

Unstructured Communication in AMR Applications

Justin “Tanner” Broaddus

This work was supported in part by the U.S. Department of Energy's National Nuclear Security Administration (NNSA) under the Predictive Science Academic Alliance Program (PSAAP-III), Award #DE-NA0003966.



Center for Understandable
Performant Exascale
Communication Systems

 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

Internship Overview



LANL Internship Goal - Establish an understanding of unstructured communication patterns present in AMR applications.

Under the guidance of Galen Shipman - LANL Computer Science

Codebase of interest - LANL's Parthenon

Research questions:

1. How do developers vary in their implementation of unstructured communication in order to perform adaptive mesh refinement (AMR)?
2. What classification of unstructured communication can be derived from inspected AMR codebases?



Center for Understandable
Performant Exascale
Communication Systems

 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

Parthenon

- AMR Infrastructure developed at Los Alamos National Laboratory (LANL), Princeton University, Michigan State University
- Performance portable task-based AMR infrastructure
- Implemented in C++
- Uses Kokkos as the shared memory parallel programming model



Center for Understandable
Performant Exascale
Communication Systems

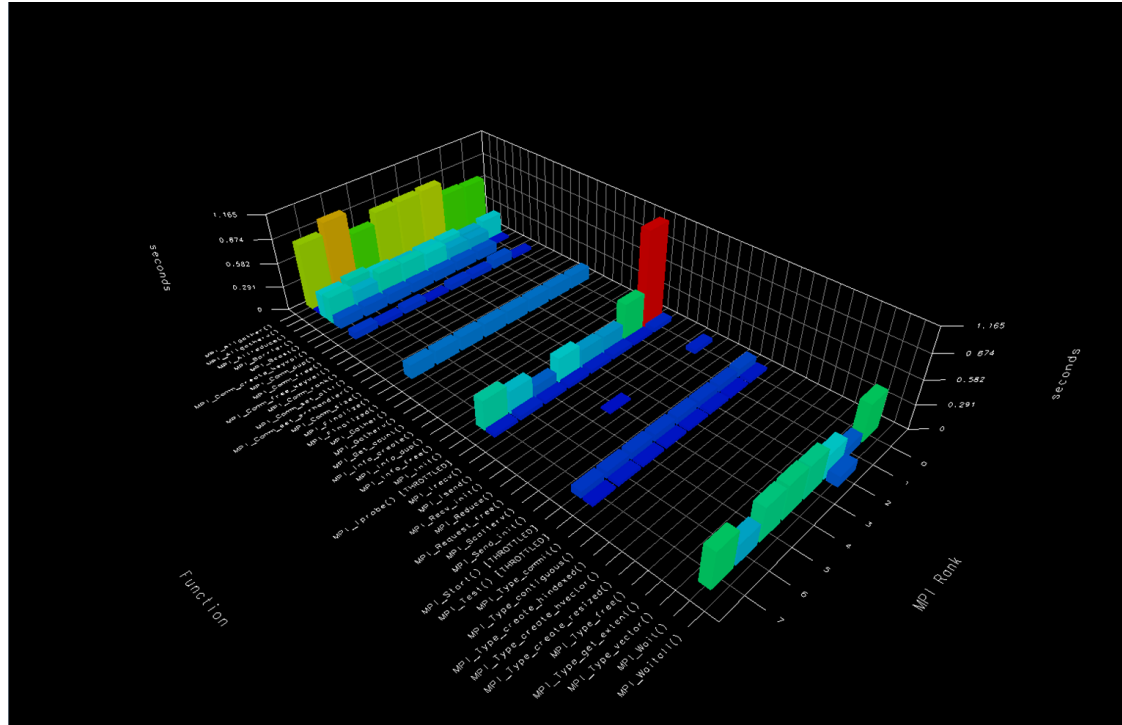
 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

Parthenon - Communication Overview

- Halo exchanges
 - Uses persistent communication; initialized by **MPI_Send_Init** and **MPI_Recv_Init**
 - Followed by **MPI_Start** and **MPI_Test**
- Load balancing
 - Considered periodically, load costs communicated by **MPI_Allgather** call
- Mesh block transfers
 - **MPI_Isend** and **MPI_Irecv**; followed by **MPI_Wait**
- Tools for analysis of communication
 - Tau's Paraprof and Jumpshot
 - Caliper and Kokkos annotations



