



# HPC SOFTWARE – DEBUGGER AND PERFORMANCE ANALYSIS TOOLS

NOVEMBER 22, 2023 | MICHAEL KNOBLOCH | M.KNOBLOCH@FZ-JUELICH.DE



# OUTLINE

- Local module setup
- Compilers
- Libraries

Make it work,  
make it right,  
make it fast.

*Kent Beck*

**Debugger and Correctness  
Tools**

**Performance Analysis Tools**

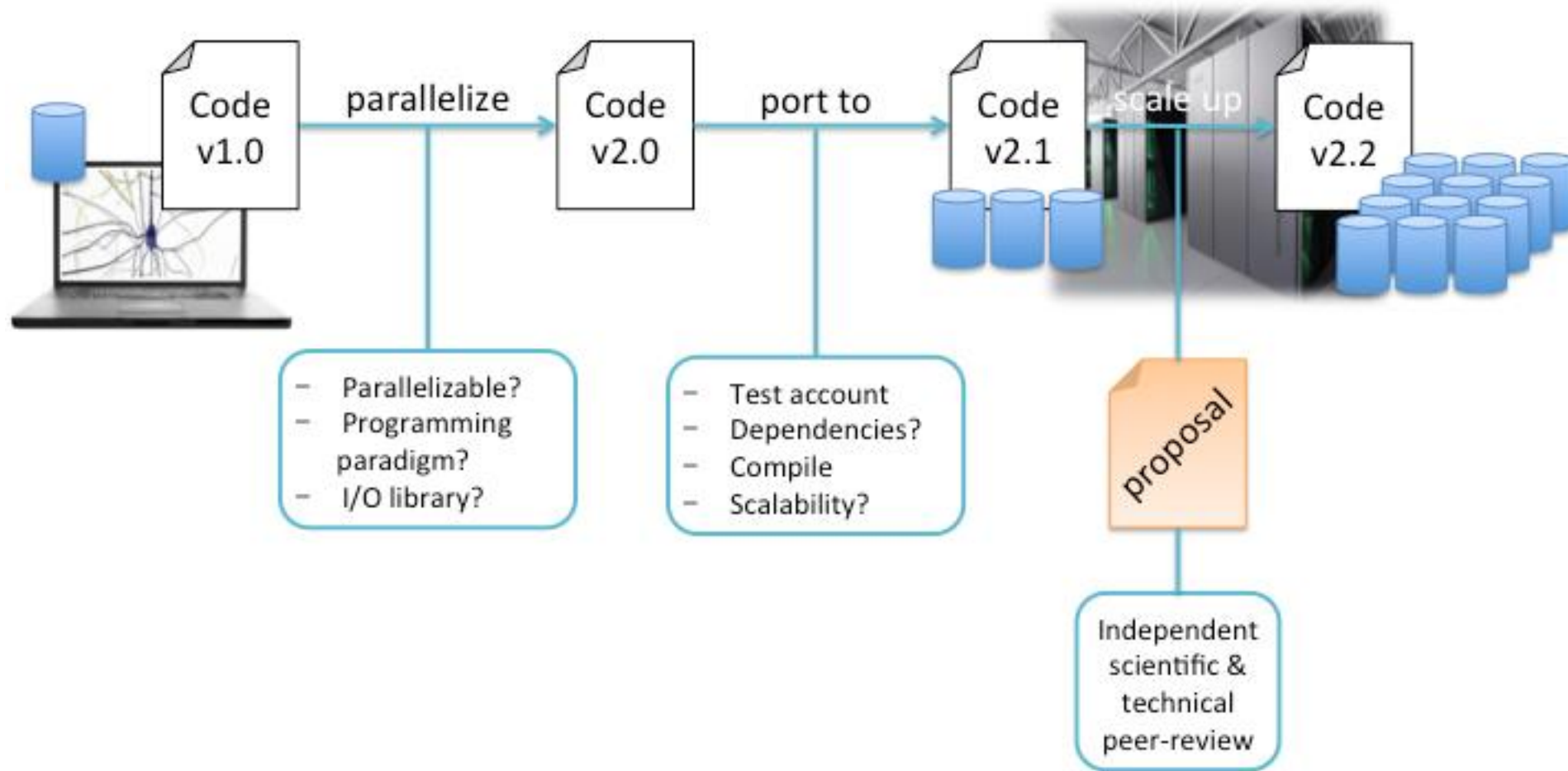


# WHY SHOULD YOU CARE ABOUT TOOLS?





# NEW APPLICATION?





# WORKING WITH LEGACY CODES?





# VETERAN HPC USER, BUT NEW TO JSC?



- Assess performance on a JSC machine



- Compare behavior on different machines



- Investigate scaling behavior





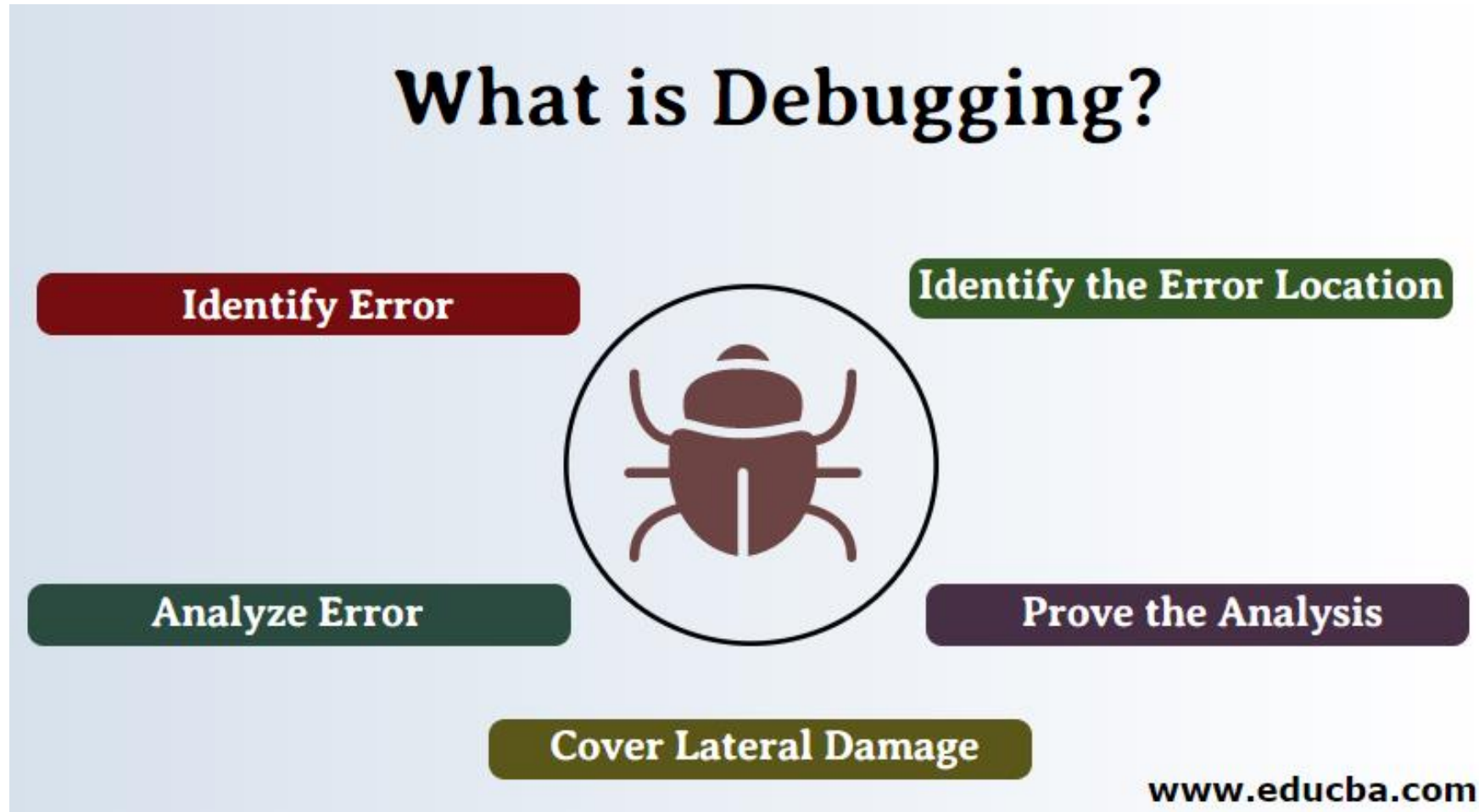
The screenshot displays a debugger interface with four main panels:

- Thread List:** Shows a list of threads. Thread 1.1 is highlighted as the 'Breakpoint' thread, with state 'Stopped' and function 'tensorflow::SoftmaxXentWithLosses'.
- Variable View:** Displays variables for the selected thread. Variables include 'nstar' (int, 0x00000006) and 'grap...' (int, 0x00000015).
- Source Code:** Shows C++ code from 'tensorflow::TF\_Run\_Setup'. Line 619 is highlighted, showing a call to 'TF\_Run\_Helper'.
- Call Stack:** Lists the sequence of function calls. The top frame is 'tensorflow::TF\_Run', which called 'tensorflow::TF\_Run\_wrapper', which in turn called 'tensorflow::TF\_Run\_Helper'.

# DEBUGGER & CORRECTNESS TOOLS



# WHAT IS DEBUGGING?





# REMINDER: DEBUGGING CAN BE FRUSTRATING





# DEBUGGING TOOLS (**STATUS: NOV 2023**)

- **Debugger:**

- CUDA-GDB
- TotalView
- LinaroForge - DDT

- **Memory Analyzer:**

- Intel Inspector

- **Correctness Checker:**

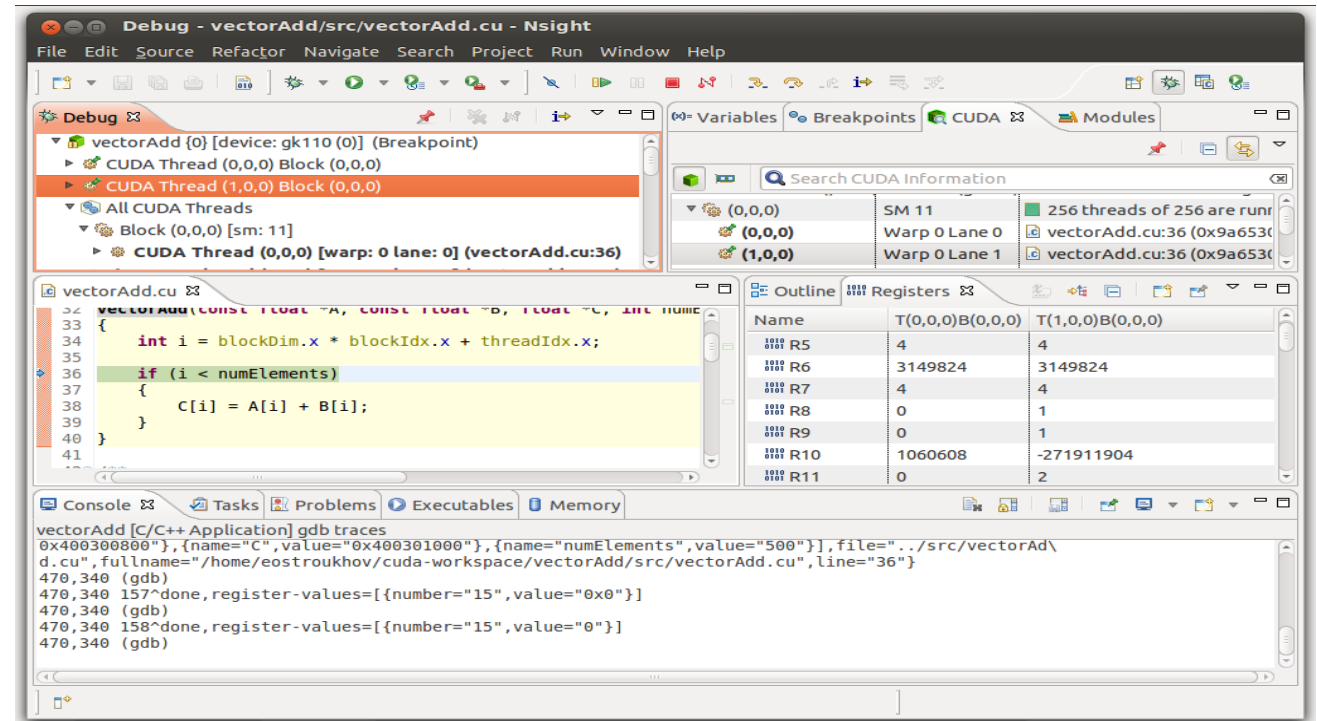
- MUST



# CUDA-GDB



- Part of the CUDA toolkit
- Extension to gdb
- CLI and GUI (Nsight)
- Simultaneously debug on the CPU and multiple GPUs
- Use conditional breakpoints or break automatically on every kernel launch
- Examine variables, read/write memory and registers
- Inspect GPU state when the application is suspended
- Identify memory access violations





- UNIX Symbolic Debugger for C/C++, Fortran, mixed Python/C++, PGI HPF, assembler programs
- JSC's “standard” debugger
- Advanced features
  - Multi-process and multi-threaded
  - Multi-dimensional array data visualization
  - Support for [parallel debugging](#) (MPI: automatic attach, message queues, OpenMP, Pthreads)
  - Scripting and [batch debugging](#)
  - Advanced memory debugging
  - Reverse debugging
  - [CUDA](#) and [OpenACC](#) support
  - Remote debugging
- **NOTE:** JSC license limited to 2048 processes (shared between all users)



# TOTALVIEW: MAIN WINDOW

The screenshot shows the TOTALVIEW IDE interface with several callouts highlighting key features:

- Thread control:** A callout points to the 'Processes & Threads' panel on the left, which lists the execution state of various threads.
- Break points:** A callout points to the 'Action Points' panel at the bottom left, showing a breakpoint set at line 91 of the `tx_cuda_matmul` process.
- Source code window:** A callout points to the central editor displaying the C++ source code for `MatMulKernel`.
- Stack trace:** A callout points to the 'Call Stack' panel on the right, showing the current stack frame.
- Local variables for selected stack frame:** A callout points to the 'Data View' panel at the bottom right, which displays the values of local variables like `A`, `width`, `height`, `stride`, and `elements`.
- Toolbar for common options:** A callout points to the top toolbar, which contains icons for running, stepping, and other debugging actions.



# LINARO FORGE - DDT

- UNIX Graphical Debugger for C/C++, Fortran, and Python programs
- Modern, easy-to-use debugger
- Advanced features
  - Multi-process and multi-threaded
  - Multi-dimesional array data visualization
  - Support for MPI parallel debugging (automatic attach, message queues)
  - Support for OpenMP (Version 2.x and later)
  - Support for CUDA and OpenACC
  - Job submission from within debugger
- <https://linaroforge.com/linaroDdt>
- **NOTE:** JSC license limited to 128 processes (shared between all users)



# DDT: MAIN WINDOW

Process controls

CUDA Thread stepping

Variables

CUDA Thread control

Source code

GPU Device information

Expression evaluator

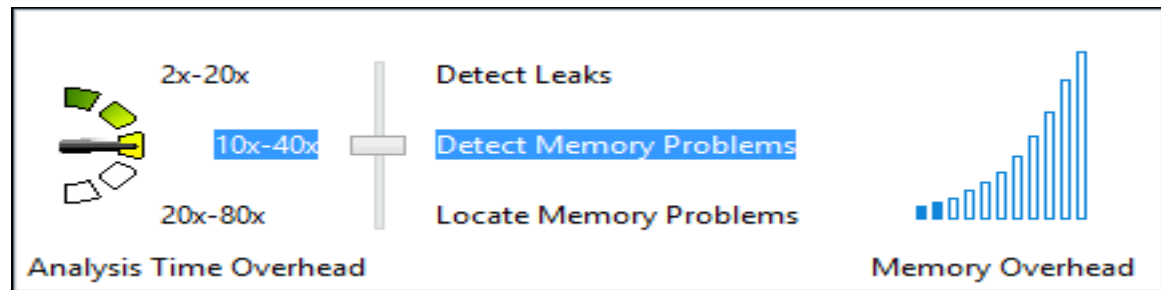
Stack trace

The screenshot displays the Arm DDT - Arm Forge 19.1.1 interface. At the top, there are process controls including 'Focus on current: Process', 'Thread', 'Step Threads Together', and 'Step CUDA threads by: Warp (default)'. Below this is the 'Threads' panel showing 'CUDA Threads (conv2d\_global)' with a grid size of 32x32x1 and block size of 16x16x1. The 'Project Files' panel on the left shows the source code structure. The central 'Source code' panel displays the C++ code for a 2D convolution filter. The 'Locals' panel on the right shows variables like 'cx', 'cy', 'output', 'x', 'y', and 'index'. The 'GPU Devices' panel shows information about the GM20B device, including Compute Capability, Number of SMs, Warps per SM, Lanes per Warp, and Registers per Lane. The 'Stacks' panel at the bottom shows the current stack trace, and the 'Expression evaluator' panel shows the evaluation of expressions like 'x+cx + (y+cy)\*width'.



