Integrating Chapel programs and MPI-Based Libraries for High-performance Graph Analysis

Trevor McCrary* tmccrary9@gatech.edu Georgia Institute for Technology Atlanta, GA, USA Karen Devine* devine.hpc@gmail.com Sandia National Laboratories, ret. Albuquerque, NM, USA

42

43

44

45

47

48

49

50

51

52

54

55

58

61

62

63

64

65

66

67

68

70

71

72

73

76

78

70

80

81

82

83

84

85

86

87

88

Andrew Younge ajyoung@sandia.gov Sandia National Laboratories Albuquerque, NM, USA

1 ABSTRACT

We identify techniques to interface Chapel programs with parallel, 2 distributed, MPI-based libraries written in C++ without storing multiple copies of shared data. This integration enables Chapel users to take advantage of the vast array of capabilities developed in parallel numerical libraries without the memory cost of duplicated data. We demonstrate two approaches to interface Chapel code with the MPI-based graph and numerical solver libraries Grafiki 8 and Trilinos. The first uses a single Chapel executable to call a C function that interacts with the MPI libraries; it requires Chapel 10 users to understand the Grafiki library interfaces and link their 11 codes with the MPI-based libraries. The second uses the Unix mmap 12 function to allow separate Chapel and Grafiki executables to read 13 and write to the same block of memory on a node; it maintains 14 greater independence between the Chapel and MPI-based codes, 15 simplifying the Chapel user's experience. We also encapsulated the 16 second approach in Docker/Singularity containers to maximize ease 17 of use by Chapel users. Comparisons of the two approaches using 18 shared and distributed memory installations of Chapel show that both approaches are feasible for sharing data between Chapel and 20 MPI-based libraries, yielding similar scalability and performance 21 with no penalty from using mmap. 22

23 KEYWORDS

24 Chapel, MPI, interoperability, mmap, container

25 ACM Reference Format:

26 Trevor McCrary, Karen Devine, and Andrew Younge. 2022. Integrating

27 Chapel programs and MPI-Based Libraries for High-performance Graph

28 Analysis. In Proceedings of The 2022 Chapel Implementers and Users Workshop

29 (CHIUW22). ACM, New York, NY, USA, 10 pages.

30 1 INTRODUCTION

We present two methods that allow data to be shared, rather 31 than copied, between partitioned global address space programs in 32 Chapel and distributed memory, parallel algorithms that are written 33 in C++ and use the Message Passing Interface (MPI) [12] library 34 to exchange data. We demonstrate our methods using applications 35 in graph analysis: a simple graph connected-component algorithm 36 in Chapel and a graph hitting-times algorithm in the graph toolkit 37 Grafiki [15]. However, the general capability to integrate Chapel 38 and MPI-based libraries is valuable in many applications, as it com-39 bines the simplicity of Chapel programming with the speed and 40 efficiency of existing high-performance MPI-based algorithms. 41

*Work completed at Sandia National Laboratories, Summer 2021. [10]

CHIUW22, June 10, 2022, Virtual

Chapel [3, 7] is a partitioned global address space (PGAS) language that is built for productive parallel programming at scale. Chapel simplifies parallel programming by providing thread-based parallel loops and managing the layout of arrays within a parallel computer's memory. Although Chapel arrays can be distributed across processors, Chapel users can access any array entry on any processor, without knowing on which processor the data is stored. Chapel manages the data movement or communication required to retrieve array entries. Thus, Chapel provides a very easy-to-use programming environment for parallel algorithm development; algorithms such as parallel graph connected-component labeling can be implemented in just a few lines of code.

In MPI-based libraries, data is also distributed across the memory spaces of parallel processors. The distribution, however, is determined by the programmer. Moreover, each processor can access only the data in its memory. Off-processor data must be communicated explicitly via message passing. This explicit control of data distribution and movement allows highly efficient parallel execution, but requires much more effort from the programmer. For example, Grafiki (the successor of TriData) [15] is an MPI-based library of high-performance parallel graph analysis algorithms with linear solvers from the Trilinos [6, 14] toolkit. Grafiki's graph computations are done via Trilinos' matrix-vector operations, with explicit inter-processor communication via MPI send/receive operations.

Integrating Chapel and Grafiki allows graph analysts to easily filter and manipulate graph data via Chapel, and then apply Grafiki's high-performance parallel algorithms to analyze the resulting data. The main challenges in the integration are avoiding duplication of data between Chapel and Grafiki, and insulating graph analysts from the complications of using external MPI-based libraries.

Sharing, rather than copying, data between Chapel and MPIbased libraries is crucial to success of the integration. Doubling the amount of memory required to couple two algorithms is often infeasible. In our demonstration, for example, Chapel users wish to analyze the largest connected component of graphs that fill much of their computer's memory; the available memory is insufficient to copy this component in memory for analysis in Grafiki. Our approaches rely on library interfaces with sufficient data abstraction to accommodate Chapel data structures (e.g., edge lists in Chapel arrays) without requiring data to be copied and/or reorganized for Grafiki (e.g., into compressed-sparse row graphs).

Chapel users should also be insulated as much as possible from the complexities of building external MPI-based libraries. One of our approaches addresses this issue by demonstrating how Chapel and Grafiki can share data through memory-mapped regions. Grafiki analysis can then be run as a separate executable or within a Docker container provided to the Chapel users.

