

Communication-Avoiding Linear Algebraic Kernel K-Means on GPUs

Julian Bellavita
Cornell University
Ithaca, NY, USA
jb2695@cornell.edu

Matthew Rubino
Cornell University
Ithaca, NY, USA
mrr245@cornell.edu

Nakul Iyer
Cornell University
Ithaca, NY, USA
npi2@cornell.edu

Andrew Chang
Cornell University
Ithaca, NY, USA
yc723@cornell.edu

Aditya Devarakonda
Wake Forest University
Winston-Salem, NC, USA
devaraa@wfu.edu

Flavio Vella
University of Trento
Trento, Italy
flavio.vella@unitn.it

Giulia Guidi
Cornell University
Ithaca, NY, USA
gguidi@cornell.edu

Abstract—Clustering algorithms are among the most important and informative data analysis tools available. Of these, the K-means algorithm is one of the most popular and powerful choices due to its simplicity and generality. However, K-means cannot compute non-linearly separable cluster boundaries, which limits its utility in certain cases. Kernel K-means is a variant of the K-means algorithm that can compute non-linearly separable cluster boundaries. Kernel K-means has significant computational and memory requirements because of the presence of a large kernel matrix which scales quadratically with dataset size. Prior work has accelerated Kernel K-means by formulating it using sparse linear algebra primitives and implementing it on a single GPU. However, that approach cannot run on datasets with more than approximately 80,000 samples due to limited GPU memory.

Clustering using Kernel K-means remains challenging because of its high computational and memory costs. In this work, we address this issue by presenting a suite of distributed-memory parallel algorithms for large-scale Kernel K-means clustering on multi-GPU systems. Our approach maps the most computationally expensive components of Kernel K-means onto communication-efficient distributed linear algebra primitives uniquely tailored for Kernel K-means, enabling highly scalable implementations that efficiently cluster million-scale datasets. Central to our work is the design of partitioning schemes that enable communication-efficient composition of the linear algebra primitives that appear in Kernel K-means.

Our 1.5D algorithm consistently achieves the highest performance, enabling Kernel K-means to scale to data one to two orders of magnitude larger than previously practical. On 256 GPUs, it achieves a geometric mean weak scaling efficiency of 79.7% and a geometric mean strong scaling speedup of $4.2\times$. Compared to our 1D algorithm, the 1.5D approach achieves up to a $3.6\times$ speedup on 256 GPUs and reduces clustering time from over an hour to under two seconds relative to a single-GPU sliding window implementation. Our results show that distributed algorithms designed with application-specific linear algebraic formulations can achieve substantial performance improvement.

I. INTRODUCTION

K-means is a widely used unsupervised learning algorithm that attempts to group similar points to identify patterns in data. It is applied in biology [40], economics [37], and computer vision [20]. However, it cannot capture non-linearly separable clusters, which can degrade clustering quality [21].

Kernel K-means is a K-means variant that can compute non-linearly separable clusters, often resulting in higher quality clustering for non-linearly separable data [23], [27], [58]. It achieves this by clustering points in a high-dimensional feature space, producing non-linear boundaries in the original space. This is enabled by the kernel trick, which computes inner products between points in feature space using a non-linear kernel applied to the original points. The resulting values are stored in a kernel matrix, \mathbf{K} , avoiding the need to explicitly project points and thus saving computation and memory.

A drawback of Kernel K-means is that its per-iteration time complexity increases quadratically with the number of points, making it costly even for medium-sized data. To address this, prior work demonstrates effective GPU parallelization by mapping the main clustering loop to sparse linear algebra and using dense linear algebra for kernel matrix computation [10]. The main limitation of this approach is that it cannot cluster datasets with more than approximately eighty thousand points due to limited GPU memory and the large size of the kernel matrix. Datasets such as MNIST8m [16] contain millions of points, and applications like molecular dynamics [23] and human activity recognition [33] often require clustering at this scale. As datasets in scientific and engineering domains continue to grow [52], single-GPU approaches to Kernel K-means are insufficient in these big data scenarios.

A common approach to scaling Kernel K-means is low-rank approximation, such as the Nyström method, which does not explicitly form the kernel matrix. The approximate approaches are generally more scalable than the exact approach, but they can yield poor outcomes for certain kernel matrices, such as those with slow spectral decay or high coherence [18], [46], [17]. Moreover, they often require extensive dataset- and k -dependent tuning, sometimes relying on access to the full kernel matrix [57]. In contrast, exact Kernel K-means can be used universally without tuning.

Here, we tackle exact Kernel K-means clustering on million-scale datasets through GPU acceleration and distributed memory, expanding the practical scale of Kernel K-means by

one to two orders of magnitude. Our approach relies on a mapping of Kernel K-means to communication-efficient sparse and dense linear algebra primitives. Through this mapping, we introduce new distributed Kernel K-means algorithms that reduce communication by explicitly exploiting domain-specific structure and effectively composing the distribution schemes of the different linear algebra primitives involved in Kernel K-means.

The linear algebraic formulation of Kernel K-means consists of a dense general matrix multiplication (GEMM) followed by multiple sparse-dense matrix multiplications (SpMM), where the dense matrix in the SpMM is the output of the GEMM. There is a rich design space of distributed SpMM and GEMM algorithms to consider, each with distinct distribution strategies and communication schedules [14], [38], [47], [56]. In this work, we develop 1D, Hybrid 1D, 1.5D, and 2D distributed Kernel K-means algorithms, each tailored to exploit domain-specific characteristics of Kernel K-means. Central to our approach is the composability, or fusion, of the GEMM and SpMM primitives. Rather than considering each operation in a vacuum, our most performant algorithm chooses partitioning schemes and communication schedules for both primitives that enable efficient communication when considered together. In addition, the linear algebraic formulation of Kernel K-means uses a sparse matrix \mathbf{V} with exactly one nonzero per column, which we leverage to design partitioning schemes and communication schedules that achieve perfect load balance while minimizing communication.

Our 1D algorithm efficiently executes the clustering loop but cannot compute the kernel matrix at scale. It is described as a baseline because its communication pattern resembles prior distributed Kernel K-means approaches [22], [55] that do not use distributed sparse primitives and are not publicly available, so a direct comparison is not possible. The Hybrid 1D algorithm addresses some scalability issues of the 1D algorithm, but it introduces significant additional memory and communication overhead, making it ineffective in some cases. The 2D algorithm is scalable and memory-efficient, but its 2D partitioning of the distance matrix introduces communication overhead when updating cluster assignments.

Our main contribution is the 1.5D algorithm, which communicates asymptotically less than other approaches and makes it feasible to use exact Kernel K-means to cluster million-scale datasets. The 1.5D algorithm is different from existing approaches to distributed Kernel K-means, which are largely similar to our 1D algorithm. A key novelty of the 1.5D algorithm is its choice of partitioning schemes for the input matrices, which enables efficient composition of GEMM and SpMM, allowing both to be computed in a communication-efficient manner. More broadly, our results highlight application-guided composition of distributed linear algebra primitives as a key research direction beyond individual scalable primitives.

Our algorithms were evaluated on three datasets and compared to a single-GPU sliding-window implementation. The 1.5D and 2D algorithms scale to 256 GPUs. The 1.5D algorithm outperforms the 1D algorithm by up to $3.6\times$ in

strong scaling and by over $2000\times$ compared to the sliding-window approach. Notably, both 1.5D and 2D can process over 1.5 million points without exhausting memory. Our algorithms are implemented in the open-source software VIVALDI (named after the composer) and are available at <https://github.com/CornellHPC/Vivaldi>. The sliding window algorithm is available at <https://github.com/aditya08/sliding-window-kernel-kmeans>.

In summary, our main contributions are:

- 1) The first open-source fully GPU-accelerated distributed memory implementation of exact Kernel K-means;
- 2) A set of new parallel algorithms for Kernel K-means based on distributed linear algebra primitives;
- 3) Our 1.5D algorithm achieves up to $3.6\times$ speedup over the baseline and makes clustering of million-scale datasets practical by increasing the dataset size limit by one to two orders of magnitude.

II. BACKGROUND

A. Kernel K-means Clustering

K-means is an unsupervised learning algorithm that partitions data into k clusters by assigning each point to the nearest centroid using a chosen distance metric, typically Euclidean [39]. Kernel K-means extends K-means to non-linearly separable data by performing clustering in a higher-dimensional feature space and mapping the resulting centroids back to the input space. Kernel K-means updates centroids iteratively by (1) computing point-centroid distances in feature space, (2) assigning each point to its nearest centroid, and (3) recomputing centroids as the mean of the assigned points. In practice, feature-space projections are never computed explicitly. Instead, Kernel K-means uses a kernel function $\kappa(x, y)$ to evaluate inner products in feature space while operating in the input space. Pairwise distances are derived from $\kappa(x, y)$, resulting in $\mathcal{O}(n^2)$ work per iteration, where n is the number of points in the dataset [21].

B. Kernel K-means with Sparse Linear Algebra

The state-of-the-art method for parallelizing Kernel K-means on a GPU formulates the problem using GEMM and the sparse linear algebra primitives SpMM and sparse matrix-vector multiplication (SpMV). The components of this formulation relevant to this work are summarized below; further details are available in [10].

Let $\mathbf{P} \in \mathbb{R}^{n \times d}$ be the (dense) point matrix, where n is the number of points and each row of d features represents one point. Then, we define the symmetric matrix $\mathbf{B} \in \mathbb{R}^{n \times n}$ as:

$$\mathbf{B} = \mathbf{P}\mathbf{P}^T \quad (1)$$

The kernel matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$ is obtained by applying an element-wise kernel function to \mathbf{B} . For example, if we define our kernel function κ to be the polynomial kernel:

$$\kappa(x, y) = (\gamma x^T y + c)^d \quad (2)$$

where $x, y \in \mathbb{R}^d$, then we have $\mathbf{K}(i, j) = \kappa(\mathbf{P}(i, :), \mathbf{P}(j, :))$, which is equivalent to $\mathbf{K}(i, j) = (\gamma \mathbf{B}(i, j) + c)^d$, therefore \mathbf{K}

can be obtained by applying κ to each entry of \mathbf{B} . In practice, \mathbf{K} is usually large, and the cost of storing or computing with \mathbf{K} is the main bottleneck in Kernel K-means.

