THE CTA PRODUCTION SYSTEM PROTOTYPE FOR LARGE-SCALE DATA PROCESSING AND SIMULATIONS
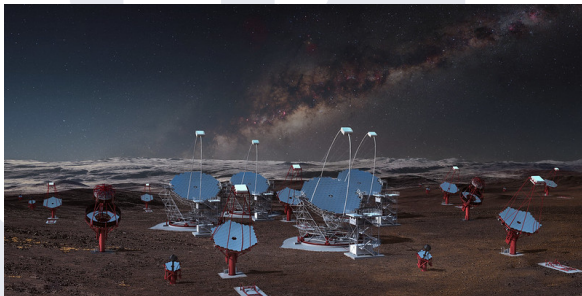
J. Bregeon, L. Arrabito, P. Maeght, M. Sanguillon for the CTA Consortium and A. Tsaregorodtsev
CNRS/IN2P3 LUPM & LPSC & CPPM

vCHEP 2021
May 20th 2021

- ▶ The next generation instrument in VHE gamma-ray astronomy (1500 participants in 31 countries)
    - ▶ Cosmic–ray origins, High Energy astrophysical phenomena, fundamental physics and cosmology
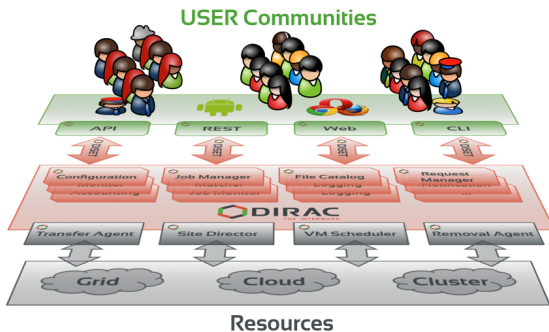


- ▶ Two arrays of Cherenkov telescopes
    - ▶ North site (La Palma, Spain): O(10) LSTs + MSTs
    - ▶ South site (Paranal, Chile): O(50) LSTs + MSTs + SSTs
- ▶ Construction starting now, operations for ∼30 years

*http://www.cta-observatory.org*

# DIRAC INTERWARE

- **D**istributed **I**nfrastructure with **Remote Agent Control**
    - all–in–one solution to access distributed resources (computing and storage) → *http://diracgrid.org*
    - started by LHCb more than 15 years ago
    - used today by many experiments: LHCb, CLIC, Belle II, CTA...
    - open source software maintained and developed by an international consortium → *http://github.com/dIRACGrid*
    - EGI instance with hundreds of users from many "small" VOs

# DIRAC for CTA

- ▶ DIRAC instance dedicated to CTA distributed at 3 sites (CC-IN2P3, PIC, DESY)
- ▶ DIRAC v7r1 deployed on 5 core servers
  - + extended with `CTADIRAC` plugin (*code on CTAO GitLab*)
  - ▶ 1 running WMS services (32 cores, 32 GB RAM)
  - ▶ 1 running WMS agents and executors (32 cores, 32 GB RAM)
  - ▶ 1 running TS and RMS (16 cores, 8GB RAM)
  - ▶ 1 running DMS + 1 DIRAC SE (16 cores, 8GB RAM, 2 TB of disk for the SE)
  - ▶ 1 running duplicated DMS, TS, RMS services (8 cores, 32 GB RAM)
  - + a VM for the web server
- ▶ Services
  - ▶ MySql databases at CC-IN2P3 (File Catalog, Transformation DB) and PIC (Accounting, Jobs DB and more...)
  - ▶ ELK instance (ElasticSearch++) at CC-IN2P3 for the DIRAC Monitoring system (new for us)
  - ▶ CVMFS for software distribution (CC-IN2P3 and DESY)
  - ▶ FTS for data movement (CERN instance)

**Transformation**

| Task Template |
| Data Filter or InputData Query |
| Output DataQuery |
| Plugin |

## Transformation System

▶ transformation: task template, input and output data

▶ 2 transformations "connected" if output data of T1 intersect input data of T2
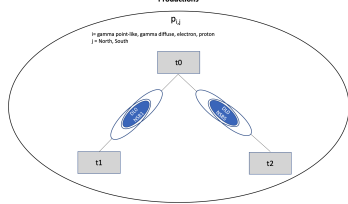
▶ workflow is data driven!