The contributions of this work include two approaches for inte-142 89 grating Chapel programs with MPI-based C++ libraries. The first 90 143 approach uses Chapel's interface to C-language functions to share 91 144 data between Chapel and the MPI-based libraries. For the second 92 approach, we developed a new Chapel data distribution Domain 146 93 that uses memory-mapped regions to share data with external pro-147 94 cesses; to our knowledge, this use of memory-mapped regions in 148 95 Chapel Domains has not been done before. We describe how to 96 149 extend the second approach for use in Docker containers, for the 97 greatest insulation of Chapel users from the details of the MPI-151 98 based libraries. We demonstrate our approaches using a simple 152 99 graph connected-components algorithm in Chapel and the high-100 153 performance graph hitting-times algorithm in Grafiki. We run on 154 101 two different Chapel environments: a single-node shared memory 102 environment and a multi-node distributed memory environment. 103 156 Our performance experiments demonstrate scalability of our meth-104 157 ods, showing that our methods can successfully integrate Chapel 105 158 and high-performance MPI-based parallel algorithms. 106 159

BACKGROUND 2 107

We describe the MPI, Chapel, Grafiki, Trilinos and Unix features 108 that are important to our integration approaches.

MPI: The Message Passing Interface (MPI) enables sharing of data 110 between distributed memory spaces. Each processor is assigned a 111 unique "rank" in [0, NumberOf Processors). Processors can access 112 only their local memory; they cannot access memory associated 113 with other processors. Instead, shared data must be communicated 114 across the interconnect network of the cluster or supercomputer. 171 115 The MPI library sends buffers of data across network links from 116 the processors owning the data to the processors needing it. 117

Chapel: Chapel uses the concept of locales to associate data with 118 individual processors for performance and scalability [7]. Chapel's 119 locales are analogous to MPI ranks. Chapel's global address space 120 allows locales to access and manipulate data that are stored on other 121 locales without explicit communication by the Chapel user; needed 122 communication is performed by the Chapel runtime environment. 123 Thus, underlying communication causes data access on a different 124 locale to be more expensive than data access on the data's own 125 locale. While Chapel can be built with an MPI back-end to do its 126 communication, an MPI back-end in Chapel is not necessary for 127 our proposed approach. In our work, we rely only on a one-to-one 128 mapping between processors' Chapel locales and MPI ranks. 129

Chapel distributed arrays are the main data structure used in 130 this work. We store lists of graph edges and vertices in Chapel 131 distributed arrays. A distributed array is implemented as a collection 132 of arrays, with one local array stored on each locale. The Chapel 133 distributed array manages the global address space indexing that 134 allows access of any array entry from any locale. Chapel's Domain 135 maps describe and manage the distribution of arrays to processes. 136 The Block Domain provides a commonly used distribution; in it, 137 the indices are partitioned evenly across the locales so that the first 138 locale has the first contiguous chunk of indices, the second locale 139 has the next chunk, and so on. An example Block array with 64 140 elements distributed across four locales in depicted in Figure 1. 141

Chapel users can create their own distributions to fit their requirements. To create a custom distribution, users create custom Domain map classes that implement the Domain map Standard Interface (DSI) [1, 2]. For our second approach below, we create a custom Domain map by modifying Chapel's Block distribution to use memory that can be shared among processes.

Chapel supports interoperability with C code. Users can access C libraries, variables, functions, structures, and constants using the extern keyword. Chapel programs can call C functions, and addresses of Chapel arrays can be passed to C functions. We use this capability in our first approach below. Chapel also has a library of definitions for C datatypes (e.g., long long int). We use these datatypes for consistency between our Chapel and C code.

Grafiki and Trilinos: Grafiki (formerly called TriData) [15] is a library of high-performance graph and hypergraph analysis algorithms written in C++. It contains algorithms for computing hitting times, spectral clustering, and eigenvector centrality. For our demonstration, we use Grafiki's hitting time algorithm, which operates on a square, symmetric matrix that may be distributed across processors. The symmetric matrix represents the adjacency matrix of a connected graph (i.e., the graph has a single connected component). Each edge (i, j) in the Chapel data corresponds to a nonzero a_{ii} in the adjacency matrix A; each vertex corresponds to a row and column of the matrix.

160

161

162

163

164

165

166

167

168

169

170

172

173

174

176

177

178

182

183

184

185

186

189

192

193

194

195

196

Grafiki's algorithms rely on linear and eigen-solvers; for these solvers as well as for abstractions of matrix and vector operations, Grafiki uses the open-source Trilinos [6, 14] framework. Trilinos has been developed and optimized for distributed memory, shared memory, and GPU parallel performance, especially in the realm of physics-based scientific simulations. Trilinos can operate with arbitrary data distributions, including two-dimensional matrix distributions favored for reducing communication in graph analysis applications. In this work, we pass the edge and vertex lists to Grafiki with the same distribution as specified by the user in Chapel; thus, the distribution of nonzeros to processors is arbitrary and matches that of the Chapel edge list distribution.

Trilinos provides an efficient compressed sparse row matrix (CrsMatrix) data structure, but creation of a CrsMatrix from Chapel data would require creating a copy of the data in CRS format - an unacceptable requirement for this project. However, Trilinos also provides a RowMatrix abstraction that allows users to implement matrix operations on their data using their own underlying data structures. The RowMatrix abstraction supports any distribution of sparse matrix entries across processors. It does not require that matrix entries be sorted in any particular manner. Performance of a user's RowMatrix strongly depends on the user's implementation and data distribution. In our work, we sacrifice some computational performance by allowing our RowMatrix to use the Chapel edge lists directly, without any reordering or reorganization, thus avoiding an additional copy of Chapel users' data.

