# Methods for performance evaluation and optimization on modern HPC systems

Felix Wolf

08-06-2011

# Objectives

- Learn about basic performance measurement and analysis methods and techniques for HPC applications
- Get to know Scalasca, a scalable and portable performance analysis tool

# Performance tuning: an old problem

"The most constant difficulty in contriving the engine has arisen from the desire to reduce the time in which the calculations were executed to the shortest which is possible."
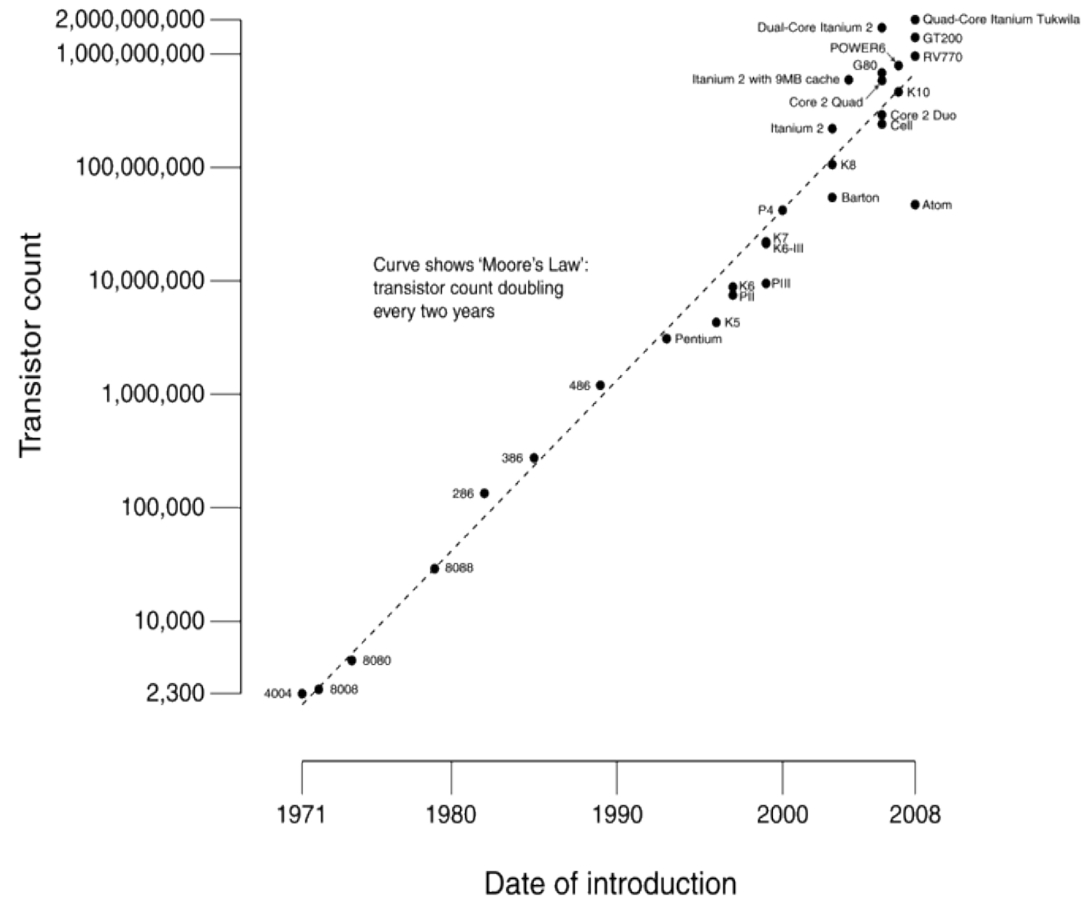
Charles Babbage
1791 - 1871

# Outline

- Principles of parallel performance
- Performance analysis techniques
- Practical performance analysis using Scalasca

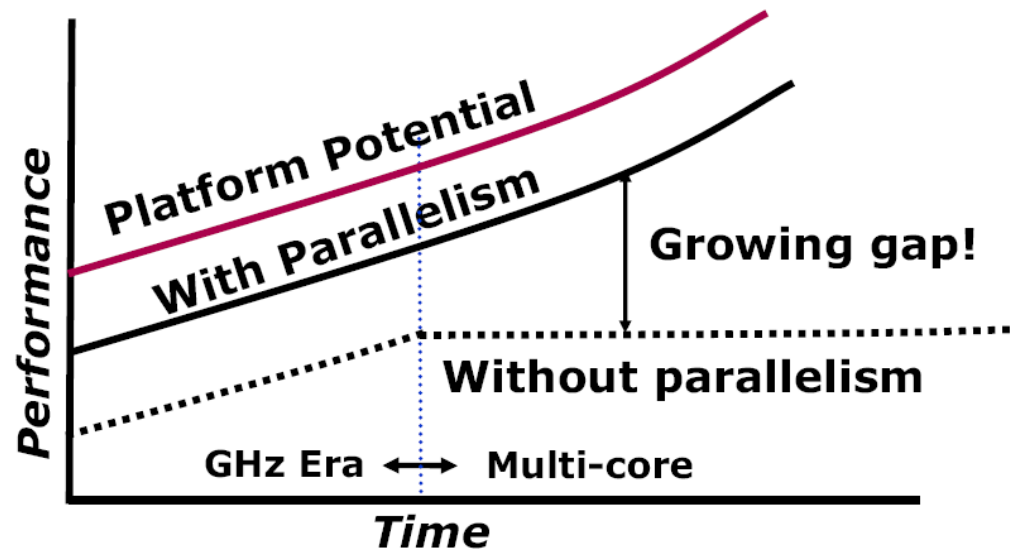# Why parallelism at all?
# Moore's Law is still in charge…



Source: Wikipedia

# Free lunch is over…

## Need for Parallelism



Parallelism is crucial for optimal performance

# Parallelism

- System/application level
  - Server throughput can be improved by spreading workload across multiple processors or disks
  - Ability to add memory, processors, and disks is called scalability

- Individual processor
  - Pipelining
  - Depends on the fact that many instructions do not depend on the results of their immediate predecessors

- Detailed digital design
  - Set-associative caches use multiple banks of memory
  - Carry-lookahead in modern ALUs

# Amdahl's Law for parallelism

- Assumption – program can be parallelized on p processors except for a sequential fraction f with

$$0 \leq f \leq 1$$

$$Speedup(p) = \frac{1}{f + \dfrac{1-f}{p}} < \frac{1}{f}$$

- Speedup limited by sequential fraction

# Available parallelism

- Overall speedup of 80 on 100 processors

$$80 = \frac{1}{f + \dfrac{1-f}{p}}$$

f $= 0.0025$

# Law of Gustafson

- Amdahl's Law ignores increasing problem size
  - Parallelism often applied to calculate bigger problems instead of calculating a given problem faster
- Fraction of sequential part may be function of problem size
- Assumption
  - Sequential part has constant runtime $\tau_f$
  - Parallel part has runtime $\tau_v(n,p)$
- Speedup

$$Speedup(n,p) = \frac{\tau_f + \tau_v(n,1)}{\tau_f + \tau_v(n,p)}$$

If parallel part can be perfectly parallelized

# Parallel efficiency

$$\text{Efficiency(p)} = \frac{\text{Speedup(p)}}{p}$$

- Metric for cost of parallelization (e.g., communication)
- Without super-linear speedup

$$\text{Efficiency(p)} \leq 1$$

- Super-linear speedup possible
  - Critical data structures may fit into the aggregate cache

