



Euclid – Feedback on IVOA data models

JC Malapert, C. Dabin

IVOA

11-13 October 2019, Groningen

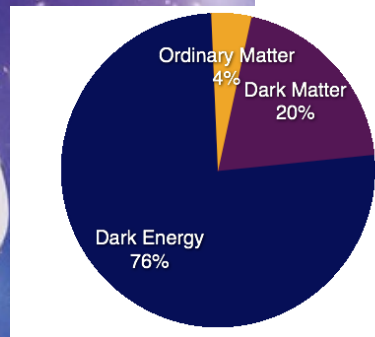
Overview

- **Context: The Euclid project**
 - In a nutshell
 - The distributed processing of the Science Ground Segment
- **The Euclid data model**
 - Presentation
 - Differences between IVOA schema and Euclid schema rules
 - Decisions about IVOA schema
 - Suggested improvements



The Euclid project

“Map the geometry and understand the nature of the **Dark Universe”**



M2 mission of ESA's Cosmic Vision Program

- Launched in 2022

1.2m space telescope + external data

- VIS instrument for visible imaging
- NISP instrument for near-IR and spectro imaging
- Ground-based surveys (KIDS, DES, LSST, J-PAS, UNIONS,...)

One of the largest organizations

- 16 countries
- 220 labs
- 1700 members

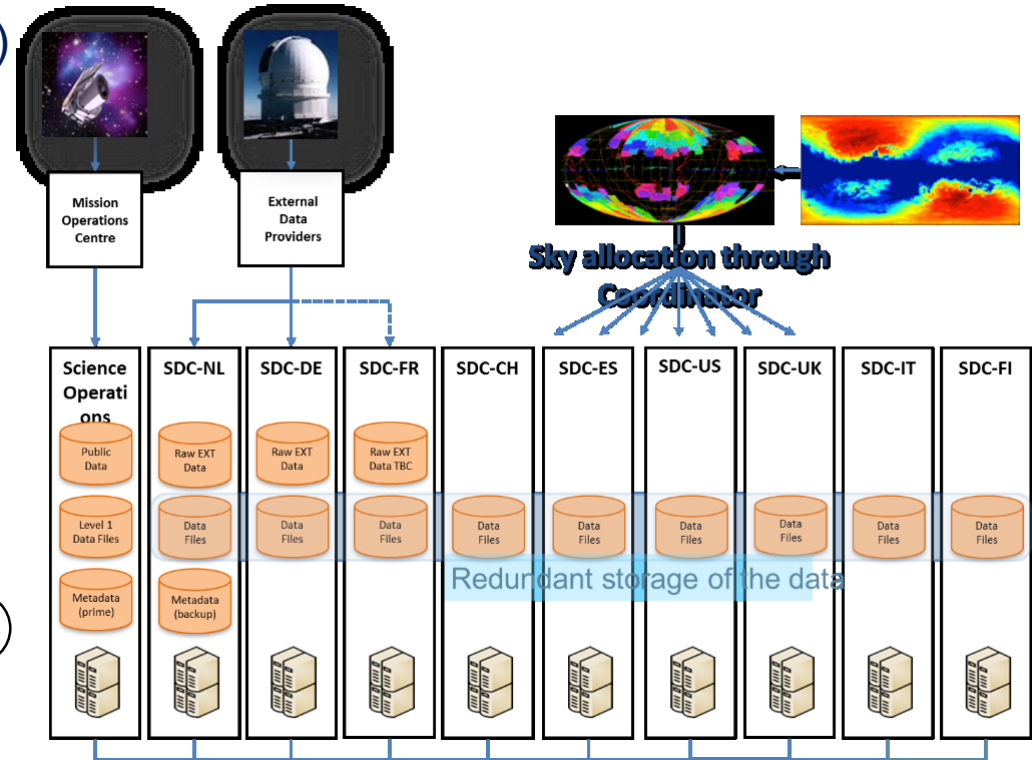
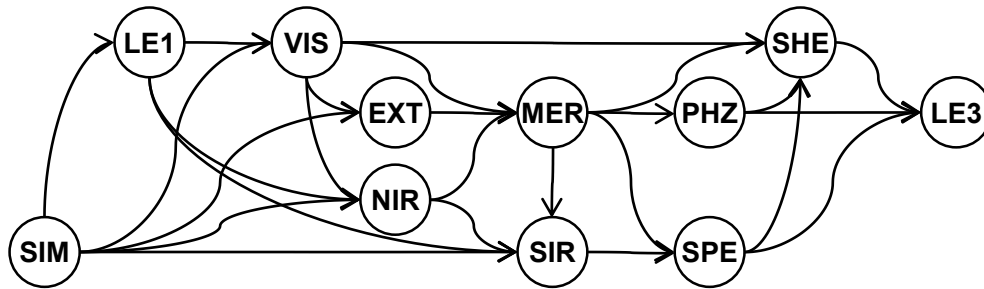
“Move the processing, not the data”

