

Performance Analysis with Vampir



Outline

Part I: Performance Analysis with Vampir & Score-P

- Introduction/Performance Engineering
- Score-P Measurement System
- Vampir Performance Data Visualization
- Part II: Vampir Analysis Exercise
 - Analysing Four Application Traces
- Part III: Vampir Hands-On
 - Visualizing and Analyzing NPB-MZ-MPI / BT





Introduction Performance Engineering





Disclaimer

It is extremely easy to waste performance!

- Bad MPI (50-90%)
- No node-level parallelism (94%)
- No vectorization (75%)
- Bad memory access pattern (99%)
- Total: $1/2 \times 1/16 \times 1/4 * 1/100$
- -~1/10000 of the peak performance

Disclaimer (2)

Performance tools will not automatically make your code run faster. They help you understand, what your code does and where to put in work.

Virtual Institute – High Productivity Supercomputing

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

by Helmholtz Association of German Research Centres

- Activities
 - Development and integration of HPC programming tools
 - Correctness checking & performance analysis
 - Academic workshops
 - Training workshops
 - Service
 - Support email lists
 - Application engagement

http://www.vi-hps.org



VI-HPS partners (founders)



Forschungszentrum Jülich

RWTH Aachen University

Jülich Supercomputing Centre





Centre for Information Services & HPC



- University of Tennessee (Knoxville)
 - Innovative Computing Laboratory









VI-HPS partners (cont.)



Arm Ltd.

Allinea Software



- Barcelona Supercomputing Center
 - Centro Nacional de Supercomputación
- Lawrence Livermore National Lab.
 - Center for Applied Scientific Computing



Leibniz Supercomputing Centre



- Technical University of Darmstadt
 - Laboratory for Parallel Programming

allinea









VI-HPS partners (cont.)





Chair for Computer Architecture



- University of Oregon
 - Performance Research Laboratory
- University of Stuttgart
 - HPC Centre



University of Versailles St-Quentin LRC ITACA







Productivity tools

• MUST & Archer

MPI & OpenMP usage correctness checking

PAPI

- Interfacing to hardware performance counters
- Periscope Tuning Framework
 - Automatic analysis and Tuning

Scalasca

Large-scale parallel performance analysis

- TAU

Integrated parallel performance system

Vampir

Interactive graphical trace visualization & analysis

Score-P

Community-developed instrumentation & measurement infrastructure

See VI-HPS Tools Guide for a overview of tools: https://www.vi-hps.org/tools



Productivity tools (cont.)

- FORGE DDT/MAP/PR: Parallel debugging, profiling & performance reports
- Extra-P: Automated performance modelling
- Kcachegrind: Callgraph-based cache analysis [x86 only]
- JUBE: Workflow execution environment
- MAQAO: Assembly instrumentation & optimization [x86-64 only]
- mpiP/mpiPview: MPI profiling tool and analysis viewer
- Open MPI: Integrated memory checking
- Open|SpeedShop: Integrated parallel performance analysis environment
- Paraver/Dimemas/Extrae: Event tracing, graphical trace visualization & analysis
- Rubik: Process mapping generation & optimization [BG only]
- SIONlib/Spindle: Optimized native parallel file I/O & shared library loading
- STAT: Stack trace analysis tools

Technologies and their integration



Recording and Studying a Program's Behavior

- Performance analysis in practice:
 - Attach measurement system Score-P to application
 - Run application with Score-P monitor
 - Result: trace/profile data
- Study the trace data with Vampir
- Repeat
 - Adapt instrumentation ("what you measure")
 - Evaluate result of a change









- Running program is periodically interrupted to take measurement
- Statistical report of program behavior
 - Not very detailed information on highly volatile metrics
 - Requires long-running applications
- Works with unmodified executables



Methodology: Instrumentation





- Measurement code is inserted such that every event of interest is captured directly
- Advantage:
 - Much more detailed information
- Disadvantage:
 - Processing of source-code / executable necessary
 - Large relative overheads for small functions

Instrumentation vs. Sampling

Instrumentation/Tracing:

- Typically more difficult than using a profiler
- Does not garuantee anything about overhead or recording size
 - It depends on the event rate
 - E.g. an MPI-only trace has very low overhead
 - Filtering required
- Sampling:
 - Does not give an absolutely accurate picture of a run
 - Cannot count function calls
 - Cannot record exact timings
 - Cannot record exact performance counters
 - It is statistical sampling
 - It cannot capture semantics of APIs, i.e. it cannot follow API usage and analyze passed arguments, e.g. transferred bytes



