# Linaro Forge

**Rudy Shand**
Field Application Engineer, Linaro Ltd

Argonne
NATIONAL LABORATORY

# Agenda

- **15 minutes Overview of DDT**
- **45 minutes DDT hands-on**
- **15 minutes Overview of MAP and Performance Reports**
- **30 minutes MAP and Performance Reports hands-on**

# DDT Supported Platforms

## Works across hardware architectures and HPC technologies

| 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 | Intel Xe-HPC | GPU Accelerator |

| Intel (x86-64) | AMD (x86-64) | arm (aarch64) | CPU Architecture |

ARGONNE
ATPESC2024
EXTREME-SCALE COMPUTING

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# DDT 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

# 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

# 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 mpiexec -n 8 python3
**%allinea_python_debug%** ./mmult.py

# 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 or 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 \
              mpiexec ./jacobi_omp_mpi_gnu.exe
```

# MAP and Performance Reports Supported Platforms

## Works across hardware architectures and HPC technologies

| 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 |

| Intel (x86-64) | AMD (x86-64) | arm (aarch64) | CPU Architecture |

# Linaro Performance tools

## Characterize and understand the performance of HPC application runs

**Commercially supported by Linaro**

Gather a rich set of data
- Analyses metric around CPU, memory, IO, hardware counters, etc.
- Possibility for users to add their own metrics

**Accurate and Astute insight**

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)

# The Performance Roadmap



**Optimizing high performance applications**

**Verification**
- Validate corrections and optimal performance

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

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

**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.
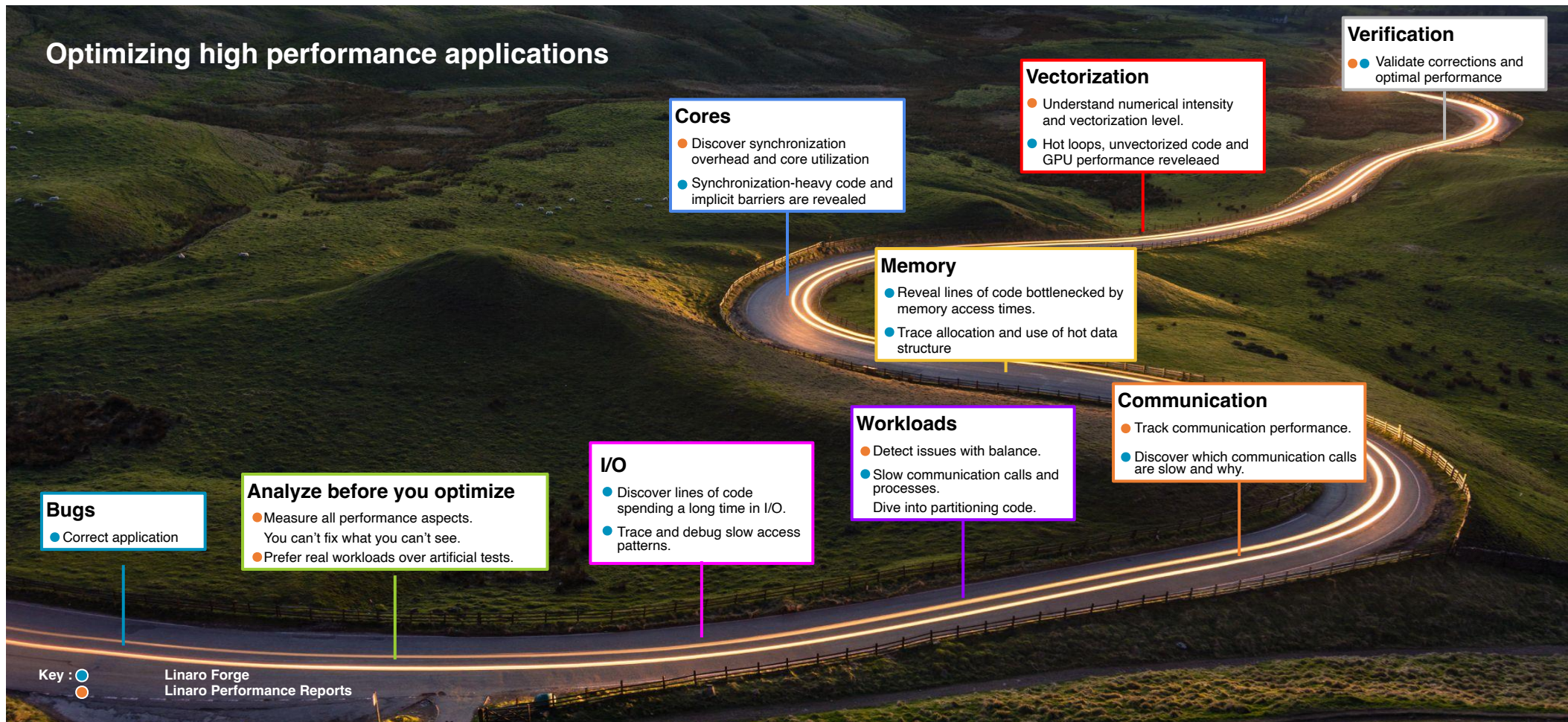  Dive into partitioning code.

**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

Argonne
NATIONAL LABORATORY

# Linaro Performance Reports

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



arm PERFORMANCE REPORTS

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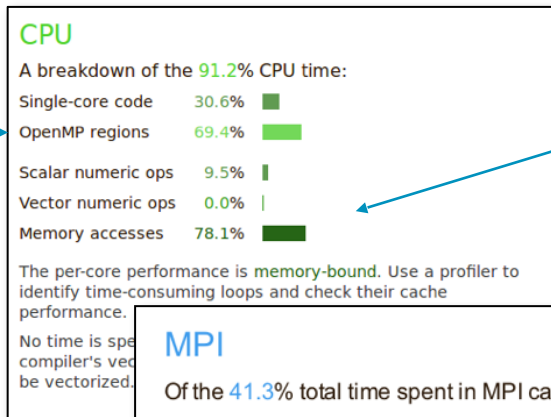
I/O   16.2%    Time spent in filesystem I/O. High values are usually bad.
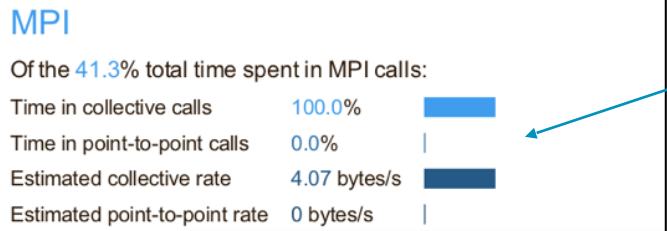This is **average**; check the I/O breakdown section for optimization advice

# Linaro Performance Reports Metrics

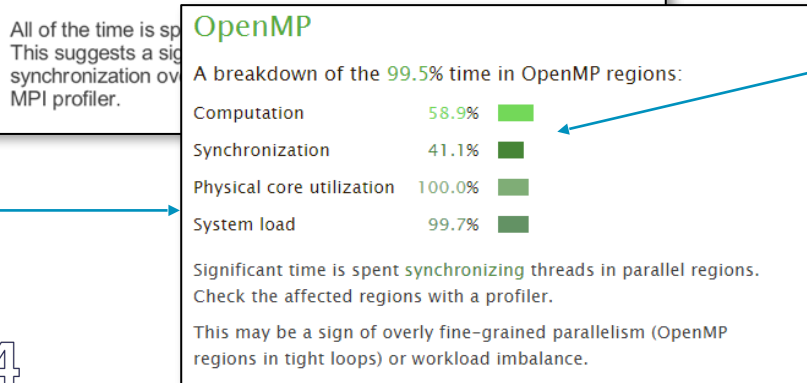## Lowers expertise requirements by explaining everything in detail right in the report
