

# Verso il Tecnopolo (e oltre....)

CCR

9 Settembre 2019

Luca dell'Agnello

---

# The CNAF Data Center in a nutshell (1/2)

---

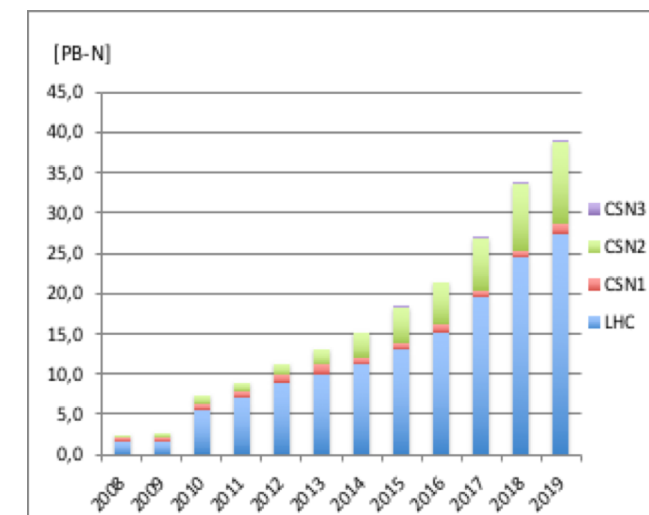
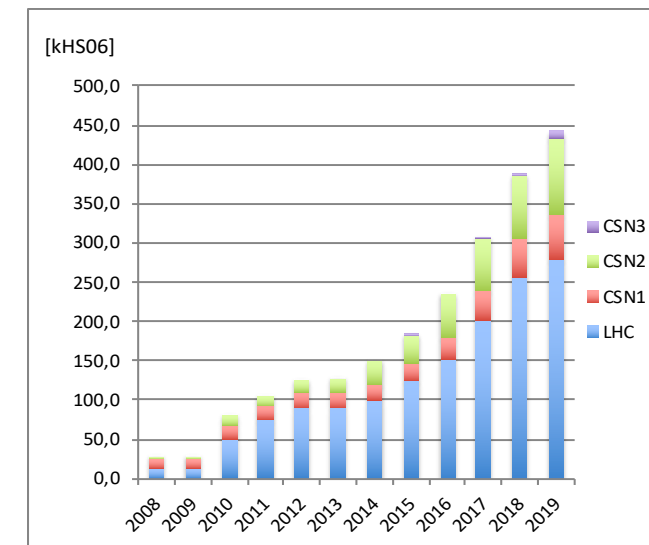
- First incarnation dates to 2003 as computing center for BaBar, CDF, Virgo and prototypical Tier-1 for LHC experiments (Alice, ATLAS, CMS, LHCb)
- Completely renewed in 2008, is now the main INFN data center
- ~400 kHS06 on production farm (~30K job slots with ~100k jobs/day on average)
  - Farm partitioned in 3 locations
    - ~210 kHS06 at CNAF data center
    - ~180 kHS06 at CINECA data center (in production since March 2018)
    - ~10 kHS06 at Bari-ReCaS data center (in production since beginning 2017)
  - Grid, local and (soon) cloud access supported
    - Possible to dynamically expand Cloud@CNAF on the Tier-1 farm
  - Migration to HtCondor ongoing
  - A secondary smaller farm manages the HPC clusters

# The CNAF Data Center in a nutshell (2/2)

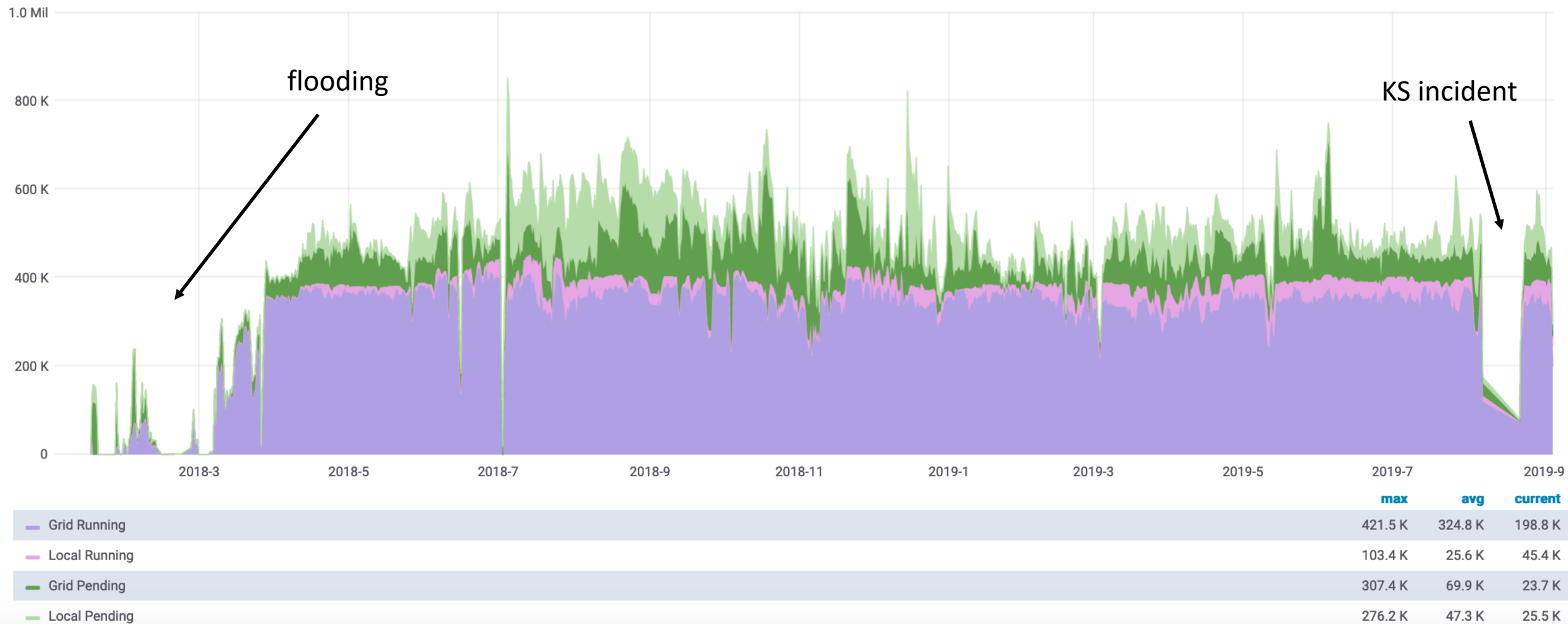
- Storage currently counts on ~34 net PB of disk and 84 PB of tapes
  - 1 tape library on line (16 drives, 8.4 TB tapes, 10000 slots)
  - A second library has been installed last week (19 drives, 20 TB tapes, ~5000 slots)
  - Data access – users can manage the data via StoRM and access the data via Posix, GridFTP, XRootD, WebDAV/HTTP
  - Completely redundant data access system
    - No single points of failure (server connections, SAN switches, controllers of the disk storage boxes)
- 640 Gbps aggregate bandwidth from CNAF
  - GPN to WAN (20 Gbps)
  - LHCOPN/ONE (2x100 Gbps)
  - 20 Gbps to Bari-RECAS
  - 400 Gbps to CINECA

# Our scientific collaborations

- Tier-1 for LHC experiments (ATLAS, CMS, ALICE and LHCb)
- Particle physics at accelerators
  - Belle2, CDF (LTDP), Compass, Kloe, LHCf, NA62, PADME (formerly also Babar and SuperB)
- Nuclear physics
  - Agata, ASFIN, Famu, FARCOS, Fazia, FOOT, n-TOF
- Astro and Space physics
  - AMS (Satellite), ARGO (Tibet), Auger (Argentina), EUCLID (Satellite), Fermi/GLAST (Satellite), LIMADOU (Satellite), LSPE (Balloon), MAGIC (Canary Islands), PAMELA (Satellite)
- Neutrino physics
  - Borexino, Cuore, Cupid, Gerda, Icarus, Opera (Gran Sasso Lab.), Enubet (accelerator), KM3NeT (underwater), Juno, Tristan
- Dark Matter search
  - DAMPE (Satellite), DarkSide, NEWS, Xenon (Gran Sasso Lab.)
- Gravitational waves physics
  - Virgo (EGO, Cascina)
- Gamma Ray Observatory
  - CTA, LHAASO
- Support offered also to other collaborations
  - CHNET, Theophys, ...



Total jobs flow over time



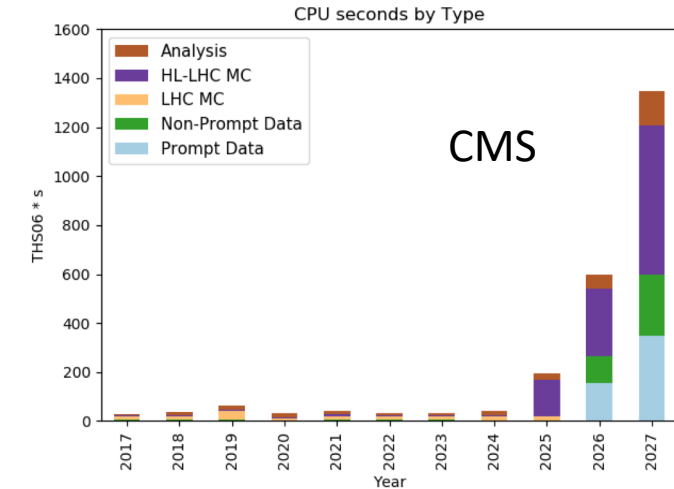
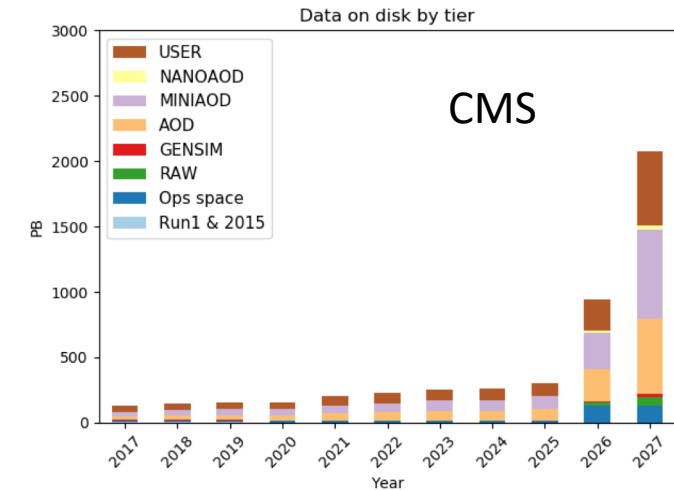
Farm usage: March 2018 – September 2019

