



# A project in the life of a data scientist

Janine C. Bennett

Sandia National Laboratories

July 30, 2015

ICERM Workshop on Mathematics in Data  
Science

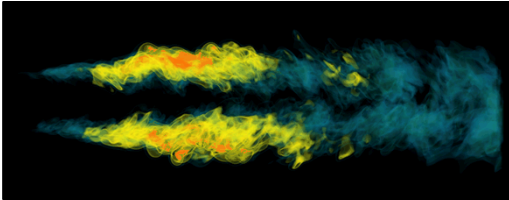


*Exceptional  
service  
in the  
national  
interest*

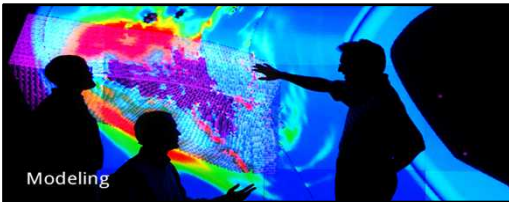


Sandia National Laboratories is a multi-program laboratory managed and operated by Sandia Corporation, a wholly owned subsidiary of Lockheed Martin Corporation, for the U.S. Department of Energy's National Nuclear Security Administration under contract DE-AC04-94AL85000. SAND NO. 2011-XXXXP

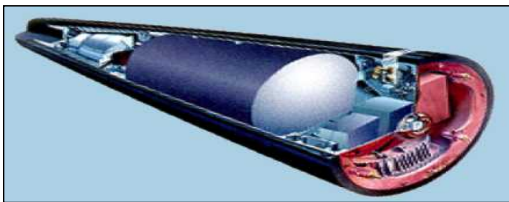
# Sandia performs scientific research in support of national policy and decision making



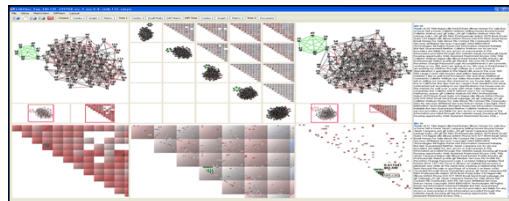
**Energy:** Reduce U.S. reliance on foreign energy, reduce carbon footprint



