# Data Management and Archiving in Astronomy:
## Status and Challenges for the Future

## Fabio Pasian

**Istituto Nazionale di Astrofisica – O.A.Trieste**

Euclid Consortium Science Ground Segment Manager

former Chair (2007-2010), International Virtual Observatory Alliance

ICTP, Trieste, 5 September 2013

# FOREWORD

## From this meeting's Web page:

«Due to the increasing usage of computer simulations, the management of large amounts of scientific data has become a challenge.»

## From the ESO Web site:

«The current total archive holding is about 65 TB, with an input rate of about 15 TB per year. This will soon drastically increase by a factor of 10 or so as the Visible and Infrared Survey Telescope for Astronomy (VISTA) with its near infrared camera will alone produce about 150 TB of data each year.»
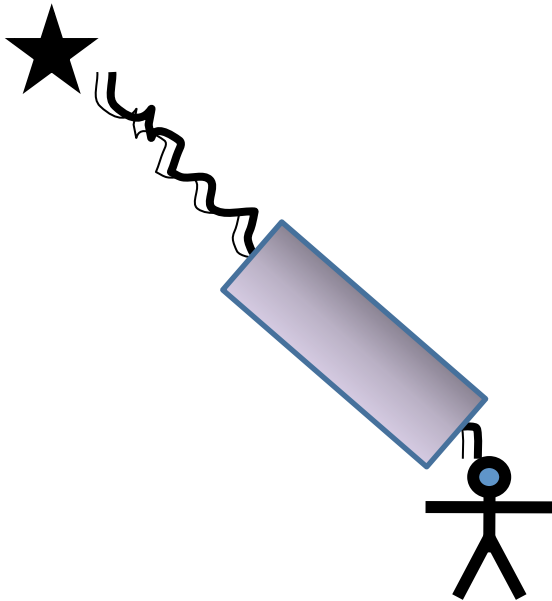
## From the «SKA: the ultimate big data challenge» paper:

«The challenge? How to transport, store and process the estimated 14 exabytes of data the antennas will gather every day.»

# SUMMARY

- Basics on astronomy

- Data sharing and standards

- The Virtual Observatory

- The big surveys
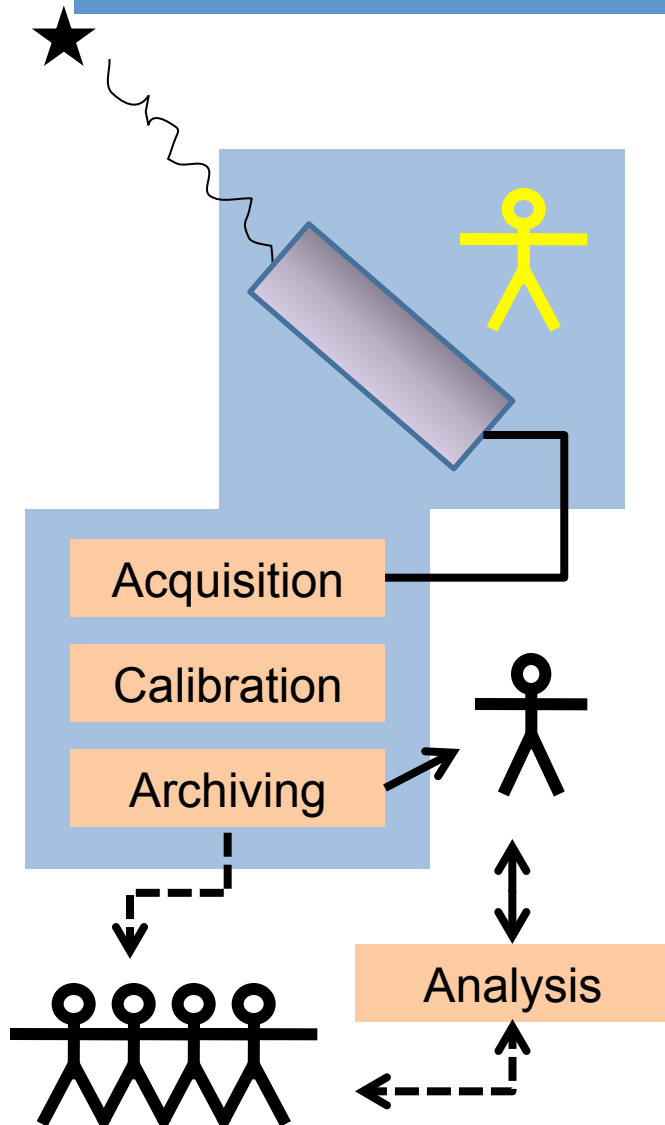  - An example: the Euclid project

- … in Trieste

- Astronomy is an <u>observational</u> science - no repeatability

- Measures $I(\lambda)$, or $I(\nu)$, or $I(E)$

- We actually observe $I'(\lambda)$ $= I(\lambda) \ * \ T(x,\lambda)$

- Most phenomena are in fact variable

- We actually observe $I''(\lambda,t) = I(\lambda,t)*T(x,\lambda,t)$

- Every single observation must be kept $\Rightarrow$ preservation!

