

Tutorial: Hands-on Practical Hybrid Parallel Application Performance Engineering

Markus Geimer (Jülich Supercomputing Centre) Sameer Shende (University of Oregon) Bert Wesarg (Technische Universität Dresden) Brian Wylie (Jülich Supercomputing Centre)



Tutorial logistics

- This is a full-day tutorial featuring hands-on exercises
 - Browser-based access available to JUWELS-Booster where the presented tools are installed
 - Exercise materials and slides are provided for you to do on your own
 - We'll be showing how the tools are used and the key functionality that they offer
 - Using provided example measurements & a simple benchmark code
- We want the tutorial to be interactive!
 - Ask questions
 - Indicate when you need assistance

Access to JUWELS-Booster

JUWELS-Booster at JSC will be used for the hands-on exercises

request an account on JSC HPC systems with the training project for tutorial

https://judoor.fz-juelich.de/projects/join/training2341

• If you don't already have an account for JSC HPC systems, you'll first need to register

- provide your E-mail address
- wait for account registration E-mail from user-services.jsc@fz-juelich.de
- follow URL in E-mail to complete registration
- provide your name, affiliation, etc
- read & accept privacy policy
- Account setup may take 30 minutes

Virtual Institute – High Productivity Supercomputing

- Goal: Improve the quality and accelerate the development process of complex simulation codes running on highly-parallel computer systems
- Start-up funding (2006–2011)

by Helmholtz Association of German Research Centres

- Activities
 - Development and integration of HPC application tools
 - Primarily correctness checking & performance analysis
 - Academic workshops: e.g. ProTools@SC23 (Sunday 12 November 2023)
 - Tools training via conference tutorials and multi-day "bring-your-own-code" Tuning Workshops
 - Face-to-face & side-by-side hands-on coaching now successfully migrated to virtual/on-line events

http://www.vi-hps.org

RCH FOR GRAND CHALLENGES

VI-HPS partners (founders)



- Forschungszentrum Jülich
 - Jülich Supercomputing Centre





- Technische Universität Dresden
 - Centre for Information Services & HPC



- University of Tennessee (Knoxville)
 - Innovative Computing Laboratory









VI-HPS partners (cont.)



Barcelona Supercomputing Center

Centro Nacional de Supercomputación



- Lawrence Livermore National Lab.
 - Center for Applied Scientific Computing
- Leibniz Supercomputing Centre









Linaro Ltd.

Allinea software



- Technical University of Darmstadt
 - Laboratory for Parallel Programming





VI-HPS partners (cont.)



Friedrich-Alexander-Universität

Technical University of Munich

Erlangen Regional Computing Center (RRZE)





 Chair for Computer Architecture University of Oregon

University of Stuttgart

Performance Research Laboratory





University of Versailles St-Quentin



HPC Centre









Jniversität Stuttgart



Productivity tools

- MUST / ARCHER
 - MPI & OpenMP usage correctness checking
- PAPI
 - Interfacing to hardware performance counters

CUBE

Analysis report exploration & processing

Scalasca

Large-scale parallel performance analysis

- TAU

Integrated parallel performance system

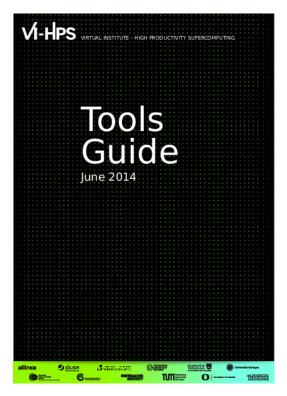
Vampir

Interactive graphical trace visualization & analysis

Score-P

Community-developed instrumentation & measurement infrastructure

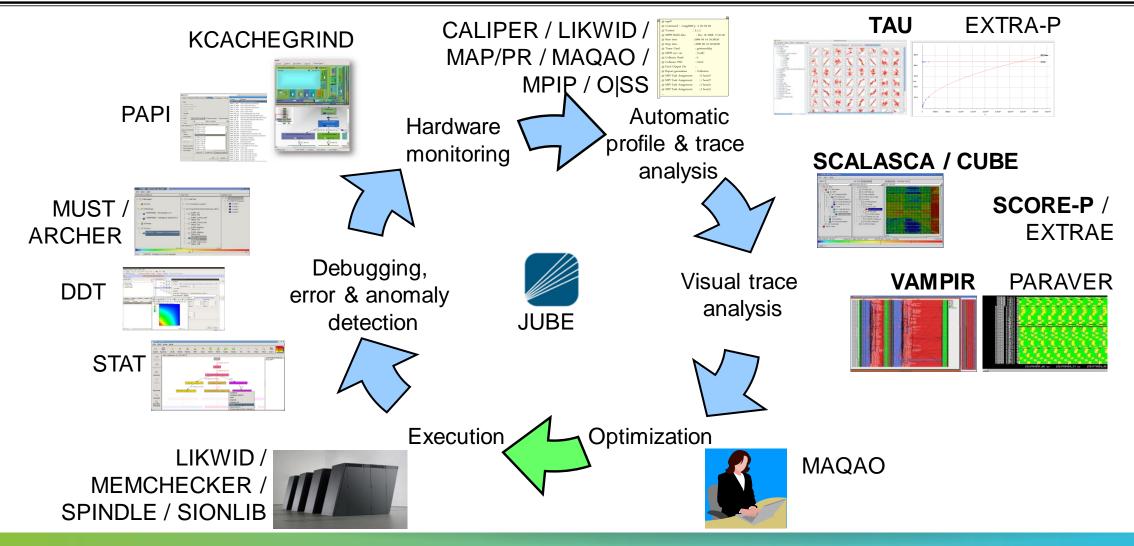
For a brief overview of tools consult the VI-HPS Tools Guide:



Productivity tools (cont.)

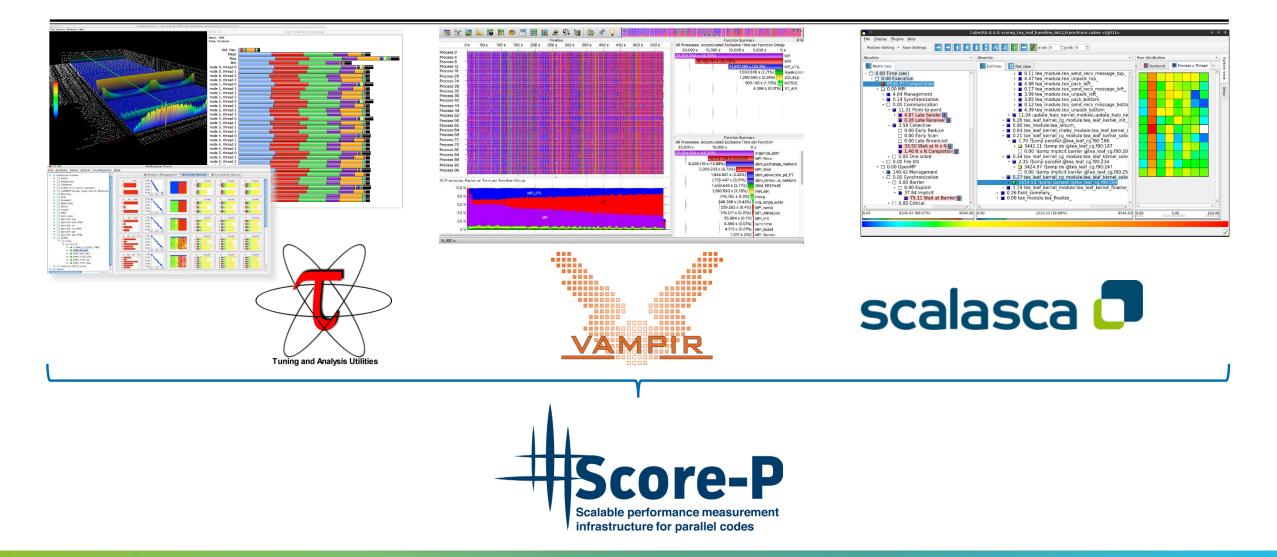
- Caliper: Library for application event annotation, logging and profiling
- Extra-P: Automated performance modelling
- FORGE DDT/MAP/PR: Parallel debugging, profiling & performance reports
- JUBE: Automatic workflow execution for benchmarking, testing & production
- Kcachegrind: Callgraph-based cache analysis [x86 only]
- LIKWID: Performance monitoring & benchmarking
- MAQAO: Assembly instrumentation & optimization [x86-64 only]
- mpiP/mpiPview: MPI profiling tool and analysis viewer
- Open MPI MemChecker: Integrated memory checking
- Open|SpeedShop: Integrated parallel performance analysis environment
- Paraver/Dimemas/Extrae: Event tracing and graphical trace visualization & analysis
- SIONlib/Spindle: Optimized native parallel file I/O & shared library loading
- STAT: Stack trace analysis tools

Technologies and their integration



SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)

Tools featured in this tutorial

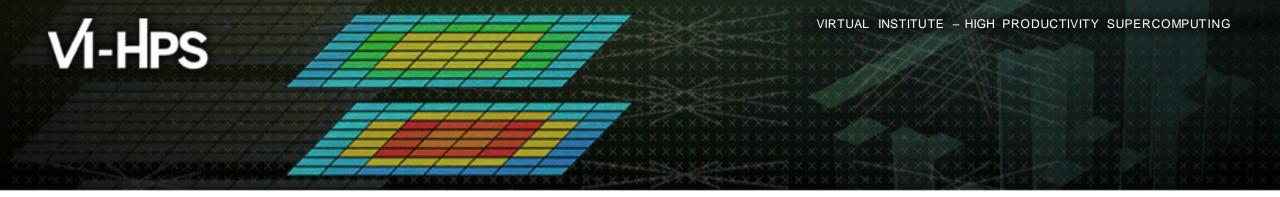


Agenda (Part 1)

Time	Торіс	Presenter
08:30	Welcome & Introduction to VI-HPS	Wylie
08:45	Introduction to parallel performance engineering	Geimer
09:15	Setup for hands-on exercises	Wylie
09:30	Instrumentation & measurement of applications with Score-P	Wesarg
10:00	Break	
10:30	Exploration of execution call-path profiles with CUBE	Wylie
11:00	Configuration & customization of Score-P measurements	Wesarg
11:30	Examination of profiles with TAU ParaProf	Shende
12:00	Lunch break	

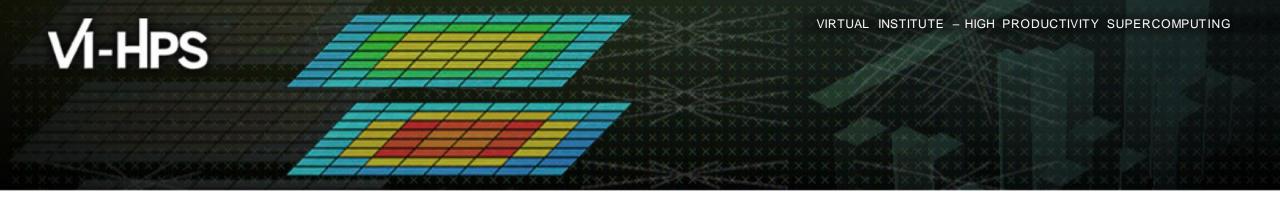
Agenda (Part 2)

Time	Торіс	Presenter
13:30	Résumé + Score-P trace collection	Wylie
13:45	Interactive trace visualization & exploration with Vampir	Wesarg
14:30	Automated analysis of traces for inefficiencies with Scalasca	Geimer
15:00	Break	
15:30	Performance data management with TAU PerfExplorer	Shende
15:45	Specialized Score-P measurements & analyses	Wesarg
16:15	Finding typical parallel performance bottlenecks	Wesarg
16:45	Review	Geimer
17:00	Adjourn	



SC23 Tutorial: Hands-on Practical Hybrid Parallel Application Performance Engineering





Introduction to Parallel Performance Engineering

Markus Geimer Jülich Supercomputing Centre

(with content used with permission from tutorials by Bernd Mohr/JSC and Luiz DeRose/Cray)



Performance: an old problem

Difference Engine

"The most constant difficulty in contriving the engine has arisen from the desire to reduce the time in which the calculations were executed to the shortest which is possible."

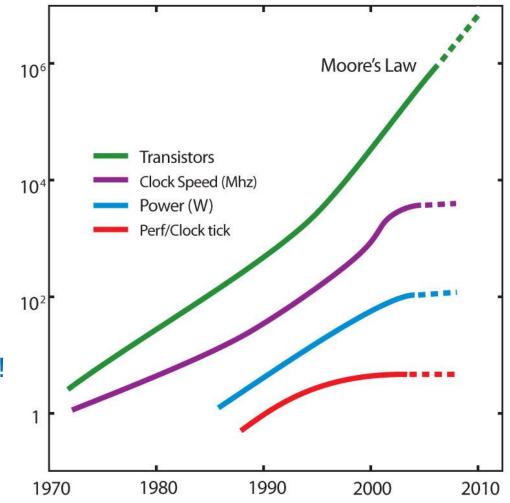
> Charles Babbage 1791 – 1871

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Today: the "free lunch" is over

- Moore's law is still in charge, but
 - Clock rates no longer increase
 - Performance gains only through increased parallelism
- Optimizations of applications more difficult
 - Increasing application complexity
 - Multi-physics
 - Multi-scale
 - Increasing machine complexity
 - Hierarchical networks / memory
 - More CPUs / multi-core

→Every doubling of scale reveals a new bottleneck!



Consequences

- Machine complexity
 - Increasing number of different architectures
 - Additional optimisation challenges related to parallelism
 - Single-core performance issues tied to increased vector length and memory hierarchy

■ → The optimisation process remains key to maintain a reasonable performance level

- Code complexity
 - Codes are harder to optimize and maintain manually
 - Optimisation is time consuming and error-prone
- → Understanding application behaviour is critical

Performance factors of parallel applications

Sequential performance factors

- Computation
 - Choose right algorithm, use optimizing compiler
- Cache and memory
 - Tough! Only limited tool support, hope compiler gets it right
- Input / output
 - Often not given enough attention

Parallel" performance factors

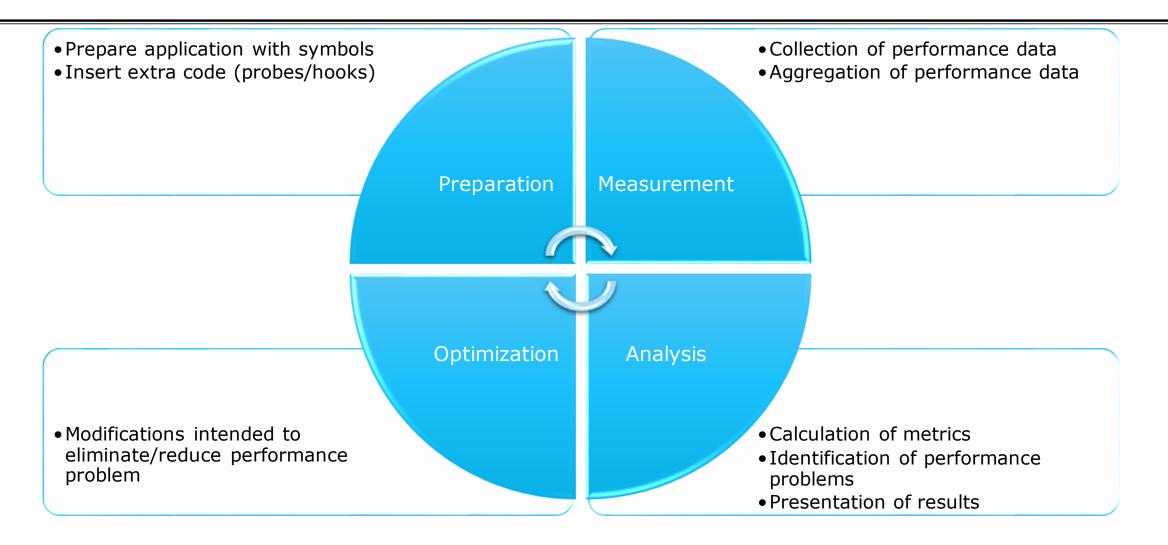
- Partitioning / decomposition
- Communication (i.e., message passing)
- Multithreading
- Synchronization / locking
 - More or less understood, good tool support

Tuning basics

- Successful engineering is a combination of
 - Careful setting of various tuning parameters
 - The right algorithms and libraries
 - Compiler flags and directives
 - ...
 - Thinking !!!
- Measurement is better than guessing
 - To determine performance bottlenecks
 - To compare alternatives
 - To validate tuning decisions and optimizations
 - After each step!
- Modeling is extremely useful but very difficult and rarely available
 - Allows to evaluate performance impact of optimization without implementing it
 - Simplifies search in large parameter space

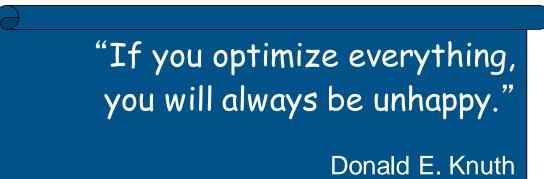
VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Performance engineering workflow



The 80/20 rule

- Programs typically spend 80% of their time in 20% of the code
- Programmers typically spend 20% of their effort to get 80% of the total speedup possible for the application
 - ► → Know when to stop!
- Don't optimize what does not matter
 - → Make the common case fast!



Metrics of performance

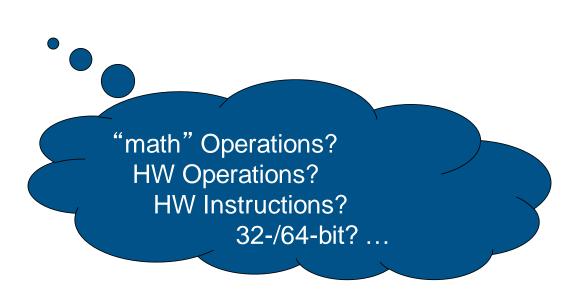
• What can be measured?

- A count of how often an event occurs
 - E.g., the number of MPI point-to-point messages sent
- The duration of some interval
 - E.g., the time spent these send calls
- The size of some parameter
 - E.g., the number of bytes transmitted by these calls
- Derived metrics
 - E.g., rates / throughput
 - Needed for normalization



Example metrics

- Execution time
- Number of function calls
- CPI
 - CPU cycles per instruction
- FLOPS
 - Floating-point operations executed per second



Execution time

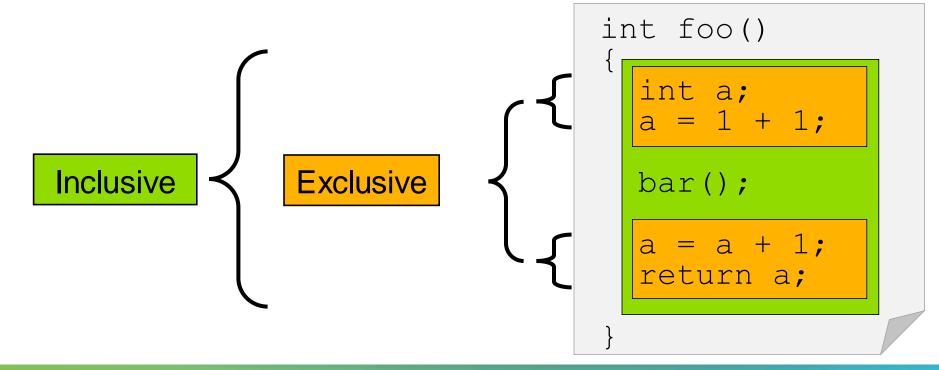
- Wall-clock time
 - Includes waiting time: I/O, memory, other system activities
 - In time-sharing environments also the time consumed by other applications

CPU time

- Time spent by the CPU to execute the application
- Does not include time the program was context-switched out
 - Problem: Does not include inherent waiting time (e.g., I/O)
 - Problem: Portability? What is user, what is system time?
- Problem: Execution time is non-deterministic
 - Use mean or minimum of several runs

Inclusive vs. Exclusive values

- Inclusive
 - Information of all sub-elements aggregated into single value
- Exclusive
 - Information cannot be subdivided further

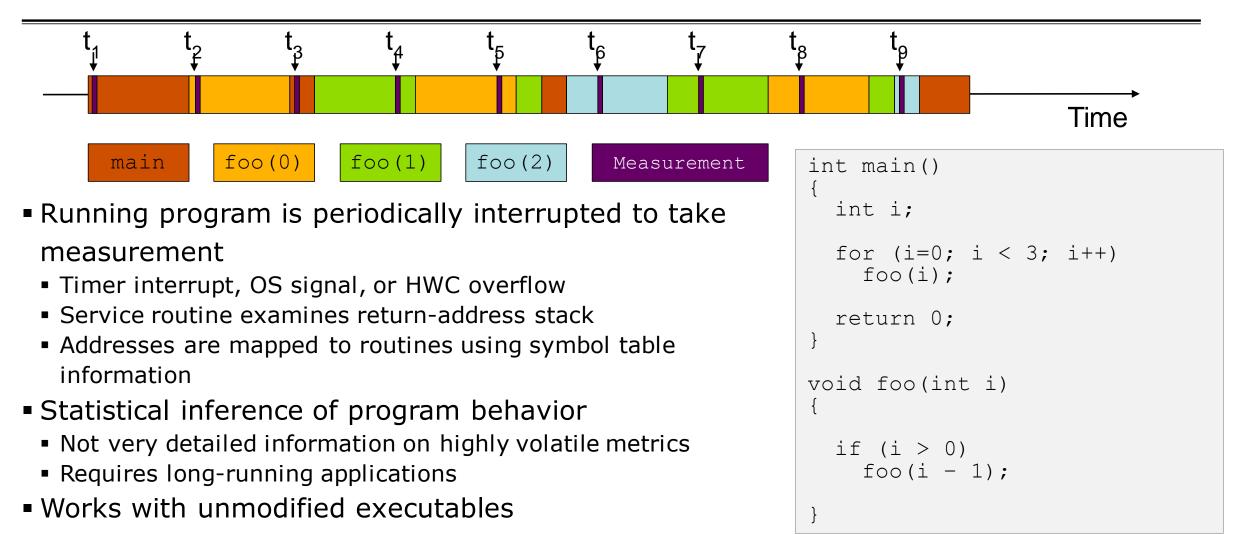


Classification of measurement techniques

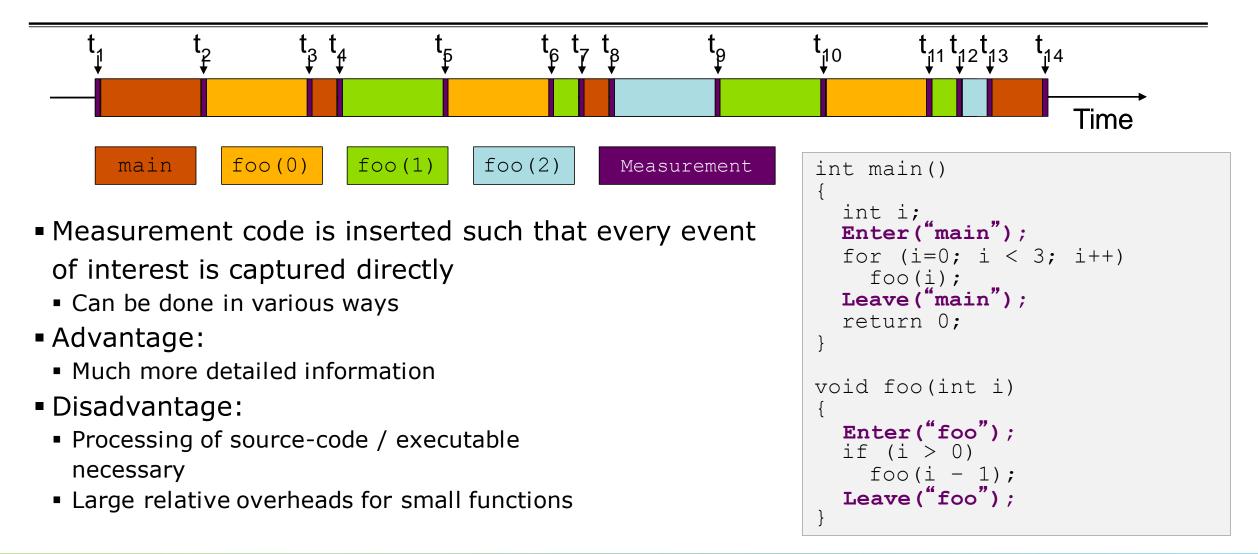
• How are performance measurements triggered?

- Sampling
- Code instrumentation
- How is performance data recorded?
 - Profiling / Runtime summarization
 - Tracing
- How is performance data analyzed?
 - Online
 - Post mortem

Sampling



Instrumentation



Instrumentation techniques

- Static instrumentation
 - Program is instrumented prior to execution
- Dynamic instrumentation
 - Program is instrumented at runtime
- Code is inserted
 - Manually
 - Automatically
 - By a preprocessor / source-to-source translation tool
 - By a compiler
 - By linking against a pre-instrumented library / runtime system
 - By binary-rewrite / dynamic instrumentation tool

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Critical issues

Accuracy

- Intrusion overhead
 - Measurement itself needs time and thus lowers performance
- Perturbation
 - Measurement alters program behaviour
 - E.g., memory access pattern
- Accuracy of timers & counters
- Granularity
 - How many measurements?
 - How much information / processing during each measurement?

Tradeoff: Accuracy vs. Expressiveness of data

Classification of measurement techniques

- How are performance measurements triggered?
 - Sampling
 - Code instrumentation
- How is performance data recorded?
 - Profiling / Runtime summarization
 - Tracing
- How is performance data analyzed?
 - Online
 - Post mortem

Profiling / Runtime summarization

- Recording of aggregated information
 - Total, maximum, minimum, ...
- For measurements
 - Time
 - Counts
 - Function calls
 - Bytes transferred
 - Hardware counters
- Over program and system entities
 - Functions, call sites, basic blocks, loops, ...
 - Processes, threads

Profile = summarization of events over execution interval

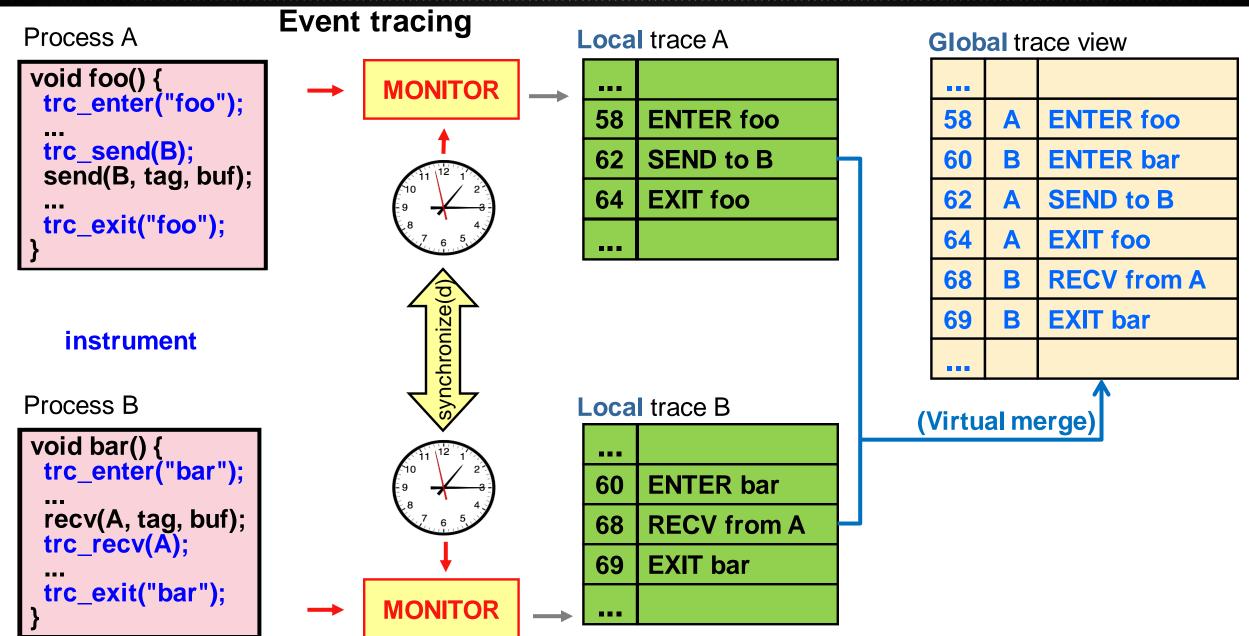
Types of profiles

- Flat profile
 - Shows distribution of metrics per routine / instrumented region
 - Calling context is not taken into account
- Call-path profile
 - Shows distribution of metrics per executed call path
 - Sometimes only distinguished by partial calling context (e.g., two levels)
- Special-purpose profiles
 - Focus on specific aspects, e.g., MPI calls or OpenMP constructs
 - Comparing processes/threads

Tracing

- Recording detailed information about significant points (events) during execution of the program
 - the program
 - Enter / leave of a region (function, loop, ...)
 - Send / receive a message, ...
- Save information in event record
 - Timestamp, location, event type
 - Plus event-specific information (e.g., communicator, sender / receiver, ...)
- Abstract execution model on level of defined events
- Event trace = Chronologically ordered sequence of event records

VI-HPS



Tracing Pros & Cons

- Tracing advantages
 - Event traces preserve the temporal and spatial relationships among individual events (@ context)
 - Allows reconstruction of dynamic application behavior on any required level of abstraction
 - Most general measurement technique
 - Profile data can be reconstructed from event traces
- Disadvantages
 - Traces can very quickly become extremely large
 - Writing events to file at runtime may causes perturbation

Classification of measurement techniques

- How are performance measurements triggered?
 - Sampling
 - Code instrumentation
- How is performance data recorded?
 - Profiling / Runtime summarization
 - Tracing
- How is performance data analyzed?
 - Online
 - Post mortem

Online analysis

- Performance data is processed during measurement run
 - Process-local profile aggregation
 - Requires formalized knowledge about performance bottlenecks
 - More sophisticated inter-process analysis using
 - "Piggyback" messages
 - Hierarchical network of analysis agents
- Online analysis often involves application steering to interrupt and re-configure the measurement

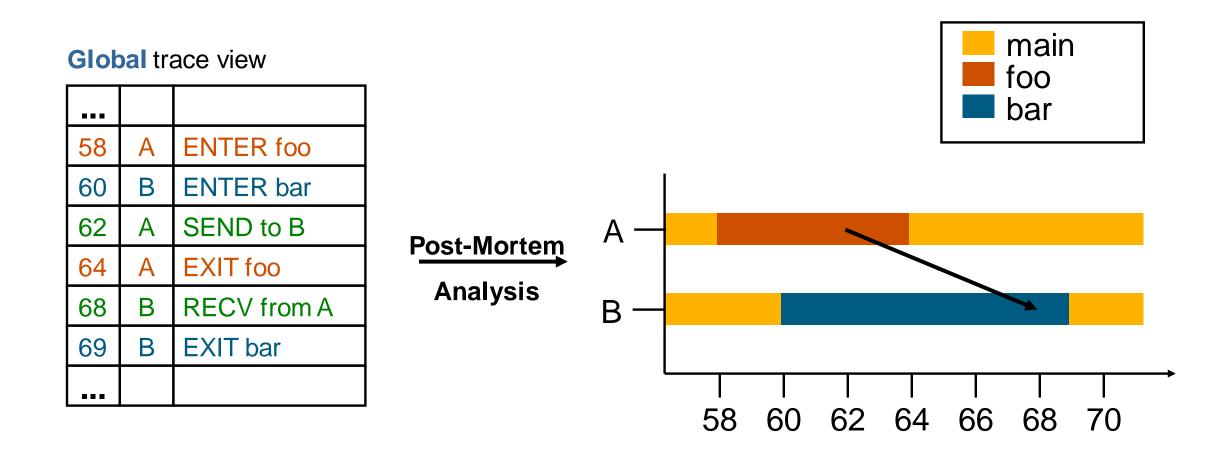


Post-mortem analysis

- Performance data is stored at end of measurement run
- Data analysis is performed afterwards
 - Automatic search for bottlenecks
 - Visual trace analysis
 - Calculation of statistics

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Example: Time-line visualization



No single solution is sufficient!

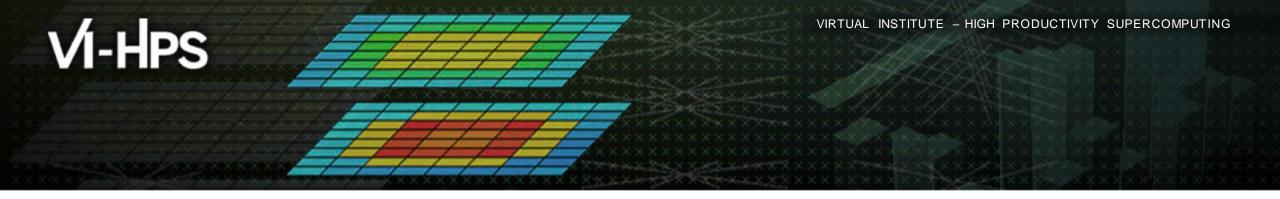


A combination of different methods, tools and techniques is typically needed!

- Analysis
 - Statistics, visualization, automatic analysis, data mining, ...
- Measurement
 - Sampling / instrumentation, profiling / tracing, ...
- Instrumentation
 - Source code / binary, manual / automatic, …

Typical performance analysis procedure

- Do I have a performance problem at all?
 - Time / speedup / scalability measurements
- What is the key bottleneck (computation / communication)?
 - MPI / OpenMP / flat profiling
- Where is the key bottleneck?
 - Call-path profiling, detailed basic block profiling
- Why is it there?
 - Hardware counter analysis, trace selected parts to keep trace size manageable
- Does the code have scalability problems?
 - Load imbalance analysis, compare profiles at various sizes function-by-function



JUWELS Booster & JupyterJSC/Xpra

Brian Wylie



JUWELS-Booster (juwels-booster@fz-juelich.de)



JUWELS-Booster system overview

936 accelerated compute nodes

- dual AMD EPYC 7402 processors, each with 24 cores, 2.8 GHz, 512 GB DDR4 RAM
- quad Nvidia A100 'Ampere' GPUs with 40GB HBM
- 4 login nodes (and additional service nodes)
 - dual AMD EPYC 7402 processors, each with 24 cores, 2.8 GHz, 512 GB DDR4 RAM
- Mellanox InfiniBand HDR full fat-tree interconnect network
- All compute nodes are diskless: OS image in RAM
- IBM Spectrum Scale (GPFS) parallel file system connection to JUST storage & HPST

Software environment

- Rocky Linux 8 OS
- Batch system workload/resource management based on SLURM from ParTec
- Programming environment
 - GNU, Intel & NVHPC compilers (for C, C++ & Fortran)
 - all supporting OpenMP and other multithreading
 - ParaStation MPI (based on MPICH3) & OpenMPI
 - Optimized mathematical libraries (Intel MKL, etc.)

Accessing software

- Hierarchical modules: `toolchain' constructed by loading compiler then MPI
- List available modules (ready to be loaded)module avail
- Search for an application/library/tool
 - module spider <name>
- Load the desired compiler
 - module load GCC
- Load the desired MPI
 - module load ParaStationMPI
- Load additional applications/libraries/tools
 - module load Score-P Scalasca Vampir

- List currently loaded modulesmodule list
- Purge all loaded modulesmodule purge
- Unload an undesired compiler
 module unload CUDA
- Save current collection of modules
 - module save [<name>]
- Restore a saved collection of modules
 - module restore [<name>]

Filesystems

- \$HOME (/p/home/\$USER)
 - private, small storage quota for each user account, poor parallel I/O performance
- \$PROJECT (/p/project)
 - shared by project members, regular backup, optimized parallel I/O performance
- \$SCRATCH (/p/scratch)
 - shared by project members, no backup, optimized parallel I/O performance, automatic purge based on last file access

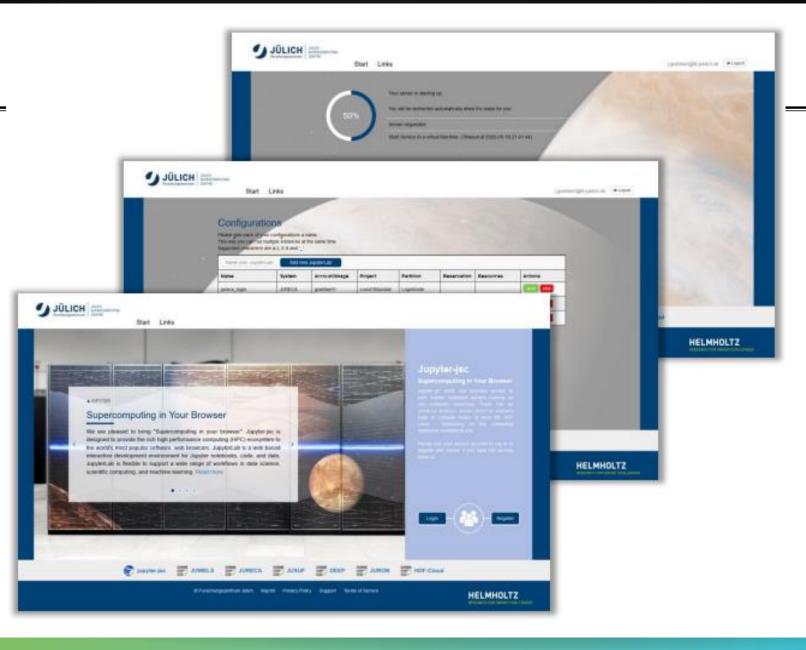
\$ARCHIVE

• ...

SLURM

- Show available partitions
 - sinfo
- Show queued jobs of user
 - squeue -u \$USER
- Cancel queued/running job
 - scancel jobid
- Submit batch script (to partition)
 - sbatch <script.sh>
 - within script **always** use srun (rather than mpiexec/mpirun) to launch application on allocated compute nodes
 - srun mpi_app.exe

Jupyter-JSC webservice



Registration

Go to training project join link in web browser

https://judoor.fz-juelich.de/projects/join/training2341

- Register for JSC HPC (JuDoor) account (if don't already have one)
 - Provide Email address and other details
 - Accept usage agreement
 - Submit the registration
 - Wait for registration Email and follow link to confirm

Portal for managing accounts, projects and resources at JSC.

Register for account

HPS

- Provide Email address
- Wait for Email with URL for account registration
- Enter name/affiliation/etc
- Accept usage agreement
- Submit the registration
- Wait for registration Email and follow link to confirm

To: Subject: J	Unity-isc@fz-juelich.de Jupyter-JSC Registration Tue, 19 May 2020 11:09:53 +0200
	r, all address was entered into the Jupyter-JSC authentication service and must be confirmed. Afterward, you have to log in again. our e-mail address.
If you did	not use your JuDoor account to log into https://jupyter-jsc.fz-juelich.de, we recommend that you change your JuDoor passwor

JuDoor Login	× +	✓ - □ X
\leftrightarrow \rightarrow C \bullet judoor.fz-jud	elich.de/login	🕶 🖄 🎓 🖬 😩 :
JuDoo	or Login	JÜLICH Supercomputing Forschungszontrum

3 3	
Login using JSC account	Login with e-mail callback
Username	Login mail address
Password	A confirmation email to confirm your identity will be sent to this address.
Login Register Reset password	Send identification mail

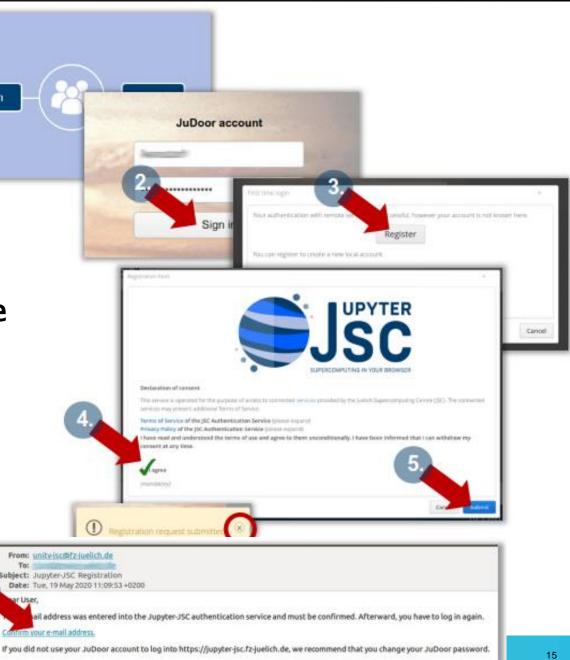
If you are stuck take a look at the 😟	JuDoor Documentation.
---------------------------------------	-----------------------

Legal Notice	Forschungszentrum Jülich, JSC	Contact Support
Privacy Policy	Forschungszehltum Julich, JSC	JuDoor Requests

Login

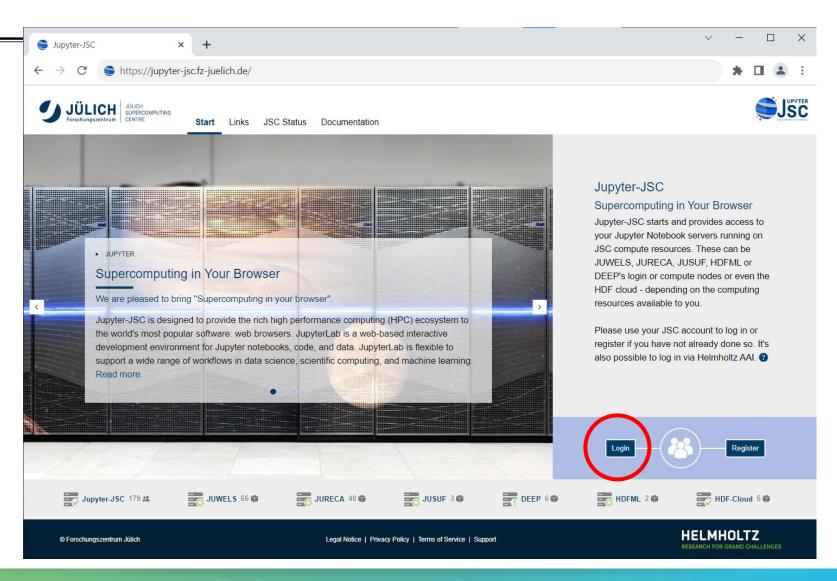
First time login

- Pre-requisite: JuDoor registration 0. with project membership and accepted systems usage agreement
- 1. Go to link in web browser https://jupyter-jsc.fz-juelich.de
- 2. Sign in with JSC HPC account (JuDoor)
- 3. Register to Jupyter-JSC
- 4. Accept usage agreement
- 5. Submit the registration
- 6. Wait for email and confirm



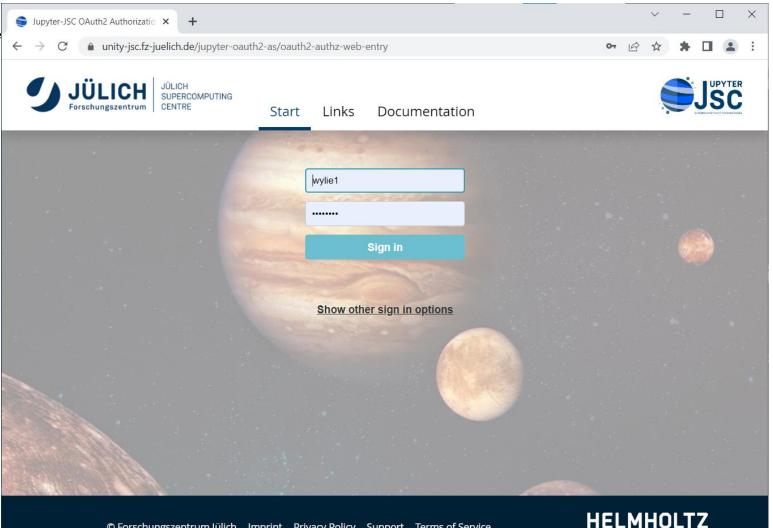
Login

• Click on Login button to enter



Sign in

• Provide JSC HPC account credentials to sign in



ND CHALLENGES

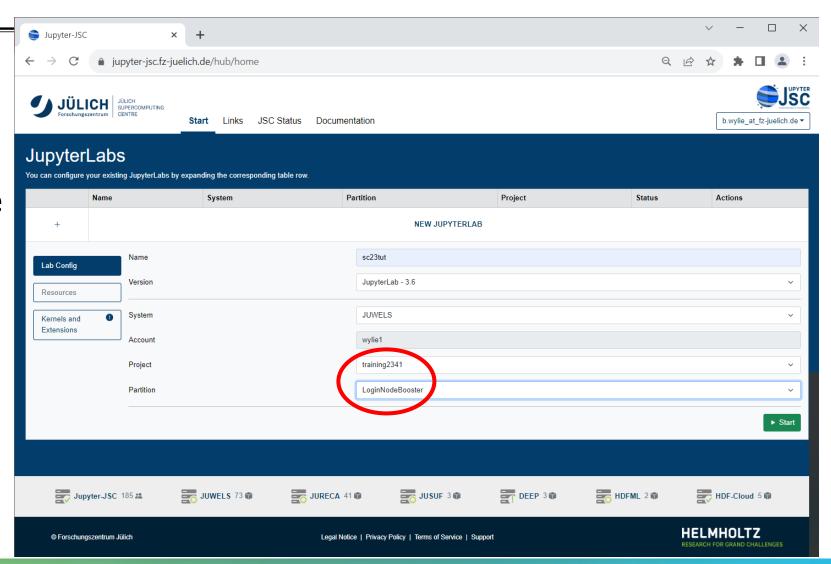
Configure

• Click on New JupyterLab configure a new service

Supyter-JSC	c ×	+					× -	
← → C		elich.de/hub/home				(2 12 12 14 1	: .
JÜL Forschungt	ICH JÜLICH SUPERCOMPUTING CENTRE	Start Links JSC S	Status Documentation				b.wylie_at_	fz-juelich.de •
Jupyter		expanding the corresponding t	able row.					
	Name	System	Partition		Project	Status	Actions	
+				NEW JUPYTER AB				
Tup	oyter-JSC 179 😃	JUWELS 66 🏟	JURECA 40 🏟	JUSUF 3 📦	DEEP 6	HDFML 2 🇊	HDF-Cloud	5 🎲
© Forschung	gszentrum Jülich		Legal Notice Privacy	Policy Terms of Service Supp	port		HELMHOLTZ RESEARCH FOR GRAND CHAI	

Define service

- Provide name of your choice for new lab configuration
- Configure from available (limited) options
 - Select partition
 LoginNodeBooster
- ... and Start



Customize extensions



- For our tutorial, we only need Xpra desktop
 - could deselect everything else

Supyter-JSC	× 🔵 JupyterLab (auto-a	a) × +			✓ – □ X
\leftrightarrow \rightarrow C $rightarrow$ jupyter-je	sc.fz-juelich.de/hub/home			Q 🖻	☆ 🛪 🖬 🏝 🗄
UILCH Forschungszentrum		Status Documentation			b.wylie_at_fz-juelich.de V
JupyterLabs You can configure your existing Jupyte	rLabs by expanding the corresponding	table row.			
Name	System	Partition	Project	Status	Actions
\frown			NEW JUPYTERLAB		
∽ sc23tut	JUWELS	LoginNode	training2341	0	► Start
Lab Config NG Resources Kern		Slurm Wrapper ()	☐ Jupyter-Al ⑦		
Extensions	DeepLearning ()) j iles	☐ Julia ① ☐ PyEarthSystem ⑦ ☐ Ruby ⑦	☐ LFortran ① ☐ PyQuantum ① ☐ Bash ①	Octave ① PyVisualization ①	
Logs	ST Desktop 🕧	🗹 Xpra 🕧			
□ Sel	ect all	Deselect all			
E Sa	ave C Reset Successfully upda	ted sc23tut. 🗙			T Delete
Jupyter-JSC 179 👪	JUWELS 70 📦	TINECA 40 🗊	JUSUF 3 🗊 🛛 📰 DEEP 6 🗊	HDFML 2 🗊	HDF-Cloud 5
© Forschungszentrum Jülich		Legal Notice Privacy f	Policy Terms of Service Support		MHOLTZ

Starting

• Now wait for server starting up ...

Supyter-JSC × Supyter-JSC × +	✓ - □ ×
$\leftarrow \rightarrow C$ a jupyter-jsc.fz-juelich.de/hub/spawn-pending/b.wylie_at_fz-juelich.de/i22e8df6732543f99ba726d1ebc845ea Q	☞ ☆ 🛊 🖬 😩 :
Start Links JSC Status Documentation	b.wylie_at_fz-juelich.de •
Your server is starting up You will be redirected automatically when it's ready for you.	
Lab Info (click to expand)	~
10% spawning	Cancel
2023_10_25 14:33:03.182: Sending request to backend service to start your service.	+
Jupyter-JSC 179 # JUWELS 67 🏶 📻 JURECA 40 🕸 📻 JUSUF 3 🏶 📻 DEEP 6 🏶 📻 HDFML 2 🏶	HDF-Cloud 5 🍘
© Forschungszentrum Jülich Legal Notice Privacy Policy Terms of Service Support	HELMHOLTZ RESEARCH FOR GRAND CHALLENGES

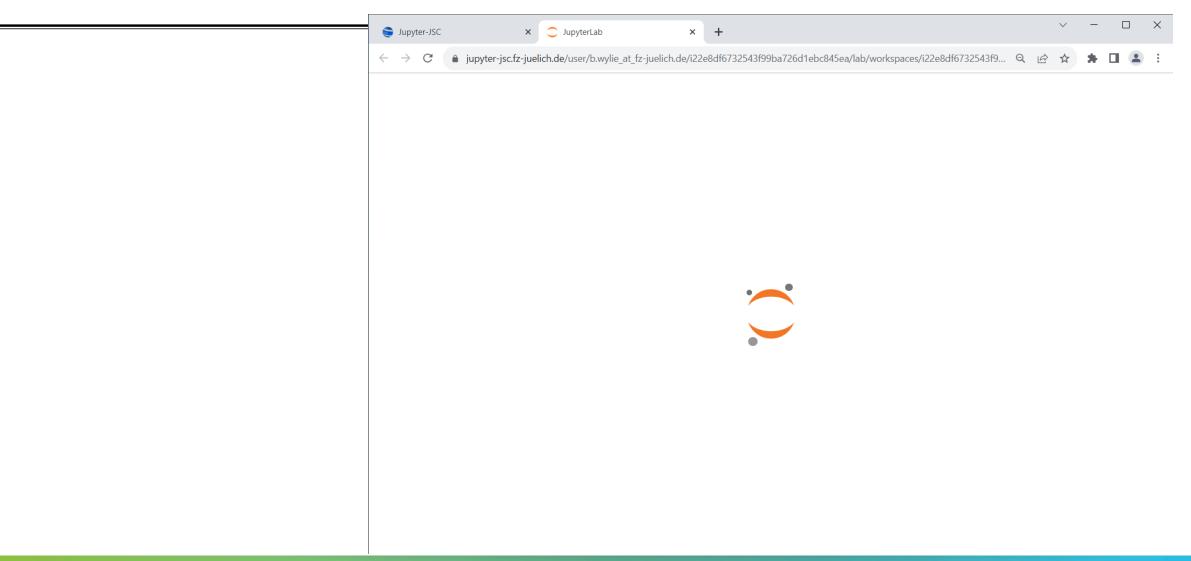
... still starting ...

• Browser tab created for new service

Supyter-JSC × Supyter-JSC	× +		✓ - □ ×
← → C 🌲 jupyter-jsc.fz-juelich.ae/indb/spann-pending/b.wy	ie_at_fz-juelich.de/i22e8df6732543f99ba726d1ebc845ea	Q	
SUBERCOMPUTING Forschungszentrum SUBERCOMPUTING CENTRE Start Links JSC Status Do	cumentation		b.wylie_at_fz-juelich.de •
Your server is starting up You will be redirected automatically when it's ready for you.			
Lab Info (click to expand)			~
	95%		
	spawning		Cancel
2023_10_25 14:43:25.465: Sending request to backend service to start your service.			+
2023_10_25 14:43:33.500: Backend communication successful.			+
2023_10_25 14:43:36.716: Setup ssh port-forwarding.			+
2023_10_25 14:43:36.837: Disk quota checked.			+
2023_10_25 14:43:36.969: Load default modules			+
2023_10_25 14:43:53.416: Load default modules done			+
2023_10_25 14:43:53.840: Add system specific configuration.			+
2023_10_25 14:43:54.872: Start JupyterLab			+
Jupyter-JSC 179 🎿 📑 JUWELS 67 🗊 📑 JUR	ECA 40 T JUSUF 3 T DEEP 6 T	THORE 2 THE PROPERTY OF THE PR	HDF-Cloud 5 🗊
© Forschungszentrum Jülich L	egal Notice Privacy Policy Terms of Service Support		HELMHOLTZ RESEARCH FOR GRAND CHALLENGES

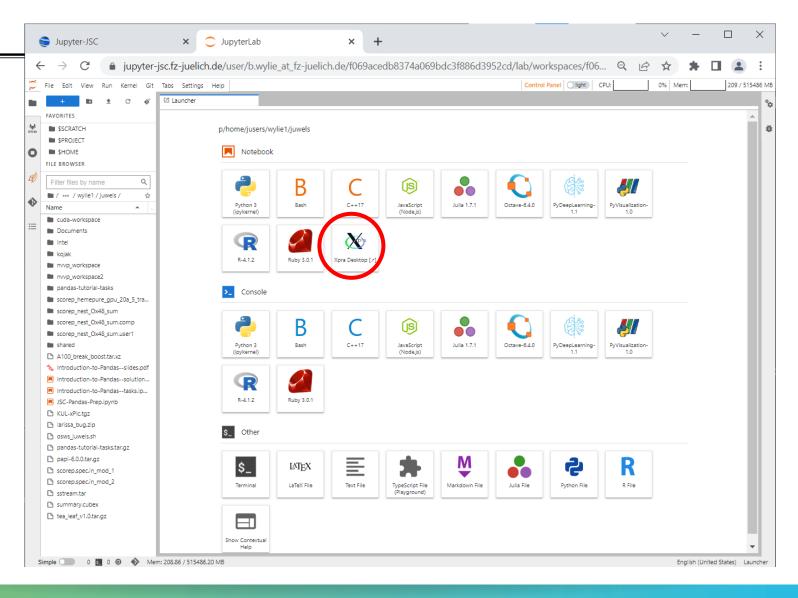
VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

... almost there



Lab service running!

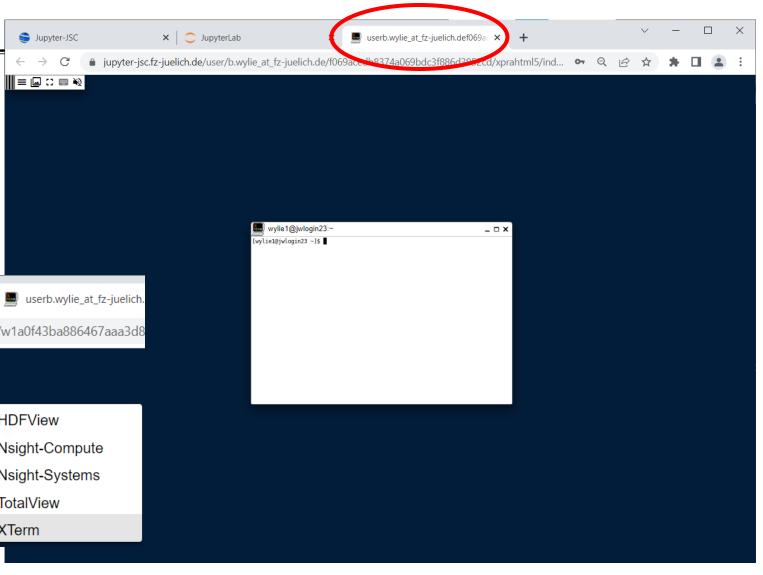
• Click on Xpra Desktop icon



Xpra desktop

- Opens new browser tab containing virtual desktop
- Xterm opened by default
 - if not, use Start menu

🤤 Jupyter-JSC	×	🔵 JupyterLab		× 📕 userb.wylie_at_fz-juelich.
← → C 🌲 jupyt	er-jsc.fz-	juelich.de/user/b.wylie_at	_fz-jue	elich.de/w1a0f43ba886467aaa3d8
= 🖬 🗆 🖦				
III Start	>	Performance Too	ls >	
E Server	>	Tools	>	HDFView
 Information 	>	Visualization Too	s >	Nsight-Compute
C Reload				Nsight-Systems
➔ Disconnect				💵 TotalView
				XTerm





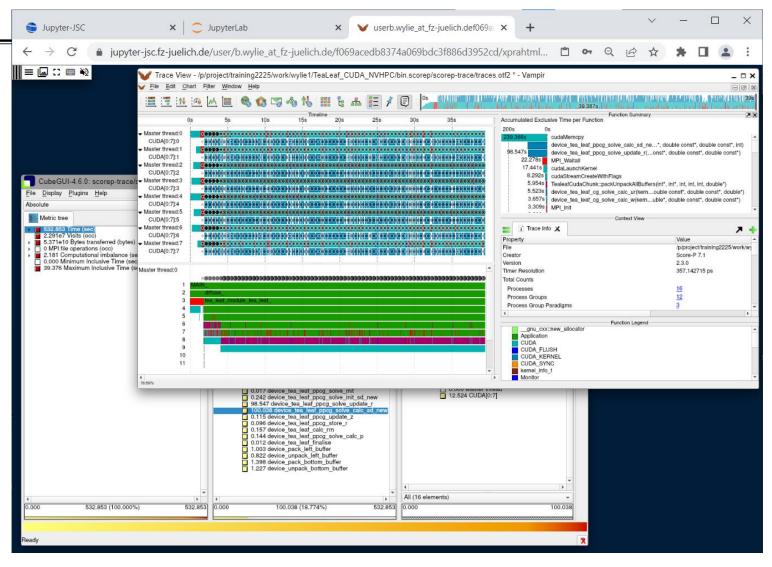
Source setup

• For tutorial, source the provided shell environment setup script

wylie1@jwlogin05:/p/scratch/training2341/wylie1	□ ×
<pre>[wyliel@jwlogin05 wyliel]\$ source /p/project/training2341/setup.sh SC23 tutorial training account setup SC23 tutorial training account wyliel set up for toolchain NVHPC(PGI)+ParaStationMP Changing to WORK directory: /p/scratch/training2341/wyliel [wyliel@jwlogin05 wyliel]\$</pre>	ч ч

Start tools

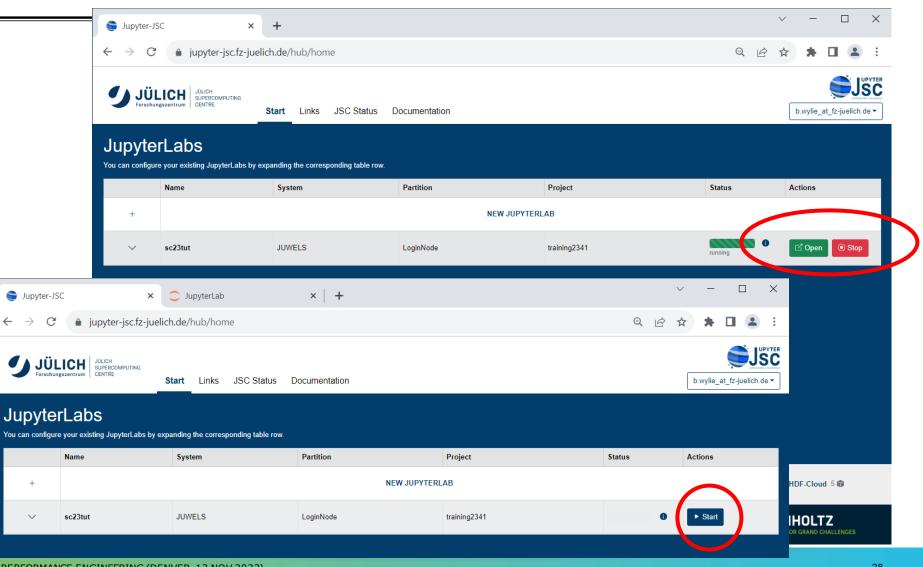
 Graphical tools are windows within the virtual desktop



Manage service

 Stop or re-open existing service

• Re-start service when desired



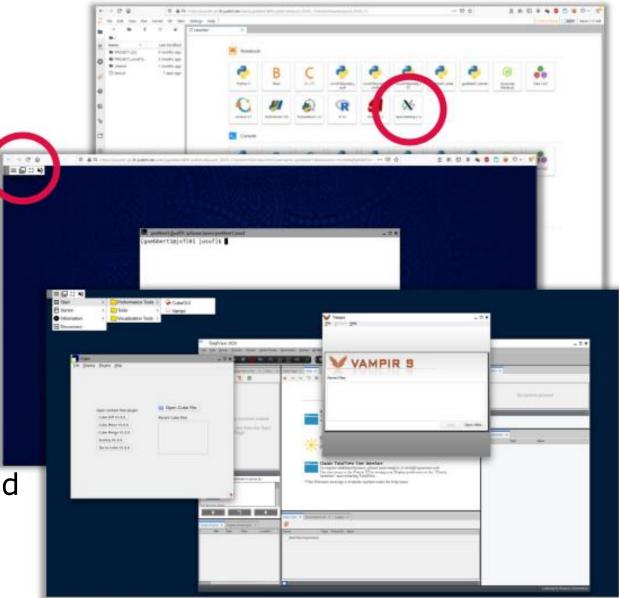
 \leftarrow

Xpra – remote desktop access

- X Persistent Remote Applications
- runs X clients on a remote host and directs display to local machine
- runs in browser
- reconnection without disrupting the forwarded application session

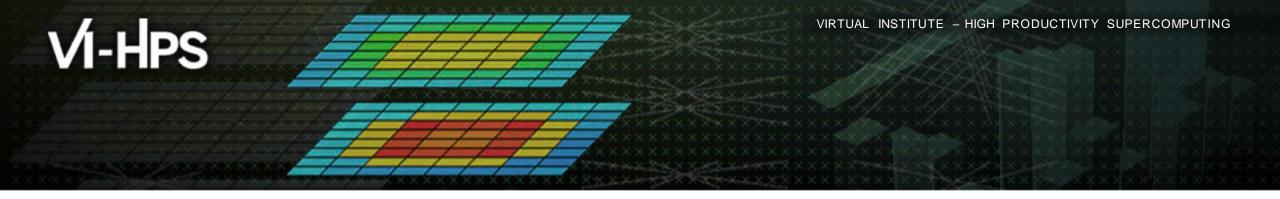
Remote desktop runs on the same node as JupyterLab does

- killed when JupyterLab session stopped
- refresh browser tab if connection lost



https://jupyter-jsc.fz-juelich.de





Hands-on: JUWELS Booster (AMD EPYC Rome + 4 x A100) TeaLeaf_CUDA

VI-HPS Team



Tutorial exercise objectives

- Familiarise with usage of VI-HPS tools
 - complementary tools' capabilities & interoperability
- Prepare to apply tools productively to your applications(s)
- Exercise is based on a small portable benchmark code
 - unlikely to have significant optimisation opportunities
- Optional (recommended) exercise extensions
 - analyse performance of alternative configurations
 - investigate effectiveness of system-specific compiler/MPI optimisations and/or placement/binding/affinity capabilities
 - investigate scalability and analyse scalability limiters
 - compare performance on different HPC platforms

• ...