# INTEL INSPECTOR

- Detects memory and threading errors
  - Memory leaks, corruption and illegal accesses
  - Data races and deadlocks
- Dynamic instrumentation requiring no recompilation
- Supports C/C++ and Fortran as well as third party libraries
- Multi-level analysis to adjust overhead and analysis capabilities
- API to limit analysis range to eliminate false positives and speed-up analysis





# INTEL INSPECTOR: GUI

The screenshot shows the Intel Inspector XE 2015 interface for detecting deadlocks and data races. The 'Problems' table lists three data race issues (P1, P2, P3) involving files like `find_and_fix_threading_errors.cpp` and `task_scheduler_init.h`. The 'Code Locations' pane shows the source code for `render_one_pixel` in `find_and_fix_threading_errors.cpp`, highlighting a write operation on `primary.scene` and a comment about a threading error. The 'Filters' pane on the right shows the severity and type of the detected issues.

ID	Type	Sources	Modules	State
P1	Data race	<code>find_and_fix_threading_errors.cpp</code> ; <code>task_scheduler_init.h</code>	<code>find_and_fix_threading_errors.exe</code>	New
P2	Data race	<code>blocked_range.h</code> ; <code>parallel_for.h</code> ; <code>partitioner.h</code> ; <code>task.h</code>	<code>find_and_fix_threading_errors.exe</code>	New
P3	Data race	<code>wmvideo.h</code>	<code>find_and_fix_threading_errors.exe</code>	New

Description	Source	Function	Module	Variable
Write	<code>find_and_fix_threading_errors.cpp:105</code>	<code>render_one_pixel</code>	<code>find_and_fix_threading_errors.exe</code>	<code>color</code>

The screenshot shows the Intel Inspector XE 2015 interface for detecting memory problems. The 'Problems' table lists four memory issues (P1, P2, P3, P4), including mismatched allocation/deallocation, invalid memory access, and memory growth. The 'Code Locations' pane shows the source code for `operator()` in `find_and_fix_memory_error...`, highlighting a loop that initializes `local_mbox` and a comment about a memory error. The 'Filters' pane on the right shows the severity and type of the detected issues.

ID	Type	Sources	State
P1	Mismatched allocation/deallocation	<code>find_and_fix_memory_errors...</code>	New
P2	Invalid memory access	<code>find_and_fix_memory_errors...</code>	New
P3	Memory growth	[Unknown]; <code>find_and_fix_me...</code>	Not fixed
P4	Memory growth	[Unknown]; <code>find_and_fix_me...</code>	Confirmed

Description	Source	Function	Module	Object Size	Offset
Write	<code>find_and_fix_memory_error...</code>	<code>operator()</code>	<code>find_and_fix_memory_error...</code>		

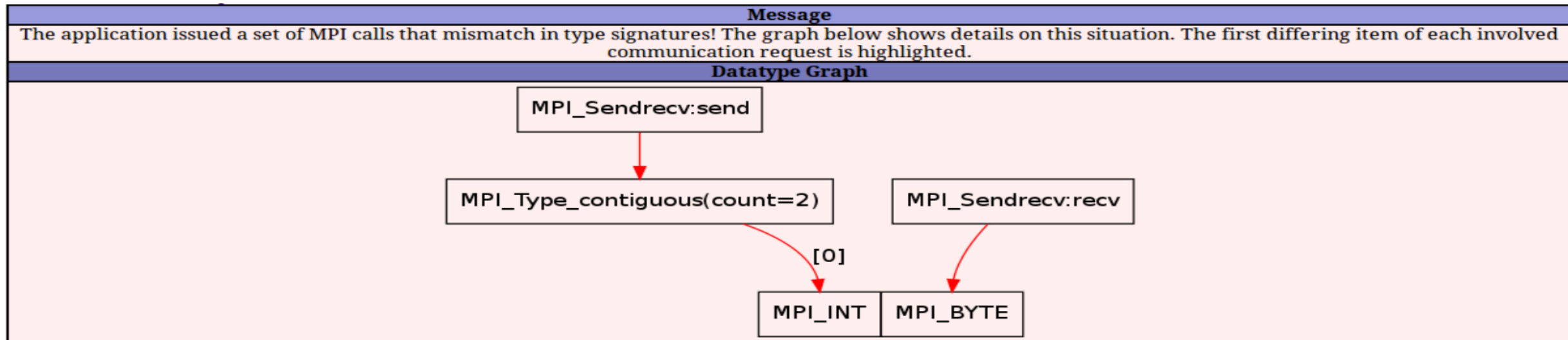


- Next generation MPI correctness and portability checker
- <https://www.i12.rwth-aachen.de/go/id/nrbe>
- MUST reports
  - Errors: violations of the MPI-standard
  - Warnings: unusual behavior or possible problems
  - Notes: harmless but remarkable behavior
  - Potential deadlock detection
- Usage
  - Relink application with `mustc`, `mustcxx`, `mustf90`, ...
  - Run application under the control of `mustrun` (requires (at least) one additional MPI process)
  - Saves output in html report



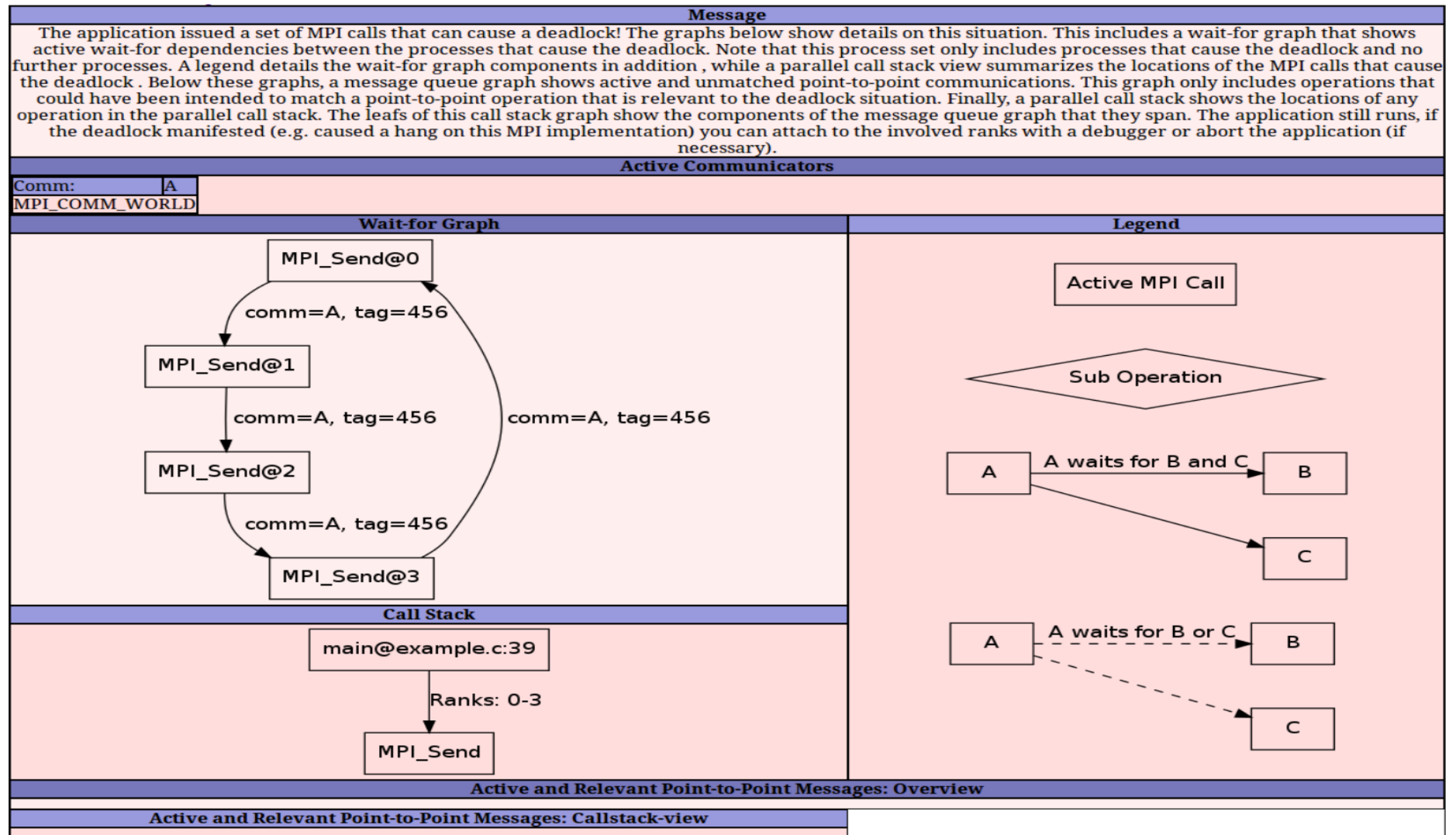
# MUST DATATYPE MISMATCH

Rank	Type	Message	From	References
0	Error	<p>A send and a receive operation use datatypes that do not match! Mismatch occurs at (contiguous) [0](MPI_INT) in the send type and at (MPI_BYTE) in the receive type (consult the MUST manual for a detailed description of datatype positions). A graphical representation of this situation is available in a <a href="#">detailed type mismatch view (MUST Output-files/MUST Typemismatch 0.html)</a>. The send operation was started at reference 1, the receive operation was started at reference 2. (Information on communicator: MPI_COMM_WORLD) (Information on send of count 1 with type:Datatype created at reference 3 is for C, committed at reference 4, based on the following type(s): { MPI_INT}Typemap = {(MPI_INT, 0), (MPI_INT, 4)}) (Information on receive of count 8 with type:MPI_BYTE)</p>	<p><b>MPI_Sendrecv</b> called from: #0 main@example.c:33</p>	<p>reference 1 rank 0: <b>MPI_Sendrecv</b> called from: #0 main@example.c:33</p> <p>reference 2 rank 1: <b>MPI_Sendrecv</b> called from: #0 main@example.c:33</p> <p>reference 3 rank 0: <b>MPI_Type_contiguous</b> called from: #0 main@example.c:29</p> <p>reference 4 rank 0: <b>MPI_Type_commit</b> called from: #0 main@example.c:30</p>





# MUST DEADLOCK DETECTION

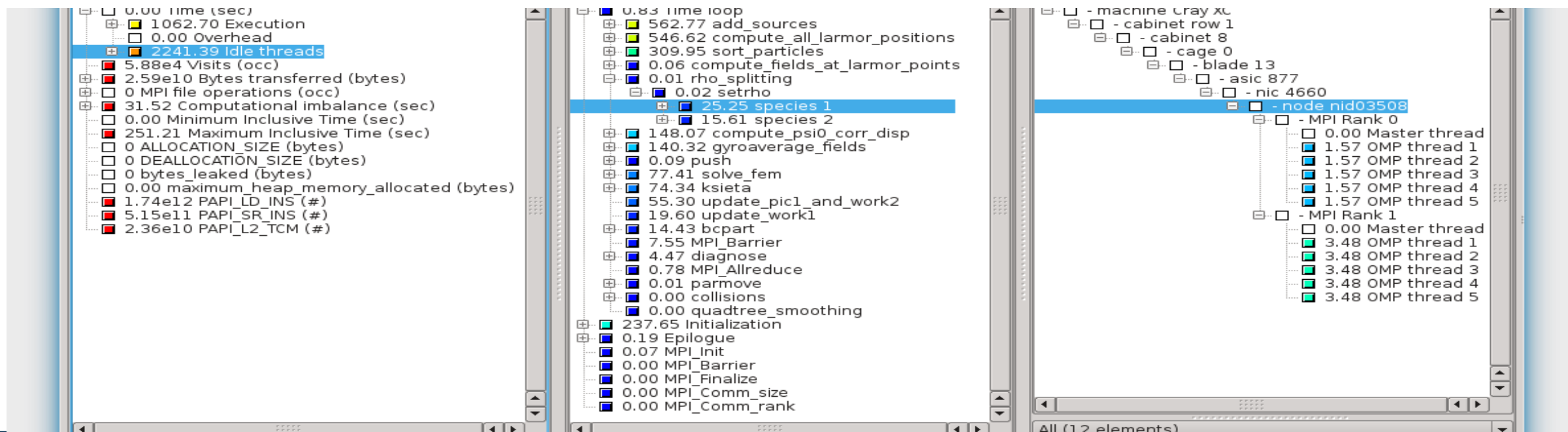




# DEBUGGING RECOMMENDATIONS

- Always debug at the lowest possible scale!
- GPU Applications:
  - Single Node / Workstation: Use CUDA-GDB
  - Multi-Node / Supercomputer: Use TotalView/DDT
- MPI Applications:
  - Check with MUST at least once
  - Use TotalView/DDT at small scale (if error occurs there), else attach to as few processes as necessary



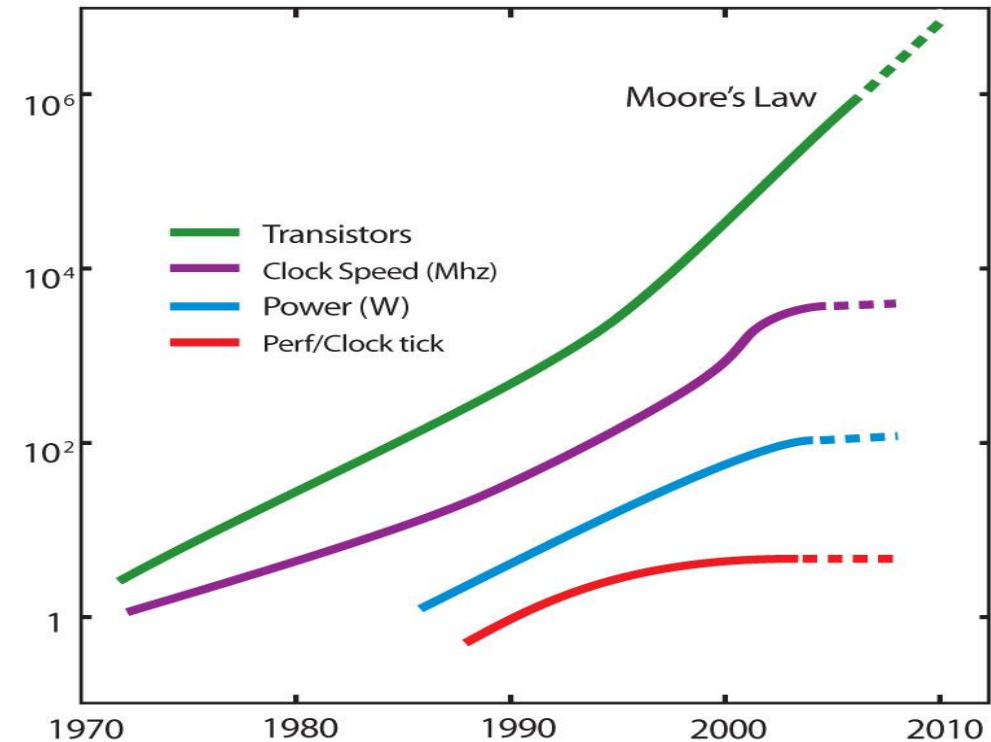


# PERFORMANCE ANALYSIS TOOLS



# 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
  - Optimization of applications more difficult
    - Increasing application complexity
      - Multi-physics
      - Multi-scale
    - Increasing machine complexity
      - Hierarchical networks / memory
      - Many-core CPUs and Accelerators
      - Modular Supercomputing Architecture
- ☞ Every doubling of scale reveals a new bottleneck!





