

# **allinea**



**Leaders in parallel software development tools**

## **Performance Profiling and Debugging at the Extreme Scale and Beyond**

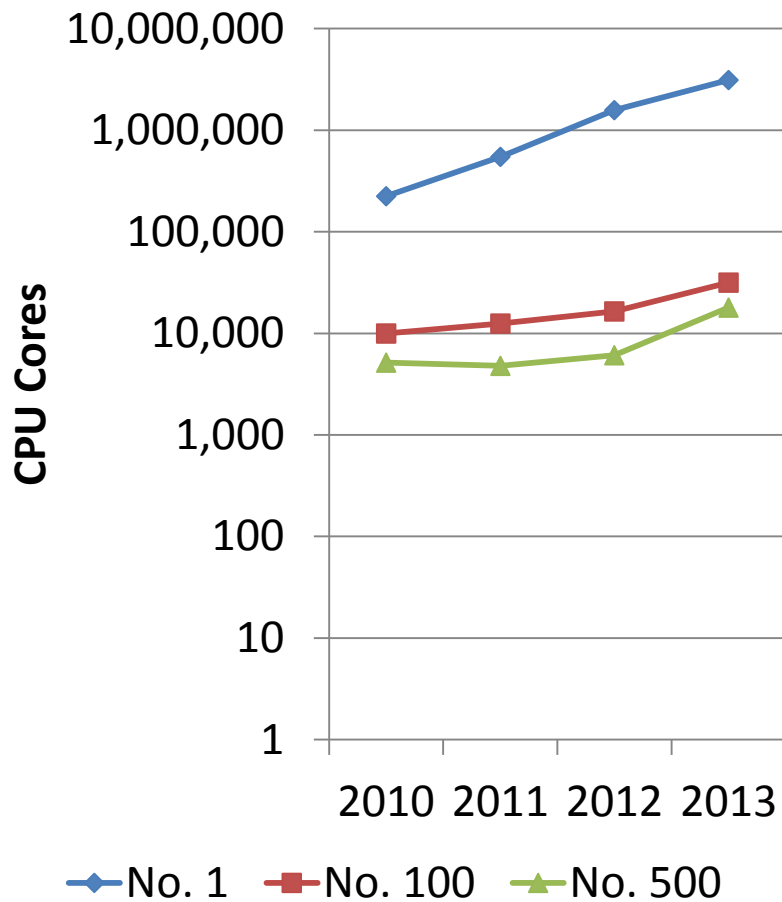
**David Lecomber  
Allinea Software  
david@allinea.com**

# Allinea Software

- **Our mission: to make HPC software development fast, simple and successful**
  - A modern integrated environment for HPC developers
  - Scalable tools for any scale of system
- **Supporting the lifecycle of application development and improvement**
  - Allinea DDT - Productively debug code
  - Allinea MAP- Enhance application performance
- **Designed for productivity**
  - Consistent integrated easy to use tools
  - Enables effective use of HPC resources and expertise



