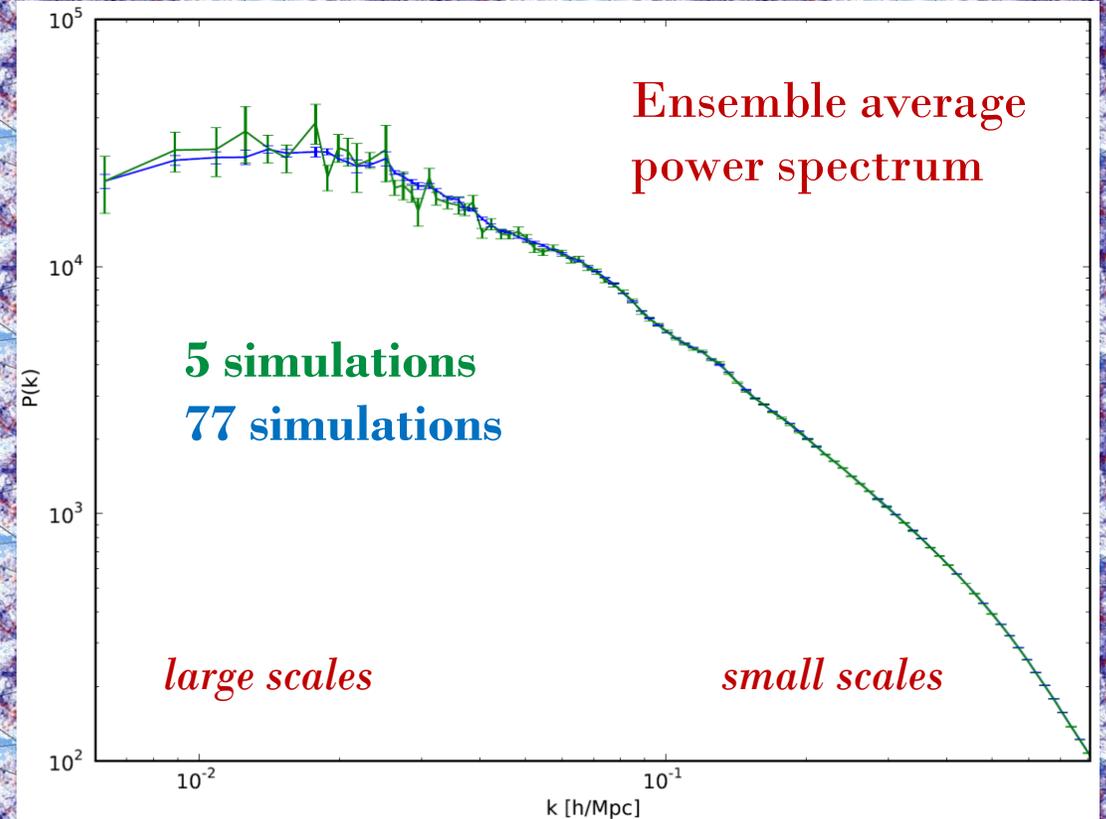# The Indra Simulations on the SciServer Science Platform

