

INSTITUT FÜR INFORMATIK

DER LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN



Master's Thesis

Graph Concepts for the DASH C++ Template Library

Stefan Effenberger

INSTITUT FÜR INFORMATIK

DER LUDWIG-MAXIMILIANS-UNIVERSITÄT MÜNCHEN



Master's Thesis

Graph Concepts for the DASH C++ Template Library

Stefan Effenberger

Aufgabensteller: Prof. Dr. Dieter Kranzlmüller

Betreuer: Tobias Fuchs

Abgabetermin: **ADD DATE**

I hereby declare that I have produced this paper without the prohibited assistance of third parties and without making use of aids other than those specified; notions taken over directly or indirectly from other sources have been identified as such.

Munich, **ADD DATE**

.....
(Signature)

Abstract

ADD ABSTRACT

Contents

1	Introduction	1
1.1	Problem statement	1
1.2	Scope and Objectives	2
2	Background	3
2.1	C++ Concepts	3
2.1.1	Language Concepts	3
2.1.2	Standard Template Library	6
2.2	High Performance Computing	8
2.3	Partitioned Global Address Space	9
2.4	DASH C++ Template Library	10
2.4.1	Dynamic memory allocation	11
3	Related Work	13
3.1	Shared Memory	13
3.2	Distributed Memory	14
3.2.1	Bulk-synchronous graph processing	14
3.2.2	Asynchronous graph processing	15
3.2.3	Linear algebra based graph processing	15
4	Graph Container Concepts	17
4.1	Fundamentals	17
4.2	Index Space	17
4.3	Memory Space	17
4.4	Element iteration	17
4.5	Semantics	17
5	Reference Implementation	19
5.1	Overview	19
5.1.1	Pointers and references	20
5.1.2	Vertices and edges	20
5.1.3	Graph types	20
5.2	Memory management	21
5.2.1	Contiguous memory	21
5.2.2	Edge list memory	22
5.2.3	Element deletion	22
5.3	Iteration	22
5.3.1	Local iteration	23
5.3.2	Global iteration	23
5.3.3	Edge iteration	23
5.4	Data access	24

Contents

6 Case studies	25
6.1 Static structure	25
6.1.1 Graph traversal	25
6.1.2 Shortest path evaluation	25
6.2 Dynamic Structure	25
6.2.1 Graph partitioning	25
6.2.2 De Bruijn Graph construction	25
7 Evaluation	27
7.1 Micro-benchmarks	27
8 Conclusion	29
8.1 Summary	29
8.2 Assessment	29
8.3 Outlook	29
List of Figures	31
Bibliography	33

1 Introduction

Many scientific projects are largely enabled by simulation. Because such simulations often require huge computational capabilities, single compute nodes with a shared-memory architecture cannot provide enough computation power and storage for numerous cases. For this reason, in High Performance Computing (HPC), work is distributed among multiple, interconnected nodes to facilitate the solving of large problems in a timely manner. Since processors cannot directly access the memory of other nodes, the traditional programming model for such systems requires programmers to explicitly distribute data between nodes via message passing. This imposes high demands on the programming skills of scientists who might not have a background in computer science.

Therefore, with the Partitioned Global Address Space (PGAS) model, a new approach is proposed: The memory space of individual nodes in a system is unified within a global address space so that each node can directly access the memory of all other nodes. Programmers are still required to keep data access between nodes to a minimum because data transfer over an interconnect is costly. To further reduce the demands on the programmer, distributed data structures that handle data distribution and load balancing are needed.

Furthermore, data-intensive tasks have been gaining a continually growing interest in the scientific community. Traditionally, applications in HPC follow a computation-centric approach by solving numerical algorithms in the fastest possible way. As “Big Data” is becoming increasingly important in scientific projects, a shift towards more data-oriented applications can be observed in recent HPC projects [ZZZ⁺14]. This trend requires distributed data structures that allow for the storage of large amounts of irregular data and cater to the needs of ever-changing dynamic data.

1.1 Problem statement

Data can be represented in numerous ways. The most generic form of data representation is enabled by *graphs*. A graph $G(V, E)$ is a pair with a set of vertices V and a set of edges E that connect the vertices. This allows for the representation of data and its relationships in regular as well as irregular patterns.

On distributed machines, graph data structures can be implemented using a variety of different characteristics. This has led to many different implementations - usually a new implementation for each algorithm - which are hardly compatible with each other. To overcome this situation, generic programming abstractions to facilitate reuse of existing code and to lower the demands on programmers are needed.

As of today, no generic graph abstractions implementing the PGAS model exist. This work therefore aims to provide a graph abstraction for C++ containers that allows for the implementation of arbitrary graph algorithms following the PGAS model on distributed memory machines.

1.2 Scope and Objectives

In this work, a C++ concept for graph containers following the PGAS model is presented. The concept is part of the DASH C++ Template Library and thus uses concepts already present in the library.

The graph concept is meant to provide a generic framework for the programming of arbitrary graph algorithms in the context of distributed machines and especially the Partitioned Global Address Space model. This means that it meets the following requirements:

- Native support for one-sided communication
- Support for the programming of synchronous graph algorithms
- Support for the programming of asynchronous graph algorithms
- Portability across platforms (**PORTABLE EFFICIENCY?**)
- Support for heterogenous systems
- **FINISH REQUIREMENTS**

Furthermore, this work provides concepts for the dynamic allocation of graph data across multiple machines with a focus on optimized data locality. **LOAD BALANCING?**

A reference implementation is then used to verify the usability, correctness and universality of the given concepts. While the concepts are designed for high performance, the reference implementation is not. This means that the evaluation of this implementation does not account for the performance of this work's concepts.

2 Background

This chapter covers some fundamental background knowledge needed for a better understanding of the following chapters of this thesis. Only explanations directly relevant to the topics of this thesis are provided.

Since the result of this work is a C++ concept, some important language expressions and concepts are firstly discussed, along with a description of the Standard Template Library on which concepts this work is built upon. The reader is then introduced to the domain of High Performance Computing which is the main application area for this work. A brief overview of the Partitioned Global Address Space programming model is then followed by a description of the DASH Library which provides core concepts used in this thesis.

2.1 C++ Concepts

The graph concepts of this work are part of the DASH C++ Template Library (see section 2.4). For this reason, the reference implementation is written in C++11 [ISO12]. This section illustrates some basic knowledge about important C++ concepts used in the implementation description of chapter 5.

2.1.1 Language Concepts

This section describes features and concepts of the C++ language the reference implementation of this thesis is based on. While C++ offers many different features for different programming styles, only a certain subset is used due to performance reasons.

Value and reference semantics

In object oriented programming, objects can either have value or reference semantics. Objects with value semantics are treated like values: An assignment operation copies the object. This way, the identity of the object is not in focus because the copied object has another identity. Operations on this copy do not affect the original object which means only the value of the object is in focus. Objects with reference semantics are referred to by references or pointers. Their identity becomes important: Multiple references point to the same object.

Many object-oriented programming languages such as Java only offer reference semantics for non-primitive types. C++ on the other hand allows the programmer to define, whether an object adheres to value semantics or reference semantics.

For an object to be able participating in value semantics, some operations like copy construction and assignment have to be implemented in a certain way. C++ compilers provide default implementations of the required operations, but depending on the object's implementation, further measures might have to be taken by programmers. For example, objects that allocate memory dynamically on the freestore have to explicitly copy the referenced

2 Background

data.

