

Data Preservation

Marcello Maggi

Coordination on Data Management with
Tommaso Boccali & Luca dell'Agnello

Definition

Digital Data are affected by digital obsolescence

Medium

Access

Semantic

Preservation means

bit preservation

access preservation

knowledge preservation

Why DP

Observations are unique:

Astrophysics, Earth Science pioneered DP:

FITS, Open Access Policies, OAIS

Experiments are reproducible:

but LEP, Tevatron, LHC, etc. will not

data exploitation by a larger community exists

Existing FP7 projects

SCIDIP-ES

(SCIENCE Data Infrastructure for Preservation - Earth Science)

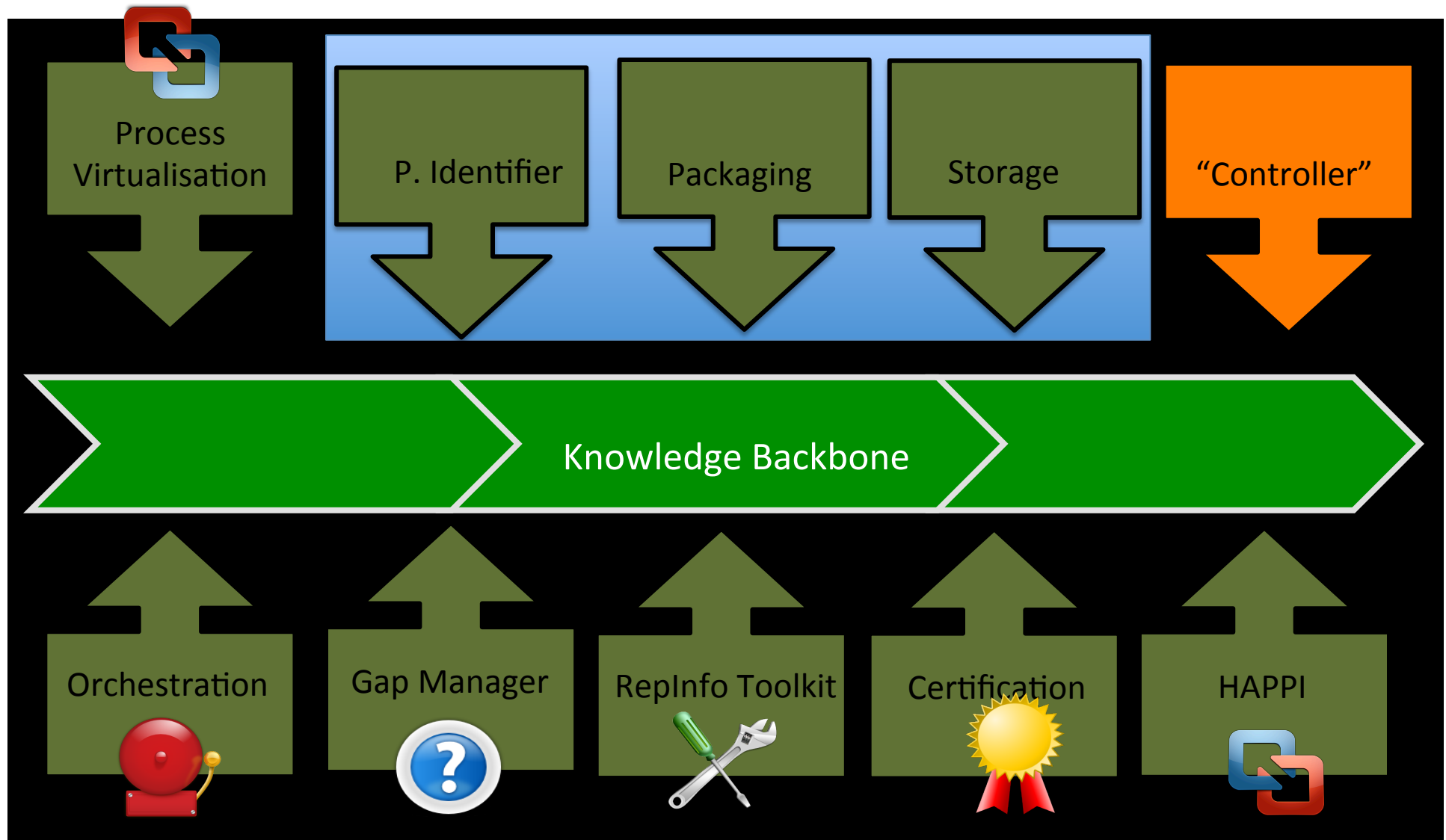
- Main Long Term Data Preservation. project Eu call INFRA-2011-1.2.2
- It address the issue of building the key information (knowledge) to allow access and understanding of experimental data in a technology independent way such that the preservation is really long term.
- The project implements components based on OAIS