9 Science Data Centers (SDCs)

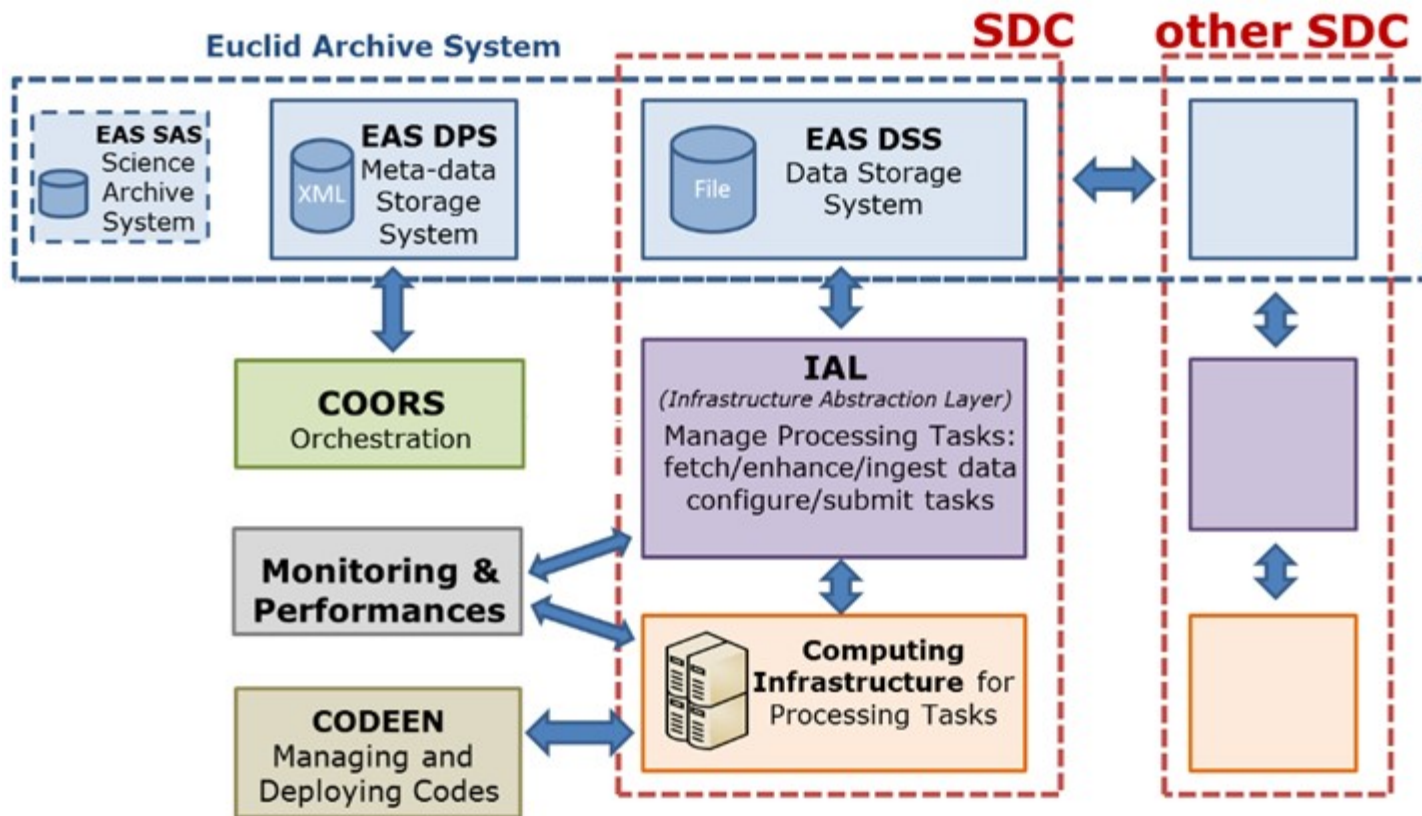
- Sky area (set of Adjacent Observations) allocation

