



# The Scalasca Scalable Parallel Performance Analysis Toolset - For POP Assessments and Beyond

Bernd Mohr

EU H2020 Centre of Excellence (CoE)



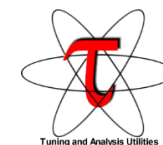
Grant Agreement No 824080

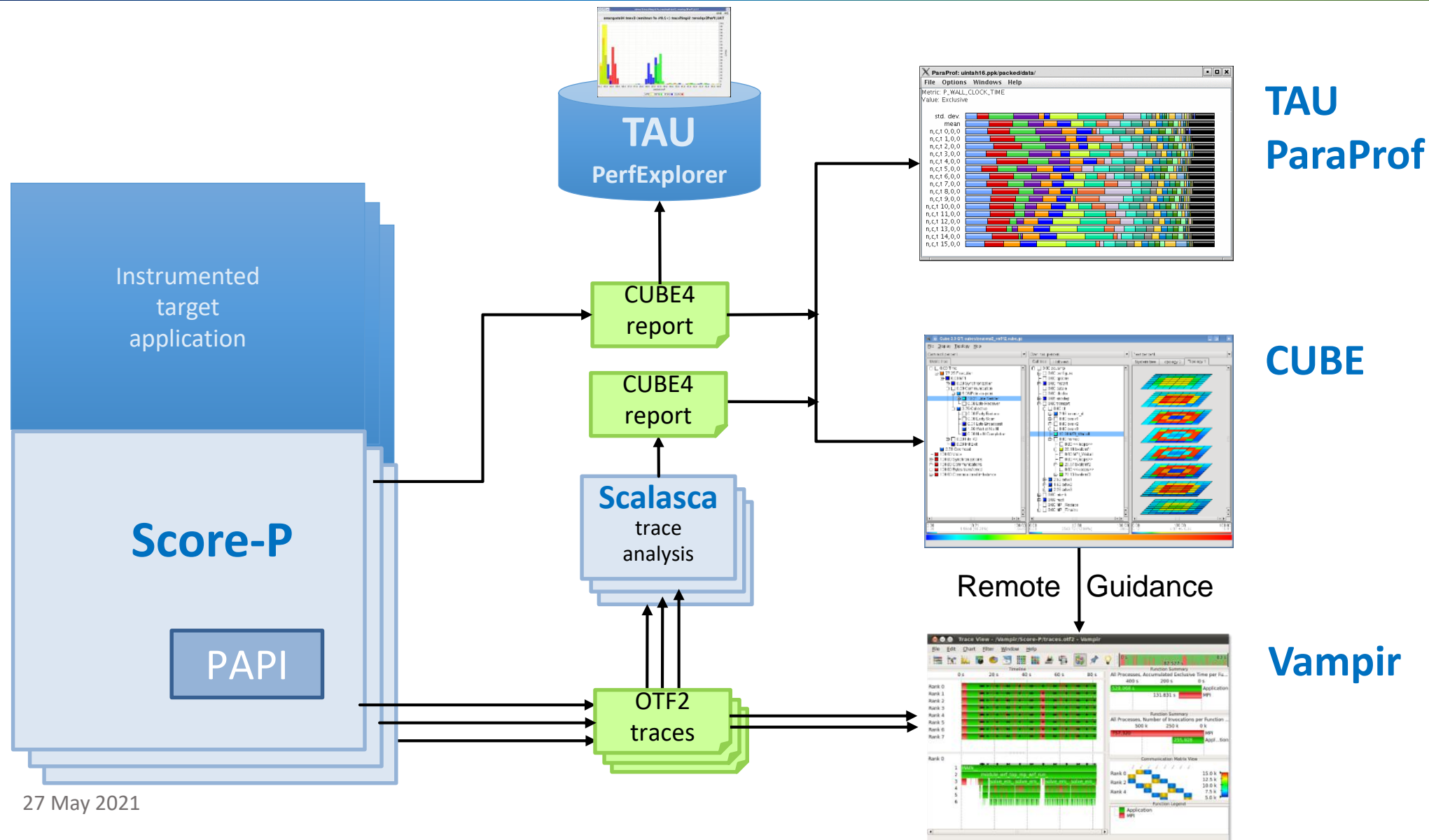
1 December 2018 – 30 November 2021



# The Score-P Tool Ecosystem

The Context



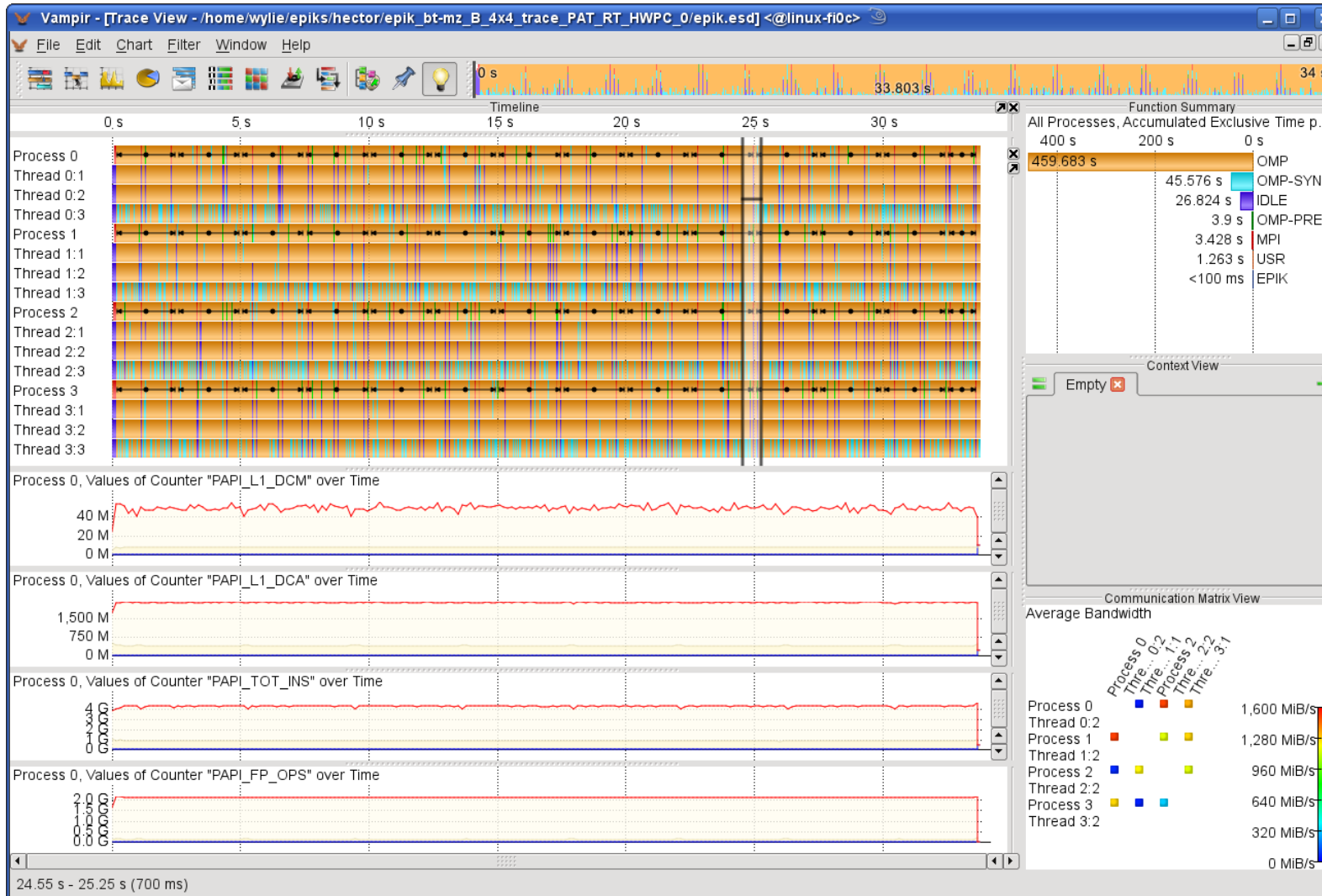




<https://score-p.org/>

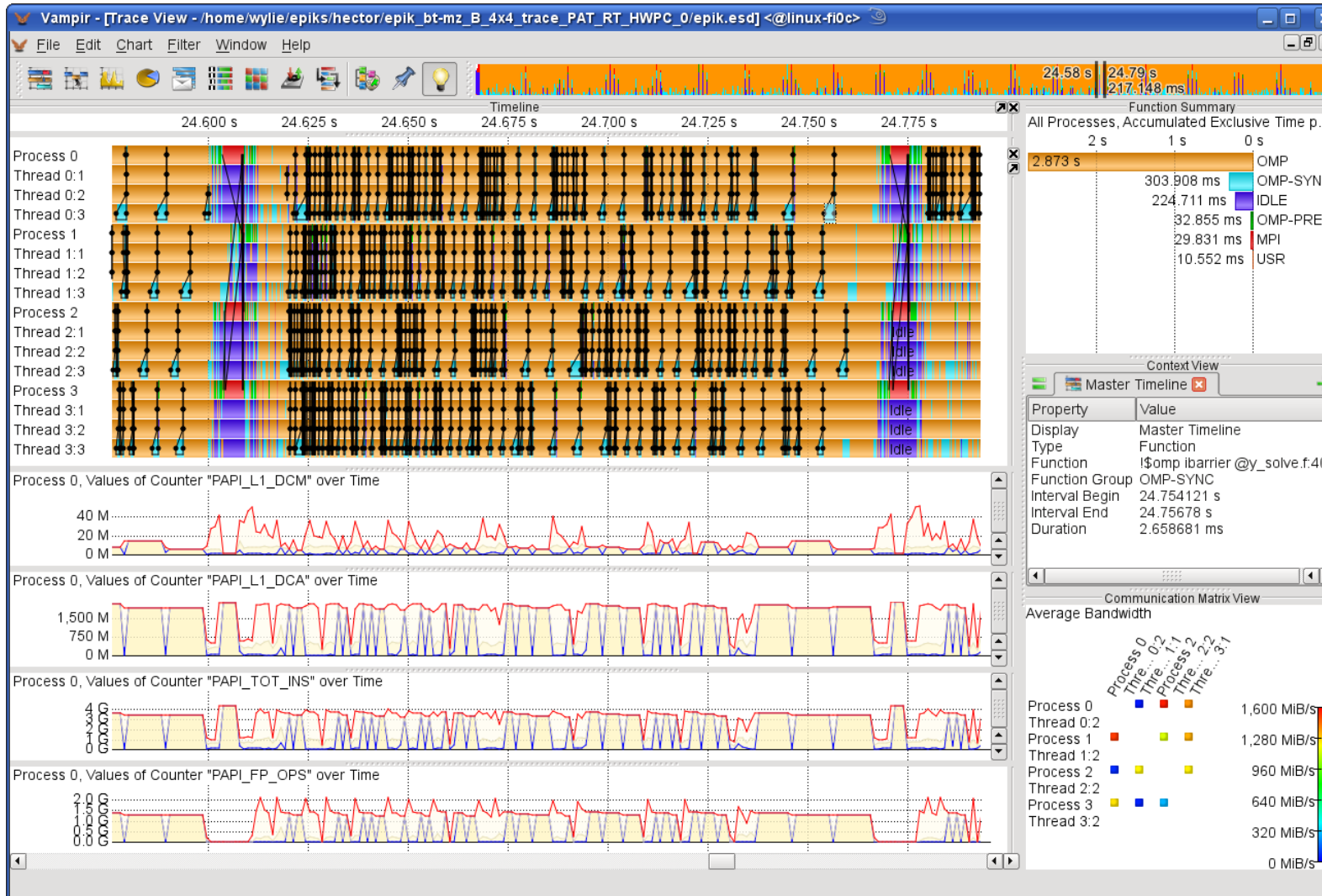
- Provide typical base functionality for HPC performance tools
- C, C++, Fortran, Python support
- **Instrumentation** (various methods)
  - Multi-process paradigms (MPI, SHMEM)
  - Thread-parallel paradigms (OpenMP, POSIX threads)
  - Accelerator-based paradigms (OpenACC, CUDA, OpenCL, Kokkos)
  - Additional execution information (HW+SW counter, I/O, memory, ...)
  - **In any combination!**
- Flexible **measurement** without re-compilation:
  - Basic and advanced **profile** generation ( $\Rightarrow$  CUBE4 format)
  - Event **trace** recording ( $\Rightarrow$  OTF2 format)
- Highly scalable measurement





