# MyMergerTree:
# A Cloud Service For Creating And Analyzing Galactic Merger Trees

Sarah Loebman, Jennifer Ortiz, Lee Lee Choo, Laurel Orr, Lauren Anderson, Daniel Halperin, Magdalena Balazinska, Thomas Quinn, and Fabio Governato

UM ASTRONOMY, UW CSE,  UW ASTRONOMY

# A Paradigm Shift in Science

Standard:

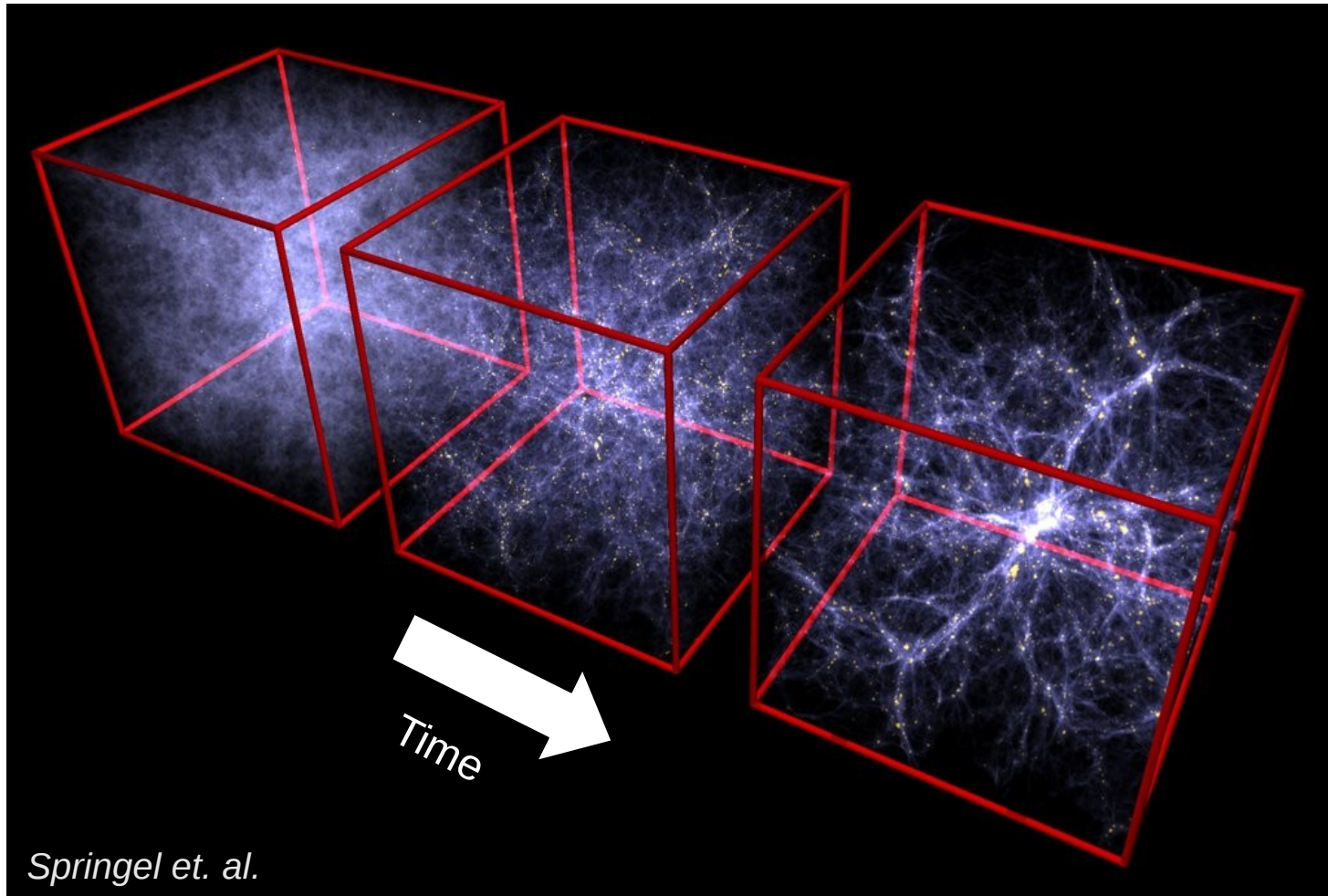What data do I have to collect to (dis)prove a theory?

