# Past and Future Challenges for Distributed Computing at the ATLAS Experiment on the Iberian Peninsula

**Helmut Wolters**
**LIP Coimbra**
on behalf of the Iberian Cloud
in the ATLAS Distributed Computing team

IBERGRID 2019
Santiago de Compostela
25 September 2019

# Overview

- **ATLAS computing model**
- **The Iberian Cloud inside ATLAS**
- **Development Activities for Atlas Distributed Computing**
- **Challenges for run 3 and HL-LHC**
- **Conclusions**

# ATLAS computing model



**Note:** "Cloud" in ATLAS means a regional setup of one Tier1 and its Tier2s in a certain geographical area!

# ATLAS computing model

## Clouds:

- **CERN, CA, DE, ES, FR, IT, ND, NL, RU, TW, UK, US**

## The Iberian Cloud (ES) inside ATLAS:

- **Tier1: PIC Barcelona**
- Tier2s:
  - Federated Spanish Tier2
    - IFIC Valencia
    - IFAE Barcelona
    - UAM Madrid
  - LIP Lisbon
  - UTFSM Santiago, Chile
  - UNLP La Paz, Argentina (inactive)

# Iberian Cloud Facilities

| Site | CPU (HEP-SPEC06) | DISK (TB) | TAPES (PB) | Availability (2018) | Reliability (2018) |
|---|---|---|---|---|---|
| **PIC-Tier1** | 40024 | 2400 | 9.6 | 98.76% | 99.60% |
| IFIC-Valencia | 26751 | 2146 | | 97.93% | 98.33% |
| IFAE-Barcelona | 10420 | 980 | | 99.22% | 99.59% |
| UAM-Madrid | 10358 | 1220 | | 99.05% | 99.66% |
| NCG-Lisbon | 4000 | 220 | | 91.00% | 92.42% |

- **Tier1: PIC-Barcelona**
- **Spanish Tier2:     IFIC-Valencia, IFAE-Barcelona, UAM- Madrid**
- **Portuguese Tier2: NCG-Lisbon**
- **Integrated in the WCLG project (World Wide LHC Computing GRID)**
  **Our cloud represents 5% of the total Tier-2 resources and 5% of the Tier-1 resources**

# Original ATLAS Computing Model

- Tier1 has associated Tier2s that are close to it in terms of network connectivity, and they form the "cloud".

- All data flow to and from Tier2s goes via its Tier1

More and more Tier2s have very good worldwide network connections and could exchange data directly between them.

This leads to the

# Nucleus ⟷ Satellite Model

- Tier2s with a big amount of storage and very good network connection get elected "Nucleus", passing job production on to smaller Tier2s (Satellites) **in any cloud**, exchanging data directly.

- **The Spanish Federated Tier2 is a Nucleus**.

# Small storage sites

- ADC has started a policy to transform Tier2 sites with "very small" storage to **diskless**
- Due to the change of the ATLAS model away from the static Tier1/Tier2 local connection, a site can use the storage of a "nearby" site if the network connection is good enough
- Lowers trouble for central ADC about solving issues on many small sites
- "Small" at the moment is defined as <400 TB, but will gradually increase over the years, probably to something like 1 PB
- Portuguese Tier2 already got an invitation,
  Studying the impact at the moment.
  It could close doors in the future.

# ADC Development Activities in Iberian Groups (1)

- Monitoring
  - Monitoring frontier-servers
  - IFIC transfers monitoring
  - Site and cloud support tools
  - ADC Live Page
- Participation in the DOMA-TPC tests and storage system performance studies (for the implementation of the tape carousel) led by ATLAS.
All of them are addressing the HL-LHC challenges.

