

Assessing the Performance Impact of an Active Global Address Space

Parsa Amini¹, Hartmut Kaiser², and Jeanine Cook³

¹ `parsa@cct.lsu.edu`,

² `hkaiser@cct.lsu.edu`

³ `jeacock@sandia.gov`

Abstract. In this research, we describe the functionality of AGAS, a subsystem of the HPX runtime system that is designed to handle data, independent of the hardware and architecture configuration. AGAS performs enables global data accesses and data migration at the expense of performing these operations at runtime. We present a method to assess the performance of AGAS and the amount of burden it imposes on application execution. With our assessment method we identified three problematic spots in the HPX version used in our experiments that need to be fixed.

Keywords: HPX, AGAS, PGAS

1 Introduction

Global address space systems attempt to boost productivity and simplify the application development cycle of distributed applications by providing a ubiquitous abstraction layer over memory spaces that are provided and managed by the operating system on each computer. SPMD-style (Single Program Multiple Data) Partitioned Global Address Space designs eliminate this layer during compilation to avoid the complexity of resolving global addresses at runtime at the cost of limiting productivity and imposing limitations on the code, while others like HPX[18, 4], HPX5[3], Charm++[19], UPC++[29], and Chapel[1] provide a runtime component that maps global addresses to virtual addresses during application execution to provide true global addresses.

The demand for models that enable applications to process massive datasets within specific time, power, and budgetary constraints continues to pose challenges for the computer science community [24, 26]. One category of such complexities includes managing a large quantity of objects across several machines and memory partitions while maximizing data locality. Several Partitioned Global Address Space systems (PGAS) [7] try to address these needs by providing control over data distribution and facilitating data accesses across processes and machines. However, initial data placement alone is not sufficient to address more complex issues such as load balancing applications that run on heterogeneous

clusters or applications that become scaling impaired over time due to increasing load imbalance such as adaptive mesh refinement (AMR), dynamic graph applications, and partial differential equation (PDE) solvers[17, 10, 12]. Active Global Address Space (AGAS) is a system that tries to address performance impaired applications. AGAS was initially proposed for the ParalleX programming model[17] and implemented in HPX. It adds an abstraction layer on top of local objects on each computing node by mapping local virtual addresses to a global address and ensuring that global addresses are valid even if the object it refers to is migrated to a different physical location. This abstraction layer needs to execute code to resolve and maintain the references and takes a portion of the execution time.

Assessing the overheads caused by AGAS is the main objective of this paper. We address this objective by developing a system that uses measurement data from AGAS and applies it to identify AGAS functions that exhibit poor scaling behavior in HPX 0.9. In the rest of this work, we discuss other pertinent research in section 2, present a general overview of AGAS functionality in section 3, explain the criteria based on which the results can be evaluated in section 4, experiments that were run and examine their results in section 5, and further analyze them and consider directions that are likely to produce more insights into improving AGAS considering our findings in section 6.

2 Related Work

In recent years, there has been a significant increase in utilization of heterogeneous clusters that use GPUs and MICs in addition to CPUs[20, 28, 21, 25]. One approach to manage such systems is using solutions[22, 28, 11] that separately use a library like MPI for explicit communication between nodes and a choice of shared-memory programming framework such as OpenMP[8], Kokkos[13, 5], or UPC[9]. Other approaches include Asynchronous PGAS runtimes[23] and Charm++[2] that use dynamic multithreading to avoid fragmenting the application development process to separately manage communication and computation while maintaining portability between various cluster configurations and providing access to heterogeneous computing resources[27].