EUDAT

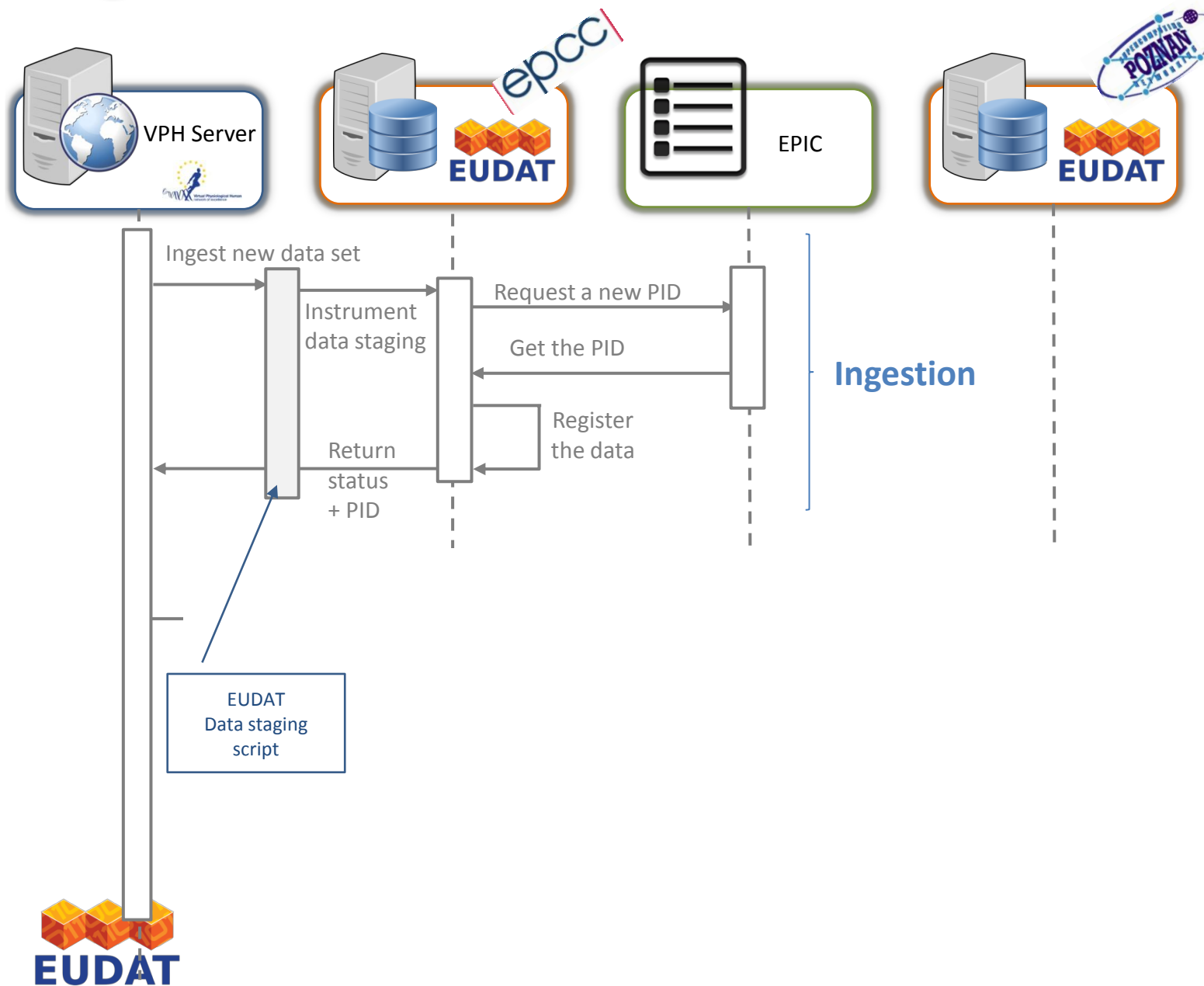
(EUROPEAN DATa infrastructure)

