# Tutorial Instruction Sheet: Analysing Parallel I/O Using Allinea MAP

Keeran Brabazon, ARM
ISC 2017, Frankfurt

## Summary

In this tutorial you will learn the following:

- How to open an Allinea MAP file and find metrics salient to file I/O
- Interpret the metrics presented to give insight into the operation of I/O routines using different I/O libraries
- How to extract information relevant to a selection of a performance profile in Allinea MAP

The exercises will give you an insight into how to get started with the analysis of I/O performance from a program run. There is extra work that can be done, and you are welcome to continue to perform an investigation into I/O with the use of the Allinea tools. If your institution does not have an installation of the Allinea tools you may request a trial from https://www.allinea.com/get-your-free-allinea-forge-and-allinea-performance-reports-trial
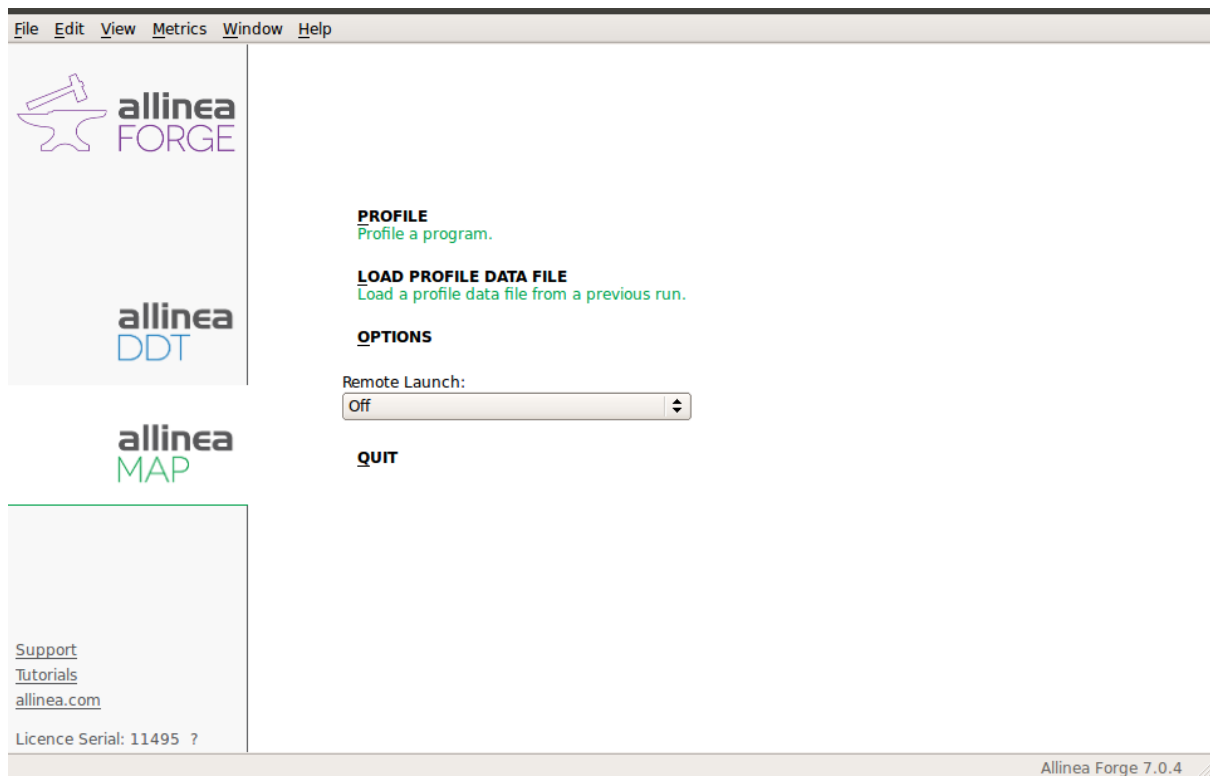
## Getting started

