

Barcelona Supercomputer Center Integration in the computing of ATLAS

Andrés Pacheco Pages

IFAE Pizza Seminar - Wednesday 29 April 2020

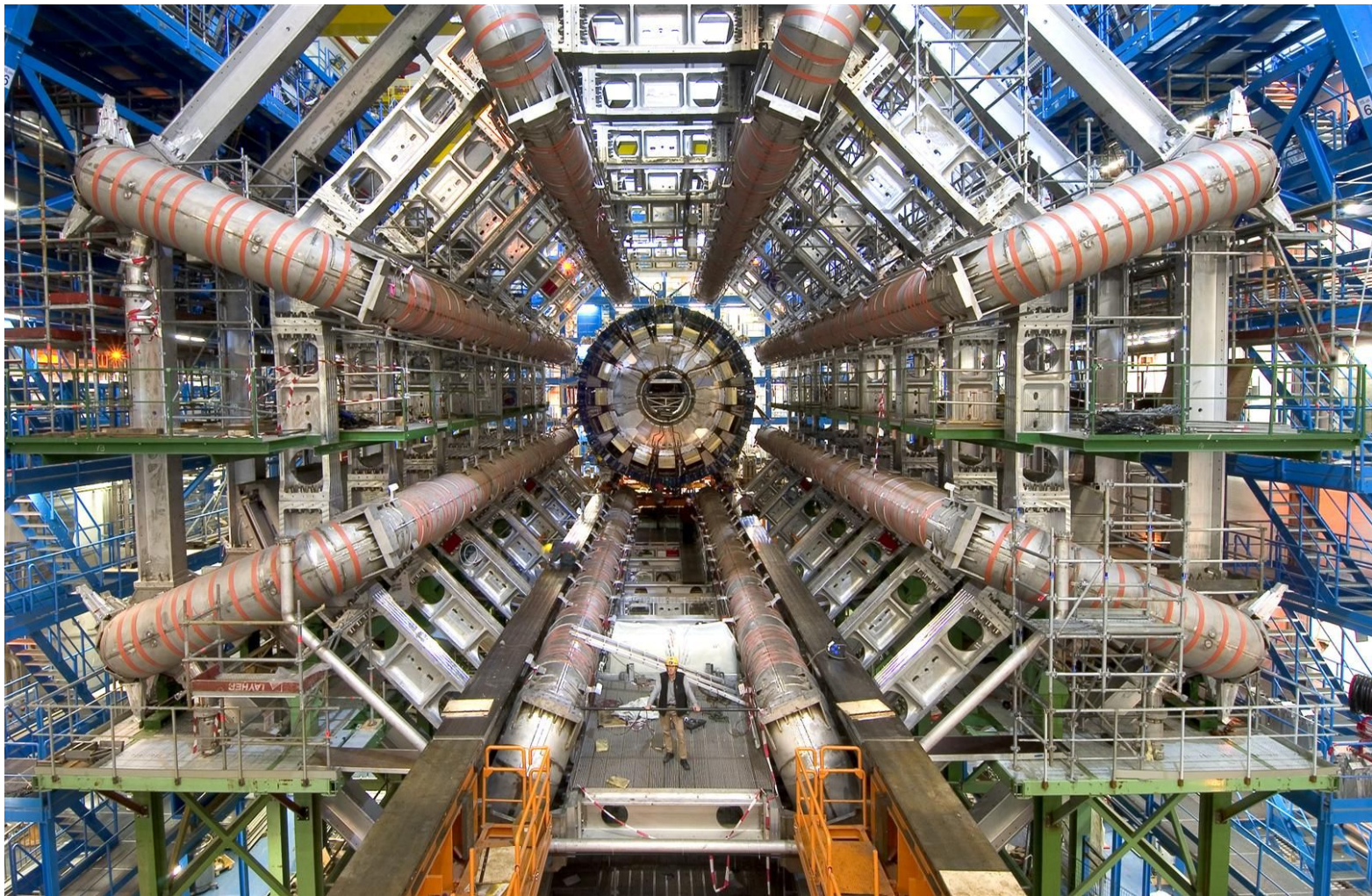


Centro de Investigaciones
Energéticas, Medioambientales
y Tecnológicas





