# Linaro Forge

Debugging and Optimising Parallel Codes

Rudy Shand - Field Application Engineer

Linaro Forge

# Agenda

- 10:00 - 11:00 Lecture on Debugging with DDT
- 11:00 - 12:30 DDT Debugging Hands-on session
- 12:30 - 13:00 Break
- 13:00 - 14:00 Lecture on Profiling with MAP
- 14:00 - 15:30 MAP Profiling Hands-on session
- 15:30 - 16:00 Break
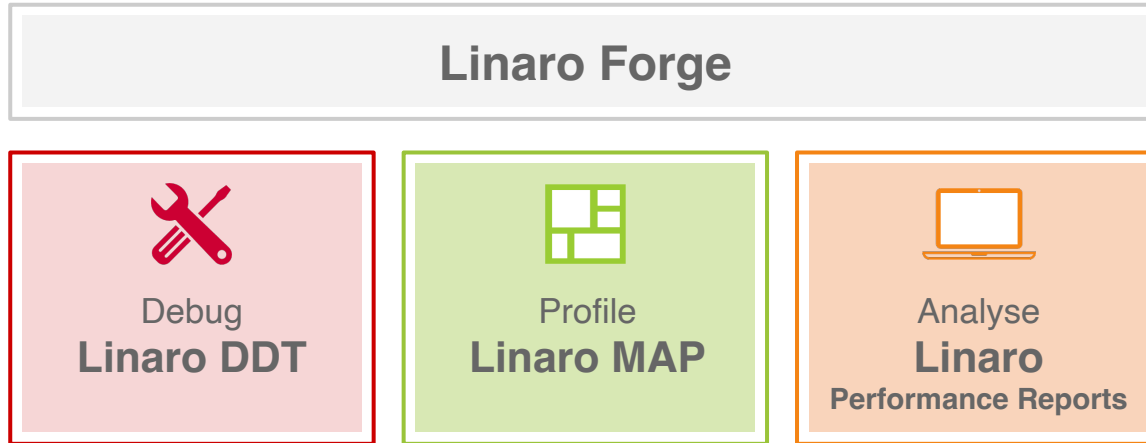- 16:00 - 17:00 Try DDT / MAP with own codes

# A little about myself

- Graduated from University of Reading
    - Cybernetics and Electronic Engineering
    - ML, Maths, Biology, Physics - No HPC

- Most working career in debuggers and performance tools
    - Arm DS-5 debugger and Streamline Performance Analyser
    - Compilers, Models, Embedded devices, mobile
    - In Embedded - both in a developer and quality role

- Joined the Arm Forge team (Now Linaro Forge)
    - Quality Lead / Field Application Engineer
    - 6 years in HPC, 11 years overall in debug and profiling tools



Linaro Forge

# HPC Development Solutions from Linaro

Best in class commercially supported tools for Linux and
high-performance computing (HPC)

| Linaro Forge |
|:---:|

| Debug<br>**Linaro DDT** | Profile<br>**Linaro MAP** | Analyse<br>**Linaro**<br>**Performance Reports** |
|:---:|:---:|:---:|

**Performance Engineering for any architecture, at any scale**

Linaro Forge

# 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

# Supported Platforms

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Intel Compiler | ROCm | CCE | ACfL | GCC | NVHPC | IBM XL | Compiler | Python |
| Intel MPI | HPE MPI | MPICH | Open MPI | … | IBM Spectrum MPI | | Slurm | PALS |
| RHEL 7+ | | SLES 15 | | Ubuntu 20.04+ | | | macOS | Windows |
| AMD ROCm | | | NVIDIA CUDA | | | | GPU Accelerator | |
| AMD/Intel (x86-64) | | Arm (AArch64) | | Power8 (ppc64le) | | | CPU Architecture | |