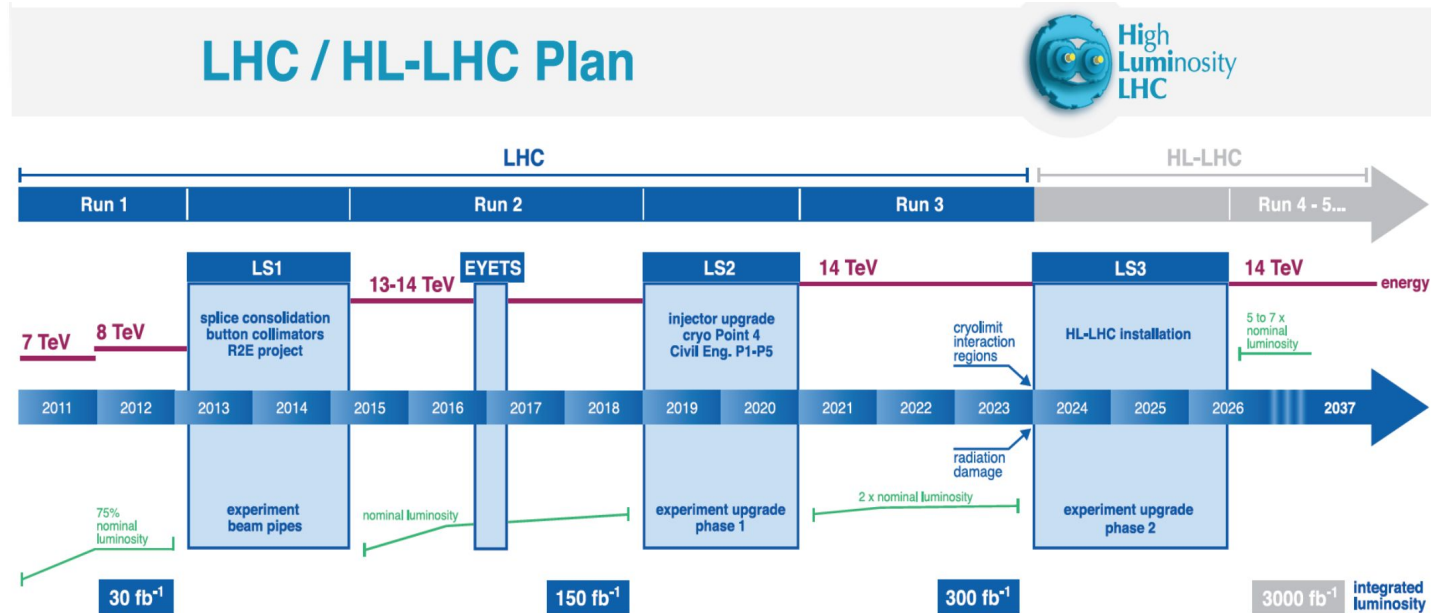# ADC Development activities in Iberian groups (2)

- **Event Index Project**
  - provide a catalog of data of all events in all processing stages needed to meet multiple use cases and search criteria.
  - Billion of events have been indexed so far (PetaBytes)!
- **Event Service**
  - Main goal: allow a more flexible and efficient usage of CPUs available when running simulation ATLAS jobs
- **Physics Case:**
  Selection of  events  with t tbar resonances (BSM)  from the SM events (background)  in  collisions pp  in the ATLAS Experiment using Machine Learning methods
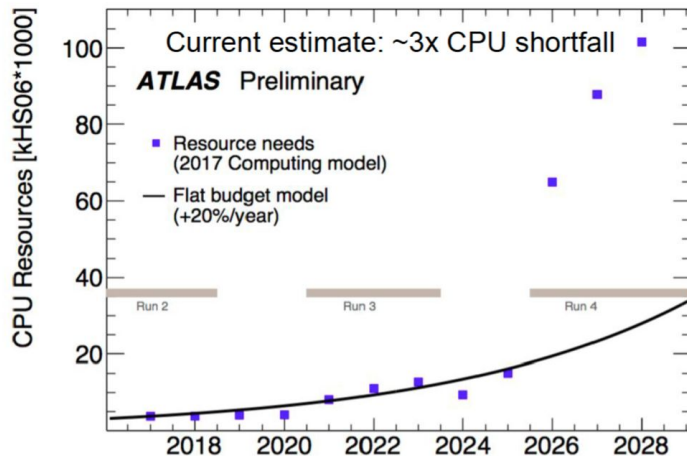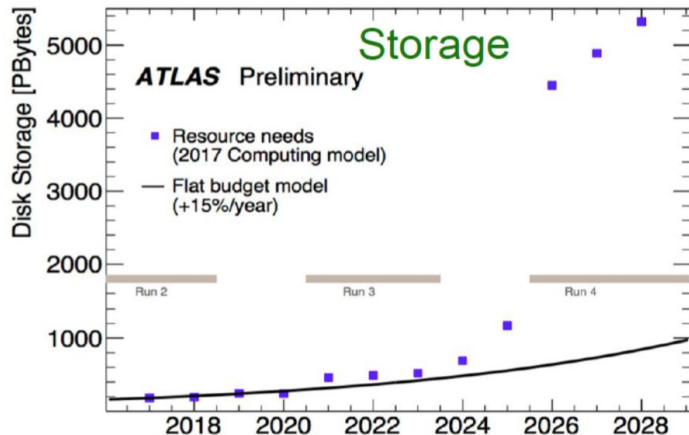
# Challenges for Run 3 and HL-LHC



