# Performance Engineering using MVAPICH and TAU

**Sameer Shende**, Srinivasan Ramesh, Allen D. Malony, Wyatt Spear, Kevin Huck
*University of Oregon*
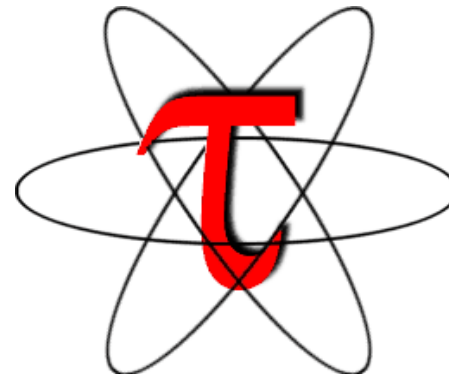
MVAPICH User's Group (MUG) Meeting
12:00pm – 12:30pm, Tuesday, August 20, 2019
Ohio Supercomputing Center, Columbus, Ohio

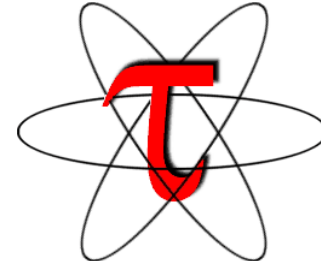THE OHIO STATE UNIVERSITY

UNIVERSITY OF OREGON

# Outline

- **Introduction**
- **The MPI Tools Interfaces and Benefits**
- **Integrating TAU and MVAPICH2 with MPI_T**

# Acknowledgments

- **The MVAPICH2 team The Ohio State University**
  - http://mvapich.cse.ohio-state.edu
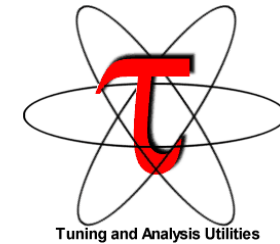- **TAU team at the University of Oregon**
  - http://tau.uoregon.edu

# TAU Performance System®

- **Tuning and Analysis Utilities (25+ year project)**
- **Comprehensive performance profiling and tracing**
  - Integrated, scalable, flexible, portable
  - Targets all parallel programming/execution paradigms

- **Integrated performance toolkit**
  - Instrumentation, measurement, analysis, visualization
  - Widely-ported performance profiling / tracing system
  - Performance data management and data mining
  - Open source (BSD-style license)
  - Uses performance and control variables to interface with MVAPICH2
- **Integrates with application frameworks**
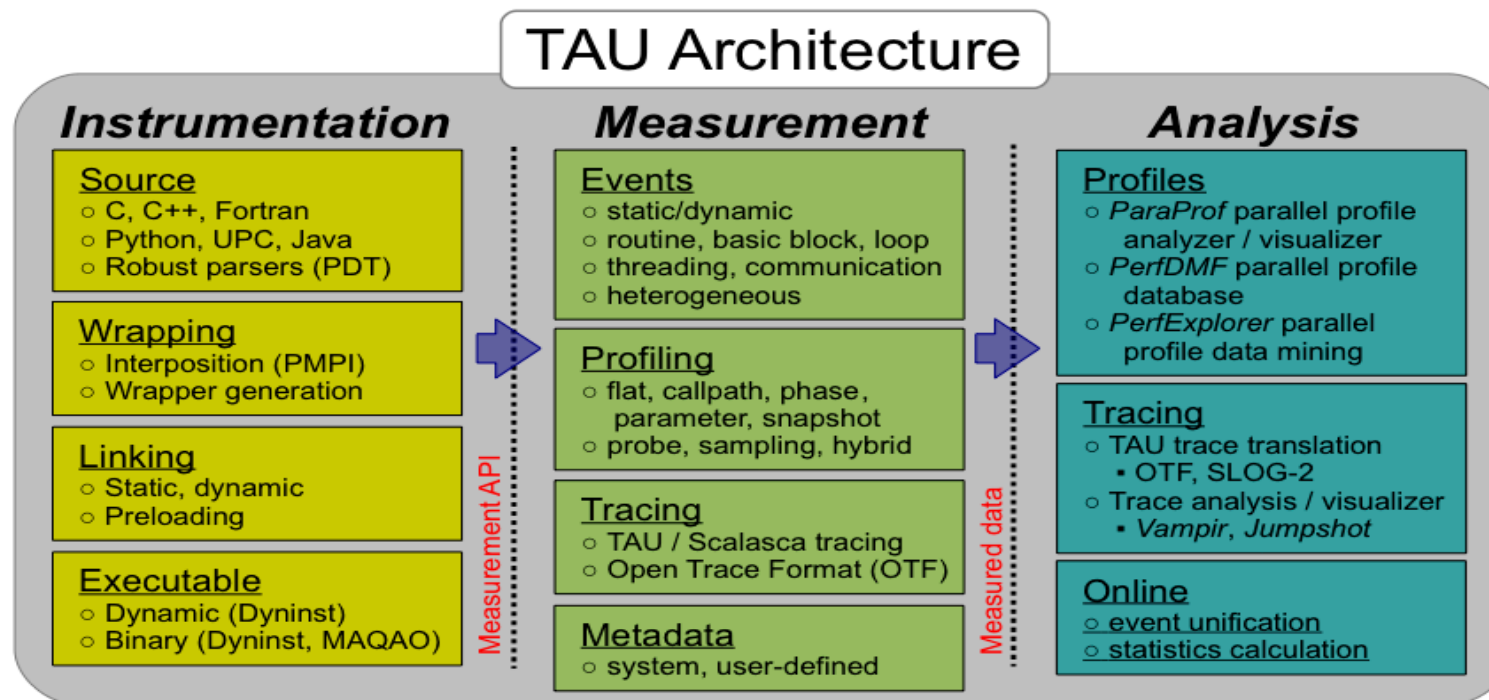- **http://tau.uoregon.edu**

