# Performance Optimization

Philip J. Mucci
Minimal Metrics

[phil@minimalmetrics.com](mailto:phil@minimalmetrics.com)

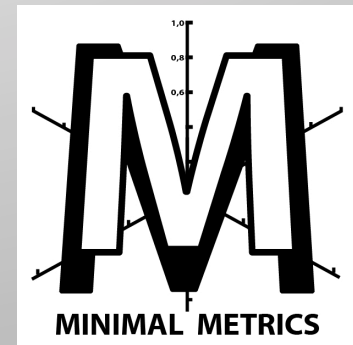04/17/2012

# Overview

o Background

o Customers

o Observations

   o Expertise

   o Programming models

   o Cloud, Virtualization and HPC

   o Tooling

   o I/O

   o Challenges for modeling

# About Me

- Consulting since 1997
  - Software and hardware performance in HPC
  - HPC system software design
  - Parallel algo. opt.
- MS from UTK at ICL under Jack Dongarra
  - PVM, PAPI
- Research Consultant since 1998.

- Software architect, BD and app. Engineering at SiCortex.
- Founded Samara Technology Group in 2008.
  - Application and hardware performance experts for hire.
- Founded Minimal Metrics in 2012.

Philip Mucci - Minimal Metrics - IDC HPC User Forum

# About Minimal Metrics

o A deep network of *the experts*.

o Evaluation, optimization and software engineering.

o Architectural evaluation.

o Moving into small-scale strategic consulting.

  o Logistic and process optimization.

  o Data collection and analytics and forecasting.

o Cofounded with Tushar Mohan

# Reasons For Services

Our favorites:

o We want to *understand* the performance of _____.

    o "Predictive vs. Reactive"