**Higher luminosity is equivalent to higher data flow.**
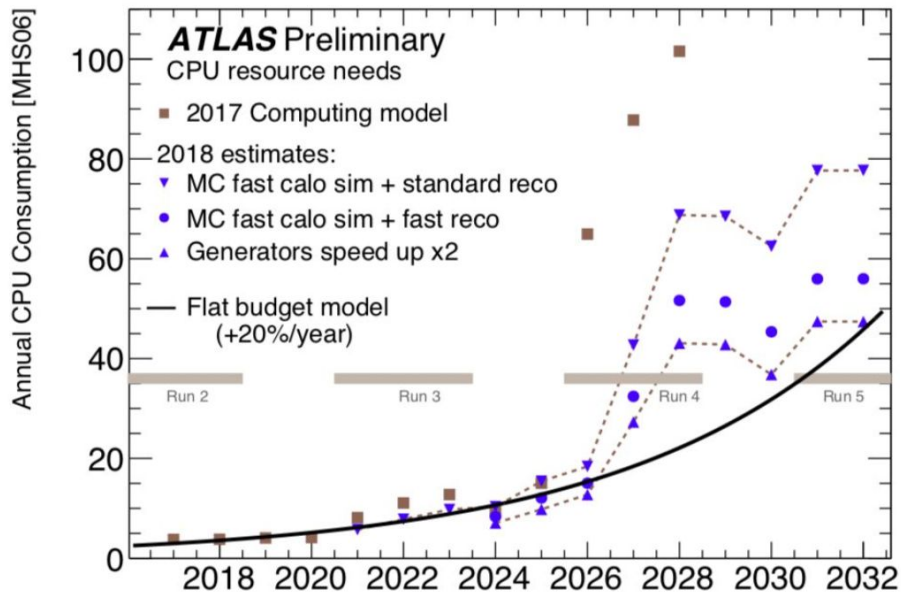**So there is an increase of one order of magnitude on the horizon!**
**Challenges on storage, network bandwidth and processing power.**

# CPU, disk storage and bandwidth requirements prediction

- **HL-LHC CPU estimations showed a ~3x shortfall with respect to the flat budget model**
- **~6x shortfall by today's estimate in Storage on Disk. Storage shortfall is our biggest problem**
- **HL-LHC will require to increase the network bandwidth by a factor 10**
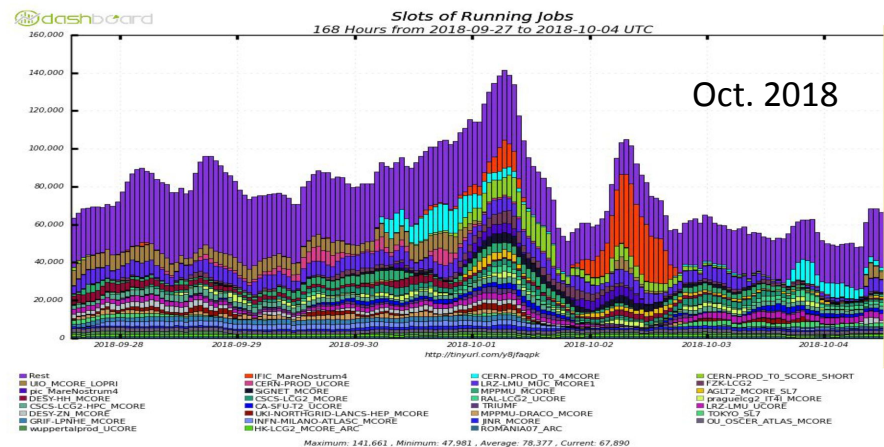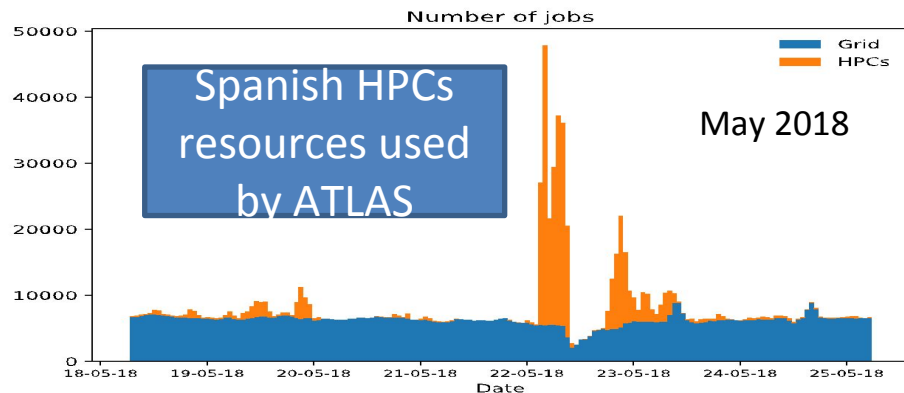
# Approaches to solve CPU shortfall



- There are a few options to face this challenge: HPC's, cloud computing and High Level Trigger Farm.

