Interoperating MPI and Charm++ for Productivity and Performance

Nikhil Jain¹, Abhinav Bhatele², Jae-Seung Yeom², Mark F. Adams³, Francesco Miniati⁴, Chao Mei⁵, and Laxmikant Kale¹ ¹University of Illinois at Urbana-Champaign, ²Lawrence Livermore National Laboratory, ³Lawrence Berkeley National Laboratory, ⁴ETH Zurich, ⁵Google Inc.

Introduction

MOTIVATION

Developing multi-physics, coupled applications using one parallel language can be challenging:

- a) For optimal performance, different modules may require features provided by distinct parallel programming languages.
- b) A good match between module requirements and language features may increase programmer's productivity if the module is implemented in the given language.
- c) Various modules may require existing software implemented in other different languages.

INTEROPERATION

- Coordinated use of multiple parallel programming languages in an application, termed interoperation, is a realistic solution to challenges described above.
- We explore interoperation between these two languages: User-driven **MPI** where the programmer explicitly defines the control flow.

System-driven**Charm++** where a runtime drives the execution based on the availability of data.

CHALLENGES

- Given the difference in control flow between Charm++ and MPI, how is control transferred and managed between them?
- How are resources (such as cores and network) and data shared between the two programming languages?
- How easy-to-use and scalable is the interoperation methodology for production applications?

Framework

charm_module1(data) {

harm_module2(data) {

// do work

CONTROL FLOW

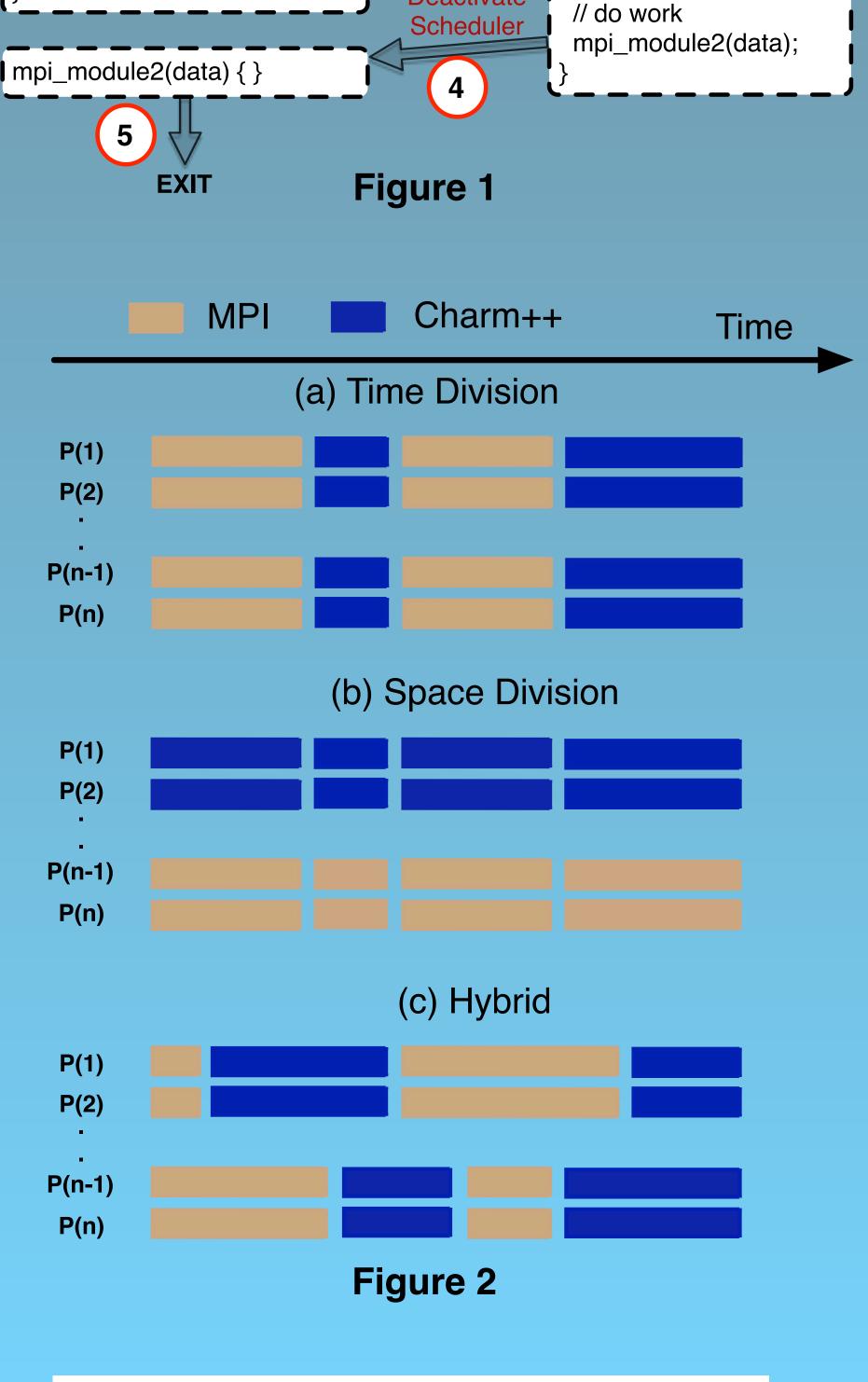
- Expose Charm++ scheduler to the programmer: new API to activate and deactivate the scheduler.
- Explicit control transfer by the programmer: begin in MPI, switch to Charm++ by invoking its scheduler, switch back to MPI by stopping the scheduler, and so on.
- Switching control infrequently is advisable for ease of programming and good performance.