# Scalability

- **Weak** scaling
  - Ability to solve a larger input problem by using more resources (here: processors)
  - Example: larger domain, more particles, higher resolution
- **Strong** scaling
  - Ability to solve the same input problem faster as more resources are used
  - Usually more challenging
  - Limited by Amdahl's Law and communication demand

# Serial vs. parallel performance

- Serial programs
  - Cache behavior and ILP

- Parallel programs
  - Amount of parallelism
  - Granularity of parallel tasks
  - Frequency and nature of inter-task communication
  - Frequency and nature of synchronization
    - Number of tasks that synchronize much higher $\rightarrow$ contention

# Goals of performance analysis

- Compare alternatives
  - Which configurations are best under which conditions?
- Determine the impact of a feature
  - Before-and-after comparison
- System tuning
  - Find parameters that produce best overall performance
- Identify relative performance
  - Which program / algorithm is faster?
- Performance debugging
  - Search for bottlenecks
- Set expectations
  - Provide information for users

# Analysis techniques (1)

- Analytical modeling
  - Mathematical description of the system
  - Quick change of parameters
  - Often requires restrictive assumptions rarely met in practice
    - Low accuracy
  - Rapid solution
  - Key insights
    - Validation of simulations / measurements
- Example
  - Memory delay $\quad t_{avg} = ht_c + (1-h)t_m$

  - Parameters obtained from manufacturer or measurement

# Analysis techniques (2)

- Simulation
  - Program written to model important features of the system being analyzed
  - Can be easily modified to study the impact of changes
  - Cost
    - Writing the program
    - Running the program
  - Impossible to model every small detail
    - Simulation refers to "ideal" system
    - Sometimes low accuracy

- Example
  - Cache simulator
  - Parameters: size, block size, associativity, relative cache and memory delays

# Analysis techniques (3)

- Measurement
  - No simplifying assumptions
  - Highest credibility
  - Information only on specific system being measured
  - Harder to change system parameters in a real system
  - Difficult and time consuming
  - Need for software tools

- Should be used in conjunction with modeling
  - Can aid the development of performance models
  - Performance models set expectations against which measurements can be compared

# Metrics of performance

- ## What can be measured?
  - A count of how many times an event occurs
    - E.g., Number of input / output requests
  - The duration of some time interval
    - E.g., duration of these requests
  - The size of some parameter
    - Number of bytes transmitted or stored

- ## Derived metrics
  - E.g., rates / throughput
  - Needed for normalization

# Primary performance metrics

- Execution time, response time
  - Time between start and completion of a program or event
  - Only consistent and reliable measure of performance
  - Wall-clock time vs. CPU time

- Throughput
  - Total amount of work done in a given time

- Performance = $\dfrac{1}{\text{Execution time}}$

- Basic principle: reproducibility

- Problem: execution time is slightly non-deterministic
  - Use mean or minimum of several runs

# Alternative performance metrics

- Clock rate

- Instructions executed
  per second

- FLOPS

  – Floating-point operations per second

- Benchmarks

  – Standard test program(s)

  – Standardized methodology

  – E.g., SPEC, Linpack

- QUIPS / HINT [Gustafson and Snell, 95]

  – Quality improvements per second

  – Quality of solution instead of effort to reach it

"Math" operations?
HW operations?
HW instructions?
Single or double
precision?

# Comparison of analysis techniques

|  | Analytical modeling | Simulation | Measurement |
|---|---|---|---|
| Flexibility | High | High | Low |
| Cost | Low | Medium | High |
| Credibility | Low | Medium | High |
| Accuracy | Low | Medium | High |

# Peak performance

- Peak performance is the performance a computer is guaranteed not to exceed



64 processors

Source: Hennessy, Patterson: Computer Architecture, 4th edition, Morgan Kaufmann

# Performance tuning cycle

# Performance measurement cycle (2)

- Instrumentation
  - Insertion of extra code (probes) into application
- Measurement
  - Collection of data relevant to performance analysis
- Analysis
  - Calculation of metrics
  - Identification of performance bottlenecks
- Presentation
  - Transformation of the results into a representation that can be easily understood by a human user
- Optimization
  - Elimination of bottlenecks

# Semantic gap



- Programmer's mental model of the program does not match the executed version
  - Performance tools needed to bridge this semantic gap

# Semantic performance mapping

- Instrumentation levels
  - Source code
  - Library
  - Runtime system
  - Object code
  - Operating system
  - Runtime image
  - Virtual machine

- Problem
  - Every level provides different information
  - Often instrumentation on multiple levels required

- Challenge
  - Mapping performance data onto application-level abstraction

# Instrumentation techniques

- Static instrumentation
  - Program is instrumented prior to execution
- Dynamic instrumentation
  - Program is instrumented at runtime
- Code is inserted
  - Manually
  - Automatically
    - By preprocessor
    - By compiler
    - By linking against preinstrumented (interposition) library
    - By binary-rewrite / dynamic instrumentation tool

# Measurement

Typical performance data include
- Counts
- Durations

inclusive duration

exclusive duration

```
int foo()
{
  int a;

  a = a + 1;

  bar();

  a = a + 1;
}
```

- Communication cost
- Synchronization cost
- IO accesses
- System calls
- Hardware events

# Critical issues

- Accuracy
  - Perturbation
    - Measurement alters program behavior
    - E.g., memory access pattern
  - Intrusion overhead
    - Measurement itself needs time and thus lowers performance
  - Accuracy of timers, counters
- Granularity
  - How many measurements
    - Pitfall: short but frequently executed functions
  - How much information / work during each measurement
- Tradeoff
  - Accuracy ⇔ expressiveness of data

# Single-node performance

- Huge gap between CPU and memory speed



Source: Hennessy, Patterson: Computer Architecture, 4th edition, Morgan Kaufmann

- Internal operation of a microprocessor potentially complex
  - Pipelining
  - Out-of-order instruction issuing
  - Branch prediction
  - Non-blocking caches

# Hardware counters

- Small set of registers that count events
- Events are signals related to the processor's internal function
- Original purpose: design verification and performance debugging for microprocessors
- Idea: use this information to analyze the performance behavior of an application as opposed to a CPU

