

Exploring energy-performance-quality tradeoffs for scientific workflows with in-situ data analyses

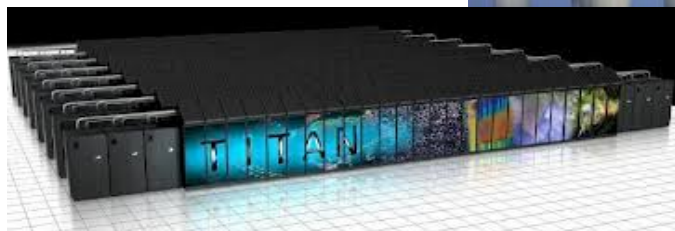
Georgiana Haldeman, **Ivan Rodero***, Manish Parashar, Sabela Ramos, Eddy Z. Zhang, Ulrich Kremer

*Rutgers Discovery Informatics Institute (RDI²)
NSF Cloud and Autonomic Computing Center (CAC)
Rutgers, The State University of New Jersey
email: irodero@rutgers.edu

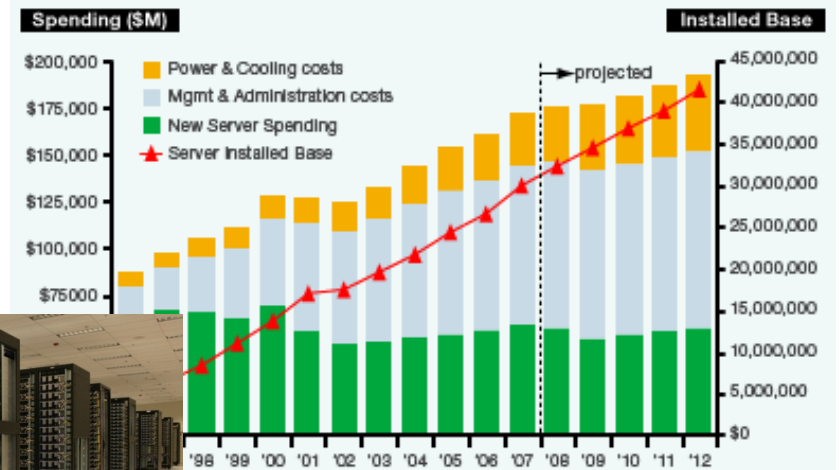
Power/Energy Challenge – Green HPC

- Power has become a critical concern for HPC/supercomputing
 - Impacts operational costs, reliability, correctness
 - End-to-end integrated power/energy management essential
- Increasing scale towards **exascale**
 - Using existing technology would require gigawatt??
 - Nuclear reactor scale??
 - > \$2.5B annual power cost

} **Target < 20MW !!**



Worldwide IT Spending on Servers, Power and Cooling, and Management/Administration



Forming the Datacenter: Consolidation, Pervasive Virtualization and Energy
DIR2009_T4_MB, Mar 2009

Modern Science & Society Transformed by Compute & Data

- New Paradigms & Practices
 - End-to-end: Seamless access, aggregation, interactions
 - Data-driven, Data/Compute-intensive; Age of Digital Observation
 - Integrative, multi-scale, online
- Multi-disciplinary collaborations
 - Individuals, groups, teams, communities, networks
 - New global science culture
- Unprecedented opportunities, challenges



Ack. M. Parashar

Clearly, modern instruments/experiments/... are producing Big Data!!

Large Hadron Collider

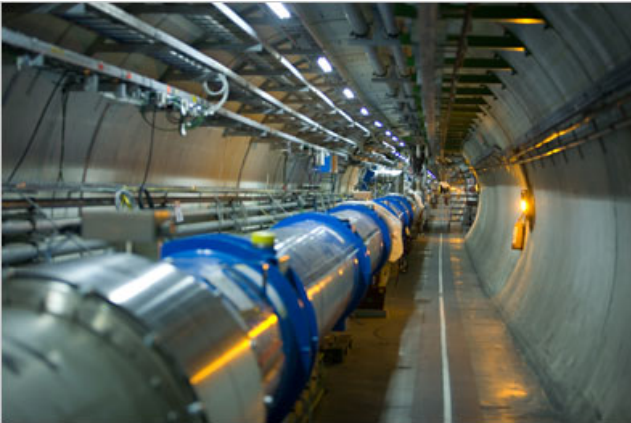


Image credit: Valerio Mezzanotti for The New York Times

Blanco 4m on Cerro Tololo



Image credit: Roger Smith/NOAO/AURA/NSF

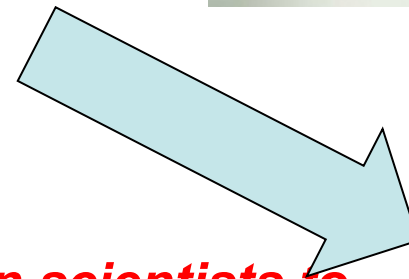
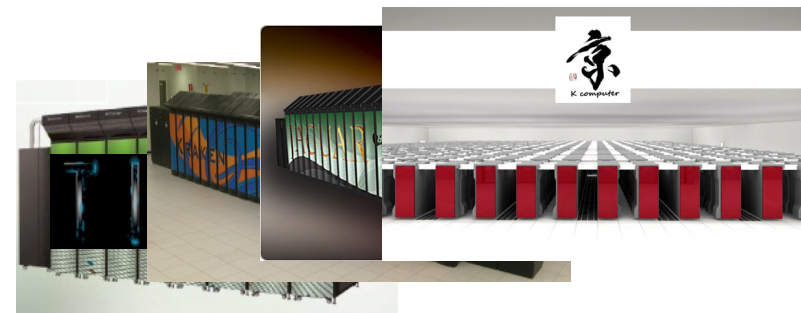
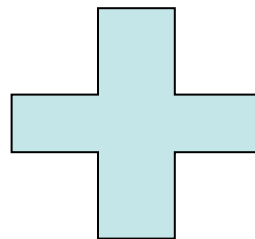
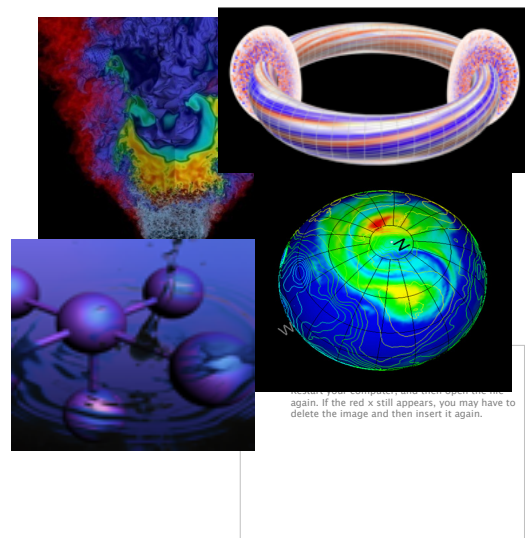
SKA project



Above is proposed image

Scientific Discovery through Simulations

- Scientific simulations running on high-end computing systems generate huge amounts of data!
 - If a single core produces 2MB/minute on average, one of these machines could generate simulation data between ~**170TB** per hour -> ~**700PB** per day -> ~**1.4EB** per year
- Successful scientific discovery depends on a comprehensive understanding of this enormous simulation data



How we enable the computation scientists to efficiently manage and explore extreme scale data: “find the needles in haystack” ??

