \in \mathcal{C}$: then Axiom $\mathcal{A}_E$ ensures $O(1)$ time message delivery, and thus $O(1)$ distance. (2) $u \in \mathcal{P}$ and $v \in \mathcal{C}$: Axiom $\mathcal{A}_C$ implies that there must be a node in $w \in \mathcal{C}$ such that $dist(u, w) = O(1)$. By Axiom $\mathcal{A}_E$, and $dist(w, v) = O(1)$, thus $dist(u, v) \leq dist(u, w) + dist(w, v) = O(1)$. (3) Both $u, v \in P$: then there must be $w, x \in \mathcal{C}$ such that $dist(u, w) = O(1)$ and $dist(x, v) = O(1)$. Since $dist(w, x) = O(1)$ by Axiom $\mathcal{A}_E$, it follows that $dist(u, v) \leq dist(u, w) + dist(w, x) + dist(x, v) = O(1)$. Hence $dist(u, v) = O(1)$ for any $u, v \in V$, so the diameter of $G(V, E)$ is constant. $\qquad \square$

The following claim shows that the above conditions are necessary.

**Claim 2.** *For $X \in \{E, C\}$, there exists a family of $n$-node partitioned networks $G_X(V, E, \mathcal{C}, \mathcal{P})$, of diameter $\Omega(n)$, which satisfy all axioms except $\mathcal{A}_X$.*

*Proof.* For $X = C$, let $G_C(V, E, \mathcal{C}, \mathcal{P})$ be the lollipop partitioned network $L_n$. As mentioned before, for this network Axiom $\mathcal{A}_C$ is violated, while the others are not. Also note that the diameter of $G_C$ is $\Omega(n)$.

For $X = E$, let $G_E(V, E, \mathcal{C}, \mathcal{P})$ be the sun partitioned network $S_n$. As mentioned before, for this network Axiom $\mathcal{A}_E$ is violated, while the others are not. Also note that the diameter of $G_E$ is $\Omega(n)$. $\qquad \square$

## 3. MST on a Core-Periphery Network

In this section we present a time-efficient randomized distributed algorithm $\mathcal{CP}$-MST for computing a *minimum-weight spanning tree* (MST) on a core-periphery network. In particular, we consider a $n$-node core periphery network $G(V, E, \mathcal{C}, \mathcal{P})$, namely, a partitioned network satisfying all three axioms, and show that a MST can be computed in a distributed manner on such a network in $O(\log^2 n)$ rounds with high probability. Upon termination, each node knows which of its edges belong to the MST. We also show that Axioms $\mathcal{A}_B$, $\mathcal{A}_E$, and $\mathcal{A}_C$ are indeed necessary, for our distributed MST algorithm to be efficient.

### 3.1. Axiom necessity

**Theorem 3.** *For each $X \in \{B, E, C\}$ there exists a family of $n$-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$, that do not satisfy Axiom $\mathcal{A}_X$, but satisfy the other two axioms; and the time complexity of* any *distributed MST algorithm on $\mathcal{F}_X$ is $\Omega(n^{\alpha_X})$, for some constant $\alpha_X > 0$.*

*Proof.* For $X = B$, consider the graph $G_B$ on Figure 2(a), in which Core is a clique of size $k$, and each node in the Core is connected to $k^3$ Periphery nodes (one node in Core is also connected to $s$, so it has $k^3 + 1$ Periphery neighbors). The number of nodes in $G_B$ is thus $n = k + k \cdot k^3 + 1 = \Theta(k^4)$. In [27], it was shown that any distributed algorithm will take at least $\Omega\left(\sqrt[4]{n}/\sqrt{\log n}\right)$ time on $G_B$. Since Core is a clique, $G_B$ satisfies Axiom $\mathcal{A}_E$. Since every node in Periphery has a direct edge to the Core, $G_B$ satisfies Axiom $\mathcal{A}_C$, i.e., it is possible to perform a convergecast in $O(1)$ time. But notice that $d_{\text{in}} = \sqrt[4]{n}$ while $d_{\text{out}} = \sqrt[4]{n^3}$ and thus $G_B$ does not satisfy Axiom $\mathcal{A}_B$.

For $X = E$, consider the graph $G_E$ on Figure 3(a), in which Core is a collection of $k$ cliques, each of size $k$, where a single node in each clique is connected to a special Core node $u$, and there are no edges between cliques. The $k^3$ Periphery nodes are arranged in $k$ columns of $k^2$ nodes each. Each node in the Core (except $u$) is connected to $k$ Periphery nodes such that the nodes in a specific clique $i$ are connected to all the Periphery nodes that reside in a specific column $i$. One Core node (from the leftmost clique) is additionally connected to $s$, and another Core node (from the rightmost clique) is connected to $r$. The number of nodes in $G_E$ is thus $n = k \cdot k + k \cdot k^2 + 2 = \Theta(k^3)$.

Assume the following weight assignment. All the edges between Core and Periphery have weight 10, except for the two edges that come from $s$ and $r$. The weights of all the edges incident to $s$ are 2, and the weights of all the edges incident to $r$ are 3. Assume also that the weights of all the rest of the edges in Periphery are 1. It's easy to see that such a weight assignment will yield an MST as illustrated in Figure 3(b). Notice that increasing the weight of some edge incident to $s$ (say, to 5), will cause this edge to be removed from the MST, and the corresponding edge incident to $r$ to be included. Thus, in order for $r$ to know which of its edges belong to the MST, it needs to receive information regarding the weights of all the edges incident to $s$, i.e., at least $k^2$ edge weights should be delivered from $s$ to $r$. In the CONGEST model, at most $O(\log n)$ edge

11

weights can be sent in a single message, hence $\Omega(k^2/\log n)$ messages must be delivered from $s$ to $r$. Next, we show that delivering $k^2/\log n$ messages from $s$ to $r$ will require at least $\Omega(k/(\log n))$ time. First, note that any path $s \to r$ that is not passing via the node $u$ has length at least $k$, thus if any of the messages avoids $u$, we are done. So now assume that all the messages take paths via $u$. Observe that the edge cut of the node $u$ (i.e., its degree) is $k$, and thus in $k/(2\log n)$ time units, it can forward at most $k^2/(2\log n)$ messages, which is not sufficient for completing the MST task. Thus, any MST algorithm on the graph $G_E$ will take at least $\Omega(k/\log n) = \Omega(\sqrt[3]{n}/\log n)$ time.

It is left to show that $G_E$ satisfies Axioms $\mathcal{A}_B$ and $\mathcal{A}_C$, but not $\mathcal{A}_E$. For every node in the Core, $d_{\text{in}} = k$ and $d_{\text{out}} = k$, except the node $u$, for which $d_{\text{out}} = 0$. So, for each node in the core $d_{\text{out}}/(d_{\text{in}}+1) = O(1)$, which means that $\mathcal{A}_B$ is satisfied. Since every node in Periphery has a direct edge to the Core, $G_E$ satisfies Axiom $\mathcal{A}_C$, i.e., it is possible to perform a convergecast in $O(1)$. It is also easy to see that the Core does not support $O(1)$-clique emulation ($\mathcal{A}_E$), since sending $k$ messages out of any clique in Core to any other clique in Core requires $k$ time, as there is only one edge connecting any clique to the node $u$.

Finally, for $X = C$, consider a graph $G_C$ on Figure 2(b), in which Core is a clique of size $k$, and each node in the Core is connected to $k/2$ Periphery nodes. One Core node is additionally connected to a cycle of size $k^2/2$ that resides in Periphery. The number of nodes in $G_C$ is thus $n = k + k \cdot k/2 + k^2/2 = \Theta(k^2)$. It is easy to see that Axioms $\mathcal{A}_B$ and $\mathcal{A}_E$ are satisfied, but Axiom $\mathcal{A}_C$ is not. For a suitable weight assignment, the decision regarding which edge of $r$ to include in the MST depends on the weights of the edges incident to $s$. The last observation implies that at least one message has to be delivered from $s$ to $r$ which will take $\Omega(k^2) = \Omega(n)$ time. $\qquad\square$

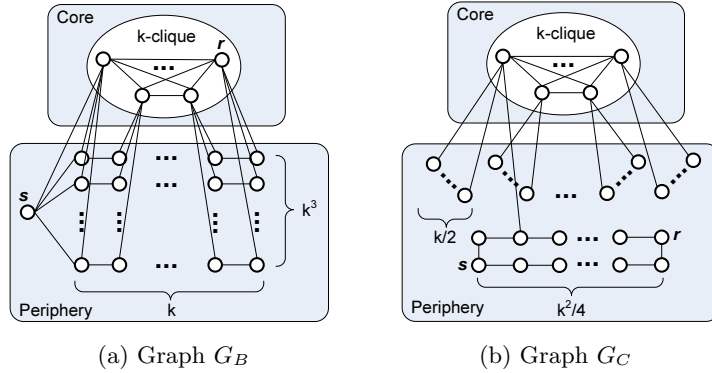

(a) Graph $G_B$        (b) Graph $G_C$

Figure 2: (a) Graph $G_B$: each node in the Core is connected to $k^3$ Periphery nodes. (b) Graph $G_C$: each node in the Core is connected to $k/2$ Periphery nodes, and one Core node is connected to cycle of length $k^2/2$ in Periphery.
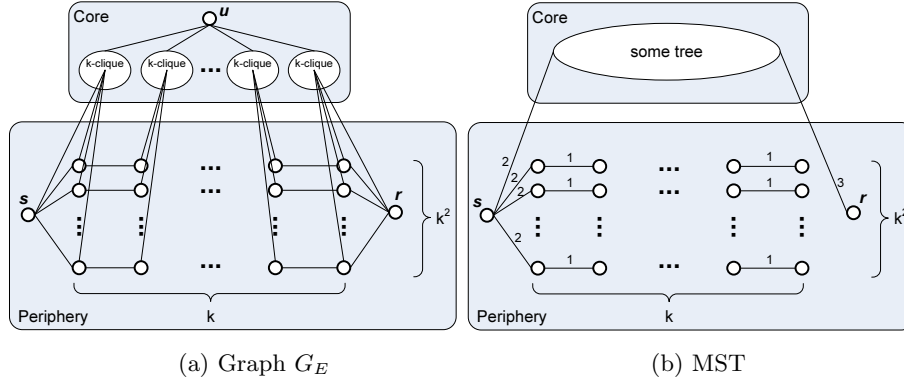
(a) Graph $G_E$        (b) MST

Figure 3: (a) Graph $G_E$: Core consists of $k$ cliques each of size $k$. Each node in the Core (except $u$) is connected to $k$ nodes in Periphery. (b) Possible MST of $G_E$.

### 3.2. Description of the $\mathcal{CP}$-MST algorithm

Let us now give a high level description of our $\mathcal{CP}$-MST algorithm. The algorithm is based on Boruvka's MST algorithm [31], and runs in $O(\log n)$ phases, each consisting of several steps. The algorithm proceeds by maintaining a forest of *tree fragments* (initially singletons) and gradually merging fragments, until the forest converges to a single tree. Throughout the execution, each node has two *officials*, namely, core nodes that represent it. In particular, recall that each node $v$ is assigned a *representative* core neighbor $r(v)$, passing information between $v$ and the core. In addition, $v$ is also managed by the *leader* $l(i)$ of its current fragment $i$. An important distinction between these two roles is that the representative of each node is fixed, while its fragment leader may change in each phase (as its fragment grows). At the beginning of each phase, every node knows the IDs of its fragment and its leader. Then, every node considers all its outgoing edges (i.e., edges with the second endpoint belonging to another fragment), and finds its minimum weight outgoing edge. This information is delivered to the core by the means of the representative nodes, which receive the information, aggregate it (as much as possible) and forward it to the leaders of the appropriate fragments. The leaders decide on the fragment merging, and inform all the nodes about new fragments IDs.