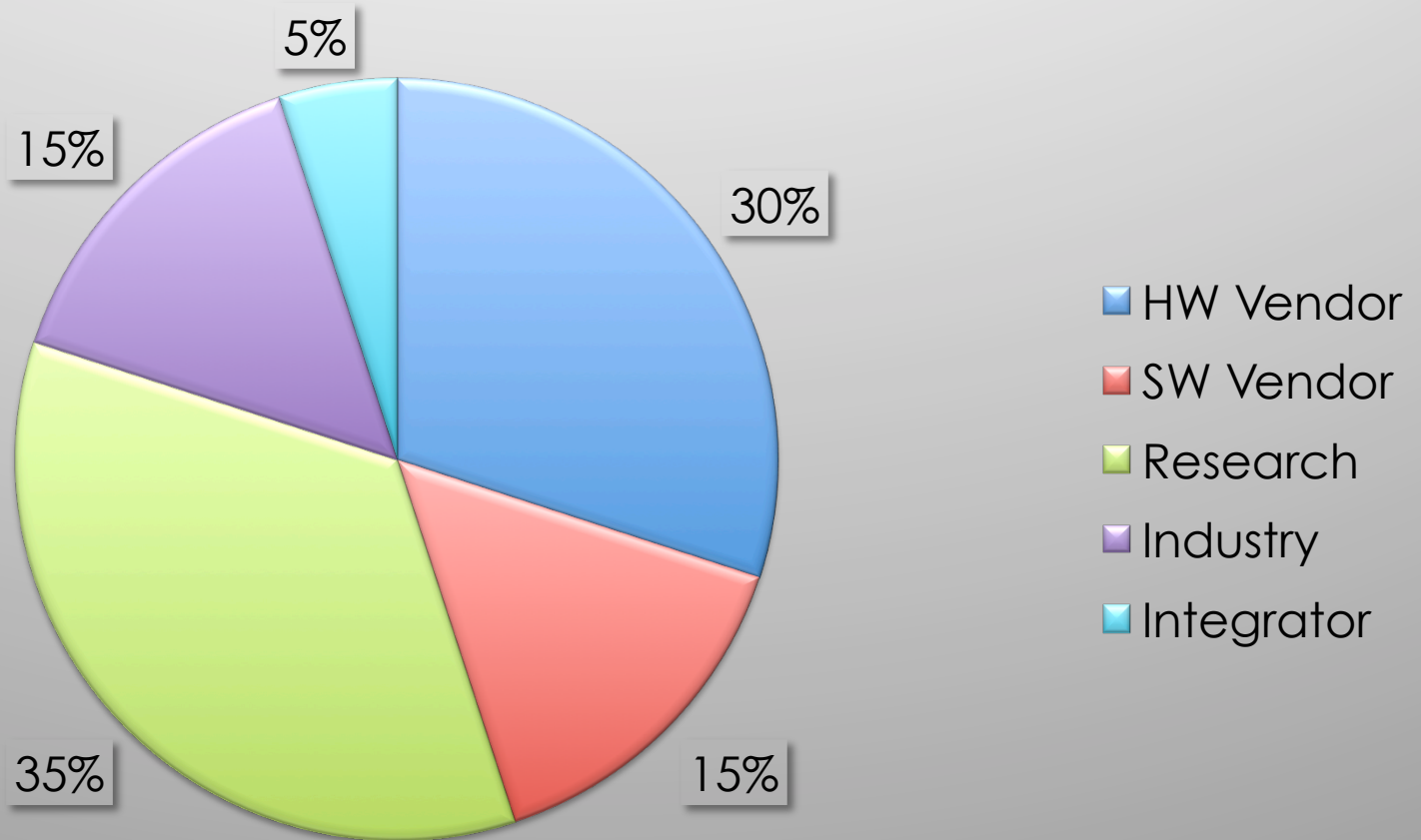o We want to implement _____ using the most experienced talent [that we can't hire].
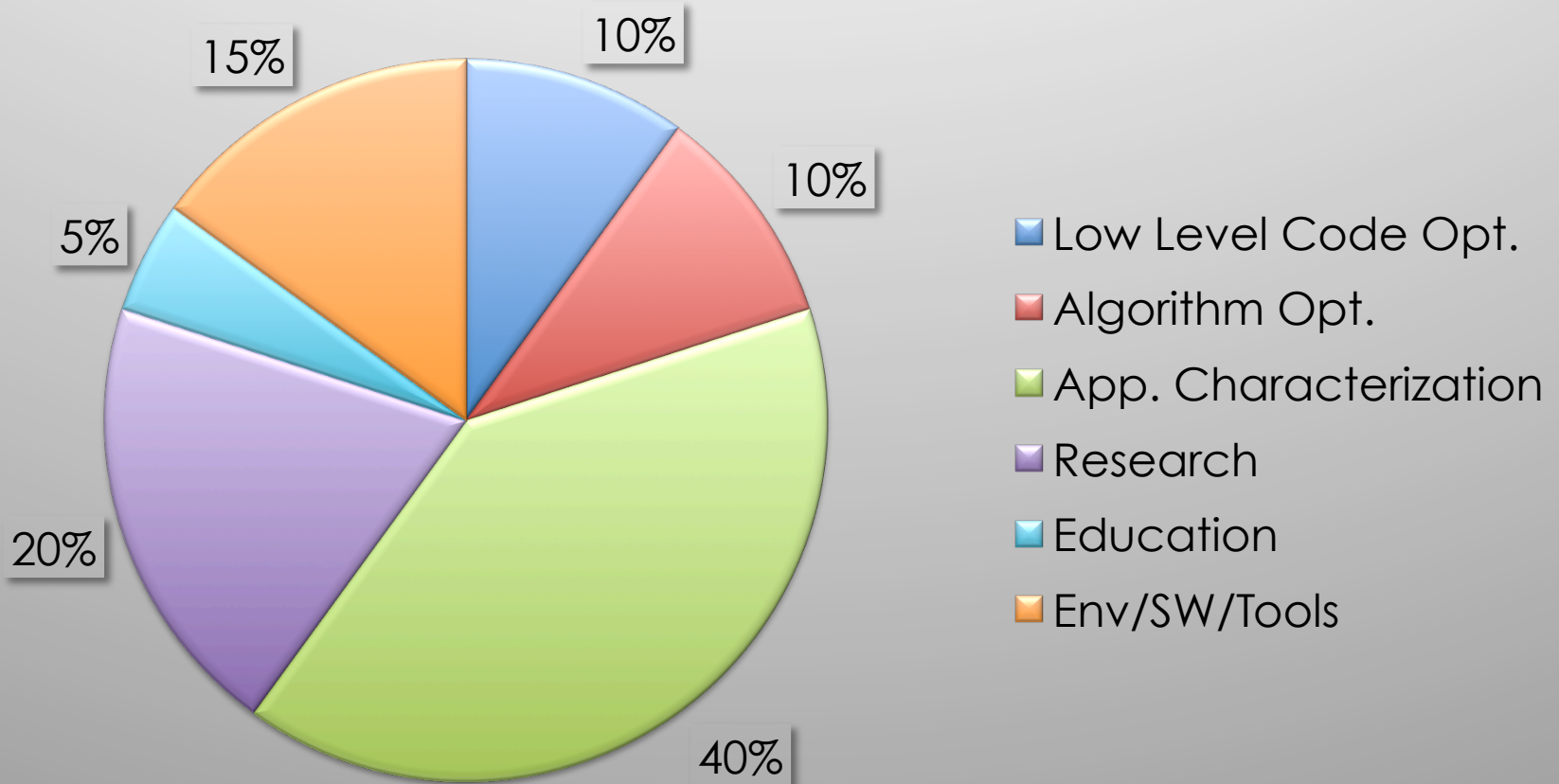


