





MatRox: Modular approach for improving data locality in Hierarchical (Mat)rix App(Rox)imation

Bangtian Liu* University of Toronto, Canada bangtian@cs.toronto.edu

Kazem Cheshmi* University of Toronto, Canada kazem@cs.toronto.edu

Saeed Soori University of Toronto, Canada sasoori@cs.toronto.edu

Michelle Mills Strout University of Arizona, USA mstrout@cs.arizona.edu

Maryam Mehri Dehnavi University of Toronto, Canada mmehride@cs.toronto.edu

Abstract

Hierarchical matrix approximations have gained significant traction in the machine learning and scientific community as they exploit available low-rank structures in kernel methods to compress the kernel matrix. The resulting compressed matrix, *HMatrix*, is used to reduce the computational complexity of operations such as HMatrix-matrix multiplications with tuneable accuracy in an evaluation phase. Existing implementations of HMatrix evaluations do not preserve locality and often lead to unbalanced parallel execution with high synchronization. Also, current solutions require the compression phase to re-execute if the kernel method or the required accuracy change. MatRox is a framework that uses novel structure analysis strategies with code specialization and a storage format to improve locality and create load-balanced parallel tasks for HMatrix-matrix multiplications. Modularization of the matrix compression phase enables the reuse of computations when there are changes to the input accuracy and the kernel function. The MatRox-generated code for matrix-matrix multiplication is 2.98×, 1.60×, and 5.98× faster than library implementations available in GOFMM, SMASH, and STRUMPACK respectively. Additionally, the ability to reuse portions of the compression computation for changes to the accuracy leads to up to 2.64× improvement with MatRox over five changes to accuracy using GOFMM.

• Software and its engineering \rightarrow Source code generation; Software performance.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

PPoPP '20, February 22-26, 2020, San Diego, CA, USA © 2020 Association for Computing Machinery. ACM ISBN 978-1-4503-6818-6/20/02...\$15.00 https://doi.org/10.1145/3332466.3374548

1 Introduction

Kernel methods in machine learning such as kernelregression [30, 50], Gaussian kernel-based learning methods [16, 38], and support vector machines [18] used in applications such as feature extraction [52] and text recognition [9, 50], as well as scientific simulations, involve computations on a dense symmetric positive definite (SPD) matrix that is obtained by computing a kernel function \mathcal{K} on pairs of points from a set of points $\{x_1, \ldots, x_N\}$. The values of the $N \times N$ kernel matrix K are given by $K(i, j) = \mathcal{K}(x_i, x_i)$, with a typically large N. For example, in Gaussian ridge regression, the kernel $exp(-||x_i - x_i||_2^2/2h^2)$, with h as bandwidth, is applied to a machine learning dataset, i.e. points. The resulting kernel matrix is used in costly matrix-matrix multiplications, with complexity $O(N^3)$, in a direct solver to minimize a loss function.

The computational complexity of kernel matrix computations is reduced significantly, leading to orders of magnitude performance improvement [5], if instead of assembling Kand operating on it, it is compressed to \widetilde{K} using the kernel function, points, and an admissibility condition [6, 28]. An admissibility condition is the value of a distance measure between points above with which the kernel value for that point pair is approximated.

Many of kernel matrices are structured (or low-rank, or data-sparse). Hierarchical matrix computations, leverage the data-sparse structure introduced by the point set distribution and the admissibility condition during a compression phase to implicitly obtain \widetilde{K} . First, points are partitioned into a hierarchy and a cluster tree (CTree) is created from the partitioning of points. Then, based on the admissibility condition, i.e. a geometric distance between points, interactions between nodes are calculated and included in the CTree to form the HTree. This structure, embedded in the HTree and the CTree, is used in the compression phase to hierarchically approximate low-rank blocks of the kernel matrix. The low-rank approximated blocks, as well as the blocks that are not approximated, are referred to as subma*trices* and form the compressed matrix \widetilde{K} , also known as the *HMatrix*. The submatrices are then operated on instead of *K* in an evaluation phase.



^{*}Both authors contributed equally to this research.

Previous work attempts to optimize hierarchical matrix computations, specifically the evaluation phase, on parallel multicore architectures [20, 32, 62]. \mathcal{H}^2 structures are amongst the most commonly used hierarchical algorithms and GOFMM [62], STRUMPACK [20], and SMASH [8] are well-known libraries that implement \mathcal{H}^2 structures. However, these implementations do not preserve locality and often lead to a load-imbalanced execution with high synchronization overheads, which limits the performance and scalability of hierarchical matrix evaluations on parallel architectures.

The order and dependency of computations during evaluation is determined by the HTree. GOFMM [62] uses the HTree as the input for dynamic task scheduling, however, their scheduling trades locality for load balance. SMASH [8] traverses the CTree level-by-level, thus, synchronization overheads increase with the length of the critical path. Also schedulers that work with the CTree do not realize the additional dependencies introduced by the admissibility condition, which leads to additional synchronization costs. Implementations such as SMASH and STRUMPACK [20] do not optimize for load balance. Finally, some libraries are optimized for a specific \mathcal{H}^2 structure, for example STRUMPACK is specialized for Hierarchical Semi-separable (HSS) [11] structures, or only support low-dimensional points, e.g. SMASH.

In this paper, we present two structure analysis algorithms to generate specialized code that improves data locality and maintains a good load balance for HMatrix evaluations. We store the submatrices in the order they are visited during evaluation to improve locality. Our work focuses on HMatrix-matrix multiplications for the evaluation phase; we use the words HMatrix evaluation and HMatrix-matrix multiplication interchangeably throughout the paper. The proposed algorithms are implemented in a framework called MatRox, which uses structure information from the points, the kernel function, as well as the admissibility and accuracy requirement. The MatRox *inspector* compresses K, analyzes structure, and generates optimized code. The *executor* computes the HMatrix-matrix multiplication.

Additionally, MatRox enables partial reuse of computations when the kernel function and/or input accuracy are modified. In scientific and machine learning simulations, often the kernel matrix has to be re-compressed because either the overall accuracy of the HMatrix-matrix multiplication is not sufficient or has to be reduced for faster evaluation, or the kernel function changes. While available libraries have to rerun the costly compression, MatRox reuses parts of the previous inspection, e.g. compression information, and also reuses the previously generated evaluation code.

The main contributions of this work include:

 A novel structure analysis strategy, based on the modularization of compression, called interaction blocking.
 Current approaches parallelize HMatrix evaluation in two parts. One of these parts visits all pairs of nodes

- that interact with a parallel reduction. The proposed interaction blocking algorithm turns reduction parallelism into embarrassingly parallel computations by organizing interactions into blocks that share data which improves locality and removes synchronization.
- The structure coarsening algorithm that is also based on the modularization of compression. The second part that current approaches parallelize in HMatrix evaluation sequentially visits all levels in the hierarchy and in-parallel computes nodes in the same level (i.e. wavefront parallelism [55].) The structure coarsening algorithm organizes the wavefront parallelization into coarser groupings to better balance data locality and load. The number of points organized into each node and a user-requested accuracy are used to approximate load.
- The Compressed Data-Sparse (CDS) storage format that organizes data based on the loop visitation order in HMatrix evaluation to improve locality.
- Implementation of the proposed strategies in a framework called MatRox. The MatRox-generated code is on average 2.98×, 1.60×, and 5.98× faster than GOFMM, SMASH, and STRUMPACK respectively.
- An approach that enables the reuse of computations in compression for when the kernel function and/or accuracy change. The approach enables deriving a family of successively-approximated HMatrices in incremental fashion, rather than starting from scratch. MatRox with reuse is 2.21× faster than GOFMM over 5 changes to the input accuracy.

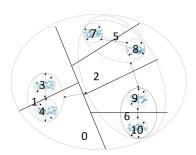
2 Approach Overview

In this section, we review the typical approach to hierarchical approximation and then provide an overview of the approach presented in this paper that is implemented in MatRox.

