



Istituto Nazionale di Fisica Nucleare

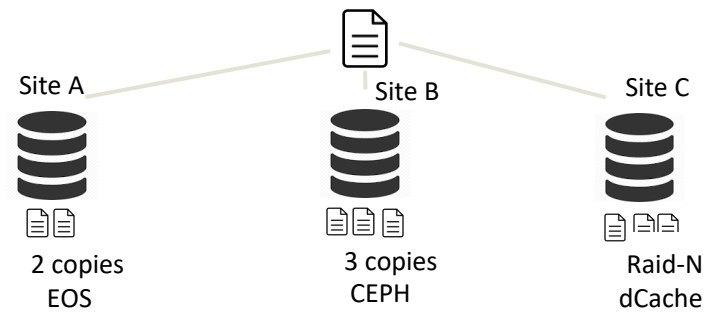


## IDDLS: Italian Distributed Data Lake for Science

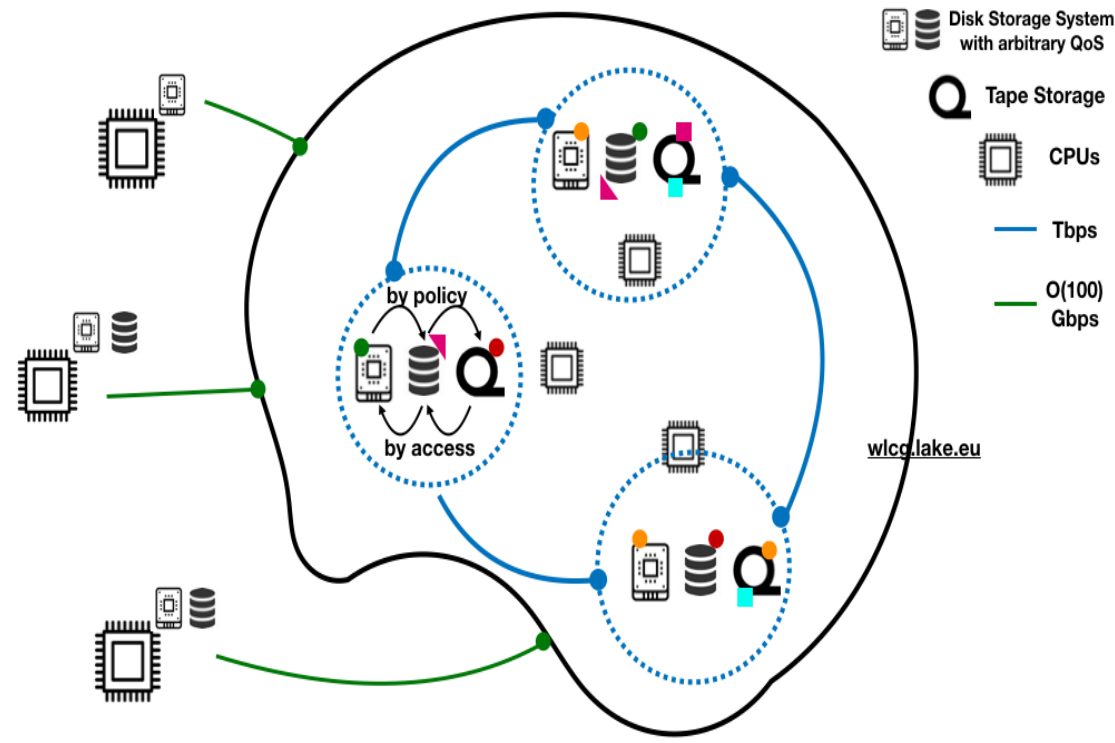
Stato progetto 13/07/2020

# + Some ideas on reducing Storage costs

- Reduce hardware cost: better exploiting the concept of QoS(Quality of Service)
  - Probably today we replicate more than we need
    - Reducing the number of copies
- Reduce Operational Cost: deploy fewer (larger) storage services maintaining high standards in availability and reliability
  - Create large storage repositories that “look like one, but it is composed of many” → **the DataLake**
- Co-location of Storage and CPU will not be guaranteed anymore
  - Need technologies for quasi-transparent data access from remote locations
    - Smart Caching