# Towards HL-LHC: resource growth

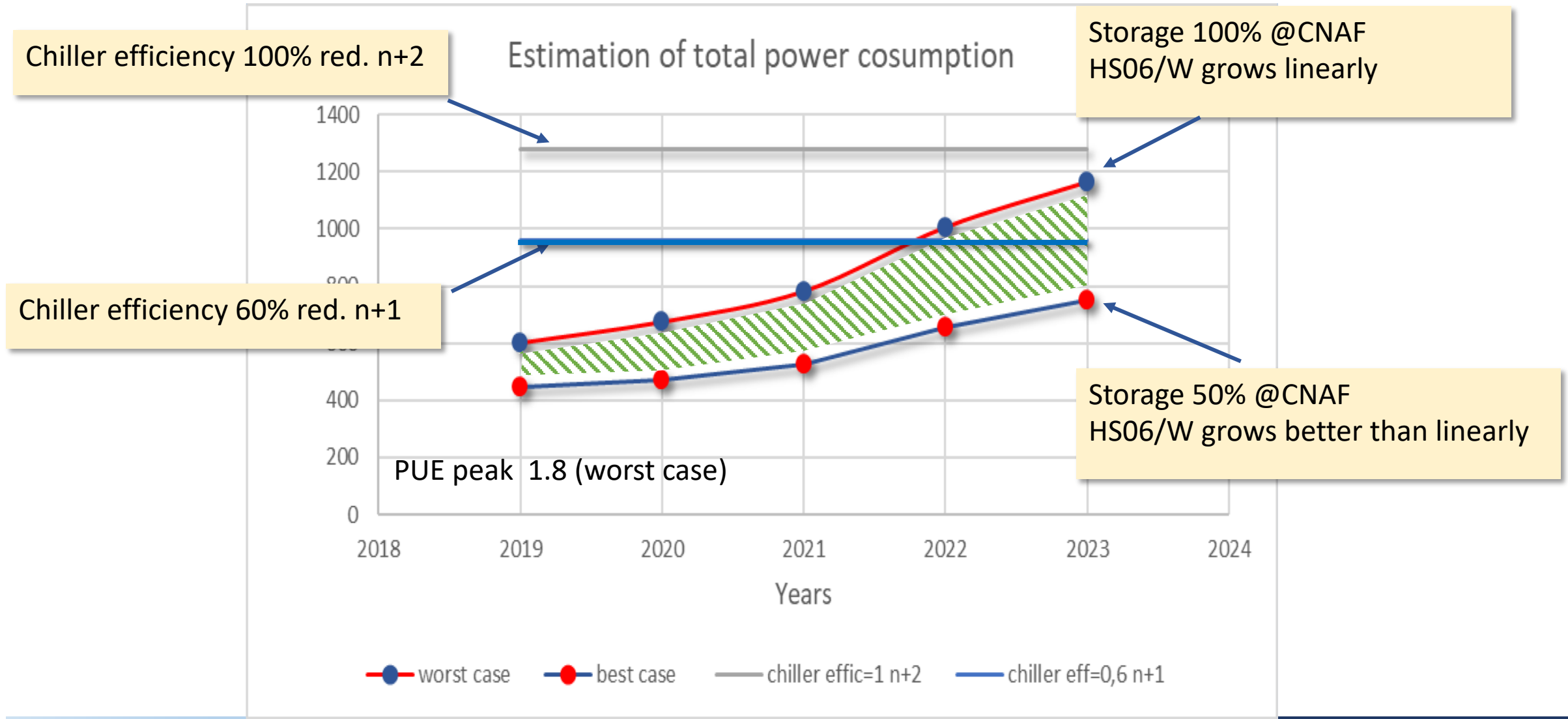
- Resource requests dominated (so far) by LHC experiments
  - But increasing requests from astro-particle experiments
- The number of CPU Cores and quantity of Storage required by LHC in RUN 4 (2026-) will be so high that won't fit in current data center.
  - At present an increase of  $\sim 10x$  is foreseen respect to 2019
- Even in RUN3 (2021-2023) LHC pledges will saturate the capacity of the current cooling system (6x320Kw chillers) of INFN Tier-1

In the short term a viable solution is to continue the strategy of having part of the farm remote.

In the medium-long term a new location for the Tier-1 data center is needed



# Data center evolution



... from the ashes of the tobacco industry the foresight of the ER Region is now creating a science park



These buildings are protected by the law as an architectural heritage



# Bologna Tecnopolo

- ECMWF's new data center will be hosted at Tecnopolo (from 2020)
- Also INFN and CINECA will locate their new data halls here
- CINECA, INFN and SISSA have won a call for a pre-exascale machine (Leonardo: 150—200 Pflops) funded by JRU EuroHPC to be hosted at Tecnopolo



**Bologna Tecnopolo (100000 sq.m.) will have one of the biggest data centers in Europe**  
These buildings are protected by the law as an architectural heritage



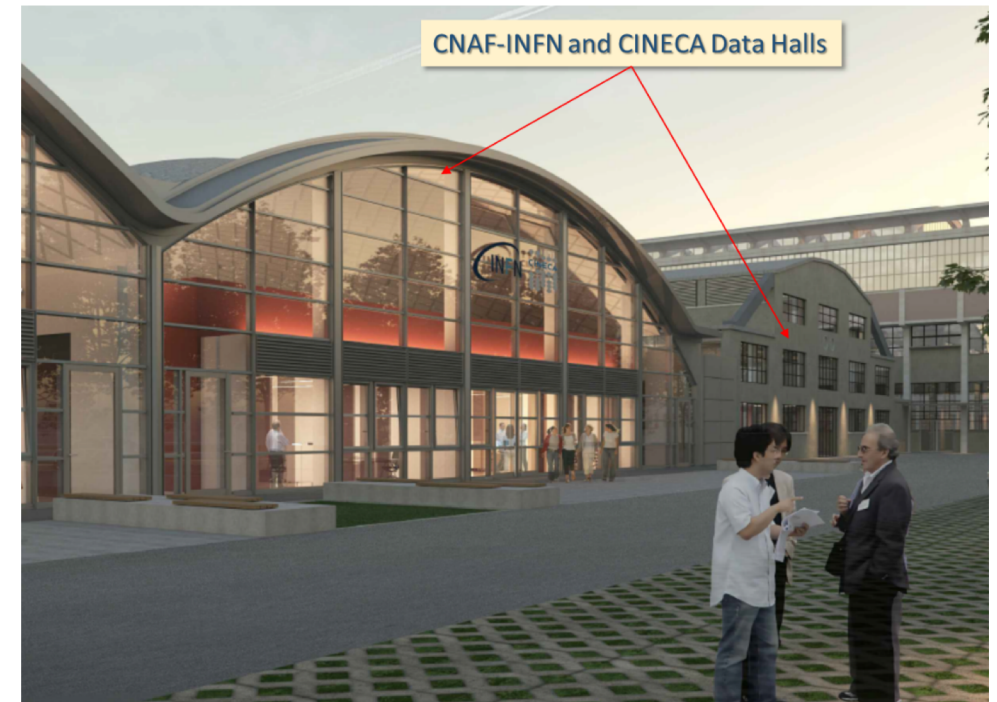
*Hall of the Pontifical Audiences (Nervi - 1963)*



