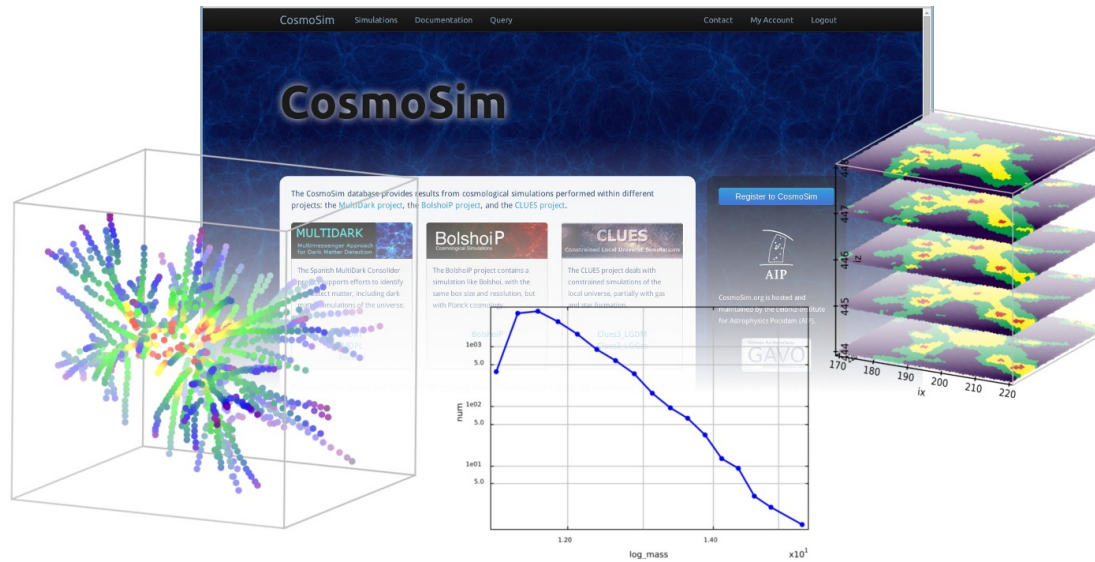# CosmoSim Database



Kristin Riebe
E-Science group @AIP
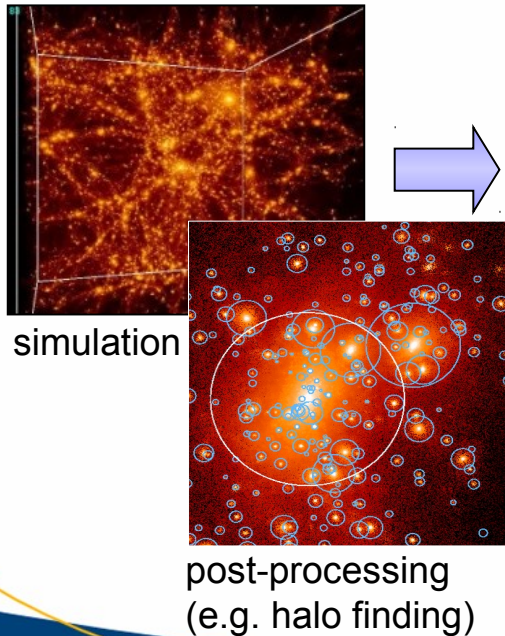
# Outline

- CosmoSim introduction
- Database technology
- Data access
- Data at CosmoSim
- Usage statistics

- Discussion

# CosmoSim

- platform for publishing data from **cosmo**logical **sim**ulations
- first version: MultiDark database, collaboration with Spanish MultiDark project
- database + web query interface
- free registration, open to anyone



simulation

post-processing
(e.g. halo finding)

data catalogue

https://www.cosmosim.org

# Why use databases?

(and not just download the files ...)