## Visual presentation of dynamic runtime behaviour

- Event timeline chart for states & interactions of processes/threads
- Communication statistics, summaries & more



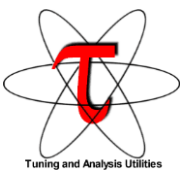
**Interactive browsing,  
zooming, selecting**

- Linked displays & statistics adapt to selected time interval

**Trace formats**

- OTF (VampirTrace)
- OTF2 (Score-P)
- EPIK (Scalasca1)



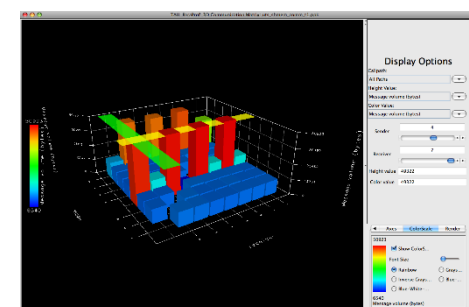
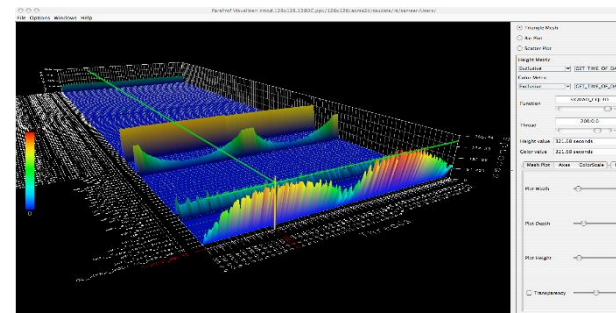
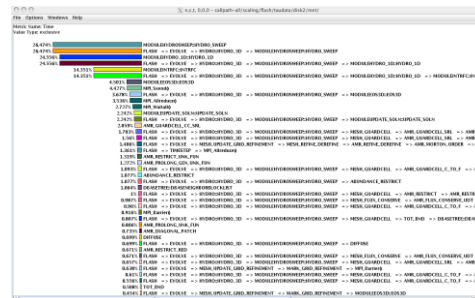
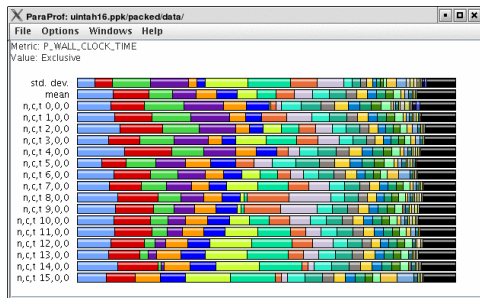


# TAU Performance System®



<http://tau.uoregon.edu/>

- **Very portable tool set** for instrumentation, measurement and analysis of parallel multi-threaded applications
- Supports
  - **Various profiling modes and tracing**
  - Various forms of **code instrumentation + sampling**
  - C, C++, Fortran, Java, Python
  - Multi-process paradigms (MPI, SHMEM, GPI, ARMCI)
  - Thread-parallel paradigms (OpenMP, POSIX threads)
  - Accelerator-based paradigms (OpenACC, CUDA, OpenCL, HIP/ROCm, OneAPI, Kokkos)