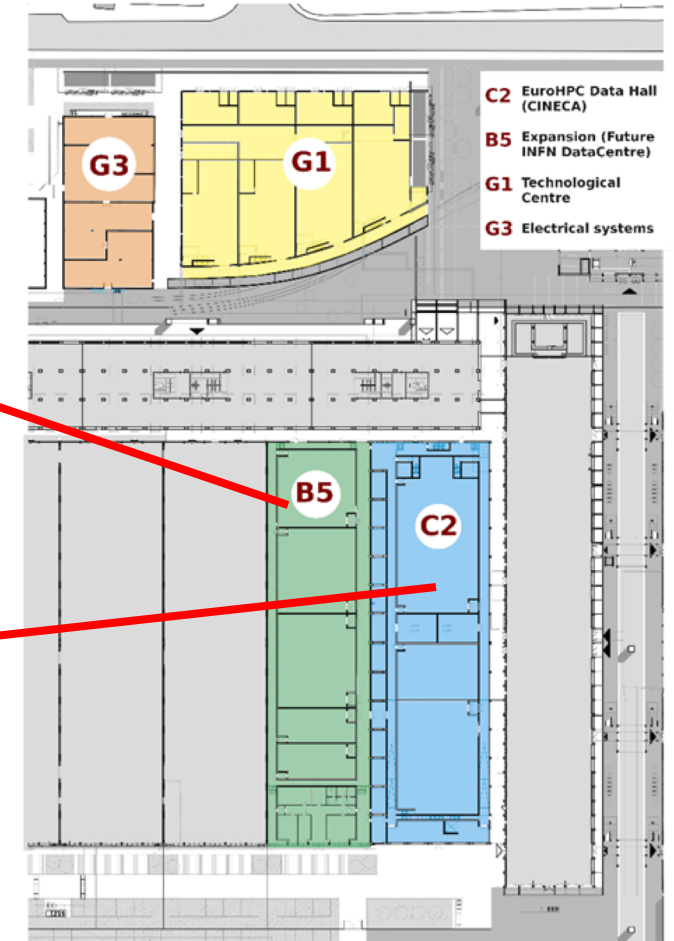
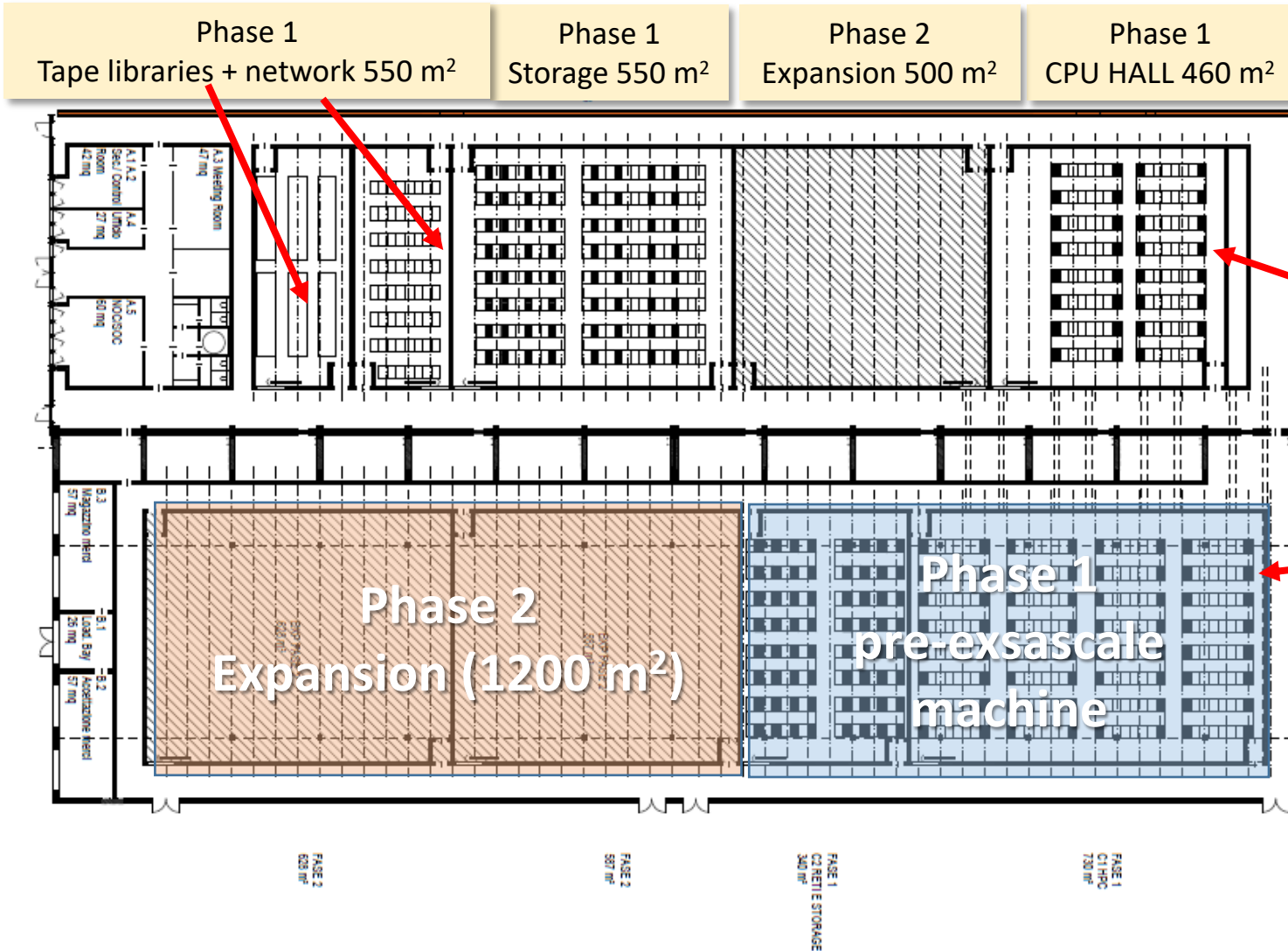
*Bologna Manifattura Tabacchi (Nervi - 1952)*

# The CNAF data center at the Bologna Tecnopolo

- CNAF datacenter will move to Bologna Tecnopolo during 2021, worst case in 2022.
- The project of the new data center is shared with CINECA.
- It is designed to be extremely energy efficient (PUE ~1.1).
- Technological plants and infrastructures are sized for 20 MW (10 MW for INFN)
- The implementation will be in two phases:
  - Phase 1: 2021-2025 up to 10-13 MW
  - Phase 2: after 2025 up to 20 MW