# PERFORMANCE FACTORS

- “Sequential” (single core) factors
  - Computation
    - ☞ Choose right algorithm, use optimizing compiler
  - Vectorization
    - ☞ Choose right algorithm, use optimizing compiler
  - Cache and memory
    - ☞ Choose the right data structures and data layout



# PERFORMANCE FACTORS

- “Parallel” (multi core/node) factors
  - Partitioning / decomposition
    - ☞ Load balancing
  - Communication (i.e., message passing)
  - Multithreading
  - Core binding / NUMA
  - Synchronization / locking
  - I/O
    - ☞ Often not given enough attention
    - ☞ Parallel I/O matters

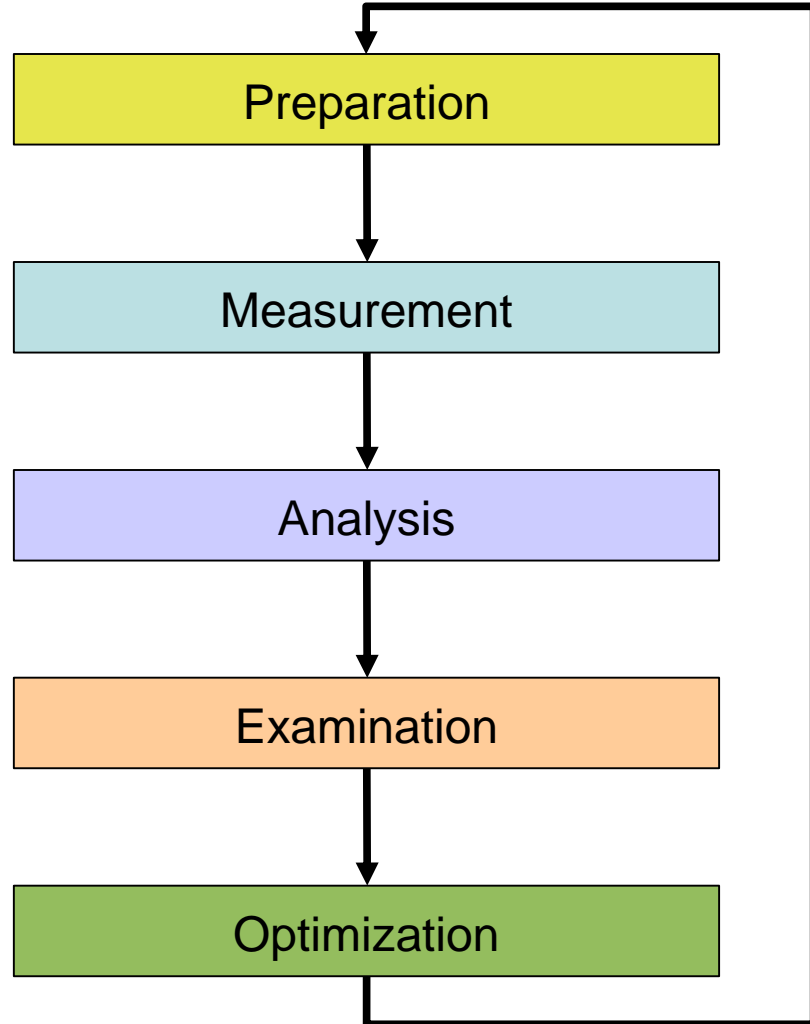


# TUNING BASICS

- Carefully set various tuning parameters
  - The right (parallel) algorithms and libraries
  - Compiler flags and directives
  - Correct machine usage (mapping and bindings)
    - 👉 Get the most performance before tuning!
- Measurement is better than guessing
  - To determine performance bottlenecks
  - To compare alternatives
  - To validate tuning decisions and optimizations
    - 👉 After each step!



# PERFORMANCE ENGINEERING WORKFLOW



- Prepare application (with symbols), insert extra code (probes/hooks)
- Collection of data relevant to execution performance analysis
- Calculation of metrics, identification of performance metrics
- Presentation of results in an intuitive/understandable form
- Modifications intended to eliminate/reduce performance problems



# THE 80/20 RULE

- Programs typically spend 80% of their time in 20% of the code

☞ *Know what matters!*

- Developers 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!*



# PERFORMANCE MEASUREMENT

## Two dimensions

**When** performance measurement is triggered

- **External trigger** (asynchronous)
  - **Sampling**
    - Trigger: Timer interrupt OR Hardware counters overflow
- **Internal trigger** (synchronous)
  - Code **instrumentation** (automatic or manual)

**How** performance data is recorded

- **Profile**
  - Summation of events over time
- **Trace**
  - Sequence of events over time



# MEASUREMENT METHODS: PROFILING

- Recording of **aggregated information**
  - Time
  - Counts
    - Calls
    - Hardware counters
- **Across program and system entities**
  - Functions, call sites, loops, basic blocks, ...
  - Processes, threads
- **Statistical information**
  - Min, max, mean and total number of values

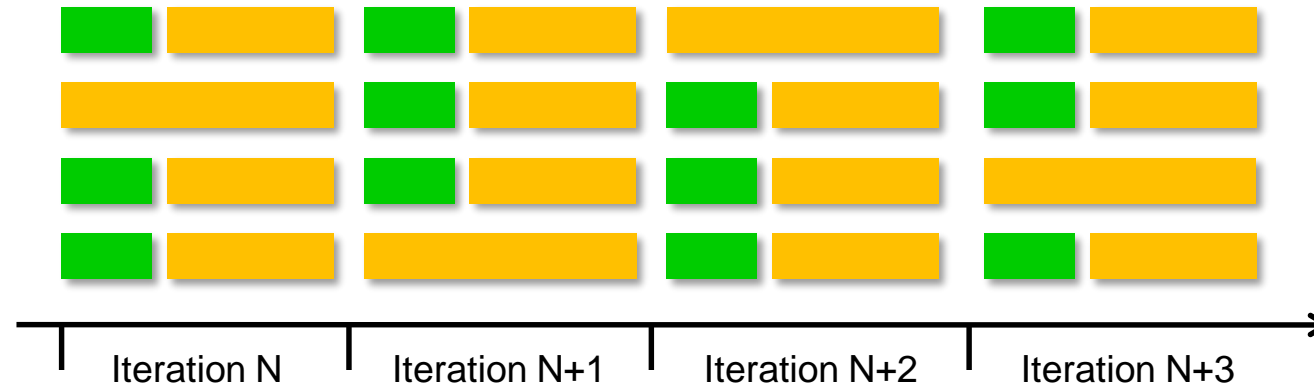
**Advantages**  
+ Works also for  
long-running programs

**Disadvantages**  
– Variations over time  
get lost

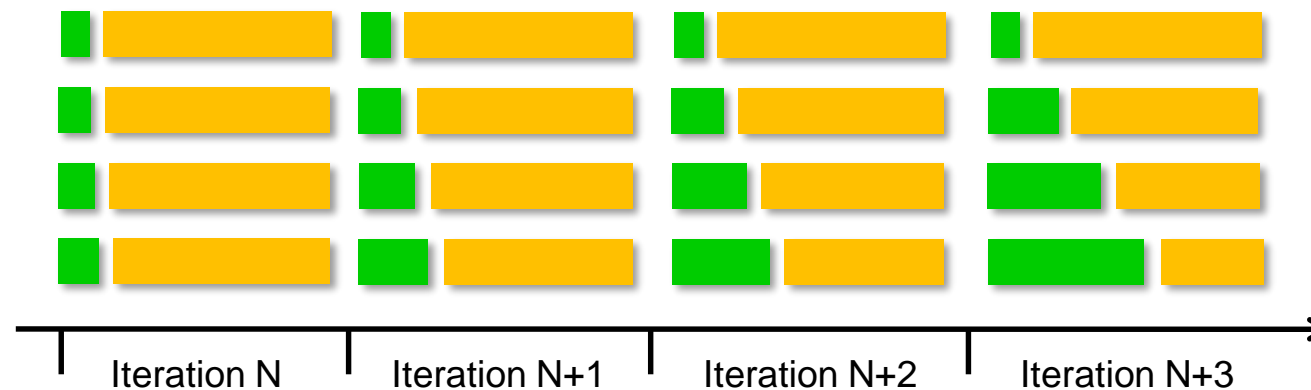


# PROFILING: ISSUES RELATED TO "AVERAGING"

- Moving bottleneck across processors can "average out" imbalances



- Imbalance changes over time  $\Rightarrow$  problem might not appear in short runs!





# MEASUREMENT METHODS: TRACING

- Recording **information about** significant points (**events**) during execution of the program
  - Enter/leave a code region (function, loop, ...)
  - Send/receive a message ...
- Save information in **event record**
  - Timestamp, location ID, event type
  - plus event specific information
- **Event trace** := stream of event records sorted by time

⇒ Abstract execution model on level of defined events

## Advantages

- + Can be used to reconstruct the dynamic behavior
- + Profiles can be calculated out of trace data

## Disadvantages

- HUGE trace files
- Can only be used for short durations or small configurations



# EVENT TRACING

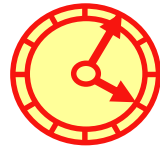
Process A

```
void foo() {  
  trc_enter("foo");  
  ...  
  trc_send(B);  
  send(B, tag, buf);  
  ...  
  trc_exit("foo");  
}
```

instrument

Process B

```
void bar() {  
  trc_enter("bar");  
  ...  
  recv(A, tag, buf);  
  trc_recv(A);  
  ...  
  trc_exit("bar");  
}
```



Local trace A

...		
58	ENTER	1
62	SEND	B
64	EXIT	1
...		

1	foo
...	

Local trace B

...		
60	ENTER	1
68	RECV	A
69	EXIT	1
...		

1	bar
...	

Global trace

...			
58	A	ENTER	1
60	B	ENTER	2
62	A	SEND	B
64	A	EXIT	1
68	B	RECV	A
69	B	EXIT	2
...			

merge

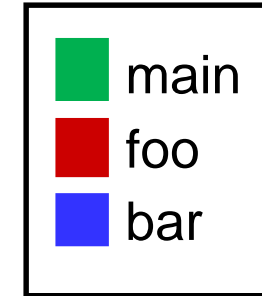
unify

1	foo
2	bar
...	

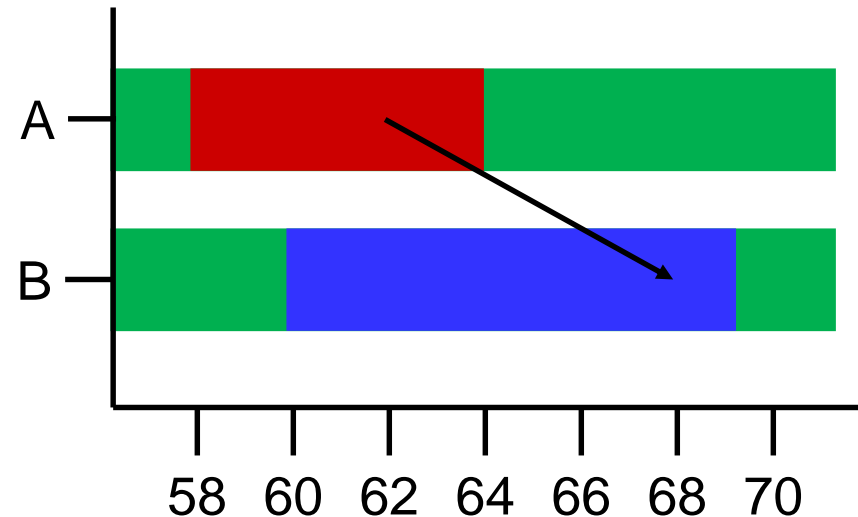


# EVENT TRACING: “TIMELINE” VISUALIZATION

1	foo
2	bar
3	...



...			
58	A	ENTER	1
60	B	ENTER	2
62	A	SEND	B
64	A	EXIT	1
68	B	RECV	A
69	B	EXIT	2
...			





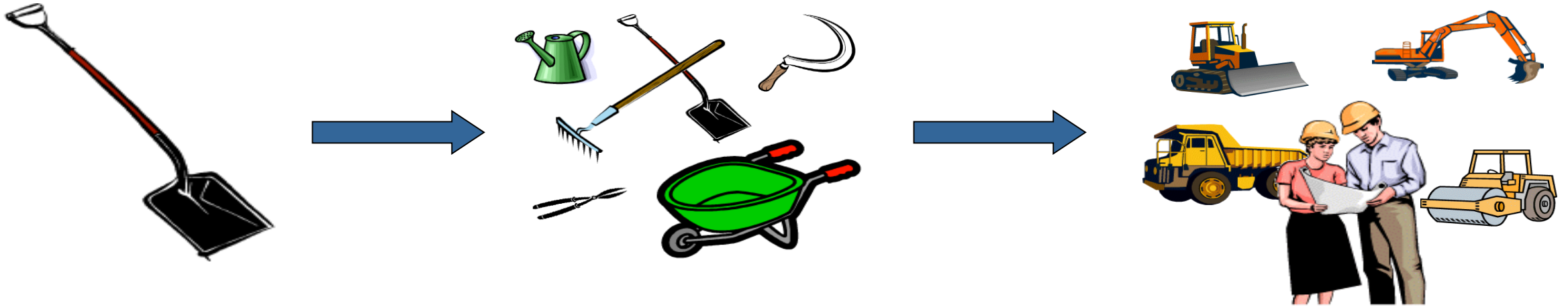
# CRITICAL ISSUES

- Accuracy
  - Intrusion overhead
    - Measurement takes 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*



# REMARK: 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, ...



# PERFORMANCE TOOLS (**STATUS: NOV 2023**)

- Score-P
- Scalasca
- Vampir[Server]
- Linaro Forge
  - Performance Reports
  - MAP
- Intel Tools
  - VTune Amplifier XE
  - Intel Advisor
- AMD uProf
- NVIDIA Tools
  - Nsight Systems
  - Nsight Compute
- Darshan
- ...



# **Score-P**

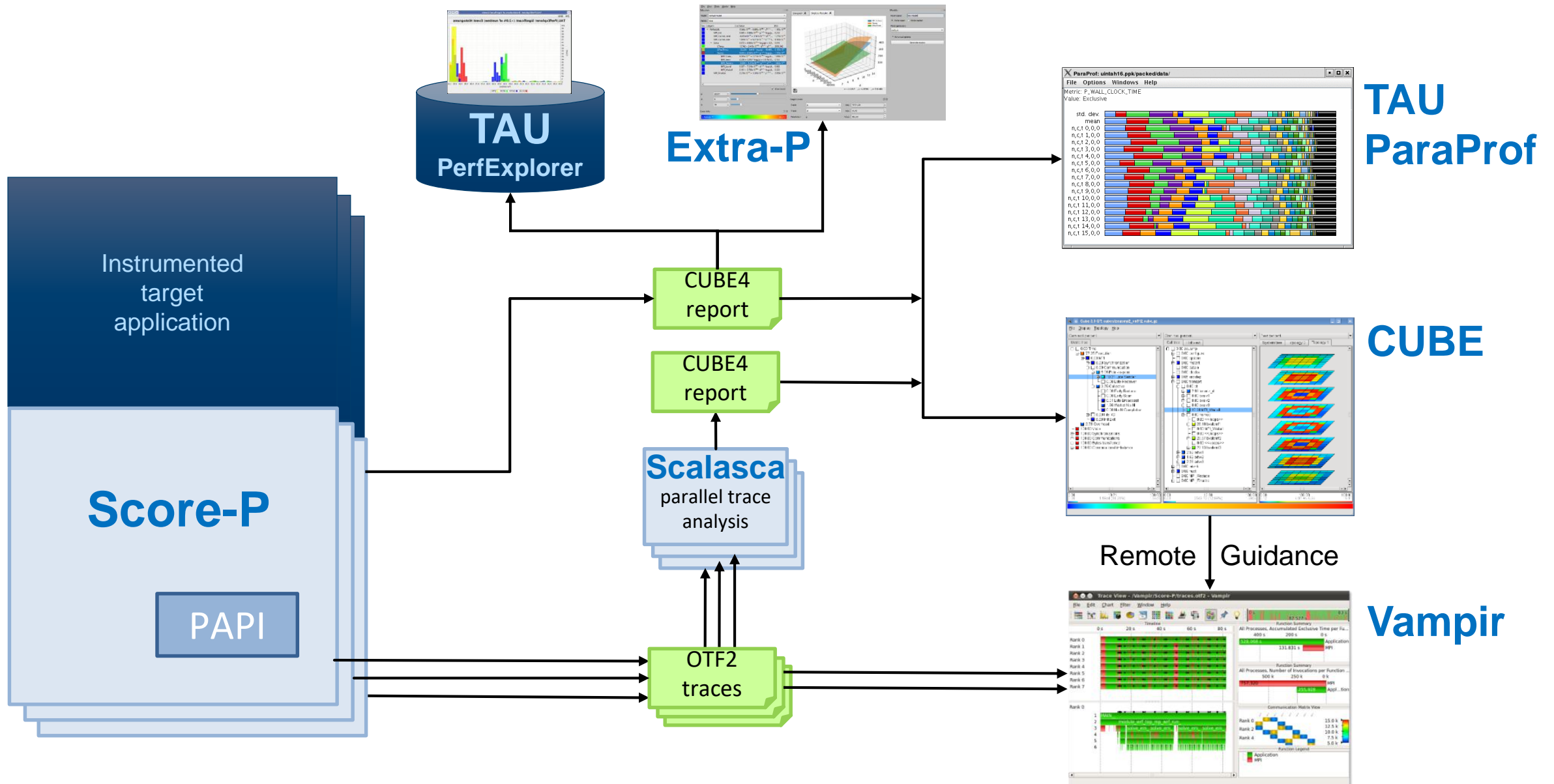
Scalable performance measurement  
infrastructure for parallel codes