Acquisition

Calibration

Archiving

Analysis

- **Service observing** due to:
  - complexity of instrumentation
  - optimisation of observational conditions
- Data handling and basic processing (calibration) **c/o observatory**
- Owner of observation gets data **from archive**
- After some time, data become **public**
- Since 1977, there is a **standard format** for data files (FITS)

5

# FITS

- **F**lexible **I**mage **T**ransportation **S**ystem
- Born in 1977 for images
- Later extended for N-dimensional data arrays, ASCII and binary tables, event lists, …
- 3 continental committees, IAU WG, permanent FITS office (in GSFC/NASA) maintaining standards and sw
- Standard format
  - ASCII header containing keywords with data description
  - standard record size: 28800 bytes
  - readers available in any astronomical application
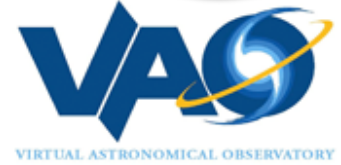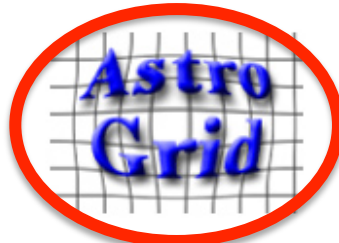- (now used by Vatican Museum for digitised books)

6

# NEED FOR ARCHIVES

- Monitor time variability of phenomena (digitization of photographic plates)

- Need to reprocess raw data given better knowledge of instrumental effects

- Compare phenomena in different bands (multi-λ astronomy)

- Increase return for investment (data re-use, educational, outreach, …)

- Statistical analysis / mining of large quantities of data

- Cope with data avalanche

Outreach Products

Calibrations

Processed Data

Raw Data

Data Headers

Observatories

Software

Logs

Instruments

Analysis

Data Products

Telescope control

vobs

Catalogues

Publications

Tables

Figures

Computing

Papers

# THE IVOA

➢ **Mission:** *"To facilitate the international coordination and collaboration necessary for the development and deployment of the tools, systems and organizational structures necessary to enable the international utilization of astronomical archives as an integrated and interoperating virtual observatory"*

- Works by telecons, "Wiki" pages, and bi-annual meetings (last one in Heidelberg [May 2013], next in Waikoloa [Sept 2013])

➢ **Needs:** standardization of data/metadata/sw, data interoperability methods, and list of available services (provided by projects)

- Structure:
    - ✓ IVOA Executive Board includes representatives from all VObs projects
    - ✓ Working and Interest Groups (including Theory)

http://www.ivoa.net

**Euro-VO : 5 EU/FP funded projects (2004-2014)**

- Findings
  - The Virtual Observatory concept is a bold community-led response to the challenges the astronomical community faces in data management and storage.  Impressive progress has been made and the momentum of the International Virtual Observatory Alliance will ensure sustained progress, provided the agency level support and funding is available.

- Recommendations
  - New projects and facilities must take the data management, storage, maintenance, and dissemination into account at the earliest planning stages, consulting potential users in the process. Agencies should recognise that this is an important long term issue and should co-ordinate plans, provide adequate funding on a long-term basis, and support development and maintenance of the needed infrastructure. *Agencies should encourage the broadening of the existing VObs collaboration into a fully representative global activity.*

# ESFRI and ASTRONET statements

- ESFRI (multi-disciplinary)
  - focus on networking, capability & throughput computing, grid architectures, software, <span style="color:red">data management and curation</span>

