

Everything you always wanted to know
about **Streaming** scientific **EXascale** data
but were afraid to ask; it's a **Yotta** bytes
(SEXY)

New York Scientific Data Summit

+ the help of almost 100 researchers!
August 15, 2016
NY, NY

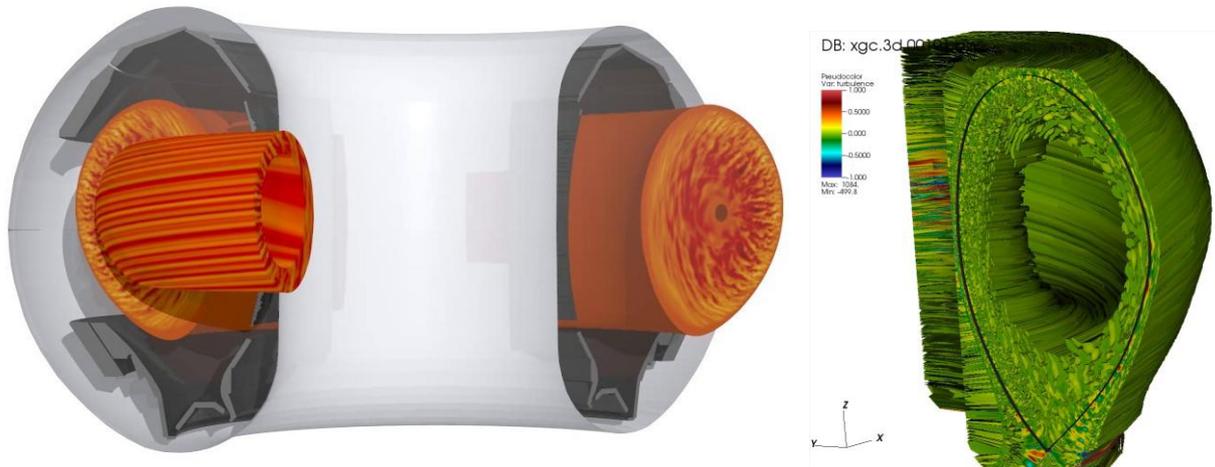


The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- eXascale Service Orient Architecture
- Movement to the Exascale
- VTK-M
- Conclusions



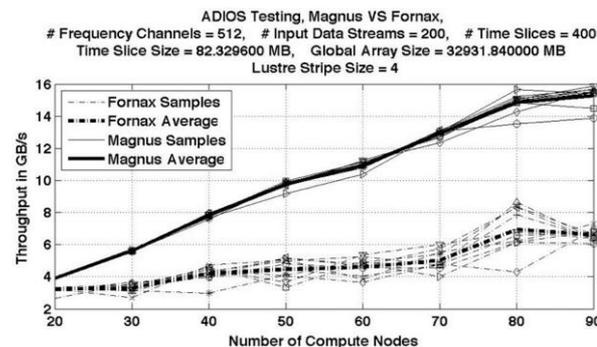
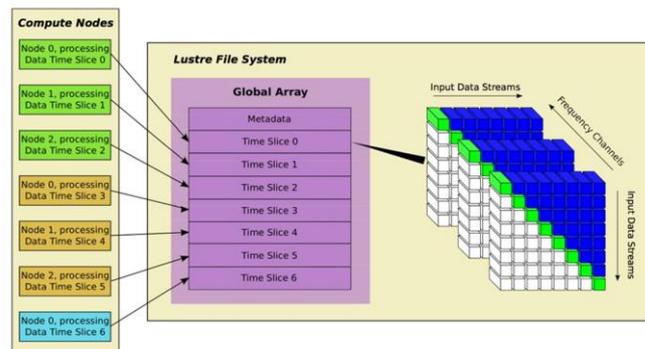
Computation: Fusion



- Develop high-fidelity Whole Device Model (WDM) of magnetically confined fusion plasmas to predict the performance of ITER
- Couple existing, well established extreme-scale gyrokinetic codes
 - GENE continuum code for the core plasma and the
 - XGC particle-in-cell (PIC) code for the edge plasma
- Data challenges
 - Couple codes (XGC to GENE) using a service framework
 - Need to save $O(10)$ PB of data from a single simulation

Observational: SKA

- The largest radio telescope in the world
- €1.5 billion project
- 11 member countries
- **2023-2030** Phase 2 constructed
- Currently conceptual design & preliminary benchmarks !
- Compute Challenge: • 100 PFLOPS
- Data Challenge: ExaBytes per day
- Challenge is to run time-division correlator and then write output data to a parallel filesystem



The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- eXascale Service Orient Architect
- Movement to the Exascale
- VTK-M
- Conclusions



1976 Cray 1. “The birth of supercomputing”

- \$8M
- The birth of vector computing
- Used Integrated Circuits
- Liquid cooled (Freon)
- 64-bit system
- 250 Mflops
- 5.5 Tons
- 16 MB of memory
- 100 MB/sec I/O system, which was a separate machine
- Later Cray YMP, C90 used SSD for absorbing the bursty I/O



2000 ASCI White

- Peak performance is 12.3 Tflops.
- Used 8,192 IBM RS6000 SP3 procs
 - Each processor ran at 375 MHz.
- 6TB of RAM
- Runs IBM's AIX operating system.
- 106 Tons.
- 160 TB of Disks.
- Located at LLNL.
- Fastest in the world until 2002.

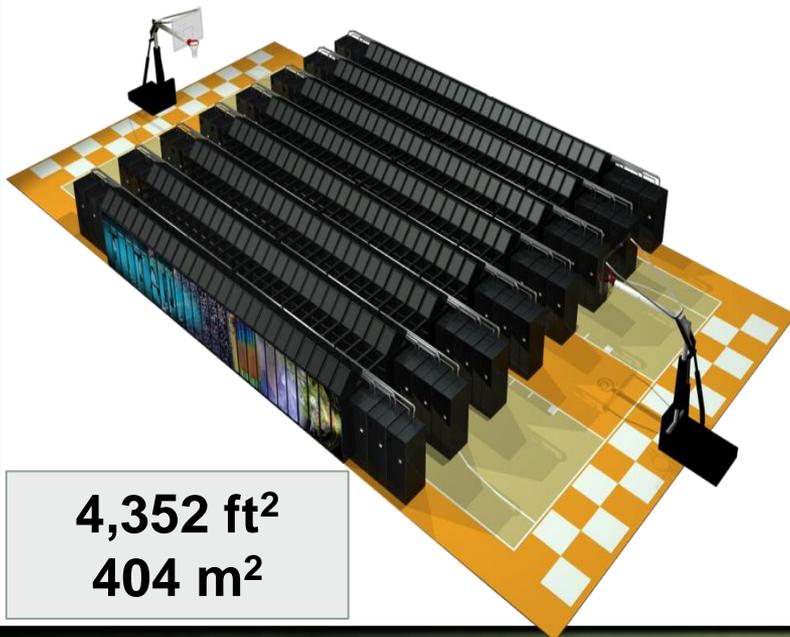


ORNL's "Titan" Hybrid System: CPU + GPGPU



SYSTEM SPECIFICATIONS:

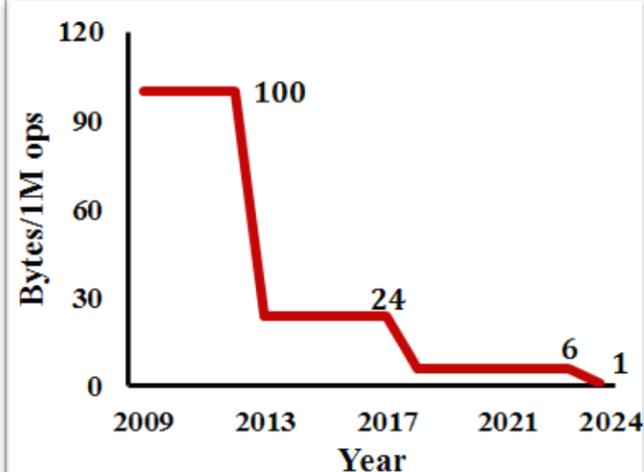
- Peak performance of 27.1 PF (24.5 & 2.6)
- 18,688 Compute Nodes each with:
- 16-Core AMD Opteron CPU (32 GB)
- NVIDIA Tesla "K20x" GPU (6 GB)
- 512 Service and I/O nodes
- 200 Cabinets
- 710 TB total system memory
- Cray Gemini 3D Torus Interconnect



4,352 ft²
404 m²

Compute-Data Gap

- File system/network bandwidth does not keep up with computing power
- Too much data to move, too little time



System attributes	NERSC Now	OLCF Now	ALCF Now	NERSC Upgrade	OLCF Upgrade	ALCF Upgrades	
Name Planned Installation	Edison	TITAN	MIRA	Cori 2016	Summit 2017-2018	Theta 2016	Aurora 2018-2019
System peak (PF)	2.6	27	10	> 30	150	>8.5	180
Peak Power (MW)	2	9	4.8	< 3.7	10	1.7	13
Total system memory	357 TB	710TB	768TB	~1 PB DDR4 + High Bandwidth Memory (HBM)+1.5PB persistent memory	> 1.74 PB DDR4 + HBM + 2.8 PB persistent memory	>480 TB DDR4 + High Bandwidth Memory (HBM)	> 7 PB High Bandwidth On- Package Memory Local Memory and Persistent Memory
Node performance (TF)	0.460	1.452	0.204	> 3	> 40	> 3	> 17 times Mira
Node processors	Intel Ivy Bridge	AMD Opteron Nvidia Kepler	64-bit PowerPC A2	Intel Knights Landing many core CPUs Intel Haswell CPU in data partition	Multiple IBM Power9 CPUs & multiple Nvidia Volta GPUs	Intel Knights Landing Xeon Phi many core CPUs	Knights Hill Xeon Phi many core CPUs
System size (nodes)	5,600 nodes	18,688 nodes	49,152	9,300 nodes 1,900 nodes in data partition	~3,500 nodes	>2,500 nodes	>50,000 nodes
System Interconnect	Aries	Gemini	5D Torus	Aries	Dual Rail EDR-IB	Aries	2 nd Generation Intel Omni-Path Architecture
File System	7.6 PB 168 GB/s, Lustre®	32 PB 1 TB/s, Lustre®	26 PB 300 GB/s GPFS™	28 PB 744 GB/s Lustre®	120 PB 1 TB/s GPFS™	10PB, 210 GB/s Lustre initial	150 PB 1 TB/s Lustre®

The proposed? Exascale System



ExFlops



DOE

Storage



The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- eXascale Service Orient Architecture
- Movement to the Exascale
- VTK-M
- Conclusions



Data driven Science: Why is this important

- Data is increasing faster than our ability to understanding it
- Big data is characterized by
 - Volume
 - Remote Sensing
 - Web (text, images, video)
 - Simulations: 1 PB for some simulations
 - Experiments (LHC)
 - Velocity
 - DOE experiments
 - Variety
 - Heterogeneous, spatial and temporal
 - Multi-resolution, Multi-sensor
 - Veracity
 - Value
- What is “Big Scientific data”
 - $O(N^p)$ algorithms, data movement will not work as $N \gg$ “Big”
 - Extracting information from the data, sharing the data