Once \mathbf{K} is computed, the clustering loop can begin. Let k be the number of clusters, initialized by some strategy. The way we initialize the clusters does not affect the clustering loop computation; it only affects convergence properties [2]. Define the sparse matrix $\mathbf{V} \in \mathbb{R}^{k \times n}$:

$$\mathbf{V}(i, j) = \begin{cases} \frac{1}{|L_i|} & \text{if point } j \text{ is in cluster } i \\ 0 & \text{otherwise} \end{cases} \quad (3)$$

Let L_i denote the points in cluster i . Notably:

- 1) \mathbf{V} is a sparse matrix with exactly n nonzeros;
- 2) \mathbf{V} has exactly 1 nonzero per column.

First, each clustering iteration computes:

$$\mathbf{E} = \mathbf{K}\mathbf{V}^T \quad (4)$$

It then initializes $\mathbf{z} \in \mathbb{R}^n$ using a masking operation on \mathbf{E} :

$$\mathbf{z}(i) = \mathbf{E}(i, cl(i)) \quad (5)$$

$cl(i)$ is a function that returns the cluster to which point i is assigned. Then, a matrix-vector product is computed:

$$\mathbf{c} = \mathbf{V}\mathbf{z} \quad (6)$$

The vector \mathbf{c} is equal to a single row of the matrix $\tilde{\mathbf{C}} \in \mathbb{R}^{n \times k}$, defined as:

$$\tilde{\mathbf{C}} = \begin{bmatrix} \|c_1\|_2^2 & \|c_2\|_2^2 & \dots & \|c_k\|_2^2 \\ \|c_1\|_2^2 & \|c_2\|_2^2 & \dots & \|c_k\|_2^2 \\ \vdots & \vdots & \vdots & \vdots \\ \|c_1\|_2^2 & \|c_2\|_2^2 & \dots & \|c_k\|_2^2 \end{bmatrix} \quad (7)$$

Given that each row of $\tilde{\mathbf{C}}$ is identical, \mathbf{c} can serve as a compact representation of $\tilde{\mathbf{C}}$. The distance matrix $\mathbf{D} \in \mathbb{R}^{n \times k}$ is then computed as:

$$\mathbf{D} = -2\mathbf{E} + \tilde{\mathbf{C}} \quad (8)$$

$\mathbf{D}(i, j)$ stores the distance between point i and cluster j in feature space. Cluster assignments are updated via a row-wise argmin and used to update \mathbf{V} , after which the next iteration can begin.

In practice, (1) is computed using GEMM, (4) using SpMM, and (6) using SpMV. Thus, the main parts of Kernel K-means reduce to linear algebra primitives with efficient parallel implementations in libraries such as cuSPARSE and cuBLAS [44], [45]. Extending Kernel K-means to distributed memory is a non-trivial task. The main challenges are avoiding communication of \mathbf{K} partitions, which require $\mathcal{O}(n^2)$ bytes, and distributing $\mathbf{P}, \mathbf{K}, \mathbf{V}, \mathbf{E}$ to perform GEMM, SpMM, and cluster updates without excessive communication or load imbalance.

C. Distributed Linear Algebra Primitives

Prior work has developed parallel algorithms for linear algebra primitives on distributed memory. For GEMM, a widely used approach is the Scalable Universal Matrix Multiply Algorithm (SUMMA) [56], which partitions the input matrices \mathbf{A} and \mathbf{B} over a 2D process grid and performs GEMM as a sequence of distributed block outer products. If P is the total number of processes, SUMMA requires \sqrt{P} rounds of communication and local computation. At iteration p , the process with row rank p broadcasts its tile of \mathbf{A} along its process row, while the process with column rank p broadcasts its tile of \mathbf{B} along its process column. Each process then multiplies the received tiles and accumulates the result into its local tile of \mathbf{C} . If we assume a tree-based broadcast, the communication cost for square matrices of dimension n under the α - β model [31] is:

$$T_{SUMMA} = \alpha \mathcal{O}(\sqrt{P} \log(\sqrt{P})) + \beta \mathcal{O}(\log(\sqrt{P}) \frac{n^2}{\sqrt{P}}) \quad (9)$$

This is within a $\log(\sqrt{P})$ factor of the standard communication lower bound for dense linear algebra, assuming no additional memory is used for replication [8].

For SpMM, two widely used algorithms are 1.5D SpMM and 2D SpMM [47]. 2D SpMM is algorithmically identical to SUMMA, with the only difference being that the first operand \mathbf{A} is sparse. The different communication schedules for 2D SpMM are called A-stationary, B-stationary, and C-stationary. These schedules differ based on which matrices are communicated: A-stationary algorithms communicate only partitions of \mathbf{B} and \mathbf{C} ; C-stationary communicates only \mathbf{B} and \mathbf{A} , and so on. 1.5D SpMM is a communication-avoiding algorithm that replicates partitions of one operand across multiple processes, reducing communication costs at the expense of higher per-process memory requirement. In the literature, only A-stationary variants of 1.5D SpMM that replicate partitions of \mathbf{B} are described [47].

For GEMM and SpMM, there are 3D algorithms that divide the matrices on a 3D process cube [53], [1] and communicate asymptotically less than 2D or 1.5D algorithms. 3D algorithms are not considered in this paper for the following reasons. Using 3D GEMM to compute \mathbf{K} would require replicating tiles of \mathbf{K} , further increasing its memory footprint. Similarly, 3D SpMM is unlikely to improve the computation of \mathbf{E} , as it generally performs worse than 2D SpMM unless a very large number of processes is used [53].

III. RELATED WORK

The standard K-means algorithm has been used for several decades [39], and prior work has studied its implementation on various parallel architectures [7], [22]. To the best of our knowledge, Kernel K-means was first introduced by Girolami [26]. Previous work on high-performance Kernel K-means has explored parallelization in shared memory [9], [10] and distributed memory [23], [55], but has not leveraged distributed linear algebra primitives. To our knowledge, these prior implementations are not open-source, and none are

fully GPU-accelerated, making direct comparison difficult. The communication schedules of these approaches closely resemble those of our 1D algorithm, which is used as a practical baseline. In Section VI, we show that the 1D algorithm is almost universally outperformed by our other approaches.

To scale Kernel K-means, previous approaches stored the kernel matrix on disk [58], while low-rank approximations such as Nyström reduced memory usage [18]. However, these approximations can degrade clustering quality for kernel matrices with slow spectral decay or high coherence [17], [18], [46]. They also introduce dataset- and k -dependent tuning parameters that are difficult to optimize at scale. Fed-KKM [60] constructed low-dimensional random feature approximations on edge devices using communication-efficient Lanczos algorithms, reducing memory usage by up to 94% but increasing the number of iterations and computational overhead, with clustering quality dependent on approximation accuracy. Choosing appropriate approximations and tuning remains challenging, whereas exact Kernel K-means eliminates both issues by computing the solution directly, making it universally effective. Prior GPU-based approaches [9], [10] are limited to single-device memory, whereas our work scales exact Kernel K-means to million-point datasets, delivering high-quality clustering without tuning.

In addition, algorithms such as spectral clustering [21], DBSCAN [25], [19], and k -NN clustering [41] can also capture nonlinear cluster structures, with k -NN graph-based clustering being particularly scalable by avoiding explicit kernel matrices. Kernel K-means is a foundational algorithm with a strong theoretical basis and widespread practical use [30], [57], [59], which is why it is the focus of this work.

Prior work has explored distributed-memory algorithms that use linear algebra, particularly sparse primitives, including graph algorithms [3], [4], [6], [15], [34], [35], [50], Markov clustering [5], GNN training [53], [54], and computational biology pipelines [11], [28], [29], [48]. The composition and optimization of sequences of linear algebra primitives has recently emerged as a key research area in high-performance computing. Bharadwaj et al. [12] develop communication-avoiding strategies for sequences of Sampled Dense-Dense Matrix Multiply (SDDMM) and SpMM operations in GNNs, showing that, in that application, performance depends on how communication schedules are composed across consecutive primitives. Compiler-based BLAS fusion [24] and DNN operator fusion [43] similarly demonstrate that optimizing primitives in isolation can leave substantial performance unrealized. In Kernel K-means, we show that adapting GEMM and SpMM sequences to sparsity and primitive composition reduces communication and improves end-to-end performance, highlighting composability as a key research direction in high-performance linear algebra.

IV. ALGORITHM DESCRIPTIONS

Here, we describe our distributed-memory Kernel K-means algorithms. Let P be the total number of processes. For simplicity, assume that P evenly divides n , \sqrt{P} evenly divides

k , and that 2D process grids are $\sqrt{P} \times \sqrt{P}$, although only the square process grid assumption is required for correctness. In a 2D process grid, the i th process in column j is denoted $P(i, j)$. \mathbf{K}_{ij} is the portion of \mathbf{K} on $P(i, j)$, and \mathbf{c}_i is the portion of \mathbf{c} on the process with global rank i . For simplicity, we set $\mathbf{B} = \mathbf{K}$, corresponding to the linear kernel (i.e., standard dot product); this simplifies notation without affecting the algorithms. Communication costs are analyzed using the standard α - β model [31], where α represents message latency and β denotes network bandwidth.

Each algorithm is defined by its strategy for partitioning \mathbf{P} , \mathbf{K} , \mathbf{V} , and \mathbf{E} , as well as its approach to distributed GEMM and SpMM. Our focus is on performing the following operations in distributed memory:

- 1) Distributed GEMM: $\mathbf{K} = \mathbf{P}\mathbf{P}^T$
- 2) Distributed SpMM: $\mathbf{E}^T = \mathbf{V}\mathbf{K}$

\mathbf{E}^T is computed instead of \mathbf{E} because SpMM libraries typically require the first operand to be sparse; this does not affect the algorithm. GEMM multiplies a tall, narrow matrix by a short, wide matrix, while SpMM multiplies a sparse matrix with one nonzero per column by a large, square dense matrix.

