

Milestone 4547

## Data Co-Processing for Extreme Scale Analysis

Preliminary Executive Summary of Milestone Report


February 8, 2013

David Rogers, Ron Oldfield, Kenneth Moreland and Nathan Fabian

Sandia National Laboratories

Sandia National Laboratories is a multi-program laboratory operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin company, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC02-04-OR21400.


1



## Summary

- Milestone 4745 "Data Co-Processing for Extreme Scale Analysis" was successfully completed on time, and demonstrated against the letter and spirit of stated Milestone.
- The Milestone Team completed over 9 million node hours of Cielo tests on both *in-situ* and *in-transit* analysis capabilities on a problem provided by a Sandia analyst.
- The results of these experiments have been detailed in a SAND report, which is published as an unclassified unlimited release document, available to the entire mod/sim community

2




## The path to Exascale

Milestone 4745 is an important step in capability development, customer engagement, and scalability development on the path to exascale. It represents significant work on the development of both *Catalyst*, an open source in-situ analysis capability, and *Nessie*, an open source data services capability.

This Milestone is part of an integrated R&D roadmap aimed at characterizing, understanding, and promoting solutions for complex analysis problems on advanced architectures.

It is an important foundation step in developing cross-cutting capabilities.

3



## Milestone 4745

SC calculations produce complex datasets that are increasingly difficult to explore and understand using traditional post-processing workflows. To advance understanding of underlying physics, uncertainties, and results of ASC codes, SNL must gather as much relevant data as possible from large simulations. This drives SNL to couple data analysis and visualization capability with a running simulation, so that high fidelity data can be extracted and written to disk. This Milestone evaluates two approaches for providing such a coupling:

- In-situ processing provides "tightly-coupled" analysis capabilities through libraries linked directly with the simulation. SNL has collaborated on developing an in-situ capability designed for this purpose.
- In-transit processing provides "loosely-coupled" analysis capabilities by performing the analysis on separate processing resources. SNL provides this capability through a "data services" capability designed for this purpose.

SNL will engineer, test and evaluate customer-driven operations on large-scale data created by a running simulation. The data operations will be performed by instrumented versions of both the in-situ and in-transit solutions, with the resulting performance data published and made available to the ASC community.

A program review will be conducted, and its results documented. A report will be submitted as a record of milestone completion.

4

## Motivation



SC calculations produce complex datasets that are increasingly difficult to explore and understand using traditional post-processing workflows. To advance understanding of underlying physics, uncertainties, and results of ASC codes, SNL must gather as much relevant data as possible from large simulations. This drives SNL to couple data analysis and visualization capability with a running simulation, so that high fidelity data can be extracted and written to disk.

- *Note: ASC program will benefit from a detailed understanding of the relationship between analyst tasks, analysis operations, and disk I/O performance.*

5

## In-situ and In-transit workflows



- In-situ processing provides "tightly-coupled" analysis capabilities through libraries linked directly with the simulation. SNL has collaborated on developing an in-situ capability designed for this purpose.



Diagram of in-situ workflow, accomplished in this Milestone through the use of Catalyst, an open source, VTK-based analysis library.

- In-transit processing provides "loosely-coupled" analysis capabilities by performing the analysis on separate processing resources. SNL provides this capability through a "data services" capability designed for this purpose.



Diagram of in-transit workflow, in which the science code communicates with data services nodes to perform analysis operations. This is accomplished in this Milestone through the use of Nessie, an open source data services library.

6

## Milestone 4745, completion criteria



SC calculations produce complex datasets that are increasingly difficult to explore and understand using traditional post-processing workflows. To advance understanding of underlying physics, uncertainties, and results of ASC codes, SNL must gather as much relevant data as possible from large simulations. This drives SNL to couple data analysis and visualization capability with a running simulation, so that high fidelity data can be extracted and written to disk. This Milestone evaluates two approaches for providing such a coupling:

- In-situ processing provides "tightly-coupled" analysis capabilities through libraries linked directly with the simulation. SNL has collaborated on developing an in-situ capability designed for this purpose.
- In-transit processing provides "loosely-coupled" analysis capabilities by performing the analysis on separate processing resources. SNL provides this capability through a "data services" capability designed for this purpose.