**A stronger integration of sites could lead to a reduction of the number of copies**



# + Sinergie in Europa

Comunità di utenti interessate all'utilizzo della tecnologia

- ESCAPE
- WLCG
- BELLE-II

Sviluppo SW per la creazione del DataLake

- XDC@INFN
- WLCG-Demonstrator@NA
- SCORES@NA
- Xcache@PG

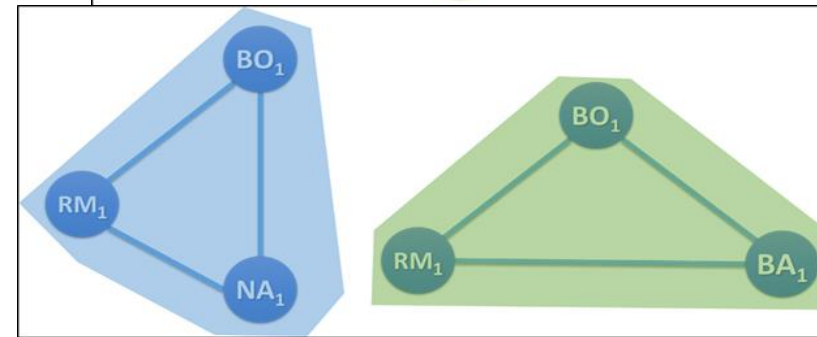
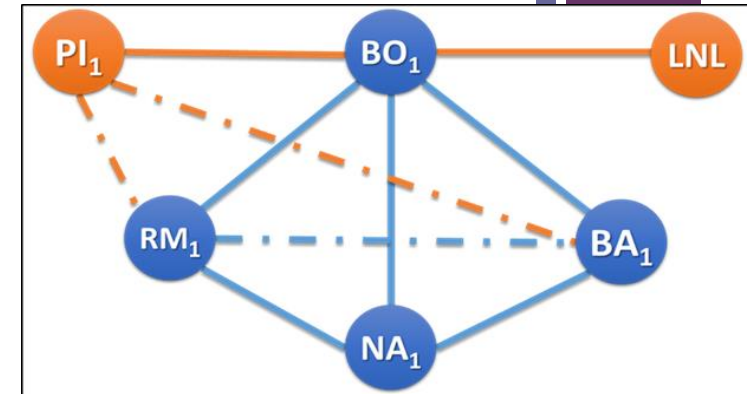
**WLCG-DOMA**

Creazione dell'infrastruttura HW

**IDDLS**

# + The project

- INFN-GARR collaboration to realize a prototype of an Italian DataLake exploiting:
  - Last generation networking technologies provided by GARR
    - DCI (Data Center Interconnection) equipment
    - SDN (Software Defined Network) deployment
  - Software for creating **scalable storage federations** provided by INFN
    - eXtreme-DataCloud project
    - SCoRES project (INFN-NA)
  - Real life use cases for testing
    - CMS
    - ATLAS
    - BELLE-II
    - Possibly involving LNGS experiments (XENON) and VIRGO



Possible topologies of the GARR Network with DCI and SDN for the DataLake

# + Timeline

## ■ 3 years project

### ■ First year

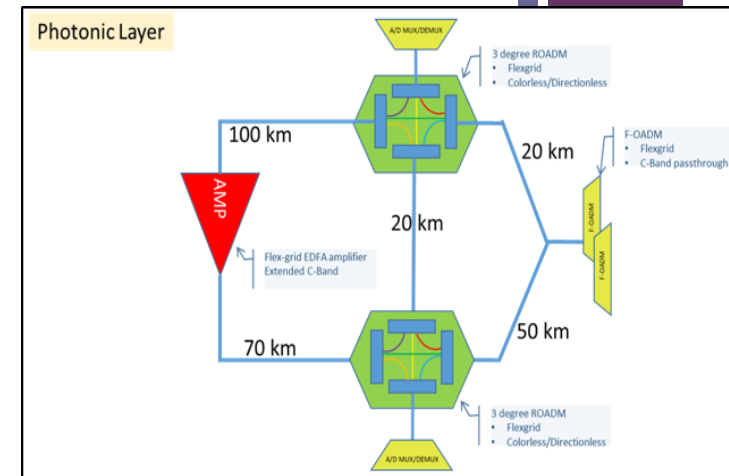
- Technology scouting for DCI equipment to be deployed by GARR
- Application (INFN) requirements analysis
- Network equipment acquisition (INFN and GARR) and Lab testing
- Deployment on mixed Lab+WAN environment of the networking equipment
- Creation of the DataLake on sites connected with standard networking and first prototype using DCI

### ■ Second year

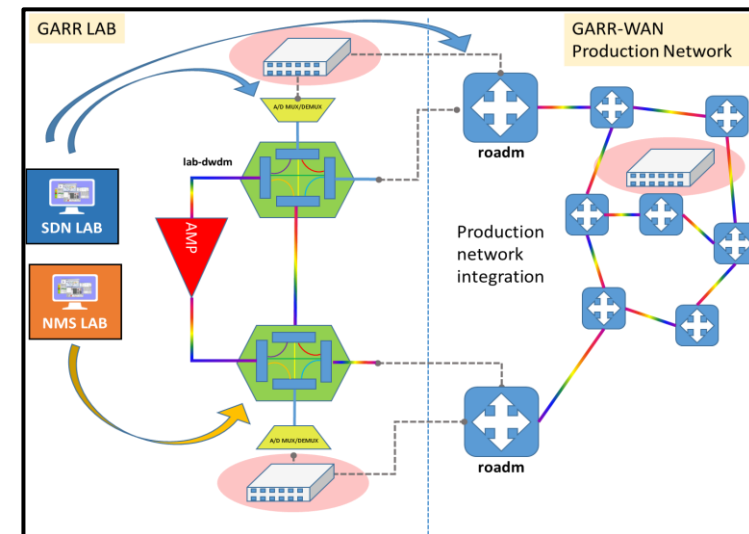
- Testing of the mixed (Lab+WAN) configuration
- Final creation of the DataLake on the 3 INFN sites with DCI systems
- Performance evaluation and comparison
- Possible acquisition of new equipment with increased performance

### ■ Third year

- Deployment only on WAN of the networking equipment
- Optimization of the DataLake
- Performance evaluation
- Final consideration



