











# MEMORIES OF A (PARALLEL) PERFORMANCE ENGINEER: 38 YEARS OF TOOLS BUILDING

21 OCT 2025 | BERND MOHR





## THE EARLY YEARS

1987 TO 1992

FRIEDRICH ALEXANDER UNIVERSITÄT ERLANGEN-NÜRNBERG

JÜLICH SUPERCOMPUTING CENTRE

## **TELEFUNKEN TR440 @ FAU**

My first computer: used in programming class in 1<sup>st</sup> semester computer science

Page 3

- Fastest computer developed in Europe at this time (end of the 1970ies)
- **Dual** processor! ⇒ BUT, FAU had the only **3-processor** version!
- Amazing huge 1.5 Mbyte main memory
- Programming via punch cards (Algol60, Fortran IV, etc)

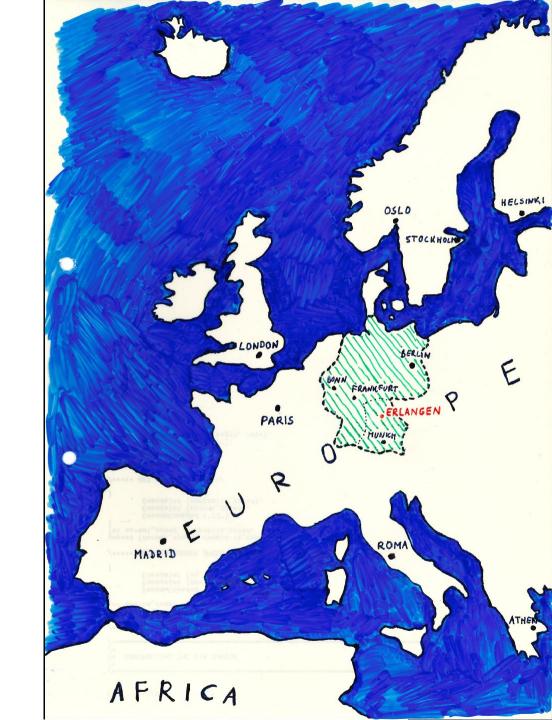


- Info
  - English: <a href="https://www.chessprogramming.org/TR\_440">https://www.chessprogramming.org/TR\_440</a>
  - German: <a href="https://de.wikipedia.org/wiki/TR\_440">https://de.wikipedia.org/wiki/TR\_440</a>



## CONTEXT

- Der SFB 182
  - "Multiprozessor- und Netzwerkkonfigurationen"
  - Four 3-year phases (1987 1998)
- Work Package C1
  - "Messung, Modellierung und Bewertung von Multiprozessoren und Rechnernetzen"
- Parallel system development @ FAU IMMD
  - EGPA (Erlangen General Purpose Array)
  - DIRMU (Distributed Reconfigurable Multiprocesser kit)
  - MEMSY (Modular Expandable Multiprocessor System)





#### Motivation (4)



- problems with monitoring non-sequential systems:
  - distributed system ↔ central monitor
  - global time (virtual, real)
  - non-reproducible behavior
  - "Probe Effect"
- other problem: many different objects
  - structure / configurations
  - operating systems
  - programming models
  - applications
- □ but: almost the same tasks
- requirement:

  object independent

Trace Monitoring and Analysis System

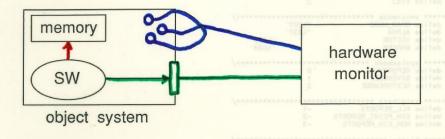




#### Monitoring (2)



- monitoring techniques:
  - software
  - hardware
  - hybrid



- hybrid monitoring:
  - <u>event definition</u> by inserting monitor instructions into the software (instrumentation)
  - event recording and timestamping with hardware

### **Event Recorder Board**

Page 6



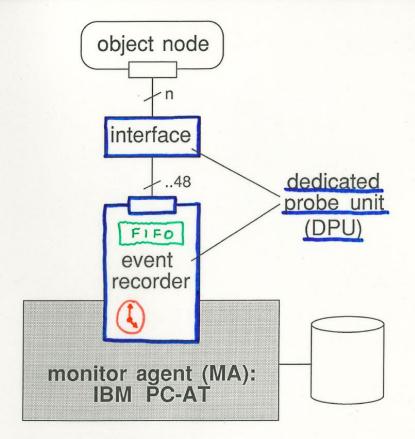


Bernd Mohr U-5



#### Monitoring (4)





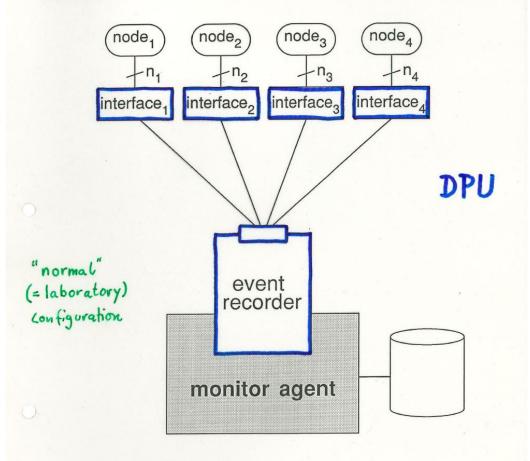
#### event recorder:

- IBM PC-AT compatible host interface
- high-resolution clock [100 ns]
- 10,000,000 events/s peak rate
- 40 timestamp 8 flags - 32K x 96 bit event buffer (FIFO) 48 event data
- 10,000 events/s mean rate (per MA)



#### Monitoring (5)





#### event recorder:

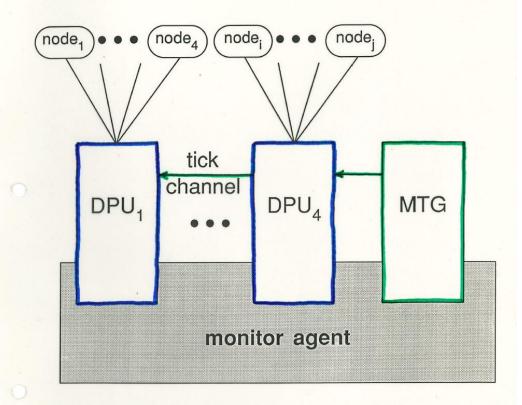
 4 independent event streams  $n_1 + n_2 + n_3 + n_4 \le 48$  bits

Bernd Mohr

Page 7

#### Monitoring (7)

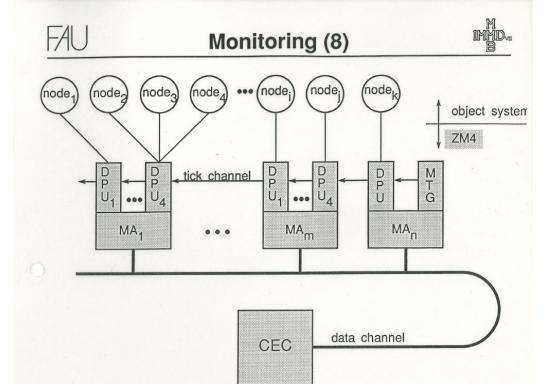




#### Measure Tick Generator:

- global time synchronization
- global start / termination
- fault-tolerant clock \_\_\_ sync events sync ok detection

> verification of clock data possible



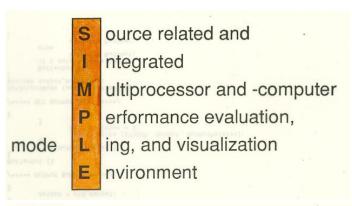
- Central Evaluation Computer
  - central control

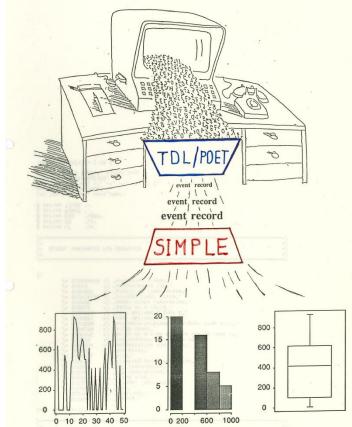
Bernd Mohr

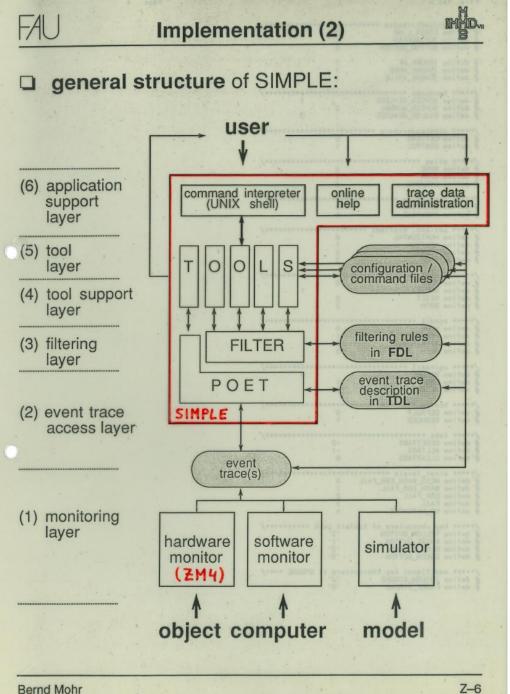
- gathers traces via data channel
- central (off-line) evaluation

#### distributed monitor system ZM4:

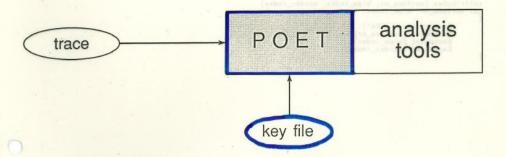
- adaptable to arbitrary computer systems
- global time with high-resolution [100 ns]







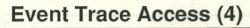




- Problem Oriented Event Trace interface
- advantages:
  - independent of event trace formats
  - standardized trace access interface
  - reusable function library in C
  - problem-oriented access