### Installing the Tools

- The tools need to be installed locally (i.e. not on ARCHER)
- Navigate to https://www.allinea.com/products/forge/download
- Choose the download for the appropriate platform you are using (Windows and Mac have front-end clients)
- Follow the installation instructions at https://www.allinea.com/user-guide/forge/Installation.html#x5-70002
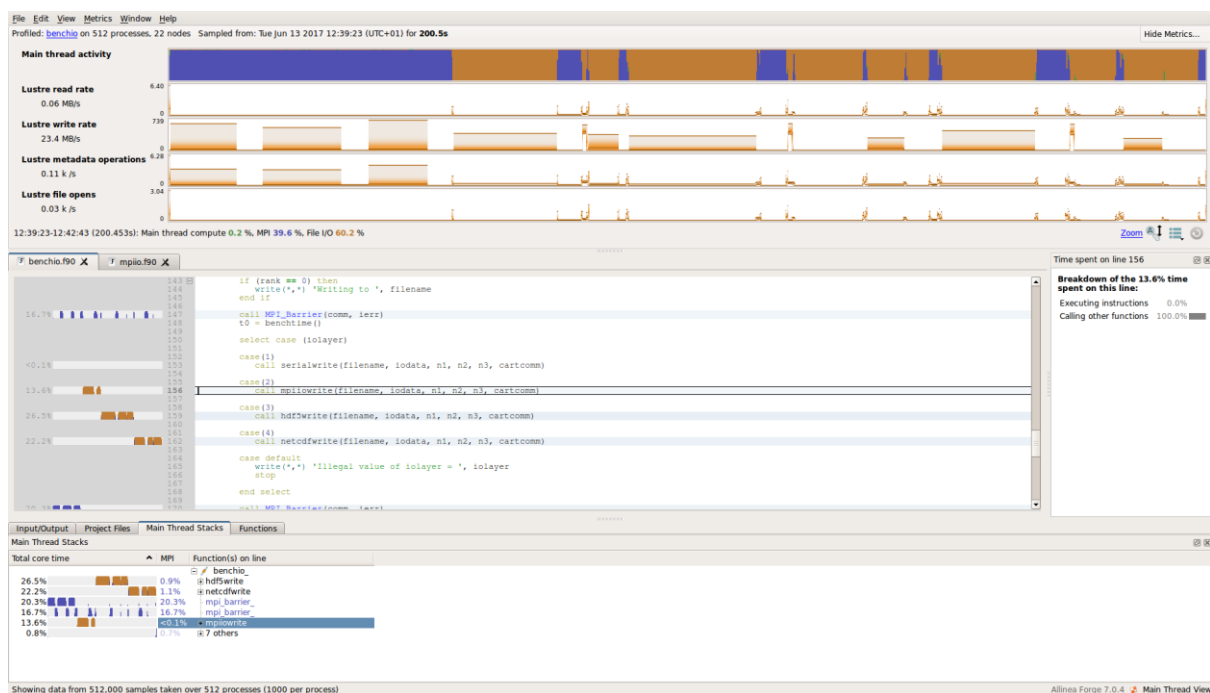
### Getting necessary resources

- Download the source files of the `benchio` program introduced in an earlier session onto your local machine
- Ensure that you have a folder called `TutorialData` with a set of MAP profiles (files ending in .map) as well as a set of output files from the benchmark runs. This file should have been provided with the folder.

## Opening a profile

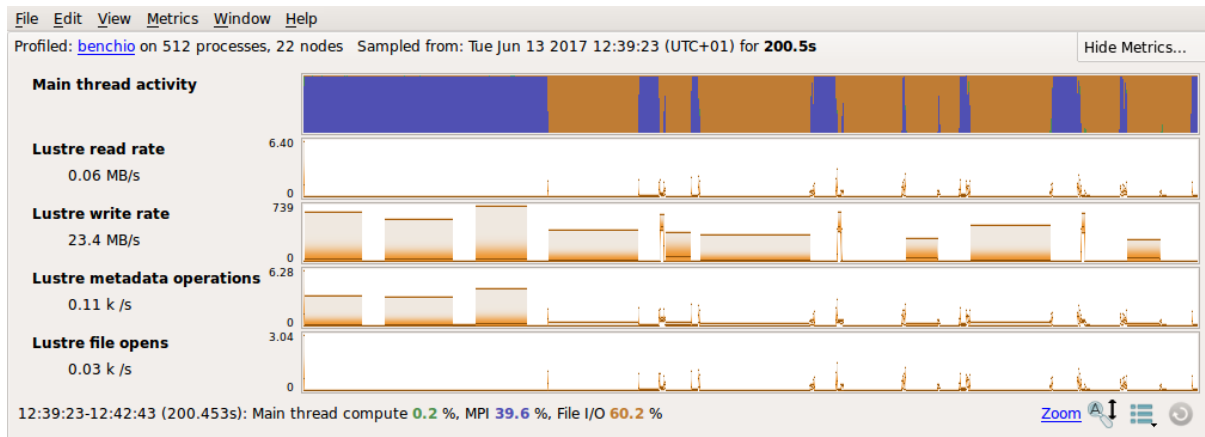- Launch the Forge GUI, and navigate to the MAP pane