"The Cleaner" – La Femme Nikita, 1990

Some are preventable:

o We implemented our own _____ solver.

o We wrote this using _____ for ease of maintenance.

o We decided that _____ was the right technology to use.

# Services by Customer



Pie chart legend:
- HW Vendor — 30%
- SW Vendor — 15%
- Research — 35%
- Industry — 15%
- Integrator — 5%

# Services by Type



- Low Level Code Opt.
- Algorithm Opt.
- App. Characterization
- Research
- Education
- Env/SW/Tools

# Expertise

o Lack throughout industry

   o Concentrated in specific verticals, with technology vendors and in research.

o Hiring for a critical task can be impossible.

   o Outsourcing is viable and cost effective.

o Educational curriculum catching up, but latency is long.

---

**Sample qualifications from a job posting:**

- Minimum 10 years related experience in a large scale R&D HPC environment.
- Expert knowledge using parallel programming techniques (e.g. MPI, OpenMP, pthreads), parallel programming languages (e.g., C, C++, F90) and scientific simulation and/or data analysis.
- Experience with parallel file systems, common data formats like NetCDF and HDF5, high-performance networking and storage systems.

# Parallel Programming Models

o We often propose (limited) library and directive-based programming.

    o Code can be easily reduced, verified and retargeted.

o Low-level technology adopted relatively quickly.

    o A bit of buyers remorse.

o Many abstractions come at a cost:

    o But good compute is available through native methods.

    o Robust data movement remains a bottleneck.

        o Between SW and HW components.

    o Limited tooling support

        o Analysis pipeline can require code transformation.
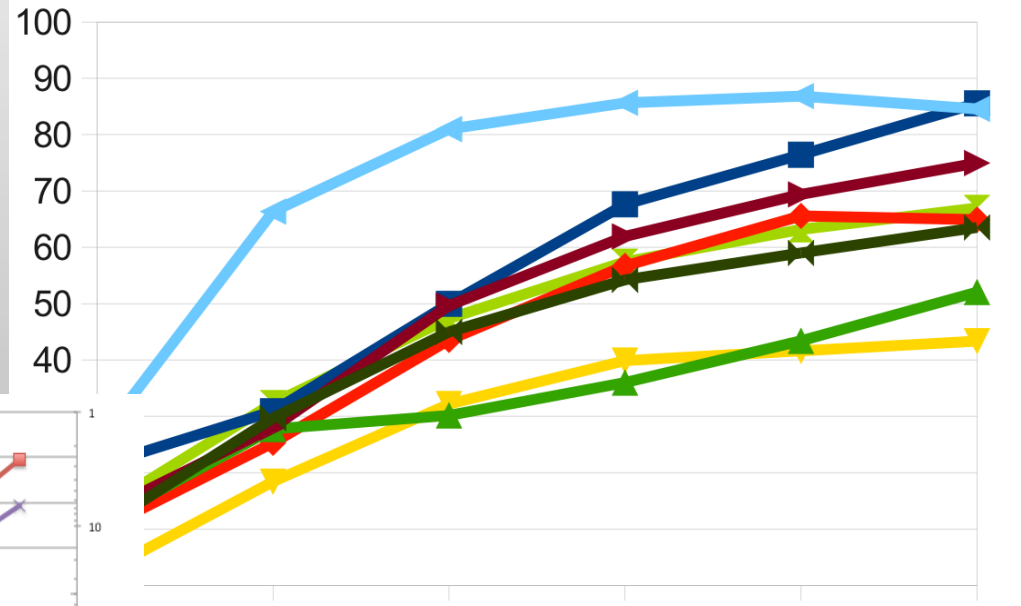
# Cloud Technology, HPC & Grid Reflux