Challenges Faced by Traditional HPC Data Pipelines

- **Data analysis challenge**

- Can current data mining, manipulation and visualization algorithms still work effectively on extreme scale machine?

- **I/O challenge**

- Increasing performance gap: disks are outpaced by computing speed

- **Data movement challenge**

- Lots of data movement between simulation and analysis machines, between coupled multi-physics simulation components -> longer latencies
- Improving data locality is critical: do work where the data resides!

- **Energy challenge**

- Future extreme systems are designed to have low-power chips – however, much greater power consumption will be due to memory and data movement!

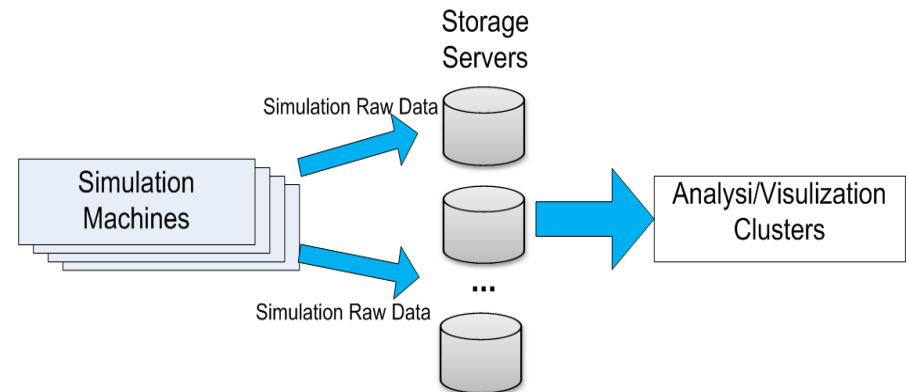


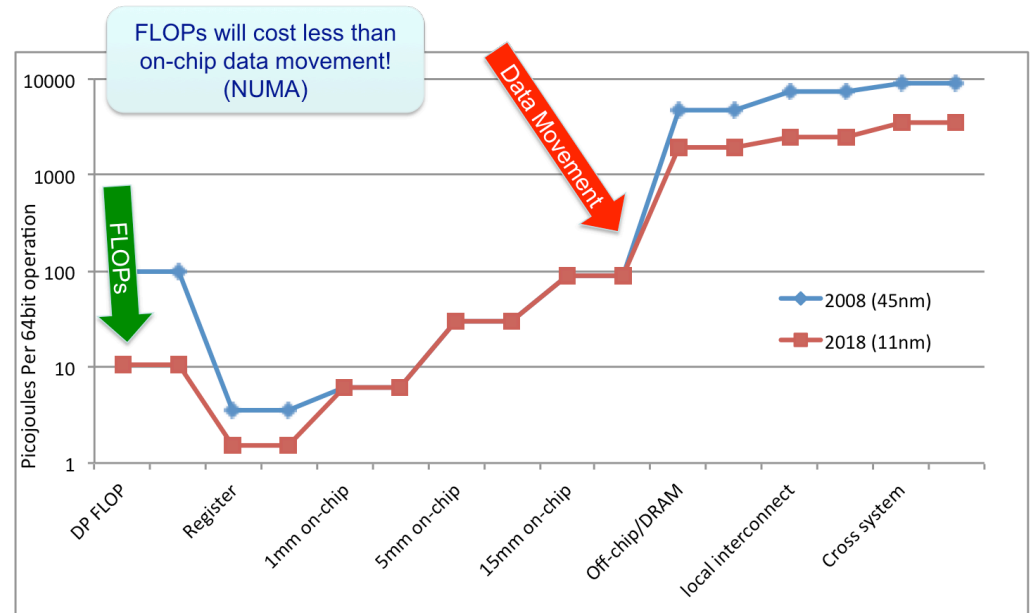
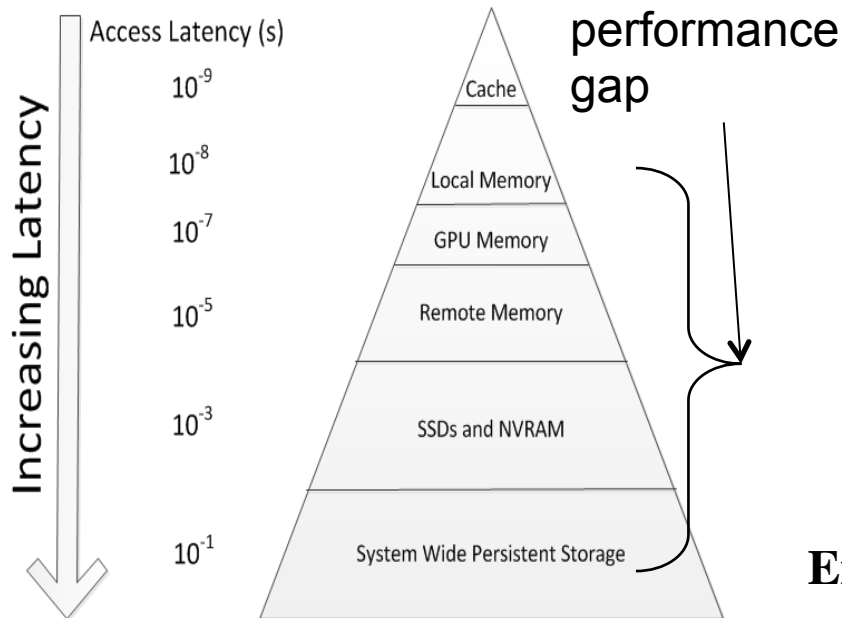
Figure. Traditional data analysis pipeline

The costs of data movement are increasing and dominating!

The Cost of Data Movement

- 40-50% energy spent in off-chip memory hierarchy!!
[Lefurgy, IEEE Computer'03]
- Moving data between node memory and persistent storage is slow!

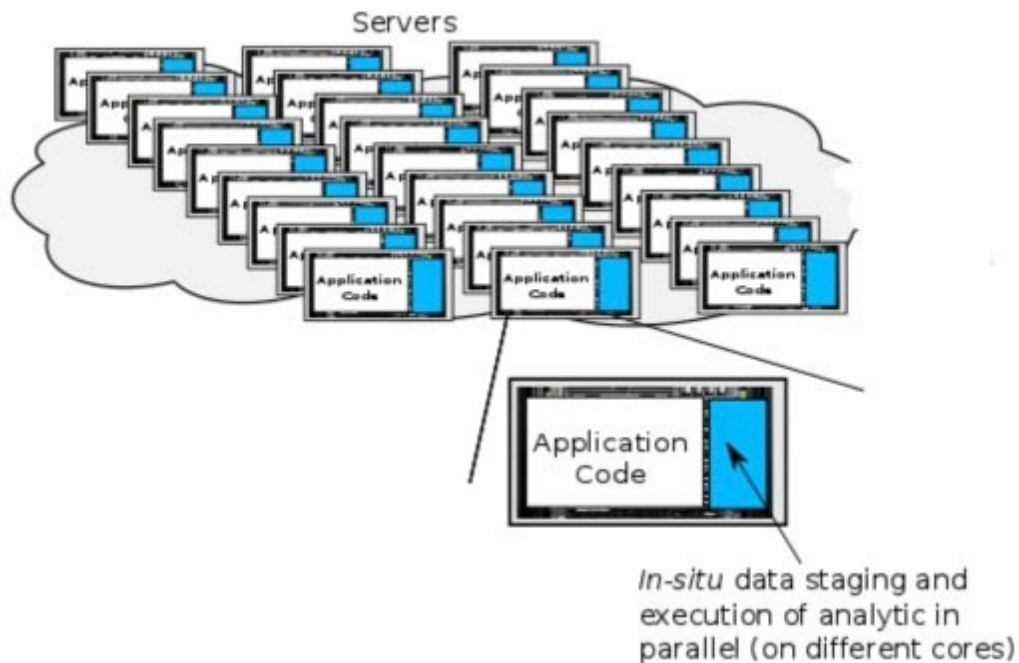
- The energy cost of moving data is a significant concern