- Select the 'Load Profile Data File' option
- Select the file called benchio_512p_128_1.map from the folder in which this document was contained and the profile will be loaded
- Show the Lustre metric information by selecting the 'Metrics -> Preset: Lustre' menu option
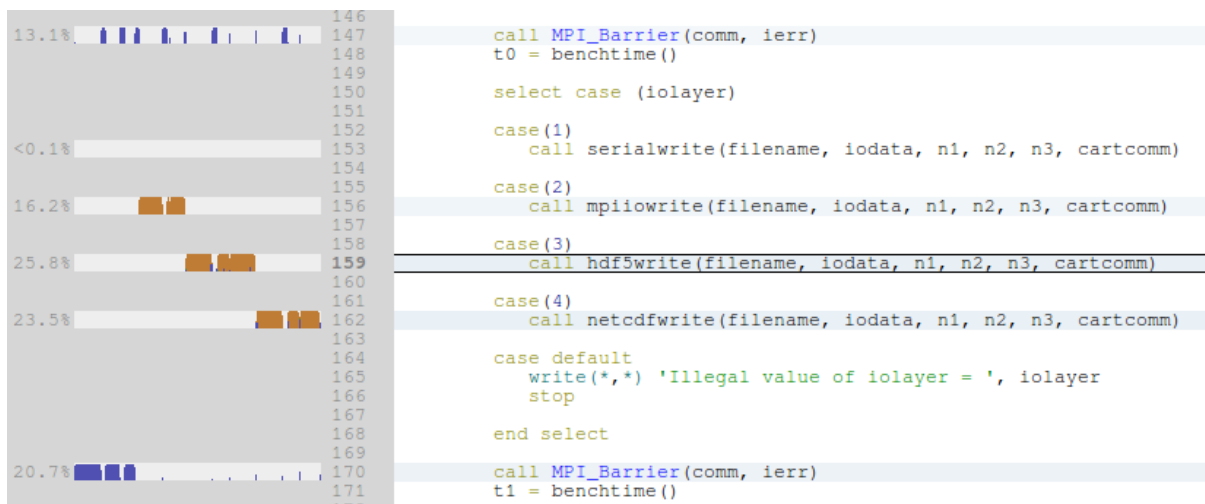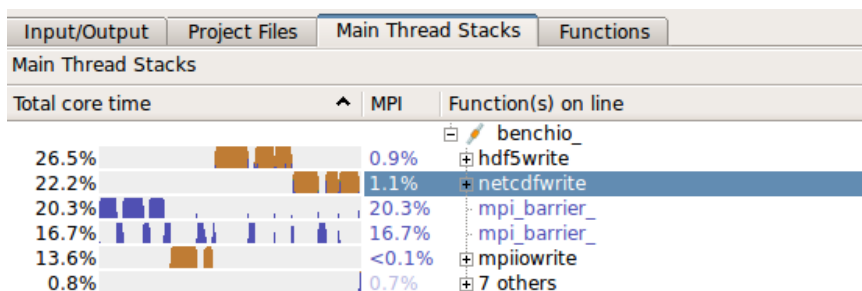
## Layout of a MAP profile

- The top-most section of a profile contains graphs over time of different metrics
- The 'Main Thread Activity' timeline is a colour-coded histogram of activity. Blue indicates MPI activity, orange is I/O. Other colours indicate activity not measured in these programs (e.g. green for compute activity such as floating point (vector) operations and memory accesses).