Listing 2.1 shows an `increment` function with both value and reference semantics. `increment_vs` takes the copy of a `VertexIndex` object as parameter, increments its offset and returns a copy of the object with incremented offset. The offset of the passed `VertexIndex` stays the same as only the offset of its copy has been changed. In contrast, `increment_rs` changes the offset of the passed object and also returns a reference to the same object.

Listing 2.1: Value and reference semantics

```
1 VertexIndex increment_vs(VertexIndex v) { // value semantics
2     return ++(v.offset);
3 }
4
5 VertexIndex & increment_rs(VertexIndex & v) { // reference semantics
6     return ++(v.offset);
7 }
```

Value semantics seem to result in a lot of copying that might hit performance. C++ compilers implement an optimization technique called *copy elision* that omits copy construction in functions returning objects with value semantics by returning the same memory location of the temporarily created object. *Copy elision* is part of the C++ standard and thus compilers are required to enforce it. For this reason, C++ value semantics in many cases have no performance drawbacks in comparison to reference semantics. Therefore, reference semantics is mainly used when identity of an object is important or when copying of an object is expensive.

Reference semantics is also important for runtime polymorphism, because objects with value semantics might be sliced (i.e. only the part of the base class is copied).

Operator Overloading

In C++, almost all existing operators can be overloaded for any operand types. Any class can therefore be handled with native operators in a completely customized way. Only four operators like the member access operator cannot be overloaded and it is not possible to create new operators that do not exist in the language itself.

Listing 2.2: Operator overloading

```
1 class Iterator {
2     public:
3         Iterator & operator++() {
4             ++position;
5             return *this;
6         }
7     private:
8         int position = 0;
9     };
10
11 Iterator it;
```

```
12 ++it;
```

Listing 2.2 shows a class `Iterator` with an overloaded version of the pre-increment operator. The operator can be directly applied to an object of `Iterator` incrementing the `position` member.

Static vs. runtime polymorphism

C++ offers two kinds of polymorphism: static and runtime. The difference lies in the way types are bound. Static polymorphism can be completely resolved at compile time, while types in runtime polymorphism have to be resolved during the runtime of a program.

In runtime polymorphism, methods of a derived class are called with a pointer of the base class by using *virtual functions* of the base class. A call to such a virtual function requires resolving which concrete derived class the pointer of the base class refers to during runtime. This is achieved by storing a *vtable* for each base class in memory and linking to this *vtable* with a pointer from all related objects. The runtime can then lookup the address of the desired class's method. Listing 2.3 shows how a method implemented in a derived class is executed calling the same method on a pointer of the base class.

Listing 2.3: Runtime polymorphism

```
1 class Base {
2     virtual void do_something() {
3         // do something
4     }
5 };
6
7 class Derived : Base {
8     virtual void do_something() {
9         // do something else
10    }
11 };
12
13 Base * b;
14 Derived d;
15 b = &d;
16 b.doSomething(); // Derived::do_something() is called
```

Because the *vtable* lookup on every call is expensive, C++ allows programmers to design polymorphic types with static polymorphism that can be completely resolved at compile time leading to better performance during runtime. This can be achieved with simple method overloading and with templates.

In C++11, templates can be defined for classes and functions. They allow the programmer to define a family of either. Listing 2.4 shows a class `Base` that accepts a template parameter which can be of any type. A call to `do_something` is delegated to an object of the type of the template parameter. Type resolving is done during compilation of the program so that a compile error would occur if no `do_something` method were available in `TypeB`.

Listing 2.4: Static polymorphism with class templates

```

1 struct TypeA {
2     void do_something() {
3         // do something
4     }
5 }
6
7 struct TypeB {
8     void do_something() {
9         // do something else
10    }
11 }
12
13 template<typename Type>
14 class Base {
15     public:
16     void do_something() {
17         type.do_something();
18     }
19     private:
20     Type type;
21 };
22
23 Base<TypeB> b;
24 b.do_something();

```

Templates actually result in code generation: For every instantiation of a class template, a new type is created which results in a larger binary file.

2.1.2 Standard Template Library

The Standard Template Library (STL) [SL95] is a C++ software library. Most parts of it have been integrated into the C++ Standard and are now part of the C++ Standard Library. The STL has been designed as a library for efficient generic programming with value semantics. Its containers are template classes that enable static polymorphism.

The STL contains a variety of containers and algorithms. Algorithms are decoupled from containers with the help of iterators so that any algorithm can work with any container that is STL-compatible. Because DASH - the library this work is part of (see section 2.4) - follows the concepts of the STL, STL algorithms can be executed with any DASH container, including the graph container of this work's reference implementation.

This section focuses on STL concepts in the C++11 standard.

Concepts

In C++11, a *concept* is a named set of requirements for a type. The C++ standard contains concepts for all components of the C++ standard library and this work provides a concept similar to these. At the time of this writing a C++ extension that allows for the formal specification and compile-time evaluation of concepts directly in code is in development [ISO15]. Because this extension is still work in progress and not part of C++11, it is not used

in this thesis. Thus, the graph concepts of this work require the programmer to manually ensure that all requirements of the respective concept are met in its implementation.

A C++ concept typically contains a description stating other concepts that are used or implemented by the concept. For example, the **SequenceContainer** concept of the C++ standard implements the **Container** concept and additionally accounts for linear arrangement of the contained elements. *Types* and *Methods/Operations* are defined along with their semantics. For some expressions, *computational complexity* is additionally defined and programmers have to keep their implementations inside of these constraints.

The graph concepts of this work try to assimilate the concepts of the standard library as far as possible.

Iterators

Iterators are the connection between *containers* and *algorithms*. Every container offers iterators with a standard interface to allow for the iteration of contained elements. Algorithms only have to comply to iterator interfaces and are agnostic of the actual interface of the container. Thus, the same algorithm can run with a variety of different containers without it being re-implemented for every existing container.

Listing 2.5 shows the computation of a sum of vector *v* using the **accumulate** algorithm of the C++ standard library. *v.begin* returns an iterator to the beginning of *v* and *v.end* an iterator past the last element of *v*. The **accumulate** algorithm then iterates over the elements of *v* without knowing the details of its underlying container.

Listing 2.5: Vector sum using standard algorithm

```

1 std::vector<int> v = { 1 , 2, 3 };
2 int sum = std::accumulate(v.begin(), v.end(), 0);

```

Because containers have different memory layouts and algorithms have different requirements on iterators, the C++11 standard defines four iterator concepts that are hierarchically organized, i.e. iterators of a higher category implement all operations of the iterators in lower categories. Table 2.1 shows these concepts starting from the lowest category:

Table 2.1: STL iterator categories

Category	operations
InputIterator	increment, read
ForwardIterator	multi-pass increment
BidirectionalIterator	decrement
RandomAccessIterator	random access

InputIterators can be incremented to iterate over elements one by one. There is no guarantee, that a value dereferenced from an **InputIterator** is still valid after it has been incremented again. This guarantee is only given by **ForwardIterators** which makes them a requirement for algorithms that have to iterate over elements more than once. **BidirectionalIterators** can be decremented and thus are able to iterate in opposing

2 Background

direction. `RandomAccessIterators` can access any element in a given range in constant time.

Any iterator can also implement the `OutputIterator` concept that enables writing operations on elements.

