# allinea

## Leaders in parallel software development tools

# Allinea DDT
# Intelligence Within
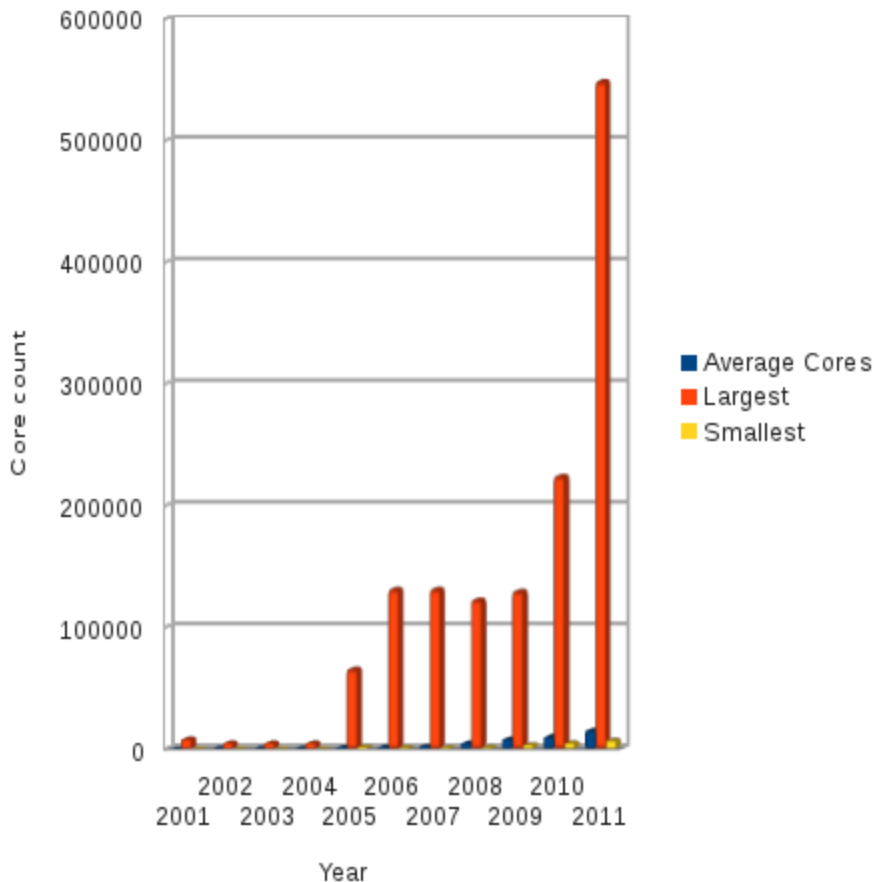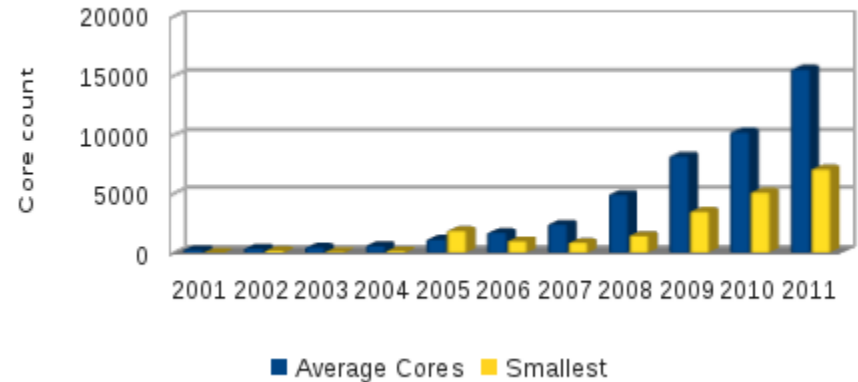
**14 September 2012**

# Agenda

- Introduction

- Allinea DDT overview

- Conclusion

# What is happening ?
# Extreme Machine Size



Growth in HPC core counts



HPC core counts

- Scientific progress requires more CPU hours
  - Maximum machine size is exploding
  - Average machine size grows exponentially