- Central section shows the source code.
- Each line on which activity was recorded has a graph similar to the 'Main thread activity' timeline which indicates when activity was recorded on that line
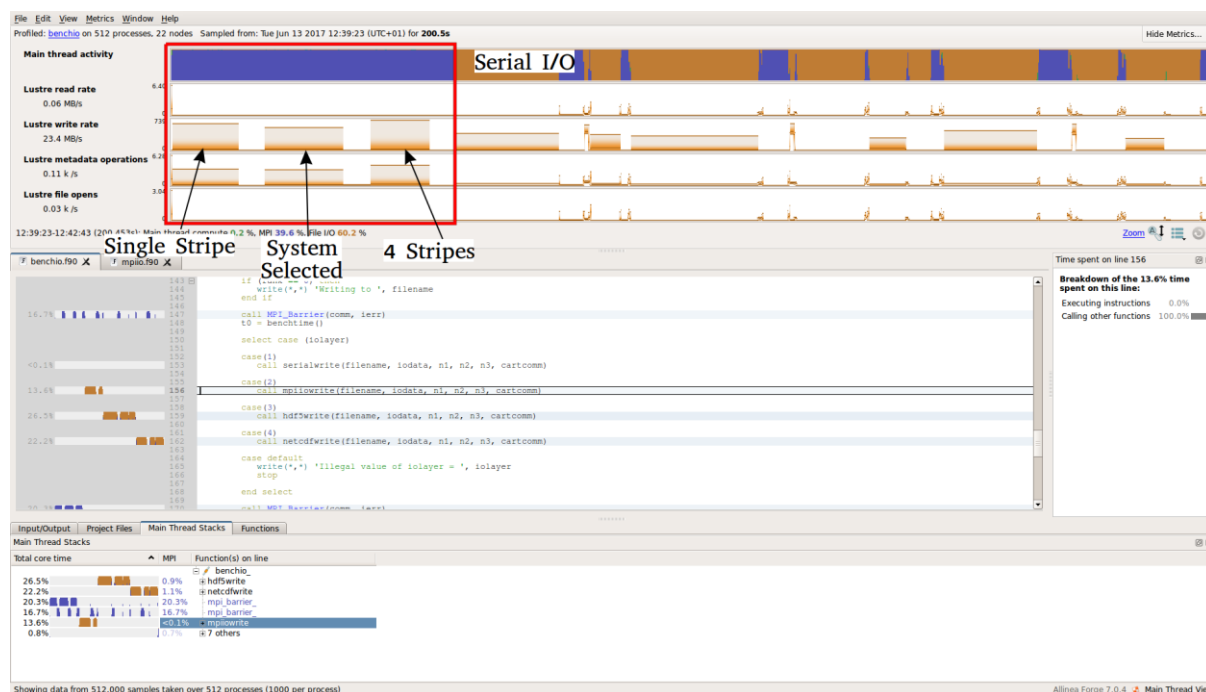


- Lower-most section has several tabs to show different data