Kernel K-means has application-specific characteristics that must be considered for efficient distributed-memory parallelization, including the unique sparsity structure of \mathbf{V} , the comparatively large size of \mathbf{K} , and the need to perform cluster updates at each iteration. To address these challenges, we present four algorithms that use different distribution strategies based on linear algebra.

In our 1D algorithm, all matrices are partitioned into 1D column blocks. This results in an algorithm that can compute the SpMM in a load-balanced manner without communicating \mathbf{K} , but the 1D layout of \mathbf{K} leads to high communication costs for both GEMM and SpMM. Two natural alternatives reduce GEMM communication: a Hybrid 1D algorithm that redistributes \mathbf{K} from 2D to 1D, and a 2D algorithm that modifies the partitioning of \mathbf{V} to be 2D. Both algorithms use SUMMA to compute \mathbf{K} , which reduces the communication costs of the GEMM. However, the Hybrid 1D algorithm incurs a high redistribution cost from 2D to 1D, while the 2D algorithm requires additional communication when updating cluster assignments.

The need to perform both a distributed GEMM and a distributed SpMM sequence motivates our main contribution: a new 1.5D algorithm. By using SUMMA for GEMM, the 1.5D algorithm preserves the natural 2D output partitioning of \mathbf{K} and combines it with a 1D partitioning of \mathbf{V} , enabling a communication-efficient algorithm for distributed SpMM without redistribution and allowing cluster updates to be computed without additional communication. This design gives 1.5D the lowest communication cost among the approaches considered.

A. 1D Kernel K-means

Our first algorithm divides matrices in a 1D columnwise manner. Each process owns a contiguous block of $\frac{n}{P}$ columns

	1D	Hybrid 1D	1.5D	2D
Kernel Matrix (\mathbf{K})	$\alpha\mathcal{O}(P) + \beta\mathcal{O}(Pnd)$	$\alpha\mathcal{O}(P) + \beta\mathcal{O}(\frac{n^2}{P} + \frac{nd}{\sqrt{P}})$	$\alpha\mathcal{O}(\sqrt{P}) + \beta\mathcal{O}(\frac{nd}{\sqrt{P}})$	Same as 1.5D
Distances Matrix (\mathbf{D}^T)	$\alpha\mathcal{O}(P) + \beta\mathcal{O}(n)$	Same as 1D	$\alpha\mathcal{O}(\sqrt{P}) + \beta\mathcal{O}(\frac{n(k+1)}{\sqrt{P}})$	$\alpha\mathcal{O}(\sqrt{P}) + \beta\mathcal{O}(\frac{n(k+1)}{\sqrt{P}} + n)$

TABLE I: Communication cost of \mathbf{K} and \mathbf{D}^T computation for each algorithm. $\log(\sqrt{P})$ terms are omitted for brevity.

Algorithm 1 1D Kernel K-means Algorithm

Require: $\mathbf{P} \in \mathbb{R}^{n \times d}$ stores points; number of clusters k .

- 1: Allgather \mathbf{P} on each process ▷ 1D GEMM
- 2: All processes p compute $\mathbf{K}_p = \mathbf{P}_p^T \mathbf{P}$
- 3: **while** not converged **do** ▷ Clustering loop
- 4: Allgather \mathbf{V} on each process
- 5: All processes p compute $\mathbf{E}_p^T = \mathbf{V} \mathbf{K}_p^T$
- 6: $\mathbf{z}_p = \text{mask}(\mathbf{E}_p^T)$
- 7: $\mathbf{c}_p = \mathbf{V}_p \mathbf{z}_p$
- 8: Allreduce \mathbf{c}_p to compute \mathbf{c} on each process
- 9: $\mathbf{D}_p^T = \text{distances}(\mathbf{E}_p^T, \mathbf{c})$
- 10: $cl_p = \text{argmin}(\mathbf{D}_p^T)$ ▷ Updated cluster assignment
- 11: Update \mathbf{V}_p using cl_p
- 12: **end while**

from each matrix. \mathbf{P}^T is initially partitioned as:

$$\mathbf{P}^T = [\mathbf{P}_1^T \quad \mathbf{P}_2^T \quad \dots \quad \mathbf{P}_P^T] \quad (10)$$

Therefore, a given tile $\mathbf{P}_i^T \in \mathbb{R}^{\frac{n}{P} \times d}$ is stored on process i . The first step of the algorithm computes $\mathbf{K} = \mathbf{P} \mathbf{P}^T$ using a 1D distributed GEMM. In this scheme, an Allgather collective replicates \mathbf{P} on each process. Each process then computes a local GEMM between its local partition of \mathbf{P}^T and the replicated \mathbf{P} , producing a block column of the kernel matrix \mathbf{K} . \mathbf{K} is partitioned columnwise as:

$$\mathbf{K} = [\mathbf{K}_1 \quad \mathbf{K}_2 \quad \dots \quad \mathbf{K}_P] \quad (11)$$

The main clustering loop then begins, with \mathbf{V} initially partitioned columnwise as:

$$\mathbf{V} = [\mathbf{V}_1 \quad \mathbf{V}_2 \quad \dots \quad \mathbf{V}_P] \quad (12)$$

Each partition of \mathbf{V} has exactly $\frac{n}{P}$ nonzeros, making communication inexpensive compared to the $\frac{n^2}{P}$ entries in each partition of \mathbf{K} . To compute \mathbf{E}^T , we use a 1D B-stationary distributed SpMM that communicates only \mathbf{V} . This requires a single Allgather to replicate \mathbf{V} across processes, followed by a local SpMM on each process to compute its corresponding partition of \mathbf{E}^T . As a result, \mathbf{E}^T is partitioned as:

$$\mathbf{E}^T = [\mathbf{E}_1^T \quad \mathbf{E}_2^T \quad \dots \quad \mathbf{E}_P^T] \quad (13)$$

Each process then performs the masking operation to create \mathbf{z}_p locally, followed by a local SpMV and a global Allreduce to obtain \mathbf{c} on every process. Finally, each process computes its partition of the distance matrix \mathbf{D}^T and updates its cluster assignments by performing a local argmin over the columns of \mathbf{D}_p^T . The 1D algorithm is summarized in Algorithm 1.

The Allgather in the 1D GEMM routine involves $\mathcal{O}(P)$ messages and sends a total of $\mathcal{O}(Pnd)$ words. Thus, the communication cost of computing \mathbf{K} is:

$$T_{K1D} = \alpha\mathcal{O}(P) + \beta\mathcal{O}(Pnd) \quad (14)$$

The Allgather used to compute \mathbf{E}^T replicates \mathbf{V} on each process. If we assume a pairwise exchange algorithm is used for the Allgather [51], this involves a total of $\mathcal{O}(P)$ messages. Since \mathbf{V} is sparse with exactly n nonzeros, the total number of words communicated is $\mathcal{O}(n)$, resulting in a total communication cost for computing \mathbf{E}^T of:

$$T_{E1D} = \alpha\mathcal{O}(P) + \beta\mathcal{O}(n) \quad (15)$$

The only other communication is the global Allreduce for \mathbf{c} , a vector of length k , which is negligible since k is small.

The main benefit of the 1D algorithm arises in the clustering loop, where columnwise partitioning of \mathbf{V} and \mathbf{K} provides perfectly load-balanced SpMM due to the uniform nonzeros in the replicated \mathbf{V} . Perfect load balance is typically difficult to achieve in distributed sparse linear algebra [14], [32], [47], making this a notable advantage of the 1D algorithm. In addition, updating cluster assignments requires no communication.

Overall, the communication cost of the 1D algorithm does not scale with the number of processes. The 1D GEMM requires an Allgather of $\mathcal{O}(Pnd)$ words, which becomes a bottleneck for large P . Similarly, the Allgather for \mathbf{E}^T sends $\mathcal{O}(n)$ words; although this amount does not increase as P increases, it may dominate when local computation becomes negligible or in a weak scaling setting. The 1D algorithm also has a large memory footprint: replicating \mathbf{P} on each process can cause out-of-memory errors when d is large, especially because partitions of \mathbf{K} must also be stored. The limitations of the 1D algorithm arise from the communication costs of both GEMM and SpMM.

B. Two Alternatives: Hybrid 1D and 2D

To reduce GEMM communication in the 1D algorithm, the 1D GEMM can be replaced with SUMMA [56], which computes distributed GEMM with a near-optimal communication cost. However, using SUMMA results in \mathbf{K} being partitioned in a 2D distribution. This prevents the use of the approach from the 1D algorithm during the main clustering loop, since \mathbf{K} is no longer distributed column-wise in 1D. Therefore, changes to the algorithm for computing the SpMM and cluster updates in the main clustering loop are necessary. Here, we consider two algorithms that implement these changes: a Hybrid 1D algorithm (H-1D) and a pure 2D algorithm.

