# Integrating AI tools in data analysis frameworks: the Vera Rubin LSST and Euclid cases

**Riccio G.**[1], Cavuoti S.[1], Angora G.[1], Brescia M.[2,1]

1. INAF - Astronomical Observatory of Capodimonte, Napoli (Italy)
1. Department of Physics "E. Pancini", University of Naples Federico II (Italy)

# Context

In the last two decades, Astronomy has been the scene of the realization of panchromatic surveys, with sophisticated instruments acquiring a huge amount of exceptional quality data.



- ESA Euclid : ~100 GB/day for 6 years → 200 TB
- Rubin/LSST : ~20 TB/night for 10 years → >60 PB
- JWST : ~30GB/day for 10 years (and more)
- GAIA : ~1 PB in 5 year
- SKA : 100 Pbytes – 3 EBytes/year
- Pan-STARRS, KiDS, DES, Herschel-ATLAS, Hi-GAL, E-ELT…

## NEEDS

- to integrate advanced **data-driven** science methodologies for the **automatic exploration** of huge and multi-dimensional **data archives**

- efficient short- and long-term **monitoring** and **diagnostics** systems
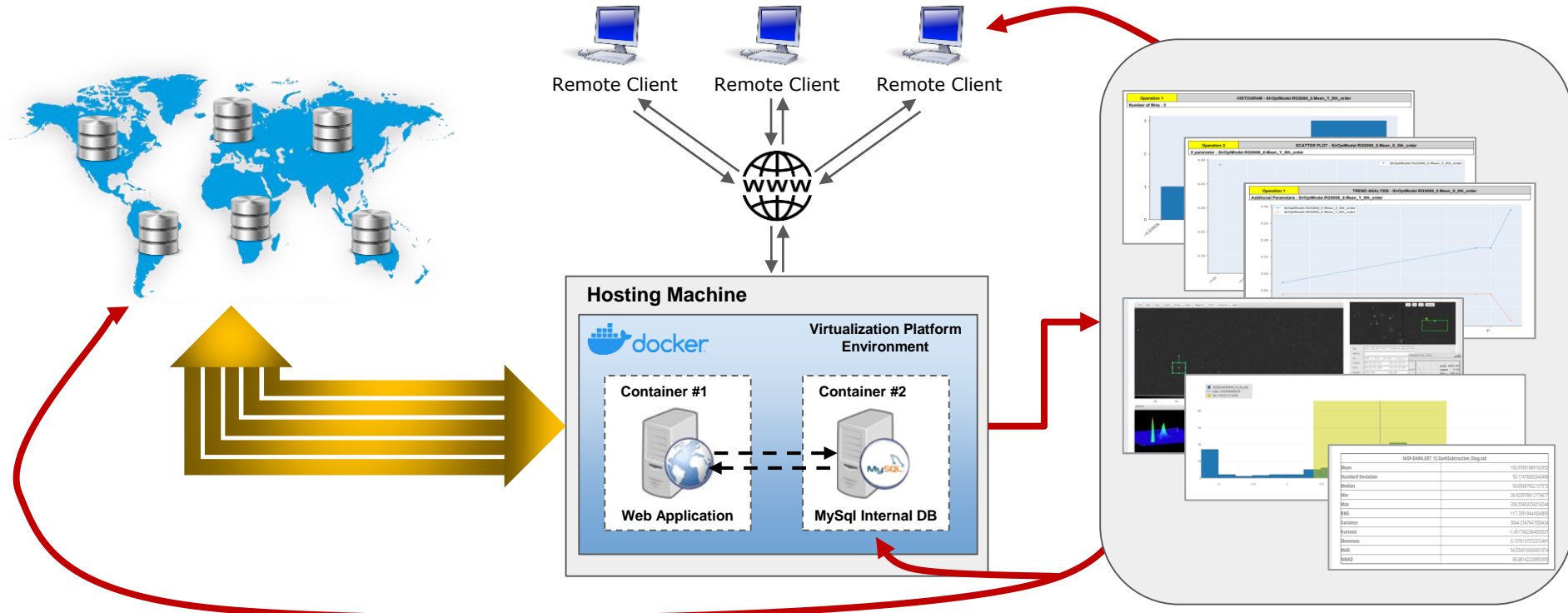
## GOALS