Containers

Standard C++ containers implement common data structures like arrays and stacks. The container library consists of sequence containers that store elements in a specific order and associative containers that allow searching for particular elements. Additionally, adaptors like `std::queue` are used to restrict the interfaces of existing containers.

All concepts of concrete implementations meet the requirements of the standard library's `Container` concept that defines element access via iterators. The `Container` concept defines `ForwardIterator` as iterator type. Containers can additionally meet the requirements of further concepts which might override certain requirements like the iterator type. The `ReversibleContainer` concept for example allows containers to either use `BidirectionalIterators` or `RandomAccessIterators`. A concrete container concept like `std::vector` that meets the requirements of `ReversibleContainer` can then define the iterator type as needed (`RandomAccessIterator` in this case).

2.2 High Performance Computing

High Performance Computing (HPC) is a broad term describing advances for the fastest possible computation of a given problem. Gustafson's Law [Gus88] suggests that a compute system can linearly grow with the problem size: A problem of two times its original size can be computed on a system with twice as many processors in the same time (best case scenario). This means that very large problems can be computed in an acceptable timeframe if there is a sufficiently large compute system available. Depending on the problem size, two different system architectures are used in HPC:

Shared Memory A shared memory system consists of a single node with multiple processors connected to the same random access memory. Memory access for the different processors can be uniform, but many systems implement a non-uniform memory access (NUMA) design where a part of the memory is assigned to each of the processors. A processor in a NUMA system can access its assigned memory faster than the memory of the other processors. Because processors can access all data at all times, communication between processors has a low cost which simplifies programming on these systems in comparison to distributed memory systems. Achieving high performance on NUMA systems is more problematic because the programmer has to take data locality into account [Lam13].

Distributed Memory Multi-processor systems in which each processor has access to its own memory space are called distributed memory systems. These systems usually consist of several shared memory nodes with the processors of one node not being able to directly access memory of other nodes. While single shared memory systems can only be scaled to a certain extent, the scalability of distributed systems is much higher [PTM96].

The nodes are connected with a network interconnect for communication between the processors. Due to the latency of the interconnect being significantly higher than the latency of a memory bus in a shared memory system, communication is much more costly. This imposes higher demands on the programmers' skills in comparison to shared memory systems.

The largest problems in science are computed on “supercomputers” like the *SuperMUC* at the *Leibniz Rechenzentrum* in Munich. These distributed memory machines consist of hundreds or even thousands of homogeneous nodes that are connected with a specialized interconnect. To this date, *message passing* is the prevalent programming model for such systems.

2.3 Partitioned Global Address Space

Shared Memory and *Message Passing* are the dominant models in HPC as of this writing. As pointed out in section 2.2 however, the usage of Message Passing requires high skills in computer architecture and programming. To ease this problem, the Partitioned Global Address Space (PGAS) model has been proposed. It unifies some of the benefits of both of these models by creating a global address space over the initially local-only address spaces of distributed machines.

Figure 2.1 a) presents the architecture of a shared-memory machine: Multiple processors share a common address space. The processors are attached to the same memory over a bus. In some systems, memory might be local to some processors which means the rest of the processors has a higher latency when trying to access the non-local memory. Still, every processor can access every part of the address space. Communication takes place *implicitly* by writing and reading shared variables. Because data written by one processor can be accessed by another processor in a fast manner, little care has to be taken regarding the decomposition of data. For this reason however, shared memory programs are typically not scalable on distributed machines [SAB⁺10].

Figure 2.1 b) shows that a distributed memory machine basically consists of several shared memory machines linked to each other via an interconnect. Since processors cannot directly access data stored in the memory of other machines, *explicit communication* is needed in order to synchronize the processors. This is typically done by two-sided communication: The *sending* of a message has to be accepted at the remote machine with a corresponding *receive* call.

Machines conduct their computations simultaneously and either synchronize in discrete time intervals or exchange data asynchronously. Either way, sending data over an interconnect imposes high latency and low throughput in comparison to the data access over a memory bus in shared memory systems. For this reason, programmers have to carefully decompose data in order to distribute the work load uniformly and minimize communication overhead.

Figure 2.1 c) illustrates the concept of Partitioned Global Address Space: The local portions of memory are unified under a global address space which allows processors to directly access data on remote machines. Data access is performed using one-sided communication:

2 Background

No *receive* call on the remote machine is needed.

Since data transferal over an interconnect is still costly, programmers have to take the same care for data locality as with the traditional message passing approach. To allow for this, the locality of a datum is directly exposed to the programmer.

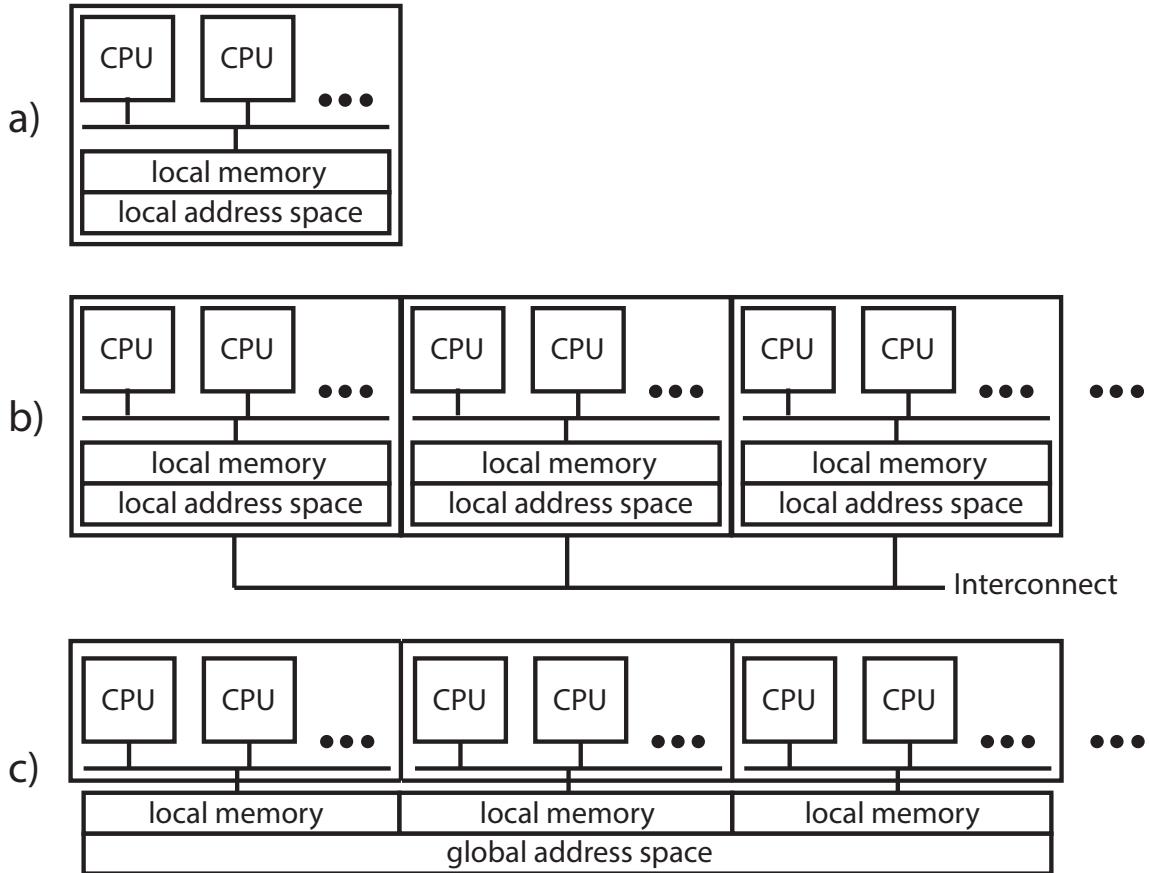


Figure 2.1: View on a) Shared Memory b) Distributed Memory c) Partitioned Global Address Space

Existing PGAS approaches are mainly comprised of dedicated programming languages such as Unified Parallel C (UPC) [C⁺05], Co-Array Fortran [NR98] or Chapel [CCZ07] that allow for compiler optimizations in respect to distributed machines but lack portability and reach. In contrast to this, efforts exist to create libraries for existing programming languages used by many HPC systems.