**Climate change:** Understand, mitigate, and adapt to the effects of global warming



**National Nuclear Security:** Maintain a safe, secure, and reliable nuclear stockpile



**Cyber:** Shore up our nation's cyber defenses, provide more fundamental understanding of cyber environment

# Interdisciplinary teams use extreme-scale experiments, modeling, and simulation to reason about complex phenomena

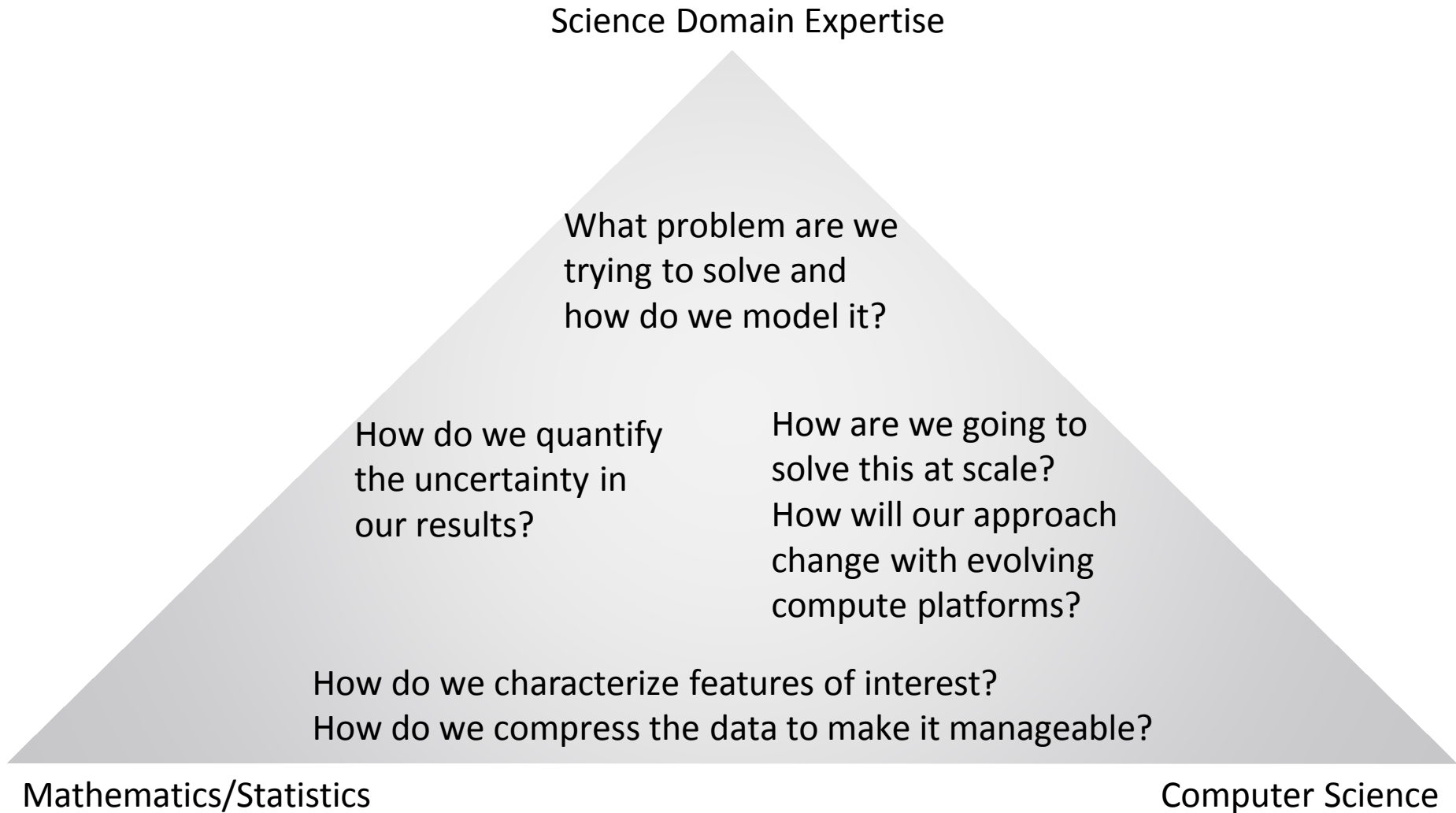
Science Domain Expertise



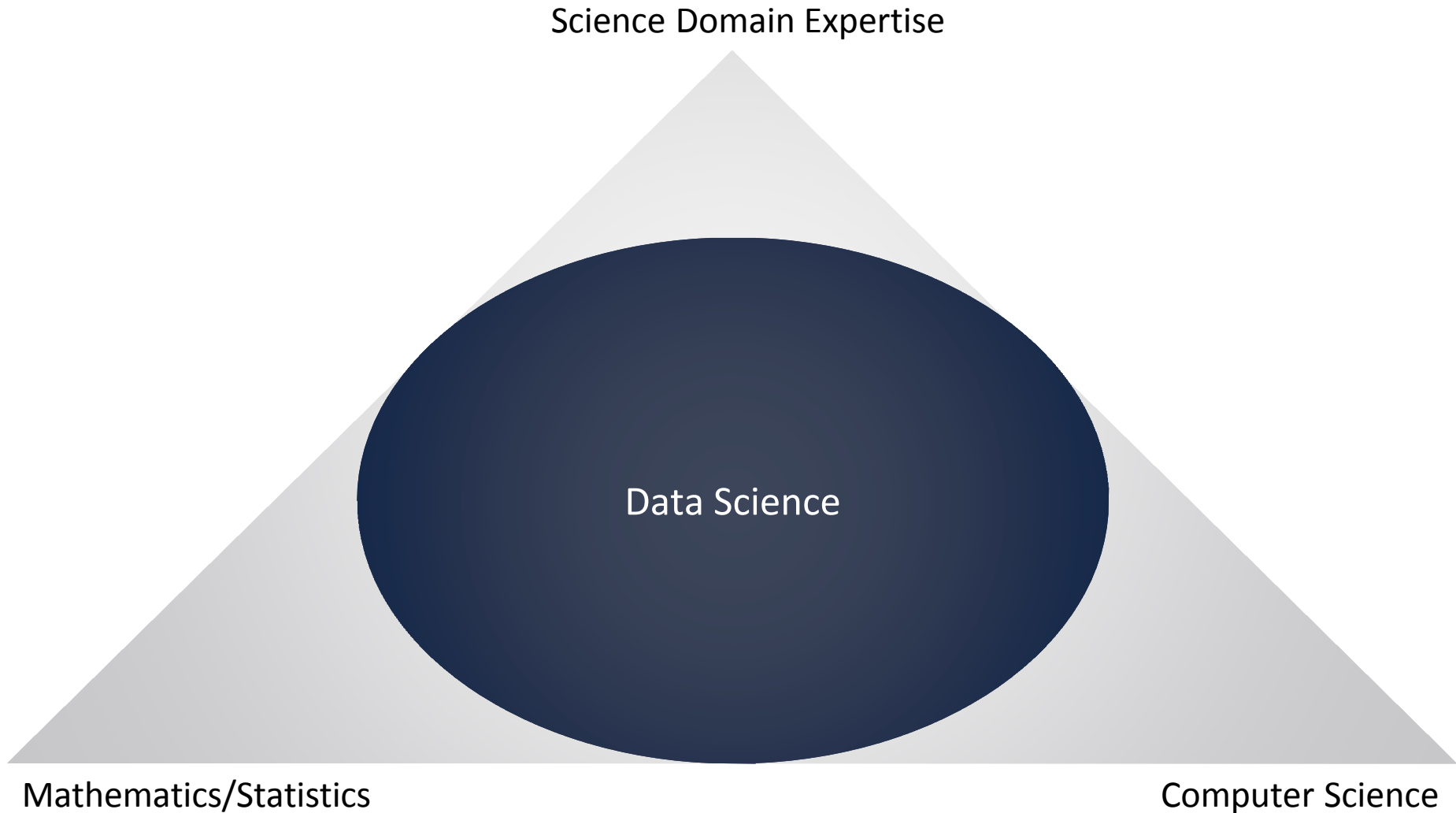
Mathematics/Statistics

Computer Science

# Interdisciplinary teams use extreme-scale experiments, modeling, and simulation to reason about complex phenomena



# Data science spans the different areas of expertise on these interdisciplinary projects



# A number of lessons learned can be gleaned from a retrospective look at a data science project

- Project details
  - Three year project
  - Relatively small team
  - Focused on enabling scientific simulations
  
- General lessons learned and math takeaways (MT)
  - Building the team
  - Scoping the work
  - Doing the work

## Acquiring funding: Data science calls for proposals often highlight the interdisciplinary nature of the work

“advance the underlying **math** and **computer science** to enable the routine use of **rigorous predictive simulation** ... more effectively use **next generation computers**”

# Acquiring funding: Data science calls for proposals often highlight the interdisciplinary nature of the work

Science Domain Expertise

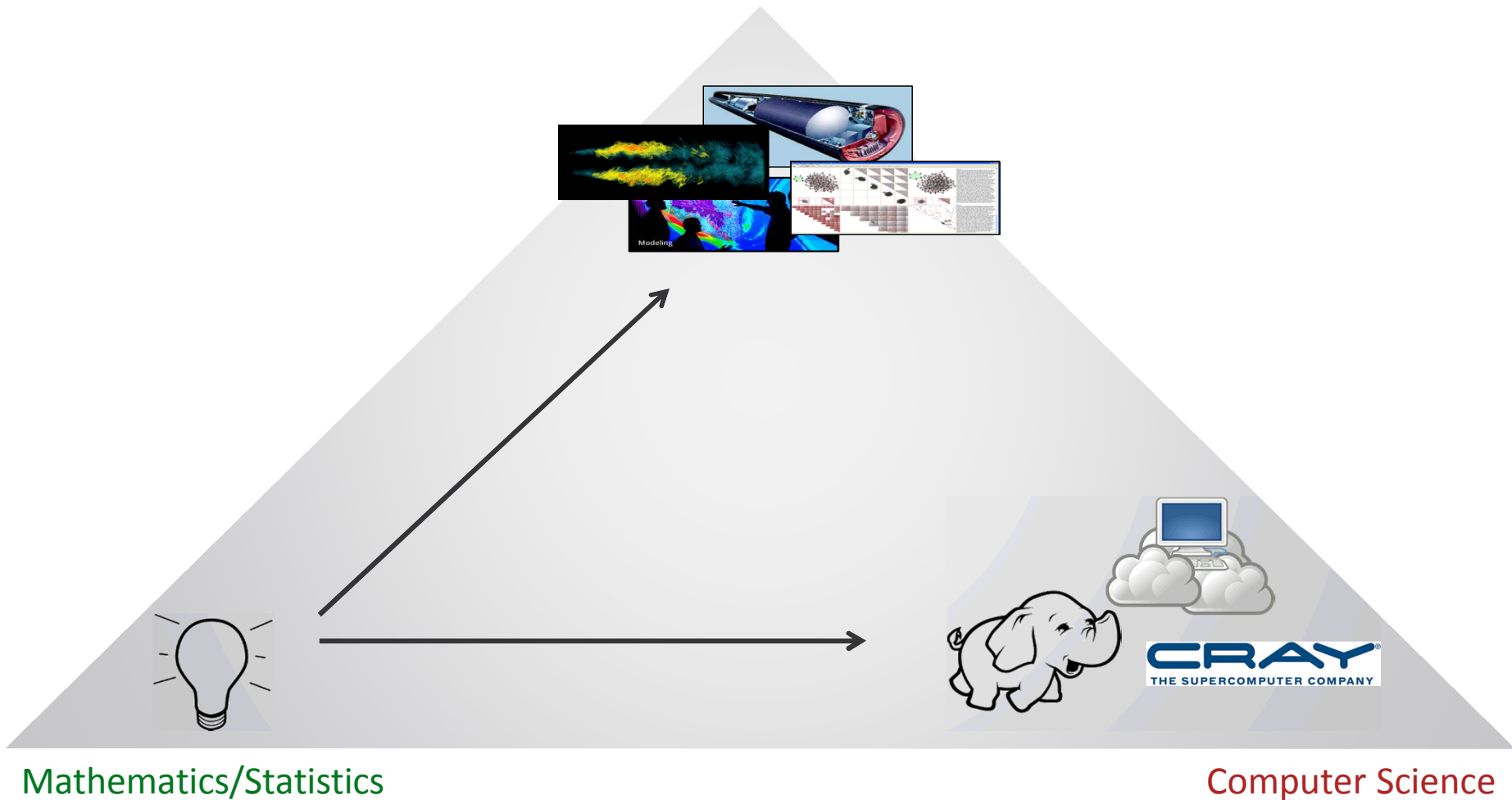
“advance the underlying **math** and **computer science** to enable the routine use of **rigorous predictive simulation** ... more effectively use **next generation computers**”

Mathematics/Statistics

Computer Science

# Successful data science projects often develop new theory and can demonstrate its applicability to mission *at scale*

Science Domain Expertise

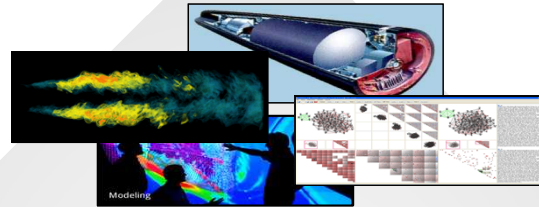


# MT0: Target a customer and compute platform for your research

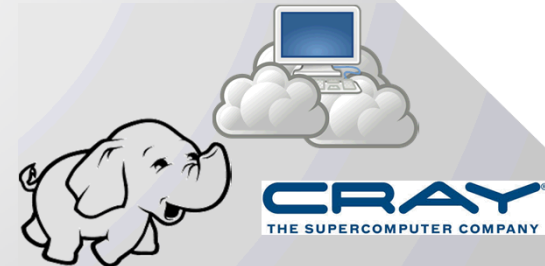
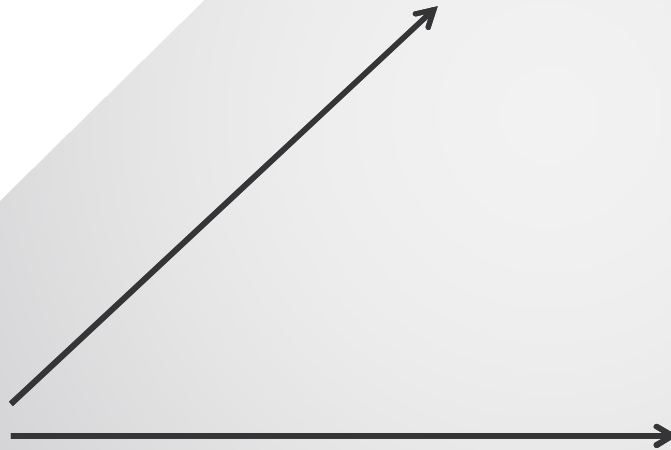
## Science Domain Expertise

Are they doing

- Exploratory science?
- Answering a yes/no question?
- Design optimization?



- How do I express my algorithm?
- How does the architecture affect an idealized mathematical model?