- Community-developed open-source
- Replaced tool-specific instrumentation and measurement components of partners
- <http://www.score-p.org>



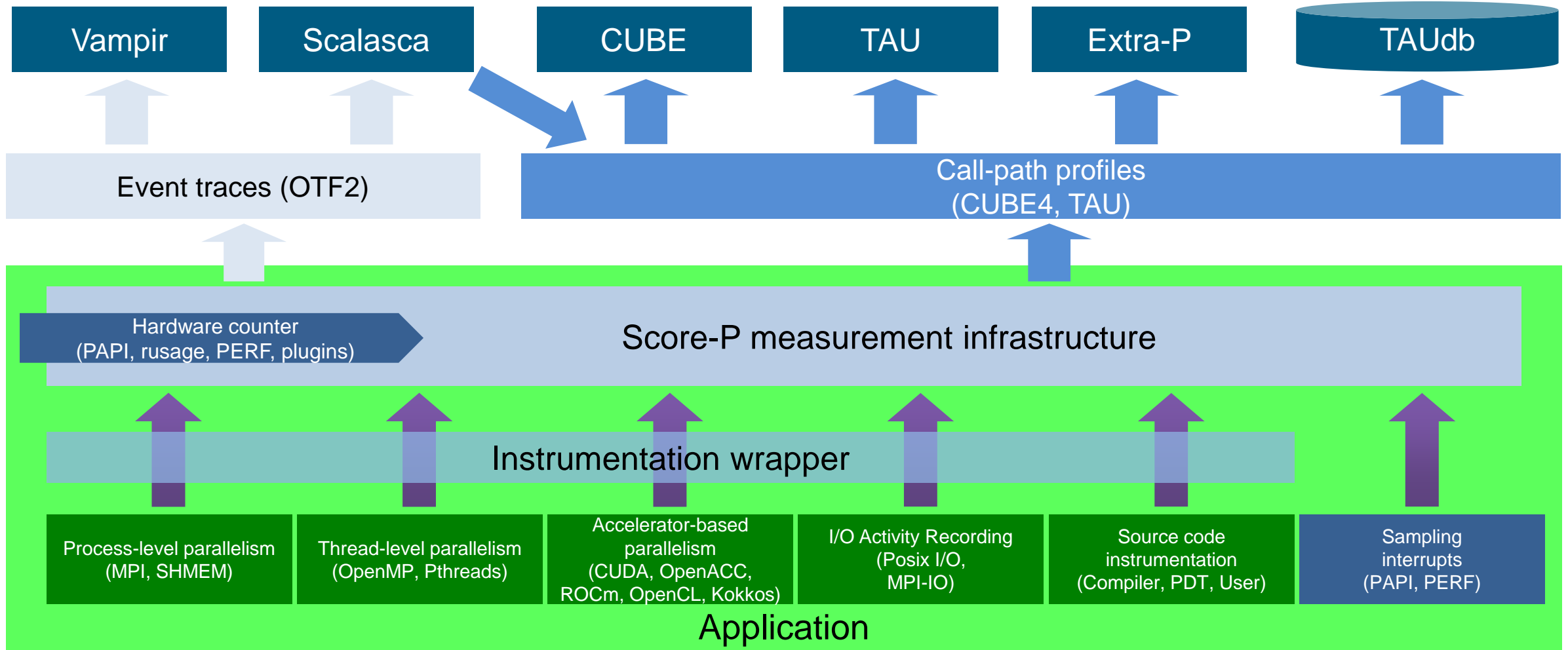


# Score-P TOOL ECOSYSTEM





# #Score-P ARCHITECTURE



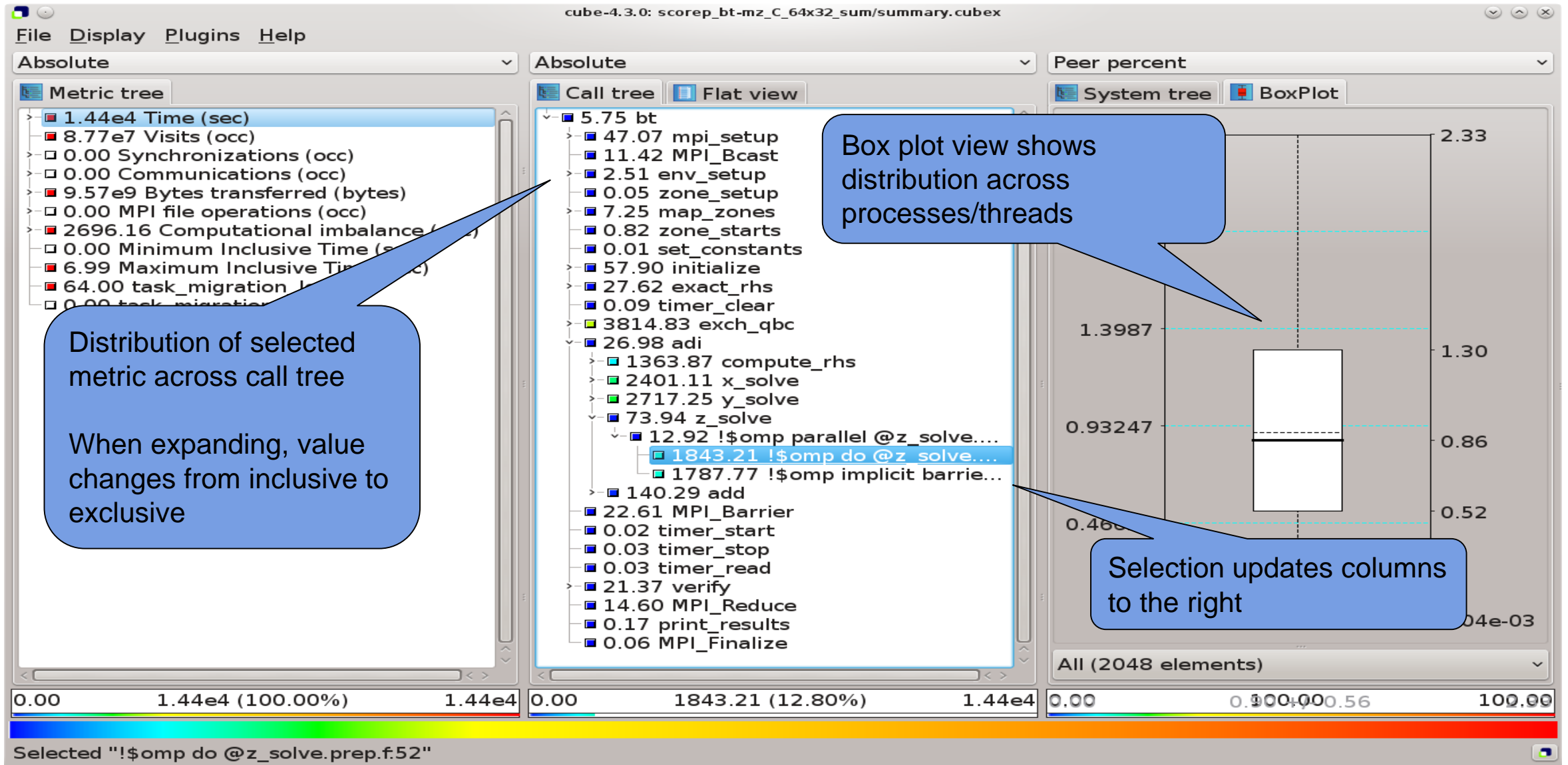


# **Score-P FUNCTIONALITY**

- Provide typical functionality for HPC performance tools
- **Instrumentation** (various methods)
  - Multi-process paradigms (MPI, SHMEM)
  - Thread-parallel paradigms (OpenMP, POSIX threads)
  - Accelerator-based paradigms (OpenACC, CUDA, OpenCL, Kokkos)
  - **In any combination!**
- Flexible **measurement** without re-compilation:
  - Basic and advanced **profile** generation (⇒ CUBE4 format)
  - Event **trace** recording (⇒ OTF2 format)
- Highly scalable I/O functionality
- Support all fundamental concepts of partner's tools



# CUBE EXAMPLE





# SCORE-P: ADVANCED FEATURES

- Measurement can be extensively configured via environment variables
- Allows for targeted measurements:
  - Selective recording
  - Phase profiling
  - Parameter-based profiling
  - ...
- GPU support: CUDA, OpenACC, OpenCL, HIP, Kokkos, ...
- Please ask us or see the user manual for details



- Scalable Analysis of Large Scale Applications

- Approach

- **Instrument** C, C++, and Fortran parallel applications (**with Score-P**)

- Option 1: scalable call-path profiling

- Option 2: scalable event trace analysis

- **Collect** event traces

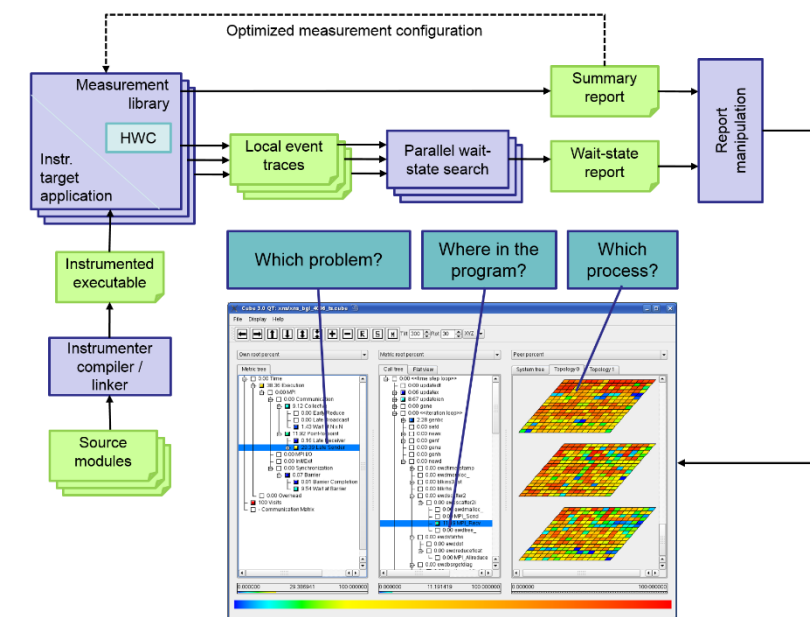
- **Process trace in parallel**

- Wait-state analysis

- Delay and root-cause analysis

- Critical path analysis

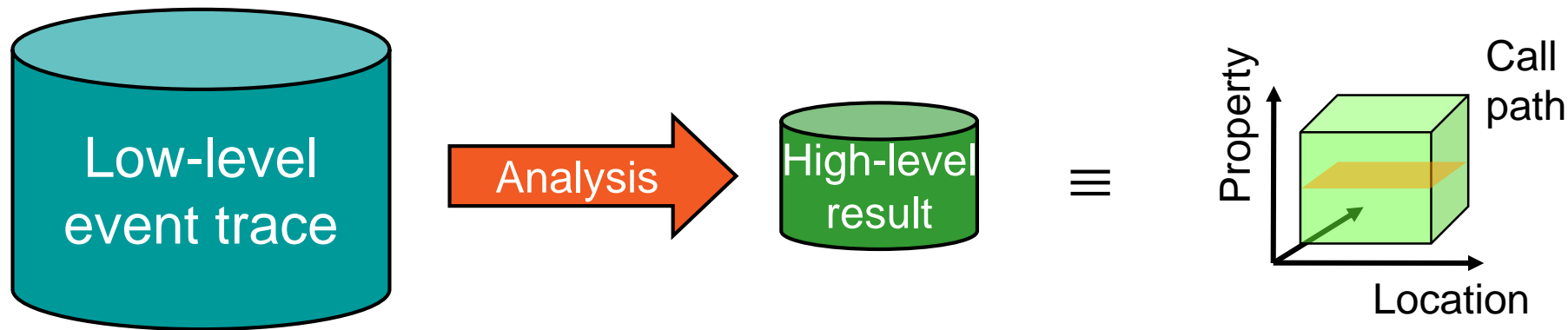
- **Categorize and rank** results





# AUTOMATIC TRACE ANALYSIS

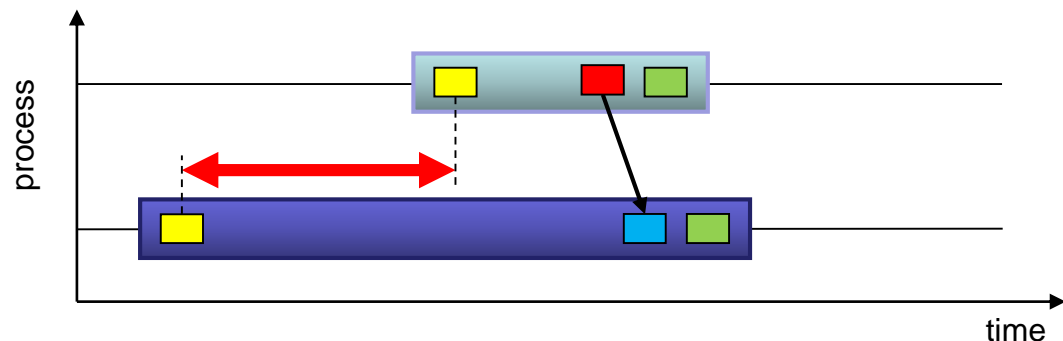
- Automatic search for patterns of inefficient behaviour
- Classification of behaviour & quantification of significance
- Identification of delays as root causes of inefficiencies



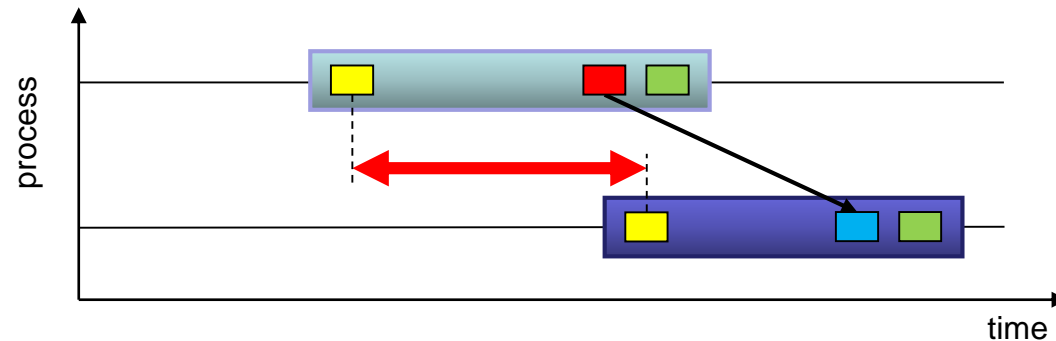
- Guaranteed to cover the entire event trace
- Quicker than manual/visual trace analysis
- Parallel replay analysis exploits available memory & processors to deliver scalability



