

COSMO HUB on Hadoop

Interactive analysis and distribution of cosmological data

J. Carretero, P. Tallada, J. Casals, M. Caubet, C. Acosta-Silva, R. Cruz, N. Tonello, F. Torradeflot, M. Eriksen, J. Delgado, C. Neissner, V. Acín, M. Delfino, S. Serrano and P. Fosalba

How it started

- CosmoHub was created to share data from cosmology projects
- It was built on top of a PostgreSQL relational DB
- **PostgreSQL wasn't scaling** as needed as data volume grew
 - Indices were not used in large datasets
 - **Most queries lasted several hours**
 - Changing the database schema was **very slow**
 - Removing data became **very inefficient**
- We knew future catalogs would grow up to 10^9 entries
 - Now we already have catalogs with up to 10^{11} !

Big Data based platform

- Apache Hadoop
 - One of the most **popular** Big Data platforms
 - Open source
 - **Distributed** storage and processing
 - Based on commodity computer clusters
 - **Scalable** from dozens up to thousands of nodes
 - ✓ Performance scales with HW
 - **Failure tolerance**
 - ✓ May use old/refurbished machines → **cheaper**
- Apache Hive
 - **Query** over **massive** data volumes using SQL

Front end

Javascript based using



Infrastructure

Big data platform managed with



Back end

REST API powered by



- 44 nodes + 3 head + UI
- 784 total cores
- 3 TiB total RAM
- 230 TiB HDD

As of February 2018

Some numbers



Experiments in CosmoHub



Conclusions

- Easy **scalability** with old servers
- **Reliability** against possible failures
- **Impressive** performance
 - A **faster system** than the previous version
- Switched focus from batch catalog generation to **interactive catalog analysis**
 - **Histogram** and **heatmap** plots
- Sampling: select a random subset of the catalog to get faster results when exploring the data