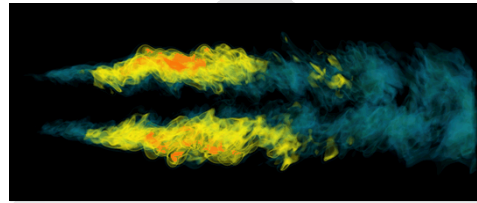


Mathematics/Statistics

Computer Science

# Proposed work: Sublinear algorithms for in-situ and in-transit data analysis at extreme scale

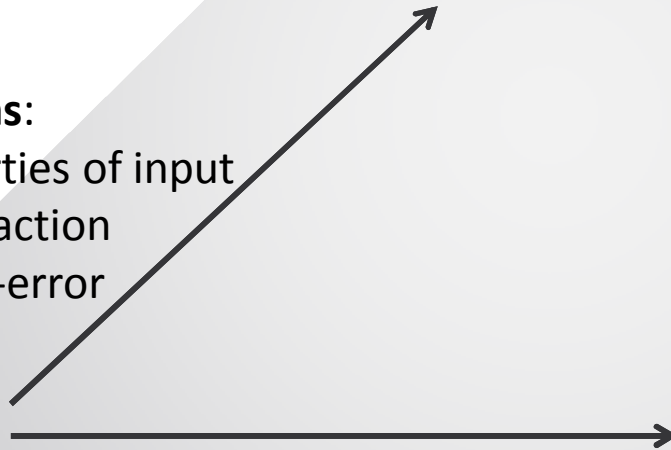
Science Domain Expertise



**Target customers:** General exploratory science (combustion use case)

## Sublinear algorithms:

- Determine properties of input with only a tiny fraction
- Quantifiable time-error tradeoffs



**Target compute platform:** Leadership-class high performance computing (HPC)



Mathematics/Statistics

Computer Science

# Our target customers perform simulations that generate large complex data sets using HPC platforms

- Case study: Direct Numerical Simulations & turbulent combustion
- Data size
  - $O(\text{Billions})$  of grid points per time step
  - $O(100\text{K})$  time steps
- Data complexity
  - Multivariate
    - $O(100)$  chemical species
    - Vector data
    - Particle data
  - Turbulence is a complex phenomenon
  - Length scales: microns to centimeters
  - Temporal scales: nanoseconds to milliseconds

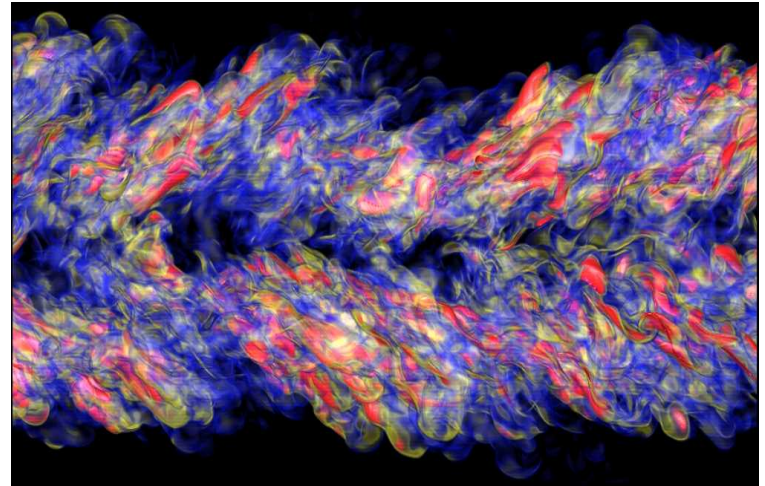
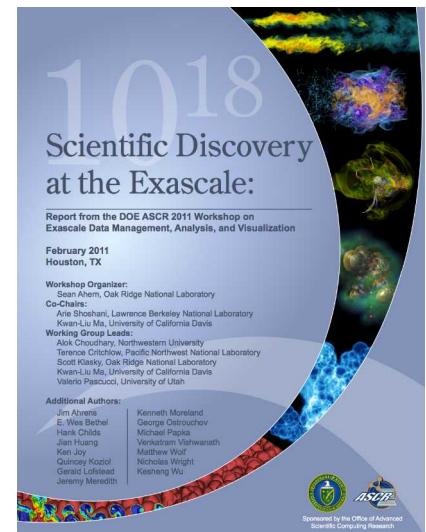
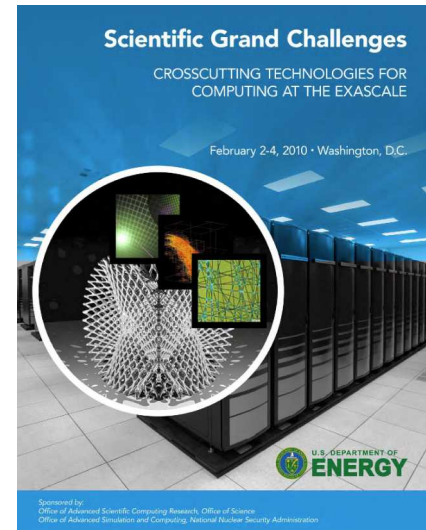


Image courtesy of Jacqueline Chen

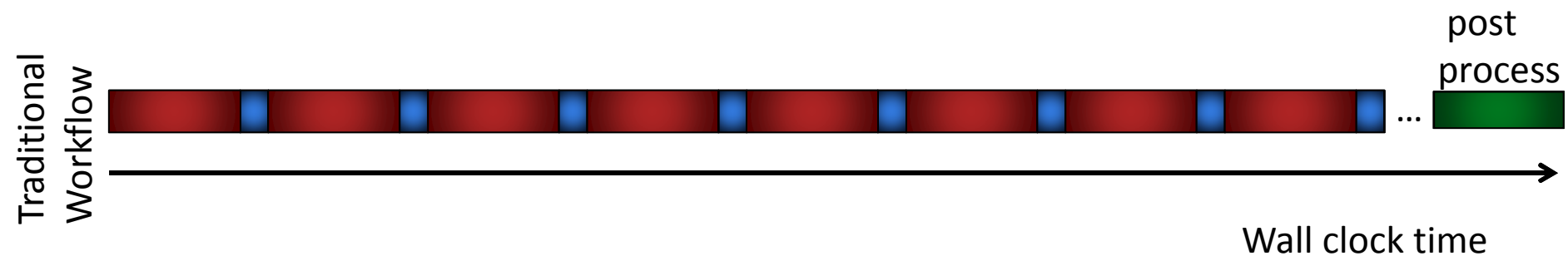
# There is a widening gap between **I/O** and **compute** capabilities on our target compute platforms

System Parameter	2011	2018		Factor Change
System Peak	2 Pf/s	1 Ef/s		500
Power	6 MW	≤20 MW		3
System Memory	0.3 PB	32-64 PB		100-200
Total Concurrency	225K	1 BX10	1B X100	40000-400000
Node Performance	125 GF	1 TF	10 TF	8-80
Node Concurrency	12	1000	10000	83-830
Network Bandwidth	1.5 GB/s	100 GB/s	1000 GB/s	66-660
System Size (nodes)	18700	1000000	100000	50-500
I/O Capacity	15 PB	30-100 PB		20-67
I/O Bandwidth	0.2 TB/s	20-60 TB/s		10-30



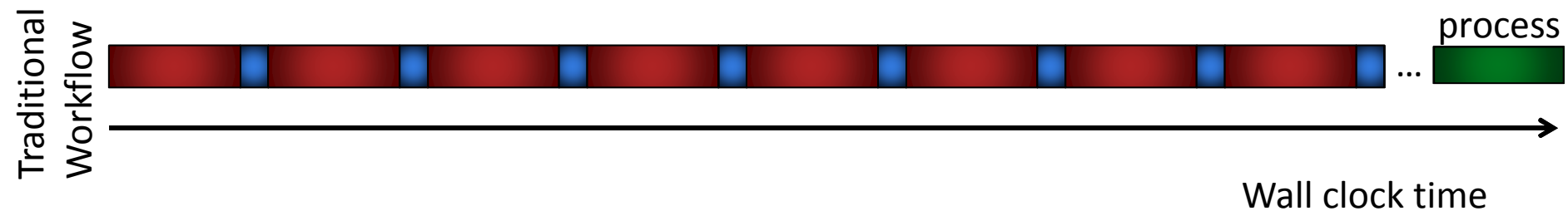
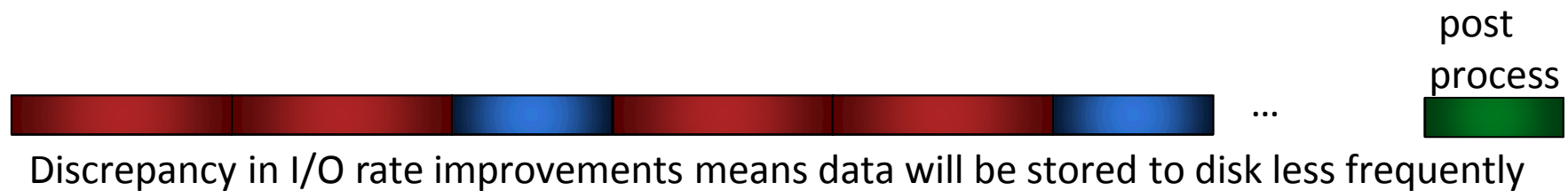
# The widening gap in compute and I/O is causing changes in the scientific workflows