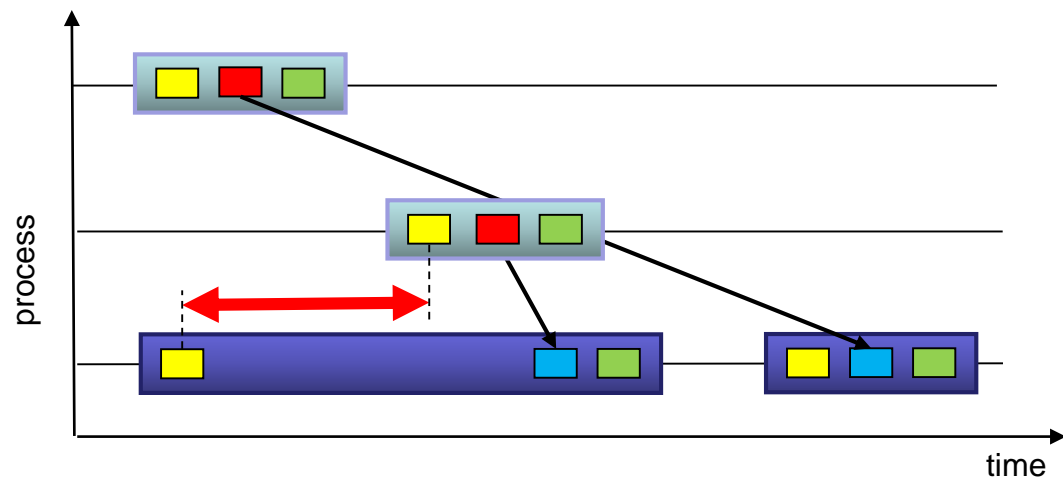
# EXAMPLE MPI WAIT STATES



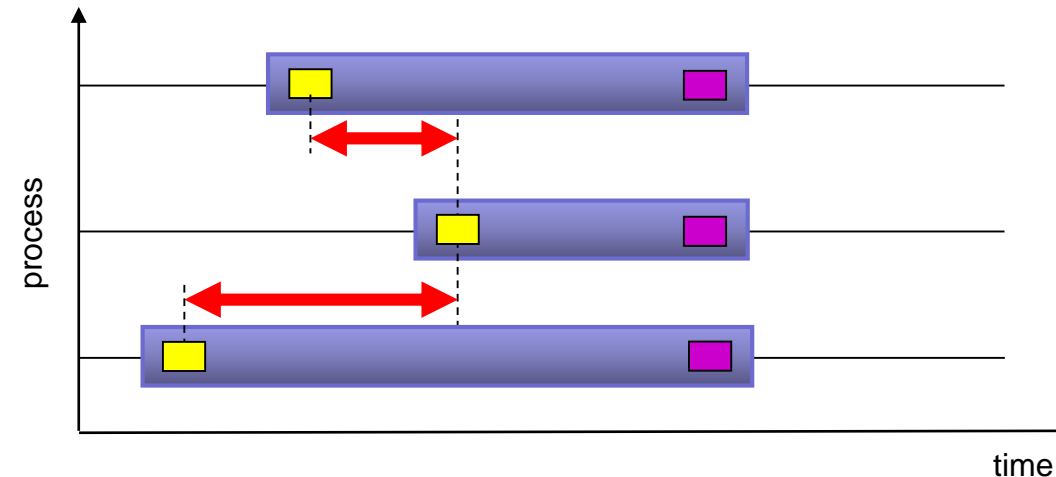
(a) Late Sender



(b) Late Receiver



(c) Late Sender / Wrong Order



(d) Wait at N x N

ENTER EXIT SEND RECV COLLEXIT



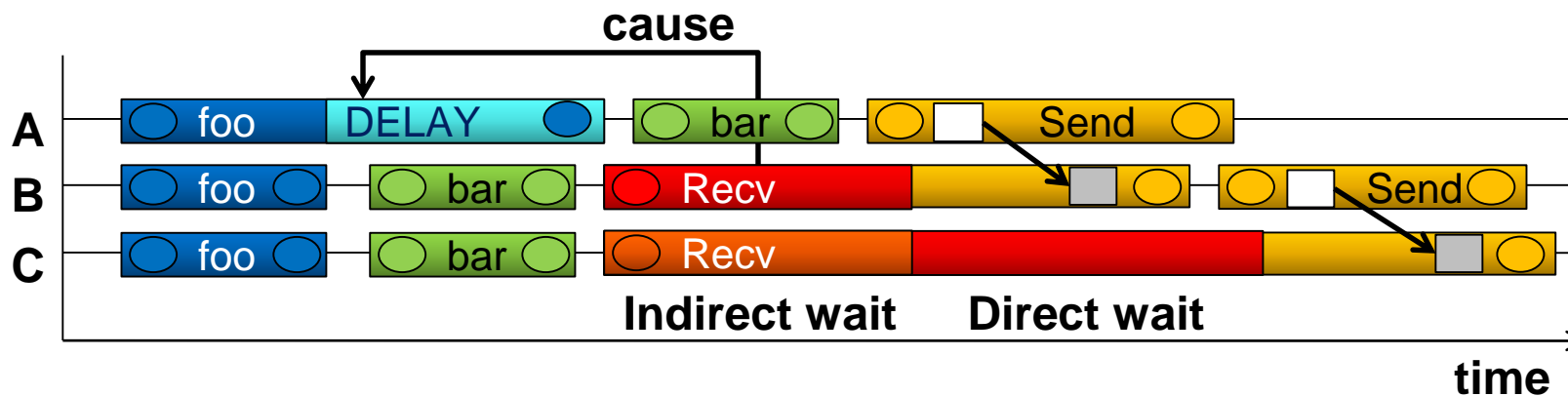
# SCALASCA ROOT CAUSE ANALYSIS

- **Root-cause analysis**

- Wait states typically caused by load or communication imbalances earlier in the program
- Waiting time can also propagate (e.g., indirect waiting time)
- Enhanced performance analysis to find the root cause of wait states

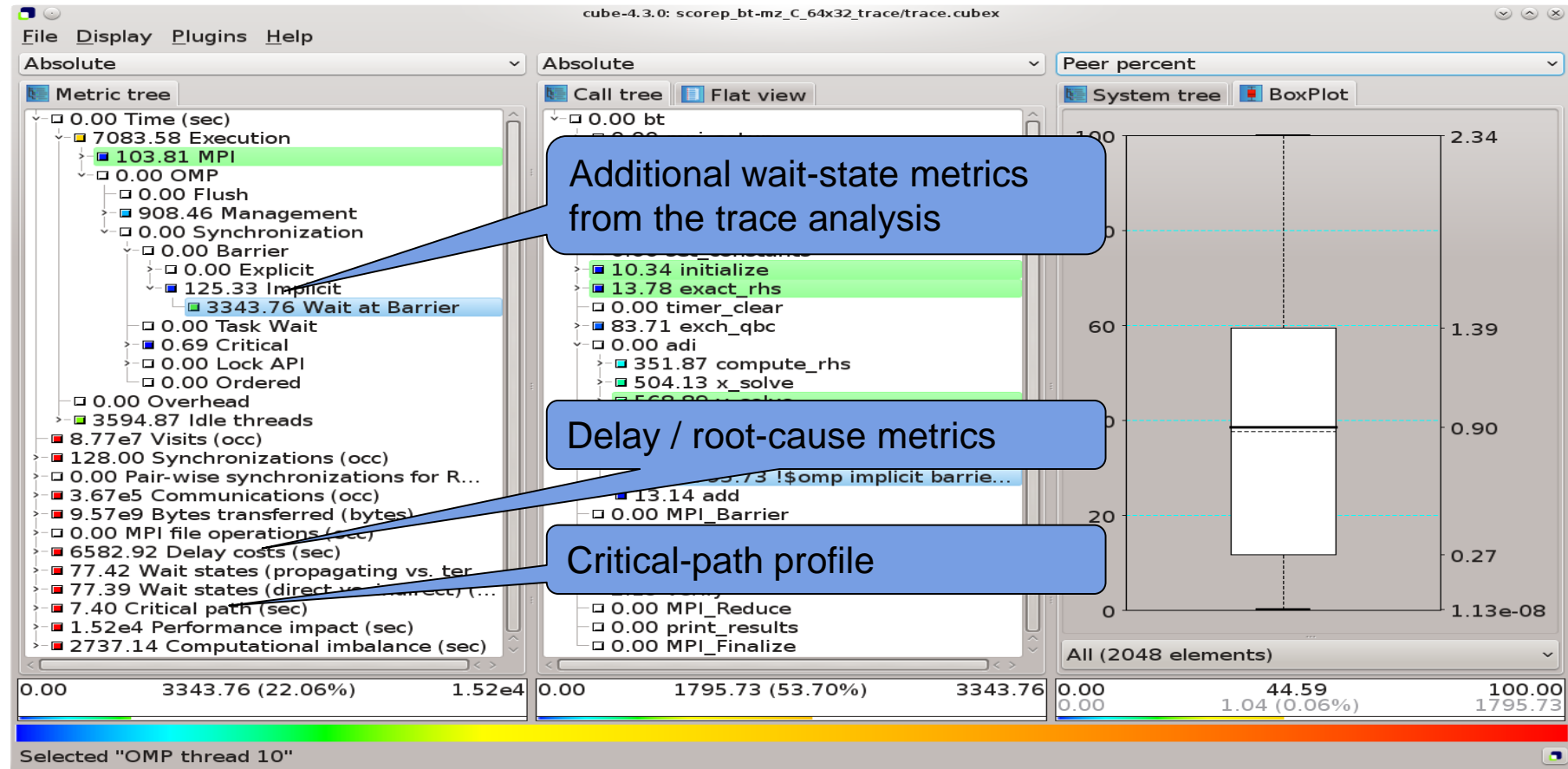
- **Approach**

- Distinguish between direct and indirect waiting time
- Identify call path/process combinations delaying other processes and causing first order waiting time
- Identify original **delay**





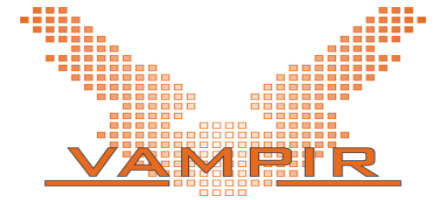
# SCALASCA TRACE ANALYSIS EXAMPLE





# VAMPIR EVENT TRACE VISUALIZER

- Offline trace visualization for Score-Ps OTF2 trace files
- Visualization of MPI, OpenMP and application events:
  - All diagrams highly customizable (through context menus)
  - Large variety of displays for ANY part of the trace
- <http://www.vampir.eu>
- Advantage:
  - Detailed view of dynamic application behavior
- Disadvantage:
  - Completely manual analysis
  - Too many details can hide the relevant parts



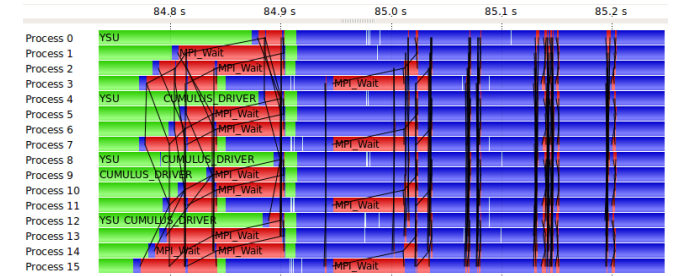


# EVENT TRACE VISUALIZATION WITH VAMPIR

- Visualization of dynamic runtime behaviour at any level of detail along with statistics and performance metrics
- Alternative and supplement to automatic analysis
- **Typical questions that Vampir helps to answer**
  - What happens in my application execution during a given time in a given process or thread?
  - How do the communication patterns of my application execute on a real system?
  - Are there any imbalances in computation, I/O or memory usage and how do they affect the parallel execution of my application?

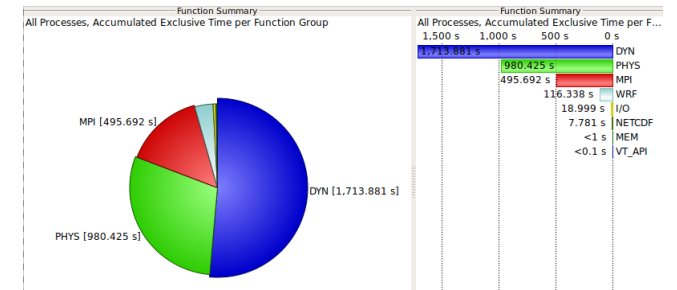
## ■ Timeline charts

- Application activities and communication along a time axis



## ■ Summary charts



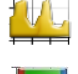

- Quantitative results for the currently selected time interval







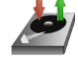



# VAMPIR PERFORMANCE CHARTS

## Timeline Charts

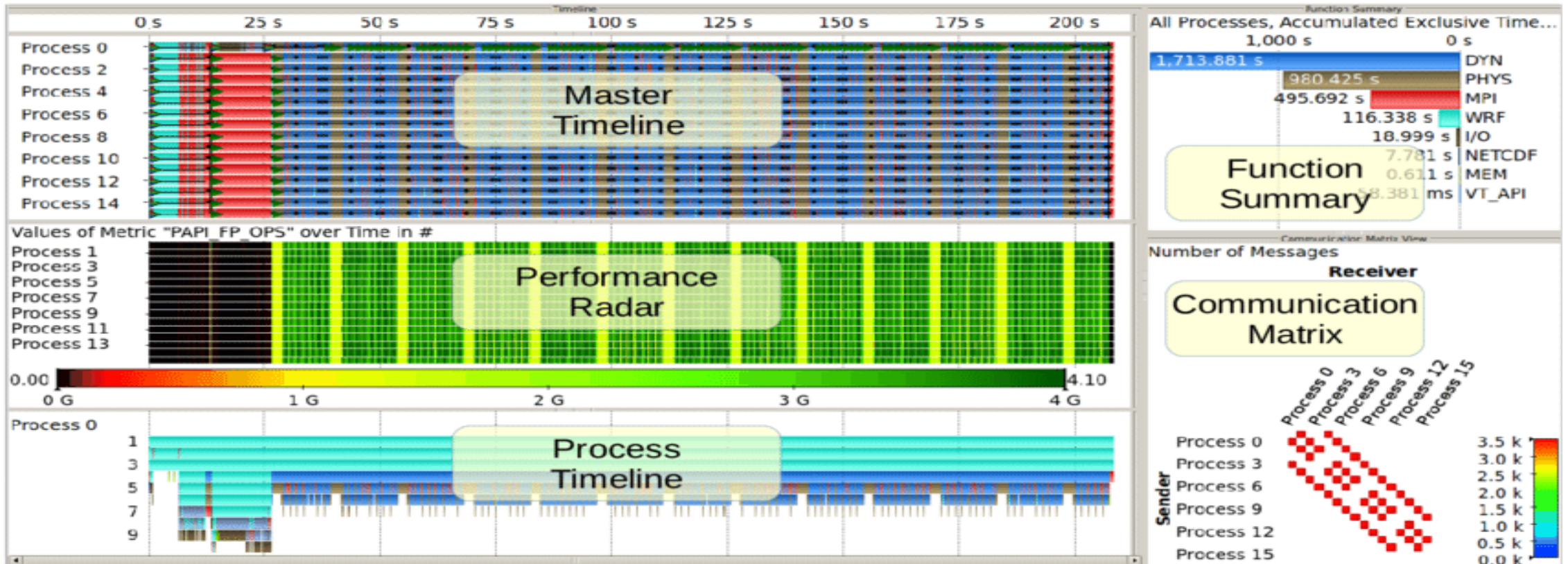
	Master Timeline	➔	<i>all threads' activities</i>
	Process Timeline	➔	<i>single thread's activities</i>
	Summary Timeline	➔	<i>all threads' function call statistics</i>
	Performance Radar	➔	<i>all threads' performance metrics</i>
	Counter Data Timeline	➔	<i>single threads' performance metrics</i>
	I/O Timeline	➔	<i>all threads' I/O activities</i>

## Summary Charts

	Function Summary		Process Summary
	Message Summary		Communication Matrix View
	I/O Summary		Call Tree



# VAMPIR DISPLAYS





# LINARO PERFORMANCE REPORTS



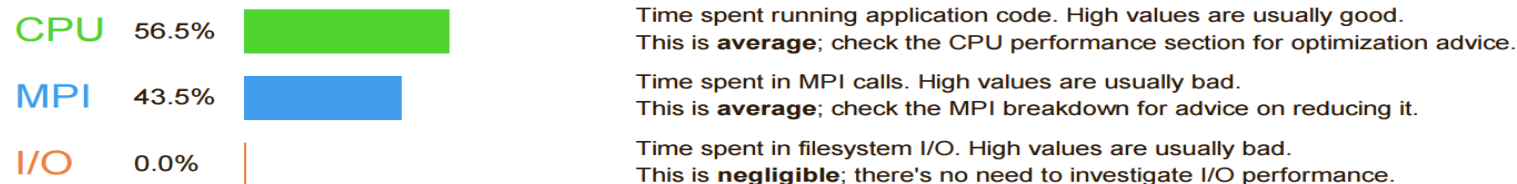
- **Single page** report provides quick overview of performance issues
- Works on unmodified, optimized executables
- Shows CPU, memory, network and I/O utilization
- Supports MPI, multi-threading and accelerators
- Saves data in HTML, CVS or text form
- <https://www.linaroforge.com/linaroPerformanceReports>
- **Note:** License limited to 128 processes (with unlimited number of threads)



# EXAMPLE PERFORMANCE REPORTS

Summary: cp2k.popt is **CPU-bound** in this configuration

The total wallclock time was spent as follows:



This application run was **CPU-bound**. A breakdown of this time and advice for investigating further is in the **CPU** section below.