Linaro Forge

# Bug classification

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

- Hangs
  - 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

# Linaro DDT Debugger Highlights



The scalable print alternative

Stop on variable change

Static analysis warnings on code errors

Detect read/write beyond array bounds

Detect stale memory allocations

Linaro Forge

# 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.

**❶ 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

Linaro Forge

# Memory debugging menu in Linaro DDT



When manual linking is used, untick "Preload" box





Linaro Forge

# Multi-dimensional Array Viewer

## What does your data look like at runtime?

### View arrays
- On a single process
- Or distributed on many ranks
- Display the array values from tables[0:11][0:11]

### 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



Linaro Forge

# DDT: Production-scale debugging

## Isolate and investigate faults at scale

### Who misbehaved?
- Merge stacks from processes and threads
- Sparklines comparing data across processes
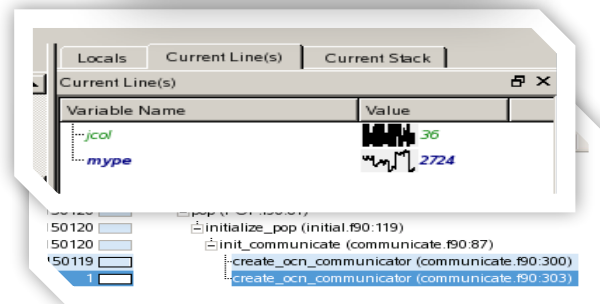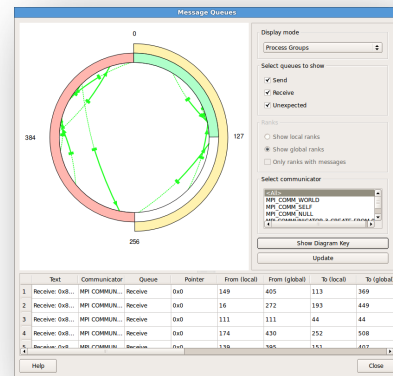- Which MPI rank

### Where is the problem?
- Integrated source code editor
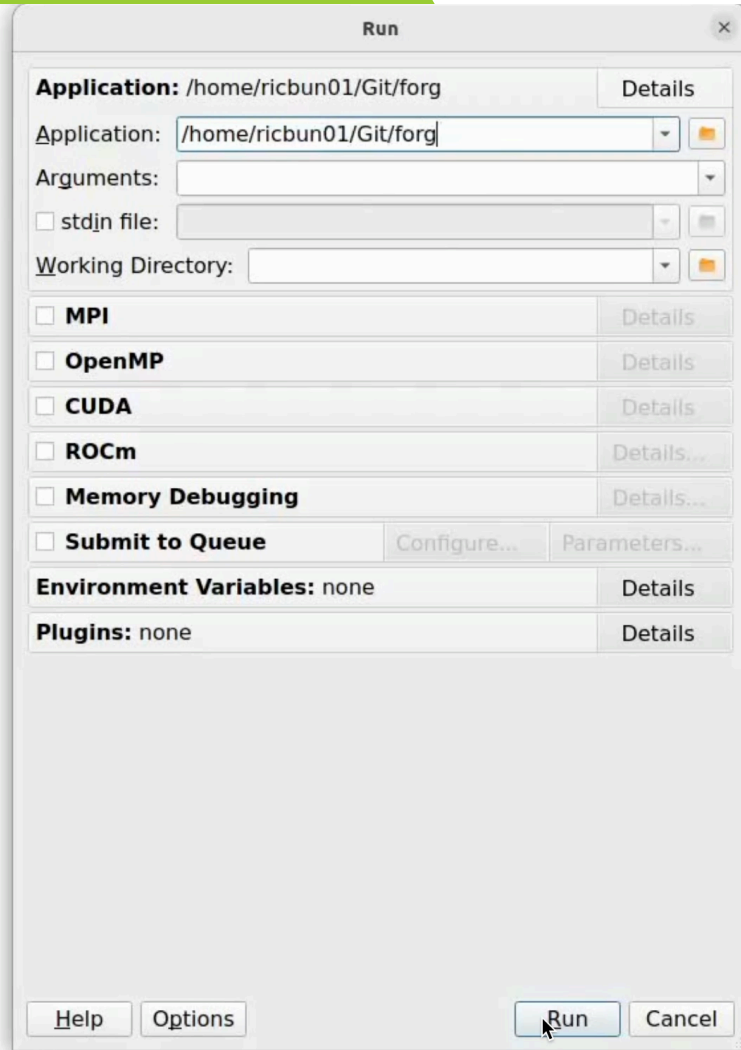- Dynamic data structure visualization

### How did it happen?
- Parse diagnostic messages
- Trace variables through execution

### Why did it happen?
- Unique "Smart Highlighting"
- Experiment with variable values



Linaro Forge

Starting a debug session

Linaro Forge

All

Create Group

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 |

Project Files | Fortran Modules

Project Files

Search (Ctrl+K)

- Application Code
  - /
  - Headers
  - Sources
    - hello.f
      - forge_constants
      - FUNC1() : INTEGER
      - hellof77
      - show_consts
      - SUB1
- External Code

hello.f

```
46        cORLD,ierr)
47            END IF
48        END IF
49
50        ALLOCATE(domain(my_rank*100000))
51        DO i = 1, SIZE(domain)
52            domain(i) = 2*i + 1
53        END DO
54
55        message="Hello From Me        !"
56
57        PRINT *,"My rank is ",my_rank,"!"
58
59        CALL SUB1()
60
61        beingwatched=1
62        tag=0
63
64        IF (my_rank.ne.0) THEN
65            PRINT *,"Greetings from process ",my_rank,"!"
66            PRINT *,"Sending message from ",my_rank,"!"
67            dest=0
68            CALL MPI_Send(message,21,MPI_CHARACTER,dest,tag,MPI_COMM_WORLD
69      c,ierr)
70            beingwatched=beingwatched-1
71        ELSE
72            message="Hello from my process"
73            DO source=1,(my_size-1)
74                PRINT *,"waiting for message from ",source
75                CALL MPI_Recv(message,21,MPI_CHARACTER,source,tag,MPI_COMM_WORLD
76      c,stat,ierr)
77                PRINT *,"Message recieved: ",message
78                ! beingwatched=beingwatched+1
79            END DO
80        END IF
81
82        beingwatched=12
83        CALL MPI_Finalize(ierr)
84        beingwatched=0
85        PRINT *,"All done...",my_rank
86        END
87
```

