

Massively Parallel Algorithms for High-Dimensional Euclidean Minimum Spanning Tree

Rajesh Jayaram
Google Research
rkjayaram@google.com

Vahab Mirrokni
Google Research
mirrokni@google.com

Shyam Narayanan
MIT*
shyamsn@mit.edu

Peilin Zhong
Google Research
peilinz@google.com

November 10, 2023

Abstract

We study the classic Euclidean Minimum Spanning Tree (MST) problem in the Massively Parallel Computation (MPC) model. Given a set $X \subset \mathbb{R}^d$ of n points, the goal is to produce a spanning tree for X with weight within a small factor of optimal. Euclidean MST is one of the most fundamental hierarchical geometric clustering algorithms, and with the proliferation of enormous high-dimensional data sets, such as massive transformer-based embeddings, there is now a critical demand for efficient distributed algorithms to cluster such data sets.

In low-dimensional space, where $d = O(1)$, Andoni, Nikolov, Onak, and Yaroslavtsev [STOC '14] gave a constant round MPC algorithm that obtains a high accuracy $(1 + \epsilon)$ -approximate solution. However, the situation is much more challenging for high-dimensional spaces: the best-known algorithm to obtain a constant approximation requires $O(\log n)$ rounds. Recently Chen, Jayaram, Levi, and Waingarten [STOC '22] gave a $\tilde{O}(\log n)$ approximation algorithm in a constant number of rounds based on embeddings into tree metrics. However, to date, no known algorithm achieves both a constant number of rounds and approximation.

In this paper, we make strong progress on this front by giving a constant factor approximation in $\tilde{O}(\log \log n)$ rounds of the MPC model. In contrast to tree-embedding-based approaches, which necessarily must pay $\Omega(\log n)$ -distortion, our algorithm is based on a new combination of graph-based distributed MST algorithms and geometric space partitions. Additionally, although the approximate MST we return can have a large depth, we show that it can be modified to obtain a $\tilde{O}(\log \log n)$ -round constant factor approximation to the Euclidean Traveling Salesman Problem (TSP) in the MPC model. Previously, only a $O(\log n)$ round was known for the problem.

*Work done as a student researcher at Google Research.

1 Introduction

The minimum spanning tree (MST) problem is one of the most fundamental problems in combinatorial optimization, whose algorithmic study dates back to the work of Boruvka in 1926 [Bor26]. Given a set of points and distances between the points, the goal is to compute a tree over the points of minimum total weight. The MST problem has received a tremendous amount of attention from the algorithm design community, leading to a large toolbox of methods for the problem [CCFC02, IT03, Ind04, FIS05, HPIM12, ANOY14, ACK⁺16, BBD⁺17a, CS09, CEF⁺05, CS04, Cha00, CRT05].

In the offline setting, where the input graph $G = (V, E)$ is known in advance, an exact randomized algorithm running in time $O(|V| + |E|)$ is known. The version of the problem where the input lies in Euclidean space has also been extensively studied, see [Epp00] for a survey. In this setting, the vertices of the graph are points in \mathbb{R}^d , and the set of weighted edges is (implicitly given by) the set of all $\binom{n}{2}$ pairs of vertices and the pairwise Euclidean distances. Despite the implicit representation, the best known runtime bound for computing an Euclidean MST exactly is n^2 , even for low-dimensional inputs. For approximations, the best known runtime bound is $O(n^{2-2/(\lceil d/2 \rceil + 1) + \varepsilon})$ for obtaining a $(1 + \varepsilon)$ -approximate solution [AESW90], and is n^{1+1/c^2} for an $O(c)$ -approximation [HIS13].

The Euclidean variant of the MST problem is particularly important in light of the tremendous success of modern *embedding* models in machine learning [MCCD13, VdMH08, HZRS16, DCLT18]. Such models encode a data point, such as an image, video, or text, into a vector in high-dimensional space, so that the semantic similarity between two data points is accurately represented by the Euclidean distance between the corresponding embedding vectors. In this setting, computing an MST is a well-known and successful technique for clustering data [LRN09, WWW09, ZMMF15, GZJ06, BBD⁺17b], which is a core component of many ML pipelines for large embedding datasets.

To deal with the sheer size of these modern embedding datasets, the typical approach is to implement algorithms in massively parallel computation systems such as MapReduce [DG04, DG08], Spark [ZCF⁺10], Hadoop [Whi12], Dryad [IBY⁺07] and others. The *Massively Parallel Computation* (MPC) model [KSV10, GSZ11, BKS17, ANOY14] is a computational model for these systems that balances accurate modeling with theoretical elegance. The MST problem in particular has been extensively studied in this model [ANOY14, KSV10, ASS⁺18, LMSV11, AAH⁺23, CJLW22]. For the non-Euclidean case, when the input is a graph with n vertices and m edges, and each machine only has n^ε space for some constant $\varepsilon < 1$, algorithms that obtain a constant approximation in $O(\log n)$ rounds and linear $O(m)$ total space are known using connected components algorithms [KSV10, ASS⁺18, BDE⁺19, CC22]. However, improving the round complexity is unlikely, as such an algorithm would refute the longstanding 1-CYCLE vs. 2-CYCLE conjecture [YV18, RVW18, LMW18, ASW19].

The Euclidean MST problem, on the other hand, is not as well understood in the MPC model as its graph-based counterpart. For low-dimensional Euclidean space, where d is a constant, [ANOY14] gave a $(1 + \varepsilon)$ approximate algorithm using only $O(1)$ MPC rounds. However, the landscape becomes much more mysterious and challenging in the high dimensional-setting.

One prevalent approach for the high-dimensional setting is the *spanner method*: one first constructs a c -approximate Euclidean spanner (i.e., a sparse graph over the points whose shortest path distance approximates the Euclidean distance to a factor of c) for some constant $c \geq 1$. Such spanners can be constructed with $O(n^{1+1/c})$ edges in $O(1)$ MPC rounds [EMMZ22, CAMZ22]. With the spanner

in hand, one can simply run the aforementioned graph-based MST algorithm to obtain a constant approximation in $O(\log n)$ rounds. However, in the Euclidean variant we have the benefit of the metric-space structure, and the one-cycle two-cycle lower bound does not apply. Thus, settling for $O(\log n)$ rounds for Euclidean MST is undesirable and perhaps unnecessary.

To date, the only known method for obtaining MPC algorithms with a $o(\log n)$ round complexity for high-dimensional MST is the *tree-embedding method* [Bar96]. In this method, one probabilistically embeds the Euclidean points $X \subset \mathbb{R}^d$ into a (log-depth) tree-metric and then computes the optimal MST in the tree metric. The latter is a simple object that can be computed in the MPC model in $O(1)$ rounds and $\tilde{O}(n)$ total space. Indyk [Ind04] gave a tree embedding that obtains a $O(d \log n)$ approximation and can be computed in the MPC model in $O(1)$ rounds and $\tilde{O}(n)$ total space. This was later improved by [CJLW22] to a $\tilde{O}(\log n)$ approximation, resulting in a $O(1)$ round $\tilde{O}(\log n)$ approximation algorithm. However, it is known that the tree-embedding method must suffer a $\Omega(\log n)$ approximation in the worst case [Bar96].

In summary, one can obtain constant approximations in $O(\log n)$ rounds via the spanner method, or a $\tilde{O}(\log n)$ approximation in constant rounds via the tree embedding method. However, both methods individually face hard barriers to removing the $O(\log n)$ factor from their round complexity or approximation (respectively). A natural question is whether this trade-off is intrinsic: namely if any algorithm for high-dimensional Euclidean MST must always use at least $\Omega(\log n)$ rounds or pay a $\Omega(\log n)$ approximation. Specifically, in this work we address the following question:

Is it possible to compute an $o(\log n)$ -approximate Euclidean minimum spanning tree in $o(\log n)$ rounds of the MPC model?

An even stronger question is whether a $\tilde{O}(1)$ approximation is possible in $O(1)$ rounds. As a positive signal in this direction, a recent result [CCAJ⁺23] demonstrated that the *cost* of the MST can be estimated to a constant factor in $O(1)$ MPC rounds. However, their algorithm is a sampling-based estimator that is far removed from a procedure that can actually compute an approximate MST. Moreover, separations between the complexity of estimating the cost of a solution and producing that solution are ubiquitous in high-dimensional geometry (e.g. for metric MST in the sublinear query model [Ind99, CS09]). For computing the MST in a distributed setting, such a result was only known in the more powerful Congested Clique [JN18] model, however implementing this algorithm in the MPC model would require $\Omega(n)$ space per machine.

In this work, we provide a positive resolution to the above question by designing a fully scalable¹ MPC algorithm in $\tilde{O}(\log \log n)$ rounds for a constant approximation of the Euclidean MST. In addition, the total space required by our algorithm is at most $O(n^{1+\varepsilon} + nd)$ where $\varepsilon > 0$ can be an arbitrarily small constant. Our result makes substantial progress towards the stronger goal of a $O(1)$ approximation in $O(1)$ rounds. Specifically, our main result is:

Theorem 1 (see Theorem 12). *Given a set $X \subset \mathbb{R}^d$ of n points, there is a MPC algorithm which outputs an $O(1)$ -approximate MST of X with probability at least 0.99 in at most $O(\log \log(n) \cdot \log \log \log(n))$ rounds. The total space required is at most $O(nd + n^{1+\varepsilon})$ and the per-machine space is $O((nd)^\varepsilon)$, where $\varepsilon > 0$ is an arbitrarily small constant.*

¹Meaning that each machine has only n^ε space for a constant $\varepsilon < 1$, see Preliminaries 1.1.

At a high-level, our algorithm bypasses the $\Omega(\log n)$ barriers intrinsic to the spanner and tree-embedding methods by combining the two approaches. Specifically, our algorithm builds a spanner and attempts to compute the MST of that spanner. However, while doing so we do *not* forget about the original metric structure. Specifically, our algorithm will exploit the metric structure by generating a spanning forest using *both* the edges of the spanner and geometric space partitions (which underpin tree-embedding methods).

Euclidian TSP. We leverage our algorithm for MST to develop the first $O(1)$ -approximation to the Euclidean traveling salesman problem (TSP) in $o(\log n)$ rounds. At a first glance, one may think that $O(1)$ -approximate TSP should directly follow from a $O(1)$ -approximate MST since any shortcutted traverse of the $O(1)$ -approximate MST gives an $O(1)$ -approximate TSP. However, the approximate MST that we computed may have diameter $\Theta(n)$, and all existing fully scalable MPC algorithms require at least $\Omega(\log(\text{diameter})) = \Omega(\log n)$ rounds to compute a traversal of the tree. To resolve this issue, we develop a new $O(\log \log(n))$ -round fully scalable MPC algorithm to compute a traverse of our approximate MST by utilizing the information of a hierarchical decomposition of the point set that we generate while computing the approximate MST. This results in the following:

Theorem 2 (see Corollary 7). *Given a set $X \subset \mathbb{R}^d$ of n points, there is a TSP algorithm which outputs an $O(1)$ -approximate TSP of X with probability at least 0.99 in at most $O(\log \log(n) \cdot \log \log \log(n))$ rounds. The total space required is at most $O(nd + n^{1+\varepsilon})$ and the per-machine space is $O((nd)^\varepsilon)$, where $\varepsilon > 0$ is an arbitrarily small constant.*

1.1 Preliminaries

MPC Model. In the Massively Parallel Computation (MPC) model, there are p machines and each machine has local memory s ; thus the total space available in the system is $p \cdot s$. The space is measured in words where each word has $O(\log(ps))$ bits. The input data has size N and is distributed arbitrarily on $O(N/s)$ machines at the beginning of the computation. If the total space satisfies $p \cdot s = O(N^{1+\gamma})$ for some $\gamma \geq 0$, and the local space satisfies $s = O(N^\varepsilon)$ for some constant $\varepsilon \in (0, 1)$, then the model is called the (γ, ε) -MPC model [ASS⁺18].

The computation in the MPC model proceeds in rounds. In every round, each machine performs arbitrary local computation on the data stored in its memory. At the end of each round, each machine sends some messages to the other machines. Since each machine only has local memory s , the total size of messages sent or received by a machine in one round can not be larger than s . For example, a machine can send a single message with size s to an arbitrary machine, or it can send a size 1 message to other s machines, but it cannot send a size s message to every machine in one round. In the next round, each machine only holds received messages in its local memory. At the end of the computation, the output is stored in a distributed way on the machines. The parallel running time (number of rounds) of an MPC algorithm is the number of above computation rounds needed to finish the computation.

We consider $\varepsilon \in (0, 1)$ to be an arbitrary constant in this paper, i.e., our algorithms can work when the memory of each machine is $s = O(N^\varepsilon)$ for any constant $\varepsilon \in (0, 1)$. Such algorithms are called *fully scalable algorithms* [ASS⁺18]. Our goal is to develop fully scalable algorithms which only require a small number of rounds and a small total space.

Basic Notation. In the remainder of the paper, we use X to denote a dataset of points that we

wish to solve either Minimum Spanning Tree or Traveling Salesman Problem over. We use n to denote the size of X and d to denote the dimensionality of X (i.e., $n = |X|$ and $X \subset \mathbb{R}^d$). Additional notation is defined in the relevant sections.

Euclidean Minimum Spanning Tree (MST) and Travelling Salesman Problem (TSP)

In the Euclidean MST problem, the input is a set of n points $X \subset \mathbb{R}^d$, where we assume that each coordinate x_i of a point $x \in \mathbb{R}^d$ can be stored in a single word of space (i.e., $O(\log ps)$ bits). The points X implicitly define a complete graph, where the vertices are X , and for any $x, y \in X$ the weight of the edge (x, y) is $\|x - y\|_2$. Our goal will be to produce a spanning tree T of this complete graph such that the weight of T , defined as $\sum_{(x_i, x_j) \in T} \|x_i - x_j\|_2$ is within a constant factor of the minimum spanning tree weight. We write $\mathbf{MST}(X)$ to denote the optimal MST weight.

In the Euclidean TSP problem, the input is the same as the MST problem. But instead of outputting a spanning tree, we want to output a Hamiltonian cycle C of X , i.e., each point appears on the cycle exactly once, such that the total length of the cycle $\sum_{(x_i, x_j) \in C} \|x_i - x_j\|_2$ is minimized up to a constant factor.

1.2 Technical Overview

1.2.1 Approximate Euclidean MST

Our starting point is a key structural fact about minimum spanning trees (observed in prior work on sublinear MST algorithms [CRT05, CS09]) that links the cost of the minimum spanning tree of a weighted graph G to the number of connected components in a sequence of auxiliary graphs. Namely, given a set of points $X \subset \mathbb{R}^d$ with pairwise distances in the range $(1, \Delta)$,² and a distance threshold $t \geq 0$, we define the t -threshold graph $G_t = (X, E_t)$ to be the graph where $(x, y) \in E_t$ if and only if $\|x - y\|_2 \leq t$. We write \mathcal{P}_t to denote the set of connected components of G_t .

Now consider the steps taken by Kruskal's MST algorithm: at the beginning, all vertices are in their own (singleton) connected component, and then at each step two connected components are merged by an edge of minimum weight. Thus, the number of edges added with weight in the range $(t, 2t]$ is precisely $|\mathcal{P}_t| - |\mathcal{P}_{2t}|$, so

$$\mathbf{MST}(X) \leq \sum_{i=0}^{\log(\Delta)-1} 2^{i+1} (|\mathcal{P}_{2^i}| - |\mathcal{P}_{2^{i+1}}|) = n - \Delta + \sum_{i=0}^{\log(\Delta)} 2^i |\mathcal{P}_{2^i}| \leq 2 \mathbf{MST}(X) \quad (1)$$

This suggests the following approach (which we call the *ideal algorithm*): for each level $t = 2^i$ where $i = 1, 2, \dots, \log(\Delta)$, compute the set \mathcal{P}_t of connected components of G_t . Then, for each t and every connected component $C \in \mathcal{P}_t$, we contract the vertices in C into a single super-node in the graph G_{2t} ; call the resulting contracted graph \bar{G}_{2t} . Next, we run a *unweighted* spanning forest algorithm on the \bar{G}_{2t} , and output every edge we found in this forest. Notice that this gives a valid spanning tree. Moreover, since each edge in G_{2t} has weight at most $2t$, by the above this spanning tree will be a 2-approximate MST. Because the sets \mathcal{P}_t are fixed (i.e., independent of the algorithm), the above procedure can be run in parallel for each value of t .

²Note that we will later be able to assume that $\Delta \leq \text{poly}(n)$.

The first challenge to this approach is that the graph G_t can be dense, namely, it may have $\Omega(n^2)$ edges. Since the total space available to our algorithm is only $O(n^{1+\varepsilon})$, we will need to compress G_t . This is precisely what is accomplished by the *spanner method*. Namely, for every t one can create a graph S_t that is an $O(1/\varepsilon)$ -approximate spanner of G_t and has at most $n^{1+\varepsilon}$ edges [HIS13]³; moreover, this construction is efficiently implementable in the MPC model (see, e.g., [EMMZ22, CAMZ22]). Such a graph S_t has the property that for every $x, y \in X$ with $\|x - y\|_2 \leq t$, there is a path of length at most 2 between x and y in S_t , and for every edge (w, z) in S_t we have $\|w - z\|_2 \leq O(t/\varepsilon)$. By adding the edges of S_t into $S_{t'}$ for each $t' \geq t$, we ensure that the edges of the graphs S_t are monotone increasing in t (since we only consider $O(\log \Delta)$ values of t , this increases the number of edges by at most a $O(\log \Delta)$ factor). Let \mathcal{P}'_t be the set of connected components in S_t . It follows from the spanner property that $|\mathcal{P}_{O(t/\varepsilon)}| \leq |\mathcal{P}'_t| \leq |\mathcal{P}_t|$. Thus, to obtain a $O(1/\varepsilon)$ approximation, it will suffice to swap out G_t with the spanner S_t in the above ideal algorithm.

The second (more serious) challenge is computing the connected components of S_t in the MPC model. Specifically, unless the one-cycle two-cycle conjecture is false, in general computing the connected components of a graph in the MPC model requires $\Omega(\log n)$ rounds. Thus, if we construct the spanner graphs S_t and then forget about the original metric space that S_t came from, then computing the connected components of S_t in the MPC model will require $\Theta(\log n)$ rounds.⁴ Instead, our goal will be to run a connectivity algorithm while crucially using the metric structure of the original points $X \subset \mathbb{R}^d$. In what follows, we describe our approach to doing this.

Approximately Computing Connected Components: Leader Compression with an Early Termination. Instead of running an MPC connectivity algorithm on S_t as a black box, we will need to open the actual inner workings of the algorithm to analyze its interplay with the underlying metric structure of the graph. At a high level, our approach will be to cut off the execution of this algorithm early and show that the intermediate solution (set of connected subgraphs) that one obtains from this partial execution is good enough to compute a constant approximate MST. To this end, we will now describe a connectivity algorithm known as *leader compression*.

The leader compression algorithm proceeds in rounds. In each round, every vertex $u \in S_t$ flips a coin; the vertices that flip heads are called “leaders” and the vertices that flip tails are called “followers”. Then, each follower vertex u merges into a uniformly random leader vertex $v \in S_t$ such that (u, v) is an edge. If no such edge to a leader vertex exists, v is untouched on that step. Each edge (u, v) that is merged in the process is contracted in the graph into a super-node, and then the process is repeated in the next round where each super-node flips a coin to be either a leader or a follower. Thus, at every time step, each super-node represents a connected sub-graph of S_t which we may subsequently grow on later steps. Ultimately, each connected component in S_t will be contracted into a single super-node. Since on each round, every vertex is merged into another vertex with probability at least $1/4$, the process will terminate after $O(\log n)$ rounds. However, since each round of leader compression takes $O(1)$ rounds in the MPC model, we cannot afford to run leader compression to completion. Instead, our approach will be to cut off the leader compression algorithm early and return the intermediate super-nodes obtained in the process.

³[HIS13] actually achieves $\sim 1/\sqrt{\varepsilon}$ approximation. But we only consider $\varepsilon = O(1)$, and optimizing such dependence is not the focus of this paper.

⁴Note that this yields the complexity of the “spanner method” described earlier in the introduction.

To analyze the early stopping of leader compression, we observe two useful properties of this algorithm: firstly, the set of edges that are merged form a spanning forest of S_t , so we can use these edges for our approximate MST.⁵ Secondly, after any $h \geq 1$ rounds of the process, for any connected component C in S_t with m vertices, we expect there to be at most $m/2^{\Omega(h)}$ super-nodes in C . Call a super-node *complete* if it is maximal, i.e. it contains its entire connected component. After h rounds of leader compression, for every connected component C we expect that either C originally had size $2^{\Omega(h)}$, or C is contained by a complete super-node.

Our main approach is then as follows: we set $h = O(\log \log n)$, and run h rounds of leader compression on S_t . We refer to a super-node remaining after h rounds as an approximate connected component, and write $\hat{\mathcal{P}}_t$ to denote the set of such approximate components at level t . We then attempt to run the *ideal algorithm* described earlier, but using the set $\hat{\mathcal{P}}_t$ instead of the true set of connected components \mathcal{P}'_t of S_t . Since we terminated leader compression early, there may be many more approximate components than true connected components. However, for every true connected component $C \in \mathcal{P}'_t$, if C was not complete then leader compression at least reduced the number of vertices in C by a factor of $2^{\Omega(h)}$. Our goal will be to use this fact to argue that for an *incomplete* component C , the actual MST cost of C is much larger than the cost we must pay for having under-merged C (i.e., splitting C into multiple approximate components). To make this argument, we will make use of ideas from the tree-embedding literature.

Using Tree-Embeddings to Handle Incomplete Components. Briefly, the idea of using tree-embeddings to generate approximate MSTs is as follows. First, for every $t = 1, 2, 4, \dots, \Delta$, one can impose a randomly shifted hypergrid over \mathbb{R}^d with side length t/\sqrt{d} . The random shift ensures that points x, y that are much closer than t/\sqrt{d} are unlikely to be split, and points x, y with $\|x - y\|_2 > t$ will deterministically be split. If n_t is the number of non-empty hyper-grid cells at level t (i.e., cells that contain at least one point in X), then results from the tree-embedding literature [Ind04, AIK08, CJLW22] imply that with good probability:⁶

$$\mathbf{MST}(X) \leq \sum_{t=2^i} t \cdot n_t \leq \text{polylog}(n) \cdot \mathbf{MST}(X) \quad (2)$$

Equation 2 is promising, as it relates the number of non-empty cells n_t to the MST cost — if we can show that leader compression returns significantly fewer approximate connected components than there were non-empty cells, then this would satisfy our earlier goal. We employ this result in the following way. First, we similarly impose a randomly shifted grid of size length t/\sqrt{d} , and perform an initial merging of all points inside of the same grid cell (i.e., points in the same grid cell are automatically merged together). Since the diameter of a grid cell is $t/\sqrt{d} \cdot \sqrt{d} = t$, these points will necessarily be in the same connected component in S_t .

After pre-merging points in the same grid cell, we next perform $h = O(\log \log n)$ rounds of leader compression, and again write $\hat{\mathcal{P}}_t$ to denote the resulting set of approximate connected components. By the above, each connected component $C \in \mathcal{P}'_t$ is either fully merged at this point, or we have

⁵If $A, B \subset S_t$ are two super-nodes that are merged together during leader-compression, then we can pick any arbitrary edge between A, B to be used for the spanning forest.

⁶Note that given a nested set of hyper-grid cells, one can easily compute an approximate spanning tree with cost at most $\sum_{t=2^i} t \cdot n_t$, see, e.g., [Ind04, CJLW22].

reduced the number of super-nodes in C by a $1/2^{\Omega(h)}$ factor. Since each grid cell was merged into a super-node before the start of leader-compression, it follows that $|\hat{\mathcal{P}}_t| \leq |\mathcal{P}'_t| + n_t/2^{\Omega(h)}$.

Now by Equation 1 and the fact that the components \mathcal{P}'_t well-approximate \mathcal{P}_t , we have $\mathbf{MST}(G) \approx \sum_{t=2^i} t \cdot |\mathcal{P}'_t|$. Thus, if we output $|\hat{\mathcal{P}}_t|$ instead of $|\mathcal{P}'_t|$ connected components, then at level t we would be paying an additional cost of at most $t \cdot n_t/2^{\Omega(h)}$ in our spanning tree (note though that it is not clear yet that we can actually achieve this cost algorithmically, since we still need to find the edges to merge the components in $\hat{\mathcal{P}}_t$). Then by Equation 2, the fact that there are n_t non-empty cells with side-length t/\sqrt{d} implies that the actual MST cost $\mathbf{MST}(G)$ is at least $n_t \cdot t/(\sqrt{d} \cdot \text{polylog}(n))$. Putting this together, the additional cost we pay is at most $\sqrt{d} \cdot \text{polylog}(n)/2^{\Omega(h)} \cdot \mathbf{MST}(G)$ at level t . By standard dimensionality reduction for ℓ_2 , we may assume $d = \Theta(\log n)$, so for $h = O(\log \log n)$, this additional cost is at most $\mathbf{MST}(G)/(\log n)^{\Theta(1)}$. Since we only incur this additive cost for $O(\log n)$ geometrically increasing values of t , the total additive error incurred from under-merging components is still a small $\mathbf{MST}(G)/(\log n)^{\Theta(1)}$.

Challenge: Maintaining Consistency of Approximate Connected Components. If we ran the above procedure and simply counted the number of approximate connected components, by the above discussion this would be sufficient to obtain a constant approximation of the *cost* of the minimum spanning tree. However, there are major issues in using this approach when trying to generate the actual tree. Specifically, to create a valid spanning tree when using the approximate components $\hat{\mathcal{P}}_t$ in the *ideal algorithm*, it is necessary that $\hat{\mathcal{P}}_t$ is a refinement of $\hat{\mathcal{P}}_{2t}$ for every t a power of 2; otherwise, we would be unable to generate the edges merging $\hat{\mathcal{P}}_t$ into $\hat{\mathcal{P}}_{2t}$ without creating cycles or leaving vertices disconnected.

One possibility is to use the components in $\hat{\mathcal{P}}_t$ as a starting point to generate the components in $\hat{\mathcal{P}}_{2t}$. However, this would require us to create $\hat{\mathcal{P}}_1, \hat{\mathcal{P}}_2, \hat{\mathcal{P}}_4, \dots$ sequentially, which would require $\log n$ rounds. Recall that the original *ideal algorithm* did not have this issue, as the sets \mathcal{P}_t of connected components did not depend on the prior execution of the algorithm. Instead, our approach is to split the set of levels $\{1, 2, 4, 8, \dots, \Delta\}$ into smaller chunks of $O(\log \log n)$ levels each, and show that the approximate components for each chunk can be computed in parallel. Specifically, we set $\alpha = (\log n)^{O(1)}$, and then define the “checkpoint” levels $\{1, \alpha, \alpha^2, \alpha^3, \dots, \Delta\}$. For every two checkpoint levels t, t' , we will compute $\hat{\mathcal{P}}_t$ and $\hat{\mathcal{P}}_{t'}$ independently and in parallel. The challenge is to do this while maintaining consistency between $\hat{\mathcal{P}}_t$ and $\hat{\mathcal{P}}_{\alpha t}$ for each checkpoint level t .