    → Research Data Alliance

- ASTRONET (Astronomy & Astrophysics)
  - recognised as must-haves to tackle the challenges of the future (<u>priority in resources</u>):
    - computing (capacity AND capability)
    - theory & simulations
    - <span style="color:red">virtual observatory</span>
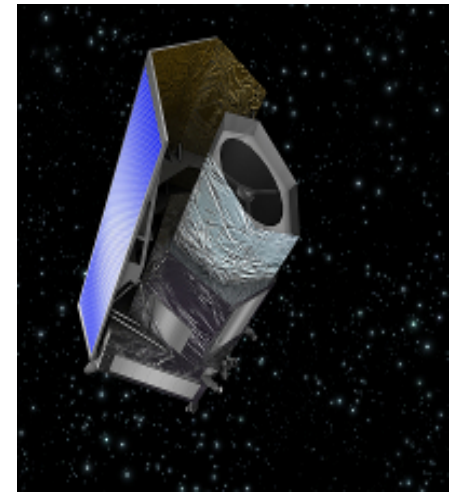    - laboratories

# BIG SURVEYS

- Digitized Sky Survey (DSS): a digital version of several photographic atlases of the sky

- Sloan Digital Sky Survey or SDSS: major multi-filter optical imaging and spectroscopic redshift survey

- PanSTARRS, DES, LSST, ...

- ROSAT All Sky Survey: 1378 distinct fields in X-rays

- Planck: whole sky in 9 $\lambda$ bands from radio to infrared

- → Euclid: whole extragalactic sky (15000 degrees$^2$) in visible and IR imaging, IR spectra

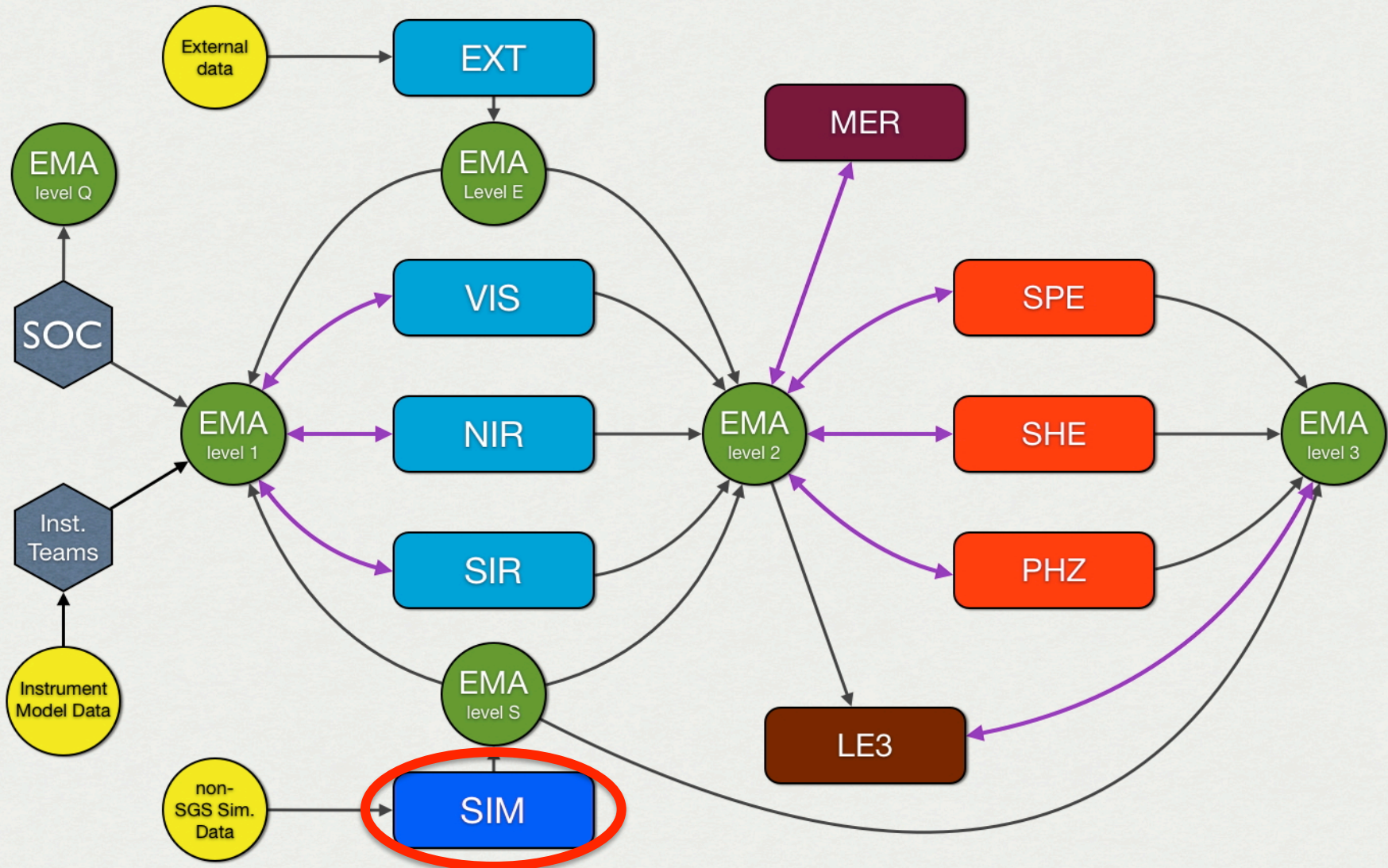- ALL DATA ARE INTEROPERABLE USING VOBS

# THE EUCLID MISSION

How did the Universe originate and what is it made of?

