# Debugging and Profiling with DDT and Map

**ALCF Simulation Data and Learning Workshop**

October 2, 2019

Ryan Hulguin

ryan.hulguin@arm.com

**Agenda**

- General Debugging and Profiling Advice
- Arm Software for Debugging and Profiling
- Debugging with DDT
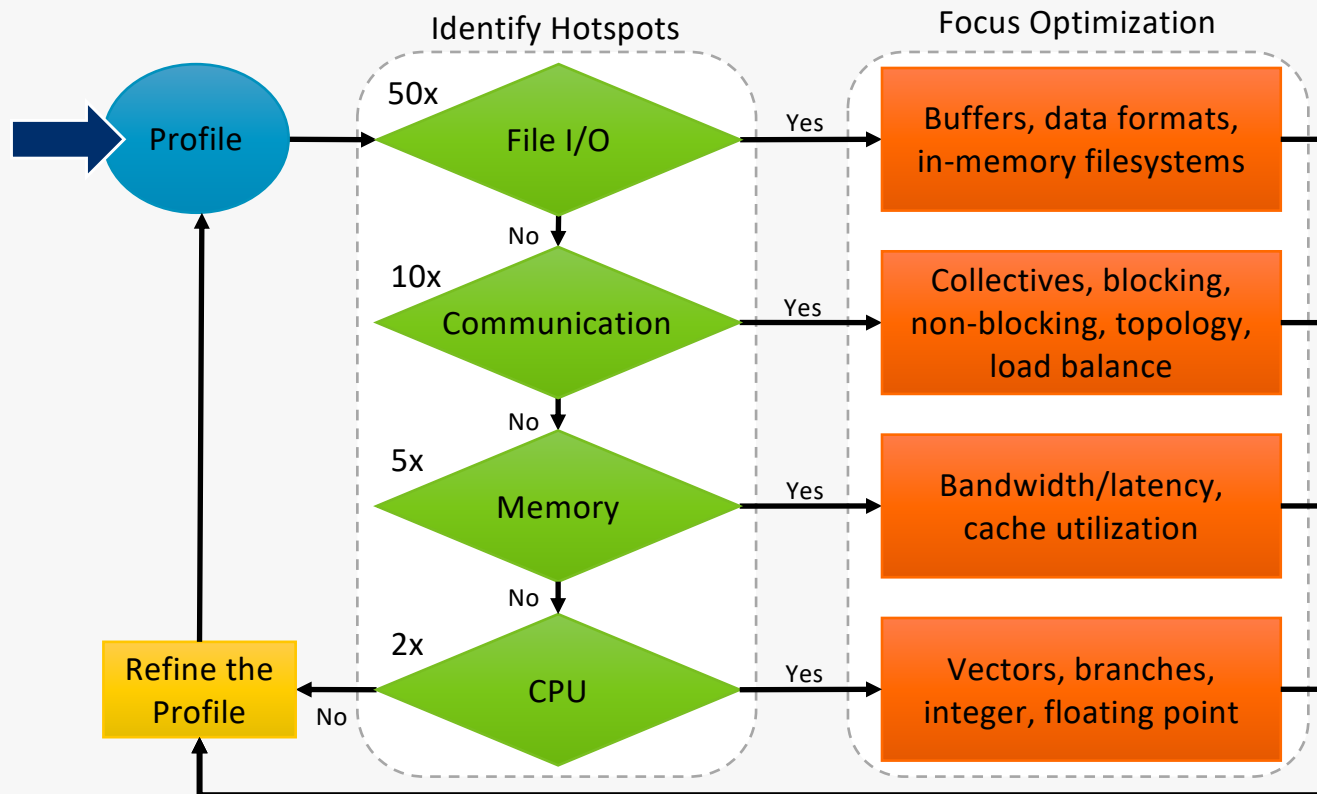- Profiling with MAP
- Theta Specific Settings

## Debugging

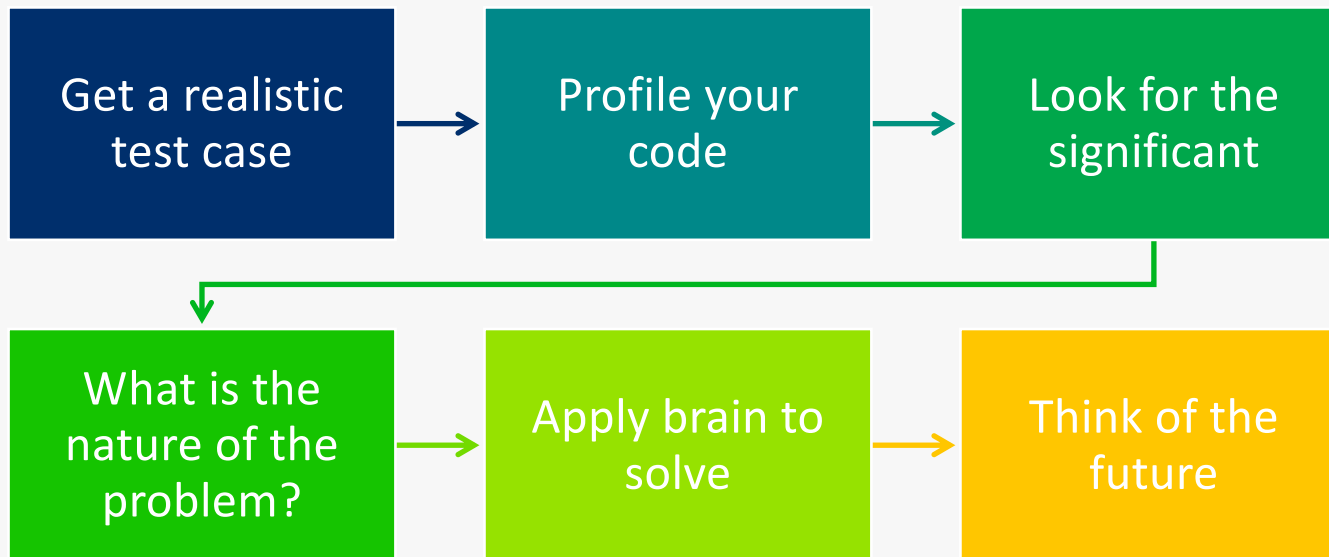Transforming a broken program to a working one
How? TRAFFIC!