To ensure consistency between $\hat{\mathcal{P}}_t$ and $\hat{\mathcal{P}}_{\alpha t}$, we need to ensure that any pair (u, v) of vertices merged in our leader compression algorithm $\hat{\mathcal{P}}_t$ was also merged in $\hat{\mathcal{P}}_{\alpha t}$. To this end, first recall that our leader compression algorithm for S_t and $S_{\alpha t}$ started by merging all points in the same grid cell with side length t/\sqrt{d} (resp., $\alpha t/\sqrt{d}$). Thus, if we enforce that our leader compression algorithm at level t *only* merges together pairs $u, v \in X$ that are in the same grid cell of side-length $\alpha t/\sqrt{d}$, then consistency will follow automatically. To enforce this condition, we simply modify the graph S_t by removing any edge $(u, v) \in S_t$ which crosses the randomly shifted grid with side-length $\alpha t/\sqrt{d}$. Since $\alpha = (\log n)^{O(1)}$ is taken sufficiently large, because of the random shift any edge (u, v) in S_t is cut with probability at most $1/(\log n)^{O(1)}$, as we must have had $\|u - v\|_2 = O(t) = O(\alpha t/(\sqrt{d} \log^{O(1)} n))$. This allows us to show that the cost of omitting these cut edges in S_t is small, and enables us to compute $\hat{\mathcal{P}}_t$ independently of $\hat{\mathcal{P}}_{\alpha t}$. Once we have computed $\hat{\mathcal{P}}_t$ and $\hat{\mathcal{P}}_{\alpha t}$ (for every checkpoint level t), we now hope to generate the intermediate sets $\hat{\mathcal{P}}_{2t}, \hat{\mathcal{P}}_{4t}, \dots, \hat{\mathcal{P}}_{\alpha t/2}$ while maintaining this consistency

property. Since there are only $O(\log \alpha) = O(\log \log n)$ intermediate sets, we can afford to compute them sequentially in $O((\log \log n)^2)$ rounds.

However, a challenge arises when attempting to compute the intermediate sets $\hat{\mathcal{P}}_{2t}, \hat{\mathcal{P}}_{4t}, \dots, \hat{\mathcal{P}}_{\alpha t/2}$. Namely, these sets still must be consistent with $\hat{\mathcal{P}}_{\alpha t}$, so for any $2t \leq \tau \leq \alpha t/4$, we cannot merge two vertices $u, v \in S_\tau$ that are not in the same approximate connected component in $\hat{\mathcal{P}}_{\alpha t}$. Unlike the situation with S_t , we cannot afford to cut edges in S_τ that cross the hyper-grid with side-length $\alpha t/\sqrt{d}$, since τ may no longer be significantly smaller than $\alpha t/\sqrt{d}$, so the probability that an edge in S_τ is cut is no longer small. Instead, one would need to cut every edge (u, v) in S_τ such that u, v were *not* in the same approximate component in $\hat{\mathcal{P}}_{\alpha t}$. Now if $\hat{\mathcal{P}}_{\alpha t} = \mathcal{P}_{\alpha t}$ were the true components in S_t , no such edges would exist. However, this will not be the case since we terminate leader compression early when constructing $\hat{\mathcal{P}}_{\alpha t}$, and leader compression does *not* have the property that each edge is equally likely to be contracted on a given round.

In fact, removing such edges from S_τ which crosses the partition $\hat{\mathcal{P}}_{\alpha t}$ can significantly increase the number of connected components in S_τ beyond what we can afford. To see this, consider the following instance of “parallel” path graphs: set $k = (2^h \cdot \alpha)^{O(1)} = (\log n)^{O(1)}$ sufficiently large, and let $X_0 = \{e_1, 2e_1, 3e_1, \dots, ke_1\}$ be a path graph, where $e_i \in \mathbb{R}^d$ is the standard basis vector for $i \in [d]$, and set $X_j = \{x + \frac{\alpha}{\sqrt{2}}e_j \mid x \in X_0\}$, for each $j \in [k]$, and let $X = \cup_{j=0}^k X_j$. Now consider running leader compression at levels $\sqrt{\alpha}$ and α ; at level α , each vertex $x \in X_j$ is connected to at $\Omega(k)$ other vertices in parallel paths X_i for $i \neq j$, but is connected to at most $O(\alpha) \ll k/(\log n)^{O(1)}$ points in the same path X_j . It follows that on every step of leader compression, the probability that x is merged with a point in the same path is $1/(\log n)^{O(1)} \ll 1/2^h$, so after h rounds of leader compression at level α each point $x \in X_j$ will be in a super-node containing no other points from X_j with probability $1 - 1/(\log n)^{O(1)}$; call such a point x totally cut. Note that there will then be $\Omega(k^2)$ totally cut vertices. Now when running leader compression at level $\sqrt{\alpha}$, the graph $S_{\sqrt{\alpha}}$ only contains edges within the same path, so if we cut edges that were not merged at level α then every isolated vertex becomes a singleton in $S_{\sqrt{\alpha}}$, thus $|\hat{\mathcal{P}}_{\sqrt{\alpha}}| = \Omega(k^2)$, so the cost of the MST produced will be at least $\sqrt{\alpha} \cdot |\hat{\mathcal{P}}_{\sqrt{\alpha}}| = \Omega(\sqrt{\alpha}k^2) > \Omega(\sqrt{\alpha}) \cdot \mathbf{MST}(X)$, which is a bad approximation. Thus, it is not possible to obtain a constant factor approximation while using under-merged clusters.

The Solution to Inconsistency: Generate Over-Merged Clusters. Our solution to the above issue is to generate approximate connected components that are *over-merged* instead of under-merged. In other words, each approximate component in the partition $\hat{\mathcal{P}}_t$ that we output will contain a full connected component in the true partition \mathcal{P}_t , and possibly more vertices as well. Consider any checkpoint level t — when running leader compression on S_t , we have that guarantee that for every true connected component C , either we fully merge C (i.e. C is complete), or we split C into at most $m/2^{\Omega(h)}$ super-nodes, where m was the number of hyper-grid cells with side length t/\sqrt{d} that intersected C . Instead of simply outputting these under-merged super-nodes as-is, we first perform an over-merging step where we arbitrarily merge together *every* “incomplete” super-node within in the same hyper-grid cell with the (larger) side length $\alpha t/\sqrt{d}$. Specifically, we can choose an arbitrary representative super-node v that is incomplete within such a hyper-grid cell, and connect every other incomplete supernode u in the same cell to v via an arbitrary edge. Since u, v were in the same hyper-grid cell, this edge will have weight at most αt . Also, recall that we modified S_t to remove edges crossing this larger hyper-grid; thus it follows that every resulting

merged cluster fully contains at least one connected component in S_t (i.e., we only over-merge). We then let $\hat{\mathcal{P}}_t$ be this set of over-merged clusters.

To argue that this over-merging of incomplete super-nodes does not increase the cost significantly, we note that for a super-node to be incomplete, it must have intersected at least $2^{\Omega(h)}$ hyper-grid cells of length t/\sqrt{d} in expectation. By Equation 2, if there are ℓ incomplete super-nodes in a cell, this means that $\mathbf{MST}(X) \geq t2^h\ell/\text{poly}(\log n)$. On the other hand, we pay a cost of $t\alpha$ to connect each incomplete super-node to the representative, for a total cost of $t\alpha\ell$. But taking $h = O(\log \log n)$ sufficiently large, it follows that $t2^h\ell/\text{polylog}(n) \gg t\alpha\ell$, thus we can afford this arbitrary over-merging step at level t . This handles the consistency between distinct checkpoint levels $1, \alpha, \alpha^2, \dots, \Delta$.

We must now consider maintaining consistency for the intermediate levels between checkpoints. Specifically, we run this over-merging algorithm to generate $\hat{\mathcal{P}}_1, \hat{\mathcal{P}}_\alpha, \hat{\mathcal{P}}_{\alpha^2}, \dots$ in parallel. Then, for each α^k , we generate the intermediate levels $\hat{\mathcal{P}}_{2\alpha^k}, \hat{\mathcal{P}}_{4\alpha^k}, \dots, \hat{\mathcal{P}}_{\alpha^{k+1}/2}$ sequentially, where $\hat{\mathcal{P}}_{2^j\alpha^k}$ is generated using $h = O(\log \log n)$ rounds of leader compression on top of the previous set of merged components in $\hat{\mathcal{P}}_{2^{j-1}\alpha^k}$. Now by construction of the over-merging procedure that we used to generate $\hat{\mathcal{P}}_{\alpha^{k+1}}$, the only way for an edge (u, v) to be in a graph S_τ , for any $\alpha^k < \tau < \alpha^{k+1}$ but *not* have been merged in $\hat{\mathcal{P}}_{\alpha^{k+1}}$ is if (u, v) crossed the hyper-grid cell with side length α^{k+2}/\sqrt{d} . But now this is acceptable, because the probability that such a cut occurs is at most $\sqrt{d}/\alpha = 1/(\log n)^{O(1)}$ over the random shift, so we can now safely remove these edges from S_τ to maintain consistency.

Note that the above algorithm result in $O((\log \log n)^2)$ rounds of MPC. However, we can further improve this by a careful application of binary-search. Namely, if we start with $\hat{\mathcal{P}}_t$ and $\hat{\mathcal{P}}_{\alpha t}$, we first generate $\hat{\mathcal{P}}_{\alpha^{1/2}t}$, then in parallel generate $\hat{\mathcal{P}}_{\alpha^{1/4}t}, \hat{\mathcal{P}}_{\alpha^{3/4}t}$, and so on. This reduces the number of sequential rounds to $O(\log \log \alpha) = O(\log \log \log n)$, and since each round needs $O(\log \log n)$ rounds of leader compression, we use only $\tilde{O}(\log \log n)$ total rounds of MPC. However, at each intermediate step, we again need to over-merge for consistency reasons. This time, if we know $\hat{\mathcal{P}}_{t/\gamma}$ and $\hat{\mathcal{P}}_{t\gamma}$ and are trying to generate $\hat{\mathcal{P}}_t$, we merge incomplete components in the same connected component of $\hat{\mathcal{P}}_{t\gamma}$, so that we ensure consistency.

Generating the edges. Given the above algorithm that maintains consistency, generating edges is now quite simple. Assuming that we have found the approximate connected components $\hat{\mathcal{P}}_1, \hat{\mathcal{P}}_2, \hat{\mathcal{P}}_4, \dots$, we can generate edges merging $\hat{\mathcal{P}}_t$ into $\hat{\mathcal{P}}_{2t}$ for every $t = 1, 2, 4, \dots$, in parallel for each t . For a fixed t , this can be done by performing $O(\log \log n)$ rounds of leader compression on S_{2t} starting with $\hat{\mathcal{P}}_t$. But this time, rather than just updating the connected components, we keep track of the edges we found. While this will not generate all of the edges for this level, the same arguments as before will imply that leader compression finds $1 - \frac{1}{(\log n)^{\Theta(1)}}$ fraction of edges needed. Finally, to fully connect $\hat{\mathcal{P}}_{2t}$, we can add arbitrary edges, and by the same argument used to analyze the tree embedding, the additional cost will not be too large.

1.2.2 Approximate Euclidean TSP

Recall that Euclidean TSP aims to find a cycle over the points such that each point is visited exactly once, and we want to minimize the total weight of the cycle. A Euclidean MST is a 2-approximation

for the Euclidean TSP problem since a shortcut Euler tour of the MST gives a valid solution for TSP; here, an Euler tour of a tree is a directed circular tour on the tree such each undirectly tree edge (u, v) appears exactly twice: once in each direction (u, v) and (v, u) . Therefore, it suffices to compute an Euler tour of the approximate MST produced by our earlier algorithm. However, this approximate MST could be a path with diameter $\Theta(n)$. Unfortunately, the best-known algorithm for computing an Euler tour algorithm requires $\Omega(\log(\text{diameter})) = \Omega(\log(n))$ rounds [ASS⁺18]. Moreover, improvements on this are unlikely as an algorithm using $o(\log n)$ rounds algorithm would refute 1-CYCLE vs. 2-CYCLE conjecture [YV18, RVW18, LMW18, ASW19].

Fortunately, in addition to the edges of approximate MST, our algorithm also outputs an $O(\log n)$ -level hierarchical decomposition $\hat{\mathcal{P}}_1, \hat{\mathcal{P}}_2, \hat{\mathcal{P}}_4, \hat{\mathcal{P}}_8, \dots$, of the point set. We will show that, though our approximate MST may have large diameter on the original point set, if we look at $\hat{\mathcal{P}}_{t/2}$ for any fixed t , and regard each cluster in $\hat{\mathcal{P}}_{t/2}$ as a node, then each cluster in $\hat{\mathcal{P}}_t$ is merged from a subset of these nodes, and the tree introduced by the edges that we selected connecting these nodes (clusters in $\hat{\mathcal{P}}_{t/2}$) has diameter at most $\text{polylog}(n)$.

To see why this is true, consider how each cluster in $\hat{\mathcal{P}}_t$ was constructed when merging clusters from $\hat{\mathcal{P}}_{t/2}$. Suppose $C \in \hat{\mathcal{P}}_t$ is merged from $\mathcal{C} = \{C_1, C_2, \dots, C_k\} \subseteq \hat{\mathcal{P}}_{t/2}$. To generate the edges, we ran $h = O(\log \log n)$ rounds of leader compression, and in each round we merged some subsets of \mathcal{C} by creating edges between pairs $(C_i, C_j) \in \mathcal{C}^2$. Let $\hat{\mathcal{P}}_t^{(i)}$ be the partition we obtained after i rounds of leader compression. Then each cluster in $\hat{\mathcal{P}}_t^{(1)}$ corresponds to a star graph over a subset of \mathcal{C} (with the leader as the center), and similarly each cluster in $\hat{\mathcal{P}}_t^{(2)}$ is a star graph over a set of clusters in $\hat{\mathcal{P}}_t^{(1)}$, and so on. Since the clusters in $\hat{\mathcal{P}}_t^{(1)}$ correspond to a tree over \mathcal{C} with diameter 2 (i.e. a star), this means that each merged tree over C_1, C_2, \dots, C_k in $\hat{\mathcal{P}}_t^{(2)}$ has diameter at most $3 \cdot 2 + 2$. Then by induction, the full tree over the vertex set \mathcal{C} in $\hat{\mathcal{P}}_t^{(i)}$ has diameter $2 \cdot 3^{i-1} - 1$. Finally, after $O(\log \log n)$ rounds of leader compression, we create a star on the clusters in C_1, C_2, \dots, C_k that were still unmerged, which will blow up the diameter by another factor of at most 2. Thus, the tree $T = (\mathcal{C}, E(T))$ of resulting from combining all these edges has diameter at most $3^{O(\log \log n)} = \text{polylog}(n)$. We call this tree the *super-node tree* over \mathcal{C} .

We will use this fact to obtain a Euler tour of this tree in $O(\log \log n)$ MPC rounds (e.g., by applying the Euler tour algorithm of [ASS⁺18], and using that the diameter is small). To this end, we now introduce a critical sub-problem, which we call the Euler Tour Join problem. Given any t and $i \geq 0$, the input to the problem is the following:

1. A cluster $C \in \hat{\mathcal{P}}_t$ that was merged from $\mathcal{C} = \{C_1, C_2, \dots, C_k\} \in \hat{\mathcal{P}}_{t/2^{2^i}}$, and an Euler tour A of the super-node tree T over \mathcal{C} .
2. For each C_i , suppose C_i was merged from clusters $C_{i,1}, C_{i,2}, \dots, C_{i,k_i} \in \hat{\mathcal{P}}_{t/2^{2^{i+1}}}$. Then we are also given the super-node tree T_i over $\{C_{i,1}, C_{i,2}, \dots, C_{i,k_i}\}$, and an Euler tour A_i of T_i .

Given the above trees T, T_1, T_2, \dots, T_k and tours A, A_1, A_2, \dots, A_k , the goal of the Euler Tour Join problem is to compute (1) the super-node tree T' over $\cup_{i=1}^k \{C_{i,1}, C_{i,2}, \dots, C_{i,k_i}\}$ which represents how C was merged from the clusters in $\hat{\mathcal{P}}_{t/2^{2^{i+1}}}$ and (2) an Euler tour A' of T' . Notice that if we can solve this problem for any $i \geq 0$, then we can run it for $i = 0, 1, 2, \dots, \log \log(\Delta)$

sequentially, by pairing up groups of levels in $\{1, 2, 4, 8, \dots, \Delta\}$, merging them, and the recursing on the merged groups. Thus, after $O(\log \log n)$ rounds, we would have a Euler tour of the full original dataset X . Thus, our main task now is to develop an $O(1)$ round fully-scalable MPC algorithm for the Euler Tour Join problem using total space $O(n)$.

We now describe the high level ideas behind our Euler Tour Join algorithm. Let $C, C_1, \dots, C_k, T, T_1, \dots, T_k$, and A, A_1, \dots, A_k be as above, and let $V = \cup_{i=1}^k \{C_{i,1}, C_{i,2}, \dots, C_{i,k_i}\}$. We will sometimes think of each C_i as being composed of the the component sets $\{C_{i,1}, C_{i,2}, \dots, C_{i,k_i}\}$, and write $x \in C_i$ to denote that $x \in \{C_{i,1}, C_{i,2}, \dots, C_{i,k_i}\}$. Now observe that for every edge $(C_i, C_j) \in E(T)$, there must have been a unique edge $(C_{i,a}, C_{j,b})$, where $C_{i,a} \in C_i, C_{j,b} \in C_j$, used by the leader compression algorithm to connect C_i, C_j , between two clusters in $\hat{P}_{t/2^{i+1}}$. Thus, we can define a function $g : E(T) \rightarrow V^2$ that specifies this mapping (i.e., in the above example $g((C_i, C_j)) = (C_{i,a}, C_{j,b})$). It is easy to see that the tree $T' = (V, E(T'))$ with edges $\{g(C_i, C_j) \mid (C_i, C_j) \in E(T)\} \cup T_1 \cup T_2 \cup \dots \cup T_k$ is a spanning tree of V . Thus, the challenge will be to compute an Euler tour of T' . In what follows, any $x, y \in V$ such that $(x, y) \in E(T')$ and $x \in C_i, y \in C_j$ and $i \neq j$, we call both x and y *terminal nodes*. In other words, terminal nodes are clusters from V that connected two clusters from $\{C_1, \dots, C_k\}$.

Given the above, one natural idea to construct A' is as follows. We follow along the Euler tour of A ; each edge (C_i, C_j) in A corresponds to an edge $(u, v) = g(C_i, C_j) \in V^2$ between two terminals $u \in C_i, v \in C_j$ which means that our tour enters cluster C_j via v after leaving C_i . Suppose the edge that follows (C_i, C_j) in A is (C_j, C_l) and the corresponding terminals are (x, y) , then it means that our tour leaves C_j via x before entering C_l . Then, it would be natural to simply plug the subtour A_j that connects v to x in our final joined tour A' to connect (u, v) and (x, y) . However, this approach fails, as it may use duplicated edges in A_j . In fact, there exists a Euler tour A such that if we visit terminals with respect to the ordering provided by A , we can never find a valid Euler tour A' . As an example, suppose we have 5 clusters

$$C_1 = \{i_1\}, \quad C_2 = \{i_2\}, \quad C_3 = \{i_3\}, \quad C_4 = \{i_4\}, \quad C_5 = \{i_{5,1}, i_{5,2}, i_{5,3}, i_{5,4}, i_{5,5}\}$$

where the inter cluster edges are $\{i_1, i_{5,1}\}, \{i_2, i_{5,2}\}, \{i_3, i_{5,3}\}, \{i_4, i_{5,4}\}$, and T_5 has edges $\{i_{5,5}, i_{5,2}\}, \{i_{5,5}, i_{5,3}\}, \{i_{5,5}, i_{5,4}\}, \{i_{5,2}, i_{5,1}\}$ with Euler tour $A_5 : i_{5,1} \rightarrow i_{5,2} \rightarrow i_{5,5} \rightarrow i_{5,3} \rightarrow i_{5,5} \rightarrow i_{5,4} \rightarrow i_{5,5} \rightarrow i_{5,2} \rightarrow i_{5,1}$. Let T be a star with Euler tour $A : C_5 \rightarrow C_1 \rightarrow C_5 \rightarrow C_4 \rightarrow C_5 \rightarrow C_2 \rightarrow C_5 \rightarrow C_3 \rightarrow C_5$. It is obvious that if one wants to follow the subpath $C_1 \rightarrow C_5 \rightarrow C_4$, the inner path of T_5 must include $i_{5,1} \rightarrow i_{5,2} \rightarrow i_{5,5} \rightarrow i_{5,4}$ as a subsequence. Similarly, if one wants to follow the subpath $C_2 \rightarrow C_5 \rightarrow C_3$, the inner path of T_5 must include $i_{5,2} \rightarrow i_{5,5} \rightarrow i_{5,3}$ as a subsequence. Thus, we must use the (directed) $i_{5,2} \rightarrow i_{5,5}$ twice.

The main issue of the above approach is that the order of visiting the neighbors (C_1, C_4, C_2, C_3) of C_5 in the Euler tour A is not consistent with the order of visiting the terminals $(i_{5,1}, i_{5,2}, i_{5,3}, i_{5,4})$ of C_5 in A_5 . Namely, A_5 suggests the order (C_1, C_2, C_3, C_4) . Therefore, we need to compute a new Euler tour \bar{A} such that the ordering of visiting the neighbors of each C_i in T is consistent with ordering of visiting terminals in A_i .

To this end, we develop a novel algorithm in the MPC model such that if each edge in T is given a weight, and we are able to compute the total weight of every path from each node to the root of T , and we are also able to compute the size of each subtree of T , then we can efficiently compute the position of each directed edge of T in the desired Euler tour \bar{A} in parallel. Fortunately, we show

that above subtree sum problem and path weight sum problem can be solved efficiently using the known Euler tour A . Therefore, can compute \bar{A} in $O(1)$ rounds and $O(n)$ total space in the fully scalable setting. Then, the remaining process of computing A' becomes simple, we chop each Euler tour A_i into paths between terminals. Then we follow the ordering of edges in \bar{A} . For a length 2 path $C_i \rightarrow C_j \rightarrow C_l$ in \bar{A} , the relevant terminals are $(u, v) = g(C_i, C_j), (x, y) = g(C_j, C_l)$, then we insert the corresponding subsequence between u and x of Euler tour A_j into the place between (u, v) and (x, y) in our final Euler tour A' . Note that this sequence insertion subroutine can be efficiently implemented in the MPC model as shown by [ASS⁺18].

1.3 Other Related Work

The MST problem has been studied extensively in multiple models of sublinear computation, including streaming, distributed algorithms, and the sublinear query model. In the sublinear query model, the implicit input is the set of $\binom{n}{2}$ distances and the goal is to estimate the weight of the MST while making a sublinear number of queries to the distances. It is known that any algorithm which actually computes an approximate MST requires $\Omega(n^2)$ queries [Ind99], hence the focus on estimating the cost. To this end, Chazelle, Rubinfeld, and Trevisan [CRT05] gave an algorithm based on estimating connected components, that gives a $(1 + \epsilon)$ -factor approximation for arbitrary graphs of maximum degree D and edge weights in $[1, W]$ using at most $O(DW\epsilon^{-3})$ queries. This result was improved by Czumaj and Sohler [CS09] for metric (e.g., Euclidean) MST, who gave a $(1 + \epsilon)$ approximation with $\tilde{O}(n\epsilon^{-8})$ queries.

In the streaming model, the points $X \subset \mathbb{R}^d$ arrive in a stream, possibly with deletions, and the goal is to estimate the cost of the MST in small space (as outputting the MST would require $\Omega(n)$ space). The first algorithm for streaming Euclidean MST was due to Indyk [Ind04], who gave a $O(\log^2 n)$ approximation in $\text{polylog}(n)$ space. This was later improved by [CJLW22] to a $\tilde{O}(\log n)$ approximation in $\text{polylog}(n)$ space, and then a $\tilde{O}(1/\epsilon^2)$ approximation in $O(n^\epsilon)$ space by [CCAJ⁺23]. The first two works employed tree embedding based approaches, whereas [CCAJ⁺23] used a connected component-based estimator similar to those used in [CRT05, CS09]. For low-dimensional space, [FIS05] gives a $(1 + \epsilon)$ approximation with space exponential in the dimension, although the techniques in this paper rely heavily on the construction of exponentially sized ϵ -nets.

The above streaming algorithms are linear sketches and therefore can be used in the MPC model to obtain constant round approximations of the cost of the MST (see e.g. [ANOY14] for a reduction from linear sketching to the MPC). However, there is a substantial gap between estimating the cost of the MST and producing the MST. For instance, the estimators in the papers [CRT05, CS09, CCAJ⁺23] are based on sampling vertices and computing the size of their connected components, thus no edges or approximate tree structure can be derived from this approach. As described in Section 1.2, substantial challenges arise when attempting to construct a tree in a consistent (i.e., no cycles) and cost-effective way.

2 Description of the MST Algorithm in the Offline Setting

In this section, we give a description for generating an $O(1)$ -approximate MST, that can be implemented in the MPC model. However, we defer the details of the MPC implementation to Section 4.2.

We first assume WLOG that $d = \Theta(\log n)$, using the Johnson-Lindenstrauss lemma [JL84]. We will define two parameters α, β such that both α/β and β are at least $(\log n)^C$ for some sufficiently large constant C . We also assume WLOG $\alpha = 2^{2^g}$ for some integer g (which can always be done by replacing α with some $\alpha' \in [\alpha, \alpha^2]$). Define G to be the complete weighted graph on X , where $(x, y) \in X$ have an edge of weight $d(x, y)$. Define the *threshold graph* G_t to connect two points x, y if $d(x, y) \leq t$. We also define \mathcal{P}_t to be the partition of X based on the connected components of G_t .

Finally, since we are only hoping for a constant-factor approximation, by Standard discretization methods (see, e.g., Proposition 1 in [CCAJ⁺23]), we may assume the aspect ratio is at most $\tilde{O}(n)$. More precisely, we may assume that the minimum distance between any two points is at least $\alpha^{100} = (\log n)^{O(1)}$, and the maximum distance between any two points is at most $O(n \cdot \alpha^{101}) \leq n^2$ (for n sufficiently large). We shift the point set such that each coordinate of each point has value in $[0, \Delta]$ where Δ is a sufficiently large power of 2 and $\Delta = O(n^2)$.

We also note a few definitions.

Definition 1. Given a dataset X , a *partition* of X is a split of X into one or more pairwise disjoint subsets of X , such that every element $x \in X$ is in exactly one of the subsets.

Given two partitions \mathcal{P}, \mathcal{Q} of X , we say that \mathcal{P} *refines* \mathcal{Q} (or equivalently, \mathcal{Q} is *refined by* \mathcal{P}) if every partition component in \mathcal{P} is a subset of some partition component in \mathcal{Q} . We use the notation $\mathcal{P} \supseteq \mathcal{Q}$ (or $\mathcal{Q} \sqsubseteq \mathcal{P}$).

Definition 2. For a positive integer n , we define the *2-adic valuation* $v_2(n)$ to be the largest nonnegative integer k such that $2^k | n$.