Locals | Current Line(s) | Current Stack

Locals

| Name | Value |
| --- | --- |
| mpi_argv_null | |
| mpi_argvs_null | |
| mpi_bottom | 0 |
| mpi_errcodes_i... | |
| mpi_in_place | 0 |
| mpi_status_ign... | |
| mpi_statuses_i... | |
| mpi_unweighted | |
| mpi_weights_e... | |
| beingwatched | 1 |
| dest | 0 |
| domain | |
| i | 1 |
| ierr | 0 |
| message | 'Hello From Me |
| my_rank | 0 |
| my_size | 20 |
| source | 0 |
| stat | 0 |
| tag | 0 |
| ompi_release_v... | 4 |
| ompi_minor_ve... | 1 |
| ompi_major_ve... | 3 |
| ompi_comm_ty... | 6 |
| ompi_comm_ty... | 7 |
| ompi_comm_ty... | 0 |
| ompi_comm_ty... | 5 |
| ompi_comm_ty... | 4 |
| ompi_comm_ty... | 3 |
| ompi_comm_ty... | 1 |
| ompi_comm_ty... | 9 |
| ompi_comm_ty... | 10 |
| ompi_comm_ty... | 2 |
| ompi_comm_ty... | 11 |
| ompi_comm_ty... | 8 |
| mpi_wtime_is_g... | 3 |
| mpi_win_unified | 0 |
| mpi_win_size | 8 |
| mpi_win_separ | 1 |

Evaluate

| Name | Value |
| --- | --- |

Input/Output | Breakpoints | Watchpoints | Stacks (All) | Tracepoints | Tracepoint Output | Logbook

Stacks (All)

| Processes | Threads | Function |
| --- | --- | --- |
| 19 | 19 | hellof77 (hello.f:65) |
| 1 | 1 | hellof77 (hello.f:72) |
| 20 | 40 | progress_engine |

Linaro Forge

# GPU Debugging



- Support both AMD and Nvidia GPUs
- Debug simultaneously on GPU and CPU

- Look and feel exactly the same
- Main Features work in GPU

- 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 -G flags

- ROCm GPU Debugging requires rocgdb to be available in your environment.

- For the hipcc compiler, kernels must be compiled with the -g flag

Linaro Forge

# 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 in "Current Line" display
  - Mpi4py, NumPy, SciPy

  ddt --connect mpirun -n 8 python3
  **%allinea_python_debug%** ./mmult.py

# 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**

```
ddt --offline -o jacobi_omp_mpi_gnu_debug.txt \
                        --trace-at _jacobi.F90:83,residual \
                        srun ./jacobi_omp_mpi_gnu.exe
```

Linaro Forge

# Report output

| 12 | ● | 0:08.188 | 0-3 | Process stopped at breakpoint in update (wave.c:216). |
|---|---|---|---|---|
| 13 | | | | Additional Information |

▼ Stacks

| Processes | Threads | Function | Source | Variables |
|---|---|---|---|---|
| 0-3 | 4 | main (wave.c:334) | ▶ iterations = update(left, right); | ▶ Rank 0, thread 1 |
| 0-3 | 4 | update (wave.c:216) | ▶ values[j] = newval[j]; | ▼ Rank 0, thread 1 |

| Name | Value |
|---|---|
| i | 0 |
| iterations | 1 |
| j | 101 |
| left | -2 (from -2 to 2) |
| now | &lt;aggregate value&gt; |
| right | 1 (from -2 to 3) |
| stop | 0 |

| Processes | Threads | Function | Source | Variables |
|---|---|---|---|---|
| 0-3 | 8 | progress_engine | | |
| 0-3 | 8 | opal_libevent2022_event_base_loop (event.c:1630) | | ▶ Rank 0, thread 2 |
| 0-3 | 4 | poll_dispatch (poll.c:165) | | ▶ Rank 0, thread 2 |
| 0-3 | 4 | poll | | |
| 0-3 | 4 | epoll_dispatch (epoll.c:407) | | ▶ Rank 0, thread 3 |
| 0-3 | 4 | epoll_wait | | |

