

COSMO **HUB** Science Portal

Jorge Carretero on behalf of the PIC Team

Red de Infraestructuras de Astronomía

Promoviendo sinergias entre grandes observatorios españoles I

Temática: Gestión de datos científicos

22- Oct - 2023, La Palma, Canarias



**Institut de Física
d'Altes Energies**



**Barcelona Institute of
Science and Technology**



**PIC
port d'informació
científica**



**Infraestructuras
Científicas y Técnicas
Singulares**



Ciemat Centro de Investigaciones
Energéticas, Medioambientales
y Tecnológicas



CFP
CIEMAT
física de partículas



**Unió Europea
Fons Europeu
Next Generation**



**GOBIERNO
DE ESPAÑA**



**Plan de Recuperación,
Transformación
y Resiliencia**



**Next Generation
Catalunya**



**Generalitat de Catalunya
Departament de Recerca
i Universitats**

Port d'Informació Científica

Google

- An advanced and multidisciplinary computing center
- Founded in 2003: collaboration between IFAE and CIEMAT
- Team of 24 people (50% scientists - 50% engineers)
 - Contact person for each experiment
 - Agile teams that embed in scientific groups
- Run production data services for large collaborations
 - Collaborative environments, distributed infrastructures
- Flexibility to adapt to evolving needs
 - Integration of data processing tools/methods
 - Trans-disciplinary cross-pollination
 - Infrastructure - experimental can-do attitude



PIC projects support

- Particle Physics:
 - LHC: Spanish primary center & ATLAS analysis facility
 - Neutrinos: T2K & DUNE
- Astrophysics:
 - Gamma-ray astronomy: MAGIC & CTA
 - Cosmology: DES, PAUS & Euclid
 - Gravitational waves: Virgo/LIGO
- Other data-intensive disciplines:
 - material sciences, bioimaging, genomics



MAPA DE INFRAESTRUCTURAS CIENTÍFICAS Y TÉCNICAS SINGULARES



Plan de Recuperación, Transformación y Resiliencia



Centro Nacional de Investigación sobre la Evolución Humana (CENIEH)
Centro de Láseres Pulsados (CLPU)
RES - Calentuda (SCAYLE)



Laboratorio Nacional de Fusión (LNIF)
Red Académica y de Investigación Española (REDIRIS)
RES - Cibeleles (UAM)
RES - Xule (CIEMAT)
NANBIOSIS - CIBER-BBN
ReDIB - BIOIMAC (UCM)
MARHIS - CEHIPAR
R-LRB - LMR
MICRONANOFABS - CT-ISOM
ELECMI - CNME
RLASB - GISA
ReDIB - TRIMA-CHIC
IABA-Centro de Microanálisis de Materiales (CMAM)



CANARIAS
Plataforma Océanica de Canarias (PLOCAN)
Observatorios de Canarias (OOC)
Gran Telescopio CANARIAS (GTC)
RES - La Palma (IAC)
MARHIS - PLOCAN-TS



GALICIA
RES - Finis Terrae (CESGA)
FLOTA - CSIC/IEO
NANBIOSIS - CIBER-BBN



CANTABRIA
RES - Altamira (UC)
MARHIS - GTIM-CCOB



NAVARRA
RES - Urderra (Nasertic)

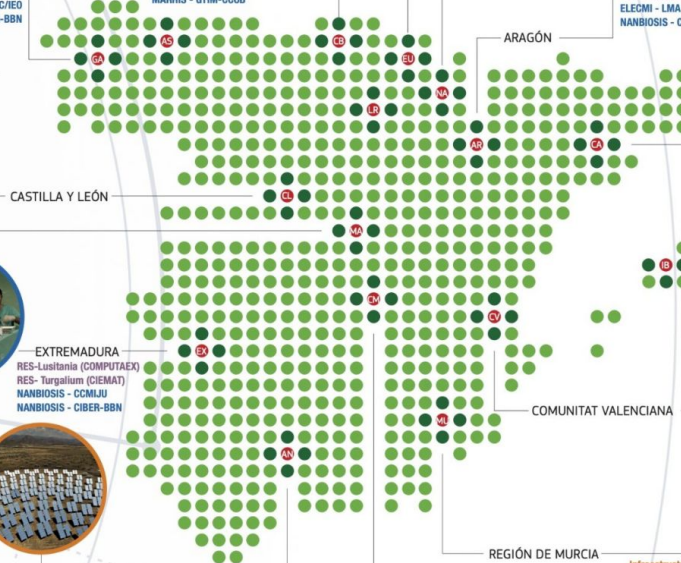


Laboratorio Subterráneo de Canfranc (LSC)
Observatorio Astronómico de Javalambre (OAJ)
RES - Caesaraugusta (UNIZAR)
ELECMI - LMA
NANBIOSIS - CIBER-BBN



Sincrotrón ALBA
RES - MareNostrum y Minotauro (BSC-CNS)
RES - PIC (CIEMAT-IFAE)

MICRONANOFABS - SECNM
MARHIS - CIEM
OmicsTech - CHAG - CRG
OmicsTech - CDS
RLASB - CReSA
NANBIOSIS - CIBER-BBN
ELECMI - UMEAP
FLOTA - CSIC



CASTILLA Y LEÓN



EXTREMADURA
RES - Lusitania (COMPUTAEX)
RES - Turgallium (CIEMAT)
NANBIOSIS - CCMIJU
NANBIOSIS - CIBER-BBN



ANDALUCÍA
Reserva Biológica de Doñana (RBD)
Plataforma Solar de Almería (PSA)
Observatorio Astronómico de Calar Alto (CAHA)
Radiotelescopio IRAM 30M
RES - Picasso (UMA)
IABA - Centro Nacional de Aceleradores (CNA)
ELECMI - DME-UCA
NANBIOSIS - BIONAND



CASTILLA - LA MANCHA
Observatorio de Yeves

ARAGÓN

CATALUÑA

ILLES BALEARS
Sistema de Observación Costero de las Illes Balears (SOCIB)
FLOTA - CSIC/IEO



RES - Tirant (UV)
MICRONANOFABS - NF-CTN
NANBIOSIS - CIBER-BBN
ReDIB - Imaging La Fe



REGIÓN DE MURCIA

Infraestructura para el Cultivo del Atún Rojo (ICAR)
FLOTA - BIO Hespérides



ANTÁRTIDA
Base Antártica Española Juan Carlos I (BAE-JCI)
Base Antártica Española Gabriel de Castilla (BAE-GdC)

"Una manera de hacer Europa"

