

Digging into MPI Performance with Caliper

Riley Shipley and Garrett Hooten



Center for Understandable, Performant Exascale Communication Systems

 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

Tools

- ExaMPI
 - UTC's C++ research-based implementation of MPI
 - Features strong progress
 - Separate progress thread in background making sure progress is always being made

Tools

- Caliper
 - LLNL program instrumentation and performance measurement library
 - Can be added to application and benchmark codes with zero overhead until activation
 - Designed for HPC applications, works for C/C++/Fortran programs on Linux/Unix

event.begin#annotation	event.end#annotation	pthread.id	time.duration	time.inclusive
BlockingProgress::progress_algorithms		35184438341872	0.000015	
	BlockingProgress::progress_algorithms	35184438341872	0.000012	0.000012
BlockingProgress::progress_transports		35184438341872	0.000017	
MPI_Send		35184372393264	2.808913	
Toolbox::generic_request		35184372393264	0.000027	
Request::set_envelope_size()		35184372393264	0.000031	
	Request::set_envelope_size()	35184372393264	0.000021	0.000021
Toolbox::start_wait_complete		35184372393264	0.000029	
BasicComm::start_request		35184372393264	0.000021	
		35184372393264	0.000019	0.000075
	BasicComm::start_request	35184372393264	0.000017	0.000115
BlockingProgress::wait_for_complete		35184372393264	0.000021	

Tools

- Hatchet
 - Python-based tool that specializes in visualizing hierarchical data (calling context trees, call graphs, nested timers, etc...)
 - Loads Caliper data as a Pandas DataFrame



Regions in Caliper

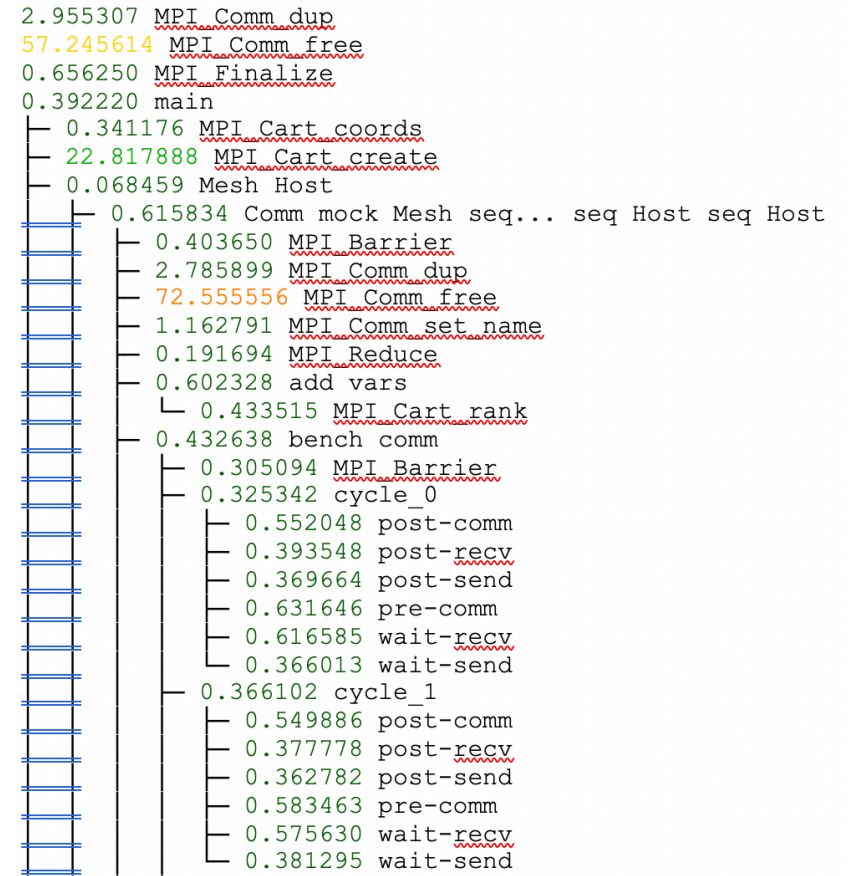
- The basis of Caliper profiling and ExaMPI message tracing
- Regions are essentially portions of code that you want timed
- Can mark an entire function for timing
 - Ex: CALI_CXX_MARK_FUNCTION
- Can mark specific portions of code
 - Ex: CALI_MARK_BEGIN("my_region") ... CALI_MARK_END("my_region")