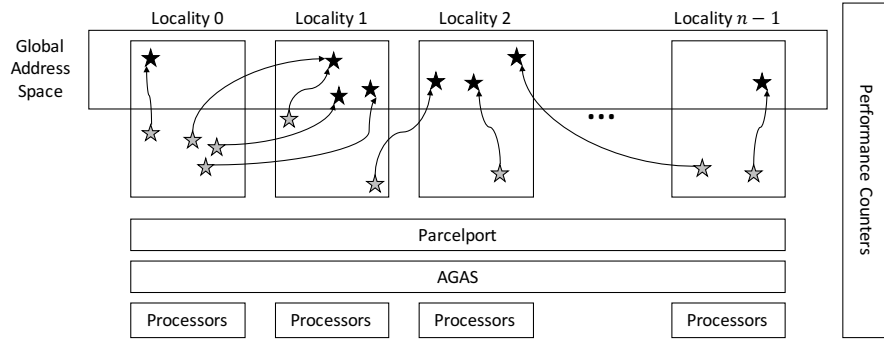
Most studies on distributed runtime systems do not include quantitative analysis of the performance of their global address space system but present the overall performance of applications using the respective model or implementation. However, amongst the runtime systems mentioned only HPX has a ubiquitous address system that is present during run time. This unique property calls for a closer look into HPX’s Active Global Address Space system behavior and performance and is the motivation behind this study.

3 Active Global Address Space

This section provides an overview of the implementation of the Active Global Address Space (AGAS) in HPX runtime system.

AGAS is a global memory addressing system that is designed to handle various memory configurations ranging from a single small embedded machine running a specialized operating system to computer clusters composed of a large number of nodes with heterogeneous computing resources. Fig. 1 illustrates AGAS's role in HPX. Typically, each part of the distributed AGAS service is hosted on a separate node of a cluster. AGAS consists of several subsystems that are depicted in Fig. 2 and are the following:

Fig. 1. AGAS provides an abstraction layer on top of virtual addresses local to each locality. Global objects are shown as black stars and gray stars indicate references to global objects. Each reference is connected to global objects by an arrow.



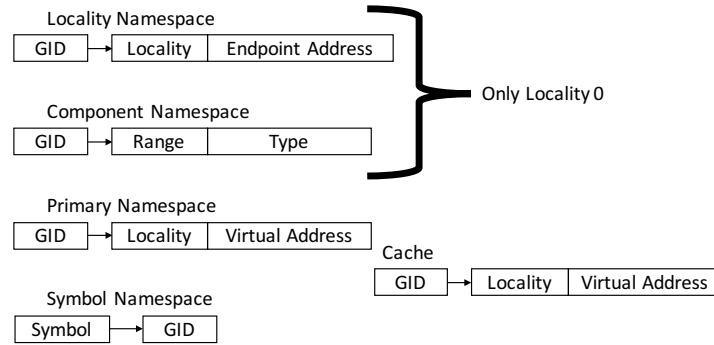
1. **Primary Namespace:** To provide uniform access to objects across the boundaries of physical partitions in a cluster, AGAS provides applications with 128-bit global identifiers (GIDs) to be used in place of virtual addresses that are local to specific nodes. Consequently, AGAS maintains mapping tables to be able to map GIDs to local virtual addresses.
2. **Locality Namespace:** Information about the nodes and computing resources allocated to each physical partition is held in the locality namespace. Each partition is called a "locality" and locality 0 is responsible for maintaining current information about all other localities.
3. **Component Namespace:** Types are registered in this namespace to facilitate resolving resource requirements during bulk memory allocations.
4. **Symbolic Namespace:** The symbolic namespace is a layer on top of the global address space that allows mapping symbolic names to global addresses for the purpose of resolving global addresses at runtime. This is useful in cases

where data about specific events needs to be collected. For example, HPX performance counter system uses the symbolic namespace for collecting performance counter data.

