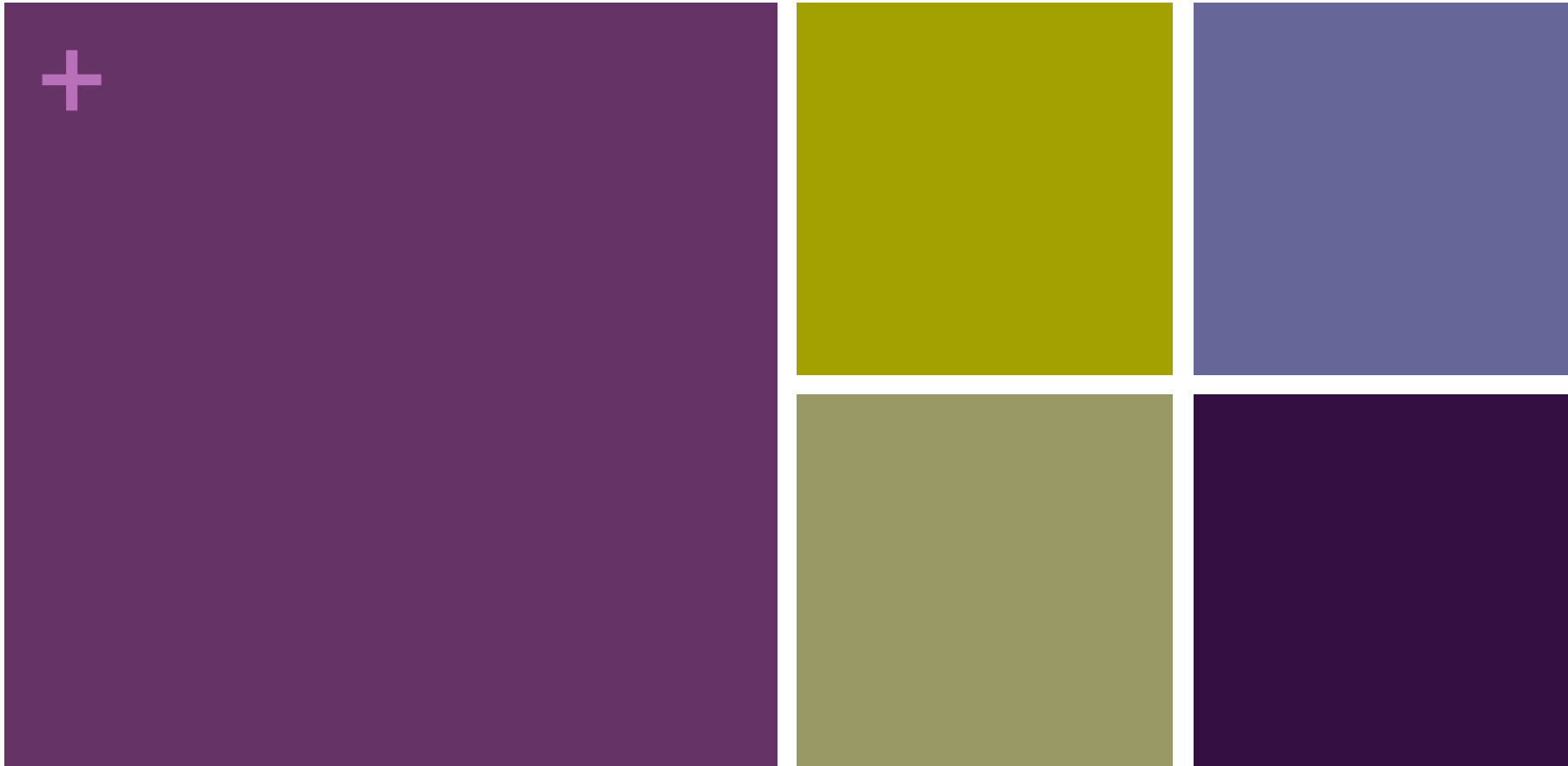


+ Obiettivi kickoff

- Presentazione del progetto e di alcune possibile tecnologie da utilizzare
- Capire i reciproci (inf-n-garr) requirement
- Organizzare gruppi di lavoro
- Definizione canali di comunicazione
- Definizione next step

