# Integration of GammaHub in ESAP

WP5 Progress Meeting, 26-27th of October 2020

J. Delgado (jordidem@pic.es), C. Nigro, L. Jouvin, A. Bruzzese, J. Rico
P. Tallada, M. Delfino, G.Merino

# Outline

- Previously on GammaHub…

- Second Use Case: Data Products

- Integration in ESAP

- Future Plans

# Previously on GammaHub… What is GammaHub?

GammaHub is intended to be an interactive science analysis platform by itself, taking previous experiences from https://cosmohub.pic.es , a platform to manage Cosmological Catalogs developed at PIC

It will provide tools to modelate large datasets on a big data environment in order to search, explore and plot billions of objects interactively using a Hadoop stack of tools

We will ingest datasets from multi-instrument astronomical gamma-ray experiments like MAGIC, HESS, VERITAS and CTA among others, but it could be applied to other science disciplines
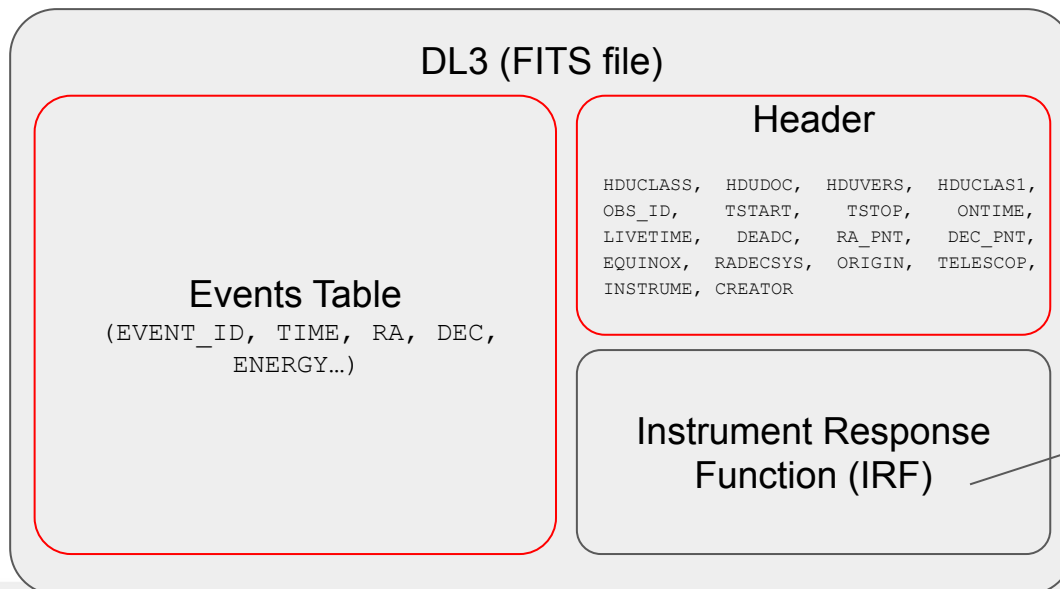
The main functionalities will be:

- <u>Interactive data selection</u>, exploration and plotting (scatter plots, histograms, heatmaps…)

- <u>Automatic Data Products</u> (Spectrum, Lightcurves…)

