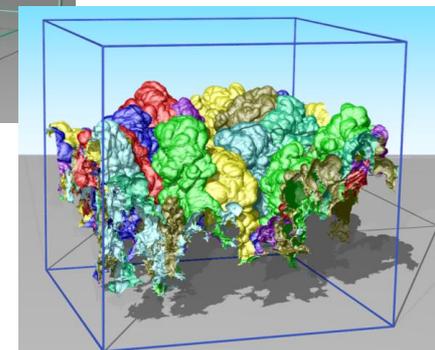
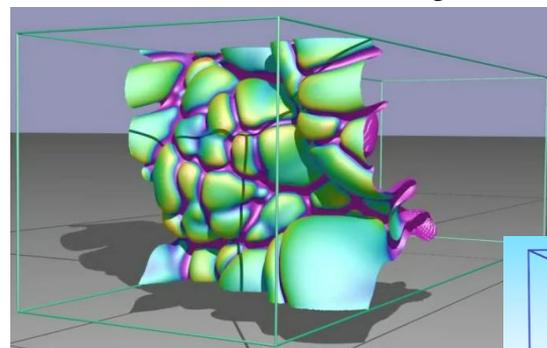
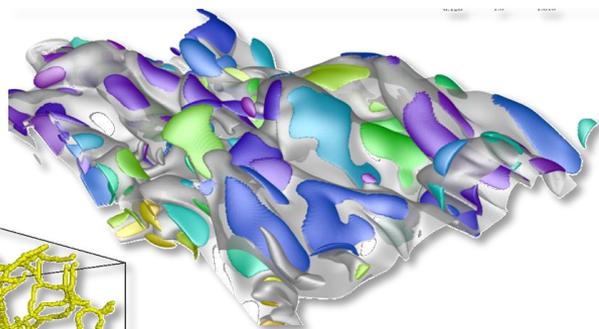
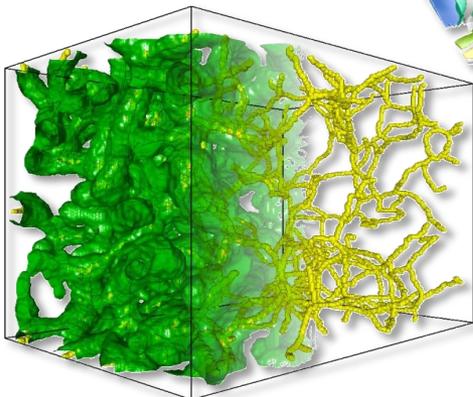


Big Data Analytics for Science Discovery



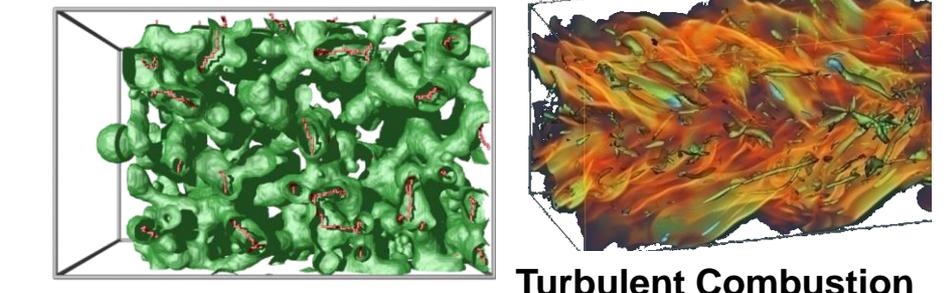
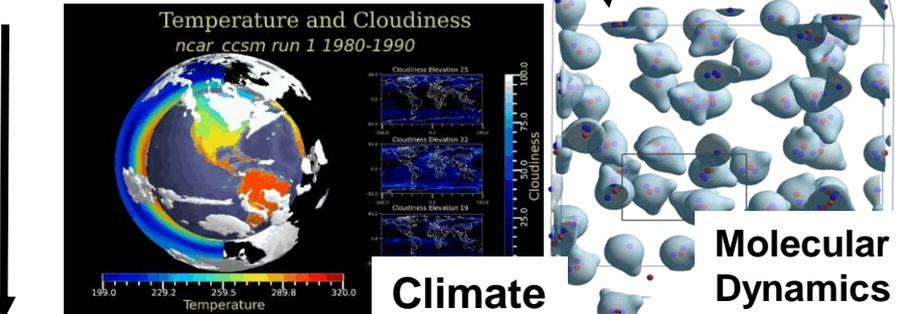
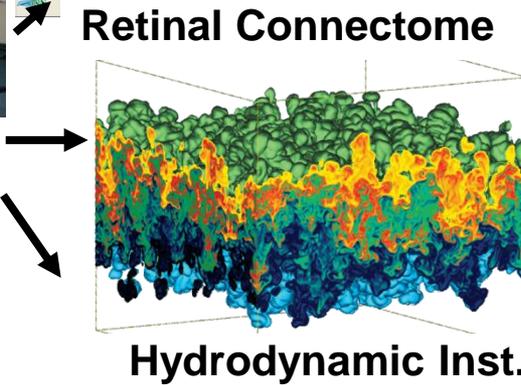
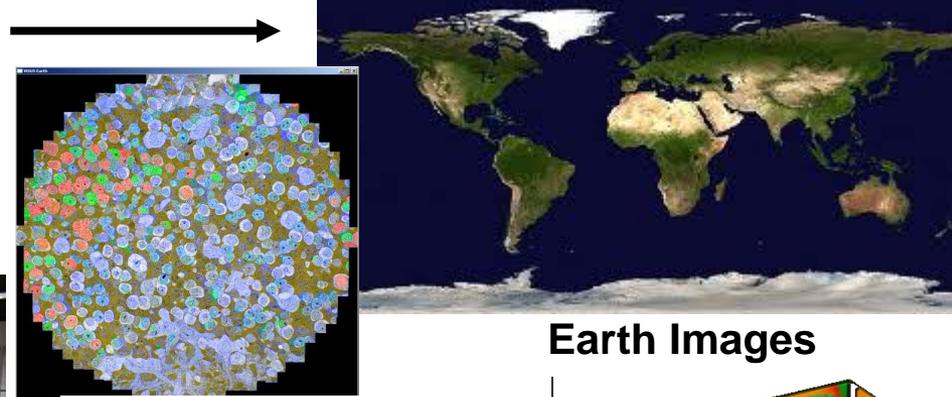
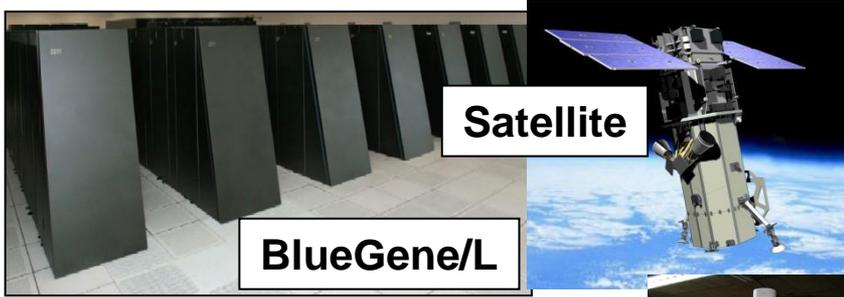
Valerio Pascucci

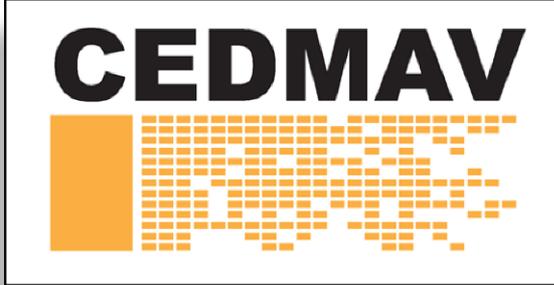
Director, Center for Extreme Data Management Analysis and Visualization

Professor, SCI institute and School of Computing, University of Utah

Laboratory Fellow, Pacific Northwest National Laboratory

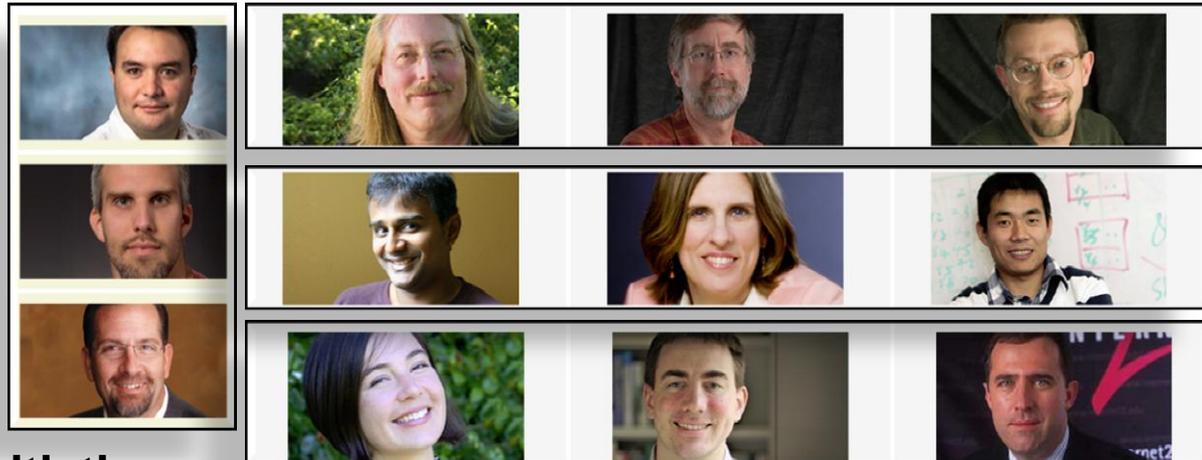
Massive Simulation and Sensing Devices Generate Great Challenges and Opportunities



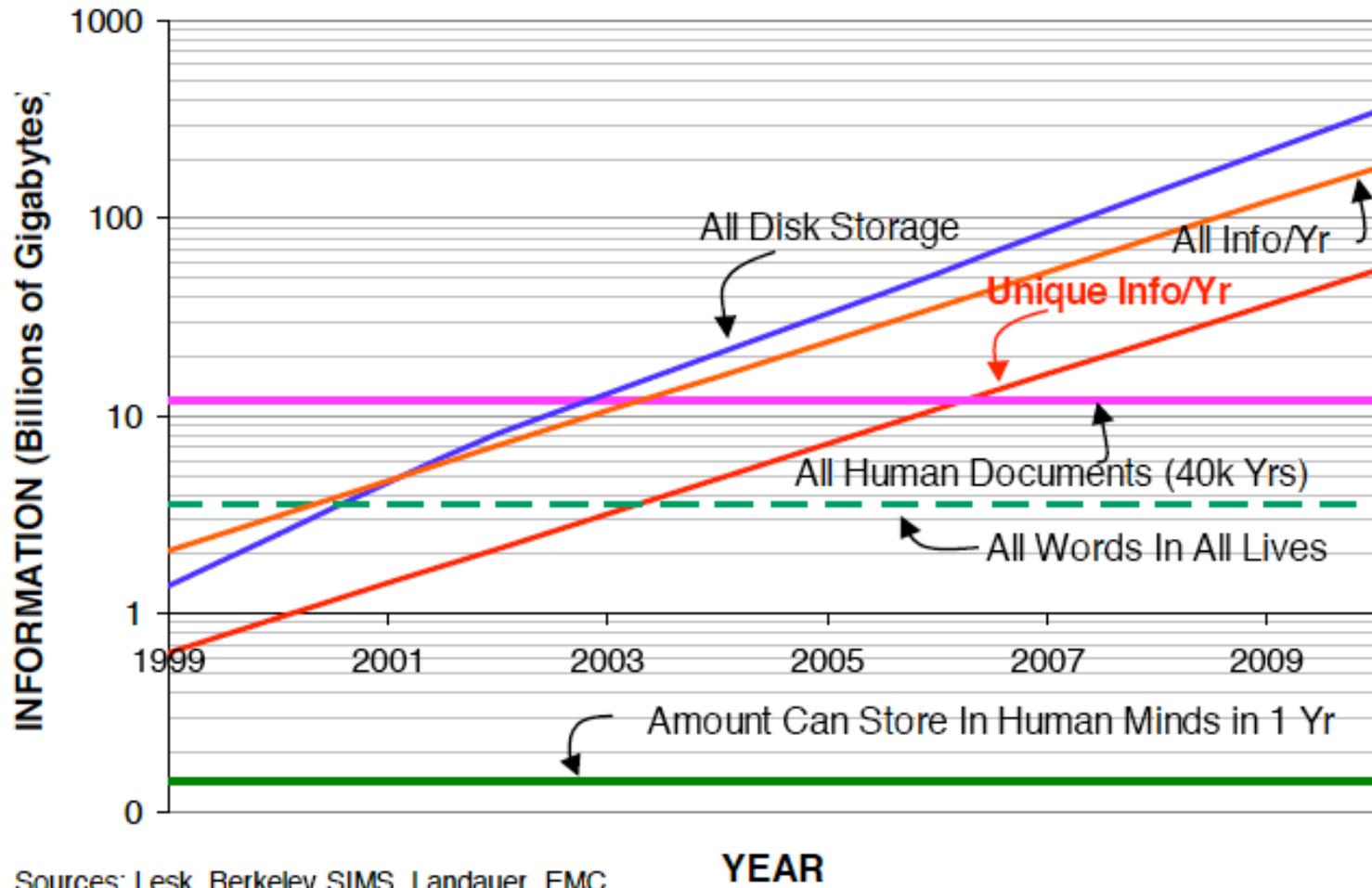


Center for Extreme Data Management, Analysis, and Visualization

- 10 Faculty + scientists, developers, students, ...
- Primary partners: UU & PNNL
- Other partnerships: NSA, INL, LLNL, ANL, Battelle,
- Involvement in national Initiatives



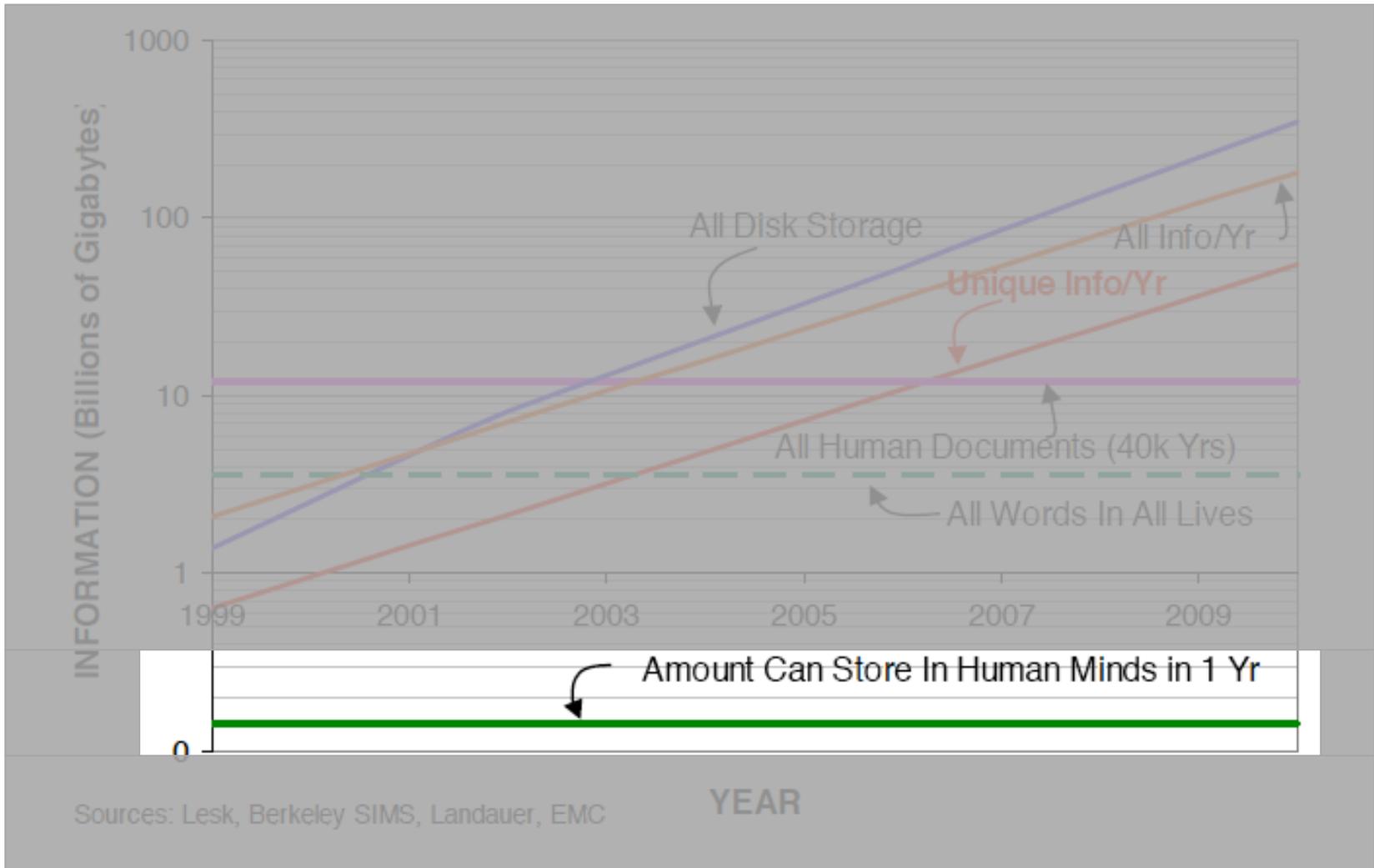
The Information Big Bang Has Come!



Sources: Lesk, Berkeley SIMS, Landauer, EMC

YEAR

Are We Hopeless ?



How Much Information Is This?

01100101

00111010

11011010

11000110

00110010

How Much Information Is This?

E

00111010

11011010

11000110

00110010

How Much Information Is This?

E

=

11011010

11000110

00110010

How Much Information Is This?

E

=

m

11000110

00110010

How Much Information Is This?

E

=

m

c

00110010

How Much Information Is This?

E

=

m

c

2

How Much Information Is This?

$$E=mc^2$$

Multiresolution vs Abstraction



Multiresolution vs Abstraction



Multiresolution vs Abstraction



Multiresolution vs Abstraction



Multiresolution vs Abstraction



Multiresolution vs Abstraction



Multiresolution vs Abstraction



Multiresolution vs Abstraction



Multiresolution vs Abstraction



Multiresolution vs Abstraction

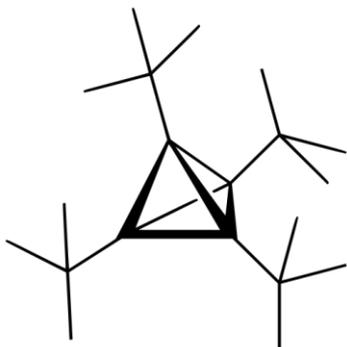
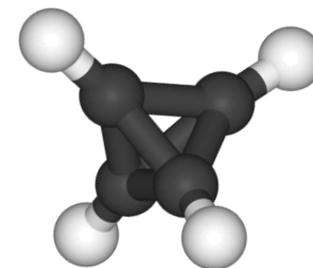
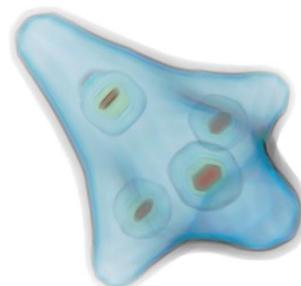


FLOWERS

Abstractions Have Been Used in Science Discovery/Communication for a Long Time

Abstraction

```
1101010100101010100101010100101010010101010101011101  
0010101010010101001010101001010101010101010111101010  
10101010010101010001110101010110101010101110101010  
10100001011010111101011001111001100111100110011001  
01010101111010101000101010101011010101010001010100  
010101001010101010101010010101010110001001010010  
11010101001010101010001010101111011011100010101010  
10101001010100101010010010010101010101010101010101
```

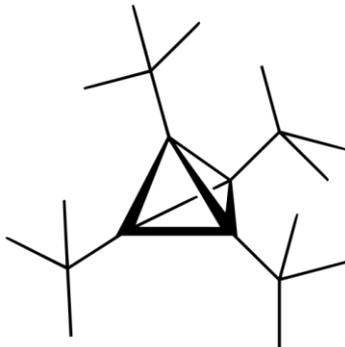
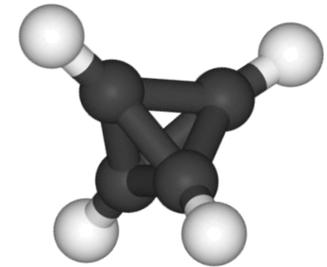
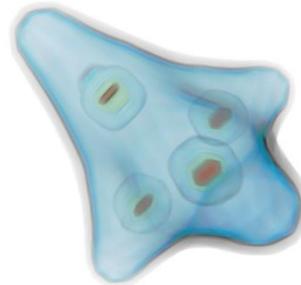


Tetrahedrane

Abstractions Have Been Used in Science Discovery/Communication for a Long Time

Abstraction

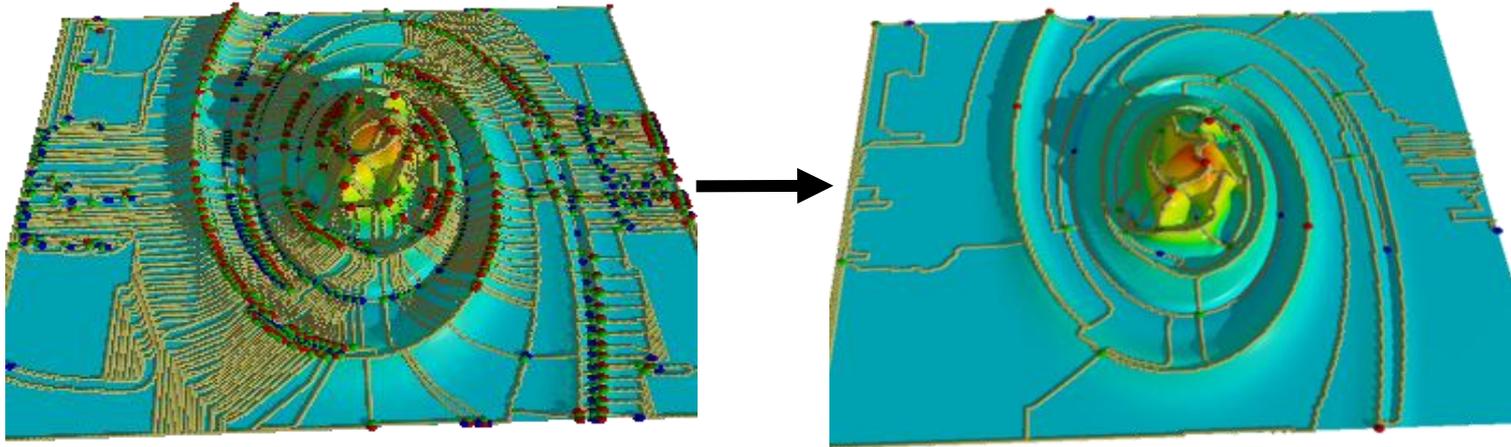
1101010100101010100101010100101010010101010101011101
0010101010010101001010101001010101010101010111101010
101010100101010100011101010101101010101110101010
10100001011010111101011001111001100111100110011001
01010101111010101000101010101011010101010001010100
010101001010101010101010010101010110001001010010
11010101001010101010001010101111011011100010101010
10101001010100101010010010010101010101010101010101



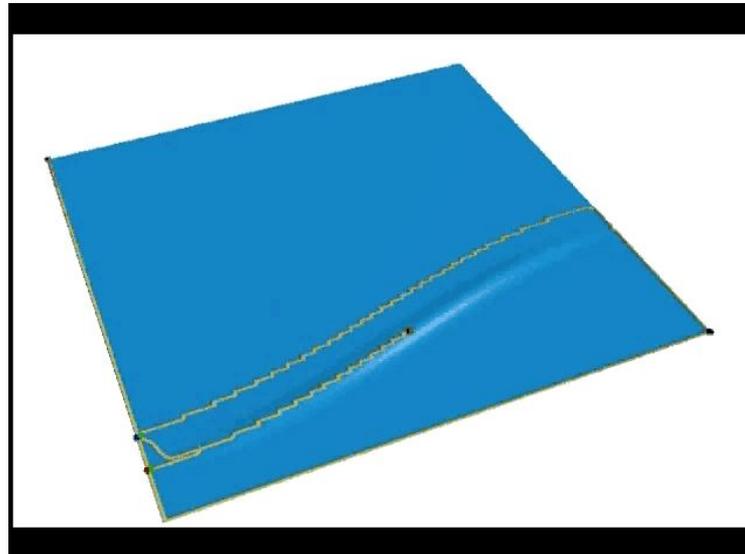
Tetrahedrane

Language

Topology is an Effective Language to Describe Abstractions of Features from Raw Data