 Simulation     Check-pointing     Analysis



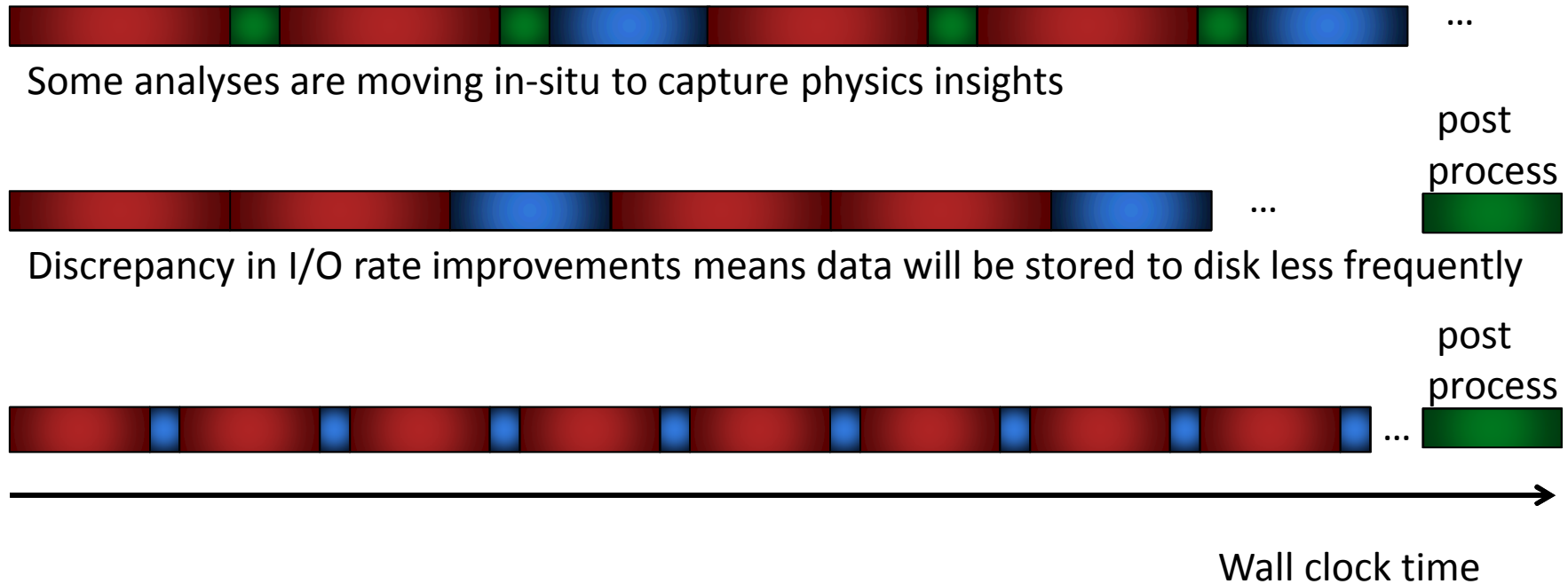
# The widening gap in compute and I/O is causing changes in the scientific workflows

 Simulation     Check-pointing     Analysis



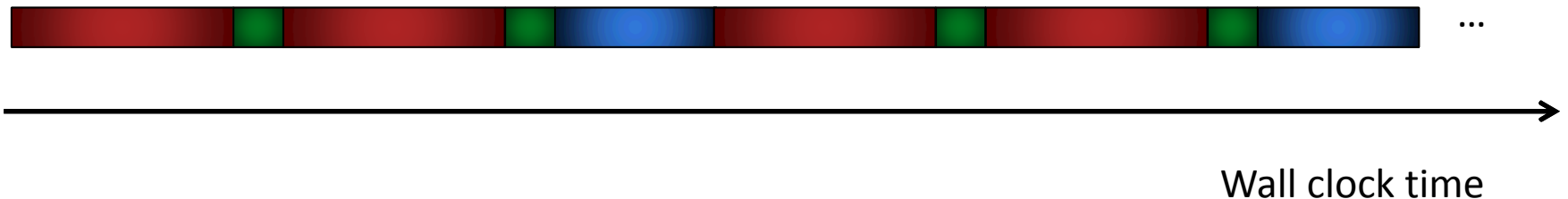
# The widening gap in compute and I/O is causing changes in the scientific workflows

 Simulation     Check-pointing     Analysis



# The change in scientific workflows introduces a number of data science challenges

 Simulation     Check-pointing     Analysis



- At what frequency should I/O or analysis be done?
- Can we make this decision in an adaptive, data-driven fashion at runtime?
  - Avoid missing interesting science
  - Avoid costly I/O when simulation state is evolving slowly
- How can we make these decisions quickly and efficiently?
- How do we design efficient analysis algorithms given in situ constraints?

# Our project aimed to apply sublinear analysis to address in situ workflow challenges

Sublinear analysis is a relatively new theoretical subfield asking:  
how to determine properties of input by seeing tiny fraction

## sublinear algorithms

- Small samples of data
- Quantifiable time-error tradeoffs
- Limited primitives for access

## in situ analysis challenges

- Too much data to move
- Constrained time budgets
- Simulation dictates data structures

There is strong alignment between theory and challenges

# Building a team with the right mix of domain, computer science and mathematics expertise is critically important

Science Domain Expertise



Mathematics/Statistics

Computer Science

# Effective communication *between* team members is as important as their individual expertise

Science Domain Expertise



Mathematics/Statistics

Computer Science

# MT1: As the breadth of project scope increases so does the social complexity – effective communication is key

Science Domain Expertise

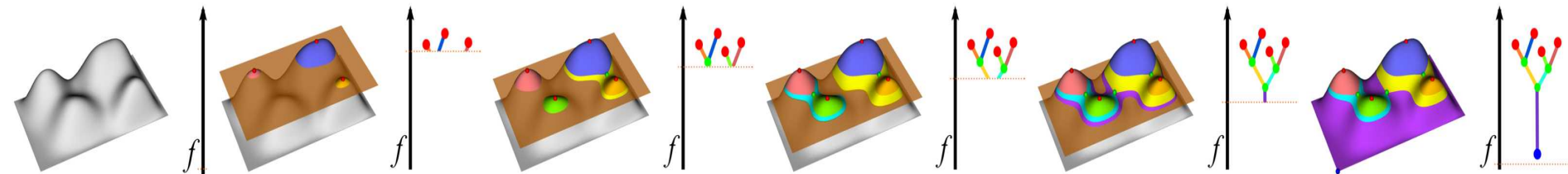


Mathematics/Statistics

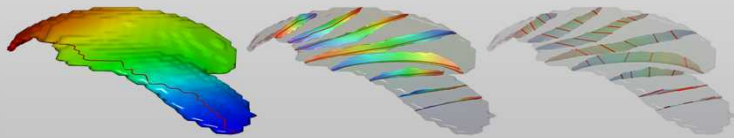
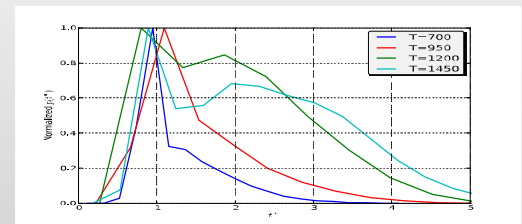
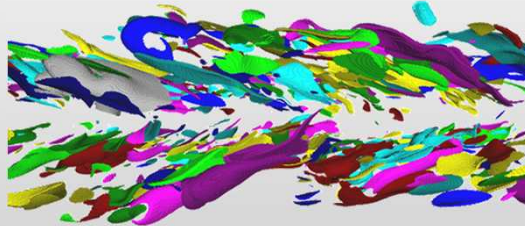
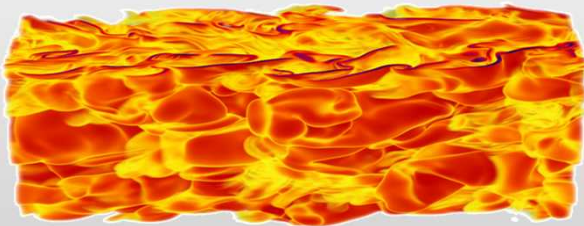
Computer Science