- To keep the **quality** of the observations under control

- To detect and circumscribe **anomalies** and **malfunctions**

- To facilitate rapid and effective **corrections**

- To ensure correct maintenance of all components and the **good health of scientific data** over time mainly crucial for space-borne observation systems, both in logistical and economic terms

# AIDA - Advanced Infrastructure for Data Analysis

*AIDA is a portable and modular web application, designed to provide an efficient and intuitive software infrastructure to support monitoring of data acquiring systems over time, diagnostics and both scientific and engineering data quality analysis, particularly suited for astronomical instruments*
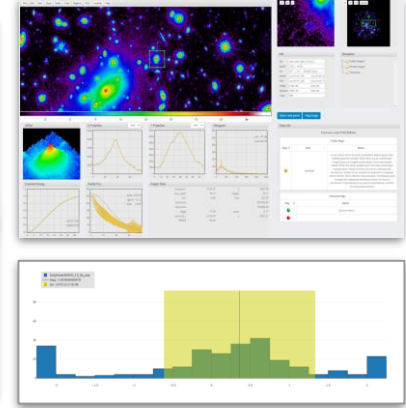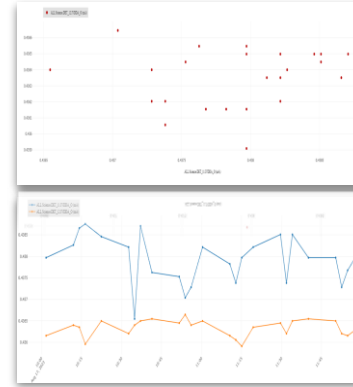
# AIDA - Main Features

## Instrument monitoring, report generation and delivery

✓ **periodic report** generation on a user-defined parameters list and delivery to remote archive

✓ **on demand customised report** generation on a user selected parameter list, locally stored

## Visualization/Exploration

✓ **series of plots** on user selected parameters/data products and ranges

✓ **pre-generated histograms** stored into remote archives

✓ **observed images** (static view, dynamic windowing, statistical characterization)

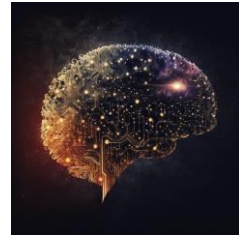## Statistics

✓ **standard** (default) estimators

✓ **special** estimations (tables/images)