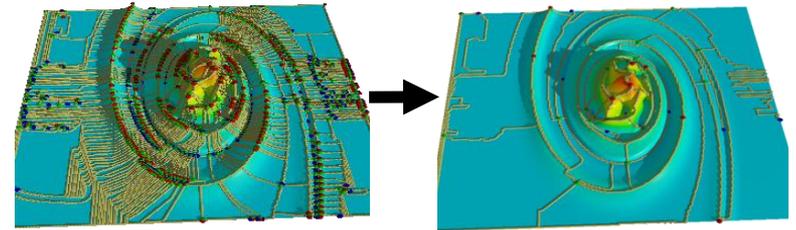


Hierarchical topology of a 2D Miranda vorticity field

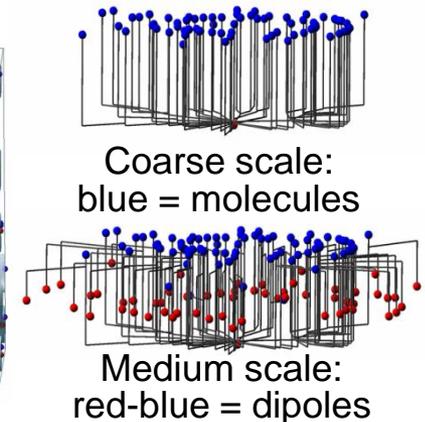
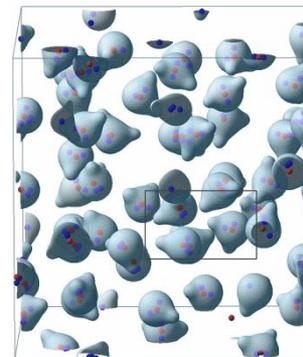


Our Framework is Based on Robust Topological Computations for Quantitative Data Analysis

- Provably robust computation
- Provably complete feature extraction and quantification
- Hierarchical topological structures used to capture multiple scales
- Error-bounded approximations associated with each scale
- Formal mathematical definition associated with each analysis
- Scalable performance in association with streaming techniques



Hierarchical topology of a 2D Miranda vorticity field

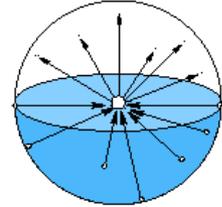


Molecular dynamics simulation (left) with abstract graph representation of its features at two scales (right)

We Adopted a Combinatorial Approach to Morse Theory for Provably Correct Computations

QED

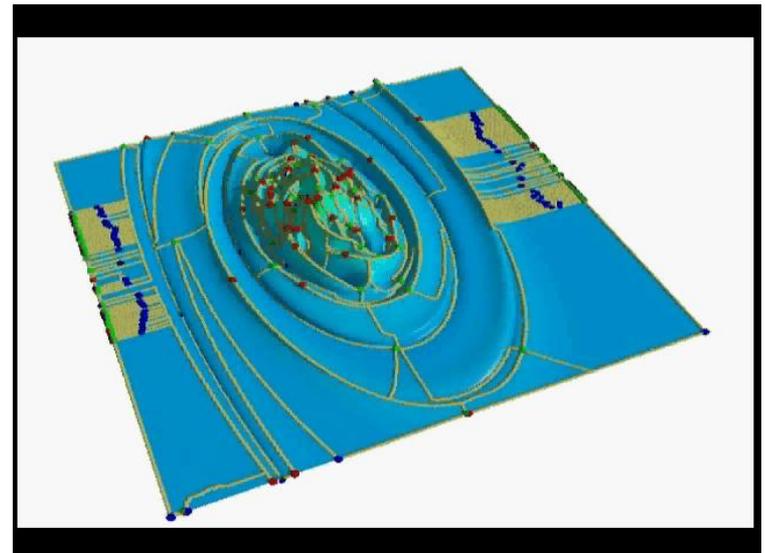
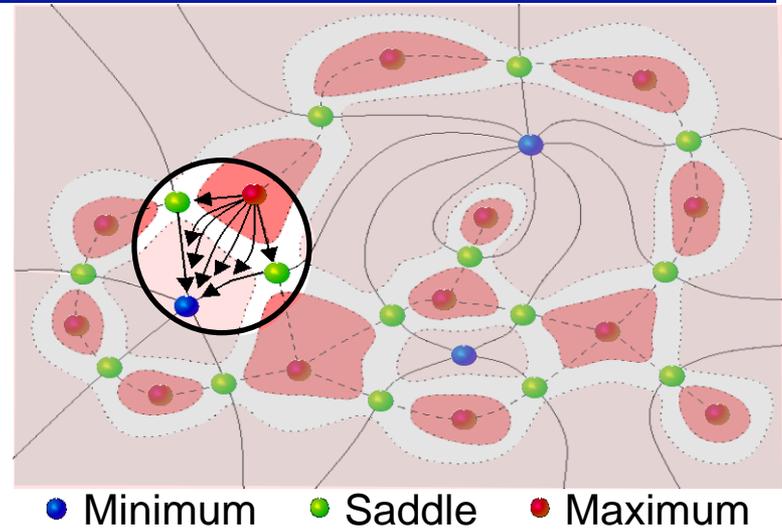
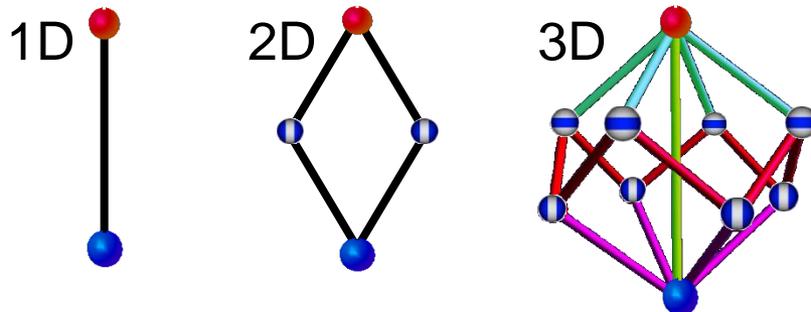
QSS

	Classical mathematical definitions	Simulation of differentiability
domain	D smooth manifold	S simplicial complex
function	f infinitely differentiable	$F(x)$ PL-extension of $f(x)$
critical point	<p><i>QED</i></p> <p>numerical</p>  <p>1D 2D</p>	<p><i>QSS</i></p> <p>combinatorial</p>  <p>3D</p>

Independent local computation yield globally consistent results

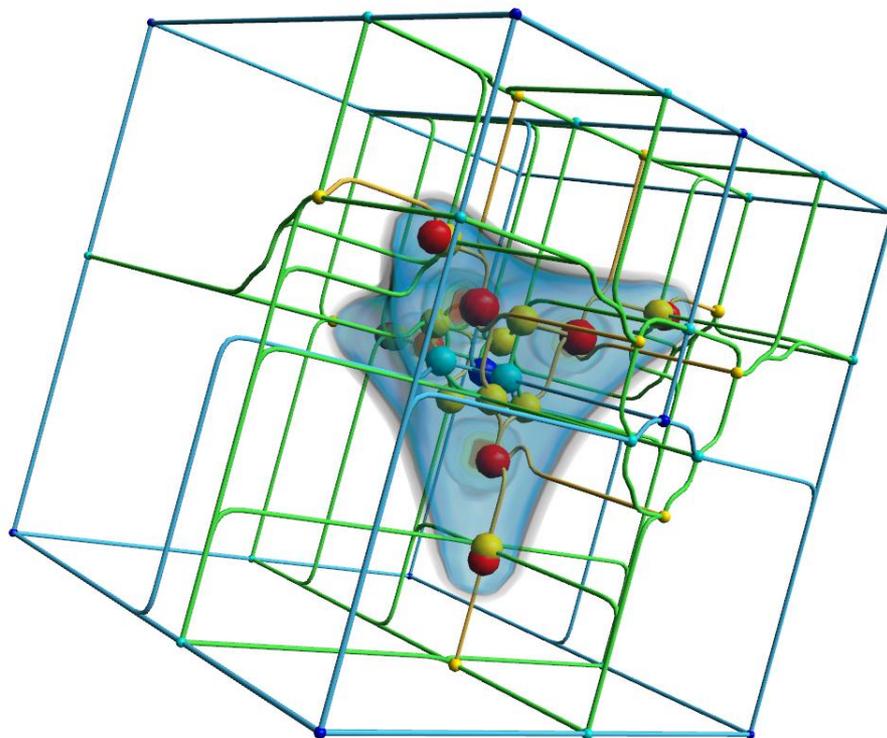
We Introduced the Morse–Smale Complex for Complete Data Analysis

- The Morse–Smale complex partitions the domain of f in regions of uniform gradient
- Generalizes the notion of monotonic interval
- Dimension of a region equal index difference of source and destination
- Remove inconsistency of local gradient evaluations



Demo C₄H₄

Morse3d Efficient Computation

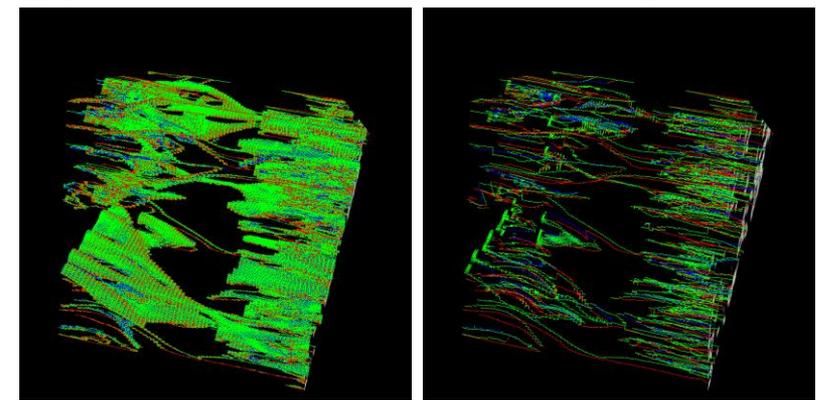
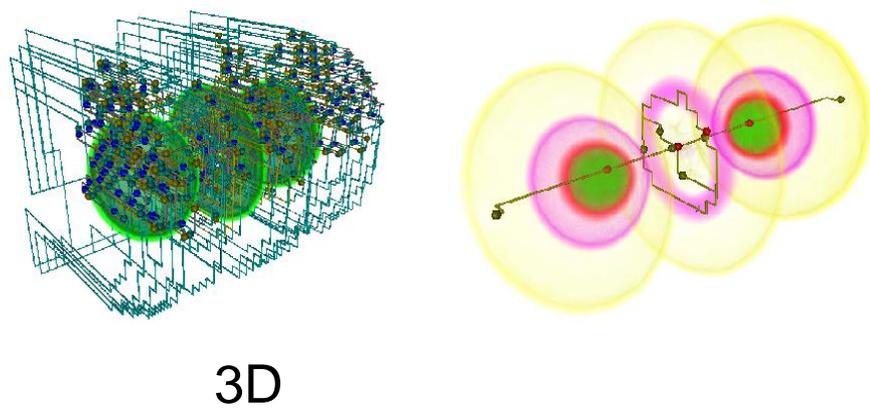
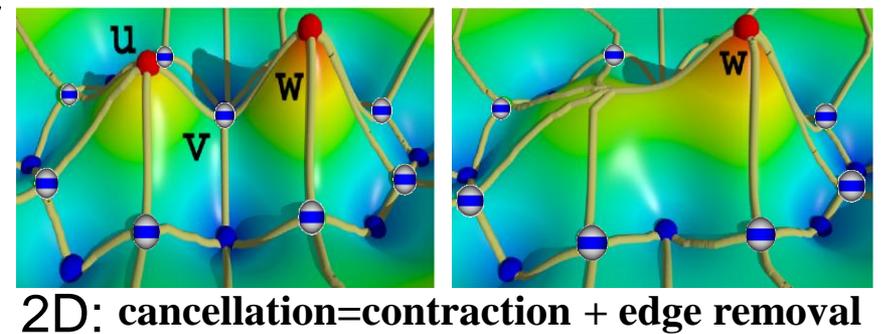
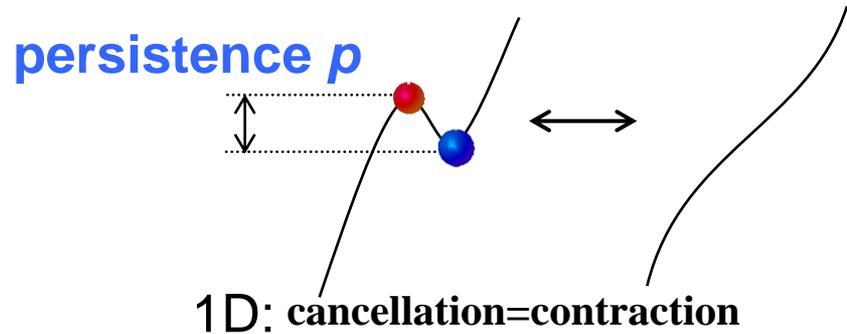


We Use Cancellations to Create a Multi-scale Representation of the Trends in the Data

Cancellations: critical points can be created or destroyed in pairs that are connected 1-manifolds

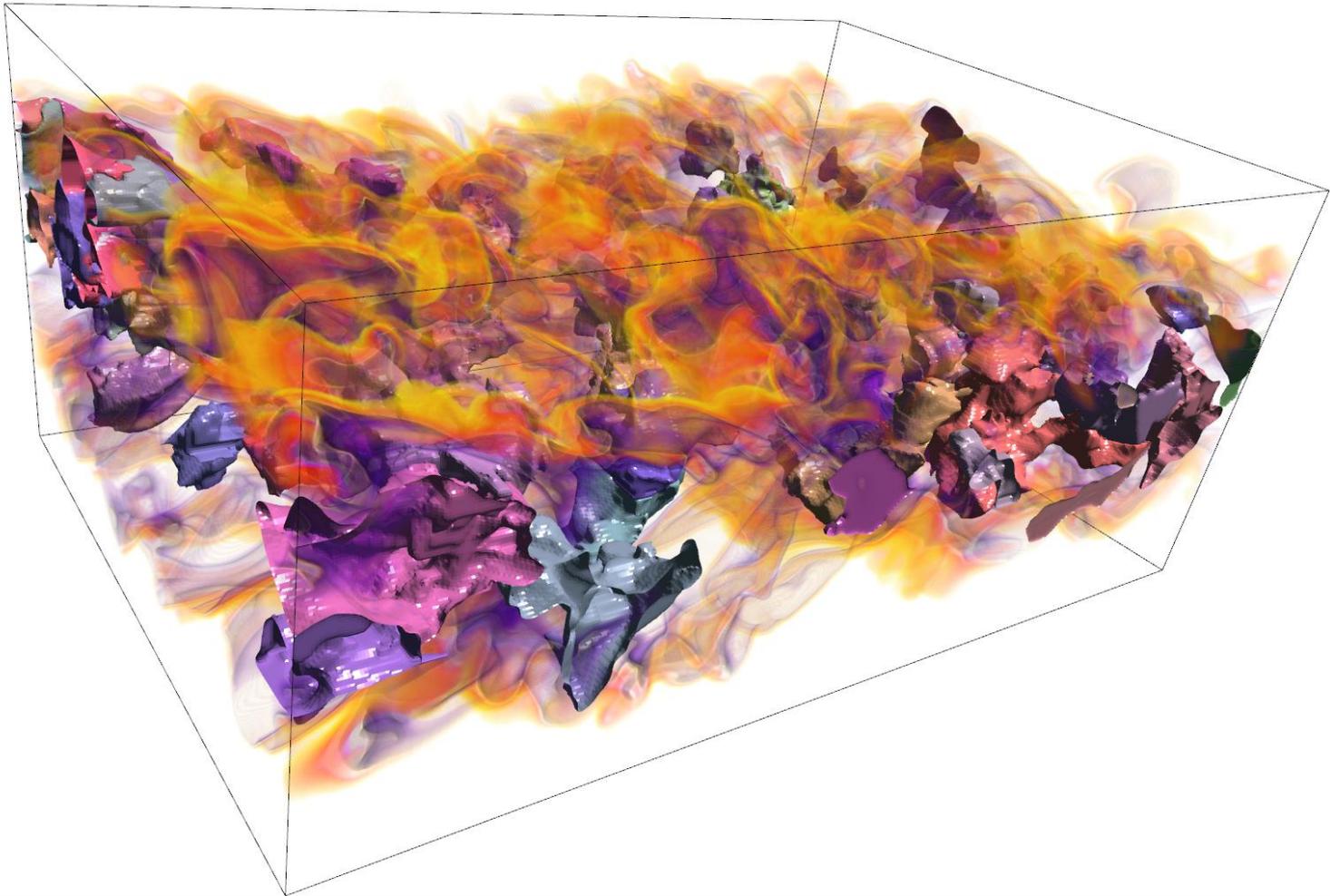
Approximation: error = persistence/2 (proven lower bound)

Multi-scale: consistent gradient segmentation at all scales

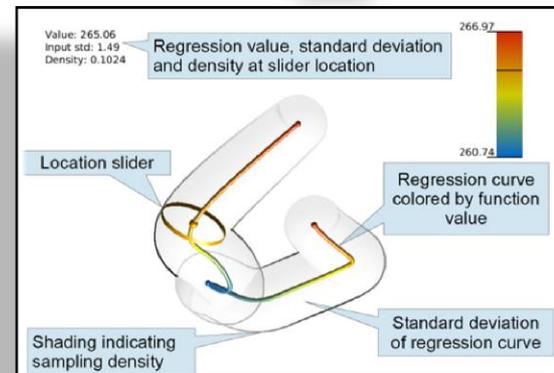
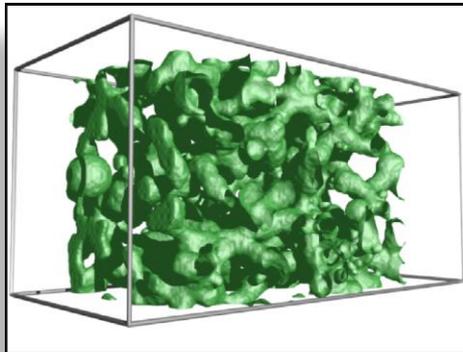
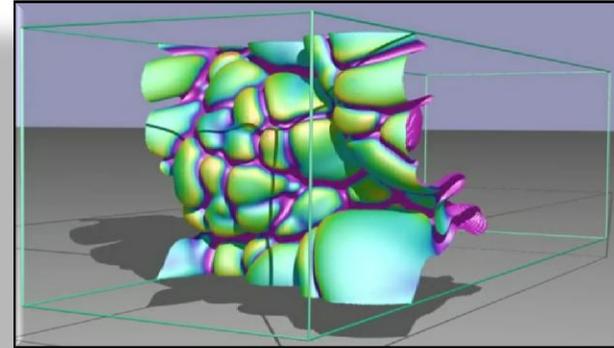
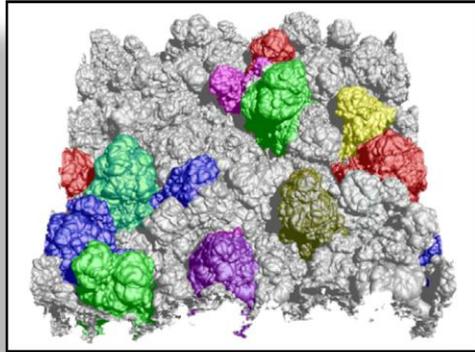


Demo S3D Combustion Simulation

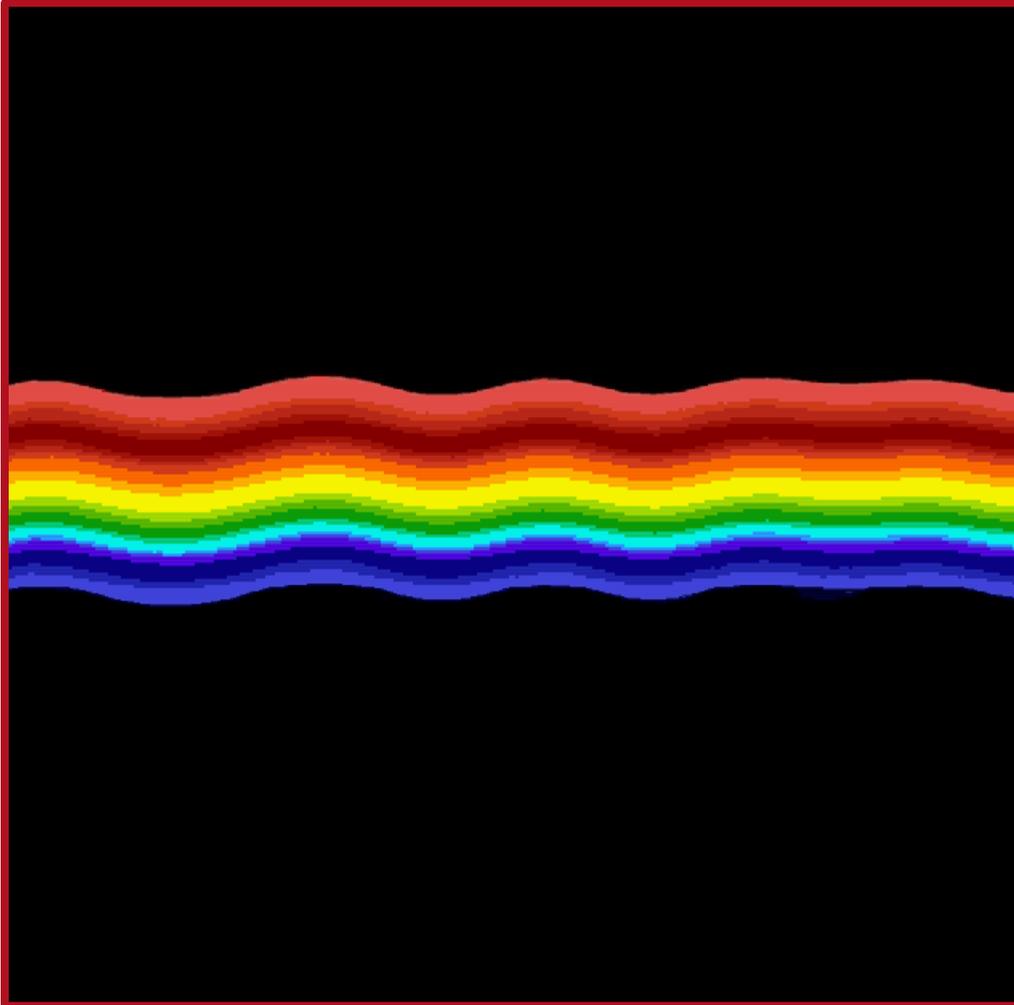
Morse3d Efficient Computation



Big Data Analytics Success Stories



Count the Number of Bubbles in a Rayleigh–Taylor Instability

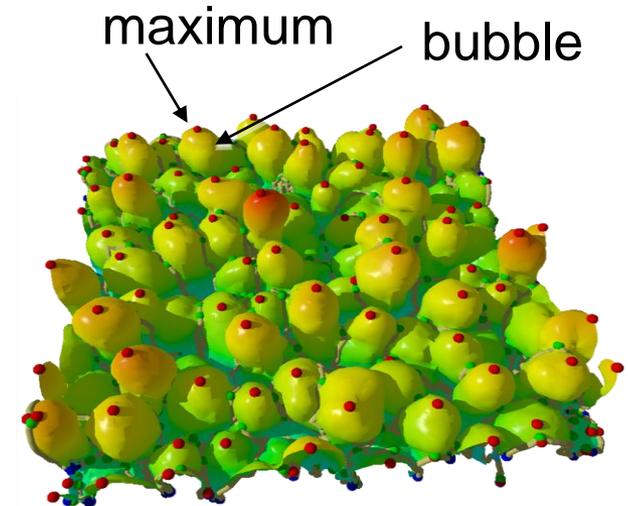


Rayleigh-Taylor instabilities arise in fusion, super-novae, and other fundamental phenomena:

- start: heavy fluid above, light fluid below
- gravity drives the mixing process
- the mixing region lies between the upper envelope surface (red) and the lower envelope surface (blue)
- 25 to 40 TB of data from simulations

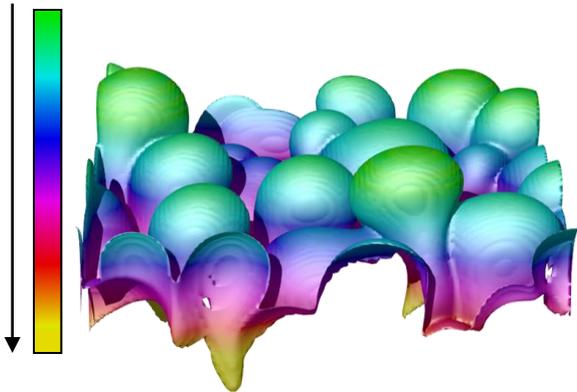
We Analyze High-Resolution Rayleigh–Taylor Instability Simulations