Start with the 3V's

- Volume

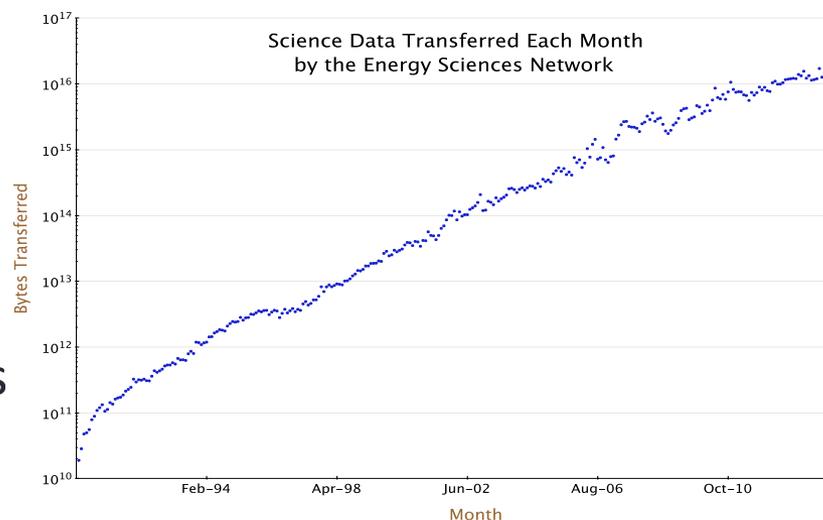
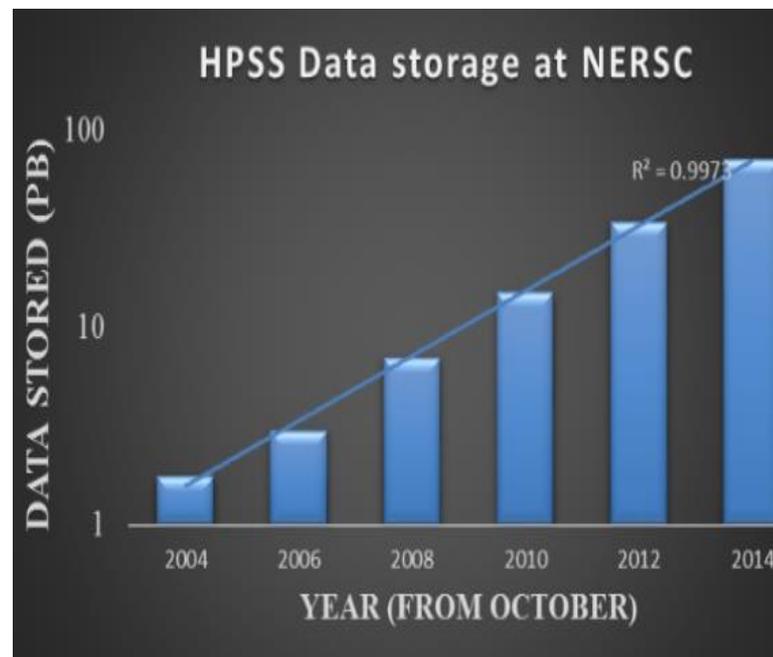
- Simulations: S3D, XGC1 > 1 PB
- Experiments: ITER: 2 PB
- Observations: Medical Slides: 100 GB/slide

- Velocity

- Simulations: QMCPack = 2.5 TB/60 s
- Experiments: ITER: 2 GB/s → 50 GB/s
- Observations SKA: 1 – 10 PB/s

- Variety

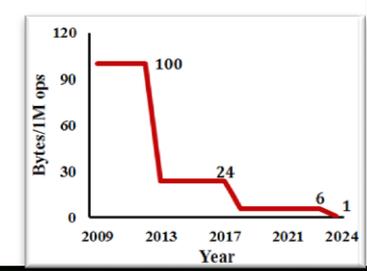
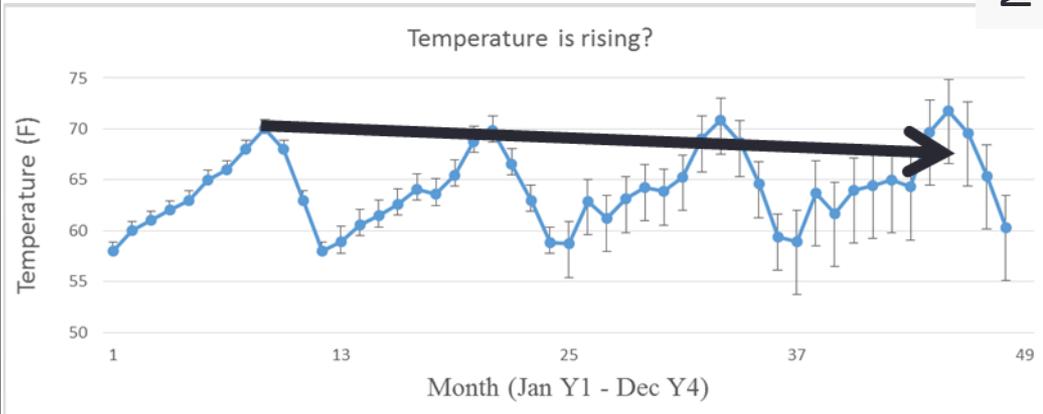
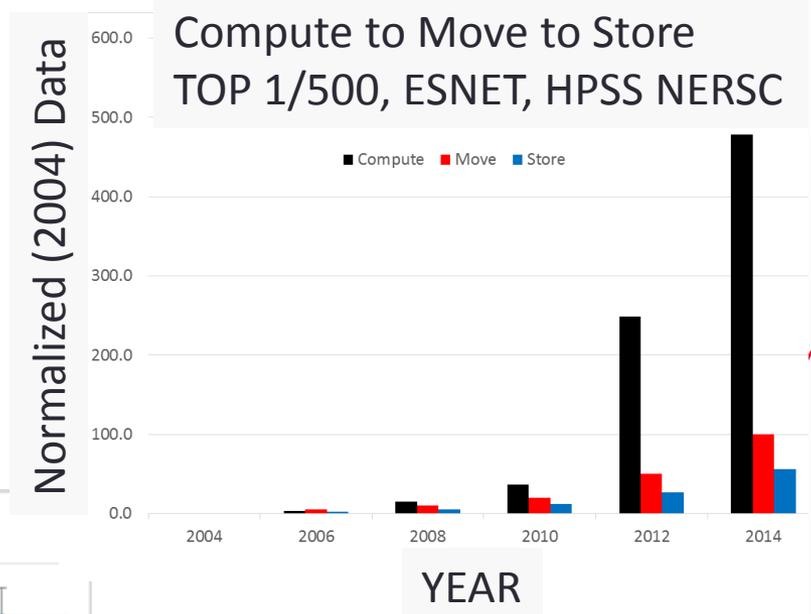
- RAMGEN: 40K variables
- QMCPack: >100 moments/flight
- XGC1: >8 moments/particle + O(100) variables for observables



End with the "other" 2 Vs

- Veracity
- How accurate is the data?
 - Diagnostics have errors
 - Simulations have errors
- Accuracy of data most be taken into account when processing data
 - UA , Quality Control, QA

- Value
 - If data has little value, or is it's too painful to access & process, then Fallback to old style



What's unique in scientific data

- Locality
 - Space-Time
- Computational Data
 - It's so easy to achieve huge orders of magnitude in reduction without losing any accuracy.
 - How do we understand all of the tradeoffs of performance vs. accuracy vs. storage vs. data movement vs. time to knowledge?
- Errors
- Data comes from our measurement from “mother nature” or our re-creation of mother nature



Community togetherness is a must

- Creating a software ecosystem to aid this process
 - Storage System
 - I/O
 - File/Stream data movement
 - Reduction
 - Queries
 - Workflows
 - Visualization
 - Analysis
 - Programming Models
 - Runtimes
 - Provenance
- Create services to plug-into an overarching framework?

The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- **Our solution**
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- eXascale Service Orient Architecture
- Movement to the Exascale
- VTK-M
- Conclusions



ADIOS 1.X

- ADIOS = An I/O abstraction framework
- Provides portable, fast, scalable, easy-to-use, metadata rich output
- Choose the I/O method at runtime
- Abstracts the API from the method
- <http://www.nccs.gov/user-support/center-projects/adios/>
- Need to provide solutions for “90% of the applications”
- We have over 2,200 citations for this work



ADIOS applications



1. Accelerator: **PIConGPU, Warp**
2. Astronomy: **SKA**
3. Astrophysics: **Chimera**
4. Combustion: **S3D**
5. CFD: **FINE/Turbo, OpenFoam**
6. Fusion: **XGC, GTC, GTC-P, M3D, M3D-C1, M3D-K, Pixie3D**
7. Geoscience: **SPECFEM3D_GLOBE, AWP-ODC, RTM**
8. Materials Science: **QMCPack, LAMMPS**
9. Medical Imaging: **Cancer pathology**
10. Quantum Turbulence: **QLG2Q**
11. Relativity: **Maya**
12. Weather: **GRAPES**
13. Visualization: **Paraview, Visit, VTK, ITK, OpenCV, VTKm**

Impact on Industry :

- **NUMECA (FINE/Turbo)** – Allowed time-varying interaction of turbomachinery-related aerodynamic phenomena
- **TOTAL (RTM)** – Allowed running of higher fidelity seismic simulations
- **FMGLOBAL (OpenFoam)** – Allowed running higher fidelity fire propagation simulations

Over 1B LCF hours from ADIOS enabled Apps 2015
Over 1,500 citations

LCF/NERSC Codes in red

2013 R&D 100 Award Winner

Adaptable I/O System for Big Data (ADIOS)

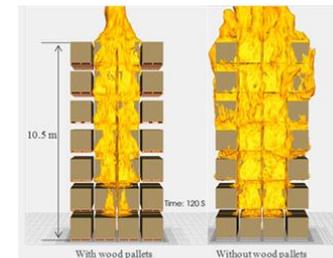
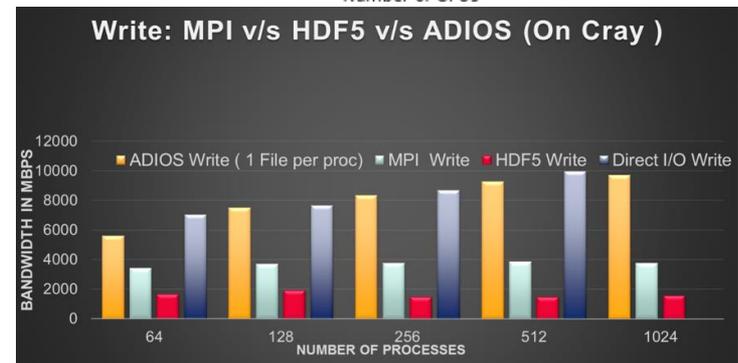
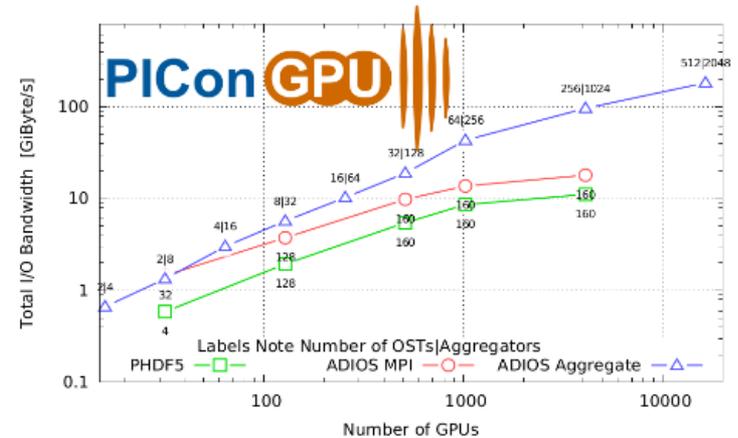
ADIOS is a portable, scalable, easy-to-use software framework conceived to solve “big data” problems. For scientists making use of high performance computers, ADIOS significantly reduces the input or output complexities typically encountered and reduces the time to solution, so researchers spend less time managing data. The software streamlines workflows and lays the foundation for exascale supercomputers to be able to run multiple tasks simultaneously.