–**T**rack the problem

–**R**eproduce

–**A**utomate - (and simplify) the test case

–**F**ind origins – where could the "infection" be from?

–**F**ocus – examine the origins

–**I**solate – narrow down the origins

–**C**orrect – fix and verify the test case is successful

# Profiling

Profiling is central to understanding and improving application performance.



Argonne Leadership Computing Facility

# Performance Improvement Workflow

```
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│                 │      │                 │      │                 │
│  Get a realistic│ ───▶ │  Profile your   │ ───▶ │  Look for the   │
│    test case    │      │     code        │      │   significant   │
│                 │      │                 │      │                 │
└─────────────────┘      └─────────────────┘      └─────────────────┘
        │◀──────────────────────────────────────────────────┘
        ▼
┌─────────────────┐      ┌─────────────────┐      ┌─────────────────┐
│                 │      │                 │      │                 │
│   What is the   │ ───▶ │  Apply brain to │ ───▶ │  Think of the   │
│  nature of the  │      │     solve       │      │     future      │
│    problem?     │      │                 │      │                 │
└─────────────────┘      └─────────────────┘      └─────────────────┘
```

Argonne NATIONAL LABORATORY

# Arm Software

arm

# Arm Forge

## An interoperable toolkit for debugging and profiling

The de-facto standard for HPC development

– Available on the vast majority of the Top500 machines in the world
– Fully supported by Arm on x86, IBM Power, Nvidia 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 parallel applications running at petascale)

## Easy to use by everyone

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

Argonne
NATIONAL LABORATORY

# Run and ensure application correctness

## Combination of debugging and re-compilation

- Ensure application correctness with
- Integrate with continuous integration system.
- Use version control to track changes and leverage Forge's built-in VCS support.

Examples:
```
$> ddt --offline aprun –n 48 ./example
$> ddt --connect aprun –n 48 ./example
```

# Visualize the performance of your application

- Measure all performance aspects with
- Identify bottlenecks and rewrite some code for better performance

Examples:
```
$> map --profile -n 48 ./example
```

# Debugging with DDT

arm

# Arm DDT – The Debugger

Who had a rogue behaviour ?

– Merges stacks from processes and threads

Where did it happen?

– leaps to source

How did it happen?

– Diagnostic messages

– Some faults evident instantly from source

Why did it happen?

– Unique "Smart Highlighting"

– Sparklines comparing data across processes

# Preparing Code for Use with DDT

As with any debugger, code must be compiled with the debug flag typically -g

It is recommended to turn off optimization flags i.e. -O0

Leaving optimizations turned on can cause the compiler to *optimize out* some variables and even functions making it more difficult to debug

# Segmentation Fault

In this example, the application crashes with a segmentation error outside of DDT.



What happens when it runs under DDT?

# Segmentation Fault in DDT



DDT takes you to the exact line where Segmentation fault occurred, and you can pause and investigate

# Invalid Memory Access



The array tab is a 13x13 array, but the application is trying to write a value to tab(4198128,0) which causes the segmentation fault.

`i` is not used, and `x` and `y` are not initialized

# It works… Well, most of the time



**SCHRODIN BUG**

A strange behaviour where the application "sometimes" crashes is a typical sign of a memory bug

Arm DDT is able to force the crash to happen

# Advanced Memory Debugging

# Heap debugging options available

**Fast**

### basic
- Detect invalid pointers passed to memory functions (e.g. malloc, free, ALLOCATE, DEALLOCATE,...)

### check-fence
- Check the end of an allocation has not been overwritten when it is freed.

### free-protect
- Protect freed memory (using hardware memory protection) so subsequent read/writes cause a fatal error.