o Options exist for turnkey HPC cloud-based environments.

  o Yet plenty of integration work remains for HPC.

o Single node performance near parity, including decent I/O.

o Communication's (and thus parallel) performance getting there:



% MPI Time

  o Per-core network bandwidth is limited.

  o Lack of low-latency, high-bandwidth comm. capability through the VM.
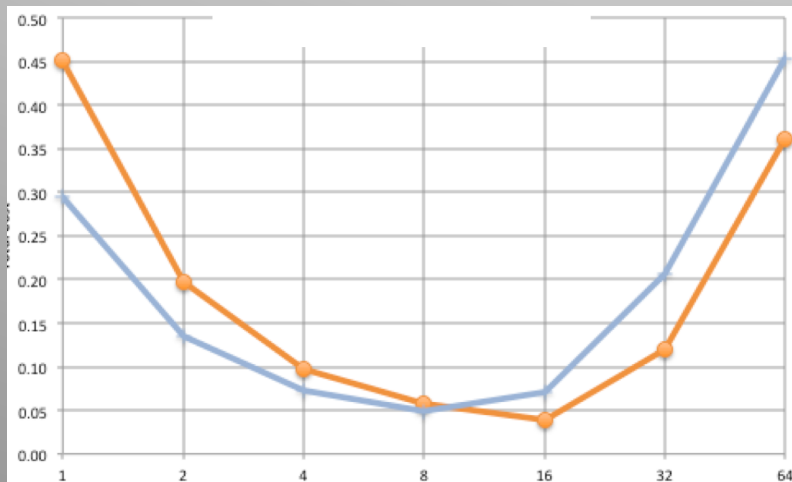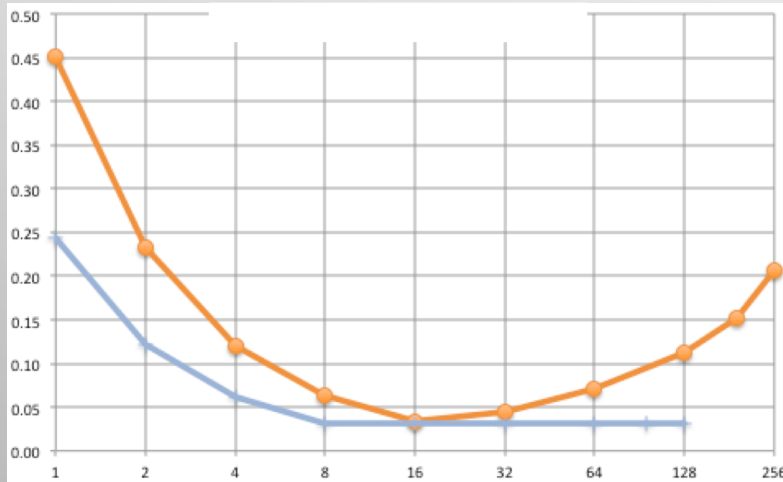
# Optimization and Virtualization