# Extreme machines are everywhere



Machine  
sizes are  
exploding

Software  
scale grows  
as machines  
grow

# Three Challenges for tools



## Scalability

- Speed and Simplification



## Heterogeneity

- Accelerators and Coprocessors

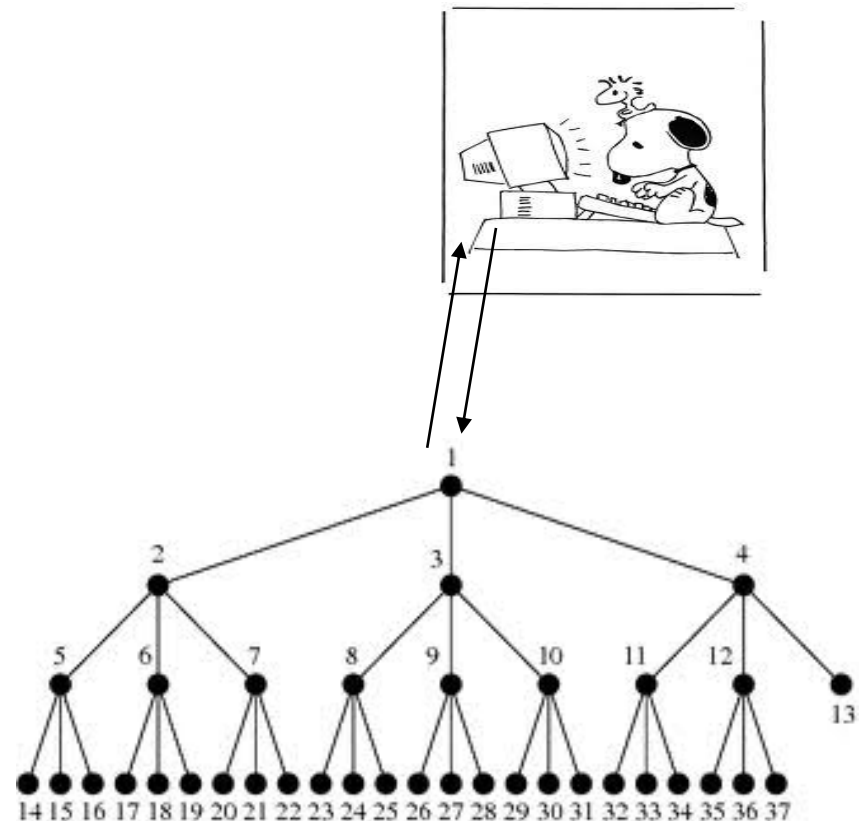


## Adoption

- Ease of Use and Education

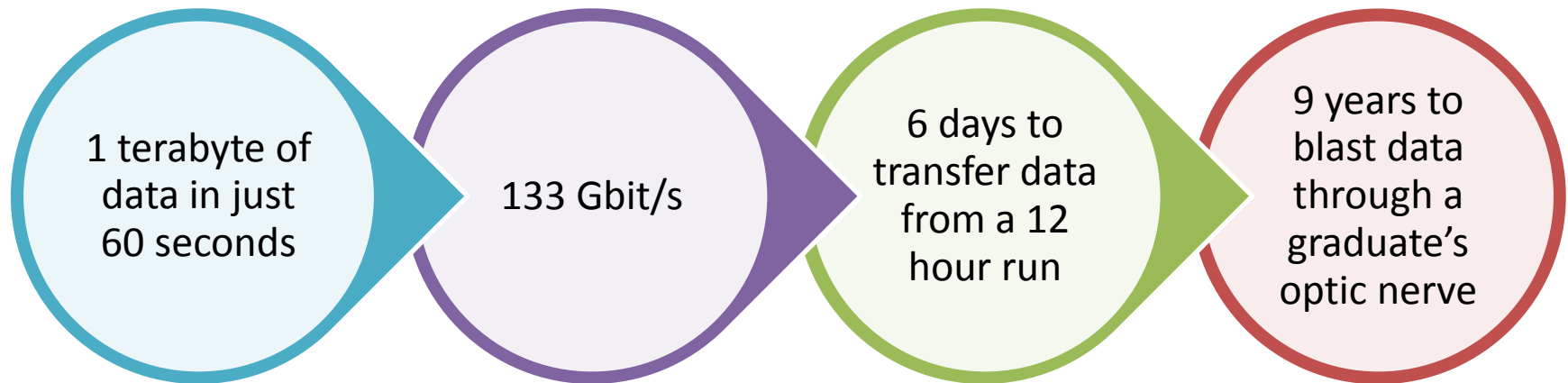
# Beneath Alinea's Petascale Tools

- Scalable tree network
  - Sends bulk commands and merges responses
  - Aggregations maintain the essence of the information
  - Don't send more data than is needed....
- Usability matters
  - The interface is as important as the speed
  - Focus on scalable components