### Added goodiness
- Memory usage, statistics, etc.

**Balanced**

### free-blank
- Overwrite the bytes of freed memory with a known value.

### alloc-blank
- Initialise the bytes of new allocations with a known value.

### check-heap
- Check for heap corruption (e.g. due to writes to invalid memory addresses).

### realloc-copy
- Always copy data to a new pointer when re-allocating a memory allocation (e.g. due to realloc)

**Thorough**

### check-blank
- Check to see if space that was blanked when a pointer was allocated/freed has been overwritten.

### check-funcs
- Check the arguments of addition functions (mostly string operations) for invalid pointers.

*See user-guide:*

*Chapter 12.3.2*

Argonne
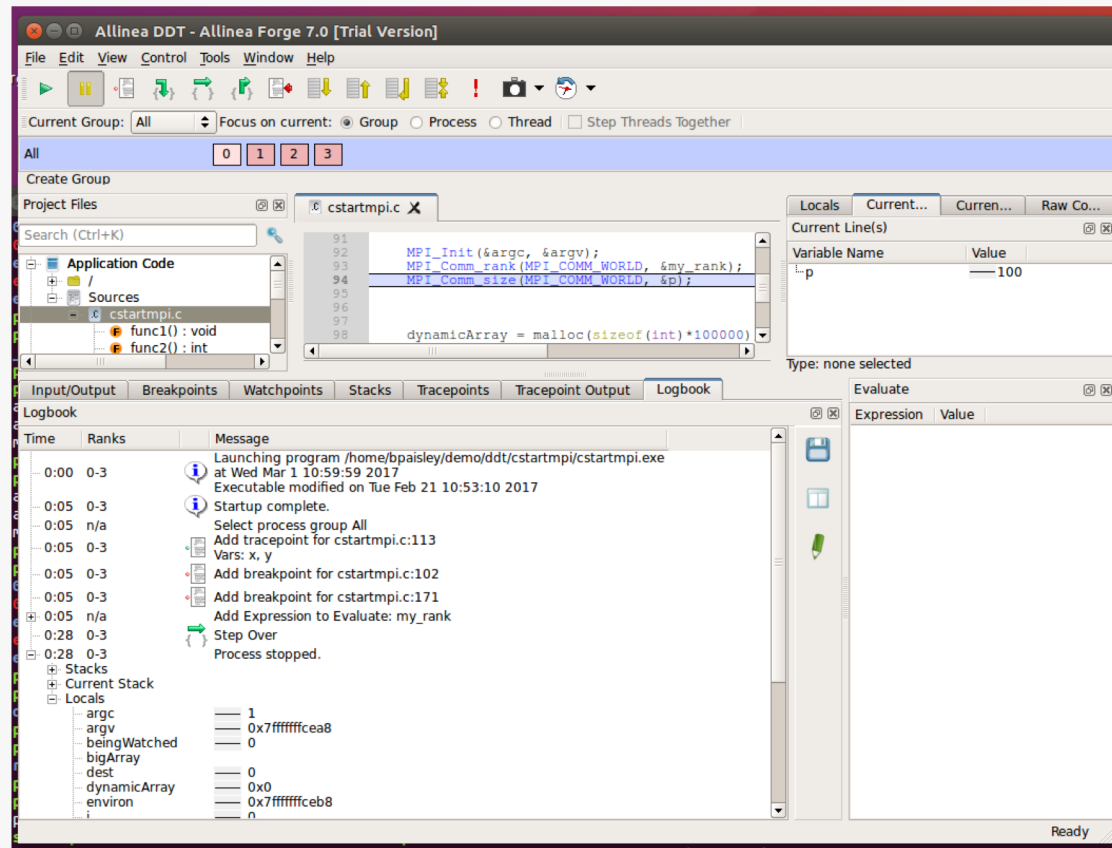NATIONAL LABORATORY

# Guard pages (aka "Electric Fences")



- **A powerful feature…:**
  - Forbids read/write on guard pages throughout the whole execution

    *(because it overrides C Standard Memory Management library)*

- **… to be used carefully:**
  - Kernel limitation: up to 32k guard pages max ( "mprotect fails" error)
  - Beware the additional memory usage cost

# New Bugs from Latest Changes

# Track Your Changes in a Logbook

# Arm DDT Demo

arm

# Five great things to try with Allinea DDT

The scalable print alternative

Stop on variable change

Static analysis warnings on code errors

