



Principles and Practice of Application Performance Measurement and Analysis on Parallel Systems

Lecture 2: Practical Performance Analysis and Tuning

1. July 2011 | Bernd Mohr

Institute for Advanced Simulation (IAS) Jülich Supercomputing Centre (JSC)

JÜLICH

CONTENT

- Fall-back
 - Simple timers
 - Hardware counter measurements
- Overview of some performance tools
 - · mpiP, TAU, Vampir + Vampirtrace
- Practical Performance Analysis and Tuning
- Challenges and open problems in performance optimization
 - Heterogeneous systems
 - · Automatic performance analysis: KOJAK/Scalasca
 - · Beyond execution time and flops
 - Extremely large HPC systems



Fall-back: Home-grown Performance Tools Silling College Control College Colleg

Timer: gettimeofday()



- UNIX function
- Returns wall-clock time in seconds and microseconds
- Actual resolution is hardware-dependent
- Base value is 00:00 UTC, January 1, 1970
- Some implementations also return the timezone

```
#include <sys/time.h>
struct timeval tv;
double walltime; /* seconds */
gettimeofday(&tv, NULL);
walltime = tv.tv_sec + tv.tv_usec * 1.0e-6;
```



Timer: clock_gettime()



- POSIX function
- For clock_id CLOCK_REALTIME returns wall-clock time in seconds and nanoseconds
- More clocks may be implemented but are not standardized
- Actual resolution is hardware-dependent

```
#include <time.h>
struct timespec tv;
double walltime; /* seconds */
Clock_gettime(CLOCK_REALTIME, &tv);
walltime = tv.tv_sec + tv.tv_nsec * 1.0e-9;
```

11-5

Timer: getrusage()



- UNIX function
- Provides a variety of different information
 - □ Including user time, system time, memory usage, page faults, etc.
 - □ Information provided system-dependent!



JÜLICH Timer: Others ■ MPI provides portable MPI wall-clock timer #include <mpi.h> double walltime: /* seconds */ walltime = MPI_Wtime(); □ Not required to be consistent/synchronized across ranks! ■ Same for OpenMP 2.0 (!) programming #include <omp.h> double walltime; /* seconds */ walltime = omp_get_wtime(); Hybrid MPI/OpenMP programming? □ Interactions between both standards (yet) undefined 11-7

Timer: Others Fortran 90 intrinsic subroutines cpu_time() system_clock() Hardware Counter Libraries Vendor APIs (PMAPI, HWPC, libhpm, libpfm, libperf, ...) PAPI



What Are Performance Counters



- Extra logic inserted in the processor to count specific events
- Updated at every cycle
- Strengths
 - Non-intrusive
 - □ Very accurate
 - □ Low overhead
- Weaknesses
 - □ Provides only hard counts
 - ☐ Specific for each processor
 - □ Access is not appropriate for the end user nor well documented
 - □ Lack of standard on what is counted

11-9

JÜLICH Hardware Counters Interface Issues Multi ■ Kernel level issues platform ☐ Handling of overflows ☐ Thread accumulation Hardware □ Thread migration counters □ State inheritance Kerne₁ ■ Multiplexing □ Overhead interface □ Atomicity Multi-platform interfaces ☐ The Performance API - PAPI University of Tennessee, USA LIKWID University of Erlangen, Germany



Hardware Measurement



- Typical measured events account for:
 - ☐ Functional units status
 - Float point operations
 - Fixed point operations
 - Load/stores
 - ☐ Access to memory hierarchy
 - ☐ Cache coherence protocol events
 - ☐ Cycles and instructions counts
 - ☐ Speculative execution information
 - Instructions dispatched
 - Branches mispredicted

II-11

Hardware Metrics



- Typical Hardware Counter
 - ☐ Cycles / Instructions
 - ☐ Floating point instructions
 - □ Integer instructions
 - □ Load/stores
 - □ Cache misses
 - □ TLB misses

- Useful derived metrics
 - □ IPC instructions per cycle
 - □ Float point rate
 - □ Computation intensity
 - ☐ Instructions per load/store
 - ☐ Load/stores per cache miss
 - □ Cache hit rate
 - □ Loads per load miss
 - □ Loads per TLB miss
- Derived metrics allow users to correlate the behavior of the application to one or more of the hardware components
- One can define threshold values acceptable for metrics and take actions regarding program optimization when values are below/above the threshold



Accuracy Issues



- Granularity of the measured code
 - ☐ If not sufficiently large enough, overhead of the counter interfaces may dominate
- Pay attention to what is **not** measured:
 - □ Out-of-order processors
 - □ Sometimes speculation is included
 - □ Lack of standard on what is counted
 - Microbenchmarks can help determine accuracy of the hardware counters

II-13

Hardware Counters Access on Linux



- Linux had not defined an out-of-the-box interface to access the hardware counters!
 - □ Linux Performance Monitoring Counters Driver (PerfCtr) by Mikael Pettersson from Uppsala X86 + X86-64
 - Needs kernel patching!
 - □ http://user.it.uu.se/~mikpe/linux/perfctr/
 - □ Perfmon by Stephane Eranian from HP IA64
 - It was being evaluated to be added to Linux
 - □ http://www.hpl.hp.com/research/linux/perfmon/
- Linux 2.6.31
 - □ Performance Counter subsystem provides an abstraction of special performance counter hardware registers



Utilities to Count Hardware Events



There are utilities that start a program and at the end of the execution provide overall event counts

 hpmcount (IBM)
 CrayPat (Cray)
 pfmon from HP (part of Perfmon for Al64)
 psrun (NCSA)
 cputrack, har (Sun)
 perf ex, ssrun (SGI)
 perf (Linux 2.6.31)

II-15

Hardware Counters: PAPI



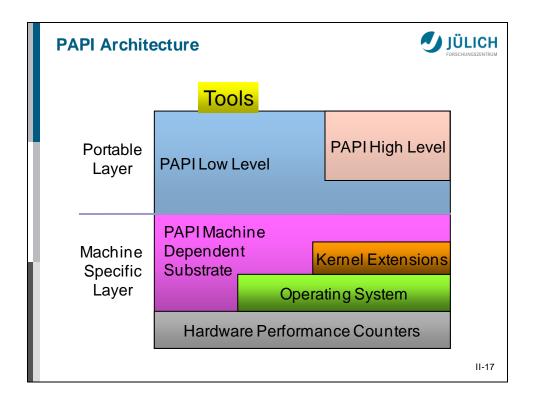
 Parallel Tools Consortium (PTools) sponsored project