# Understanding Application Performance using TAU

- **How much time** is spent in each application routine and outer *loops*? Within loops, what is the contribution of each *statement*?

- **How many instructions** are executed in these code regions?
Floating point, Level 1 and 2 *data cache misses*, hits, branches taken?

- **What is the memory usage** of the code? When and where is memory allocated/de-allocated? Are there any memory leaks?

- **What are the I/O characteristics** of the code?  What is the peak read and write *bandwidth* of individual calls, total volume?

- **What is the contribution of each *phase*** of the program? What is the time wasted/spent waiting for collectives, and I/O operations in Initialization, Computation, I/O phases?

- **How does the application *scale*?** What is the efficiency, runtime breakdown of performance across different core counts?

- **How can I tune MPI for better performance?** What performance and control does MVAPICH2 export to observe and control its performance?

# TAU Performance System®

## Parallel performance framework and toolkit

- Supports all HPC platforms, compilers, runtime system
- Provides portable instrumentation, measurement, analysis



TAU Architecture

**Instrumentation**

**Source**
- C, C++, Fortran
- Python, UPC, Java
- Robust parsers (PDT)

**Wrapping**
- Interposition (PMPI)
- Wrapper generation

**Linking**
- Static, dynamic
- Preloading

**Executable**
- Dynamic (Dyninst)
- Binary (Dyninst, MAQAO)

Measurement API

**Measurement**

**Events**
- static/dynamic
- routine, basic block, loop
- threading, communication
- heterogeneous

**Profiling**
- flat, callpath, phase, parameter, snapshot
- probe, sampling, hybrid

**Tracing**
- TAU / Scalasca tracing
- Open Trace Format (OTF)

**Metadata**
- system, user-defined

Measured data

**Analysis**

**Profiles**
- *ParaProf* parallel profile analyzer / visualizer
- *PerfDMF* parallel profile database
- *PerfExplorer* parallel profile data mining

**Tracing**
- TAU trace translation
  - OTF, SLOG-2
- Trace analysis / visualizer
  - *Vampir, Jumpshot*

**Online**
- event unification
- statistics calculation

THE OHIO STATE UNIVERSITY

UNIVERSITY OF OREGON

# TAU Instrumentation Approach

**Supports both direct and indirect performance observation**

- Direct instrumentation of program (system) code (probes)

- Instrumentation invokes performance measurement

- Event measurement: performance data, meta-data, context

- Indirect mode supports sampling based on periodic timer or hardware performance counter overflow based interrupts

**Support for user-defined events**

- *Interval* (Start/Stop) events to measure exclusive & inclusive duration

- *Atomic events* (Trigger at a single point with data, e.g., heap memory)
  - Measures total, samples, min/max/mean/std. deviation statistics

- *Context events* (are atomic events with executing context)
  - Measures above statistics for a given calling path

# Direct Observation: Events

## Event types

- Interval events (begin/end events)
  - Measures exclusive & inclusive durations between events
  - Metrics monotonically increase
- Atomic events (trigger with data value)
  - Used to capture performance data state
  - Shows extent of variation of triggered values (min/max/mean)

## Code events

- Routines, classes, templates
- Statement-level blocks, loops

# Inclusive and Exclusive Profiles

- **Performance with respect to code regions**
- **Exclusive measurements for region only**
- **Inclusive measurements includes child regions**



```
int foo()
{
    int a;
    a =a + 1;

bar();

    a =a + 1;
    return a;
}
```

exclusive duration

inclusive duration

# How much data do you want?

Limited Profile  Loop Profile  Callpath Profile

O(KB) ⟷ O(TB)

**Flat Profile**  Callsite Profile  **Trace**

# Types of Performance Profiles

***Flat* profiles**
- Metric (e.g., time) spent in an event
- Exclusive/inclusive, # of calls, child calls, …

***Callpath* profiles**
- Time spent along a calling path (edges in callgraph)
- *"main=> f1 => f2 => MPI_Send"*
- Set the TAU_CALLPATH and TAU_CALLPATH_DEPTH environment variables

***Callsite* profiles**
- Time spent along in an event at a given source location
- Set the TAU_CALLSITE environment variable

***Phase* profiles**
- Flat profiles under a phase (nested phases allowed)
- Default "main" phase
- Supports static or dynamic (e.g. per-iteration) phases

## Instrumentation

### Add hooks in the code to perform measurements

**Source instrumentation using a preprocessor**

- Add timer start/stop calls in a copy of the source code.
- Use Program Database Toolkit (PDT) for parsing source code.
- Requires recompiling the code using TAU shell scripts (tau_cc.sh, tau_f90.sh)
- Selective instrumentation (filter file) can reduce runtime overhead and  narrow instrumentation focus.

