



High performance tools to debug, profile, and analyze your applications

Debugging and Profiling your HPC Applications

David Lecomber, CEO and Co-founder

david@allinea.com



About this talk



allinea
FORGE
DDT + MAP

- Learn how to debug and profile your code
 - Techniques to take home
- Tools we will use: Allinea Forge
 - Debugging with Allinea DDT
 - Profiling with Allinea MAP
 - NB. Allinea MAP is not supported on BG/Q
- Where to find Allinea's tools
 - > 70% of Top 500 have at least one Allinea tool

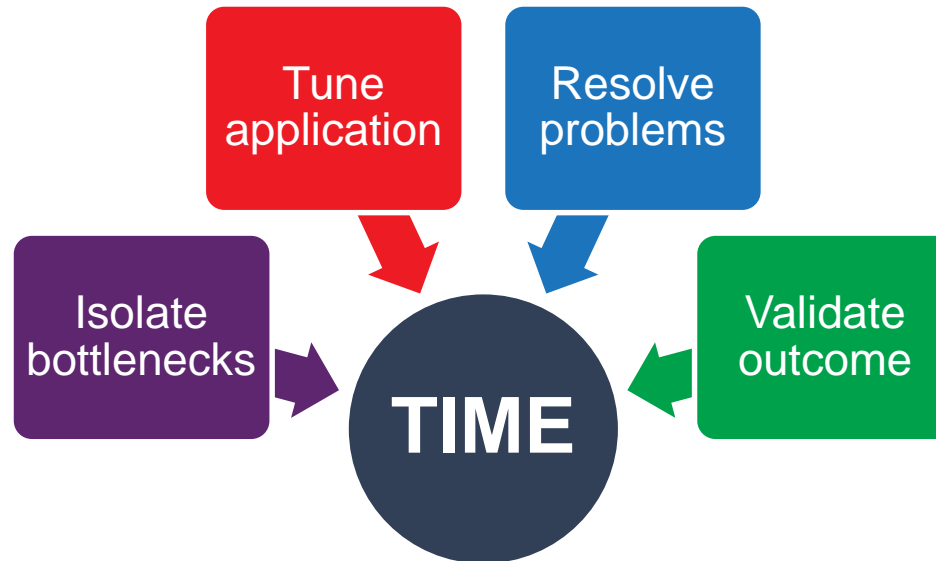


allinea

Motivation

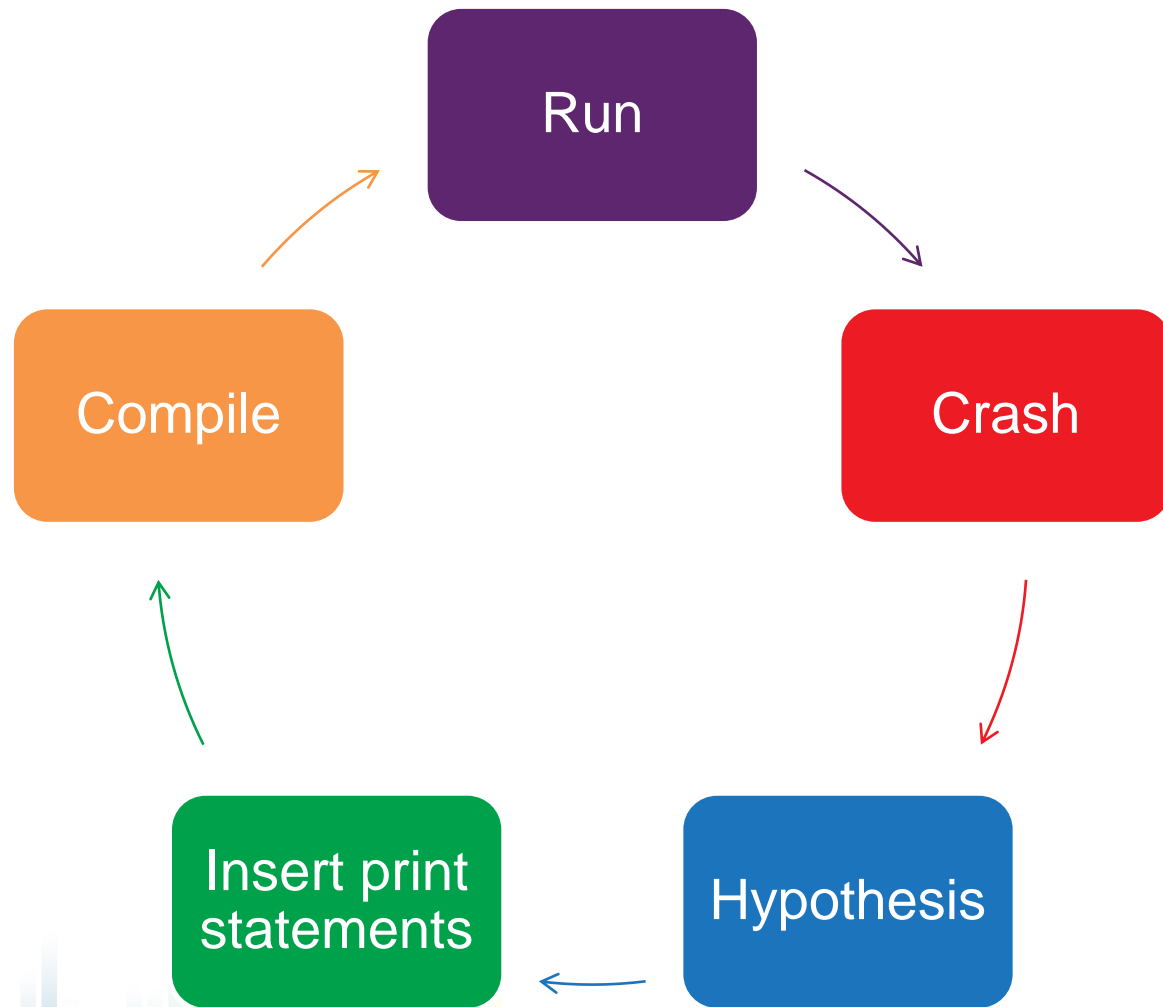
- HPC systems are finite
 - Limited lifetime to achieve most science possible
 - Sharing a precious resource means your limited allocation needs to be used well
- Your time is finite
 - PhD to submit
 - Project to complete
 - Paper to write
 - Career to develop
- Doing good things with HPC means creating better software, faster
 - Unrivalled productive and easy-to-use development environment...
 - ... To help reach the highest level of performance and scalability
- High performance parallel code needs tools designed for the challenge

Use the right software tool to be faster

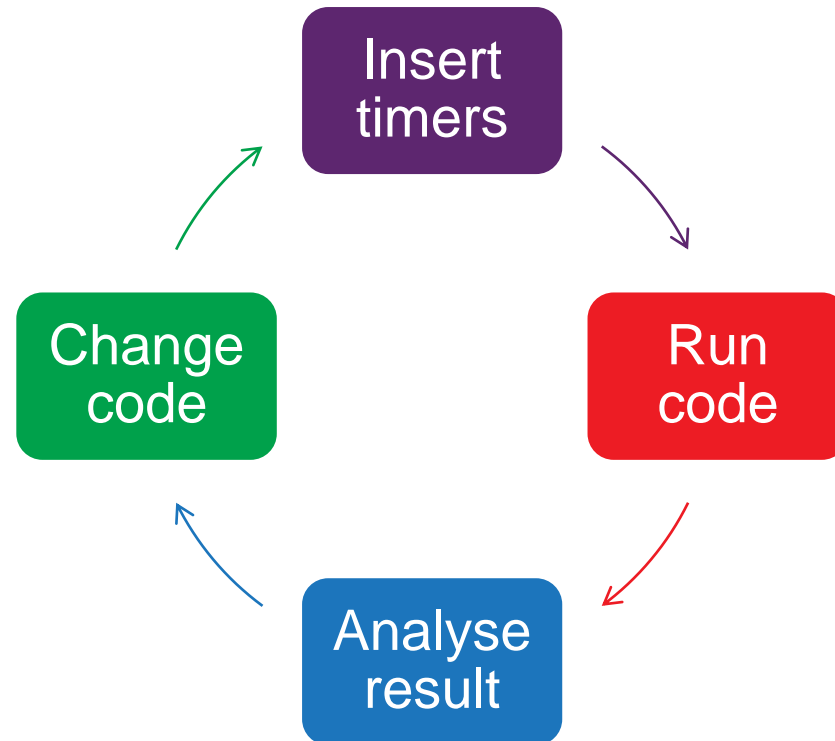


- What parts of the code would benefit most from being rewritten?
- How should I modify a code to make it better (or work at all)?

