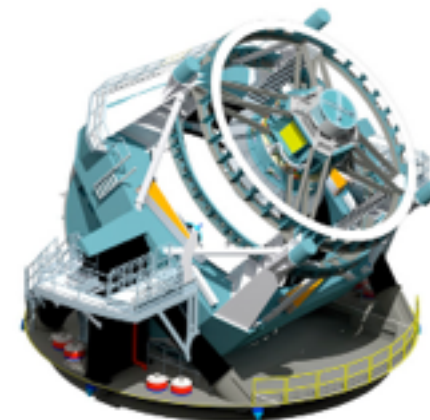# Computation-Driven Discovery for the Dark Universe

Salman Habib
HEP and MCS Divisions
Argonne National Laboratory

PIs: K. Heitmann (ANL), A. Slozar (BNL), S. Dodelson (FNAL),
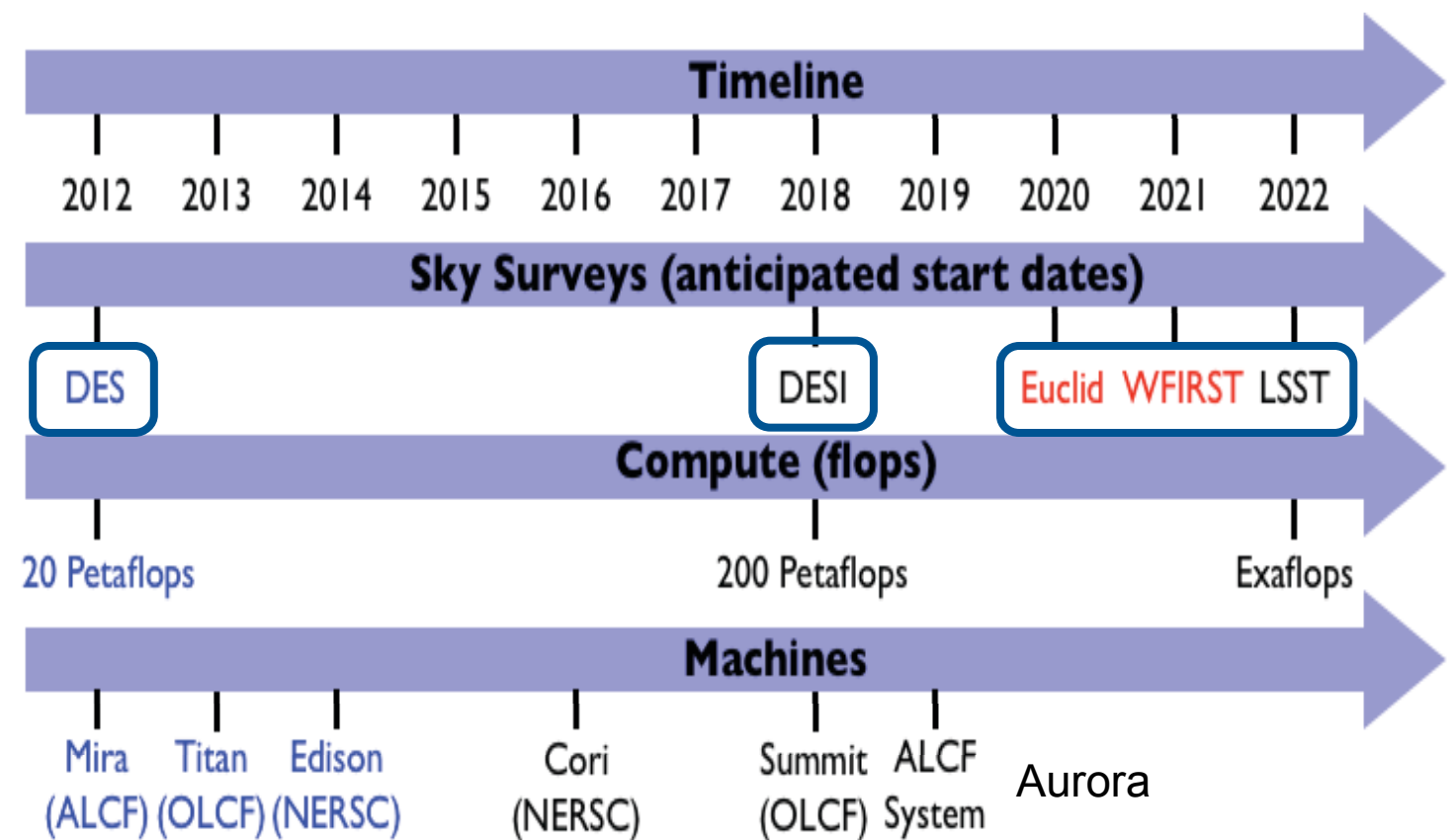P. Nugent (LBNL), J. Ahrens (LANL), R. Wechsler (SLAC)

*ASCR*
*HEP*

# Science at the Cosmic Frontier: 'Precision' Cosmology

- **Instrumentation Advances:** Wide/deep sky coverage, multi-wavelength observations (near-term target, ~1% errors on cosmological parameters or better)

- **Precision Cosmology Science:**

  - Nature of cosmic acceleration (physics of structure formation)

  - Nature and interactions of dark matter

  - Primordial fluctuations and Early universe

  - Probes of fundamental physics (e.g., neutrino sector)



- **Theory/Modeling:** Predictive theory and modeling becoming critically important

- **Computational Cosmology:** Primary theoretical/modeling approach to bridge the the **"theory/modeling gap"**

# Data 'Overload': Observations of Cosmic Structure

- Cosmology=Physics+Statistics
  - Mapping the sky with large-area surveys across multiple wave-bands, at remarkably low levels of statistical error

**CMB temperature anisotropy: theory meets observations**

Galaxies in a moon-sized patch (Deep Lens Survey). LSST will cover 50,000 times this size (~400PB of data)

**LSST**

**SDSS BOSS**

**The same signal in the galaxy distribution**

# The Precision Cosmology Revolution



Compilation for SH by E. Gawiser (1999)

sCDM

Planck (2013)

Compilation (1999)

CMB

LSS

Equivalent to one modern GPU

4 orders of magnitude!

Concurrent Supercomputing Progress

BOSS (2013)

post-recon

# Connecting Theory and Observations: Challenges & Opportunities



Supercomputer    Dark matter    Galaxies

Theory

LSST weak lensing shear power spectrum

z<0.7

0.7<z<1.2

1.2<z<3

LCDM w=-0.9

LSST galaxies    Sloan Digital Sky Survey

- **Error bars are shrinking dramatically**
  - ▸ Many predictions have to be accurate at the sub-percent level
  - ▸ Modeling and understanding of systematics is becoming ever more important (e.g. baryonic effects)
  - ▸ We can go beyond LCDM and explore new fundamental physics: neutrinos, modified gravity, dynamical dark energy, self-interacting dark matter ...
- **Surveys are going deeper and will target/resolve fainter/different galaxies**
  - ▸ Synthetic sky map making becomes more difficult, more physics needed
  - ▸ Significantly higher resolution simulations will be required
  - ▸ New cosmological probes, cross-correlations across multiple probes

# Example: Analytics/Workflow Complexity



**Gaussian Random Field Initial Conditions** → **High-Resolution N-Body Code (HACC)** → **Multiple Outputs Halo/Sub-Halo Identification** → **Halo Merger Trees** → **Semi-Analytic Modeling Code (Galacticus)** → **Galaxy Catalog** → **Realistic Image Catalog** → **Atmosphere and Instrument Modeling** → **Data Management Pipeline** → **Data Analysis Pipeline** → **Scientific Inference Framework**

- **Simulation Campaigns:** Statistics of virtual universes; construction of emulators

- **Modeling the Measurement:** End-to-End modeling necessary to understand crucial systematic errors

- **PDACS:** Custom workflow system under development

- **Data-Intensive Computing:** New architectures needed

# Large Scale Structure: Vlasov-Poisson Equation

$$\frac{\partial f_i}{\partial t} + \dot{\mathbf{x}}\frac{\partial f_i}{\partial \mathbf{x}} - \nabla\phi\frac{\partial f_i}{\partial \mathbf{p}} = 0, \qquad \mathbf{p} = a^2\dot{\mathbf{x}},$$

$$\nabla^2\phi = 4\pi G a^2(\rho(\mathbf{x},t) - \langle\rho_{\mathrm{dm}}(t)\rangle) = 4\pi G a^2 \Omega_{\mathrm{dm}}\delta_{\mathrm{dm}}\rho_{\mathrm{cr}},$$

$$\delta_{\mathrm{dm}}(\mathbf{x},t) = (\rho_{\mathrm{dm}} - \langle\rho_{\mathrm{dm}}\rangle)/\langle\rho_{\mathrm{dm}}\rangle),$$

$$\rho_{\mathrm{dm}}(\mathbf{x},t) = a^{-3}\sum_i m_i \int d^3\mathbf{p}\, f_i(\mathbf{x},\dot{\mathbf{x}},t).$$

**Cosmological Vlasov-Poisson Equation**

- **Properties of the Cosmological Vlasov-Poisson Equation:**
  - **6-D PDE with long-range interactions, no shielding, all scales matter, models gravity-only, collisionless evolution**
  - **Extreme dynamic range in space and mass (in many applications, million to one, 'everywhere')**
  - **Jeans instability drives structure formation at all scales from smooth Gaussian random field initial conditions**

# Large Scale Structure Simulation Requirements
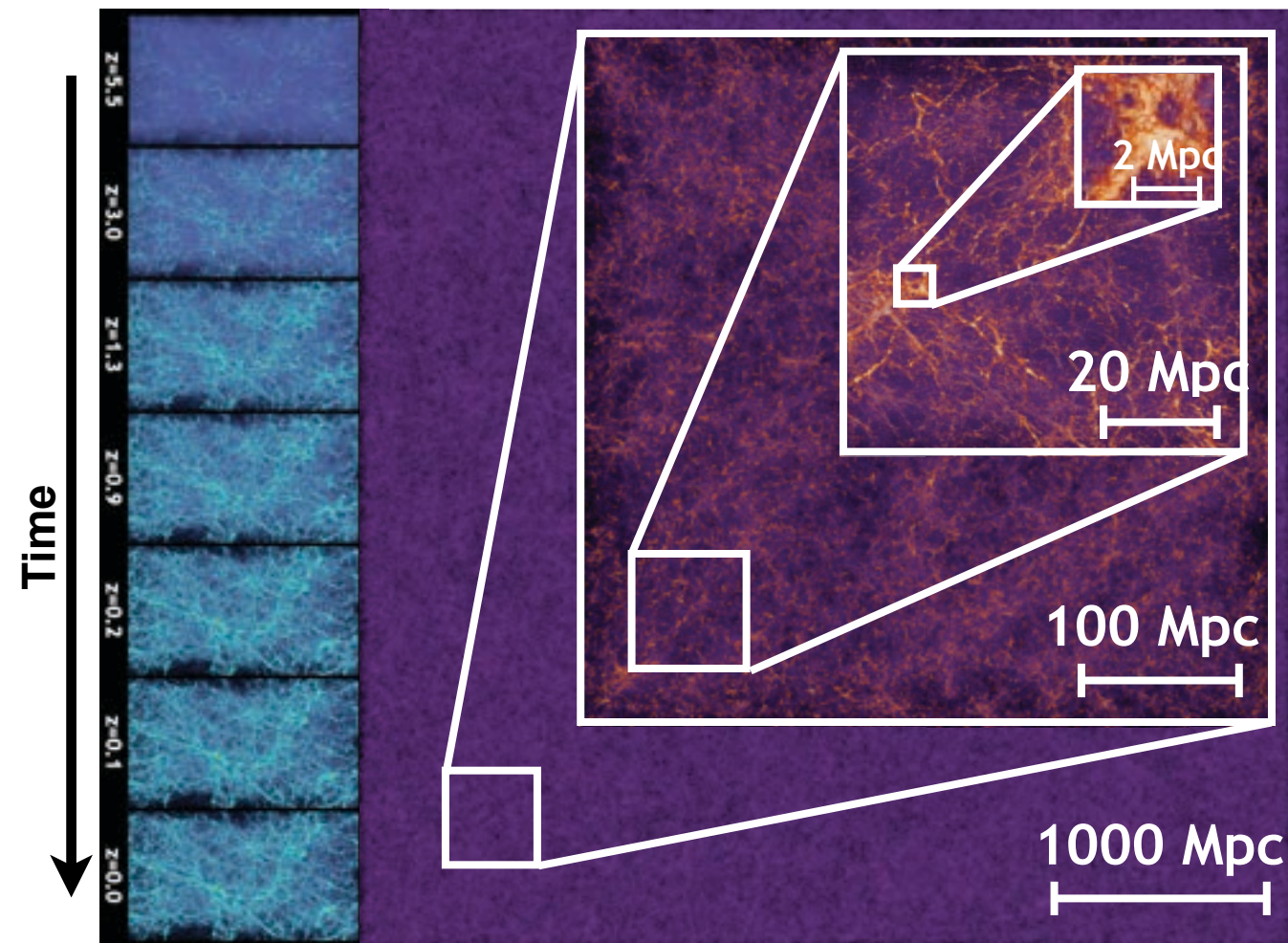
- **Force and Mass Resolution:**

  - Galaxy halos ~100kpc, hence force resolution has to be ~kpc; with Gpc box-sizes, a dynamic range of a million to one

  - Ratio of largest object mass to lightest is ~10000:1

- **Physics:**

  - Gravity dominates at scales greater than ~0.1 Mpc

  - Small scales: galaxy modeling, semi-analytic methods to incorporate gas physics/feedback/star formation