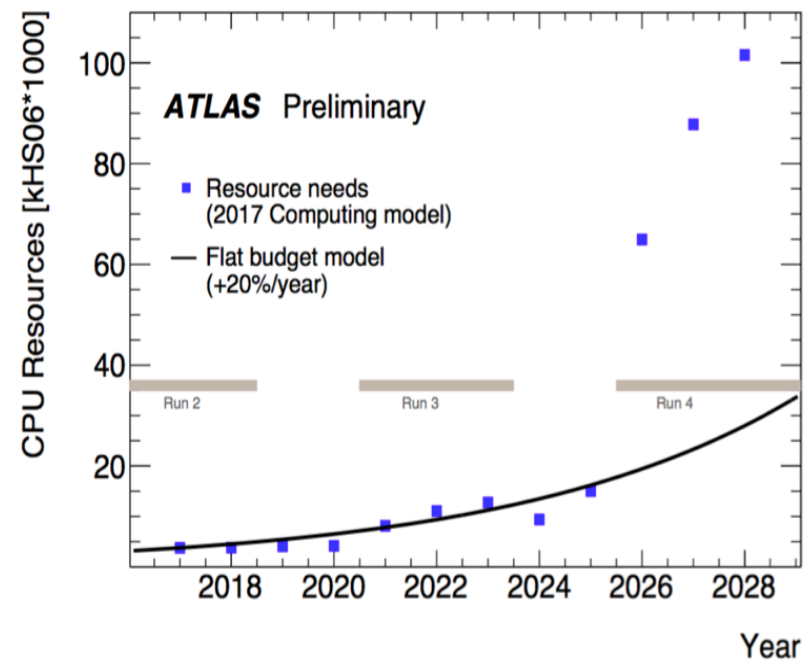
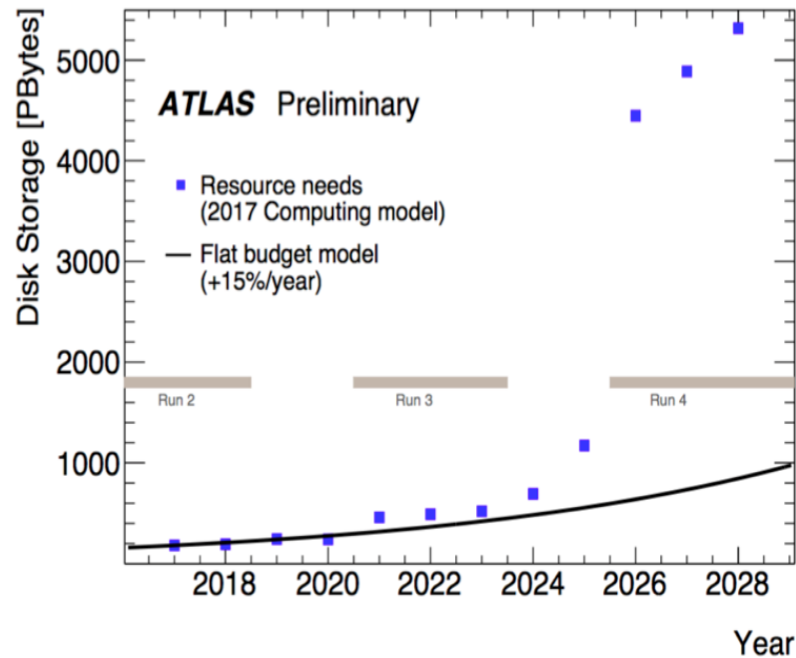




IDDLS: Italian Distributed
Data Lake for Science

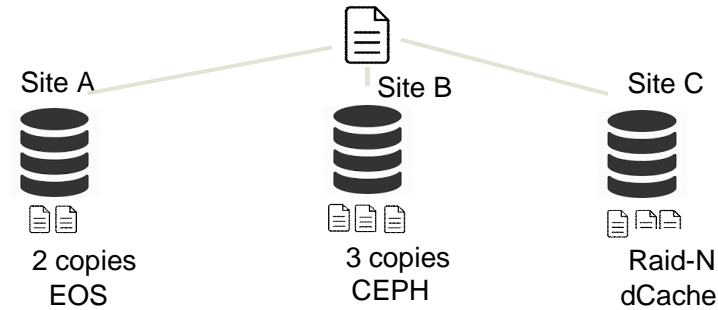
+ Motivation

- HL-LHC computing needs are above the expected technology evolution (15%/yr) and funding (flat)
- We need to optimize hardware usage and operational costs
- High fraction of the infrastructure and operation costs is due to Storage