Debugging in practice...

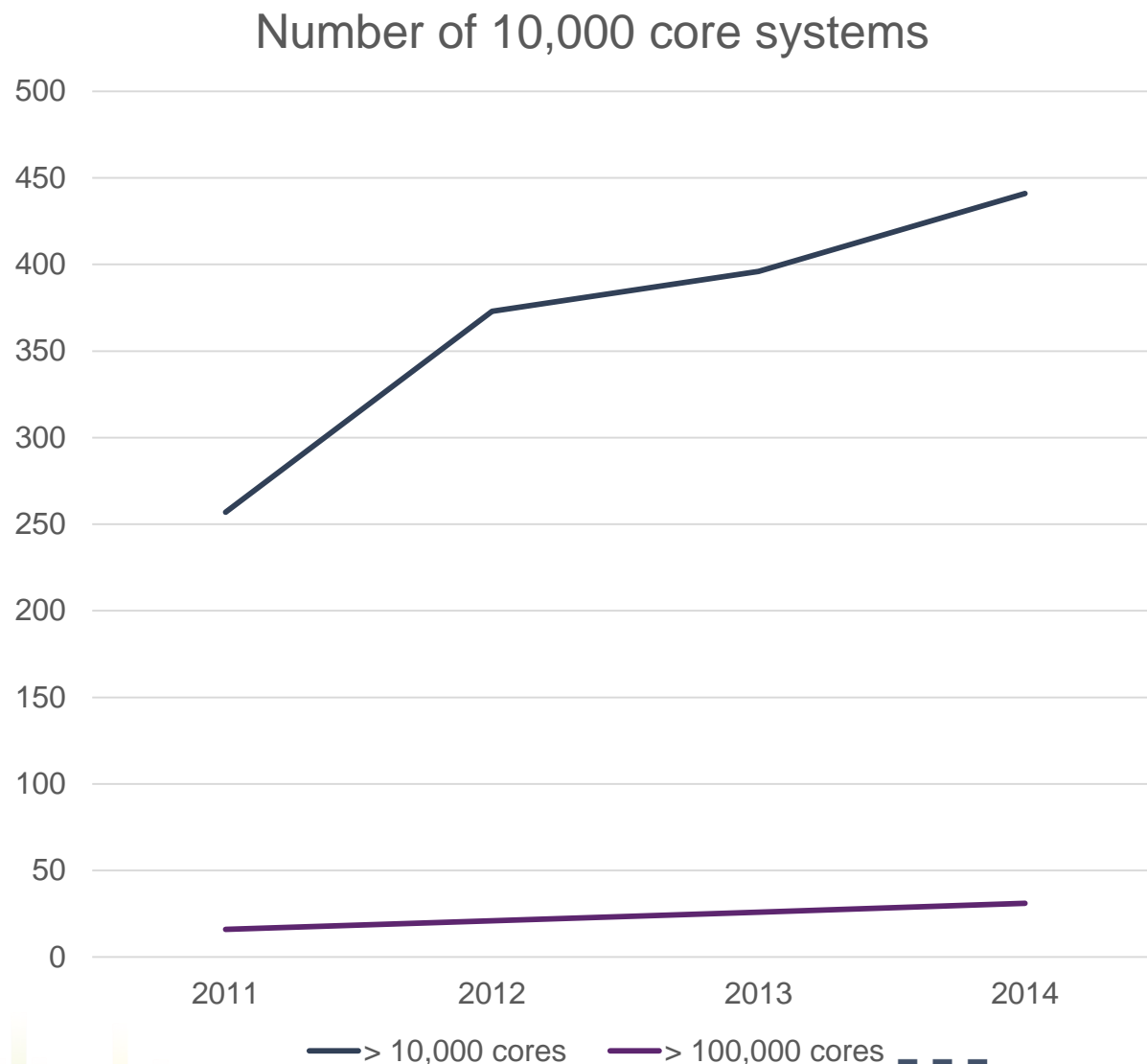


Optimization in Practice

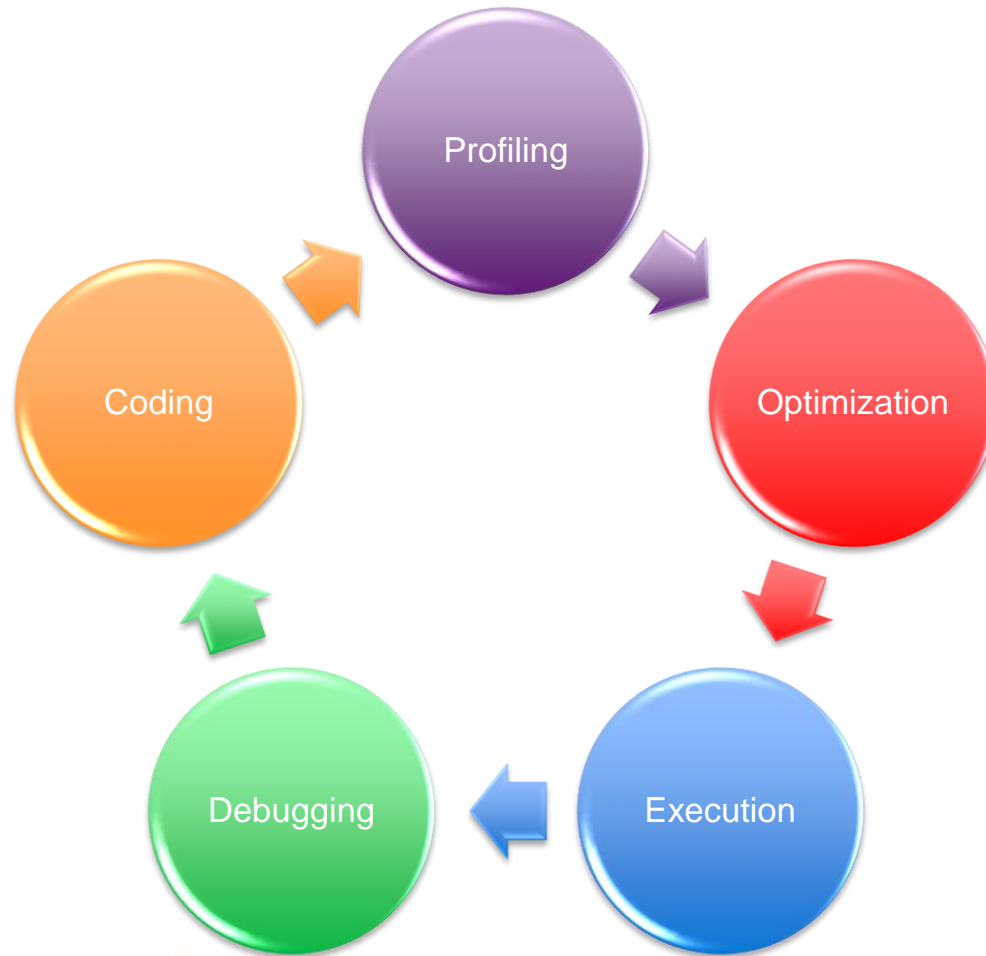


Motivation

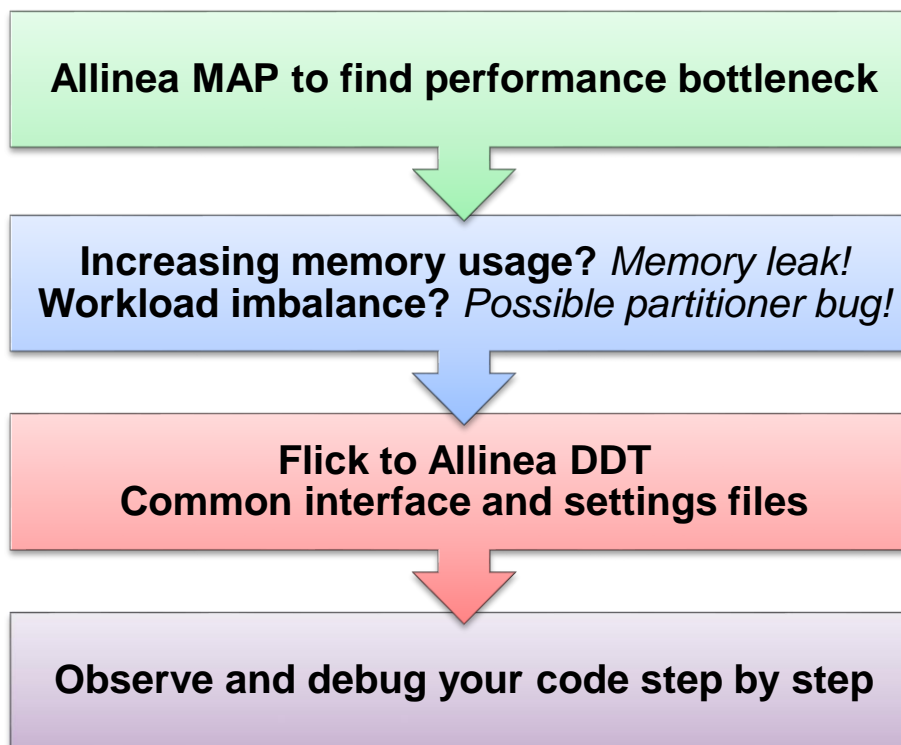
- “Without capable highly parallel software, large supercomputers are less useful”
 - Council on Competitiveness
- “1% of HPC application codes can exploit 10,000 cores”
 - IDC, 2011