2.4 DASH C++ Template Library

DASH [FFK16] is a compiler-free PGAS approach: It consists of a simple C++ library that can be compiled with any C++ compiler and thus can be used out-of-the-box on most HPC systems. The library is part of the Priority Programme “Software for Exascale Computing” (SPPEXA)¹ which supports research on computing systems achieving 10^{18} floating point

¹<http://www.sppexa.de>

operations per second and above. While PGAS languages require existing programs to be completely rewritten from scratch, DASH allows the applications to be incrementally ported and thus facilitates wider adoption of the PGAS model in the HPC community.

DASH operates on top of the *DASH Runtime* (DART) which is a PGAS memory allocation and communication abstraction written in C. DART enables global memory allocation, pointers to remote memory locations and one-sided communication on top of existing libraries like MPI [For12] or GASPI [GS13]. With DART-MPI [ZMI⁺14], a fully functional DART abstraction on top of MPI-3 is used in DASH releases at the time of this writing.

In DASH, processing elements are referred to as *units*. Units can be any processing element such as threads or processes. DASH programs are implemented using the Single Program Multiple Data (SPMD) model: The data is partitioned onto the participating units and each unit executes the same code on its part of the data. Furthermore, units form *teams* that can be created at runtime. Because HPC hardware topologies become more complex over time (e.g. [KDSA08]), DASH supports hierarchical team creation to allow for a more fine-grained exploitation of data locality compared to the typical local-remote distinction of the PGAS model.

Data is referred to in terms of global pointers and references. A `GlobPtr<T>` object holds information about the unit and local memory location of the referenced datum. It can be dereferenced to a `GlobRef<T>` object which behaves like a C++ reference and can be converted to an object of type T. This type conversion triggers a one-sided get operation transferring the data from its remote source to the caller. Similarly, data can be written into the referenced memory location of a `GlobRef<T>` object.

DASH provides a set of containers for distributed data storage. Aside from the static data structures Array and Matrix, dynamic data structures are available. Since the graph concepts of this work belong into the latter category, details of it are discussed in the following.

2.4.1 Dynamic memory allocation

Dynamic allocation in DASH is encapsulated in the `GlobHeapMem` concept. `GlobHeapMem` offers two basic operations to dynamically allocate memory during runtime: `grow` and `shrink`. These operations increase or decrease the local size of the memory allocated on the respective unit. Changes in memory space are not reflected in global address space until the operation `commit` is called which publishes the changes across all units.

A dynamic container in DASH pre-allocates some memory during its initialization. When the memory is completely used, further additions of elements result in `GlobHeapMem.grow` operations. A call to the `barrier` operation of the container results in all newly added elements of the container to be publicly available on all units.

Multiple `grow` operations result in a scattered memory space: Each call to `grow` creates a new *bucket* - a contiguous memory region in the local freestore. A class implementing the `GlobHeapMem` concept keeps track of each bucket on each unit so that element indices can be translated to concrete memory locations. Figure 2.2 illustrates the memory spaces of two units after two `grow` operations for six elements each have been called. The buckets are allocated at different memory locations. Data access to an element therefore requires the memory location of the bucket and the offset of the element inside the bucket. For this reason, bucket locations and sizes are exchanged between all units during a `commit` operation

2 Background

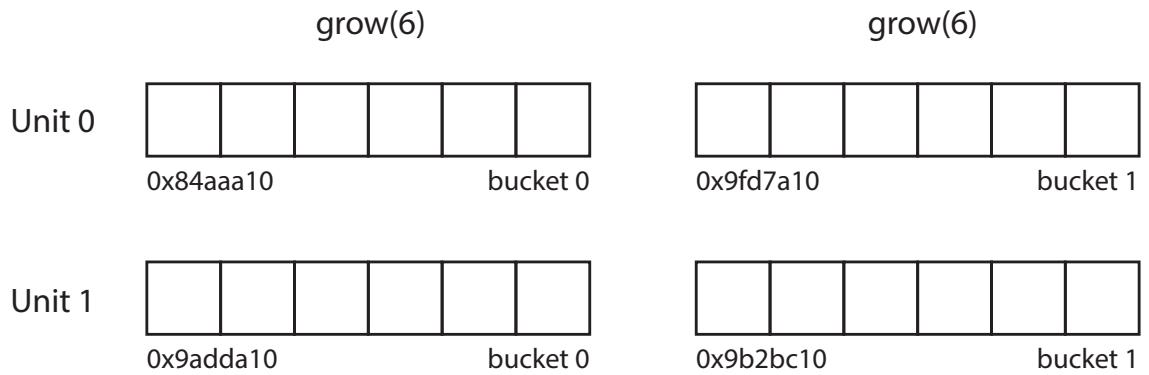


Figure 2.2: Memory space of two units after two `GlobHeapMem.grow` operations

and a reference to the object holding the bucket data is needed in every pointer and iterator.

3 Related Work

3.1 Shared Memory

Graph processing on distributed machines tends to result in low processor utilization due to bad locality [LGHB07]. Shared memory systems on the other hand allow for simpler graph algorithms with better processor utilization [SB13]. Because some shared memory systems with several terabytes of memory exist (e.g. the LRZ teramem¹), graphs up to a certain size can also be computed on such systems. While many libraries and frameworks for graph processing on shared memory systems exist, this section only focuses on work relevant to the HPC community.

Spatio-Temporal Interaction Networks and Graphs Extensible Representation (STINGER) [BBAB⁺09] is a graph data structure for shared memory systems. STINGER tries to provide a standardized and extensible data structure that allows portability across different graph processing frameworks. The data structure consists of a set of standards that have to be abided by actual implementations.

At its core, STINGER is an adjacency list [LINK TO EXPLANATION?](#) with a fixed structure. Vertices are stored in a contiguous memory region and point to blocks of edges. The set of all edges of a single vertex can be distributed between multiple blocks with each block containing a pointer to the next block. Vertices and edges both store a weight value and edges additionally store timestamps. This allows edges to be automatically removed after they have reached a certain age. STINGER also accounts for fault tolerance by defining standards for checkpointing functionality.

The **MultiThreaded Graph Library (MTGL)** [BHKK07] focuses on graph processing using specialized multithreading systems that contain processors supporting a large amount of thread contexts. The MTGL is specifically designed for the Cray MTA architecture which limits its use significantly because a Cray MTA system has to be available. Brian W. Barrett et al. have implemented an MTGL version for commodity hardware [BBMW09] but its performance is significantly lower compared to the original version for serial machines.