- **Computing 'Boundary Conditions':**

  - Total memory in the PB+ class
  - Performance in the 10 PFlops+ class
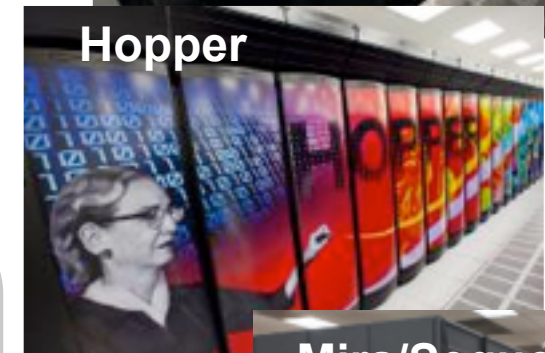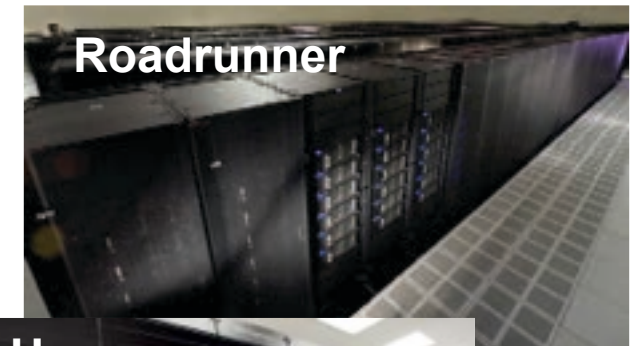  - Wall-clock of ~days/week, in situ analysis



Gravitational Jeans Instability: 'Outer Rim' run with 1.1 trillion particles

Key motivation for HACC: Can the Universe be run as a short computational 'experiment'?
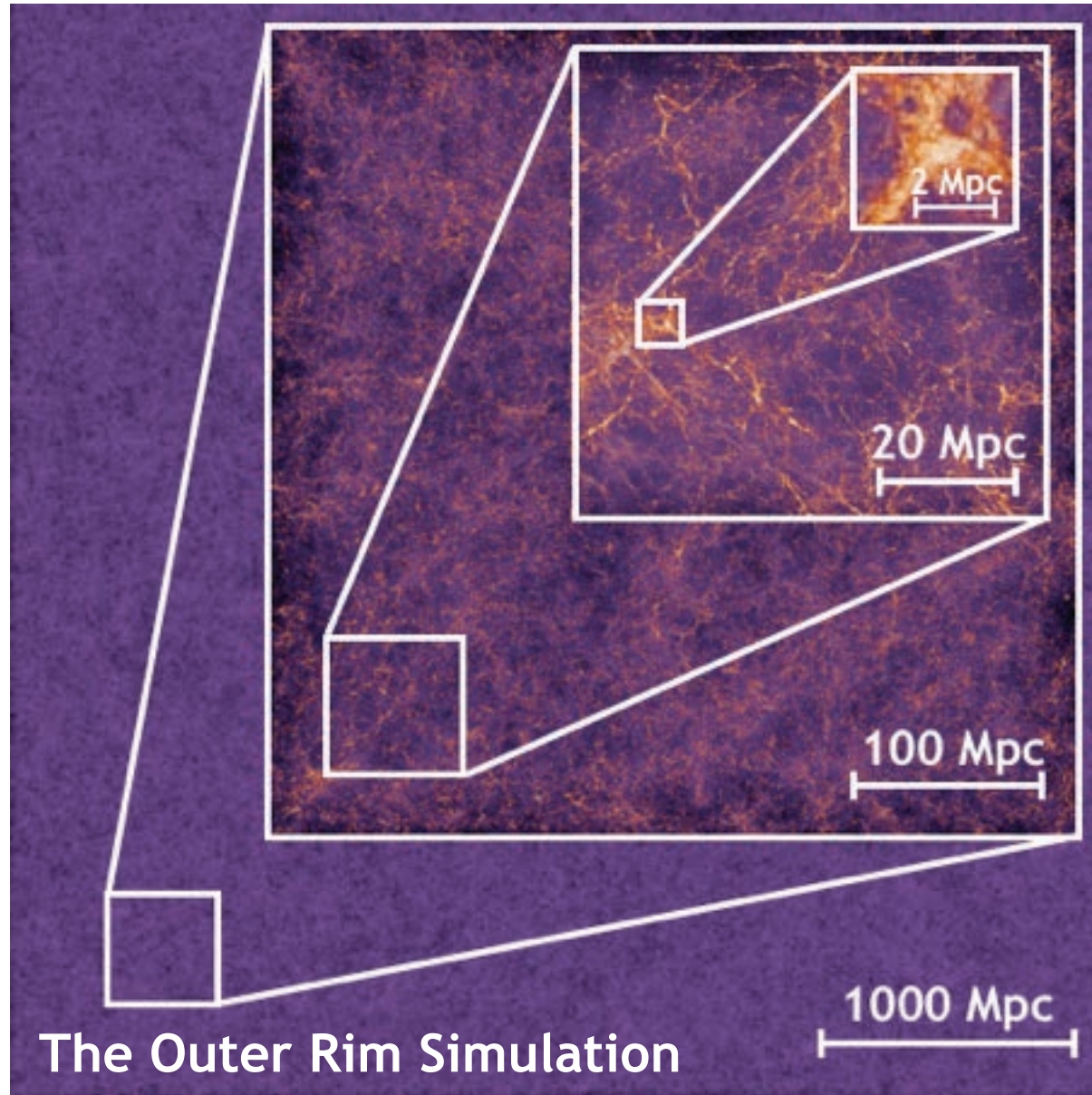
# Combating Architectural Diversity with HACC

- **Architecture-independent performance/scalability:** 'Universal' top layer + 'plug in' node-level components; minimize data structure complexity and data motion

- **Programming model:** 'C++/MPI + X' where X = OpenMP, Cell SDK, OpenCL, CUDA, --

- **Algorithm Co-Design:** Multiple algorithm options, stresses accuracy, low memory overhead, no external libraries in simulation path

- **Analysis tools:** Major analysis framework, tools deployed in stand-alone and in situ modes

- **Performance:** First production science code to break 10PFlops sustained, runs at full scale on all current DOE supercomputing systems (Gordon Bell Finalist 2012/2013, benchmark code for CORAL procurement)

- **Load Balancing:** New, highly efficient, task-based load balancing implemented on Titan
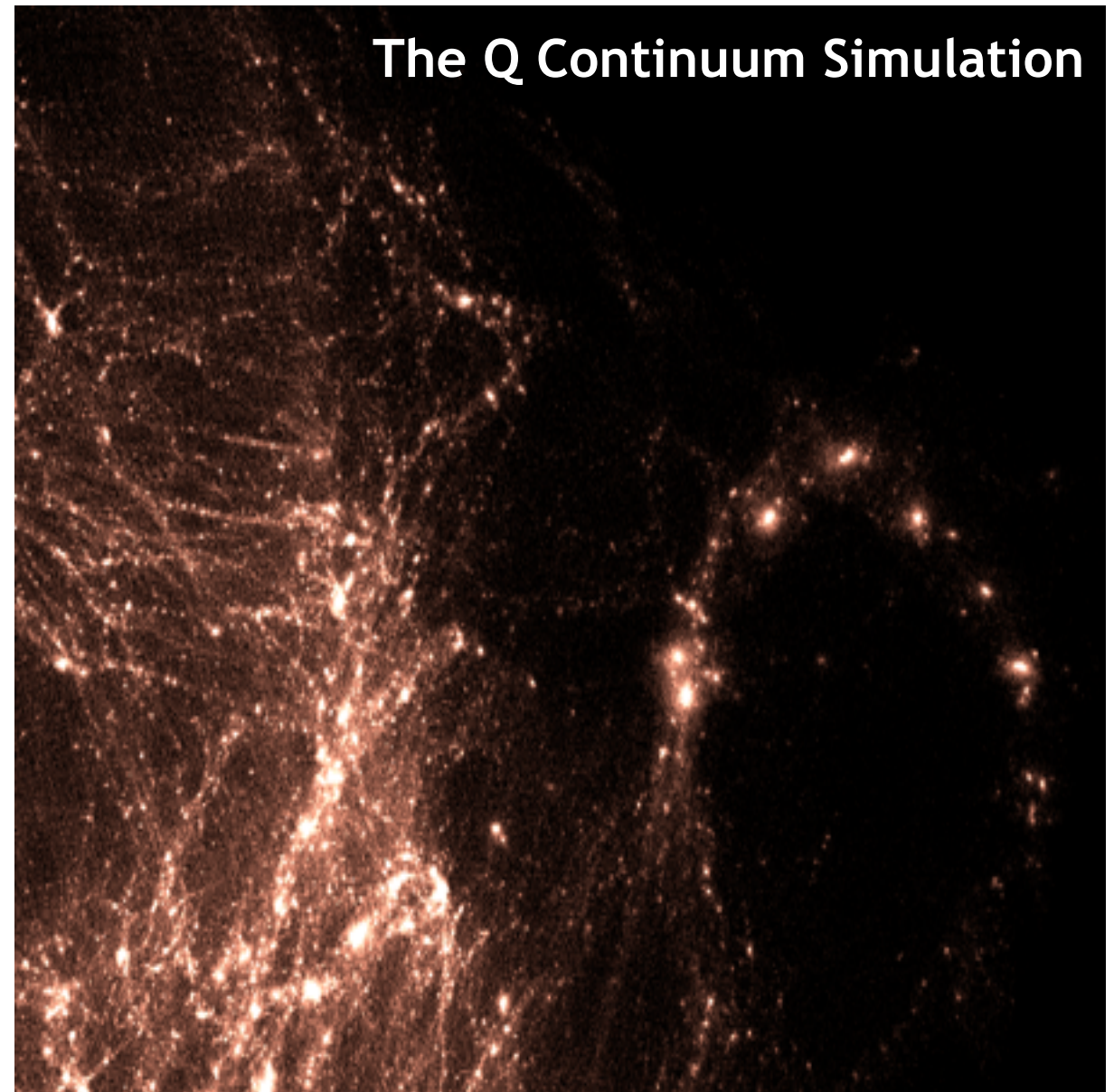
Roadrunner