## Bridget Falck

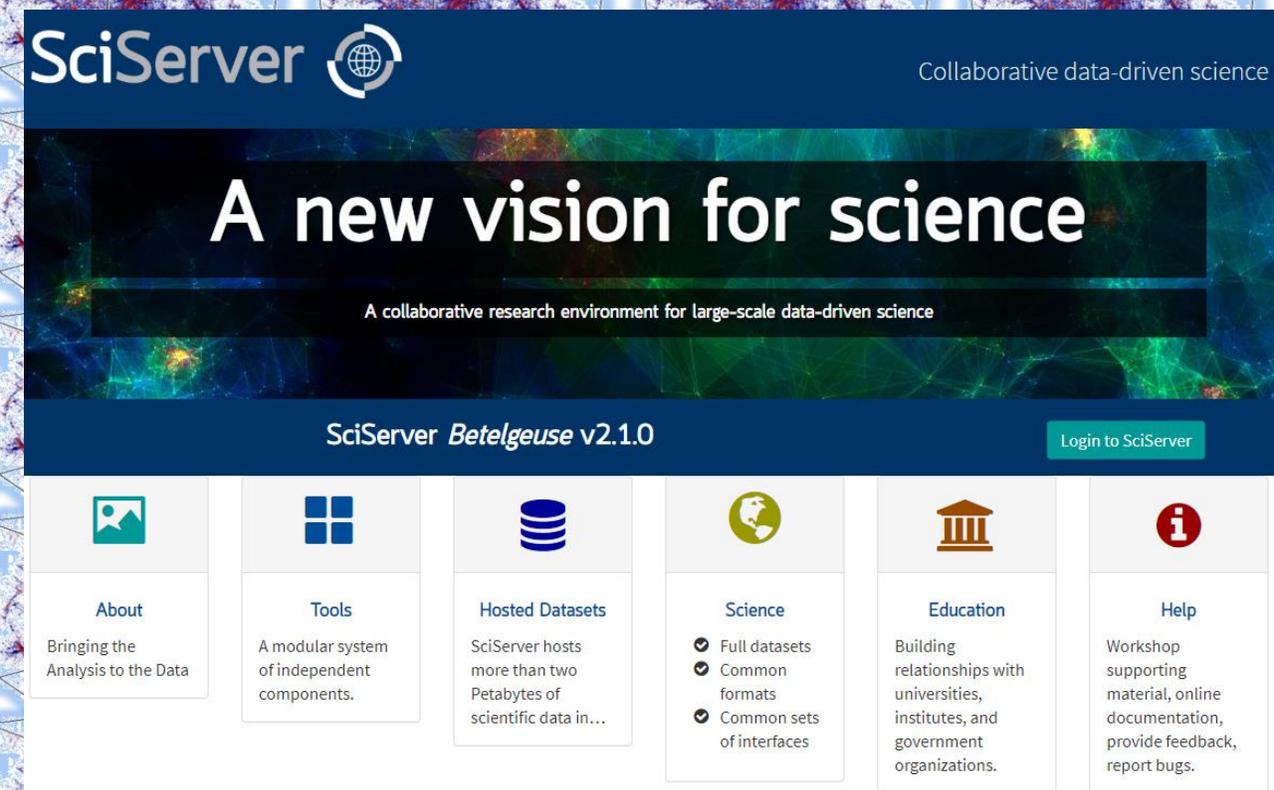### Johns Hopkins University

# Motivation: Cosmic Variance

- Theoretical predictions of large-scale structure require numerical simulations, but we can't simulate our observable Universe exactly, only its statistical properties

- Both simulations and observations have good statistics on small scales but poor statistics on large scales – the cosmic variance limit

- Need to run many simulations with different initial conditions but *same model and parameters*

Ensemble average power spectrum

**5 simulations**
**77 simulations**

*large scales*          *small scales*

# The Indra Simulations

- Suite of 384 simulations with the same cosmology

  - Each a 1 Gpc/h-sided box with $1024^3$ dark matter particles, run with L-Gadget code
  - Output: 64 snapshots of particle positions and velocities, FOF/SUBFIND halo catalogs, and 505 time-steps of Fourier-space density grids

- 750 TB of data, available to the public and *computationally-accessible* via the SciServer

  - Ensemble averages and covariances, conditional and extreme statistics, mock galaxy catalogs and lightcones, etc.
  - Test-bed for new data architectures and analysis tools