In H-1D, SUMMA is used to compute \mathbf{K} . Then, \mathbf{K} is redistributed from SUMMA's 2D layout to a 1D column-wise

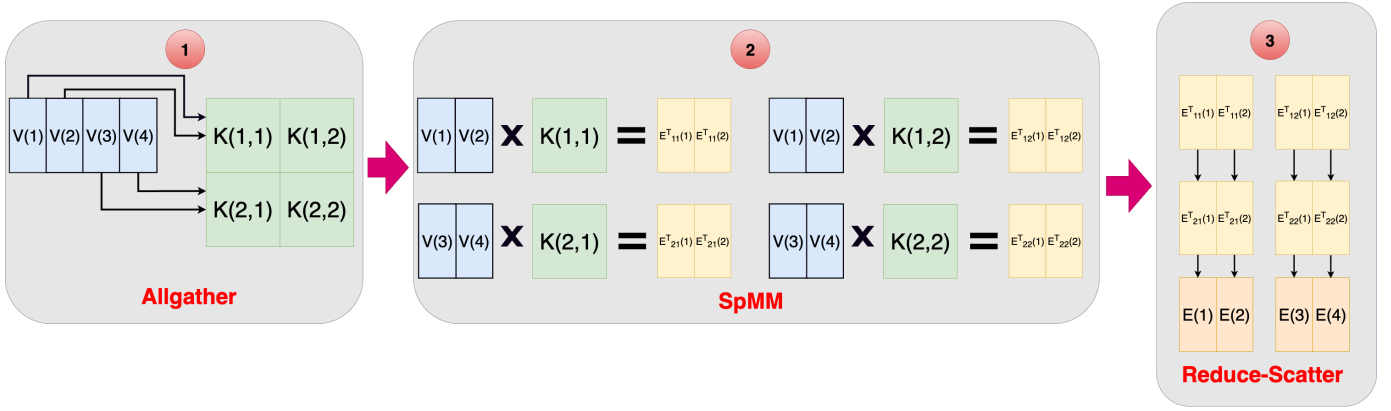


Fig. 1: 1.5D SpMM algorithm on $P = 4$ processes. \mathbf{V} is partitioned 1D columnwise and \mathbf{K} in 2D. (1) The nonzeros of each \mathbf{V} partition are replicated along the corresponding process row. (2) Each process performs a local SpMM with its \mathbf{V} replicas and local \mathbf{K} partition, producing partial sums of \mathbf{E}^T . (3) The partial sums are split along columns and reduced across process columns, resulting in a 1D columnwise partitioning of \mathbf{E}^T .

distribution, after which the main clustering loop is identical to that of the 1D algorithm. This redistribution is performed using an Alltoallv. Computing \mathbf{K} using SUMMA incurs the following communication cost:

$$T_{K-SUMMA} = \alpha \mathcal{O}(\sqrt{P} \log(\sqrt{P})) + \beta \mathcal{O}\left(\log(\sqrt{P}) \frac{nd}{\sqrt{P}}\right) \quad (16)$$

H-1D enables a more communication-efficient algorithm for performing GEMM, but it incurs additional communication costs due to the redistribution of \mathbf{K} . In particular, during the redistribution step, each process sends $\mathcal{O}(P)$ messages and communicates at most $\mathcal{O}(n^2/P^{3/2})$ words to up to \sqrt{P} processes, resulting in a worst-case redistribution cost of:

$$T_{Redist} = \alpha \mathcal{O}(P) + \beta \mathcal{O}\left(\frac{n^2}{P}\right) \quad (17)$$

This $\mathcal{O}(n^2/P)$ redistribution volume can make the H-1D algorithm communication costly unless a very large number of processes are used. Furthermore, the scalability issues of SpMM in the 1D algorithm remain present in H-1D.

The second alternative we consider is a pure 2D approach, where both \mathbf{V} and \mathbf{K} are partitioned across a 2D process grid. This enables the use of SUMMA to compute \mathbf{K} without needing to redistribute \mathbf{K} from a 1D to a 2D layout before starting the clustering loop. The computation of \mathbf{E}^T uses a B-stationary SpMM [47], communicating only \mathbf{V} entries and partial sums of \mathbf{E}^T . This results in a communication cost of:

$$T_{E2D} = \alpha \mathcal{O}(\sqrt{P}) + \beta \mathcal{O}\left(\frac{n(k+1)}{\sqrt{P}}\right) \quad (18)$$

The communication for computing \mathbf{E}^T now scales with the processes; however, the 2D partitioning of \mathbf{E}^T introduces extra communication for cluster updates. In particular, computing each point's nearest centroid requires an argmin reduction along process columns, and an Allreduce along process rows

to compute cluster sizes. This results in a total cluster update communication of:

$$T_{update} = \alpha \mathcal{O}(\log(\sqrt{P})) + \beta \mathcal{O}\left(\log(\sqrt{P})n\right) \quad (19)$$

The pure 2D algorithm has an asymptotically smaller communication volume than either 1D algorithm, but it requires communication during cluster updates, resulting in an overall communication cost higher than that of the 1.5D algorithm.

C. 1.5D Kernel K-means

Here, we present our 1.5D distributed Kernel K-means algorithm, which can compute both the GEMM and SpMM in a scalable fashion. The key idea behind our 1.5D algorithm is to use SUMMA to compute \mathbf{K} while keeping \mathbf{V} in a 1D distribution. The subsequent distributed SpMM, $\mathbf{E}^T = \mathbf{V}\mathbf{K}$, is performed with \mathbf{V} 1D-distributed and \mathbf{K} 2D-distributed, enabling communication costs that decrease with process count and avoiding redistribution of \mathbf{K} to 1D. The SpMM is further structured to produce \mathbf{E}^T in a 1D columnwise distribution, allowing cluster updates to proceed without additional communication. This approach achieves scalable GEMM and SpMM while avoiding the limitations of the other three algorithms.

For simplicity, the 1.5D approach is described using $P = 4$, although Algorithm 2 describes the general case. To avoid communicating large partitions of \mathbf{K} during SpMM, we use a B-stationary communication schedule that exchanges only entries of \mathbf{V} and partial sums of the output \mathbf{E}^T . The distributed SpMM begins with an Allgather that replicates subsets of the 1D-partitioned \mathbf{V} along each process row. In particular, the first two partitions of \mathbf{V} are replicated along the first process row, and the last two along the second process row: $[\mathbf{V}_1 \ \mathbf{V}_2]$ are replicated on processes $P(1,1)$ and $P(1,2)$, while $[\mathbf{V}_3 \ \mathbf{V}_4]$ are replicated on processes $P(2,1)$ and $P(2,2)$. Each process then performs a local SpMM:

$$\mathbf{E}_{ij}^T = \mathbf{K}_{ij} [\mathbf{V}_{2(i-1)} \ \mathbf{V}_{2(i-1)+1}] \quad (20)$$

Algorithm 2 1.5D Kernel K-means Algorithm

Require: $\mathbf{P} \in \mathbb{R}^{n \times d}$ stores points; number of clusters k .

- 1: Use SUMMA to compute \mathbf{K} , partitioned across 2D grid
- 2: **while** not converged **do** ▷ Clustering loop
- 3: Allgather $[\mathbf{V}_{(i-1)\sqrt{P}} \cdots \mathbf{V}_{(i-1)\sqrt{P}+\sqrt{P}}]$ on $P(i, :)$, store in \mathbf{V}_{gather}
- 4: Compute $\mathbf{E}_{ij}^T = \mathbf{V}_{gather} \mathbf{K}_{ij}$
- 5: Partition \mathbf{E}_{ij}^T into $[\mathbf{E}_{ij}^T(1) \cdots \mathbf{E}_{ij}^T(\sqrt{P})]$
- 6: Reduce-scatter along process columns. $\mathbf{E}_{ij}^T(l)$ is reduced on process with global rank $p = l + (j-1)\sqrt{P}$
- 7: \mathbf{E}^T is now partitioned along columns in a 1D fashion.
- 8: $\mathbf{z}_p = \text{mask}(\mathbf{E}_p^T)$
- 9: $\mathbf{c}_p = \mathbf{V}_p \mathbf{z}_p$
- 10: Allreduce \mathbf{c}_p to compute \mathbf{c} on each process
- 11: $\mathbf{D}_p^T = \text{distances}(\mathbf{E}_p^T, \mathbf{c})$
- 12: $cl_p = \text{argmin}(\mathbf{D}_p^T)$ ▷ Updated cluster assignment
- 13: Update \mathbf{V}_p using cl_p
- 14: **end while**

Each partition of \mathbf{V} contains $\frac{n}{P}$ nonzeros, and the same number of \mathbf{V} partitions are replicated on each process. Thus, the local SpMMs are load balanced, as in the 1D algorithm.

Each \mathbf{E}_{ij}^T has k rows and $\frac{n}{\sqrt{P}}$ columns and represents a partial sum of the output. To form the complete \mathbf{E}^T , partial sums \mathbf{E}_{ij}^T associated with the same process column are accumulated. For example, when $P = 4$, $\mathbf{E}_{11}^T + \mathbf{E}_{21}^T$ produces the first $\frac{n}{2}$ columns of \mathbf{E}^T , while $\mathbf{E}_{12}^T + \mathbf{E}_{22}^T$ produces the remaining $\frac{n}{2}$ columns. The resulting sums can then be partitioned along either rows or columns, which determines the final layout of \mathbf{E}^T . Prior 1.5D SpMM approaches [47] partition along rows via a Reduce-Scatter across process columns, resulting in the following partitioning of each \mathbf{E}_{ij}^T :

$$\mathbf{E}_{ij}^T = \begin{bmatrix} \mathbf{E}_{ij}^T(1) \\ \mathbf{E}_{ij}^T(2) \end{bmatrix} \quad (21)$$

Partitions labeled (1) are reduced on the first process in column j , and those labeled (2) on the second. This results in a 2D partitioning of \mathbf{E}^T , which is undesirable because it increases communication for both cluster updates and \mathbf{V} updates. A 1D columnwise partitioning of \mathbf{E}^T would be preferable, as it would avoid this additional communication. To achieve this, we modify the Reduce-Scatter operation so that each \mathbf{E}_{ij}^T is partitioned along columns instead of rows:

$$\mathbf{E}_{ij}^T = [\mathbf{E}_{ij}^T(1) \quad \mathbf{E}_{ij}^T(2)] \quad (22)$$

Each $\mathbf{E}_{ij}^T(l)$ has k rows and $\frac{n}{P}$ columns. A Reduce-Scatter is performed so that process p reduces $\mathbf{E}_{ij}^T(l)$ with $2(j-1) + l = p$. \mathbf{E}^T is fully computed and naturally 1D columnwise partitioned, which enables cluster updates and \mathbf{V} updates to be computed without any communication, as in the 1D algorithm. Our 1.5D SpMM algorithm is shown in Figure 1.

