# Analyzing large astronomical datasets: The Science Portal solution

Angelo Fausti Neto
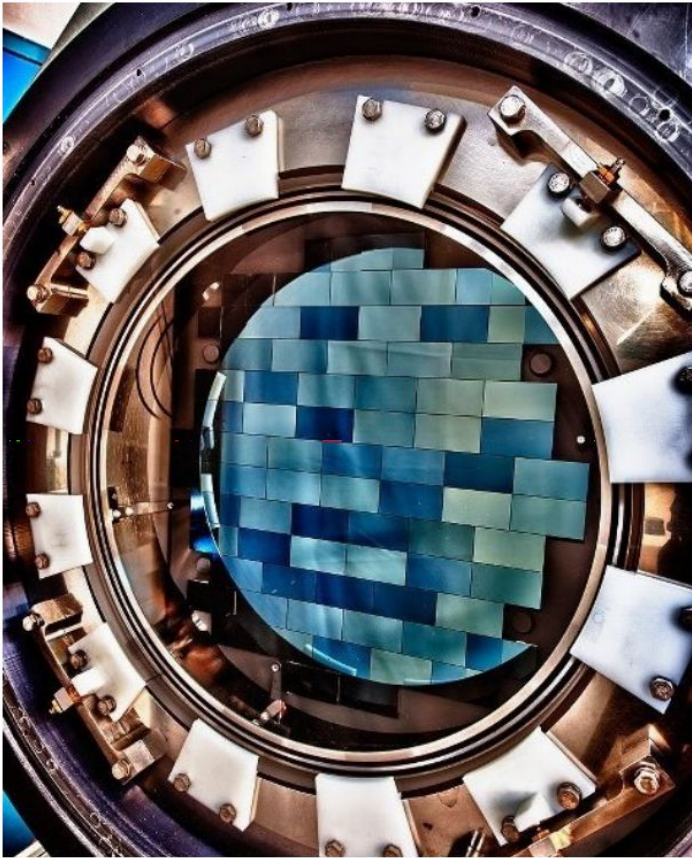
Laboratório Interinstitucional de e-Astronomia (LIneA)
Large Synoptic Survey Telescope (LSST)

# Outline

- Brazilian participation in Large Surveys
  - The role of LIneA  (see Ricardo Ogando's presentation)
  - DES-Brazil, BPG-SDSS, DESI, BPG-LSST

- Science Portal
  - Applications in DES
  - Challenges

- Conclusions and perspectives

# Dark Energy Survey (DES)

https://www.darkenergysurvey.org/



Dark Energy Camera, in operation since 2012 at Blanco (CTIO)

- Photometric Survey (grizY)

- DECam 570 Mpixels (62 CCDs)

- Blanco 4m  (CTIO)

- ~ 300 exposures each night (500GB)

- 5 years (~100 nights/year)

- 5.000 sq deg

- 4th year of operations (Ago 2016)

- First Public Data Release - DR1 (2017)

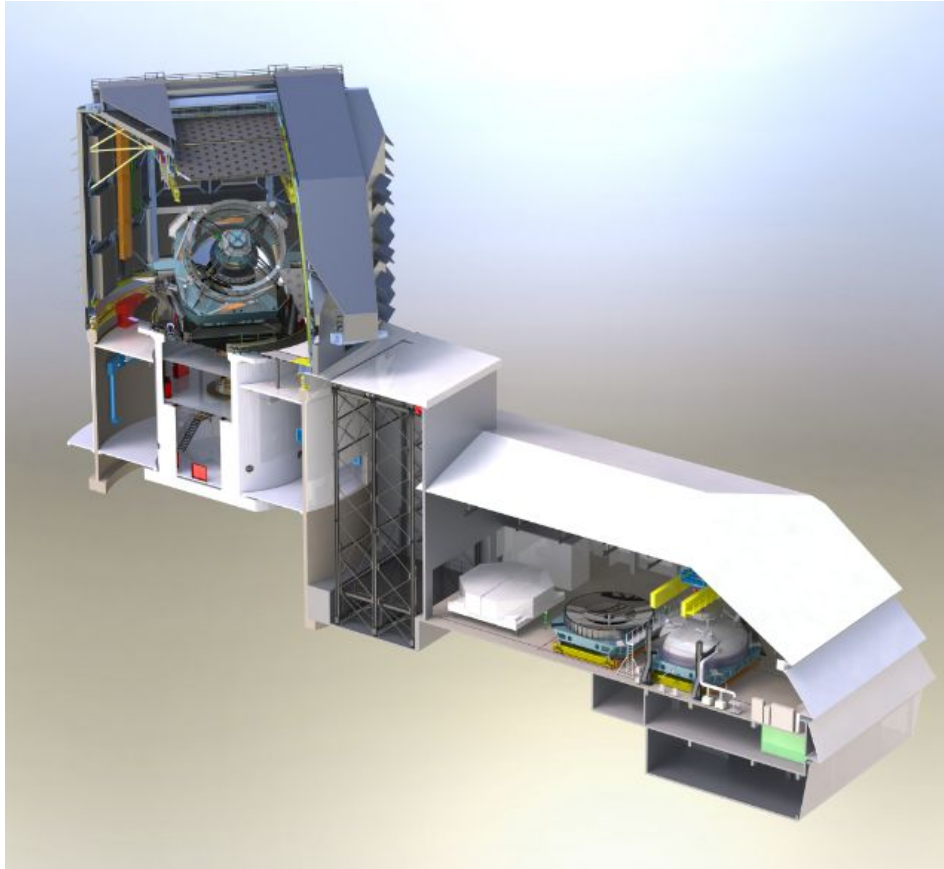- Objects catalog ~1Billion

- DES-Brazil  (2007)



