



Big Data, Legacy & Preservation (for science and technology)

Mario G. Lattanzi
INAF-OATo



1- I Big Data sono fra noi e sono (anche) tanto nostri (INAF)!

Gaia HTC@DPCT – Primo Big Data store del cielo

- **Piattaforma Operativa e di Test & Sviluppo**
- **Acquisizioni HW e SW 'incrementali' eseguiti nel tempo attraverso una attenta pianificazione delle necessita' della missione.**
- ✓ **INTERNET LINK PER RICEZIONE DATI GREZZI GAIA DA ESRIN (Spagna) : 1Gbps (300 Mbps guaranteed) via GARR.**
- ✓ **SPAZIO DISCO: 2.5 PB di spazio distribuito tra 2 unita' di 'storage' HP P7400 ed una HP P8400.**
- ✓ **HPC PIATTAFORMA OPERATIVA: 14 servers HP DL580 G7/G9 per un totale di oltre 600 CPU cores e 4.5 TB di memoria RAM.**
- ✓ **PIATTAFORMA DEV & TEST: 7 servers HP (per altri 300 CPU cores)**
- ✓ **DATA BASE SERVERS: 3 servers HP DL580 G7 (32 CPU cores, 256 MB RAM ognuno) basati sulla tecnologia Oracle RAC (DBMS Oracle).**
- ✓ **INFRA MONITORING AND MANAGEMENT: services based on VMWare virtual environment configured with two HP DL 580 G7 servers clustered and managed by vCenter Server.**
- ✓ **BACKUP SERVERS: HP DL580 G7 dedicat al back-up del DB ed ai backups del filesyste via 'snapshots' dei 'data volumes'.**
- ✓ **3 LIVELLI DI BACK-UP : L1 on primary storage array, L2 on disks (StoreOnce 6600) and L3 on tape libraries (HP ESL G3).**
- ✓ **COLEGAMENTO DIRETTO AI SISTEMI HPC DEL CINECA: accesso ai sistemi HPC della classe Marconi CINECA per la ricostruzione della sfera celeste di Gaia (via MOU dedicato CINECA-INAF)**



Piattaforma operativa del DPCT presso ICT ALTEC



DPCT-Operation Room @ALTEC

Sezione Piattaforma Test e Sviluppo I DPCT@ALTEC

DECT will OPERATE AS SUCH
TO 2025 → 2030 (?)



to 4+ PBy

2- Anche la legacy Gaia e' gia' qui (vedi DR3 e prossima DR4)

GAIÀ HAS ALREADY
ENTERED ITS LEGACY ERA

WHAT IS IT?

α WHAT TO PRESERVE
("CONSERVATION")
AND HOW TO PRESERVE

α WHAT INFRASTRUCTURE

☐ From DBMS-based pipelines



{ File system - based pipelines
+ Metadata DB

↳ This was the key

☐ From "pbin" → FITS
(sequendo la BAV)

N. B. BAV= Biblioteca Apostolica Vaticana

THE GAIA LEGACY IS NOW!

IN ITS ORIGINAL FORM PRESENTED AS
PART OF "NITIC" (GARICCI, PI)
→ produced the TLS
experiment or

« The Living Sky Prototype »