- Further options: use fast simulation instead of full one. And speed up the MC generators by a factor two.

- Running on GPU's is also feasible, but needs significantly time and effort to adapt our software to new architecture
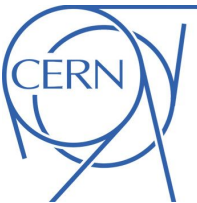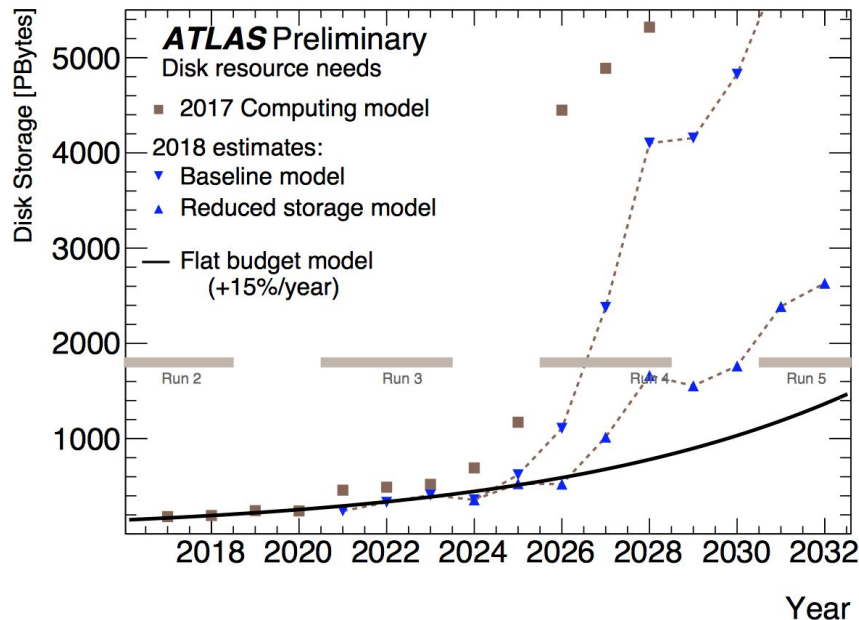
# Use of HPC resources

- Opportunistic resources turn out to be a meaningful way to face the future HL-LHC challenges in terms of CPU requirements
- High Performance Centers (HPC's) have been tested recently. Further work is required to use these resources since they aren't available from the ATLAS GRID.
- In a collaboration between IFIC and IFAE, we received hours in Spanish HPC's (RES and PRACE):
  **Mare Nostrum:** 4 million hours (PIC and IFIC)
  **Lusitania:**        2 million hours (IFIC)
- **Cibeles** in Madrid, opportunistic
- We plan to try a similar exercise at the **BOB supercomputer in Minho**, a new Portuguese HPC infrastructure, as opportunistic resource

- IFIC led ATLAS simulation when profiting of opportunistic HPC resources

# Approaches to solve Storage shortfall

**No opportunistic storage…so far**



- Increase investment in computing

- New file formats (to reduce filesize, many data formats for physics analysis)

- "Less data"

- Use of tapes. But this option slows down the workflow

- Data Lakes / DOMA

# Conclusions

- The Iberian Cloud contributes around 5% of the total resources deployed in the Tier-1 and Tier-2 sites.

- Spanish Tier-2 has the so-called "nucleus" status in ATLAS.
  ⇨ Major responsibilities and larger work volume!

- Portuguese Tier-2 getting stable and reliable again after storage migration to new hardware

- Not only deployment of CPU and storage resources for the ATLAS experiment but also several researching activities are carried out by the teams in the Iberian Cloud.

- Run-3 and HL-LHC defy the current ATLAS computing model. CPU resources, disk storage and bandwidth shortfalls have to be assessed and faced asap.

  - Usage of opportunistic resources (HPC)
  - Data format, data processing and storage..
  - Bandwidth requirements seem to be satisfied in time

# Thanks for your attention!

Helmut Wolters (LIP)

Santiago González de la Hoz (IFIC)

Carlos Acosta-Silva (IFAE,PIC)

Javier Aparisi Pozo (IFIC)

Mário David (LIP)

Jose Del Peso (UAM)

Álvaro Fernández Casani (IFIC)

José Flix Molina (PIC,CIEMAT)

Esteban Fullana Torregrosa (IFIC)

Carlos García Montoro (IFIC)

Jorge Gomes (LIP)

Julio Lozano Bahilo (IFIC)

João Paulo Martins (LIP)

Gonzalo Merino (CIEMAT,PIC)

Almudena del Rocio Montiel (UAM)