11 Processing Functions (PFs)

- Developed by the Science Ground Segment (SGS) labs
- Each PF runs on all SDCs



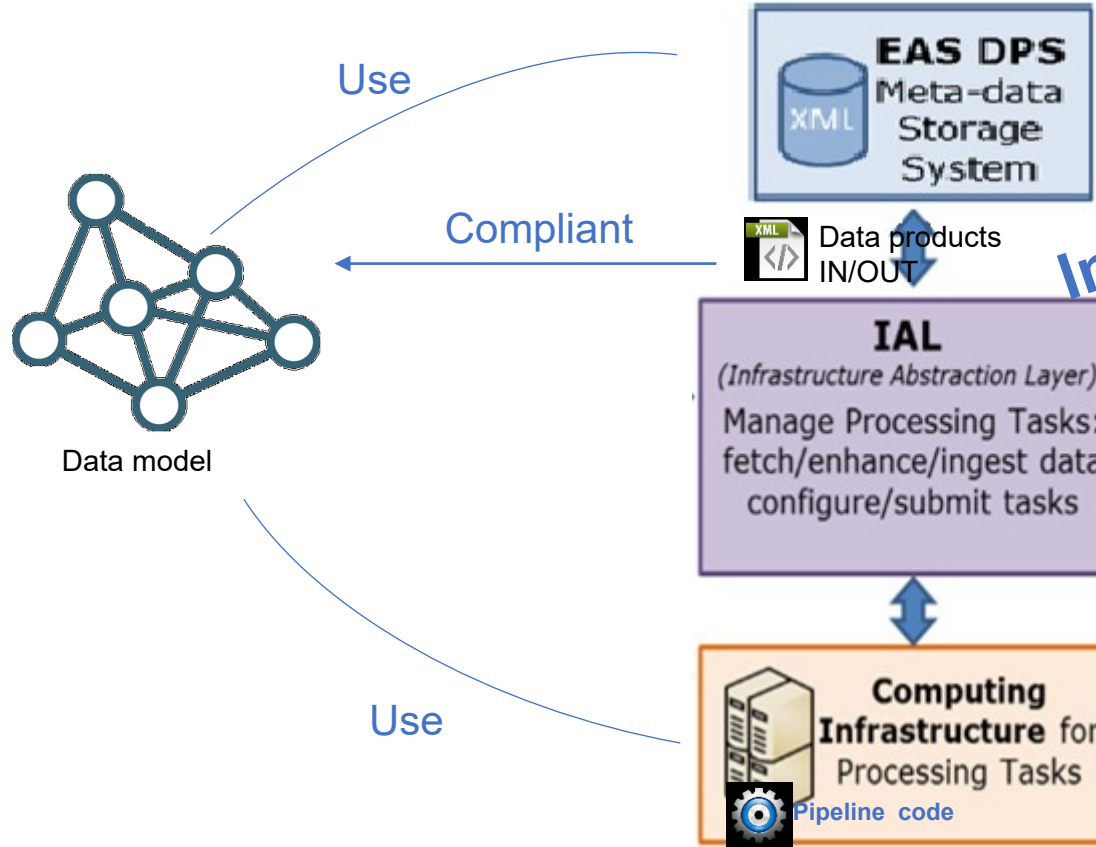