# The Company

- Development tools company
  - Leading in HPC software tools market worldwide
  - Global customer base
    - Blue-chip engineering, government and academic research
- Allinea DDT
  - The leading debugger in parallel computing
  - Production use from desktop to extreme scale
  - World's only scalable debugger
    - Record holder for debugging software on largest machines
    - First at Petascale – and first for GPUs and ARM support!
- Allinea OPT
  - Profiling tool for parallel applications

allinea
www.allinea.com

# Collaboration – National Labs

Partnership to develop Petascale debugger with NVIDIA support, with Cray and Caps Entreprise

Partnership on Full Scale debugging on IBM Blue Gene /P & /Q

Partnership with CEA French Atomic Energy Authority on scalable programming, CUDA and Allinea OPT

European partnership to develop techniques and solutions which address the exascale challenges
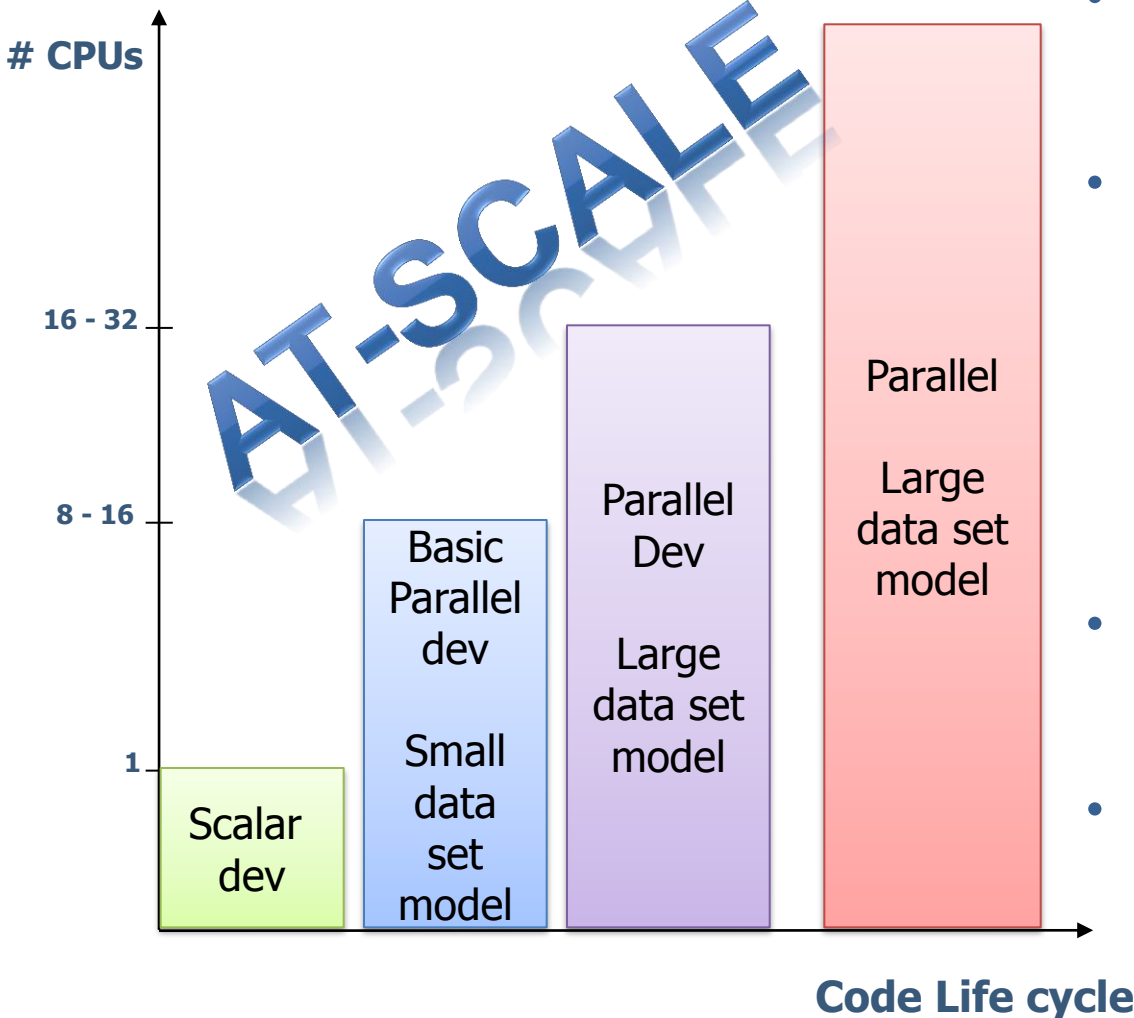
Partnership with BSC through its utilization of ARM technology to develop energy efficient HPC systems