The Cray MTA architecture [CFS99] utilizes processors with relatively slow speed and without data caches. Memory access is always blocking but each processor features a large amount of threads so that there are always free threads available for computation in most cases. This results in a simple and fine-grained synchronization model that is not available on commodity hardware, but also creates subtle concurrency and performance issues. These issues are addressed by the MTGL in order to allow straight-forward programming of graph algorithms on these platforms.

Because of the architecture of the Cray MTA, locality is not an issue in the MTGL: Data is partitioned by the runtime system of the MTA. The interface of the MTGL is - like many

¹https://www.lrz.de/services/compute/special_systems/teramem/

3 Related Work

graph processing frameworks - based on the Boost Graph Library (BGL) [SLL01]. It is less generic though: Only *directed graphs* are supported and data is always stored in an *adjacency list* structure.

3.2 Distributed Memory

Processing of large graphs that do not fit into the memory of a single machine takes place on distributed memory machines in HPC. Shared memory approaches that hold parts of the graph in persistent memory (e.g. GraphChi [KBG12]) do not have sufficient performance in order to be useful in the field of high performance graph processing. However, because some graphs like scale-free social networks lack locality, graph partitioning on distributed memory machines is a challenge that can result in unpredictable performance due to low processor utilization [BHKK07]. While performance and scalability are key capabilities of the work mentioned in this section, they can not always be guaranteed.

3.2.1 Bulk-synchronous graph processing

In *Parallel Computing*, Valiant's *Bulk Synchronous Parallel (BSP)* model [Val90] describes a way to exploit locality: Data is partitioned across processors in a way that allows each processor to perform independent computations. BSP algorithms run a series of *supersteps*, each consisting of three phases:

1. Local computation on each processor.
2. Global exchange of computed data between processors.
3. Synchronization of all processors.

The performance of BSP algorithms relies heavily on a well-suited decomposition of the data. Because of the explicit barrier synchronization at the end of each superstep, processor time may be wasted on processors that finish their work faster than others. Deadlocks and data races are however prevented by it.

The *Boost Graph Library (BGL)* [SLL01] is a widely used sequential graph processing framework. With the **Parallel Boost Graph Library (PBGL)** [GL05], the BGL has been ported to the environment of distributed machines using the BSP model. It aims at providing a general-purpose library for graph processing on distributed machines. The PBGL data structures adjacency list and adjacency matrix therefore are designed as generic as possible and the library includes a wide variety of different graph algorithms that can be extended with the *Visitor Pattern* [PJ98]. Properties of graph components are not restricted to simple edge weights: The PBGL allows programmers to define arbitrary property data structures for vertices and edges. Supported graph types in the PBGL include directed graphs, undirected graphs and bidirectional graphs.

Pregel [MAB⁺10] is a framework for large-scale graph processing developed by Google. Like the PBGL, it partitions graphs by vertex and uses the BSP model for computation. Vertices are statically partitioned by a user-defined hash function. Pregel additionally accounts for fault tolerance with checkpointing mechanisms but only supports directed graphs. Programming in Pregel is restrictive because it only allows the definition of vertex visitors that have access to messages sent by other vertices in the previous superstep.

3.2.2 Asynchronous graph processing

3.2.3 Linear algebra based graph processing

Many graph algorithms can be expressed by a set of linear algebra operations [FR11]. The graph is represented by an adjacency matrix storing edge weights and additional vectors are used as auxiliary data structures.

Because an adjacency matrix can be partitioned into blocks with each block being placed into the memory of a single processor, the graph can be partitioned in a 2-dimensional way. As research has shown, breadth-first-search and related algorithms can be implemented in a way that exploits the 2-D partitioning for higher performance in comparison to a 1-dimensional partitioning (i.e. partitioning based on vertices instead of edges) [BM11].

The **Combinatorial BLAS** [BG11] exploits these findings with a library containing *Basic Linear Algebra Subroutines* (BLAS) specifically targeted at graph processing. It provides a subset of typical BLAS operations like matrix-matrix multiplication that are sufficient for expressing most graph algorithms. Graph data is stored in a sparse matrix data structure. Vertices are therefore not mutable which results in graphs being static. Also, algorithms like Dijkstra's shortest path, that are based on priority queues, tend to have lower performance.

The Combinatorial BLAS is also used as a computational engine for a high-level graph processing framework: the **Knowledge Discovery Toolbox (KDT)** [LAB⁺12]. The KDT aims at easing use of graph algorithms for domain experts that are not experts in computer science at the same time. To achieve this, it provides an API that allows to run algorithms on graphs with only few lines of code. Algorithm developers can extend the library by creating algorithms with the underlying linear algebra primitives. While KDT can be easily used by data scientists without background in computer sciences, the performance and scalability are still promising: In comparison to the PBGL, speedups from 3 to 12 can be observed depending on the problem size and core count. Because it uses the Combinatorial BLAS as underlying engine, KDT also inherits its restrictions explained above.

4 Graph Container Concepts

4.1 Fundamentals

4.2 Index Space

4.3 Memory Space

4.4 Element iteration

4.5 Semantics

5 Reference Implementation

This section explains a concrete reference implementation of the concepts in chapter 4. The implementation is part of the DASH C++ Template Library and thus written in C++11. It is based on basic C++ concepts illustrated in section 2.1. The reference implementation will be referred to as `dash::Graph` in this chapter.

5.1 Overview

`dash::Graph` is part of the dynamic data containers of the DASH Library. As such, it is interacting with existing components of the library. Figure 5.1 depicts the components and their main interactions.