Detect read/write beyond array bounds

Detect stale memory allocations

Argonne
NATIONAL LABORATORY

# Arm DDT cheat sheet

Load the environment module
- – $ module load **forge/19.1.2**
- – $ module unload **xalt**

Prepare the code
- – $ cc **-O0 -g** myapp.c -o myapp.exe

Start Arm DDT in interactive mode
- – $ **ddt** aprun -n 8 ./myapp.exe arg1 arg2

Or use the reverse connect mechanism
- – On the login node:
  - $ ddt &
- – (or use the remote client) **<- Preferred method**
- – Then, edit the job script to run the following command and submit:
  - **ddt --connect** aprun -n 8 ./myapp.exe arg1 arg2

Argonne ▲
NATIONAL LABORATORY

# Profiling with MAP

arm

# Arm MAP – The Profiler

Small data files

<5% slowdown

No instrumentation

No recompilation

# Glean Deep Insight from our Source-Level Profiler

| | |
|---|---|
| **Memory usage** (M)<br>6.1  -  62.6  ( 29.3 avg ) | |
| **MPI call duration** (ms)<br>0  -  727.9  ( 136.2 avg ) | |
| **MPI point-to-point** (/s)<br>0  -  217  ( 1.0 avg ) | |
| **MPI collectives** (/s)<br>0  -  172  ( 0.6 avg ) | |
| **CPU memory access** (%)<br>0  -  100  ( 18.7 avg ) | |
| **CPU floating-point** (%)<br>0  -  100  ( 10.9 avg ) | |
| **CPU vector** (%)<br>0  -  100  ( 37.6 avg ) | |
| **CPU branch** (%)<br>0  -  60  ( 0.5 avg ) | |

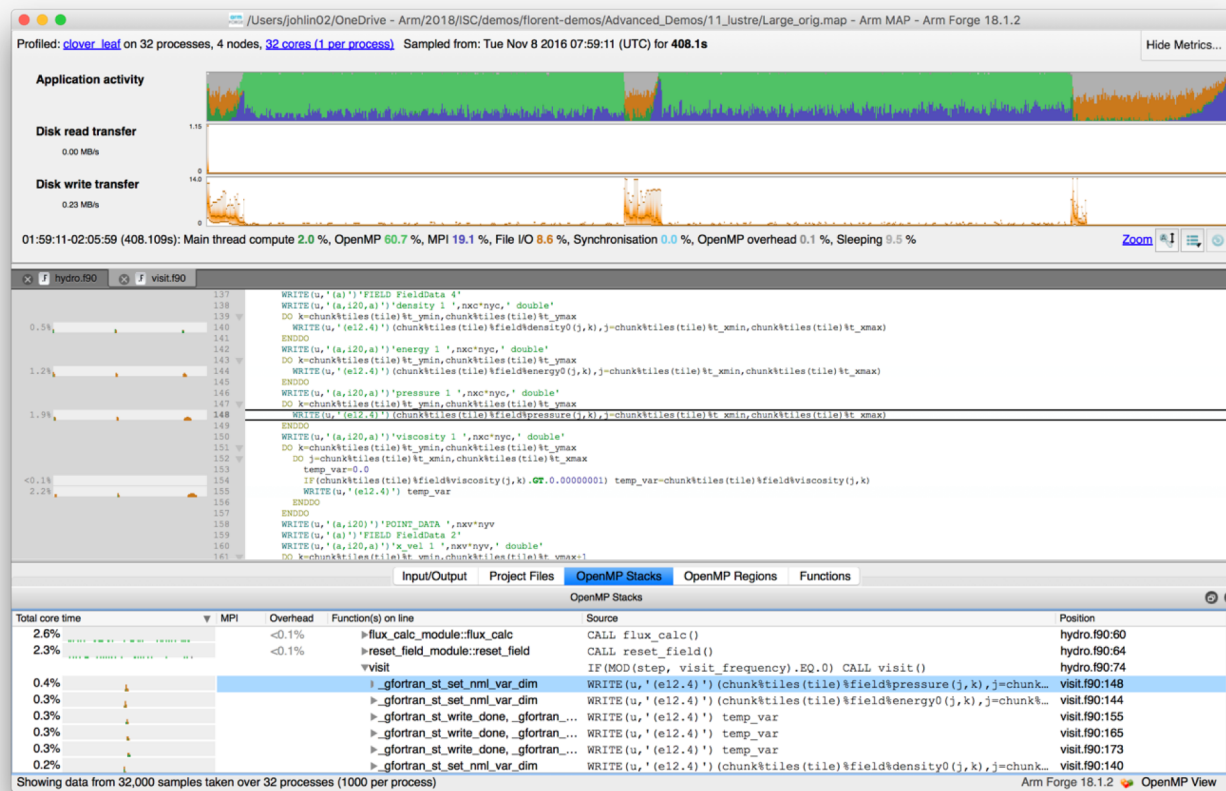**Track memory usage across the entire application over time**