Data-driven:

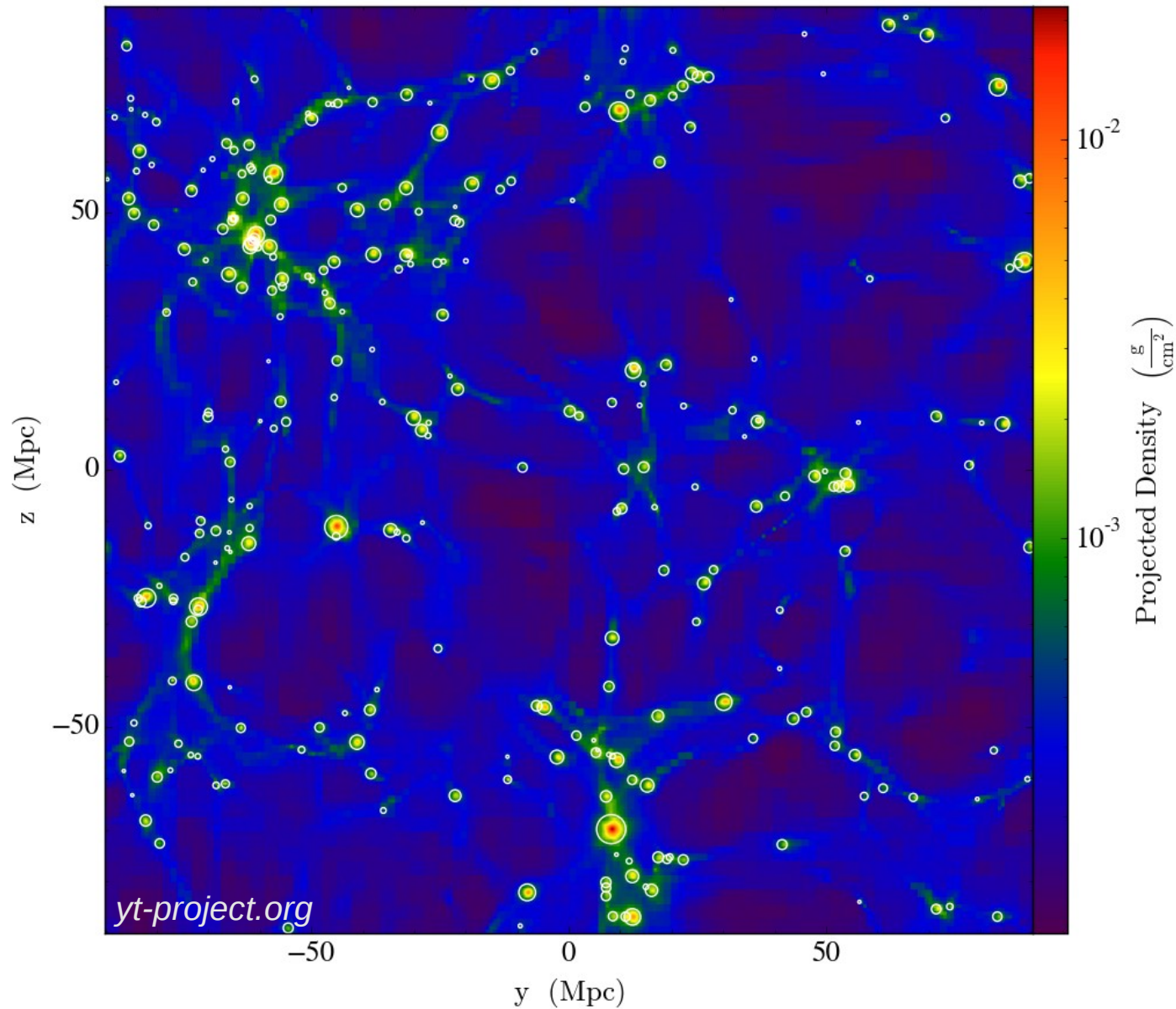What theories can I test given the data that I already have?

