



HPX Applications and Performance Adaptation

Alice Koniges¹, Jayashree Ajay Candadai², Hartmut Kaiser³, Kevin Huck⁴, Jeremy Kemp⁵, Thomas Heller⁶, Matthew Anderson², Andrew Lumsdaine², Adrian Serio³, Michael Wolf⁷, Bryce Adelstein Lelbach¹, Ron Brightwell⁷, Thomas Sterling²

¹Berkeley Lab, ²Indiana University, ³Louisiana State University, ⁴University of Oregon, ⁵University of Houston, ⁶Friedrich-Alexander University, ⁷Sandia National Laboratories



The HPX runtime system is a critical component of the DOE XPRESS (eXascale Programming Environment and System Software) project and other projects world-wide. We are exploring a set of innovations in execution models, programming models and methods, runtime and operating system software, adaptive scheduling and resource management algorithms, and instrumentation techniques to achieve unprecedented efficiency, scalability, and programmability in the context of billion-wide parallelism. A number of applications have been implemented to drive system development and quantitative evaluation of the HPX system implementation details and operational efficiencies and scalabilities.

Applications and Characteristic Behavior

- LULESH:** (Livermore Unstructured Lagrangian Explicit Shock Hydrodynamics) For details see Deep Dive to right.
- Mini-ghost:** A mini-app for exploring boundary exchange strategies using stencil computations in scientific parallel computing. Implemented by decomposing the spatial domain, inducing a "halo exchange" of process-owned boundary data.
- N-Body Code:** An event driven constraint based execution model using the Barnes-Hut algorithm where the particles are grouped by a hierarchy of cube structures using a recursive algorithm. It uses an adaptive octree data structure to compute center of mass and force on each of the cubes with resultant $O(N \log N)$ computational complexity making use of LibGeoDecomp, an auto-parallelizing library.
- PIC:** Two 3D particle-in-cell (PIC) codes – GTC and PICAR. The gyrokinetic toroidal code (GTC) was developed to study turbulent transport in magnetic confinement fusion plasmas. It models the interactions between fields and particles by solving the 5D gyro-averaged kinetic equation coupled to the Poisson equation. PICAR is a mini-app with the key functionalities of PIC accelerator codes, including a Maxwell solver using an arbitrary order finite-difference scheme (staggered/centered), a particle pusher using the Boris algorithm, and an energy conserving field gathering with high order particle shape factors.
- miniTri:** A newly developed triangle enumeration-based data analytics minapp. miniTri mimics the computation requirements of an important set of data science applications, not well represented by traditional graph search benchmarks such as Graph500. An asynchronous HPX-based approach enables our linear algebra-based implementation of miniTri to be significantly more memory efficient, allowing us to process much larger graphs.
- CMA:** (Climate Mini-App) For details see Deep Dive to right.
- Kernels:** Various computational kernels, such as matrix transpose and fast multiple algorithms, which are used to explore features of HPX and compare to other approaches.

LULESH (Deep Dive)

LULESH solves the Sedov blast wave problem. In three dimensions, the problem is spherically-symmetric and the code solves the problem in a parallelepiped region. In the figure, symmetric boundary conditions are imposed on the colored faces such that the normal components of the velocities are always zero. Free boundary conditions are imposed on the remaining boundaries.



The LULESH application is implemented as a hexahedral mesh-based code with two centerings. Element centering stores thermodynamic variables such as energy and pressure. Node centering stores kinematic variables such as positions and velocities. The simulation is run via time integration using a Lagrangian leapfrog algorithm. There are three main computational phases within each time step: advance node quantities, advance element quantities, and calculate time constraints. There are three communication patterns, each regular, static, and uniform: face adjacent, 26 neighbor, and 13 neighbor communications, illustrated below:



[1] LULESH is available from: <http://webpages.tdt.edu/~lshayley/>

CMA (Deep Dive)

The climate mini-app (CMA) models the performance profile of an atmospheric "dynamic core" (dcycore) for non-hydrostatic flows. The code uses a conservative finite-volume discretization on an adaptively-refined cubed-sphere grid. An implicit-explicit (IMEX) time integrator combines a vertical implicit operator (which is FLOP-bound) with a horizontal explicit operator (which is bandwidth-bound). There are three major sources of load imbalance in the code:

- The number of iterations required to perform the non-linear vertical solves will vary across the grid.
- Exchanges across the boundaries of the six panels of the cubed sphere are more computationally expensive than intra-panel exchanges.
- The adaptively refined mesh will add and removed refined regions as the simulation evolves.

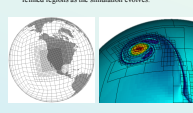
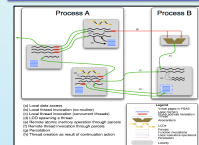


Figure: Left image shows an example of an adaptively refined cubed-sphere grid used in climate codes. Right image shows vorticity dynamics for a climate test problem with AMR.

The mini app is implemented using the Chombo adaptive mesh refinement (AMR) framework and has been built as MPI+OMP and HPX backend. The mini-app is being used to explore performance on multi-core architectures (e.g. Xeon Phi) and to explore the benefits of using HPX for finite-volume AMR codes to combat dynamic load imbalance.