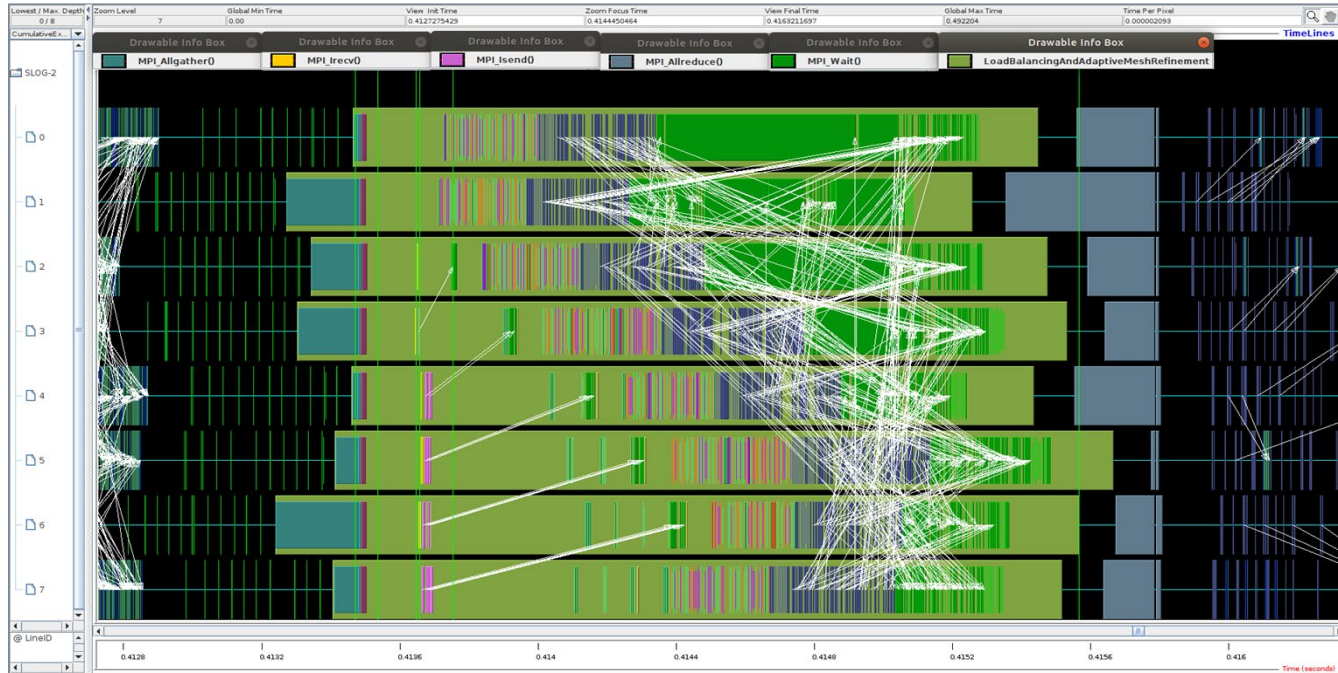
Parthenon - MPI Profiling



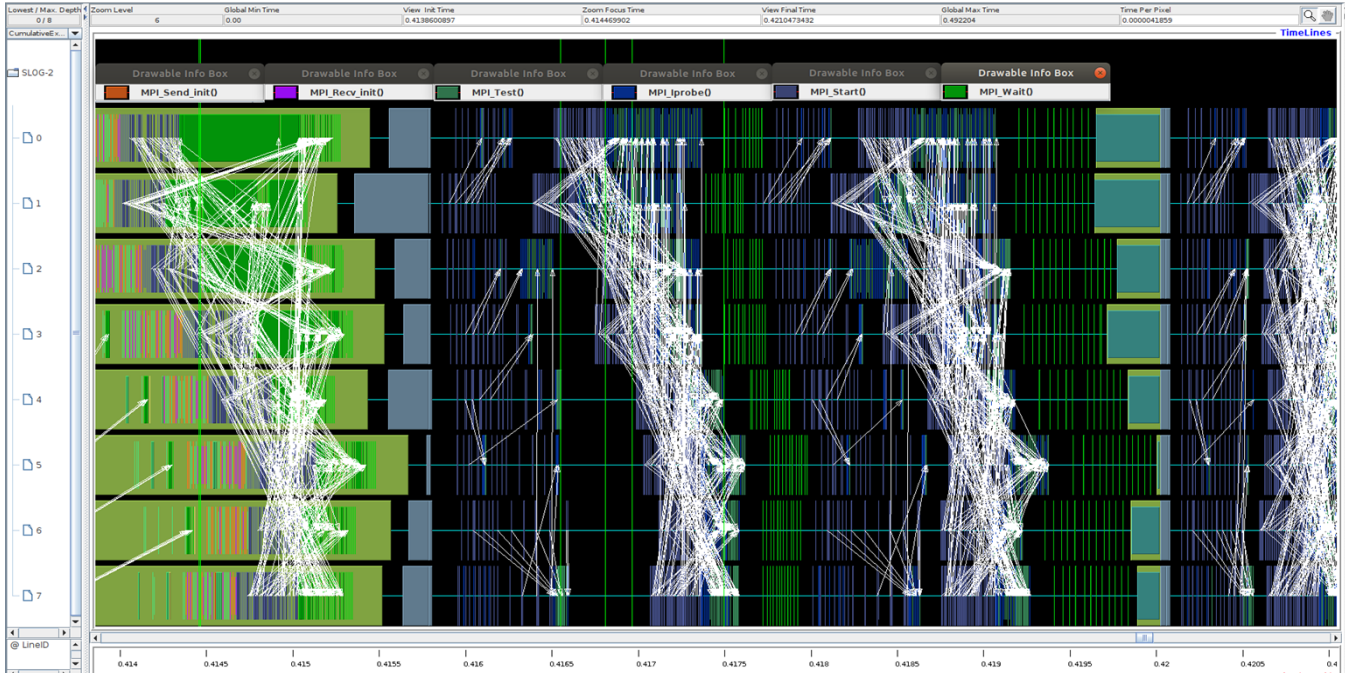
Center for Understandable
Performant Exascale
Communication Systems

 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

Parthenon - Load Balancing Trace



Parthenon - Halo Exchange Trace



Center for Understandable
Performant Exascale
Communication Systems

THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

CLAMR Overview

- Collection of cell-based AMR mini-apps developed at LANL
- Tests algorithms to be used in heterogeneous computing environments
- OpenCL as the shared memory parallel programming model



Center for Understandable
Performant Exascale
Communication Systems

 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

CLAMR Observations

- Difference between the min and max times for neighbor boundary updates in MVAPICH2 compiled binaries
- MPI processes stuck in **MPI_Waitall**
- **MPI_Barrier** after **MPI_Waitall** call improves performance for MVAPICH2 compiled binaries
- Issue present L7 function

```

CPU: state_timer_refine_potential 2.2446 4.2383 4.2379 s min/median/max
CPU: state_timer_calc_mpot 1.1493 2.9786 3.0291 s min/median/max
CPU: mesh_timer_refine_smooth 1.1085 1.1493 1.2259 s min/median/max
CPU: state_timer_rezone_all 3.9166 3.9117 3.9117 s min/median/max
CPU: mesh_timer_partition 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_calc_neighbors 33.2815 33.2977 33.3073 s min/median/max