Unix: The Unix mmap function can be used to create memorymapped regions that can be shared by independent Unix processes. With this function, two separate processes can read and write to the same block of memory by using the Unix MAP_SHARED flag. The mmap function works with the function shm_open, which takes



Figure 1: Example of a Block distribution across four locales of a Chapel distributed array with 64 elements; elements are distributed evenly across locales, with a contiguous chunk of indices assigned to each locale.

a backing file name (C string) for mmap. This backing file name 197 connects mmap calls on separate processes to the same block of 198 memory. Chapel has its own mmap function, sys_mmap, which takes 100 the same arguments as the Unix function, invokes the system mmap 200 function and returns an error code. Because we used the Unix mmap 201 function in sys/mman.h in our early explorations, we continued 202 using it with Chapel, specifying the extern keyword to refer to 203 the Unix mmap function. However, we expect Chapel's sys_mmap 204 would work identically. In our second approach below, we use 205 mmap capability to share edge and vertex lists in mapped memory 206 between separate Chapel and C++ processes or containers. 207

208 3 METHODOLOGY

In the general use case that we wish to support, graph analysts load
graph data into Chapel edge lists. They then filter the graph data in
some way to identify vertices and edges of interest. The resulting
subgraph is shared with an MPI-based library for further analysis,
and results are shared back to Chapel. Our goal is to accomplish
this workflow without copying or reformatting data that is to be
shared between Chapel and the MPI-based library.

Our demonstration follows the pattern of this general use 216 case. We load a graph from files that contain coordinate pairs 217 $(source_k, target_k)$ and, optionally, a weight w_k for all edges k in 218 the graph (i.e., nonzeros in an adjacency matrix). The edges are 219 stored in Chapel arrays *source* and *target* distributed across locales 220 243 with Chapel's Block distribution. We filter the graph by identifying 221 244 its largest connected component via a parallel label-propagation 222 algorithm written in a few lines of Chapel code. We then share the 245 223 246 edge lists, along with an array identifying vertices in the largest 224 connected component, with the MPI-based Grafiki library. Grafiki 225 computes vertex hitting times on the largest connected component. ²⁴⁷ 226 Both of our integration approaches create a Trilinos RowMatrix 248 227 (in C++) to describe the matrix that is passed to Grafiki. The 249 228 RowMatrix class directly uses our shared, distributed edge lists 250 229 to answer queries about the matrix and perform matrix operations. 251 230 The most important method of RowMatrix is sparse-matrix vector 252 231 multiplication (SpMV): $y = \alpha Ax + \beta y$ for matrix *A*, vectors *x* and *y*, ²⁵³ 232 and constants α and β . In parallel with distributed data, SpMV re- ²⁵⁴ 233 quires communication of off-processor vector entries x_i to be used 234 in multiplication with local matrix entries a_{ij} , and of subproducts 256 235 $a_{ij}x_i$ to be accumulated into vector entries y_i . This communication 257 236 is done via MPI using the communication operations from Trili-258 237 nos' CrsMatrix implementation. However, the localApply – the 259 238

²³⁹ on-processor SpMV operation — differs from Trilinos' CrsMatrix ²⁶⁰

240 since our RowMatrix uses coordinate-formatted edge lists from 26

241 Chapel rather than Trilinos' compressed-sparse row format. Our 262





localApply uses a straightforward loop over the edge lists, multiplying each edge value for edge ($source_k, target_k$) with the appropriate, possibly communicated, x vector entry x_{target_k} . For our demonstration, we did not attempt to optimize this operation, but several optimizations (e.g., threading, GPU parallelism) are possible.

3.1 Method One: Direct Chapel-to-C calls

For our first method, we use a straightforward approach in which a single Chapel executable reads the graph data file, performs the connected-component label propagation, and then calls directly to a C function that then calls a C++ function that calls Grafiki. A high-level flow chart of its execution is in Figure 2.

This method takes advantage of Chapel's C interoperability to directly call C code. An overview of the code is shown in Figure 3. We start by having Chapel call a short C function glueC wrapped in extern "C" controls to allow it to be compiled by a C++ compiler. The C function then calls a C++ function glueChapelGrafiki with the same arguments. Function glueChapelGrafiki instantiates a RowMatrix object using the Chapel edge lists, and calls Grafiki. The C and C++ code are in a file separate from the Chapel code. The file is compiled by a C++ compiler and its object file is linked with the Chapel object file at build time. Integrating Chapel programs and MPI-Based Libraries for High-performance Graph Analysis

```
// In Chapel source code
                                                           285
require "glue.h";
                                                           286
extern proc glueC(nEdges, *src, *tgt, ...);
                                                           287
                                                           288
coforall loc in Locales {
                                                           289
  on loc {
                                                           290
    var mySrc = source.localSubdomain();
                                                           291
    var myTgt = target.localSubdomain();
                                                           292
    glueC(mySrc.size:size t,
                                                           293
           c_ptrTo(source[mySrc.low]),
                                                           294
           c_ptrTo(target[myTgt.low]), ...);
                                                           295
  }
                                                           296
}
                                                           297
                                                           298
// In glue.h file
                                                           299
extern "C" {
                                                           300
  int glueC(nEdge, *src, *tgt, ...);
                                                           301
}
                                                           302
                                                           303
// In glue.cpp C++ file
                                                           304
int glueChapelGrafiki(nEdge, *src, *tgt, ...) {
                                                           305
  Build RowMatrix (nEdge, src, tgt, ...);
                                                            306
  return Grafiki_hittingTimesDriver(RowMatrix, ...);<sub>307</sub>
}
                                                           309
extern "C" {
                                                           310
  int glueC(nEdge, *src, *tgt, ...) {
                                                           311
    return glueChapelGrafiki(nEdge, src, tgt, ...);
                                                           312
  }
                                                           313
}
                                                           314
```

Figure 3: Method One: Chapel's path to call C++ code directly.