allinea
www.allinea.com

# Application Development Life Cycles



# CPUs

16 - 32

8 - 16

1

Scalar dev

Basic Parallel dev

Small data set model

Parallel Dev

Large data set model

Parallel

Large data set model

AT-SCALE

Code Life cycle

- Bugs are present at any stage

- Bugs get more and more tricky
  - More threads/processes
  - More data
  - Easy bugs already found

- Time to debug increases

- Debugging at scale : this is what you mostly do !

allinea
www.allinea.com

# Example

```
$> mpirun –np 4 ./cpi.exe
      Process 0 on dopey
      Process 2 on dopey
      Process 1 on dopey
      Process 3 on dopey
      pi is approximately 3.1416009869231249
      Error is 0.0000083333333318
      wall clock time = 0.012800
```

**Happy !!**

----------------------------------------------------------------

```
$> mpirun –np 8 ./cpi.exe
      Process 5 on dopey
      Process 1 on dopey
      Process 4 on dopey
      Process 6 on dopey
      Process 3 on dopey
      Process 0 on dopey
      Process 2 on dopey
      Process 7 on dopey

^Cmpirun: killing job...
```

**Deadlock !!**

allinea
www.allinea.com

# How to debug At-Scale ?

- Natural reflex : printf !

- No fully integrated solution : time consuming

  - Problems do not appear at smaller scale

  - Takes time to move the problems to a smaller size

  - Need to use multiple tools for multiple needs

  - More cores means more bugs...

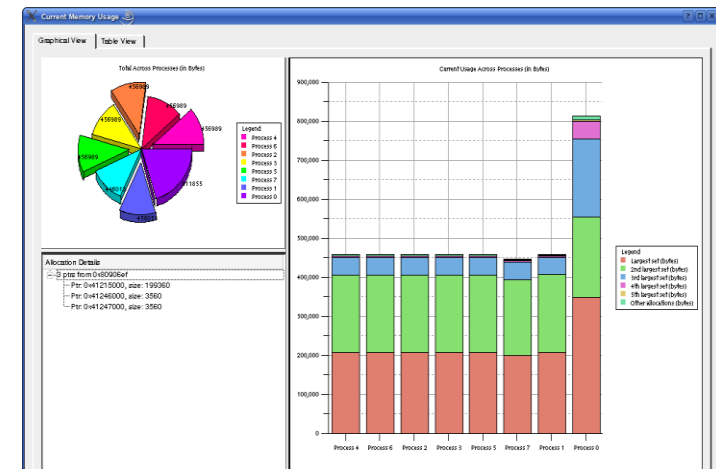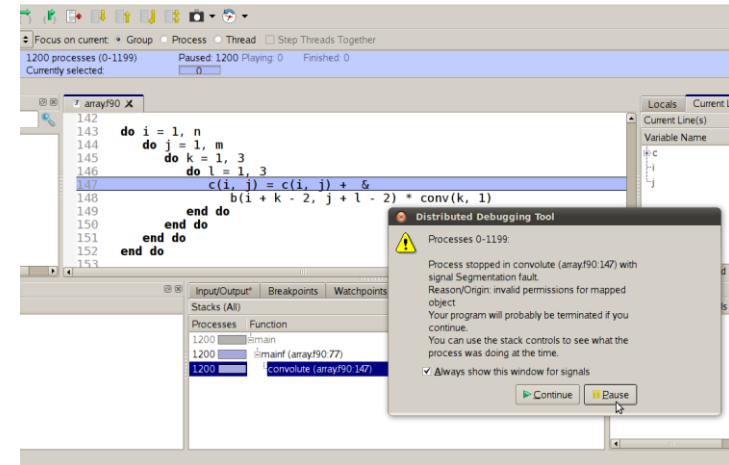  - ... And more debugging information available