5. AGAS Cache: AGAS cache stores mapping between the most recently used global addresses to localities and local virtual addresses. If a task requires an object that does not live on the same locality then the task is sent to the locality where the object currently is. However, in order to do so, AGAS has to determine the current location of the queried object. If AGAS does not know the current location then it forwards the query to the locality where the object was originally created. The locality on which an object is created stays responsible for maintaining the current location of the object during its entire lifetime. If an object is likely to be accessed again then the locality that forwards a task will also request the original locality to report the object's current location. This information is stored in the AGAS cache. The AGAS cache is small since it is designed for speed and therefore, it does not know all local objects.
6. Garbage Collection: A distributed garbage collection system tracks objects during their lifetime and frees the consumed memory when an object goes out of scope and therefore can no longer be accessed in the program. HPX complies to the C++ standard specifications and hence, it provides structures similar to what the modern C++ standard requires. Additionally, GIDs in HPX can be managed or unmanaged and in case of the former AGAS tracks that GID until the reference is lost so that it can free up the memory space when possible. AGAS uses reference counting to determine if there are existing references to an object and a credit-based scheme for remote references. The local counter is updated when a new reference is created and when a reference goes out of scope on the same locality. As for remote references, the credit system works as follows:
 - An object is created, AGAS sets the object's credit to a large number.
 - When an object is referenced by another locality the object's credit is split in half and a copy of the reference is kept at both localities, and a flag is set on the original reference to indicate that the object is referenced globally.
 - When a copy runs out of credit, it will ask the lender AGAS instance for more. If the original reference itself does not have enough credit, the AGAS instance responsible for that reference will grant it more credits.
 - When a reference goes out of scope AGAS returns all borrowed credits to the original reference.
 - If there are no local references and all credits have been returned then AGAS can remove the object from memory during garbage collection.

Garbage collection at runtime requires execution of code that is otherwise not present and consumes computing resources. Performing garbage collection requires executing code that is not the user's application. AGAS tries to minimize garbage collection sweeps by performing it when the volume of garbage reaches a certain threshold that can be specified by the users. It is also possible to manually initiate it inside applications by developers.

Fig. 2. The four namespaces inside AGAS. Primary namespace on each AGAS instance contains GID to local address mappings. Locality namespace holds information about all AGAS instances. Component namespace tracks bulk memory allocations dedicated to types. Symbolic namespace contains mappings between GIDs and special strings that can be used for various purposes such as facilitating the collection of data about application execution.



3.1 AGAS Performance Counters

HPX is a runtime system that includes novel abstractions on top of ordinary operating systems and hardware that are more difficult to benchmark using traditional performance measuring tools such as hardware performance counters included in Intel processors. HPX introduces a set of performance counters to let developers monitor the performance of its subsystems, including AGAS, during execution.

Similar to hardware performance counters, HPX performance counters are designed to expose performance data on the underlying function calls and their durations.

3.2 HPX bootstrap and teardown

Before an HPX application can start executing, HPX has to initialize. This process is called bootstrap and includes registering runtime services, data types, performance counters[14], and symbols. Similarly during teardown, after applications built with HPX complete, HPX removes all objects, frees all allocated hardware resources and may perform additional tasks such as collecting performance counter data when applicable.

In this research, we exclude data from the bootstrap and teardown phases since they do not directly provide useful data on how AGAS might impact the performance or scaling behavior during runtime.

4 Quantifying AGAS Performance

The main objective of any HPC memory addressing system, including AGAS, is to handle massive amount of data that is used by applications that use it. In HPX this task is performed by AGAS. To study AGAS, we use OctoTiger[15, 6] as our application and run it on the cluster provided by the Center of Computation and Technology at Louisiana State University. In the rest of this section we discuss the applications we use to benchmark AGAS and the performance metrics that we use judge AGAS based on.

4.1 OctoTiger

LSU OctoTiger is a state of the art multi-physics AMR HPX application that simulates the merger of double white dwarf binaries. Its computations include hydrodynamics, gravitational, and radiation transportation solvers. It fully utilizes the capabilities of AGAS in HPX since it dynamically changes the computational resolution based on the actual needs of the simulation and thus, it is an inherently balance impaired application.

We use OctoTiger to study the behavior of AGAS when used by applications running on large machines. We ran our experiments on SuperMIC at Louisiana State University, a hybrid cluster composed of Xeon, Xeon Phi Knights Corner, and NVIDIA Tesla processors. Due to the complexity of the computations performed by OctoTiger, adjusting the size of the problem with reasonable accuracy to demonstrate weak scaling is not practical and therefore we presenting the impacts of strong scaling on AGAS behavior.

4.2 System and Environment Setup

We run our experiments on SuperMIC, a hybrid Xeon/Xeon Phi cluster that comprises 360 compute nodes located at Louisiana State University. More details about the configuration of SuperMIC is included in Table 1. We run OctoTiger on SuperMIC from two nodes to 200 nodes only using the Xeon processors on each machine using HPX 0.99[16].

