# Data validation beyond Big Data

Edwin A. Valentijn

Kapteyn Astronomical Institute

6 June 2018  VST in the era of large sky surveys-  Napoli

# STORY LINES

- processing/archiving/distribution:

     - AstroWISE- KiDs - Ou-Ext – Euclid


- data validation:

     - lineage  -  OU-Ext  - Euclid-  Facts and Fakes
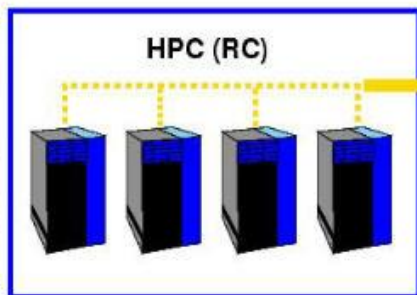

Sequence of hypes:

GRID -  Big Data  - Machine learning ->  data validation

# The Datacentric approach
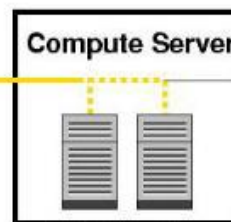## local networks and distributed



2003
RUG-CIT

OmegaCEN & HPC

© 2003 Astro–Wise

# Astro-WISE – Data federations

## Distributed Information Systems - handling surveys
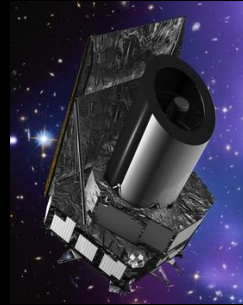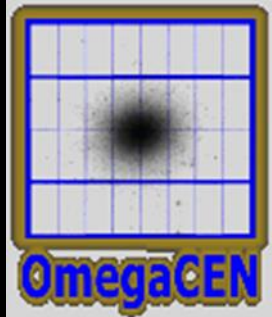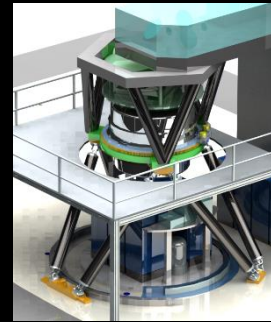## since 2003  -   it works
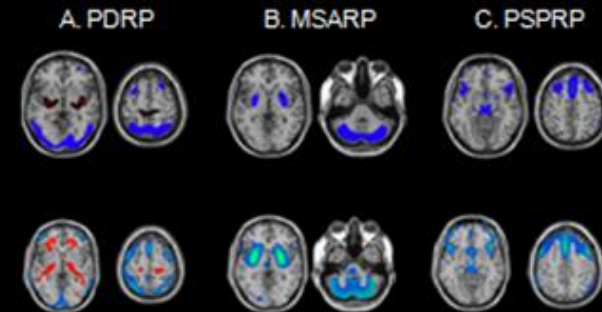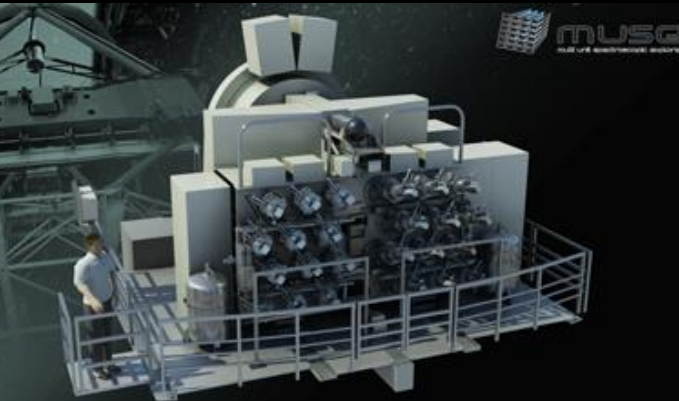### OmegaCEN@Kapteyn datacenter ~15-20 fte

KiDS                    -  ESO – OmegaCAM@VST
MUSE                  -  ESO - VLT
Lofar - LTA          -  Astron
Glimps -  AI Handwritten text  – Lifelines DNA
Target Holding
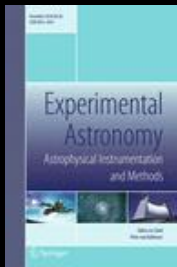
-> Euclid           -  ESA
-> Micado          -  ESO - ELT

# all published

http://www.astro-wise.org     Manuals & tutorials

http://www.rug.nl/target     Target Consortium

Experimental Astronomy - Vol. 35, 2013

All papers are online

## Target and (Astro-)WISE technologies
## Data federations and its applications

E. A. Valentijn[1], K. Begeman[1], A. Belikov[1], D. R. Boxhoorn[1],
J. Brinchmann[2], J. McFarland[1], H. Holties[3], K. H. Kuijken[2],
G. Verdoes Kleijn[1], W-J. Vriend[1], O. R. Williams[4],
J. B. T. M. Roerdink[5], L. R. B. Schomaker[6], M. A. Swertz[7],
A. Tsyganov[4] and G. J. W. van Dijk[8]

[1]Kapteyn Astronomical Institute, University of Groningen,
email: valentyn@astro.rug.nl
[2]Leiden Observatory, Leiden University
[3]ASTRON, Dwingeloo
[4]Center for Information Technology, University of Groningen
[5]Johann Bernoulli Institute, University of Groningen
[6]ALICE, University of Groningen
[7]University Medical Center Groningen, University of Groningen
[8]Target Holding, Groningen

Astroinformatics 2016
IAU symposium 325
Datafederations
Valentijn et al. 2017

# Links as workhorse in data federations

- Distributed  Information Systems
    - Users, computers, storage
- Processing   and Quality control
- Reproducable    ( re-processing)

2018: Open Science  - **FAIR** principles

**F**indable  **A**ccessable  **I**nteroperable **R**eproducable

# The universe as a spreadsheet
## Target Diagram/Data lineage /backward chaining ++ programming - dependencies

Entangled I/O

raw=RawScienceFrame processed=ReducedScienceFrame

process.astro-wise.org/Process

Euclid conference Bonn

OCam   OCen   EV   Wiki   Euclid   Printers   Mail   Routers   AA   TARGET   Nieuwe map   Infoversum   Target   NOS   GIC   NOS   DVHN   Astrop