Comparing IBM Spectrum MPI and ExaMPI

- Summer 2022 Project -- Comparing Spectrum MPI performance vs ExaMPI in Comb (LLNL's HPC benchmark)
- Used Caliper's built-in *profile.mpi* configuration to automatically record time spent in MPI functions
- Ran tests on LLNL's Lassen supercomputer

Comparing IBM Spectrum MPI and ExaMPI

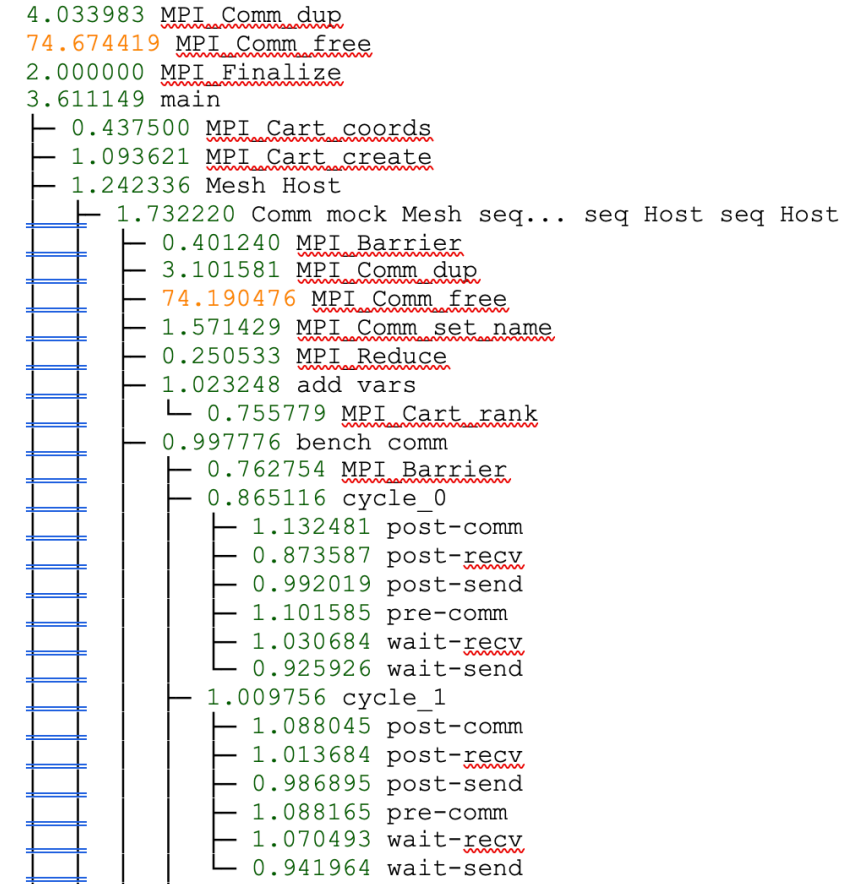
- Initial results were disappointing
 - Comb with ExaMPI was slower than Spectrum, even when there was no MPI happening
 - Hatchet calculation (right) shows Spectrum was usually significantly faster (values < 1)



Ratio of Spectrum completion times to ExaMPI completion times in our initial Comb assessment

Comparing IBM Spectrum MPI and ExaMPI

- Narrowed down the problem to 2 areas:
 - Build/Compile time
 - Re-built tests to ensure they were the same
 - Running them again produced the same results
 - Launch/Runtime
 - Suspected that runtime configurations did not match
 - Found process/core affinities between ExaMPI and Spectrum were different
- New core affinity code added
- Updated comparison shows ExaMPI and Spectrum performance are near equal (value ≈ 1)



Ratio of ExaMPI completion times to Spectrum times in our updated Comb assessment