4.3 Performance Metrics

Every globally accessible object in HPX has a global id that AGAS manages and resolves to local virtual addresses at runtime. However, address resolution at runtime is an overhead that consumes computational resources that application developers expect to be used by the applications rather than the runtime system. To quantify and study the overhead of AGAS, we look at the following fundamental operations, the number of calls, and the amount of time that is spent performing these operations.

Table 1. SuperMIC Configuration

Nodes available	360
Architecture	Ivy Bridge, Knights Corner
Cores/Node	2×10 cores
Processor	Intel Xeon E5-2680 @ 2.8 GHz, Intel Xeon Phi 7120P @1.238 GHz
Memory	64 GB DDR3 @ 1,866 MHz
Connection	56 Gbps Infiniband
Operating System	Red Hat Enterprise Linux 6.8

Bind, Unbind, Object Lookup Bind and Unbind operations occur when a global object is created and deleted, respectively. An object lookup operation takes place each time AGAS attempts to resolve a global Id.

Locality Lookup Each AGAS instance only knows about itself and locality 0. When an AGAS instance needs to communicate with another locality and does not have information about the appropriate communication endpoint to do so it needs to query that information from locality 0.

Parcel Routing HPX implements active messages in the form of parcels. A parcel is packaged information that triggers an operation upon reception by an AGAS instance. For example, when an HPX task needs to operate on information that resides on another locality, AGAS serializes the task along with its arguments and state information and forwards it to the locality on which data resides.

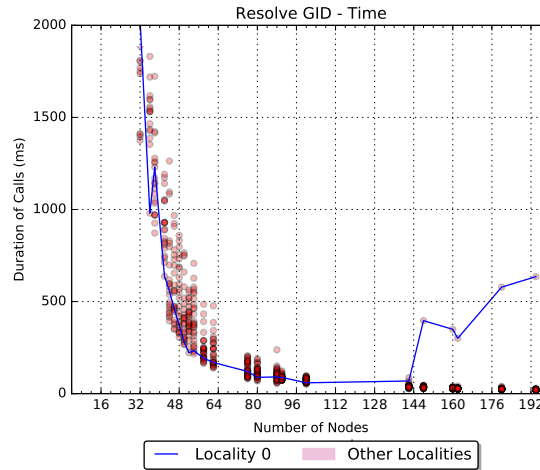
Cache Whenever AGAS decides that an address is likely to be queried again, it stores that information in the local AGAS cache of that locality. To serve its purpose of actually improving performance the AGAS cache is designed to hold a limited number of entries that is defined by user configuration.

Garbage Collection HPX provides managed objects whose lifetime is controlled by the garbage collection mechanism instead of the programmer having to free the memory consumed by each object. HPX uses reference counting to track references local to each AGAS instance and a credit scheme to manage global references. Once an object no longer has any local references and/or all its global reference credits are returned it is deemed garbage and is collected when the garbage collection operation is triggered by reaching a certain threshold or by the application developer. Registration and removal of global references is done by calling `increment_credit` and `decrement_credit` functions, respectively.

5 Performance Results

As mentioned before, the main objective of this study is to identify the AGAS operations that need performance improvement and to study the effects of strong scaling on AGAS. In this section, we present the results of this study.

Fig. 3. Total time spent to resolve GID calls while running OctoTiger on 2 to 200 nodes of SuperMIC with fixed problem size (Strong Scaling). Each (red) circle refers to the measured amount of time a locality has spent executing resolve GID calls. Higher intensities indicate overlapping measurements. The line highlights locality 0's behavior.

