Andrey Prokopenko prokopenkoav@ornl.gov Oak Ridge National Laboratory Oak Ridge, Tennessee, USA Piyush Sao Oak Ridge National Laboratory Oak Ridge, Tennessee, USA

ABSTRACT

Computing the Euclidean minimum spanning tree (EMST) is a computationally demanding step of many algorithms. While workefficient serial and multithreaded algorithms for computing EMST are known, designing an efficient GPU algorithm is challenging due to a complex branching structure, data dependencies, and load imbalances. In this paper, we propose a single-tree Borůvka-based algorithm for computing EMST on GPUs. We use an efficient nearest neighbor algorithm and reduce the number of the required distance calculations by avoiding traversing subtrees with leaf nodes in the same component. The developed algorithms are implemented in a performance portable way using ArborX, an open-source geometric search library based on the Kokkos framework. We evaluate the proposed algorithm on various 2D and 3D datasets, show and compare it with the current state-of-the-art open-source CPU implementations. We demonstrate 4-24× speedup over the fastest multi-threaded implementation. We prove the portability of our implementation by providing results on a variety of hardware: AMD EPYC 7763, Nvidia A100 and AMD MI250X. We show scalability of the implementation, computing EMST for 37 million 3D cosmological dataset in under a 0.5 second on a single A100 Nvidia GPU.

CCS CONCEPTS

• **Computing methodologies** → *Parallel algorithms*.

KEYWORDS

Euclidean minimum spanning tree, parallel algorithm, GPU

ACM Reference Format:

Andrey Prokopenko, Piyush Sao, and Damien Lebrun-Grandié. 2022. A single-tree algorithm to compute the Euclidean minimum spanning tree on GPUs. In 51st International Conference on Parallel Processing (ICPP '22), August 29-September 1, 2022, Bordeaux, France. ACM, New York, NY, USA, 10 pages. https://doi.org/10.1145/3545008.3546185

https://doi.org/10.1145/3545008.3546185

Damien Lebrun-Grandié Oak Ridge National Laboratory Oak Ridge, Tennessee, USA



Figure 1: Performance (in *MFeatures/sec*) for the dualtree [18] (MLPACK), WSPD [27] (MEMOGFK), and single-tree (this work using ARBORX) approaches on AMD EPYC 7763 CPU (sequential and multi-threaded), and Nvidia A100 and AMD MI250X (single GCD) GPU architectures for a 3D cosmological dataset.

1 INTRODUCTION

Given a set of n points in a d-dimensional space, the *Euclidean* minimum spanning tree (EMST) problem determines the minimum spanning tree (MST) of the distance graph of the set, i.e., a graph where each pair of vertices are connected by an edge of weight equal to the Euclidean distance between them. Computing EMST is an important task in a variety of applications, including data clustering [7], Euclidean traveling salesman problem [12], cosmology [22], wireless network connectivity [17], computational fluid dynamics [25], and many others.

Most algorithms for computing an MST on a general graph are variants of the three classical algorithms: 1926 Borůvka's algorithm [5], 1956 Kruskal's algorithm [15] and 1957 Prim's algorithm [24]. These algorithms share the same general idea, constructing an MST iteratively. At any instant during the computation, an algorithm maintains a set of non-overlapping sets of vertices called components. Initially, all components consist of a single vertex. On each step of an algorithm, some components are merged using a subset of the graph edges. The algorithm terminates when there remain no edges connecting separate components.

The fundamental difference of the EMST problem and the MST one lies in the graph structure. MST operates on a sparse graph, where the number of edges is a small fraction of all possible edges with the same vertices. On the other hand, EMST uses the distance graph, which is *complete*, with each pair of vertices connected by an edge, for a total of n(n - 1)/2 edges. It is prohibitively expensive to both construct and store a complete graph for large problems,

Publication rights licensed to ACM. ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or affiliate of the United States government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only.

ICPP '22, August 29-September 1, 2022, Bordeaux, France

^{© 2022} Copyright held by the owner/author(s). Publication rights licensed to ACM. ACM ISBN 978-1-4503-9733-9/22/08...\$15.00

This manuscript has been authored by UT-Battelle, LLC, under contract DE-AC05-00OR22725 with the U.S. Department of Energy. The United States Government retains and the publisher, by accepting the article for publication, acknowledges that the United States Government retains a nonexclusive, paid-up, irrevocable, world-wide license to publish or reproduce the published form of this manuscript, or allow others to do so, for United States Government purposes.

ICPP '22, August 29-September 1, 2022, Bordeaux, France

Prokopenko et al.



Figure 2: Borůvka's algorithm. (a) Initial state (each component having a single vertex). (b) The state after a few Borůvka iterations. (c) Closest neighbors from a different component for each vertex. (d) The shortest outgoing edge for each component. (e) The new components after the merge (the initial state of the next Borůvka iteration).

as well as run classical MsT algorithms on such a graph. Thus, the distance graph is typically used implicitly.

To solve the EMST problem, Bentley and Friedman [4] proposed combining classical MST algorithms with a nearest-neighbor search algorithm using a spatial indexing structure (e.g., a *k*-d tree). However, a straightforward implementation of this approach performs poorly. The main bottleneck of the algorithm is in the excessive number of distance calculations in the later iterations of an MST algorithm. Since, the key challenge in designing an efficient EMST algorithm has been in careful pruning of distance calculations during the runtime.

Two popular approaches have emerged: one based on the *well* separated pair decomposition (WSPD) [6, 23, 27], and the other using the *dual-tree* framework [18]. Under some mild assumptions on the distribution of the data points, the dual-tree method provides the sharpest worst-case bounds for any dimensional space. The two approaches have been shown to perform well on CPUs, including multi-threaded parallelization.

However, the existing approaches are limited in scalability and performance for an efficient GPU computation. The dual-tree algorithm, in general, is outperformed by the best-known WSPD-based approach. And while WSPD is asymptotically O(n)[2], the hidden constants of the algorithm are very high. In fact, we observed that the WSPD computation dominates the overall time in the EMST computation and is comparable with a single-tree implementation in the sequential case. Thus, an EMST algorithm or implementation that is sequentially efficient, is scalable with respect to problem size, and is amenable to GPU parallelism remains an open problem.