- Largely environmental.
  - OS and software stack
  - I/O and MPI
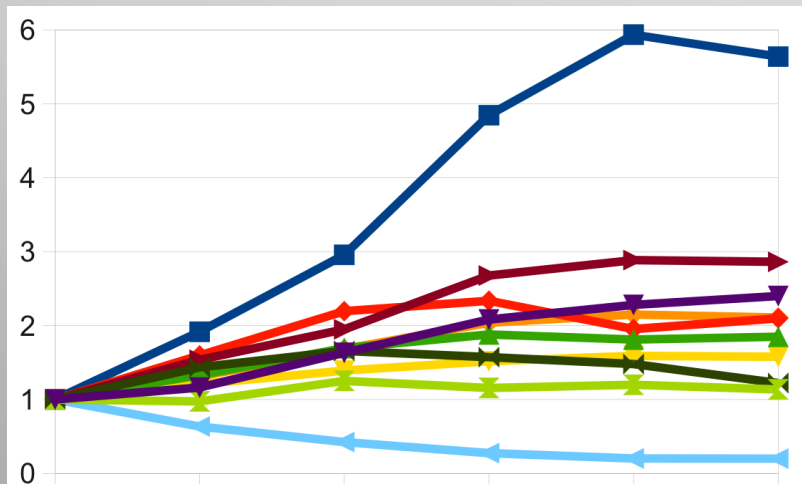- Intra-node MPI still quite good.
- Not so with off-node





- Scaling is lost due to high MPI latencies for un-accelerated comm. in VM's for MPI

Philip Mucci - Minimal Metrics - IDC
HPC User Forum

# Economics of HPC in the Cloud



- Pricing requires very good scaling to be cost effective.
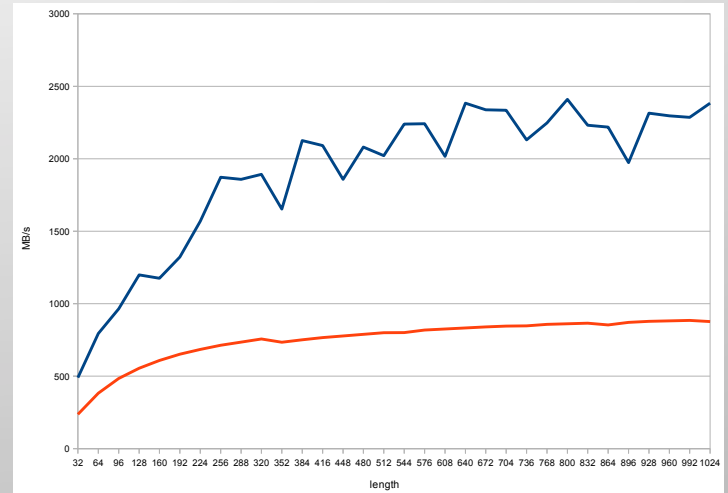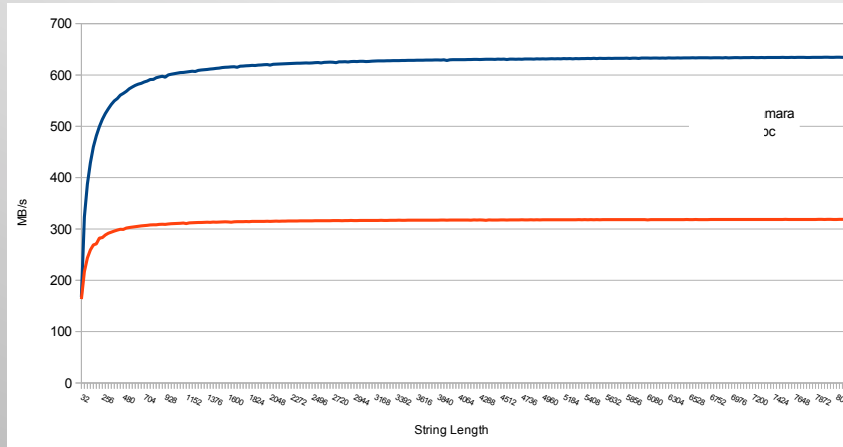- Scaling's worth is related to the importance of the problem.
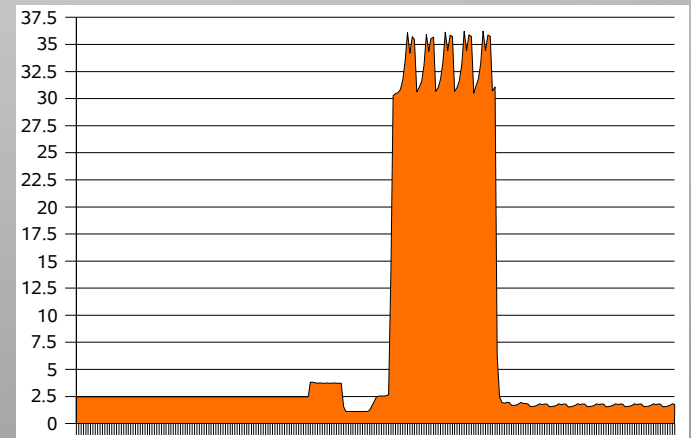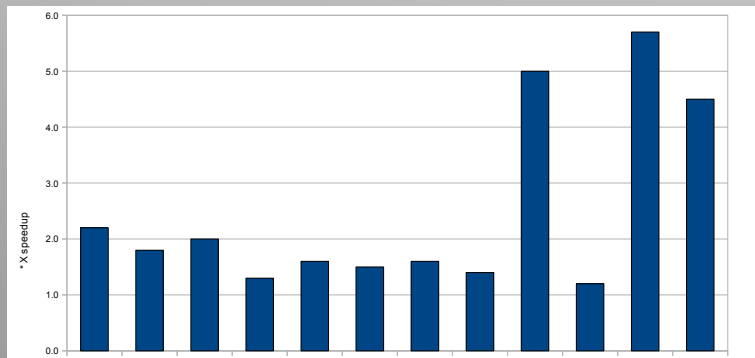
# Optimizing ISV Applications