High level approach: First, we will approximately generate \mathcal{P}_t , the connected components of G_t , for $t = 1, \alpha, \alpha^2, \dots, \alpha^H$, for $H = \Theta(\log n / \log \log n)$. We will approximately compute each of these partitions \mathcal{P}_t in parallel, using $h = O(\log \log n)$ rounds. Next, given the approximate connected components for G_{α^k} and $G_{\alpha^{k+1}}$, we attempt to generate approximate connected components for $G_{2\alpha^k}, G_{4\alpha^k}, \dots, G_{\alpha^{k+1}/2}, G_{\alpha^{k+1}}$. Finally, we generate edges to form an approximately minimal spanning tree.

Quadtree and Spanner: We start off knowing a randomly shifted *quadtree* Q . We recall the definition of a quadtree, along with some relevant notation.

Definition 3. A *randomly shifted Quadtree* is constructed as follows. First, we choose a random vector $a = (a_1, a_2, \dots, a_d) \in \mathbb{R}^d$, where each coordinate a_i is drawn uniformly from $[0, \Delta]$. For each t that is a power of 2 between 1 and 2Δ (where $\Delta = \Theta(n^2)$ is a power of 2), we use a to generate a grid of side length t , which we call the *grid at level* t . Specifically, there is a one-to-one correspondence between each grid cell in level t and each vector $(c_1, c_2, \dots, c_d) \in \mathbb{Z}^d$, i.e., the corresponding cell denotes the set of points: $\{(x_1, x_2, \dots, x_d) \in \mathbb{R}^d \mid \forall i \in [d], x_i \in [c_i \cdot t - a_i, (c_i + 1) \cdot t - a_i]\} \subset \mathbb{R}^d$. Finally, for any cell c at level t , define X_c to be the set of points in X contained in the cell c .

Note that the grids are nested, i.e., for every t , every cell of side length t in the Quadtree is contained in a cell of side length $2t$. Finally, for the largest grid length 2Δ , it is easy to see that there is a unique cell c in level Δ such that $X_c = X$. This is because every $a_i \in [0, \Delta]$ and because we assume all data points have coordinates in $[0, \Delta]$.

Next, we will also assume we have a series of 2-hop Euclidean spanners. We recall the definition of a 2-hop Euclidean spanner.

Definition 4. Given a dataset $Y \in \mathbb{R}^d$, a C -approximate 2-hop Euclidean spanner of length t is a graph on Y with edge set E , with the following two properties.

1. For any $p, q \in X$ such that $\|p - q\|_2 \leq t$, either $(p, q) \in E$ or there exists $r \in X$ such that $(p, r), (r, q) \in E$.
2. For any two points $p, q \in Y$ with $\|p - q\|_2 > t$, we must have $(p, q) \notin E$.

In other words, any two points within distance at most t are connected by a path of length at most 2, but any two points of distance more than t cannot be directly connected in the graph.

We will use the fact that for any dataset Y , any level t , and any constant $0 < \varepsilon < 1$, there exists a $O(1/\varepsilon)$ -approximate 2-hop Euclidean spanner of Y with at most $O(|Y|^{1+\varepsilon})$ edges. More specifically, for each t a power of 2, let k be such that $\alpha^{k-1} < t \leq \alpha^k$. We generate $\tilde{G}_t(X_c)$ to be a 2-hop Euclidean spanner of length t generated on X_c , where c is a cell at level α^{k+1}/β in the quadtree, where $\tilde{G}_t(X_c)$ has at most $O(|X_c|^{1+\varepsilon})$ edges. In addition, let \tilde{G}_t be the union of $\tilde{G}_t(X_c)$ across all cells c at level α^{k+1}/β . We also make sure that $\tilde{G}_t \subseteq \tilde{G}_{2t}$ for all t , simply by adding \tilde{G}_t to $\tilde{G}_{2t}, \tilde{G}_{4t}, \dots$: it is clear that this does not violate the definition of an $O(1/\varepsilon)$ -approximate spanner. Since there are at most $O(\log n)$ such levels, and since $\sum |X_c|^{1+\varepsilon} \leq (\sum |X_c|)^{1+\varepsilon} = n^{1+\varepsilon}$, every \tilde{G}_t has at most $O(n^{1+\varepsilon})$ edges. Finally, we let $\tilde{\mathcal{P}}_t$ be the partitioning of X based on the connected components of \tilde{G}_t . Note that $\tilde{\mathcal{P}}_t \supseteq \tilde{\mathcal{P}}_{2t}$ for all t , since $\tilde{G}_t \subseteq \tilde{G}_{2t}$.

We note that the assumptions on dimensionality reduction, aspect ratio, and our assumption that we have a Quadtree Q and 2-hop $O(1/\varepsilon)$ -approximate Euclidean spanners \tilde{G}_t can all be achieved using $O(1)$ rounds of MPC. We discuss this in Section 4.

Leader Compression Algorithm We need to first explain the leader compression algorithm. We assume that we start out with some starting partition \mathcal{P} of a dataset X , and are given some graph H on X . The goal is to connect components that have edges between them in H , to form larger connected components.

Formally, we will define a single round of leader compression on \mathcal{P} with respect to H as follows. For each connected component $C \in \mathcal{P}$, we assign every $x \in C$ some leader, which is a vertex $x^* \in C$. A round of leader compression works as follows. First, each leader will generate a random bit (uniformly 0 or 1), and will broadcast the value to all of its descendants (i.e., the rest of the vertices in C). So, every x_i now has a value of 0 or 1, matching that of its leader. Next, every x_i with a value of 1 will send a message to all of its neighbors in the graph H . Importantly, after these messages, every x_j that has value 0 knows the set of x_i that have value 1 (which means they do not have the same leader, since x_j has value 0) and are connected to x_j in H . For the purpose of MPC implementation, it will suffice for each such x_j to know a single such x_i (if one exists). Simultaneously, every x_j that has value 0 will choose such an x_i (assuming such an x_i exists), and x_j then sends the x_i value to its current leader x_k , in $O(1)$ rounds of MPC. Each leader x_k with value 0 will choose a single descendant x_j with such a x_i (if such an x_j exists: note j could equal k). Then, x_k will update its leader (as well as the leader of all of x_k 's descendants) to be the leader of x_i .

Finally, we can define h rounds of leader compression on \mathcal{P} with respect to H as follows. For the first round, we run a round of leader compression on \mathcal{P} with respect to H . We now have a potentially coarser partition of connected components, and for the second round we run leader compression on this new partition with respect to H . We repeat this for h rounds, and at each round we use the updated partition and assignments of leaders.

We include pseudocode for a round of leader compression in Algorithm 1.

Part 1: Approximately generating \mathcal{P}_t , $t = \alpha^k$. We start with a randomly shifted grid from Q of level t/β . We will automatically connect all points that are in the same cell. (Note that all such points must have distance at most $t\sqrt{d}/\beta$ from each other, so they should be in the same connected component at this level anyway). This forms an initial partitioning $\mathcal{P}_t^{(0)}$ of X .

Now, each (nonempty) component $\mathcal{P}_t^{(0)}$ starts with a “leader” point. We now perform $h = O(\log \log n)$ rounds of leader compression on $\mathcal{P}_t^{(0)}$ with respect to \tilde{G}_t . At the end of these rounds, we will have connected components and thus have $\bar{\mathcal{P}}_t$, our preliminary estimate for \mathcal{P}_t .

We will next convert $\bar{\mathcal{P}}_t$ into $\hat{\mathcal{P}}_t$ by doing some additional merging. For each component $S \in \bar{\mathcal{P}}_t$, we will check whether S is a complete component. To do so, we look at all of the edges in \tilde{G}_t : for each edge we check whether its endpoints are in the same connected component in $\bar{\mathcal{P}}_t$. If not, we send a message to both endpoints. Every vertex will then know if there is an edge leaving it into another connected component. Thus, every leader of a connected component will know if the connected component has any edges leaving it. We call such a component *incomplete*. To form $\hat{\mathcal{P}}_t$, we group together all incomplete components in $\bar{\mathcal{P}}_t$ with their leaders in the same cell c at level α^{k+1}/β .

We include pseudocode for this step in Algorithm 2.

Part 2: Approximately generating \mathcal{P}_t , $\alpha^k < t < \alpha^{k+1}$, t a power of 2. Recall that $\alpha = 2^{2^g}$. We will generate $\hat{\mathcal{P}}_t$ in decreasing order of $v_2(\log_2 t)$ (recall Definition 2 for $v_2(\cdot)$). In other words, we first generate $\hat{\mathcal{P}}_t$ for $t = \alpha^{k+1/2} = \alpha^k \cdot 2^{2^{g-1}}$, then for $t = \alpha^{k+1/4} = \alpha^k \cdot 2^{2^{g-2}}$ and $t = \alpha^{k+3/4} = \alpha^k \cdot 2^{3 \cdot 2^{g-2}}$ simultaneously, and so on. This will result in $g = \log_2 \log_2 \alpha$ iterations. When generating some $\alpha^k < t < \alpha^{k+1}$, we define $\kappa = 2^{2^{v_2(\log_2 t)}}$. We note that $v_2(\log_2(t/\kappa)), v_2(\log_2(t \cdot \kappa)) > v_2(\log_2 t)$, so we have already generated $\hat{\mathcal{P}}_{t/\kappa}$ and $\hat{\mathcal{P}}_{t \cdot \kappa}$. We will use these two partitions to generate $\hat{\mathcal{P}}_t$.

To do so, we start with the connected components from $\hat{\mathcal{P}}_{t/\kappa}$. Then, we perform $h = O(\log \log n)$ rounds of leader compression on $\hat{\mathcal{P}}_{t/\kappa}$ with respect to \tilde{G}_t , to create $\bar{\mathcal{P}}_t$. Finally, we convert $\bar{\mathcal{P}}_t$ to $\hat{\mathcal{P}}_t$, by performing merging in a similar way as in Part 1. Specifically, we check each edge in \tilde{G}_t and see if the endpoints are in the same connected component in $\bar{\mathcal{P}}_t$, and use this information to determine whether each component in $\bar{\mathcal{P}}_t$ is complete or incomplete. (Namely, a connected component C in $\bar{\mathcal{P}}_t$ is *incomplete* iff there exists an edge \tilde{G}_t with exactly one vertex in C .) Finally, we group together all incomplete connected components with their leaders in the same connected component in $\hat{\mathcal{P}}_{t \cdot \kappa}$.

We include pseudocode for this step in Algorithm 3.

Part 3: Generating the edges. We have listed out the approximate connected components $\hat{\mathcal{P}}_t$ for each t a power of 2. Now, for any such t , we want to generate edges connecting $\hat{\mathcal{P}}_{t/2}$ into $\hat{\mathcal{P}}_t$.

Algorithm 1 LEADERCOMPRESSION($\mathcal{P}, \ell, H, \text{Edges}$): A single round of leader compression on \mathcal{P} with respect to H , where ℓ is the function mapping each node to its leader in the partition. Edges is a boolean value that denotes whether we want to generate edges (which is only needed for Step 3).

```

1:  $S = \{\ell(x) : x \in X\}$  { $S$  is the set of leader nodes}
2: for  $x \in S$  do
3:    $b_x \leftarrow \text{Unif}(\{0, 1\})$ . { $b_x$  is the random bit for leader node  $x$ }
4: for  $x \in X$  do
5:    $b_x \leftarrow b_{\ell(x)}$ .
6: for  $x \in X$  with  $b_x = 1$  do
7:   for  $e = (x, y) \in H$  do
8:      $x$  sends message “ $(x, \ell(x))$ ” to  $y$ 
9:   for  $y \in X$  with  $b_y = 0$  do
10:    Select any message  $(x, \ell(x))$  sent to  $y$  (if it exists)
11:     $y$  sends message “ $(x, y, \ell(x))$ ” to  $\ell(y)$ .
12: for  $z \in S$  with  $b_z = 0$  do
13:   Select any message  $(x, y, \ell(x))$  sent to  $z$  (if it exists)
14:   if Edges then
15:      $E \leftarrow E \cup \{(x, y)\}$ . { $E$  is a set of edges, initialized to  $\emptyset$ }
16:    $z$  broadcasts message “ $\ell(x)$ ” to all descendants, all descendants (including  $z$ ) update their
     leader to be  $\ell(x)$ .
17: Update partition  $\mathcal{P}$  based on leader function.
18: Return  $(\mathcal{P}, \ell, H, E)$ . {If Edges is False, we return False instead of  $E$ .}

```

We will again perform $h = O(\log \log n)$ rounds of leader compression on $\hat{\mathcal{P}}_{t/2}$ with respect to \tilde{G}_t , but this time we keep track of all edges that we added. We will show that the partition we generate after doing leader compression still refines $\hat{\mathcal{P}}_t$. So, to finish generating $\hat{\mathcal{P}}_t$ along with the necessary edges, we connect all disconnected components in each partition of $\hat{\mathcal{P}}_t$, using arbitrary edges.

We include pseudocode for this step in Algorithm 4.

3 Analysis of the Offline Algorithm

3.1 Comparison to the Spanner Graph

Given the spanners \tilde{G}_t for each t a power of 2 (as described in Section 2), we create the graph \tilde{G} as follows. For each pair of distinct vertices (y, z) , we assign (y, z) the weight t for t the smallest power of 2 such that $(y, z) \in \tilde{G}_t$ (if such a t exists). Otherwise, we do not connect (y, z) .

We first note that any spanning tree in \tilde{G} does not increase in cost by too much after converting to the corresponding Euclidean spanning tree.

Proposition 1. *Suppose that T is a spanning tree in \tilde{G} . Then, the cost of T in X (i.e., over the true Euclidean distance) is at most $O(1/\varepsilon)$ times the cost of T in \tilde{G} .*

Algorithm 2 PART1(X, t, Q, \tilde{G}_t): Generating $\hat{\mathcal{P}}_t$ for $t = \alpha^k$.

- 1: Generate partition $\mathcal{P}_t^{(0)}$ of X based on the quadtree at level t/β , i.e., two points in X are in the same connected component in $\mathcal{P}_t^{(0)}$ if and only if they are in the same cell at level t/β .
 - 2: Choose an arbitrary leader for each (nonempty) component in $\mathcal{P}_t^{(0)}$, and let the corresponding leader mapping with respect to $\mathcal{P}_t^{(0)}$ to be $\ell_t^{(0)}$.
 - 3: **for** $i = 1$ to $h = O(\log \log n)$ **do**
 - 4: $(\mathcal{P}_t^{(i)}, \ell_t^{(i)}, \tilde{G}_t, \text{False}) \leftarrow \text{LEADERCOMPRESSION}(\mathcal{P}_t^{(i-1)}, \ell_t^{(i-1)}, \tilde{G}_t, \text{False})$.
 - 5: **for** $e = (x, y) \in \tilde{G}_t$ **do**
 - 6: **if** $\ell_t^{(h)}(x) \neq \ell_t^{(h)}(y)$ **then**
 - 7: (x, y) sends INCOMPLETE to both x and y
 - 8: **for** $x \in X$ that received INCOMPLETE message **do**
 - 9: x sends INCOMPLETE to $\ell_t^{(h)}(x)$.
 - 10: Merge all Incomplete components with leaders in the same cell in Q at level α^{k+1}/β , to form $\hat{\mathcal{P}}_t$.
-

Algorithm 3 PART2($X, t, \tilde{G}_t, \hat{\mathcal{P}}_{t/\kappa}, \hat{\mathcal{P}}_{t\cdot\kappa}$): Generating $\hat{\mathcal{P}}_t$, given $\hat{\mathcal{P}}_{t/\kappa}$ and $\hat{\mathcal{P}}_{t\cdot\kappa}$.

- 1: Initialize $\mathcal{P}_t^{(0)} \leftarrow \hat{\mathcal{P}}_{t/\kappa}$, and choose an arbitrary leader for each component in $\mathcal{P}_t^{(0)}$. Let the corresponding leader mapping with respect to $\mathcal{P}_t^{(0)}$ to be $\ell_t^{(0)}$
 - 2: **for** $i = 1$ to $h = O(\log \log n)$ **do**
 - 3: $(\mathcal{P}_t^{(i)}, \ell_t^{(i)}, \tilde{G}_t, \text{False}) \leftarrow \text{LEADERCOMPRESSION}(\mathcal{P}_t^{(i-1)}, \ell_t^{(i-1)}, \tilde{G}_t, \text{False})$.
 - 4: **for** $e = (x, y) \in \tilde{G}_t$ **do**
 - 5: **if** $\ell_t^{(h)}(x) \neq \ell_t^{(h)}(y)$ **then**
 - 6: (x, y) sends INCOMPLETE to both x and y
 - 7: **for** $x \in X$ that received INCOMPLETE message **do**
 - 8: x sends INCOMPLETE to $\ell_t^{(h)}(x)$.
 - 9: Merge all Incomplete components with leaders in the same component in $\hat{\mathcal{P}}_{t\cdot\kappa}$, to form $\hat{\mathcal{P}}_t$.
-

Proof. We recall that $(x, y) \in \tilde{G}_t$ only if $\|x - y\|_2 \leq O(t/\varepsilon)$. Therefore, if $T = \{(y_i, z_i)\}_{i=1}^{n-1}$, the cost of T in X is $\sum_{i=1}^{n-1} \|y_i - z_i\|$. However, if the edge (y_i, z_i) had weight t_i , then $\|y_i - z_i\| \leq O(t_i/\varepsilon)$, so $t_i \geq \Omega(\varepsilon) \cdot \|y_i - z_i\|$. Thus, the cost of T in \tilde{G} is at least $\sum_{i=1}^{n-1} \Omega(\varepsilon) \cdot \|y_i - z_i\|$. This completes the proof. \square

Next, we show that the minimum spanning tree cost in \tilde{G} is not much more than the Euclidean spanning tree cost. To do so, we make use of the following two simple but important propositions.

Proposition 2 (Folklore). *For any edge e connecting two points $p, q \in [0, \Delta]^d$ of length $w = \|p - q\|_2$, the probability that a randomly shifted grid of side length $L \leq \Delta$ splits the edge (i.e., p, q are not in the same cell in this grid) is at most $\frac{w \cdot \sqrt{d}}{L}$*

Proposition 3. [CS09, restated] *Let G be any weighted graph on a dataset X , with all edges having weight at most Δ and at least some sufficiently large constant. Let $\text{MST}(G)$ denote the weight of the minimum spanning tree of G . For each t , let \mathcal{P}_t represent the partition of X representing the*

Algorithm 4 PART3($X, t, \tilde{G}_t, \hat{\mathcal{P}}_{t/2}, \hat{\mathcal{P}}_t$): Generating the edges merging $\hat{\mathcal{P}}_{t/2}$ into $\hat{\mathcal{P}}_t$.

- 1: Initialize $\mathcal{P}_t^{(0)} \leftarrow \hat{\mathcal{P}}_{t/2}$, and choose an arbitrary leader for each component in $\mathcal{P}_t^{(0)}$. Let the corresponding leader mapping with respect to $\mathcal{P}_t^{(0)}$ to be $\ell_t^{(0)}$.
 - 2: **for** $i = 1$ to $h = O(\log \log n)$ **do**
 - 3: $(\mathcal{P}_t^{(i)}, \ell_t^{(i)}, \tilde{G}_t, E_t^{(i)}) \leftarrow \text{LEADERCOMPRESSION}(\mathcal{P}_t^{(i-1)}, \ell_t^{(i-1)}, \tilde{G}_t, \text{True})$.
 - 4: Create an arbitrary star on the set of leaders $\{\ell_t^{(h)}(x)\}$ in the same component in $\hat{\mathcal{P}}_t$, to create a forest of stars F_t .
 - 5: **Return** $F_t \cup \bigcup_{i=1}^h E_t^{(i)}$.
-

connected components of the threshold graph G_t of edges with weight at most t . Then,

$$\mathbf{MST}(G) = \Theta(1) \cdot \left(\sum_{t=1}^{\Delta} (|\mathcal{P}_t| - 1) \cdot t \right),$$

where the sum ranges over all $1 \leq t \leq \Delta$ where t is a power of 2.

Let \mathbf{MST} represent the true minimum spanning tree cost of X , and for any weighted graph G , let $\mathbf{MST}(G)$ represent the minimum spanning tree cost of G . First, we prove the following proposition.

Proposition 4. Suppose that $\alpha^{k-1} < t \leq \alpha^k$. Then, $\mathbb{E}[|\tilde{\mathcal{P}}_t|] \leq |\mathcal{P}_t| + \mathbf{MST} \cdot \frac{\sqrt{d} \cdot \beta}{\alpha^{k+1}}$.

Proof. Let $L := \alpha^{k+1}/\beta$. For each component $S \in \mathcal{P}_t$, let M_S represent the minimum spanning tree of X_S , and $\mathbf{MST}(S)$ represent the cost of M_S . Note that $\mathbf{MST} \geq \sum_{S \in \mathcal{P}_t} \mathbf{MST}(S)$, by Kruskal's algorithm. Now, for each edge $e \in M_S$ of weight $w(e)$, the probability that the randomly shifted grid at level L splits e is at most $\frac{w(e) \cdot \sqrt{d}}{L}$, by Proposition 2. This means that the expected number of edges in M_S across all $S \in \mathcal{P}_t$ that are cut is at most $\sum_{S \in \mathcal{P}_t} \frac{\mathbf{MST}(S) \cdot \sqrt{d}}{L} \leq \mathbf{MST} \cdot \frac{\sqrt{d}}{L} = \mathbf{MST} \cdot \frac{\sqrt{d} \cdot \beta}{\alpha^{k+1}}$.

Next, for any piece of a tree that has not been cut, all of the points in this piece will be in the same connected component in the spanner \tilde{G}_t . Therefore, the additional number of connected components is at most the number of cut edges, which completes the proof. \square

Hence, we have the following corollary.

Corollary 1. The minimum spanning tree cost in \tilde{G} , in expectation, is at most $O(1) \cdot \mathbf{MST}(X)$.

Proof. By Proposition 3, we can write $\mathbf{MST}(X) = \Theta(1) \cdot (\sum_t (|\mathcal{P}_t| - 1) \cdot t)$, where t ranges as powers of 2 from 1 to Δ . Likewise, $\mathbf{MST}(\tilde{G}) = \Theta(1) \cdot (\sum_t (|\tilde{\mathcal{P}}_t| - 1) \cdot t)$. So,

$$\mathbb{E}[\mathbf{MST}(\tilde{G})] \leq O(1) \cdot \left(\mathbf{MST}(X) + \sum_t t \cdot (\mathbb{E}[|\tilde{\mathcal{P}}_t|] - |\mathcal{P}_t|) \right).$$

Using Proposition 4, and the fact that $t \leq \alpha^k$ in Proposition 4, this is at most

$$\begin{aligned} O(1) \cdot \left(\mathbf{MST}(X) + \sum_t t \cdot \mathbf{MST}(X) \cdot \frac{\sqrt{d} \cdot \beta}{t \cdot \alpha} \right) &= O(1) \cdot \left(\mathbf{MST}(X) + \mathbf{MST}(X) \cdot \sum_t \frac{\sqrt{d} \cdot \beta}{\alpha} \right) \\ &= O(1) \cdot \mathbf{MST}(X), \end{aligned}$$

as long as $\alpha \geq \sqrt{d} \cdot \beta \cdot \log n$. \square

Hence, it suffices to find an $O(1)$ -approximate MST in the graph \tilde{G} . This tree will have cost at most $O(1) \cdot \mathbf{MST}(X)$ in \tilde{G} by Corollary 1, so by Proposition 1, it also has Euclidean cost at most $O(1/\varepsilon) \cdot \mathbf{MST}(X)$. The rest of the analysis will go into showing the algorithm finds an $O(1)$ -approximate MST in \tilde{G} .

3.2 Important Properties of Leader Compression with Early Termination

Here, we note some simple but important properties of leader compression, and some general properties of the approximate connected components we form.

Definition 5. Given two partitions \mathcal{P} and \mathcal{Q} on X , we define $\mathcal{P} \oplus \mathcal{Q}$ to be the finest partition \mathcal{R} such that $\mathcal{R} \subseteq \mathcal{P}$ and $\mathcal{R} \subseteq \mathcal{Q}$. Equivalently, it is the partition generated by merging a spanning forest of \mathcal{P} and of \mathcal{Q} , and taking the connected components.

Given a graph H on X , we define \mathcal{P}_H as the set of connected components of the graph H . We abuse notation and write $\mathcal{P} \oplus H$ to mean $\mathcal{P} \oplus \mathcal{P}_H$.

First, we note the following basic proposition.

Proposition 5. Let $\mathcal{P}^{(0)}$ be a starting partition, with a graph H . After h rounds of leader compression, let $\mathcal{P}^{(h)}$ be the set of connected components. Then, $\mathcal{P}^{(0)} \oplus H \subseteq \mathcal{P}^{(h)} \subseteq \mathcal{P}^{(0)}$.

Proof. Since we are only connecting connected components together, we trivially have that $\mathcal{P}^{(h)} \subseteq \mathcal{P}^{(0)}$. To prove that $\mathcal{P}^{(0)} \oplus H \subseteq \mathcal{P}^{(h)}$, we first consider the case that $h = 1$. In this case, we never merge two connected components in $\mathcal{P}^{(0)}$ unless they had an edge in H , so the proof is clear.

For general $h \geq 2$, we proceed by induction (base case $h = 1$ is already done). We know that $\mathcal{P}^{(0)} \oplus H \subseteq \mathcal{P}^{(h-1)} \subseteq \mathcal{P}^{(0)}$. Since $\mathcal{P}^{(h-1)} \subseteq \mathcal{P}^{(0)}$, this implies that $\mathcal{P}^{(h-1)} \oplus H \subseteq \mathcal{P}^{(0)} \oplus H$. However, we also know that $\mathcal{P}^{(h-1)} \supseteq \mathcal{P}^{(0)} \oplus H$, so this clearly implies that $\mathcal{P}^{(h-1)} \oplus H \supseteq \mathcal{P}^{(0)} \oplus H$. So, in fact $\mathcal{P}^{(h-1)} \oplus H = \mathcal{P}^{(0)} \oplus H$. Therefore, by the base case, if we performing a single round of leader compression on $\mathcal{P}^{(h-1)}$ to obtain $\mathcal{P}^{(h)}$, we have that $\mathcal{P}^{(0)} \oplus H = \mathcal{P}^{(h-1)} \oplus H \subseteq \mathcal{P}^{(h)}$, which completes the proof. \square

We next note a simple proposition about $\bar{\mathcal{P}}_t$.

Proposition 6. For $t = \alpha^k$, we have that $\bar{\mathcal{P}}_t \supseteq \tilde{\mathcal{P}}_t$. Hence, every connected component $C \in \bar{\mathcal{P}}_t$ is in a single cell at level α^{k+1}/β .