Hopper

Mira/Sequoia

Titan

Edison

S. Habib et al., arXiv:1410.2805 (in press)

# The Q Continuum and the Outer Rim Simulations

## Simulating the LCDM Universe with Unprecedented Volume and Resolution



The Outer Rim Simulation
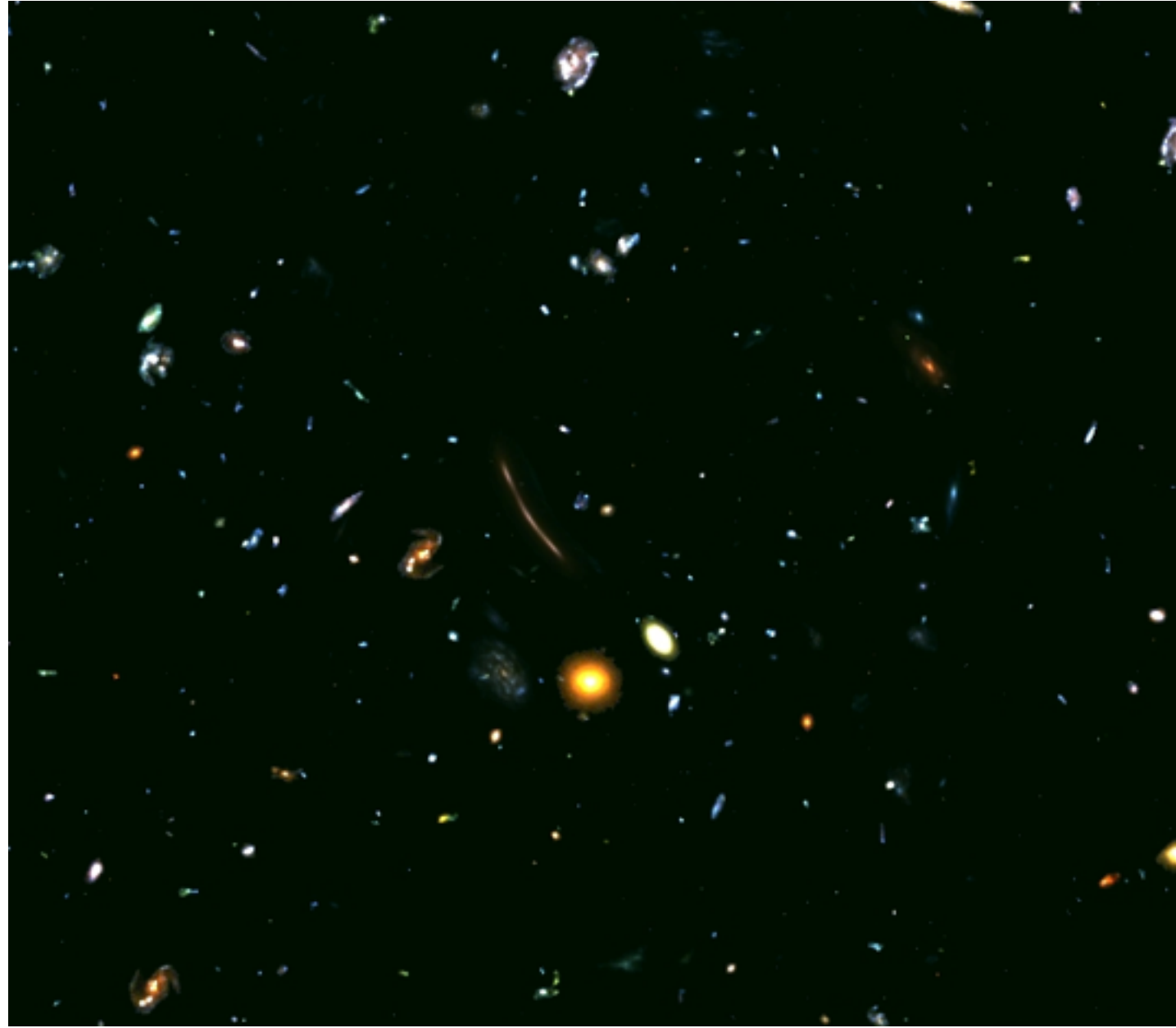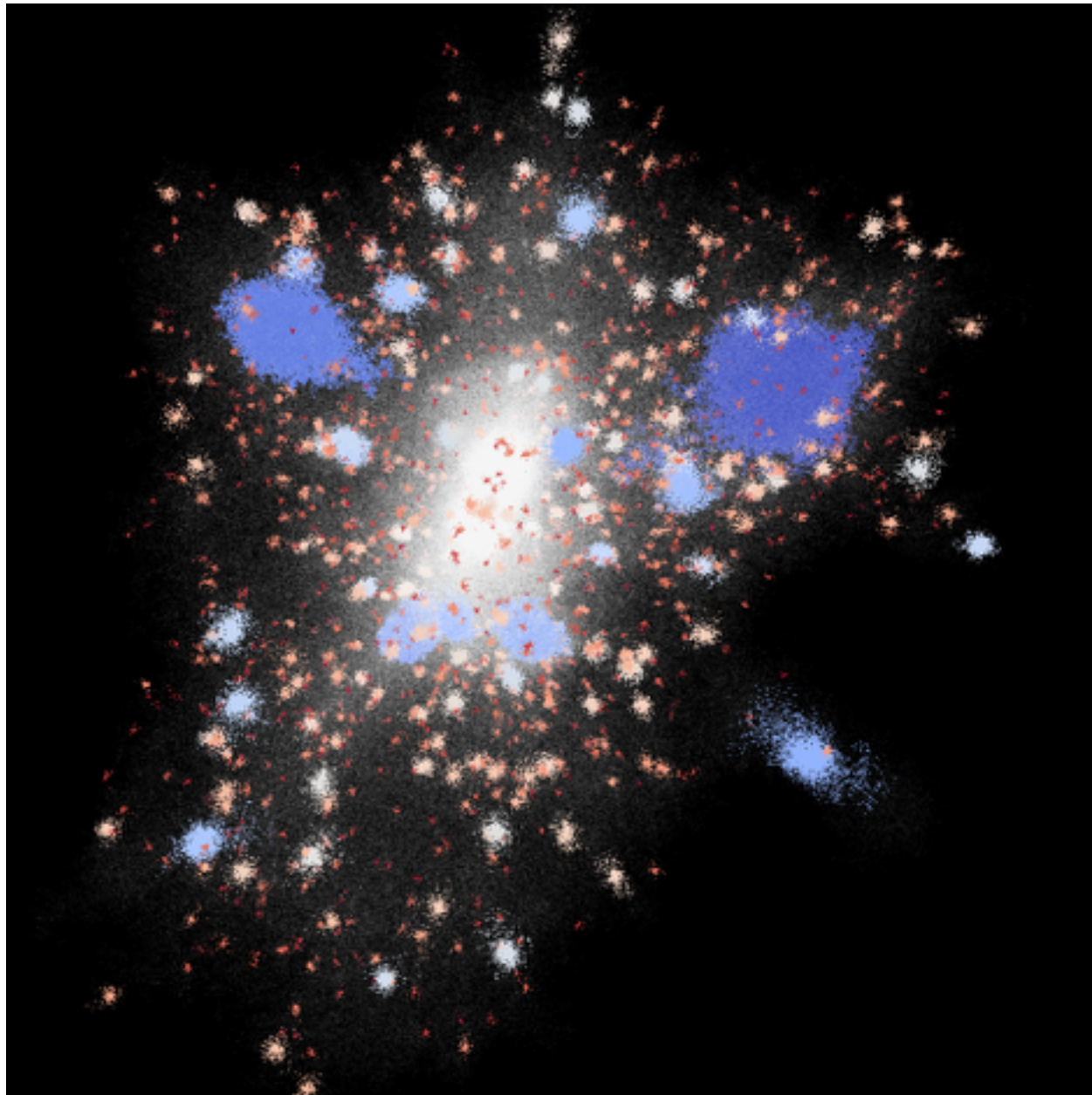


The Q Continuum Simulation

$(4225 \text{ Mpc})^3$ volume, 1.07 trillion particles carried out on ~67% of Mira at Argonne, 4PB of data, 216x Millennium simulation

$(1300 \text{ Mpc})^3$ volume, 0.55 trillion particles carried out on ~90% of Titan at Oak Ridge, 2PB of data, 64x Bolshoi simulation
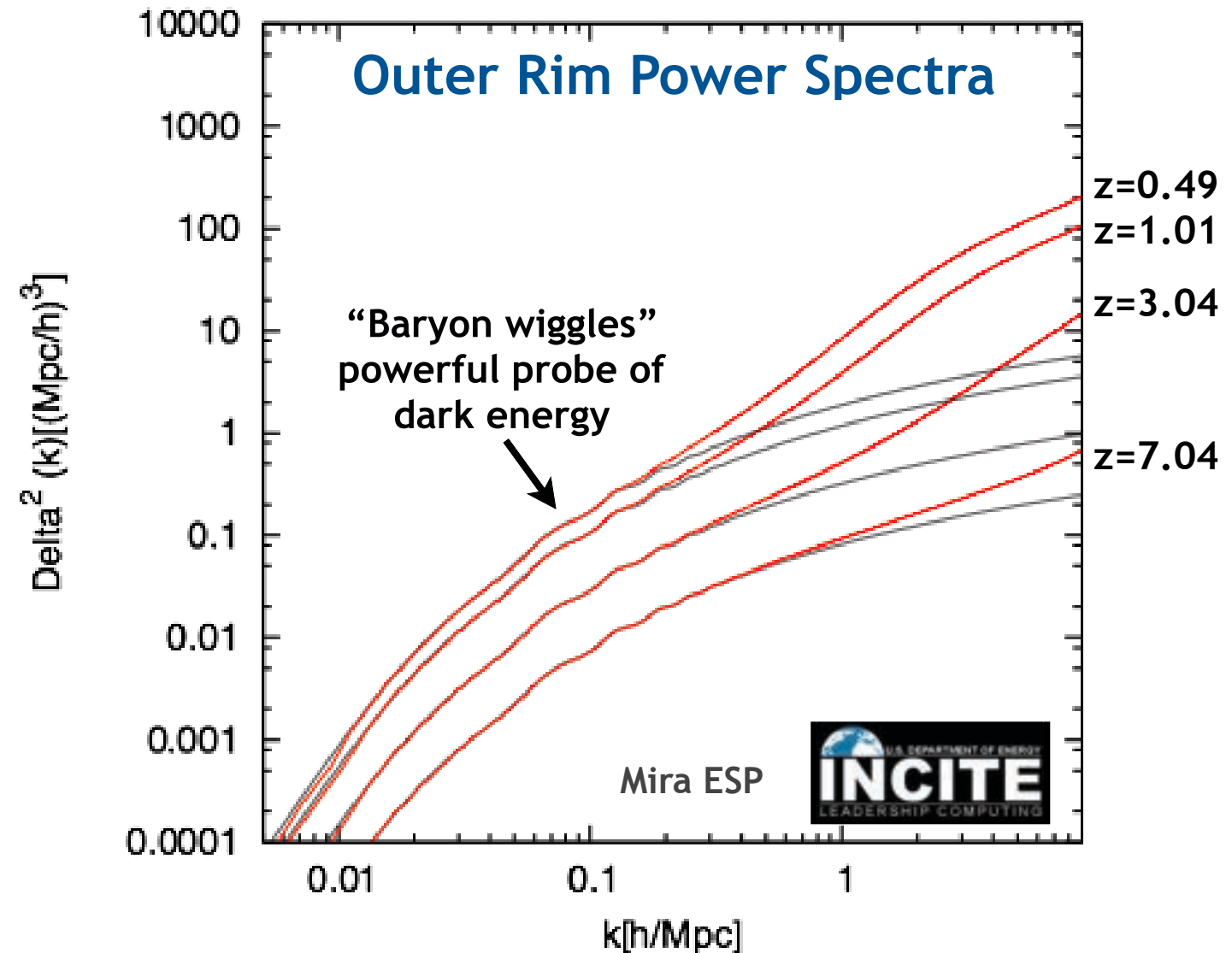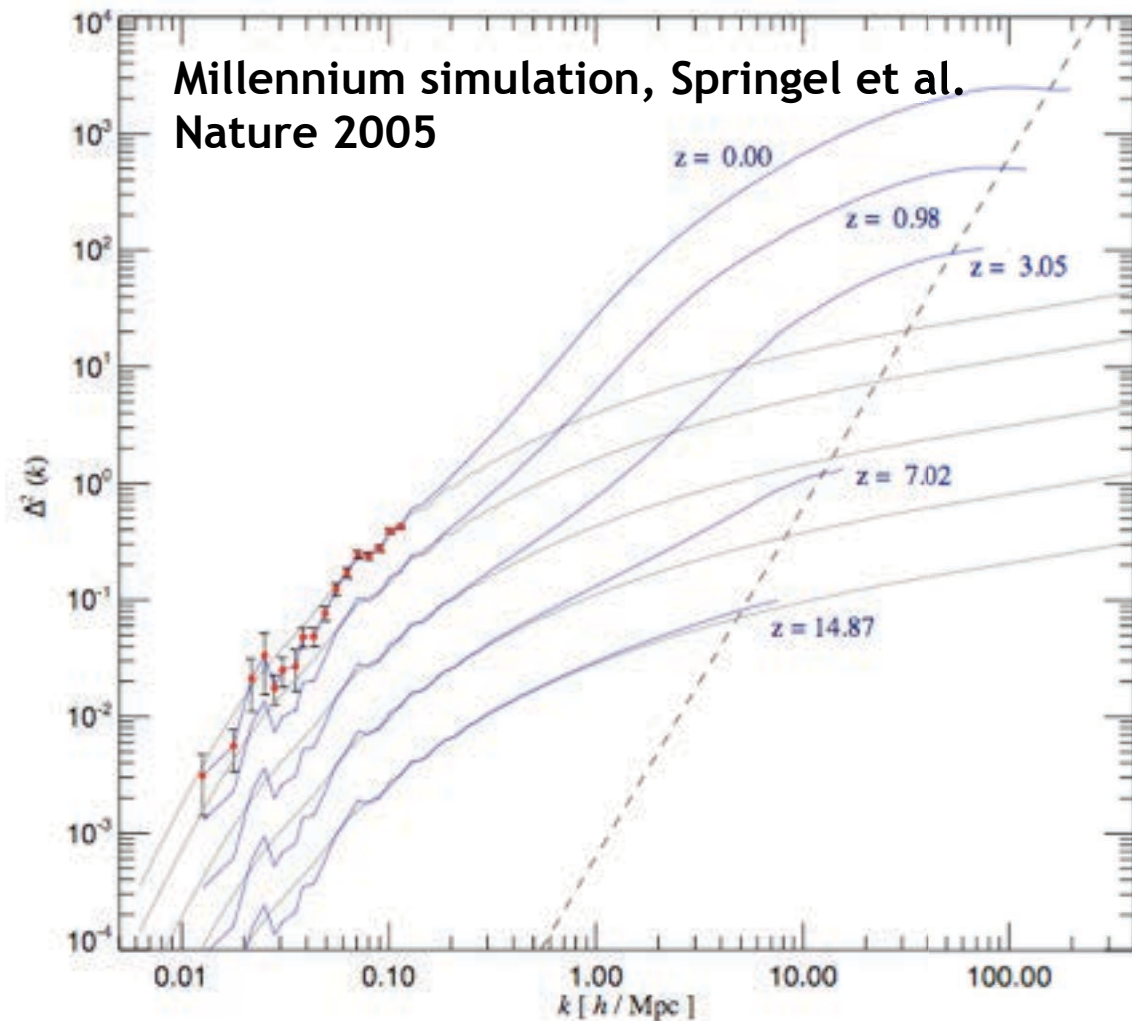
# Cosmology with the Q Continuum Run



- Many cosmological statistics available at high resolution
- Highly resolved cluster-scale halos used for strong lensing predictions (left, halo with ~1400 subhalos, right, background galaxies lensed by a simulated cluster)