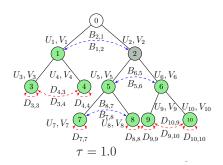
2.1 Background

Current library implementations of hierarchical matrix approximations typically have an interface as follows. The user provides to the library, the pointset shown in Figure 1a, an admissibility parameter τ , a kernel function, and a desired block approximation accuracy (bacc). The compression phase approximates the kernel matrix, implicitly created using the points and the kernel function. The admissibility parameter τ dictates how pairs of points are determined to be far or close, and the block approximation accuracy bacc indicates how closely submatrices need to be approximated. A representation of the compressed matrix is the input to the evaluation code that multiplies the HMatrix with another matrix or vector. Compression is typically expensive, however, the objective is to compress the kernel matrix once and reuse over many evaluations, e.g. multiple matrix-vector multiplications or a matrix-matrix multiplication.

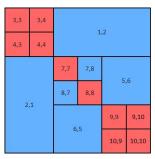




(a) Example pointset that has been hierarchically partitioned. For example partitions 3 and 4 are both included in partition 1.



(b) Cluster tree and Htree. Each node represents a partition. Edges between nodes indicate an interaction.



(c) Conceptual diagram of approximation matrix. Red submatrices are not approximated, blue ones are.

```
1 // Loops with reduction
2 #pragma omp parallel for
3 for (every near interaction (i,j)){
4 #pragma omp atomic
5 Y_i += D_{i,j} * W_j;
6 }
7
8 
9 // Loops with carried dependencies
10 for(l=depth; l>=1;l--)
11 #pragma omp parallel for
12 for(every node i \in level-1)
13 if(i is leaf node)
14 S_i = V_i^T * W_i
15 else
16 S_i += V_i^T * S_{lc}
17 S_i += V_i^T * S_{rc}
```

```
(d) Library code for evaluation for D and V matrices
```

```
1 // Blocked loop
2 #pragma omp parallel for schedule(dynamic)
3 for(every block b)
     for(every d \in blockset[b])
       i = d[0]; j = d[1];
       Y_i += D_{i,j} * W_i;
 8 // Loop over the CTree
9 for(every coarsen level cl)
10 #pragma omp parallel for schedule(static)
     for(every sub-tree st)
       for(every i \in coarsenset[c1][st])
         if (i is a leaf node)
           S_i = V_i^T * W_i
           S_i += V_i^T * S_{lc}
S_i += V_i^T * S_{rc}
16
17
```

(e) MatRox code for evaluation for D and V matrices.

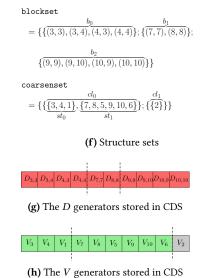


Figure 1. With HMatrix approximation, Y = KW is approximately computed with $Y \approx \widetilde{K}W$. Figures 1a and 1b demonstrate how \widetilde{K} is implicitly obtained and Figure 1c shows its conceptual diagram. Figure 1b illustrates the hierarchical organization of subpartitions of the pointset and the submatrices U_i , V_i , B_i and D_i (where i is a CTree node) that will be generated to approximate interactions between points in various partitions. \widetilde{K} is implicitly obtained from these submatrices. Figure 1d demonstrates a typical library approach and Figures 1e to 1h show MatRox's approach for computing $Y \approx \widetilde{K}W$, i.e. HMatrix evaluation. In a simplified form (not accounting for hierarchy, i.e. submatrices per node), we have $\widetilde{K} = UBV^T + D$ and thus the operations to perform $Y = \widetilde{K}W$ are, OP1: Y += D*W, OP2: $S = V^T*W$, OP3: L += B*S, OP4: Y += U*L; S is an intermediate matrix. The codes in Figures 1d and 1e perform these operations in hierarchical form; i is a node in the tree, i and i are the left and right children of node i respectively. Lines 3-6 in Figures 1d and 1e apply OP1 on the D submatrices; Lines 10-17 do a CTree bottom-up traversal for OP2 on V submatrices; OP3 and OP4 are the same loop type as OP1 and OP2 respectively and thus are not shown for space. Figure 1d is a typical library implementation of HMatrix evaluation using the CTree, Line 10 in the figure visits the levels of CTree for OP2. MatRox groups the i submatrices into blocks based on interactions and the i matrices into coarsen level sets (Figure 1f) and also stores the submatrices in the CDS format as shown in Figures 1g and 1h. The MatRox code in Figure 1e uses the sets and CDS to optimize HMatrix evaluation; i is the index for coarsen levels, and i is the index for sub-trees in the i-th coarsen level. The figure only shows blocking of near interactions in MatRox.

Compression. To approximate K, points are first clustered into hierarchically-organized sub-domains. Figure 1a shows a clustering that creates 10 sub-domains. Figure 1b shows a cluster tree for this clustering with each of the tree nodes representing a sub-domain. If two sub-domains i and j are

Far from each other, their interaction, (i,j), is approximated. Interactions between sub-domains Near to each other are not approximated. The admissibility condition $\tau dist(\alpha,\beta) > (diam(\alpha) + diam(\beta))$ in which $dist(\alpha,\beta)$ is the geometrical distance between the two sub-domains α and β , $diam(\alpha)$ is



the diameter of α , and τ is the input admissibility parameter, determines near-far interactions. The added dashed edges in blue and red to the cluster tree in Figure 1b represent these interactions and together form the HTree. The red edges are near interactions and blue edges show the far interactions.

Figure 1c shows an example approximated *K* matrix with colored sub-blocks. The blue-colored blocks are matrix blocks that are approximated during the compression phase¹. The degree to which each block is approximated is determined by the *submatrix-rank* (srank). The srank is adaptively tuned to meet the user-requested *block approximation accuracy* (*bacc*). The red blocks are not approximated.

Evaluation. The approximated matrix \widetilde{K} is never explicitly assembled, instead the HTree is used during evaluation to compute the desired HMatrix-matrix multiplication.

Figure 1d shows a typical library implementation of the evaluation phase. Existing library implementations of evaluation are classified into (i) Loops with reduction that operate on the near and far interacting nodes in the HTree; (ii) Loops with sparse dependencies between parents and children that do a bottom-up or top-down, level-by-level traversal of the cluster tree to generate other portions of the approximated matrix. Lines 3-6 in Figure 1d show the reduction loop computing near interactions by operating on the D submatrices; operation on the B submatrices is of the same loop type. Lines 10-17 in Figure 1d show the loop over the CTree that operates on the V matrices with bottom-up traversal; the U submatrices are operated on with the same loop type but using top-down traversal. Some library implementations perform the level-by-level traversal with a synchronization between levels and others place tasks into a dynamic task queue to enable run-time load balancing.

2.2 Approach implemented in MatRox

MatRox is composed of an inspector and an executor. The inspector is a modularized implementation of compression. It analyzes structure to generate optimized evaluation code and to store the submatrices associated with nodes in the cluster tree into an optimized data structure called Compressed Data-Sparse (CDS). Together the optimization of the code and reorganization of the data lead to faster evaluation compared to libraries. Additionally, when the inputs to the inspector do not change, the inspection can be conducted at compile-time.

The approach presented here improves data locality and reduces synchronization costs by blocking computations and associated data when the computations share data and are dependent on each other. For example, the computations involving the submatrices (9,9), (9,10), (10,9), and (10,10) are blocked together using *interaction blocking*. The interaction blocking algorithm analyzes the HTree to create a *blockset* shown in Figure 1f that creates an order for computation.

```
1 /// MatRox Inspector Code
2 #include <matrox.h>
      // Inputs declaration
     Points points("path/to/load/points");
     Float(64) tau = 0.65:
     Float(64) bacc = 1e-5;
     Ker kfunc = GAUSSIAN;
      // Outputs declaration
     HMatrix H("path/to/store/hmat.cds");
10
      Op HMatMul("path/to/store/matmul.h");
      inspector(points, tau, kfun, bacc, &H, &HMatMul);
12 /// MatRox Executor Code
13 #include "path/to/matmul.h"
     HMatrix H("path/to/load/hmat.cds");
     Dense W("path/to/load/matrix/W");
     Dense Y(Float(64), H.dim * W.cols);
      Y = matmul(H, W);
```

Figure 2. How a user provides parameters and calls the MatRox inspector for compression and executor for evaluation.