Proof. Initially, $\mathcal{P}_t^{(0)}$ is based on the quadtree at level α^k/β , and every two points in the same cell have distance at most $\alpha^k \cdot \sqrt{d}/\beta \leq t$ and are also in the same cell at level α^{k+1}/β , which means they must be in the same connected component in $\tilde{\mathcal{P}}_t$. Hence $\tilde{\mathcal{P}}_t \subseteq \mathcal{P}_t^{(0)}$. Then, we perform h rounds of leader compression on $\mathcal{P}_t^{(0)}$ with respect to \tilde{G}_t , to obtain $\bar{\mathcal{P}}_t \supseteq \mathcal{P}_t^{(0)} \oplus \tilde{G}_t$, by Proposition 5. However, since $\tilde{\mathcal{P}}_t \subseteq \mathcal{P}_t^{(0)}$, this means $\mathcal{P}_t^{(0)} \oplus \tilde{G}_t = \mathcal{P}_t^{(0)} \oplus \tilde{\mathcal{P}}_t = \tilde{\mathcal{P}}_t$.

Finally, for $t = \alpha^k$, \tilde{G}_t never crosses a cell at level α^{k+1}/β . Therefore, any component in $\bar{\mathcal{P}}_t$ is contained in a single cell at level α^{k+1}/β . \square

The following important lemma helps us understand the approximate connected components $\tilde{\mathcal{P}}_t$.

Lemma 1. *For every t a power of 2, the following hold.*

1. *If $t \leq \alpha^k$, then every connected component in $\hat{\mathcal{P}}_t$ is contained in the same cell at level α^{k+1}/β .*
2. *$\hat{\mathcal{P}}_t \subseteq \tilde{\mathcal{P}}_t$.*
3. *$\hat{\mathcal{P}}_t \subseteq \hat{\mathcal{P}}_{t/2}$.*

Proof. We shall prove these three properties in an inductive manner. In the base case, we prove the first two claims for $t = \alpha^k$, and that $\hat{\mathcal{P}}_{\alpha^{k+1}} \subseteq \hat{\mathcal{P}}_{\alpha^k}$. For the inductive step, we consider creating $\hat{\mathcal{P}}_t$, given $\hat{\mathcal{P}}_{t/\kappa}$ and $\hat{\mathcal{P}}_{t \cdot \kappa}$, where $\alpha^k \leq t/\kappa$ and $t \cdot \kappa \leq \alpha^{k+1}$. We inductively assume that every connected component in $\hat{\mathcal{P}}_{t/\kappa}$ and $\hat{\mathcal{P}}_{t \cdot \kappa}$ is in a single cell at level α^{k+2}/β , that $\hat{\mathcal{P}}_{t/\kappa} \subseteq \tilde{\mathcal{P}}_{t/\kappa}$ and $\hat{\mathcal{P}}_{t \cdot \kappa} \subseteq \tilde{\mathcal{P}}_{t \cdot \kappa}$, and finally, $\hat{\mathcal{P}}_{t \cdot \kappa} \subseteq \hat{\mathcal{P}}_{t/\kappa}$. Then, the inductive step proves that every connected component in $\hat{\mathcal{P}}_t$ is in a single cell at level α^{k+2}/β , that $\hat{\mathcal{P}}_t \subseteq \tilde{\mathcal{P}}_t$, and that $\hat{\mathcal{P}}_{t \cdot \kappa} \subseteq \hat{\mathcal{P}}_t \subseteq \hat{\mathcal{P}}_{t/\kappa}$.

Base case: $g = 0$, i.e., $t = \alpha^k$. Suppose that $x, y \in X$ are in the same connected component in \tilde{G}_t . Then, either x, y are in the same component in $\bar{\mathcal{P}}_t$ (and thus in $\hat{\mathcal{P}}_t$), or x, y are both in incomplete connected components in $\bar{\mathcal{P}}_t$. Since $\bar{\mathcal{P}}_t \supseteq \tilde{\mathcal{P}}_t$ by Proposition 6, this means the leaders of x and y (in $\bar{\mathcal{P}}_t$) are in the same connected component in \tilde{G}_t , which means they are in the same cell at level α^{k+1}/β . So, these two components are merged, and thus x, y are in the same component in $\hat{\mathcal{P}}_t$. Hence, $\hat{\mathcal{P}}_t \subseteq \tilde{\mathcal{P}}_t$. In addition, every connected component in $\bar{\mathcal{P}}_t$ is contained in a single cell at level α^{k+1}/β , by Proposition 6. Therefore, every connected component in $\hat{\mathcal{P}}_t$ is fully contained in a single cell at level α^{k+1}/β . Finally, all points in a single cell at level α^{k+1}/β start off connected as $\mathcal{P}_{\alpha^{k+1}}^{(0)}$, so this also implies that $\hat{\mathcal{P}}_{\alpha^{k+1}} \subseteq \hat{\mathcal{P}}_{\alpha^k}$.

Inductive step: First, we show that $\hat{\mathcal{P}}_{t \cdot \kappa} \subseteq \hat{\mathcal{P}}_t \subseteq \hat{\mathcal{P}}_{t/\kappa}$. The claim that $\hat{\mathcal{P}}_t \subseteq \hat{\mathcal{P}}_{t/\kappa}$ is immediate, because we generate $\hat{\mathcal{P}}_t$ by performing leader compression on $\hat{\mathcal{P}}_{t/\kappa}$, and then doing additional merging. So, we prove that $\hat{\mathcal{P}}_{t \cdot \kappa} \subseteq \hat{\mathcal{P}}_t$. By our inductive hypothesis, $\hat{\mathcal{P}}_{t/\kappa} \supseteq \hat{\mathcal{P}}_{t \cdot \kappa}$. To generate $\bar{\mathcal{P}}_t$, we perform rounds of leader compression on $\hat{\mathcal{P}}_{t/\kappa}$ with respect to \tilde{G}_t . But, we know that $\bar{\mathcal{P}}_t \supseteq \tilde{\mathcal{P}}_{t \cdot \kappa} \supseteq \hat{\mathcal{P}}_{t \cdot \kappa}$ (the last part following from our inductive hypothesis), so $\bar{\mathcal{P}}_t \supseteq \hat{\mathcal{P}}_{t \cdot \kappa}$. Finally, we only connect components of $\bar{\mathcal{P}}_t$ to form $\hat{\mathcal{P}}_t$ if they are incomplete and their leaders are in the same component of $\hat{\mathcal{P}}_{t \cdot \kappa}$, and since $\bar{\mathcal{P}}_t \supseteq \hat{\mathcal{P}}_{t \cdot \kappa}$, this implies that $\hat{\mathcal{P}}_t \supseteq \hat{\mathcal{P}}_{t \cdot \kappa}$, as desired.

Next, by the inductive hypothesis, we have that $\hat{\mathcal{P}}_{t,\kappa} \sqsupseteq \hat{\mathcal{P}}_{\alpha^{k+1}}$, and by the base case, every connected component in $\hat{\mathcal{P}}_{\alpha^{k+1}}$ is contained in a single cell at level α^{k+2}/β . Therefore, the same holds for $\hat{\mathcal{P}}_t$, since $\hat{\mathcal{P}}_t \sqsupseteq \hat{\mathcal{P}}_{t,\kappa} \sqsupseteq \hat{\mathcal{P}}_{\alpha^{k+1}}$.

To finish the proof, we show that $\hat{\mathcal{P}}_t \sqsubseteq \tilde{\mathcal{P}}_t$. Consider any x, y in the same component in $\tilde{\mathcal{P}}_t$. If they are in the same component in $\bar{\mathcal{P}}_t$, then they are also in the same component in $\hat{\mathcal{P}}_t$. Otherwise, x and y are in two different incomplete components in $\bar{\mathcal{P}}_t$, but are in the same connected component in $\hat{\mathcal{P}}_{t,\kappa}$, since $\hat{\mathcal{P}}_t \sqsupseteq \hat{\mathcal{P}}_{t,\kappa} \sqsupseteq \hat{\mathcal{P}}_{\alpha^{k+1}}$ by the inductive hypothesis. Recall that we showed in the first paragraph of the proof that $\bar{\mathcal{P}}_t \sqsupseteq \hat{\mathcal{P}}_{t,\kappa}$. Therefore, since x and y are in the same connected component in $\hat{\mathcal{P}}_{t,\kappa}$, so are the leaders of x and y in $\bar{\mathcal{P}}_t$. Therefore, the components get merged together, so x and y are in the same connected component in $\hat{\mathcal{P}}_t$. \square

Next, we prove an probabilistic result about leader compression.

Lemma 2. *Consider any partitioning $\mathcal{P}^{(0)}$ of X , and a graph H . After h rounds of leader compression, $\mathbb{E}[|\mathcal{P}^{(h)}| - |\mathcal{P}^{(0)} \oplus H|] \leq (3/4)^h \cdot (|\mathcal{P}^{(0)}| - |\mathcal{P}^{(0)} \oplus H|)$.*

Proof. We first prove the lemma for $h = 1$. Fix a connected component $C \in \mathcal{P}^{(0)} \oplus H$, and suppose that C splits into C_1, \dots, C_r in $\mathcal{P}^{(0)}$, for some $r \geq 1$. If $r = 1$, then C still splits into 1 connected component in $\mathcal{P}^{(1)}$, since $\mathcal{P}^{(0)} \oplus H \sqsubseteq \mathcal{P}^{(1)} \sqsubseteq \mathcal{P}^{(0)}$. Otherwise, for every C_j , there must exist some C_i , $i \neq j$, such that there is an edge in H connecting C_i to C_j (or else C_j is not connected to the rest of C). Hence, with probability at least $1/4$, C_j has value 0, C_i has value 1, and there must be a message sent from some point in C_i to some point in C_j . This means that C_j will become merged to some connected component (not necessarily C_i) and will have a new head. Hence, if $r > 1$, the value of r multiplies by at most $3/4$ in expectation after a single round of compression.

Hence, in all cases, the expectation of $r - 1$ multiplies by at most $3/4$. However, $|\mathcal{P}^{(0)}| - |\mathcal{P}^{(0)} \oplus H|$ precisely equals the sum of $r - 1$ across all connected components in $\mathcal{P}^{(0)} \oplus H$ at the beginning, and $|\mathcal{P}^{(1)}| - |\mathcal{P}^{(0)} \oplus H|$ precisely equals the sum of $r - 1$ across all connected components in $\mathcal{P}^{(0)} \oplus H$ at the end. Therefore, $\mathbb{E}[|\mathcal{P}^{(1)}| - |\mathcal{P}^{(0)} \oplus H|] \leq (3/4) \cdot (|\mathcal{P}^{(0)}| - |\mathcal{P}^{(0)} \oplus H|)$.

For general values of h , note that $\mathcal{P}^{(0)} \oplus H = \mathcal{P}^{(h-1)} \oplus H$. Therefore, we have $\mathbb{E}[|\mathcal{P}^{(h)}| - |\mathcal{P}^{(0)} \oplus H|] \leq (3/4) \cdot (|\mathcal{P}^{(h-1)}| - |\mathcal{P}^{(0)} \oplus H|)$, which, after an inductive argument, completes the proof. \square

3.3 Part 1 of the Algorithm

In this subsection, we analyze Part 1 of the algorithm for some fixed $t = \alpha^k$.

First, we recall the following, which is immediate by combining Proposition 6 and Lemma 1.

Proposition 7. *For every t a power of α , $\bar{\mathcal{P}}_t \sqsupseteq \tilde{\mathcal{P}}_t \sqsupseteq \hat{\mathcal{P}}_t$.*

Recall that by our algorithmic construction (Part 1), and because $\beta > \sqrt{d}$, all points in a cell at level t/β must be in the same connected component for any of $\tilde{\mathcal{P}}_t, \bar{\mathcal{P}}_t, \hat{\mathcal{P}}_t$.

Lemma 3. *For some fixed $t = \alpha^k$, let $\mathcal{P}_t^{(h)}$ be the partitioning after h rounds of leader compression on $\mathcal{P}_t^{(0)}$ with respect to \tilde{G}_t . Then, the expected number of incomplete connected components in $\bar{\mathcal{P}}_t$ is at most $(3/4)^h \cdot \alpha \cdot \frac{\text{MST}(\tilde{G})}{t}$.*

Proof. First, note that $\mathcal{P}_t^{(0)} \oplus \tilde{G}_t = \tilde{\mathcal{P}}_t$, since $\mathcal{P}_t^{(0)} \supseteq \tilde{\mathcal{P}}_t$. Therefore, by Lemma 2, $\mathbb{E}[|\mathcal{P}_t^{(h)}| - |\tilde{\mathcal{P}}_t|] \leq (3/4)^h \cdot (|\mathcal{P}_t^{(0)}| - |\tilde{\mathcal{P}}_t|) \leq (3/4)^h \cdot (|\mathcal{P}_t^{(0)}| - 1)$. Next, note that $|\mathcal{P}_t^{(0)}|$ equals the number of distinct nonempty cells in the quadtree Q at side length t/β . Any points in distinct cells can only be connected by an edge in the spanner of length at least $2\alpha^{k-1}$. Hence, $\mathbf{MST}(\tilde{G}) \geq 2\alpha^{k-1} \cdot (|\mathcal{P}_t^{(0)}| - 1)$, which means that $|\mathcal{P}_t^{(0)}| - 1 \leq \frac{\mathbf{MST}(\tilde{G})}{2\alpha^{k-1}}$. Hence, $\mathbb{E}[|\mathcal{P}_t^{(h)}| - |\tilde{\mathcal{P}}_t|] \leq (3/4)^h \cdot \frac{\mathbf{MST}(\tilde{G})}{2\alpha^{k-1}} = (3/4)^h \cdot \alpha \cdot \frac{\mathbf{MST}(\tilde{G})}{2t}$.

Next, we note that the number of incomplete connected components in $\mathcal{P}_t^{(h)}$ is at most $2 \cdot (|\mathcal{P}_t^{(h)}| - |\tilde{\mathcal{P}}_t|)$. To see why, if any connected component $C \in \tilde{\mathcal{P}}_t$ is split into r_C connected components in $\mathcal{P}_t^{(h)}$, then $|\mathcal{P}_t^{(h)}| - |\tilde{\mathcal{P}}_t| = \sum_{C \in \tilde{\mathcal{P}}_t} (r_C - 1)$, but the number of incomplete connected components is $\sum_{C \in \tilde{\mathcal{P}}_t} r_C \cdot \mathbf{1}(r_C \geq 2)$. Since $2(r_C - 1) \geq r_C$ whenever $r_C \geq 2$, this means that $2 \cdot (|\mathcal{P}_t^{(h)}| - |\tilde{\mathcal{P}}_t|)$ is at least the number of incomplete connected components. Hence, the expected number of incomplete connected components is at most $(3/4)^h \cdot \alpha \cdot \frac{\mathbf{MST}(\tilde{G})}{t}$. \square

Hence, the following holds.

Corollary 2. *For some fixed $t = \alpha^k$, $\mathbb{E}[|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t|] \leq (3/4)^h \cdot \frac{\alpha}{t} \cdot \mathbf{MST}(\tilde{G})$.*

Proof. By Proposition 7, we have that $|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t| \leq |\bar{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t|$. However, to form $\hat{\mathcal{P}}_t$ from $\bar{\mathcal{P}}_t$, we only merge incomplete connected components together. So, $|\bar{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t|$ is at most the number of incomplete connected components, minus 1. Hence, by Lemma 3, the expectation is at most $(3/4)^h \cdot \frac{\alpha}{t} \cdot \mathbf{MST}(\tilde{G})$. \square

3.4 Part 2 of the Algorithm

In this section, we consider constructing $\hat{\mathcal{P}}_t$ where $\alpha^k < t < \alpha^{k+1}$.

Proposition 8. *Let \mathcal{P} be any partition, and $\mathcal{P}_1, \mathcal{P}_2 \subseteq \mathcal{P}$ be coarser partitions. Then, $|\mathcal{P}| + |\mathcal{P}_1 \oplus \mathcal{P}_2| \geq |\mathcal{P}_1| + |\mathcal{P}_2|$.*

Proof. Let $M = |\mathcal{P}|$. Consider graphs on $[M]$, where each vertex represents a partition in \mathcal{P} . Let G_1 be the spanning forest on $[M]$ that generates the partition \mathcal{P}_1 (via the connected components of G_1) and G_2 be the graph that generates the partition \mathcal{P}_2 on $[M]$. Then, $G_1 \cup G_2$ generates the partition $\mathcal{P}_1 \oplus \mathcal{P}_2$. Since G_1, G_2 are forests, the number of edges $|G_1|, |G_2|$ equal $M - |\mathcal{P}_1|, M - |\mathcal{P}_2|$, respectively. Hence, $|G_1 \cup G_2| \leq 2M - |\mathcal{P}_1| - |\mathcal{P}_2|$. The number of connected components in $G_1 \cup G_2$ is at least $M - |G_1 \cup G_2|$, with equality if and only if $G_1 \cup G_2$ is also a forest. Therefore, $|\mathcal{P}_1 \oplus \mathcal{P}_2| \geq M - (2M - |\mathcal{P}_1| - |\mathcal{P}_2|) = |\mathcal{P}_1| + |\mathcal{P}_2| - M$, which completes the proof. \square

Define $T := |\tilde{\mathcal{P}}_{t/\kappa}| - |\hat{\mathcal{P}}_{t/\kappa}|$. Since $\tilde{\mathcal{P}}_t, \hat{\mathcal{P}}_{t/\kappa} \subseteq \tilde{\mathcal{P}}_{t/\kappa}$ by Lemma 1, we have the following corollary.

Corollary 3. *We have that $|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t| \leq T$.*

We recall that $\tilde{\mathcal{P}}_t \supseteq \hat{\mathcal{P}}_t$. The rest of this section is devoted to bounding $\mathbb{E}[|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t|]$, which generalizes Corollary 2 to t not necessarily a power of α .

Next, we bound the number of incomplete connected components, similar to Lemma 3.

Lemma 4. For some fixed $\alpha^k < t < \alpha^{k+1}$, the expected number of incomplete connected components in $\bar{\mathcal{P}}_t$ is at most $(3/4)^h \cdot \frac{\alpha^2}{t} \cdot \mathbf{MST}(\tilde{G})$.

Proof. Define $\mathcal{P}_t^{(0)} := \hat{\mathcal{P}}_{t/\kappa}$ as the starting partition before leader compression. Let $\mathcal{P}_t^{(h)}$ be the partition after performing h rounds of leader compression on $\mathcal{P}_t^{(0)}$ with respect to \tilde{G}_t . By Lemma 2, we know that $\mathbb{E} \left[|\mathcal{P}_t^{(h)}| - |\mathcal{P}_t^{(0)} \oplus \tilde{\mathcal{P}}_t| \mid \mathcal{P}_t^{(0)} \right] \leq (3/4)^h \cdot \left(|\mathcal{P}_t^{(0)}| - |\mathcal{P}_t^{(0)} \oplus \tilde{\mathcal{P}}_t| \right)$. Since $\mathcal{P}_t^{(0)} = \hat{\mathcal{P}}_{t/\kappa}$, this means that

$$\mathbb{E} \left[|\mathcal{P}_t^{(h)}| - |\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t| \mid \hat{\mathcal{P}}_{t/\kappa} \right] \leq (3/4)^h \cdot \left(|\hat{\mathcal{P}}_{t/\kappa}| - |\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t| \right).$$

By Part 3 of Lemma 1, know that $\hat{\mathcal{P}}_{t/\kappa} \subseteq \hat{\mathcal{P}}_{\alpha^k}$, so we can bound $|\hat{\mathcal{P}}_{t/\kappa}|$ as at most the number of nonempty cells at level α^k/β . Moreover, $|\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t| \geq 1$, which means that $|\hat{\mathcal{P}}_{t/\kappa}| - |\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t|$ is at most the number of nonempty cells at level α^k/β minus 1, which as we saw in Lemma 3 is at most $\frac{\mathbf{MST}(\tilde{G})}{2\alpha^{k-1}} \leq \alpha^2 \cdot \frac{\mathbf{MST}(\tilde{G})}{2t}$, because $t \leq \alpha^{k+1}$. So, by removing the conditioning on $\hat{\mathcal{P}}_{t/\kappa}$, we have that $\mathbb{E} \left[|\mathcal{P}_t^{(h)}| - |\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t| \right] \leq (3/4)^h \cdot \alpha^2 \cdot \frac{\mathbf{MST}(\tilde{G})}{2t}$.

Next, recall that $\mathcal{P}_t^{(h)} \supseteq \hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t$, and every incomplete connected component in $\mathcal{P}_t^{(h)}$ has an edge leaving it in \tilde{G}_t . So, if any component $C \in \hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t$ is split into r_C connected components in $\mathcal{P}_t^{(h)}$, then $|\mathcal{P}_t^{(h)}| - |\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t| = \sum_{C \in \hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t} (r_C - 1)$, but the number of incomplete connected components is $\sum_{C \in \hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t} r_C \cdot \mathbf{1}(r_C \geq 2)$. Since $2(r_C - 1) \geq r_C$ whenever $r_C \geq 2$, this means that the number of incomplete connected components is at most $2 \cdot (|\mathcal{P}_t^{(h)}| - |\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t|)$. Hence, the expected number of incomplete connected components is at most $(3/4)^h \cdot \alpha^2 \cdot \frac{\mathbf{MST}(\tilde{G})}{t}$. \square

We next show that $|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t|$ is small, in expectation.

Lemma 5. Suppose t is a power of 2 with $\alpha^k < t < \alpha^{k+1}$, and let $\kappa = 2^{2^{v_2(\log_2 t)}}$. Then, $\mathbb{E} \left[|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t| \right] \leq \mathbb{E} \left[|\tilde{\mathcal{P}}_{t/\kappa}| - |\hat{\mathcal{P}}_{t/\kappa}| \right] + (3/4)^h \cdot \frac{\alpha^2}{t} \cdot \mathbf{MST}(\tilde{G})$.

Proof. Since we only merge incomplete connected components together to go from $\bar{\mathcal{P}}_t = \mathcal{P}_t^{(h)}$ to $\hat{\mathcal{P}}_t$, we know that

$$\mathbb{E} \left[|\mathcal{P}_t^{(h)}| - |\hat{\mathcal{P}}_t| \right] \leq (3/4)^h \cdot \frac{\alpha^2}{t} \cdot \mathbf{MST}(\tilde{G}), \quad (3)$$

by Lemma 4. Next,

$$\mathbb{E} \left[|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t| \right] \leq \mathbb{E} \left[|\tilde{\mathcal{P}}_{t/\kappa}| - |\hat{\mathcal{P}}_{t/\kappa}| \right], \quad (4)$$

by Corollary 3. Finally, $\mathcal{P}_t^{(h)} \supseteq \hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t$, by Proposition 5. Therefore,

$$\mathbb{E} \left[|\hat{\mathcal{P}}_{t/\kappa} \oplus \tilde{\mathcal{P}}_t| - |\mathcal{P}_t^{(h)}| \right] \leq 0. \quad (5)$$

Adding Equations (3), (4), and (5) together completes the proof. \square

By a simple inductive argument, we have $|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t|$ is small in expectation. Specifically, we have the following corollary.

Corollary 4. For any t a power of 2, $\hat{\mathcal{P}}_t \sqsubseteq \tilde{\mathcal{P}}_t$, and $\mathbb{E} \left[|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t| \right] \leq (3/4)^h \cdot \frac{\alpha^3}{t} \cdot \mathbf{MST}(\tilde{G})$.

Proof. Suppose $\alpha^k \leq t < \alpha^{k+1}$. We can construct a sequence $t = t_0 > \dots > t_r = \alpha^k$, where $t_{i+1} = t_i / 2^{2^{v_2(\log_2 t_i)}}$. Then, we can form a telescoping sum

$$\mathbb{E} \left[|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t| \right] = \mathbb{E} \left[|\tilde{\mathcal{P}}_{t_r}| - |\hat{\mathcal{P}}_{t_r}| \right] + \sum_{i=0}^{r-1} \left(\mathbb{E} \left[|\tilde{\mathcal{P}}_{t_i}| - |\hat{\mathcal{P}}_{t_i}| \right] - \mathbb{E} \left[|\tilde{\mathcal{P}}_{t_{i+1}}| - |\hat{\mathcal{P}}_{t_{i+1}}| \right] \right).$$

Using Corollary 2 and Lemma 5, this is at most

$$\begin{aligned} &\leq \left(\frac{3}{4} \right)^h \cdot \frac{\alpha}{t_r} \cdot \mathbf{MST}(\tilde{G}) + \sum_{i=0}^{r-1} \left(\frac{3}{4} \right)^h \cdot \frac{\alpha^2}{t_i} \cdot \mathbf{MST}(\tilde{G}) \\ &\leq \left(\frac{3}{4} \right)^h \cdot \alpha^2 \cdot \mathbf{MST}(\tilde{G}) \cdot \sum_{i=0}^r \frac{1}{t_i} \\ &\leq \left(\frac{3}{4} \right)^h \cdot \alpha^2 \cdot \mathbf{MST}(\tilde{G}) \cdot \left(\frac{1}{t} + \frac{2}{t} + \frac{4}{t} + \dots + \frac{1}{\alpha^k} \right) \\ &\leq \left(\frac{3}{4} \right)^h \cdot \alpha^2 \cdot \mathbf{MST}(\tilde{G}) \cdot \frac{2}{\alpha^k} \\ &\leq \left(\frac{3}{4} \right)^h \cdot \frac{\alpha^3}{t} \cdot \mathbf{MST}(\tilde{G}), \end{aligned}$$

since $t \leq \alpha^{k+1}/2$ so $\frac{2}{\alpha^k} \leq \frac{\alpha}{t}$. □

3.5 Part 3 of the Algorithm

From Lemma 1, we have generated a nested set of partitions $\hat{\mathcal{P}}_1 \sqsupseteq \hat{\mathcal{P}}_2 \sqsupseteq \hat{\mathcal{P}}_4 \sqsupseteq \dots$. In Part 3, we generate the edges: we must show that the cost of the edges we produce is at most $O(1) \cdot \mathbf{MST}(\tilde{G})$, and that the edges we produce generate a spanning tree.

When connecting $\hat{\mathcal{P}}_{t/2}$ into $\hat{\mathcal{P}}_t$, we generate a new partition $\mathcal{P}'_t^{(h)}$ after h rounds of leader compression on $\hat{\mathcal{P}}_{t/2}$ with respect to \tilde{G}_t . First, we note that $\mathcal{P}'_t^{(h)}$ still refines $\hat{\mathcal{P}}_t$.

Proposition 9. We have $\mathcal{P}'_t^{(h)} \sqsupseteq \hat{\mathcal{P}}_t$.

Proof. By Proposition 5, $\mathcal{P}'_t^{(h)} \sqsupseteq \hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t$. Moreover, by Lemma 1, $\hat{\mathcal{P}}_{t/2} \sqsupseteq \hat{\mathcal{P}}_t$ and $\tilde{\mathcal{P}}_t \sqsupseteq \hat{\mathcal{P}}_t$, which means that $\hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t \sqsupseteq \hat{\mathcal{P}}_t$. In summary, $\mathcal{P}'_t^{(h)} \sqsupseteq \hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t \sqsupseteq \hat{\mathcal{P}}_t$. □

Let \hat{G}_t represent the edges generated when connecting $\hat{\mathcal{P}}_{t/2}$ into $\hat{\mathcal{P}}_t$ (and $\hat{G}_1 = \emptyset$).

Lemma 6. $\bigcup_{t' \leq \alpha^H} \hat{G}_{t'}$ is a spanning tree of X .