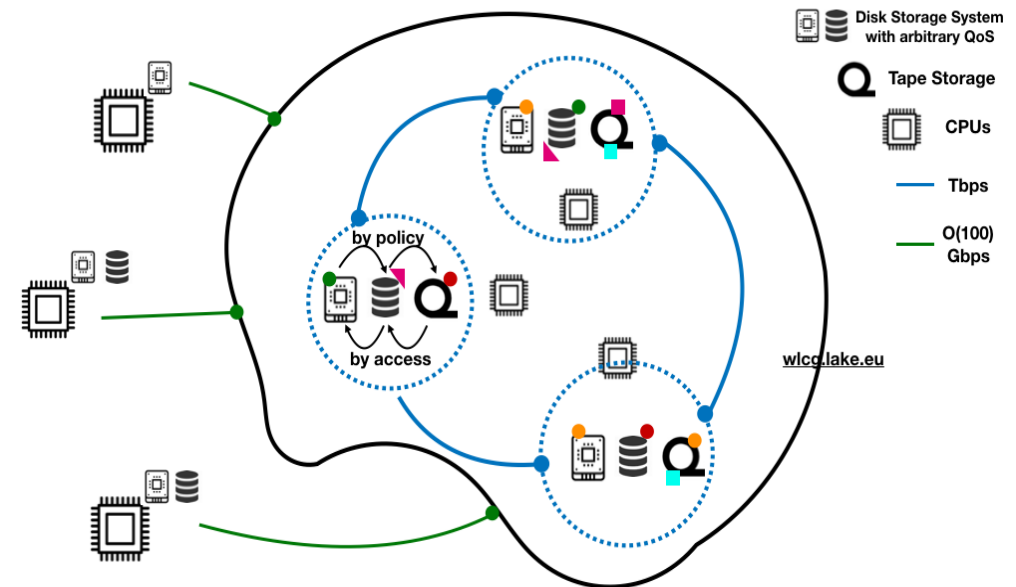
+ Some ideas on reducing Storage costs

- 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



A stronger integration of sites could lead to a reduction of 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