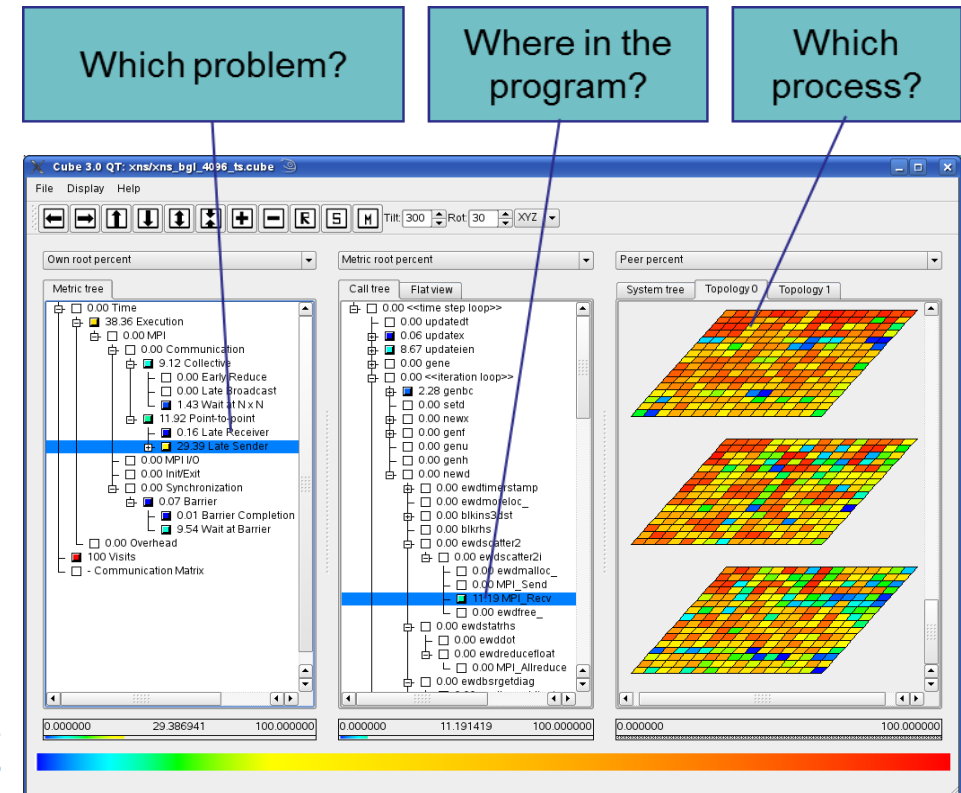


- Scalable Analysis of Large Scale Applications

<https://scalasca.org/>

- 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



CUBE



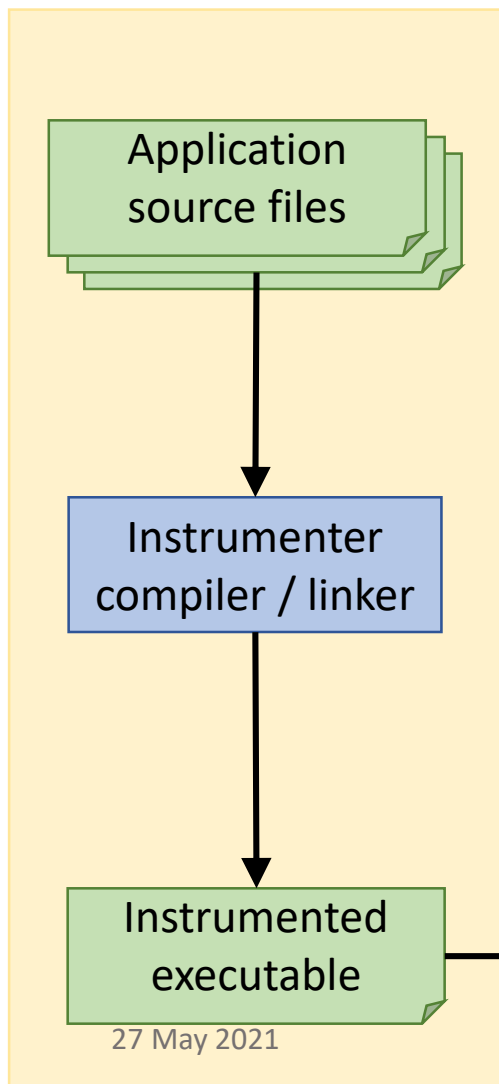


# The Scalasca Workflow

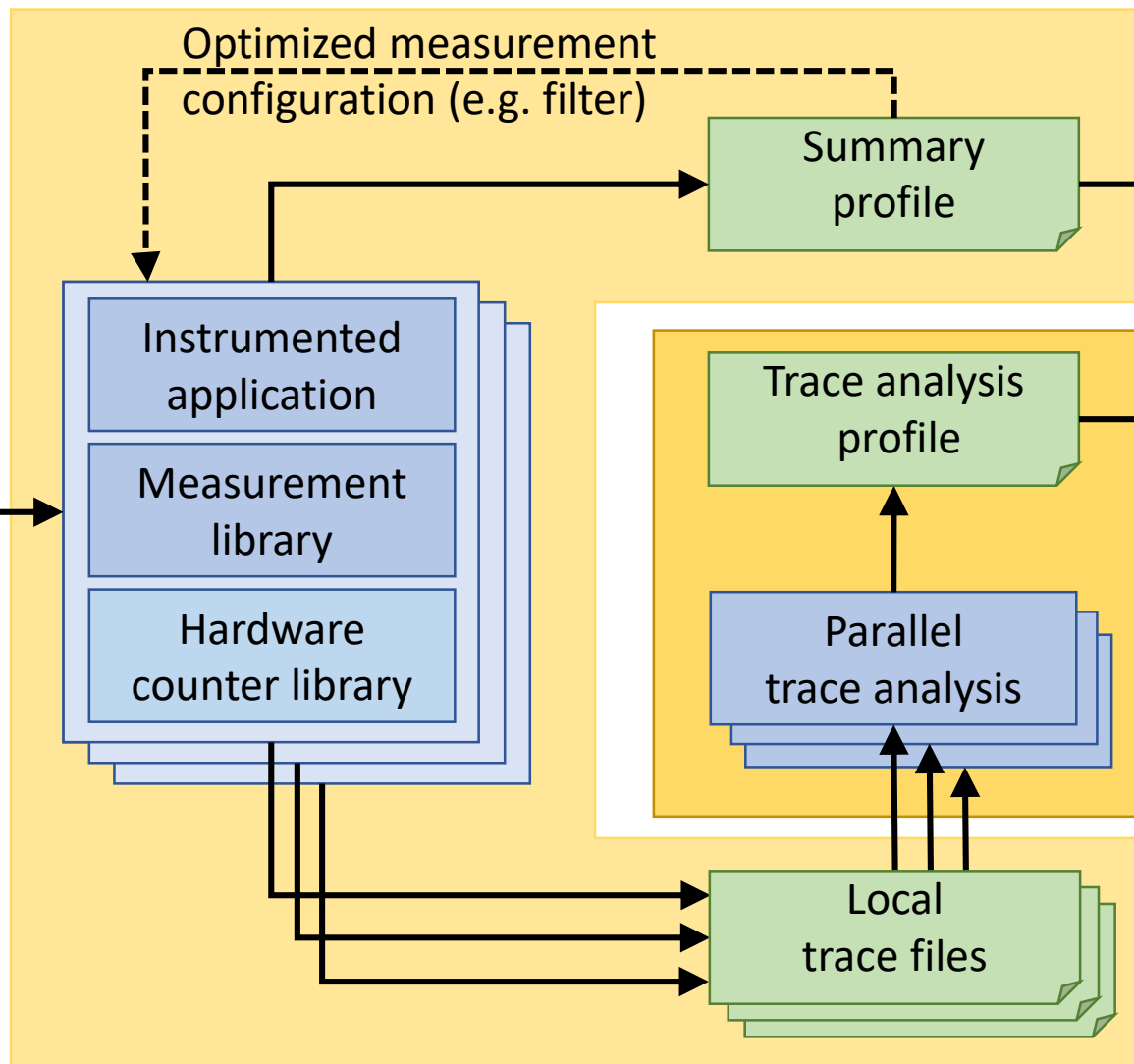
The Basics



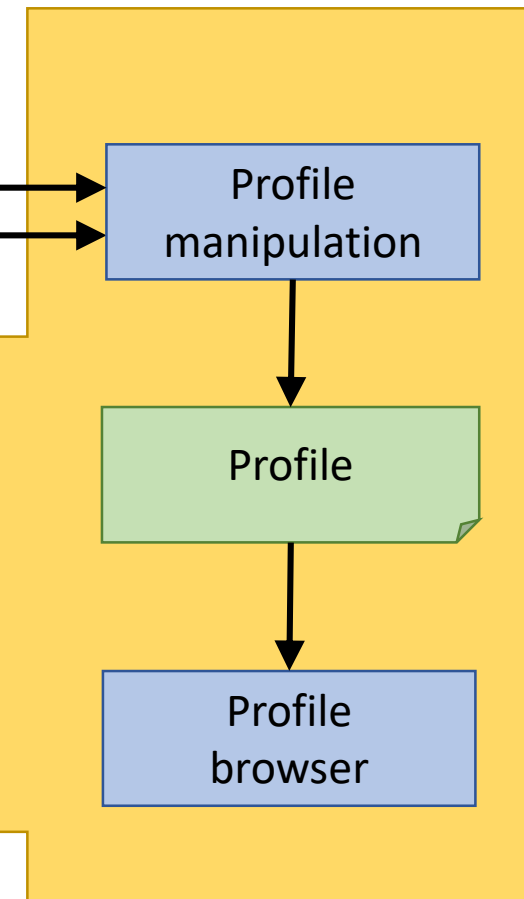
## 1. Instrumentation



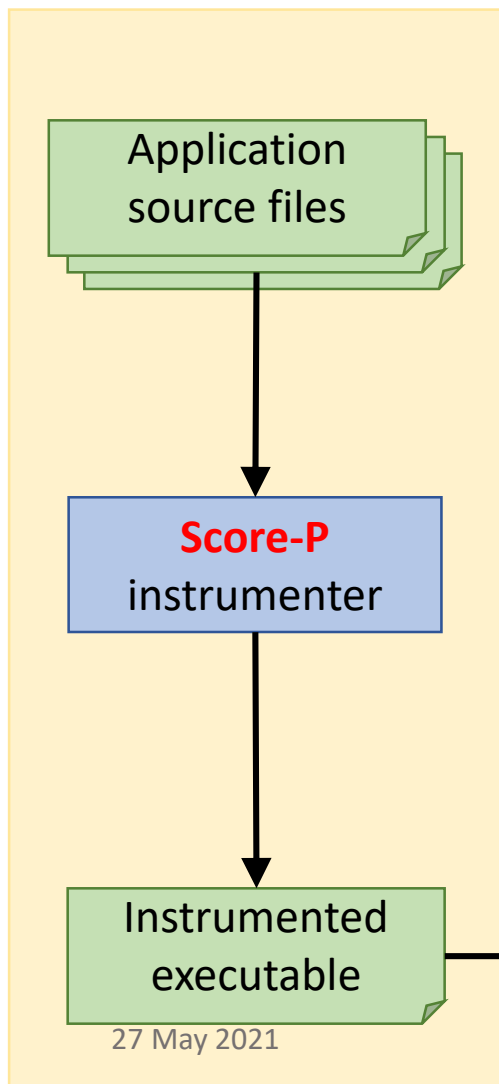
## 2. Measurement



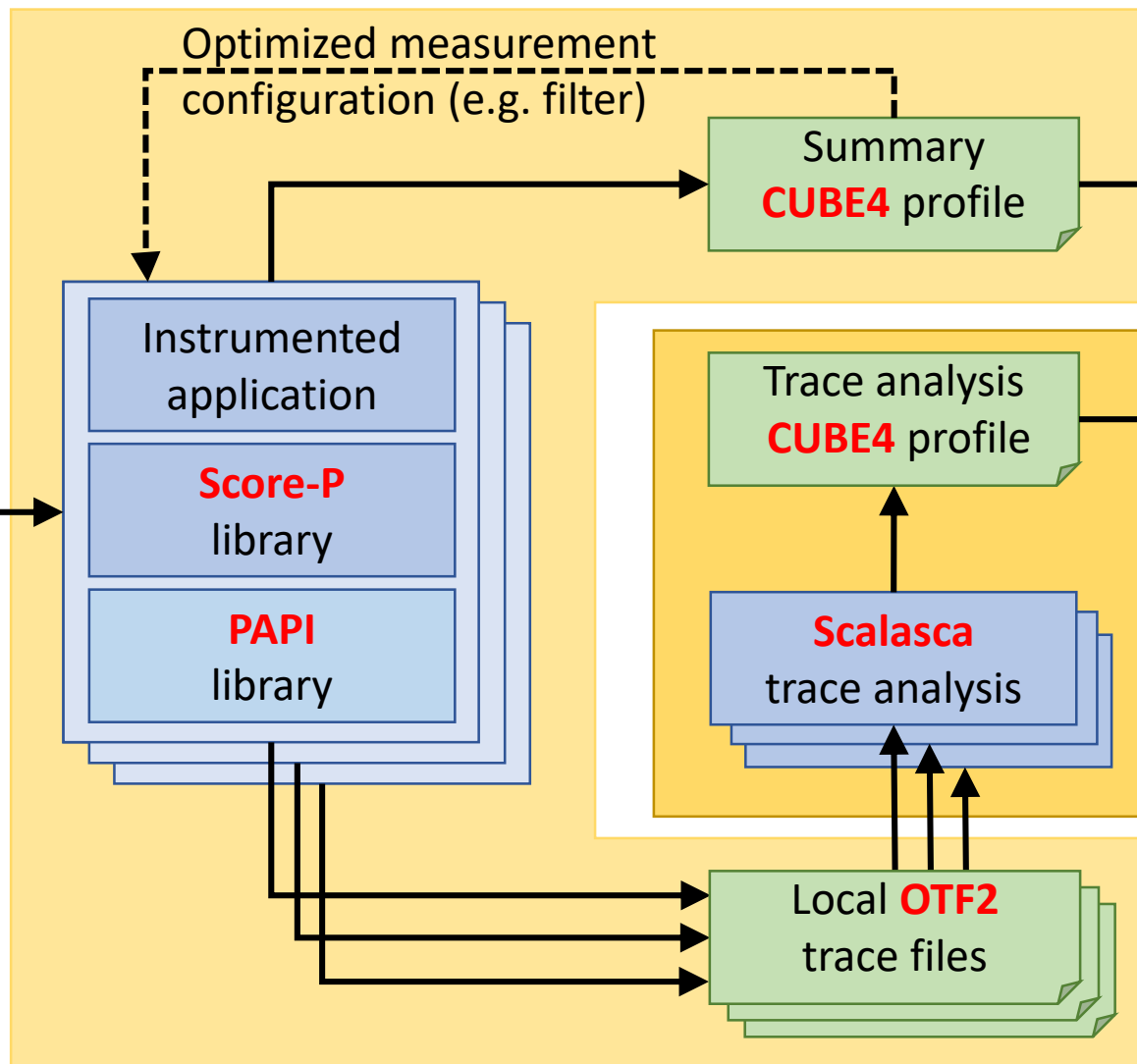
## 3. Analysis



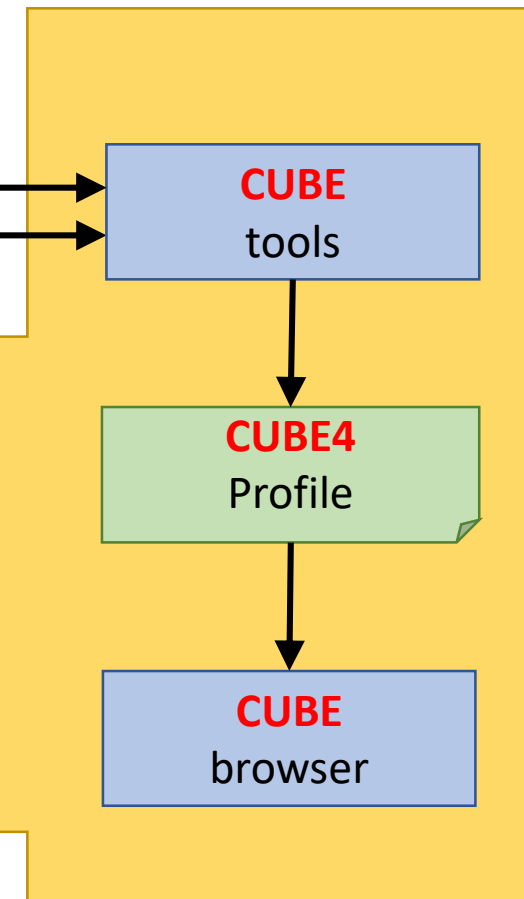
## 1. Instrumentation



## 2. Measurement



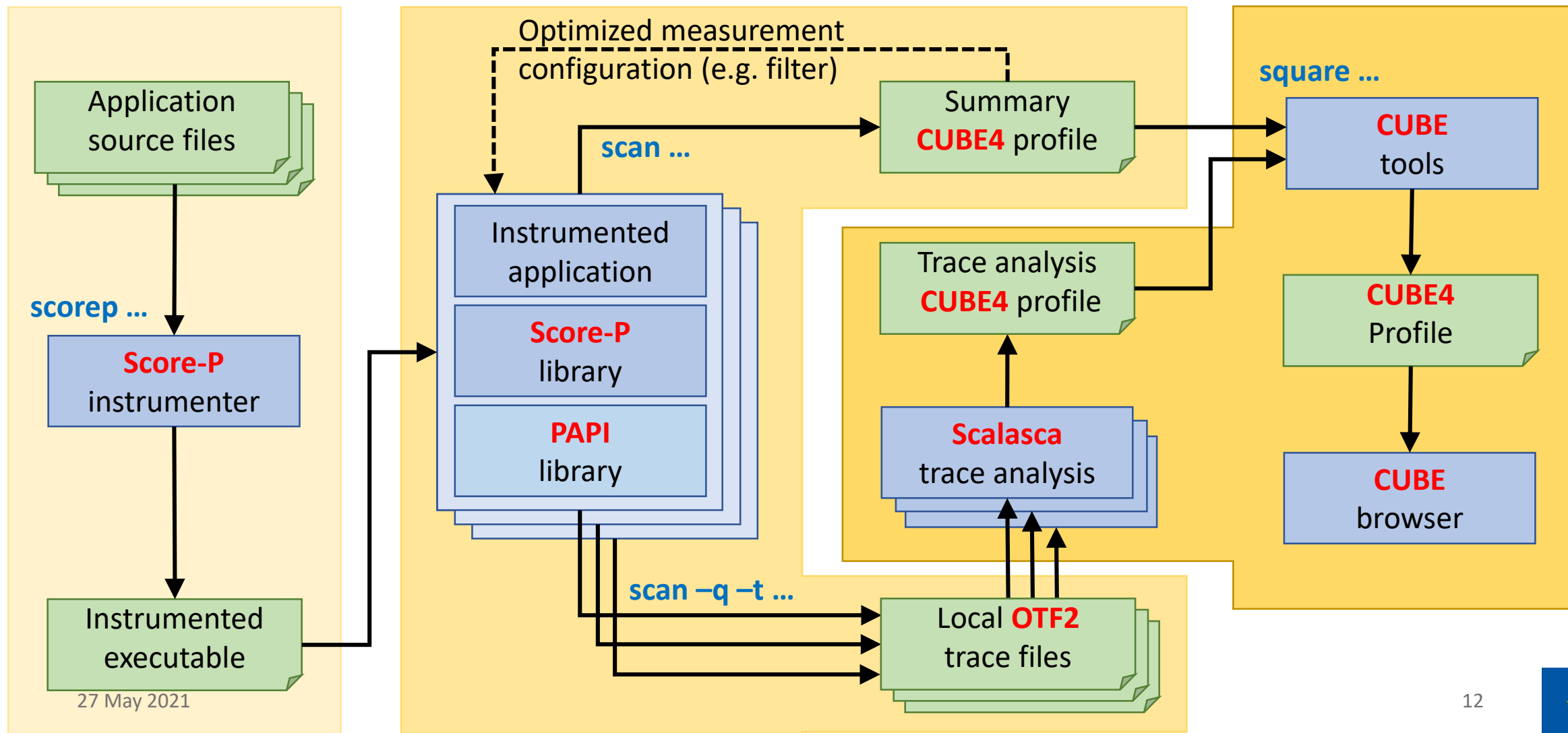
## 3. Analysis



## 1. Instrumentation

## 2. Measurement

## 3. Analysis





# POP Assessments with Scalasca

This is the way



<https://pop-coe.eu/further-information/learning-material>

- **Original (POP1) Metrics**

- [Article](#) explaining the POP Standard Metrics for Parallel Performance Analysis
- [Presentation](#) summarizing the POP Standard Metrics for Parallel Performance Analysis