– Timelines



Terms Used and How They Connect





Score-P Measurement System



Fragmentation of tools landscape

- Several performance tools co-exist
 - Separate measurement systems and output formats
- Complementary features and overlapping functionality
- Redundant effort for development and maintenance
 - Limited or expensive interoperability
- Complications for user experience, support, training

Vampir	Scalasca	TAU	Periscope
VampirTrace	EPILOG /	TAU native	Online
OTF	CUBE	formats	measurement

Score-P project idea

- Start a community effort for a common infrastructure
 - Score-P instrumentation and measurement system
 - Common data formats OTF2 and CUBE4

Partners:

- TU Dresden, FZ Jülich, TU München, University of Oregon, RWTH Aachen, TU Darmstadt
- Developers save manpower by sharing development resources
- Users need only a single installation, resulting in single learning curve
- Ensure a single official release version at all times which will always work with the tools
- Commitment to joint long-term cooperation
 - Development based on meritocratic governance model
 - Open for contributions and new partners



Bundesministerium für Bildung und Forschung

GEFÖRDERT VON





Data Collection with Score-P



Instrument (Compile & Link)





Data Collection with Score-P



Measurements are configured via environment variables

```
$ scorep-info config-vars --full
SCOREP_ENABLE_PROFILING
[...]
SCOREP_ENABLE_TRACING
[...]
SCOREP_TOTAL_MEMORY
Description: Total memory in bytes for the measurement system
[...]
SCOREP_EXPERIMENT_DIRECTORY
Description: Name of the experiment directory
[...]
```



Data Collection with Score-P



Example for generating and loading a trace:

<pre>\$ export SCOREP_ENABLE_PROFILING=false</pre>
\$ export SCOREP_ENABLE_TRACING=true
\$ export SCOREP_METRIC_PAPI=PAPI_TOT_INS, PAPI_TOT_CIC
\$ mpirun -np 4 ./a.out
\$ ls -R
scorep-20170323_1309_7243761919249966 a.out
./scorep-20170323 1309 7243761919249966:
scorep.cfg traces/ traces.def traces.otf2
<pre>\$ vampir scorep-20170323_1309_7243761919249966/traces.otf2</pre>



Data Collection with Score-P



- Score-P supports a combination of instrumentation/sampling:
 - Instrumentation for MPI/OpenMP events
 - Sampling for everything else

Simple configuration, e.g.:

% export SCOREP_ENABLE_TRACING=true

% export SCOREP_ENABLE_UNWINDING=true

% export SCOREP_SAMPLING_EVENTS=perf_cycles@2000000

Score-P: Workflow / Filtering



Use scorep-score to define a filter

- Exclude short frequently called functions from measurement
 - For profiling: reduce measurement overhead (if necessary)
 - For tracing: reduce measurement overhead and total trace size



Filter file:

\$ vim scorep.filt
SCOREP_REGION_NAMES_BEGIN EXCLUDE
matmul_sub
matvec_sub
binvcrhs

Good Tracing Practices



- Traces can become large; Their handling challenging
 - Trace size proportional to: number of processes/threads (width), duration (length), and measurement detail (depth)
- Intermediate flushes => High degree of perturbation
 - Either use less detail or larger trace buffers
- Traces should be written to a parallel file systems
 - E.g. "/work", "/scratch", or "/p/lscratchc"
- Moving large traces can become challenging
 - However, systems with more memory/core can analyze larger traces
 - Alternatively, undersubscribe for increased memory/core ratios



Vampir Performance Data Visualization



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



Event Trace Visualization with Vampir

The value of seeing how an application executes on the machine

- Application code computing coulomb forces
- The workload was distributed evenly across available processes
- The user expected perfect parallelized code
- However the underlying algorithm worked differently than expected

Visualization of the application execution instantly shows a problem in the parallelization approach



Main Performance Charts of Vampir

Timeline Charts



- Master Timeline
- Process Timeline
- Summary Timeline
- Performance Radar
- Counter Data Timeline
- I/O Timeline

- all threads' activities
- *single thread's activities*
- all threads' function call statistics
- all threads' performance metrics
- single threads' performance metrics
- all threads' I/O activities

Summary Charts



- Function Summary Message Summary
- I/O Summary

Process Summary Communication Matrix View Call Tree

İ 🔜

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

% vampir



Visualization Modes (2)