Andreu Pacheco Pages (IFAE,PIC)

João Pina (LIP)

Javier Sánchez Martínez (IFIC)
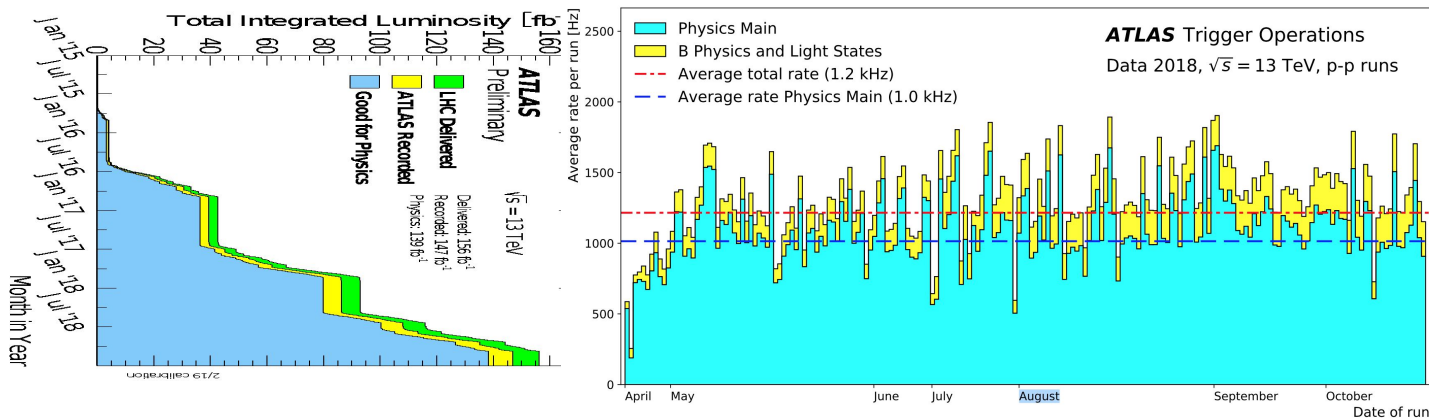
José Salt (IFIC)

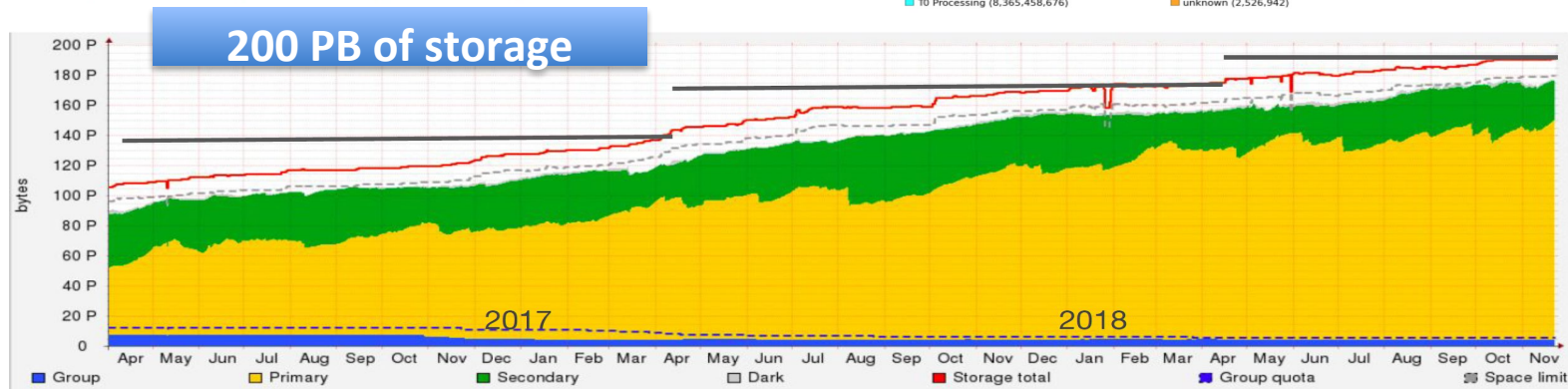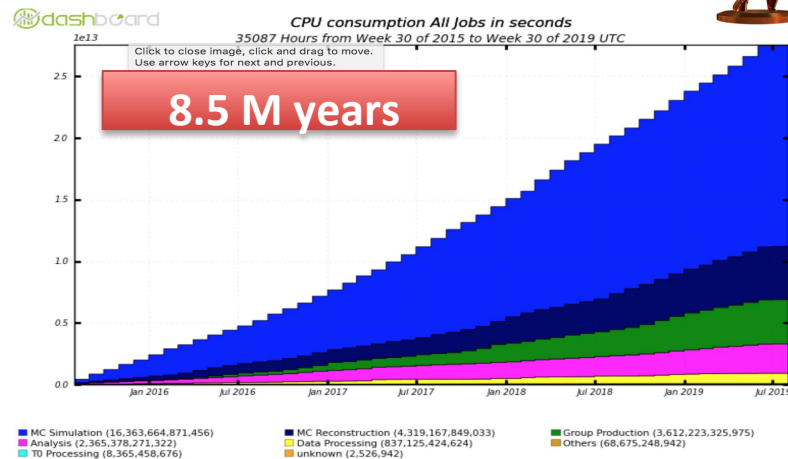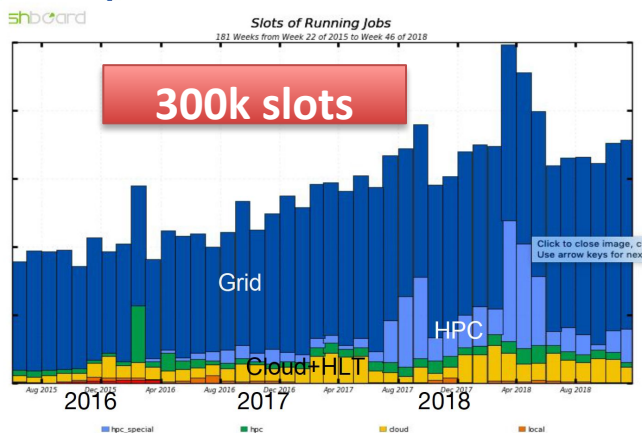Aresh Vedaee (IFAE,PIC)