- **New (POP2) Hybrid Metrics**

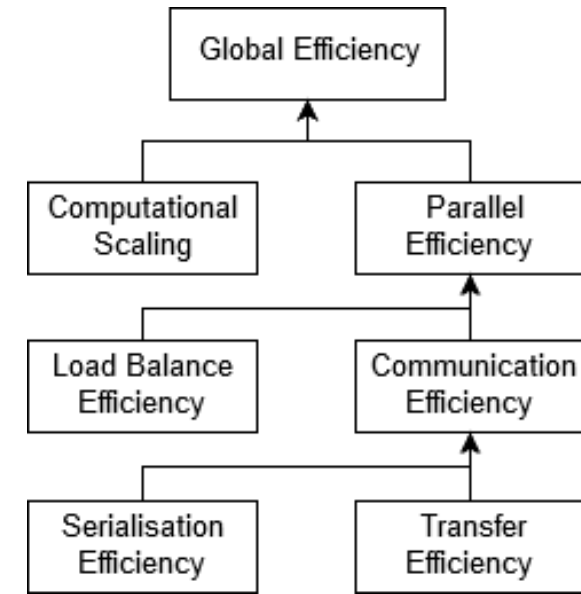
- • [Introduction](#) explaining the POP2 Standard Metrics for Performance Analysis of Hybrid Parallel Applications
- [Cheat sheet](#) for Additive Hybrid Metrics
- [Cheat sheet](#) for Multiplicative Hybrid Metrics
- [In-depth explanation](#) of the Additive Hybrid Metrics
- • [Webinar](#) Identifying Performance Bottlenecks in Hybrid MPI + OpenMP Software



# Recap: POP Phase1 (MPI) Metrics



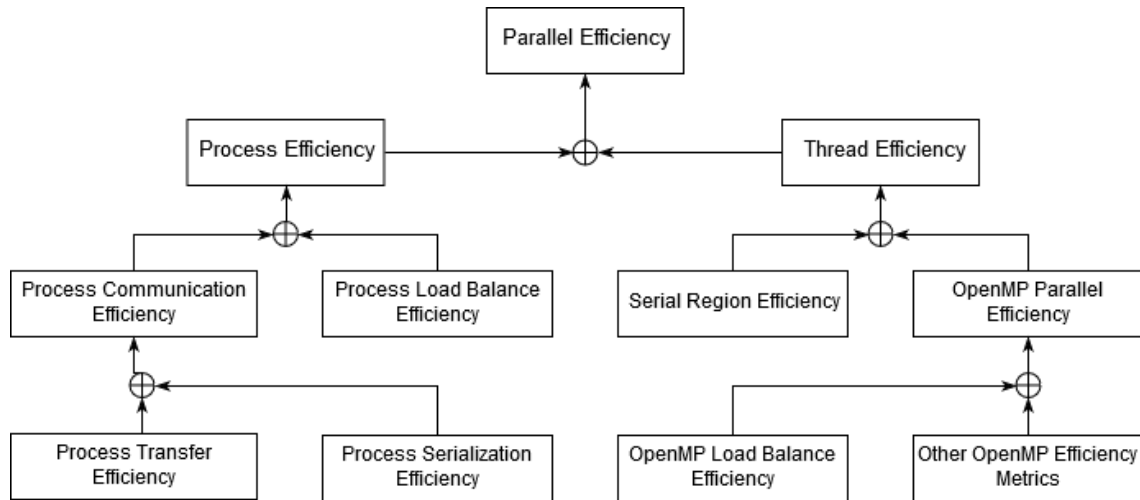
- The following metrics are used in a POP Performance Audit:
- Global Efficiency (GE):  $GE = PE * CompE$ 
  - Parallel Efficiency (PE):  $PE = LB * CommE$ 
    - **Load Balance** Efficiency (LB):  $LB = avg(CT)/max(CT)$
    - **Communication** Efficiency (CommE):  $CommE = SerE * TE$ 
      - Serialization Efficiency (SerE):  
 $SerE = max(CT / TT \text{ on ideal network})$
      - Transfer Efficiency (TE):  $TE = TT \text{ on ideal network} / TT$
  - **Computation** Scaling (CompS)
    - Computed out of IPC Scaling and Instruction Scaling
    - For strong scaling: ideal scaling -> efficiency of 1.0



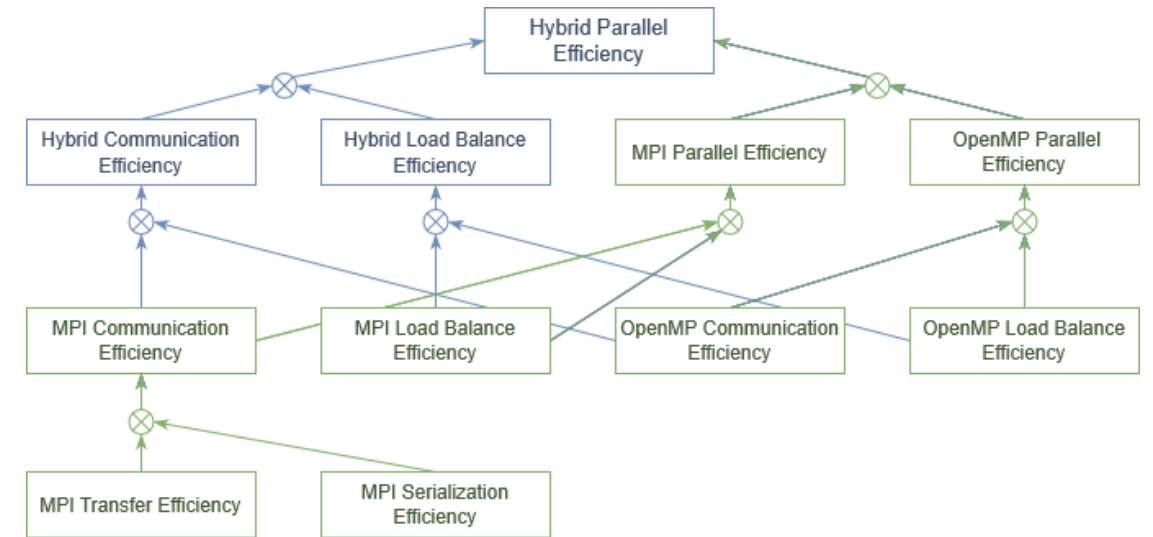
CT = Computational time  
TT = Total time



# Recap: POP Phase2 (Hybrid) Metrics



Additive Version



Multiplicative Version



1. Instrument application and setup measurement parameters (e.g. filtering)
  - `scorep <comp+link+cmds>`
  - `scan <exec+cmd> ...`
2. For parallel efficiency: perform trace measurement and analysis
3. For computational scaling: perform profile measurement with suitable HW counters
  - `scan -P pop <exec+cmd>`
4. Merge profile and trace measurement
5. Post-process measurement
6. Analyze POP metrics with **Cube Advisor**
  - `square <measurement+archive>`

for scaling, repeat

Requires

- Scalasca V2.6
- Cube 4.6



# POP Metrics + Scalasca



**Select POP metric set**

**Select Region of Interest**

**Find Advisor on General Tab**

CubeGUI-4.6.0: scorep\_tea\_leaf\_16p32x8\_multi-run\_c2\trace+su

File Display Plugins Help

NAG POP Hybrid Assessment (Additive) Runtime threshold: 2.5 % runtime

Metric tree

- 0.00 Time (sec)
  - 2460.67 Execution
    - 0.00 Overhead
    - 1325.66 Idle threads
    - 3.64e8 Visits (occ)
    - 288 MPI synchronizations (occ)
    - 0 MPI pair-wise one-sided synchronizations (occ)
    - 6.26e6 MPI communications (occ)
    - 0 MPI file operations (occ)
    - 6.01e10 MPI bytes transferred (bytes)
    - 1484.10 Delay costs (sec)
    - 6.21 MPI point-to-point wait states (propagating vs. terminal) (sec)
    - 6.21 MPI point-to-point wait states (direct vs. indirect) (sec)
    - 14.86 Critical path (sec)
    - 3804.82 Performance impact (sec)
    - 95.68 Computational imbalance (sec)
    - 0.00 Minimum Inclusive Time (sec)
    - 18.88 Maximum Inclusive Time (sec)
    - 0 io\_bytes\_read (bytes)
    - 0 io\_bytes\_written (bytes)
    - 6.08e12 PAPI\_TOT\_INS (#)
    - 6.59e12 PAPI\_TOT\_CYC (#)

Absolute

Call tree Flat tree

- 0.02 tea\_leaf.inst
  - 0.00 tea\_leaf
    - 73.59 tea\_init\_comms
    - 0.57 !\$omp parallel @tea\_leaf.f90:45
    - 5.21 initialise
    - 0.00 diffuse
      - 0.00 timer
      - 0.06 set\_field
      - 0.03 timeslep
      - 1.82 tea\_leaf
        - 1.00 timer
        - 540.69 update\_halo
        - 3.96 tea\_leaf\_kernel\_init\_cg\_fortran
        - 58.70 tea\_allsum
        - 0.93 tea\_leaf\_kernel\_cheby\_copy\_u
        - 779.35 tea\_leaf\_kernel\_solve\_cg\_fortran\_calc\_w
        - 645.25 tea\_leaf\_kernel\_solve\_cg\_fortran\_calc\_ur
        - 347.08 tea\_leaf\_kernel\_solve\_cg\_fortran\_calc\_p
        - 2.07 tea\_leaf\_kernel\_finalise
      - 0.30 field\_summary
      - 0.02 tea\_allgather
      - 0.01 tea\_finalize

Advisor Source

Recalculate ☐ automatic ☐ direct calculation

POP Hybrid Assessment ( Additive ) : tea\_leaf\_module::tea\_leaf

Parallel Efficiency	0.61	<div></div>	Good	?
+ Process Efficiency	0.67	<div></div>	Good	?
++ Computation Load Balance	0.78	<div></div>	Good	?
++ Communication Efficiency	0.89	<div></div>	Very good	?
+++ Serialisation Efficiency	0.91	<div></div>	Very good	?
+++ Transfer Efficiency	0.98	<div></div>	Very good	?
+ Thread Efficiency	0.94	<div></div>	Very good	?
+++ Amdahl Efficiency	0.96	<div></div>	Very good	?
+++ OpenMP Region Efficiency	0.99	<div></div>	Very good	?
Resource stall cycles				!
IPC	0.93		Value	?
Instructions (only computation)	5.14e12		Value	?
Computation time	1733.75		Value	?