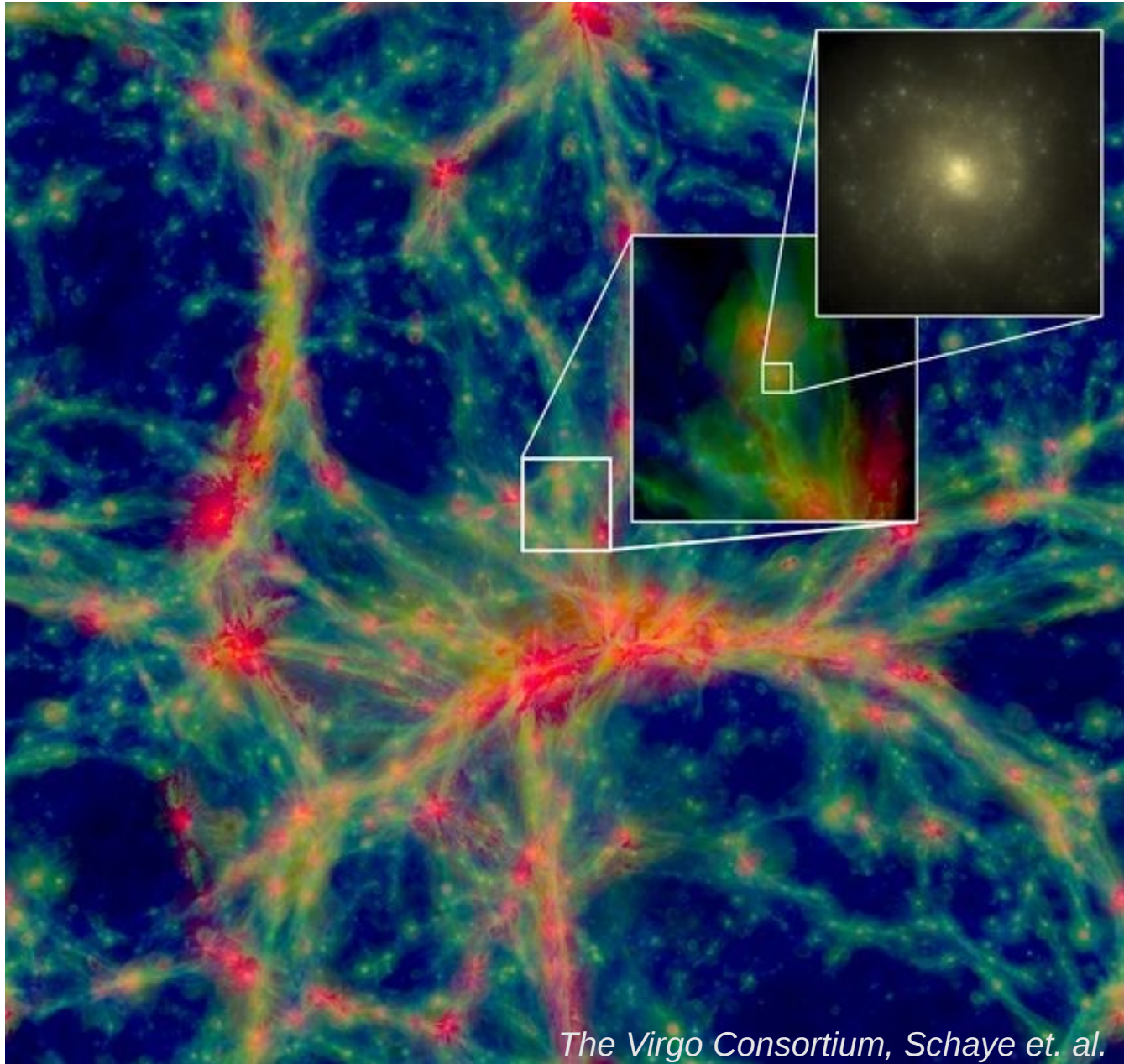Ex: The Sloan Digital Sky Survey (SDSS)

The Millennium Simulation

## Simulated Cosmological Volume



Time

*Springel et. al.*

yt-project.org

**4**

The Virgo Consortium, Schaye et. al.