- Builds data e-Infrastructure for storing, preserving and sharing through standard tools and services
- Data preservation is considered only for the bit preservation part

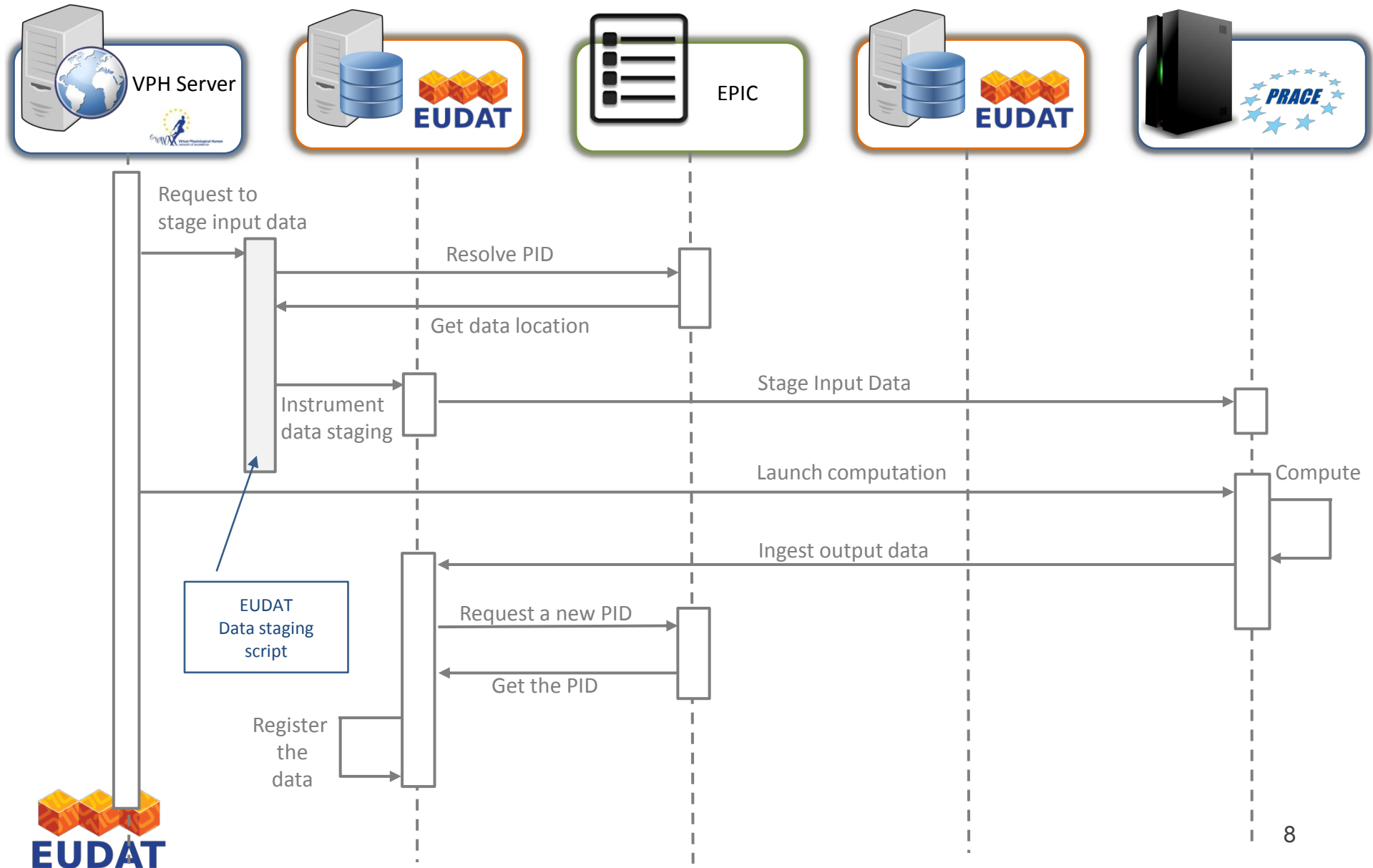
SCIDIP-ES



Ingestion



Staging

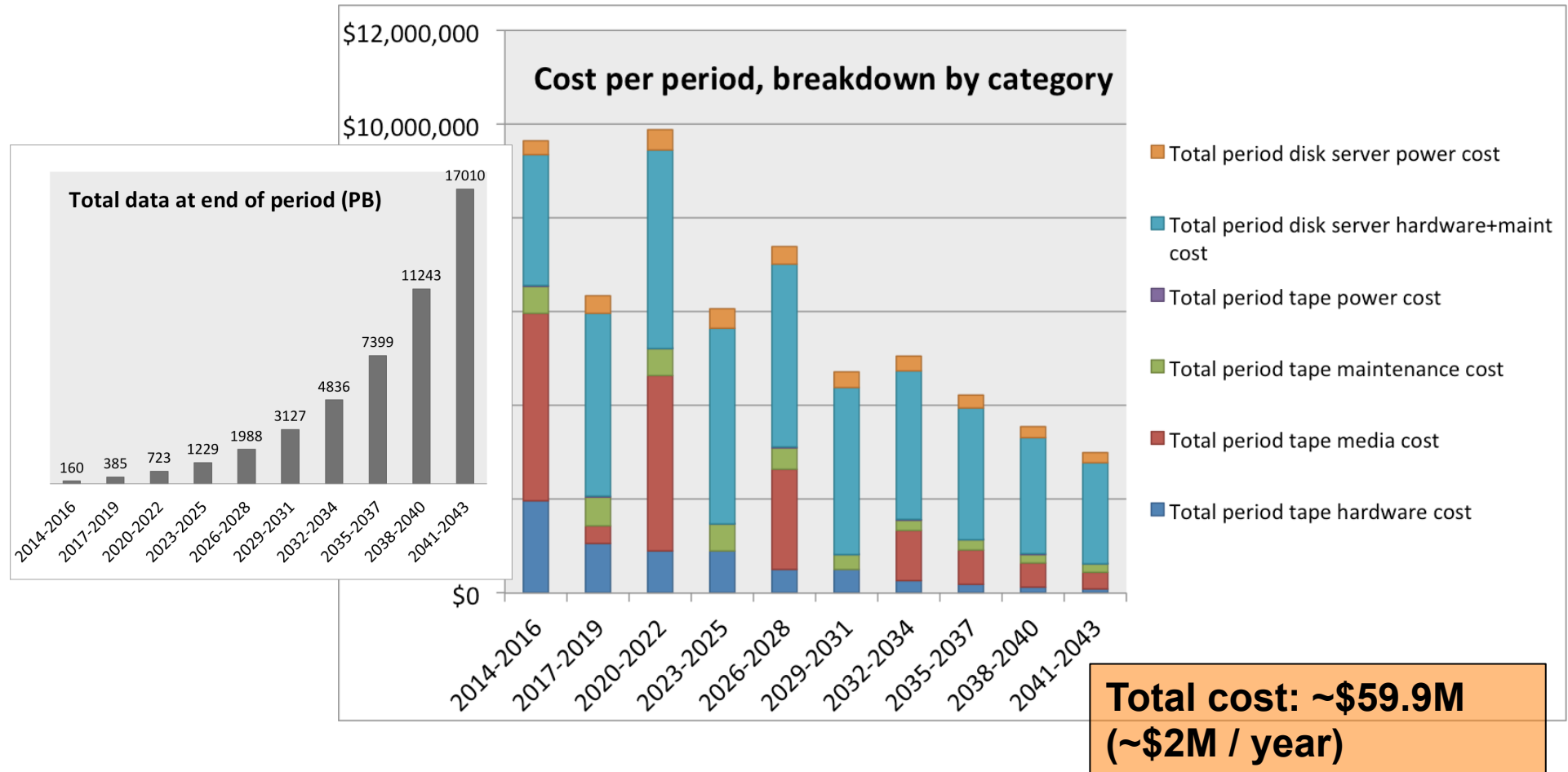




International Collaboration for Data Preservation and Long Term Analysis in High Energy Physics

- ICFA panel since 2009
- DPHEP Project Manager appointed in October 2012 Jamie Shiers, CERN
- Intense activity in the last year:
 - Two DPHEP workshops (CERN)
 - Monthly Implementation Board meetings
 - Participation to many conferences and workshops: visibility and new opportunities
- CA agreement proposed (ICFA statement 2013)
 - Declare interest and possible areas of contributions (light commitment)
 - Strong support for further funding programs (H2020)
 - Signed by CERN, positive echo and possible imminent signature from DESY, INFN, discussions ongoing in US