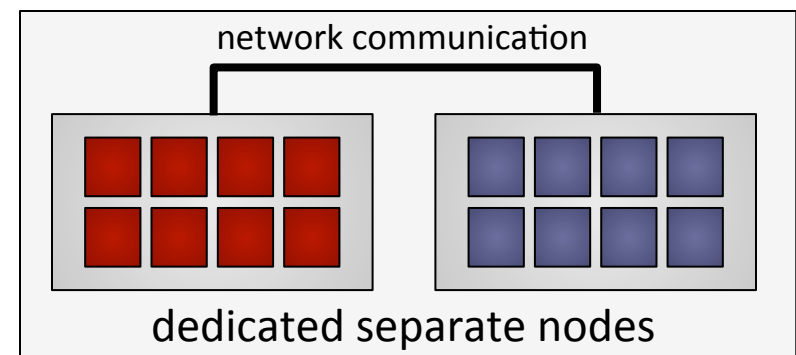
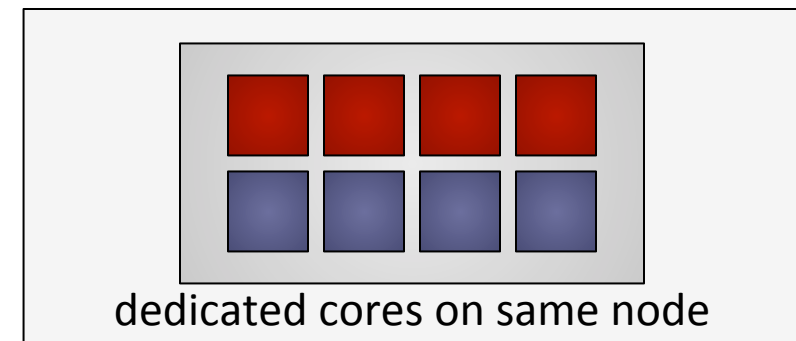
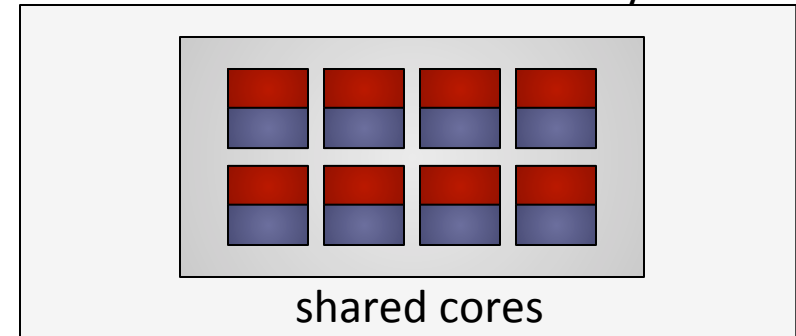
$$\text{Energy_move_data} = \frac{\text{bitrate} * \text{length}^2}{\text{cross_section_area_of_wire}}$$

Rethinking the Data Management Pipeline – In-Situ Data Analytics

- Location of analysis compute resources
 - Same cores as the simulation
 - Dedicated cores on the same node