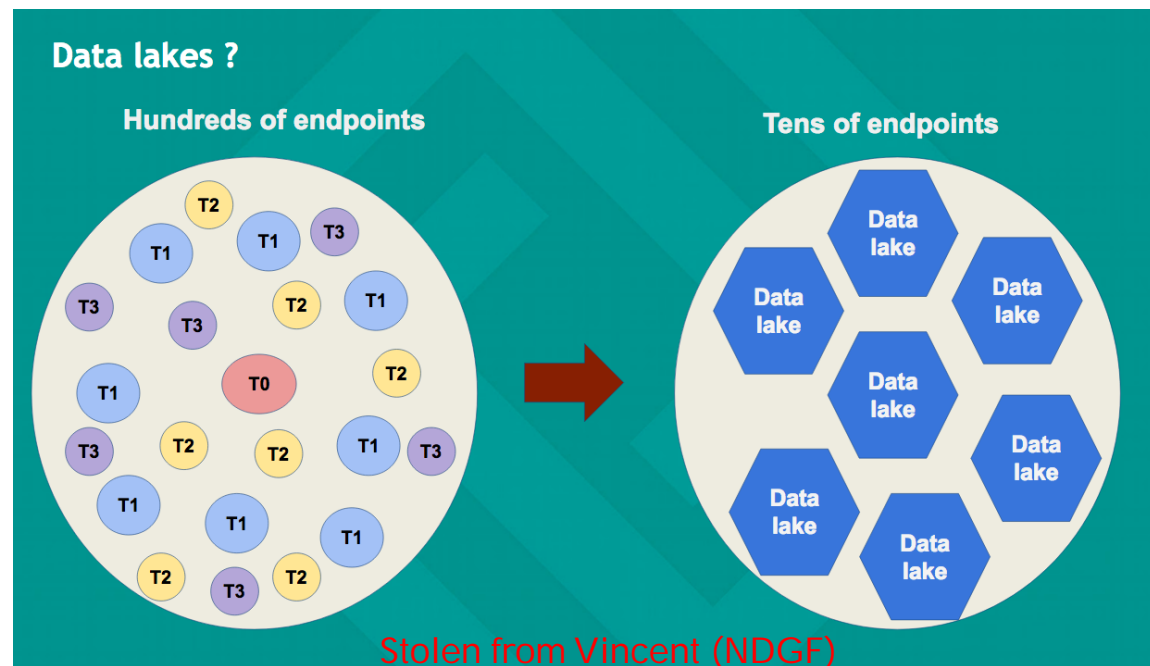


+ On the definition of the *Data Lake*

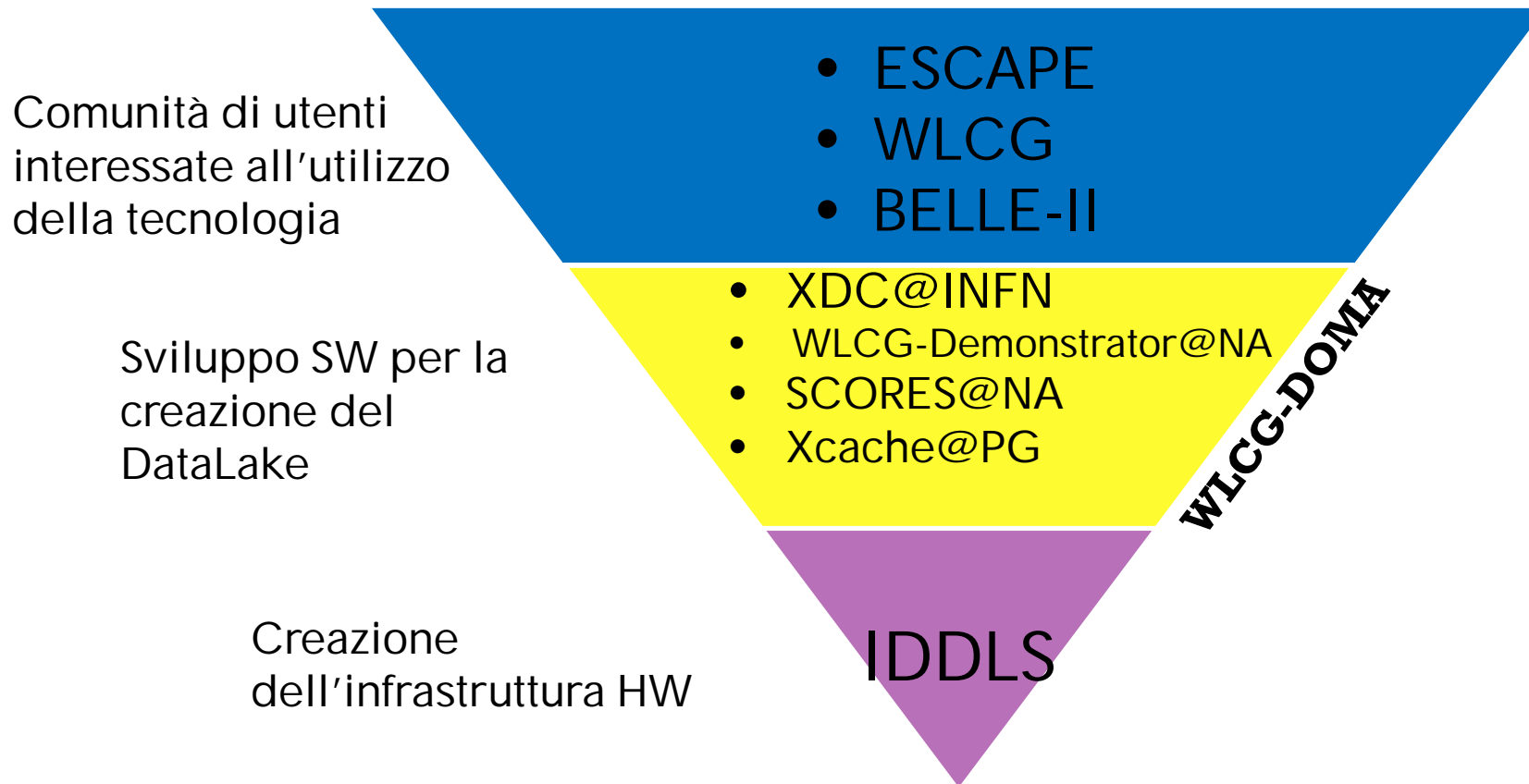
5

- Very diverse understanding of the expression *data lake*.
 - Attempts to define it by 'name space' or 'region' or 'country' all failed.

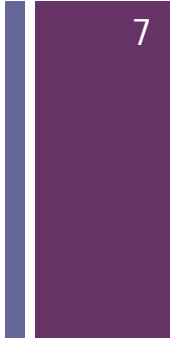
✗ **Looks like one thing, but is composed of many**



+ Sinergie in Europa



+ Sinergie in Europa



Comunità di utenti interessate all'utilizzo della tecnologia

- ESCAPE
- WLCG
- BELLE-II

Sviluppo SW per la creazione del DataLake