Setup for exercises

- Connect to your training account on JUWELS Booster (with X11-forwarding)
 - % ssh -X <yourid>@juwels-booster.fz-juelich.de

Not needed with JupyterJSC!

Set account and default environment (NVHPC + ParaStationMPI) via helper script

% source /p/project/training2341/setup.sh

Needed in each new shell

Copy tutorial sources to your WORK directory

% cd \$WORK
% tar zxvf \$PROJECT/examples/tea_leaf.tar.gz
% cd TeaLeaf CUDA

WORK=/p/scratch/training2341/\$USER PROJECT=/p/project/training2341

Case study: TeaLeaf_CUDA

- HPC mini-app developed by the UK Mini-App Consortium
 - Solves the linear 2D heat conduction equation on a spatially decomposed regular grid using a 5 point stencil with implicit solvers

- Part of the Mantevo 3.0 suite
- Available on GitHub: https://uk-mac.github.io/TeaLeaf/
- CUDA-enabled MPI version written in Fortran90, run using default testcase
 - Optional OpenMP (only used during initialization): vary the number of threads for each MPI process
 - Run with 4 MPI tasks-per-node so that each has a dedicated GPU
 - Run on 1 or more nodes to see its strong scaling performance: 2 is sufficient for exercise
 - Experiment with different MPI libraries, compilers & compiler optimisations
 - Experiment with different bindings/affinities of MPI processes and OpenMP threads
 - Experiment with different solvers (and other testcases)
- Provided version of Makefile and sources customized for this tutorial
 - builds tea_leaf executable in separate directory when using instrumentation

TeaLeaf_CUDA source directory

% **ls**

Benchmarks/ Makefile

README.md build_field.f90 calc_dt.f90 config/ cuda_common.hpp cuda_errors.cu cuda_strings.cu cuda_strings.hpp data.f90 definitions.f90 diffuse.f90 field_summary.f90 field_summary_kernel_cuda.cu ftocmacros.h

generate chunk.f90 generate_chunk_kernel_cuda.cu global mpi.f90 host_reductions kernel cuda.cu init cuda.cu initialise.f90 initialise chunk.f90 initialise chunk kernel cuda.cu jobscripts/ kernel files/ makefile.deps pack kernel cuda.cu parse.f90 read input.f90 report.f90 set field.f90

set field kernels cuda.cu start.f90 tea.f90 tea.in tea leaf.f90 tea leaf cq.f90 tea leaf cheby.f90 tea leaf common.f90 tea leaf kernel cuda.cu tea leaf ppcq.f90 tea solve.f90 timer.f90 timer c.c timestep.f90 update halo.f90 update halo kernel cuda.cu

25 Fortran90 modules, 1 C module, 10 CUDA modules

TeaLeaf_CUDA: Makefile

```
#Crown Copyright 2014 AWE
 This file is part of TeaLeaf.
#
 TeaLeaf is free software...
 Agnostic, platform independent Makefile for the TeaLeaf benchmark code.
 It is not meant to be clever in any way, just a simple build script.
  this works as well:-
#
                                                                                 Specify the suite of compilers
#
 make COMPILER=GNU [OPENMP=1]
#
                                                                                  (and optionally OpenMP)
. . .
#PREP=scorep --cuda
                                                                                 No instrumentation by default
MPI COMPILER=$(PREP) mpifort
C MPI COMPILER=$(PREP) mpicc
# No preposition for CXX MPI COMPILER!
CXX MPI COMPILER=mpic++
NVC\overline{C}=$(\overline{P}REP) nvcc -ccbin $(CXX MPI COMPILER)
. . .
```

Building tea_leaf

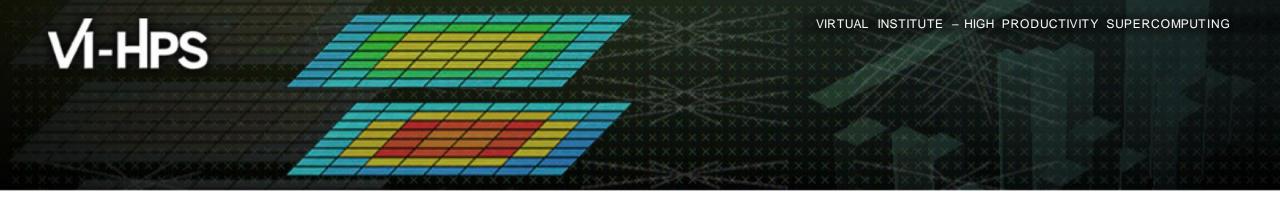
```
% make COMPILER=PGI
mpif90 -fastsse -qopt -q -c data.f90 -o data.o
[ ... ]
mpif90 -fastsse -gopt -g -c tea leaf.f90 -o tea leaf.o
mpif90 -fastsse -gopt -g -c diffuse.f90 -o diffuse.o
mpicc -fastsse -gopt -g -c timer c.c -o timer c.o
nvcc -ccbin mpicxx -std=c++14 -I/p/software/juwelsbooster/stages/2024/software/CUDA/12/include \
 -gencode arch=compute 80, code=sm 80 -restrict -Xcompiler "-fastsse -gopt -c -g" -DNO ERR CHK \
 -O3 -c cuda errors.cu -o cuda errors.o
[ ... ]
mpif90 -fastsse -qopt -q \
 data.o definitions.o global mpi.o tea.o report.o timer.o parse.o read input.o initialise chunk.o \
 build field.o update halo.o start.o generate chunk.o initialise.o field summary.o calc dt.o \setminus
 timestep.o set field.o tea leaf common.o tea leaf cg.o tea leaf cheby.o tea leaf ppcg.o \
 tea leaf jacobi.o tea solve.o tea leaf.o diffuse.o \
 timer c.o
 cuda errors.o cuda strings.o field summary kernel cuda.o generate chunk kernel cuda.o init cuda.o \
 initialise chunk kernel cuda.o pack kernel cuda.o set field kernel cuda.o tea leaf kernel cuda.o \
 update halo kernel cuda.o \setminus
 -L/p/software/juwelsbooster/stages/2024/software/CUDA/12/lib64 \
 -1stdc++ -1cudart \setminus
 -o bin/tea leaf
```

TeaLeaf_CUDA jobscript for reference execution

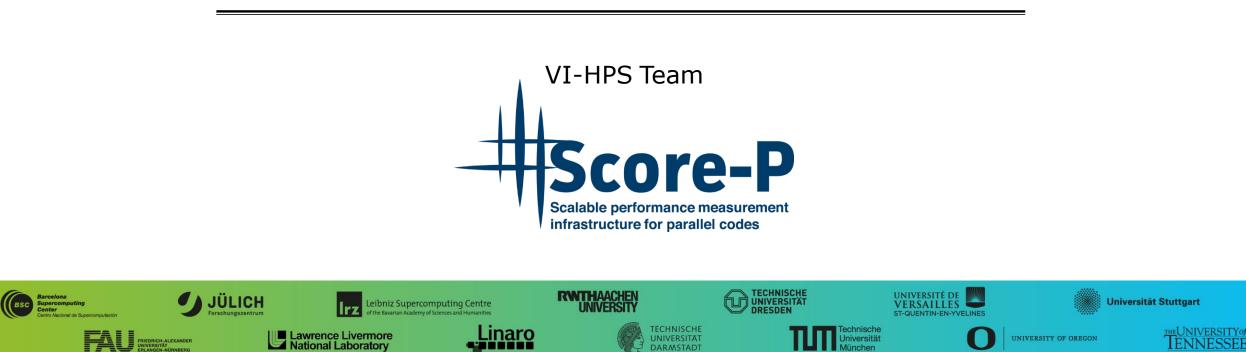
 Copy jobscript and % cd bin % cp ../jobscripts/juwelsbooster/run.sbatch . check/modify its ⁹ cat run.sbatch contents #!/bin/bash #SBATCH --job-name=TeaLeaf then submit # Job name #SBATCH --nodes=2 # Total number of nodes requested # Number of GPUs per node #SBATCH --gres=gpu:4 #SBATCH --ntasks-per-node=4 # Number of MPI tasks per node (one per GPU) #SBATCH --time=00:05:00 # Max. wall-clock time (hh:mm:ss) #SBATCH --account=training2341 # Project account to be charged #SBATCH --output=%x.%j.out # Output files #SBATCH --error=%x.%j.out #SBATCH --partition=develbooster # Job partition srun ./tea leaf % sbatch run.sbatch

TeaLeaf_CUDA Reference Execution

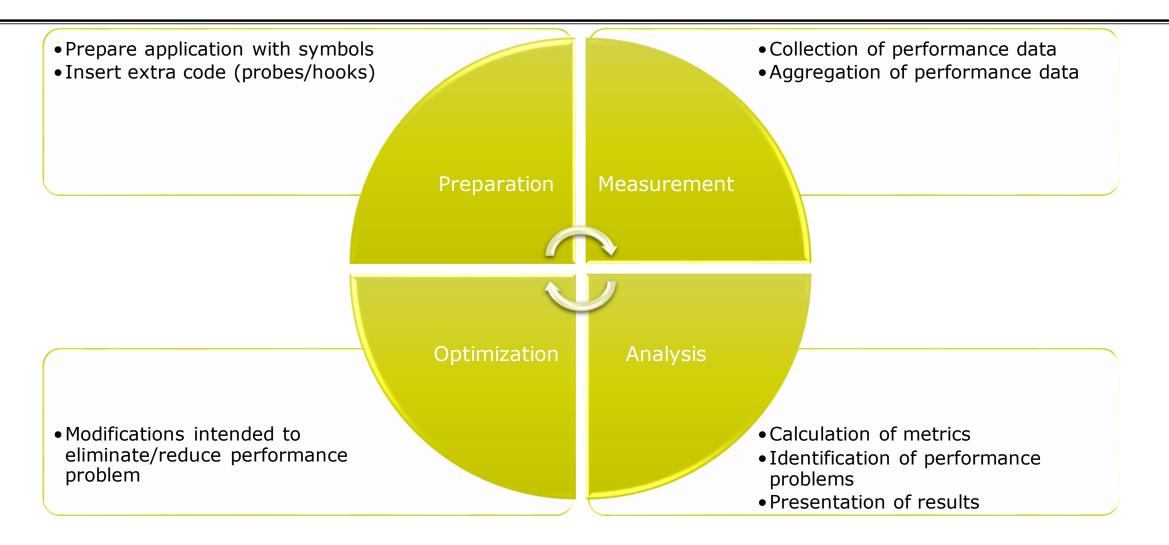
<pre>% cat TeaLeaf.<job_id>.out</job_id></pre>	Verify the reported
Tea Version 1.400	execution
MPI Version Task Count 8	configuration and
Output file tea.out opened. All output will go there.	that the test
CUDA in rank 1 using jwb0149.juwels GPU device 0 NVIDIA A100-SXM4-40GB (0:68:0) CUDA in rank 3 using jwb0149.juwels GPU device 0 NVIDIA A100-SXM4-40GB (0:196:0)	execution passed
CUDA in rank 2 using jwb0149.juwels GPU device 0 NVIDIA A100-SXM4-40GB (0:132:0)	
CUDA in rank 0 using jwb0149.juwels GPU device 0 NVIDIA A100-SXM4-40GB (0:3:0) Solver to use: PPCG	
Preconditioner to use: None	
Step 1 time 0.0000000 timestep 4.00E-03	
CUDA in rank 5 using jwb0150.juwels GPU device 0 NVIDIA A100-SXM4-40GB (0:68:0)	
CUDA in rank 7 using jwb0150.juwels GPU device 0 NVIDIA A100-SXM4-40GB (0:196:0)	
CUDA in rank 6 using jwb0150.juwels GPU device 0 NVIDIA A100-SXM4-40GB (0:132:0) CUDA in rank 4 using jwb0150.juwels GPU device 0 NVIDIA A100-SXM4-40GB (0:3:0) Hint save	the henchmark output
	the benchmark output
This test is considered PASSED (Or note the	e run time) to be able to
First step overhead -0.3670756816864014 refer to it la	ater
Wall clock 33.43737912178040	



Score-P – A Joint Performance Measurement Run-Time Infrastructure for Scalasca, TAU, and Vampir



Performance engineering workflow



Score-P



- Infrastructure for instrumentation and performance measurements
- Instrumented application can be used to produce several results:
 - Call-path profiling: CUBE4 data format used for data exchange
 - Event-based tracing: OTF2 data format used for data exchange
 - Online profiling: In conjunction with the Periscope Tuning Framework
- Supported parallel paradigms:
 - Multi-process: MPI, SHMEM
 - Thread-parallel: **OpenMP**, Pthreads
 - CUDA, OpenCL, OpenACC, HIP Accelerator-based:
- Open Source; portable and scalable to all major HPC systems
- Initial project funded by BMBF
- Close collaboration with PRIMA project funded by DOE

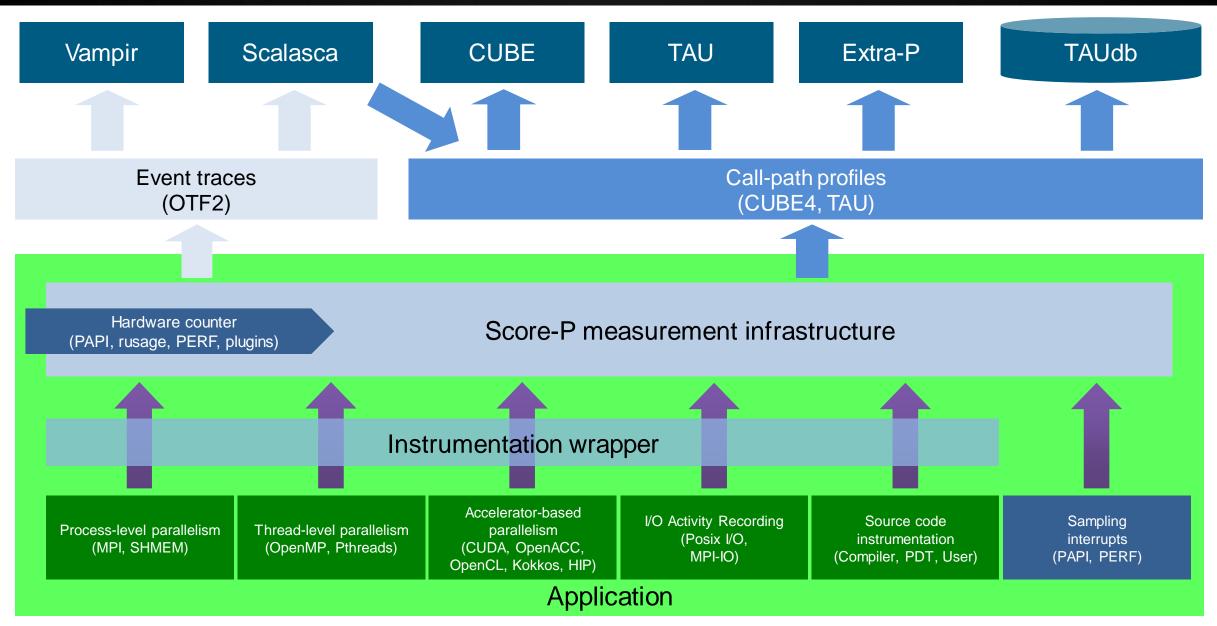
GEFÖRDERT VOM

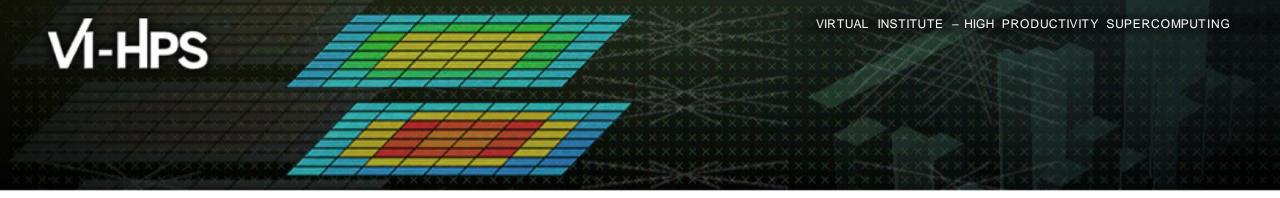


Bundesministerium für Bildung und Forschung



VI-HPS





Hands-on: TeaLeaf_CUDA





Performance analysis steps

• 0.0 Reference preparation for validation

- 1.0 Program instrumentation
- 1.1 Summary measurement collection
- 1.2 Summary analysis report examination
- 2.0 Summary experiment scoring
- 2.1 Summary measurement collection with filtering
- 2.2 Filtered summary analysis report examination

3.0 Event trace collection

3.1 Event trace examination & analysis

Local installation (JUWELS Booster)

Set account and default environment (NVHPC + ParaStationMPI) via helper script:

% source /p/project/training2341/setup.sh

Load the modules for the tool environment:

% module load Score-P CubeGUI

Copy tutorial sources to your WORK directory (or your personal workspace)

Only required if not done already (for opening exercise)

```
% cd $WORK
% tar xf $PROJECT/examples/tea_leaf.tar.gz
% cd TeaLeaf CUDA
```

TeaLeaf_CUDA: Makefile

```
#Crown Copyright 2014 AWE
 This file is part of TeaLeaf.
 TeaLeaf is free software...
 Agnostic, platform independent Makefile for the TeaLeaf benchmark code.
 It is not meant to be clever in any way, just a simple build script.
 this works as well:-
                                                                              Specify the suite of compilers
#
 make COMPILER=PGI [OPENMP=1]
#
                                                                              (and optionally OpenMP)
. . .
#PREP=scorep --cuda
                                                                              No instrumentation by default
MPI COMPILER=$(PREP) mpifort
C MPI COMPILER=$(PREP) mpicc
# No preposition for CXX MPI COMPILER!
                                                                               Uncomment or set PREP
CXX MPI COMPILER=mpic++
NVCC=$(PREP) nvcc -ccbin $(CXX MPI COMPILER)
                                                                               to instrumenter preposition
. . .
```

Instrumenting tea_leaf

```
% make COMPILER=PGI PREP="scorep --cuda"
scorep --cuda mpif90 -fastsse -gopt -g -c data.f90 -o data.o
[...]
mpicc -fastsse -gopt -g -c timer c.c -o timer c.o
scorep --cuda nvcc -ccbin mpicxx -I/p/software/juwelsbooster/stages/2024/software/CUDA/12/include \
 -std=c++14 -gencode arch=compute 80, code=sm 80 -restrict -Xcompiler "-fastsse -gopt -c -g" \
 -DNO ERR CHK -O3 -c cuda errors.cu -o cuda errors.o
[...]
scorep --cuda mpif90 -fastsse -gopt
                                      -a /
data.o definitions.o global mpi.o tea.o report.o timer.o parse.o read input.o initialise chunk.o \
build field.o update halo.o start.o generate chunk.o initialise.o field summary.o calc dt.o \
timestep.o set field.o tea leaf common.o tea leaf cg.o tea leaf cheby.o tea leaf ppcg.o \
tea leaf jacobi.o tea solve.o tea leaf.o diffuse.o timer c.o
cuda errors.o cuda strings.o field summary kernel cuda.o generate chunk kernel cuda.o init cuda.o \
initialise chunk kernel cuda.o pack kernel cuda.o set field kernel cuda.o tea leaf kernel cuda.o \
update halo kernel cuda.o \setminus
-L/p/software/juwelsbooster/stages/2024/software/CUDA/12/lib64 \
-lstdc++ -lcudart \setminus
-o bin.scorep/tea leaf
```

Measurement configuration: scorep-info

```
% scorep-info config-vars --full
SCOREP ENABLE PROFILING
 Description: Enable profiling
 [...]
SCOREP ENABLE TRACING
 Description: Enable tracing
 [...]
SCOREP TOTAL MEMORY
 Description: Total memory in bytes for the measurement system
 [...]
SCOREP EXPERIMENT DIRECTORY
 Description: Name of the experiment directory
 [...]
SCOREP FILTERING FILE
 Description: A file name which contain the filter rules
 [...]
SCOREP METRIC PAPI
 Description: PAPI metric names to measure
 [...]
SCOREP METRIC RUSAGE
 Description: Resource usage metric names to measure
 [...]
SCOREP CUDA ENABLE
 Description: CUDA measurement features
 [... More configuration variables ...]
```

 Score-P measurements are configured via environmental variables

Required for CUDA measurements. [yes|default] recommended to start with.

Summary measurement collection

% cd bin.scorep

% cp ../jobscripts/juwelsbooster/scorep.sbatch .
% cat scorep.sbatch

Score-P measurement configuration
export SCOREP_CUDA_ENABLE=default
export SCOREP_CUDA_BUFFER=48M

export SCOREP_EXPERIMENT_DIRECTORY=scorep-tea_leaf-8
#export SCOREP_FILTERING_FILE=../config/scorep.filter

#export SCOREP_ENABLE_TRACING=true
#export SCOREP TOTAL MEMORY=250M

Run the application
srun ./tea_leaf

\$ sbatch scorep.sbatch

- Change to the directory containing the new executable before running it with the desired configuration
- Check settings

Leave these lines commented out for the moment

Submit job

VIRTUAL VIRTUAL

TeaLeaf_CUDA Reference Execution

Verify the reported % cat TeaLeaf scorep.<job id>.out execution Tea Version 1.400 MPT Version configuration and OpenMP Version that the test Task Count: 8 execution passed Input read finished. Using CUDA Kernels [...] Solver to use: PPCG Preconditioner to use: None Test problem 6 is within 0.1397839E-05% of the expected solution This test is considered PASSED First step overhead -0.5253252983093262 Compare to previous reference Wall clock 38.01989197731018 execution without instrumentation

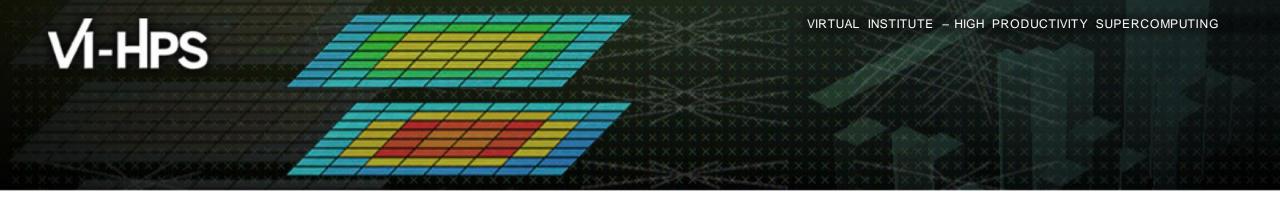
TeaLeaf summary analysis report examination

% ls tea_leaf tea.in TeaLeaf_scorep.<job id>.out scorep.sbatch scorep-tea leaf-8/ % ls scorep-tea leaf-8 MANIFEST.md profile.cubex scorep.cfg % cube scorep-tea leaf-8/profile.cubex [CUBE GUI showing summary analysis report] Hint: Copy 'profile.cubex' to local system (laptop) using 'scp' to improve responsiveness of GUI

- Creates experiment directory including
 - A brief content overview (MANIFEST.md)
 - A record of the measurement configuration (scorep.cfg)
 - The analysis report that was collated after measurement (profile.cubex)
- Interactive exploration with Cube

Further information

- Community instrumentation & measurement infrastructure
 - Instrumentation (various methods)
 - Basic and advanced profile generation
 - Event trace recording
 - Online access to profiling data
- Available under 3-clause BSD open-source license
- Documentation & Sources:
 - <u>https://www.score-p.org</u>
- User guide also part of installation:
 - orefix>/share/doc/scorep/{pdf,html}/
- Support and feedback: support@score-p.org
- Subscribe to news@score-p.org, to be up to date



Analysis report examination with Cube

Brian Wylie Jülich Supercomputing Centre





Cube

 CubeLib
 DOI
 10.5281/zenodo.7737408

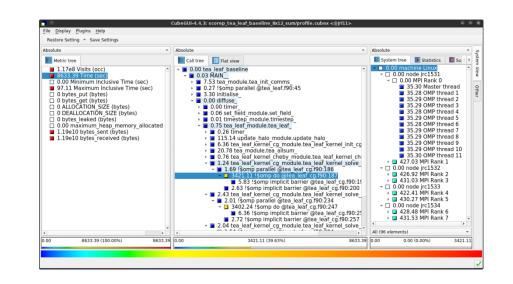
 CubeGUI
 DOI
 10.5281/zenodo.7737411

- Parallel program analysis report exploration tools
 - Libraries for XML+binary report reading & writing
 - Algebra utilities for report processing
 - GUI for interactive analysis exploration
 - Requires $Qt \ge 5$
- Originally developed as part of the Scalasca toolset



- Can be installed independently of Score-P, e.g., on laptop or desktop
- Latest release: Cube v4.8.2 (Sept 2023)

Note: source distribution tarballs for Linux, as well as binary packages provided for Windows & MacOS, from **www.scalasca.org** website in software/Cube-4x



VI-HPS

VIRTUAL INSTITUTE -- HIGH PRODUCTIVITY SUPERCOMPUTING

Cube GUI

mailto: scalasca@fz-juelich.de



Run remote (Jupyter-JSC)

- start Jupyter-JSC and then start Xpra desktop
- load cube module and start cube

[jwlogin~]\$ module load CubeGUI
[jwlogin~]\$ cube ./scorep-20221114_*/profile.cubex

Run remote (ssh)

- start X server (e.g., Xming) locally
- connect to juwels with X forwarding enabled

desk\$ ssh -X <yourid>@juwels-booster.fz-juelich.de Welcome to ...

```
[jwlogin~]$ module load CubeGUI
```

```
[jwlogin~]$ cube ./scorep-20221114_*/profile.cubex
```

Install & run *local*

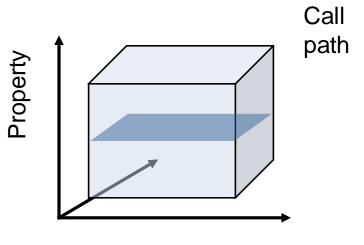
- install Cube GUI locally on desktop
 - binary packages available for MacOS & Windows and externally provided by OpenHPC and various Linux distributions
 - source package available for Linux, requires Qt
 - configure/build/install manually or use your favourite framework (e.g. Spack or EasyBuild)
- copy .cubex file (or entire scorep directory) to desktop from remote system
 OR locally mount remote filesystem
- start cube locally

```
desk$ mkdir $HOME/mnt
desk$ sshfs [user@]remote.sys:[dir] $HOME/mnt
desk$ cd $HOME/mnt
desk$ cube ./scorep-20221114_*/profile.cubex
```

VICTOR VI

Analysis presentation and exploration

- Representation of values (severity matrix) on three hierarchical axes
 - Performance property (metric)
 - Call path (program location)
 - System location (process/thread)
- Three coupled tree browsers
- Cube displays severities
 - As value: for precise comparison
 - As *colour*: for easy identification of hotspots
 - Inclusive value when closed & exclusive value when expanded
 - Customizable via display modes

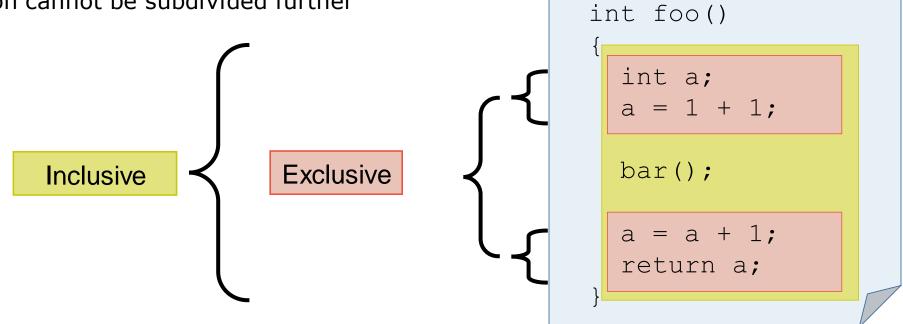


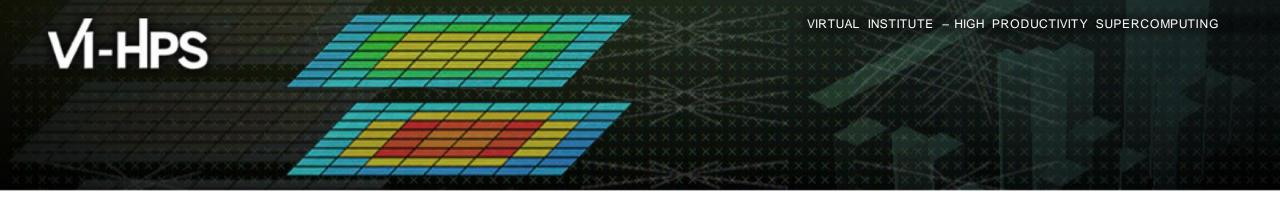


Inclusive vs. exclusive values



- Inclusive
 - Information of all sub-elements aggregated into single value
- Exclusive
 - Information cannot be subdivided further





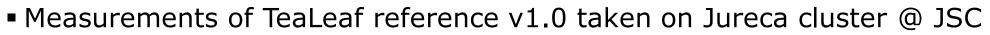
Demo: TeaLeaf case study





Case study: TeaLeaf

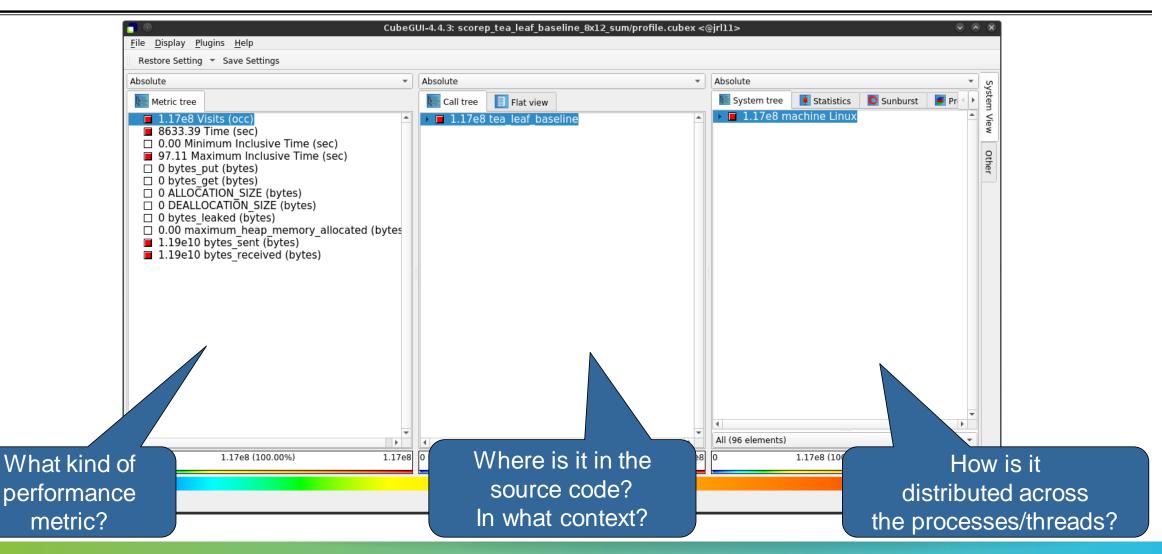
- HPC mini-app developed by the UK Mini-App Consortium
 - Solves the linear 2D heat conduction equation on a spatially decomposed regular grid using a 5 point stencil with implicit solvers
 - Part of the Mantevo 3.0 suite
 - Available on GitHub: http://uk-mac.github.io/TeaLeaf/



- Using Intel 19.0.3 compilers, Intel MPI 2019.3, and Score-P 5.0
- Run configuration
 - 8 MPI ranks with 12 OpenMP threads each

```
% cd ~/workshop-vihps/Experiments
% cube scorep_tea_leaf_baseline_8x12_sum/profile.cubex
[GUI showing summary analysis report]
```

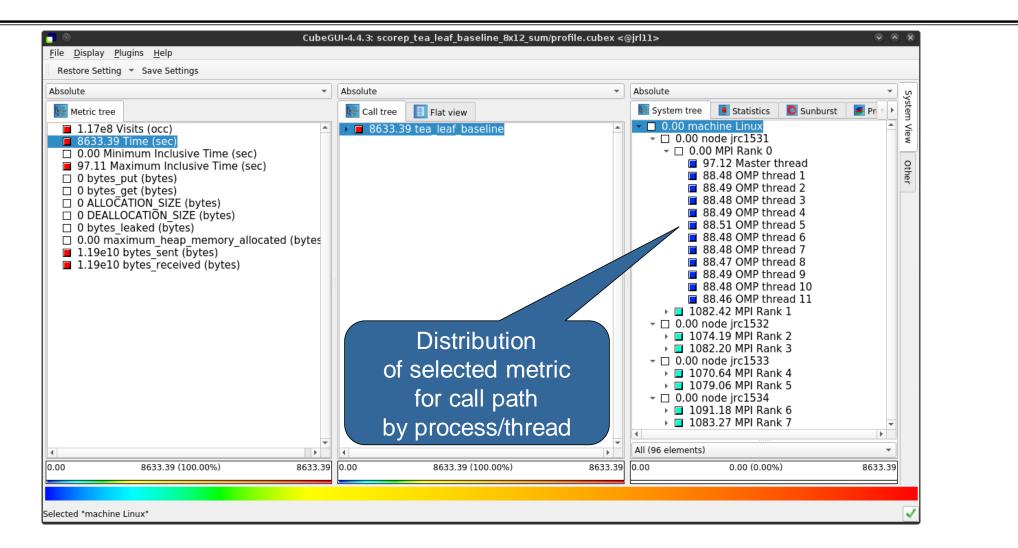
Score-P analysis report exploration (opening view)



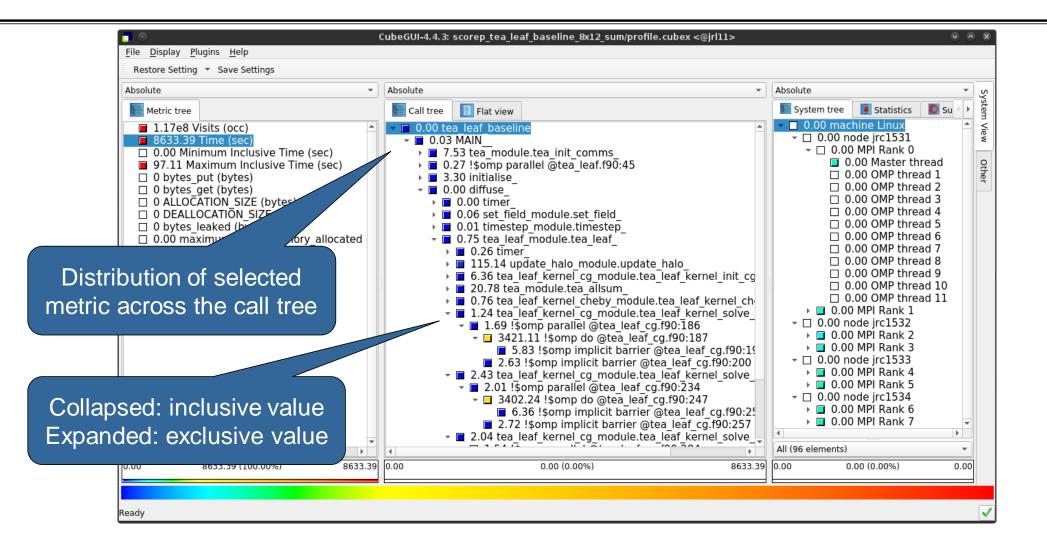
Metric selection



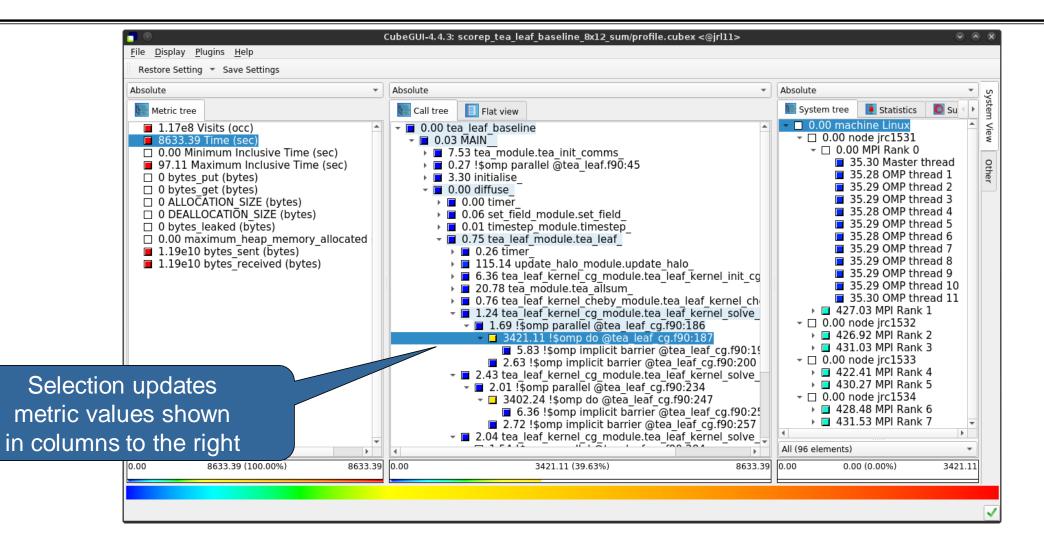
Expanding the system tree



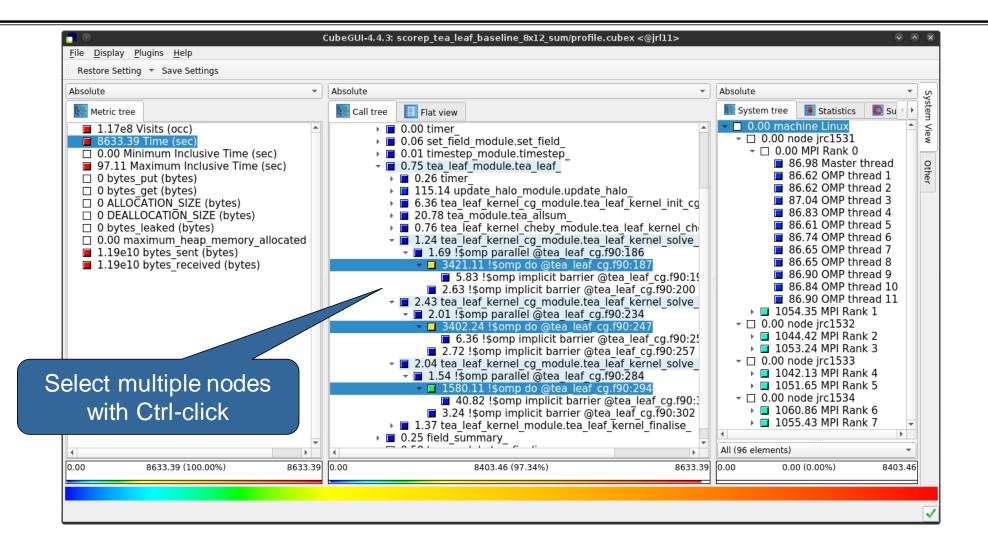
Expanding the call tree



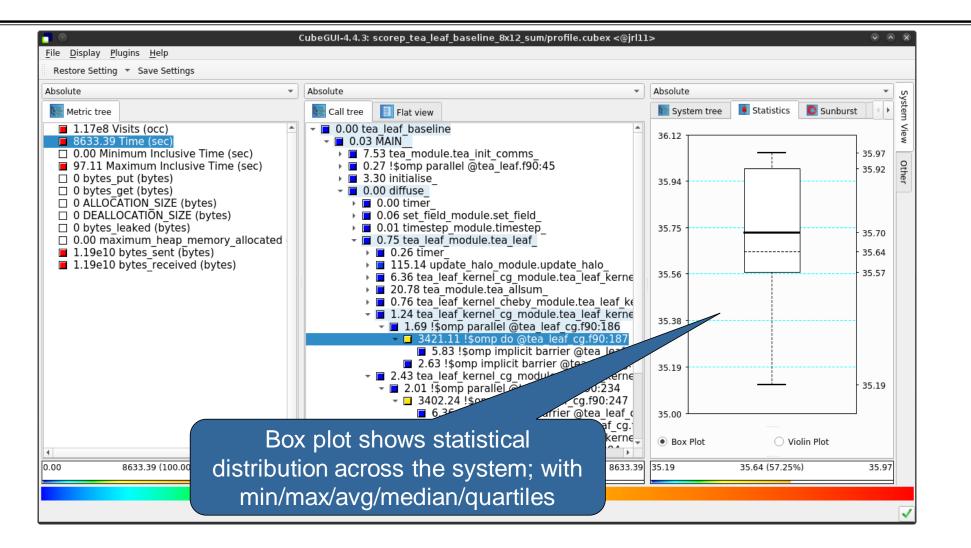
Selecting a call path



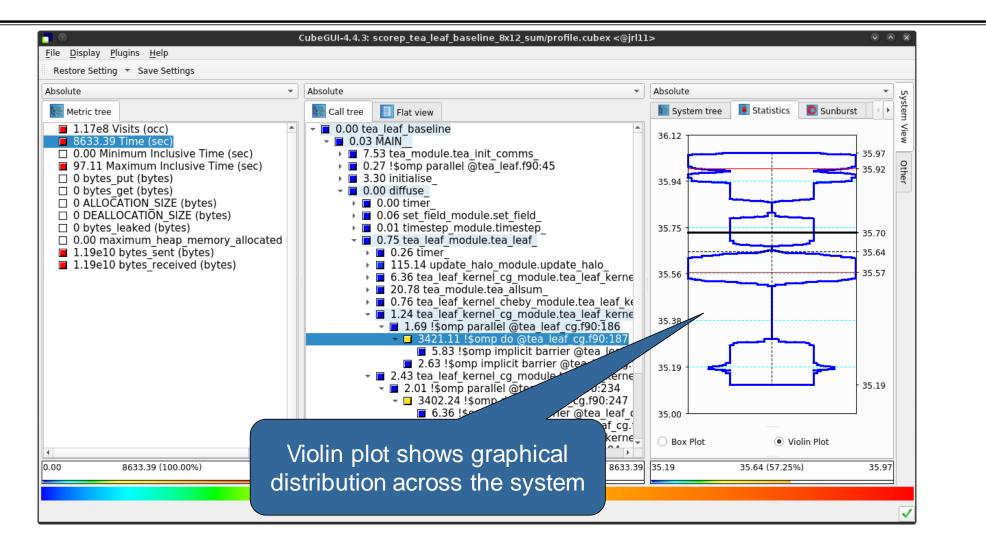
Multiple selection



Box plot view



Violin plot view



Topology view

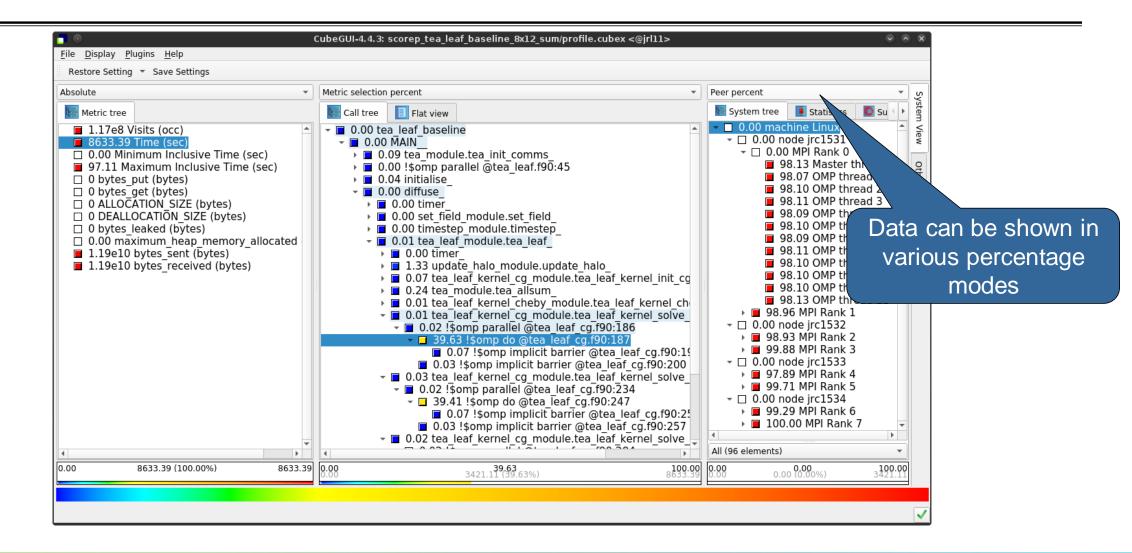
Absolute 🗸	Absolute	*	Peer percent			~ v
Metric tree	Call tree Flat view		Statistics	🖸 Sunburst	🕖 Process x Threa	System
 1.17e8 Visits (occ) 8633.39 Time (sec) 0.00 Minimum Inclusive Time (sec) 97.11 Maximum Inclusive Time (sec) 0 bytes_put (bytes) 0 bytes get (bytes) 0 ALLOCATION SIZE (bytes) 0 DEALLOCATION_SIZE (bytes) 0 bytes_leaked (bytes) 0.00 maximum_heap_memory_allocated 1.19e10 bytes_received (bytes) 1.19e10 bytes_received (bytes) 	 0.00 tea leaf baseline 0.03 MAIN 7.53 tea_module.tea init_comms 0.27 !\$omp parallel @tea_leaf.f90:45 3.30 initialise_ 0.00 diffuse_ 0.00 timer 0.06 set_field_module.set_field 0.01 timestep_module.timestep_ 0.75 tea_leaf_module.tea_leaf_ 0.26 timer 115.14 update_halo_module.update_halo 6.36 tea_leaf_kernel_cg_module.tea_lea 0.76 tea_leaf_kernel_cg_module.tea_lea 1.24 tea_leaf_kernel_cg_module.tea_lea 1.69 !\$omp parallel @tea_leaf_cg.f9 3421.11 !\$omp do @tea_leaf_cg.f9 2.63 !\$omp implicit barrier @tea 2.63 !\$omp implicit barrier @tea 3402.24 !\$om	a leaf ke a leaf ke af kerne 00:186 f90:187 ea 0:234 f90:247 ca leaf c				View Other
Sh	lows topological distribution	if_cg. kerne≖	4			
0.00 8633.39 (100.00%)	across the system	8633.39	0.00	0.00		100.00

VIRTUAL INSTITUTE -- HIGH PRODUCTIVITY SUPERCOMPUTING

Topology view (cont.)

Restore Setting	↓ ↓	Peer percent 🔹
 Metric tree 1.17e8 Visits (occ) 8633.39 Time (sec) 0.00 Minimum Inclusive Time (sec) 97.11 Maximum Inclusive Time (sec) 0 bytes_put (bytes) 0 bytes get (bytes) 0 ALLOCATION SIZE (bytes) 0 DEALLOCATION_SIZE (bytes) 0 bytes_leaked (bytes) 0.00 maximum heap_memory_allocated 1.19e10 bytes_received (bytes) 1.19e10 bytes_received (bytes) 	 Call tree Flat view 0.00 tea leaf baseline 0.03 MAIN 7.53 tea_module.tea_init_comms 0.27 !\$omp parallel @tea_leaf.f90:45 3.30 initialise 0.00 diffuse 0.00 timer 0.06 set_field_module.set_field 0.01 timestep_module.timestep 0.75 tea_leaf_module.tea_leaf 0.26 timer 0.76 tea_leaf_kernel_cg_module.tea_leaf_kerne 20.78 tea_module.tea_allsum 0.76 tea_leaf_kernel_cheby_module.tea_leaf 1.24 tea_leaf_kernel_cg_module.tea_leaf 3.421.11 !\$omp do @tea_leaf_c	Statistics Sunburst Process x Thread
	Selection & right-click	<u>ا</u>
.00 8633.39 (100.00%) 8633.39	shows details	0.00 85.46 100.00 0.00 0.06 (1.10%) 5.83

Alternative display modes



Important display modes

Absolute

Absolute value shown in seconds/bytes/counts

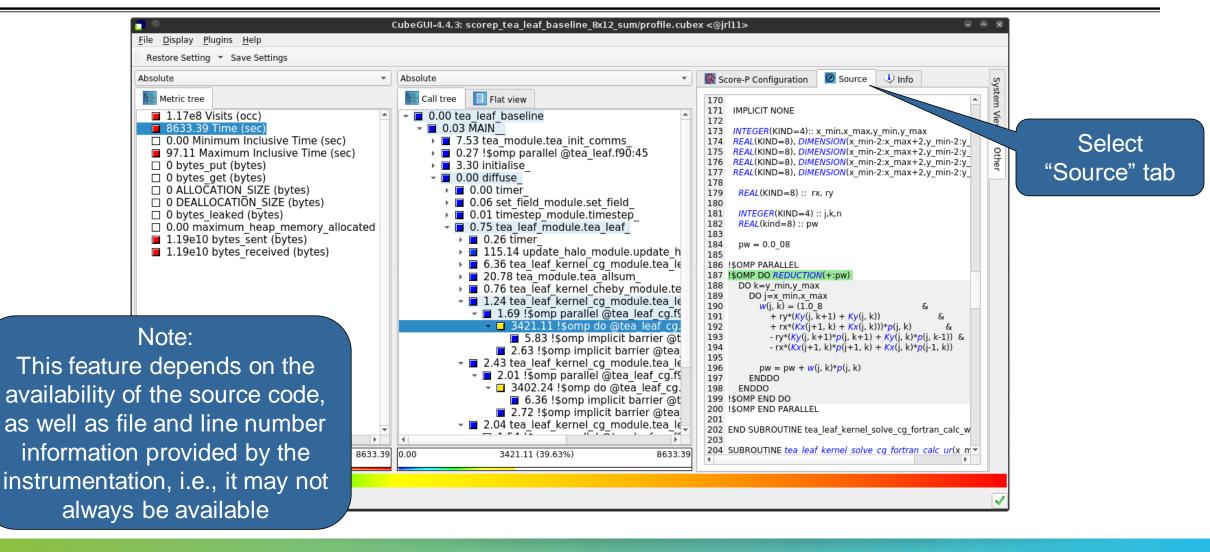
Selection percent

- Value shown as percentage w.r.t. the selected node "on the left" (metric/call path)
- Peer percent (system tree only)
 - Value shown as percentage relative to the maximum peer value

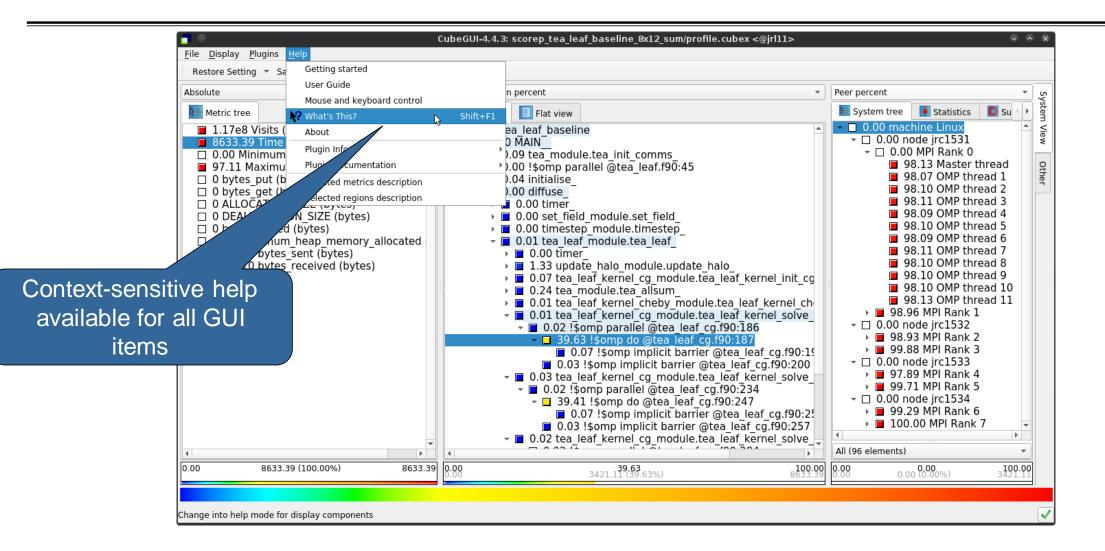
Source-code view via context menu

Absolute	Absolute	-	Absolute			
				Statistics 🛛 🖸 Su 📢		
Metric tree	Call tree Flat view		System tree		Su	
 1.17e8 Visits (occ) 8633.39 Time (sec) 0.00 Minimum Inclusive Time (sec) 97.11 Maximum Inclusive Time (sec) 0 bytes_put (bytes) 0 bytes_get (bytes) 0 ALLOCATION_SIZE (bytes) 0 DEALLOCATION_SIZE (bytes) 0 bytes_leaked (bytes) 0.00 maximum_heap_memory_allocated 1.19e10 bytes_received (bytes) 1.19e10 bytes_received (bytes) 	 0.00 tea_leaf_baseline 0.03 MAIN	rnel ch	 0.00 node 0.00 M 35.3 35.2 35.3 	e jrc1531 PI Rank 0 30 Master threa 29 OMP threac 29 OMP threac 29 OMP threac 29 OMP threac 29 OMP threac 29 OMP threac 29 OMP threac 30 OMP threac 30 OMP threac 30 OMP threac 40 OMP t	1 2 3 4 5 6 7 8 9 10	
t-click opens ntext menu	 2.01 !\$omp parallel @tea_leaf_cg.f90:234 3402.24 !\$omp do @tea_leaf_cg.f90:24 6.36 !\$omp implicit barrier @tea_lea 	Documentation Set as loop Expand/collap Hiding	pse)	MPI Rank 3 jrc1533 MPI Rank 4 MPI Rank 5 jrc1534 MPI Rank 6 MPI Rank 7		
		Cut call tree Find items	,		•	
4		- Clear found it	tems		*	
	2 22 2421 24 (22 22)			00%)	3421.11	
0.00 8633.39 (100.00%) 8633.39	0.00 3421.11 (39.63%)	Sort tree iten	ns)	00%)	5421.11	

Source-code view



Context-sensitive help



Scalasca report post-processing

- Scalasca's report post-processing derives additional metrics and generates a structured metric hierarchy
- Automatically run (if needed) when using the square convenience command:

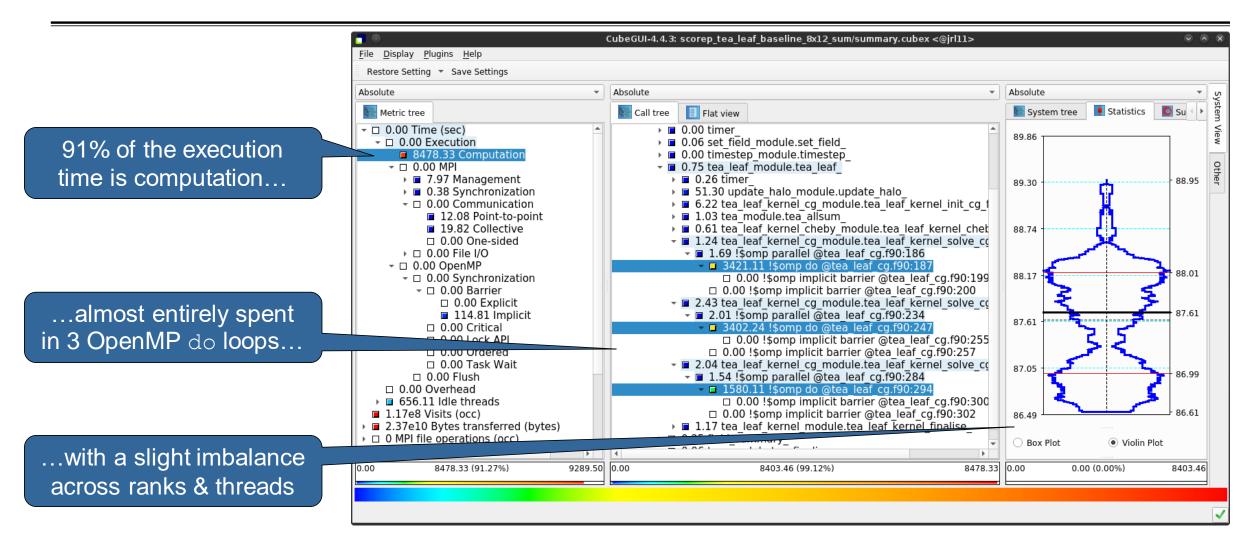
% square scorep_tea_leaf_baseline_8x12_sum
INFO: Post-processing runtime summarization report (profile.cubex)...
INFO: Displaying ./scorep_tea_leaf_baseline_8x12_sum/summary.cubex...