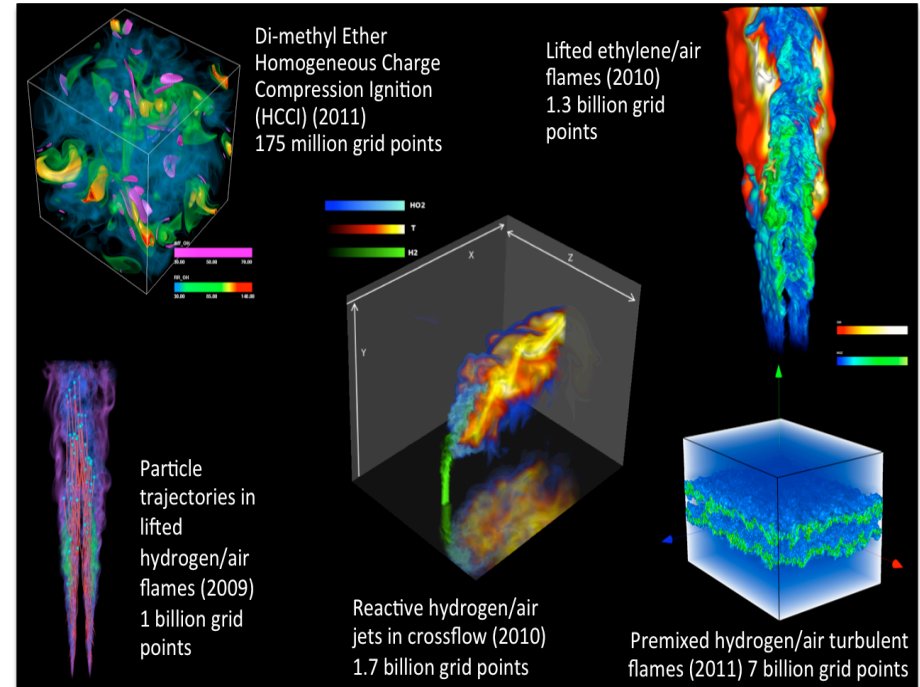


■ simulation ■ analysis



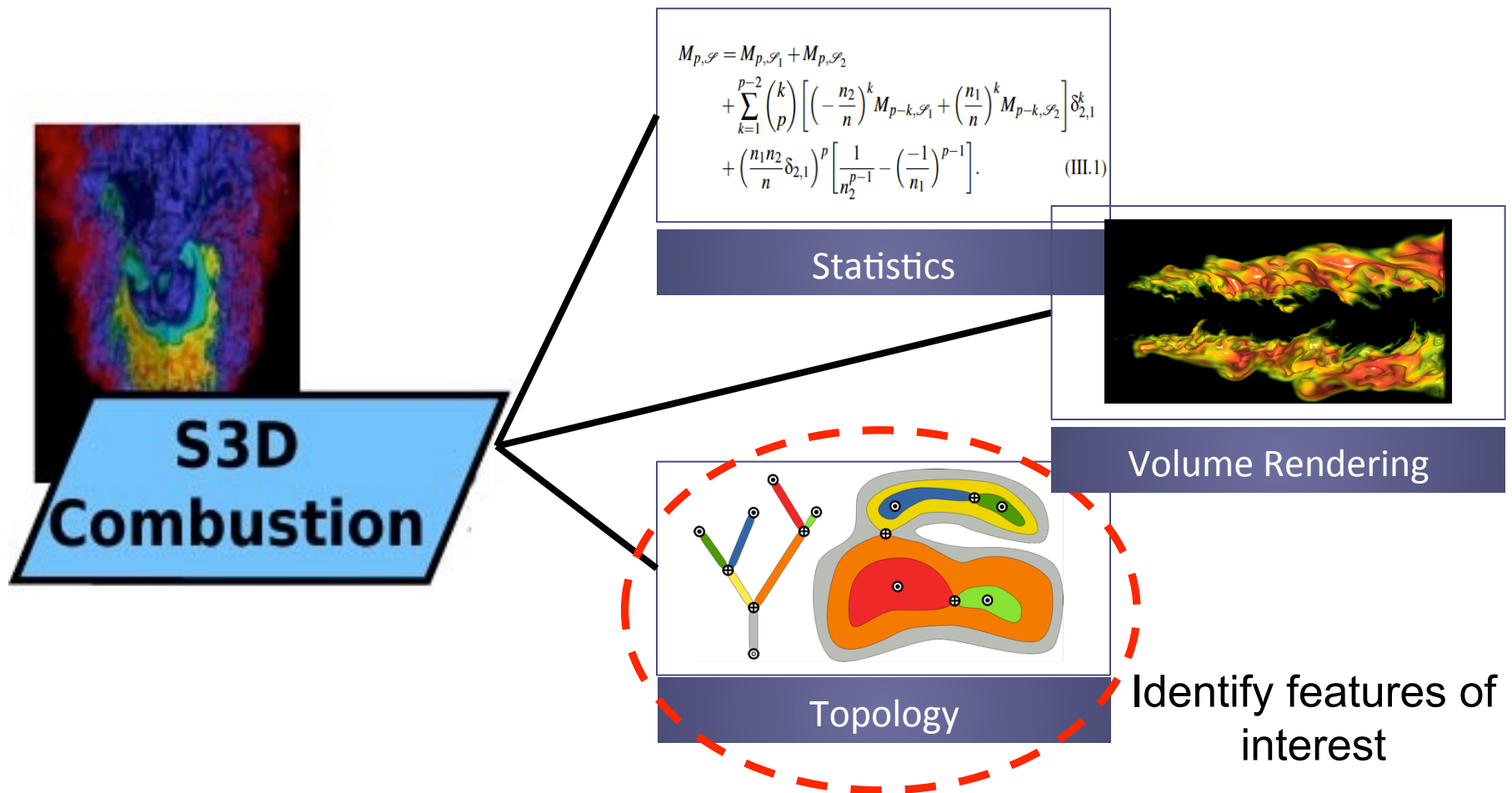
Power Behavior of In-situ Analytics Pipeline

- Combustion simulation workflow with an **in-situ data analytics** pipeline
- With research groups from

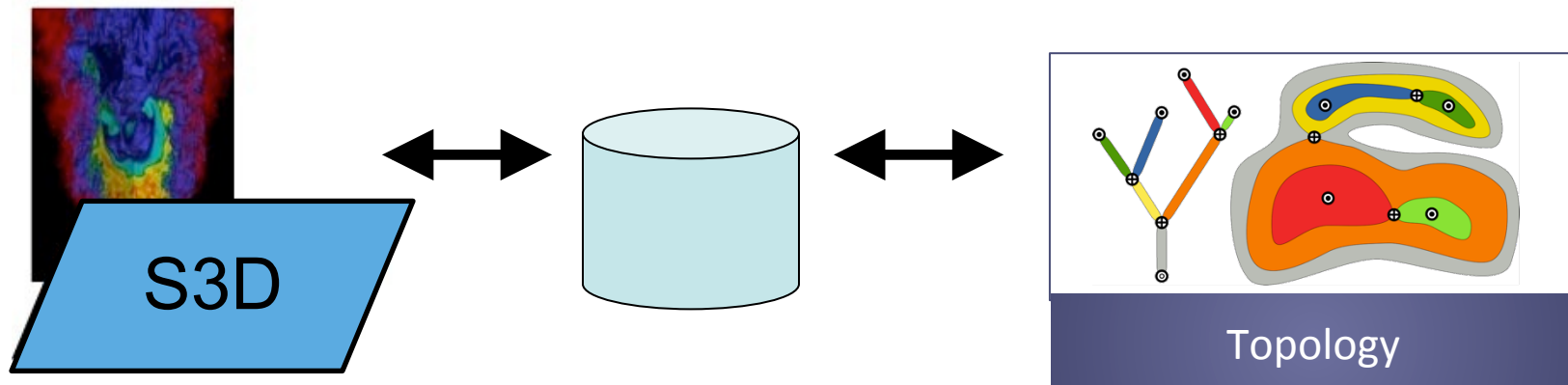


Recent data sets generated by S3D, developed at the Combustion Research Facility, Sandia National Laboratories

In-situ DataAnalysis as Part of S3D



Example: Simulation + Data Analysis Workflow



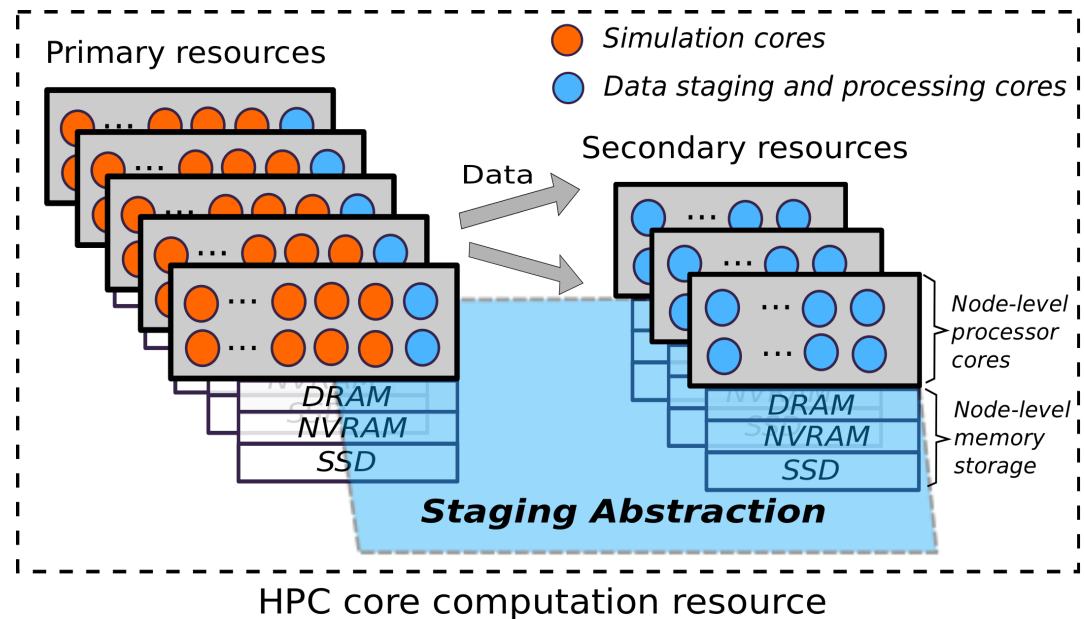
- Modeling data placement and data paths
 - Deep memory hierarchy **1**
- Modeling in-situ analysis choices (cores sharing)
- Opportunities for speculation **2**