✓ statistical analysis on **image pixels**

## Machine Learning

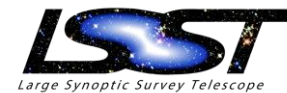✓ **Regression/classification** experiments on available data
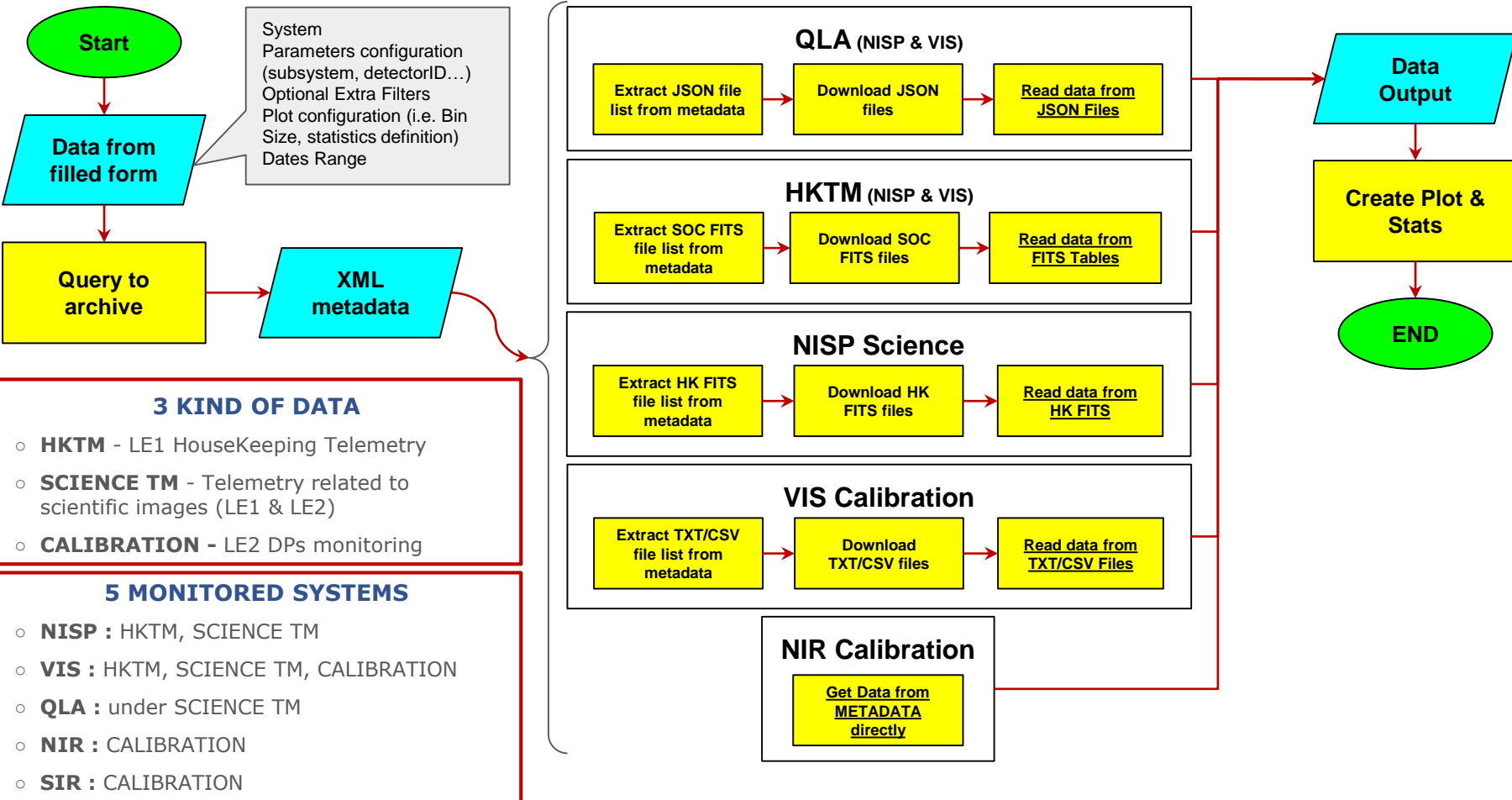
# Code Name : Flexibility

AIDA has been designed as a **modular system**, based on **Object-Oriented Programming** and specific information on DB, so it is possible to **extend its functionalities**, by integrating and customizing monitoring and diagnostics systems, as well as scientific data analysis solutions, including machine/deep learning and data mining methods

➢ **Available plots and statistics are defined as classes/functions** linked to a specific table in DB. To add a new operation, it is sufficient to implement the related class/function and add it to the local DB;

➢ **A JSON configuration file is associated to every system monitored by AIDA**. It includes info about the instrument and connection to the related data and metadata archives;