ICTS CON LOCALIZACIÓN ÚNICA

ICTS DISTRIBUIDAS

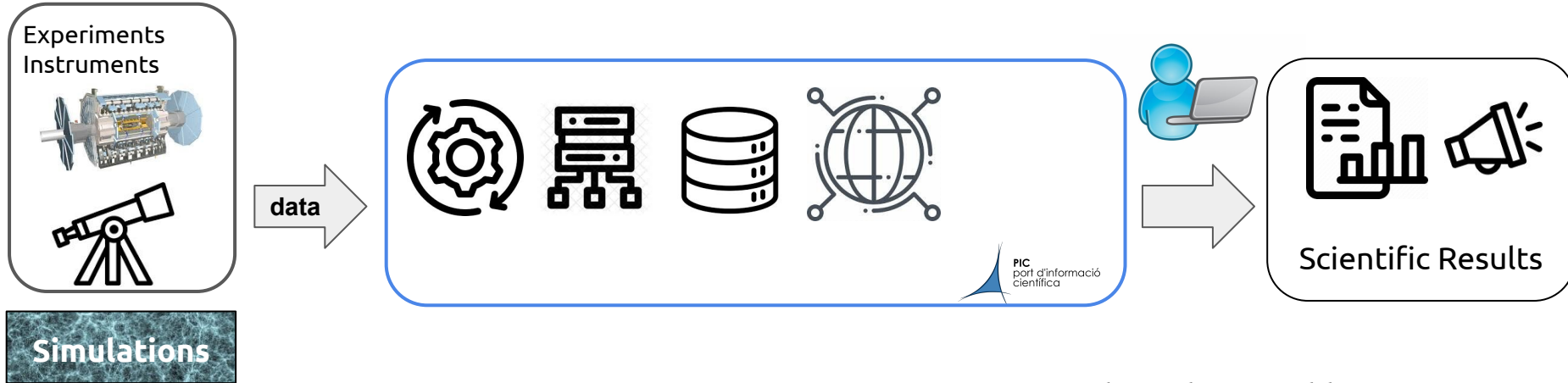
RED DE ICTS

Red Española de Supercomputación

- ICTS distribuida por toda la geografía española
 - 14 nodos interconectados
- Ofrece recursos y servicios de supercomputación, gestión de datos e inteligencia artificial a proyectos científicos y tecnológicos innovadores y de alta calidad
 - Convocatorias competitivas



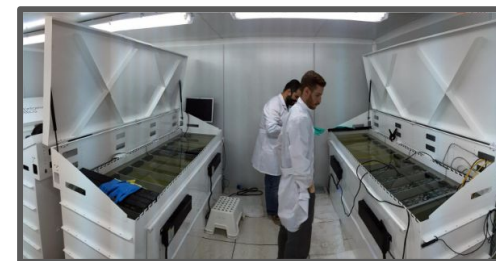
What do we do at PIC?



- Data transfer
- User support so they can process and analyze the data
- Store data and/or results
- Data distribution
- Provide tools to enable science
- Promote the dissemination of the science we work with

- Facilities, ~150 kW IT
 - ~120 kW in 150 m² air-cooled room
 - high efficiency, PUE 1.44
 - ~30 kW in 25 m² liquid immersion cooling system
 - PUE 1.1
- Data processing services
 - Disk - dCache: **20 PB**
 - Tape - Enstore: **63 PB**
 - Computing - HTCondor: **12000 cores**, 18 GPUs
 - Computing - Hadoop: 720 cores, 2.5 PB net storage
- Connectivity
 - 2x100 Gbps to Academic Network
 - **Largest data mover in Spanish academic network: 100PB** in+out per year

IBM TS4500



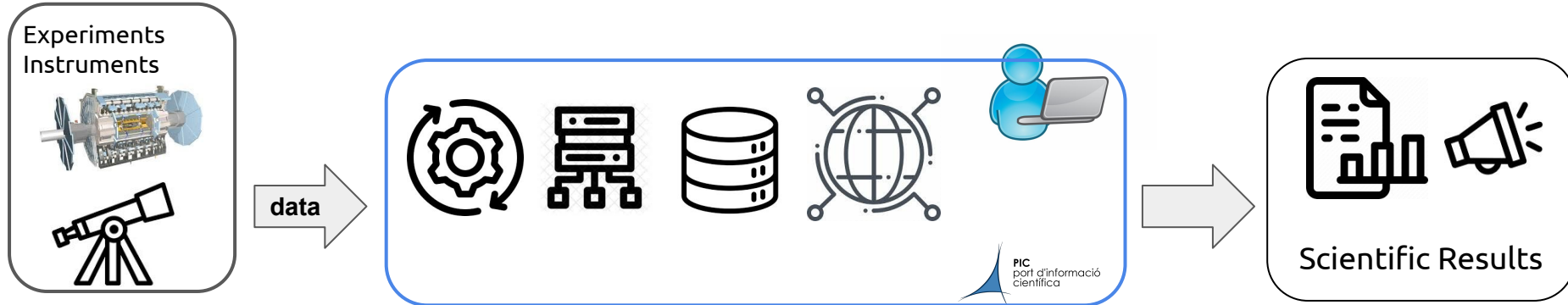
- Computing, data processing and analysis
 - Batch processing through HTCondor
 - Notebooks ecosystem: Jupyter (Dask / *Spark*)
 - CosmoHub
- Mass storage
 - Tape
 - Disk
- Web Services (Gitlab, Wiki, Redmine, Webdav, Monitoring, etc.)
- Consulting support

What do we do at PIC?



- Data transfer
- User support so they can process and analyze the data
- Store data and/or results
- Data distribution
- Provide tools to enable science
- Promote the dissemination of the science we work with

What do we do at PIC?



- Data transfer
- User support so they can process and analyze the data
- Store data and/or results
- Data distribution
- Provide tools to enable science
- Promote the dissemination of the science we work with

Computing key recommendations

- Define and plan the data lifecycle
- Favor distributed architectures
- Ecological awareness
- **Bring user to the data**
 - **open & collaborative environments**
- People are very valuable assets



Jupyter notebook

Server Options

Select custom options for your profile

Memory (RSS)

2 GB

CPUS

1

GPUS

0

User options

Experiment Select your experiment

Start

Notebook

Python 3 (ipykernel) common dcache-tutorial desktop [1] hmf_env VS Code (IDE) [1]

Console

Python 3 (ipykernel) common dcache-tutorial hmf_env master

Other

Terminal Text File Markdown File Python File Show Contextual Help

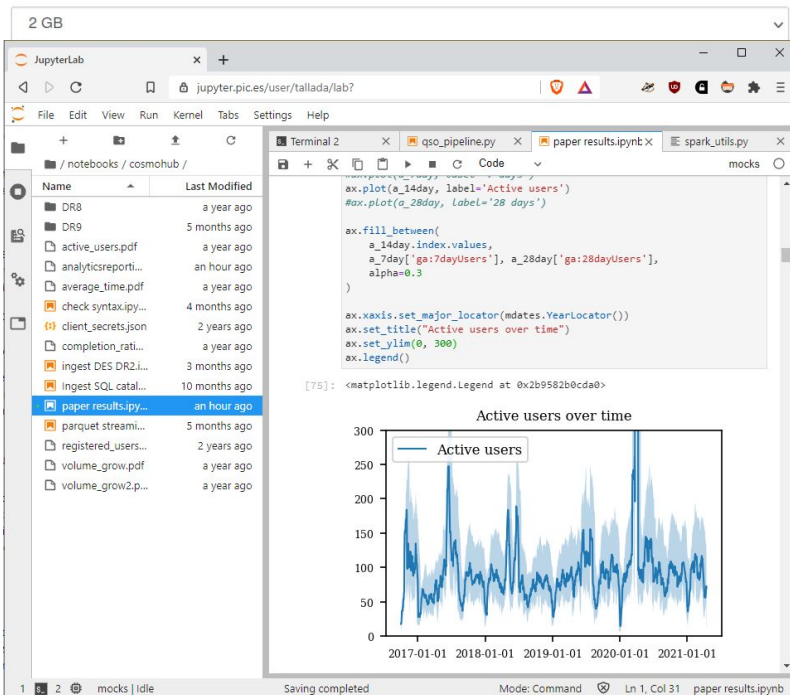
<https://jupyter.pic.es/>

Server Options

Select custom options for your profile

Memory (RSS)

