



Multicore Performance Analysis at Scale: From two Embedded Cores to one Million HPC Cores

Bernd Mohr Jülich Supercomputing Centre, Germany



Forschungszentrum Jülich

JÜLICH SUPERCOMPUTING CENTRE

Forschungszentrum Jülich GmbH





- Germany's largest national laboratory
- About 5600 employees
- Research areas
 - Information technology
 - Health (Neuroscience / brain research)
 - Energy and environment

Jülich Supercomputing Centre (JSC)





HPC Centre for

- Forschungszentrum Jülich
- Jülich Aachen Research Alliance (JARA)
- Germany as GCS (1 of 3 German National Centres)
- Europe (1st European Centre inside PRACE)



How and Why

MULTICORE PERFORMANCE ANALYSIS

HPC Architectures: State of the Art





Performance Challenges for HPC Systems



- HPC systems consist of
 - Complex configurations
 - With a huge number of components
 - Very likely heterogeneous
 - With never enough memory
 - Dynamically changing configuration due to fault recovery + power saving
- Deep software hierarchies of large, complex software components are needed to make use of such systems
- Sophisticated integrated performance measurement, analysis, and optimization capabilities are required to efficiently operate an HPC system

Performance Challenges for Multicore Systems **J**Ü

- Multicore Systems | consist of
 - Complex configurations
 - With a huge number of components
 - Very likely heterogeneous
 - With never enough memory



- Dynamically changing configuration due to fault recovery + power saving
- Deep software hierarchies of large, complex software components are needed to make use of such systems
- Sophisticated integrated performance measurement, analysis, and optimization capabilities are required to efficiently operate an Multicore Systems

Measurement Methods: Profiling



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

Measurement Methods: Tracing



- Recording information about significant points (events) during execution of the program
 - Enter/leave a code region (function, loop, ...)
 - Send/receive a message ...
- Save information in event record
 - Timestamp, location ID, event type
 - plus event specific information
- **Event trace** := stream of event records sorted by time
- Can be used to reconstruct the dynamic behavior
 ⇒ Abstract execution model on level of defined events



Event Tracing: "Timeline" Visualization







More than numbers and diagrams

INSIGHTFULNESS

Interactive Event Trace Analysis: Vampir





Visual presentation of dynamic runtime behaviour

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



http://www.vampir.eu/

Vampir GUI (zoom)





Interactive browsing, zooming, selecting

 Linked displays & statistics adapt to selected time interval

Trace formats

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

"A picture is worth 1000 words..."



10	10.2 mc 10.2	5 - 19.693 III5 = 0 me 10	1.927 IIIS) Ame	10 G me		Proces
ocess 0 683 MPI_Recv	5 15.21	115 15	.4 1113	13.0 113	MPI	Proces Proces Proces
ocess 1 6 80 83 MPI_Finali	ze				Application	Proces Proces Proces
cess 2 6 80 83 MPI_	Finalize					Proces Proces Proces
cess 3 6 MPI_Recv 83	MPI_Finalize					Proces Proces Proces
cess 4 6 MPI_Recv	83 MPI_Finalize					Proces Proces Proces
cess 5 6 MPI_Recv	83 MPI_Final	lize				Proces Proces Proces
cess 6 6 MPI_Recv	83 MPI	Finalize				Proces Proces Proces Proces
cess 7 6 MPI_Recv	83	MPI_Finalize				Proces Proces Proces Proces
cess 8 6 MPI_Recv		83 MPI_Final	ize			Proces Proces Proces Proces
cess 9 6 MPI_Recv		83 MPI_	Finalize			Proces Proces Proces
cess 10 <mark>6 MPI_Recv</mark>		83	MPI_Finalize			Proces Proces Proces
cess 11 6 MPI_Recv			83 MPI_Fin	alize		Proces Proces Proces
cess 12 6 MPI_Recv			83 MP	I_Finalize		Proces Proces Proces
cess 13 6 MPI_Recv			8:	B MPI_Fina	lize	Proces Proces Proces
cess 14 <mark>6 MPI_Recv</mark>				83 54		Proces Proces Proces
cess 15 6 MPI_Recv				8:	3	Proces Proces Proces Proces



• MPI ring program

"Real world" example

"What about 1000's of pictures?" (with 100's of menu options)







Example Automatic Analysis: Late Sender





(a) Late Sender

process



ENTER

(c) Late Sender / Wrong Order

Example MPI Wait States





Scalasca

- Scalable Analysis of Large Scale Applications
- Approach





http://www.scalasca.org/

- Instrument C, C++, and Fortran parallel applications
- Option 1: <u>scalable</u> call-path profiling
- Option 2: <u>scalable</u> 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



Late Sender Analysis

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







Late Sender Analysis + Application Topology

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



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





Direct Wait Time Analysis

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



JSC



Indirect Wait Time Analysis

 Indirect waits occurs for ranks processing warmer areas







Delay Costs Analysis

 Delays NOT caused on ranks processing ice!







