

Debugging and Optimizing Parallel Codes with Linaro Forge

Rudy Shand Field Application Engineer

linaro**forge**

Agenda

- 10:00 12:00 Session 1: Ensuring Program Correctness with Linaro DDT
- 13:00 15:00 Session 2: Performance Engineering with Linaro Performance Reports and Linaro MAP

Bug classification

- Crashes
 - One or more processes in application terminates
 - Most common and generally easiest to solve



- Deadlocks Stuck waiting for something that never happens
- Livelocks Making local progress, but no global progress

Race conditions

- One or more threads accessing the same data at the same time in non deterministic way
- Shows up as incorrect answer or sometimes crashes



Linaro Forge

An interoperable toolkit for debugging and profiling



The de-facto standard for HPC development

- Most widely-used debugging and profiling suite in HPC
- Fully supported by Linaro on Intel, AMD, Arm, Nvidia, AMD GPUs, etc.



State-of-the art debugging and profiling capabilities

- Powerful and in-depth error detection mechanisms (including memory debugging)
- Sampling-based profiler to identify and understand bottlenecks
- Available at any scale (from serial to exascale applications)



Easy to use by everyone

- Unique capabilities to simplify remote interactive sessions
- Innovative approach to present quintessential information to users

Linaro Forge: Where Most of Top Supercomputers turn for Performance Excellence

Build reliable and optimized code on multiple Server and HPC architectures

Linaro Forge combines



Linaro DDT

Market leading, simple to use HPC debugger for C/C++, Fortran and Python applications.



Linaro MAP

Effortless performance analysis for experts and novices alike.

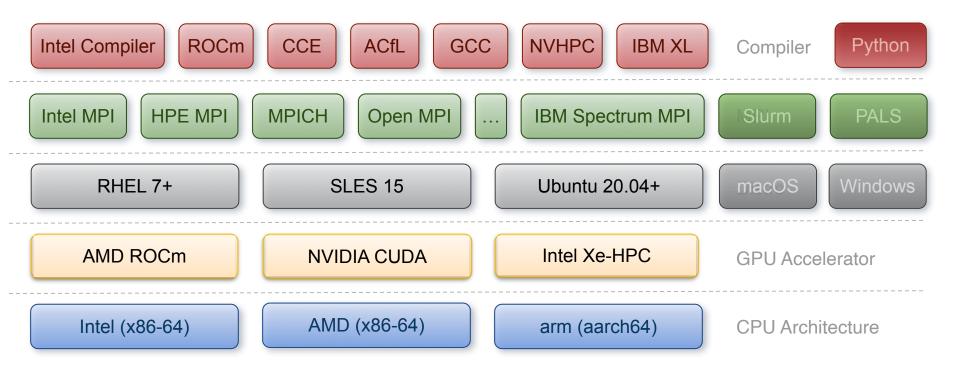


LinaroPerformance Reports

At a glance, single-page, application performance summary.

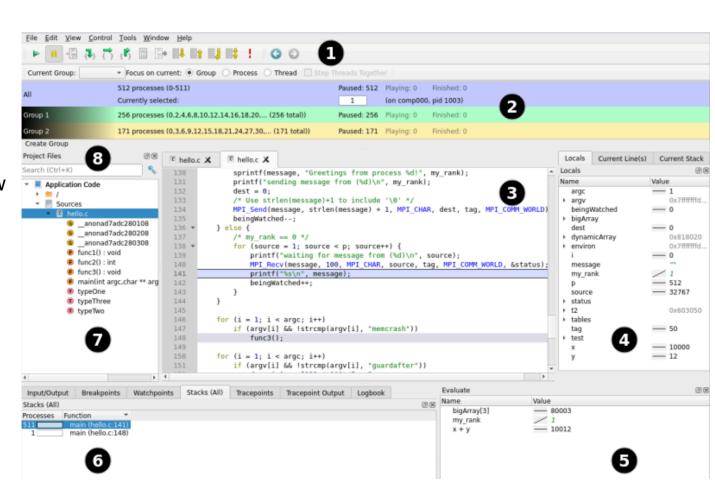
Performance
Engineering for
any architecture,
at any scale