Application Development Workflow



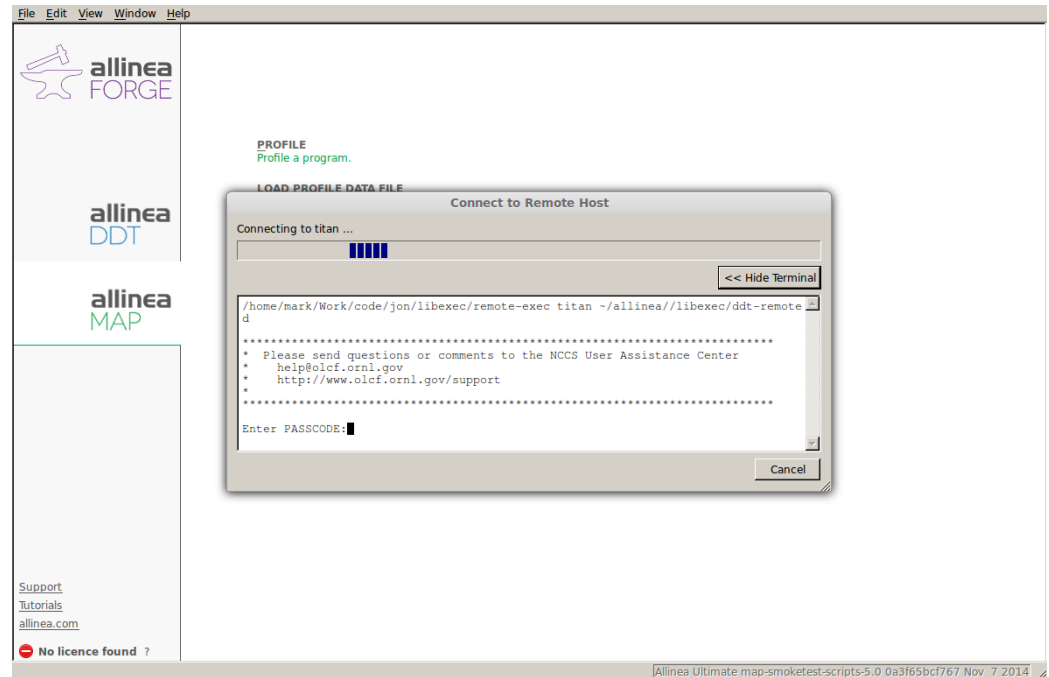
Hello Alinea Forge!

The Alinea logo, consisting of the word "allinea" in a bold, dark blue, sans-serif font. In the bottom left corner of the slide, there is a decorative graphic of vertical bars of varying heights and colors, ranging from purple to yellow.

HPC means being productive on remote machines



- ✓ Linux
- ✓ OS/X
- ✓ Windows
- ✓ Multiple hop SSH
- ✓ RSA + Cryptocard
- ✓ Uses server license



allinea

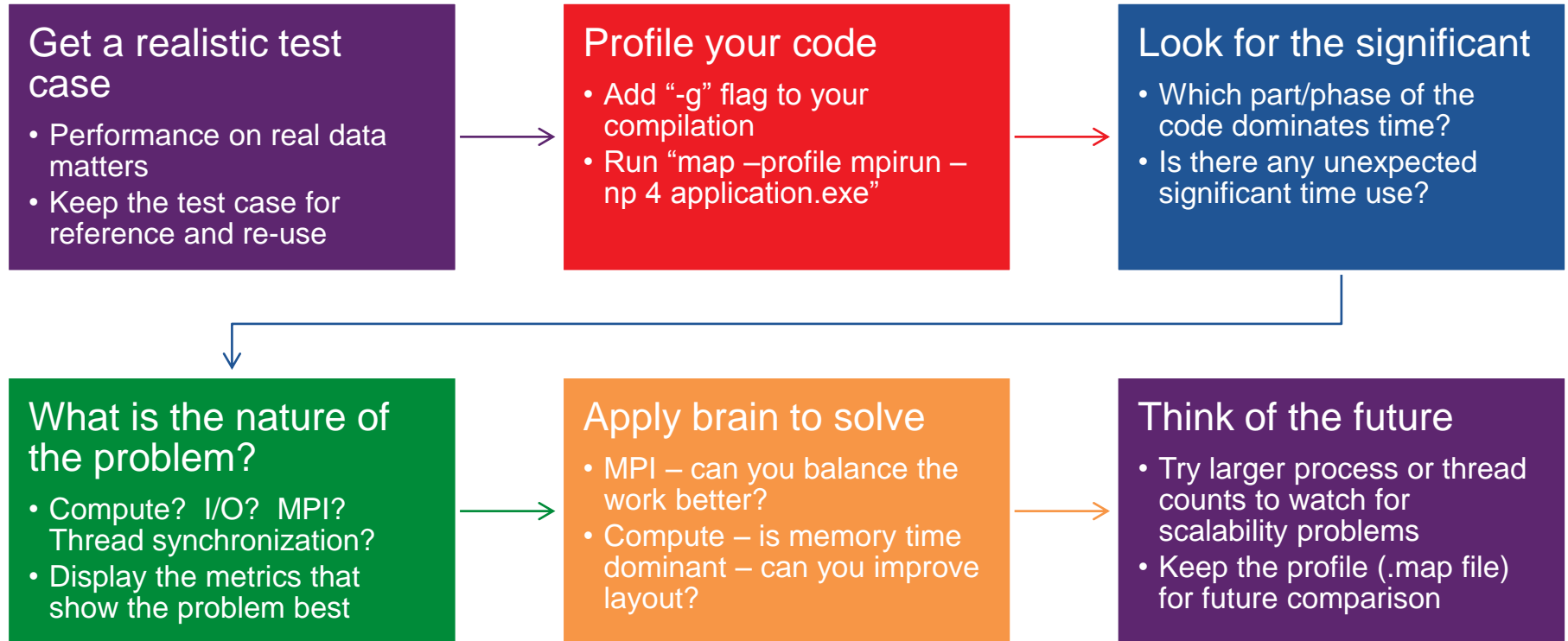
Profiling for performance

- Code optimisation can be time-consuming...



– (image courtesy of xkcd.com)

