



<http://ls.st/fqn>

# The Challenge of Large Dataset Analysis

---

*or why I believe exciting times are ahead*

**Prof. Mario Juric**

DIRAC Institute | eScience Institute | UW Astronomy

DATA INTENSIVE RESEARCH IN ASTROPHYSICS AND COSMOLOGY  
COLLEGE OF ARTS & SCIENCES | UNIVERSITY of WASHINGTON



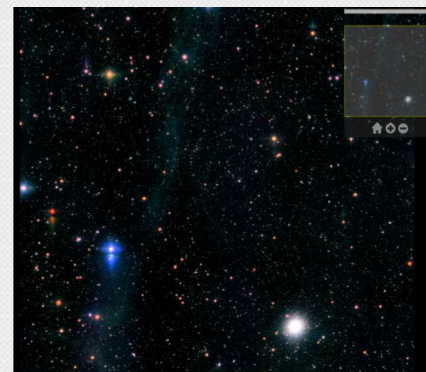
# Synopsis

- How to think about survey datasets
- What do surveys deliver (with LSST as an example)
- How will we analyze PB-scale datasets
- The possibilities beyond the current paradigm

# Telescopes as just (Expensive) Peripherals



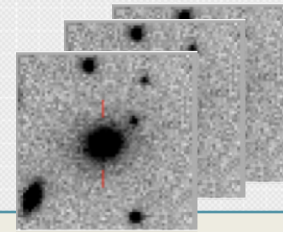
Peripherals



Datasets

# This is the Important Bit: From Data to Knowledge

Computationally (and cognitively) expensive, science-case specific



And metadata!

*Scientists* **Model** ← *inference* – **Data** *Projects*

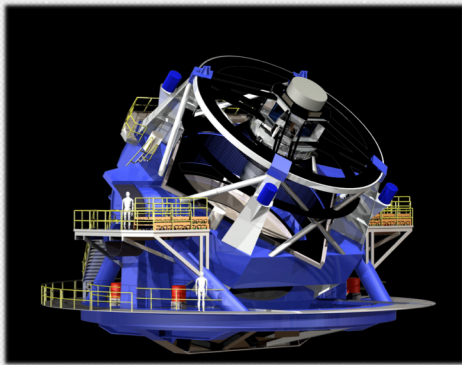
*Scientists* **Model** ← *inference* – *Scientists* **Catalog** ← *Projects* **Data Processing** – *Projects* **Data**

Computationally cheaper,  
Easier to understand,  
Science-case specific

- Computationally expensive, general
- Reprojection; may or may not involve compression
- Almost always introduces some information loss
- Data Processing == Instrumental Calibration + Measurement

# Surveys: Turning the Sky into a Databases

- The ultimate deliverable of a survey is not the telescope, nor the instruments; it is the fully reduced data.
- All science comes from catalogs and images
- The telescope is still an (expensive) data collection peripheral



Peripherals



Code & Machines



Databases

Table 4: Level 2 Catalog Object Table

Name	Type	Unit	Description
psRadcTai	double	time	Point source model: Time at which the object was at position <code>radec</code> .
psPm	float[2]	mas/yr	Point source model: Proper motion vector.
psParallax	float	mas	Point source model: Parallax.
psFlux	float[ugrizy]	nmgy	Point source model fluxes <sup>58</sup> .
psCov	float[66]	various	Point-source model covariance matrix <sup>59</sup> .
psLnL	float		Natural <i>log</i> likelihood of the observed data given the point source model.
bdRadc	double[2]	degrees	B+D model <sup>60</sup> : $(\alpha, \delta)$ position of the object at time <code>radecTai</code> , in each band.

# ***#1 Challenge:***

*General purpose processing while  
minimizing information loss.*

# Guiding Principles for LSST Data Products

- There are virtually infinite options on what quantities (features) one can measure on images. But if catalog generation is understood as a (generalized) cost reduction tool, the guiding principles become easier to define:
  - 1. Maximize science enabled by the catalogs**
    - Working with images takes time and resources; a large fraction of science cases should be enabled by just the catalog.
    - Be considerate to the user: provide even sub-optimal measurements if they will enable leveraging of existing experience and tools
  - 2. Minimize information loss**
    - Choose good models
    - Provide (as much as possible) estimates of likelihood surfaces, not just single point estimators
  - 3. Provide and document the transformation (the software)**
    - Measurements are becoming increasingly complex and systematics limited; need to be maximally transparent about how they're done

# What LSST will Deliver:

## A Data Stream and a Database

- A stream of  $\sim 10$  million time-domain events per night, detected and transmitted to event distribution networks within 60 seconds of observation.
- A catalog of orbits for  $\sim 6$  million bodies in the Solar System.
- A catalog of  $\sim 37$  billion objects (20B galaxies, 17B stars),  $\sim 7$  trillion single-epoch detections (“sources”), and  $\sim 30$  trillion forced sources, produced annually, accessible through online databases.
- Deep co-added images.

