

# allinea

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

## Large scale heterogeneous applications made easy with Alinea

ECMWF, 27/10/2016

Florent Lebeau  
[flebeau@allinea.com](mailto:flebeau@allinea.com)



# Weather and Forecasting models today



Scalability

Efficiency

Simplicity

**allinea**

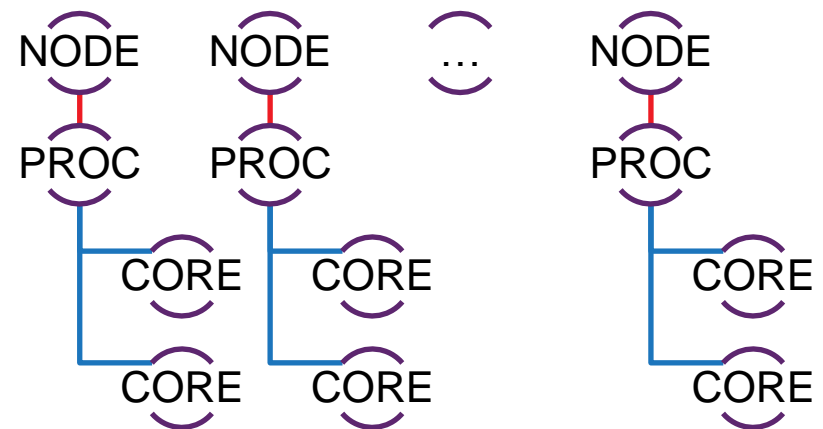
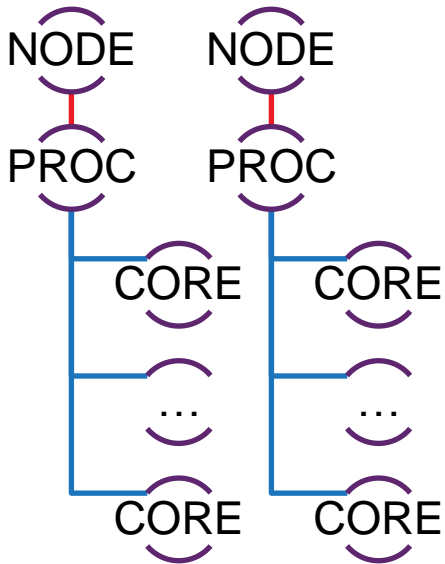
# Towards more vertical & horizontal scalability

INTEL  
KNIGHTS  
LANDING

NVIDIA  
GPUS

NEXT-GEN  
INTEL  
XEON

ARM v8



# Heterogeneous Systems

- Use many processor architectures
  - x86\_64 + GPUs
  - x86\_64 + KNL
  - OpenPower + GPUs
  - ARMv8 + GPUs
  - ...
- Goal:
  - Increase compute power with specialised processors
  - Improve energy efficiency
- Parallel Programming Languages:
  - MPI
  - OpenMP
  - CUDA
  - OpenACC
  - ...

# Alinea's vision

- **Helping maximize HPC production**



- Reduce HPC systems operating costs
- Resolve cutting-edge challenges
- Promote Efficiency (as opposed to Utilization)
- Transfer knowledge to HPC communities

- **Helping the HPC community design the best applications**



- Reach highest levels of performance and scalability
- Improve scientific code quality and accuracy