This enables the blocked loop in Figure 1e to be fully parallel, because the blocks are selected to eliminate reduction dependencies between block computations. MatRox uses the OpenMP dynamic scheduler to load-balance based on the blocksets. Interaction blocking also improves temporal locality with reusing W_j and Y_i in Figure 1e line 6. Spatial locality is also improved through consecutive storage of submatrices D_i in CDS.

The loop over the CTree is reorganized into coarsen levels and load-balanced sub-trees within those coarsen levels. The structure coarsening algorithm analyzes the cluster tree to create a coarsenset shown in Figure 1f that contains the coarsen levels and sub-trees. In Figure 1b, there are two coarsen levels (cl_0, cl_1) . The green nodes are in coarsen level 0 and node 2 is in coarsen level 1 by itself (node 0 is not involved in any computation). The coarsened loop in Figure 1e has a sequential loop over the coarsen levels and then a parallel loop over all of the sub-trees within that coarsenset. The sub-trees are load-balanced based on the srank. MatRox uses the static OpenMP scheduler for loops on coarsensets. V sub-matrices associated with all of the nodes in a coarsenset are organized together in the CDS format to improve spatial locality as shown in Figure 1g and 1h. Structure coarsening also improves temporal locality by keeping subtrees in one coarsenset, so that matrix S in lines 14-17 in Figure 1e is

An example of how MaRox is used is shown in Figure 2 in which the user provides the points, the kernel function, the admissibility condition, and block accuracy to the inspector. The output of the inspector is used by the executor to complete evaluation. The CDS stored submatrices, shown with H in Figure 2, as well as the generated HMatrix-matrix multiplication code are used by the executor.

In addition to the interaction blocking and structure coarsening optimization, MatRox also specializes the evaluation



 $^{^1\}mbox{We}$ use interpolative decomposition [42] for approximations.

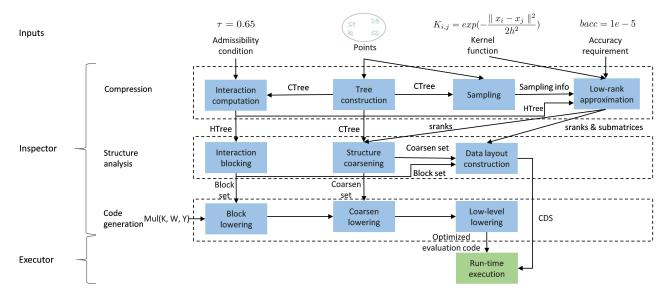


Figure 3. MatRox takes admissibility, points, kernel function, and accuracy as inputs and generates a storage format and an optimized code for HMatrix-matrix multiplication. It first compresses the matrix in the compression phase and then inspects the output of the compression phase in structure analysis. MatRox then uses the result of structure analysis, i.e. structure sets, to generate an optimized code and a storage format CDS. The MatRox executor runs the generated code with CDS.

code for a given matrix block. For example, since the parallelism in the HTree is less when we get closer the root, MatRox peels the last iteration of the nested computation to exploit block-level parallelism, e.g. with parallel BLAS. With all these changes, MatRox obtains 9.06× speedup compared to GEMM and 2.11× compared to GOFMM for covtype dataset on a Haswell processor.

3 Modular HMatrix Approximation

MatRox consists of an inspector that generates specialized code and an efficient storage of the compressed data to improve locality and load balance in HMatrix-matrix multiplications. Figure 3 shows the overview of MatRox. The inspector is separated into three phases of modular compression, structure analysis, and code generation. The user-provided inputs are separately passed to their respective modules in compression. Compression generates the submatrices, sranks, as well as the CTree and HTree to be used by different components of structure analysis. Information from structure analysis, i.e. the structure sets, is used along with an internal representation of the HMatrix-matrix multiplication for code lowering and specialization in the code generation stage. Finally the generated code and the submatrices stored in CDS are used by the executor for an efficient HMatrix-matrix multiplication. This section describes the MatRox internals.

3.1 Modularizing compression

MatRox provides a modularized design for the compression phase by defining four well-separated modules, i.e. interaction computation, tree construction, sampling, and low-rank approximation. Each module has well-defined inputs and creates and stores one or more of the outputs HTree, CTree, the *sranks*, and the submatrices, which we call the *structure information*.

By modularizing compression we divide it into smaller pieces, each of which will take only the required user-provided inputs and/or inputs from another piece. With the tree construction, interaction computation, low-rank approximation, and sampling modules, the structure information are separately stored and are passed to the structure analysis phase or to other parts in compression. The following discusses each module in the compression phase.

Tree construction and interaction computation. The points are inputs to the tree construction module and the output is the CTree. The CTree is constructed recursively using a partitioning algorithm with the tree root as the entire pointset. The partitioning terminates when the number of points in the leaf node is less than a pre-defined constant m, i.e. leaf size. MatRox uses two partitioning algorithms, kd-tree [47] and two-means [50], which are respectively used for low ($d \le 3$) and high (d > 3) dimensional data. The interaction computation module takes as input the CTree and the admissibility parameter to find near and far subdomains using the admissibility condition. It then adds the interactions to the CTree to create and store the HTree.

Low-rank approximation and sampling. The inputs to the low-rank approximation module are the HTree, kernel function, block-accuracy, and the sampling information and the outputs are the sranks and the submatrices. Interpolative Decomposition (ID) [42] is used to create the U, V, and



B submatrices with low-rank approximation and the full-rank D submatrices, i.e. the near blocks, are stored without approximation. Each low-rank block in that is compressed with a rank that is adaptively tuned to meet the input block-accuracy specified by the user. The rank with which a block is approximated with is stored in the sranks vector. ID can be expensive for larger block sizes [34], thus, sampling techniques are used to reduce the overhead of ID [39].

Sampling is a separate module in MatRox compression that takes only the points and the CTree as inputs and generates the sampling information for each sub-domain to be used by the low-rank approximation module. MatRox uses nearest-neighbour sampling [35] to reduce the overhead of low-rank approximation. The sampling module first takes the unclustered points to generate the k-nearest-neighbour list for each point [15]. k is the number of sampled points, sampling size [62], and is a predefined constant. Finding the exact k-nearest-neighbours of all points can be costly (points with high dimensions.) To reduce this overhead, we use a greedy search based on random projection trees that recursively partitions the points along a random direction [15]. The lists are then combined for each block using the clustering in CTree to form a nearest-neighbour list for the corresponding sub-domain/block. Finally, importance sampling [39] selects from the nearest-neighbour list of a block and generates the sampling information for that block.

3.2 Structure analysis

As shown in Figure 3, after finishing compression, all structure information is known, and MatRox analyzes this information using the interaction blocking and structure coarsening algorithms to create the coarsenset and blockset that are later used to generate specialize code for HMatrix-matrix multiplication. The submatrices and sranks from the modular compression phase are used with the sets to store the HMatrix in the CDS format. In this section, we describe this structure analysis.

Interaction blocking. As shown in Algorithm , the interaction blocking algorithm takes the HTree and an additional parameter called blocksize, as inputs and creates the blockset. We only show the interaction blocking algorithm for near interactions; far interactions follow a similar algorithm. The interaction blocking algorithm in lines 3-9, maps a near interaction between nodes i and j to the location of (i/blocksize, j/blocksize) in the blocks array. This mapping increases the possibility that interactions that involve the same node are in the same block which increases locality. However, as shown in line 5 of Figure 1d all interactions (i,j) will write to the location i of j, thus, these interactions have to be put in the same block of blockset to eliminate synchronization; this is implemented in Lines 10-16 of Algorithm .