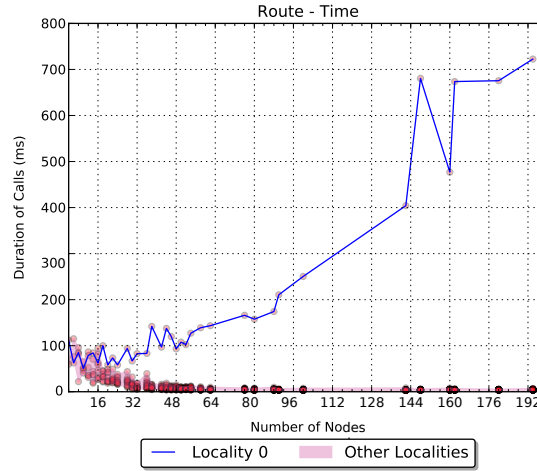


Appropriate distribution of work is of particular interest to computer scientists. In case of AGAS, ideally we expect the number of specific operations and the amount of time spent to perform them to be similar. Address resolution is one of the main functions of AGAS that translates a HPX GID to the actual virtual address on a machine. Its operation is complicated if the object is not currently located on the locality address resolution is taking place and in such case AGAS tries to determine which locality the object currently lives on, serialize the task that asked to access the foreign GID in an HPX parcel, and send it to the locality object is living in. If an AGAS instance has no knowledge of the locality that is holding an object then it takes advantage of the fact that a locality an object is created on stays responsible for it during the object's entire lifetime and uses the metadata in the GID to determine the locality the object

was originally created on and sends the query there. Fig. 3 shows the amount of time GID resolution takes while running and OctoTiger. By looking at the graph we notice that while the the amount of work performed by each locality is relatively similar up to 146 nodes, locality 0 starts to perform more work than other localities. This may indicate that most objects globally referenced were created on locality 0 and it requires further investigation to determine if this effect is caused by application behavior, non-optimal data distribution, or HPX itself.

Parcels are the basic communication block in HPX. Parcels are active messages that trigger an operation on the target AGAS instance that opens them and usually contain a serialized task and its arguments. Parcel routing is the operation in which the a parcel is sent to a different locality. Fig. 4 shows that locality 0 ends up handling a disproportionate and increasing number of routing requests regardless of the application profile and problem size and hence, designate the routing operation as candidate for optimization.

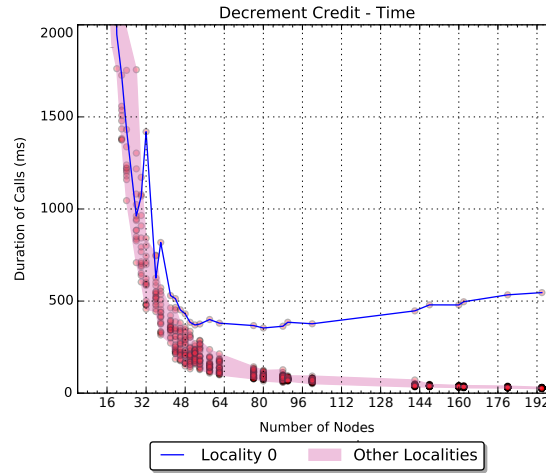
Fig. 4. Total time spent in route calls while running OctoTiger on 2 to 200 nodes of SuperMIC with fixed problem size (Strong Scaling). Each (red) circle refers to measured amount of time a locality has spent executing resolve GID calls. Higher intensities indicate overlapping measurements. The line highlights locality 0's behavior.



For a garbage collection system to function the HPX runtime system must be able to determine if an object is currently being used. HPX uses reference

counting for local references and a credit-based system to keep track of global references. The operations that lend and return credits are called decrement credit and increment credit, respectively. Fig 5 depicts the behavior of decrement_credit calls taken place while running OctoTiger with the same number of objects as the number of nodes are increased. This figure shows that a linear increase in number of objects results in an exponential growth in the amount of time spent on managing global references. Although this behavior is likely to be different for applications that do not have as many global references per each object and therefore decrement_credit is a good candidate for further optimizations.

Fig. 5. Total time spent in decrement calls while running OctoTiger on 2 to 200 nodes of SuperMIC with fixed problem size (Strong Scaling). Each (red) circle refers to measured amount of time a locality has spent executing resolve GID calls. Higher intensities indicate overlapping measurements. The line highlights locality 0's behavior.



6 Conclusions and Implications for Future Work

In this work we introduce AGAS, its subsystems, and a method to study how the amount of time that is spent executing AGAS code, an overhead that does not exist in PGAS systems that statically resolve global references during compilation, is affected as the number of nodes increase.