- **Large eddy simulation run on Linux cluster: 1152 x 1152 x 1152**
 - ~ 40 G / dump
 - 759 dumps, about 25 TB
- **Direct numerical simulation run on BlueGene/L: 3072 x 3072 x Z**
 - Z depends on width of mixing layer
 - More than 40 TB
- **Bubble-like structures are observed in laboratory and simulations**
- **Bubble dynamics are considered an important way to characterize the mixing process**
 - Mixing rate = $\partial(\#bubbles) / \partial t$.
- **There is no prevalent formal definition of bubbles**

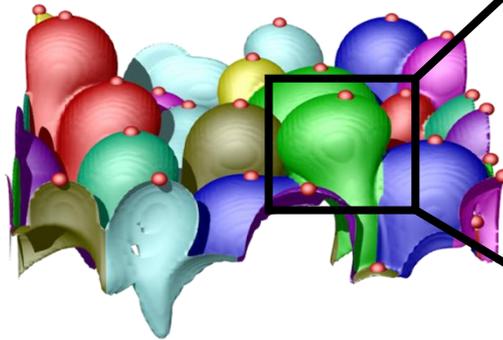


We Compute the Morse–Smale Complex of the Upper Envelope Surface

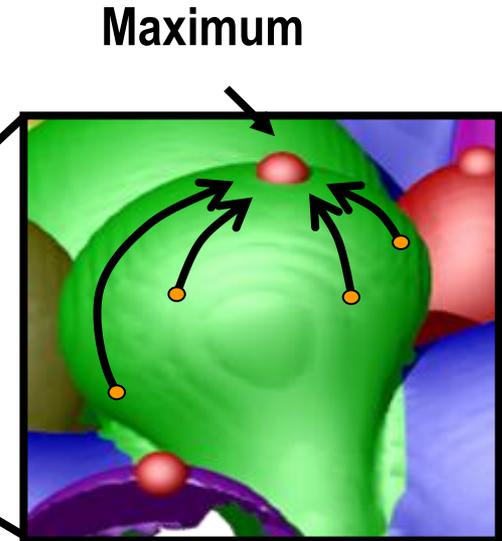
$$F(x) = z$$



$F(x)$ on the surface is aligned against the direction of gravity which drives the flow



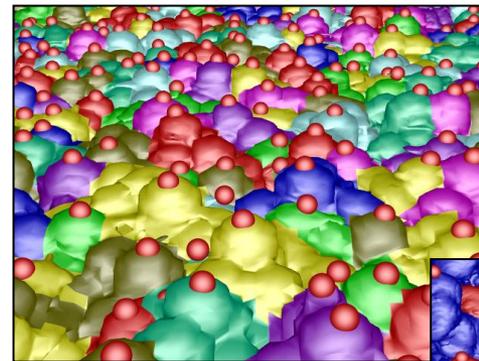
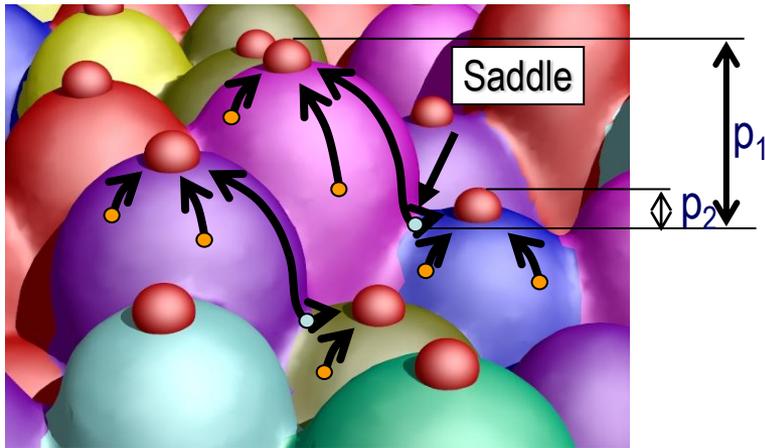
Morse complex cells drawn in distinct colors



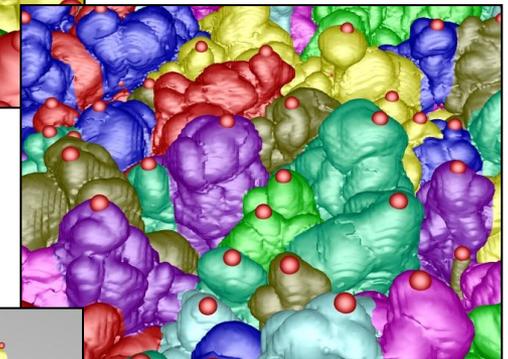
In each Morse complex cell, all steepest ascending lines converge to one maximum

A Hierarchical Model is Generated by Simplification of Critical Points

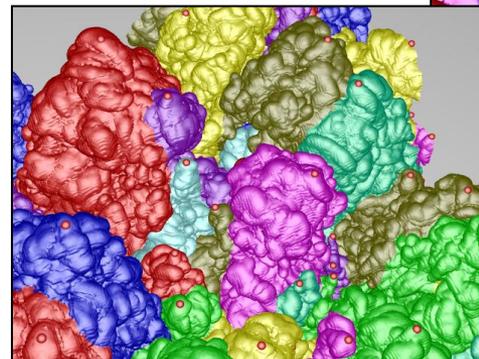
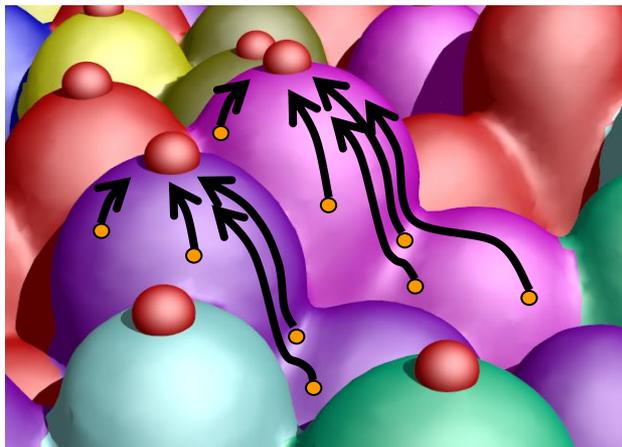
- Persistence is used to produce coarse segmentations
- Coarse scales preserve high-persistence critical points



$T=100$

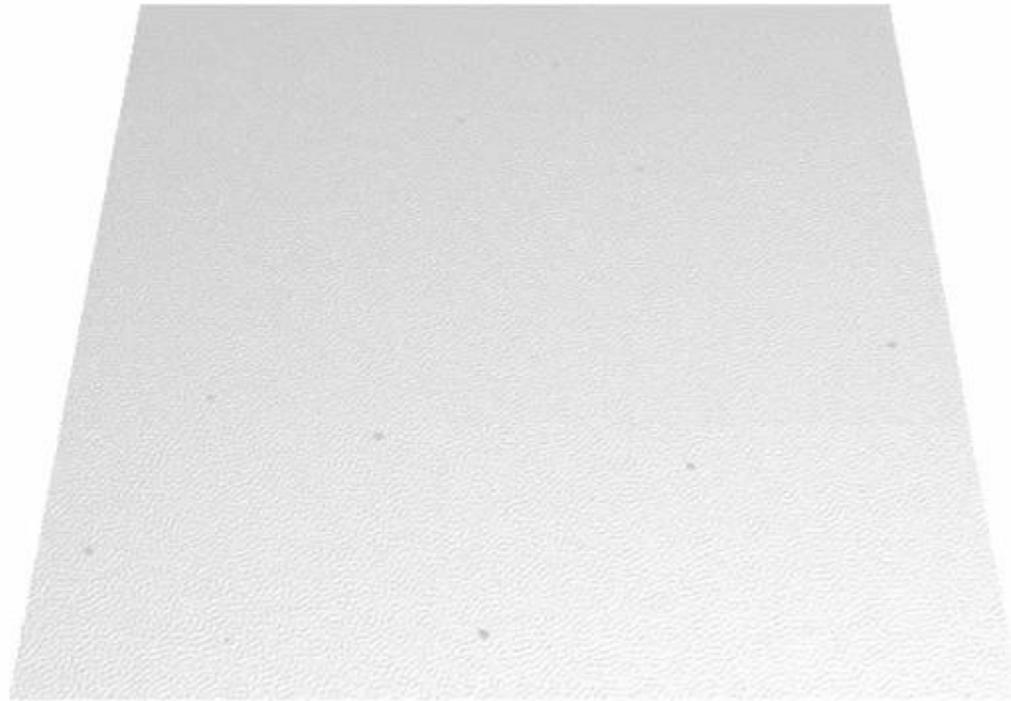


$T=353$

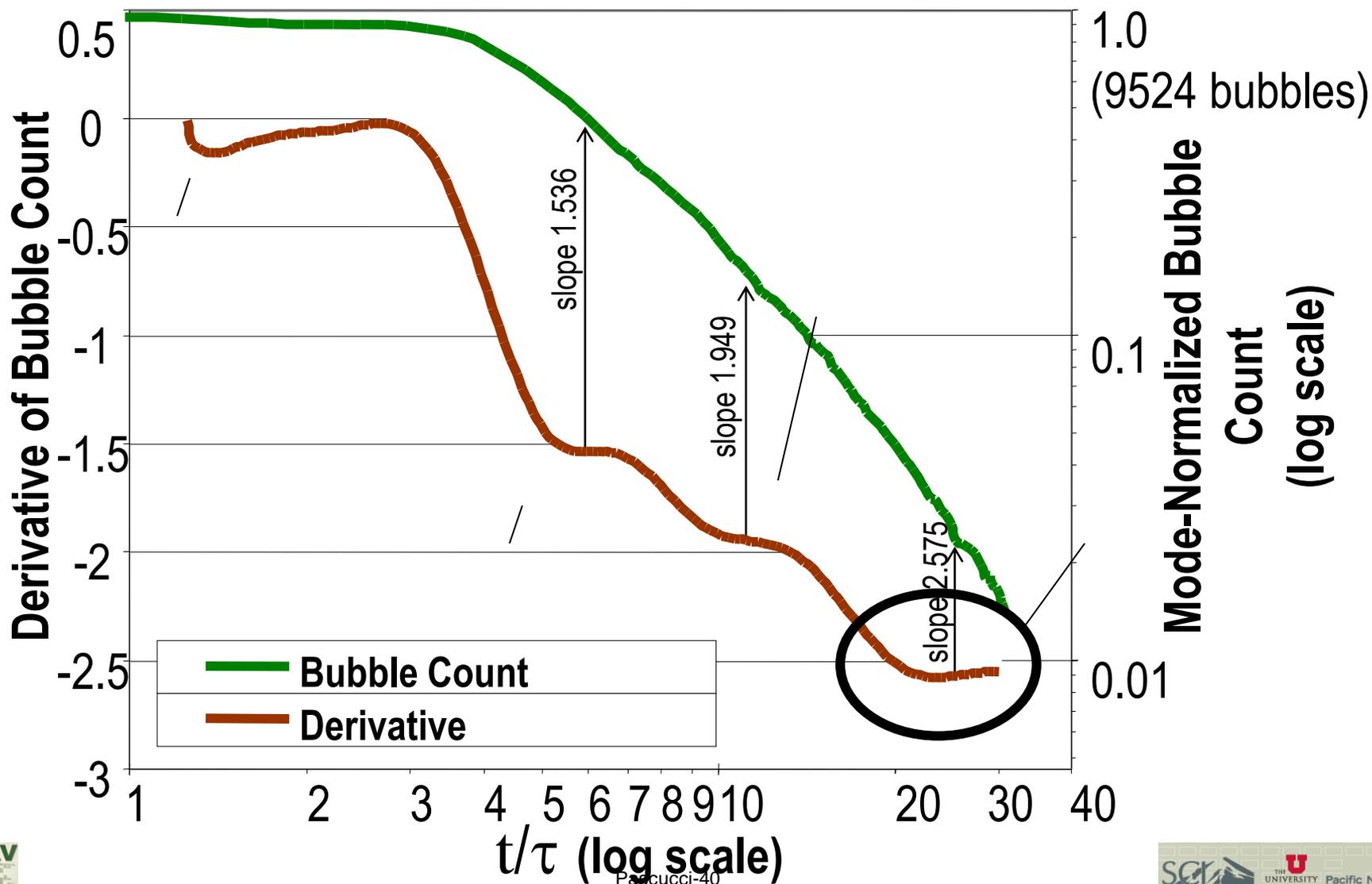


$T=700$

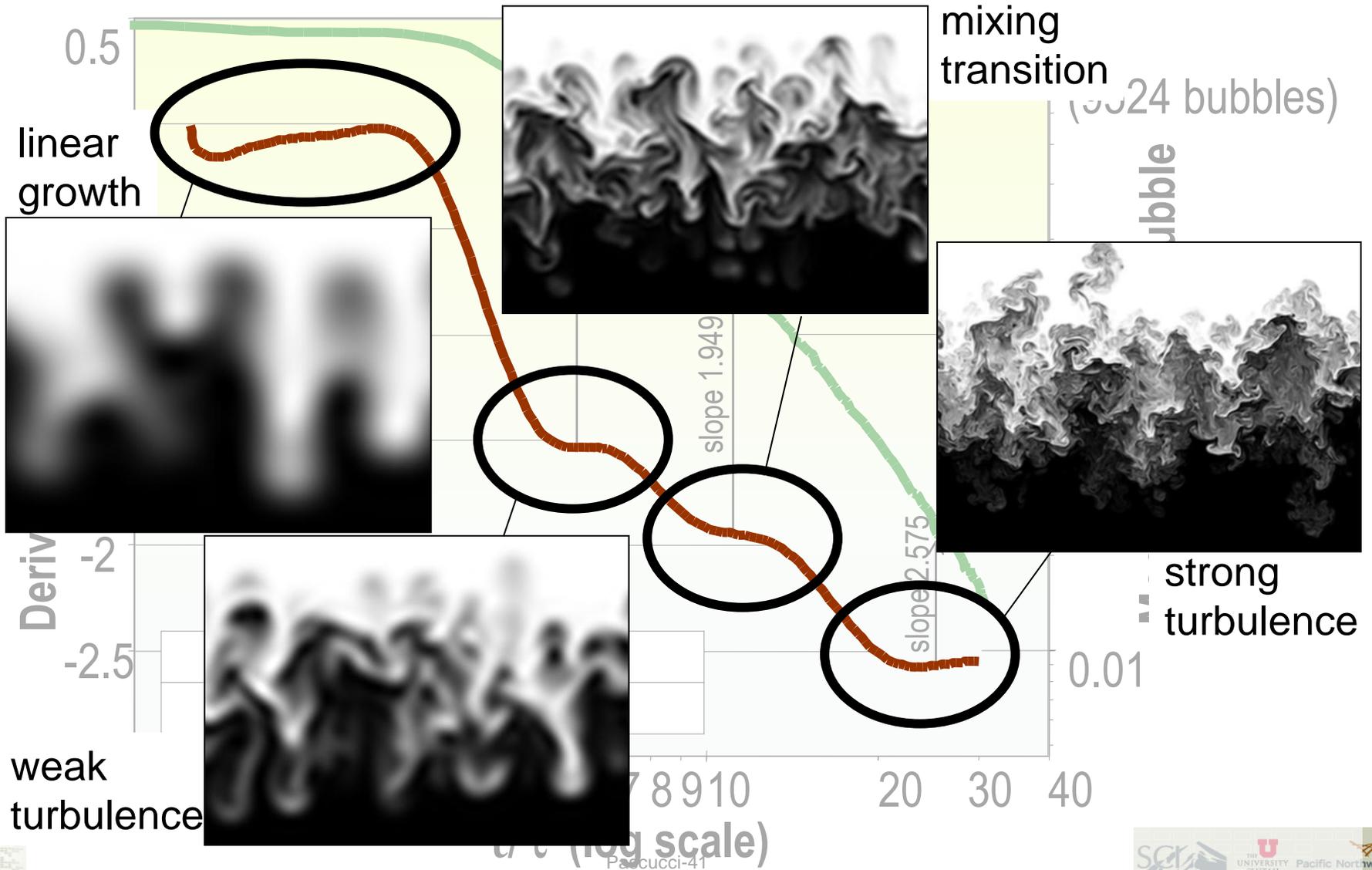
First Robust Bubble Tracking From Beginning to Late Turbulent Stages



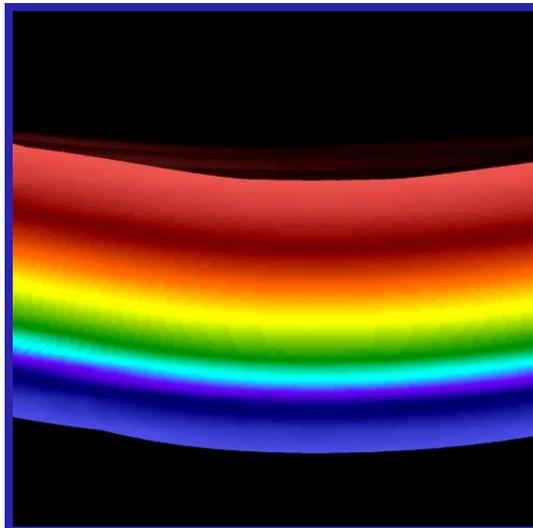
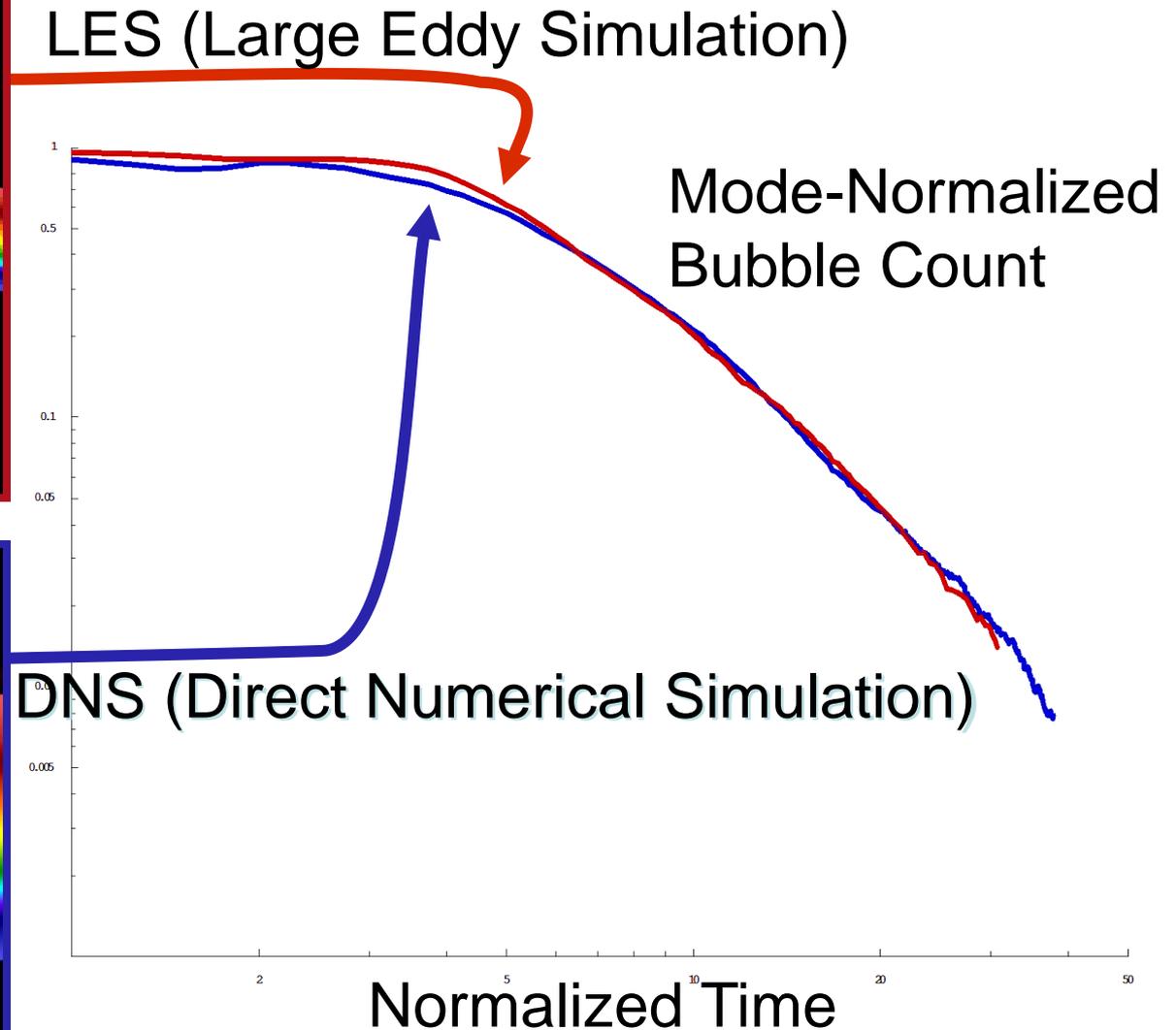
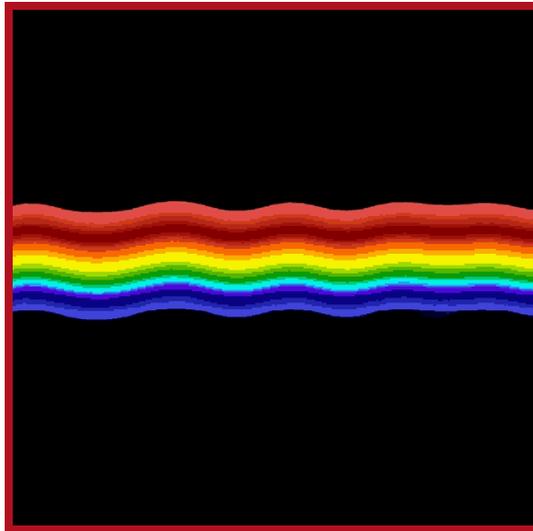
First Time Scientists Can Quantify Robustly Mixing Rates by Bubble Count



We Provide the First Quantification of Known Stages of the Mixing Process

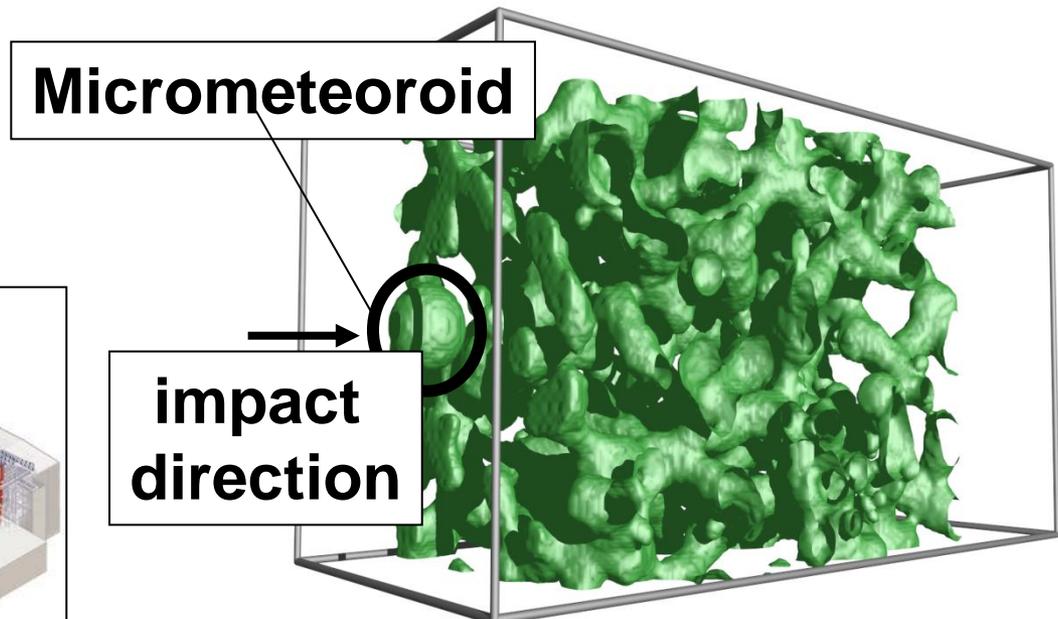
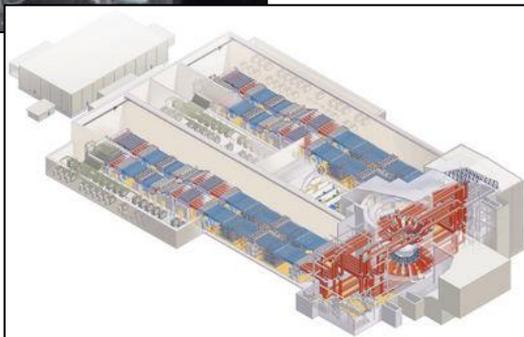
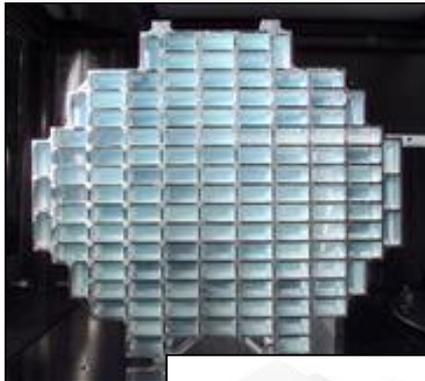
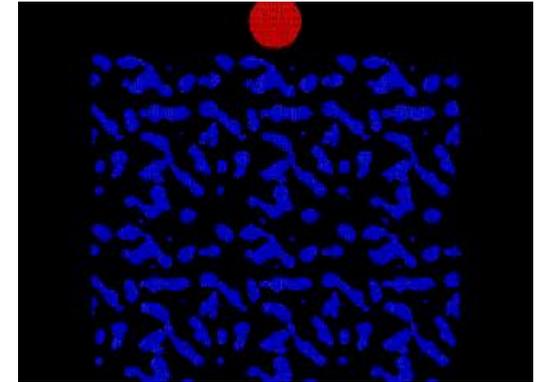


We Provided the First Feature-Based Validation of a LES with Respect to a DNS



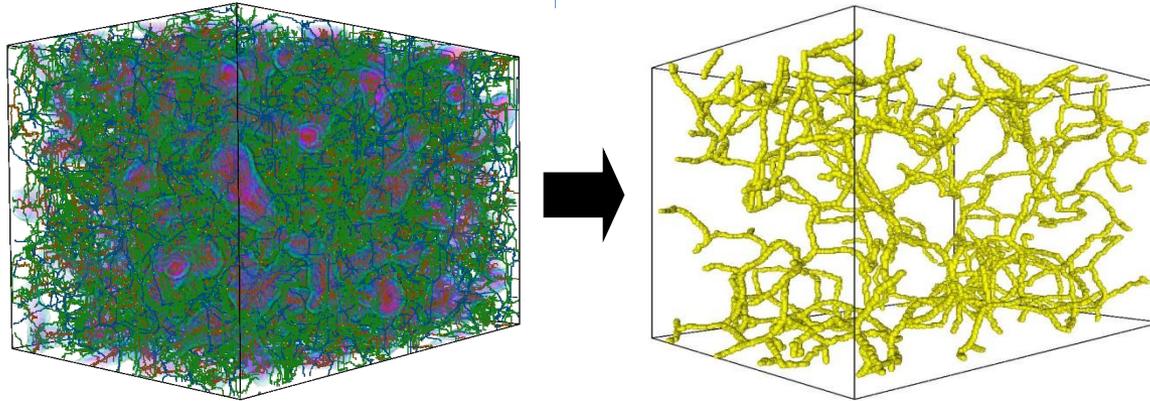
Quantitative Analysis of the Impact of a Micrometeoroid in a Porous Medium

- Many possible applications:
 - NASA's Stardust Spacecraft
 - National Ignition Facility Targets
 - Light and Robust Materials
 - many more...

