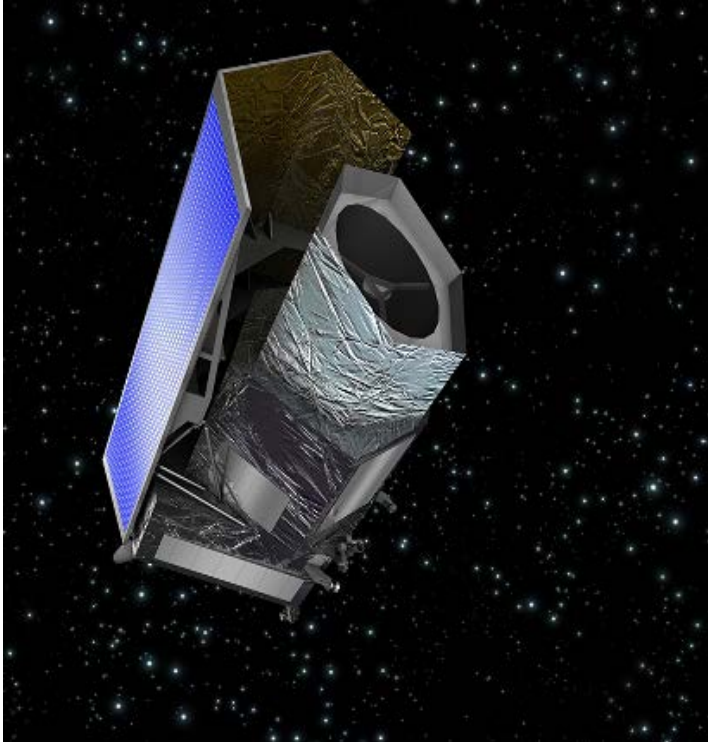# The Euclid Archive Data Processing and Storage Systems: a distributed infrastructure for Euclid

Rees Williams

University of Groningen

on behalf of the Euclid Archive Development Team

# Euclid



- ➢ ESA mission to map dark matter & dark energy
- ➢ Launch on Soyuz in 2022

- ➢ Wide field survey
- ➢ 1.2m telescope
- ➢ 3 Optical & IR instruments

- ➢ 15 countries
- ➢ 200+ institutes
- ➢ 2000+ consortium members

## Unprecedented data volumes for astronomy satellite mission

1 Peta byte images from Euclid

10 Peta byte images from earth observatories
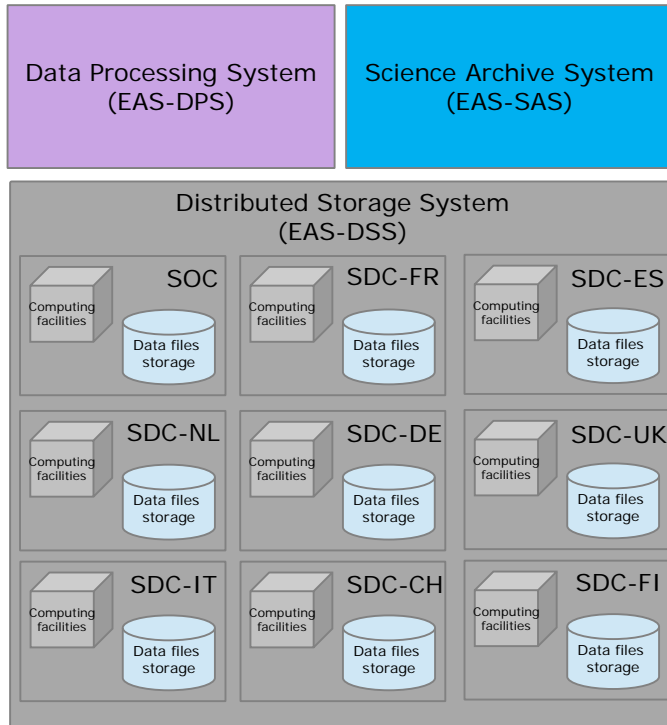
Massive simulation efforts required

### Key features

➢ Very strict requirements on quality and calibration
➢ heavy (re)processing needed from raw data to science products
➢ mission could generate 25 Pbytes/year
➢ 10 billion objects in catalogue
➢ 10 national science data centres (SDC)
➢ distributed data processing and storage