**Spot MPI and OpenMP imbalance and overhead**

**Optimize CPU memory and vectorization in loops**

**Detect and diagnose I/O bottlenecks at real scale**
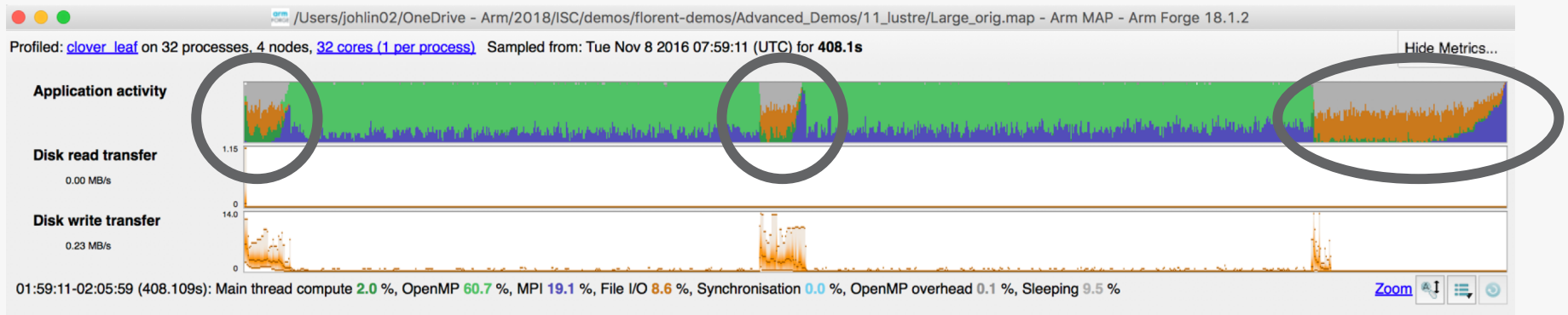
Argonne
NATIONAL LABORATORY

# Initial profile of CloverLeaf shows surprisingly unequal I/O

Each I/O operation should take about the same time, but it's not the case.

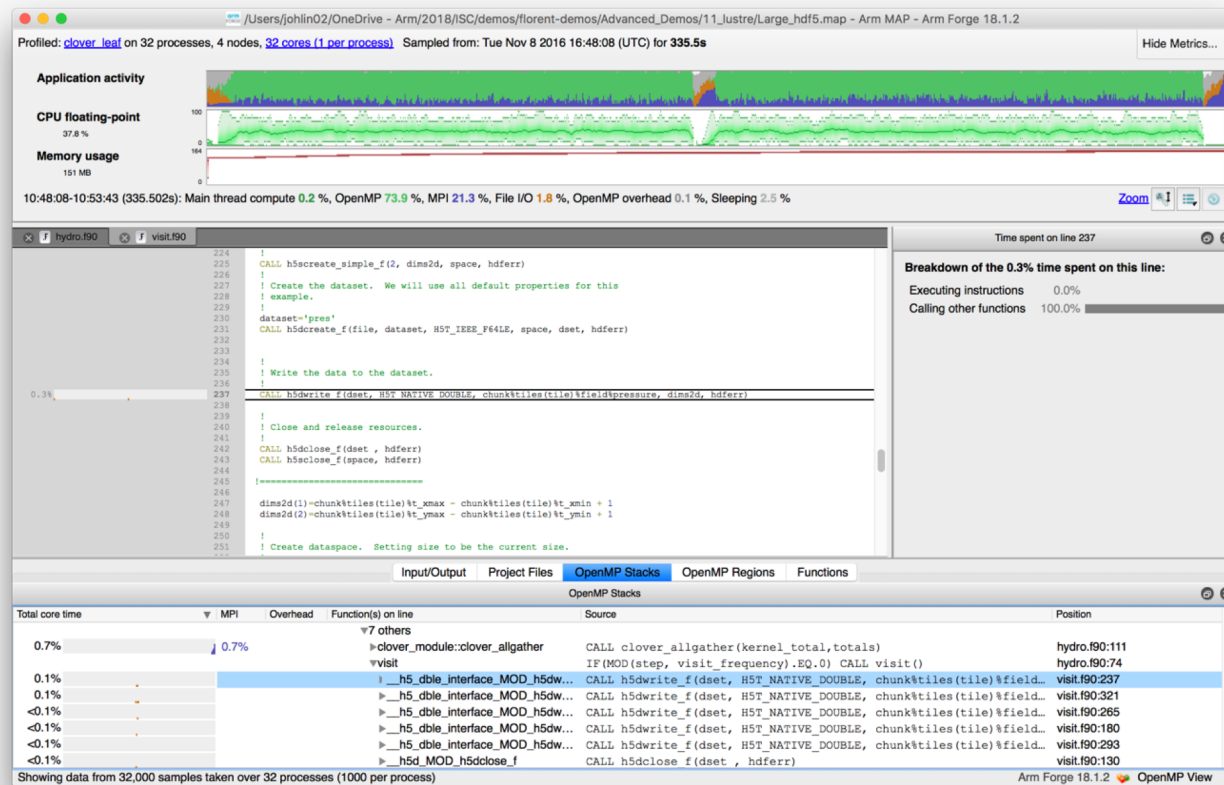# Symptoms and causes of the I/O issues

Sub-optimal file format and surprise buffering.