# Our initial plan aimed to make in-situ analysis algorithms more efficient using sampling

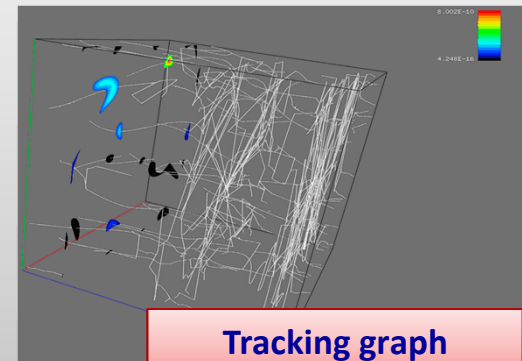
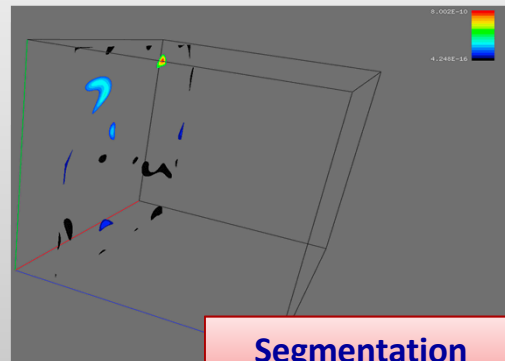
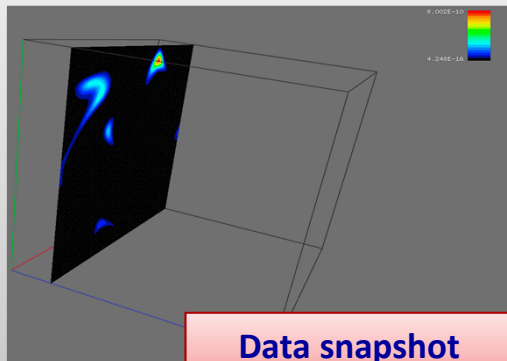
- Merge tree algorithm for encoding level-set behavior of a function defined on a mesh
- Once computed, provides a compact representation of domain segmentation
- Enables efficient queries of feature-based quantities of interest



# Merge trees enable a variety of feature-based exploratory techniques as a post-process



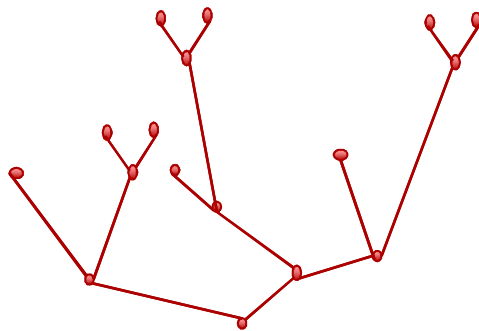
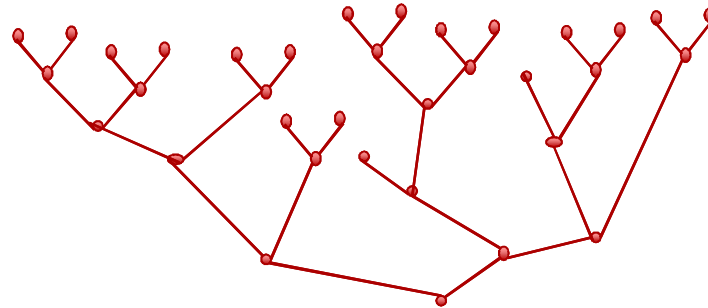
Identify features, characterize their shapes and analyze the behavior of other variables within these features



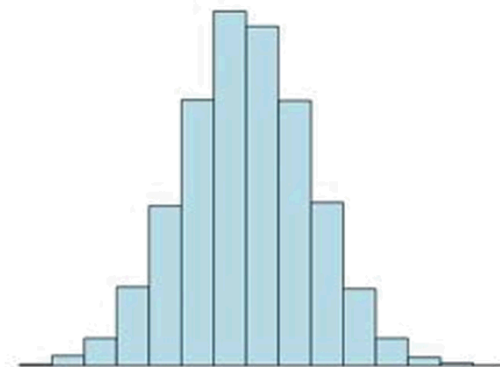
Tracking features in space and time

# Our strategy was to use sampling on each processing element and merge results into a final tree

Approximate full merge tree



Merge tree for global  
connectivity



Sampled histogram for local  
on-node data

## Our finding: a sampling-based approach doesn't provide a computational win

- Mesh and associated data are distributed across processing elements (PE)
- Data decomposition is optimized for the simulation
- Communication costs end up dominating run time
- You can compute the full tree on a PE in approximately the same amount of time as the sampled histogram!
- Better to focus on a fully distributed implementation, rather than use sampling

## Our finding: a sampling-based approach doesn't provide a computational win

- Mesh and associated data are distributed across processing elements (PE)
- Data decomposition is optimized for the simulation
- Communication costs end up dominating run time
- You can compute the full tree on a PE in approximately the same amount of time as the sampled histogram!
- Better to focus on a fully distributed implementation, rather than use sampling

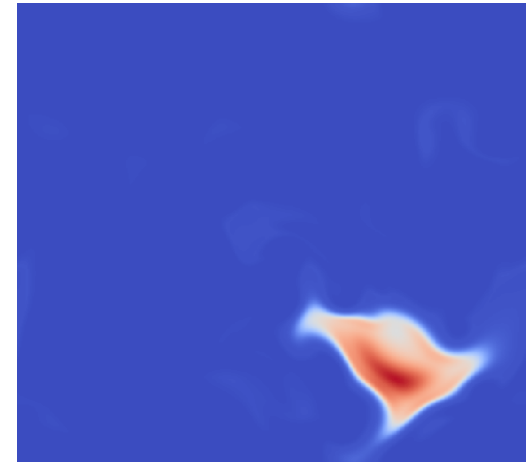
**MT2: Algorithms can look beautiful in theory, but may not be worth deploying in practice due to constraints of your target user and architecture**

## MT3: Question your underlying algorithmic assumptions to mitigate technical constraints

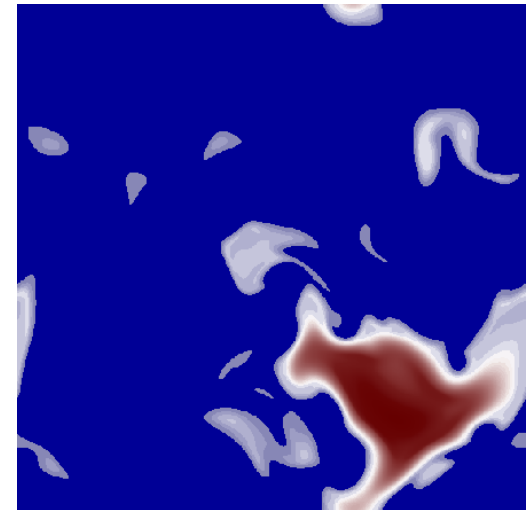
- What is performant in serial may not work well in a distributed fashion
  - All to all communication is expensive!
  - In-situ data layouts will likely not be optimal for your algorithms
- Are you exposing the maximum amount of parallelism?
  - Data parallelism: same or different tasks operating on different data
  - Task parallelism: different tasks operating on the same data
  - Pipeline parallelism: Overlapping communication and computation
    - How asynchronous is your approach?

# We had better success applying the sampling based technique to visualization

- Mapping function values to colors is key to visualization
- Many tools map linearly by default
  - User interaction to refine the map
- Developed fast, efficient algorithm to build color-maps using cumulative density function-based sketches
- Our distribution-based sketch enables better feature identification
- Available open source as a ParaView plugin
- Published in IEEE Symposium on Large-Scale Data Analysis & Visualization 2013