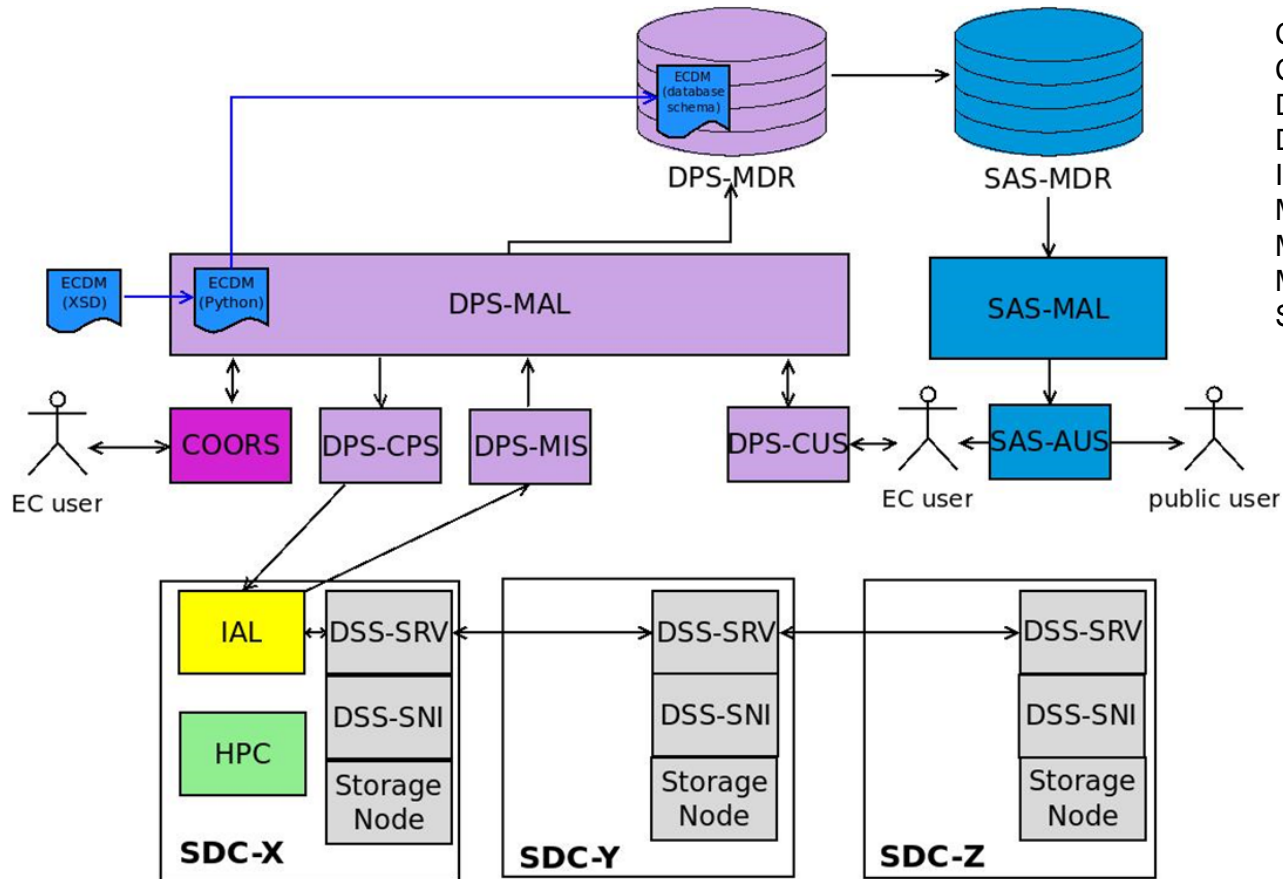# The Euclid Archive System: Data Centric Approach

- The Euclid Archive System (EAS) has a central role in Euclid data processing
- Traceability and data lineage are strong requirements
  - EAS acts as interface between all ground system components
  - EAS data distributed across the 10 national data centres
  - EAS metadata is centrally stored
  - EAS metadata contains all information except images and spectra
  - EAS stores dependencies of data products

- Builds on experience with the Astro-WISE information system used by
  - OmegaCAM (KiDS), MUSE

# Euclid Archive System: Components



- EAS-DPS: Functionality needed to support operations of the SGS, most notably Data Processing. Implements the complex Euclid Common Data Model (poster by T. Nutma).