- Performance Application Programming Interface
- Two interfaces to the underlying counter hardware:
 - ☐ The high-level interface simply provides the ability to start, stop and read the counters for a specified list of events
 - ☐ The low-level interface manages hardware events in user defined groups called EventSets
- Timers and system information
- C and Fortran bindings
- Experimental PAPI interface to performance counters support in the linux 2.6.31 kernel
- http://icl.cs.utk.edu/papi/





PAPI Predefined Events



- Common set of events deemed relevant and useful for application performance tuning (wish list)
 - □ papiStdEventDefs.h
 - □ Accesses to the memory hierarchy, cache coherence protocol events, cycle and instruction counts, functional unit and pipeline status
 - □ Run PAPI papi_avail utility to determine which predefined events are available on a given platform
 - ☐ Semantics may differ on different platforms!
- PAPI also provides access to native events on all supported platforms through the low-level interface
 - □ Run PAPI papi_native_avail utility to determine which predefined events are available on a given platform



```
JÜLICH
PAPI Preset Listing
  (derose@jaguar1) 184% papi_avail
  LibLustre: NAL NID: 0005dc02 (2)
Lustre: OBD class driver Build Version: 1, info@clusterfs.com
  Test case avail.c: Available events and hardware information.
  Vendor string and code
  Model string and code : AMD K8 (13)
CPU Revision : 1.000000
  CPU Revision
                              : 2400.000000
  CPU Megahertz
  CPU's in this Node
Nodes in this System
  Total CPU's
   Number Hardware Counters : 4
  Max Multiplex Counters : 32
  Name
                    Code
                                      Avail Deriv
                                                       Description (Note)
                    0x80000000
  PAPI L1 DCM
                                               Yes
                                                       Level 1 data cache misses ()
Level 1 instruction cache misses ()
  PAPI_L1_ICM
                    0x80000001
                                      Yes
                                               Yes
  PAPI_L2_DCM
                    0x80000002
                                      Yes
                                                       Level 2 data cache misses ()
                                               No
  PAPI_L2_ICM
                    0x80000003
                                               No
No
                                                        Level 2 instruction cache misses ()
                                      Yes
                                                       Level 3 data cache misses ()
  PAPI_L3_DCM
                    0x80000004
  PAPI_L3_ICM
                    0x80000005
                                                        Level 3 instruction cache misses ()
                                                       Level 1 cache misses ()
Level 2 cache misses ()
  PAPI_L1_TCM
                    0x80000006
                                      Yes
                                               Yes
                    0x80000007
  PAPI_L2_TCM
                                      Yes
                                               Yes
                    0x80000008
  PAPI_L3_TCM
                                                       Level 3 cache misses ()
                                      No
                                               No
                                                                                                 11-20
```



```
Example: papi_avail -e PAPI_L1_TCM (AMD Opteron JÜLICH
Event name:
                                                            PAPI_L1_TCM
 Event Code:
Number of Native Events:
                                                            0x80000006
                                                             L1 cache misses
 Short Description:
                                                             Level 1 cache misses
Long Description:
Developer's Notes:
Developer's Notes: ||
Derived Type: | DERIVED_ADD|
Postfix Processing String: ||
| Native Code[0]: 0x40001e1c DC_SYS_REFILL_MOES|
| Number of Register Values: 2|
| Register[0]: 0x20f P3 Ctr Mask|
| Register[1]: 0x1e43 P3 Ctr Code|
| Native Event Description: | Refill from system. Cache bits: Modified Owner Exclusive Shared|
    Native Code[1]: 0x40000037
                                                       IC_SYS_REFILL
   Number of Register Values: 2|
Register[0]: 0xf P3 Ctr Mask|
Register[1]: 0x83 P3 Ctr Code|
   | Native Code[2]: 0x40000036 | IC_L2_REFILL | Number of Register Values: 2| Register[0]: 0xf | P3 Ctr Mask | Register[1]: 0x82 | P3 Ctr Code |
   |Register[1]: 0x82 P3 Ctr Code
|Native Event Description: |Refill from L2|
   | Native Code[3]: 0x40001e1b | DC_L2_REFILL_MOES|
| Number of Register Values: 2|
| Register[0]: 0x20f | P3 Ctr Mask|
| Register[1]: 0x1e42 | P3 Ctr Code|
| Native Event Description: | Refill from L2. Cache bits: Modified Owner Exclusive Shared
                                                                                                                                                                  11-21
```

```
JÜLICH
PAPI papi_native_avail Utility (AMD Opteron)
(derose@sleet) 187% papi_native_avail |more
Test case NATIVE_AVAIL: Available native events and hardware information.
Vendor string and code : AuthenticAMD (2)
Model string and code : AMD K8 Revision : 10.000000
                                  : AMD K8 Revision C (15)
: 10.000000
CPU Megahertz
CPU's in this Node
Nodes in this System
Total CPU's
Number Hardware Counters
                                   : 2193.406982
 Max Multiplex Counters
The following correspond to fields in the PAPI_event_info_t structure.

Symbol Event Code Count
| Short Description|
| Long Description|
| Derived|
   PostFix
 The count field indicates whether it is a) available (count >= 1) and b) derived
(count > 1)
FP_ADD_PIPE
                      0x40000000
  | Dispatched FPU ops - Revision B and later revisions - Speculative multiply pipe ops excluding junk ops | Register Value[0]: 0xf P3 Ctr Mask | Register Value[1]: 0x200 P3 Ctr Code|
                                                                                                                          .
II-22
```



High Level API



- Meant for application programmers wanting simple but accurate measurements
- Calls the lower level API
- Allows only PAPI preset events
- Eight functions:
 - □ PAPI_num_counters
 - □ PAPI_start_counters, PAPI_stop_counters
 - □ PAPI_read_counters
 - □ PAPI_accum_counters
 - □ PAPI_flops
 - □ PAPI_flips, PAPI_ipc (New in Version 3.x)
- Not thread-safe (Version 2.x)

11-23

Example: Quick and Easy Mflop/s



```
program papiMflops
         parameter (N=1024) include "f77papi.h"
         integer*8 fpins
         real*4
                      realtm, cputime, mflops
         integer
                       ierr
                       a(N,N)
         real*4
         call random_number(a)
         call PAPIF_flops(realtm, cputime, fpins, mflops, ierr)
         do j=1,N
           do i=1,N
a(i,j)=a(i,j)*a(i,j)
            end do
         end do
         call PAPIF_flops(realtm, cputime, fpins, mflops, ierr)
print *,' realtime: ', realtm, ' cputime: ', cputime
print *,' papi_flops: ', mflops, ' Mflop/s'
% ./papiMflops
                     3.640159
                                                          3.630502
  realtime:
                                         cputime:
  papi_flops:
                       29.67809
                                           MFlops.
                                                                                       11-24
```



```
JÜLICH
General Events
        program papicount
            parameter (N=1024)
include "f77papi.h"
            integer*8 values(2)
            integer events(2), ierr
real*4 a(N,N)
            call random_number(a)
            events(1) = PAPI_L1_DCM
            events(2) = PAPI_L1_DCA
            call PAPIF_start_counters(events, 2, ierr)
            do j=1,N
do i=1,N
              a(i,j)=a(i,j)*a(i,j) end do
            end do
            call PAPIF_read_counters(values, 2, ierr)
print *,' L1 data misses : ', values(1)
print *,' L1 data accesses: ', values(2)
        end
 % ./papi count
                                                  13140168
    L1 data misses :
                                                 500877001
    L1 data accesses:
                                                                                       11-25
```