- Write rate is less than 14MB/s.

- Writing an ASCII output file.

- Writes not being flushed until buffer is full.

  - Some ranks have much less buffered data than others.

  - Ranks with small buffers wait in barrier for other ranks to finish flushing their buffers.
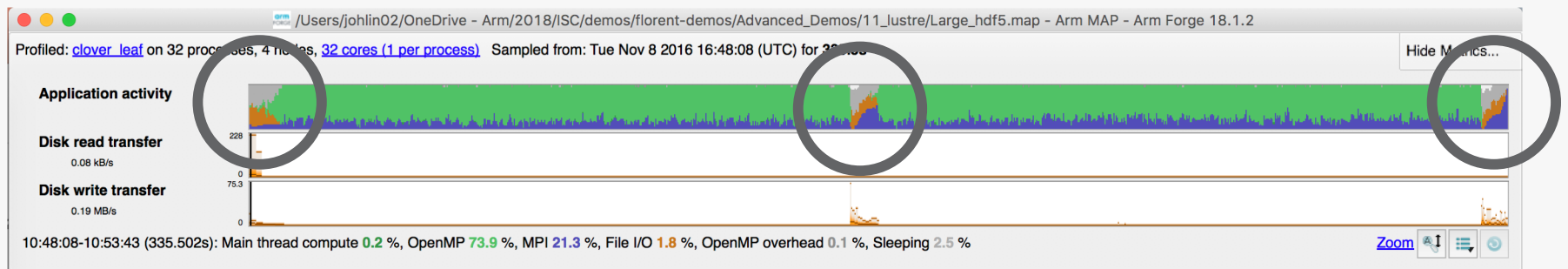
# Solution: use HDF5 to write binary files

Using a library optimized for HPC I/O improves performance and portability.



Argonne Leadership Computing Facility

# Solution: use HDF5 to write binary files

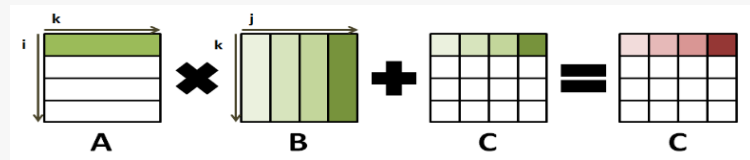Using a library optimized for HPC I/O improves performance and portability.



- Replace Fortran write statements with HDF5 library calls.
  - Binary format reduces write volume and can improve data precision.
  - Maximum transfer rate now 75.3 MB/s, over 5x faster.
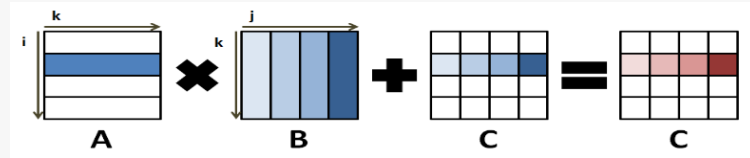- Note MPI costs (blue) in the I/O region, so room for improvement.

# Matrix Multiplication Example
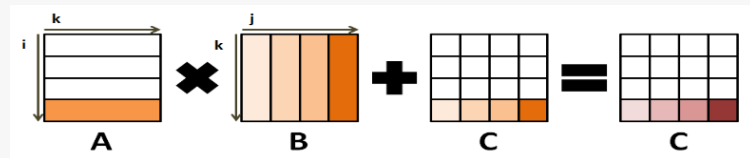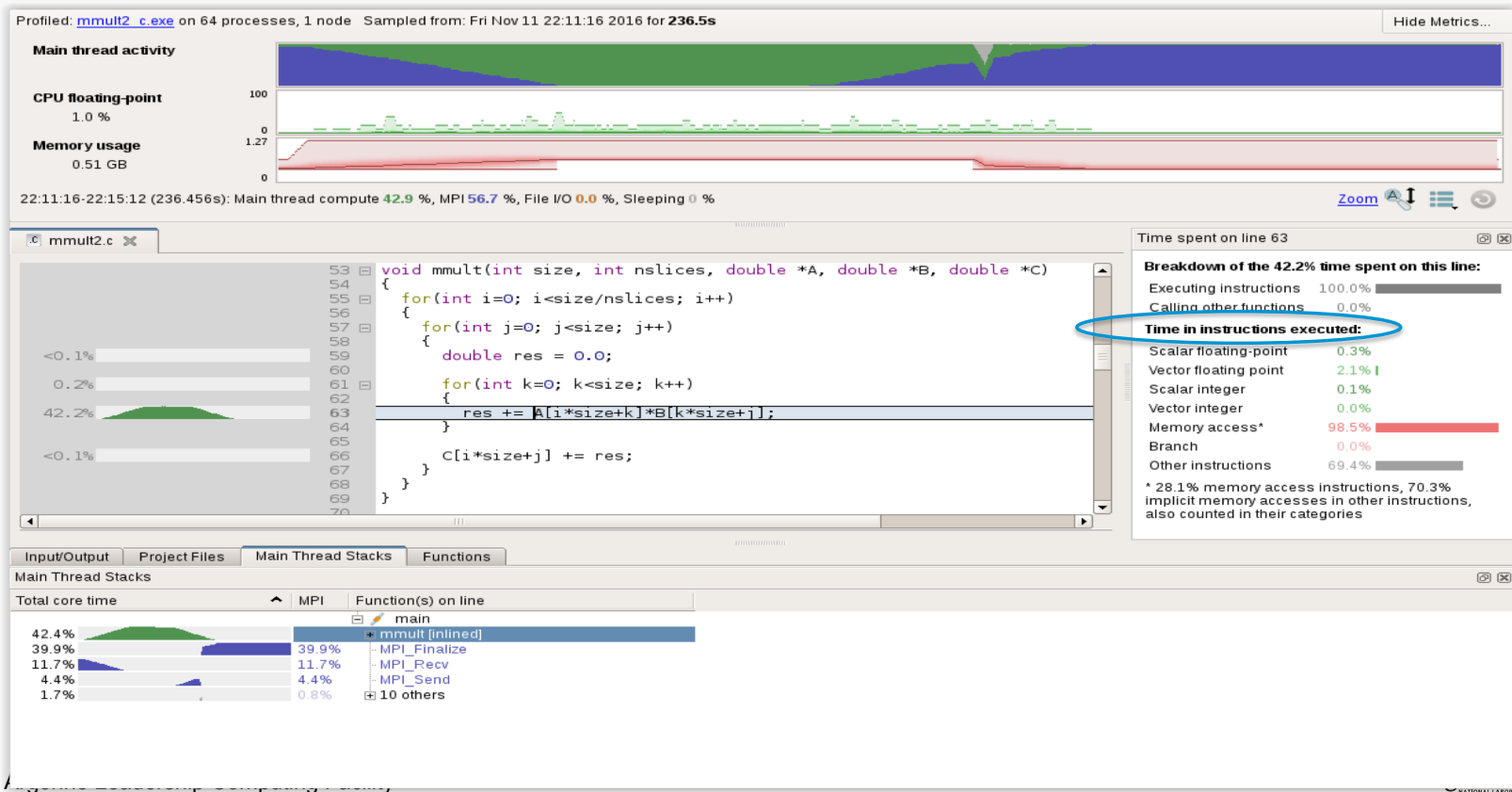
$$C = A \times B + C$$