Começa o quarto ano de observações do Dark Energy Survey

19 de agosto de 2016   http://www.linea.gov.br/Noticias/

# Large Synoptic Survey Telescope (LSST)

https://www.lsst.org/scientists



- Photometric Survey (ugrizy)
- LSST Cam 3.2 Gpixels (189 CCDs)
- M1/M3 primary 8.4m  (CTIO)
- 1 "visit" = 2 exposures in 30s
- 1000 "visits" each night (~15 TB)
- 18.000 sq deg - twice a week for 10 years
- Start operations in 2022
- 10 Million alerts/night
- Objects Catalog ~37 Billion
- BPG-LSST (2016)

http://bpg-lsst.linea.gov.br/

# Common challenges

- ● Big science questions
  - ○ Active participation in international collaborations
  - ○ "M-shaped" person (astronomy, scientific computing, statistics)
- ● Big collaborations
  - ○ Cultural change  (**collaborative tools)**
  - ○ Communication challenges (**telecons, distributed information**)
  - ○ Understand the projects and their opportunities
- ● Big data volume and variety
  - ○ Infrastructure to support science (data transfer, storage and processing)
  - ○ Efficient data preparation and analysis
  - ○ Data science techniques

# The "new astronomy"

**IN THE PAST**

Data  >  Inference  > Model

(astronomer)

**THESE DAYS**

Data  > Processing > Catalog  > Inference > Model

("Big data" problem)

("New astronomer")

"M-shaped" person

Adapted from Andrew Connolly - LSST

# The role of LIneA in Brazil

http://www.linea.gov.br/

- ## Support science
  - DES-Brazil (2007), BPG-SDSS-III (2008), BPG-SDSS-IV (2014), DESI (2016) and BPG-LSST (2016)
  - 71 members  (ON, UFRJ, USP, UNESP, UNICAMP, UFABC, UFRGS, UFSM)
  - Scientific and technical education
  - Public Outreach

- ## Data center operation
  - Hardware and services maintenance (external contract)
    - helpdesk, e-mail, twiki, git, doc-db, slack, etc
  - Data transfer (RNP/Brazil)
  - Database and processing consultant (LNCC/Brazil)

- ## Software Development
  - R&D
  - Software Development team (9 FTEs)
  - 3 PhDs in astronomy + external contributions

# What is the Science Portal?

- It is the collection of software tools developed by LIneA to **assess the data quality** and to **explore large astronomical datasets.**

- It is also a **web framework** that facilitates the integration of the **science analysis codes** into a hardware and software infrastructure that provides **efficient data preparation and analysis**.

# Science Portal applications in DES

- Data Quality Assessment in real-time
  - Quick Reduce @ CTIO (2012)
  - http://quick1.ctio.noao.edu:8080/

- Data validation and exploration
  - Science Server @ Fermilab (2014)
  - Science Server @ NCSA (2016)
  - https://des-portal.fnal.gov/
  - http://desportal.cosmology.illinois.edu/

A new production system each 2 years

- Preparation of science-ready catalogs and science workflows
  - Science Portal @ LIneA (2016)
  - http://des-portal.linea.gov.br/

# DES Quick Reduce @ CTIO

Data reduction in real-time ~250k DECam exposures (~2.5M CCDs) since 2012
Developed for DES but also available for other programs

# Monitoring DES observations
## Daily transfer of QR results to Fermilab, access to DES collaboration
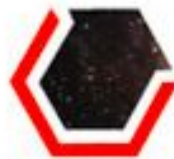
# THE DARK ENERGY SURVEY

Welcome to the DES Science Portal @ FNAL

## Services available

- Data Upload
- Visualization tools
- Cutout service
- User query

### Login

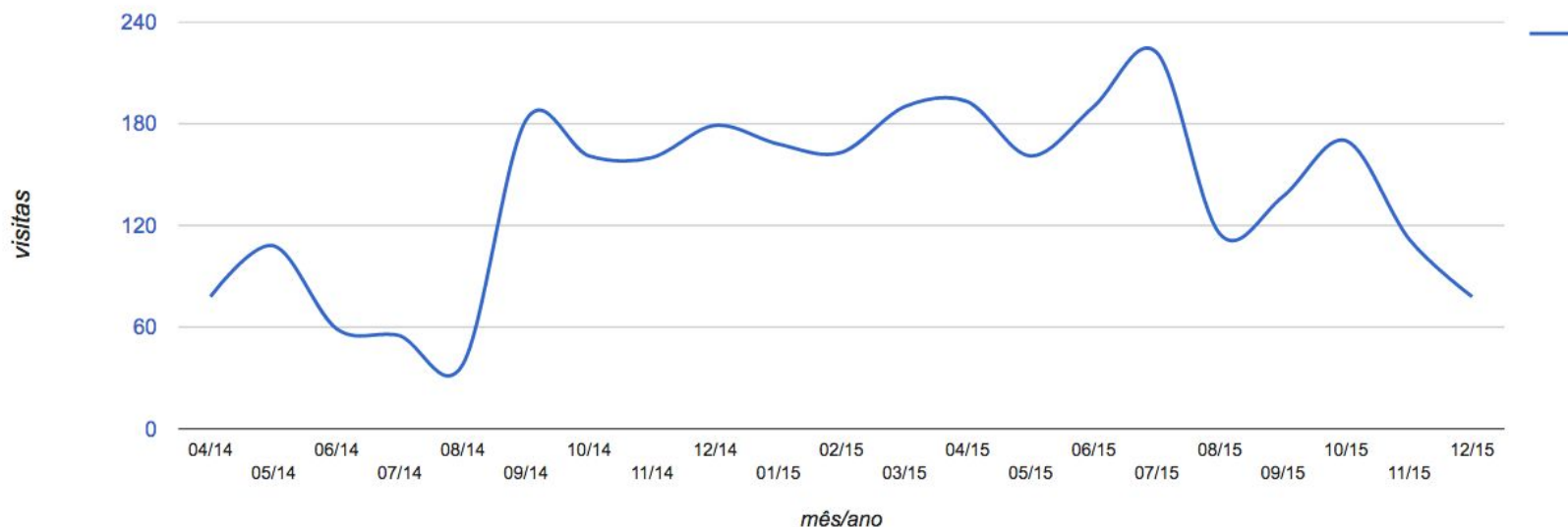Use your FNAL services username and password.