2 GB



```

ax.plot(a_14day, label='Active users')
#ax.plot(a_28day, label='28 days')

ax.fill_between(
    a_14day.index.values,
    a_7day['ga:7dayUsers'], a_28day['ga:28dayUsers'],
    alpha=0.3
)

ax.xaxis.set_major_locator(mdates.YearLocator())
ax.set_title("Active users over time")
ax.set_ylim(0, 300)
ax.legend()
    
```

Active users over time

Active users

2017-01-01 2018-01-01 2019-01-01 2020-01-01 2021-01-01

Notebook

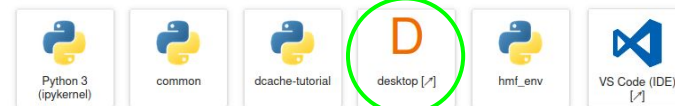
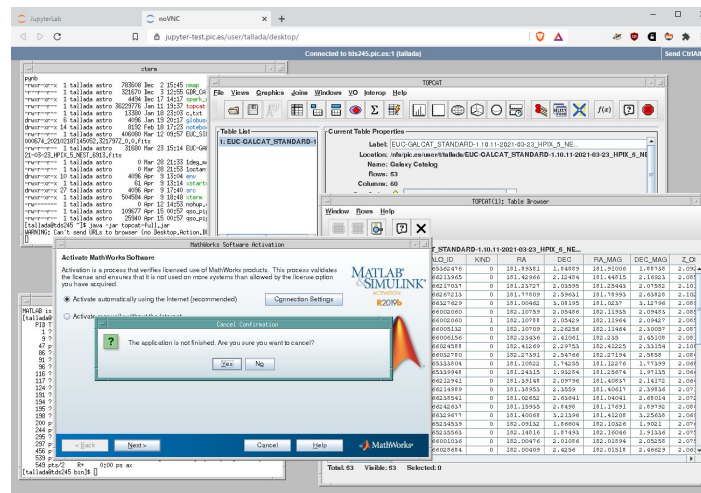



Table List: EU-CALCAT_STANDAR1

Current Table Properties: EU-CALCAT_STANDAR1-1:0-11-0201-03-23_HPR_5_M

Name: Galaxy Catalog

Format: #3

Columns: 40

NO_ID	RING	PA	DEC	PA_MAG	DEC_MAG	Z_O
0050476	0	181.13381	1.94879	181.20109	1.81793	2.096
0051305	0	181.40266	2.12484	181.44815	2.16063	2.098
0021703	0	181.23721	2.03559	181.23443	2.07582	2.110
0047133	0	181.17693	2.02613	181.17693	2.05021	2.110
0021203	0	181.06482	1.98159	181.06337	1.91794	2.108
0020188	0	182.10789	2.08489	182.11365	2.07465	2.110
0010000	1	182.10781	1.05403	182.11944	2.07462	2.085
0010132	0	182.10789	2.16238	182.11444	2.10097	2.088
0000156	0	182.13460	1.43061	182.1300	2.43060	2.088
0010458	0	182.43269	2.27153	182.43225	2.33354	2.100
0002160	0	182.27361	2.94768	182.27194	2.81859	2.098
0033304	0	181.18922	1.74235	181.18276	1.71789	2.068
0033384	0	181.12413	1.02824	181.12094	1.07130	2.068
0021241	0	181.13141	2.09719	181.14041	2.14372	2.064
0014689	0	181.18083	1.26509	181.18011	1.26931	2.071
0021541	0	181.02482	2.62441	181.04041	2.49104	2.072
0044264	0	181.12928	2.14619	181.11962	2.14021	2.074
0026181	0	181.46061	1.33396	181.43201	1.25650	2.061
0031000	0	182.07362	1.74606	182.10628	1.80021	2.074
0025583	0	182.18160	1.81493	182.18068	1.73331	2.071
0020056	0	182.08476	2.02088	182.01094	2.02820	2.071
0020044	0	182.08063	1.46436	182.03311	1.46602	2.068

Matlab Software Activation

Activate MathWorks Software

Activation is a process that verifies licensed use of MathWorks products. This process validates the license and ensures that it is not used on more systems than allowed by the license option you have specified.

Activate automatically using the Internet (recommended) Connection Settings

Cancel Confirmation

The application is not finished. Are you sure you want to cancel?

Yes No

<https://jupyter.pic.es/>

- **Interactive exploration** (visualization)

- **Very fast (85% plots < 30s)**
- **Full dataset plots** (over all rows)
 - May use sampling
 - Cone search tool
 - 1D histogram & 2D heatmap
- **Guided process (no SQL knowledge required)**
 - Expert mode
 - Custom UDFs: healpix, geometric...

- **Distribution**

- **Parquet, CSV, FITS format** (custom SerDe)
- Email with a link to download dataset

- Based on Python / Flask + Hadoop / Hive

- **Data**

- 40 TiB catalogued data
- >100 catalogs (simulated and observed)
- Supporting multiple projects
 - DES, PAU, Euclid, MICE and Gaia

- **Users**

- >1500 usuarios registrados
 - ~150 active users worldwide
- >13.5K custom catalogs generated
- >20k interactive queries

- **Performance**

- **>75% custom catalogs finish in <3 min**
- Resource queues with reservation
- Preemption to keep interactive response times

Science Portal: service improvements

- Multi-instrument, -frequency, -messenger
 - GW, optical, Gammas, Neutrinos, etc.
- Integrate with Jupyter notebooks
- Guided data products:
 - Spatial / time cross-match
- Additional UDFs functions:
 - geometric, array, footprint, MOCs
- Federated identity with EduGAIN
 - Per-user storage, quotas & accounting
 - Each user has its own schema (MyDB)
- VO protocols:
 - TAP (ADQL, VOTable, ...)
 - UWS, VOspace, ...
- New frontend
 - more plot types (healpix maps)
 - improved density/performance
- Social features
 - Catalogs and plots can be linked/shared
 - Public permissions to edit metadata (with manual moderation)