Low Level API



- Increased efficiency and functionality over the high level PAPI interface
- 54 functions
- Access to native events
- Obtain information about the executable, the hardware, and memory
- Set options for multiplexing and overflow handling
- System V style sampling (profil())
- Thread safe





CONTENT

- Fall-back
 - Simple timers
 - · Hardware counter measurements
- Overview of some performance tools
 - · mpiP, TAU, Vampir + Vampirtrace
- Practical Performance Analysis and Tuning
- Challenges and open problems in performance optimization
 - Heterogeneous systems
 - Automatic performance analysis: KOJAK/Scalasca
 - Beyond execution time and flops
 - Extremely large HPC systems

11-27

MPI Profiling: mpiP



■ Scalable, light-weight MPI profiling library



- Generates detailed text summary of MPI behavior
 - ☐ Time spent at each MPI function callsite
 - ☐ Bytes sent by each MPI function callsite (where applicable)
 - MPI I/O statistics
 - ☐ Configurable traceback depth for function callsites
- Controllable from program using MPI_Pcontrol
 - ☐ Allows you to profile just one code module or cycle
 - ☐ Allows mpiP profile dumps mid-run
- Uses PMPI interface ⇒ only re-link of application necessary
- http://mpip.sourceforge.net/



```
JÜLICH
mpiP Text Output Example
@ Version: 3.1.1
// 10 lines of mpiP and experiment configuration options
// 8192 lines of task assignment to BlueGene topology information
@--- MPI Time (seconds) -----
Task AppTime MPITime
                                                 MPI%
                  37.7
                                                 66.89
     0
                                    25.2
8191
        37.6
3.09e+05
                           26
2.04e+05
@--- Callsites: 26 -----
ID Lev File/Address Line Parent_Funct
1 0 coarsen.c 542 hypre_StructCoarsen
                                                                                         MPI_Call Waitall
 // 25 similiar lines
      - Aggregate Time (top twenty, descending, milliseconds)
l Site Time App% MPI% Contail 21 1.03e+08 33.27 50.49 0.10
tall 1 2.88e+07 9.34 14.17 0.10
Call
                                                                                            ĆOV
                                                                                           0.11
 Waitall
 Waitall
 // 18 similiar lines
                                                                                                         11-29
```

```
JÜLICH
mpiP Text Output Example (cont.)
      Aggregate Sent Message Size (top twenty, descending, bytes)
Site Count Total Avrg Sent%
d 11 845594460 7.71e+11 912 59.92
Call
Isend
Allreduce
// 6 similiar lines
                              10
                                         49152
                                                    3.93e+05
                                                                                0.00
@--- Callsite Time statistics (all, milliseconds): 212992 ----
                                             Max Mean Min App% MPI%
275 0.1 0.000707 29.61 44.27
           Site Rank Count
21 0 111096
Name
Waitall
Waitall 21 8191 65799 882 0.24 0.000707 41.98 60.66 Waitall 21 * 577806664 882 0.178 0.000703 33.27 50.49 // 213,042 similar lines
 @--- Callsite Message Sent statistics (all, sent bytes) -
        Site Rank Count Max
11 0 72917 2.621e+05
                                               Max Mean
Name
                                                                      8 6.206e+07
Isend
                                                           851.1
Isend 11 8191 46651 2.621e+05 11 * 845594460 2.621e+05 // 65,550 similiar lines
                                                            1029
                                                                            4.801e+07
                                                                        8 7.708e+11
                                                           911.5
Tsend
                                                                                     11-30
```



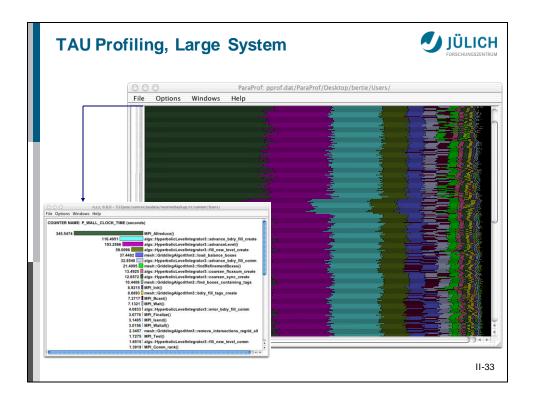
"Swiss Army Knife" of **JÜLICH Performance Analysis: TAU** Very portable tool set for instrumentation, measurement and analysis of parallel multi-threaded applications Instrumentation API supports choice □ between profiling and tracing □ of metrics (i.e., time, HW Counter (PAPI)) Uses Program Database Toolkit (PDT) for C, C++, Fortran source code instrumentation Supports □ Languages: C, C++, Fortran 77/90, HPF, HPC++, Java, Python ☐ Threads: pthreads, Tulip, SMARTS, Java, Win32, OpenMP ☐ Systems: same as KOJAK + Windows + MacOS + ... http://tau.uoregon.edu/ http://www.cs.uoregon.edu/research/pdt/ II-31

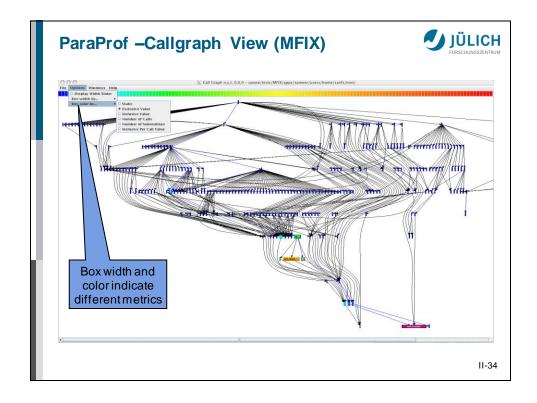
TAU Instrumentation



- Flexible instrumentation mechanisms at multiple levels
 - □ Source code
 - manual
 - automatic
 - □ C, C++, F77/90/95 (Program Database Toolkit (PDT))
 - □ OpenMP (directive rewriting with Opari)
 - □ Object code
 - pre-instrumented libraries (e.g., MPI using PMPI)
 - statically-linked and dynamically-loaded (e.g., Python)
 - □ Executable code
 - dynamic instrumentation (pre-execution) (DynInst)
 - virtual machine instrumentation (e.g., Java using JVMPI)
- Support for performance mapping
- Support for object-oriented and generic programming