In this paper, we propose an efficient algorithm for EMST suitable for both CPUs and GPUs. Our algorithm is based on Borůvka's algorithm and uses a *single* tree to perform the nearest-neighbor queries. We use a bounding volume hierarchy (BVH) as our tree structure as it is very efficient for unstructured low-dimensional data on GPU (we note, however, that the described algorithms are general and are applicable to other tree structures such as k-d tree). To reduce the number of the distance calculations when finding the closest outgoing edge for a given component (the most expensive part of each Borůvka iteration), we keep track of the component membership of the children of the internal tree nodes. This allows nearest neighbor queries to bypass subtrees where all leaf nodes lie in the same component. Our motivation for this work comes from an astronomy application which requires high performance to analyze the data from a cosmological simulation. We show the results from one such dataset in Figure 1. Additionally, we show that our GPU implementation achieves 270 *MFeatures/sec* (million features, the product of the number of points and dimensions, processed per second) on nVidia A100, which is 17× faster than best known multithreaded implementation.

In our implementation, we used ArborX [16], a performance portable geometric search library using Kokkos framework [26]. This allows us to study the algorithm on both CPU and multiple GPU architectures (e.g., Nvidia A100, AMD MI250X). We evaluate the proposed algorithm on various 2D and 3D datasets, and compare it with the current state-of-the-art open-source CPU implementations.

Our key contributions are:

- We provide the first performance portable algorithm and implementation for the EMST problem. Our algorithm is efficient in the sequential case and outperforms the best publicly known sequential algorithm in a best-case by 50% (see Section 4.1).
- Compared to the best available multithreaded algorithm, our GPU implementation is up to 24× faster, and our multithreaded implementation is within 0.5-2×. This results in a best-case performance of 270 *MFeatures/sec*. In contrast, the best sequential algorithm achieves 1.2 *MFeatures/sec* on the latest AMD EPYC 7763 and 16 *MFeatures/sec* by an efficientmultithreaded implementation(See section 4.2).
- We show that our proposed algorithm gracefully handles certain non-Euclidean distances. Specifically, we show that our algorithm is efficient when used with *mutual reachability distance*, a variant of Euclidean distance used in a popular clustering algorithm HDBSCAN^{*} [20].
- We provide a comprehensive set of experiments on three architectures to establish or provide empirical evidence for several properties of our algorithm and implementation including performance portability, asymptotic linear cost growth with problem size and lower threshold problem size to achieve performance saturation relative to CPU.

Algorithm 1 Borůvka's algorithm		
1:	procedure BORUVKA($G = \{V, E, W\}$)	
2:	$T \leftarrow (V, \emptyset) \triangleright$ initialize graph T with vertices from V and no edges	
3:	while <i>T</i> has more than one connected component do	
4:	for all components C of T do	
5:	$S \leftarrow \emptyset$	
6:	for all vertices v in C do	
7:	$D \leftarrow \{a \in E \mid a = (v, w), w \notin C\}$	
8:	$e \leftarrow \text{minimum}$ weight edge in D	
9:	$S \leftarrow S \cup e$	
10:	$e \leftarrow \text{minimum}$ weight edge in S	
11:	Add e to the graph T	

The remainder of the paper is organized as follows. Section 2 introduces the Borůvka algorithm and gives an overview of the related work. Section 3 describes the proposed algorithm. We compare our implementation with the state-of-the-art CPU implementations and demonstrate the algorithm's performance on GPUs in Section 4.

2 BACKGROUND

In this Section, we will briefly review the background and the related work relevant to subsequent discussion. **Notations:** Let $G = \{V, E, W\}$ be a weighted undirected connected graph. Here, V is the set of vertices of size n, E is the set of edges of size m, and W are the weights of the edges. We will use a *component* to describe a subset of of V with its meaning clear from the context.

MST computation. For a connected graph, the MsT is the tree subgraph with the least sum of edge weights. MsT is unique if all edge weights are distinct. Many algorithms to compute MsT are based on a greedy approach, using the fact that the minimum weight edge in any edge cut will be in MsT if it is unique (if it is not unique, any one of the edges with minimum weight can be chosen). The three most popular algorithms for computing MsT are Borůvka's algorithm [5], Prim's algorithm [24], and Kruskal's algorithm [15].

Prim's algorithm operates on a single component. At the start, the component is assigned a single vertex. On each step of the algorithm, the component is expanded by adding a vertex connected by an edge of the minimum weight in the component's cut. Prim's algorithm has $O(m \log n)$ complexity, and is inherently sequential.

In the Kruskal's algorithm, each vertex is initially assigned to its own component, and all edges are sorted by weights. On each step, the edge with the minimum weight is chosen among all edges that have the vertices in different components. The components of the vertices of that edge are then merged together. Kruskal's algorithm has $O(m \log n)$ complexity. It allows for a limited parallelism which is insufficient for a GPU.

Borůvka's algorithm. Borůvka's algorithm was one of the first published algorithms to compute an MsT. Similarly to the Kruskal's algorithm, it maintains a set of components, each initially containing a single vertex. Similarly to the Prim's algorithm, each component is expanded through finding the minimum weight edge in its cut. Unlike Prim's algorithm, however, the computation of the minimum weight edges can be done in parallel, and components are expanded through merging components together, rather than adding a single vertex. At the start of the Borůvka's algorithm (Algorithm 1), each component is initialized with an individual vertex, $C_i = \{v_i\}$. On each step, the algorithm determines the edge with the minimum weight in the cut of each active component. In other words, for a component C_i we find an edge $e_i = (v, u) \in E$, $v \in C_i$ and $u \in C_j$, $j \neq i$, with the minimum weight. We will call such edge e_i the *smallest outgoing edge* for the component C_i , and denote $C_i \rightarrow C_j$. The found edge is added to the list of edges in the MsT, and the two components C_i and C_j are merged, $C_i \leftarrow C_i \cup C_j$.

It is not guaranteed that the smallest outgoing edge for C_i with an end in C_j would be the smallest outgoing edge for C_j . Instead, the found edges result will typically produce a chain of components $C_{i_1} \rightarrow \cdots \rightarrow C_{i_{s-1}} \leftrightarrow C_{i_s}$. Each such chain terminates in a pair of components with their smallest outgoing edges pointing to each other. All components belonging to the same chain can be merged together in the same Borůvka iteration. In practice, this results in the Borůvka's algorithm requiring far fewer iterations compared to its theoretical upper bound of $\lceil \log_2(n) \rceil$.