- data becomes huge, not quick to download
- let server do most of the calculations
- retrieve only subsets/results, **not** complete catalogues
- Structured Query Language: SQL, quite easy to learn
- examples: sort/filter halos, calculate mass functions, merger trees, follow mass growth of stellar disk, ...

```sql
SELECT * FROM MDR1.FOF
WHERE snapnum=85
ORDER BY mass DESC LIMIT 10
```

10 most massive FOF groups at z=0

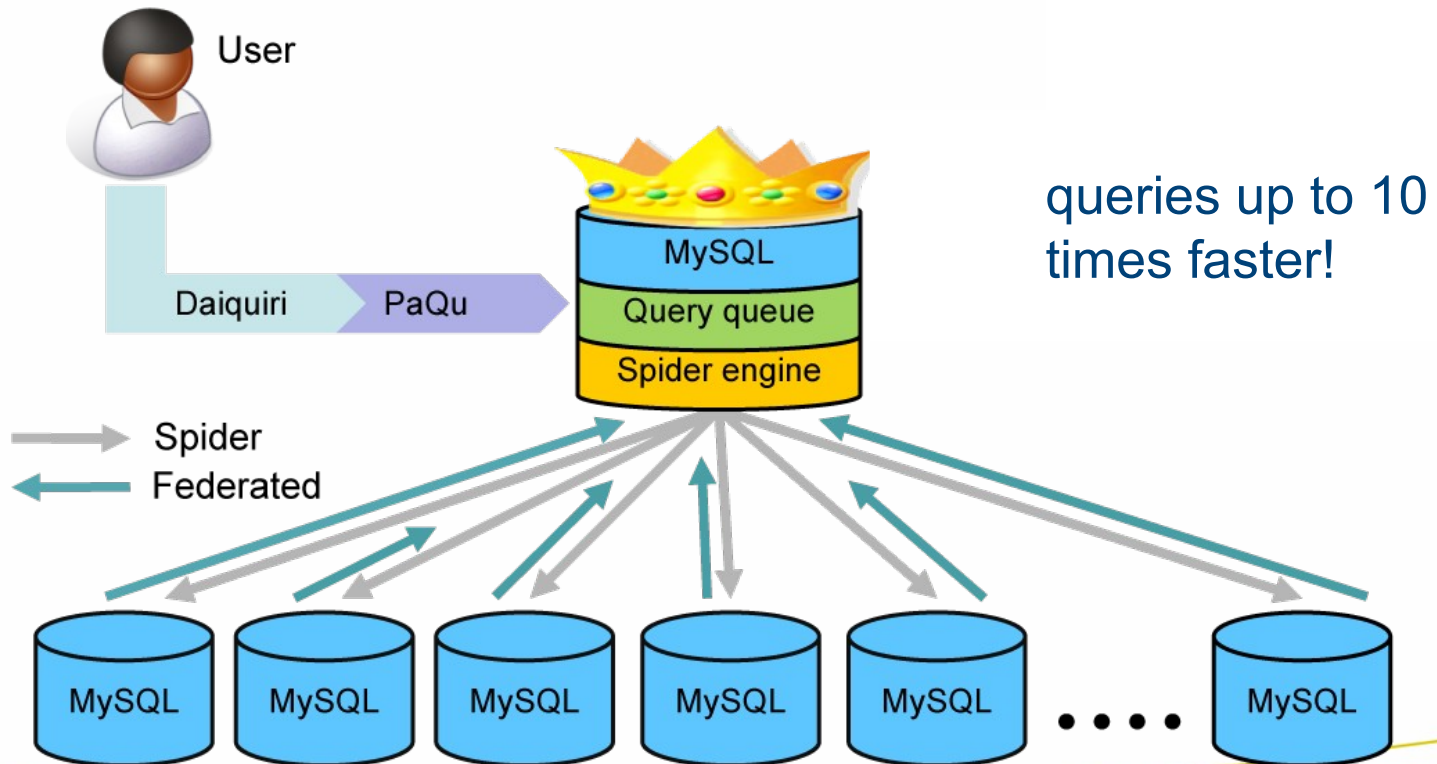=> Just share the query, don't need to share the data!