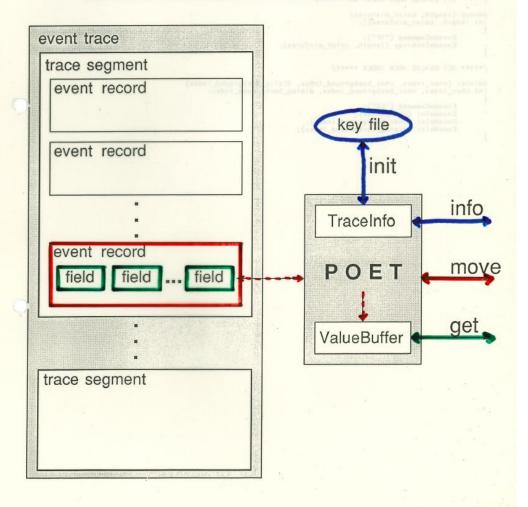


Bernd Mohr





- event trace handling with POET:
  - generic abstract data type

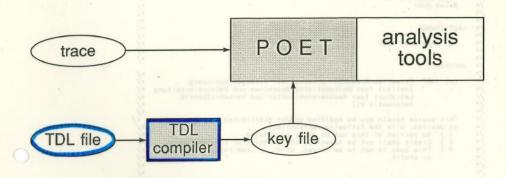


Z-9



#### **Event Trace Access (5)**





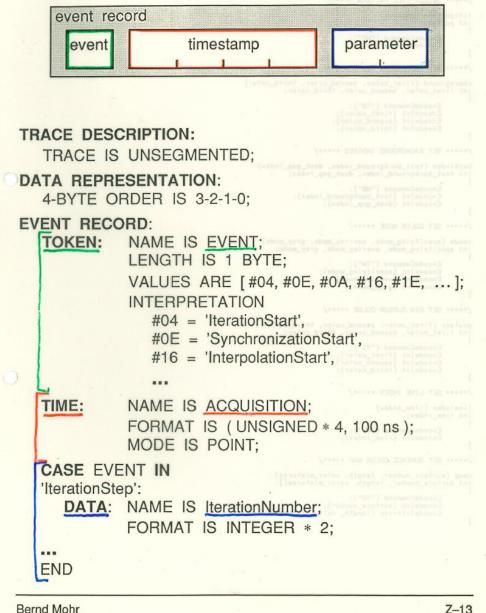
- event Trace Description Language
- differences to other systems:
  - problem-oriented description
  - description of representation possible
  - ASCII + binary trace formats
- advantages:

  - syntax and consistency checks



#### **Event Trace Access (6)**





Bernd Mohr

Z-11



#### **Event Trace Analysis (5)**



☐ example: FDL description

```
INIT < global variables>;
START FILTERMODULE:
   NAME IS < name>;
   ACTIONS
      IF < condition>
      DO
         PASS RECORD;
SET <variable> = <expr>;
         SWITCH TO < name>:
      END
FILTERMODULE:
   <name>
  ACTIONS
      PASS RECORD;
      IF < condition>
         INCREMENT < variable >;
      IF < condition>
         DECREMENT < variable>;
      IF < condition> EXIT:
```



Bernd Mohr

#### **Event Trace Analysis (5)**



**example**: FDL description

```
INIT COUNTER no of proc;
START FILTERMODULE:
  NAME IS search_begin;
  ACTIONS
     IF EVENT=='PrgStart' AND
        PROCESSOR=='Master'
     DO
        PASS RECORD;
        SET no of proc = 1;
        SWITCH TO program only;
     END
FILTERMODULE:
  NAME IS program only;
  ACTIONS
     PASS RECORD:
     IF EVENT == 'PrgStart'
        INCREMENT no of proc;
     IF EVENT == 'PrgEnd'
        DECREMENT no of proc;
     IF no of proc == 0 EXIT;
```

C-16



#### **Event Trace Analysis (6)**



#### trace validation:

- test whether measurement OK
- confirm expected behavior
- find unexpected behavior
- ☐ tool CHECKTRACE
  - for standard tests
- ☐ tool VARUS
  - VAlidation RUles checking System
  - allows user specified and problem-oriented assertions:

#### **ASSERT**

NUMBER (EVENT=='IterationStart') == NUMBER (EVENT=='IterationEnd')
ELSE "Iteration counting error";

#### **ASSERT**

IterationNumber INCREASING BY 1 ELSE "IterationNumber sequence error";



- example: trace statistics
- recommand language ⇔ tool features

  FREQUENCY <field name>

  DISTANCE <expr>
- identifiers, values ⇔ trace data

  FREQUENCY EVENT

  DISTANCE EVENT == 'receive' AND

  NODE == 'server'
- for application-dependent features:
  - predefined <u>standard names</u>
    - EVENT
    - ACQUISITION
    - NODE
    - PROCESS
  - predefined standard event types
    - send / receive
    - ...

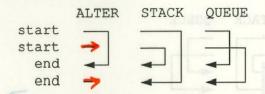
C-17

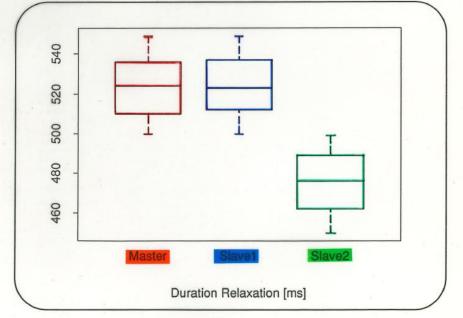
## **Trcstat Description Files**

```
□ general:
```

```
    FREQUENCY <fieldname>
        [ NAME = <identifier> ]
    DISTANCE <situation>
        [ NAME = <identifier> ]
    DURATION <situation> <situation>
        [ NAME = <identifier> ]
        MODE = ALTER | STACK | QUEUE
    <situation> := <event> [ @ <fieldname> ]
```

#### duration modes:





```
#01H DURATION 'RelaxB' 'RelaxE' PROCESSOR
#01S 461 [ms] on Slave2
#01S 523 [ms] on Master
#01S 527 [ms] on Slave1
#01S 470 [ms] on Slave2
#01S 502 [ms] on Master
#01S 520 [ms] on Slave1
            Master Slave1
                          Slave2
#01T
               231
                                    693
#01T
                      231
                              231
       no:
                                    450 [ms]
#01T
       min:
               500
                      500
                              450
                                    549 [ms]
#01T
      max:
               549
                      549
                                    508 [ms]
#01T mean:
               524
                      524
#01T
               207
                                    748 [ms]
                      215
       var:
```





#### **Event Trace Analysis (8)**



- □ activity:
  - interval in the dynamic behavior
  - defined by sequences of events
- ☐ inspired by **EDL** [Bates]
- □ tool FACT
  - Find ACTtivities

ments the To employed (2)

 allows user-specified and problem-oriented activity definitions as regular event expressions

#### **ACTIVITY** Iteration IS

IterationStart → IterationStep \*

→ Sync ?

→ IterationEnd

**END** 





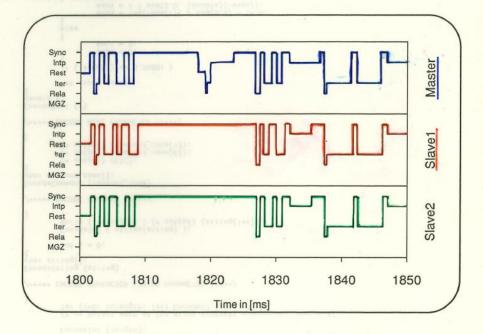
#### **Event Trace Analysis (3)**



#### ☐ Gantt diagram:

- program activities versus time
- functional dynamic behavior
- dependencies of activities
- duration of activities

#### ☐ tool GANTT



```
• GANTT DIAGRAM [<name>] WITH
                                           phase1
     <phase> : <evlist> ;
                                           phase2
                                           phase3
     ...
 END
• ACTIVITY GANTT DIAGRAM [<name>] WITH
                                           phase1
     <phase> : <evlist> TO <evtlist> ;
                                          phase2
                                           phase3
 END
• NUMBER GANTT DIAGRAM [<name>] WITH
     RANGE IS <num> TO <num>
     INCREASING BY <evlist> ;
     DECREASING BY <evlist> ;
 END
• FOR <tokenname> <interpretlist> DO
     <ganttdescriptions>
 END
```





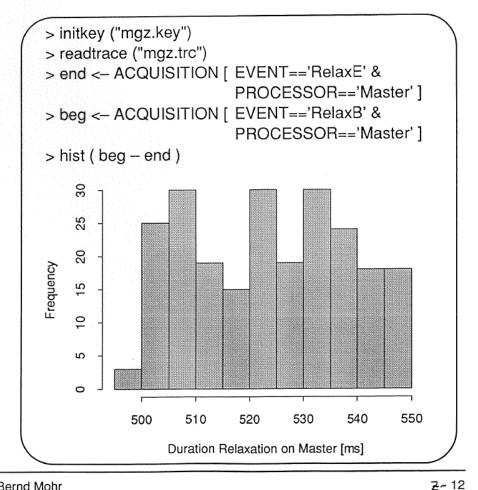
Bernd Mohr Z–17



#### example: trace analysis



- data analysis and graphics package **S** (AT&T)
  - high-level programming environment
  - interactive query language
  - S-POET interface



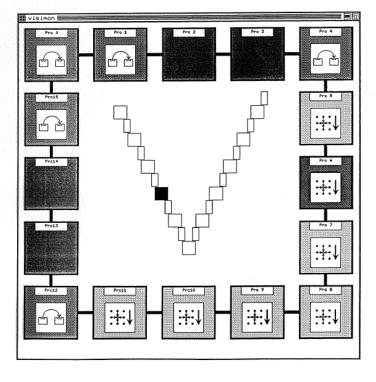


Bernd Mohr

#### example: animation



- animation tool VISIMON
  - based on X-Windows (X10!)
  - user specified animation description
  - program and data animation
  - "slow-motion" and "event-by-event" mode



#### **Applications**



object system / operating system	monitor / interface	application
DIRMU / DIRMOS	logic analyzer ZM4 / parallel port	numerical application simulation program
Transputer	ZM4 / link adapter ZM4 / bus adapter	communication system TRACOS
SUPRENUM / PEACE	ZM4 / 7 segment display	ray tracing
IBM-PC / OS/2, MSDOS	ZM4 / Centronics	protocol software B-ISDN, FDDI
IBM-PC / XENIX	ZM4 / Centronics	protocol software
SUN4 / SunOS	ZM4 / VME bus	X-Windows
SIEMENS robot control	ZM4 / SMP bus	robot control software
CCC3280 / XELOS	software monitor	multiprocessor UNIX
IBM-PC network	ZM4 / Centronics	Electrical Load Supervision Control System

IBM Zurich Research Lab

IBM ENC Heidelberg

Siemens Erlangen

Siemens Munich

> Fudan Univ. Shanghai



DIRMU (FAU, 1985-1988)



## SUPRENUM

- German: SUPerREchner für NUMerische Anwendungen
- English: super-computer for numerical applications



- German research project to develop a parallel computer from 1985 through 1990.
- Major effort aimed at developing a national expertise in HPC both at HW and SW level

Page 20

- Developed 256-node Suprenum-1
  - Fastest massively parallel MIMD computer in the world during a period in 1992
  - Still, project was canceled and considered a commercial failure
- https://en.wikipedia.org/wiki/SUPRENUM



#### What worked

- Hybrid instrumentation
- Sophisticated and highly advanced hardware monitor (ZM4)
- Fully flexible trace analysis framework (SIMPLE)
  - Generated by different sources (SW tracing, HW monitoring, LANalyzer, log files, simulation outputs, ....)
  - Of diverse applications (MP Unix, network stack, simulations, robot control software)

#### What didn't

- ZM4 too expensive for really large configurations
- SIMPLE unusable by non-expert (highly complex programming required)
- IFF issue: what's wrong? Investigated target? Description?