Cost of Curation

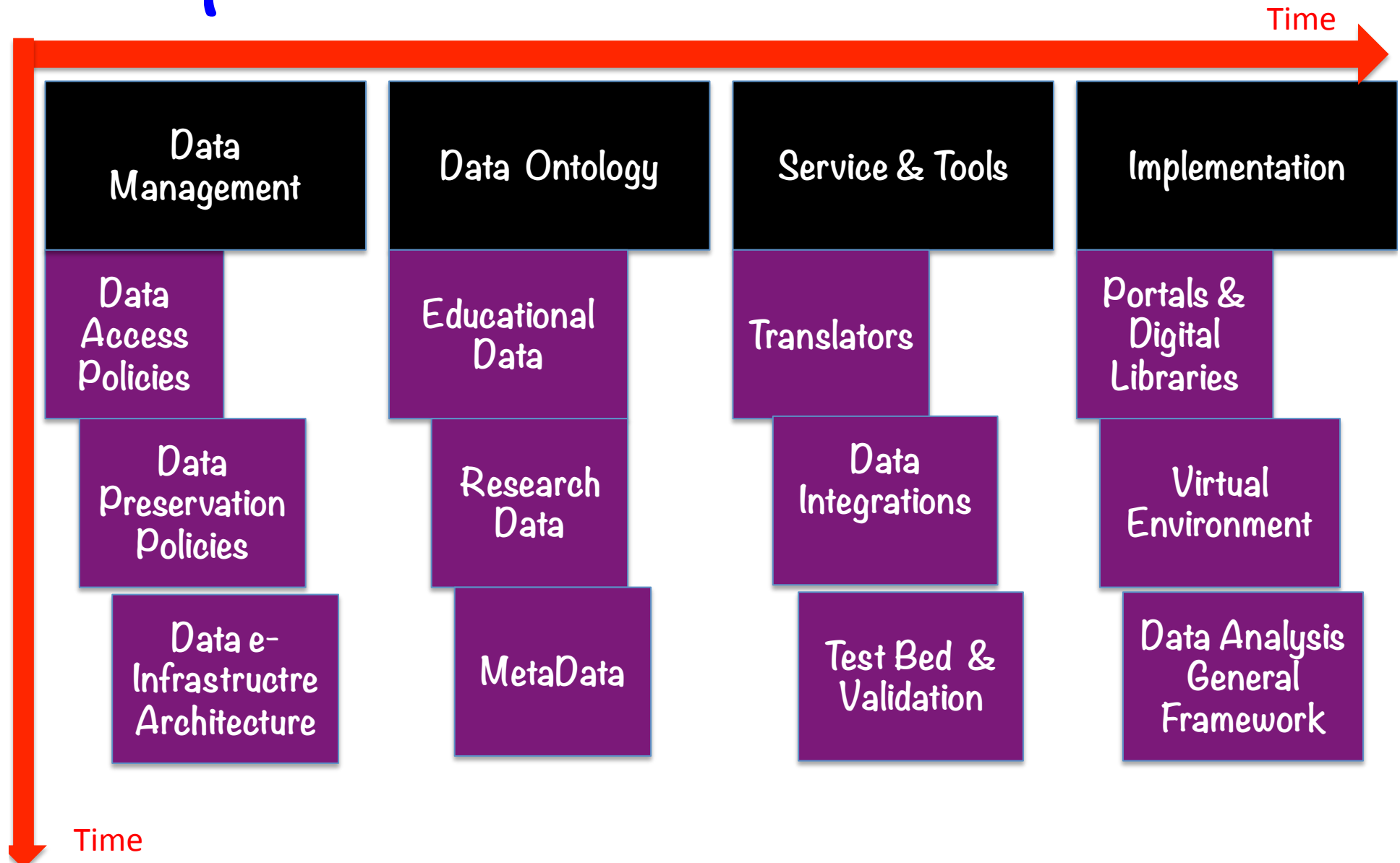


The annual cost of WLCG
(infrastructure, operations, services)
is ~100 M€

DPHEP Portal

- **Digital library** tools (Invenio) & services (CDS, INSPIRE, ZENODO) + related tools (HepData, RIVET, ...)
- **Sustainable software**, coupled with advanced **virtualization** techniques, “snap-shotting” (incl. CernVM[FS]) **and** **validation** frameworks
- Proven bit preservation at the 100PB scale, together with a **sustainable** funding model with an outlook to 2040/50 (+HEPiX coordination + RDA WG + H2020?)
- **Open Data** (“Open everything”)

Open Data Work Plan Structure



The enemy

Each single experiments design proprietary SW
and data format for resource optimization...