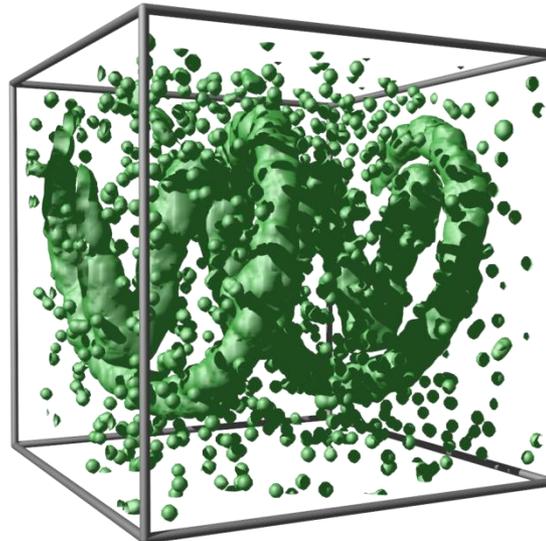


The Topological Reconstruction Method is Validated with a Controlled Test Shape

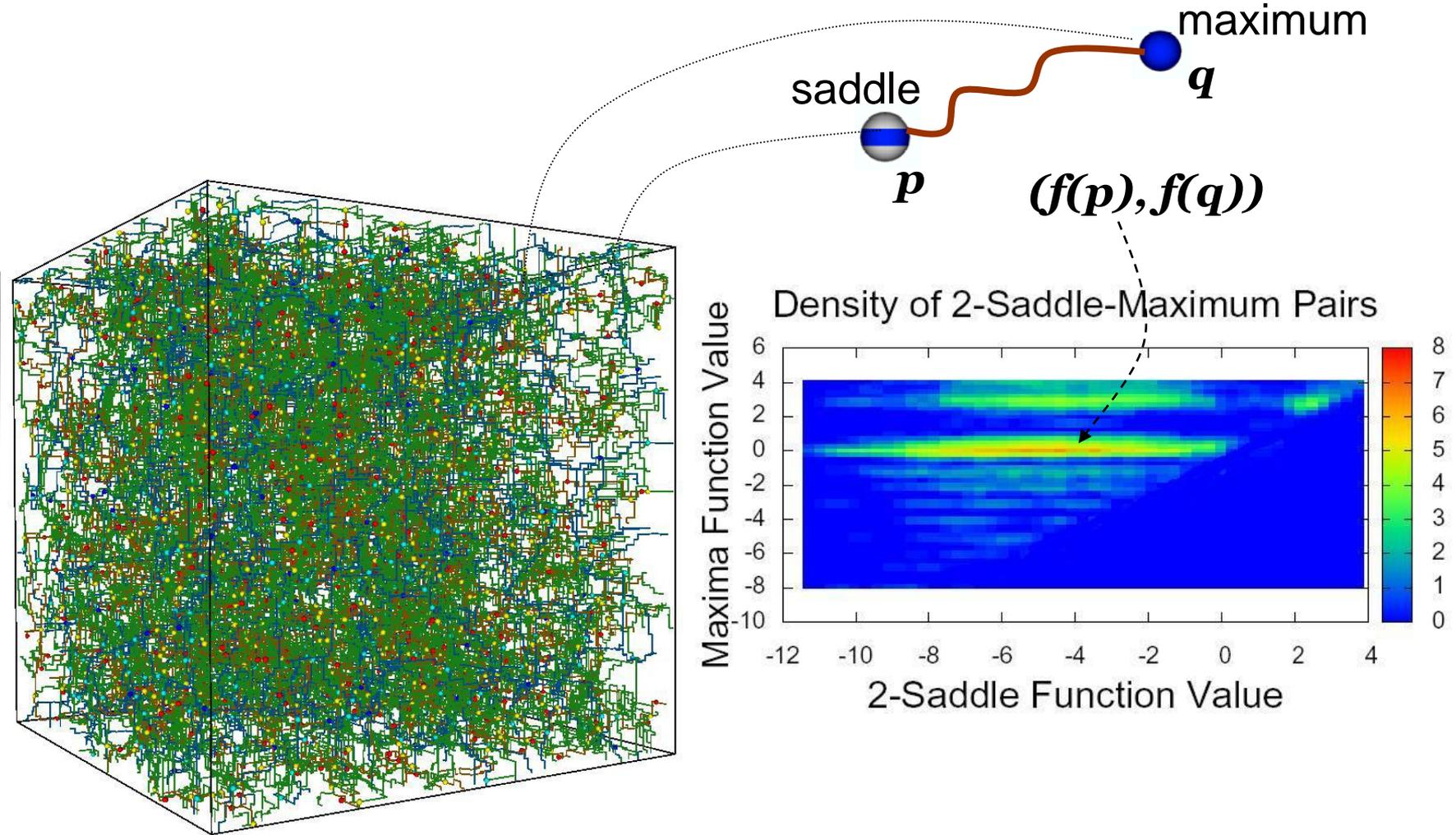
Challenge: robust reconstruction of the structure of a porous medium



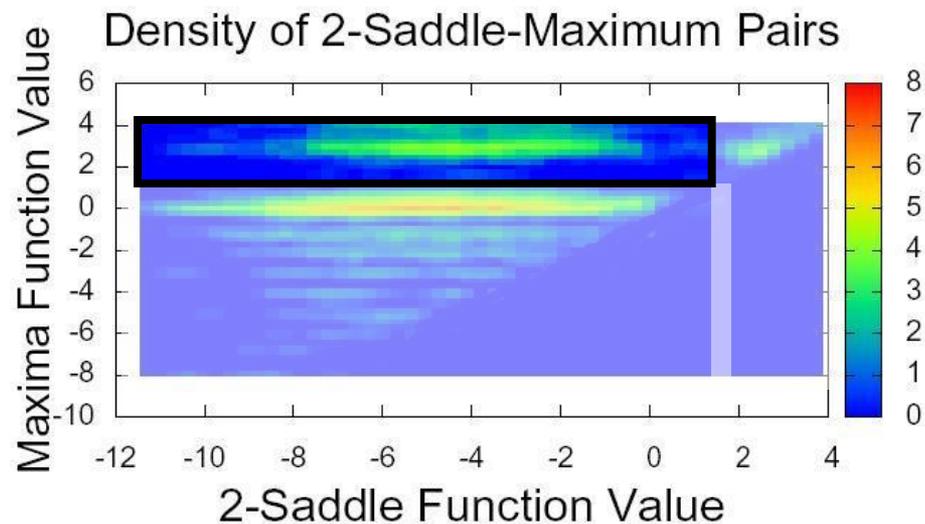
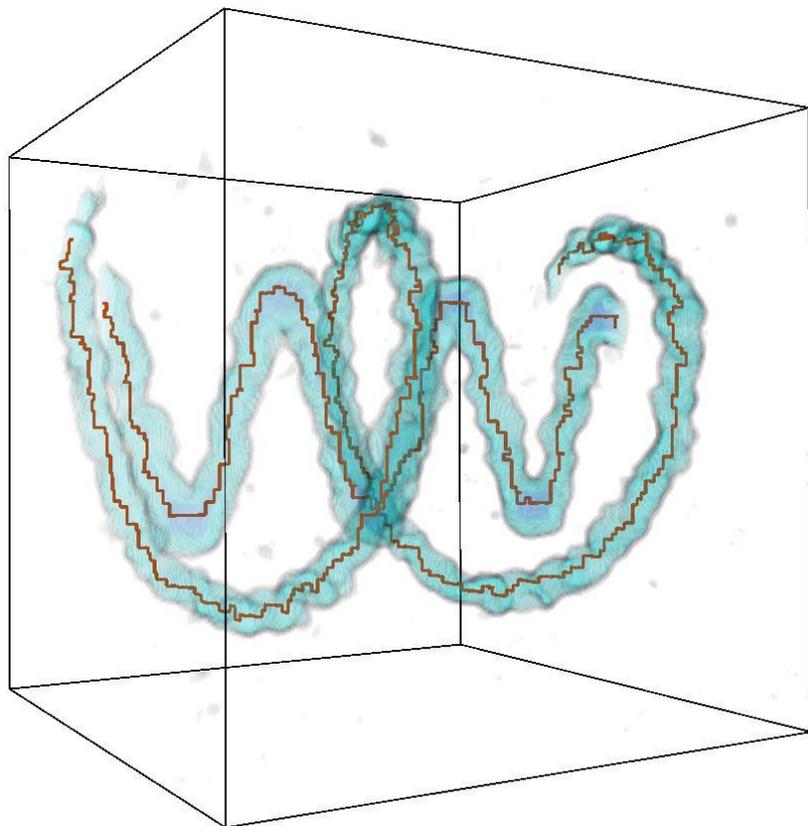
Preparation: we develop control test data to validate the approach



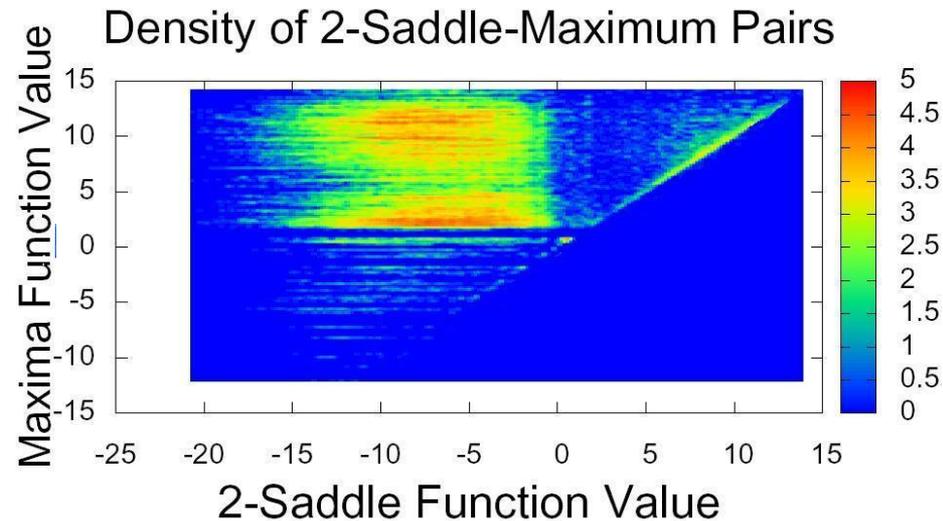
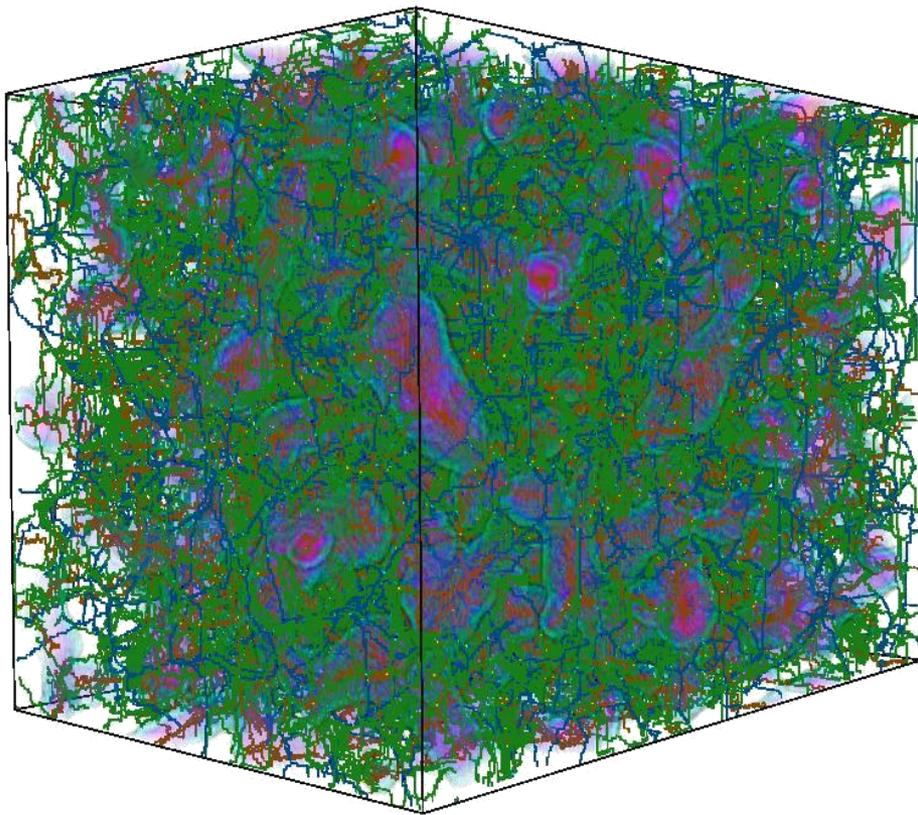
We Report the Distribution of Topological Features in the Full Resolution Data



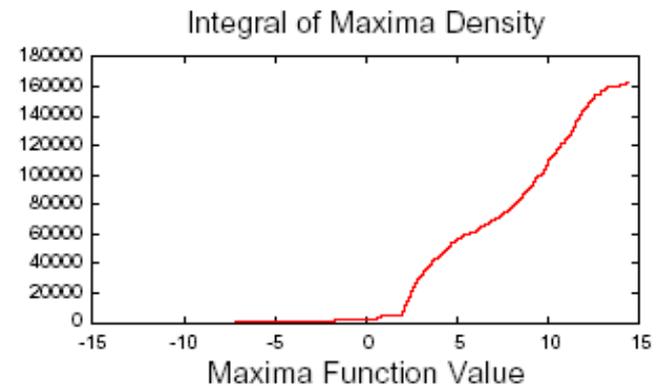
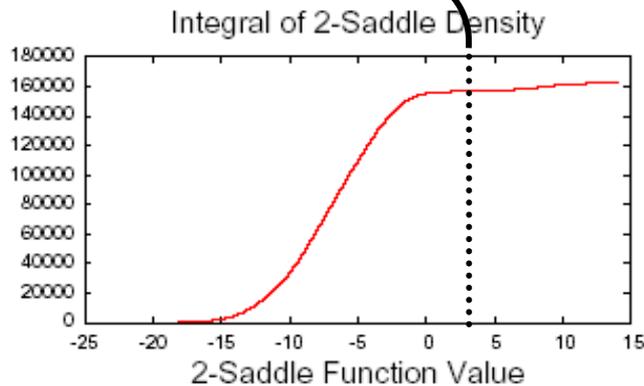
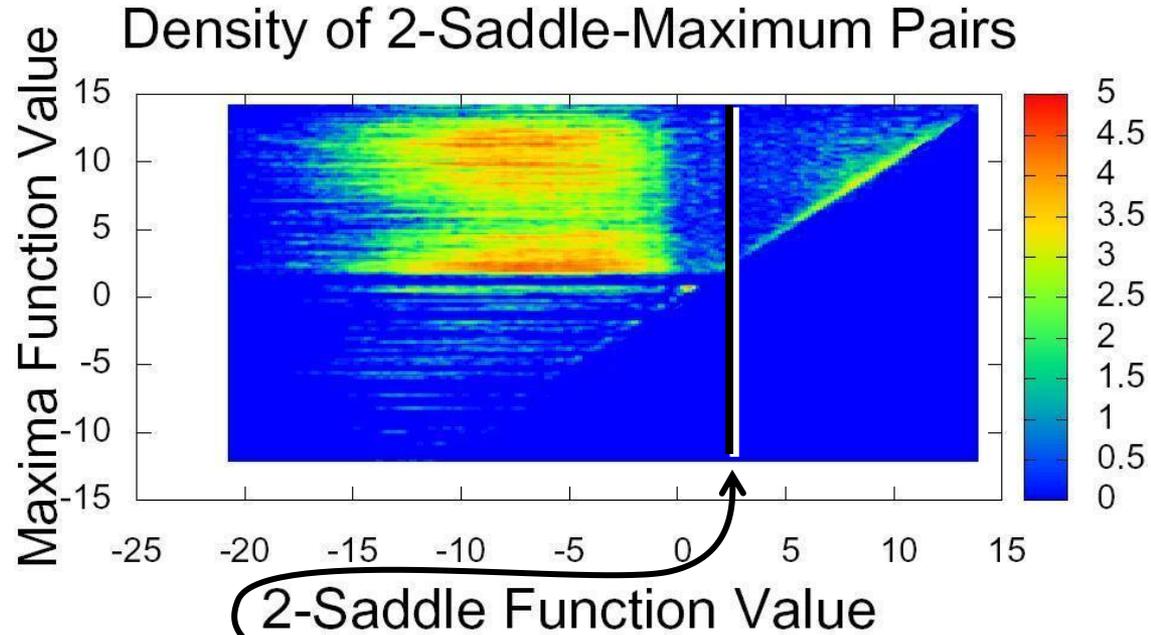
The Hierarchical Morse-Smale Complex Has Very Good Reconstruction Properties



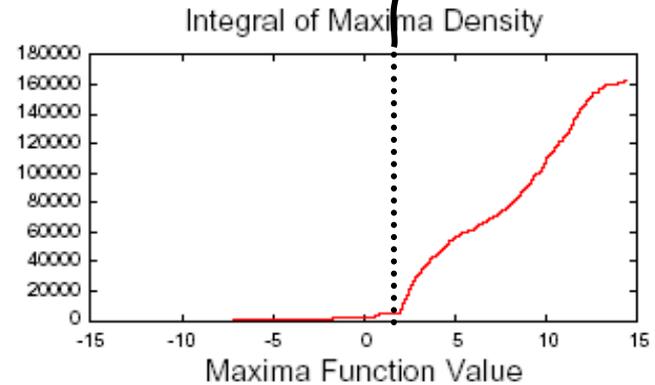
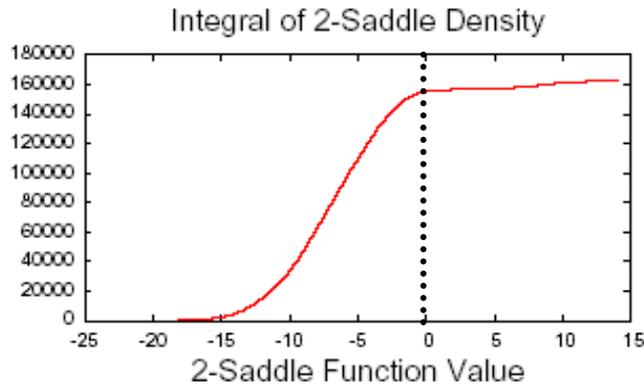
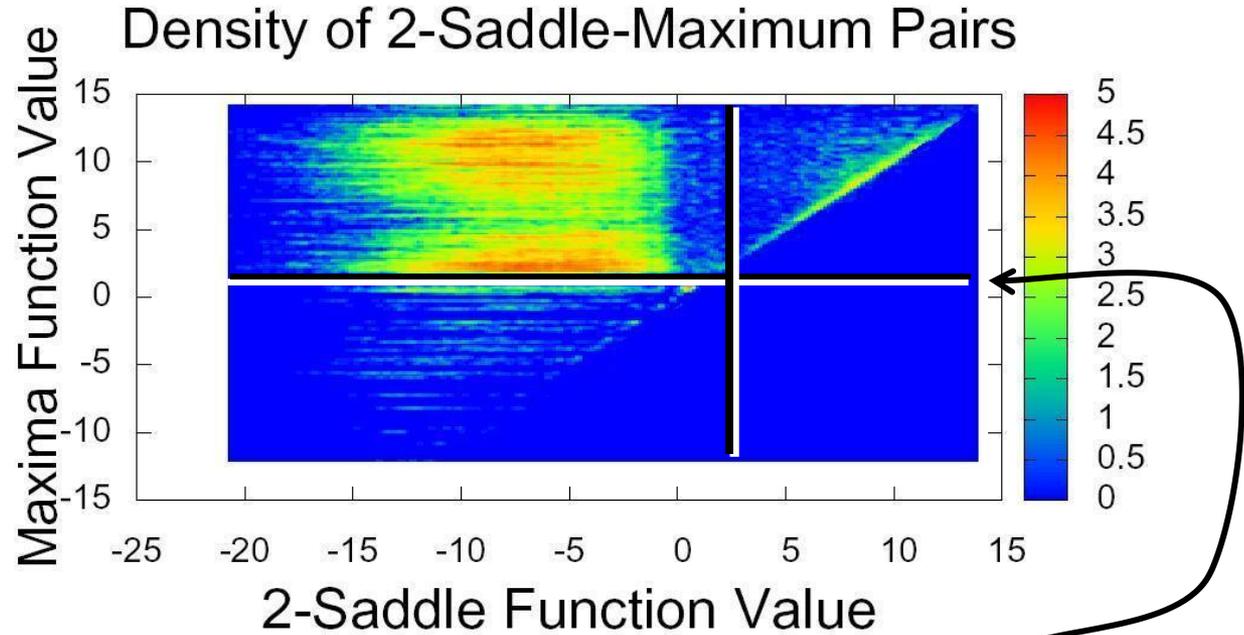
We Compute the Complete Morse-Smale Complex for the Porous Medium