# Database server

- 10 shard nodes with distributed data, 1 head node
- MariaDB (MySQL variant)
  - MyISAM engine (no transactions => fast)
  - Spider engine for distributed queries
- open source
- own developments in E-Science group:
  - see http://github.com/aipescience
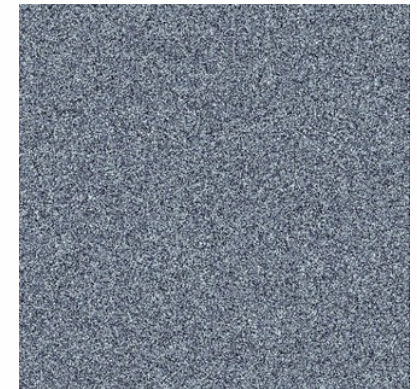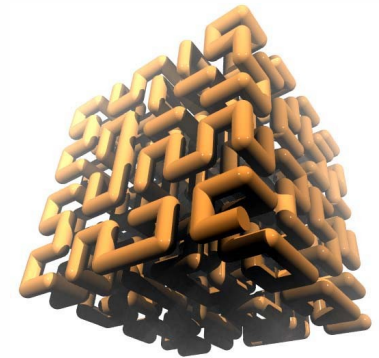    and http://github.com/adrpar

# Spider engine

- data tables partitioned, distributed over 10 nodes engine
- PaQu reformulates queries, head node sends them to nodes
- head node collects data via federated table

queries up to 10 times faster!
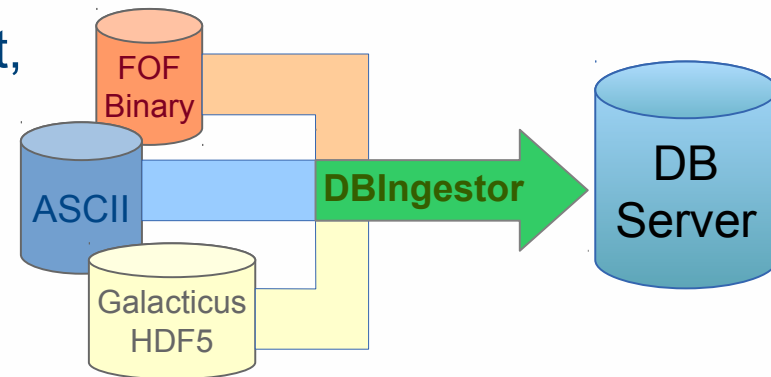
# Further MySQL plugins

- **C-library libhilbert**
  - For creating indexes of space-filling Peano-Hilbert curve in up to 20 dimensions

- **MySQL sprng**
  - Based on The Scalable Parallel Random Number Generators Library (SPRNG, www.sprng.org)
  - Implements several random number generators
  - Better random sampling for large numbers than with built-in function
  - useful for e.g. extracting a random sample of halos/galaxies

see https://github.com/adrpar/

# Upload: DBIngestor

- Data has variety of formats, need transformations
- DBIngestor library: one tool to load them all
  - by Adrian Partl, open source: https://github.com/aipescience/DBIngestor
  - adjustable to any database server
  - write own file readers (e.g. FofIngest, see https://github.com/kristinriebe)
  - apply converters during ingestion
    - e.g. unit conversion, type conversion (int/real), adding identifiers, grid indexes
  - apply asserters (not nan, inf, null etc.)
  - => transform and upload in one go
  - => easier to preserve the workflow for later reference