Lab deployment at GARR for testing





Mixed Lab+WAN deployment



# Milestone 2019 e stato



- **Gennaio 2019: Kickoff meeting @GARR**
  - <https://agenda.infn.it/event/17957/>
- **Gennaio – Giugno:** test apparati trasmissivi in Lab GARR (slide successive)
- **30/06/2019: Scelta degli apparati di networking per la creazione del DataLake**
  - Apparati GARR testati
    - Infinera Groove G30 
    - Transponder TP100GOTN XTM Infinera
    - Juniper ACX6360 (and MX240)
  - Switch con porte a 100Gb per connessione lato datacenter valutati:
    - Arista 7050SX3-48YC8, Arista716032CQ, DCS-7050CX3-32S-F 
    - Extreme Networks SLX930
    - NEXSUS 93180
    - CISCO Cisco Catalyst9500-32C-A
- **31/12/2019: Gara assegnata** Gara assegnata: Cisco Catalyst9500-32C-A
- **31/12/2019 Deployment degli apparati** di rete in una configurazione Lab+WAN
  - primi portotipi DataLake su apparati DCI e standard

Delayed, tests only in the Lab

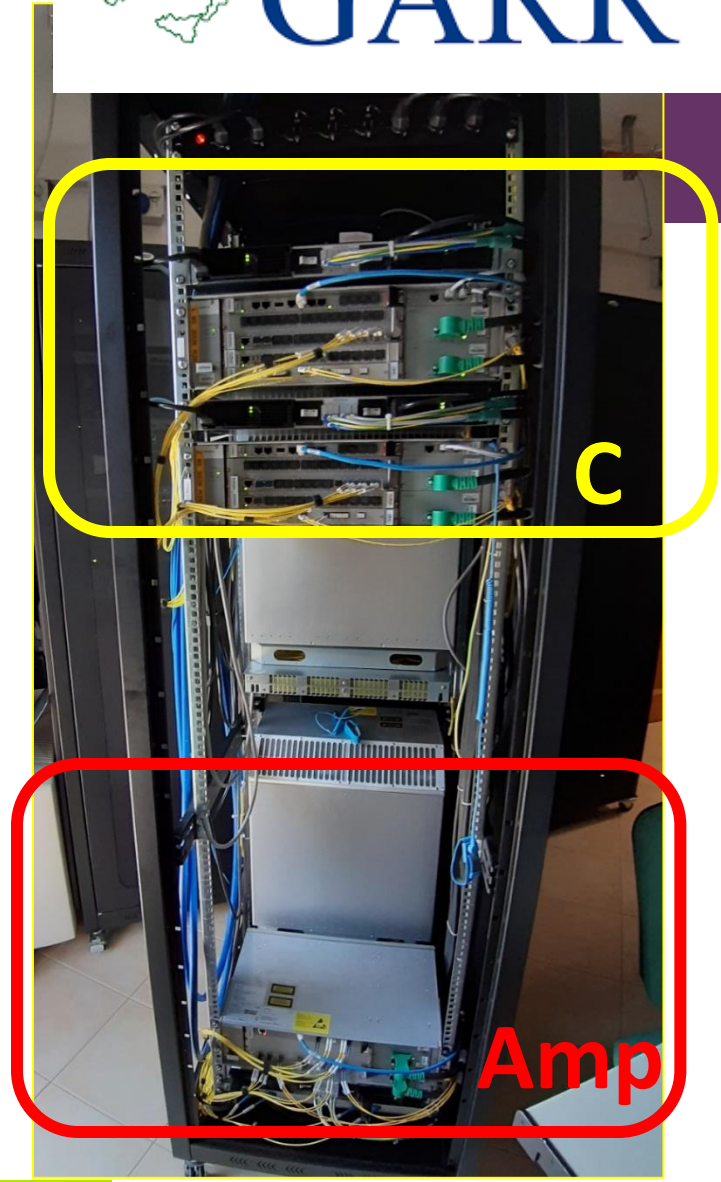
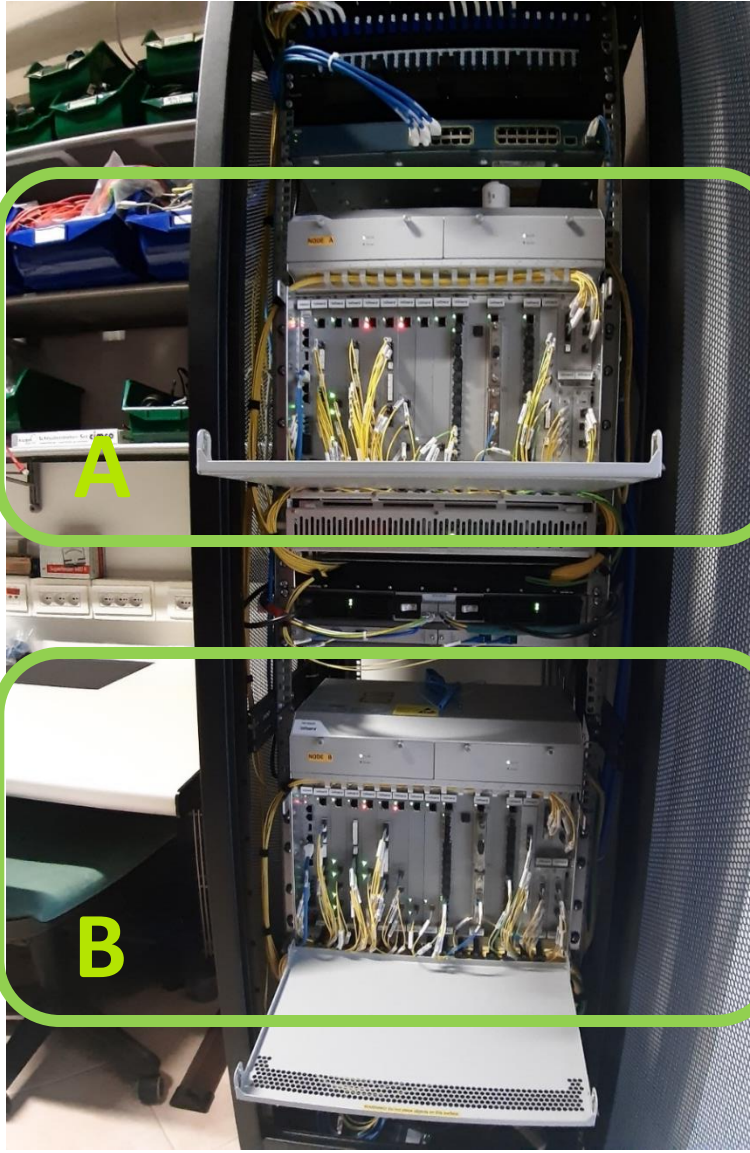
- **31/12/2019: Primi run per la valutazione delle performance sui prototipi**