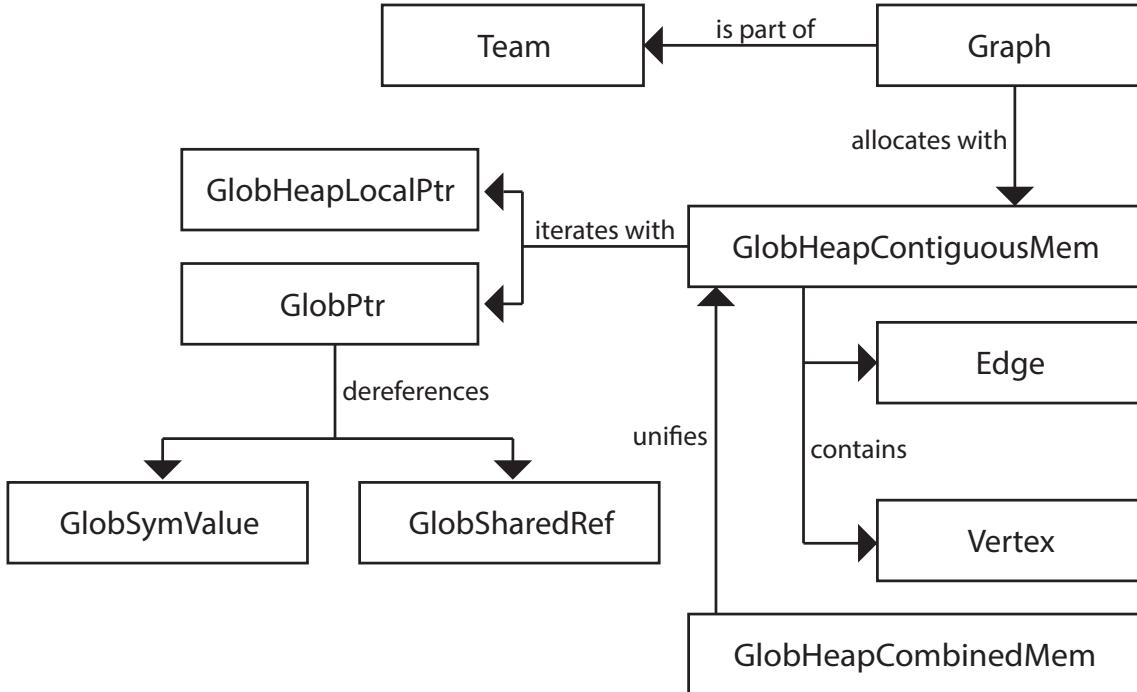


Figure 5.1: `dash::Graph` component overview

The `Graph` gets initialized with a reference to an existing `Team`. By default, this is `Team::All()` which includes all units DASH has been initialized with. The `Graph` creates three instances of `GlobHeapContiguousMem` for vertices, inbound edges and outbound edges. These instances are used to globally allocate memory for the respective elements. Since the `Graph` also allows for the iteration of all (inbound and outbound) edges, `GlobHeapCombinedMem`

unifies the memory spaces of the two `GlobHeapContiguousMem` instances. Both `GlobHeapContiguousMem` and `GlobHeapCombinedMem` use a specialized template version of `GlobPtr` to iterate over the memory space. Each `GlobPtr` object can then be dereferenced to a `GlobSharedRef` object which enables direct access to the referenced element.

5.1.1 Pointers and references

To allow `GlobPtr` and `GlobSharedRef` act like real pointers and references respectively, operators like the increment and dereference operators are overloaded. This results in usage analogous to native pointers and references:

Listing 5.1: Operator overloading in `GlobPtr` and `GlobSharedRef`

```

1 typedef GlobHeapContiguousMem<std::vector<int>> g_mem_type;
2 GlobPtr<int, g_mem_type> ptr(mem, 0); // position 0 in index space
3 ++ptr; // go to position 1 in index space
4 auto ref = *ptr; // dereference ptr to GlobSharedRef object
5 int val = ref; // convert reference to value

```

In this case, `mem` is an instance of `GlobHeapContiguousMem` holding at least two globally available elements.

5.1.2 Vertices and edges

Vertices and edges are modeled as individual classes: `Vertex` and `Edge`. Each instance of these classes contains objects of the *property* classes defined by the user as template parameters of `dash::Graph`. The property classes can only contain variables of static size. These variables are filled with default values upon initialization as per *default initialization* of the C++ standard.

Additionally, vertices and edges contain their respective index in global index space. The index is not publicly available and can only be read by the `Graph` class and some of its internal classes. Storing the indices internally is required because of global iteration and remote memory access: After a unit has globally iterated over the elements and dereferenced an element, its index is lost and the user has no way to recover it. The index, however, is needed for various use cases such as iterating over the outbound edges of a certain vertex.

Edges also contain indices of their source and target vertices.

5.1.3 Graph types

`dash::Graph` currently supports *directed* and *undirected* graph types. A bidirectional graph type is not implemented but can be integrated later on. As described in [LINK CORRECT SECTION](#), directed graphs can be instantiated in two different variants depending on the need for inbound edge iteration. This is due to the fact that this iteration is not necessary for most algorithms but requires additional communication that lowers the overall performance of the graph, even if not used at all.

For undirected graphs, edges are replicated to the edge lists of both participating vertices. For directed graphs, the same mechanism is used when inbound edge iteration is needed. Edges in directed graphs without inbound edge iteration are not replicated in any way.

5.2 Memory management

`dash::Graph`'s memory space is handled by `GlobHeapContiguousMem`. This class follows the concept described in section 4.3 but adds an additional feature: Fully contiguous global memory regions. While the memory concept only demands single edge lists to be contiguous, `GlobHeapContiguousMem` allocates a contiguous memory space for all vertices as well as all edges for better locality exploitation. However, this comes at the expense of additional memory re-allocations in each epoch.

5.2.1 Contiguous memory

Figure 5.2 illustrates the basic scheme of the contiguous memory allocation. *Region 1* is a publicly available contiguous memory region. The memory location of *Region 1* is known by other units and thus cannot be changed outside of a commit operation. Because of this, *Region 2* is allocated at another memory location that might not be contiguous to *Region 1*. *Region 2* contains elements that have been added in the current epoch - they are available only locally and cannot be seen by other units.

The `commit` operation starts a new epoch by packing the elements of the two memory regions into another, contiguous memory region and notifying other units about the changed location and size of the region. Traversing the elements in this region is now as simple as incrementing a pointer. Local iteration over the elements however requires a hop between the publicly available region and the local-only region because they are not contiguous. `GlobHeapLocalPtr` can iterate over an arbitrary number of buckets with contiguous memory (see section 5.3). In this case, always two buckets (the two mentioned memory regions) are used.

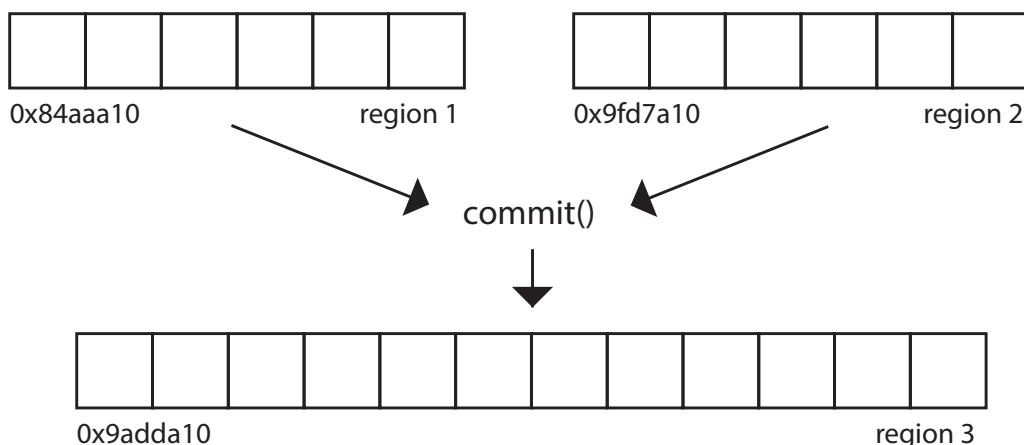


Figure 5.2: Contiguous memory allocation

5.2.2 Edge list memory

The mechanism in subsection 5.2.1 is used for all vertices created at a certain unit. The same mechanism cannot be used for the entirety of all edges created at a single unit, because edge lists have to be contiguous. Appending an element to an edge list would invalidate the offsets of all edges following behind and would also require a re-allocation invalidating the memory location. Figure 5.3 depicts the problem: The new edge inserted at the end of *edge list 1* would be placed at offset 6 in the memory region, incrementing the offsets of all following elements.