# Typical hardware counters

| Cycle count | |
|---|---|
| Instruction count | All instructions<br>Floating point<br>Integer<br>Load / store |
| Branches | Taken / not taken<br>Mispredictions |
| Pipeline stalls due to | Memory subsystem<br>Resource conflicts |
| Cache | I/D cache misses for different levels<br>Invalidations |
| TLB | Misses<br>Invalidations |

# Profiling

- Mapping of aggregated information
  - Time
  - Counts
    - Calls
    - Hardware counters
- Onto program and system entities
  - Functions, loops, call paths
  - Processes, threads

- Methods to create a profile
  - PC sampling (statistical approach)
  - Interval timer / direct measurement (deterministic approach)

# Profiling (2)

- Sampling
  - Statistical measurement technique
    - Based on the assumption that a subset of a population being examined is representative for the whole population
    - Requires long-running programs
  - Periodic operating system signal interrupts the running program
  - Interrupt service routine examines return-address stack to find address of instruction being executed when interrupt occurred
  - Using symbol-table information this address is mapped onto specific subroutine

- Interval timing
  - Time measurement at beginning and end of a code region
  - Requires high-resolution / low-overhead clock

# Call-path profiling

- Behavior of a function may depend on caller (i.e., parameters)
- Flat function profile often not sufficient
- How to determine call path at runtime?
  - Runtime stack walk
  - Maintain shadow stack
    - Requires tracking of function calls

```
main()
{
   A( );
   B( );
}


A( )      B( )
{         {
   X();      Y();
   Y();   }
}
```

# Event tracing



- Typical events
  - Entering and leaving a function
  - Sending and receiving a message

# Why tracing?

- High level of detail
- Allows in-depth post-mortem analysis of program behavior
  - Time-line visualization
  - Automatic pattern search
- Identification of wait states

# Obstacle: trace size



- Problem: width and length of event trace

# Tracing vs. profiling

- Advantages of tracing
  - Event traces preserve the temporal and spatial relationships among individual events
  - Allows reconstruction of dynamic behavior of application on any required abstraction level
  - Most general measurement technique
    - Profile data can be constructed from event traces
- Disadvantages
  - Traces can become very large
  - Writing events to a file at runtime can cause perturbation
  - Writing tracing software is complicated
    - Event buffering, clock synchronization, …

**scalasca**

- Scalable performance-analysis toolset for parallel codes
  - Focus on communication & synchronization
- Integrated performance analysis process
  - Performance overview on call-path level via call-path profiling
  - In-depth study of application behavior via event tracing
- Supported programming models
  - MPI-1, MPI-2 one-sided communication
  - OpenMP (basic features)
- Available for all major HPC platforms

# Joint project of



JÜLICH
FORSCHUNGSZENTRUM

German Research School
for Simulation Sciences

# The team

# www.scalasca.org

# Installations and users

- Companies
    - Bull (France)
    - Dassault Aviation (France)
    - Efield Solutions (Sweden)
    - GNS (Germany)
    - INTES (Germany)
    - MAGMA (Germany)
    - RECOM (Germany)
    - SciLab (France)
    - Shell (Netherlands)
    - Sun Microsystems (USA, Singapore, India)
    - Qontix (UK)
- Research/supercomputing centers
    - ANL (USA)BSC (Spain)
    - CASPUR (Italy)
    - CEA (France)
    - CERFACS (France)
    - CINECA (Italy)
    - CSC (Finland)
    - CSCS (Switzerland)
    - DLR (Germany)
    - DKRZ (Germany)
    - EPCC (UK)
    - FZJ (Germany)
    - HLRN (Germany)
    - HLRS (Germany)
    - ICHEC (Ireland)
    - IDRIS (France)
    - KIT (Germany)
    - LLNL (USA)

- Research/supercomputing centers (cont.)
    - LRZ (Germany)
    - MCH (Switzerland)
    - NCAR (USA)
    - NCSA (USA)
    - ORNL (USA)
    - PIK (Germany)
    - PSC (USA)
    - RZG (Germany)
    - SARA (Netherlands)
    - SAITC (Bulgaria)
    - TACC (USA)
- Universities
    - Lund University (Sweden)
    - MSU (Russia)
    - RPI (USA)
    - RWTH (Germany)
    - TUD (Germany)
    - UOregon (USA)
    - UTK (USA)
- DoD/MoD computing centers
    - AFRL DSRC (USA)
    - ARL DSRC (USA)
    - ARSC DSRC (USA)
    - AWE (UK)
    - ERDC DSRC (USA)
    - Navy DSRC (USA)
    - MHPCC DSRC (USA)
    - SSC-Pacific (USA)
    - MetOffice (UK)

# Wait-state analysis

- Classification
- Quantification



(a) Late Sender

(b) Late Sender / Wrong Order

(c) Late Receiver
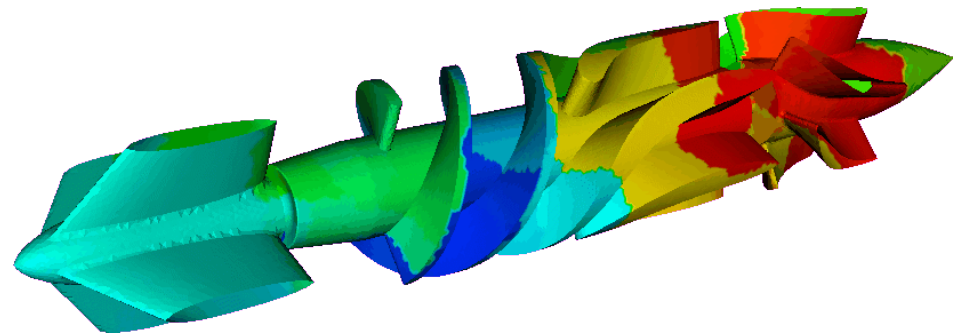
# XNS CFD simulation application

- Computational fluid dynamics code
  - Developed by Chair for Computational Analysis of Technical Systems, RWTH Aachen University
  - Finite-element method on unstructured 3D meshes
  - Parallel implementation based on message passing
  - >40,000 lines of Fortran & C
  - DeBakey blood pump test case
    - Scalability of original version limited <1024 CPUs