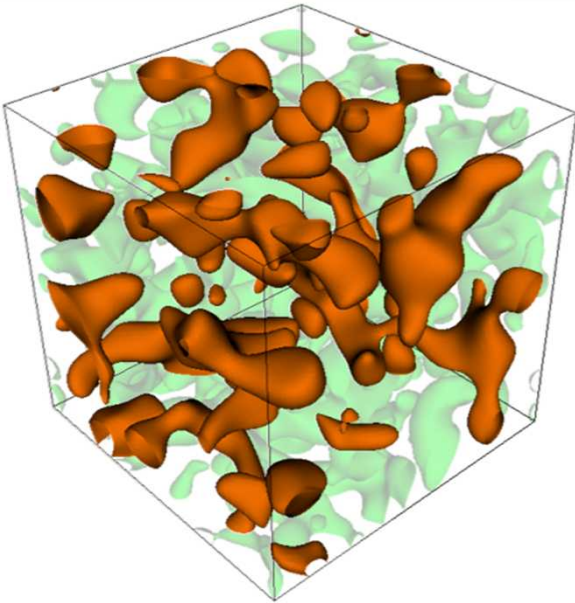
Linear color map



CDF-based color map

# Our most impactful results stemmed from prolonged discussions with domain scientists regarding workflows

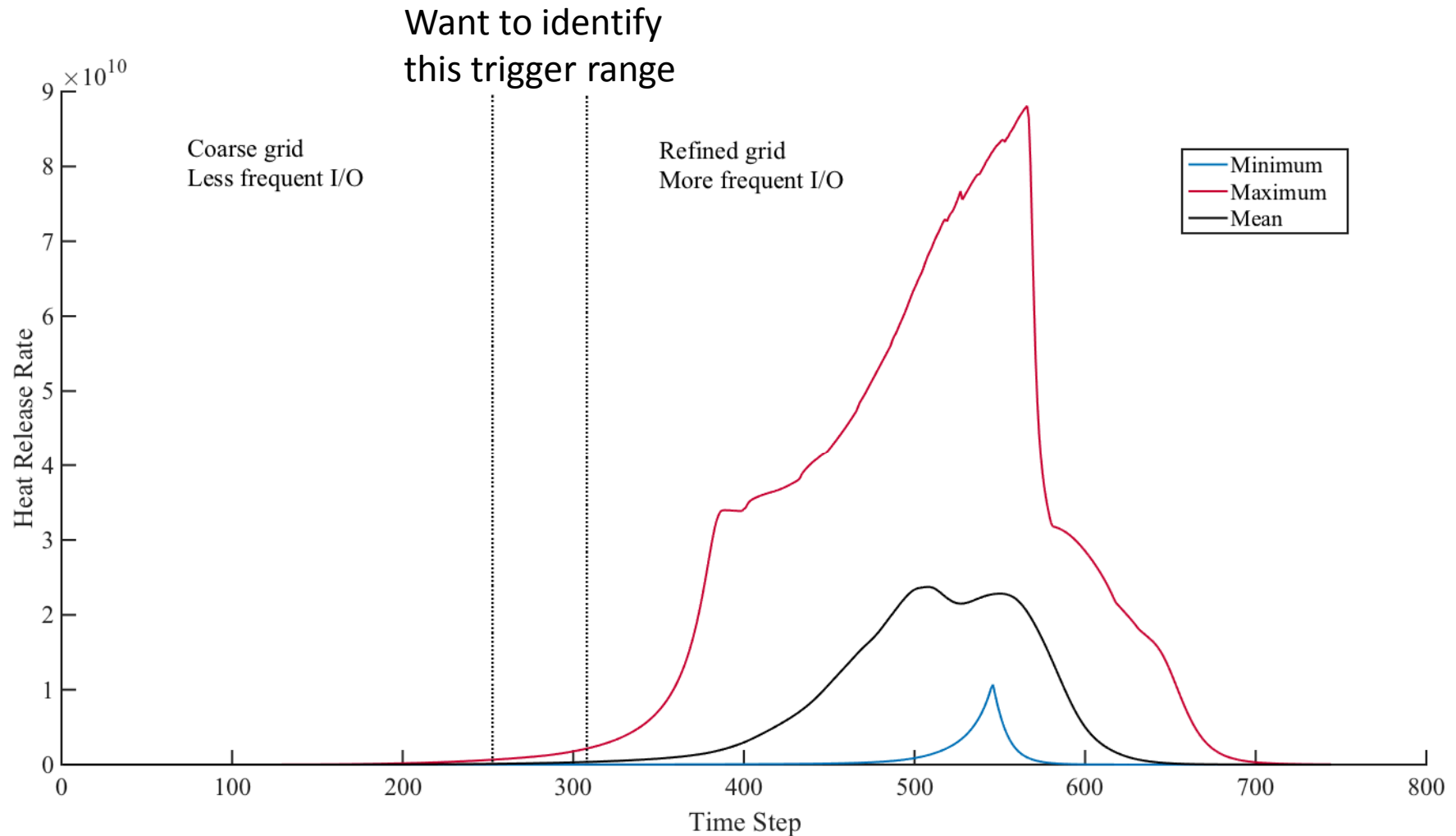
- At what frequency should I/O or analysis be done?
- Can we make this decision in an adaptive, data-driven fashion at runtime?
  - Avoid missing interesting science
  - Avoid costly I/O when simulation state is evolving slowly
- How can we make these decisions quickly and efficiently?



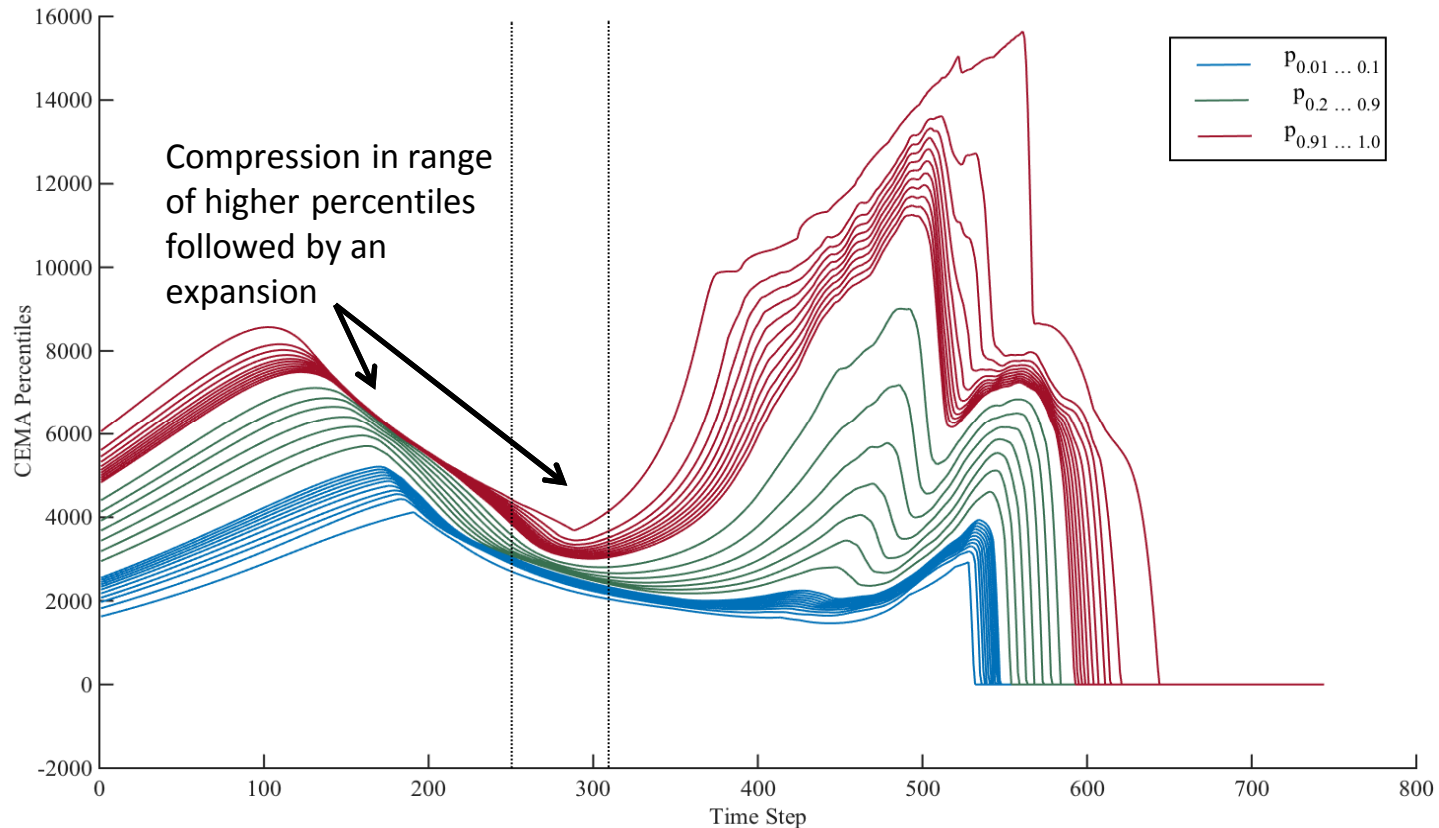
HCCI: Homogenous Compression Charge Ignition