315

316

317

318

319

320

321

322

323

324

325

330

331

332

333

334

335

336

338

339

Each Chapel locale must call the C function in parallel, using Chapel's coforall parallel-for loop and its built-in *Locales* array as in Figure 3. The coforall loop creates a parallel task *loc* for each locale, and the subsequent on *loc* directive ensures that each locale calls glueC independently. This parallelism is needed to prevent Grafiki from hanging in collective MPI communications; all locales (i.e., all ranks) must participate in the call to glueC.

The short C function passes pointers from Chapel to the MPIbased library. Each locale can provide, for example, the number of edges and vertices in the locale, and pointers to its local *source* and *target* arrays. The appropriate way to get the pointer to the local arrays is shown in Figure 3: Chapel's c_ptrTo function obtains the C pointer to the lowest-indexed value of the array in the locale.

This approach is a simple path for integrating MPI-based C++ 276 libraries with Chapel applications, but it has several drawbacks. It 277 is intrusive to Chapel algorithm development. In order for Chapel 278 to directly call to Grafiki, Chapel users must understand how to 279 call C functions from Chapel. They must ensure that their Chapel 280 installation, Grafiki library, and Trilinos library were built all with 281 compatible compilers and MPI libraries. They also need to deter-282 mine how to link the Grafiki and Trilinos libraries with their Chapel 283 executable. Our next approach reduces this burden for Chapel users. 284

3.2 Method Two: Separate Chapel and C++ Processes

In our second method, a Chapel process interacts with an independent, concurrently running, C++ process, as outlined in Figure 4. The C++ process waits for a semaphore to be posted by Chapel before starting computation. The Chapel program reads edges into arrays that use a new shareBlock Domain that allows them to be shared with the C++ program through a mmap memory region. It filters the data (in our case, labeling connected components). It then writes metadata about the edge lists to a file and sets the semaphore indicating the C++ program can begin. The C++ program reads the metadata file to get information on accessing the shared data. It then creates a RowMatrix with the shared data, and calls Grafiki as in our first method. Upon completion, it resets the semaphore, indicating to the Chapel program that the Grafiki result is ready and Chapel can proceed. Details of each step follow.

To enable the Chapel data to be shared with the C++ program, we created a modified version of Chapel's Block Domain (from Chapel v1.23) in which the local arrays are built using mmap memory. We refer to this modified Domain as a shareBlock Domain.

Chapel code for our new shareBlock Domain is shown in Figure 5. The Chapel Block Domain stores the local arrays in a variable named myElems; this variable is assigned with a call to Chapel's domain.buildArray(). In our shareBlock Domain, we replace domain.buildArray() with a function createMmap() (shown in red) that calls the Unix mmap function to allocate memory of the requested size (the element's data type size times the domain size).

Each locale's mmap local array uses a separate backing file. The backing file name is generated by a new function getBackingFile(); it is a string consisting of the locale's ID and a counter to indicate which shareBlock is stored there. For example, the first shareBlock array allocated across four locales uses backing files /share.bak0-1, /share.bak1-1, /share.bak2-1, and /share.bak3-1. After all of the locales' local arrays are allocated with mmap, we increment a counter *Status* that is used to differentiate between separate shareBlock arrays' backing files.

To use the mmap region for the local array, we pass its pointer to modified versions of Chapel's makeArrayFromExternArray makeArrayFromPtr and (see Figure 6). Function makeArrayFromPtr is a short function calls makeArrayFromExternArray. that Function makeArrayFromExternArray creates a Domain for the new array and returns a call to Chapel's _newArray(), returning a new Chapel local array. We modified makeArrayFromPtr and makeArrayFromExternArray to accept a Domain as an argument. The Domain was needed because the original functions created a new Domain from 0 to the size of the array minus one, while the Block Domain expected indices matching the global indices of the array. Without the Domain argument, Chapel would assign a copy of the mmap region to the myElems array, rather than use the mmap region directly. Our modified shareMakeArrayFromExternArray uses the same Domain as myElems, so we avoid making a copy, and myElems refers directly to the mmap memory.

For the C++ program to access the mmap regions, it needs to know some metadata about the arrays using shareBlock Domains. Specifically, it needs the names of the backing files associated with

Integrating Chapel programs and MPI-Based Libraries for High-performance Graph Analysis



Figure 4: Method Two: Separate, concurrently running Chapel and C++ executables share graph data in mmap memory; execution is coordinated by a semaphore and a small metadata file.

389

390

393

each locale's local arrays and the size of those arrays. The Chapel 368 341 program writes these fields to a metadata file, as in Figure 7. (While 342 369 we used a regular file for the metadata, one could easily use a mmap 343 370 memory space to share the metadata instead. But because the file 371 344 is very small, using a regular file is feasible and straightforward.) 345 This example file shows a matrix with 25 vertices and 105 edges 372 346 distributed across four locales, with backing arrays for the source 347 and *target* edge lists and the *component* labels. Locale 0 has 27 edges 348 and 7 vertices, while other locales each have 26 edges and 6 vertices. 349 While the Chapel program runs, the C++ executable waits for 350 376 a semaphore to indicate that the data is ready for use. After the 351 377 Chapel program posts the semaphore to flag the C++ program to 352 378 begin, the processor with rank 0 in the C++ program reads the entire 353 379 metadata file into a buffer and broadcasts it to the other processors. 354 380 Each processor then parses the buffer to extract the information 355 381 (backing file names and data sizes) associated with its rank. The 356 382 processors then set up pointers to their mmap data and call the MPI-357 383 based libraries with the data. In our demonstration, the processors 358 384 creates the *source*, *target*, and *component* arrays using the backing 359 385 files, collectively create a distributed RowMatrix wrapping the data, 360 and call Grafiki's hitting times algorithm. 361 387 The key advantage of this method is that it is less intrusive to 362 388

Chapel developers. Developers can declare distributed arrays using shareBlock, and the shareBlock class handles creating the mmap arrays with the required size. In addition, this method allows a C++ process to be built and run separately from the Chapel program, eliminating the need for Chapel programmers to link their code with MPI-based libraries. The disadvantage, for now, is that this method is restricted to the block distribution; extensions to other Chapel distributions are feasible but not yet done. Also, currently, we cannot reshape a shareBlock Domain after it is initialized.