# Where to find Alinea's tools

Over 65% of Top 100 HPC systems

- From small to very large tools provision

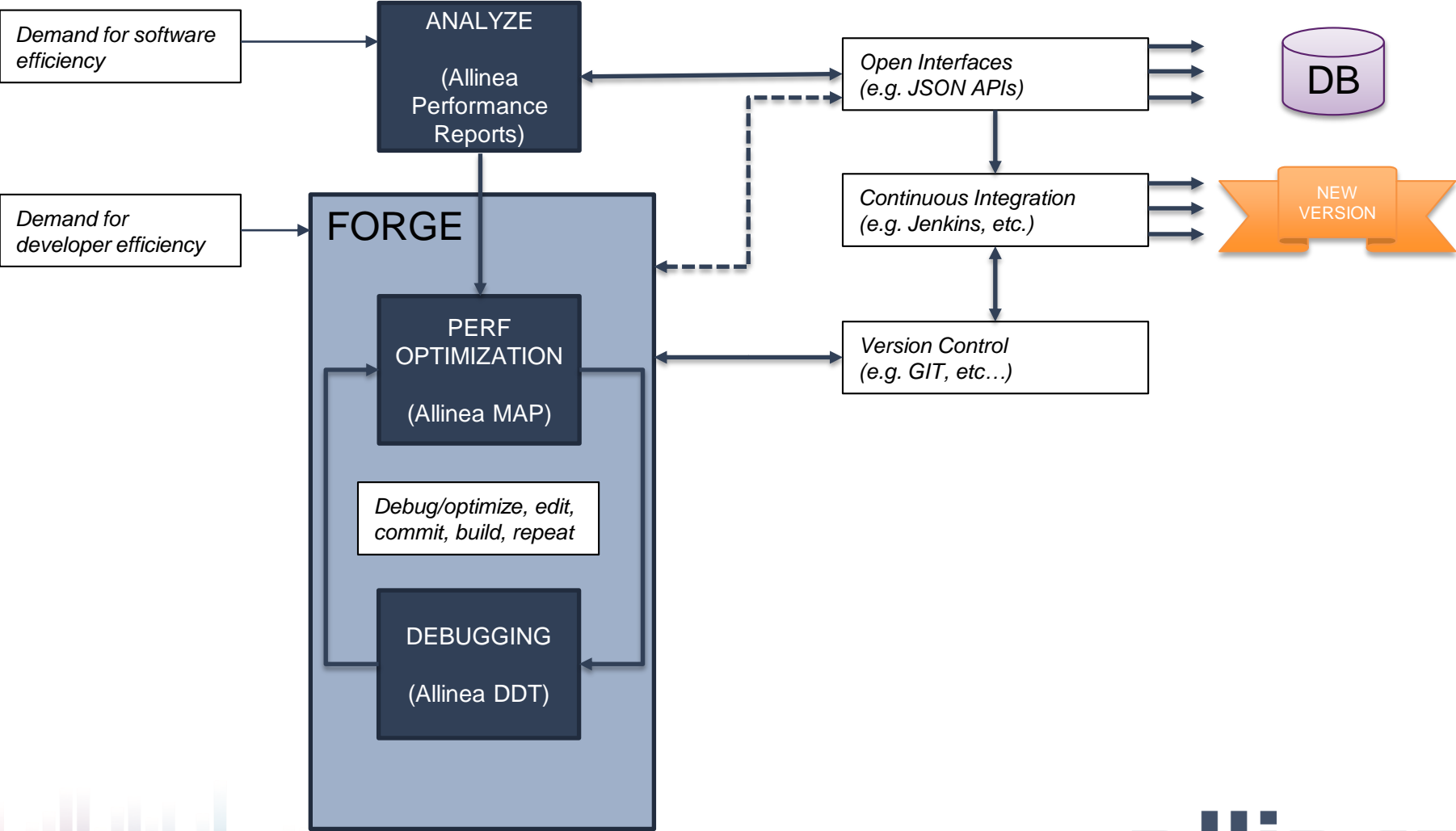
8 of the Top 10 HPC systems

- From 1,000 to 700,000 core tools usage

Future leadership systems

- Millions of cores usage

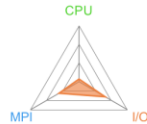
# Development process workflow



# Analyse with Alinea Performance Reports

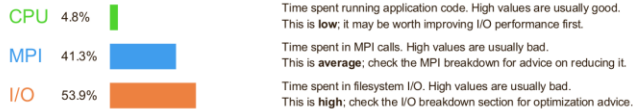


Executable: MADbench2  
Resources: 16 processes, 1 node  
Machine: sandybridge2  
Start time: Mon Nov 4 12:27:50 2013  
Total time: 109 seconds (2 minutes)  
Full path: /tmp/MADbench2  
Notes: 12-core server / HDD / 16 readers + writers



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 **4.8%** 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

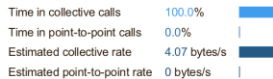
A breakdown of how the **53.9%** total I/O time was spent:



Most of the time is spent in **write operations**, which have a very 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 **41.3%** total time spent in MPI calls:



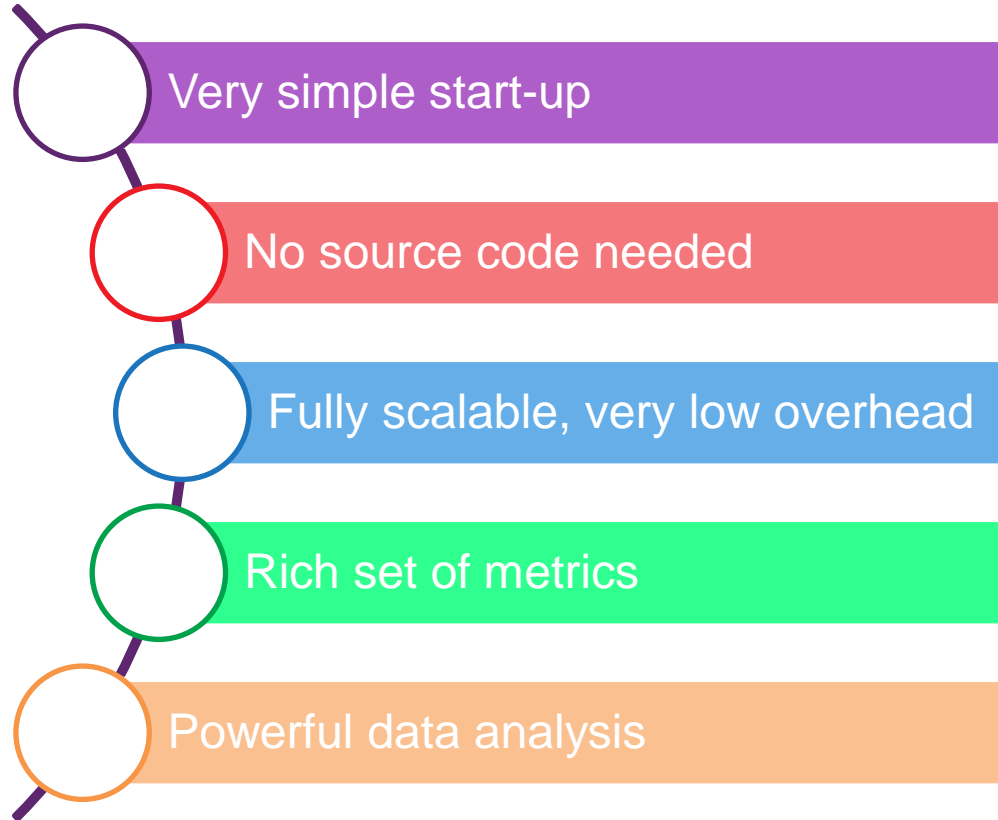
All of the time is spent in **collective calls** with a very low transfer rate. This suggests a significant load imbalance is causing synchronization overhead. You can investigate this further with an MPI profiler.

## 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.





# Increase the efficiency of your jobs

## Accelerators

A breakdown of how accelerators were used:

GPU utilization	78.3%	<div style="width: 78.3%;"></div>
Global memory accesses	70.9%	<div style="width: 70.9%;"></div>
Mean GPU memory usage	31.5%	<div style="width: 31.5%;"></div>
Peak GPU memory usage	38.7%	<div style="width: 38.7%;"></div>

Significant time is spent in **global memory accesses**. Try modifying kernels to use shared memory instead and check for bad striding patterns.

**GPU utilization** is acceptable. Use the NVIDIA Visual Profiler to improve kernel performance.

## Energy

A breakdown of how the 3.6 Wh was used:

CPU	62.9%	<div style="width: 62.9%;"></div>
System	37.1%	<div style="width: 37.1%;"></div>
Mean node power	92.4 W	<div style="width: 92.4%;"></div>
Peak node power	94 W	<div style="width: 94%;"></div>

Significant energy is wasted during MPI communications. It may be more efficient to use fewer nodes with more data on each node.

Significant time is spent waiting for memory accesses. Reducing the **CPU** clock frequency could reduce overall energy usage.