Candidates

Callpath	Issue
----------	-------

System View Topologies General

0.00 2460.67 (64.99%) 3786.33

0.00 1.82 (0.07%) 2460.67

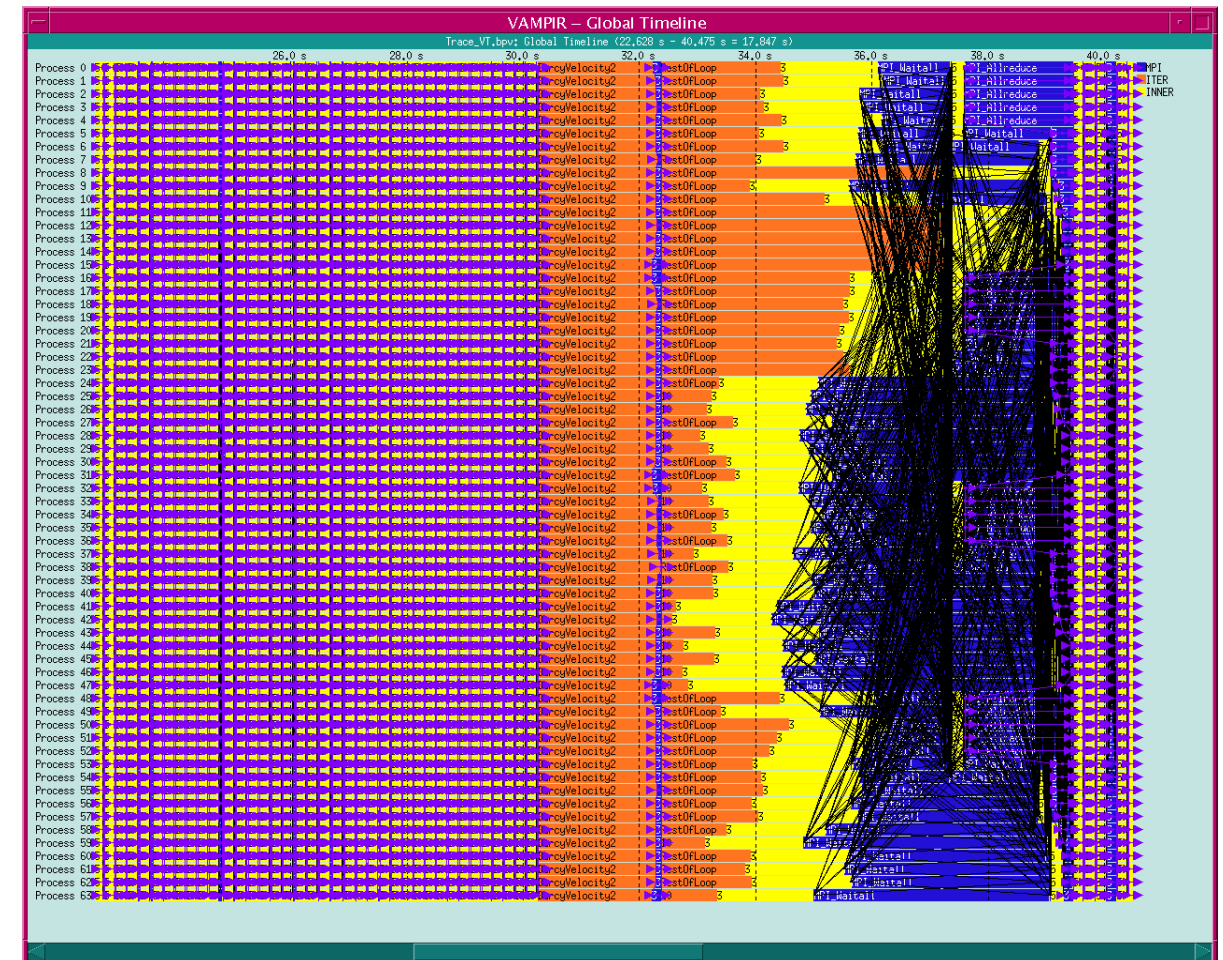


# (Detailed) Performance Analysis with Scalasca

Finding the cause(s)



# “A picture is worth 1000 words...”



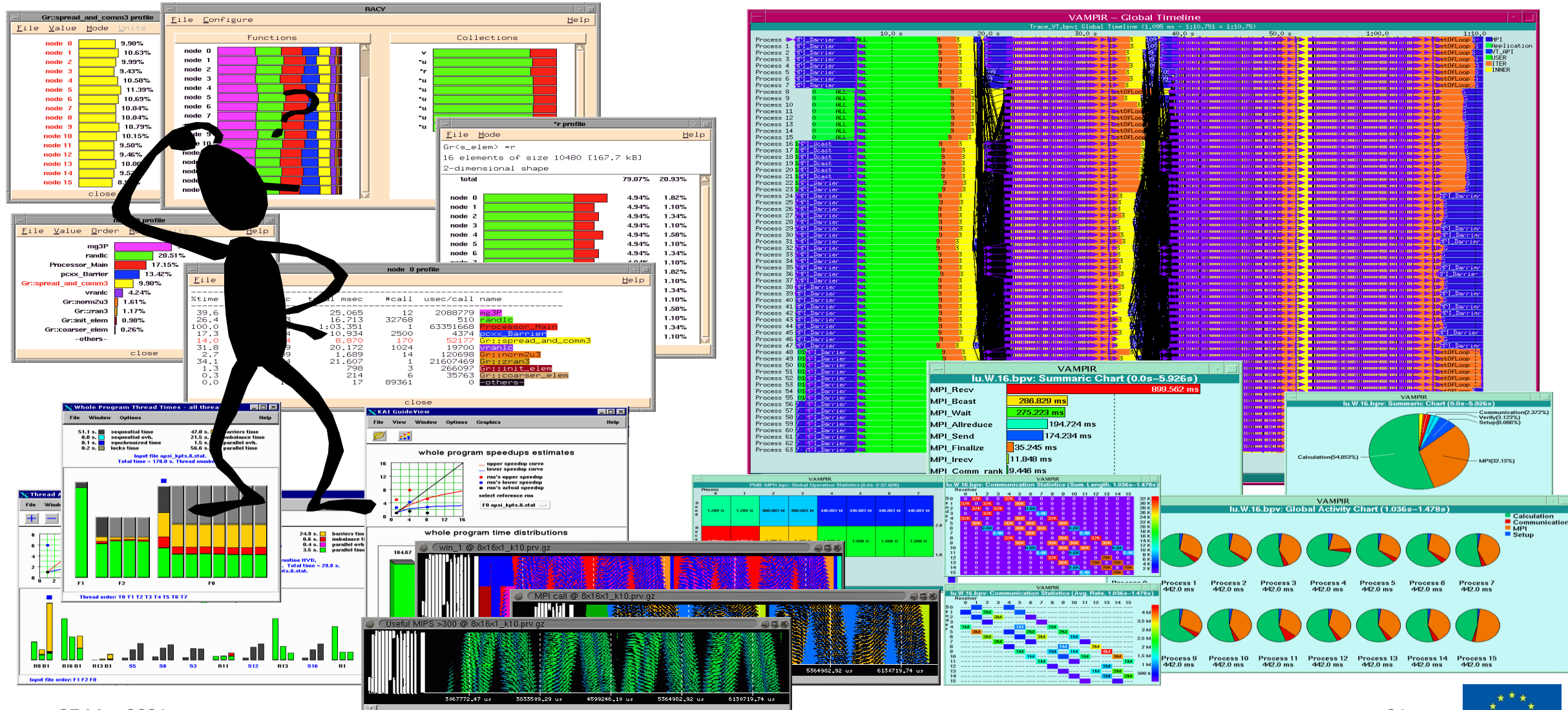
- MPI ring program

- “Real world” example

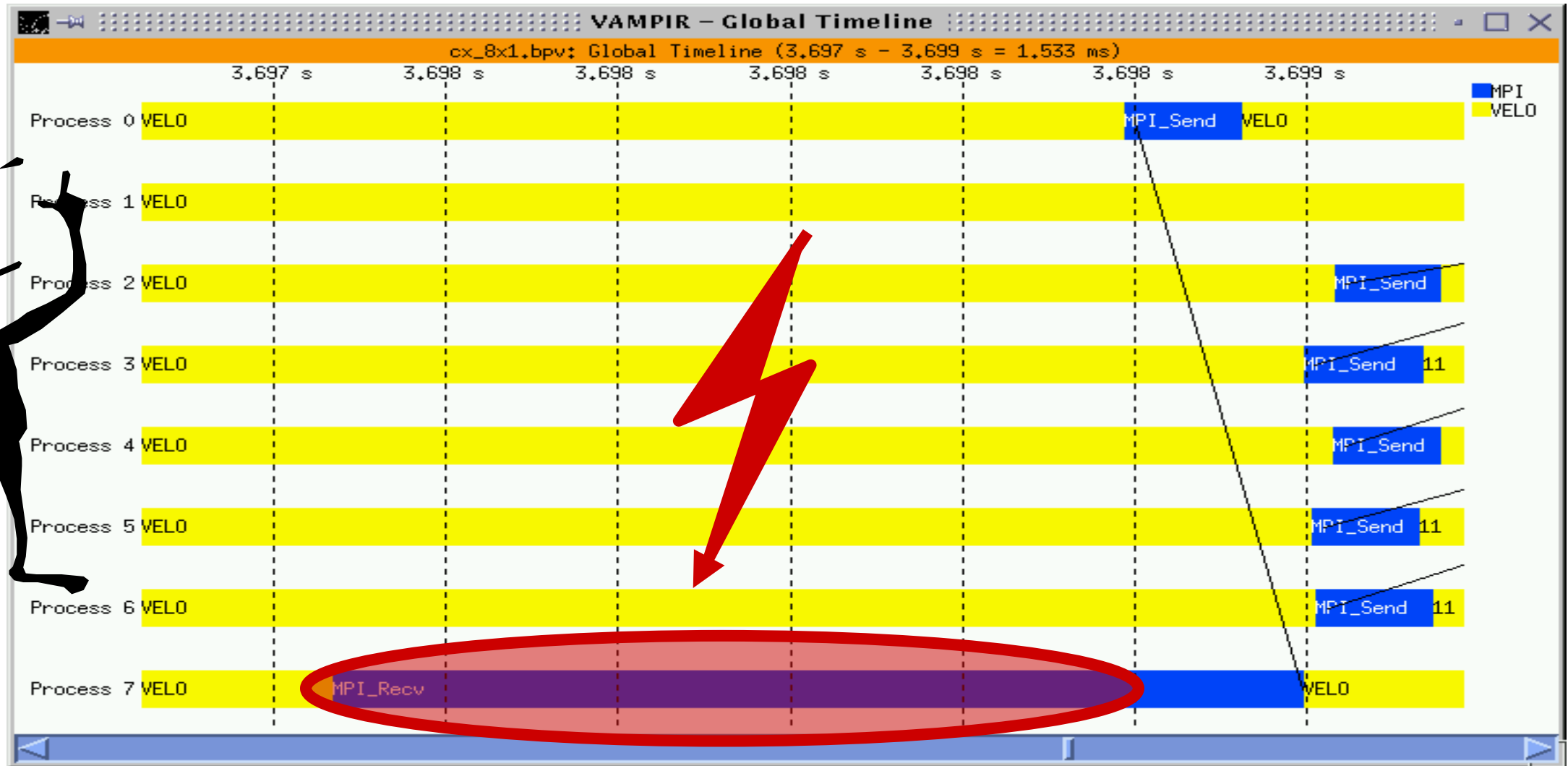
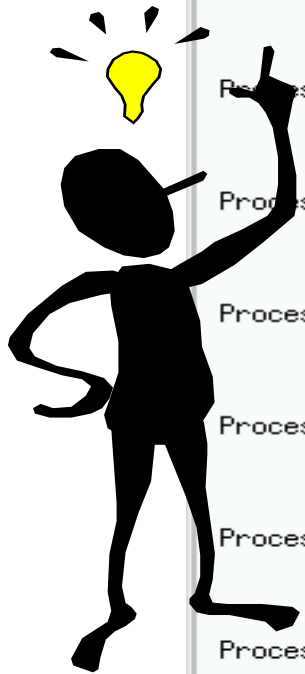




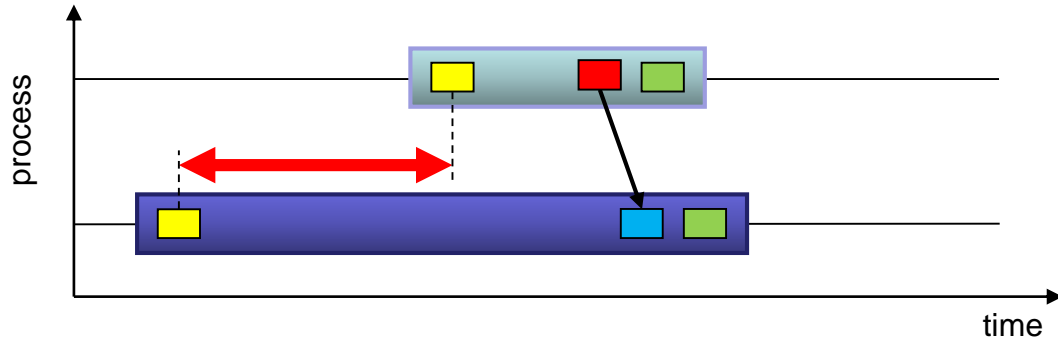
# “What about 1000’s of pictures?” (with 100’s of menu options)