- Data access from <u>Python notebooks for complex analysis</u>

# Previously on GammaHub… The Data Model

Data Level 3 (DL3) v.0.2 (2018) contained in <u>FITS files.</u> The DL3 format is being promoted as a new standard in the context of building CTA
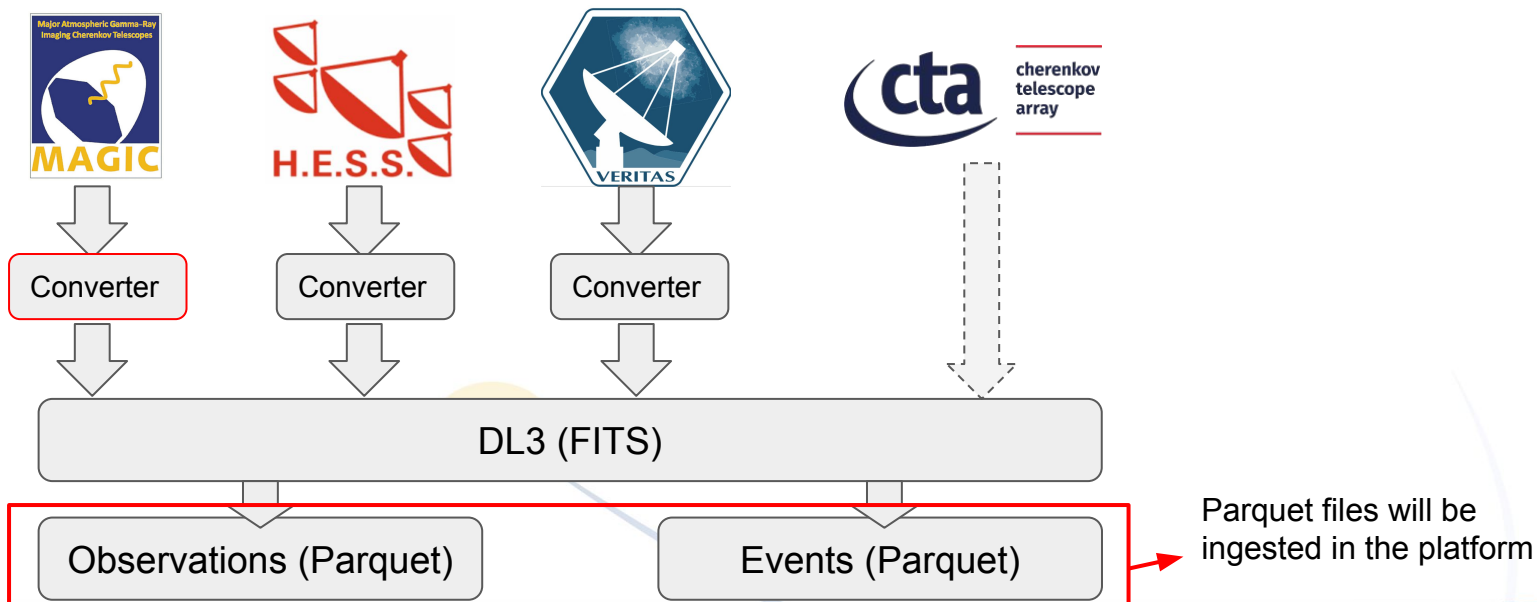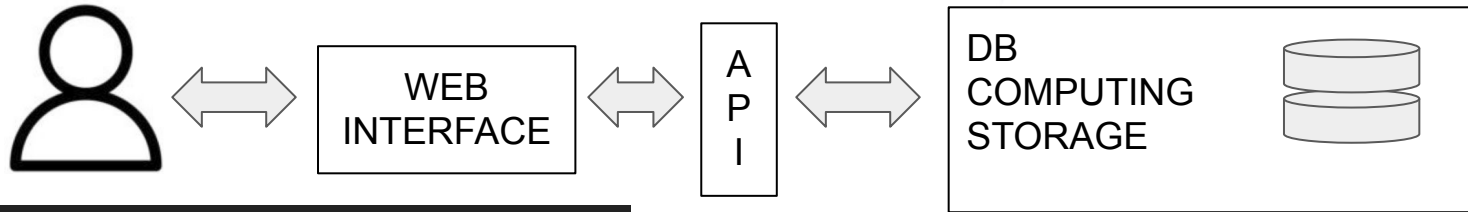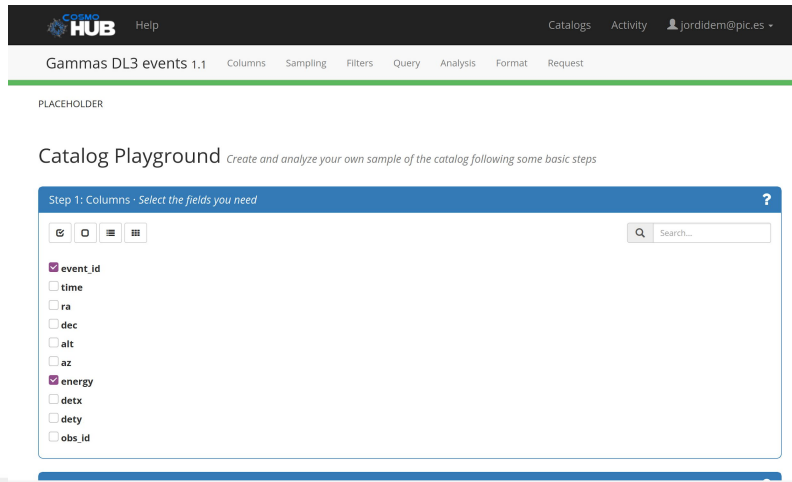