3.3 Containerization of the C++ Library

To make Grafiki even easier to use by Chapel programmers, we have containerized our C++ glue program with Grafiki, and shown that a Chapel program and the Grafiki container can share memory through the same mmap array mechanism. This containerization allows Grafiki developers to provide tools that are fully encapsulated, sparing Chapel users from have to compile the Grafiki executable. Our Docker container encapsulates our C++ glue program, Grafiki and all of the libraries on which Grafiki depends (Trilinos, Kokkos, MPI, BLAS). We use Singularity [9], a container system designed for high-performance parallel computing, to instantiate our container with MPI parallelism. Our workflow then proceeds as in Method Two (Figure 4), with the user's Chapel program and and the Grafiki C++ container accessing the same mmap regions.

One challenge in integrating Chapel with containers was our use of semaphores. While mmap memory works across containers, POSIX semaphores do not. Semaphores are created on the parent process' stack. They are unusable in the container after namespace creation, which copies the semaphore rather than addressing it. Thus, the semaphore never unlocks in the C++ program. Our solution is to implement a "pseudo-semaphore" signaling mechanism. Since mmap works across containers, we use a single integer in mmap

CHIUW22, June 10, 2022, Virtual

Integrating Chapel programs and MPI-Based Libraries for High-performance Graph Analysis

```
class LocBlockArr {
proc init() {
    // Chapel's Block domain allocated local array as follows
    // this.myElems = this.locDom.myBlock.buildArray(...);
    // New shareBlock domain allocates local array with mmap memory
    if(this.locDom.myBlock.size) {
      var myPtr = createMmap(eltType, this.locDom.myBlock.size:uint);
      this.myElems = shareMakeArrayFromPtr(myPtr, this.locDom.myBlock.size:uint, this.locDom.myBlock);
      allLocalesBarrier.barrier();
      if (here.id == 0) then
        Status += 1;
    } else {
      this.myElems = this.locDom.myBlock.buildArray(...);
    }
}
// Chapel functions to create mmap array
proc createMmap(type eltType, size:uint):c_ptr {
    var myBytes = (size * c_sizeof(eltType)): uint;
    var fd = shm_open(getBackingFile(), O_RDWR | O_CREAT | O_EXCL, accessPerms): c_int;
    ftruncate(fd, myBytes:off_t);
    var region;
    region = mmap(nil, myBytes:size_t, PROT_READ | PROT_WRITE, MAP_SHARED, fd:fd_t, 0:off_t);
    return region:c_ptr(eltType);
}
proc getBackingFile():c_string {
    return ("/share.bak" + here.id:string + "\-" + Status:string).c_str();
}
```

Figure 5: Method Two: The new shareBlock Domain allocates mmap memory (shown in red) in place of the memory allocated in Chapel's Block Domain (Chapel v1.23). If the size of the domain is nonzero, an mmap array is created to fit the size requirements of the Domain and datatype. Backing files are tracked by incrementing *Status* after all local arrays have been allocated.

419

421

477

423

474

425

426

427

120

430

394memory to mimic semaphore functionality. The value of this flag412395indicates which process should be working and which should be413396idling. Currently, we use a spin lock for this functionality.414397With containerization, the Chapel user is spared the chore of415398building Grafiki and everything on which it depends. This approach416399delivers the highest ease of use to graph analysts.417

400 4 EXPERIMENTS AND RESULTS

We demonstrate our approaches with a simple Chapel graph analy-401 sis application that shares the graph with the C++ library Grafiki. 402 Specifically, our Chapel program reads a matrix from a Matrix Mar-403 ket file, symmetrizing the matrix if necessary, and constructs lists 404 of edges corresponding to the matrix nonzeros. It then identifies 405 the largest connected component of the graph using a simple, com-406 monly used, label propagation algorithm. For label propagation, 407 each vertex's label is initialized to its vertex number. Then the 408 propagation algorithm loops over edges, giving each vertex of an 409 edge the lower-valued label of the edge's two vertices. Iteration 410 over edges is done in parallel with respect to locales and continues 411

until no vertex labels change. The Chapel code then counts the number of vertices with each label and identifies the label with the most vertices. The label of the largest component is passed, along with the edge lists and vertex component labels, to Grafiki, which computes hitting times within the largest component. While vertices that are not in the largest component remain in the edge lists passed to Grafiki, they are ignored in the RowMatrix sparse matrix-vector multiplication operation.