Lots of little heat kernels slowly develop until ignition occurs

# Their goal was to optimize I/O frequency and mesh resolution in a data driven manner based on heat release



# Domain expert hypothesis: Chemical Explosive Mode Analysis (CEMA) is a good predictor but too expensive

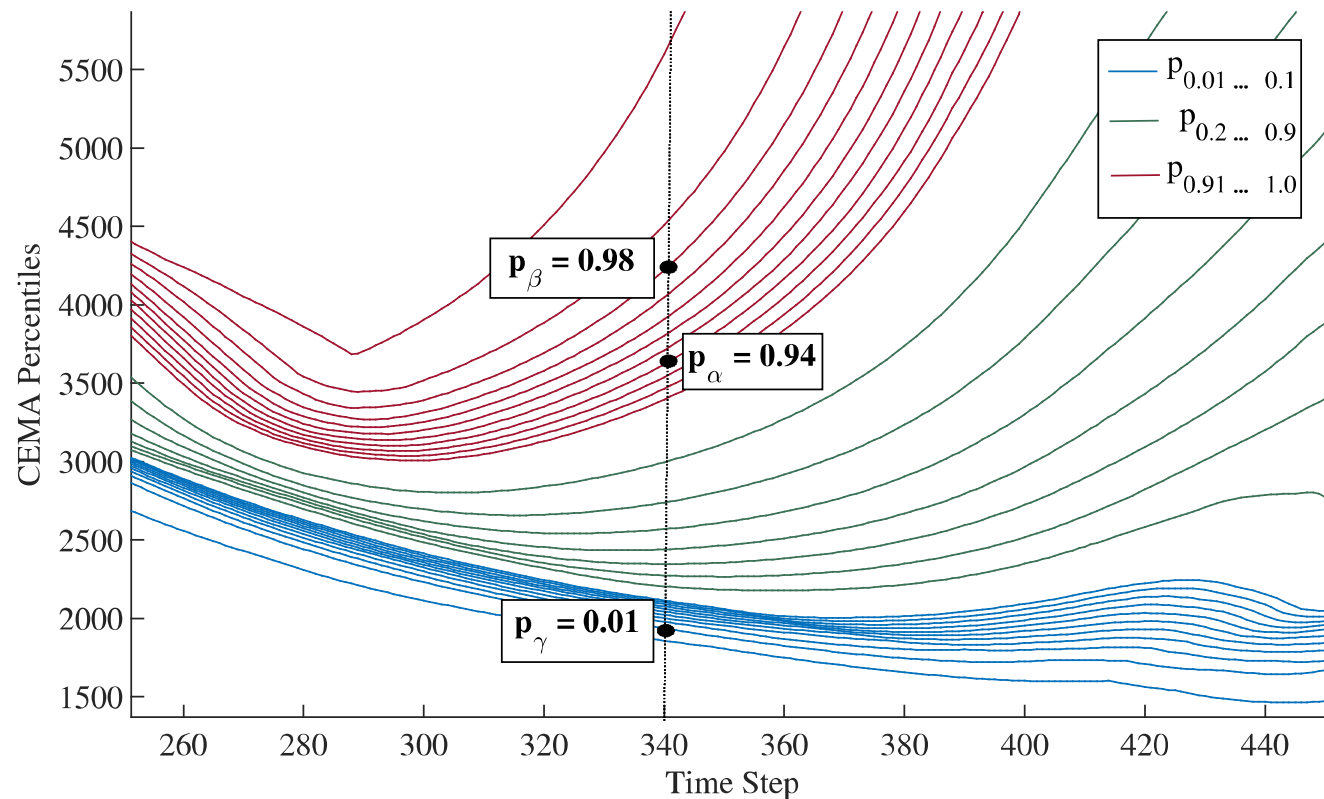


- Point-wise Jacobian of chemical species
- Cost is prohibitively expensive – up to 60 times the cost of a simulation timestep

# We defined a simple noise-resistant indicator function to characterize the spread of CEMA percentiles

$$P_{\alpha,\beta,\gamma}(t) = \frac{p_{\alpha}(t) - p_{\gamma}(t)}{p_{\beta}(t) - p_{\gamma}(t)}$$

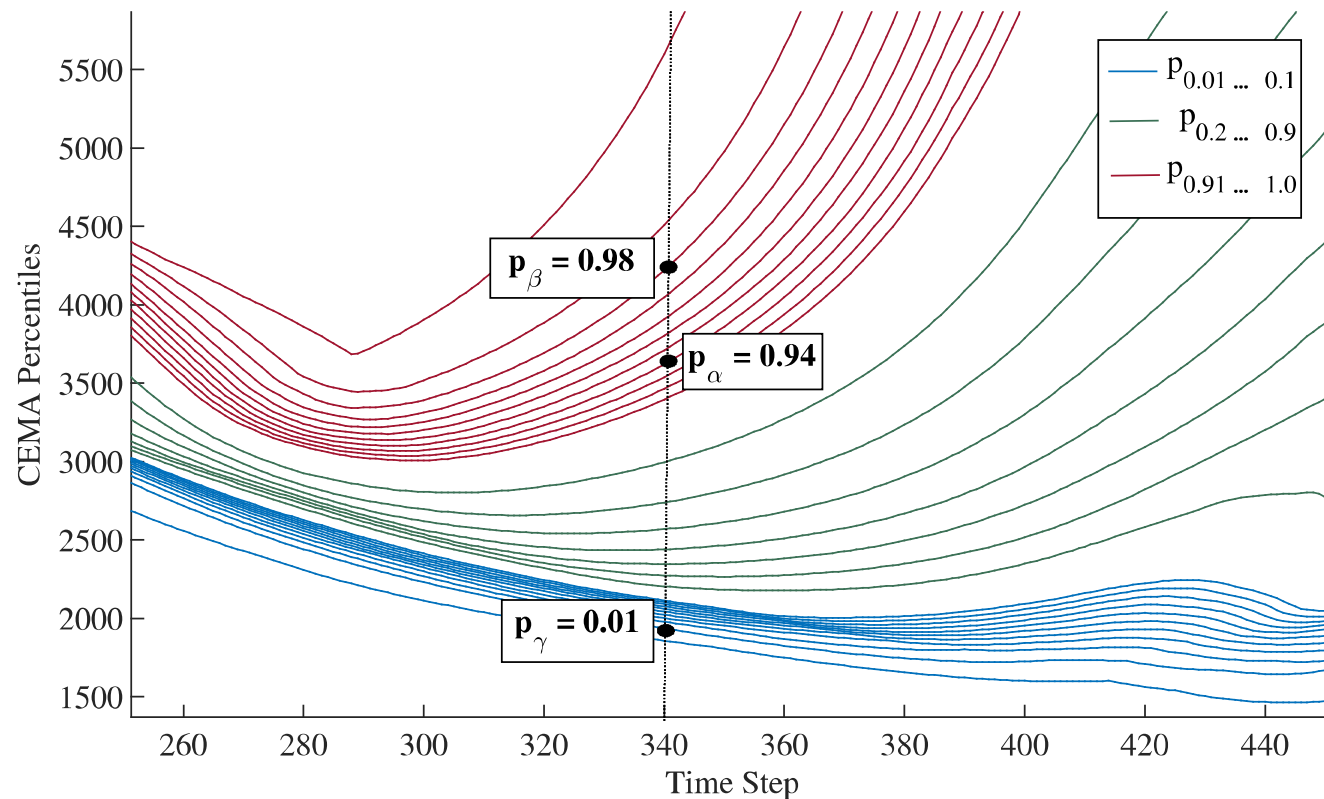
Avoid use of minimum and maximum percentiles as these are outliers



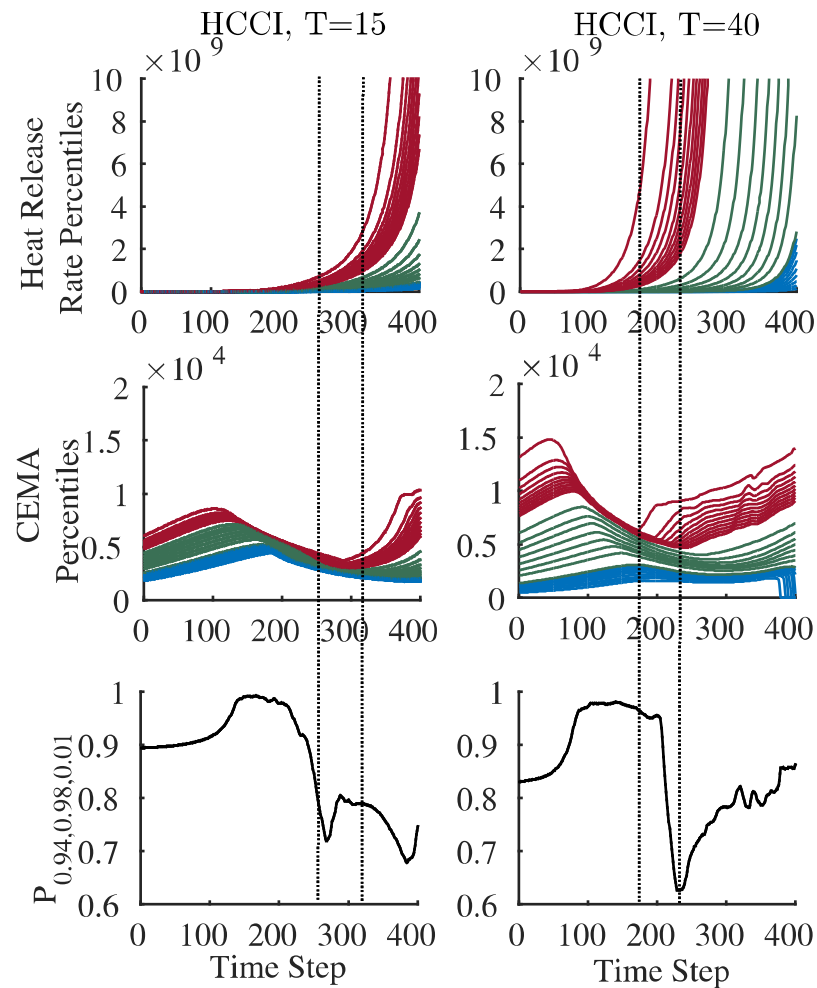
# We defined a simple noise-resistant indicator function to characterize the spread of CEMA percentiles

$$P_{\alpha,\beta,\gamma}(t) = \frac{p_{\alpha}(t) - p_{\gamma}(t)}{p_{\beta}(t) - p_{\gamma}(t)}$$