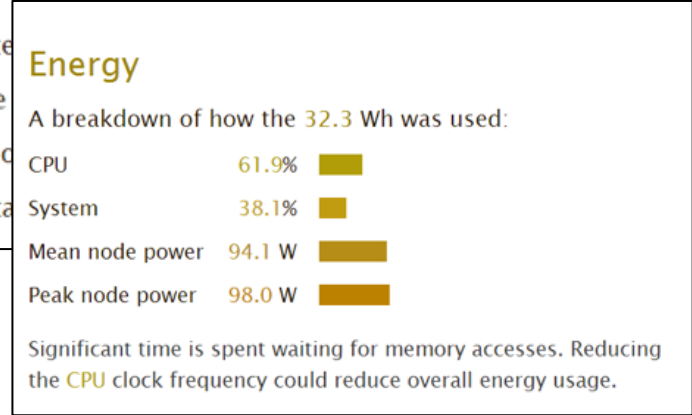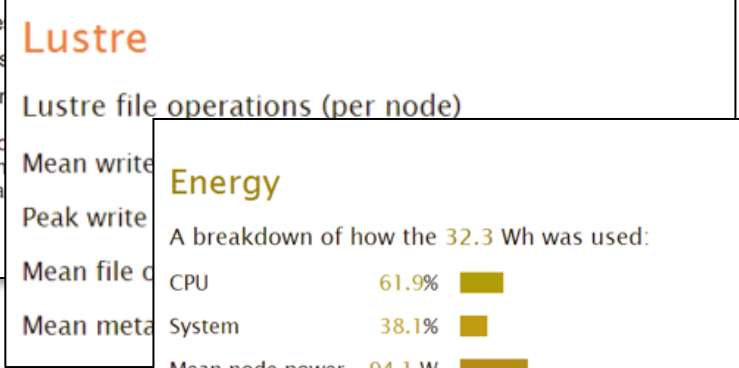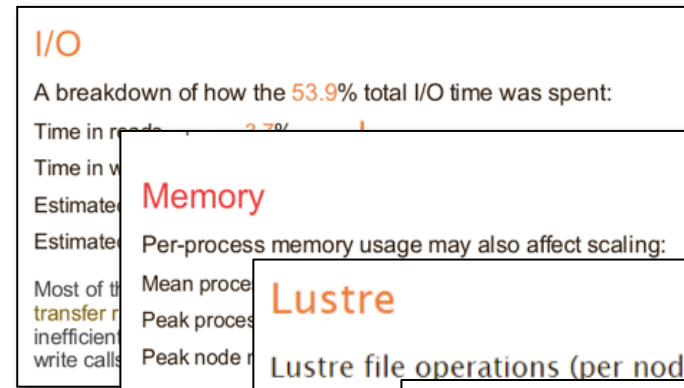
Multi-threaded parallelism

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 performance is memory-bound. Use a profiler to identify time-consuming loops and check their cache performance.

No time is spe... compiler's ve... 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 sp...
This suggests a sig...
synchronization ov...
MPI profiler.

### 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.

### I/O

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

Time in read...
Time in w...
Estimate...
Estimate...

Most of th...
transfer r...
inefficien...
write calls...

### Memory

Per-process memory usage may also affect scaling:

Mean proce...
Peak proce...
Peak node...

The peak no...
the total nu...
processes a...

### Lustre

Lustre file operations (per node)

Mean write...
Peak write...
Mean file o...
Mean meta...

### 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.

ARGONNE
ATPESC 2024
EXTREME-SCALE COMPUTING

extremecomputingtraining.anl.gov

Argonne
NATIONAL LABORATORY

# 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
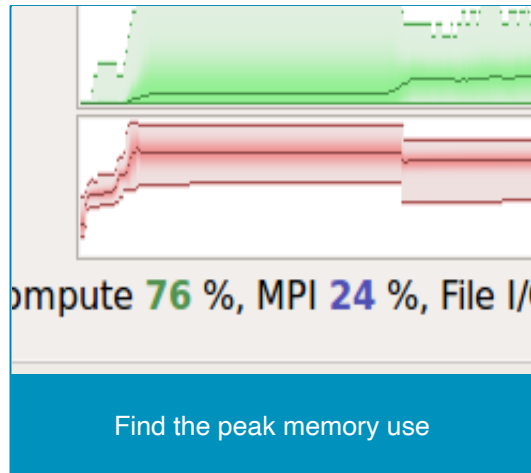
Designed for 'hot-spot' analysis