The main challenges in obtaining the proof were in bounding the running time, which required careful analysis. There are two major sources of potential delays in the algorithm. The first involves sending information between officials (representatives to leaders and vice versa). Note that there are only $O(\sqrt{m})$ officials, but they may need to send information about $m$ edges, which can lead to congestion. For example, if more than $\alpha \cdot \sqrt{m}$ messages need to be sent to an official of degree $\sqrt{m}$, then this will take at least $\alpha$ rounds. We use randomization of leaders and the property of clique emulation to avoid this situation, and make sure that officials do not have to send, or receive, more than

$O(\sqrt{m} \log m)$ messages in a phase.

The second source for delays is the fragment merging procedure. This further splits into two types of problems. The first is that a chain of fragments that need to be merged could be long, and in the basic distributed Boruvka's algorithm will take long time (up to $n$) to resolve. This problem is overcome by using a modified pointer jumping technique, similar to [27]. The second problem is that the number of fragments that need to be merged could be large, resulting in a large number of *merging* messages that contain, for example, the new fragment ID. This problem is overcome by using randomization, and by reducing the number of messages needed for resolving a merge.

Before describing the subsequent phases of the algorithm, a few definitions are in place. At any point throughout the execution, let $f(u)$ denote the fragment that $u$ belongs to. Dually, let $V^i$ denote the set of nodes in fragment $i$, and let $V^i(w)$ denote the subset of $V^i$ consisting of the nodes that are represented by $w$. For a representative $w \in \mathcal{C}$, let $F_{rep}(w)$ be the set of fragments that $w$ represents, namely, $F_{rep}(w) = \{i \mid V^i(w) \neq \emptyset\}$, and let $F_{lead}(w)$ be the set of fragments that $w$ leads, namely, $F_{lead}(w) = \{i \mid l(i) = w\}$. For a set of nodes $S^i$ belonging to the same fragment $i$, an *outgoing edge* is one whose second endpoint belongs to a different fragment. Let $\text{mwoe}(S^i)$ be the *minimum weight outgoing edge* of $S^i$. For a node $u$, a fragment $i$ and a representative $w$, we may occasionally refer to the fragment's $\text{mwoe}$ as either $\text{mwoe}(u)$, $\text{mwoe}(V^i)$ or $\text{mwoe}(V^i(w))$. The *merge-partner* of fragment $i$, denoted $\text{mp}(i)$, is the fragment of the second endpoint of the edge $\text{mwoe}(V^i)$. Define $F_{lead}^j(w) \subseteq F_{lead}(w)$ to be the set of fragments led by $w$ that attempt to merge with the fragment $j$, i.e., $F_{lead}^j(w) = \{i \mid i \in F_{lead}(w) \text{ and } \text{mp}(i) = j\}$. Define a *speaker* fragment $\text{spk}^j(w) = \min F_{lead}^j(w)$, that is responsible for sending merge-requests on behalf of all the fragments in $F_{lead}^j(w)$, and updating them upon the reception of merge-replies.

We now proceed with the description of the algorithm.

*Phase 0 – Initialization.*

1. **Obtaining a Representative.** Each node $u \in V$ obtains a representative $r(u) \in \mathcal{C}$ in the core. In particular, if $u \in \mathcal{C}$, it represents itself, i.e., $r(u) = u$. Each periphery node $u \in \mathcal{P}$ sends a "representative-request" message towards the core $\mathcal{C}$ with its ID. This is done in parallel, using a $\gamma$-convergecast protocol on $\mathcal{P}$ and $\mathcal{C}$, which ensures that each such message is received by some node in $\mathcal{C}$. Once a node $w \in \mathcal{C}$ receives such a message, it replies to $u$ on the same route, and $u$ sets $r(u) = w$.

2. **Renaming.** Each node $u \in V$ receives a unique ID, $\text{id}(u) \in [1, \ldots, n]$. This step can be performed in the following simple way: each node sends to its representative its ID, and each representative sends its own ID and the number of nodes it represents, to all core members. Now, every core member can sort the core IDs and reserve a sufficiently large range of IDs for each representative. Each node in the core can now set its own new ID,

14

and send unique new IDs in the range $[1 \ldots n]$ to the nodes it represents. We assume nodes in the core $\mathcal{C}$ take IDs $[1 \ldots n_c]$.

3. **Fragment ID Initialization.** Each node $u \in V$ forms a singleton fragment with its unique $\text{id}(u)$.

4. **Obtaining a Leader.** Each initial fragment $i = f(u)$ (which is a singleton at this phase) obtains a leader by asking the representative $r(u)$ of node $u$, to select a random Core member $w$ uniformly at random, and declare it as a leader of $i$, $l(i) = w$. This is done, in a balanced way, by picking a random permutation and assigning leaders according to it, hence every node in $\mathcal{C}$ becomes the leader of $O(n_c)$ fragments.

5. **Fragment State Initialization.** Each leader keeps a state (*active / frozen / root / waiting*) for each of its fragments. The initial state of all fragments is *active*.

*Phase $b \in \{1 \ldots B\}$ (similar to Boruvka's phases).*

1. **Finding mwoe.** Each $u \in V$ finds an edge $(u, v) = \texttt{mwoe}(u)$, and obtains $f(v)$ and $l(f(v))$.

2. **Periphery to Representatives.** Each node $u \in V$ sends $(u, v) = \texttt{mwoe}(u)$, $f(u)$, $l(f(u))$, $f(v)$ and $l(f(v))$ to its representative $r(u) \in \mathcal{C}$.

3. **Representatives to Leaders.** Each representative $w \in \mathcal{C}$, for each fragment $i \in F_{rep}(w)$, sends $(u, v) = \texttt{mwoe}(V^i(w))$, $i$, $f(v)$, and $l(f(v))$ to the leader $l(i)$ of $i$.

4. **Leaders Merge Fragments.** Each leader $w \in \mathcal{C}$, for each fragment $i \in F_{lead}(w)$, finds $(u, v) = \texttt{mwoe}(V^i)$ and $\texttt{mp}(i) = f(v)$, and then executes MERGEFRAGS($i$).

5. **Leaders to Representatives.** Each leader $w \in \mathcal{C}$, for each *active* fragment $i \in F_{lead}(w)$, sends an update message with the new fragment name $newID(i)$, the new leader node $l(newID(i))$ and the edge to add to the MST, to all the representatives of the nodes in $V^i$. If $w \neq l(newID(i))$, then the fragment $i$ is removed from $F_{lead}(w)$.

6. **Representatives to Periphery.** Each representative $w \in \mathcal{C}$, for each $i \in F_{rep}(w)$, for which the update message with $newID(i)$ and $l(newID(i))$ was received, forwards it to all the nodes of $V^i(w)$.

### 3.3. MERGEFRAGS *procedure*

The MERGEFRAGS procedure is the essential part of our algorithm, executed at each phase $b$. The procedure is executed by each leader $w \in \mathcal{C}$, for each fragment $i \in F_{lead}(w)$. For a fragment $i$, its leader maintains a state parameter $state(i) \in \{active, frozen, root, waiting\}$. Each fragment $i$ attempts to merge

with some other fragment $\mathtt{mp}(i)$. Towards that, the leader of $i$ initiates a merge-request to a leader of the fragment $\mathtt{mp}(i)$ (the fragment at the other end of $\mathtt{mwoe}(i)$). Since these requests do not have to be reciprocal, merge requests usually form a *merge-tree*, whose nodes are fragments and whose directed edges represent merge-requests (see Figure 4 for illustration). In order to minimize the number of merge-request messages sent by fragment leaders, we propose to designate, for each set of fragments sharing the same leader that attempt to merge with the same target fragment, a *speaker* fragment, that will act on behalf of all the fragments in the set, and update all of them upon reception of merge-replies.
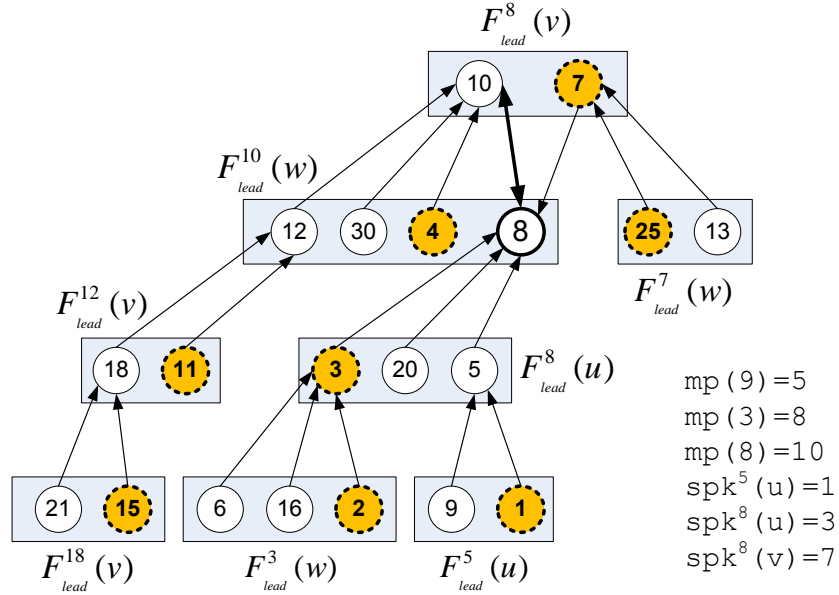


$$\mathtt{mp}(9)=5$$
$$\mathtt{mp}(3)=8$$
$$\mathtt{mp}(8)=10$$
$$\mathtt{spk}^5(u)=1$$
$$\mathtt{spk}^8(u)=3$$
$$\mathtt{spk}^8(v)=7$$

Figure 4: Illustration of a fragments merge-tree. An arrow $i \to j$ means that fragment $i$ attempts to merge with fragment $j$, i.e., $j = \mathtt{mp}(i)$. The root of the merge-tree is fragment 8, since it has a reciprocal arrow with fragment 10 and $8 < 10$.

The root of that tree is a fragment that received a reciprocal merge-request (actually, there are two such fragments, so the one with the smaller ID is selected as a root). However, since merge-requests are not sent by every fragment, only by *speakers*, the root node is detected by the *speaker*, and not the root fragment itself (except for the case when the root is the *speaker*). For example, in Figure 4, fragment 4 sends a merge request to fragment 10, and gets a merge-reply with the *next* pointer of 10, which is 8 (the *next* pointer of fragment $i$ is always set initially to $\mathtt{mp}(i)$). Fragment 4 then realizes that 8 belongs to $F_{lead}^{10}(w)$, and thus identifies the reciprocity between 8 and 10. Fragment 4 (the *speaker*) then notifies 8 that it should be the root (7 does not notify 10 since $8 < 10$). For a

detailed description of the root finding procedure see Algorithm FINDROOT($i$) in Appendix A.

When a fragment $i$ that is led by $w$ is in the *active* state and attempts to merge with another fragment ($\mathtt{mp}(i)$), it first tries to find the root using the procedure FINDROOT (see Appendix A for the pseudocode). By the end of the FINDROOT procedure, $i$ may not find a root, in which case its state will become $frozen$; $i$ may find that the root is another fragment $k$ in $F_{lead}^{\mathtt{mp}(i)}(w)$, and then $i$ will notify $k$, but $i$'s state will become $frozen$; $i$ may find that it is a root by itself, in which case its state will become $root$; and finally, $i$ may be notified by a *speaker* of $F_{lead}^{\mathtt{mp}(i)}(w)$ and $i$'s state will become $root$.