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

Creazione dell'infrastruttura HW

IDDLS

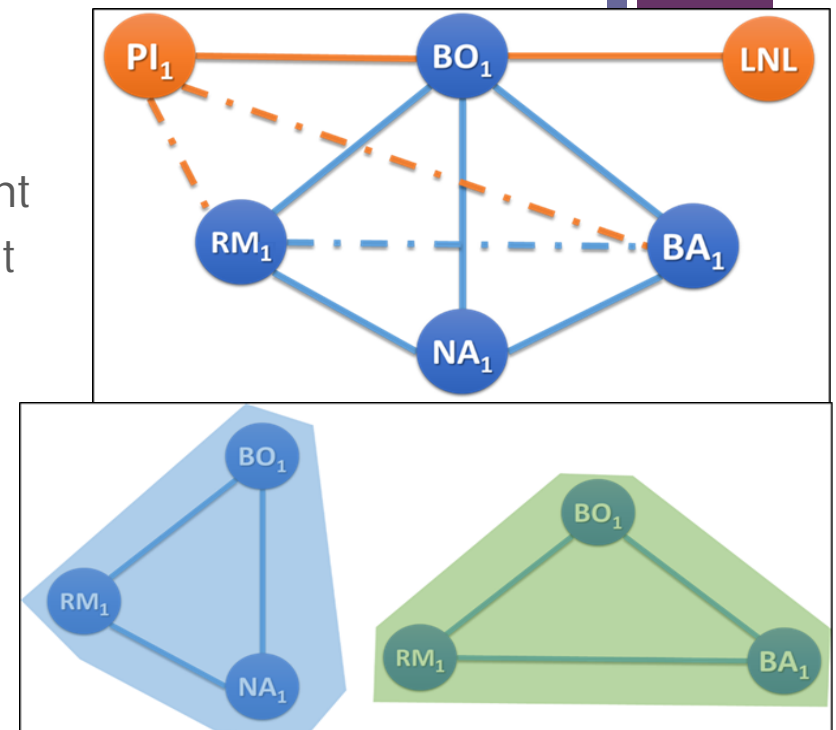
- ESCAPE: Astro-particle cluster (including SKA)
- 32 partner (CNRS lead)
 - 16Meuro
 - 800keuro INFN

The scientific drive towards combining and aligning data from different facilities online and offline will open-up the way **towards the implementation of a *data-lake* infrastructure for astronomy and physics** and it will be offered as a pillar infrastructure to be connected to EOSC for the next decades' data challenges.