DL3 (FITS file)

**Events Table**
(EVENT_ID, TIME, RA, DEC, ENERGY…)

**Header**

```
HDUCLASS,  HDUDOC,   HDUVERS,  HDUCLAS1,
OBS_ID,    TSTART,    TSTOP,    ONTIME,
LIVETIME,  DEADC,    RA_PNT,   DEC_PNT,
EQUINOX,  RADECSYS,  ORIGIN,  TELESCOP,
INSTRUME, CREATOR
```

Instrument Response Function (IRF)

Data is modeled in the Hive DB

Not searchable, it will be accessed directly from disk

# Previously on GammaHub… The Data Model

Data Level 3 (DL3) v.0.2 (2018) contained in FITS files. The DL3 format is being promoted as a new standard in the context of building CTA



Parquet files will be ingested in the platform

# Previously on GammaHub… First use case

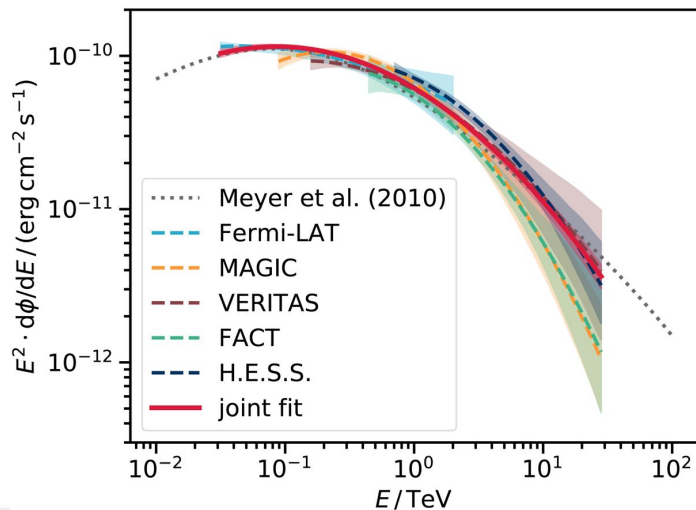1. Web interface for interactive/on-line exploration and visualization of the selected datasets.



✓ Data Searching (single and multi-instrument exploration)
✓ Data Plotting (Scatter, Histogram 1D, Heatmaps 2D…)
✓ Batch custom subsets

✓ User guidance / User friendly look and feel
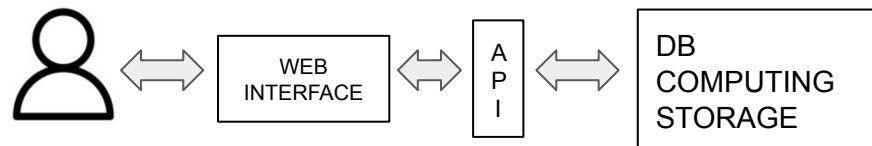✓ No SQL knowledge needed (but Expert mode included)