**Compiler-based instrumentation**

- Use system compiler to add a special flag to insert hooks at routine entry/exit.
- Requires recompiling using TAU compiler scripts (tau_cc.sh, tau_f90.sh…)

**Runtime preloading of TAU's Dynamic Shared Object (DSO)**

- No need to recompile code! Use **mpirun tau_exec ./app**  with options.
- Requires dynamic executable (link using **–dynamic** on Cray systems).

# Outline

- **Introduction**

- **The MPI Tools Interfaces and Benefits**

- **Integrating TAU and MVAPICH2 with MPI_T**

# Overview of the MVAPICH2 Project

**High Performance open-source MPI Library for InfiniBand, Omni-Path, Ethernet/iWARP, and RDMA over Converged Ethernet (RoCE)**

- MVAPICH (MPI-1), MVAPICH2 (MPI-2.2 and MPI-3.1), Started in 2001, First version available in 2002

- MVAPICH2-X (MPI + PGAS), Available since 2011

- Support for GPGPUs (MVAPICH2-GDR) and MIC (MVAPICH2-MIC), Available since 2014

- Support for Virtualization (MVAPICH2-Virt), Available since 2015

- Support for Energy-Awareness (MVAPICH2-EA), Available since 2015

- Support for InfiniBand Network Analysis and Monitoring (OSU INAM) since 2015

- **Used by more than 3,025 organizations in 89 countries**

- **More than 562,000 (> 0.5 million) downloads from the OSU site directly**

- Empowering many TOP500 clusters (Nov '18 ranking)

  - 3[rd] ranked 10,649,640-core cluster (Sunway TaihuLight) at NSC, Wuxi, China

  - 5[th], 448,448 cores (Frontera) at TACC

  - 8[th], 391,680 cores (ABCI) in Japan

  - 15[th], 570,020 cores (Neurion) in S. Korea and many others

- Available with software stacks of many vendors and Linux Distros (RedHat, SuSE, and OpenHPC)

- **http://mvapich.cse.ohio-state.edu**

**Empowering Top500 systems for over a decade**

**Partner in TACC Frontera System**

18 Years & Counting!

2001-2019

THE OHIO STATE UNIVERSITY

UNIVERSITY OF OREGON

# MVAPICH2 and TAU



- **TAU and MVAPICH2 are enhanced with the ability to generate recommendations and engineering performance report**
- **MPI libraries like MVAPICH2 are now "reconfigurable" at runtime**
- **TAU and MVAPICH2 communicate using the MPI-T interface**

# Why PMPI is not good enough?



- Takes a "black box" view of the MPI library

# MPI_T usage semantics



```
int MPI_T_cvar_get_info(int cvar_index, char *name, int*name_len, int *verbosity,
     int MPI_T_pvar_handle_alloc(MPI_T_pvar_session session, int pvar_index,
          char *desc, int *desc_len, int *bind, int *scope);
```

# MPI_T support with MVAPICH2

- Support performance variables (PVAR)

  - Variables to track different components within the MPI library

- Initial support for Control Variables  (CVAR)

  - Variables to modify the behavior of MPI Library

| | | |
|---|---|---|
| **Memory Usage:**<br>- current level<br>- maximum watermark | **InfiniBand N/W:**<br>- #control packets<br>- #out-of-order packets | **Pt-to-pt messages:**<br>- unexpected queue length<br>- unexp. match attempts<br>- recvq. length |
| **Registration cache:**<br>- hits<br>- misses | **Shared-memory:**<br>- limic/ CMA<br>- buffer pool size & usage | **Collective ops:**<br>- comm. creation<br>- #algorithm invocations<br>[Bcast – 8; Gather – 10]<br>… |

# Co-designing Applications to use MPI-T

Example Pseudo-code: Optimizing the eager limit dynamically:

```
MPI_T_init_thread(..)
MPI_T_cvar_get_info(MV2_EAGER_THRESHOLD)
if (msg_size < MV2_EAGER_THRESHOLD + 1KB)
    MPI_T_cvar_write(MV2_EAGER_THRESHOLD, +1024)
MPI_Send(..)
MPI_T_finalize(..)
```

# Outline

- **Introduction**
- **The MPI Tools Interfaces and Benefits**
- **Integrating TAU and MVAPICH2 with MPI_T**

# Integrating TAU with MVAPICH2 through MPI_T Interface



- **Enhance existing support for MPI_T in MVAPICH2 to expose a richer set of performance and control variables**