TAU Usage



- Usage:
 - □ Specify programming model by setting TAU_MAKEFILE to one of \$<TAU_ROOT>/<arch>/lib/Makefile.tau-* Examples from Linux cluster with Intel compiler and PAPI:
 - MPI: Makefile.tau-icpc-papi-mpi-pdt
 - OMP: Makefile.tau-icpc-papi-pdt-openmp-opari
 - OMPI: Makefile.tau-icpc-papi-mpi-pdt-openmp-opari
 - □ Compile and link with
 - tau_cc.sh file.c ...
 - tau_cxx.sh file.cxx...
 - tau_f90.sh file.f90 ...
 - □ Execute with real input data
 Environment variables control measurement mode
 - TAU_PROFILE, TAU_TRACE, TAU_CALLPATH, ...
 - □ Examine results with paraprof

Vampirtrace MPI Tracing Tool



- Library for Tracing of MPI and Application Events
- □ Records MPI-1 point-to-point and collective communication
 - □ Records MPI-2 I/O operations and RMA operations
 - □ Records user subroutines
- Uses the standard MPI profiling interface
- Usage:
 - □ Compile and link with
 - vtcc -vt:cc mpicc file.c ...
 - vtcxx -vt:cxx mpicc file.cxx...
 - vtf90 -vt:f90 mpif90 file.f90 ...
 - □ Execute with real input data (⇒ generates <exe>.otf)



Vampirtrace MPI Tracing Tool



 Versions up to 4.0 until 2003 were commercially distributed by PALLAS as VampirTrace



- Current status
 - □ Commercially distributed by Intel
 - Version 5 and up distributed as Intel Trace Collector



- For Intel based platforms only
- http://software.intel.com/en-us/articles/intel-cluster-toolkit/
- □ New open-source VampirTrace version 5 distributed by Technical University Dresden
 - Based on KOJAK's measurement system with OTF backend
 - http://www.tu-dresden.de/zih/vampirtrace/
 - Is also distributed as part of Open MPI



11-37

Visualization and Analysis of MPI Programs Commercial product Originally developed by Forschungszentrum Jülich Originally developed by Forschungszentrum Jülich



Vampir Event Trace Visualizer



 Versions up to 4.0 until 2003 were commercially distributed by PALLAS as Vampir



- Current status
 - □ Commercially distributed by Intel
 - Version 4 distributed as Intel Trace Analyzer



- For Intel based platforms only
- Intel meanwhile released own new version (ITA V6 and up)
- http://software.intel.com/en-us/articles/intel-cluster-toolkit/
- □ Original Vampir (and new VNG) commercially distributed by Technical University Dresden
 - http://www.vampir.eu

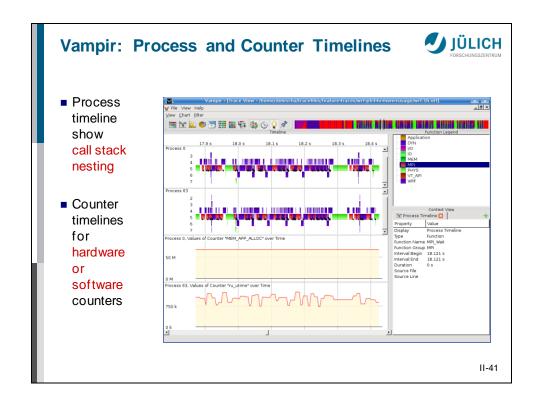


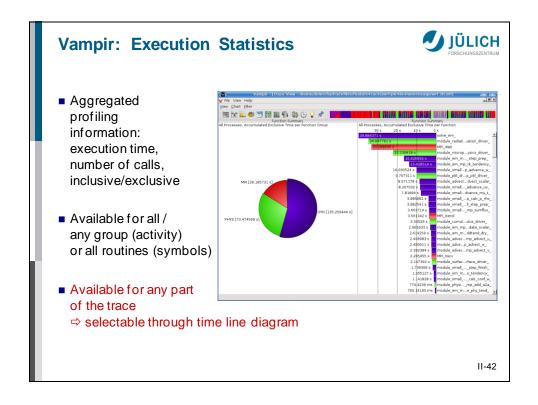
11-39

JÜLICH Vampir: Time Line Diagram Functions organized into groups coloring ocess 13 by group ocess 22 cess 25 cess 31 Message ocess 37 lines can ocess 43 be colored ocess 49 by tag or ocess 55 size Information about states, messages, collective and I/O operations

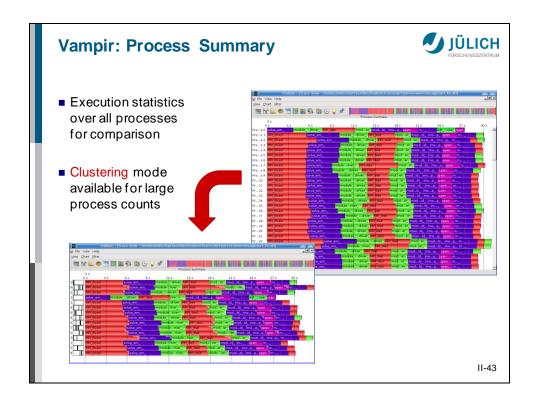
available through clicking on the representation

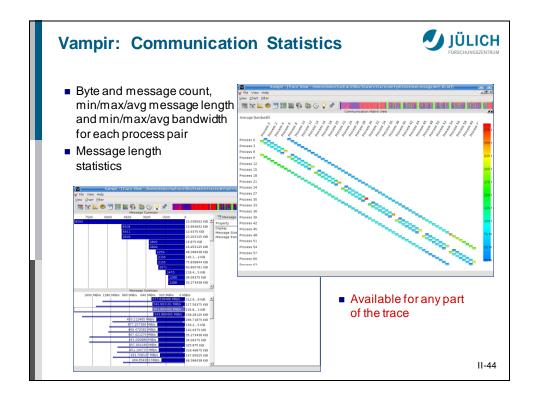














Other Profiling Tools



- gprof
 - □ Available on many systems
 - □ Compiler instrumentation (invoked via –g –pg)
- FPMPI-2 (ANL)
 - □ http://www.mcs.anl.gov/fpmpi/
 - ☐ MPI profiler (invoked via re-linking)
 - □ special: Optionally identifies synchronization time
 - □ single output file: count, sum, avg, min, max over ranks
- ThreadSpotter (Rogue Wave) [commercial product]
 - □ http://www.roguewave.com/products/threadspotter.aspx
 - □ Sampling-based memory and thread performance analyzer
 - ☐ Works on un-modified, optimized executables