Once a fragment enters the *root* state, it starts waiting for all the tree fragments to send it merge-requests. These merge-requests are sent by each fragment, using the pointer-jumping procedure PJ (see Appendix A for the pseudocode), while it is in the $frozen$ state. Once the requests of all the tree fragments reach the *root* (using pointer-jumping), it chooses a new random ID ($newID$) for the fragment, among all the fragments in the tree, and a random Core node to be the new leader ($newLead$) for this fragment, and sends this information back to all of them. At this point, the merge-tree is considered to be resolved, and all its fragments (including the *root*) change their state to *active*. The simple state diagram of Algorithm MERGEFRAGS can be found in Figure 5, and a detailed pseudocode in Appendix A.
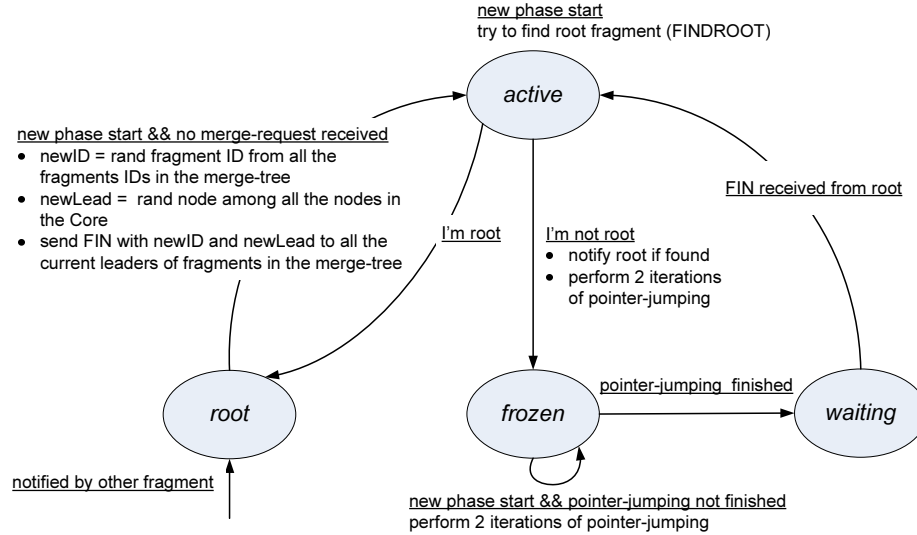


Figure 5: State diagram of Algorithm MERGEFRAGS. The algorithm is executed by every leader for every fragment it leads. Each title indicates an event, and the text below it is the action performed upon that event.

Now we briefly describe the pointer-jumping approach used to resolve fragment *merge-trees*. Pointer-jumping is a technique (developed in [37] and often

17

used in parallel algorithms) for contracting a given linked list of length $k$, causing all its elements to point to the last element, in $O(\log k)$ steps. We use the pointer-jumping approach for resolving merge-trees, viewing the fragments as the nodes in a linked list, with each fragment $i$ initially pointing at $\mathtt{mp}(i)$. Each fragment chain can be of length at most $O(n)$, and thus can be contracted in $\log n$ rounds, resulting in a $\log n$ time overhead per phase. In order to overcome this, we use a technique first introduced in [27], called "amortized pointer-jumping", in which the contraction of long chains is deferred to later phases, while in each phase only a small constant number of pointer-jumps is performed. The argument for the correctness of this approach is that if the chain (or tree) is large, then the resulting fragment, once resolved, is large enough to satisfy the fragment growth rate needed to complete the algorithm in $B = O(\log n)$ phases (see Claim 6).

*3.4. Correctness of the $\mathcal{CP}$-MST Algorithm*

We now show that our $\mathcal{CP}$-MST algorithm is correct, i.e., it results in an MST. The following claim shows that the MERGEFRAGS algorithm indeed resolves a merge-tree.

**Claim 3.** *Once a merge-tree becomes* active, *all the (old) fragments in the tree have the same (new) fragment ID.*

*Proof.* The claim follows directly from the description above, and the observation that in the pointer-jumping procedure, at every step, at least one more node points to the *root*. Thus, if at some phase the *root* of the merge-tree does not receive any merge-request, then every other fragment in the tree is in *waiting* state, i.e., points to the *root*. Consequently, the *root* knows all the fragments in the tree, and can inform their leaders about the new fragment ID, $newID$, and the new leader node $l(newID)$. $\qquad\square$

**Claim 4.** *The $\mathcal{CP}$-MST algorithm emulates Boruvka's MST algorithm and thus constructs an MST for the network.*

*Proof.* In Boruvka's algorithm, fragment merges can be performed in any order. What's important is that a merge between any two fragments will occur if, and only if, they share an edge that is an $\mathtt{mwoe}$, for at least one of the fragments. Since our algorithm satisfies this property, it results in an MST. $\qquad\square$

*3.5. Running time analysis of the $\mathcal{CP}$-MST algorithm*

We now analyze the running time of Algorithm $\mathcal{CP}$-MST in a Core-Periphery network. To do that, we analyze each part of the algorithm separately, and prove the following Claims 5, 6 and 7. We start with the initialization phase.

**Claim 5.** *The initialization phase (Phase 0) takes $O(1)$ rounds.*

*Proof.* Due to Axiom $\mathcal{A}_C$, the convergecast process employed in step 1 of Phase 0 requires $O(1)$ rounds. The renaming step (step 2 of Phase 0) can be done in $O(1)$ rounds, due to Axiom $\mathcal{A}_E$. The operation of obtaining a leader (step 3 of Phase 0) requires $O(1)$ rounds due to the Axiom $\mathcal{A}_C$. □

Now we show that the number of Boruvka phases needed in algorithm $\mathcal{CP}$-MST is $B = O(\log n)$.

**Claim 6.** *Algorithm $\mathcal{CP}$-MST takes $O(\log n)$ phases, i.e., $B = O(\log n)$.*

*Proof.* The proof is by induction. Assume that every active fragment $f$ at phase $x \le i$ has size (in nodes) $|f| > \min(2^x, n)$. We show that in the phase $j > i$ at which $f$ becomes active again, its size will be at least $\min(2^j, n)$. In phase $i$, $f$ joins a merge-tree that was created at some phase $k \le i$, and according to induction assumption, every fragment in this tree has size at least $\min(2^k, n)$. That tree will be resolved in phase $j$, i.e., after $j - k$ phases. Let $D$ be the diameter of the tree in phase $j$. Since the algorithm uses pointer jumping with two iterations at each phase, it follows that $j - k \le \lceil \log D \rceil / 2$. The size of the resolved tree is at least $\min(2^k, D)$, since it comprises of at least $D$ fragments, each of size at least $\min(2^k, n)$. Clearly,

$$2^k D \;=\; 2^{k + \log D} \;\ge\; 2^{k + \frac{\lceil \log D \rceil}{2}} \;\ge\; 2^j,$$

and thus $|f| \ge \min(2^j, n)$. So each active fragment at phase $j$ is of size at least $\min(2^j, n)$. If in phase $\lceil \log n \rceil$ there are no active fragments, then the algorithm waits for at most $\log n$ time, which is sufficient to resolve any fragments tree, and then, the size of the fragment is $\min(2^{2 \log n}, n) = n$, which means that the algorithm has terminated. □

Finally, we analyze the steps performed in phases $b \in [1, \ldots, B]$. First, we give the following auxiliary lemma. The result of this lemma is well known, and its proof is analogous to the proof of Lemma 5.1 in [30].

**Lemma 1.** *For every real $x > 0$, when up to $k$ balls are thrown independently and uniformly at random into at least $w$ bins, the maximum loaded bin has at most $O(k/w + \log x)$ with probability at least $1 - w/x^c$, where $c$ is an arbitrary constant.*

*Proof.* Let $X_i$ be the random variable representing the number of balls in bin $i$. For integer $l \ge 0$,

$$
\begin{aligned}
\Pr(X_i \ge l) \;&\le\; \binom{k}{l} \cdot \left(\frac{1}{w}\right)^l \le \frac{k^l}{l!} \cdot \frac{1}{w^l} = \left(\frac{k}{w}\right)^l \cdot \frac{1}{l!} \\
&\le\; \left(\frac{k}{w}\right)^l \cdot \left(\frac{e}{l}\right)^l = \left(\frac{ek}{wl}\right)^l
\end{aligned}
$$

For $l = c_1(k/w + \log x)$, we obtain

$$\Pr(X_i \geq c_1(k/w + \log x)) \ \leq \ \left(\frac{e}{c_1}\right)^{c_1(k/w+\log x)} \ \leq \ \frac{1}{x^c} \ ,$$

where $c = c(c_1)$ is an arbitrary constant.

By taking union bound over all the $w$ bins, we obtain that the probability that any bin has at most $O(k/w + \log x)$ balls with probability of at least $1 - w/x^c$. $\qquad\square$

**Claim 7.** *For every phase $b \in [1, \ldots, B]$, the running times of the main steps $1, 2$ and $6$ are bounded by $O(1)$, and of steps $3, 4$ and $5$ by $O(\log n)$. Thus, every phase $b$ takes $O(\log n)$ rounds.*

*Proof.* In step 1, every node sends a single message to all its neighbors, so the running time is $O(1)$. In step 2, each node $u \in V$ sends $\mathtt{mwoe}(u)$ to $r(u) \in \mathcal{C}$ using $\gamma$-convergecast. By Axiom $\mathcal{A}_C$, the running time is $O(\gamma) = O(1)$. Next, consider step 3. Since the network satisfies Axiom $\mathcal{A}_E$, one may assume that $\mathcal{C}$ is a clique. To derive the running time of this step we have to calculate how many messages are sent between a representative $u$ and a leader $v$ in $\mathcal{C}$. It suffices to look only at Core edges, since this step involves communication only between nodes in $\mathcal{C}$ (representatives and leaders). By Theorem 2(2), $d_{out}(u) = O(n_c)$, and since $\mathcal{P}$ and $\mathcal{C}$ form a $\Theta(1)$-convergecaster, it follows that on each edge towards $\mathcal{P}$, $u \in \mathcal{C}$ receives a constant number of "representative-requests" at the initialization phase. The last claim implies that $u$ represents $O(n_c)$ nodes, and thus at most $O(n_c)$ fragments. Hence, every representative node has to send $O(n_c)$ messages, each destined to a leader of a specific fragment. Since the "leadership" on a fragment is assigned independently at random to the nodes in $\mathcal{C}$, sending messages from representative to leaders is analogous to throwing $O(n_c)$ balls into $n_c$ bins. Hence by Lemma 1 with an appropriate constant $c$, the most loaded edge (bin) from a representative $u$ to some leader $v$ handles $O(\log n_c)$ messages (balls), with probability at least $1 - 1/n_c^8$. Applying the union bound over all $O(n_c)$ representative nodes and all $O(\log n) = O(\log n_c)$ phases of the algorithm, we get that the most loaded edge in step 3, is at most $O(\log n)$ with probability at least $1 - 1/n$. Thus, this step takes $O(\log n)$ time.

In step 4, every execution of Procedure MERGEFRAGS requires sending / receiving merge-request/reply messages from every leader $u \in \mathcal{C}$, for each fragment $i \in F_{lead}(u)$. For each merge-request there is exactly one merge-reply, so it suffices to count only merge-requests. Moreover, if there are multiple fragments that have the same leader node and need to send a merge-request to the same other fragment, only one message will actually be sent by the *speaker* fragment. The last observation implies that every request message sent by a leader is destined to a different fragment (i.e., to its leader). As in the analysis of the previous step, since "leadership" is assigned independently at random to the nodes in $\mathcal{C}$, sending messages from leaders to leaders is analogous to throwing $A$ balls into $n_c$ bins, where $A$ is the number of fragments that the node

$u$ leads. Using Lemma 1 with an appropriate constant $c$, $A$ can be bounded with high probability by $O(\sqrt{n})$, since up to $n$ fragments (balls) are assigned to $n_c = \Omega(\sqrt{n})$ Core nodes (bins).

We now apply Lemma 1 with $O(\sqrt{n})$ balls (fragments led by a node), and $n_c = \Omega(\sqrt{n})$ bins (edges towards other Core nodes), and conclude that the most loaded edge from a leader $u$ to some other leader $v$ carries $O(\log n)$ messages, with probability at least $1 - 1/(\sqrt{n})^8 = 1 - 1/n^4$. Applying the union bound over all the $O(\sqrt{m})$ leaders, and all $O(\log n)$ phases of the algorithm, we get that the most loaded edge in the process of sending merge-requests carries at most $O(\log n)$ messages, with probability at least $1 - 1/n^c$.