Delayed, tests only in the Lab





# Slides from GARR...





# Fiber spools

- 16 spools of single G.652d fiber
- (vendor: Fujikura)
  - 8 spools 50km -> 4 pairs
  - 8 spools 25 km -> 4 pairs

In total 600km single fiber

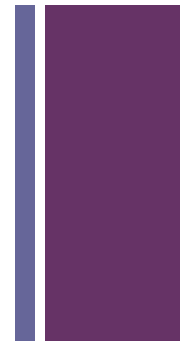
The **flexible topology** changes through a patch panel on the rack top



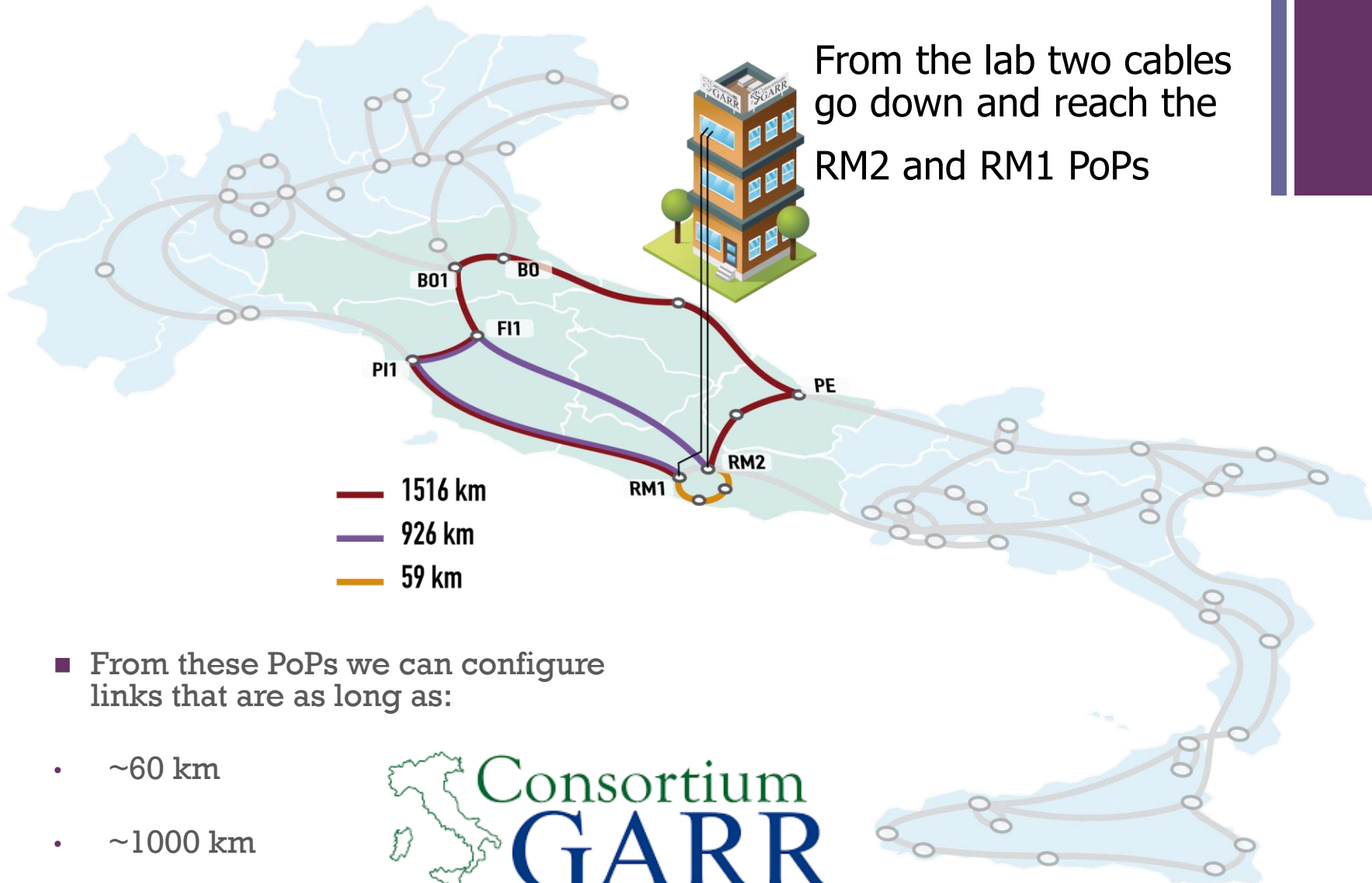




# Links towards the production environment



From the lab two cables go down and reach the RM2 and RM1 PoPs



■ From these PoPs we can configure links that are as long as:

- ~60 km
- ~1000 km
- ~1500 km



# + Test transponders 100G

1. Transponder TP100GOTN XTM Infinera  
Client: 100G Ethernet  
Line: CFP1 100G ACO - QPSK coherent



2. Juniper ACX6360 (and MX240)  
Client: 100-Gbps (QSFP28) pluggable interfaces  