In the 1.5D algorithm, computing \mathbf{K} has the same communication cost as H-1D, as it relies on SUMMA. The only

difference in communication arises when computing \mathbf{E}^T . The Allgather that replicates \mathbf{V} partitions along process rows, sending $\mathcal{O}(\sqrt{P})$ messages and $\mathcal{O}(n/P)$ words to \sqrt{P} processes, resulting in a total communication cost of:

$$T_{allgather} = \alpha \mathcal{O}(\sqrt{P}) + \beta \mathcal{O}\left(\frac{n}{\sqrt{P}}\right) \quad (23)$$

The bandwidth scales with the number of processes, unlike T_{E1D} from the 1D algorithm, while the latency term is reduced by a factor of \sqrt{P} . The Reduce-Scatter used to finalize \mathbf{E}^T sends $\mathcal{O}(\log(\sqrt{P}))$ messages and communicates $\mathcal{O}(\frac{nk}{P})$ words to \sqrt{P} processes, yielding a total communication cost of:

$$T_{reduce-scatter} = \alpha \mathcal{O}(\log(\sqrt{P})) + \beta \mathcal{O}\left(\frac{nk}{\sqrt{P}}\right) \quad (24)$$

The overall communication cost of computing \mathbf{E}^T is:

$$T_{E1.5D} = \alpha \mathcal{O}(\sqrt{P}) + \beta \mathcal{O}\left(\frac{n(k+1)}{\sqrt{P}}\right) \quad (25)$$

Our 1.5D algorithm reduces communication compared to 2D because cluster assignment does not require any communication. The $\frac{nk}{\sqrt{P}}$ bandwidth term is larger than that in 1D approaches for small P , but it scales with the process count: for large P , it is less than the $\mathcal{O}(n)$ bandwidth term for 1D. The 1.5D approach uses only $\mathcal{O}(\sqrt{P})$ messages, resulting in lower latency costs than 1D, and leverages SUMMA for \mathbf{K} without the redistribution overhead of H-1D, further reducing communication.

The 1.5D algorithm is the most effective and demonstrates the importance of considering the composability of primitives in distributed algorithms that use linear algebra. By design, keeping \mathbf{V} 1D-partitioned enables SpMM to efficiently consume the 2D-partitioned \mathbf{K} produced by SUMMA without redistribution. Moreover, keeping \mathbf{E}^T 1D-partitioned eliminates communication for cluster updates, highlighting the importance of application-aware partitioning schemes.

Our approach, developed for distributed Kernel K-means, could also benefit other clustering algorithms that use GEMM and SpMM, such as K-means [36] and spectral clustering [21]. In distributed memory, these algorithms could similarly benefit from communication schedules that emphasize composability.

V. IMPLEMENTATION

This section describes the implementation of each proposed algorithm. The algorithms are written in C++, using dense matrices in row-major order and local \mathbf{V} partitions in compressed sparse column (CSC) format. Storing dense matrices in row major order is known to improve the performance of cuSPARSE's SpMM routine [45]. The implementation is GPU-capable, storing matrices and vectors on the device and using single-precision floating-point numbers and 32-bit integers for indices. The local sparse operations use cuSPARSE [45], distributed GEMM uses SLATE [49], and communication is managed through GPU-aware MPI.

For the 1D and 1.5D algorithms, communication of \mathbf{V} partitions involves only their local row indices. This is sufficient to initialize the partition of \mathbf{V} on the receiving process, since the column pointers array contains the columns in the partition, and values can be computed using the global cluster sizes from an Allreduce at the end of each iteration.

In the 2D algorithm, the column pointers array is sent along with the local row indices, but values are initialized locally as in the 1D and 1.5D algorithms. \mathbf{V} is initialized by assigning points to clusters in a round-robin fashion. Other initialization strategies that provide stronger convergence guarantees, such as K-Means++ [2], are left for future work.

A. 1D and Hybrid 1D Implementation

The 1D implementation begins by initializing \mathbf{P}^T and \mathbf{P} on the devices. In the pure 1D implementation, both matrices are partitioned in a 1D fashion, while in the H-1D implementation, both matrices are partitioned in a 2D fashion. From here, SLATE’s GEMM routine computes \mathbf{K} . For the H-1D algorithm only, \mathbf{K} is then redistributed to a 1D columnwise partitioning using Alltoallv. From this point, the 1D and H-1D implementations proceed identically. \mathbf{V} is initialized using the round-robin assignment strategy.

\mathbf{V} is communicated via MPI_Allgather, followed by partitioned \mathbf{E}^T computation using cuSPARSE SpMM. The masking for \mathbf{z} uses a custom kernel, and \mathbf{c} is computed with cuSPARSE SpMV and finalized with a global MPI_Allreduce. A hand-written kernel sums $\tilde{\mathbf{C}}$ using \mathbf{c} and \mathbf{E}^T to form \mathbf{D}^T . Then, local argmin values are computed with another kernel, producing partitions of the cluster assignments array, which also serve as row indices for local \mathbf{V} partitions. Finally, a global Allreduce computes cluster sizes for updating \mathbf{V} .

B. 2D Implementation

In our 2D approach, \mathbf{P}^T , \mathbf{P} , and \mathbf{K} are partitioned in a 2D fashion. Our B-stationary 2D SpMM uses MPI_Allgather to replicate \mathbf{V} tiles along each process row, avoiding MPI_Bcast, and uses MPI_Reduce to sum partial results of \mathbf{E}^T block rows. This single Allgather approach is preferred over the typical \sqrt{P} Broadcast method, as prior work [47], [53] has shown that Broadcast can result in low arithmetic intensity in local SpMM computation and thus poorer performance. In addition, because partitions of \mathbf{V} have varying numbers of nonzeros in the 2D algorithm, the Broadcast approach would cause load imbalance in local SpMM. Using Allgather avoids this imbalance by ensuring each process column contains exactly $\frac{n}{\sqrt{P}}$ nonzeros. This choice does not significantly affect the communication cost analysis in Section IV, except that the $\log(\sqrt{P})$ factors in the latency and bandwidth terms are removed. \mathbf{c} is computed as in the 1D routine, except the Allreduce is performed along process rows. Cluster updates use a local argmin on each \mathbf{D}^T partition, followed by an MPI_Allreduce with MPI_MINLOC along process columns. Finally, an MPI_Allreduce along process rows computes cluster sizes, enabling a local kernel to update \mathbf{V} .

C. 1.5D Implementation

In the 1.5D implementation, \mathbf{P}^T and \mathbf{P} are 2D-partitioned, and \mathbf{K} is computed using SLATE GEMM. \mathbf{V} is partitioned as in the 1D algorithms. The main difference from the 1D algorithms occurs in computing \mathbf{E}^T . For each process row i , partitions of \mathbf{V} :

$$[\mathbf{V}_{(i-1)\sqrt{P}} \quad \mathbf{V}_{(i-1)\sqrt{P}+1} \quad \cdots \quad \mathbf{V}_{(i-1)\sqrt{P}+\sqrt{P}}],$$

are gathered on the i -th diagonal process $P(i, i)$ using MPI_Gather. Each diagonal process $P(i, i)$ then uses MPI_Bcast to replicate the partitions of \mathbf{V} along row i , and cuSPARSE SpMM computes $\mathbf{E}_{ij}^T(l)$. This is equivalent to performing an MPI_Allgather to replicate the partitions of \mathbf{V} shown above along process row i , and has an equivalent communication cost while simplifying the implementation. Finally, MPI_Reduce_scatter_block is used to sum the partial results of the 1D-partitioned \mathbf{E}^T along each process column while splitting the summed outputs along columns, yielding \mathbf{E}^T partitioned in a 1D fashion. Because each local SpMM produces output in row-major order, each $\mathbf{E}_{ij}^T(l)$ must be converted to column-major order before the MPI_Reduce_scatter_block to ensure that the portion sent to each process is stored contiguously. The additional time required for this conversion was negligible. Processes in the 2D grid are arranged in column-major order. This ensures that MPI_Reduce_scatter_block along process columns naturally places the fully reduced partitions of \mathbf{E}^T on contiguous processes, which is necessary for the 1D partitioning of \mathbf{E}^T . Computing \mathbf{D}^T and updating cluster assignments are performed in the same way as in the 1D implementation.

VI. RESULTS

In this section, we evaluate the performance of the four Kernel K-means algorithms presented in this paper on three real-world libSVM datasets [16], and compare the 1.5D algorithm’s runtime to a single-GPU implementation of Kernel K-means that uses a sliding-window approach to handle data that does not fit on a single GPU. To the best of our knowledge, there are no open-source GPU-accelerated distributed Kernel K-means implementations, so a direct comparison to a distributed baseline is not possible. However, since the 1D algorithm has a communication pattern similar to the non-linear algebraic formulations of distributed Kernel K-means described in the literature [23], [55], we use our 1D algorithm as a baseline.

A. Experiment Information

The experiments were run on the GPU partition of NERSC’s Perlmutter supercomputer [42]. Each node has four 80 GB NVIDIA A100 GPUs¹ connected by NVLink 3.0, with nodes interconnected via a dragonfly network and four Cassini-11 NICs. Codes were compiled with NVCC 12.9.41 (-O3, C++17) and used cuSPARSE 12.9, Cray MPICH 8.1.30, and SLATE 2025.05.28. The datasets from libSVM used to

¹On Perlmutter, there are also nodes with 40 GB A100s, but we run on the 80 GB nodes.

Dataset	n	d	Domain
KDD-sampled	8,407,752	10,000	Education
HIGGS	11,000,000	28	Physics
MNIST8m	8,100,000	784	Vision

TABLE II: The libSVM [16] datasets used for evaluation.

benchmark our algorithms are shown in Table II. The datasets were selected from different scientific domains and have varying numbers of features. For KDD, 10,000 features were randomly sampled to keep \mathbf{P} at a manageable size. This is a common practice for high-dimensional data, where dimensionality reduction or feature sampling is typically applied [13]. Unless otherwise noted, the experiments run the clustering loop for 100 iterations to ensure that runtime differences arise from performance, not convergence rate. Benchmarks use $k \in \{16, 32, 64\}$ and the polynomial kernel with $\gamma = 1, c = 1, d = 2$. Here, G denotes the total number of GPUs.