# Database access: webinterface

- Daiquiri web application
  - http://escience.aip.de/daiquiri
  - developed by J. Klar and A. Partl
  - modular, highly customizable
  - using PHP, Zend-framework
  - authentication, query interface
  - wordpress integration for documentation
  - open source
  - also used for other projects (databases in Madrid, Gaia at AIP)
  - use SAMP for sending results to VO clients (Topcat) directly from webpage
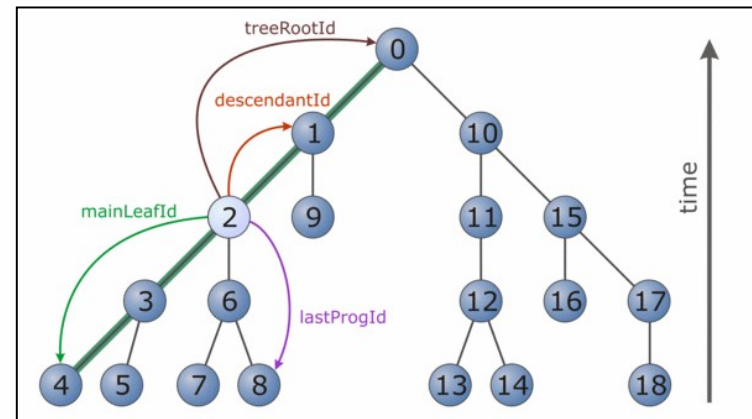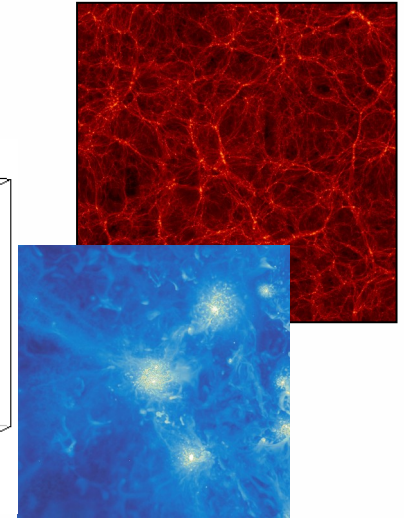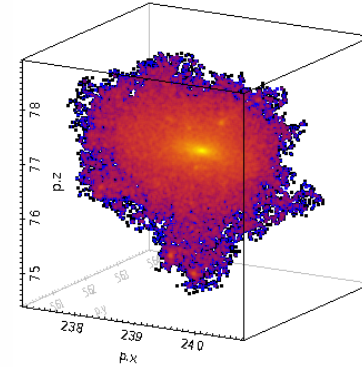
# Database access

- UWS
  - universal worker service, virtual observatory standard
  - interfaces to create, execute, abort or delete query jobs
  - write scripts for submitting many jobs at once
- different tools available:
  - **httpie/curl/wget:**

    ```
    http --auth <username>:<password> --print b GET
    https://www.cosmosim.org/uws/query
    ```
  - **uws-client** (https://github.com/aipescience/uws-client):
    - python client, supports latest features