Username: [          ]

Password: [          ]

[ Login ]

Need Help?

~180 DES collaboration members each month

"If you build they will come..."

## Visitas Fermilab



*visitas* (y-axis): 0, 60, 120, 180, 240

*mês/ano* (x-axis): 04/14, 05/14, 06/14, 07/14, 08/14, 09/14, 10/14, 11/14, 12/14, 01/15, 02/15, 03/15, 04/15, 05/15, 06/15, 07/15, 08/15, 09/15, 10/15, 11/15, 12/15

12

# Upload, cutout and visual inspection
## e.g Strong Lensing, Galaxy Clusters

# New Science Server @ NCSA*

**Science Server**

- 🏠 Home
- ✔ Releases
- ⭐ Sky Viewer
- ▦ Tile Viewer
- ◎ Target Viewer
- 🗄 Sky Query
- ⬆ Upload
- 🖼 Cutout Server
- 📖 Science Products

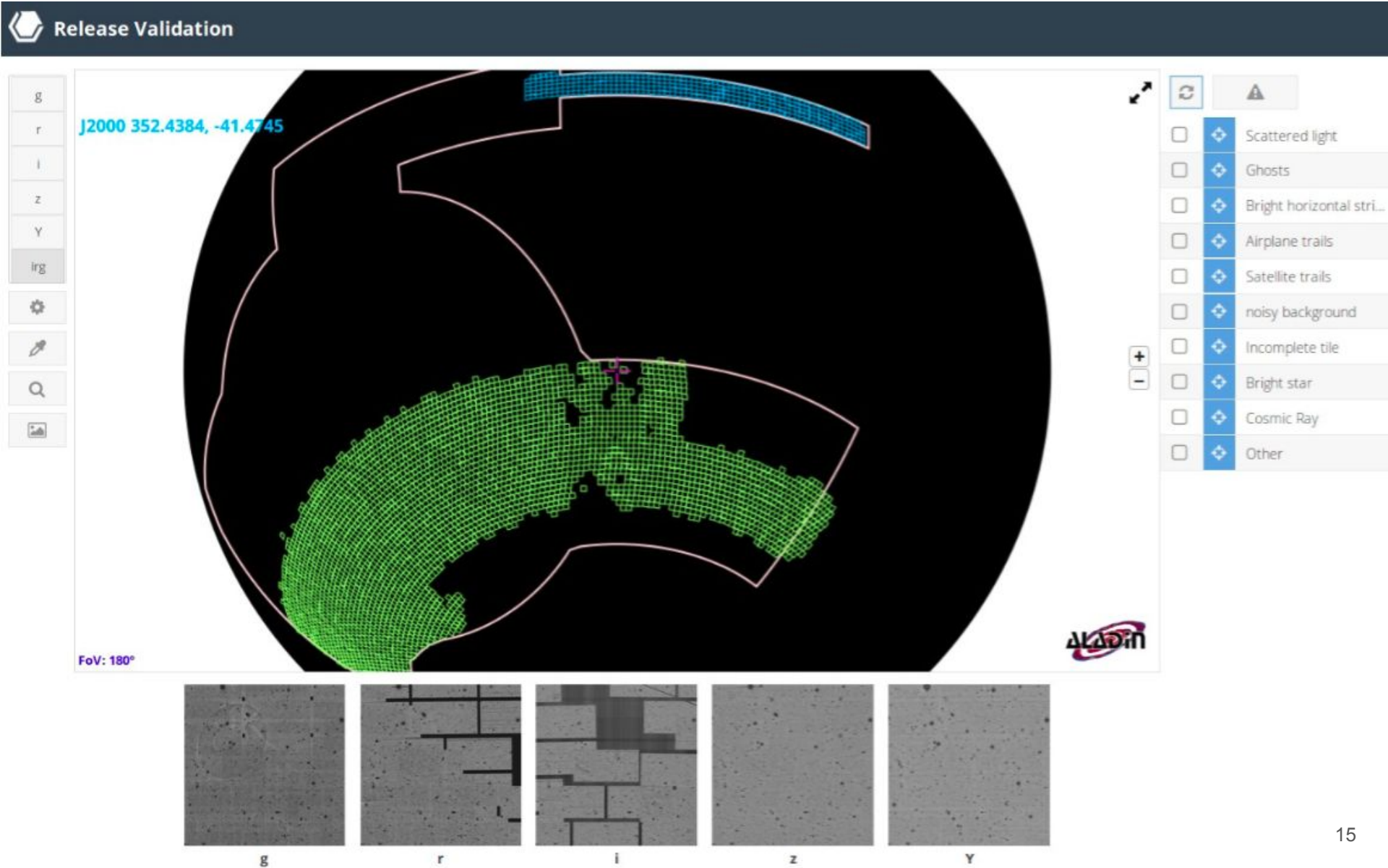| | |
|---|---|
| **Releases** — Summary information of DES releases and validation | **Sky Viewer** — All-sky visualization of DES releases in grizY and RGB with overlay of tile grid and objects |
| **Tile Viewer** — Inspect DES tiles using visiOmatic tools | **Target Viewer** — Manage lists of targets with image display, cutouts, ranking and reject functionalities. |
| **Sky Query** — Query catalogs using sample queries or keep your own query library | **Upload** — Upload external data to the Science Server |
| **Cutout Server** — Create co-added or single epoch cutouts from a list of coordinates | **Science Products** — Serve catalogs created by the collaboration |

New technologies:

- Python 3
- Django 1.9
- ExtJS 6
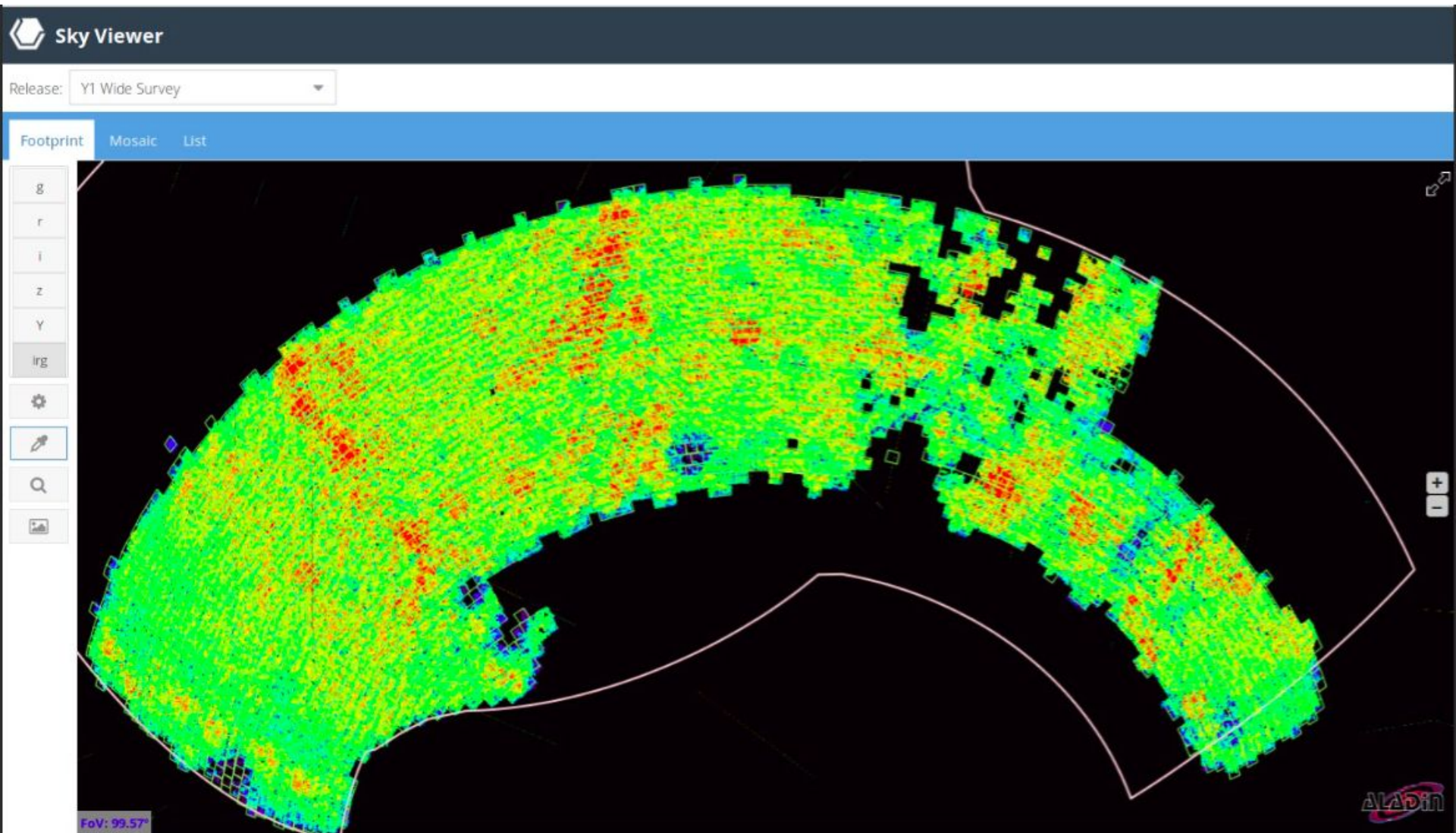- Aladin, VisiOmatic

\* DES DR1 prototype (2017)

14

# Data Release Validation
## Being used to evaluate DES internal data releases

# Visualization of the survey properties
## HEALPix map of the survey magnitude limit

# Detailed image visualization in the browser
## (visiOmatic  E. Bertin)



**Release Validation**

Identification of defects

- Scattered light
- Ghosts
- Bright horizontal stri..
- Airplane trails
- Satellite trails
- noisy background
- Incomplete tile
- Bright star
- Cosmic Ray
- Other