Transformation 1 → DL0 → Transformation 2 → DL1 → Transformation 3 → DL2

## Production System — introduced in version v7r0

▶ production as a set of "linked" transformation"

▶ design to handle the workflow at high level

▶ meant to be used by "operators" in *production*

**Connected transformations**

t1 → OutputData Query → data → InputData Query → t2

**Productions**

$p_{ij}$

i= gamma point-like, gamma diffuse, electron, proton
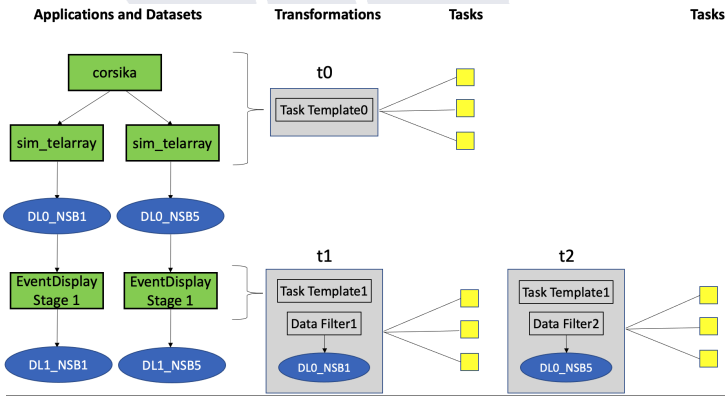j = North, South

t0

t1    t2

## Prepare transformations

- ▶ setup workload into DIRAC Job API
- ▶ setup Job to be used through transformation

## Build Productions

- ▶ define input and output datasets, as meta-queries
- ▶ assemble Productions from Transformation definition and datasets

## Run Productions

- ▶ assign simulation jobs to production
- ▶ activate production and monitor workflow processing

## Run Productions

▶ monitoring mostly through web interface

▶ failed processing jobs are automatically retried (FailOverRequest module)

▶ stop/resume productions in case of issues (scripts)

▶ monitor datasets content (scripts)

## Tools examples

▶ in vanilla DIRAC
  ▶ `dirac-prod-get-all`
  ▶ `dirac-prod-start` / `stop`
▶ in CTADIRAC extension
  ▶ `cta-prod-extend-sim`
  ▶ `cta-prod-monitor`
  ▶ `cta-prod-create-dataset`
  ▶ `cta-prod-show-dataset`

▶ `PROD5(b)` goals for CTA
- ▶ improved telescope and camera descriptions
- ▶ fine tuning of telescope positions

▶ A few numbers
- ▶ ran "almost" non stop for 9 months
- ▶ 200 MHS06.hours (20% by users) in 2.8 millions of jobs
- ▶ 15 different sites across Europe
- ▶ 1.7 PB on disk in 4.7 millions of files



Running jobs by JobType
44 Weeks from Week 26 of 2020 to Week 18 of 2021

Max: 10,532, Average: 2,049, Current: 165



Total Number of Jobs by Site
44 Weeks from Week 26 of 2020 to Week 18 of 2021

# DATA MANAGEMENT

▶ CTA now handles 6 PB on disk and 2 PB on tape
  ▶ use DIRAC for data management operations
  ▶ data removal
  ▶ moving data from disk to tapes
▶ DIRAC framework
  ▶ Dirac File Catalog with file and replica meta data
  ▶ Request Management System (asynchronous tasks)
  ▶ Transformation System (combine with Production System!)
▶ Example: Move (1 PB) of PROD3 to tape
  ▶ recently setup CERN FTS to run under the hood
  ▶ one Transformation per dataset — meta–data driven!



Number of Transfers
~ 5 K per hour

Throughput
500 MB/s

- ▶ ELK stack: Elasticsearch, Logstash and Kibana
  - ▶ CTA relies on the CC-IN2P3 instance
- ▶ need proper configuration of the DIRAC instance (and some work to support latest ELK version)
- ▶ gives access to a lot of information in a very efficient way
  - ▶ a few examples below
  - ▶ now need to dig more and extract specific information for our needs
  - ▶ in particular job parameters like memory usage

▶ How do we fully automatize this kind of workflow?
  ▶ issue: the look-up table/ML model/BDT are "part" of the software, how do we know that these are ready and that the next step can be run ?
  → different solutions envisaged, need to try!



Monte Carlo Simulation Setup
DL0 → DL1 → DL2

TS 3
Recon Training

TS 1a

MC DL0 Train

TS 2
Image Recon

DL1 Train

ML Model

TS 1b

MC DL0 Test

DL1 Test

**Software**

**Input Data**

How to wait for the ML model
to be ready and trigger the
launch of the TS 3 Reconstruction ?

TS 3
Event Recon

DL2

# Conclusion

- ▶ DIRAC is a modular open source tool for distributed computing, driven by an open minded consortium that offers a place for help and discussions.
- ▶ CTADIRAC is the DIRAC instance developed for the construction of CTA, it's been up and running for more than 8 years: hundreds of millions of jobs, Petabytes of data.
- ▶ CTADIRAC also largely benefits from HEP computing ecosystem: VOMS, CVMFS, FTS…
- ▶ CTADIRAC is mostly ready to become the CTA computing resources and workflow management system: in–kind contributions to the CTA Observatory being discussed now, hopefully finalized by end of the year.
- ▶ Learn more about CTA performance and see the full list of computing centers that provide resources and support (Thanks to them!) at: *https://www.cta-observatory.org/science/cta-performance/*.