## OpenMP

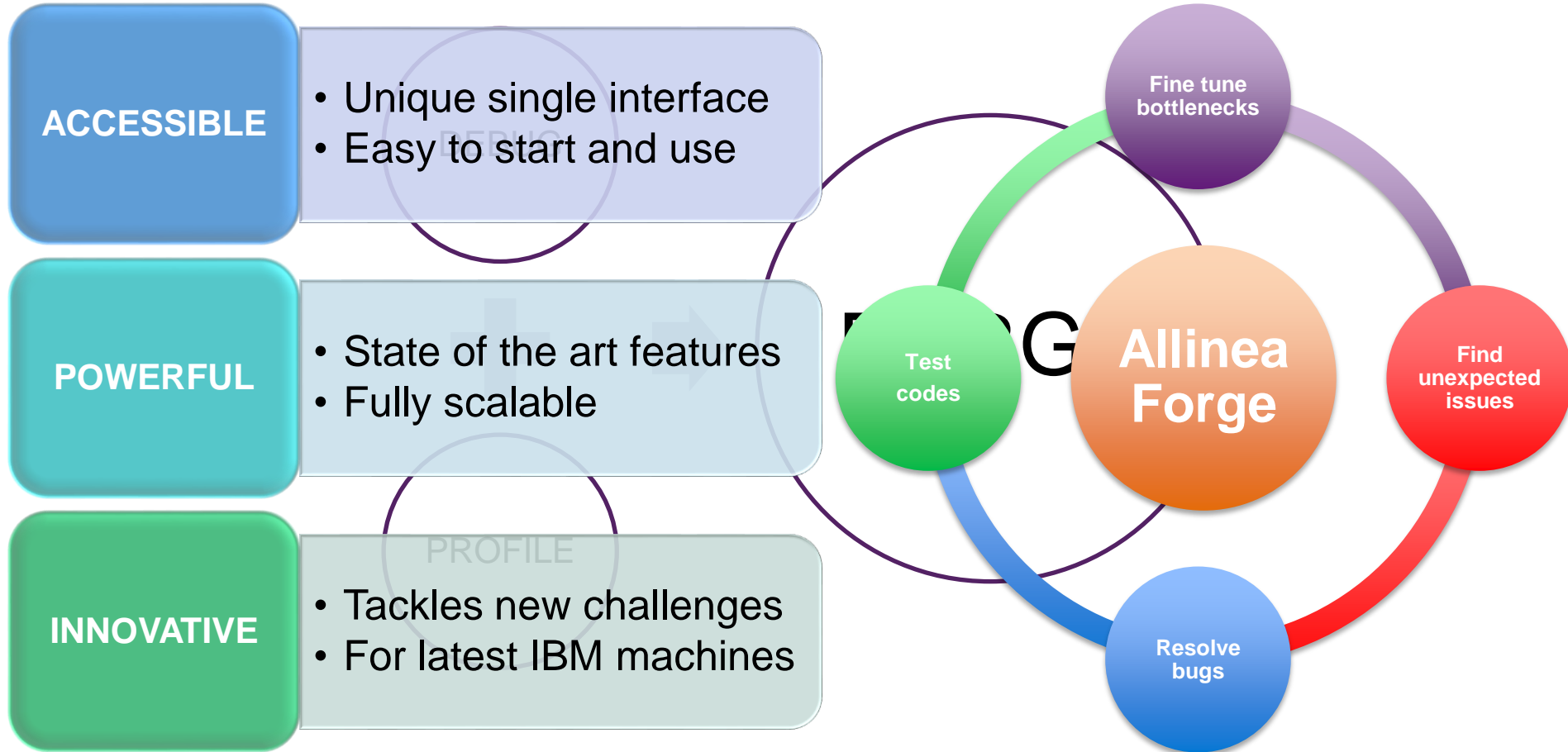
A breakdown of the 99.5% time in OpenMP regions:

Computation	58.9%	<div style="width: 58.9%;"></div>
Synchronization	41.1%	<div style="width: 41.1%;"></div>
Physical core utilization	100.0%	<div style="width: 100.0%;"></div>
System load	99.7%	<div style="width: 99.7%;"></div>

Significant time is spent **synchronizing** threads in parallel regions. Check the affected regions with a profiler.

This may be a sign of overly fine-grained parallelism (OpenMP regions in tight loops) or workload imbalance.