A. Pacheco Pages - Pizza Seminar - Wednesday 29 April 2020

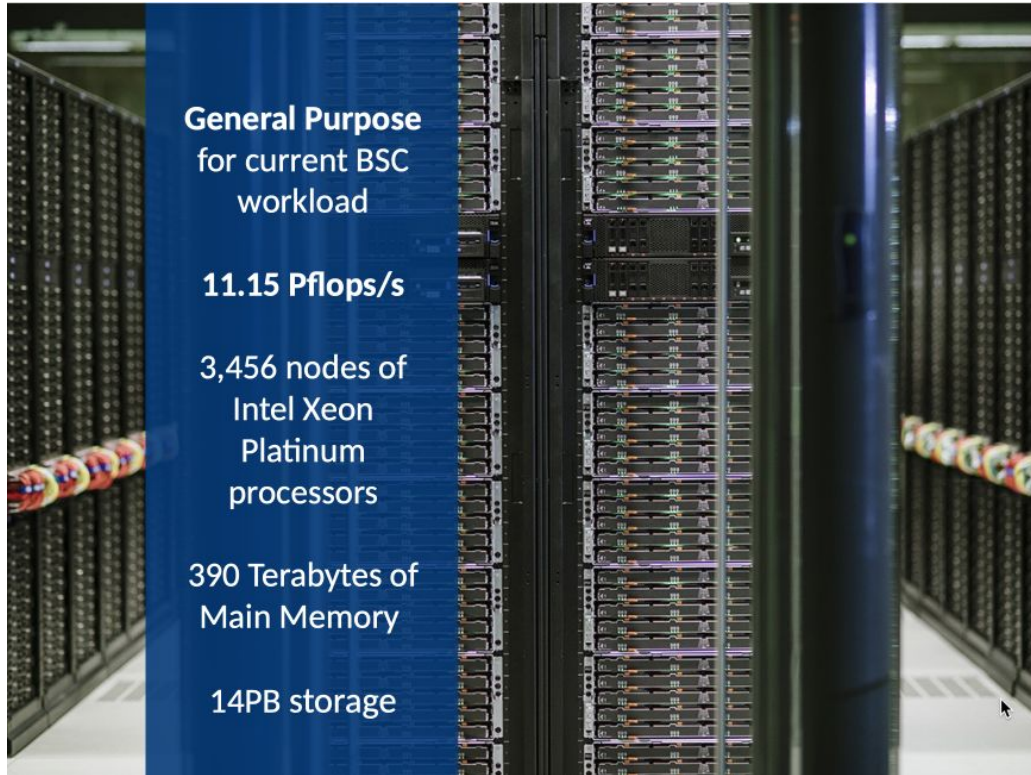


A. Pacheco Pages - Pizza Seminar - Wednesday 29 April 2020



A. Pacheco Pages - Pizza Seminar - Wednesday 29 April 2020

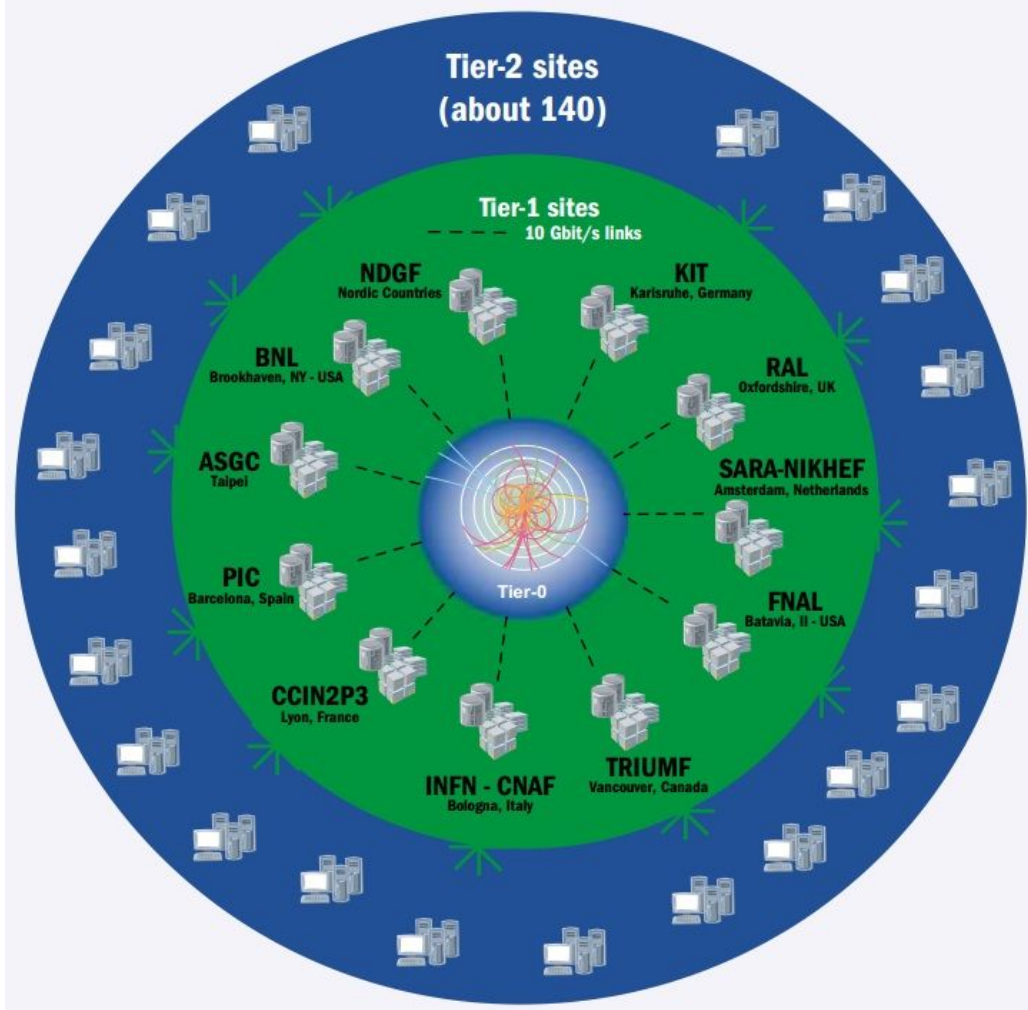
MareNostrum4 Picture

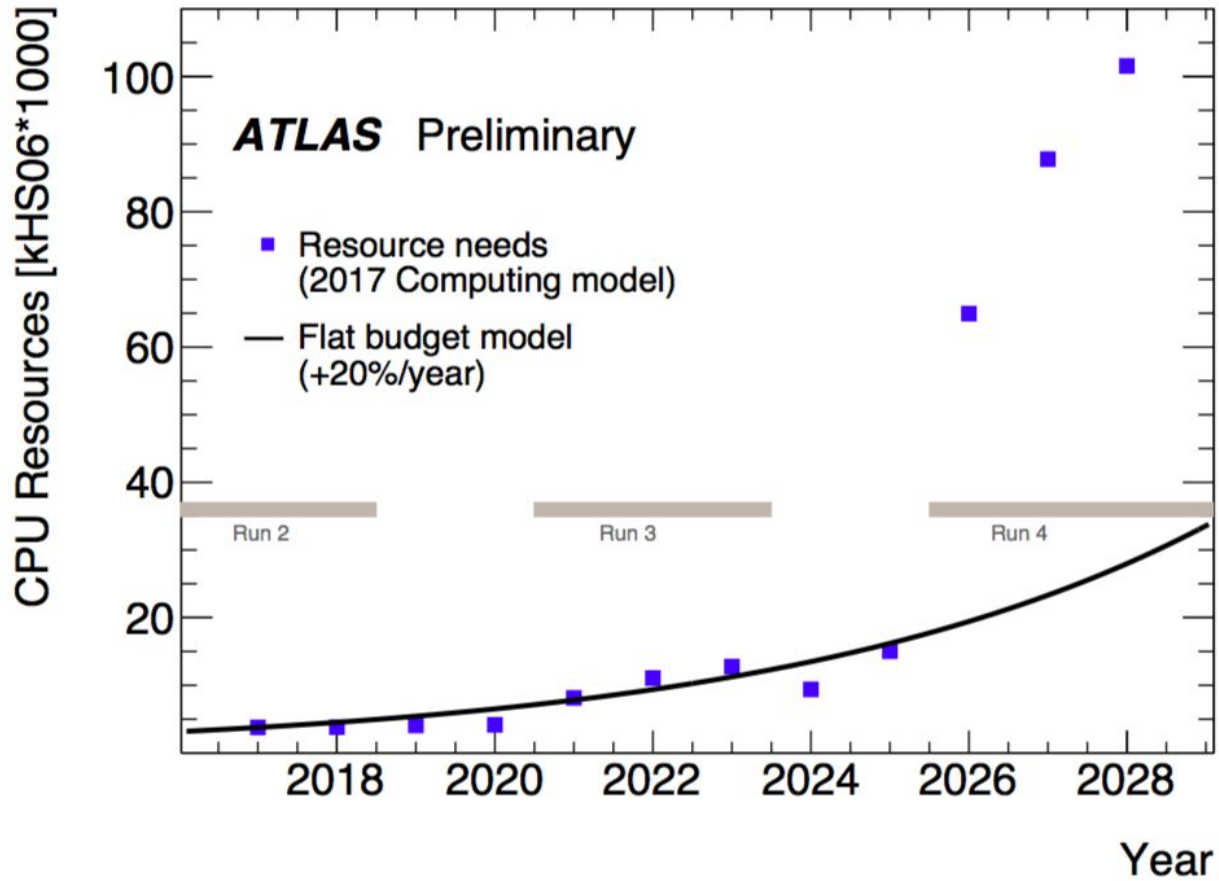


- Each node has two Intel Xeon Platinum chips, each with 24 processors, amounting to a **total of 165,888 processors** and a main memory of **2 GB RAM per processor**.
- Batch system: **SLURM**
- Operating system: **SUSE Linux 6**
- Shared file system: **GPFS**