# BACKUP

# ATLAS event recording performance



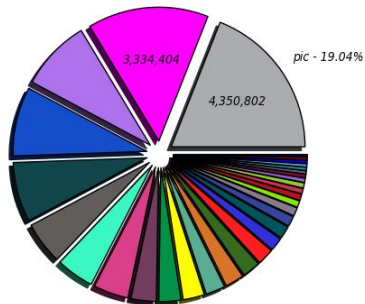**94% of the luminosity delivered by LHC is collected!!**

**Event Rate in Run I (HLT readout): 300 Hz**
**Event Rate in Run II (HLT readout): 1.2 kHz**
**Expected Event Rate in Run III (HLT readout): 5-10 kHz**
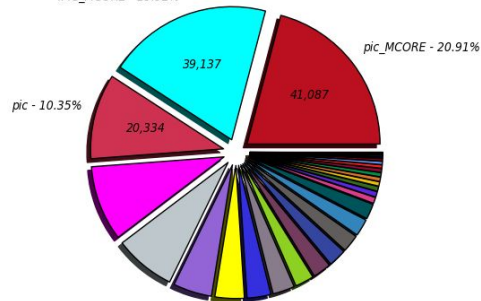
# ATLAS GRID performance



**300k slots**

**8.5 M years**

**200 PB of storage**

# 3.1 Spanish Cloud performance in Run II



Completed jobs (Sum: 22,854,366)
IFIC - 14.59%
3,334,404
pic - 19.04%
4,350,802

**More than 22 million finished jobs**



NEvents Processed in MEvents (Million Events) (Sum: 196,466)
IFIC_MCORE - 19.92%
39,137
pic_MCORE - 20.91%
41,087
pic - 10.35%
20,334

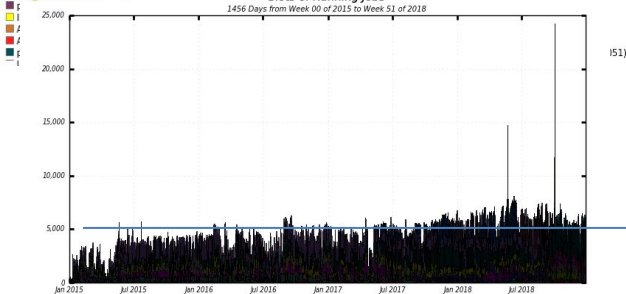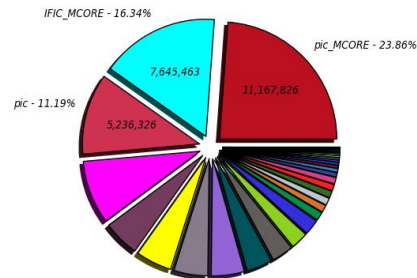**More than 196 million events proccessed**