Structure coarsening. The structure coarsening algorithm creates a coarsenset that optimizes the loops over the

```
Algorithm: Interaction blocking for near interac-
 tions
               :HTree, blocksize
   Input
   Output
               :blockset
1 blockDim = (HTree.numNodes - 1 + blocksize) / blocksize
2 blocks[blockDim,blockDim] = (0,0);
   /* Find blocks based on near interactions
3 for every node i \in HTree and i != root do
       iid = (i-1) / blocksize
4
       for node j \in HTree.near[i] do
5
           jid = (j-1) / blocksize
 7
           blocks(iid,jid).append(i,j)
       end
9 end
   /* Add blocks into blockset
10 for i=0; i<blockDim; i++ do
       for j=0; j < blockDim; j++ do
           if blocks(i,j).size() > 0 then
12
               blockset[i].append(blocks(i,j))
13
14
           end
       end
15
16 end
```

CTree by improving locality while maintaining load balance. The algorithm is an adaptation of the Load-Balanced level Coarsening (LBC) method [13] with the difference that here the algorithm is designed for binary trees and a different cost model based on sranks is used to balance load. As shown in Algorithm, structure coarsening takes the CTree, the sranks, number of sub-trees p, and a tuning parameter aqq as inputs and generates a coarsenset. In lines 2-7 the levels of the CTree are coarsened to build the coarsened levels. A level of a tree refers to nodes with the same height. *Tree[lb:ub]* shows a coarsen level that includes nodes with levels in the range of lb-ub. Algorithm builds all disjoint trees in a coarsen level, line 5, and stores them in coarsenset in post-order. For example in Figure 1b, the disjoint trees of HTree[0:1] are sub-trees with a root node in 1, 5, and 6. This ensures all nodes with dependency are assigned to the same thread to improve locality. The structure coarsening algorithm computes the cost of each sub-tree using sranks in lines 8-14. The subtree cost is related to the size of submatrices associated with the subtree nodes and is determined by sranks. The computed costs are used in lines 15-19 of Algorithm to merge the initial disjoint sub-trees with a first-fit bin-packing algorithm [14] and to create p new sub-trees that are load balanced. These sub-trees will execute in parallel.

Data layout construction. In the final phase of structure analysis, MatRox uses the structure sets to store the generated submatrices in a format, which we call the compressed data-sparse (CDS) format, that improves locality during the HMatrix-matrix multiplication. CDS follows the order of computations in the blocked and coarsened loops



Algorithm: Structure coarsening Input :CTree, p, agg, sranks Output :coarsenset 1 $l = \lceil CTree.height/agg \rceil$ 2 **for** i=0; i< l; i++ **do** $lb = i^*agg;$ ub = (i+1)*agg;cl = disjoint_subtrees(CTree[lb:ub]); coarsenset.append(cl); 6 7 end /* Cost estimation for each node 8 for $node x \in CTree do$ **if** $x \in CTree.leafnodes$ **then** CTree[x].cost = cost(sranks(x))10 11 else CTree[x].cost = cost(sranks(x),12 $\operatorname{srank}(\operatorname{lchild}(x)) + \operatorname{sranks}(\operatorname{rchild}(x)))$ end 13 14 end /* Merge subtrees in each coarsen level 15 **for** i=0; i< l; i++ **do** cl = coarsenset[i] 16 nPart = cl.size() > p ? p : cl.size()/217 coarsenSet[i] = bin_pack(cl, nPart); 18 19 end

which is obtained from the structure sets. More specifically, the U, V submatrices are stored in the order specified by the coarsenset and the B, D submatrices are stored by the order specified by the near and far blocksets. The size of each submatrix is known with sranks and is used as the offsets in CDS (see Figure 1f for an example.)

3.3 Code generation

Code generation in MatRox uses structure information to lower an internally generated abstract syntax tree (AST) to an optimized evaluation code. Figure 3 shows different components of code generation. MatRox lowers the AST in either the block or the coarsen lowering stages or both. The resulting lowered code from these stages iterates over the structure set. Figure 1e shows an example lowered code where the blocked loop iterates over the blockset and the coarsen loop goes over coarsenset. The number of blocks and number of levels are used to determine whether the block and/or coarsen lowering should be applied. If the number of blocks is larger than an architecture-related threshold, block-threshold, MatRox applies block lowering. Similarly the number of levels and a coarsen-threshold is used to determine the application of coarsen lowering. These thresholds are defined to ensure there are enough parallel workloads that amortize the initial cost of launching threads. MatRox further optimizes the lowered code with low-level optimizations, using structure information, as shown in Figure 3.

4 Experimental Results

We compare the inspector and executor performance of MatRox to the corresponding parts from STRUMPACK [20], GOFMM [62], and SMASH [8], which are well-known libraries for HMatrix approximation. The inspector performance for MatRox is quite similar to the existing libraries. The resulting matrix-matrix multiply performed by the executor is much improved over existing library implementations due to improvements in data locality and parallelism.

4.1 Methodology

We select a set of datasets, i.e. points, used also in prior work and shown in Table 1 from real-world machine learning and scientific applications. Problem IDs 1-8 are machine learning datasets from the UCI repository [2] and are high-dimensional points. Problem IDs 9-13 are scientific computing points that are low-dimensional [8]. STRUMPACK only runs for small datasets, i.e. problem IDs 5, 6, 8, 13. We use a Gaussian kernel [57] with bandwidth of 5 when comparing to GOFMM and STRUMPACK. For comparisons to SMASH we use their default settings of kernel function (1/||x-y||) and the scientific pointsets ID 9-13 (SMASH only supports 1-3 dimensional points); MatRox uses the same setting when compared to SMASH. The HMatrix is multiplied with a randomly generated dense matrix W of size $N \times Q$.

Testbed architectures are Haswell with Xeon™ E5-2680v3, 12 cores, 2.5 GHz, 30MB L3, and KNL with Xeon™ Phi 7250, 68 cores, 1.4 GHz, and 34MB L3. All tools are compiled with icc/icpc 18.0.1 with −03. For BLAS/LAPACK routines we use MKL [56]. MatRox is implemented in C++ in double precision, available at: https://github.com/kobeliu85/MatRox_RU. The median of 5 executions is reported for each experiment.

When comparing to libraries we use their default settings and use the same in MatRox, e.g sampling size=32, maxi $mum\ rank=256$. MatRox-specific parameters are agg=2, p = number of physical cores, blocksize = 2 for near and blocksize = 4 for far interactions, coarsen-threshold=4, and block-threshold= number of leaf nodes. bacc is set to 1e-5 for all experiments with MatRox and the libraries and the overall accuracy, i.e. accuracy of the HMatrix-matrix multiplication, is the same in MatRox and the libraries. We choose the admissibility condition to match the library's default setting. STRUMPACK only supports HMatrix structures with a very large admissibility condition in which all off-diagonal blocks are low-rank approximated; also known as HSS. GOFMM uses a budget parameter instead of admissibility, which we also implement in MatRox. Recommended budget settings in GOFMM are 0.03 and 0, in our results we refer to the former \mathcal{H}^2 -b and the later as HSS as its structure is HSS. SMASH's default admissibility is 0.65, which we also use. To be consistent across libraries we use the Gaussian kernel for SMASH experiments; the SMASH paper [8] uses a different kernel function.



Table 1. Dataset: N is number of points, d is point dimension.

ID	1	2	3	4	5	6	7
Data	covtype	higgs	mnist	susy	letter	pen	hepmass
N	100k	100k	60k	100k	20k	11k	100k
d	54	28	780	18	16	16	28
ID	8	9	10	11	12	13	
ID Data	8 gas				12 sunflower		
		grid			sunflower		

4.2 Performance of the MatRox inspector

MatRox's inspector time is similar to that of libraries since the time for structure analysis and code generation are negligible compared to the compression time as shown in Figure 4. Structure analysis and code generation in MatRox is on average 8.7 percentage of inspection time. The compression time of STRUMPACK is larger than MatRox and GOFMM because of using a different compression method.

Figure 4 also shows that the inspector time is amortized with increasing Q (number of columns in matrix approximated K is multiplied by) because the evaluation time grows with Q. The figure compares the MatRox overall time, includes inspector and executor times, with the overall time of libraries, i.e. compression and evaluation, for all datasets and Q sizes of 1, 1K, 2K, and 4K on Haswell. The compression time for both \mathcal{H}^2 -b and HSS and for all tools will not change for Q=1 and a larger Q. For example, for susy with \mathcal{H}^2 -b, MatRox's overall Speedup vs GOFMM is $1.56\times$ for Q=1K and $2.02\times$ for Q=4K. Figure 4 does not include SMASH because SMASH only supports matrix-vector multiplication (Q=1). We compared MatRox with SMASH for Q=1 and our results show that the overall time of MatRox and its evaluation time is on average $1.1\times$ and $1.6\times$ faster.