Avoid use of minimum and maximum percentiles as these are outliers

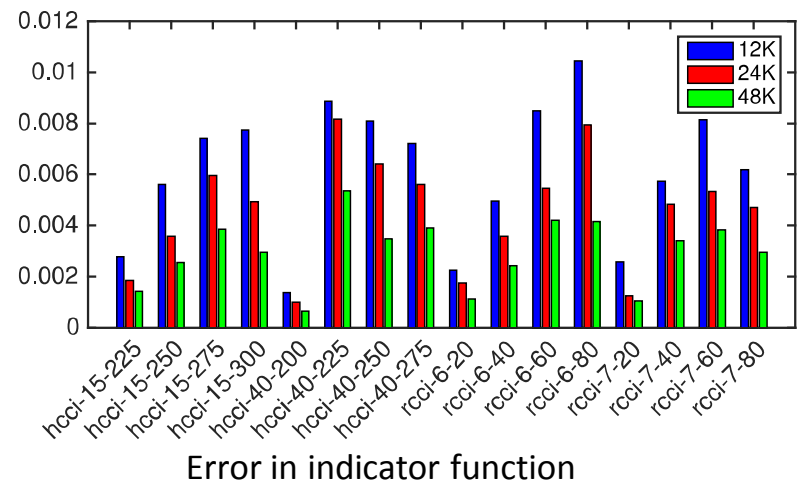
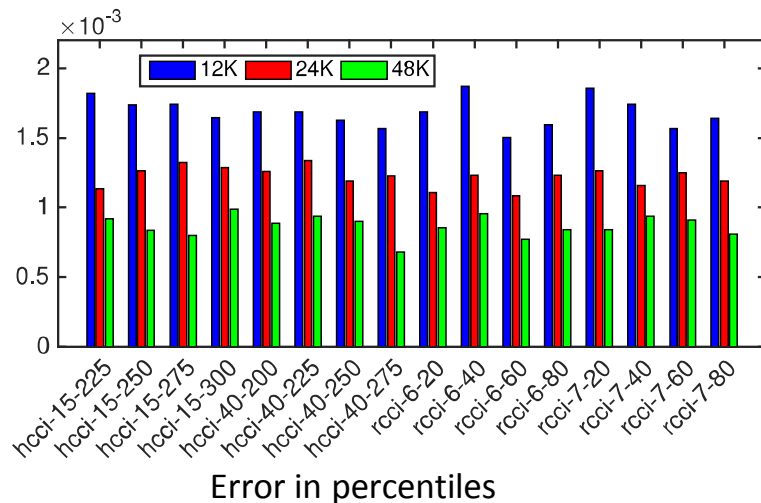


We found that the “ideal” trigger range corresponded to the indicator function passing a threshold from above

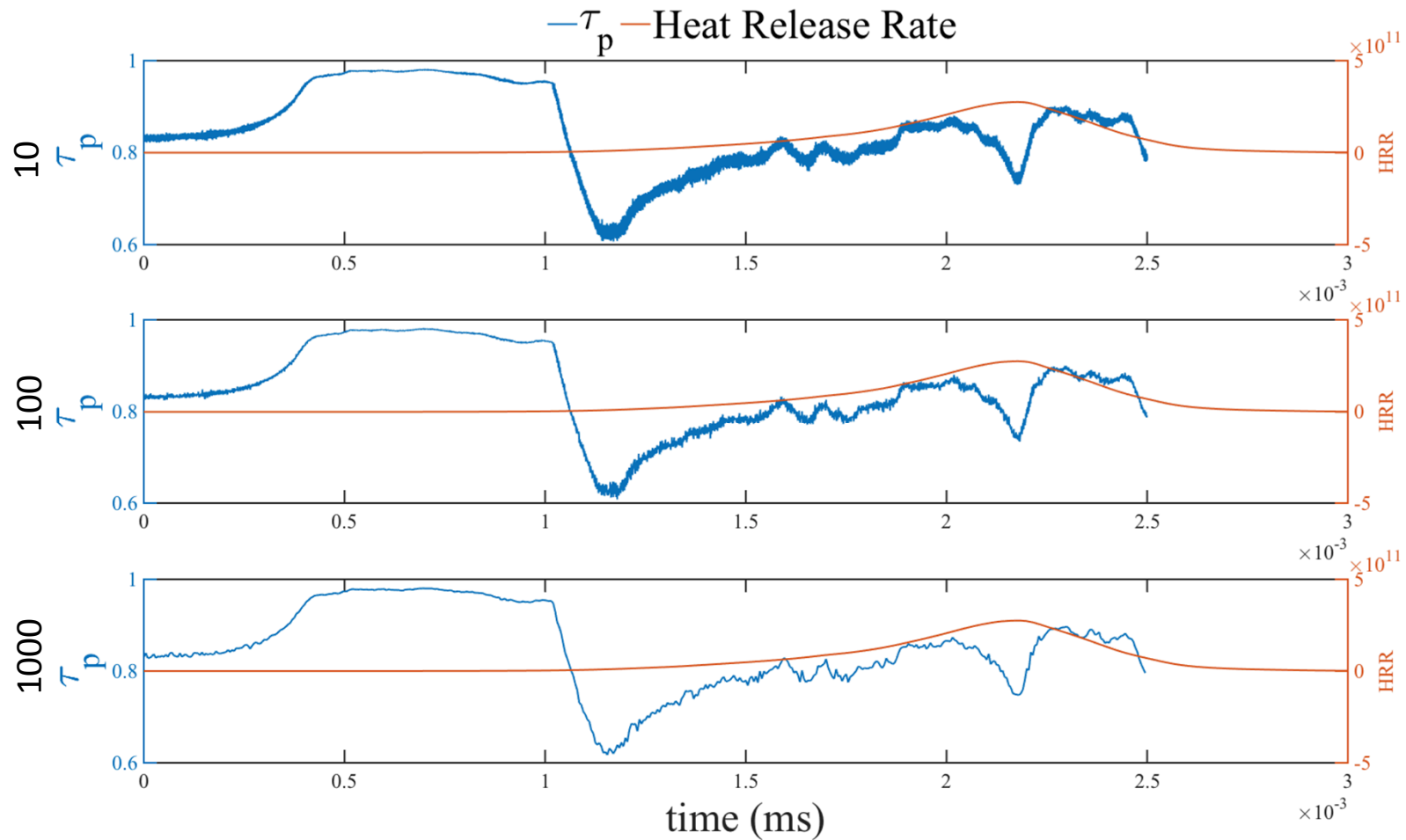


# We developed a simple strategy for sampling percentiles that scales nicely

1. Sample  $k$  independent, uniform indices  $r_1, r_2, \dots, r_k$  in  $\{1, 2, \dots, N\}$ .  
Denote by  $\hat{A}$  the sorted array  $[A(r_1), A(r_2), \dots, A(r_k)]$ .
  2. Output the  $\alpha$ -percentile of  $\hat{A}$  as the estimate,  $\hat{p}_\alpha$ .
- Number of samples depends only on the required accuracy
  - Provide theoretical bounds on error
  - Provide empirical results showing error in practice



# Our approach is being deployed now in-situ to enable dynamic workflows



## MT4: Work with domain experts to understand social constraints that impede adoption of your techniques

- Operating in-situ means you are deploying on production runs
- Scientists are often risk-averse to deploy new technologies in these settings
  - Month-long simulation runs using codes that have evolved over decades
- Sampling was a win in some use cases, not others
  - Final quantities of interest: we encountered hesitance
  - In situ visualization for debugging: positive feedback
  - Control flow decisions: positive feedback
- Proofs and empirical tests together provided the confidence levels the scientists needed to deploy in-situ
  - Neither was sufficient on its own

# Math takeaways for large interdisciplinary research

- MT0: Target a customer and compute platform for your research
- MT1: As the breadth of scope of a project increases so does the social complexity – effective communication is key!
- MT2: Algorithms can look beautiful in theory, but may not be worth deploying in practice due to constraints of your target user and architecture
- MT3: Question your underlying algorithmic assumptions to mitigate technical constraints
- MT4: Work with domain experts to understand social constraints that would impede adoption (often unanticipated)
- MT5: Seeing your theory deployed in practice on real problems at scale is incredibly rewarding and makes MT0-MT4 worth it!