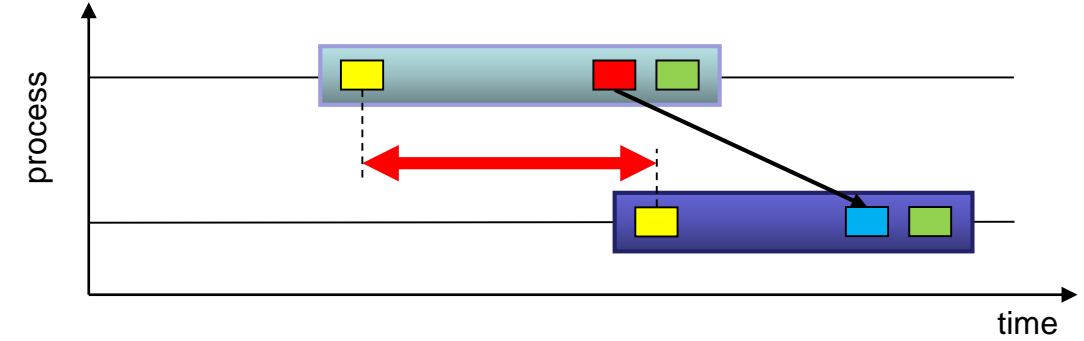
# Example Automatic Analysis: Late Sender



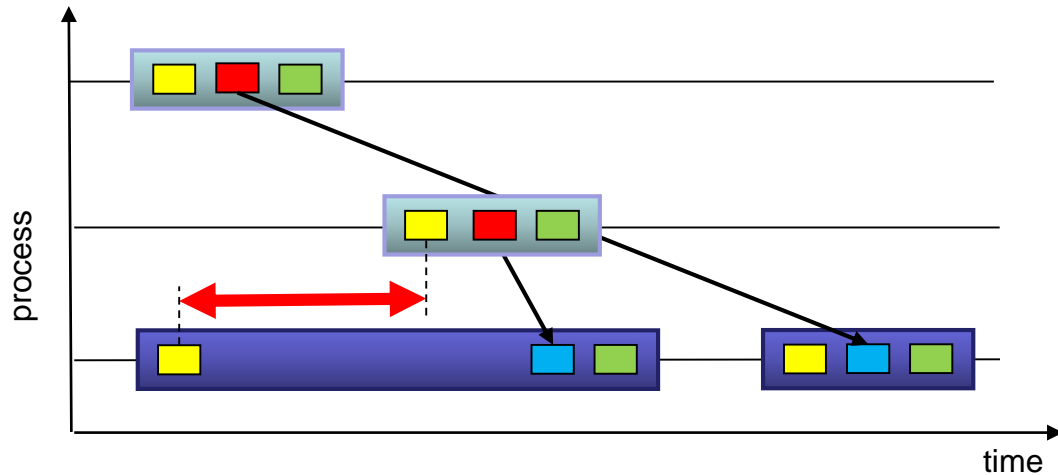
# Example MPI Wait States



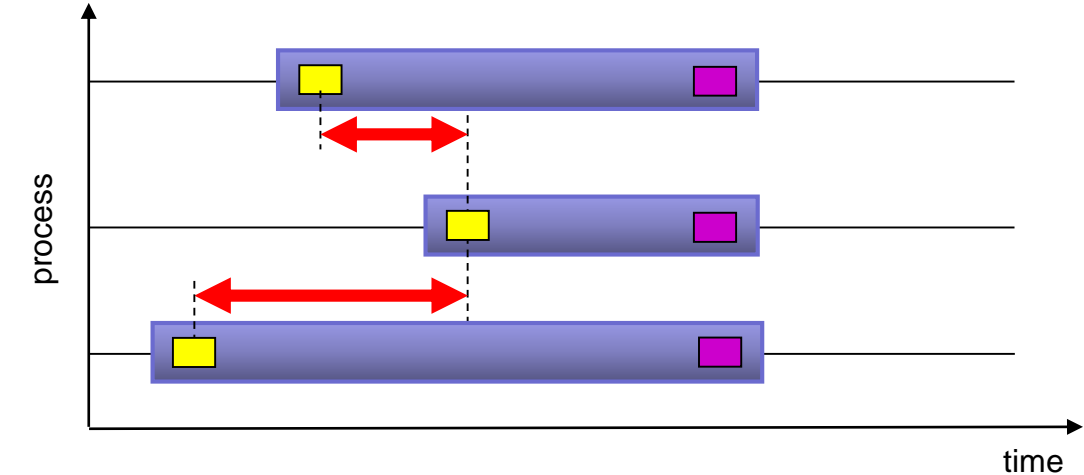
(a) Late Sender



(b) Late Receiver



(c) Late Sender / Wrong Order



(d) Wait at N x N

ENTER EXIT SEND RECV COLLEXIT

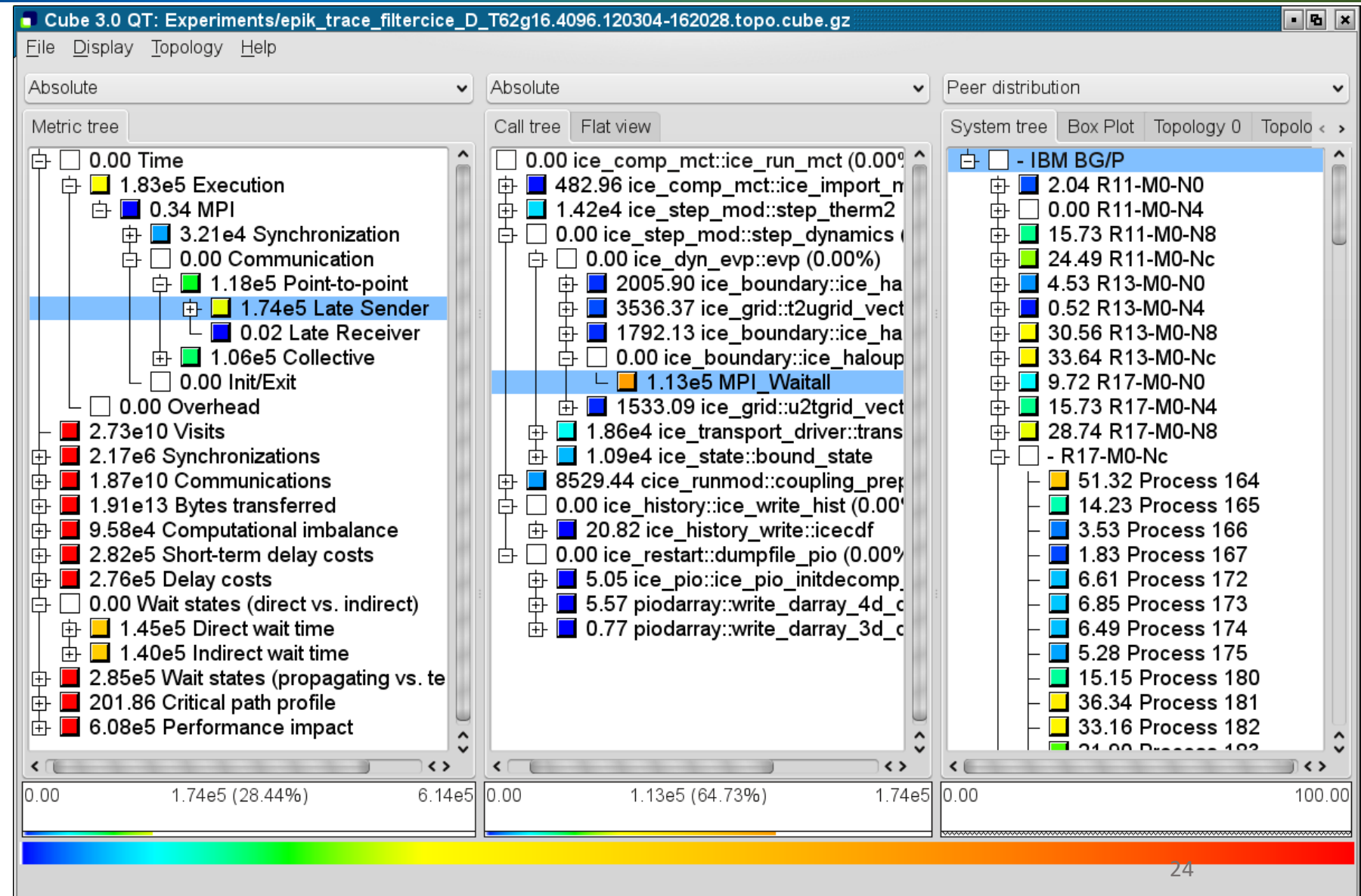


# Scalasca Example: CESM Sea Ice Module



## Late Sender Analysis

- Finds waiting at MPI\_Waitall() inside ice boundary halo update
- Shows distribution of imbalance across system and ranks

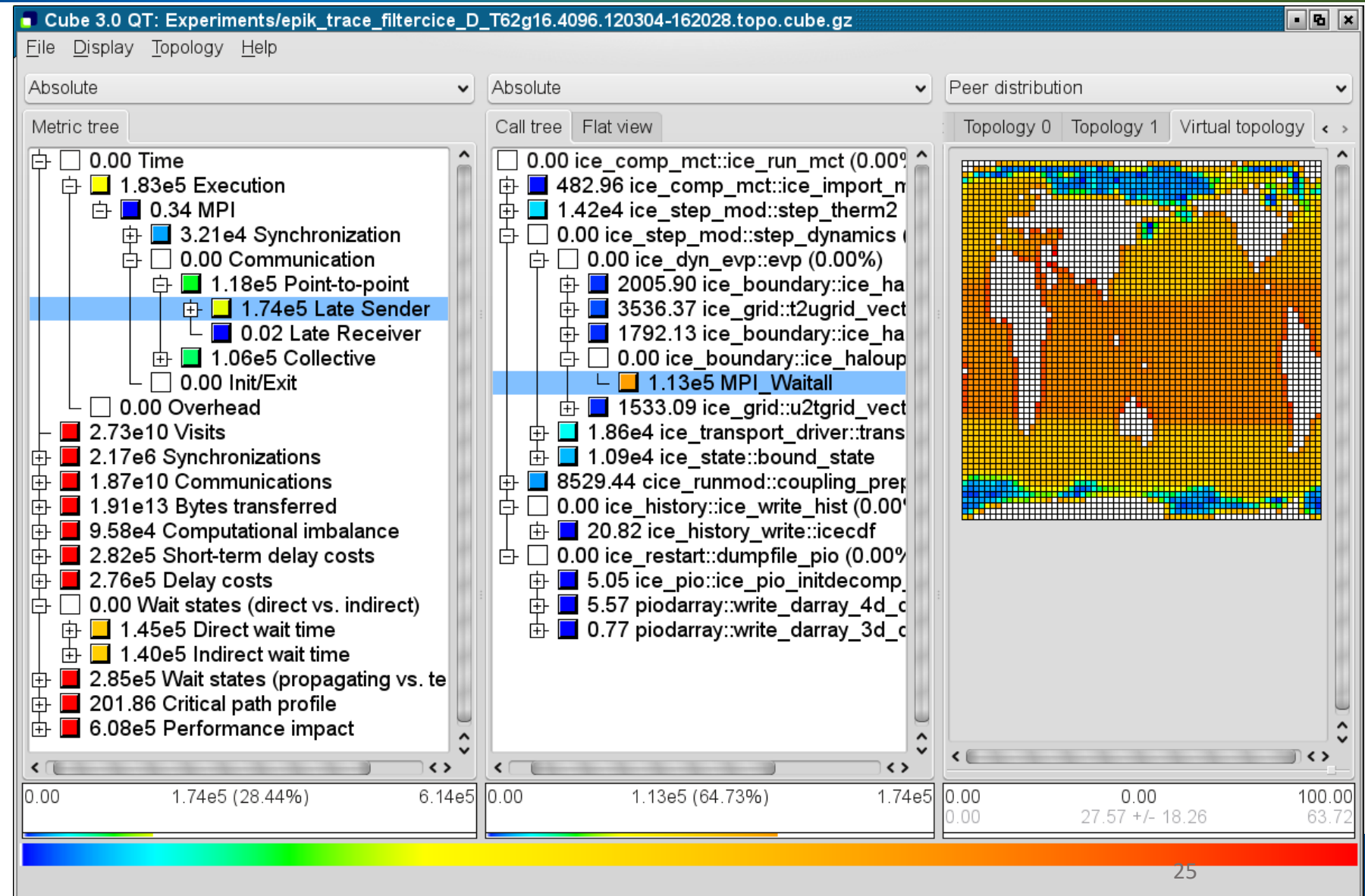


# Scalasca Example: CESM Sea Ice Module



## Late Sender Analysis + Application Topology

- Shows distribution of imbalance over topology
- MPI topologies are automatically captured
- Also: topology Process x Threads