Not needed in Long Term

LEP decade: Factor 300 increase of CPU power

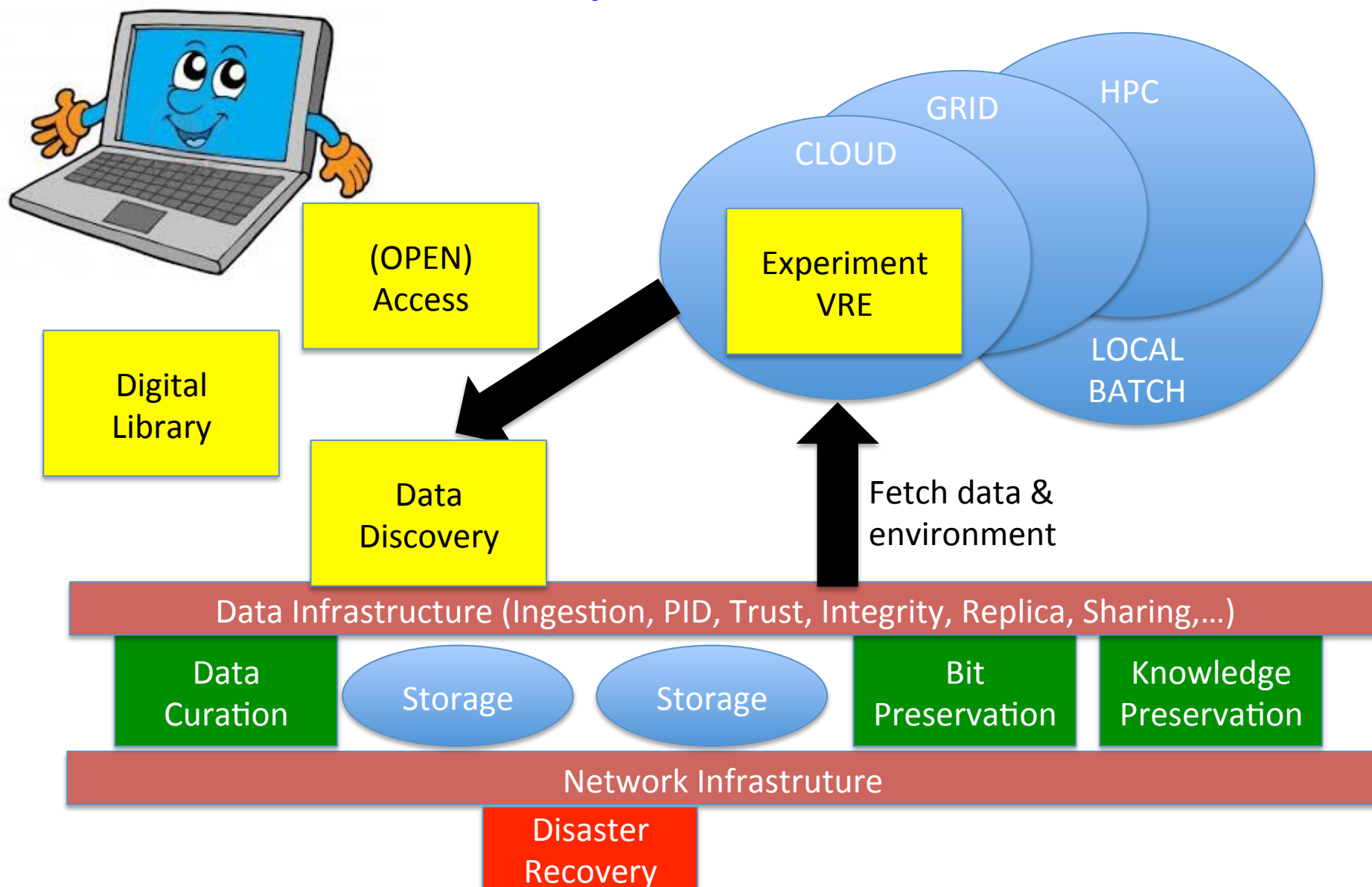
SHIFT50 DEC Alpha had 320 Cern Units = 2.5 SpecHEP

Today 1 machine is 64 core & 560 SpecHEP

Projects..?