RESOURCE SHARING

- ✓ Automatic shared use of low-level network resources via distinct communication domains, e.g. create clients using PAMI_Client_create on Blue Gene/Q.
- \checkmark Three schemes available for sharing cores: 1. *Time Division*: execute either Charm++ or MPI on all cores at a given time during execution. 2. Space Division: statically divide cores into subsets that execute either Charm++ or MPI. 3. Hybrid Division: divide cores into subsets that alternate between MPI and Charm++ together.
- \checkmark API for division of cores based on MPI communicators.

DATA SHARING

- ✓ Method 1: A generic *data transfer repository* to which each language module deposits (or queries) the data. Under the hood, communication performed by the repository.
- ✓ Method 2: *Pointer-based* data sharing via reserved memory within nodes assisted by explicit communication by programmer code.



int main(int argc, char **argv) {

// Initialization

mpi_module1(data)

mpi_module1(data) {

charm module1(da

// do work

_ _ _ _ _ _ _ _ _ _ _ _

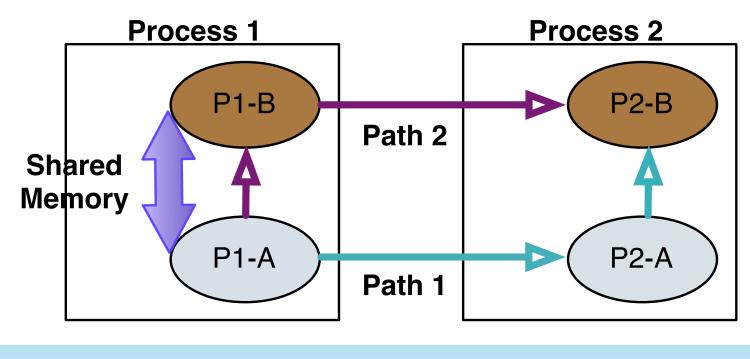
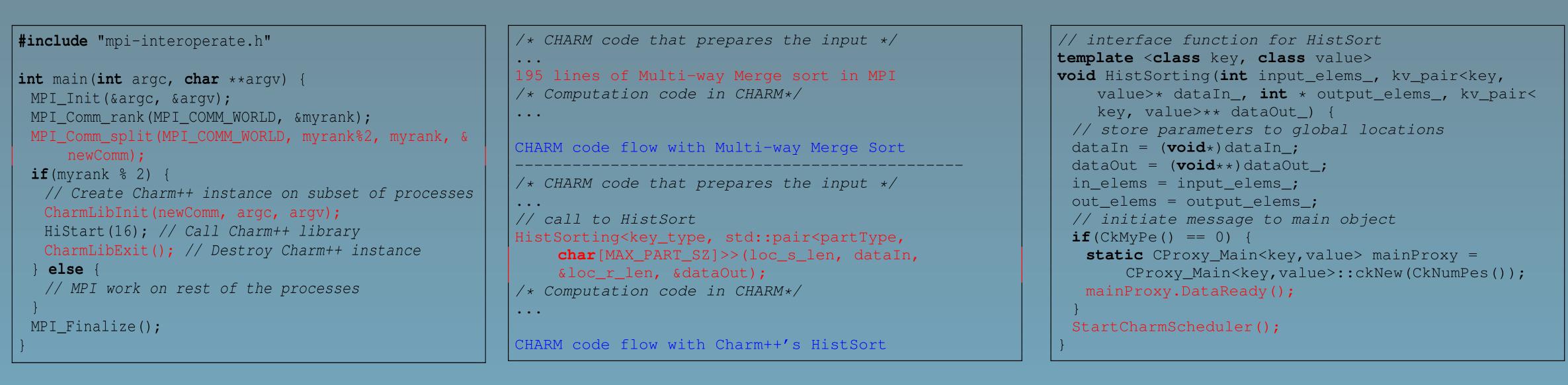
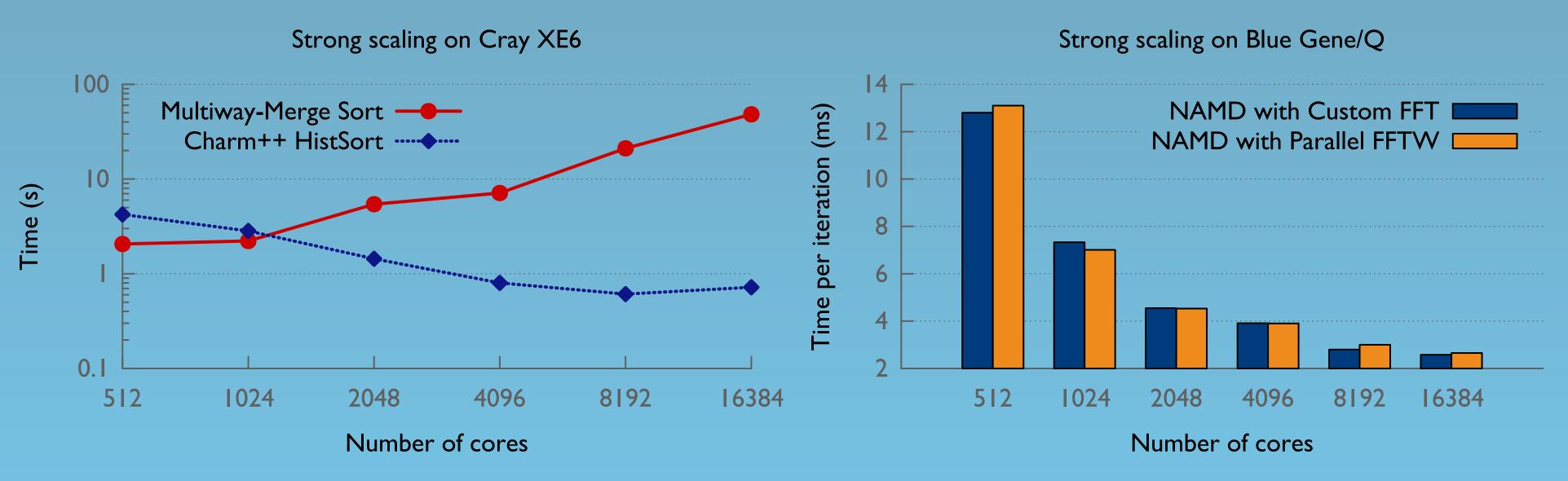