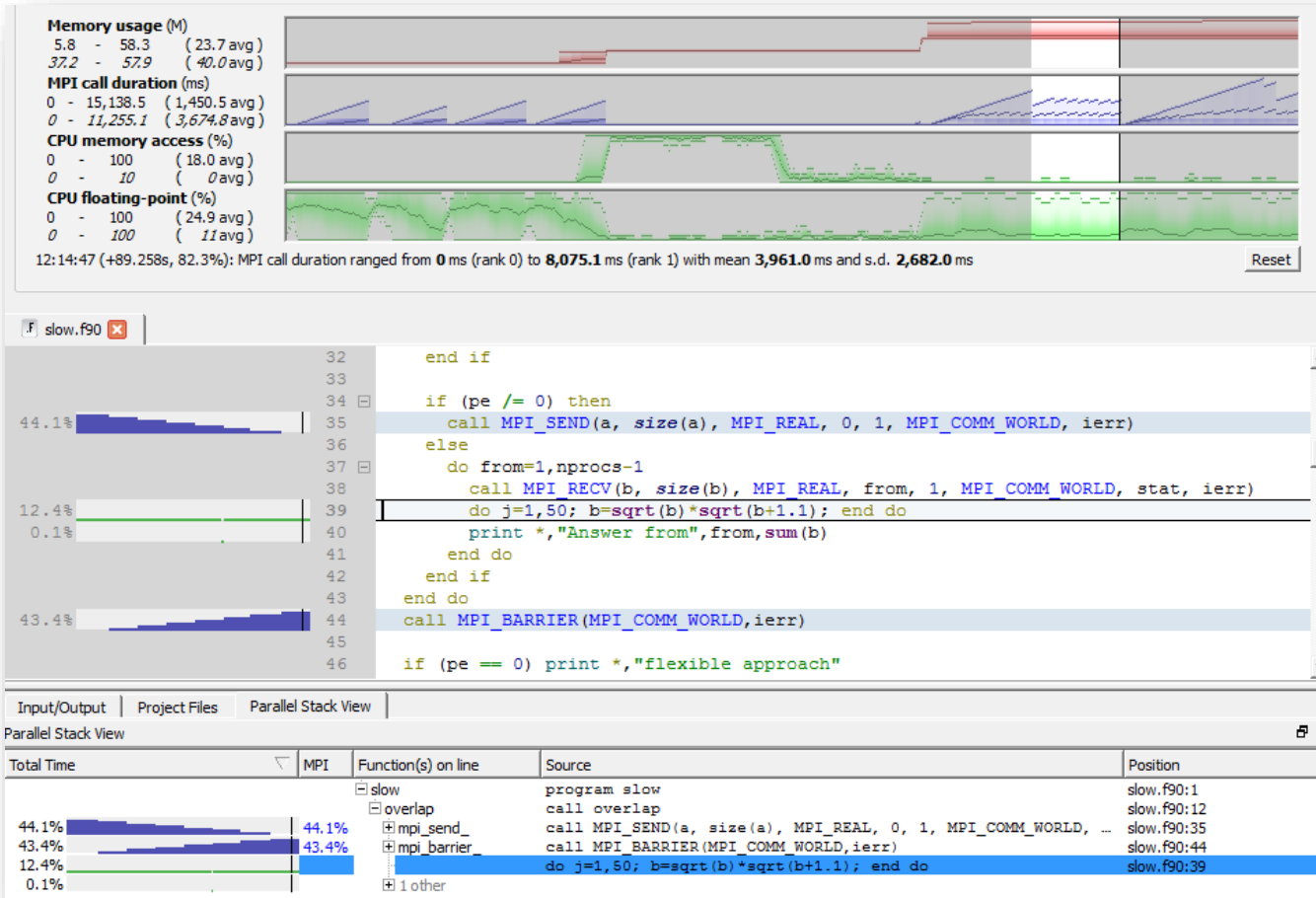


# Beware Exploding Bandwidth Needs...

Trivial 16,000 process wave equation code



# Yield Focussed Example



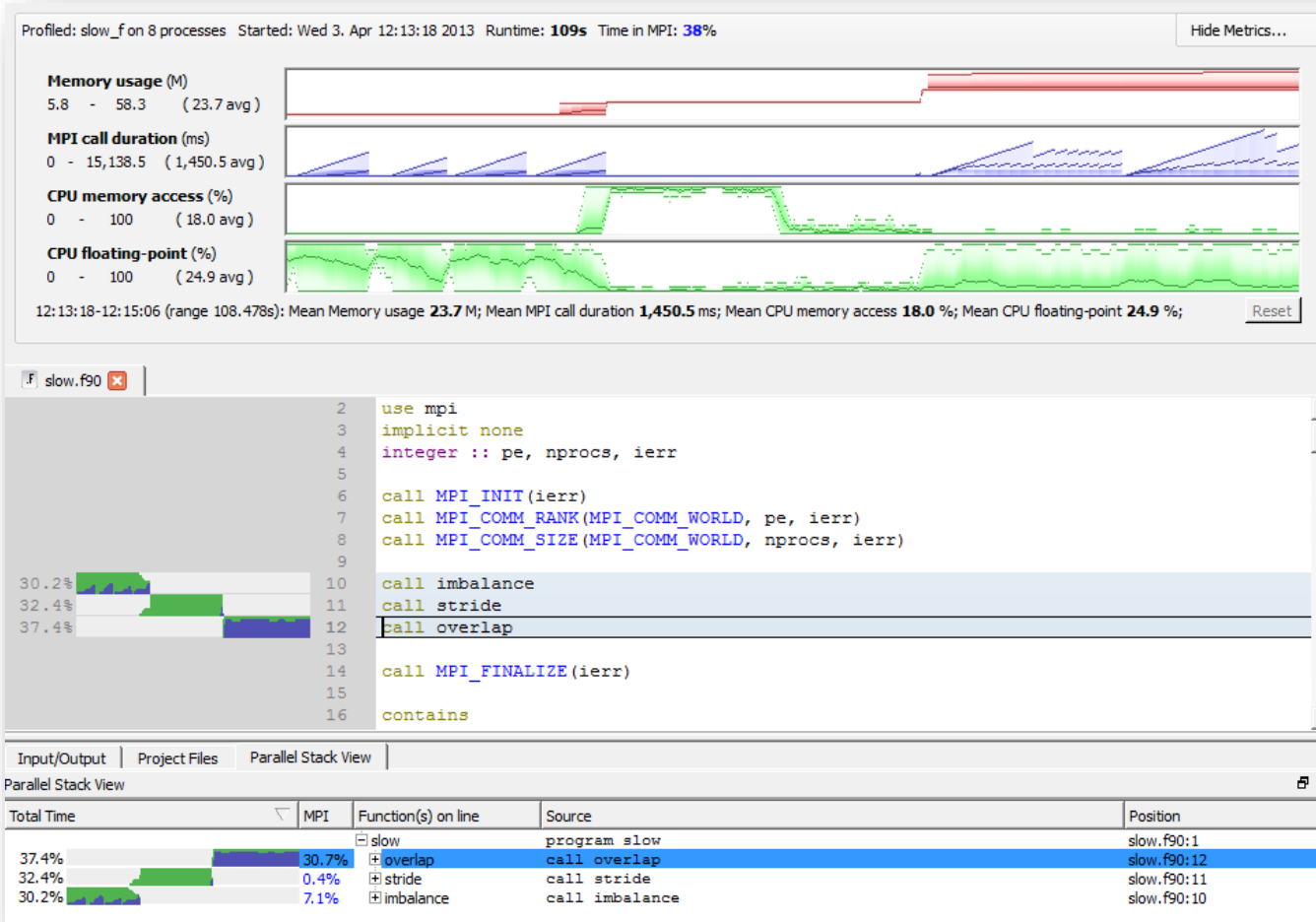
Show distributions and min/max ranks

Show per-line information

Focus on movement through code

Process 0 busy computing on line 39

# Attacking Visual Scalability



Common  