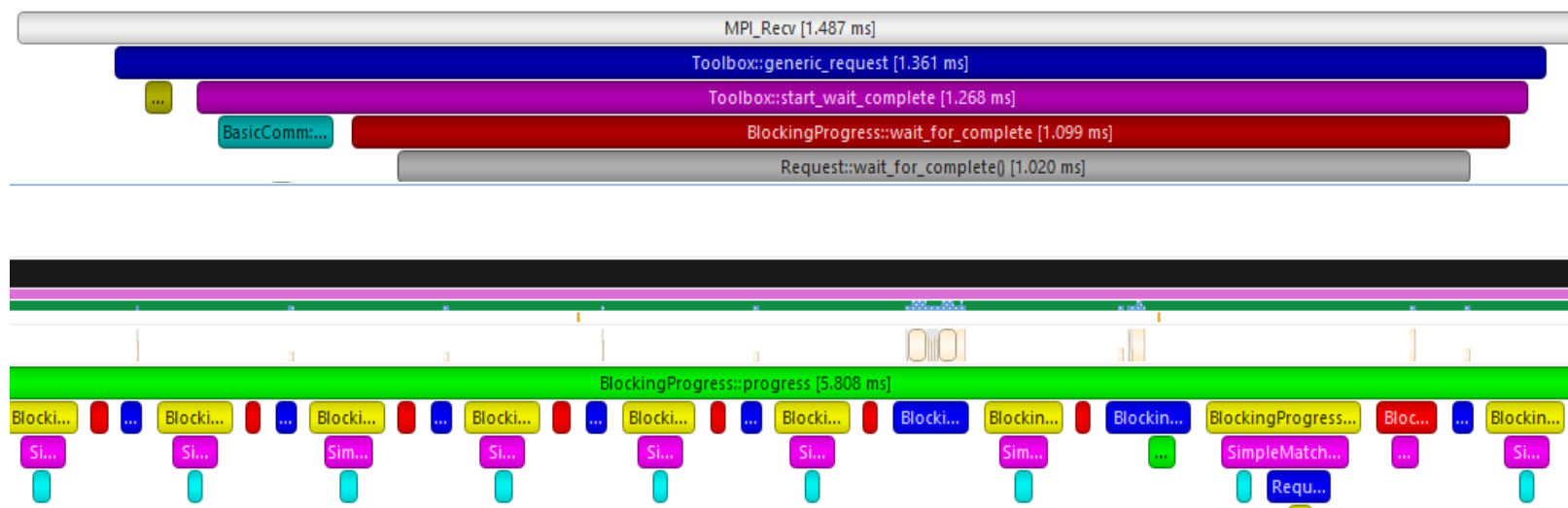


Multiple MPI Profiles

- External
 - Adding Caliper to a benchmark to compare MPI implementations
 - MPI is linked into Caliper at build time
 - Caliper automatically times each MPI binding
- Internal
 - Caliper is built without linking MPI
 - Regions are added *inside* of ExaMPI
 - ExaMPI is then built with Caliper
 - Creates a statistical profile of each part of ExaMPI

Visualization & Analysis

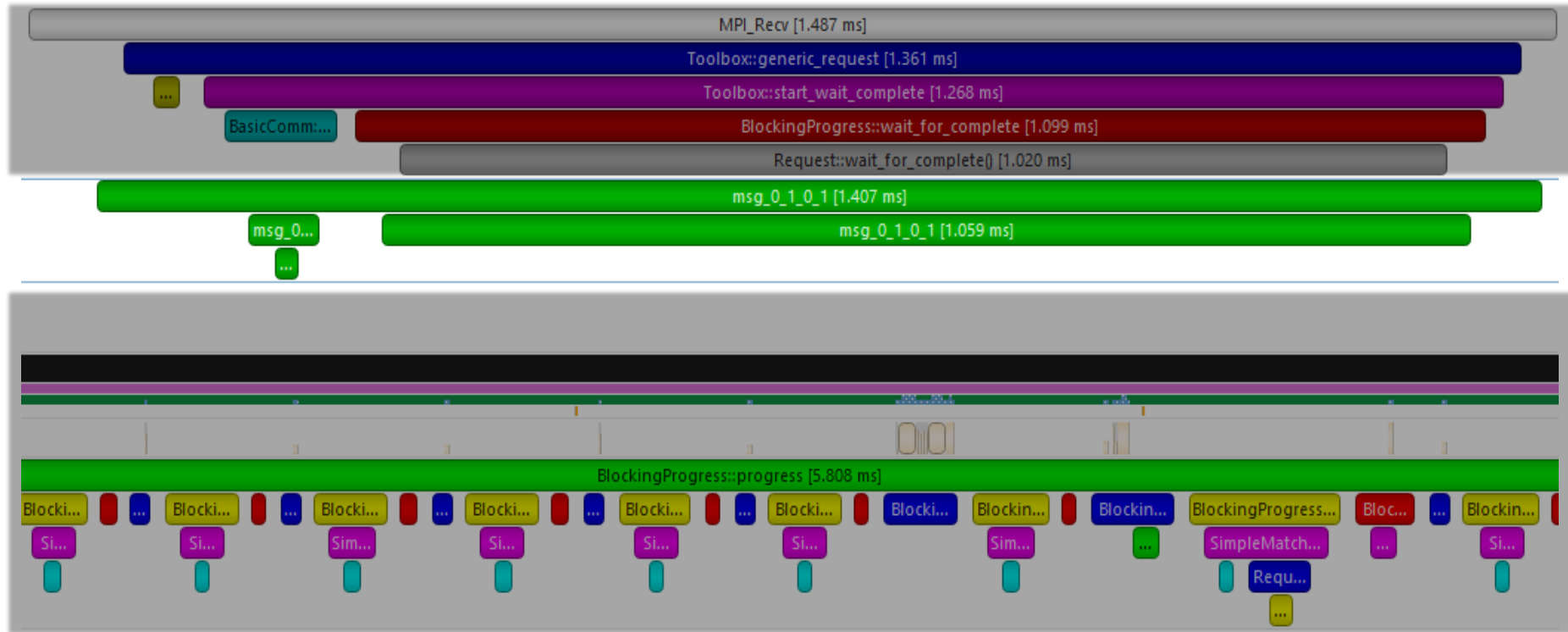
- NVIDIA NSight
 - Caliper can export regions to NVTX
 - Requires Caliper to be built with CUDA
 - Allows us to see exact start/end of each part of ExaMPI