To study AGAS's behavior we chose a multiphysics AMR application called OctoTiger that generates and works with a significant number of objects. We identify the performance metrics that reveal AGAS's performance and used the corresponding counters in HPX's Performance Counter framework to collect performance data from our strong scaling experiments.

Our observations show that in the cases of three basic AGAS operations, resolve GID, route, and decrement credit, the AGAS instance on locality 0 performs an increasing amount of work as the problem is strongly scaled. Fig. 4 shows parcel routing as the operation that is affected the most with an aggregated sum of 722 milliseconds (0.0016% of execution time) spent on routing parcels by AGAS-invoked tasks on locality 0 on 200 nodes. Using our methods we were able to identify these issues which pose possible performance bottlenecks that are being investigated and preliminary results appear to show that this scaling issue has been addressed. Inclusion of a copy of the locality namespace to each locality's AGAS cache in the current version of HPX is one such optimization that reduces traffic forwarding to locality 0 due to not knowing the endpoint address for a locality.

References

1. Cascade High Productivity Language - Cray. <http://chapel.cray.com/>
2. Charm++. <http://charm.cs.illinois.edu>
3. HPX-5. <http://hpx.crest.iu.edu/>
4. HPX on Github. <https://github.com/STELLAR-GROUP/hpx>
5. Kokkos on Github. <https://github.com/kokkos/kokkos>
6. OctoTiger on Github. <https://github.com/STELLAR-GROUP/octotiger>
7. PGAS - Partitioned Global Address Space Languages. <http://www.pgas.org>
8. The OpenMP API specification for parallel programming. <http://www.openmp.org/specifications/>
9. Unified Parallel C. <https://upc-lang.org>
10. Anderson, M., Brodowicz, M., Kaiser, H., Adelstein-Lelbach, B., Sterling, T.: Adaptive Mesh Refinement for Astrophysics Applications with ParalleX (oct 2011), <http://arxiv.org/abs/1110.1131>
11. Chorley, M.J., Walker, D.W.: Performance analysis of a hybrid MPI/OpenMP application on multi-core clusters. *Journal of Computational Science* 1(3), 168–174 (aug 2010), <http://linkinghub.elsevier.com/retrieve/pii/S1877750310000396>
12. Dekate, C.: Extreme Scale Parallel N-Body Algorithm with Event-Driven Constraint-Based Execution Model. Ph.D. thesis, Louisiana State University, Baton Rouge, LA, USA (2011)
13. Edwards, H.C., Trott, C.R., Sunderland, D.: Kokkos: Enabling manycore performance portability through polymorphic memory access patterns. *Journal of Parallel and Distributed Computing* 74(12), 3202 – 3216 (2014), <http://www.sciencedirect.com/science/article/pii/S0743731514001257>, domain-Specific Languages and High-Level Frameworks for High-Performance Computing
14. Grubel, P., Kaiser, H., Cook, J., Serio, A.: The performance implication of task size for applications on the hpx runtime system. In: *Cluster Computing (CLUSTER)*, 2015 IEEE International Conference on. pp. 682–689. IEEE (2015)