Challenge 1: How to handle massive data

- Data volume evolution: from a laptop, to a data center, to the cloud?
- Goals
 - Scalable storage
 - Short iteration cycles
 - Fast access
 - Efficient analysis
 - Reproducibility
 - Traceability
- Our approach
 - Hadoop, Hive, Spark (code to the data!)
 - Multidisciplinary teams
 - so that scientists and software engineers understand each other
 - stand behind developers / Tutoring / Teaching (Bootcamps, MOOCs)
 - Jupyter+Spark notebooks, jupytertext, git+LFS
 - Optimize algorithms
 - Custom algorithms (spark + treecorr)

Challenge 2: Gaining (more) independent users

- **Goal:**
 - Reach a broader community
 - User should not notice the transition from laptop, to cluster, to grid
 - Simple, usable interfaces
 - Guided processes
 - Avoid configuration files and terminals
 - While still providing access to advanced features for expert users
- **Our approach**
 - **Interfaces**
 - CosmoHub
 - JupyterHub (+ extensions: i.e. VNC)
 - HTTP/Webdav
 - VO protocols
 - **Training:**
 - Documentation / Quick start guides (e.g. HTCondor / Spark)
 - Bootcamps / MOOC

Thanks for your attention Questions?

Credits to: E. Acción, V. Acín, C. Acosta, A. Alou, A. Bruzzese, J. Carretero, J. Casals, R. Cruz, M. Delfino, J. Delgado, M. Eriksen, D. Graña, J. Flix, E. Johana, G. Merino, C. Neissner, A. Pacheco, C. Pérez, A. Pérez-Calero, E. Planas, M.C. Porto, J. Priego, P. Tallada, F. Torradeflot

www.pic.es

Questions

- How many users prefer the query builder to the expert SQL mode?
 - We don't know. We did not take the time to measure it.
- How do you see what the users are doing and get their feedback?
 - We can see the queries done
 - Batch historic
 - Interactives are a bit more difficult to get (Tez)
 - Many users directly know us and they send emails
 - Contact form (a bit hidden)
- Are sub-queries allowed?
 - Yes
- What happens if a user tries a very-long running query, or one generates huge output?
 - We don't have a limit. Up to know nothing horrible happened, except very large custom catalogs that they cannot be used. We remove them after one month.
 - We could add a time limit
 - Disk quota for the future
- What is the database backend and why did you choose it?
 - Hive
 - Free software
 - SQL interface (transparent migration from Postgres)
 - Workflow that perfectly fits
 - Storage and processing scale linearly, parallel queries, great performance (low latency "not allowed")
 - Expandable
- Does CosmoHub still run when the database is being backup?
 - Yes. File system snapshot
- What is the codebase for the online plot generator
 - Custom SQL wrapper to aggregate plot values in Hive
 - Backed transfers values in csv+json via websockets to the browser
 - Plots in frontend use plotly

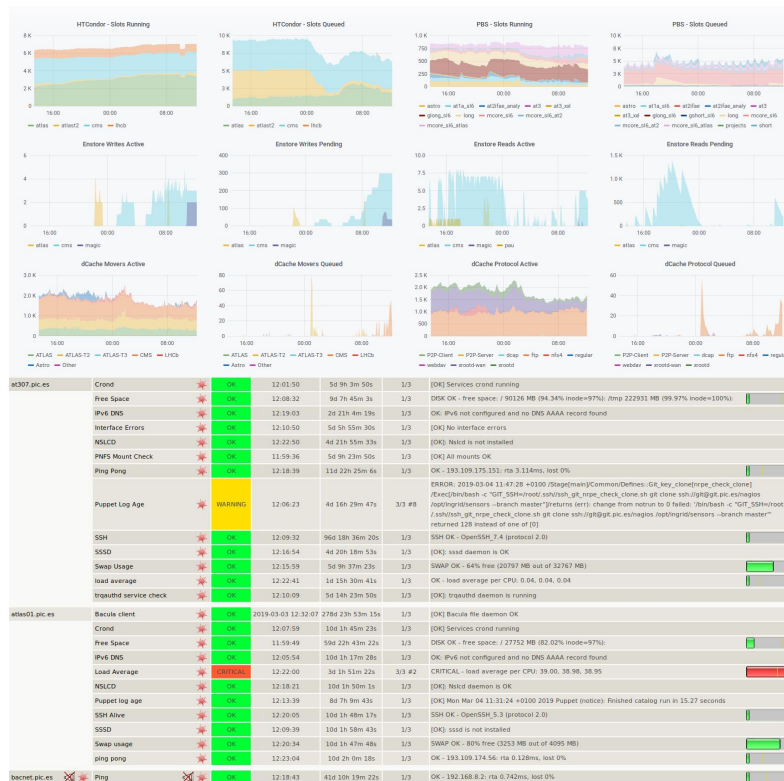
- Administration

- Puppet
- LDAP
- Wiki
- Git


- Monitoring

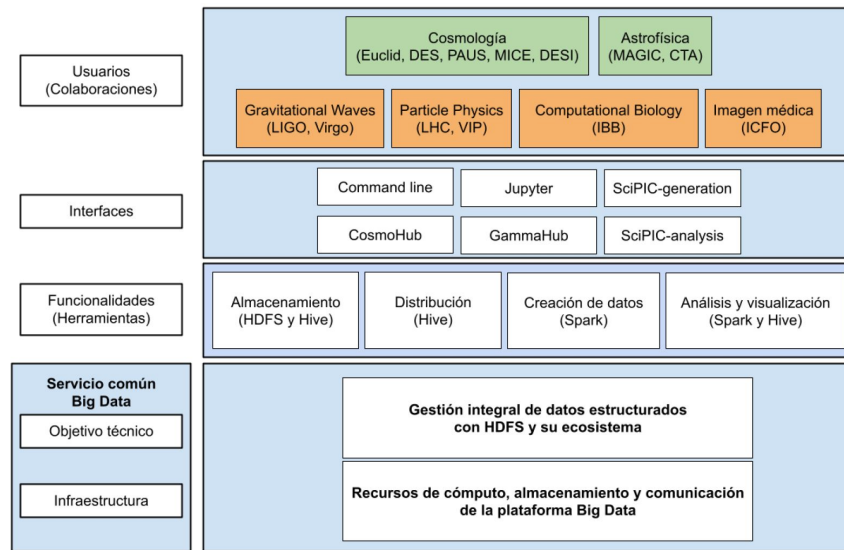
- Nagios + Thruk
- Elasticsearch + Logstash + Kibana
- Collectd + Graphite + Grafana
- Pakiti
- Cacti

- Manager on Duty (MoD) (24/7)