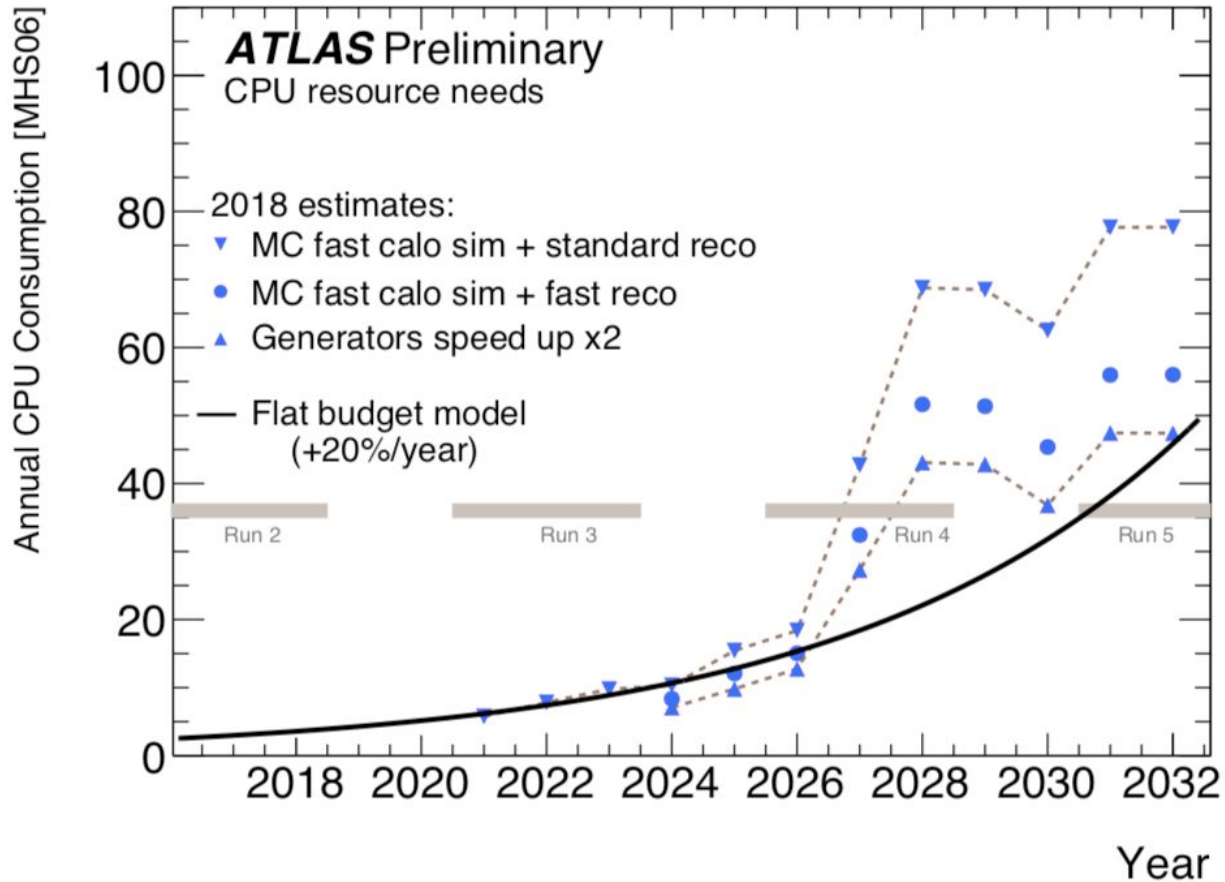


Minotauro at BSC





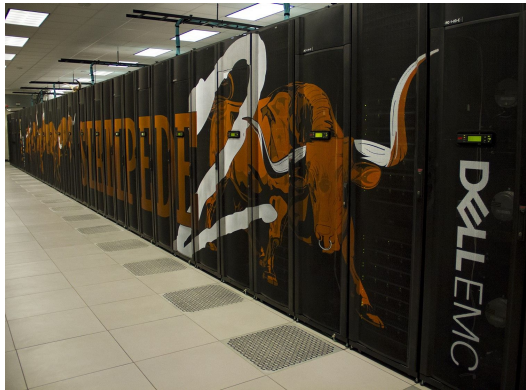
OLD ESTIMATES



NEW ESTIMATES



A. Pacheco Pages - Pizza Seminar - Wednesday 29 April 2020

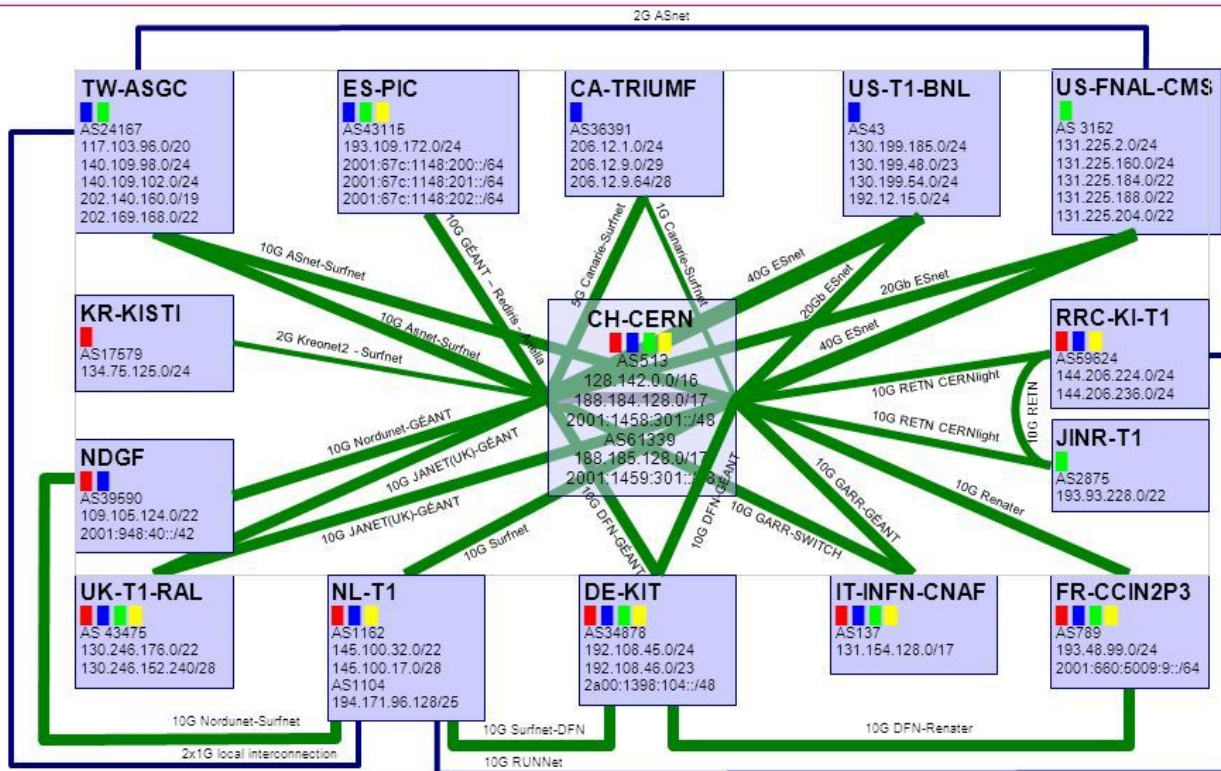


A. Pacheco Pages - Pizza Seminar - Wednesday 29 April 2020

— T0-T1 and T1-T1 traffic
— GEANT provided links
— T1-T1 traffic only
- - - Not deployed yet
 (thick) $\geq 10\text{Gbps}$
 (thin) $< 10\text{Gbps}$