15. Kadam, K., Clayton, G.C., Motl, P.M., Marcello, D., Frank, J.: Numerical simulations of close and contact binary systems having bipolytropic equation of state. In: American Astronomical Society Meeting Abstracts. vol. 229 (2017)
16. Kaiser, H., Adelstein-Lelbach, B., Heller, T., Berg, A., Bikineev, A., Mercer, G., Habraken, J., Anderson, M., Serio, A., Biddiscombe, J., Stumpf, M., Schfer, A., Brandt, S.R., Bourgeois, D., Grubel, P., Amatya, V., Huck, K., Bacharwar, D., Viklund, L., Yang, S., Schnetter, E., Bcorde5, Brodowicz, M., Byerly, Z., bghimire, Jakhola, P., Bross, C., andreasbuhr, Guo, C., atrantan: hpx: HPX V0.9.11: A general purpose C++ runtime system for parallel and distributed applications of any scale (Nov 2015), <https://doi.org/10.5281/zenodo.33656>
17. Kaiser, H., Brodowicz, M., Sterling, T.: ParalleX An Advanced Parallel Execution Model for Scaling-Impaired Applications. In: Barolli, L.F., Feng, W.c.V.T. (eds.) 2009 International Conference on Parallel Processing Workshops. pp. 394–401. ICPPW, IEEE, Vienna, Austria (sep 2009), <http://stellar.cct.lsu.edu/pubs/icpp09.pdf> <http://ieeexplore.ieee.org/document/5364511/>
18. Kaiser, H., Heller, T., Adelstein-Lelbach, B., Serio, A., Fey, D.: HPX - A Task Based Programming Model in a Global Address Space. In: Proceedings of the 8th International Conference on Partitioned Global Address Space Programming Models. p. 6. ACM, Eugene, OR, USA (2014), <http://stellar.cct.lsu.edu/pubs/pgas14.pdf> <http://doi.acm.org/10.1145/2676870.2676883>
19. Kale, L.V., Krishnan, S.: CHARM++. In: Proceedings of the eighth annual conference on Object-oriented programming systems, languages, and applications - OOPSLA '93. OOPSLA '93, vol. 28, pp. 91–108. ACM Press, New York, New York, USA (oct 1993), <http://portal.acm.org/citation.cfm?doid=165854.165874>
20. Lena, O.: GPI2 for GPUs: A PGAS framework for efficient communication in hybrid clusters. Parallel Computing: Accelerating Computational Science and Engineering (CSE) 25, 461 (2014)
21. Potluri, S.: Enabling Efficient Use of MPI and PGAS Programming Models on Heterogeneous Clusters with High Performance Interconnects. Ph.D. thesis, Ohio State University (2014), https://etd.ohiolink.edu/pg_10?0::NO:10:P10_ACCESSION_NUM:osu1397797221
22. Rabenseifner, R., Hager, G., Jost, G.: Hybrid MPI/OpenMP Parallel Programming on Clusters of Multi-Core SMP Nodes. In: 2009 17th Euromicro International Conference on Parallel, Distributed and Network-based Processing. pp. 427–436. PDND '09, IEEE, Weimar, Germany (feb 2009), <http://ieeexplore.ieee.org/document/4912964/>
23. Saraswat, V., Almasi, G., Bikshandi, G., Cascaval, C., Cunningham, D., Grove, D., Kodali, S., Peshansky, I., Tardieu, O.: The Asynchronous Partitioned Global Address Space Model. In: The First Workshop on Advances in Message Passing. pp. 1–8 (2010)
24. Sarkar, V., Harrod, W., Snavely, A.E.: Software challenges in extreme scale systems. Journal of Physics: Conference Series 180, 012045 (jul 2009), <http://stacks.iop.org/1742-6596/180/i=1/a=012045?key=crossref.f67cd168995105550454fe5e8f7d8cf2>
25. Sidelnik, A., Chamberlain, B.L., Garzaran, M.J., Padua, D.: Using the High Productivity Language Chapel to Target GPGPU Architectures. Tech. rep. (2011), <http://hdl.handle.net/2142/18874>
26. Sterling, T., Stark, D.: A High-Performance Computing Forecast: Partly Cloudy. Computing in Science & Engineering 11(4), 42–49 (jul 2009), <http://ieeexplore.ieee.org/document/5076318/>

27. Wong, M., Kaiser, H., Heller, T.: Towards Massive Parallelism (aka Heterogeneous Devices/Accelerator/GPGPU) support in C++ with HPX. Tech. rep. (2015), <http://www.open-std.org/jtc1/sc22/wg21/docs/papers/2016/p0234r0.pdf>
28. Yang, C.T., Huang, C.L., Lin, C.F.: Hybrid CUDA, OpenMP, and MPI parallel programming on multicore GPU clusters. *Computer Physics Communications* 182(1), 266–269 (jan 2011), <http://linkinghub.elsevier.com/retrieve/pii/S0010465510002262>
29. Zheng, Y., Kamil, A., Driscoll, M.B., Shan, H., Yelick, K.: UPC++: A PGAS Extension for C++. In: 2014 IEEE 28th International Parallel and Distributed Processing Symposium. pp. 1105–1114. IPDPS 14, IEEE, Phoenix, AZ, USA (may 2014), <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6877339>