6 steps to improve performance



MAP in a nutshell



Small data files



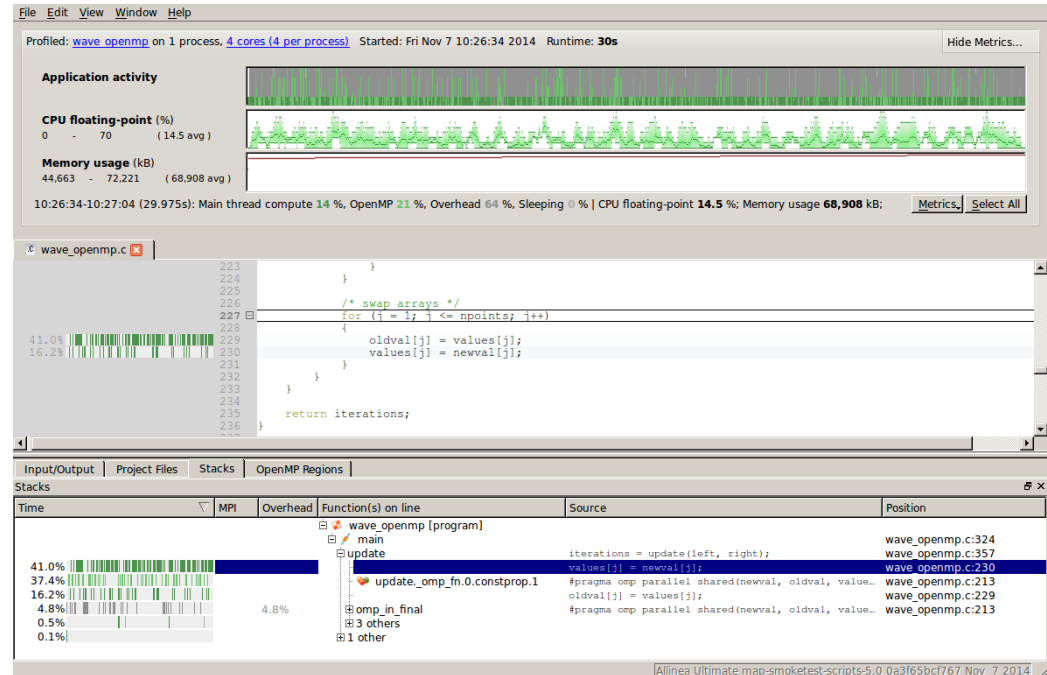
<5% slowdown



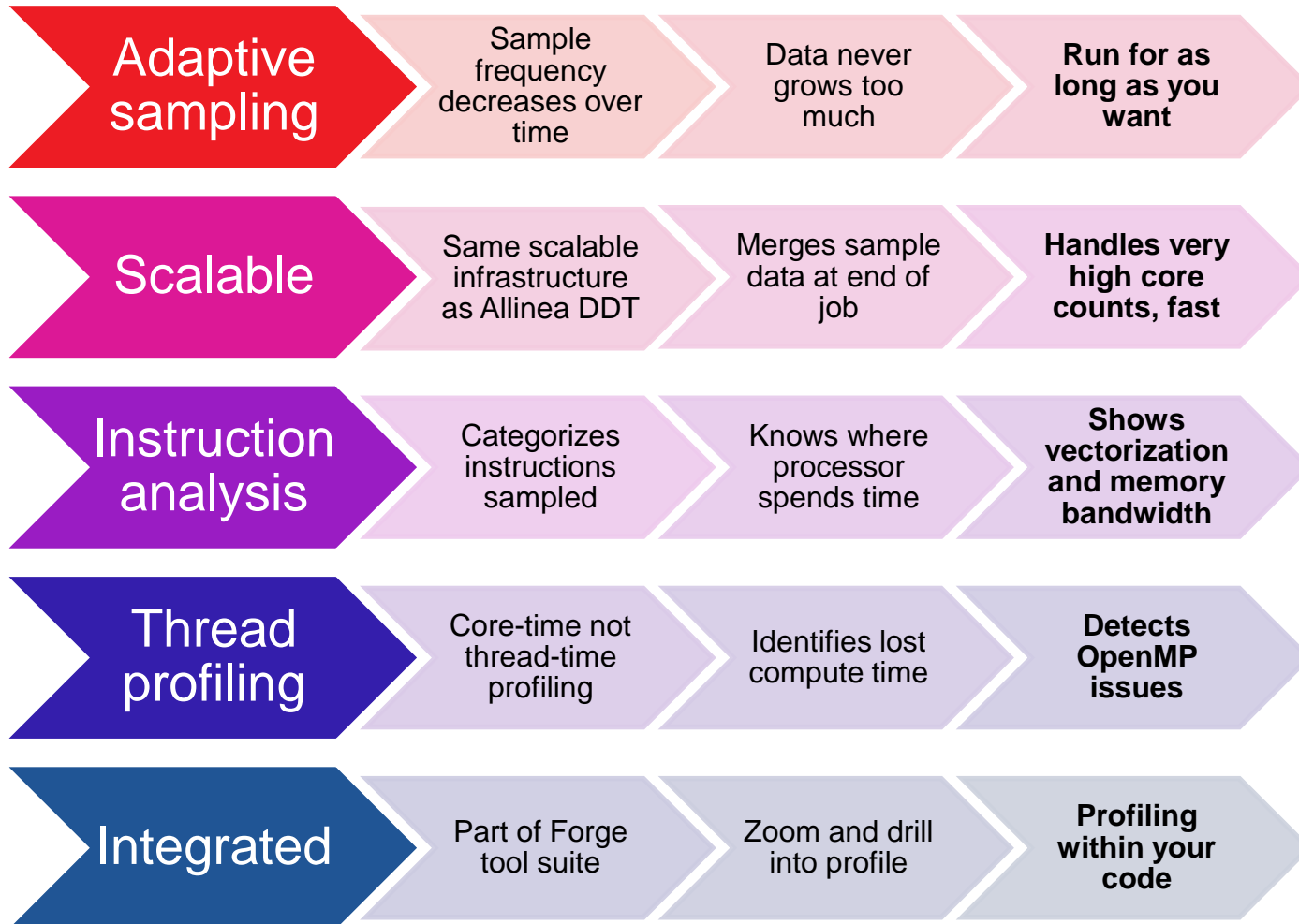
No instrumentation



No recompilation



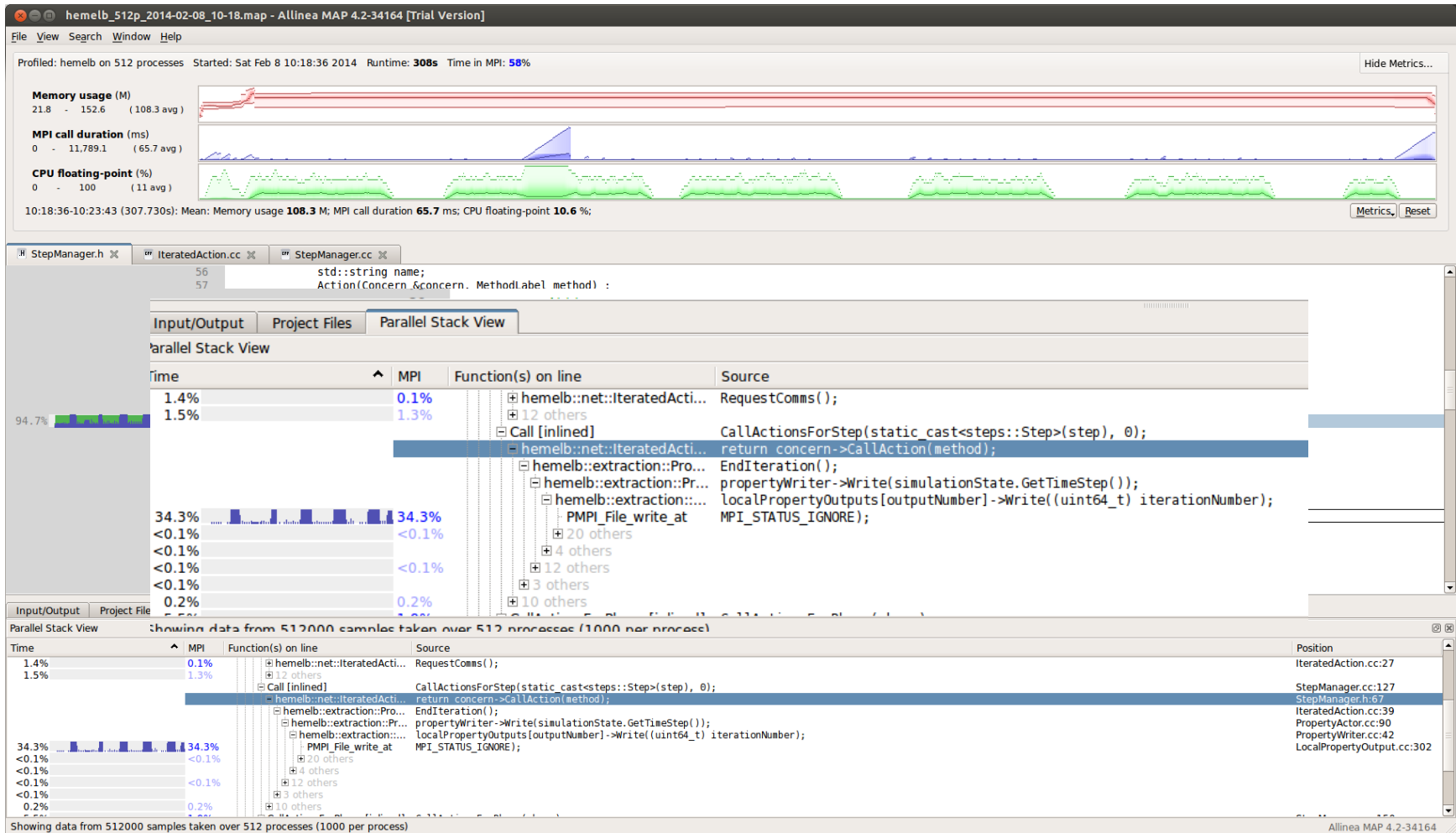
How Allinea MAP is different



Above all...