**Astro-WISE Homepage**

**Target Processor**

**Contact**
Willem-Jan Vriend

**DB User**
awevalentyn

**Help**
Getting Started

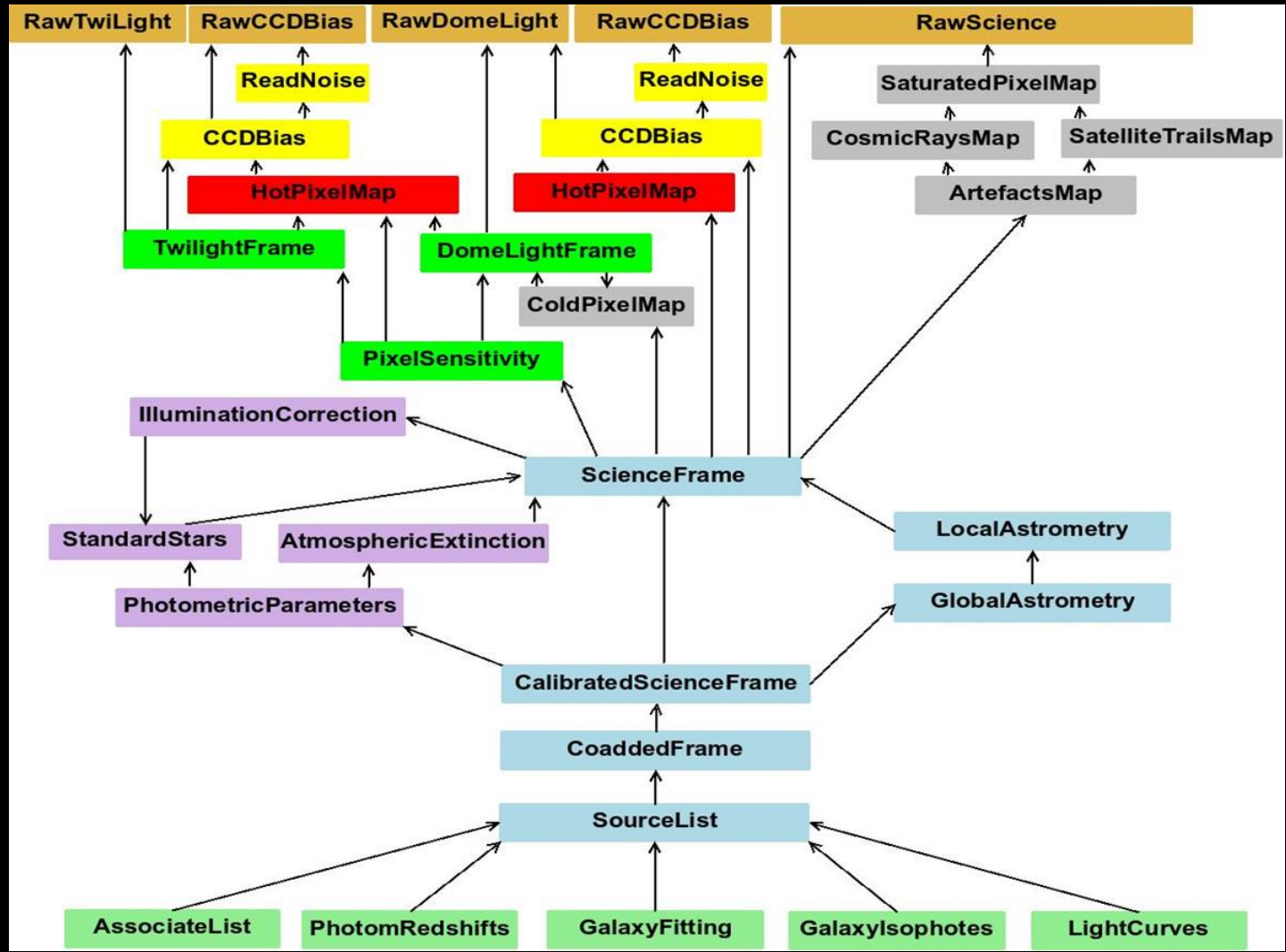**Project**
KIDS

**Instrument**
OMEGACAM

**State**
1. Preselect Target
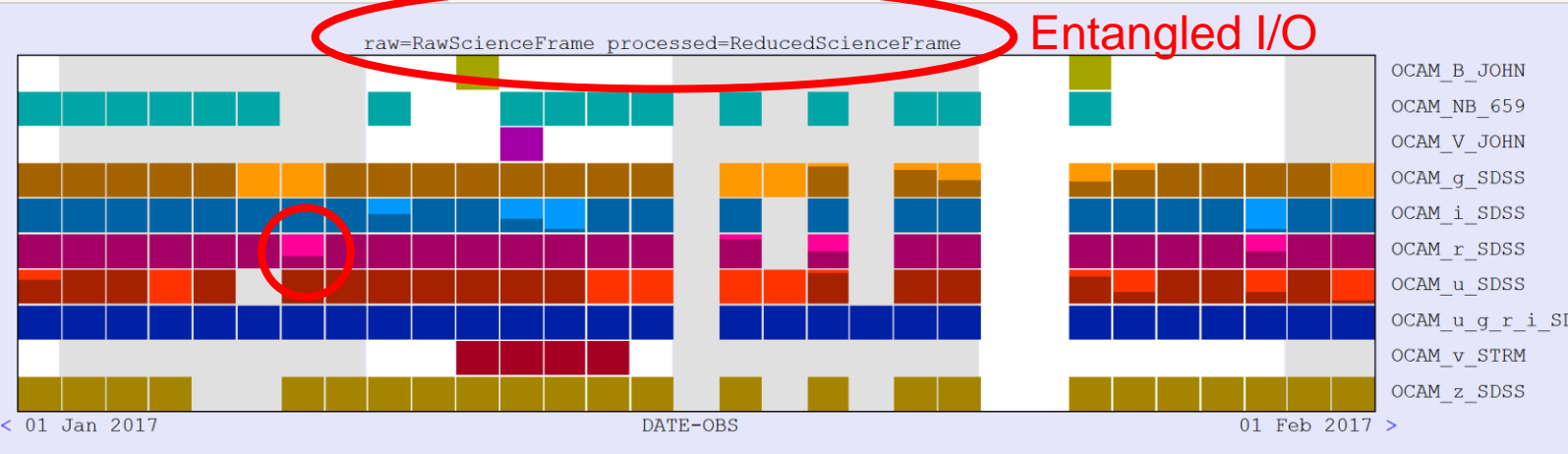2. Specify Target
3. Select Target(s)
4. Process or Query

**Options**
Preferences
Process Parameters
Upload Code
Job overview

OCAM_B_JOHN
OCAM_NB_659
OCAM_V_JOHN
OCAM_g_SDSS
OCAM_i_SDSS
OCAM_r_SDSS
OCAM_u_SDSS
OCAM_u_g_r_i_SD
OCAM_v_STRM
OCAM_z_SDSS

< 01 Jan 2017     DATE-OBS     01 Feb 2017 >

## Specify Target

Specify a period and click show. For the selected period all available observations will be shown in the above view. Each block corresponds to one or a set of observations with a specific filter or observing block. Click on a block to get an overview of the possible targets. You can also use the extended query form.

