Enabling Graph-Based Profiling Analysis using Hatchet

Student: Ian Lumsden

Mentors: Stephanie Brink, Michael R. Wyatt II, Todd Gamblin, and Michela Taufer



LLNL-PRES-815496

This work was performed under the auspices of the U.S. Department of Energy by Lawrence Livermore National Laboratory under contract DE-AC52-07NA27344. Lawrence Livermore National Security, LLC



[2] L. Adhianto, S. Banerjee, M. Fagan, M. Krentel, G. Marin, J. Mellor-Crummey, and N. Tallent. 2010. HPCToolkit: Tools for Performance Analysis of Optimized Parallel Programs. *Concurrency and Computation: Practice and Experience* 22, 6 (2010), 285-701

2





3





[2] L. Adhianto, S. Banerjee, M. Fagan, M. Krentel, G. Marin, J. Mellor-Crummey, and N. Tallent. 2010. HPCToolkit: Tools for Performance Analysis of Optimized Parallel Programs. *Concurrency and Computation: Practice and Experience* 22, 6 (2010), 285-701

4

TEN



[2] L. Adhianto, S. Banerjee, M. Fagan, M. Krentel, G. Marin, J. Mellor-Crummey, and N. Tallent. 2010. HPCToolkit: Tools for Performance Analysis of Optimized Parallel Programs. *Concurrency and Computation: Practice and Experience* 22, 6 (2010), 285-701

5



[2] L. Adhianto, S. Banerjee, M. Fagan, M. Krentel, G. Marin, J. Mellor-Crummey, and N. Tallent. 2010. HPCToolkit: Tools for Performance Analysis of Optimized Parallel Programs. *Concurrency and Computation: Practice and Experience* 22, 6 (2010), 285-701

6



- Profiler output uses a custom, unique data format
- Visualization and analysis tools are GUI-based with limited analysis capabilities
- Both issues above limit user analysis options



Hatchet



- Python-based library
- A new, **general** data analysis tool that can read HPC profiling data from **different profilers**
 - Store the raw performance data into a *pandas* DataFrame
 - Represent the relational *caller-callee* data with directed acyclic graph (DAG)



Hatchet



- Python-based library
- A new, **general** data analysis tool that can read HPC profiling data from **different profilers**
 - Store the raw performance data into a *pandas* DataFrame
 - Represent the relational *caller-callee* data with directed acyclic graph (DAG)





HPCToolkit with Hatchet



[2] L. Adhianto, S. Banerjee, M. Fagan, M. Krentel, G. Marin, J. Mellor-Crummey, and N. Tallent. 2010. HPCToolkit: Tools for Performance Analysis of Optimized Parallel Programs. *Concurrency and Computation: Practice and Experience* 22, 6 (2010), 285-701

10

HPCToolkit with Hatchet



Hatchet



- Python-based library
- A new, **general** data analysis tool that can read HPC profiling data from **different profilers**
 - Store the raw performance data into a *pandas* DataFrame
 - Represent the relational *caller-callee* data with directed acyclic graph (DAG)
- Hatchet restricts users to table-based analysis of the raw performance data
- Hatchet does NOT support analysis of the relational data collected by HPC profilers

[1] A. Bhatele, S. Brink, and T. Gamblin. 2019. Hatchet: Pruning the Overgrowth in Parallel Profiles. In *Proc. of the International Conference of High Performance Computing, Networking, Storage, and Analysis (SC19).*



Research Goals

- Design new graph-based filtering query language to enable the use of relational data collected by profilers in analysis
- Implement this query language and integrate it into Hatchet analysis
- Use the augmented Hatchet to analyze the performance of different MPI calls in HPC benchmark applications



Research Goals

- Design new graph-based filtering query language to enable the use of relational data collected by profilers in analysis
- Implement this query language and integrate it into Hatchet analysis
- Use the augmented Hatchet to analyze the performance of different MPI calls in HPC benchmark applications



• Filter nodes in a graph using a user-provided path pattern (called a *query path*)



- Filter nodes in a graph using a user-provided path pattern (called a *query path*)
- User Input:
 - A *Query Path* represented as a list of *abstract graph nodes*



- Filter nodes in a graph using a user-provided path pattern (called a *query path*)
- User Input:
 - A *Query Path* represented as a list of *abstract graph nodes*
 - Two parts to an *abstract graph node*:



- Filter nodes in a graph using a user-provided path pattern (called a *query path*)
- User Input:
 - A *Query Path* represented as a list of *abstract graph nodes*
 - Two parts to an *abstract graph node*:
 - 1. A wildcard specifying the number of real graph nodes to match to the *abstract* graph node (default is 1)

- Filter nodes in a graph using a user-provided path pattern (called a *query path*)
- User Input:
 - A *Query Path* represented as a list of *abstract graph nodes*
 - Two parts to an *abstract graph node*:
 - 1. A wildcard specifying the number of real graph nodes to match to the *abstract* graph node (default is 1)
 - 2. A filter used to determine whether a real graph node matches the *abstract* graph node (default is an "always true" filter)



- Algorithm:
 - 1. Read and parse the user's query path
 - 2. Match real nodes in the graph being filtered to the abstract nodes in the query path
 - 3. Collect all graph paths that match the full query path
 - 4. Create a new graph (GraphFrame) containing only the nodes in the matched paths

Research Goals

- Design new graph-based filtering query language to enable the use of relational data collected by profilers in analysis
- Implement this query language and integrate it into Hatchet analysis
- Use the augmented Hatchet to analyze the performance of different MPI calls in HPC benchmark applications

Hatchet Query Language: APIs



Hatchet Query Language: Example





Hatchet Query Language: Example





Hatchet Query Language: Example





Publicly available starting with Hatchet v1.2.0



Research Goals

- Design new graph-based filtering query language to enable the use of relational data collected by profilers in analysis
- Implement this query language and integrate it into Hatchet analysis
- Use the augmented Hatchet to analyze the performance of different MPI calls in HPC benchmark applications



Case Study: Collect Data

- Run each of the following benchmarks with MVAPICH (M) and Spectrum-MPI (S) using 64, 128, 256, and 512 MPI ranks, and profile the runs with HPCToolkit
 - AMG2013
 - Kripke
 - Lammps
- Load the generated profiles into Hatchet



Case Study: Collect Data

- Perform the benchmarking on LLNL's Lassen supercomputer
 - LLNL CZ (Public Collaboration) CORAL Supercomputer
 - 795 AC922 Nodes
 - 2 IBM Power9 CPUs per node (20 cores per node)
 - 256 GB Memory per node
 - Infiniband EDR interconnect







Case Study: Extract MPI Layer

- Filter the call graphs to get subgraphs rooted at standard MPI function calls
- Query Path used to extract MPI Layer:



Case Study: Calculate Percent MPI Time for each MPI Function



AMG2013

Only MPI functions that contribute at least 5% of total MPI time are shown

| MPI Function Calls | | | | | | |
|--------------------|--|---------------|--|--|--|--|
| MPI_Finalize | | MPI_Allreduce | | | | |
| MPI_Allgather | | MPI_Waitall | | | | |
| Remaining MPI Time | | | | | | |



Case Study: Calculate Percent MPI Time for each MPI Function



AMG2013

Only MPI functions that contribute at least 5% of total MPI time are shown

| MPI Function Calls | | | | | | |
|--------------------|--|---------------|--|--|--|--|
| MPI_Finalize | | MPI_Allreduce | | | | |
| MPI_Allgather | | MPI_Waitall | | | | |
| Remaining MPI Time | | | | | | |



Case Study: Calculate Percent MPI Time for each MPI Function



AMG2013

Only MPI functions that contribute at least 5% of total MPI time are shown

| MPI Function Calls | | | | | |
|--------------------|--|---------------|--|--|--|
| MPI_Finalize | | MPI_Allreduce | | | |
| MPI_Allgather | | MPI_Waitall | | | |
| Remaining MPI Time | | | | | |



Case Study: Calculate Percent MPI Time for Child Calls of MPI Functions



AMG2013

Only children functions that contribute at least 10% of total MPI time are shown

| Child Function Calls | | | | |
|-----------------------------|--|-----------------------------|--|--|
| <unknown file=""></unknown> | | <unknown file=""></unknown> | | |
| pthread_spin_lock.c:26 | | memset.S:1133 | | |
| malloc.c:0 | | Geometry.h:0 | | |
| stl_vector.h:0 | | Remaining MPI Time | | |