- **Get and display MPI Performance Variables (PVARs) made available by the runtime in TAU**

- **Control the runtime's behavior via MPI Control Variables (CVARs)**

- **Add support to MVAPICH2 and TAU for interactive performance engineering sessions**

# Three Scenarios for Integration



Scenario 1: Non-interactive mode



Scenario 2: User-interactive mode



Scenario 3: Policy driven mode

# TAU Performance Measurement Model



enter/exit events
are "interval" events

(in shared memory)

application-wide
performance data

# TAU Plugin Architecture

**Extend TAU *event* interface for plugins**

- Events: *interval, atomic*
- Specialized on event ID
- Synchronous operation

**Create TAU interface for *trigger* plugins**

- Named trigger
- Pass application data
- Synchronous
- Asynchronous using agent plugin

# TAU Plugin Architecture

- **Both event and trigger plugins are synchronous**
  - Directly called from the application
  - Execute inline with the application
  - Use an application's thread of execution
- **Consider utilizing a separate thread of execution to perform performance analysis functions**
  - For instance, periodic daemon to sample performace
- **Design an *agent* plugin mechanism**
  - Create an execution thread to execute plugin
  - Register plugin with this execution thread

THE OHIO STATE UNIVERSITY

UNIVERSITY OF OREGON

# TAU Plugin Architecture

- **Parallel performance systems do not typically do runtime analytics when making measurements**

- **Want to extend a performance system with additional analytics functionality WITHOUT building it directly into the performance system**

- **Apply a plugin architecture approach**

  - Develop analytics plugins (common, application)

  - Register (load) them with the performance system

- **Plugins have access to performance data state**

- **Plugins can utilize the parallel execution context**

# Plugin-based Infrastructure for Non-Interactive Tuning

- **TAU supports a *fully-customizable* plugin infrastructure based on callback event handler registration for salient states inside TAU:**
  - Function Registration / Entry / Exit
  - Phase Entry / Exit
  - Atomic Event Registration / Trigger
  - Init / Finalize Profiling
  - Interrupt Handler
  - *MPI_T*
- **Application can define its own "trigger" states and associated plugins**
  - Pass arbitrary data to trigger state plugins

# TAU Customization

- **TAU states can be *named* or *generic***
- **TAU distinguishes named states in a way that allows for separation of occurrence of a state from the action associated with it**
  - Function entry for "foo" and "bar" represent distinguishable states in TAU
- **TAU maintains an internal map of a list of plugins associated with each state**

# TAU Runtime Control of Plugin

- **TAU defines a plugin API to deliver access control to the internal plugin map**
- **User can specify a regular expression to control plugins executed for a class of named states at runtime**
  - Access to map on a process is serialized: application is expected to access map through main thread

## TAU Phase Based Recommendations

- **MiniAMR: Benefits from hardware offloading using SHArP hardware offload protocol supported by MVAPICH2 for MPI_Allreduce operation**
- **Recommendation Plugin:**
  - Registers callback for *"Phase Exit"* event
  - Monitors message size through PMPI interface
  - If message size is low and execution time inside MPI_Allreduce is significant, a recommendation is generated on ParaProf (TAU's GUI) for the user to set the CVAR enabling SHArP

# TAU Per-Phase Recommendations in ParaProf

# Enhancing MPI_T Support

- **Introduced support for new MPI_T based CVARs to MVAPICH2**
  - MPIR_CVAR_MAX_INLINE_MSG_SZ
    - Controls the message size up to which "inline" transmission of data is supported by MVAPICH2
  - MPIR_CVAR_VBUF_POOL_SIZE
    - Controls the number of internal communication buffers (VBUFs) MVAPICH2 allocates initially. Also, MPIR_CVAR_VBUF_POOL_REDUCED_VALUE[1] ([2...n])
  - MPIR_CVAR_VBUF_SECONDARY_POOL_SIZE
    - Controls the number of VBUFs MVAPICH2 allocates when there are no more free VBUFs available
  - MPIR_CVAR_IBA_EAGER_THRESHOLD
    - Controls the message size where MVAPICH2 switches from eager to rendezvous protocol for large messages
- **TAU enhanced with support for setting MPI_T CVARs in a non-interactive mode for uninstrumented applications**

THE OHIO STATE UNIVERSITY

UNIVERSITY OF OREGON

# MVAPICH2

- **Several new MPI_T based PVARs added to MVAPICH2**
  - mv2_vbuf_max_use, mv2_total_vbuf_memory etc
- **Enhanced TAU with support for tracking of MPI_T PVARs and CVARs for uninstrumented applications**
  - ParaProf, TAU's visualization front end, enhanced with support for displaying PVARs and CVARs
  - TAU provides tau_exec, a tool to transparently instrument MPI routines
    - Uninstrumented:
      % mpirun –np 1024 ./a.out
    - Instrumented:
      - % export TAU_TRACK_MPI_T_PVARS=1
      - % export TAU_MPI_T_CVAR_METRICS=MPIR_CVAR_VBUF_POOL_SIZE
      - % export TAU_MPI_T_CVAR_VALUES=16
      - % mpirun -np 1024 *tau_exec -T mvapich2,mpit* ./a.out