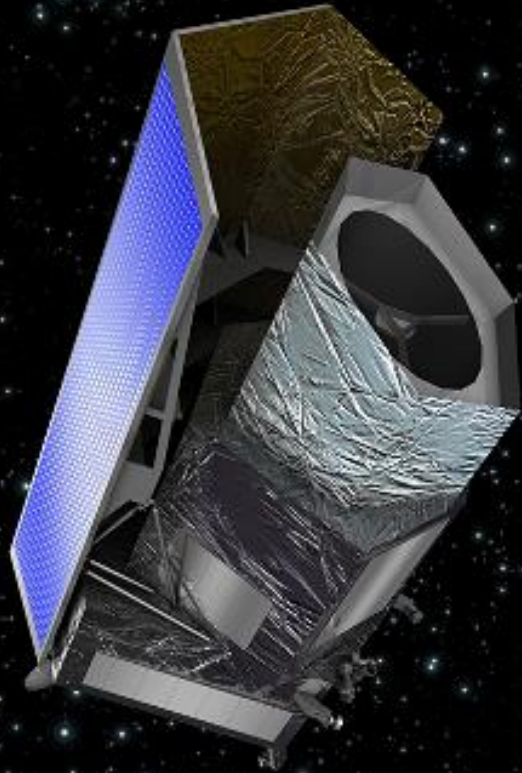### Period Selection (DATE-OBS)

| Year | Quarter | Month | Week |
|------|---------|-------|------|
| 2017 | <none> | 1 jan | <none> |

### Optional Settings

| Name | Value |
|------|-------|
| Filter | <none> |
| Group by | ⦿ Filter   ◯ Observing Block   ◯ Template |
| Filtering | ☑ Flagged data   ☐ Project only |

Show

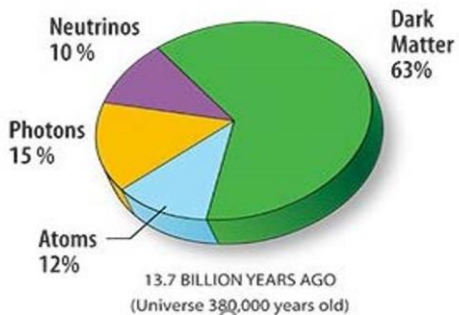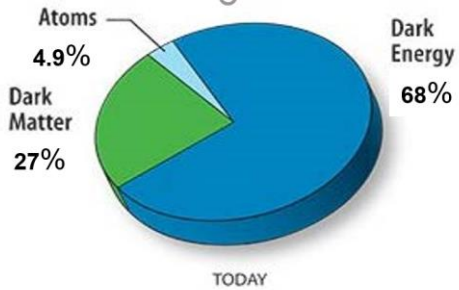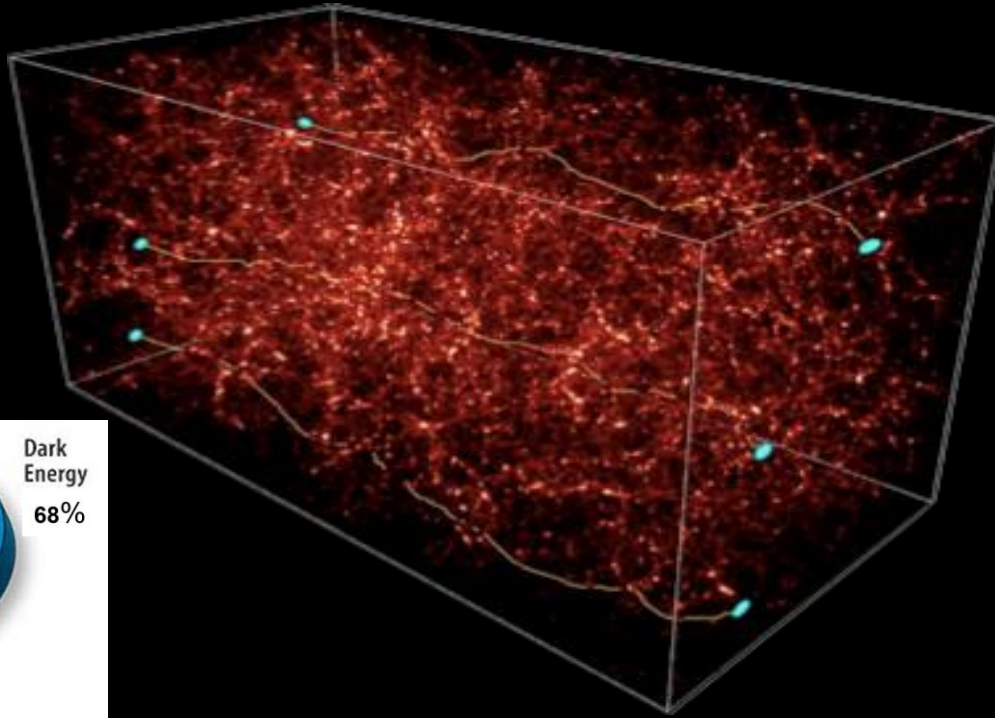| | raw | | processed | | | |
|---|---|---|---|---|---|---|
| | 192 | | 0 | OCAM_B_JOHN | JohnsonB |
| | 9184 | | 0 | OCAM_NB_659 | UnknownNB659 |
| | 32 | | 0 | OCAM_V_JOHN | JohnsonV |
| | 6624 | | 2400 | OCAM_g_SDSS | SloanG |
| | 10624 | | 2048 | OCAM_i_SDSS | SloanI |
| | 11008 | | 640 | OCAM_r_SDSS | SloanR |
| | 7808 | | 2595 | OCAM_u_SDSS | SloanU |
| | 2976 | | 0 | OCAM_u_g_r_i_SDSS | SloanUGR |
| | 128 | | 0 | OCAM_v_STRM | StromgrenV |
| | 1376 | | 0 | OCAM_z_SDSS | SloanZ |

# Euclid



ESA   launch in May 2021

Euclid Archive System (EAS)