**Heitmann et al. 2015 (in press)**

# Cosmology with HACC: Exquisite Statistics



Millennium simulation, Springel et al. Nature 2005

Outer Rim Power Spectra

"Baryon wiggles" powerful probe of dark energy
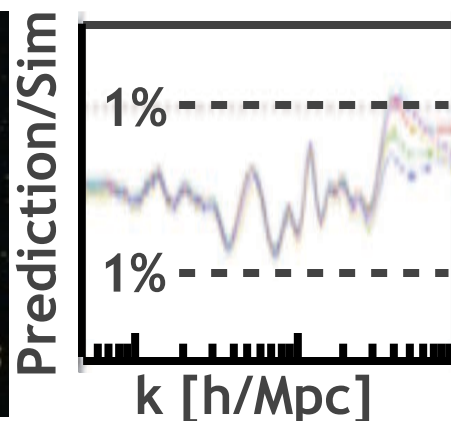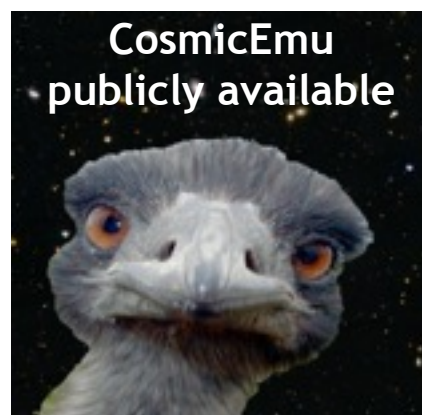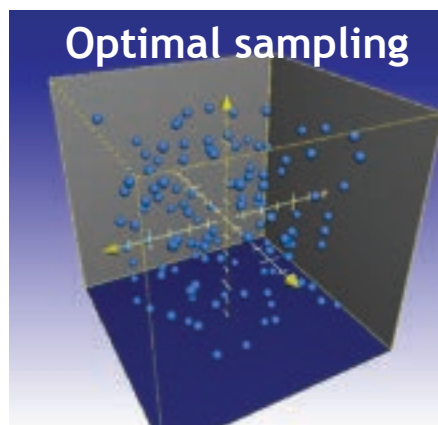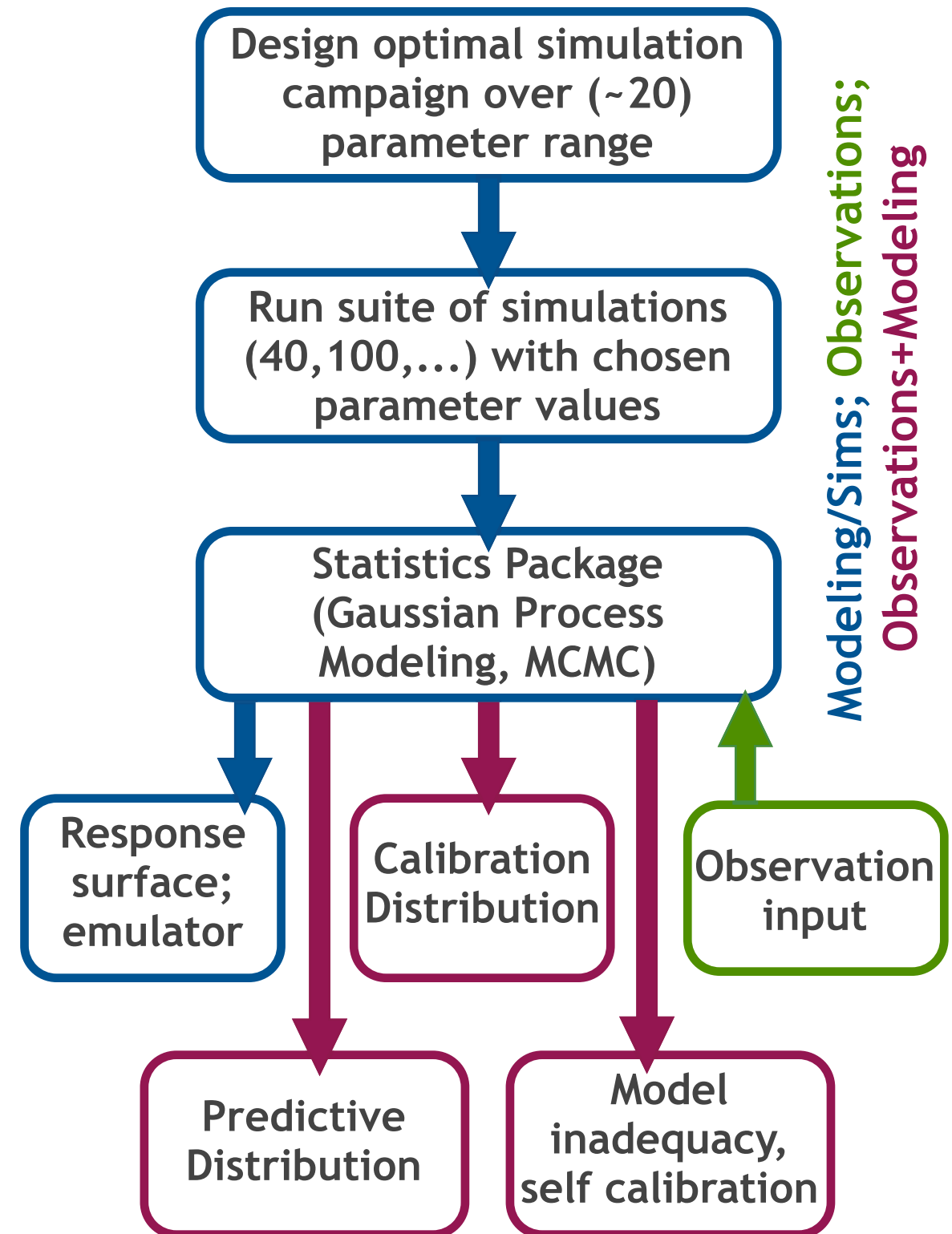
z=0.49
z=1.01
z=3.04
z=7.04

Mira ESP

- Mass resolution of Millennium simulation and Outer Rim run very similar (~ $10^9 M_\odot$ particle mass), but volume different by a factor of 216 (Outer Rim volume = Millennium XXL, but with 7 times higher mass resolution)

- Exceptional statistics at high resolution enable many science projects

**Habib et al. 2014 (in press)**

# Cosmic Calibration: Solving the Inverse Problem

- **Challenge:** To extract cosmological constraints from observations in non-linear regime, need to run Markov Chain Monte Carlo code; input: 10,000 - 100,000 different models

- **Current strategy:** Fitting functions for e.g. P(k), accurate at 10% level, not good enough!
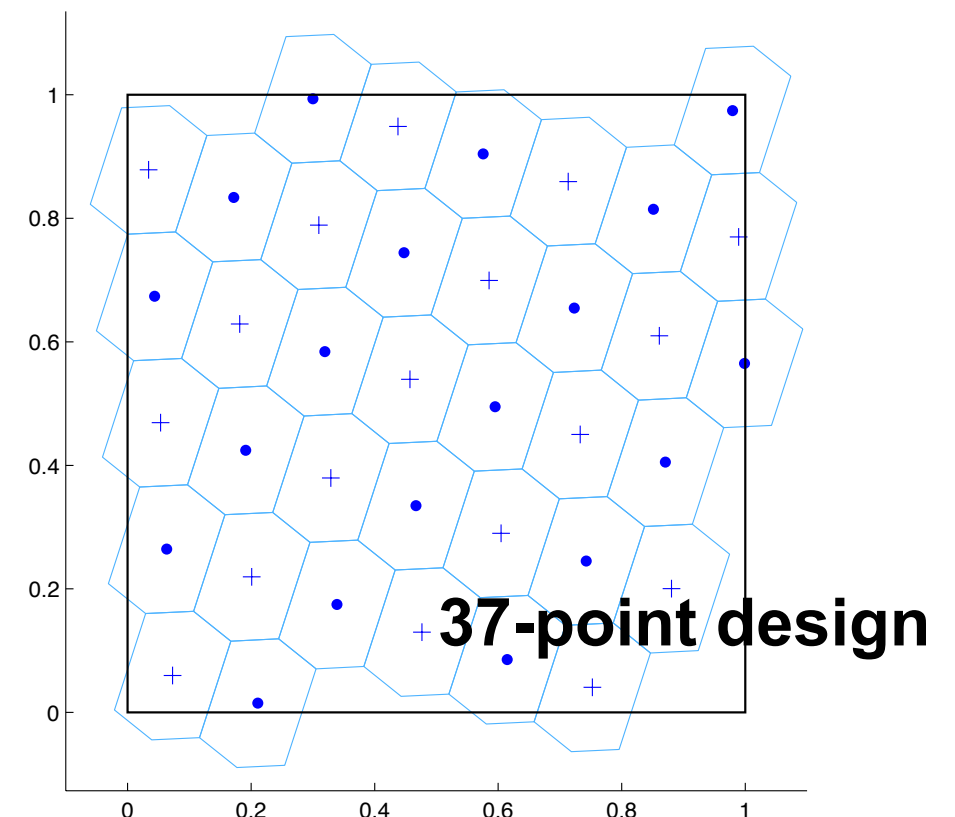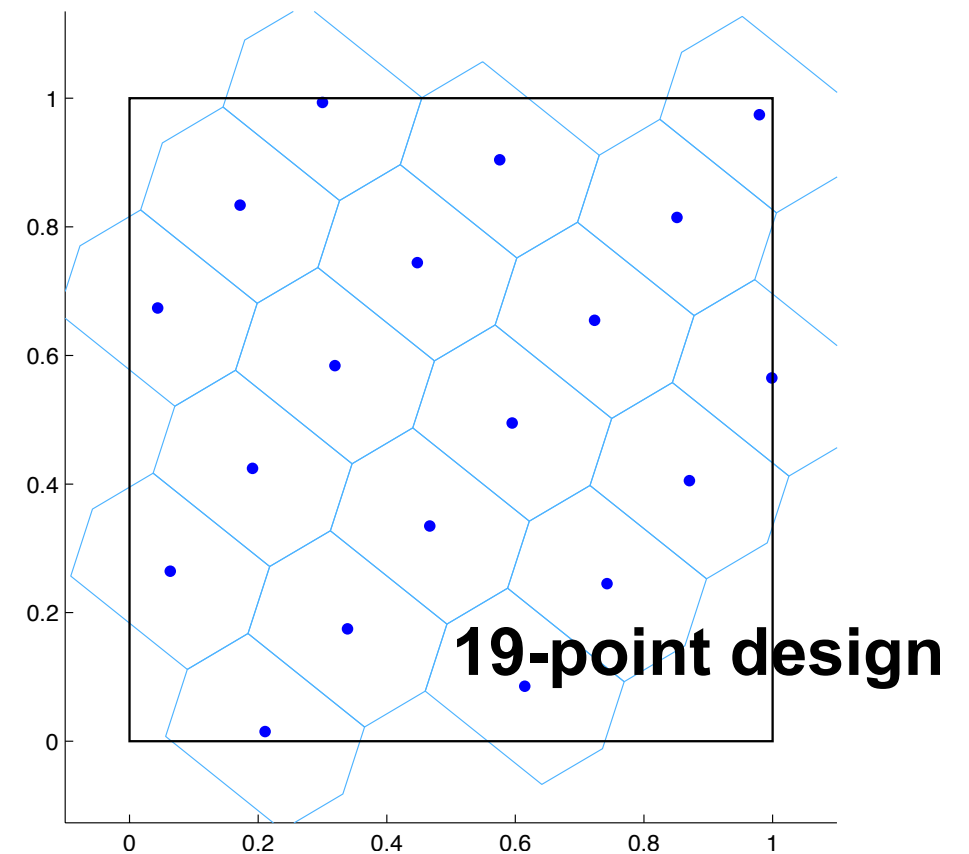
- **Brute force:** Hopeless —

- **Solution:** Emulators



Optimal sampling

CosmicEmu publicly available

Prediction/Sim

1% - - - - - - - - - - -

1% - - - - - - - - - - -

k [h/Mpc]

**Heitmann et al. 2006, Habib et al. 2007**

Design optimal simulation campaign over (~20) parameter range

↓

Run suite of simulations (40,100,...) with chosen parameter values

↓

Statistics Package (Gaussian Process Modeling, MCMC)