# PVARs Exposed by MVAPICH2

# CVARs Exposed by MVAPICH2

# Using MVAPICH2 and TAU with Multiple CVARs

● To set CVARs or read PVARs using TAU for an uninstrumented binary:

```
% export TAU_TRACK_MPI_T_PVARS=1
% export TAU_MPI_T_CVAR_METRICS=
      MPIR_CVAR_VBUF_POOL_REDUCED_VALUE[1],
      MPIR_CVAR_IBA_EAGER_THRESHOLD
% export TAU_MPI_T_CVAR_VALUES=32,64000
% export PATH=/path/to/tau/x86_64/bin:$PATH
% mpirun -np 1024 tau_exec -T mvapich2,mpit   ./a.out
% paraprof
```

# VBUF usage without CVARs

TAU: ParaProf: Context Events for: node 0 - mpit_withoutcvar_bt.C.1k.ppk

| Name △ | MaxValue | MinValue | MeanValue | Std. Dev. | NumSamples | Total |
|---|---|---|---|---|---|---|
| mv2_total_vbuf_memory (Total amount of memory in bytes used for VBUFs) | 3,313,056 | 3,313,056 | 3,313,056 | 0 | 1 | 3,313,056 |
| mv2_ud_vbuf_allocated (Number of UD VBUFs allocated) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_ud_vbuf_available (Number of UD VBUFs available) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_ud_vbuf_freed (Number of UD VBUFs freed) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_ud_vbuf_inuse (Number of UD VBUFs inuse) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_ud_vbuf_max_use (Maximum number of UD VBUFs used) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_vbuf_allocated (Number of VBUFs allocated) | 320 | 320 | 320 | 0 | 1 | 320 |
| mv2_vbuf_available (Number of VBUFs available) | 255 | 255 | 255 | 0 | 1 | 255 |
| mv2_vbuf_freed (Number of VBUFs freed) | 25,545 | 25,545 | 25,545 | 0 | 1 | 25,545 |
| mv2_vbuf_inuse (Number of VBUFs inuse) | 65 | 65 | 65 | 0 | 1 | 65 |
| mv2_vbuf_max_use (Maximum number of VBUFs used) | 65 | 65 | 65 | 0 | 1 | 65 |
| num_calloc_calls (Number of MPIT_calloc calls) | 89 | 89 | 89 | 0 | 1 | 89 |
| num_free_calls (Number of MPIT_free calls) | 47,801 | 47,801 | 47,801 | 0 | 1 | 47,801 |
| num_malloc_calls (Number of MPIT_malloc calls) | 49,258 | 49,258 | 49,258 | 0 | 1 | 49,258 |
| num_memalign_calls (Number of MPIT_memalign calls) | 34 | 34 | 34 | 0 | 1 | 34 |
| num_memalign_free_calls (Number of MPIT_memalign_free calls) | 0 | 0 | 0 | 0 | 0 | 0 |

THE OHIO STATE UNIVERSITY

UNIVERSITY OF OREGON

# VBUF usage with CVARs



TAU: ParaProf: Context Events for: node 0 - bt-mz.E.vbuf_pool_16.1k.ppk

| Name △ | MaxValue | MinValue | MeanValue | Std. Dev. | NumSamp... | Total |
|---|---|---|---|---|---|---|
| mv2_total_vbuf_memory (Total amount of memory in bytes used for VBUFs) | 1,815,056 | 1,815,056 | 1,815,056 | 0 | 1 | 1,815,056 |
| mv2_ud_vbuf_allocated (Number of UD VBUFs allocated) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_ud_vbuf_available (Number of UD VBUFs available) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_ud_vbuf_freed (Number of UD VBUFs freed) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_ud_vbuf_inuse (Number of UD VBUFs inuse) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_ud_vbuf_max_use (Maximum number of UD VBUFs used) | 0 | 0 | 0 | 0 | 0 | 0 |
| mv2_vbuf_allocated (Number of VBUFs allocated) | 160 | 160 | 160 | 0 | 1 | 160 |
| mv2_vbuf_available (Number of VBUFs available) | 94 | 94 | 94 | 0 | 1 | 94 |
| mv2_vbuf_freed (Number of VBUFs freed) | 5,479 | 5,479 | 5,479 | 0 | 1 | 5,479 |
| mv2_vbuf_inuse (Number of VBUFs inuse) | 66 | 66 | 66 | 0 | 1 | 66 |
| mv2_vbuf_max_use (Maximum number of VBUFs used) | 66 | 66 | 66 | 0 | 1 | 66 |
| num_calloc_calls (Number of MPIT_calloc calls) | 89 | 89 | 89 | 0 | 1 | 89 |
| num_free_calls (Number of MPIT_free calls) | 130 | 130 | 130 | 0 | 1 | 130 |
| num_malloc_calls (Number of MPIT_malloc calls) | 1,625 | 1,625 | 1,625 | 0 | 1 | 1,625 |
| num_memalign_calls (Number of MPIT_memalign calls) | 56 | 56 | 56 | 0 | 1 | 56 |
| num_memalign_free_calls (Number of MPIT_memalign_free calls) | 0 | 0 | 0 | 0 | 0 | 0 |