SNL will engineer, test and evaluate customer-driven operations on large-scale data created by a running simulation. The data operations will be performed by instrumented versions of both the in-situ and in-transit solutions, with the resulting performance data published and made available to the ASC community.

A program review will be conducted, and its results documented. A report will be submitted as a record of milestone completion.

7

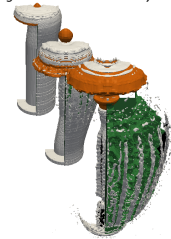
## Experiment Driver



Milestone focused on "customer-driven operations on large-scale data created by a running simulation"

Customer driver use case: characterize fragments in an explosion simulation, an analysis step critical for understanding shock physics

- Partner: Jason Wilke
- Critical steps
  - Find fragments (multiple operations required)
  - Characterize fragments (mass, velocity, etc.)
  - Retain relevant data



Milestone experiments focused on identifying the fragments. This operation is a significantly complex part the analysis, so it serves as a useful way to characterize the operations in the driver use case.

Full range of data experiments run at 32k cores on Cielo (half the machine). Partial experiments done at 64k cores on Cielo (the entire machine). This report presents results from the 32k runs.

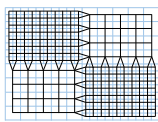
8

## Fragment detection

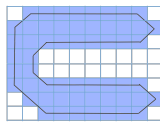


- Operations required for fragment detection (requires a watertight surface)

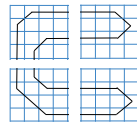
1. Find block neighbors
2. Fix AMR boundaries by inserting degenerate cells
3. Create surfaces
4. Find components across neighbors



Step 2



Step 3



Step 4

9

## Implemented Workflows



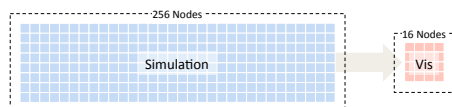
- In situ:** A CTH job that directly runs *in situ* analysis
  - Baseline:** Basic algorithm with somewhat redundant step of global communication to find AMR block neighbors
  - Refined:** Improved algorithm that gets AMR block neighbors from CTH
- In transit:** CTH transfers data to separate server job
  - Extra nodes:** CTH job size same as other runs, extra nodes are used to allocate the VDA service
  - Internal nodes:** CTH job given less nodes that are assigned to VDA service so that together both jobs use the same nodes as other runs
- Spyplot file:** Write Spyplot files from CTH with no analysis

10

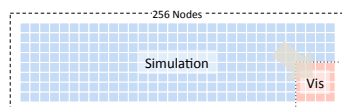
## In Transit Allocations



"Extra Nodes" allocated for VDA services



"Internal Nodes" included in job allocation



## Experiment Configurations



- All experiments performed on Cielo supercomputer at LANL, jointly managed by Los Alamos National Laboratory and Sandia National Laboratories
  - 1,840 node Cray XE6
  - Node: 2 AMD Opteron 6100 (Magny-Cours) 8-way processor chips
    - Total of 16 cores/node
    - 2.6 GHz peak computation speed per core
  - Peak of 1.37 Petaflops
  - 32 GB memory/core

## Experiment, cont'd



- All applications complete 500 cycles (i.e., timestep calculations) of the CTH code.
- The first four applications execute an analysis operation once every 10 cycles
- Spyplot file application outputs spyplot data at a fixed interval in simulated time, calculated so that the application executed the same number of analysis operations performed by the in-situ and in-transit applications

13

## Experiment, cont'd



- For each application, we ran strong scaling experiments for three different datasets.
  - Each data set comes from the same initial conditions but with a different maximum level of refinement
  - Measurements of different job sizes with different data set sizes provides a weak scaling overview.