# Merger Trees



halo
(galaxy)

PRESENT
DAY

PAST

TIMESTEP

4

3

2

1

What tool can generate this structure from the data?

Need a data service can match billions of particles across time

Visualization fast to load, platform independent, easy to use/share

Many Data Services          Many Cloud Services



Myria

**Myria**

- ***In the Cloud***: A RESTful Query-as-a-Service platform

- ***Expressive***: A compiler framework for multiple iterative RA-based languages

- ***Efficient***: A parallel, shared-nothing, iterative execution engine

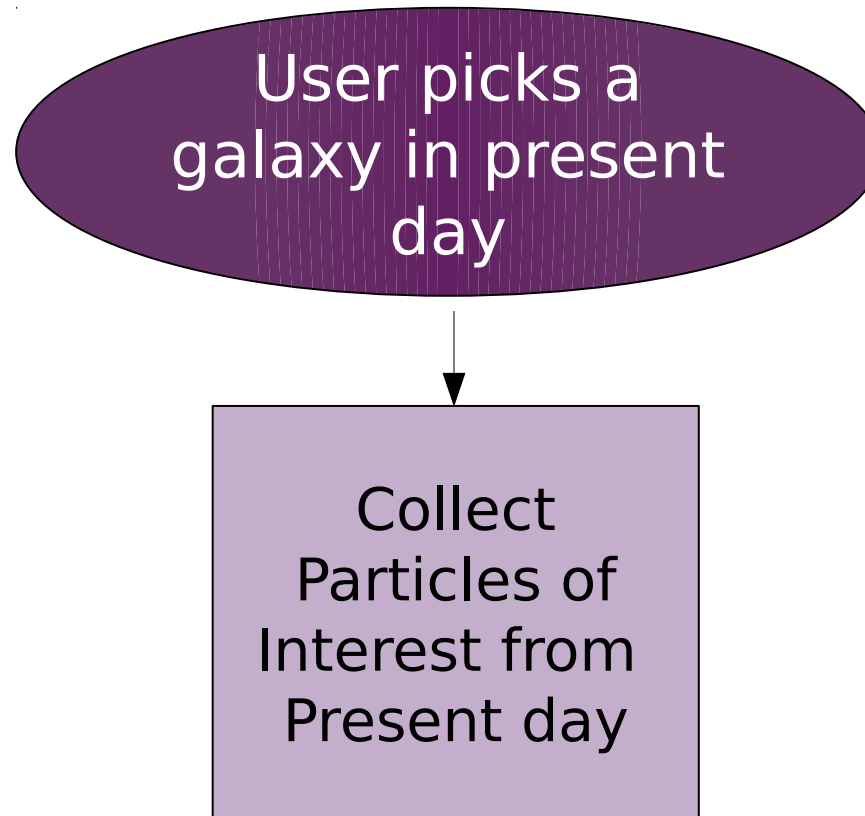# *Myria in the browser*

# Generating a Merger Tree in Myria

# Demo

Challenges:

   - Expressing scientific problems declaratively

   - Physical Tuning for high performance

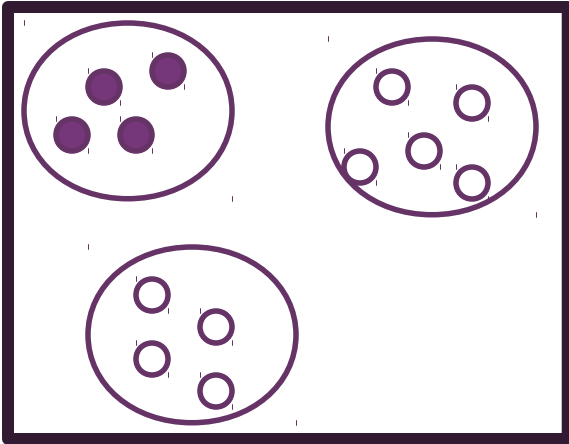   - Visualizing naturally and easily

# How to generate trees?

User picks a galaxy in present day

Collect Particles of Interest from Present day

Look earlier in the simulation….
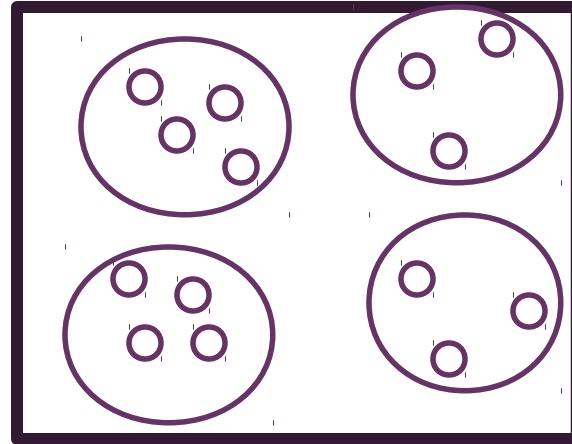Where are the particles at previous timesteps?

# Select *ParticlesOfInterest*

**SELECT** s.iOrder, s.mass, s.type, s.grp

**FROM** Snapshot1818 s  *-- present day snapshot*

**WHERE** s.grp = 'user selection'
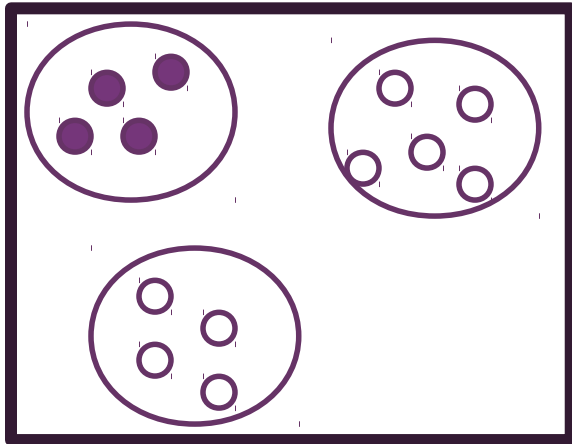
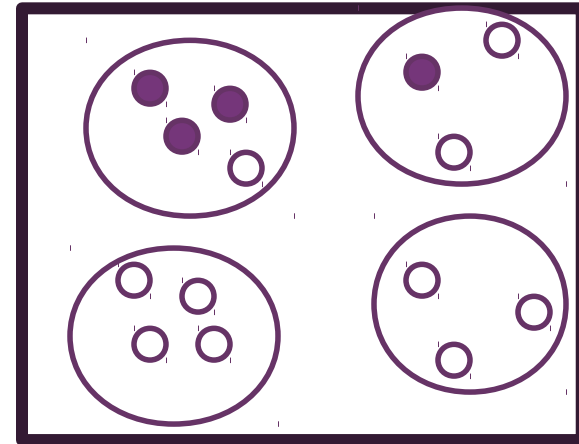Snapshot 1818