- Default tab shows the call stacks in which activity was recorded.
- Activity timelines in the left-hand margin indicate when activity was recorded in the call stacks
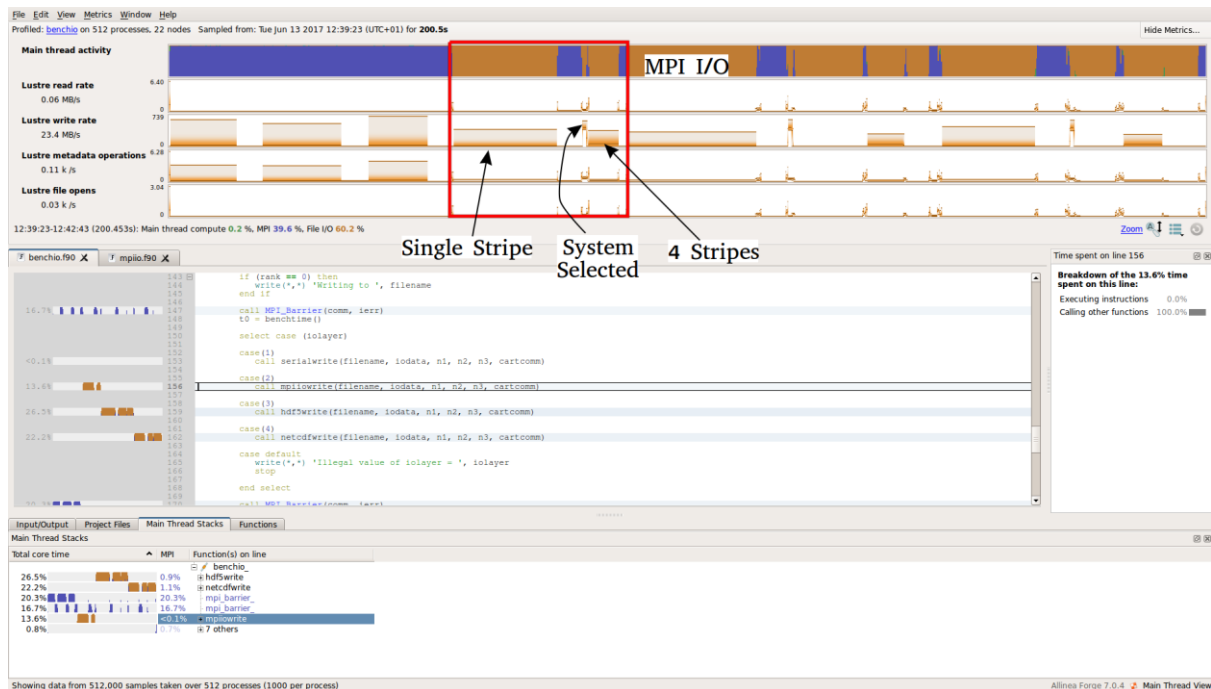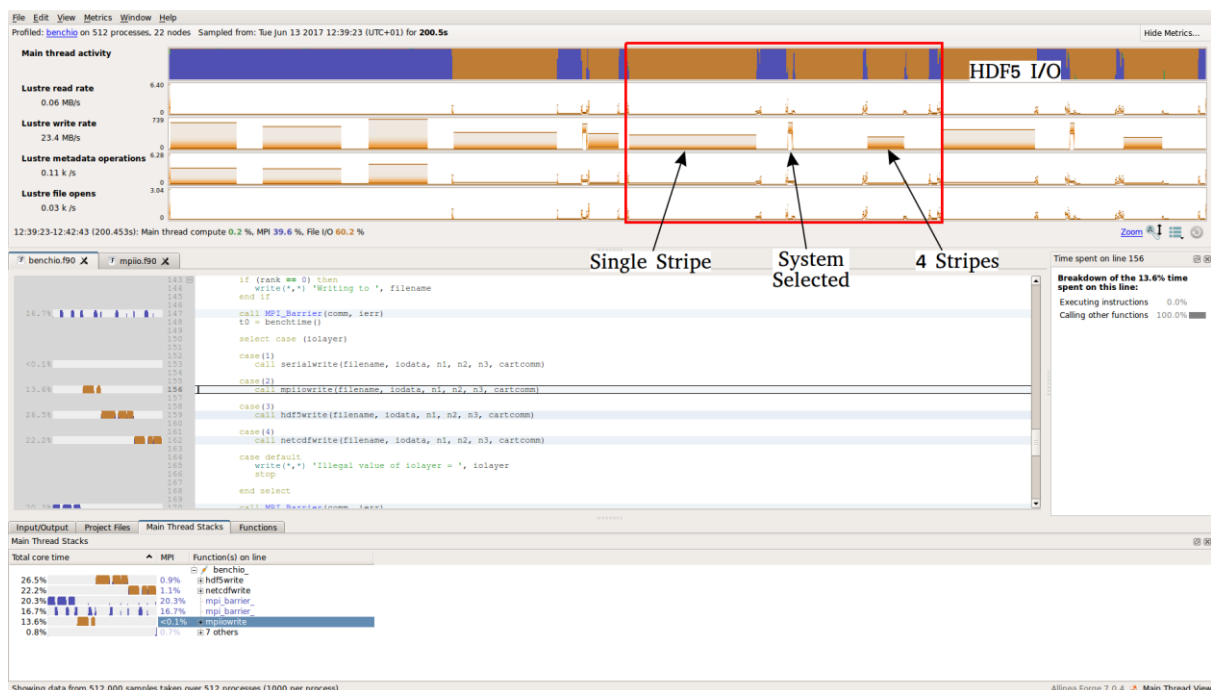
## Matching a profile to a program