Proof. We prove inductively that $\bigcup_{t' \leq t} \hat{G}_{t'}$ forms a spanning forest for $\hat{\mathcal{P}}_t$.

For the base case $t = 1$, recall that we assume the minimum distance between any two points in X is at least α^{100} , which means by Part 1 of Lemma 1, every connected component in $\hat{\mathcal{P}}_1$ consists of a single element. Since \hat{G}_1 has no edges, the base case follows.

Define $\hat{G}_{\leq t} := \bigcup_{t' \leq t} \hat{G}_{t'}$. For the inductive step, assume that $\hat{G}_{\leq t/2}$ generates a spanning forest for $\hat{\mathcal{P}}_{t/2}$. Then, leader compression only adds edges to connect two distinct connected components together, so when adding the edges generated by h rounds of leader compression on $\hat{\mathcal{P}}_{t/2}$ with respect to \tilde{G}_t , we obtain a spanning forest for $\mathcal{P}'_t{}^{(h)}$, where $\mathcal{P}'_t{}^{(h)} \supseteq \hat{\mathcal{P}}_t$ by Proposition 9. Finally, we connect all disconnected components in each partition of $\hat{\mathcal{P}}_t$, which generates a spanning forest for $\hat{\mathcal{P}}_t$.

Due to Part 2 of Lemma 1, $\hat{\mathcal{P}}_t$ refines $\tilde{\mathcal{P}}_t$. Since $\tilde{\mathcal{P}}_t$ contains only one part X for $t = \alpha^H$, $\bigcup_{t' \leq \alpha^H} \hat{G}_{t'}$ is a spanning tree of X . \square

Now, we fix a value t , and consider the edges generated to connect $\hat{\mathcal{P}}_{t/2}$ into $\hat{\mathcal{P}}_t$. We will bound the sum of the weights of these edges.

Lemma 7. *The sum of the weights of all edges in \hat{G}_t , in expectation, is at most $t \cdot (|\tilde{\mathcal{P}}_{t/2}| - |\tilde{\mathcal{P}}_t|) + (3/4)^h \cdot \left(\alpha^3 + \frac{2\alpha^4\sqrt{d}}{\beta}\right) \cdot \mathbf{MST}(\tilde{G})$.*

Proof. By Proposition 5, we have $\hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t \subseteq \mathcal{P}'_t{}^{(h)} \subseteq \hat{\mathcal{P}}_{t/2}$, and by Lemma 2, $\mathbb{E} \left[|\mathcal{P}'_t{}^{(h)}| - |\hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t| \right] \leq (3/4)^h \cdot (|\hat{\mathcal{P}}_{t/2}| - |\hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t|)$. As in Lemma 4, we know that $|\hat{\mathcal{P}}_{t/2}| - |\hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t| \leq \alpha^2 \cdot \frac{\mathbf{MST}(\tilde{G})}{t}$, so $\mathbb{E} \left[|\mathcal{P}'_t{}^{(h)}| - |\hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t| \right] \leq (3/4)^h \cdot \alpha^2 \cdot \frac{\mathbf{MST}(\tilde{G})}{t}$. In addition, $|\hat{\mathcal{P}}_{t/2} \oplus \tilde{\mathcal{P}}_t| - |\tilde{\mathcal{P}}_t| \leq 0$, and $\mathbb{E} \left[|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t| \right] \leq (3/4)^h \cdot \frac{\alpha^3}{t} \cdot \mathbf{MST}(\tilde{G})$ by Corollary 4. So, $\mathbb{E} \left[|\mathcal{P}'_t{}^{(h)}| - |\hat{\mathcal{P}}_t| \right] \leq (3/4)^h \cdot \frac{2\alpha^3}{t} \cdot \mathbf{MST}(\tilde{G})$.

Assume $\alpha^k \leq t < \alpha^{k+1}$. At level t , we produce $|\hat{\mathcal{P}}_{t/2}| - |\mathcal{P}'_t{}^{(h)}|$ edges of weight t , and then $|\mathcal{P}'_t{}^{(h)}| - |\hat{\mathcal{P}}_t|$ edges to combine all of $\hat{\mathcal{P}}_t$ together. But by Part 1 of Lemma 1, the diameter of any cell in $\hat{\mathcal{P}}_t$ is at most $\alpha^{k+1} \cdot \sqrt{d}/\beta$, so the total weight of all the edges, in expectation, is at most

$$\begin{aligned} t \cdot \mathbb{E} \left[|\hat{\mathcal{P}}_{t/2}| - |\mathcal{P}'_t{}^{(h)}| \right] + \frac{\alpha^{k+1}\sqrt{d}}{\beta} \cdot \mathbb{E} \left[|\mathcal{P}'_t{}^{(h)}| - |\hat{\mathcal{P}}_t| \right] &\leq t \cdot \mathbb{E} \left[|\hat{\mathcal{P}}_{t/2}| - |\mathcal{P}'_t{}^{(h)}| \right] + \frac{\alpha^{k+1}\sqrt{d}}{\beta} \cdot \left(\frac{3}{4}\right)^h \cdot \frac{2\alpha^3}{t} \cdot \mathbf{MST}(\tilde{G}) \\ &\leq t \cdot \mathbb{E} \left[|\hat{\mathcal{P}}_{t/2}| - |\mathcal{P}'_t{}^{(h)}| \right] + \frac{2\alpha^4\sqrt{d}}{\beta} \cdot \left(\frac{3}{4}\right)^h \cdot \mathbf{MST}(\tilde{G}) \\ &\leq t \cdot \mathbb{E} \left[|\tilde{\mathcal{P}}_{t/2}| - |\hat{\mathcal{P}}_t| \right] + \frac{2\alpha^4\sqrt{d}}{\beta} \cdot \left(\frac{3}{4}\right)^h \cdot \mathbf{MST}(\tilde{G}). \end{aligned}$$

Moreover, $\mathbb{E}[|\tilde{\mathcal{P}}_t| - |\hat{\mathcal{P}}_t|] \leq (3/4)^h \cdot \frac{\alpha^3}{t} \cdot \mathbf{MST}(\tilde{G})$ by Corollary 4, so $\mathbb{E}[|\tilde{\mathcal{P}}_{t/2}| - |\hat{\mathcal{P}}_t|] \leq |\tilde{\mathcal{P}}_{t/2}| - |\tilde{\mathcal{P}}_t| + (3/4)^h \cdot \frac{\alpha^3}{t} \cdot \mathbf{MST}(\tilde{G})$. So overall, the weight of the edges we added, in expectation is at most

$$t \cdot (|\tilde{\mathcal{P}}_{t/2}| - |\tilde{\mathcal{P}}_t|) + (3/4)^h \cdot \left(\alpha^3 + \frac{2\alpha^4\sqrt{d}}{\beta}\right) \cdot \mathbf{MST}(\tilde{G}). \quad \square$$

Adding this over all levels t , we get this is at most $\mathbf{MST}(\tilde{G}) + (3/4)^h \cdot \log n \cdot \left(\alpha^3 + \frac{2\alpha^4\sqrt{d}}{\beta} \right) \cdot \mathbf{MST}(\tilde{G})$. So, we just need $h = O(\log \log n)$.

Lemma 8. *If h is a sufficiently large $\Theta(\log \log n)$, $\bigcup_{t' \leq \alpha^h} \hat{G}_{t'}$ is an $O(1)$ -approximate \mathbf{MST} for X in expectation.*

Proof. According to Lemma 7, the expected total weights of $\bigcup_{t' \leq \alpha^h} \hat{G}_{t'}$ is at most

$$\mathbf{MST}(\tilde{G}) + (3/4)^h \cdot \log n \cdot \left(\alpha^3 + \frac{2\alpha^4\sqrt{d}}{\beta} \right) \cdot \mathbf{MST}(\tilde{G}) = O(1) \cdot \mathbf{MST}(\tilde{G}).$$

According to Corollary 1 and Proposition 1, $\bigcup_{t' \leq \alpha^h} \hat{G}_{t'}$ is a $O(1)$ -approximate \mathbf{MST} for X in expectation. \square

4 Euclidean MST in the MPC Model

In this section, we describe how to implement our algorithms in the Massively Parallel Computation (MPC) model. In Section 4.1, we briefly introduce some useful existing algorithmic primitives in the MPC model. In Section 4.2, we show the detailed implementation of our Euclidean MST algorithm in the MPC model.

4.1 Existing Algorithmic Primitives in the MPC model

One of the most basic MPC algorithmic primitive is sorting.

Theorem 3 ([Goo99, GSZ11]). *Given size N input data where each data item has size at most $N^{o(1)}$, there is a fully scalable MPC algorithm to sort and index (rank) these data items in $O(1)$ rounds and using $O(N)$ total space.*

Sorting can be used to build other MPC subroutines. One important use case of sorting is to simulate a PRAM algorithm in the MPC model. The PRAM model is a classic model of parallel computation. Roughly speaking, there is a shared memory and multiple processors in the PRAM model. The computation proceeds in rounds, and each processor does at most one algorithmic operation in one round. In addition, processors can simultaneously read and write the shared memory cells in a round. The efficiency of a PRAM algorithm is measured by *depth* (longest chain of dependencies of the computation) and *work* (total running time over processors). A CRCW (Concurrent Read Concurrent Write) PRAM algorithm means that processors can simultaneously read and write the same shared memory cell in a round. It was shown that any PRAM algorithm can be simulated in the MPC model (see e.g. [GSZ11, ASS⁺18]). We refer readers to Appendix E of [ASS⁺18] for detailed implementations. We formally state the guarantees of the simulation as the following.

Theorem 4 ([GSZ11, ASS⁺18]). *Given a CRCW PRAM algorithm with $O(D)$ depth and $O(W)$ work, it can be simulated by a fully scalable MPC algorithm with $O(D)$ rounds and $O(W)$ space.*

One application of Theorem 4 is the following simultaneous access problem: several arrays with total size $O(m)$ are stored distributedly on the machines, and p machines want to simultaneously access some entries of these arrays where each machine has s non-adaptive queries. This operation can be done in $O(1)$ rounds and total space $O(m + p \cdot s)$, and it is fully scalable. One application of simultaneous access is duplicating and broadcasting a small size message.

Theorem 5 (See e.g., [ASS⁺18]). *Given a data item with size $N^{o(1)}$ where N is the target output size, there is a fully scalable MPC algorithm⁷ uses $O(1)$ rounds and $O(N)$ space that duplicates the data item such that the total size over all duplications is N .*

Another important application of sorting is to duplicate the input data small number of times and rearrange them. This allows us to run multiple tasks which requires the same input data in parallel.

Theorem 6 (See e.g., [ASS⁺18]). *Given size N input data and a parameter $m = N^{o(1)}$, there is a fully scalable algorithm which duplicates the input data m times. The algorithm takes $O(1)$ rounds and $O(N \cdot m)$ total space.*

The following theorem gives a subroutine of applying dimensionality reduction efficiently in the MPC model. Note that each input data point is not necessarily small enough to fit into the memory of a single machine and the entries of each data point can distributed arbitrarily over machines.

Theorem 7 (Theorem 3 of [AAH⁺23]). *Given a set of points $x_1, x_2, \dots, x_n \in \mathbb{R}^d$ whose entries are distributed arbitrarily over the machines in the MPC model, there is a fully scalable MPC algorithm which outputs $x'_1, x'_2, \dots, x'_n \in \mathbb{R}^{d'}$ where $d' = O(\log n)$ and with probability at least $1 - 1/\text{poly}(n)$, $\forall i, j \in [n], 0.99\|x_i - x_j\|_2^2 \leq \|x'_i - x'_j\|_2^2 \leq 1.01\|x_i - x_j\|_2^2$. In addition, $\forall i \in [n], x'_i$ is held entirely by a single machine. The algorithm takes $O(1)$ rounds, and the total space needed is $O(nd + n \log^3 n)$.*

The following theorem states that there is an efficient MPC algorithm which computes a constant approximate Euclidean spanner. Note that a graph can be represented by a set of edges, and we store the edges in an arbitrarily distributed way over machines in the MPC model.

Theorem 8 (See e.g., [EMMZ22]). *Given parameters $C > 1, t > 0$ and a set of points $X = \{x_1, x_2, \dots, x_n\} \subset \mathbb{R}^d$ distributed arbitrarily on several machines in the MPC model with $d = O(\log n)$, there is a fully scalable MPC algorithm which outputs a graph $G = (X, E)$ such that with probability at least $1 - 1/\text{poly}(n)$:*

1. *For any $x, y \in X$, if there is an edge in G between x and y , $\|x - y\|_2 \leq C \cdot t$.*
2. *For any $x, y \in X$ with $\|x - y\|_2 \leq t$, there is a path in G connecting x, y using at most 2 edges.*

In addition, G has at most $n^{1+1/C^2+o(1)}$ number of edges. The algorithm uses $O(1)$ rounds and needs total space $n^{1+1/C^2+o(1)}$.

Another algorithmic primitive in the MPC model is to find predecessors. In particular, given pairs $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ where x_i have a total ordering and $y_i \in \{0, 1\}$. The total size of all pairs is N and each x_i has size $N^{o(1)}$. The goal is to output $(x_1, i_1), (x_2, i_2), \dots, (x_n, i_n)$ where $\forall j \in [n], i_j$ satisfies that (1) $y_{i_j} = 1$, (2) x_{i_j} is the largest one such that $x_{i_j} \leq x_j$.

⁷Here the “fully scalable” means that the algorithm works for local memory size N^δ for any constant $\delta > 0$.

Theorem 9 ([GSZ11, ASS⁺18]). *There is a fully scalable MPC algorithm that solves the predecessor problem in $O(1)$ rounds and $O(N)$ total space.*

We also state an MPC algorithmic primitive to index elements in sets.

Theorem 10 ([ASS⁺18]). *Given sets S_1, S_2, \dots, S_k of items with total size N , there is a fully scalable MPC algorithm which runs in $O(1)$ rounds and uses total space $O(N)$ to output the size of each set S_i , and for each $x \in S_i$, outputs its ranking / index within the set S_i .*

Similarly, the above primitive can be easily extended to compute a prefix sum for each element in each set.

Theorem 11 ([ASS⁺18]). *Given sets S_1, S_2, \dots, S_k of ranked items with total size N , where each item has a weight, there is a fully scalable MPC algorithm which runs in $O(1)$ rounds and uses total space $O(N)$ to output a prefix sum $\sum_{j \in S_l: \text{item } j \text{ has smaller rank than item } i} (\text{weight of } j)$ for each item $i \in S_l$ for each $S_l, l \in [k]$.*

4.2 Implementation of Euclidean MST in the MPC Model

In this section, we provide details of the MPC implementation of our Euclidean MST algorithm.

Theorem 12 (Constant approximate Euclidean MST in the MPC model). *Given n points from \mathbb{R}^d , there is a fully scalable MPC algorithm which outputs an $O(1)$ -approximate MST with probability at least 0.99. The number of rounds of the algorithm is $O(\log \log(n) \cdot \log \log \log(n))$. The total space required is at most $O(nd + n^{1+\varepsilon})$ where $\varepsilon > 0$ is an arbitrary small constant.*

Proof. The approximation guarantee is shown by Lemma 8 in Section 3. In the remaining of the proof, we show how to implement the algorithm described in Section 2 in the MPC model.

Preprocessing. We firstly need to run dimensionality reduction. By applying the algorithm mentioned in Theorem 7, we are able to reduce the dimension of each point to $d' = \Theta(\log n)$ while preserve all pairwise distances up to a constant factor with probability at least 0.999. This step requires $O(1)$ rounds and $O(nd + n \log^3 n)$ total space. In addition, the algorithm is fully scalable. Then we can use sorting (Theorem 3) to sort all values in each dimension separately and parallelly to learn the maximum and the minimum values in each dimension. Then according to Theorem 5, we are able to use $O(1)$ rounds to make every machine learn the maximum and the minimum value of each dimension. Thus, we are able to shift the point set such that the minimum value in each dimension is 0. Note that base on the maximum value and minimum value of each coordinate, we can verify whether all points are at the same location. If all points are at the same location, the MST cost is 0 and we can output an arbitrary spanning tree. Otherwise we linearly scale the entire dataset such that the maximum value of the point set $\Delta = \Theta(n^2 / \log n)$.

Next, we need to build the randomly shifted grid and round the coordinates of each point. To do this, we let one machine generate a randomly shifted vector $(a_1, a_2, \dots, a_{d'}) \in \mathbb{R}^{d'}$ where each coordinate is drawn uniformly from $[0, \Delta)$ (this is the same vector as described in Section 2). Since $d' = \Theta(\log n)$, we apply Theorem 5 again, we can use $O(1)$ rounds to make every machine learns

$(a_1, a_2, \dots, a_{d'})$. Therefore, all machines agree on the same randomly shifted grid. Each machine checks each point stored in its memory, and moves the point to the closest center of a cell at level $t = \alpha^{100}$. Note that each movement might change the cost of MST by $\sqrt{d'} \cdot \alpha^{100} = \text{polylog } n$, and because we have n points, the total change will be at most $n \text{polylog } n$. However, since the MST cost is at least the maximum distance which is at least $\Delta = \Theta(n^2 / \log n)$, the MST cost is changed by at most a constant factor by above movements. After doing above steps, our point set satisfies that the maximum distance is at most n^2 and the minimum distance is at least α^{100} which are the same as mentioned in Section 2.

Now we describe how to generate \tilde{G}_t for $t = 1, 2, 4, \dots, \alpha^H$. Firstly, we apply Theorem 6 to duplicate input data $\log(\alpha^H)$ times. This step is fully scalable and takes $O(1)$ rounds and $O(n \log n) \cdot \log(\alpha^H)$ total space. Since we duplicated the input data, we can handle each $t = 1, 2, 4, \dots, \alpha^H$ in parallel. Consider a fixed t and let k be such $\alpha^{k-1} < t \leq \alpha^k$. As discussed above, all machines agree on the same randomly shifted grid, thus each cell at level α^{k+1}/β can be uniquely determined by a vector in $\mathbb{Z}^{d'}$ (see the discussion in Section 2). Let $c(x)$ be the cell at level α^{k+1}/β containing x . We create $(c(x), x)$ for each point x . These steps only require local computations. Let $\varepsilon > 0$ be an arbitrary constant. Then we sort all $(c(x), x)$ (Theorem 3) but we restrict that each machine only uses memory $\Theta(s^{1/(1+\varepsilon/2)})$ where s is the local memory size of each machine. The algorithm is fully scalable, takes $O(1)$ rounds and $O(n^{1+\varepsilon/2} \log n)$ total space. Then the points X_c in the same cell c are rearranged into consecutive machines. We set C used in Theorem 8 to be $\xi/\sqrt{\varepsilon}$ where $\xi \geq 1$ is a sufficiently large constant. If X_c is entirely contained by a single machine, then $\tilde{G}_t(X_c)$ can be computed locally and with probability at least $1 - 1/\text{poly}(n)$,

1. For any $p, q \in X_c$, if there is an edge in $\tilde{G}_t'(X_c)$ between p and q , $\|p - q\|_2 \leq C \cdot t$.
2. For any $p, q \in X_c$ with $\|p - q\|_2 \leq t$, there is a path in $\tilde{G}_t'(X_c)$ connecting p, q using at most 2 edges.

Since $\tilde{G}_t'(X_c)$ only has $O(|X_c|^{1+\varepsilon/2})$ edges, $\tilde{G}_t'(X_c)$ can still be stored locally (note each machine only consumes $s^{1+\varepsilon/2}$ space during the sorting). Otherwise, if X_c is distributed over multiple machines, then only the first machine containing points from X_c and the last machine containing points from X_c might contain points that are not in X_c . Then there are two cases: (1) If X_c is only on two consecutive machines, we can send them into an arbitrary machine and handle them locally as discussed above. (2) If X_c is on more than two consecutive machines, the first machine containing points in X_c can move them to the second machine containing points in X_c , and the last machine containing points in X_c can move them to the second to the last machine that contains points in X_c . By applying the predecessor algorithm (Theorem 9), if a machine holds points in X_c , it learns the first machine contains points in X_c as well as the last machine that contains points in X_c . Finding predecessor processor requires $O(1)$ rounds and $O(n \log n)$ space, and it is fully scalable in the MPC model. Then we are able to apply Theorem 8 for each X_c in parallel. Therefore, we obtain graphs $\tilde{G}_t'(X_c)$ for all X_c in $O(1)$ rounds and the total space required is $O(n^{1+\varepsilon/2})$. This step is also fully scalable. In addition, with probability at least $1 - 1/\text{poly}(n)$, for every X_c , $G_t'(X_c)$ satisfies the properties mentioned in Theorem 8 with respect to X_c . Then we let G_t' be the union of all $G_t'(X_c)$. This step can be done by only local computation, i.e., renaming the graph.

Then, we duplicate $\tilde{G}_1', \tilde{G}_2', \tilde{G}_4', \dots, \tilde{G}_{\alpha^H}'$ $\log(\alpha^H)$ times. According to Theorem 6, this is fully scalable and only takes $O(1)$ rounds and $O(n^{1+\varepsilon/2} \log n)$ total space. For each edge of the i -th copy

of \tilde{G}'_t , if $t \leq 2^i$, we add this edge into \tilde{G}_{2^i} . Otherwise, we delete the edge. Therefore, we obtain that $\tilde{G}_t = \bigcup_{t' \leq t} \tilde{G}'_{t'}$, and thus $\tilde{G}_1 \subseteq \tilde{G}_2 \subseteq \tilde{G}_4 \subseteq \dots \subseteq \tilde{G}_{\alpha^H}$. The above step can be done by renaming the edges which only requires local computation. Finally, we apply sorting (Theorem 3) again to remove duplicated edges in each \tilde{G}_t .

To conclude, the entire preprocessing stage takes $O(1)$ rounds and uses $O(nd + n^{1+\varepsilon/2} \log n)$ total space. In addition, these steps are fully scalable in the MPC model.

Leader compression. Now we describe how to implement leader compression algorithm (Algorithm 1) in the MPC model. We firstly describe the format of the input data of leader compression stored in the system. Note that the partition \mathcal{P} is implied by the leader mapping function ℓ : $\forall x, y \in X$, x and y are in the same part if and only if $\ell(x) = \ell(y)$. Therefore, we store ℓ in the system to denote \mathcal{P} . In particular we have $|X|$ tuples $(x, \ell(x))$ distributed on the machines. The graph H is represented a set of edges, and the edges (x, y) are also distributed on the machines. The boolean value Edges is known by all machines (this can be done by Theorem 5 for broadcasting).

We firstly need to compute the set of leaders S . This step can be done by looking at each tuple $(x, \ell(x))$. If $\ell(x) = x$, we generate the random bit b_x for the leader node x (e.g., create a tuple $(\text{"b"}, x, b_x)$). Above steps only require local computation. The next step is to let $b_x \leftarrow b_{\ell(x)}$ for every $x \in X$. We can regard it as a simultaneous access problem in the PRAM model. In particular, the vector b is stored in the shared memory. Each processor corresponds to a tuple $(x, \ell(x))$. If $x \neq \ell(x)$, then the processor reads the value of $b_{\ell(x)}$ and writes the value of b_x . This operation has $O(1)$ depth and $O(|X|)$ work in the PRAM model. Therefore, according to the simulation algorithm (Theorem 4), we can use $O(1)$ rounds and total space $O(n)$ to compute b_x for all $x \in X$. In addition, this step is fully scalable. Then the next step is to find all edges $(x, y) \in H$ such that $b_x = 1$ and $b_y = 0$, and for each such edge, create a tuple $(\ell(y), (x, y, \ell(x)))$. This step can also be regarded as a simultaneous access problem in the PRAM model. In particular, the vector b and the vector $\ell(\cdot)$ is stored in the shared memory. Each processor corresponds to each edge $(x, y) \in H$. The process that all processors simultaneously read corresponding $b_x, b_y, \ell(x), \ell(y)$ has $O(1)$ depth and $O(|H| + |X|)$ work. Therefore, by applying the simulation algorithm again (Theorem 4), we can use $O(1)$ rounds and $O(|H| + |X|)$ total space in the MPC model to create all tuples $(\ell(y), (x, y, \ell(x)))$ for edges $(x, y) \in H$ whose $b_x = 1, b_y = 0$. This step is also fully scalable. Then for each $z = \ell(y)$, if there are some tuples $(z, (\cdot, \cdot, \cdot))$ created, we only keep an arbitrary one $(\ell(y), (x, y, \ell(x)))$. This deduplication step can be done by sorting (Theorem 3). Thus it is fully scalable, takes $O(1)$ rounds and uses total space at most $O(|H|)$. For every tuple $(\ell(y), (x, y, \ell(x)))$ that is kept, if Edges is true, we add edge (x, y) into the output edge set E . This operation can be done by local computation only (e.g., create a tuple $(\text{"E"}, (x, y))$). Finally, we need to update $\ell(y)$ for all $y \in X$. In particular, if there is a tuple $(\ell(y), (x, y, \ell(x)))$ kept, we need to update $\ell(y) \leftarrow \ell(x)$. Otherwise we keep $\ell(y)$ unchanged. This operation can also be seen as a simultaneous access in the PRAM model. In particular, we have a size $|X|$ array stored in the shared memory, if there is a tuple $(\ell(y), (x, y, \ell(x)))$ kept, then the entry of the array with index $\ell(y)$ has value $(x, y, \ell(x))$. Otherwise, the entry of the array with index $\ell(y)$ has value empty. Each processor corresponds to each tuple $(y, \ell(y))$. Each processor simultaneously accesses the array and see whether the entry with index $\ell(y)$ is empty or not. If the entry is empty, the processor keeps $(y, \ell(y))$ unchanged. Otherwise the entry is $(x, y, \ell(x))$, the processor updates $(y, \ell(y))$ to $(y, \ell(x))$. These PRAM operations take $O(1)$ depth and $O(|X|)$ total work. Therefore, according to the simulation theorem again (Theorem 4), the above steps can also

be done using $O(1)$ rounds and $O(|X|)$ total space, and the algorithm is fully scalable in the MPC model.

To conclude, the one leader compression step can be implemented in the MPC model using $O(1)$ rounds and $O(|X| + |H|)$ total space. In addition, the algorithm is fully scalable.

Part 1 of the algorithm. We consider how to implement Algorithm 2 in the MPC model. Notice that we will run Algorithm 2 for $t = \alpha, \alpha^2, \alpha^3, \dots, \alpha^H$ simutenously in parallel. Firstly, we can apply Theorem 6 to duplicate the point set H times. This operation is fully scalable and takes $O(1)$ rounds and $O(H \cdot n \log n)$ total space. Note that the graphs $\tilde{G}_\alpha, \tilde{G}_{\alpha^2}, \tilde{G}_{\alpha^3}, \dots, \tilde{G}_{\alpha^H}$ are already computed and stored distributedly in the system. We can use sorting (Theorem 3) to send each \tilde{G}_t and a duplicated copy of the input point set to a group of machines. Thus we can handle $t = \alpha, \alpha^2, \alpha^3, \dots, \alpha^H$ simutenously in parallel. This operation is fully scalable and takes $O(1)$ rounds and $O(n^{1+\varepsilon/2} \log n)$ total space.