# Allinea Forge: the toolkit for HPC developers



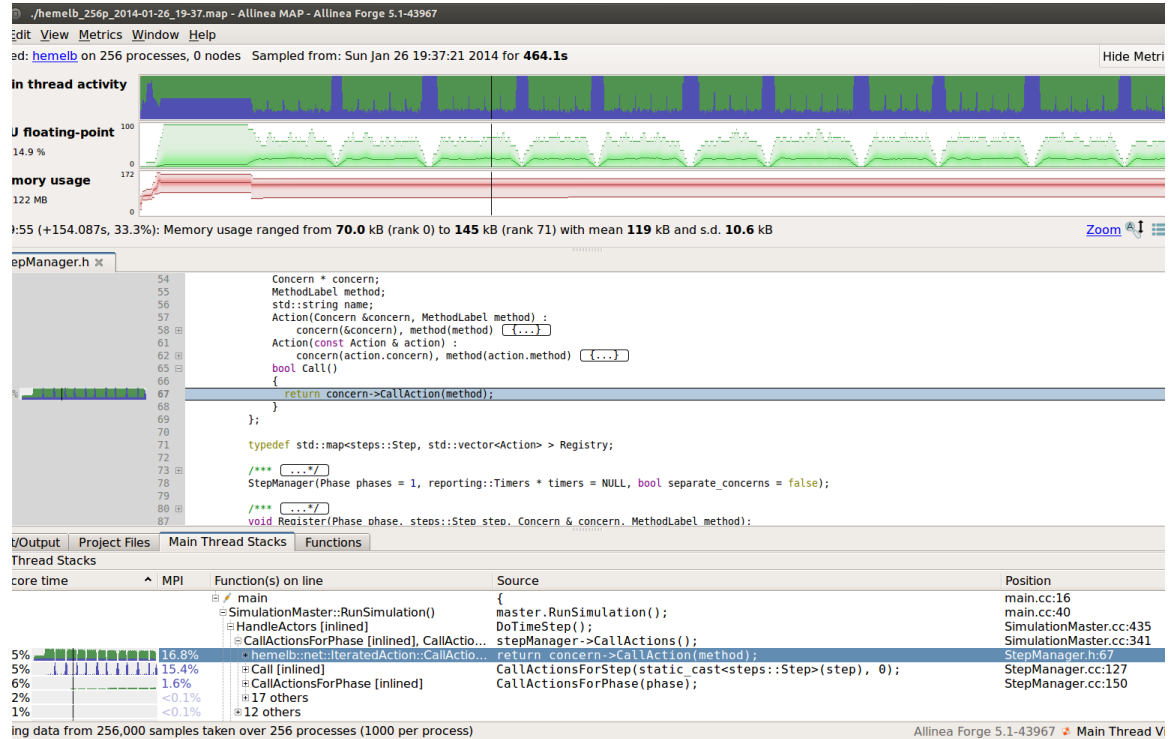
# Allinea MAP – the profiler

 No instrumentation

 Low overhead

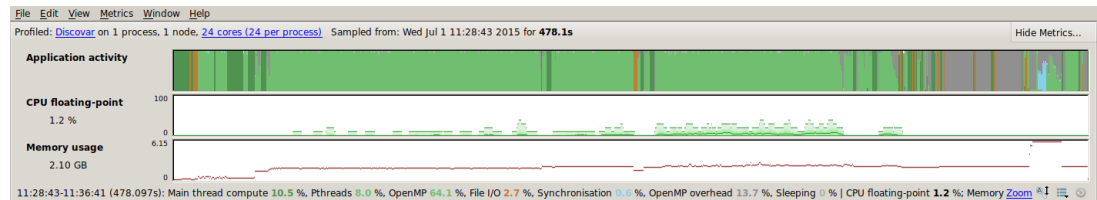
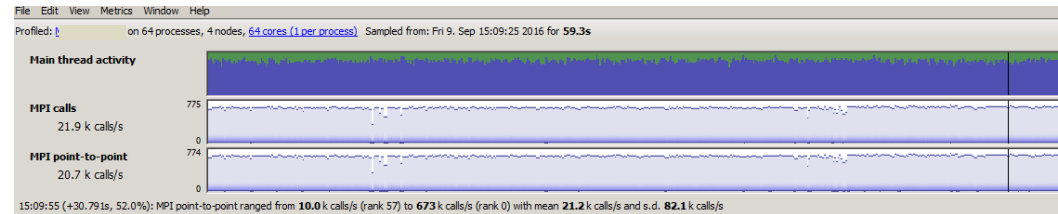
 Scalable