Line: 200-Gbps CFP2-DCO coherent DWDM pluggable interfaces, which support 100-Gbps QPSK and 200-Gbps 8QAM, and 16QAM modulation options



3. Infinera Groove G30  
Client: up to 4x100GbE QSFP28  
Line: CHM1 sled – 400G 2xCFP2-ACO (100G QPSK, 150G 8QAM, 200G 16QAM)



# + Apparati lato GARR

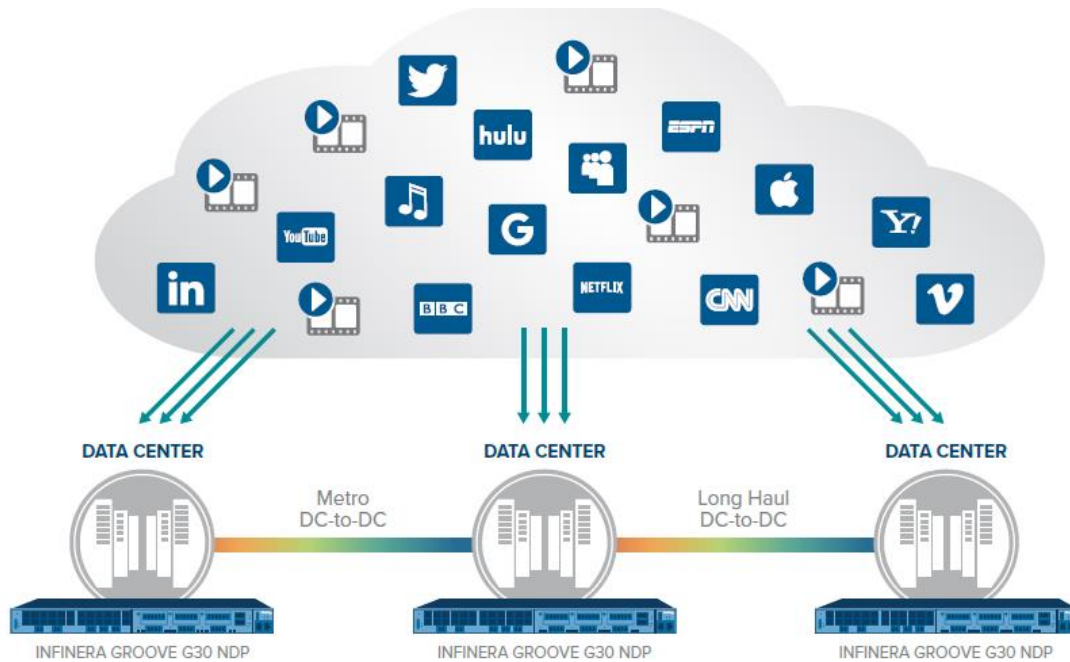


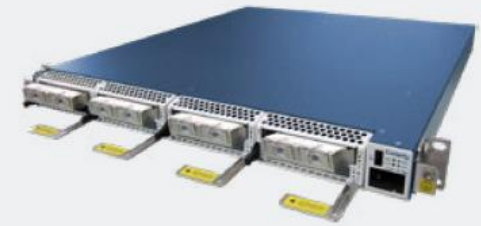
Figure 1: Powering High Performance, Cost-efficient Data Center Connectivity

## THE PURPOSE-BUILT INFINERA GROOVE G30 NETWORK DISAGGREGATION PLATFORM

The Infinera Groove G30 Network Disaggregation Platform (NDP) is an innovative 1RU modular open transport solution for cloud and data center networks that can be equipped as a muxponder terminal solution and as an Open Line System (OLS) optical layer solution. Purpose-built for interconnectivity applications, the disaggregated Groove G30 delivers industry-leading density, flexibility, and low power consumption.