Supported Platforms

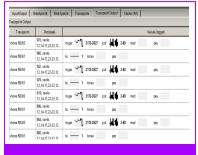


DDT UI

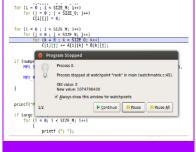
- 1 Process controls
- 2 Process groups
- 3 Source Code view
- 4 Variables
- 5 Evaluate window
- 6 Parallel Stack
- 7 Project files
- 8 Find a file or function



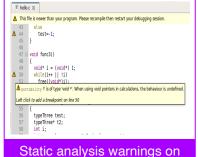
Linaro DDT Debugger Highlights



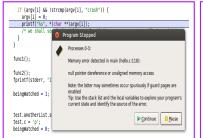
The scalable print alternative



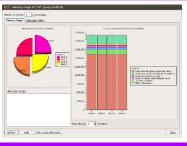
Stop on variable change



code errors



Detect read/write beyond array bounds



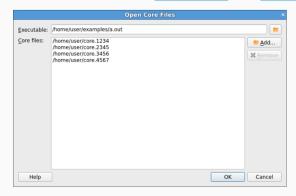
Detect stale memory allocations

Core files

You can open and debug one or more core files generated by your application.

Procedure

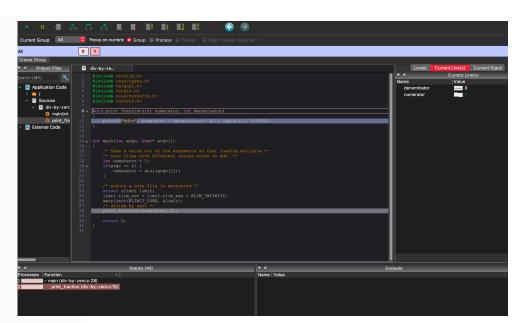
1. On the Welcome page click Open Core Files . The Open Core Files window opens.



2. Select an executable and a set of core files, then click **OK** to open the core files and start debugging them.

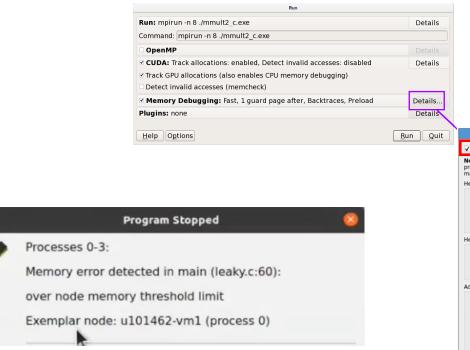
Note

While Linaro DDT is in this mode, you cannot play, pause, or step, because there is no process active. You are, however, able to evaluate expressions and browse the variables and stack frames saved in the core files.



- View core files for CPU's
- View core files for GPU's

Memory debugging menu in Linaro DDT



When manual linking is used, untick "Preload" box

✓ Preload the memory debugging library Language: Recommended Note: Preloading only works for programs linked against shared libraries. If your program is statically linked, you must relink it against the dmalloc library Heap Debugging Fast Balanced Thorough Custom Enabled Checks: basic More Information Heap Overflow/Underflow Detection Add guard pages to detect out of bounds heap access Advanced Set node memory threshold at 90 percent Check heap consistency every 100 heap operations ▼ Store stack backtraces for memory allocations Only enable for these processes: Help OK Cancel

Multi-dimensional Array Viewer

What does your data look like at runtime?

View arrays

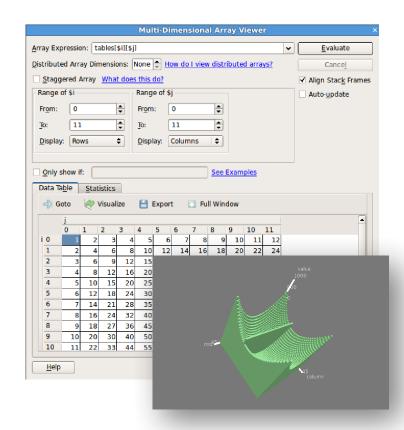
- On a single process
- Or distributed on many ranks

Use metavariables to browse the array

- Example: \$i and \$j
- Metavariables are unrelated to the variables in your program
- The bounds to view can be specified
- Visualise draws a 3D representation of the array

Data can also be filtered

 "Only show if": \$value>0 for example \$value being a specific element of the array



DDT: Production-scale debugging

Isolate and investigate faults at scale

Who misbehaved?

