# allinea

Now part of **ARM**

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

# Analyzing I/O Profiles

## I/O Profiling at Scale

Keeran Brabazon, ARM

DKRZ, UIOP Workshop
23rd March 2017

allinea FORGE   allinea DDT   allinea MAP   allinea PERFORMANCE REPORTS

# Acknowledgements

KTH:
Stefano Markidis, Sergio Rivas Gomez,
Bo Peng

# Allinea – What is it?

- HPC Tools company since 2002

- Help the HPC community develop and design the best applications and make the most use of HPC clusters

- Part of ARM since December 2016
  - Continue to improve tools for new uses in HPC
  - Support for all HPC applications and hardwares

allinea
Now part of ARM

# Products

- ## Allinea Forge
  - Combined debugging and profiling in same interface
  - Designed for application developers

- ## Allinea Performance Reports
  - Summary of application performance
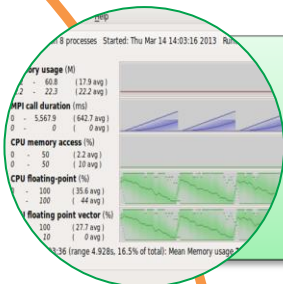  - Designed for system administrators
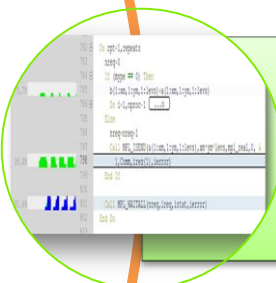
# Profiling with Allinea MAP

# Allinea MAP

- Use of performance analysis tools can help to focus attention on the parts of a program with worst performance

- Allinea MAP can do so for applications running with 100k+ processes
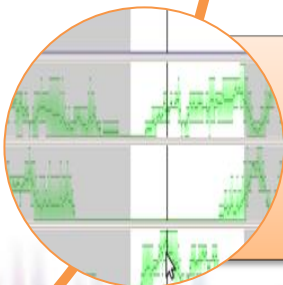
# Allinea MAP



## Low overhead measurement

- Accurate, non-intrusive application performance profiling
- Seamless – no recompilation or relinking required

## Easy to use

- Source code viewer pinpoints bottleneck locations
- Zoom in to explore iterations, functions and loops

## Deep

- Measures CPU, communication, I/O and memory to identify problem causes
- Identifies vectorization and cache performance

**allinea**
Now part of **ARM**
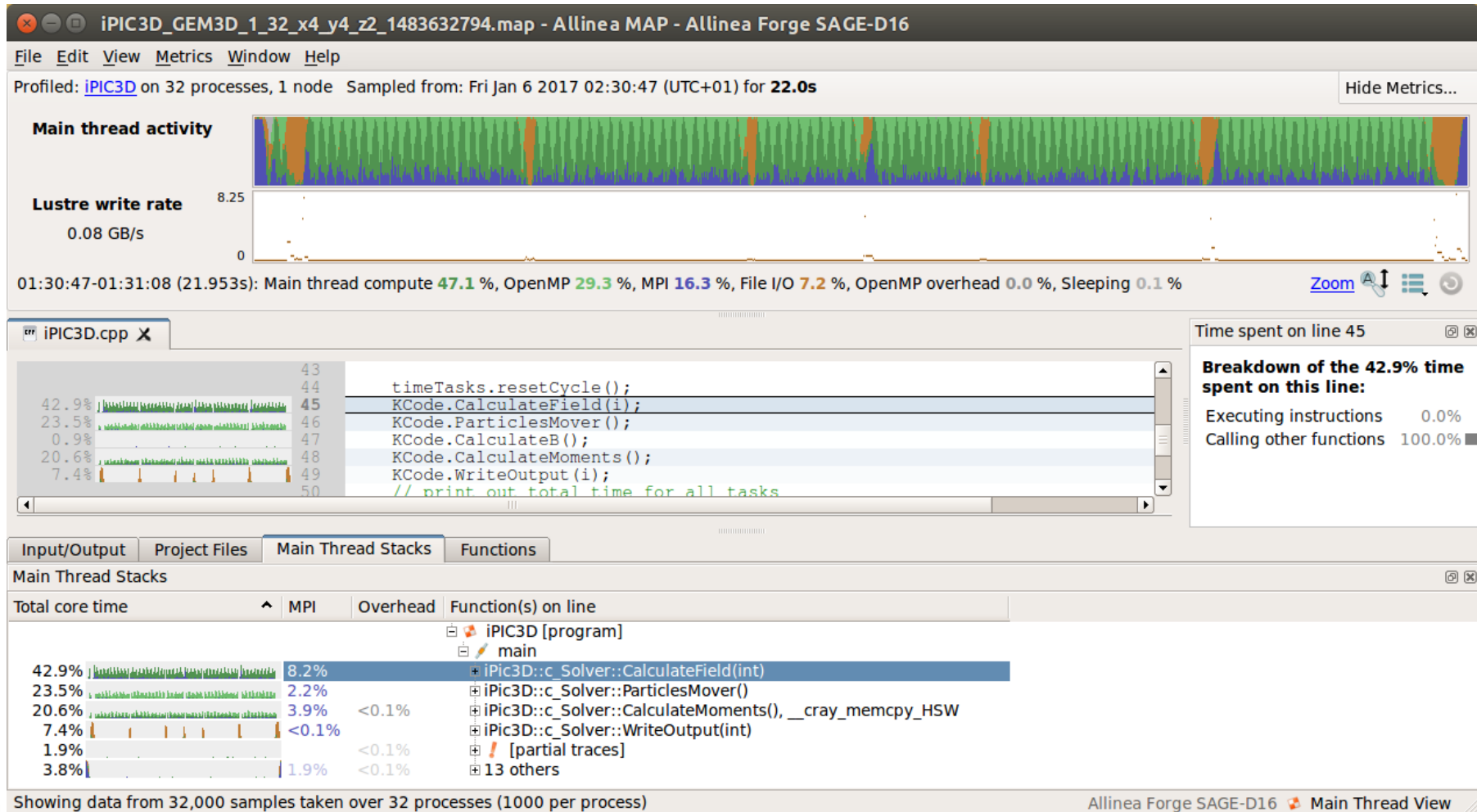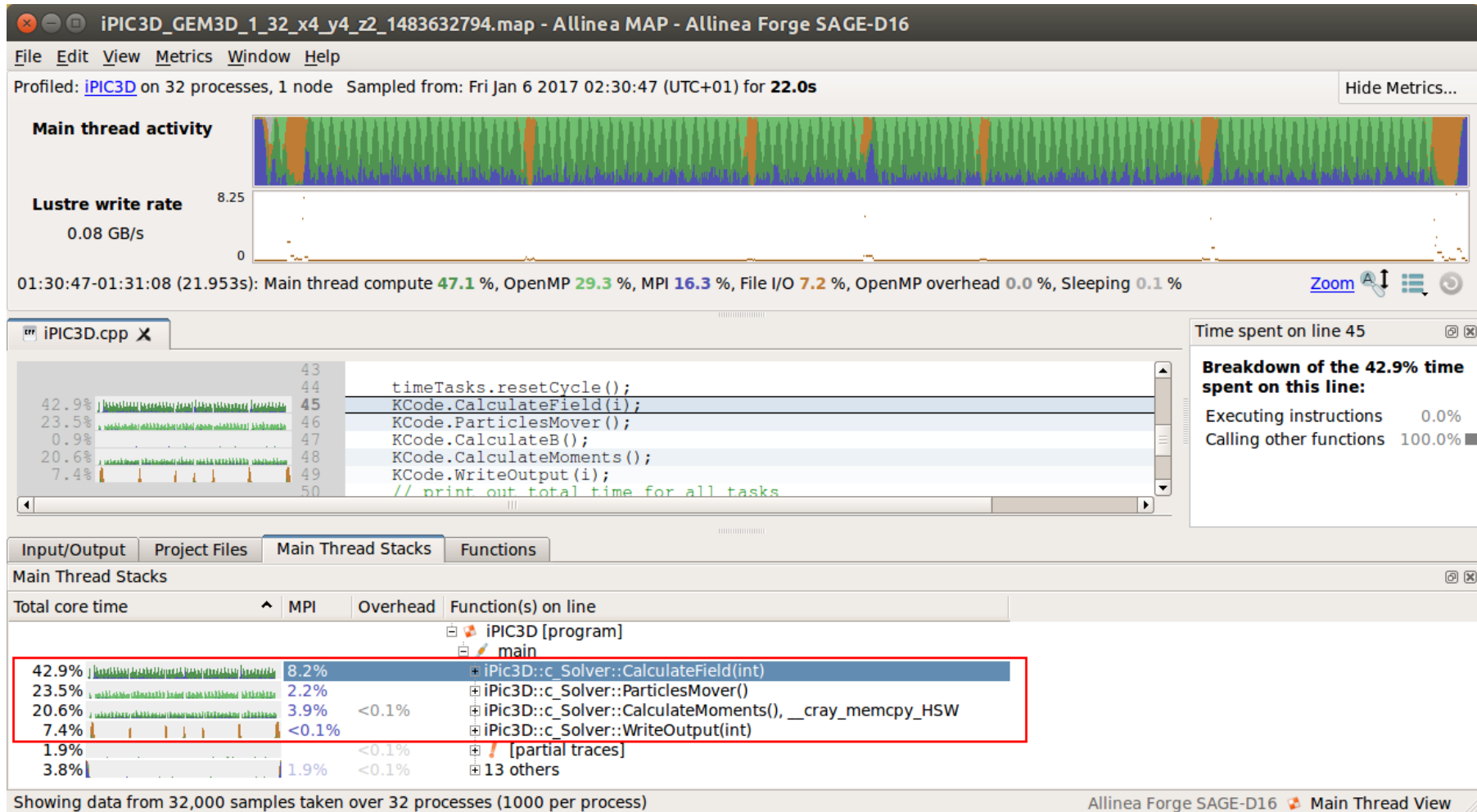
# I/O Profiling – Worked Example