## THE POSTDOC YEARS

1993 TO 1995

UNIVERSITY OF OREGON, EUGENE



## CONTEXT

- ARPA funded project "pC++"
  - Programming Environments, Compiler Technology and Runtime Systems for Object-Oriented Parallel Processing
    - Dennis Gannon, Indiana University
    - Postdocs: Pete Beckman, Francois Bodin
    - pC++ compiler and runtime system, Sage++ toolkit
  - Languages, Libraries and Performance Evaluation Tools for Scalable Parallel Systems

Page 23

- A. Malony, J. Cuny, University of Oregon
- Postdoc: Bernd Mohr
- TAU program analysis tools



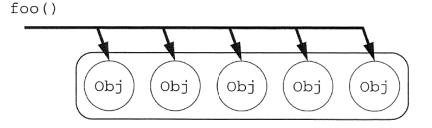
#### pC++ - The Programming Language Ideas

regular C++: programmers apply operators and functions to objects as "messages"

```
class Obj {
   int x;
   void foo();
};
Obj myObj;
myObj.foo();
```

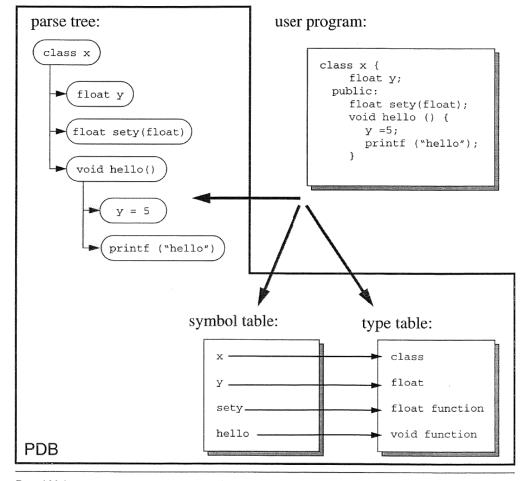
pC++: this concept is extended so that an operator or function can be applied to a large set, grid, array (:= collection) of objects (:= elements) in parallel

```
Collection Vector { ... };
Vector<Obj> paraObj(AlignObj, DistrObj);
paraObj.foo();
```

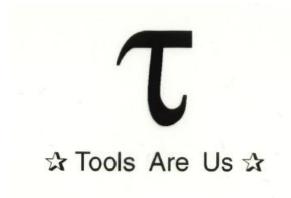


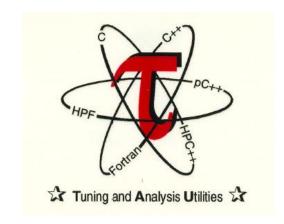
#### Sage++

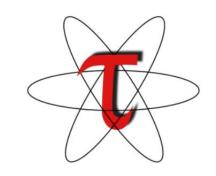
- ☐ C++ class library for building program analysis and transformation tools
- contains parsers for Fortran 77 / 90, ANSI C, and C++
- organized as a class hierarchy for accessing and modifying the parse tree, symbol table, and type table



## **TAU LOGO EVALUATION**







TAU Performance System ®

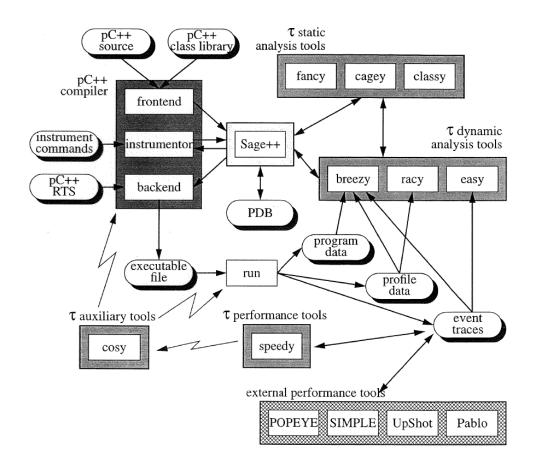


https://www.cs.uoregon.edu/research/paraducks/



## The pC++ Programming Environment







#### ☆ Tools Are Us ☆

- $\Box$  Currently available  $\tau$  tools:
  - Cosy (COmpile manager Status displa Y)
  - O fancy (File ANd Class display)
  - O cagey (CAII Graph Extended displa Y)
  - O classy (CLASS hierarch Y browser)
  - O racy (Routine and data ACcess profile displa Y)
  - speedy (Speedup and Parallel Execution Extrapolation Displa Y)
  - O **breezy** (**BR**eakpoint **E**xecutive **E**nvironment for visuali **Z**ation and data displa **Y**)
- □ Prototypes:
  - O easy (Event And State displa Y)
  - O dandy (Distributed Array Navigator Displa Y)
  - O crafty (ContRol flow And FuncTion displaY)
  - geeky (GEeky Editing and symbol loo Kup displa Y)
  - O POPEYE, DAQV (data and performance visualization)

University of Oregon

- $\Box$   $\tau$  can work with a local or remote pC++ language system
- □ T originally designed for C++/pC++ programs



- ☐ Providing a user (program-level) view
  - T graphical interface objects represent pC++ language level objects: collections, classes, methods, functions
- Support for high-level, parallel programming languages
  - T designed and implemented in concert with pC++ system
  - T translates low-level performance data  $\Leftrightarrow$  language level
- Integration with compilers and runtime systems
  - → T integrated with pC++ runtime system.
  - $\mathsf{T}$  uses Sage++ for access to **P**rogram **D**ata **B**ase (PDB)
- Portability, extensibility, and retargetability
  - au implemented using C/C++ and Tcl/Tk for portability
  - o implemented as <u>hypertools</u> for extensibility
  - au uses Sage++ for retargetability
- Usability

Bernd Mohr

- tool objects represent hyperlinks
- wo has on-line hypertext help



#### Hypertools

- tools are distinct tools, but they act in concert as if they were a single, monolithic application
- implemented using hyperlinks and global features

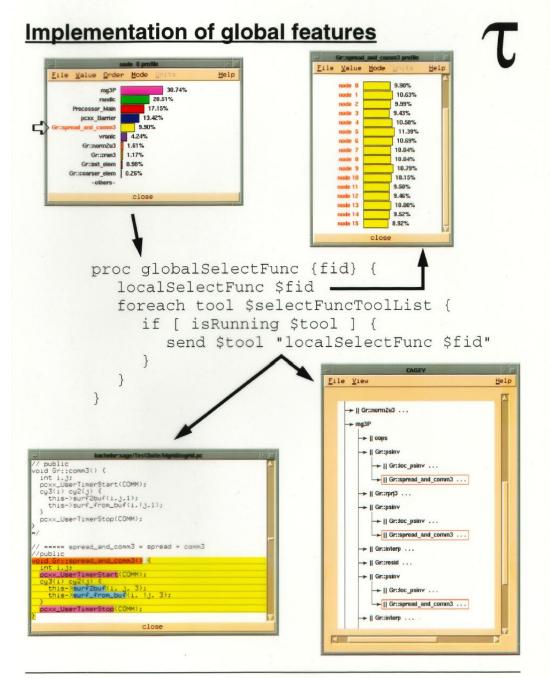
#### ☐ Hyperlinks

- T graphical interface objects act like in hypermedia documents
- selecting an object of interest brings up a more detailed or related view of the object
  - e.g., selecting a class in the class hierarchy graph displays a table of class members

#### □ global features

- synchronized hyperlinks: execution of a global feature in one tool automatically updates views in the other tools
- currently implemented:
  - O select-function O load-depfile
  - O select-class O select-callsite

Mar '95

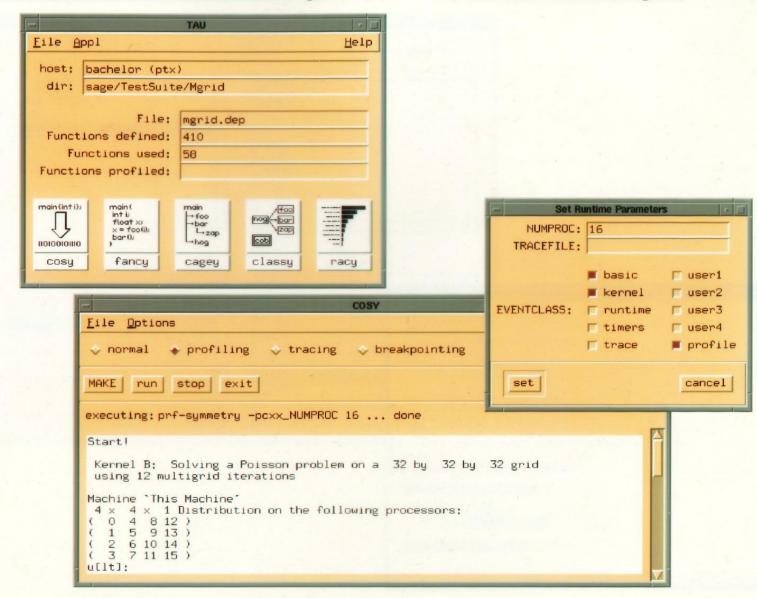




### TAU / COSY

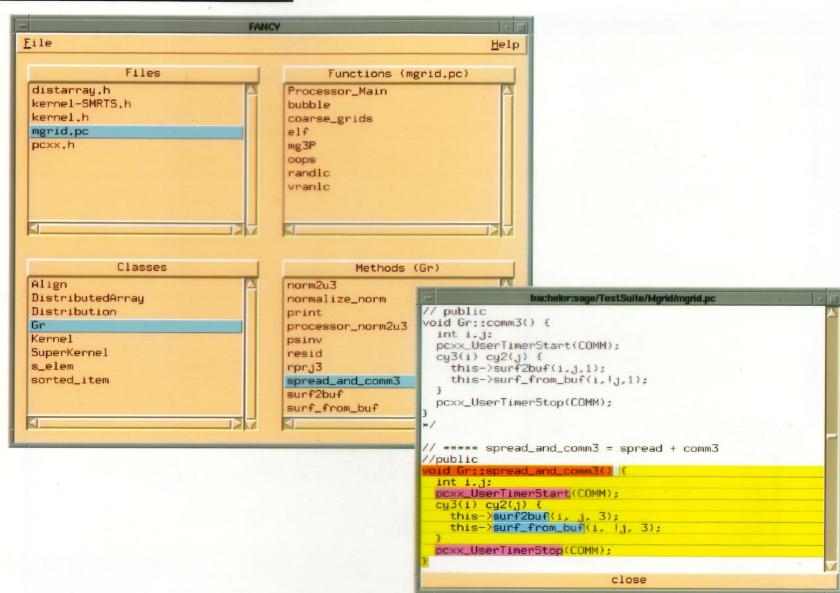
### Main Control Panel / Compile and Execute Manager