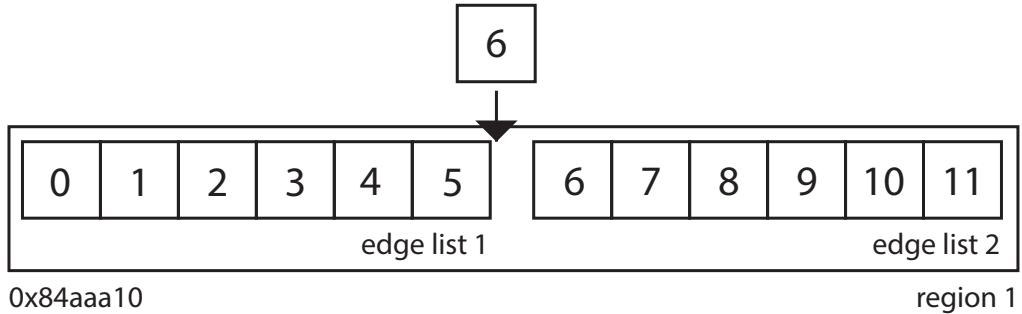


Figure 5.3: Element insertion into edge list invalidating offsets

For this reason, each edge list has to be maintained with individual second regions for element insertions: The mechanism in subsection 5.2.1 is used for each edge list individually and the `commit` operation packs all edge lists into one contiguous memory region.

Vertices contain the indices of their respective edge lists and edge list locations and sizes are maintained in `GlobHeapContiguousMem` so that edge lists of a given vertex can be easily iterated.

5.2.3 Element deletion

Deleting an element in contiguous memory regions requires its memory location to get invalidated instead of removed to avoid shifting of elements and offset invalidation. If many delete operations occur, the memory space will get scattered. For this reason, it is necessary to take measures that reduce the scattering of the memory to a minimum. `GlobHeapContiguousMem` uses a *free list*: Deleting an element results in its memory location being added to the free list. If another element is added, a memory location from the back of the free list is used to store the element. Only if the free list is empty, new memory is allocated.

Because invalidated elements are not part of the memory space anymore, iterators have access to the free list so that they can skip the respective elements.

5.3 Iteration

The four different iteration spaces of the graph concept (see section 4.4) are handled by the same iterator classes. Since the memory space of `dash::Graph` is epoch-based, iteration

space of local iterators can be different to the local part of the iteration space of global iterators.

5.3.1 Local iteration

Local iteration in dynamic containers of DASH is handled by `GlobHeapLocalPtr` that can iterate over multiple non-contiguous memory buckets. `GlobHeapContiguousMem` holds a list of objects containing bucket meta-data including size and a local native pointer to its beginning memory location. The bucket list is passed to a `GlobHeapLocalPtr` along with the position the pointer currently holds in the index space of the buckets. The buckets are equal to the allocated memory regions of `GlobHeapContiguousMem` as described in section 5.2.

Iteration is done by `increment/decrement` operations which result in `GlobHeapLocalPtr` calculating the bucket number the current position belongs to. A `dereference` operation then accesses the local pointer of the bucket object.

5.3.2 Global iteration

For global iteration, a template specialization of `GlobPtr` for `GlobHeapContiguousMem` is used. It is similar to the template specialization of `GlobPtr` for `GlobHeapMem`, which is used for other dynamic DASH containers like `dash::List`.

While the bucket meta-data used by `GlobHeapLocalPtr` contains only information about local buckets, `GlobPtr` requires information about all buckets on all units. Therefore, `GlobHeapContiguousMem` distributes bucket information between all units during a `commit` operation. This information includes bucket sizes as well as DART global pointers to the memory locations of the buckets in global address space. **EXPLAIN DART POINTERS?**

With the bucket size information and the DART abstraction of global pointers, `GlobPtr` can iterate over global memory space the same way as `GlobHeapLocalPtr` iterates over local memory space. The iteration order is given by the canonical order of units.

5.3.3 Edge iteration

Inbound and outbound edges are handled by separate `GlobHeapContiguousMem` objects. Iteration of either is therefore handled as described in subsection 5.3.1 and subsection 5.3.2. The graph concept however requires iteration over all existing edges. For this reason, `GlobHeapCombinedMem` is used to unify the iteration spaces of multiple `GlobHeapContiguousMem` objects.

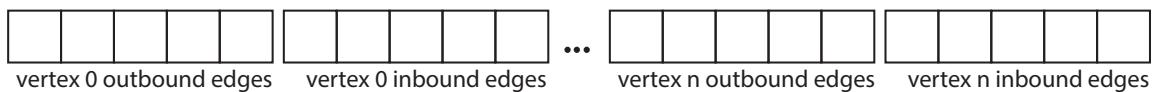


Figure 5.4: Combined edge iteration space

Figure 5.4 illustrates the order in which the combination of inbound and outbound edges are iterated over.

5.4 Data access

Remote memory data access is strictly handled via `GlobSharedRef` objects. These objects issue DART `get`/`put` operations to access data on remote machines. If referenced data resides locally, `GlobSharedRef` directly accesses it in the memory.

Because edges are replicated for some graph types (see subsection 5.1.3), writing data to an element with a `GlobSharedRef` object would introduce consistency problems because it is restricted to referencing single memory locations. `GlobSymValue` **ADAPT NAME ONCE DEFINED** is set up with two DART global pointers to account for this problem. An assignment to a `GlobSymValue` object results in two memory locations being updated.

6 Case studies

6.1 Static structure

6.1.1 Graph traversal

6.1.2 Shortest path evaluation

6.2 Dynamic Structure

6.2.1 Graph partitioning

6.2.2 De Bruijn Graph construction

7 Evaluation

7.1 Micro-benchmarks

8 Conclusion

8.1 Summary

8.2 Assessment

8.3 Outlook

List of Figures

2.1	View on a) Shared Memory b) Distributed Memory c) Partitioned Global Address Space	10
2.2	Memory space of two units after two <code>GlobHeapMem.grow</code> operations	12
5.1	dash::Graph component overview	19
5.2	Contiguous memory allocation	21
5.3	Element insertion into edge list invalidating offsets	22
5.4	Combined edge iteration space	23

Bibliography

- [BBAB⁺09] BADER, David A. ; BERRY, Jonathan ; AMOS-BINKS, Adam ; CHAVARRÍA-MIRANDA, Daniel ; HASTINGS, Charles ; MADDURI, Kamesh ; Poulos, Steven C.: Stinger: Spatio-temporal interaction networks and graphs (sting) extensible representation. In: *Georgia Institute of Technology, Tech. Rep* (2009)
- [BBMW09] BARRETT, Brian W. ; BERRY, Jonathan W. ; MURPHY, Richard C. ; WHEELER, Kyle B.: Implementing a portable multi-threaded graph library: The MTGL on Qthreads. In: *Parallel & Distributed Processing, 2009. IPDPS 2009. IEEE International Symposium on* IEEE, 2009, S. 1–8
- [BG11] BULUÇ, Aydin ; GILBERT, John R.: The Combinatorial BLAS: Design, implementation, and applications. In: *The International Journal of High Performance Computing Applications* 25 (2011), Nr. 4, S. 496–509
- [BHKK07] BERRY, Jonathan W. ; HENDRICKSON, Bruce ; KAHAN, Simon ; KONECNY, Petr: Software and algorithms for graph queries on multithreaded architectures. In: *Parallel and Distributed Processing Symposium, 2007. IPDPS 2007. IEEE International* IEEE, 2007, S. 1–14
- [BM11] BULUÇ, Aydin ; MADDURI, Kamesh: Parallel breadth-first search on distributed memory systems. In: *Proceedings of 2011 International Conference for High Performance Computing, Networking, Storage and Analysis* ACM, 2011, S. 65
- [C⁺05] CONSORTIUM, UPC u. a.: UPC language specifications v1. 2. In: *Lawrence Berkeley National Laboratory* (2005)
- [CCZ07] CHAMBERLAIN, Bradford L. ; CALLAHAN, David ; ZIMA, Hans P.: Parallel programmability and the chapel language. In: *The International Journal of High Performance Computing Applications* 21 (2007), Nr. 3, S. 291–312
- [CFS99] CARTER, Larry ; FEO, John ; SNAVELY, Allan: Performance and Programming Experience on the Tera MTA. In: *PPSC*, 1999
- [FFK16] FÜRLINGER, Karl ; FUCHS, Tobias ; KOWALEWSKI, Roger: DASH: a C++ PGAS library for distributed data structures and parallel algorithms. In: *High Performance Computing and Communications; IEEE 14th International Conference on Smart City; IEEE 2nd International Conference on Data Science and Systems (HPCC/SmartCity/DSS), 2016 IEEE 18th International Conference on* IEEE, 2016, S. 983–990
- [For12] FORUM, Message Passing I: *MPI: A Message-Passing Interface Standard Version 3.0.* 09 2012. – Chapter author for Collective Communication, Process Topologies, and One Sided Communications