Case Study: Calculate Percent MPI Time for Child Calls of MPI Functions



AMG2013

Only children functions that contribute at least 10% of total MPI time are shown

| Child Function Calls | | | | |
|-----------------------------|--|-----------------------------|--|--|
| <unknown file=""></unknown> | | <unknown file=""></unknown> | | |
| pthread_spin_lock.c:26 | | memset.S:1133 | | |
| malloc.c:0 | | Geometry.h:0 | | |
| stl_vector.h:0 | | Remaining MPI Time | | |

Number of MPI Ranks

Case Study: Calculate Percent MPI Time for Child Calls of MPI Functions



AMG2013

Only children functions that contribute at least 10% of total MPI time are shown

| Child Function Calls | | | | | |
|--|--|---|--|--|--|
| <unknown file=""> [libmlx5.so.1.0.0]:0</unknown> | | <unknown file=""> [libmlx5.so.1.0.0]:1133</unknown> | | | |
| pthread_spin_lock.c:26 | | memset.S:1133 | | | |
| malloc.c:0 | | Geometry.h:0 | | | |
| stl_vector.h:0 | | Remaining MPI Time | | | |

Zoom into specific benchmarks (**AMG2013**) and examine the children of specific MPI functions (**MPI_Allgather**)

| MPI Function Calls | | | | | |
|--------------------|--------------------|--|---------------|--|--|
| | MPI_Finalize | | MPI_Allreduce | | |
| | MPI_Allgather | | MPI_Waitall | | |
| | Remaining MPI Time | | | | |



Children Functions



Number of MPI Ranks

| Child Function Calls | | | | | | |
|--|--|--------------------|--|--|--------------|--|
| <unknown file=""> [libmlx5.so.1.0.0]:0</unknown> | libmlx5.so.1.0.0]:0 <unknown file=""> Geomet [libmlx5.so.1.0.0]:1133</unknown> | | | | Geometry.h:0 | |
| pthread_spin_lock.c:26 | | memset.S:1133 | | | malloc.c:0 | |
| stl_vector.h:0 | | Remaining MPI Time | | | Time | |



Zoom into specific benchmarks (AMG2013) and examine the children of specific MPI functions (MPI_Allgather)

| MPI Function Calls | | | | | |
|--------------------|--------------------|--|---------------|--|--|
| | MPI_Finalize | | MPI_Allreduce | | |
| | MPI_Allgather | | MPI_Waitall | | |
| | Remaining MPI Time | | | | |



Number of MPI Ranks

| Child Function Calls | | | | | | | |
|---|--|--|--------------------|--|------------|--|--|
| <unknown file=""> [libmlx5.so.1.0.0]:0 <unknown file=""> [libmlx5.so.1.0.0]:1133 Geometry.h:0</unknown></unknown> | | | | | | | |
| pthread_spin_lock.c:26 | | | memset.S:1133 | | malloc.c:0 | | |
| stl vector.h:0 | | | Remaining MPI Time | | | | |



128

Number of MPI Ranks

64

256

512

Zoom into specific benchmarks (**AMG2013**) and examine the children of specific MPI functions (**MPI_Allgather**)

| MPI Function Calls | | | | | |
|--------------------|--------------------|--|---------------|--|--|
| | MPI_Finalize | | MPI_Allreduce | | |
| | MPI_Allgather | | MPI_Waitall | | |
| | Remaining MPI Time | | | | |



Number of MPI Ranks

| 11.93 s | 4.05 s | 12.15 s | 4.57 s | 13.58 s | 4.67 s | 0011 |
|---------|--------|---------|--------|---------|--------|------|
| | | | | | | |

Children Functions



100

Number of MPI Ranks

| Child Function Calls | | | | | | |
|--|--|---|--------------------|--|--------------|--|
| <unknown file=""> [libmlx5.so.1.0.0]:0</unknown> | | <unknown file=""> [libmlx5.so.1.0.0]:1133</unknown> | | | Geometry.h:0 | |
| pthread_spin_lock.c:26 | | | memset.S:1133 | | malloc.c:0 | |
| stl_vector.h:0 | | | Remaining MPI Time | | | |



Zoom into specific benchmarks (**AMG2013**) and examine the children of specific MPI functions (**MPI_Allgather**)