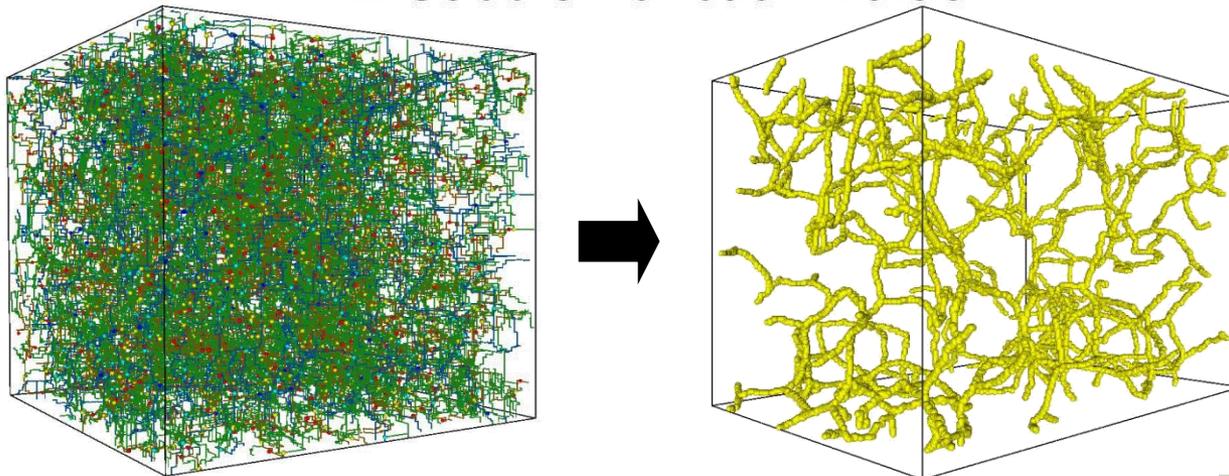
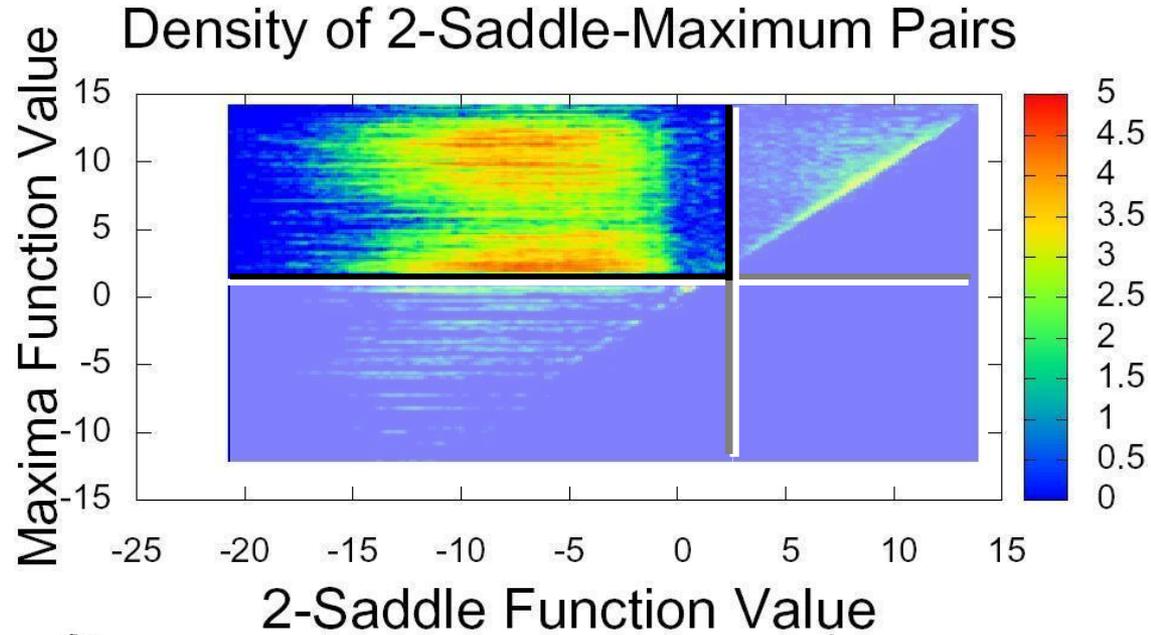
Need to Find Proper Threshold Values and Characterize the Stability of the Solution



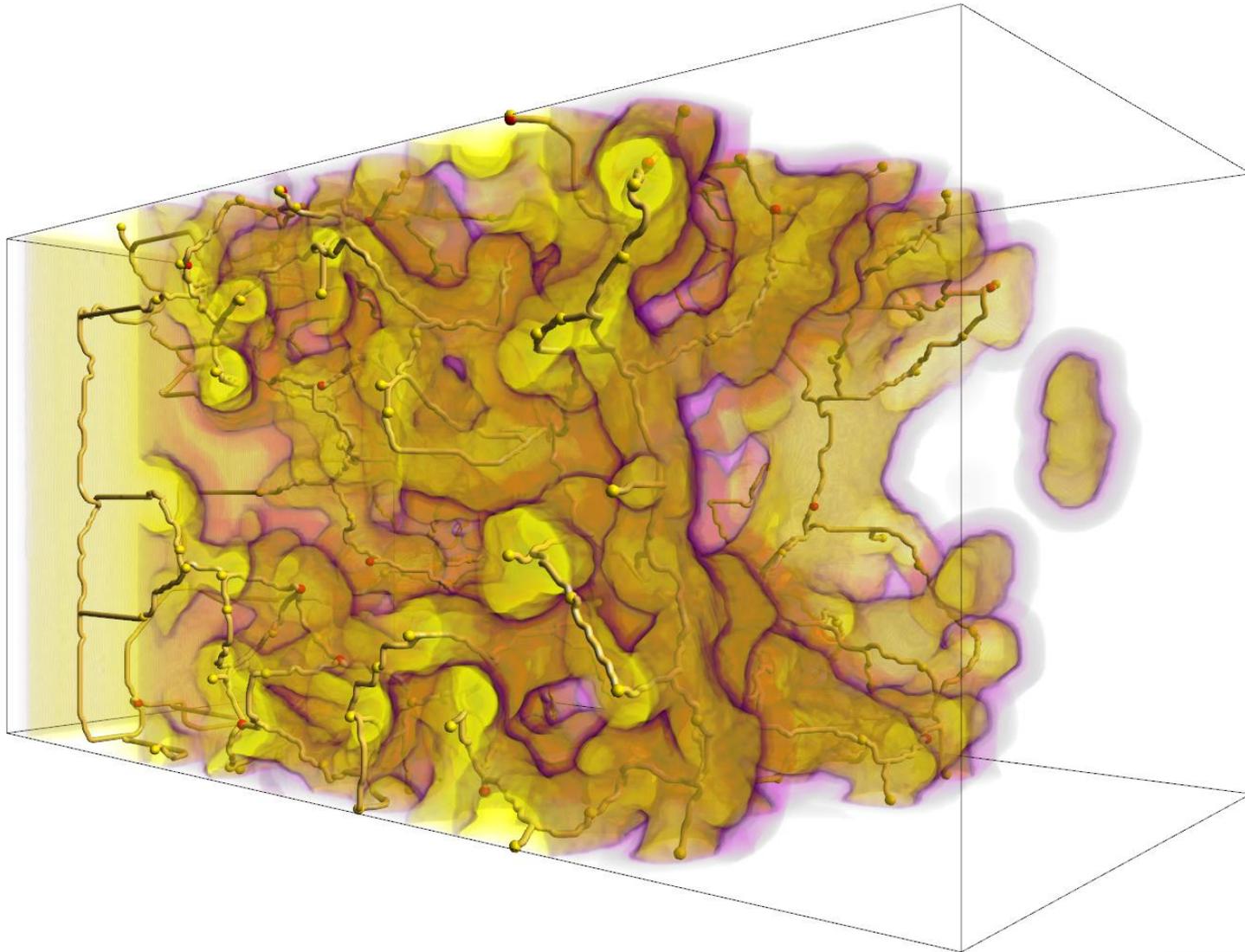
Need to Find Proper Threshold Values and Characterize the Stability of the Solution



We Obtain a Robust Reconstruction of the Filament Structures in the Material

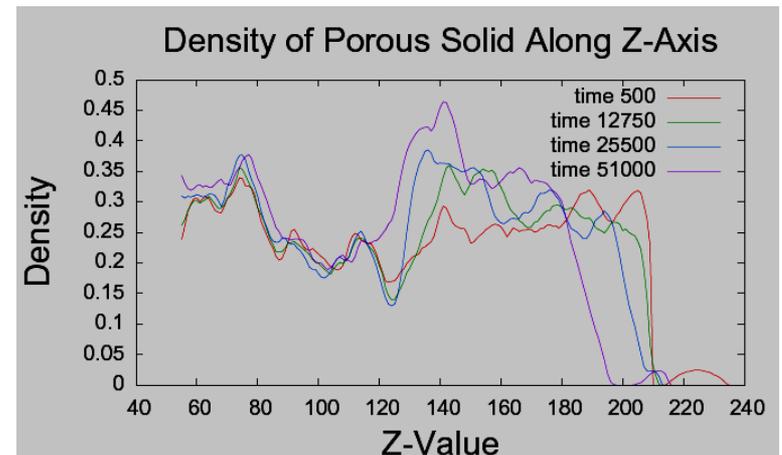
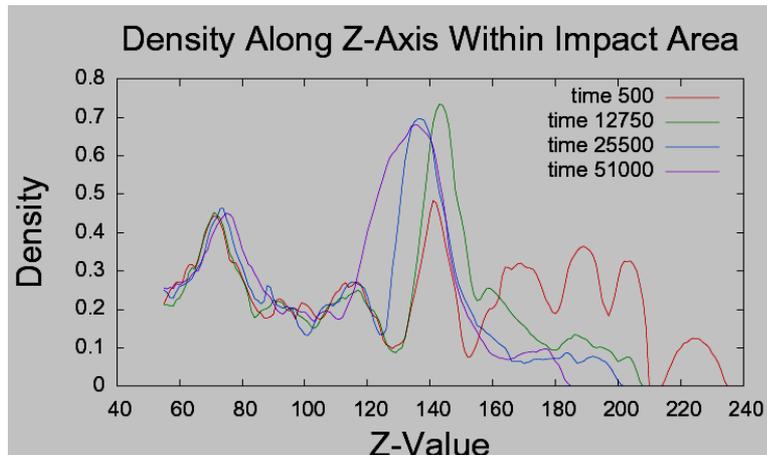


Demo Porous Medium



The Extracted Structures Allow to Quantify the Change in Porosity of the Material

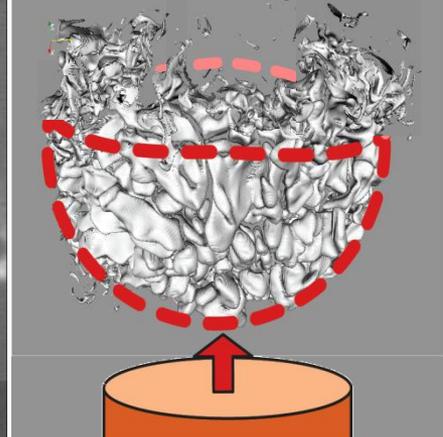
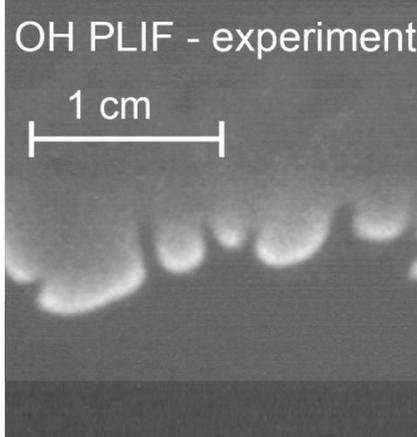
Density profiles



Decay in porosity of the material

Metric	t=500	t=12750	t=25500	t=51000
# Cycles	762	340	372	256
Total Length	34756	24316	23798	18912

Understanding Turbulence for Low Emission, High Efficiency Combustion



Experiment

Simulation

- Lean premixed H₂ flames
- Low Swirl Combustion (LSC) Burners
- Low pollution in energy production
- High Efficiency in fuel consumption
- Scalable from residential to industrial use
- Each variable 3.9-4.5 TB

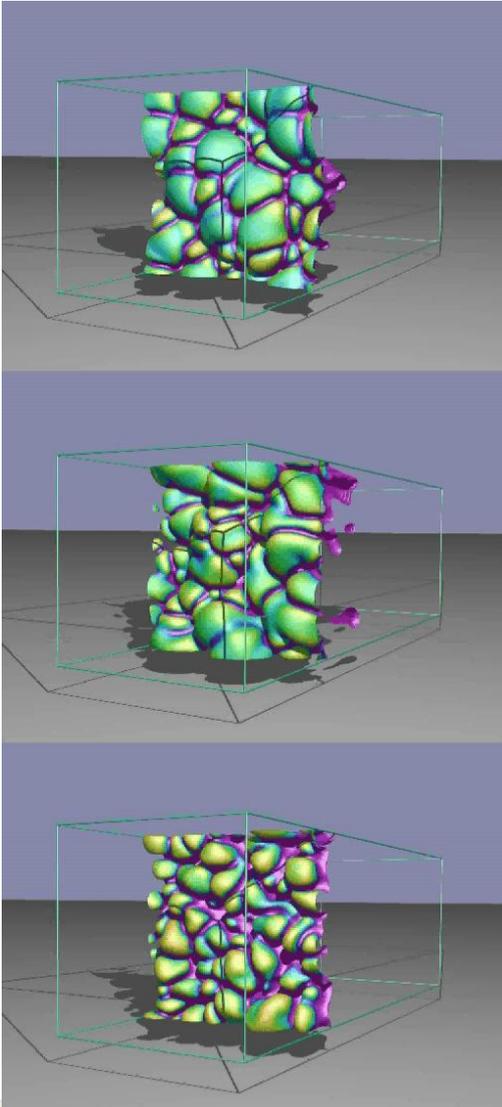


1" burner (5 kW, 17 KBtu/hr)



28" burner (44 MW, 150 MBtu/hr)

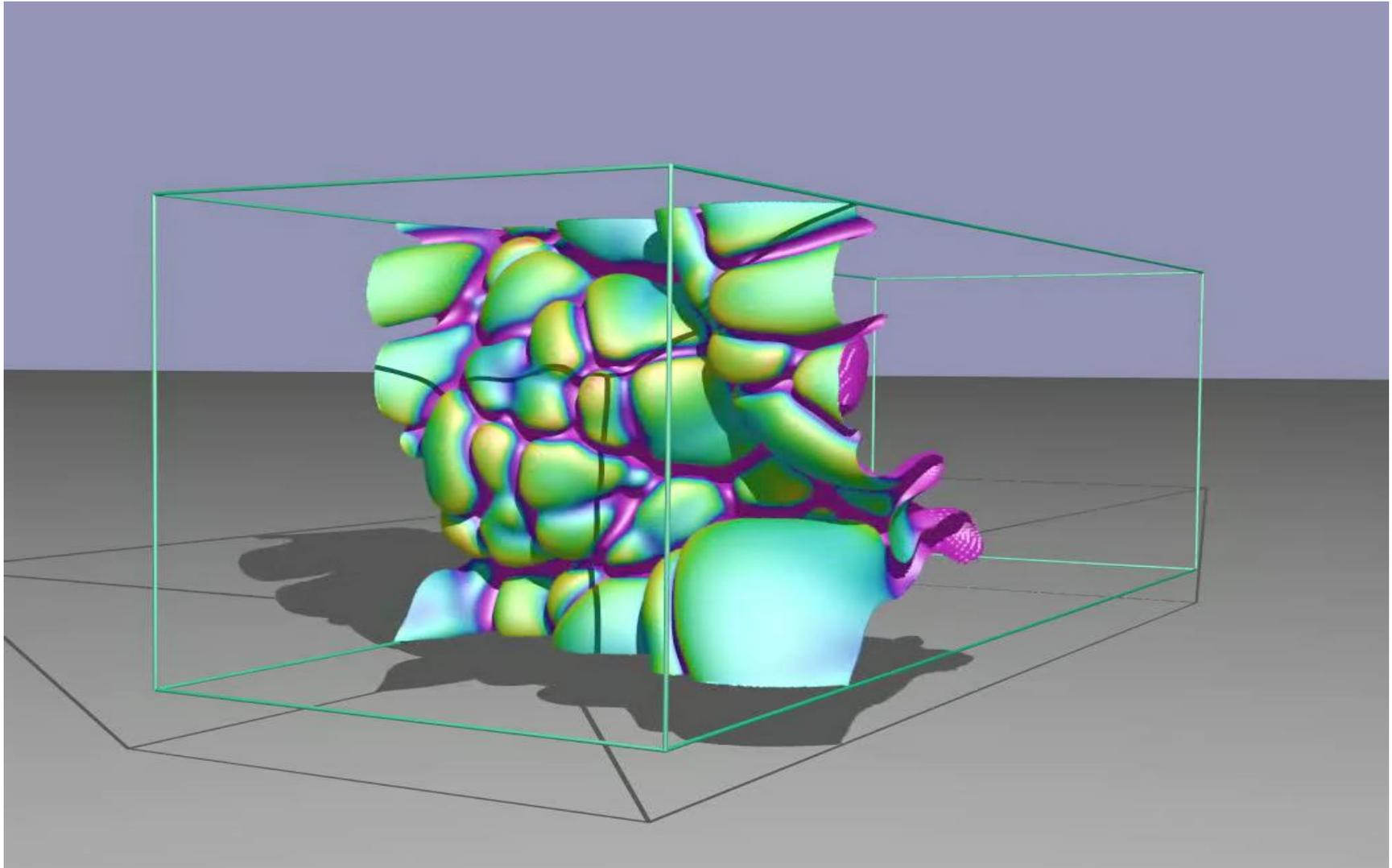
We Take on the Challenge of Developing a Quantitative Analysis Detecting Turbulence



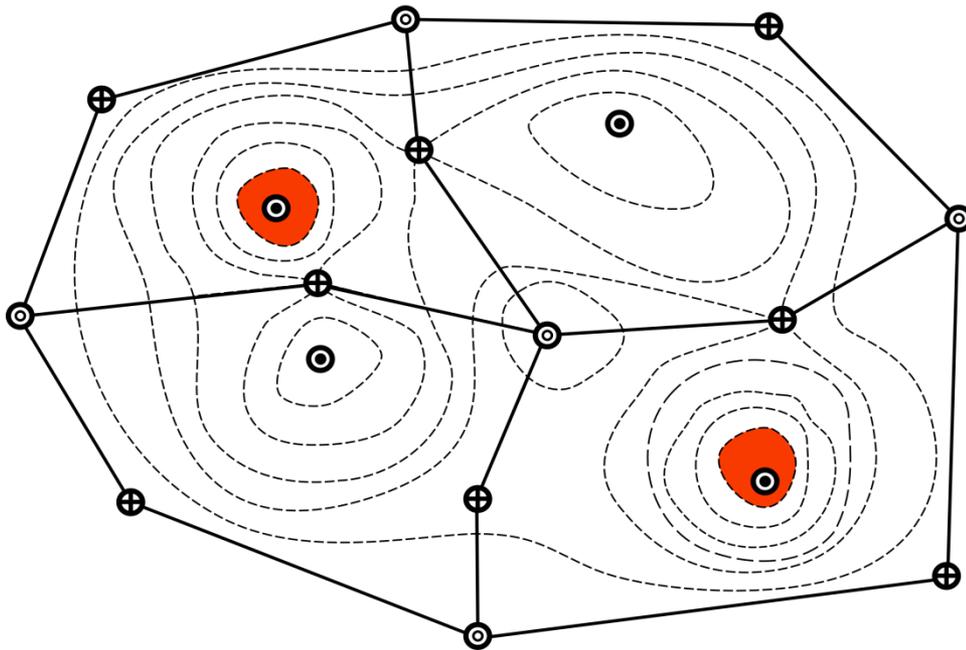
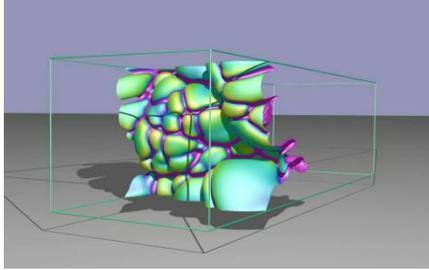
Understanding combustion processes over a broad range of burning conditions is an important problem for designing engines and power plants.

- Simulation with AMR mesh.
- Simulations of lean premixed hydrogen flames with three degrees of turbulence.
- Can we identify precisely and track in time burning regions?
- Can we discriminate the degree of turbulence from a quantitative analysis?

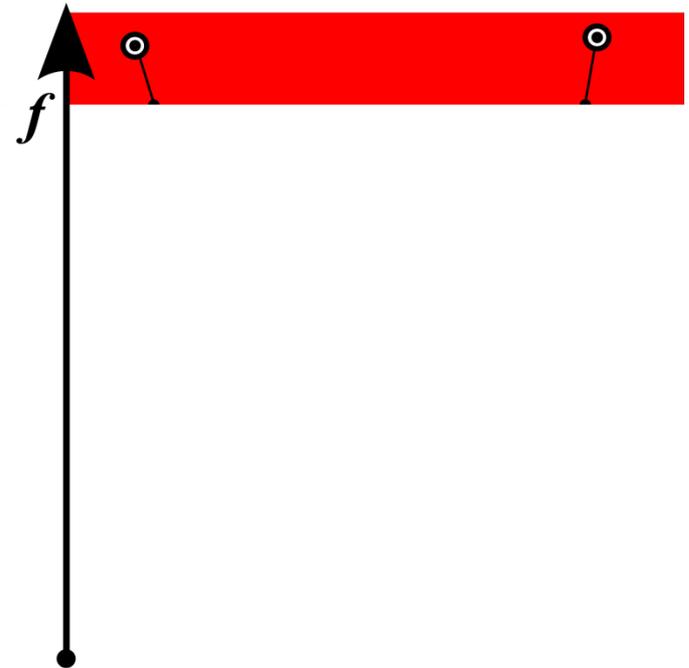
We Build a Reduced Topological Model of H₂ Consumption on an Isothermal Surfaces



We Build a Reduced Topological Model of H2 Consumption on an Isothermal Surfaces

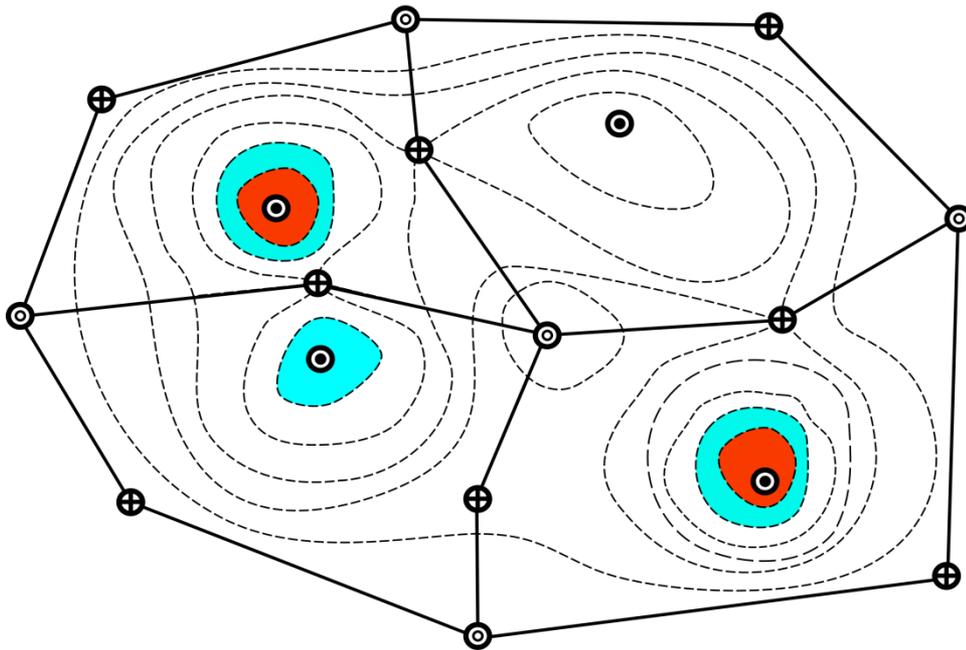
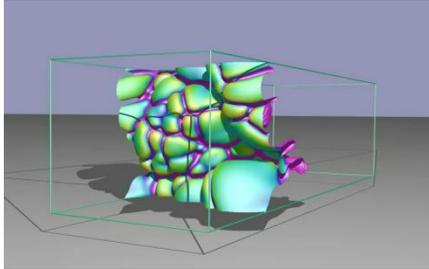