- Aimed at any performance problem that matters
 - MAP focuses on time
- Does not prejudge the problem
 - Doesn't assume it's MPI messages, threads or I/O
- If there's a problem..
 - MAP shows you it, next to your code

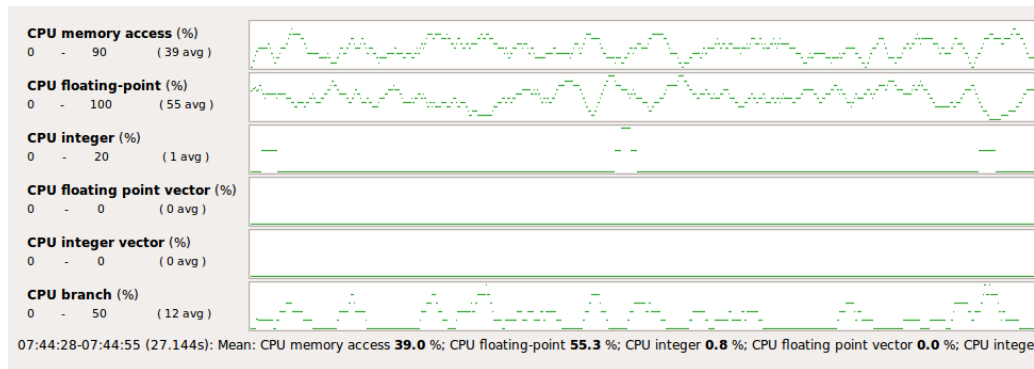
Scaling issue – 512 processes



Simple fix... reduce periodicity of output

Deeper insight into CPU usage

- Runtime of application still unusually slow



- Allinea MAP identifies vectorization close to zero
- Why? Time to switch to a debugger!

Some types of bug

Bohrbug

Steady, dependable bug

Heisenbug

Vanishes when you try to debug (observe)

Mandelbug

Complexity and obscurity of the cause is so great that it appears chaotic

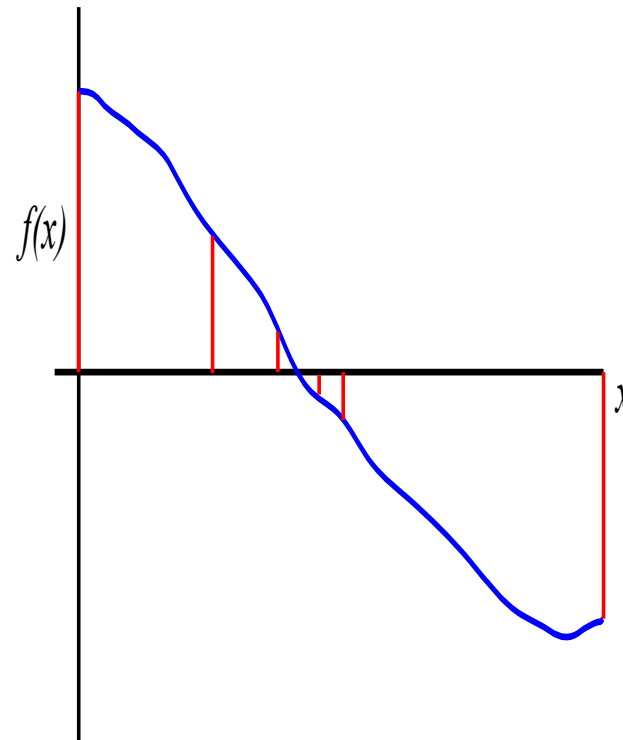
Schroedinbug

First occurs after someone reads the source file and deduces that it never worked, after which the program ceases to work

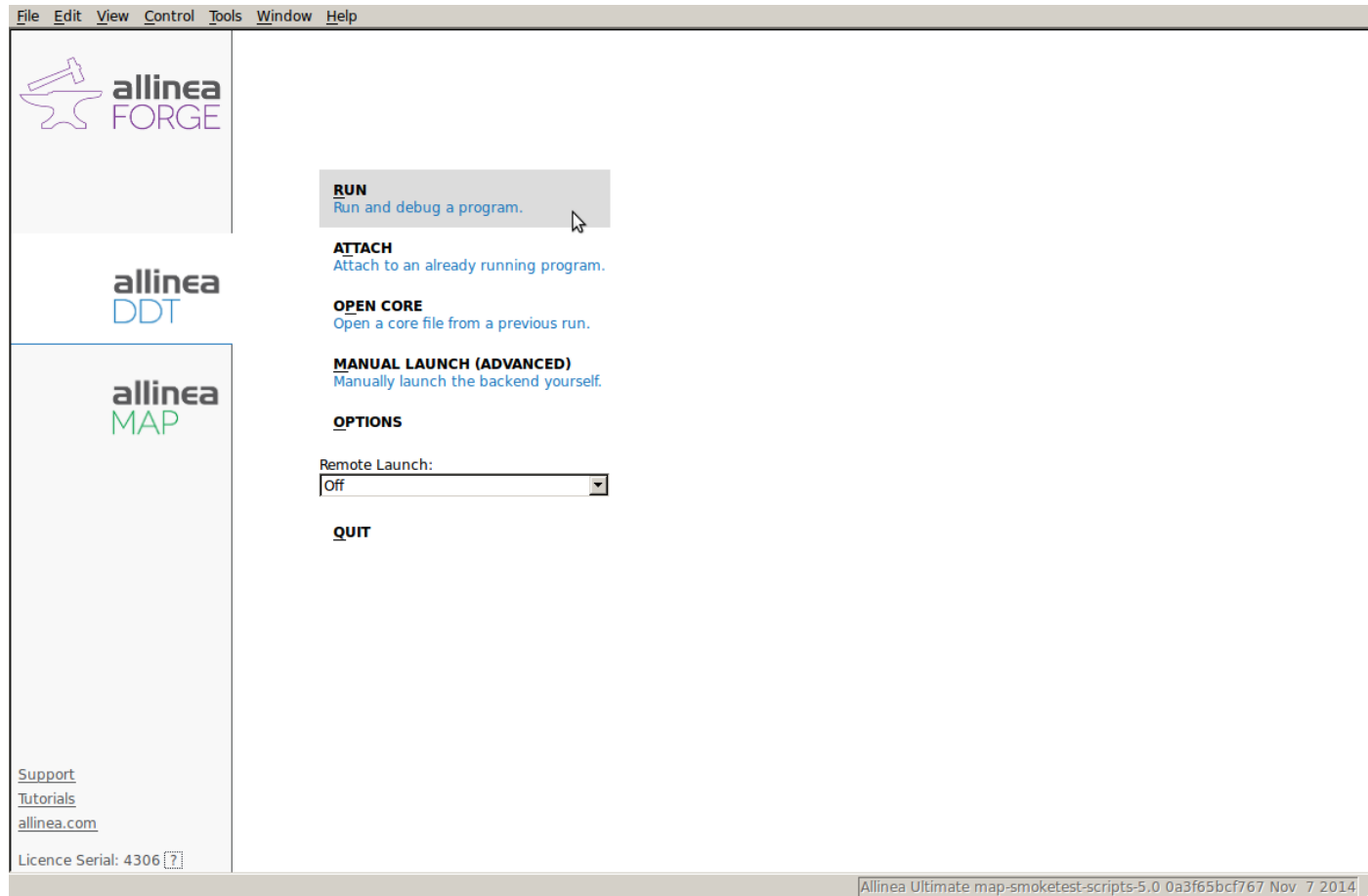
Debugging

- Transforming a broken program to a working one
- How?
 - Track the problem
 - Reproduce
 - Automate - (and simplify) the test case
 - Find origins – where could the “infection” be from?
 - Focus – examine the origins
 - Isolate – narrow down the origins
 - Correct – fix and verify the testcase is successful
- Suggested Reading:
 - Zeller A., “Why Programs Fail”, 2nd Edition, 2009
 - Zen and the Art of Motorcycle Maintenance, Robert M. Pirsig