- Ensure the profile benchio_512p_128_1.map is open
- Recall the program profiled performs I/O using the paradigms
  - Serial I/O
  - MPI I/O
  - HDF5 I/O
  - NetCDF I/O

  and that it writes data to files which have
  - 1 stripe
  - System selected striping
  - Default system striping (this is 4 on ARCHER)
- From the write rate identify the writes using Serial I/O

- The writes using MPI I/O



- And the writes using HDF I/O and NetCDF I/O

# Investigation of HDF5 I/O

This section starts demonstrates how to investigate the performance of I/O in the applications that have been considered. Only the performance of the HDF5 output data is considered. Interpreting the performance of the other methods of I/O presented is left as an exercise.

- Look at the bandwidth reported in the output file benchio_512p_128_1.o which is for the MAP profile benchio_512p_128_1.map
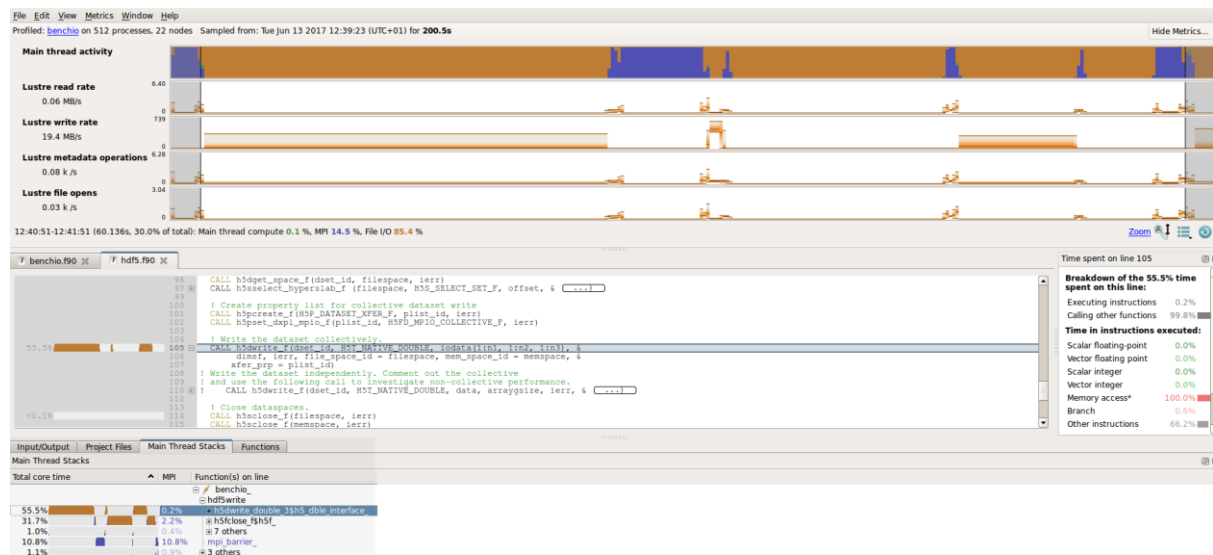- Output for HDF5 is as follows

```
------
HDF5
------

Writing to unstriped/hdf5.dat
time =  26.017146230908111 , rate =  314.86927610331014  MiB/s
Deleting: unstriped/hdf5.dat

Writing to striped/hdf5.dat
time =  15.180744504323229 , rate =  539.63097776048153  MiB/s
Deleting: striped/hdf5.dat

Writing to defstriped/hdf5.dat
time =  12.465544771635905 , rate =  657.17143936140462  MiB/s
Deleting: defstriped/hdf5.dat
```
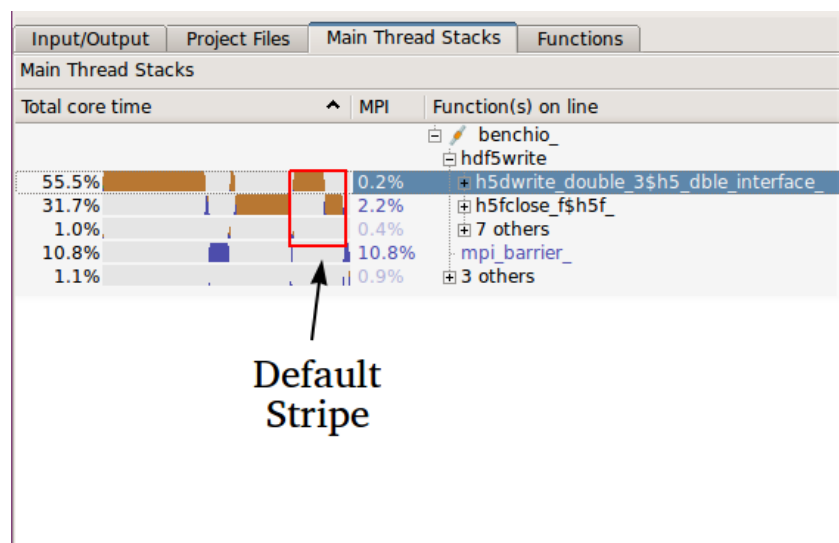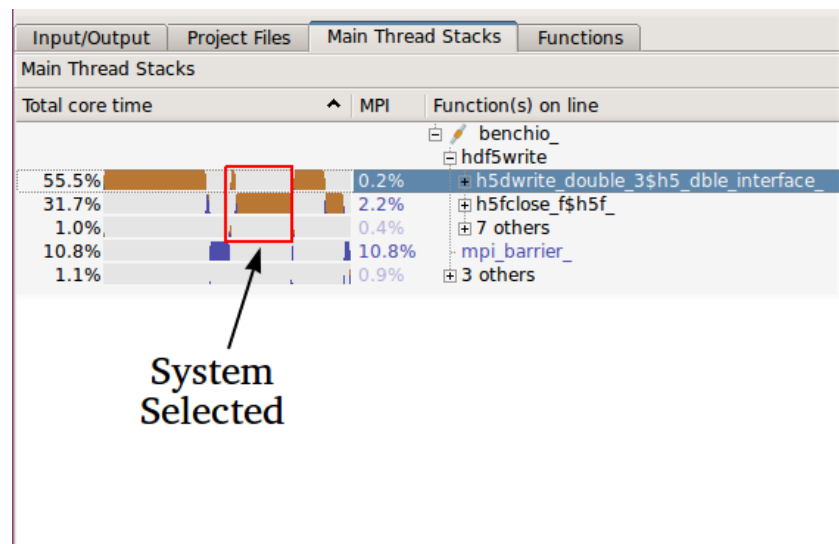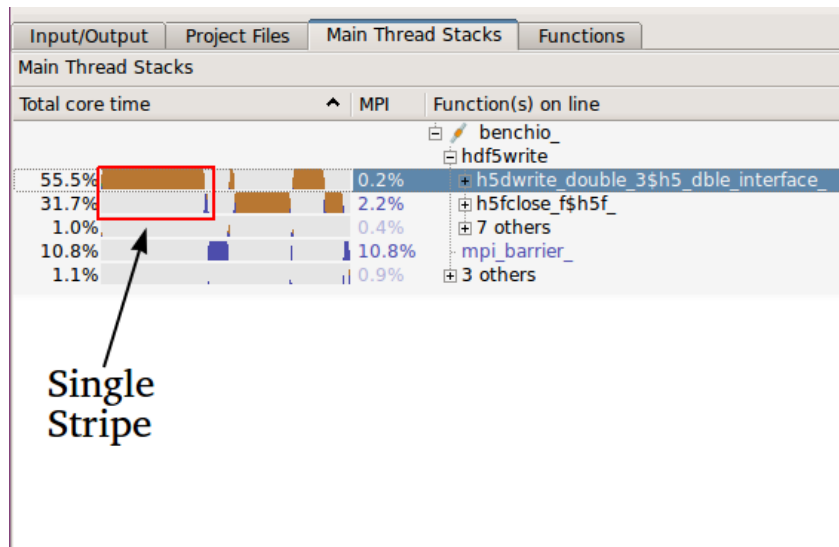
- Identify when writing to HDF5 files was taking place by looking at the activity timeline in the stack view in the lower-most pane of the GUI
- Left click on the metrics timeline and drag to select when the `hdf5write` procedure was active. Left click to zoom in
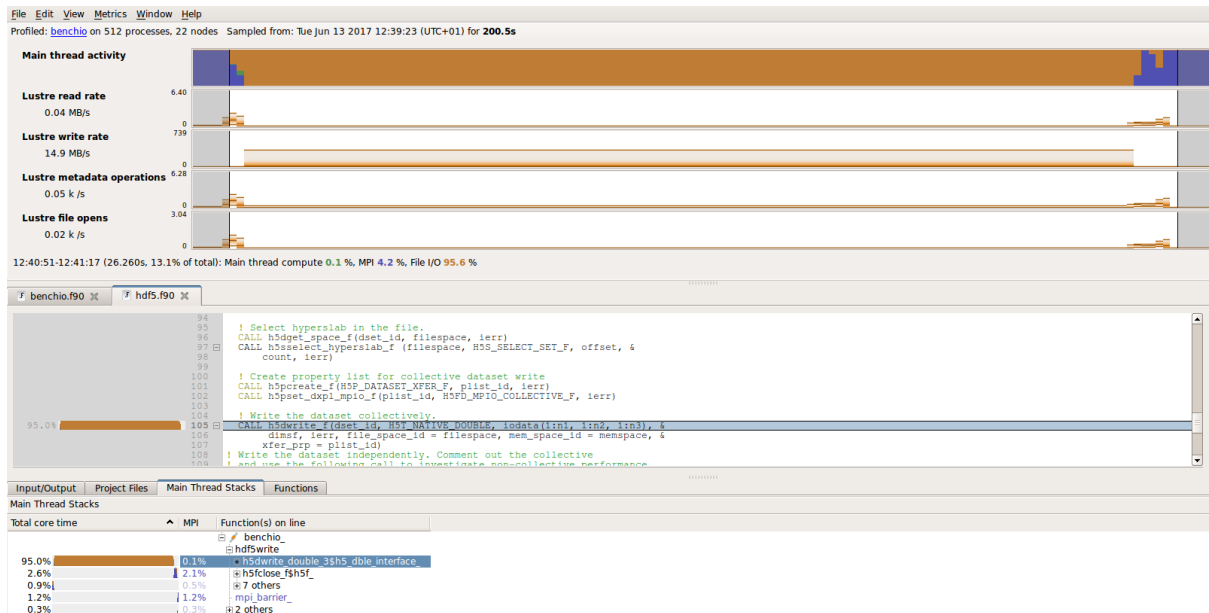