1 Data Staging over Deep Memory Hierarchy

- Small DRAM capacity per core – even aggregated memory on dedicated nodes can hardly keep all coupled data (given the ratio of resource allocations for compute nodes and dedicated nodes)

Hybrid Staging

- Spans horizontally across compute nodes
- Spans vertically across the multi-level memory hierarchy, e.g. DRAM/NVRAM/SSD, to extend the capacity of in-memory data staging

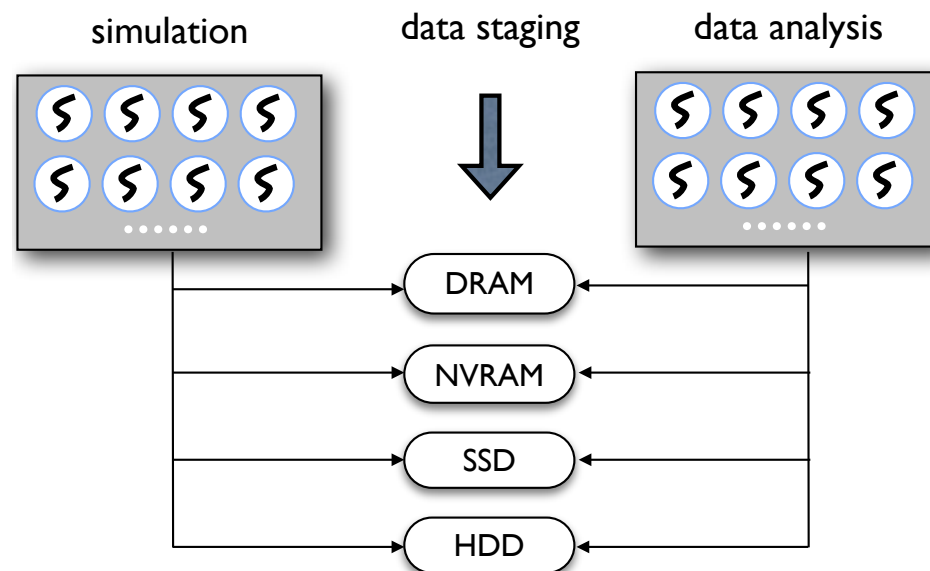


1 Data Staging over Deep Memory Hierarchy

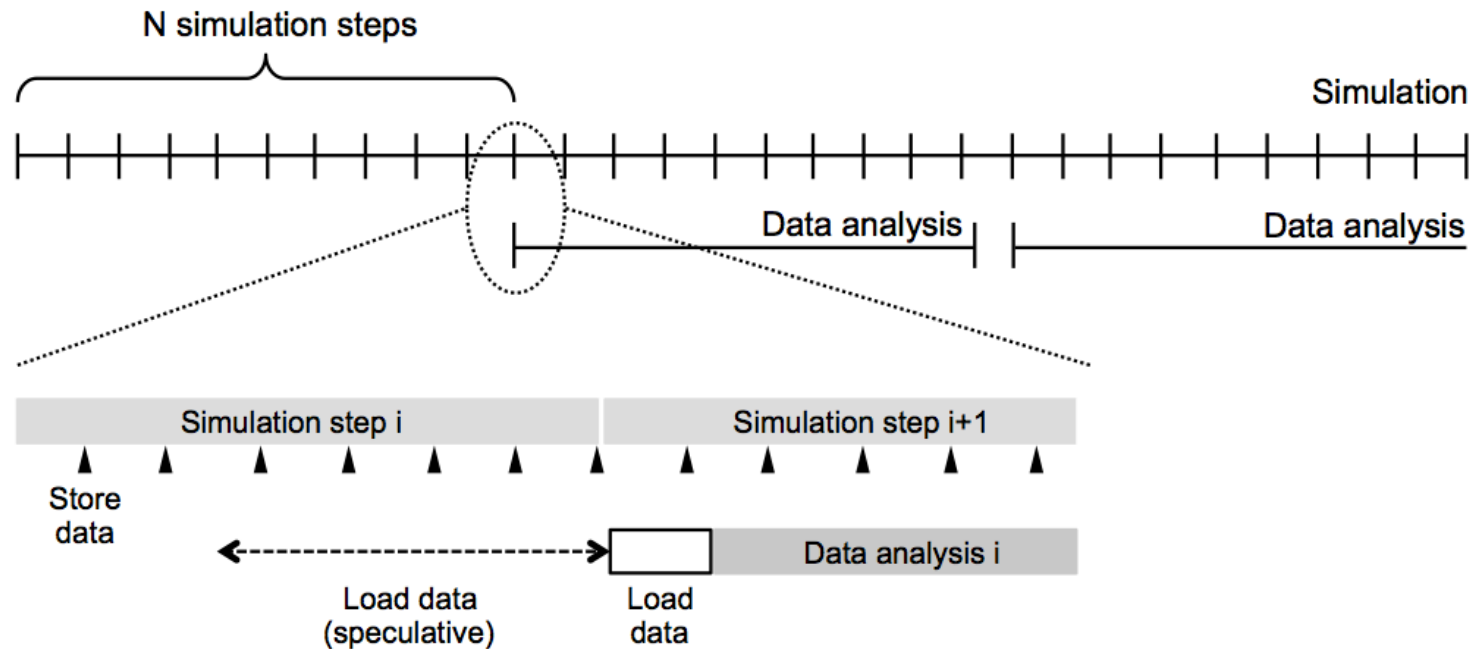
- Small DRAM capacity per core – even aggregated memory on dedicated nodes can hardly keep all coupled data (given the ratio of resource allocations for compute nodes and dedicated nodes)

Hybrid Staging

- Spans horizontally across compute nodes
- Spans vertically across the multi-level memory hierarchy, e.g. DRAM/NVRAM/SSD, to extend the capacity of in-memory data staging



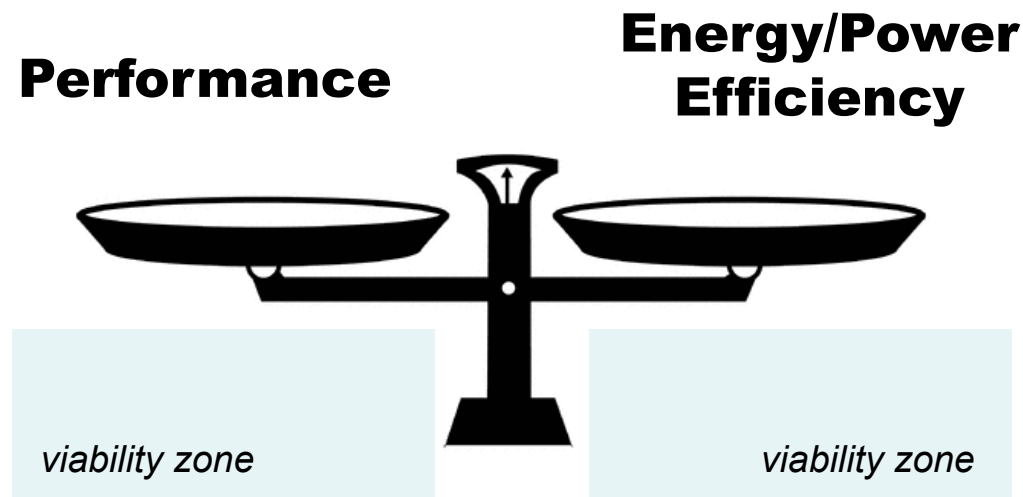
2 Synthetic Workflow for Understand Relative Behaviors