Partitioned finite-element mesh

# Call-path profile: Computation

# Call-path profile: P2P messaging

# Call-path profile: P2P sync. ops.
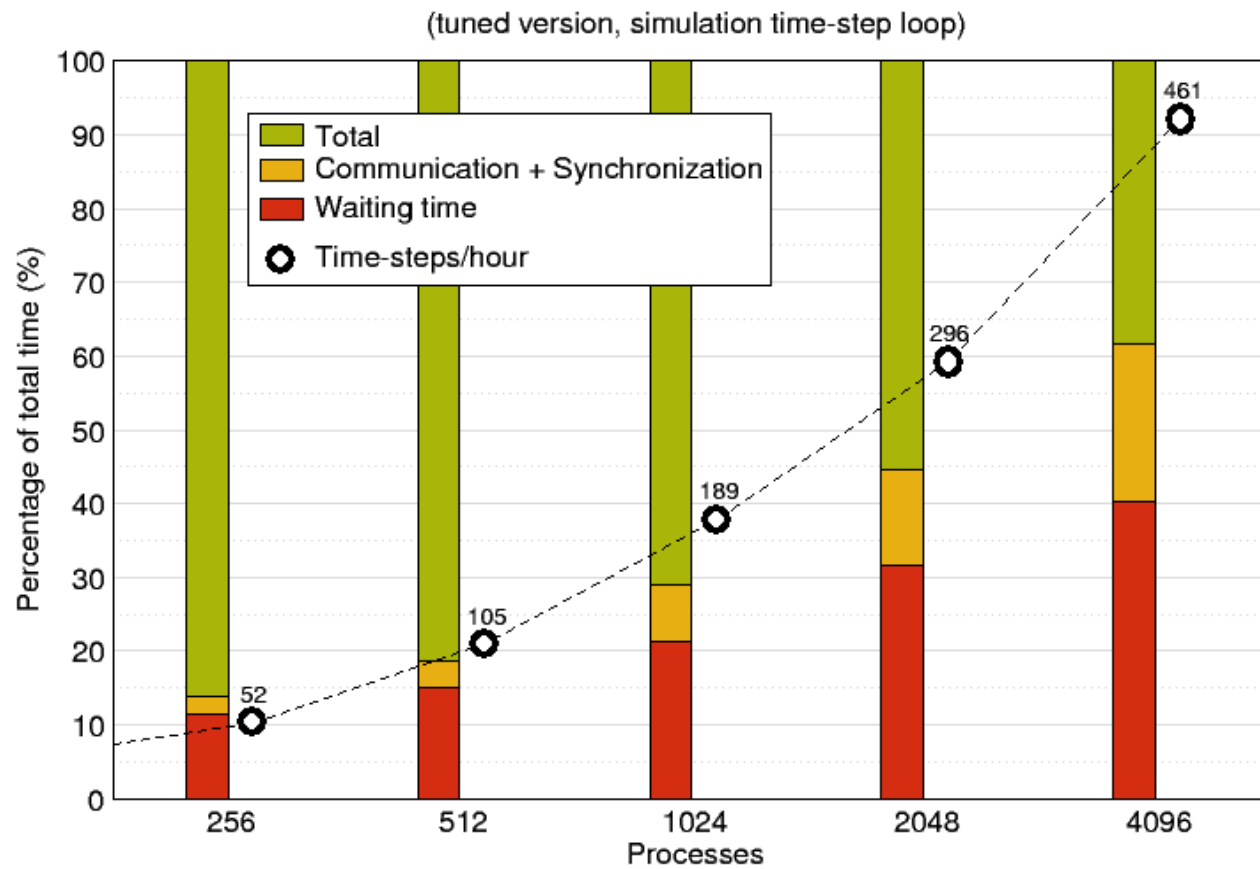
# Trace analysis: Late sender

# XNS scalability remediation

- Review of original XNS
  - Computation is well balanced
  - Real communication is very imbalanced
  - Huge amounts of P2P synchronisations
    - Grow exponentially with number of processes

- Elimination of redundant messages
  - Relevant neighbor partitions known in advance from static mesh partitioning
  - Most transfers still required at small scale while connectivity is relatively dense
  - Growing benefits at larger scales (>512)
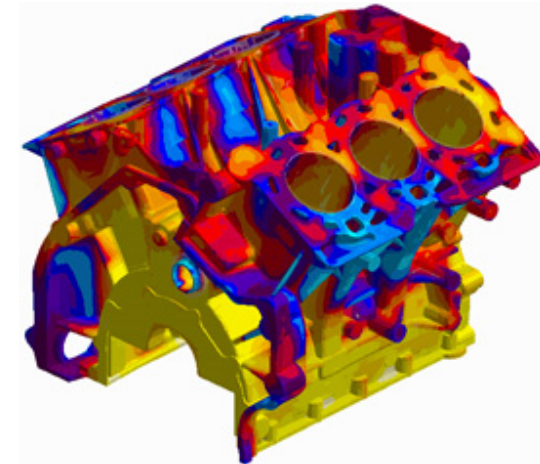
# After removal of redundant messages

# XNS wait-state analysis of tuned version



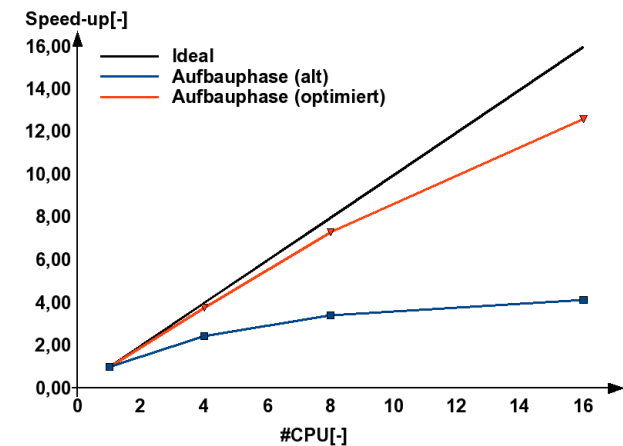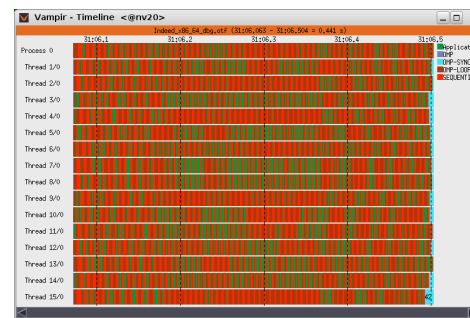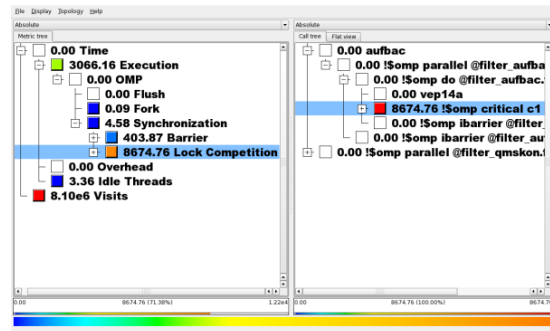(tuned version, simulation time-step loop)

# MAGMAfill by MAGMASOFT® GmbH

- Simulates mold-filling in casting processes

- Scalasca used
  - To identify communication bottleneck
  - To compare alternatives using performance algebra utility