- data centric information system

- many of the WISE concepts

- prototype uses Astro-WISE

- db hosted in the Euclid SDC-NL in Groningen

# Weak gravitational lensing as probe of dark matter
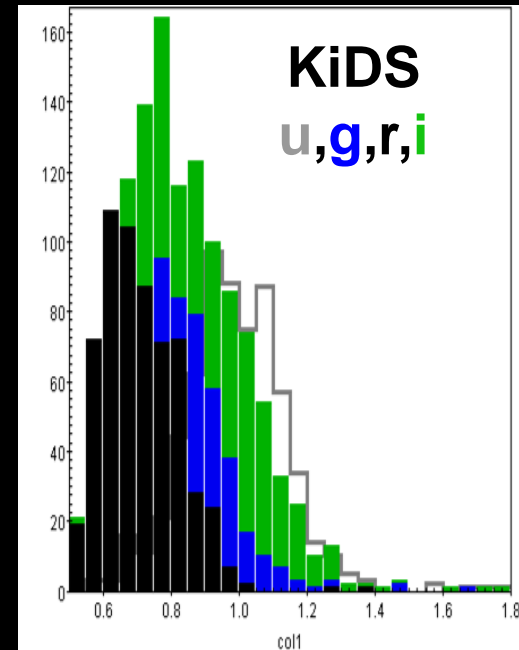


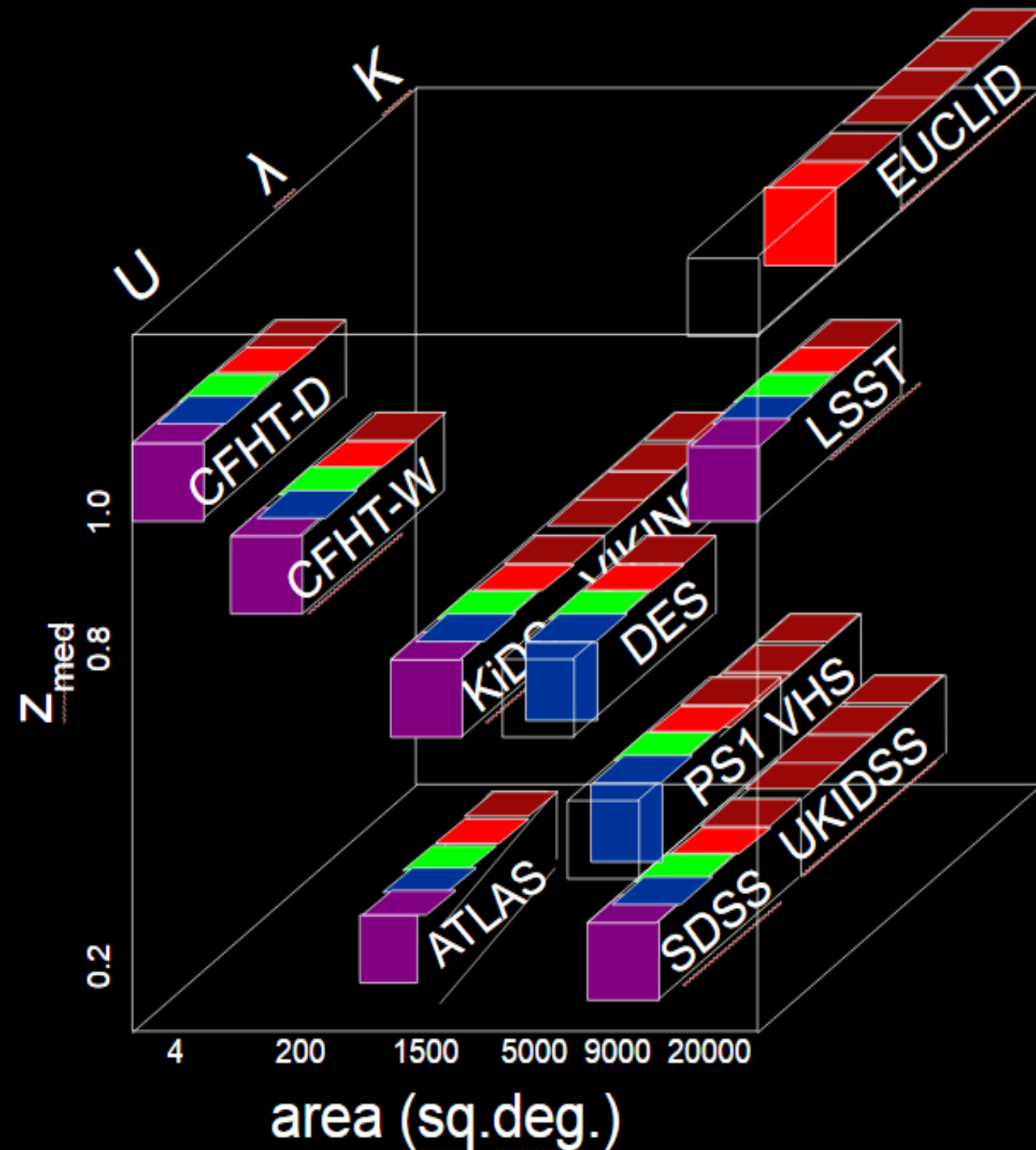KiDS:          < 100 $10^6$ redshifts
EUCLID:        1.5 $10^9$ redshifts  - phot- z
                          Ground based data – OU-Ext
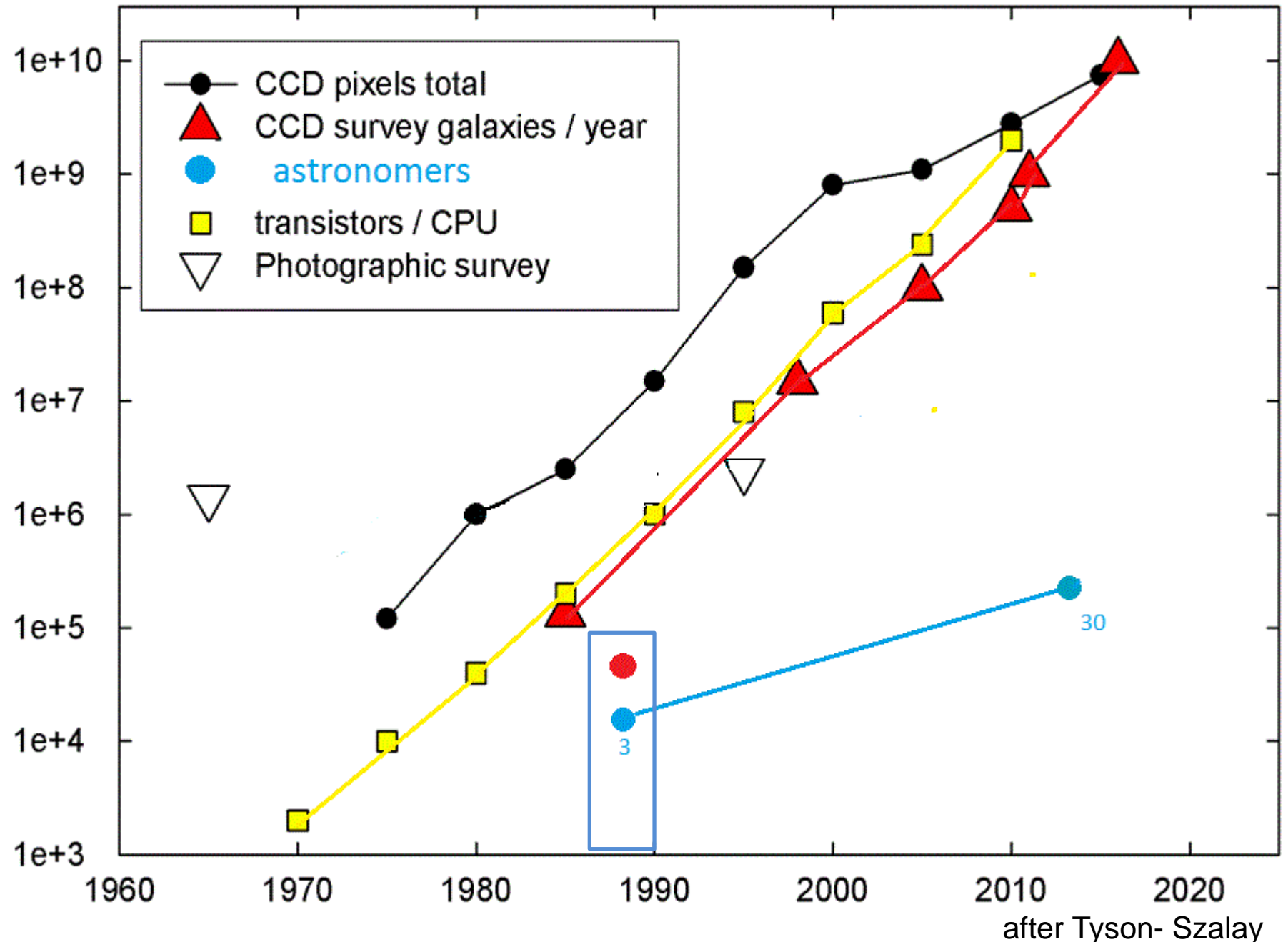Every galaxy has its own 4 PSFs
QC- bias – re-processing
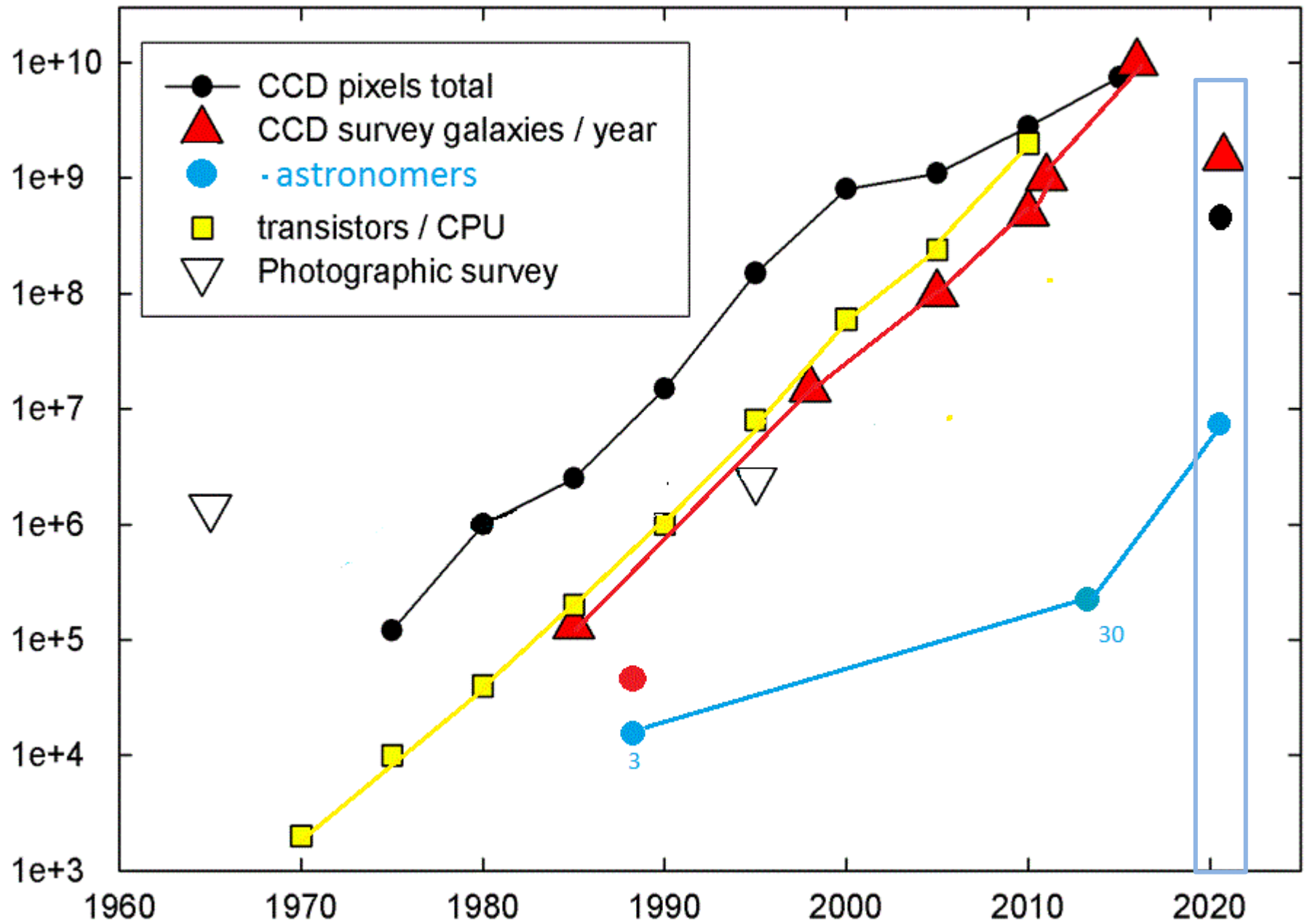
KiDS/VIKING