CPU: mesh_timer_hash_setup 7.1314 7.3969 7.3570 s min/median/max
CPU: mesh_timer_hash_query 3.9666 4.1236 4.3581 s min/median/max
CPU: mesh_timer_find_boundary 1.4795 1.8370 2.2330 s min/median/max
CPU: mesh_timer_push_setup 0.6429 0.6889 0.6889 s min/median/max
CPU: mesh_timer_push_boundary 0.4611 1.0801 9.3026 s min/median/max
CPU: mesh_timer_locat_list 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_layer1 0.5561 2.1061 3.2465 s min/median/max
CPU: mesh_timer_layer2 0.3581 1.7464 3.0779 s min/median/max
CPU: mesh_timer_layer_list 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_copy_mesh_data 0.3132 0.3458 0.4462 s min/median/max
CPU: mesh_timer_fill_mesh_ghost 0.0050 0.0068 0.0068 s min/median/max
CPU: mesh_timer_fill_neigh_ghost 0.5816 3.3174 3.5235 s min/median/max
CPU: mesh_timer_set_corner_neigh 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_neigh_adjust 0.2602 0.2641 0.2661 s min/median/max
CPU: mesh_timer_setup_comm 0.7076 0.7836 15.8901 s min/median/max
CPU: state_timer_mess_sum 0.0094 0.0295 0.0330 s min/median/max
CPU: mesh_timer_load_balance 0.9030 9.3173 9.7593 s min/median/max
CPU: mesh_timer_calc_spatial_coordi 0.0000 0.0000 0.0000 s min/median/max
-----
Profiling: Total CPU time was 62.8369 62.8547 62.8919 s min/median/max
-----
Mesh Ops (Neigh+rezone+smooth+balance) 47.2573 47.6461 48.1677 s min/median/max
Mesh Ops Percentage 75.1697 75.8151 76.5403 percent min/median/max
-----
Profiling: Total time was 64.0156 64.0156 64.0157 s min/median/max