Next, we focus on how to implement Algorithm 2 for a fixed t . The first step is to generate partition $\mathcal{P}_t^{(0)}$. As we discussed before, we only need to generate the leader mapping function $\ell_t^{(0)}(\cdot)$ such that x and y are in the same part of $\mathcal{P}_t^{(0)}$ if and only if $\ell_t^{(0)}(x) = \ell_t^{(0)}(y)$. Recall the definition $\mathcal{P}_t^{(0)}$, two points x, y are in the same part if and only if they are in the same cell in the grid of level t/β . As discussed in the preprocessing step, all machines agree on the same randomly shifted grid, thus each cell at level t/β can be uniquely determined by a vector in $\mathbb{Z}^{d'}$ (see the discussion in Section 2). Let $c(x) \in \mathbb{Z}^{d'}$ denote the cell which contains point x in level t/β . Then we can compute $(c(x), x)$ for every point $x \in X$. This step only requires local computation. We can sort (Theorem 3) all $(c(x), x)$ such that each point x learns whether it is the first point in the cell $c(x)$. For each $(c(x), x)$, we create a tuple $((c(x), x), 1)$ if x is the first point in the cell $c(x)$ and set $\ell_t^{(0)}(x) = x$. Otherwise we create a tuple $((c(x), x), 0)$. Then, by applying the predecessor algorithm (Theorem 9), each point x learns the index of the point y which is the first point within the same cell $c(x) = c(y)$. By using the index of the point y and simulating the simultaneous access operation in PRAM (Theorem 4), each point x learns y which is the first point in the same cell. Then we set $\ell_t^{(0)}(x) = y$. The above operations only require $O(1)$ rounds and $O(n \log n)$ total space.

Once we obtained $\ell_t^{(0)}$, we are able to run $h = O(\log \log n)$ rounds of leader compression process. As we discussed previously, running h rounds of leader compression takes $O(\log \log n)$ rounds and $O(|X| + |\tilde{G}_t|)$ space. The algorithm is fully scalable as well. For each $(x, y) \in \tilde{G}_t$, we want to access $\ell_t^{(h)}(x)$ and $\ell_t^{(h)}(y)$, this again can be regarded as a simultaneous access problem in the PRAM model where $\ell_t^{(h)}$ is stored in the shared memory and each processor corresponds to each edge $(x, y) \in \tilde{G}_t$. Each processor requires to read $\ell_t^{(h)}(x)$ and $\ell_t^{(h)}(y)$, and if they are not equal, the processor writes INCOMPLETE to the $\ell_t^{(h)}(x)$ -th entry and the $\ell_t^{(h)}(y)$ -th entry of an array in the shared memory. According to the simulation theorem (Theorem 4), these steps can be done in $O(1)$ rounds in the MPC model and $O(|X| + |\tilde{G}_t|)$ space, and it is fully scalable. Then for each point x , we can check whether it is a leader, i.e., $\ell_t^{(h)}(x) = x$ and whether it is marked as INCOMPLETE. This operation can be done by simulating simultaneous access in PRAM (Theorem 4) which takes $O(1)$ rounds and $O(n)$ total space, and is fully scalable. Let I be the set of all incomplete leaders. Then for each $x \in I$, we can locally compute $(c(x), x)$ where $c(x) \in \mathbb{Z}^{d'}$ indicates the cell at the level α^{k+1}/β containing x . Then we can sort all such $(c(x), x)$, and each point in I learns whether it is the

first point among I and in the cell $c(x)$. By following the similar approach of computing $\ell_t^{(0)}$, we compute $\ell_t^{(h+1)}(x) = y$ for each $x \in I$ where $y \in I$ is the first point in the cell $c(x)$. Similar as computing $\ell_t^{(0)}$, this step takes $O(1)$ MPC rounds and requires $O(n \log n)$ total space, and it is fully scalable. For each leader which is not INCOMPLETE, i.e., $x \notin I$, we set $\ell_t^{(h+1)}(x) = x$. This can be done by simulating simultaneous access in PRAM (Theorem 4). Thus, it is fully scalable, has $O(1)$ rounds, and uses space $O(n)$. Finally, for each tuple $(x, \ell_t^{(h)}(x))$, we access $\ell_t^{(h+1)}(\ell_t^{(h)}(x))$ and create a new tuple $(x, \ell_t^{(h+1)}(\ell_t^{(h)}(x)))$ to indicate $\ell_t^{(h+1)}(x) \leftarrow \ell_t^{(h+1)}(\ell_t^{(h)}(x))$. This can also be done by simulating simultaneous access in PRAM (Theorem 4). Thus, it is fully scalable, has $O(1)$ rounds, and uses space $O(n)$. The output $\hat{\mathcal{P}}_t$ is represented by $\ell_t(\cdot) \equiv \ell_t^{(h+1)}(\cdot)$.

To conclude, Part 1 of the algorithm (Algorithm 2) is fully scalable. It requires $O(\log \log n)$ rounds and requires total space $O(n^{1+\varepsilon/2} \log n)$.

Part 2 of the algorithm. We consider how to implement Algorithm 3 for all $t = 1, 2, 4, 8, \dots, \alpha^H$. Firstly, we can duplicate our input data $\log(\alpha^H)$ times (Theorem 6) which takes $O(1)$ rounds and $O(n \log n \cdot \log(\alpha^H))$ total space, and the algorithm is fully scalable. Then we can use disjoint set of machines to handle each t separately. Note that \tilde{G}_t are already stored in the system due to the preprocessing stage. We also send \tilde{G}_t to the group of machines which will be used to compute $\hat{\mathcal{P}}_t$. This is also fully scalable and takes only $O(1)$ rounds and total space $O(n^{1+\varepsilon/2} \log n)$.

As described in Section 2, instead of computing all t at the same time, we handle t in decreasing order of $v_2(\log t)$ (recall Definition 2 for $v_2(\cdot)$). In particular, we already obtained $\hat{\mathcal{P}}_t$ for all $t = 1, \alpha, \alpha^2, \dots, \alpha^H$ by running Part 1 of the algorithm. More precisely, we obtained $\ell_t(\cdot)$ for these t . In the first iteration, $\kappa = \sqrt{\alpha}$, we will handle $t = \sqrt{\alpha}, \alpha^{1.5}, \dots, \alpha^{H-0.5}$ at the same time in parallel, i.e., we will have $\hat{\mathcal{P}}_t$ for all $t = 1, \alpha^{0.5}, \alpha, \alpha^{1.5}, \dots, \alpha^H$ at the end of the first iteration. In the second iteration, $\kappa = \alpha^{0.25}$, we will have $\hat{\mathcal{P}}_t$ for all $t = 1, \alpha^{0.25}, \alpha^{0.5}, \alpha, \dots, \alpha^H$ at the end of the second iteration. In the i -th iteration we will have $\kappa = \alpha^{1/2^i}$, and we will have $\hat{\mathcal{P}}_t$ for all $t = \kappa^0, \kappa, \kappa^2, \dots, \alpha^H$ at the end of the i -th iteration. Thus, we will have $O(\log \log(\alpha)) = O(\log \log \log(n))$ iterations in total. Note that if $\hat{\mathcal{P}}_t$ is computed in the i -th iteration, then $\forall j \geq 1$, $\hat{\mathcal{P}}_t$ (more precisely, $\ell_t(\cdot)$) will be the input for computing $\hat{\mathcal{P}}_{t \cdot \alpha^{1/2^{i+j}}}$ and $\hat{\mathcal{P}}_{t/\alpha^{1/2^{i+j}}}$. Therefore, when all $\ell_t(\cdot)$ are computed after the i -th iteration, we can use duplication process (Theorem 6) to duplicate all $\ell_t(\cdot)$ at most $O(\log \log \log(n))$ times and send each duplicated copy of $\hat{\mathcal{P}}_t$ (more precisely, $\ell_t(\cdot)$) to the groups of machines that will be used to compute $\hat{\mathcal{P}}_{t \cdot \alpha^{1/2^j}}$ or $\hat{\mathcal{P}}_{t/\alpha^{1/2^j}}$ for $j \geq 1$.

When a group of machines receives all required inputs $X, \tilde{G}_t, \ell_{t/\kappa}(\cdot), \ell_{t \cdot \kappa}(\cdot)$, it starts to run Algorithm 3 to compute $\hat{\mathcal{P}}_t$ (i.e., $\ell_t(\cdot)$). We firstly initialize $\ell_t^{(0)}(\cdot)$ to be $\ell_{t/\kappa}(\cdot)$. Then, similar as Part 1 of the algorithm, we run $O(\log \log n)$ rounds of leader compression. This step takes $O(\log \log n)$ MPC rounds and $O(|X| + |\tilde{G}_t|)$ space. This step is fully scalable. Then by using the same process described in how to implement Part 1 of the algorithm, we can simulate PRAM operations (Theorem 4) to find all $x \in X$ that are marked as INCOMPLETE. For each leader x ($\ell_t^{(h)}(x) = x$) which is marked as INCOMPLETE, let $\ell_t^{(h+1)}(x) = \ell_{t \cdot \kappa}(x)$. For each leader x ($\ell_t^{(h)}(x) = x$) which is not marked as INCOMPLETE, let $\ell_t^{(h+1)}(x) = x$. Finally, for each $x \in X$, let $\ell_t^{(h+1)}(x) \leftarrow \ell_t^{(h+1)}(\ell_t^{(h)}(x))$. Similar as before, above operations used to compute $\ell_t^{(h+1)}(\cdot)$ can be seen as simultaneous accesses in the PRAM model which requires $O(1)$ depth and $O(n)$ total work. By applying PRAM simulation

again (Theorem 4), the above procedure to compute $\ell_t^{(h+1)}(\cdot)$ requires $O(1)$ MPC rounds and $O(n)$ total space, and it is fully scalable. Note that the output $\hat{\mathcal{P}}_t$ is represented by $\ell_t(\cdot) \equiv \ell_t^{(h+1)}(\cdot)$.

To conclude, to compute $\hat{\mathcal{P}}_t$ (more precisely, $\ell_t(\cdot)$) for all $t = 1, 2, 4, \dots, \alpha^H$, our algorithm takes $O(\log \log(n) \cdot \log \log \log(n))$ rounds. The total space required is at most $n^{1+\varepsilon/2} \cdot \text{polylog}(n)$. In addition, the algorithm is fully scalable.

Part 3 of the algorithm. Finally, let us consider how to implement Algorithm 4 in parallel for all t in the MPC model.

Similar as before, we use a disjoint group of machines to compute $\text{PART3}(X, t, \tilde{G}_t, \hat{\mathcal{P}}_{t/2}, \hat{\mathcal{P}}_t)$ for each t . To do this, we need to duplicate $X \log(\alpha^H)$ times and duplicate $\hat{\mathcal{P}}_t$ (i.e., $\ell_t(\cdot)$) 2 times ($\hat{\mathcal{P}}_t$ is used to compute $\text{PART3}(X, t, \tilde{G}_t, \hat{\mathcal{P}}_{t/2}, \hat{\mathcal{P}}_t)$ and $\text{PART3}(X, 2t, \tilde{G}_{2t}, \hat{\mathcal{P}}_t, \hat{\mathcal{P}}_{2t})$). Then we send each copy of the data and \tilde{G}_t to the corresponding group of machines. According to Theorem 6, above process can be done using $O(1)$ rounds and $n^{1+\varepsilon/2} \text{polylog}(n)$ space, and these steps are fully scalable.

In the following, we describe how to implement Algorithm 4 for any fixed t . We initialize $\ell_t^{(0)}(\cdot)$ to be $\ell_{t/2}(\cdot)$. This only requires local computation. Then, as we discussed previously, $O(\log \log n)$ rounds of leader compression can be implemented in $O(\log \log n)$ MPC rounds and $O(n^{1+\varepsilon/2} \log(n))$ total space, and the process is fully scalable. The i -th round of leader compression outputs a set of edges $E_t^{(i)}$. At the end of the last leader compression process, let the obtained leader mapping be $\ell_t^{(h)}(\cdot)$. For each $x \in X$, we check whether $\ell_t^{(h)}(x) = \ell_t(x)$, if not, we create a set $S_{\ell_t(x)}$ and add $\ell_t^{(h)}(x)$ into the set $S_{\ell_t(x)}$. Note that this operation can be regarded as simultaneous access of $\ell_t^{(h)}(\cdot)$ and $\ell_t(\cdot)$ under the PRAM model, which only has $O(1)$ depth and $O(n)$ work. Therefore, according to the simulation theorem (Theorem 4), the computation of all $S_{\ell_t(x)}$ only requires $O(1)$ MPC rounds and $O(n)$ total space, and it is fully scalable. Then, for each set S_c if z is in S_c , we add an edge $\{z, c\}$ into F_t . Thus F_t will be a forest of stars. Then we output $F_t \cup \bigcup_{i=1}^h E_t^{(i)}$.

To conclude, to compute $F_t \cup \bigcup_{i=1}^h E_t^{(i)}$ for all $t = 1, 2, 4, \dots, \alpha^H$, our algorithm takes $O(\log \log(n))$ rounds. The total space required is at most $n^{1+\varepsilon/2} \cdot \text{polylog}(n)$. In addition, the algorithm is fully scalable.

Put them together. All parts of our algorithm are fully scalable. Thus, the overall algorithm is fully scalable as well. The preprocessing stage takes $O(1)$ rounds and $O(nd + n^{1+\varepsilon/2} \log n)$ total space. Running Algorithm 2 for all $t = 1, \alpha, \alpha^2, \dots, \alpha^H$ takes $O(\log \log(n))$ rounds and $O(nd + n^{1+\varepsilon/2} \log n)$ total space. Running Algorithm 3 for all remaining t takes $O(\log \log(n) \cdot \log \log \log(n))$ rounds and $n^{1+\varepsilon/2} \cdot \text{polylog } n$ total space. Running Algorithm 4 for all t takes $O(\log \log(n))$ rounds and $n^{1+\varepsilon/2} \cdot \text{polylog } n$ total space. Therefore, the entire algorithm takes $O(\log \log(n) \log \log \log(n))$ rounds and $O(nd + n^{1+\varepsilon})$ total space. \square

5 Euler Tour of Approximate Euclidean MST in MPC

In this section, we show how to compute an Euler tour of the approximate MST that we obtained. Note that the actual diameter of the tree that we obtained can be very large (for example, the tree that we obtained can be a path which has diameter $\Theta(n)$). Therefore, we cannot use the algorithm of [ASS⁺18] which requires $\Omega(\log(\text{diameter}))$ rounds. We propose a new method. In particular, given the approximate MST and the hierarchy of clusters that we used to obtain the tree, we are able to use linear total space and $O(1)$ MPC rounds to compute an Euler tour of the tree, and our algorithm is fully scalable.

5.1 Additional Existing Algorithmic Primitives in the MPC Model

A standard way to store a sequence $A = (a_1, a_2, \dots, a_n)$ in the MPC model is that we store tuples $(1, a_1), (2, a_2), (3, a_3), \dots, (n, a_n)$ arbitrarily on the machines in a distributed manner (see e.g., [ASS⁺18]). In this way, one can easily use sorting (Theorem 3) to reorder the elements such that the sequence is stored in consecutive machines and each machine stores corresponding consecutive elements in order. Similarly, a standard way to store a mapping $f : [n] \rightarrow [m]$ in the MPC model is that we store tuples $(1, f(1)), (2, f(2)), \dots, (n, f(n))$.

The sequence insertion problem is stated as the following: Given $k+1$ sequences $A = (a_1, a_2, \dots, a_m)$, A_1, A_2, \dots, A_k and a mapping $f : [k] \rightarrow \{0\} \cup [m]$, the goal is to output a sequence A' which is obtained by inserting every A_i , ($i \in [k]$) into A such that A_i is between the element $a_{f(i)}$ and $a_{f(i)+1}$. Let the total size of the input sequences is $|A| + |A_1| + |A_2| + \dots + |A_k| = N$. The size of A' is also N , and the space to store the input mapping f is $O(k) = O(N)$ as well.

Theorem 13 ([ASS⁺18]). *There is a fully scalable MPC algorithm solving the sequence insertion problem in $O(1)$ rounds and $O(N)$ space.*

Theorem 14 ([ASS⁺18]). *Given a tree of n nodes with diameter at most Λ , there is a fully scalable which takes $O(\log(\Lambda))$ rounds and $O(n^{1+\varepsilon})$ total space to output an Euler tour (see Definition 6) of the tree, where $\varepsilon > 0$ is an arbitrarily small constant.*

5.2 Join Euler Tours of Spanning Tree of Sub-clusters

Given a tree, we treat each edge $\{x, y\}$ as two directed edges (x, y) and (y, x) . An Euler tour of a tree is a traverse sequence of tree edges such that each directed edge (x, y) appeared in the sequence exactly once. The following gives a formal definition of an Euler tour of a tree.

Definition 6 (Euler tour of a tree). Given a tree $T = (V, E)$ with $|V| = n$ nodes, the Euler tour of T is a sequence $((v_1, v_2), (v_2, v_3), \dots, (v_{2n-2}, v_1))$ of length $2n - 2$ where $\forall i \in [2n - 2], \{v_i, v_{(i \bmod (2n-2))+1}\} \in E$, and $\forall \{x, y\} \in E$, there exists exactly one $i \in [2n - 2]$ and exactly one $i' \in [2n - 2]$ such that $(v_i, v_{(i \bmod (2n-2))+1}) = (x, y)$ and $(v_{i'}, v_{(i' \bmod (2n-2))+1}) = (y, x)$.

5.2.1 Properties of An Euler Tour of a Tree

In this section, we give some formal definitions for Euler tour of a tree. Most observations stated in this section can also be found in the textbooks of algorithms such as [CLRS22].

Definition 7 (Parent pointers). Given a rooted tree T over V , we use $\text{par} : V \rightarrow V$ to denote a set of parent pointers which is a mapping from each node to its parent, and we set $\text{par}(v) = v$ for the root v of T .

Note that $\text{par}(\cdot)$ contains the information of all tree edges and the root that we select.

Definition 8 (Size of subtree). Given parent pointers $\text{par}(\cdot)$ of tree with root v , let $\text{size}_{\text{par}}(u)$ denote the number of nodes in the subtree rooted at u . If $\text{par}(\cdot)$ is clear in the context, we omit the subscript par and use $\text{size}(u)$ for short.

Definition 9 (Ordering of children). Given a rooted tree represented by parent pointers par , and an ordering of children of each node, we use $\text{child}_{\text{par}}(v)$ to denote the set of children of v , $\text{child}_{\text{par}}(v, i)$, ($i \in [|\text{child}_{\text{par}}(v)|]$) to denote the i -th child of v , and $\text{rank}_{\text{par}}(v)$ to denote its rank among its siblings, i.e., $\text{rank}_{\text{par}}(v)$ satisfies $\text{child}_{\text{par}}(\text{par}(v), \text{rank}_{\text{par}}(v)) = v$ if v is not a root and $\text{rank}_{\text{par}}(v) = 1$ otherwise.

If par is clear in the context, we use $\text{child}(\cdot)$ and $\text{rank}(\cdot)$ without subscript par for short.

Observation 1 (Euler tour for different root). Let $A = ((v_1, v_2), (v_2, v_3), (v_3, v_4), \dots, (v_{2n-2}, v_1))$ be any Euler tour of a tree T , then any circular shift of A , $((v_j, v_{j+1}), (v_{j+1}, v_{j+2}), \dots, (v_{2n-2}, v_1), (v_1, v_2), (v_2, v_3), \dots, (v_{j-1}, v_j))$, for $j \in [2n-2]$, is still an Euler tour of T .

Definition 10. If $A = ((v_1, v_2), (v_2, v_3), (v_3, v_4), \dots, (v_{2n-2}, v_1))$ is an Euler tour of a tree $T = (V, E)$ and $v_1 = u \in V$, then A is an Euler tour with respect to the root u .

Definition 11. Let $A = (e_1, e_2, \dots, e_{2n-2})$ be an Euler tour of a tree T where $e_i = (v_i, v_{i+1})$. Given any node v , we say that the first appearance of v is at $\text{first}_A(v) = i$ of A if i is the smallest value such that $v_i = v$, and we say that the last appearance of v is at $\text{last}_A(v) = i$ if i is the largest value such that $v_{i+1} = v$. If the Euler tour is clear in the context, we omit the subscript A of **first** and **last**.

Observation 2 (Parent via Euler tour). Let $A = ((v_1, v_2), (v_2, v_3), (v_3, v_4), \dots, (v_{2n-2}, v_1))$ be an Euler tour of a tree T with root v . If u is not a root, then $v_{\text{first}(u)-1}$ is the parent of u .

Observation 3 (Euler tour of subtree). Let $A = (e_1, e_2, \dots, e_{2n-2})$ be an Euler tour of a tree T with respect to the root u . Let v be any node in tree T , let i be the first appearance of v in A , and let j be the last appearance of v in A . Then $(e_i, e_{i+1}, \dots, e_j)$ is an Euler tour of the subtree rooted at v (note that $j = i - 1$ when v is a leaf and thus the tour is an empty sequence), and thus $(j - i + 1)/2 + 1$ is the size of the subtree rooted at v .

Lemma 9 (Euler tour from the ordering of children). Given a tree T rooted at v with parent pointers par , and given a mapping $\text{child}_{\text{par}}(\cdot, \cdot)$ (Definition 9), then we are able to construct an Euler tour A with respect to the root v such that for any node u in the tree, we have $\text{first}(\text{child}(u, 1)) <$

$\mathbf{first}(\text{child}(u, 2)) < \mathbf{first}(\text{child}(u, 3)) < \dots < \mathbf{first}(\text{child}(u, |\text{child}(u)|))$. In addition, for any non-root node u , let the path $((x_1, x_2), (x_2, x_3), \dots, (x_{k-1}, x_k))$ be the unique simple path from the root v to u (i.e., $x_1 = v, x_k = u$), then the position of the (directed) edge $(\text{par}(u), u)$ in A is:

$$p = \sum_{l=2}^k \left(1 + \sum_{i=1}^{\text{rank}(x_l)-1} \mathbf{size}(\text{child}(x_{l-1}, i)) \right),$$

and the position of the (directed) edge $(u, \text{par}(u))$ is $p + 2 \cdot (\mathbf{size}(x_k)) - 1$

Proof. The proof is by induction and construction. If T only has one node, then $\text{child}(v) = \emptyset$ and we just output A to be an empty sequence.

Suppose the statement is true for any tree whose depth is at most $l-1$, now we consider how to construct A for a tree with depth l . Suppose T has root v and $\forall i \in [|\text{child}(v)|], \text{child}(v, i) = v_i$. We construct A to be following:

$$(v, v_1), A_1, (v_1, v), (v, v_2), A_2, (v_2, v), \dots, (v, v_{|\text{child}(v)|}), A_{|\text{child}(v)|}, (v_{|\text{child}(v)|}, v),$$

where A_i is an Euler tour of the subtree rooted at v_i via our induction hypothesis. Then, if u is in a subtree of v_i , then we still have $\mathbf{first}(\text{child}(u, 1)) < \mathbf{first}(\text{child}(u, 2)) < \mathbf{first}(\text{child}(u, 3)) < \dots < \mathbf{first}(\text{child}(u, |\text{child}(u)|))$ by our induction hypothesis. Otherwise, if $u = v$, it is clear that $\mathbf{first}(v_1) < \mathbf{first}(v_2) < \mathbf{first}(v_3) < \dots < \mathbf{first}(v_{|\text{child}(v)|})$. In the remaining of the proof, we need to determine the position of an edge. If an edge is (v, v_i) , then its position is clearly $1 + \sum_{j=1}^{i-1} 2 \cdot \mathbf{size}(v_j)$. If an edge is (v_i, v) , then its position is $\sum_{j=1}^i 2 \cdot \mathbf{size}(v_j)$. Consider an edge $(\text{par}(u), u)$ is in the subtree rooted at v_i . Let $((x_1, x_2), (x_2, x_3), \dots, (x_{k-1}, x_k))$ be the path from v to u , i.e., $x_1 = v, x_2 = v_i, x_{k-1} = \text{par}(u), x_k = u$. Then by our induction hypothesis, the position of $(\text{par}(u), u)$ is

$$\begin{aligned} & 1 + \left(\sum_{j=1}^{\text{rank}(v_i)-1} 2 \cdot \mathbf{size}(v_j) \right) + \sum_{l=3}^k \left(1 + \sum_{j=1}^{\text{rank}(x_l)-1} 2 \cdot \mathbf{size}(\text{child}(x_{l-1}, j)) \right) \\ &= \sum_{l=2}^k \left(1 + \sum_{j=1}^{\text{rank}(x_l)-1} 2 \cdot \mathbf{size}(\text{child}(x_{l-1}, j)) \right). \end{aligned}$$

Similarly, the position of $(u, \text{par}(u))$ is

$$\begin{aligned} & 1 + \left(\sum_{j=1}^{\text{rank}(v_i)-1} 2 \cdot \mathbf{size}(v_j) \right) + \sum_{l=3}^k \left(1 + \sum_{j=1}^{\text{rank}(x_l)-1} 2 \cdot \mathbf{size}(\text{child}(x_{l-1}, j)) \right) + 2 \cdot \mathbf{size}(x_k) - 1 \\ &= \sum_{l=2}^k \left(1 + \sum_{j=1}^{\text{rank}(x_l)-1} 2 \cdot \mathbf{size}(\text{child}(x_{l-1}, j)) \right) + 2 \cdot \mathbf{size}(x_k) - 1. \end{aligned}$$

Therefore, we conclude the proof of the induction. \square

Observation 4 (Euler tour and a path from the root). Let $A = (e_1, e_2, \dots, e_{2n-2})$ be an Euler tour of a tree T with respect to the root u . Consider any node v . If i is the first appearance of v , and $B = ((x_1, x_2), (x_2, x_3), \dots, (x_{k-1}, x_k))$ is a simple directed path from u to v (i.e., $x_1 = u, x_k = v$), then we have:

1. $\forall j \in [k-1]$, (x_j, x_{j+1}) appears in (e_1, \dots, e_{i-1}) in A , and (x_{j+1}, x_j) appears in (e_i, \dots, e_{2n-2}) .
2. If $\{p, q\}$ is a tree edge but neither (p, q) nor (q, p) appears in the path B , then both (p, q) and (q, p) appear simultaneously in (e_1, \dots, e_{i-1}) or appear simultaneously in (e_i, \dots, e_{2n-2}) .