➢ **Repositories and systems have a dedicated classes** which implement methods for interfacing AIDA with the data repository. To add a new system/repository, it is sufficient to create its own configuration file (only for systems), implement the related class and methods, and fill DB with required information.
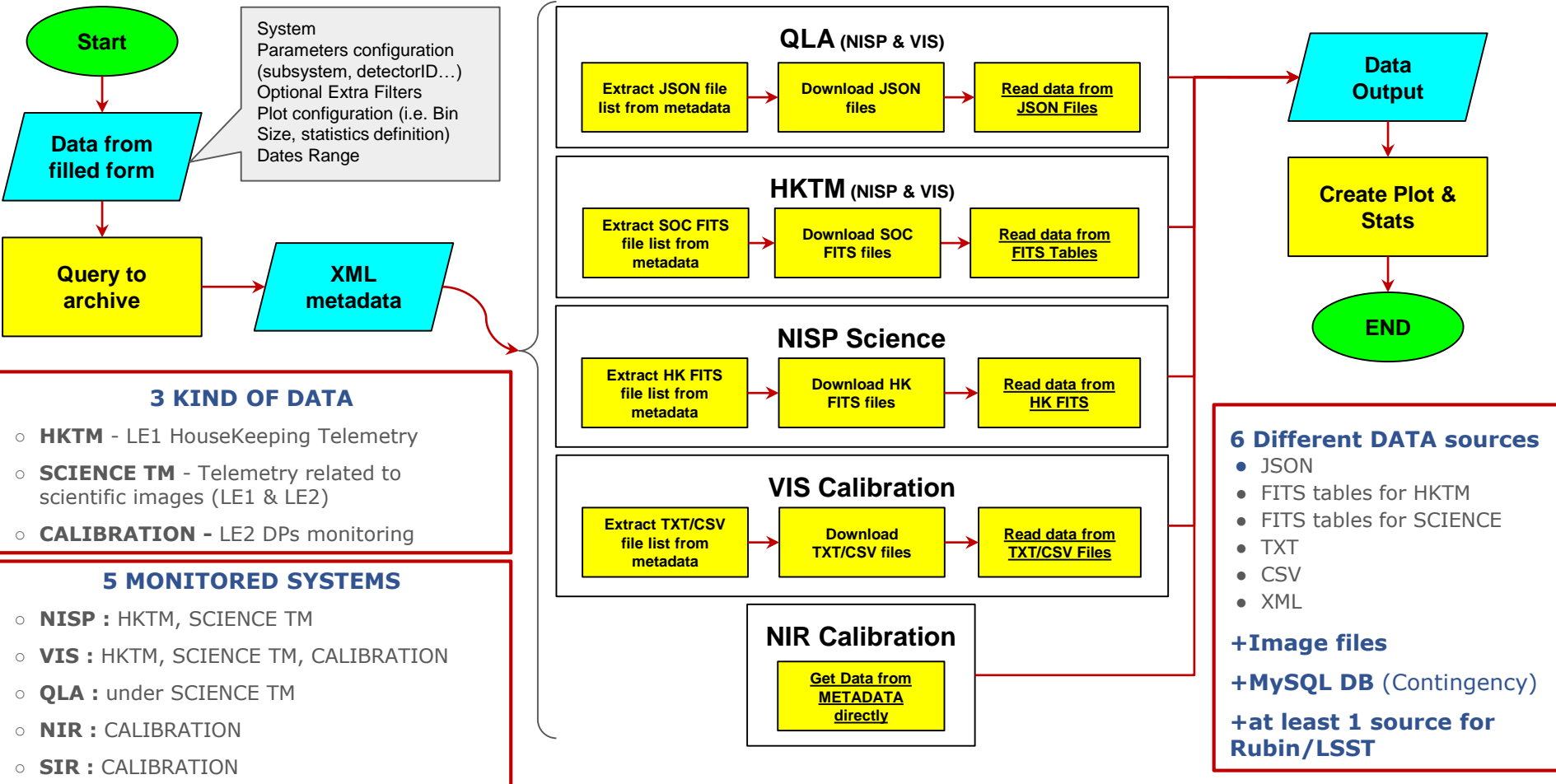