Other Profiling Tools II



- ompP (UC Berkeley)
 - □ http://www.ompp-tool.com
 - □ OpenMP profiler (invoked via OPARI source instrumentation)
- HPCToolkit (Rice University)
 - □ http://www.hpctoolkit.org
 - ☐ Multi-platform sampling-based callpath profiler
 - ☐ Works on un-modified, optimized executables
- Open|SpeedShop (Krell Insitute with support of LANL, SNL, LLNL)
 - □ http://www.openspeedshop.org
 - ☐ Comprehensive performance analysis environment
 - ☐ Uses binary instrumentation (via DynInst)



Other Tracing Tools MPE / Jumpshot (ANL) Part of MPICH2 Invoked via re-linking Only supports MPI P2P and collectives; SLOG2 trace format Extrae / Paraver (BSC/UPC) http://www.bsc.es/paraver Measurement system (Extrae) and visualizer (Paraver) Powerful filter and summarization features Very configurable visualization

JÜLICH

CONTENT

- Fall-back
 - Simple timers
 - · Hardware counter measurements
- Overview of some performance tools
 - mpiP, TAU, Vampir + Vampirtrace
- Practical Performance Analysis and Tuning
- Challenges and open problems in performance optimization
 - Heterogeneous systems
 - Automatic performance analysis: KOJAK/Scalasca
 - · Beyond execution time and flops
 - · Extremely large HPC systems



Practical Performance Analysis and Tuning JÜLICH Successful tuning is combination of □ Right algorithm and libraries ☐ Compiler flags and pragmas / directives (Learn and use them) □ THINKING Measurement is better than reasoning / intuition (= guessing) ☐ To determine performance problems ☐ To validate tuning decisions / optimizations (after each step!)

Practical Performance Analysis and Tuning JÜLICH



11-49

- It is easier to optimize a slow correct program than to debug a fast incorrect one
 - ⇒ Debugging before Tuning
 - ⇒ Nobody really cares how fast you can compute the wrong answer
- The 80/20 rule
 - ☐ Program spents 80% time in 20% of code
 - □ Programmer spends 20% effort to get 80% of the total speedup possible in the code
 - ⇒ Know when to stop!
- Don't optimize what doesn't matter
 - ⇒ Make the common case fast



Typical Performance Analysis Procedure



- 1. Do I have a performance problem at all?
 - ⇒ Time / hardware counter measurements
 - ⇒ Speedup and scalability measurements
- 2. What is the main bottleneck (computation/communication...)?
 - ⇒ Flat profiling (sampling / prof)
- 3. Where is the main bottleneck?
 - ⇒ Call graph profiling (gprof)
 - ⇒ Detailed (basic block) profiling
- 4. Why is it there?
 - ⇒ Hardware counters analysis
 - ⇒ Trace <u>selected</u> parts to keep trace files manageable
- 5. Does my code have scalability problems?
 - ⇒ Load Imbalance analysis
 - ⇒ Profile code for typical small and large processor count
 - ⇒ Compare profiles function-by-function

II-51

CONTENT



- Fall-back
 - Simple timers
 - · Hardware counter measurements
- Overview of some performance tools
 - · mpiP, TAU, Vampir + Vampirtrace
- Practical Performance Analysis and Tuning
- Challenges and open problems in performance optimization
 - Heterogeneous systems
 - Automatic performance analysis: KOJAK/Scalasca
 - Beyond execution time and flops
 - Extremely large HPC systems



Heterogenous Systems

••	
JÜL	ICH
, – –	

- Current trend to use hardware acceleration to speedup calculations
 - ☐ IBM Cell (e.g. LANL Roadrunner)
 - □ Clearspead
 - □ FPGA-based acceleration
 - ☐ GPU-based acceleration
- In the long run: new programming models needed
- Very little to not existing tool support
- ⇒Tool research opportunities!

11-53

GPU Performance Tools



- CUDA profiler
- TAU
 - □ CUDA + OpenCL profiling (host side)
- VampirTrace
 - □ CUDA tracing (host side)



Current Related Projects



 Hybrid Programming for Heterogeneous Architectures



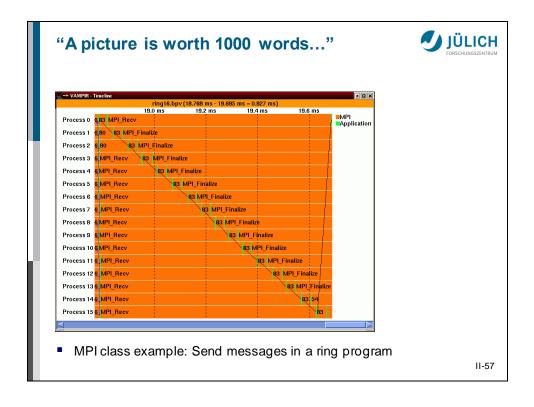
- □ EU ITEA2 funded project 2010 2013
- ☐ Successor of successful ParMA project
- 25 partners for France, Germany, Spain, Sweden including Bull, TUD, BSC, USQV, JSC, Acumem
- □ Develop programming models and tools
- ☐ Evaluation with large set of industrial codes

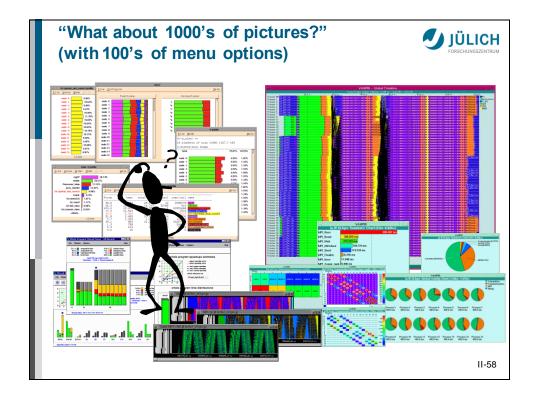
JÜLICH FORSCHUNGSZENTRUM

CONTENT

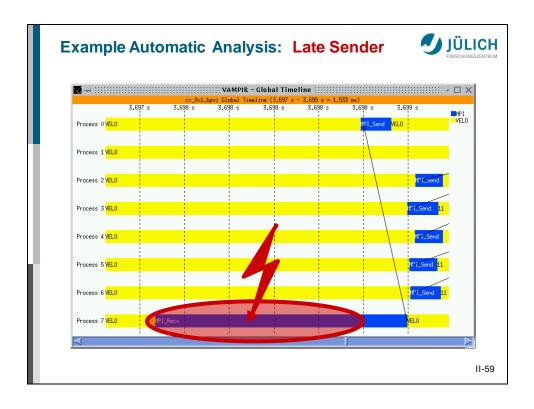
- Fall-back
 - Simple timers
 - · Hardware counter measurements
- Overview of some performance tools
 - mpiP, TAU, Vampir + Vampirtrace
- Practical Performance Analysis and Tuning
- Challenges and open problems in performance optimization
 - Heterogeneous systems
 - · Automatic performance analysis: KOJAK/Scalasca
 - Beyond execution time and flops
 - · Extremely large HPC systems