B. Weak Scaling

In this section, we evaluate the weak scaling of our Kernel K-means algorithms. For each GPU count G , $n = \sqrt{G} \times 96,000$ points are sampled so that \mathbf{K} fits within the aggregate GPU memory and the per-GPU workload for computing \mathbf{K} and \mathbf{E}^T remains constant as n increases. Figure 2 illustrates weak scaling for the four algorithms across datasets and k values. The 1.5D algorithm consistently achieves the best weak scaling, while the 2D algorithm outperforms both the 1D and H-1D approaches. For all k values and datasets, our 1.5D algorithm achieves a geometric mean weak scaling efficiency of 86.87% at 64 GPUs and 79.72% at 256 GPUs. Our 1.5D algorithm can cluster datasets with more than 1.5 million points, representing an increase of approximately $19\times$ over the largest dataset size clustered using exact Kernel K-means in prior work [10]. The H-1D algorithm cannot run on more than 16 GPUs because of the extra memory needed to redistribute \mathbf{K} from 2D to 1D. For KDD, the 1D algorithm fails on more than 4 GPUs, highlighting its memory inefficiency: the 1D GEMM routine requires replicating \mathbf{P} on each GPU. For KDD, with $d = 10,000$, the global \mathbf{P} matrix reaches several GB as G increases, making it impossible to store both a local \mathbf{K} partition and the replicated \mathbf{P} on a single GPU. In contrast, the 1.5D and 2D algorithms handle all problem sizes without memory issues, as they avoid both \mathbf{K} redistribution and \mathbf{P} replication.

Figure 3 illustrates the runtime breakdown of all four algorithms across GPU counts for MNIST8m and HIGGS with $k = 64$, highlighting their performance differences. The 1D algorithm scales poorly at higher GPU counts because the \mathbf{K} computation time increases with G , as discussed in Section IV. In contrast, our 1.5D algorithm benefits from SUMMA, enabling better scalability, especially for MNIST8m, where the larger d amplifies this advantage. The communication cost of computing \mathbf{E}^T in the 1.5D algorithm is comparable to that of the 1D algorithm for large G , despite the additional Reduce-Scatter. This aligns with our analysis in Section IV, which predicted that the extra cost would be negligible, espe-

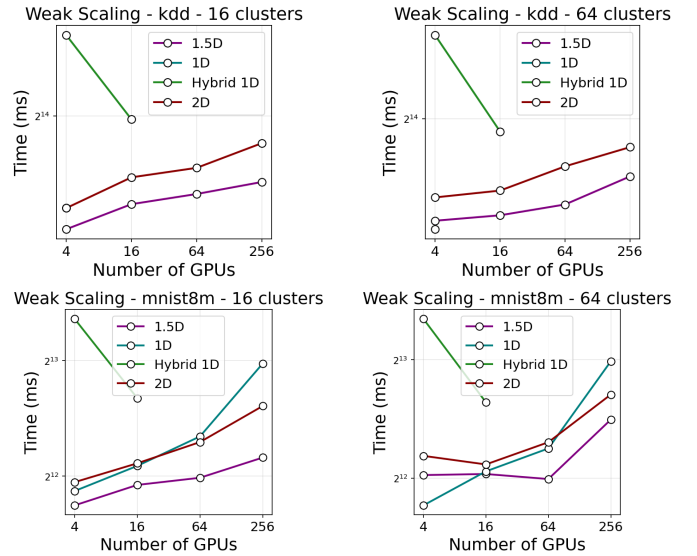


Fig. 2: The weak scaling evaluation on three datasets and $k \in \{16, 64\}$. Results for the HIGGS dataset and $k = 32$ are omitted for clarity.

cially at scale. The H-1D algorithm is the slowest due to the high cost of computing \mathbf{K} , which is dominated by redistribution overhead. Consequently, computing \mathbf{K} is more expensive than in the 1D case, even though a more communication-efficient algorithm is used.

The 2D algorithm scales similarly to the 1.5D algorithm but generally performs worse due to additional communication during cluster updates, particularly the MPI_Allreduce with MPI_MINLOC. For large G , this overhead becomes significant. Our analysis showed that MPI_Allreduce communicates $\mathcal{O}(\log(\sqrt{P})n)$ words. The MPI_MINLOC operator also doubles the buffer size to store an additional integer, further increasing the communication volume.

C. Strong Scaling

This section evaluates the strong scaling of our Kernel K-means algorithms using $n = 192,000$ sampled points per dataset, resulting in a kernel matrix near the single-node memory limit. Figure 4 illustrates strong scaling results for the four algorithms across datasets and k values. The 1.5D algorithm consistently scales best. Its advantage is most evident at lower k values (e.g., $k = 16$), where it scales efficiently to 256 GPUs on KDD. For MNIST8m and larger k values, the 1.5D algorithm generally does not scale beyond 64 GPUs but still outperforms others in scalability and runtime. Overall, across all values of k and all datasets, the 1.5D algorithm achieves a geometric mean strong scaling speedup of $4.65\times$ at 64 GPUs and $4.16\times$ at 256 GPUs. Both the 2D and H-1D algorithms consistently scale better than the 1D algorithm.

Figure 5 illustrates the runtime breakdown for strong scaling on MNIST8m and KDD with $k = 64$. Other datasets and other values of k show similar patterns to the ones in these plots. As in the weak scaling case, the 1D algorithm is limited by

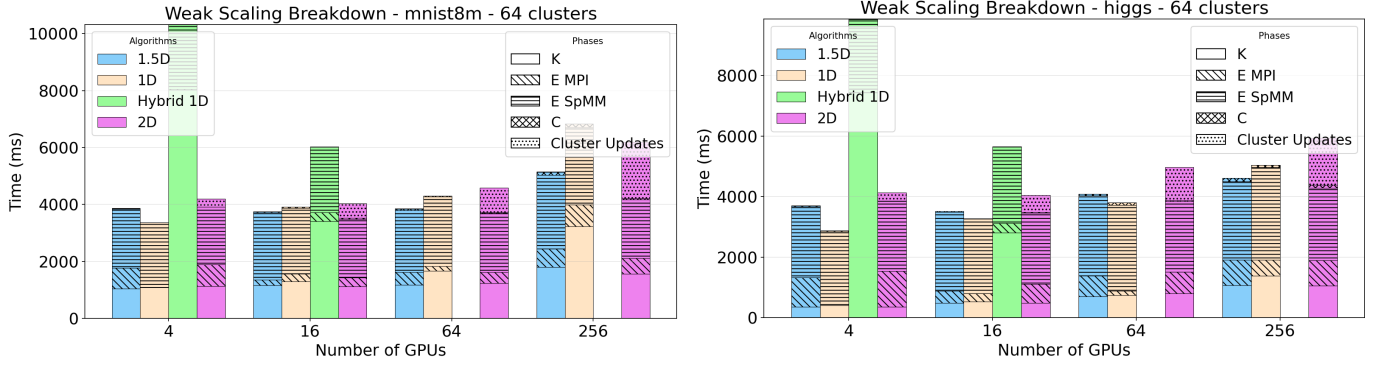


Fig. 3: Weak scaling runtime breakdown for MNIST8m and HIGGS for $k = 64$.

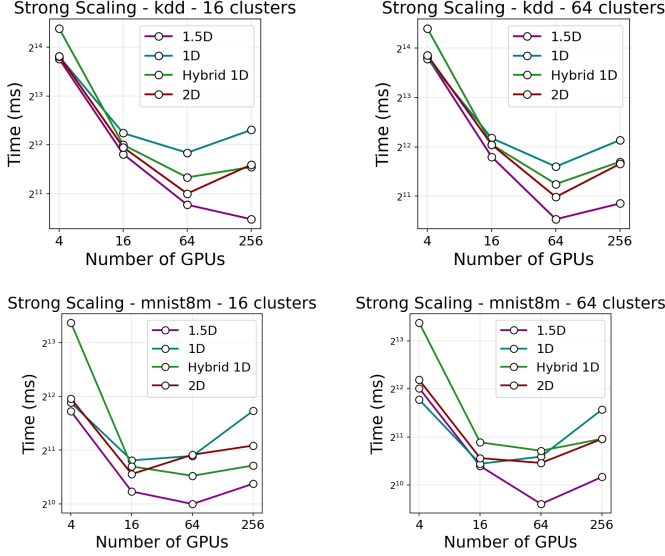


Fig. 4: Strong scaling evaluated on three datasets with $k \in \{16, 64\}$. Results for the HIGGS dataset and $k = 32$ are omitted for clarity.

poor \mathbf{K} scalability, while the 1.5D algorithm can avoid that bottleneck. The additional communication required to compute \mathbf{E}^T in 1.5D is minimal and scales well, eventually matching the 1D algorithm’s communication cost. In strong scaling, the $\mathcal{O}(\frac{n^2}{\sqrt{P}})$ redistribution cost in H-1D scales with GPU count, but the $\mathcal{O}(P)$ message count still creates a latency bottleneck, keeping it behind 1.5D. The 2D algorithm faces the same issue as in weak scaling: the MPI_Allreduce for argmin does not scale with GPU count and eventually becomes a bottleneck.

D. Comparison to Single GPU

Finally, we compare our 1.5D algorithm to a single-GPU sliding-window baseline for cases where \mathbf{K} exceeds GPU memory, using $n = 192,000$. The 1.5D approach is emphasized because it consistently outperforms other strategies.

Our sliding window algorithm is based on the approach to Kernel K-means described in [58], which stores \mathbf{K} on disk and loads it in blocks. In contrast, we recompute blocks of \mathbf{K} on the fly, trading increased computation for reduced disk I/O and host-device data movement, both of which are

more costly than computation. The sliding window algorithm computes a $b \times n$ block row of \mathbf{K} and updates b rows of \mathbf{D} at each step. Once $\lceil n/b \rceil$ steps are completed, new cluster assignments (cl_p and \mathbf{V}) are computed before the next Kernel K-means iteration. We set $b = 8192$, which yielded the best empirical performance. Computing \mathbf{K}_p is the main bottleneck in the sliding window algorithm because of the cost of GEMM and potential nonlinear operations, such as exponentiation for the polynomial kernel. Therefore, we store \mathbf{V} as a dense matrix in our experiments, since the time required to recompute \mathbf{K}_p dominates the computation of \mathbf{E}^T . Popcorn [10] was considered as a single-GPU baseline, but it failed to cluster datasets larger than approximately 80,000 points, making it unsuitable for our large-scale focus. To the best of our knowledge, no open-source distributed Kernel K-means implementations with GPU support exist. Therefore, the sliding window algorithm serves as the most relevant baseline outside our 1D algorithm.