iPIC3D

# iPIC3D

- Particle-in-cell code to model interaction between solar wind and Earth's magnetic field

- Practical problem sizes have billions of particles with velocity, current and charge density

- I/O performed for visualization (every 20 iterations) and checkpointing (every 50 iterations)

- Run on Beskow – Cray XC40, 32 Broadwell per Node
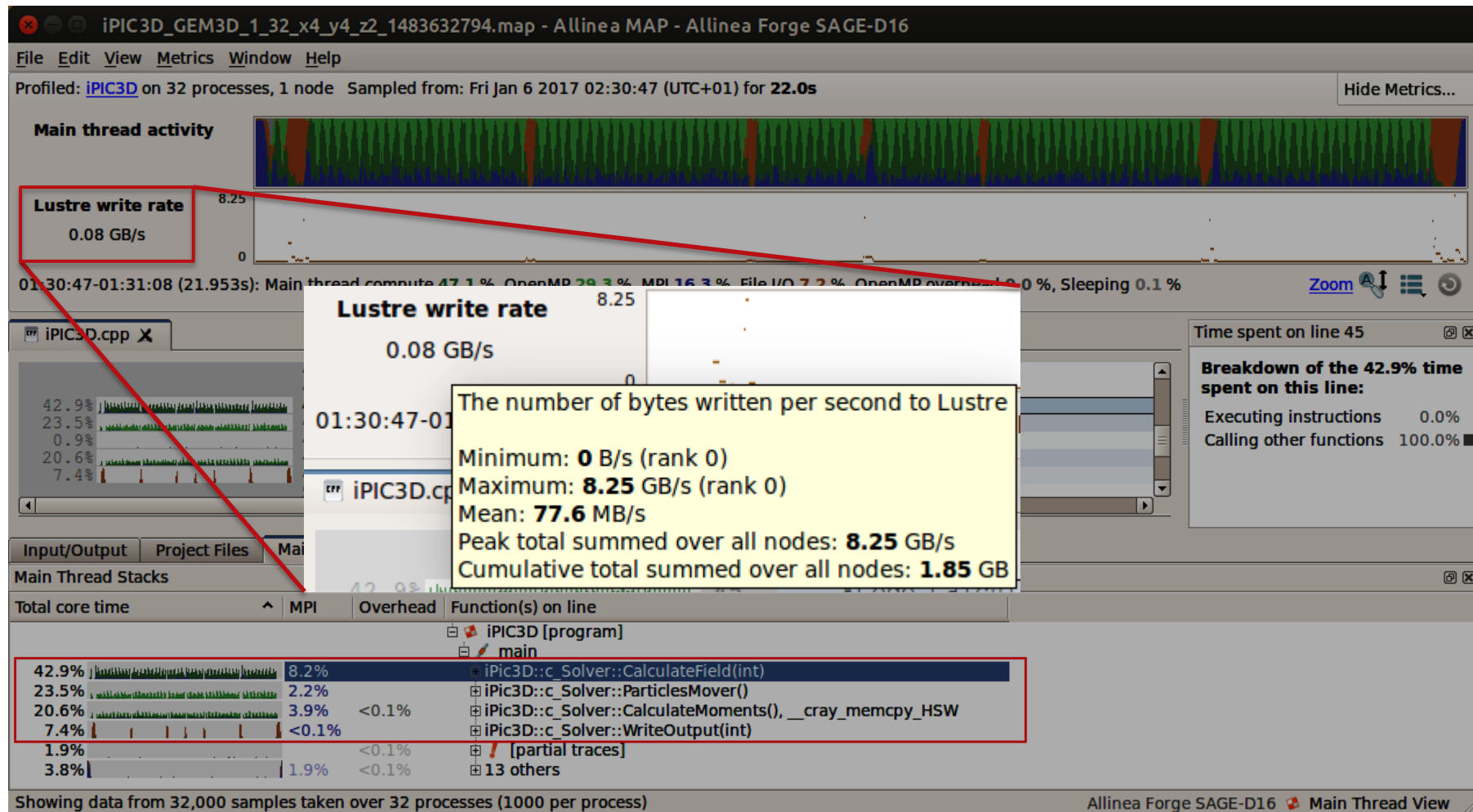
allinea
Now part of ARM