- Other simulation suites beyond Indra are needed and being produced

  - Vary cosmological models and parameters, include hydrodynamical effects, etc.
  - Variety of scientific questions and codes means no standard outputs
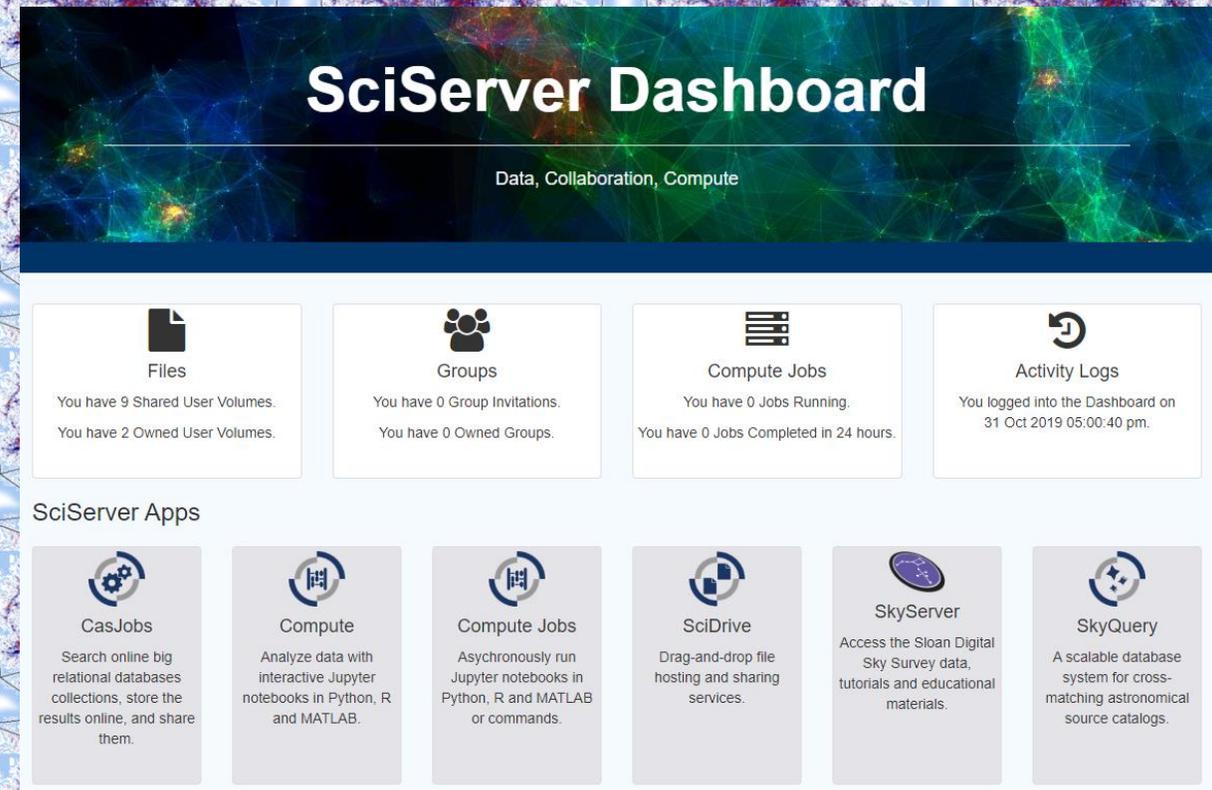
# SciServer Science Platform

# Indra Infrastructure



SciServer

Compute UI | Groups | Files | Apps (SkyServer, CasJobs, …)

Personal Databases and Filesystems

Indra Databases and Files

Compute Domains
Docker containers on virtual machines
Interactive and Batch

Indra Databases

Indra Data Files

FileDB
Distributed filesystem and compute cluster
Dask parallel python
~200 TB of Indra