The research was funded by DOE's Oak Ridge Leadership Computing Facility, the Office of Advanced Scientific Computing Research, the Office of Fusion Energy Science, and the National Science Foundation.

The ORNL team consisted of (seated) Norbert Pochossat, Gu Yuan Tian; (standing) Jong Youl Choi, Hasan Abbas, Jeremy Logan, Scott Klasky; and (not pictured) Roselyne Tchoua. Also pictured are Karsten Schwab and Matthew Wolf (Joint Faculty, Georgia Institute of Technology), Manish Parashar (Rutgers University), Nazwa Samatova (Joint Facility, North Carolina)

Impact to other LCF applications

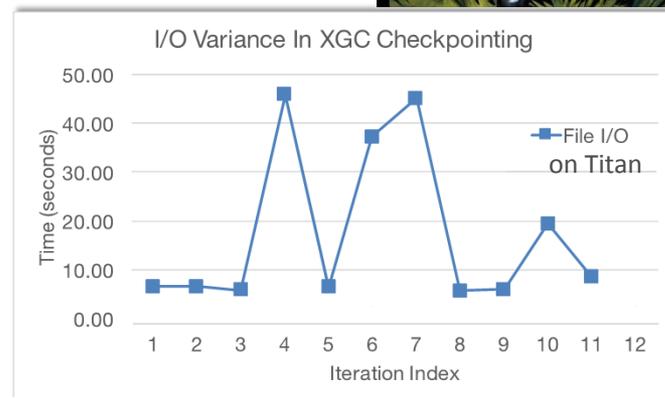
- Accelerators – PIconGPU
 - M. Bussmann, et al. - HZDR
 - Study laser-driven acceleration of ion beams and its use for therapy of cancer
 - Over 184 GB/s on 16K nodes on Titan
- Seismic Imaging – RTM by Total Inc.
 - Pierre-Yves Aquilanti, TOTAL E&P
TBs as inputs, outputs PBs of results along with intermediate data
- FMGLOBAL using OpenFOAM with ADIOS
 - “ADIOS was able to achieve a 12X performance improvement from the original IO “
 - <https://www.olcf.ornl.gov/2016/01/05/fighthing-fire-with-firefoam/>



Stacking commodities on wood pallets slows horizontal fire spread, versus absence of pallets

I/O Variability on LCF systems

- **Static** assumptions about I/O provisioning are **sensitive** to point contention
 - Aggregation with write-behind strategy
 - Stripe alignment: to avoid contention
 - Slowdown on a single node can stall I/O



IO PERFORMANCE VARIABILITY DUE TO EXTERNAL INTERFERENCE

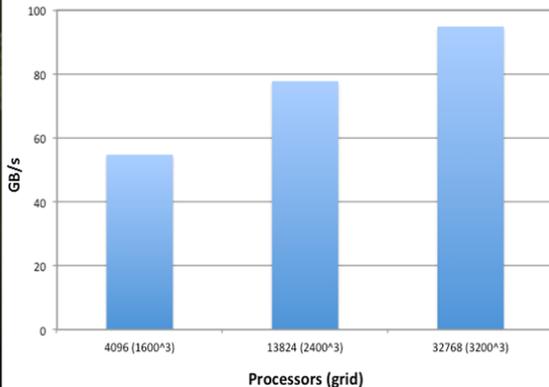
Machine	Number of Samples	Avg. IO Bandwidth (MB/sec)	Std. Deviation	Covariance
Jaguar	469	1.78e+4	1.07e+4	60.09%
Franklin	2581	6.22e+3	2.50e+3	40.22%
XTP(with Int.)	400	7.89e+2	3.44e+2	43.68%
XTP(without Int.)	320	1.44e+3	1.28e+2	8.86%

- J. Lofstead, F. Zheng, Q. Liu, S. Klasky, R. Oldfield, T. Kordenbrock, K. Schwan, M. Wolf, *Managing variability in the IO performance of petascale storage systems* in *Proceedings of the 2010 ACM/IEEE International Conference for High Performance Computing, Networking, Storage and Analysis*, IEEE Computer Society, pp. 1–12.
- Q. Liu, N. Podhorski, J. Logan, S. Klasky, *Runtime I/O ReRouting + Throttling on HPC Storage* in 5th USENIX Workshop on Hot Topics in Storage and File Systems, USENIX, Berkeley, CA. <https://www.usenix.org/conference/hotstorage13/workshop-program/presentation/Liu>.

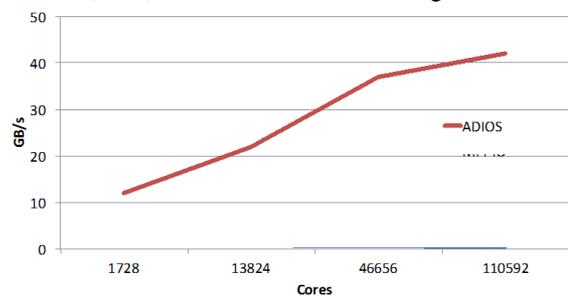
Performance Portability

- ADIOS has methods optimized for different platforms
 - BGQ-GPFS, Cray XK-Lustre, IB clusters
- Different methods can be changed for R/W optimizations
- E.g. Quantum Physics QLG2Q code

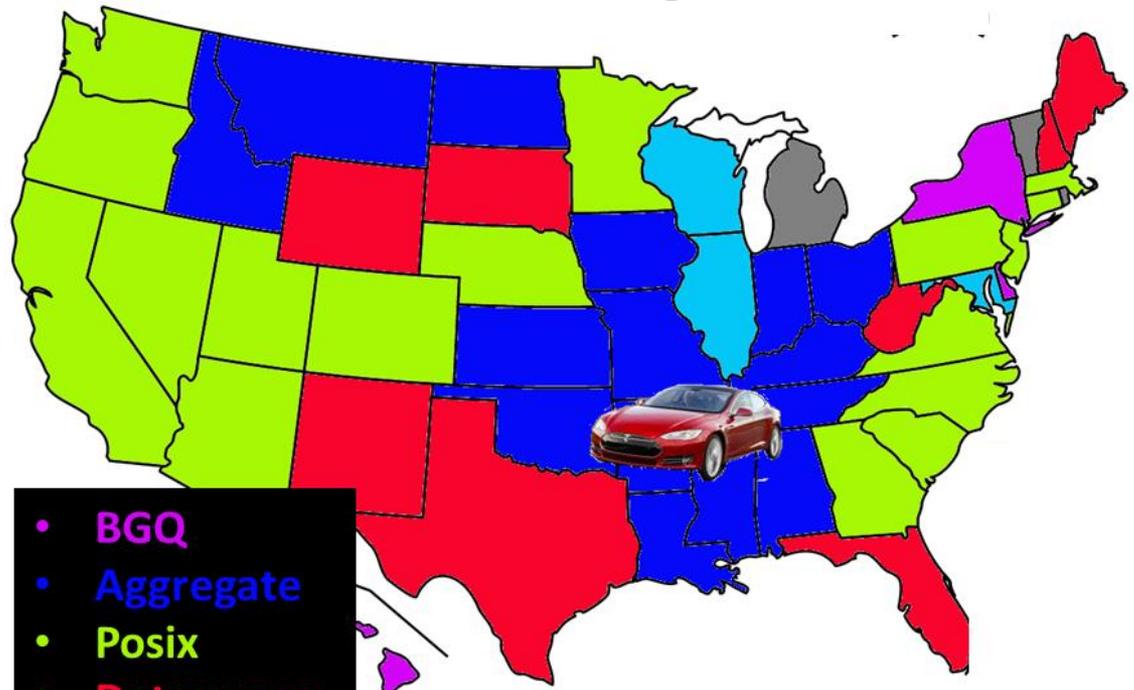
QLG2Q (Quantum Physics) on Titan



QLG2Q with ADIOS vs. MPI-IO on JaguarPF



Performance Island navigation with ADIOS



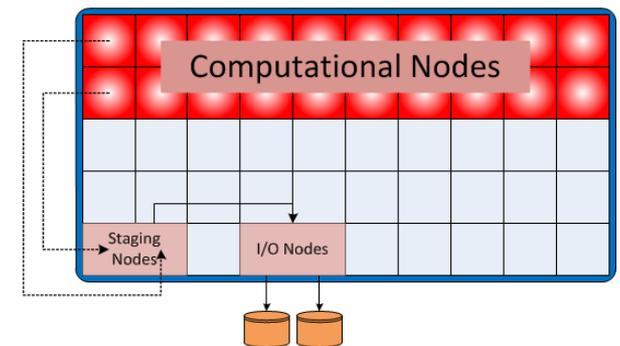
The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- **Staging for Data-in-Motion**
- Refactoring via reduction, queries
- eXascale Service Orient Architecture
- Movement to the Exascale
- VTK-M
- Conclusions



Staging to address I/O performance/variability

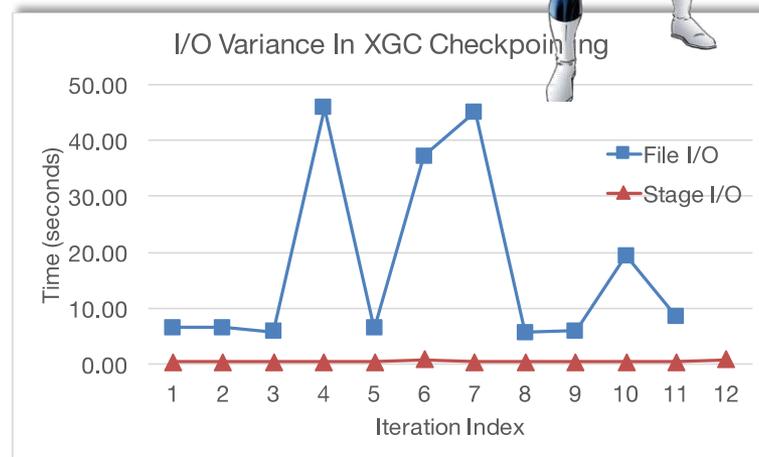
- Simplistic approach to staging
 - Decouple application performance from storage performance (burst buffer)
 - Move data directly to remote memory in a “staging” area
 - Write data to disk from staging
- Built on past work with threaded buffered I/O
 - Buffered asynchronous data movement with a single memory copy for networks which support RDMA
 - Application blocks for a very short time to copy data to outbound buffer
 - Data is moved asynchronously using server-directed remote reads
- Exploits network hardware for fast data transfer to remote memory



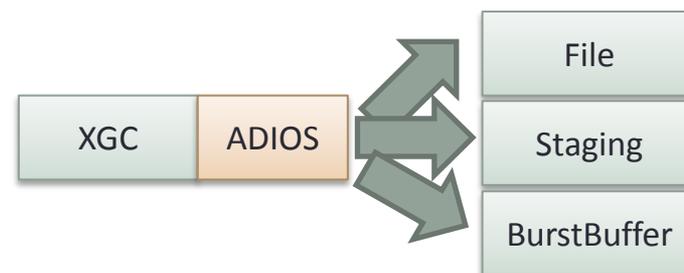
XGC I/O Variability



- Staging of XGC output
 - Observed high I/O variations with files on XGC check-pointing
 - Research on Staging I/O with Rutgers and GA Tech: Decoupling app and I/O
 - Reduced I/O variations with staging methods as well as high I/O throughput
 - Creates new research opportunities with BurstBuffer (ADIOS + Staging + BurstBuffer)
- Impact: I/O **performance** is much more **predictable**

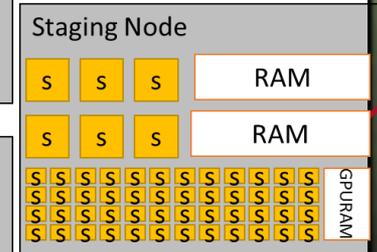


XGC I/O variations on checkpoint writing. Compared File I/O and Staging I/O on Titan.