ISV scaling well below state-of-the-art

- Code is immutable.
  - And rarely changes.
- System optimization.
  - Uptime
  - System configuration
  - Libraries*
  - Parallel run-time
  - Storage
  - CPU availability
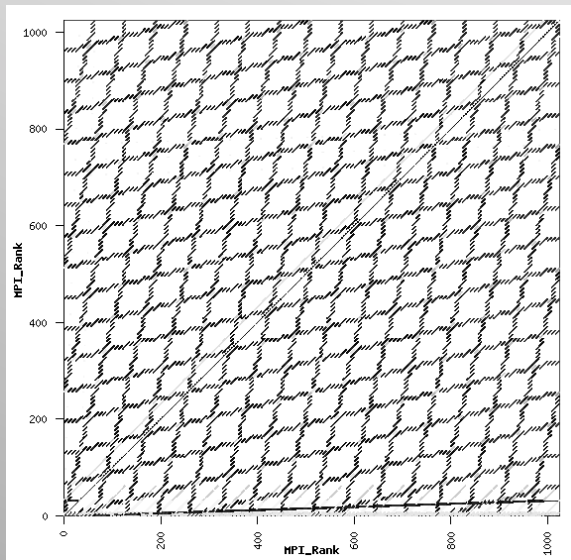- 80/20 rule.

# Performance of GNU/Linux





libscm ceil_sp Performance

- Ain't what you think it is for emerging architectures.



st_memset, 1k to 64k performance

Philip Mucci - Minimal Metrics - IDC
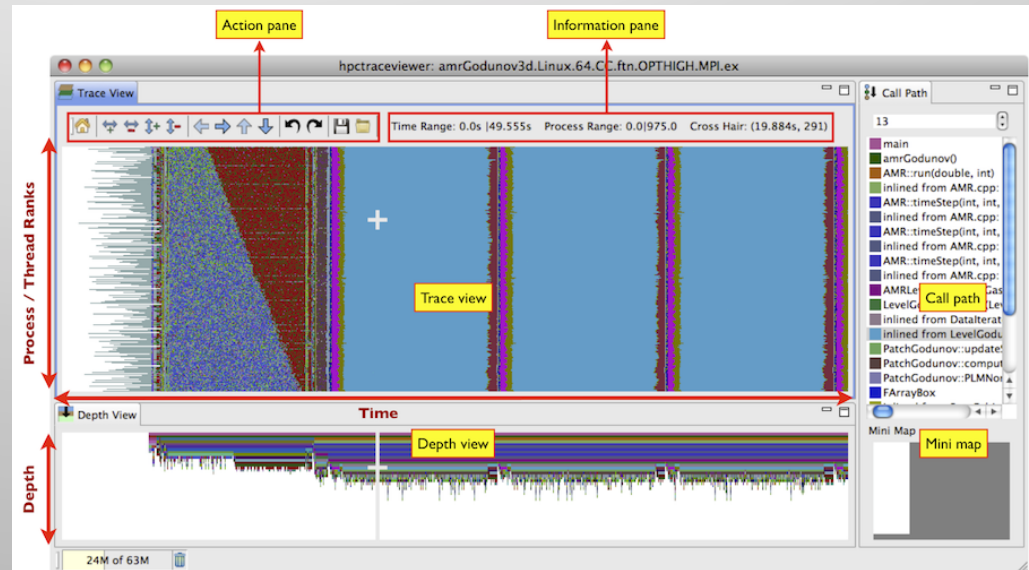HPC User Forum