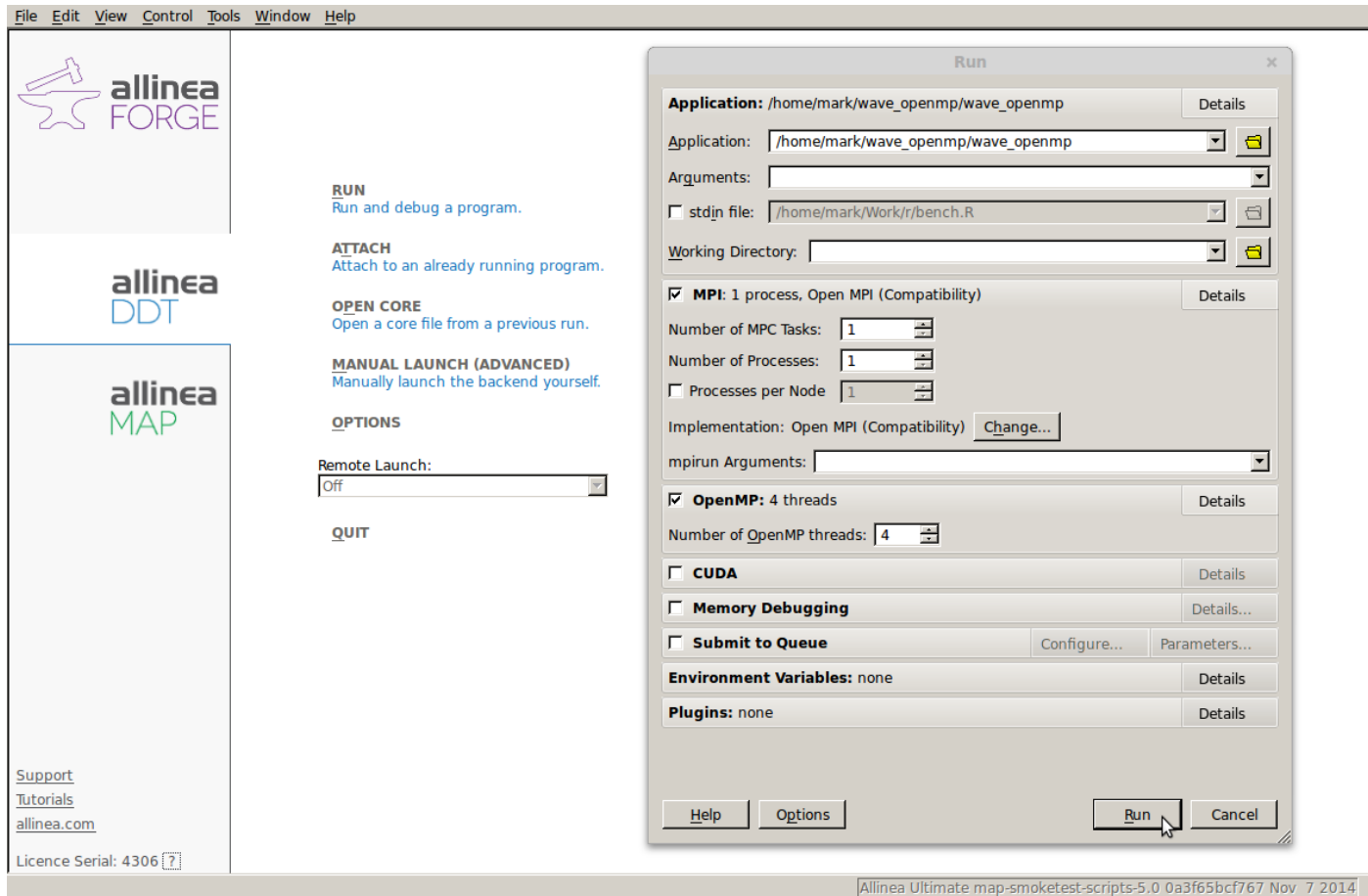
- The first debugger: print statements
 - Each process prints a message or value at defined locations
 - Diagnose the problem from evidence and intuition
- A long slow process
 - Analogous to bisection root finding
- Broken at modest scale
 - Too much output – too many log files



While still connected to the server we switch to the debugger

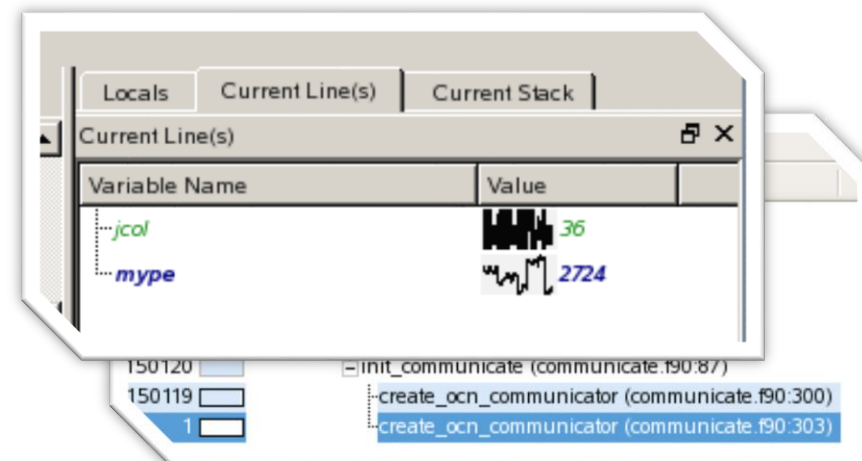
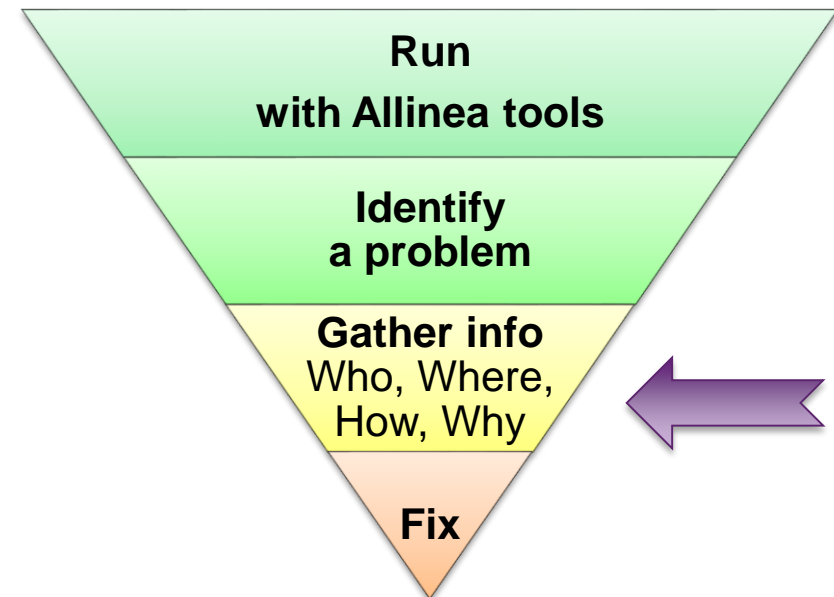


It's already configured to reproduce the profiling run



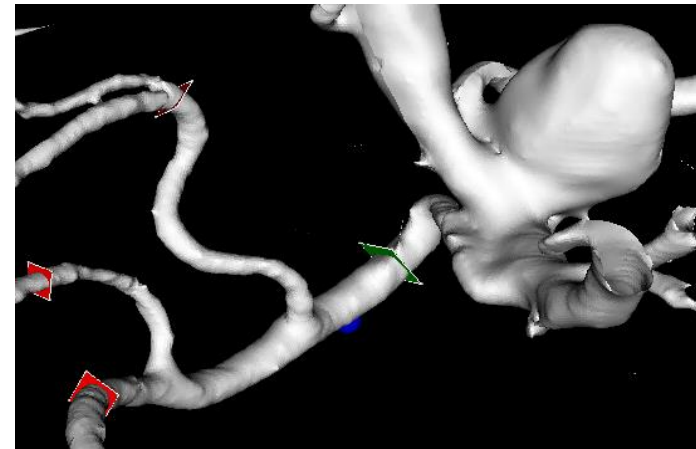
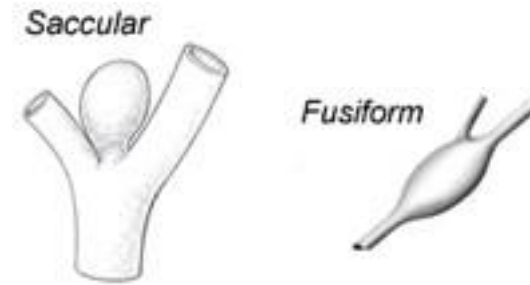
Solving Software Defects

- Who had a rogue behavior ?
 - Merges stacks from processes and threads
- Where did it happen?
 - leaps to source
- How did it happen?
 - Diagnostic messages
 - Some faults evident instantly from source
- Why did it happen?
 - Unique “Smart Highlighting”
 - Sparklines comparing data across processes