Response surface; emulator

Calibration Distribution

Observation input

Predictive Distribution

Model inadequacy, self calibration

**Modeling/Sims; Observations; Observations+Modeling**
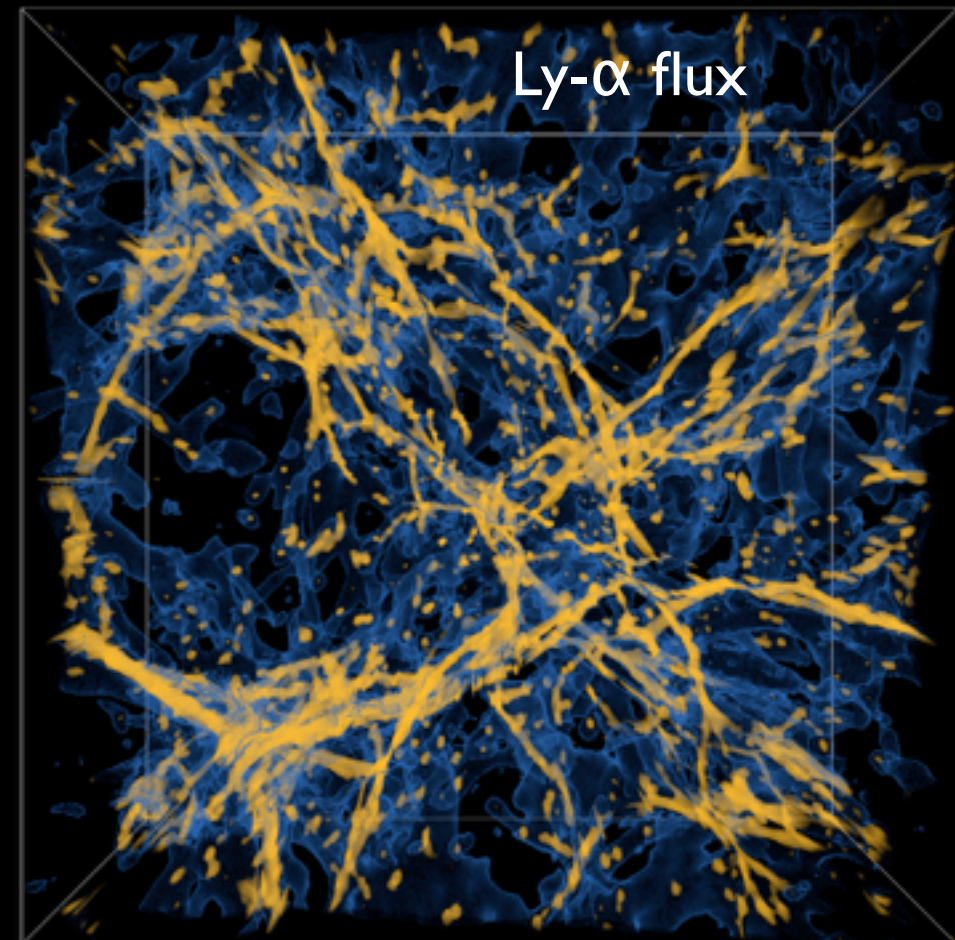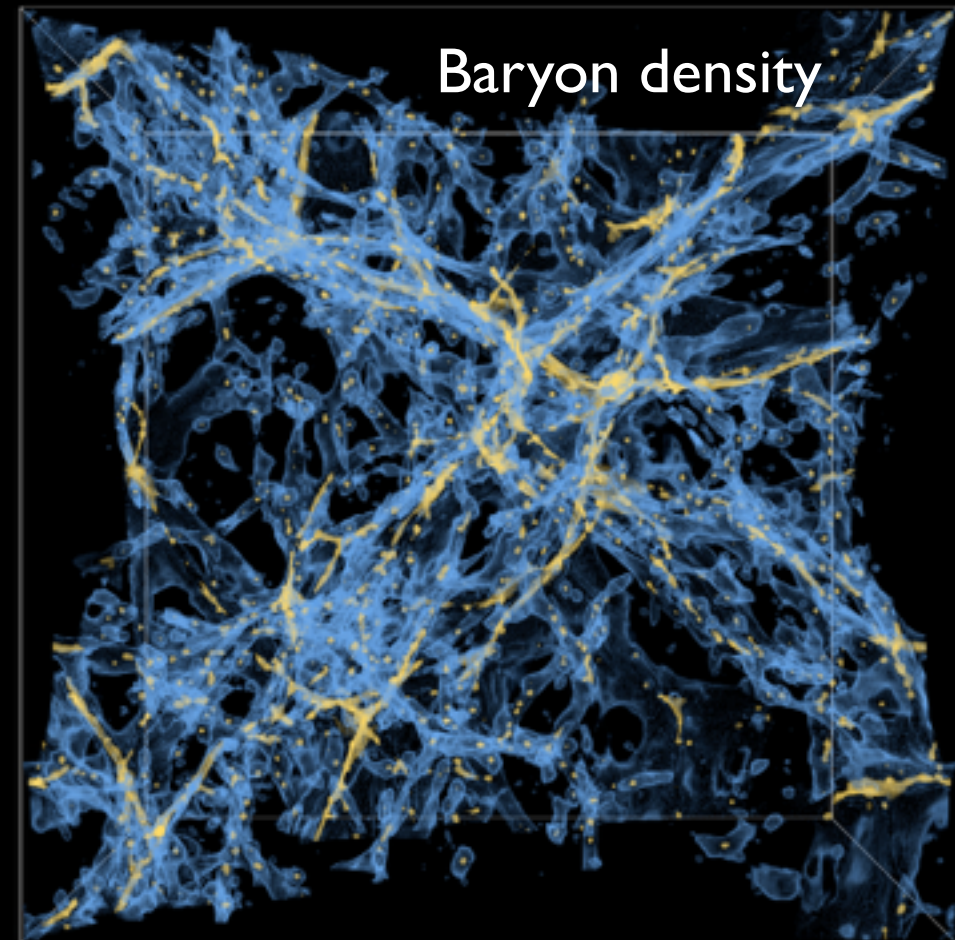
# Emulator Science

- **Previous/Current Results:**

  - Emulator for the matter power spectrum (current state of the art); Heitmann et al. 2014

  - Emulators for galaxy statistics, Kwan et al. 2015 (in press)

  - Emulators for halo profiles (c(M) relation), Kwan et al. 2013

  - Emulators for covariances, in prep.

- **Titan-Mira Universe Suite:**

  - Increase number of dimensions to 8, adding neutrinos and dynamical dark energy

  - Introduce new nested lattice sampling method for increase of accuracy; good accuracy with only 26 simulation runs
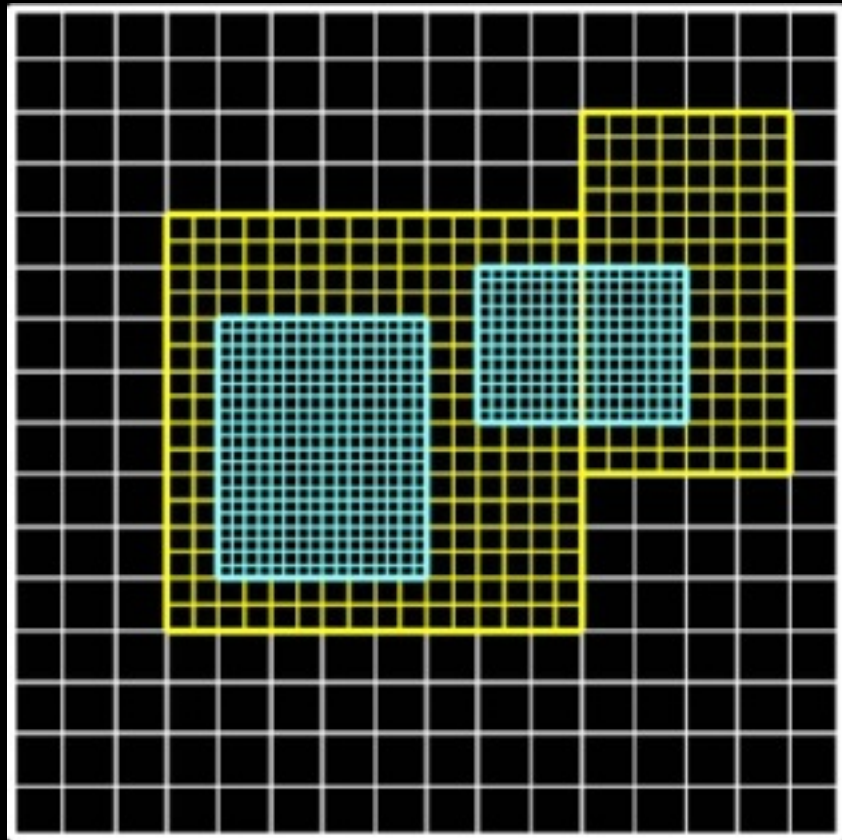
**19-point design**

**37-point design**

**Heitmann et al. 2015 (in prep)**

# Nyx

- 3-D Cartesian grid, finite volume representation

- Evolve dark matter as collisionless Lagrangian fluid

- Evolve baryons as ideal gas using unsplit, Godunov-type methodology

- Adaptive mesh refinement (AMR) to extend dynamic range

- Uses BoxLib software framework developed at LBL
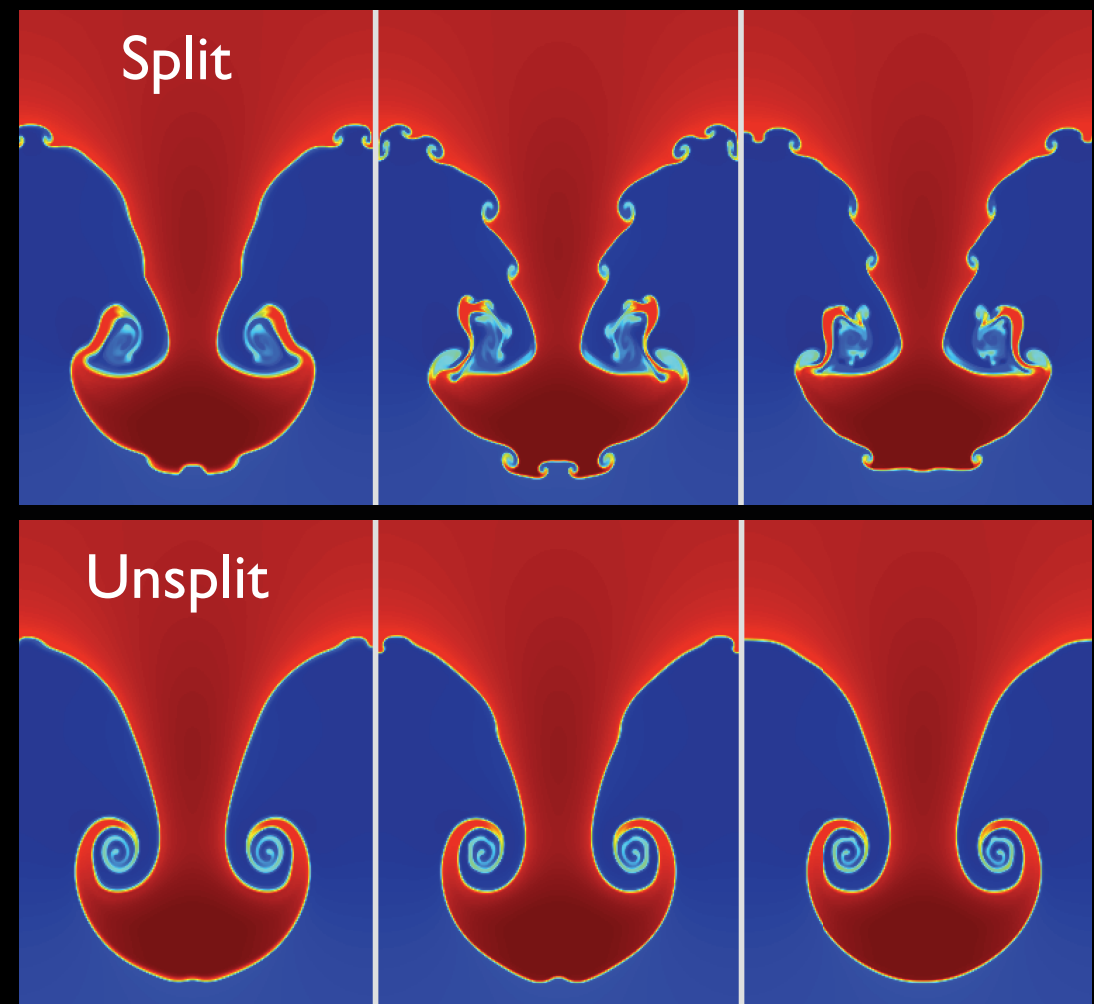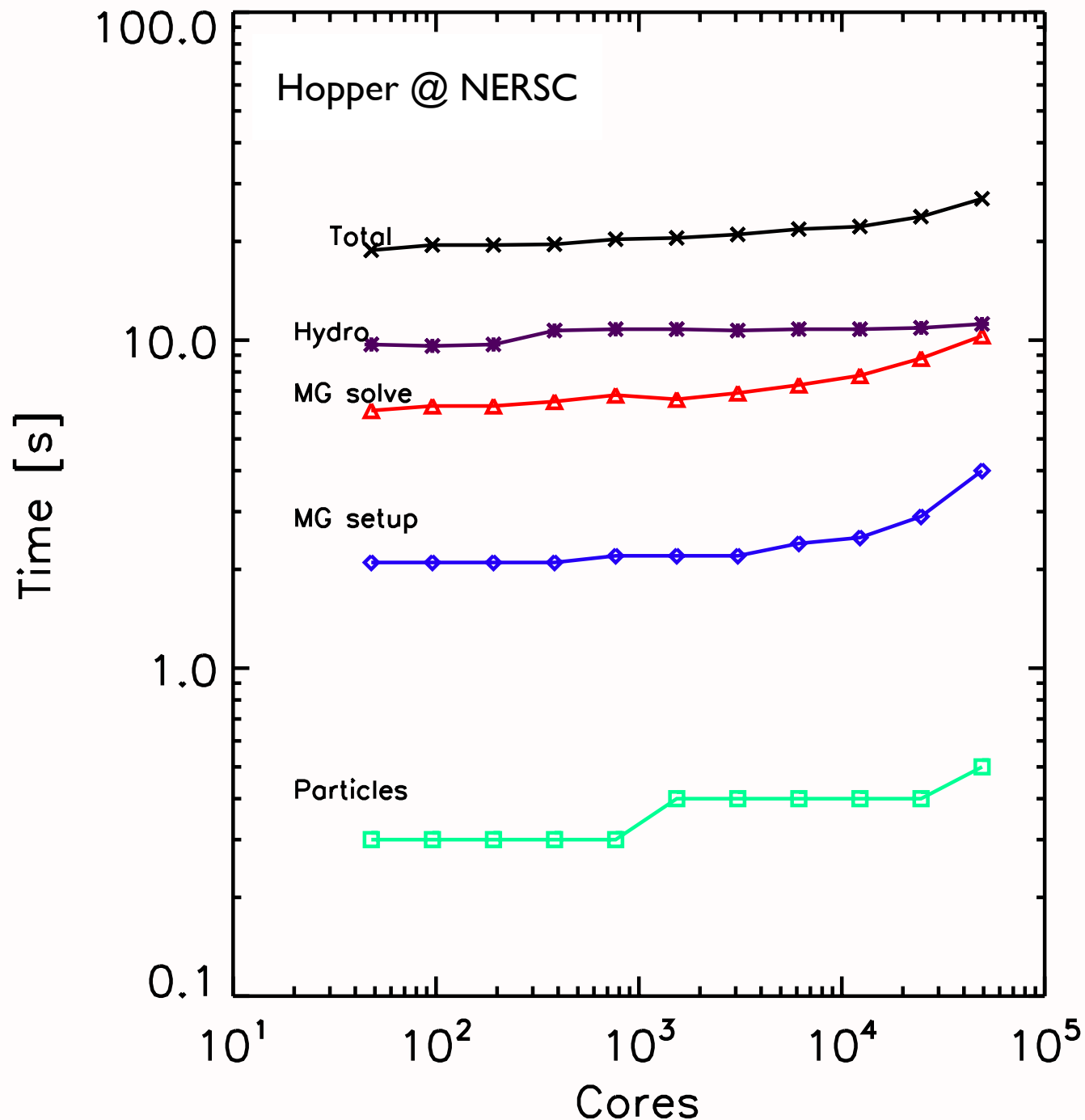
- Code paper: ApJ, 765, 39 (2013)



Baryon density

Ly-$\alpha$ flux

- **AMR:** patch-based refinement, with jump up to a factor of 4.