	CTH			In transit Server		
	Host Core Nodes	In transit Core Nodes	Internal Core Nodes	External Nodes Core Nodes	Internal Nodes Core Nodes	External Nodes Core Nodes
33K Blocks - 5 levels						
128	8	96	6	16	2	2
256	16	224	14	16	2	2
512	32	496	30	16	2	2
1,024	64	992	62	16	2	2
226K Blocks - 6 levels						
1,024	64	768	48	128	16	128
2,048	128	1,536	112	128	16	128
4,096	256	3,072	240	128	16	128
8,192	512	6,144	496	128	16	128
1.5M Blocks - 7 levels						
8,192	512	4,096	256	1,024	128	896
16,384	1,024	8,192	512	1,024	128	896
32,768	2,048	16,384	1,024	1,024	128	896
65,536	4,096	32,768	2,048	1,024	128	896

Table shows the range of core sizes used for the various experiments. For every application we used the maximum 16 cores per node for the CTH client, since CTH is primarily bound by computation and scales very well.

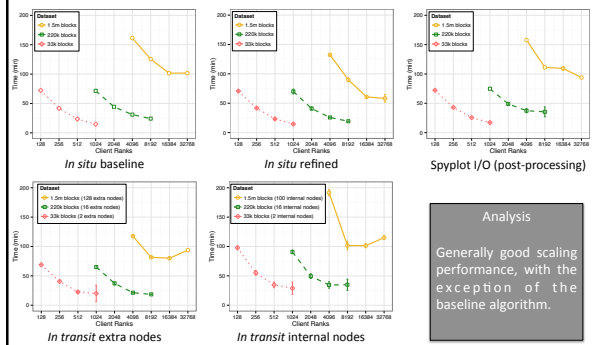
14

## Results



15

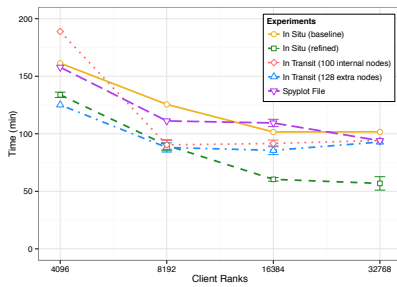
## Basic Timing



Analysis  
Generally good scaling performance, with the exception of the baseline algorithm.

16

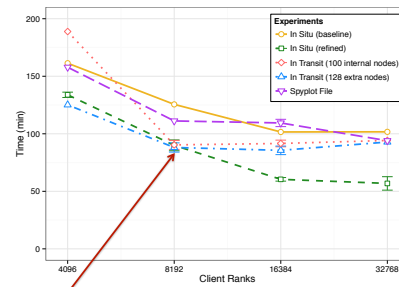
## Pipeline Summary Timing



Generally good scaling performance, with the exception of the baseline algorithm.

17

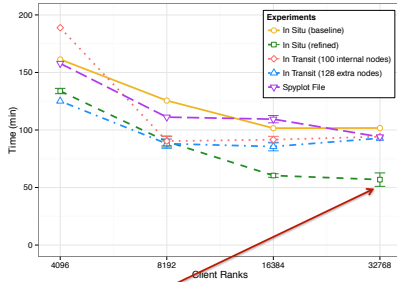
## Pipeline Summary Timing



Sweet spot at 8K cores: *in transit* with unrefined algorithm beats *in situ* with refined algorithm.

18

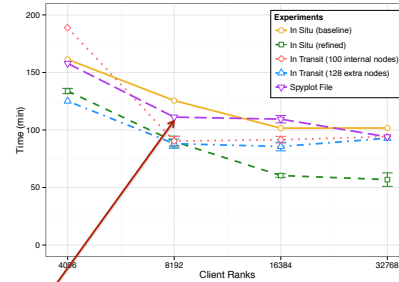
## Pipeline Summary Timing



No significant improvement at 32K cores. Probably insufficient work for analysis (only 45 blocks per process).