# iPIC3D Profiling – 32 Processes (1 node)

# iPIC3D Profiling – 32 Processes (1 node)

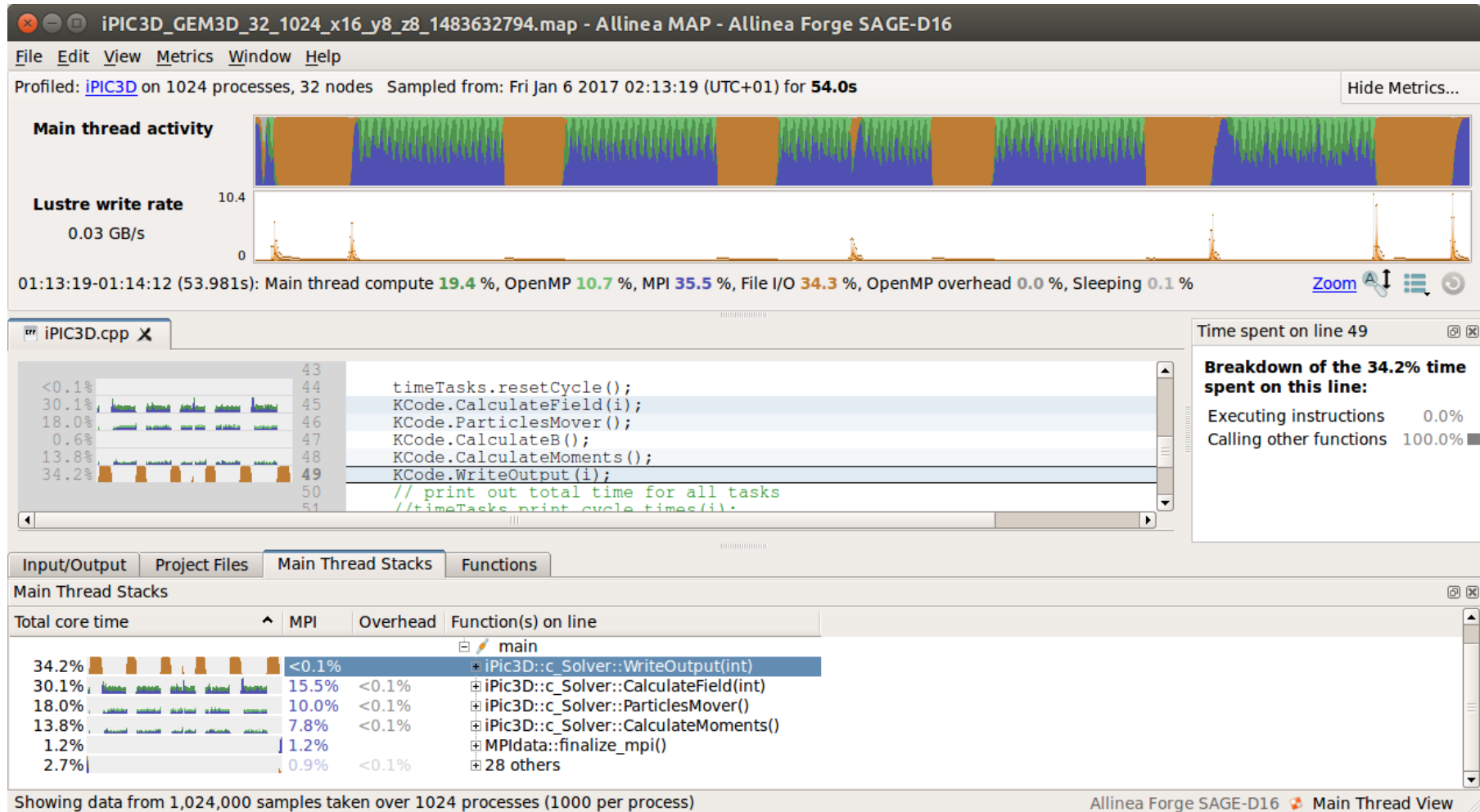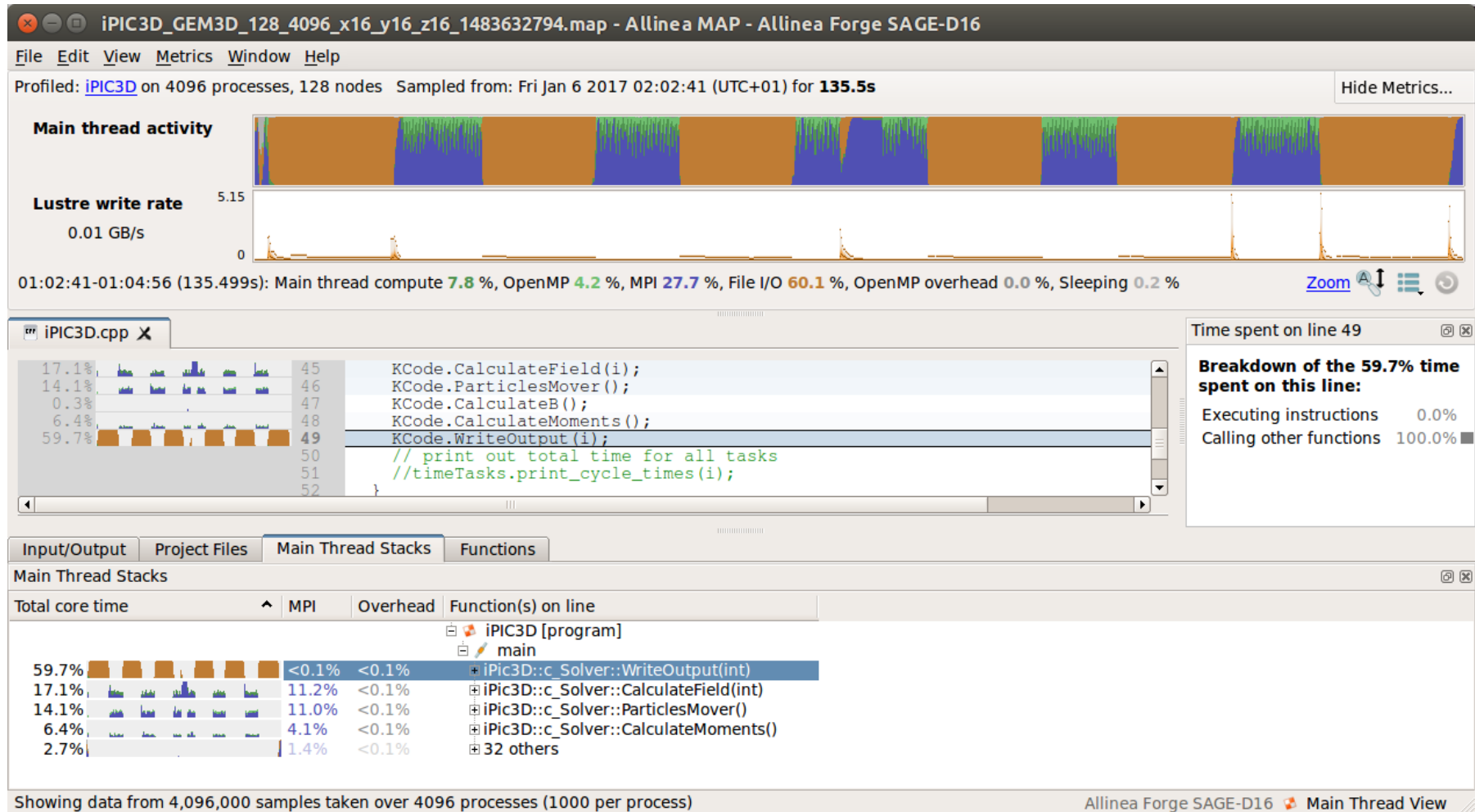# iPIC3D Profiling: 32 Processes (1 node)