Staging

Petascale Staging: 6% I/O overhead



The Story

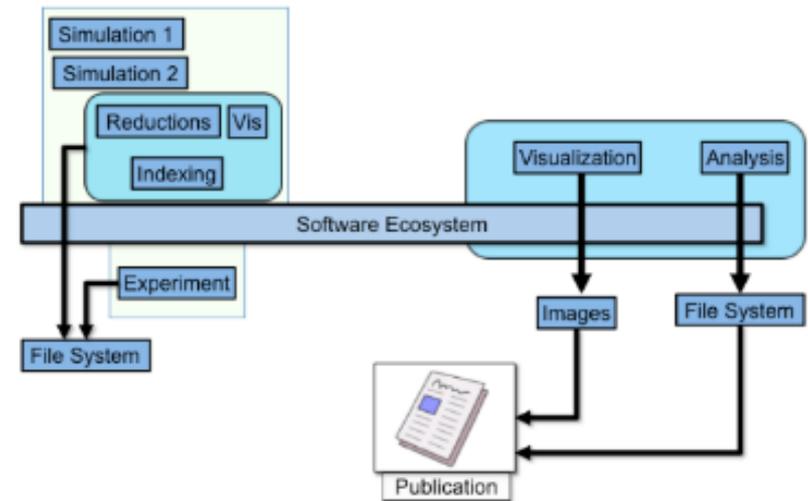
- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- eXascale Service Orient Architecture
- Movement to the Exascale
- VTK-M
- Conclusions



Where Do We Spend Our Time in Science?

- Goals

- Make the process **predictable**
- Make the process **adaptable**
- Make the process **scalable**
- Make the software **easy-to-use**



- Observation

- Too much time is spent in managing, moving, storing, retrieving , and turning the science data into **knowledge**

- Our specific Goals

- Refactor data so that we can efficiently store, find, retrieve data in predictable times set by the users, negotiated with the system

How to fill the gap

- Scientific data can be modelled and refactored
 - Exploit structure to optimize data storage and processing
 - Split data into blocks with varying precision
 - Remember how data was originally created to regenerate on demand
- **Move Information**
 - Data compression
 - Possible loss, but within acceptable bounds from our understanding of the physics model of the simulation and/or measurement errors
 - Change data \rightarrow data = model + information
 - $F = F + \delta F \approx F + \delta F_1 + \delta F_2 + \varepsilon$



Emphasize utility + information

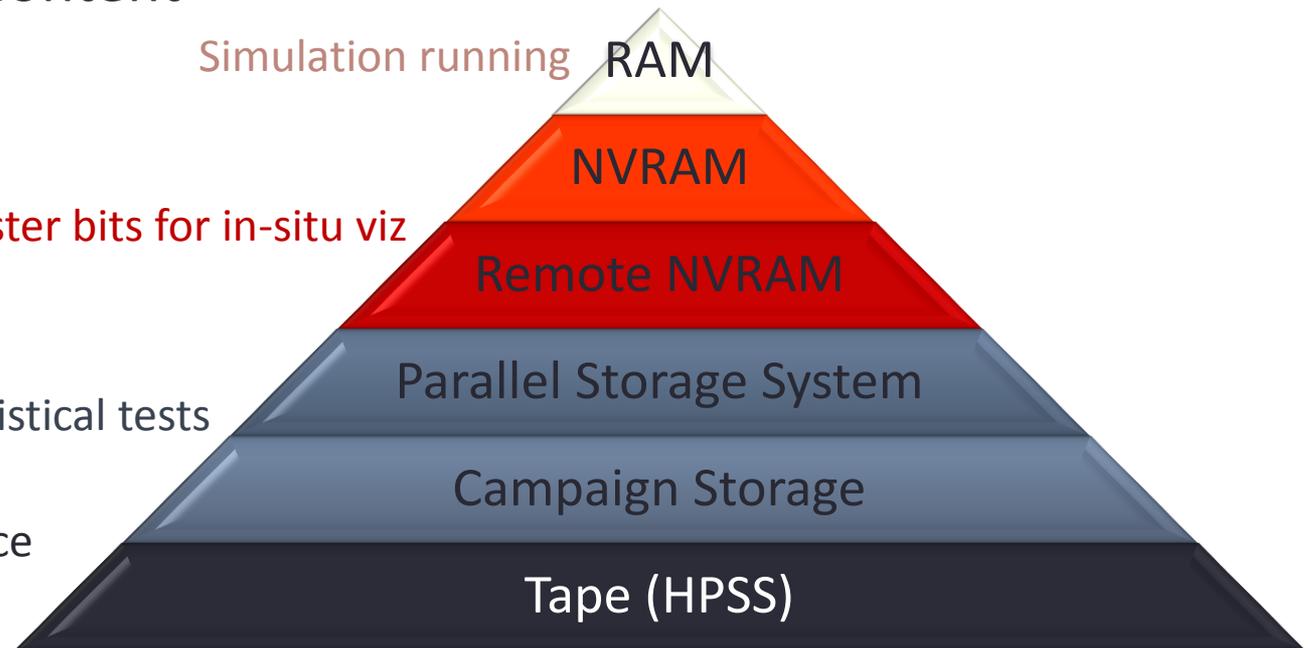
- What do I need where? When / how fast do I need it? How precisely do I need it?
- Utilize data storage hierarchy to address variation in demands
- Couple with data compression algorithms, which reduce the movement + storage cost (size and time), but retain most information content



Fewer, faster bits for in-situ viz

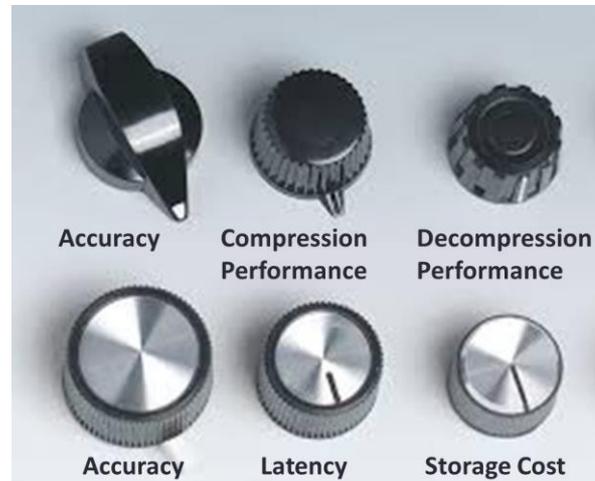
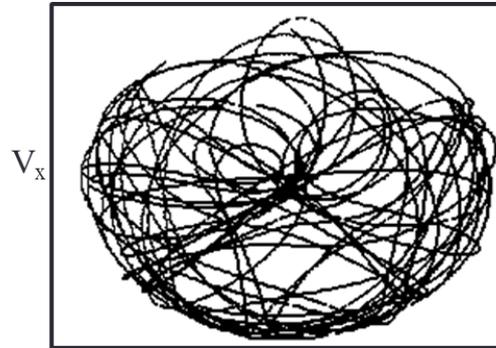
Extra precision for statistical tests

Longest persistence



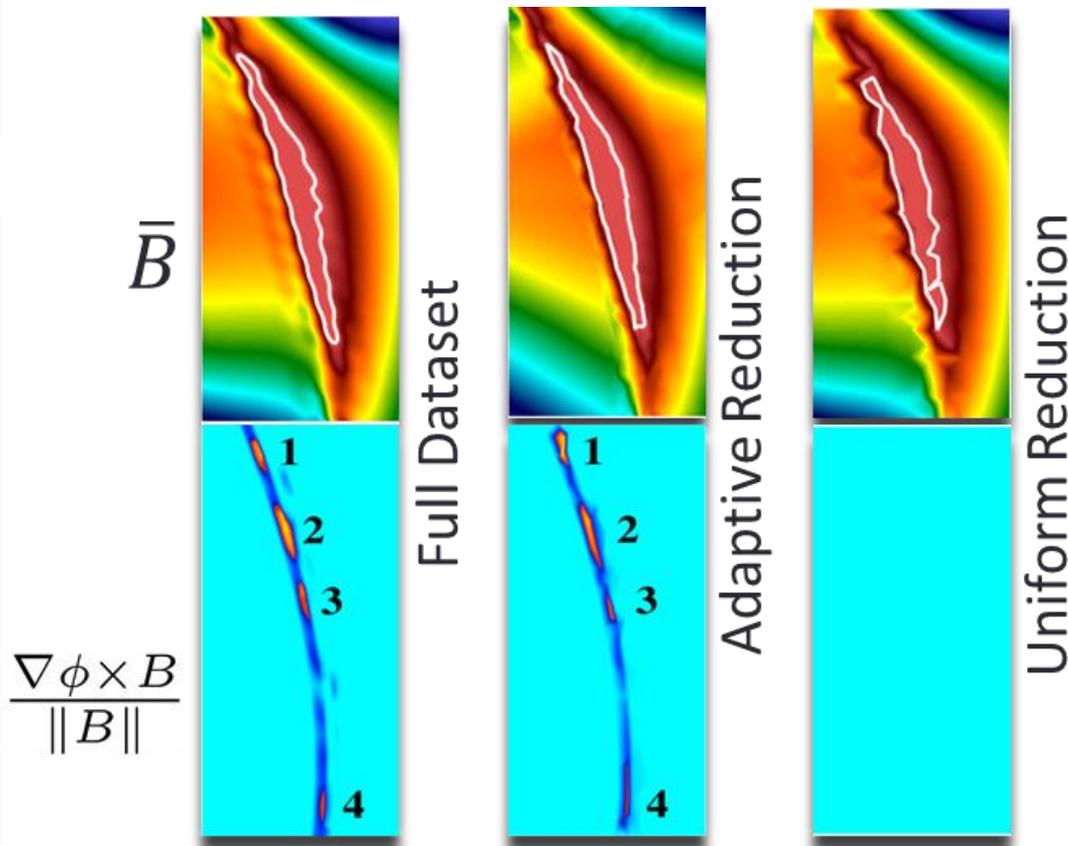
Question

- How can we move, reduce, and process this dataset, $f(x, v_x)$, which is 10 GB per timestep and has 100 time-steps, to < 1MB, moving it in <10 s over the WAN, with 100% accuracy?
 - And then performing different analysis on this data?