- DPHEP portal: build in collaboration with other disciplines, including RDA IG and the APA...
Proposal being prepared
- Digital libraries: continue existing collaborations & projects (also H2020 funding...)
Services in place and in use
- Sustainable “bit preservation” – certified repositories as part of EINFRA-1-2014 – based on RDA WG “RFCs” for interfaces, functionality, federation etc.
Also a service! Certification and standardisation desirable...
- “Knowledge capture & preservation”: collaborate with SCIDIP-ES, APA(RSEN) and many others
- Open “Big Data”: a Big Opportunity (for RDA?)
Proposal exists – to be discussed in DPHEP shortly & with EU JRC / EUDAT / RECODE

The Picture



The proposed implementation in EUDAT

- Eudat competitive call for data preservation

ALEPH as “model experiment”

30 TB divided in 150000 files with a typical size of 200 MB.

for

Data Volume Collection

Data Discovery System

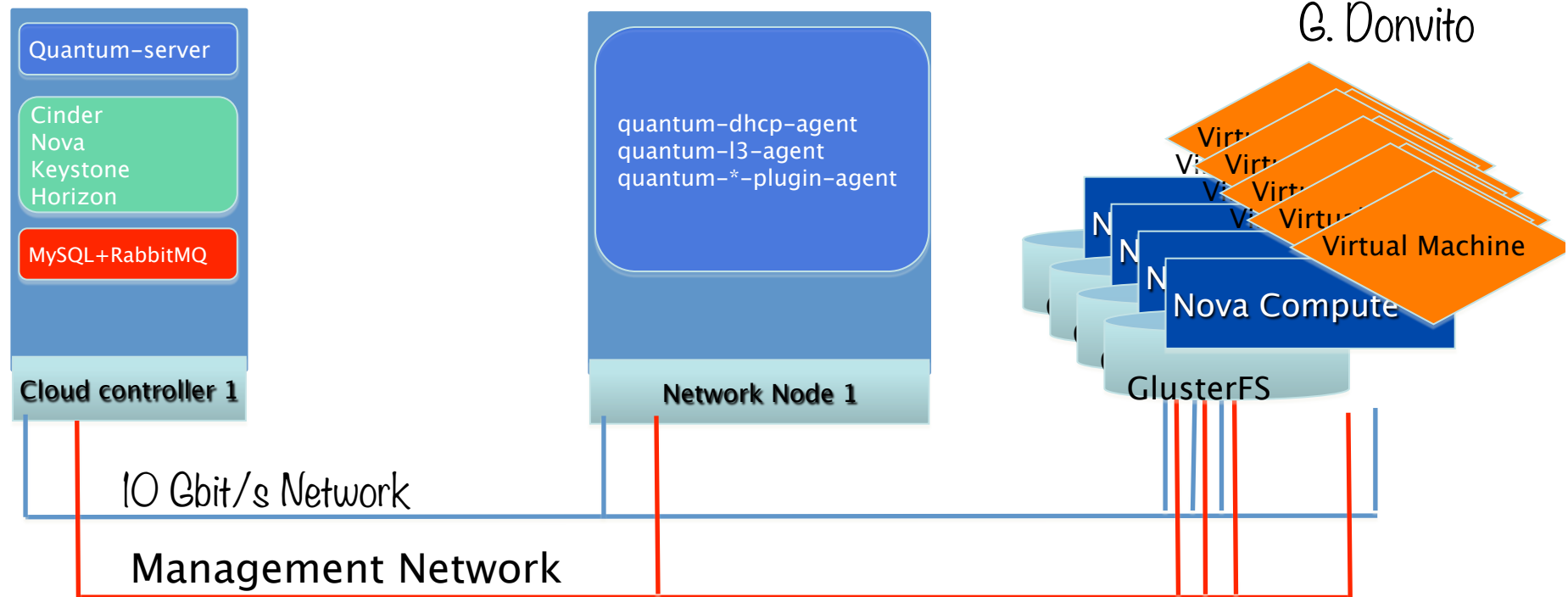
Data Analysis Cycle

Interplay between

a tailored Virtual Research Environment

and the use EUDAT tools to access the data

Implementation in Prisma Cloud



- 11 server
- 264 CPU/Core
- 880 GB di RAM
- 66TB HDD (DAS) 7.2K rpm

- 1 Manager Node
- 1 Network Node
- 10 Gbit/s on Wide-Area-Network
- In production with KVM nodes

Summary

- In several key areas, there are **services** and **solutions** in place that needs to be made interoperable and be integrated (“lego” project)
- Other requires further R&D (Scidip-research?)
- Preservation costs are negligible in a long term