# iPIC3D Profiling

- I/O does not take up a large amount of run time

- 32 processes is rather small – go to larger core counts with more I/O performed

# iPIC3D Profiling: 1024 Processes (32 Nodes)

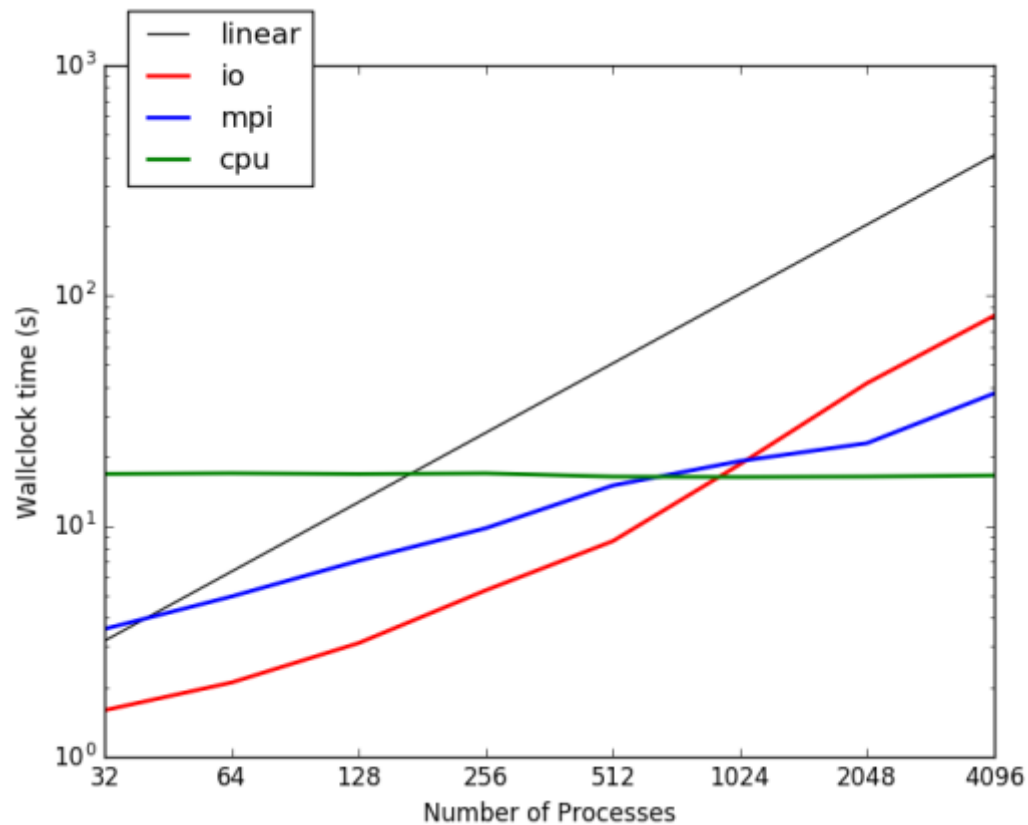# iPIC3D Profiling: 4096 Processes (128 Nodes)