Prompt

Data Rel.



## Prompt: Time-Domain Event Alerts

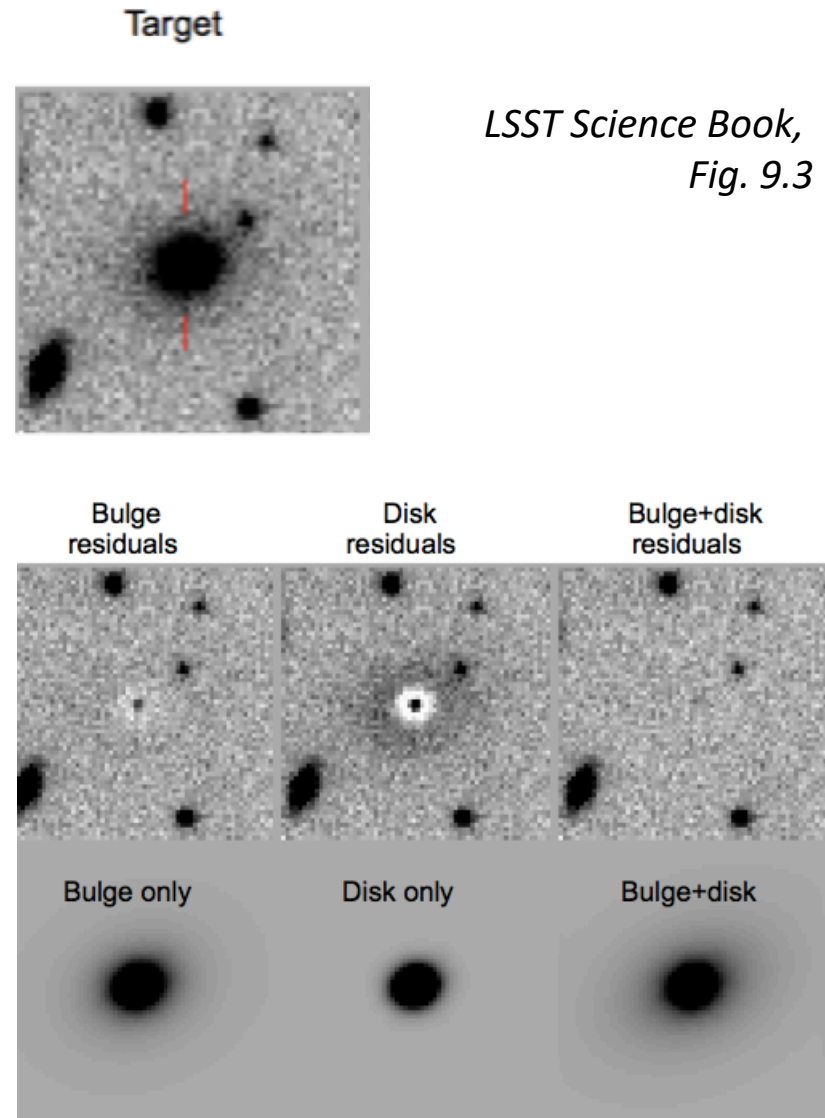
- We expect a high rate of alerts, **approaching 10 million per night**. We'll also provide an ***alert filtering service***, to select subsets of alerts, as well as serve the full stream to external ***event brokers***.
- **Each alert will include the following:**
  - **Alert and database ID**: IDs uniquely identifying this alert.
  - The photometric, astrometric, and shape characterization of the detected source
  - 30x30 pixel (on average) **cut-out of the difference image** (FITS)
  - 30x30 pixel (on average) **cut-out of the template image** (FITS)
  - The time series (up to a year) of all previous detections of this source
  - Various summary statistics (“features”) computed of the time series
- **The goal is to quickly transmit nearly everything LSST knows about any given event, enabling downstream classification and decision making**

# Annual Data Releases

- **Made available in *Data Releases***
  - Annually, except for Year 1
    - Two DRs for the first year of data
- **Well calibrated, consistently processed, catalogs and images**
  - Catalogs of objects, detections, detections in difference images, etc.
- **Complete reprocessing of all data, for each release**
  - Every DR will reprocess *all* data taken up to the beginning of that DR
- **Projected catalog sizes:**
  - **18 billion objects** (DR1) → **37 billion** (DR11)
  - **750 billion observations** (DR1) → **30 trillion** (DR11)