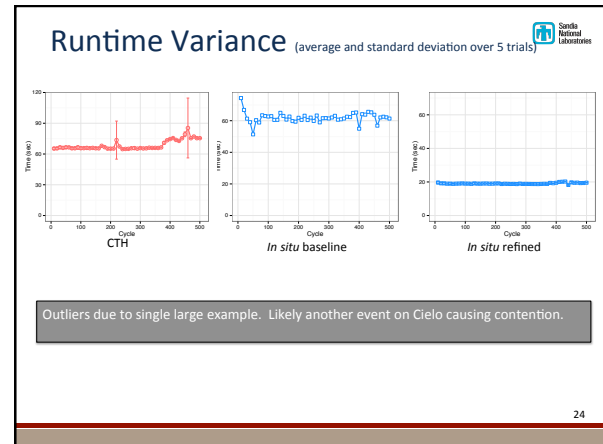
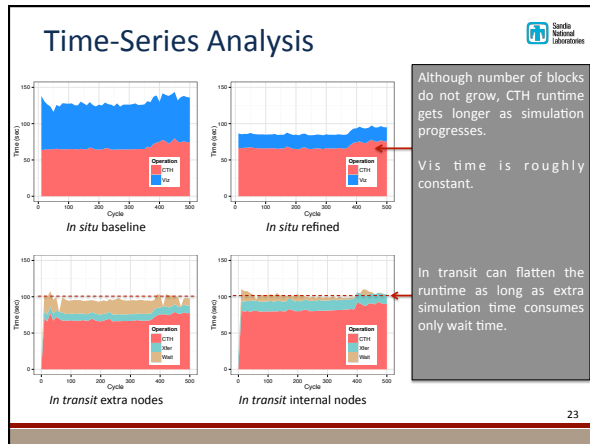
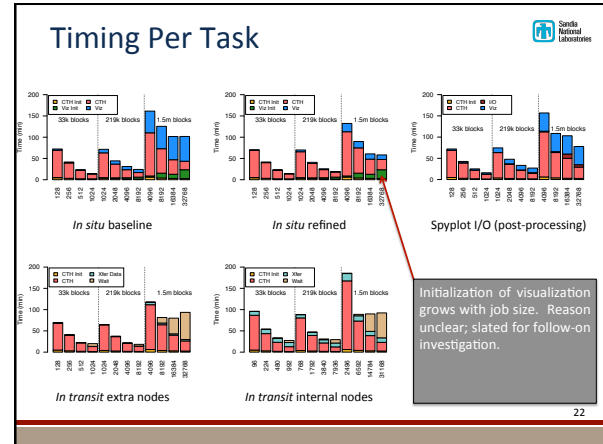
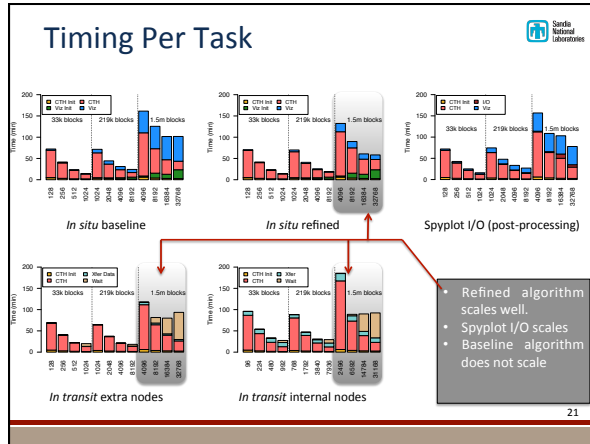
19