Figure 2 demonstrates the steps in a single Borůvka iteration. At the beginning of the k-th iteration, we have five components (Figure 2b). First, each point finds the closest neighbor belonging to a different component than its own (Figure 2c). This forms a set of candidate edges for each component. Then, we choose the shortest candidate edge for each component (Figure 2d) and add it to a set of found MsT edges. Finally, the newly found MsT edges connect previously disconnected components. To merge the components, we compute the new component label for each point. We can see that one of the new components was formed by merging three components.

Borůvka's algorithm is guaranteed to converge (i.e., produce a correct MsT) only when all edge weights are distinct. Otherwise, the found edges may result in a cycle. This situation may be avoided by a suitable tie-breaking resolution when selecting the smallest outgoing edges. One of the ways to achieve that is by using indices of the vertices for the comparison of the edges. For example, given two edges $e_1 = (v_1, w_1)$ and $e_2 = (v_2, w_2)$ of the same weight, one could define $e_1 < e_2$ if min $(v_1, w_1) < \min(v_2, w_2)$, or min $(v_1, w_1) = \min(v_2, w_2)$ and max $(v_1, w_1) < \max(v_2, w_2)$.

The parallel nature of the Borůvka's algorithm make it well suited for a GPU implementation.

EMST computation. Given a set of points *X* in a *d*-dimensional space, Euclidean minimum spanning tree (EMST) is defined as an MST of its *distance* graph. The distance graph \mathcal{D} of *X* is a complete graph, with each vertex corresponding to a point in *P*, and each edge $e_{ij} = (p_i, p_j)$ having the weight $w_{ij} = ||p_i - p_j||_2$. Explicitly computing and storing \mathcal{D} is undesirable as it requires $O(n^2d)$ operations and $O(n^2)$ storage, which is prohibitively expensive for large datasets. For that reason, it is usually used implicitly.

As the complexity of the MST algorithms is at least linear in the number of edges, regular MST algorithms are not suitable for the EMST problem as they would have quadratic complexity with respect to the number of points $O(n^2)$.

For the two-dimensional case, MsT calculation can be performed on a *Delaunay* triangulation of the points, which only has O(n)edges. However, Delaunay triangulation worst-case complexity grows from $O(n \log n)$ in the two-dimensional case to $\Theta(n^2)$ for higher dimensions.

Instead, EMST algorithms combine a general MST algorithm with a data structure to accelerate the search for the nearest neighbors. Bentley and Friedman [4] proposed the first such EMST algorithm using a k-d tree-based nearest neighbor searches together with Prim's algorithm. The authors estimated $O(n \log n)$ operations for most distributions of points, albeit not rigorously. A key limitation of this approach is that it will often perform many redundant distance computations. This stems from the iterative nature of MST algorithms, where the nearest-neighbor queries can be run multiple times for the same points.

Pruning the number of the redundant distance computations for EMST was explored in many works. The two popular strategies emerged: the *well-separated pair decomposition* (WSPD) [6], and the *dual-tree* algorithms [18].

A pair of sets of points (P, Q) is called *well-separated* if the shortest distance between any point in P to any point in O is greater than the diameter of both of the sets. For a given set of points, WSPD is defined as a sequence of well-separated pairs (P_i, Q_i) such that for any pair of points $p, q \in X$ there exists a well-separated pair (P_k, Q_k) with $p \in P_k$ and $q \in Q_k$. With WSPD, EMST computation can be reformulated as a computation of the bichromatic closest pair (BCP) [2] between the well-separated pairs, and performing an MST computation using the found BCP edges. The first algorithm based on this approach was proposed in Agrawal et. al. [2]. Narasimhan [23] proposed GeoMST which combined WSPD and BCP with the Kruskal's algorithm. The algorithm was improved further by computing some BCP lazily or avoiding them altogether. Recently, Wang et.al [27] developed a parallel shared-memory variant based on this approach. To our knowledge, this is currently the fastest sequential and multithreaded parallel open-source implementation. The algorithm proposed in [27] algorithm was also shown to work with certain non-Euclidean distance metrics, such as the mutual-reachability distance for computing HDBSCAN^{*} [7].

March et al [18] proposed an EMST algorithm based on the dualtree framework. Unlike the single tree algorithm of Bentley and Friedman, where the nearest neighbor queries are performed separately for every point, the dual-tree algorithm performs such a query for a subtree in the spatial search tree. The algorithm used the component-wide upper and lower distance bounds during the tree traversal to avoid unnecessary distance computations. Under certain assumptions on the distribution of points, dual-tree has the best worst-case asymptotic complexity. In [19], researchers used the algorithm for non-Euclidean mutual-reachability distance of the HDBSCAN* algorithm.

Kokkos. Kokkos [9, 26] is a performance-portable programming model. It provides abstractions for expressing several parallel execution patterns such as parallel_{for,reduce,scan}. These patterns take function objects (*e.g.*, C++ lambdas) as arguments to execute for a given kernel index. While fairly restricted, this programming models allows maximum flexibility for mapping the patterns to an execution model. To this end, Kokkos provides an *execution space* abstraction that represents an execution resource, and a *memory space* that represents an abstract memory resource. A user is required to make sure that an execution space has access to the memory space that the data is in. For example, if the execution space is Kokkos::Cuda and the data is on the host (Kokkos::HostSpace memory space), an explicit data transfer is required to put the data on the device (Kokkos::CudaSpace memory space).

Kokkos also provides an abstraction for a multi-dimensional array data structure called View. It is a polymorphic structure, whose layout depends on the memory the data resides in (host or device). For example, a one-dimensional view on a GPU would automatically result in a coalesced data access pattern.

Together, these abstractions are implemented in a C++ library¹. Kokkos supports multiple backends, allowing the code written in Kokkos to run on a variety of hardware. Pertinent to this work, Kokkos supports Nvidia GPUs through the Cuda backend, AMD GPUs through the HIP backend, serial host through the Serial backend, and parallel host through the OpenMP backend.

ArborX. ARBORX [16] is a performance-portable geometric search library based on Kokkos. At its core, ARBORX implements a highly efficient parallel data structure, bounding volume hierarchy (BVH), to allow fast computation of the two types of the search queries: spatial (*e.g.*, searching for all objects within a certain distance of an object of interest) and nearest (*e.g.*, searching for a certain number of the closest objects regardless of their distance from an object of interest).