## CPU

A breakdown of how the **56.5%** total CPU time was spent:

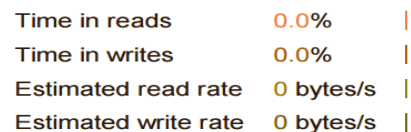


The per-core performance is **memory-bound**. Use a profiler to identify time-consuming loops and check their cache performance.

Little time is spent in **vectorized instructions**. Check the compiler's vectorization advice to see why key loops could not be vectorized.

## I/O

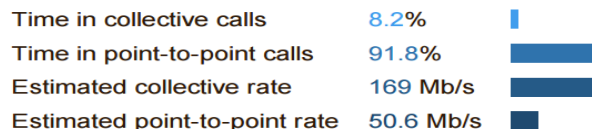
A breakdown of how the **0.0%** total I/O time was spent:



No time is spent in **I/O operations**. There's nothing to optimize here!

## MPI

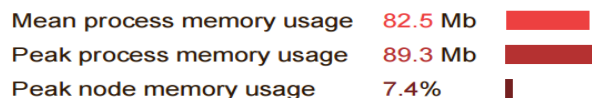
Of the **43.5%** total time spent in MPI calls:



The **point-to-point** transfer rate is low. This can be caused by inefficient message sizes, such as many small messages, or by imbalanced workloads causing processes to wait. Use an MPI profiler to identify the problematic calls and ranks.

## Memory

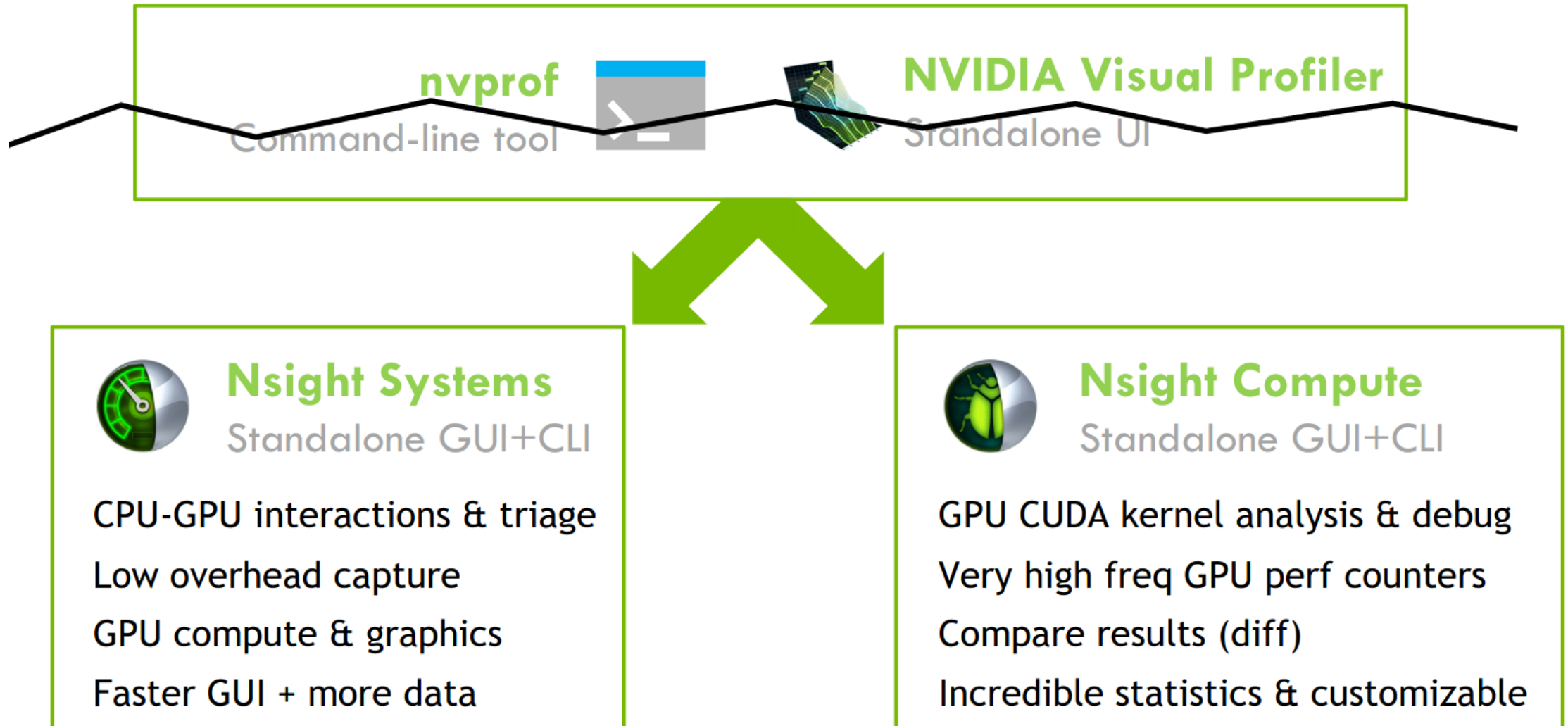
Per-process memory usage may also affect scaling:



The **peak node memory usage** is low. You may be able to reduce the total number of CPU hours used by running with fewer MPI processes and more data on each process.



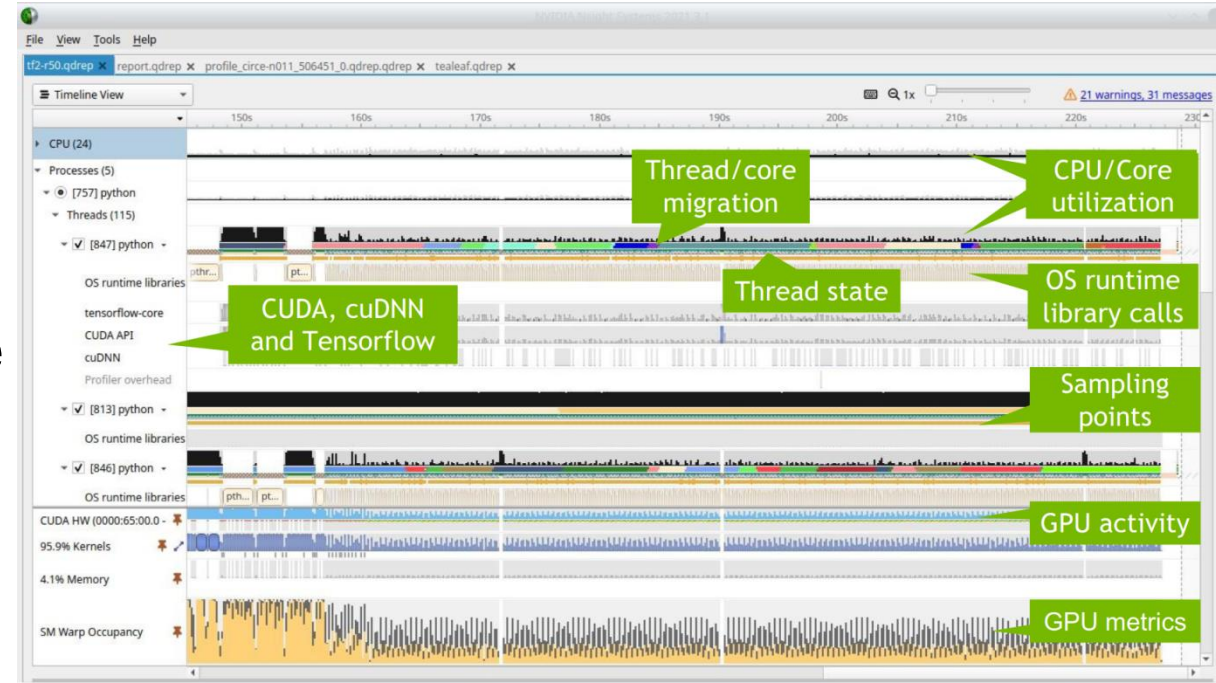
# NVIDIA TOOLS -- LEGACY TRANSITION





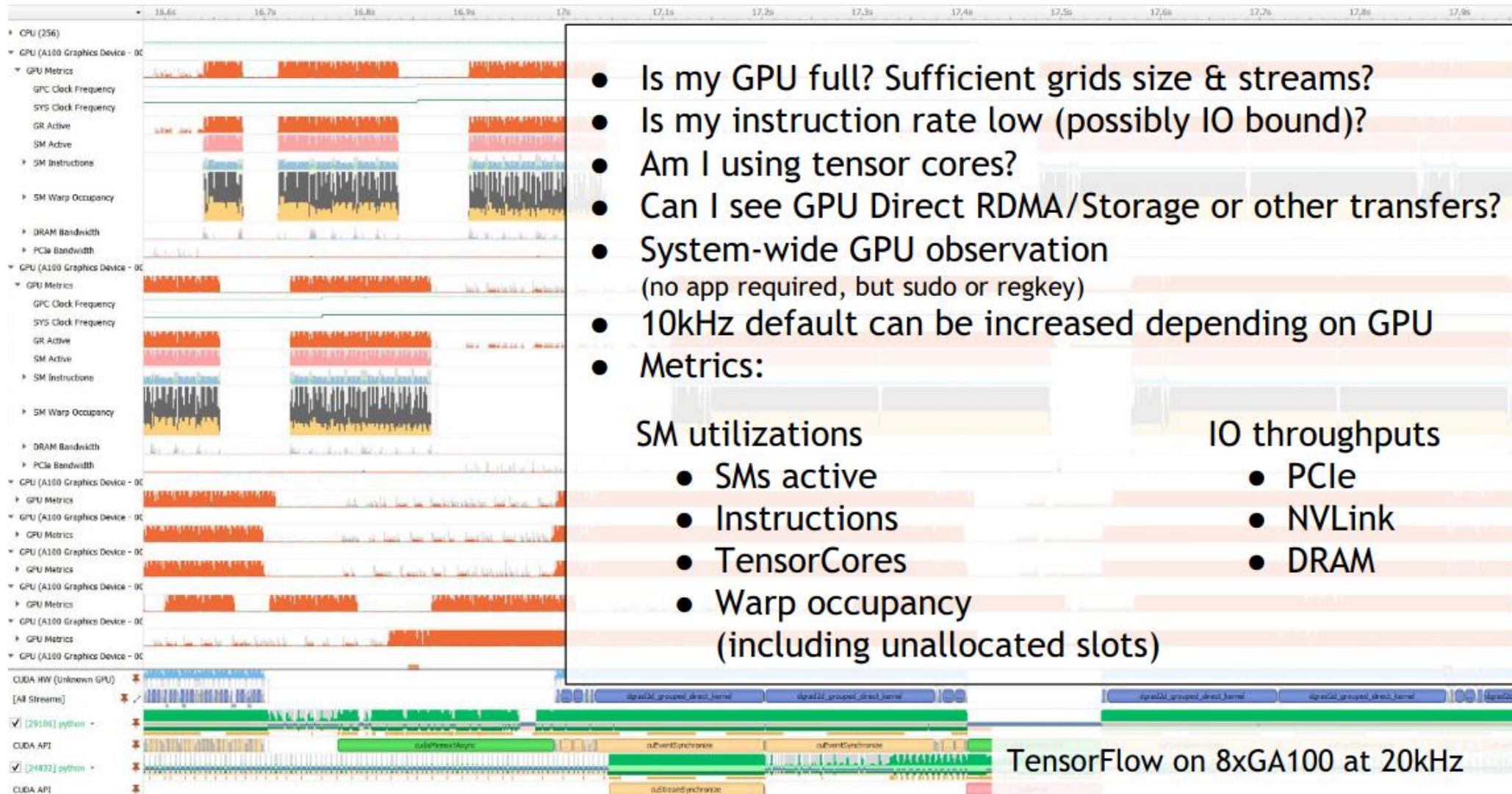
# NSIGHT SYSTEM

- System-wide application tuning
- Locate optimization opportunities
  - Visualize millions of events on a timeline
  - See gaps of unused CPU and GPU time
- Balance workloads across multiple CPUs and GPUs
  - CPU utilization and thread state
  - GPU streams, kernels, memory transfers, etc.
- Multi-platform support
  - Linux, Windows and Mac OS X (host-only)
  - x86-64, Power9, ARM server, Tegra (Linux & QNX)