TAU: ParaProf Manager

- Applications
  - ▼ Standard Applications
    - ▼ Default App
      - ▼ Default Exp
        - ▼ bt-mz.E.vbuf_pool_16.1k.pp
          - TIME

| TrialField | Value |
|---|---|
| MPI Processor Name | c526-502.stampede.tacc.utexas.edu |
| MPIR_CVAR_VBUF_POOL_SIZE | 0 (old) -> 16 (new), This set the size of the VBUF pool |

Total memory used by VBUFs is reduced from 3,313,056 to 1,815,056

# VBUF Memory Usage Without CVAR

# VBUF Memory Usage With CVAR



```
% export TAU_TRACK_MPI_T_PVARS=1
% export TAU_MPI_T_CVAR_METRICS=MPIR_CVAR_VBUF_POOL_SIZE
% export TAU_MPI_T_CVAR_VALUES=16
% mpirun -np 1024 tau_exec -T mvapich2  ./a.out
```

# TAU: Extending Control Variables on a Per-Communicator Basis

- **Based on named communicators (MPI_Comm_set_name) in an application, TAU allows a user to specify triples to set MPI_T cvars for each communicator:**
  - Communicator name
  - MPI_T CVAR name
  - MPI_T CVAR value
    - % ./configure –mpit –mpi –c++=mpicxx –cc=mpicc –fortran=mpif90 …
    - % make install
    - % export TAU_MPI_T_COMM_METRIC_VALUES=<comm, cvar, value>,…
    - % mpirun –np 64 tau_exec –T mpit   ./a.out
    - % paraprof

# COMB LLNL App MPI_T Tuning for COMB_MPI_CART_COMM

bash-4.2$ TAU_MPI_T_COMM_METRIC_VALUES=COMB_MPI_CART_COMM,MPIR_CVAR_GPUDIRECT_LIMIT,2097152,COMB_MPI_CART_COMM,MPIR_CVAR_USE_GPUDIRECT_RECEIVE_LIMIT, 2097152,COMB_MPI_CART_COMM,MPIR_CVAR_CUDA_IPC_THRESHOLD,16384 MV2_USE_CUDA=1 mpirun -np 8 tau_exec -ebs -T mvapich2,mpit,cuda9,cupti,communicators,gnu -cupti ./ comb -comm post_recv wait_all -comm post_send wait_all -comm wait_recv wait_all -comm wait_send wait_all 200_200_200 -divide 2_2_2 -periodic 1_1_1 -ghost 1_1_1 -vars 3 -cycles 100 - comm cutoff 250 -omp_threads 1

Started rank 0 of 8

Node lassen710

Compiler COMB_COMPILER

Cuda compiler COMB_CUDA_COMPILER

GPU 0 visible undefined

Not built with openmp, ignoring -omp_threads 1.