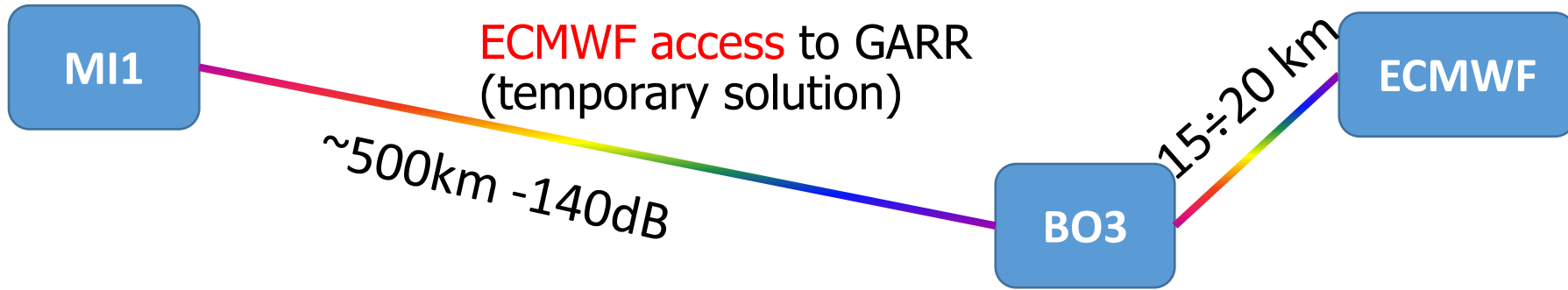


Infinera Groove G30 Open Line System (OLS)



Infinera Groove G30 Muxponder (MUX)

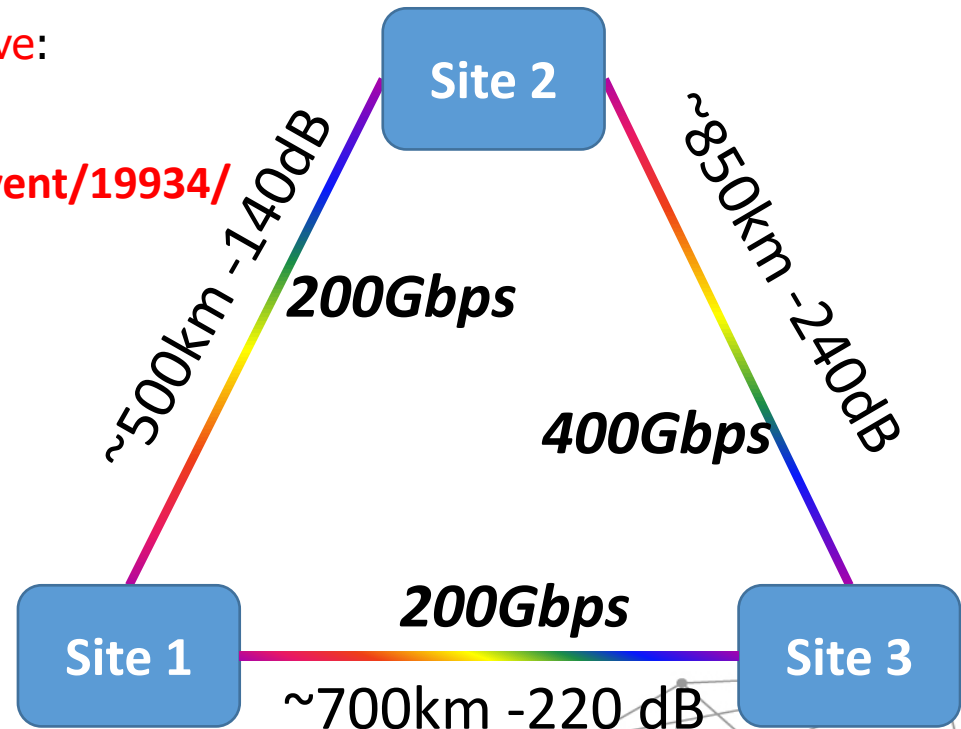
# Which transponder for 100GE



For the **Italian Distributed Data Lake initiative**:  
 Large bandwidth over hundreds of kms