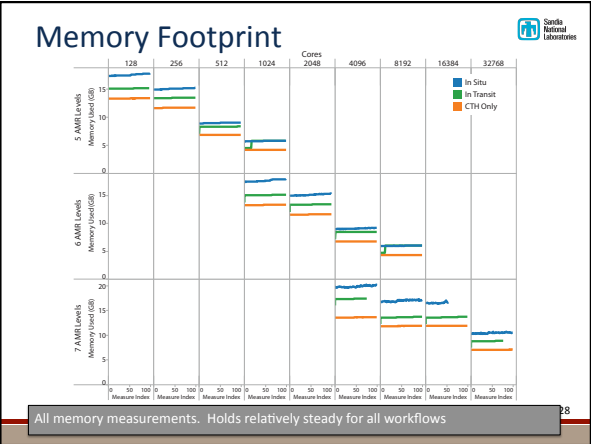
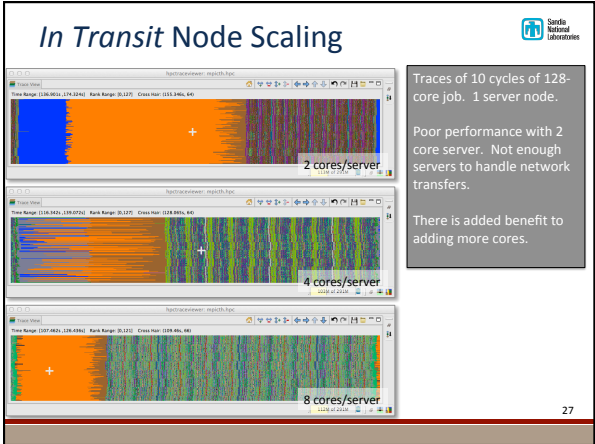
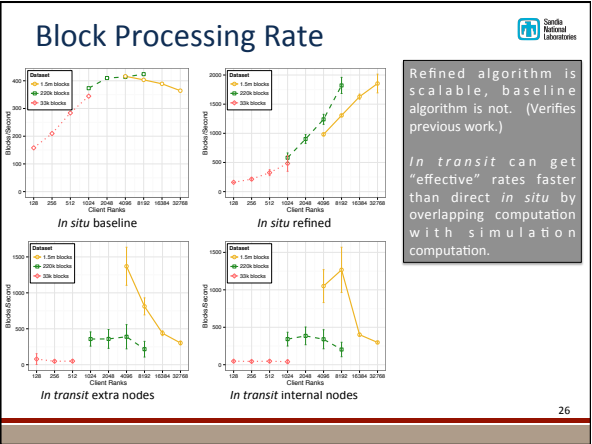
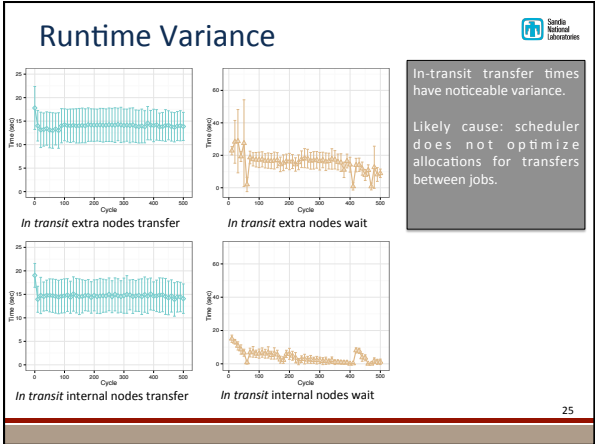
## Pipeline Summary Timing

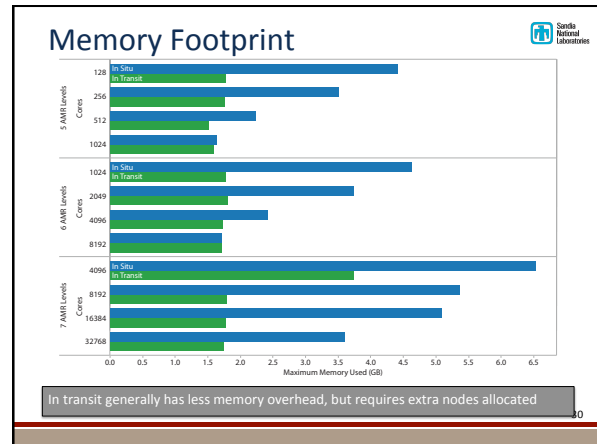
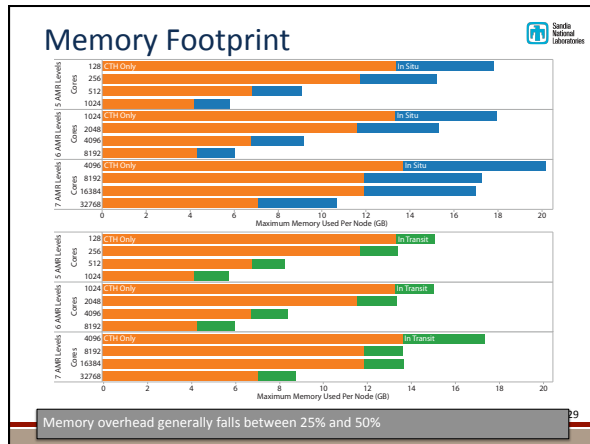