On local machine with remote VampirServer



Visualization: After Tracing













Example Analysis: Visualizing and Analyzing NPB-MZ-MPI / BT



Vampir - [Trac	e View - /home/dolescha/tracefiles/feature	-traces/wrf-p64-io-mem-rus	age/wrf.1h.otf]	×
¥ File ∨iew Help			-	_ 8 ×
⊻iew <u>C</u> hart <u>F</u> ilter				
) 🚟 🔛 💹 🌑 🔁 🛄 🚻		E TELEVILLE MUTATINE AUG	Function Legend	
			Application	
26.860 s 26.865 s	26.870 s 26.875 s 26.880 s 26.8	885 s 26.890 s	DYN	
Process 0	MPI Wait	solve_em_	10	
Process 1	Wait MRI Wait	solve_em_	MEM	
Process 2	MPI_Wait	solve_em_	- MPI DHVS	
Process 3	MPI_Wait MPI_Wait	solve_em_		
Process 4 module_em_mp_rk_s	step_prep	solve_em_	WRF	
Process 5	MPEWait	solve_em_		
Process 6 module_em_mp_rk_s	step_prep_	solve_em_		
Process 7	ait soh	ve_em_		
Process 8 module_em_mp_rk_s	step_prepMRI_Wait	solve_em_	Context View	
Process 9 MPI_Wait		solve_em_	🚟 Master Timeline 🔯	- +
Process 10		solve_em_	Property Value	
Process 11	Wile wait	solve_em_	Display Master Timeline	
Process 12	MPI Wait	solve_em_	Function Name MPI Wait	
Process 13	MPLWait	solve_em_	Function Group MPI	
Process 14 module_em_mp_rk_s	step_prep_	solve_em_	Interval Begin 26.872264 s	
Process 15 module_em_mp_rk_s	step_prep_	solve_em_	Interval End 26.885236 s	
Process 16 module_em_mp_rk_s	step_prep_	solve_em_	Source File	
Process 17 module_em_mp_rk_s	step_prep_	solve_em_	Source Line	
Process 18	MPI_Wait	solve_em_		
Process 19 module_em_mp_rk_s	step_prep_	solve_em_		
Process 20 module_em_mp_rk_s	step_prep_	solve_em_		
Process 21	MPI_Wait MPI_Wait	solve_em_		
Process 22	PI_Wait	solve_em_		
Process 23 module_em_mp_rk_s	step_prepMPI_Wait	solve_em_		
Process 24	MPI_Wait	MPI_Wait		

: 🗮 🐹 🌌	<u>11</u>	- 🔜 🧉) 🔄		1	5 19	1	🄊 💡	34.14	912 340	8 34,14946 s 607 µs				
34 149119	s +50	us +1	00 us	+150 u	e s +2	00 us	+250 us	+300 µ	<u>،</u>	XX	All Processes	Accumulate	Function S ed Exclusive	Summary Time per Fun	ction
04.140110	:	μο · ·	:	100 μ			. 200 μο		-	[2 ms	/ loounnaide	1 ms	0	ms
Process 0 MPI_	Wait			_						×	2.37 ms				MPI_Wait
Process 1	MP	I_Isend		MPI	_Wait		MPI_Wai	t		[1.708 ms			SOLVE_EM
Process 2 MPI_	Wait							SOLVE_EN	1				550.796 μ	s	MPI_lsend
Process 3 SOLV	VE_EM												445.604	μs	MPI_Irecv
Process 4 MPL						_							340.60	07 μs	COUPLEOMENTUM
Process 5 SOLV	VE_EM			PI_Irecv			MPI_IS	end	\leq					25.423 µs	CALCULATE_FULL
Process 6 MPL	_irecv			IVIPI_ISC	na									9.35 µs	RK_STEP_PREP
Process 7 MPL	_wait				9										
Process 8 MPL	Wait				3						i				
Process 10	MPE	Irecy		M	PL Isenc								Context	View	a
Process 11 MPL	Wait				1 1_100110				_		Descentu			1	
Process 12 COU	JPLE MON	MENTUM									Property	value			
Process 13 SOLV	VE EM										Display	Master II	meline		
Process 14 SOLV	VE_EM										Туре	Message			
Process 15 SOLV	VE_EM										Message Typ	e Point to p	oint		
											Origin	Process 1			
											Destination	Process 2			
											Communicato	r MPI Comm	unicator 0		
											Tag	2			
											Start Time	34.14916	1 s		
											Arrival Time	34,14937	9 s		
											Duration	218.05 us			
											Size	100 3350	38 KiB		
											Data Data	002 7402			
											Data Rate	092.7493	VO WID/S		

Visualization of the NPB-MZ-MPI / BT trace



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





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

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





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

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



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





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





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



Context View: Detailed information about function "initialize_".

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



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



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



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



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





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

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

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





I/O Features Operations over Time



I/O Features Data Rates over Time



I/O Features Summaries with Totals





I/O Features Summaries per File



Aggregated data for specific resource

I/O Features Operations per File





Summary and Conclusion



Summary

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



http://www.vampir.eu

vampirsupport@zih.tu-dresden.de