The last part of step 4 is when the root fragment sends the FIN ("finish") message to all the merge-tree members. at the beginning of phase $j$, the size of each active fragment is at least $2^j$ (see Claim 6) and at most $2^{j+1}$ (as the root does not release a tree at phase $j - 1$ if it is too large). Thus, the number of merge-trees resolved at phase $i$ is at most $n/2^{i+1}$ (every resolved tree becomes an active fragment at the next phase). In case $2^{i+1} \leq \sqrt{n}$, it follows from the Lemma 1 that a leader node $u \in \mathcal{C}$ has at most $O(\sqrt{n}/2^{i+1} + \log n)$ roots (at most $n/2^{i+1}$ balls into at least $\sqrt{n}$ bins). For each root, a leader has to send a message, for each fragment in its tree. The number of fragments in the tree is bounded by the number of nodes in the tree, which is $2^{i+2}$ (this is because the tree becomes an active fragment at the beginning of the next phase $j = i + 1$, and its size is limited by $2^{j+1}$). Thus, a leader has to send $O(\sqrt{n}/2^{i+1} + \log n) \cdot 2^{i+2} = O(\sqrt{n} \log n)$ messages. Each message is destined to a leader of some fragment, which is located at the randomly chosen node in $\mathcal{C}$. So, again, by Lemma 1, we have that the most loaded edge carries $O(\log n)$ messages with high probability.

In case $2^{i+1} > \sqrt{n}$, Lemma 1 yields that a leader $u \in \mathcal{C}$ has at most $O(\log n)$ roots (at most $n/2^{i+1}$ balls into at least $\sqrt{n}$ bins). Since a root has to send at most one message to each leader (even if the node is a leader of multiple fragments of the tree), the total number of messages needed to be sent by a leader is $O(n_c \log n)$. Since every message is destined to a random leader, by Lemma 1 we obtain a bound of $O(\log n)$ on the maximum edge load, with high probability.

Overall, step 4 takes $O(\log n)$ time. Step 5 obviously takes the same time ($O(\log n)$) as step 3, since it involves the transfer of the same amount of information (except in the opposite direction). Step 6 takes the same time ($O(1)$) as step 2, since again it involves transferring the same amount of information (in the opposite direction). □

We have established the following theorem.

**Theorem 4.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, Algorithm $\mathcal{CP}$-MST constructs a MST in $O(\log^2 n)$ rounds, with high probability.*

## 4. Additional Algorithms in Core-Periphery Networks

In addition to MST, we have considered a number of other distributed problems of different types, and developed algorithms for these problems that can be efficiently executed on core-periphery networks. In particular, we dealt with the following set of tasks, related to matrix operations: (M1) Sparse matrix transposition. (M2) Multiplication of a sparse matrix by a vector. (M3) Multiplication of two sparse matrices.

We then considered problems related to calculating aggregate functions of initial values, initially stored one at each node in $V$. In particular, we developed efficient algorithms for the following problems: (A1) Finding the rank of each value, assuming the values are ordered. (As output, each node should know the rank of the element it stores.) (A2) Finding the median of the values. (A3) Finding the (statistical) mode, namely, the most frequent value. (A4) Finding the number of distinct values stored in the network. Each of these problems requires $\Omega(D)$ rounds on general networks of diameter $D$. We show that on a $\mathcal{CP}$-network these tasks can be performed in $O(1)$ rounds.

An additional interesting task is defined in a setting where the initial values are split into disjoint groups, and requires finding the $r$ largest values of each group. This task can be used, for example, for finding the most popular headlines in each area of news. Here, there is an $O(r)$ round solution on a $\mathcal{CP}$-network. (The diameter is a lower bound for this task in general networks.)

In all of these problems, we also establish the necessity of all 3 axioms, by showing that there are network families satisfying 2 of the 3 axioms, for which the general lower bound holds.

### 4.1. Technical preliminaries

A few definitions are in place. Let $A$ be a matrix, in which each entry $A(i, j)$ can be represented by $O(\log n)$ bits (i.e., it fits in a single message in the CONGEST model). Denote by $A_{i,*}$ (respectively, $A_{*,i}$) the $i$th row (resp., column) of $A$. Denote the $i$th entry of a vector $s$ by $s(i)$. We assume that the nodes in $\mathcal{C}$ have IDs $[1, \ldots, n_c]$, and this is known to all of them. A square $n \times n$ matrix $A$ with $O(k)$ nonzero entries in each row and each column is hereafter referred to as an $O(k)$-*sparse matrix*.

Our algorithms make extensive use of the following theorem of [24].

**Theorem 5** ([24])**.** *Consider a fully connected network of $n_c$ nodes, where each node is given up to $M_s$ messages to send, and each node is the destination of at most $M_r$ messages. There exists an algorithm* SendMsg *that delivers all the messages to their destinations in the CONGEST model in $O((M_s + M_r)/n_c)$ rounds with high probability.*

This theorem provides a time-efficient procedure for messages delivery in a core that satisfies Axiom $\mathcal{A}_E$. Note that the result of the theorem holds with high probability, which implies that it exploits a randomized algorithm. Nevertheless, our algorithms, presented below, can be considered "mostly deterministic", in the sense that all the *decisions* they make are deterministic. That

is, the choices concerning which messages should be send where during different stages of our algorithms are made deterministically, and randomness is used only in the information delivery algorithm `SendMsg` of Theorem 5, which is used as a low-level procedure for routing messages from sources to destinations over a complete network. Hence the time bound of each of our algorithms holds with the same probability as those of the calls to `SendMsg` used in that algorithm.

In some of our algorithms, we use the following result on distributed sorting in a complete network, presented in [23].

**Theorem 6** ([23])**.** *Consider a complete network $G(V,E)$ with node ID's $[1,\dots,n]$. Each node is given $n$ values. For simplicity assume all $n^2$ values are distinct[3]. The following tasks can be performed deterministically in $\Theta(1)$ rounds.*

1. *Value learning (VL): Node $i$ needs to learn the values with indices $[i(n-1)+1,\dots,in]$ according to the total order of all values.*

2. *Index Learning (IL): Node $i$ needs to determine the indices of its input (initial) values in the total order of all values.*

**Observation 1.** *Theorem 6 can be naturally extended to the case where each node initially holds $O(n)$ keys (instead of exactly $n$).*

*4.2. Matrix transposition (MT)*

Initially, each node in $V$ holds one row of an $O(k)$-sparse matrix $A$ (along with its index). The *matrix transposition (MT)* task is to distributively calculate the matrix $A^T$, and store its rows in such a way that the node that stores row $A_{i,*}$ will eventually store row $A_{i,*}^T$. We start with a lower bound.

**Theorem 7.** *Any algorithm for transposing an $O(k)$-sparse matrix on an arbitrary network of diameter $D$ requires $\Omega(D)$ rounds. On a $\mathcal{CP}$-network, $\Omega(k)$ rounds are required.*

*Proof.* Consider a nonzero entry $A(i,j)$, where $j \neq i$. Consider the nodes $u$ and $v$ that initially store $A_{i,*}$ and $A_{j,*}$ respectively. Clearly, in any algorithm for MT, $A(i,j)$ should be delivered to the node $v$ (which is required to eventually obtain $A_{*,j} = A_{j,*}^T$). Since the distance $dist(u,v)$ may be as large as the diameter, the lower bound is $\Omega(D)$ rounds.

For a $\mathcal{CP}$-network, the lower bound on MT is $\Omega(k)$, since there are inputs for which row $A_{i,*}^T$ has $k$ nonzero values, which must be delivered to the node that initially has row $A_{i,*}$. There are $\mathcal{CP}$-networks whose minimum degree is 1 (see Figure 1(I) for an illustration) and hence delivering $\Omega(k)$ messages will require $\Omega(k)$ communication rounds. $\qquad\square$

---

[3]This limitation can be eliminated by chaining each value, with the node ID and its order at the node. Thus, each input value becomes unique.

**Algorithm 1.** MT generating $A^T$ on a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$.
(1) Each $u \in V$ sends its row (all the nonzero values with their indices in $A$) to its representative $r(u) \in \mathcal{C}$. Now, each representative has $O(n_c)$ rows of $A$ (or, $O(kn_c)$ entries of $A$).
(2) Each representative sends each entry $A(i, j)$ it has to the node in $\mathcal{C}$ that is responsible for the row $A^T_{j,*}$. Every node in $\mathcal{C}$ is responsible for the rows of $A^T$ indexed $1 + (n/n_c)(i-1), \ldots, (n/n_c)i$. (Assume $n/n_c$ is integral.)
[$*$ Now, each node in $\mathcal{C}$ stores $O(n/n_c)$ rows of $A^T$. $*$]
(3) Each node $u \in V$ that initially stored the row $i$ of $A$, requests $A^T_{i,*}$ from its representative. The representative gets the row from the corresponding node in $\mathcal{C}$, and sends it back to $u$.

**Theorem 8.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, transposing an $O(k)$-sparse matrix can be completed in $O(k)$ rounds with high probability.*

*Proof.* Consider Algorithm 1 and the $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$. Step 1 of the algorithm will take $O(k)$ rounds, since each row has up to $k$ nonzero entries and sending one entry takes $O(1)$ due to Axiom $\mathcal{A}_C$. Now each representative has $O(kn_c)$ values, since it represents up to $O(n_c)$ nodes in $\mathcal{P}$ (due to Axiom $\mathcal{A}_B$).

In the beginning of Step 2, each representative knows the destination for each of the $A(i, j)$ entries it has (since, by agreement, each node in $\mathcal{C}$ is responsible for collecting entries for specific rows of $A^T$). So, it will send $O(kn_c)$ messages, each one to a specific single destination. Since each node in $\mathcal{C}$ is responsible for $O(n/n_c)$ rows of $A^T$, it will receive $O(kn/n_c)$ messages. Thus, using Axiom $\mathcal{A}_E$ and Theorem 5, the running time is $O(k)$.

At Step 3, a single row (with $O(k)$ nonzero entries) is sent, by each node, to its representative (which takes $O(k)$ time, due to the Axiom $\mathcal{A}_C$). Then the requests are delivered to the appropriate nodes in $\mathcal{C}$, and the replies with the appropriate rows of $A^T$ are received back by the representatives. All this takes $O(k)$ rounds, due to Axiom $\mathcal{A}_E$ and Theorem 5. Then the rows of $A^T$ are delivered to the nodes that have requested them. Due to the Axiom $\mathcal{A}_C$ this will also take $O(k)$ rounds. $\square$

We now show the necessity of the Axioms $\mathcal{A}_B$, $\mathcal{A}_E$ and $\mathcal{A}_C$ for achieving $O(k)$ running time.

**Theorem 9.** *For each $X \in \{B, E, C\}$ there exist a family of $n$-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$ that satisfy all axioms except $\mathcal{A}_X$, and input matrices of size $n \times n$ for every $n$, such that the time complexity of any matrix transposition (MT) algorithm on the networks of $\mathcal{F}_i$, with the corresponding inputs, is $\Omega(n)$.*

*Proof.* Consider the following cases where in each case, one of the axioms is not satisfied, while the other two are satisfied.
Necessity of $\mathcal{A}_B$: Consider the family of dumbbell partitioned networks $D_n$. As discussed earlier, Axiom $\mathcal{A}_B$ is violated, while the others hold. Let $A$ be a matrix where, at least, the following $n/2$ entries are nonzero (assume $n/2$

is even): $A(n/2 + 1, 1), A(n/2 + 2, 2), \ldots, A(n, n/2)$. Then we input the rows $A_{1,*} - A_{n/2,*}$ to the nodes in the first star of of $D_n$, and rows $A_{n/2+1,*} - A_{n,*}$ to the nodes in the second star of $D_n$. Clearly, the entries we specified before are initially located in the second star, but they all must be delivered to the first star (by the problem definition, an entry $A(i, j)$ should eventually be stored in a node that initially has row $A_{j,*}$). Since there is only one edge connecting the stars, any distributed algorithm for the specified task will take $\Omega(n)$ rounds.

Necessity of $\mathcal{A}_E$: Consider the family of sun partitioned networks $S_n$. As discussed earlier, Axiom $\mathcal{A}_E$ is violated, while the others hold. The diameter of $S_n$ is $\Omega(n)$, hence, any distributed MT algorithm will require $\Omega(n)$ communication rounds (due to the lower bound discussed earlier).