- Expand the `hdf5write` procedure in the call stack in the lower-most section of the screen
- The activity is split into three parts – first for write to file with single stripe, then system selected stripe, then default striping.

Single Stripe



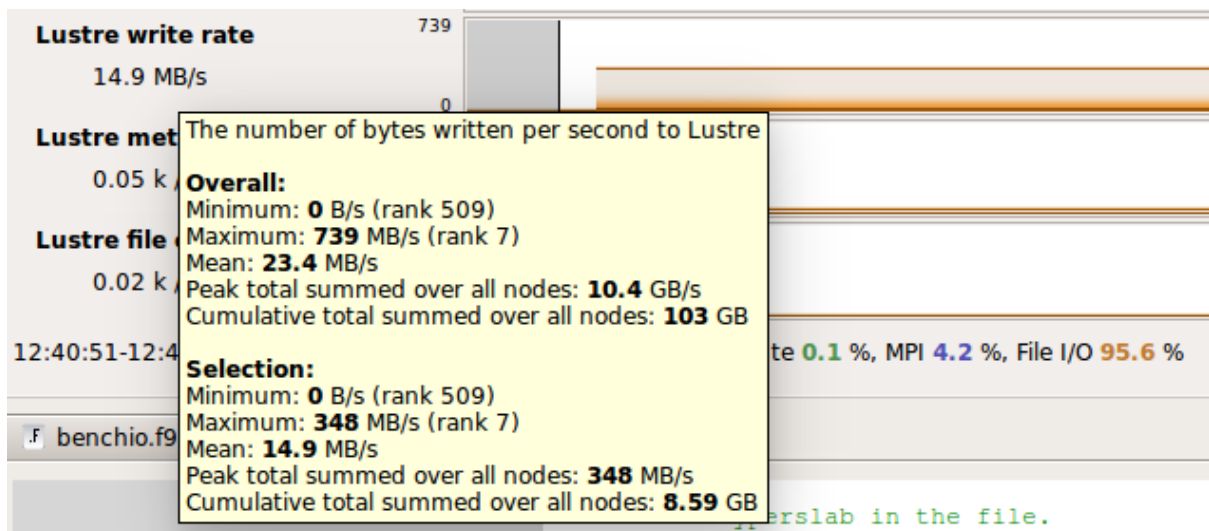System Selected



Default Stripe

- Activity recorded in writing doubles and closing the file is different for each of the files

- To zoom out at any point right click in the metrics view in the top-most pane of the GUI
- Are bytes transferred to Lustre when closing the file? (Hint: look in the 'Lustre write rate' metric in the top-most section of the GUI)
- Look at the code. Does the bandwidth reported take into account only the time when the program is writing data? (Hint: Search in benchio.f90 for where procedure `benchtime` is called)
- Estimate what the bandwidth of the I/O was across the whole of a call to hdf5write. We demonstrate how to do this for the write to file with a single stripe.
- Select the section of the profile from the first HDF5 write to the end of the first `h5fclose`



- Hover the mouse over the 'Lustre write rate' label to get a tool-tip up



- The label states that the metric is collected over all *nodes,* not over the processes. As there are 3 nodes used in this run (this information is displayed near the top of the GUI window) what is the combined bandwidth over all the nodes?
- How does this compare to the bandwidth reported in the benchmark output?

- How long does the closing of the file take for this file?
- How long does the write and the closing of the file take for the other HDF5 files?
- Does the time taken to close the file correlate to the number of stripes of the file? (Hint: file with the system defined striping has 48 stripes for this run, and default is 4 stripes)

## Further investigation

The following questions may be answered in your own time, using the profiles provided in TutorialData. You may also perform you own investigation and ask questions about other parts of the MAP profiles that have been provided. The questions posed here should make you think about what the profiles are displaying, and how to draw conclusions about the operation of the I/O routines on the ARCHER system at EPCC.

- Which is the fastest method of performing I/O?
  - For 8 processes
  - For 64 processes
  - For 512 processes
- Which method of I/O seems to be the most scalable (i.e. has the highest bandwidth going to scale)?
- Which method of I/O has the least amount of overhead associated with it? In this context I define overhead to be the amount of time spent in I/O routines in which no data is being written to the file-system.
- Does the overhead associated with using libraries (such as MPI I/O, HDF5 and NetCDF) scale with the number of processes used, the file size or the number of stripes contained in a file?
- What is the best method of parallel I/O to use in general? (Hint: This is a very tricky question to answer, and if you find a good answer then please publish it!)

## Resources for Allinea tools

If you are interested in getting to know the Allinea tool suite you may download our trial package. This is available from
https://s3-eu-west-1.amazonaws.com/com.allinea.resources/trial_package.tar.gz

Visit https://www.allinea.com/products/downloads/evaluation-resources for more resources getting started. If you would like to use the tools on your own applications you are free to do so.

The licence for the Allinea tools with which you have been supplied is valid until Tuesday 20th June 2017. If you would like to continue testing the Allinea tools after this date please request a trial from https://www.allinea.com/get-your-free-allinea-forge-and-allinea-performance-reports-trial