- Data analysis pipeline (e.g., S3D+Topology analysis)
- Synthetic kernels to evaluation relative behaviors
- Potential use of **speculative** data movement
 - Out-of-the-core data movement vs. traditional speculation at CPU level

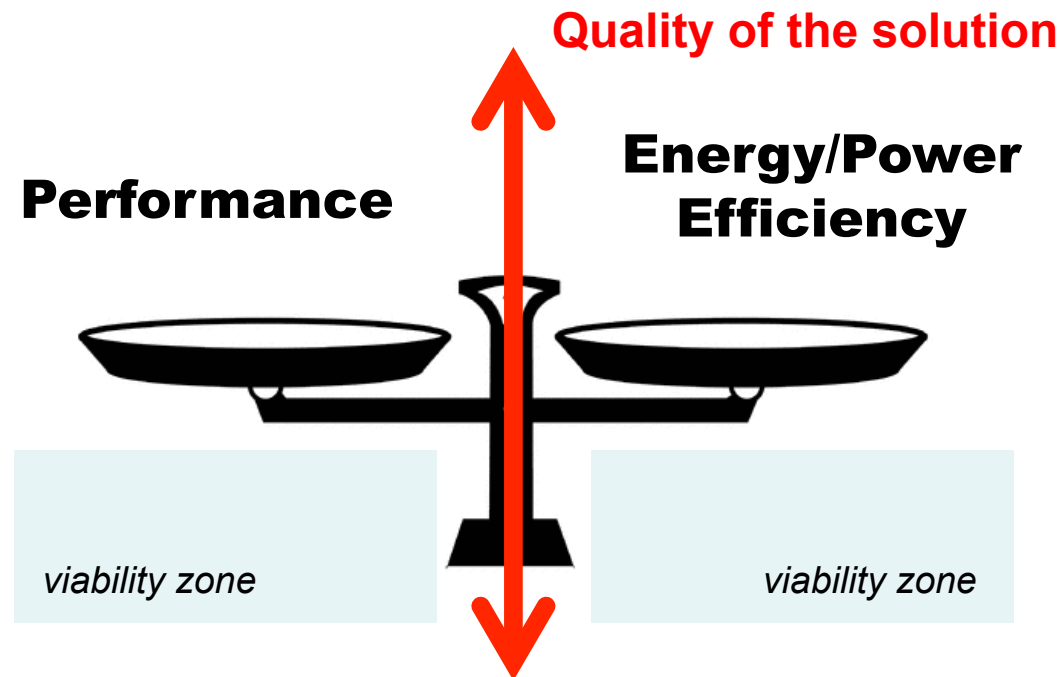
Understanding Behaviors and Tradeoffs

- Performance and Energy/Power
- Limitations when “viability zones” are exclusive



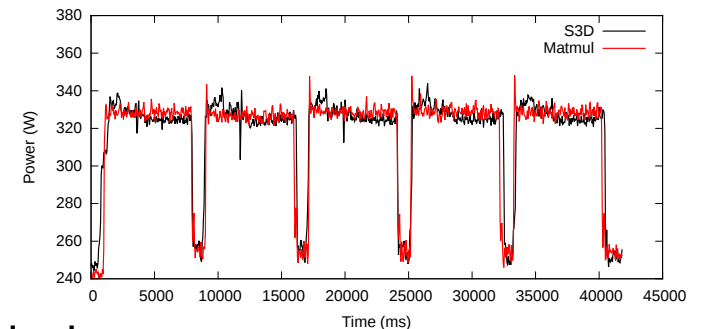
Understanding Behaviors and Tradeoffs

- Quality/accuracy of the solution
 - E.g., single/double precision, convergence values, AMR codes, etc.
- **Frequency of analysis** to represent quality of the solution



Evaluation Methodology

- Evaluation framework (single node)
 - Customizable multi-threaded framework which takes care of the workflow synchronization
 - Can run different kernel/applications as the simulation and analysis are customizable
- Synthetic kernels to evaluation relative behaviors
 - Matrix multiplication in simulation steps
 - Word finding in analysis steps
- Multiple customizable input parameters
 - Number of cores assigned for simulation/analysis
 - Data path (HDD, SSD, etc.)
 - Frequency of analysis, i.e. every x number of steps



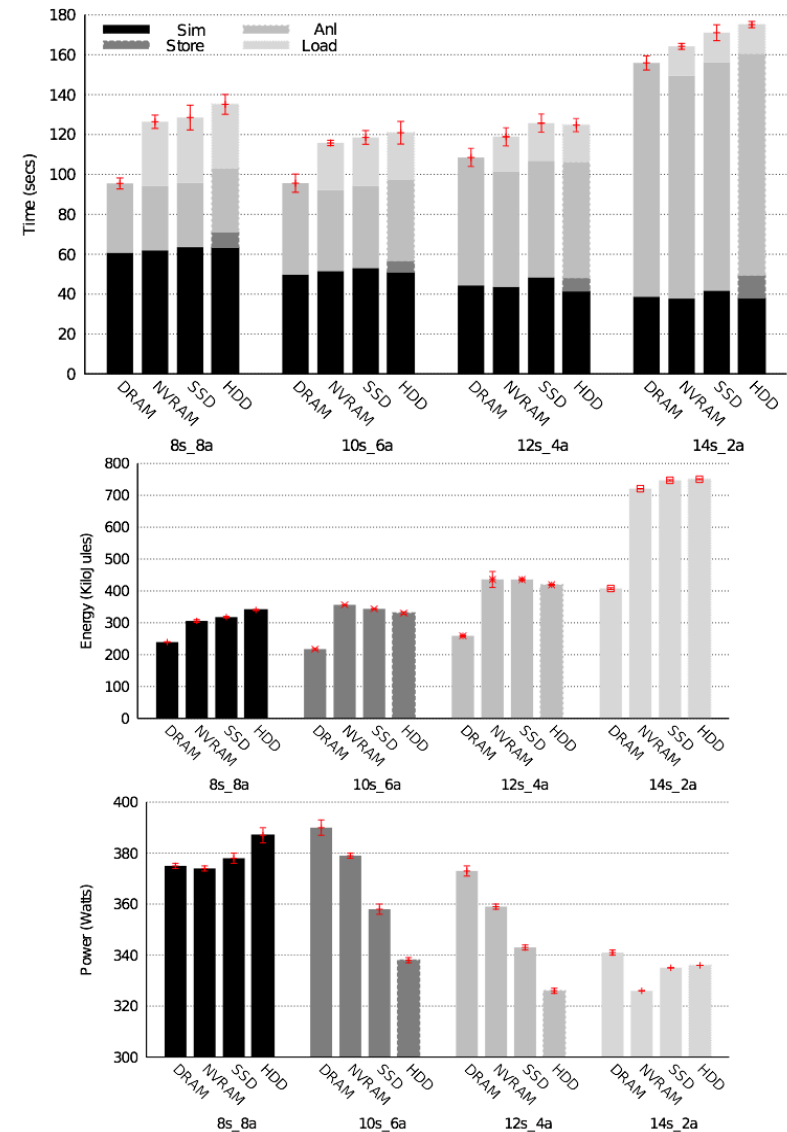
Computational And data-enabled Platform for Energy efficiency Research (CAPER)