▶ Current Stack

▼ Evaluate

| Name | Value |
|---|---|
| 3*j*j | 30603 |
| j | 101 |

| 14 | ▶ | 0:11.009 | 0-3 | Play |
|---|---|---|---|---|

Linaro Forge

# The Forge GUI and where to run it

DDT provides a powerful GUIs that can be run in a variety of configurations.



mydesktop

SSH

mycluster-login

qsub

Compute Nodes

Linaro Forge

# Hands on Setup

## Remote System

Host coolmuc2
    Hostname lxlogin1.lrz.de
    user <username>

/lrz/sys/courses/hlin1w23/linaro/linaro-forge-training.tar.gz

module load ddt/23.1.1

## Local Machine

Install Forge *https://www.linaroforge.com/downloadForge*

*Forge userguide*

Linaro Forge

# Remote connection to CoolMUC-2

# Explore a core file

# Hands on session

## System Info

*https://doku.lrz.de/coolmuc-2-11484376.html*
**CoolMUC-2**: 812 nodes:
- 28-core Intel Hazwell processor per node
- 64GB DDR4 memory per node
- cm2_tiny partition

*https://doku.lrz.de/running-parallel-jobs-on-the-linux-cluster-11484078.html*
*Interactive Session:*
- module load salloc_conf/cm2_tiny
- salloc -J linaro-hands-on --partition=cm2_tiny --time 00:30:00 --reservation=hlin1w23

*Scripting:*
- <linaro-forge-training>/slurm-coolmuc2.qtf
- <linaro-forge-training>/submit-job.sh

# Hands on session

## Build and run debug examples

# Use default Intel modules

*# build deadlock, simple and split programs*
cd <linaro-forge-training>/correctness/debug
make -f Makefile

*# run simple example with ddt*
ddt --connect mpiexec -n 4 ./simple

*# offline-debugging*
sbatch submit-job.sh

# Linaro Performance tools

## 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

**Commercially supported by Linaro**

### 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

**Accurate and Astute insight**

### 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)

**Relevant advice to avoid pitfalls**

Linaro Forge

# Linaro Performance Reports

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



Command: mpiexec.hydra –host node–1,node–2 –map–by socket –n 16 –ppn 8 ./Bin/low_freq/../../Src//hydro –i ./Bin/low_freq/../../../Input/input_250x125_corner.nml
Resources: 2 nodes (8 physical, 8 logical cores per node)
Memory: 15 GiB per node
Tasks: 16 processes, OMP_NUM_THREADS was 1
Machine: node–1
Start time: Thu Jul 9 2015 10:32:13
Total time: 165 seconds (about 3 minutes)
Full path: Bin/../Src

**I/O**

A breakdown of the 16.2% I/O time:

Time in reads                    0.0%
Time in writes                   100.0%
Effective process read rate      0.00 bytes/s
Effective process write rate     1.38 MB/s

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 investigate which write calls are affected.

Summary: hydro is MPI–bound in this configuration

Compute  20.6%
Time spent running application code. High values are usually good.
This is **very low**; focus on improving MPI or I/O performance first

MPI  63.2%
Time spent in MPI calls. High values are usually bad.
This is **high**; check the MPI breakdown for advice on reducing it

I/O  16.2%
Time spent in filesystem I/O. High values are usually bad.
This is **average**; check the I/O breakdown section for optimization advice

Linaro Forge

# Linaro Performance Reports Metrics

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

Multi-threaded parallelism

SIMD parallelism

Load imbalance

OMP efficiency System usage

**CPU**

A breakdown of the 91.2% CPU time:

| | | |
|---|---|---|
| Single-core code | 30.6% | |
| OpenMP regions | 69.4% | |
| Scalar numeric ops | 9.5% | |
| Vector numeric ops | 0.0% | |
| Memory accesses | 78.1% | |

The per-core perform... identify time-consum... performance.

No time is spent in v... compiler's vectoriza... be vectorized.

**MPI**

Of the 41.3% total time spent in MPI calls:

| | | |
|---|---|---|
| Time in collective calls | 100.0% | |
| Time in point-to-point calls | 0.0% | |
| Estimated collective rate | 4.07 bytes/s | |
| Estimated point-to-point rate | 0 bytes/s | |

All of the time is spent in collective calls with a very low transfer rate. This suggests a significa... synchronization overhe... MPI profiler.

**I/O**

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

| | |
|---|---|
| Time in reads | |
| Time in writes | |
| Estimated read rate | |
| Estimated write rate | |

Most of the time is s... transfer rate. This m... inefficient access pa... write calls are affect...

**Memory**

Per-process memory usage may also affect scaling:

| | | |
|---|---|---|
| Mean process memory usage | 160 Mb | |
| Peak process mem... | | |
| Peak node memory... | | |

The peak node mem... the total number of ... processes and more ...