- Merge stacks from processes and threads
- Sparklines comparing data across processes
- Which MPI rank

Where is the problem?

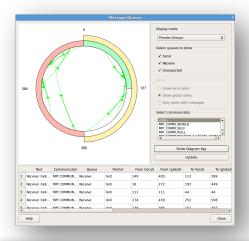
- Integrated source code editor
- Dynamic data structure visualization

How did it happen?

- Parse diagnostic messages
- Trace variables through execution

Why did it happen?

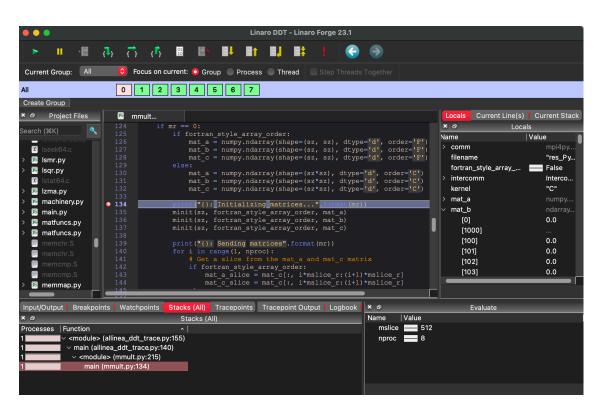
- Unique "Smart Highlighting"
- Experiment with variable values





Python Debugging

- Debug Features
 - Sparklines for Python variables
 - Tracepoints
 - MDA viewer
 - Mixed language support
- Improved Evaluations:
 - Matrix objects
 - Array objects
 - Pandas DataFrame
 - Series objects
- Python Specific:
 - Stop on uncaught Python exception
 - Show F-string variables
 - Mpi4py, NumPy, SciPy



Debugging Nvidia GPUs

Using Linaro DDT

Debug code simultaneously on the A100 GPU and the CPU

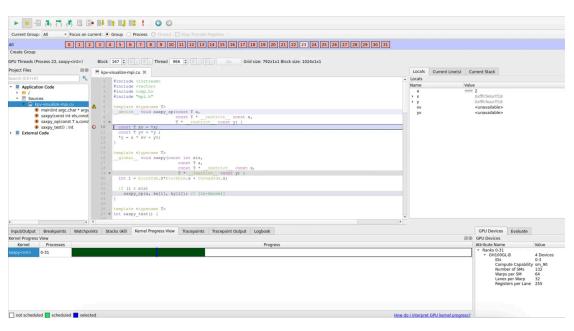
Controlling the GPU execution:

- All active threads in a warp will execute in lockstep.
 Therefore, DDT will step 32 threads at a time.
- Play/Continue runs all GPU threads
- Pause will pause a running kernel

Key (additional) GPU features:

- Kernel Progress View
- GPU thread in parallel stack view
- GPU Thread Selector
- GPU Device Pane

For NVIDIA's nvcc compiler, kernels must be compiled with the -g and -G flags



Run DDT in offline mode

Run the application under DDT and halt or report when a failure occurs

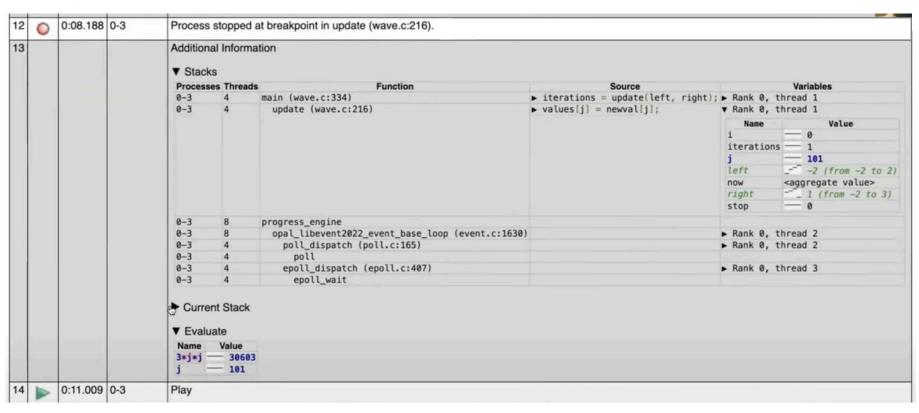
You can run the debugger in non-interactive mode

- For long-running jobs / debugging at very high scale
- For automated testing, continuous integration...