- **Hydro:** unsplit finite volume scheme better characterizes fluid flow.
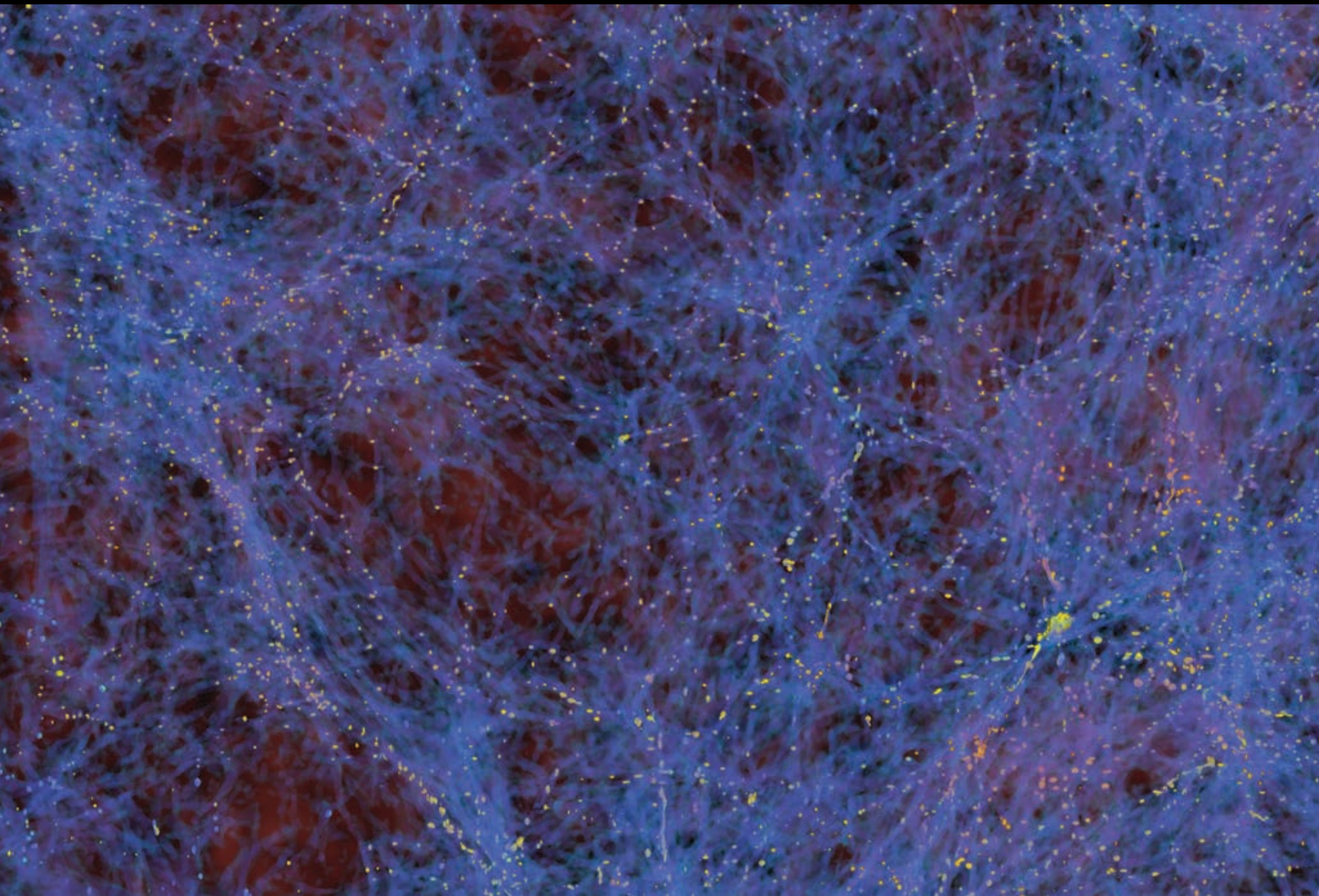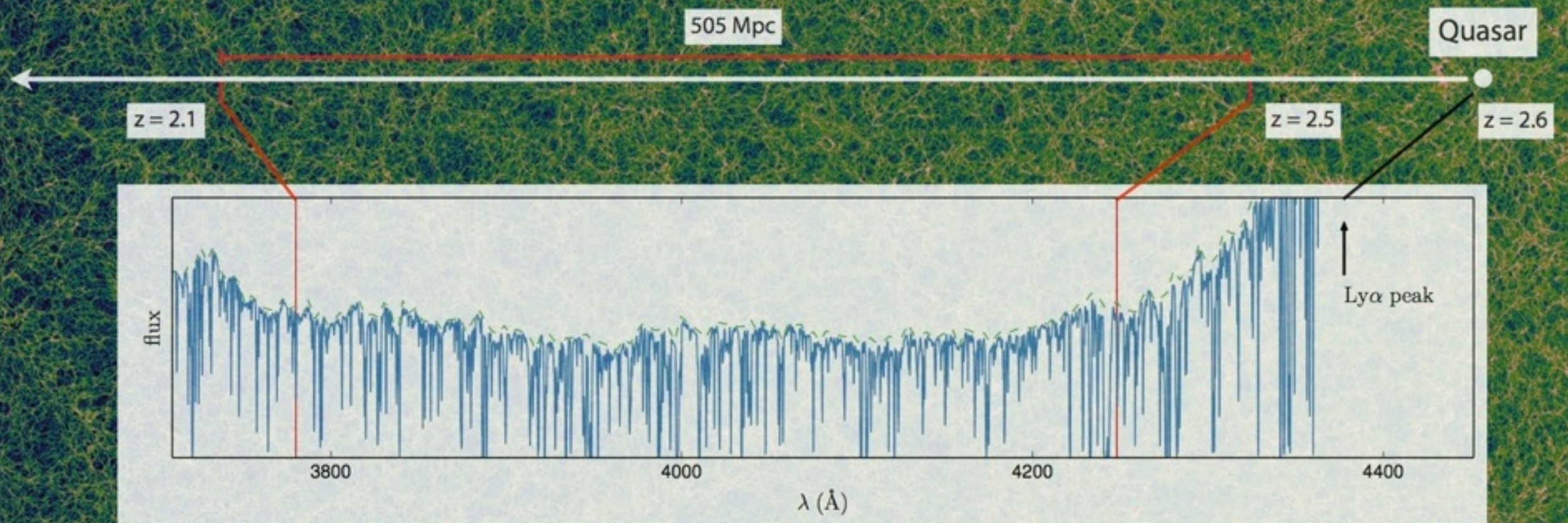
Split

Unsplit

# Excellent scaling



- Currently we are using NERSC resources under ALCC allocation.

- Mostly running $2048^3$ and $4096^3$ runs.

- Hopper/Edison: standard cluster architecture, 24 cores on a node, 32/64GB per node, ~5,000 nodes.