**Lustre**

Lustre file operations (per node)

| | |
|---|---|
| Mean write r... | |
| Peak write ra... | |
| Mean file op... | |
| Mean metad... | |

**OpenMP**

A breakdown of the 99.5% time in OpenMP regions:

| | | |
|---|---|---|
| Computation | 58.9% | |
| Synchronization | 41.1% | |
| Physical core utilization | 100.0% | |
| System load | 99.7% | |

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.

**Energy**

A breakdown of how the 32.3 Wh was used:

| | | |
|---|---|---|
| CPU | 61.9% | |
| System | 38.1% | |
| Mean node power | 94.1 W | |
| Peak node power | 98.0 W | |

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

Linaro Forge

# The Performance Roadmap

**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.

## Verification
- Validate corrections and optimal performance

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

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

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

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

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

## I/O
- Discover lines of code spending a long time in I/O.
- Trace and debug slow access patterns.

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

## Bugs
- Correct application

Key : 
- Linaro Forge
- Linaro Performance Reports

LinaroForge

# Performance Improvement

### i, j, k



In-memory layout  Excellent spatial locality

Good spatial locality

Poor spatial locality

4096 elements apart

### i, k, j



In-memory layout

Think,

code,

run, run, run…

…to test and measure many
different implementations

| Loop order (outer to inner) | Running time (s) |
|---|---|
| i, j, k | 1155.77 |
| i, k, j | 177.68 |
| j, i, k | 1080.61 |
| j, k, i | 3056.63 |
| k, i, j | 179.21 |
| 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];
        }
    }
}
```

### 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];
        }
    }
}
```

Linaro Forge

# MAP Capabilities

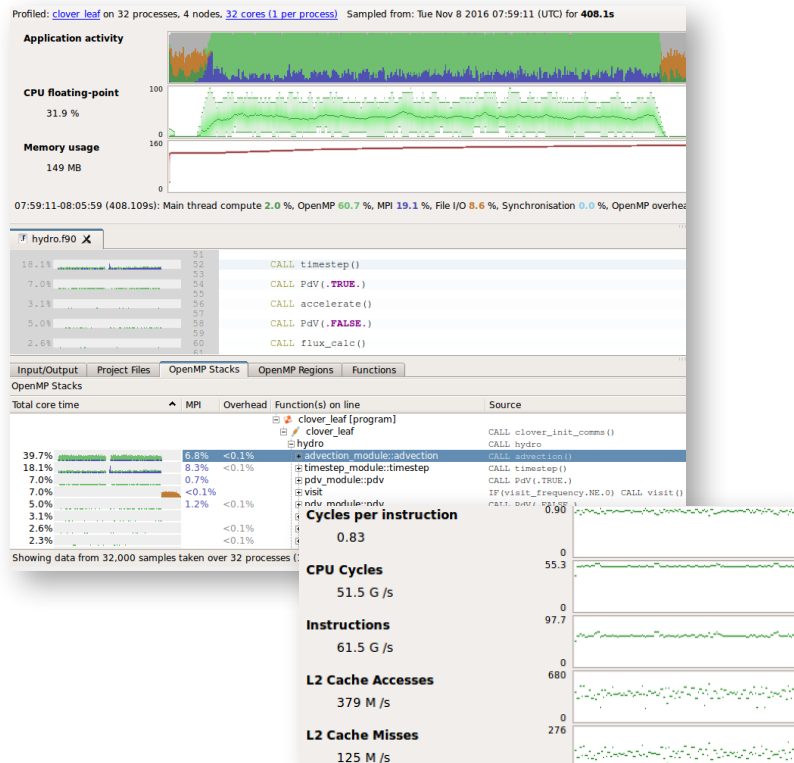MAP is a sampling based scalable profiler

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

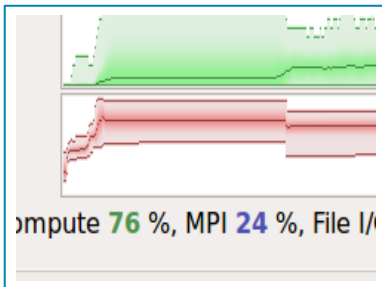Designed for 'hot-spot' analysis

- Stack traces
- Augmented with performance metrics
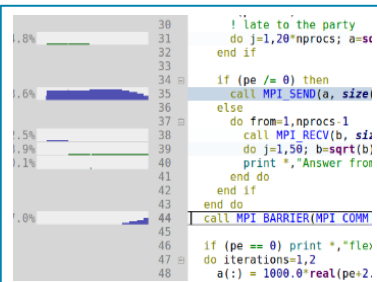
Adaptive sampling rate