# Scalasca Root Cause Analysis

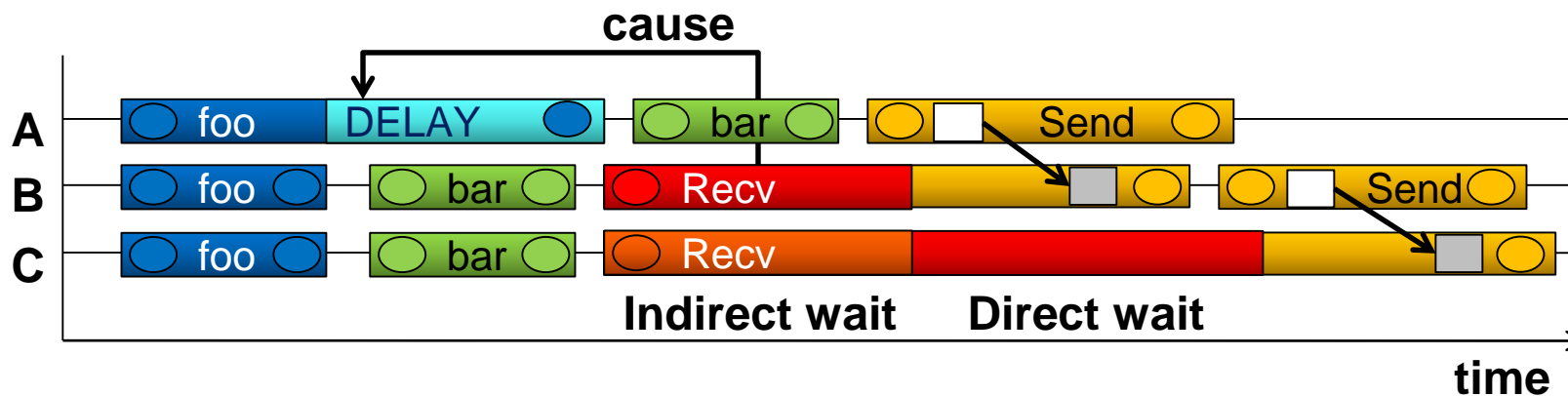


- **Root-cause analysis**

- 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**



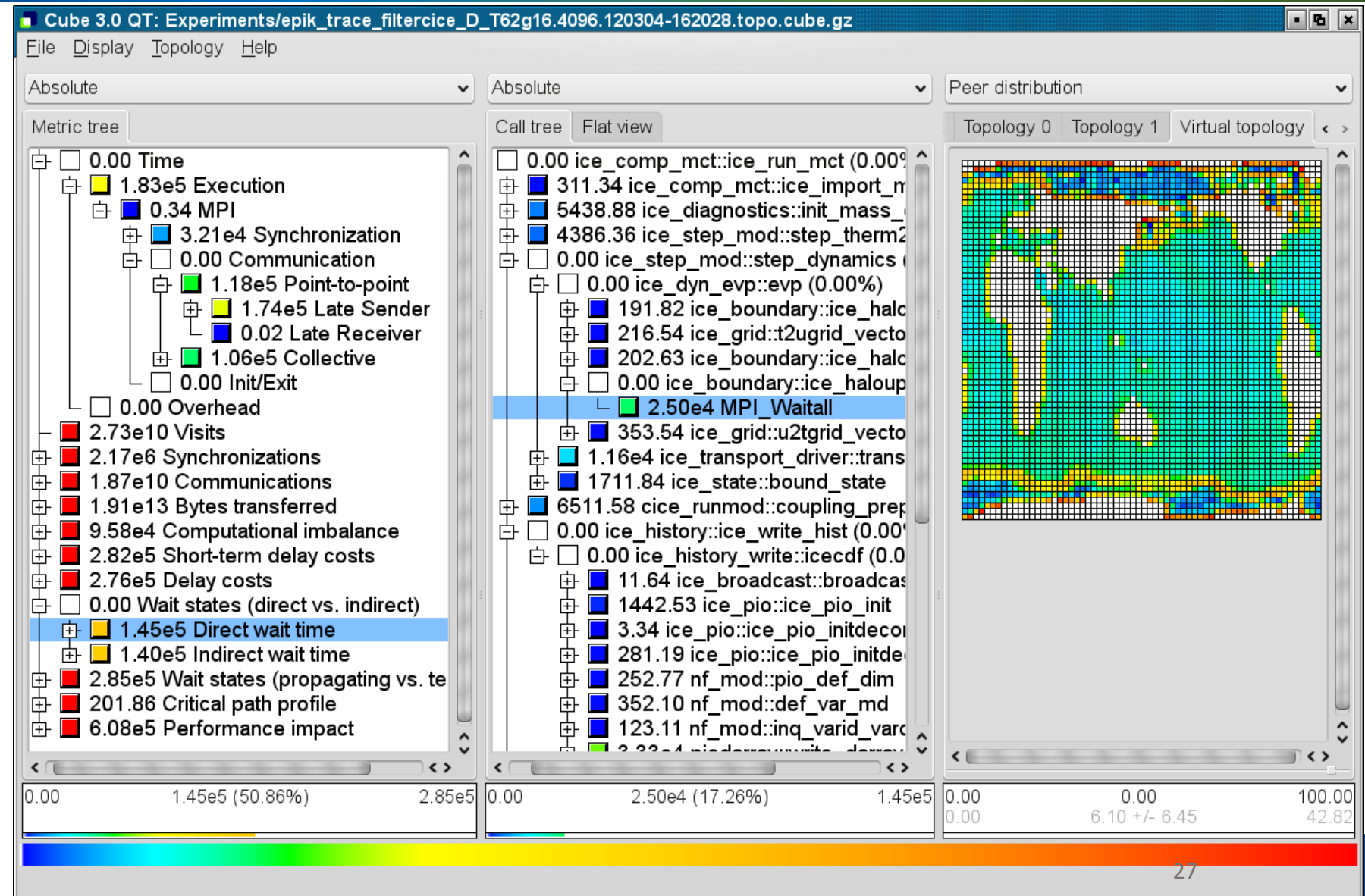


# Scalasca Example: CESM Sea Ice Module



## Direct Wait Time Analysis

- Direct wait caused by ranks processing areas near the north and south ice borders