Necessity of $\mathcal{A}_C$: Consider the family of lollipop partitioned networks $L_n$. As discussed earlier, Axiom $\mathcal{A}_C$ is violated, while the others hold. Again, the diameter of $L_n$ is $\Omega(n)$, hence, any distributed MT algorithm requires $\Omega(n)$ rounds. □

### 4.3. Vector by matrix multiplication (VMM)

Let $s$ be a vector of size $n$, and $A$ be a square $n \times n$ $O(k)$-sparse matrix. Initially, each node in $V$ holds one entry of $s$ (along with its index), and one row of $A$ (along with its index). The *vector by matrix multiplication (VMM)* task is to distributively calculate the vector $s' = sA$, and store its entries at the corresponding nodes in $V$, such that the node that initially stored $s(i)$ will store $s'(i)$. We start with the lower bound.

**Theorem 10.** *Any algorithm for the multiplication of a vector by an $O(k)$-sparse matrix (VMM) on any network requires $\Omega(D)$ rounds. On a $\mathcal{CP}$-network, $\Omega(k/\log n)$ rounds are required.*

*Proof.* The $\Omega(D)$ time lower bound for VMM on an arbitrary network follows since in order to obtain $s'(1)$, we need, at least, to multiply $s(1)$ by $A(1, 1)$ (assuming $s(1) \neq 0$ and $A(1, 1) \neq 0$), which might take $\Omega(D)$ rounds in case $s(1)$ and $A(1, 1)$ are located at different nodes $u$ and $v$ at distance $dist(u, v) = D$.

Now we show that there exists a $\mathcal{CP}$-network for which the lower bound on VMM is $\Omega(k/\log n)$ rounds. Consider a $\mathcal{CP}$-network, as in Figure 1(I). Let $u$ be a node in $\mathcal{P}$ (whose degree is 1). Let $v$ be any other node in $V$. Assume that $u$ initially stores the row $A_{1,*}$,, and the entry $s(1)$, while $v$ stores row $A_{2,*}$ and the entry $s(2)$.

Next, we show a reduction from the well-known *equality problem* (EQ), in which two parties are required to determine whether their input vectors $x, y \in \{0, 1\}^k$ are equal. Assuming the existence of a procedure $P$ for our VMM problem, we use it to solve the EQ problem. Given input vectors $x, y$ for the EQ problem (at $u$ and $v$ respectively), we create an input for the VMM problem in the following way. Node $u$ assigns $A(1, i) = x(i)$, for every $i \in [1, \ldots, k]$ and $s(1) = 1$; while node $v$ assigns $A(2, i) = y(i)$ for every $i \in [1, \ldots, k]$ and $s(2) = 1$. All the other entries of $A$ and $s$ are initialized to 0. It follows that $s'(i) = \sum_{j=1}^{n} s(j)A(j, i) = A(1, i) + A(2, i) = x(i) + y(i)$ for every

$i \in [1, \dots, k]$. Given the value of $s'(i)$, one can decide whether $x(i) = y(i)$ for every $i \in [1, \dots, k]$, since clearly, $x(i) = y(i)$ if'f $s'(i) \in \{0, 2\}$ (and otherwise $s'(i) = 1$). Notice that the vector $s'$ is stored distributedly in the network, one entry in each node. But the indication to $v$ and $u$, whether all the entries are in $\{0, 2\}$, can be delivered in $O(1)$ rounds in the following way. Each node in $\mathcal{P}$ sends its entry of $s'$ to its representative, who checks all the received entries and sends an indication bit to all the other nodes in $\mathcal{C}$. So, every node in $\mathcal{C}$ knows now whether all the entries in $s'$ are in $\{0, 2\}$ (actually, we are interested only in the first $k$ entries). Representatives can now inform the nodes in $\mathcal{P}$ they represent in $O(1)$ rounds. It follows that using procedure $P$, one can solve the EQ problem, which is known to require at least $k$ bits of communication. Therefore, assuming that each message has $O(\log n)$ bits, our problem requires $\Omega(k/\log n)$ communication rounds. $\qquad\square$

**Algorithm 2.** VMM on a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$.
(1) Each $u \in V$ sends the entry of $s$ it has (along with its index) to its representative $r(u) \in \mathcal{C}$ (recall that if $u \in \mathcal{C}$ then $r(u) = u$).
(2) $\mathcal{C}$ nodes redistribute the $s$ entries among them, so that the node with ID $i$ stores indices $[1 + (n/n_c)(i-1), \dots, (n/n_c)i]$ (assume $n/n_c$ is integral).
(3) Each $u \in V$ sends the index of the row of $A$ it has to its representative $r(u) \in \mathcal{C}$.
(4) Each representative requests the $s(i)$ entries corresponding to rows $A_{i,*}$ that it represents, from the $\mathcal{C}$ node storing it.
(5) Once getting them, it sends them to the nodes in $\mathcal{P}$ it represents.
(6) Each $u \in V$ sends the products $\{A(i,j)s(i)\}_{j=1}^{n}$ to its representative.
(7) Each representative sends each nonzero value $A(i,j)s(i)$ it has (up to $O(kn_c)$ values) to the representative responsible for $s(j)$, so it can calculate $s'(j)$.
(8) Each node $u \in V$, that initially stored $s(i)$, requests $s'(i)$ from its representative. The representative gets the entry from the corresponding node in $\mathcal{C}$ and sends it back to $u$.

**Theorem 11.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, the multiplication of a vector by an $O(k)$-sparse matrix (VMM) can be completed in $O(k)$ rounds with high probability.*

*Proof.* Consider Algorithm 2 and the $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$. At Step 1, due to $\mathcal{A}_B$ and $\mathcal{A}_C$, each representative will obtain $O(n_c)$ entries of $s$ in $O(1)$ rounds. For Step 2, we use Theorem 5 with the parameters $M_s = O(n_c)$ and $M_r = O(n/n_c)$, and thus such a redistribution will take $O((n_c + n/n_c)/n_c) = O(1)$ rounds. At Step 3, due to $\mathcal{A}_B$ and $\mathcal{A}_C$ each representative will obtain $O(n_c)$ row indices of $A$ in $O(1)$ rounds.

For Step 4, we again use Theorem 5 with the parameters $M_s = O(n_c)$ (indices of rows each representative has), $M_r = O(n/n_c)$ (number of entries of $s$ stored in each node in $\mathcal{C}$), and obtain a running time of $O((n_c + n/n_c)/n_c) = O(1)$ rounds for this step. At Step 5, each representative gets the required elements of $s$, which takes $O(1)$ rounds due to Theorem 5, and then sends them

to the nodes in $\mathcal{P}$ it represents, which also takes $O(1)$ rounds due to $\mathcal{A}_C$. Step 6 takes $O(k)$ rounds, since $A$ has up to $k$ nonzero entries in each row. Step 7 again uses Theorem 5 with parameters $M_s = O(kn_c)$, $M_r = O(n/n_c)$, and thus its running time is $O(kn/n_c^2) = O(k)$.

At Step 8, a single message is sent by each node to its representative (which takes $O(1)$ rounds due to $\mathcal{A}_C$), then the requests are delivered to the appropriate nodes in $\mathcal{C}$, and the replies with the appropriate entries of $s'$ are received back by the representatives. All this takes $O(1)$ rounds, due to Axiom $\mathcal{A}_E$ and Theorem 5. Then the entries of $s'$ are delivered to the nodes that have requested them. By $\mathcal{A}_C$ this also takes $O(1)$ rounds. □

The following theorem shows the necessity of the axioms for achieving $O(k)$ running time.

**Theorem 12.** *For each $X \in \{B, E, C\}$ there exists a family of $n$-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$, that satisfy all axioms except $\mathcal{A}_X$, and input matrices of size $n \times n$ and vectors of size $n$, for every $n$, such that the time complexity of* any *algorithm for VMM on the networks of $\mathcal{F}_X$ with the corresponding-size inputs is $\Omega(n/\log n)$.*

*Proof.* Necessity of $\mathcal{A}_B$: Consider the family of dumbbell partitioned networks $D_n$ (which violates only Axiom $\mathcal{A}_B$). Denote the core's nodes as $u$ and $v$. In $O(k)$ rounds $u$ (resp. $v$) can collect all the rows of $A$ and entries of $s$ stored at the nodes of $\mathcal{P}$ connected to $u$ (resp., $v$). So, assuming $n/2$ is an integer, after $O(k)$ rounds, $u$ and $v$ have each $n/2$ rows of $A$ and $n/2$ entries of $s$. Assume also an input, for which $u$ has rows $A_{1,*}, A_{2,*}, \ldots, A_{n/2,*}$, and entries $s(n/2 + 1), s(n/2 + 2), \ldots, s(n)$, and $v$ has all the remaining rows of $A$ and entries of $s$. We again show a reduction from the EQ problem. Assuming the existence of a procedure $P$ for VMM, we use it to solve the EQ problem. Given input vectors $x, y$ for the EQ problem (at $u$ and $v$ respectively), we create an input for VMM in the following way. Node $u$ assigns $A(i, i) = x(i) + 1$ for every $i \in [1, \ldots, n/2]$, and $s(i) = x(i) + 1$ for every $i \in [n/2 + 1, \ldots, n]$; while node $v$ assigns $A(i, i) = y(i) + 1$, for every $i \in [n/2 + 1, \ldots, n]$ and $s(i) = y(i) + 1$ for every $i \in [1, \ldots, n/2]$. All the other entries of $A$ are initialized to 0, thus $A$ is a diagonal matrix. It follows that $s'(i) = \sum_{j=1}^{n} s(j)A(j, i) = s(i)A(i, i) = (x(i)+1)(y(i)+1)$ for every $i$. Given the value of $s'(i)$, one can decide whether $x(i) = y(i)$, since clearly, $x(i) = y(i)$ if'f $s'(i) \in \{1, 4\}$ (and otherwise $s'(i) = 2$). It follows that, using procedure $P$, one can solve the EQ problem, which in this case requires at least $n$ bits of communication. Hence the VMM problem requires $\Omega(n/\log n)$ communication rounds.

The proof of the necessity of $\mathcal{A}_E$ and $\mathcal{A}_C$ is the same as in the proof of Theorem 9, and is based on the diameter argument. □

*4.4. Matrix multiplication (MM)*

Let $A$ and $B$ be square $n \times n$ matrices, with $O(k)$ nonzero entries in each row and each column. Initially, each node in $V$ holds one row of $A$ (along with

its index), and one row of $B$ (along with its index). The *matrix multiplication (MM)* task is to distributively calculate $C = AB$, and store its rows at the corresponding nodes in $V$, such that the node that initially stored row $i$ of $B$ will store row $i$ of $C$.

**Theorem 13.** *Any algorithm for $O(k)$-sparse matrix multiplication on any network requires $\Omega(D)$ rounds. On a $\mathcal{CP}$-network, $\Omega(k^2)$ rounds are required.*

*Proof.* The lower bound for general networks follows since in order to obtain $C(1,1)$ we need, at least, to multiply $A(1,1)$ by $B(1,1)$ (assuming an input in which $A(1,1) \neq 0$ and $B(1,1) \neq 0$), which might take $\Omega(D)$ rounds in case $A(1,1)$ and $B(1,1)$ are located at nodes $u$ and $v$, at distance $dist(u,v) = D$.