Background Info

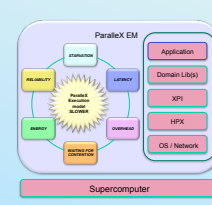
Parallel Execution Model



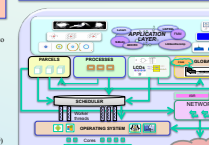
- Lightweight multi-threading**
 - Divides work into smaller tasks
 - Increases concurrency
- Message-driven computation**
 - Move work to data
 - Keeps work local, steps blocking
- Constraint-based synchronization**
 - Declarative criteria for work
 - Event driven
 - Eliminates global barriers
- Data-directed execution**
 - Merge of flow control and data structure
 - Shared name space
 - Global address space
 - Simplifies random gathers

High Performance ParallelX (HPX)

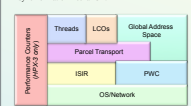
- The HPX runtime system refines the ParallelX execution model to support large-scale irregular applications:
 - Localities
 - Active Global Address Space (AGAS)
 - ParallelX Processes
 - Complexes (ParallelX Threads and Thread Management)
 - Parcel Transport and Parcel Management
 - Local Control Objects (LCOs)
- Sits between the application and OS
- Portable interface: C++11/14 (HPX-3 only), XPI
- Comprehensive set of parallel C++ algorithms (HPX-2 only)
- Automatic distributed garbage collection in AGAS (HPX-3 only)
- Flexible set of execution and scheduling policies
- Performance counter framework (HPX-2 only)



HPX Architecture



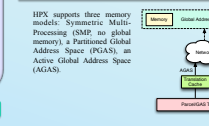
HPX avoids the use of locks and/or barriers in parallel computation through the use of LCOs, which are lightweight synchronization objects used by threads in a control mechanism. Reads and writes on LCOs are globally atomic and require no other synchronization mechanism.



HPX has optimized transports built on top of Photon (HPX-3 only) and other communication libraries

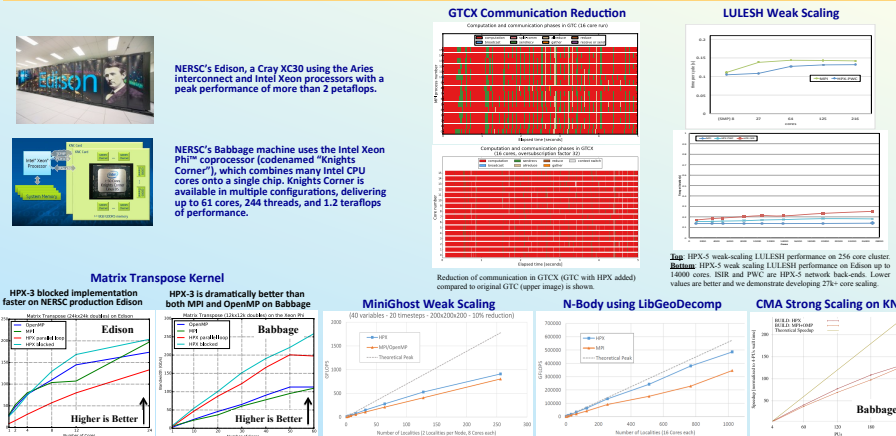
- Two-sided handshaking transport (ISR)
- Pre-points receive to reduce probe overhead
- One-sided Put-With-Command completion (PWC)
- Local remote modifications for RDMA operations
- RDMA communication optimizations using Photon: turn large puts into gets, buffer coalescing

Memory Models and Transport



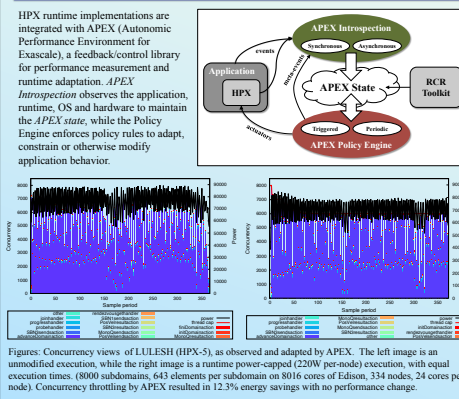
- [1] Thomas Sterling, Daniel Kogler, Matthew Anderson, and Marco Budnik. SLURM: A performance model for parallel computing. *Supercomputing Frontiers and Innovations*, 1:42–57, September 2014.
- [2] Hartmut Kaiser, Thomas Heller, Bryce Adelstein-Lelbach, Adam Serio, Dorian Fox, HPX: A Task-Based Programming Model in a Global Address Space. *PLAS 2014: The 18th International Conference on Partitioned Global Address Space Programming Models* (2014).

Results Showing Benefits of HPX



Performance Adaptation, Legacy Applications, and Summary

APEX : Performance Adaptation

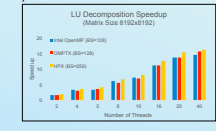


Legacy Application Support

OMPTX is an HPX implementation of the Intel OpenMP runtime, enabling existing OpenMP applications to execute with HPX.



Within a node, performance achieved for OpenMP programs using OMPTX is comparable to the Intel OpenMP Runtime. For example, below we show speedup for a blocking LU decomposition benchmark using the optimal block size for each implementation. These results were obtained with a dual-socket Ivy Bridge processor at LSU.



Summary

Exascale programming models and runtime systems are at a critical juncture in development. Systems based on light-weight tasks and data dependence are an excellent method for extracting parallelism and achieving performance. HPX is emerging as an important new path with support from US Department of Energy, the National Science Foundation, the Bavarian Research Foundation, and the European Horizon 2020 Programme. Application performance of HPX codes on very recent architectures, including current and prototypical next-generation Cray-Intel machines, is very good. For some applications, the performance using HPX is significantly better than standard MPI + OpenMP implementations. Performance adaptation using APEX provides significant energy savings with no performance change. We have shown that legacy applications using OpenMP can run under the HPX runtime system effectively.