We tested our approaches on Sandia's Kahuna high-performance data analytics cluster. Kahuna has Dual Socket Intel E5-2683v3 2.00GHz CPUs with 28 cores and 256 GB of memory. Kahuna uses a shared-memory build of Chapel 1.23; that is, Kahuna's Chapel did not enable an MPI backend. Thus, all experiments on Kahuna could use only a single multi-core node. For Method One, in which Chapel calls Grafiki directly, only one MPI rank (and, thus, one locale) can be used. Method Two, in which separate processes are used for Grafiki and Chapel, allows more MPI parallelism, as the MPI parallelism is not tied to Chapel's shared-memory implementation; we can run Grafiki and Chapel with up to 28 ranks (locales). Integrating Chapel programs and MPI-Based Libraries for High-performance Graph Analysis

```
// Chapel's makeArrayFromPtr
                                                       //New shareMakeArrayFromPtr
proc makeArrayFromPtr(
                                                      proc shareMakeArrayFromPtr(
     value: c_ptr , num_elts: uint) {
                                                            value:c_ptr , num_elts:uint , dom:domain) {
    var data = chpl_make_external_array_ptr(
                                                           var data = chpl_make_external_array_ptr(
                     value , num_elts);
                                                                            value , num_elts);
                                                           return shareMakeArrayFromExternArray(
    return makeArrayFromExternArray(
                     data , value.eltType);
                                                                           data , value.eltType , dom);
}
                                                      }
// Chapel's makeArrayFromExternArray
                                                       //New shareMakeArrayFromExternArray
proc makeArrayFromExternArray(
                                                       proc shareMakeArrayFromExternArray(
     value: chpl_external_array, type eltType) {
                                                            value:chpl_external_array, type eltType,
    var dom = 0..number_elements \setminus -1;
                                                            dom:domain) {
    var arr = new unmanaged
                                                           var arr = new unmanaged
              DefaultRectangularArr(dom=dom, ...);
                                                                          DefaultRectangularArr (dom=dom._value, ...);
    dom.add_arr(arr, locking = false);
                                                           dom.add_arr(arr, locking=false);
    return newArray(arr);
                                                           return newArray(arr);
}
                                                      }
```

Figure 6: Method Two: Chapel's original (left) vs. our new (right) shared makeArrayFromPtr: in the shared version, a shareBlock Domain is passed as an argument to shareMakeArrayFromExternArray and DefaultRectangularArr.

source /share.bak0-1 /share.bak1-1 /share.bak2-1 /share.bak3-1 target /share.bak0-2 /share.bak1-2 /share.bak2-2 /share.bak3-2 component / share . bak0 -3 / share . bak1 -3 / share . bak2 -3 / share . bak3 -3 nEdge 27 26 26 26 nVtx 7 6 6 6

Figure 7: Method Two: Chapel shares matrix data with the separate C++ executable that calls Grafiki via a small metadata file containing per-locale backing file names and data sizes.

463

465

466

We also tested on Sandia's Mutrino Cray computing system with 455 431 100 Intel Haswell nodes with 96 GB of memory per node. Mutrino 456 432 has a distributed memory build of Chapel 1.24 with an MPI backend. 457 433 Thus, we could run both approaches with multiple ranks and locales. 458 434 To evaluate the effect of inter-node communication and provide 459 435 the greatest contrast to the shared-memory Kahuna experiments, 460 436 we chose to assign one rank (locale) per node on Mutrino. 461 437

Method One uses only one executable; Chapel and Grafiki exe- 462 438 cute on the same cores of allocated nodes. Method Two launches 439 two executables that run concurrently on the nodes (using SLURM's 440 "srun -overlap" for each executable); cores may or may not be 441 oversubscribed, depending on the number of ranks per node. For 442 Method Two, we show results with a C++ executable for Grafiki; 467 443 past experiences saw no performance degradation using containers. 468 444 We tested the performance of our methods with small and large 469 445 graphs from the SuiteSparse matrix collection [4]. The small graph, 470 446 bcsstk29.mtx, has 28 strongly connected components and a largest 471 447 component with 13.8K vertices and 620K edges. The large graph, 472 448 GAP-kron.mtx, has 78M strongly connected components (many of 473 449 them singleton vertices), and a largest component with 63M vertices 474 450 and 4.2B edges. The large test case uses over half of a Mutrino node's 475 451 memory for the source and target arrays alone (> 67 GB), making 476 452 the graph too large to be copied; on a single node, data sharing 453 477 between Chapel and Grafiki is the only way to solve the problem. 478 454

In Table 1, we show the execution times for label propagation in our connected-component algorithm written in Chapel, and the linear solve time in Grafiki's hitting time algorithm using both integration methods and the small graph bcsstk29.mtx. Results are shown for both the Kahuna and Mutrino systems. Execution times for Methods One and Two are comparable, with no significant loss of performance caused by using mapped memory in Method Two.

On both platforms, we see that adding MPI ranks accelerates the linear solve in Grafiki, with reasonable scaling even for this small graph. Adding locales also accelerates each iteration of label propagation in the Chapel-based connected component algorithm. Adding locales can increase the number of iterations required for connected component labeling, as each locale operates on a subset of the edges, slowing propagation across the full graph. This effect was seen in this small graph, in which the ordering (sorted by target vertex) of the input graph was optimal for label propagation with one locale; increasing the number of locales from one to 16 increased the number of propagation iterations from two to nine.