Data-Scope
Peta-scale storage with high throughput
Permanently store full PB of Indra
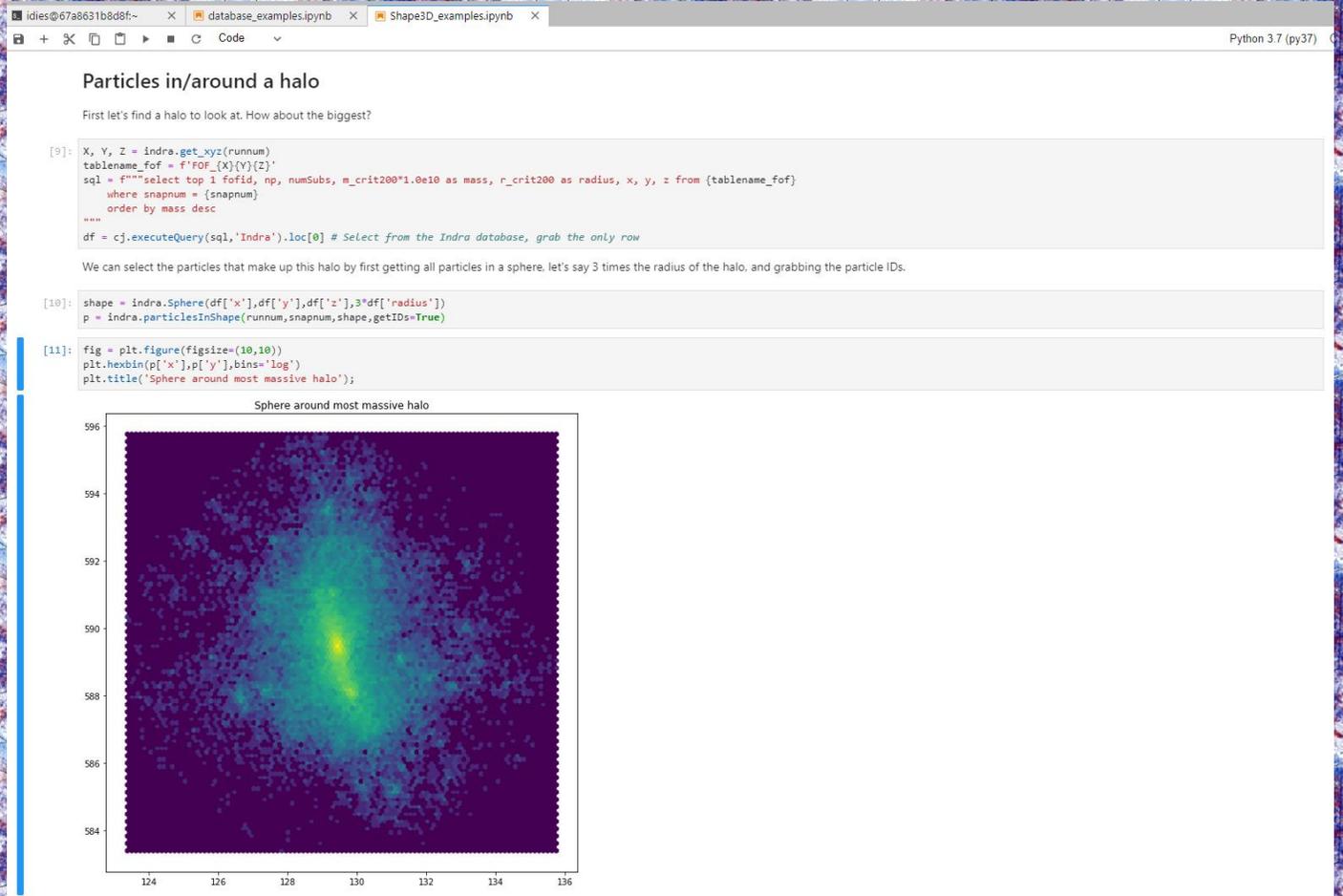
- Permanent read-only storage connected to compute domain