- Throws data away - 1,000 samples per process
- Low overhead, scalable and small file size
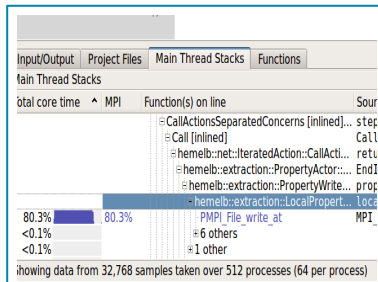
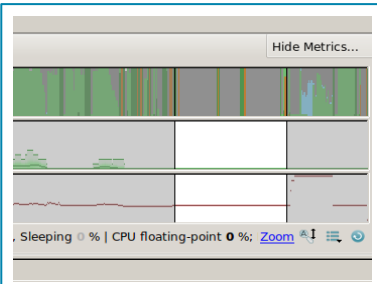# Linaro MAP Source Code Profiler Highlights
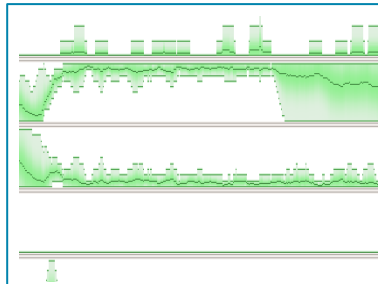

Find the peak memory use


Fix an MPI imbalance


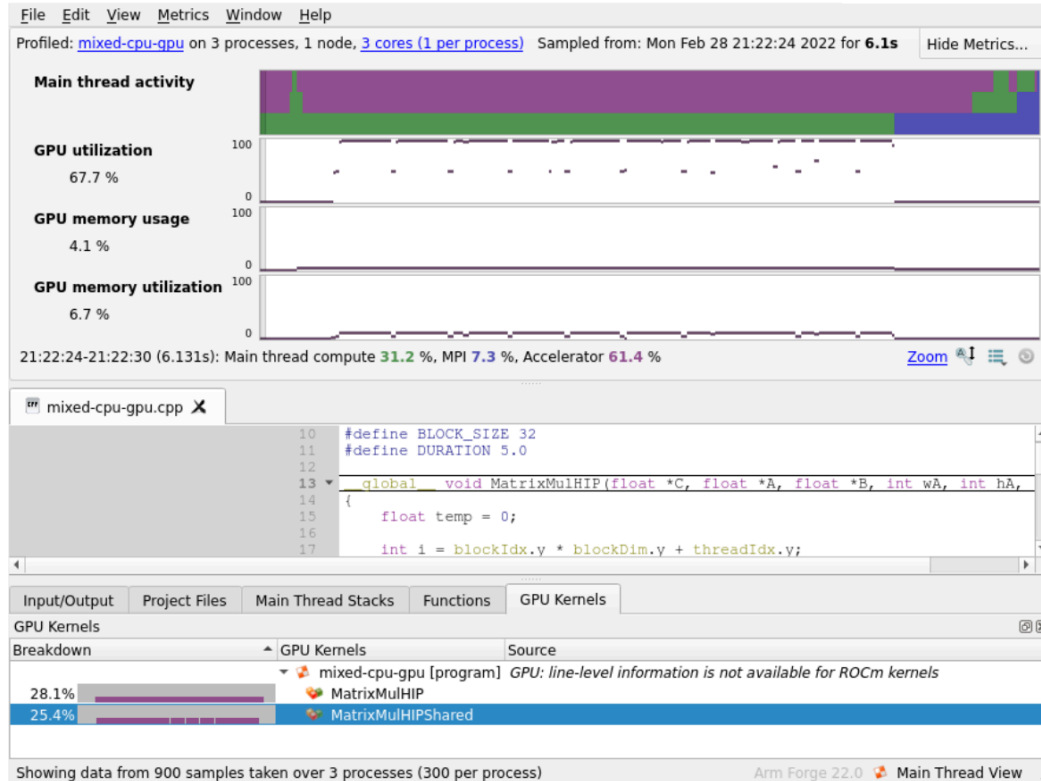Remove I?O bottleneck


Make sure OpenMP regions make sense


Improve memory access


Restructure for vectorization

Linaro Forge

# ROCm AMD GPU Profiling



## Profile

- Ran for 6s, taking 300 samples per process
- 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