■ = Alice ■ = Atlas
■ = CMS ■ = LHCb

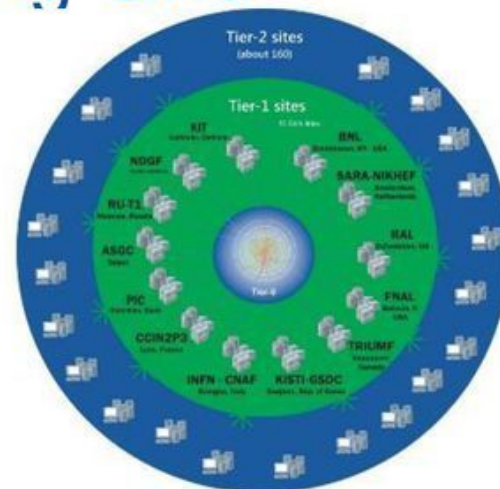
p2p prefix: 192.16.166.0/24 - 2001:1458:302::/48
 edoardo.martelli@cern.ch 20141212



LHC Data Distribution: WLCG

Worldwide LHC Computing Grid

- The Worldwide LHC Computing Grid (WLCG) is a global collaboration of **170 data centres around the world, in 42 countries**
- The CERN data centre (Tier-0) distributes the LHC data worldwide to the other WLCG sites (Tier-1 and Tier-2)
- WLCG provides global computing resources to store, distribute and analyse the LHC data
- The resources are distributed – for funding and sociological reasons



Tier-0 (CERN):

- Initial data reconstruction
- Data distribution
- **Data recording & archiving**

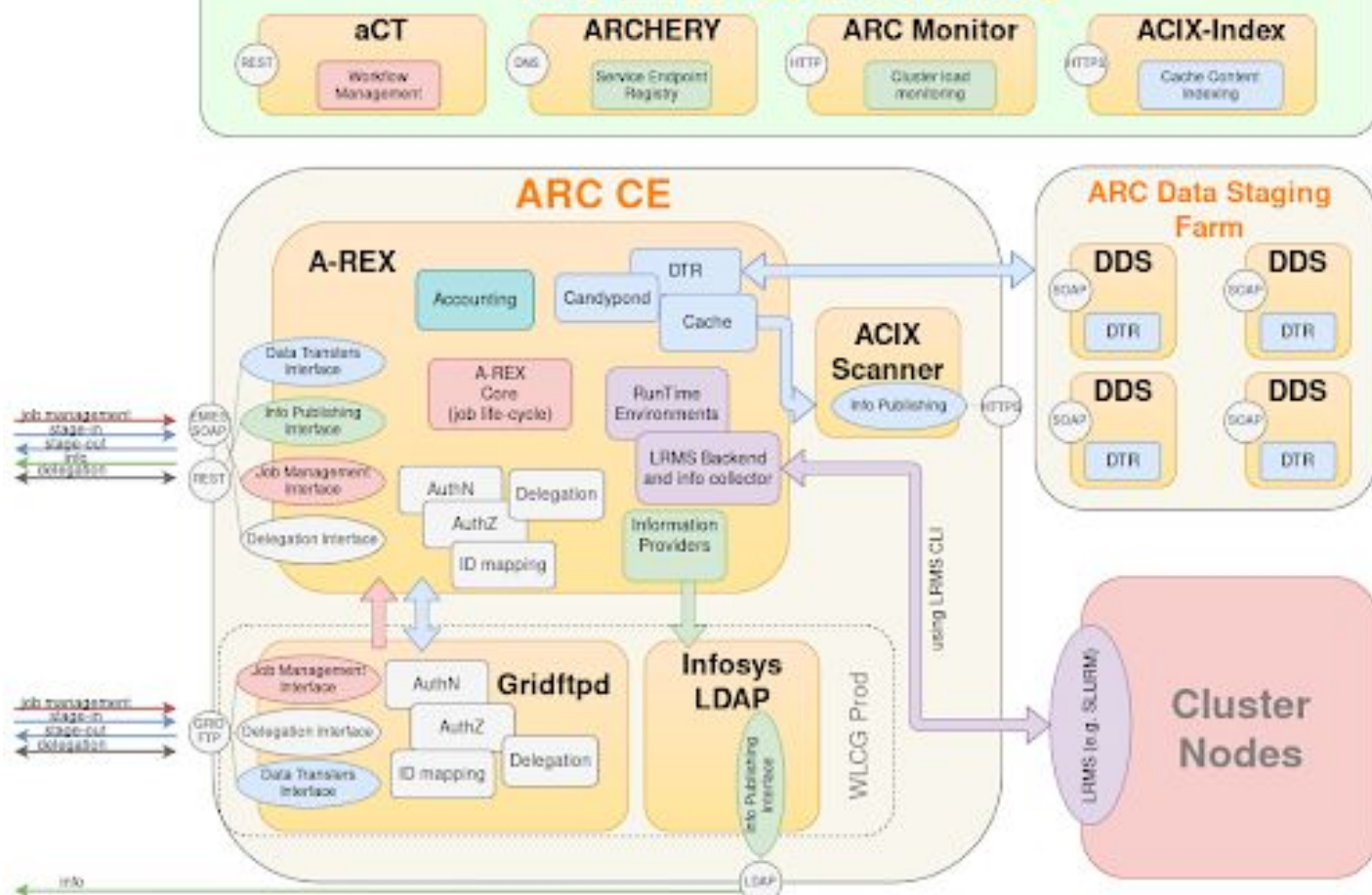
Tier-1 (13 centres):

- **Permanent storage**
- Re-processing
- Analysis

Tier-2 (~140 centres):