[GUI showing post-processed summary analysis report]

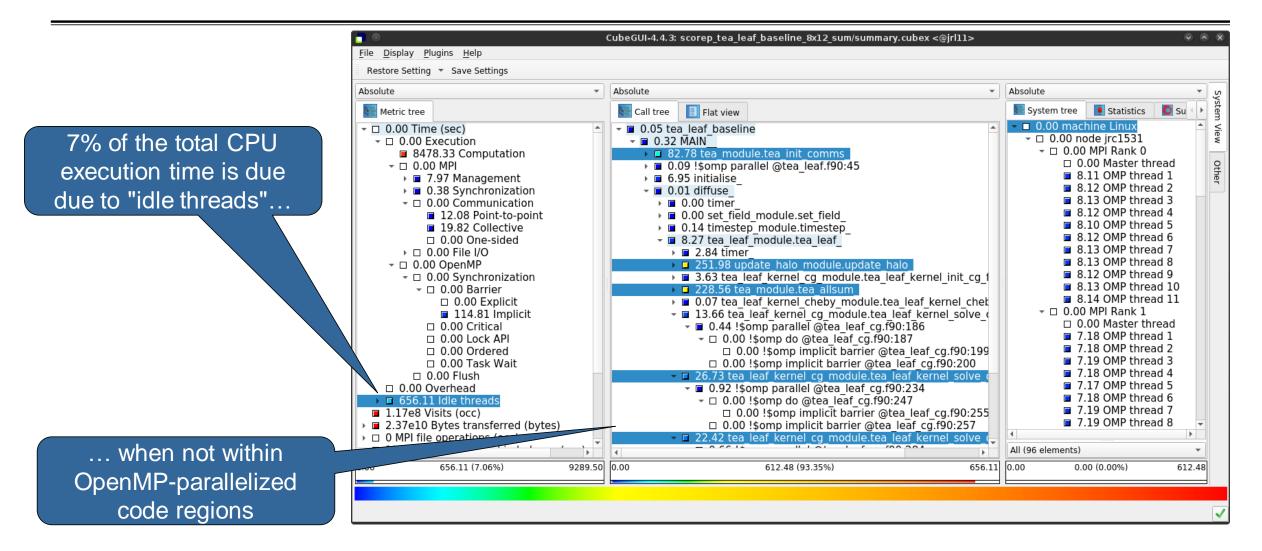
Post-processed summary analysis report

CubeGUI-4.4.3: scorep tea leaf baseline 8x12 sum/summary.cubex <@jrl11> File Display Plugins Help Restore Setting
 Save Settings Absolute Absolute Absolute ¥ I Ŧ Split base metrics into Syste Statistics 🖸 Sunburst System tree Flat view Metric tree 🚛 Call tree more specific metrics, 0.00 machine Linux □ 0.00 Time (sec) 0.00 tea leaf baseline View - 0.00 node jrc1531 0.03 MAIN 0.00 Execution e.g. computation vs - 0.00 MPI Rank 0 8478.33 Computation I 0.00 tea module.tea init comms ▶ ■ 0.00 !\$omp parallel @tea leaf.f90:45 35.30 Master thread Othe - 0.00 MPI 35.28 OMP thread 1 parallelization costs 7.97 Management 2.15 initialise 35.29 OMP thread 2 0.38 Synchronization 0.00 diffuse 35.29 OMP thread 3 - 0.00 Communication 0.00 timer 35.28 OMP thread 4 ▶ ■ 0.06 set field module.set field 12.08 Point-to-point 35.29 OMP thread 5 ▶ ■ 0.00 timestep module.timestep 19.82 Collective 35.28 OMP thread 6 0.00 One-sided 35.29 OMP thread 7 • 0.26 timer □ 0.00 File I/O
 □ 35.29 OMP thread 8 -
 0.00 OpenMP
 0.00
 0 35.29 OMP thread 9 ▶ ■ 6.22 tea leaf kernel cg module.tea □ 0.00 Synchronization I.03 tea module.tea allsum 35.29 OMP thread 10 - 0.00 Barrier ▶ ■ 0.61 tea leaf kernel cheby module. 35.30 OMP thread 11 □ 0.00 Explicit -
0.00 MPI Rank 1 I.24 tea leaf kernel cg module.tea 114.81 Implicit 🝷 🖬 1.69 !\$omp parallel @tea leaf cg. 35.59 Master thread 0.00 Critical 35.58 OMP thread 1 ✓ ■ 3421.11 !\$omp do @tea leaf c 0.00 Lock API □ 0.00 !\$omp implicit barrier @ 35.58 OMP thread 2 0.00 Ordered 35.58 OMP thread 3 0.00 Task Wait 0.00 !\$omp implicit barrier @te 35.58 OMP thread 4 2.43 tea leaf kernel cg module.tea 0.00 Flush 35.58 OMP thread 5 ▼ ■ 2.01 !\$omp parallel @tea leaf cg. 0.00 Overhead 35.59 OMP thread 6 - 3402.24 !\$omp do @tea leaf co G56.11 Idle threads 35.59 OMP thread 7 0.00 !somp implicit barrier a 1.17e8 Visits (occ) 35.58 OMP thread 8 2.37e10 Bytes transferred (bytes) 0.00 !\$omp implicit barrier @te Þ 2.04 tea leaf kernel cg module.tea O MPI file operations (occ) All (96 elements) Ŧ F 4 9289.50 0.00 3421.11 (40.35%) 8478.33 0.00 0.00 (0.00%) 00.0 8478.33 (91.27%) 3421.11

TeaLeaf summary report analysis (I)



TeaLeaf summary report analysis (II)



VIRTUAL INSTITUTE -- HIGH PRODUCTIVITY SUPERCOMPUTING

TeaLeaf summary report analysis (III)

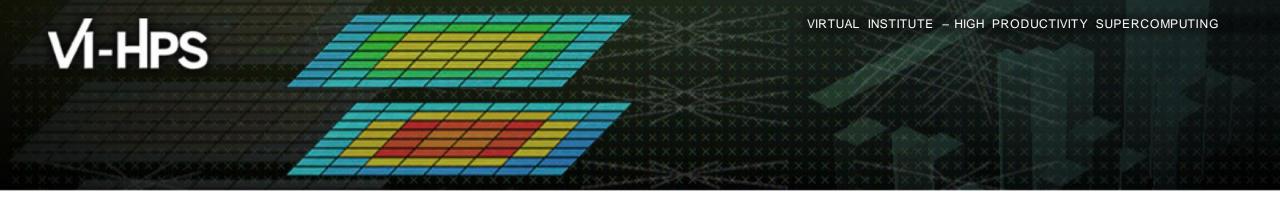
MPI communication time is negligible (0.34%); but communication is only on the master threads (MPI_THREAD_FUNNELED)

Absolute	Absolute	Absolute
Absolute +	Absolute	
Metric tree	Call tree 🔲 Flat view	🔚 System tree 🚺 Statistics 🚺 Su 🕐
- 🗆 0.00 Time (sec)		0.00 machine Linux
 0.00 Execution 	□ 0.00 MAIN □	
8478.33 Computation	• □ 0.00 tea module.tea init comms	
	▶ □ 0.00 !\$omp parallel @tea leaf.f90:45	4.88 Master thread
7.97 Management	•	0.00 OMP thread 1
0.38 Synchronization		0.00 OMP thread 2
 0.00 Communication 	→ □ 0.00 timer	0.00 OMP thread 3
12.08 Point-to-point	▷ □ 0.00 set field module.set field	0.00 OMP thread 4
19.82 Collective	•	0.00 OMP thread 5
0.00 One-sided	- 🗆 0.00 tea leaf module.tea leaf	0.00 OMP thread 6
▶ □ 0.00 File I/O	→ □ 0.00 timer	0.00 OMP thread 7
	12.03 update halo module.update halo	0.00 OMP thread 8
 0.00 Synchronization 	D 0.00 tea leaf kernel cg module.tea leaf kernel init cg 1	0.00 OMP thread 9
	19.74 tea module.tea allsum	0.00 OMP thread 10
0.00 Explicit	D 0.00 tea leaf kernel cheby module.tea leaf kernel cheb	0.00 OMP thread 11
114.81 Implicit	→ □ 0.00 tea leaf kernel cg module.tea leaf kernel solve cg	
□ 0.00 Critical	¬ □ 0.00 !\$omp parallel @tea leaf cg.f90:186 ¬	3.97 Master thread
0.00 Lock API		0.00 OMP thread 1
0.00 Ordered	□ 0.00 !\$omp implicit barrier @tea leaf cg.f90:199	0.00 OMP thread 2
0.00 Task Wait	□ 0.00 !\$omp implicit barrier @tea leaf cg.f90:200	0.00 OMP thread 3
0.00 Flush		□ 0.00 OMP thread 4
0.00 Overhead	¬ □ 0.00 !\$omp parallel @tea leaf cg.f90:234	0.00 OMP thread 5
656.11 Idle threads	0.00 !\$omp do @tea leaf cg.f90:247	0.00 OMP thread 6
1.17e8 Visits (occ)	□ 0.00 !\$omp implicit barrier @tea leaf cg.f90:255	0.00 OMP thread 7
	0.00 !\$omp implicit barrier @tea leaf cg.f90:257	0.00 OMP thread 8
□ 0 MPI file operations (occ)	□ 0.00 tea_leaf_kernel_cg_module.tea_leaf_kernel_solve_cc □ 0.00 tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_solve_cc □ 0.00 tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.tea_leaf_kernel_cg_module.te	
		All (96 elements)
.00 31.90 (0.34%) 9289.5	0.00 31.78 (99.62%) 31.90	0 0.00 0.00 (0.00%) 31.78

Cube: Further information

- Parallel program analysis report exploration tools
 - Libraries for Cube report reading & writing
 - Algebra utilities for report processing
 - GUI for interactive analysis exploration
- Available under 3-clause BSD open-source license
- Documentation & sources:
 - https://www.scalasca.org
- User guide also part of installation:
 - <prefix>/share/doc/cubegui/CubeUserGuide.pdf
- Contact:
 - mailto: scalasca@fz-juelich.de







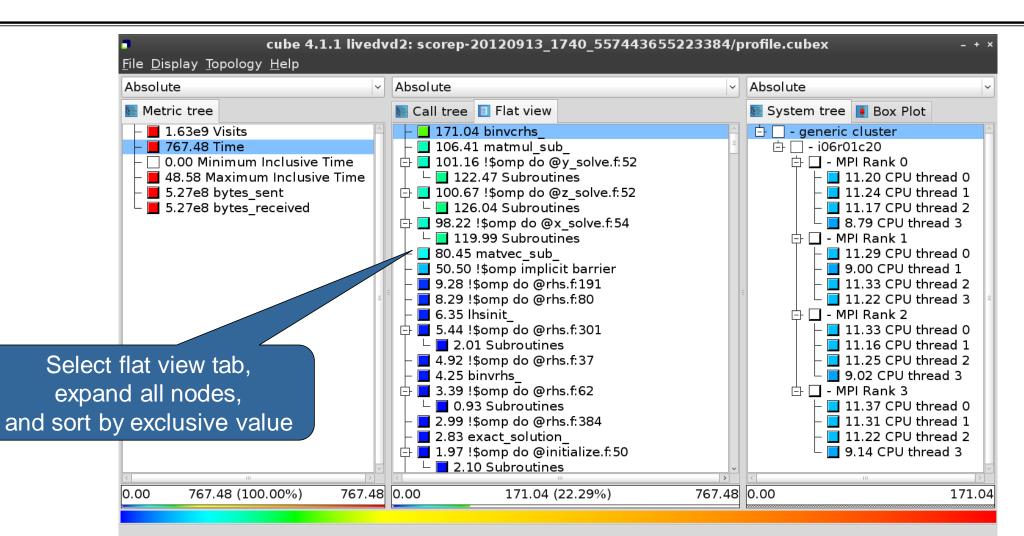




VI-HPS

Flat profile view





Derived metrics



Derived metrics are defined using CubePL expressions, e.g.:

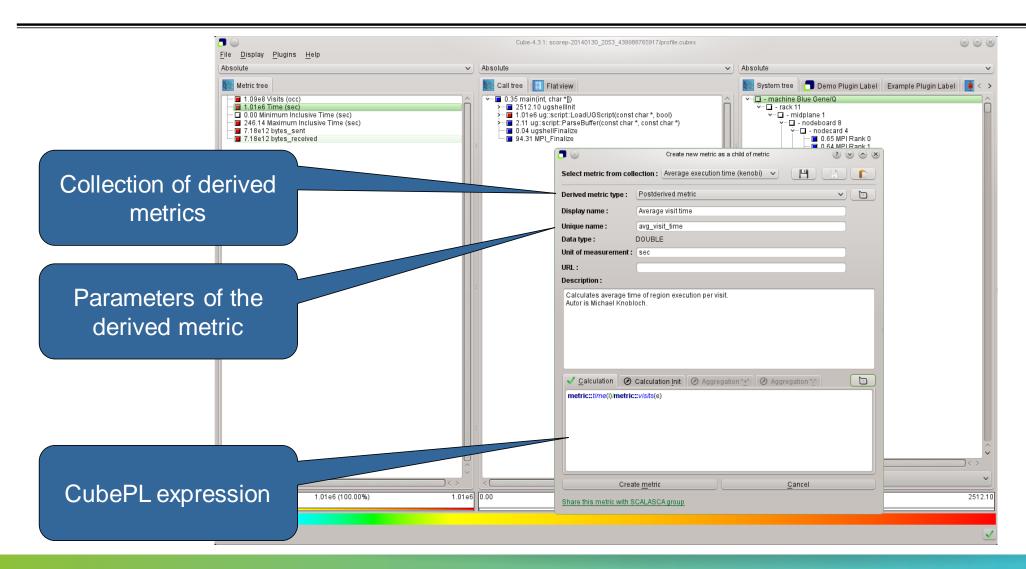
metric::time(i)/metric::visits(e)

- Values of derived metrics are not stored, but calculated on-the-fly
- Types of derived metrics:
 - Prederived: evaluation of the CubePL expression is performed before aggregation
 - Postderived: evaluation of the CubePL expression is performed after aggregation
- Examples:
 - "Average execution time": Postderived metric with expression

metric::time(i)/metric::visits(e)

Derived metrics in Cube GUI





Select metric from collection : --- please sele

✓ <u>C</u>alculation Ø Calculation Init Ø Aggreg

metric::PAPI FP OPS()/metric::time()

Edit <u>m</u>etric

Share this metric with SCALASCA group

FLOPS

flops

Derived metric type : Postderived metric

Display name :

Unique name : Data type :

URL : Description :

Unit of measurement :

Example: FLOPS based on PAPI_FP_OPS and time

	Ci	ıbe=4.3.1: scorep_8x4_sum/profile.cubex (on froggy1)	_ D X			
	<u>F</u> ile <u>D</u> isplay <u>P</u> lugins <u>H</u> elp					
	☐ Restore Setting ▼ Save Settings					
froggy1)	Absolute	Absolute	Absolute			
	keric tree	🔚 Call tree 📘 Flat view	🔚 System tree 🛛 Barplot 🔰 Heatmap 🚺 🚺 Boy 4 🕨			
	□ 1.17e7 Visits (occ)	■ 3.17e5 MAIN	🖻 = - machine Linux 📃			
	□ 1148.49 Time (sec)		e □ - node frog6			
	□ 0.00 Minimum Inclusive Time (sec)	■ 6.34e4 MPI Bcast	i⇒□ - MPI Rank 0			
	■ 41.57 Maximum Inclusive Time (₽ ■ 2.05e5 env setup	□ 1.17e9 Master thread			
	□ 0 bytes_put (bytes)	■ 7.39e5 zone_setup_	9.43e8 OMP thread 1			
	□ 0 bytes_get (bytes)	■ ■ 9.31e5 map_zones_	■ 9.47e8 OMP thread 2			
	■ 5.75e12 PAPI TOT INS (#)	9.39e4 zone_starts_	■ 9.47e8 OMP thread 3			
	2.69e12 PAPI_TOT_CYC (#)	■ 6.16e5 set constants	🖶 🗆 - MPI Rank 1			
	■ 2.12e12 PAPI FP_OPS (#)	🗐 🖶 🖬 5.91e8 initialize	🛛 🖬 1.17e9 Master thread			
	■ 3.12e9 bytes sent (bytes)	□ 0.00 exact rhs	■ 9.87e8 OMP thread 1			
	3.12e9 bytes_received (bytes)	□ □ 145.62 !\$omp parallel @exac	- ■ 9.68e8 OMP thread 2			
	■ 1.84e9 FLOPS	■ ■ 2.54e4 !\$omp do @exact_r	9.72e8 OMP thread 3			
		□ 9.65e8 !\$omp do @exact_r	🖶 🗆 - MPI Rank 2			
		∎ 9.62e8 !\$omp do @exact_r	□ 1.10e9 Master thread			
		■ ■ 8.14e8 !\$omp do @exact_r	- ■ 8.97e8 OMP thread 1			
" <u>+</u> " 🖉 Aggregation " <u>-</u> " 🛅		□ 1.21e5 !\$omp do @exact_r	■ 8.77e8 OMP thread 2			
		□ 0.00 !\$omp implicit barrier	■ 8.76e8 OMP thread 3			
			🖻 🗆 - MPI Rank 3			
		🗉 🖻 🖬 1.94e9 adi_	🗖 🗖 1.09e9 Master thread			
		□ 2.19e5 MPI_Barrier	■ 9.06e8 OMP thread 1			
		■ ■ 1.92e9 < <bt_iter>> (200 itera</bt_iter>	■ 9.04e8 OMP thread 2			
		■ ■ 1.98e8 verify_	■ 9.02e8 OMP thread 3			
Cancel		□ 1.05e5 MPI_Reduce				

9.65e8 (-0.00%)

Þ

1.84e9 0.00

Selected "!\$omp do @exact_rhs.f:46"

1.84e9 (100.00%)

0.00

•

All (32 elements)

0.00... -179769313486231570814527423731704356798070.

-12858016489314434.00

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

CUBE algebra utilities



Extracting solver sub-tree from analysis report

% cube_cut -r '<<ITERATION>>' scorep_bt-mz_C_32x4_sum/profile.cubex Writing cut.cubex... done.

Calculating difference of two reports

% cube_diff scorep_bt-mz_C_32x4_sum/profile.cubex cut.cubex
Writing diff.cubex... done.

- Additional utilities for merging, calculating mean, etc.
- Default output of cube_utility is a new report utility.cubex
- Further utilities for report scoring & statistics
- Run utility with `-h' (or no arguments) for brief usage info

Iteration profiling



- Show time dependent behavior by "unrolling" iterations
- Preparations:
 - Mark loop body by using Score-P instrumentation API in your source code

```
SCOREP_USER_REGION_DEFINE( scorep_bt_loop )
SCOREP_USER_REGION_BEGIN( scorep_bt_loop, "<<bt_iter>>", SCOREP_USER_REGION_TYPE_DYNAMIC )
SCOREP_USER_REGION_END( scorep_bt_loop )
```

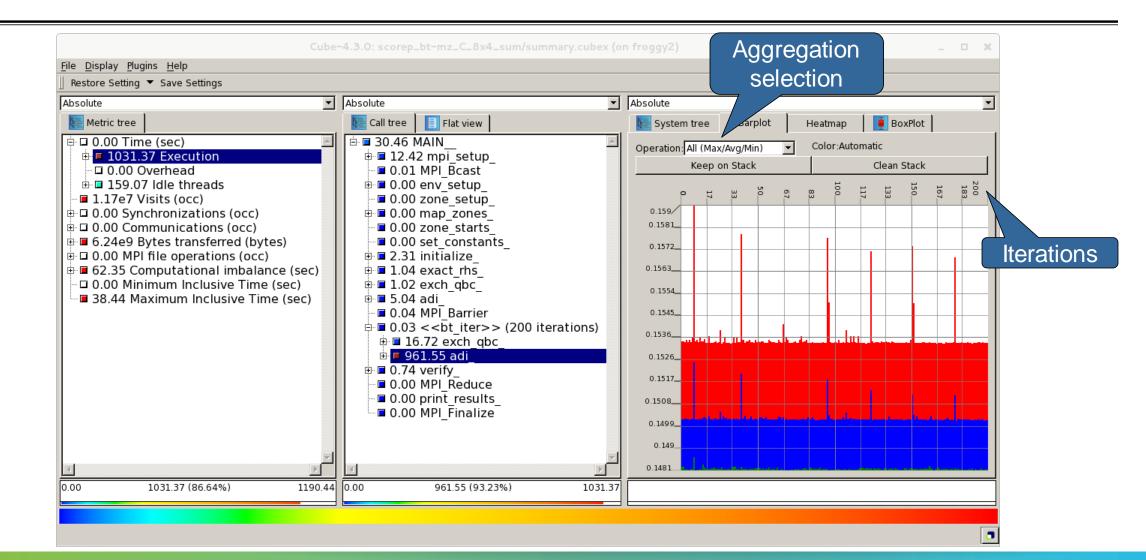
- Result in the Cube profile:
 - Iterations shown as separate call trees
 - Useful for checking results for specific iterations

or

- Select your user-instrumented region and mark it as loop
- Choose "Hide iterations"
- \succ View the Barplot statistics or the (thread x iterations) Heatmap

Iteration profiling: Barplot

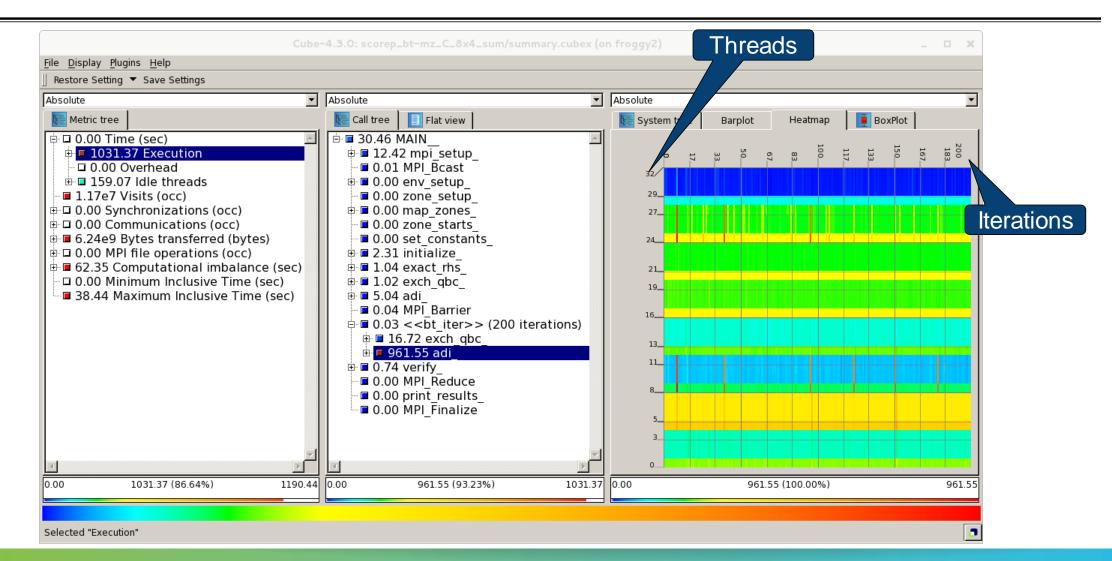


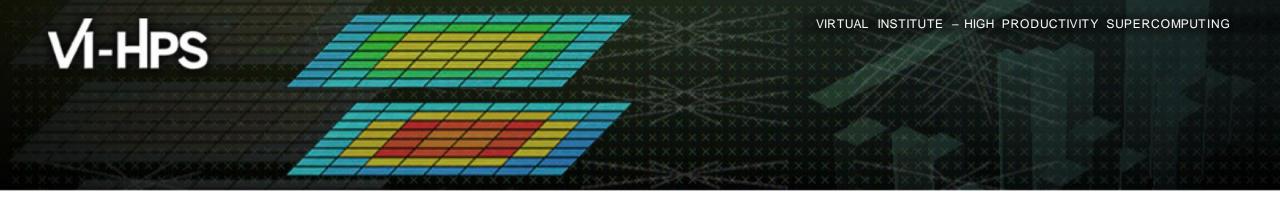


VIRTUAL INSTITUTE -- HIGH PRODUCTIVITY SUPERCOMPUTING

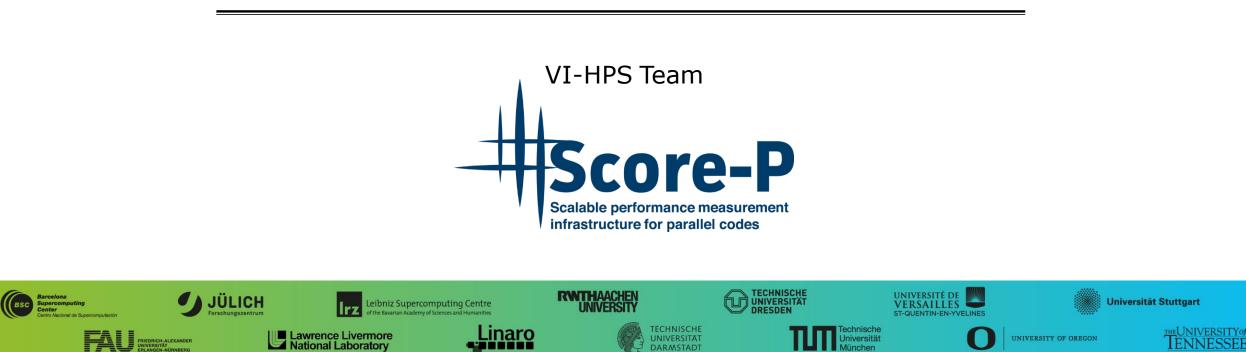
Iteration profiling: Heatmap







Score-P – A Joint Performance Measurement Run-Time Infrastructure for Scalasca, TAU, and Vampir



Congratulations!?

- If you made it this far, you successfully used Score-P to
 - instrument the application
 - analyze its execution with a summary measurement, and
 - examine it with one of the interactive analysis report explorer GUIs
- In revealing the call-path profile annotated with
 - the "Time" metric
 - Visit counts
 - MPI message statistics (bytes sent/received)
- ... but how good was the measurement?
 - The measured execution produced the desired valid result
 - but there wasn't much GPU-related performance information
 - and the execution took rather longer than expected!
 - even when ignoring measurement start-up/completion, therefore
 - it was probably dilated by instrumentation/measurement overhead

Performance analysis steps

- 0.0 Reference preparation for validation
- 1.0 Program instrumentation
- 1.1 Summary measurement collection
- 1.2 Summary analysis report examination
- 2.0 Summary experiment customisation & scoring
- 2.1 Summary measurement collection with filtering
- 2.2 Filtered summary analysis report examination

3.0 Event trace collection

3.1 Event trace examination & analysis

Mastering heterogeneous applications

Record CUDA application events and device activities

% export SCOREP_CUDA_ENABLE=default

- Record OpenCL application events and device activities
 - % export SCOREP_OPENCL_ENABLE=api,kernel
- Record OpenACC application events
 - % export SCOREP_OPENACC_ENABLE=yes
 - Can be combined with CUDA if it is a NVIDIA device
 - % export SCOREP_CUDA_ENABLE=kernel

Up to now: using default values

For all available options check: scorep-info config-vars --full

TeaLeaf CUDA extended summary measurement

```
% edit scorep.sbatch
% cat scorep.sbatch
```

```
# Score-P measurement configuration
export SCOREP_EXPERIMENT_DIRECTORY=scorep-tea_leaf-8.extended
export SCOREP_CUDA_ENABLE=default,driver,sync
export SCOREP_CUDA_BUFFER=48MB
#export SCOREP_FILTERING_FILE=../config/scorep.filter
```

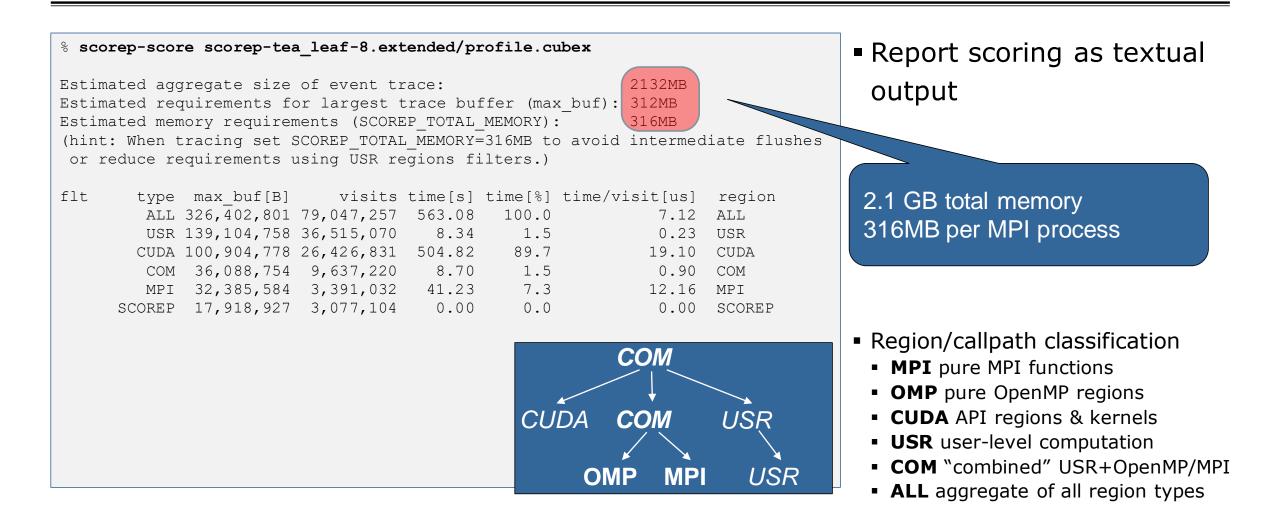
```
# Run the application
srun ./tea leaf
```

 $\frac{1}{2}$ sbatch scorep.sbatch

 Set new experiment directory and re-run measurement with extended CUDA event configuration

```
Submit job
```

TeaLeaf summary analysis result scoring



Score-P filtering: Automatic Generation of Filter Files

Basic usage: scorep-score -g

default heuristic targets:

- Buffer usage: relevancy
- Time per visits: overhead
- Creates annotated filter file:
 - initial_scorep.filter
 - Repeated calls create backups
 - Usage with -f <file> results in inclusion

• Objective:

- Starting point for filtering
- Syntax introduction

-g [<list>]

Generation of an initial filter file with the name 'initial_scorep.filter'. A valid parameter list has the form KEY=VALUE[,KEY=VALUE]*. By **default**, uses the following control parameters:

`bufferpercent=1,timepervisit=1`

A region is included in the filter file (i.e., excluded from measurement) if it matches all of the given conditions, with the following keys:

- `bufferpercent`

- `bufferabsolute`

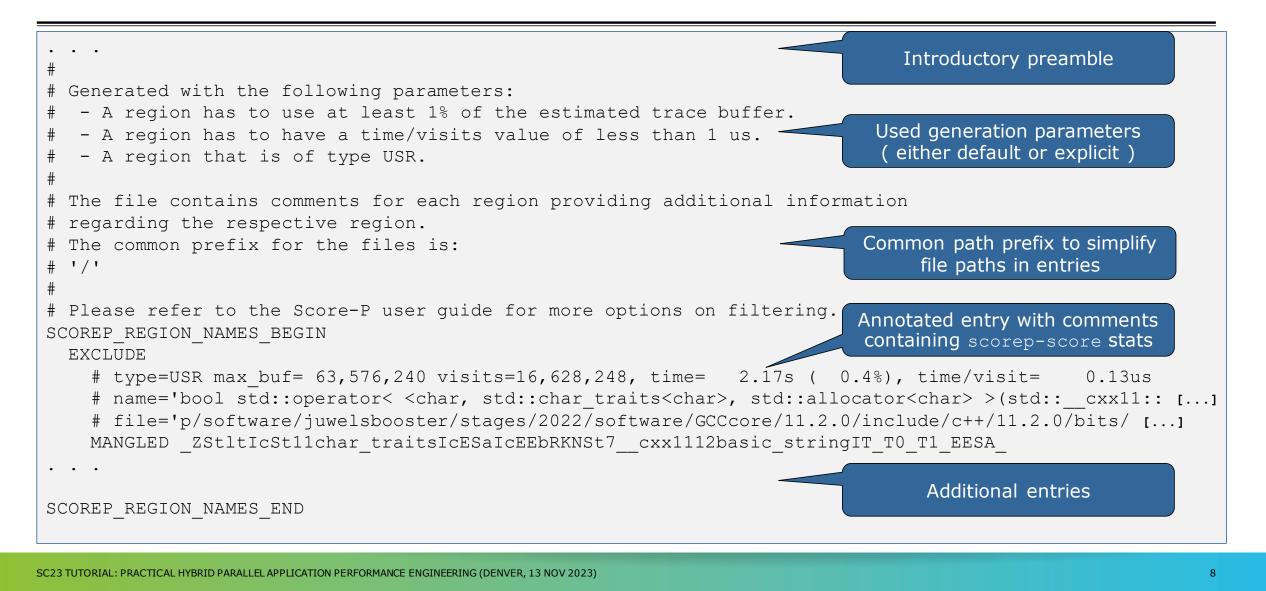
- `timepervisit`

- `visits`

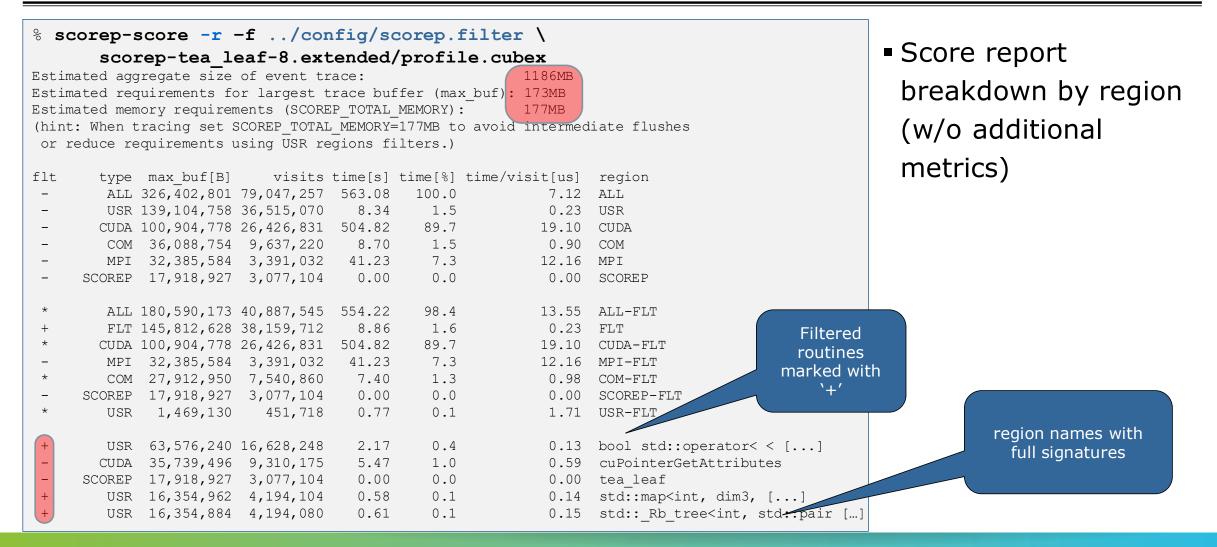
- `type`

- : estimated memory requirements exceed the given threshold in percent of the total estimated trace buffer requirements
- : estimated memory requirements exceed the given absolute threshold in MB
- : number of visits exceeds the given threshold
- : time per visit value is below the given threshold in microseconds
- : region type matches the given value
 (allowed: 'usr', 'com', 'both')

Score-P filtering: Automatic Generation of Filter Files



TeaLeaf summary analysis report filtering – Preview mode



TeaLeaf summary analysis report filtering – Preview incl. Counters

° SC	-	score <mark>-f</mark> cep-tea_lea	-	_					eport scoring with rospective filter
		regate size of			<i>,</i> , , , , , , , , , , , , , , , , , ,	3486MB			
		uirements for ory requirement				ouf): 506MB 510MB	~		
		racing set SCO							
		-				regions filters			
flt	type	max_buf[B]	visits	time[s]	time[%]	time/visit[us]	region		
-	ALL	1,006,819,003							3.5 GB of memory in total,
-	USR	454,765,555	36,515,070	8.34	1.5		USR		
-	CUDA	329,881,005	26,426,831	504.82	89.7	19.10	CUDA		506 MB per rank!
-	COM	117,982,465	9,637,220	8.70	1.5	0.90	COM		
-	MPI	60,485,278	3,391,032	41.23	7.3	12.16	MPI		
-	SCOREP	43,704,700	3,077,104	0.00	0.0	0.00	SCOREP		(Including 2 metric values)
*	ALL	530,123,873					ALL-FLT		
+	FLT	476,695,130	38,159,712	8.86	1.6	0.23	FLT		
*	CUDA	329,881,005	26,426,831	504.82	89.7	19.10	CUDA-FLT		
*	COM	91,253,875	7,540,860	7.40	1.3		COM-FLT		
-	MPI	60,485,278	3,391,032	41.23	7.3	12.16	MPI-FLT		
-	SCOREP	43,704,700	3,077,104	0.00	0.0	0.00	SCOREP-FLT		
*	USR	4,802,925	451 , 718	0.77	0.1	1.71	USR-FLT		

TeaLeaf filtered summary measurement

```
% edit scorep.sbatch
% cat scorep.sbatch
# Score-P measurement configuration
export SCOREP_EXPERIMENT_DIRECTORY=scorep_tea_leaf_sum.filtered
export SCOREP_CUDA_ENABLE=default,driver,sync
export SCOREP_CUDA_BUFFER=48MB
export SCOREP_FILTERING_FILE=../config/scorep.filter
# Run the application
srun ./tea_leaf
```

% sbatch scorep.sbatch

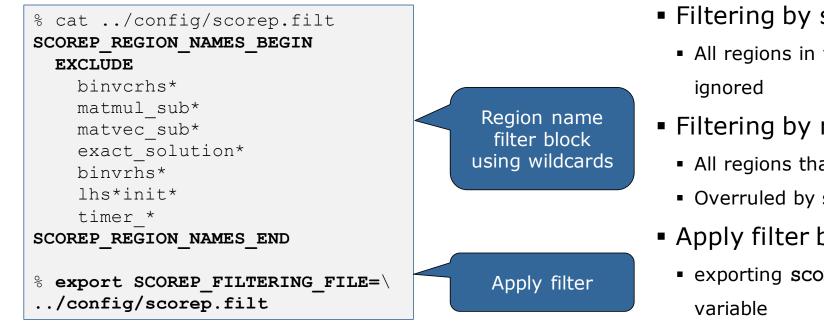
 Set new experiment directory and re-run measurement also with new filter configuration

```
Submit job
```

× × × × × × × × × × × VIRTUAL×INSTITUTE – HIGH_PRODUCTIVITY_SUPERCOMPUTING ×××××××××××××××××××××

Score-P filtering





- Filtering by source file name
 - All regions in files that are excluded by the filter are
- Filtering by region name
 - All regions that are excluded by the filter are ignored
 - Overruled by source file filter for excluded files
- Apply filter by
 - exporting scorep_filtering_file environment
- Apply filter at
 - Run-time
 - Compile-time (GCC-plugin, Intel compiler)
 - Add cmd-line option --instrument-filter
 - No overhead for filtered regions but recompilation

Source file name filter block



- Keywords
 - Case-sensitive
 - SCOREP_FILE_NAMES_BEGIN, SCOREP_FILE_NAMES_END
 - Define the source file name filter block
 - Block contains EXCLUDE, INCLUDE rules
 - EXCLUDE, INCLUDE rules
 - Followed by one or multiple white-space separated source file names
 - Names can contain bash-like wildcards *, ?, []
 - Unlike bash, * may match a string that contains slashes
- EXCLUDE, INCLUDE rules are applied in sequential order
- Regions in source files that are excluded after all rules are evaluated, get filtered

```
# This is a comment
SCOREP_FILE_NAMES_BEGIN
    # by default, everything is included
    EXCLUDE */foo/bar*
    INCLUDE */filter_test.c
SCOREP_FILE_NAMES_END
```

Region name filter block



- Keywords
 - Case-sensitive
 - SCOREP_REGION_NAMES_BEGIN,

SCOREP_REGION_NAMES_END

- Define the region name filter block
- Block contains EXCLUDE, INCLUDE rules
- EXCLUDE, INCLUDE rules
 - Followed by one or multiple white-space separated region names
 - Names can contain bash-like wildcards *, ?, []
- EXCLUDE, INCLUDE rules are applied in sequential order
- Regions that are excluded after all rules are evaluated, get filtered

```
# This is a comment
SCOREP_REGION_NAMES_BEGIN
    # by default, everything is included
    EXCLUDE *
    INCLUDE bar foo
        baz
        main
SCOREP_REGION_NAMES_END
```

VICTOR VICT

Region name filter block, mangling

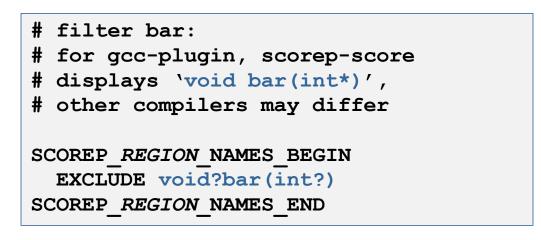


- Name mangling
 - Filtering based on names seen by the measurement system
 - Dependent on compiler
 - Actual name may be mangled
- scorep-score names as starting point

(e.g. matvec_sub_)

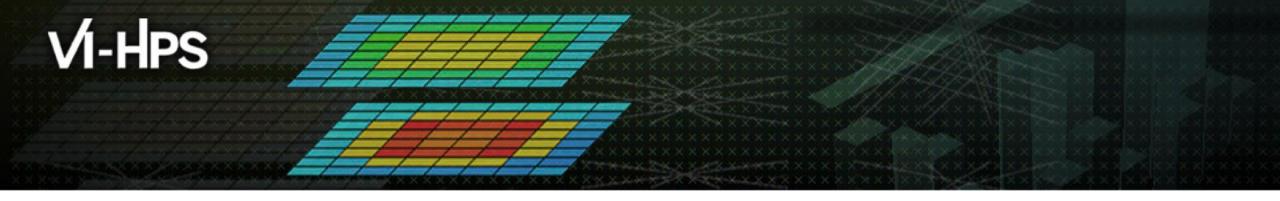
- Use * for Fortran trailing underscore(s) for portability
- Use ? and * as needed for full signatures or overloading

```
void bar(int* a) {
    *a++;
}
int main() {
    int i = 42;
    bar(&i);
    return 0;
}
```

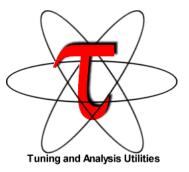


Further information

- Community instrumentation & measurement infrastructure
 - Instrumentation (various methods)
 - Basic and advanced profile generation
 - Event trace recording
 - Online access to profiling data
- Available under 3-clause BSD open-source license
- Documentation & Sources:
 - https://www.score-p.org
- User guide also part of installation:
 - orefix>/share/doc/scorep/{pdf,html}/
- Support and feedback: support@score-p.org
- Subscribe to news@score-p.org, to be up to date



Examination of profiles with TAU ParaProf



Sameer Shende sameer@cs.uoregon.edu

University of Oregon

http://tau.uoregon.edu/TAU_ParaProf_SC23.pdf



Application Performance Engineering using TAU

- How much time is spent in each application routine and outer *loops*? Within loops, what is the contribution of each *statement*? What is the time spent in OpenMP loops?
- How many instructions are executed in these code regions? Using Likwid or PAPI, TAU measures floating point, Level 1 and 2 *data cache misses*, hits, branches taken.
- What is the time taken in OS routines for thread scheduling? How much time is wasted?
- What is the memory usage of the code? When and where is memory allocated/de-allocated? Are there any memory leaks? What is the memory footprint of the application? What is the memory high water mark?
- 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?

TAU: Quickstart Guide

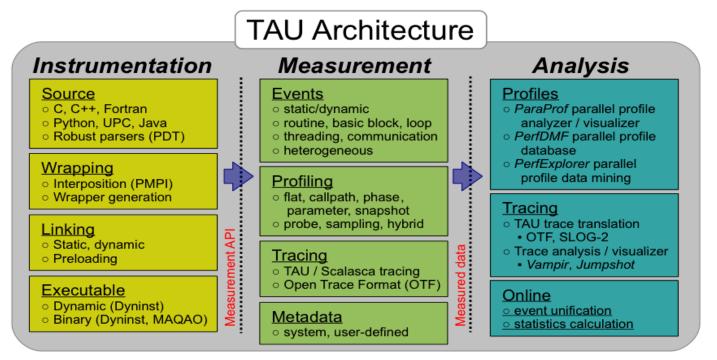
Profiling:

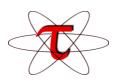
MPI: % mpirun -np 16 tau_exec -ebs ./a.out
 MPI+OpenMP: % export TAU_OMPT_SUPPORT_LEVEL=full;
 % mpirun -np 16 tau_exec -T ompt -ompt ./a.out
 Pthread: % mpirun -np 16 tau_exec -T mpi,pthread -ebs ./a.out
 CUDA: % mpirun -np 4 tau_exec -T cupti -cupti -ebs ./a.out
 Analysis: % pprof -a -m | more; % paraprof (GUI)
Tracing:
 Norming MDL: % expondent TAU_EDAGE=1; expondent TAU_EDAGE_FORMEret F0

- Vampir: MPI: % export TAU_TRACE=1; export TAU_TRACE_FORMAT=otf2 % mpirun -np 16 tau_exec ./a.out; vampir traces.otf2 &
- Chrome/Jumpshot: % export TAU_TRACE=1; mpirun -np 64 tau_exec ./a.out
 % tau_treemerge.pl;
- Chrome: % tau_trace2json tau.trc tau.edf -chrome -ignoreatomic -o app.json Chrome browser: chrome://tracing (Load -> app.json) or Perfetto.dev
- Jumpshot: tau2slog2 tau.trc tau.edf -o app.slog2; jumpshot app.slog2

TAU Performance System[®]

- Parallel performance framework and toolkit
 - Supports all HPC platforms, compilers, runtime system
 - Provides portable instrumentation, measurement, analysis





TAU Performance System

- Instrumentation
 - Fortran, C++, C, UPC, Java, Python, Chapel
 - Automatic instrumentation
- Measurement and analysis support
 - MPI, OpenSHMEM, ARMCI, PGAS, DMAPP
 - pthreads, OpenMP, OMPT interface, hybrid, other thread models
 - GPU, CUDA, OpenCL, OpenACC, ROCm, HIP
 - Parallel profiling and tracing
 - Use of Score-P for native OTF2 and CUBEX generation
 - Efficient callpath profiles and trace generation using Score-P
- Analysis
 - Parallel profile analysis (ParaProf), data mining (PerfExplorer)
 - Performance database technology (TAUdb)
 - 3D profile browser

TAU's Support for Runtime Systems

• MPI

- PMPI profiling interface
- MPI_T tools interface using performance and control variables

Pthread

Captures time spent in routines per thread of execution

OpenMP

- OMPT tools interface to track salient OpenMP runtime events
- Opari source rewriter
- Preloading wrapper OpenMP runtime library when OMPT is not supported

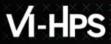
OpenACC

- OpenACC instrumentation API
- Track data transfers between host and device (per-variable)
- Track time spent in kernels

TAU's Support for Runtime Systems (contd.)

OpenCL

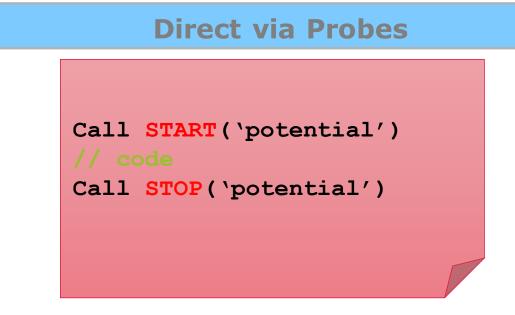
- OpenCL profiling interface
- Track timings of kernels
- Intel[®] OneAPI
 - Level Zero
 - Track time spent in kernels executing on GPU
 - Track time spent in OneAPI runtime calls
- CUDA
 - Cuda Profiling Tools Interface (CUPTI)
 - Track data transfers between host and GPU
 - Track access to uniform shared memory between host and GPU
- ROCm
 - Rocprofiler and Roctracer instrumentation interfaces
 - Track data transfers and kernel execution between host and GPU
- Kokkos
 - Kokkos profiling API
 - Push/pop interface for region, kernel execution interface
- Python
 - Python interpreter instrumentation API
 - Tracks Python routine transitions as well as Python to C transitions



Examples of Multi-Level Instrumentation

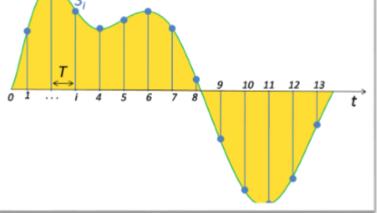
- MPI + OpenMP
 - MPI_T + PMPI + OMPT may be used to track MPI and OpenMP
- MPI + CUDA
 - PMPI + CUPTI interfaces
- Kokkos + OpenMP
 - Kokkos profiling API + OMPT to transparently track events
- Kokkos + pthread + MPI
 - Kokkos + pthread wrapper interposition library + PMPI layer
- Python + CUDA + MPI
 - Python + CUPTI + pthread profiling interfaces (e.g., Tensorflow, PyTorch) + MPI
- MPI + OpenCL
 - PMPI + OpenCL profiling interfaces

Performance Data Measurement



- Exact measurement
- Fine-grain control
- Calls inserted into code





- No code modification
- Minimal effort
- Relies on debug symbols (-g)

SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)

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

Using TAU's Runtime Preloading Tool: tau_exec

- Preload a wrapper that intercepts the runtime system call and substitutes with another
 - •MPI
 - ■OpenMP
 - POSIX I/O
 - Memory allocation/deallocation routines
 - Wrapper library for an external package
- No modification to the binary executable!
- Enable other TAU options (communication matrix, OTF2, event-based sampling)
- Add tau_exec before the name of the binary
 - ■mpirun –np 64 tau_exec ./a.out
 - mpirun tau_exec -T ompt,v5,mpi,papi -ompt ./a.out

VI-HPS

tau_exec

\$ tau_e	xec				
Usage:	tau_exec [opti	.ons] [] <exe> <exe options=""></exe></exe>	Tau_exec preloads		
Usage: Options	-v -s -qsub -io -memory_debug -cuda -cupti -opencl -openacc -ompt -armci -ebs -ebs_period=<	<pre>Verbose mode Show what will be done but don't actually do anything (dryrun) Use qsub mode (BG/P only, see below) Track I/O Track memory allocation/deallocation g Enable memory debugger Track GPU events via CUDA Track GPU events via CUDA Track GPU events via CUPTI (Also see env. variable TAU_CUPTI_API) Track GPU events via OpenCL Track GPU events via OpenACC (currently PGI only) Track OpenMP events via OMPT interface Track ARMCI events via PARMCI Enable event-based sampling Scount> Sampling period (default 1000) Scounter> Counter (default itimer) Enable Unified Memory events via CUPTI</pre>	Iau_exec preloads the TAU wrapper libraries and performs measurements.		
	-T <disable,g< th=""><th>SNU,ICPC,MPI,OMPT,OPENMP,PAPI,PDT,PROFILE,PTHREAD,SCOREP,SERIAL> : Specify TAU tag .e.so> : Specify additional load library</th><th>S</th></disable,g<>	SNU,ICPC,MPI,OMPT,OPENMP,PAPI,PDT,PROFILE,PTHREAD,SCOREP,SERIAL> : Specify TAU tag .e.so> : Specify additional load library	S		
		<pre>options> : Specify TAU library directly</pre>			
Notes:	-gdb	Run program in the gdb debugger No need to	recompile the application!		
	Defaults if unspecified: -T MPI				
	MPI is ass	sumed unless SERIAL is specified			

tau_exec Example (continued)

```
Example:
    mpirun -np 2 tau exec -T icpc,ompt,mpi -ompt ./a.out
    mpirun -np 2 tau exec -io ./a.out
Example - event-based sampling with samples taken every 1,000,000 FP instructions
    mpirun -np 8 tau exec -ebs period=1000000 -ebs source=PAPI FP INS ./ring
Examples - GPU:
    tau exec -T serial, cupti -cupti ./matmult (Preferred for CUDA 4.1 or later)
   tau exec -openacc ./a.out
   tau exec -T serial -opencl ./a.out (OPENCL)
   mpirun -np 2 tau exec -T mpi, cupti, papi -cupti -um ./a.out (Unified Virtual Memory in CUDA 6.0+)
qsub mode (IBM BG/Q only):
    Original:
      gsub -n 1 --mode smp -t 10 ./a.out
    With TAU:
      tau exec -qsub -io -memory -- qsub -n 1 ... -t 10 ./a.out
Memory Debugging:
    -memory option:
      Tracks heap allocation/deallocation and memory leaks.
    -memory debug option:
      Detects memory leaks, checks for invalid alignment, and checks for
      array overflow. This is exactly like setting TAU TRACK MEMORY LEAKS=1
      and TAU MEMDBG PROTECT ABOVE=1 and running with -memory
```

 tau_exec can enable event based sampling while launching the executable using the -ebs flag!