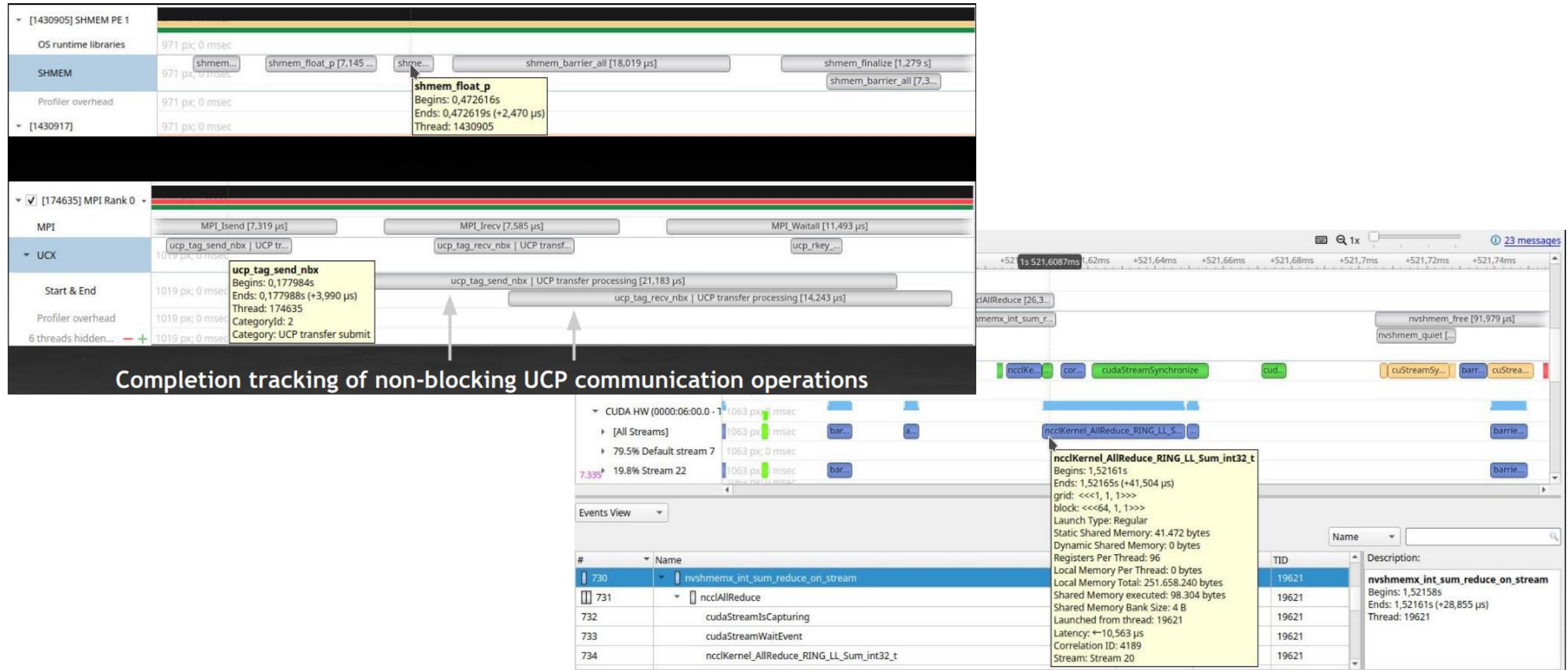


# GPU METRIC SAMPLING



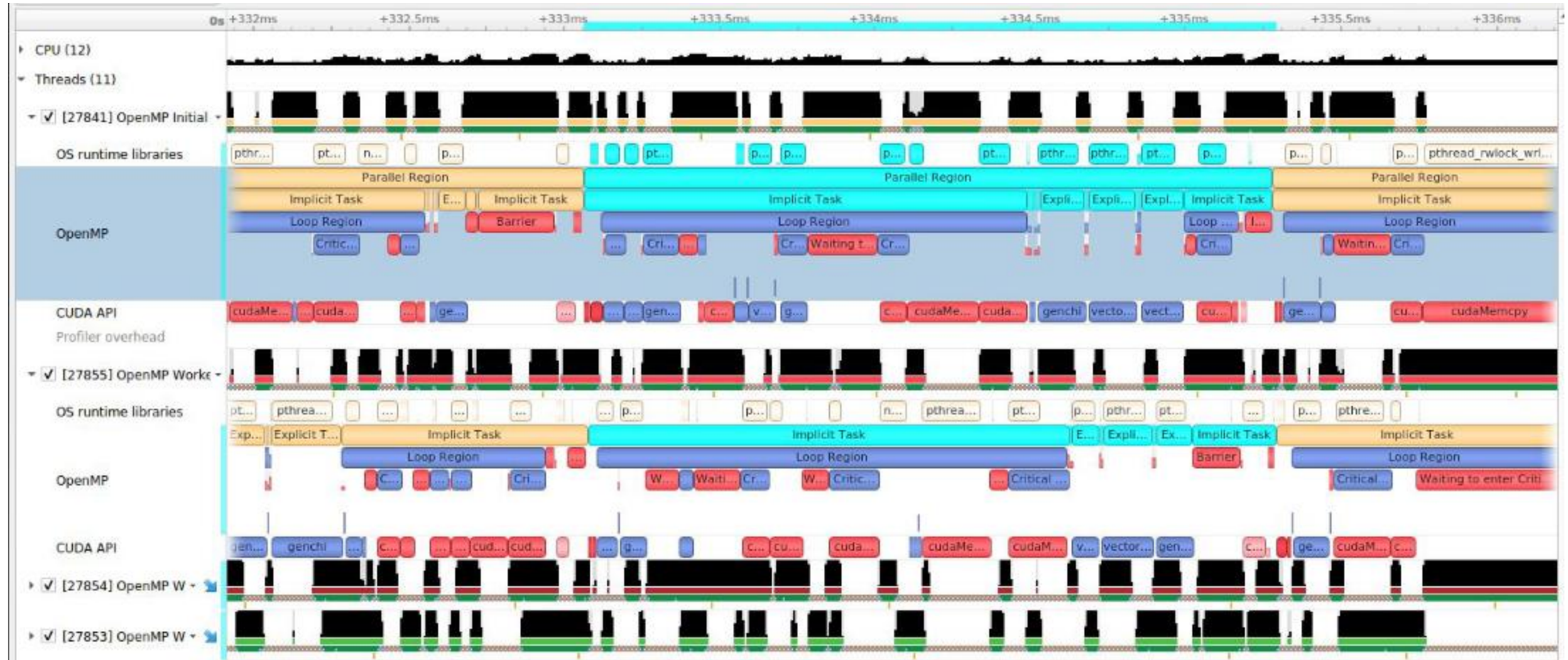


# MULTI NODE SUPPORT – SHMEM, MPI, UCX, AND NCCL





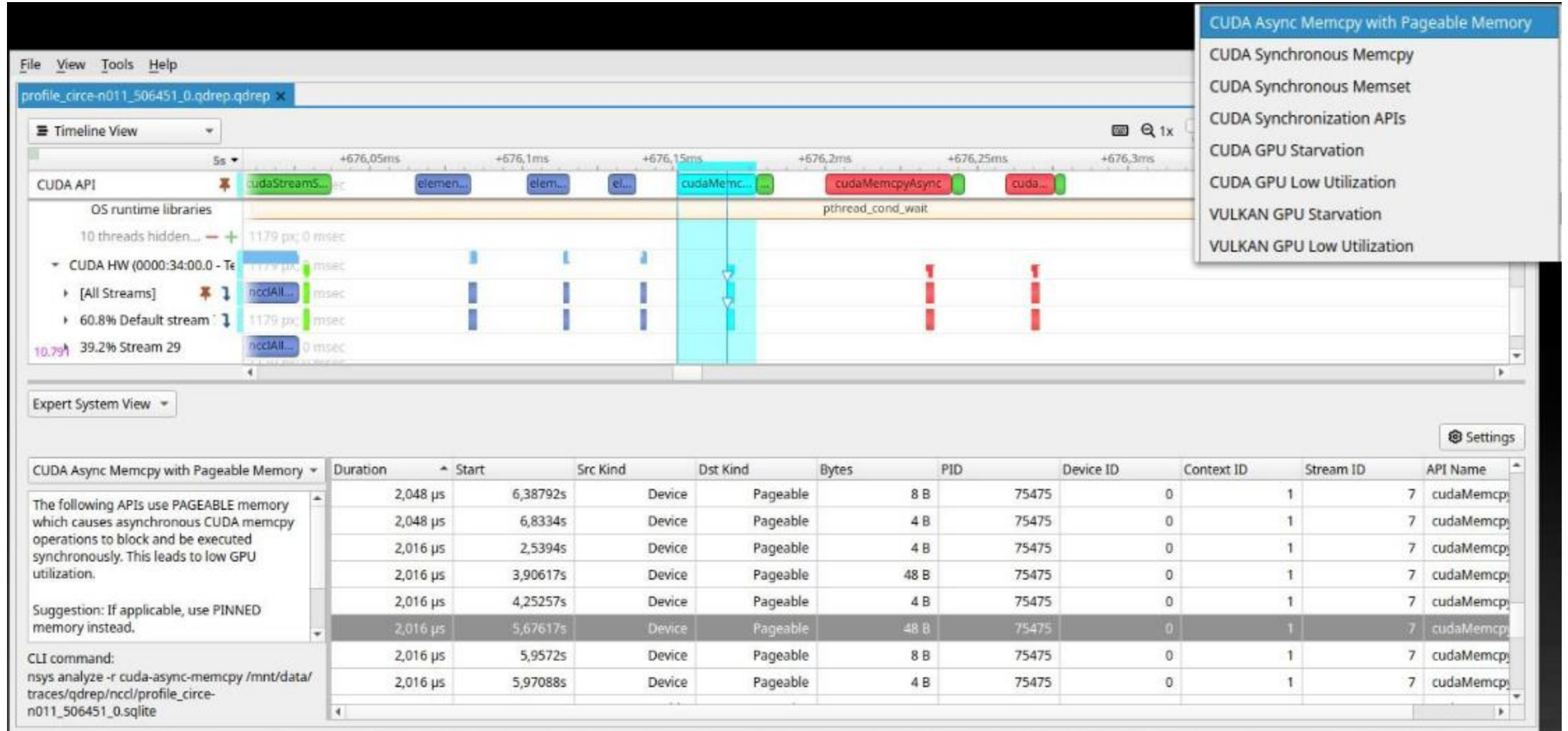
# OPENMP



OMPT-capable OpenMP runtime required



# EXPERT SYSTEM





# NSIGHT COMPUTE

- Interactive CUDA kernel profiler
- Targeted metric sections for various performance aspects
- Customizable data collection and presentation (tables, charts, ...)
- GUI and CLI
- Python-based API for guided analysis and post-processing
- Support for remote profiling across machines and platforms

The screenshot displays the NVIDIA Nsight Compute interface. The top section shows a summary table of kernel launches. The bottom section provides a detailed view of GPU throughput metrics, including Compute (SM) Throughput, Memory Throughput, L1/TEX Cache Throughput, L2 Cache Throughput, and DRAM Throughput. It also includes a 'High Compute Throughput' warning and a 'FP64/32 Utilization' section.

Page:	Summary	Launch:	0 - 43843 - device_tea_leaf_ppcg_sol	▼	▼	Add Baseline	▼	Apply Rules	▼	Occupancy Calculator	Copy as Image
	Launch	Time	Cycles	Regs	GPU	SM Frequency	CC	Process			
Current	43843 - device_tea_leaf_ppcg_solve_init (126, 1001, 1)x(32, 4, 1)	217.63 usecond	297,114	40	0 - NVIDIA GeForce RTX 2080 Ti	1.36 cycle/nsecond	7.5	[15958] tea_leaf			

ID	Time	API Call ID	Function Name	Demangled N:	Process	Device Name	Grid Size	Block Size	Cycles [cycle]	Duration [msecond]	Compute Throughput [%]	Memc
0	2021-Dec-1...	43843	device_tea_leaf_ppcg...	device_te...	[15958] tea_leaf	NVIDIA GeForce...	126, 1001, 1	32, 4, 1	297,114	0.22	77.89	
1	2021-Dec-1...	43857	device_tea_leaf_ppcg_sol...	device_tea_J...	[15958] tea_leaf	NVIDIA GeForce RT...	126, 1001, 1	32, 4, 1	1,264,921	0.94	54.78	
2	2021-Dec-1...	43860	device_tea_leaf_ppcg_sol...	device_tea_J...	[15958] tea_leaf	NVIDIA GeForce RT...	126, 1001, 1	32, 4, 1	1,462,446	1.07	86.22	
3	2021-Dec-1...	43863	device_tea_leaf_ppcg_sol...	device_tea_J...	[15958] tea_leaf	NVIDIA GeForce RT...	126, 1001, 1	32, 4, 1	1,443,836	1.06	23.81	

Page:	Details	Launch:	0 - 43843 - device_tea_leaf_ppcg_sol	▼	▼	Add Baseline	▼	Apply Rules	▼	Occupancy Calculator	Copy as Image
	Launch	Time	Cycles	Regs	GPU	SM Frequency	CC	Process			
Current	43843 - device_tea_leaf_ppcg_solve_init (126, 1001, 1)x(32, 4, 1)	217.63 usecond	297,114	40	0 - NVIDIA GeForce RTX 2080 Ti	1.36 cycle/nsecond	7.5	[15958] tea_leaf			

GPU Speed Of Light Throughput		
Compute (SM) Throughput [%]	77.89	Duration [usecond]
Memory Throughput [%]	45.03	Elapsed Cycles [cycle]
L1/TEX Cache Throughput [%]	68.22	SM Active Cycles [cycle]
L2 Cache Throughput [%]	2.30	SM Frequency [cycle/nsecond]
DRAM Throughput [%]	0.12	DRAM Frequency [cycle/nsecond]

**High Compute Throughput** Compute is more heavily utilized than Memory: Look at the [Compute Workload Analysis](#) report section to see what the compute pipelines are spending their time doing. Also, consider whether any computation is redundant and could be reduced or moved to look-up tables.

**FP64/32 Utilization** The ratio of peak float (fp32) to double (fp64) performance on this device is 32:1. The kernel achieved 0% of this device's fp32 peak performance and 19% of its fp64 peak performance. If [Compute Workload Analysis](#) determines that this kernel is fp64 bound, consider using 32-bit precision floating point operations to improve its performance. See the [Kernel Profiling Guide](#) for more details on runtime analysis.

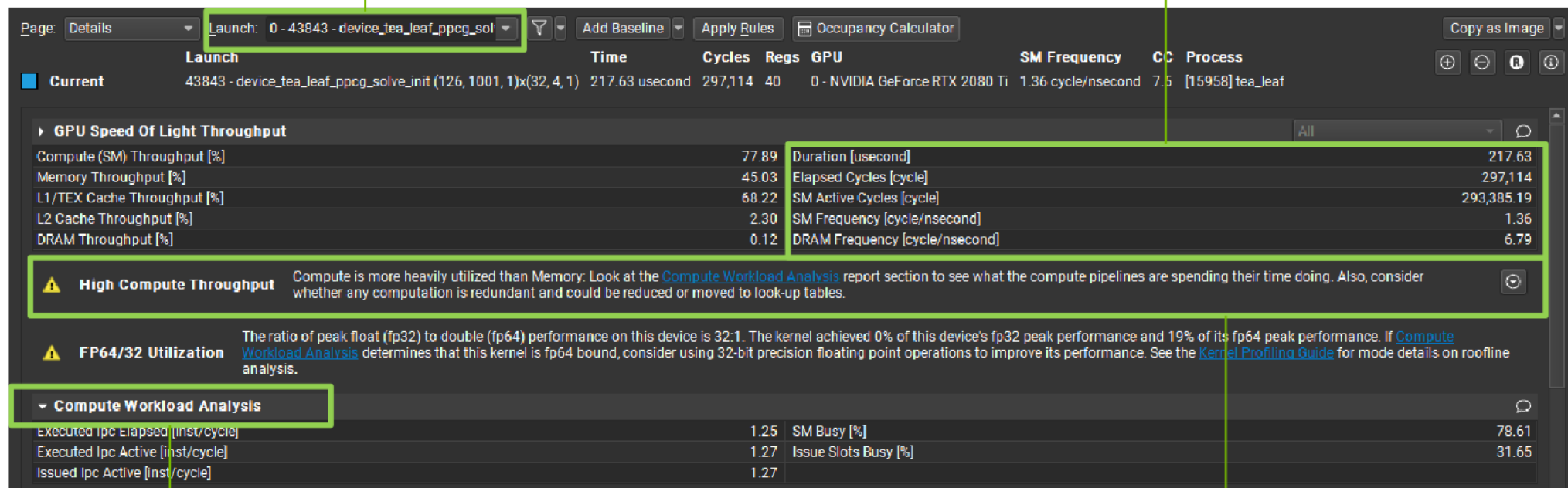
Compute Workload Analysis		
Executed lpc Elapsed [inst/cycle]	1.25	SM Busy [%]
Executed lpc Active [inst/cycle]	1.27	Issue Slots Busy [%]
Issued lpc Active [inst/cycle]	1.27	



# PROFILER REPORT

Selected result

Metric values



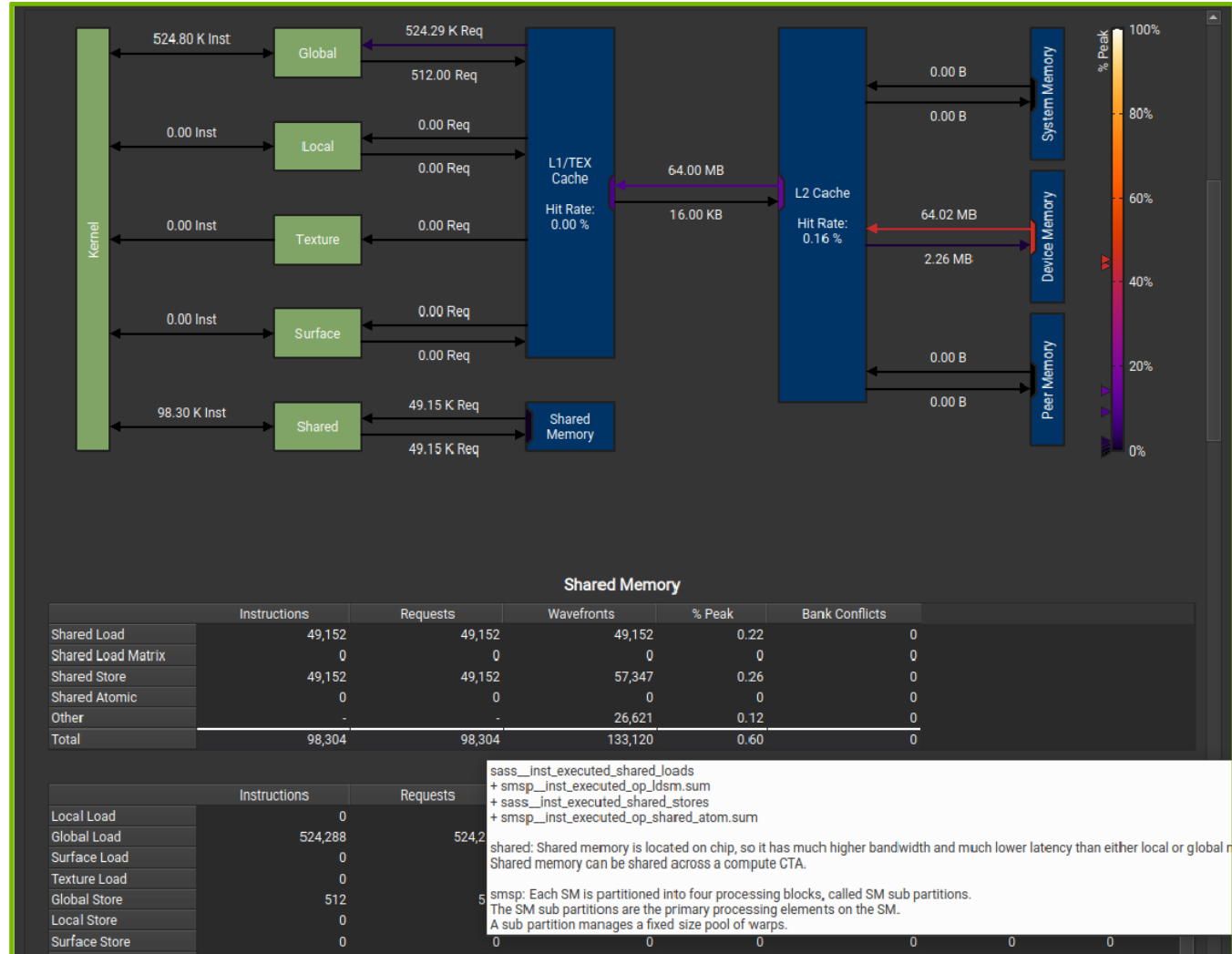
Expandable  
Sections

Expert Analysis  
(Rules)