- 23% overall runtime improvement

# INDEED by GNS® mbh

- Finite-element code for the simulation of material-forming processes
  - Focus on creation of element-stiffness matrix
- Tool workflow
  - Scalasca identified serialization in critical section as bottleneck
  - In-depth analysis using Vampir
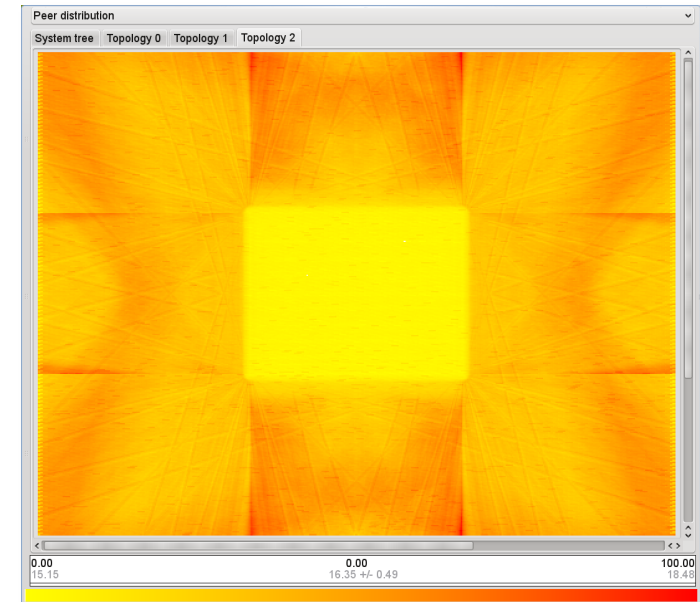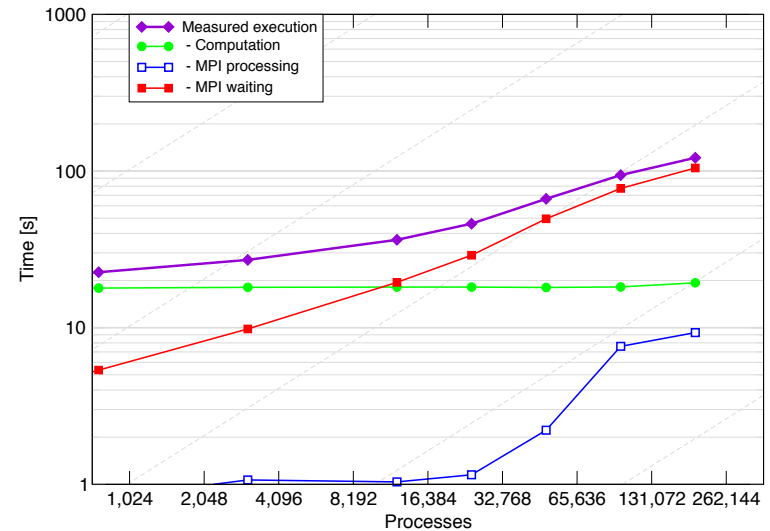- Speedup of 30-40% after optimization

# Scalability in terms of the number of cores

- Application study of ASCI Sweep3D benchmark
- Identified MPI waiting time correlating with computational imbalance
- Measurements & analyses demonstrated on
  - Jaguar with up to 192k cores
  - Jugene with up to 288k cores

Brian J.N. Wylie et al.: Large-scale performance analysis of Sweep3D with the Scalasca toolset. Parallel Processing Letters, 20(4):397-414, December 2010.
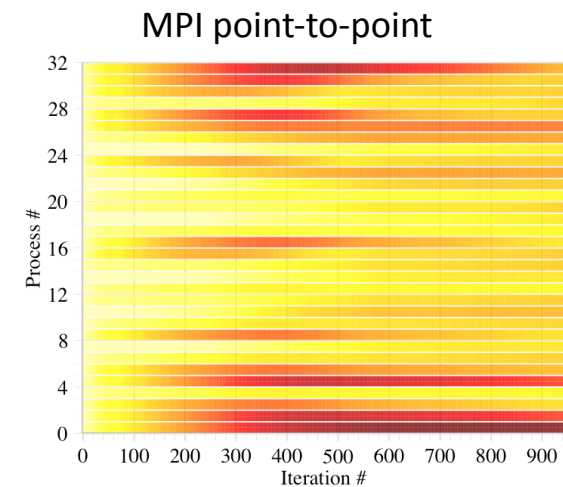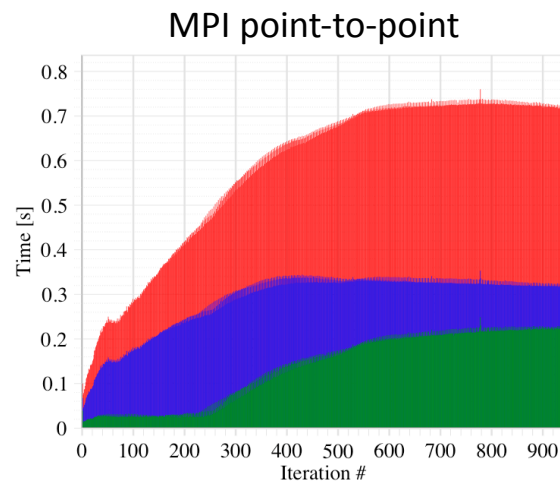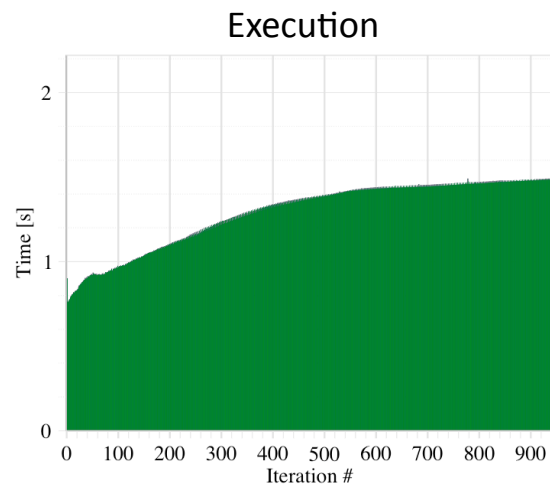