Slots of Running Jobs
1456 Days from Week 00 of 2015 to Week 51 of 2018

**On average, 5000 slots occupied by running jobs daily**



NFiles Produced (Pie Graph) (Sum: 46,801,836)
IFIC_MCORE - 16.34%
7,645,463
pic_MCORE - 23.86%
11,167,826
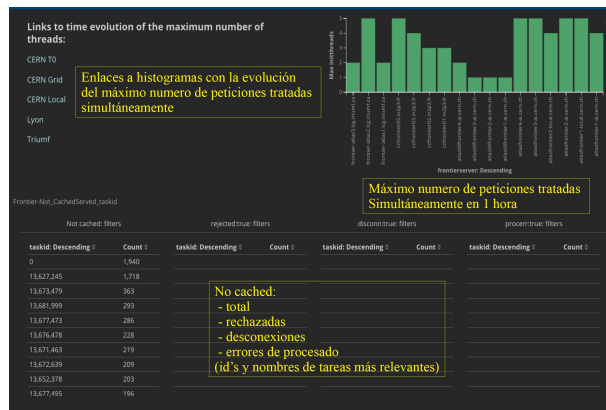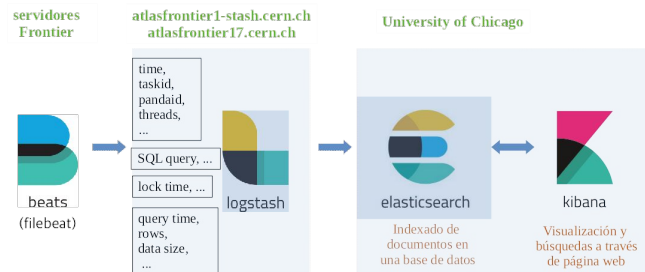pic - 11.19%
5,236,326

**More than 46 million files produced**
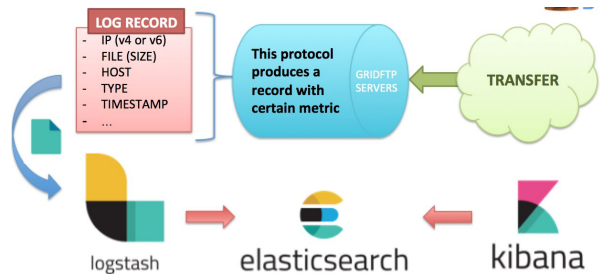
# IFIC activities
# Monitoring frontier-servers

- *Frontier servers* optimize the access to the so-called "conditions database" (variables of the ATLAS detector), needed to run simulation or production jobs.
  - A '*squid*' server provides '*caching*'of data
  - A '*servlet*' Tomcat connects to the Oracle database when needed
- The monitoring system collects information about the *queries* from '*log*' files:
- System operates in a steady and stable way processing 12M of *queries daily on average*
- Allows the visualization of meaningful information by means of *Dashboard. Namely, summary tables and histograms*
- Incorporates resources to send e-mail warning when a site performance deteriorates
- Thanks to its versatility and the close relation between ATLAS and CMS on this project, some CMS servers are
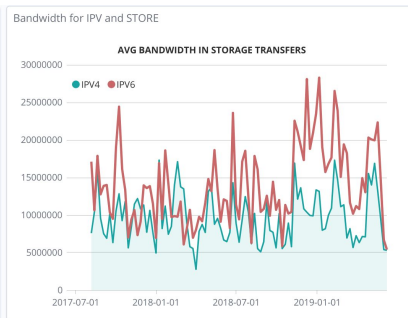
# IFIC transfers monitoring



By means of ELK stack, we can filter, store, analyze and display huge amount of information of the transfers made from and to our Tier-2 center

# Event Index Project

- The EventIndex Project aims to provide a catalog of data of all events in all processing stages needed to meet multiple use cases and search criteria. Only meaningful information is indexed.

- Billion of events have been indexed so far (PetaBytes)!

- Some use cases:
    1. Event picking
    2. Duplicate event checking
    3. Overlap detection
    4. Trigger checks and event skimming
    5. Trigger Counter

# Event Service

- Main goal: a more flexible and efficient usage of the CPU available when running simulation ATLAS jobs. Improving performance in HPC's
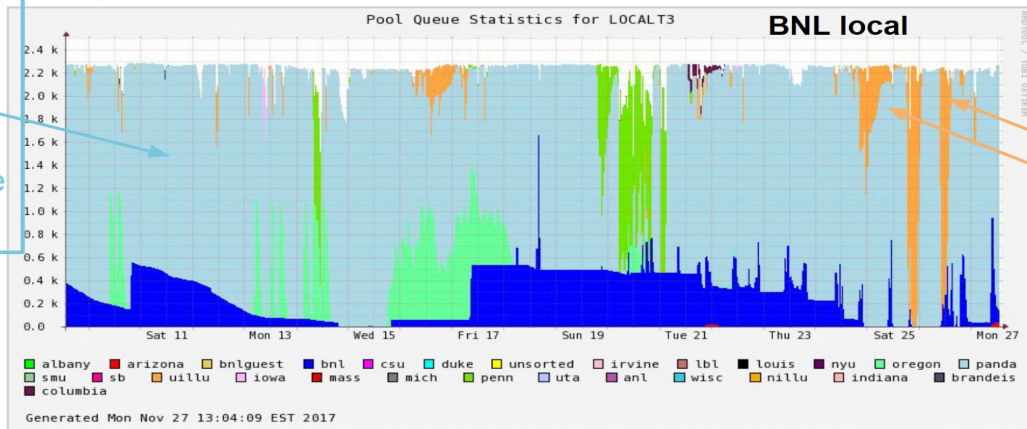  1. Splitting event bunch in GRID jobs: from 1000 to 1!!!
  2. In tier-3 facilities, when no user analysis are running, ES uses this CPU