domain

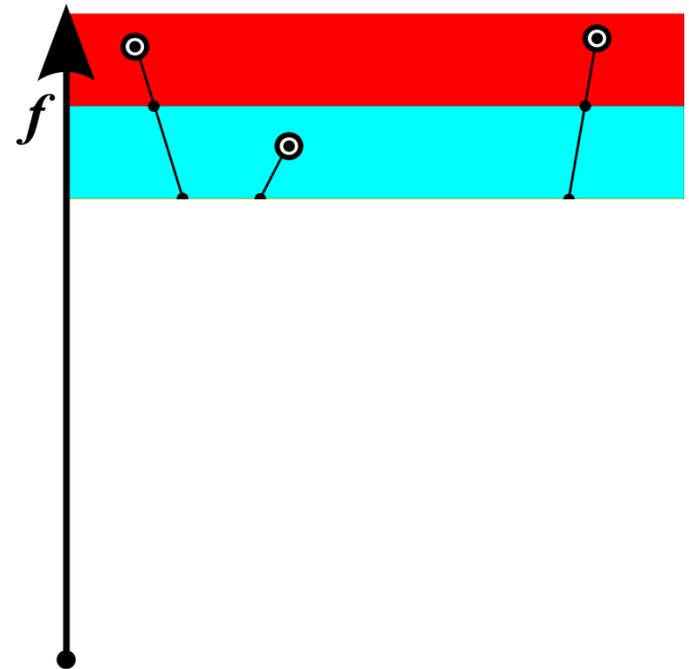


hierarchy

We Build a Reduced Topological Model of H2 Consumption on an Isothermal Surfaces

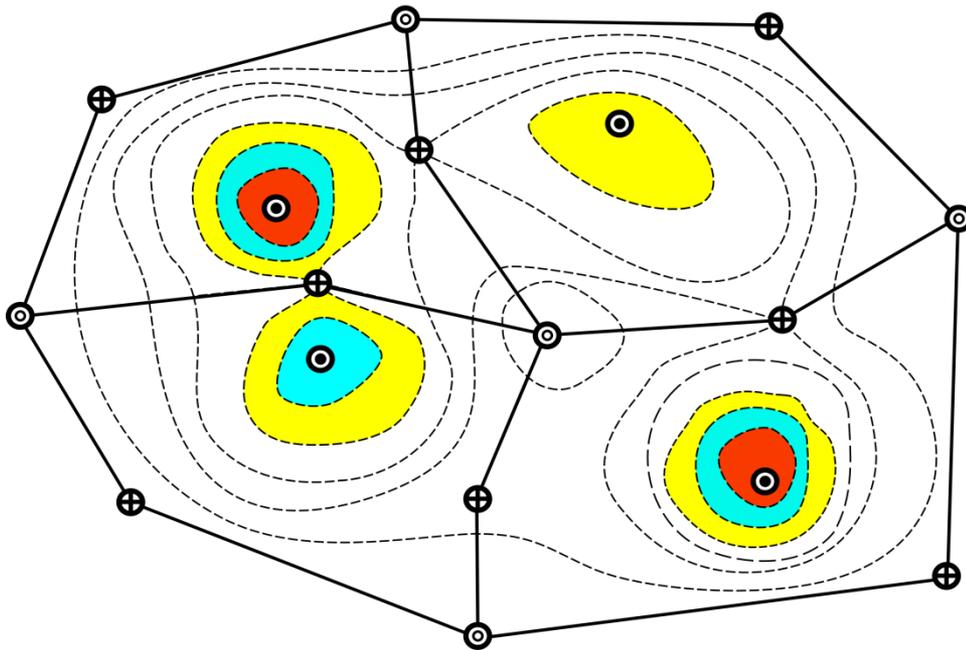
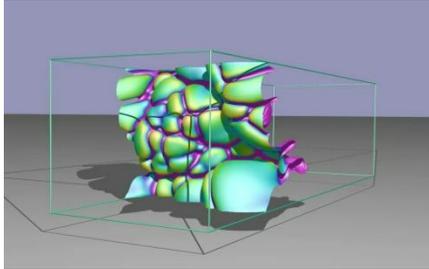


domain

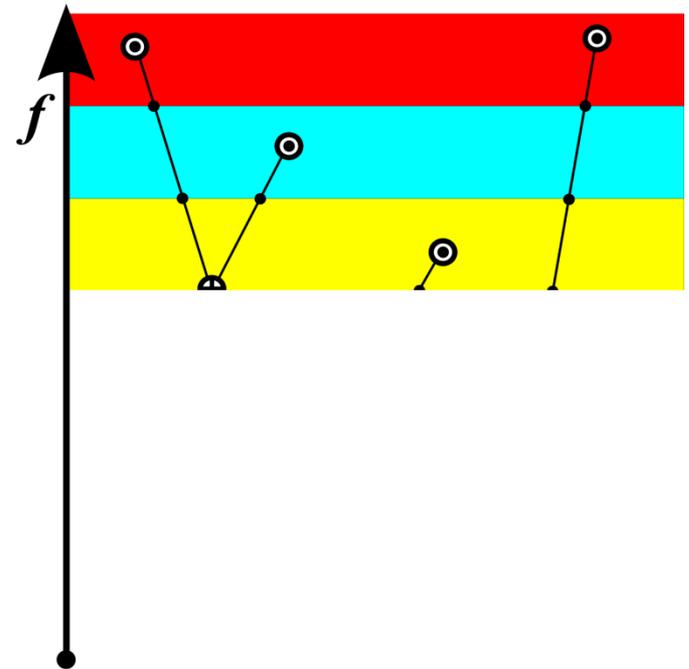


hierarchy

We Build a Reduced Topological Model of H2 Consumption on an Isothermal Surfaces

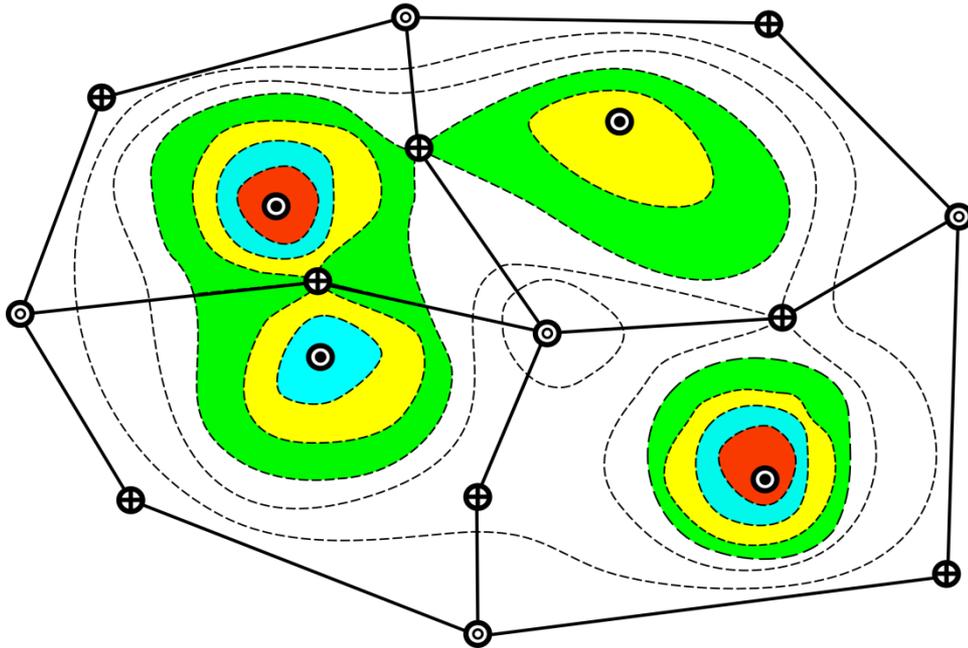
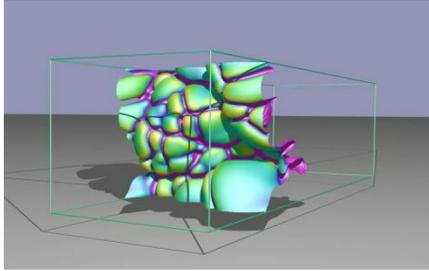


domain

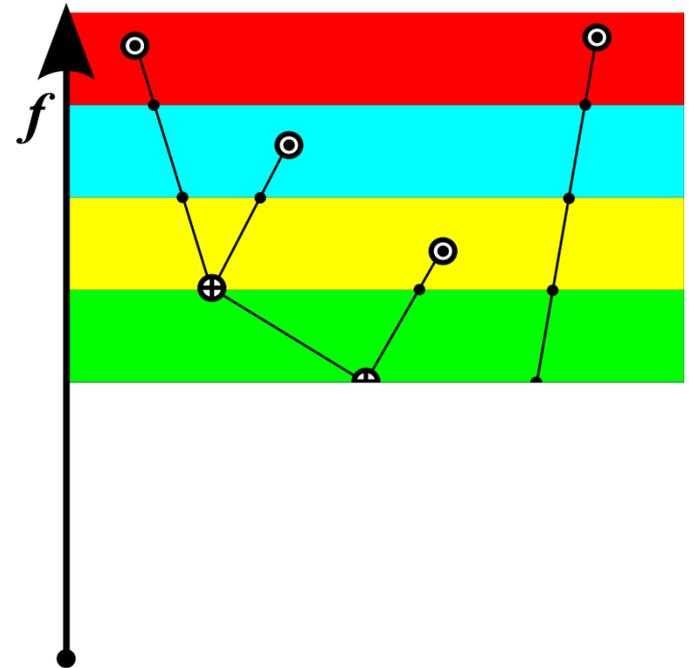


hierarchy

We Build a Reduced Topological Model of H2 Consumption on an Isothermal Surfaces

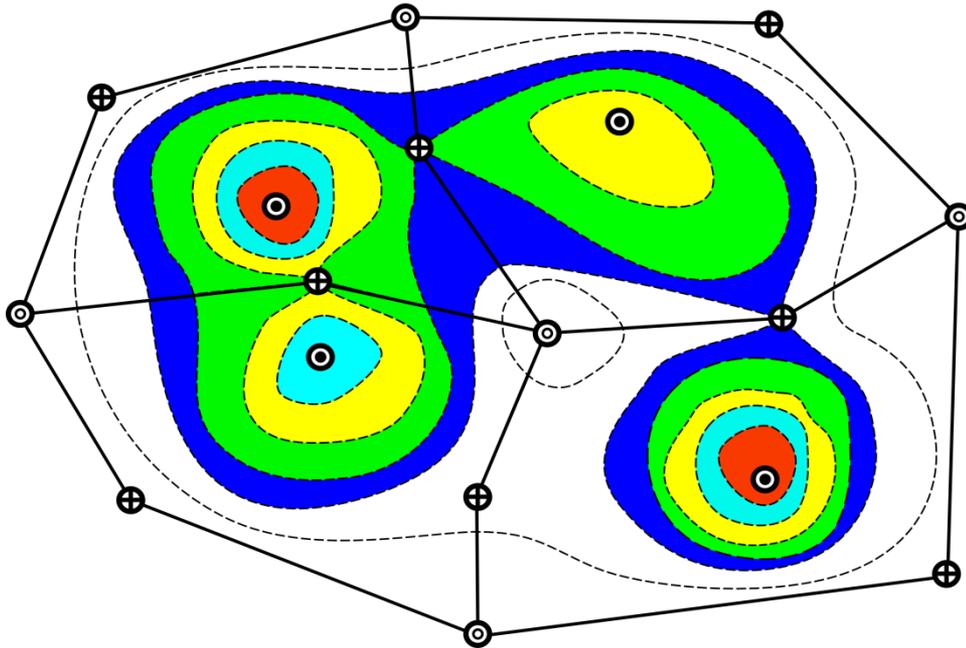
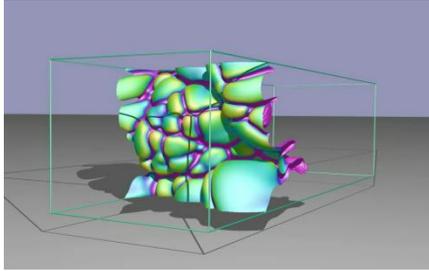


domain

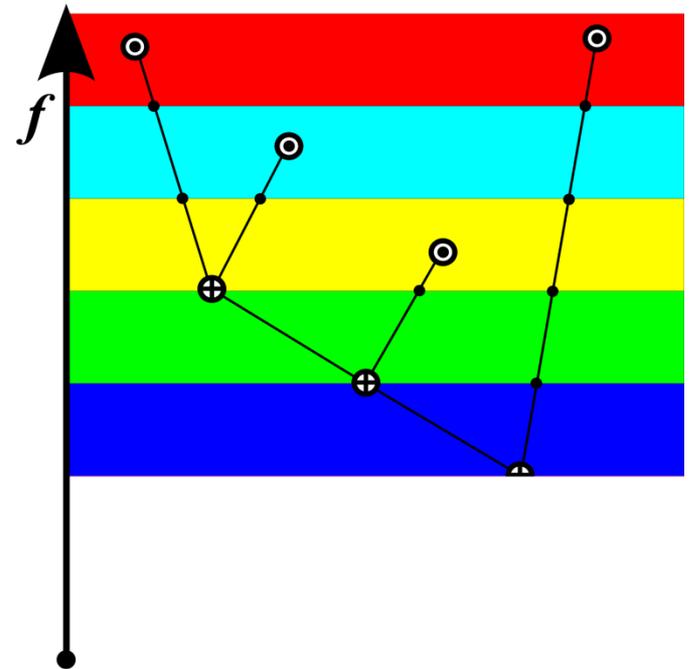


hierarchy

We Build a Reduced Topological Model of H2 Consumption on an Isothermal Surfaces

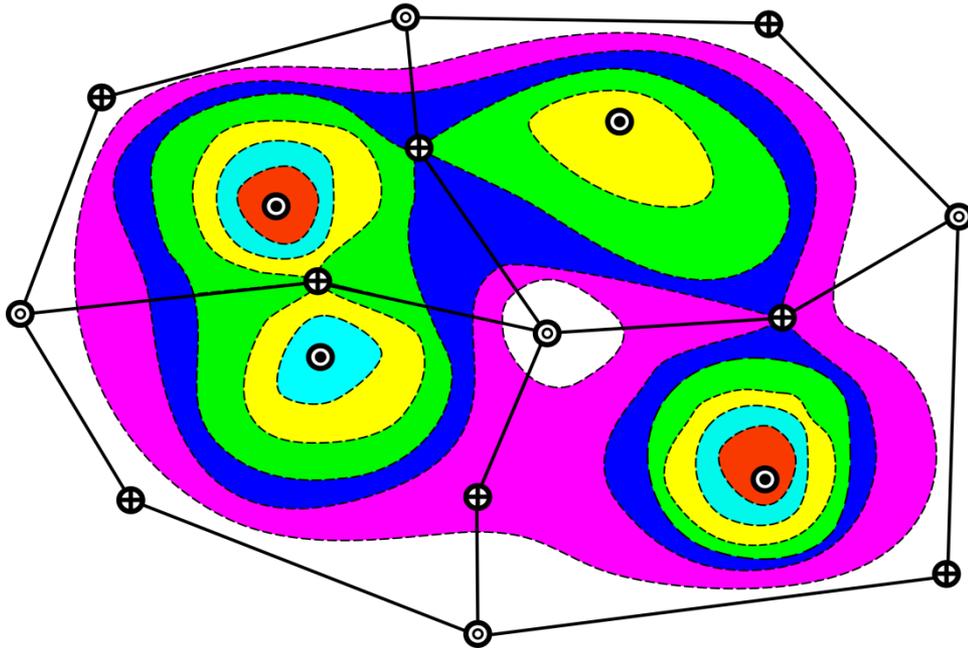
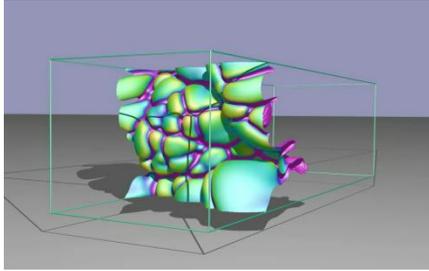


domain

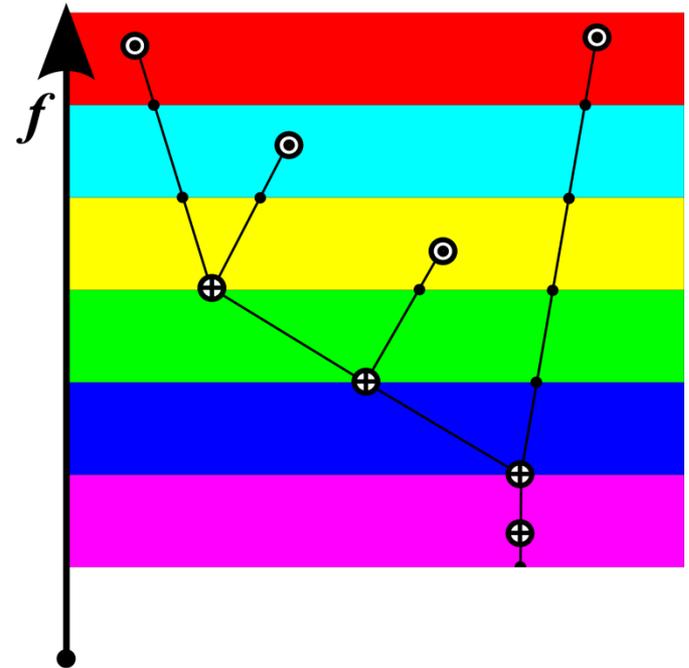


hierarchy

We Build a Reduced Topological Model of H2 Consumption on an Isothermal Surfaces

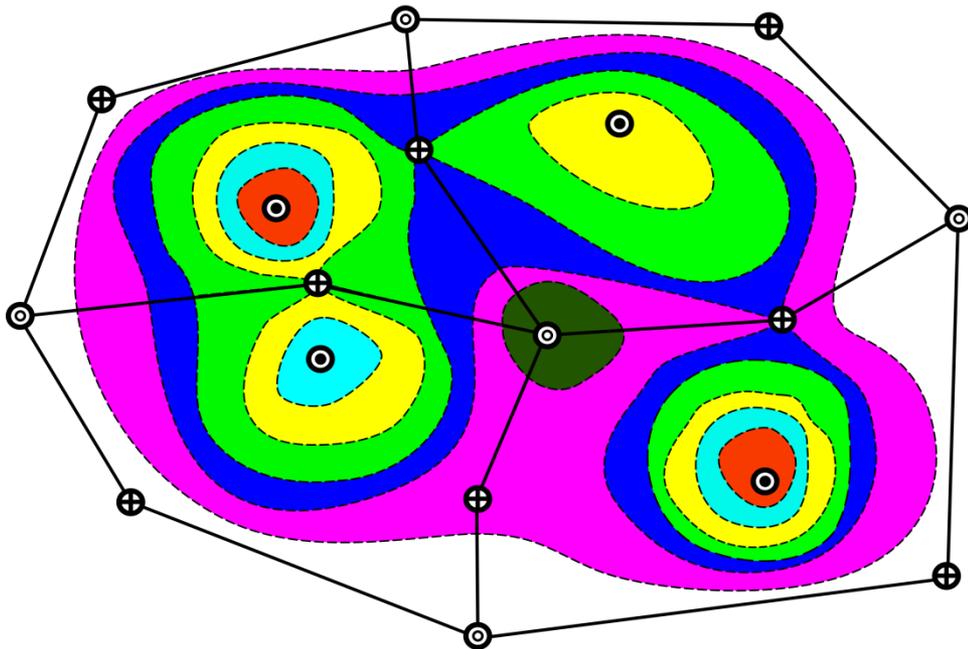
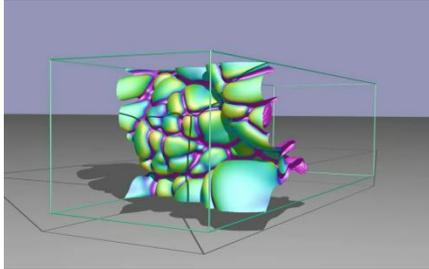