For a $\mathcal{CP}$-network, consider a network illustrated on Figure 1(I), where the degree of a node $u \in \mathcal{P}$ is 1. Assume that initially $u$ has row $A_{1,*}$ and $B_{1,*}$ and thus, by problem definition, it has to eventually receive the row $C_{1,*}$. We show now that $\Omega(k^2)$ messages are required to allow $u$ to obtain $C_{1,*}$. Assume that $\{b_i^j\}$ for $i,j \in [1,\dots,k]$ are $k^2$ distinct values. Consider an input in which $A(1,i) = 1$ for every $i \in [2,\dots,k+1]$, $B(2,i) = b_i^1$ for every $i \in [1,\dots,k]$, $B(3,i) = b_i^2$ for every $i \in [k+1,\dots,2k]$; and so on until $B(k+1,i) = b_i^k$ for every $i \in [k(k-1)+1,\dots,k^2]$. All other entries of $A$ and $B$ are set to 0. It is easy to see that $C(1,i) = b_i^1$ for every $i \in [1,\dots,k]$, $C(1,i) = b_i^2$ for every $i \in [k+1,\dots,2k]$; and so on until $C(1,i) = b_i^k$ for every $i \in [k(k-1)+1,\dots,k^2]$. So, in order to obtain the row $C_{1,*}$, $u$ must receive all the $k^2$ values $\{b_i^j\}$, which will take $\Omega(k^2)$ communication rounds. $\square$

**Algorithm 3.** $O(k)$-sparse matrix multiplication on a $\mathcal{CP}$-network $G(V,E,\mathcal{C},\mathcal{P})$.
(1) Each node in $V$ send its row of $B$ to its representative.
[∗ Now, each node in $\mathcal{C}$ has $O(kn_c)$ entries of $B$. ∗]
(2) Nodes in $\mathcal{C}$ redistribute the rows of $B$ among themselves, so that the node with ID $i$ will store rows $1 + (n/n_c)(i-1),\dots(n/n_c)i$ (assuming $n/n_c$ is an integer).
[∗ Now, each $u \in \mathcal{C}$ has $O(n/n_c)$ rows of $B$. ∗]
(3) Each node in $V$ sends its row of $A$ to its representative. Notice that the row $i$ of $B$ needed to be multiplied only by values of the column $i$ of $A$.
(4) Each $u \in \mathcal{C}$ sends the values of $A$ it has to the nodes in $\mathcal{C}$, which hold the corresponding rows of $B$. I.e., the value $A(i,j)$ will be sent to the node in $\mathcal{C}$, which holds the row $B_{j,*}$.
[∗ Now all the summands are prepared and distributed across all the nodes in $\mathcal{C}$, and it is left to combine corresponding summands. ∗]
(5) Each $u \in \mathcal{C}$ sends each of its values to the corresponding node in $\mathcal{C}$, that is responsible for gathering the summands for the values of specific $O(n/n_c)$ rows of the resulting matrix $C$.
(6) Each node $u \in V$ that initially stored row $i$ of $B$, requests row $i$ of $C$ from its representative. The representative gets the row from the corresponding node in $\mathcal{C}$, and sends it back to $u$.

**Theorem 14.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, the multiplication of two $O(k)$-sparse matrices can be completed in $O(k^2)$ rounds with high probability.*

*Proof.* Consider Algorithm 3 and the network $G(E, V)$, with a partition $\langle \mathcal{C}, \mathcal{P} \rangle$ that satisfies Axioms $\mathcal{A}_B$, $\mathcal{A}_E$, $\mathcal{A}_C$.

Step 1 takes $O(k)$ rounds, due to the Axiom $\mathcal{A}_C$ and the fact that the number of nonzero entries in each row is bounded by $k$. Each node in $\mathcal{C}$ will have $O(kn_c)$ entries of $B$, since it represents $O(n_c)$ nodes in $\mathcal{P}$ due to the Axiom $\mathcal{A}_B$. Using Theorem 5 with the parameters $M_s = O(kn_c)$ and $M_r = O(kn/n_c)$, we show that the redistribution, performed at Step 2, takes $O((kn_c + kn/n_c)/n_c) = O(k)$ rounds.

For Step 4, we again use Theorem 5, with the parameters $M_s = O(k)O(n_c)$ ($O(k)$ values per row), $M_r = O(n/n_c)O(k)$ ($O(k)$ values per column of $A$), and obtain that the running time is $O(k)$. Note that each node in $\mathcal{C}$ has $O(n/n_c)$ rows (of $B$), with $O(k)$ elements in each. Each row was multiplied by $O(k)$ values received in the previous step ($O(k)$, since each column of $A$ has $O(k)$ nonzero values). Thus, each $u \in \mathcal{C}$ has $O(k^2 n/n_c)$ values that needed to be sent to the appropriate nodes in the next step.

At Step 5, each $u \in \mathcal{C}$ sends each of its values to the corresponding node in $\mathcal{C}$ that is responsible for gathering the summands for the values of specific $O(n/n_c)$ rows of the resulting matrix $C$. Clearly, $M_s = O(k^2 n/n_c)$ since each row of $B$ has $O(k)$ different entries, and was multiplied by $O(k)$ different entries of $A$. Now let's find $M_r$. Each $u \in \mathcal{C}$ is responsible for $O(n/n_c)$ rows of $C$. Thus, e.g., for a row $C_{1,*}$ it needs to receive the following summands: $C(1, 1) = \sum A(1, i)B(i, 1)$, $C(1, 2) = \sum A(1, i)B(i, 2)$, ..., $C(1, n) = \sum A(1, i)B(i, n)$. Since the number of nonzero entries in each row of $A$ and $B$ is $O(k)$, the number of nonzero entries in each row of $C$ is bounded by $O(k^2)$. Thus, for each row of $C$, a node in $C$ will receive $O(k^2)$ messages. So, $M_r = O(k^2 n/n_c)$ and thus the running time of this step is $O(k^2)$.

At the last step, each node $u \in V$ sends a single message (request for a row) to its representative. This takes $O(1)$ rounds due to $\mathcal{A}_C$. Then, the representative gets the row from the corresponding node in $\mathcal{C}$, and sends it back to $u$. Using Axiom $\mathcal{A}_E$ and Theorem 5 with $M_s = O(n_c)$ and $M_r = O(n/n_c)$, delivering the request inside the core takes $O(1)$ rounds. In a similar way, sending the rows inside the core takes $O(k^2)$ rounds. The same amount of time is required to deliver those rows to the nodes in $\mathcal{P}$ that requested them ($O(1)$ per row entry, due to $\mathcal{A}_C$, and $O(k^2)$ nonzero entries per row). $\square$

**Theorem 15.** *For each $X \in \{B, E, C\}$ there exist a family of $n$-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$, that satisfy all axioms except $\mathcal{A}_X$, and input matrices of size $n \times n$, for every $n$, such that the time complexity of any MM algorithm for the multiplication for $O(k)$-sparse matrices on the networks of $\mathcal{F}_X$ with the corresponding inputs is $\Omega(n/\log n)$.*

*Proof.* Necessity of $\mathcal{A}_B$: Again consider the dumbbell partitioned networks $D_n$. As done in the proof of Theorem 12, we can show a reduction from the EQ

problem to the MM problem for sparse matrices. Here we initialize $A$ and $B$ to be diagonal matrices, and the first core node, $u$, will have first half of $A$'s rows and second half of $B$'s rows. Due to the initialization, each entry of the resulting matrix $C$ will be a multiplication of the corresponding entries of $x$ and $y$. Thus, obtaining $C$ will allow us to determine whether $x$ and $y$ are equal. Hence the lower bound for the MM problem is again $\Omega(n/\log n)$ communication rounds. The proof for $\mathcal{A}_E$ and $\mathcal{A}_C$ is the same as in the proof of Theorem 9. $\square$

*4.5. Rank Finding (RF)*

Next, we show another set of algorithms that deals with aggregate functions. Let each node $v \in V$ hold some initial value. The *rank finding (RF)* task requires each node to know the position of its value in the ordered list of all the values.

**Theorem 16.** *Any algorithm for finding the rank of each value on any network requires $\Omega(D)$ rounds.*

*Proof.* For two nodes $u, v \in V$ to decide whose value is larger, at least one bit of information must travel between them. Thus, the lower bound for this task on any graph is $\Omega(D)$. $\square$

**Algorithm 4.** RF on a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$.
(1) Each node in $V$ sends its initial value to its representative in $\mathcal{C}$.
(2) Nodes in $\mathcal{C}$ sort all the values they have, and obtain the ranks of those values.
[∗ Now, each node $u \in \mathcal{C}$ knows the ranks of the values it had received from the nodes it represents. ∗]
(3) The ranks are delivered to the appropriate nodes in $\mathcal{P}$.

**Theorem 17.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, the RF task can be completed in $O(1)$ rounds with high probability.*

*Proof.* Consider Algorithm 4 and the $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$. The first step takes $O(1)$, due to Axiom $\mathcal{A}_C$. Now, each representative $u \in \mathcal{C}$ has $O(n_c)$ values, due to Axiom $\mathcal{A}_B$. Step 2 is performed in $O(1)$ rounds, due to Theorem 6 and Axiom $\mathcal{A}_E$. The last step takes $O(1)$ rounds, due to $\mathcal{A}_C$. $\square$

**Theorem 18.** *For each $X \in \{B, E, C\}$ there exist a family of $n$-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$ that satisfy all axioms except $\mathcal{A}_X$, and input matrices of size $n \times n$ and vectors of size $n$, for every $n$, such that the time complexity of any RF algorithm on the networks of $\mathcal{F}_X$ with the corresponding inputs is $\Omega(n)$.*

*Proof.* Necessity of $\mathcal{A}_B$: Consider the family of dumbbell partitioned networks $D_n$. Consider a sorted list of $n$ values ($a_i \leq a_{i+1}$): $a_1, a_2, a_3, ..., a_n$. Clearly, every two adjacent values must be compared in any sorting procedure. For example, assuming that some sorting procedure does not compare $a_2$ and $a_3$,

and also $a_4 = a_3 + 2$, we can replace $a_2$ with $a_3 + 1$ and get the same output $a_1, a_2, a_3, ..., a_n$, which is now incorrect since $a_2 > a_3$. Denote the two core nodes by $u$ and $v$. Node $u$ and all the nodes in $\mathcal{P}$ that are connected to it, will be assigned values (one per node) with odd indices in the ordered list (i.e., $a_1, a_3, \ldots$). The other $n/2$ values will be assigned to the remaining $n/2$ nodes (one value to one node). Note that in order to obtain a sorted list, at least the following comparisons must take place: $a_1$ with $a_2$, $a_3$ with $a_4$, and so on. Hence about $n/2$ pairs of values have to be brought together, while initially they are located at the different sides of the link $(u, v)$. Thus, $\Omega(n)$ messages must be sent over the link $(u, v)$ in order to perform the sorting task. This takes $\Omega(n)$ communication rounds.

The proof for $\mathcal{A}_E$ and $\mathcal{A}_C$ is the same as in Thm. 9. $\square$

*4.6. Median Finding (MedF)*

Let each node $v \in V$ hold some initial value. The *median finding (MedF)* task requires each node to know the value, which is the median in the ordered list of all the values in the network.

**Theorem 19.** *Any median finding algorithm on any network requires $\Omega(D)$ rounds.*

*Proof.* For two nodes $u, v \in V$, if $u$ wants to obtain the median of the initial values, then at least one bit of information must travel between $u$ and $v$ (otherwise $u$ will never know that $v$ exists, while $v$ may even be the median). Since the distance $dist(u, v)$ can be as large as the diameter, the lower bound for this task is $\Omega(D)$ communication rounds. $\square$

**Algorithm 5.** MedF on a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$.
(1) Each node in $V$ sends its initial value to its representative in $\mathcal{C}$.
(2) Nodes in $\mathcal{C}$ sort all their values and obtain the ranks of those values.
[* Now, each node $i \in \mathcal{C}$ now knows values with indices $i(n_c - 1) + 1, \ldots in_c$ according to the total order of all values. *]
(3) The node in $\mathcal{C}$ that has the value with index $n^2/2$ (namely, the median value) sends it to all the nodes in $\mathcal{C}$.
(4) The median value is delivered to the nodes in $\mathcal{P}$.