Fattibilità
conclusa
nel 2021

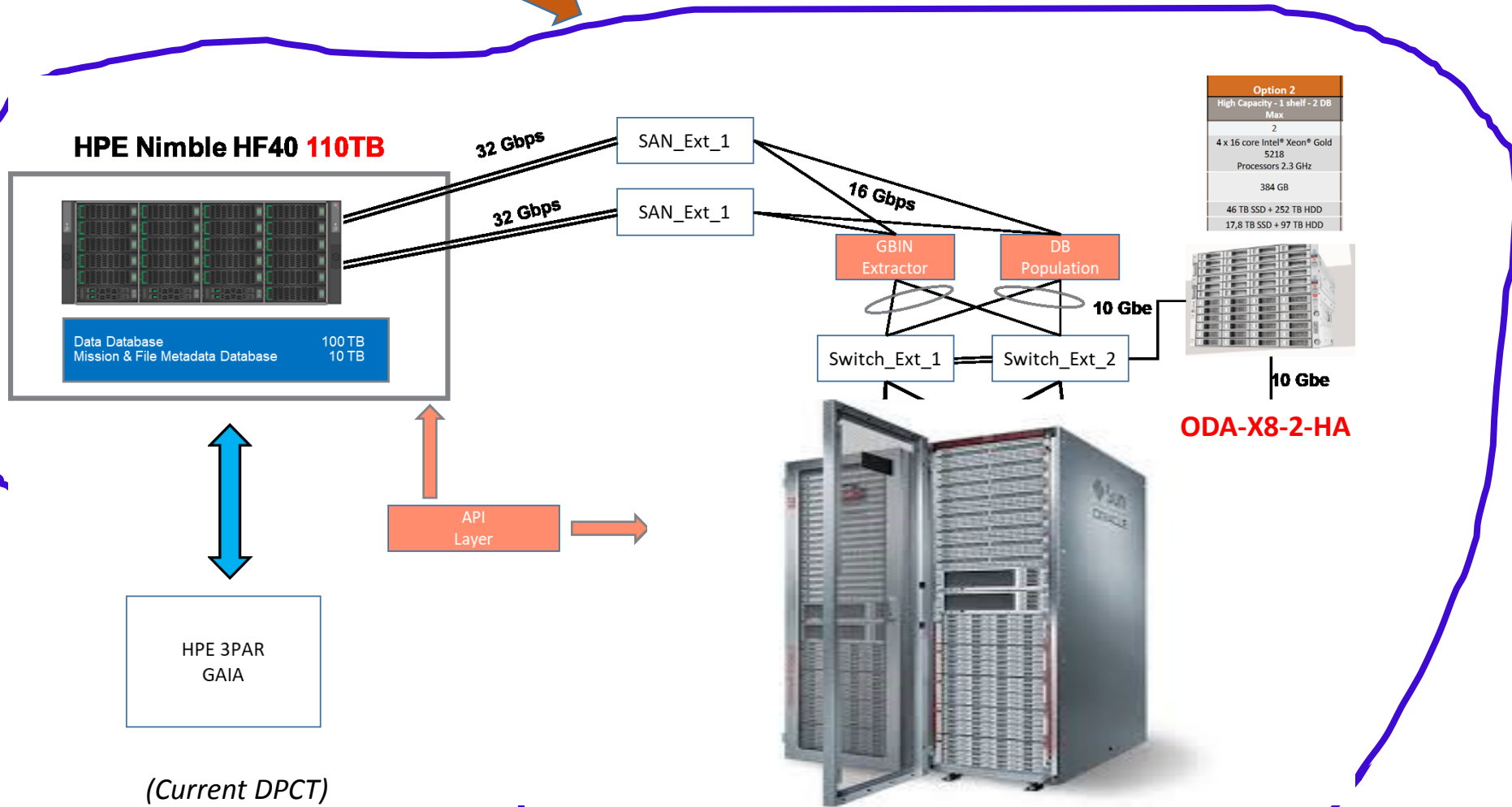
100 TBY NAS + small
ODA from OATS
Fully integrated at
DPCT

TLS prototype produced
3 POC'S : FEASIBILITY { File system !!!
+ Meta Data DB ...
SERVICE TO OPERATIONS, LEGACY
(→ SERVING SCIENCE... & possibly
more)

Generic
Data
Requests

CHANGE OF PARADIGM
POSSIBLE

Full scale realisation: OPS4 expansion of DPCT (with ASI resources)



Option 2
High Capacity - 1 shelf - 2 DB Max
2
4 x 16 core Intel® Xeon® Gold 5218 Processors 2.3 GHz
384 GB
46 TB SSD + 252 TB HDD
17,8 TB SSD + 97 TB HDD

ODA-X8-2-HA
 64 CPU
 18 TBy SSD
 97 TBy HDD

ZFS ZS9-2
 480 TBy (RAID2)

Oracle ZFS file system compatible con DPCT.