Figure 6 illustrates the speedup of the 1.5D algorithm over the sliding window baseline across datasets and k values. At 256 GPUs, 1.5D is more than $10\times$ faster in all cases, with the largest gain on KDD at $k = 16$, achieving a $2749.8\times$ speedup and reducing runtime from over an hour to under 2 seconds. The speedups are especially pronounced for datasets with large d , where recomputing \mathbf{K}_p in the sliding window approach is more costly. Being able to cluster large datasets with Kernel K-means in seconds instead of hours significantly expands the scope of potential applications of Kernel K-means, since large-scale clustering in reasonable time is now possible.

VII. CONCLUSION AND FUTURE WORK

This work presents GPU-accelerated distributed-memory algorithms for Kernel K-means based on sparse and dense linear algebra primitives. By tailoring communication strategies and partitioning schemes to the structure of this clustering algorithm, particularly the sparsity of \mathbf{V} and the interaction between GEMM and SpMM, we make large-scale exact Kernel K-means clustering practical. In contrast to previous single-GPU exact approaches, our method enables clustering of datasets that are one to two orders of magnitude larger within reasonable time.

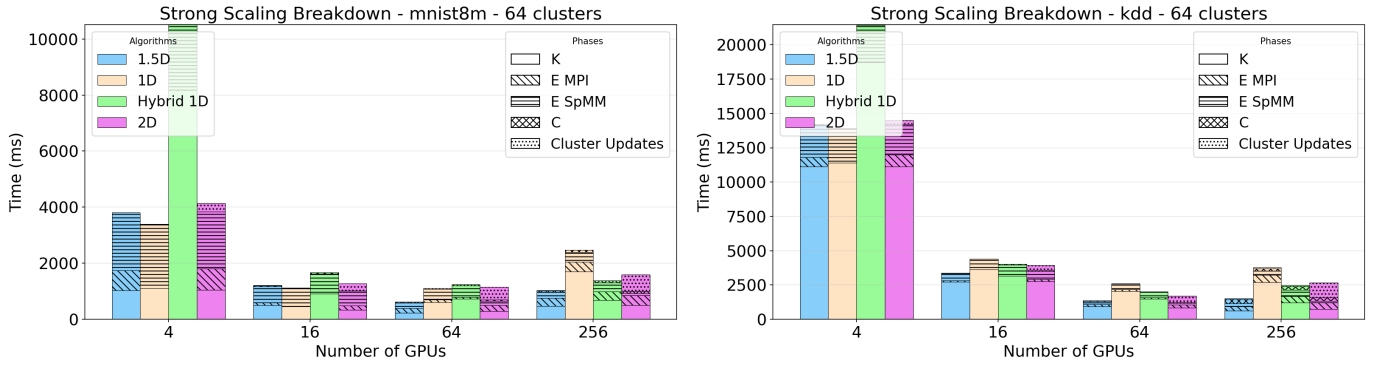


Fig. 5: The strong scaling runtime breakdown on MNIST8m and KDD for $k = 64$.

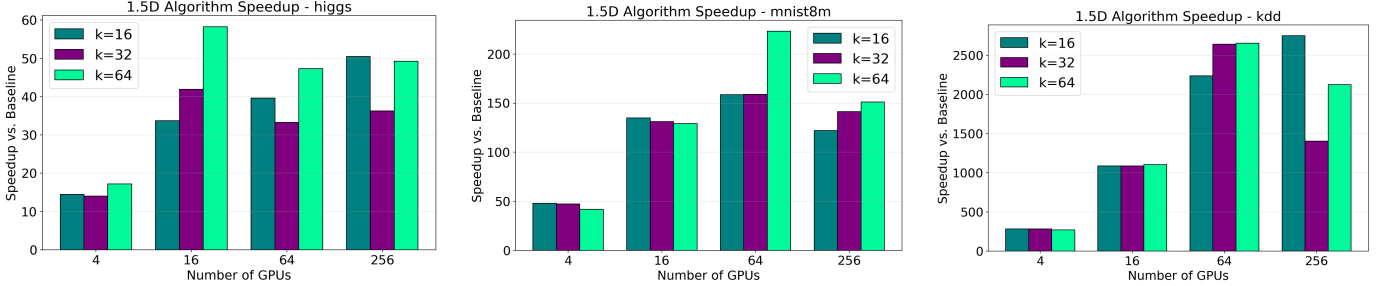


Fig. 6: The evaluation of strong scaling speedup over sliding window approach.

In this work, we study four distributed formulations: a 1D algorithm, a hybrid 1D algorithm, a 2D algorithm, and a 1.5D algorithm. The proposed 1.5D algorithm minimizes communication volume and achieves better scalability, with communication decreasing as the number of processes increases, unlike existing approaches where communication remains constant or increases more rapidly. Its effectiveness comes from the application-aware composability of linear algebra primitives: 1D-partitioned V enables SpMM to consume the 2D-partitioned K from SUMMA without redistribution, and 1D-partitioned E^T eliminates communication for cluster updates. As a result, our 1.5D algorithm demonstrates strong and weak scaling up to 256 GPUs on large real-world datasets, achieving substantial speedups over a single-GPU sliding-window baseline. Most importantly, it enables exact Kernel K-means clustering of datasets with more than 1.5 million points, representing a dramatic increase in tractable problem size compared to prior work.

For future work, we plan to develop additional clustering algorithms, such as spectral clustering and mean-shift, using sparse linear algebra, and implement them on distributed memory systems. In this scenario, we plan to explore mixed-precision techniques for Kernel K-means, standard K-means, and related clustering algorithms to ensure portability to newer architectures favoring lower precision. More broadly, we aim to develop methodologies that leverage domain-specific sparsity structures and sequences of linear algebra primitives to select communication-efficient data distribution and implementation strategies. One possible direction is an autotuning framework that analyzes end-to-end sequences of sparse and dense linear algebra primitives to automatically select optimal

data distributions and communication schedules for the entire workflow.

ACKNOWLEDGMENT

This research used resources of the National Energy Research Scientific Computing Center, a DOE Office of Science User Facility supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-AC02-05CH11231, using NERSC award ASCR-ERCAP0030076. This material is based upon work supported by the U.S. Department of Energy, Office of Science, Office of Advanced Scientific Computing Research, Department of Energy Computational Science Graduate Fellowship under Award Number DE-SC0025528. The authors acknowledge financial support from *ICSC – Centro Nazionale di Ricerca in High-Performance Computing, Big Data and Quantum Computing*, funded by the European Union – NextGenerationEU. The authors wish to disclose that generative AI and editing assistants, such as Microsoft Copilot and InstaText, were used to assist with grammar checking and to improve the clarity of the writing.

REFERENCES

- [1] Ramesh C Agarwal, Susanne M Balle, Fred G Gustavson, Mahesh Joshi, and Prasad Palkar. A three-dimensional approach to parallel matrix multiplication. *IBM Journal of Research and Development*, 39(5):575–582, 1995.
- [2] David Arthur and Sergei Vassilvitskii. k-means++: the advantages of careful seeding. In *Proceedings of the Eighteenth Annual ACM-SIAM Symposium on Discrete Algorithms*, SODA '07, page 1027–1035, USA, 2007. Society for Industrial and Applied Mathematics.
- [3] Ariful Azad and Aydın Buluç. Lacc: A linear-algebraic algorithm for finding connected components in distributed memory. In *2019 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 2–12. IEEE, 2019.