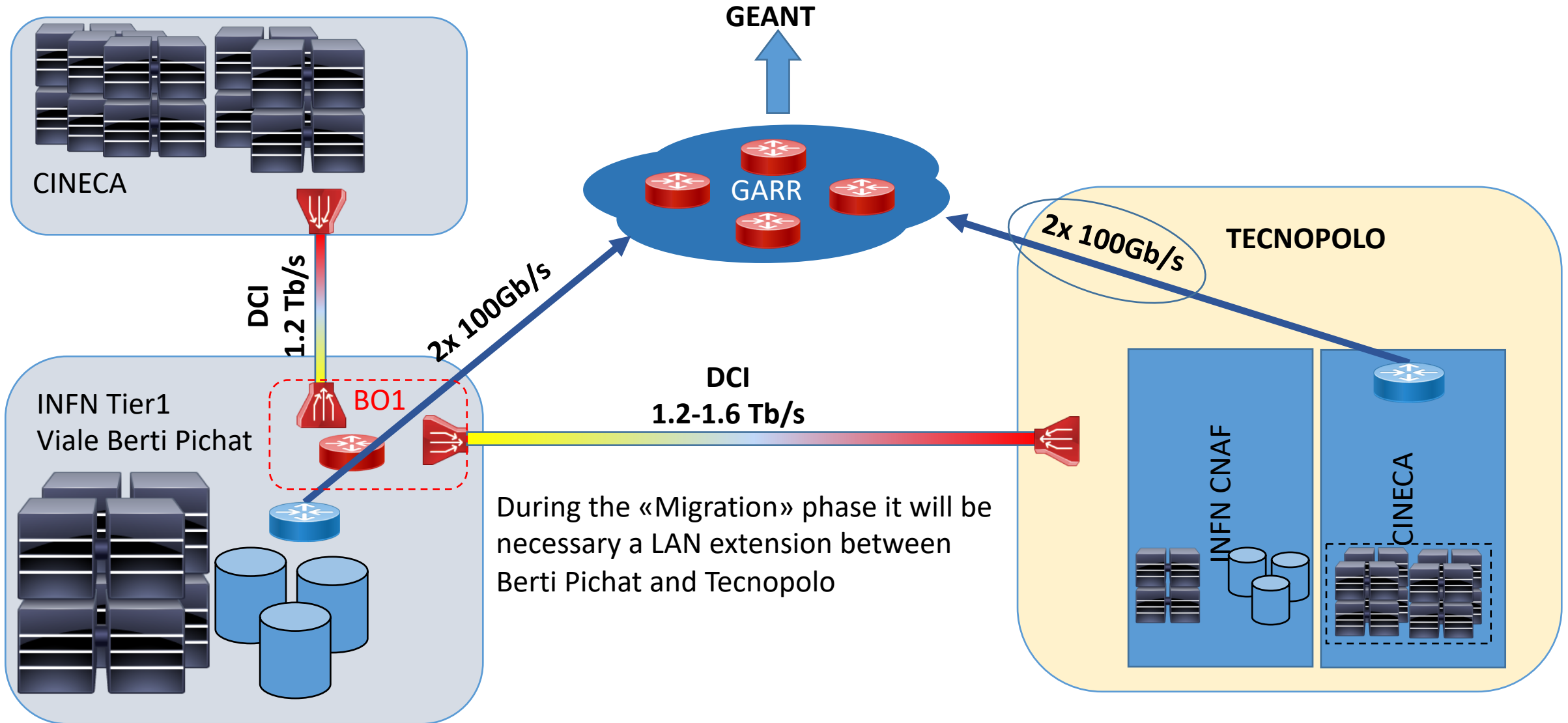


# Preliminary Layout data halls

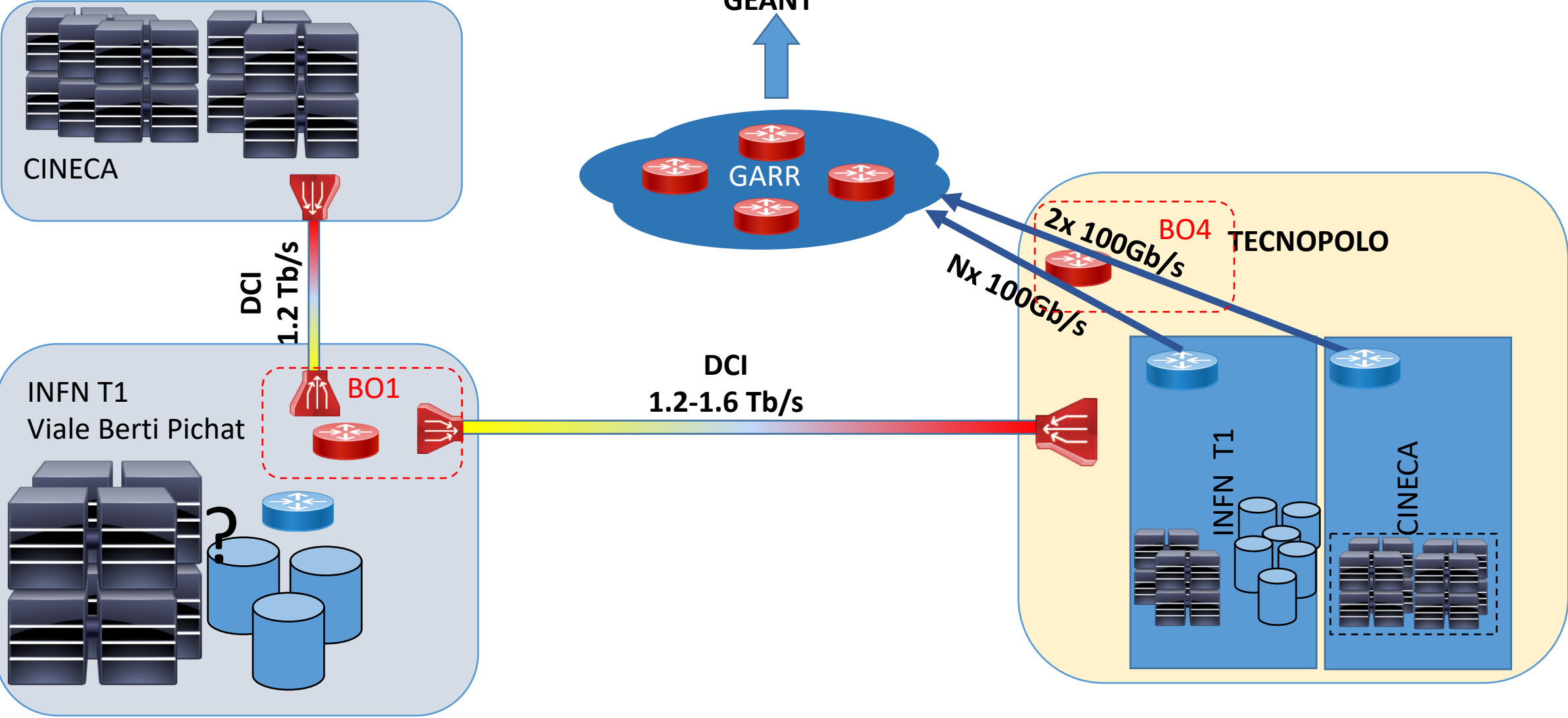


- Work in progress for a detailed strategy
  - Move will happen during LHC RUN 3
  - In a transitional phase the Tier-1 resources will be hosted part at CNAF data center and at Tecnopolo
- The transition phase will be based on the following approach
  - LAN extension (DCI) between CNAF and Tecnopolo (as currently with CINECA)
  - GPN and LHCONE/OPN uplinks to be switched from CNAF to Tecnopolo
  - Tape libraries and data to be moved asap
  - Disk buffer at Tecnopolo to move data first and then storage systems
  - Delivery of tenders to be phased depending on the precise date of entry into production of new data center
  - CPU at CNAF could be phased out naturally

# Connectivity (migration phase)



# Connectivity



1. Project Management
    - Overall coordination, timing, documentation, tenders, etc.
  2. In-box facilities
    - Power distribution, cooling, rack placement, network, etc.
  3. Pre-exascale machine (Leonardo)
    - Architectural requirements for HTC and interface with Tier1 (including the use of storage), procurement, etc.
  4. Services
    - Definition of portfolio of cloud services, farming virtualization, storage architecture, ITC services (SI+SSNN), operations, ISO27001, security, etc.
- *On top of this, a coordination with the INFN computing infrastructure would be needed*
    - *Integration within data lakes (WLCG and non WLCG), tighter integration with INFN Tier2s, cloud@INFN*



- HEP experiments are evaluating the “Data Lake “ model based on few very big datacenters
  - INFN Tier-1 will be part of the WLCG “Data Lake“
- For non WLCG experiments, some sort of this is needed too.
  - Ensure replication of non-reproducible data at infrastructural level w/o user intervention to a second site (→ see IDDLS)
  - Unique access point to data
- This could (should) be complemented by an INFN cloud

*“Data lakes are an extension of storage consolidation, where geographically distributed storage centers, potentially deploying different storage technologies, **are operated and accessed as a single entity**”.*

# Backup

---

- CNAF is a quite solid National Center operating in the ICT area, provides computing resources to most of the INFN experiments, has a strong R&D program and an innovative approach to the technology transfer, but ....
- ... the flooding showed the strong weakness of the CNAF site and forced us, to give a future to the Center, to immediately think of an alternative solution.
- Almost at the same time it has been possible to make synergy with the CINECA with the aim to transfer both data centers to the tecnopolo
- This gives to INFN and CNAF two great opportunities:
  - Make the most efficient use of a large HPC machine that can provide a significant part of our computer power needs
  - Be a strong candidate as data center for the HL-LHC data lake era, the future we would like for CNAF

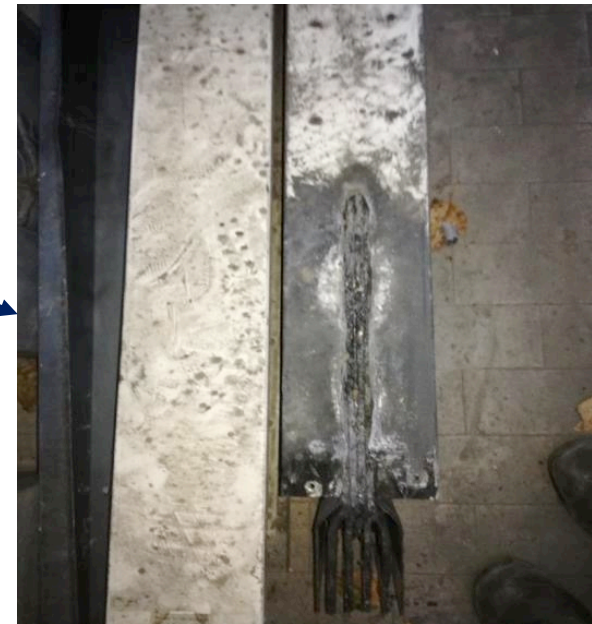
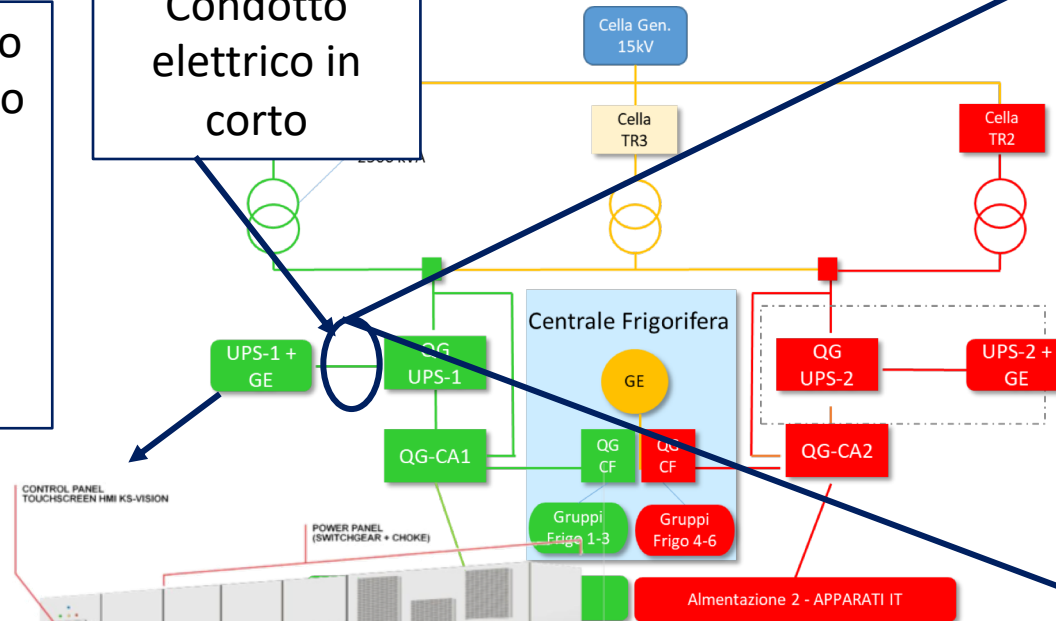
# Cause del guasto

---

1. numerosi indizi indicano che il fermo del sistema sia stato provocato da un corto circuito sul condotto sbarre che collega il gruppo rotante al quadro elettrico di controllo. Il corto circuito è avvenuto in prossimità di due giunti di connessione tra le sbarre ed è stato causato da un deterioramento della connessione con aumento della resistenza di contatto e conseguente surriscaldamento. Secondo questa interpretazione il guasto dell'UPS rotante (KS) è un effetto e non la causa dell'incidente.
2. l'alternatore del KS ha quindi subito un importante stress che ha causato la rottura della interconnessione meccanica con l'albero del volano. Tale interconnessione è stata ripristinata cambiando le componenti meccaniche rotte.
3. Il condotto elettrico (blindo sbarra) in questione (16 m) è stato completamente sostituito (cannibalizzando quello dell'UPS 2 che andrà ordinato quanto prima) nei giorni scorsi e il suo isolamento testato dal nostro personale. Vista la peculiarità di questo condotto è stato deciso che, prima di rimettere in produzione il KS, tale condotto deve essere certificato nuovamente dalla ditta che l'ha prodotto (Schneider). Oggi alle 17 interverrà il tecnico di Schneider per la certificazione. Se non verranno riscontrati problemi si tenterà nuovamente di far ripartire il KS, prima a vuoto e poi finalmente in produzione. Anche nelle migliori delle ipotesi servirà almeno tutta giovedì.
4. Nel frattempo (pochi minuti fa) l'allacciamento dell'UPS statico da 400 kW per fornire continuità ai nostri carichi più critici è stato completato (collaudo compreso). I carichi critici sulla linea rossa sono ora protetti da UPS.  
**Sono in atto le procedure per far ripartire il centro.**

In conseguenza del corto l'alternatore si è bloccato mentre il volano continuava a girare. L'interconnessione meccanica tra i due è saltata.

Condotto elettrico in corto



Volano

Alternatore

- Progressiva riapertura dei servizi
  - Sistema disco
  - Interfacce sistemi disco (SRM, gridfp, Xrootd, ecc. )
  - Tape library
  - Farm + HPC Farm
- A disposizione degli esperimenti circa 4000 slot (40000 HS06) di calcolo oltre al pledge complessivo
- Overpledge assegnato a esperimenti più in difficoltà
  - Richieste da soli due esperimenti al momento

# Altre azioni per mitigare gli effetti del down

---

- Assistenza per trasferire i dati o parte di essi in RECAS-Bari per avere ulteriore pledge di calcolo e disco (anche qui fino ad esaurimento risorse)
- Assistenza per recuperare eventuali file con problemi e/o recupero da tape