- Goal: To understand the nature of dark energy and dark matter by accurate measurement of the accelerated expansion of the Universe

- Targets are galaxies and clusters of galaxies out to z~2, in a wide extragalactic survey covering 15 000 deg², plus a deep survey covering an area of 40 deg²

- Two instruments (NIR imager and spectrograph + visible imager) + need for external data

- Approved by the European Space Agency (Oct 2011)

- Date of launch: March 2020



15

# EUCLID DATA FLOW



SOC

$100 \text{ GB/day} \rightarrow 10^{14} \text{ B/3years} = 100 \text{ TB/3years}$

simulations
excluded !

(compressed data everywhere!)

60%    30%    10%

VIS    x5x3    1 PB
NIR    x5x3    0.5 PB
SIR    x2x3    60 TB

$1 \text{ Gb/s} \rightarrow 10 \text{ TB/day}$
off-line transfer?

100 Mb/s    50 Mb/s    7 Mb/s

MER
2xEuclid + Ground    13 PB

EXT    10 PB

100 Mb/s

PHZ    x2    < 0.1 PB

SHE    x1    1 PB

SPE    x1    0.4 PB

LE3    x1    0.1 PB

# EUCLID GROUND SEGMENT: OVERALL VIEW



Euclid

scientific community

**VObs**

EA is built jointly by EC and SOC,
and is managed by SOC
ELA is an EA function, allowing public access
to a subset of EA data

An organization based on the
decomposition in Organization
Units (OU), corresponding to a
subset of overall EUCLID Data
Processing.

Each OU produces
algorithms which are
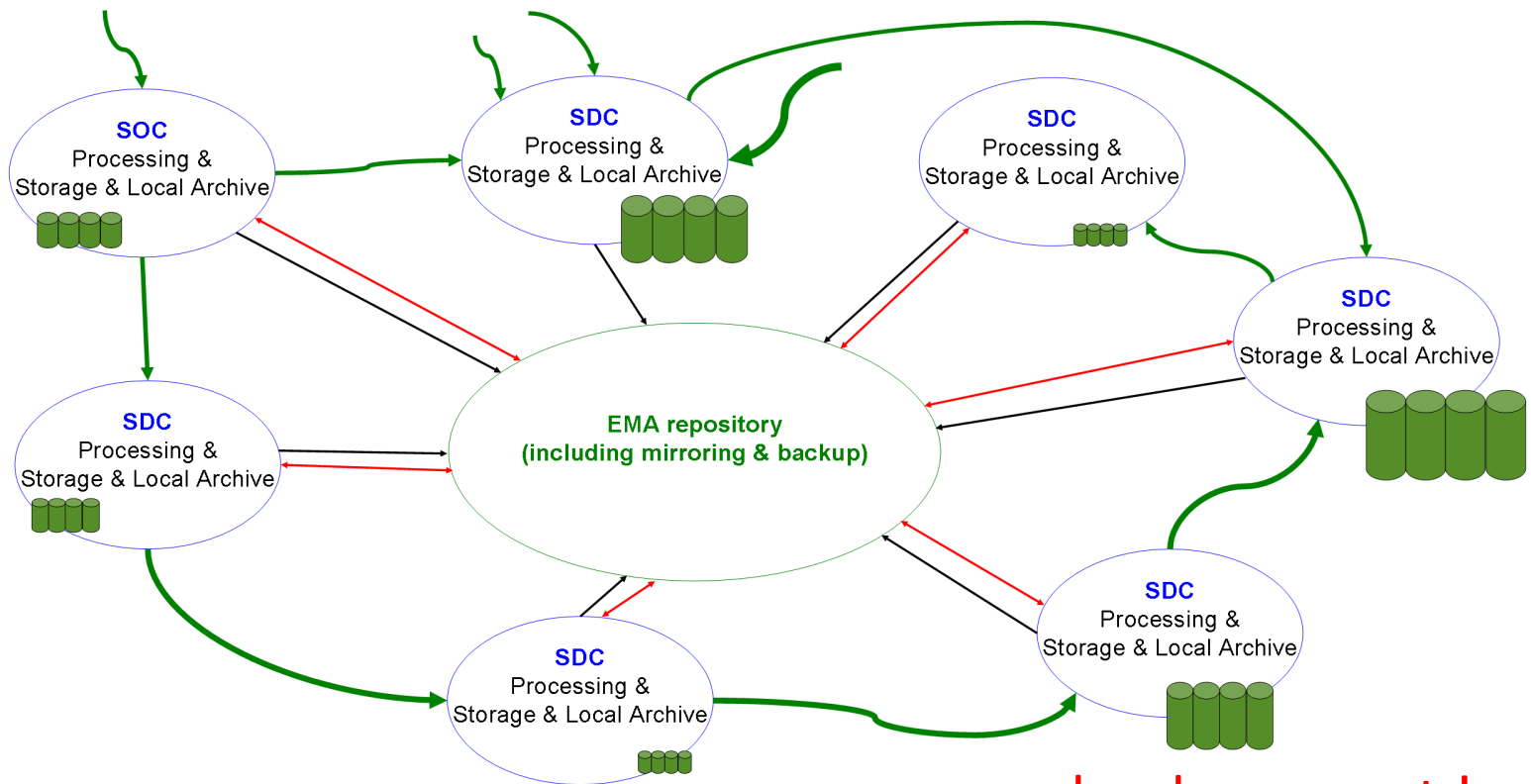integrated and executed
in an SDC (Science
Data Center)