Simplifying TAU's usage (tau_exec)

- Uninstrumented execution linked with –dynamic (dynamic executables only!)
- % mpirun -np 16 ./a.out
- Track MPI performance
 - % mpirun -np 16 tau_exec ./a.out
- Track OpenMP, and MPI performance (MPI enabled by default; OMPT in Clang 9+, Intel 19+) % export TAU_OMPT_SUPPORT_LEVEL=full;
 - % mpirun -np 16 tau_exec -T mpi,pdt,ompt,papi -ompt ./a.out
- Track memory operations
 - % export TAU_TRACK_MEMORY_LEAKS=1
 - % mpirun -np 16 tau_exec -memory_debug ./a.out (bounds check)
- Use event based sampling (compile with -g)
 - % mpirun -np 16 tau_exec -ebs ./a.out
 - Also -ebs_source=<PAPI_COUNTER> -ebs_period=<overflow_count> -ebs_resolution=<file|function|line>
- Load wrapper interposition library
- % mpirun -np 16 tau_exec -loadlib=<path/libwrapper.so> ./a.out
- Track GPGPU operations (-rocm, -I0, -opencl, -cupti, -cupti –um, -openacc):
 - % mpirun -np 16 tau_exec -cupti ./a.out

Configuration tags for tau_exec

```
% ./configure -pdt=<dir> -mpi -papi=<dir>; make install
Creates in $TAU:
Makefile.tau-papi-mpi-pdt(Configuration parameters in stub makefile)
shared-papi-mpi-pdt/libTAU.so
% ./configure -pdt=<dir> -mpi; make install creates
Makefile.tau-mpi-pdt
shared-mpi-pdt/libTAU.so
To explicitly choose preloading of shared-<options>/libTAU.so change:
% mpirun -np 256 ./a.out to
% mpirun -np 256 tau exec -T <comma separated options> ./a.out
% mpirun -np 256 tau exec -T papi,mpi,pdt ./a.out
Preloads $TAU/shared-papi-mpi-pdt/libTAU.so
% mpirun -np 256 tau exec -T papi ./a.out
Preloads $TAU/shared-papi-mpi-pdt/libTAU.so by matching.
% aprun -n 256 tau exec -T papi,mpi,pdt -s ./a.out
Does not execute the program. Just displays the library that it will preload if executed without the -s option.
NOTE: -mpi configuration is selected by default. Use -T serial for
Sequential programs.
```

TAU Execution Command (tau_exec)

Uninstrumented execution

• % mpirun -np 256 ./a.out

Track GPU operations

- % mpirun –np 256 tau_exec –rocm ./a.out
- % mpirun –np 256 tau_exec –cupti ./a.out
- % mpirun –np 256 tau_exec –opencl ./a.out
- % mpirun –np 256 tau_exec –openacc ./a.out
- % mpirun –np 256 tau_exec –l0 ./a.out
- Track MPI performance
 - % mpirun -np 256 tau_exec ./a.out

Track I/O, and MPI performance (MPI enabled by default)

- % mpirun -np 256 tau_exec -io ./a.out
- Track OpenMP and MPI execution (using OMPT for Intel v19+ or Clang 8+)
 - % export TAU_OMPT_SUPPORT_LEVEL=full;
 - % mpirun –np 256 tau_exec –T ompt,mpi -ompt ./a.out

Track memory operations

- % export TAU_TRACK_MEMORY_LEAKS=1
- % mpirun –np 256 tau_exec –memory_debug ./a.out (bounds check)
- ■Use event based sampling (compile with -g)
 - % mpirun –np 256 tau_exec –ebs ./a.out
 - Also -ebs_source=<PAPI_COUNTER> -ebs_period=<overflow_count> -ebs_resolution=<file | function | line>

Hands-On Exercises for ParaProf

% source /p/project/training2341/setup.sh
% wget <u>http://tau.uoregon.edu/demo.ppk</u>
% paraprof demo.ppk &

% wget http://tau.uoregon.edu/data.tgz % tar zxf data.tgz; cd data/tau; % paraprof *.ppk &

TAU's Runtime Environment Variables

Environment Variable	Default	Description	
TAU_TRACE	0	Setting to 1 turns on tracing	
TAU_CALLPATH	0	Setting to 1 turns on callpath profiling	
TAU_TRACK_MEMORY_FOOTPRINT 0		Setting to 1 turns on tracking memory usage by sampling periodically the resident set size and high water mark of memory usage	
TAU_TRACK_POWER	0	Tracks power usage by sampling periodically.	
TAU_CALLPATH_DEPTH	2	Specifies depth of callpath. Setting to 0 generates no callpath or routine information, setting to 1 generates flat profile and context events have just parent information (e.g., Heap Entry: foo)	
TAU_SAMPLING	1	Setting to 1 enables event-based sampling.	
TAU_TRACK_SIGNALS	0	Setting to 1 generate debugging callstack info when a program crashes	
TAU_COMM_MATRIX	0	Setting to 1 generates communication matrix display using context events	
TAU_THROTTLE	1	Setting to 0 turns off throttling. Throttles instrumentation in lightweight routines that are called frequently	
TAU_THROTTLE_NUMCALLS	100000	Specifies the number of calls before testing for throttling	
TAU_THROTTLE_PERCALL	10	Specifies value in microseconds. Throttle a routine if it is called over 100000 times and takes less than 10 usec of inclusive time per call	
TAU_CALLSITE	0	Setting to 1 enables callsite profiling that shows where an instrumented function was called. Also compatible with tracing.	
TAU_PROFILE_FORMAT	Profile	Setting to "merged" generates a single file. "snapshot" generates xml format	

Runtime Environment Variables

Environment Variable	Default	Description
TAU_METRICS	TIME	Setting to a comma separated list generates other metrics. (e.g., ENERGY,TIME,P_VIRTUAL_TIME,PAPI_FP_INS,PAPI_NATIVE_ <event>:<subevent>)</subevent></event>
TAU_TRACE	0	Setting to 1 turns on tracing
TAU_TRACE_FORMAT	Default	Setting to "otf2" turns on TAU's native OTF2 trace generation (configure with –otf=download)
TAU_EBS_UNWIND	0	Setting to 1 turns on unwinding the callstack during sampling (use with tau_exec –ebs or TAU_SAMPLING=1)
TAU_EBS_RESOLUTION	line	Setting to "function" or "file" changes the sampling resolution to function or file level respectively.
TAU_TRACK_LOAD	0	Setting to 1 tracks system load on the node
TAU_SELECT_FILE	Default	Setting to a file name, enables selective instrumentation based on exclude/include lists specified in the file.
TAU_OMPT_SUPPORT_LEVEL	basic	Setting to "full" improves resolution of OMPT TR6 regions on threads 1 N-1. Also, "lowoverhead" option is available.
TAU_OMPT_RESOLVE_ADDRESS_EAGERLY	1	Setting to 1 is necessary for event based sampling to resolve addresses with OMPT. Setting to 0 allows the user to do offline address translation.
TAU_EVENT_THRESHOLD	0.5	Define a threshold value (e.g., .25 is 25%) to trigger marker events for min/max

Runtime Environment Variables

Environment Variable	Default	Description
TAU_TRACK_MEMORY_LEAKS	0	Tracks allocates that were not de-allocated (needs –optMemDbg or tau_exec –memory)
TAU_EBS_SOURCE	TIME	Allows using PAPI hardware counters for periodic interrupts for EBS (e.g., TAU_EBS_SOURCE=PAPI_TOT_INS when TAU_SAMPLING=1)
TAU_EBS_PERIOD	100000	Specifies the overflow count for interrupts
TAU_MEMDBG_ALLOC_MIN/MAX	0	Byte size minimum and maximum subject to bounds checking (used with TAU_MEMDBG_PROTECT_*)
TAU_MEMDBG_OVERHEAD	0	Specifies the number of bytes for TAU's memory overhead for memory debugging.
TAU_MEMDBG_PROTECT_BELOW/ABOVE	0	Setting to 1 enables tracking runtime bounds checking below or above the array bounds (requires – optMemDbg while building or tau_exec –memory)
TAU_MEMDBG_ZERO_MALLOC	0	Setting to 1 enables tracking zero byte allocations as invalid memory allocations.
TAU_MEMDBG_PROTECT_FREE	0	Setting to 1 detects invalid accesses to deallocated memory that should not be referenced until it is reallocated (requires –optMemDbg or tau_exec –memory)
TAU_MEMDBG_ATTEMPT_CONTINUE	0	Setting to 1 allows TAU to record and continue execution when a memory error occurs at runtime.
TAU_MEMDBG_FILL_GAP	Undefined	Initial value for gap bytes
TAU_MEMDBG_ALINGMENT	Sizeof(int)	Byte alignment for memory allocations

Installing and Configuring TAU

Installing PDT:

- wget http://tau.uoregon.edu/pdt.tgz
- ./configure; make ; make install

Installing TAU :

- wget http://tau.uoregon.edu/tau.tgz
- ./configure -ompt -c++=mpiicpc -cc=mpiicc -fortran=mpiifort -mpi -bfd=download -pdt=<dir> -papi=<dir>
- make install; export PATH=<taudir>/x86_64/bin:\$PATH
- All configurations are stored in <taudir>/.all_configs if you wish to see how TAU was configured!
- •Using TAU for source instrumentation:
 - export TAU_MAKEFILE=<taudir>/x86_64/lib/Makefile.tau-<TAGS>
 - make CC=tau_cc.sh CXX=tau_cxx.sh F90=tau_f90.sh
 - Use tau_exec with uninstrumented binaries instead of recompiling the source code.

Installing TAU on your laptop for paraprof (GUI)

Microsoft Windows

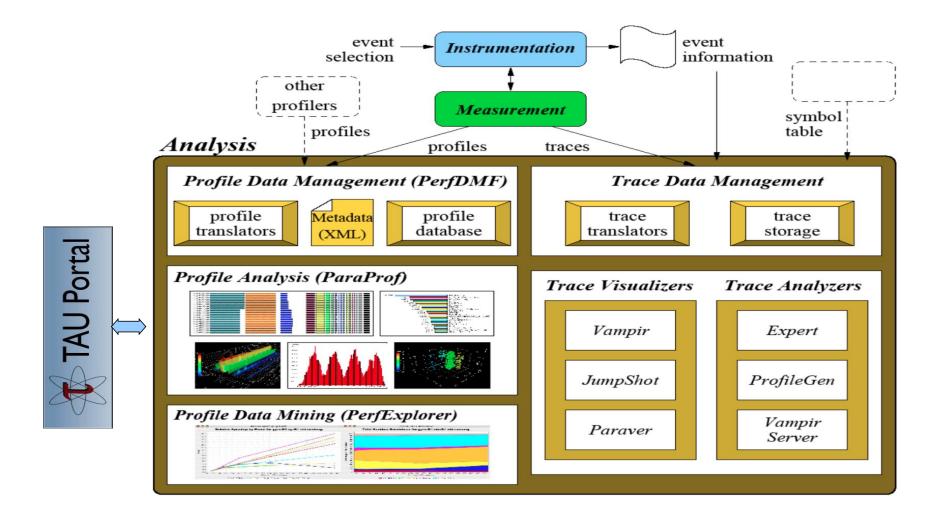
- Install Java from Oracle.com
- http://tau.uoregon.edu/tau.exe
- Install, click on a ppk file to launch paraprof

macOS

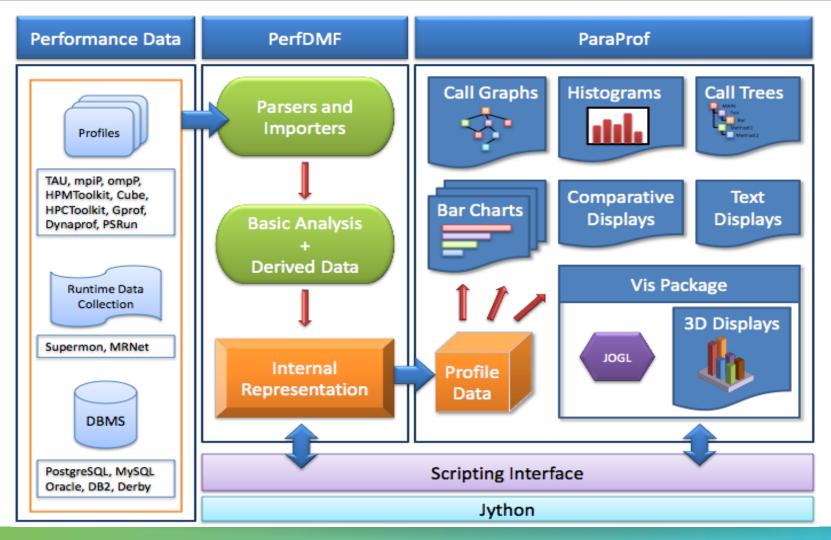
- Install Java 11.0.3:
 - Download <u>http://tau.uoregon.edu/java.dmg</u>
 - If you have multiple Java installations, add to your ~/.zshrc (or ~/.bashrc as appropriate):
 - export PATH=/Library/Java/JavaVirtualMachines/jdk-11.0.3.jdk/Contents/Home/bin:\$PATH
 - java -version
- Download and install TAU (copy to /Applications from dmg):
 - http://tau.uoregon.edu/tau.dmg
 - export PATH=/Applications/TAU/tau/apple/bin:\$PATH
 - paraprof app.ppk &
- macOS (arm64, M1/M2)
 - <u>http://tau.uoregon.edu/java_arm64.dmg</u>
 - http://tau.uoregon.edu/tau_arm64.dmg
- Linux (http://tau.uoregon.edu/tau.tgz)
 - ./configure; make install; export PATH=<taudir>/x86_64/bin:\$PATH
 - paraprof app.ppk &

TAU's Analysis Tools: ParaProf

TAU Analysis



ParaProf Profile Analysis Framework

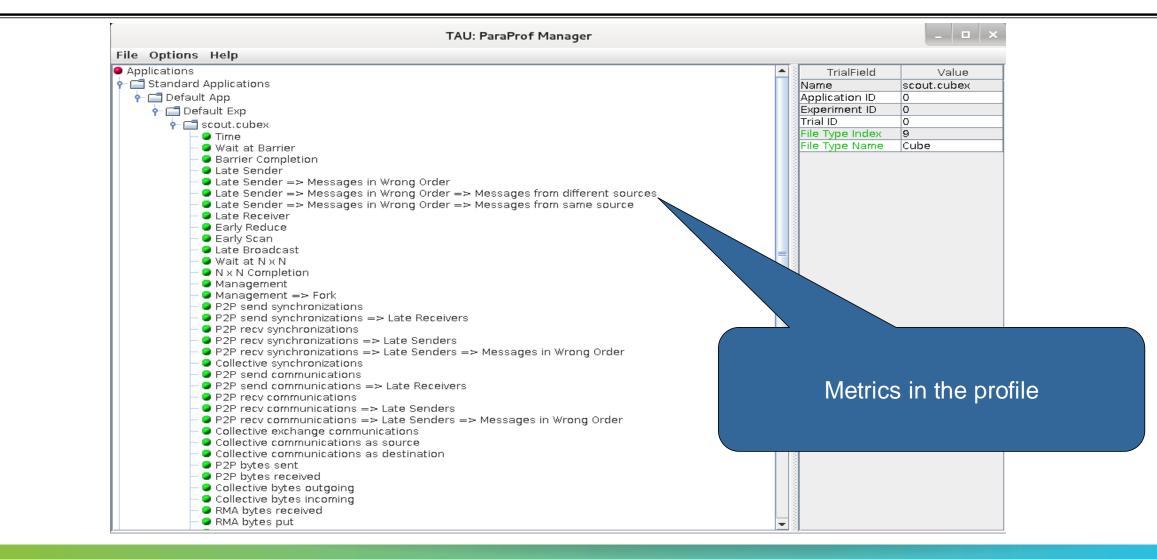


SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)

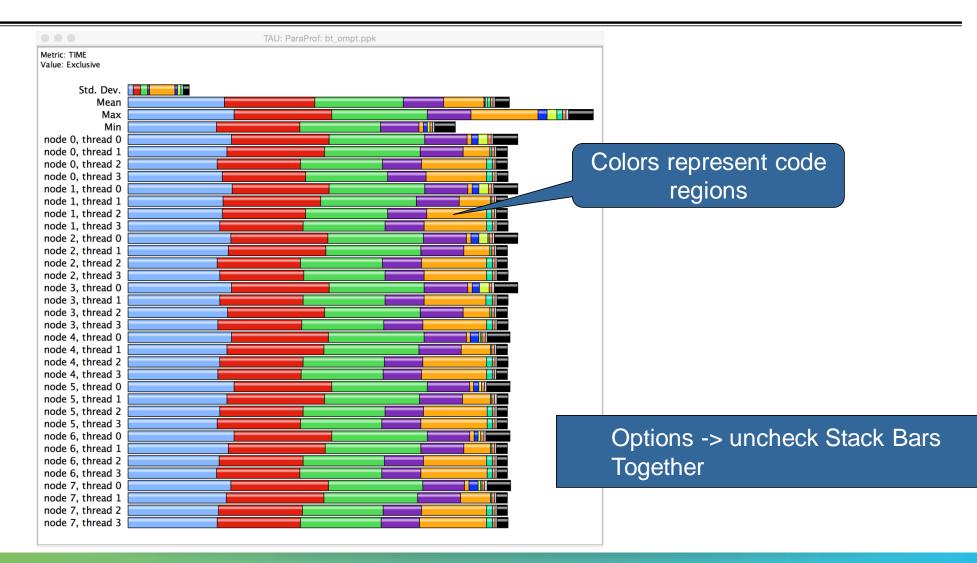
TAU Analysis Tools: paraprof

	Applications	TrialField	Value
Launch paraprof		Name	bt_ompt.ppk
Edditeri parapior	Standard Applications	Application ID	0
	🔻 📩 Default App	Experiment ID	0
	🔻 🚞 Default Exp	Trial ID	0
% paraprof	🔻 🥥 bt_ompt.ppk	CPU Cores	8
	TIME	CPU MHz	2600.000
	Default (jdbc:h2:/Users/sameer/.ParaProf/perfdmf/perfdmf;AUTO_SERVER=TRUE)	CPU Type	Intel(R) Xeon(R) CPU E5-2670 0 @ 2.60GHz
		CPU Vendor	GenuineIntel
		CWD	/scratch/sameer/NPB3.3-MZ-MPI/bin
		Cache Size	20480 КВ
		Command Line	./bt-mz_C.8
		Executable	/scratch/sameer/NPB3.3-MZ-MPI/bin/bt-mz_C.8
		File Type Index	0
		File Type Name	ParaProf Packed Profile
Metric		Hostname	frog9
Wiethe		Local Time	2015-05-18T00:37:38+02:00
		MPI Processor Name	frog9
		Memory Size	65944056 kB
		Node Name	frog9
		OMP_CHUNK_SIZE	1
		OMP_DYNAMIC	off
		OMP_MAX_THREADS	4
		OMP_NESTED	off
		OMP_NUM_PROCS	4
		OMP_SCHEDULE	UNKNOWN
		OS Machine	x86_64
		OS Name	Linux
		OS Release	2.6.32-279.5.2.bl6.Bull.33.x86_64
		OS Version	#1 SMP Sat Nov 10 01:48:00 CET 2012

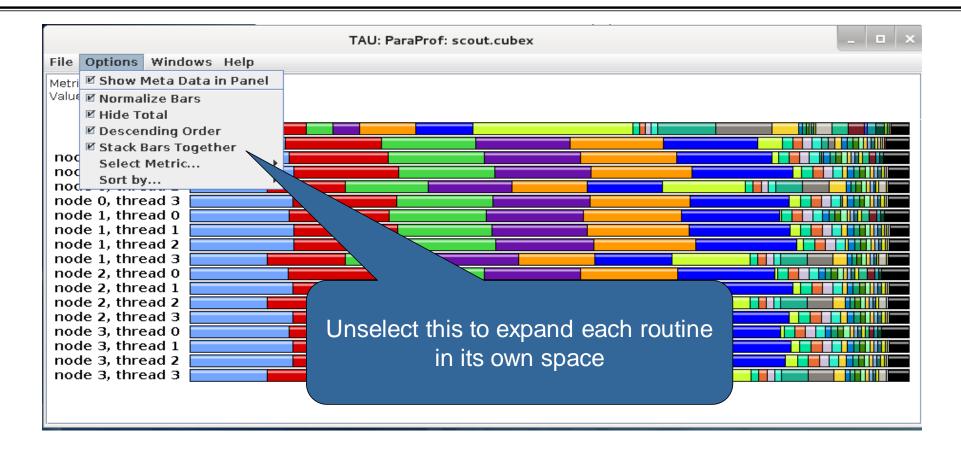
ParaProf Manager Widow: scout.cubex



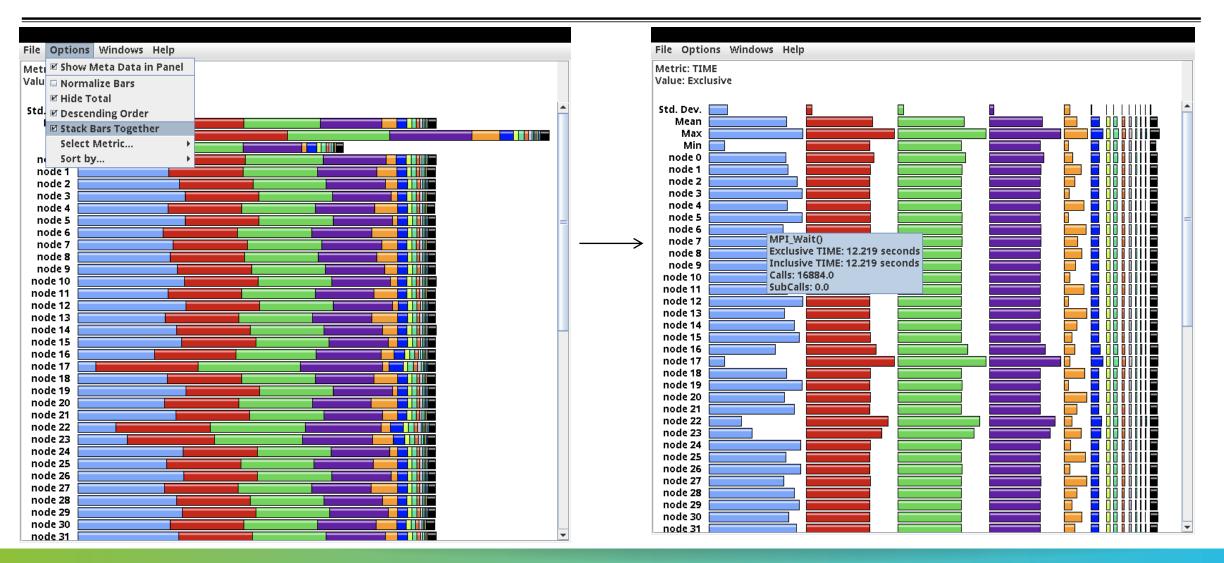
Paraprof main window



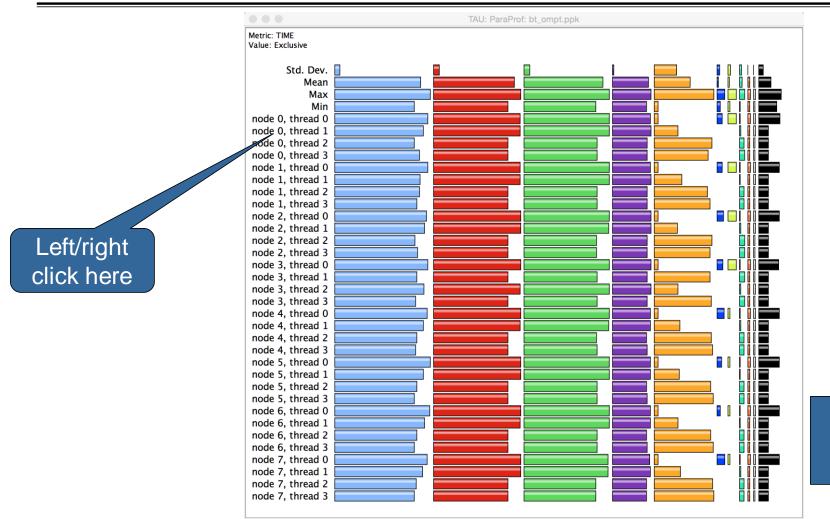
Paraprof main window



ParaProf Profile Browser



Paraprof main window



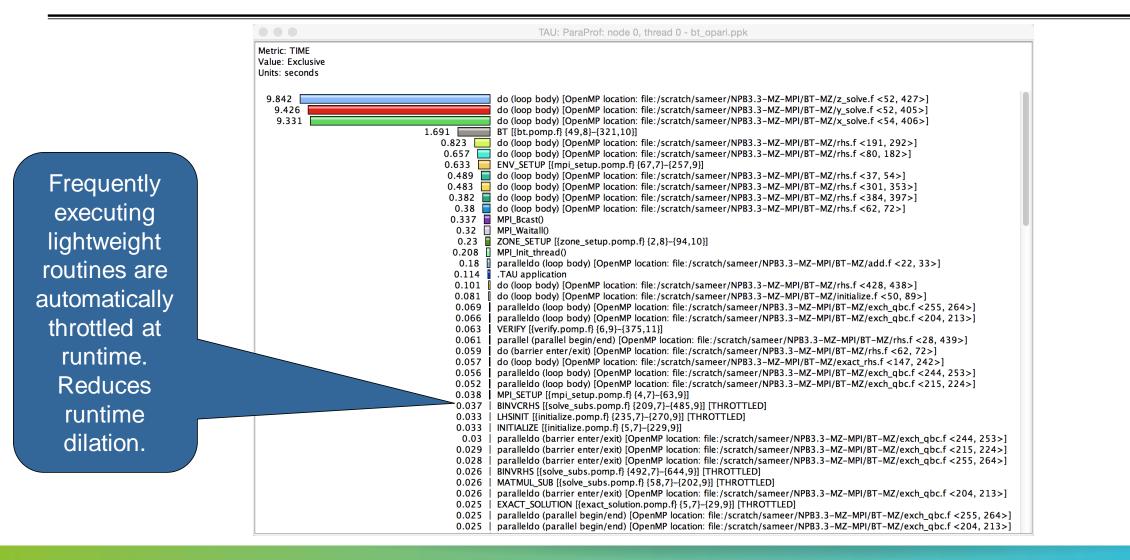
Each routine occupies its own space. Can see the extent of imbalance across all threads.

XXXXXXXXXXXXX × × × × × × VIRTUAL×INSTITUTE ~ HIGH PRODUCTIVITY SUPERCOMPUTING

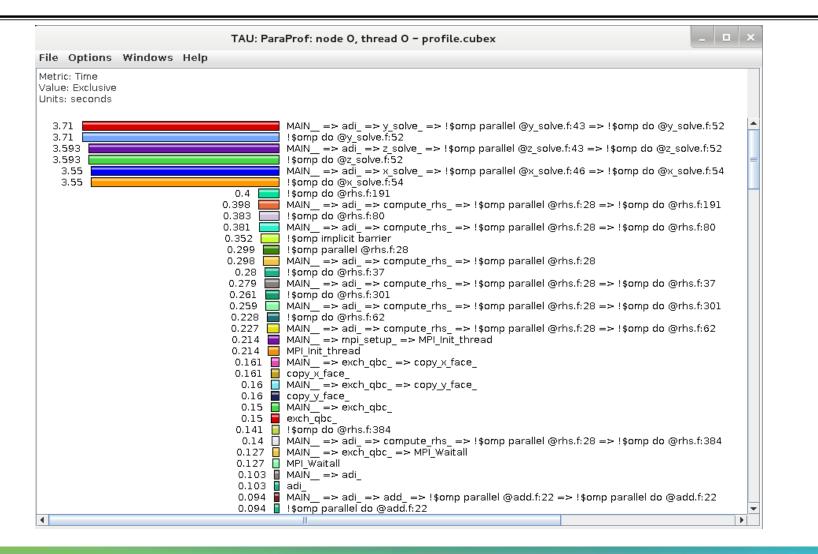
Paraprof node window (function barchart window)

	TAU: ParaProf: node 0, thread 1 - bt_ompt.ppk		
	Metric: TIME Value: Exclusive Units: seconds		
Exclusive time spent in each	8.214 8.038 7.899 3.549 2.223 0.206 0.184 0.147 0.098 0.085 0.075 0.066 0.065 0.066	 OpenMP_LOOP: L_add_22_par_loop0_2_19 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/add.f} {22,0}] OpenMP_WAIT_BARRIER: L_compute_rhs_28_par_region0_2_306 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/rhs.f} {28,0}] OpenMP_BARRIER: L_compute_rhs_28_par_region0_2_306 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/rhs.f} {28,0}] OpenMP_LOOP: L_copy_x_face_255_par_loop1_2_87 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {255,0}] OpenMP_LOOP: L_copy_y_face_204_par_loop0_2_176 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {204,0}] OpenMP_LOOP: L_copy_y_face_215_par_loop1_2_211 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {215,0}] OpenMP_LOOP: L_copy_x_face_244_par_loop0_2_54 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {214,0}] 	
code region (OpenMP loop) is shown here for MPI rank 0	0.049 0.049 0.048 0.047 0.02 0.029 0.029 0.021 0.021	 OpenMP_IMPLICIT_TASK: L_copy_x_face_255_par_loop1_2_87 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {255,0}] OpenMP_IMPLICIT_TASK: L_copy_y_face_215_par_loop1_2_211 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {215,0}] OpenMP_IMPLICIT_TASK: L_copy_x_face_244_par_loop0_2_54 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {244,0}] OpenMP_LOOP: L_initialize_28_par_region0_2_193 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/y_solve_f} {43,0}] OpenMP_IMPLICIT_TASK: L_y_solve_43_par_region0_2_44 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve_f} {43,0}] OpenMP_IMPLICIT_TASK: L_z_solve_46_par_region0_2_43 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve_f} {43,0}] OpenMP_IMPLICIT_TASK: L_x_solve_46_par_region0_2_43 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve_f} {43,0}] OpenMP_IMPLICIT_TASK: L_x_solve_146_par_region0_2_43 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve_f} {46,0}] OpenMP_IMPLICIT_TASK: L_add_22_par_loop0_2_19 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve_f} {42,0}] 	
thread 1	0.02 0.02 0.02 0.02 0.01 0.01 0.01 0.01	 5 OpenMP_BARRIER: L_copy_y_face_215_par_loop1_2211 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {215,0}] 5 OpenMP_BARRIER: L_copy_x_face_255_par_loop1_2.87 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {255,0}] 5 OpenMP_BARRIER: L_copy_x_face_244_par_loop0_2.54 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {244,0}] 5 OpenMP_BARRIER: L_copy_y_face_204_par_loop0_2.176 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {204,0}] 6 OpenMP_BARRIER: L_y_solve_43_par_region0_2_43 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/exch_qbc.f} {204,0}] 8 OpenMP_BARRIER: L_y_solve_44_par_region0_2_43 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/y_solve.f} {43,0}] 9 OpenMP_BARRIER: L_z_solve_43_par_region0_2_44 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {46,0}] 9 OpenMP_BARRIER: L_z_solve_43_par_region0_2_44 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43,0}] 9 OpenMP_BARRIER: L_z_solve_43_par_region0_2_19 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43,0}] 	

Instrumenting Source Code with PDT and Opari

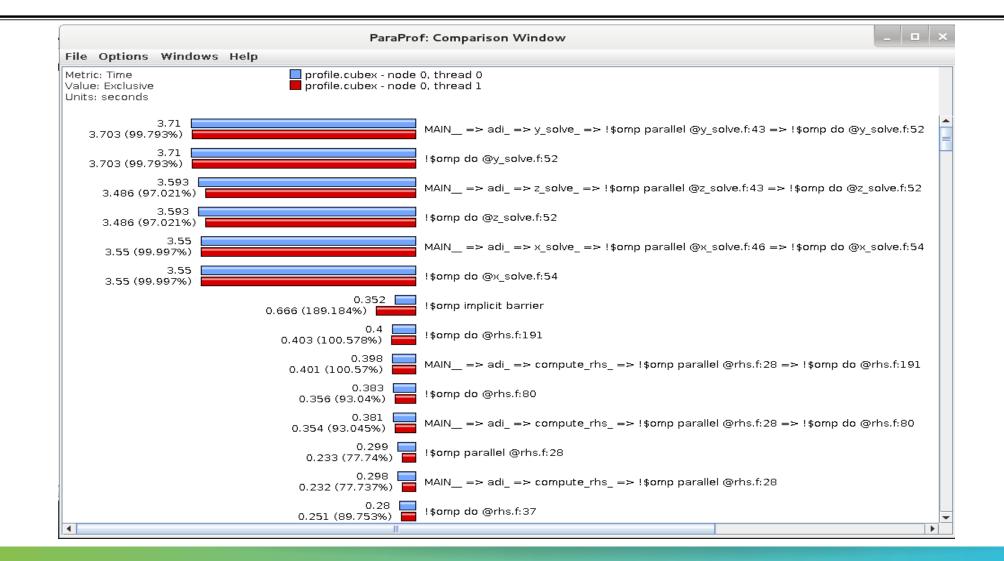


ParaProf: Node view in a callpath profile

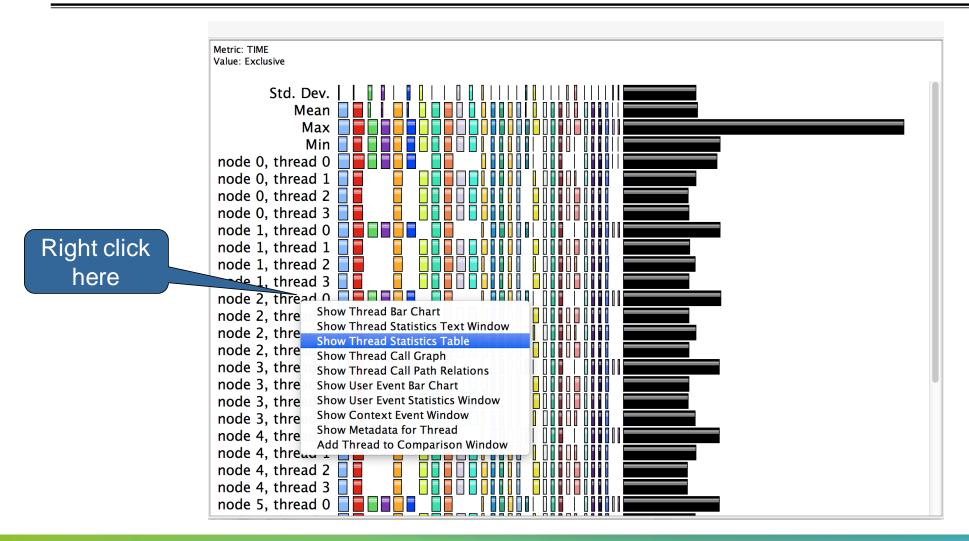


VICTOR COMPUTING

ParaProf: Add thread to comparison window



Paraprof Thread Statistics Table with TAU_SAMPLING=1



ParaProf: Thread Statistics Table

TAU: ParaProf: Sta	tistics for: node 0, thread 0 - scout.c	ubex		_ 🗆 ×
File Options Windows Help				
Time	▼			
Name	Exclusive Time 🗸 🥿	Inclusive Time	Calls	Child Calls
\$00 minimum solve.f:52	5.817	5.817	3,216	0
-solve.f:52	5.657	5.657	3,216	0
- 🔄 !\$omp do @x_solve.f:54	5.609	5,609	3,216	0
- _ !\$omp do @rhs.f:191	0.609	20	3,232	0 _
-somp do @rhs.f:80	0.583		3,232	0
– MPI_Waitall	0.402	6	603	0
-somp implicit barrier	0.402		_	
• 🗖 !\$omp do @rhs.f:301	0.36			
	0.026	Click to	sort by a	a given metric, d
\$omp implicit barrier	0			
- 🔤 !\$omp do @rhs.f:37	0.343	and m	iove to r	earrange columr
<mark>⊱ </mark>]\$omp do @rhs.f:62	0.225			
	0.004	0.004	3,210	U
\$omp implicit barrier	0	0	16	0
– <mark>–</mark> MPI_Init_thread	0.218	0.218	1	0
- <mark>-</mark> !\$omp do @rhs.f:384	0.199	0.199	3,232	0
🗠 🗖 !\$omp parallel do @add.f:22	0.099	0.111	3,216	3,216
- <mark>-</mark> !\$omp do @rhs.f:428	0.069	0.069	3,232	0
– MPI_Isend	0.043	0.043	603	0
	0.04	0.04	32	0
🗠 🗖 !\$omp parallel @rhs.f:28	0.03	2.536	3,232	51,712
\$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$	0.021	0.029	6,432	6,432
🗠 🗖 !\$omp parallel do @exch_qbc.f:255	0.02	0.033	6,432	6,432
\$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$ \$\$	0.02	0.053	6,432	6,432
\$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$ \$	000		FinderScreen	nSnapz003.png

ParaProf

- Click on Columns:
- to sort by incl time
- Open binvcrhs
- Click on Sample

TAU: ParaProf: Statistics for: node 0 - /rwthfs/rz/cluster/	work/hpclab17/NPB3.3	3-MZ-MPI/bin		
ile Options Windows Help				
Name	Exclusive TIME	Inclusive TIME V	Calls	Child Calls
TAU application	9.167	9.308	1	2,4
- CONTEXT] .TAU application	0	9.019	901	
[SUMMARY] binvcrhs [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ,	2.89	2.89	288	
SUMMARY] matmul_sub_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT	1.27	1.27	127	
SUMMARY] x_solve [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/x	1.16	1.16	116	
SUMMARY] z_solve_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/z	1.08	1.08	108	
SUMMARY] y_solve_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/y	1.08	1.08	108	
SUMMARY] compute_rhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/B	0.83	0.83	83	
SUMMARY] matvec_sub_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT	0.49	0.49	49	
SUMMARY] lhsinit_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/in	0.08	0.08	8	
SAMPLE] add [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/add.f}	0.05	0.05	5	
SUMMARY] binvrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/s	0.04	0.04	4	
SUMMARY] exact_solution_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/	0.02	0.02	2	
SAMPLE] copy_x_face [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ	0.01	0.01	1	
SUMMARY] exact_rhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-M.	0.01	0.01	1	
[SAMPLE] initialize_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/in	0.009	0.009	1	
MPI_Init_thread()	0.155	0.155	1	
MPI_Finalize()	0.022	0.022	1	
MPI_Waitall()	0.018	0.018	804	
MPI_Irecv()	0.004	0.004	804	
MPI_Isend()	0.001	0.001	804	
MPI_Comm_split()	0	0	1	
MPI_Bcast()	0	0	9	
MPI_Reduce()	0	0	3	
MPI_Barrier()	0	0	2	
MPI_Comm_size()	0	0	1	
MPI_Comm_rank()	0	0	2	

Paraprof Thread Statistics Table

Name TAU application DenMP_PARALLEL_REGION: L_z_solve_43_par_region0_2_44 DenMP_IMPLICIT_TASK: L_z_solve_43_par_region0_2_44 [{ COPENDP_LOOP: L_z_solve_43_par_region0_2_44 [{ context] OpenMP_LOOP: L_z_solve_43_par_region0_ v = [CONTEXT] OpenMP_LOOP: L_z_solve_43_par_region0_ v = [SUMMARY] L_z_solve_43_par_region0_2_44 [{ context] Context] Context] context]	[{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {4 /scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43 :h/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43,0}] _2_44 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve	,0}]	re TIME 1.754 0.061 0.04 8.528	Inclusive TIME ⊽ 36.26 8.692 8.568	Calls 1 6,432 6,432	Child Calls 88,049 12,864
 OpenMP_PARALLEL_REGION: L_z_solve_43_par_region0_2_44 OpenMP_IMPLICIT_TASK: L_z_solve_43_par_region0_2_44 [OpenMP_LOOP: L_z_solve_43_par_region0_2_44 [(CONTEXT] OpenMP_LOOP: L_z_solve_43_par_region0_2_44 [SUMMARY] L_z_solve_43_par_region0_2_44 [<pre>//scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43 h/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43,0}] 2_44 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve</pre>	,0}]	0.061 0.04	8.692 8.568		12,864
 OpenMP_IMPLICIT_TASK: L_z_solve_43_par_region0_2_44 [OpenMP_LOOP: L_z_solve_43_par_region0_2_44 [CONTEXT] OpenMP_LOOP: L_z_solve_43_par_region0_ [SUMMARY] L_z_solve_43_par_region0_2_44 [<pre>//scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43 h/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43,0}] 2_44 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve</pre>	,0}]	0.04	8.568		
 OpenMP_LOOP: L_z_solve_43_par_region0_2_44 [{/scrate [CONTEXT] OpenMP_LOOP: L_z_solve_43_par_region0_ [SUMMARY] L_z_solve_43_par_region0_2_44 [{/scrate 	h/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {43,0}] 2_44 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve				6,432	
 CONTEXT] OpenMP_LOOP: L_z_solve_43_par_region0_ [SUMMARY] L_z_solve_43_par_region0_2_44 [{/scrat 	2_44 [{/scratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve	e fl (43 0)]	8.528			6,432
SUMMARY] L_z_solve_43_par_region0_2_44 [{/scrat		o fl (43 01)		8.528	6,432	0
	ch/sameer/NPB3 3-M7-MPI/BT-M7/z solve f3]	e.ij [+J,0]]	0	9.23	847	0
			3.67	3.67	340	0
	ch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f}]		3.67	3.67	340	0
	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {419}]	Show Source Code	0.22	0.22	21	0
	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} [58]	Show Function Bar Chart	0.17	0.17	16	0
	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} [418]	Show Function Histogram Assign Function Color	0.16	0.16 0.11	12 11	0
	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {123}] atch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {193}]	Reset to Default Color	0.11	0.11	5	0
	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/Z_solve.f} [193]		0.03	0.07	7	0
	atch/sameer/NPB3.3–MZ–MPI/BT–MZ/Z_solve.f} {247}]		0.07	0.07	6	0
SAMPLE] L z solve 43 par region0 2 44 [/sc	atch/sameer/NPB3.3–MZ–MPI/BT–MZ/z_solve.f} {158}]		0.06	0.06	5	0
	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {313}]		0.06	0.06	4	0
SAMPLE] L_z_solve_43_par_region0_2_44 [{/sc	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {230}]		0.06	0.06	4	0
Choose SAMPLE] L_z_solve_43_par_region0_2_44 [{/scr	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {308}]		0.05	0.05	3	0
SAMPLE] L_z_solve_43_par_region0_2_44 [{/sci	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {191}]		0.05	0.05	3	0
"Show	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {81}]		0.05	0.05	4	0
■ [SAMPLE] L_Z_solve_43_par_region0_2_44 [{/sci	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {301}]		0.05	0.05	5	0
	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {67}]		0.05	0.05	5	0
[SAMPLE] L_Z_SOIVE_43_par_region0_2_44 [{/Sci	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {175}]		0.04	0.04	4	0
	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {89}]		0.04	0.04	4	0
SAMPLE] L_Z_SOIVE_43_par_regionU_Z_44 [/SCI	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {55}]		0.04	0.04	4	0
	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {275}]		0.04	0.04	4	0
	ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {129}] ratch/sameer/NPB3.3-MZ-MPI/BT-MZ/z_solve.f} {168}]		0.04 0.04	0.04 0.04	4	0
	atch/sameer/NPB3.3-MZ-MPI/BT-MZ/Z_solve.1} [108]] atch/sameer/NPB3.3-MZ-MPI/BT-MZ/Z_solve.f} [238]]		0.04	0.04	4	0

ParaProf

X TAU: ParaProf: Statistics for: node 0 - /rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/bin

File Options Windows Help

Name	Exclusive TIME	Inclusive TIME 🗸	Calls	Child Calls
P-■.TAU application	9.1	67 9.368	1	2,432 📤
- CONTEXT] .TAU application		0 9.019	901	0
[SUMMARY] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f}]	2	89 2.89	288	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f}_{228}]	0	14 0.14	14	0
SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} Show Sou	urce Code 0	09 0.09	9	0
SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} Show In S		09 0.09	9	0
SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} Show Fur	nction Histogram 0	06 0.06	6	0
SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} Show Fur	nction Bar Chart 0	06 0.06	6	0
- [SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} Assign Fu	unction Color 0	06 0.06	6	0
	Default Color 0	06 0.06	6	0
SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} - [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f}	0	05 0.05	5	0
SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {332}]	0	05 0.05	5	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {275}]	0	05 0.05	5	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {331}]	0	04 0.04	4	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {445}]	0	04 0.04	4	0
SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {254}]	0	04 0.04	4	0
[SAMPLE] binvcrhs_ [{/wthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {314}]	0	04 0.04	4	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {343}]		04 0.04	4	0
[SAMPLE] binvcrhs_ [{/wthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {403}]	0	04 0.04	4	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {389}]	0	03 0.03	3	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {415}]	0	03 0.03	3	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {247}]	0	03 0.03	3	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {300}]	0	03 0.03	3	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {309}]	0	03 0.03	3	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {444}]	0	03 0.03	3	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {468}]	0	03 0.03	3	0
[SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {242}]		03 0.03	3	0
SAMPLE] binvcrhs_ [{/rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f} {407}]		03 0.03	3	0
SAMPLE1 binvcrhs_l{/wthfs/rz/cluster/work/hoclab17/NPB3.3-MZ-MPI/BT-MZ/solve_subs.f}_{412}	0	030.03	3	0 //.

san

in

Statement Level Profiling with TAU

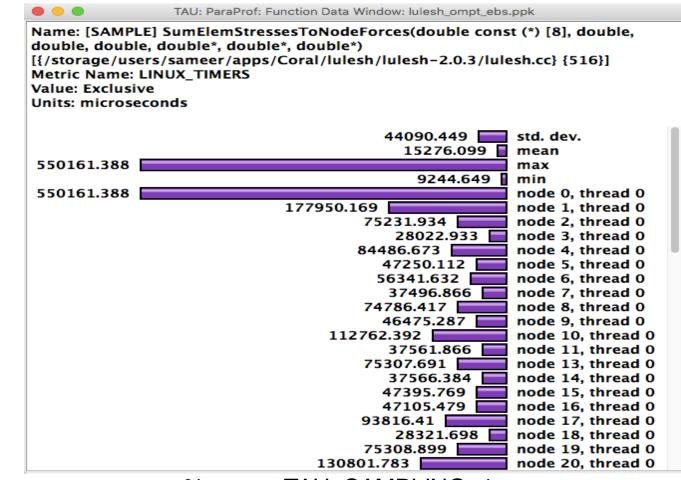
	TAU: ParaProf: Source Browser: /scratch/sameer/NPB3.3	-MZ-MPI/BT-MZ/x_solve.f
	File Help	
	353 call matmul_sub(lhs(1,1,aa,i)	
	354 > lhs(1,1,cc,i- base > lhs(1,1,bb,i)	
	555)
	356	
	357	
	358 C 359 C multiply c(i,j,k) by b_inverse and cop	w back to c
	a $multiply mba(1 + k) by b inverse(1 + k)$	
) and copy to this
	sol	
	rha(1 + i + k)	
	554	
	365 enddo	
	566	
	367	
	368 C	1)
	370 C	
	371 call matvec_sub(lhs(1,1,aa,isize	:),
	372 > rhs(1,isize-1	.,j,k),rhs(1,isize,j,k))
	373	
ource	374 C	
ouroc	375 C B(isize) = B(isize) - C(isize-1)*A(isi	.ze)
	376 C	、
cation	377 call matmul_sub(lhs(1,1,aa,isize	
	378 > lhs(1,1,cc,is 379 > lhs(1,1,bb,is	
ubara		126))
vhere	380	
	381 C	to the
ples are		
ipies ale	383 Ccall binvrhs(lhs(1,1,bb,isize),	
· .	384 Call Dinvrns(Lns(1,1,DD,1size), 385 > rhs(1,isize,j,k	
aken.	386	, ,
	387	
	388 C	
ompute	389 c back solve: if last cell, then generat	e U(isize)=rhs(isize)
	390 c else assume U(isize) is loaded in un p	ack backsub_info
	₃₉₁ c so just use it	
ensive 📃 🗌	₃₉₂ c after call u(istart) will be sent to n	ext cell
	393 C	
	394	
egion.	do i=isize-1,0,-1	
	do m=1,BLOCK_SIZE	
	do n=1, BLOCK_SIZE	3 6 3
	rhs(m,i,j,k) = rhs(m,i,j,k)],K)
	399 > - lhs(m,n,cc,i)*rh	s(n,1+1,j,k)
	400 en ddo	
	401 enddo	
	402 enddo	

ParaProf Comparison Window



SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)

TAU – Event Based Sampling (EBS)



% export TAU_SAMPLING=1

Examples: Callstack Sampling in TAU

TAU: ParaProf: Statistics for: n,c,t 2,0,0 - gamess_unw_call_ebs.ppk		
Name		Calls
■ Name	Inclusive TIME 79.592	
▼ ■ MPI_Recv()	75.607	6,870
<pre>The second second</pre>	74.848	1,497
[UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@] MAIN_ [{/gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@] MAIN_ [{/gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@]		524
[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_fortran.c.67 [@] beging_ [{/gpfs/mira-home/sameer/gamess/theta-tau/ddi/src/ddi_fortran.c.67 [@] beging_ [{/gpfs/mira-home/sameer/gamess/theta-tau/ddi/src/ddi_fortran.c.67 [@] beging_ []		434
UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/gamess.f.538 [@] main [{/gpfs/mira-home/sameer/gamess-theta-tau/object/gamess.f.538 [@]		237
[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi init.c.113 [@] ddi init [{/gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi init.c.113 [@] ddi init [{/gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/gamess-theta-tau/gamess-theta-tau/gamess-theta-tau/gamess-theta-tau/gamess-theta-tau/gamess-theta-tau/gamess-theta-tau/gamess-theta-tau/		174
IUNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [@] DDI_Init [{/gpfs/mira-home/yuri/dist/C		115
UNWIND] /lib64/libc-2.22.so.0 [@] _start [{/home/abuild/rpmbuild/BUILD/glibc-2.22/csu//sysdeps/x86_64/start.S} {118}]	0.2	4
[SAMPLE] GNII_DIaProgress [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	0.2	4
[UNWIND] [/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0] [@] UNRESOLVED UNKNOWN	0.15	3
[SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	0.051	1
[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt/7	0.05	1
MPI_Finalize()	3.601	1
MPI_Send()	0.122	6,866
MPI_Init_thread()	0.112	1
[CONTEXT].TAU application	0.05	1
MPI_Bcast()	0.014	6
MPI_Allgather()	0.004	3
MPI_Barrier()	0.003	7
MPI_Comm_create()	0.002	4
MPI_Gather()	0.002	1
MPI_Comm_split()	0.002	1
MPI_Group_intersection()	0.001	1
MPI_Comm_group()	0.001	1
MPI_Group_incl()	0	3
MPI_Comm_rank()	0	6
MPI_Comm_size()	0	2
% export TAU_SAMPLING=1; export TAU_EBS_UNWIND=1		

UNWINDING CALLSTACKS

TAU: ParaProf: Statistics for: n,c,t 2,0,0 - gamess_unw_call_ebs.ppk		
Name	Inclusive TIME V	Calls
.TAU application	79.592	
MPI_Recv()	75.607	6,87
CONTEXT] MPI_Recv()	74.848	1,49
UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@] MAIN_ [{/gpfs/mira-home/sameer/gamess-theta		52
🔻 🗖 [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_fortran.c.67 [@] beging_ [{/gpfs/mira-home/sameer/	g 21.7	43
🔻 🗖 [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [@] ddi_init_ [{/gpfs/mira-home/yuri/dis		43
🔻 🗖 [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [@] DDI_Init [{/gpfs/mira-home/yuri	/ 21.7	43
🔻 🗖 [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_recv.c.65 [@] DDI_Server [{/gpfs/mira-home/y	/ 21.7	43
🔻 🗖 [UNWIND] /lus/theta-fs0/software/perftools/tau/tau-2.26.3/src/Profile/TauMpi.c.2371 [@] DDI_Recv_request [{/gpfs/mir	a 21.7	43
🔻 🗖 [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPI_Recv [{/lus/theta-fs0/so	fi 21.7	43
🔻 🗖 [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] PMPI_Recv [{/opt/cray/pe/	n 21.7	43
🔻 🗖 [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPIDI_CH3I_Progress [{/	21.45	42
🔻 🗖 [UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [@] MPID_nem_gni_poll [{	/ 15.95	31
[SAMPLE] GNI_SmsgGetNextWTag [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0}	10.349	20
[SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	5.6	11
[UNWIND] gni_poll.c.0 [@] MPID_nem_gni_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel	e 5.25	10
[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPID_nem_gni_poll [{	/ 0.25	
[UNWIND] UNRESOLVED [@] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel/16.0/lib/lib/libmpich_intel/16.0/lib/lib/libmpich_intel/16.0/lib/lib/libmpich_intel/16.0/lib/lib/lib/lib/lib/lib/lib/lib/lib/lib	t 0.25	
🕨 🔲 [UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/gamess.f.538 [@] main [{/gpfs/mira-home/sameer/gamess-theta-t	a 11.85	2
UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [@] ddi_init_ [{/gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [@]	G 8.701	17
[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [@] DDI_Init [{/gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [@] DDI_Init [{/gpfs/mira-home/yuri/dist/Server.c.99 [@] Server.c.99 [@] DDI_Init [{/gpfs/mira-home/yuri/dist/Server.c.99 [@] Server.c.99 [@]	5.75	1
UNWIND] /lib64/libc-2.22.so.0 [@] _start [{/home/abuild/rpmbuild/BUILD/glibc-2.22/csu//sysdeps/x86_64/start.S} {118}]	0.2	
[SAMPLE] GNII_DIaProgress [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	0.2	
[UNWIND] [/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0] [@] UNRESOLVED UNKNOWN	0.15	
[SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	0.051	
[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt/	0.05	
MPI_Finalize()	3.601	
▶ ■ MPI_Send()	0.122	6,8
MPI_Init_thread()	0.112	
CONTEXT] .TAU application	0.05	

% export TAU_SAMPLING=1; export TAU_EBS_UNWIND=1

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

UNWINDING CALLSTACKS

TAU: ParaProf: Statistics for: n,c,t 2,0,0 - gamess_unw_call_ebs.ppk		
Name	Inclusive TIME V	Calls
TAU application	79.592	
MPI_Recv()	75.607	6,87
CONTEXT] MPI_Recv()	74.848	1,49
[UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@] MAIN_ [{/gpfs/mira-home/sameer/gamess-theta	- 26.196	52
[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_fortran.c.67 [@] beging_ [{/gpfs/mira-home/sameer/	g 21.7	43
🔻 🗖 [UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/gamess.f.538 [@] main [{/gpfs/mira-home/sameer/gamess-theta-t	ta 11.85	23
🔻 🗖 [UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@] MAIN [{/gpfs/mira-home/sameer/gamess-the	et 11.85	23
🔻 🗖 [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_fortran.c.67 [@] beging_ [{/gpfs/mira-home/san	n 11.85	23
🔻 🗖 [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [@] ddi_init_ [{/gpfs/mira-home/yu	r 11.85	23
🔻 🗖 [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [@] DDI_Init [{/gpfs/mira-home	./ 11.85	23
🔻 🗖 [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_recv.c.65 [@] DDI_Server [{/gpfs/mira-ho	11.85	23
🔻 🗖 [UNWIND] /lus/theta-fs0/software/perftools/tau/tau-2.26.3/src/Profile/TauMpi.c.2371 [@] DDI_Recv_request [{/gpf	s 11.85	23
🔻 🗖 [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPI_Recv [{/lus/theta-fs	s(11.85	23
[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] PMPI_Recv [{/opt/cray	y, 11.7	23
[SAMPLE] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1} {	0 11.3	22
[SAMPLE] MPIDU_Sched_are_pending [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.	.0 0.2	
[SAMPLE] MPID_nem_gni_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1} {0}] 0.15	
[SAMPLE] MPID_nem_network_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.	1 0.05	
[UNWIND] ch3_progress.c.0 [@] PMPI_Recv [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.sc	o. 0.15	
UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [@] ddi_init_ [{/gpfs/mira-home/yuri/dist/	G 8.701	17
UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [@] DDI_Init [{/gpfs/mira-home/yuri/dist	/ 5.75	1:
UNWIND] /lib64/libc-2.22.so.0 [@] _start [{/home/abuild/rpmbuild/BUILD/glibc-2.22/csu//sysdeps/x86_64/start.S} {118}]	0.2	
[SAMPLE] GNII_DlaProgress [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	0.2	
UNWIND] [/opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so.0.6.0.0] [@] UNRESOLVED UNKNOWN	0.15	
[SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0} {0}]	0.051	
[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt	/: 0.05	
MPI_Finalize()	3.601	
MPI_Send()	0.122	6,86
MPI_Init_thread()	0.112	
CONTEXT] .TAU application	0.05	

Deep Learning: Tensorflow

TAU: ParaProf: Statistics for: node 0, thread 8 - nt3_baseline_keras2.ppk		
Name	Inclusiv	Calls ⊽
 TAU application 	519.211	1
[CONTEXT].TAU application	509.222	50,915
SAMPLE] Eigen::internal::gebp_kernel <float, 0="" 0,="" eigen::internal::blas_data_mapper<float,="" float,="" long,="">,</float,>	240.632	24,089
[SAMPLE]pthread_cond_wait [{} {0}]	86.384	8,634
[SAMPLE] Eigen::internal::gemm_pack_rhs <float, eigen::internal::tensorcontractionsubmapper<float,="" long,="" lor<="" p=""></float,>	51.345	5,135
[SAMPLE] Eigen::internal::gemm_pack_rhs <float, eigen::internal::tensorcontractionsubmapper<float,="" long,="" lor<="" p=""></float,>	24.375	2,416
[SAMPLE] void tensorflow::SpatialMaxPoolWithArgMaxHelper <eigen::threadpooldevice, float="">(tensorflow::Opk</eigen::threadpooldevice,>	16.301	1,630
[SAMPLE]memset_sse2 [{} {0}]	13.446	1,336
[SAMPLE] Eigen::TensorEvaluator <eigen::tensorcontractionop<eigen::array<eigen::indexpair<long>, 1ul> co</eigen::tensorcontractionop<eigen::array<eigen::indexpair<long>	5.99	599
[SAMPLE] long Eigen::internal::operator/ <long, false="">(long const&, Eigen::internal::TensorIntDivisor<long, false)<="" p=""></long,></long,>	5.843	585
[SAMPLE] std::_Function_handler <void (long,="" eigen::internal::tensorexecutor<eigen::tensorassignop<<="" long),="" p=""></void>	5.377	538
[SAMPLE] floatvector Eigen::TensorEvaluator <eigen::tensorbroadcastingop<eigen::indexlist<int, eigen::typ<="" p=""></eigen::tensorbroadcastingop<eigen::indexlist<int,>	4.862	487
[SAMPLE] Eigen::TensorEvaluator <eigen::tensorcontractionop<eigen::array<eigen::indexpair<long>, 1ul> co</eigen::tensorcontractionop<eigen::array<eigen::indexpair<long>	4.775	478
[SAMPLE] Eigen::TensorEvaluator <eigen::tensorassignop<eigen::tensormap<eigen::tensor<float, 1,="" long=""></eigen::tensorassignop<eigen::tensormap<eigen::tensor<float,>	4.037	404
[SAMPLE] Eigen::internal::gemm_pack_lhs <float, eigen::internal::tensorcontractionsubmapper<float,="" lon<="" long,="" p=""></float,>	3.679	367
[SAMPLE] Eigen::internal::EvalRange <eigen::tensorevaluator<eigen::tensorassignop<eigen::tensormap<eigen< p=""></eigen::tensorevaluator<eigen::tensorassignop<eigen::tensormap<eigen<>	2.981	298
[SAMPLE] tensorflow::MaxPoolingOp <eigen::threadpooldevice, float="">::SpatialMaxPool(tensorflow::OpKernelCo</eigen::threadpooldevice,>	2.915	295
[SAMPLE] std::_Function_handler <void (long,="" eigen::internal::tensorexecutor<eigen::tensorassignop<<="" long),="" p=""></void>	2.91	291
[SAMPLE] std::_Function_handler <void (long,="" eigen::internal::tensorexecutor<eigen::tensorassignop<<="" long),="" p=""></void>	2.772	277
[SAMPLE] Eigen::internal::gemm_pack_lhs <float, eigen::internal::tensorcontractionsubmapper<float,="" long)<="" long,="" p=""></float,>	2.481	248
[SAMPLE] std::_Function_handler <void (long,="" eigen::internal::tensorexecutor<eigen::tensorassignop<<="" long),="" p=""></void>	2.148	215
[SAMPLE] void Eigen::internal::call_dense_assignment_loop <eigen::map<eigen::matrix<float, -1,="" -1;<="" 0,="" p=""></eigen::map<eigen::matrix<float,>	2.008	197
[SAMPLE] Eigen::NonBlockingThreadPoolTempl <tensorflow::thread::eigenenvironment>::WorkerLoop(int) [{/hc</tensorflow::thread::eigenenvironment>	1.999	200
[SAMPLE] Eigen::internal::ptranspose(Eigen::internal::PacketBlock <floatvector, 4="">&) [{crtstuff.c} {0}]</floatvector,>	1.919	192
[SAMPLE] Eigen::internal::gemm_pack_rhs <float, eigen::internal::tensorcontractionsubmapper<float,="" long,="" lor<="" p=""></float,>	1.607	160
[SAMPLE] Eigen::TensorEvaluator <eigen::tensorcontractionop<eigen::array<eigen::indexpair<long>, 1ul> co</eigen::tensorcontractionop<eigen::array<eigen::indexpair<long>	1.518	152

% tau_python -ebs nt3_baseline_keras2.py (CANDLE)