AMG Children Functions

Allgather Children Functions



| Child Function Calls | | | | | | | |
|----------------------|--|--|-----------------------------|---------------------|--|--------------|--|
| | <unknown file=""> [libmlx5.so.1.0.0]:0</unknown> | | <unknown file=""></unknown> | | | Geometry.h:0 | |
| | | | [lib | mlx5.so.1.0.0]:1133 | | | |
| | pthread_spin_lock.c:26 | | memset.S:1133 | | | malloc.c:0 | |
| | stl_vector.h:0 | | | Remaining MPI Time | | | |



Lessons Learned: MPI

The

pthread_spin_lock function is consistently a major contributor to MPI runtime (i.e., 20% or more of MPI time)



AMG Children Functions

Allgather Children Functions



| Child Function Calls | | | | | | | |
|----------------------|--|--|-----------------------------|---------------------|--|--------------|--|
| | <unknown file=""> [libmlx5.so.1.0.0]:0</unknown> | | <unknown file=""></unknown> | | | Geometry.h:0 | |
| | | | [lib | mlx5.so.1.0.0]:1133 | | | |
| | pthread_spin_lock.c:26 | | memset.S:1133 | | | malloc.c:0 | |
| | stl_vector.h:0 | | | Remaining MPI Time | | Time | |



Lessons Learned: MPI

In AMG2013, the inferior performance of MPI_Allgather in Spectrum-MPI can possibly be attributed to the use of pthread_spin_lock



AMG Children Functions

Allgather Children Functions



| Child Function Calls | | | | | | |
|--|----------------|---|--------------------|--|--------------|--|
| <unknown file=""> [libmlx5.so.1.0.0]:0</unknown> | | <unknown file=""> [libmlx5.so.1.0.0]:1133</unknown> | | | Geometry.h:0 | |
| pthread_spin_lock.c:26 | | | memset.S:1133 | | malloc.c:0 | |
| stl_vector.h:0 | stl_vector.h:0 | | Remaining MPI Time | | Time | |



Lessons Learned: Hatchet

- Using the query language and Hatchet, we are able to:
 - Extract all call paths specific to a given library
 - Determine the performance contributions of function calls used by these libraries
 - Correlate children function calls to specific important library API calls in an application
 - Use this correlation to determine children function calls that contribute the most to the performance of the targeted library API call
 - Compare the correlation of children and API calls across libraries to determine possible causes for performance differences



Future Work

- Use and expand on the techniques and tools from this work to study large scale scalability and replicability problems
 - Use information gained from studying these problems to propose mitigation strategies
 - Develop tools to help others study and mitigate these problems



Additional Resources on Hatchet

Hatchet SC19 Paper



Hatchet SC20 Paper



Hatchet GitHub







Disclaimer

This document was prepared as an account of work sponsored by an agency of the United States government. Neither the United States government nor Lawrence Livermore National Security, LLC, nor any of their employees makes any warranty, expressed or implied, or assumes any legal liability or responsibility for the accuracy, completeness, or usefulness of any information, apparatus, product, or process disclosed, or represents that its use would not infringe privately owned rights. Reference herein to any specific commercial product, process, or service by trade name, trademark, manufacturer, or otherwise does not necessarily constitute or imply its endorsement, recommendation, or favoring by the United States government or Lawrence Livermore National Security, LLC. The views and opinions of authors expressed herein do not necessarily state or reflect those of the United States government or Lawrence Livermore National Security, LLC, and shall not be used for advertising or product endorsement purposes.

Supplemental Slides



Case Study: Benchmarks

- Kripke (LLNL) [5]
 - Mini-App Benchmark
 - Deterministic
 - Runs S_N Neutron Transport simulation



Case Study: Benchmarks

- Lammps (SNL) [6]
 - Popular classical molecular dynamics simulator
 - Also contains a set of benchmarks



Case Study: Benchmarks

- AMG2013 (LLNL) [7]
 - Parallel algebraic multigrid solver for linear systems
 - Proxy App for hypre's BoomerAMG solver
 - Highly synchronous code

Case Study: MPIs

- MVAPICH (v2.3.3)
 - Project led by the Network-Based Computing Laboratory at Ohio State University
 - Based on MPICH (Argonne National Lab)
 - Conforms with MPI 3.1 standard



Case Study: MPIs

- Spectrum-MPI (v10.2)
 - Created by IBM
 - Based on OpenMPI (first created by LANL, Indiana University, and the University of Tennessee)
 - Conforms with MPI 3.1 standard