PIC Big Data common service

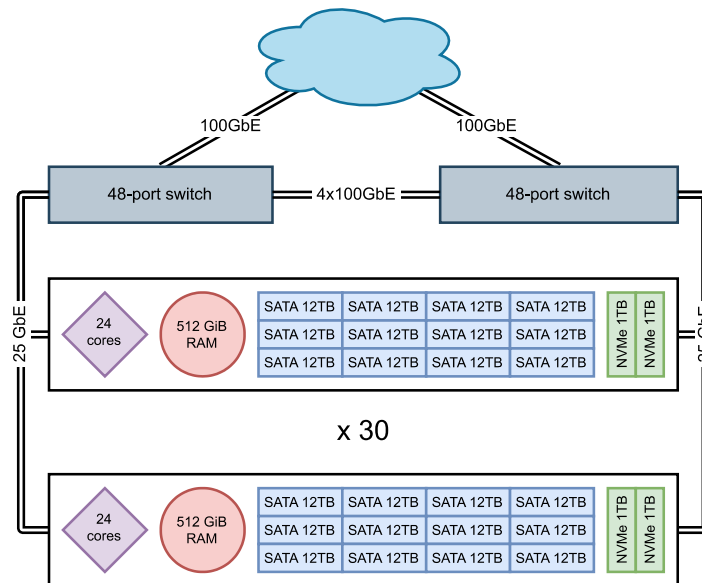
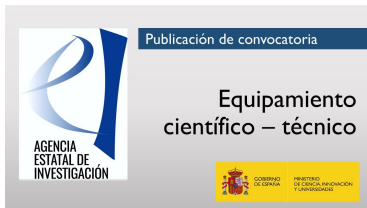
- Based on  **hadoop**
 - Open source Big Data Platform
 - Distributed storage and processing
 - Runs on commodity computer clusters
 - Scalable from dozens up to thousands of nodes
 - *Performance scales with HW*
 - Fault tolerant
 - Simple machines working together
 - *no single point of failure*
- Last update (summer 2020):
 - Custom DIY nodes
 - 12 nodes AMD Threadripper 1920X
 - 128 GB RAM, 12 x 3 TB SATA HDD hot-swap
 - 2x1 TB NVMe SSD i 2x10-GBASE-T LAN
 - Hortonworks HDP 3.1.4
 - Hadoop 3.1.0
 - Hive 3.1.0
 - Spark 2.3.2



Organización servicio común Big Data. Las disciplinas que hacen uso habitual de la plataforma aparecen en color verde mientras que los potenciales usuarios aparecen en color naranja.

Hadoop platform expansion

- “Convocatoria de ayudas para la adquisición de equipamiento científico-técnico” (EQC2021-007479-P) (2021-2024)
- Near future update (Dec 2023):
 - **30 nodes, 720 cores**
 - **15 TiB RAM, 2.0 PB, 60TB NVMe**
 - Custom Hadoop distribution



Catalog data volume

Project	Date	volume / night	Number of objects (catalog)
SDSS	2000 - now	variable	2×10^6
MICE GC (sims)	2013	NA	5×10^8
DES	2013 - 2018	2.5 TiB	4×10^8
GAIA*	2014 - 2019	40 GiB	1.8×10^9
Euclid Flagship (sims)	2016 - now	NA	8.7×10^9
Euclid	2023 - 2029	100 GiB	1.5×10^9
LSST	2024 - 2034	15 TiB	1×10^{10}

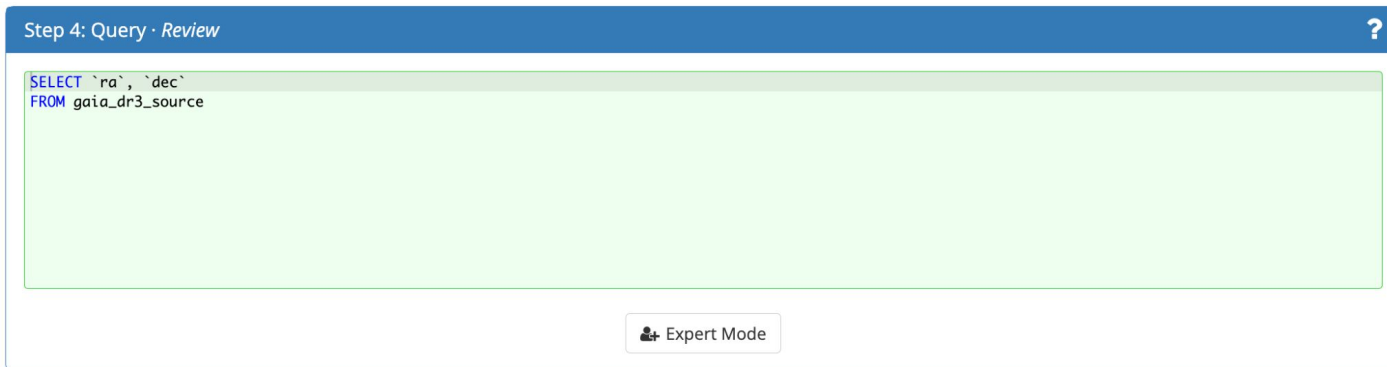
* DR3 full sky star catalog

Guided process

- Catalog description
- Value Added Data · Directly download useful or necessary files to analyse the catalog
- Step 1: Columns · Select the fields you need
- Step 2: Sampling · Select a subset and get faster results
- Step 3: Filters · Add conditions to refine your search
- Step 3a: Cone Filter · Restrict results around a sky position
- Step 4: Query · Review

Guided process

- Catalog description
- Value Added Data · Directly download useful or necessary files to analyse the catalog
- Step 1: Columns · Select the fields you need
- Step 2: Sampling · Select a subset and get faster results
- Step 3: Filters · Add conditions to refine your search
- Step 3a: Cone Filter · Restrict results around a sky position
- Step 4: Query · Review



The screenshot shows a web interface for reviewing a query. At the top, a blue header bar contains the text "Step 4: Query · Review" on the left and a question mark icon on the right. Below the header is a large, light green rectangular area containing a SQL query: `SELECT `ra`, `dec`
FROM gaia_dr3_source`. At the bottom center of the interface, there is a button with a person icon and the text "Expert Mode".