# Data Release Catalog Contents

- **Object characterization (models):**
  - Moving Point Source model
  - Double Sérsic model (bulge+disk)
    - Maximum likelihood peak
    - Samples of the posterior (hundreds)
- **Object characterization (non-parametric):**
  - Centroid: ( $\alpha$ ,  $\delta$ ), per band
  - Adaptive moments and ellipticity measures (per band)
  - Aperture fluxes and Petrosian and Kron fluxes and radii (per band)
- **Colors:**
  - Seeing-independent measure of object color
- **Variability statistics:**
  - Period, low-order light-curve moments, etc.



# Analysis: Subset – Download – Analyze

**SDSS Query / CasJobs**

Help Tools Query History MyDB Import Groups Output Profile Queues SkyServer Logout myunc

**'My Query' Details**

JobID	TaskName	Context	Queue	Submitted	Started	Finished	Status
4930248	My Query	DR7	600	7/28/2010 8:04:02 AM	7/28/2010 8:04:10 AM	7/28/2010 8:09:21 AM	Finished

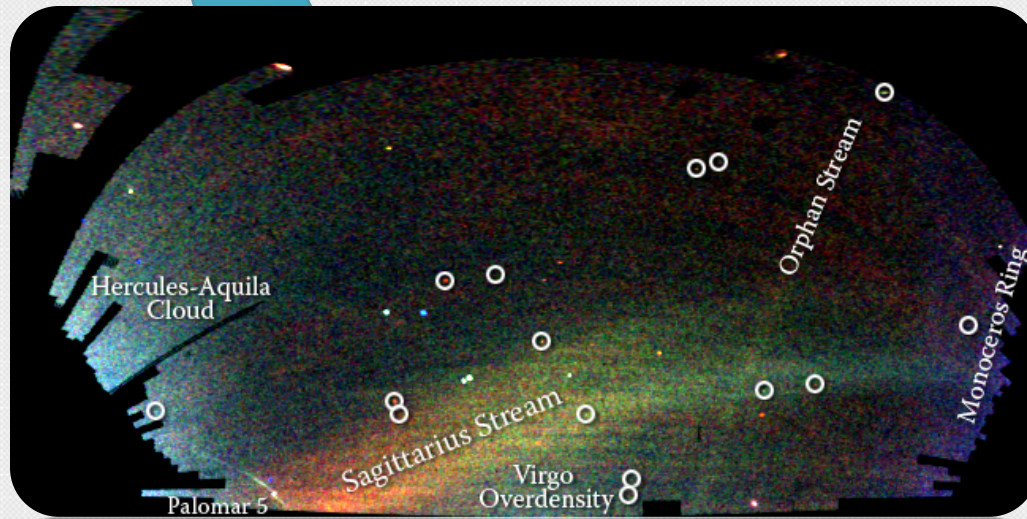
Executed on Rows: 600 Message: Query Complete

```

Query
-----
select
  p.objid,
  p.ra,
  p.dec,
  p.wavelength_id as W,
  p.dered_u as U, p.dered_g as G, p.dered_r as R, p.dered_i as I, p.dered_z as Z,
  p.err_u, p.err_g, p.err_r, p.err_i, p.err_z,
  p.flags
from
  Star s
into
  mpdb.north1cc3
where
  b > 40
  and 0 <= c - r and r - i < 0.4
  and 21 <= z and z < 21.1
    
```

Contact: casjobs@casjobs.org (Revision: 1.23.9, Last modified: Wednesday, May 14, 2008 at 1:02:25 AM)