- [4] Ariful Azad, Aydin Buluç, and John Gilbert. Parallel triangle counting and enumeration using matrix algebra. In *2015 IEEE International Parallel and Distributed Processing Symposium Workshop*, pages 804–811. IEEE, 2015.
- [5] Ariful Azad, Georgios A Pavlopoulos, Christos A Ouzounis, Nikos C Kyrpides, and Aydin Buluç. Hipmcl: a high-performance parallel implementation of the markov clustering algorithm for large-scale networks. *Nucleic acids research*, 46(6):e33–e33, 2018.
- [6] Ariful Azad, Oguz Selvitopi, Md Taufique Hussain, John R. Gilbert, and Aydin Buluc. Combinatorial blas 2.0: Scaling combinatorial algorithms on distributed-memory systems, 2021.
- [7] Maria Fiorina Balcan, Steven Ehrlich, and Yingyu Liang. Distributed k-means and k-median clustering on general topologies. NIPS’13, page 1995–2003, Red Hook, NY, USA, 2013. Curran Associates Inc.
- [8] Grey Ballard, James Demmel, Olga Holtz, and Oded Schwartz. Minimizing communication in numerical linear algebra. *SIAM Journal on Matrix Analysis and Applications*, 32(3):866–901, 2011.
- [9] Mohammed Baydoun, Hassan Ghaziri, and Mohammed Al-Husseini. Cpu and gpu parallelized kernel k-means. 74(8):3975–3998, August 2018.
- [10] Julian Bellavita, Thomas Pasquali, Laura Del Rio Martin, Flavio Vella, and Giulia Guidi. Popcorn: Accelerating kernel k-means on gpus through sparse linear algebra. In *Proceedings of the 30th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming*, pages 426–440, 2025.
- [11] Maciej Besta, Raghavendra Kanakagiri, Harun Mustafa, Mikhail Karasikov, Gunnar Rätsch, Torsten Hoefer, and Edgar Solomonik. Communication-efficient jaccard similarity for high-performance distributed genome comparisons. In *2020 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 1122–1132. IEEE, 2020.
- [12] Vivek Bharadwaj, Aydin Buluç, and James Demmel. Distributed-memory sparse kernels for machine learning, 2022.
- [13] Christos Boutsidis, Anastasios Zouzias, Michael W. Mahoney, and Petros Drineas. Randomized dimensionality reduction for k-means clustering, 2014.
- [14] Benjamin Brock, Aydin Buluç, and Katherine Yelick. Rdma-based algorithms for sparse matrix multiplication on gpus. In *Proceedings of the 38th ACM International Conference on Supercomputing*, pages 225–235, 2024.
- [15] Aydin Buluç and John R Gilbert. The combinatorial blas: design, implementation, and applications. *Int. J. High Perform. Comput. Appl.*, 25(4):496–509, November 2011.
- [16] Chih-Chung Chang and Chih-Jen Lin. Libsvm: a library for support vector machines. *ACM transactions on intelligent systems and technology (TIST)*, 2(3):1–27, 2011.
- [17] Yifan Chen, Ethan N. Epperly, Joel A. Tropp, and Robert J. Webber. Randomly pivoted cholesky: Practical approximation of a kernel matrix with few entry evaluations, 2024.
- [18] Radha Chitta, Rong Jin, Timothy C. Havens, and Anil K. Jain. Approximate kernel k-means: solution to large scale kernel clustering. In *Proceedings of the 17th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD ’11, page 895–903, New York, NY, USA, 2011. Association for Computing Machinery.
- [19] Dingsheng Deng. DbSCAN clustering algorithm based on density. In *2020 7th international forum on electrical engineering and automation (IFEEA)*, pages 949–953. IEEE, 2020.
- [20] Nameirakpam Dhanachandra, Khumanthem Manglem, and Yambem Jina Chanu. Image segmentation using k -means clustering algorithm and subtractive clustering algorithm. *Procedia Computer Science*, 54:764–771, 2015. Eleventh International Conference on Communication Networks, ICCN 2015, August 21–23, 2015, Bangalore, India Eleventh International Conference on Data Mining and Warehousing, ICDMW 2015, August 21–23, 2015, Bangalore, India Eleventh International Conference on Image and Signal Processing, ICISP 2015, August 21–23, 2015, Bangalore, India.
- [21] Inderjit S. Dhillon, Yuqiang Guan, and Brian Kulis. Kernel k-means: Spectral clustering and normalized cuts. In *Proceedings of the Tenth ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages 551–556, New York, NY, USA, 2004. Association for Computing Machinery.
- [22] Reza Farivar, Daniel Rebolledo, Ellick Chan, and Roy Campbell. A parallel implementation of k-means clustering on gpus. In *Proceedings of the 2008 International Conference on Parallel and Distributed Processing Techniques and Applications (PDPTA)*, pages 340–345, Las Vegas, NV, USA, 2008.
- [23] Marco Jacopo Ferrarotti, Sergio Decherchi, and Walter Rocchia. Distributed kernel k-means for large scale clustering. *CoRR*, abs/1710.03013, 2017.
- [24] Jiří Filipovič, Matúš Madzin, Jan Fousek, and Luděk Matyska. Optimizing cuda code by kernel fusion: application on blas. *The Journal of Supercomputing*, 71(10):3934–3957, 2015.
- [25] Nahid Gholizadeh, Hamid Saadatfar, and Nooshin Hanafi. K-dbscan: An improved dbSCAN algorithm for big data. *The Journal of supercomputing*, 77(6):6214–6235, 2021.
- [26] Mark Girolami. Mercer kernel-based clustering in feature space. *IEEE transactions on neural networks*, 13(3):780–784, 2002.
- [27] Mehmet Gönen and Adam A Margolin. Localized data fusion for kernel k-means clustering with application to cancer biology. *Advances in neural information processing systems*, 27, 2014.
- [28] Giulia Guidi, Gabriel Raulet, Daniel Rokhsar, Leonid Oliker, Katherine Yelick, and Aydin Buluc. Distributed-memory parallel contig generation for de novo long-read genome assembly. In *Proceedings of the 51st International Conference on Parallel Processing*, pages 1–11, 2022.
- [29] Giulia Guidi, Oguz Selvitopi, Marquita Ellis, Leonid Oliker, Katherine Yelick, and Aydin Buluç. Parallel string graph construction and transitive reduction for de novo genome assembly. In *2021 IEEE International Parallel and Distributed Processing Symposium (IPDPS)*, pages 517–526. IEEE, 2021.
- [30] Teny Handhayani and Lely Hiryanto. Intelligent kernel k-means for clustering gene expression. *Procedia Computer Science*, 59:171–177, 2015.
- [31] Roger W. Hockney. The communication challenge for mpp: Intel paragon and meiko cs-2. *Parallel Computing*, 20(3):389–398, 1994.
- [32] Yuxi Hong and Aydin Buluc. A sparsity-aware distributed-memory algorithm for sparse-sparse matrix multiplication, 2024.
- [33] Ahmed AM Jamel and Bahriye Akay. Human activity recognition based on parallel approximation kernel k-means algorithm. *Computer Systems Science & Engineering*, 35(6), 2020.
- [34] Jeremy Kepner and John Gilbert. *Graph algorithms in the language of linear algebra*. SIAM, 2011.
- [35] Raye Kimmerer, Timothy G Mattson, Scott McMillan, Benjamin Brock, Erik Welch, Michel Pelletier, and José E Moreira. The graphblas 3.0 project. In *2024 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*, pages 478–481. IEEE, 2024.
- [36] Martin Kruliš and Miroslav Kratochvíl. Detailed analysis and optimization of cuda k-means algorithm. In *Proceedings of the 49th International Conference on Parallel Processing*, ICPP ’20, New York, NY, USA, 2020. Association for Computing Machinery.
- [37] R.J. Kuo, L.M. Ho, and C.M. Hu. Integration of self-organizing feature map and k-means algorithm for market segmentation. *Computers & Operations Research*, 29(11):1475–1493, 2002.
- [38] Grzegorz Kwasniewski, Marko Kabić, Maciej Besta, Joost VandeVondele, Raffaele Solcà, and Torsten Hoefer. Red-blue pebbling revisited: near optimal parallel matrix-matrix multiplication, 2019.
- [39] S. Lloyd. Least squares quantization in pcm. *IEEE Transactions on Information Theory*, 28(2):129–137, 1982.
- [40] Yi Lu, Shiyong Lu, Farshad Fotouhi, Youping Deng, and Susan J Brown. Incremental genetic k-means algorithm and its application in gene expression data analysis. *BMC bioinformatics*, 5:1–10, 2004.
- [41] Alessandro Lulli, Thibault Debatty, Matteo Dell’Amico, Pietro Michiardi, and Laura Ricci. Scalable k-nn based text clustering. In *2015 IEEE International Conference on Big Data (Big Data)*, pages 958–963. IEEE, 2015.
- [42] National Energy Research Scientific Computing Center (NERSC). Perlmuter supercomputer. <https://www.nersc.gov/systems/perlmutter/>, 2021. Accessed: 2025-04-12.
- [43] Wei Niu, Jiexiong Guan, Yanzhi Wang, Gagan Agrawal, and Bin Ren. Dnnfusion: accelerating deep neural networks execution with advanced operator fusion. In *Proceedings of the 42nd ACM SIGPLAN International Conference on Programming Language Design and Implementation*, pages 883–898, 2021.
- [44] NVIDIA. cublas 12.2. <https://docs.nvidia.com/cuda/archive/12.2.2/cublas/>, 2023.
- [45] NVIDIA. cusparse 12.2. <https://docs.nvidia.com/cuda/archive/12.2.2/cusparse/>, 2023.

- [46] Elizaveta Rebrova, Gustavo Chavez, Yang Liu, Pieter Ghysels, and Xiaoye Sherry Li. A study of clustering techniques and hierarchical matrix formats for kernel ridge regression, 2018.
- [47] Oguz Selvitopi, Benjamin Brock, Israt Nisa, Alok Tripathy, Katherine Yelick, and Aydın Buluç. Distributed-memory parallel algorithms for sparse times tall-skinny-dense matrix multiplication. In *Proceedings of the 35th ACM International Conference on Supercomputing, ICS '21*, page 431–442, New York, NY, USA, 2021. Association for Computing Machinery.
- [48] Oguz Selvitopi, Saliya Ekanayake, Giulia Guidi, Muaaz G Awan, Georgios A Pavlopoulos, Ariful Azad, Nikos Kyrpides, Leonid Oliker, Katherine Yelick, and Aydın Buluç. Extreme-scale many-against-many protein similarity search. In *SC22: International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–12. IEEE, 2022.
- [49] SLATE. Slate users’ guide. <https://icl.utk.edu/files/publications/2020/icl-utk-1664-2020.pdf>, 2020.
- [50] Edgar Solomonik, Maciej Besta, Flavio Vella, and Torsten Hoefer. Scaling betweenness centrality using communication-efficient sparse matrix multiplication. In *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, pages 1–14, 2017.
- [51] Rajeev Thakur, Rolf Rabenseifner, and William Gropp. Optimization of collective communication operations in mpich. *The International Journal of High Performance Computing Applications*, 19(1):49–66, 2005.
- [52] Davide Tosi, Redon Kokaj, and Marco Rocchetti. 15 years of big data: a systematic literature review. *Journal of Big Data*, 11(1):73, 2024.
- [53] Alok Tripathy, Katherine Yelick, and Aydın Buluç. Reducing communication in graph neural network training, 2020.
- [54] Alok Tripathy, Katherine Yelick, and Aydın Buluç. Distributed matrix-based sampling for graph neural network training. *Proceedings of Machine Learning and Systems*, 6:253–265, 2024.
- [55] Nikolaos Tsapanos, Anastasios Tefas, Nikolaos Nikolaidis, and Ioannis Pitas. A distributed framework for trimmed kernel k-means clustering. *Pattern recognition*, 48(8):2685–2698, 2015.
- [56] Robert A. van de Geijn and Jerrell Watts. Summa: Scalable universal matrix multiplication algorithm. Technical report, USA, 1995.
- [57] Shusen Wang, Alex Gittens, and Michael W. Mahoney. Scalable kernel k-means clustering with nystrom approximation: Relative-error bounds, 2019.
- [58] Rong Zhang and Alexander I Rudnicky. A large scale clustering scheme for kernel k-means. In *2002 International conference on pattern recognition*, volume 4, pages 289–292. IEEE, 2002.
- [59] Yi Zhang, Jie Lu, Feng Liu, Qian Liu, Alan Porter, Hongshu Chen, and Guangquan Zhang. Does deep learning help topic extraction? a kernel k-means clustering method with word embedding. *Journal of Informetrics*, 12(4):1099–1117, 2018.
- [60] Xiaochen Zhou and Xudong Wang. Memory and communication efficient federated kernel k-means. *IEEE Transactions on Neural Networks and Learning Systems*, 35(5):7114–7125, 2022.