- Distributed filesystem for parallel computation

- Leverage relational databases

- Accessible through SciServer with its collaboration tools and hosted astrophysical datasets
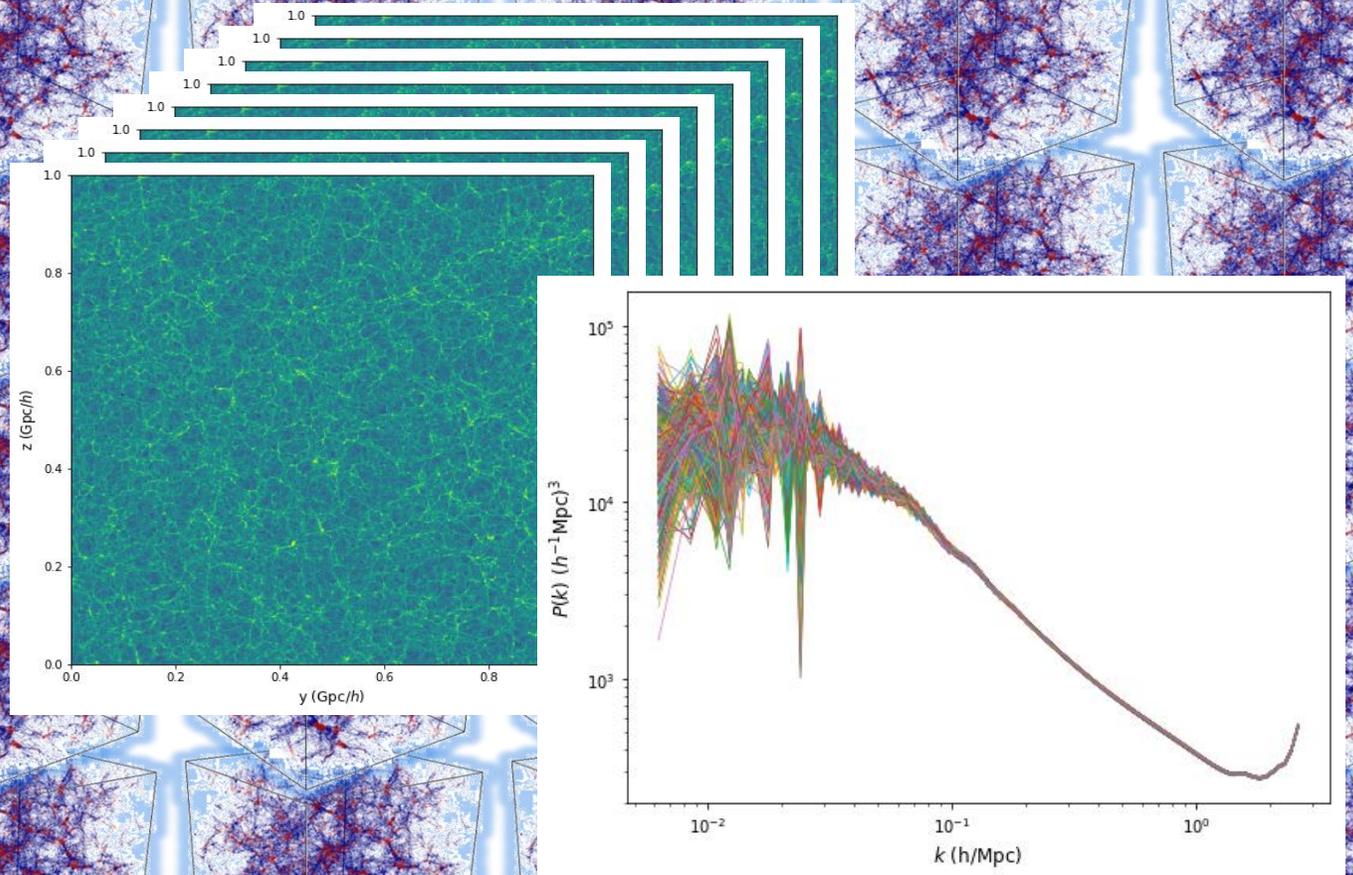
# The indra-tools software library

- Python library to read and interact with data
  - Don't assume expert users
  - Make it easier for experts
  - Hide the file system

- Example notebooks
  - Make it easy to get started
  - Show sample database queries
  - Explain advanced features (e.g., Shape3D)

# Infrastructure that enables heavy computation

- FileDB distributed filesystem with a Dask parallel python cluster

- 448 Cloud-In-Cell density grids calculated in 2 hours!

- 481 billion particles total

- Still testing different use cases and building job submission capability

# Discussion Questions

- What hardware and technologies are required to host large public data sets and make them computationally accessible?

- What are the unique requirements or challenges of hosting simulated data vs. observational archives?

- When we plan archives for large missions, how do we ensure that we don't leave theory behind?

- …