A specialized version of AIDA is already the **official monitoring** and analysis tool for the **ESA Euclid space mission** and another one is going to be used for the commissioning of the **V. Rubin Telescope, suitable also for LSST survey data**
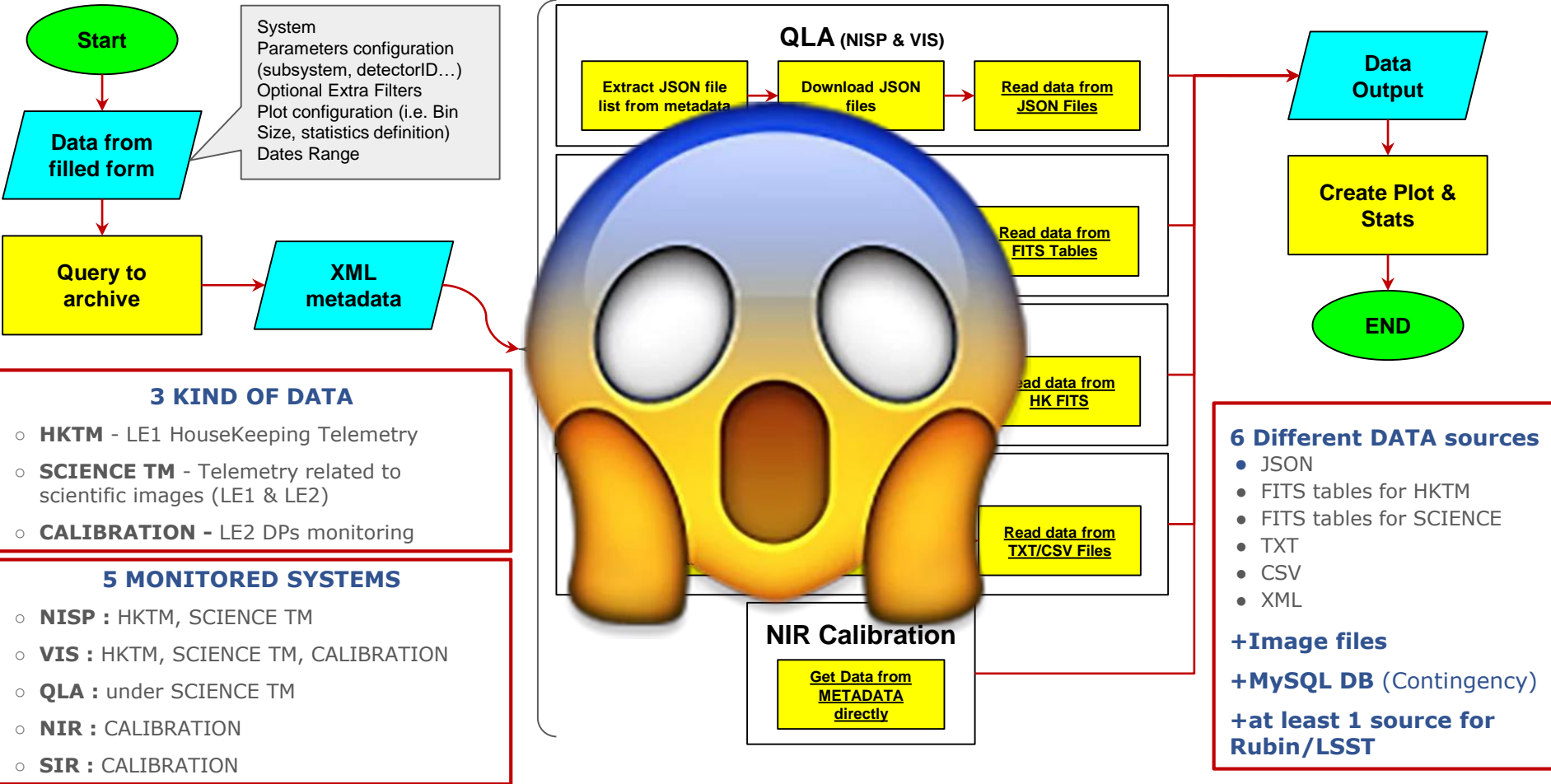
# AIDA/IODA for Euclid Data

```
Start
```