      ```
      uws --host https://www.cosmosim.org/uws/query --user
      <username> --password <password> list
      ```
  - **cosmosim**-package for „**astroquery**" by Austen Groener:
    - https://github.com/astropy/astroquery/cosmosim
      (Beware: needs to be updated!)
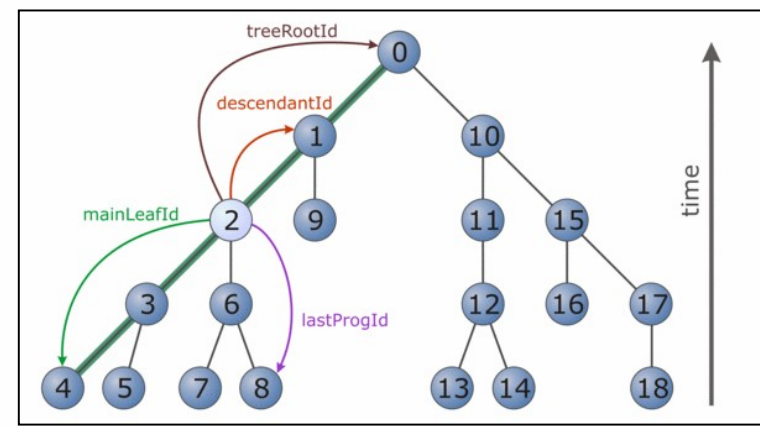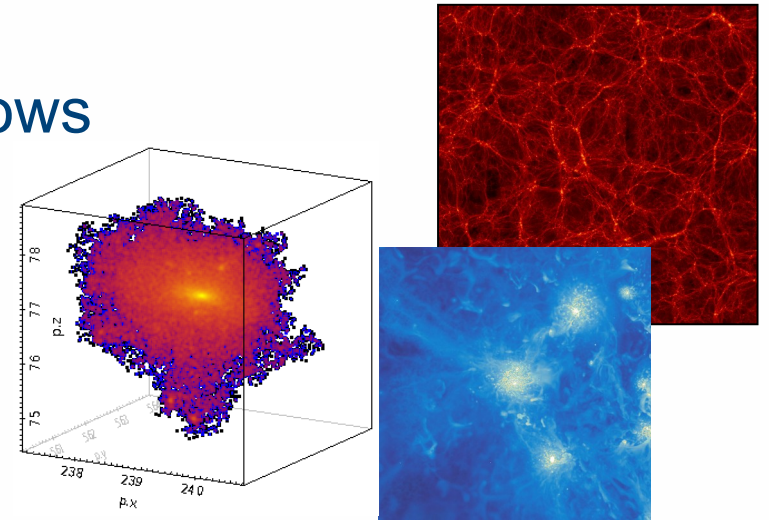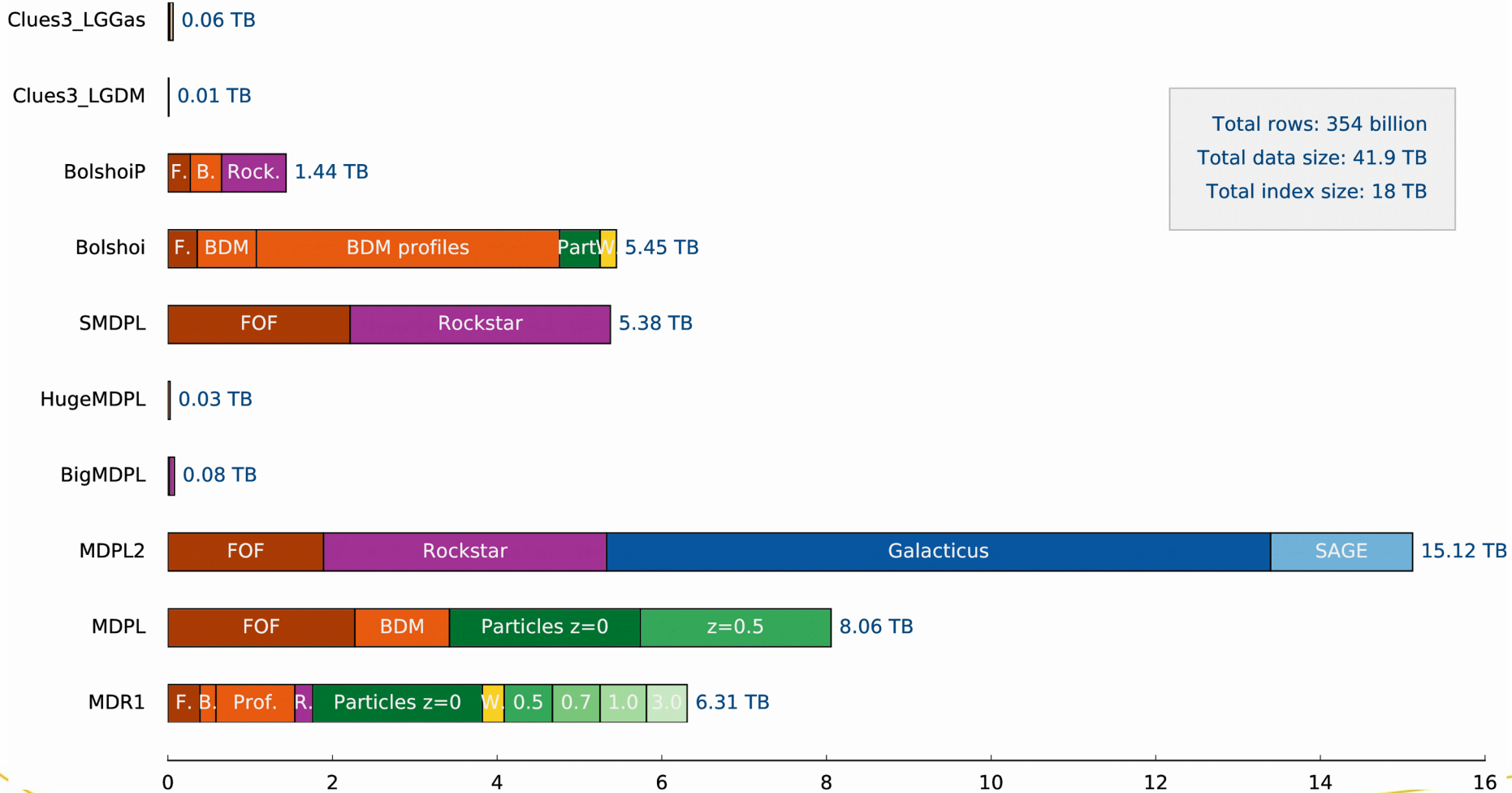
# Data at CosmoSim