**What can we do... ?**

# HPC's current challenge

- GPUs – a rival to traditional processors
  - AMD, Intel, ARM
  - NVIDIA, OpenCL, CUDA
- A big challenge for HPC developers
  - Data transfer
  - Several memory levels
  - Grid/block layout and thread scheduling
  - Synchronization
- New languages, compilers, potential standards



Processing flow on CUDA

allinea
www.allinea.com

# Allinea DDT
# In a Nutshell



- Graphical debugger designed for:
  - C/C++, Fortran, UPC, CUDA
  - Multithreaded code
    - Single address space
  - Multiprocess code
    - Interdependent or independent processes
  - GPU codes
    - Hybrid software
  - Any mix of the above
- Managing concurrency
  - Emphasizing differences
  - Collective control
- Strong feature set
  - Memory debugging
  - Data analysis



allinea
www.allinea.com

# Allinea DDT

## Gather, Sort and Display information

- **User and administrator friendly**

- **Flow control**

- **Data monitoring**

- **Many environments**



allinea
www.allinea.com

# Allinea DDT

## Gather, Sort and Display information

- **User and administrator friendly**
  - Get started easily
  - Fast, reliable, simple and intuitive GUI interface
  - Offline debugging

- Flow control

- Data monitoring

- Many environments



allinea
www.allinea.com

# User and administrator friendly
## Getting started

- Quick creation of runs
- Well integrated in workload schedulers
  - SLURM, LSF, SGE, PBS,....

# User and administrator friendly Intuitive GUI interface

- Intelligent GUI that adapts to the environment
  - From workstation to large scale clusters



www.allinea.com

# User and administrator friendly Intuitive GUI interface

- Intelligent GUI that adapts to the environment
  - Even in CUDA environments

# User and administrator friendly Offline debugging

- Using a workload scheduler
  - Machines are available when the scheduler decides (by night ?)
  - Can be tricky to get a big cluster exactly when the developer wants it

- Offline debugging : printf replacement
  - Tracepoints and offline debugging
  - Job runs without debugger interface and record variables

- Worlds first scalable batch debugger
  - Set tracepoints, breakpoints, and run !
  - Memory debugging errors, crashes
  - Reports in HTML or plain text



**Allinea DDT Off-line Log**

| Messages | Tracepoints | Output |

**Messages**

| Type | Time | Processes | Message |
|---|---|---|---|
| 💡 | 08:53:45.437 | n/a | Launching program /home/david/Work/HEAD/code/ddt/libexec/ddt.bin. |
| ⚠ | 08:53:48.154 | n/a | DDT could not find valid debug information in one or more of your processes. Source files, local variables and other features may be unavailable or inaccurate. Please check you are using the correct debug interface and have compiled your code with debug information. |
| ✎ | 08:53:48.436 | 0 | Memory error detected in FcPatternReference from /usr/lib/x86_64-linux-gnu/libfontconfig.so.1. Thread 1 attempted to dereference a null pointer or execute an SSE instruction with an incorrectly aligned memor Tip: Use the stack list and the local variables to explore your program's current state and identify the source of the Threads,Function 1,main 1, QApplication::QApplication 1, QApplicationPrivate::construct 1, QApplicationPrivate::x11_apply_settings 1, FcInit 1, FcInitLoadConfigAndFonts 1, FcInitLoadConfig 1, FcConfigParseAndLoad 1, XML_ParseBuffer 1, 0x00007ffff0808e6b 1, 0x00007ffff08072e2 1, 0x00007ffff080a77e 1, 0x00007ffff08098e5 1, FcConfigParseAndLoad 1, FcConfigParseAndLoad 1, FcConfigParseAndLoad 1, XML_ParseBuffer 1, 0x00007ffff0808e6b 1, 0x00007ffff08072e2 1, 0x00007ffff080a77e 1, 0x00007ffff08098e5 1, FcConfigParseAndLoad 1, FcValueSave 1, FcValueSave 1, FcPatternReference |
| 💡 | 08:53:48.477 | n/a | Every process in your program has terminated. |

| Messages | Tracepoints | Output |

**Tracepoints**

www.allinea.com

# Allinea DDT and Offline Debugging Using the CLI

- General Options :

```
-offline <YYY>    : activate offline mode
-n <XXX>          : the number of MPI processes
-memory           : enable memory debugging
-ddtsession       : the session file to use (if exists)
```