To do so, use following arguments:

- \$ ddt --offline --output=report.html mpirun ./jacobi_omp_mpi_gnu.exe
 - --offline enable non-interactive debugging
 - --output specifies the name and output of the non-interactive debugging session
 - Html
 - Txt
 - Add --mem-debug to enable memory debugging and memory leak detection

Report output



9 Step Guide

Optimizing high performance applications

Improving the efficiency of your parallel software holds the key to solving more complex research problems faster.

This pragmatic, 9 Step best practice guide, will help you identify and focus on application readiness, bottlenecks and optimizations one step at a time.

Cores

1/0

Discover lines of code

patterns.

spending a long time in I/O.

Trace and debug slow access

- Discover synchronization overhead and core utilization
- Synchronization-heavy code and implicit barriers are revealed

Validate corrections and optimal performance

Verification

- Understand numerical intensity and vectorization level.
 Hot loops, unvectorized code and
- Hot loops, unvectorized code an GPU performance reveleaed

Vectorization

Memory

- Reveal lines of code bottlenecked by memory access times.
- Trace allocation and use of hot data structure

Workloads

- Detect issues with balance.
- Slow communication calls and processes.
- Dive into partitioning code.

Communication

- Track communication performance.
- Discover which communication calls are slow and why.

Bugs

Correct application

Analyze before you optimize

- Measure all performance aspects.
 You can't fix what you can't see.
- Prefer real workloads over artificial tests.

Linaro Performance Reports

Characterize and understand the performance of HPC application runs



Gather a rich set of data

- Analyses metric around CPU, memory, IO, hardware counters, etc.
- Possibility for users to add their own metrics



Build a culture of application performance & efficiency awareness

- Analyses data and reports the information that matters to users
- Provides simple guidance to help improve workloads' efficiency



Relevant advice to avoid pitfalls

Adds value to typical users' workflows

- Define application behaviour and performance expectations
- Integrate outputs to various systems for validation (eg. continuous integration)
- Can be automated completely (no user intervention)

MPI

1/0

63.2%

16.2%

Linaro Performance Reports

A high-level view of application performance with "plain English" insights

Time spent in MPI calls. High values are usually bad.

This is high; check the MPI breakdown for advice on reducing it Time spent in filesystem I/O. High values are usually bad.