Euclid SGS components





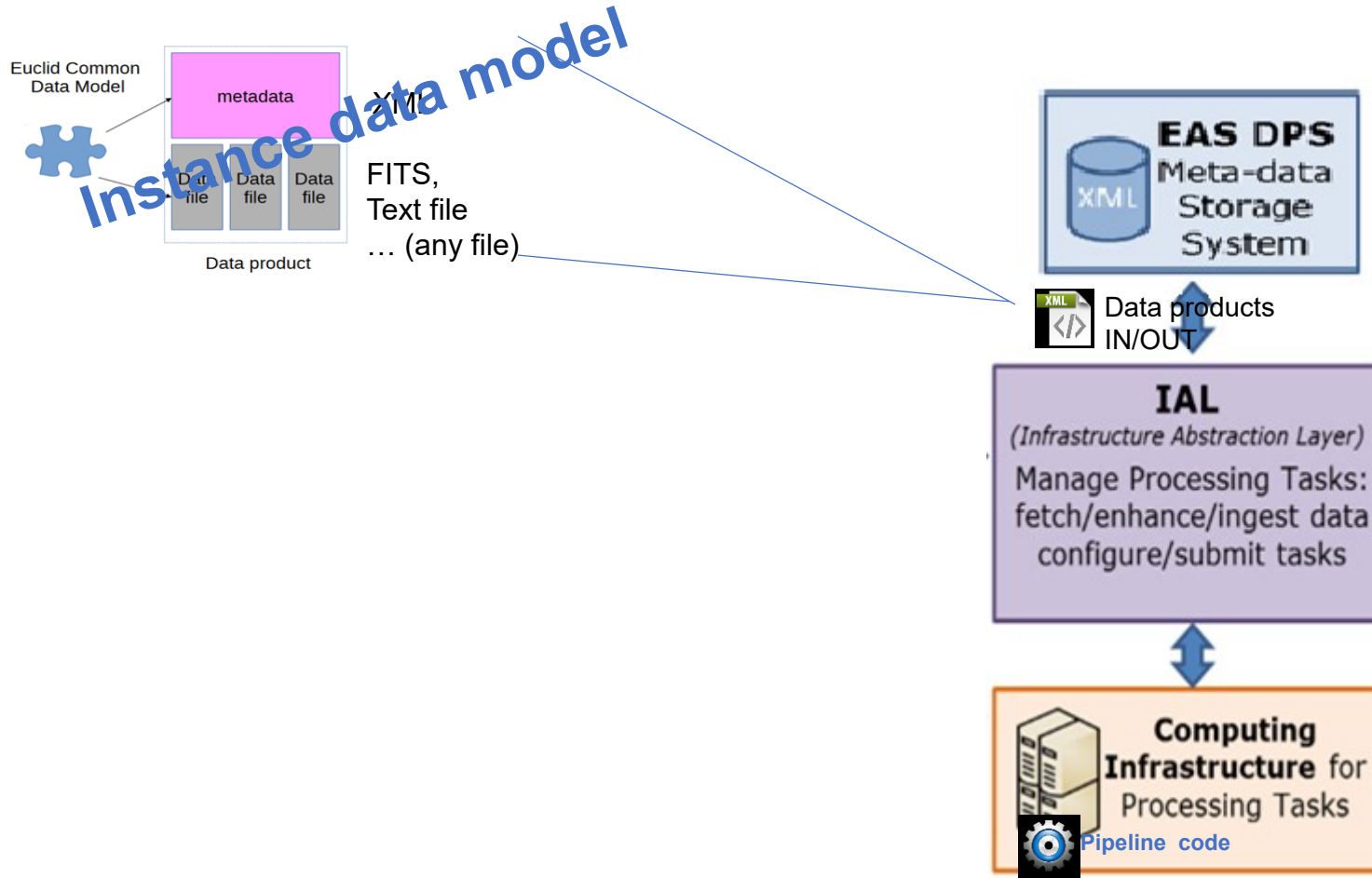
Euclid Data Model

What is the data model used for ?