Cart coords        0       0        0

Message policy cutoff 250

Post Recv using wait_all method

Post Send using wait_all method

Wait Recv using wait_all method

Wait Send using wait_all method

Num cycles      100

Num vars            3

ghost_widths      1       1       1

sizes           200      200      200

divisions        2      2      2

periodic          1       1       1

division map

map               0        0        0

map             100      100      100

map             200      200      200

Starting test memcpy seq dst Host src Host

Starting test Comm mock Mesh seq Host Buffers seq Host seq Host

Starting test Comm mpi Mesh seq Host Buffers seq Host seq Host



Default

With MPI_T CVARs

# COMB Profile

| Name △ | Exclusive TAUGP... | Inclusive TAUGP... | Calls | Child Calls |
|---|---|---|---|---|
| ▼ ◻ .TAU application | 3.114 | 6.855 | 1 | 6,806 |
| ▼ ◼ [CONTEXT] .TAU application | 0 | 3.09 | 103 | 0 |
| ◼ [SAMPLE] COMB::detail::reset_1::operator()(int, int, int, int) const [{/usr/global/tools/tau/tr | 0.57 | 0.57 | 19 | 0 |
| ◼ [SAMPLE] COMB::detail::set_1::operator()(int, int, int, int) const [{/usr/global/tools/tau/trair | 0.42 | 0.42 | 14 | 0 |
| ◼ [SAMPLE] COMB::detail::set_copy::operator()(int, int) const [{/usr/global/tools/tau/training | 0.06 | 0.06 | 2 | 0 |
| ◼ [SAMPLE] COMB::detail::set_copy::operator()(int, int) const [{/usr/global/tools/tau/training | 0.45 | 0.45 | 15 | 0 |
| ◼ [SAMPLE] COMB::detail::set_n1::operator()(int, int) const [{/usr/global/tools/tau/training/a | 0.06 | 0.06 | 2 | 0 |
| ◼ [SAMPLE] __nv_hdl_wrapper_t<false, false, __nv_dl_tag<void (*)(CommContext<mock_pol | 0.03 | 0.03 | 1 | 0 |
| ◼ [SAMPLE] syscall [{/usr/lib64/libc–2.17.so} {0}] | 0.03 | 0.03 | 1 | 0 |
| ◼ [SAMPLE] void detail::copy_idxr_idxr<double const, detail::indexer_list_idx, double, detail: | 0.03 | 0.03 | 1 | 0 |
| ▶ ◼ [SUMMARY] void COMB::do_cycles<mock_pol, seq_pol, seq_pol, seq_pol>(CommContext< | 0.36 | 0.36 | 12 | 0 |
| ▶ ◼ [SUMMARY] void COMB::do_cycles<mock_pol, seq_pol, seq_pol, seq_pol>(CommContext< | 0.33 | 0.33 | 11 | 0 |
| ▶ ◼ [SUMMARY] void COMB::do_cycles<mpi_pol, seq_pol, seq_pol, seq_pol>(CommContext<n | 0.39 | 0.39 | 13 | 0 |
| ▶ ◼ [SUMMARY] void COMB::do_cycles<mpi_pol, seq_pol, seq_pol, seq_pol>(CommContext<n | 0.36 | 0.36 | 12 | 0 |
| ▶ ◼ MPI_Barrier() | 0.292 | 0.292 | 8 | 0 |
| ◼ MPI_Barrier() [ <comm> = <COMB_MPI_CART_COMM> ] | 0.292 | 0.292 | 8 | 0 |

Name: .TAU application => [CONTEXT] .TAU application => [SAMPLE]
COMB::detail::reset_1::operator()(int, int, int, int) const
[{/usr/global/tools/tau/training/apps/COMB_LLNL/Comb/include/comb.hpp} {121}]
Metric Name: TAUGPU_TIME
Value: Exclusive
Units: seconds

| | |
|---|---|
| 0.712 | max |
| 0.51 | min |
| 0.081 | std. dev. |
| 0.595 | mean |
| 0.57 | node 0 |
| 0.69 | node 1 |

Name: .TAU application => [CONTEXT] .TAU application => [SAMPLE]
COMB::detail::set_1::operator()(int, int, int, int) const
[{/usr/global/tools/tau/training/apps/COMB_LLNL/Comb/include/comb.hpp} {90}]
Metric Name: TAUGPU_TIME
Value: Exclusive
Units: seconds

| | |
|---|---|
| 0.6 | max |
| 0.361 | min |
| 0.068 | std. dev. |
| 0.436 | mean |
| 0.42 | node 0 |
| 0.45 | node 1 |