YOU KNOW YOU MADE IT ...

... IF 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:

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

Copyright © 2014, Intel Corporation. All rights reserved. *Other names and brands may be claimed as the property of others.

Optimization Notice

Which Performance Issues are automatically identified?

Point-to-point exchange problems:

Late Sender



wait time

send

Late Receiver

Late Receiver

P1

Problems with global collective operation performance:

Wait at Barrier



Early Reduce



Σ

Late Broadcast





Source:

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

Copyright © 2014, Intel Corporation. All rights reserved. *Other names and brands may be claimed as the property of others.

receive

time

Optimization Notice

(intel)

14



Together we are strong

INTEGRATION

Integration



- Need integrated tool (environment) for all levels of parallelization
 - Inter-node (MPI, PGAS, SHMEM)
 - Intra-node (OpenMP, multi-threading, multi-tasking)
 - Accelerators (OpenACC, CUDA, OpenCL)
- Integration with performance modeling and prediction
- No tool fits all requirements
 - Interoperability of tools
 - Integration via open interfaces



Scalable performance measurement infrastructure for parallel codes









Technische Universität München



UNIVERSITY OF OREGON

Score-P Tool Ecosystem





February 2016

JSC

Connect to Vampir Trace Browser





Show most Severe Pattern Instances



cube 4.1.1 livedvd2: scorep_bt-mz_B_4x4_trace/trace.cubex							
<u>F</u> ile <u>D</u> isplay <u>T</u> opology <u>H</u> elp							
Absolute	Absolute ~ Absolute	~					
Metric tree	💽 Call tree 🔲 Flat view 💽 System	tree 🚺 Box Plot					
 0.00 Time 300.91 Execution 0.00 MPI 0.01 Synchronization 0.00 Communication 0.39 Point-to-point 1.38 Late Sender 0.00 Collective 0.00 Collective 0.00 Early Reduce 0.00 Early Scan 0.00 Kat Broadcast 0.00 Na N Completion 0.00 File I/O 0.87 Init/Exit 0.00 OMP 0.00 Filush 2.17 Management 0.00 Synchronization 22.99 Barrier 	O.00 MAIN_ Call site O.00 mpi_set Call site O.00 mPi_Set Call site O.00 env_set Called region O.00 zone_set Expand/collapse O.00 zone_sta Cut call tree O.00 set_cont Cut call tree O.00 set_cont Find items O.00 copy Clear found items O.00 copy O.00 copy O.00 copy O.00 copy O.00 copy O.00 onl O.00 onl Min/max values O.00 onl Min/max values O.00 onl Min/max values O.00 onl Min/max values O.00 onl Select "Max severity in trace	neric cluster * i06r01c20 - - MPI Rank 0 - 0.34 Thread 0 - 0.00 Thread 1 - 0.00 Thread 3 - MPI Rank 1 - 0.00 Thread 3 - MPI Rank 1 - 0.00 Thread 1 - 0.00 Thread 1 - 0.00 Thread 3 - MPI Rank 2 - 0.00 Thread 1 - 0.00 Thread 1 - 0.00 Thread 3 - MPI Rank 3 - 0.21 Thread 0 - 0.00 Thread 1 - 0.00 Thread 3					
0.00 1.38 (0.41%) 337.45	^{0.00} browser" from context mer	1.38					
	of call paths marked with	a					
Shows the most severe instance of pattern in trace brows							
	red frame						

Investigate most Severe Instance in Vampir







To infinity and beyond

EXTREME CONCURRENCY

Typical HPC System Size





Personal Motivation



JuQueen

IBM BlueGene/Q 28 racks 458,752 cores 1,835,008 HW threads

2012/2013: Most powerful machine in Europe!

Most **parallel** machine in Europe!



VampirServer: Trace Visualization S3D@200,448



 OTF trace 4.5 TB

 Vampir Server running with 20,000 analysis processes



JSC

Scalasca Trace Analysis sweep3D@294,912 BGP

JSC

- 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.