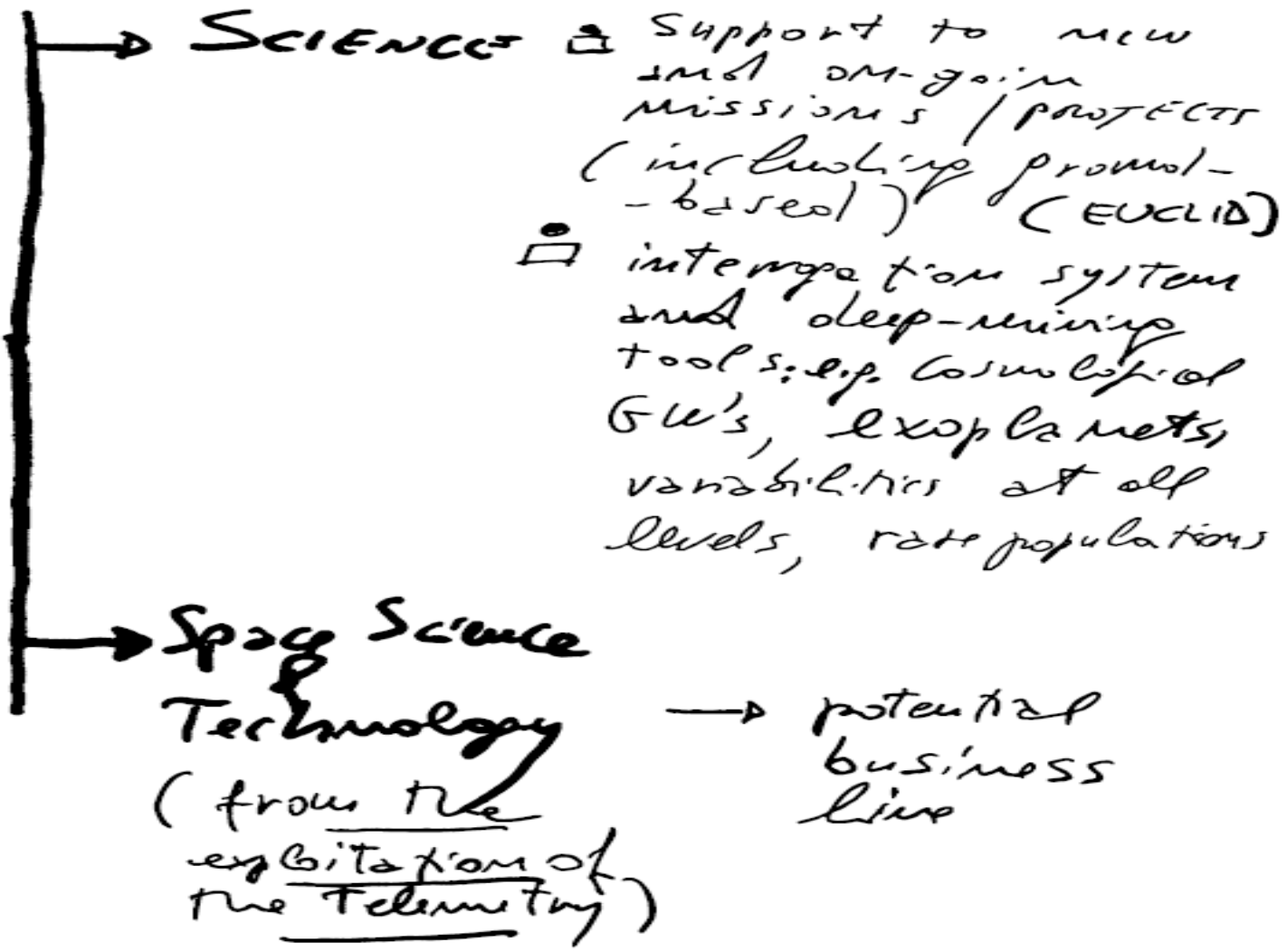
^{GAIA}
WITH OPERATION RESOURCES
(AND THE DEDICATION OF PEOPLE
→ THE HUMAN FACTOR!)
WE HAVE INITIATED THE
IMPLEMENTATION OF THIS
'HYBRID OBJECT' WHICH
CAN BECOME THE CORE/ENGINE
OF THE INAF TLS OR
"la biblioteca INAF del
cosmos (la cosmoteca)".



OPS4

3- Toward Exascale HTC with Big Data

STRATEGIC FUNCTIONS



Supporting thousands of simultaneous complex queries/ transactions

Oracle Exadata Database Machine X9M-2: quarter rack (minimum)

OPS4+TLS MUST BE
CONCERNED WITH VERY LONG
TERM PRESERVATION &

INTEROPERABILITY!
PRESERVATION: not only ^{PERSISTENT} data
accessible at any time, but also
PRESERVE PROCESSING CAPABILITY
AND EXPAND.

ALL THIS MUST BE IMPLEMENTED BY A
DEDICATED TEAM WITH THE TALENT
OF MULTIDISCIPLINARITY AND THE
FIGURE OF DATA SCIENTIST FORMED AND
"OWNED" BY IMAF.

WITHIN AN INFRASTRUCTURE THAT
IS ECONOMICALLY AND ENVIRONMENTALLY
SUSTAINABLE.

Conclusions

The Gaia Legacy is a leading project for INAF in the context of its spoke of the National Center for HPC&Big Data HTC and is awaiting PNRR resources: clearly defined 3-year implementation plan and sustainability prospects.

Needs:

- Oracle adequate Exadata Machine
- 3 dedicated positions (data scientist, DB specialist, SW specialist)