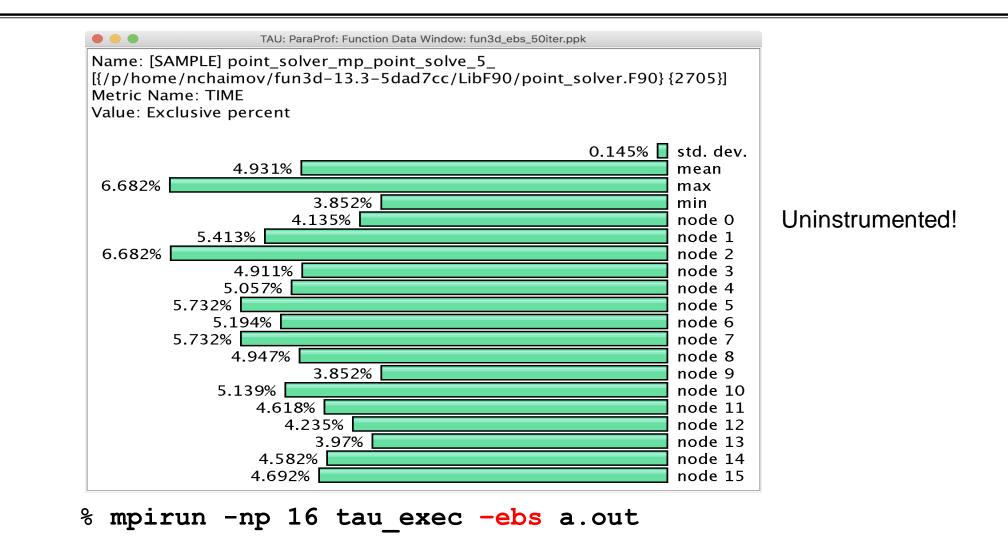


VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Sampling Tensorflow

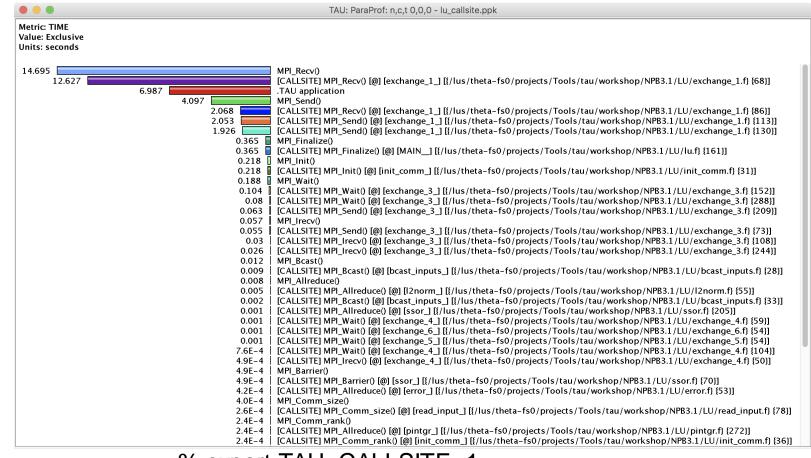
	TAU: ParaProf: Function Data Window: nt3_baseline_keras2.ppk
Eigen::interna Eigen::interna Eigen::Tenso Eigen::MakeP Eigen::array < false>::opera Eigen::Tenso Eigen::MakeP	TIME ve
53.463 50.094 53.463 50.193 52.872 51.145 52.442 52.618 51.345	15.44 std. dev. 5.144 mean max max 0.02 min node 0, thread 1 node 0, thread 2 node 0, thread 3 node 0, thread 3 node 0, thread 4 node 0, thread 4 node 0, thread 5 node 0, thread 5 node 0, thread 6 node 0, thread 7 node 0, thread 47 node 0, thread 47 node 0, thread 55 node 0, thread 55 0.05 node 0, thread 55 0.06 node 0, thread 55

Event Based Sampling (EBS)



VICTOR VICT

Callsite Profiling and Tracing





CALLPATH THREAD RELATIONS WINDOW

Sorted	Name: TIME By: Inclusive seconds			
	Exclusive	Inclusive	Calls/Tot.Calls	Name[id]
>	0.121	79.592	1	.TAU application
	0.002	0.002	1/1	MPI Gather()
	0.004	0.004	3/3	MPI Allgather()
	0.122	0.122	6866/6866	MPI Send()
	0.002	0.002	1/1	MPI Comm split()
	8.9E-5	8.9E-5	2/2	MPI Comm size()
	4.6E-4	4.6E-4	3/3	MPI Group incl()
	75.607	75.607	6870/6870	MPI Recv()
	0.002	0.002	4/4	MPI Comm create()
	9.5E-5	9.5E-5	6/6	MPI Comm rank()
	5.4E-4	5.4E-4	1/1	MPI Comm group()
	0.003	0.003	7/7	MPI_Barrier()
	0.112	0.112	1/1	MPI Init_thread()
	6.3E-4	6.3E-4	1/1	MPI_Group_intersection()
	0	0.05	1/1	[CONTEXT] .TAU application
	3.601	3.601	1/1	MPI_Finalize()
	0.014	0.014	6/6	MPI_Bcast()
	75.607	75.607	6870/6870	.TAU application
>	75.607	75.607	6870	MPI_Recv()
	0	74.848	1497/1497	[CONTEXT] MPI_Recv()
	0	74.848	1497/1497	MPI_Recv()
>	0	74.848	1497	[CONTEXT] MPI_Recv()
	0	8.701	174/1371	[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_init.c.113 [@] ddi_i
	0	26.196	524/763	[UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@] MAIN_ [/gpfs/mira-home/sameer/gamess-theta-tau/object/unport.f.410 [@] MAIN_ [/gpfs/mira-home/sameer/gamess-tau/object/unport.f.410 [@] MAIN_ [/gpfs/mira-home/sameer/gamess-tau/object/unp
	0.2	0.2	4/138	[SAMPLE] GNII_DlaProgress [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.
	0	5.75	115/1484	[UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi_server.c.99 [0] DDI_
	0	0.2	4/5	[UNWIND] /lib64/libc-2.22.so.0 [@] _start [{/home/abuild/rpmbuild/BUILD/glibc-2.22/csu//s
	-	11.85	237/239	[UNWIND] /gpfs/mira-home/sameer/gamess-theta-tau/object/gamess.f.538 [@] main [{/gpfs/mira- [SAMPLE] GNL GECCHEvent [/ent/gray/gray/gray] /6 0 14 6 0 4 0 14 1 go72b1a2 ari/lib(/liburgi co
	0.051	0.051 0.05	1/273 1/1197	[SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.sc [UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich intel.so.3.0.1.0 [@] MPIE
	0	0.05	3/7	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0 14.1 ge7db4a2.ari/lib64/libugni.so.0.6.0.0] [@] UN
	0	21.7	434/1197	[UNWIND] [/opt/cray/ugn1/6.0.14-6.0.4.0_14.1_ge/db4a2.ar1/11664/116ugn1.so.0.6.0.0] [0] 0N [UNWIND] /gpfs/mira-home/yuri/dist/Github/gamess-theta-tau/ddi/src/ddi fortran.c.67 [0] beg

CALLPATH THREAD RELATIONS WINDOW

orted	Name: TIME By: Exclusive seconds			
	Exclusive	Inclusive	Calls/Tot.Calls	Name[id]
>	75.607 75.607	75.607 75.607	6870/6870 6870	.TAU application MPI Recv()
	0	74.848	1497/1497	[CONTEXT] MPI_Recv()
	0.15	0.15	3/444	[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] PMPI_Re
>	22.046 22.196	22.046 22.196	441/444 444	[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] MPIDI_C [SAMPLE] MPID_nem_gni_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so
	5.6	5.6	112/273	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [@] MPID_ne
	0.051	0.051	1/273	[CONTEXT] MPI_Recv()
	7.651 0.35	7.651 0.35	153/273 7/273	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [@] MPID_n [UNWIND] [/opt/cray/ugni/6.0.14-6.0.4.0_14.1_ge7db4a2.ari/lib64/libugni.so.0.6.0.0] [@] UNRE
>	13.652	13.652	273	[SAMPLE] GNI_CqGetEvent [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so.0.
	11.3	11.3	226/226	[UNWIND] /opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel.so.3.0.1.0 [@] PMPI_R(
>	11.3	11.3	226	[SAMPLE] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib/libmpich_intel
	10.349 10.349	10.349 10.349	207/207	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [@] MPID_nd
>	10.349	10.349	207	[SAMPLE] GNI_SmsgGetNextWTag [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni
	0.2	0.2	4/138	[CONTEXT] MPI_Recv()
	6.701	6.701	134/138	[UNWIND] /opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so.0.6.0.0 [@] GNI_Cq
>	6.901	6.901	138	[SAMPLE] GNII_DlaProgress [{/opt/cray/ugni/6.0.14-6.0.4.0_14.1ge7db4a2.ari/lib64/libugni.so
	5.25	5.25	105/109	[UNWIND] gni_poll.c.0 [0] MPID_nem_gni_poll [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/lib
>	0.2 5.45	0.2 5.45	4/109 109	[UNWIND] gni_poll.c.0 [@] MPIDI_CH3I_Progress [{/opt/cray/pe/mpt/7.6.3/gni/mpich-intel/16.0/1; [SAMPLE] MPID nem gni check localCQ [{gni poll.c} {0}]
	5.45	5.45	107	[BANELD] MEID_Nem_GHI_CHECK_IOCAICY [{GHI_DOIICS {0}]
	3.601	3.601	1/1	.TAU application
>	3.601	3.601	1	MPI Finalize()

ParaProf: Callpath Thread Relations Window

	ons Windows H				
etric N	ame: Time				
orted E	By: Exclusive				
nits: se	econds				
	0.04	0.04	32/32	!\$omp parallel @initialize.f:28	
>	0.04	0.04	32	!\$omp do @initialize.f:50	
	0.03	2.536	3232/3232	compute_rhs_	
>	0.03	2,536	3232	!\$omp parallel @rhs.f:28	
	9.8E-4	9.8E-4	3232/3232	!\$omp master @rhs.f:424	=
	0.225	0.228	3232/3232	!\$omp_do_@rhs.f:62	
	0.002	0.002	3232/3232	!\$omp master @rhs.f:74	
	0.002	0.002	3232/3232	!\$omp master @rhs.f:293	
	0.199	0.199	3232/3232	!\$omp_do_@rhs.f:384	
	0.002	0.002	3232/3232	!\$omp master @rhs.f:183	
	0.343	0.343	3232/3232	!\$omp_do_@rhs.f:37	
	0.016	0.016	3232/3232	!\$omp_do_@rhs.f:372	
	0.014	0.027	3232/3232	!\$omp_do_@rhs.f:413	
	0.609	0.609	3232/3232	!\$omp do @rhs.f:191	
	0.36	0.386	3232/3232	!\$omp_do_@rhs.f:301	
	0.583	0.583	3232/3232	!\$omp_do_@rhs.f:80	
	0.019	0.019	3232/3232	!\$omp do @rhs.f:400	
	0.006	0.006	3232/51680	!\$omp implicit barrier	
	0.069	0.069	3232/3232	!\$omp_do_@rhs.f:428	
	0.015	0.015	3232/3232	!\$omp_do_@rhs.f:359	
	0.001	0.000	6 400 (6 400	Idama accellel Grach she from	
	0.021	0.029	6432/6432	!\$omp parallel @exch_qbc.f:215	
>	0.021	0.029	6432	!\$omp parallel do @exch_qbc.f:215	
	0.007	0.007	6432/51680	!\$omp implicit barrier	
	0.02	0.033	6432/6432	!\$omp parallel @exch qbc.f:255	
>	0.02	0.033	6432	!\$omp_parallel_do_@exch_qbc.f:255	
	0.013	0.013	6432/51680	!\$omp implicit barrier	

Callsite Profiling and Tracing (TAU_CALLSITE=1)

Metric: TIME Value: Exclusive	TAU: ParaProf: n,c,t 0,0,0 - lu_callsite.ppk	
Units: seconds		
14.695	MPI_Recv()	
12.627	[CALLSITE] MPI_Recv() [@] [exchange_1_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_1.f] {68	-}]
	6.987	
	4.097 MPI_Send()	
	2.068 [CALLSITE] MPI_Recv() [@] [exchange_1_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_1.f} {86 2.053 [CALLSITE] MPI_Send() [@] [exchange_1_1]{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_1.f}	
	2.053 [CALLSITE] MPI_Send() [@] [exchange_1_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_1.f} {11 1.926 [CALLSITE] MPI_Send() [@] [exchange_1_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_1.f} {13	
	1.928 [CALLSTE] MPI_send() [@] [exchange_1_] [{/lus/theta=is0/projects/ioois/tau/workshop/NPB3.1/LU/exchange_1.] {13	JO}]
	0.365 [[CALLSITE] MPI_Finalize() [@] [MAIN_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/lu.f} {161}]	
	0.218 [CALLSITE] MPI Init() [@] [init_comm_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/init_comm.f} {31}]	
	0.188 MPI Wait()	
	0.104 [] [CALLSITE] MPI_Wait() [@] [exchange_3_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_3.f} {15	231
	0.08 [CALLSTE] MP] Wait() [@] [exchange_3] [[/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_3.1] [2	
	0.063 [CALSITE] MPI Send() [@] [exchange 3] [[/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange 3, [] 20	
	0.057 MPI Irecv()	
	0.055 [CALLSITE] MPI Send() [@] [exchange_3] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_3.f} {73	331
	0.03 CALLSITE] MPI Irecv() [@] [exchange 3] [//lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange 3.f] [10	
	0.026 CALLSITE] MPI_Irecv() [@] [exchange_3_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_3.f}	
	0.012 MPI Bcast()	
	0.009 [CALLSITE] MPI_Bcast() [@] [bcast_inputs_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/bcast_inputs.f}	{28}]
	0.008 MPI_Allreduce()	
	0.005 [CALLSITE] MPI_Allreduce() [@] [l2norm_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/l2norm.f} {55}]	
	0.002 [CALLSITE] MPI_Bcast() [@] [bcast_inputs_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/bcast_inputs.f}	{33}]
	0.001 [CALLSITE] MPI_Allreduce() [@] [ssor_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/ssor.f} {205}]	
	0.001 [CALLSITE] MPI_Wait() [@] [exchange_4_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_4.f] {59	
	0.001 [CALLSITE] MPI_Wait() [@] [exchange_6_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_6.f} {54	
	0.001 [CALLSITE] MPI_Wait() [@] [exchange_5_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_5.f} {54	
	7.6E-4 [CALLSITE] MPI_Wait() [@] [exchange_4_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_4.f} {10	
	4.9E-4 [CALLSITE] MPI_Irecv() [@] [exchange_4_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_4.f} {50	J}]
	4.9E-4 MPI_Barrier()	
	4.9E-4 [CALLSITE] MPI_Barrier() [@] [ssor_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/ssor.f} {70}]	
	4.2E-4 [CALLSITE] MPI_Allreduce() [@] [error_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/error.f} {53}]	
	4.0E-4 MPI_Comm_size()	
	2.6E-4 [CALLSITE] MPI_Comm_size() [@] [read_input_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/read_input.	t} {78
	2.4E-4 MPI_Comm_rank()	
	2.4E-4 [CALLSITE] MPI_Allreduce() [@] [pintgr_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/pintgr.f} {272}]	n (24
	2.4E-4 [CALLSITE] MPI_Comm_rank() [@] [init_comm_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/init_comm.	.1} {36

Identifying MPI Collective Sync Wait in Thread Callpath Relations

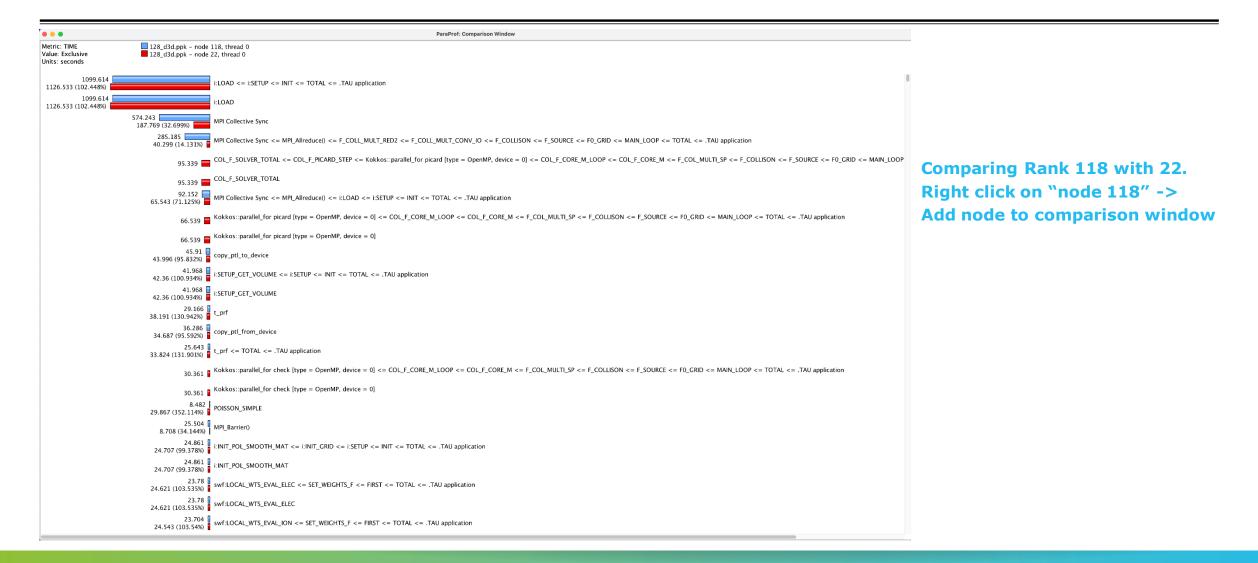
TAU: ParaProf: Call Path Data n,c,t, 118,0,0 - 128_d3d.ppk Metric Name: TIME Sorted By: Exclusive Units: seconds Inclusive Calls/Tot.Calls Name[id] Exclusive 1191.772 1/1 i:SETUP 1099.614 1099.614 1191.772 1 i:LOAD --> 0.006 92.158 3/9543 MPI Allreduce() 9.8E-4 9.8E-4 11/15177 MPI_Gatherv() 1.448 43/15177 MPI Gather() 1.448 MPI_Alltoall() 15.353 15.353 46/15177 89.821 89.821 MPI_Bcast() 4311/15177 6.777 6.777 195/15177 MPI Allgather() MPI Reduce() 68.678 68.678 991/15177 9.179 9.179 12/15177 MPI Comm dup() MPI_Allgatherv() 0.125 0.125 25/15177 382.861 382.861 MPI Allreduce() 9543/15177 574.243 574.243 15177 MPI Collective Sync --> 2.507 2.508 10/186 DISTRIBUTE_F0G 2.433 2.434 10/186 F_UPD_F0_SP 5.156 5.158 20/186 F0_CHARGE_SEARCH_INDEX 5.505 5.507 22/186 PULLBACK WEIGHT 24.86 UPDATE PTL WEIGHT 24.872 102/186 0.473 0.473 2/186 MAIN LOOP 4.975 4.977 20/186 DIAG fØ PORT1 PTL 45.91 45.93 186 copy ptl to device --> Kokkos::parallel_for set_buffer_particles_d [type = Cuda, device = 0] 0.02 0.02 186/272

MPI Collective Sync is the time spent in a barrier operation inside a collective

SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)

VICTOR VICT

Thread Comparison Window

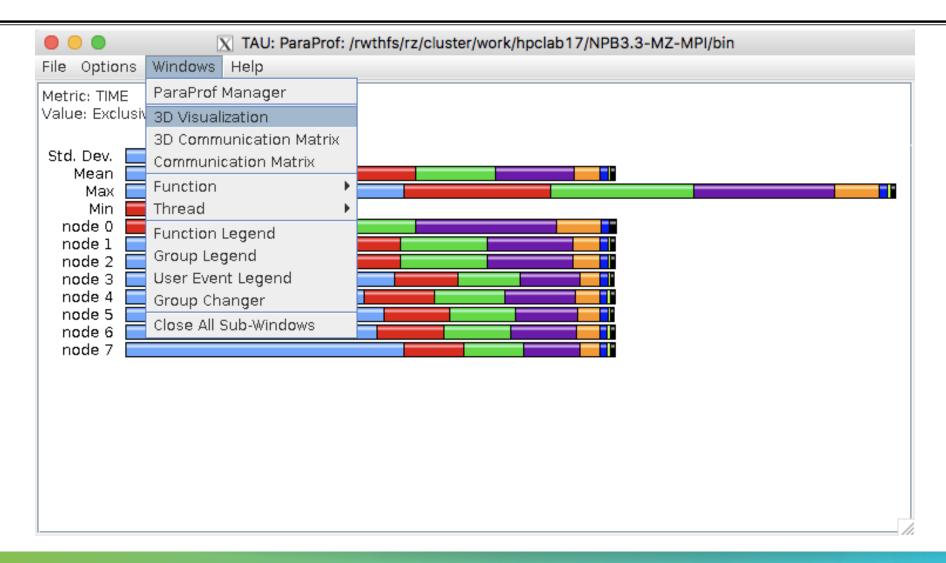


V VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

TAU – Context Events

	TAU: ParaProf: Context Events f	or thread: n,c,t, 1,0,0 -	samarc_obe_4p	_iomem_cp.ppk				
Name 🗸			Total	MeanValue N	umSamples M	inValue	MaxValue	Std. Dev.
 .TAU application 								
▶ read()			<i>(</i>	•• •				
▶ fopen64()	Write ba	ndwidth	ner t					
▶ fclose()		i a vi a ci i						
OurMain()			-					
malloc size			25,235	1,097.174	23	11	12,032	2,851.143
free size			22,707	1,746.692	13	11	12,032	3,660.642
OurMain [{wrapper.py}{3}]								
▶ read()								
malloc size			3,877	323.083	12	32	981	252.72
free size								122
▶ fopen64()		Dyte				f :1		
▶ fclose()		Byle	s wrii	llen ll	o each	ше		
<pre><module> [{obe.py}{8}]</module></pre>		-,	• • • • • •					
writeRestartData [{samarcInterface.py}{145}]								
samarcWriteRestartData								
vrite()	↓							
WRITE Bandwidth (MB/s) <file="sam< td=""><td>narc/restore.00002/nodes.000</td><td>004/proc.00001"></td><td></td><td>74.565</td><td>117</td><td>0</td><td>2,156.889</td><td>246.386</td></file="sam<>	narc/restore.00002/nodes.000	004/proc.00001">		74.565	117	0	2,156.889	246.386
WRITE Bandwidth (MB/s) <file="sam< td=""><td>narc/restore.00001/nodes.000</td><td>004/proc.00001"></td><td>↓</td><td>77.594</td><td>117</td><td>0</td><td>1,941.2</td><td>228.366</td></file="sam<>	narc/restore.00001/nodes.000	004/proc.00001">	↓	77.594	117	0	1,941.2	228.366
WRITE Bandwidth (MB/s)				76.08	234	0	2,156.889	237.551
Bytes Written <file="samarc restore<="" td=""><td>e.00002/nodes.00004/proc.00</td><td>0001"></td><td>2,097,552</td><td>17,927.795</td><td>117</td><td>1</td><td>1,048,576</td><td>133,362.946</td></file="samarc>	e.00002/nodes.00004/proc.00	0001">	2,097,552	17,927.795	117	1	1,048,576	133,362.946
Bytes Written <file="samarc restore<="" td=""><td>e.00001/nodes.00004/proc.00</td><td>0001"></td><td>2,097,552</td><td>17,927.795</td><td>117</td><td>1</td><td>1,048,576</td><td>133,362.946</td></file="samarc>	e.00001/nodes.00004/proc.00	0001">	2,097,552	17,927.795	117	1	1,048,576	133,362.946
Bytes Written			4,195,104	17,927.795	234	1	1,048,576	133,362.946
▶ open64()								

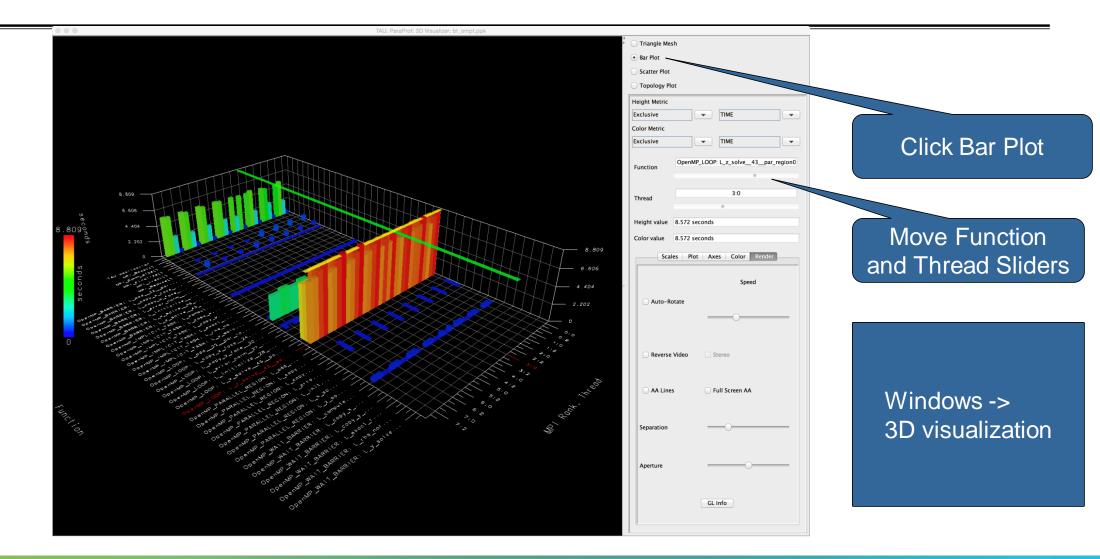
ParaProf with Optimized Instrumentation



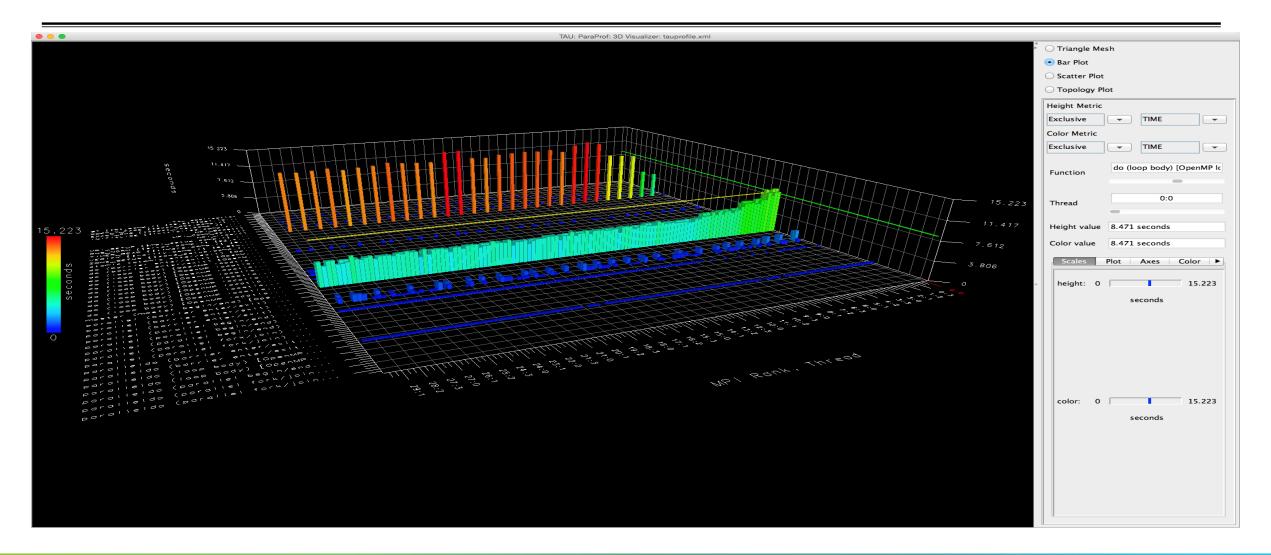
Create a Selective Instrumentation File, Re-instrument, Re-run

C TAU: ParaProf: /rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/bin	TAU: ParaProf: Selective Instrumentation File Generator	
File Options Windows Help Export Profile	Output File: /rwthfs/rz/cluster/work/hpclab17/NPB3.3-MZ-MPI/bin/select.tau	
Convert to Phase Profile Create Selective Instrumentation File Add Mean to Comparison Window	Exclude Throttled Routines	
Save Preferences	✓ Exclude Lightweight Routines	
Print Close This Window	Lightweight Routine Exclusion Rules	
Exit ParaProf!	Microseconds per call: 10	
node 4 Image: Contract of the contract	Number of calls: 100000	
node 7	Excluded Routines	
	lhsinit_ exact_solution_ matvec_sub_	
	matmul_sub	
	save Merge close	

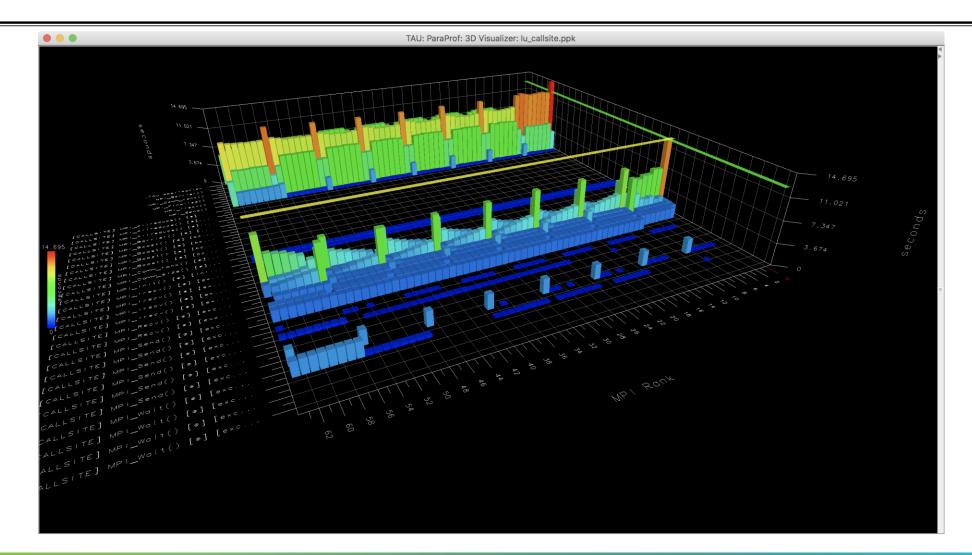
Paraprof 3D visualization window



ParaProf: 3D Visualization Window Showing Entire Profile



Callsite Profiling and Tracing

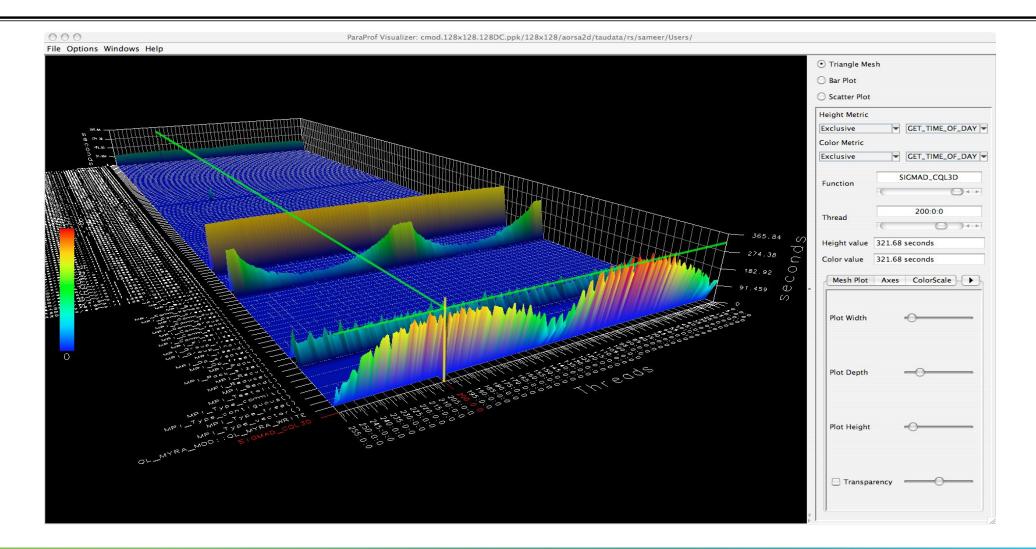


Callsite Profiling and Tracing

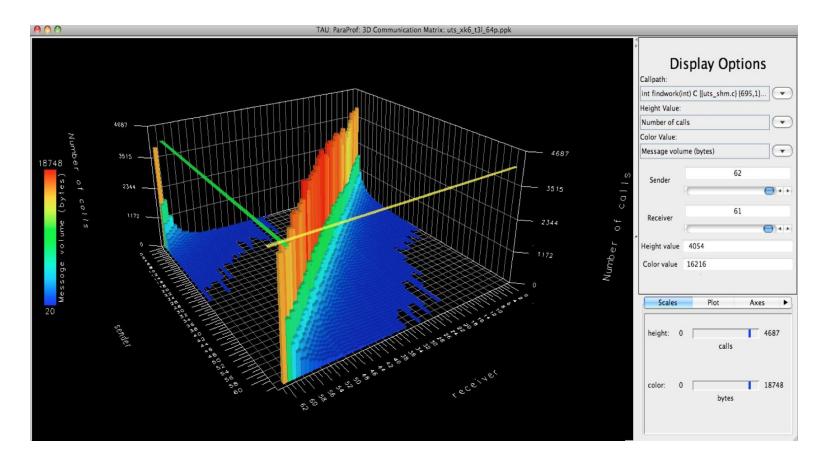
•••	TAU	J: ParaProf: 3D Visualizer: lu_callsite.ppk	
r Þ	Triangle Me	sh	
	🗿 Bar Plot		
	O Scatter Plot		
	O Topology P	lot	
	Height Metric		
	Exclusive	ТІМЕ	•
1 (D) 90	Color Metric		
9 40 40 40 40 40 40 40 40 40 40 40 40 40	Exclusive	ТІМЕ	•
		[CALLSITE] MPI_Recv() [@] [exchange_1_] [{/lus/theta-fs0/projects/Tools/tau/workshop/NPB3.1/LU/exchange_1.f} {6	68}]
	Function		
	Thread	0	-11
en e			
and and the second s Second second	Height value	12.627 seconds	- 1
-40 ⁹⁷	Color value	12.627 seconds	
		Scales Plot Axes Color Render	
	height:	0 14.695 seconds	
		seconds	
	color:	0 14.695	
		seconds	

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Parallel Profile Visualization: ParaProf



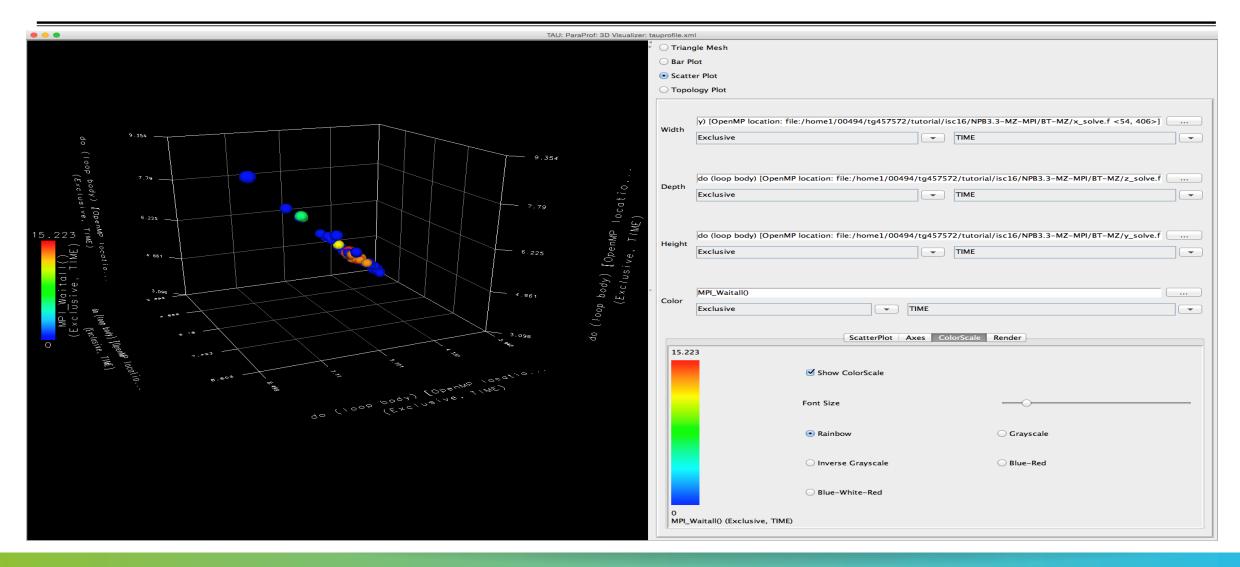
ParaProf 3D Communication Matrix



% export TAU_COMM_MATRIX=1

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

ParaProf: 3D Scatter Plot



ParaProf: Score-P Profile Files, Database

Applications	TrialField	Value
Standard Applications	Name	profile.cubex
- Call Default App	Application ID	0
e C Default Exp	Experiment ID	0
	Trial ID	0
	File Type Index	9
Minimum Inclusive Time	File Type Name	Cube
- • PAPI TOT CYC		
- PAPI FP INS		
– 9 ru_utime		
— 🥯 ru_stime		
– • ru_maxrss		
— 🤍 ru_ixrss		
- • ru_idrss		
— ❷ ru_isrss — ❷ ru minflt		
- 9 ru maifit		
- 9 ru nswap		
- • ru_inblock		
– 🔍 ru oublock		
— ● ru msgsnd		
– • ru_msgrcv		
– 🥯 ru_nsignals		
- • ru_nvcsw		
— • ru_nivesw		
- • bytes_sent		
bytes_received		
Default (jdbc:h2:/home/livetau/.ParaProf//perfdmf;AUTO_SERVER=TRUE) Default (jdbc:h2:/home/livetau/.ParaProf/perfexplorer_working_(jdbc:h2:/home/livetau/.ParaProf/perfexplorer_working Add Application		
a benevered in a manual data and a manual benevered and in a more that the more the more that the more that the more the more that the more that the more the m		
Add Experiment		
Add Trial		

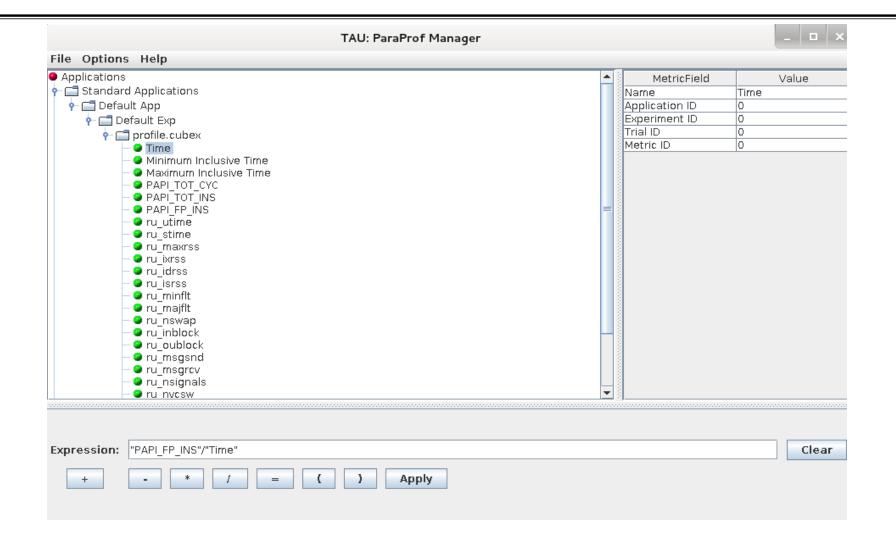
ParaProf: File Preferences Window

	ParaProf Preferences	_ 🗆 ×
File		
Font SansSerif Bold Size	n,c,t 0,0,0 n,c,t 0,0,1 n,c,t 0,0,2	
Italic 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1		
Window defaults	Settings	
Units Seconds 💌	 Reverse Call Paths Interpret threads that do not call a given function as a 0 value for statistics computation 	
Show Values as Percent	Generate data for reverse calltree (requires lots of memory) (does not apply to currently loaded profiles) Show Source Locations Auto label node/context/threads	
Restore Defaults	Apply	Cancel

ParaProf: Group Changer Window

	TAU: Pa	raProf: Group Changer: prot	file.cubex	
	Region	Current		Available
filter:				new group
!\$omp atomic !\$omp atomic	_	CUBE_DEFAULT	CUBE_	CALLPATH
!\$omp do @er	ror.f:33			
!\$omp do @er	ror.f:91			
!\$omp do @ex	act_rhs.f:147		<	
!\$omp do @ex	act_rhs.f:247			
!\$omp do @ex	act_rhs.f:31			
!\$omp do @ex	act_rhs.f:346			
!\$omp do @ex	(act_rhs.f:46			
!\$omp do @ini	itialize.f:100			
!\$omp do @ini	itialize.f:119			
!\$omp do @ini	itialize.f:137			
!\$omp do @ini	itialize.f:156			
!\$omp do @ini	itialize.f:174		>	
!\$omp do @ini	itialize.f:192			
!\$omp do @ini	itialize.f:31	T		
•				

ParaProf: Derived Metric Panel in Manager Window



VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Sorting Derived FLOPS metric by Exclusive Time

TAU: Par	aProf: node 0, thread 0 – profile.cubex	_ 0 ×
ile Options Windows Help		
Metric: (PAPI_FP_INS / Time) Value: Exclusive Units: Derived metric shown in seconds format Sorted By: Exclusive (Time)		
3.0217E9 3.0217E9	MAIN => adi_ => y_solve_ => !\$omp parallel @y_solve.f:43 => !\$omp do @y_solve.f:52 !\$omp do @y solve.f:52	2
3.2421E9 3.2421E9	MAIN_ => adi_ => z_solve_ => !\$omp parallel @z_solve.f:43 => !\$omp do @z_solve.f:52 !\$omp do @z_solve.f:52	2
3.0673E9 3.0673E9	MAIN_ => adi_ => x_solve_ => !\$omp parallel @x_solve.f:46 => !\$omp do @x_solve.f:54 !\$omp do @x solve.f:54	4
3.3299E9 3.3298E9 3.5138E9	\$omp do @rhs.f:191 MAIN_ => adi_ => compute_rhs_ => !\$omp parallel @rhs.f:28 => !\$omp do @rhs.f:191 \$omp do @rhs.f:80	
3.514E9 1965740.083	MAIN_ => adi_ => compute_rhs_ => !\$omp parallel @rhs.f:28 => !\$omp do @rhs.f:80	
2518815.107	!\$omp parallel @rhs.f:28	
2518981.066 3.502E8	MAIN => adi_ => compute_rhs_ => !\$omp parallel @rhs.f:28 \$omp do @rhs.f:37	
3.4975E8	MAIN_ => adi_ => compute_rhs_ => !\$omp parallel @rhs.f:28 => !\$omp do @rhs.f:37 !\$omp do @rhs.f:301	
4.0205E9	MAIN_ => adi_ => compute_rhs_ => !\$omp parallel @rhs.f:28 => !\$omp do @rhs.f:301	
393146.074	!\$omp do @rhs.f:62	
393024.443 60.754	<pre> MAIN_ => adi_ => compute_rhs_ => !\$omp parallel @rhs.f:28 => !\$omp do @rhs.f:62 MAIN => mpi setup => MPI Init thread</pre>	
60.754	MANU_ => mpl_setup_ => MPl_mt_thread	
2218222.902	MAIN_ => exch_qbc_ => copy_x_face_	
2218222.902	copy x face	
2217983.431	MAIN_=> exch_qbc_ => copy_y_face_	
2217983.431	copy_face	
2691052.918	MAIN_=> exch_qbc_	
2691052.918	exch_qbc	
1.5944E9	!\$omp do @rhs.f:384	
1.5944E9] MAIN_ => adi_ => compute_rhs_ => !\$omp parallel @rhs.f:28 => !\$omp do @rhs.f:384	t I
65007.137	MAIN_ => exch_qbc_ => MPI_Waitall	
		•

ParaProf Hands-On

SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)

Login to Xpra Desktop via https://jupyter-jsc.fz-juelich.de

```
$ source /p/project/training2341/setup.sh
$ tar xf $PROJECT/examples/tea leaf.tar.gz
 cd TeaLeaf CUDA
Ş
$ make
$ cd bin
$ cp ../jobscripts/juwelsbooster/tau.sbatch .
$ cat tau.sbatch
....
srun tau exec -T cupti, mpi -cupti -ebs ./tea leaf
$ sbatch tau.sbatch
$ pprof -a | more
$ paraprof &
```

ParaProf: TeaLeaf_CUDA

Control Contro Control Control Control Control Control Control Control Control Co	nde1/juwels/workshop/SC22/TeaLeaf_CUDA/bir	ı		_ 🗆 ×	TAU: ParaProf: /p/home/jusers/shende1/juwels/workshop/SC2: File Options Windows Help	2/TeaLeaf_CUE T TAU: ParaProf: Function Data File Options Windows Help	الWindow: /p/home/jusers/shende1/jt المالية المالية المالية المالية المالية المالية المالية المالية المالية ال
Name	Exclusive TAUGPU TIME Inclusive		Calls	Child Calls	Metric: TAUGPU_TIME Value: Exclusive	Name:	
TAU application	0.003	47.092	Calls		value: Exclusive	device_tea_leaf_ppcg_solve_ double const*, double*, double	
taupreload_main	6.154	47.089	1	676,842	Std. Dev.		ble const*, double const*, doul
cudaMemcpy	29.198	29.198	214,620	010,042	Mean Mean	const*, double const*, doubl	
- MPI_Waitall()	8.663	8.663	104,774	0		Metric Name: TAUGPU_TIM	E
MPI Init()	0.177	1.075	1	12	Min	Value: Exclusive	
cudaStreamCreateWithFlags	0.898	0.898	1	0	node 0, thread 0	Units: seconds	
cudaLaunchKernel [THROTTLED]	0.705	0.705	100,001	0	node 0, thread 1		
MPI_Allreduce()	0.131	0.551	4,752	23,760	node 1, thread 0	6.252	std. dev.
MPI Collective Sync	0.384	0.384	4,764	0	node 1, thread 1	6.252	mean
MPI_Testall()	0.238	0.238	52,387	0	node 2, thread 0	12.526	max
MPI_Finalize()	0.142	0.142	1	3	node 2, thread 1	12.466	min
MPI_Isend() [THROTTLED]	0.119	0.119	100,001	0	node 3, thread 0	12.484	node 0, threa
MPI_Irecv() [THROTTLED]	0.098	0.098	100,001	0	node 3, thread 1	12.518	node 1, thread node 2, thread
MPI_Cart_create()	0.073	0.073	100,001	0	node 4, thread 1	12.466	node 2, threa node 3, threa
MPI_Barrier()	0.039	0.039	7	0	node 5, thread 0	12.522	node 3, threa
cudaPointerGetAttributes	0.037	0.037	19,056		node 5, thread 1	12.499	node 5, threa
cudaMalloc	0.029	0.029	48	0	node 6, thread 0	12.526	node 6, threa
cudaGetDeviceProperties	0.003	0.003	1	0	node 6, thread 1	12.523	node 7, threa
cudaDeviceSynchronize	0.002	0.002	132	0	node 7, thread 0		
MPI_Reduce()	0	0.001	12	60	node 7. thread 1		
cudaFree	0.001	0.001	8	0	TAU: ParaProf: 3D Visualizer: /p/home/jusers/shende1/juwels	/workshop/SC22/TeaLeaf_CUDA/bin	_ 0
cudaMemset	0	0	38	0	File Options Windows Help		
cudaGetLastError	0	0	46	0		Triangl	Mosh
cudaStreamDestroy	0	0	1	0			
cudaSetDevice	0	0	4	0		Bar Plo	t
cudaGetDeviceCount	0	0	9	0			Plot
MPI_Info_delete()	0	0	1	0		- Senter	FIG
MPI_Cart_shift()	0	0	2	0		 Topolo 	gy Plot
MPI_Cart_coords()	0	0	1	0			*-1-
MPI_Dims_create()	0	0	2	0		Height Me	
MPI_Comm_size()	0	0	2	0		Exclusive	TAUGPU_TIME
MPI_Comm_rank()	0	0	2	0		Color Met	ric
cudaGetDevice	0	0	1	0		Exclusive	
						Function	device_unpack_top_buffer(kern
TAU: ParaProf: node 0, thread 1 - /p/home/jusers/shende1/juwels/wo	orkshop/SC22/TeaLeaf_CUDA/bin			_ 0 X		Function	()
le Options Windows Help							0:1
etric: TAUGPU_TIME						Thread	
lue: Exclusive							
its: seconds						Height v	alue 0.202 seconds
						Color va	ue 0.202 seconds
0.154 .TAU application				^	and the set of the set		

Login to Xpra Desktop via https://jupyter-jsc.fz-juelich.de

\$ source /p/project/training2341/setup.sh

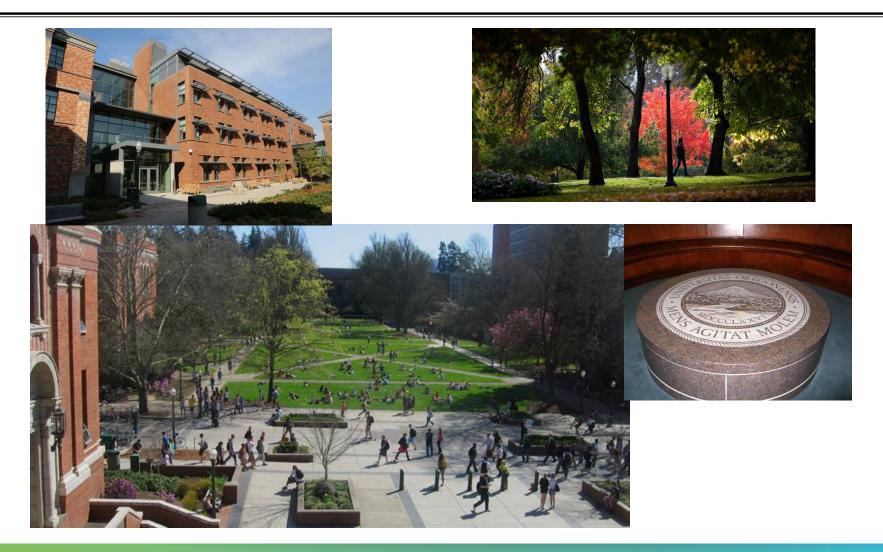
\$ wget <u>http://tau.uoregon.edu/demo.ppk</u>

\$ paraprof demo.ppk &

Windows -> 3D Visualization -> Bar Chart -> Function and Thread slider

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Performance Research Lab, University of Oregon, Eugene, USA



Support Acknowledgments





Acknowledgement

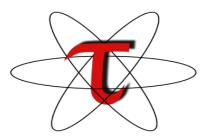
This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of two U.S. Department of Energy organizations (Office of Science and the National Nuclear Security Administration) responsible for the planning and preparation of a capable exascale ecosystem, including software, applications, hardware, advanced system engineering, and early testbed platforms, in support of the nation's exascale computing imperative.





VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Download TAU from U. Oregon



http://tau.uoregon.edu

http://www.hpclinux.com [LiveDVD, OVA] <u>https://e4s.io</u> [Containers for Extreme-Scale Scientific Software Stack]

Free download, open source, BSD license



Hands-on resume: TeaLeaf MPI+CUDA



Recap: Setup for exercises

- Connect to your training account on JUWELS Booster (with X11-forwarding)
 - % ssh -X <yourid>@juwels-booster.fz-juelich.de

Use Jupyter-JSC/Xpra instead!

Set account and default environment (NVHPC + ParaStationMPI) via helper script

% source /p/project/training2341/setup.sh

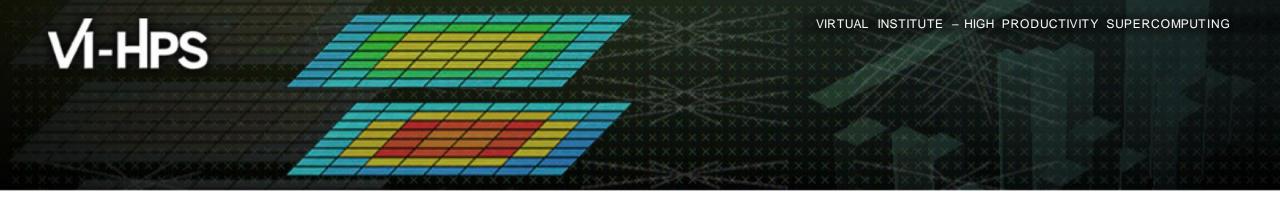
- Change to WORK directory containing TeaLeaf_CUDA sources & bin.scorep directory
 - Existing executable instrumented by Score-P can be reused
 - $\frac{9}{8}$ cd \$WORK
 - % cd TeaLeaf_CUDA
 - [⊗] cd bin.scorep
- Load Score-P module

% module load Score-P CubeGUI

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

TeaLeaf_CUDA trace measurement collection for Vampir...

```
Change to
% edit scorep.sbatch
                                                                       directory with the
% cat scorep.sbatch
                                                                       Score-P
# Score-P measurement configuration
                                                                       instrumented
export SCOREP CUDA ENABLE=default
                                                                       executable and
export SCOREP CUDA BUFFER=48M
export SCOREP EXPERIMENT DIRECTORY=scorep-tea leaf-8 vampir
                                                                       edit the job script
export SCOREP FILTERING FILE=../config/scorep.filter
export SCOREP ENABLE TRACING=true
                                                                   Uncomment lines to enable
export SCOREP TOTAL MEMORY=250M
                                                                   Score-P trace collection
# Run the application
srun ./tea leaf
                                                                      Submit the job
% sbatch scorep.sbatch
```



Interactive visualization and time-interval statistics with Vampir



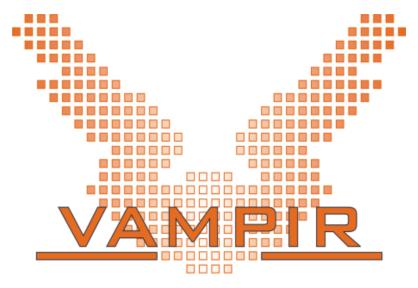
Outline

Part I: Welcome to the Vampir Tool Suite

- Mission
- Event Trace Visualization
- Vampir & VampirServer
- The Vampir Displays

Part II: Vampir Hands-On

Visualizing and analyzing NPB-MZ-MPI / BT



VIRTUAL×INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Event Trace Visualization with Vampir

- Alternative and supplement to automatic analysis
- Show dynamic run-time behavior graphically at any level of detail
- Provide statistics and performance metrics

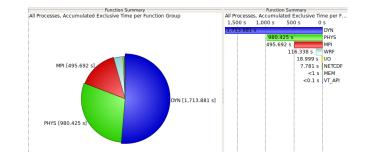
Timeline charts

Show application activities and communication along a time axis

Summary charts

Provide quantitative results for the currently selected time interval

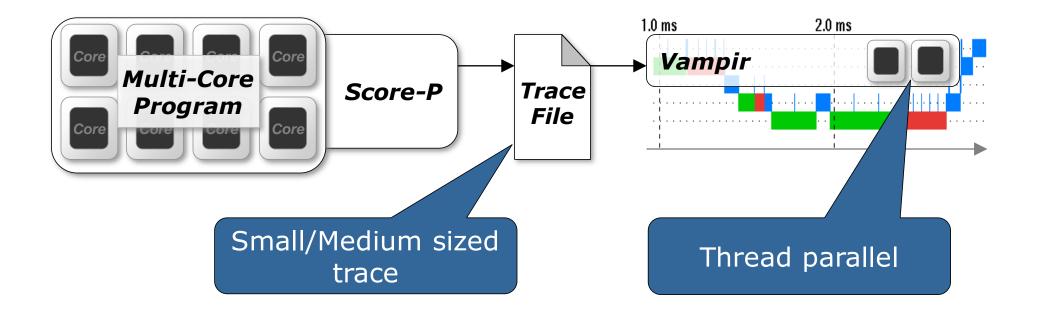
	84.8 s	84.9 s	85.0 s	85.1 s	85.2 s
	1	-		1	
Process 0	YSU				
Process 1	MPI_Wait				
Process 2	MPI	Wait			
Process 3			MPI_Wait		
Process 4	YSU CUMULAS_DRI	VER			
Process 5	MPT	Wait			
Process 6	MPI	Wait			
Process 7			-MPI_Wait		
Process 8		VER			v.
Process 9		Wait			
Process 10	MPT	Wait			
Process 11		/	-MPI_Wait		
Process 12	YSU CUMULUS CRIVER				
Process 13	MPI	Wait			
Process 14	MPI_Wait MPI	Wait			
Process 15			MPI_Wait		14



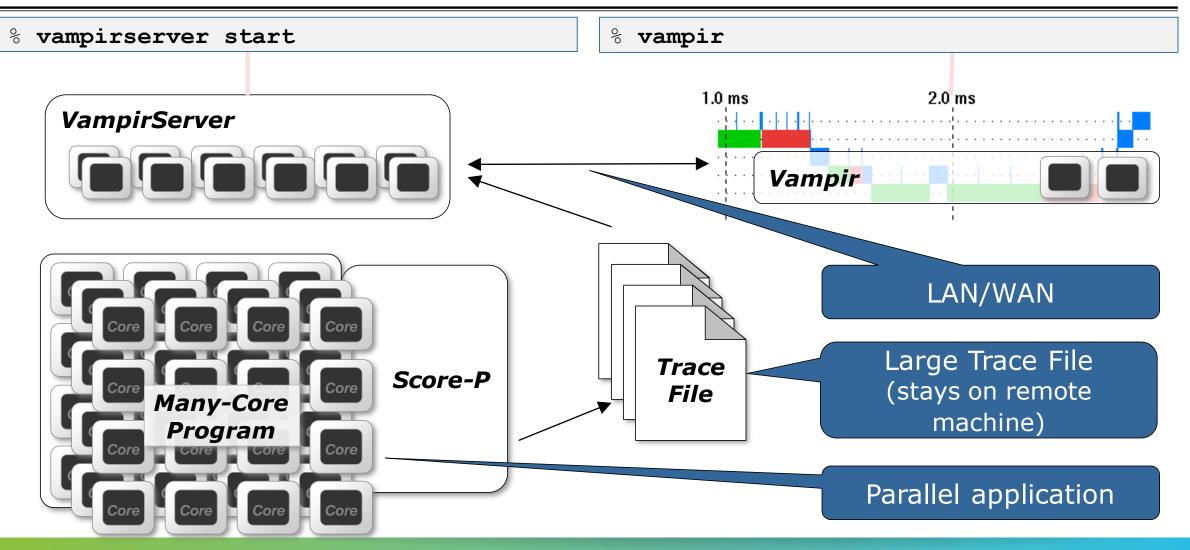
VICTOR VICT

Visualization Modes (1) Directly on front end or local machine

% vampir



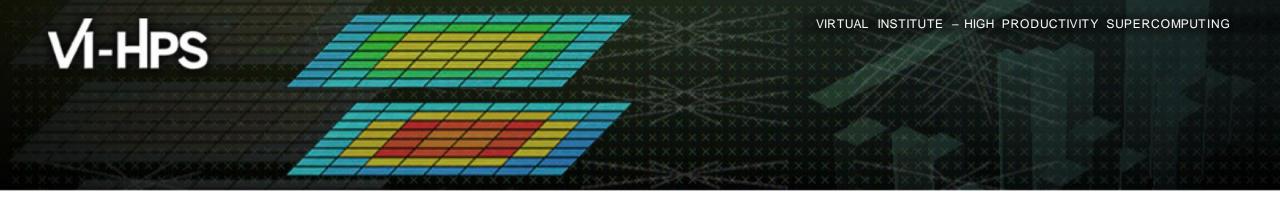
Visualization Modes (2) On local machine with remote VampirServer



VICTOR VICT

The main displays of Vampir

- Timeline Charts:
 - 🚎 Master Timeline
 - 📷 Process Timeline
 - Counter Data Timeline
 - 🛚 🐻 Performance Radar
- Summary Charts:
 - Sunction Summary
 - Message Summary
 - Process Summary
 - Communication Matrix View



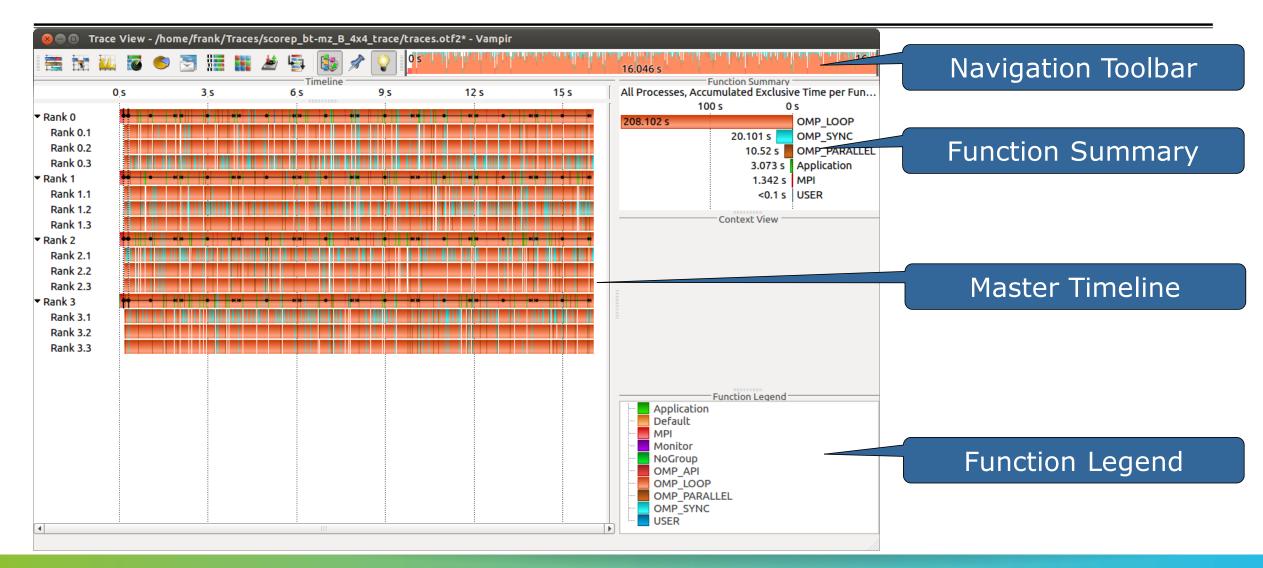
Hands-on: Visualizing and analyzing NPB-MZ-MPI / BT



Start Vampir inside Xpra Desktop

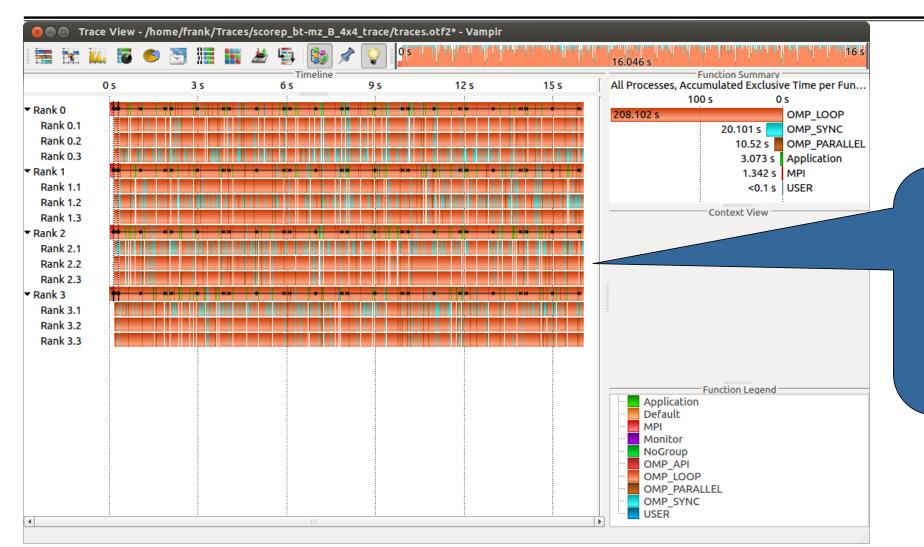
= 🖃 🖸 💷 🔌				
= 🔛 :: 📟 📢	>	Performance Tools >	袶 CubeGUI	
Server	>		Vampir	
Information	>	Visualization Tools >	vanpii	
Reload	Í			xterm -hold -e 'launch_vampir.sh'
Disconnect				
				wesarg1@jwlogin04:~ _ C X [wesarg1@jwlogin04 ~]\$
				[wesarg1@jwlogin04 ~]\$

Visualization of the NPB-MZ-MPI / BT trace



Visualization of the NPB-MZ-MPI / BT trace Master Timeline

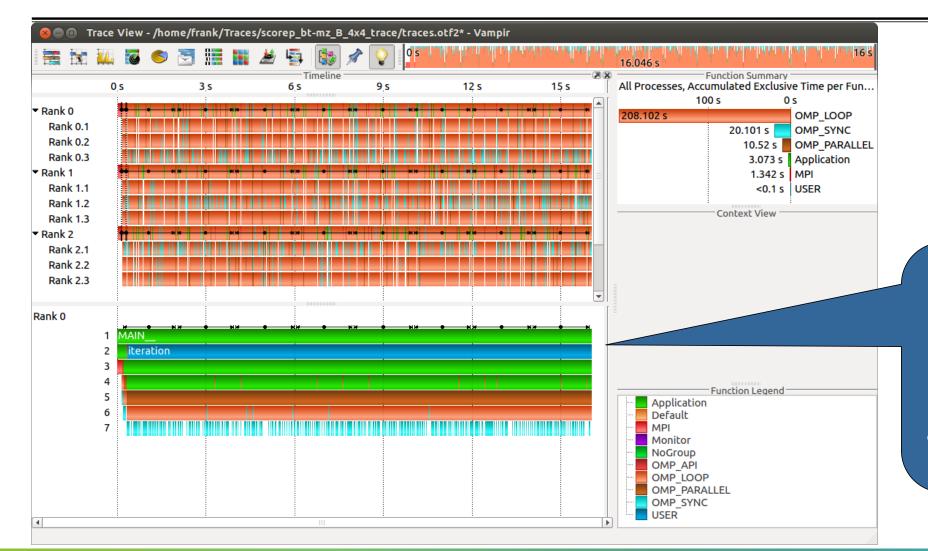




Detailed information about functions, communication and synchronization events for collection of processes.