The benefits of improving evaluation with MatRox are more for larger Qs. In scientific and machine learning applications, Q is typically large and often close to N, shown in Table 1. Examples include multigrid methods in which the coefficient matrix is multiplied by a large matrix [7], Schur complement methods in hybrid solvers [61], highorder finite-elements [17], as well as direct solvers [31]. We also ran the un-approximated matrix-matrix multiplication KW with GEMM. For the tested datasets on average MatRox's overall time is $18\times$ faster than GEMM for Q=2K; the speedups obtained with HMatrix evaluation are significantly higher than GEMM for larger Qs. We use a Q=2K for the rest of our experiments unless stated otherwise.

4.3 The performance of the MatRox executor

Figures 5 and 6 show the performance breakdown of MatRox's executor (evaluation) versus GOFMM and STRUMPACK for HSS and \mathcal{H}^2 -b on both Haswell and KNL.

The theoretical floating point operations (FLOP) in HMatrix evaluation is computed for each dataset and the same number is used for each dataset across all tools. As shown in Figure 5, MatRox's evaluation is on average 3.41× and 5.98× faster than GOFMM and STRUMPACK respectively on Haswell. For KNL, MatRox is on average 8.37× and 21.25× faster than GOFMM and STRUMPACK in order. Figure 6 shows the MatRox's executor time of \mathcal{H}^2 -b versus GOFMM in which MatRox is on average 2.98× and 4.20× faster than GOFMM on Haswell and KNL respectively.

The performance breakdown shows the effect of CDS as well as the interaction blocking and structure coarsening algorithms. To show the effect of CDS, we run both the MatRox executor, that uses CDS, and GOFMM and STRUMPACK, that use a tree-based storage format, with a single thread and label in order with CDS (seq) and TB (seq). Coarsen, block, and low-level lowering applied in MatRox are labeled with coarsen, block, and low-level. The parallel version of GOFMM and STRUMPACK use a dynamic scheduling labeled with DS in Figures 5 and 6.

The different admissibility conditions in HSS and \mathcal{H}^2 -b allows us to demonstrate MatRox's performance for different HMatrix structures. HMatrix structures differ by the number of blocks that they low-rank approximate, which changes the ratio of loops over the CTree to loops with reduction. Because in HSS no off-diagonal blocks are full-rank, loops over the CTree dominate its execution time. As a result, from Figures 5 and 6, structure coarsening contributes to a performance improvement of on average 79.2% for HSS on Haswell, which is more that the 46.8% on average improvement from structure coarsening for \mathcal{H}^2 -b. A similar trend is observed on KNL where structure coarsening improves performance on average 81.2% for HSS and on average 65.1% for \mathcal{H}^2 b. Also, while interaction blocking contributes on average 38.3% and 25.2% to the performance of the MatRox generated code for \mathcal{H}^2 -b on Haswell and KNL respectively, block lowering is never activated for HSS since the number of loops with reduction, i.e. near-far interactions, never exceeds the block-threshold. MatRox peels the last iteration of loops over the CTree. Low-level transformations lead to on average 6.28% and 4.24% performance improvement of the MatRox code in HSS and \mathcal{H}^2 -b respectively on Haswell. Since HSS is dominated by loops over the CTree, the effect of low-level transformation is also more prominent in HSS. Similar effect is visible on KNL where, low-level transformations lead to on average 12.20% and 7.72% performance improvement of the MatRox code in HSS and \mathcal{H}^2 -b respectively.

One major goal of inspection and lowering in MatRox is to improve locality. Figure 7 shows the correlation between the performance of MatRox generated code versus the cost of average memory access latency among all datasets for both HSS and \mathcal{H}^2 -b on Haswell. We measure the average memory access latency [29] that is computed based on the number of memory accesses, miss-ratio of different cache



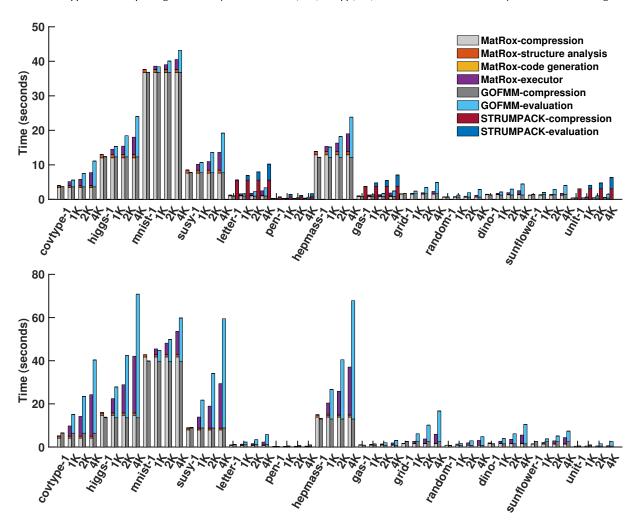


Figure 4. The overall time in MatRox, GOFMM, and STRUMPACK for different values of Q i.e., 1, 1K, 2K and 4K for HSS (top) and \mathcal{H}^2 -b (bottom) on Haswell. STRUMPACK only supports HSS and can execute for the letter, pen, gas, and unit datasets.

levels, and TLB, and use it as a proxy for locality. We use the PAPI [53] library to collect L1, LLC (Last-level Cache), TLB access and misses and number of memory accesses. The coefficient of determination or R² is 0.81 that shows a good correlation between speedup and memory access latency.

Figure 8 compares the wait time to CPU time in MatRox and GOFMM on Haswell. To obtain these times, we are using Intel's VTune Amplifier similar to [13]. CPU time is the time that the CPU takes to do HMatrix evaluation and wait time measures the time that a thread is stalled because of APIs that cause synchronization. As shown in Figure 8, the wait times in MatRox are less than 12.2% and are lower than that of GOFMM's across all datasets which demonstrates better load balance in MatRox. The MatRox generated code utilizes on average 90.5% of CPU time (a ratio of 0.095) across all datasets while GOFMM on average gives 77.6% CPU utilization.

4.4 Scalability of the MatRox executor

Figure 9 shows the scalability of MatRox executor vs. GOFMM, STRUMPACK, and SMASH for two datasets on Haswell and KNL; other datasets follow a similar trend. SMASH does not support *covtype* and in the figure, MatRox-Skernel is MatRox with SMASH settings. We select KNL in addition to Haswell to demonstrate MatRox's strong scalability on more cores, i.e. 68 cores of KNL. MatRox scales well on both architectures while the libraries show poor scalability with increasing number of cores. For example, GOFMM's performance reduces from 34 to 68 cores. MatRox's strong scaling is because structure coarsening and interaction blocking improve locality and reduce synchronization while maintaining load-balance.

5 Reusing Inspection

The modular design in MatRox enables the reuse of specific outputs of the inspector when parts of the input change. In libraries, any change to the inputs results in re-running the



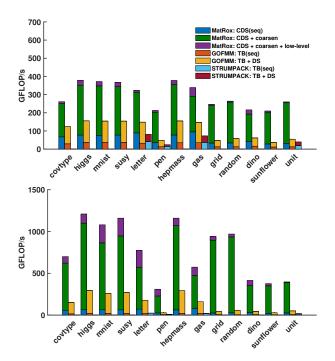


Figure 5. The performance of executor/evaluation in MatRox vs. GOFMM for HSS on Haswell (top) and KNL (bottom). Labels seq, TB, and DS are sequential, tree-based format, and dynamic scheduling respectively. Effects of CDS, structure coarsening, and low-level transformations are shown separately. Missing bars for STRUMPACK mean it could not run that dataset.

entire compression and evaluation phases. However, when the kernel function and/or the input accuracy change in MatRox, the modules and components in the inspector that do not rely on these inputs can execute only once and be reused. MatRox enables this reuse by separating the inspector into two phases, i.e. *inspector-p1* and *inspector-p2*. The inputs to inspector-p2 are the kernel function and the input accuracy and it is composed of the low-rank approximation, structure coarsening, and data layout construction modules in Figure 3. The remaining parts of the inspector in Figure 3 belong to inspector-p1. A change in the kernel function and/or the accuracy only requires inspector-p2 and the executor to be re-ran. Figure 10 shows an example code that allows for the reuse of inspector-p1 when the *bacc* changes.