Master process



Slave process 1



Slave process n-1

# Matrix Multiplication Profile

# Enabling Vectorization

The compiler is unable to vectorize efficiently because of the following line
in C:
**`res += A[i*size+k]*B[k*size+j];`**
and in F90:
**`res=A(i*size+k)*B(k*size+j)+res`**

rewrite mmult to have
in C:
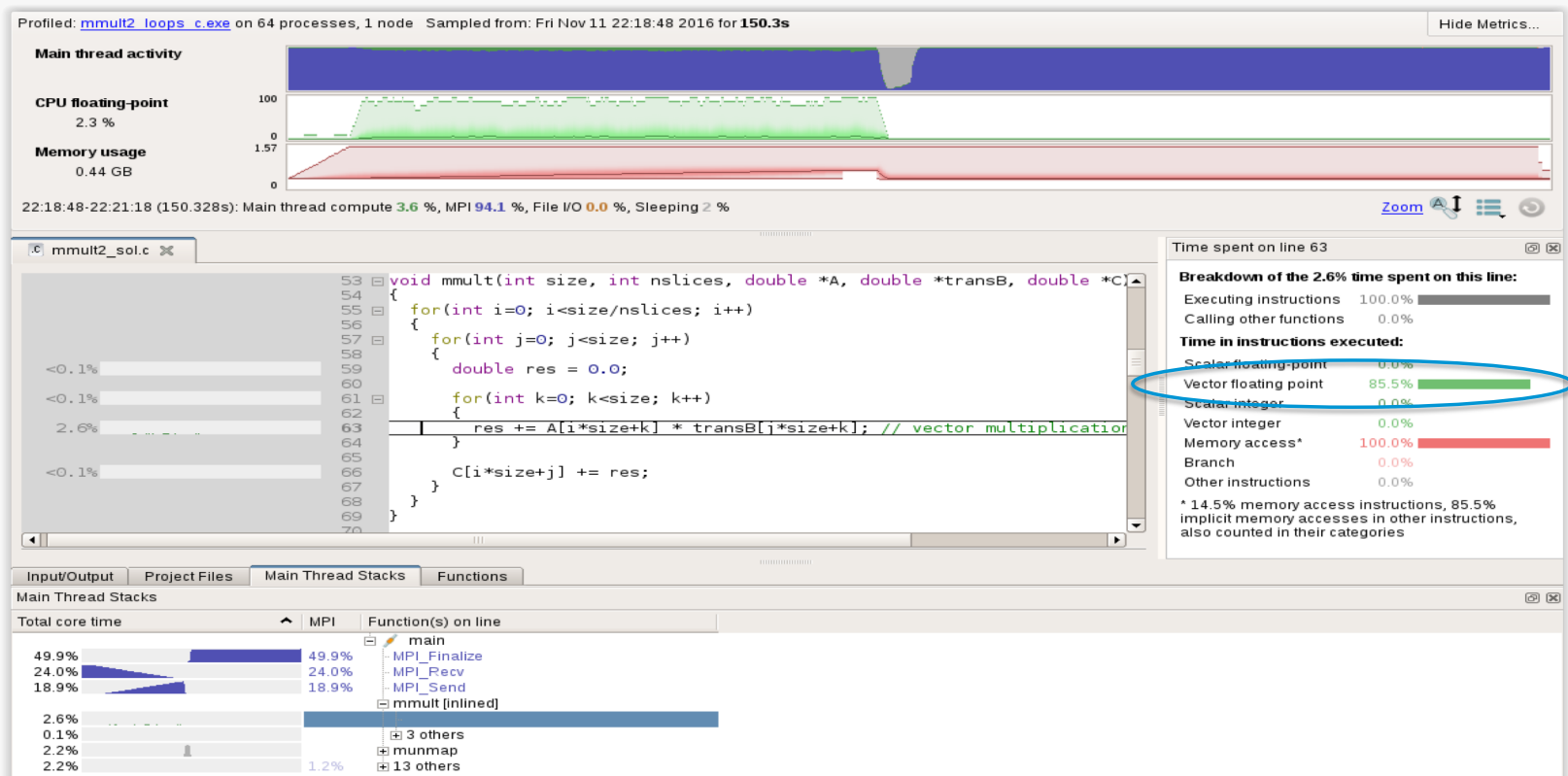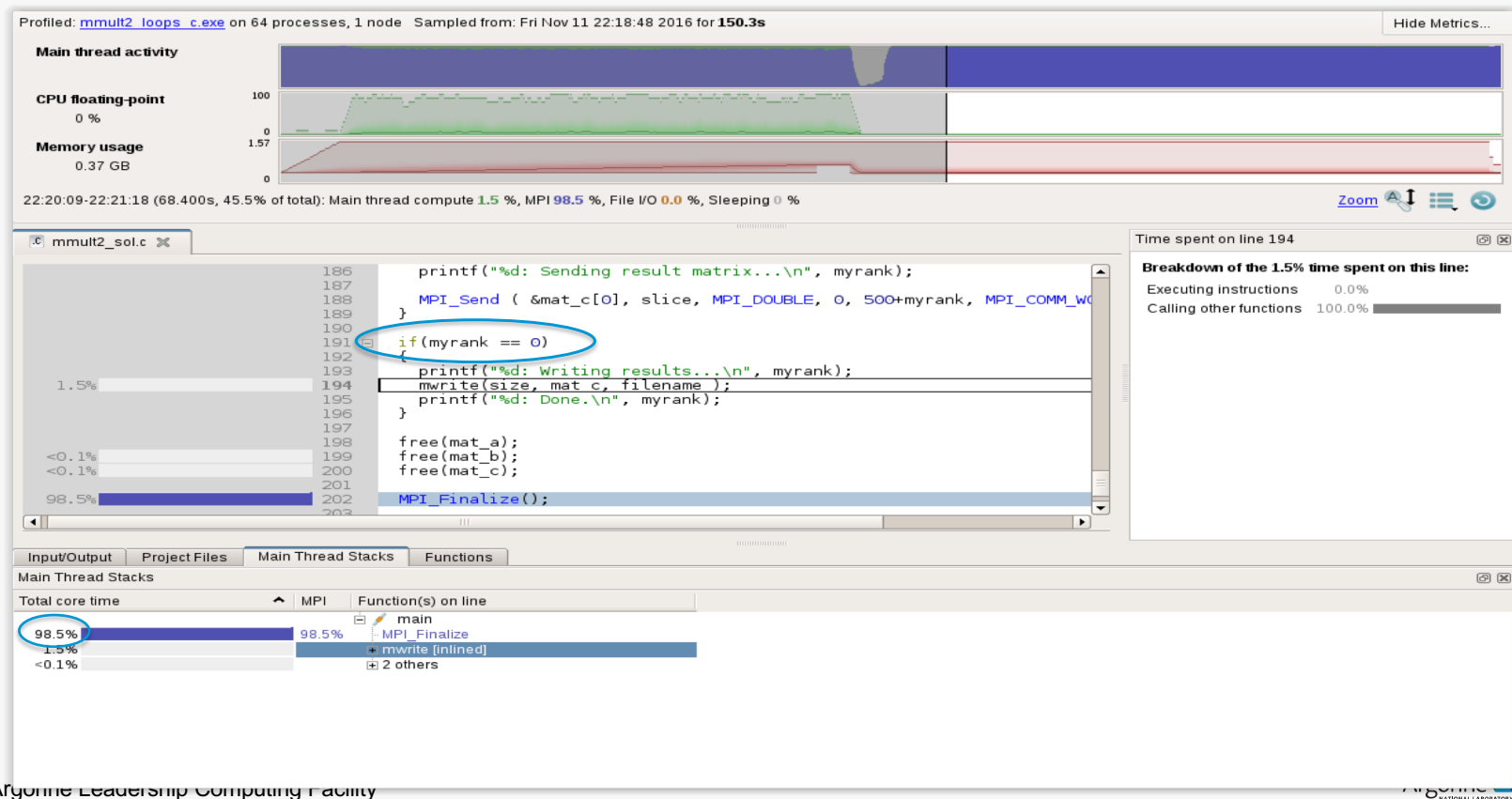**`res += A[i*size+k]*transB[j*size+k];`**
and in F90:
**`res=A(i*size+k)*transB(j*size+k)+res`**

Argonne
NATIONAL LABORATORY

# Improving Data Layout and Access Pattern

# Serial Bottleneck

# Inefficient I/O

```
if(myrank == 0)
{
  printf("%d: Receiving result matrix...\n", myrank);
   […]
}
else
{
  printf("%d: Sending result matrix...\n", myrank);
   […]
}


if(myrank == 0)
{
  printf("%d: Writing results...\n", myrank);
  mwrite(size, mat_c, filename);
}
```