```

CLAMR compiled w/ MVAPICH2

```

CPU: state_timer_rezone_all 3.9078 3.9091 3.9098 s min/median/max
CPU: mesh_timer_partition 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_calc_neighbors 25.8390 25.8526 25.8639 s min/median/max
CPU: mesh_timer_hash_setup 7.1589 7.5878 7.9701 s min/median/max
CPU: mesh_timer_hash_query 4.0785 4.1575 4.4782 s min/median/max
CPU: mesh_timer_find_boundary 1.4645 1.8823 2.2533 s min/median/max
CPU: mesh_timer_push_setup 0.0502 0.0070 1.0069 s min/median/max
CPU: mesh_timer_push_boundary 0.2565 0.4076 0.5316 s min/median/max
CPU: mesh_timer_locat_list 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_layer1 0.6030 2.2101 3.4291 s min/median/max
CPU: mesh_timer_layer2 0.3595 2.0652 3.3769 s min/median/max
CPU: mesh_timer_layer_list 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_copy_mesh_data 0.3334 0.3702 0.4889 s min/median/max
CPU: mesh_timer_fill_mesh_ghost 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_fill_neigh_ghost 0.5799 3.3697 3.4734 s min/median/max
CPU: mesh_timer_set_corner_neigh 0.0000 0.0000 0.0000 s min/median/max
CPU: mesh_timer_neigh_adjust 0.2611 0.2660 0.2724 s min/median/max
CPU: mesh_timer_setup_comm 0.1128 2.3837 8.4010 s min/median/max
CPU: state_timer_mess_sum 0.0042 0.0341 0.0342 s min/median/max
CPU: mesh_timer_load_balance 9.0555 9.3503 9.7628 s min/median/max
CPU: mesh_timer_calc_spatial_coordi 0.0000 0.0000 0.0000 s min/median/max
-----
Profiling: Total CPU time was 55.5333 55.5521 55.5745 s min/median/max
-----
Mesh Ops (Neigh+rezone+smooth+balance) 39.9882 40.2425 40.6864 s min/median/max
Mesh Ops Percentage 71.9815 72.4353 73.2280 percent min/median/max

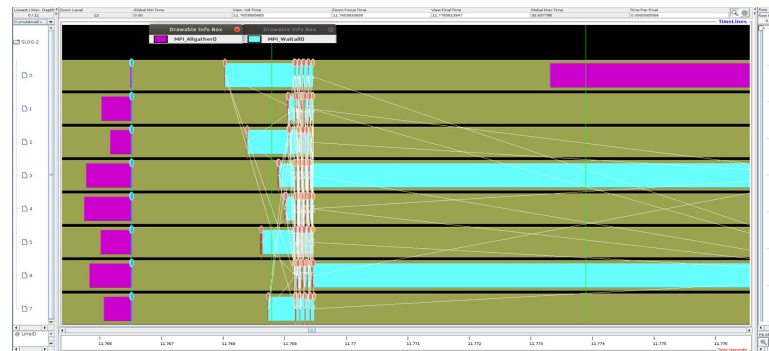
```