ARBORX implements a linear BVH structure following the works [3, 13], which has been shown to perform well for low-dimensional data on GPUs. The user data is linearized using a space-filling curve (Z-curve) to improve the locality of the geometric objects during the construction. It is then followed by a fully parallel bottom-up construction algorithm to produce a binary tree structure (hierarchy). Given *n* data points, the resulting tree would have n - 1 internal nodes and *n* leaf nodes, for a total of 2n - 1 nodes. This very fast construction algorithm produces a tree of sufficient quality in most situations.

During the search (also called a traversal), each thread is assigned a single query, and all the traversals are performed independently in parallel in a top-down manner. To reduce the data and thread divergence, the queries are pre-sorted with the goal to assign neighboring threads the queries that are geometrically close.

3 ALGORITHM

Our algorithm follows the general approach of combining a classical MsT algorithm with an efficient data structure for finding nearest neighbors for the components. We use the Borůvka's algorithm as it exposes the most parallelism out of all classical algorithms (see Section 2).

In this work, we use a single tree algorithm for two reasons. First, a parallel implementation of the dual-tree algorithms on GPU accelerators is an open research problem, and a high-performance implementation is a significant challenge. Second, a single-tree approach is much easier to implement by reusing efficient parallel geometric search algorithms, allowing to tap into existing efficient GPU implementations. As we will demonstrate in Section 4, a single tree implementation works well in practice.

The Borůvka's algorithm is iterative in nature (see Section 2 for an overview). Figure 3 provides a high level overview of our implementation, with a detailed description provided later in this

¹https://github.com/kokkos/Kokkos

ICPP '22, August 29-September 1, 2022, Bordeaux, France

```
// ExecutionSpace is the Kokkos execution space
// (where a kernel is executed). MemorySpace is
// the Kokkos memory space (where the data resides).
ExecutionSpace exec_space;
Kokkos::View<int*, MemorySpace> labels("labels", n);
// Initialize labels by placing each vertex into a
// separate component
Kokkos::parallel_for(
  Kokkos::RangePolicy<ExecutionSpace>(exec_space, 0, n),
  KOKKOS_LAMBDA(int i) { labels(i) = i; });
// Construct BVH
ArborX::BVH<MemorySpace> bvh(exec_space, data);
// Perform Boruvka iterations
int num_components = n;
do {
  // Propagate leaf node labels to the internal nodes
  // [parallel_for]
  reduceLabels(exec_space, bvh, labels);
  // Compute upper bounds on the length of the shortest
  // outgoing edge for each component [parallel_for]
  Kokkos::View<float*, MemorySpace> upper_bounds =
    computeUpperBounds(exec_space, bvh, labels);
  // Find the shortest outgoing edge for each component
  // [parallel for]
  Kokkos::View<Kokkos::Pair<int,int>*, MemorySpace>
    component_out_edges =
        findComponentsOutgoingEdges(exec_space, bvh,
                                    labels, upper_bounds);
  // Merge components using the found edges through
  // updating the labels [parallel_for]
  num_components =
    mergeComponents(space, component_out_edges, labels);
```

} while (num_components > 1);

Figure 3: The single-tree EMST algorithm C++ implementation using ArborX and Kokkos.

Section. Each iteration consists of two phases. In the first phase, we find the shortest outgoing edge for each component

(findComponentsOutgoingEdges). Using these edges, the components are merged in the second phase (mergeComponents). We will now describe the algorithms for both phases.

Finding the shortest outgoing edge

We will denote by C_i^k the *i*th component on the *k*th Borůvka iteration. At the start of the Borůvka's algorithm, each component is initialized with an individual vertex, $C_i^0 = \{v_i\}$. As the algorithm proceeds, the components are merged together using the found edges.

Let $C^k = \{C_i^k\}_{i=1}^{s_k}$ be the set of components on iteration $k, C_i^k \cap C_j^k = \emptyset$ for $i \neq j$. The goal of this phase of the algorithm is to find edges $e_i^k, i = 1, \ldots, s_k$ such that $e_i^k = \arg\min\{\|(u_i^k, v_i^k)\| \mid u_i^k \in C_i^k$ and $v_i^k \in C_j^k, j \neq i\}$. The component C_j^k is the closest component to C_i^k , and we denote this relationship by $C_i^k \to C_j^k$.



Figure 4: Propagation of leaf node labels to the internal nodes. Gray denotes invalid labels.

This problem can be seen as the nearest neighbor problem with an additional constraint that the nearest neighbor of a point must belong to a component different from the one the point belongs to, followed by choosing the shortest edge for all points in a component. Thus, we can follow a general approach to solving nearest neighbor problem on GPUs. Using ARBORX, this is done by assigning each point to a single thread and executing the neighbor searches in bulk (i.e., with all threads launching at the same time). Each thread executes a stack-less top-down traversal.

One of the challenges in designing an efficient algorithm lies in that the components grow in size with each Borůvka iteration. Examining all nearest points regardless of their component membership becomes progressively more expensive. Without trimming the number of the distance computations, this leads to $O(n^2)$ cost on the later Borůvka iterations.

Thus, we propose two optimization procedures to maintain a moderate cost of each tree traversal regardless of the component size.

Optimization 1: subtree skipping. We focus on reducing the number of the tree nodes encountered during the traversal by each thread. Specifically, individual thread skips the subtrees where each leaf node belongs to the same component as the point assigned to the thread. A similar approach was proposed in [19] in the context of dual-trees. While the benefit of this approach is limited on the earlier iterations of the algorithm, when the components are small, it is critical on the later iterations. In our experience, the cost of Borůvka's iterations tends to progressively decrease, with later iterations typically taking a small fraction of the earlier ones.

Our implementation uses a flat array of size n, called *labels*, to indicate a membership of a point in a component². As each point in the dataset is also a leaf in the constructed tree, we can associate each leaf node with a label of its component.

Before running the nearest neighbor algorithm, we propagate the labels from the leaf nodes to the internal nodes (reduceLabels in Figure 3). For a binary tree-based index, such as our case, this is done in a single bottom-up traversal algorithm (parallel_for). Each thread is assigned a leaf node, and traverses up the tree. The first thread accessing an internal node stores its label and terminates, while the second thread combines both labels, updates the internal node's label and continues upwards. If the labels of the children of an internal node are the same, the same label is assigned to that internal node. Otherwise, the internal node is assigned an invalid

 $^{^2\}mathrm{The}$ content of the array changes on each Borůvka iteration, as the components are merged together.