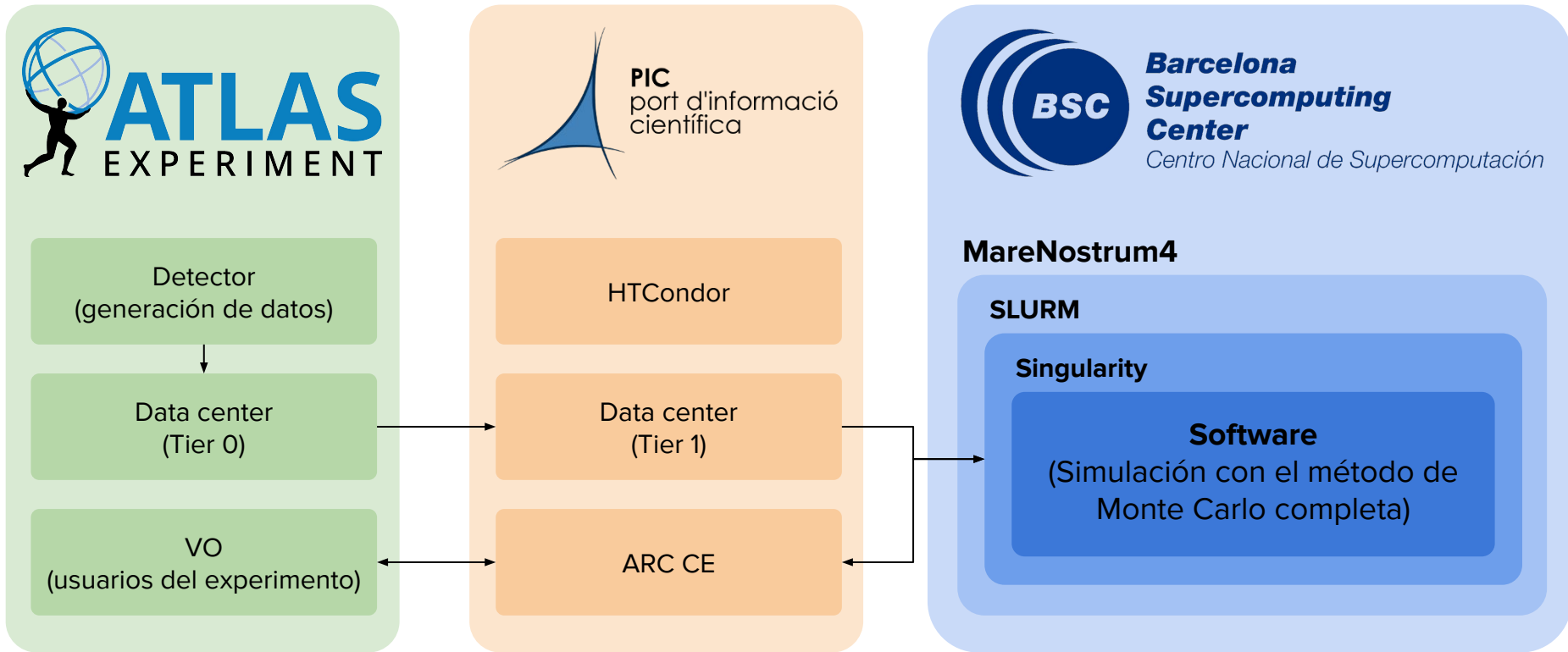
- Simulation
- End-user analysis

ARC Infrastructure Services



Workflow at MareNostrum4?

- We must **copy all the input files** using the DTN by mounting a sshfs file system between PIC and BSC.
- We must **submit the jobs** using the login nodes and running on validated Singularity images with all the software preloaded.
- We must **check the status of the jobs** using the login nodes.
- We must **retrieve the output files** using the sshfs filesystem.