Linaro Forge

# Python Profiling

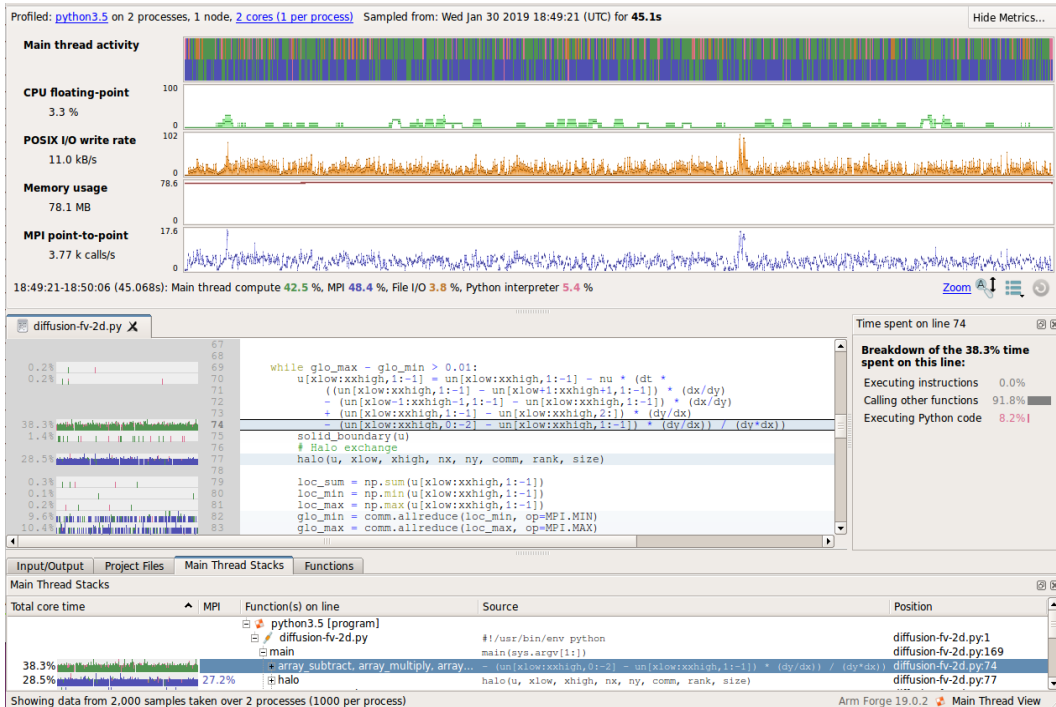19.0 adds support for Python
- Call stacks
- Time in interpreter

Works with MPI4PY
- Usual MAP metrics

Source code view
- Mixed language support

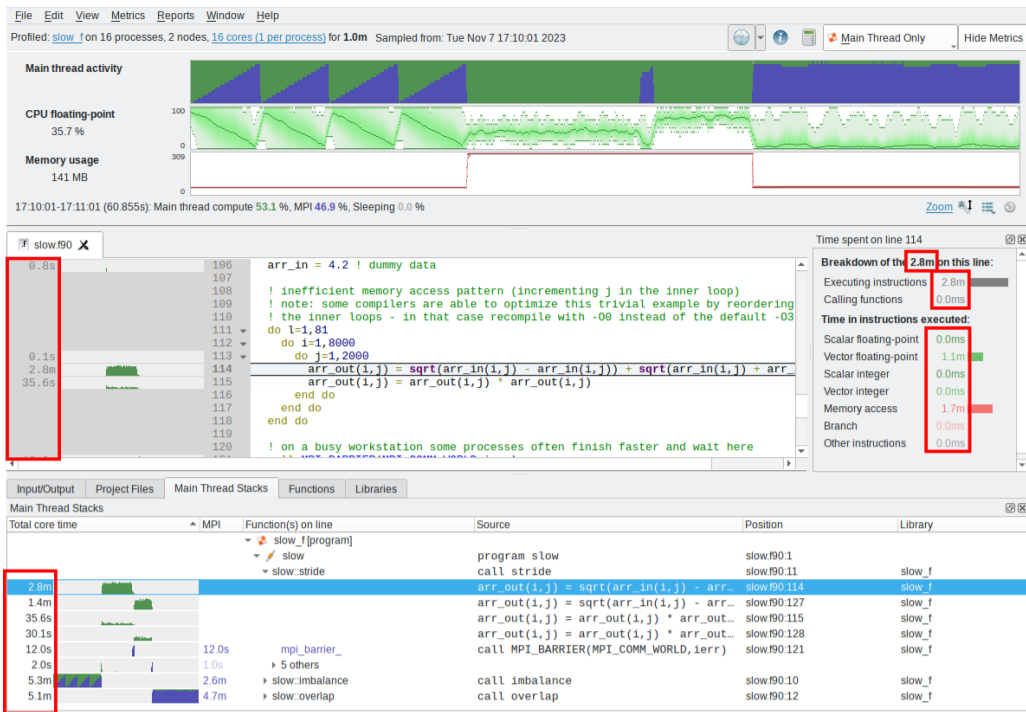**Note:  Green as operation is on numpy array, so backed by C routine, not Python (which would be pink)**



```
map --profile jsrun -n 2 python3 ./diffusion-fv-2d.py
```
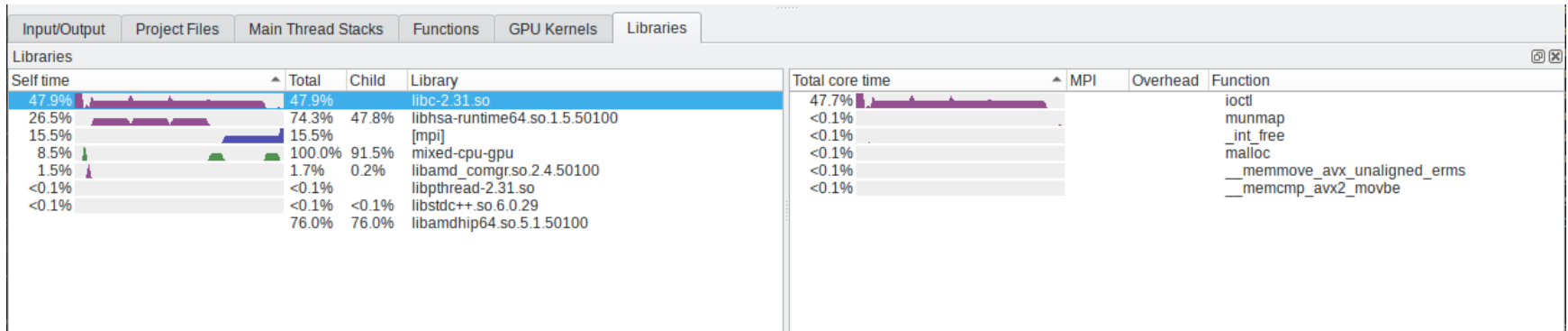
# Toggle percentage-time and core-time in MAP