Prokopenko et al.

Algorithm 2 Optimized single-thread nearest neighbor traversal algorithm for a given thread with index *i*.

1:	$point \leftarrow data(i)$ \triangleright data point assigned to a thread
2:	$component \leftarrow labels(i) \qquad \triangleright$ component that the point belongs to
3:	$radius \leftarrow upper_bounds(component)$
4:	$shortest_distance \leftarrow \infty$ \triangleright the best found distance
5:	$closest_neighbor \leftarrow \emptyset$
6:	Initialize stack with the root node
7:	while stack is not empty do
8:	Pop the stack and assign it to <i>node</i>
9:	<pre>if distance(point, node) > radius then</pre>
10:	continue
11:	for all children <i>child</i> of <i>node</i> do
12:	$d \leftarrow distance(point, child)$
13:	if $labels(child) \neq component$ and $d \leq radius$ then
14:	if <i>child</i> is a leaf node then
15:	if d < shortest_distance then
16:	$shortest_distance \leftarrow d$
17:	$closest_neighbor \leftarrow child$
18:	$radius \leftarrow d$
19:	else
20:	Insert <i>child</i> into the stack
21:	Update the component's shortest outgoing edge if necessary

label to indicate that the corresponding subtree has leaf nodes from multiple components. Figure 4 shows an example of internal node labels based on the labels of the leaf nodes. We see that the leaf nodes belong to three different components, and that there are two subtrees (green and orange) containing the leaf nodes belonging to a single component. A thread performing the nearest neighbor search for a point with a green label will skip the subtree with the root at internal node 3.

Optimization 2: upper bounds for the outgoing edges. Further improvements may be achieved by using the fact that we are looking for the closest neighbor for *all* points in a component. The distance to an encountered neighbor of one point automatically provides an upper bound on the shortest outgoing edge for the full component. In the extreme case, this upper bound can be updated each time a thread encounters a leaf node.

However, we take a more moderate approach. We observe that if two points belong to different components, we can use the distance between them as an upper bound for the shortest outgoing edge for both of the components. For this bound to be useful, it is also desirable that these two points are close to each other. In general, it is not a trivial task to find such pairs. However, one of the steps in constructing a linear BVH is sorting data along a space-filling curve (typically, Z-curve using Morton indices). We then use any neighboring pair of points on the curve with different labels to initialize the upper bounds for the components (computeUpperBounds in Figure 3). This works well in practice as a pair of points with close Morton indices are likely to be close geometrically.

Traversal algorithm. Algorithm 2 shows the pseudo-code of the nearest neighbor algorithm executed by an individual thread. On line 3, the cutoff radius is set to the upper bound distance for the component. If the distance from a data point to a bounding volume of a tree node less than the current value of the cutoff radius, the children of the node are examined (line 11). We check that there is

at least one node belonging to a different component in the subtree with *child* as root, and that the bounding volume of the *child* node is within the cutoff distance (line 13). If a child node is a leaf node closer than the closest neighbor found so far (line 15), we update the cutoff radius value and the closest neighbor values. Otherwise, if the child is an internal node, it is inserted into the stack for the later examination (line 20). Finally, once the closest neighbor this point is found, we compare and update the component's shortest outgoing edge if necessary (line 21).

Non-Euclidean metrics. We also note that while we described the procedure for Euclidean distance, it will also work for certain other metrics. In particular, the mutual reachability metric used in HDBSCAN* can be integrated with a regular nearest neighbor traversal. The only change to the algorithm is that the cutoff radius during the traversal is set to the mutual reachability distance instead of the regular Euclidean distance. This is made possible by the fact that the mutual reachability distance is always greater or equal to the Euclidean one, and thus nodes truncated by the mutual reachability distance.

Merging components together

In the second phase, we use the edges found in the first phase to merge components together. As mentioned in Section 2, a single iteration of the Borůvka's algorithm results in chains of components. The merge procedure is straightforward and is embarrassingly parallel. For every point, we follow the chain until reaching the terminal pair of components with their shortest outgoing edges pointing to each other, and update the value of the labels array to be the component with the smallest index of that pair.

4 EXPERIMENTAL RESULTS

In our implementation, we used ArborX [16], an open-source library for the tree-based implementations, and Kokkos library [26] for a device-independent programming model. The implemented algorithm is available in the main ArborX repository³.

For our rate metric, we used the number of features processed per second, nd/t, where *n* is the number of points in the dataset, *d* is the dimension, and *t* is the time taken. We also denote by *MFeatures/sec* the number representing millions of features processed per second. We chose to include the dimension in our rate metric to allow cross-dimensional comparison of the datasets.

Testing environment. The numerical studies presented in the paper were performed using AMD EPYC 7763 (64 cores), Nvidia A100 and a single GCD (Graphics Compute Die) of AMD MI250X⁴. The chips are based on TSMC's N7+, N7 and N6 technology, respectively, and can be considered to belong to the same generation.

We used GCC 11.2.0 compiler for AMD EPYC 7763, NVCC 11.5 for Nvidia A100, and ROCm 4.5 for AMD MI250X.

Datasets. For our experiments, we used a combination of artificial and real-world datasets:

³https://github.com/arborx/ArborX

⁴Currently, HIP (Heterogeneous-computing Interface for Portability) – the programming interface provided by AMD – only allows the use of each GCD as an independent GPU.

ICPP '22, August 29-September 1, 2022, Bordeaux, France



Figure 5: Performance comparison of the sequential EMST implementations on AMD EPYC 7763.



Figure 6: Performance comparison of the parallel EMST implementations using AMD EPYC 7763, Nvidia A100 and AMD MI250X (single GCD).