- NSF funded research instrument
- SyperMicro SYS-4027GR-TRT (support up to 8 GPUs concurrently)
- Phase 1
 - 8 servers with 2 Intel Xeon Ivybridge E5-2650v2 (16 cores)
 - 128GB of DRAM
 - 1TB of PCIe Flash-based NVRAM (Fusion-io iodrive2)
 - 2TB of SSD (RAID)
 - 4TB of hard disk (RAID)
 - Intel Xeon Phi 7120P
 - Infiniband FDR and 10G Ethernet
- Phase 2
 - NVIDIA K40
- Instrumentation
 - Coarse grain: PDU (1Hz)
 - **Fine grain: Yokogawa DL850E ScopeCorder (1KHz) – from modules at 10Ms/s**

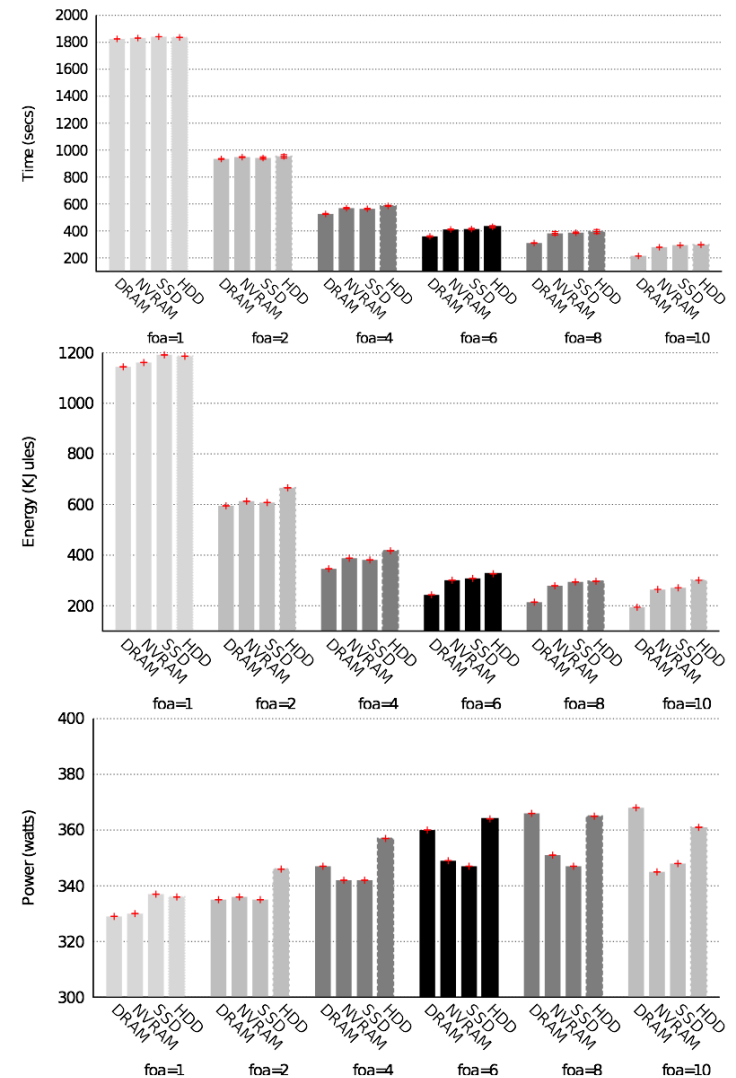
1 Data Staging over (Deep) Memory Hierarchy

- Empirical evaluation on CAPER
- Execution time, energy consumption, and average power of the workflow's execution using different configurations and devices for data staging
- Each group of columns represents one configuration (number of cores for simulation/analysis)
- **Goal: finding sweet spots for in-situ data analysis**
 - In this example 12 core for simulation and 4 for analysis
 - Using very few cores for analysis delays simulation tasks
 - Higher power with power demanding devices (e.g., DRAM)



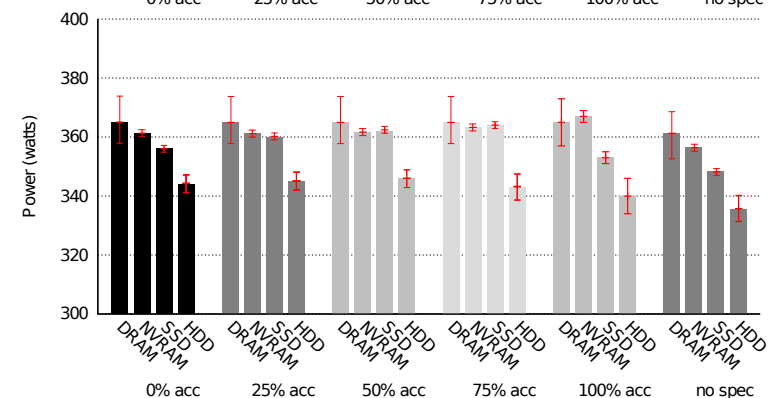
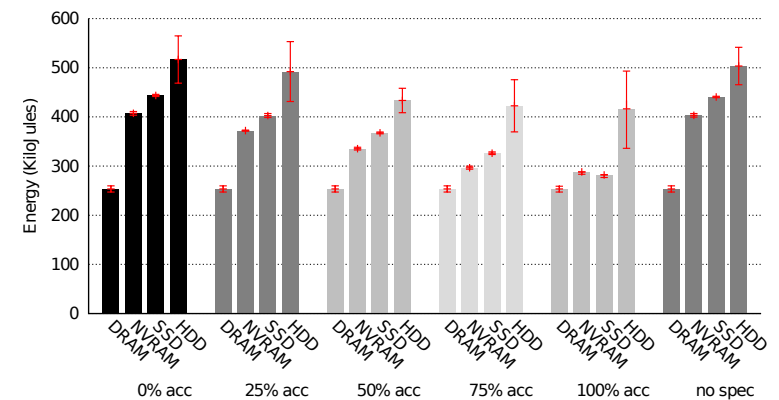
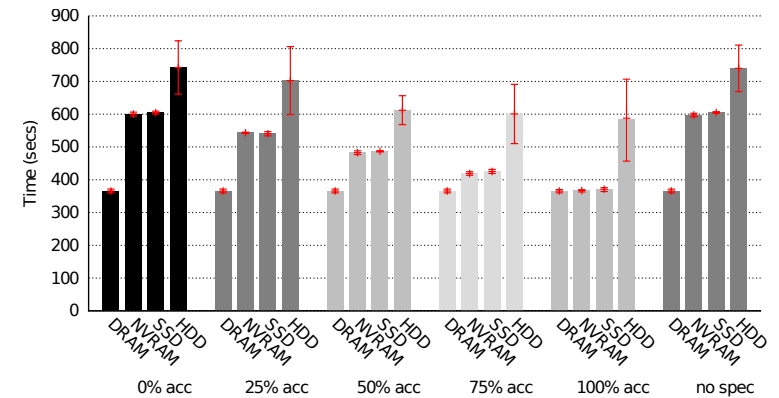
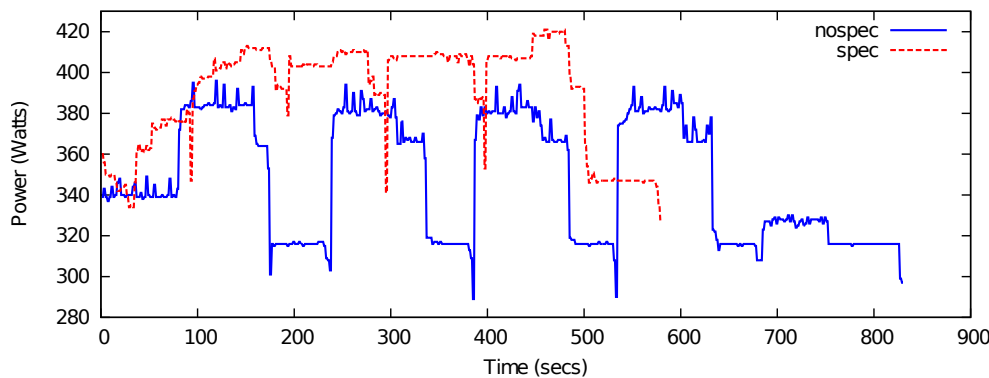
Tradeoffs with the Quality of Solution (Freq. of analysis)