# Strategia e medio termine

---

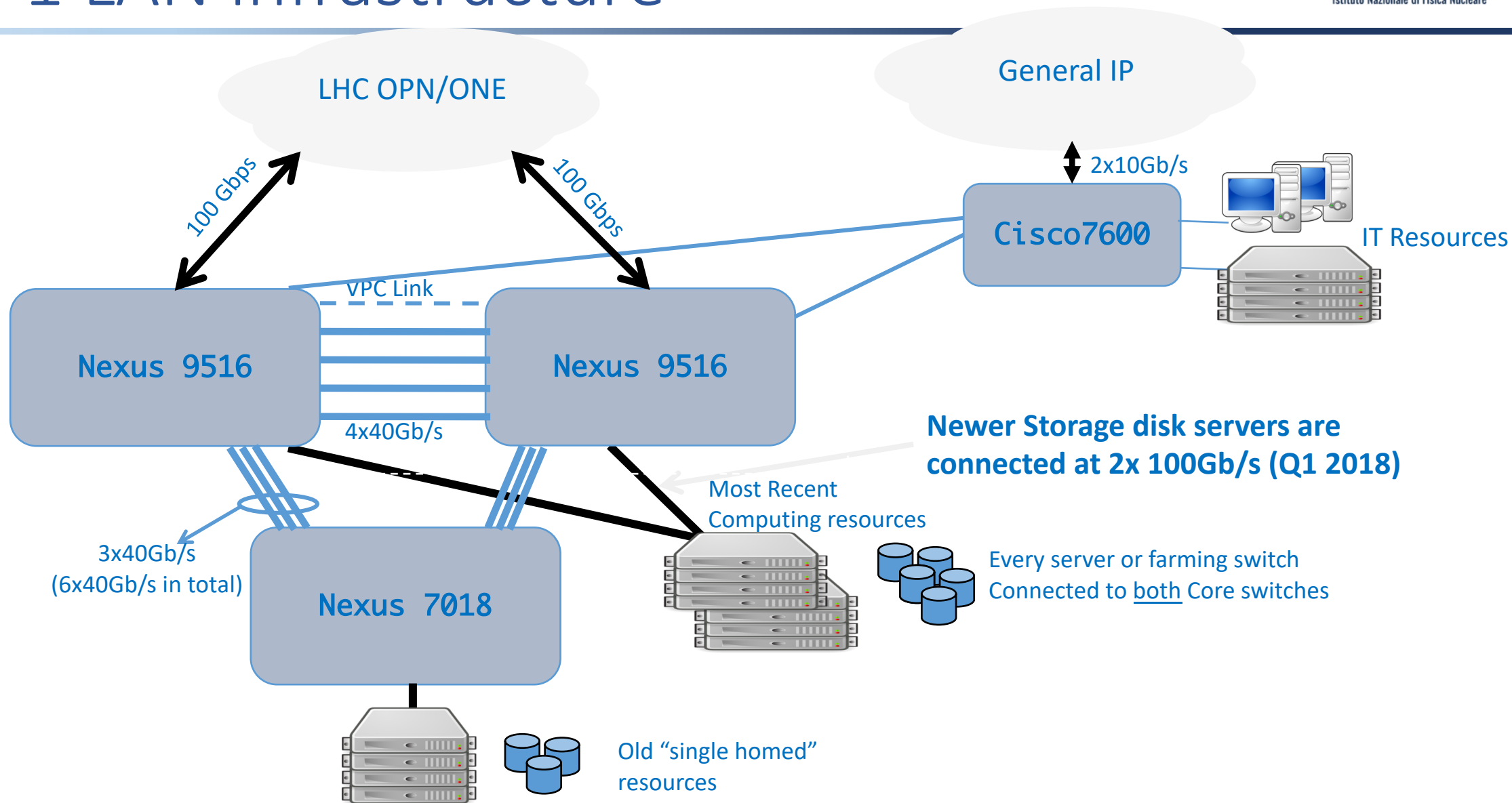
- Continueremo con la strategia di continuità mostrata (1 UPS rotante + 1 statico) fino a manutenzione straordinaria effettuata. Probabilmente fino al trasferimento al Tecnopolo (2021)
- Le infrastrutture elettriche hanno circa 10 anni, sono in contratto di manutenzione, ma guasti agli apparati come quello successo potrebbero aumentare di frequenza in futuro.
- Dobbiamo comunque avere una strategia che ci garantisca di poter continuare a lavorare anche in presenza di guasti di questo tipo
- Va implementata una replica automatica (comunque non a carico dell'utente) dei dati più preziosi su un altro sito
- L'idea è quella di anticipare il modello data lake che implementeremo per HL-LHC, magari in forma semplificata, tra il CNAF e un sito del PON IBISCO. L'idea era già stata presentata circa un anno fa e ora a PON approvato cercheremo di implementarla.



Back-up slides

---

# Tier-1 LAN infrastructure

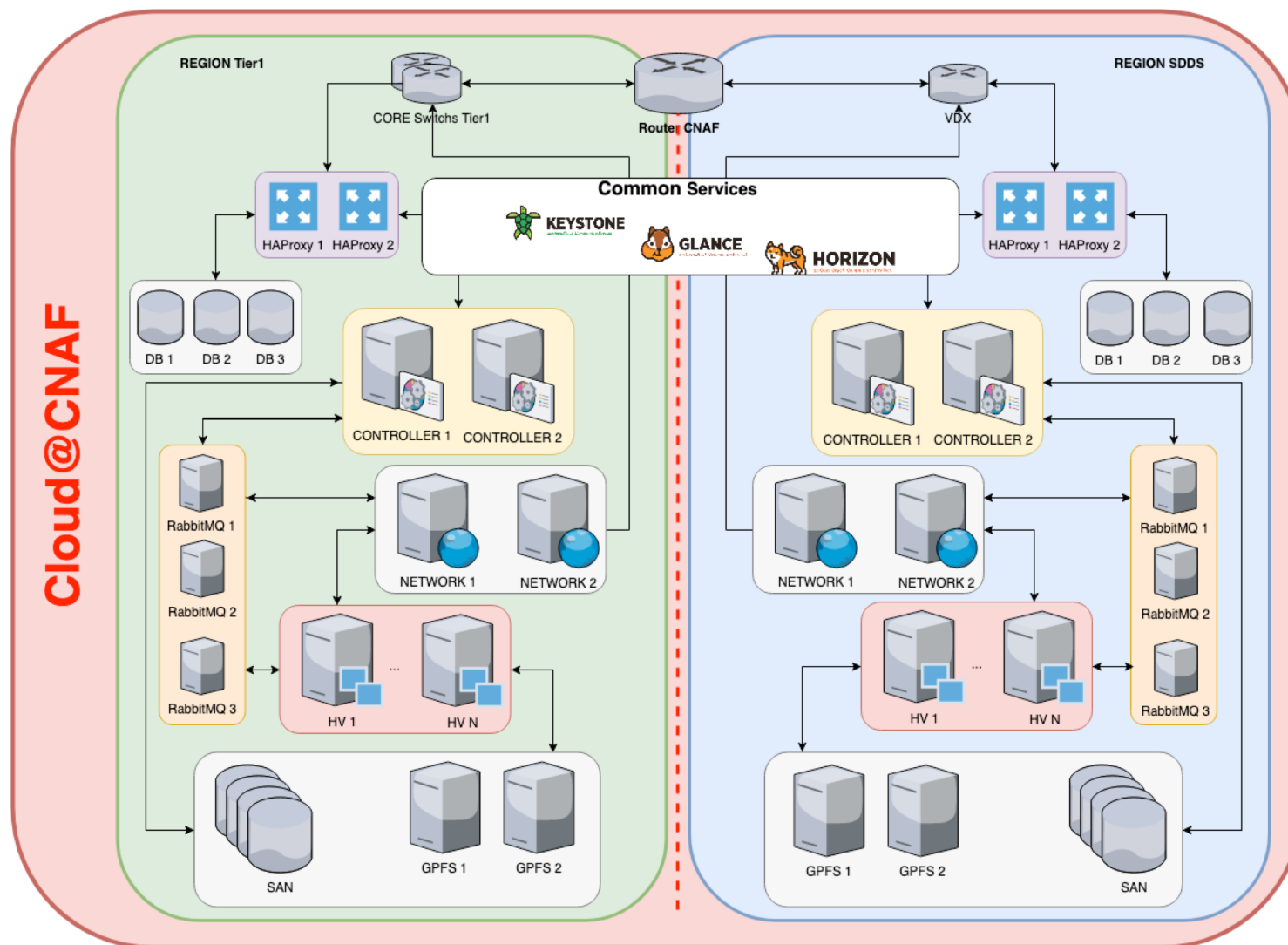


# Cloud@CNAF (1)

---

- First installation dates back to 2014
- Up to now addressing the needs of small groups w/o pledged resources
- Since the beginning of 2019 started the migration of Cloud@CNAF to a more reliable infrastructure based on the latest OpenStack release
  - Co-managed by SDDS and farm groups
- Aiming to extend the cloud-user community also to major experiments (including WLCG ones) with direct access to the Tier-1 resources and implementing a dynamic resource scaling between Grid and Cloud infrastructures
  - Mechanism for the dynamic partitioning of the Tier-1 farm is already in place

# Cloud@CNAF (2)



Cloud@CNAF

# Plan for tests on CINECA-HPC

- Test for LHC on the CINECA Marconi A2 partition (KNL based)
  - 64x4 cores per node, 96 GB
  - Currently no virtualization, no external network access
- Production phase (near future):
  - Use parasitic + grant based allocation
  - Install CVMFS + Singularity
  - Submit from CNAF via CREAM tunneling to CINECA's Slurm
  - Tunnel all external access via the Infinera DCI link

*Kick-off meeting June 14<sup>th</sup>*



## Marconi- A2

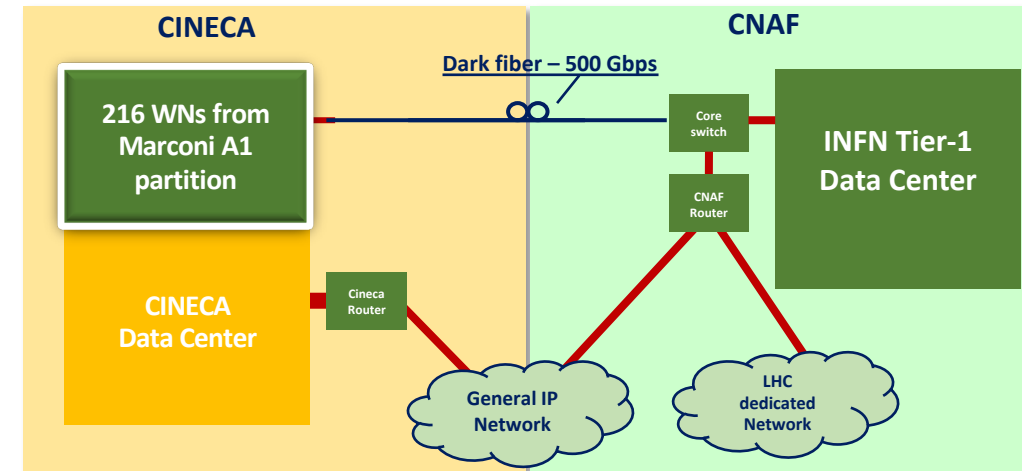
Model: Lenovo Adam Pass  
Architecture: Intel OmniPath Cluster