Corollary 5 (Total weight on the path via Euler tour). *Let $A = (e_1, e_2, \dots, e_{2n-2})$ be an Euler tour of a tree T with respect to the root u . Let w be a weight function such that every edge $\{p, q\}$ has a weight $w(\{p, q\})$ in the tree. Let w' be the weight function of directed edge such that $\forall \text{edge } \{p, q\}$, if $p = \text{par}(q)$, then $w'(p, q) = w(\{p, q\})$ and $w'(q, p) = -(w(\{p, q\}))$. Let $B = ((x_1, x_2), (x_2, x_3), \dots, (x_{k-1}, x_k))$ be a simple directed path from u to v , then the total edge weight on the path $\sum_{i=2}^k w(\{x_{i-1}, x_i\}) = \sum_{j=1}^{\text{first}(x_k)-1} w'(e_j)$, i.e., the prefix sum over A with weights w' .*

Proof. Due to Observation 4, for the prefix of A , if an edge is on the path, we will count its edge weight once, otherwise, the non-path edge weight is either cancelled since both directions appeared, or none of the directions is appeared so we did not count it. \square

5.2.2 MPC Algorithms for Trees via Euler Tour

Lemma 10 (Change root of the Euler tour). *Given an arbitrary Euler tour A of a tree T over n nodes V , and given any node $v \in V$, there is a fully scalable MPC algorithm which outputs an Euler tour A' of tree T with respect to the root v in $O(1)$ rounds and $O(n)$ total space.*

Proof. We use sorting (Theorem 3) to find the first appearance of v in $A = \{e_1, e_2, \dots, e_{2n-2}\}$, and send $i = \text{first}_A(v)$ to all machines by broadcasting algorithm (Theorem 5). These steps only take $O(1)$ rounds and $O(n)$ total space, and they are fully scalable. Then compute the circular shift $A' = (e_i, e_{i+1}, \dots, e_{2n-2}, e_1, \dots, e_{i-1})$ of A . Since every machine learns $\text{first}_A(v)$, computing the circular shift only requires local computation. According to Observation 1, the circular shift A' is also an Euler tour of T . Since A' starts from node v , A' is an Euler tour of T with respects to the root v . \square

Lemma 11 (Subtree sizes). *Given an arbitrary Euler tour A of a tree T over n nodes V with root v , there is a fully scalable MPC algorithm which outputs the size of each subtree u in $O(1)$ rounds and $O(n)$ total space.*

Proof. We suppose A is an Euler tour of T with respect to the root v . Otherwise, we can apply Lemma 10 to make A satisfy above condition, and the operation is fully scalable, takes $O(1)$ rounds and $O(n)$ space.

We use sorting (Theorem 3) to compute $\text{first}(u)$ and $\text{last}(u)$ for each $u \in V$. It is fully scalable and only takes $O(1)$ rounds and $O(n)$ total space. Then, for each u in parallel, we simultaneously query both $\text{first}(u)$ and $\text{last}(u)$, which can be done in the fully scalable setting by Theorem 4 in $O(1)$ rounds and $O(n)$ total space. According to Observation 3, $\text{size}(u) = (\text{last}(u) - \text{first}(u) + 1)/2 + 1$. \square

Lemma 12 (Euler tour in consistent with ordering of children). *Given a tree T rooted at v over n nodes V with parent pointers par , and given a mapping $\text{child}_{\text{par}}(\cdot, \cdot)$ (Definition 9), if we also have an arbitrary Euler tour A' of T , there is an MPC algorithm which constructs an Euler tour A with respect to the root v such that for any node u in the tree, it satisfies $\text{first}(\text{child}(u, 1)) < \text{first}(\text{child}(u, 2)) < \text{first}(\text{child}(u, 3)) < \dots < \text{first}(\text{child}(u, |\text{child}(u)|))$. In addition, for any non-root node u , let the path $((x_1, x_2), (x_2, x_3), \dots, (x_{k-1}, x_k))$ be the unique simple path from the root v to u (i.e., $x_1 = v, x_k = u$), then the position of the (directed) edge $(\text{par}(u), u)$ in A is:*

$$p = \sum_{l=2}^k \left(1 + \sum_{i=1}^{\text{rank}(x_l)-1} \text{size}(\text{child}(x_{l-1}, i)) \right),$$

and the position of the (directed) edge $(u, \text{par}(u))$ is $p + 2 \cdot (\text{size}(x_k)) - 1$. In addition, the algorithm is fully scalable, and it takes $O(1)$ rounds and $O(n)$ total space.

Proof. Note that the formulation of the position of each directed edge $(\text{par}(u), u), (u, \text{par}(u))$ is given, and Lemma 9 already showed that the resulting sequence is indeed an Euler tour with stated properties. Therefore, our goal is to provide an MPC algorithm to efficiently compute the position of each directed edge in the resulting Euler tour A .

We can assume A' is an Euler tour of T with respect to the root v . Otherwise, we apply Lemma 10 to make A' be with respect to the root v which takes $O(1)$ rounds, $O(n)$ total space in the fully scalable setting. We apply Lemma 11 on A' to compute the size of each subtree in $O(1)$ rounds, $O(n)$ total space in the fully scalable setting. Next, for each node u , we want to compute $w(\{u, \text{par}(u)\}) = 1 + \sum_{i=1}^{\text{rank}(u)-1} \text{size}(\text{child}(\text{par}(u), i))$ for each non-root u . To achieve this, for each mapping $y = \text{child}(x, i)$ in parallel, we create $\text{rank}(y) = i$, add y into a set S_x , and simultaneously access the value $\text{size}(y)$. The simultaneous access operation can be done by Theorem 4. Then we apply the subset prefix sum algorithm (Theorem 11) over $\{S_x\}$ and thus each non-root u learns the value $\sum_{i=1}^{\text{rank}(u)-1} \text{size}(\text{child}(\text{par}(u), i))$, and it sets the edge weight $w(\{u, \text{par}(u)\}) = 1 + \sum_{i=1}^{\text{rank}(u)-1} \text{size}(\text{child}(\text{par}(u), i))$. The overall number of rounds needed to compute $w(\cdot)$ is $O(1)$, and it requires $O(n)$ total space, and is fully scalable.

Finally, we want to compute the position of each directed edge $(\text{par}(u), u), (u, \text{par}(u))$ in the final sequence A . Note that the position of $(\text{par}(u), u)$ is $\sum_{x \in V: x \neq v \text{ and } x \text{ is on the path from the root } v \text{ to } u} w(x)$. Suppose $A' = (e'_1, e'_2, \dots, e'_{2n-2})$. If $e'_i = (v'_i, v'_{i+1})$ satisfies $v'_i = \text{par}(v'_{i+1})$, let $w'_i = w(\{v'_i, v'_{i+1}\})$, otherwise, let $w'_i = -w(\{v'_i, v'_{i+1}\})$. Then according to Corollary 5, the position of $(\text{par}(u), u)$ is a prefix sum $\sum_{i=1}^{\text{first}_{A'}(u)} w'_i$. Note that computing w' only requires simultaneous accesses which can be done in $O(1)$ rounds, $O(n)$ total space, and in fully scalable setting according to Theorem 4. Then computing prefix sum $\sum_{i=1}^j w'_i$ can be done using Theorem 11. It takes $O(1)$ rounds, $O(n)$ total space, and is fully scalable. Note that the position of $(u, \text{par}(u))$ is the position of $(\text{par}(u), u)$ plus $2 \cdot \text{size}(u) - 1$. Since $\text{size}(u)$ is already computed, all computations of positions $(\text{par}(u), u)$ can be done in parallel with an additional call of simultaneous access (Theorem 4).

Therefore, the overall algorithm only takes $O(1)$ rounds, $O(n)$ total space. The algorithm is fully scalable. \square

5.2.3 Euler Tour Join

Now consider the following problem. We define the problem of joining Euler tours as following:

- **Inputs:**

1. A partition $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ of a node set V , represented by a leader mapping $\ell : V \rightarrow V$ where two nodes u, v are in the same component iff $\ell(u) = \ell(v)$. Each component C_i is represented by its leader node, i.e., $u \in C_i$ and $\ell(u) = u$.
2. A spanning tree $T = (\mathcal{C}, E)$ over components C_1, C_2, \dots, C_k where each node in T corresponds to a component C_i . An Euler tour A of T is also given.
3. An edge mapping g , where $\forall e = \{C_i, C_j\} \in E$, $g(e) = \{x, y\} \subseteq V$ and $x \in C_i, y \in C_j$, we sometimes also abuse the notation and use directed version $g(C_i, C_j) = (x, y)$ and $g(C_j, C_i) = (y, x)$. We use g^{-1} to indicate the inverse mapping.
4. A spanning tree $T_i = (C_i, E_i)$ for each component $i \in [k]$, and an Euler tour A_i of T_i .

- **Outputs:**

1. A spanning tree $T' = \left(V, \{g(e) \mid e \in E\} \cup \bigcup_{i \in [k]} E_i \right)$ over V , and an Euler tour A' of T' .

In Algorithm 5, we show a novel MPC algorithm that solves the Euler tour join problem efficiently.

Lemma 13 (Correctness of Algorithm 5). *The output A' of Algorithm 5 is an Euler tour of the output T' , and T' is a spanning tree of V , where the edges of T' is the union of edges in T_1, T_2, \dots, T_k and edges $\{g(e) \mid e \in E\}$.*

Proof. Since each T_i is a spanning tree of C_i , and T is a spanning tree of $\{C_1, C_2, \dots, C_k\}$ and \hat{E} are inter-cluster edges with respect to T , we know that T' is a spanning tree of $\bigcup_{i \in [k]} C_i = V$, and the edges of T' is union of edges in T_1, T_2, \dots, T_k and $\hat{E} = \{g(e) \mid e \in E\}$.

Since \bar{A} is an Euler tour of T which means that both directions of each edge in \hat{E} must appear in \hat{A} , which means that every node in \hat{V} also appears in \hat{A} . According to our construction of A' , both directions of every edge in T' must appear in A' and $|A'| = |A| + \sum_{x \in \hat{V}} |A_x| = |A| + \sum_{i \in [k]} |A_i| = 2|V| - 2$ which implies that each direction of each edge in T' also appeared exactly once.

In the remaining of the proof, we only need to show that A' is indeed a cycle. Consider an arbitrary (u, v) appeared in \hat{A} . By our construction of \hat{A} , $\{u, v\}$ must be an inter-cluster edge. There are several cases. Let u be in the cluster C_i and v be in the cluster C_j .

Case 1.1: $C_j = \text{par}_T(C_i)$ and there is some (p, q) after (u, v) in \hat{A} satisfying $q = v$. Suppose $C_i = \text{child}_{\text{par}_T}(C_j, l)$ for some l . If there is some (p, q) after (u, v) in \hat{A} satisfying $q = v$, then there must be a C_r which contains p , and $\text{par}_T(C_r) = C_j$, and $C_r = \text{child}_{\text{par}_T}(C_j, l')$ for some $l' > l$. This means that $(x, y) = g(C_j, \text{child}_{\text{par}_T}(C_j, l + 1))$ satisfies that $x = v$. Since in \bar{A} , $(C_j, \text{child}_{\text{par}_T}(C_j, l + 1))$ directly follows (C_i, C_j) , (x, y) directly follows (u, v) in \hat{A} , and thus (x, y) directly follows (u, v) in A' as well. Since $x = v$, $(u, v), (x, y)$ is a connected path in A' .

Algorithm 5 EULERTOURJOIN($\mathcal{C}, T, A, g, \{T_i\}, \{A_i\}$): Generating a spanning tree of V and the Euler tour of the spanning tree.

- 1: **Inputs:** $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ is a partition over a node set V , $T = (\mathcal{C}, E)$ is a spanning tree over \mathcal{C} , A is an arbitrary Euler tour of T , g is a mapping such that $\forall e = \{C_i, C_j\} \in E, g(e) = \{x, y\}$ where $x \in C_i, y \in C_j$, and $T_i = (C_i, E_i)$ is a spanning tree of C_i and A_i is an Euler tour of T_i .
 - 2: Let $\hat{E} = \{g(e) \mid e \in E\}$. Let $T' = (V, \hat{E} \cup \bigcup_{i \in [k]} E_i)$ be the spanning tree over V .
 - 3: Let $\hat{V} = \{x \in V \mid \exists y, \{x, y\} \in \hat{E}\}$.
 - 4: Choose an arbitrary root $C \in \mathcal{C}$ for T , and shift A to make it be an Euler tour of T with respect to the root C .
 - 5: For each $C_i, i \in [k]$, compute its parent $\text{par}_T(C_i)$ in T when root is C .
 - 6: For each $i \in [k]$, let $(x, y) = g(C_i, \text{par}_T(C_i))$, and circular shift A_i to make A_i be an Euler tour of T_i with root x . If C_i is a root, we circular shift A_i to make A_i be an Euler tour of T_i with an arbitrary node from \hat{V} as a root.
 - 7: For each $i \in [k]$, sort all children of C_i in T : Consider two children C_j and C_q of C_i , if $(x, y) = g(C_i, C_j), (x', y') = g(C_i, C_q)$ and $\text{first}_{A_i}(x) < \text{first}_{A_i}(x')$, then we regard that C_j has smaller rank than C_q among children of C_i . If $x = x'$, we can given an arbitrary ordering for C_j and C_q .
 - 8: Compute a new Euler tour \bar{A} for T with respect to the root C , satisfying that $\forall C_i \in \mathcal{C}, \text{first}_{\bar{A}}(\text{child}(C_i, 1)) < \text{first}_{\bar{A}}(\text{child}(C_i, 2)) < \dots < \text{first}_{\bar{A}}(\text{child}(C_i, |\text{child}(C_i)|))$.
 - 9: Compute \hat{A} : for each (C_p, C_q) appeared in the sequence of \bar{A} , replace it with $g(C_p, C_q)$.
 - 10: Compute A_x for each $x \in \hat{V}$: Suppose x is in C_i , and $A_i = ((v_1, v_2), (v_2, v_3), \dots, (v_{2 \cdot |C_i| - 2}, v_1))$, let $A_x = ((v_{\text{first}_{A_i}(x)}, v_{\text{first}_{A_i}(x)+1}), (v_{\text{first}_{A_i}(x)+1}, v_{\text{first}_{A_i}(x)+2}), \dots, (v_j, v_{(j \bmod (2 \cdot |C_i| - 2)) + 1}))$ where $j > \text{first}_{A_i}(x)$ is the smallest value such that either $j = 2 \cdot |C_i| - 2$ or $j = \text{first}_{A_i}(x')$ for some $x' \in \hat{V}$.
 - 11: Compute a mapping $f : \hat{V} \rightarrow [2 \cdot |\mathcal{C}| - 2]$. Suppose $\hat{A} = ((v'_1, v'_2), (v'_2, v'_3), \dots, (v'_{2 \cdot |\mathcal{C}| - 2}, v'_1))$. For each $x \in \hat{V}$, compute $f(x) = j$ where $j \in [2 \cdot |\mathcal{C}| - 2]$ is the largest value such that $v'_{(j \bmod (2 \cdot |\mathcal{C}| - 2)) + 1} = x$.
 - 12: Compute a sequence A' by plugging A_x for each $x \in \hat{V}$ into \hat{A} , where A_x should be inserted directly after $(v'_{f(x)}, v'_{(f(x) \bmod (2 \cdot |\mathcal{C}| - 2)) + 1})$ in \hat{A} .
 - 13: **Output:** T' and its Euler tour A' .
-

Case 1.2: $C_j = \text{par}_T(C_i)$ and there is no (p, q) after (u, v) in \hat{A} satisfying $q = v$. Suppose C_i is the last child of C_j . Then A_v must be a suffix of A_j . If C_j is the root of T , then A_v is a suffix of A' , and since A' starts with a node which is also the first node in A_j , and A_j itself is a cycle, we know that it is valid to put A_v at the end of A' . Otherwise, the edge follow (C_i, C_j) in \bar{A} is $(C_j, \text{par}_T(C_j))$. Let $(x, y) = g(\text{par}_T(C_j), C_j)$. Since we circularly shifted A_j to make A_j be an Euler tour of T_j with respect to the root y , it means that the first entry in A_j should be (y, \cdot) (i.e., started from y), and the last entry in A_j should be (\cdot, y) (i.e., ended with y). Since $(C_j, \text{par}_T(C_j))$ follows (C_i, C_j) in \bar{A} , we have (y, x) follows (u, v) in \hat{A} . By our construction of A' , we inserted A_v between (u, v) and (y, x) . Since A_v starts from v (by our construction) and ends in y (since A_v is a suffix of A_j), $(u, v), A_v, (y, x)$ is a connected path in A' .

Suppose C_i is not the last child of C_j and $C_i = \text{child}_{\text{par}_T}(C_j, l)$ for some l . Then the edge that follows (C_i, C_j) in \bar{A} must be $(C_j, \text{child}_{\text{par}_T}(C_j, l + 1))$. Let $(x, y) = g(C_j, \text{child}_{\text{par}_T}(C_j, l + 1))$. By our construction of A_v , we know that A_v must end at some node in \hat{V} . In the following, we show that A_v must end in x . Firstly, we have $\text{first}_{A_j}(x) > \text{first}_{A_j}(v)$ since this is how we used to sort all children of C_j in T . Then, if A_x ends in $z \in \hat{V}$ which is between $\text{first}_{A_j}(x)$ and $\text{first}_{A_j}(v)$, then it means that there is some cluster C_r containing z such that $\text{par}_T(C_r) = C_j$ and $g(C_j, C_r) = (z, z')$, and in addition, we have $\text{first}_{A_j}(x) < \text{first}_{A_j}(z) < \text{first}_{A_j}(v)$. This implies that the rank of C_r

among children of C_j should be between l and $l + 1$ which leads to a contradiction. Therefore, $(u, v), A_v, (x, y)$ is a connected path in A' .

Case 2.1: $C_i = \text{par}_T(C_j)$ and there is some (p, q) after (u, v) in \hat{A} satisfying $q = v$. In this case, it implies that C_j cannot be a leaf because we should be able to find C_r such that C_r contains p , $g(C_r, C_j) = (p, q)$ and $\text{par}_T(C_r) = C_j$. Note that we circular shifted A_j such that A_j starts from v . Based on how we sorted the children of C_j , we now that $g(C_j, \text{child}(C_j, 1))$ should be some (x, y) where $x = v$. In \bar{A} , $(C_j, \text{child}(C_j, 1))$ should follow (C_i, C_j) immediately, which means (x, y) follows (u, v) in \hat{A} as well as A' . Since $v = x$, $(u, v), (x, y)$ is a connected path in A' .

Case 2.2: $C_i = \text{par}_T(C_j)$ and there is no (p, q) after (u, v) in \hat{A} satisfying $q = v$. Note that we circular shifted A_j such that A_j starts from v . If C_j is a leaf in T , then the edge in \bar{A} following (C_i, C_j) is (C_j, C_i) and thus, the edge in \hat{A} following (u, v) is (v, u) . Since A_j is a cycle starting from v , we have $(u, v), A_j, (v, u)$ in A' which is a connected path (or cycle).

If C_j is not a leaf in T , the edge $(C_j, \text{child}_{\text{par}_T}(C_j, 1))$ should directly follow (C_i, C_j) in \bar{A} and thus $(x, y) = g(C_j, \text{child}_{\text{par}_T}(C_j, 1))$ directly follows (u, v) in \hat{A} . Since A_j starts from v , we have A_v is a prefix of A_j . If A_v ends in some $z \in \hat{V}$ which is before $\text{first}_{A_j}(x)$, then there must be a cluster C_r such that $g(C_j, C_r) = (z, z')$ which implies that C_r should have smaller rank than $\text{child}_{\text{par}_T}(C_j, 1)$ which contradicts to the definition of $\text{child}_{\text{par}_T}(C_j, 1)$. Therefore A_v must ends in x which implies that $(u, v), A_j, (x, y)$ is a connected path in A' .

Put them together. Since each (u, v) in \hat{A} can find a valid path in A' to the following (v, p) in \hat{A} , A' is indeed a cycle and thus A' is Euler tour of T' . \square

Theorem 15 (Euler tour join). *Given (1) a partition $\mathcal{C} = \{C_1, C_2, \dots, C_k\}$ over a set of n nodes V , (2) a spanning tree $T = (\mathcal{C}, E)$ over \mathcal{C} and an arbitrary Euler tour A of T , (3) an edge mapping g where $\forall \{C_i, C_j\} \in E, g(e) = \{x, y\}$ and $x \in C_i, y \in C_j$, and (4) spanning trees T_1, T_2, \dots, T_k and corresponding Euler tours A_1, A_2, \dots, A_k for C_1, C_2, \dots, C_k respectively, there is an MPC algorithm (Algorithm 5) that outputs a spanning tree T' over V and an Euler tour A' of T' using $O(1)$ rounds and $O(n)$ total space where the edges of T' is the union of edges in T_1, T_2, \dots, T_k and edges $\{g(e) \mid e \in E\}$. In addition, the algorithm is fully scalable.*

Proof. The correctness of Algorithm 5 is proved by Lemma 13. In the remaining of the proof, we show how to implement Algorithm 5 in the MPC model using $O(1)$ rounds, $O(n)$ total space, and make the implemented algorithm fully scalable.

Computing \hat{E} only require simultaneously access $g(e)$ for each $e \in E$. According to Theorem 4, this step only takes $O(1)$ rounds and $O(n)$ total space. Computing T' only requires local computation (i.e., renaming the tuples). To compute \hat{V} , for each $\{x, y\} \in E$, it generates x and y . This only requires local computation. Then we run sorting (Theorem 3) to remove duplicates. To choose an arbitrary $C \in \mathcal{C}$ to be a root, we can sort (Theorem 3) all C_1, C_2, \dots, C_k via an arbitrary sorting key and broadcast (Theorem 5) the name of the cluster with the smallest sorting key. We apply lemma 10 to circular shift A to make A be an Euler tour of T with respect to the root C . These steps take $O(1)$ rounds, $O(n)$ total space, and are fully scalable. To compute $\text{par}_T(C_i)$ for all C_i

simultaneously, we use sorting (Theorem 3) to compute $\mathbf{first}_A(C_i)$, then we use Observation 2 to compute $\text{par}_T(C_i)$ which only requires another call of sorting. Therefore computing $\text{par}_T(C_i)$ for all C_i can be done using $O(1)$ rounds, $O(n)$ total space, and in fully scalable setting. Computing $g(C_i, \text{par}_T(C_i))$ only requires simultaneous accesses (Theorem 4), and circular shifting all A_i at the same time can be done in parallel using Lemma 10. These steps take $O(1)$ rounds, $O(n)$ total space, and are fully scalable. To sort the children of each C_i , we do the following:

1. For all $i \in [k], x \in C_i$, simultaneously compute $\mathbf{first}_{A_i}(x)$. This step can be done via sorting (Theorem 3) which takes $O(1)$ rounds, $O(n)$ total space and is fully scalable.
2. For all $j \in [k]$, compute $(x, y) = g(C_j, C_i)$ where $C_i = \text{par}_T(C_j)$, add C_j into a set S_{C_i} and give C_j a sorting key $\mathbf{first}_{A_i}(y)$. These steps only require simultaneous accesses (Theorem 4), and thus they can be finished in $O(1)$ rounds and $O(n)$ total space. The computation is also fully scalable.
3. We use Theorem 10 to simultaneously compute the rank of each C_j in the set S_{C_i} where $C_i = \text{par}_T(C_j)$. Suppose C_j is the l -th element in S_{C_i} , then we set $\text{child}(C_i, l) = C_j$. This step takes $O(1)$ rounds, $O(n)$ total space, and is fully scalable.

Now we have parent pointers $\text{par}_T(\cdot)$ for T , a mapping $\text{child}(\cdot, \cdot)$ and an Euler tour A of T . We apply Lemma 12 to compute a new Euler tour \bar{A} of T with respect to the root C , such that $\forall C_i \in \mathcal{C}$, $\mathbf{first}_{\bar{A}}(\text{child}(C_i, 1)) < \mathbf{first}_{\bar{A}}(\text{child}(C_i, 2)) < \dots < \mathbf{first}_{\bar{A}}(\text{child}(C_i, |\text{child}(C_i)|))$. According to Lemma 12, this step only takes $O(1)$ rounds, $O(n)$ total space, and it is fully scalable. For each (C_p, C_q) appeared in \bar{A} , we simultaneously access $g(C_p, C_q)$ and replace (C_p, C_q) with $g(C_p, C_q)$ to obtain \hat{A} . This step can be done by Theorem 4, and it takes $O(1)$ rounds, $O(n)$ total space, and it is fully scalable. Next, we need to compute A_x for $x \in \hat{V}$. To do it, for each (u, v) appeared in each A_i , we check whether its position is equal to $\mathbf{first}_{A_i}(u)$, if it is, we mark (u, v) as 1 otherwise, we mark (u, v) as 0. Then, for each (u, v) we find the closest (u', v') appeared in A_i such that (u', v') is marked as 1 and is appeared before (u, v) . Then we put (u, v) into $A_{u'}$ and its index can be derived from the distance between (u, v) and (u', v') in A_i . Above steps can be implemented by simultaneous accesses (Theorem 4) and predecessor algorithm (Theorem 9). Therefore, computing A_x for all $x \in \hat{V}$ simultaneously only requires $O(1)$ rounds and $O(n)$ total space. Computing f can be done using sorting (Theorem 3) which can be done in $O(1)$ rounds and $O(n)$ total space and is fully scalable. Finally, we run sequence insertion algorithm (Theorem 13) on $\hat{A}, f, \{A_x\}$, to obtain A' , and it takes $O(1)$ rounds and $O(n)$ total space and is fully scalable as well. \square

5.3 Euler Tour via Hierarchical Decomposition

Let V be a set of nodes. Let $\mathcal{C}_0 \supseteq \mathcal{C}_1 \supseteq \dots \supseteq \mathcal{C}_L$ be a hierarchy of partitions on V where $\mathcal{C}_0 = \{\{v\} \mid v \in V\}$ and $\mathcal{C}_L = \{V\}$. Let E_1, E_2, \dots, E_L be arbitrary sets of edges between nodes in V such that $\forall l \in [L], \mathcal{C}_{l-1} \oplus E_l = \mathcal{C}_l$ and $|\mathcal{C}_{l-1}| - |E_l| = |\mathcal{C}_l|$, i.e., any edge in E_l only connects two different components in \mathcal{C}_{l-1} . Then it is easy to see that $T = (V, \bigcup_{l \in [L]} E_l)$ is a spanning tree of V . If L is small, and for every $l \in [L]$, the tree obtained by regarding each component in \mathcal{C}_{l-1} as a node and regarding each pair of components in \mathcal{C}_{l-1} that are connected by an edge in E_l as an edge has a low diameter, then we provide a novel MPC algorithm to efficiently compute an Euler tour of T .

Algorithm 6 EULERTOURVIAHIERARCHICALDECOMPOSITION($\{\mathcal{C}_l\}$ (represented by ℓ_l), $\{E_l\}$):
Generating an Euler tour of the spanning tree $T = (V, E_1 \cup E_2 \cup \dots \cup E_L)$.