- Execution time, energy consumption, and average power of the workflow's execution for different frequency of analysis ("foa") and different devices for data staging
- Frequency of analysis foa = k means that the data analysis is performed every k simulation steps
- **Execution time and energy consumption decreases as the frequency of analysis decreases**
- **However, average power increases as more computation/data movement is performed in parallel**



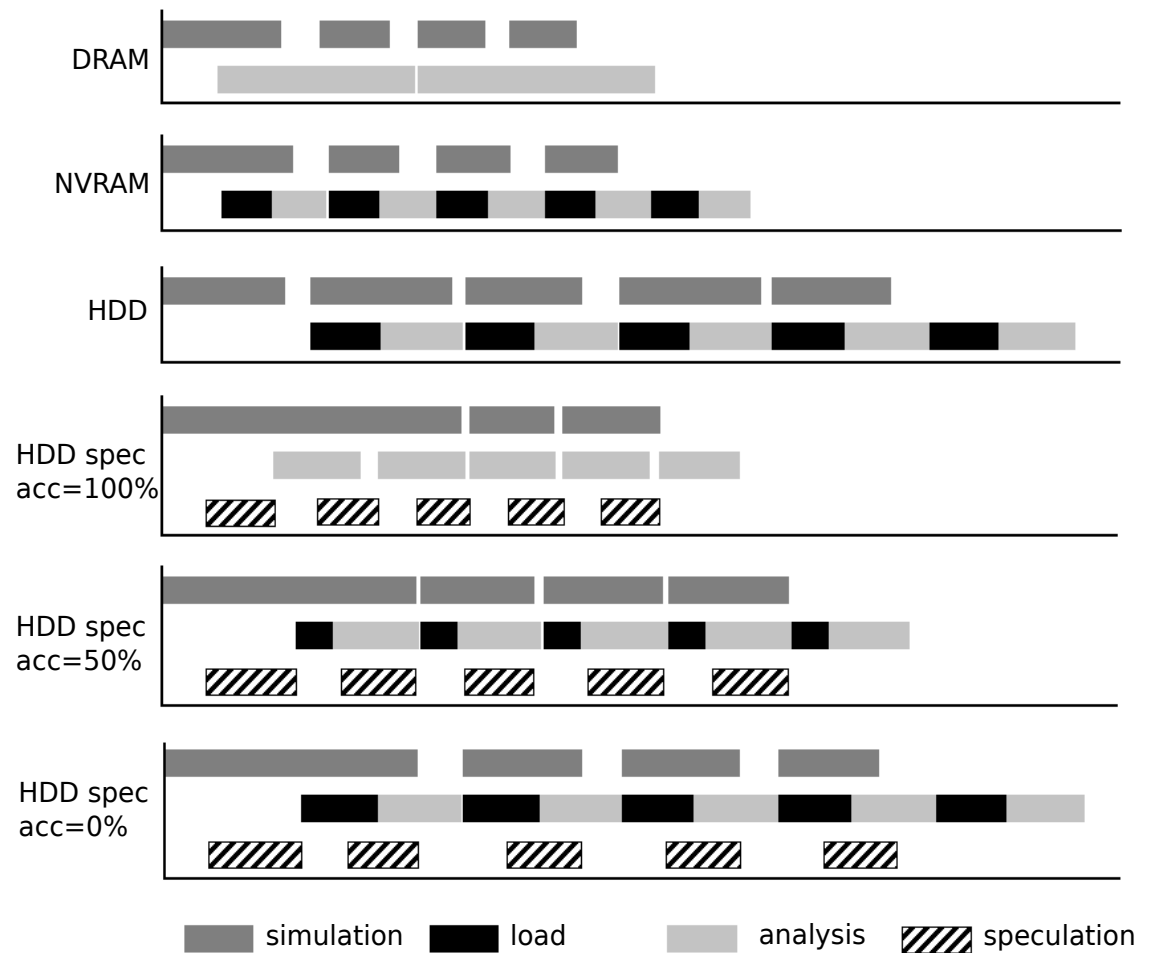
2 Speculative Data Movement

- Data speculation incurs little overhead when it is 0% accurate vs. no speculation both in terms of time and energy consumption
- The average power is higher when performing data speculation because it shares resources with the simulation and the analysis



Speculative Data Movement (cont.)

- Results present tradeoffs that can be exploited at runtime
 - E.g., Execution behavior with NVRAM is similar with HDD when data speculation is accurate
 - However, overall energy consumption is a bit higher with NVRAM (device power requirements)



Conclusions

- Costs (energy, latency) related to transporting, processing and analyzing increasing **data volumes and rates** are limiting the insights from extreme scale applications
- Energy/power-efficiency in combination with other objectives – understanding tradeoffs are important
 - **Quality of solution**, Performance, Resiliency, etc.
- Using **data speculation** in data-intensive workflows can positively impact energy consumption without much negative impact on performance or the quality of the solution
- Co-design is essential along multiple dimensions
 - E.g., **runtime system** to balance these tradeoffs

Thank You!

Ivan Rodero, Ph.D.

Rutgers Discovery Informatics Institute
NSF Cloud and Autonomic Computing Center
Rutgers, The State University of New Jersey

Email: irodero@rutgers.edu

WWW: <http://nsfcac.rutgers.edu/people/irodero/>