Nodes: 3.600  
Processors: 1 x 68-cores Intel Xeon Phi 7250 CPU (Knights Landing), 1.40 GHz  
Cores: 68 cores/node (272 with HyperThreading), 244.800 cores in total  
RAM: 16 GB/node MCDRAM + 96 GB/node DDR4  
Internal Network: Intel OmniPath Architecture 2:1

Peak Performance: 11 PFlop/s

# Farm extension to CINECA

- ~180 kHS06 leased from CINECA (from Marconi A1 partition)
  - 216 WNs (72 cores HT each), 3.5 GB/core
  - 10 Gbit connection to rack switch
  - 4x40 from rack switch to router aggregator
- Dedicated fiber couple directly connecting Tier-1 core switches to our aggregation router at CINECA (upgradable to 1.2 Tbps) via Infinera DCI
- Quasi-LAN situation (RTT: 0.48 ms vs. 0.28 ms on T1 LAN)
- No disk cache, direct access to CNAF storage
- Efficiency comparable to partition @CNAF



12 x 100 Gb Ethernet QSFP28

# INFN Tier-1 extension off premises

---

- Farm partitioned in 3 locations
  - ~210 kHS06 at CNAF data center
  - ~180 kHS06 at CINECA data center (in production since March 2018)
  - ~10 kHS06 at Bari-ReCaS data center (in production since beginning 2017)
- Able also to extend on commercial clouds
  - Proved by small scale tests on Aruba and Azure and scalability tests with hybrid cloud in the context of EU project HNSciCloud
- Requirements
  - Access to remote WNs must be transparent for users
    - i.e. INFN Tier-1 CEs (or LBS) as unique access points also to these resources
  - Key issue is data access (i.e. remote or cache)

# Tape library

---

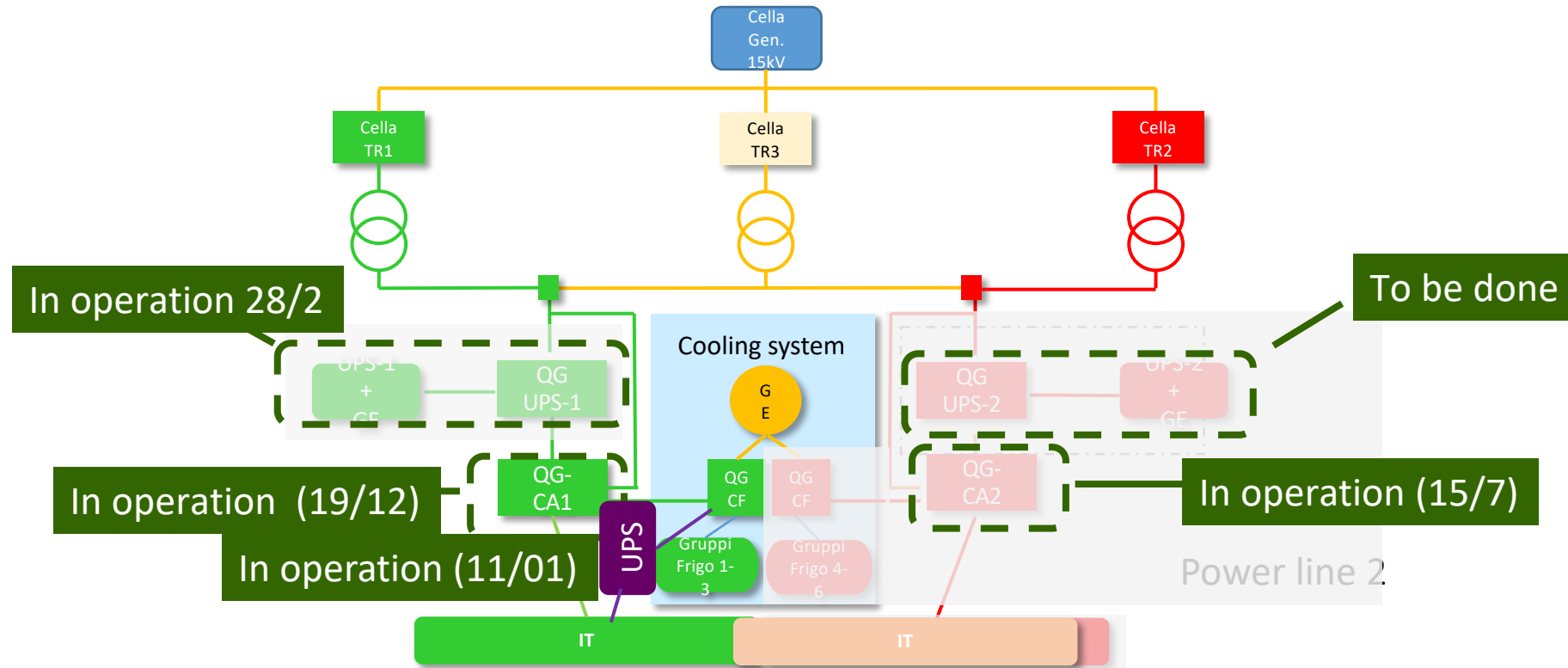
- A tape library SL8500 (10000 slots) is installed
  - 16 T10kD drives (all interconnect via 16 Gbps FC to the TAN)
  - Almost fully populated with tape cartridges (8.4 TB each)
    - Last bunch of tapes (1600) to be delivered this month
- Because of the flooding, several tapes were damaged
  - Recovered in a specialized laboratory or data retransferred
- Complete repack performed (2014-2016)
  - Performing statistical checks on tapes after the flooding
- Tape drives up to now statically allocated to experiments
  - Dynamic allocation of drives in pre-production mode
- A second library should be delivered in a month
  - Based IBM technology (18 TB tapes, ~5000 slots)



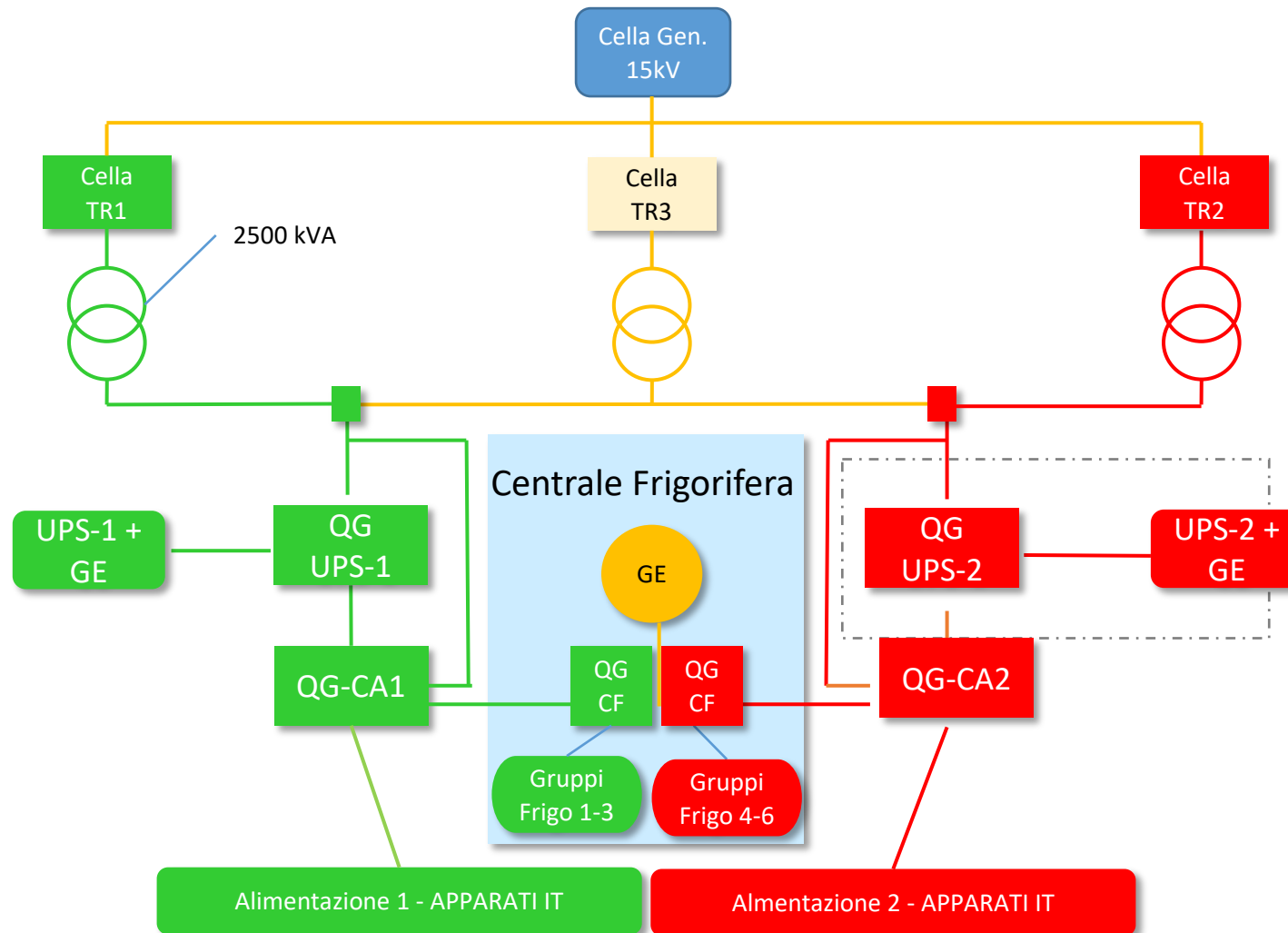
# La strategia di continuità elettrica dopo l'allagamento