System
Parameters configuration
(subsystem, detectorID…)
Optional Extra Filters
Plot configuration (i.e. Bin Size, statistics definition)
Dates Range

```
Data from filled form
```

```
Query to archive
```

```
XML metadata
```

## 3 KIND OF DATA

- ○ **HKTM** - LE1 HouseKeeping Telemetry
- ○ **SCIENCE TM** - Telemetry related to scientific images (LE1 & LE2)
- ○ **CALIBRATION -** LE2 DPs monitoring

## 5 MONITORED SYSTEMS

- ○ **NISP :** HKTM, SCIENCE TM
- ○ **VIS :** HKTM, SCIENCE TM, CALIBRATION
- ○ **QLA :** under SCIENCE TM
- ○ **NIR :** CALIBRATION
- ○ **SIR :** CALIBRATION

### QLA (NISP & VIS)

```
Extract JSON file list from metadata → Download JSON files → Read data from JSON Files
```

### HKTM (NISP & VIS)

```
Extract SOC FITS file list from metadata → Download SOC FITS files → Read data from FITS Tables
```

### NISP Science

```
Extract HK FITS file list from metadata → Download HK FITS files → Read data from HK FITS
```

### VIS Calibration

```
Extract TXT/CSV file list from metadata → Download TXT/CSV files → Read data from TXT/CSV Files
```

### NIR Calibration

```
Get Data from METADATA directly
```

```
Data Output
```

```
Create Plot & Stats
```

```
END
```

# AIDA/IODA for Euclid Data

# AIDA/IODA for Euclid Data

**Start**

**Data from filled form**

System
Parameters configuration (subsystem, detectorID…)
Optional Extra Filters
Plot configuration (i.e. Bin Size, statistics definition)
Dates Range

**Query to archive**

**XML metadata**

## QLA (NISP & VIS)

**Extract JSON file list from metadata**

**Download JSON files**

**Read data from JSON Files**

**Read data from FITS Tables**

**Read data from HK FITS**

**Read data from TXT/CSV Files**

### NIR Calibration

**Get Data from METADATA directly**

**Data Output**

**Create Plot & Stats**

**END**

### 3 KIND OF DATA

○ **HKTM** - LE1 HouseKeeping Telemetry

○ **SCIENCE TM** - Telemetry related to scientific images (LE1 & LE2)

○ **CALIBRATION -** LE2 DPs monitoring

### 5 MONITORED SYSTEMS

○ **NISP :** HKTM, SCIENCE TM

○ **VIS :** HKTM, SCIENCE TM, CALIBRATION

○ **QLA :** under SCIENCE TM

○ **NIR :** CALIBRATION

○ **SIR :** CALIBRATION

### 6 Different DATA sources
- JSON
- FITS tables for HKTM
- FITS tables for SCIENCE
- TXT
- CSV
- XML

**+Image files**

**+MySQL DB** (Contingency)

**+at least 1 source for Rubin/LSST**

# Machine Learning Tools



The tool includes **more than 100 prediction**, **classification and regression models** based on **Machine Learning** to apply on available tabular data, useful in this case to identify operating anomalies or correlations between instrumental information.
**Deep Learning** methods coding is on going

# Summary

✓ The **AIDA web application** has been designed to provide an **efficient and intuitive software infrastructure** to support **monitoring** of data acquisition systems over time, **diagnostics** and both scientific and engineering **data quality analysis**, in particular for astronomical instruments

✓ It provides **a number of tools** for data analysis & system diagnostics

- ❑ **Instrument monitoring, report generation and delivery**
- ❑ **Visualization Exploration**
- ❑ **Statistics**
- ❑ **Machine/Deep Learning**

✓ a specific version of AIDA is already the **official monitoring** and analysis tool for the **ESA Euclid space mission** and another one is going to be used for the commissioning of the **V. Rubin Telescope, suitable also for LSST survey data**

✓ Being designed as a modular system, **it is possible to integrate and customize** monitoring and diagnostics systems, as well as scientific data analysis solutions

An **high level of standardization** for data and tools is crucial to easily customize AIDA to have a **general infrastructure** for as many astronomical projects as possible

# Ideas for next AI tools (1)

**REPORT**

## Periodic automatic or on-demand generation



Why not **a pre-trained LLM** to automatically create configuration files?

# Ideas for next AI tools (2)



A very useful tool could be a function, runnable from the Image Explorer panel, to **automatically** generate **thumbnails** from images, to be used by Deep Learning methods

A standard and automatic thumbnail extractor would be very useful for astronomical community in general

**Based on JS9 Library : https://js9.si.edu/**