# Optimization Tools

o Some excellent commercial and open source tools now available.

- o Some require far more knowledge than others to be effective.
- o Tools for MPI, OpenMP, I/O, GPUs and processors down to the instruction level.
  - o Many now include *time* as a dimension of measurement.
- o Focus is on bottom-up view: explain global performance through local observations.
- o Much more robust collection, visualization (and some prediction) capabilities.

o Still lacking full job performance accounting.

# Advanced Performance Visualization



**IPM**
Point to point data flow
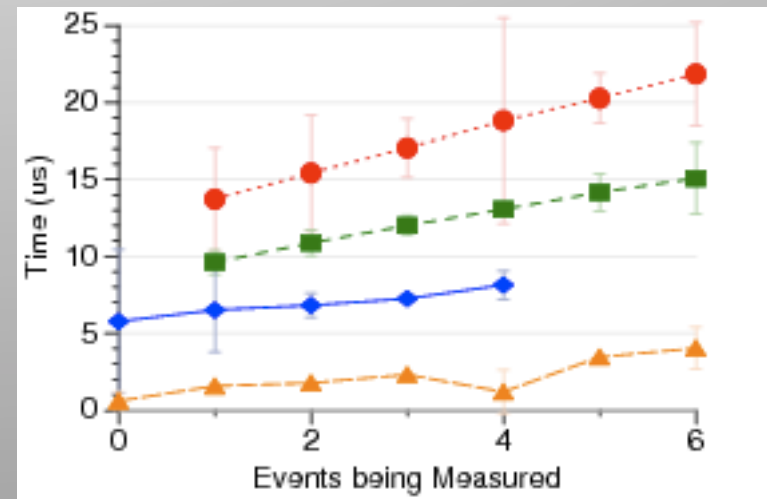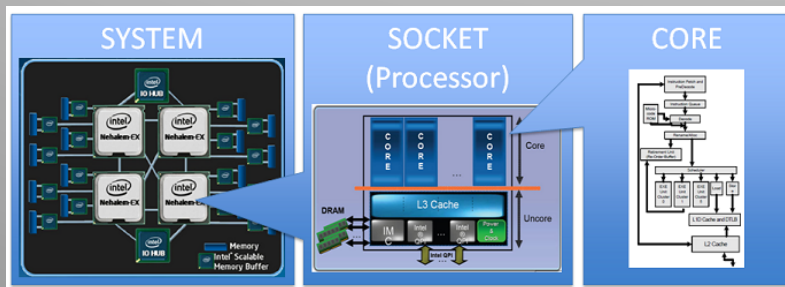
**HPCToolkit**
Metric vs. Task vs. Function (and Depth)

# Tools Workflow

o Naïve methods:
  o Instrument and recompile.
    o But changes characteristics of original code.
  o Measure time only.
    o But answers "where", but not "why" and "by how much"
o Methods now are largely passive and in-situ.
  o Instrumentation is inserted in binary form at run-time.
    o Or by the compiler with knowledge that this code is special.
  o Measure application, operating system and hardware performance events **that are relevant and actionable**.
  o Do so with minimal intrusion.