Data reduction challenges

- Most publications refer to error as $E \equiv \|f - \tilde{f}\|$
- But $E \equiv \left\| (g_i^t(f(\vec{x}, t)) - \widetilde{g_i^t(f_i^t(\vec{x}, t))}) \right\|$
- What's the impact of errors on derived quantities?



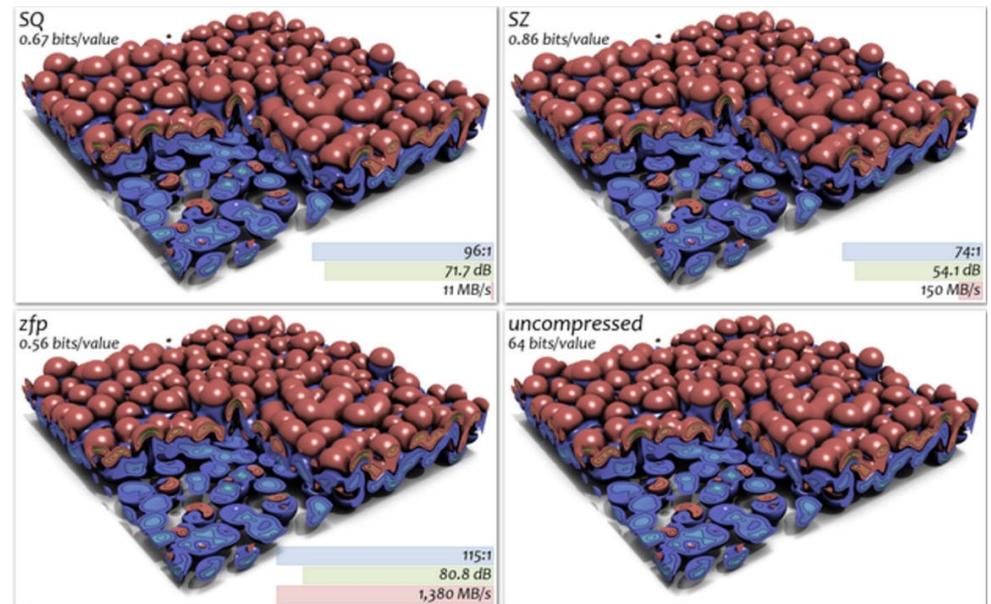
Auditing data during the streaming phase

- Streaming Information
 - Too much data to move in too little time
 - Storage sizes/speed doesn't keep up with making NRT decisions
- Fit the data with a model
 - From our physics understanding
 - From our understanding of the entropy in the data
 - From our understanding of the error in the data
 - Change data into $\text{data} = \text{model} + \text{information}$ ($f = F + \delta f$)
- Streaming Information
 - Reconstruct “on-the-fly”
- Query data from the remote source



ZFP: lossy JPEG-like floating-point array compression

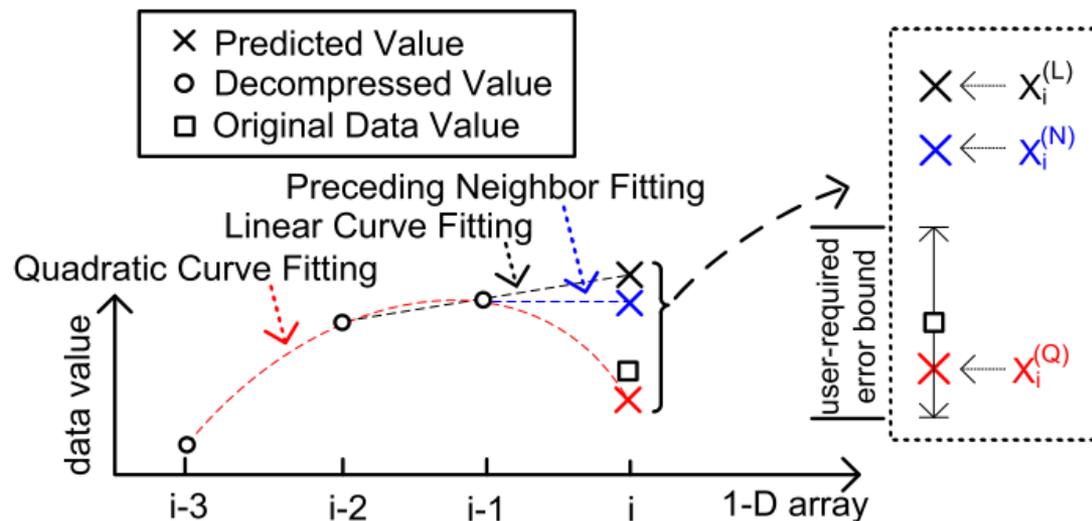
- Discrete Cosine Transform transform, adjustable precision of coefficients kept
- Good for regular grids, “smooth” data
- Integrated with ADIOS
- Natural bit prioritizing
- Fast
- No hard-bound on errors



Peter Lindstrom, Fixed-Rate Compressed Floating-Point Arrays, *IEEE Transactions on Visualization and Computer Graphics* **20**, 2674 (2014).

SZ: curve fitting ++ reduction of multiple snapshots

- Predictable data replace with selected best-fit function
- Unpredictable data compressed with lossy binary representation
- Adjustable hard-bounds on error
- Very fast decompression
- Coming to ADIOS



EVALUATION USING ATM WITH FAIRLY LARGE DATA SIZE (1.5TB)

	error bound $\delta=10^{-4}$			error bound $\delta=10^{-6}$		
	CR	Cmpr time	Decmpr time	CR	Cmpr time	Decmpr time
ZFP	3	27166 sec	30395 sec	2.3	31627 sec	36254 sec
SZ	5.4	43980 sec	6598 sec	4.02	51951 sec	7788 sec

S. Di and F. Cappello, Fast Error-Bounded Lossy HPC Data Compression with SZ, in 2016 IEEE International Parallel and Distributed Processing Symposium (IPDPS), pages 730–739, 2016.

The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- **eXascale Service Orient Architecture**
- Movement to the Exascale
- VTK-M
- Conclusions



XSSA: eXtreme Scale Service Architecture

- Philosophy based on **Service-Oriented Architecture**
 - System management
 - Changing requirements
 - Evolving target platforms
 - Diverse, distributed teams
- Applications built by assembling services
 - Universal view of functionality
 - Well defined API
- Implementations can be easily modified and assembled
- **Manage complexity while maintaining performance, scalability**
 - Scientific problems and codes
 - Underlying disruptive infrastructure
 - Coordination across codes and research teams
 - End-to-end workflows



Services

- Execution Environments
 - On simulation nodes
 - Shared resources with simulation
 - Dedicated resources
 - On dedicated nodes (staging)
 - Post hoc services
- Data Models
 - Codes are different
 - Limit data copy and conversion
 - Support data reductions
 - Integration with SIRIUS

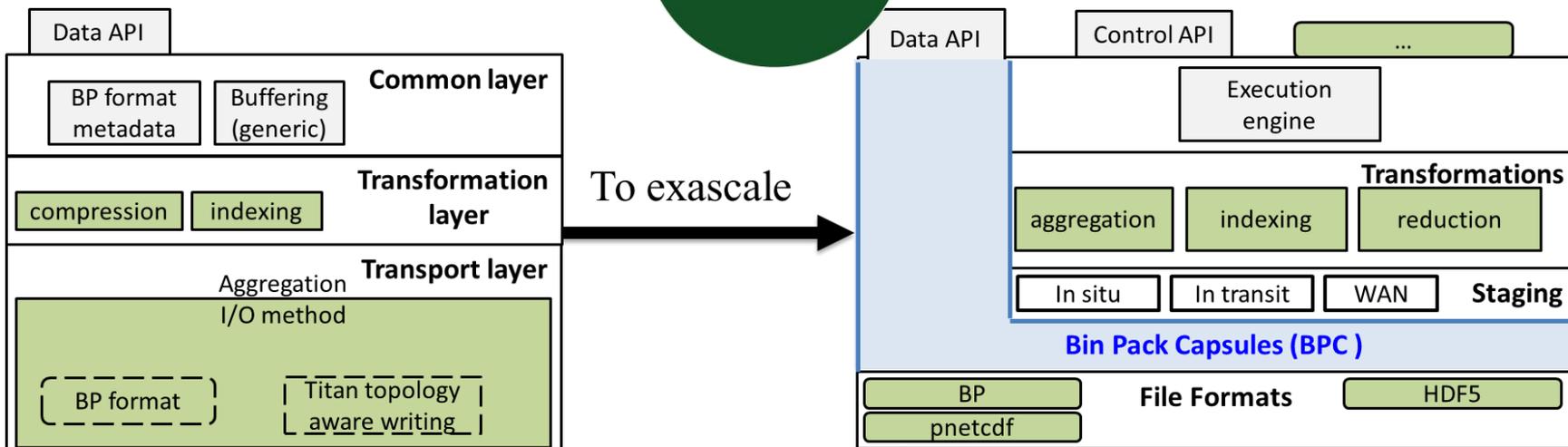
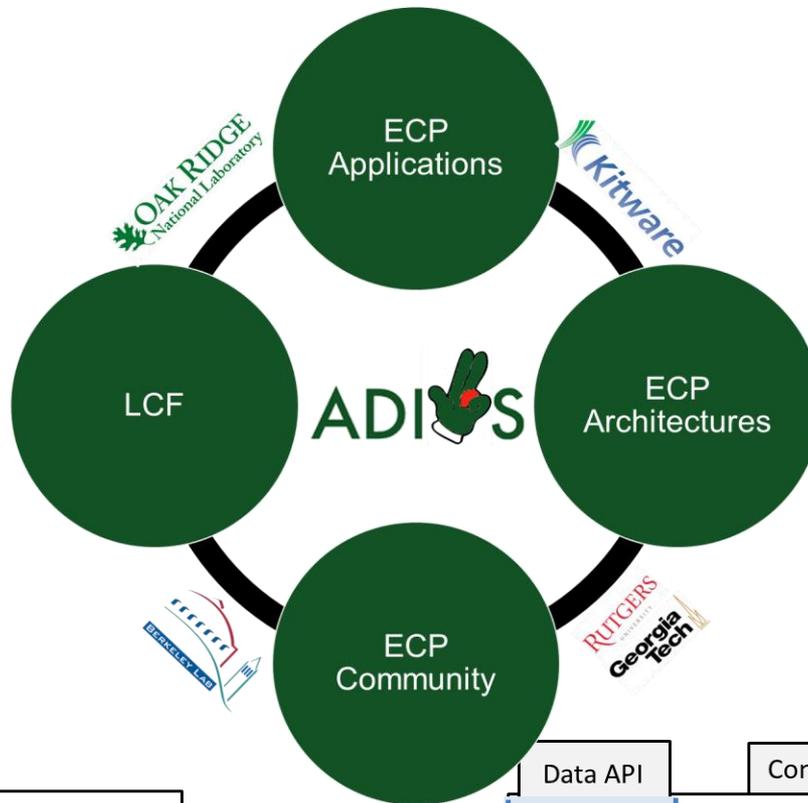