- 1: **Inputs:** Leader mappings $\ell_0(\cdot), \ell_1(\cdot), \dots, \ell_L(\cdot)$ representing a hierarchical decomposition $\mathcal{C}_0 \supseteq \mathcal{C}_1 \supseteq \dots \supseteq \mathcal{C}_L$ of n nodes V where $\mathcal{C}_0 = \{\{v\} \mid v \in V\}$ and $\mathcal{C}_L = \{V\}$, and sets of edges E_1, E_2, \dots, E_L between nodes in V such that $\forall l \in [L], \mathcal{C}_{l-1} \oplus E_l = \mathcal{C}_l$ and $|\mathcal{C}_{l-1}| - |E_l| = |\mathcal{C}_l|$.
- 2: Let $R = \log L$. {We suppose L is a power of 2.}
- 3: {Notation: $\forall l \in [L] \cup \{0\}, \forall C \in \mathcal{C}_l, \forall l' \leq l$, we denote

$$C^{(l')} = \{C' \in \mathcal{C}_{l'} \mid C' \subseteq C\},$$

i.e., interpret C as a set of clusters at level l' . Similarly, $\forall l \in [L] \cup \{0\}, \forall E \subseteq \bigcup_{j=l}^L E_j, \forall l' < l$, we denote

$$E^{(l')} = \{\{C_x, C_y\} \mid \{x, y\} \in E, x \in C_x, y \in C_y, \text{ and } C_x, C_y \in \mathcal{C}_{l'}\},$$

i.e., interpret E as a set of edges between nodes at level l' , where each node denotes a cluster at level l' . Let $E(C)$ denotes the subset of edges where both end nodes are in C . For $0 \leq l' \leq l \leq L$, we denote $g^{(l \rightarrow l')}$ as following, if there is an edge $\{x, y\} \in \bigcup_{j=l}^L E_j$, then $g^{(l \rightarrow l')}(\{C_x, C_y\}) = \{C'_x, C'_y\}$ where $C_x, C_y \in \mathcal{C}_l, C'_x, C'_y \in \mathcal{C}_{l'}, x \in C'_x \subseteq C_x, y \in C'_y \subseteq C_y$. }

- 4: **for** $l \in [L]$ **in parallel do**
 - 5: **for** $C \in \mathcal{C}_l$ **in parallel do**
 - 6: Compute the Euler tour $A_l^{(0)}(C)$ of the tree $T_l^{(0)}(C) = (C^{(l-1)}, E_l(C)^{(l-1)})$
 - 7: **for** $r := 1 \rightarrow R$ **do**
 - 8: **for** $l \in \{1 \cdot 2^r, 2 \cdot 2^r, 3 \cdot 2^r, 4 \cdot 2^r, \dots, L\}$ **in parallel do**
 - 9: **for** $C \in \mathcal{C}_l$ **in parallel do**
 - 10: $T_l^{(r)}(C), A_l^{(r)}(C) \leftarrow \text{EULERTOURJOIN}(\mathcal{C}, T, A, g, \{T_i\}, \{A_i\})$ {Algorithm 5} where:
 $C = C^{(l-2^{r-1})}, T = T_l^{(r-1)}(C), A = A_l^{(r-1)}(C), g = g^{(l-2^{r-1} \rightarrow l-2^r)},$
 $\{T_i\} = \{T_{l-2^{r-1}}^{(r-1)}(C') \mid C' \in C^{(l-2^{r-1})}\}, \{A_i\} = \{A_{l-2^{r-1}}^{(r-1)}(C') \mid C' \in C^{(l-2^{r-1})}\}$
 - 11: **Output:** Euler tour $A_L^{(R)}(V)$ of spanning tree $T_L^{(R)}(V)$.
-

In the MPC model, we use a leader mapping $\ell_l : V \rightarrow V$ to denote the partitioning \mathcal{C}_l , i.e., x, y are in the same component in \mathcal{C}_l iff $\ell_l(x) = \ell_l(y)$. The edges E_l are distributed arbitrarily on the machines.

Lemma 14. (*correctness of Algortihm 6*) For $r \in \{0\} \cup [R], l \in [L]$ with $l \bmod 2^r = 0, \forall C \in \mathcal{C}_l$, $T_l^{(r)}(C)$ and $A_l^{(r)}(C)$ satisfy following properties:

1. $T_l^{(r)}(C)$ is a spanning tree where each node in the tree denotes a cluster $C' \in \mathcal{C}_{l-2^r}$ and $C' \subseteq C$.
2. The edge set of $T_l^{(r)}(C)$ is $\{\{C_x, C_y\} \mid \{x, y\} \in \bigcup_{j=l-2^r+1}^l E_j, x \in C_x, y \in C_y, C_x, C_y \in \mathcal{C}_{l-2^r}, C_x, C_y \subseteq C\}$.
3. $A_l^{(r)}(C)$ is an Euler tour of $T_l^{(r)}(C)$.

Proof. Our proof is by induction. The base case is $r = 0$, the claimed properties of $T_l^{(0)}(C), A_l^{(0)}(C)$ obviously hold for all $l \in [L], C \in \mathcal{C}_l$.

Now consider the case that the claimed properties hold for $T_l^{(r-1)}(C)$, $A_l^{(r-1)}(C)$ for all $l \in [L]$ with $l \bmod 2^{r-1} = 0$ and all $C \in \mathcal{C}_l$. Consider a cluster $C \in \mathcal{C}_l$ and how we compute $T_l^{(r)}(C)$ and $A_l^{(r)}(C)$. By our induction hypothesis, $T_l^{(r-1)}(C)$ is a spanning tree where each node in the tree denotes a cluster $C' \in \mathcal{C}_{l-2^{r-1}}$ and $C' \subseteq C$. $A_l^{(r-1)}(C)$ is an Euler tour of $T_l^{(r-1)}(C)$. In addition, the edge set of $T_l^{(r-1)}(C)$ is $\{\{C_x, C_y\} \mid \{x, y\} \in \bigcup_{j=l-2^{r-1}+1}^l E_j, x \in C_x, y \in C_y, C_x, C_y \in \mathcal{C}_{l-2^{r-1}}, C_x, C_y \subseteq C\}$. Then for each $C' \in \mathcal{C}_{l-2^{r-1}}$ and $C' \subseteq C$, we have that $T_{l-2^{r-1}}^{(r-1)}(C')$ is a spanning tree where each node in the tree denotes a cluster $C'' \in \mathcal{C}_{l-2^r}$ and $C'' \subseteq C'$. $A_{l-2^{r-1}}^{(r-1)}(C')$ is an Euler tour of $T_{l-2^{r-1}}^{(r-1)}(C')$. The edges of $T_{l-2^{r-1}}^{(r-1)}(C')$ is $\{\{C_x, C_y\} \mid \{x, y\} \in \bigcup_{j=l-2^r+1}^{l-2^{r-1}} E_j, x \in C_x, y \in C_y, C_x, C_y \in \mathcal{C}_{l-2^r}, C_x, C_y \subseteq C'\}$. Since $g^{(l-2^{r-1} \rightarrow l-2^r)}$ maps each edge in $\{\{C_x, C_y\} \mid \{x, y\} \in \bigcup_{j=l-2^{r-1}+1}^l E_j, x \in C_x, y \in C_y, C_x, C_y \in \mathcal{C}_{l-2^{r-1}}, C_x, C_y \subseteq C\}$ to a corresponding edge in $\{\{C_x, C_y\} \mid \{x, y\} \in \bigcup_{j=l-2^r+1}^l E_j, x \in C_x, y \in C_y, C_x, C_y \in \mathcal{C}_{l-2^r}, C_x, C_y \subseteq C\}$. By applying Theorem 15, we have that $T_l^{(r)}(C)$ is a spanning tree where each node in the tree denotes a cluster $C' \in \mathcal{C}_{l-2^r}$ and $C' \subseteq C$. $A_l^{(r)}(C)$ is an Euler tour of $T_l^{(r)}(C)$. In addition, The edge set of $T_l^{(r)}(C)$ is $\{\{C_x, C_y\} \mid \{x, y\} \in \bigcup_{j=l-2^r+1}^l E_j, x \in C_x, y \in C_y, C_x, C_y \in \mathcal{C}_{l-2^r}, C_x, C_y \subseteq C\}$. Therefore, the stated properties also hold for $T_l^{(r)}(C)$ and $A_l^{(r)}(C)$. \square

If we regard each singleton cluster as the node itself, we get the following corollary.

Corollary 6 (Correctness of final output of Algorithm 6). *At the end of Algorithm 6, $T_L^{(R)}(V)$ is a spanning tree over V . The edge set of $T_L^{(R)}(V)$ is $\bigcup_{j=1}^L E_j$. $A_L^{(R)}(V)$ is an Euler tour of $T_L^{(R)}(V)$.*

Theorem 16 (Euler tour via hierarchical decomposition). *Consider leader mappings $\ell_0(\cdot), \ell_1(\cdot), \dots, \ell_L(\cdot)$ representing a hierarchical decomposition $\mathcal{C}_0 \supseteq \mathcal{C}_1 \supseteq \dots \supseteq \mathcal{C}_L$ of n nodes V where $\mathcal{C}_0 = \{\{v\} \mid v \in V\}$ and $\mathcal{C}_L = \{V\}$, and arbitrary sets of edges E_1, E_2, \dots, E_L between nodes in V such that $\forall l \in [L], \mathcal{C}_{l-1} \oplus E_l = \mathcal{C}_l$ and $|\mathcal{C}_{l-1}| - |E_l| = |\mathcal{C}_l|$. Let Λ be an upper bound such that $\forall l \in [L]$, if we regard each cluster in \mathcal{C}_{l-1} as a node and each pair of clusters with an edge in E_l connecting them as an edge, the diameter of the tree is at most Λ . Then there is a fully scalable MPC algorithm (Algorithm 6) which takes $O(\log(L) + \log(\Lambda))$ rounds and $O(nL + n^{1+\varepsilon})$ total space outputting an Euler tour of the spanning tree $T = \left(V, \bigcup_{l \in [L]} E_l\right)$ where $\varepsilon > 0$ is an arbitrarily small constant.*

Proof. The correctness is proved by Corollary 6. In the following, we show how to implement Algorithm 6 in the MPC model. Note that we store all leader mappings on machines and thus it takes total space $O(nL)$. The inputs for computing $T_l^{(0)}(C)$ over all $l \in [L], C \in \mathcal{C}_l$ are disjoint, thus we can use sorting (Theorem 3) to assign each Euler tour computation to a disjoint group of machines. Then we independently apply Corollary 14 for each task in parallel. Therefore, we use $O(\log \Lambda)$ rounds and $O(n^{1+\varepsilon})$ total space to compute $T_l^{(0)}(C)$ for all $l \in [L], C \in \mathcal{C}_l$. The algorithm is fully scalable.

Then, we run R iterations. In each iteration, we run for all considered l and all $C \in \mathcal{C}_l$ in parallel. We need to call EULERTOURJOIN (Algorithm 5) for each l and C . Due to Lemma 14, the inputs for subroutines of EULERTOURJOIN in one iteration are disjoint, and the total size of inputs are at most $O(n)$. Therefore, we can use sorting (Theorem 3) to assign each subroutine to a disjoint group of machines. Note that, we only need to compute one $g^{(l-2^{r-1} \rightarrow l-2^r)}$ for all $C \in \mathcal{C}_l$. We

only need simultaneous accesses (Theorem 4) to $\ell_{l-2^{r-1}}$ and ℓ_{l-2^r} to compute $g^{(l-2^{r-1} \rightarrow l-2^r)}$. Therefore, to prepare the inputs for each subroutine and assign machines for each subtask, we need $O(1)$ rounds, $O(nL)$ total space and the process is fully scalable. Then, for each subtask of calling EULERTOURJOIN, we apply Theorem 15, which takes $O(n)$ space in total. It takes $O(1)$ rounds and is fully scalable. Since we have $O(R) = O(\log L)$ iterations, it takes $O(R)$ rounds.

Therefore, the total number of rounds is $O(\log R + \log \Lambda)$, and the total space required is at most $O(n^{1+\varepsilon} + nL)$ where $\varepsilon > 0$ is an arbitrarily small constant. The entire algorithm is fully scalable. \square

5.4 Euler Tour of Approximate Euclidean MST

Then, by applying Algorithm 6 on the approximate MST and the corresponding hierarchical decomposition $\hat{\mathcal{P}}_1, \hat{\mathcal{P}}_2, \hat{\mathcal{P}}_4, \dots, \hat{\mathcal{P}}_{\alpha_H}$ that we obtained in Section 2 and Section 4, we are able to output an Euler tour of our approximate MST.

Theorem 17 (Approximate Euclidean MST with Euler tour). *Given n points from \mathbb{R}^d , there is a fully scalable MPC algorithm which outputs an $O(1)$ -approximate MST with probability at least 0.99. In addition, the algorithm also outputs an Euler tour of the outputted approximate MST. The number of rounds of the algorithm is $O(\log \log(n) \cdot \log \log \log(n))$. The total space required is at most $O(nd + n^{1+\varepsilon})$ where $\varepsilon > 0$ is an arbitrary small constant.*

Proof. The approximate Euclidean MST is shown by Theorem 12. Note that at the end of Part 3 (Algorithm 4) of our algorithm, we also obtain a hierarchical decomposition: $\hat{\mathcal{P}}_1 \supseteq \hat{\mathcal{P}}_2 \supseteq \hat{\mathcal{P}}_4 \supseteq \dots \supseteq \hat{\mathcal{P}}_{\alpha_H}$ which has at most $\text{polylog}(n)$ levels. In addition, we showed that $\hat{\mathcal{P}}_{t/2} \oplus (F_t \cup \bigcup_{i \in [h]} E_t^{(i)}) = \hat{\mathcal{P}}_t$ and $|\hat{\mathcal{P}}_{t/2}| - |F_t \cup \bigcup_{i \in [h]} E_t^{(i)}| = |\hat{\mathcal{P}}_t|$.

Claim 1. *For any t , if we regard each cluster in $\hat{\mathcal{P}}_{t/2}$ as a node and each pair of clusters in $\hat{\mathcal{P}}_{t/2}$ that is connected by an edge in $F_t \cup \bigcup_{i \in [h]} E_t^{(i)}$ as an edge, then each tree has diameter at most $2^{O(h)} = \text{polylog}(n)$.*

Proof. If we just run one round leader compression, each tree can have diameter at most 2. Since in every round of leader compression, we merge clusters using a star which may blow up the diameter by at most a factor of 5. After leader compression, we create stars to merge incomplete components. This operation can blow up the diameter by a factor of 5 as well. Therefore, the diameter can be at most $5^{O(h)} = \text{polylog}(n)$. \square

Then, by plugging $\{\hat{\mathcal{P}}_t \mid t = 2^j, j \in [O(\log n)]\}$ and $\{F_t \cup \bigcup_{i \in [h]} E_t^{(i)} \mid t = 2^j, j \in [O(\log n)]\}$ into Theorem 16, we can use additional $O(n^{1+\varepsilon})$ total space and $O(\log \log(n))$ number of rounds to compute an Euler tour of our approximate MST. The algorithm is fully scalable. \square

The Euclidean travelling salesman problem (TSP) is stated as the following: Given n points, the goal is to output a cycle over points such that each point is visited exactly once such that the total length of the cycle is minimized.

Corollary 7 (Approximate Euclidean TSP). *Given n points from \mathbb{R}^d , there is a fully scalable MPC algorithm which outputs an $O(1)$ -approximate TSP solution with probability at least 0.99. The number of rounds of the algorithm is $O(\log \log(n) \cdot \log \log \log(n))$. The total space required is at most $O(nd + n^{1+\varepsilon})$ where $\varepsilon > 0$ is an arbitrary small constant.*

Proof. It is easy to observe that the MST cost and the length of Euler tour of the MST are the same up to a factor of 2. Since optimal TSP cost is greater than MST cost and less than the cost of Euler tour, we only need to output a shortcut Euler tour of a constant approximate MST to get a constant approximate TSP solution.

By applying Theorem 17, we obtain an Euler tour of $O(1)$ approximate MST. Then we can use sorting and (re)indexing (Theorem 3) to firstly only keep the first appearance of each point on the Euler tour and then recompute the index of each point in the deduplicated tour sequence to provide an $O(1)$ -approximate TSP solution. Above steps only takes $O(1)$ additional MPC rounds and $O(n)$ additional total space. These steps are fully scalable. \square

References

- [AAH⁺23] AmirMohsen Ahanchi, Alexandr Andoni, MohammadTaghi Hajiaghayi, Marina Kinitel, and Peilin Zhong. Massively parallel tree embeddings for high dimensional spaces. In *Proceedings of the 35th ACM Symposium on Parallelism in Algorithms and Architectures*, pages 77–88, 2023.
- [ACK⁺16] Alexandr Andoni, Jiecao Chen, Robert Krauthgamer, Bo Qin, David P Woodruff, and Qin Zhang. On sketching quadratic forms. In *Proceedings of the 2016 ACM Conference on Innovations in Theoretical Computer Science*, pages 311–319. ACM, 2016.
- [AESW90] Pankaj K Agarwal, Herbert Edelsbrunner, Otfried Schwarzkopf, and Emo Welzl. Euclidean minimum spanning trees and bichromatic closest pairs. In *Proceedings of the sixth annual symposium on Computational geometry*, pages 203–210, 1990.
- [AIK08] Alexandr Andoni, Piotr Indyk, and Robert Krauthgamer. Earth mover distance over high-dimensional spaces. In *SODA*, volume 8, pages 343–352. Citeseer, 2008.
- [ANOY14] Alexandr Andoni, Aleksandar Nikolov, Krzysztof Onak, and Grigory Yaroslavtsev. Parallel algorithms for geometric graph problems. In *Proceedings of the Forty-Sixth Annual ACM Symposium on Theory of Computing, STOC '14*, page 574–583, New York, NY, USA, 2014. Association for Computing Machinery.
- [ASS⁺18] Alexandr Andoni, Zhao Song, Clifford Stein, Zhengyu Wang, and Peilin Zhong. Parallel graph connectivity in log diameter rounds. In *2018 IEEE 59th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 674–685. IEEE, 2018.
- [ASW19] Sepehr Assadi, Xiaorui Sun, and Omri Weinstein. Massively parallel algorithms for finding well-connected components in sparse graphs. In *Proceedings of the 2019 ACM Symposium on principles of distributed computing*, pages 461–470, 2019.

- [Bar96] Yair Bartal. Probabilistic approximation of metric spaces and its algorithmic applications. In *Proceedings of 37th Conference on Foundations of Computer Science*, pages 184–193. IEEE, 1996.
- [BBD⁺17a] Mohammadhossein Bateni, Soheil Behnezhad, Mahsa Derakhshan, MohammadTaghi Hajiaghayi, Raimondas Kiveris, Silvio Lattanzi, and Vahab Mirrokni. Affinity clustering: Hierarchical clustering at scale. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [BBD⁺17b] MohammadHossein Bateni, Soheil Behnezhad, Mahsa Derakhshan, MohammadTaghi Hajiaghayi, Raimondas Kiveris, Silvio Lattanzi, and Vahab Mirrokni. Affinity clustering: Hierarchical clustering at scale. *Advances in Neural Information Processing Systems*, 30, 2017.
- [BDE⁺19] Soheil Behnezhad, Laxman Dhulipala, Hossein Esfandiari, Jakub Lacki, and Vahab Mirrokni. Near-optimal massively parallel graph connectivity. In *2019 IEEE 60th Annual Symposium on Foundations of Computer Science (FOCS)*, pages 1615–1636. IEEE, 2019.
- [BKS17] Paul Beame, Paraschos Koutris, and Dan Suciu. Communication steps for parallel query processing. *Journal of the ACM (JACM)*, 64(6):1–58, 2017.
- [Bor26] Otakar Boruvka. O jistém problému minimálním. 1926.
- [CAMZ22] Vincent Cohen-Addad, Vahab Mirrokni, and Peilin Zhong. Massively parallel k -means clustering for perturbation resilient instances. In *International Conference on Machine Learning*, pages 4180–4201. PMLR, 2022.
- [CC22] Sam Coy and Artur Czumaj. Deterministic massively parallel connectivity. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 162–175, 2022.
- [CCAJ⁺23] Xi Chen, Vincent Cohen-Addad, Rajesh Jayaram, Amit Levi, and Erik Waingarten. Streaming euclidean mst to a constant factor. In *Proceedings of the 55th Annual ACM Symposium on Theory of Computing*, STOC 2023, page 156–169, New York, NY, USA, 2023. Association for Computing Machinery.
- [CCFC02] Moses Charikar, Kevin Chen, and Martin Farach-Colton. Finding frequent items in data streams. In *Proceedings of the 29th International Colloquium on Automata, Languages and Programming*, ICALP ’02, pages 693–703, London, UK, UK, 2002. Springer-Verlag.
- [CEF⁺05] Artur Czumaj, Funda Ergün, Lance Fortnow, Avner Magen, Ilan Newman, Ronitt Rubinfeld, and Christian Sohler. Approximating the weight of the euclidean minimum spanning tree in sublinear time. *SIAM J. Comput.*, 35(1):91–109, 2005.
- [Cha00] Bernard Chazelle. A minimum spanning tree algorithm with inverse-ackermann type complexity. *J. ACM*, 47(6):1028–1047, nov 2000.

- [CJLW22] Xi Chen, Rajesh Jayaram, Amit Levi, and Erik Waingarten. New streaming algorithms for high dimensional emd and mst. In *Proceedings of the 54th Annual ACM SIGACT Symposium on Theory of Computing*, pages 222–233, 2022.
- [CLRS22] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. *Introduction to algorithms*. MIT press, 2022.
- [CRT05] Bernard Chazelle, Ronitt Rubinfeld, and Luca Trevisan. Approximating the minimum spanning tree weight in sublinear time. *SIAM Journal on Computing*, 34(6):1370–1379, 2005.
- [CS04] Artur Czumaj and Christian Sohler. Estimating the weight of metric minimum spanning trees in sublinear-time. In László Babai, editor, *Proceedings of the 36th Annual ACM Symposium on Theory of Computing, Chicago, IL, USA, June 13-16, 2004*, pages 175–183. ACM, 2004.
- [CS09] Artur Czumaj and Christian Sohler. Estimating the weight of metric minimum spanning trees in sublinear time. *SIAM Journal on Computing*, 39(3):904–922, 2009.
- [DCLT18] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- [DG04] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: Simplified data processing on large clusters. 2004.
- [DG08] Jeffrey Dean and Sanjay Ghemawat. Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113, 2008.
- [EMMZ22] Alessandro Epasto, Mohammad Mahdian, Vahab Mirrokni, and Peilin Zhong. Massively parallel and dynamic algorithms for minimum size clustering. In *Proceedings of the 2022 Annual ACM-SIAM Symposium on Discrete Algorithms (SODA)*, pages 1613–1660. SIAM, 2022.
- [Epp00] David Eppstein. Spanning trees and spanners., 2000.
- [FIS05] Gereon Frahling, Piotr Indyk, and Christian Sohler. Sampling in dynamic data streams and applications. In *Proceedings of the Twenty-First Annual Symposium on Computational Geometry, SCG ’05*, page 142–149, New York, NY, USA, 2005. Association for Computing Machinery.
- [Goo99] Michael T Goodrich. Communication-efficient parallel sorting. *SIAM Journal on Computing*, 29(2):416–432, 1999.
- [GSZ11] Michael T Goodrich, Nodari Sitchinava, and Qin Zhang. Sorting, searching, and simulation in the mapreduce framework. In *International Symposium on Algorithms and Computation*, pages 374–383. Springer, 2011.
- [GZJ06] Oleksandr Grygorash, Yan Zhou, and Zach Jorgensen. Minimum spanning tree based clustering algorithms. In *2006 18th IEEE International Conference on Tools with Artificial Intelligence (ICTAI’06)*, pages 73–81. IEEE, 2006.

- [HIS13] Sarel Har-Peled, Piotr Indyk, and Anastasios Sidiropoulos. Euclidean spanners in high dimensions. In Sanjeev Khanna, editor, *Proceedings of the Twenty-Fourth Annual ACM-SIAM Symposium on Discrete Algorithms, SODA 2013, New Orleans, Louisiana, USA, January 6-8, 2013*, pages 804–809. SIAM, 2013.
- [HPIM12] Sarel Har-Peled, Piotr Indyk, and Rajeev Motwani. Approximate nearest neighbor: Towards removing the curse of dimensionality. *Theory of computing*, 8(1):321–350, 2012.
- [HZRS16] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [IBY⁺07] Michael Isard, Mihai Budiu, Yuan Yu, Andrew Birrell, and Dennis Fetterly. Dryad: distributed data-parallel programs from sequential building blocks. In *Proceedings of the 2nd ACM SIGOPS/EuroSys European Conference on Computer Systems 2007*, pages 59–72, 2007.
- [Ind99] Piotr Indyk. Sublinear time algorithms for metric space problems. In *Proceedings of the thirty-first annual ACM symposium on Theory of computing*, pages 428–434, 1999.
- [Ind04] Piotr Indyk. Algorithms for dynamic geometric problems over data streams. In *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*, pages 373–380. ACM, 2004.
- [IT03] Piotr Indyk and Nitin Thaper. Fast color image retrieval via embeddings. In *Workshop on Statistical and Computational Theories of Vision (at ICCV)*, 2003.
- [JL84] William B. Johnson and Joram Lindenstrauss. Extensions of lipschitz mappings into a hilbert space. *Contemporary Mathematics*, 26:189–206, 1984.
- [JN18] Tomasz Jurdziński and Krzysztof Nowicki. Mst in $o(1)$ rounds of congested clique. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 2620–2632. SIAM, 2018.
- [KSV10] Howard Karloff, Siddharth Suri, and Sergei Vassilvitskii. A model of computation for mapreduce. In *Proceedings of the twenty-first annual ACM-SIAM symposium on Discrete Algorithms*, pages 938–948. SIAM, 2010.
- [LMSV11] Silvio Lattanzi, Benjamin Moseley, Siddharth Suri, and Sergei Vassilvitskii. Filtering: a method for solving graph problems in mapreduce. In *Proceedings of the twenty-third annual ACM symposium on Parallelism in algorithms and architectures*, pages 85–94, 2011.
- [LMW18] Jakub Lacki, Vahab Mirrokni, and Michał Włodarczyk. Connected components at scale via local contractions. *arXiv preprint arXiv:1807.10727*, 2018.
- [LRN09] Chih Lai, Taras Rafea, and Dwight E Nelson. Approximate minimum spanning tree clustering in high-dimensional space. *Intelligent Data Analysis*, 13(4):575–597, 2009.

- [MCCD13] Tomas Mikolov, Kai Chen, Greg Corrado, and Jeffrey Dean. Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*, 2013.
- [RVW18] Tim Roughgarden, Sergei Vassilvitskii, and Joshua R Wang. Shuffles and circuits (on lower bounds for modern parallel computation). *Journal of the ACM (JACM)*, 65(6):1–24, 2018.
- [VdMH08] Laurens Van der Maaten and Geoffrey Hinton. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [Whi12] Tom White. *Hadoop: The definitive guide*. ” O’Reilly Media, Inc.”, 2012.
- [WWW09] Xiaochun Wang, Xiali Wang, and D Mitchell Wilkes. A divide-and-conquer approach for minimum spanning tree-based clustering. *IEEE Transactions on Knowledge and Data Engineering*, 21(7):945–958, 2009.
- [YV18] Grigory Yaroslavtsev and Adithya Vadapalli. Massively parallel algorithms and hardness for single-linkage clustering under lp distances. In *International Conference on Machine Learning*, pages 5600–5609. PMLR, 2018.
- [ZCF⁺10] Matei Zaharia, Mosharaf Chowdhury, Michael J Franklin, Scott Shenker, and Ion Stoica. Spark: Cluster computing with working sets. In *2nd USENIX Workshop on Hot Topics in Cloud Computing (HotCloud 10)*, 2010.
- [ZMMF15] Caiming Zhong, Mikko Malinen, Duoqian Miao, and Pasi Fränti. A fast minimum spanning tree algorithm based on k-means. *Information Sciences*, 295:1–17, 2015.