- Ngsim (2D) [1] consists of ~12M 2D points corresponding to car trajectories on three highways. We also use one of these highways as a separate dataset Ngsimlocation3.
- **PortoTaxi** (2D) [21] consists of 1,710,000+ trajectories with ~81M 2D points in total, corresponding to the trajectories of several hundred taxis operating in the city of Porto, Portugal.
- RoadNetwork3D (2D) [14] consists of ~400K 2D points of the road network of the North Jutland province in Denmark.
- GeoLife24M3D (3D) [28] consists of ~24M 3D points corresponding to a user location data (longitude, latitude, altitude), and has a very skewed distribution.
- Hacc37M (3D) and Hacc497M (3D) consist of the 3D data taken from a single rank of a cosmology simulation performed with HACC [11]. *Hacc37M* was taken from a 1024³ particles simulation, and has ~37M points. *Hacc497M* was taken from a 3072³ particles simulations, and has ~497M points.
- VisualVar10M2D (2D) and VisualVar10M3D (3D) were produced by the generator of [10]. Both datasets are of size 10M.
- Normal100M2 (2D), Normal300M2 (2D), Normal100M3 (3D), (Normal200M3 (3D) consist of randomly generated points with zero mean and one standard deviation in all the dimensions. The dataset sizes are 100M, 300M, 100M, 300M, respectively.
- Uniform100M2 (2D) and Uniform100M3 (3D) are randomly generated datasets where all the points are distributed

uniformly inside a unit square (cube) in 2D (3D), both centered at the origin. Both datasets are of size 100M.

Competing Algorithms: We compare the performance of our algorithm to MLPACK [8] implementation of the dual-tree algorithm [18] available at https://github.com/mlpack/mlpack, and to MEMOGFK implementation of the [27] available at https://github.com/wangyiqiu/hdbscan.

4.1 Sequential performance

Our first goal is to compare the sequential performance of the implementations. Figure 5 shows the results comparing MLPACK, MEMOGFK and ARBORX on a variety of datasets using a single thread on AMD EPYC 7763.

We observe that MLPACK is slower than MEMOGFK for all the datasets. The sequential performance of our algorithm is competitive for most datasets, and is $1.5 \times$ faster than MEMOGFK for the Ngsimlocation3. The only outlier is the GeoLife24M3D. Our investigation showed that the properties of that dataset make it challenging to construct a high quality BVH. Specifically, the extremely high density of certain regions is under-resolved by the space-filling curve, resulting in significant bounding volume overlaps among nodes of certain subtrees. We believe that this issue can be addressed by increasing the resolution of the Z-curve grid, e.g., by using 128-bit Morton codes instead of 64-bit ones.

An interesting observation is that the performance of all implementations seem to be dimension-agnostic, as the rates are similar between 2D and 3D datasets.



Figure 7: Effect of the dataset size on the parallel performance using AMD EPYC 7763 and Nvidia A100.

4.2 Parallel performance

We now compare the parallel performance of the best multi-threaded implementation MEMOGFK using AMD EPYC 7763 with the parallel CPU and GPU implementations of ARBORX run AMD EPYC 7763, Nvidia A100 and AMD MI250X (single GCD). The results are presented on Figure 6.

We observe that our ARBORX implementation achieves 45-270 *MFeatures/sec* on an Nvidia A100, and is faster by 4-24× than MEM-OGFK. The relative performance between different datasets observed on AMD MI250X is qualitatively similar to observed performance in Nvidia A100. For both AMD MI250X and Nvidia A100, we achieve the best performance for *Hacc37M*, and the worst performance for *GeoLife24M3D*. The ARBORX on a single GCD of AMD MI250X is faster by 2-12× than the multithreaded MEMOGFK on AMD EPYC 7763.

The good and bad cases are the similar between MEMOGFK and ARBORX. Both implementations achieve the best performance on *Hacc37M*, and the worst on *GeoLife24M3D* (see the discussion in the previous Section). We also observe lower performance of ARBORX on the *RoadNetwork3D*. This is caused by the smaller size of that dataset, which is not enough to fully saturate a GPU.

In general, there is little qualitative differences in performance between 2D and 3D datasets. In other words, performance has little variability with respect to the dimension of the data, but is more dependent on the distribution of points. One exception to that are the uniform datasets, where we see up to 20% reduced performance for the 3D datasets with respect to the 2D dataset.

We find that our ARBORX multi-threaded implementation achieves 10-17 *MFeatures/sec* on AMD EPYC 7763 (with an exception of the *GeoLife24M3D* dataset), which puts it within factor 0.5-2× of the MEMOGFK. A currently known limitation of the multi-threaded implementation is the poor scaling of the sort algorithm. The native multi-threaded Kokkos::BinSort showed very poor performance on some of the datasets, and was replaced by an std::sort, a serial sort from the standard C++ library. For larger datasets, the serial nature of this sort becomes a dominant cost. We look to replace it with a robust multi-threaded sort implementation in the future.

We also note that we have not used any architecture-specific optimization for any device, and that we do not attempt to study the impact of architectural differences in ARBORX performance. Nevertheless, we would like to make a qualifying remark about relative performance on AMD MI250X and Nvidia A100. We primarily used the Nvidia A100 for algorithm and software development, debugging and profiling. Doing so may result in *performance bias* for Nvidia A100 since our algorithmic design process was guided by performance hotspots observed on the Nvidia A100.

4.3 Scaling performance

We now explore the performance of the algorithms with respect to the number of points in a dataset. As all algorithms are sensitive to the distribution of points in a dataset, we try to maintain a given distribution by randomly sampling a large dataset a specified number of times, producing a subset with the same data distribution.

We show the results of the sampling experiment for three datasets in Figure 7. The performance of each algorithm increases with the number of samples until it reaches saturation. This empirically demonstrates the asymptotic linear complexity of the two algorithms. Otherwise, if the complexity was higher than linear, our metric would decrease with increasing number of samples.

We also observe that ARBORX seems to start peaking around 10⁶ mark. In contrast, MEMOGFK achieves its peak performance at much higher number of points. This is counter intuitive, as typically CPU algorithms reach peak performance at lower problem sizes compared to similar GPU algorithms.

4.4 Analysis of computational phases

We show the relative cost and scaling of the different computation phases for MEMOGFK and ARBORX algorithms.

МемоGFK algorithm consists of four phases: tree construction (T_{tree}) , WSPD calculation (T_{wspd}) , Kruskal's MST algorithm (T_{mst}) , and auxiliary routines (T_{mark}) . Figure 8a shows the breakdown for MEMOGFK. The lower portion of each bar corresponds to the multi-threaded performance, while the full bar is the sequential performance. The numbers indicate the ratio between the two.