# Second Use Case: Data Products

1. Data Products: implement interactive common data analysis workflows using Gamma-ray data (spectrum, lightcurves…)

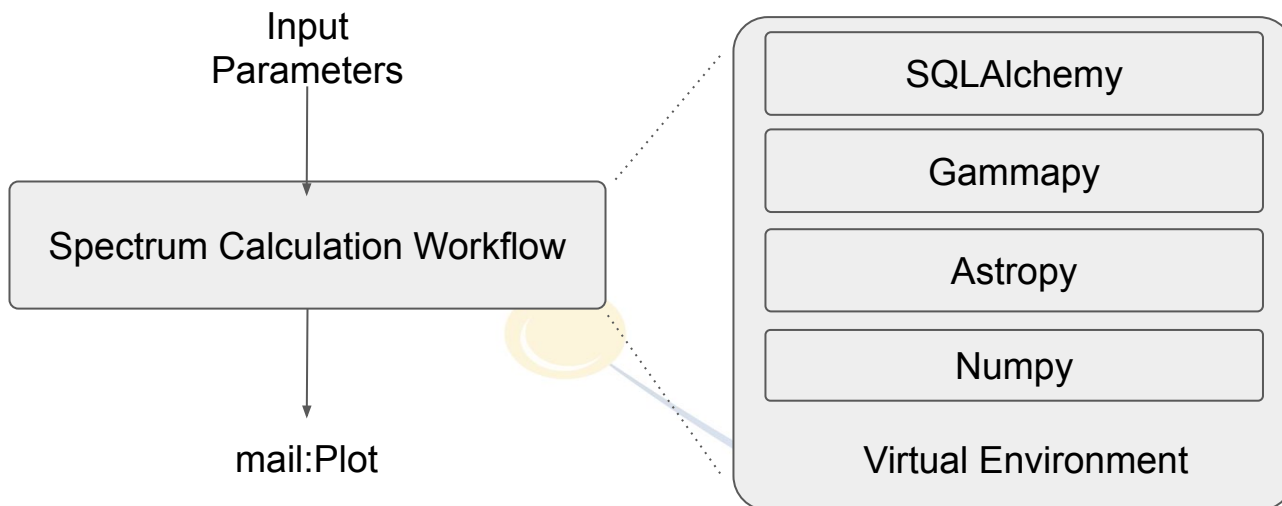Example: *Spectrum calculation*



← Crab Nebula SED.
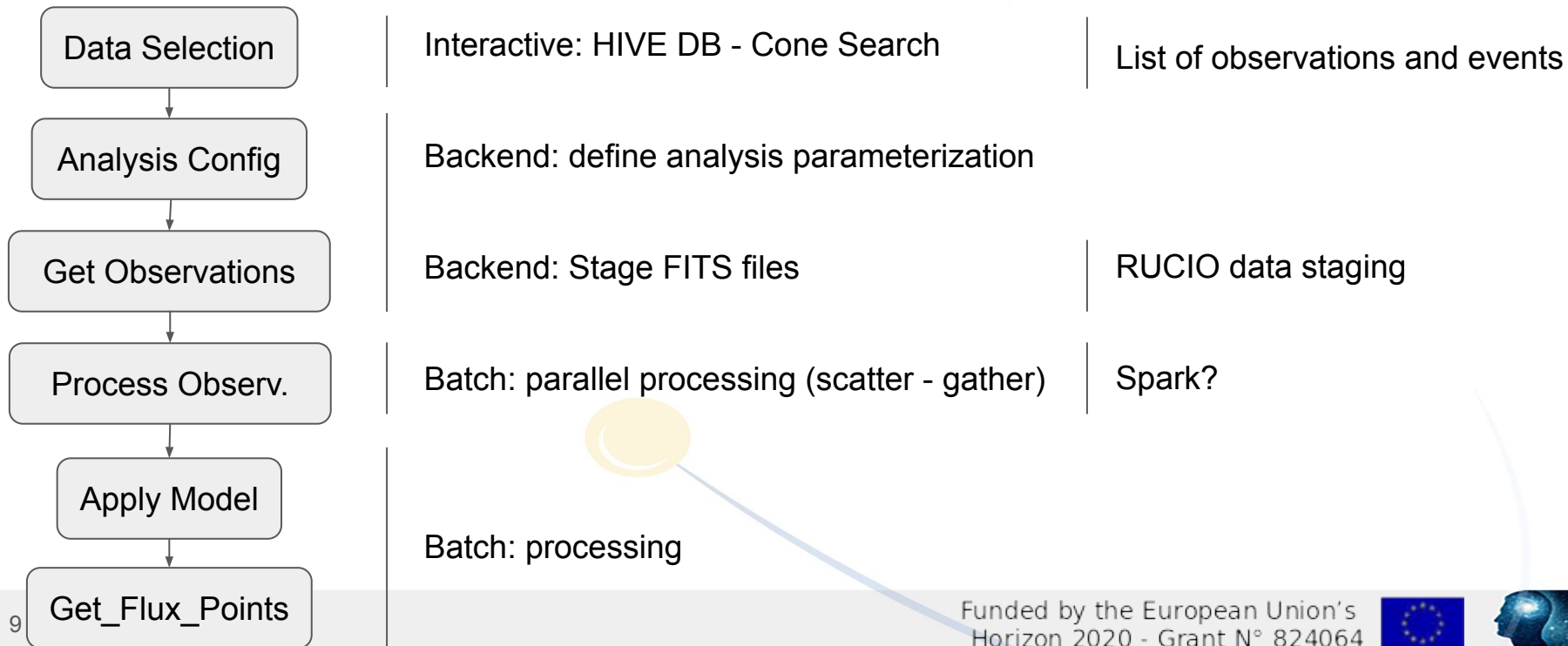C.Nigro et Al (2019)