- Analysis pipeline on par with simulations.
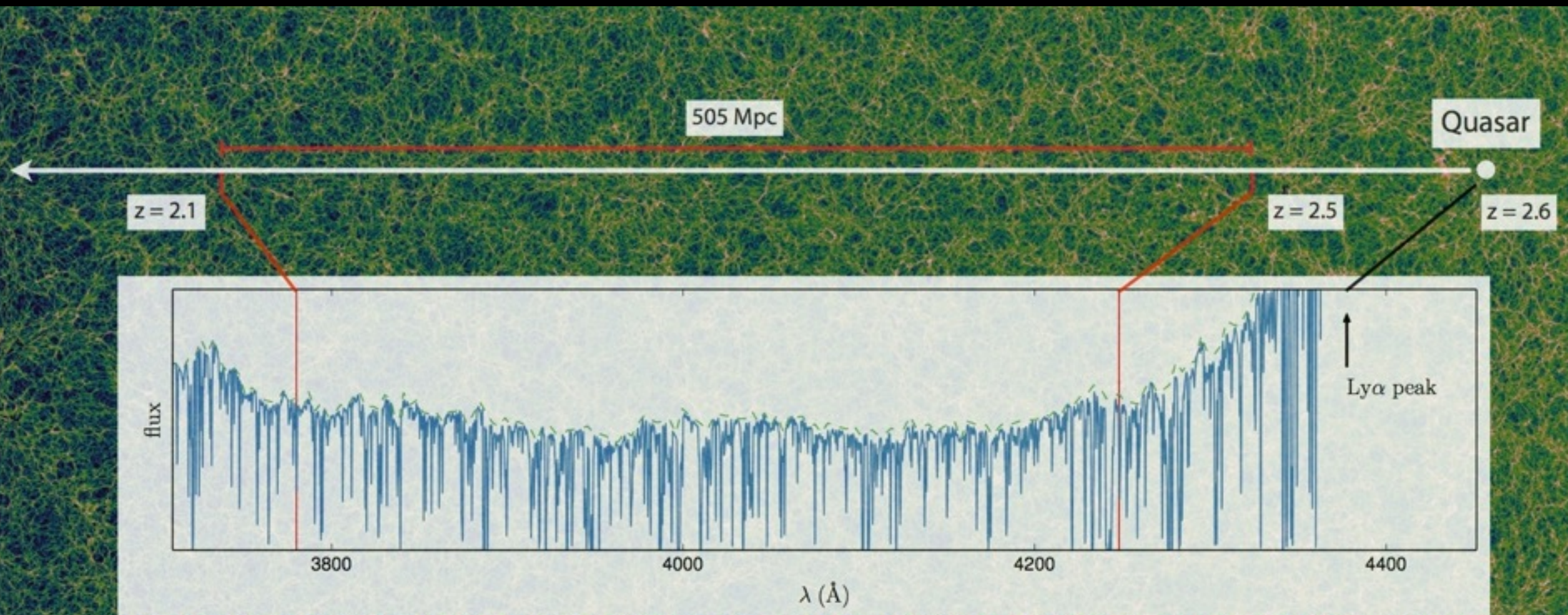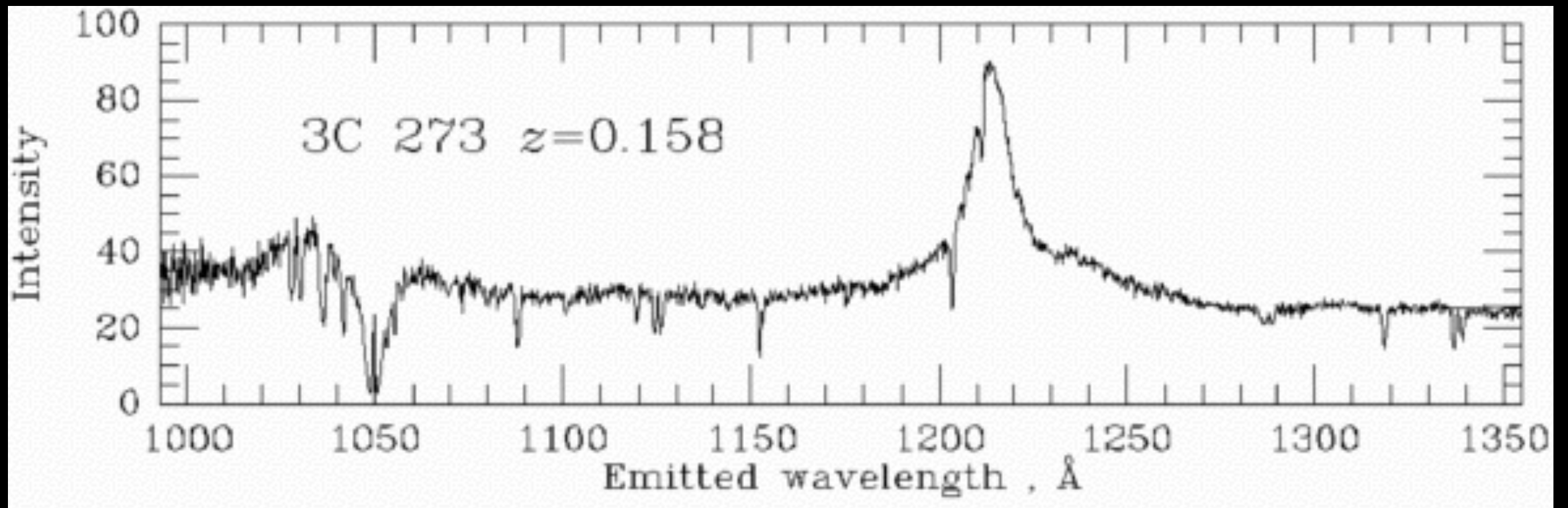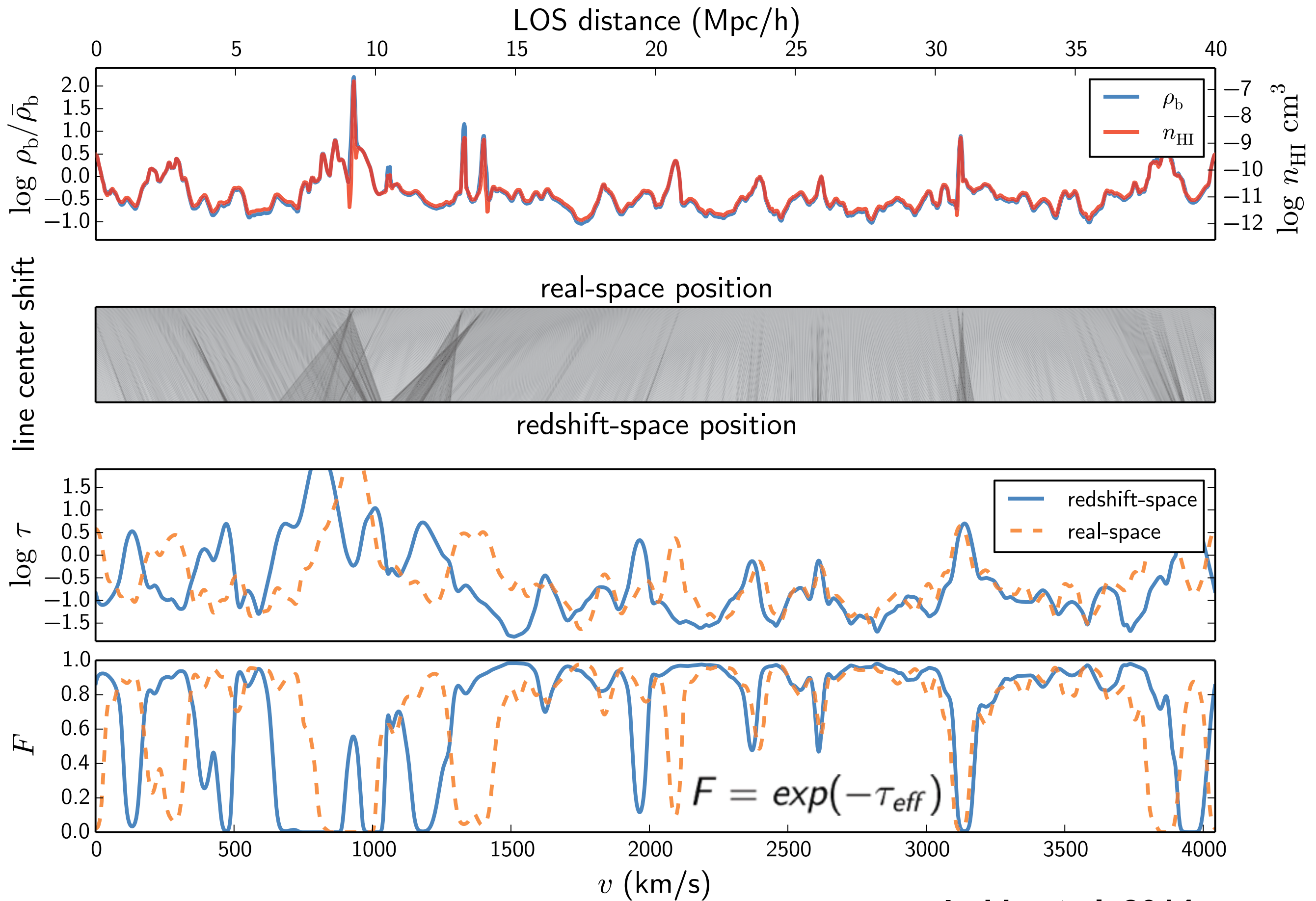
4096³ hydro simulation
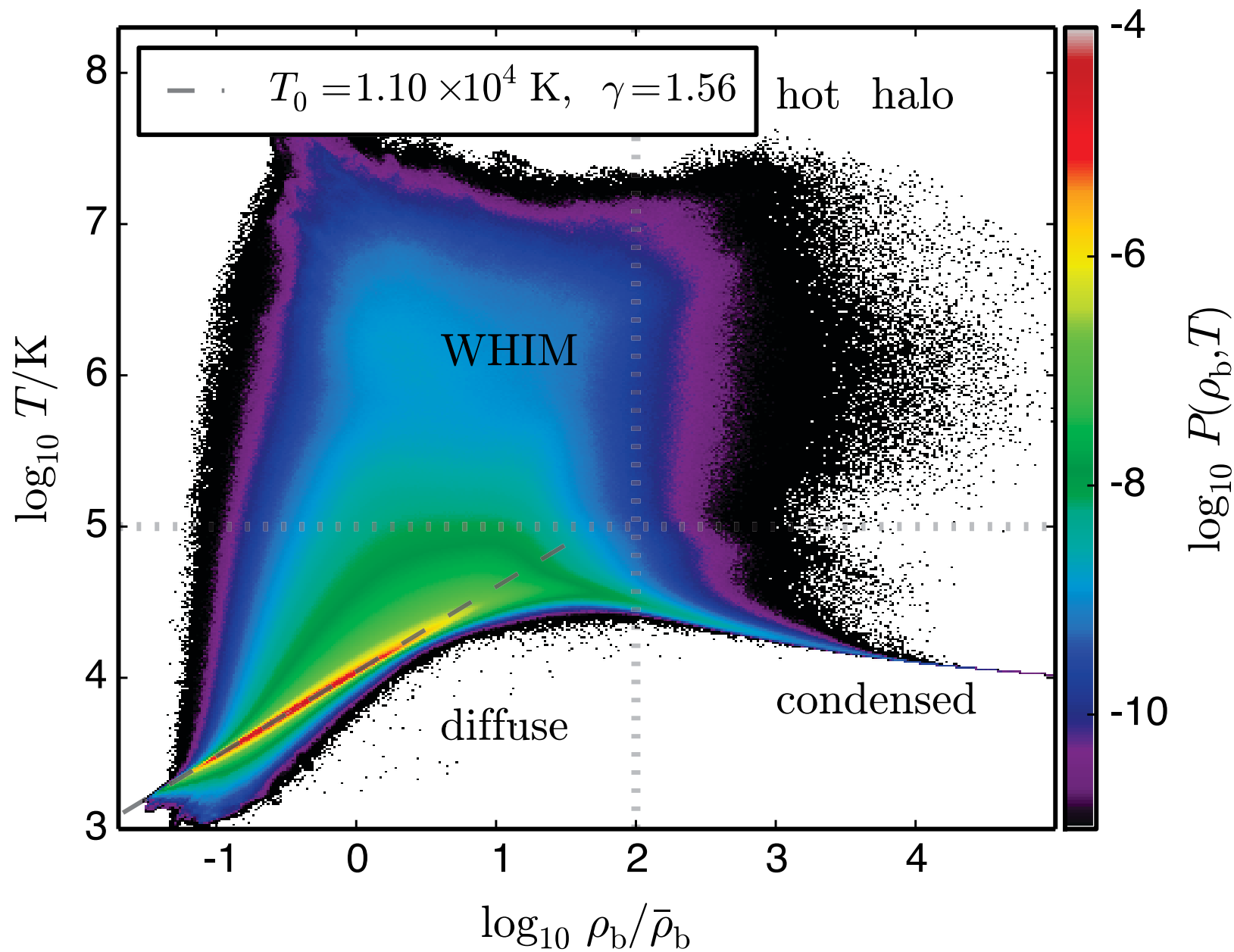
Blue: F~0; Red: F~1

# The Lyman-α forest in optically-thin hydro simulations

Zarija Lukić, Lawrence Berkeley National Laboratory

(Casey Stark, Peter Nugent, Martin White, Avery Meiksin, Ann Almgren)

3C 273 z=0.158



505 Mpc

Quasar

z = 2.1          z = 2.5          z = 2.6

Lyα peak

flux

λ (Å)

LOS distance (Mpc/h)

real-space position

redshift-space position

line center shift

$F = exp(-\tau_{eff})$

$v$ (km/s)

**Lukic et al. 2014**
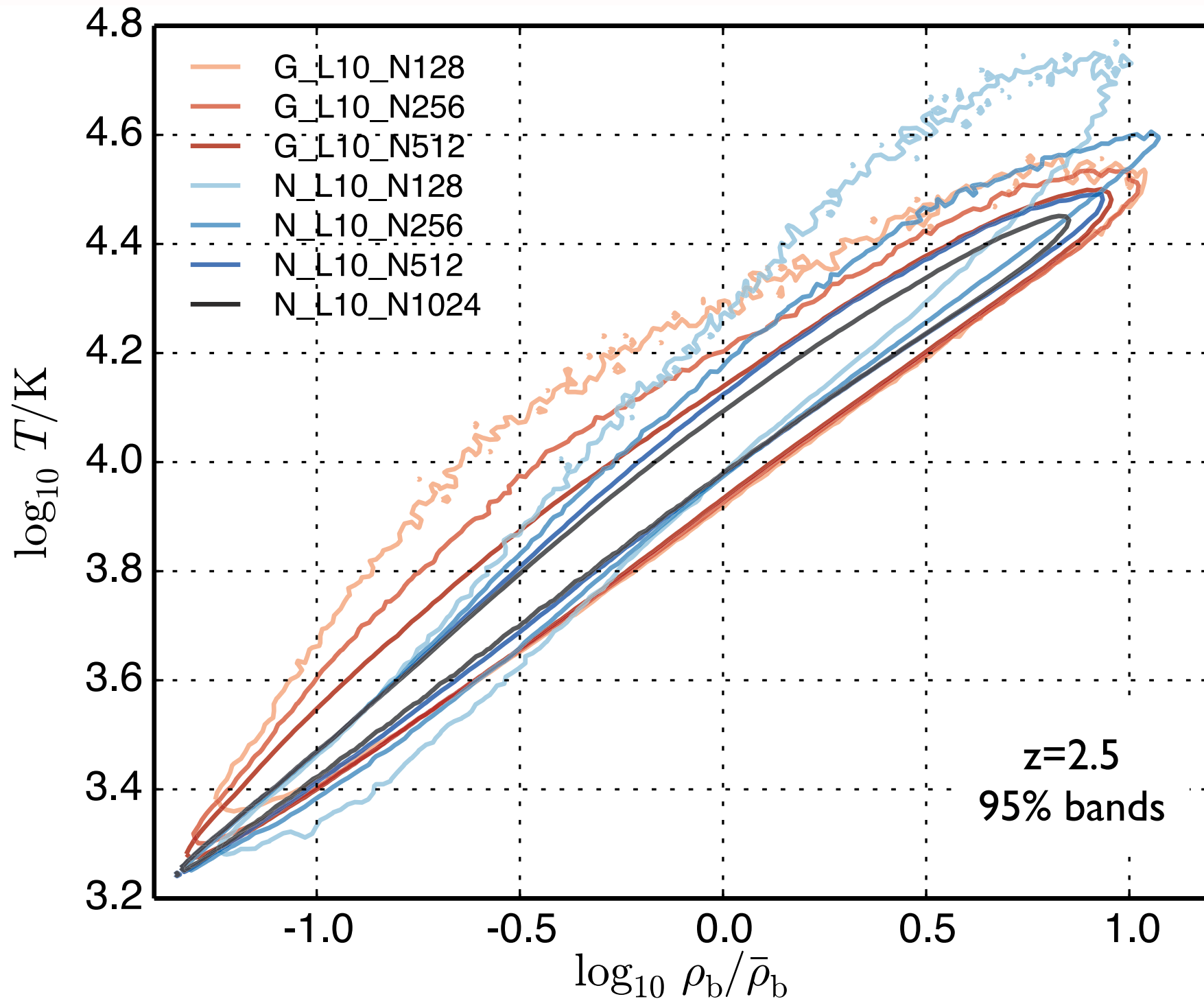
# "Equation of state"



$$T = T_0 \left( \frac{\rho}{\rho_0} \right)^{\gamma - 1}$$

- 4 phases of gas in simulations: "diffuse" relevant for the forest. (Tight density-temperature relation in this regime.)