Full control over image properties (contrast, RGB levels) JPEG is generated in real-time at the server and sent to the client.
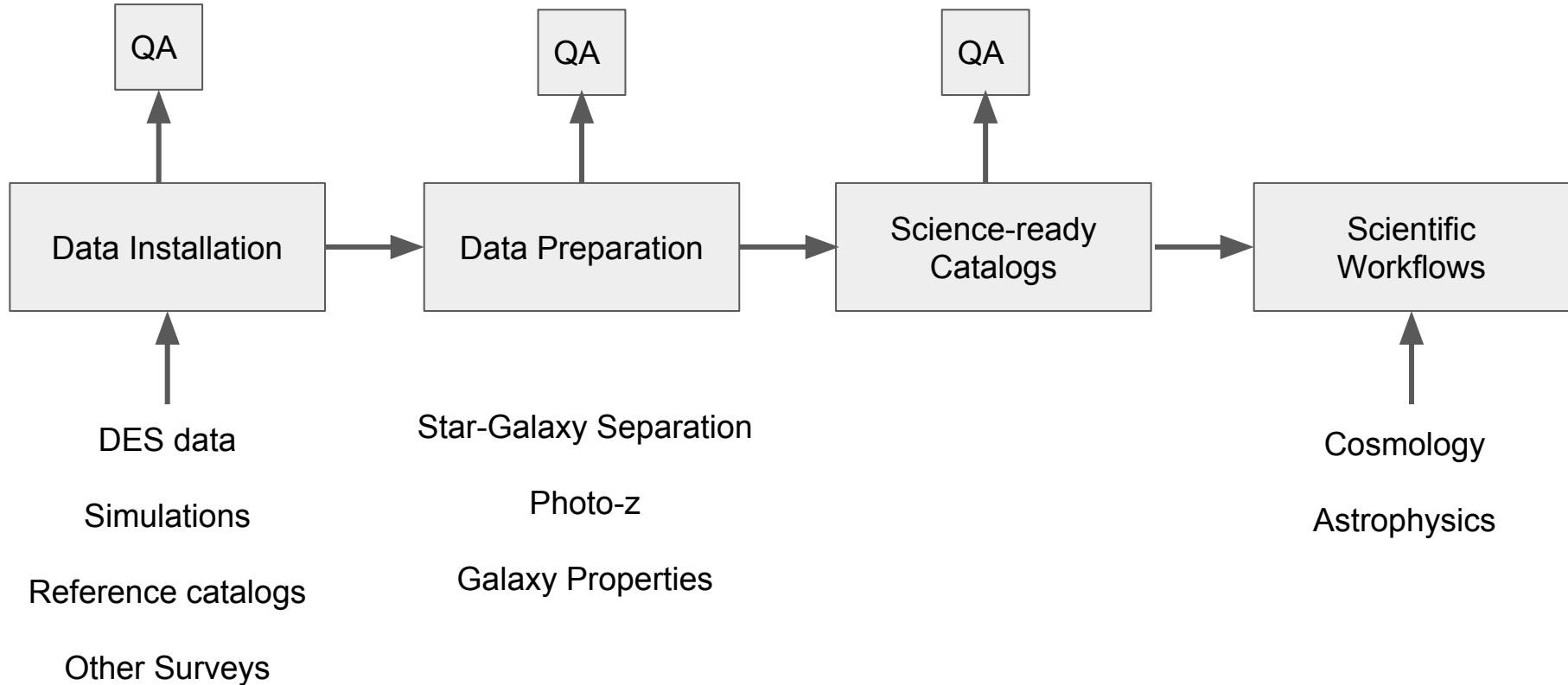
Dark Energy Survey @ 2016 NCSA/LIneA

grizY thumbnails

g    r    i    z    Y

# Science Portal: Science Analysis Framework



QA        QA        QA

Data Installation → Data Preparation → Science-ready Catalogs → Scientific Workflows

DES data

Simulations

Reference catalogs

Other Surveys

Star-Galaxy Separation

Photo-z

Galaxy Properties

Cosmology

Astrophysics

# Preparation of science-ready catalogs

(Fausti et. al 2016 in prep)



**Four steps:**

(a) region selection
(b) object selection
(c) column selection
(d) addition of value-added quantities  (sg separation, photo-z, galaxy properties)
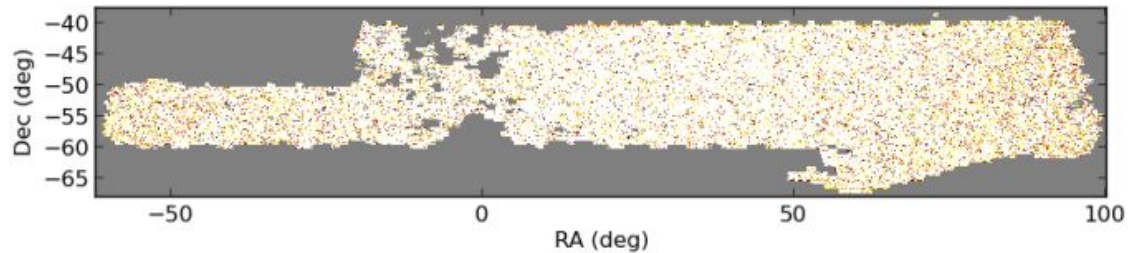
# Why building this infrastructure?

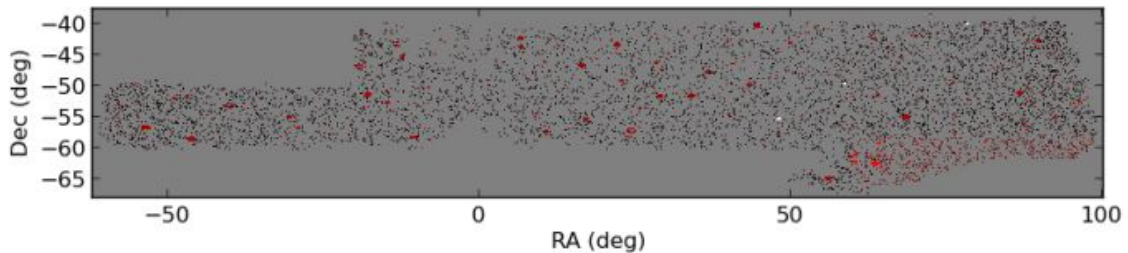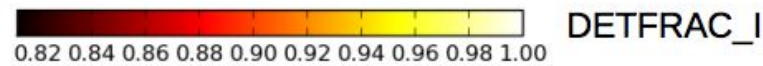(Fausti et. al 2016 in prep)

- <u>Preservation</u> of the codes developed by DES Collaboration to create "ancillary products"