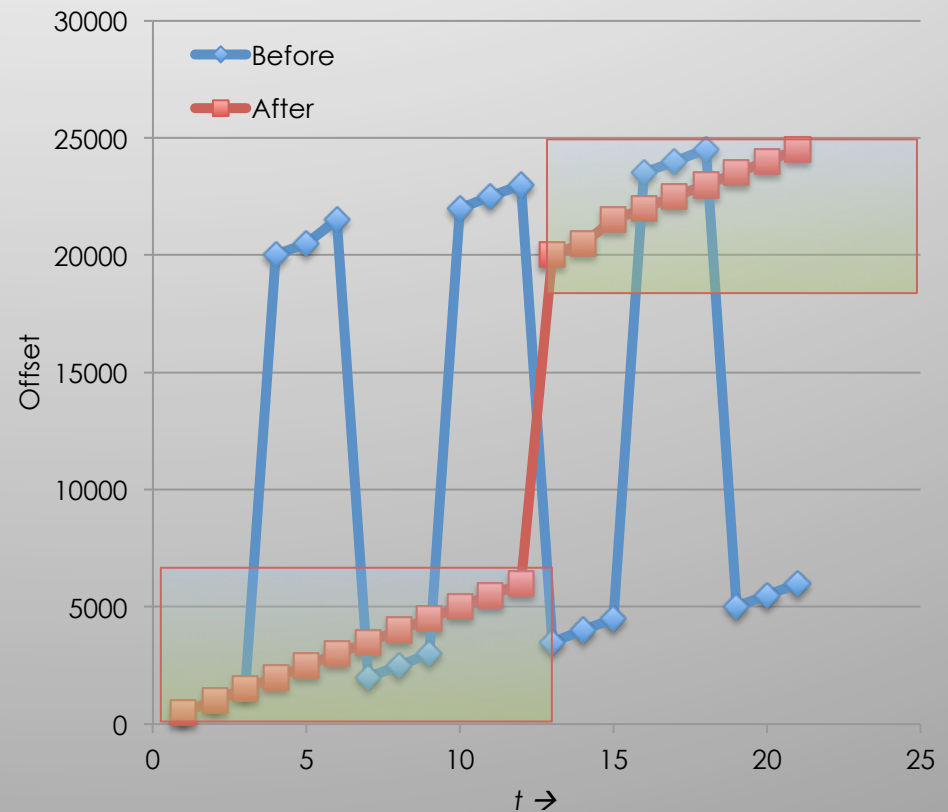# Performance Monitoring

o Hardware PMU's
  o Logic capable of counting and sampling events of interest.
  o Now both on and off-core and in many devices.
o Software
  o System events with significant performance penalties.
o OS support maturing slowly, often regressing.
  o Low-latency, non-privileged access.
o Access often accomplished through PAPI.
  o Only as good as OS support.

# Tools for I/O

o I/O bounds abound.

o Lack of bottom-up tooling.

  o System-level tools provide device level statistics.

  o Good for capacity & fault diagnosis, not tuning.

Access reordering and caching

# Challenges Related to Modeling

o How will my application run on a new platform?

o Anything other than a kernel is non-trivial.

  o HW monitoring and tools allow us to precisely analyze and predict execution traces, not arbitrary code segments.

o Application performance is now largely data-set dependent.

  o Problems are often irregular and/or sparse.

  o Algorithms may be highly configurable.

  o Convergence criteria may be different.

o The data-set needs to be part of the input vector for any model.

o PMaC @ SDSC reflects the state of the art.

# Software Systems

o Software and knowledge are well behind exploiting what the hardware is capable of.

o Quotes from this morning:
- o "Software hurdles are rising to the top for most users"
- o "Software leadership will become the new battleground"
- o "HPC experts often have a narrow view of a new applied user world"
- o "We require ease of everything and just want it to work"
- o "[Engineers become] too hyped about the tools and not about the problem being solved."

# Thanks Dad.



- o **John Francis Mucci**
  - o 5/19/1942 – 2/7/2010
- o From Ridgway, PA
  - o PhD in High Energy Physics from Carnegie Mellon.
- o Career
  - o Director GSG @ Digital
  - o VP of Sales, Marketing and Technical Research at Thinking Machines
  - o Cofounder and CEO of Topical Net, Links2Go, Continuum Software and SiCortex
- o Married Patricia A. Mucci in 1967. Two sons, Philip and David.

Philip Mucci - Minimal Metrics - IDC HPC User Forum

# Thank You

[phil@minimalmetrics.com](mailto:phil@minimalmetrics.com)

Philip Mucci - Minimal Metrics - IDC
HPC User Forum