# DATA TRANSFER ANALYSIS

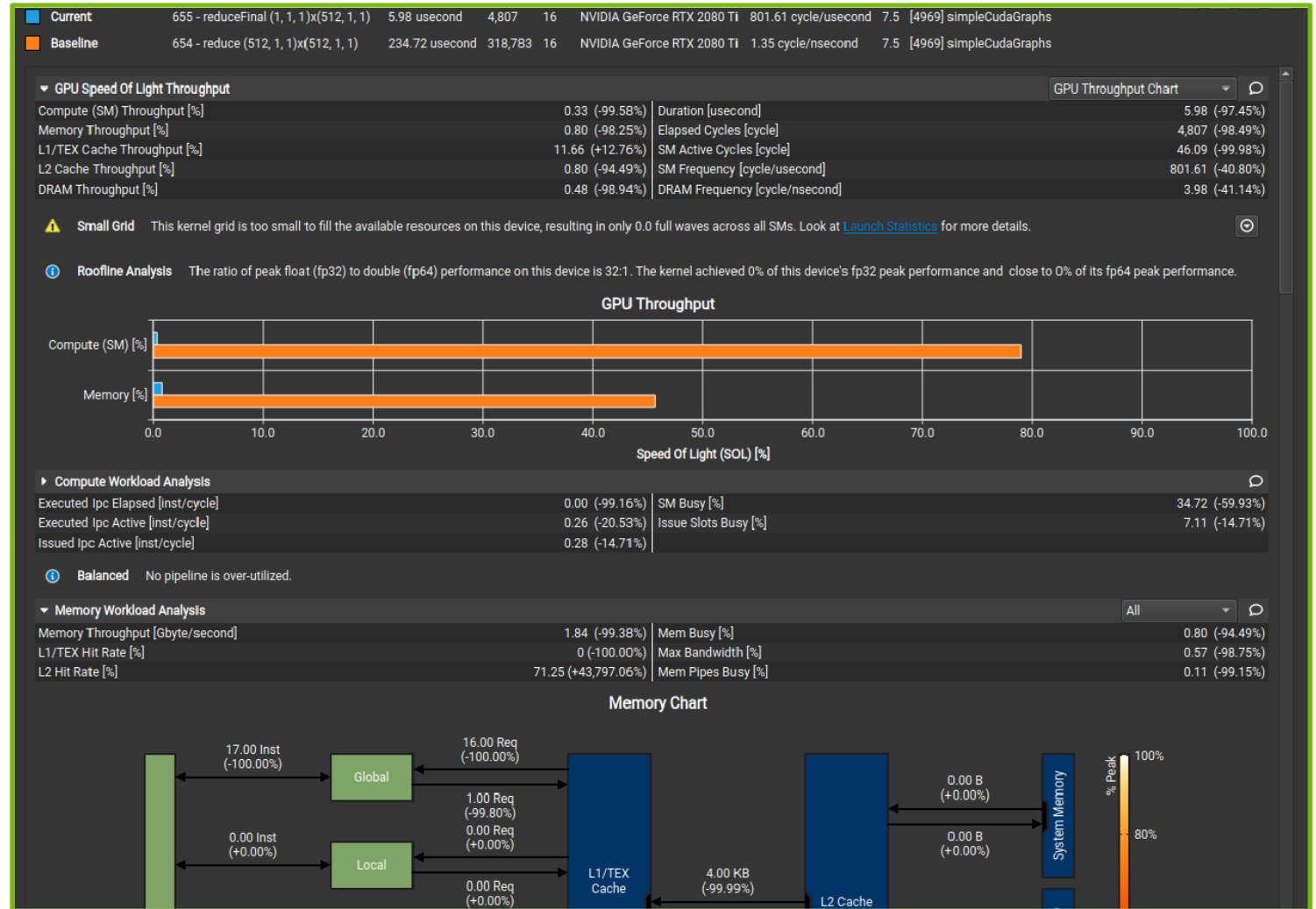
- Detailed memory workload analysis chart and tables
- Shows transferred data or throughputs
- Tooltips provide metric names, calculation formulas and detailed background info





# BASELINE COMPARISON

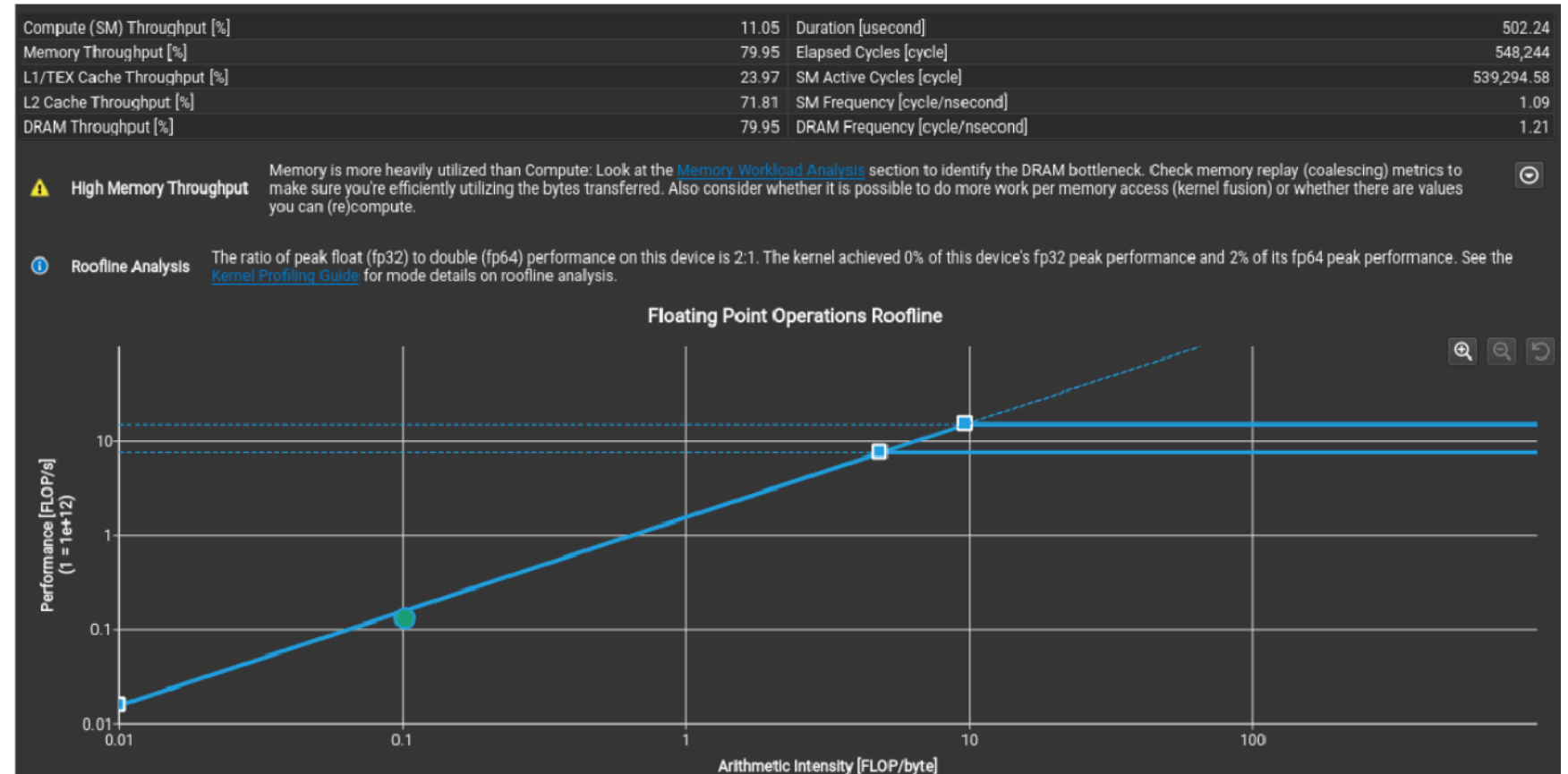
- Comparison of results directly within the tool with "Baselines"
- Supported across kernels, reports, and GPU architectures





# ROOFLINE ANALYSIS

- Determine whether the application is memory bound or compute bound
- Guided analysis points to detailed analysis of the most severe problem





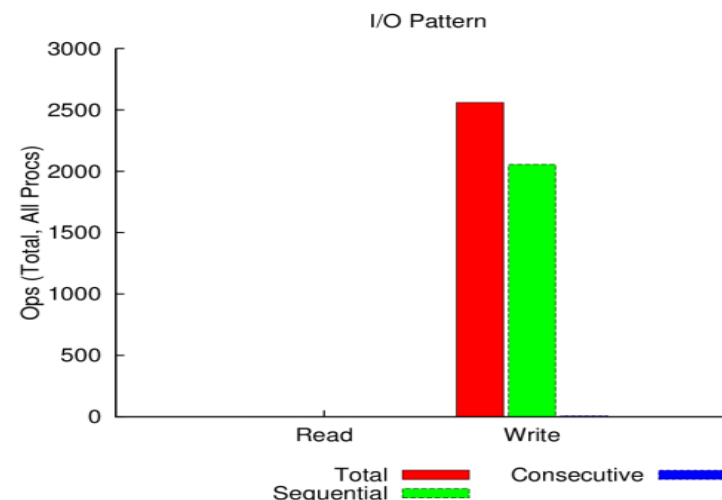
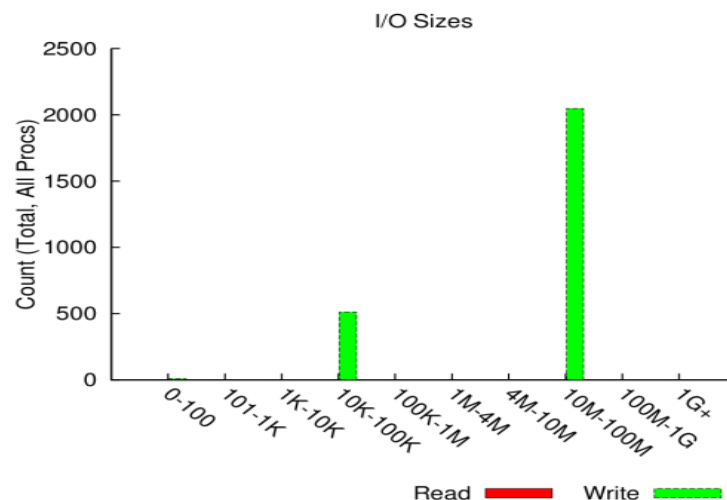
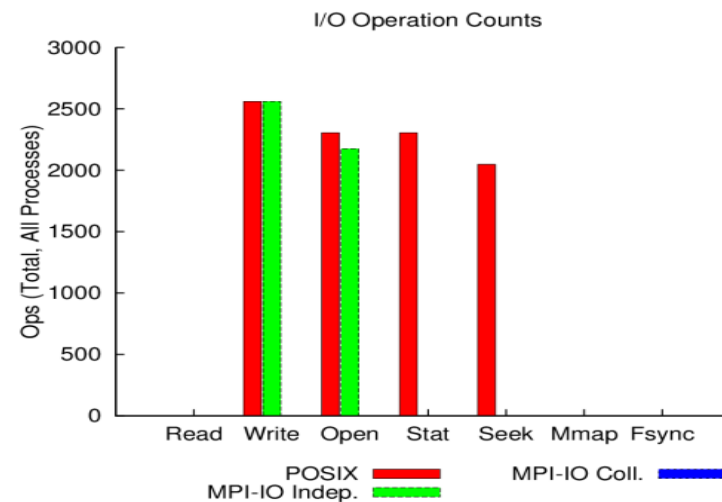
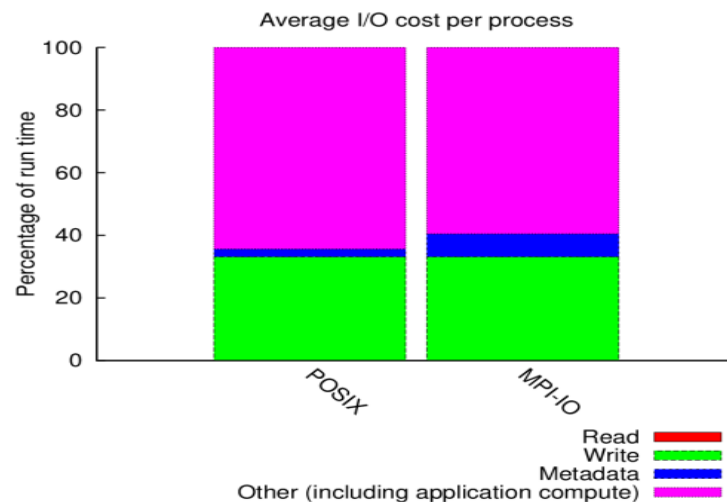
# DARSHAN

- I/O characterization tool logging parallel application file access
- Summary report provides quick overview of performance issues
- Works on unmodified, optimized executables
- Shows counts of file access operations, times for key operations, histograms of accesses, etc.
- Supports POSIX, MPI-IO, HDF5, PnetCDF, ...
- Binary log file written at exit post-processed into PDF report
- <http://www.mcs.anl.gov/research/projects/darshan/>
- Open Source: installed on many HPC systems



# EXAMPLE DARSHAN REPORT EXTRACT

jobid: | uid: | nprocs: 4096 | runtime: 175 seconds





# PERFORMANCE ANALYSIS RECOMMENDATIONS

- Measure and analyze at the desired scale (once you have a reasonable measurement setup)
- Get performance overview with Performance Reports
  - CPU Issues:
    - Use Vtune (on Intel nodes) or uProf (on AMD nodes)
    - Use perf / LIKWID / PAPI
  - MPI Issues: Use Scalasca/Vampir
  - GPU Issues: Use NVIDIA tools
  - I/O Issues: Use DARSHAN
- OR: Do it all with Score-P/Scalasca/Vampir

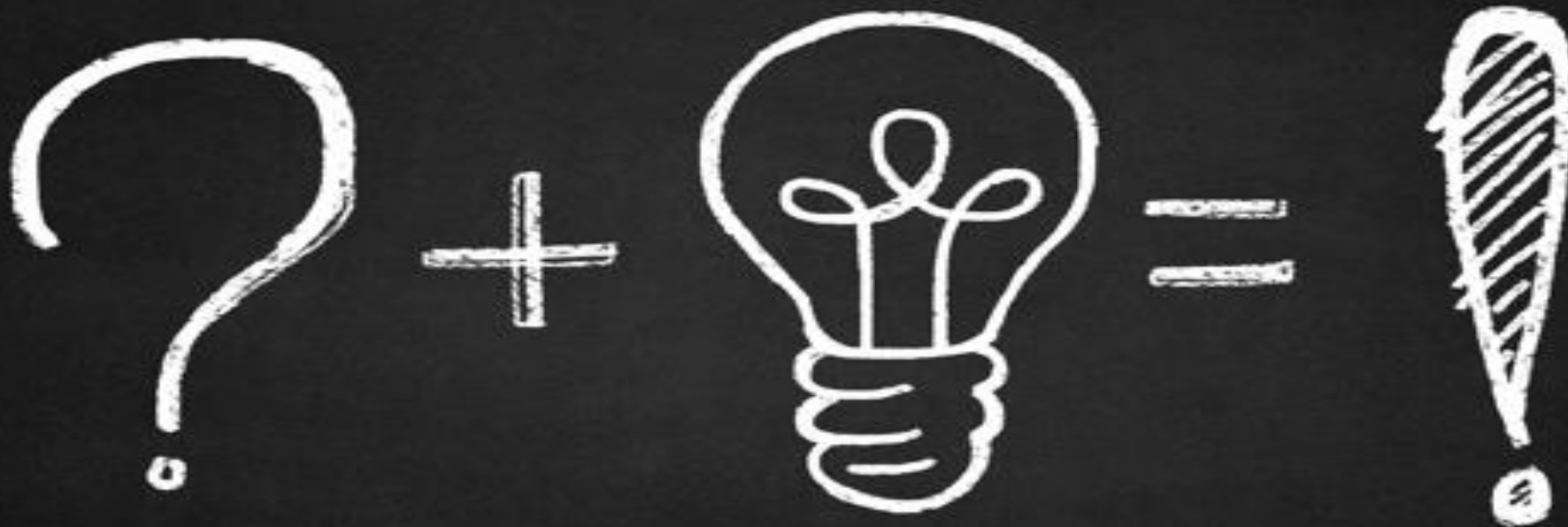


# NEED HELP?

- Talk to the experts
  - Use local 1<sup>st</sup>-level support, e.g. SC support, SimLab
  - Use mailing lists
  - JSC/NVIDIA Application Lab
  - ATML Parallel Performance
  - ATML Application Optimization and User Service Tools

👉 Successful performance engineering often is a collaborative effort





# QUESTIONS