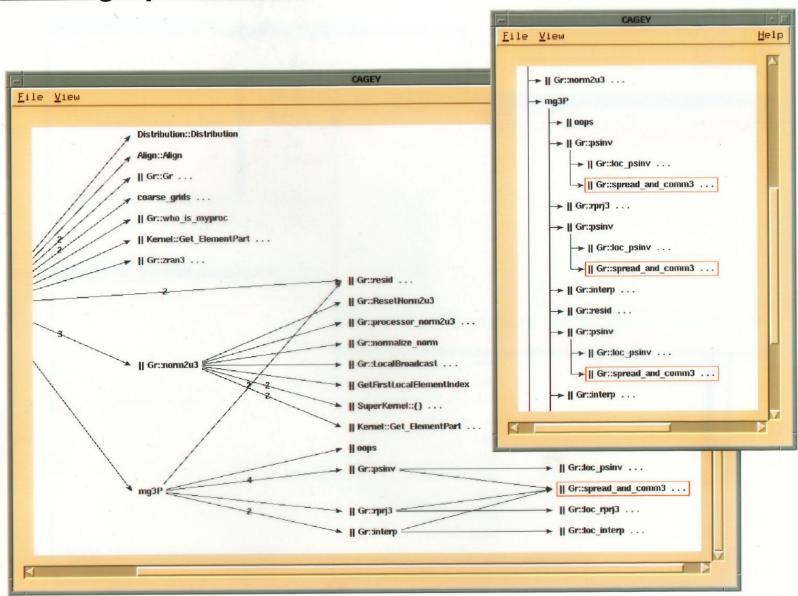








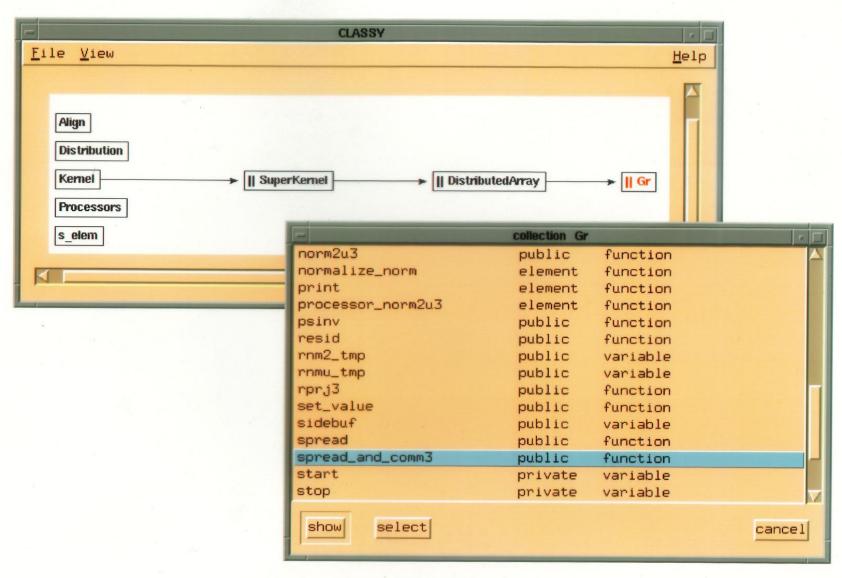






## **CLASSY Class Hierarchy Browser**







## PPROF / RACY Parallel Program Profiler



uses instr tool to dynamically instrument		
O	entry and exit of	
	functions, class member functions, constructors, and destructors	
	in the user application and the pC++ kernel	
uses "profile-instrumented" runtime system		
data gathered for each profiled function:		
O	time spent in function including and excluding its children	
0	number of calls	
data	data gathered for each pC++ collection:	
0	number of local and remote accesses per node	
profile data is stored in one file per node in a machine-independent format		
ppr	pprof: prints ASCII profile report (like UNIX's prof)	
rac	racy: graphical profile data browser	

#### ☐ Example pprof Output

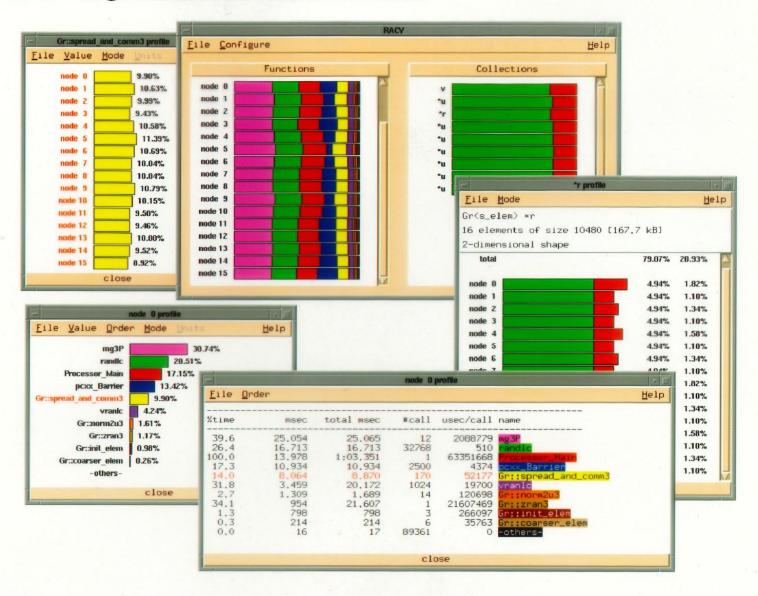
```
NODE 0:
%time msec total msec #call usec/call name
                            1 28869157 Processor Main
100.0 15.761
                 28.869
                                4172459 initw
 14.5 4.166
                 4.172
 12.2 3.512
                 3.512 560567
                                      6 Vector::get
  6.9
                 2.004
                               1002201 DistVector::DistVector
  6.6 1.919
                 1.919
                           44
                                  43628 Barrier
  5.9 1.694
                                  52958 Vector::allocData
                 1.694
  5.5
         132
                 1.578
                          1 1578414 poisson_solve
  2.9
        1
                   835
                           32
                                  26117 ParSinTransform
  2.9
                   833
         638
                                  26055 Vector::SinTransform
   local remote collection
accesses accesses num
   16619
             496
                   0
                      F
    9789
               5
                   1 U
// similar output from other 31 processors deleted
FUNCTION SUMMARY (mean):
       msec total msec #call usec/call name
100.0 15.312
                            1 28301704 Processor_Main
                 28.301
 14.4 4.058
                 4.064
                                4064819 initw
                            1
 12.4 3.519
                 3.519 563042
                                      6 Vector::get
 7.9 2.246
                 2.246
                                  51065 Barrier
  7.1
          9
                 2.004
                           2 1002249 DistVector::DistVector
  5.6
        116
                 1.578
                           1 1578486 poisson_solve
  4.8
      1.348
                 1.348
                              42138 Vector::allocData
  3.0
                   845
                           3.2
                                 26418 ParSinTransform
 3.0
        640
                   843
                                  26356 Vector::SinTransform
COLLECTION SUMMARY:
DistVector<Vector> F, collection #0
      513 elements of size 4216, 1-dimensional
      532996 local / 15841 remote accesses
DistVector<Vector> U, collection #1
      513 elements of size 4216, 1-dimensional
      282720 local / 423 remote accesses
```

Mar '95

### **RACY**

### Parallel Program Profile Visualizer







#### **Event Tracing Support Tools**

#### pcxx\_merge:

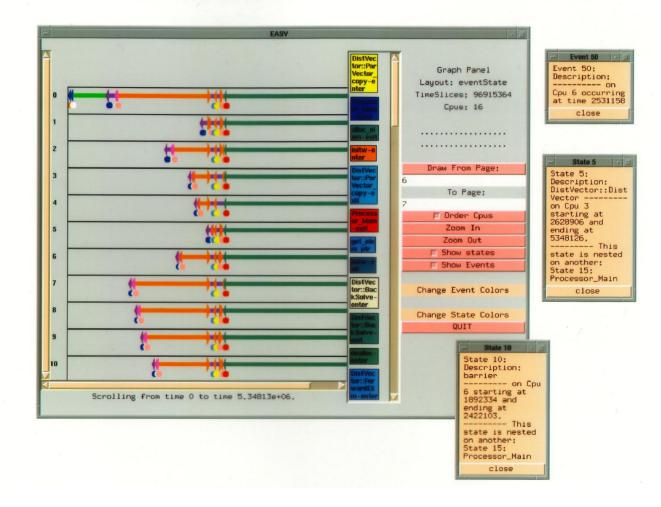
- merges node traces to global system trace according to timestamps
- establishes global timescale, if target system doesn't have global clock
- pcxx\_convert: trace format converter
  - O generic ASCII
  - O alog
  - O SDDF

#### **Event Trace Tools**

#### easy

- $\circ$   $\circ$   $\tau$  event trace browser for alog format
- SIMPLE (University of Erlangen, Germany)
  - O trace format independent event trace analysis environment
- ☐ **Upshot** (Argonne National Laboratory)
  - O simple+portable X11 event trace browser for alog format
- ☐ Pablo (University of Illinois)
  - sophisticated event trace analysis environment based on SDDF format

## **EASY Event Trace and State Browser**





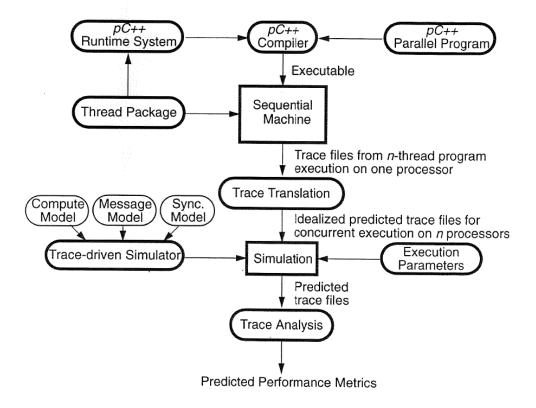
August'95

## ExtraP Performance Extrapolation and Analysis



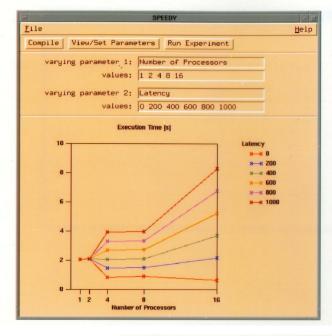
#### ☐ ExtraP

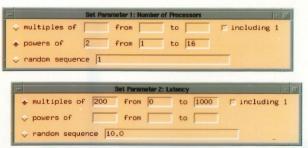
- O high-level event tracing of a *n*-thread pC++ program on a uniprocessor workstation
- O trace-driven simulation for prediction of performance on n-processor parallel machine