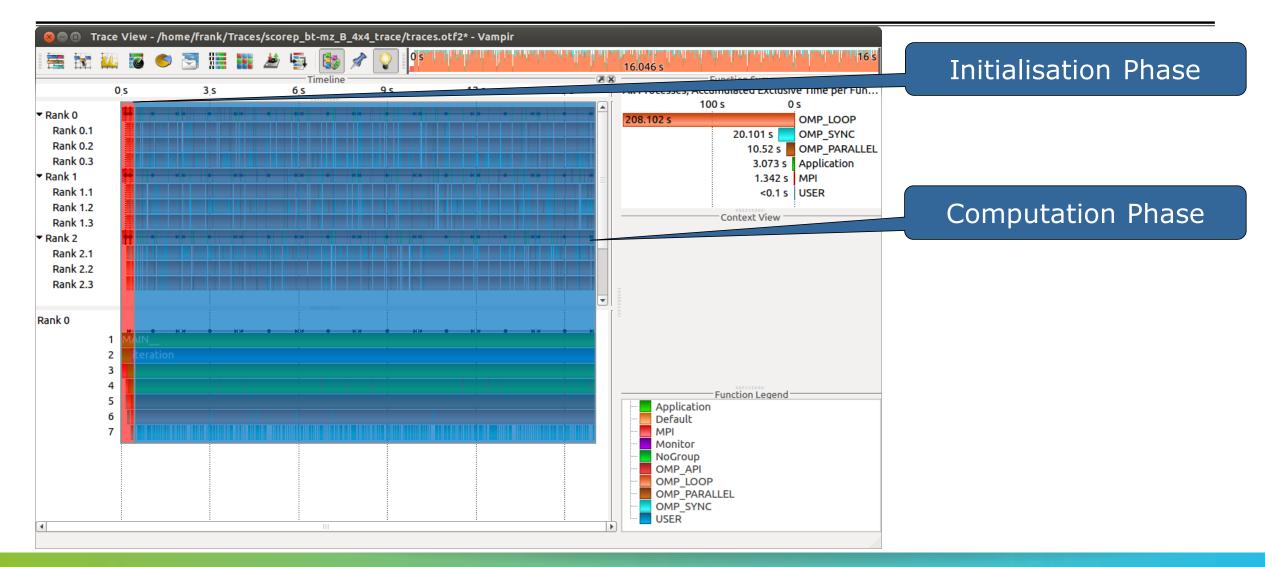
Visualization of the NPB-MZ-MPI / BT trace Process Timeline





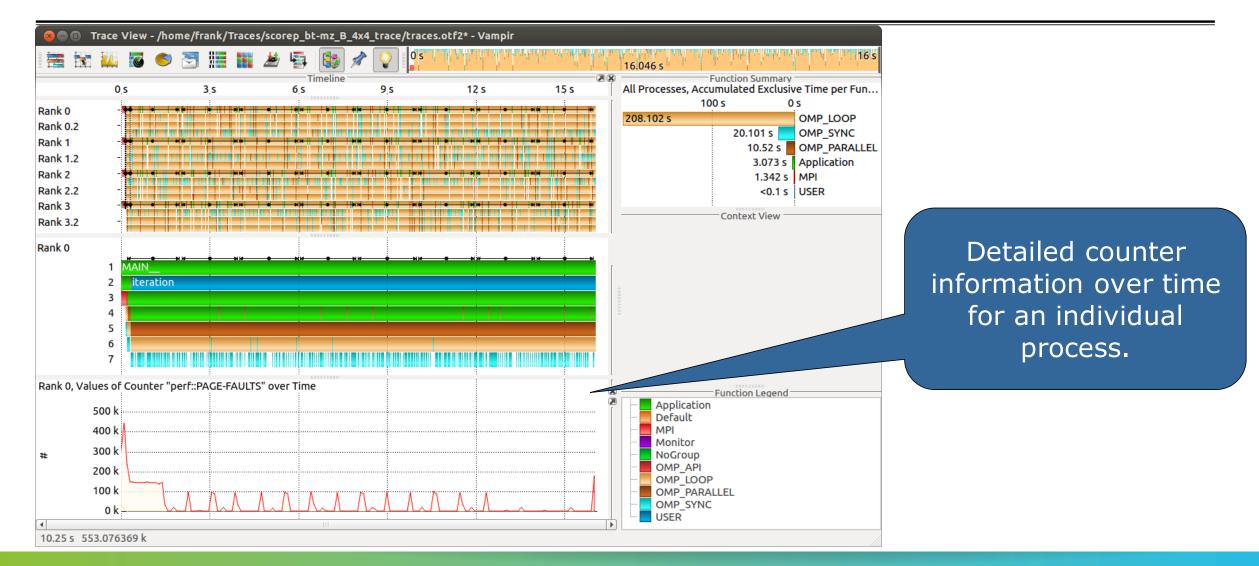
Detailed information about different levels of function calls in a stacked bar chart for an individual process.

Visualization of the NPB-MZ-MPI / BT trace Typical program phases



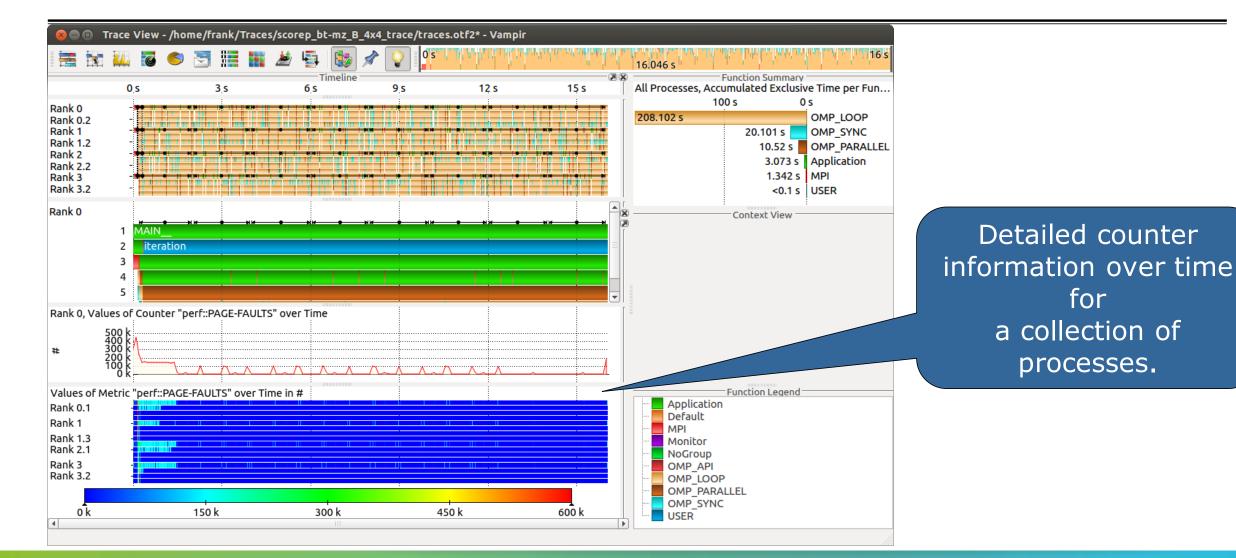
Visualization of the NPB-MZ-MPI / BT trace Counter Data Timeline



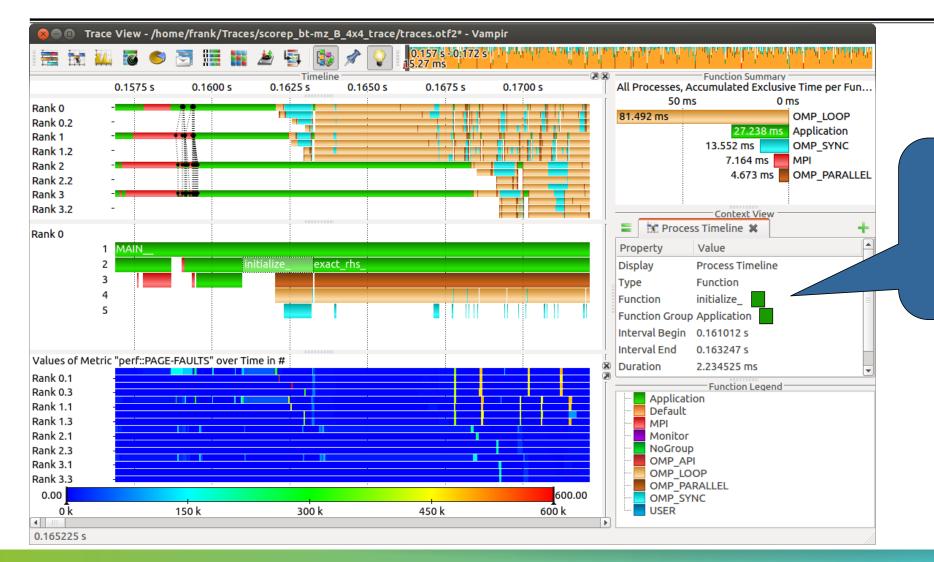


Visualization of the NPB-MZ-MPI / BT trace Performance Radar



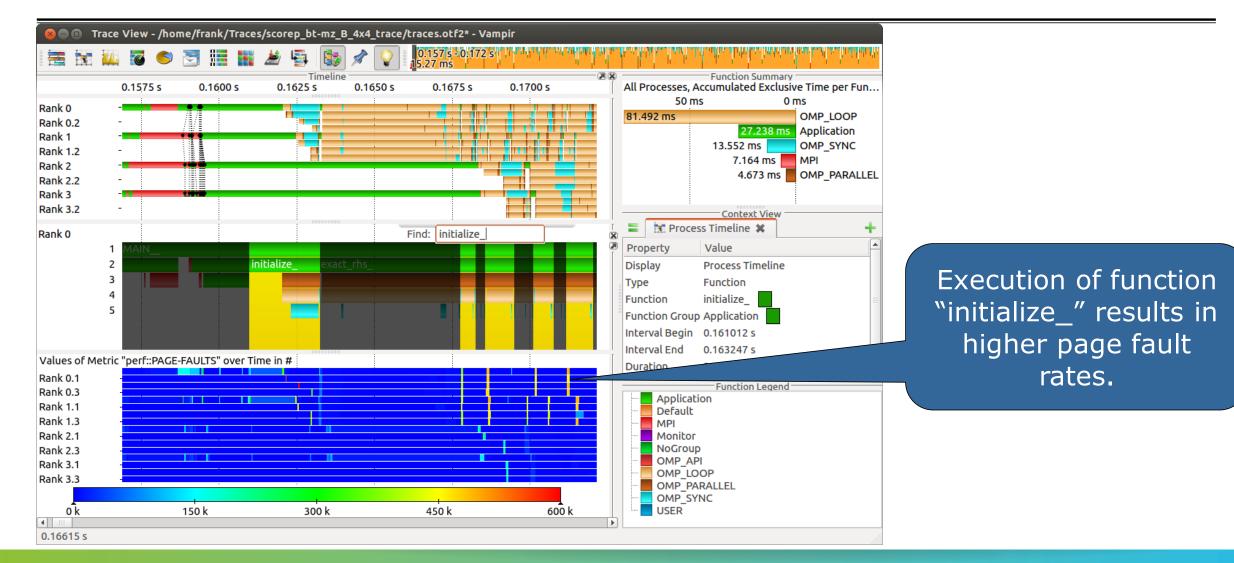


Visualization of the NPB-MZ-MPI / BT trace Zoom in: Inititialisation Phase

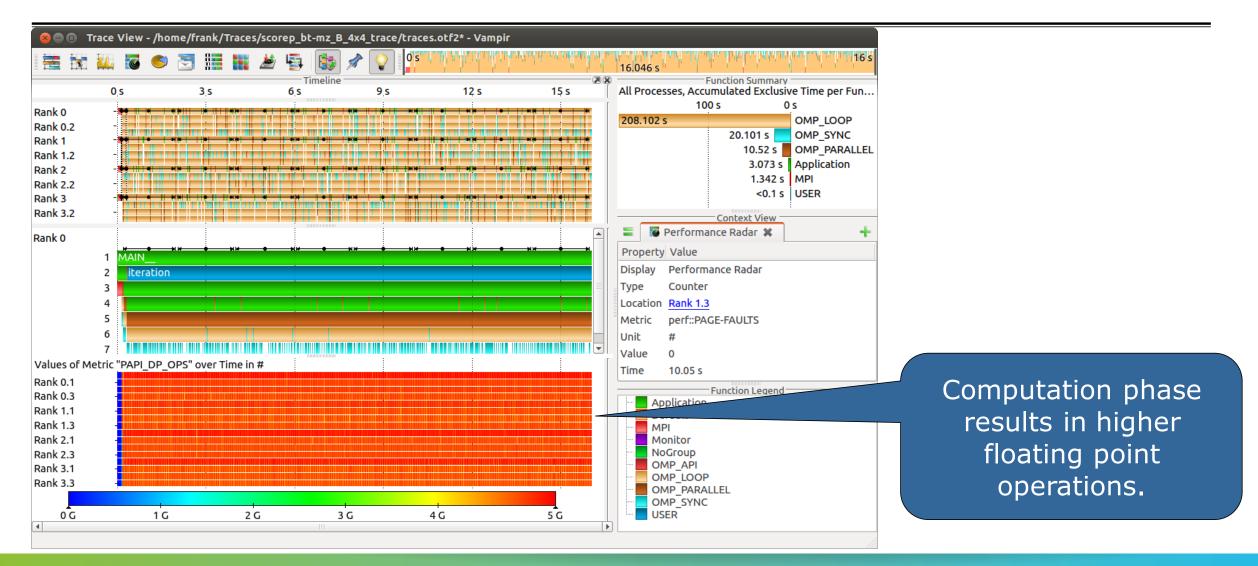


Context View: Detailed information about function "initialize_".

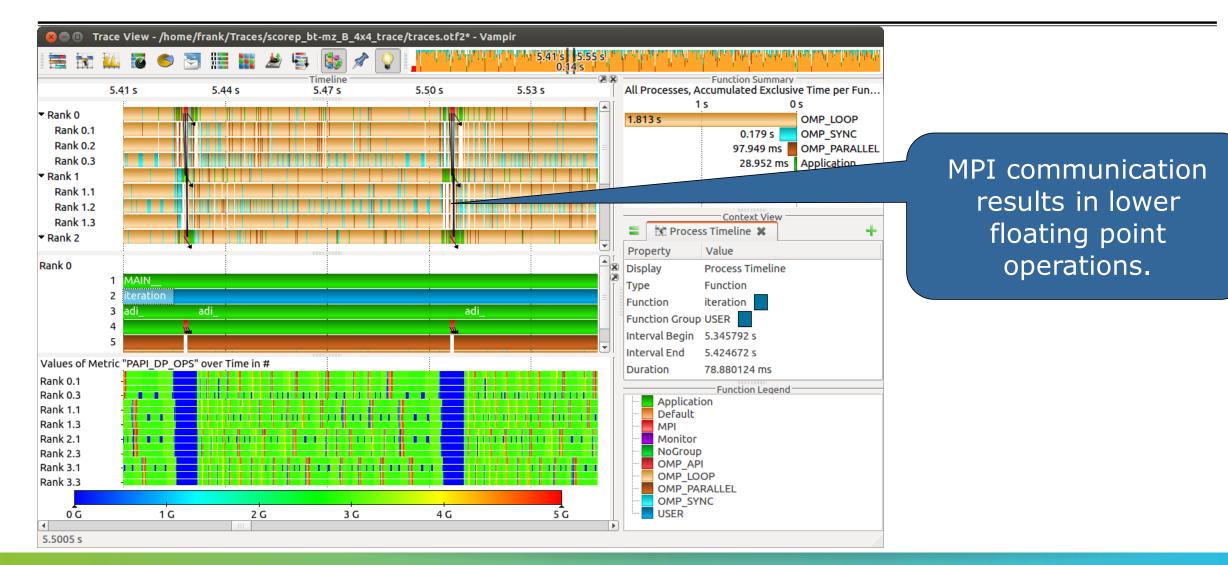
Visualization of the NPB-MZ-MPI / BT trace Find Function



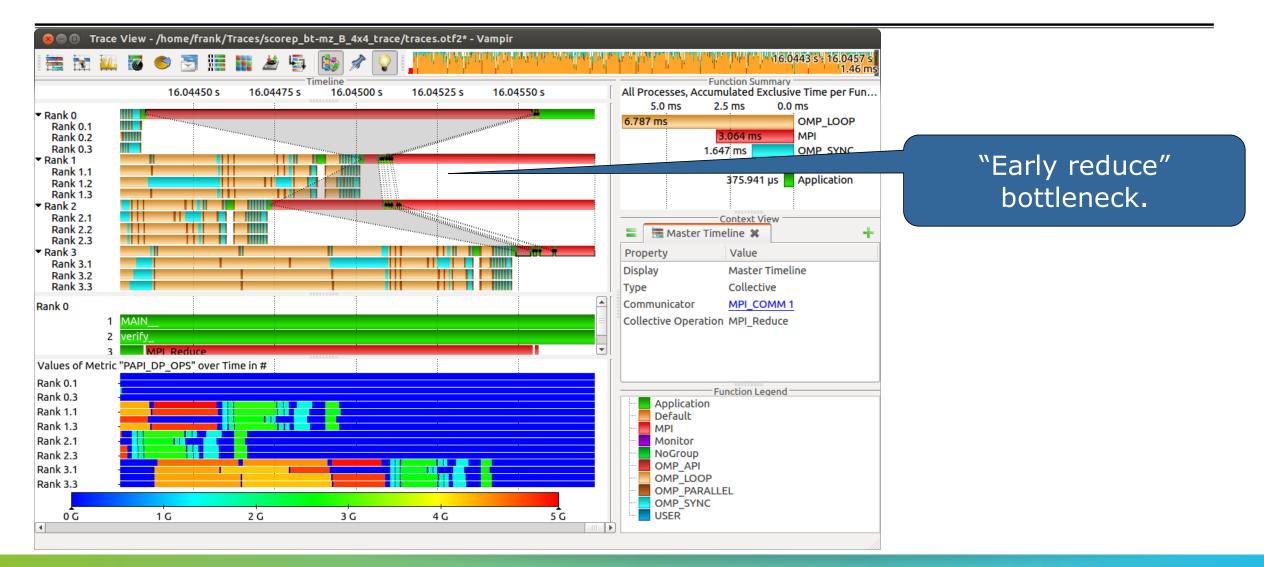
Visualization of the NPB-MZ-MPI / BT trace Computation Phase



Visualization of the NPB-MZ-MPI / BT trace Zoom in: Computation Phase

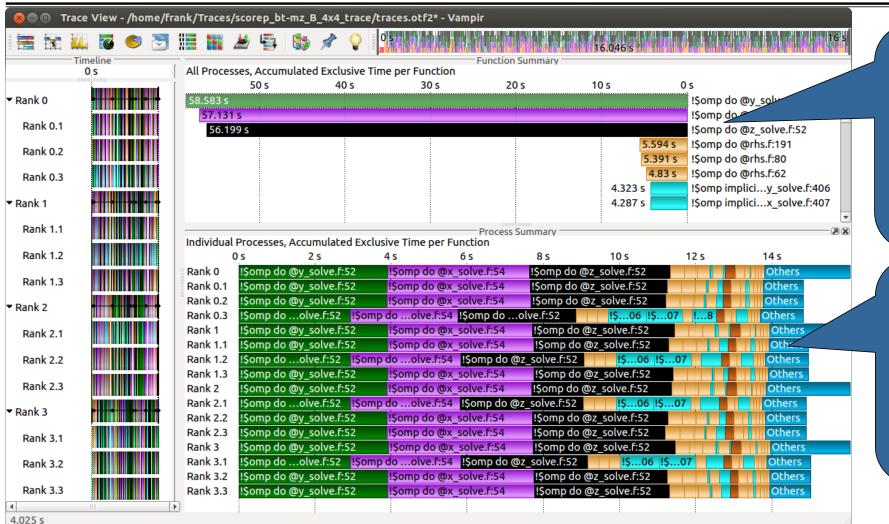


Visualization of the NPB-MZ-MPI / BT trace Zoom in: Finalisation Phase



Visualization of the NPB-MZ-MPI / BT trace Process Summary



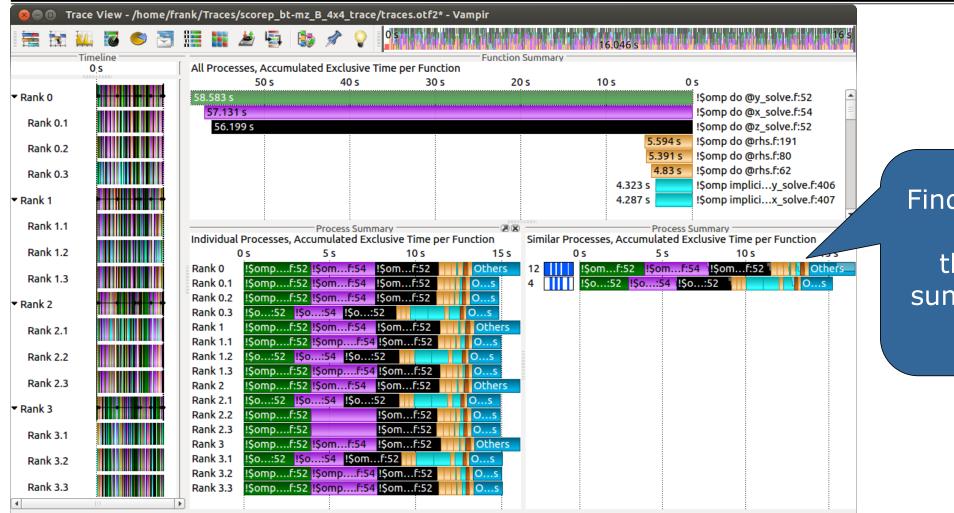


Function Summary: Overview of the accumulated information across all functions and for a collection of processes.

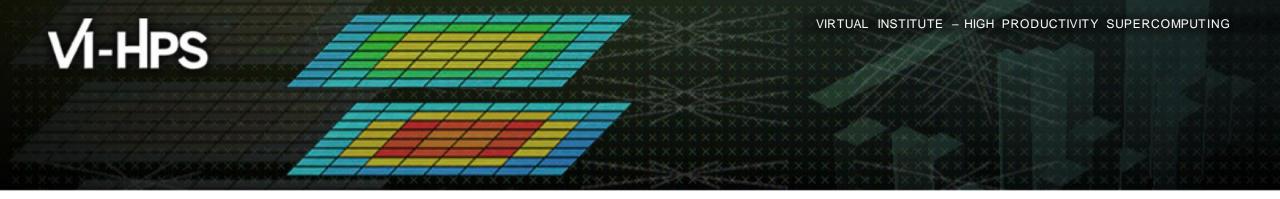
Process Summary: Overview of the accumulated information across all functions and for every process independently.

Visualization of the NPB-MZ-MPI / BT trace Process Summary





Find groups of similar processes and threads by using summarized function information.



Summary and Conclusion

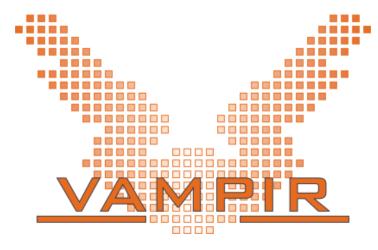


Summary

- Vampir & VampirServer
 - Interactive trace visualization and analysis
 - Intuitive browsing and zooming
 - Scalable to large trace data sizes (20 TiByte)
 - Scalable to high parallelism (200,000 processes)
- Vampir for Linux, Windows, and Mac OS X

VIRTUAL VI





Visit us at ZIH booth #1533



service@vampir.eu

SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)



Automatic trace analysis with the Scalasca Trace Tools

Markus Geimer Jülich Supercomputing Centre

trace tools scalasca



Scalasca Trace Tools

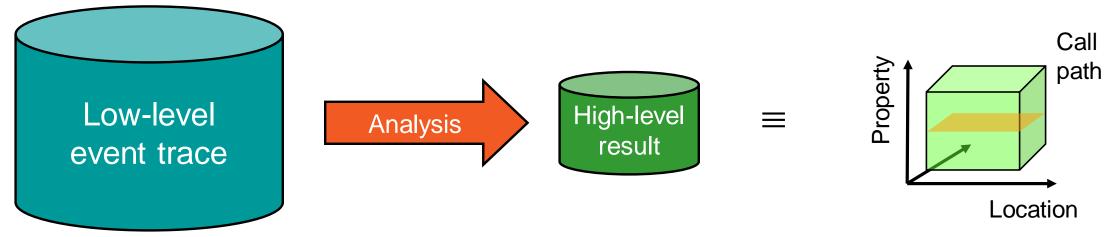
DOI 10.5281/zenodo.4103922

- Scalable trace-based performance analysis toolset for the most popular parallel programming paradigms
 - Current focus: MPI, OpenMP, and (to a limited extend) POSIX threads
 - Analysis of traces including only host-side events from applications using CUDA, OpenCL, or OpenACC (also in combination with MPI and/or OpenMP) is possible, but results need to be interpreted with some care
- Specifically targeting large-scale parallel applications
 - Demonstrated scalability up to 1.8 million parallel threads
 - Of course also works at small/medium scale
- Latest release:
 - Scalasca Trace Tools v2.6.1 (Dec 2022)

Automatic trace analysis

Idea

- Automatic search for patterns of inefficient behaviour
- Classification of behaviour & quantification of significance
- Identification of delays as root causes of inefficiencies



- Guaranteed to cover the entire event trace
- Quicker than manual/visual trace analysis
- Parallel replay analysis exploits available memory & processors to deliver scalability

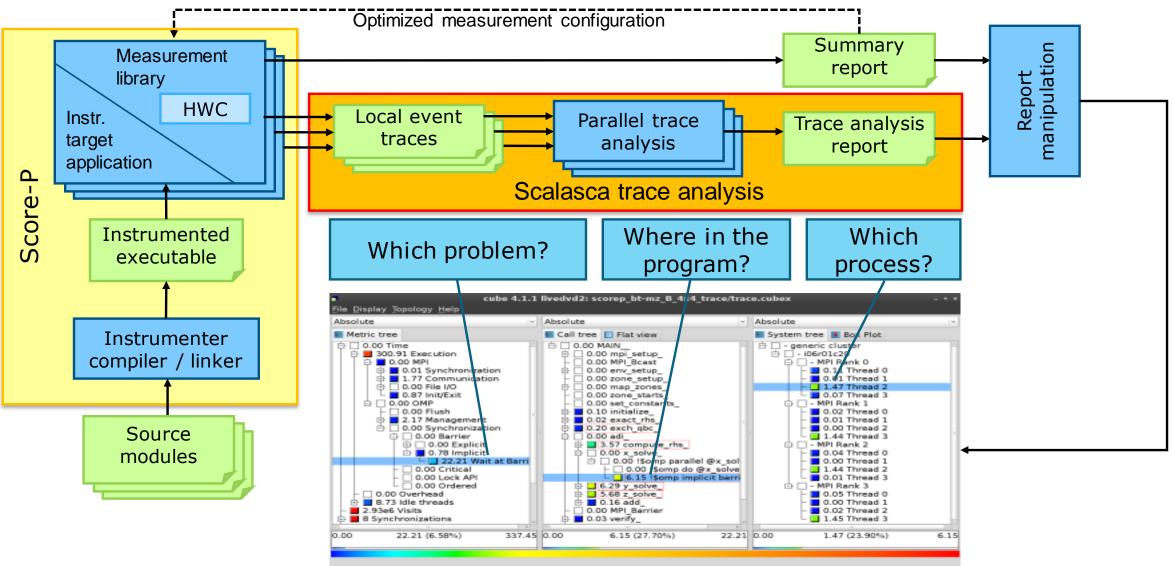
Scalasca Trace Tools: Features

- Open source, 3-clause BSD license
- Supports all major HPC platforms
- Uses Score-P instrumenter & measurement libraries
 - Scalasca v2 core package focuses on trace-based analyses
 - Provides convenience commands for measurement, analysis, and post-processing
 - Supports common data formats
 - Reads event traces in OTF2 format
 - Writes analysis reports in CUBE4 format

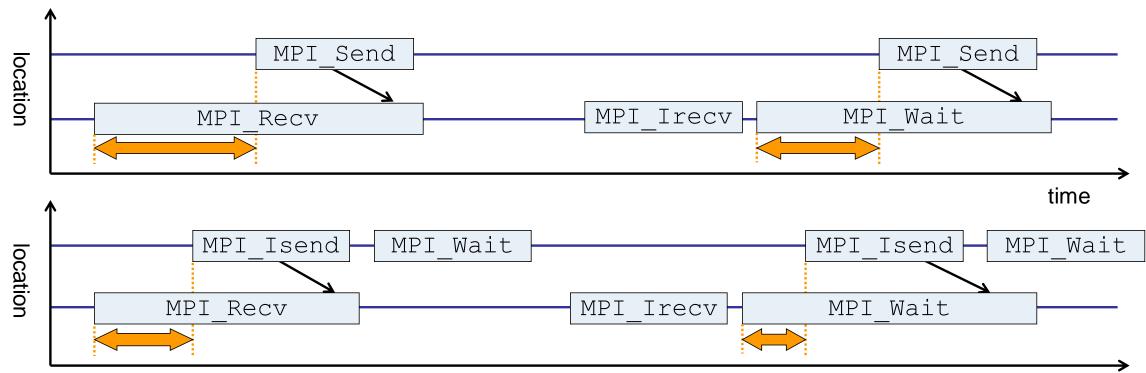
Current limitations:

- Unable to handle traces ...
 - with MPI thread level exceeding MPI_THREAD_FUNNELED
 - containing memory events, CUDA/OpenCL device events (kernel, memcpy), SHMEM, or OpenMP nested parallelism
- PAPI/rusage metrics for trace events are ignored

Scalasca workflow



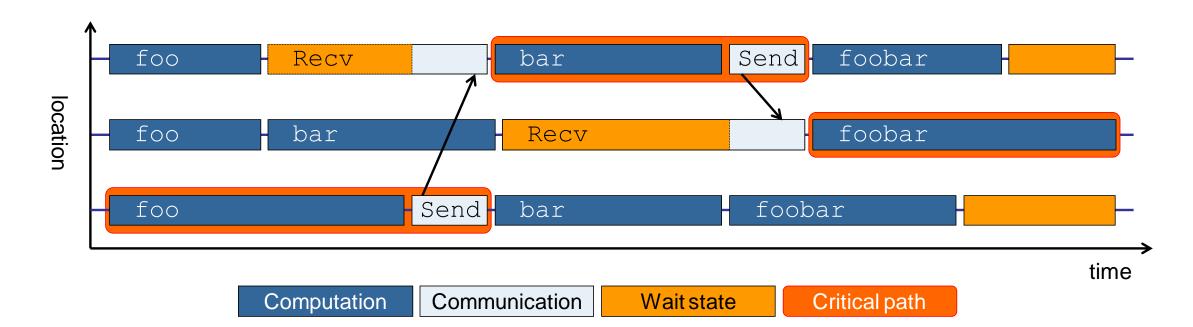
Example: "Late Sender" wait state



time

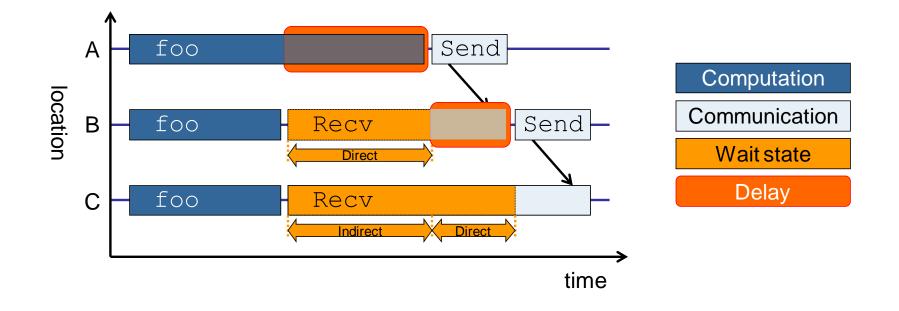
- Waiting time caused by a blocking receive operation posted earlier than the corresponding send
- Applies to blocking as well as non-blocking communication

Example: Critical path



- Shows call paths and processes/threads that are responsible for the program's wall-clock runtime
- Identifies good optimization candidates and parallelization bottlenecks

Example: Root-cause analysis



- Classifies wait states into direct and indirect (i.e., caused by other wait states)
- Identifies delays (excess computation/communication) as root causes of wait states
- Attributes wait states as *delay costs*



Hands-on: TeaLeaf MPI+CUDA





Recap: Setup for exercises

- Connect to your training account on JUWELS Booster (with X11-forwarding)
 - % ssh -X <yourid>@juwels-booster.fz-juelich.de
- Set account and default environment (NVHPC + ParaStationMPI) via helper script
 - % source /p/project/training2341/setup.sh
- Change to directory containing TeaLeaf_CUDA sources
 - Existing instrumented executable can be reused
 - $\frac{9}{6}$ cd \$WORK
 - **% cd TeaLeaf_CUDA**
- Load Scalasca module
 - Depends on (i.e., implicitly loads) Score-P & CubeGUI
 - % module load Scalasca

TeaLeaf_CUDA summary measurement collection...



% cd bin.scorep % cp ../jobscripts/juwelsbooster/scalasca.sbatch . % cat scalasca.sbatch # Score-P measurement configuration export SCOREP_CUDA_ENABLE=default export SCOREP_CUDA_BUFFER=48M #export SCOREP_CUDA_BUFFER=48M #export SCOREP_EXPERIMENT_DIRECTORY=scorep-tea_leaf-8 export SCOREP_FILTERING_FILE=../config/scorep.filt #export SCOREP_ENABLE_TRACING=true #export SCOREP_TOTAL_MEMORY=250M

Scalasca configuration
export SCAN_ANALYZE_OPTS="--time-correct"

Run the application
scan -s srun ./tea_leaf

% sbatch scalasca.sbatch

 Change to directory with the Score-P instrumented executable and edit the job script

Hint:

scan = scalasca -analyze
-s = profile/summary (default)

Submit the job

TeaLeaf_CUDA summary measurement



S=C=A=N: Scalasca 2.6.1 runtime summarization
S=C=A=N: ./scorep_tea_leaf_8_sum experiment archive
S=C=A=N: Wed Oct 25 14:51:20 2023: Collect start
srun ./tea leaf

Tea version 1.400

[... More application output ...]

S=C=A=N: Wed Oct 25 14:52:09 2023: Collect done (status=0) 49s S=C=A=N: ./scorep_tea_leaf_8_sum complete. Run the application using the Scalasca measurement collection & analysis nexus prefixed to launch command

 Creates experiment directory: scorep_tea_leaf_8_sum

TeaLeaf_CUDA summary analysis report examination



Score summary analysis report

% square -s scorep_tea_leaf_8_sum
INFO: Post-processing runtime summarization result (profile.cubex)...
INFO: Score report written to ./scorep_tea_leaf_8_sum/scorep.score

Post-processing and interactive exploration with Cube

% square scorep_tea_leaf_8_sum INFO: Displaying ./scorep_tea_leaf_8_sum/summary.cubex... Hint:

Copy 'summary.cubex' to local system (laptop) using 'scp' to improve responsiveness of GUI

[GUI showing summary analysis report]

 The post-processing derives additional metrics and generates a structured metric hierarchy

TeaLeaf_CUDA trace measurement collection...



% cp ../jobscripts/juwelsbooster/scalasca.sbatch .

 ${\,\rm \$}\,$ cat scalasca.sbatch

Score-P measurement configuration
export SCOREP_CUDA_ENABLE=runtime
export SCOREP_CUDA_BUFFER=48M
#export SCOREP_EXPERIMENT_DIRECTORY=scorep-tea_leaf-8
export SCOREP_FILTERING_FILE=../config/scorep.filt
#export SCOREP_ENABLE_TRACING=true
export SCOREP_TOTAL_MEMORY=250M

```
# Scalasca configuration
export SCAN ANALYZE OPTS="--time-correct"
```

```
# Run the application
scan -t srun ./tea leaf
```

 Change to directory with the Score-P instrumented executable and edit the job script

Hint:

scan = scalasca -analyze
-t = trace collection & analysis

Submit the job

TeaLeaf_CUDA trace measurement ... collection



S=C=A=N: Scalasca 2.6.1 trace collection and analysis S=C=A=N: Wed Oct 25 14:58:52 2023: Collect start srun ./tea leaf

Tea version 1.400

[... More application output ...]

S=C=A=N: Wed Oct 25 14:59:33 2023: Collect done (status=0) 41s

 Starts measurement with collection of trace files ...

TeaLeaf_CUDA trace measurement ... analysis



S=C=A=N: Wed Oct 25 14:59:33 2023: Analyze start srun scout.mpi --time-correct ./scorep tea leaf 8 trace/traces.otf2 SCOUT (Scalasca 2.6.1) Analyzing experiment archive ./scorep tea leaf 8 trace/traces.otf2 Opening experiment archive ... done (0.009s). Reading definition data ... done (0.008s). Reading event trace data... done (0.0005).Preprocessing... done (0.695s).Timestamp correction... done (0.721s).Analyzing trace data... done (10.108s).Writing analysis report... done (0.137s). : 426.422MB Max. memory usage # passes : 1
violated : 0 Total processing time : 12.489s S=C=A=N: Wed Oct 25 14:59:49 2023: Analyze done (status=0) 16s

 Continues with automatic (parallel) analysis of trace files

TeaLeaf CUDA trace analysis report exploration



 Produces trace analysis report in the experiment directory containing trace-based wait-state metrics

% square scorep_tea_leaf_8_trace INFO: Post-processing runtime summarization report (profile.cubex)... INFO: Post-processing trace analysis report (scout.cubex)... INFO: Displaying ./scorep_tea_leaf_8_trace/trace.cubex...

[GUI showing trace analysis report]

Hint:

Run 'square -s' first and then copy 'trace.cubex' to local system (laptop) using 'scp' to improve responsiveness of GUI



Demo: TeaLeaf MPI+OpenMP case study





Case study: TeaLeaf MPI+OpenMP

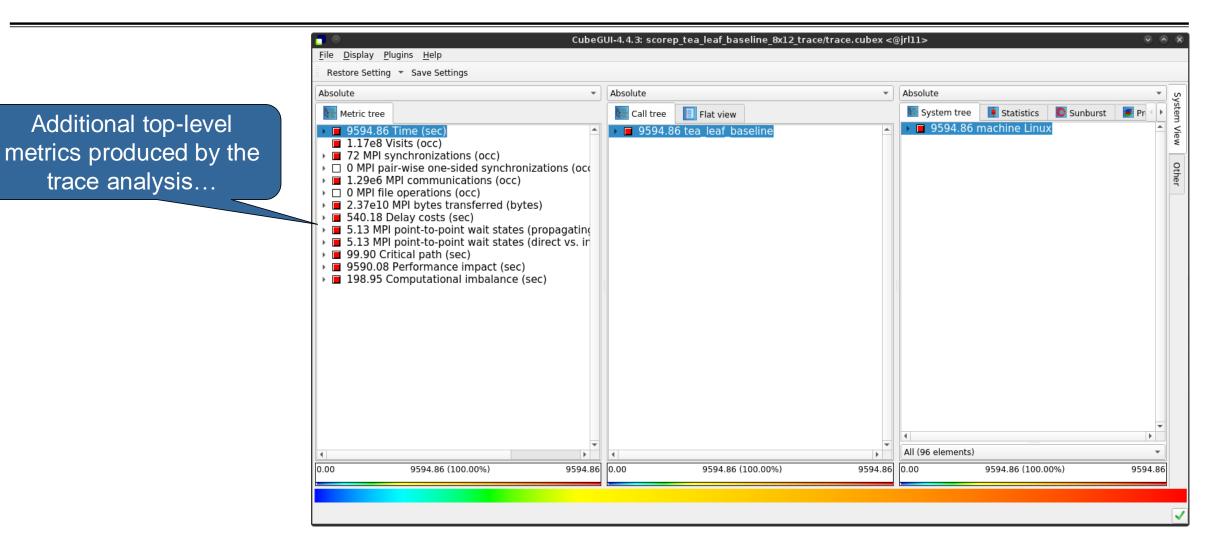
- HPC mini-app developed by the UK Mini-App Consortium
 - Solves the linear 2D heat conduction equation on a spatially decomposed regular grid using a 5 point stencil with implicit solvers
 - Part of the Mantevo 3.0 suite
 - Available on GitHub: https://uk-mac.github.io/TeaLeaf/
- Measurements of TeaLeaf reference v1.0 taken on Jureca cluster @ JSC
 - Using Intel 19.0.3 compilers, Intel MPI 2019.3, Score-P 5.0, and Scalasca 2.5
 - Run configuration
 - 8 MPI ranks with 12 OpenMP threads each
 - Distributed across 4 compute nodes (2 ranks per node)
 - Test problem "5": 4000 × 4000 cells, CG solver

% cube scorep_tea_leaf_baseline_8x12_trace/trace.cubex

[GUI showing post-processed trace analysis report]

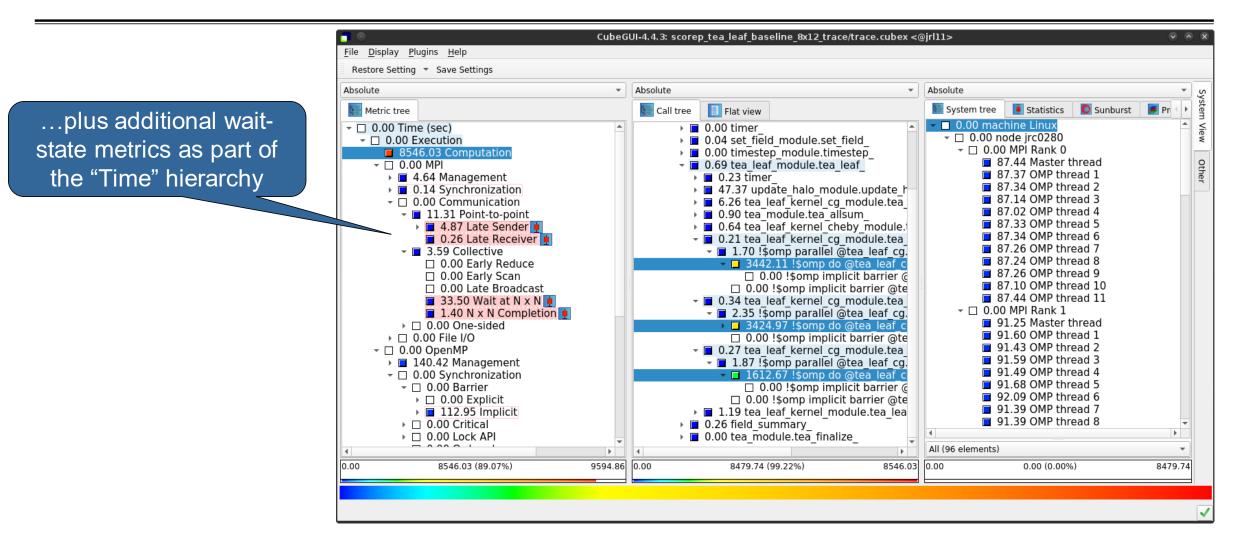


Scalasca analysis report exploration (opening view)



Scalasca wait-state metrics

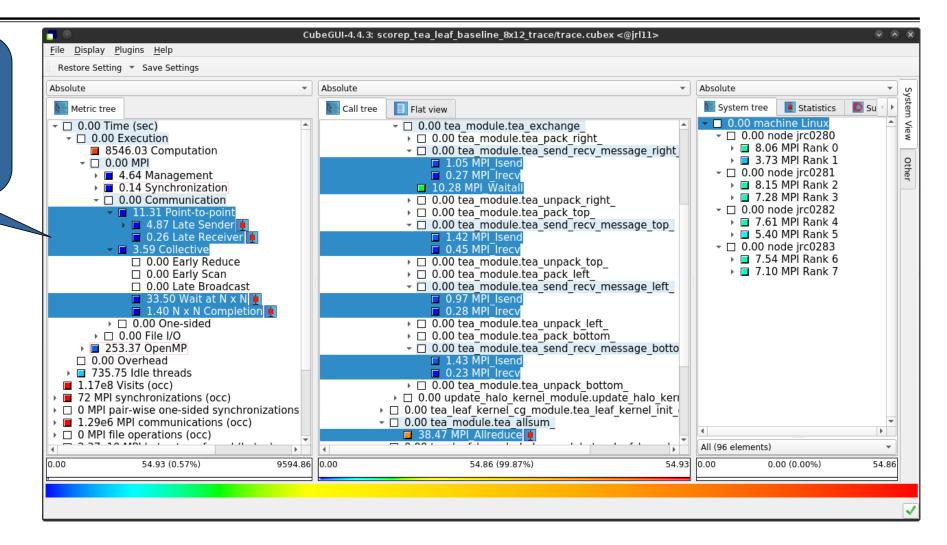




TeaLeaf Scalasca report analysis (I)



While MPI communication time and wait states are small (~0.6% of the total execution time)...



TeaLeaf Scalasca report analysis (II)



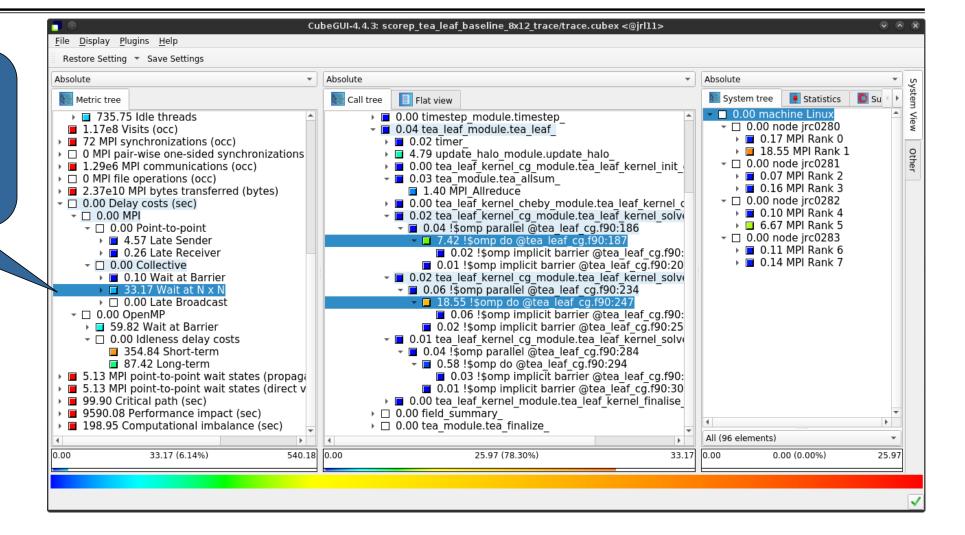
CubeGUI-4.4.3: scorep tea leaf baseline 8x12 trace/trace.cubex <@jrl11> File Display Plugins Help Restore Setting Save Settings Absolute Absolute Absolute * Flat view 🔚 System tree Statistics O Su 🔚 Call tree Metric tree 0.00 machine Linux 735.75 Idle threads 7.00 tea module.tea exchange 1.17e8 Visits (occ) 24.77 MPI Rank 0 72 MPI synchronizations (occ) 1.98 tea module.tea send recv message right 24.21 MPI Rank 1 Othe 0 MPI pair-wise one-sided synchronizations 11.56 MPI Isend 1.29e6 MPI communications (occ) 2.99 MPI Trecv 20.93 MPI Rank 2 D MPI file operations (occ) 56.82 MPI Waita 21.55 MPI Rank 3 2.37e10 MPI bytes transferred (bytes) 4.79 tea module.tea unpack right - 0.00 node jrc0282 6.85 tea module.tea pack top 0.00 Delay costs (sec) 23.46 MPI Rank 4 I.25 tea module.tea send recv message top → □ 0.00 MPI 24.15 MPI Rank 5 15.65 MPI Isend 0.00 Point-to-point 4.57 Late Sender 4.95 MPI Irecv 19.39 MPI Rank 6 7.10 tea module.tea unpack top 0.26 Late Receiver 20.40 MPI Rank 7 4.87 tea module.tea pack left □ 0.00 Collective 0.10 Wait at Barrier 1.92 tea module.tea send recv message left 33.17 Wait at N x N 10.63 MPI Isend → □ 0.00 Late Broadcast 3.13 MPI Trecv 4.59 tea module.tea unpack left 0.00 OpenMP 6.98 tea module.tea pack bottom 59.82 Wait at Barrier 1.34 tea module.tea send recv message botto 0.00 Idleness delay costs 15.69 MPI Isend 354.84 Short-term 87.42 Long-term 2.55 MPI Irecv 6.96 tea module.tea unpack bottom 5.13 MPI point-to-point wait states (propaga 3.83 update halo kernel module.update halo keri 5.13 MPI point-to-point wait states (direct v → ■ 3.55 tea leaf kernel cg module.tea leaf kernel init 99.90 Critical path (sec) 9590.08 Performance impact (sec) 9.87 tea module.tea allsum Þ. 198.95 Computational imbalance (sec) 54.90 MPI Allreduce All (96 elements) F 4 354.84 0.00 540.18 0.00 0.00 (0.00%) 178.86 0.00 354.84 (65.69%) 178.86 (50.41%)

...they directly cause a significant amount of the OpenMP thread idleness

TeaLeaf Scalasca report analysis (III)



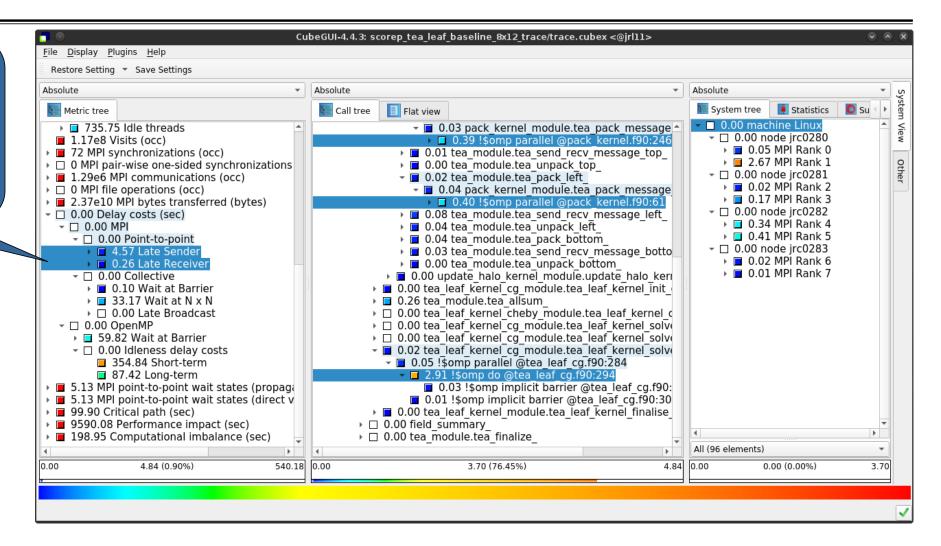
The "Wait at NxN" collective wait states are mostly caused by the first 2 OpenMP do loops of the solver (on ranks 5 & 1, resp.)...