horizontal axis



Aggregate across  
all processes



Highlight  
imbalance visually



Always refer to  
source code



# Complimentary Approaches

## Allinea MAP

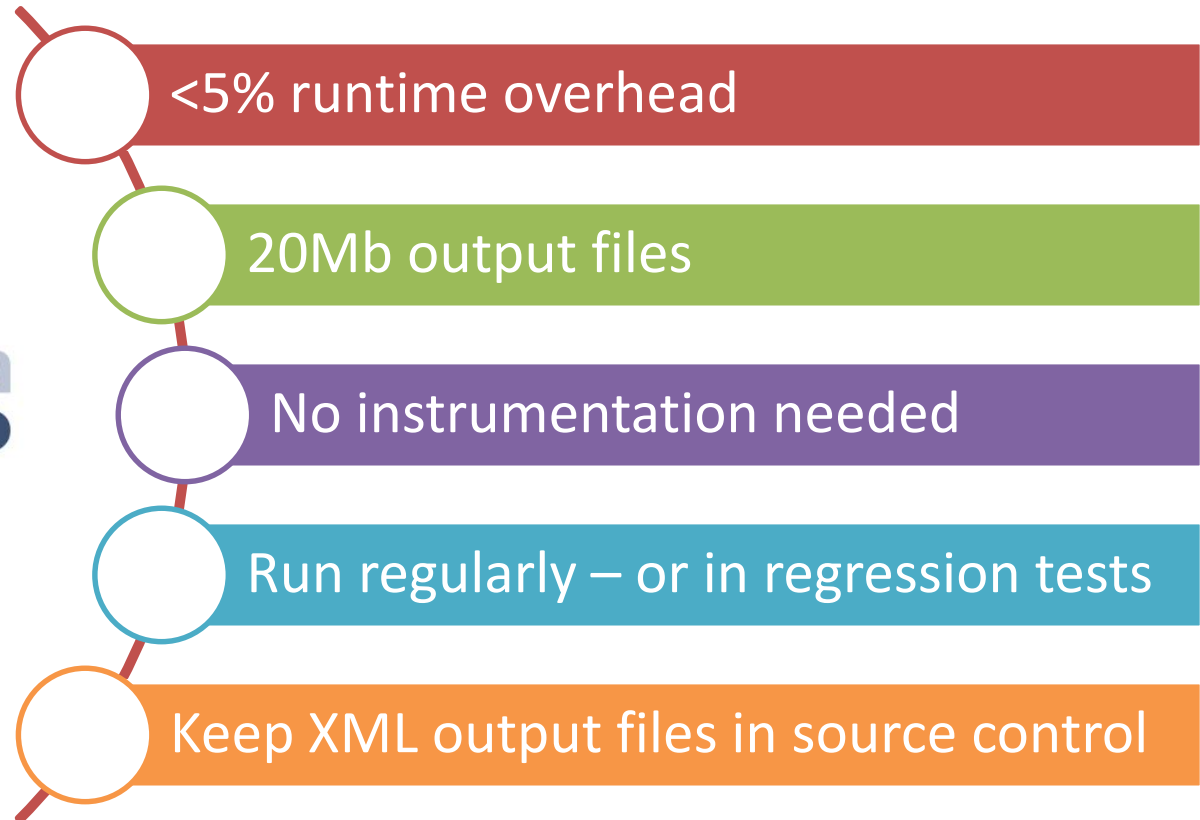
- Quick, low-overhead way to characterize performance
- See which lines of code are hotspots
- Identify common problems at once