- <u>Reproducibility</u> of the catalogs

- <u>Control</u> the parameters used in the creation of the catalogs

- <u>Provenance</u> of the input data products

- <u>Documentation</u> of the catalog and its properties
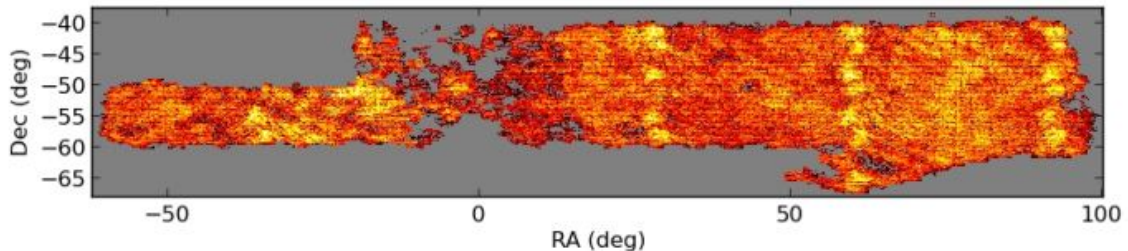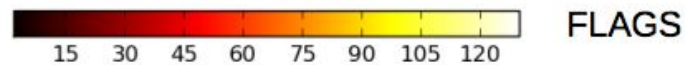
# Example: galaxy magnitude-limited catalog
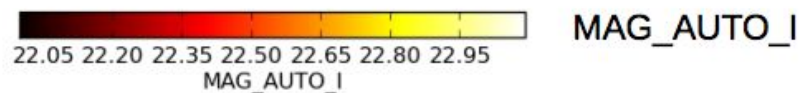
(Fausti et. al 2016 in prep)
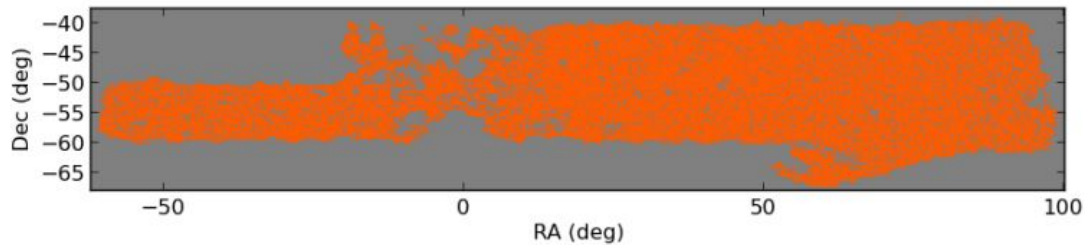


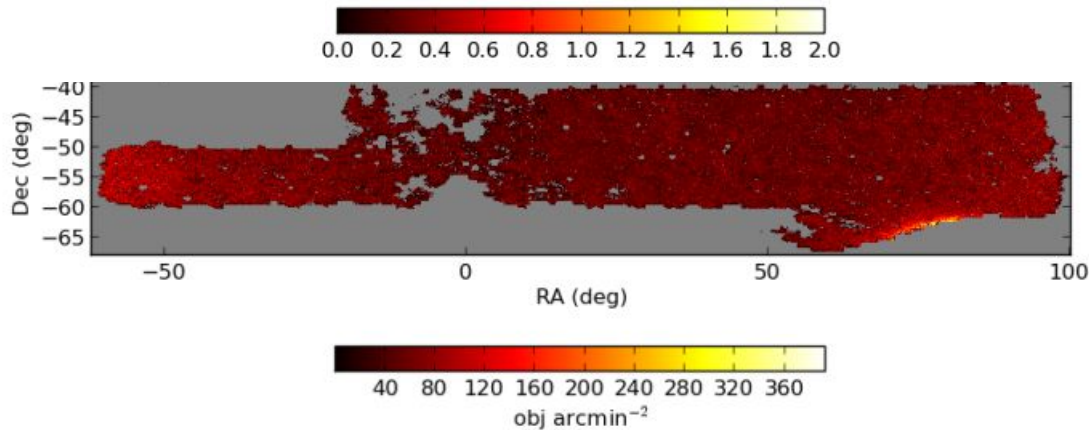"Good" regions for science

Foreground objects mask

Magnitude limit map

21

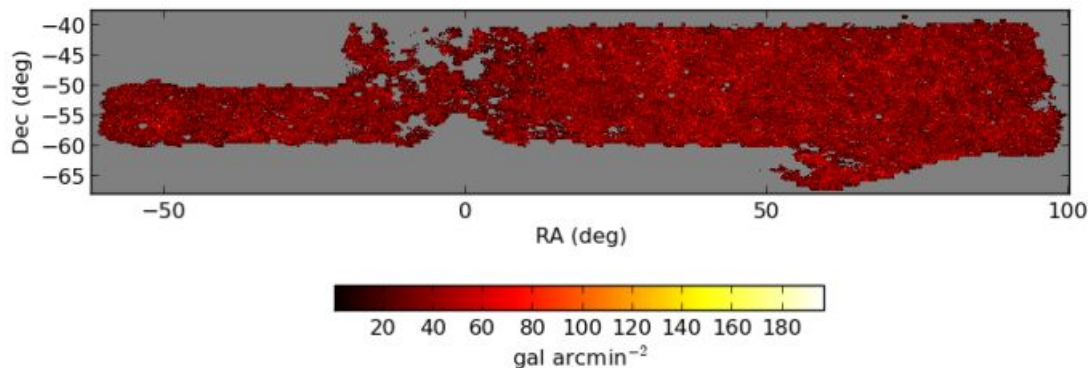# Preparation of science-ready catalogs

(Fausti et. al 2016 in prep)



Footprint map  (after region selection)

Object Selection

Star-galaxy separation, photo-z

**Science-ready catalog**

# Configuration Interface
## All decisions about the catalog content are made here

Input Data | **Configuration** | Summary

Selected config: **System default**

**Cluster Catalog**

☑ Query Builder

☑ Catalog Properties

Configuration

Save | Select | Share with users

Share with groups | Reset | Set as default

~50 configuration parameters!