- Create breakpoints :   `-break-at LOCATION[,N:M:P]`

## Example:

```
$> ddt –offline myjob.html –break-at main.c:22,-:2:-,x myapp
```

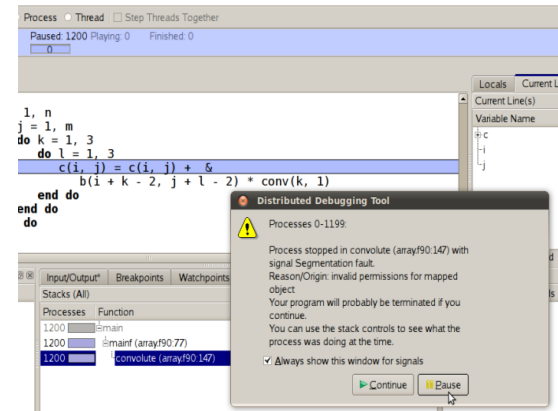- Create tracepoints :   `-trace-at LOCATION[,N:M:P],VAR1,VAR2...`

## Example:

```
$> ddt –offline myjob.html –trace-at trisol.F90:14, -:-:3,x myapp
```

**allinea**
www.allinea.com

# Allinea DDT

## Gather, Sort and Display information

- User and administrator friendly

- Flow control
  - Static analysis
  - Control progress at scale
  - Understand deadlocks
  - Start investigation

- Data monitoring

- Many environments



**allinea**
www.allinea.com

# Flow control
# Static analysis

- Fix those errors before they bite !
- Static analysis
  - Integrated with cppcheck
  - Also includes ftncheck

```
29
30      threads = calloc(sizeof(pthread_t), nthreads);
31      ids = calloc(sizeof(int), nthreads);
32
33      init_mutex();
34
35      pthread_mutex_lock(mutley);
36      for (i = 0; i < nthreads; ++i) {
37          ids[i] = i;
38          pthread_create (threads + i, NULL, &thread,
39      }
40      pthread_mutex_unlock(mutley);
41      for (i = 0; i < nthreads; ++i)
42          pthread_join (threads[i], NULL);
43
44      return 0;
45
```
error Memory leak: threads
error Memory leak: ids     oid *q)
```
49      volatile int busy = 0;
50      volatile int locker = 0;     /* to be amended by
51      int i, j;
52      double k = 1;
53      int tid = *(int*) q;
54
55      usleep(rand() % 31);
56
```

# Flow control
# Controlling progress at scale

- Bulk control is essential for multicore debugging
  - Group processes together
  - Play, step, reach breakpoints... Based on groups
  - Change interleaving order by stepping/playing selectively

# Flow control
# Resolving MPI issues



- We can see messages
  - MPI standard exists for debugging message queues
- Integrates with MPI correctness tools
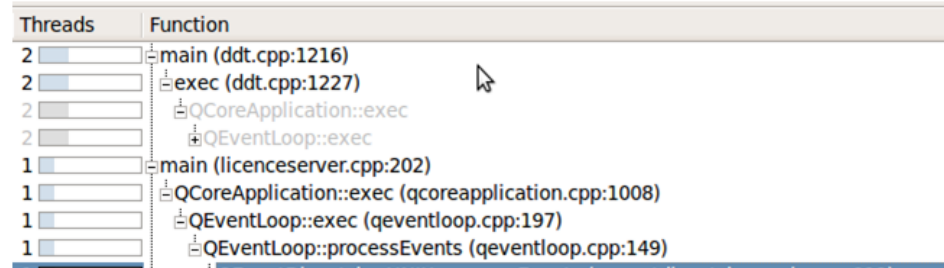  - Check correctness of the messages
  - Intel MPI Checker ; Marmot

- Find cycles or blockage in message queue display
  - Parallel stacks
  - Variables
- More details than examining variables and processes alone

# Flow control
# Understanding what happens

- Application crashes
  - Threads/processes can be anywhere
  - Impossible to scroll through them individually