February 2016

ÜLICH

Performance Tool Scaling: Scalasca



- Latest full application test case
 - Granular Dynamics Simulation
 - Based on Physics Engine (PE) Framework (Erlangen), MPI only
 - PRACE @ ISC Award winner
- Scalasca 1.x experiments on JUQUEEN
 - Full machine experiment: 28,672 nodes x 32 MPI ranks
 - 917,504 processes [Limit: Memory / System metadata]
 - Largest number of threads: 20,480 nodes x 64 MPI ranks
 - 1,310,720 processes
 [Limit: Memory / System metadata]
- Scalasca 2.x / Score-P 1.4.1 Nekbone (CORAL benchmark) on JUQUEEN
 - Profiles: 28,672 x 64 = 1,835,008 threads !!!
 - Traces: 10,240 x 64 = 655,360 threads

[Limit: OTF2]

Scalasca: 1,835,008 Threads Test Case



• Nekbone

CORAL
 benchmark

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







USE CASES

Success Story: TerrSysMP

- Scale-consistent highly modular integrated multi-physics sub-surface/surface hydrology-vegetation atmosphere modelling system
- fully-coupled MPMD simulation consisting of
 - COSMO (Weather prediction
 - CLM (Community Land Model
 - ParFlow (Parallel Watershed Flow)
 - OASIS coupler

JSC





Success Story: TerrSysMP



- Identified several sub-components bottlenecks:
 - Inefficient communication patterns
 - Unnecessary/inefficient code blocks
 - Inefficient data structures
- Performance of sub-components improved by factor of 2!
- Scaling improved from 512 to 32768 cores!







Optimize Industrial HPC Applications on Heterogeneous Architectures

H4H Project Review#3, Repsol, Madrid, Spain 19th September 2013



Company	Area	Performance Analysis	Programming Tools	High level tools
ATEME	Video compression	ThreadSpotter™ Scalasca, VampirTrace	HMPP	
Dassault	Simulation of aircraft	ThreadSpotter™	HMPP	SAMG
Aviation	design	Scalasca, Vampir		Scilab
		PAS2P		LAMA
Efield	Electromagnetic	ThreadSpotter™	HMPP	SAMG
	fields modeling and	Scalasca, Vampir		CuBLAS
	simulation	Marmot		LAMA
GNS	Metal forming	Scalasca/Score-P	HMPP	SAMG
	processes simulation	VampirTrace		
INTES	Implicit finite element	VampirTrace		
	analysis system			
MAGMA	Casting process	ThreadSpotter™		SAMG, LAMA
	simulation	Scalasca		
RECOM	3D combustion	VampirTrace	OpenACC	
Repsol	Seismic imaging and	PAS2P		Scilab
	reservoir simulation			



- Significant performance improvements for Applications through performance analysis, code restructuring / porting on GPU
- DASSAULT:
 - LAMA ⇒ sparse linear equations solver GPU-MPI **↑3x to 4x faster**
 - Scilab/GPGPU ⇒ processing electromagnetic measurements **↑3x faster**
 - Scilab/MPI ⇒ very good scalability reached in inverse design code
- RECOM:
 - Optimized particle deposition algorithm : **↑7.9x** GPU vs CPU
- MAGMA:
 - Sparse linear equations solver on GPU with SAMG 12x to 2.65x faster
 - Sparse linear equations solver on GPU with LAMA CG solver **^3x** faster
 - MPI-GPU solver kernel version with CUDA yielded a 3x speedup (8 CPU MPI vs 8 CPU / 1 GPU version)

||||||4||||

- GWT:
 - IUMD code → 13x speedup on GPU compared 16 OpenMP Threads Intel Xeon E5 CPU
- INTES:
 - XPU (SMP+GPU) **12.5x faster than pre-existing** 16-core SMP version
 - Combined with additional DMP parallelization, a static analysis of a motor block benchmark even showed a speedup of 68.6 against the single core time (total job) → Turnaround time for typical simulation cycles reduced from 1 week to about 2 hours (SMP+XPU+DMP).
- REPSOL:
 - **1.9x** Speedup for seismic code with GPUs and linear scaling with MPI
 - **12.6x** Speedup for Reservoir Simulator with GPUs speed up of reservoir





Runtime Analysis of Parallel applications for Industrial software Development

- Collaboration between
 - Corporate Technology Multicore Expert Center of Siemens AG
 - Jülich Supercomputing Centre
- Results
 - Tool support (Score-P, Scalasca, Vampir) for
 - Threading models (POSIX, QT, ACE, Windows)
 - MTAPI (Multicore Association Task API)
 - Windows port of basic Score-P and Cube



Corporate Technology | Siemens AG

ScoreP, Scalasca and Vampir @ Siemens AG



Unrestricted © Siemens AG 2015. All rights reserved

Unrestricted © Siemens AG 2015. All rights reserved

Targeted Analysis using adapted ScoreP allows relevant focussing

Via ScoreP, Cube, Vampir

- Effective (sub-)hotspot localization (80%)
- Identified legacy threads
- Gained insight into call hierarchy/stack





SIEMENS





Questions?





scalasca 🗖

- http://www.scalasca.org
- scalasca@fz-juelich.de





- http://www.score-p.org
- support@score-p.org