We see that in the sequential case, the costliest step is the computation of WSPD. However, WSPD calculation scales well with the number of cores and achieves the best case speed-up of $57 \times$ on 64 CPU cores. On the other hand, tree construction is not a bottleneck in the sequential case, but its poor scaling makes it the slowest phase of the EMST computation for many datasets.

ICPP '22, August 29-September 1, 2022, Bordeaux, France



(a) Breakdown of different phases of MEMOGFK and their speed-up over sequential on AMD EPYC 7763.



(b) Breakdown of different phases of ARBORX and their speed-up over sequential on Nvidia A100.

Figure 8: Breakdown of different phases of МемоGFK and ARBORX

ARBORX algorithm consists of only two phases: tree construction (T_{tree}) , and Borůvka's MST algorithm (T_{mst}) . Except for *RoadNetwork3D*, which is of small size, both phases scale well on GPU, and achieve the best speed-up of 360× and 350×, respectively.

4.5 Mutual reachability distance

HDBSCAN* [7] is a popular unsupervised clustering algorithm. Similarly to EMST, it seeks to construct an MST on a complete graph of a set of points. The main difference is that instead of using Euclidean distance, it uses the mutual reachability distance (m.r.d.). Given two points u and v, m.r.d. is defined as

$$d_{mreach}(u, v) = \max \{ d_{core}(u), d_{core}(v), \|u - v\|_2 \}.$$

Here, $d_{core}(u)$ is the *core distance*, defined as the distance to the k_{pts} th nearest neighbor (including the point itself), where k_{pts} is an input parameter to HDBSCAN*. When run with $k_{pts} = 1$, d_{mreach} is equivalent to the regular Euclidean distance.

Computing an MST in this scenario requires two changes to the regular EMST calculations. First, core distances have to be determined prior to running an MST algorithm. Second, an EMST algorithm must be modified to allow for a non-Euclidean distance metric. Both MEMOGFK and ARBORX (see Section 3) allow use of the m.r.d. metric.

In this Section, we would like to explore the effect of using m.r.d. with different values of k_{pts} on both the runtime and relative speedup of MEMOGFK and ARBORX implementations. Figure 9 shows the effect of varying values of k_{pts} on the runtime of the implementations for two datasets. T_{core} and T_{emst} denote the time to compute core-distances and the total time to compute MST with m.r.d, respectively.

We first observe that increasing values of k_{pts} results in growth of T_{core} . This is entirely expected as more neighbors are to be found. However, the kernel cost grows faster in the ARBORX implementation on GPU compared to the MEMOGFK on CPU. For example, for the *Hacc37M* the speedup of ARBORX over MEMOGFK drops from 20 at $k_{pts} = 2$ to only 12.7 at $k_{pts} = 16$. This is likely caused by the cost of thread divergence when maintaining priority queues for every thread.

The increase in T_{emst} is partially caused by the increase in T_{core} . The cost of the Borůvka iterations kernel is less clear. The difference between m.r.d. and Euclidean distance only affects earlier Borůvka iterations, when the distances to the closest neighbors are smaller than their core distances. Thus, many neighbors for a given point will all have the same m.r.d. distance to it, resulting in more expensive neighbor searches. This effect disappears on the later Borůvka iterations when the Euclidean distance dominates. We have also observed that increasing k_{pts} will result in more components getting merged on the earlier Borůvka iterations. In general, the cost of that kernel does not increase much with k_{pts} , staying within 30% of $k_{pts} = 2$.

5 CONCLUSION

We presented a single-tree algorithm for the EMST problem designed to exploit the massively threaded parallelism available on GPUs. The key strength of our approach is its simplicity through the use of a single-tree traversal with certain optimizations to prune the neighbor search. We evaluated the sequential, multithreaded, and GPU versions of our approach using a variety of datasets on multiple hardware architectures including Nvidia and AMD GPUs. We demonstrated that it was performance portable across these platforms and its excellent performance on GPUs compared to the best multi-threaded implementation. We conclude that our approach is efficient for a low-dimensional data. It remains to be seen if our findings hold for the data of higher dimension, which we plan to explore in our future work.

ACKNOWLEDGEMENTS

This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of the U.S. Department of Energy Office of Science and the National Nuclear Security Administration.

This research used resources of the Oak Ridge Leadership Computing Facility at the Oak Ridge National Laboratory, which is supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC05-00OR22725.

REFERENCES

- 2018. Next Generation Simulation (NGSIM) Vehicle Trajectories and Supporting Data. Available online: https://catalog.data.gov/dataset/next-generationsimulation-ngsim-vehicle-trajectories-and-supporting-data. Accessed: 2021-03-06.
- [2] Pankaj K Agarwal, Herbert Edelsbrunner, Otfried Schwarzkopf, and Emo Welzl. 1991. Euclidean minimum spanning trees and bichromatic closest pairs. *Discrete & Computational Geometry* 6, 3 (1991), 407–422.

ICPP '22, August 29-September 1, 2022, Bordeaux, France

Prokopenko et al.