**Theorem 20.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, the MedF task can be completed in $O(1)$ rounds with high probability.*

*Proof.* Consider Algorithm 5 and the $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$. The first step takes $O(1)$ due to Axiom $\mathcal{A}_C$. Now, each representative $u \in \mathcal{C}$ has $O(n_c)$ values due to Axiom $\mathcal{A}_B$. Step 2, is performed in $O(1)$ rounds, due to Theorem 6 and Axiom $\mathcal{A}_E$. The last step will take $O(1)$ rounds, due to $\mathcal{A}_C$. $\square$

**Theorem 21.** *For each $X \in \{B, E, C\}$ there exist a family of $n$-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$, that satisfy all axioms except $\mathcal{A}_X$, and input matrices of size $n \times n$, and vectors of size $n$, for every $n$, such that the time complexity of* any *algorithm for finding median of all the initial values on the networks of $\mathcal{F}_X$, with the corresponding inputs is $\Omega(\log n)$.*

*Proof.* Necessity of $\mathcal{A}_B$: Consider the family of dumbbell partitioned networks $D_n$. As discussed earlier, Axiom $\mathcal{A}_B$ is violated, while the others hold. In constant time, each of the two centers, A and B, can learn the inputs of its star. Now, the problem becomes for A and B to learn the median of the union of their sets. This operation is known to require at least $\Omega(\log n)$ communication rounds. More formally, it is shown in [36], that the Median function does not admit a deterministic protocol using $O(\log^{1-\epsilon} n)$ rounds, and a logarithmic amount of communication at each round, for any $\epsilon > 0$ (even though the total communication complexity of the problem is known to be only $O(\log n)$ bits). This allows us to conclude the same lower bound in our case too.

The proof for $\mathcal{A}_E$ and $\mathcal{A}_C$ is the same as in Thm. 9. $\qquad\square$

*4.7. Mode Finding (ModF)*

Let each node $v \in V$ hold some initial value. The *mode finding (ModeF)* task requires each node to know the value (values) that appears most frequently.

**Theorem 22.** *Any mode finding algorithm on any network requires $\Omega(D)$ rounds.*

*Proof.* For three nodes $u, v, w \in V$, assume an input, for which the most frequent value appears with frequency 2, and is initially located at $v$ and $w$, while all the other nodes have other distinct values. Obviously, if $u$ needs to learn the mode, at least one bit of information must travel from $v$ to $u$, since otherwise $u$ will not be aware of the $v$'s existence. Since the distance $dist(u, v)$ can be as large as diameter, the lower bound for this task is $\Omega(D)$ communication rounds. $\qquad\square$

**Algorithm 6.** ModF on a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$.
(1) Each node in $V$ sends its initial value to its representative in $\mathcal{C}$.
(2) Nodes in $\mathcal{C}$ perform sorting of all the values they have, and obtain the ranks of those values. So, each node $i \in \mathcal{C}$ now knows values with indices $i(n_c - 1) + 1, \ldots in_c$, according to the total order of all values.
(3) Each node in $\mathcal{C}$ sends to each other node in $\mathcal{C}$: (i) its most frequent value (values) and its frequency, (ii) its minimum value and its frequency, (iii) its maximum value and its frequency.
[∗ The minimum and maximum are needed in order capture the most frequent values that appear at the boundaries of the ranges. ∗]
(4) Each node in $\mathcal{C}$ finds the most frequent value (values), and delivers it to the nodes in $\mathcal{P}$ it represents.

**Theorem 23.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, the ModF task can be completed in $O(k)$ rounds with high probability.*

*Proof.* Consider Algorithm 6 and the $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$. The first step takes $O(1)$ due to Axiom $\mathcal{A}_C$. Now, each representative $u \in \mathcal{C}$ has $O(n_c)$ values due to Axiom $\mathcal{A}_B$. Step 2, is performed in $O(1)$ rounds, due the Theorem 6 and Axiom $\mathcal{A}_E$. Step 3 takes $O(1)$, due to Axiom $\mathcal{A}_E$ since each node in $\mathcal{C}$ needs to send $O(1)$ values to all the other nodes in $\mathcal{C}$. The last step will take $O(1)$ rounds due to $\mathcal{A}_C$ (and assuming there are $O(1)$ most frequent values). $\square$

**Theorem 24.** *For each $X \in \{B, E, C\}$ there exist a family of $n$-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$ that satisfy all axioms except $\mathcal{A}_X$, and input matrices of size $n \times n$ and vectors of size $n$, for every $n$, such that the time complexity of any algorithm for the finding mode on the networks of $\mathcal{F}_X$, with the corresponding inputs is $\Omega(n/\log n)$.*

*Proof.* Necessity of $\mathcal{A}_B$: Consider the family of dumbbell partitioned networks $D_n$. As discussed earlier, Axiom $\mathcal{A}_B$ is violated while the other hold. Assume such an input that every element appears exactly once or twice on each side (for simplicity, assume there are $n/4$ types of elements altogether, and some nodes do not have any element). Hence, the most frequent element will appear 2, 3 or 4 times in the graph. The case where the answer is 2 occurs only when every element appears exactly once on every side. This case is easy to check in a constant amount of communication between the two centers, so we assume we do this check first, and it remains to consider the case where this does not happen. It thus remains to decide whether the answer is 3 or 4. To do that, A (the center of the first star) defines a set $S_A$ of all the elements that appear twice in its star, and B defines a set $S_B$ similarly for its side. Now the answer is 4 if'f the sets $S_A$ and $S_B$ intersect. Set disjointness has communication complexity $n$, so A and B must exchange at least $n$ bits, or, at least $\Omega(n/\log n)$ messages. Since these messages all cross the single edge connecting the two centers, they require this much time.

The proof for $\mathcal{A}_E$ and $\mathcal{A}_C$ is the same as in Thm. 9. $\square$

*4.8. Finding the number of distinct values (DF)*

Let each node $v \in V$ hold some initial value. The *number of distinct values finding (DF)* requires each node to know the number of distinct values present in the network.

**Theorem 25.** *Any DF algorithm on any network requires $\Omega(D)$ rounds.*

*Proof.* For two nodes $u, v \in V$, assume all input values in the network are distinct. Obviously, that if $u$ needs to learn the number of distinct values, at least one bit of information must travel from $v$ to $u$, since otherwise, $u$ will not be aware of the $v$'s existence. Since the distance $dist(u, v)$ can be as large as diameter, the lower bound for this task is $\Omega(D)$ communication rounds. $\square$

**Algorithm 7.** DF on a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$.

(1) Each node in $V$ sends its initial value to its representative in $\mathcal{C}$.

(2) Nodes in $\mathcal{C}$ perform sorting of all the values they have, and obtain the ranks of those values.

[∗ Now, each node $i \in \mathcal{C}$ now knows values with indices $i(n_c - 1) + 1, \ldots in_c$ according to the total order of all values. ∗]

(3) Each node in $\mathcal{C}$ sends to every other node in $\mathcal{C}$, the number of distinct values and the two border values (min,max).

[∗ Now, every node in $\mathcal{C}$ is able to find the number of distinct values (for each repeated border value, decrease 1 from the total count). ∗]

(4) Each representative delivers the number of distinct values to the nodes in $\mathcal{P}$ it represents.

**Theorem 26.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, the DF task requires $O(1)$ rounds with high probability.*

*Proof.* Consider Algorithm 7 and the $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$. The first step takes $O(1)$ due to Axiom $\mathcal{A}_C$. Now, each representative $u \in \mathcal{C}$ has $O(n_c)$ values due to the Axiom $\mathcal{A}_B$. Step 2, is performed in $O(1)$ rounds, due the Theorem 6 and Axiom $\mathcal{A}_E$. Step 3 takes $O(1)$, due to Axiom $\mathcal{A}_E$ since each node in $\mathcal{C}$ needs to send $O(1)$ values to all the other nodes in $\mathcal{C}$. The last step will take $O(1)$ rounds, due to $\mathcal{A}_C$. $\square$

**Theorem 27.** *For each $X \in \{B, E, C\}$ there exist a family of n-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$, that satisfy all axioms except $\mathcal{A}_X$, and input matrices of size $n \times n$ and vectors of size $n$, for every $n$, such that the time complexity of* any *algorithm for finding the number of distinct values on the networks of $\mathcal{F}_X$, with the corresponding inputs is $\Omega(n/\log n)$.*

*Proof.* Necessity of $\mathcal{A}_B$: Consider the family of dumbbell partitioned networks $D_n$. As discussed earlier, Axiom $\mathcal{A}_B$ is violated, while the others hold. Assume that the inputs are taken out of a range of $m$ distinct possible values. In constant time, the two star centers A and B can collect $m$-bit vectors $x$ and $y$ respectively, representing the values that exist in their respective stars. The goal is for A and B to decide the number of distinct values in the graph, i.e., the number of 1's in the vector $x \vee y$. We show a reduction from set disjointness to this problem, hence, the lower bound for set disjointness holds for it. Assume we are given a procedure $P$ for our problem. We use it to solve set disjointness as follows. (1) A computes $|x|$ and informs B. (2) B computes $|y|$ and informs A. (3) A and B invoke procedure $P$ and compute $|x \vee y|$. (4) The answer is "yes" (the sets are disjoint) iff $|x \vee y| = |x| + |y|$. It is easy to verify that the reduction is correct, hence, we get the desired lower bound of $\Omega(m/\log n)$ on the number of round required for finding the number of distinct values.

The proof for $\mathcal{A}_E$ and $\mathcal{A}_C$ is the same as in Thm. 9. $\square$

*4.9. Get the top k ranked by areas (TopK)*

Let each node $v \in V$ hold some initial value. Assume that each value is assigned to a specific area of a total $\sqrt{n}$ areas. E.g., values may represent news and areas topics, so that each news belongs to a specific topic. Assume also, that each node in $V$ is interested in one specific area. The Getting the top $k$ ranked values by areas (TopK) task is to deliver to each node in $V$ the largest $k$ values, from the area it is interested in.

**Theorem 28.** *Any TopK algorithm on any network requires $\Omega(D)$ rounds. On a $\mathcal{CP}$-network $\Omega(k)$ rounds are required.*

*Proof.* First, let us show the lower bound for any network. For two nodes $u, v \in V$, assume input in which $u$ is interested in the value initially stored at $v$. Obviously, delivering the value from $v$ to $u$ will take at least $dist(u, v)$. Since the distance $dist(u, v)$ may be as large as diameter, the lower bound on the running time is $\Omega(D)$.

For a $\mathcal{CP}$-network, the lower bound is $\Omega(k)$, since obviously, there are inputs for which $k$ values must be delivered to a node in $\mathcal{P}$. There are $\mathcal{CP}$-networks, in which minimum degree is 1 (see Figure 1(I) for an illustration) and hence delivering $\Omega(k)$ messages will require $\Omega(k)$ communication rounds. $\square$

**Algorithm 8.** TopK on a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$.

Without loss of generality, assume that all the values are taken from the range $[1, \ldots, n]$, and each area has its own range for its values (e.g., politics $[1, \ldots, 100]$, sports $[101, \ldots, 200]$, etc.).
(1) Each node in $V$ sends its initial value to its representative in $\mathcal{C}$.
(2) Perform sorting using Theorem 6 and Axiom $\mathcal{A}_E$.
[∗ Now, each node $i \in \mathcal{C}$ knows values with indices $i(n_c - 1) + 1, \ldots in_c$ according to the total order of all values. ∗]
(3) Each node in $\mathcal{C}$ sends the largest $k$ values from each area, to the appropriate node in $\mathcal{C}$, so that each node in $\mathcal{C}$ will be responsible for at most one area (recall that we have $\sqrt{n}$ areas and $n_c = \Omega(\sqrt{n})$).
(4) Each representative sends requests for areas (up to $O(n_c)$ areas) requested by nodes it represents. Each request is destined to the specific node in $\mathcal{C}$ responsible for the requested area. Upon request, each node in $\mathcal{C}$ returns the $k$ largest values, for the area it is responsible for, to the requesting nodes.
(5) Each representative $u \in \mathcal{C}$ delivers the values to the nodes in $\mathcal{P}$ it represents.