We ran the same experiments with the large GAP-Kron.mtx graph; execution times are in Table 2. Again, we see that Grafiki's linear solve time scales reasonably well with the number of MPI ranks, as does Chapel's label propagation time with the number of locales. On Mutrino, we see a difference in the label propagation time for our two approaches, with per-iteration time significantly longer using Method Two. This time difference is not seen for the 494

495 496

497

498

499

500

501 502

503

504

505

506

507

508

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

537

538

539

540

5/11

542

543

544

545

546

547

548

| | | Chapel | | Grafiki | |
|----------|---------|--------------------|------|--------------------|--------|
| | | LabelPropagation | | Hitting Times | |
| | Number | Time per iteration | | Time per iteration | |
| | of | (seconds) | | (seconds) | |
| Platform | Locales | One | Two | One | Two |
| Kahuna | 1 | 1.24 | 1.22 | 0.0285 | 0.0287 |
| | 2 | NA | 1.91 | NA | 0.0162 |
| | 4 | NA | 1.56 | NA | 0.0087 |
| | 8 | NA | 1.00 | NA | 0.0045 |
| | 16 | NA | 0.65 | NA | 0.0051 |
| Mutrino | 1 | 0.82 | 0.89 | 0.0215 | 0.0230 |
| | 2 | 0.65 | 0.77 | 0.0113 | 0.0120 |
| | 4 | 0.44 | 0.47 | 0.0058 | 0.0062 |
| | 8 | 0.22 | 0.29 | 0.0031 | 0.0033 |
| | 16 | 0.11 | 0.15 | 0.0018 | 0.0018 |

Table 1: Runtime of Chapel connected-component label prop-509agation and Grafiki hitting times on Kahuna and Mutrino510using Method One (Direct calls from Chapel to Grafiki, Sec-511tion 3.1) and Method Two (Separate Chapel and Grafiki exe-512cutables, Section 3.2) with small graph bcsstk29.mtx513

| | | Chapel | | Grafiki | |
|----------|---------|--------------------|-------|--------------------|------|
| | | LabelPropagation | | Hitting Times | |
| | Number | Time per iteration | | Time per iteration | |
| | of | (seconds) | | (seconds) | |
| Platform | Locales | One | Two | One | Two |
| Kahuna | 1 | 8626 | 7476 | 1187 | 1235 |
| | 4 | NA | 9103 | NA | 364 |
| | 16 | NA | 3749 | NA | 188 |
| Mutrino | 4 | OOM | 10333 | OOM | 371 |
| | 16 | 1274 | 3079 | 103 | 105 |
| | 64 | 320 | 824 | 30 | 26 |

Table 2: Runtime of Chapel connected-component label propagation and Grafiki hitting times on Kahuna and Mutrino using Method One (Direct calls from Chapel to Grafiki, Section 3.1) and Method Two (Separate Chapel and Grafiki executables, Section 3.2) with large graph *GAP-Kron.mtx*

Grafiki linear solve; Grafiki linear solve times for both approaches 480 are nearly identical. Thus, the degredation in Chapel's label propa-481 gation must be due to our implementation of shareBlock rather 482 than some inherent mapped-memory access issue on Mutrino. We 483 suspect it arises due to differences between Chapel versions 1.23 484 and 1.24. Our shareBlock Domain was built from the Block Do-485 main in Chapel 1.23; it may not exploit performance enhancements in Chapel version 1.24. More investigation and a transfer of our 487 memory-mapped approach to newer versions of Chapel are needed. 488

489 5 CONCLUSIONS AND FUTURE WORK

We have described new methods for integrating Chapel algorithms
with MPI-based numerical libraries, enabling Chapel users to take
advantage of existing, optimized parallel algorithms while maintaining the simplicity of Chapel programming. Our approaches enabled

integration of Chapel algorithms with MPI-based libraries *without requiring additional copies of the user data*, allowing Chapel users to solve problems that fill the computer's memory. We demonstrated our approaches with shared-memory and distributed-memory versions of Chapel, as well as with Docker/Singularity containers.

While our demonstration used a simple Chapel program for finding connected components of a graph and the MPI-based graph libraries Grafiki and Trilinos for computing hitting times within the largest component, our methodology extends beyond this single demonstration. A matrix abstraction in Trilinos (RowMatrix) allowed Trilinos to perform matrix operations using the Chapelformatted edge lists directly (i.e., without copying and/or reformatting the Chapel data). With our methods, any MPI-based library could be used, provided it either has native data structures identical to Chapel's or has abstraction capabilities similar to Trilinos'. To share data from any Chapel program with MPI-based libraries, Chapel users simply need to substitute mmap-enabled data structures (e.g., shareBlock) for standard Chapel arrays in their programs.

This project raises many opportunities for future work.

Arkouda [11] is often used for large-scale graph analysis because of its elegant NumPy-like interface. Arkouda is built on Chapel, so using our methods to integrate Arkouda and MPI-based libraries would be a natural extension of this work.

An option in Chapel to use a user-defined allocator when creating arrays would simplify extension of this work to other Chapel data objects and versions. To enable Chapel to use mmap memory, we refactored several Chapel functions to use the mmap memory (see Figures 5 and 6). These functions were based on the Block Domain in a specific version of Chapel (version 1.23), and would need to be updated for new versions of Chapel and other Domain types. An option for a user-defined allocator would allow us to provide an allocator based on mmap that could be used in many Chapel objects.

Redistributing graphs' edge lists to increase locality and reduce off-processor data accesses would speed matrix-vector multiplication and other operations in both Chapel and Grafiki. Load balancing tools such as Zoltan [5] or ParMETIS [8] balance computational work while reducing off-processor data dependencies. Even simple sorting of edge lists has been shown to reduce cache misses and speed execution [13]. Our current work relies on the users' data layout; users would need to pre-process the data to increase locality.