- **Efficiently and flexibly exploit any CPUs available**
  - Efficient use of **opportunistic** and volatile resources
    - Dynamically making use of available CPUs as they appear



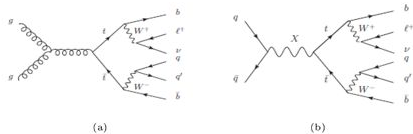When CPU is not used by local users, EventService can use it

BNL local

ES jobs killed, restarted, re-killed, etc...
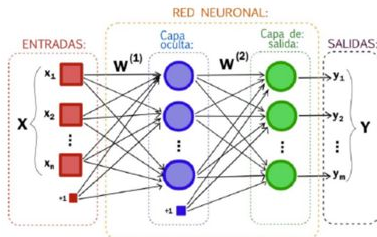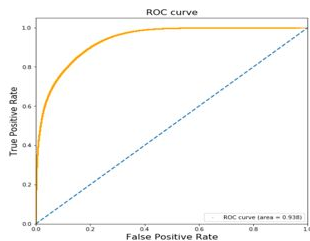
See other example by Rod

# Summary: ML @ T2-ATLAS-IFIC project (J. Lozano and J. Salt)

- **Physics Case** 1: Selection of events with t tbar resonances (BSM) from the SM events (background) in collisions pp in the ATLAS Experiment using ML methods



a) Diagrama de Feynman del modelo SM (Fondo). b) Diagrama de Feynman del modelo BSM (Señal).

**21 low-level** + **5 high-level variables**
(kinematic variables) (Invariant masses of decay schema processes)



| Dataset | NN_MLP | NN_Shallow | NN_keras | NN-sklearn | R.F. | R. LIN | R. LOG |
|---------|--------|------------|----------|-----------|------|--------|--------|
| 500000  | 0.931  | 0.924      | ---      | 0.935     | 0.934| 0.909  | 0.909  |
| 730000  | 0.938  | 0.937      | 0.939    | -----     | -----| 0.910  | 0.910  |

NN_MLP: Neural Network – Multilayer perceptron      NN_sktlearn: Neural Network- using sktlearn

NN_Shallow: Neural Network- Shallow          R.F. : Random Forest
NN_Keras: Neural Network_ using Keras         R.LIN: Linear Regression
                              R. LOG. Logistic Regression



AUC vs mass of the signal (ttbar resonance)
LR: Logistic regression
ETC: Extra Trees Classifiers
RFC: Random fFores classifiers

- **Physics Case** 2: Searching for DM in ATLAS experiment by applying ML methods to detect Outliers          just starting

## Use of ARTEMISA facility:

**ARTificial Environment for ML and Innovation in Scientific Advanced computing ( computing resources based on GPU's)**

Project "Application of ML methods for studies on New Physics in ATLAS"
Use Case: ttbar resonances

It will be used to find the optimal configuration of the neural network or the random forest. There are algorithms in Keras and in Scikit-learn that you give them some ranges for the typical parameters: number of layers, number of neurons, activation function, etc. and they alone test all the permutations and give you the best according to several criteria.

# HPC Usage 2019

**PIC:**

| MareNostrum | Scheduled (kHours) | Used (kHours) | % | |
|---|---|---|---|---|
| RES FI-2019-1-0035 | 700 | 822.08 | 117 | completed |
| RES FI-2019-2-0030 | 2000 | 703.55 | 35 | grant valid until end of October |
| PRACE 2010PA5027 | 50 | 54.18 | 108 | completed |

Around 2.75 Mhours with 1.2 Mhours granted to IFIC gives 4 millions hours granted in MareNostrum4 this year.

**IFIC:** this year 1.2 Mhours in **MareNostrum** , and 2 Mhours in **Lusitania** both of them already consumed

**UAM:** **Cibeles** in Madrid, opportunistic

**LIP:** **BOB** (Minho), planned, opportunistic