TeaLeaf Scalasca report analysis (IV)



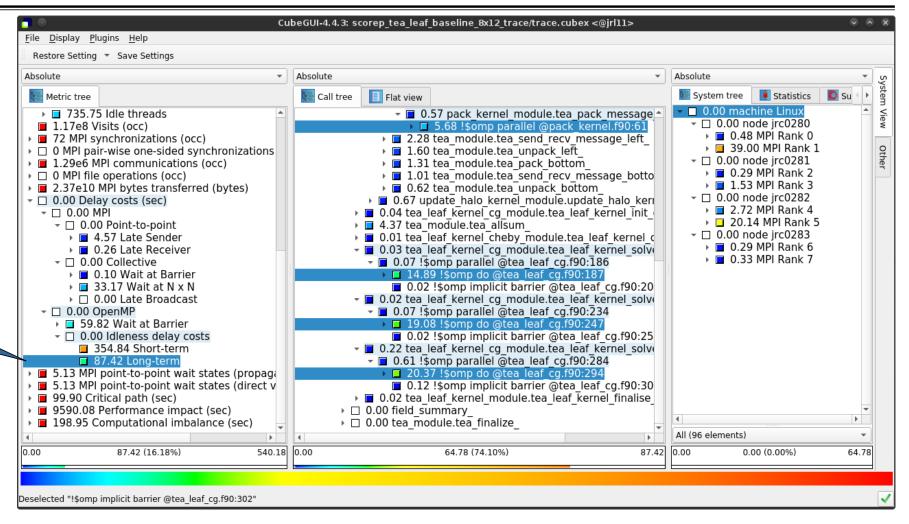
...while the MPI pointto-point wait states are caused by the 3rd solver do loop (on rank 1) and two loops in the halo exchange



TeaLeaf Scalasca report analysis (V)

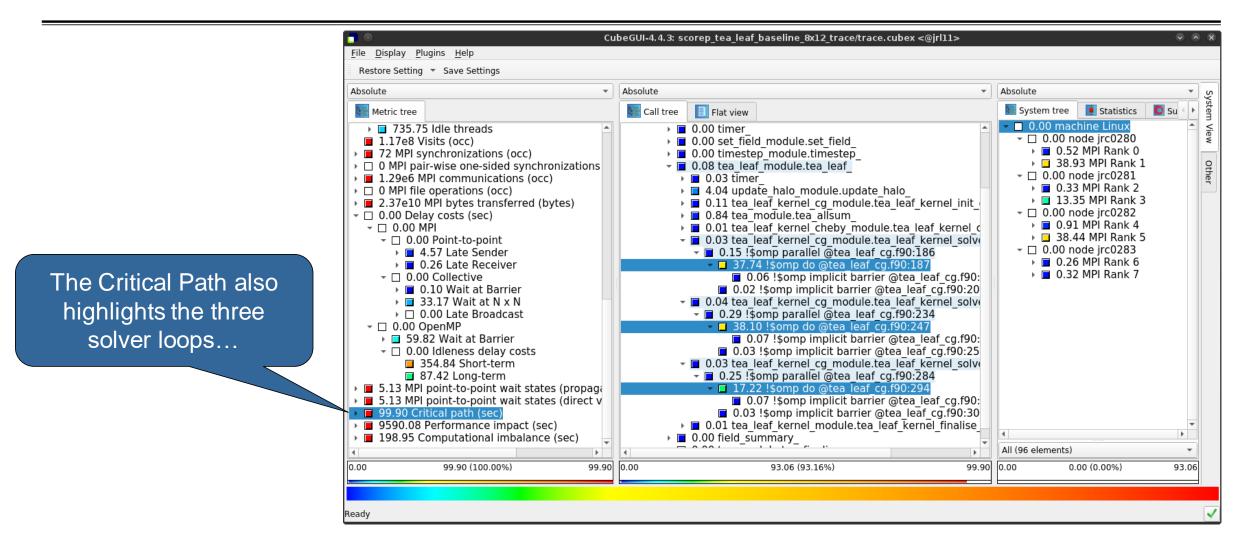


Various OpenMP do loops (incl. the solver loops) also cause OpenMP thread idleness on other ranks via propagation



TeaLeaf Scalasca report analysis (VI)





TeaLeaf Scalasca report analysis (VII)

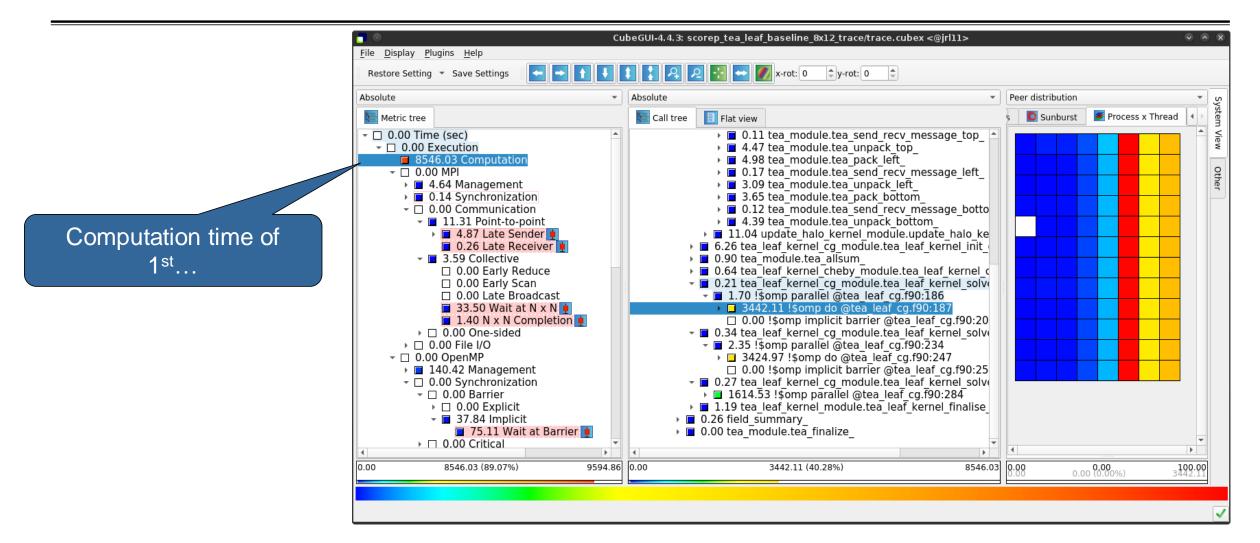


CubeGUI-4.4.3: scorep tea leaf baseline 8x12 trace/trace.cubex <@jrl11> File Display Plugins Help Restore Setting
 Save Settings Absolute Absolute Absolute * Statistics 📕 Flat view 🔚 System tree 🔘 Su Metric tree 🔚 Call tree 0.00 machine Linux 735.75 Idle threads 0.72 MPI Waitall ▶ □ 0.00 tea module.tea unpack right 1.17e8 Visits (occ) 0.03 MPI Rank 0 I 0.20 tea module.tea pack top 72 MPI synchronizations (occ) 3.07 MPI Rank 1 Othe D 0 MPI pair-wise one-sided synchronizations I 0.20 tea module.tea send recv message top 1.29e6 MPI communications (occ) • 0.21 tea module.tea unpack top Image: 0.01 MPI Rank 2 D MPI file operations (occ) I 0.28 tea module.tea pack left 0.28 MPI Rank 3 ▶ ■ 0.32 tea module.tea send recv message left 2.37e10 MPI bytes transferred (bytes) - 0.00 node jrc0282 I 0.23 tea module.tea unpack left 0.02 MPI Rank 4 I 0.31 tea module.tea pack bottom ¬ □ 0.00 MPI 2.30 MPI Rank 5 • 0.28 tea module.tea send recv message botto
 0.00 Point-to-point
 - 🗆 0.00 node jrc0283 4.57 Late Sender I 0.18 tea module.tea unpack bottom 0.01 MPI Rank 6 0.12 update halo kernel module.update halo keri 0.26 Late Receiver 0.02 MPI Rank 7 → ■ 0.02 tea leaf kernel cg module.tea leaf kernel init
 □ 0.00 Collective
 0.10 Wait at Barrier • 0.09 tea module.tea allsum 0.68 MPI Allreduce 33.17 Wait at N x N 🕨 🔲 0.00 tea leaf kernel cheby module.tea leaf kernel o → □ 0.00 Late Broadcast - 0.00 tea leaf kernel cg module.tea leaf kernel solve 🝷 🔲 0.01 !\$omp parallel @tea leaf cg.f90:186 59.82 Wait at Barrier ▶ 🗖 1.90 !\$omp do @tea leaf cg.f90:187 0.00 Idleness delay costs 0.01 !\$omp implicit barrier @tea leaf cg.f90:20 354.84 Short-term - 🖬 0.00 tea leaf kernel cg module.tea leaf kernel solv 87.42 Long-term 0.02 !\$omp parallel @tea_leaf_cg.f90:234
 2.45 !\$omp do @tea_leaf_cg.f90:247 5.13 MPI point-to-point wait states (propaga 5.13 MPI point-to-point wait states (direct v 0.01 !\$omp implicit barrier @tea leaf cg.f90:25 90.49 Critical path (sec) 9.41 Imbalance 🗝 🔲 0.00 tea leaf kernel cg module.tea leaf kernel solv 0.47 !\$omp parallel @tea leaf cg.f90:284 Þ. 9590.08 Performance impact (sec) All (96 elements) Ŧ 9.41 0.00 99.90 0.00 0.00 (0.00%) 0.00 9.41 (9.42%) 5.75 (61.12%) 5.75

...with imbalance (time on critical path above average) mostly in the first two loops and MPI communication

TeaLeaf Scalasca report analysis (VIII)





3424.97 (40.08%)

TeaLeaf Scalasca report analysis (IX)



CubeGUI-4.4.3: scorep tea leaf baseline 8x12 trace/trace.cubex <@jrl11> File Display Plugins Help Restore Setting
 Save Settings x-rot: 0 v-rot: 0 ٢ Peer distribution Absolute Absolute System Process x Thread 🔚 Call tree Flat view Sunburst Metric tree • 0.11 tea module.tea send recv message top View 0.00 Execution 4.47 tea module.tea unpack top 8546.03 Computation ▶ ■ 0.17 tea module.tea send recv message left Othe • □ 0.00 MPI 3.09 tea module.tea unpack left 4.64 Management 3.65 tea module.tea pack bottom 0.14 Synchronization 0.00 Communication Image: 0.12 tea_module.tea_send_recv_message_botto 11.31 Point-to-point ...and 2nd do loop 11.04 update halo kernel module.update halo ke 4.87 Late Sender 0.26 Late Receiver ▶ ■ 6.26 tea leaf kernel cg module.tea leaf kernel init mostly balanced within ▶ ■ 0.90 tea module.tea allsum 3.59 Collective • 0.64 tea leaf kernel cheby module.tea leaf kernel o 0.00 Early Reduce - 0.21 tea leaf kernel cg module.tea leaf kernel solv 0.00 Early Scan each rank, but vary → ■ 1.70 !\$omp parallel @tea leaf cg.f90:186 0.00 Late Broadcast 3442.11 !somp do @tea leaf cg.f90:187 33.50 Wait at N x N considerably across 0.00 !\$omp implicit barrier @tea leaf cg.f90:20 1.40 N x N Completion 🝷 🖬 0.34 tea leaf kernel cg module tea leaf kernel solv 0.00 One-sided ranks... ▶ □ 0.00 File I/O 🝷 🔲 2.35 !\$omp parallel @tea leaf cg.f90:234 3424.97 !\$omp do @tea leaf cg.f90:247 - 0.00 OpenMP 0.00 !\$omp implicit barrier @tea leaf cg.f90:25 140.42 Management • 0.27 tea leaf kernel cg module.tea leaf kernel solve ■ 1614.53 !\$omp parallel @tea leaf cg.f90:284 - 🗆 0.00 Barrier 🗖 1.19 tea leaf kernel module.tea leaf kernel finalise ▶ □ 0.00 Explicit 37.84 Implicit 🕨 🗖 0.26 field summary 75.11 Wait at Barrier ▶ ■ 0.00 tea module.tea finalize D 0.00 Critical • F 4

9594.86

0.00

0.00

8546.03 (89.07%)

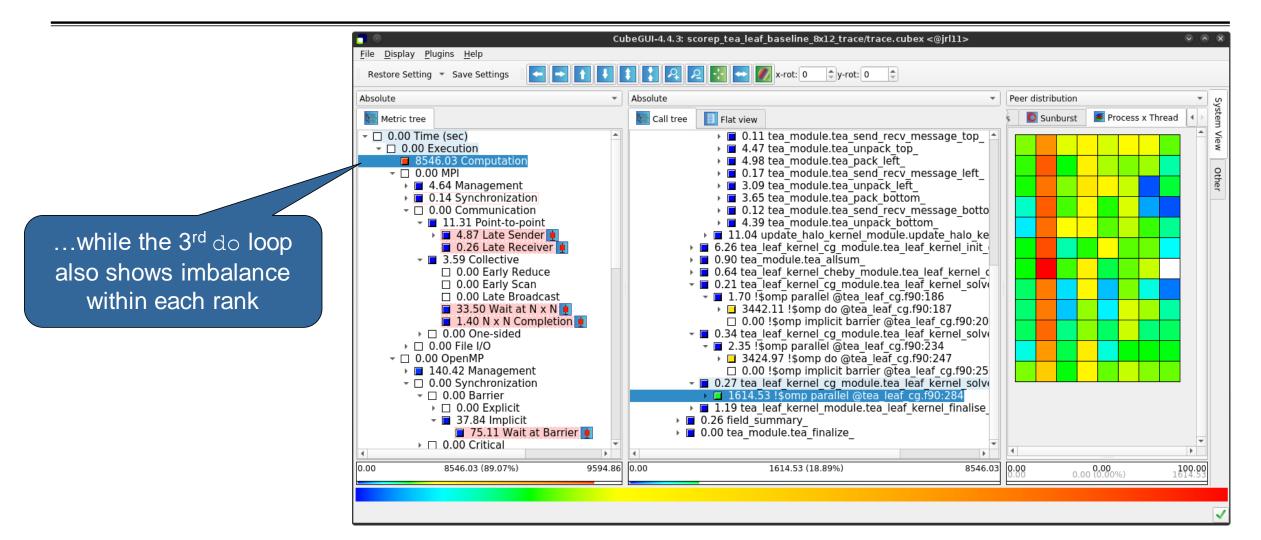
100.00

8546.03 0.00

0.00 (0.00

TeaLeaf Scalasca report analysis (X)





SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)

TeaLeaf analysis summary

- The first two OpenMP do loops of the solver are well balanced within a rank, but are imbalanced across ranks
 - → Requires a global load balancing strategy
- The third OpenMP do loop, however, is imbalanced within ranks,
 - causing direct "Wait at OpenMP Barrier" wait states,
 - which cause indirect MPI point-to-point wait states,
 - which in turn cause OpenMP thread idleness
 - → Low-hanging fruit
- Adding a SCHEDULE (guided) clause reduced
 - the MPI point-to-point wait states by ~66%
 - the MPI collective wait states by ~50%
 - the OpenMP "Wait at Barrier" wait states by ~55%
 - the OpenMP thread idleness by ~11%
 - → Overall runtime (wall-clock) reduction by ~5%

Scalasca Trace Tools: Further information

- Collection of trace-based performance tools
 - Specifically designed for large-scale systems
 - Features an automatic trace analyzer providing wait-state, critical-path, and delay analysis
 - Supports MPI, OpenMP, POSIX threads, and hybrid MPI+OpenMP/Pthreads
- Available under 3-clause BSD open-source license
- Documentation & sources:
 - https://www.scalasca.org
- Contact:
 - mailto: scalasca@fz-juelich.de



Acknowledgement



This tutorial is sponsored by the DEEP-SEA project.





The DEEP Projects have received funding from the European Commission's FP7, H2020, and EuroHPC Programmes, under Grant Agreements n° 287530, 610476, 754304, and 955606.

The EuroHPC Joint Undertaking (JU) receives support from the European Union's Horizon 2020 research and innovation programme and Germany, France, Spain, Greece, Belgium, Sweden, United Kingdom, Switzerland.



Reference material





Scalasca command – One command for (almost) everything



```
<sup>9</sup> scalasca
Scalasca 2.6.1
Toolset for scalable performance analysis of large-scale parallel applications
usage: scalasca [OPTION]... ACTION <argument>...
    1. prepare application objects and executable for measurement:
       scalasca -instrument <compile-or-link-command> # skin (using scorep)
    2. run application under control of measurement system:
       scalasca -analyze <application-launch-command> # scan
    3. interactively explore measurement analysis report:
       scalasca -examine <experiment-archive|report> # square
Options:
   -c, --show-config
                         show configuration summary and exit
                         show this help and exit
   -h, --help
   -n, --dry-run
                         show actions without taking them
       --quickref
                         show quick reference quide and exit
       --remap-specfile
                         show path to remapper specification file and exit
   -v, --verbose
                         enable verbose commentary
                         show version information and exit
   -V, --version
```

■ The `scalasca -instrument' command is deprecated and will be remove in the next major release
⇒ use Score-P instrumenter directly

Scalasca convenience command: scan / scalasca -analyze



% scan		
Scalasca 2.6.1: measurement collection & analysis nexus		
<pre>usage: scan {options} [launchcmd [launchargs]] target [targetargs]</pre>		
where {options} may include:		
-h		: show this brief usage message and exit.
-v		: increase verbosity.
		: show command(s) to be launched but don't execute.
-q	Quiescent	: execution with neither summarization nor tracing.
-s	Summary	: enable runtime summarization. [Default]
-t	Tracing	: enable trace collection and analysis.
-a	Analyze	: skip measurement to (re-)analyze an existing trace.
-e	exptdir	: Experiment archive to generate and/or analyze.
		(overrides default experiment archive title)
-f	filtfile	: File specifying measurement filter.
-1	lockfile	: File that blocks start of measurement.
-R	#runs	: Specify the number of measurement runs per config.
-M	cfgfile	: Specify a config file for a multi-run measurement.
-P	preset	: Specify a preset for a multi-run measurement, e.g., 'pop'.
-L	-	: List available multi-run presets.
-D	cfgfile	: Check a multi-run config file for validity and dump
		: the processed configuration for comparison.

Scalasca measurement collection & analysis nexus

Automatic measurement configuration



- scan configures Score-P measurement by automatically setting some environment variables and exporting them
 - E.g., experiment title, profiling/tracing mode, filter file, ...
 - Precedence order:
 - Command-line arguments
 - Environment variables already set
 - Automatically determined values
- Also, scan includes consistency checks and prevents corrupting existing experiment directories
- For tracing experiments, after trace collection completes then automatic parallel trace analysis is initiated
 - Uses identical launch configuration to that used for measurement (i.e., the same allocated compute resources)

Scalasca convenience command: square / scalasca -examine

<pre>% square Scalasca 2.6.1: analysis</pre>	roport ovploror						
usage: square [OPTIONS] <experiment archive="" cube="" file="" =""></experiment>							
-C <none full<="" quick="" th="" =""><th>> : Level of sanity checks for newly created reports</th></none>	> : Level of sanity checks for newly created reports						
-c <number></number>	: Consider number of counters when doing scoring (-s)						
-F	: Force remapping of already existing reports						
-f filtfile	: Use specified filter file when doing scoring (-s)						
-s	: Skip display and output textual score report						
-v	: Enable verbose mode						
-n	: Do not include idle thread metric						
-S <mean merge="" =""></mean>	: Aggregation method for summarization results of						
	each configuration (default: merge)						
-T <mean merge="" =""></mean>	: Aggregation method for trace analysis results of						
	each configuration (default: merge)						
-A	: Post-process every step of a multi-run experiment						
-I	: Ignore structural sanity checks and force aggregation						
	of measurements in a multi-run experiment						
-x <scorep-score opt=""></scorep-score>	: Pass option(s) to scorep-score						

Scalasca analysis report explorer (Cube)

Scalasca advanced command: scout - Scalasca automatic trace analyzer



```
% scout.hyb --help
SCOUT (Scalasca 2.6.1)
Copyright (c) 1998-2022 Forschungszentrum Juelich GmbH
Copyright (c) 2014-2021 RWTH Aachen University
Copyright (c) 2009-2014 German Research School for Simulation Sciences GmbH
Usage: <launchcmd> scout.hyb [OPTION]... <ANCHORFILE | EPIK DIRECTORY>
Options:
  --statistics
                    Enables instance tracking and statistics [default]
                    Disables instance tracking and statistics
  --no-statistics
  --critical-path
                    Enables critical-path analysis [default]
  --no-critical-path Disables critical-path analysis
                    Enables root-cause analysis [default]
  --rootcause
                    Disables root-cause analysis
  --no-rootcause
  --single-pass
                    Single-pass forward analysis only
  --time-correct
                    Enables enhanced timestamp correction
  --no-time-correct
                    Disables enhanced timestamp correction [default]
  --verbose, -v
                    Increase verbosity
  --help
                    Display this information and exit
```

Provided in serial (.ser), OpenMP (.omp), MPI (.mpi) and MPI+OpenMP (.hyb) variants

Scalasca advanced command: clc_synchronize



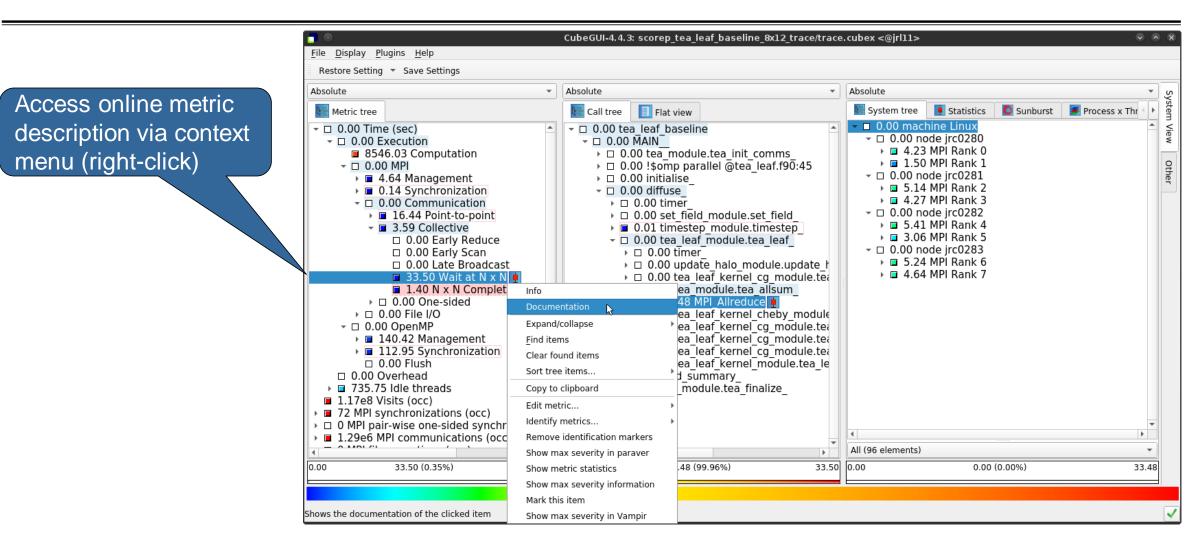
Scalasca trace event timestamp consistency correction

Usage: <launchcmd> clc_synchronize.hyb <ANCHORFILE | EPIK_DIRECTORY>

- Provided in MPI (.mpi) and MPI+OpenMP (.hyb) variants
- Takes as input a trace experiment archive where the events may have timestamp inconsistencies
 E.g., multi-node measurements on systems without adequately synchronized clocks on each compute node
- Generates a new experiment archive (always called ./clc_sync) containing a trace with event timestamp inconsistencies resolved
 - E.g., suitable for detailed examination with a time-line visualizer

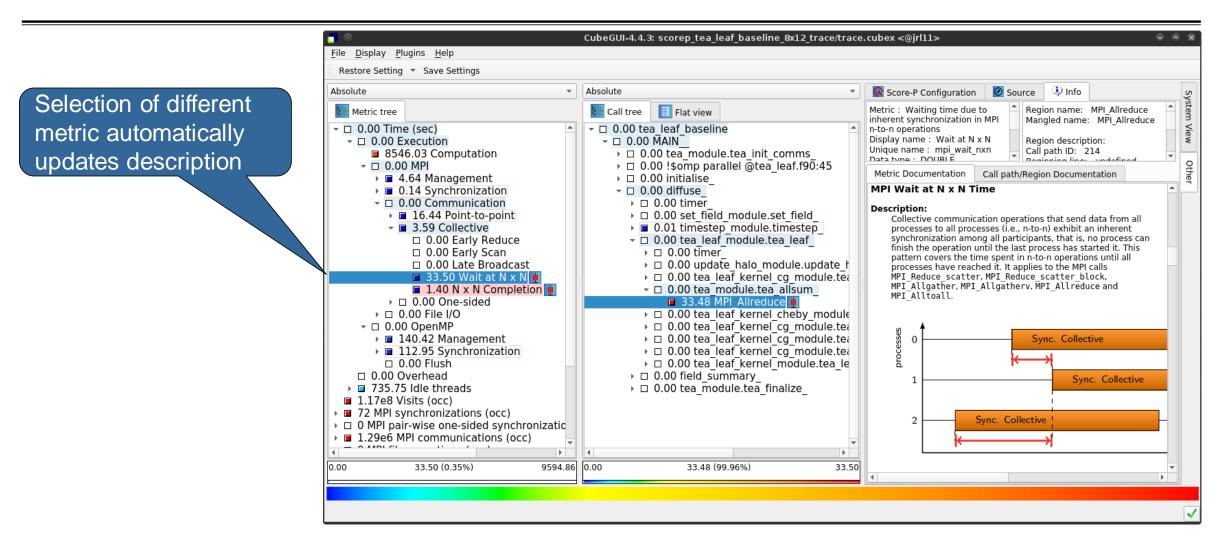
Online metric description





Online metric description (cont.)





menu

VIRTUAL INSTITUTE – HIGH PRODUCTIVITY SUPERCOMPUTING

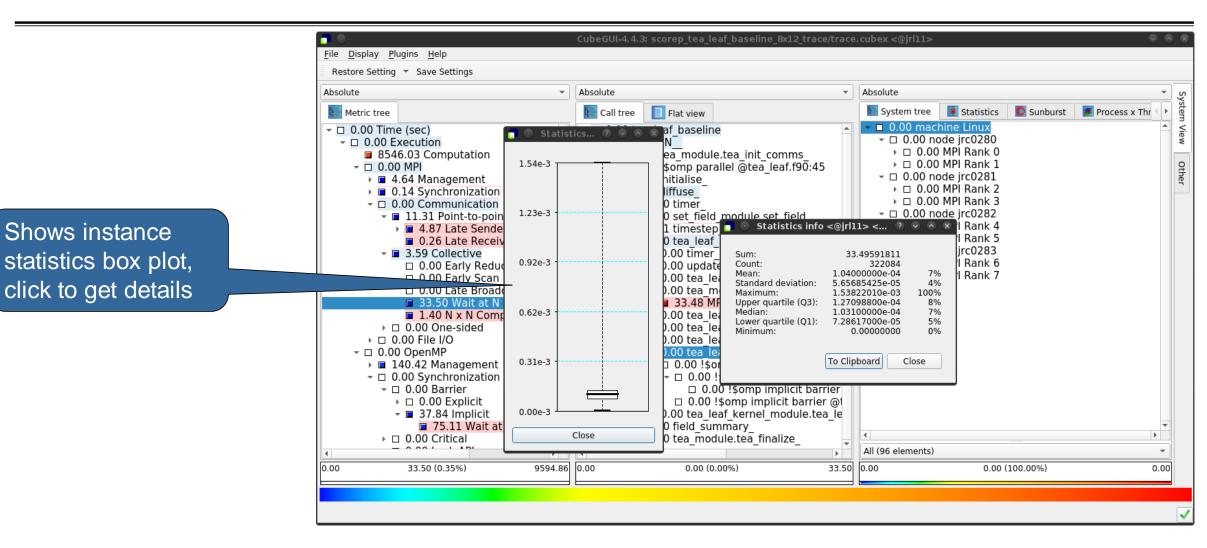
Metric statistics



CubeGUI-4.4.3: scorep tea leaf baseline 8x12 trace/trace.cubex <@jrl11> File Display Plugins Help Restore Setting
 Save Settings Absolute Absolute Absolute Ŧ Syste 🚺 Sunburst 🔚 System tree Statistics 🕖 Process x Thr Call tree Flat view Metric tree 🔺 🖃 0.00 machine Linux 0.00 Time (sec) 0.00 tea leaf baseline View - D 0.00 MAIN 0.00 Execution 4.23 MPI Rank 0 8546.03 Computation D 0.00 tea module.tea init comms → □ 0.00 !\$omp parallel @tea leaf.f90:45 I 1.50 MPI Rank 1 Othe 0.00 MPI - 0.00 node jrc0281 ▶ □ 0.00 initialise 5.14 MPI Rank 2 0.14 Synchronization □ 0.00 diffuse
 Access metric statistics 4.27 MPI Rank 3 0.00 Communication ▶ □ 0.00 timer - 0.00 node jrc0282 I1.31 Point-to-point ▶ □ 0.00 set field module.set field 5.41 MPI Rank 4 Image: 0.01 timestep module.timestep for metrics marked with 4.87 Late Sender 3.06 MPI Rank 5 0.26 Late Receiver - 0.00 node jrc0283 3.59 Collective ▶ □ 0.00 timer box plot icon from context 5.24 MPI Rank 6 0.00 Early Reduce D 0.00 update halo module.update ł 4.64 MPI Rank 7 → □ 0.00 tea leaf kernel cg module.tea 0.00 Early Scan 0.00 Late Broadcast 48 MPI Allreduce 33.50 Wait at N x N Info 1.40 N x N Complet ea leaf kernel cheby module Documentation ea leaf kernel cg module.tea D 0.00 One-sided ea_leaf_kernel_cg_module.tea Expand/collapse ea_leaf_kernel_cg_module.tea - □ 0.00 OpenMP Find items 0 !somp parallel @tea leaf c ▶ ■ 140.42 Management Clear found items 0.00 !\$omp do @tea leaf cg.f □ 0.00 Synchronization → □ 0.00 !\$omp implicit barrier □ 0.00 Barrier Sort tree items ... 0.00 !\$omp implicit barrier @I 0.00 Explicit Copy to clipboard 37.84 Implicit ea leaf kernel module.tea le Edit metric... 75.11 Wait at Ba [▶]d summary module.tea finalize Þ ▶ □ 0.00 ical Identify metrics... All (96 elements) Ŧ Remove identification markers Þ. 33.50 (0.35%) .48 (99.96%) 33.50 0.00 0.00 (0.00%) 33.48 0.00 Show max severity in paraver Show metric statistics Show max severity information Mark this item Show max severity in Vampir

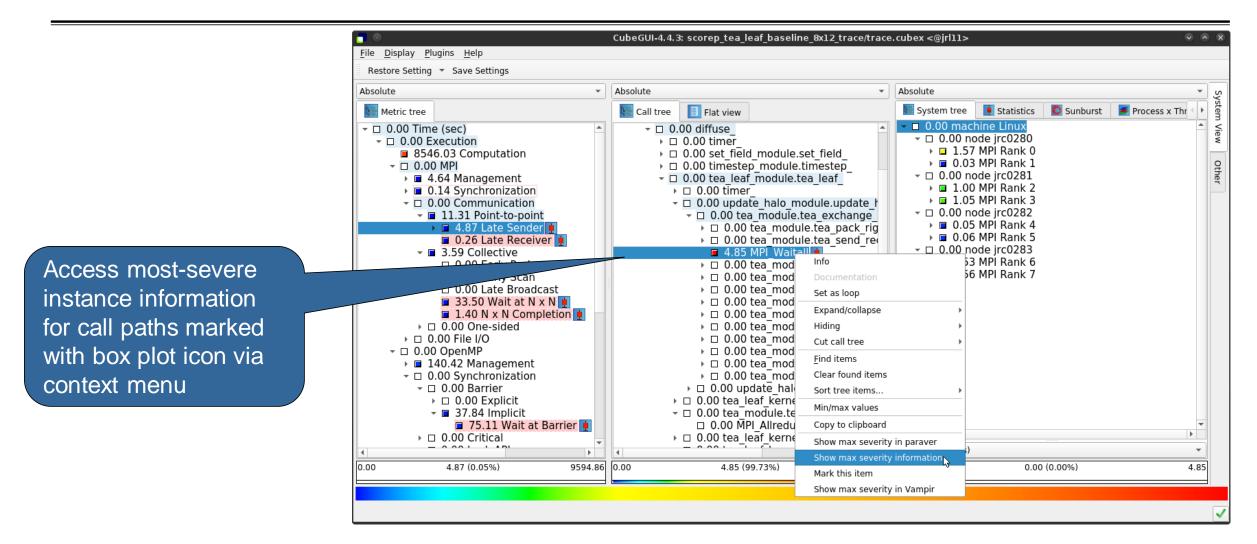
Metric statistics (cont.)



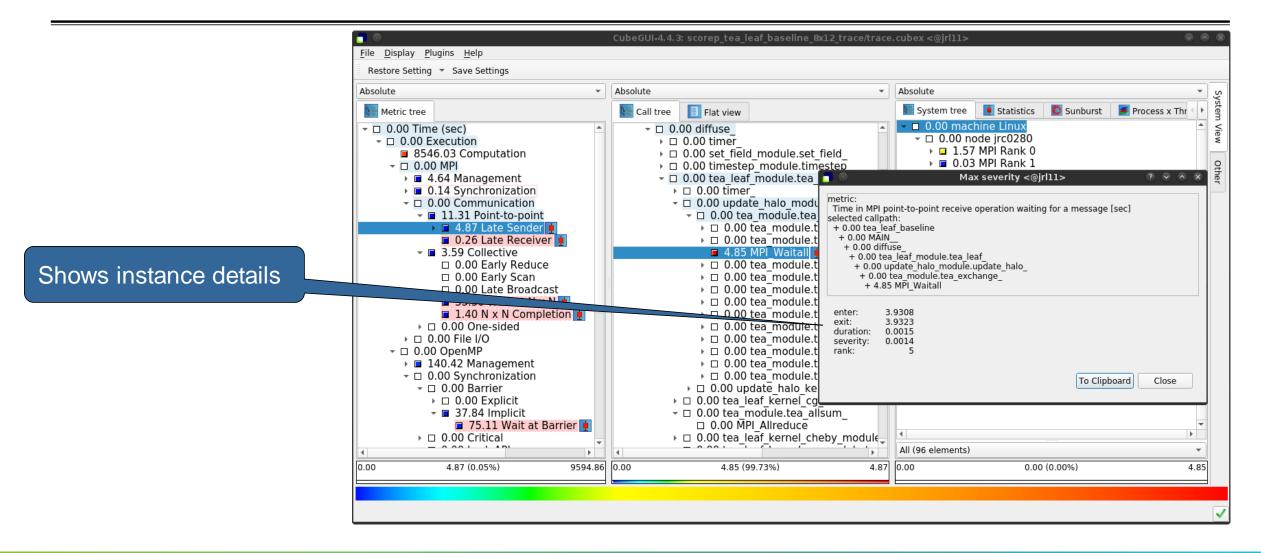


Metric instance statistics



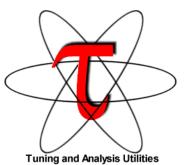


Metric instance statistics (cont.)





Performance Data Management with TAU PerfExplorer



Sameer Shende <u>sameer@cs.uoregon.edu</u>

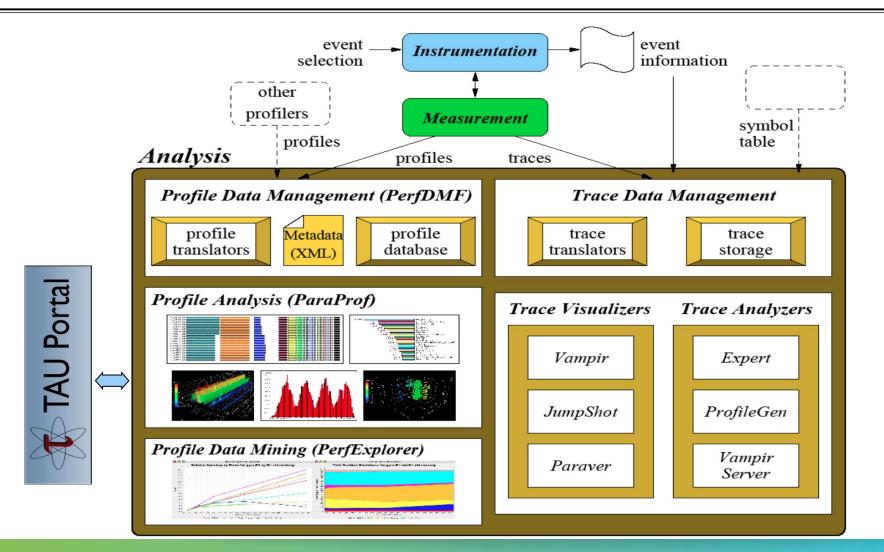
University of Oregon

http://tau.uoregon.edu/TAU_PerfExplorer_SC23.pdf

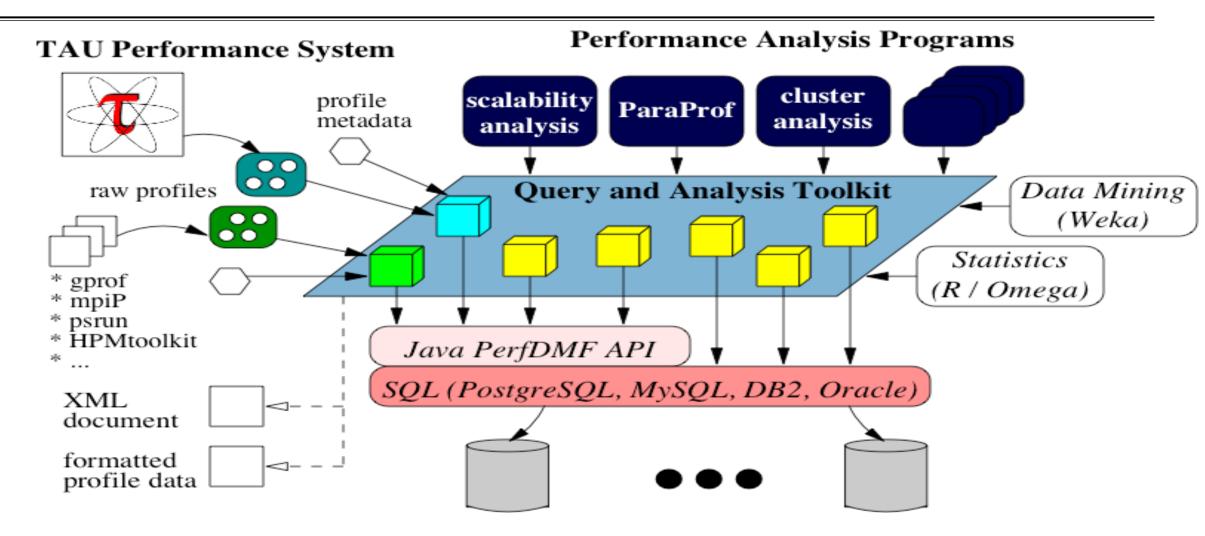


TAU's Analysis Tools: PerfExplorer

TAU Analysis



TAUdb: Performance Data Management Framework



Using TAUdb

Configure TAUdb (Done by each user)

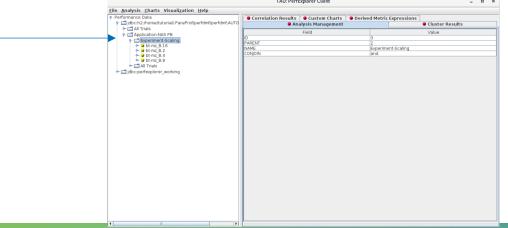
% taudb_configure --create-default

- Choose derby, PostgreSQL, MySQL, Oracle or DB2
- Hostname
- Username
- Password
- Say yes to downloading required drivers (we are not allowed to distribute these)
- Stores parameters in your ~/.ParaProf/taudb.cfg file
- Configure PerfExplorer (Done by each user)
 - % perfexplorer_configure
- Execute PerfExplorer
 - % perfexplorer

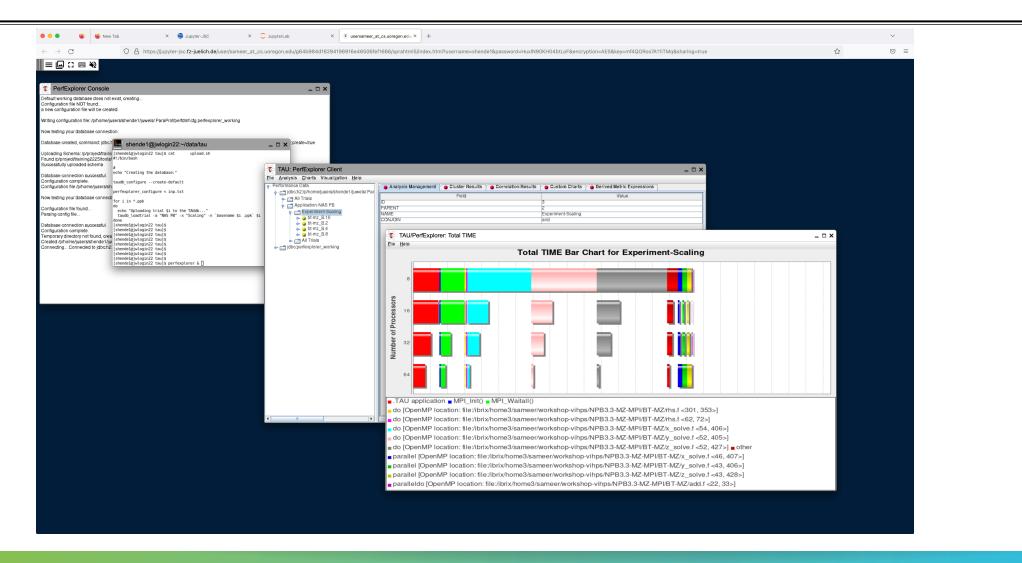


PerfExplorer Hands-on

% source /p/project/training2341/setup.sh % wget http://tau.uoregon.edu/data.tgz % tar xf data.tgz; cd data/tau % cat ./upload.sh % ./upload.sh % perfexplorer # Click and expand first database jdbc:/home/tutorial/.ParaProf/perfdmf.cfg # Click and expand Application-NAS PB and Experiment-Scaling # Select by clicking Experiment-Scaling # Select by clicking Experiment-Scaling # Go to Charts -> Total Execution Time, then select other charts



PerfExplorer Hands-on



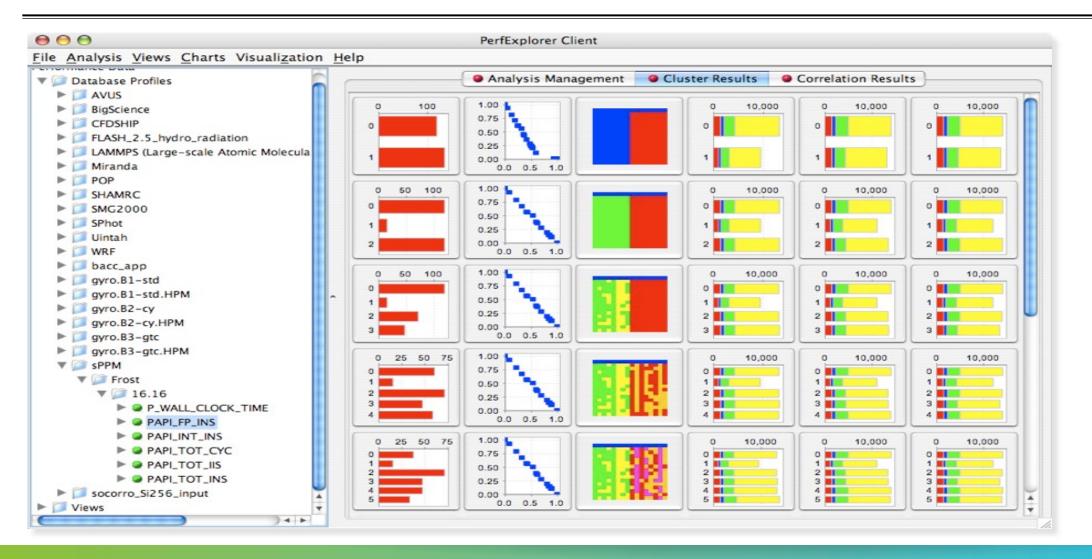
Performance Data Mining (PerfExplorer)

- Performance knowledge discovery framework
 - Data mining analysis applied to parallel performance data
 - comparative, clustering, correlation, dimension reduction, ...
 - Use the existing TAU infrastructure
 - TAU performance profiles, taudb
 - Client-server based system architecture
- Technology integration
 - Java API and toolkit for portability
 - taudb
 - R-project/Omegahat, Octave/Matlab statistical analysis
 - WEKA data mining package
 - JFreeChart for visualization, vector output (EPS, SVG)

PerfExplorer: Using Cluster Analysis

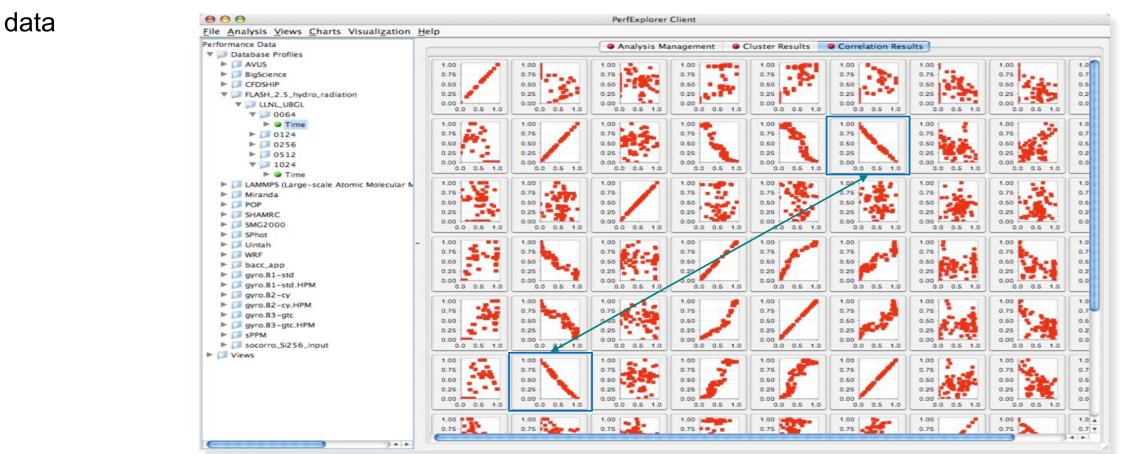
- Performance data represented as vectors each dimension is the cumulative time for an event
- k-means: k random centers are selected and instances are grouped with the "closest" (Euclidean) center
- New centers are calculated and the process repeated until stabilization or max iterations
- Dimension reduction necessary for meaningful results
- Virtual topology, summaries constructed

PerfExplorer - Cluster Analysis (sPPM)



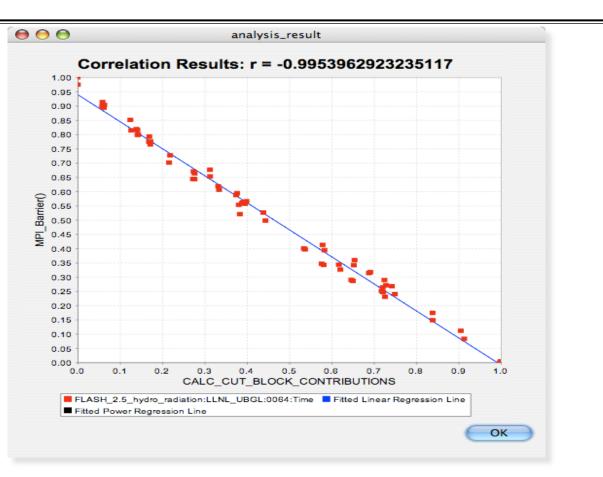
PerfExplorer - Correlation Analysis (Flash)

• Describes strength and direction of a linear relationship between two variables (events) in the



PerfExplorer - Correlation Analysis (Flash)

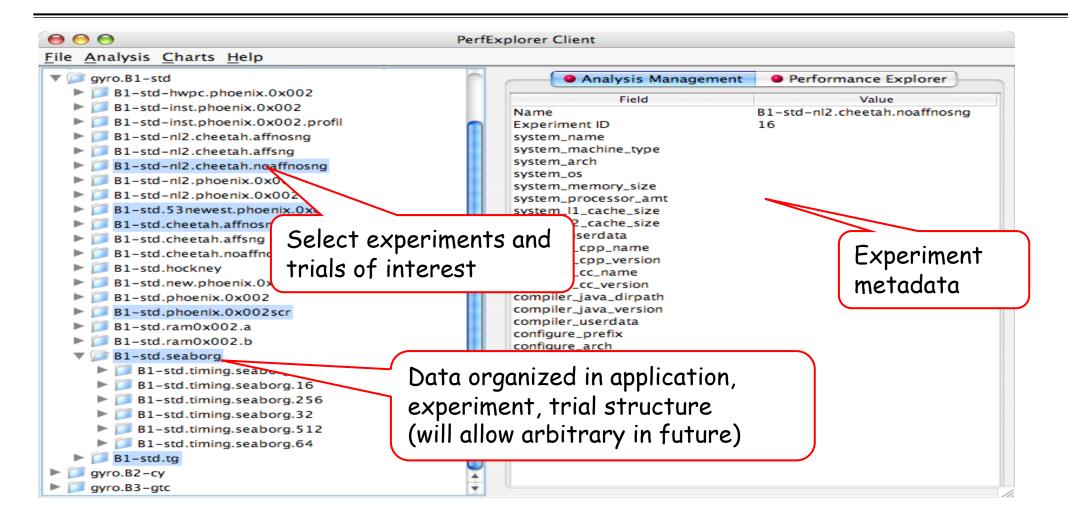
 -0.995 indicates strong, negative relationship
 As CALC_CUT_ BLOCK_CONTRIBUTIONS() increases in execution time, MPI_Barrier() decreases



PerfExplorer - Comparative Analysis

- Relative speedup, efficiency
 - total runtime, by event, one event, by phase
- Breakdown of total runtime
- Group fraction of total runtime
- Correlating events to total runtime
- Timesteps per second

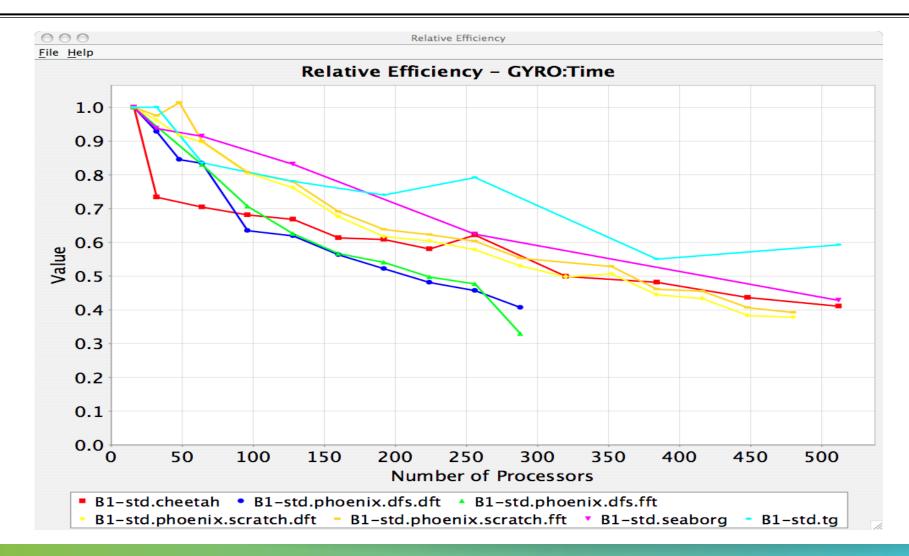
PerfExplorer - Interface



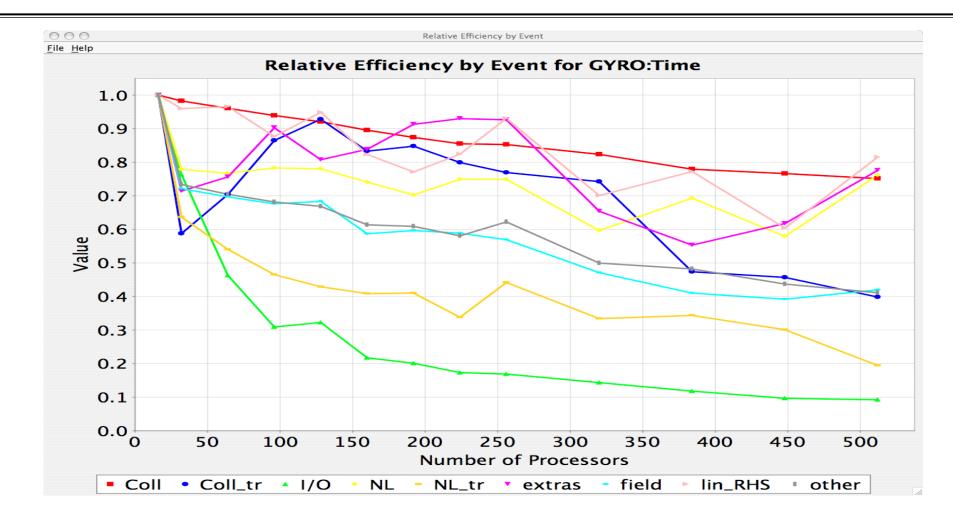
PerfExplorer - Interface

000		PerfEx	plorer Client		
ile <u>A</u> nalysis	<u>C</u> harts <u>H</u> elp				
V 📁 gyro.B1-s	Set <u>G</u> roup Name		Analysis Management	Performance Explorer	
🕨 📁 B1-stc	Set Metric of Interest		Field	Value	
B1-stc	B1-stc Set Total Number of Timesteps B1-stc Timesteps Per Second B1-stc Relative Efficiency B1-stc Relative Efficiency by Event B1-stc Relative Efficiency for One Event		Name	B1-std-nl2.cheetah.noaffnosng	
B1-stc			Experiment ID 16 system_name system_machine_type		
▶ 📁 B1-stc					
B1-stc					
B1-stc			system_arch		
► 📁 B1-stc			system_os		
B1-stc Relative Speedup					
B1-stc	Relative Speedup by Event				
B1-stc	Relative Speedup for One Event	Selec	t analysis		
► B1-str	<u>Communication Time / Total Ru</u>		A		
▶ 🗐 B1-stc.	Runtime Breakdown		compiler_cpp_name		
 B1-std.hockney B1-std.new.phoenix.0x002 B1-std.phoenix.0x002 B1-std.phoenix.0x002scr B1-std.ram0x002.a B1-std.ram0x002.b B1-std.seaborg B1-std.timing.seaborg.128 		_	compiler_cpp_version		
		-	compiler_cc_name compiler_cc_version		
			compiler_java_dirpath		
			compiler_java_version		
			compiler_userdata		
			configure_prefix		
			configure_arch configure_cpp		
			configure_cc		
▶ 2 B1-std.timing.seaborg.16			configure_jdk		
	std.timing.seaborg.256		configure_profile		
	std.timing.seaborg.32		configure_userdata		
	std.timing.seaborg.512		userdata		
	std.timing.seaborg.64				
B1-std.					
gyro.B2-cy	-	Ų			
gyro.B2-cy		-			
p gyro.b3-gt		<u> </u>			

PerfExplorer - Relative Efficiency Plots

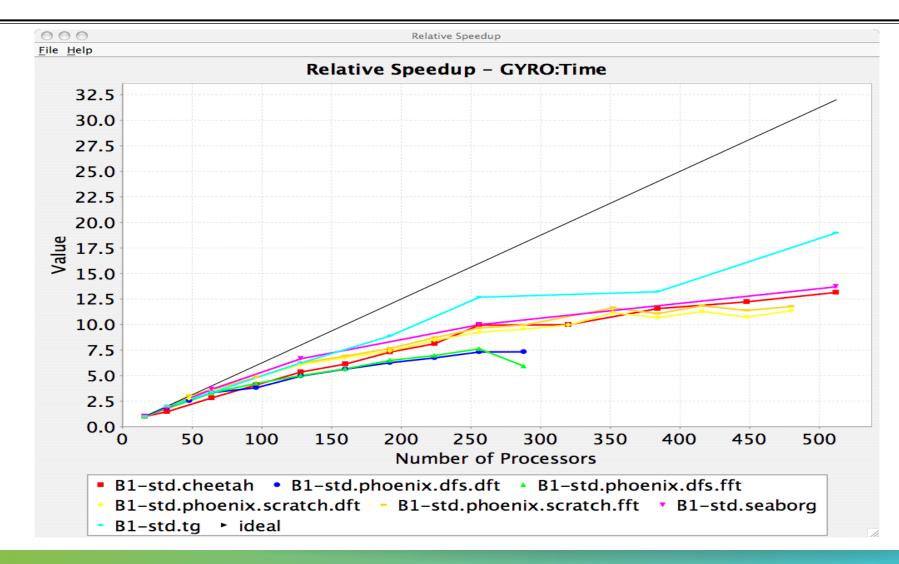


PerfExplorer - Relative Efficiency by Routine

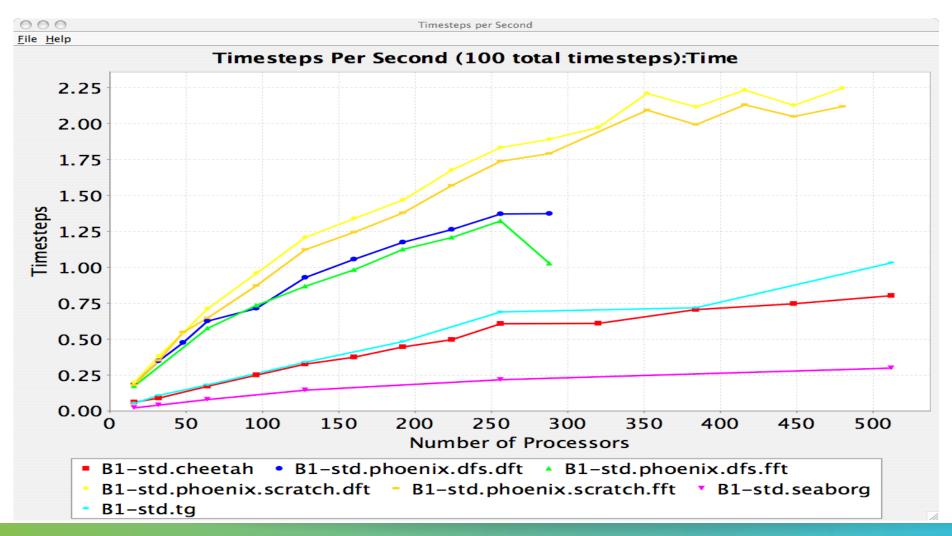


 \times

PerfExplorer - Relative Speedup



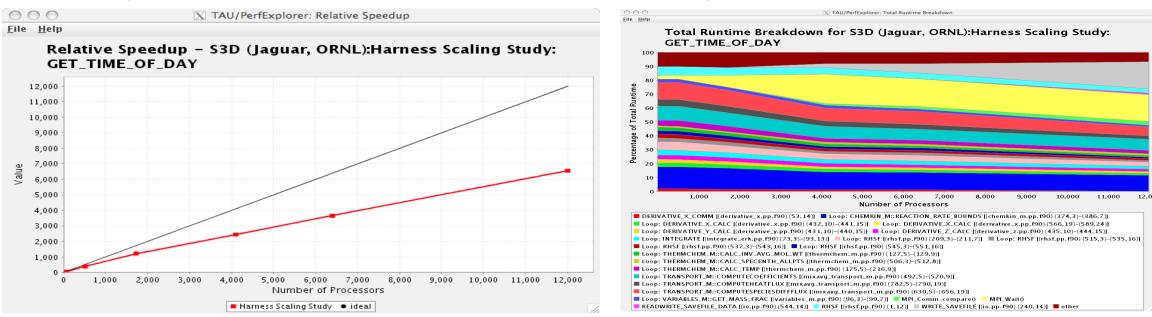
PerfExplorer - Timesteps Per Second



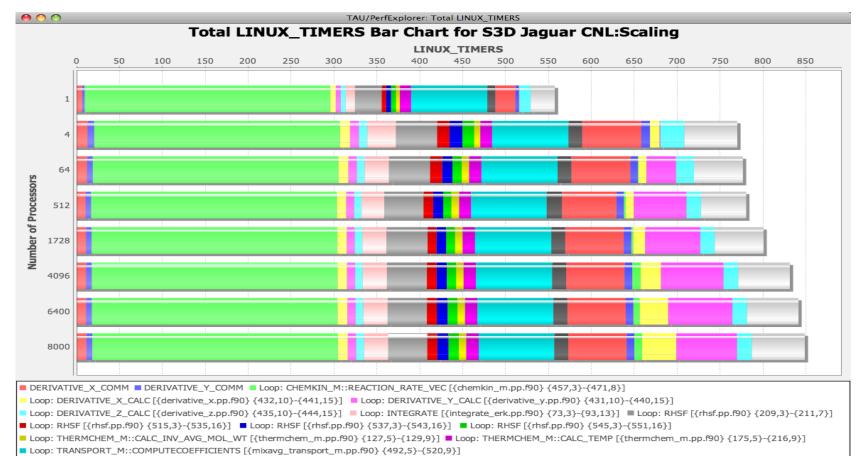


Evaluate Scalability

- Goal: How does my application scale? What bottlenecks occur at what core counts?
- Load profiles in taudb database and examine with PerfExplorer



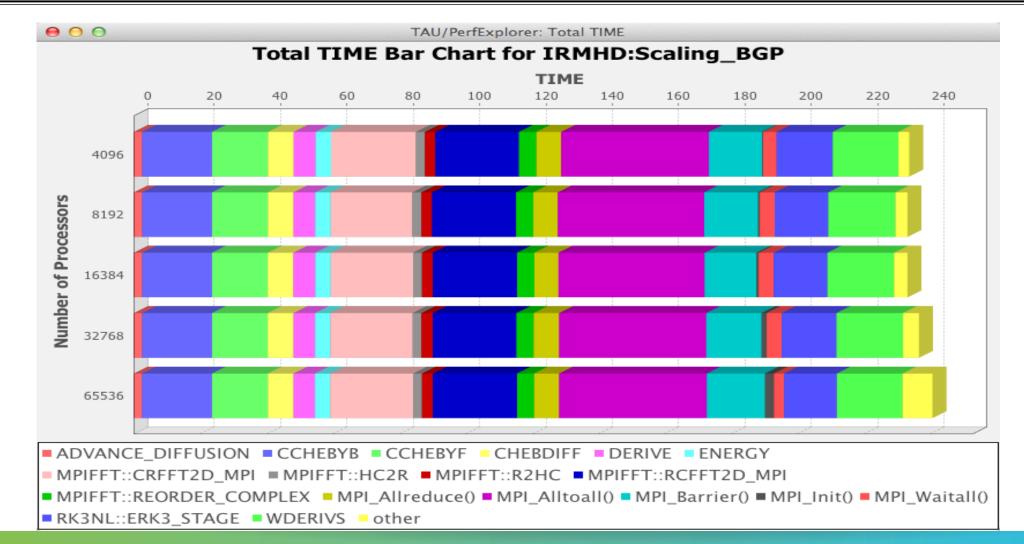
Evaluate Scalability



Loop: TRANSPORT_M::COMPUTEHEATFLUX [{mixavg_transport_m.pp.f90} {782,5}-{790,19}]

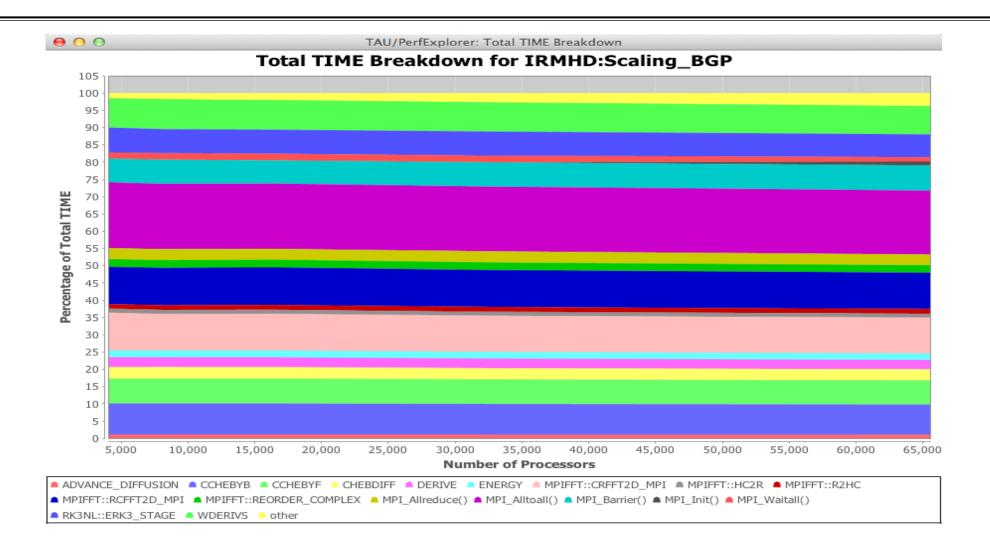
Loop: TRANSPORT_M::COMPUTESPECIESDIFFFLUX [{mixavg_transport_m.pp.f90} {630,5}-{656,19}] Loop: VARIABLES_M::GET_MASS_FRAC [{variables_m.pp.f90} {96,3}-{99,7}]
MPI_Barrier() MPI_Isend() MPI_Wait() RHSF other

PerfExplorer

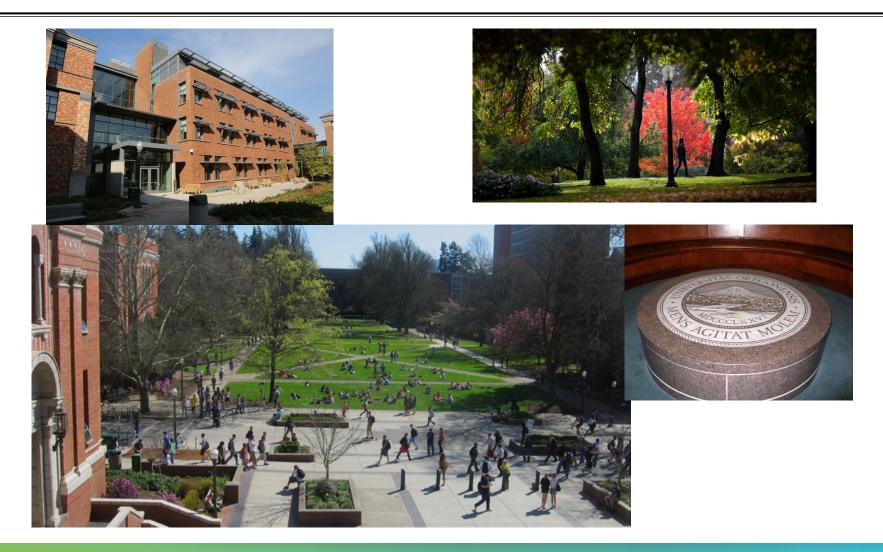


VI-HPS

PerfExplorer



Performance Research Lab, University of Oregon, Eugene, USA



Support Acknowledgments





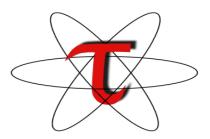
Acknowledgement

This research was supported by the Exascale Computing Project (17-SC-20-SC), a collaborative effort of two U.S. Department of Energy organizations (Office of Science and the National Nuclear Security Administration) responsible for the planning and preparation of a capable exascale ecosystem, including software, applications, hardware, advanced system engineering, and early testbed platforms, in support of the nation's exascale computing imperative.