This work allows Chapel users to leverage years of effort in MPI-based parallel computing for physics-based and graph-based applications by using libraries like Trilinos and Grafiki. Comparisons between Chapel and MPI in terms of performance and easeof-use are beyond the scope of this work, but are important for evaluating the productivity versus performance trade-off in using Chapel. Also important is the sharing of hardware resources between Chapel and MPI processes in modern computer architectures. Our experiments showed that running separate processes (Method Two) on the same cores of a single node did not slow performance relative to running a single process (Method One). Further experiments using Trilinos' multithreaded implementation would help in understanding interaction between Chapel and multithreaded libraries. Moreover, given the extensive recent efforts to adapt MPI-based libraries to accelerator architectures, experiments using Chapel+MPI+GPU would be of great interest and potential benefit to the graph analysis community.

551 6 ACKNOWLEDGEMENTS

- 552 We thank several people for enabling this work. Jon Berry provided
- motivation and use cases for this project. Samuel Knight, Omar
- 554 Aaziz and Connor Brown provided system support on Sandia's
- 555 computers. Chapel developers Elliot Ronaghan, Brad Chamber-
- ⁵⁵⁶ lain, Michael Ferguson, and Lydia Duncan in the Chapel Discourse
- ⁵⁵⁷ Group provided technical advice and got us over several hurdles.
- Danny Dunlavy and the conference reviewers provided valuable
 feedback on this paper.
- 560 Sandia National Laboratories is a multimission laboratory man-
- aged and operated by National Technology and Engineering So-
- ⁵⁶² lutions of Sandia, LLC, a wholly owned subsidiary of Honeywell
- ⁵⁶³ International Inc., for the U.S. Department of Energy's National
- ⁵⁶⁴ Nuclear Security Administration under contract DE-NA0003525.

565 **REFERENCES**

577 578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

- [1] Bradford Chamberlain, Sung-Eun Choi, Steven Deitz, David Iten, and Vassily
 Litvinov. 2010. User-Defined Distributions and Layouts in Chapel: Philosophy
 and Framework. In 2nd USENIX Workshop on Hot Topics in Parallelism. Cray Inc,
 1–12. https://chapel-lang.org/publications/hotpar10-final.pdf
- [2] Bradford Chamberlain, Steven Deitz, David Iten, and Sung-Eun Choi. 2011. Au thoring User-Defined Domain Maps in Chapel. In *Chapel User's Group 2011*. Cray
 Inc, 1–6. https://chapel-lang.org/publications/cug11-final.pdf
- [3] B. L. Chamberlain, E. Ronaghan, B. Albrecht, L. Duncan, M. Ferguson, B. Harshbarger, D. Iten, D. Keaton, V. Litvinov, and P. Sahabu. 2018. Chapel comes of age: Making scalable programming productive. In *Cray User Group 2018*. https: //cug.org/proceedings/cug2018_proceedings/includes/files/pap130s2-file1.pdf
 - [4] Timothy A. Davis and Yifan Hu. 2011. The University of Florida sparse matrix collection. ACM Trans. Math. Software 38, 1 (2011), 1:1 – 1:25. https://sparse. tamu.edu/
 - [5] K.D. Devine, E.G. Boman, R.T. Heaphy, R.H. Bisseling, and U.V. Catalyurek. 2006. Parallel Hypergraph Partitioning for Scientific Computing. In Proc. of 20th International Parallel and Distributed Processing Symposium (IPDPS'06). IEEE.
 - [6] Michael A Heroux, Roscoe A Bartlett, Vicki E Howle, Robert J Hoekstra, Jonathan J Hu, Tamara G Kolda, Richard B Lehoucq, Kevin R Long, Roger P Pawlowski, Eric T Phipps, et al. 2005. An overview of the Trilinos project. ACM Trans. Math. Software 31, 3 (2005), 397–423. http://dx.doi.org/10.1145/1089014.1089021
 - [7] Hewlett Packard Enterprise Development LP. 2021. Chapel Documentation. https://chapel-lang.org/docs/
 - [8] G. Karypis and V. Kumar. 1997. ParMETIS: Parallel Graph Partitioning and Sparse Matrix Ordering Library. Technical Report 97-060. Dept. Computer Science, University of Minnesota. http://glaros.dtc.umn.edu/gkhome/metis/parmetis/ overview
 - [9] Gregory M. Kurtzer, Vanessa Sochat, and Michael W. Bauer. 2017. Singularity: Scientific containers for mobility of compute. PLOS ONE 12, 5 (05 2017), 1–20.
 - [10] Trevor M. McCrary, Karen D. Devine, and Andrew J. Younge. 2022. Integrating PGAS and MPI-Based Graph Analysis. Technical Report SAND2022-0653R. Sandia National Laboratories. Computer Science Research Institute Summer Proceedings 2021, J.D. Smith and E. Galvan, eds., 235–249.
- [11] Michael Merrill, William Reus, and Timothy Neumann. 2019. Arkouda: Interactive
 Data Exploration Backed by Chapel. In *Proceedings of the ACM SIGPLAN 6th Chapel Implementers and Users Workshop.* 28. https://doi.org/10.1145/3329722.
 3330148
- [12] Message Passing Interface Forum. 1994. MPI: A Message-Passing Interface Standard.
 Technical Report. University of Tennessee, USA.
- [13] Eric T. Phipps and Tamara G. Kolda. 2019. Software for Sparse Tensor Decomposition on Emerging Computing Architectures. SIAM Journal on Scientific
 Computing 41, 3 (2019), C269–C290.
- [14] The Trilinos Project Team. 2021. The Trilinos Project Website. https://trilinos.
 github.io/
- [15] Michael Wolf, Daniel Dunlavy, Richard B Lehoucq, Jonathan W Berry, and Daniel
 Bourgeois. 2018. *TriData: High Performance Linear Algebra-Based Data Analytics*.
 Technical Report. Sandia National Laboratories.