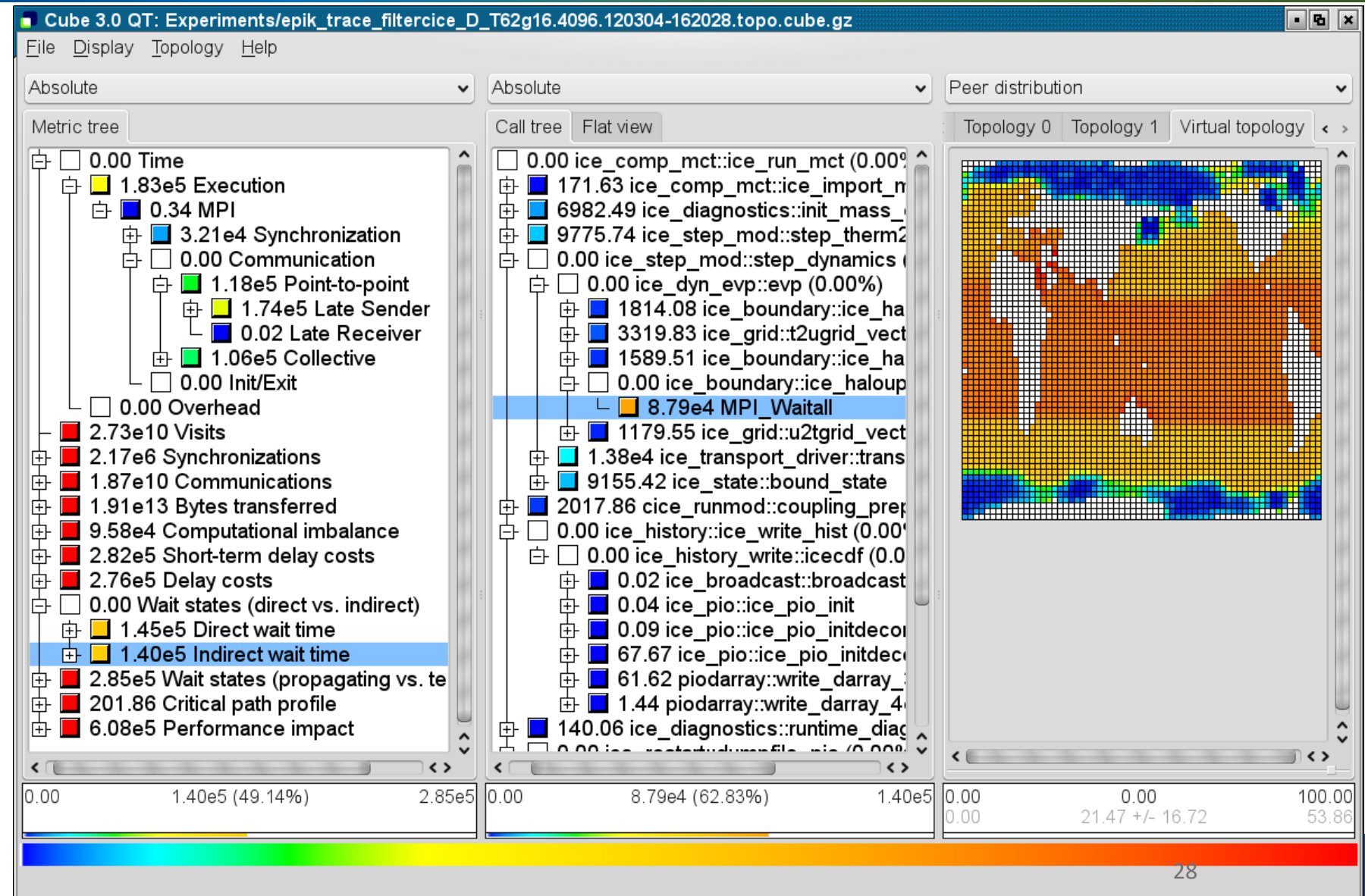


# Scalasca Example: CESM Sea Ice Module



## Indirect Wait Time Analysis

- Indirect waits occurs for ranks processing warmer areas

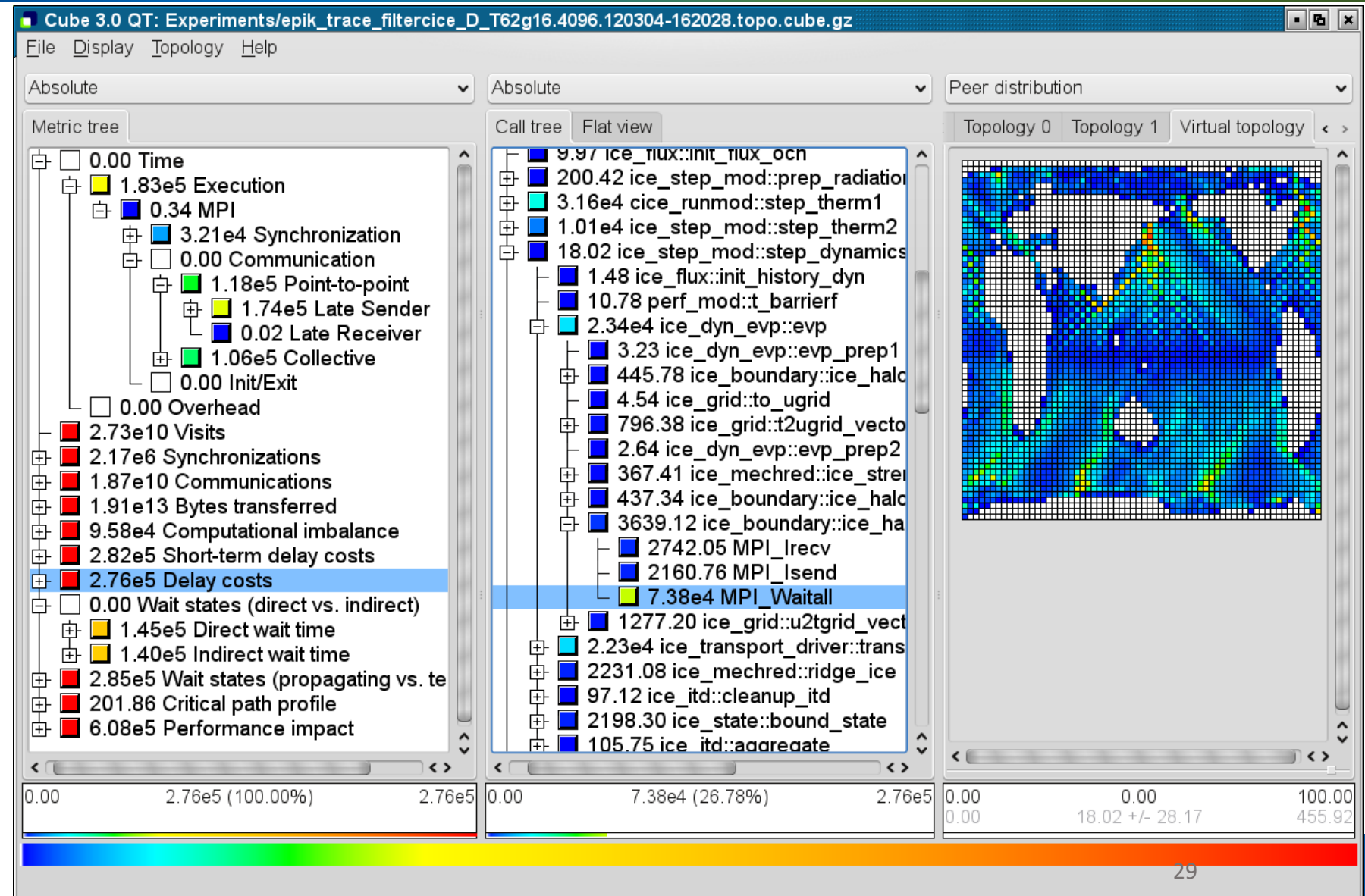


# Scalasca Example: CESM Sea Ice Module



## Delay Costs Analysis

- Delays **NOT** caused on ranks processing ice!





# Scalability

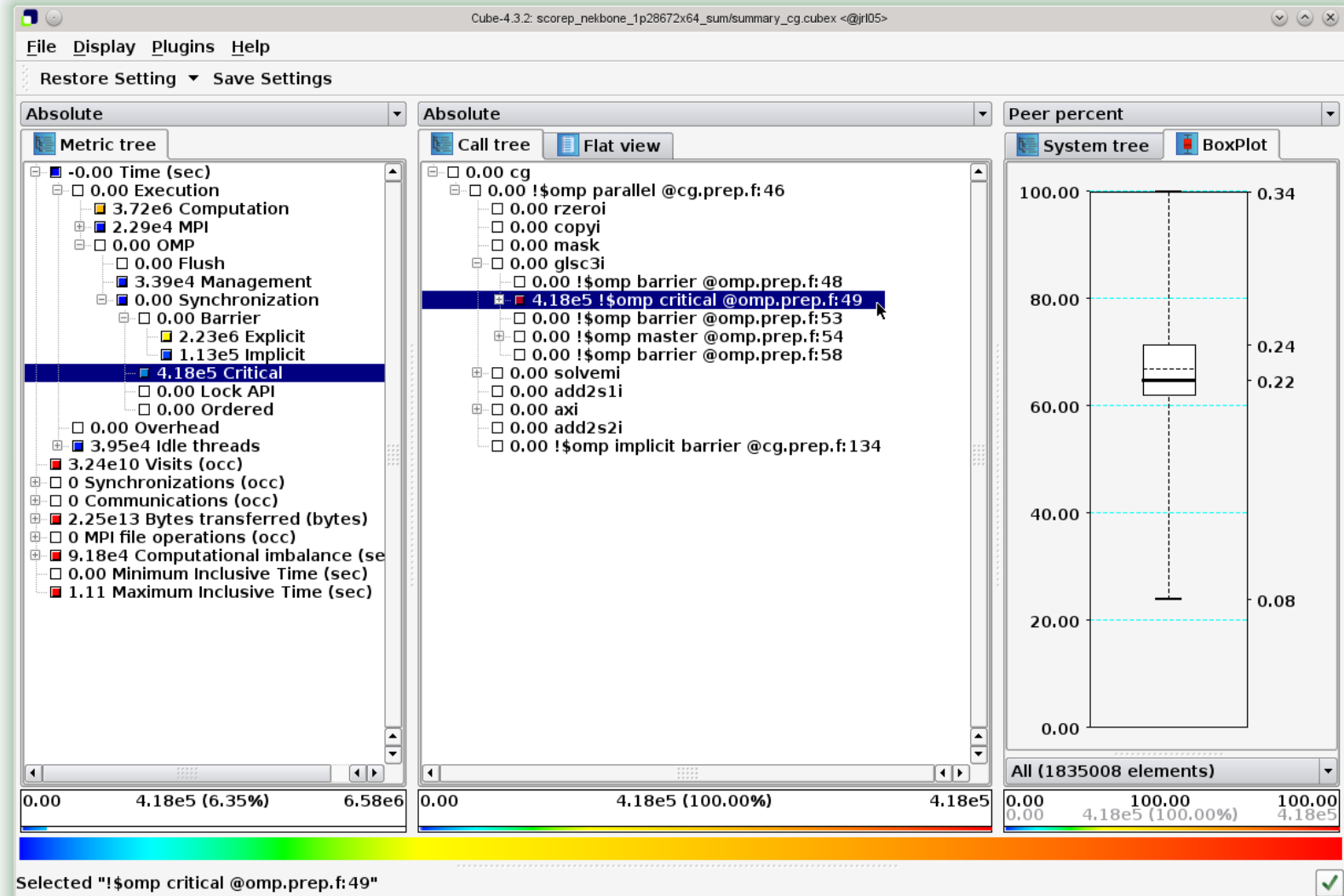
Towards Exascale



# Scalasca: 1,835,008 Threads Test Case



- Nekbone  
(CORAL benchmark)
- Measurement of **full system BGQ JuQueen computer run**
  - $28,672 \times 64 = 1,835,008$  threads
  - A few TB trace data!
  - A few Million events!
- Possible because:
  - **Scalable parallel measurement and trace analysis**

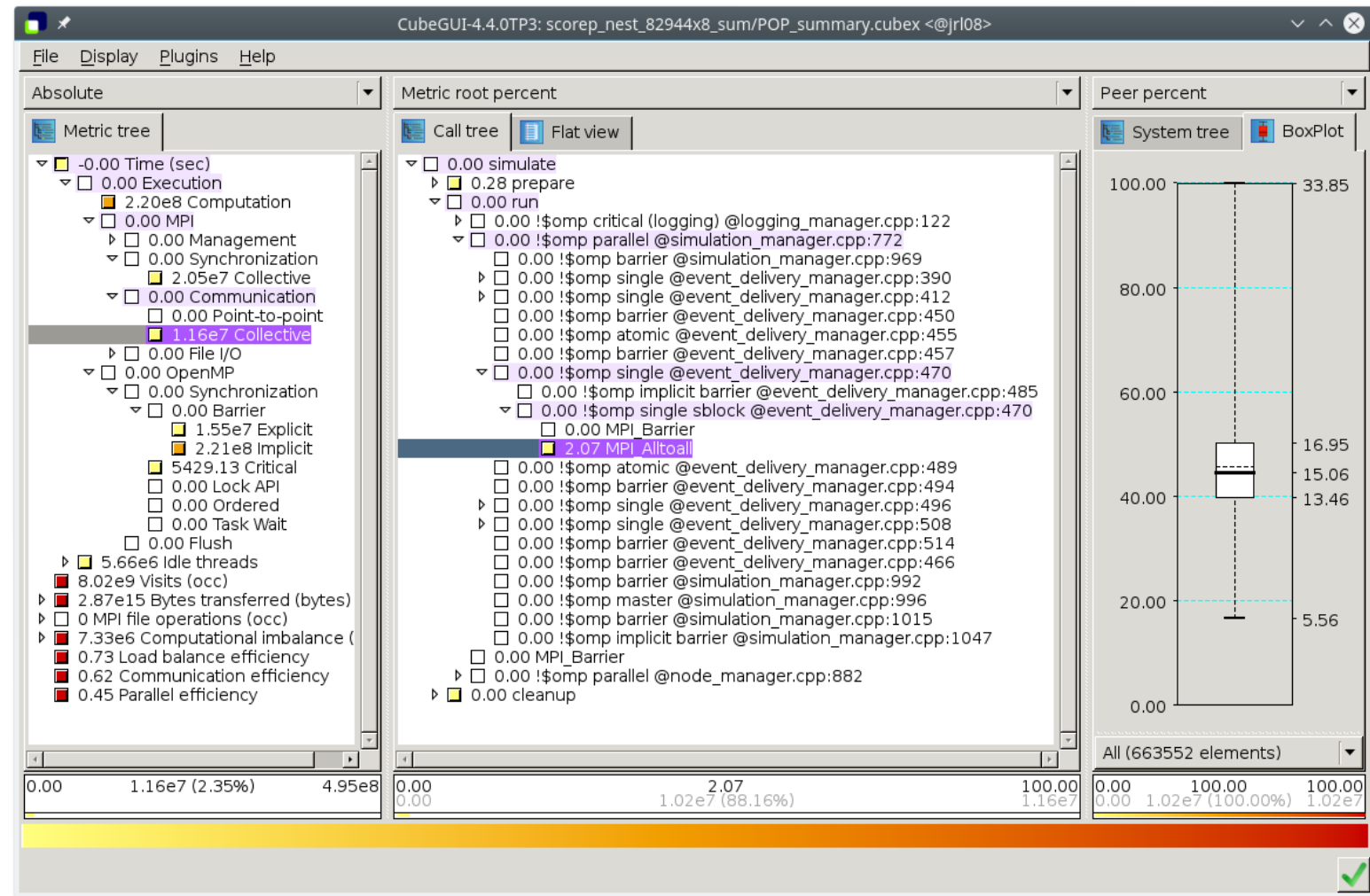




# Scalasca: User Analysis of NEST on K Computer



- Jülich **nest::** neural network simulator code
- Measurement of **full system K computer run**
  - 82,944 nodes
  - 663,552 threads
- Performance analyst
  - Itaru Kitayama (RIKEN)
- Analysis of MPI and OpenMP communication and synchronization at large scale







# Wrapup

What I hope you remember



- **Scalasca** facilitates
  - **Simple POP assessments** based on POP metrics!
  - **Advanced performance analysis!**
- Even on **extreme scales!**

<https://pop-coe.eu/further-information/online-training>

- POP Online Training modules



- [Installing POP Tools: Score-P, Scalasca, Cube](#)
- [Using POP Tools: Score-P and Scalasca](#)
- [Using POP Tools: Cube](#)

<https://pop-coe.eu/blog/19th-pop-webinar-identifying-performance-bottlenecks-in-hybrid-mpi-openmp-software>



- POP Hybrid Metrics Webinar

<https://www.vi-hps.org/training/tws/tw40.html>

- 40<sup>th</sup> VI-HPS Tuning Workshop
  - 14<sup>th</sup> to 18<sup>th</sup> June 2021
  - Online (ZOOM)

<https://app.swapcard.com/widget/event/isc-high-performance-2021-digital/planning/UGxhbm5pbmdfNDUzNTlw>

- ISC 2021 Tutorial “Determining Parallel Application Execution Efficiency and Scaling using the POP Methodology”
  - 24<sup>th</sup> June 2021, 14:00 – 18:00 CET
  - Judit Gimenz (BSC), Brian Wylie (JSC)
  - Online (SwapCard)



# Performance Optimisation and Productivity

A Centre of Excellence in HPC

## Contact:

 <https://www.pop-coe.eu>

 [pop@bsc.es](mailto:pop@bsc.es)

 [@POP\\_HPC](https://twitter.com/POP_HPC)

 [youtube.com/POPHPC](https://youtube.com/POPHPC)