In scientific and machine learning simulations, typically the input accuracy and the kernel function change more frequently than the input points and the admissibility condition. The reuse of inspector-p1 in MatRox reduces the overhead of these changes. For example in finite-elements the discretization, i.e. points, is often fixed [22], in statistical learning the training samples, i.e. points, are reused during offline training [30], and in N-body problems the CTree is only reconstructed during rebuild intervals [4, 43]. The admissibility condition also often remains the same in simulations and is known by the domain practitioner as it depends on

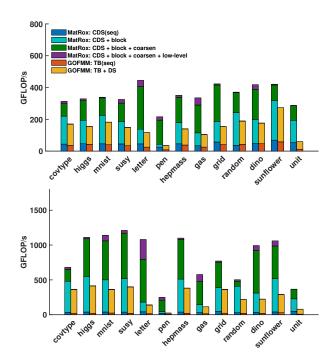


Figure 6. The performance of executor/evaluation in MatRox vs. GOFMM for \mathcal{H}^2 -b on Haswell (top) and KNL (bottom). Labels seq, TB, and DS are sequential, tree-based format, and dynamic scheduling respectively. Effects of CDS, structure coarsening, interaction blocking, and low-level transformations are shown separately.

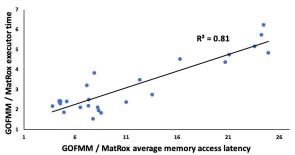


Figure 7. Effect of improving locality on MatRox's speedup vs. GOFMM on Haswell. Average memory access latency shows the average cost of accessing memory.

the problem structure. However, users often need to tune the parameters in a kernel function specially in machine learning simulations. For example, the bandwidth h in the Gaussian kernel [57] is typically tuned with cross-validation to avoid overfitting [45]. Also, often the practitioner needs to tune the input accuracy (bacc) because the *overall accuracy* of the HMatrix-matrix multiplication is not sufficient, i.e. bacc is correlated with the overall accuracy with a loose upper bound [38], or the user might decide to trade accuracy for faster evaluation (or vise-vesa) with re-compression.

Figure 11, shows the correlation between *bacc* and overall accuracy ϵ_f obtained from $\epsilon_f = \|\widetilde{K}W - KW\|_F / \|KW\|_F$



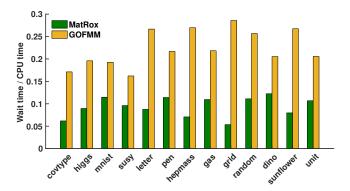


Figure 8. The ratio of wait time to the total execution time of executor/evaluation in MatRox vs. GOFMM for \mathcal{H}^2 -b on Haswell. HSS also follows a similar trend.

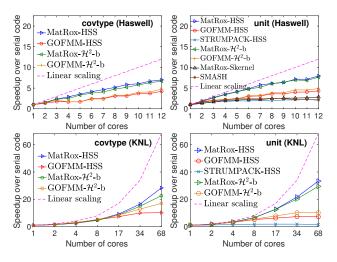
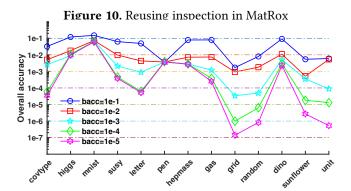


Figure 9. Scalability result on Haswell (top two) and KNL (bottom two) for datasets covtype (left) and unit (right).

for \mathcal{H}^2 -b. As demonstrated, with a *bacc* of 1e-3 more than 50% of the datasets do not reach an overall accuracy of 1e−3 and thus the user has to retune. The tuning becomes more important when more accurate results are required and also depends on the spectrum (i.e., eigenvalues) of the kernel Matrix. Figure 12 shows the MatRox's overall time compared to GOFMM for \mathcal{H}^2 -b with 5 changes to the input accuracy bacc, 1e-1 to 1e-5, with reusing inspector-p1. We do not include STRUMPACK for space but it follows a similar trend. As shown in the Figure, MatRox's overall time is on average 2.21× faster than GOFMM. For high dimensional datasets such as mnist sampling is expensive, 89.2% of the compression time in mnist, and thus the reuse of sampling in inspector-p1 leads to a speedup of 2.64× for minst compared to GOFMM. MatRox with inspector reuse vs SMASH leads to on average $1.37 \times$ speedup with up to $2.4 \times$ for sunflower.

```
1 ///MatRox inspector code for re-using inspection
  #include <matrox.h>
     //Inputs declaration
    Points points("path/to/load/points");
    Float(64) tau = 0.65;
    //Outputs declaration
    Op HMatMul("path/to/save/mat_mul.c");
    Tree ctree("path/to/store/ctree");
    Set GroupSet("path/to/store/groupset");
10
    Vector<Int(32)> sampling("path/to/store/sampling");
    Op HMatMul("path/to/store/matmul.h");
11
    inspector_p1(points,tau,&HMatMul,&ctree,&GroupSet,&
12
          sampling):
  ///MatRox executor Code for re-using inspection
15 #include "path/to/matmul.h"
    Tree ctree("path/to/load/ctree");
    Set GroupSet("path/to/load/groupset");
17
    Vector<Int(32)> sampling("path/to/load/sampling");
    Float(64) tau; Ker kfunc = GAUSSIAN;
19
    Dense W("path/to/load/matrix/W");
20
    Dense Y(Float(64), ctree[0].num_points * W.cols);
22
    // Accuracy tuning
    for(auto acc : {1e-1,1e-2,1e-3,1e-4,1e-5}){
23
      inspector_p2(kfunc,acc,sampling,ctree,GroupSet,&H);
24
      Y = matmul(H, W);
25
26
    }
```



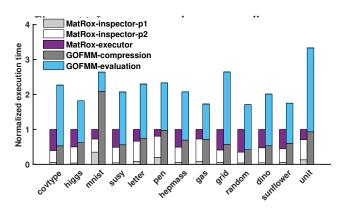


Figure 12. MatRox reusing inspection for \mathcal{H}^2 -b on Haswell.



6 Related Work

The presented approach applies a modular inspectorexecutor strategy to HMatrix Approximation. This section summarizes previous approaches to improving the performance of HMatrix approximation and related inspectorexecutor strategies that were used in other contexts.

Hierarchical matrices. Hierarchical matrices are used to approximate matrix computations in almost linear complexity. Hackbusch first introduced \mathcal{H} -matrices [6, 24], to generalize fast multipole methods [23], where the matrix is partitioned hierarchically with a cluster tree and then parts of the off-diagonal blocks are approximated. Later, \mathcal{H}^2 -structures were introduced [27] which use nested basis matrices [25], to further reduce the computational complexity of dense matrix computations. \mathcal{H}^2 has gained significant traction in recent years [6, 26]. Hierarchical semi-separable (HSS) [11, 59, 60] are a specific class of \mathcal{H}^2 structures. Matrox supports HSS and other classes of \mathcal{H}^2 using a binary cluster tree; we abbreviate \mathcal{H}^2 matrix with HMatrix.

HMatrix approximations have a compression and an evaluation phase. Numerous algorithms have been studied for HMatrix compression [3, 10, 58] including interpolative decomposition (ID) [40]. MatRox uses ID in its compression phase and contributes on improving the performance of the evaluation phase, which is the focus of many of the recent works on HMatrix computations [62, 63].

Specialized libraries for HMatrix computations. Numerous specialized libraries implement HMatrix evaluations on different platforms and for different evaluation operations. HMatrix algorithms have been implemented on platforms ranging from shared memory [20, 32, 62], distributed memory [35, 37, 51], and many-core such as GPUs [36]. Hierarchical matrices have been studied to accelerate matrix factorization [1, 11, 60]. Ghysels et. al. [19] introduces an algebraic preconditioner based on HSS. Other work has improved matrix inversion [41] and matrix-vector/matrix multiplication [11]. STRUMPACK, GOFMM, and SMASH are the most well-known libraries that support HMatrix-matrix/vector multiplications. SMASH [8] supports 1-3D datasets while GOFMM [62] and STRUMPACK [20] also support datasets of higher dimension. MatRox generates code for HMatrixmatrix/vector multiplications for datasets of all dimensions on multicore platforms.