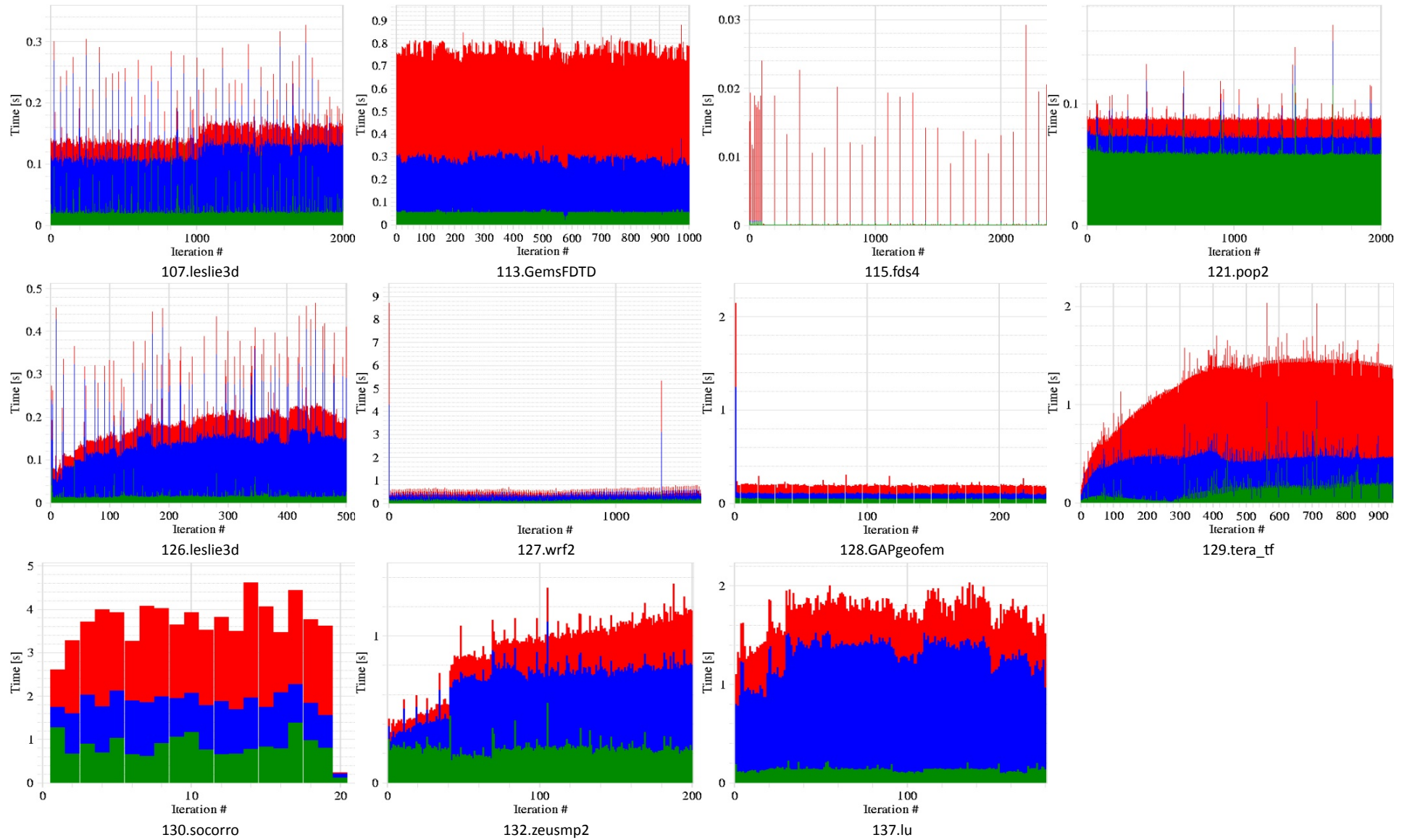Jaguar, MK = 10 (default)



Computation

# Performance dynamics

- Most simulation codes work iteratively
- Growing complexity of codes makes performance behavior more dynamic – even in the absence of failures
  - Periodic extra activities
  - Adaptation to changing state of computation
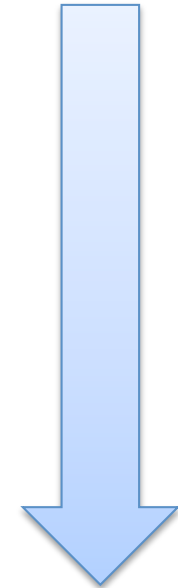- External influence (e.g., dynamic reconfiguration)



129.tera_tf

# P2P communication in SPEC MPI 2007 suite



107.leslie3d  113.GemsFDTD  115.fds4  121.pop2

126.leslie3d  127.wrf2  128.GAPgeofem  129.tera_tf

130.socorro  132.zeusmp2  137.lu

# Scalasca's approach to performance dynamics

- Capture overview of performance dynamics via time-series profiling
  - Time and count-based metrics
- Identify pivotal iterations
  - If reproducible
- In-depth analysis of these iterations via tracing
  - Analysis of wait-state formation including root cause analysis    **New**
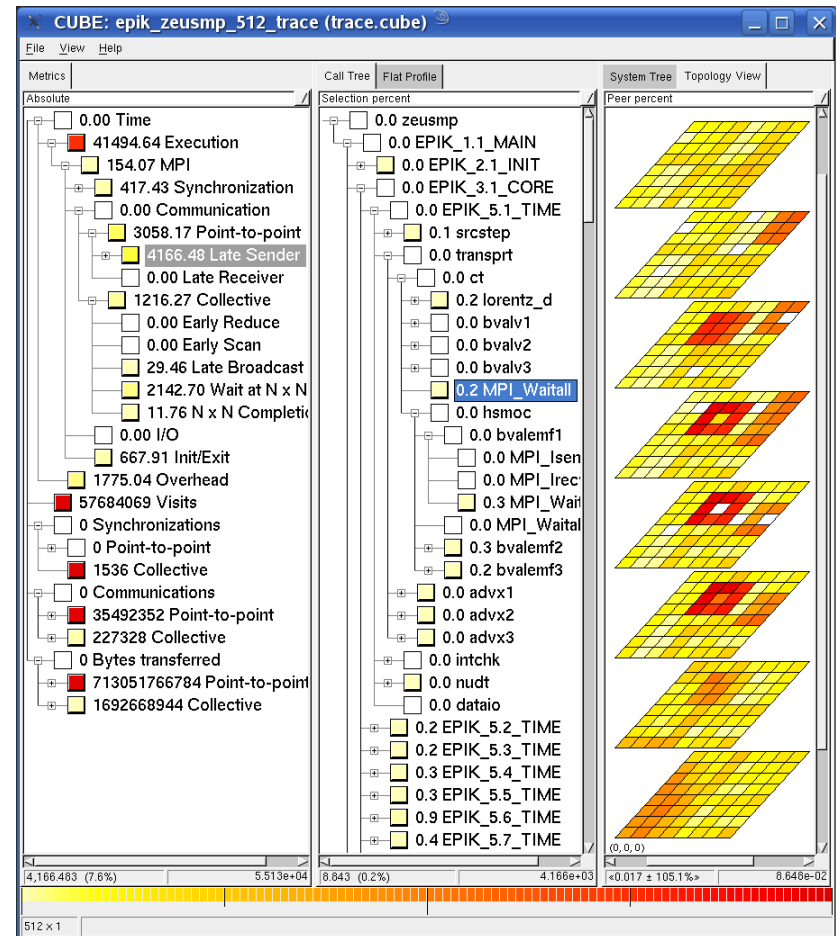  - Tracing restricted to iterations of interest

# Scalable time-series call-path profiling

- Instrumentation of the main loop to distinguish individual iterations

- Complete call-tree recorded for each iteration
  - With multiple metrics collected for every call-path