KiDS
u,**g**,r,i

Seeing (")

# Trends in Optical Astronomy Survey Data

Legend:
- ● CCD pixels total
- ▲ CCD survey galaxies / year
- ● astronomers
- ■ transistors / CPU
- ▽ Photographic survey

after Tyson- Szalay

# Trends in Optical Astronomy Survey Data



**Legend:**
- ● CCD pixels total
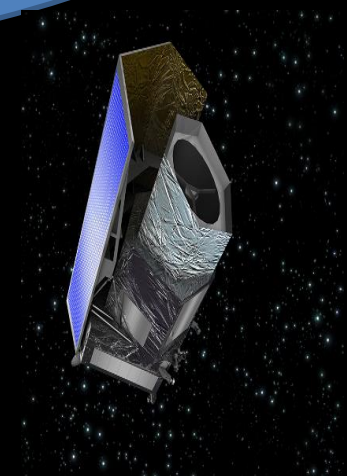- ▲ CCD survey galaxies / year
- ● astronomers
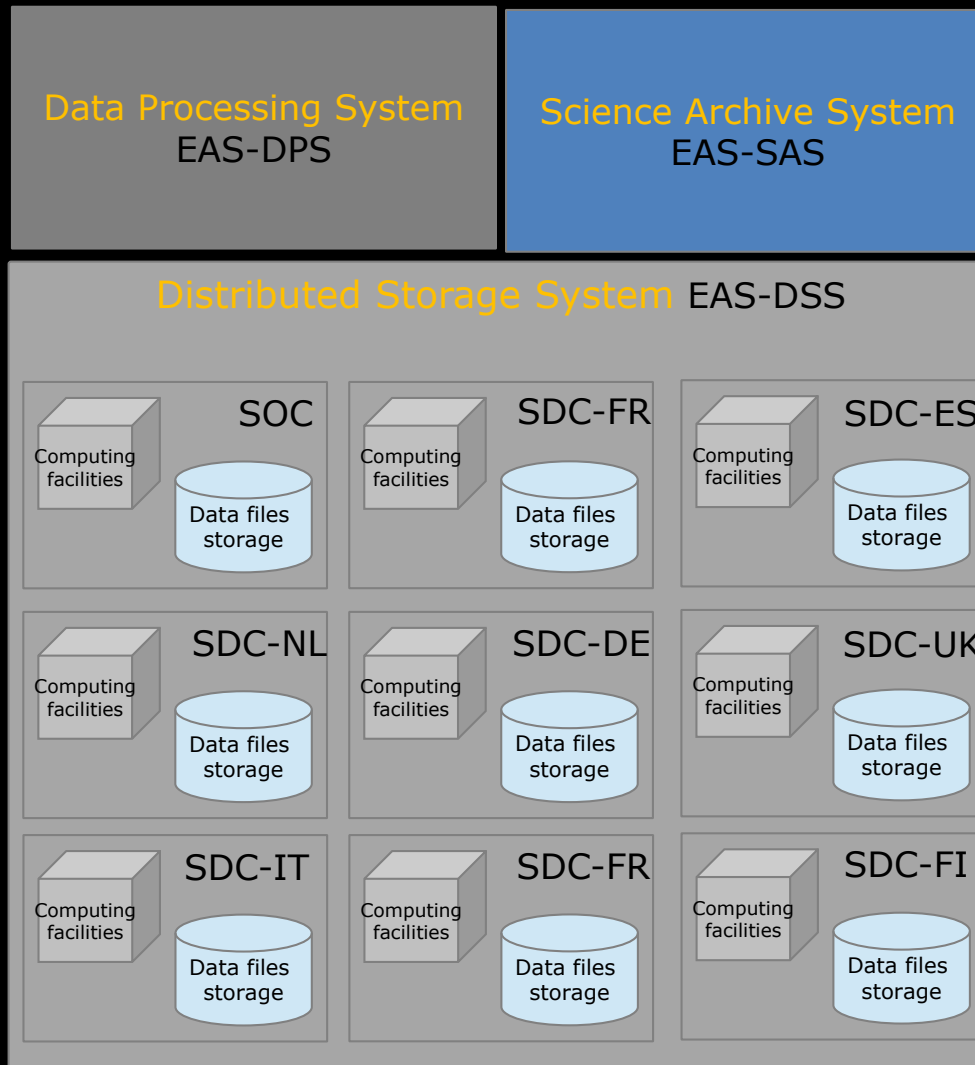- ■ transistors / CPU
- ▽ Photographic survey