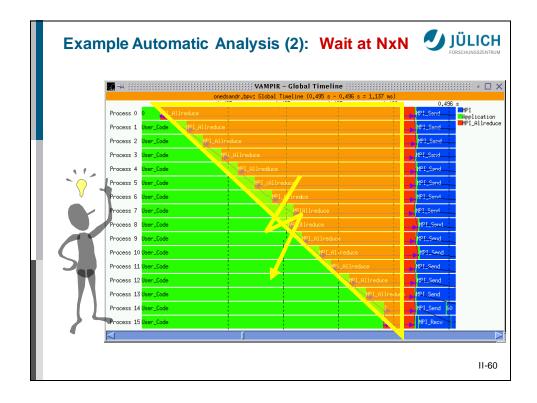




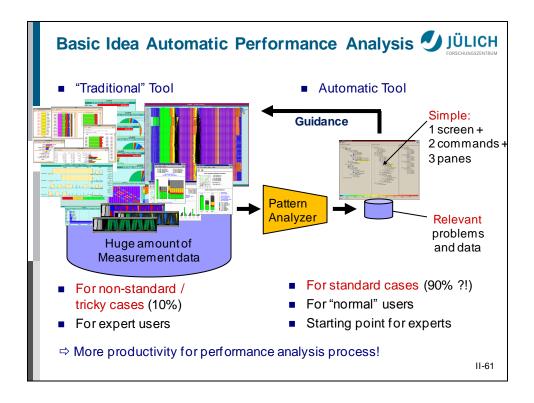












The KOJAK Project

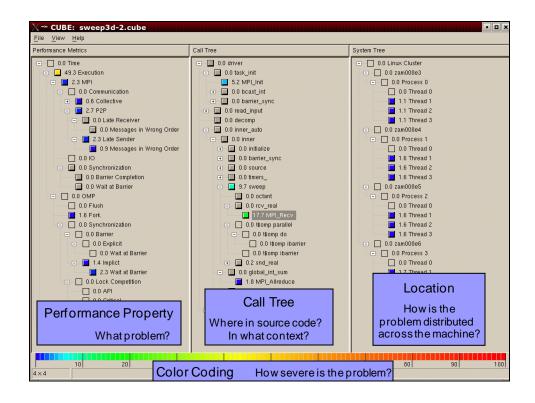


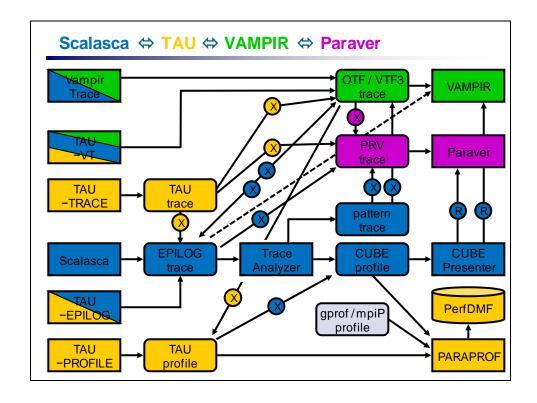
- Kit for Objective Judgement and Automatic Knowledge-based detection of bottlenecks
- Forschungszentrum Jülich
- Innovative Computing Laboratory, TN
- Started 1998



- □ Instrument C, C++, and Fortran parallel applications
 - Based on MPI, OpenMP, SHMEM, or hybrid
- □ Collect event traces
- □ Search trace for event patterns representing inefficiencies
- □ Categorize and rank inefficiencies found
- http://www.fz-juelich.de/jsc/kojak/











CONTENT

- Fall-back
 - Simple timers
 - · Hardware counter measurements
- Overview of some performance tools
 - · mpiP, TAU, Vampir + Vampirtrace
- Practical Performance Analysis and Tuning
- Challenges and open problems in performance optimization
 - Heterogeneous systems
 - · Automatic performance analysis: KOJAK/Scalasca
 - Beyond execution time and flops
 - Extremely large HPC systems

II-65

Beyond Execution Time and Flops



- High performance on today's architectures requires extremely effective use of the cache and memory hierarchy of the system
 - □ Very complex designs
 - □ Will get "worse" with advanced multi- and many core chips
- tools for memory performance analysis needed
 - □ Current approach:
 - Measure cache, memory, TLB events (per execution unit)
 - Rarely tell what data object(s) are cause of the problem
 - Main problem:
 - if cause known, what is the fix?
 - Or worse, what is the portable fix?
 - ⇒ Should be really better left to the compiler!



Beyond Execution Time and Flops II

JÜLICH	
EORSCHIINGSZENTRII	١.

- Do we need tools for I/O performance analysis?
 - □ Currently only very few tools available (CrayPat, Pablo I/O, ...)
 - □ Perhaps scientific programmers only need training in parallel I/0 facilities already available today?
- Tools for new programming paradigms?
 - ☐ Obvious next candidate: one-sided communication
 - SHMEM, MPI-2 RMA, ARMCI, ...
 - □ "Easy" to monitor; intercept calls e.g., using wrapper funcs
 - CAF, UPC
 - □ "Harder"; probably need compiler support
 - □ Other candidates?

11-67

Beyond Execution Time and Flops III



- Tool (sets and environments) must be able to handle "hybrid" cases!
 - ☐ Message passing (MPI)
 - ☐ Multi-threading (OpenMP, pthreads, ...)
 - □ One-sided communication
 - □ (Parallel) I/O
 - □...
- BIGGEST challenge:
 - Many tools for instrumentation, measurement, analysis, and visualization
 - ☐ What about (automatic) optimization tools?