domain

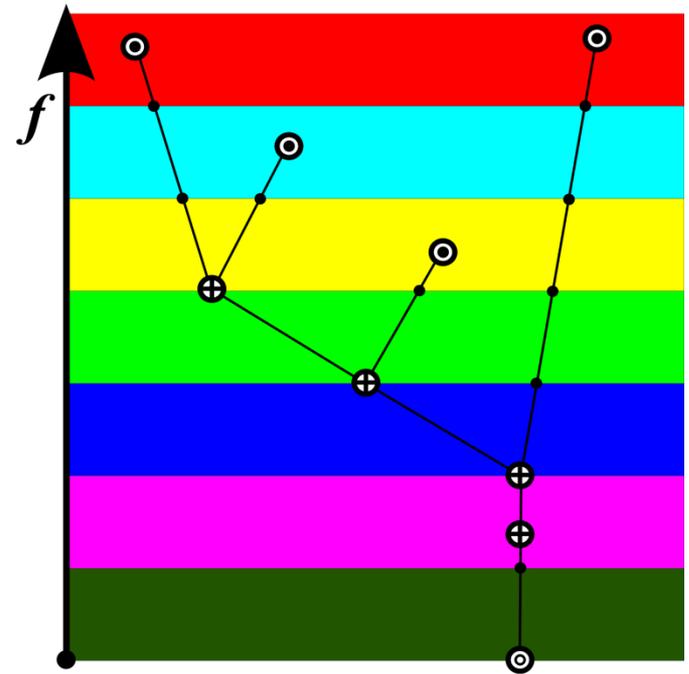


hierarchy

We Build a Reduced Topological Model of H2 Consumption on an Isothermal Surfaces

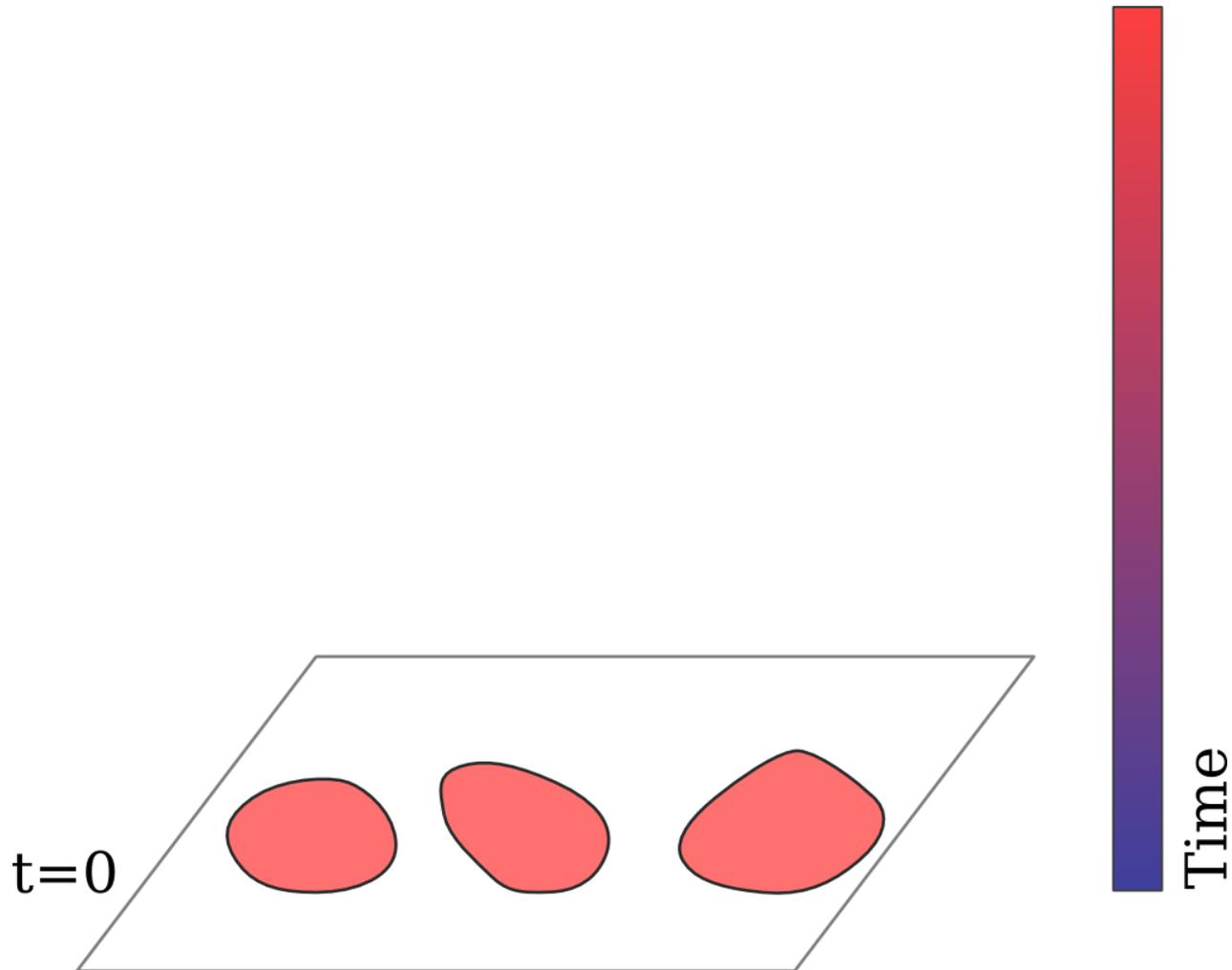


domain

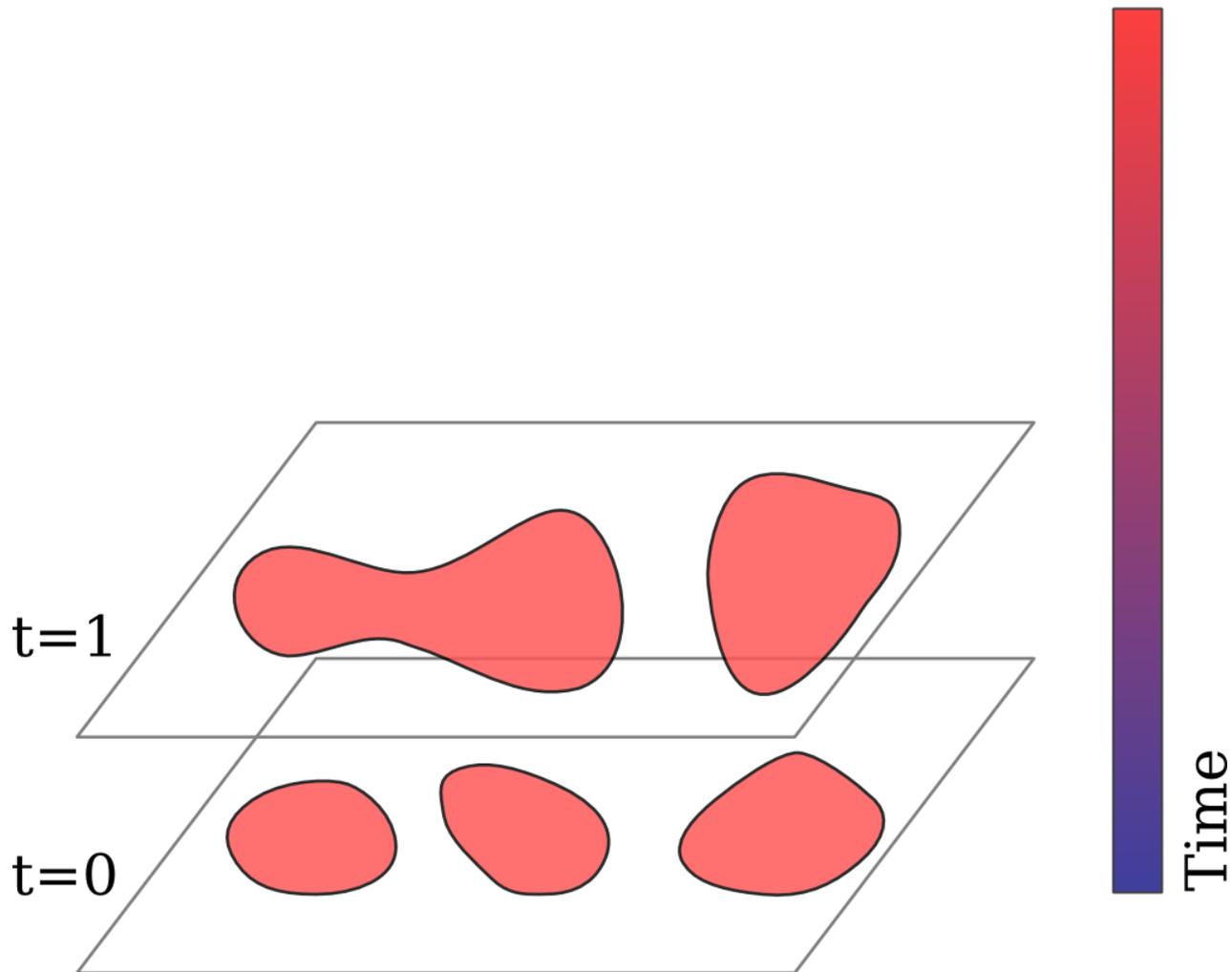


hierarchy

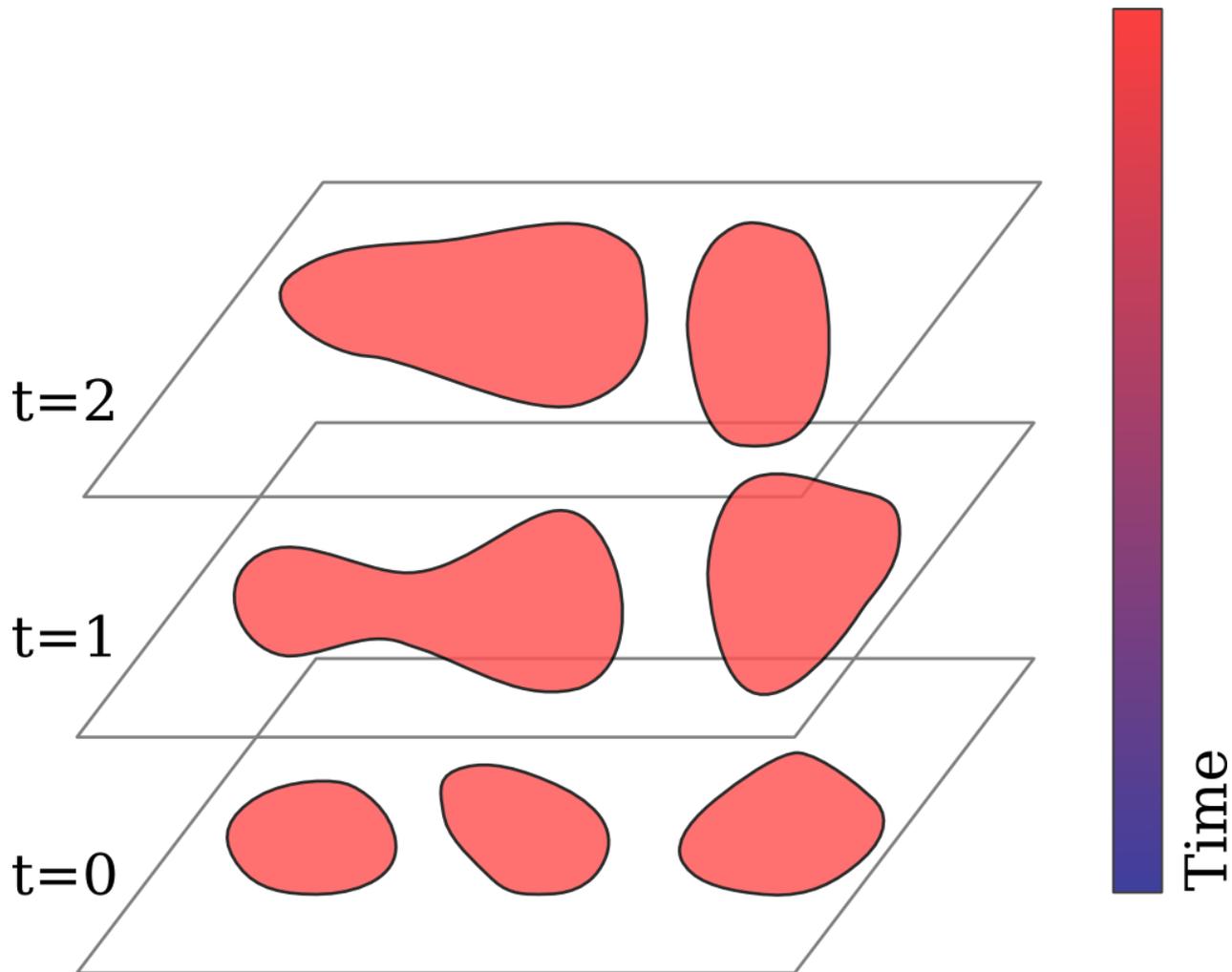
We track in time by interpolation in 4D and contraction to a Reeb Graph



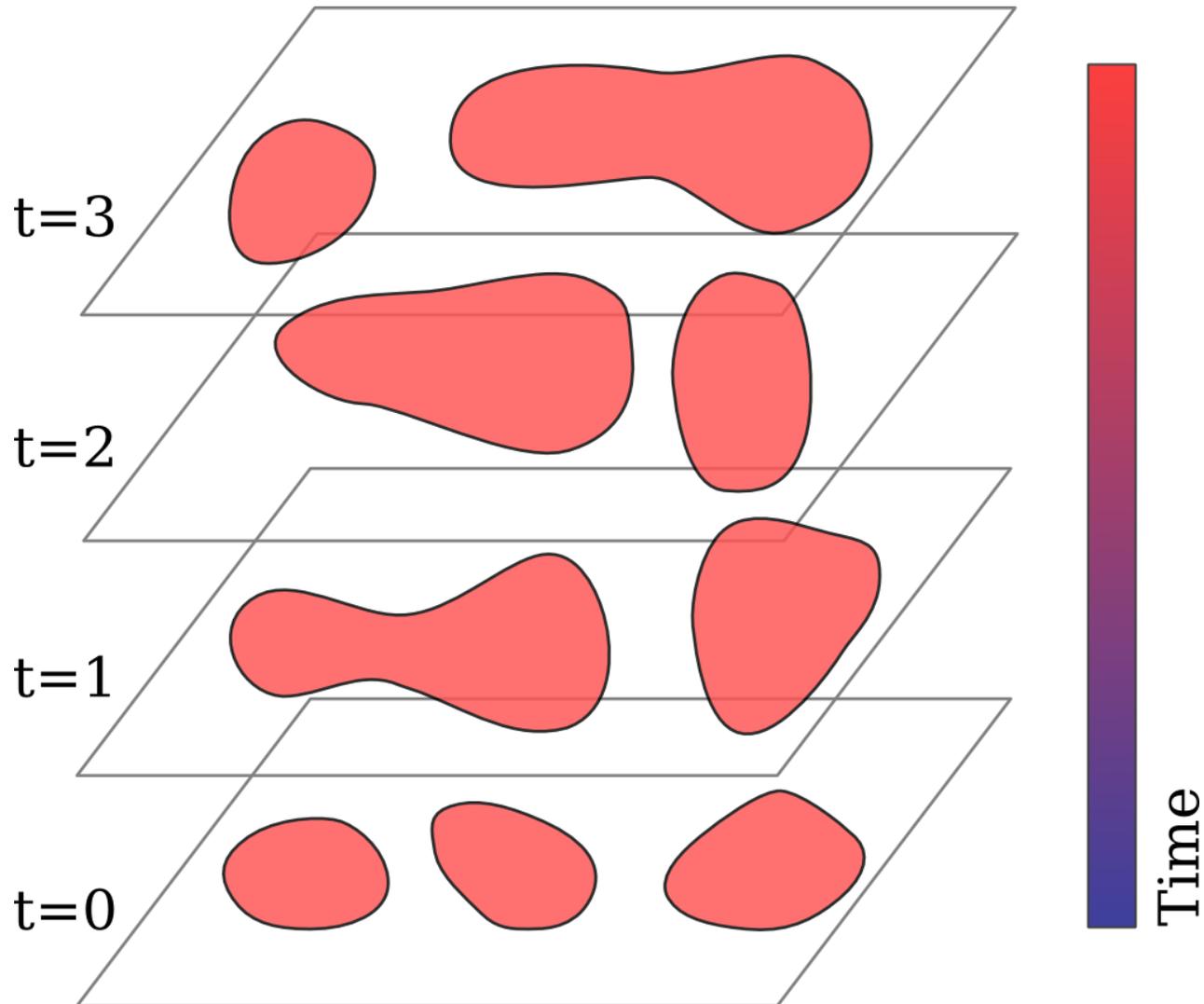
We track in time by interpolation in 4D and contraction to a Reeb Graph



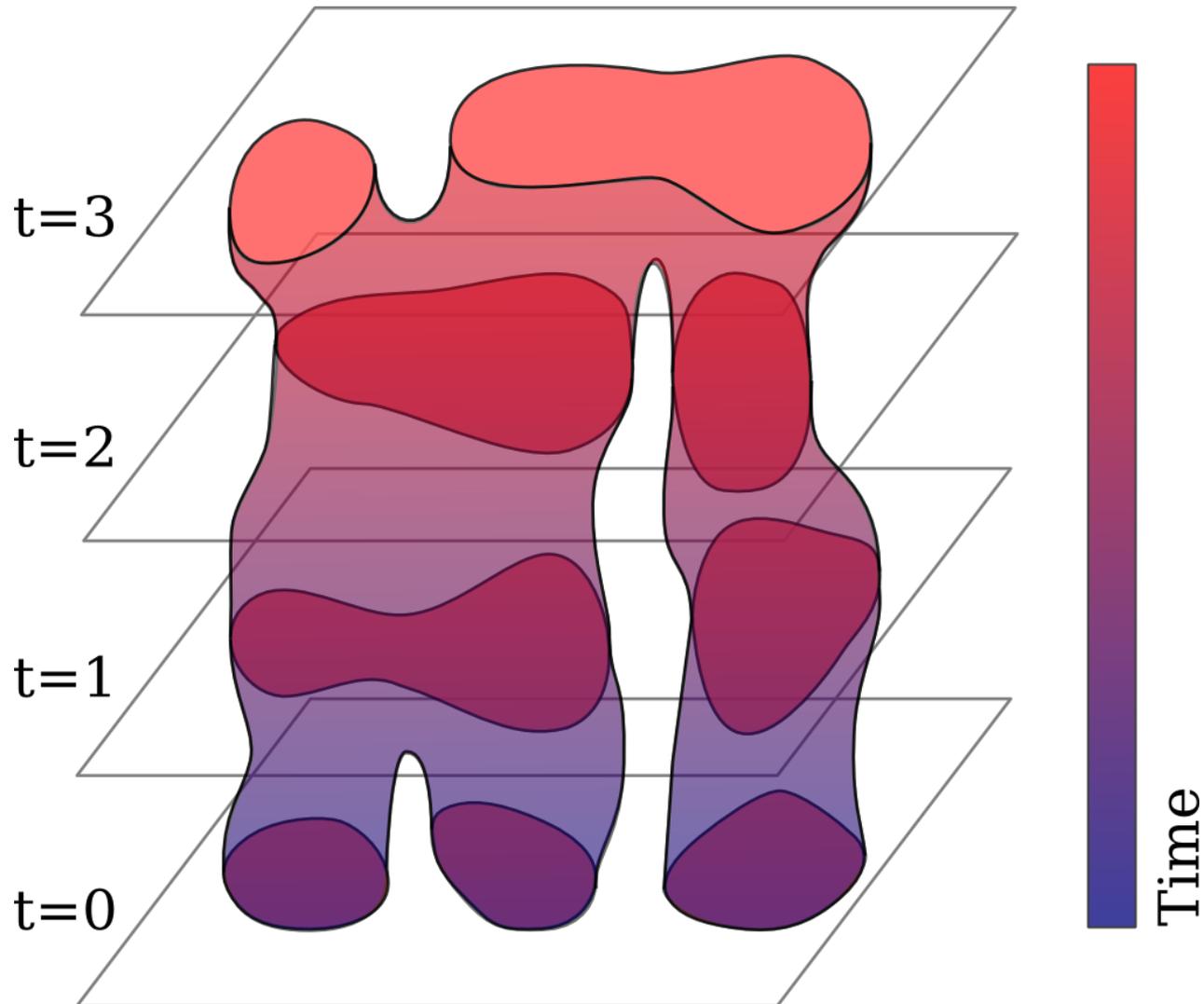
We track in time by interpolation in 4D and contraction to a Reeb Graph



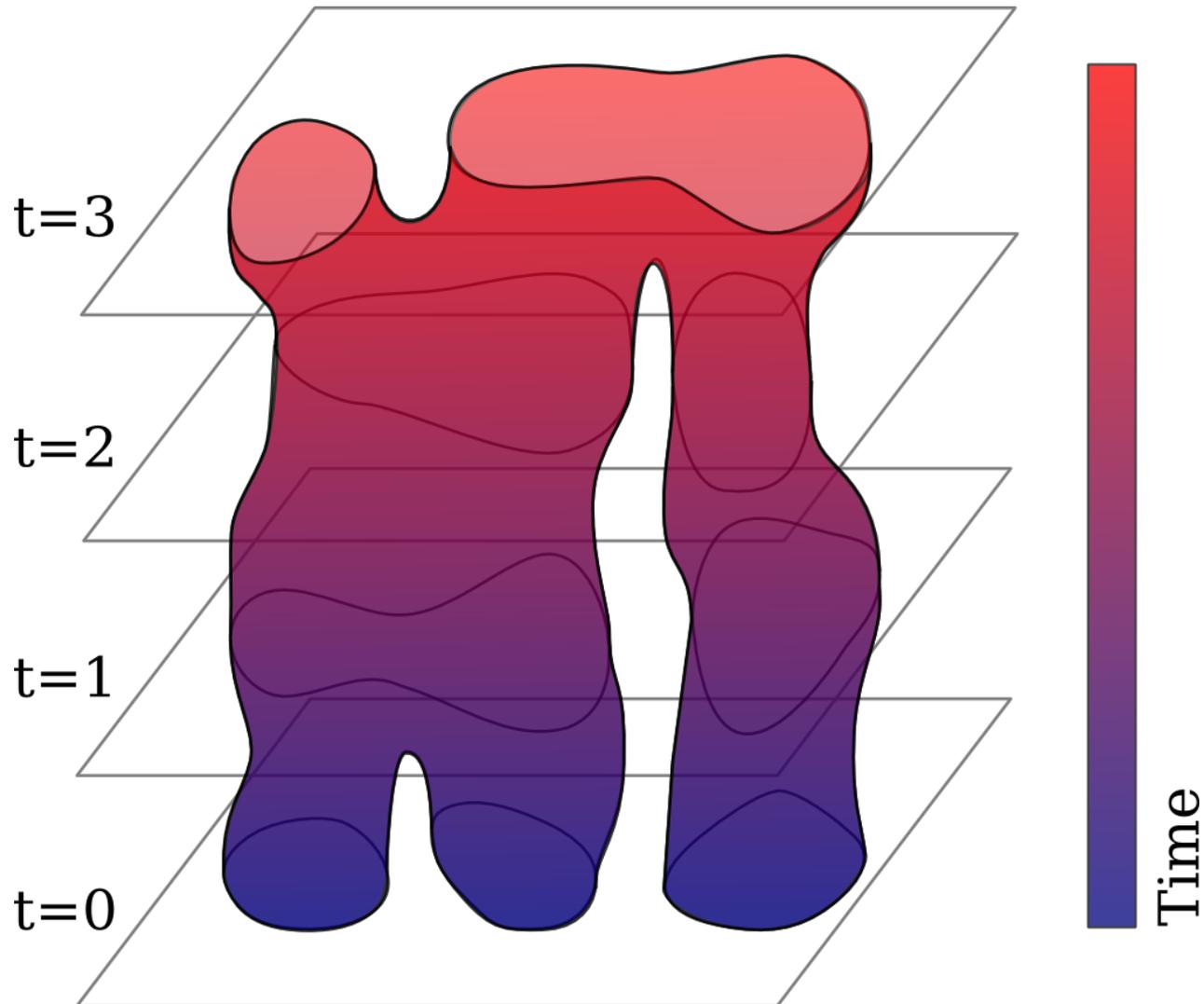
We track in time by interpolation in 4D and contraction to a Reeb Graph



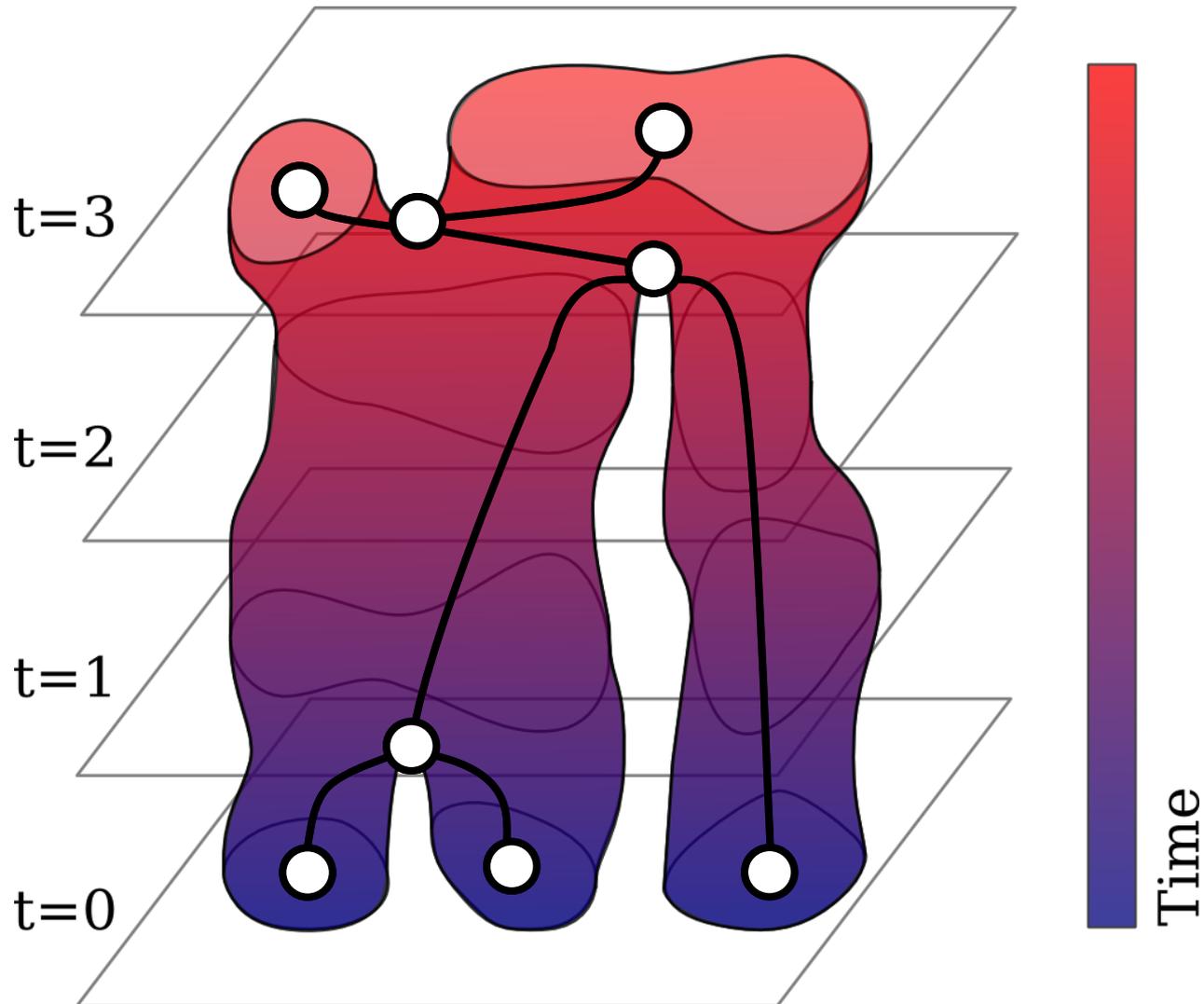
We track in time by interpolation in 4D and contraction to a Reeb Graph



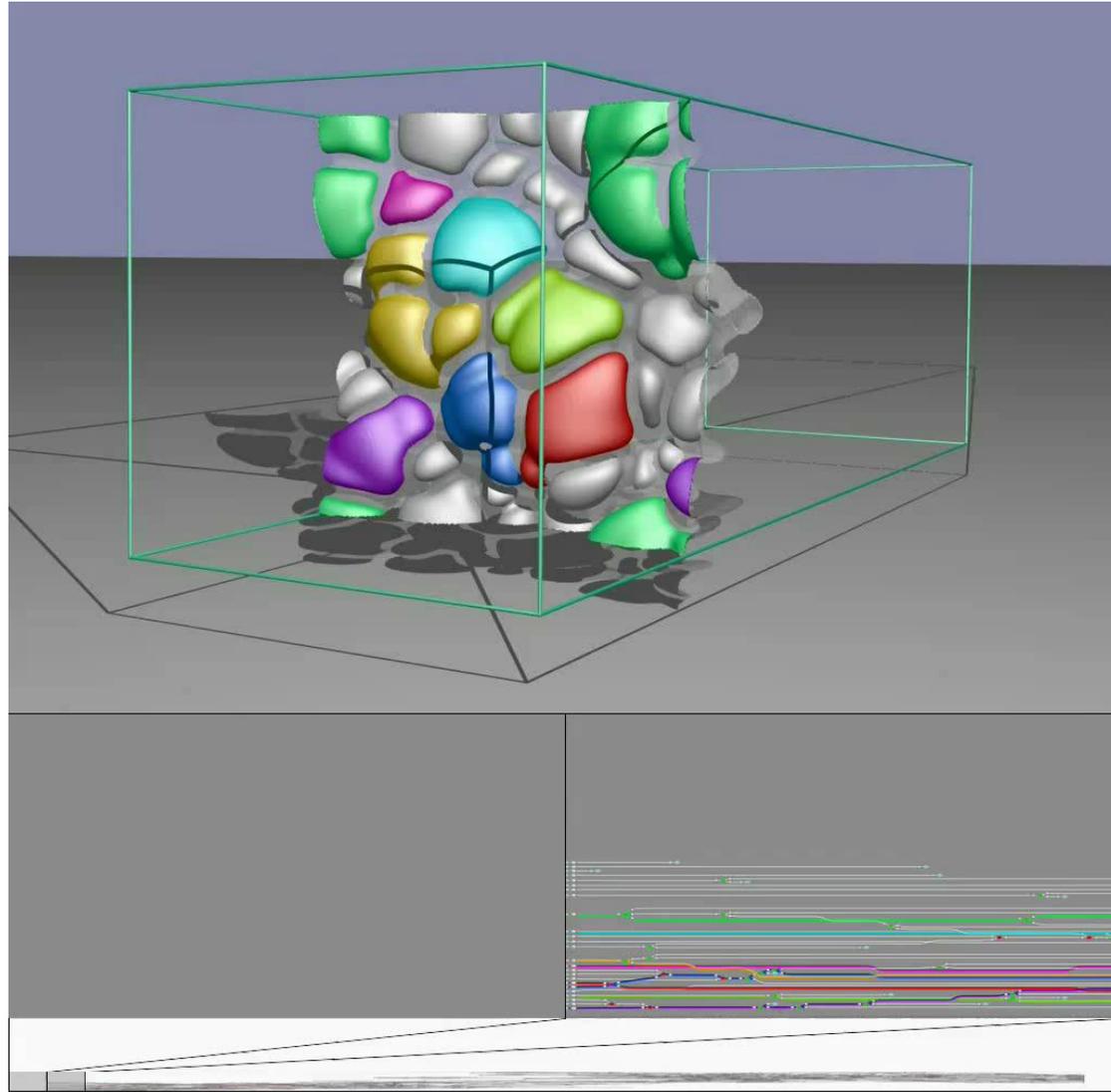
We track in time by interpolation in 4D and contraction to a Reeb Graph