- Low-overhead online compression of iteration profiles
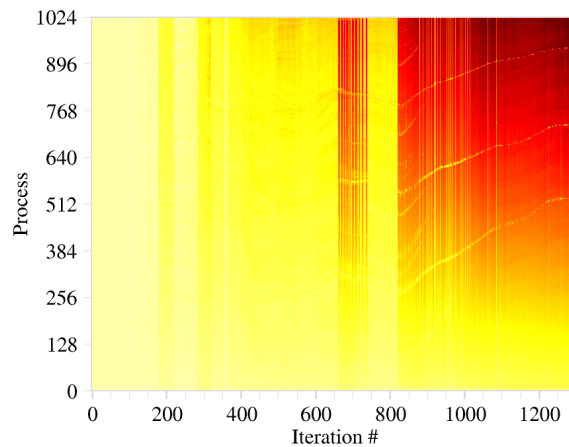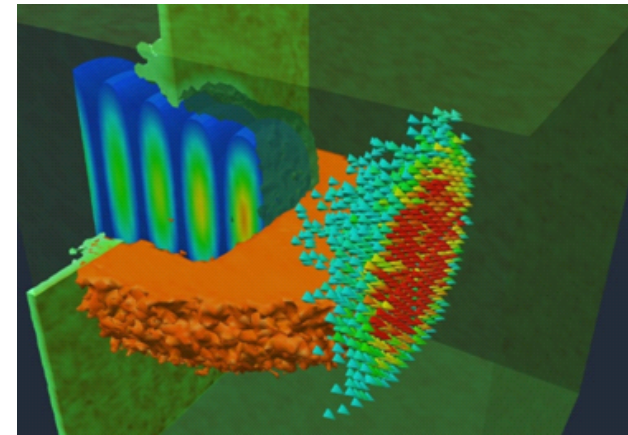  - Reduces memory requirements

Zoltán Szebenyi et al.: Space-Efficient Time-Series Call-Path Profiling of Parallel Applications. In Proc. of the SC09 Conference, Portland, Oregon, ACM, November 2009.
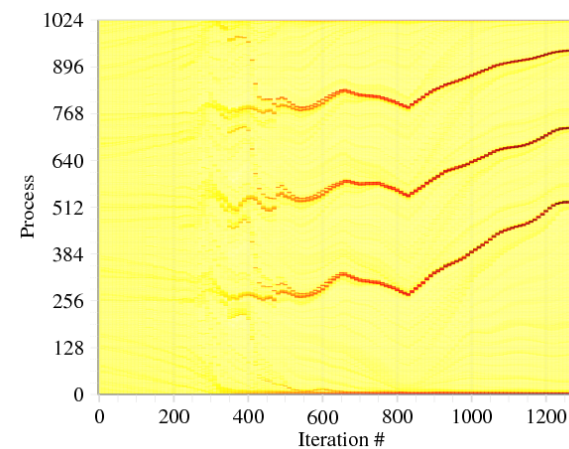
# Pretty Efficient Parallel Coulomb-solver (PEPC)

- Multi-purpose parallel tree code
  - Molecular dynamics
  - Laser-plasma interactions
- Developed at JSC





Late Sender



# particles owned by a process

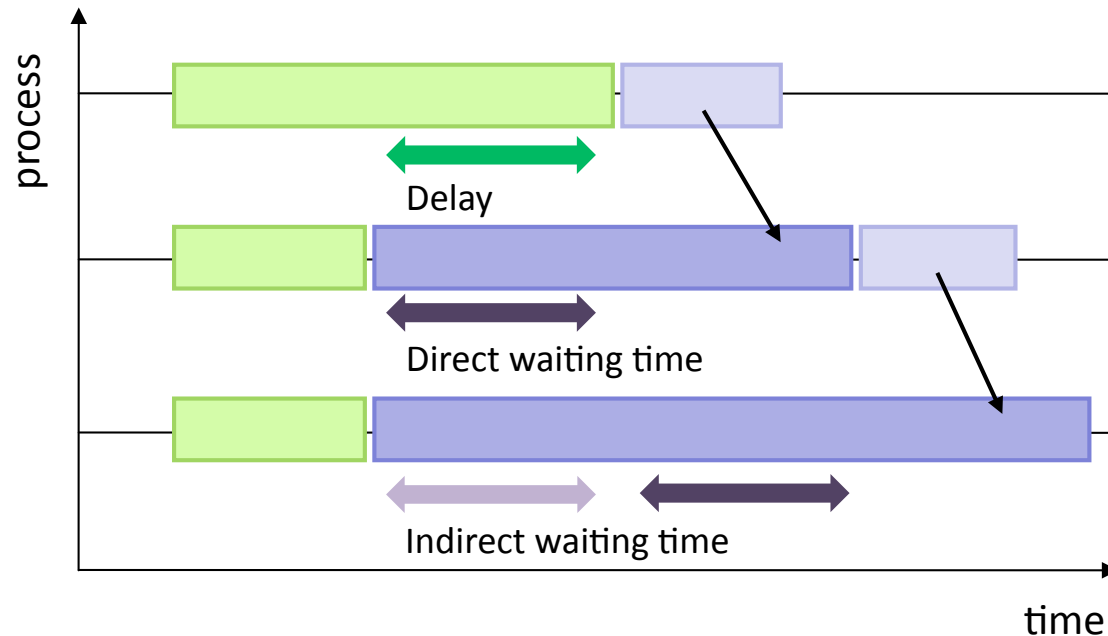# Reconciling sampling and direct instrumentation

- Semantic compression needs direct instrumentation to capture communication metrics and to track the call path
- Direct instrumentation may result in excessive overhead
- New hybrid approach
  - Applies low-overhead sampling to user code
  - Intercepts MPI calls via direct instrumentation
  - Relies on efficient stack unwinding
  - Integrates measurements in statistically sound manner

Joint work with **Lawrence Livermore National Laboratory**

Zoltan Szebenyi et al.: Reconciling sampling and direct instrumentation for unintrusive call-path profiling of MPI programs. In Proc. of IPDPS, Anchorage, AK, USA. IEEE Computer Society, May 2011. (to appear)
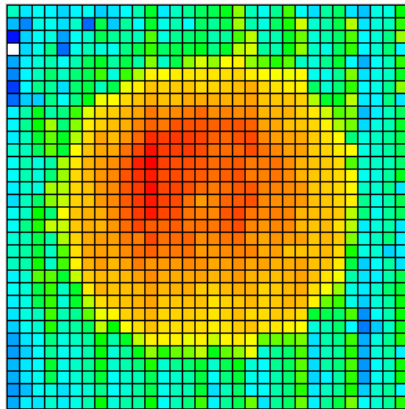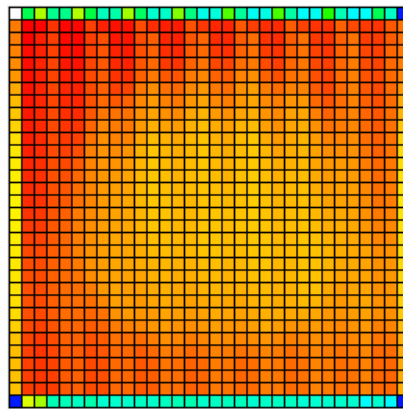
# Delay analysis