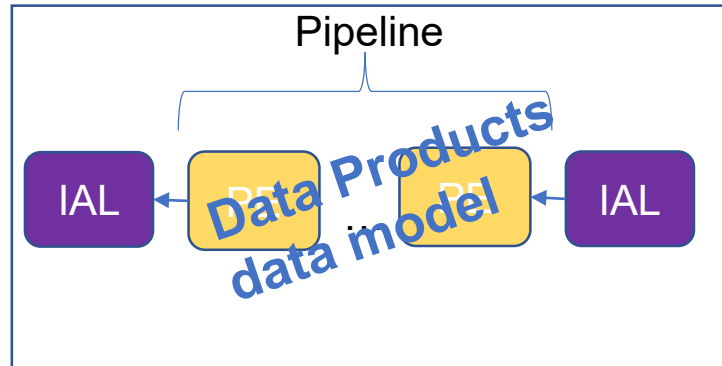
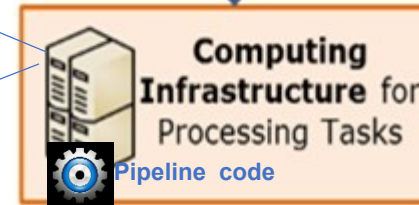
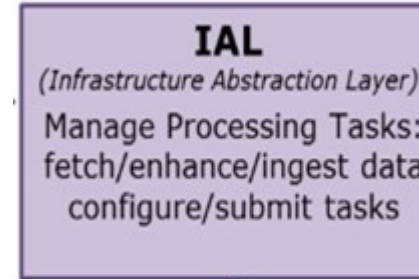


Interface data model

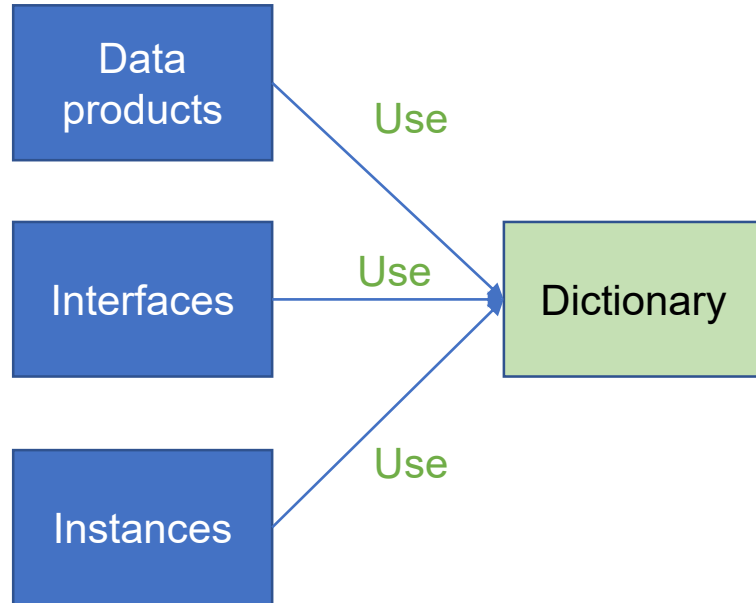
What is the data model used for ?



What is the data model used for ?



How is the data model structured ?



4 different themes (each theme has its own namespace) :

- **bas** : common definitions shared by everyone
- **ins** : instrument specific definitions
- **pro** : processing function specific definitions
- **sys** : system specific definitions (storage, processing, survey, production plan,...)

Dictionary : ins namespace

this branch (acronym of Instruments) is related to the specific data types describing the instruments characteristics

(ex: position of the observatory, filter used)

Dictionary : pro namespace

this branch (acronym of Processing) is related to the different processing functions of the overall Euclid processing pipeline. This branch is broken down into as many processing functions and sub processing functions as implemented: le1, vis, nir, sir, ... and could be completed as much as desired.

Dictionary : sys namespace

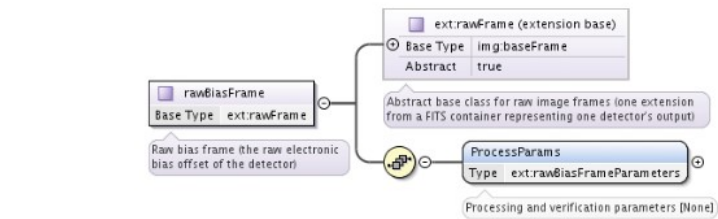
this branch (acronym of system) is the common namespace for processing and operational data types used in the Common Data Model and that are shared by transverse components and non PFs components of Euclid SGS such as Infrastructure Abstraction Layer, Orchestration component in charge of the data and processing distribution on SDCs and EC-SURV that is the component in charge of the Survey definition.