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
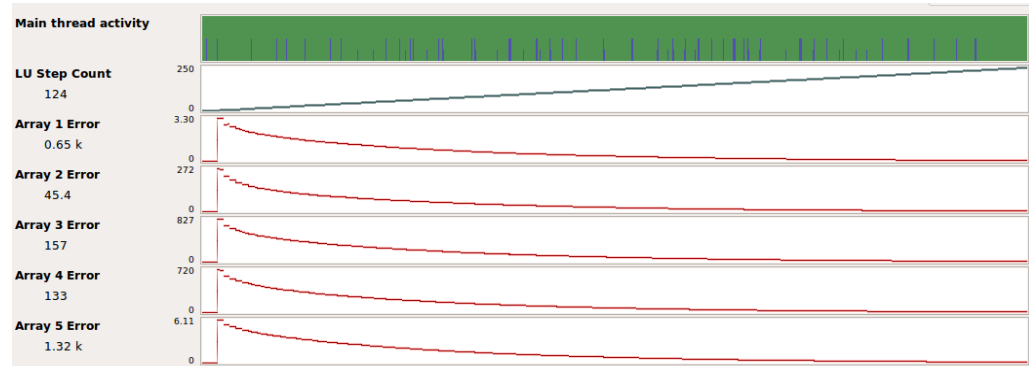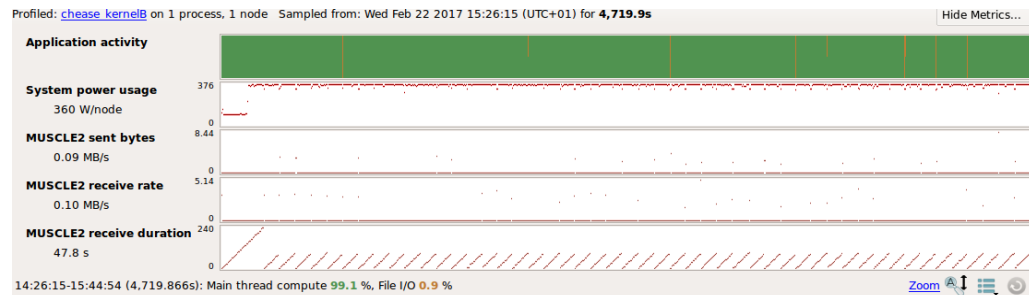
# Libraries tab in MAP

- 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).

# Custom metric example: MUSCLE2 & LU error terms

https://github.com/arm-hpc/custom-metrics



- Customized application instrumentation, eg, NPB LU
- Record error terms of solve
- Plot over time and step count for optimisation

Linaro Forge

# Matrix Multiplication example

## Build and run matrix multiplication example

https://docs.linaroforge.com/23.1.1/html/forge/worked_examples_appendix/mmult/analyze.html

```
# Build C and Fortran Examples
  export MPIF90=mpif90
  make -f mmult.makefile

# Build Python Examples
  module load python
  python -m venv run-mmult
  . run-mmult/bin/activate
  pip3 install numpy=='1.23.5' scipy mpi4py
  make -f mmult_py.makefile

# Debug using UI
  ddt --connect mpirun -n 8 ./mmult_c -s 3072
  ddt --connect python3 %allinea_python_debug% ./mmult.py -s 3072

# Offline profile
  sbatch submit-job.sh
```

Linaro Forge

# Thank you

~

[rudy.shand@linaro.org](mailto:rudy.shand@linaro.org)

Linaro Forge