CLAMR compiled w/ OpenMPI

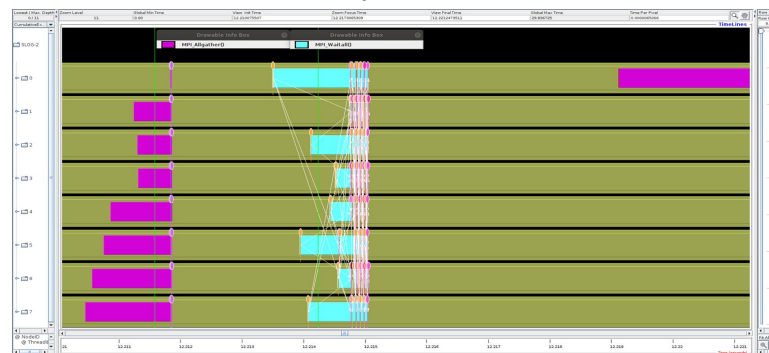


CLAMR Observations

- Difference between the min and max times for neighbor boundary updates in MVAPICH2 compiled binaries
- MPI processes stuck in **MPI_Waitall**
- **MPI_Barrier** after **MPI_Waitall** call improves performance for MVAPICH2 compiled binaries
- Issue present L7 function



CLAMR compiled w/ MVAPICH2



CLAMR compiled w/ OpenMPI

CLAMR Observations

- Difference between the min and max times for neighbor boundary updates in MVAPICH2 compiled binaries
- MPI processes stuck in **MPI_Waitall**
- **MPI_Barrier** after **MPI_Waitall** call improves performance for MVAPICH2 compiled binaries
- Issue present L7 function

```

CPU: state timer refine potential      2.2446      4.2383      4.2279      $ min/median/max
CPU: state timer calc mpot            1.1493      2.9786      3.0291      $ min/median/max
CPU: mesh timer refine_smooth         1.0805      1.1493      1.2259      $ min/median/max
CPU: state timer rezone_all           3.9166      3.9117      3.9124      $ min/median/max
CPU: mesh timer partition              0.0000      0.0000      0.0000      $ min/median/max
CPU: mesh timer calc neighbors        33.2815     33.2977     33.3073      $ min/median/max
CPU: mesh timer hash setup            7.1334      7.3969      7.3570      $ min/median/max
CPU: mesh timer hash query            3.9666      4.1236      4.3581      $ min/median/max
CPU: mesh timer find boundary         1.4795      1.8370      2.2330      $ min/median/max
CPU: mesh timer push setup            0.6889      0.6889      0.6889      $ min/median/max
CPU: mesh timer push boundary         0.4611      1.0801      9.3026      $ min/median/max
CPU: mesh timer local_list            0.0000      0.0000      0.0000      $ min/median/max
CPU: mesh timer layer1                2.1061      2.1061      2.1061      $ min/median/max
CPU: mesh timer layer2                1.7464      3.0779      3.0779      $ min/median/max
CPU: mesh timer layer list            0.0000      0.0000      0.0000      $ min/median/max
CPU: mesh timer copy mesh data       0.3132      0.3458      0.3458      $ min/median/max
CPU: mesh timer fill mesh_ghost       0.0050      0.0068      0.0080      $ min/median/max
CPU: mesh timer fill neigh_ghost      0.5816      0.3174      3.5235      $ min/median/max
CPU: mesh timer set corner neigh      0.0000      0.0000      0.0000      $ min/median/max
CPU: mesh timer neigh adjust          0.2602      0.2641      0.2641      $ min/median/max
CPU: mesh timer setup_comm            0.7076      0.7836      15.9801     $ min/median/max
CPU: state timer mess sum             0.0094      0.0295      0.0330      $ min/median/max
CPU: mesh timer load balance          0.9030      0.9173      0.9793      $ min/median/max
CPU: mesh timer calc spatial_coordi   0.0000      0.0000      0.0000      $ min/median/max
-----
Profiling: Total CPU time was 62.8369 62.8547 62.8919 $ min/median/max
-----
Mesh Ops (Neigh+rezone+smooth+balance) 47.2573 47.6461 48.1677 $ min/median/max
Mesh Ops Percentage 75.1697 75.8151 76.5403 percent min/median/max
-----
Profiling: Total time was 64.0156 64.0156 64.0157 $ min/median/max
  