Message Tracing

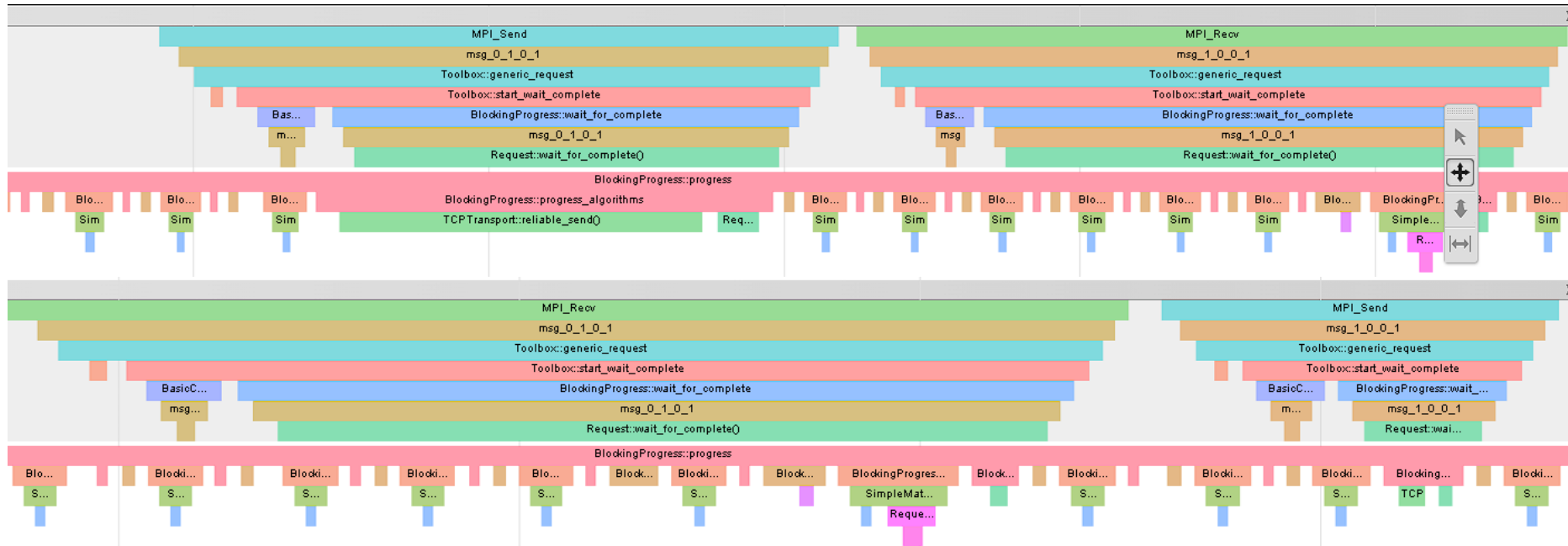
- Goal is to measure end-to-end communication times
 - Use existing regions as **annotation** regions
 - Add new **message** regions when using message data
 - Message regions are stored separate from annotation regions
 - Allows us to find a message at any point, including across ranks
- Message identification
 - Message code – send process rank, destination process rank, tag
 - Each thread of each process rank stores message codes
 - Code + count = unique ID

NSight & Message Tracing



Combining multiple ranks

- JSON trace visualizer (://tracing)
 - Caliper can combine traces from multiple ranks into one file
 - Result can be viewed in Chrome-based browsers



Future Work

- Timeline Synchronization
 - Caliper does not account for different rank start times
 - Aligning message regions should show full message timelines
- Further Analysis
 - Hatchet becomes more complex to use with custom regions
 - We are continuing to refine our profiling models by working with LLNL
- Continuing to use profiling results to make improvements to ExaMPI