# iPIC3D Profiling

- Experiments were run to show weak scaling

- Scaling behaviour cumbersome to view in MAP files

- Export to JSON of profile (new in Forge version 7.0+) allows user to post-process and visualise data
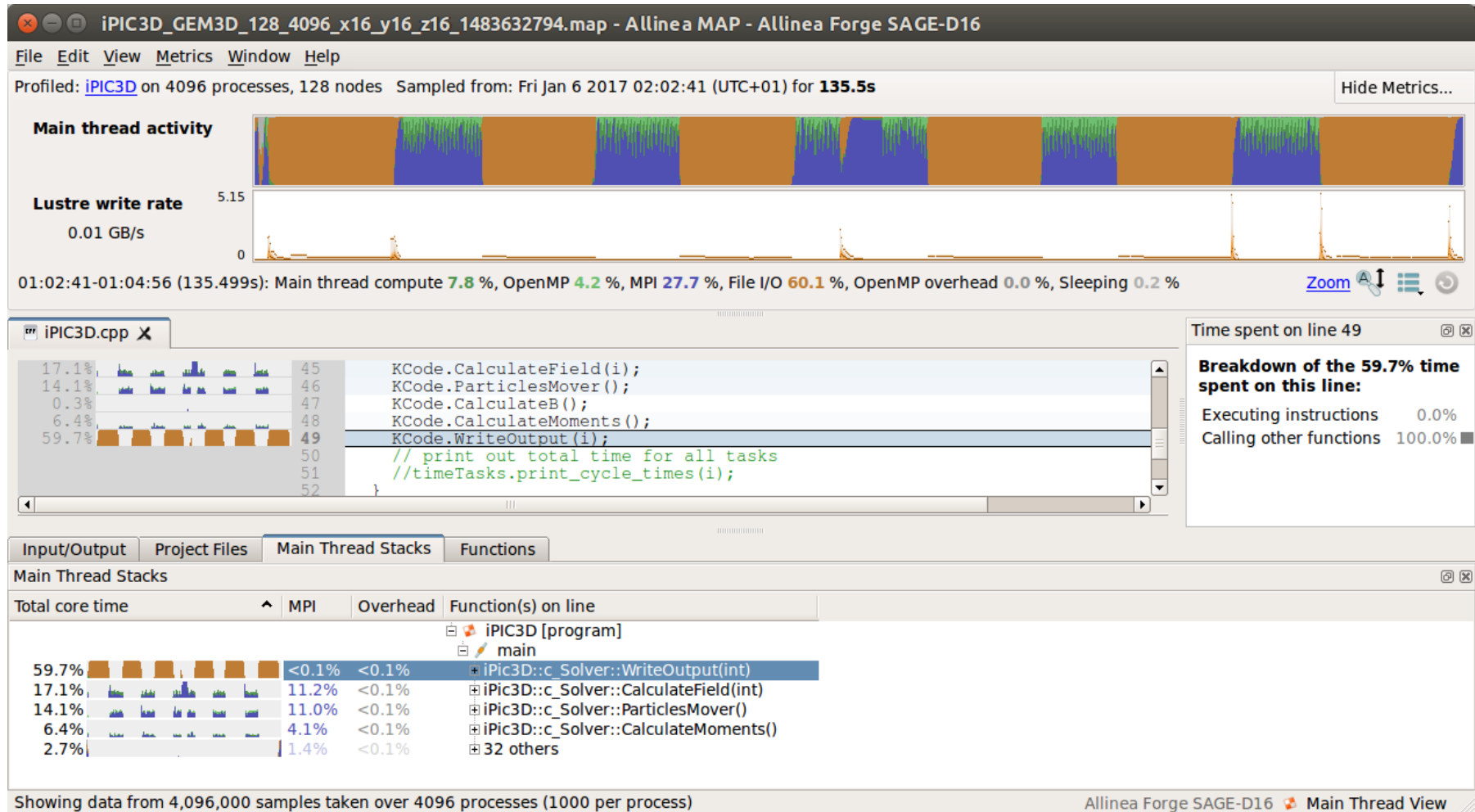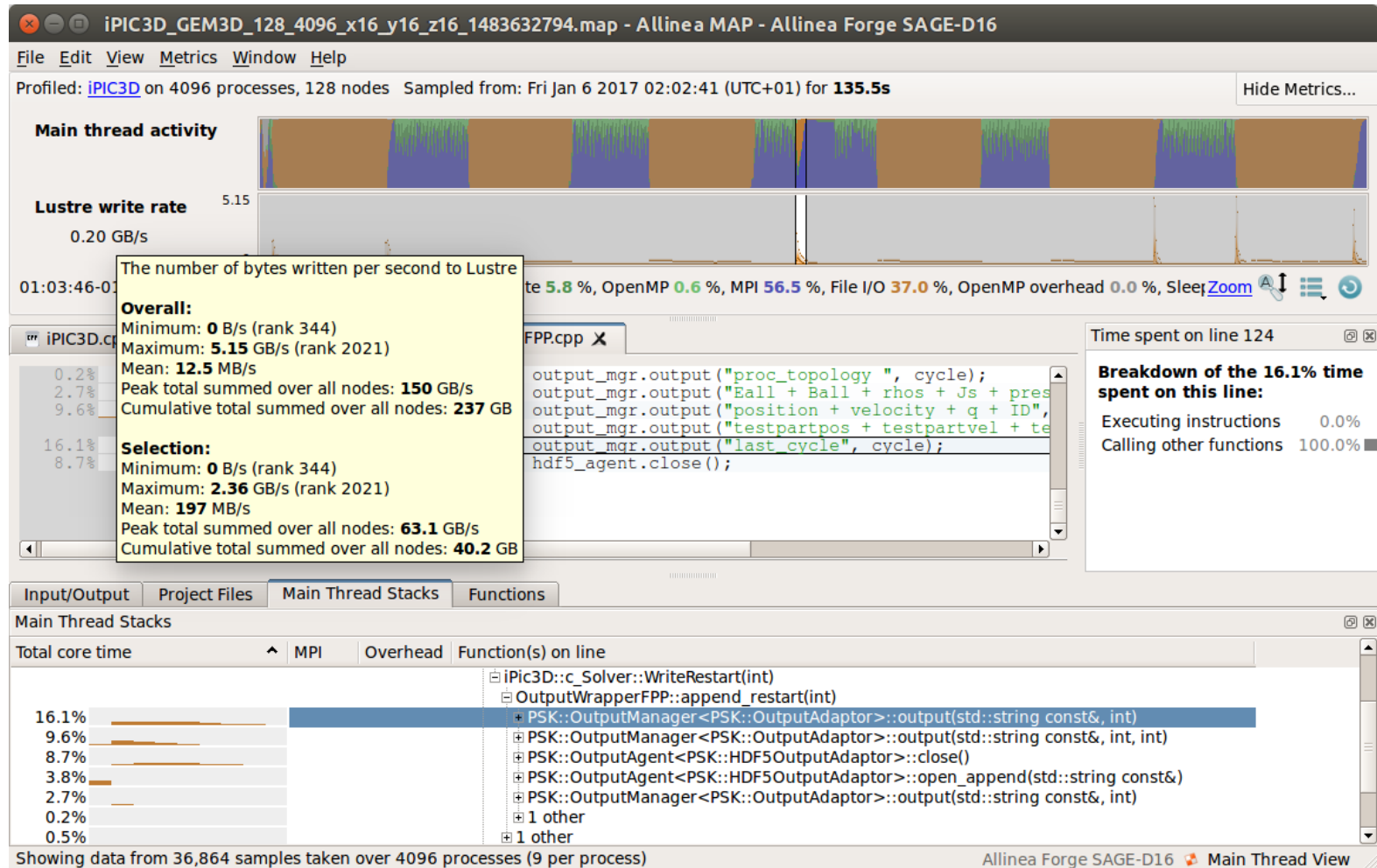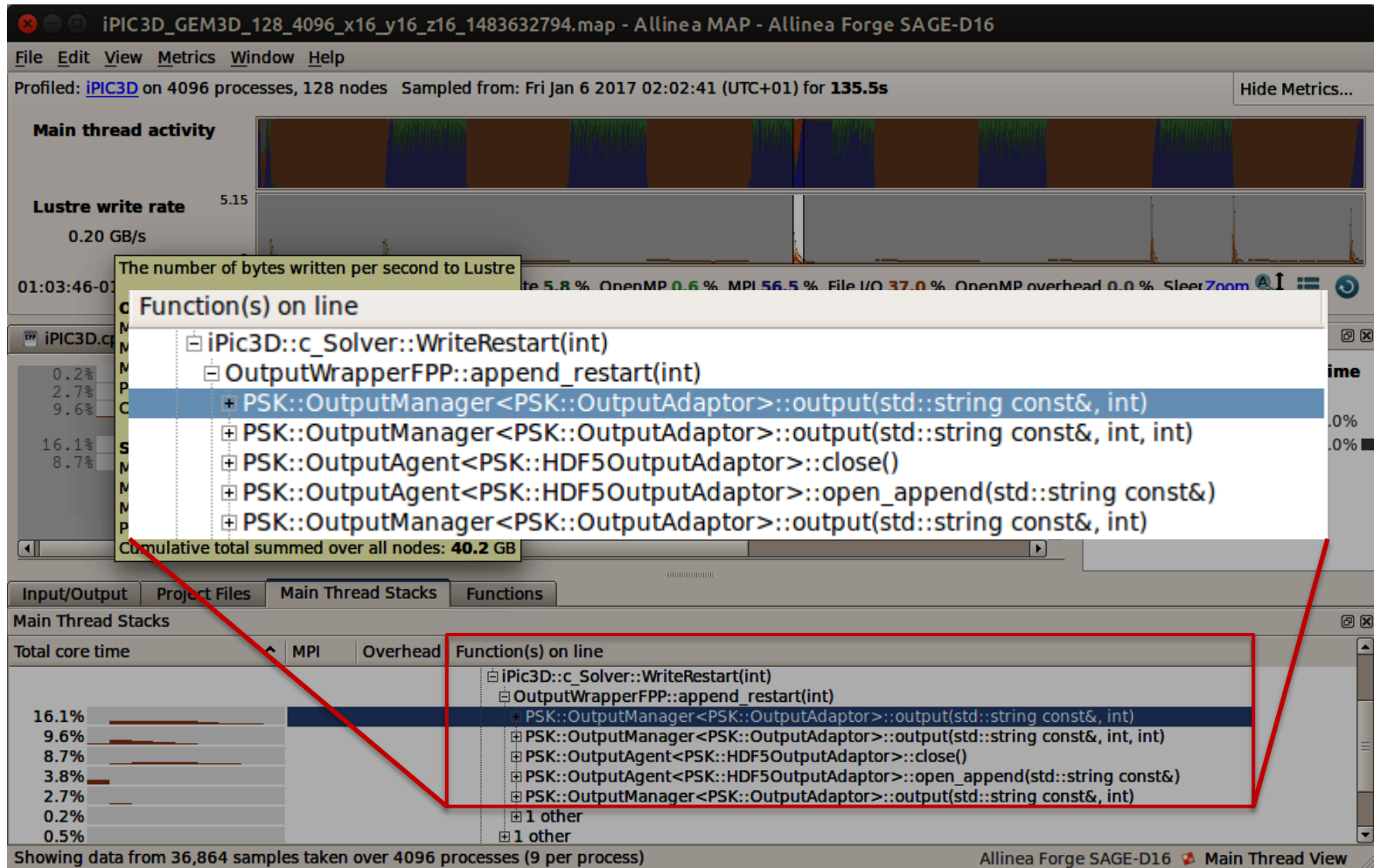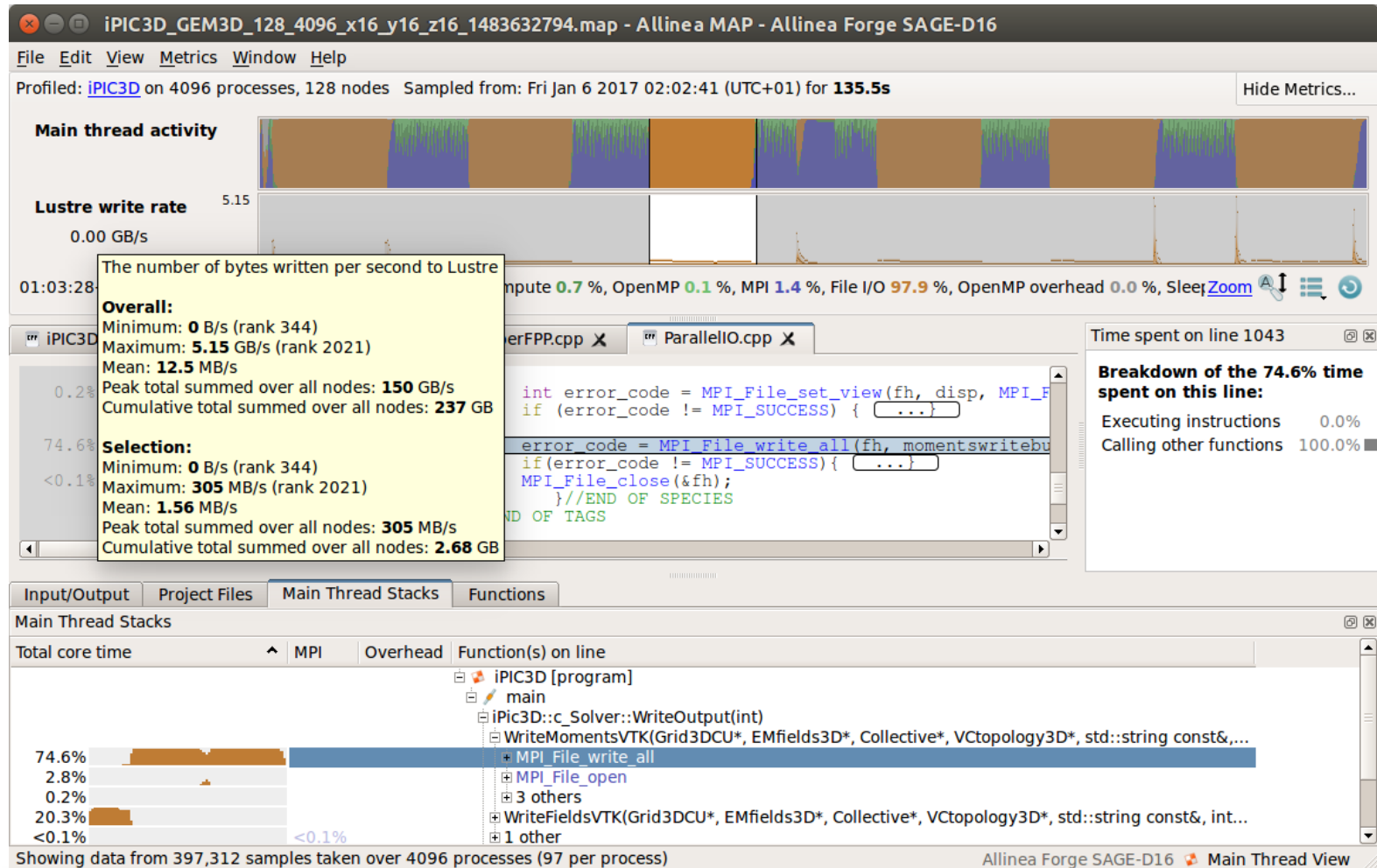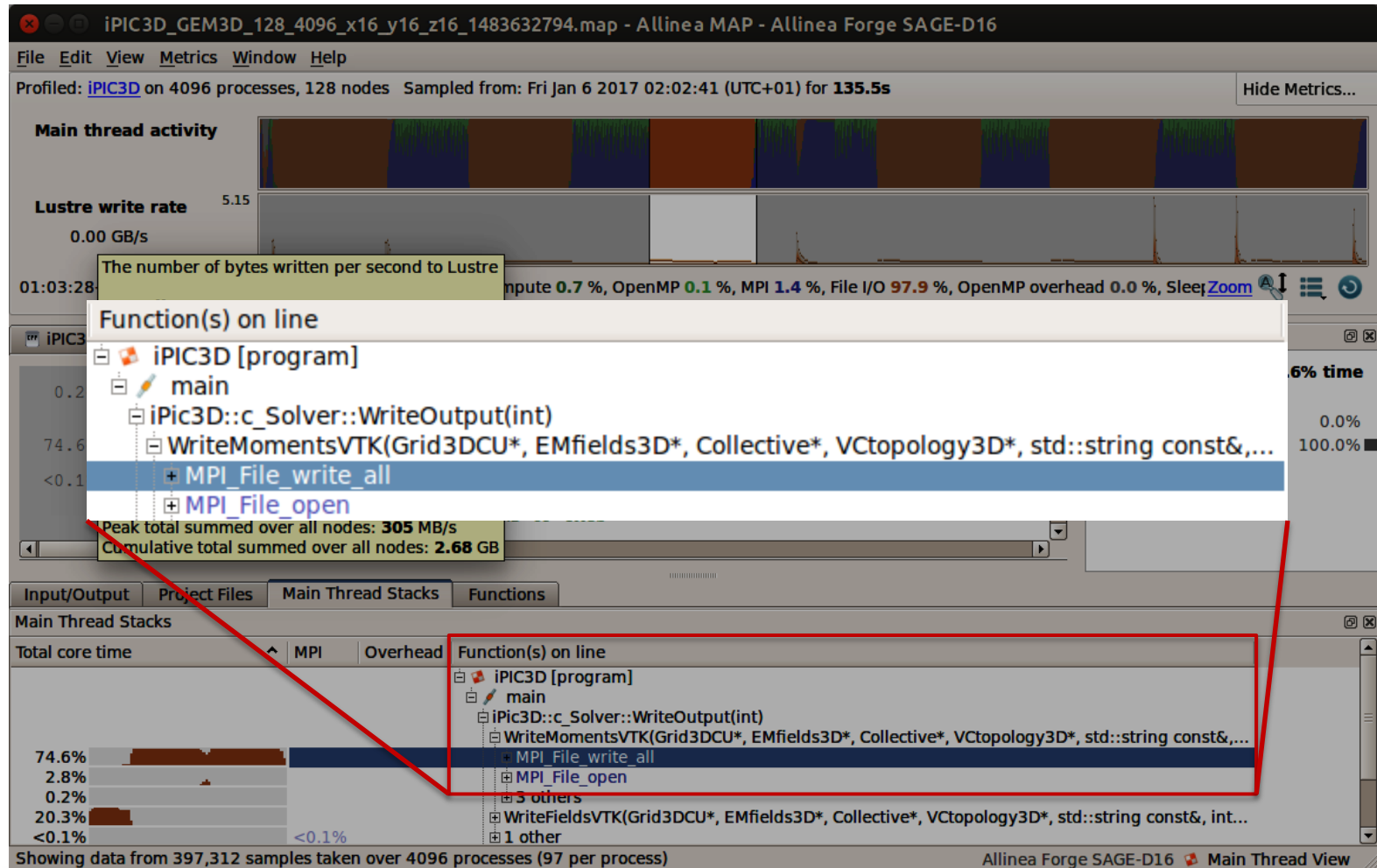
# iPIC3D – Weak Scaling