# Second Use Case: Data Products

1. Data Products: implement interactive common data analysis workflows using Gamma-ray data (spectrum, lightcurves…)
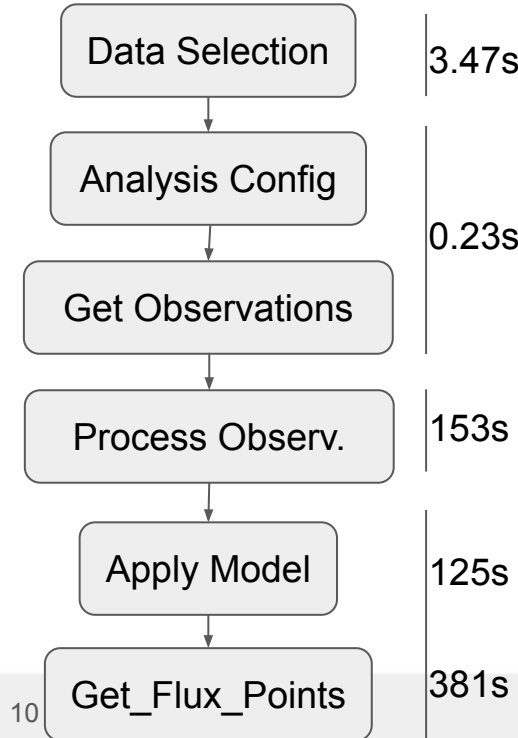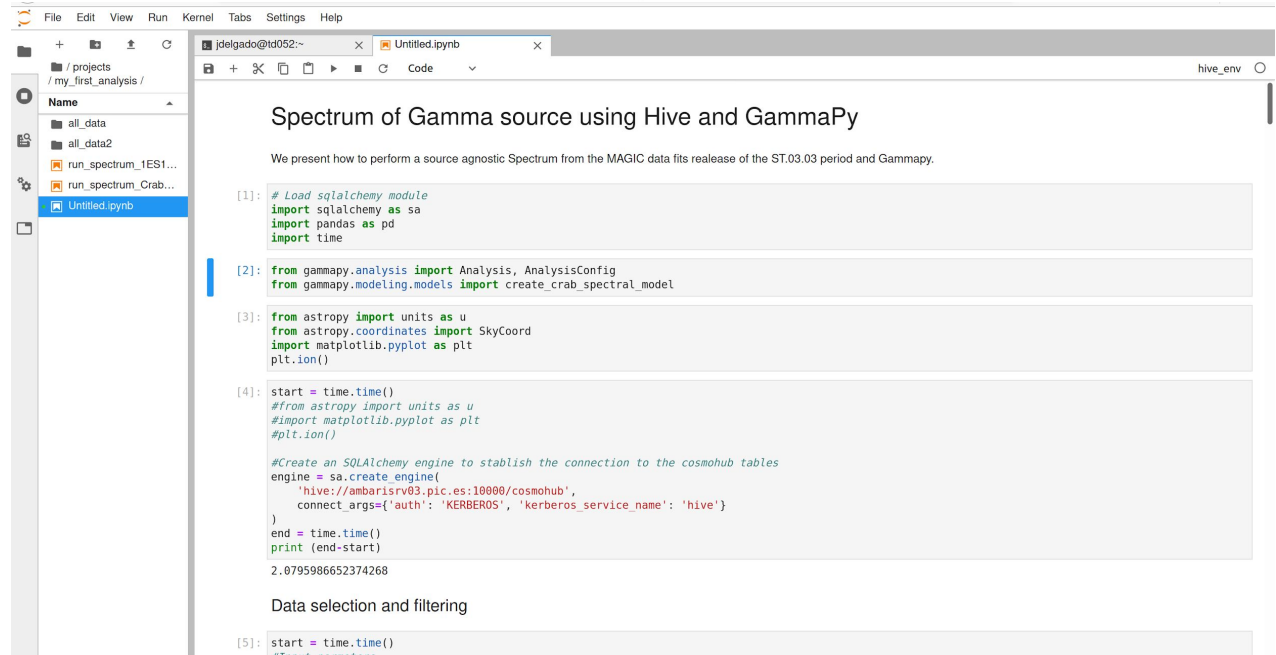
Example: *Spectrum calculation*