CONTENT

- Fall-back
 - Simple timers
 - · Hardware counter measurements
- Overview of some performance tools
 - · mpiP, TAU, Vampir + Vampirtrace
- Practical Performance Analysis and Tuning
- Challenges and open problems in performance optimization
 - · Heterogeneous systems
 - · Automatic performance analysis: KOJAK/Scalasca
 - · Beyond execution time and flops
 - Extremely large HPC systems

II-69

Increasing Importance of Scaling



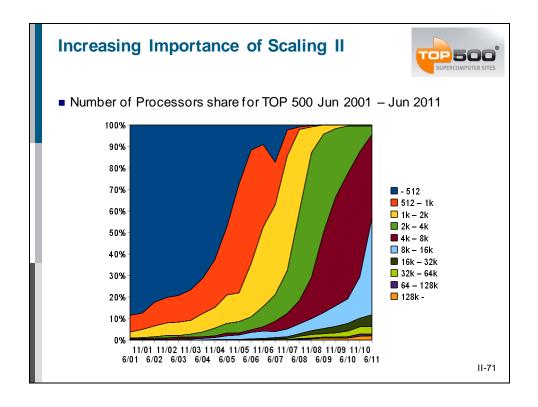
Number of Processors share for TOP 500 Jun 2011

NProc	Count	Share	∑Rmax	Share	∑NProc
1025-2048	2	0.4%	168 TF	0.3%	2,632
2049-4096	20	4.0%	1,177 TF	2.0%	71,734
4097-8192	195	39.0%	9,759 TF	16.6%	1,262,738
8193-16384	224	44.8%	15,216 TF	25.8%	2,337,998
> 16384	59	11.8%	32,556 TF	55.3%	4,100,234
Total	500	100%	58,876 TF	100%	7,775.336

■ Average system size: 15,551 cores

■ Median system size: 8,556 cores





System attributes	2010	"2015"		"2018"		Difference 2010 & 2018
System p eak	2 Pflop/s	200 Pflop/s		1 Eflop/sec		O(1000)
Power	6 MW	15 MW		~20	~20 MW	
System memory	0.3 PB	5	PB	32-64 PB		O(100)
Node performance	125 GF	0.5 TF	7TF	1 TF	10 TF	O(10) – O(100)
Node memory BW	25 GB/s	0.1 TB/sec	1 TB/sec	0.4 TB/sec	4 TB/sec	O(100)
Node concurrency	12	O(100)	O(1,000)	O(1,000)	O(10,000)	O(100) – O(1000)
Total Concurrency	225,000	O(O(10 ⁸) O(10 ⁹)		10º)	O(10,000)
Total Node Interconnect BW	1.5 GB/s	20 GB/sec O(1day)		200 GB/sec		O(100)
MTTI	days			O(1 day)		- O(10)



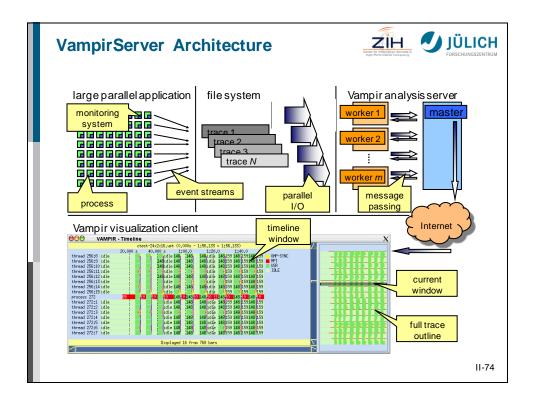
Extremely large HPC systems



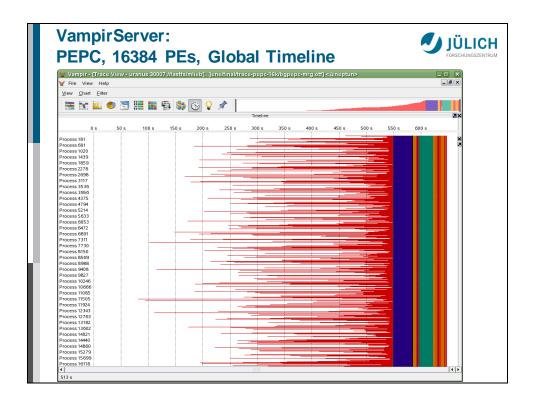
- Many HPC computing centers have systems with some 1000s of PEs
 Today's tools (e.g., gprof, Vampir, TAU) can handle 500–1000 PEs
 Measurements give enough insight to optimize for 2000-4000 PEs
- Now: IBM BlueGene, Cray XT5 ⇒ some 100,000s PEs
 - □ Small customer base
 - □ Very unique performance problems:
 - Why does it work on 16.000 PEs but not for 32.000 PEs?
 - □ Load balance problems amplify performance degradation
 - □ Interesting challenging problem,

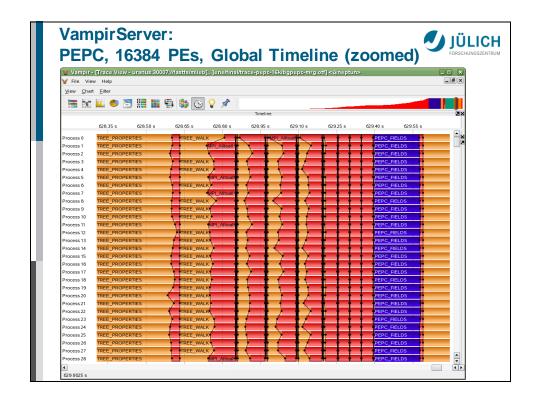
but given scarce resources in tools research? ...

- Lots of data ⇒ need better data mgmt, parallel analysis!?
- Bigger problem: scalable result displays / visualizations

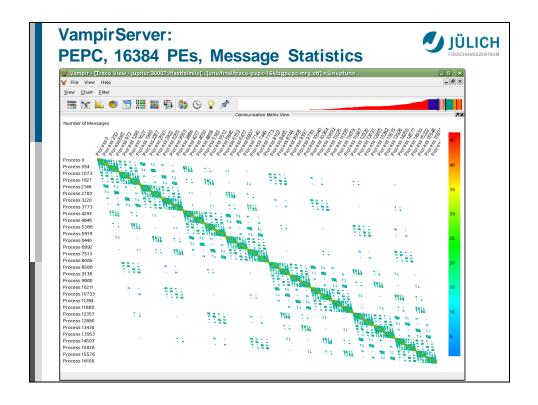


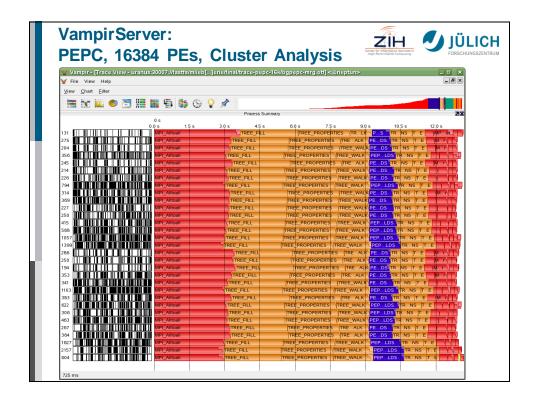














Vampirserver





- Parallel client/server version of Vampir
 - □ Can handle much larger trace files
 - □ Remote visualization possible
- Usage
 - ☐ Start VampirServer daemon
 - vngd -n <t> # on frontend, uses <t> threads
 - mpiexec -n vngd # cluster, uses processes
 - ☐ Start the Vampir visualizer (vampir)
 - on local system if available
 - on frontend in 2nd shell