- [3] Ciprian Apetrei. 2014. Fast and Simple Agglomerative LBVH Construction. In Computer Graphics and Visual Computing (CGVC), Rita Borgo and Wen Tang (Eds.). The Eurographics Association. https://doi.org/10.2312/cgvc.20141206 ZSCC: NoCitationData[s0].
- [4] Bentley and Friedman. 1978. Fast Algorithms for Constructing Minimal Spanning Trees in Coordinate Spaces. *IEEE Trans. Comput.* C-27, 2 (Feb. 1978), 97–105. https://doi.org/10.1109/TC.1978.1675043 Conference Name: IEEE Transactions on Computers.
- [5] Otakar Borůvka. 1926. O jistém problému minimálním. Práce Mor. Prirodved. Spol. v Brne (Acta Societ. Scienc. Natur. Moravicae) 3, 3 (1926), 37–58.
- [6] Paul B Callahan and S Rao Kosaraju. 1995. A decomposition of multidimensional point sets with applications to k-nearest-neighbors and n-body potential fields. *Journal of the ACM (JACM)* 42, 1 (1995), 67–90.
- [7] Ricardo J. G. B. Campello, Davoud Moulavi, Arthur Zimek, and Jörg Sander. 2015. Hierarchical Density Estimates for Data Clustering, Visualization, and Outlier Detection. ACM Transactions on Knowledge Discovery from Data 10, 1 (July 2015), 5:1–5:51. https://doi.org/10.1145/2733381
- [8] Ryan R. Curtin, Marcus Edel, Mikhail Lozhnikov, Yannis Mentekidis, Sumedh Ghaisas, and Shangtong Zhang. 2018. mlpack 3: a fast, flexible machine learning library. *Journal of Open Source Software* 3, 26 (June 2018), 726. https://doi.org/ 10.21105/joss.00726
- [9] H. Carter Edwards, Christian R. Trott, and Thomas Sunderland. 2014. Kokkos: Enabling manycore performance portability through polymorphic memory access patterns. J. Parallel and Distrib. Comput. 74, 12 (Dec. 2014), 3202–3216. https: //doi.org/10.1016/j.jpdc.2014.07.003 Publisher: Academic Press.
- [10] Junhao Gan and Yufei Tao. 2017. On the Hardness and Approximation of Euclidean DBSCAN. ACM Transactions on Database Systems 42, 3 (July 2017), 14:1-14:45. https://doi.org/10.1145/3083897
- [11] Salman Habib, Adrian Pope, Hal Finkel, Nicholas Frontiere, Katrin Heitmann, David Daniel, Patricia Fasel, Vitali Morozov, George Zagaris, Tom Peterka, et al. 2016. HACC: Simulating sky surveys on state-of-the-art supercomputing architectures. *New Astronomy* 42 (2016), 49–65.
- [12] Michael Held and Richard M Karp. 1970. The traveling-salesman problem and minimum spanning trees. Operations Research 18, 6 (1970), 1138–1162.
- [13] T. Karras. 2012. Maximizing Parallelism in the Construction of BVHs, Octrees, and K-d Trees. In Proceedings of the Fourth ACM SIGGRAPH / Eurographics Conference on High-Performance Graphics (EGGH-HPG'12). Eurographics Association, Goslar Germany, Germany, 33–37. https://doi.org/10.2312/EGGH/HPG12/033-037
- [14] Manohar Kaul, Bin Yang, and Christian S Jensen. 2013. Building accurate 3d spatial networks to enable next generation intelligent transportation systems. In 2013 IEEE 14th International Conference on Mobile Data Management, Vol. 1. IEEE, 137–146.
- [15] Joseph B. Kruskal. 1956. On the Shortest Spanning Subtree of a Graph and the Traveling Salesman Problem. *Proc. Amer. Math. Soc.* 7, 1 (1956), 48–50. https: //doi.org/10.2307/2033241 Publisher: American Mathematical Society.
- [16] D. Lebrun-Grandié, A. Prokopenko, B. Turcksin, and S. R. Slattery. 2020. ArborX: A Performance Portable Geometric Search Library. ACM Trans. Math. Softw. 47, 1, Article 2 (Dec. 2020), 15 pages. https://doi.org/10.1145/3412558
- [17] Xiang-Yang Li and Peng-Jun Wan. 2001. Constructing minimum energy mobile wireless networks. ACM SIGMOBILE Mobile Computing and Communications Review 5, 4 (2001), 55–67.
- [18] William B March, Parikshit Ram, and Alexander G Gray. 2010. Fast euclidean minimum spanning tree: algorithm, analysis, and applications. In Proceedings of the 16th ACM SIGKDD international conference on Knowledge discovery and data mining. 603–612.

- [19] Leland McInnes and John Healy. 2017. Accelerated Hierarchical Density Based Clustering. In 2017 IEEE International Conference on Data Mining Workshops (ICDMW). 33–42. https://doi.org/10.1109/ICDMW.2017.12 ISSN: 2375-9259.
- [20] Leland McInnes, John Healy, and Steve Astels. 2017. hdbscan: Hierarchical density based clustering. J. Open Source Softw. 2, 11 (2017), 205.
- [21] Luis Moreira-Matias, Joao Gama, Michel Ferreira, Joao Mendes-Moreira, and Luis Damas. 2013. Predicting taxi-passenger demand using streaming data. *IEEE Transactions on Intelligent Transportation Systems* 14, 3 (2013), 1393–1402.
- [22] Krishna Naidoo, Lorne Whiteway, Elena Massara, Davide Gualdi, Ofer Lahav, Matteo Viel, Héctor Gil-Marín, and Andreu Font-Ribera. 2020. Beyond two-point statistics: using the minimum spanning tree as a tool for cosmology. *Monthly Notices of the Royal Astronomical Society* 491, 2 (2020), 1709–1726.
- [23] Giri Narasimhan, Jianlin Zhu, and Martin Zachariasen. 2000. Experiments with computing geometric minimum spanning trees. In *Proceedings of ALENEX'00*. Citeseer, 183–196.
- [24] R. C. Prim. 1957. Shortest connection networks and some generalizations. *The Bell System Technical Journal* 36, 6 (Nov. 1957), 1389–1401. https://doi.org/10.1002/j.1538-7305.1957.tb01515.x Conference Name: The Bell System Technical Journal.
- [25] S Subramaniam and SB Pope. 1998. A mixing model for turbulent reactive flows based on Euclidean minimum spanning trees. *Combustion and Flame* 115, 4 (1998), 487–514.
- [26] Christian R. Trott, Damien Lebrun-Grandié, Daniel Arndt, Jan Ciesko, Vinh Dang, Nathan Ellingwood, Rahulkumar Gayatri, Evan Harvey, Daisy S. Hollman, Dan Ibanez, Nevin Liber, Jonathan Madsen, Jeff Miles, David Poliakoff, Amy Powell, Sivasankaran Rajamanickam, Mikael Simberg, Dan Sunderland, Bruno Turcksin, and Jeremiah Wilke. 2022. Kokkos 3: Programming Model Extensions for the Exascale Era. *IEEE Transactions on Parallel and Distributed Systems* 33, 4 (April 2022), 805–817. https://doi.org/10.1109/TPDS.2021.3097283 Conference Name: IEEE Transactions on Parallel and Distributed Systems.
- [27] Yiqiu Wang, Shangdi Yu, Yan Gu, and Julian Shun. 2021. Fast parallel algorithms for euclidean minimum spanning tree and hierarchical spatial clustering. In Proceedings of the 2021 International Conference on Management of Data. 1982– 1995.
- [28] Yu Zheng, Like Liu, Longhao Wang, and Xing Xie. 2008. Learning transportation mode from raw gps data for geographic applications on the web. In *Proceedings* of the 17th international conference on World Wide Web. 247–256.