- Dopo l'allagamento del novembre 2017, insieme con il CIAC (CNAF Infrastructure Advisory Committee) è stata decisa una nuova strategia per la continuità elettrica del Tier1 che teneva conto anche della necessità di mettere in manutenzione straordinaria i due sistemi di UPS rotanti (UPS + Generatore Diesel). La manutenzione straordinaria comporta di spedire in Belgio per alcuni mesi ciascun sistema. La strategia prevedeva:
  - di far ripartire immediatamente UPS-1 verde
  - Di mettere UPS-2 rosso in manutenzione straordinaria e spedirlo in Belgio per alcuni mesi. Stessa sorte sarebbe poi successe a UPS-1 quando il 2 sarebbe tornato
  - Di acquisire un UPS statico da 400 kW per proteggere i carichi critici (in primis storage) in assenza di UPS-2 (e poi UPS-1)
- UPS-1 verde è ripartito nei tempi previsti
- Per ordinare la manutenzione straordinaria di UPS-2 e il rifacimento del relativo quadro elettrico finito a bagno abbiamo impegnato più di 1 anno a causa di problemi di «burocrazia» e contrattuali. Il contratto è stato firmato il 1 agosto 2019 !!!
- Anche l'UPS 400 kW è stato consegnato in ritardo e la messa in servizio era previsto per i primi giorni di settembre
- Il 6 agosto eravamo quindi con la continuità solo sulla linea verde e il guasto ci ha lasciato senza continuità su tutto il centro