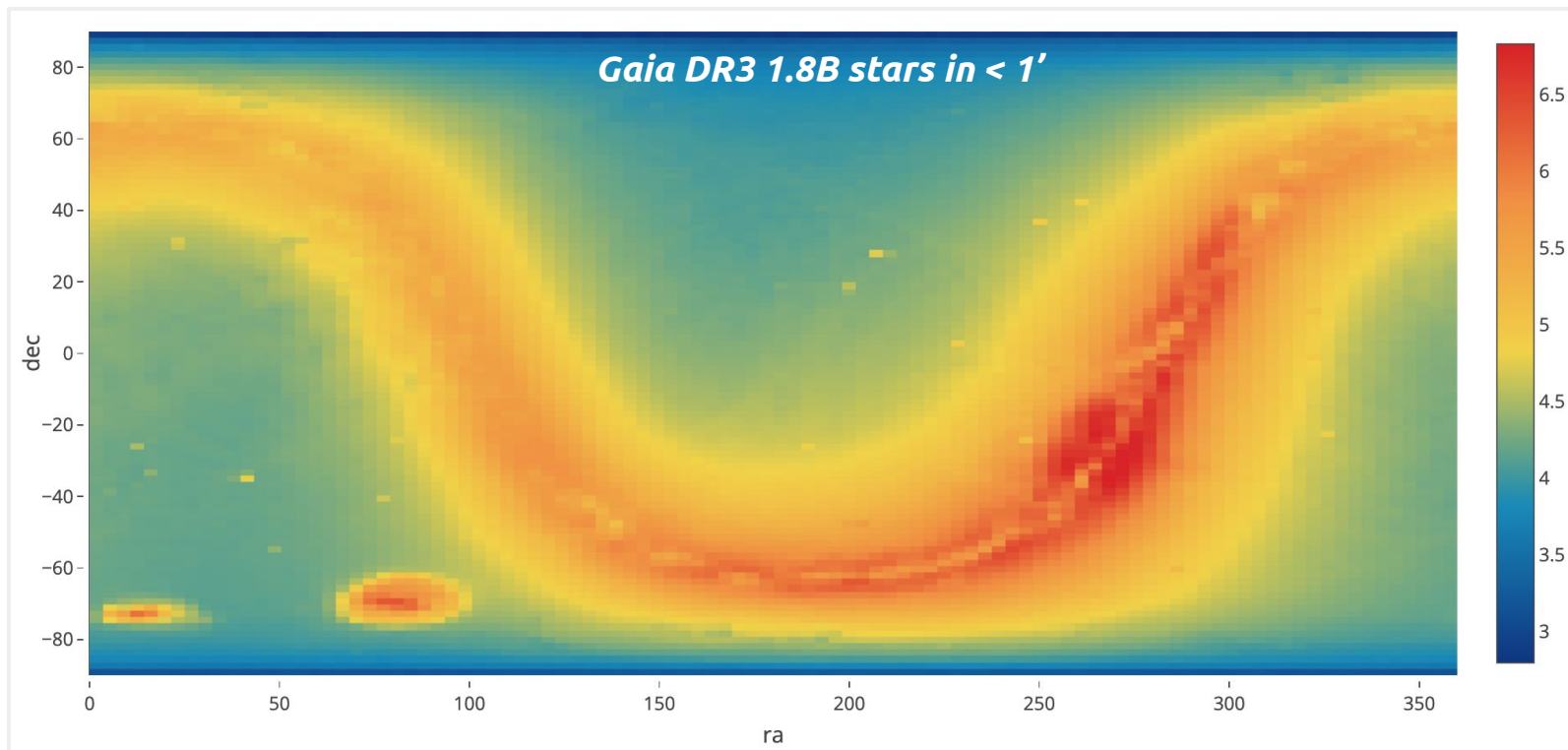
Guided process

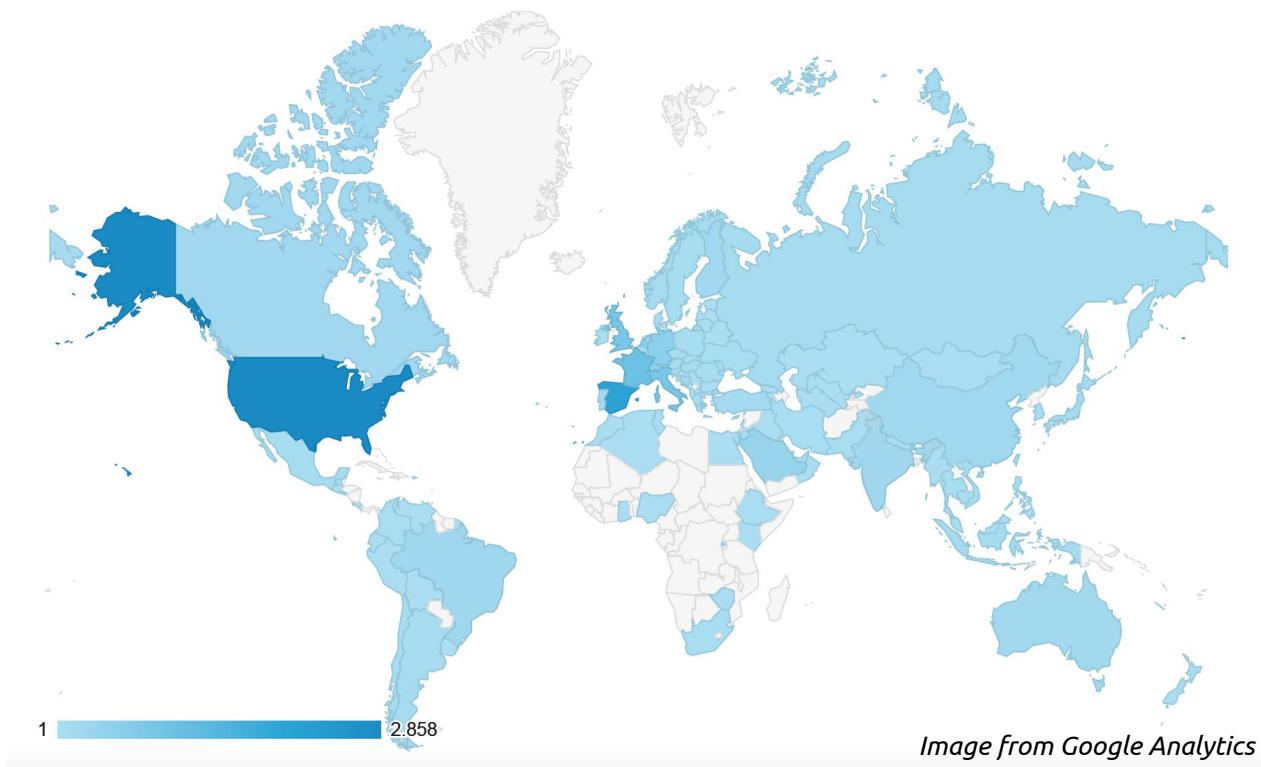
- Catalog description
- Value Added Data · Directly download useful or necessary files to analyse the catalog
- Step 1: Columns · Select the fields you need
- Step 2: Sampling · Select a subset and get faster results
- Step 3: Filters · Add conditions to refine your search
- Step 3a: Cone Filter · Restrict results around a sky position
- Step 4: Query · Review
- Step 5: Analysis · Explore the selected data (Table, Scatter, Histogram, **Heatmap**)

Guided process

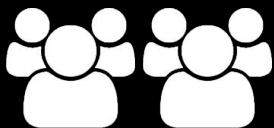
- Catalog description
- Value Added Data · Directly download useful or necessary files to analyse the catalog
- Step 1: Columns · Select the fields you need
- Step 2: Sampling · Select a subset and get faster results
- Step 3: Filters · Add conditions to refine your search
- Step 3a: Cone Filter · Restrict results around a sky position
- Step 4: Query · Review
- Step 5: Analysis · Explore the selected data (Table, Scatter, Histogram, Heatmap)
- Step 6: Format · Select a file type (csv.bz2, FITS, parquet, ASDF)
- Step 7: Request · Review citation guides