- Finding where processes crashed is essential
  - Allinea DDT merges stacks from processors and threads into a tree
  - Common faults patterns instantly evident
    - Divergence, deadlocks...
  - Information scalable without overload





allinea
www.allinea.com

# Allinea DDT

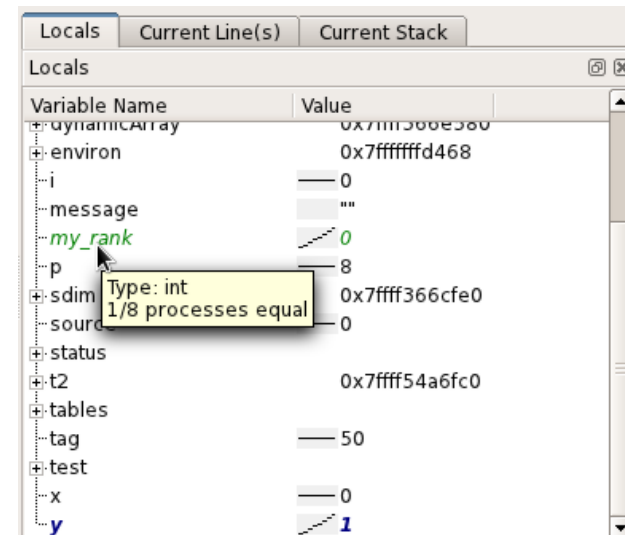## Gather, Sort and Display information

- User and administrator friendly

- Flow control

- Data monitoring
  - Monitor variables
  - Detect memory errors
  - Check the calculation data

- Many environments



allinea
www.allinea.com

# Data monitoring
# Monitor variables

- Developers need to see data
  - Too many variables to trawl them manually
  - Too many tasks or thread to display them at the same time

- Intelligent data management inside the debugger
  - Automatic monitoring of the data
  - Subtle highlights differences
  - Sparklines and smart display

- Even more detailed analysis if needed
  - Cross process comparison
  - Historical values of variables



| Locals | Current Line(s) | Current Stack |
| --- | --- | --- |

| Locals | | |
| --- | --- | --- |
| Variable Name | | Value |
| dynamicArray | | 0x7fff366e380 |
| environ | | 0x7fffffffd468 |
| i | | 0 |
| message | | "" |
| my_rank | | 0 |
| p | | 8 |
| sdim | Type: int | 0x7ffff366cfe0 |
| source | 1/8 processes equal | 0 |
| status | | |
| t2 | | 0x7ffff54a6fc0 |
| tables | | |
| tag | | 50 |
| test | | |
| x | | 0 |
| y | | 1 |

allinea
www.allinea.com

# Data monitoring
# Smart displays : tracepoints

- "printf" is still mostly used but serious drawbacks:
  - Need to recompile the code
  - Information randomly printed on screen (depending on interleaving)
  - One line per process

- Scalable and advanced printf :
  - First step to reconcile printf and GUI
  - Information sorted by steps
  - Merged by groups of processes
  - No information overload
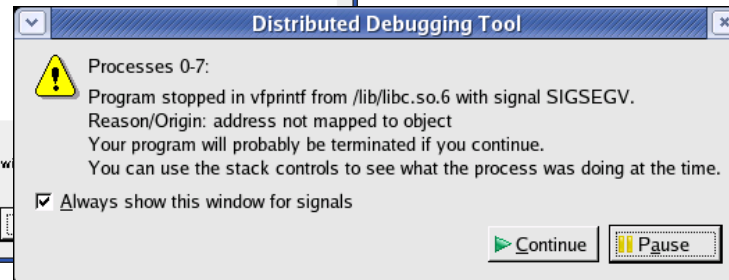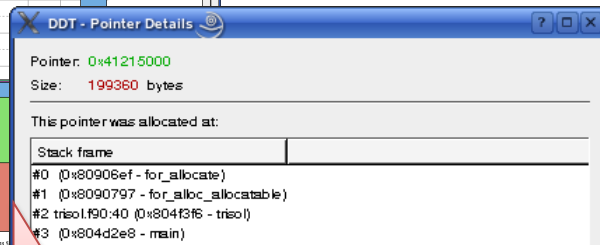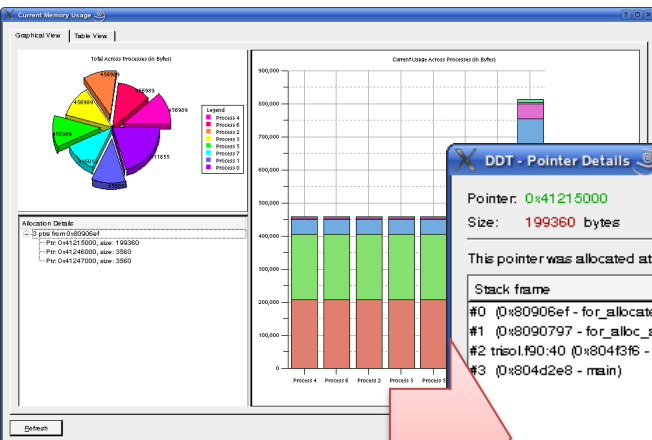  - Possibility to filter the printed values
  - Save output for offline analysis