# Distributed communities acces-proces-calibrate-analyse publish
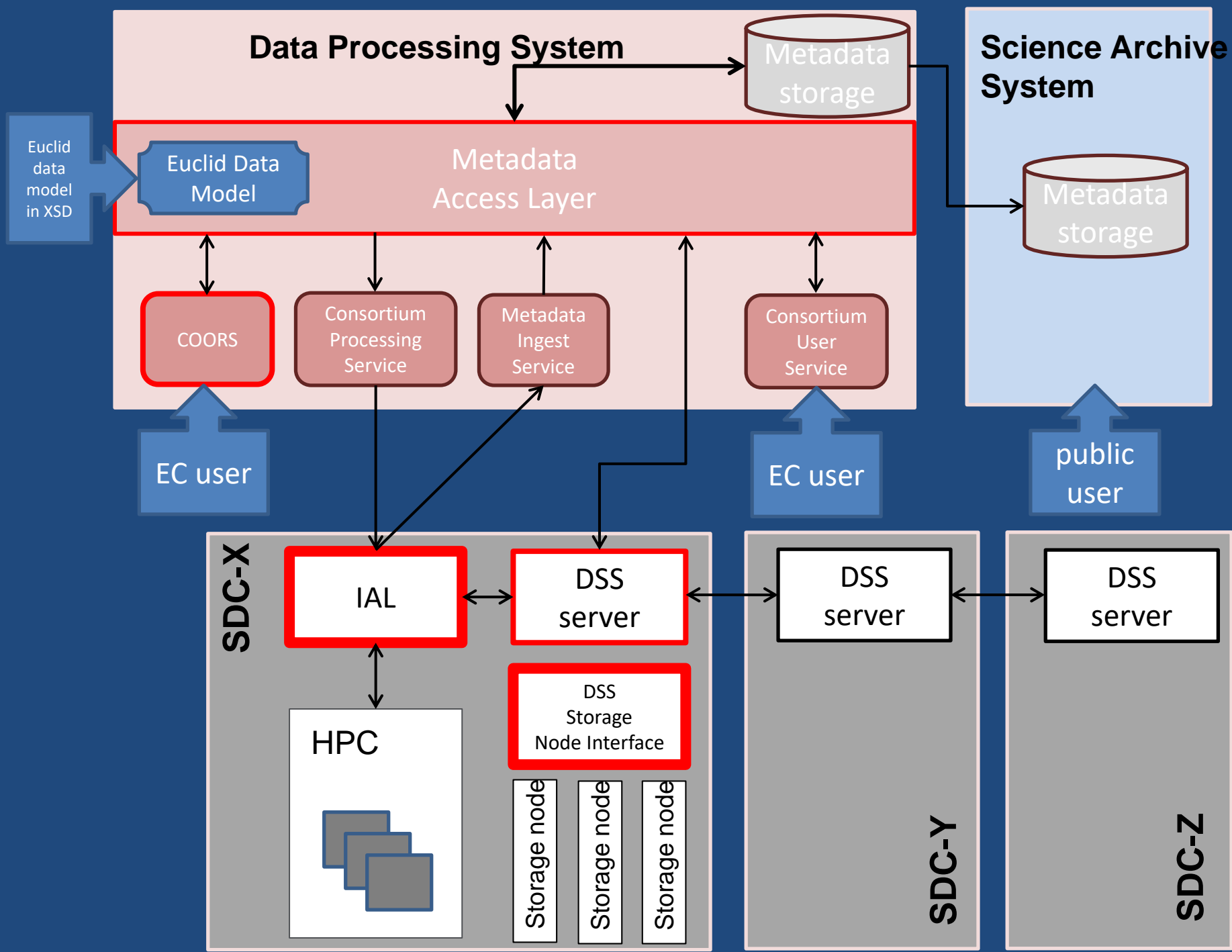
Euclid:

- o 1500 registered members and growing
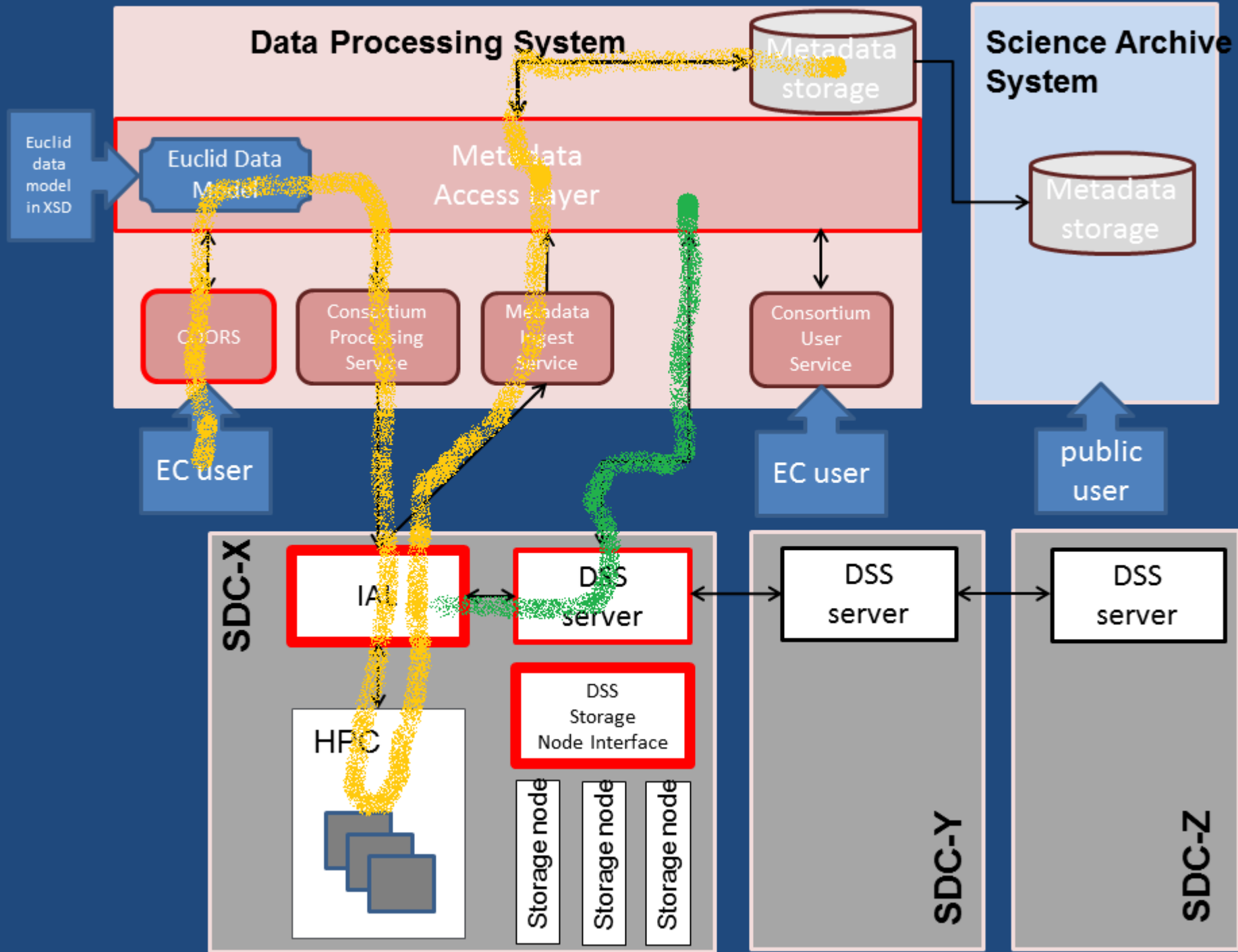- o 200 laboratories/departments
- o 16 countries contributing
- o NASA/US: provides the IR detectors.