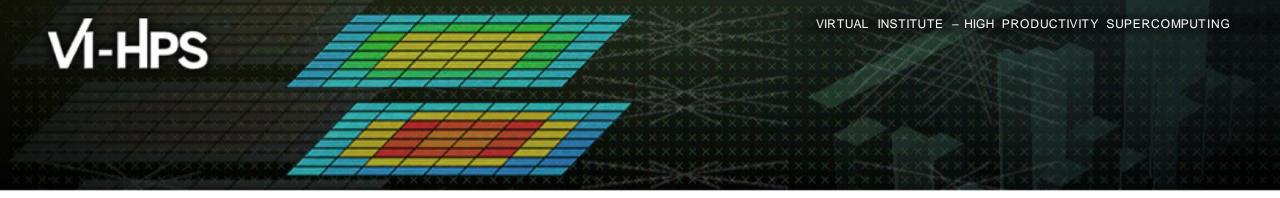
Download TAU from U. Oregon



http://tau.uoregon.edu

http://www.hpclinux.com [LiveDVD, OVA] https://e4s.io [Containers for Extreme-Scale Scientific Software Stack]

Free download, open source, BSD license



Score-P: Specialized Measurements and Analyses



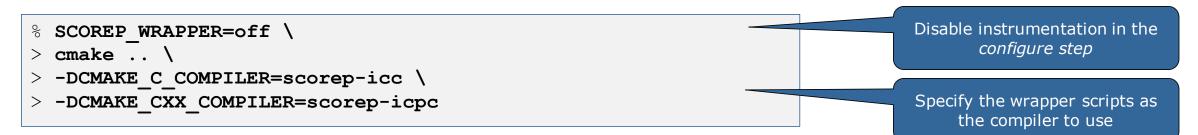




Mastering build systems



- Hooking up the Score-P instrumenter scorep into complex build environments like Autotools or CMake was always challenging
- Score-P provides convenience wrapper scripts to simplify this
- Autotools and CMake need the used compiler already in the configure step, but instrumentation should not happen in this step, only in the build step



- Allows to pass addition options to the Score-P instrumenter and the compiler via environment variables without modifying the *Makefiles*
- Run scorep-wrapper --help for a detailed description and the available wrapper scripts of the Score-P installation

Score-P user instrumentation API



- No replacement for automatic compiler instrumentation
- Can be used to further subdivide functions
 - E.g., multiple loops inside a function
- Can be used to partition application into coarse grain phases
 - E.g., initialization, solver, & finalization
- Enabled with --user flag to Score-P instrumenter
- Available for Fortran / C / C++

Score-P user instrumentation API (Fortran)



```
#include "scorep/SCOREP User.inc"
subroutine foo(...)
  ! Declarations
  SCOREP USER REGION DEFINE ( solve )
  ! Some code...
  SCOREP USER REGION BEGIN( solve, "<solver>", \setminus
                             SCOREP USER REGION TYPE LOOP )
  do i=1,100
   end do
  SCOREP USER REGION END( solve )
  ! Some more code...
end subroutine
```

- Requires processing by the C preprocessor
 - For most compilers, this can be automatically achieved by having an uppercase file extension, e.g., main.F or main.F90

Score-P user instrumentation API (C/C++)



```
#include "scorep/SCOREP User.h"
void foo()
 /* Declarations */
 SCOREP USER REGION DEFINE ( solve )
 /* Some code... */
  SCOREP USER REGION BEGIN( solve, "<solver>",
                             SCOREP USER REGION TYPE LOOP )
 for (i = 0; i < 100; i++)
    [...]
  SCOREP USER REGION END( solve )
  /* Some more code... */
```

Score-P user instrumentation API (C++)



```
#include "scorep/SCOREP User.h"
void foo()
  // Declarations
  // Some code ...
    SCOREP USER REGION( "<solver>",
                         SCOREP USER REGION TYPE LOOP )
    for (i = 0; i < 100; i++)
      [...]
  // Some more code...
```

Score-P measurement control API



Can be used to temporarily disable measurement for certain intervals

- Annotation macros ignored by default
- Enabled with --user flag

```
#include "scorep/SCOREP_User.inc"
subroutine foo(...)
! Some code...
SCOREP_RECORDING_OFF()
! Loop will not be measured
do i=1,100
[...]
end do
SCOREP_RECORDING_ON()
! Some more code...
end subroutine
```

```
#include "scorep/SCOREP_User.h"
void foo(...) {
   /* Some code... */
   SCOREP_RECORDING_OFF()
   /* Loop will not be measured */
   for (i = 0; i < 100; i++) {
      [...]
   }
   SCOREP_RECORDING_ON()
   /* Some more code... */</pre>
```

Fortran (requires C preprocessor)

C / C++

Enriching measurements with performance counters



Record metrics from PAPI:

```
% export SCOREP_METRIC_PAPI=PAPI_TOT_CYC
```

```
% export SCOREP_METRIC_PAPI_PER_PROCESS=PAPI_L3_TCM
```

• Use PAPI tools to get available metrics and valid combinations:

```
% papi_avail
```

% papi_native_avail

Record metrics from Linux perf:

- % export SCOREP_METRIC_PERF=cpu-cycles
- % export SCOREP_METRIC_PERF_PER_PROCESS=LLC-load-misses
- Use the perf tool to get available metrics and valid combinations:

% perf list

- Write your own metric plugin
 - Repository of available plugins: https://github.com/score-p

Only the master thread records the metric (assuming all threads of the process access the same L3 cache)

Mastering application memory usage



- Determine the maximum heap usage per process
- Find high frequent small allocation patterns
- Find memory leaks
- Support for:
 - C, C++, MPI, and SHMEM (Fortran only for GNU Compilers)
 - Profile and trace generation (profile recommended)
 - Memory leaks are recorded only in the profile
 - Resulting traces are not supported by Scalasca yet

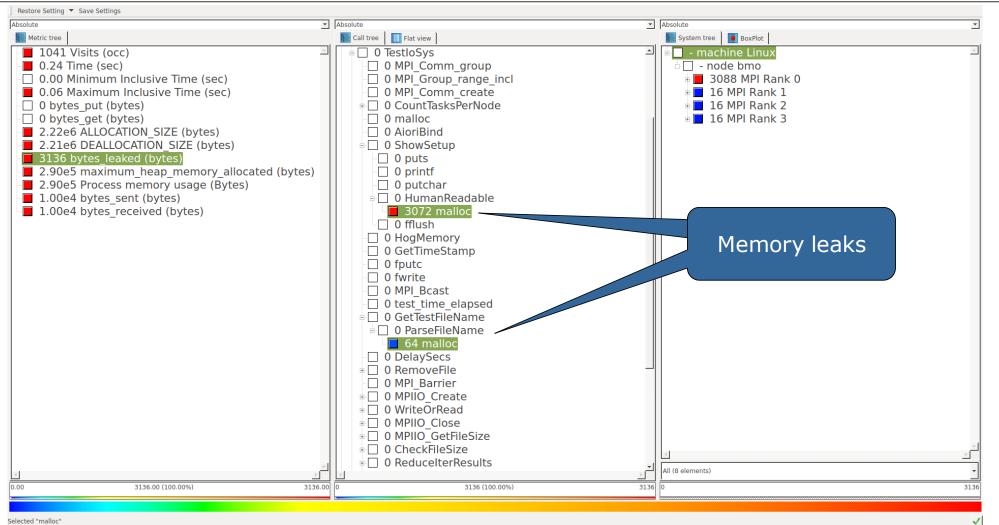
% export SCOREP_MEMORY_RECORDING=true % export SCOREP_MPI_MEMORY_RECORDING=true Set new configuration variable to enable memory recording

Mastering application memory usage



Restore Setting 🔻 Save Settings				
Absolute	*		Absolute	
Wetric tree	1	Call tree Flat view	System tree BoxPlot	
1041 Visits (occ)	<u>^</u>	🛚 🗖 2.90e5 main	🖻 📃 – machine Linux 🖆	
0.24 Time (sec)		- PER PROCESS METRICS	🗄 🗌 - node bmo	
- 0.00 Minimum Inclusive Time (sec)			🛛 🔲 2.90e5 MPI Rank 0	
- 0.06 Maximum Inclusive Time (sec)			🗉 🧧 2.87e5 MPI Rank 1	
- 0 bytes_put (bytes)			■ 2 .87e5 MPI Rank 2	
0 bytes_get (bytes)				
2.22e6 ALLOCATION_SIZE (bytes)				
2.21e6 DEALLOCATION_SIZE (bytes)				
 3136 bytes_leaked (bytes) 2.90e5 maximum_heap_memory_allocated (by 	(toc)			
 2.90e5 maximum_neap_memory_anocated (by 2.90e5 Process memory usage (Bytes) 	(les)			
 1.00e4 bytes_sent (bytes) 				
1.00e4 bytes_sent (bytes)				
1.00e4 bytes_received (bytes)				
			Different maxim	
			Different maxin	num
			hoop ucogoe r	aar
			heap usages p	Jei
			ranks	
			TUTINS	,
			× >	
	-		All (8 elements)	
0 2.90e5 (100.00%)	2 90e5		0.00 2.90e5	
2.5053 (100.0070)	2.5025	21901 51902 (0.00 MIN - TI 3103 T 34005 T 1100 T 45 1 3 T 104 2 0 1020 1 3 20 1 4 3 0 3 0 3 1 1 10	2.5065	
Selected "main"			1	

Mastering application memory usage



Selected "malloc"

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Hybrid measurement with sampling

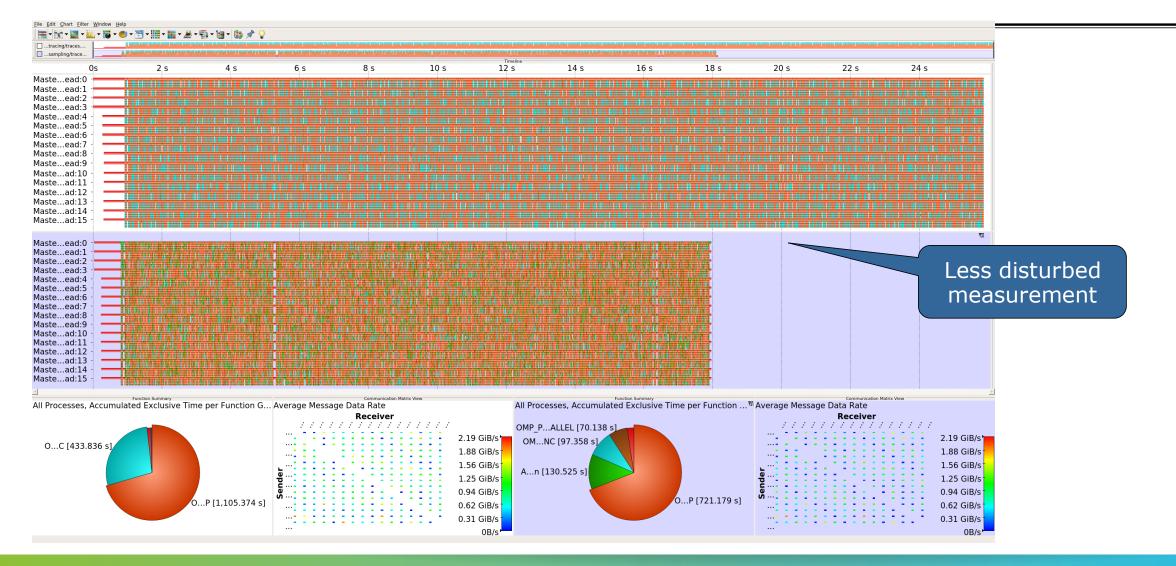


- Automatic compiler instrumentation greatly disturbs C++ applications because of frequent/short function calls => Use sampling instead
- Novel combination of sampling events and instrumentation of MPI, OpenMP, ...
 - Sampling replaces only compiler instrumentation (use --nocompiler)
 - Instrumentation is still used for parallel activities (MPI, OpenMP, CUDA, I/O)
- Supports profile and trace generation
 - % export SCOREP_ENABLE_UNWINDING=true
 - % # use the default sampling frequency
 - % #export SCOREP_SAMPLING_EVENTS=perf_cycles@2000000

 Set new configuration variable to enable sampling

Mastering C++ applications





Wrapping calls to 3rd party libraries

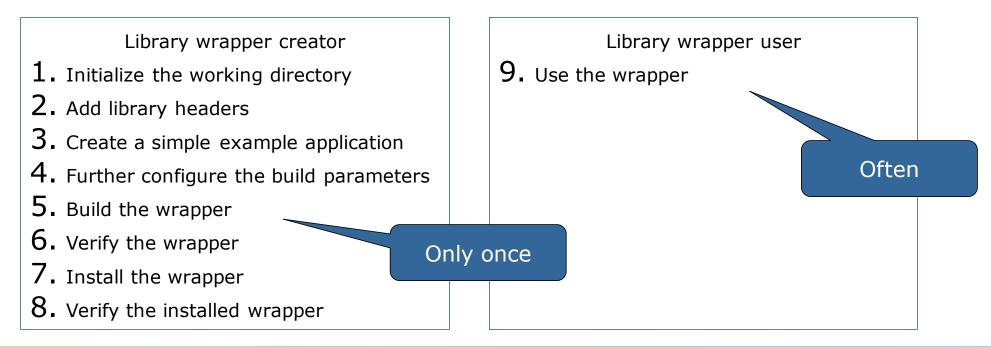


- Score-P does not record function calls to non-instrumented external libraries
- Increase insight into the behavior of the application
 - How does the application use the external library?
 - How does this compares to the usage of other libraries?
- Manual user instrumentation of the application using the library should be avoided
- Vendor provided libraries cannot be instrumented, but API provided in headers

Wrapping calls to 3rd party libraries: Library wrapper generator



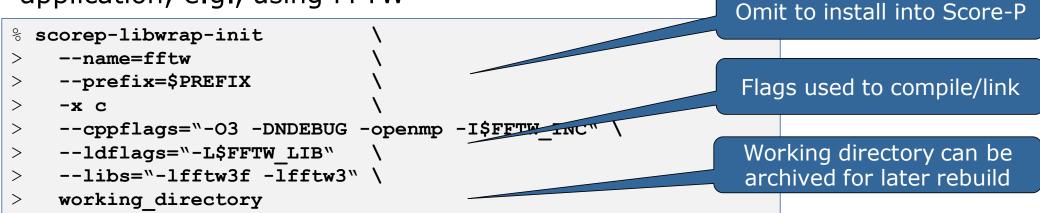
- Workflow to generate library wrappers for most C/C++ libraries
- Tailored towards user of the external library, not users of Score-P
- Results can be shared by multiple users
- Workflow driver scorep-libwrap-init --help provides instructions



Wrapping calls to 3rd party libraries: Workflow



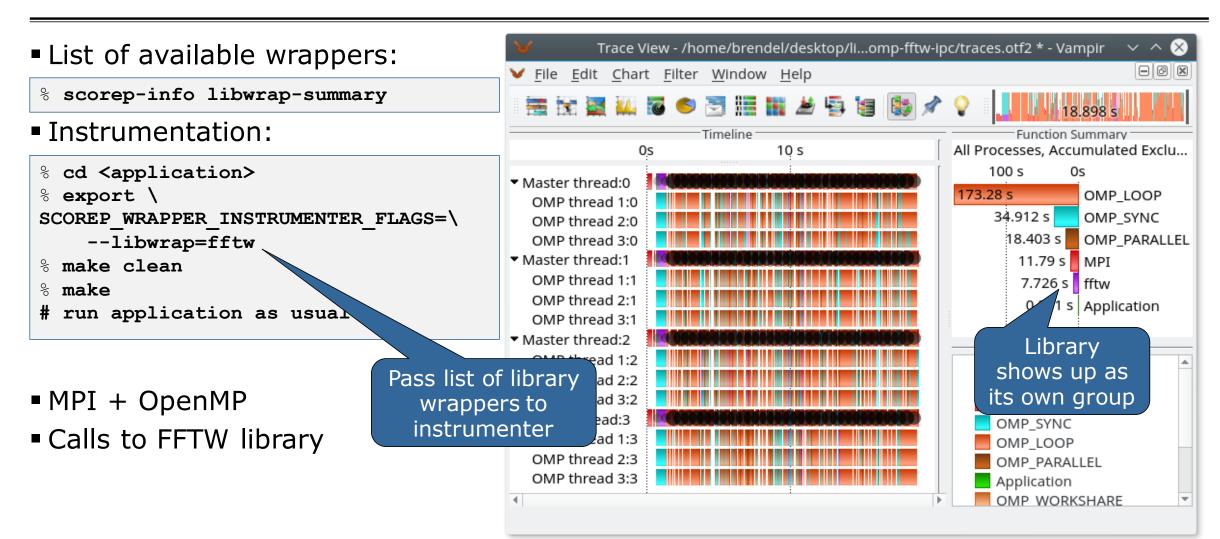
Start workflow by telling scorep-libwrap-init how you would compile and link an application, e.g., using FFTW



Generate and build wrapper

0/0	cd working directory		
			ells you how to use the
0/0	make #	Generate and build wrapper	wrapper with Score-P
00	make check #	See if header analysis matches symbols	
00	make install #		
010	<pre>make installcheck #</pre>	More checks: Linking etc.	

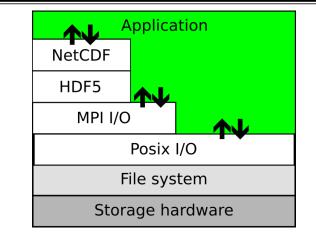
Wrapping calls to 3rd party libraries: Usage and result

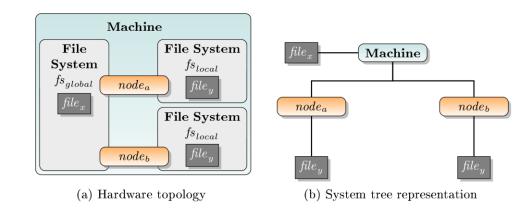


File I/O recording



- Omnipresent in todays HPC applications
- Record interaction between multiple layers
 - MPI I/O (MPI_File_open)
 - ISO C I/O (fopen)
 - POSIX I/O (open, interface to OS)
- System tree information determine whether file resides in a shared filesystem
- High level of detail
 => Trace data might increase dramatically





B_EFF I/O

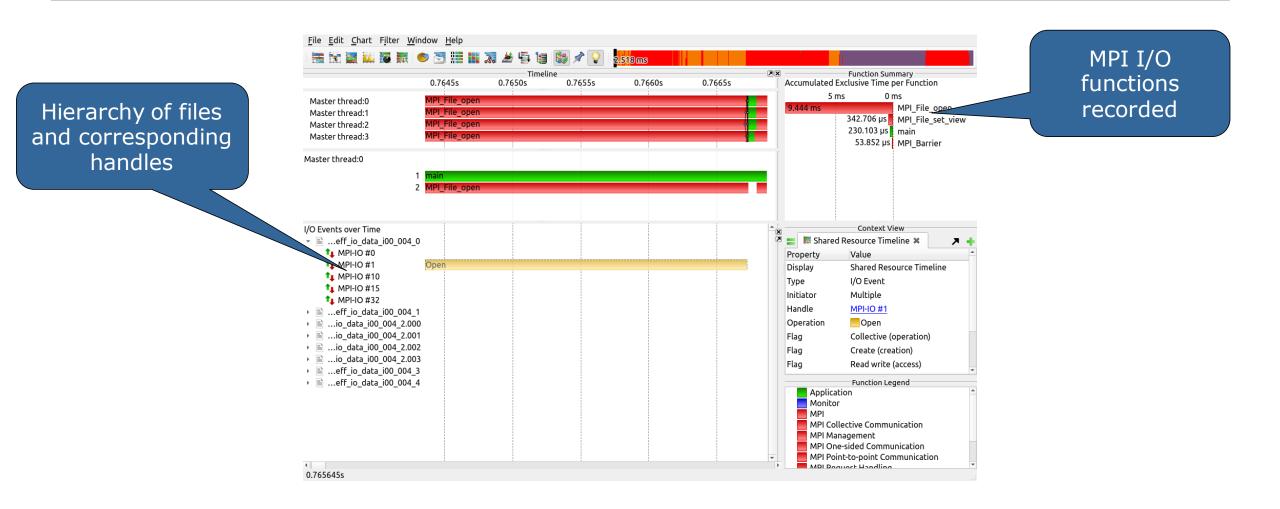
- MPI I/O benchmark
- <u>fs.hlrs.de/projects/par/mpi//b_eff_io/</u>

MPI I/O is enabled by default

```
% scorep-mpicc -o b_eff_io b_eff_io.c
% export SCOREP_EXPERIMENT_DIRECTORY=scorep-b_eff_io-4-profile
% mpirun -n 4 -c 6 ./b_eff_io -MB 2048 -MT 98304 -rewrite -N 4 -T 60
% scorep-scorep -g scorep-b_eff_io-4-profile/profile.cubex
% export SCOREP_EXPERIMENT_DIRECTORY=scorep-b_eff_io-4-tracing
% export SCOREP_FILTERING_FILE=initial_scorep.filter
% export SCOREP_ENABLE_TRACING=true
% export SCOREP_ENABLE_TRACING=true
% export SCOREP_TOTAL_MEMORY=31MB
% mpirun -n 4 -c 6 ./b_eff_io -MB 2048 -MT 98304 -rewrite -N 4 -T 60
```

VIRTUAL INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

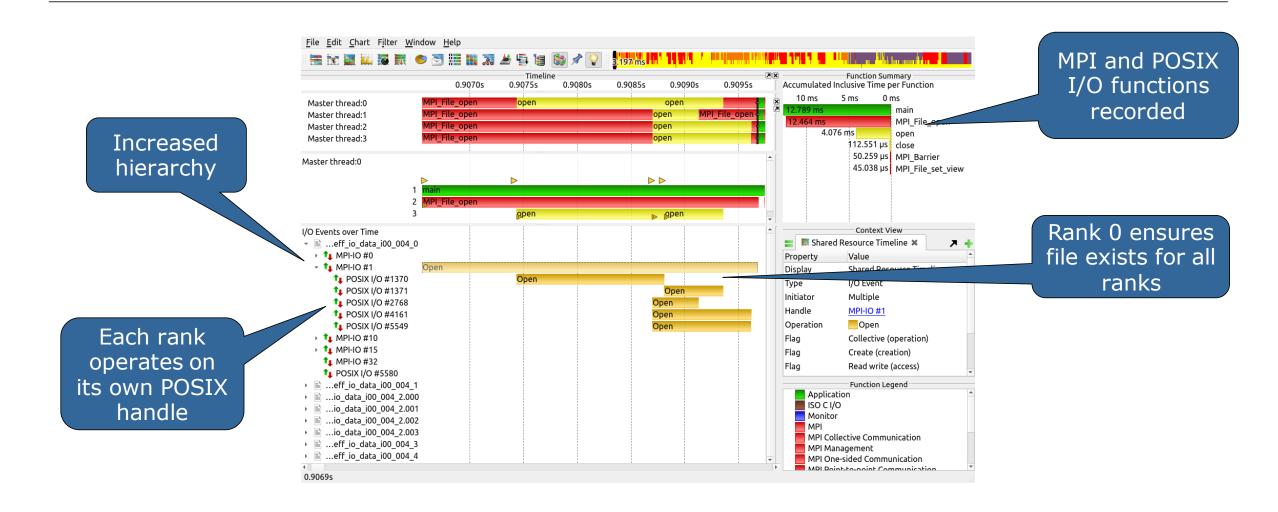
Result visualization



B_EFF I/O

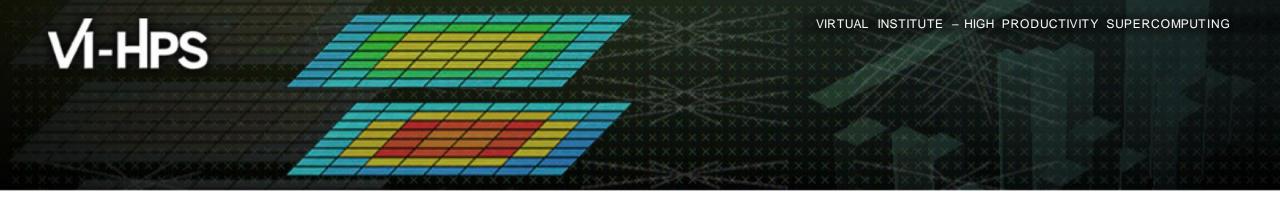
- Enabling ISO C and POSIX I/O instrumentation
- Instrumentation might require threading support

Result visualization



Further information

- Community instrumentation & measurement infrastructure
 - Instrumentation (various methods) and sampling
 - Basic and advanced profile generation
 - Event trace recording
 - Online access to profiling data
- Available under 3-clause BSD open-source license
- Documentation & Sources:
 - https://www.score-p.org
- User guide also part of installation:
 - <prefix>/share/doc/scorep/{pdf,html}/
- Support and feedback: support@score-p.org
- Subscribe to news@score-p.org, to be up to date



Exercise: Typical performance bottlenecks and how they can be identified

The VI-HPS Team



VIRTUAL VIRTUAL

Sparse matrix vector multiplication

$$\begin{pmatrix} y_1 \\ \vdots \\ y_m \end{pmatrix} = \begin{pmatrix} a_{11} & \cdots & a_{n1} \\ \vdots & \ddots & \vdots \\ a_{m1} & \cdots & a_{mn} \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

- A sparse matrix is a matrix populated primarily with zeros
- Only non-zero elements of a_{ij} are saved efficiently in memory

Algorithm

```
foreach row r in A
y[r.x] = 0
foreach non-zero element e in row
y[r.x] += e.value * x[e.y]
```

VIRTUAL×INSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Sparse matrix vector multiplication

Naive OpenMP algorithm

```
#pragma omp parallel for
foreach row r in A
 y[r.x] = 0
 foreach non-zero element e in row
 y[r.x] += e.value * x[e.y]
```

- Distributes the rows of A evenly across the threads in the parallel region
- The distribution of the non-zero elements may influence the load balance in the parallel application

Local installation (JUWELS Booster)

Set account and default environment (NVHPC + ParaStationMPI) via helper script:

% source /p/project/training2341/setup.sh

Load the modules for the tool environment:

% module load Score-P CubeGUI

Copy tutorial sources to your WORK directory (or your personal workspace)

Only required if not done already (for opening exercise)

[%] cd \$₩ORK

% tar xf \$PROJECT/examples/smxv.tar.gz

 $\frac{9}{6}$ cd SMXV

VI-HPS

SMxV: Makefile

```
% cat Makefile
CC=pqcc
OMP CFLAGS=-fopenmp
LIB\overline{S} = -lm
TARGETS=\
     smxv-omp
all: $(TARGETS)
clean: $(RM) *.o $(TARGETS) smxv-omp.scorep
CPPFLAGS=-DLITTLE ENDIAN -DITERLIMIT=500
smxv-omp: smxv.c y Ax.h
         $(CC) $(OMP CFLAGS) $(CPPFLAGS) $(CFLAGS) -0 $@ $< $(LIBS)
scorep: smxv-omp.scorep
smxv-omp.scorep: smxv.c y Ax.h
         scorep $(CC) $(OMP CFLAGS) $(CPPFLAGS) $(CFLAGS) -o $@ $< $(LIBS)</pre>
```

Load imbalances in OpenMP codes

• Exercise I: Determine scaling behavior for number of threads between 1 and 12

% make
% sbatch -c 1 run.sbatch
% sbatch -c 2 run.sbatch
...
% sbatch -c 12 run.sbatch

Load imbalances in OpenMP codes

• Exercise II: Collection and analyze measurement

% make scorep % sbatch -c 12 scorep.sbatch

Load imbalances in OpenMP codes: Profile examination

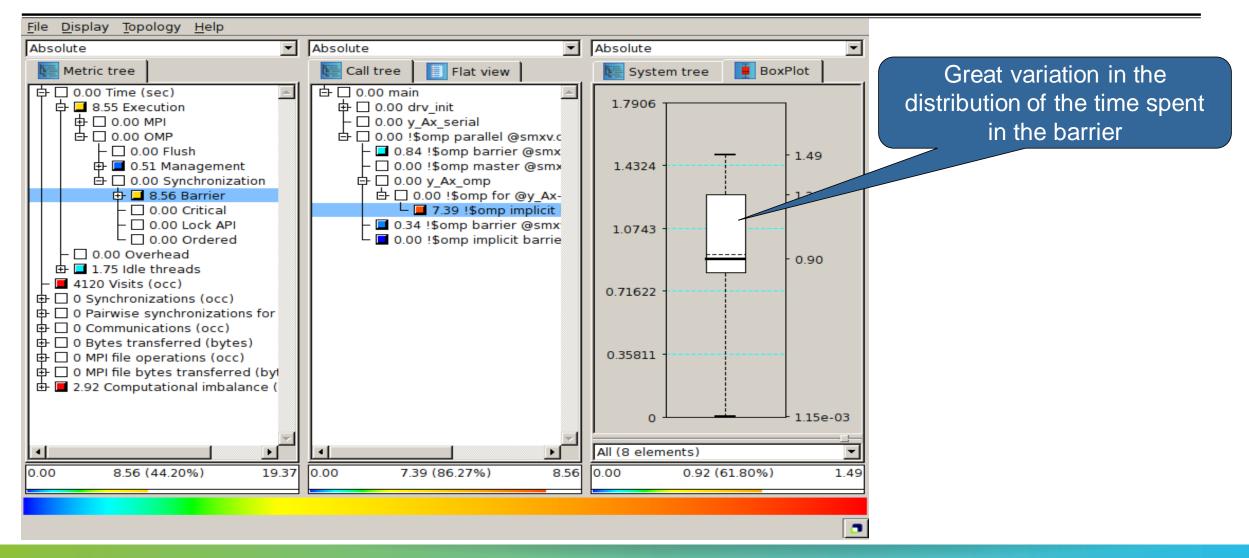
- Two metrics which indicate load imbalances:
 - Time spent in OpenMP barriers
 - Computational imbalance

• Open measurement with Cube

% cube scorep-smxv-12/trace.cubex

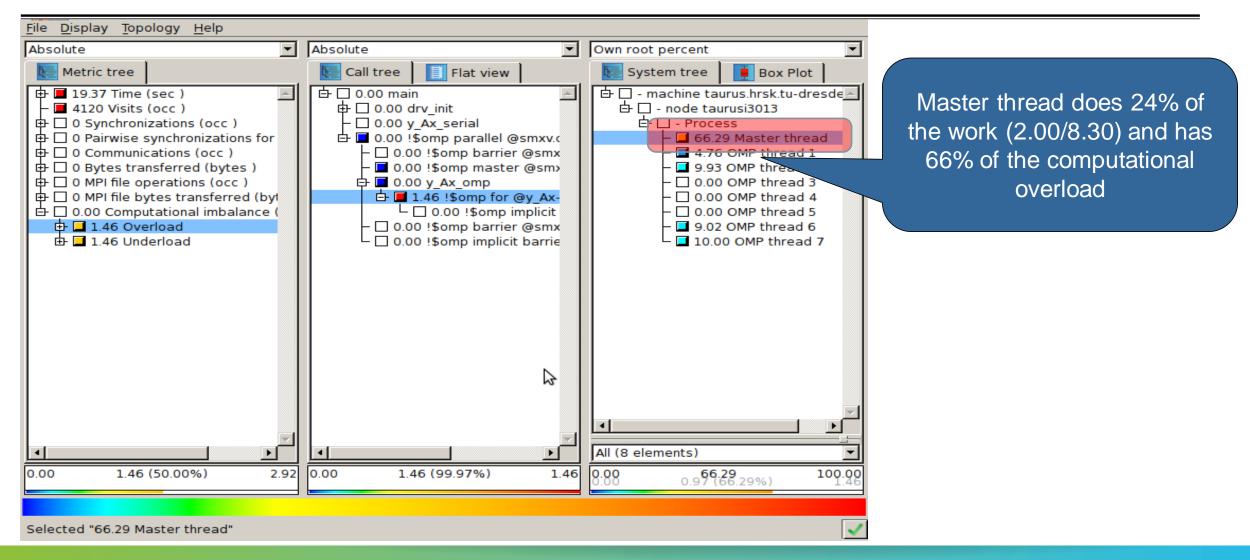
[CUBE GUI showing trace analysis report]

Time spent in OpenMP barriers



VICTOR VICT

Computational imbalance



Sparse matrix vector multiplication

- Scheduling policy determined via environmental variable OMP_SCHEDULE:
 - 1. static[,n]
 - Schedule is fix at start of work share
 - 2. dynamic[,n]
 - Work is split into packages of n-items. Each thread requests new package

Exercise III: Measure improved algorithm

```
% cat scorep.sbatch
...
export SCOREP_EXPERIMENT_DIRECTORY=scorep-smxv-$SLURM_CPUS_PER_TASK-dynamic
export OMP_SCHEDULE=...
scan -t srun ./smxv-omp.scorep yax_large.bin
% sbatch -c 12 scorep.sbatch
```

Profile Analysis

- Two metrics which indicate load imbalances
 - Time spent in OpenMP barriers
 - Computational imbalance

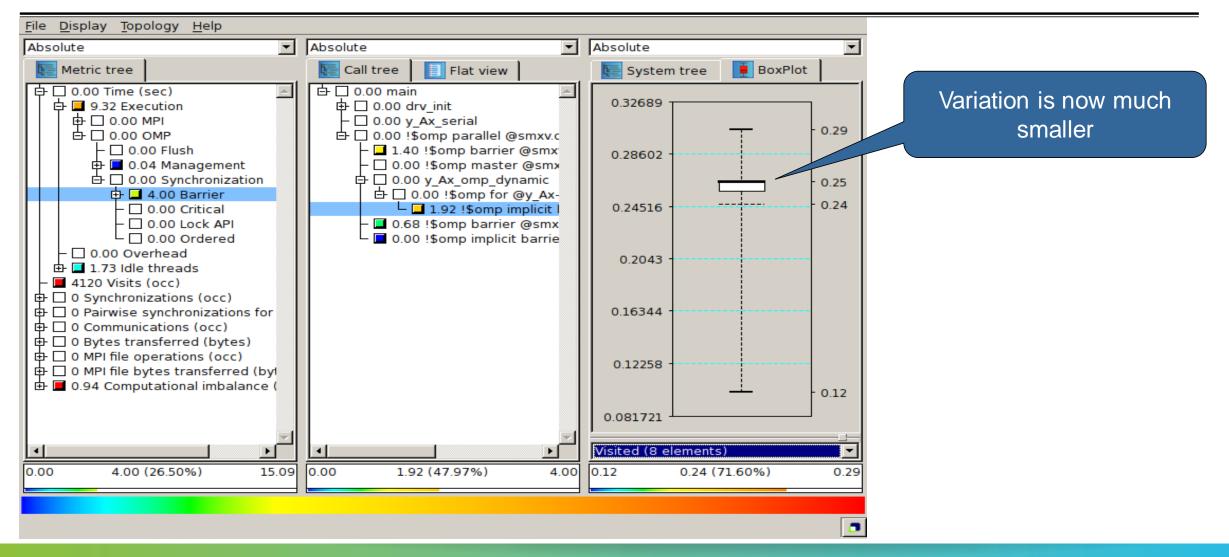
Open measurement with Cube

% cube scorep-smxv-12-dynamic/trace.cubex

[CUBE GUI showing trace analysis report]

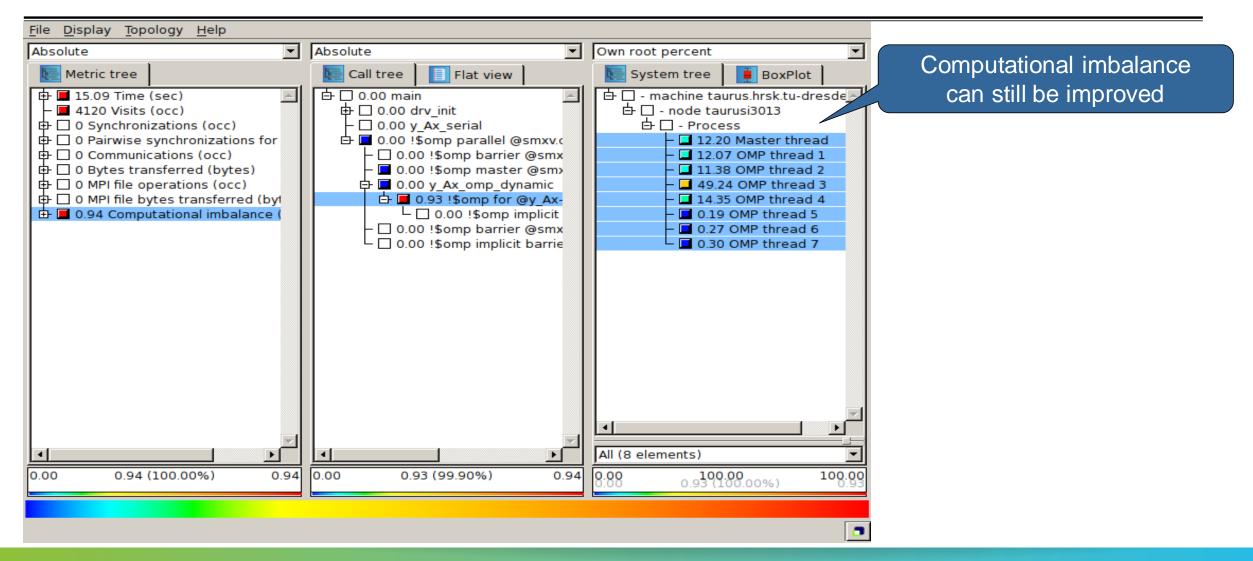
VIRTUAL VINSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Time spent in OpenMP barriers



VICTOR VICT

Computational imbalance





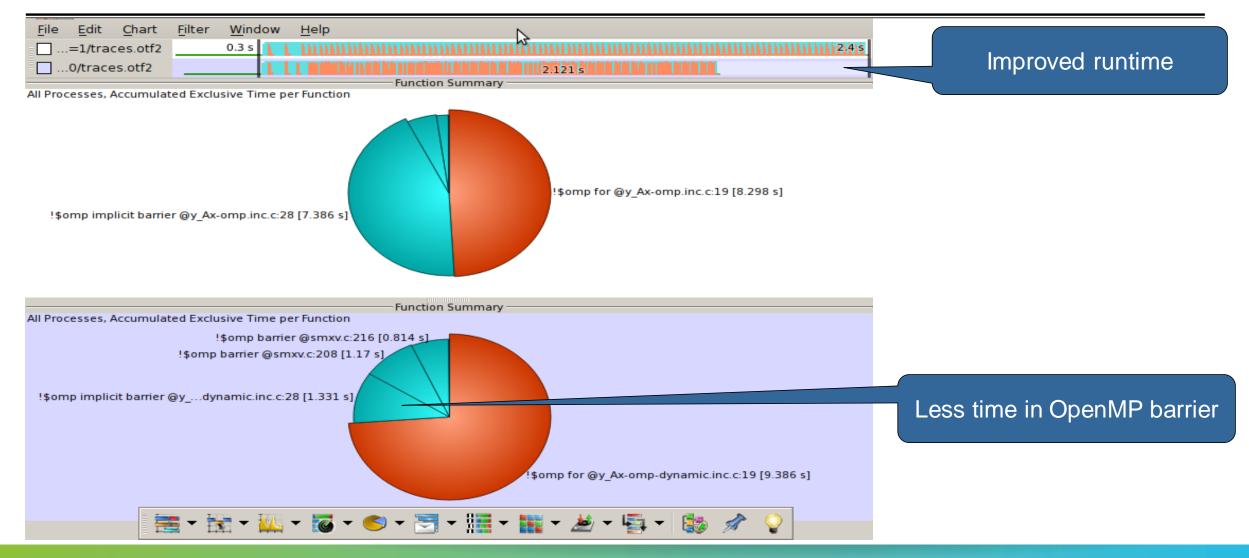
VIRTUAL VINSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Trace comparison

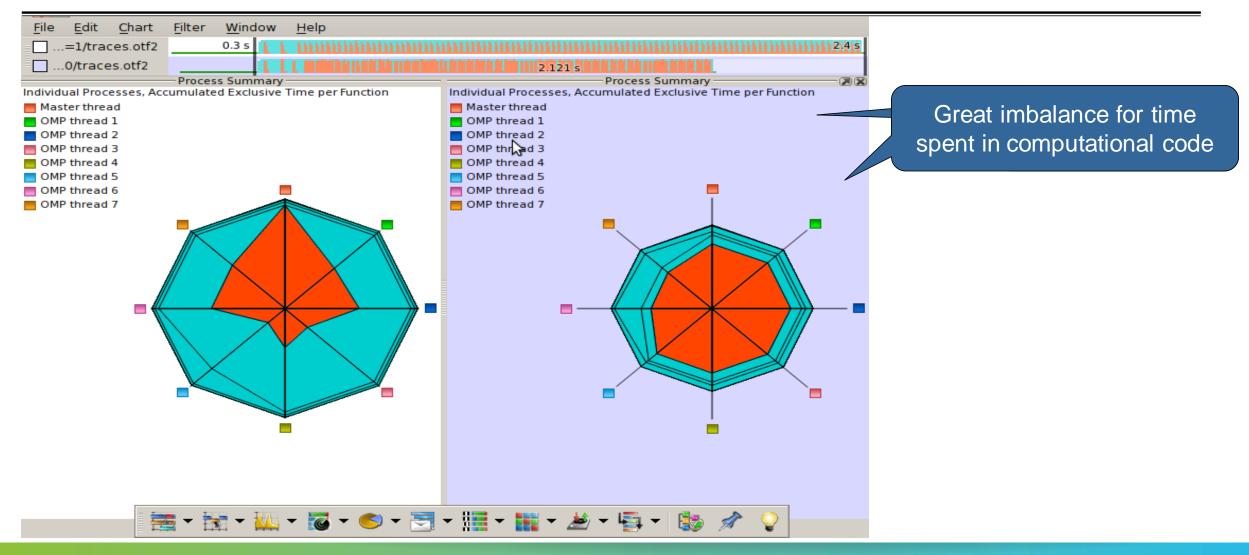
Open both traces in Vampir

VIRTUAL VIRTUAL VINSTITUTE - HIGH PRODUCTIVITY SUPERCOMPUTING

Time spent in OpenMP barriers



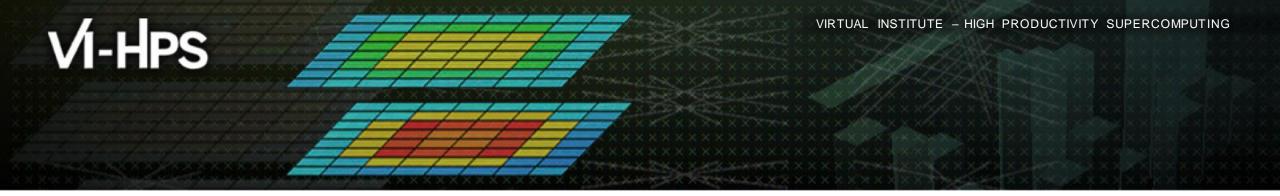
Computational imbalance



Validate optimization

Exercise IV

```
% cat run.sbatch
...
export OMP_SCHEDULE=...
srun ./smxv-omp yax_large.bin
% sbatch -c 12 run.sbatch
```



Review

Markus Geimer Jülich Supercomputing Centre



Summary

You've been introduced to a variety of toolswith hints to apply and use the tools effectively

Tools provide complementary capabilities
computational kernel & processor analyses
communication/synchronization analyses

Ioad-balance, scheduling, scaling, ...

Tools are designed with various trade-offs

- general-purpose versus specialized
- platform-specific versus agnostic
- simple/basic versus complex/powerful

Tool selection

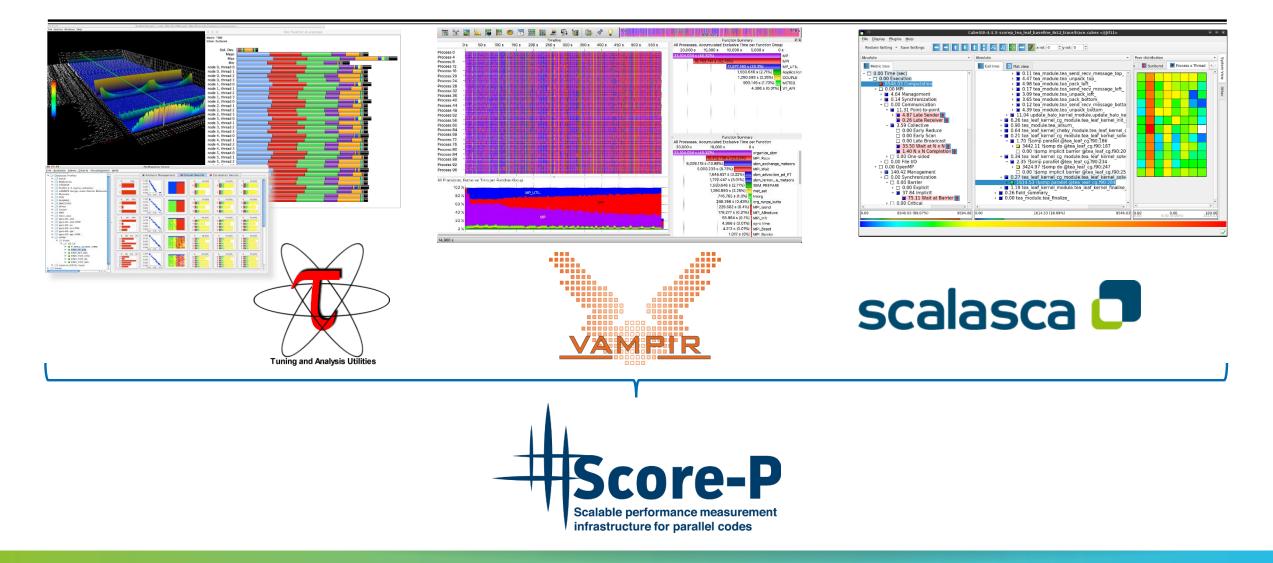
Which tools you use and when you use them likely to depend on the situation

- which are available on (or for) your computer system
- which support your programming paradigms and languages
- which you are familiar (comfortable) with using
- which type of issue you suspect
- which question you want to have answered

Being aware of (potentially) available tools and their capabilities can help finding the most appropriate tools

Portable multi-platforms tools complement platform-specific/vendor tools

Tools featured in this tutorial



Workflow (getting started)

First ensure that the parallel application runs correctly

- no-one will care how quickly you can get invalid answers or produce a set of corefiles
- parallel debuggers help isolate known problems
- correctness checking tools can identify other issues
 - (that might not cause problems right now, but will eventually)
 - e.g., race conditions, invalid/non-compliant usage

Best to start with an overview of execution performance

- fraction of time spent in computation (CPU & GPU) vs comm/synch vs I/O
- which sections of the application/library code are most costly
- Example profilers: Score-P + Cube/ParaProf, TAU

and how it changes with scale or different configurationsprocesses vs threads, mappings, bindings

Workflow (communication/synchronization)

Communication issues generally apply to every computer system (to different extents) and typically grow with the number of processes/threads

- Weak scaling: fixed computation per thread, and perhaps fixed localities, but increasingly distributed
- Strong scaling: constant total computation, increasingly divided amongst threads, while communication grows
- Collective communication (particularly of type "all-to-all") result in increasing data movement
- Synchronizations of larger groups are increasingly costly
- Load-balancing becomes increasingly challenging, and imbalances more expensive
 - generally manifests as waiting time at following collective ops

Workflow (wasted time waiting)

Waiting times are difficult to determine in basic profiles

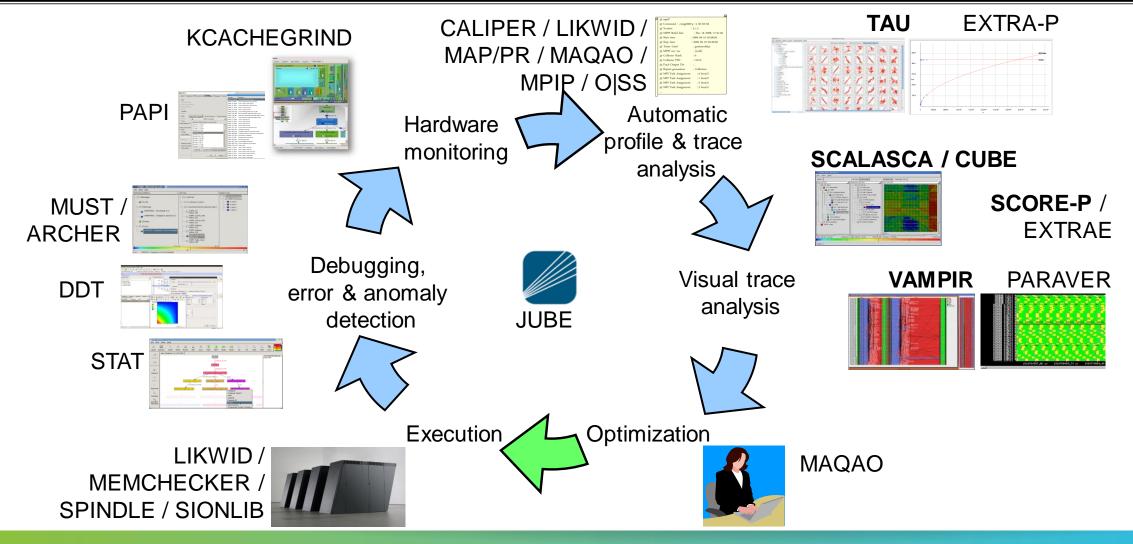
- Part of the time each process/thread spends in communication & synchronization operations may be wasted waiting time
- Need to correlate event times between processes/threads
 - Post-mortem event trace analysis avoids interference and provides a complete history
 - Scalasca automates trace analysis and ensures waiting times are completely quantified
 - Vampir allows interactive exploration and detailed examination of reasons for inefficiencies

Workflow (core computation)

Effective computation within processors/cores (and GPUs) is also vital

- Optimized libraries may already be available
- Optimizing compilers can also do a lot
 - provided the code is clearly written and not too complex
 - appropriate directives and other hints can also help
- Processor hardware counters can also provide insight
 - although hardware-specific interpretation required
- Tools available from processor and system vendors help navigate and interpret processor-specific performance issues

Technologies and their integration



SC23 TUTORIAL: HANDS-ON PRACTICAL HYBRID PARALLEL APPLICATION PERFORMANCE ENGINEERING (DENVER, 13 NOV 2023)

Further information

Website

- Introductory information about the VI-HPS portfolio of tools for high-productivity parallel application development
 - VI-HPS Tools Guide
 - links to individual tools sites for details and download
- Training material
 - tutorial slides
 - user guides and reference manuals for tools
- News of upcoming events
 - tutorials and workshops
 - mailing-list sign-up for announcements

https://www.vi-hps.org