then connect client to server, load and analyze trace file <exe>.otf

finally: vngd-shutdown

The Scalasca Project



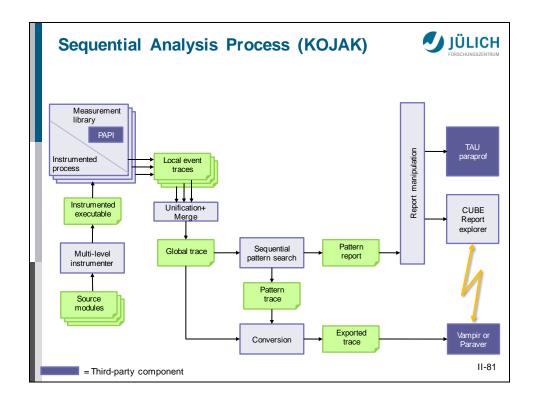
- Scalable Analysis of Large Scale Applications
- scalasca 🗖
- Follow-up project to KOJAK
- http://www.scalasca.org/

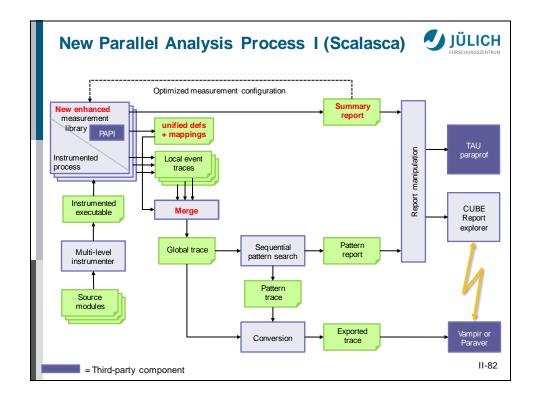
- Approach
 - □ **Instrument** C, C++, and Fortran parallel applications
 - Based on MPI, OpenMP, SHMEM, or hybrid
 - ☐ Option 1: scalable call-path profiling
 - □ Option 2: scalable event trace analysis
 - Collect event traces
 - Search trace for event patterns representing inefficiencies
 - Categorize and rank inefficiencies found
- Supports MPI 2.2 and basic OpenMP



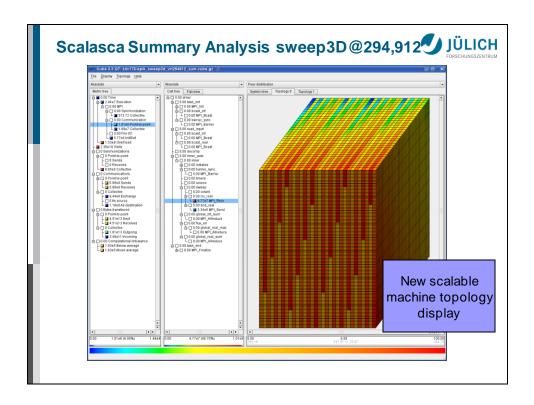


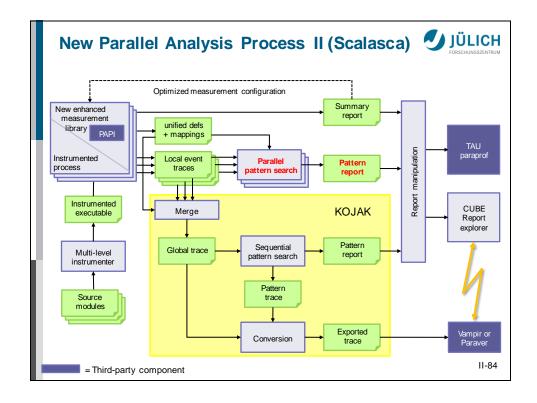














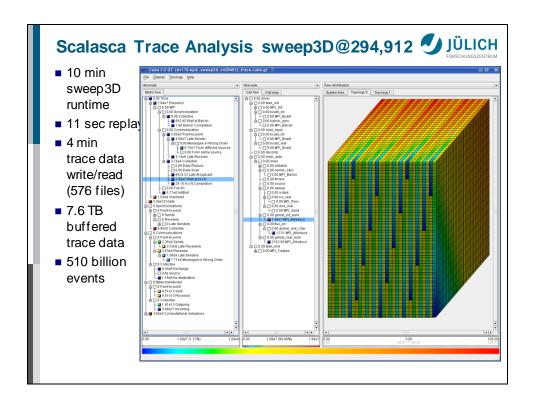
11-85

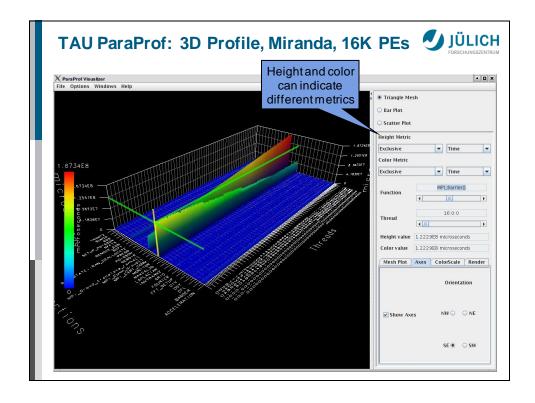
Parallel pattern search to address wide traces As many processes / threads as used to run the application ⇒ Can run in same batch job!! Each process / thread responsible for its "own" local trace data Idea: "parallel replay" of application Analysis uses communication mechanism that is being analyzed Use MPI P2P operation to analyze MPI P2P communication, use MPI collective operation to analyze MPI collectives, ... Communication requirements not significantly higher and (often lower) than requirements of target application In-memory trace analysis Available memory scales with number of processors used

□ Local memory usually large enough to hold local trace

JÜLICH Example: Late Sender Sequential approach (EXPERT) ☐ Scan the trace in sequential order waiting □ Watch out for receive event □ Use links to access other constituents □ Calculate waiting time time New parallel approach (SCOUT) ☐ Each process identifies local constituents Sender sends local constituents to receiver □ Receiver calculates waiting time 11-86









Performance Tuning: Still a Problem?	JÜLICH FORSCHUNGSZENTRUM
[Slide left intentionally blank]	
	11-89