- Classification of waiting times into
  - Direct vs. indirect
  - Propagating vs. terminal
- Attributes costs of wait states to delay intervals
  - Scalable through parallel forward and backward replay of traces

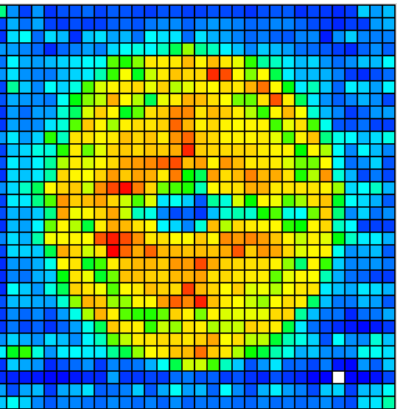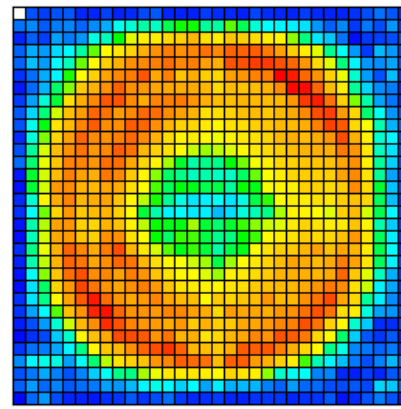# Delay analysis of code Illumination

- Particle physics code (laser-plasma interaction)
- Delay analysis identified inefficient communication behavior as cause of wait states
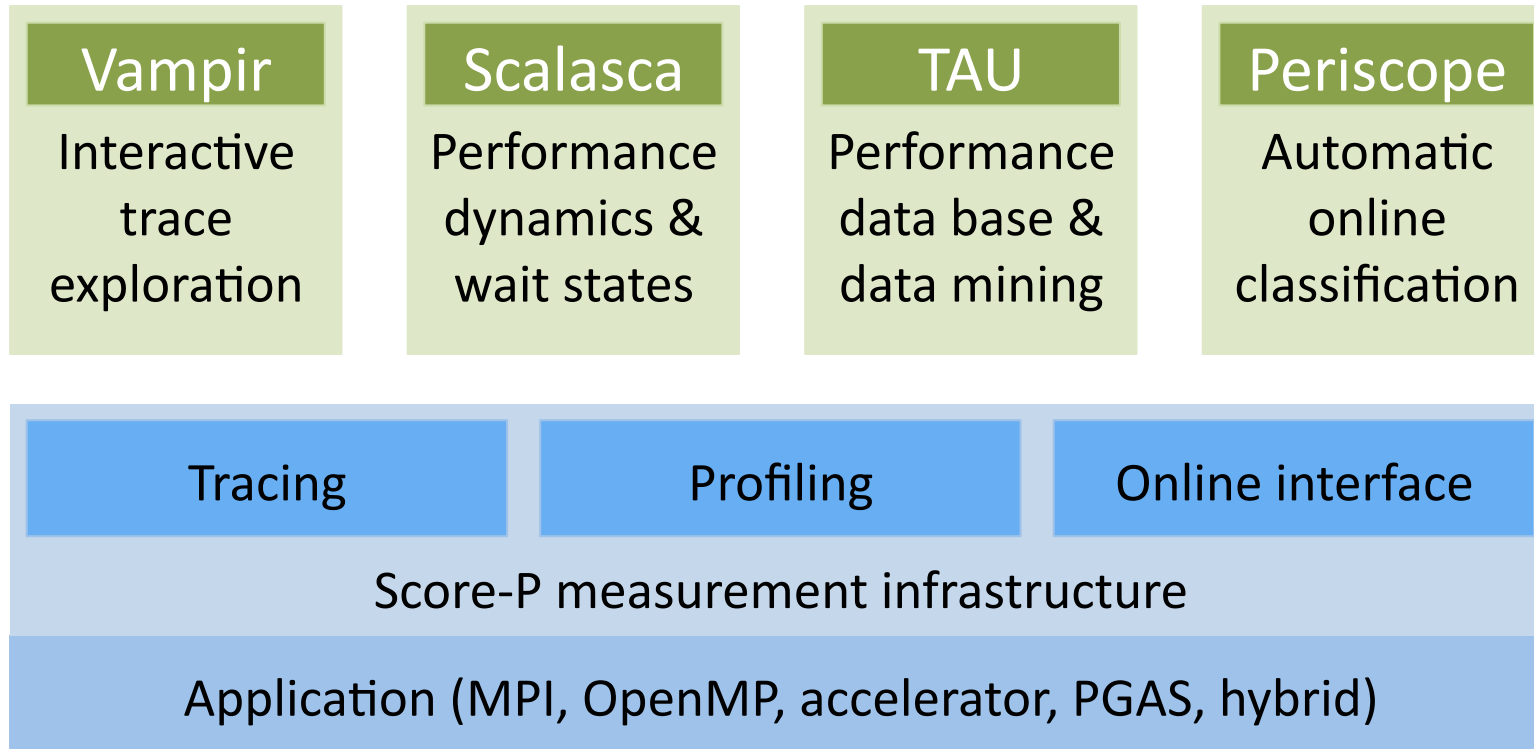


Computation

Propagating wait states:
Original vs. optimized code

Costs of direct delay
in optimized code

# Score-P measurement system

| Vampir | Scalasca | TAU | Periscope |
|---|---|---|---|
| Interactive trace exploration | Performance dynamics & wait states | Performance data base & data mining | Automatic online classification |

| Tracing | Profiling | Online interface |
|---|---|---|

Score-P measurement infrastructure

Application (MPI, OpenMP, accelerator, PGAS, hybrid)

JÜLICH FORSCHUNGSZENTRUM

German Research School for Simulation Sciences

gns

TUM Technische Universität München

RWTH AACHEN UNIVERSITY

TECHNISCHE UNIVERSITÄT DRESDEN

O UNIVERSITY OF OREGON

# Future work

- Further scalability improvements
- Emerging architectures and programming models
  - PGAS languages
  - Accelerator architectures
- Interoperability with 3$^{rd}$-party tools
  - Common measurement library for several performance tools

## The virtual institute in a…



- Partnership to develop advanced programming tools for complex simulation codes
- Goals
  - Improve code quality
  - Speed up development
- Activities
  - Tool development and integration
  - **Training**
  - Support
  - Academic workshops
- www.vi-hps.org

# Thank you!