```

MVAPICH2 CLAMR w/o MPI_Barrier

```

CPU: state timer rezone_all           3.9346      3.9356      3.9360      $ min/median/max
CPU: mesh timer partition              0.0000      0.0000      0.0000      $ min/median/max
CPU: mesh timer calc neighbors        25.9420     25.9515     25.9595      $ min/median/max
CPU: mesh timer hash setup            7.1185      7.5718      7.9349      $ min/median/max
CPU: mesh timer hash query            4.0759      4.2055      4.4323      $ min/median/max
CPU: mesh timer find boundary         1.4763      1.8480      2.2263      $ min/median/max
CPU: mesh timer push setup            0.6000      0.6000      0.6000      $ min/median/max
CPU: mesh timer push boundary         0.3218      0.4595      0.5910      $ min/median/max
CPU: mesh timer local_list            0.0000      0.0000      0.0000      $ min/median/max
CPU: mesh timer layer1                0.5893      2.1708      3.3478      $ min/median/max
CPU: mesh timer layer2                2.0940      2.0940      2.0940      $ min/median/max
CPU: mesh timer layer list            0.0000      0.0000      0.0000      $ min/median/max
CPU: mesh timer copy mesh data       0.3374      0.3658      0.4930      $ min/median/max
CPU: mesh timer fill mesh_ghost       0.0050      0.0069      0.0081      $ min/median/max
CPU: mesh timer fill neigh_ghost     3.4727      3.4727      3.4727      $ min/median/max
CPU: mesh timer set corner neigh      0.0000      0.0000      0.0000      $ min/median/max
CPU: mesh timer neigh adjust          0.2609      0.2649      0.2676      $ min/median/max
CPU: mesh timer setup_comm            0.0755      2.3668      6.4141      $ min/median/max
CPU: state timer mess sum             0.0053      0.0290      0.0324      $ min/median/max
CPU: mesh timer load balance          0.9842      9.3409      9.7475      $ min/median/max
CPU: mesh timer calc spatial_coordi   0.0000      0.0000      0.0000      $ min/median/max
-----
Profiling: Total CPU time was 55.5448 55.5502 55.5855 $ min/median/max
-----
Mesh Ops (Neigh+rezone+smooth+balance) 40.0000 40.3497 40.7863 $ min/median/max
Mesh Ops Percentage 71.9882 72.6276 73.3951 percent min/median/max
  
```

MVAPICH2 CLAMR w/ MPI_Barrier



CLAMR Research Questions

1. How do methods for performing weak progress in MPI implementations affect the performance of an application in different runtime environments?
2. What can be done to optimize the effectiveness of a weak progress engine in order to avoid performance anomalies present amongst different MPI implementations?



Future Work

- Expand the investigation of unstructured communication to different AMR codebases (HIGRAD, xRAGE, CLAMR)
- Begin formulating a classification of unstructured communication
- Begin investigating methods of weak progress
- Assist L7 Benchmark Development (Savannah Camp's work)
- Provide feedback on Dr. Ryan Marshall's reproducibility project



Center for Understandable
Performant Exascale
Communication Systems

 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

Acknowledgements

- Galen Shipman - LANL Computer Science
- Anthony Skjellum, Patrick Bridges, Amanda Bienz, Puri Bangalore, and Savannah Camp (PSAAP)
- Ryan Marshall - University of Alabama (PSAAP)
- The Parthenon development team (especially Jonah Miller)
- Bob Robey - LANL



Center for Understandable
Performant Exascale
Communication Systems

 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA

References

- CLAMR, <https://github.com/lanl/CLAMR>
- Parthenon, <https://github.com/lanl/parthenon>
- Caliper, <https://github.com/LLNL/Caliper>, <https://software.llnl.gov/Caliper/>
- Tau, <https://www.cs.uoregon.edu/research/tau/home.php>



Center for Understandable
Performant Exascale
Communication Systems

 THE UNIVERSITY OF TENNESSEE
CHATTANOOGA