# iPIC3D Profiling: 4096 Processes (128 Nodes)

# iPIC3D Profiling: 4096 Processes (128 nodes)

# iPIC3D Profiling: 4096 Processes (128 nodes)

# iPIC3D Profiling: 4096 Processes (128 nodes)
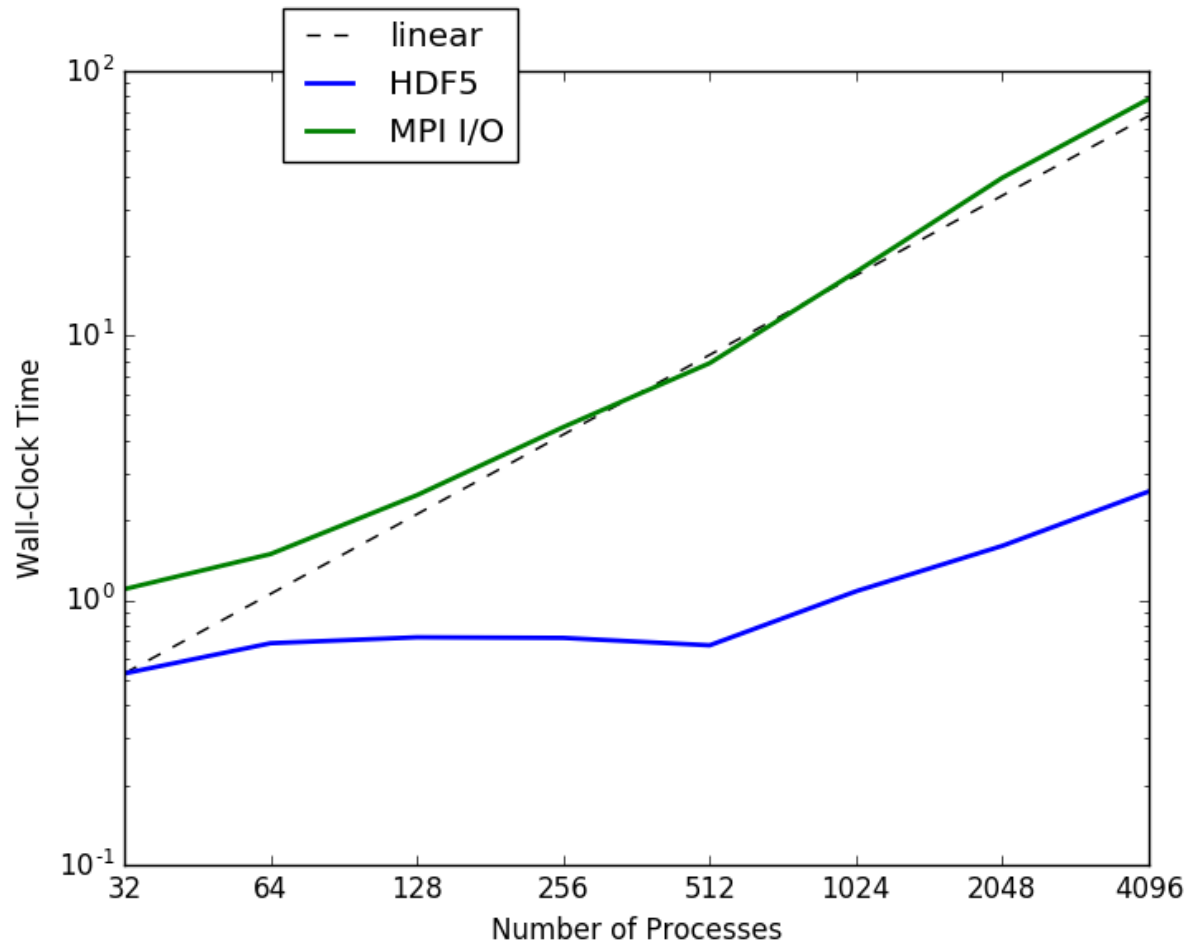
# iPIC3D Profiling: 4096 Processes (128 nodes)

# iPIC3D Profiling

- HDF5 and MPI identified as I/O libraries for fast and slow phases of I/O, respectively

- Allinea MAP can show percentage time spent in different functions

- Calculate wallclock time spent in MPI I/O and HDF5 I/O and plot over weak scaling runs