## Record Everything

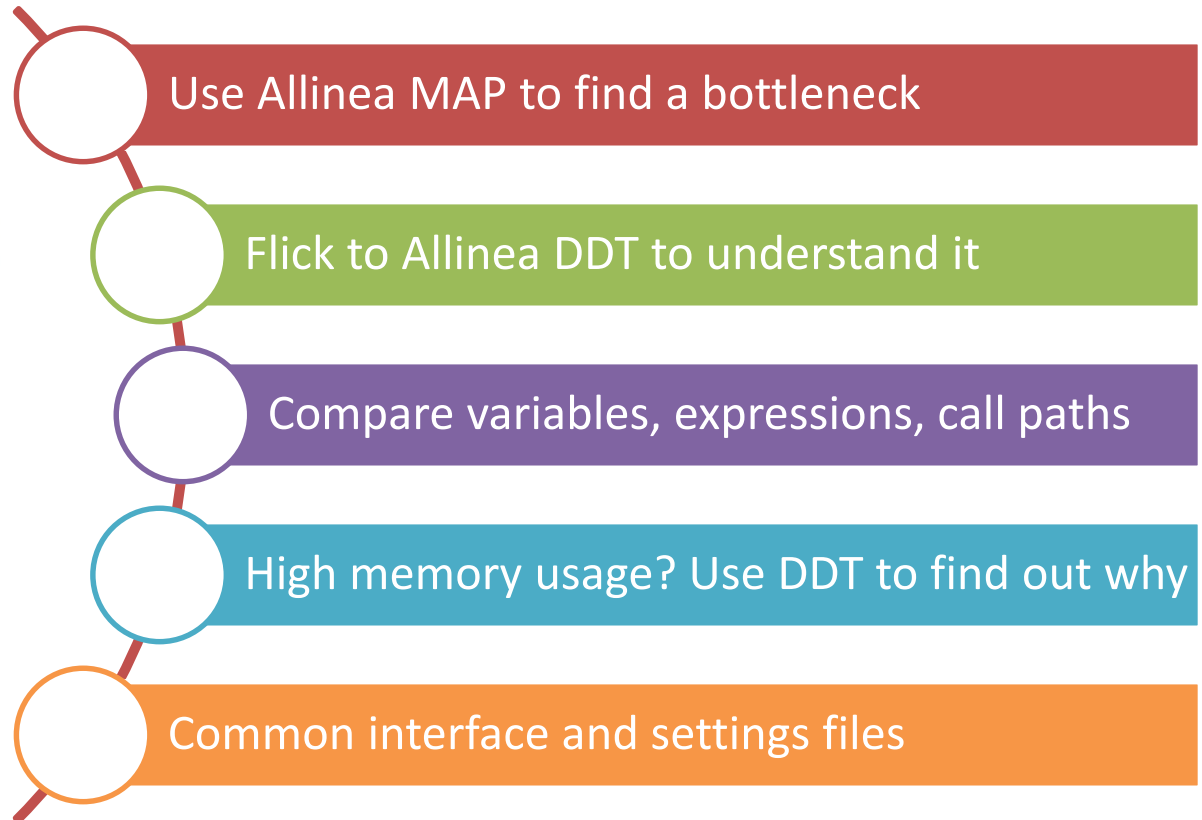
- Pass more obscure problems to an expert
- Now know which loop to instrument and which performance counters should be recorded