The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- eXascale Service Orient Architecture
- **Movement to the Exascale**
- VTK-M
- Conclusions

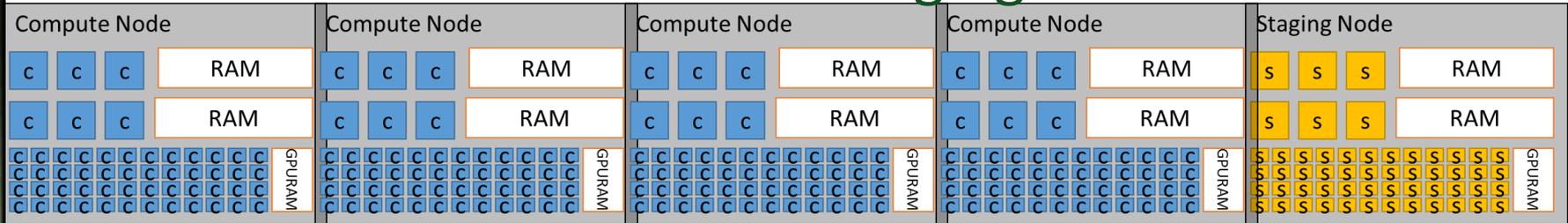


ADIOS 2

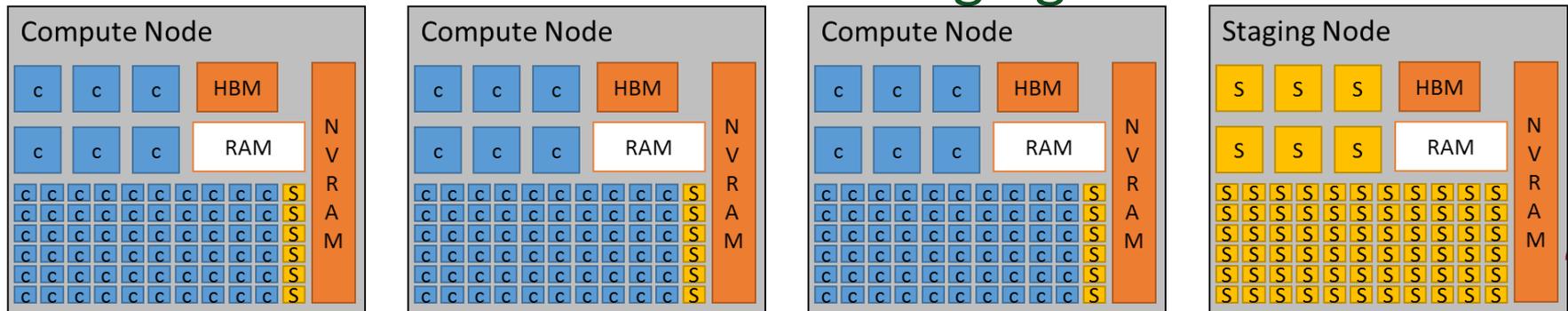


Staging

Petascale Staging

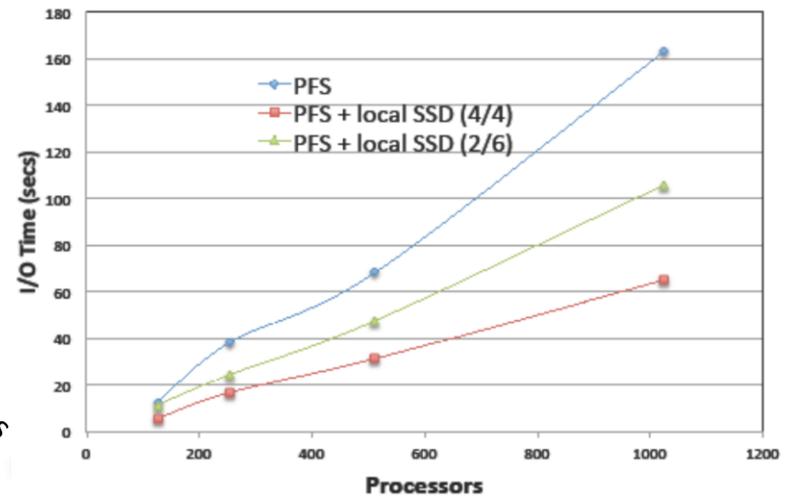
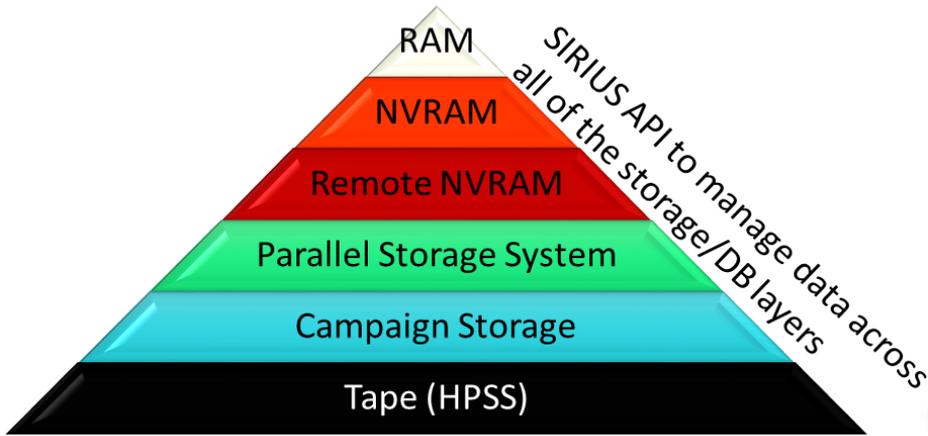


To Exascale Staging

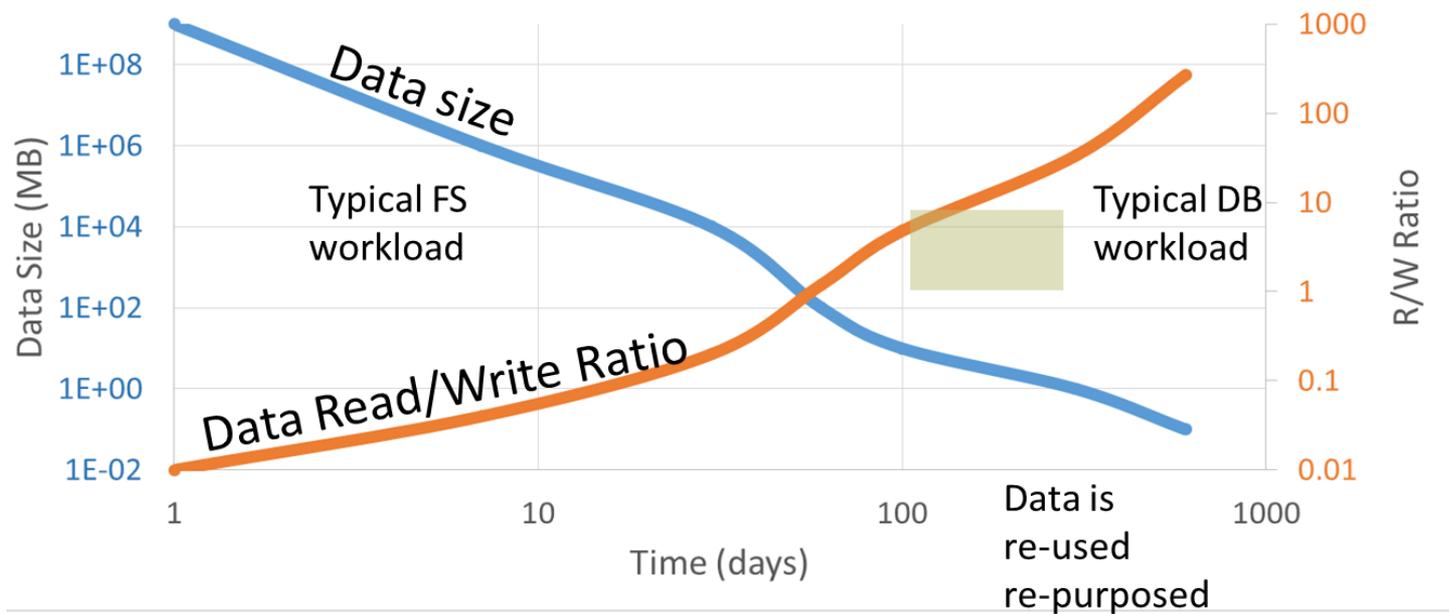


- Use compute and deep-memory hierarchies to optimize overall workflow for power vs. performance tradeoffs
- Abstract complex/deep memory hierarchy access
- Placement of analysis and visualization tasks in a complex system
- Impact of network data movement compared to memory movement
- Abstraction allows staging
 - On-same core
 - On different cores
 - On different nodes
 - On different machines
 - Through the storage system