HPC could be brain surgery

- Brain aneurysms
 - 2-5% of population – most are undiagnosed
 - 30,000 rupture in US each year – 40% fatal
 - Early discovery and treatment increases survival rates
- Neurosurgery as HPC
 - MRI provides the blood vessel structure
 - Intra-cranial blood flow and pressures is just complex CFD
 - Full brain 3D model is 2-10GB geometry
- Individualized HPC
 - Patient's MRI scan enables surgical decision: whether to operate, how to operate, ...
 - Circle of Willis requires super-Petascale ~~machine~~ software
 - Need answer in minutes or hours
- ... but it crashes at 49152 cores



Allinea DDT 4.2.1-36484

File View Control Search Tools Window Help

Current Group: All Focus on current: ☒ Group ☐ Process ☐ Thread ☐ Step Threads Together

24576 processes (0-24575) Paused: 17223 Playing: 7353 Finished: 0

Currently selected: 260 (on nid00194, pid 9481, main thread IWP 9481)

Create Group

Project Files

Search (Ctrl+K)

VolumeTrav...
wave.c
weird.c
WholeGeom...
Writer.cc
wspace.c
XdrFileWrite...
XdrMemRea...
XdrMemWrit...
XdrReader.c...
XdrWriter.cc
XmlAbstract...
xyzpart.c
External Code

MpiEnvironment.cc xyzpart.c

```
551 ikvsortii(ntsamples, allpicks);
552
553
554 /* Select the final splitters. Set the boundaries to
555 for (i=1; i<npes; i++)
556     mypicks[i] = allpicks[i*ntsamples/npes];
557 mypicks[0].key = IDX_MIN;
558 mypicks[npes].key = IDX_MAX;
559
560
561 WCOREPOP; /* free allpicks */
562
563 STOPTIMER(ctrl, ctrl->AuxTmr2);
564 STARTTIMER(ctrl, ctrl->AuxTmr3);
565
```

Locals Current Line(s) Current Stack

Current Line(s)

Variable Name	Value
allpicks	0x2aab8055e010
i	2245
mypicks	0x2a6f8f0
npes	24575
ntsamples	1818550

Type: none selected

Input/Output Breakpoints Watchpoints Stacks Tracepoints Tracepoint Output Logbook

Stacks

Processes	Threads	Function
17223	17223	main (main.cc:37)
17223	17223	SimulationMaster::SimulationMaster (SimulationMaster.cc:63)
17223	17223	SimulationMaster::Initialise (SimulationMaster.cc:154)
17223	17223	hemelb::geometry::GeometryReader::LoadAndDecompose (GeometryReader.cc:188)
17223	17223	hemelb::geometry::GeometryReader::OptimiseDomainDecomposition (GeometryReader.cc:188)
17223	17223	hemelb::geometry::decomposition::OptimisedDecomposition::OptimisedDecomposition (OptimisedDecomposition.cc:188)
17223	17223	hemelb::geometry::decomposition::OptimisedDecomposition::CallParmetis (OptimisedDecomposition.cc:188)
17223	17223	ParMETIS_V3_PartGeomKway (gkmetis.c:90)
17223	17223	libparmetis_Coordinate_Partition (xyzpart.c:58)
17223	17223	libparmetis_PseudoSampleSort (xyzpart.c:556)

Evaluate

Expression	Value
i * ntsamples	-212322546

Type: int
Range: from -2147259746 to -12282046
49/17223 processes equal

7353 processes playing

Computer titan-ext7 Allinea DDT 4.2.1-36484 Sun Aug 10, 7:50 PM

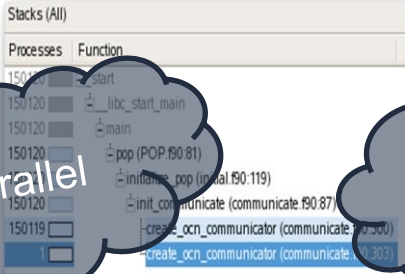
Debugging by Inspiration

- Some errors are harder than others to diagnose
 - A bug that occurs 1% of the time is many times harder to fix than one that occurs 100% of the time
 - Often caused by incorrect memory usage
- Get help from Allinea DDT's memory debugging
 - Checks double frees
 - Checks use of dangling (freed) pointers
 - Can force O/S to check for read/write beyond bounds
 - Make some random bugs deterministic and occur 100% of the time ("Guard Pages")
 - showing memory leaks (group by allocation point)
- Other tips – get inspiration from a colleague
 - try explaining your code
 - The very act of explaining your thinking in the code can
 - don't have a colleague to hand?
 - Follow <http://www.rubberduckdebugging.com>



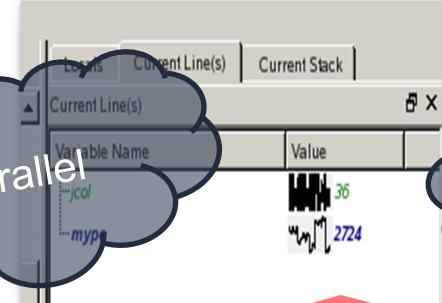
Favorite Allinea DDT Features for Scale

Parallel



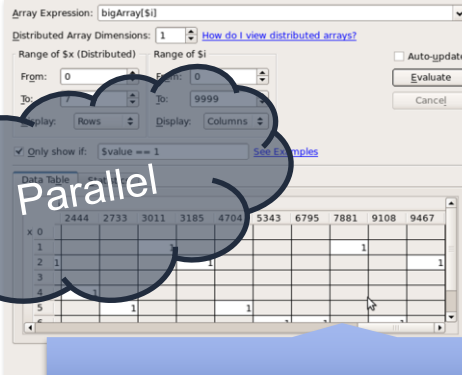
Parallel stack view

Parallel




Automated data comparison: sparklines

Parallel



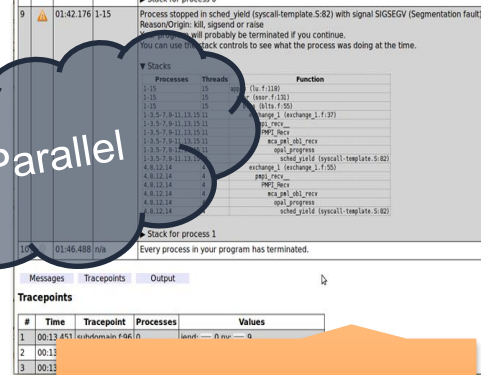
Parallel array searching

Parallel



Step, play, and breakpoints

Parallel



Offline debugging

Today's Status on Scalability

- Debugging and profiling
 - Active users at 100,000+ cores debugging
 - 50,000 cores is largest profiling tried to date (and was Very successful)
 - ... and active users with just 1 process too
- Deployed on
 - ORNL's Titan, NCSA Blue Waters, ANL Mira etc.
 - Hundreds of much smaller systems – academic, research, oil and gas, genomics, etc.
- Tools help the full range of programmer ambition
 - Very small slow down with either tool (< 5%)

The scalable print alternative

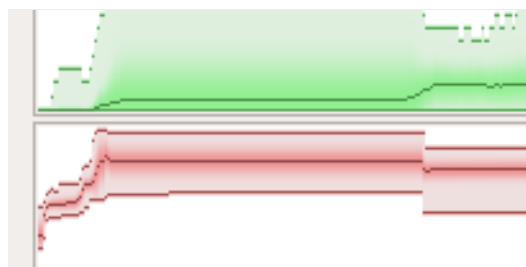
Stop on variable change

Static analysis warnings on code errors

Detect read/write beyond array bounds

Detect stale memory allocations

Six Great Things to Try with Allinea MAP



pute 76 % MPI 24 % File

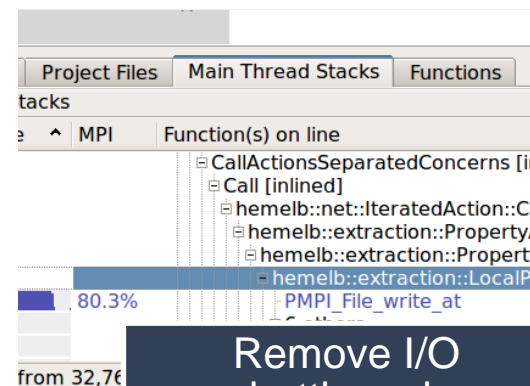
Find the peak memory use

```

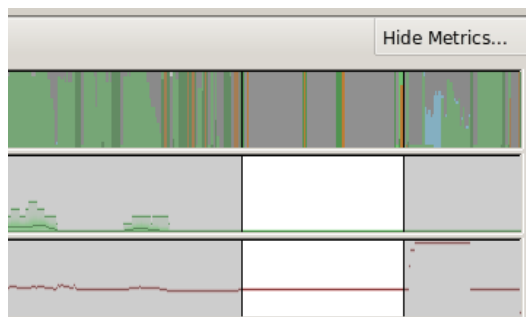
30      ! late to the party
31      do j=1,20*nprocs; a
32      end if
33
34      if (pe /= 0) then
35      call MPI_SEND(a, si
36      else
37      do from=1,nprocs-1
38      call MPI_RECV(b,
39      do j=1,50; b=sqrt
40      print *, "Answer f
41      end do
42      end if
43      end do
44      call MPI_BARRIER(MPI CO