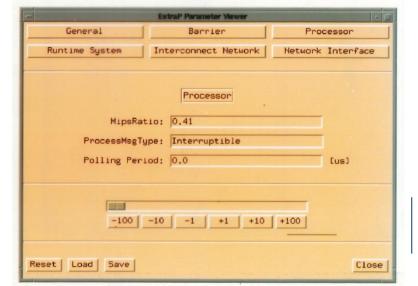
## SPEEDY Performance Extrapolation and Analysis







### **ExtraP**



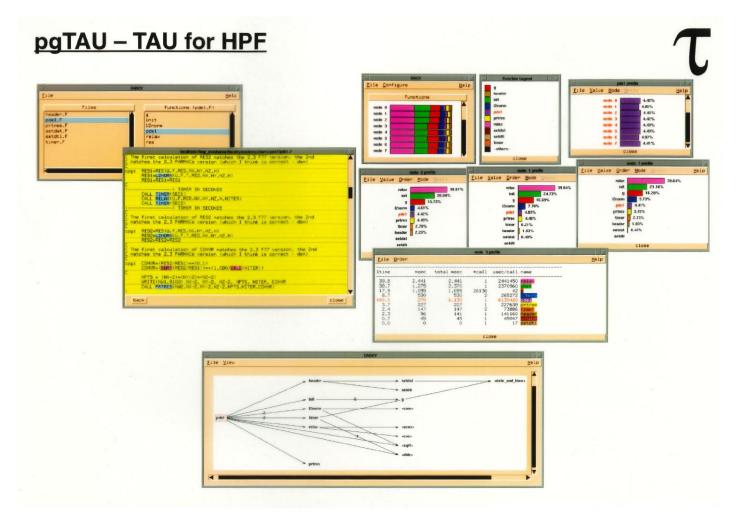
JÜLICH SUPERCOMPUTING CENTRE

#### pC++ - Supported Systems

- □ Shared memory systems
  - O Kendall Square KSR-1 / KSR-2
  - O Sequent Symmetry (under Dynix + PTX)
  - O SGI (Power) Challenge + Onyx
  - O Convex SPP-1
  - O (BBN TC2000)
- Distributed memory systems
  - O TMC CM-5
  - O Intel Paragon
  - O IBM SP-1 / SP-2
  - O Cray T3D / T3E
  - O Meiko-CS2
  - O Workstation Clusters with PVM + MPI (homogeneous)
- ☐ UNIX Workstations (SUN, HP, DEC, IBM, SGI, ...)
  - O serialized
  - O thread-based (Pthreads, LWP, AT&T tasks, Awesime)
- The <u>same</u> pC++ program will run <u>without modification</u> on all platforms



#### **BEYOND PC++**



- prototype port of
  - (function browser) fancy
  - (callgraph browser)
  - racy (profile data browser)

to HPF compiler system of The Portland Group, Inc.

- changes needed for static browsers
  - pgdep
    - tool for generating HPF PDB (Program Data Base)
    - generated from intermediate f77 sources
  - om
    - new object manager which provides the TAU standard static browser interface to HPF PDB
- changes needed for profiling
  - HPF compiler already supports instrumentation
  - rewrite of the profiler runtime system functions to output pprof / racy compatible profile data files



Page 38

#### **ASSESSMENT**

#### What worked

- Early fully featured parallel programming environment (pC++ and TAU)
  - Easy to use (build, run, analyze)
  - Global features (hyper tools)
  - Although build for / integrated into pC++, TAU was easy to retarget

#### Undecided

Most innovative or worst configuration system (before GNU configure or Cmake)

#### What didn't

C++ parsing is just to complicated to be implemented within an University project