The SGS System Team
provides support and
tools for the whole of the
SGS (SOC + EC-SGS)

Ground
Station

MOC
ESOC

DDS

SOC
ESAC

ELA

EA

External data
(PanStarrs, DES, ...)

EC-SGS
Project Office

Simulation

OU-SIM

SDC   SDC   SDC   SDC   SDC   SDC   SDC   SDC

OU-PHZ   OU-SHE   OU-VIS   OU-LE3   OU-SPE   OU-SIR   OU-MER   OU-NIR   OU-EXT

Phot Red Sh.   Morpho & Shear   VIS Imag   Level 3   Spectro Meas   Nir Spectro   Euclidisation   Nir Imag   Ext Data
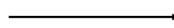
**OUs are transnational**

OU coordinator
OU Deputy Coordinator

18

# EUCLID ARCHIVE SYSTEM

- For internal use within the Euclid Consortium:
    - Centralised metadata repository (at ESA) based on DBMS
    - Distributed data files (images, spectra, catalogues)
    - At the core of the (data-centric) Euclid processing:
        - Processing «pipelines» work on fractions of the sky stored on local computing centre
        - Data replicated at some other centre as necessary

- Public data are replicated at the European Space Astronomy Centre of ESA (ESAC) and from there are made accessible to the external community through VObs-compliant interfaces.

# EUCLID DATA + PROCESSING LOGICAL ARCHITECTURE



Data Product (image) transfer between SDC's

Metadata update

Query/Metadata

«cloud» concept !
full pipeline is run on sky tiles
virtualisation – Hadoop?

## IN ITALY (… actually, led by TRIESTE)

- Coordination of INAF ICT Unit
- INAF centre for Astronomical Archives (IA2) and VObs.it (~ 150 TB)
- Planck- LFI Data Processing Centre (~ 250 TB)
- In Euclid:
  - management of Science Ground Segment (Project Office)
  - design, development and operations of SDC-Italy
  - development of NIR, SIR, MER, SPE, LE3, SIM processing
  - coordination of the infrastructure needed for simulations
- Participation in MIUR-approved DHTCS-IT project (Distributed High-Throughput Computing System) → cloud for research
- Submitted projects to MIUR («premiali»):
  - AstroCloud (INAF-led, with UniNA and SISSA)
  - participation in PIDES (multi-disciplinary data preservation)
- Participation in RDA and CINECA «Big Data» initiatives

21

# CONCLUSIONS

- Increasing detector performance + astronomical all-sky surveys: from Big Data to **Bigger Data** (VObs-compliant)

- Theory data and modelling add up to the data volume – there is no limit to the data that can be produced, the physics to model is already there – (VObs-compliant)

- Individual projects may be able to cope through a cloud-like approach

- «What can we afford to throw away?»

- Long tradition in data sharing (VObs) → standards (they are a blessing, but also a constraint)

- There is political rather then technical issue: how to perform massive data analysis within the VObs?

# THANK YOU
# FOR YOUR ATTENTION !

**fabio.pasian@inaf.it**
**riccardo.smareglia@inaf.it**