General Information | **Region Selection** | Object Selection | Column Selection

▸ Mangle Detrac Map

▾ Bad Regions Mask

☐ 1 - Regions with bad astrometric colors

☑ 2 - Fainter 2MASS star region (8 < J < 12)

☑ 4 - Large nearby object (R3C catalog)

☑ 8 - Bright 2MASS star region (5 < J < 8)

☐ 16 - Near the LMC

☑ 32 - Yale Bright Star region

☐ 64 - High density of crazy colors

☑ 128 - Globular Clusters (William et al. 2010)

▸ Depth Map

▸ Systematic Maps

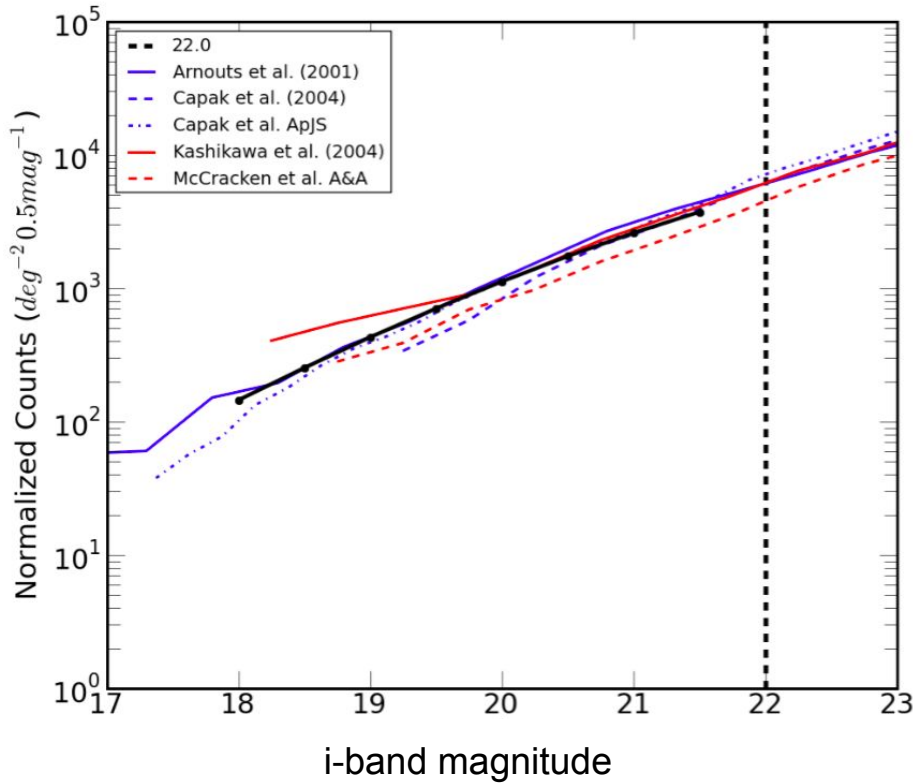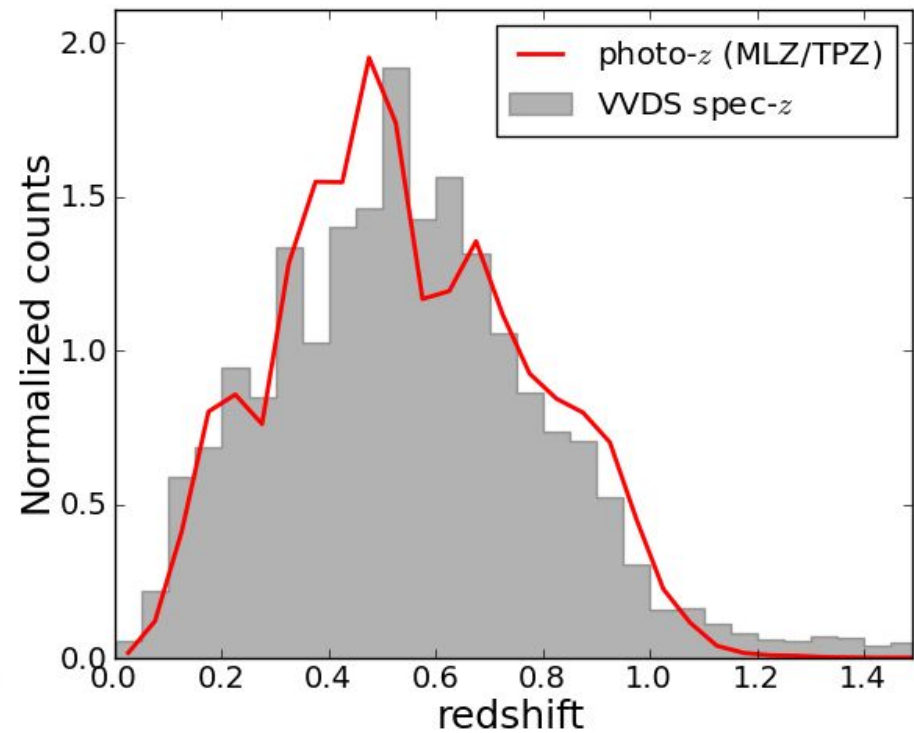▸ Additional Mask

Next

# Catalog properties

Number counts

photo-z and spec-z distribution

# Scientific Workflows @ LIneA

# Example of a cluster finder workflow

**WAZP**

Process ID: 10024456

**Process Summary** | **Results** | **Comments**

| Input Data | |
|---|---|
| Data Release | Y1A1 |
| Data Set | STRIPE82 |
| Value-Added Catalogs | Cluster Catalog |

| Output Data | |
|---|---|
| Targets | Cluster Members 1 |
| Targets | Galaxy Clusters 1 |

| Process Information | |
|---|---|
| Stage | None |
| Process ID | 10024456 |
| User | Cristiano Singulani |
| Start | 2016-09-01 17:22:07 |
| End | 2016-09-01 20:41:43 |
| Execution Time | 03:19:36 |
| Expiration Date | 2016-09-08 20:41:43 |
| Size | 29169420 |
| Status | Success |
| Overall Success Rate | 100% |
| Total Number of Jobs | 288 |
| Time Profiler | 🔍 |

| Module | Duration | Config | Error Log | Pipeline Out | Log | Condor Log | NC | Success Rate | Status |
|---|---|---|---|---|---|---|---|---|---|
| Slicing | 0:00:05 | 🔧 | - | - | - | - | ▫ | 100% | ✔ |
| Split Area | 0:00:16 | - | - | - | - | - | ▫ | 100% | ✔ |
| Visibility Maps | 0:01:36 | - | - | - | - | - | ▫ | 100% | ✔ |
| Background Model | 0:09:47 | - | - | - | - | - | ▫ | 100% | ✔ |
| WAZP per tile | 2:28:07 | - | - | - | - | - | ▫ | 100% | ✔ |
| Merge Results | 0:00:10 | - | - | - | - | - | ▫ | 100% | ✔ |