allinea
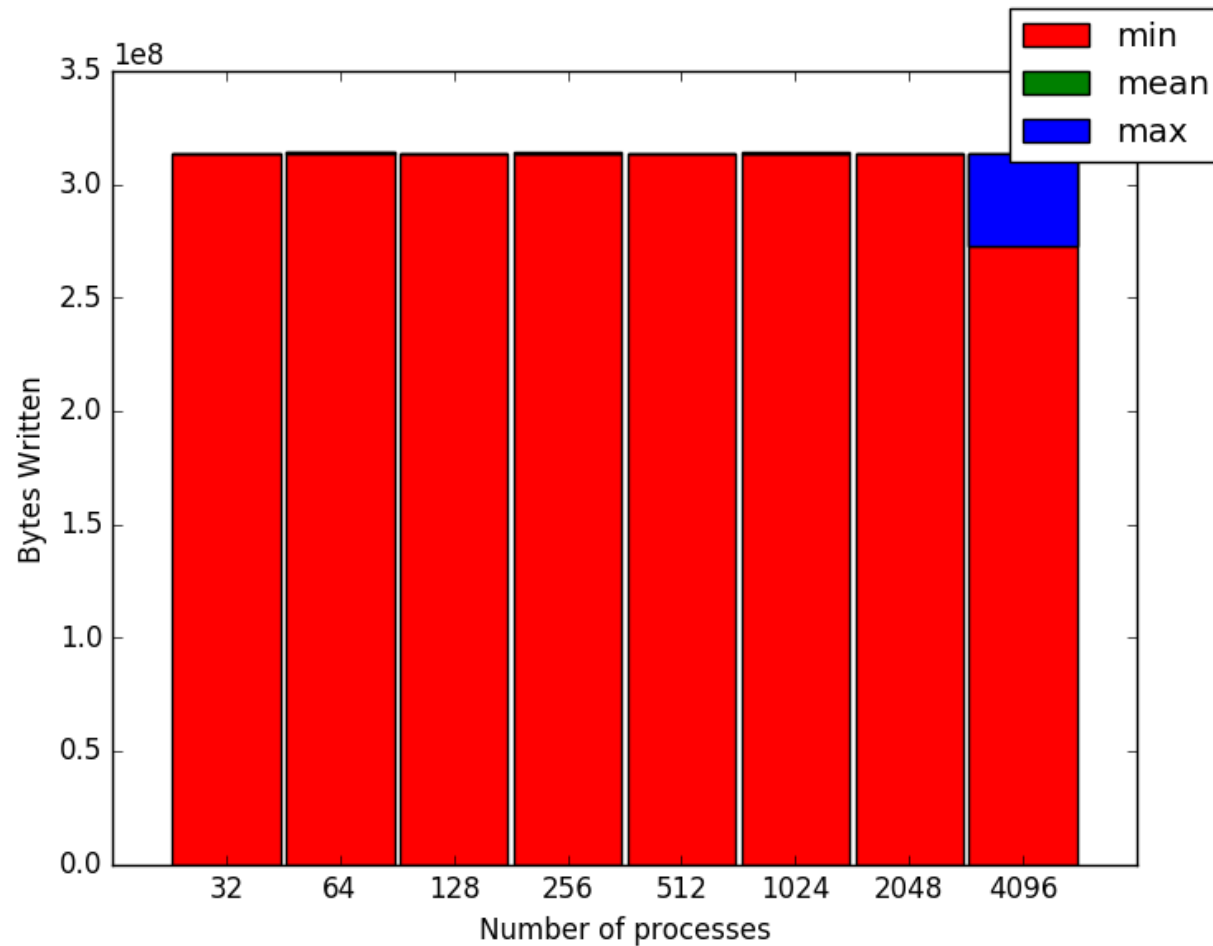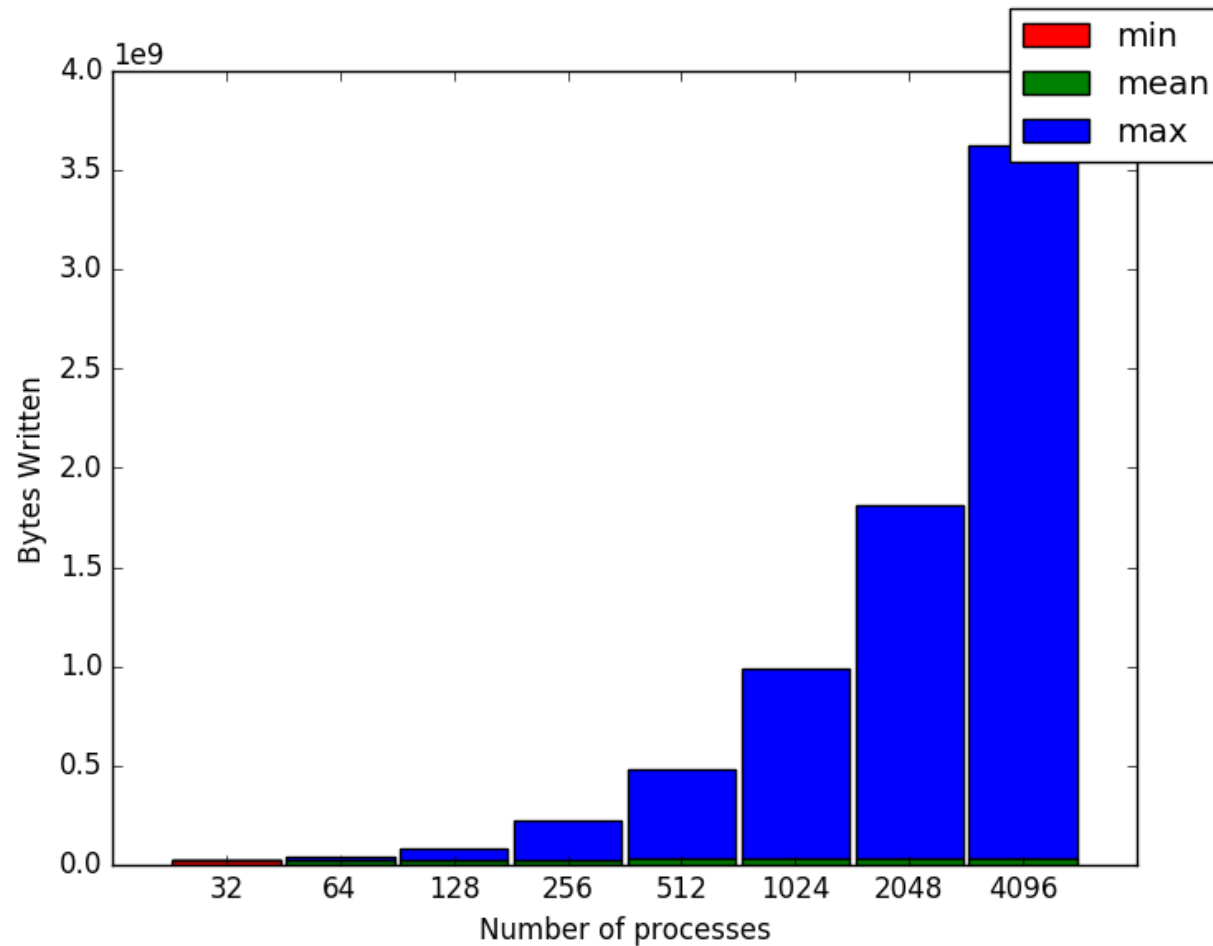Now part of **ARM**

# iPIC3D I/O time

# iPIC3D Profiling

- MPI I/O performs collective write from all processes. 512kB are written per process. Scaling of time is linear with amount of data to write

- HDF5 performs small file write from all processes (10MB per process)

- How much data is written from each node?

# iPIC3D I/O Volume – HDF5

# iPIC3D I/O Volume – MPI I/O

# iPIC3D I/O

- Cray collective MPI I/O accumulates data to a single node


- Possible improvements
  - Dedicated I/O nodes
  - Asynchronous I/O (e.g. burst buffers, function offloading - SAGE)
  - Limiting number of writers per file

# I/O Areas of Interest

- Ideal Performance – is it achievable or desirable to achieve for real applications?

  - How sensitive to system performance would an application be which achieves the maximum I/O throughput?

  - Can I/O bandwidth be maximised in an application through better I/O management?

allinea
Now part of ARM

# I/O Areas of Interest

- System I/O – how does application I/O tie in with the system?
  - View application I/O alongside system I/O

  - What about I/O in a group of related programs (i.e. workflows)?

  - Should I/O bandwidth be maximised at a system or workflow level rather than at an application level?

allinea
Now part of ARM

# Last Words - Profiling

- Profiling at small scales may not show the whole picture – measurement at larger scales show problems related to scale

- MAP provides low overhead measurement with rich information

- Export to JSON of performance data allows for post-processing and analysis across many runs