**Nyx result from Lukic et al. 2014**
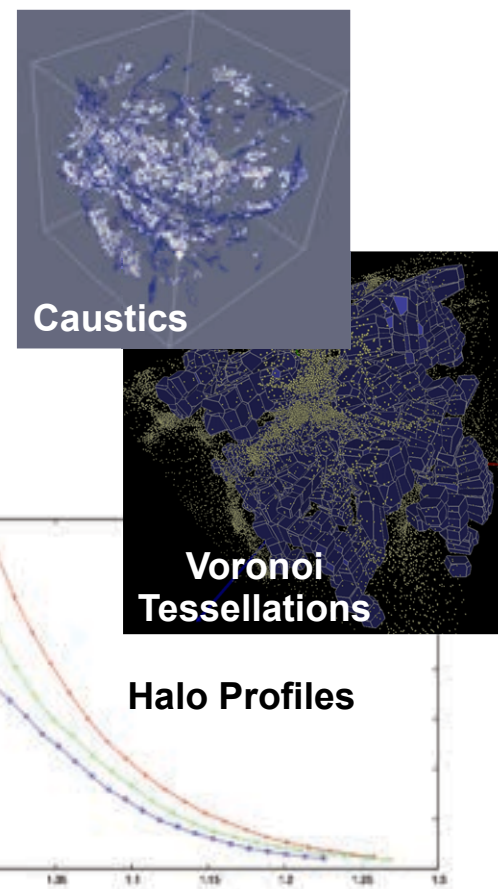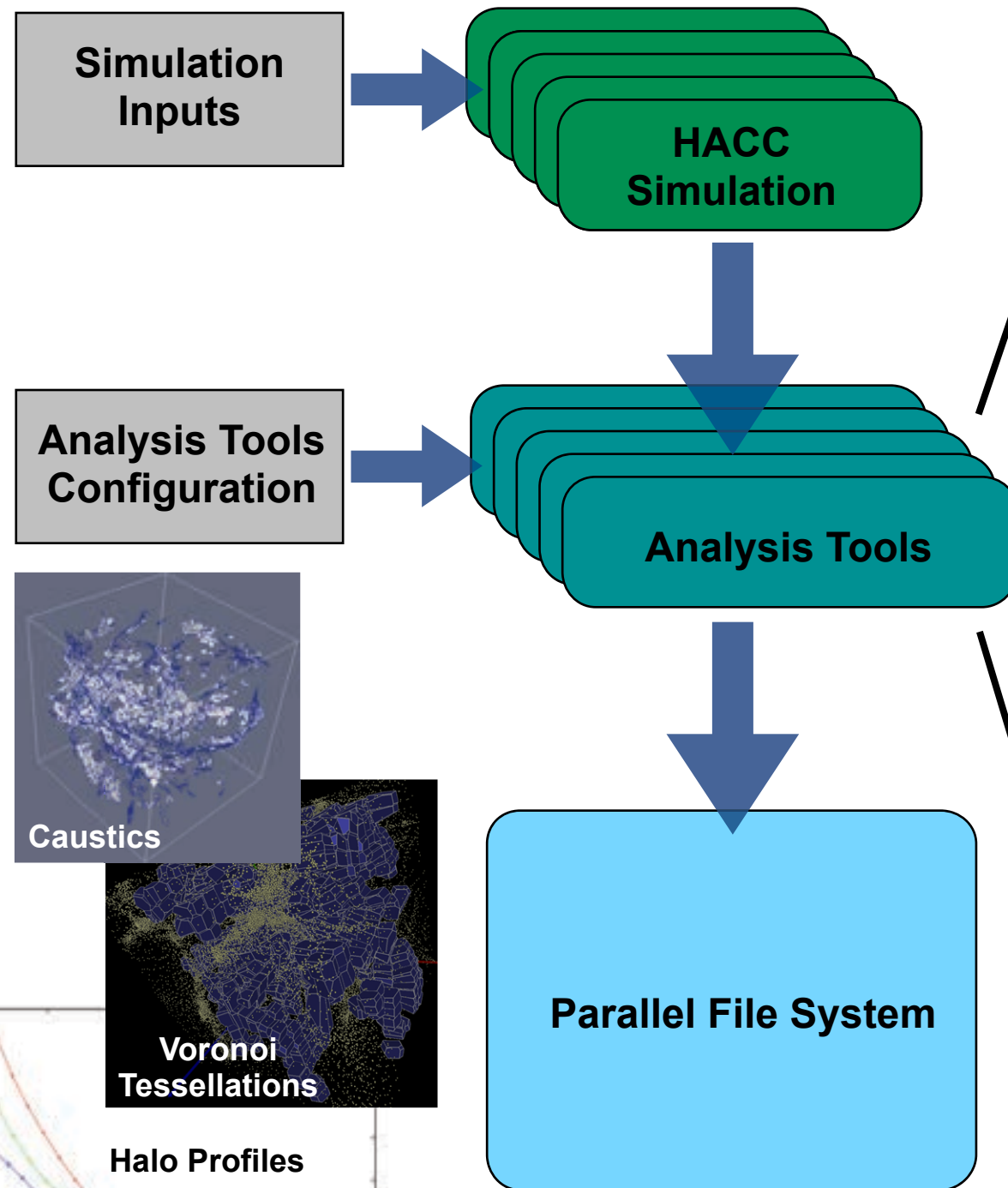
# Density - temperature



- SPH (Gadget) vs. Eulerian (Nyx) code.

Stark et al. in prep.

# In Situ Analysis

- **Data Reduction:** A trillion particle simulation with 100 analysis steps has a storage requirement of ~4 PB -- in situ analysis reduces it to ~200 TB

- **I/O Chokepoints:** Large data analyses difficult because I/O time > analysis time, plus scheduling overhead

- **Fast Algorithms:** Analysis time is only a fraction of a full simulation timestep

- **Ease of Workflow:** Large analyses difficult to manage in post-processing

Simulation Inputs → HACC Simulation

Analysis Tools Configuration → Analysis Tools

Caustics

Voronoi Tessellations

Halo Profiles

Parallel File System

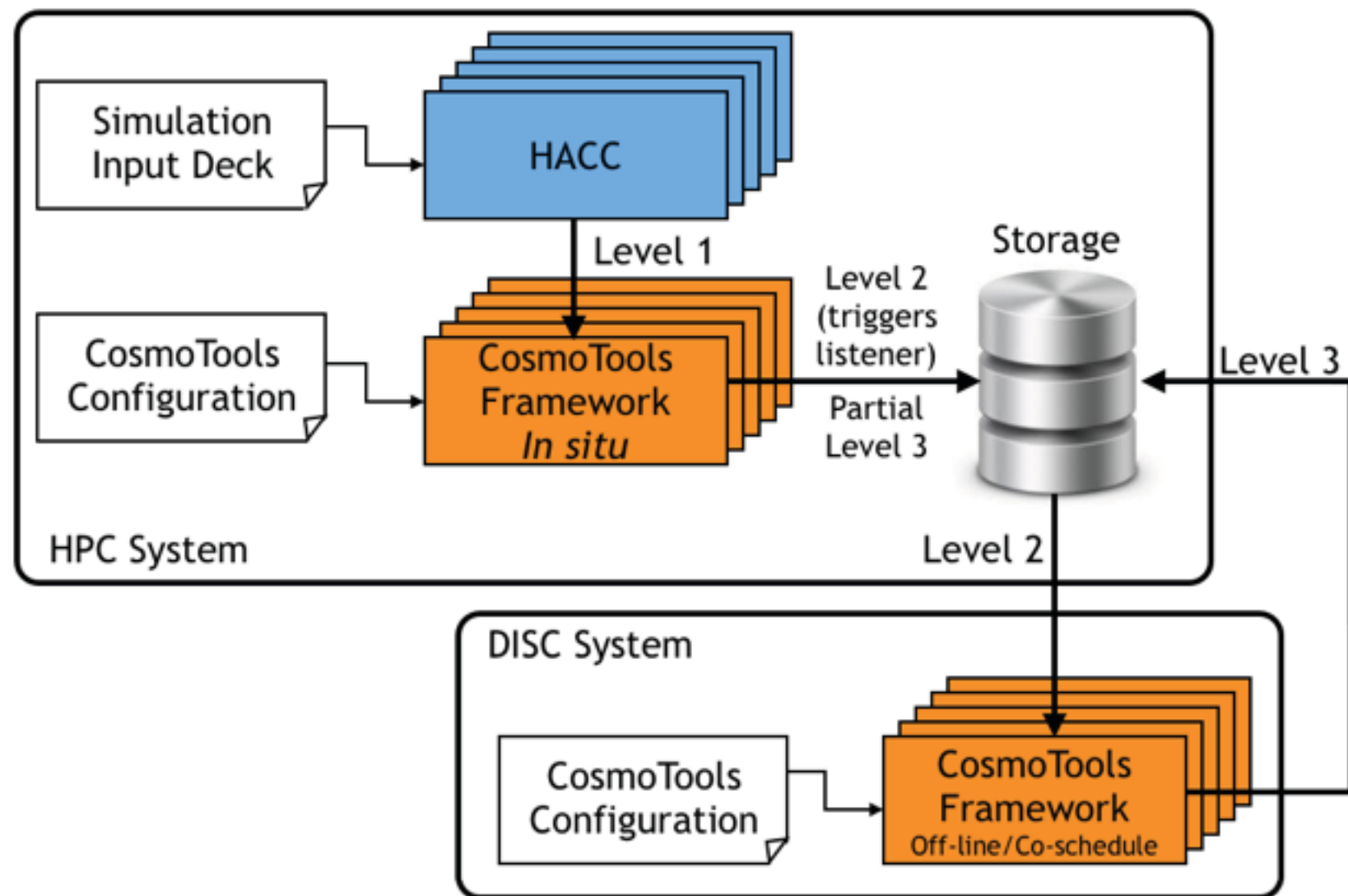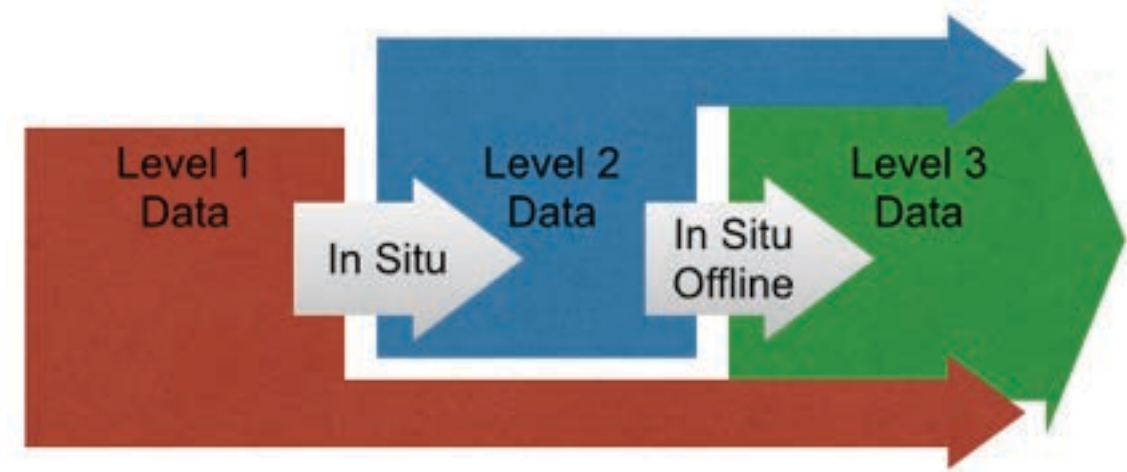*k*-d Tree Halo Finders

Voronoi Tesselation

Merger Trees

N-point Functions

Predictions go into Cosmic Calibration Framework to solve the Cosmic Inverse Problem
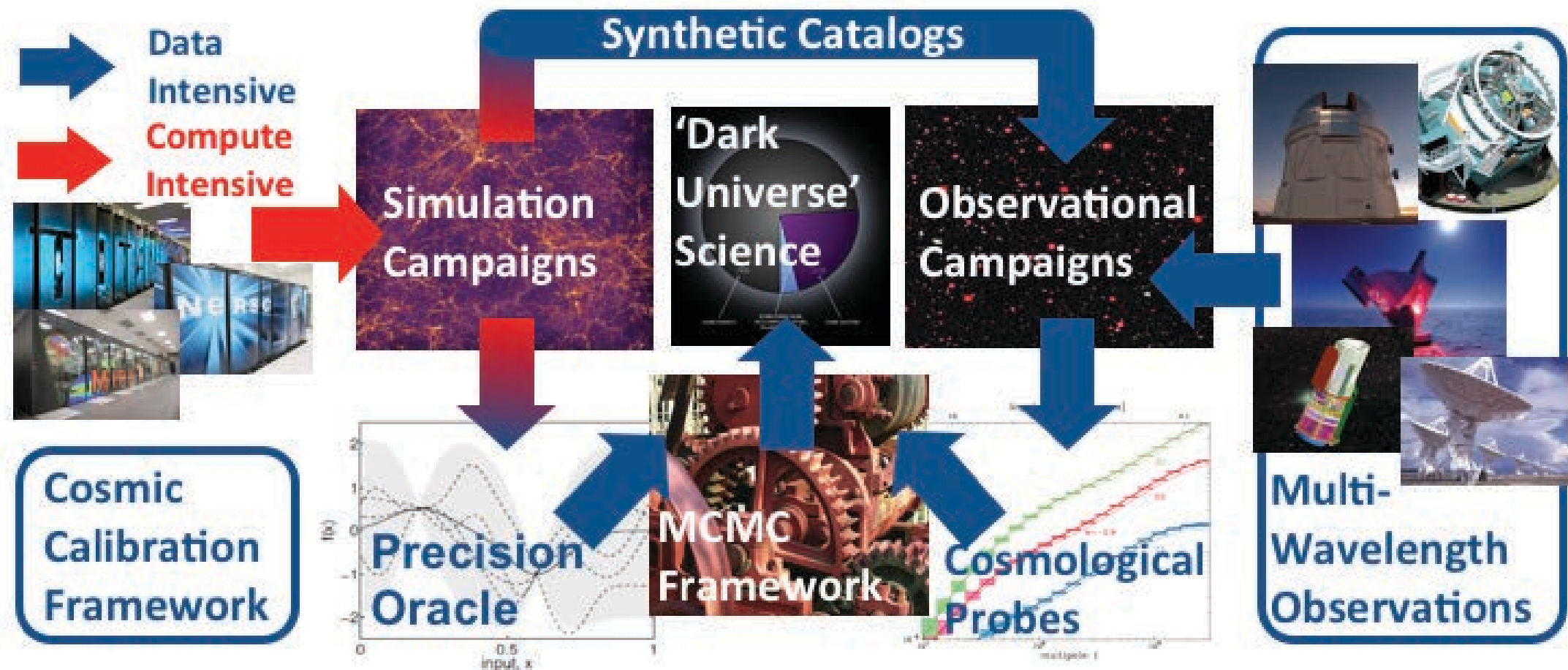
# In Situ Analysis and Co-Scheduling

- **Analysis Dataflows:** Analysis dataflows are complex and any future strategy must combine elements of in situ and off-line approaches

- **CosmoTools Test:** Test of coordinated off-line analysis ("co-scheduling")

- **Portability:** Analysis routines implemented using PISTON (part of VTK-m, built on NVIDIA's Thrust library)

- **Example Case (Titan):** Large halo analysis (strong scaling bottleneck) offloaded to alternative resource using a listener script that looks for appropriate output files



**Sewell et al. 2015, SC15 (to appear)**

# SciDAC-3++: Computing the Sky — Simulation and Analysis for Cosmological Surveys



- **Highlights:**
  - **Next-generation emulators (including covariances), in situ visualization framework, merger trees in CosmoTools, new HACC algorithms for Summit**
  - **Add neutrino capability to Nyx, implement approaches for next-gen architectures, add in situ capability to Nyx**
  - **In situ data reduction schemes for ART, refactoring of ART I/O framework**