# Surprising Benefits

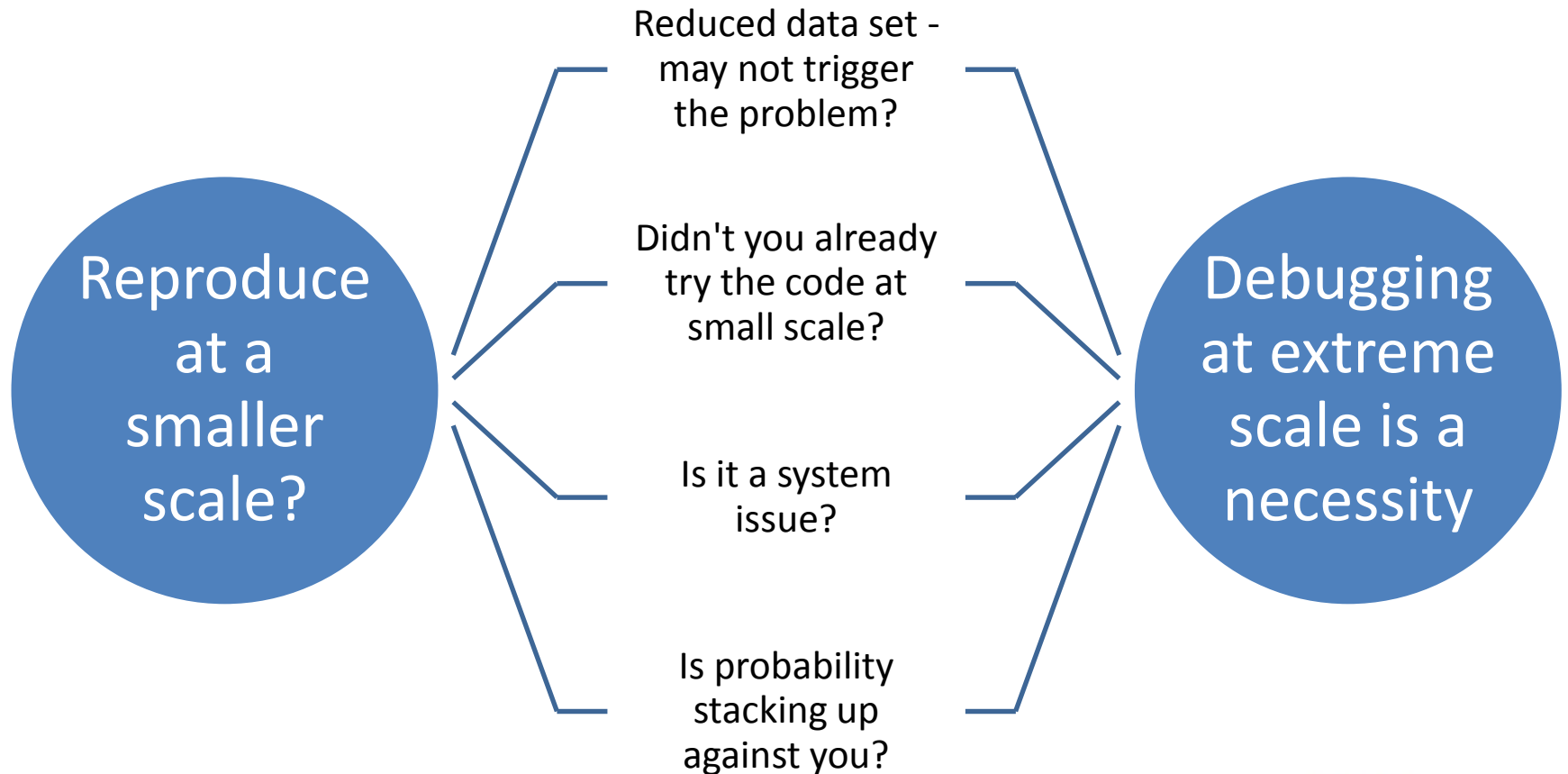
---



# Integrated with Allinea DDT

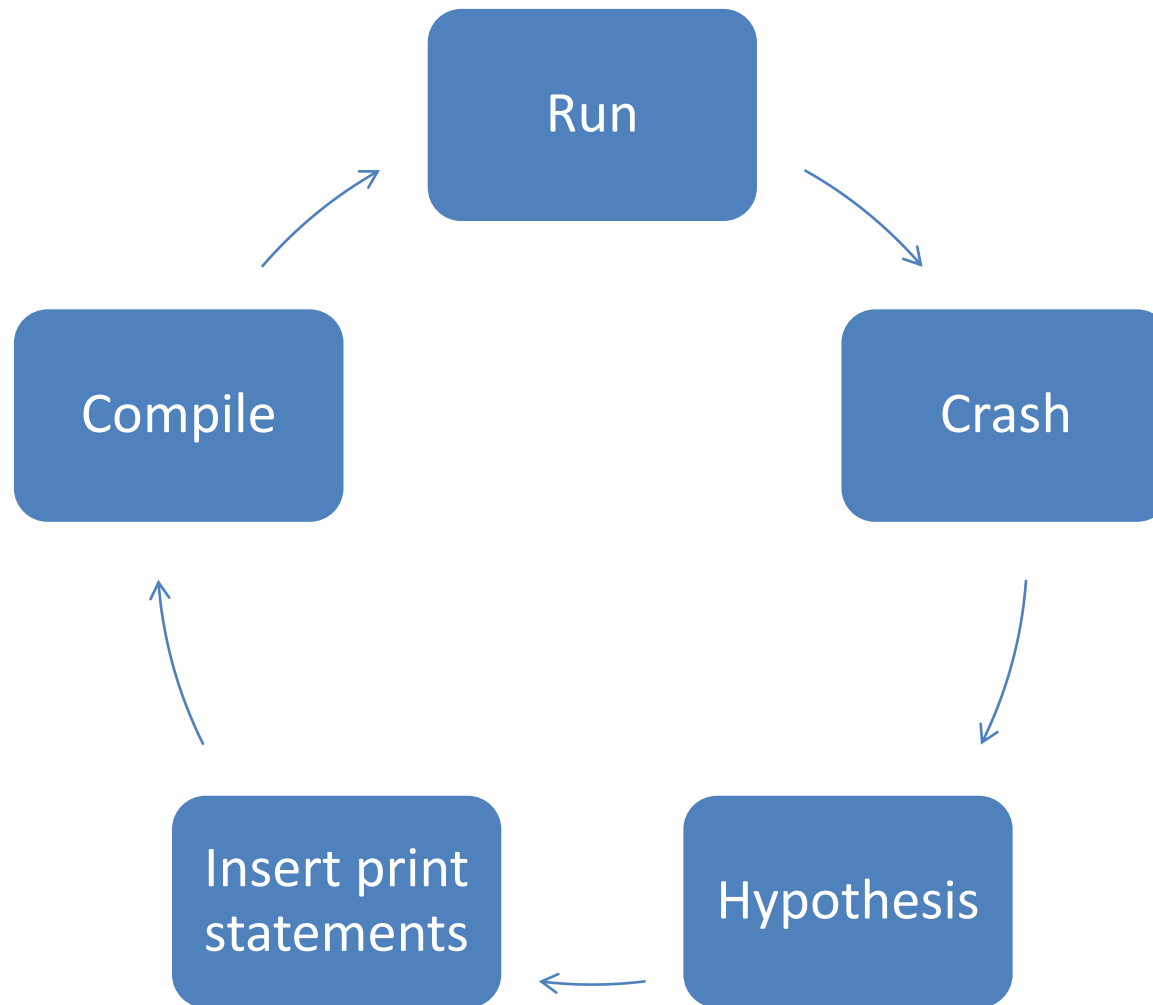


# Bug fixing as scale increases



# Debugging in practice...

---



# Titan and Mira

## Titan

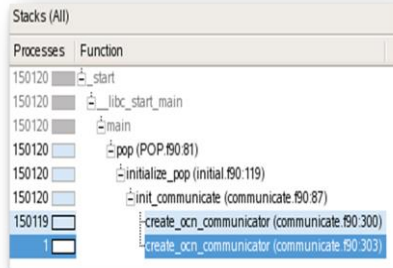
- 18,688 nodes
- 18,688 NVIDIA Kepler K20 GPUs
- 299,008 CPU cores
- 50,233,344 CUDA cores

## Mira

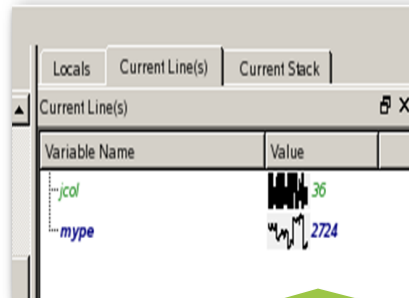
- 49,152 nodes
- 786,432 cores
- 3,145,728 hardware threads

Does the printf workflow “work”?