```

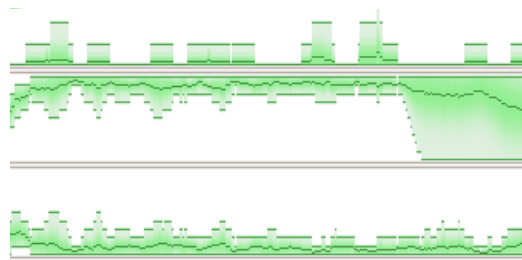
Fix an MPI imbalance



Remove I/O bottleneck



Make sure OpenMP regions make sense



Improve memory access

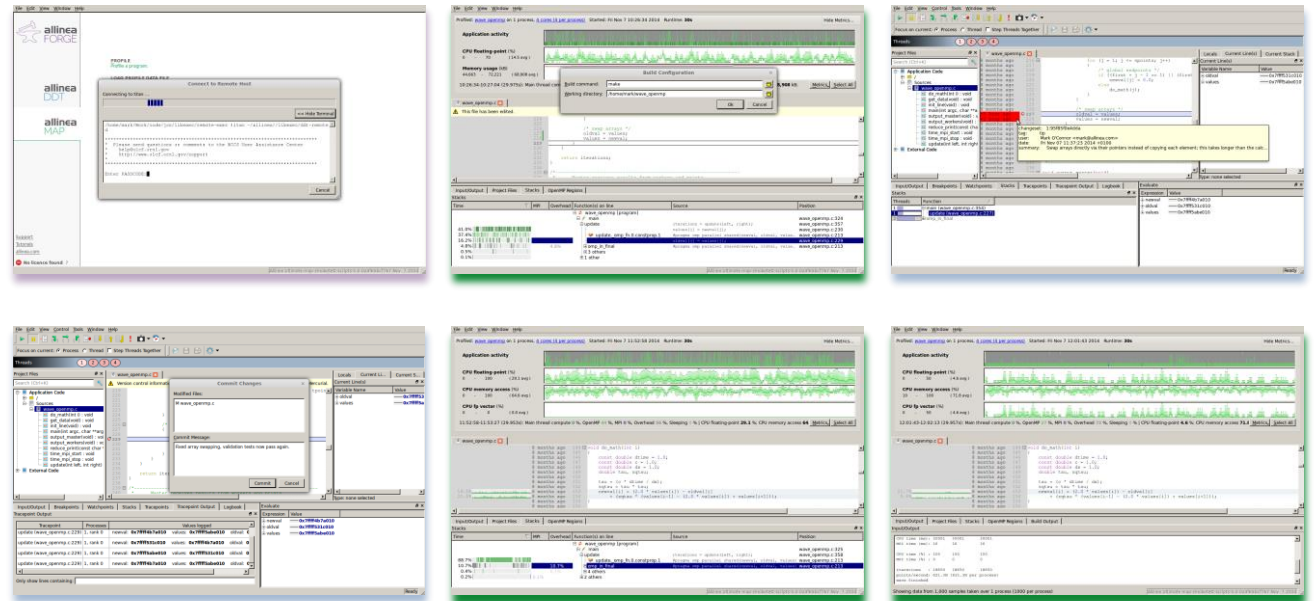
```

size, nproc, mat a
A[i*size+k]*B[k*s

```

Restructure for vectorization

A Productive HPC Development Workflow



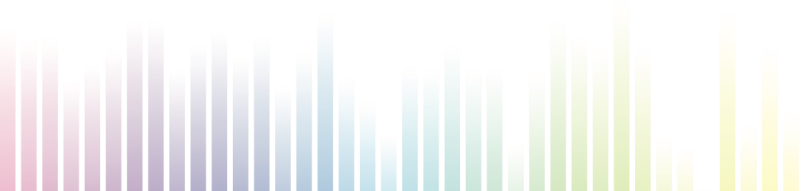
allinea

Tour/Demo of Forge

- Debugging with DDT
 - Walkthrough example – cstartmpi.exe
- Profiling with MAP
 - Matrix-multiplication

Hands on Session

- Use Allinea DDT on your favorite system to debug your code – or example codes
- Use Allinea MAP on NERSC Edison or ANL Cooley to see your code performance
- Use Allinea DDT and Allinea MAP together to improve our test code
- Can you beat a 50% speed up?



Getting started on Mira/Cooley

- Install local client on your laptop
 - www.allinea.com/products/forge/downloads
 - Linux – installs full set of tools
 - Windows, Mac – just a remote client to the remote system
 - Run the installation and software
 - “Connect to remote host”
 - Hostname:
 - username@cetus.alcf.anl.gov
 - username@cooley.alcf.anl.gov
 - Remote installation directory: /soft/debuggers/ddt
 - Click Test
- Congratulations you are now ready to debug on Mira/Vesta/Cetus – or debug and profile on Cooley.

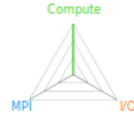
Performance and debugging challenge

- Obtain the code:
 - `git clone https://github.com/estrabd/2d-heat.git`
- The challenge
 - Improve run time by 50%
 - Parameters: `-h 1000 -w 1000` set the dimensions
- Submit/run
 - `mpirun -np 16 2d-heat.x -h 1000 -w 1000`
- Use MAP to find the lowest hanging fruits
- Use DDT to look at behaviour - confirm what is safe to change
- Edit and build in the GUI: File/Configure Build/ and re-run

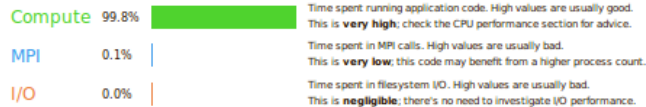
Initial Performance Report



Command: oshrun -n 4 -mca sshmem mmap
./shmem_2dheat.x -h 1000 -w 1000
Resources: 1 node (2 physical, 4 logical cores per node)
Memory: 7 GB per node
Tasks: 4 processes
Machine: ip-172-30-0-163
Start time: Mon Aug 3 10:58:52 2015
Total time: 42 seconds (1 minute)
Full path: /home/training/openshmem-release-1.0d/examples
Input file:
Notes:



Summary: shmem_2dheat.x is **Compute-bound** in this configuration

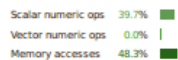


This application run was **Compute-bound**. A breakdown of this time and advice for investigating further is in the **CPU** section below.

As very little time is spent in **MPI** calls, this code may also benefit from running at larger scales.

CPU

A breakdown of the 99.8% CPU time:

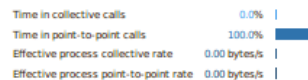


The per-core performance is memory-bound. Use a profiler to identify time-consuming loops and check their cache performance.

No time is spent in **vectorized instructions**. Check the compiler's vectorization advice to see why key loops could not be vectorized.

MPI

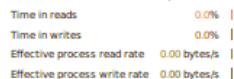
A breakdown of the 0.1% MPI time:



Most of the time is spent in **point-to-point calls** with a very low transfer rate. This suggests load imbalance is causing synchronization overhead; use an MPI profiler to investigate.

I/O

A breakdown of the 0.0% I/O time:



No time is spent in **I/O** operations. There's nothing to optimize here!

Threads

A breakdown of how multiple threads were used:

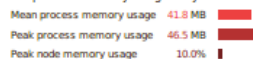


No measurable time is spent in **multithreaded code**.

The system load is high. Check that other jobs or system processes are not running on the same nodes.

Memory

Per-process memory usage may also affect scaling:



The peak node memory usage is very low. Running with fewer MPI processes and more data on each process may be more efficient.

- Using Forge – profile, debug, edit, build, repeat!
 - profiler to find targets
 - debugger to reveal usage
 - Make changes
 - Configure the Build command
 - make CFLAGS='-O3 -g'
 - **NB: single quotes..**
 - Build, repeat!