Page 39

Didn't support traces very well





### THE LATER YEARS

1996 TO NOW

/

JÜLICH SUPERCOMPUTING CENTRE



#### **CONTEXT: 25 YEARS OF AUTOMATIC TRACE ANALYSIS**

1999 – 2004



- EU ESPRIT + IST Working Group
- http://www.fz-juelich.de/apart/\*



- Sequential analyzer EXPERT
- http://www.fz-juelich.de/zam/kojak/\*

• 2006 – now



- Helmholtz Virtual Institute
- http://www.vi-hps.org/

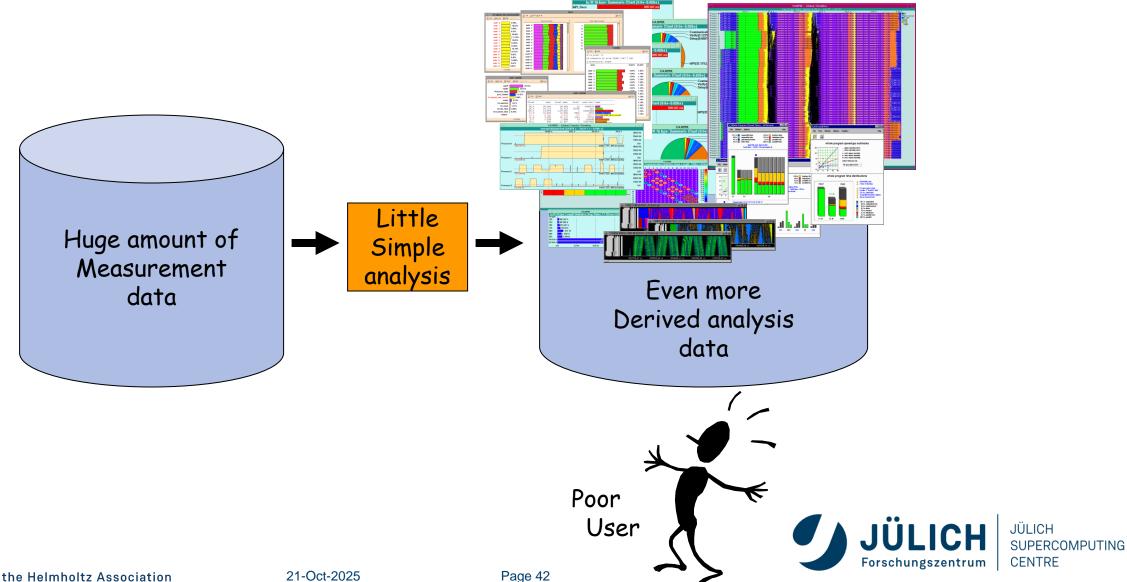


- Parallel analyzer SCOUT
- http://www.scalasca.org



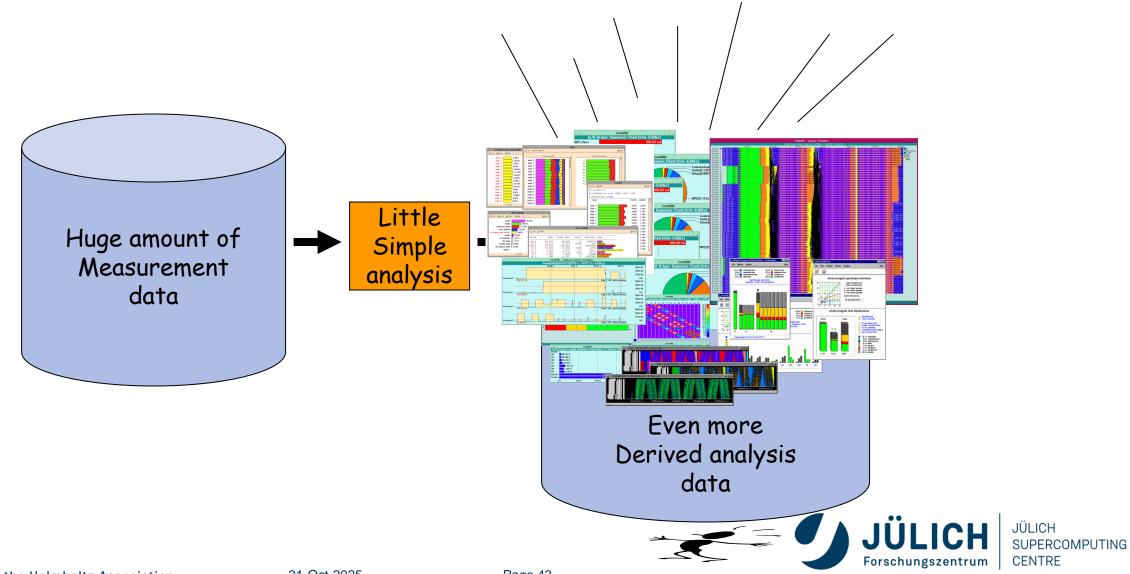
#### TRADITIONAL PERFORMANCE TOOLS





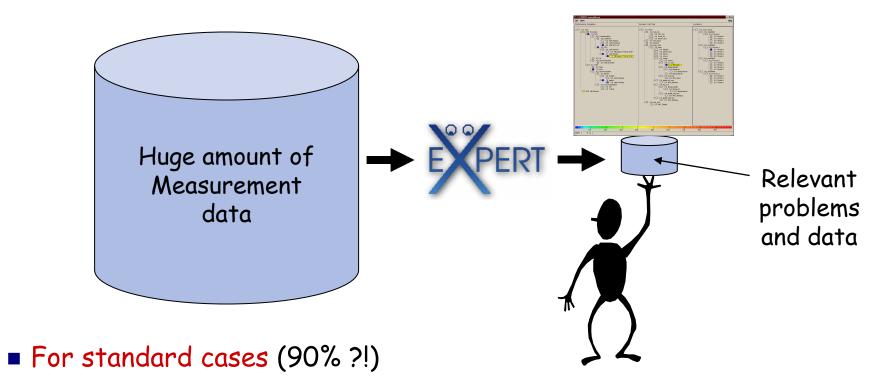
#### TRADITIONAL PERFORMANCE TOOLS





#### **SOLUTION PART 1: AUTOMATIC TOOL**



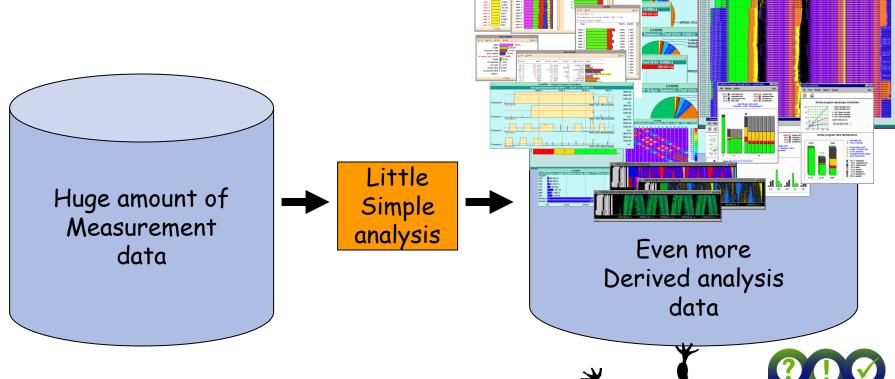


- For "normal" users
- Starting point for experts



#### **SOLUTION PART 2: EXPERT TOOLS + EXPERT**



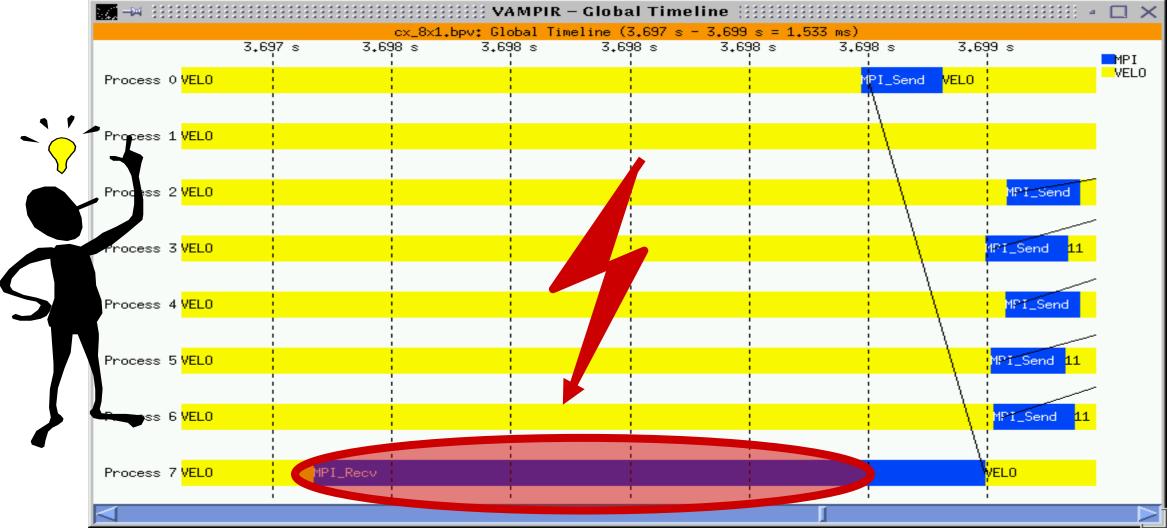


Page 45

- For non-standard / tricky cases (10%)
- ⇒ More productivity for performance analysis process!



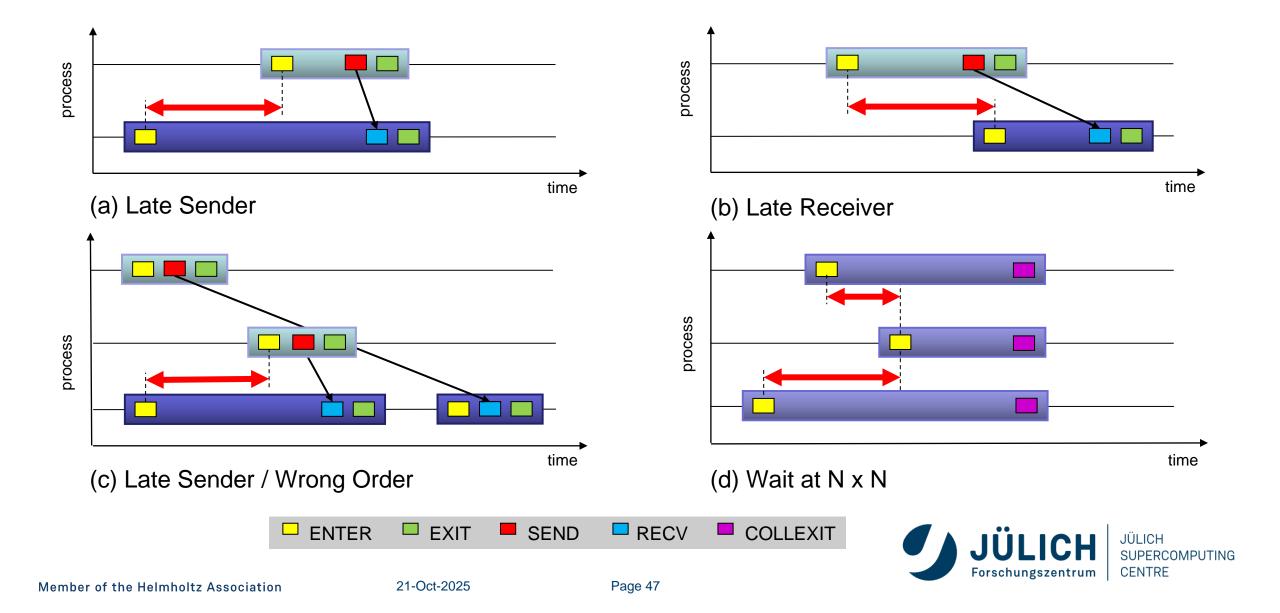
#### **EXAMPLE AUTOMATIC ANALYSIS: LATE SENDER**



Page 46

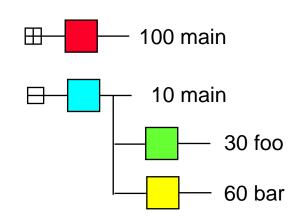
Member of the Helmholtz Association

#### **EXAMPLE MPI WAIT STATES**



[2001]

- Performance behavior
  - 3 dimensional matrix
  - Hierarchical dimensions
- Weighted tree
  - Tree browser
  - Each node has weight
    - \* Percentage of CPU allocation time
    - \* E.g. time spent in subtree of call tree
  - Displayed weight depends on state of node
    - \* Collapsed (including weight of descendants)
    - Expanded (without weight of descendants)
  - Displayed using
    - \* Color
      - Allows to easily identify hot spots (bottlenecks)
    - \* Numerical value
      - Detailed comparison



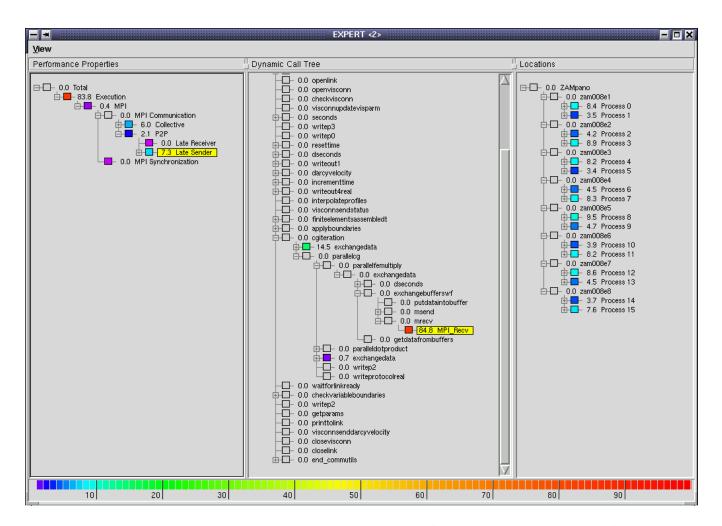


#### PRESENTATION OF PERFORMANCE BEHAVIOR (2)

Page 49



- Three views
  - Performance property
  - Call tree
  - Locations
- Interconnected
  - View refers to selection in left neighbor
- Two modes
  - Absolute: percent of total **CPU** allocation time
  - Relative: percent of selection in left neighbor
- Collapsing/expanding of nodes
  - Analysis on all hierarchy levels



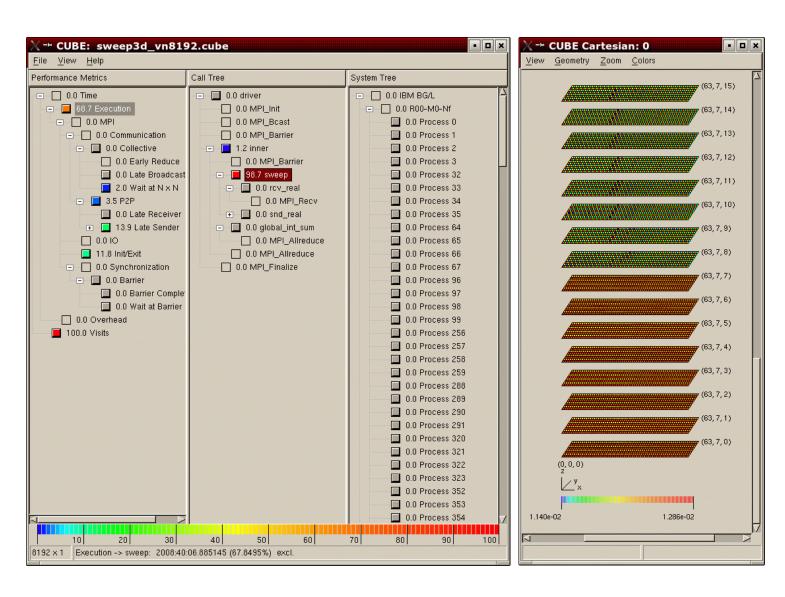


[2003]