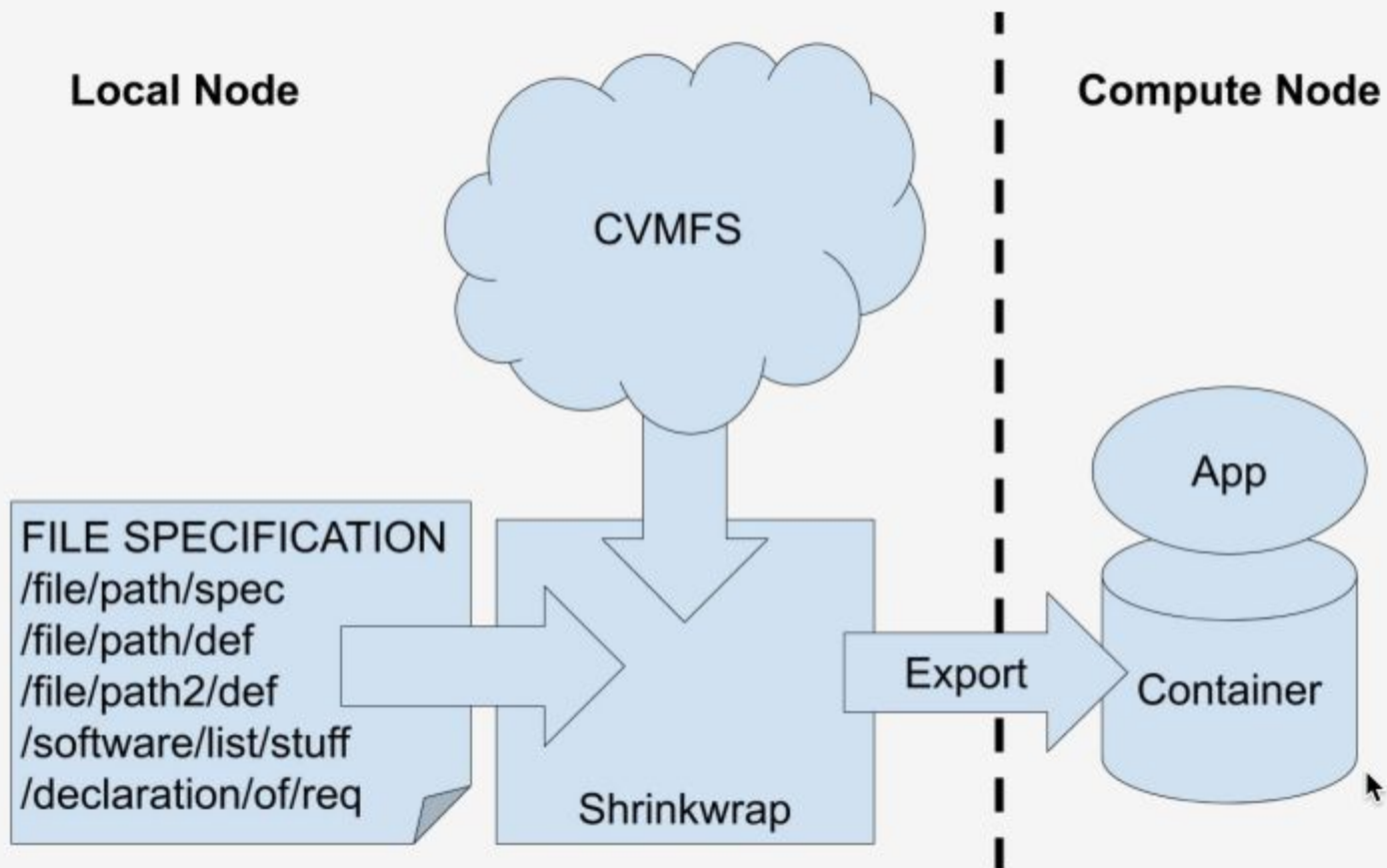


Pipeline

How to solve the problem of running on isolated worker nodes?

- The **working solution** now is to **create a filesystem with a partial copy of the ATLAS CVMFS filesystem repository** and including files containing **detector conditions**. The latest tool is called Shrinkwrap.
- This works because the **releases used for simulation are very few**.
- Then the filesystem is copied inside a **Singularity image** running a validated operating system (**CC7**).
- The **“problem”** is to **find the right list of files to be copied** to the image and the balance of the number of images to maintain: one image per ATLAS release, per workflow,... just run the parrot utility on a workflow to get an idea of the list of files accessed from cvmfs... thousands.

Use Case: HPC Environments



How do we get grants at MareNostrum4? RES

- The main **source of allocation of cpu hours** come from the “**Red Española de Supercomputación**” (**RES**) competitive program.
- Web is www.bsc.es/res
- You enter, you register, you request the time and then **you get approved or denied every 4 months**.
- You can **get hours** allocated in **any center of the RES**.

Home

IMPORTANT: For security reasons, we recommend you to change your password once a year. If you have not changed your RES password in the last 12 months, please click [here](#).

NEWS: New [Frequently Asked Questions \(FAQ\)](#) for RES applicants and users



Applications and activities

Below you have a list of your available applications and activities.

[+ New Application](#)

The deadline for next period applications is 12/05/2020 11:00:00, CEST.

Time left: **2 weeks 0 day 09:50:33**

[Add Publications](#)

Please, if you have any new publication add it to your dissemination information.

[Next Period](#) [Current Period](#) [Past Periods](#)

[Click here](#) to show / hide meaning of icons and colors area.

Current Period Applications and Activities

2020, March 1st - 2020, June 30th **2020-1**

FI-2020-1-0027

Monte Carlo Simulation for the ATLAS Experiment at the CERN LHC at the MareNostrum

OTA

[View Application](#)

[Technical Information](#)

[Reviews](#)

[View Reports](#)

[Add Report](#)

[View CPU Usage](#)

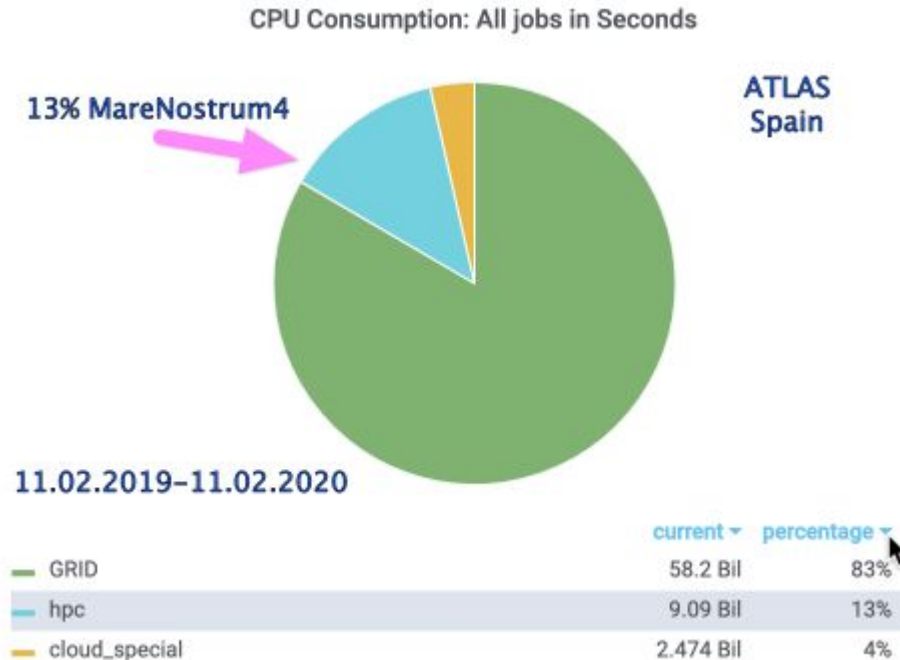
[Manage Users](#)