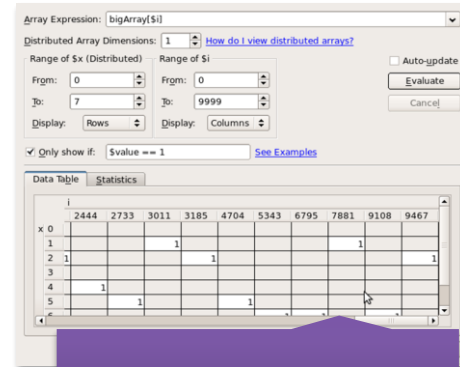
# Top 5 features at scale



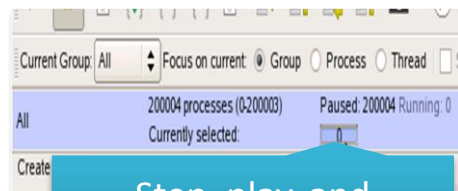
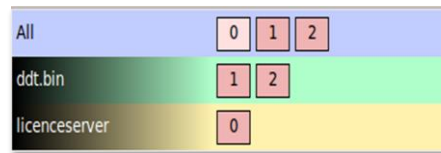
Parallel stack view



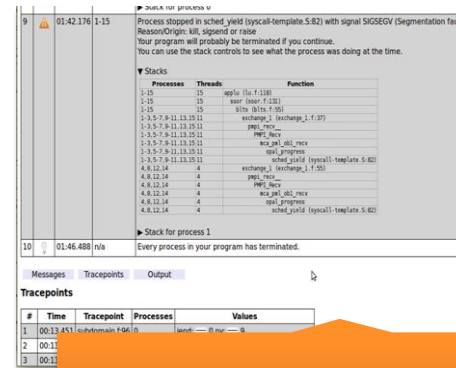
Automated data comparison: sparklines



Parallel array searching



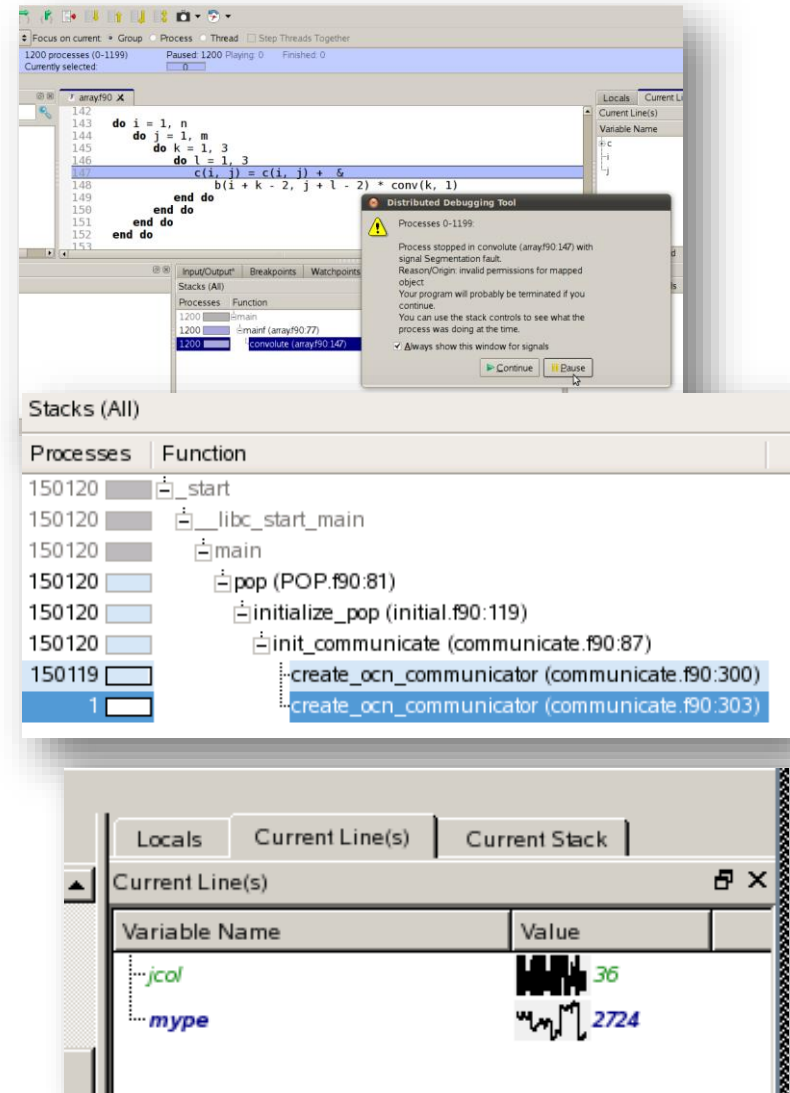
Step, play, and breakpoints



Offline debugging

# Allinea DDT: Scalable debugging by design

- **Where did it happen?**
  - Allinea DDT leaps to source automatically
  - Merges stacks from processes and threads
- **How did it happen?**
  - Some faults evident instantly from source
- **Why did it happen?**
  - Real-time data comparison and consolidation
  - Unique “Smart Highlighting” – colouring differences and changes
  - Sparklines comparing data across processes
- **Force crashes to happen?**
  - Memory debugging makes many random bugs appear every time