| Tracepoint Output | | |
|---|---|---|
| Tracepoint | Processes | |
| cstartmpi.c:98 | 4, ranks 0-3 | x: —— 0 |
| cstartmpi.c:98 | 4, ranks 0-3 | x: —— 10 |
| cstartmpi.c:98 | 4, ranks 0-3 | x: —— 20 |
| cstartmpi.c:98 | 4, ranks 0-3 | x: —— 30 |
| cstartmpi.c:98 | 4, ranks 0-3 | x: —— 40 |
| cstartmpi.c:98 | 4, ranks 0-3 | x: —— 50 |

Only show lines containing:

allinea
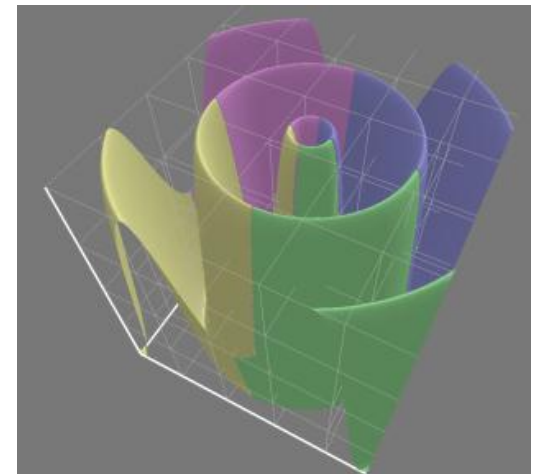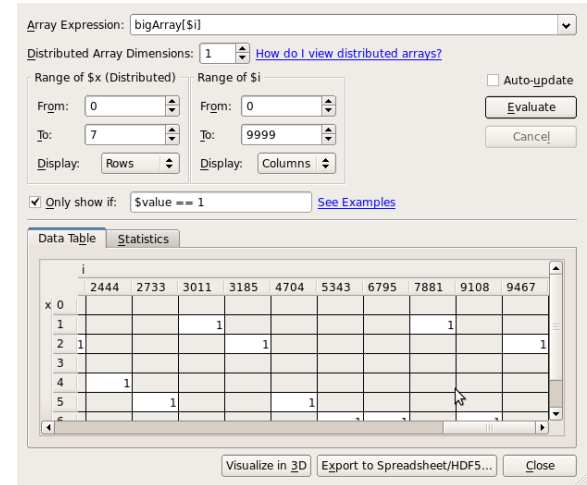
# Data monitoring
# Memory debugging

- Random errors are sneaky
  - Can't fix a bug that doesn't repeat
  - Often caused by memory issues
- Memory debugging can force the bug
  - Better to happen every time, than only during product demos

- Memory debugging :
  - Places agent between memory library and user process
  - Communicates problems to the debuggers
  - Monitors usage : detect memory leaks
  - Automatically protects ends of arrays
  - Trigger instant stop on touching memory
  - Many classes of errors can be checked
  - Also has CUDA support

# Data monitoring
# Searching haystacks

- Arrays are the building blocks of HPC
  - Largest jobs accumulate TB of data
  - Usually 2GB or 4GB per core