Dictionary : bas namespace

- **cat**: definition of catalog of astrophysical sources (transients, galaxies, stars, ...)
- **cot**: **coordinates-time domain such as time, space, spectral etc. not covered by standard IVOA STC scheme (see: imp/stc)**
- **dqc**: definition of Data Quality Common types shared by all Processing Functions
- **dtd**: data types introduced by Euclid common data model or redefined from standard XML data types: array, matrices, angles, lists, ...
- **fit**: description of fits files format used in Euclid (each fits file used by Processing Functions shall be referenced in this section)
- **img**: basic definitions for Euclid (and external) images (frames)
- **imp**: **import of external schema derived from IVOA (stc), CCSDS, ESO or FITS standards (elementary data types derived directly from the FITS specification)**
- **mat**: mathematical definitions for polynomial computation, Gaussian,...
- **msk**: definitions of the binary masks related to images
- **ppr**: definitions related to the Parameters of some common processing steps of the overall Euclid processing pipeline
- **psf**: definitions related to the Point Spread Function object and its processing
- **utd**: **Units used by Euclid, we have to mention that the units used are conform to the IVOA standards referenced here: <http://www.ivoa.net/internal/IVOA/UnitsDesc/Units.html>**

Data model implementation

Use binding generation in pipeline code



```

class RawBiasFrame(RawFrame):
    """ Raw bias frame (the raw electronic bias offset of the detector) """
    ProcessParams=persistent('Processing and verification parameters [None]',RawBiasFrameParameters,None)
  
```

```

CREATE TABLE "AMQPST"."RAWBIASFRAME"
(
  "OBJECT_ID" RAW(16) DEFAULT SYS_GUID(), "CREATOR" NUMBER DEFAULT 1, "PROJECT"
  "PRIVILEGES" NUMBER DEFAULT 1, "CREATED" TIMESTAMP (6) DEFAULT SYS_EXTRACT_UTCTIME
  "MODIFIED" TIMESTAMP (6) DEFAULT SYS_EXTRACT_UTCTIME,"VERSION" NUMBER DEFAULT
  "CHIP" RAW(16), "CHIPS" VARCHAR2(30 BYTE),
  "CREATEDATE" TIMESTAMP (6) DEFAULT TO_TIMESTAMP('1998-01-01 00:00:00','YYYY-MM-DD HH24:MI:SS'),
  "DATE" TIMESTAMP (6) DEFAULT TO_TIMESTAMP('1998-01-01 00:00:00','YYYY-MM-DD HH24:MI:SS'),
  "DATETIME" TIMESTAMP (6) DEFAULT TO_TIMESTAMP('1998-01-01 00:00:00','YYYY-MM-DD HH24:MI:SS'),
  "EXTOBJECTID" RAW(16), "EXTOBJECTS" VARCHAR2(30 BYTE), "EXTENSION" NUMBER(,0) DEFAULT 0,
  "INSTAT" RAW(16), "INSTATS" VARCHAR2(30 BYTE), "INSTRUMENT" RAW(16), "INSTRUMENTS" VARCHAR2(30 BYTE),
  "INVALID" NUMBER(,0) DEFAULT 0, "LIST" BINARY_DOUBLE DEFAULT 0.0, "PIDWORDS" BINARY_DOUBLE DEFAULT 0.0,
  "INDEX" NUMBER(,0) DEFAULT 0, "INDEXES" NUMBER(,0) DEFAULT 0,
  "OBJECTS" RAW(16), "OBJECTS" VARCHAR2(30 BYTE), "OBSLOCK" RAW(16), "OBSLOCKS" VARCHAR2(30 BYTE),
  "OBSERVER" RAW(16), "OBSERVERS" VARCHAR2(30 BYTE), "OVERSCANSTAT" RAW(16), "OVERSCANSTATS" VARCHAR2(30 BYTE),
  "OVERSCANSTAT" RAW(16), "OVERSCANSTATS" VARCHAR2(30 BYTE),
  "OVSCK" NUMBER(,0) DEFAULT 0, "OVSCKPR" NUMBER(,0) DEFAULT 0, "OVSCKPST" NUMBER(,0) DEFAULT 0,
  "OVSCKY" NUMBER(,0) DEFAULT 0, "OVSCKYPR" NUMBER(,0) DEFAULT 0, "OVSCKYPT" NUMBER(,0) DEFAULT 0,
  "PRSCANSTAT" RAW(16), "PRSCANSTATS" VARCHAR2(30 BYTE),
  "PRSCANSTAT" RAW(16), "PRSCANSTATS" VARCHAR2(30 BYTE),
  "PROCESSPARAMS" RAW(16), "PROCESSPARAMS" VARCHAR2(30 BYTE),
  "PRSCY" NUMBER(,0) DEFAULT 0, "PRSCYPR" NUMBER(,0) DEFAULT 0, "PRSCYPT" NUMBER(,0) DEFAULT 0,
  "PRSCY" NUMBER(,0) DEFAULT 0, "PRSCYPR" NUMBER(,0) DEFAULT 0, "PRSCYPT" NUMBER(,0) DEFAULT 0,
  "QUALITYFLAG" NUMBER(,0) DEFAULT 0,
  "RAWBIAS" RAW(16), "RAWBIAS" VARCHAR2(30 BYTE), "STORAGE" RAW(16), "STORAGE" VARCHAR2(30 BYTE),
  "TEMPLAT" RAW(16), "TEMPLATS" VARCHAR2(30 BYTE), "UTC" BINARY_DOUBLE DEFAULT 0.0
)
  