Writing files surprisingly fast. Although slower than most alternatives, still a viable option.

20







## Conclusions

31

## Conclusions

**In transit can provide a performance improvement over in situ in some circumstances, but the window is narrower than we initially expected it would be.**

In transit analysis has an added overhead above embedded in situ analysis involving transferring data between parallel jobs. Given an analysis with perfect linear scalability, we suspect in transit workflows will always have an added cost, and our results support this. With an analysis that does not scale perfectly, possibly due to communication overhead, it is theoretically possible for in transit to be faster by reducing the size of the analysis job. This is one of the motivations for choosing an analysis task that requires significant communication. In our results, we do find instances where in transit is faster, but by a smaller margin and for fewer configurations than we initially anticipated. So although in transit has several other positive features, we do not anticipate performance to be the main motivations for using it.

32



## Conclusions



### Memory overhead will be an important trade-off space.

The baseline amount of memory added to the CTH job to perform in situ processing is roughly 100MB per core. Considering that our embedded in situ library is a fully featured visualization toolkit containing over 2 million lines of code and algorithms developed over almost 2 decades, this overhead is not unreasonable. Nevertheless, this footprint can be problematic for simulations already tight on memory. Because of this, efforts are already underway to improve our memory footprint by making finer modules and being more selective on the available algorithms. This, of course, requires a compromise between the size of the library and the algorithms that are dynamically available. We also note that our algorithm has the potential to generate sizable meshes of its own. Thus, it may be fruitful to pursue and support incremental algorithms where possible.

33

## Conclusions



### Initialization time matters

Our scaling efforts to date focus on the scalability of the algorithms invoked during the run of a simulation. The initialization cost, a one-time penalty, has yet to be seriously considered. However, based on our HPCToolkit measurements, initialization becomes a significant cost at high process counts.

### Disk-based I/O is not dead . . . yet.

Our initial assumption was that it would not be feasible to output full results at a fine enough temporal resolution from CTH to disk storage to perform our high fidelity analysis. However, our control workflow shows that although the overall time to write data to disk and then read back again incurs a large cost, it is still realistic to do so. Thus, users may still choose to incur the extra overhead to use a traditional offline post-processing visualization and data analysis workflow.

34

## Conclusions



### Better job scheduling is important

One of the more complicated parts of running an in transit workflow is scheduling the simulation job and service job to run in tandem. Frankly, the capabilities of the scheduler are inadequate for our needs. We cannot start and stop jobs independently and make reconnections dynamically. Another experiment we would like to do but is challenging to schedule is to allow simulation and service to share nodes. Since each node has 16 cores, perhaps we could get better transfer performance by allocating one core per node for service and the rest for simulation. A similar scheduling scheme will be important to take advantage of burst buffers in future architectures.

35

## Summary



- Milestone 4745 “Data Co-Processing for Extreme Scale Analysis” was successfully completed on time, and demonstrated against the letter and spirit of stated Milestone.
- The Milestone Team completed over 9 million node hours of Cielo tests on both *in-situ* and *in-transit* analysis capabilities on a problem provided by a Sandia analyst.
- The results of these experiments have been detailed in a SAND report, which is published as an unclassified unlimited release document, available to the entire mod/sim community

36