# Euclid Archive system – EAS – lay out

# Euclid-EXT: massive pixel volumes - distributed archives
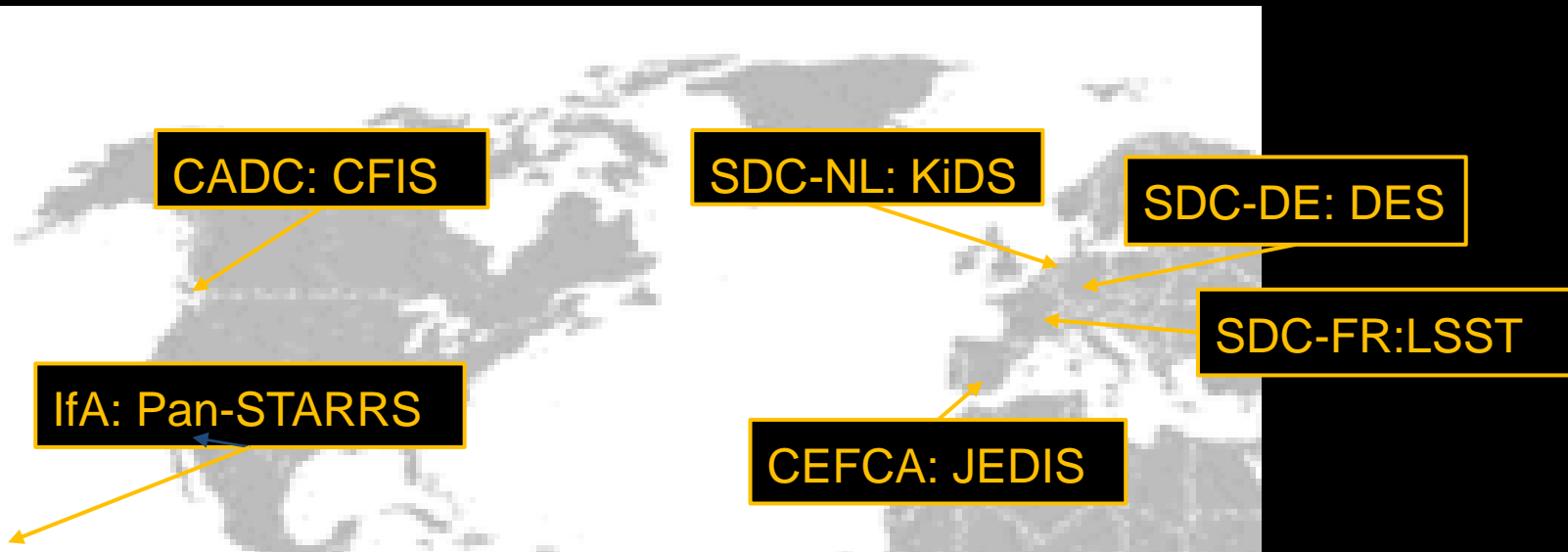


CADC: CFIS

SDC-NL: KiDS

SDC-DE: DES

SDC-FR:LSST

IfA: Pan-STARRS

CEFCA: JEDIS

**~6E5 EXPOSURES RAW DATA (TBYTE)**

- Euclid, 219
- DES, 156
- CFIS, 62.4
- JEDIS, 31.2
- KiDSVIKING, 40
- LSST-Euclid, 2152
- PanSTARRS-Euclid, 1076
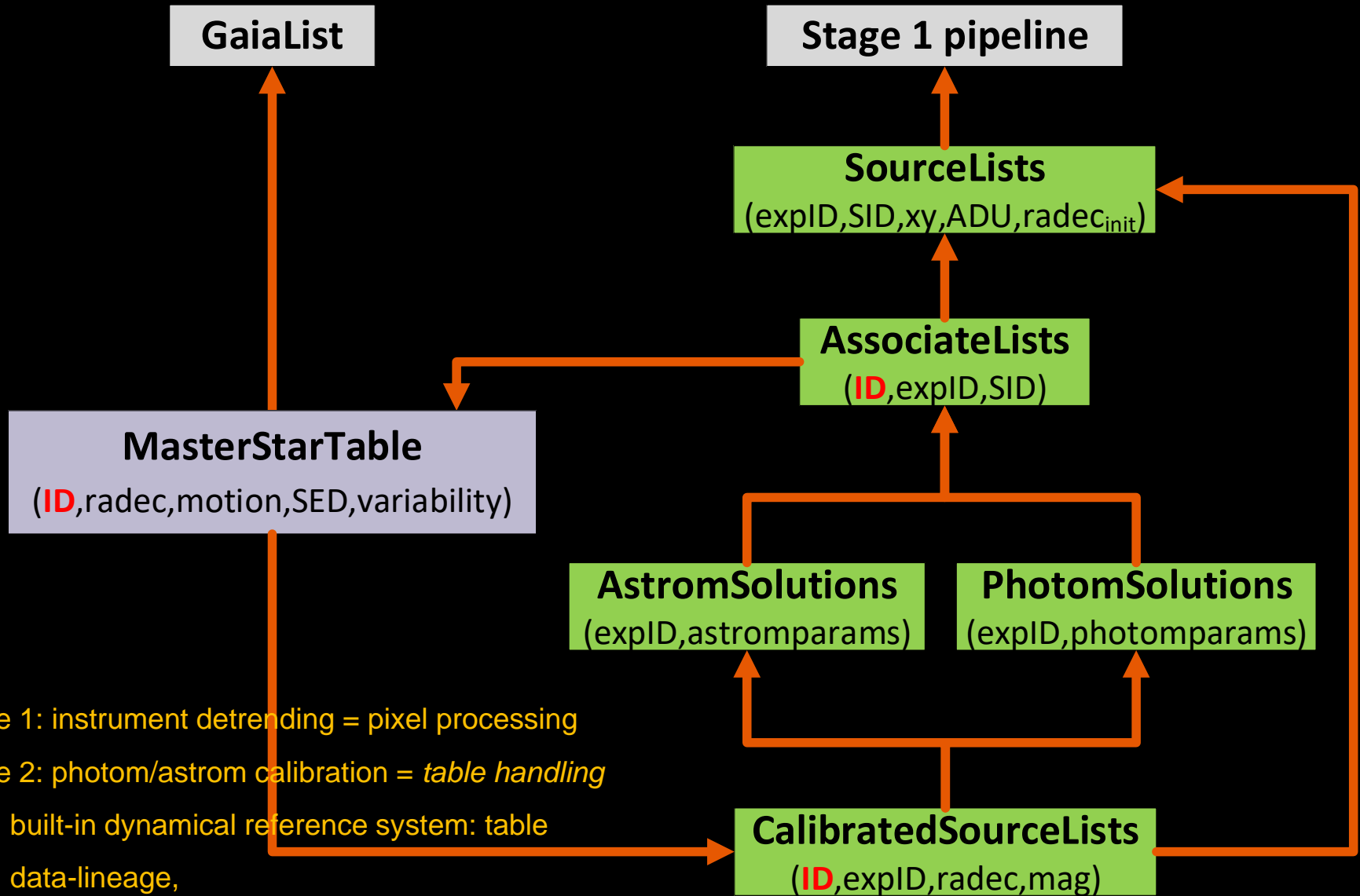
# From KiDS to Euclid-EXT