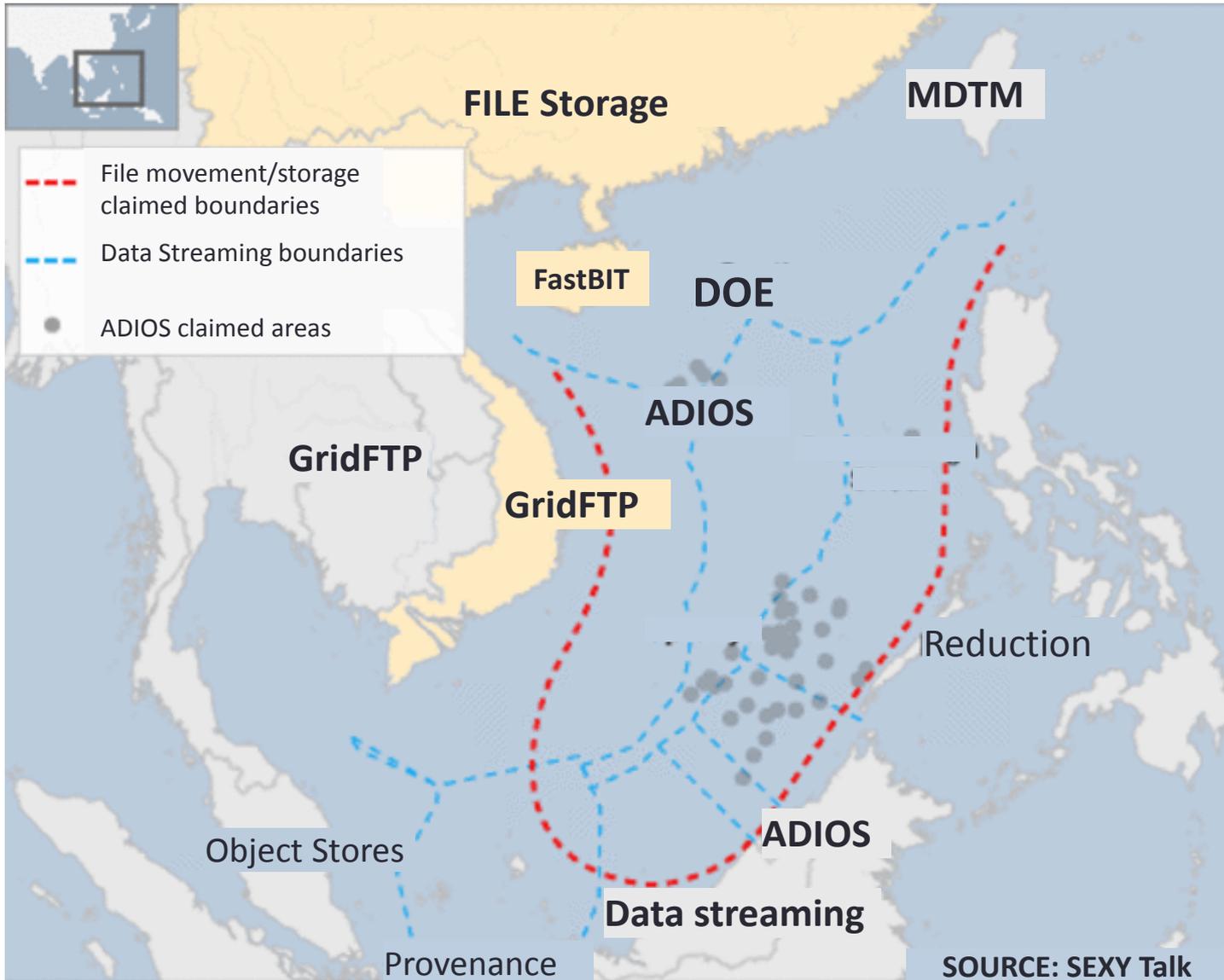
Abstractions across File System to DB



Evolution of Data to Information



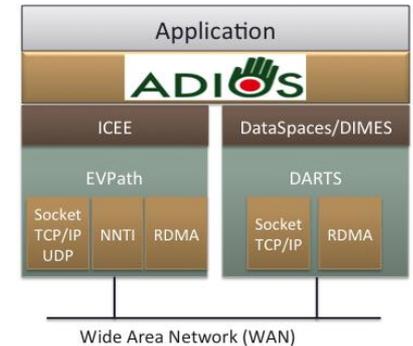
Files vs. Streams



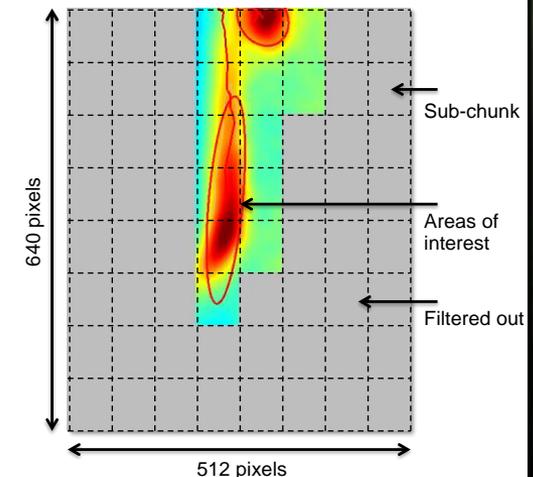
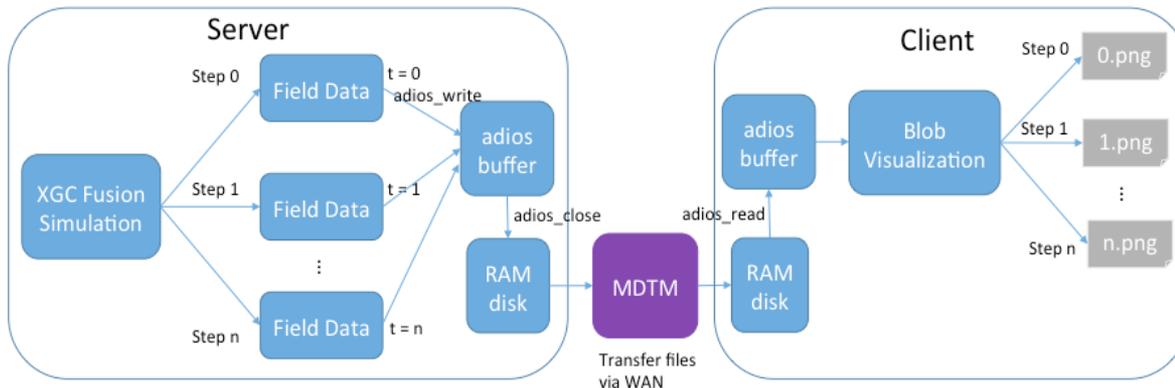
Data Staging for EOD data

- ICEE

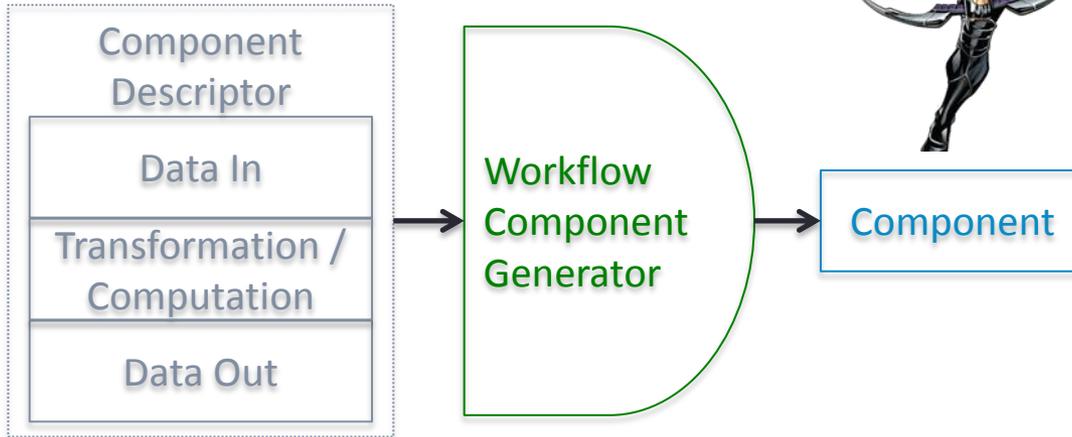
- Using EVPath package (GATech)
- Support uniform network interface for TCP & RDMA
- Allows us to stage data over the LAN/WAN
- Dataspaces + sockets/IB
- MDTM now connected for fast data transfer tool, that fully takes advantage of multi-core on DTNs (<https://web.fnal.gov/project/mdtm>)



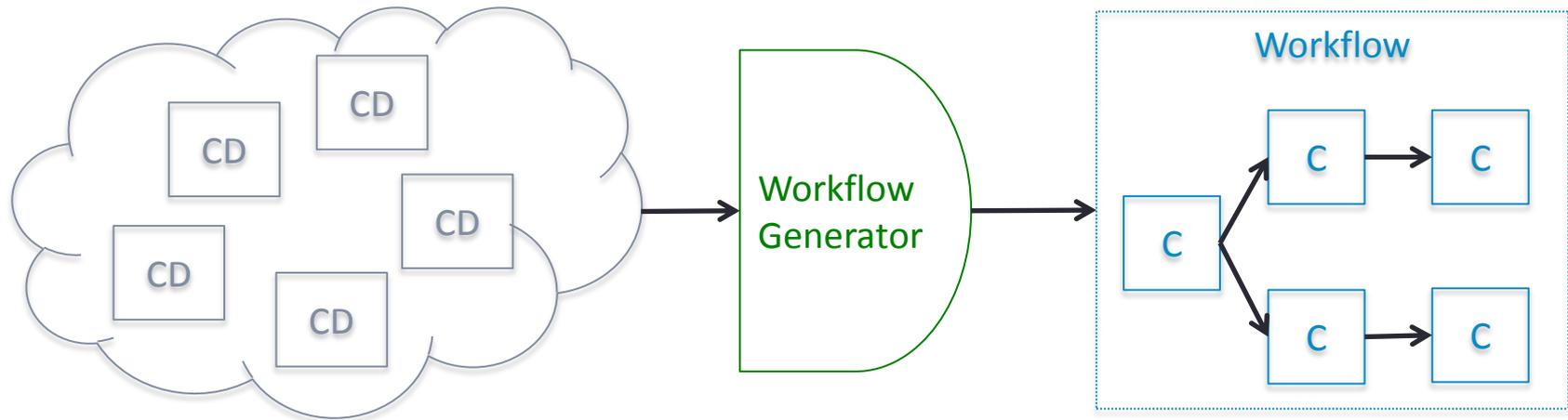
- Select only areas of interest and send
- Reduce payload on average by about 5X



Mini-Application Workflows



- These skeletal workflows are used to analyze system performance.
- Each workflow component is generated from descriptions of computation and data usage
- A complete workflow is generated from a set of component descriptors along with their associated data dependencies
- Allows the system to be developed and tested independently from simulations, analysis and visualization codes



The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- eXascale Service Orient Architect
- Movement to the Exascale
- VTK-M
- Conclusions