How do we get grants at MareNostrum4? PRACE

- **Another program** we can apply for resources at BSC is **PRACE (Partnership for Advanced Computing in Europe)**
- Web is: <http://www.prace-ri.eu/how-to-apply/>
- There are **several types of calls from 2 months till 1 year**. You can get the allocation at the MareNostrum4 or at any of the HPCs in Europe. You select which you want explicitly.
- The **smallest grant is 2 months and 50 khours**. PRACE Preparatory Access type A.

	Distribution	Description		Access Criteria
Peer Reviewed evaluation	40%	PRACE		PRACE Access committee
	40%	RES	General access	RES Access Committee
			Time in advance	
			Unexperienced users	
Strategic projects (up to 7%)		Board of Trustees + RES annual evaluation		
Expost evaluation	5%	Director time		Director BSC
	15%	BSC		Internal BSC

CPU from ATLAS jobs in Spanish sites: 13% correspond to jobs in MN4



- On the left, we have the CPU consumption pie chart of ATLAS jobs by resource type 1 year to date.
- ATLAS has already got off-pledge 13% of the Spanish contribution to the CPU from MareNostrum4 using queues at IFIC and PIC.

Plans and the next move

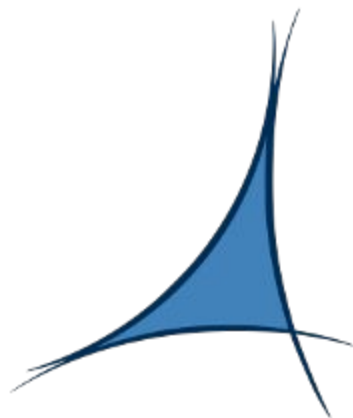
- Current plans is to increase the use of BSC thanks to the strategic program.
- We plan at PIC to run 1 million hours per month and increase each quadrimester.
- We need some work to increase the types of simulations we can run.
- After simulation the next target are the analysis jobs in containerized images.
 - Useful for analysis using GPUs

Can we replace the LHC computer centers?

- The answer is not.
- We need at least grid centers to receive the data from the experiment, store it on disk and **tape**, distribute, and **reprocess** the data. As well as simulate and analyze.
- The same is valid for simulated data once is produced, needs to be archived.
- The reconstruction of the data needs access to the databases of detector information, which is hard to upload to any supercomputer center.

Summary and conclusions

- We have managed to integrate the ATLAS Simulation jobs into the MareNostrum 4.
- The BSC has included the LHC computing in the list of strategic projects.
- We expect that the transition to MareNostrum 5 can be straightforward with 17 times more computing power in 2021.
- We still need grid computing for the LHC
 - Still many workflows cannot run in the BSC due to the lack of connectivity
 - We need to store, distribute and archive to tape the data.
- Thanks to the work of Carlos Acosta (PIC) and Elvis Diaz (UAB Student), all the PIC team and the collaboration with IFIC.



PIC
port d'informació
científica

FASES DE LA DESESCALADA

MAYO							JUNIO						
L	M	X	J	V	S	D	L	M	X	J	V	S	D
	28	29	30	1	2	3	1	2	3	4	5	6	7
4	5	6	7	8	9	10	8	9	10	11	12	13	14
11	12	13	14	15	16	17	15	16	17	18	19	20	21
18	19	20	21	22	23	24	22	23	24	25	26	27	28
25	26	27	28	29	30	31	29	30					

FASE 0
 FASE 1
 FASE 2
 FASE 3

Salida controlada de menores	Apertura de comercio excepto centros comerciales.	Apertura del interior de locales a 1/2 de su aforo	Flexibilización de la movilidad general.
Deporte individual al aire libre	Apertura de restauración con un 30% de ocupación en terrazas.	Apertura excepcional de centros escolares para clases de refuerzo o la selectividad.	Centros comerciales a un 50% de su capacidad y con una distancia de 2 metros.
Locales con cita previa para tener llevar comida a domicilio	Apertura de hoteles y alojamientos turísticos excluyendo zonas comunes.	Cines, teatros y similares a 1/2 de su aforo.	Mayor aforo en restauración, preservando las distancias de seguridad.
Entrenamiento individual de deportistas profesionales	Lugares de culto a 1/2 de su capacidad	Equipamientos culturales a 1/2 de su aforo. Si es al aire libre, máximo 400 personas sentadas.	
Preparación de todos los locales públicos con medidas de protección	Apertura de centros de alto rendimiento deportivo.	Lugares de culto al 50% de su aforo.	

Conceptos generales:

- La provincia es la unidad territorial de medición.
- La duración de las fases será como mínimo de dos semanas cada una y el avance irá condicionado a indicadores de salud pública y a la evolución de los datos.
- La desescalada será asimétrica según la evolución de cada provincia.
- Inicio de la nueva normalidad como mínimo a partir del 25 de junio: Se permite el desplazamiento entre provincias y se siguen manteniendo las normas de seguridad y distancia social.