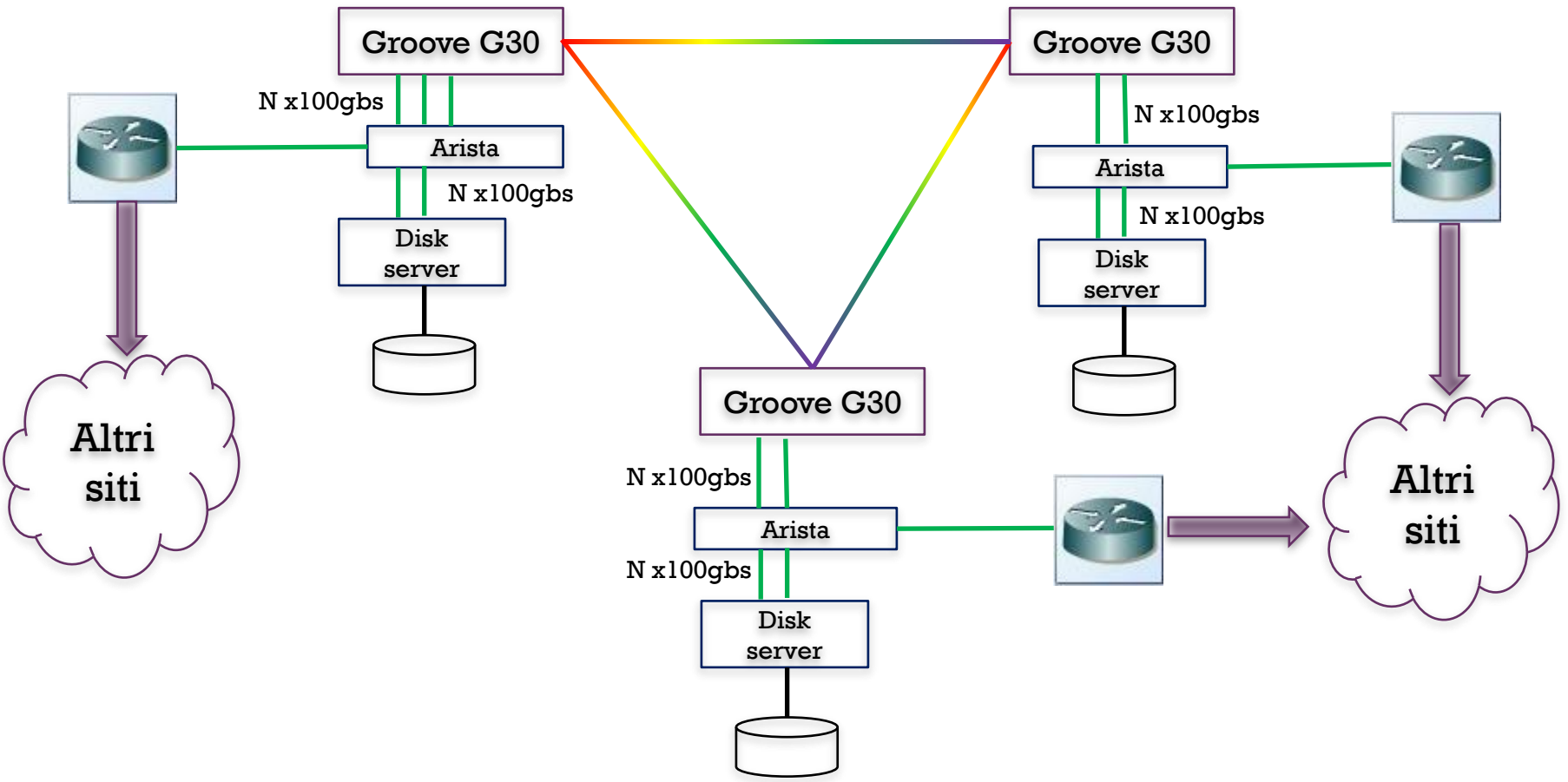
See full test details: <https://agenda.infn.it/event/19934/>

	ECMWF	IDDLs
INFINERA XTM	OK	NOK
Juniper ACX	OK	NOK
INIFNERA GROOVE	OK	OK





# IDDLS Network Architecture





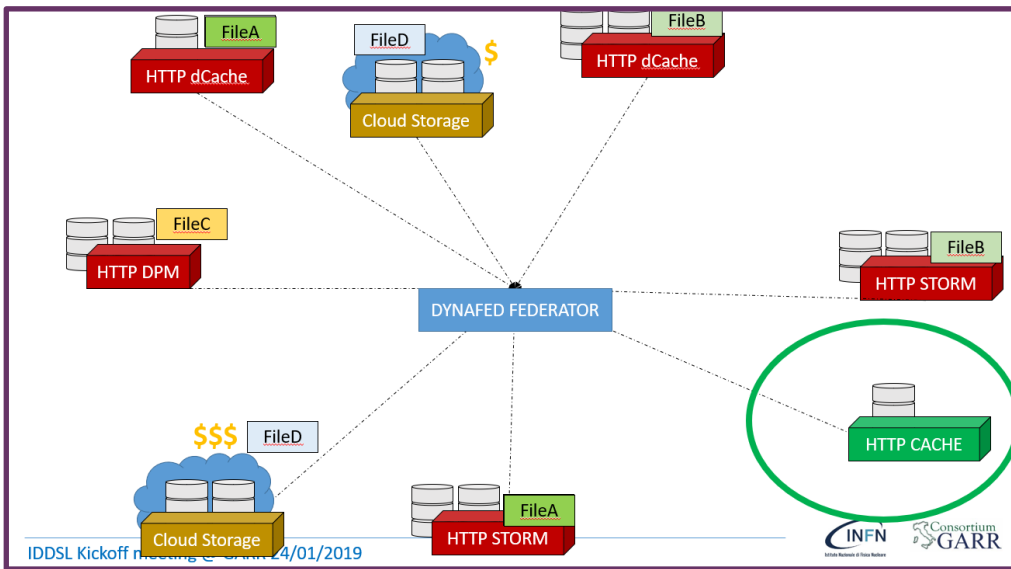
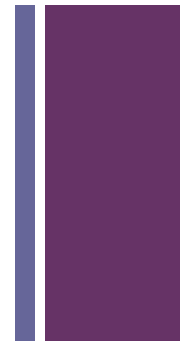