- Integrated visualization tool
  - Search data across all tasks or threads
  - Data displayed on a picture as requested by the user
  - Export at runtime





allinea
www.allinea.com

# Data Monitoring
# Unified Parallel C

- ## New programming models
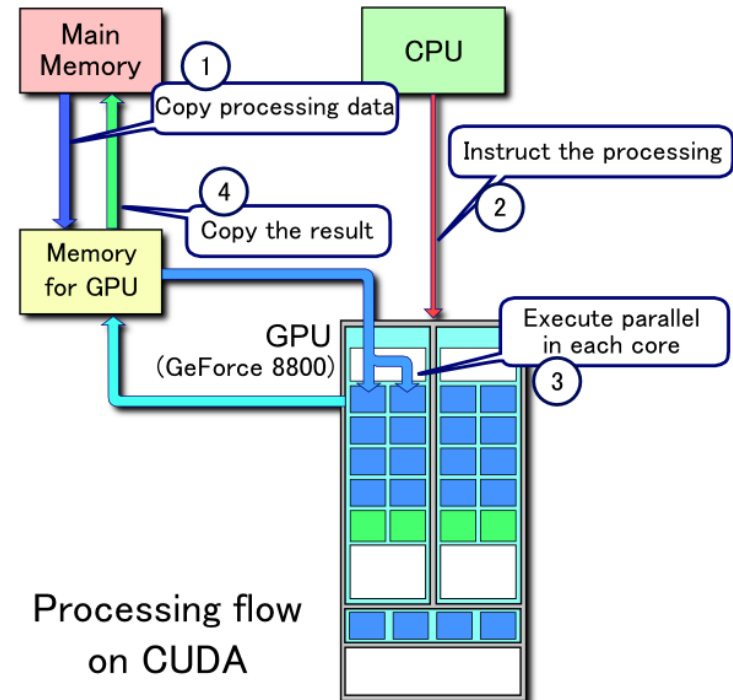  - – Support for Cray UPC
  - – Support for Cray Co-Array Fortran

# Allinea DDT

## Gather, Sort and Display information

- User and administrator friendly

- Flow control

- Data monitoring

- Many environments
  - Debugging at scale
  - GPU Debugging



Processing flow on CUDA

allinea
www.allinea.com

# Allinea DDT and Scalability
# Simplifies debugging for everyone

- At scale : new problems appear
  - Small and medium environments : multiple tools
  - More cores means more bugs...
  - ... And more debugging information available

DDT 3.0 Performance Figures



- How can the design of Allinea DDT help your debugging whatever the size of your cluster (workstation or cluster) ?

- High-end architecture (in any Allinea DDT version)
  - Data consolidation (merged and sorted)
  - Highlight differences
  - Very low footprint

allinea
www.allinea.com

# Allinea DDT and GPU
# Successful add-on

- Built on vendors low level efforts
  - Nvidia cuda-gdb, compiler
  - Cray compiler

- Execution model is unusual
  - GUI supports blocks and grids
  - Support 32 thread units (warps)

- Mixed GPU/CPU in one interface
  - Interaction with CPUs
  - Easy to switch between contexts (stacks, threads, data...)
  - Support multiple nodes



allinea
www.allinea.com

# Many environments
# Directives support



- Wide range of partnerships
  - Support the environments you use for hybrid development
  - You swiftly benefit from the latest updates
  - Read our latest white papers available with PGI and CAPS

# Many environments
# Directives support - OpenACC



- OpenACC : New parallel programming standard led by
  - CAPS Enterprise, Cray, Nvidia and PGI

- Fully supported within Allinea DDT
  - CAPS and PGI : still adapting their software to the new standard
  - Same mechanisms and features for CUDA programs
  - First implementation to be supported: Cray OpenACC

- You benefit from our partnerships


www.allinea.com

# Summary

- Bugs happen at all stage of the development and during the entire life of the code

- An intelligent debugger is now available to fix bugs quickly
  - Other methods have limited success and issues at scale
  - Intelligence makes debugging easier and faster for the developer

- Allinea DDT scales in both performance and interface
  - Breaking all records and making problems now manageable

**allinea**
www.allinea.com

![allinea — Leaders in parallel software development tools]

# Thank you

Your contacts :
- Technical Support team :  support@allinea.com
- Sales team :  sales@allinea.com