# How do scientists and developers interact?

Adding a new workflow to the Science Portal

Example 1: GE Science Workflow
    1 Description
    2 Contact points
    3 Science Code
    4 Pipeline definition
    5 Input Data
    6 Configuration Parameters
    7 Output Data
    8 Design of the Process Log
    9 Schedule
    10 Comunication tools

Example of specification document

# Challenges I - Data Processing

- DES Year 1 (~1500 sq deg)  objects catalog ~400G
- ~140M objects e ~600 attributes
- Data partitioned in 3,703 files ~100-150M  (DES tile)
- Data access during processing
  - PostgreSQL DB  ✕
  - Lustre File System  ✕
  - Hadoop File System  (local data processing)  ✓
- DES Year 3 release (Setembro 2016)
  - Objects catalog > 1TB partitioned in 10k files

# Challenges I - Data processing
## Parallel and distributed processing (similar to Map-Reduce)



Implementation problems
- Moving data during processing
- Consolidation of results in a single step

# Challenges II - How to use other resources available?

- SDumont (LNCC/Brazil), FermiGrid , Blue Waters, NERSC
  - Different environments:  PBS, Condor, SLURM, Condor-g
  - Data movement
  - "Big software" complex software and dependencies
- **Science-as-a-service**
  - Science APIs (iPlant/CyVerse)
  - Science Gateways (NERSC)
- **Portability**
  - cloud processing (private or public clouds)
- Federation of private clouds
  - UFCG Distributed systems lab (Francisco Brasileiro)
- OpenStack, AWS
  - LSST/SQuaRE (Frossie Economou)

# Challenges III - Data access and distribution

- Optimized data transfer (RNP/Brazil)

- Science portal integration with DES Science DB (NCSA)

  - Data access interface (large variety of products)

  - Documentation

- Science Server  DES DR1

- LSST/DM Data Acess Center (DAC) prototype at  LIneA

# Conclusions and perspectives

- LIneA: support brazilian participation in DES, SDSS, DESI and LSST

- Science Portal: necessary infrastructure for efficient science analysis

- DES Public DR1 (2017)

- LSST first light in 2020 start operations in 2022

- DES as a "prototype" for LSST

# Extra Slides

# The Science Portal and DES

http://www.linea.gov.br/

- Software development started in 2007*
- 9 years! about 56 FTEs
- 8 international reviews

| Emphasis | When | Where |
|---|---|---|
| Introduction, Science Workflows | Oct 2010 | Fermilab |
| Precam, Quick Reduce, Science Workflows | Oct 2011 | UPenn |
| Quick Reduce | May 2012 | MPA |
| End-to-end vision, data preparation | Jul 2013 | Fermilab |
| Data validation | Nov 2013 | Fermilab |
| Data validation and exploration | Ago 2014 | Fermilab |
| Data validation, exploration and science-ready catalogs | Nov 2014 | NCSA |
| Science-ready catalogs | May 2015 | Fermilab |

\* https://youtu.be/1Qv8HOoeUF4

# Monitoring the execution of all processes involved

Release: Y1A1 ▼    Dataset: SPT ▼

## Data Instalation

| Pipeline | Start | Duration | Runs | Status |
|---|---|---|---|---|
| Install Catalogs | 2016-03-08 15:40:13 | 01:51:56 | 1 | 🟢 |
| Install Mangle Mask | 2016-06-10 10:21:14 | 05:49:15 | 3 | 🟢 |
| Install Bright Mask | 2016-06-27 13:20:37 | 00:01:22 | 4 | 🟢 |
| Install Depth Maps | 2016-06-10 10:24:11 | 01:09:19 | 2 | 🟢 |
| Systematic Maps | 2016-06-13 12:47:35 | 12:43:31 | 4 | 🟢 |
| Zeropoint Correction | 2016-08-11 13:13:51 | 05:32:55 | 5 | 🔴 |
| QA Coadd | | | | ⚪ |
| | | Total: 27:8:17 | | |

## Data Preparation

| Pipeline | Start | Duration | Runs | Status |
|---|---|---|---|---|
| SG Separation | 2016-05-25 13:35:42 | 02:37:35 | 3 | 🟢 |
| Spectroscopic Sample | 2016-08-08 10:19:51 | 00:03:47 | 27 | 🟢 |
| Training Set Maker | 2016-07-20 10:40:48 | 01:35:41 | 6 | 🟢 |
| Photo-z Training | 2016-06-27 10:17:59 | 03:26:39 | 2 | 🟢 |
| Photo-z Compute | 2016-06-14 16:29:09 | 02:36:11 | 13 | 🟢 |
| Galaxy Properties | 2016-07-13 15:16:10 | 10:38:08 | 2 | 🔴 |
| | | Total: 20:57:0 | | |

## Science-ready Catalogs

| Pipeline | Start | Duration | Runs | Status |
|---|---|---|---|---|
| Cluster | 2016-08-07 17:38:21 | 02:45:44 | 25 | 🟢 |
| GE | 2016-05-17 14:40:45 | 01:52:37 | 1 | 🟢 |
| GA | 2016-05-24 10:58:30 | 01:15:09 | 2 | 🟢 |
| | | Total: 5:53:30 | | |

## Stages

- Data Installation
- Data Preparation
- Science-ready catalogs

- 16 workflows
- 64 data products