- about 40 TB data, ~ 350 billion rows
- 10 simulations
- box sizes: 64 Mpc/h -- 4 Gpc/h
- data products:
  - simulation particles
  - density fields
  - dark matter halo catalogues
  - dark matter profiles
  - merger trees



depth first order for fast retrieval of merger trees

# Data at CosmoSim

- about 40 TB data, ~ 350 billion rows
- 10 simulations
- box sizes: 64 Mpc/h -- 4 Gpc/h
- data products:
  - simulation particles
  - density fields
  - dark matter halo catalogues
  - dark matter profiles
  - merger trees

  + galaxy data
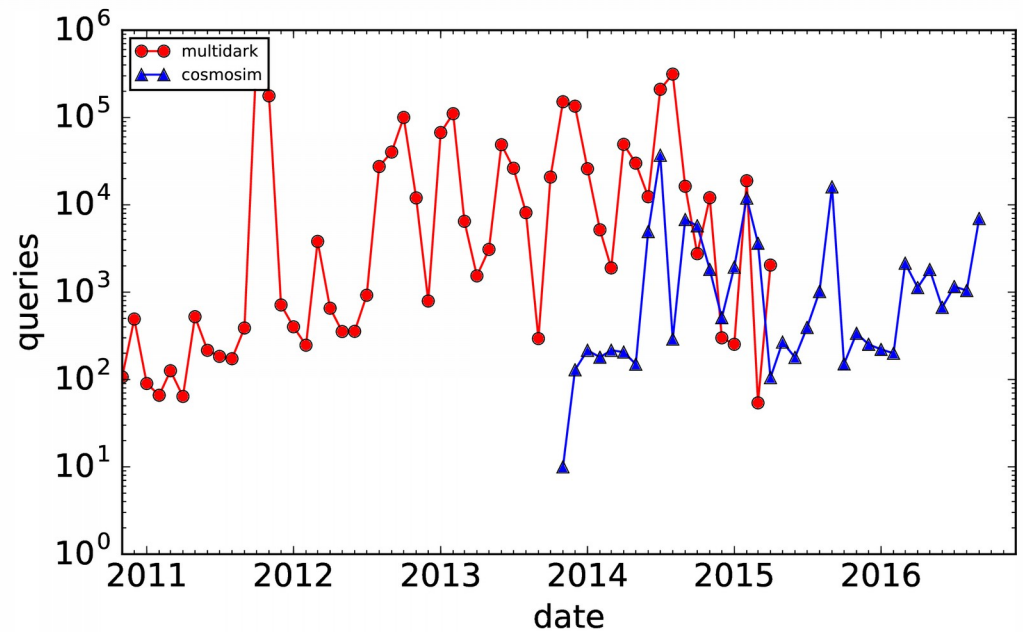


depth first order for fast retrieval of merger trees
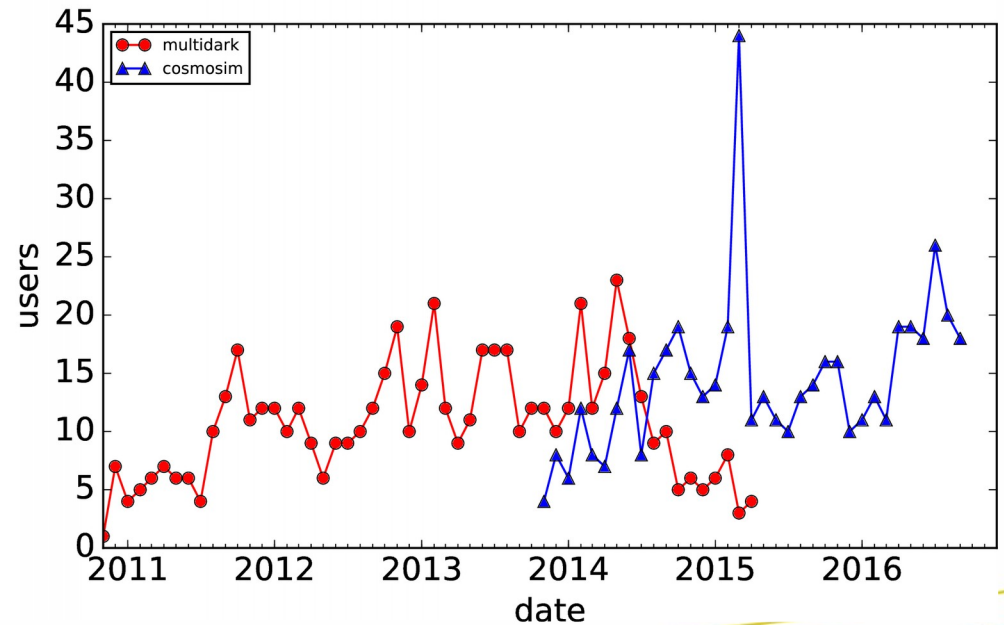
# Data at CosmoSim



Clues3_LGGas — 0.06 TB

Clues3_LGDM — 0.01 TB

BolshoiP — F. B. Rock. — 1.44 TB

Bolshoi — F. BDM BDM profiles PartW. — 5.45 TB

SMDPL — FOF Rockstar — 5.38 TB

HugeMDPL — 0.03 TB

BigMDPL — 0.08 TB

MDPL2 — FOF Rockstar Galacticus SAGE — 15.12 TB

MDPL — FOF BDM Particles z=0 z=0.5 — 8.06 TB

MDR1 — F. B. Prof. R. Particles z=0 W. 0.5 0.7 1.0 3.0 — 6.31 TB

0   2   4   6   8   10   12   14   16

Total rows: 354 billion
Total data size: 41.9 TB
Total index size: 18 TB

13

# User statistics

- Number of queries per month

- Number of (unique) users per month

  400 users registered, 178 with at least 1 query

# Most wanted data on CosmoSim

- MDR1: 640,000 queries; 154 users
- Bolshoi: 34,000 queries;   89 users
- MDPL2:   6,500 queries;   34 users

- Particles:       79,000 queries; 82 users
- FOFParticles: 37,500 queries; 43 users
- FOF:             53,000 queries; 107 users
- BDM:             23,000 queries;   86 users

# Demo: querying data from CosmoSim

# Discussion

- Galaxy data:
  - When will final SAG version be available?
  - Which other formats? (Galform?)
  - Time plan for publishing galaxy data? Papers?
  - Generate DOIs for citations?

- Which simulations next?
  - SMDPL?
  - Time frame?

- New data:
  - need sample data with description + sample read routine
  - estimated data volume
  - time scale for intended publication
  - updates on database are slow (depending on data volume ...)
  - => upload data already into finalized format