# Apparati rete e storage

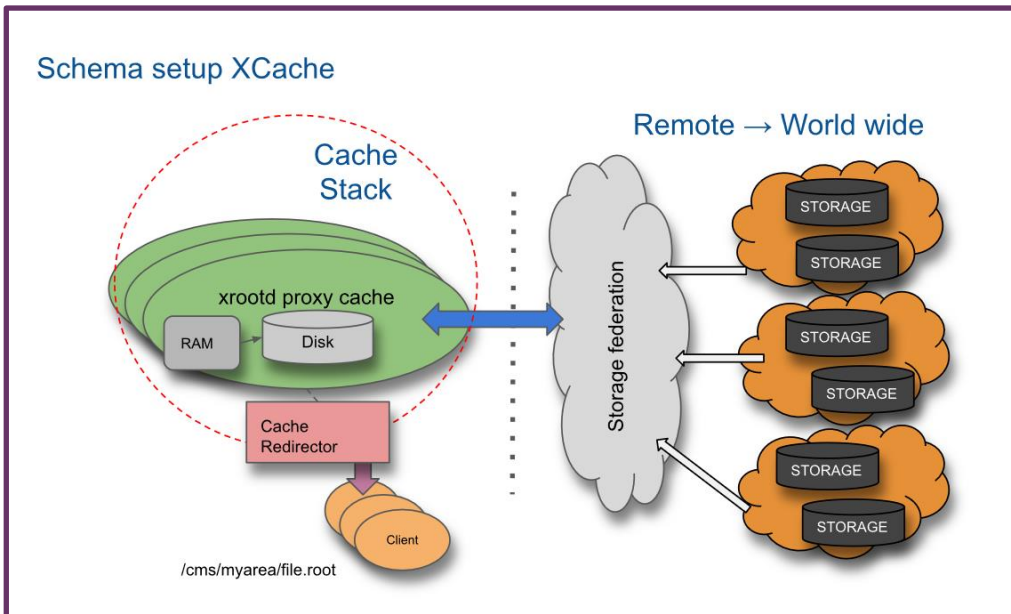


- Acquistato 3xCisco Catalyst9500-32C-A
  - 32x100Gb/s QSFP+ + 2x10Gb/s SFP+
  - + 9 ottiche 100Gb/s
  - manutenzione
- Da definire apparati di storage per test scala seconda parte da progetto
  - 3 disk server finanziati dal progetto
  - 1 server con scheda a 100gbs già acquistato @CNAF
  - 2 server da acquistare su fondi 2020 (15k già sbloccati dal SJ)
- Storage da definire
  - Piccoli sistemi SSD
  - Recuperarlo da dismissioni/produzione se disponibile

# + Lake Architecture(s)



HTTP federation



XROOTD federation + cache

+ QoS support: RUCIO

# + Preventivi 2021

- Finalizzazione preventivi 2021 17/07/2020
  - 1.5k richiesti per ciascuna sede + [1.5k@cnafe](mailto:1.5k@cnafe) per Osservatore
  - Tasca 10k at cnafe per potenziamento storage nei disk server
    - SSD anche consumer
  - Tasca al CNAF per server nelle sedi con CPU: 15k
  - Afferenze CNAF: 1.6FTE

Nome	%
Zani	10
Chiarelli	10
Dell'Agnello	10
Falabella	10
Fattibene	20
Sapunenko	20
Vistoli	10
Giacomini	10
Soares	10
Vianello	10
Cesini	40

# + Gruppi di lavoro

- WP1 – Management
  - Coordinamento, rapporti CSN5 e referee, organizzazione meeting
  - Progress report periodici
  - Procedure acquisti
- WP2 - Studio, definizione e implementazione dei link ad alta velocità
  - Scouting tecnologico delle soluzioni Data Centre Interconnect (DCI)
  - Identificazione dei requisiti degli esperimenti INFN
  - Integrazioni delle componenti HW e SW delle tecnologie DCI
  - Sperimentazione mista laboratorio e infrastruttura di rete geografica
  - Condivisione dello spettro in ambiente protetto;
  - Modelli di provisioning
  - Modelli di gestione e controllo
  - Sperimentazione su infrastruttura geografica su 3 siti

# + Gruppi di lavoro

## ■ WP3 – Creazione del DataLake

- Definizione dello stato dell'arte delle tecnologie esistenti
- Implementazione del DataLake con tecnologie basate su protocollo HTTD/XROOTD con e senza sistemi di caching
- Implementazione del DataLake con tecnologie differenti (eventuali)

## ■ WP4 – Testing del DataLake

- Definizione della testsuite del progetto basata sul software degli esperimenti rappresentati, almeno CMS, ATLAS e BELLEII
- Esecuzione della testsuite sul DataLake sfruttando sia i siti interconnessi con tecnologie di tipo DCI che di tipo legacy
- Interazione con sedi INFN o legate all'ente produttrici di dati (i.e. LNGS, CASCINA) che possano essere interessate a testare le soluzioni del progetto
- Valutazione delle performance ottenute