Guided process





Ubicación física de los usuarios de CosmoHub desde su creación



~ 150 active users



~ 12K custom catalogs



~ 40 TiB hosted data



> 100 catalogs

Public catalogs

- Gaia (DR3, Mean Spectrum, EDR3, DR2 & DR1)
- DESI Legacy Survey DR8 PZ
- DESI Legacy Survey with Photoz (DR8)
- COSMOS 2020 (Classic | Farmer)
- COSMOS 2015 Laigle (v2.1)
- LSST DESC DC2 (Truth-match | Object table)
- DES DR2
- DES Y1A1 Morphological catalog (v1.0)
- DES Y1A1 Gold Data (v1.0)
- GLADE (v2.3, v2.4) & GLADE +
- VIPERS photometry and spectroscopy (PDR2)
- KiDS (DR4)
- CANDELS Bulge-Disk decomposition (2018)
- CFHTLenS (good fields) (v1.2)
- Alhambra photometric redshifts (v1.0)
- ALHAMBRA S/G CLASSIFIED (v1.0)
- PAUS+COSMOS photo-z catalog (v0.4)
- PAUS-COSMOS Early Data Release (v1.0)
- PAU.MillGas Lightcone (2016-07-18)
- DEEP2 Redshift catalog (DR4)
- MICE halo properties
- MICECAT (v2.0, v1.0)

Hive tuning

- We have set the platform so that queries over large tables are really fast:
 - Apache Tez execution engine instead of the venerable Map-reduce engine
 - ORCfile: a new table (column based) storage format
 - Vectorized query technique: batches of 1024 rows at once

Load balancing

- Set up two different queues given the two different profiles:
- 'Interactive': real-time analysis (low latency)
- 'Batch': custom catalogs (high latency)
- Configure queue shares and preemption:
- batch jobs take idle resources to maximize efficiency (10-90)
- interactive jobs can take resources from batch queue (90-100)

- ReST API powered by Flask:
- flask-restful - ReST framework
- sqlalchemy - database ORM
- websockets - bidirectional communications
- gevent - asynchronous framework
- pyhive - hive connection library
- pyhdfs - hdfs bindings

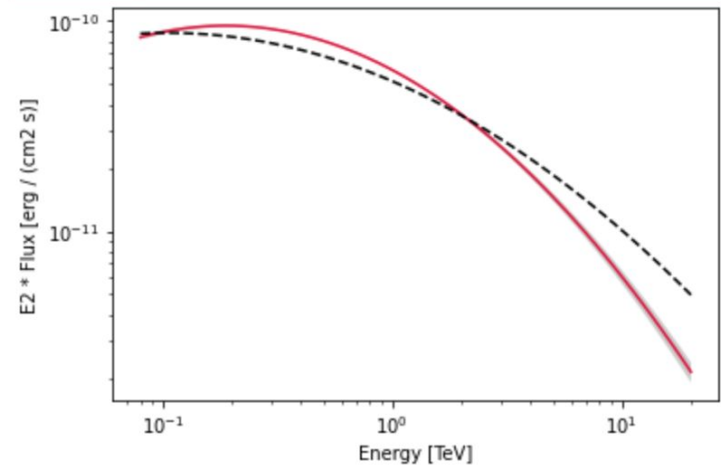
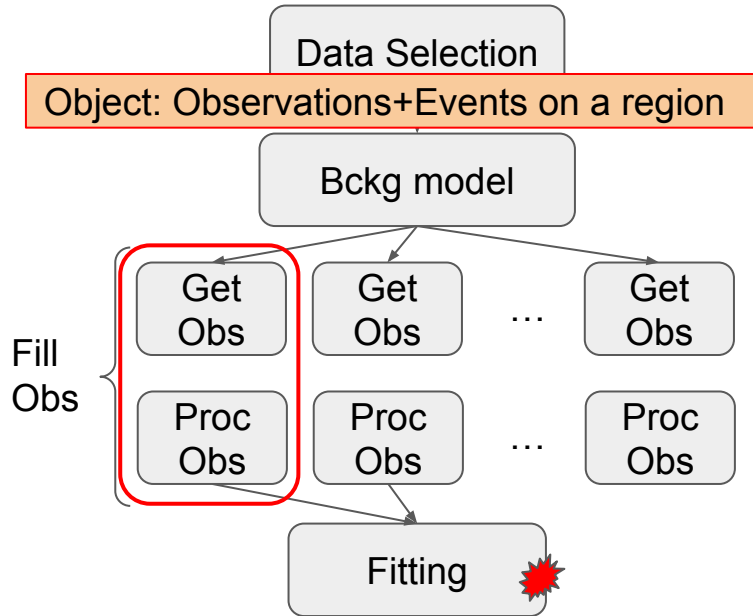
- Responsive Web interface powered by:
- Angular JS - web app oriented HTML framework
- Bootstrap - responsive frontend framework
- Plot.ly for plotting
- Wordpress as backend to edit "static" content

- Plan Complementario: Astrofísica y Física de Altas Energías (2022-2025)
- Línea 8: Computación, big data e inteligencia artificial:
 - “Dar un salto cualitativo en la participación española en la siguiente generación de proyectos internacionales líderes en el área de la Astrofísica y Física de Altas Energías, con un énfasis particular en sus aspectos más tecnológicos” (MICINN)
 - “Aprovechar al máximo las capacidades actuales de las infraestructuras existentes de análisis de big data, expandiendo su capacidad y su ambición, hasta tener en España un hub de datos de astronomía multi-mensajero único en Europa” (BOE 2023)