- EAS-DSS: Functionality to manage and transfer data stored across many locations.

- EAS-SAS: Functionality needed to support the scientific community (talk by S. Nieto).

COORS- Coordination & orchestration
CPS – Consortium processing service
DPS – Data Processing System
DSS - Distributed Storage System
IAL - Infrastructure abstraction layer
MAL- Metadata access layer
MDR - Metadata Repository
MIS – Metadata Ingestion service
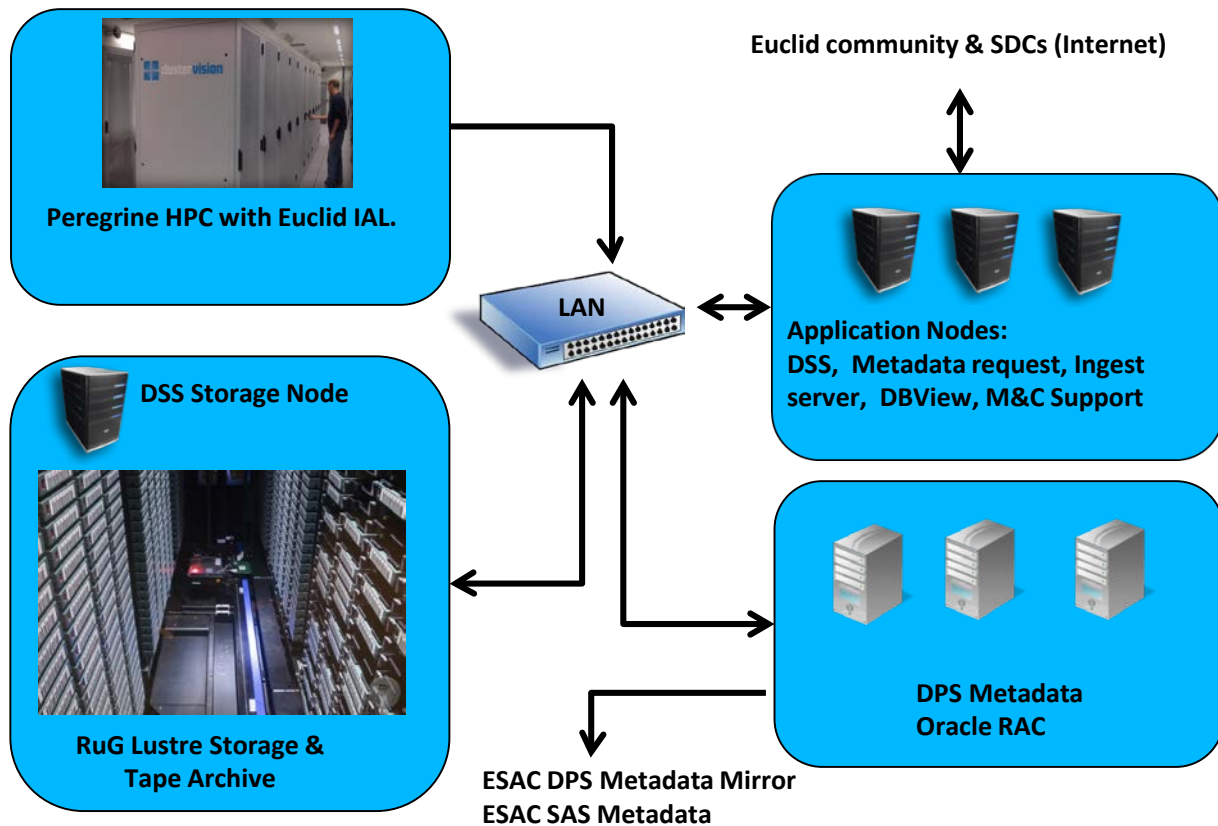SAS – Science Archive System

# Distributed Storage System Servers: I

- File location stored in central metadata system - provides a global "file system"
  - Process coordination system can run jobs at cluster closest to data

- DSS server installed at each SDC can access data to other SDCs.
  - Users can run jobs locally without knowing the location of data

- Common code base with Astro-WISE, MUSE-WISE, NOVA-LTA
  - Decreases maintenance overhead