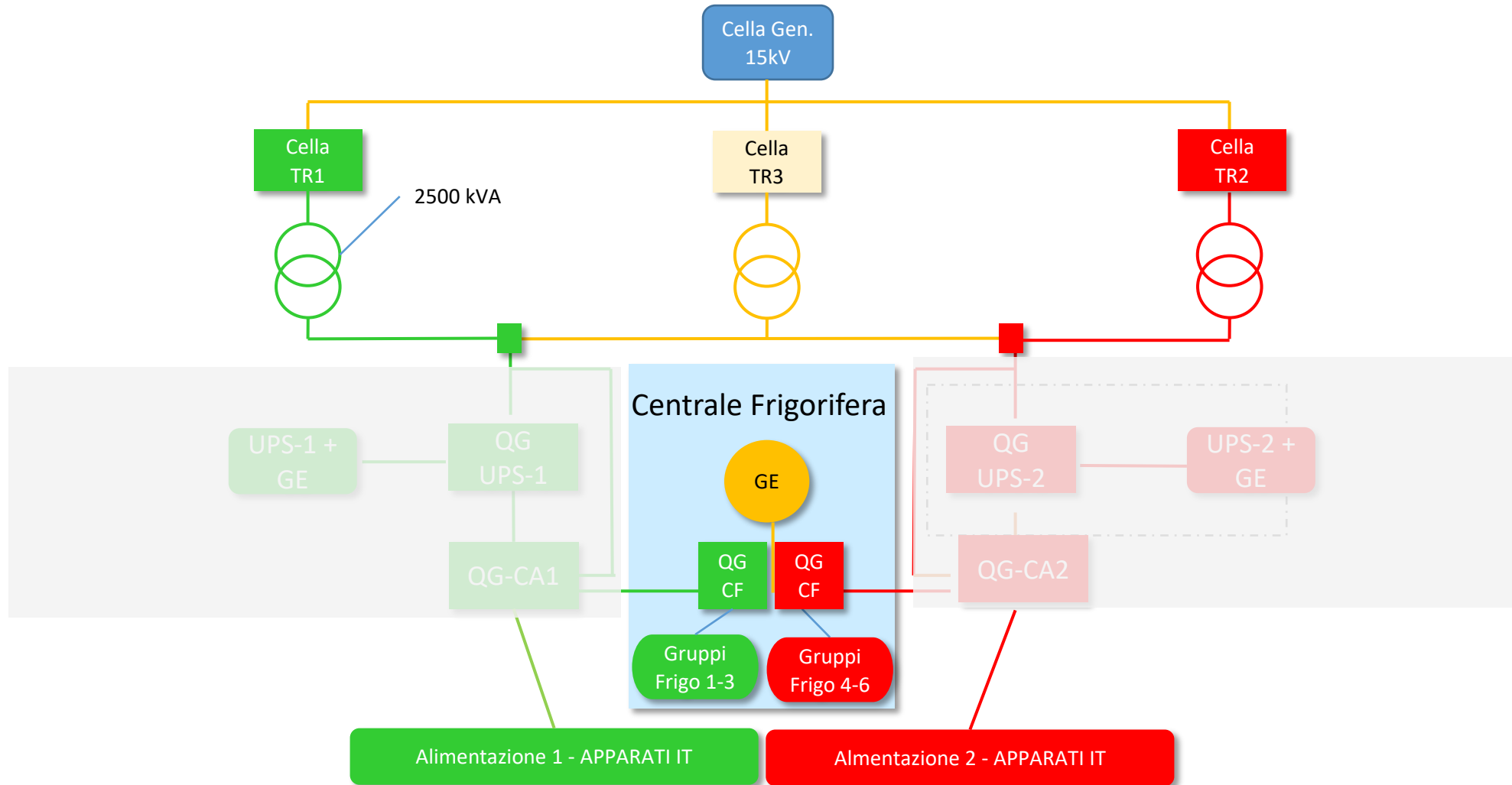
# Present Power Center status



# Infrastruttura elettrica: come era prima dell'allagamento del nov 2017

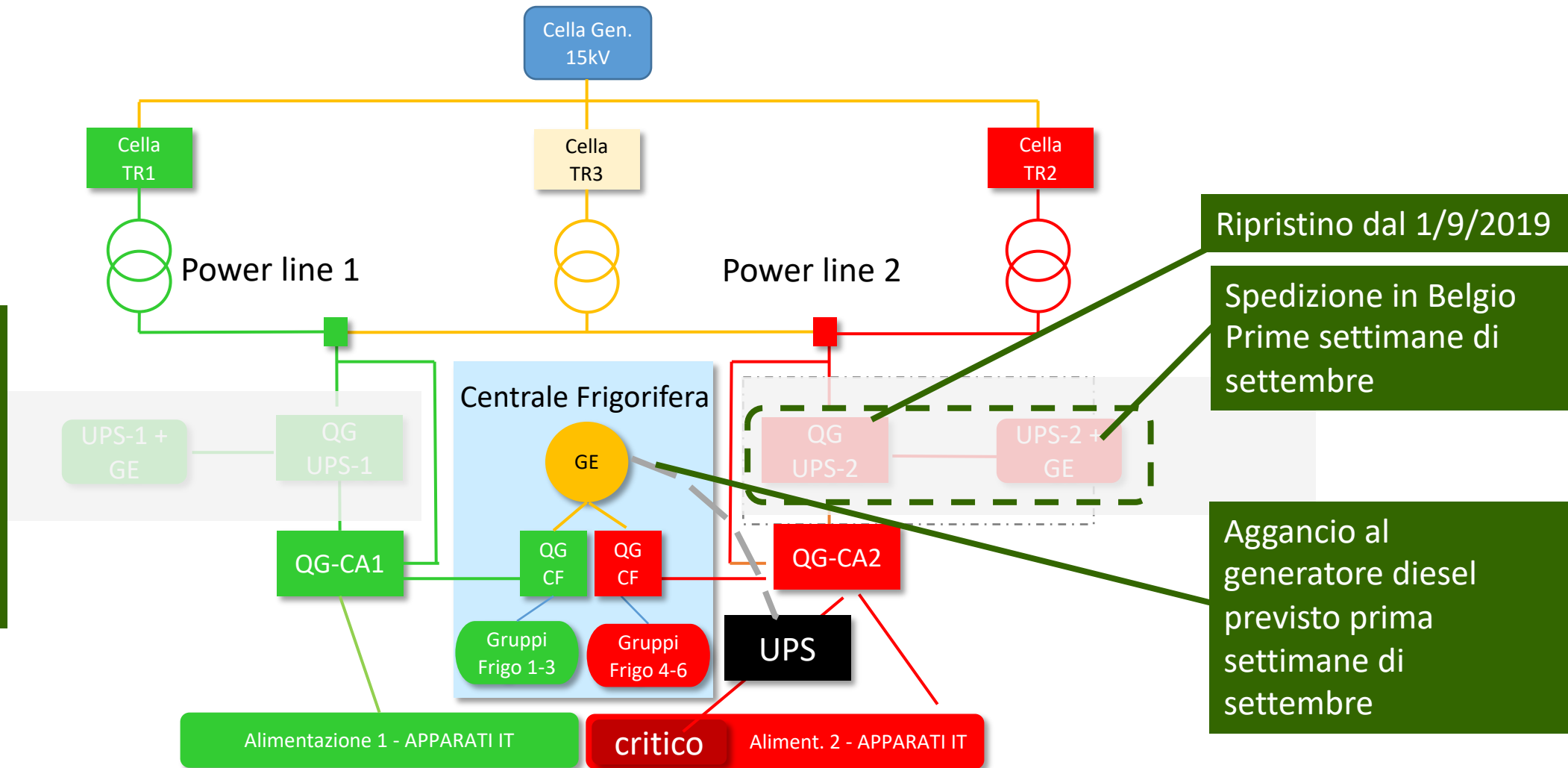


# Effetti dell'allagamento sulla parte elettrica



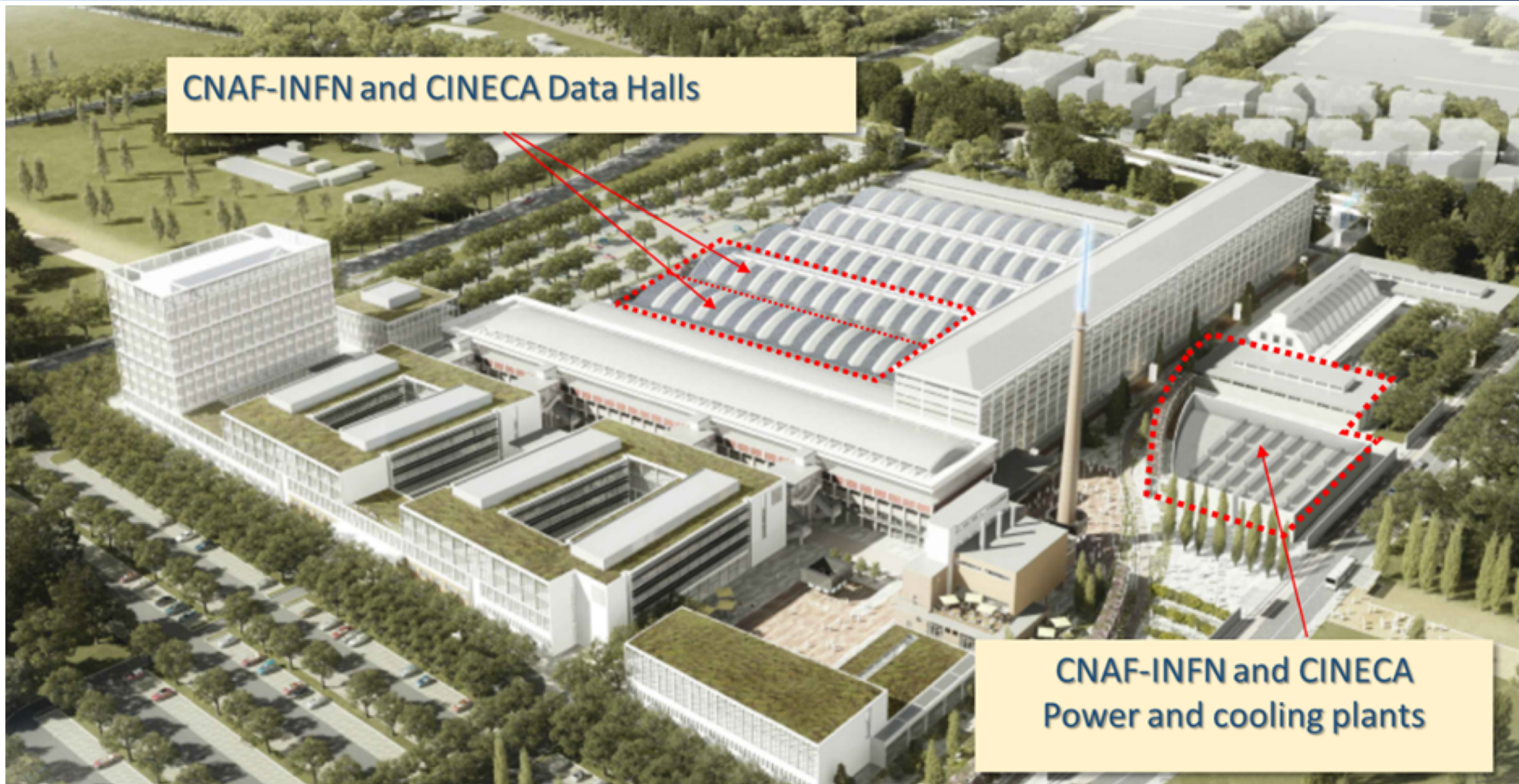
# Situazione attuale:

- Certificazione condotto elettrico oggi dalle ore 17
- Avviamento UPS - 1 a seguire se il condotto è ok



# History pills of the site of Manifattura Tabacchi

<b>1949-1963</b>	<b>proposal, design, project implementation</b>
1963-2004	production
2004-2008	Site closed, no longer productive
2009	The Emilia Romagna (ER) Region acquires the site
2011	ER Region launches an international tender to transform the site in a Tecnopolo/Science Park) to gather the research institutions, university labs, etc. of the Bologna area. Won by the German Company GMP.
2017	ECMWF (European Centre for Medium-Range Weather Forecasts) approves the proposal by Italian Government and ER Region to host ECMWF's new data center at Bologna Science Park
2017	Two main actors in HPC and Big Data handling, CINECA e INFN, join the Bologna Science Park
	Bologna Science Park becomes, with ECMWF, CINECA and INFN, one the bigger HPC and Big Data site in Europe

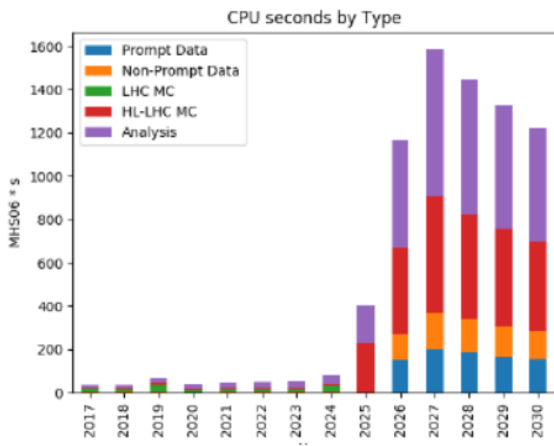
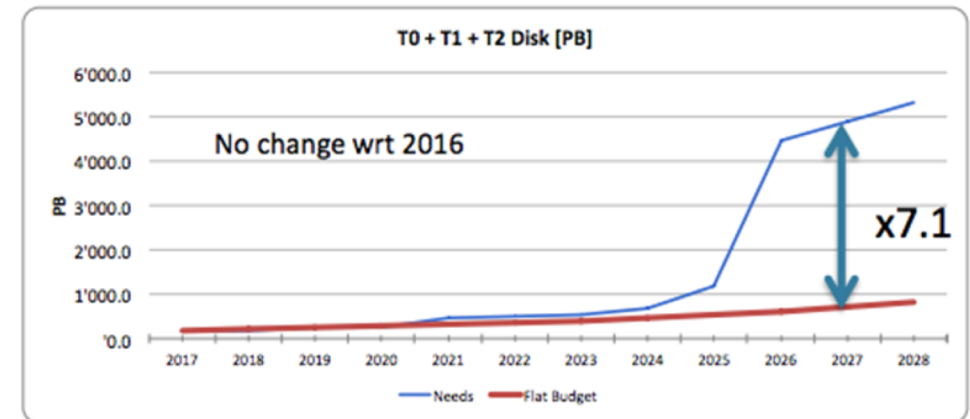
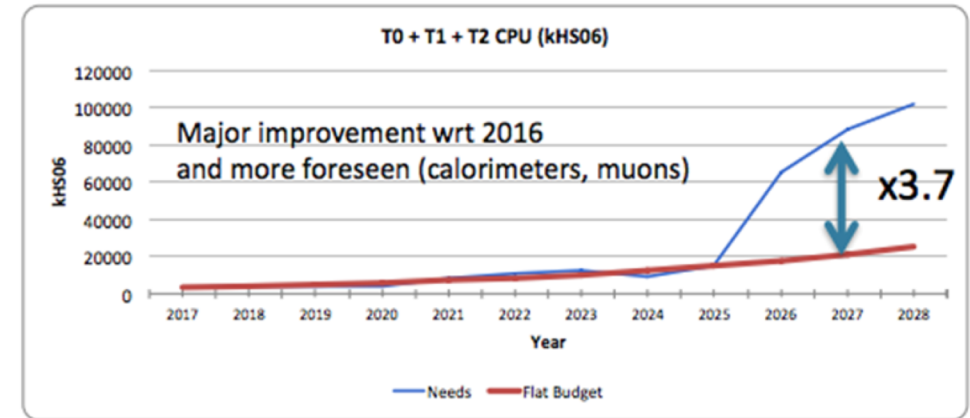
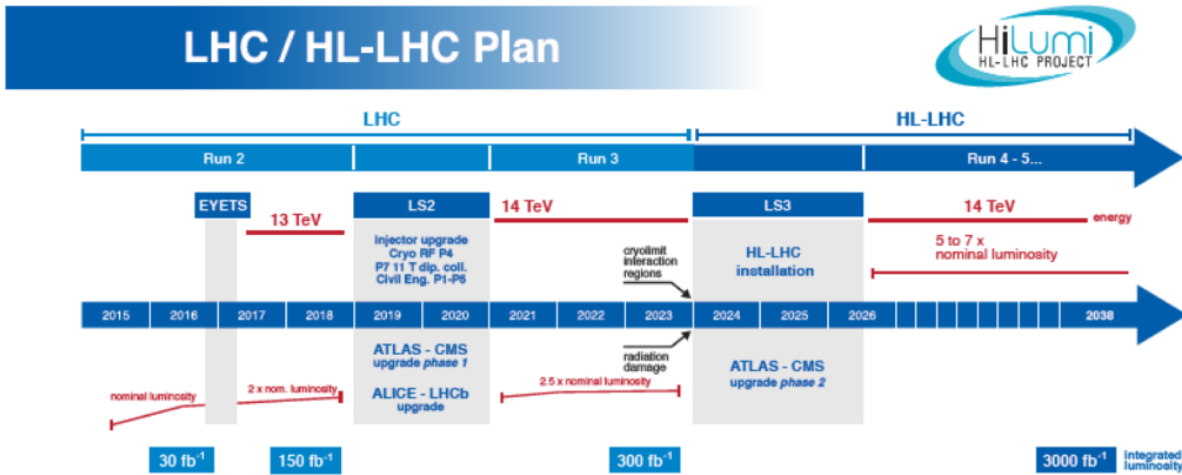


# Towards HL-LHC: data center infrastructure

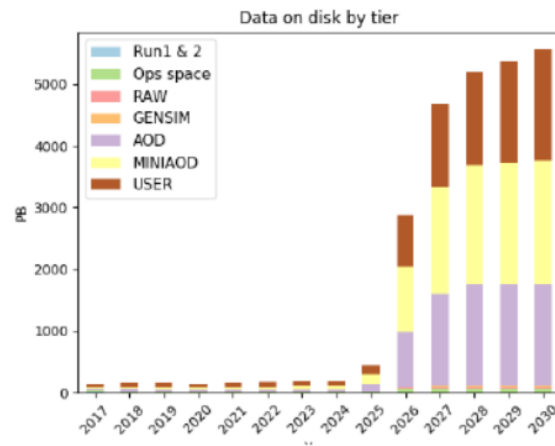
- In principle, the power center can cope with the load up to 2023 but this does not hold for the cooling infrastructure
  - The efficiency of chillers is reduced respect to the nominal one
  - During Summer redundancy could reduce to n+1 with 750 kW of installed IT
  - Not sustainable with 1 MW of installed IT (likely scenario in 2022-23)
- This is one of the reasons why part of the load (i.e. part of the farm) is remote
- Moreover we have had 2 serious incidents (flooding and KS break) in the last 2 years (since infrastructure renovation in 2008)
  - The first one hints to look for a better location
  - The second one could be a sign of obsolescence



# Long term evolution: HL-LHC



(a)



(b)