- Stack traces
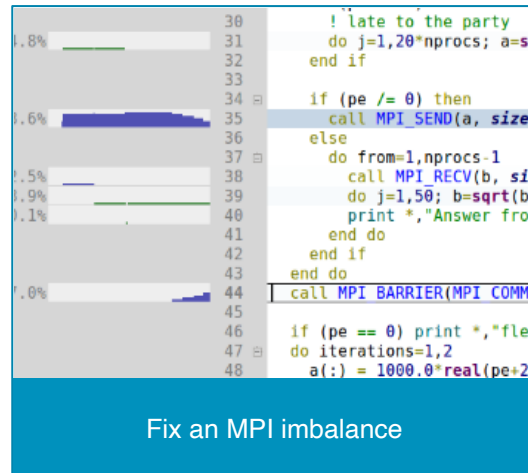- Augmented with performance metrics

Adaptive sampling rate
- Throws data away - 1,000 samples per process
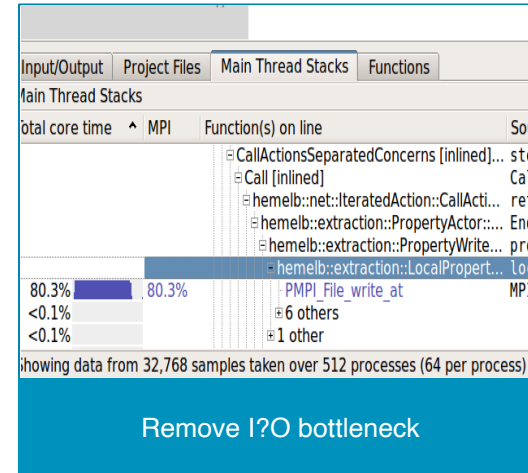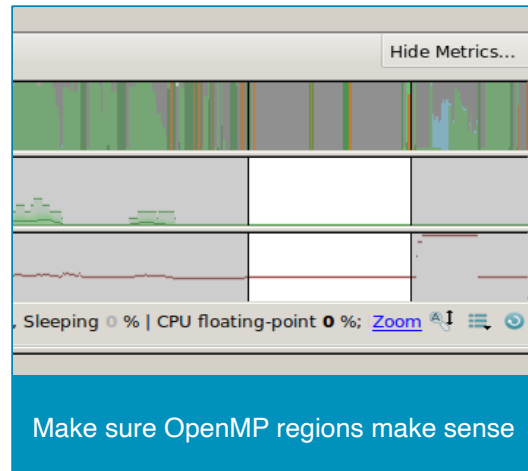- Low overhead, scalable and small file size
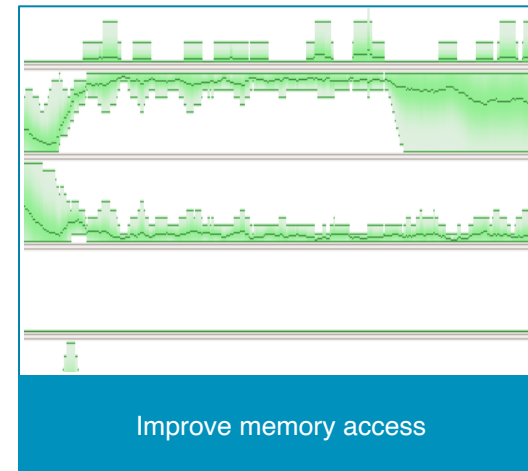
# MAP Highlights


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

# 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

extremecomputingtraining.anl.gov

# 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 mpiexec -n 2 python ./diffusion-fv-2d.py
```

# Compiler Remarks

Annotates source code with compiler remarks
- Remarks are extracted from the compiler optimisation report
- Compiler remarks are displayed as annotations next to your source code

Colour coded
- Their colour indicates the type of remark present in the following priority order:
  1. Red: failed or missed optimisations
  2. Green: successful or passed optimisations
  3. White: information or analysis notes

Compiler Remarks menu.
- Specify build directories for non-trivial build systems
- Filter out remarks

# Thank you

***Linaro Website***

www.linaro.org

***Linaro Forge Website***

www.linaroforge.com

***Contacts***

Sales: sales@forge.linaro.com

Support: support@forge.linaro.com