- Simple interface:
  - store, retrieve, copy, delete, check, make_local

- Extended interface: cut-out services

# Distributed Storage System Servers : II

- No constraint on storage systems at data centres

- DSS server currently supports the access to data using the following protocols
  - Posix file system
  - sftp
  - https
  - gridftp,
  - Astro-Wise dataserver
  - iRods
  - XRootD
  - Openstack

# Euclid infrastructure at SDC-NL



**Peregrine HPC with Euclid IAL.**

**Euclid community & SDCs (Internet)**

DPS – Data Processing System
DSS - Distributed Storage System
IAL - Infrastructure abstraction layer
SAS – Science Archive System

**LAN**

**Application Nodes:**
**DSS,  Metadata request, Ingest**
**server,  DBView, M&C Support**

**DSS Storage Node**

**RuG Lustre Storage &**
**Tape Archive**

**ESAC DPS Metadata Mirror**
**ESAC SAS Metadata**

**DPS Metadata**
**Oracle RAC**

# EAS Data Processing Services

EAS-DPS services:
- consortium processing service (CPS)
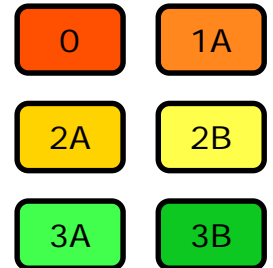- consortium user service (CUS)

**CPS:**
- IAL interface
- COORS interface
- COORS plugin
- archive function service
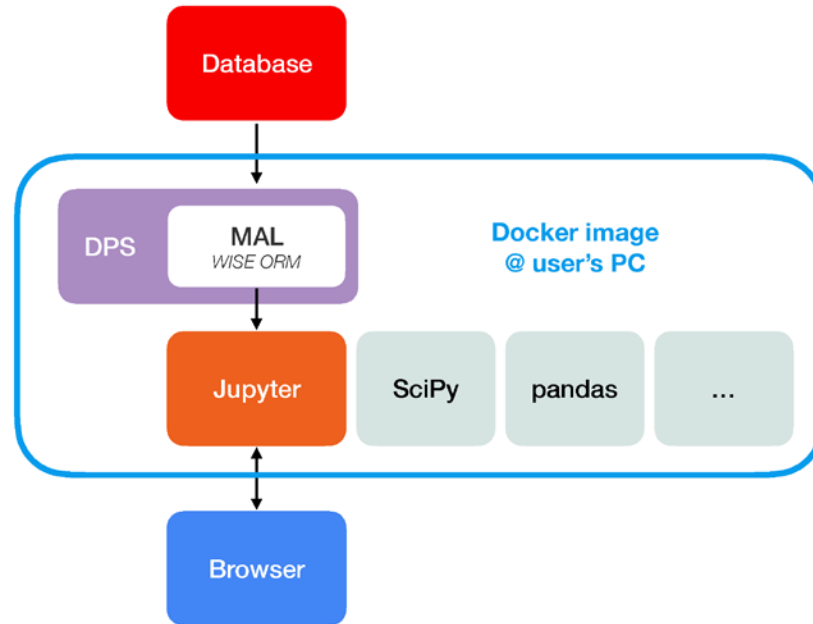- data distribution service

**CUS:**
- metadata explorer (MEX, novice user)
- parametric plot (novice user)
- DBview (expert user)
- metadata access layer (MAL, expert user)
- calibration service
- quality view service
- data lineage viewer

Maturity Levels
see: MLs definition
(restricted access)

| 0 | 1A |
| 2A | 2B |
| 3A | 3B |

Full data model review
Data products/definitions grouped

SC456 software stack

# Euclid Science Challenges

- EAS is currently supporting the so called Euclid Science Challenge #4, #5 & #6.

  - Distributed processing across all 10 data centres
  - Two small wide field regions (7 observations)
  - One large wide field region (38 observations c.f. 40,000 observations in total)
  - Challenge due to finish Nov 2019

- Issues found

  - Management of data processing too manpower intensive (needs more tools)
  - Spatial queries need to be faster

# Conclusions

- Basic functionality for the Euclid archive is in place.

- WISE technology shown to scale to Euclid data scales
  - datacentric distributed processing and storage

- Performance tests show minimum levels reached
  - Except for complex spatial queries

- More work on user interface required for
  - Calibration scientists
  - Data quality investigations