This is average; check the I/O breakdown section for optimization advice

0.0%

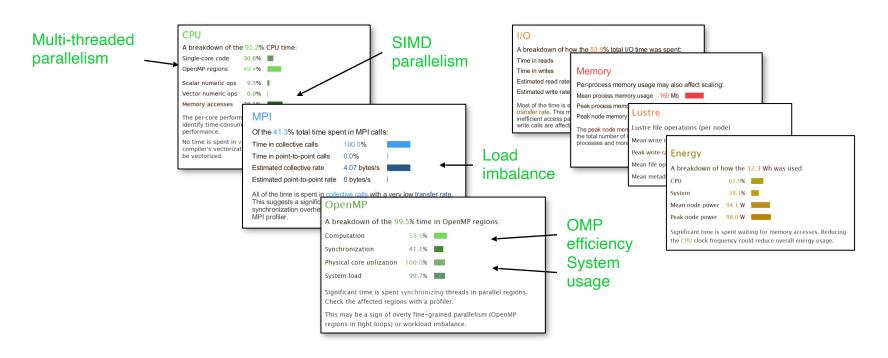
100.0%

1.38 MB/s

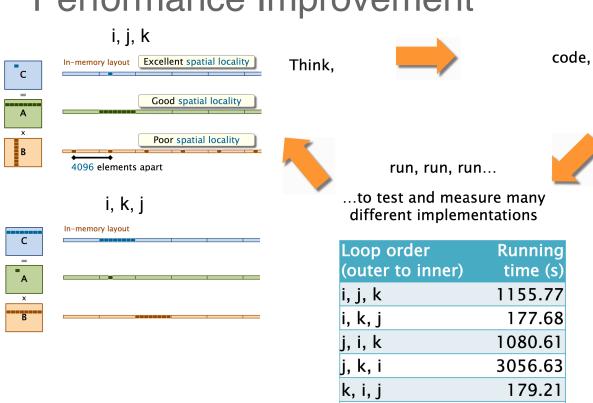
mpiexec.hvdra -host node-1.node-2 -map-by 1/0 socket -n 16 -ppn 8 ./Bin/low freq/../../Src//hydro Command: A breakdown of the 16.2% I/O time: ./Bin/low_freq/../../../Input/input_250x125_corner.nml 2 nodes (8 physical, 8 logical cores per node) Resources: Time in reads Memory: 15 GiB per node 16 processes, OMP_NUM_THREADS was 1 Tasks: Time in writes Machine: node-1 Effective process read rate 0.00 bytes/s Thu Jul 9 2015 10:32:13 Start time: 165 seconds (about 3 minutes) Total time: Effective process write rate Bin/../Src Most of the time is spent in write operations with a very low effective transfer rate. This may be caused by contention for the filesystem or inefficient access patterns. Use an I/O profiler to Summary: hydro is MPI-bound in this configuration investigate which write calls are affected. Time spent running application code. High values are usually good. Compute 20.6% This is very low; focus on improving MPI or I/O performance first

Linaro Performance Reports Metrics

Lowers expertise requirements by explaining everything in detail right in the report



Performance Improvement



k, j, i

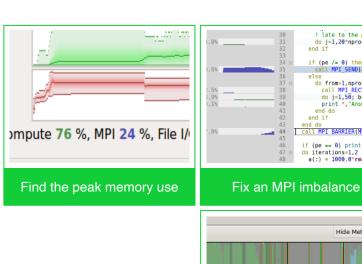
3032.82

```
i, j, k

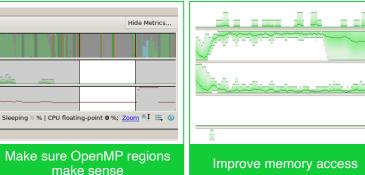
for (int i = 0; i < n; ++i) {
  for (int j = 0; j < n; ++j) {
    for (int k = 0; k < n; ++k) {
        C[i][j] += A[i][k] * B[k][j];
    }</pre>
```

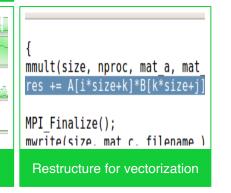
i, k, j for (int i = 0; i < n; ++i) { for (int k = 0; k < n; ++k) { for (int j = 0; j < n; ++j) { C[i][j] += A[i][k] * B[k][j];</pre>

Linaro MAP Source Code Profiler Highlights









MAP Capabilities

MAP is a sampling based scalable profiler

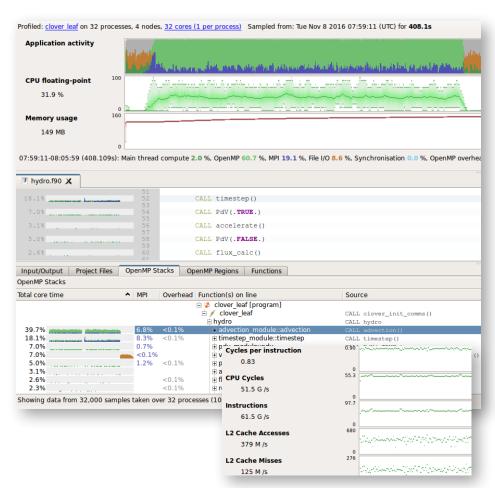
- Built on same framework as DDT
- Parallel support for MPI, OpenMP, CUDA
- Designed for C/C++/Fortran