 Small data files

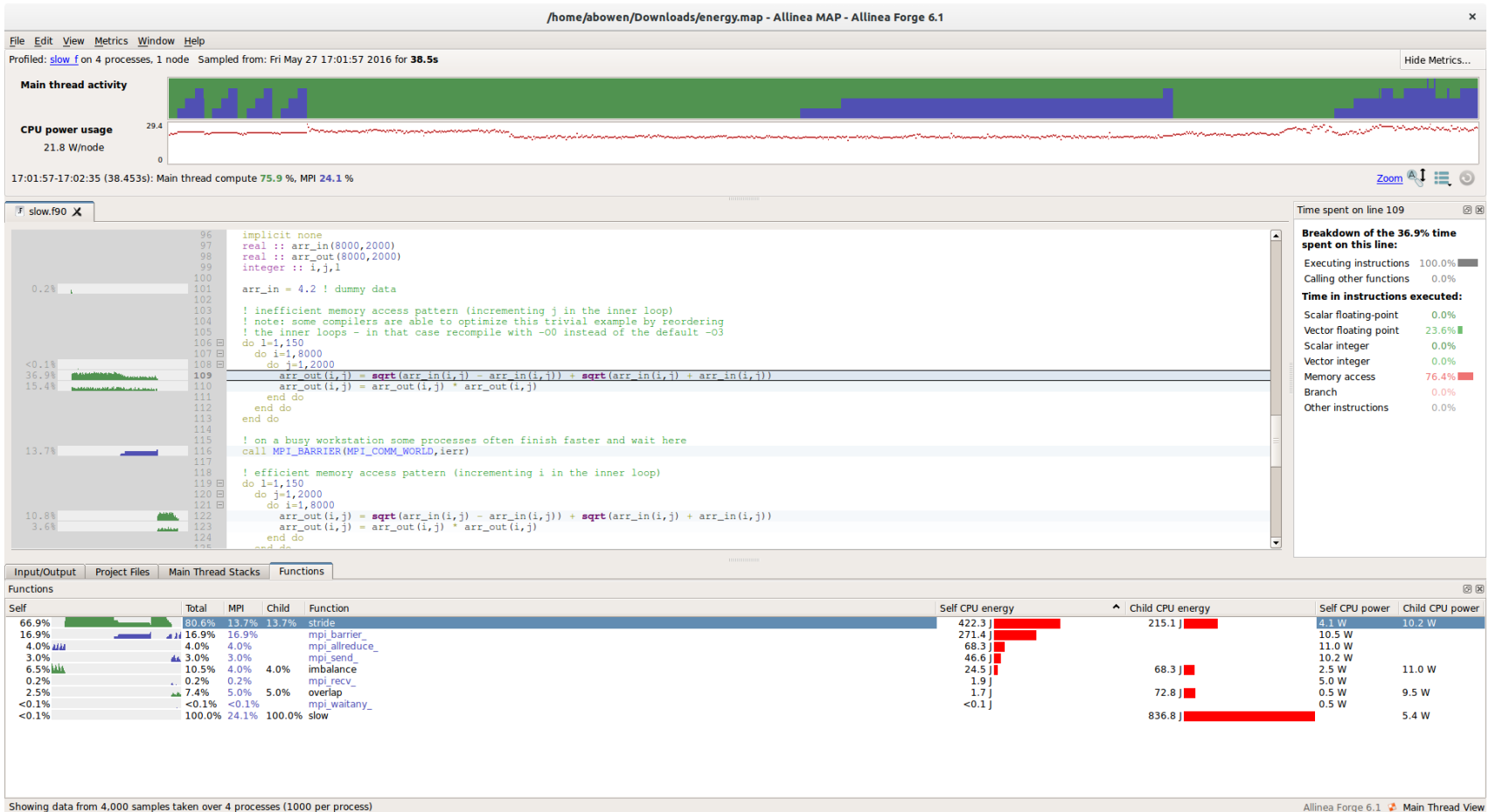


# Quickly spot application bottlenecks

- Find patterns of MPI and OpenMP imbalance
- Offload compute intensive parallel regions to be offloaded to a coprocessor or accelerator