Inspector-executor approaches. MatRox uses a domain-specific inspector-executor approach to generate code for HMatrix evaluation. Recent work [12, 21, 33, 44, 46, 48] have proposed inspector-executors that inspect the data dependency graphs in sparse matrix computations to apply code optimizations that general compilers cannot apply. Amongst them, inspectors based on level-by-level wavefront parallelism [49, 55] are the most well-known, but do not optimize for locality and load-balance. Cheshmi *et. al.* [13] present an approach to

improve wavefront inspectors, with the LBC algorithm by coarsening levels for better locality and creating balanced partitions. However, LBC only works for DAG from a specific class of sparse matrix. MatRox improves the data locality and parallelism in reduction and tree-based loops for HMatrix approximation with a novel structure coarsening method that uses a cost model of submatrix ranks and uses a specialized partition for binary trees.

7 Conclusion and Future Work

We demonstrate a novel structure analysis approach based on modular compression to generate specialized code and an efficient storage that improves the performance HMatrix approximations on multicore architectures. The proposed interaction blocking and structure coarsening algorithms, improve locality in HMatrix evaluations while maintaining load-balance. The modular approach used in MatRox, allows parts of the inspector to be reused when the kernel function and accuracy change. MatRox outperforms state-of-the-art libraries for HMatrix-matrix multiplications on different multicore processors. In future work, we plan to extend MatRox to support other kernels such as matrix factorization with support for GPUs and distributed systems, and ultimately build a domain-specific language for a large class of approximate linear algebra computations.

8 Acknowledgments

This research is supported in part by NSERC Discovery Grants (RGPIN-06516, DGECR-00303), the Canada Research Chairs program, and by the U.S. National Science Foundation (NSF) Award Numbers CCF-1657175 and CCF-1563732. The experimental platform is provided by the Extreme Science and Engineering Discovery Environment (XSEDE) [54], which is supported by the NSF grant number ACI-1548562. The research was enabled in part by Compute Canada and Scinet (http://www.computecanada.ca).

References

- Amirhossein Aminfar, Sivaram Ambikasaran, and Eric Darve. 2016. A fast block low-rank dense solver with applications to finite-element matrices. J. Comput. Phys. 304 (2016), 170–188.
- [2] Kevin Bache and Moshe Lichman. 2013. UCI Machine Learning Repository [http://archive. ics. uci. edu/ml]. Irvine, CA: University of California. School of information and computer science 28 (2013).
- [3] Mario Bebendorf and Sergej Rjasanow. 2003. Adaptive low-rank approximation of collocation matrices. *Computing* 70, 1 (2003), 1–24.
- [4] Jeroen Bédorf, Evghenii Gaburov, and Simon Portegies Zwart. 2012. A sparse octree gravitational N-body code that runs entirely on the GPU processor. J. Comput. Phys. 231, 7 (2012), 2825–2839.
- [5] Steffen Börm and Jochen Garcke. 2007. Approximating Gaussian Processes with H²-Matrices. In European Conference on Machine Learning. Springer 42–53
- [6] Steffen Börm, Lars Grasedyck, and Wolfgang Hackbusch. 2003. Introduction to hierarchical matrices with applications. *Engineering analysis with boundary elements* 27, 5 (2003), 405–422.



- [7] William L Briggs, Steve F McCormick, et al. 2000. A multigrid tutorial. Vol. 72. Siam.
- [8] Difeng Cai, Edmond Chow, Lucas Erlandson, Yousef Saad, and Yuanzhe Xi. 2018. SMASH: Structured matrix approximation by separation and hierarchy. Numerical Linear Algebra with Applications 25, 6 (2018), e2204.
- [9] Nicola Cancedda, Eric Gaussier, Cyril Goutte, and Jean-Michel Renders. 2003. Word-sequence kernels. Journal of machine learning research 3, Feb (2003), 1059–1082.
- [10] Tony F Chan. 1987. Rank revealing QR factorizations. Linear algebra and its applications 88 (1987), 67–82.
- [11] Shiv Chandrasekaran, Ming Gu, and Timothy Pals. 2006. A fast ULV decomposition solver for hierarchically semiseparable representations. SIAM J. Matrix Anal. Appl. 28, 3 (2006), 603–622.
- [12] Kazem Cheshmi, Shoaib Kamil, Michelle Mills Strout, and Maryam Mehri Dehnavi. 2017. Sympiler: transforming sparse matrix codes by decoupling symbolic analysis. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. ACM, 13.
- [13] Kazem Cheshmi, Shoaib Kamil, Michelle Mills Strout, and Maryam Mehri Dehnavi. 2018. ParSy: inspection and transformation of sparse matrix computations for parallelism. In Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis. IEEE Press, 62.
- [14] Edward G Coffman, Jr, Michael R Garey, and David S Johnson. 1978. An application of bin-packing to multiprocessor scheduling. SIAM J. Comput. 7, 1 (1978), 1–17.
- [15] Sanjoy Dasgupta and Yoav Freund. 2008. Random projection trees and low dimensional manifolds.. In *STOC*, Vol. 8. Citeseer, 537–546.
- [16] Yi Ding, Risi Kondor, and Jonathan Eskreis-Winkler. 2017. Multiresolution kernel approximation for Gaussian process regression. In Advances in Neural Information Processing Systems. 3740–3748.
- [17] Tingxing Dong, Veselin Dobrev, Tzanio Kolev, Robert Rieben, Stanimire Tomov, and Jack Dongarra. 2014. A step towards energy efficient computing: Redesigning a hydrodynamic application on CPU-GPU. In 2014 IEEE 28th International Parallel and Distributed Processing Symposium. IEEE. 972–981.
- [18] Shai Fine and Katya Scheinberg. 2001. Efficient SVM training using low-rank kernel representations. *Journal of Machine Learning Research* 2, Dec (2001), 243–264.
- [19] Pieter Ghysels, Xiaoye Sherry Li, Christopher Gorman, and François-Henry Rouet. 2017. A robust parallel preconditioner for indefinite systems using hierarchical matrices and randomized sampling. In Parallel and Distributed Processing Symposium (IPDPS), 2017 IEEE International. IEEE, 897–906.
- [20] Pieter Ghysels, Xiaoye S Li, François-Henry Rouet, Samuel Williams, and Artem Napov. 2016. An efficient multicore implementation of a novel HSS-structured multifrontal solver using randomized sampling. SIAM Journal on Scientific Computing 38, 5 (2016), S358–S384.
- [21] R Govindarajan and Jayvant Anantpur. 2013. Runtime dependence computation and execution of loops on heterogeneous systems. In Proceedings of the 2013 IEEE/ACM International Symposium on Code Generation and Optimization (CGO). IEEE Computer Society, 1–10.
- [22] Lars Grasedyck, Ronald Kriemann, and Sabine Le Borne. 2008. Parallel black box H-LU preconditioning for elliptic boundary value problems. Computing and visualization in science 11, 4-6 (2008), 273–291.
- [23] Leslie Greengard and Vladimir Rokhlin. 1987. A fast algorithm for particle simulations. Journal of computational physics 73, 2 (1987), 325–348.
- [24] Wolfgang Hackbusch. 1999. A Sparse Matrix Arithmetic Based on H-Matrices. Part I: Introduction to H-Matrices. Computing 62, 2 (1999), 89–108.
- [25] Wolfgang Hackbusch. 2015. Hierarchical matrices: algorithms and analysis. Vol. 49. Springer.