## Example – ORNL's Jaguar

---

- HPC code fails on 98,304 cores
- Random processes crashing
- Printf? Which processes and where?
- Too costly to repeat
- Allinea DDT finds cause first time

# Can Allinea MAP help with other tools?

## Profiles are gzipped XML

- Not currently publicly documented
- Provides entire MAP GUI data
- Samples and metrics, source file locations

## *Could* process by scripts

- Auto-instrument for other tools
- Auto-pick tool for next step (VAMPIR or Vtune!)

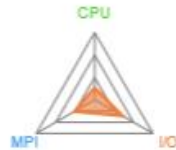
## *Could* have a Plug-in metric API?

- Shared libraries to record metrics and supply to MAP API
- Allinea preloads – almost every MPI
- Take burden of platform matrix explosion away?

# Allinea Performance Reports

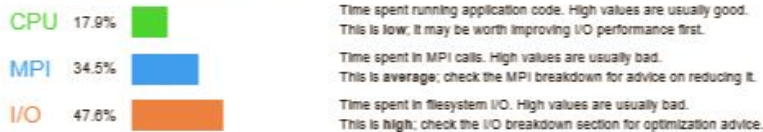


Executable: MADbench2  
Resources: 9 processes, 1 node  
Machine: sandybridge2  
Start time: Mon Nov 4 12:28:45 2013  
Total time: 11 seconds (0 minutes)  
Full path: /tmp/MADbench2  
Notes: 12-core server / HDD / 9 readers + writers



Summary: MADbench2 is **I/O-bound** in this configuration

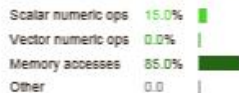
The total wallclock time was spent as follows:



This application run was **I/O-bound**. A breakdown of this time and advice for investigating further is in the **I/O** section below.

## CPU

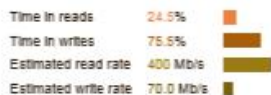
A breakdown of how the 17.9% total CPU time was spent:



The per-core performance is **memory-bound**. Use a profiler to identify time-consuming loops and check their cache performance. No time was spent in **vectorized instructions**. Check the compiler's vectorization advice to see why key loops could not be vectorized.

## I/O

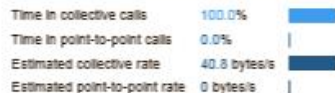
A breakdown of how the 47.6% total I/O time was spent:



Most of the time is spent in **write operations**, which have a low transfer rate. This may be caused by contention for the filesystem or inefficient access patterns. Use an I/O profiler to investigate which write calls are affected.

## MPI

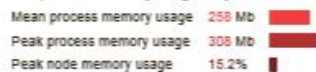
Of the 34.5% total time spent in MPI calls:



All of the time is spent in **collective calls** with a very low transfer rate. This suggests a significant load imbalance is causing synchronization overhead. You can investigate this further with an MPI profiler.

## Memory

Per-process memory usage may also affect scaling:



The **peak node memory usage** is low. You may be able to reduce the total number of CPU hours used by running with fewer MPI processes and more data on each process.

- How well do your **applications** match your **hardware**?
  - How well is application X optimized for **this system**?
  - Does it benefit from running at this **scale**?
  - Are there **I/O** or **networking** bottlenecks affecting performance?
  - Which **hardware**, **software** or **configuration** changes will improve performance?

## This week..

---

- Meet us at booth #1719
  - Get a demo
  - Talk to the team
  - Ask a question
- Enter the draw
  - Can you detect a performance problem?
  - Daily draw – win a Kindle Fire HDX

# More talks during the week

---

- Monday
  - 4.10pm – Extreme Scale Performance Tools Workshop – Room 501
    - *Performance Profiling and Debugging at the Extreme Scale and Beyond*
- Tuesday
  - 11.00am – Intel booth - #2701
    - *Discovering bottlenecks without pain: Get performance on Intel Xeon Phi with Allinea MAP and Allinea DDT*
- Wednesday
  - 11.30am – DoE booth #1327
    - *OpenSHMEM tools*
  - 1.30pm – Fujitsu booth - #2718
    - *Develop efficient HPC applications at scale on FX10*
- Thursday
  - 1.15pm – VisIt/Intelligent Light booth #4216
    - *Allinea DDT and VisIt: debugging HPC applications using a visualization tool*
  - 3.30pm – Exhibitor Forum – Room 501/502
    - *Pick your battles: Getting results faster with Intel Xeon Phi and NVIDIA CUDA*