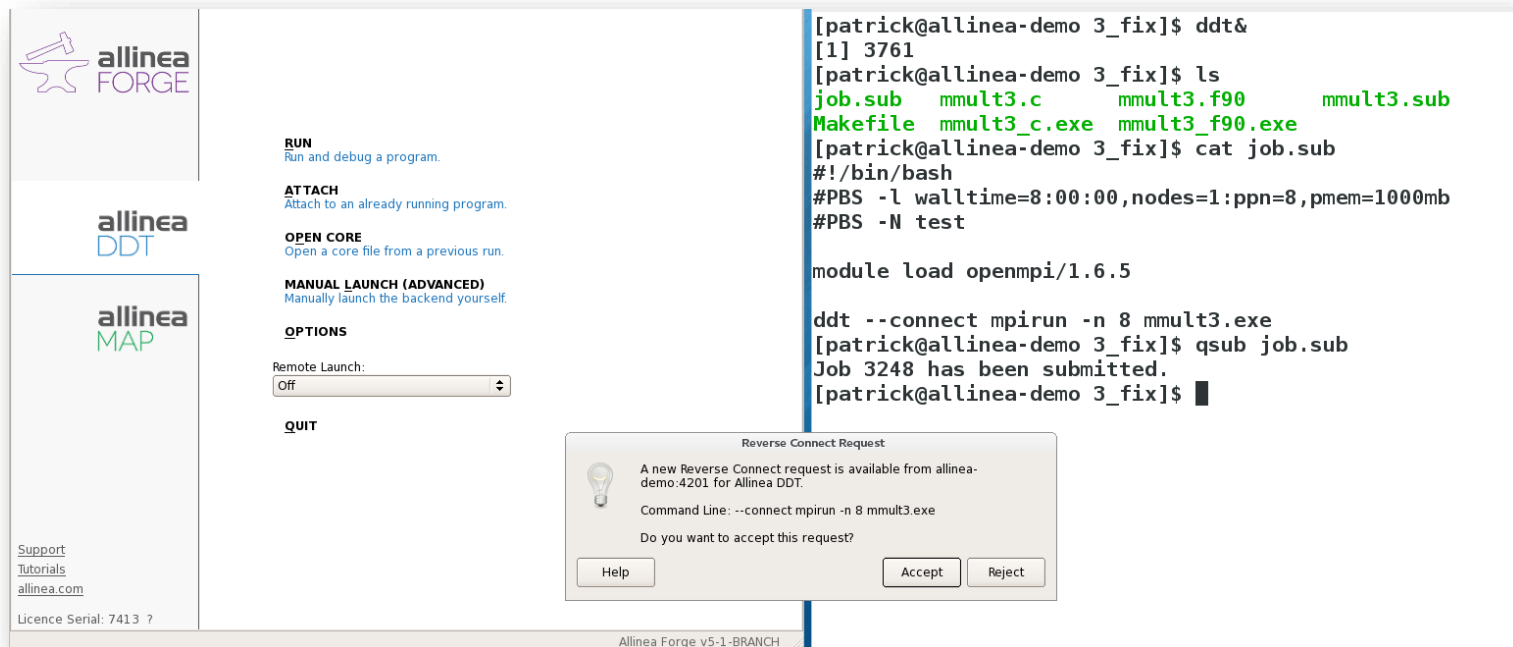


# Energy analytics



# Allinea DDT – the debugger

- Easily debug executable in your workflow with the “reverse connect” mechanism



The screenshot displays the Allinea Forge v5-1-BRANCH interface. On the left, there are logos for Allinea Forge, Allinea DDT, and Allinea MAP. The main area shows the Allinea DDT configuration options: RUN (Run and debug a program.), ATTACH (Attach to an already running program.), OPEN CORE (Open a core file from a previous run.), MANUAL LAUNCH (ADVANCED) (Manually launch the backend yourself.), OPTIONS (Remote Launch: Off), and QUIT. A terminal window on the right shows the following commands and output:

```
[patrick@allinea-demo 3_fix]$ ddt&
[1] 3761
[patrick@allinea-demo 3_fix]$ ls
job.sub  mmult3.c  mmult3.f90  mmult3.sub
Makefile mmult3_c.exe mmult3_f90.exe
[patrick@allinea-demo 3_fix]$ cat job.sub
#!/bin/bash
#PBS -l walltime=8:00:00,nodes=1:ppn=8,pmem=1000mb
#PBS -N test

module load openmpi/1.6.5

ddt --connect mpirun -n 8 mmult3.exe
[patrick@allinea-demo 3_fix]$ qsub job.sub
Job 3248 has been submitted.
[patrick@allinea-demo 3_fix]$
```

A 'Reverse Connect Request' dialog box is open, displaying a lightbulb icon and the following text:

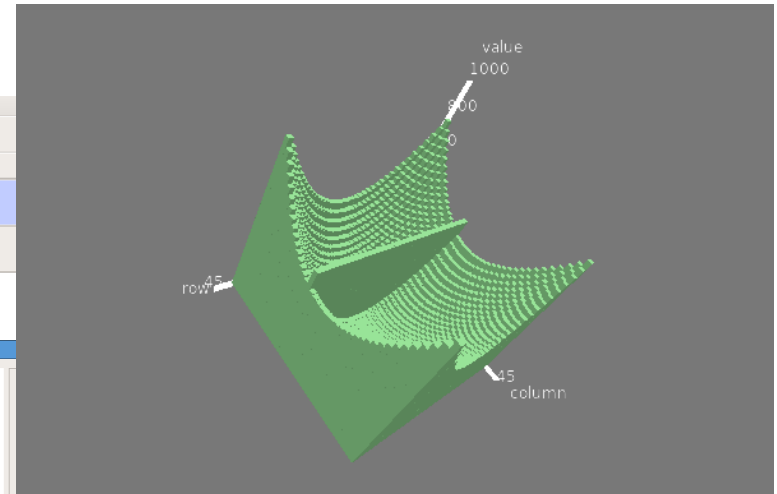
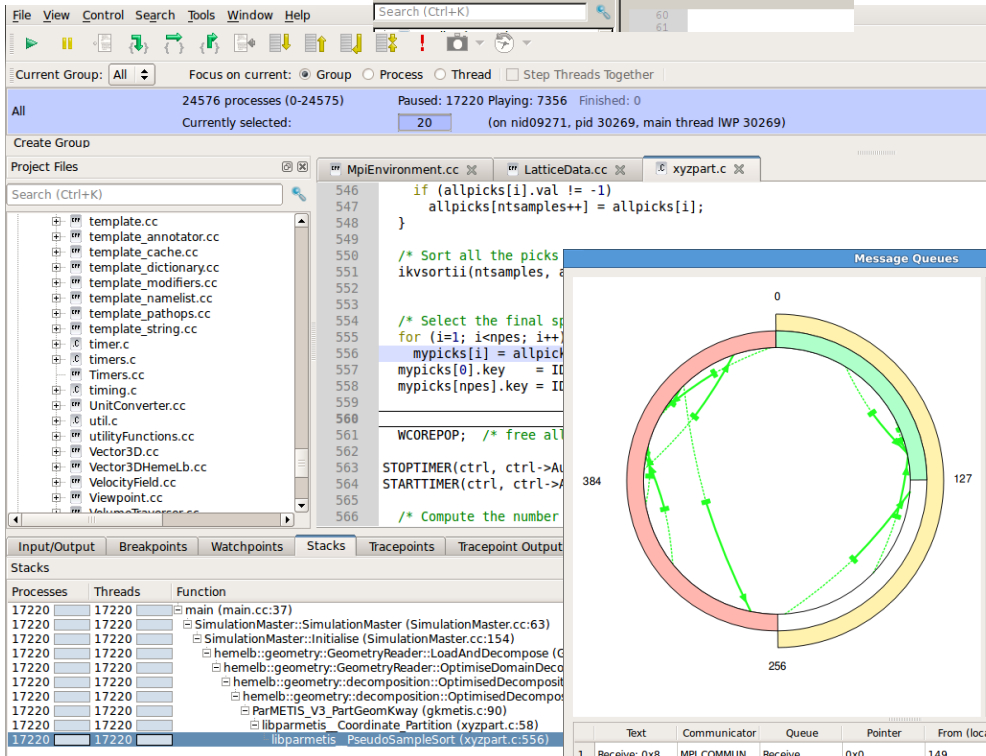
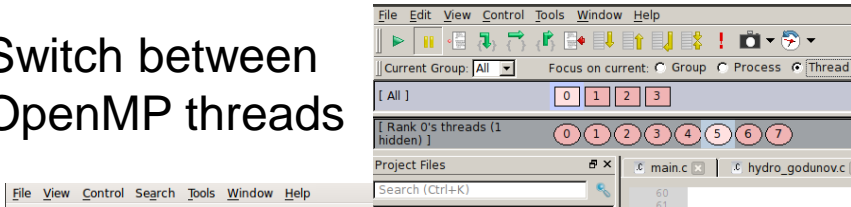
Reverse Connect Request  
A new Reverse Connect request is available from allinea-demo:4201 for Allinea DDT.  
Command Line: --connect mpirun -n 8 mmult3.exe  
Do you want to accept this request?