- [26] Wolfgang Hackbusch and Steffen Börm. 2002. Data-sparse approximation by adaptive H2-matrices. Computing 69, 1 (2002), 1–35.
- [27] W Hackbusch, B Khoromskij, and SA Sauter. 2000. On H2-matrices: Lectures on applied mathematics.
- [28] Wolfgang Hackbusch, Boris N Khoromskij, and Ronald Kriemann. 2004. Hierarchical matrices based on a weak admissibility criterion. Computing 73, 3 (2004), 207–243.
- [29] John L Hennessy and David A Patterson. 2017. *Computer architecture:* a quantitative approach. Elsevier.
- [30] Thomas Hofmann, Bernhard Schölkopf, and Alexander J Smola. 2008. Kernel methods in machine learning. The annals of statistics (2008), 1171–1220.
- [31] Eun-Jin Im, Katherine Yelick, and Richard Vuduc. 2004. Sparsity: Optimization framework for sparse matrix kernels. *The International Journal of High Performance Computing Applications* 18, 1 (2004), 135–158.
- [32] Ronald Kriemann. 2005. Parallel-matrix arithmetics on shared memory systems. Computing 74, 3 (2005), 273–297.
- [33] Weifeng Liu, Ang Li, Jonathan Hogg, Iain S Duff, and Brian Vinter. 2016. A synchronization-free algorithm for parallel sparse triangular solves. In European Conference on Parallel Processing. Springer, 617–630.
- [34] William B March and George Biros. 2017. Far-field compression for fast kernel summation methods in high dimensions. Applied and Computational Harmonic Analysis 43, 1 (2017), 39–75.
- [35] William B March, Bo Xiao, and George Biros. 2015. ASKIT: Approximate skeletonization kernel-independent treecode in high dimensions. SIAM Journal on Scientific Computing 37, 2 (2015), A1089–A1110.
- [36] William B March, Bo Xiao, D Yu Chenhan, and George Biros. 2015. An algebraic parallel treecode in arbitrary dimensions. In Parallel and Distributed Processing Symposium (IPDPS), 2015 IEEE International. IEEE, 571–580
- [37] William B March, Bo Xiao, Sameer Tharakan, D Yu Chenhan, and George Biros. 2015. A kernel-independent FMM in general dimensions. In High Performance Computing, Networking, Storage and Analysis, 2015 SC-International Conference for. IEEE, 1–12.
- [38] William B March, Bo Xiao, Sameer Tharakan, Chenhan D Yu, and George Biros. 2015. Robust treecode approximation for kernel machines. In Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM, 775–784.
- [39] William B March, Bo Xiao, Chenhan D Yu, and George Biros. 2016. ASKIT: an efficient, parallel library for high-dimensional kernel summations. SIAM Journal on Scientific Computing 38, 5 (2016), S720–S749.
- [40] Per-Gunnar Martinsson. 2011. A fast randomized algorithm for computing a hierarchically semiseparable representation of a matrix. SIAM 7. Matrix Anal. Appl. 32, 4 (2011), 1251–1274.
- [41] Per-Gunnar Martinsson and Vladimir Rokhlin. 2005. A fast direct solver for boundary integral equations in two dimensions. J. Comput. Phys. 205, 1 (2005), 1–23.
- [42] Per-Gunnar Martinsson, Vladimir Rokhlin, and Mark Tygert. 2011. A randomized algorithm for the decomposition of matrices. Applied and Computational Harmonic Analysis 30, 1 (2011), 47–68.
- [43] Yohei Miki and Masayuki Umemura. 2017. GOTHIC: Gravitational oct-tree code accelerated by hierarchical time step controlling. New Astronomy 52 (2017), 65–81.
- [44] Mahdi Soltan Mohammadi, Tomofumi Yuki, Kazem Cheshmi, Eddie C Davis, Mary Hall, Maryam Mehri Dehnavi, Payal Nandy, Catherine Olschanowsky, Anand Venkat, and Michelle Mills Strout. 2019. Sparse computation data dependence simplification for efficient compilergenerated inspectors. In Proceedings of the 40th ACM SIGPLAN Conference on Programming Language Design and Implementation. ACM, 594–609.
- [45] Vlad I Morariu, Balaji V Srinivasan, Vikas C Raykar, Ramani Duraiswami, and Larry S Davis. 2009. Automatic online tuning for fast Gaussian summation. In Advances in neural information processing



- systems. 1113-1120.
- [46] Maxim Naumov. 2011. Parallel solution of sparse triangular linear systems in the preconditioned iterative methods on the GPU. NVIDIA Corp., Westford, MA, USA, Tech. Rep. NVR-2011 1 (2011).
- [47] Stephen M Omohundro. 1989. Five balltree construction algorithms. International Computer Science Institute Berkeley.
- [48] Jongsoo Park, Mikhail Smelyanskiy, Narayanan Sundaram, and Pradeep Dubey. 2014. Sparsifying synchronization for highperformance shared-memory sparse triangular solver. In *International Supercomputing Conference*. Springer, 124–140.
- [49] Lawrence Rauchwerger, Nancy M Amato, and David A Padua. 1995. Run-time methods for parallelizing partially parallel loops. In Proceedings of the 9th international conference on Supercomputing. ACM, 137–146.
- [50] Elizaveta Rebrova, Gustavo Chávez, Yang Liu, Pieter Ghysels, and Xiaoye Sherry Li. 2018. A study of clustering techniques and hierarchical matrix formats for kernel ridge regression. In 2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW). IEEE, 883–892.
- [51] François-Henry Rouet, Xiaoye S Li, Pieter Ghysels, and Artem Napov. 2016. A distributed-memory package for dense hierarchically semiseparable matrix computations using randomization. ACM Transactions on Mathematical Software (TOMS) 42, 4 (2016), 27.
- [52] Ana R Teixeira, Ana Maria Tomé, and Elmar Wolfgang Lang. 2008. Feature extraction using low-rank approximations of the kernel matrix. In *International Conference Image Analysis and Recognition*. Springer, 404–412.
- [53] Dan Terpstra, Heike Jagode, Haihang You, and Jack Dongarra. 2010. Collecting performance data with PAPI-C. In Tools for High Performance Computing 2009. Springer, 157–173.
- [54] J. Towns, T. Cockerill, M. Dahan, I. Foster, K. Gaither, A. Grimshaw, V. Hazlewood, S. Lathrop, D. Lifka, G. D. Peterson, R. Roskies, J. R. Scott, and N. Wilkins-Diehr. 2014. XSEDE: Accelerating Scientific Discovery. Computing in Science and Engineering 16, 5 (Sept.-Oct. 2014), 62–74. https://doi.org/10.1109/MCSE.2014.80
- [55] Anand Venkat, Mahdi Soltan Mohammadi, Jongsoo Park, Hongbo Rong, Rajkishore Barik, Michelle Mills Strout, and Mary Hall. 2016. Automating wavefront parallelization for sparse matrix computations. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. IEEE Press, 41.
- [56] Endong Wang, Qing Zhang, Bo Shen, Guangyong Zhang, Xiaowei Lu, Qing Wu, and Yajuan Wang. 2014. Intel math kernel library. In High-Performance Computing on the Intel® Xeon PhiTM. Springer, 167–188.
- [57] Christopher KI Williams and Carl Edward Rasmussen. 1996. Gaussian processes for regression. In Advances in neural information processing systems. 514–520.
- [58] Christopher KI Williams and Matthias Seeger. 2001. Using the Nyström method to speed up kernel machines. In Advances in neural information processing systems. 682–688.
- [59] Yuanzhe Xi and Jianlin Xia. 2016. On the stability of some hierarchical rank structured matrix algorithms. SIAM J. Matrix Anal. Appl. 37, 3 (2016), 1279–1303.
- [60] Jianlin Xia, Shivkumar Chandrasekaran, Ming Gu, and Xiaoye S Li. 2010. Fast algorithms for hierarchically semiseparable matrices. Numerical Linear Algebra with Applications 17, 6 (2010), 953–976.
- [61] Ichitaro Yamazaki and Xiaoye S Li. 2010. On techniques to improve robustness and scalability of a parallel hybrid linear solver. In *Interna*tional Conference on High Performance Computing for Computational Science. Springer, 421–434.
- [62] Chenhan D Yu, James Levitt, Severin Reiz, and George Biros. 2017. Geometry-oblivious FMM for compressing dense SPD matrices. In Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis. ACM, 53.

[63] Chenhan D Yu, Severin Reiz, and George Biros. 2018. Distributed-memory hierarchical compression of dense SPD matrices. In Proceedings of the International Conference for High Performance Computing, Networking, Storage, and Analysis. IEEE Press, 15.