Bibliography

- [FR11] FINEMAN, Jeremy T. ; ROBINSON, Eric: Fundamental graph algorithms. In: *Graph Algorithms in the Language of Linear Algebra* 22 (2011), S. 45
- [GL05] GREGOR, Douglas ; LUMSDAINE, Andrew: The parallel BGL: A generic library for distributed graph computations. In: *Parallel Object-Oriented Scientific Computing (POOSC)* 2 (2005), S. 1–18
- [GS13] GRÜNEWALD, Daniel ; SIMMENDINGER, Christian: The GASPI API specification and its implementation GPI 2.0. In: *7th International Conference on PGAS Programming Models* Bd. 243, 2013
- [Gus88] GUSTAFSON, John L.: Reevaluating Amdahl's Law. In: *Commun. ACM* 31 (1988), Mai, Nr. 5, 532–533. <http://dx.doi.org/10.1145/42411.42415>. – DOI 10.1145/42411.42415. – ISSN 0001–0782
- [ISO12] ISO: *ISO/IEC 14882:2011 Information technology — Programming languages — C++*. Geneva, Switzerland : International Organization for Standardization, 2012. – 1338 (est.) S. http://www.iso.org/iso/iso_catalogue/catalogue_tc/catalogue_detail.htm?csnumber=50372
- [ISO15] ISO: *C++ Extensions for Concepts*. Geneva, Switzerland : International Organization for Standardization, 2015 <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2015/n4553.pdf>
- [KBG12] KYROLA, Aapo ; BLELLOCH, Guy E. ; GUESTRIN, Carlos: Graphchi: Large-scale graph computation on just a pc USENIX, 2012
- [KDSA08] KIM, John ; DALLY, William J. ; SCOTT, Steve ; ABTS, Dennis: Technology-driven, highly-scalable dragonfly topology. In: *ACM SIGARCH Computer Architecture News* Bd. 36 IEEE Computer Society, 2008, S. 77–88
- [LAB⁺12] LUGOWSKI, Adam ; ALBER, David ; BULUÇ, Aydm ; GILBERT, John R. ; REINHARDT, Steve ; TENG, Yun ; WARANIS, Andrew: A flexible open-source toolbox for scalable complex graph analysis. In: *Proceedings of the 2012 SIAM International Conference on Data Mining* SIAM, 2012, S. 930–941
- [Lam13] LAMETER, Christoph: NUMA (Non-Uniform Memory Access): An Overview. In: *Queue* 11 (2013), Juli, Nr. 7, 40:40–40:51. <http://dx.doi.org/10.1145/2508834.2513149>. – DOI 10.1145/2508834.2513149. – ISSN 1542–7730
- [LGHB07] LUMSDAINE, Andrew ; GREGOR, Douglas ; HENDRICKSON, Bruce ; BERRY, Jonathan: Challenges in parallel graph processing. In: *Parallel Processing Letters* 17 (2007), Nr. 01, S. 5–20
- [MAB⁺10] MALEWICZ, Grzegorz ; AUSTERN, Matthew H. ; BIK, Aart J. ; DEHNERT, James C. ; HORN, Ilan ; LEISER, Naty ; CZAJKOWSKI, Grzegorz: Pregel: a system for large-scale graph processing. In: *Proceedings of the 2010 ACM SIGMOD International Conference on Management of data* ACM, 2010, S. 135–146
- [NR98] NUMRICH, Robert W. ; REID, John: Co-Array Fortran for parallel programming. In: *ACM Sigplan Fortran Forum* Bd. 17 ACM, 1998, S. 1–31

- [PJ98] PALSBERG, Jens ; JAY, C B.: The essence of the visitor pattern. In: *Computer Software and Applications Conference, 1998. COMPSAC'98. Proceedings. The Twenty-Second Annual International IEEE*, 1998, S. 9–15
- [PTM96] PROTIC, J. ; TOMASEVIC, M. ; MILUTINOVIC, V.: Distributed shared memory: concepts and systems. In: *IEEE Parallel Distributed Technology: Systems Applications* 4 (1996), Summer, Nr. 2, S. 63–71. <http://dx.doi.org/10.1109/88.494605>. – DOI 10.1109/88.494605. – ISSN 1063–6552
- [SAB⁺10] SARASWAT, Vijay ; ALMASI, George ; BIKSHANDI, Ganesh ; CASCAVAL, Calin ; CUNNINGHAM, David ; GROVE, David ; KODALI, Sreedhar ; PESHANSKY, Igor ; TARDIEU, Olivier: The asynchronous partitioned global address space model. In: *The First Workshop on Advances in Message Passing*, 2010, S. 1–8
- [SB13] SHUN, Julian ; BLELLOCH, Guy E.: Ligra: a lightweight graph processing framework for shared memory. In: *ACM Sigplan Notices* Bd. 48 ACM, 2013, S. 135–146
- [SL95] STEPANOV, Alexander ; LEE, Meng: *The standard template library*. Bd. 1501. Hewlett Packard Laboratories 1501 Page Mill Road, Palo Alto, CA 94304, 1995
- [SLL01] SIEK, Jeremy G. ; LEE, Lie-Quan ; LUMSDAINE, Andrew: *The Boost Graph Library: User Guide and Reference Manual, Portable Documents*. Pearson Education, 2001
- [Val90] VALIANT, Leslie G.: A bridging model for parallel computation. In: *Communications of the ACM* 33 (1990), Nr. 8, S. 103–111
- [ZMI⁺14] ZHOU, Huan ; MHEDHEB, Yousri ; IDREES, Kamran ; GLASS, Colin W. ; GRACIA, José ; FÜRLINGER, Karl: DART-MPI: an MPI-based implementation of a PGAS runtime system. In: *Proceedings of the 8th International Conference on Partitioned Global Address Space Programming Models* ACM, 2014, S. 3
- [ZZZ⁺14] ZHAO, D. ; ZHANG, Z. ; ZHOU, X. ; LI, T. ; WANG, K. ; KIMPE, D. ; CARNS, P. ; ROSS, R. ; RAICU, I.: FusionFS: Toward supporting data-intensive scientific applications on extreme-scale high-performance computing systems. In: *2014 IEEE International Conference on Big Data (Big Data)*, 2014, S. 61–70