VTKm

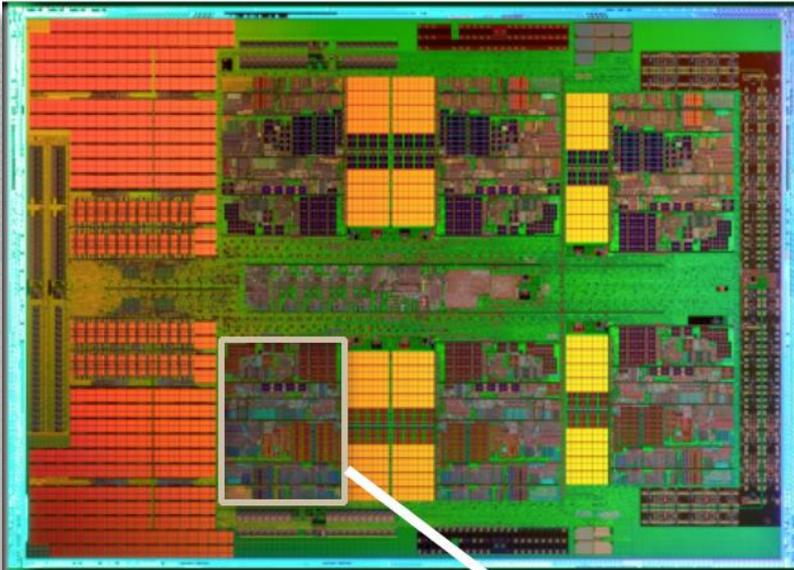


EAVL

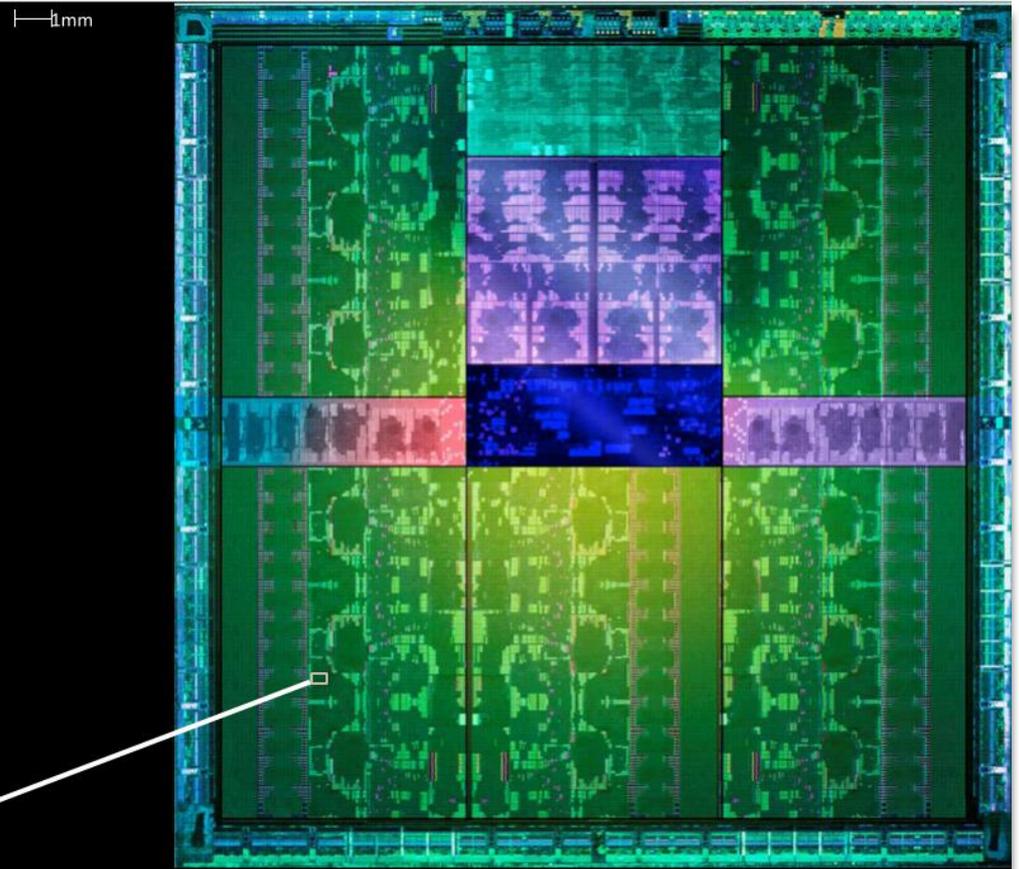
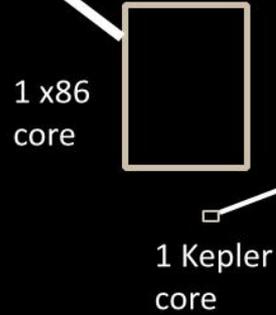
DAX

PISTON

VTKm



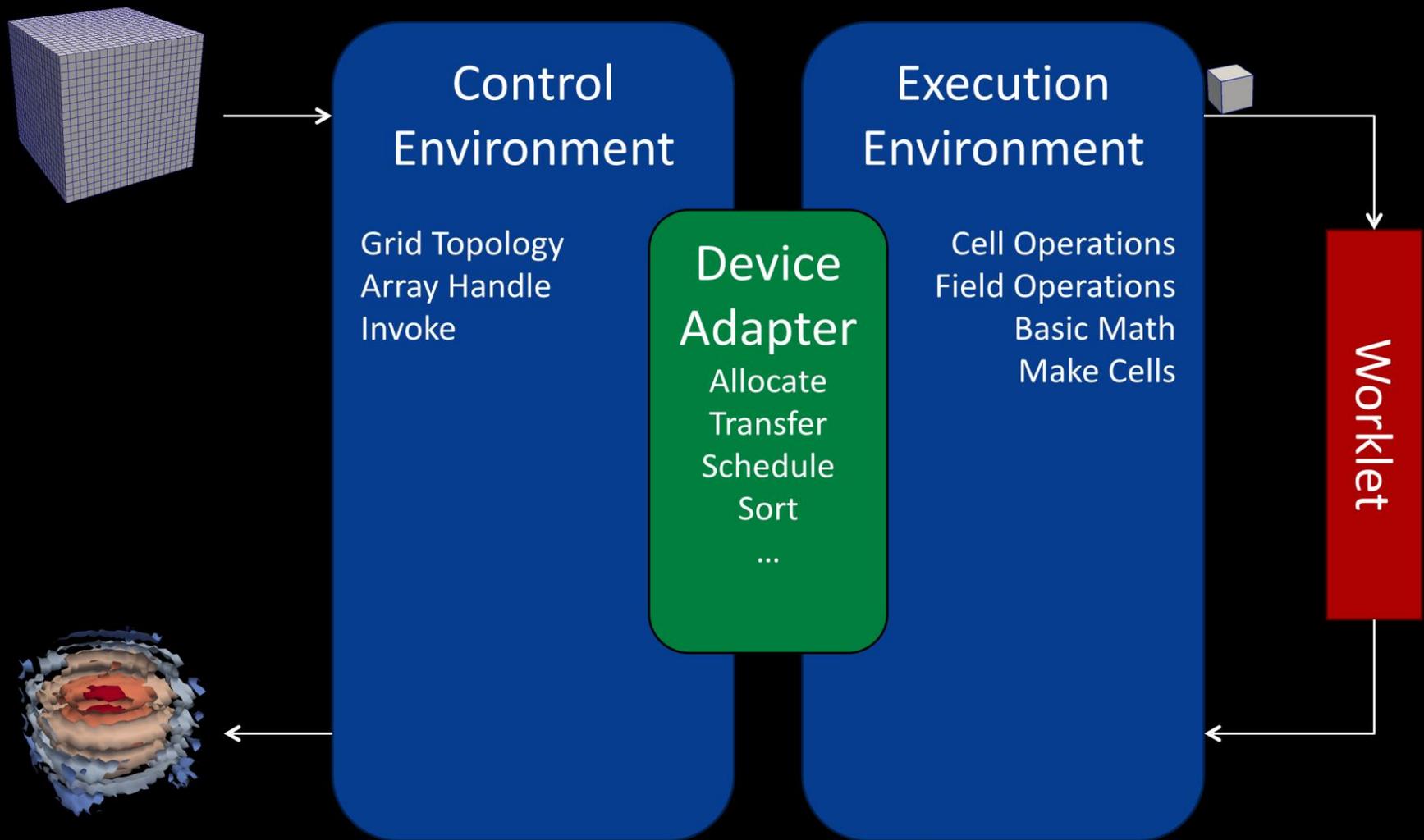
AMD x86
Full x86 Core
+ Associated Cache
8 cores per die
MPI-Only feasible



NVIDIA GPU

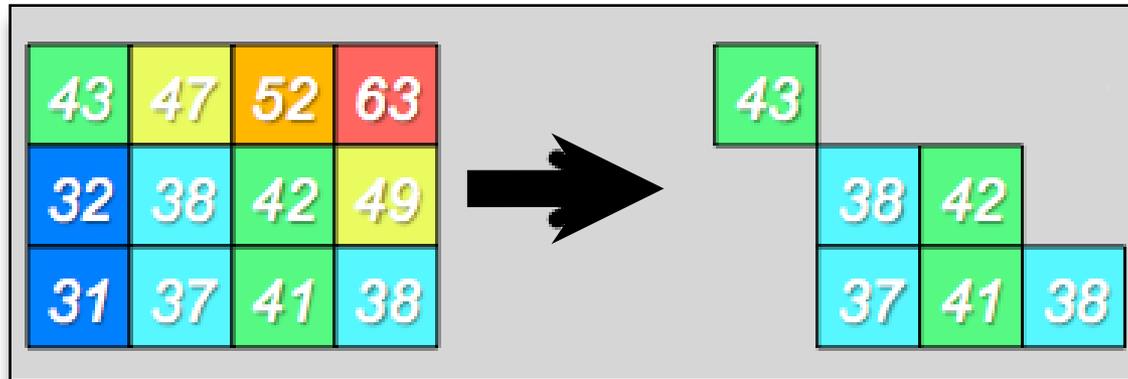
2,880 cores collected in 15 SMX
Shared PC, Cache, Mem Fetches
Reduced control logic
MPI-Only not feasible

VTKm Framework



Example: Memory and Algorithmic Efficiency

Threshold regular grid: $35 < \text{density} < 45$



Traditional Data Model

Fully unstructured grid

- Explicit points
- Explicit cells

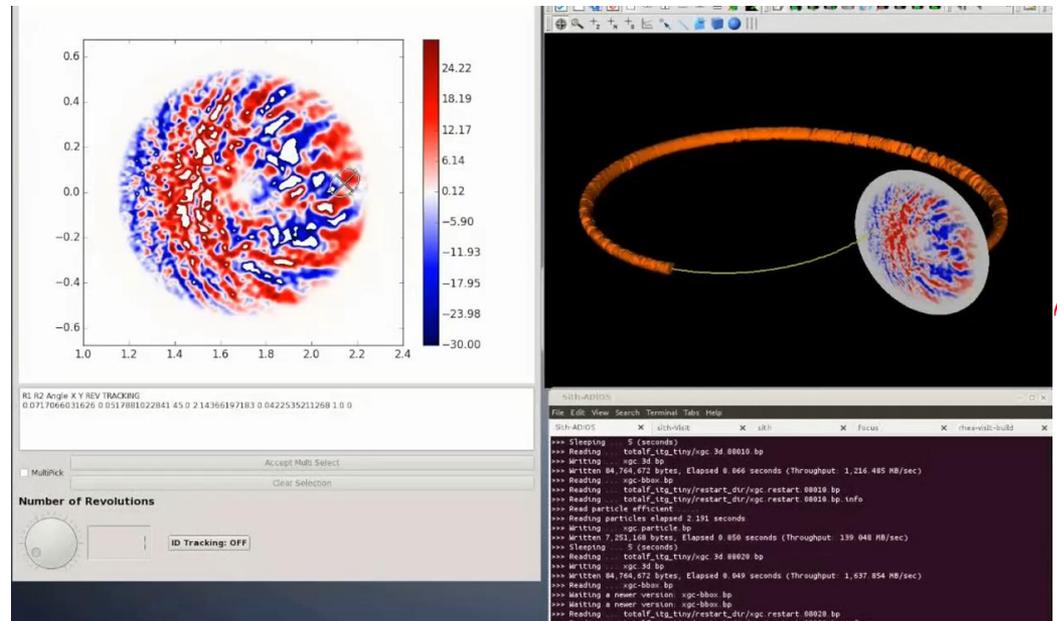
VTK-m Data Model

Hybrid implicit/explicit grid

- Implicit points
- Explicit cells

ADIOS and VTKm

- VTKm and XViz are efforts to prepare for the increasing complexity of extreme-scale visualization
 - Addresses: Emerging Processors, In Situ Visualization, and Usability
 - Minimizes resources required for analysis and visualization
 - Processor/Accelerator awareness provides execution flexibility
- Idea is to incorporate VTK-M into the ADIOS software eco-system



The Story

- Scientific data
- Evolution of HPC hardware
- Community driven software eco-system for the 5Vs
- Our solution
- Staging for Data-in-Motion
- Refactoring via reduction, queries
- eXascale Service Orient Architecture
- Movement to the Exascale
- VTK-M
- **Conclusions**



Conclusions

- Applications from Computation, Experiments, and Observations are pulling us to create new technologies to aid in their data processing
- Technology is pushing us to optimize our solutions on new hardware
- Collaboration in the community on common infrastructure(s) is important
- Movement is from data-at-rest technologies to data-in-motion
 - Is it in situ? (in-place) – no memory movement?
- Computational Science technologies should help applications producing/consuming “BIG” scientific Data

Questions



Acknowledgements

- EPSI
- ICEE
- OLCF
- RSVP
- MONA
- Big Data Express
- Xviz
- VTM-M
- SDAV
- Exact
- SENSEI
- Sirius