THE OHIO STATE UNIVERSITY

UNIVERSITY OF OREGON

# CVARs Exposed by MVAPICH2

Metadata for n,c,t 0,0,0

| Name | Value |
| --- | --- |
| MPI Processor Name | lassen710 |
| MPIR_CVAR_CUDA_IPC_THRESHOLD | 16384 |
| MPIR_CVAR_GPUDIRECT_LIMIT | 2097152 |
| MPIR_CVAR_USE_GPUDIRECT_RECEIVE_LIMIT | 2097152 |
| MPI_T CVAR: MPIR_CVAR_ABORT_ON_LEAKED_HANDLES | If true, MPI will call MPI_Abort at MPI_Finalize if any MPI object handles ha... |
| MPI_T CVAR: MPIR_CVAR_ALLGATHERV_PIPELINE_MSG_SIZE | The smallest message size that will be used for the pipelined, large–mes... |
| MPI_T CVAR: MPIR_CVAR_ALLGATHER_COLLECTIVE_ALGORITHM | This CVAR selects proper collective algorithm for allgather operation. |
| MPI_T CVAR: MPIR_CVAR_ALLGATHER_LONG_MSG_SIZE | For MPI_Allgather and MPI_Allgatherv, the long message algorithm will be... |
| MPI_T CVAR: MPIR_CVAR_ALLGATHER_SHORT_MSG_SIZE | For MPI_Allgather and MPI_Allgatherv, the short message algorithm will b... |
| MPI_T CVAR: MPIR_CVAR_ALLREDUCE_COLLECTIVE_ALGORITHM | This CVAR selects proper collective algorithm for allreduce operation. |
| MPI_T CVAR: MPIR_CVAR_ALLREDUCE_SHORT_MSG_SIZE | the short message algorithm will be used if the send buffer size is <= th... |
| MPI_T CVAR: MPIR_CVAR_ALLTOALLV_COLLECTIVE_ALGORITHM | This CVAR selects proper collective algorithm for alltoallv operation. |
| MPI_T CVAR: MPIR_CVAR_ALLTOALL_COLLECTIVE_ALGORITHM | This CVAR selects proper collective algorithm for alltoall operation. |
| MPI_T CVAR: MPIR_CVAR_ALLTOALL_MEDIUM_MSG_SIZE | the medium message algorithm will be used if the per–destination messa... |
| MPI_T CVAR: MPIR_CVAR_ALLTOALL_SHORT_MSG_SIZE | the short message algorithm will be used if the per–destination message... |
| MPI_T CVAR: MPIR_CVAR_ALLTOALL_THROTTLE | max no. of irecvs/isends posted at a time in some alltoall algorithms. Set... |
| MPI_T CVAR: MPIR_CVAR_ASYNC_PROGRESS | If set to true, MPICH will initiate an additional thread to make asynchrono... |
| MPI_T CVAR: MPIR_CVAR_BCAST_COLLECTIVE_ALGORITHM | This CVAR selects proper collective algorithm for broadcast operation. |
| MPI_T CVAR: MPIR_CVAR_BCAST_LONG_MSG_SIZE | Let's define short messages as messages with size < MPIR_CVAR_BCAST_... |
| MPI_T CVAR: MPIR_CVAR_BCAST_MIN_PROCS | Let's define short messages as messages with size < MPIR_CVAR_BCAST_... |
| MPI_T CVAR: MPIR_CVAR_BCAST_SHORT_MSG_SIZE | Let's define short messages as messages with size < MPIR_CVAR_BCAST_... |
| MPI_T CVAR: MPIR_CVAR_CH3_EAGER_MAX_MSG_SIZE | This cvar controls the message size at which CH3 switches from eager to... |
| MPI_T CVAR: MPIR_CVAR_CH3_ENABLE_HCOLL | If true, enable HCOLL collectives. |
| MPI_T CVAR: MPIR_CVAR_CH3_INTERFACE_HOSTNAME | If non–NULL, this cvar specifies the IP address that other processes shoul... |
| MPI_T CVAR: MPIR_CVAR_CH3_NOLOCAL | If true, force all processes to operate as though all processes are located... |
| MPI_T CVAR: MPIR_CVAR_CH3_ODD_EVEN_CLIQUES | If true, odd procs on a node are seen as local to each other, and even pr... |
| MPI_T CVAR: MPIR_CVAR_CH3_PORT_RANGE | The MPIR_CVAR_CH3_PORT_RANGE environment variable allows you to s... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_ACTIVE_REQ_THRESHOLD | Threshold of number of active requests to trigger blocking waiting in op... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_DELAY_ISSUING_FOR_PIGGYBACKING | Specify if delay issuing of RMA operations for piggybacking LOCK/UNLOC... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_OP_GLOBAL_POOL_SIZE | Size of the Global RMA operations pool (in number of operations) that st... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_OP_PIGGYBACK_LOCK_DATA_SIZE | Specify the threshold of data size of a RMA operation which can be piggy... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_OP_WIN_POOL_SIZE | Size of the window–private RMA operations pool (in number of operation... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_POKE_PROGRESS_REQ_THRESHOLD | Threshold at which the RMA implementation attempts to complete reque... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_SCALABLE_FENCE_PROCESS_NUM | Specify the threshold of switching the algorithm used in FENCE from the ... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_SLOTS_SIZE | Number of RMA slots during window creation. Each slot contains a linked... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_TARGET_GLOBAL_POOL_SIZE | Size of the Global RMA targets pool (in number of targets) that stores inf... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_TARGET_LOCK_DATA_BYTES | Size (in bytes) of available lock data this window can provided. If current ... |
| MPI_T CVAR: MPIR_CVAR_CH3_RMA_TARGET_LOCK_ENTRY_WIN_POOL_SIZE | Size of the window–private RMA lock entries pool (in number of lock entr... |

# TAU's ParaProf 3D Browser

# Download TAU from U. Oregon



Tuning and Analysis Utilities

http://tau.uoregon.edu

http://taucommander.com

http://www.hpclinux.com  [OVA for VirtualBox]

https://e4s.io  [ Extreme-Scale Scientific Software Stack, Containers for HPC]

for more information

Free download, open source, BSD license

# PRL, OACISS, University of Oregon, Eugene



www.uoregon.edu

# Support Acknowledgements

**US Department of Energy (DOE)**
- ANL
- Office of Science contracts, ECP
- SciDAC, LBL contracts
- LLNL-LANL-SNL ASC/NNSA contract
- Battelle, PNNL and ORNL contract

**Department of Defense (DoD)**
- PETTT, HPCMP

**National Science Foundation (NSF)**
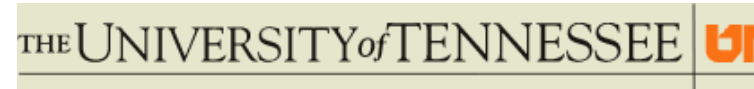- SI2-SSI, Glassbox

**NASA**

**CEA, France**

**Partners:**
- University of Oregon
- The Ohio State University
- ParaTools, Inc.
- University of Tennessee, Knoxville
- T.U. Dresden, GWT
- Jülich Supercomputing Center

# Acknowledgment

THE OHIO STATE UNIVERSITY

UNIVERSITY OF OREGON