The dialog box has three buttons: Help, Accept, and Reject.

# Debug large scale application

Switch between OpenMP threads

Visualise data structures



Debugger message queue details and variable values. The 'Message Queues' window shows a circular diagram with nodes and arrows. Below it is a table of message details:

Text	Communicator	Queue	Pointer	From (local)	From (global)	To (local)	To (global)
1 Receive: 0x8...	MPI COMMUN...	Receive	0x0	149	405	113	369
2 Receive: 0x8...	MPI COMMUN...	Receive	0x0	16	272	193	449
3 Receive: 0x8...	MPI COMMUN...	Receive	0x0	111	111	44	44
4 Receive: 0x8...	MPI COMMUN...	Receive	0x0	174	430	252	508
5 Receive: 0x8...	MPI COMMUN...	Receive	0x0	130	305	151	407

Below the table are checkboxes for 'Unexpected', 'Show local ranks', 'Show global ranks', and 'Only ranks with messages'. A 'Select communicator' dropdown is also present.

Variable values window shows:

- `ntsamples`: <value optimized out>
- `nps`: <value optimized out>
- `nrecv`: 1065353216
- `ntsamples`: <value optimized out>
- `nvtxs`: 9224

Message queue debugging

# Automate debugging with offline mode

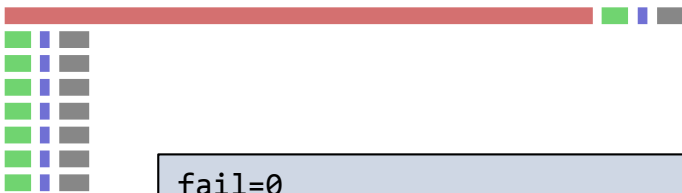
#	Time	Tracepoint	Processes	Values
1	21:18.172	jacobi_mpi_omp_gnu.exe (_jacobi.f90:83)	0-127	residual: 2.57

## Memory Leak Report

This report shows unfreed memory allocations when the program finished executing. Clicking an item in the bar chart below will show additional details about the allocations, including where they were allocated.

All 8 ranks:

Rank 0: 583.11 kB  
Rank 1: 58.71 kB  
Rank 2: 58.71 kB  
Rank 3: 58.71 kB  
Rank 4: 58.71 kB  
Rank 5: 58.71 kB  
Rank 6: 58.71 kB  
Rank 7: 58.71 kB



Legend

main (mmult3.c:139)  
ompi\_free\_list\_grow  
event\_del\_internal (minheap-internal.h)  
Other



```
fail=0
# --- check DDT tracepoint (residual)
f=jacobi_omp_mpi_gnu_debug.txt
resid=`grep ^tracepoint $f |awk -Fresidual: '{print $2}' |tail -1 |cut -c2-5`
if [ "$resid" != "2.57" ] ; then
    ((fail++))
    echo "Test has failed resid=$resid"
else
    echo "Test has succeeded"
```



# Conclusion

- Analyse application efficiency and understand behaviour with Allinea Performance Reports
- Develop faster by debugging and optimising large-scale applications with Allinea Forge
- Available for latest architectures:
  - x86\_64
  - KNL
  - CUDA 8.0
  - ARMv8
  - OpenPower

# allinea

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

## Thank you!

Any question, please ask.

Technical Support team : [flebeau@allinea.com](mailto:flebeau@allinea.com) or [support@allinea.com](mailto:support@allinea.com)  
Sales team : [sales@allinea.com](mailto:sales@allinea.com) or [marcin@allinea.com](mailto:marcin@allinea.com)