# Second Use Case: Data Products

*Spectrum calculation workflow using Gammapy*

| | | |
|---|---|---|
| **Data Selection** | Interactive: HIVE DB - Cone Search | List of observations and events |
| ↓ | | |
| **Analysis Config** | Backend: define analysis parameterization | |
| ↓ | | |
| **Get Observations** | Backend: Stage FITS files | RUCIO data staging |
| ↓ | | |
| **Process Observ.** | Batch: parallel processing (scatter - gather) | Spark? |
| ↓ | | |
| **Apply Model** | | |
| ↓ | Batch: processing | |
| **Get_Flux_Points** | | |

9

# Second Use Case: Data Products

*Spectrum calculation workflow using Gammapy*

# Integration in ESAP

## ESAP Architectural Design
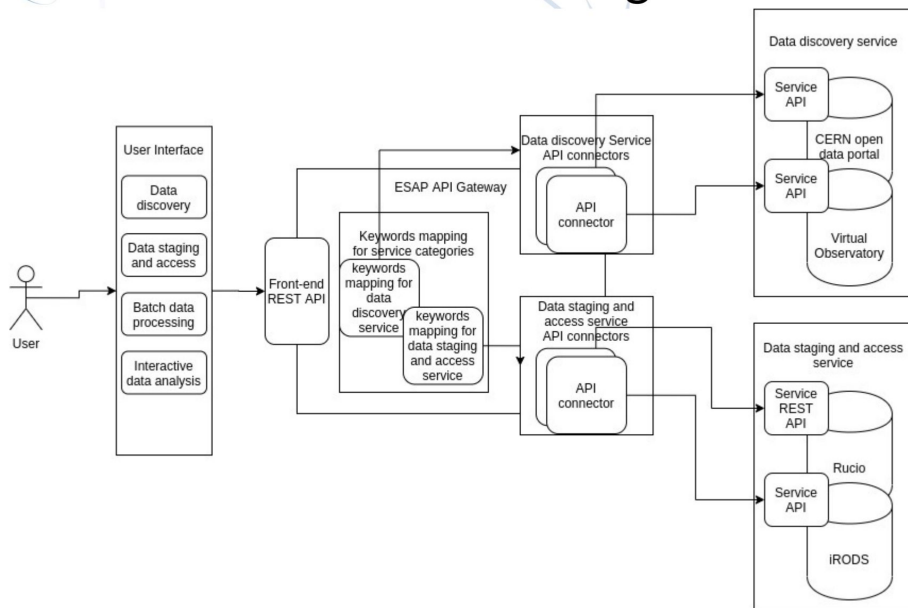


Open questions to solve:

- Data selection:
    - is HIVE/HADOOP compatible with ESAP?
    - How to use IVOA tools to publish and get data?

- Data staging:
    - We will upload DL3 files in the Data Lake
    - Install Rucio Client (Rob's demo)

- Computing: how we manage the parallelism?

# Future plans

Next steps in the developing of GammaHub

- Continue the data ingestion from other instruments, other MAGIC periods

- Upload DL3 data in the Data Lake and test the Rucio Client to stage and read data

- Implement parallelism using Spark (from jupyter notebooks to standalone)

  - Perform a multi-instrument spectral analysis of the Crab Nebula with Gammapy (reproduce C.Nigro et al *Towards open and reproducible multi-instrument analysis in gamma-ray astronomy*)

- Open to contribute to the resource federation, workflows publication...