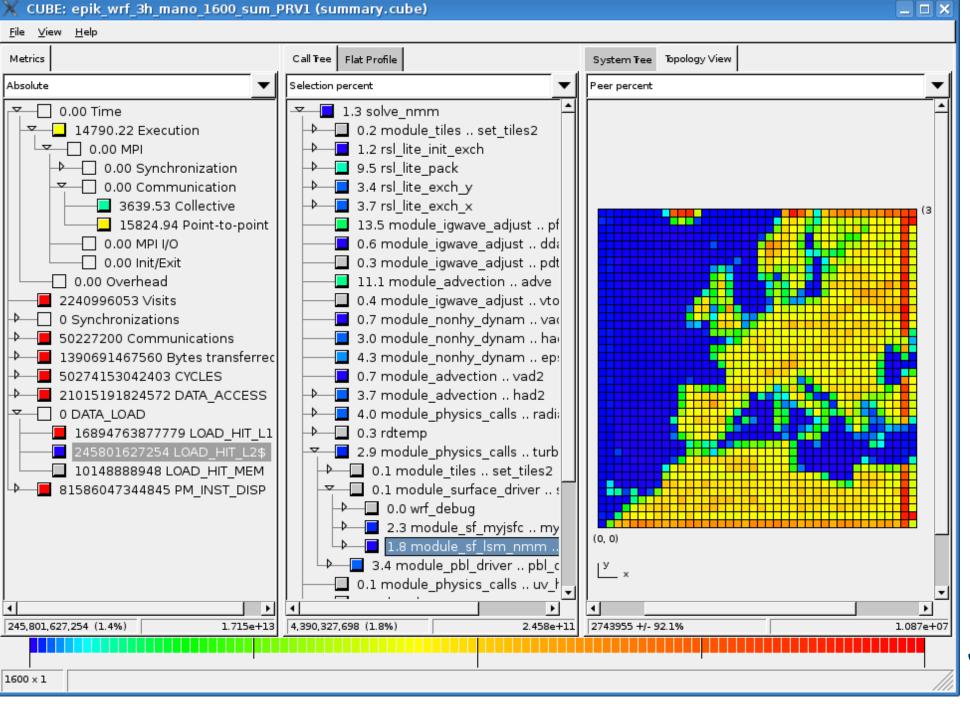
#### **EXAMPLE: SWEEP3D ON 8192 BG/L PES**





- New topology display
- Shows distribution of pattern over HW topology
- Scales to larger systems





[2009]

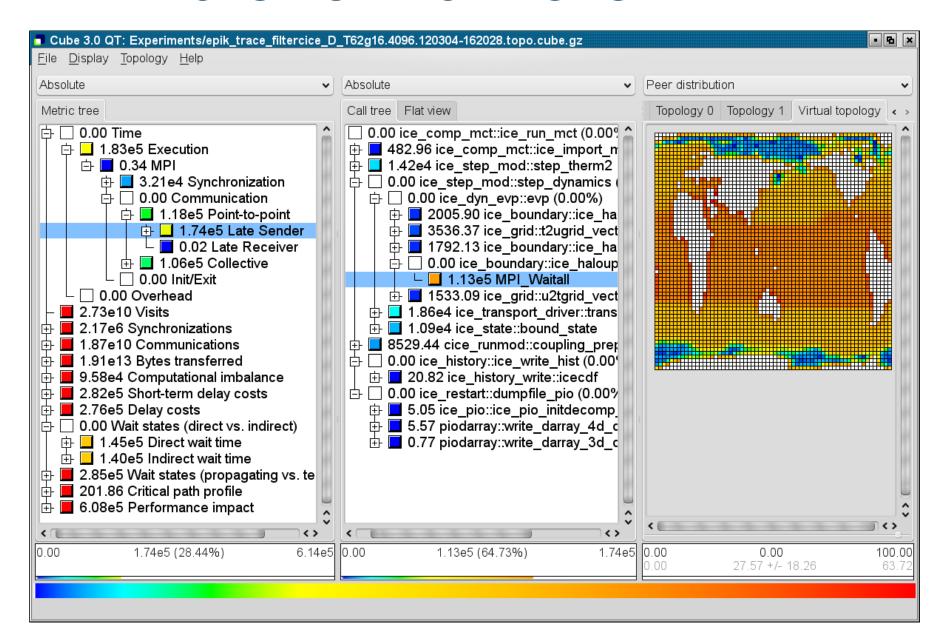
Application topologies



#### SCALASCA EXAMPLE: CESM SEA ICE MODULE

# Late Sender Analysis + Application Topology

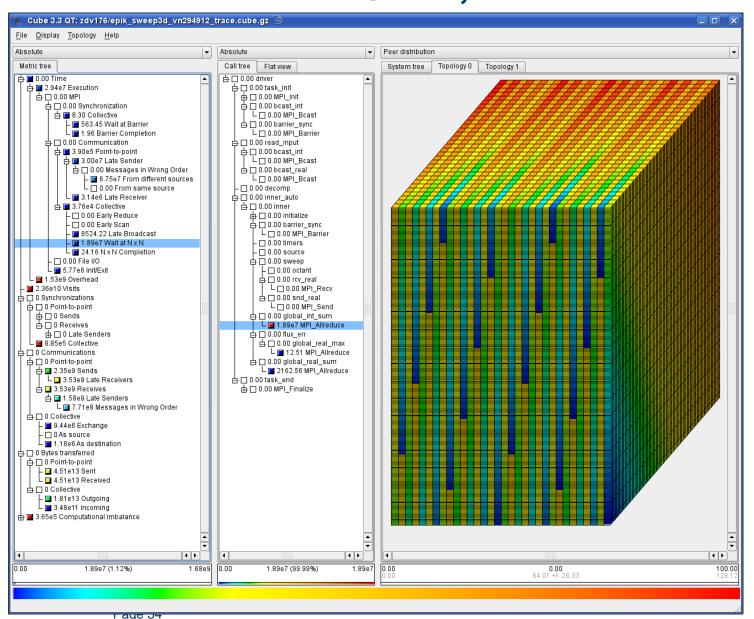
- Shows distribution of imbalance over topology
- MPI topologies are automatically captured



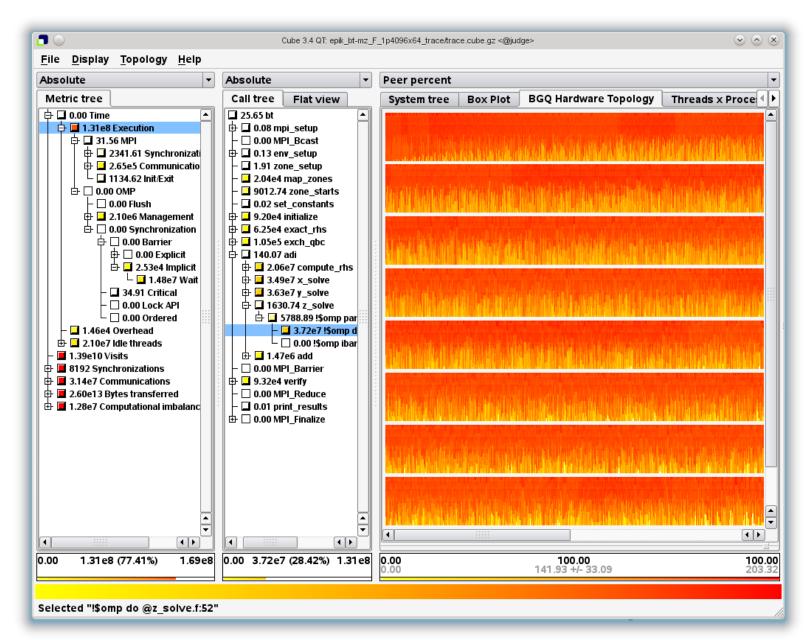
#### SCALASCA TRACE ANALYSIS SWEEP3D@294,912 BGP [2010]

- 10 min sweep3D runtime
- 11 sec analysis
- 4 min trace data write/read (576 files)
- 7.6 TB buffered trace data
- •510 billion events

B. J. N. Wylie, M. Geimer, B. Mohr, D. Böhme, Z.Szebenyi, F. Wolf: Large-scale performance analysis of Sweep3D with the Scalasca toolset. Parallel Processing Letters, 20(4):397-414, 2010.



#### SCALASCA TRACE ANALYSIS BT-MZ@1,048,704 BGQ [2013]

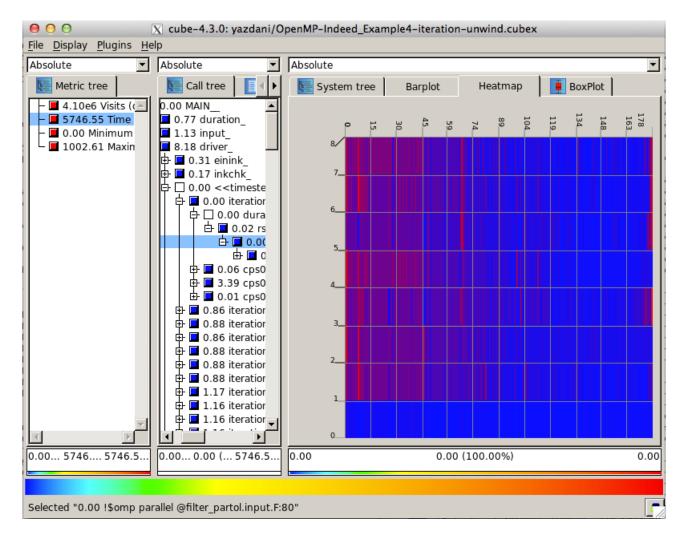




Phase profiling

 Collects data for each instance of phases marked in program instead of aggregating it

 Shows data over "time" (phase instances) for each rank/thread

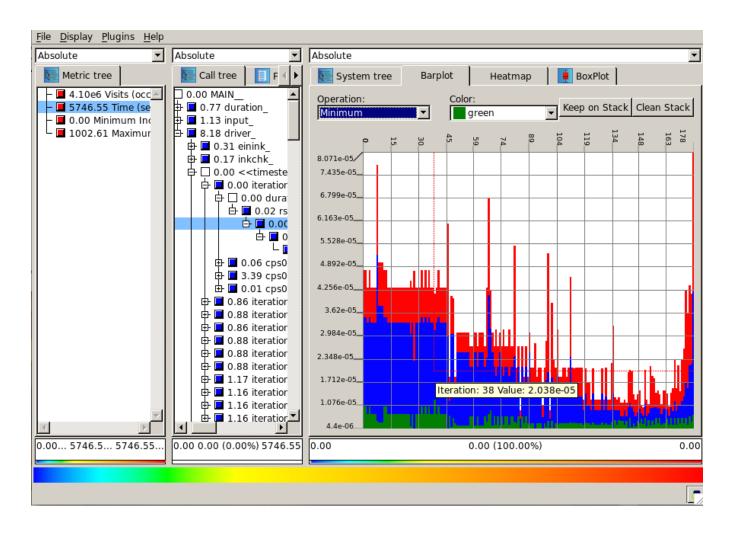




#### **CUBE VIZ PLUGINS: PHASE BARPLOT**

#### Phase profiling

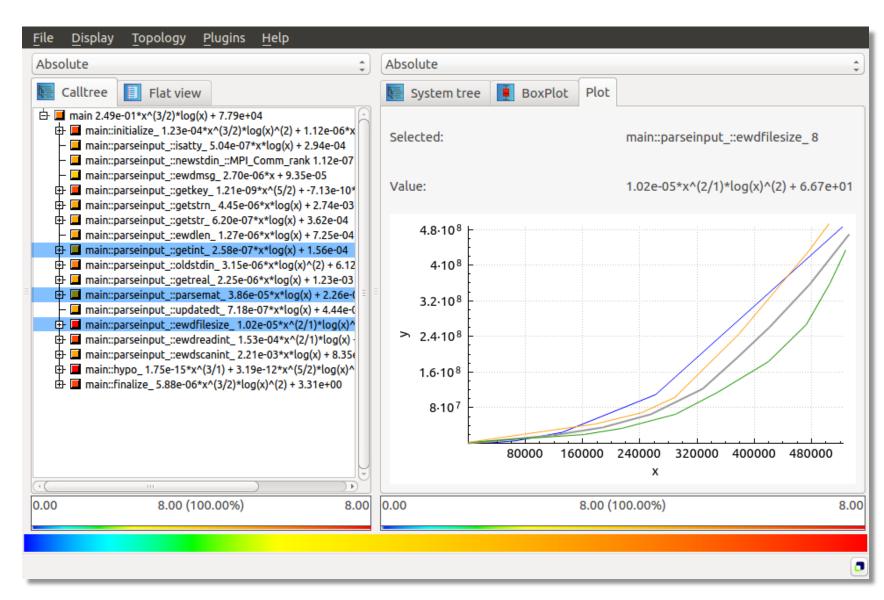
- Collects data for each instance of phases marked in program instead of aggregating it
- Shows data over "time" (phase instances) for each rank/thread