+ 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 (H2020 - INFN lead)
 - SCoRES project (INFN-NA)
 - Real life use cases for testing
 - CMS
 - ATLAS
 - BELLE-II
 - Possibly involving LNGS experiments (XENON) and VIRGO

8



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

+ Timeline

■ 3 years project

■ First year

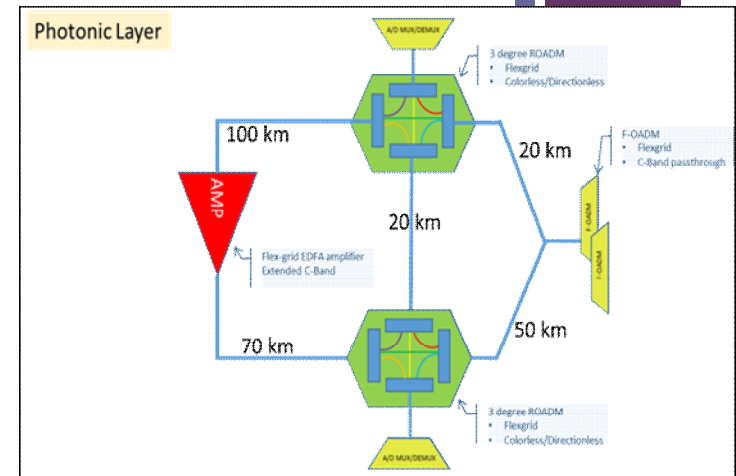
- Technology scouting for DCI equipment to be deployed by GARR
- Application (INFN) requirements analysis
- Network equipment acquisition (INFN and GARR) and Lab testing
- Deployment on mixed Lab+WAN environment of the networking equipment
- Creation of the DataLake on sites connected with standard networking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

■ Second year

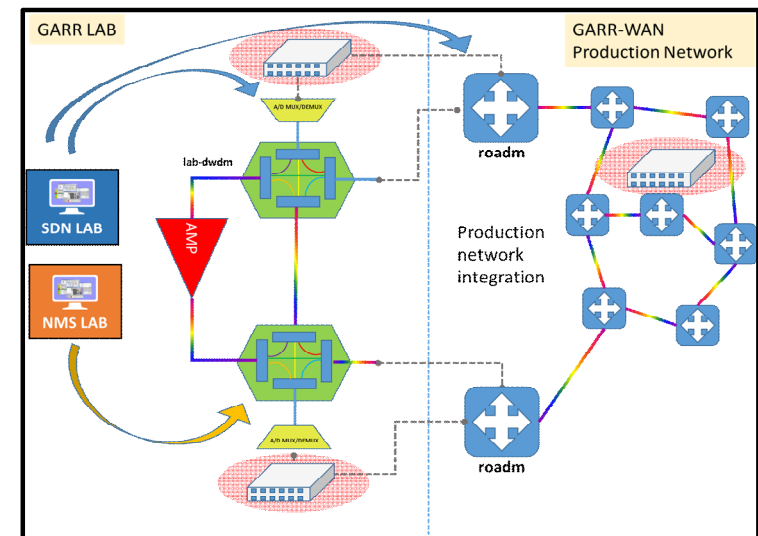
- Testing of the mixed (Lab+WAN) configuration
- Final creation of the DataLake on the 3 INFN sites with DCI systems
- Performance evaluation and comparison
- Possible acquisition of new equipment with increased performance