Euclidization
Changing reference systems
Astrometry- photometry

# Target diagram ( ++ dependencies) for OU-EXT – Euclid external data - stage 2- dynamic Euclidization

**GaiaList**

**Stage 1 pipeline**

**SourceLists**
$(\text{expID},\text{SID},\text{xy},\text{ADU},\text{radec}_{\text{init}})$

**AssociateLists**
(**ID**,expID,SID)

**MasterStarTable**
(**ID**,radec,motion,SED,variability)

**AstromSolutions**
(expID,astromparams)

**PhotomSolutions**
(expID,photomparams)

**CalibratedSourceLists**
(**ID**,expID,radec,mag)

Stage 1: instrument detrending = pixel processing

Stage 2: photom/astrom calibration = *table handling*

- built-in dynamical reference system: table
  data-lineage,
- QC, re-processing

# Beyond Big Data

- QC and re-processing  – Kids Euclid    **FAIR**
- OU EXT >  Billion – dynamic tables

All techniques  go back to the source

Scientists  and journalists- > Fact and Fakes

Structured data   and unstructured data

# TARGET Fieldlab

**Fact or Fake**
- News items tracking
- Open Science Applications
- Data lineage

**Sensor Grids**
- Timeseries : trend prediction
- Open Seismic Sensor Grid
- Wearables

**VR Valley**
- $360^0$ imaging
- VR editing Platform
- Social applications
- Medical applications

## Proeftuin gebruikers

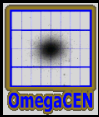Demo project
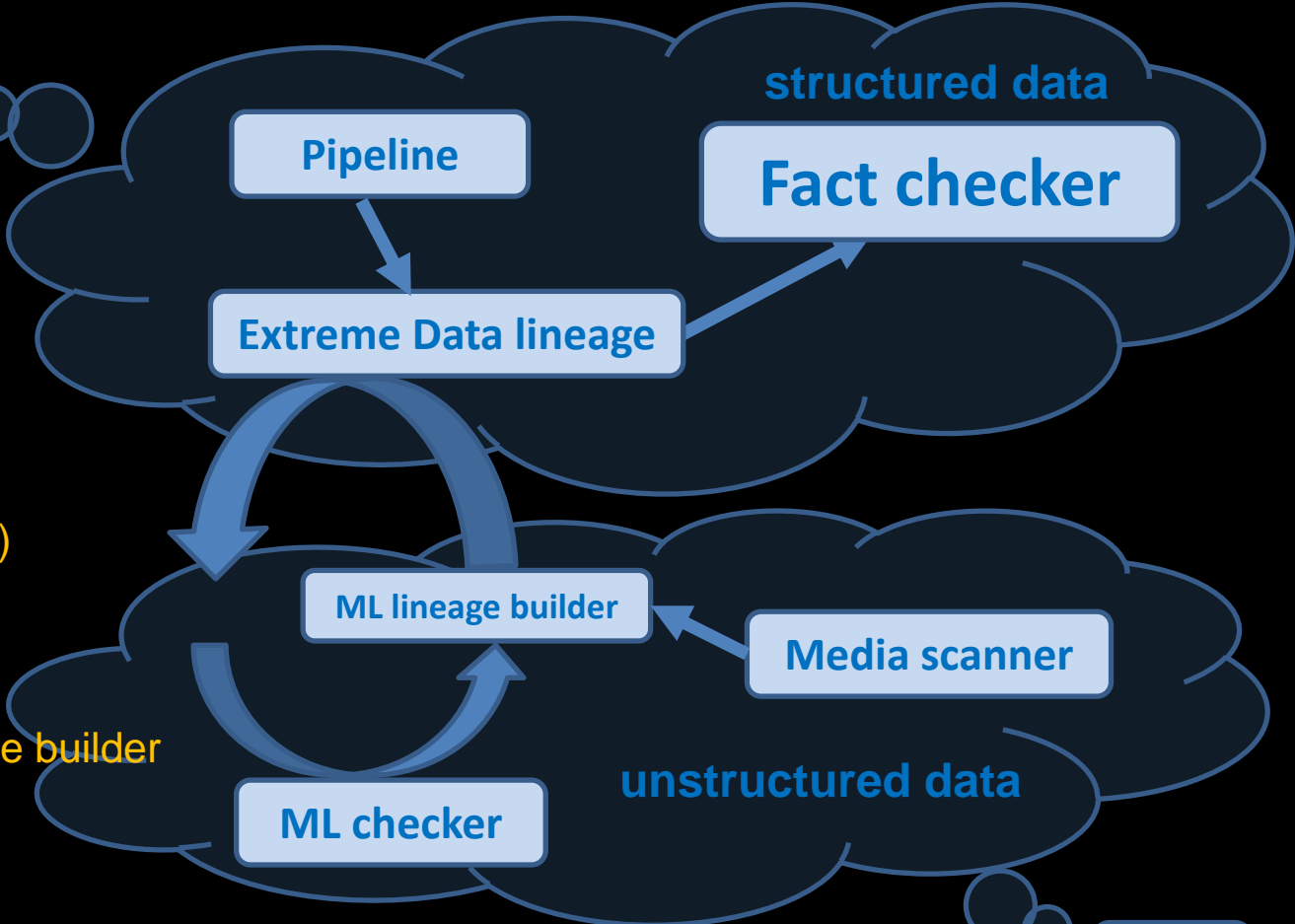(Crowdy News
TRAIN AIAAS BV)

Demo project
(Target Holding)

Tender

Demo project
(Horus VR
Yellowbird)

Andere & nieuwe klanten

# conclusions

Next level is all about Data validation

- check ML

- QC

- systematics  in data sets

- OU-ext  dynamic Euclidization

- unstructured data:  ML + lineage

Almost all about going back to the source

Facts and Fakes