Designed for 'hot-spot' analysis

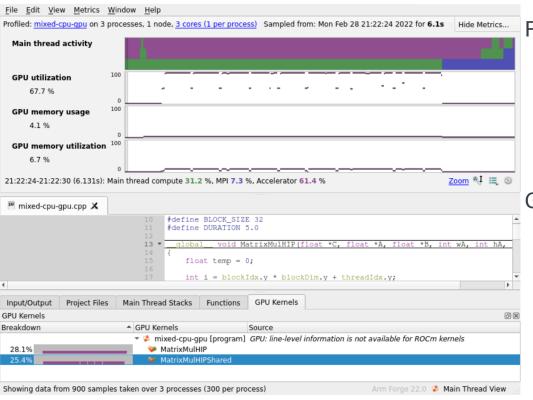
- Stack traces
- Augmented with performance metrics

Adaptive sampling rate

- Throws data away 1,000 samples per process
- Low overhead, scalable and small file size



GPU Profiling



Profile

- Supports both AMD and Nvidia GPUs
- Able to bring up metadata of the profile
- Mixed CPU [green] / GPU [purple] application
- CPU time waiting for GPU Kernels [purple]
- GPU Kernels graph indicating Kernel activity

GUI information

- GUI is consistent across platforms
- Zoom into main thread activity
- Ranked by highest contributors to app time

Python Profiling

19.0 adds support for Python

- Call stacks
- Time in interpreter

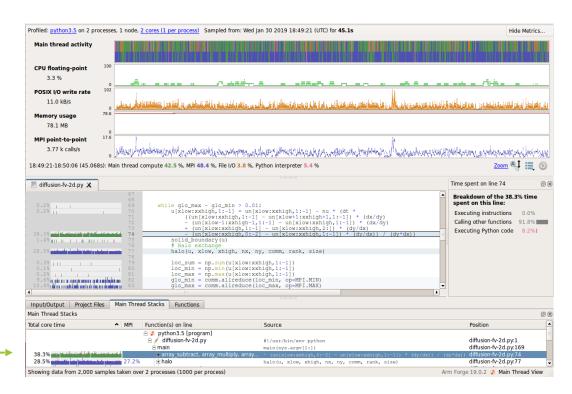
Works with MPI4PY

Usual MAP metrics

Source code view

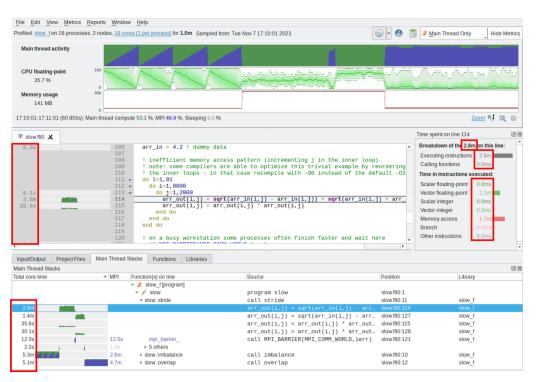
Mixed language support

Note: Green as operation is on numpy array, so backed by C routine, not ——Python (which would be pink)



map --profile mpirun -n 2 python ./diffusion-fv-2d.py

Toggle percentage-time and core-time in MAP

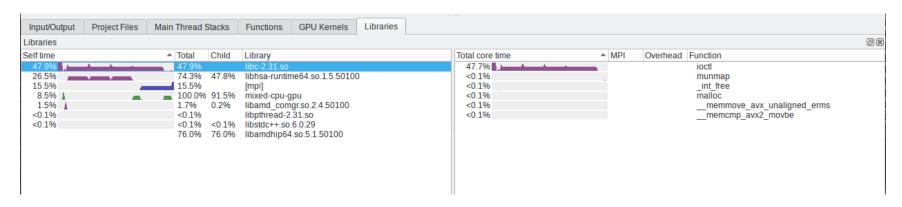


Use for direct comparisons between runs at the same scale (process/core counts).