**Theorem 29.** *On a $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$, the TopK task requires $O(k)$ rounds with high probability.*

*Proof.* Consider Algorithm 8 and the $\mathcal{CP}$-network $G(V, E, \mathcal{C}, \mathcal{P})$. From Theorem 2, we know that $n_c = \Omega(\sqrt{n})$, thus we can say that the number of areas is $O(n_c)$. The first step takes $O(1)$ due to Axiom $\mathcal{A}_C$. Now, each representative $u \in \mathcal{C}$ has $O(n_c)$ values due to the Axiom $\mathcal{A}_B$. At Step 2, each node has $O(n_c)$ values (each destined to a specific single node), so $M_s = n_c$. Each node has

to receive $M_r = k$ values (more precisely: $M_r = 2k$, since it is possible that after the initial sorting, an area is split across two nodes, and each of these two nodes will send up to $k$ values from that area, and the receiving node will have to select the correct $k$). Thus, this step will take (according to Theorem 5) $O((n_c + k)/n_c) = O(k/n_c)$. Step 3 takes $O(k)$, since sending requests will take $O(1)$ (due to the Axiom $\mathcal{A}_E$ and Theorem 5 with $M_s = n_c$, $M_r = n_c$), and receiving the desired values will take $O(k)$ (since $M_s = kn_c$ and $M_r = kn_c$). The last step will take $O(k)$, since each node in $\mathcal{P}$ needs $k$ values, and delivering a single value from $r(u) \in \mathcal{C}$ to $u \in \mathcal{P}$ takes $O(1)$, due to the Axiom $\mathcal{A}_C$. $\square$

**Theorem 30.** *For each $X \in \{B, E, C\}$ there exist a family of $n$-node partitioned networks $\mathcal{F}_X = \{G_X(V, E, \mathcal{C}, \mathcal{P})(n)\}$, that satisfy all axioms except $\mathcal{A}_X$, and input matrices of size $n \times n$ and vectors of size $n$, for every $n$, such that the time complexity of* any *algorithm for finding the $k$ top ranked values from a specific area on the networks of $\mathcal{F}_X$, with the corresponding inputs is $\Omega(kn_c)$.*

*Proof.* Necessity of $\mathcal{A}_B$: Consider the family of dumbbell partitioned networks $D_n$. As discussed earlier, Axiom $\mathcal{A}_B$ is violated, while the others hold. Consider $\sqrt{n}/2$ areas, and assume that all the $k\sqrt{n}/2$ values belonging to these areas are initially located at the nodes of the first star of $D_n$. Consider a subset of nodes of the second star. Let the subset size be $\sqrt{n}/2$ and each node in this subset is interested in different area, of the areas stored in the first star we mentioned. We can see that in order to complete the task, all the $k\sqrt{n}/2$ values have to be sent from the first to the second star, which are interconnected by a single edge. Thus, the running time will be $\Omega(k\sqrt{n})$ rounds.

The proof for $\mathcal{A}_E$ and $\mathcal{A}_C$ is the same as in Thm. 9. $\square$

### References

### References

[1] Adamic, L. (1999). The small world web. *Research and Advanced Technology for Digital Libraries*, pages 852–852.

[2] Avin, C., Lotker, Z., Pignolet, Y. A., and Turkel, I. (2012). From caesar to twitter: An axiomatic approach to elites of social networks. *CoRR*, abs/1111.3374.

[3] Awerbuch, B. (1987). Optimal distributed algorithms for minimum weight spanning tree, counting, leader election and related problems. In *Proc. STOC*, pages 230–240.

[4] Baset, S. and Schulzrinne, H. (2006). An analysis of the skype peer-to-peer internet telephony protocol. In *Proc. INFOCOM*, pages 1–11.

[5] Berns, A., Hegeman, J., and Pemmaraju, S. (2012). Super-fast distributed algorithms for metric facility location. In *Proc. ICALP*, pages 428–439.

[6] Bonanno, G., Caldarelli, G., Lillo, F., and Mantegna, R. (2003). Topology of correlation-based minimal spanning trees in real and model markets. *Phys. Rev. E*, 68.

[7] Bonato, A. (2008). *A course on the web graph*, volume 89. AMS.

[8] Borgatti, S. and Everett, M. (2000). Models of core/periphery structures. *Social networks*, 21(4):375–395.

[9] Chen, C. and Morris, S. (2003). Visualizing evolving networks: Minimum spanning trees versus pathfinder networks. In *Proc. INFOVIS*, pages 67–74.

[10] Chung, F. and Lu, L. (2006). *Complex graphs and networks*. Number 107. AMS.

[11] Dolev, D., Lenzen, C., and Peled, S. (2012). Tri, tri again": Finding triangles and small subgraphs in a distributed setting. *arXiv preprint arXiv:1201.6652*.

[12] Easley, D. A. and Kleinberg, J. M. (2010). *Networks, Crowds, and Markets - Reasoning About a Highly Connected World*. Cambridge Univ. Press.

[13] Elkin, M. (2006). An unconditional lower bound on the time-approximation trade-off for the distributed minimum spanning tree problem. *SIAM J. Computing*, 36(2):433–456.

[14] Feamster, N., Rexford, J., and Zegura, E. (2013). The road to sdn. *Queue*, 11(12):20:20–20:40.

[15] Fujita, M., Krugman, P. R., and Venables, A. J. (2001). *The spatial economy: Cities, regions, and international trade*. MIT press.

[16] Gallager, R., Humblet, P., and Spira, P. (1983). A distributed algorithm for minimum-weight spanning trees. *ACM Trans. on Programming Lang. and Syst.*, 5:66–77.

[17] Garay, J. A., Kutten, S., and Peleg, D. (1998). A sublinear time distributed algorithm for minimum-weight spanning trees. *SIAM J. Computing*, 27:302–316.

[18] Hojman, D. and Szeidl, A. (2008). Core and periphery in networks. *J. Economic Theory*, 139(1):295–309.

[19] Holme, P. (2005). Core-periphery organization of complex networks. *Physical Review E*, 72:046111.

[20] Jung, K., Kim, B., and Vojnovic, M. (2012). Distributed ranking in networks with limited memory and communication. In *Proc. ISIT*, pages 980–984.

[21] Krugman, P. (1991). Increasing Returns and Economic Geography. *J. Political Economy*, 99(3):483–499.

[22] Kutten, S. and Peleg, D. (1998). Fast distributed construction of small k-dominating sets and applications. *J. Algorithms*, 28:40–66.

[23] Lenzen, C. (2013). Optimal deterministic routing and sorting on the congested clique. In *Proc. PODC*, pages 42–50.

[24] Lenzen, C. and Wattenhofer, R. (2011). Tight bounds for parallel randomized load balancing. In *Proc. STOC*, pages 11–20.

[25] Liang, J., Kumar, R., and Ross, K. W. (2006). The fasttrack overlay: A measurement study. *Computer Networks*, 50:842 – 858.

[26] Lotker, Z., Patt-Shamir, B., Pavlov, E., and Peleg, D. (2005). Minimum-weight spanning tree construction in o(log log n) communication rounds. *SIAM J. Computing*, 35(1):120–131.

[27] Lotker, Z., Patt-Shamir, B., and Peleg, D. (2006). Distributed MST for constant diameter graphs. *Distributed Computing*, 18(6):453–460.

[28] Lynch, N. (1995). *Distributed Algorithms*. Morgan Kaufmann.

[29] MacCormack, A. (2010). The architecture of complex systems: do core-periphery structures dominate? In *Proc. Acad. Management*, pages 1–6.

[30] Mitzenmacher, M. and Upfal, E. (2005). *Probability and Computing: Randomized Algorithms and Probabilistic Analysis*. Cambridge Univ. Press.

[31] Nesetril, J., Milkova, E., and Nesetrilova, H. (2001). Otakar boruvka on minimum spanning tree problem translation of both the 1926 papers, comments, history. *Discrete Mathematics*, 233(1 - 3):3 – 36.

[32] Newman, M. (2010). *Networks: an introduction*. Oxford Univ. Press.

[33] Patt-Shamir, B. and Teplitsky, M. (2011). The round complexity of distributed sorting. In *Proc. PODC*, pages 249–256.

[34] Peleg, D. (2000). *Distributed Computing: A Locality-Sensitive Approach*. SIAM.

[35] Peleg, D. and Rubinovich, V. (2000). A near-tight lower bound on the time complexity of distributed minimum-weight spanning tree construction. *SIAM J. Computing*, 30:1427–1442.

[36] Soltys, K. (2011). The hardness of median in the synchronized bit communication model. In *Proc. TAMC*, pages 409–415.

[37] Wyllie, J. (1979). *The complexity of parallel computations*. Technical report. Dept. of Computer Science, Cornell University.

# APPENDIX

## Appendix A. Pseudocodes for Section 3

---

**Algorithm** MERGEFRAGS($i$)

---

Executed every phase by every leader $w \in \mathcal{C}$, for each fragment $i \in F_{lead}(w)$.

1: **if** $state(i) = active$ **then**
2:     $next \leftarrow \text{mp}(i)$
3:     $state(i) \leftarrow frozen$
4:     $[isFound, rootFrag] \leftarrow \text{FINDROOT}(i)$
5:     **if** $isFound = \text{TRUE}$ **then**
6:         $state(rootFrag) \leftarrow root$

7: **if** $state(i) = frozen$ **then**
8:     $[isFinished, next] \leftarrow \text{PJ}(i, next, 2)$
9:     $next(F_{lead}^{\text{mp}(i)}(w)) \leftarrow next$
10:     **if** $isFinished = \text{TRUE}$ **then**
11:         $state(F_{lead}^{\text{mp}(i)}(w)) \leftarrow waiting$

12: **if** $state(i) = waiting$ **then**
13:     receive merge-requests
14:     send merge-replies with $next$ (which now points to the $root$)
15:     wait for FIN msg from $root$ with $newID$ and $newLead$
16:     **if** FIN msg received **then**
17:         $newID(F_{lead}^{\text{mp}(i)}(w)) \leftarrow newID$
18:         $state(F_{lead}^{\text{mp}(i)}(w)) \leftarrow active$

19: **if** $state(i) = root$ **then**
20:     wait for incoming merge-requests
21:     store the sources of the requests
22:     reply on all requests with NULL
23:     **if** num of requests $= 0$ and size of merge-tree $\leq 2^{2+phase}$ **then**
24:         $newID \leftarrow$ random ID among all fragments in the merge-tree
25:         $newLead \leftarrow$ random node in $\mathcal{C}$
26:         send FIN msg with $newID$ and $newLead$ to all the stored sources
27:         $state(i) \leftarrow active$

---

---

**Algorithm** PJ($i, next, iter$) (pointer jumping)

---

Executed by each fragment $i$ in the $frozen$ state
**Input:** $next$ – first fragment to try, $iter$ – how many pointer-jumps to perform

**Output 1:** indication whether the $root$ was reached
**Output 2:** pointer to the $root$ or to the next fragment in the chain

  1: **while** $iter > 0$ **do**
  2:    **if** $i = \text{spk}^{\text{mp}(i)}(w)$ **then**
  3:      send merge-request to $next$

  4:    receive merge-requests
  5:    send merge-replies with $next$

  6:    **if** $i = \text{spk}^{\text{mp}(i)}(w)$ **then**
  7:      receive merge-reply with $next'$
  8:      **if** $next' = \text{NULL}$ **then**
  9:        **return** [TRUE, $next$]
10:      $next \leftarrow next'$
11:      $iter \leftarrow iter - 1$

12: **return** [FALSE, $next$]

---

**Algorithm** FINDROOT($i$)

---

  1: **if** $i = \text{spk}^{\text{mp}(i)}(w)$ **then**
  2:    send merge-request to $\text{mp}(i)$

  3: receive merge-requests
  4: send merge-replies with $\text{mp}(i)$

  5: **if** $i = \text{spk}^{\text{mp}(i)}(w)$ **then**
  6:    receive merge-reply with $next'$
  7:    **if** $next' \in F_{lead}^{\text{mp}(i)}(w)$ **then**
  8:      **if** $next' \leq \text{mp}(i)$ **then**
  9:        **return** [TRUE, $next'$]

10: **return** [FALSE, NULL]

---