#### **CATWALK: MODELING RESULT VISUALIZATION**

[2015]

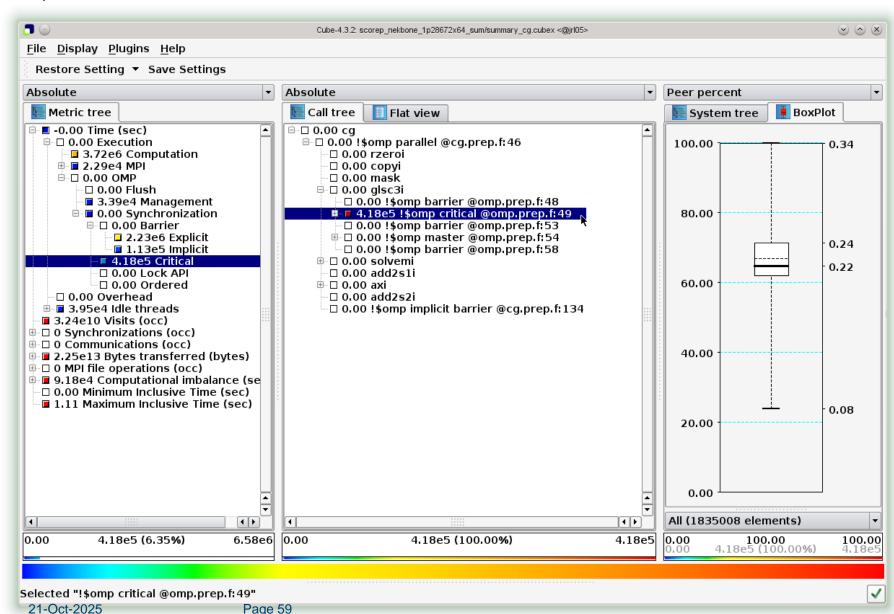




#### SCALASCA: 1,835,008 THREADS TEST CASE

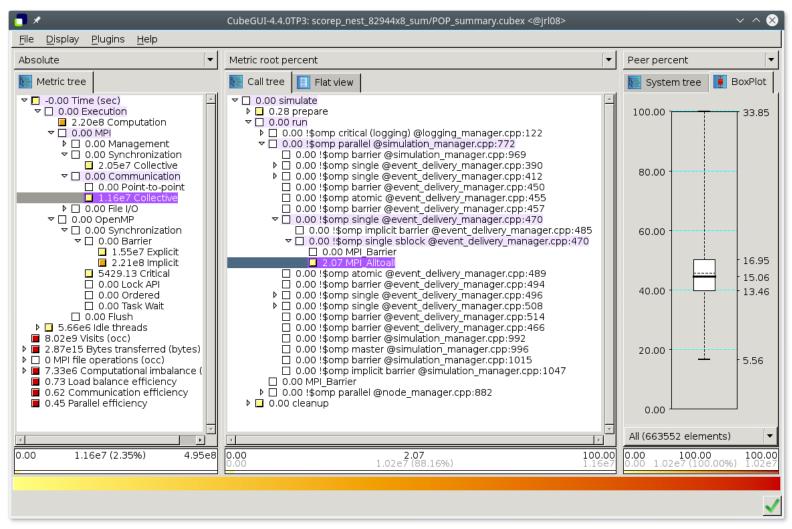
[2016]

- Nekbone
- CORAL benchmark
- JuQueen experiment
- 28,672 x 64 = 1,835,008 threads
- Load imbalance at OpenMP critical
- section



#### 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





Page 60

You KNOW YOU made it ...

# ... WHEN LARGE COMPANIES "COPY" YOUR STUFF



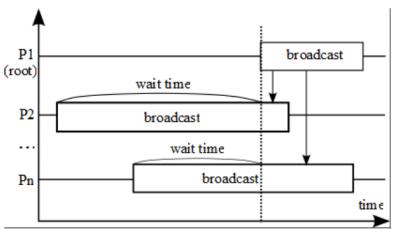
## Introducing the Intel® Trace Analyzer and Collector Performance Assistant

Motivation: Improve method of performance analysis via the GUI

#### Solution:

- Define common/known performance problems
- Automate detection via the Intel® Trace Analyzer

Example: A "Late Broadcast" is not easy to identify with existing views

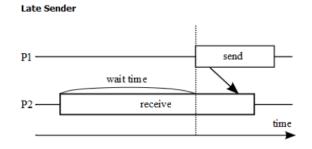


Source (no longer exists):

## Which Performance Issues are automatically identified?

Point-to-point exchange problems:

Late Sender

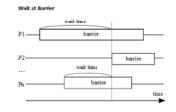


Late Receiver

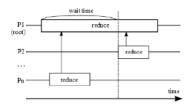
P1 wait time
P2 receive

Problems with global collective operation performance:

Wait at Barrier

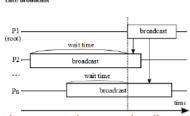


Early Reduce



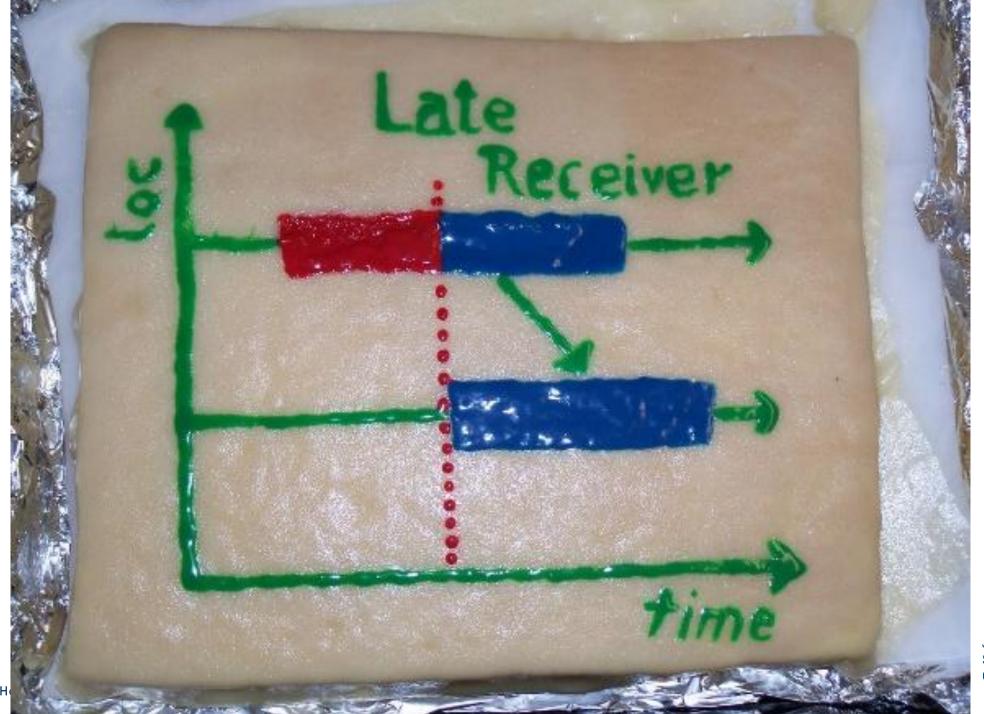
Optimization Notice

Late Broadcast



Source (no longer exists):

https://software.intel.com/en-us/videos/quickly-discover-performance-issues-with-the-intel-trace-analyzer-and-collector-90-beta



JÜLICH SUPERCOMPUTING CENTRE

#### **ASSESSMENT**

- Score-P / Scalasca installed on many HPC sites world-wide
- Used in daily work by performance analysts (e.g. POP CoE)
- User interface consistent for 24 years (but many enhancements)
- Support by vendors: Intel, AMD, (Siemens)
- Lots of work behind the scenes
  - Scalasca1 (Epilog) ⇒ scalasca2 (Score-P)
  - Constant bug fixing
  - Constant scaling improvements
  - Lately: MPI 3+4+5, OpenMP 3+4+5, OMPT, F2008 MPI interface, Pthreads, ...
  - GPU support: OpenMP target, OpenACC, CUDA HIP/ROCm, Kokkos, ...



#### **FUTURE WORK / OPPORTUNITIES (I)**

- Memory and vectorization performance analysis
  - Hard to capture performance data
    - Only possible if suitable hardware counters are provided
    - VERY processor specific ⇒ hard for open-source portable tools
    - ⇒ hard to attribute to data structures
- Trend towards task-based / asynchronous programming models
  - Very dynamic execution might be non reproducible ⇒ off-line tools fail
  - Timelines work on the node level, but hard to get the "big picture"
    - ⇒ good high-level metrics still missing here
  - 3-pane Cube display shows its limitations here ... ?!



#### **FUTURE WORK / OPPORTUNITIES (II)**

- Trend towards more modern programming languages (Python, C++)
  - How to automatically instrument template-based frameworks and programming styles?
  - Measurement overhead for small methods
  - How to present the data on Python level (and not on the interpreter low-level)?
- Performance assessment of data analytics codes
- Performance assessment of workflows
- Use of AI methods for performance analysis?!



### **QUESTIONS?**