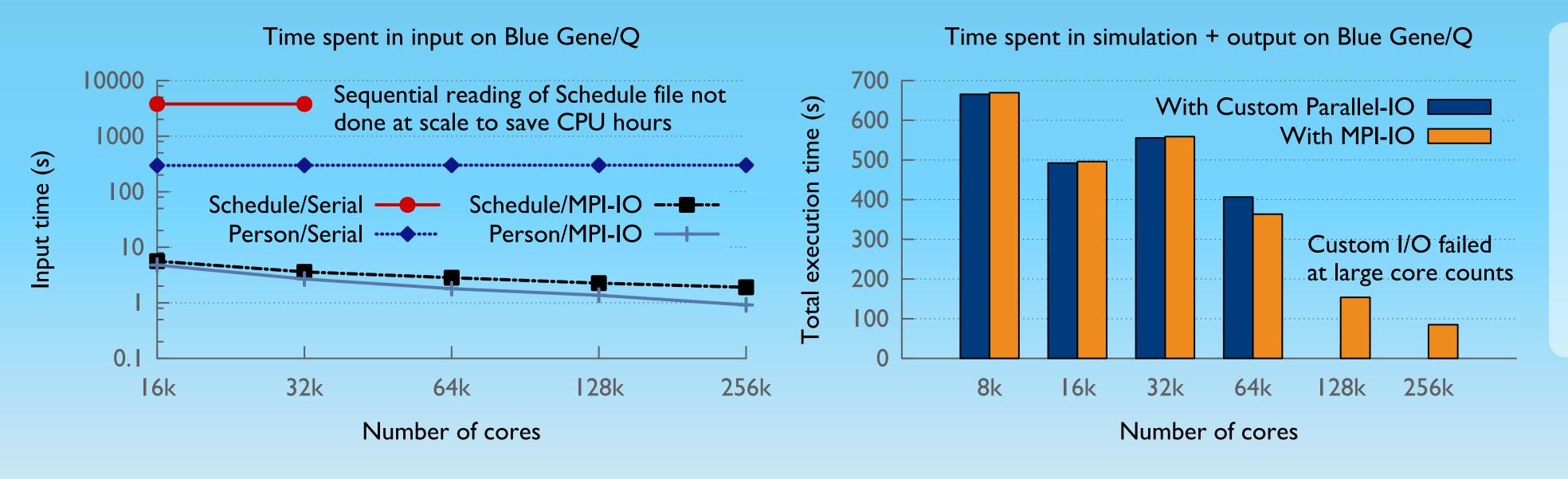
Figure 3







MPI-based Chombo using Charm++-based HistSorting with *time*division: the performance of global sorting, which is the scaling bottleneck, is improved by 48x with minimal changes to the code (as shown in the figures at the top)



Examples and Results

Left figure: basic MPI example to interoperate with Charm++; lines in red are the only additional code required. Centre and right figures: all changes required for calling HistSorting library in Charm++ from a production MPI code.s

> Charm++-based NAMD using MPI-based FFTW with space-division: similar performance obtained without the need to maintain a custom FFT library in Charm++.

Charm++-based EpiSimdemics using MPI-IO with *hybrid-division*: scaling bottleneck due to file reading is eliminated, while file output is enabled at scale.