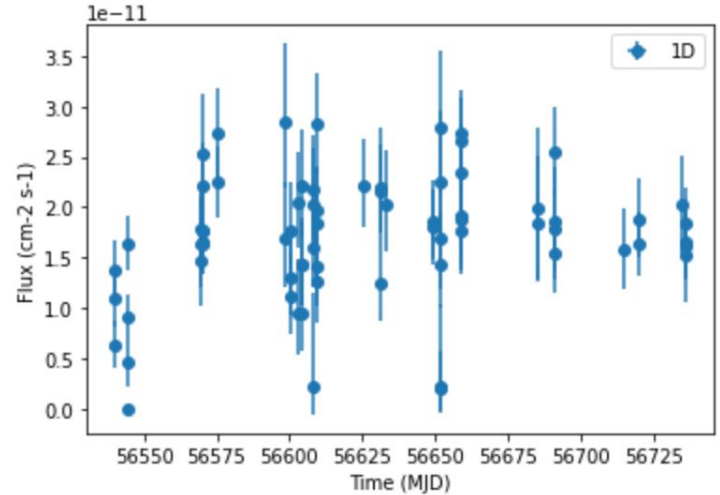
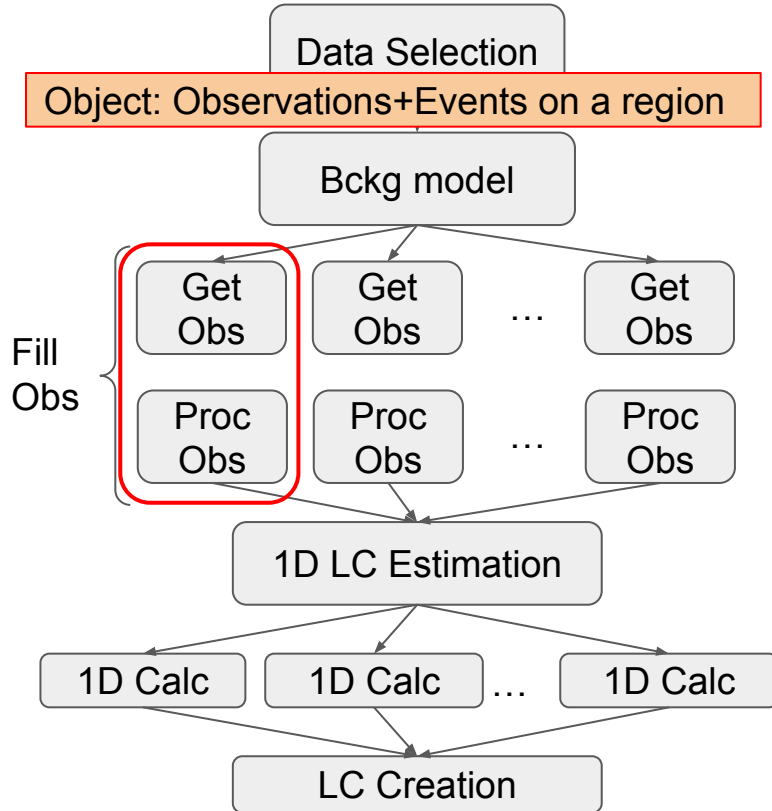
GammaHub Data Products: introducing parallelism with Spark

Gammapy low level SED workflow + parallelism (v.1)



GammaHub Data Products: introducing parallelism with Spark

Gammapy low level LC workflow + parallelism (v.1)



Multinstrument GammaHub Data Products

Objective: allow to perform multi-instrument data selection and analysis

Login

Data Selection

Analysis Parameters

Source Name

DEC (deg)

RA (deg)

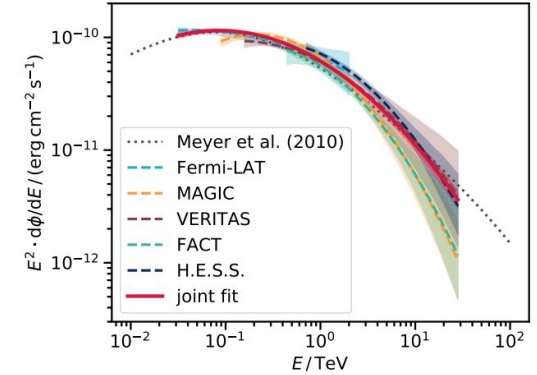
DEC (deg)

Radius (deg)

Select instr...
 All
 MAGIC
 HESS
 Veritas
 CTA

Enable multiple instruments selection
according to the data stored in Hive

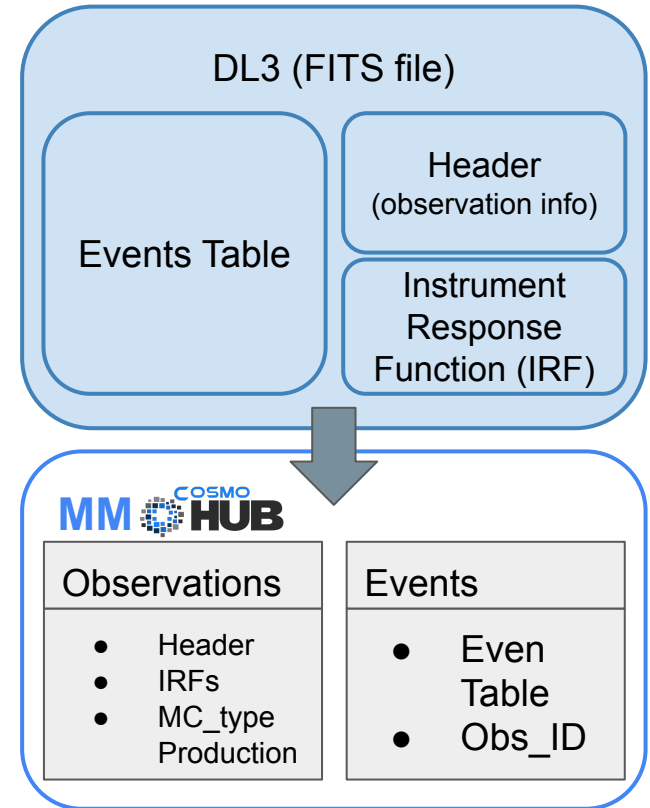
←



Example: joint crab paper. Scale the calculation using parallel calculations to obtain each curve and then combine all together on a single plot

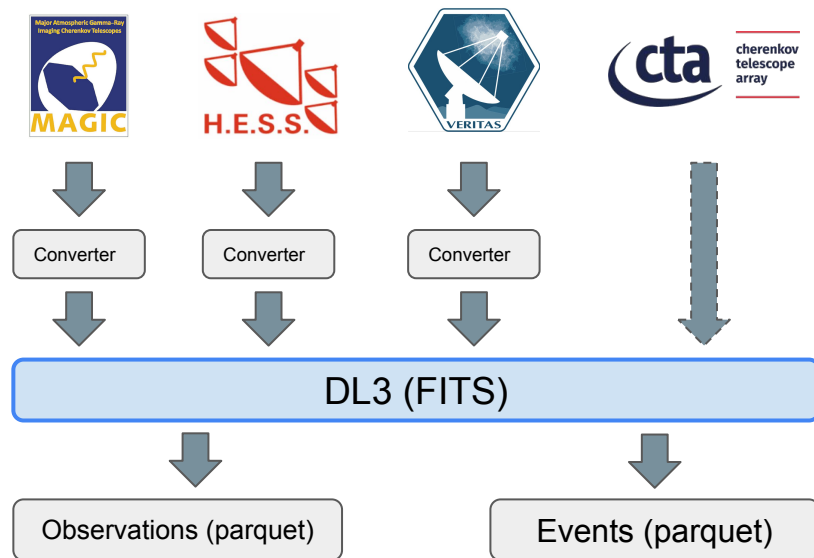
Science Portal: service improvements

- Multi-instrument, -frequency, -messenger
 - GW, optical, Gammas, Neutrinos, etc.
- Guided data products:
 - **Spatial / time cross-match**
 - ad-Hoc Value Added Data
- Integrate with **Jupyter notebooks**
- Additional UDFs functions:
 - geometric, array, footprint, MOCs



Science Portal: service improvements

- Multi-instrument, -frequency, -messenger
 - GW, optical, Gammas, Neutrinos, etc.
- Guided data products:
 - **Spatial / time cross-match**
 - ad-Hoc Value Added Data
- Integrate with **Jupyter notebooks**
- Additional UDFs functions:
 - **geometric, array, footprint, MOCs**



parquet file format: columnar data format + useful metadata for data ingestion