- Implement a Data Product to calculate Light Curves

# Thank you!

WP5 Progress Meeting, 27th of October 200

J. Delgado (jordidem@pic.es), C. Nigro, L. Jouvin, A. Bruzzese, J. Rico,
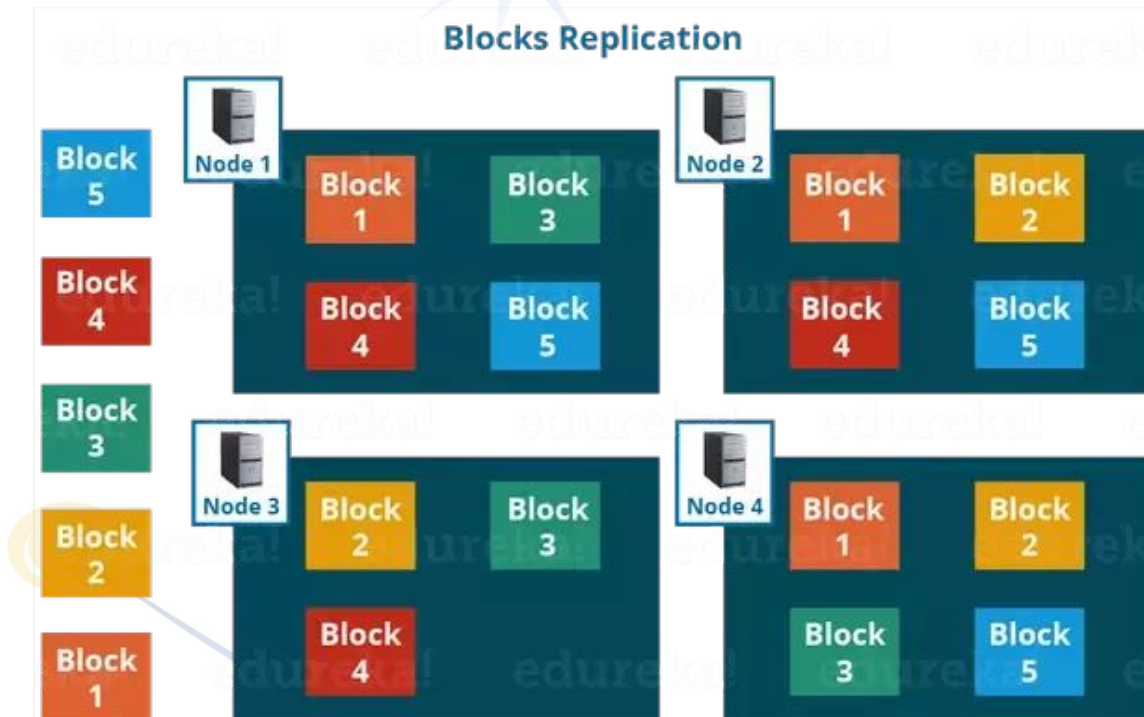P. Tallada, M. Delfino, G.Merino

# GammaHub components: Infrastructure

Data replication on HDFS:
- File is divided in blocks of data
- Default replication factor x3

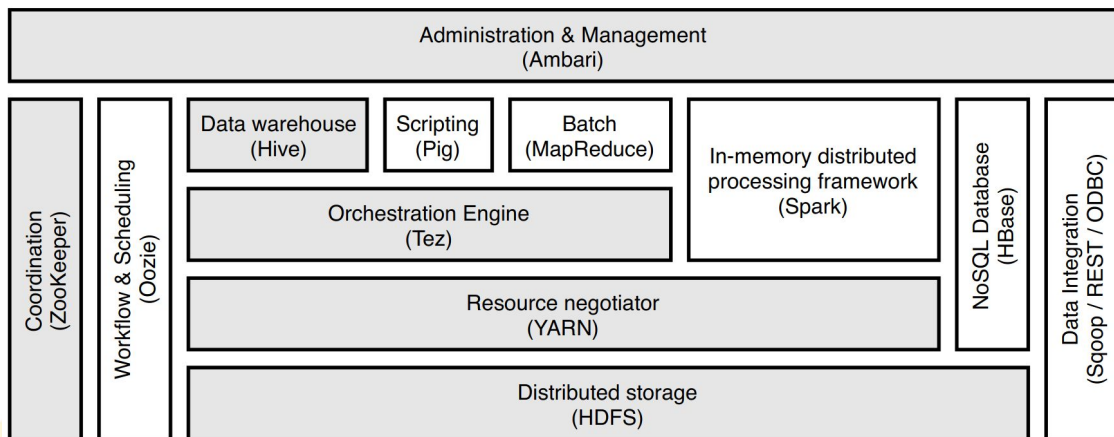File of 128MB -> 384MB stored in the platform

- Hadoop nodes includes access

- Rack awareness algorithm, to distribute data on different data nodes to improve network performance and fault tolerance



**Blocks Replication**

# GammaHub components: Infrastructure

Apache Hadoop:
- Distributed processing
  - Map-reduce model

- Dedicated computing cluster
  - Each node is a computing and data node

- Easy to scale from single server to thousand machines, each one offering local computation and storage

- Fault tolerance



Hadoop Layered ecosystem applied to CosmoHub (P.Tallada et. al. 2020 arxiv )