■ 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

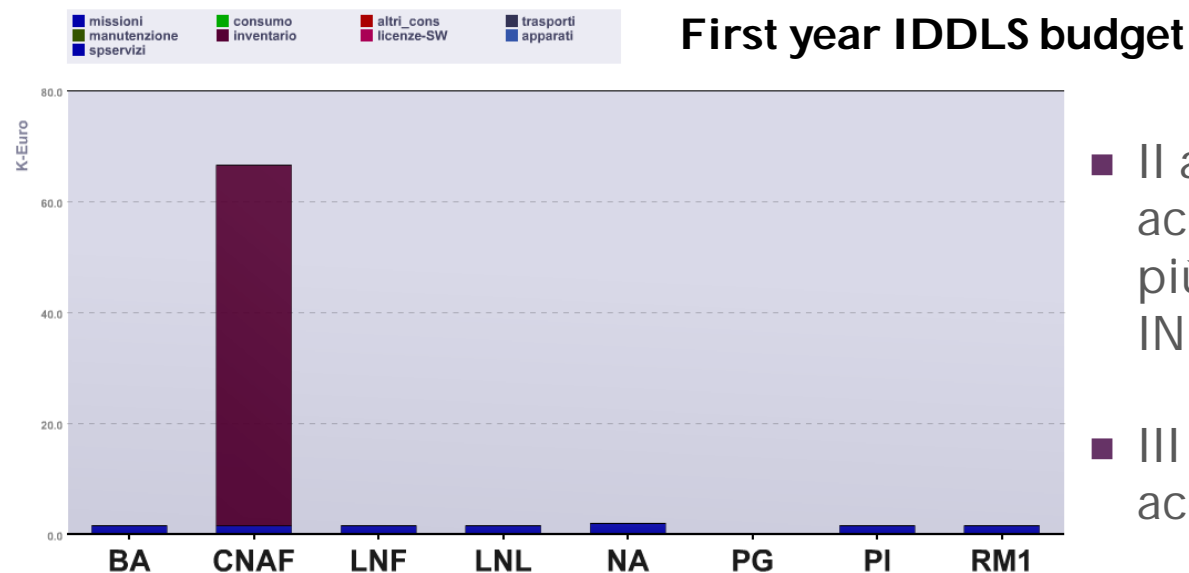
+ 2019 Milestone (simplified)

- 30/06/2019: Scelta degli apparati di networking per la creazione del Datalake
- 31/12/2019: Deployment degli apparati di rete in una configurazione Lab+WAN – primi prototipi DataLake su apparati DCI e standard
- 31/12/2019: Primi run per la valutazione delle performance sui prototipi

+ Budget richiesto

- 200k - contributo in kind del GARR con sistemi DCI per tre sedi INFN
- I anno richiesta tasca di 65k al CNAF per 3 apparati di rete per connettività verso sistemi GARR a 100gbs
 - 3 switch con uplink a 100gbs + 3 NIC + ottiche per 3 sedi
 - gli apparati saranno distribuiti alle sede opportune quando decise
 - studio in fase di progetto per decidere quale sia la migliore topologia
- Missioni: 1.5 per ciascuna sede che ne ha fatto richiesta

11



- Il anno: possibili nuove acquisizione di apparati di rete più performanti da parte di INFN, budget similare
- III anno non sono previsti acquisti

+ Budget Assegnato

12

Sez. & Suf.	MISS			CON			ALTRICONS			TRA			SEM			PUB			MAN			INV			LIC-SW			APP			SPSERVIZI			TOTALE		
	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.	Sj	Dot.	Ant.			
BA	1.5																															1.5				
	0.5																																0			
CNAF	3.0																															68				
	1.5																															6.5	40.0	0		
LNF	1.5																															1.5				
	0.0																																	0		
LNL	1.5																																1.5			
	0.5																																0.5		0	
NA	1.5	0.5																														1.5	0.5			
	0.5	0.0																														0.5			0	
PI	1.5																																1.5			
	0.5																															0.5			0	
RM1	1.5																																1.5			
	0.0																																		0	
TOTALE	12	0.5																														77	0.5			
				12.5			0			0					0																			77.5		
		3.5		0	0																												8.5	40.0	0.0	0.0
				3.5			0.0			0.0					0.0						0.0												0.0		48.5	

+ Sedi e personale

Sede	Personale			FTE			
	Tecnologi	Tecnici	Ricerc./Borse	Tecnologi	Tecnici	Ricerc./Borse	TOT
CNAF	7	2	0	1.4	0.2	0	1.6
BARI	3	2	0	0.2	0.2	0 +0.9	1.3
LNF	1	0	0	0.1	0	0	0.1
LNL	4	0		0.6 +0.4	0	0	1.0
NA	3	0	2	0.3 +0.1	0	0.1 +0.5	1.0
PG	1	0	3	0.1 +0.1	0	0.2 +0.5	0.9
PI	2	0	3	0.1 +0.45	0	0.2 +0.3	1.05
RM1	1	0	0	0.1	0	0	0.1
TOT.	21 +1	4	6 +2	2.9 +1.05	0.4	0.5 +2.2	7.05

13

GARR: 3x15%

+ Afferenze

- Vedi excel



+ Gruppi di lavoro

15

- 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

16

- 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