- Easily determine if a change has made a portion of code faster, slower, or largely unchanged.
- Performance report automatically includes both percentage-time and core time
- Core-time is an estimation, but should be very close to the application run time

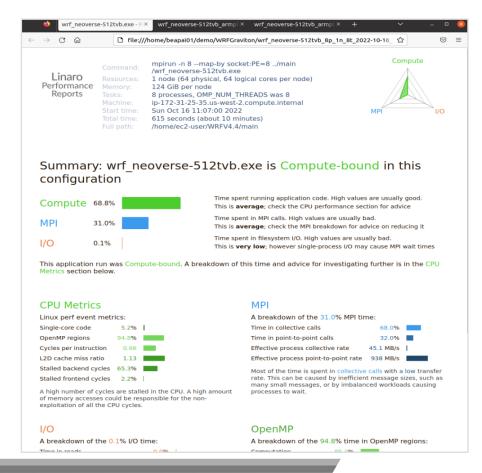
Libraries tab in MAP

- List time spent in shared libraries (left)
- List entry point functions into the selected library (right)

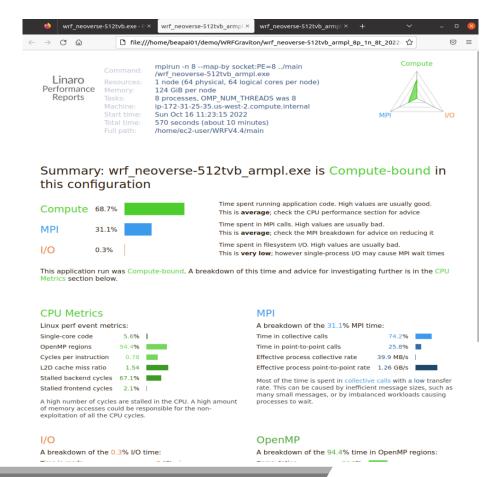


Use to identify the libraries that would benefit the most from optimisation or replacement (e.g. alternative maths library or memory management implementation).

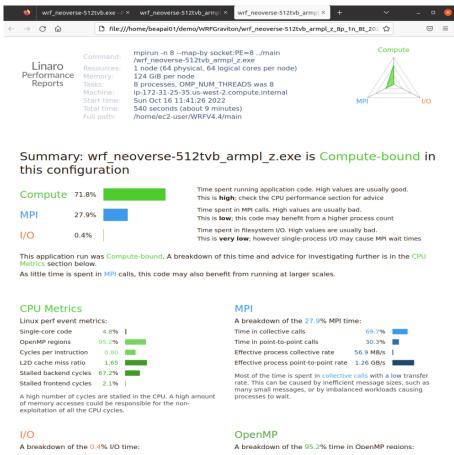
WRF build without enhancements



WRF with Arm Performance Libraries



WRF with Arm Performance Libraries and IO compression libraries

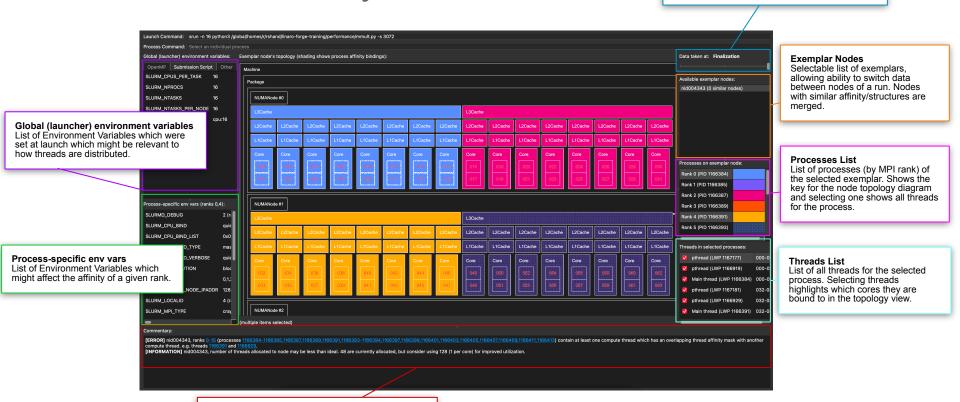


Linaro Forge

MAP Thread Affinity Advisor

Snapshot Selector

Change at which point of a run the Affinity data is shown (*Library Load, Initialisation, Finalization*).



Commentary

A list of commentary, providing information and advice on Memory Imbalance. Core Utilization etc.