```

Database schema

Original Euclid Data Model (XSD)

Python classes

User Interface and services

ROWNUM	project_id	PRIVILEGES	object_id	CreationDate	Date	EuclidObs	Extension bit
1	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
2	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
3	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
4	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
5	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
6	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
7	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
8	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
9	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
10	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
11	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
12	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
13	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
14	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
15	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
16	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
17	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
18	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
19	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
20	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
21	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
22	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
23	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
24	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
25	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
26	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0
27	1	2	object view	2011-08-01 01:18:08	2011-08-01 12:01:04	2011-08-01 12:01:04	0

- Detailed
- ECDM objects implemented in full in the EAS
- Provides lineage & traceability

Choosing XML schema as the formal language of data description brings numerous advantages; some of them are particularly desirable:

- **XML documents conform with xsd are structured**
- **XML documents conform with xsd are portable and readable without proprietary tools**
- **XML documents conform with xsd are easily customised**

A major shortcoming of this choice is the verbosity and the volume of XML documents. This drawback is addressed via specific rules on xsd.

Concerns related to IVOA schema vs Euclid rules

1/ Lack of maintainability and reusability

- Every schema is self-contained and redefine its own types with a lot of similarity with other schemas
- coordSysType is defined in spectrumDM
- Some types 'redshiftFrameType' are defined twice (in STC and spectrumDM) with different definitions
- No namespace, no import, no versioning of schema inside

2/ Complexity :

- STC make an intensive use of Substitution Group. Not recommended
- Type « abstract » : lack of visibility
- No « mixed content » : complexity
- No « anyType » : complexity

3/ Understandability :

- Elements and types are mixed : it may be confusing

Decisions about IVOA schema

Facing above incompatibilities between ivoa schema defined in: <http://www.ivoa.net/xml/> and EUCLID rules , the following decisions were made:

- Preferred STC metadata linear string implementation,
- Huge data sets such as catalogs that describe similar objects with same attributes, and same reference frames should be packed as far as possible in a specific header (to minimize the volume of the instance of schema),
- Extract types (simple and complex) from IVOA xsd schema relevant for Euclid,
- Dictionary (thematic tree) to set up from extraction of these types: basic, units, misc, projection, timescale, coordinates, photo, spectrum, astroobjects,
- Keep track of the IVOA source through dedicated annotation, keep IVOA syntax for a better identification of reused IVOA types,
- Re use of ref Frame enumeration: projection, time, code projection,...
- Simplification of attribute Group

Conclusion

Suggested improvements on schema :

- Define common types in IVOA data models
- Regroup types by namespace and manage versionning
- Reuse types in scientific data models (characterisation, spectrum, ...)
- Use only one root element
- Keep it as simple as possible
- Have a look on 'similar' standards from EO : OGC
 - Spatial objects have simpler definitions