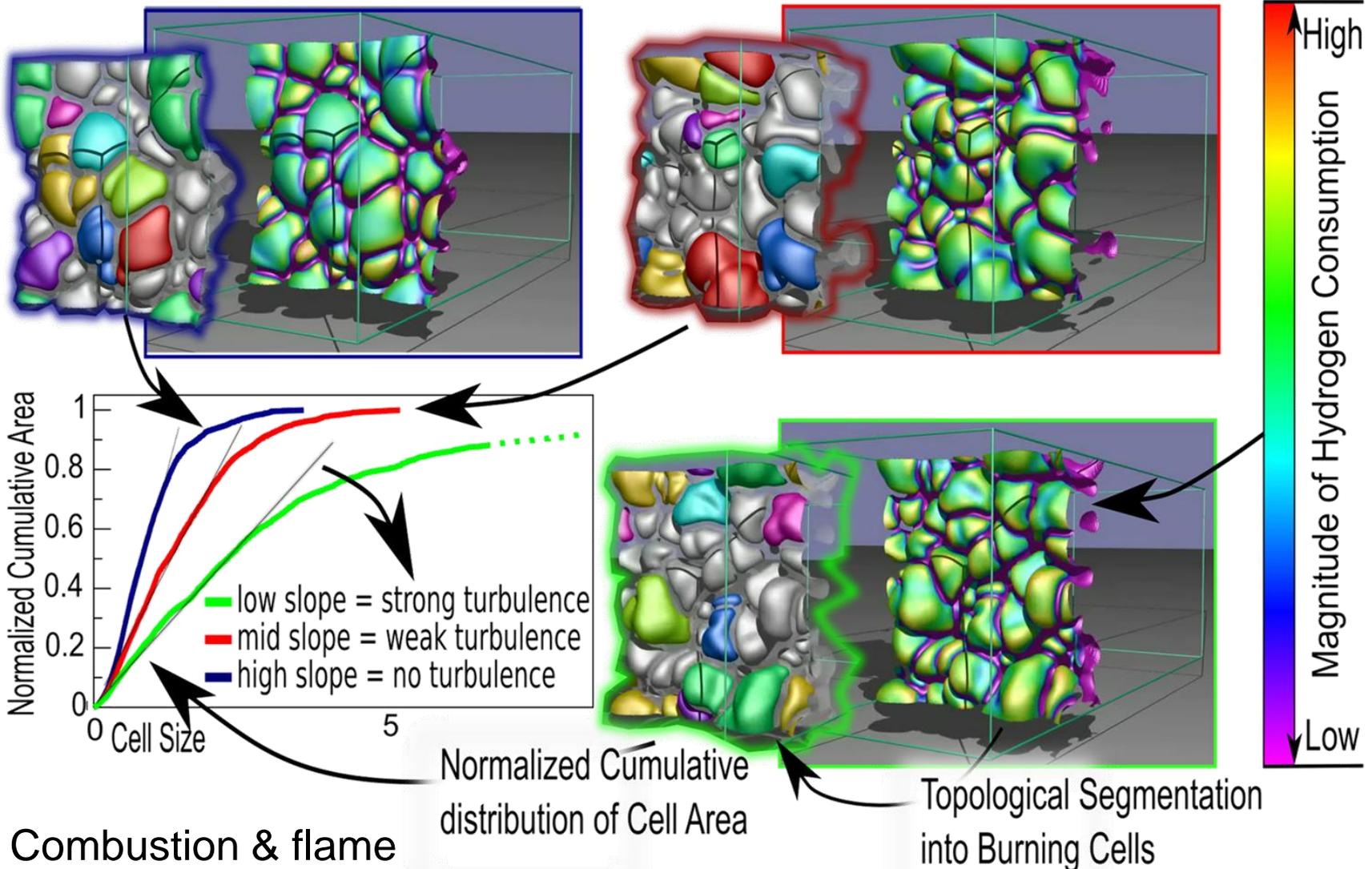
We track in time by interpolation in 4D and contraction to a Reeb Graph



Each Set of Parameters Results in a Robust Segmentation and Tracking of Burning Cells



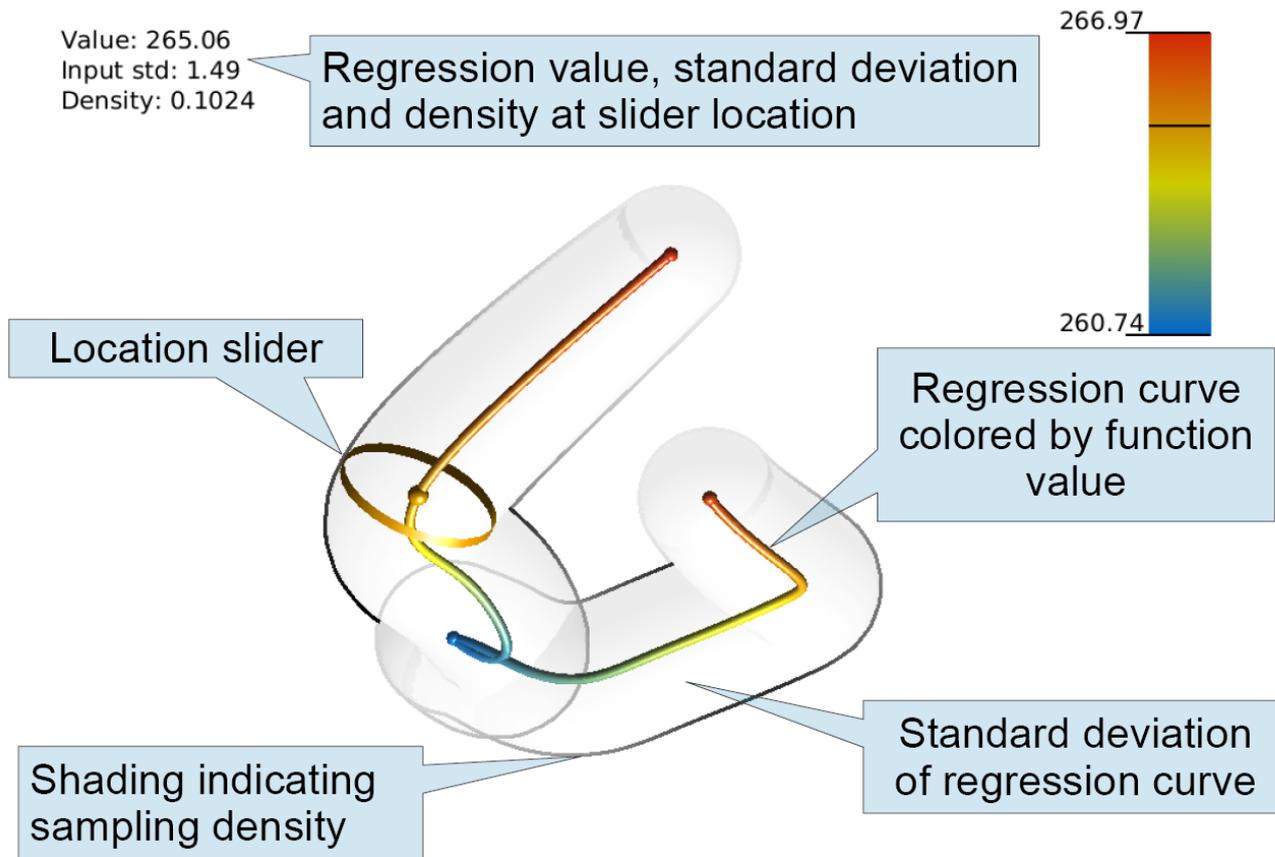
Topological Segmentation Allows to Quantify Turbulence as Slope of the Area Distributions



- Combustion & flame

Exploration of High Dimensional Functions for Sensitivity Analysis

Integrated presentation of statistics and topology



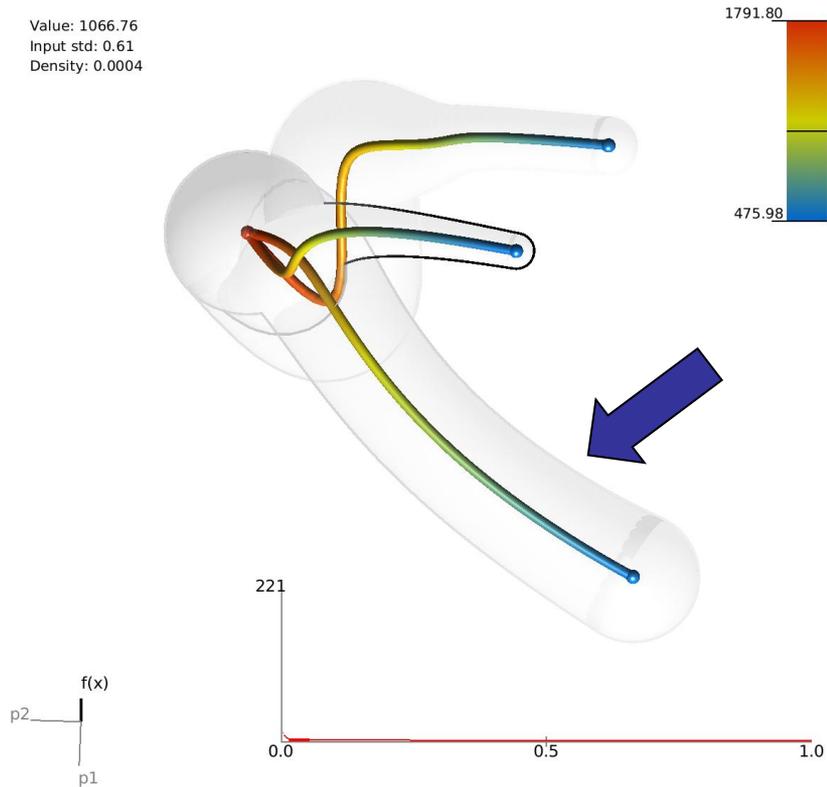
Analysis of Combustion Simulations

Combustion Simulation of Jet CO/H₂-Air Flames

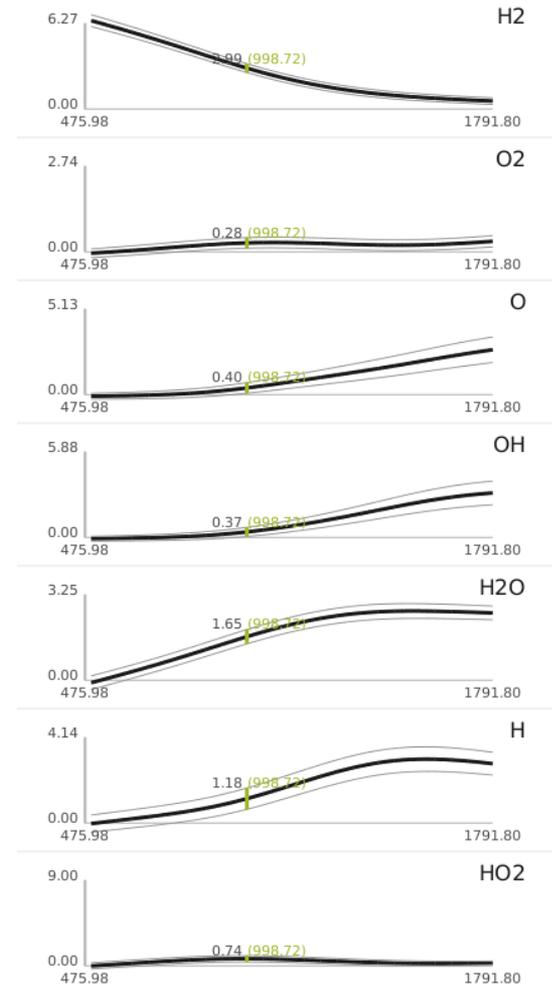
Input: Composition of 10 chemical species

Output: Temperature

The Framework Allows Detailed Visualization and Analysis of High Dimensional Functions

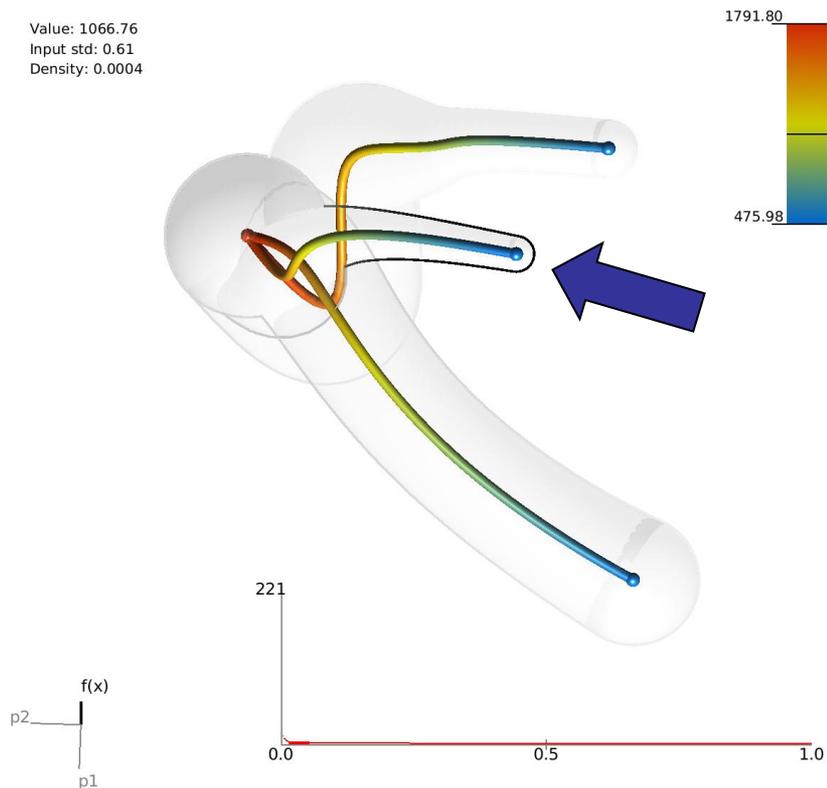


10 dimensional data set describing the heat release wrt. to various chemical species in a combustion simulation

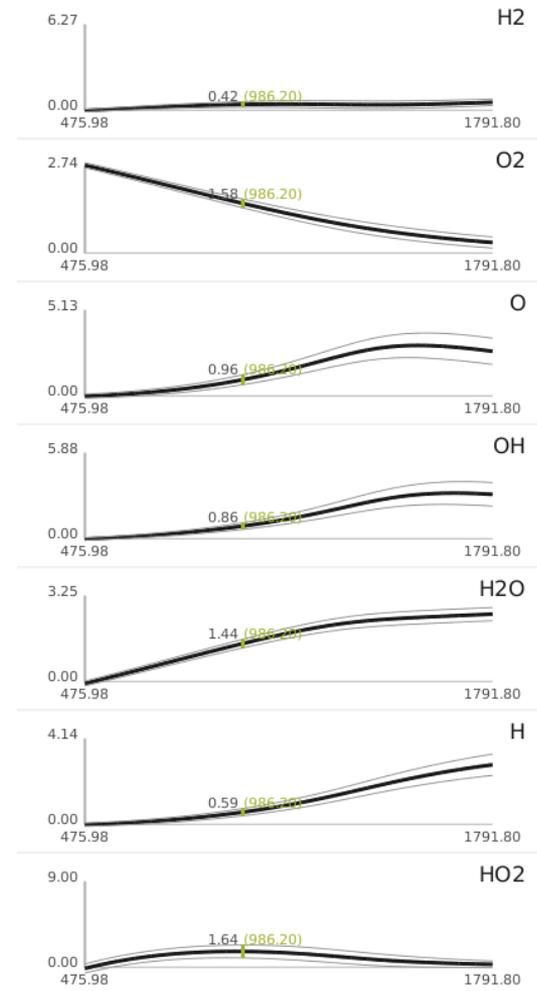


Pure fuel

The Framework Allows Detailed Visualization and Analysis of High Dimensional Functions



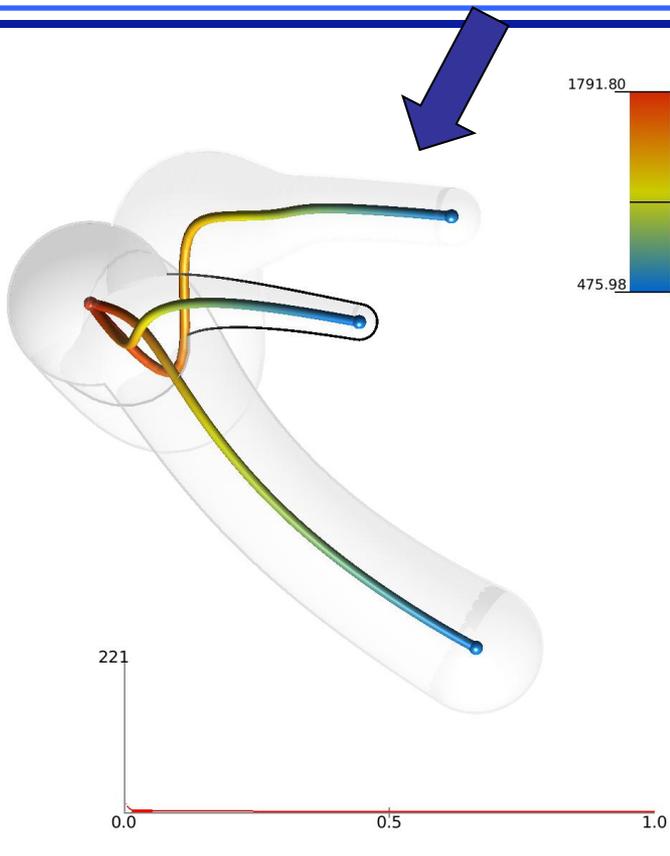
10 dimensional data set describing the heat release wrt. to various chemical species in a combustion simulation



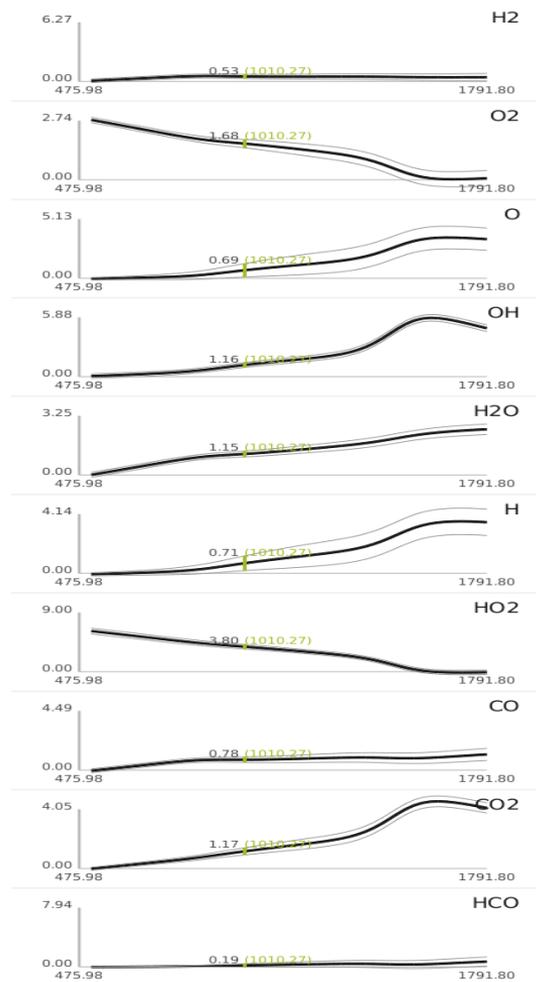
Pure oxidizer

The Framework Allows Detailed Visualization and Analysis of High Dimensional Functions

Value: 1066.76
Input std: 0.61
Density: 0.0004



10 dimensional data set describing the heat release wrt. to various chemical species in a combustion simulation



Local extinction

Combustion Simulation of Jet CO/H₂-Air Flames

Input: Composition of 10 chemical species

Output: Temperature

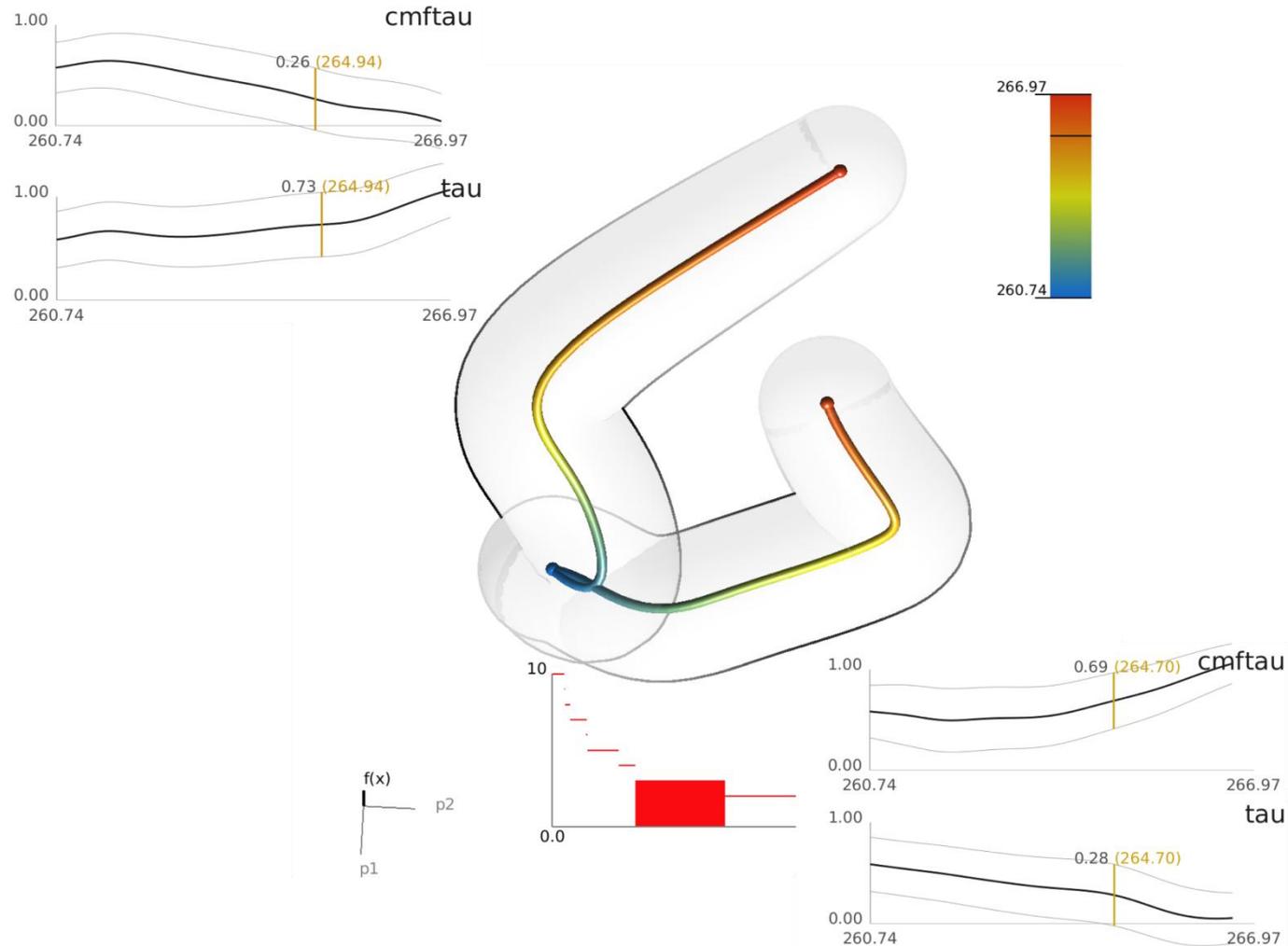
Analysis of Climate Data

Community Atmosphere Climate Model

Input: 21 parameter settings

Output: Net long wave flux (thermal radiation)

The Framework Reveals Relationship Between Convection and Global Long Wave Flux



Community Atmosphere Climate Model



Input: 21 parameter settings

Output: Net long wave flux (thermal radiation)

Data Analysis and Visualization Center is a Catalyst for a Virtuous Cycle of Collaborative Activities



- **Tight cycle of :**
 - basic research,
 - software deployment
 - user support
- **Coordination among many projects:**
 - unified techniques for several applications
- **Strong University-Lab-Industry collaboration**
- **Focused technical approach:**
 - performance tools for fast data access
 - general purpose data exploration
 - error bounded quantitative analysis
 - feature extraction and tracking
- **Interdisciplinary collaboration with domain scientists (from math to physics):**
 - motivating the work
 - formal theoretical approaches
 - feedback to specific disciplines