# Improve Scalability of I/O Routines

```
printf("%d: Writing results...\n", myrank);


MPI_File_open(MPI_COMM_WORLD, filename, MPI_MODE_CREATE+MPI_MODE_WRONLY, MPI_INFO_NULL,
&fh);
MPI_File_set_view(fh, 0, MPI_DOUBLE, MPI_DOUBLE, "native", MPI_INFO_NULL);
MPI_File_write_at(fh, slice*myrank*sizeof(double), &mat_c[0], slice, MPI_DOUBLE, &st);
MPI_Barrier(MPI_COMM_WORLD);
MPI_File_close(&fh);
```
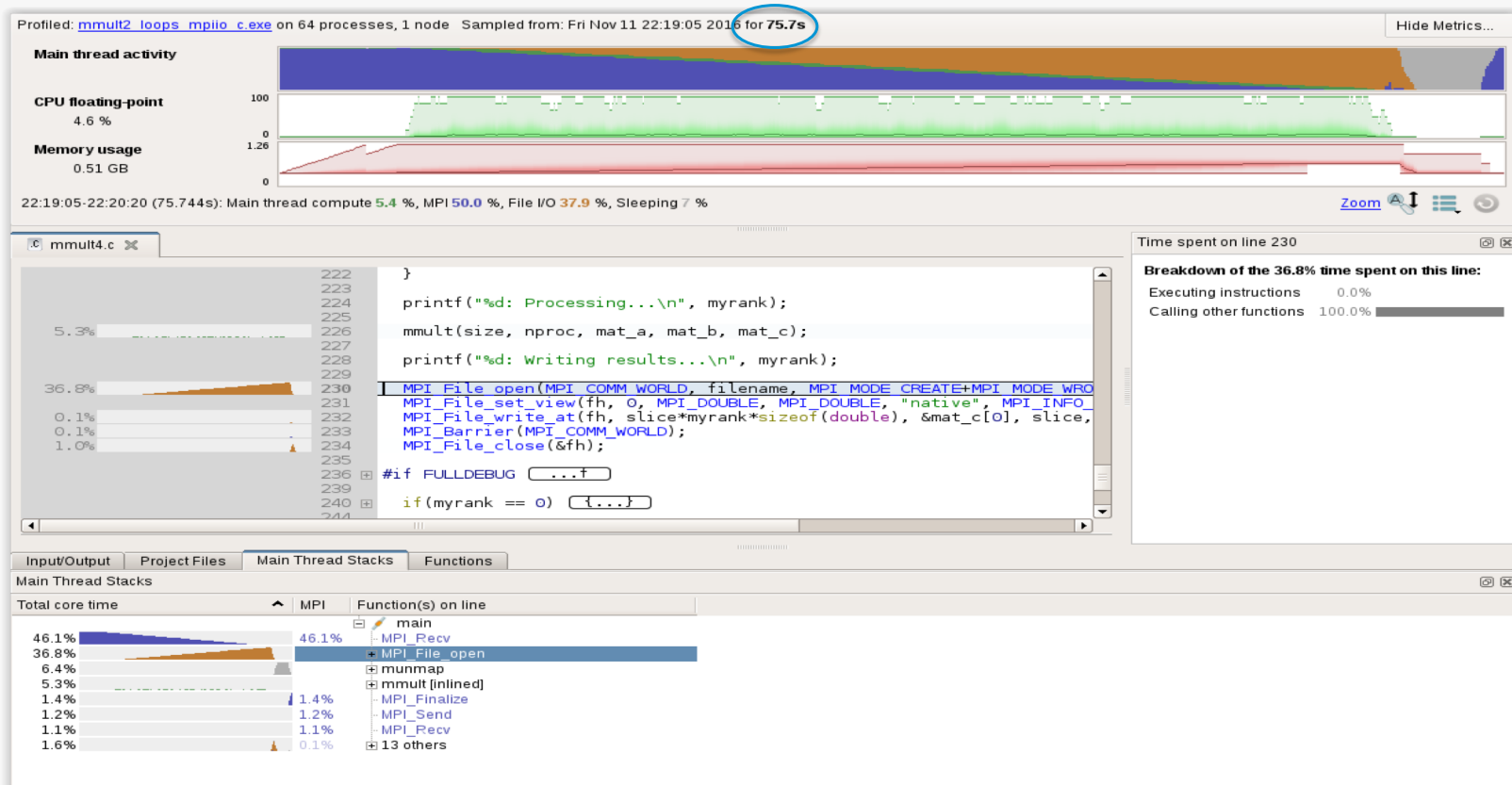
# 3x Speedup from Original Code

# Arm MAP cheat sheet

Load the environment module (manually specify version)
- – $ module load **forge/19.1.2**

Generate the wrapper libraries (static is default on Theta)
- – $ make-profiler-libraries --platform=default --lib-type=static

Unload Darshan module (It wraps MPI calls which cannot be used with MAP)
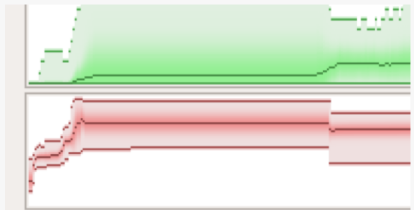- – $ module unload darshan

Follow the instructions displayed to prepare the code
- – $ cc -O3 **-g** myapp.c -o myapp.exe **-Wl,@/path/to/profiler_wrapper_libraries/allinea-profiler.ld**
- – Edit the job script to run Arm MAP in "profile" mode
- – $ **map --profile** -n 8 ./myapp.exe arg1 arg2

Open the results
- – On the login node:
  - • $ map myapp_Xp_Yn_YYYY-MM-DD_HH-MM.map
- – (or load the corresponding file using the remote client connected to the remote system or locally)

Argonne
NATIONAL LABORATORY

# Six Great Things to Try with Allinea MAP


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

Argonne
NATIONAL LABORATORY

# Theta Specific Settings

arm

# Configure the remote client

**Install the Arm Remote Client**

- Go to : https://developer.arm.com/products/software-development-tools/hpc/downloads/download-arm-forge

**Connect to the cluster with the remote client**

– Open your Remote Client

– Create a new connection: Remote Launch ➔ Configure ➔ Add

- Hostname: `<username>@theta.alcf.anl.gov`

- Remote installation directory:

  `/soft/debuggers/forge-19.1.2-2019-08-06`

Argonne
NATIONAL LABORATORY

# Static Linking Extra Steps

To enable advanced memory debugging features in DDT, you must link explicitly against our memory libraries
Simply add the link flags to your Makefile, or however appropriate

lflags = -L/soft/debuggers/ddt/lib/64 -Wl,--undefined=malloc -ldmalloc -Wl,--allow-multiple-definition

In order to profile with MAP, static profiler libraries must be created with the command
make-profiler-libraries --lib-type=static

Instructions to link the libraries will be provided after running the above command

Argonne
NATIONAL LABORATORY

# Questions?

arm

Thank You!

Danke!

Merci!

谢谢!

ありがとう!

Gracias!

Kiitos!

감사합니다

धन्यवाद

arm