Snapshot 1745

# Join across time *AllParticlesTable*

**SELECT** i.iOrder, i.mass, i.type, i.time, i.grp

**FROM** *ParticlesOfInterest* i,  Snapshot1745 s
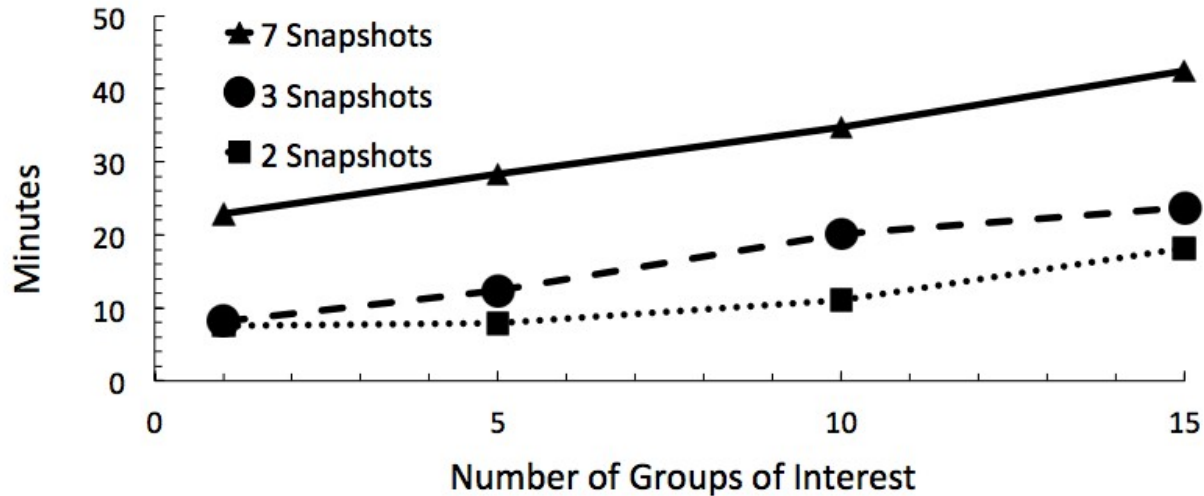
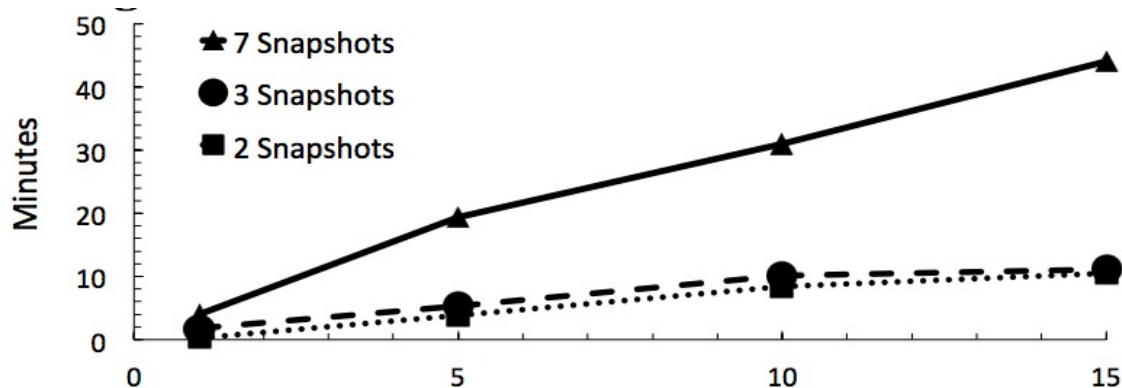**WHERE** i.iOrder = s.iOrder

Snapshot 1818

Snapshot 1745

How data was arranged had a huge impact on query speed!

Not Hashed on Ingest



Hashed on Ingest



16

**d3js.org**

# Conclusions

- Translation of problem easier than predicted
- Physical tuning required most effort
- Specialized visualization tools existed
- Ingesting and validation still needed

How efficiently can we solve other problems?

Exploring dynamic arrangement of data.

# Questions?

# Myria Team

**Magda Balazinska**, Bill Howe, Dan Suciu (faculty)

**Dan Halperin** (postdoc, technical lead)
Victor Almeida (postdoc)
Andrew Whitaker (research scientist)

PhD Students
Shumo Chu
Eric Gribkoff
Jeremy Hyrkas
Paris Koutris
Ryan Maas
Dominik Moritz
**Laurel Orr**
**Jennifer Ortiz**
Emad Soroush
Jingjing Wang
ShengLiang Xu

Undergraduate Students
Lee Lee Choo
Vaspol Ruamviboonsuk