JobID	TaskName	Context	Queue	Submitted	Started	Finished	Status		
7757_301	1.74	186	6	8.12783867256709	26.62724537921	17.37402	37.92875	0.03894481	0.02558813
7757_301	1.74	188	3	8.12732322524192	26.6251199416623	20.1466	21.35297	0.3803744	0.3302762
4288_301	1.39	682	3	24.5161178422305	-1.16579446393527	22.97832	24.3259	0.2672399	0.5248437
4288_301	1.39	683	3	24.5179406515354	-1.1792069022485	22.62052	25.09109	0.1850479	0.6585805
4288_301	1.39	684	6	24.5189452933148	-1.1691508180898	21.4247	23.04125	0.06608535	0.1968172
4136_301	1.61	935	6	36.471592759892	-1.06093938828308	27.71762	23.14112	0.1580014	0.1796867
4136_301	1.61	936	3	36.4717583011316	-1.1378448207726	22.81683	23.88123	0.1742272	0.3269605
4136_301	1.61	937	3	36.4717582434391	-1.13784497192974	22.81147	23.87586	0.1734457	0.3247895
4288_301	1.48	311	3	24.689203338022	-1.22631696217547	21.2002	21.67521	0.05694564	0.06316777
4288_301	1.48	312	3	24.68468692692246	-1.21704918362007	20.30207	21.04976	0.02927161	0.04032911
4288_301	1.48	313	6	24.6848016690377	-1.08292772289886	24.92263	25.72778	0.68427	0.5471938
5598_301	1.61	792	3	35.1787950113407	6.14573538435867	22.43574	23.83793	0.1125768	0.2867745
5598_301	1.61	793	3	35.1787950113434	6.14573538316393	22.43573	23.77753	0.1174541	0.2741985
5598_301	1.61	794	6	35.1787349107839	6.14481612145222	24.6701	24.8507	0.4675894	0.4049984
2699_301	1.48	597	3	12.0768019816408	-3.32677418219699	20.18116	23.27577	0.126546	0.2278369
2699_301	1.48	598	6	12.0770027529666	-3.32913243320258	22.12757	23.79366	0.1215217	0.3403472
2699_301	1.48	599	3	12.0832728187538	-3.52539818226738	22.29741	23.32808	0.1377919	0.2253349
94_301	1.38	279	6	348.524768659138	-0.043098883870374	20.75028	21.13888	0.04839022	0.04444739
94_301	1.38	281	6	348.53257887691	-1.02365035629542	21.10668	21.15248	0.06605724	0.04672166
4288_301	1.76	766	3	30.0899167300738	-1.25189466355601	22.57054	22.91144	0.180731	0.1984752
4288_301	1.76	767	3	30.0899962731195	-1.14314301954175	21.99419	23.78533	0.1094794	0.4162939
4288_301	1.76	768	6	30.0899156247035	-1.19236549812758	22.64196	22.91776	0.1927483	0.2022232
7937_301	1.84	354	3	5.9312226470183	26.61206380856974	22.63659	23.76799	0.3950225	0.5983699
7937_301	1.84	355	3	5.93084744314193	26.6157353187021	21.27009	22.13992	0.08000753	0.1490507
7937_301	1.84	356	6	5.93087089167375	26.5703384153732	24.52393	26.33231	0.066107	0.6528299
3996_301	1.81	615	6	216.228187662076	11.9472286471793	19.93735	21.0716	0.0294876	0.03470354
3996_301	1.81	616	3	216.23835672189	12.082424512804	20.23574	21.74497	0.02392315	0.03942162
3996_301	1.81	617	3	216.2193792328008	11.9171482588482	22.68762	23.06279	0.1416738	0.1555459
6354_301	1.29	2997	6	332.302798261878	41.2179102291468	20.95895	22.51458	0.04065189	0.1163793
6354_301	1.29	2998	6	332.355206585822	41.2447081198007	20.78153	22.35575	0.03523029	0.09754675
6354_301	1.29	2999	3	332.30785409452	41.2613114989108	21.08566	22.75478	0.04217106	0.1333208
2986_301	1.59	251	6	127.586395164415	2.97627035997022	24.42951	26.40598	0.053675	0.0810285
2986_301	1.59	252	6	127.585992707459	2.98207493069978	25.42876	25.06479	0.5993609	0.8693208
2986_301	1.59	253	6	127.588133654789	2.98377339888094	22.25157	25.16562	0.0720802	0.7840060
5598_301	1.31	1254	6	347.289019727991	6.00342438896053	24.31121	25.55253	0.4780855	0.4634795
5598_301	1.31	1255	6	347.292686860266	6.04280417791226	25.21082	25.48026	0.4824083	0.473972
5598_301	1.31	1256	6	347.294399617648	6.08023292496797	24.812	25.48199	0.5272527	0.4741035



# Data Volumes

	ZTF	LSST
Number of detections	1 trillion	7 trillion
Number of objects	1 billion	37 billion
Nightly alert rate	1 million	10 million
Nightly data rate	1.4 TB	15 TB
Alert latency	< 20 minutes	60 seconds

Science analysis code

~50kb

***If the data is big...***

***... bring the code to the data.***

# What LSST will Deliver:

## A Data Stream, a Database, and a (small) Cloud

- A stream of  $\sim 10$  million time-domain events per night, detected and transmitted to event distribution networks within 60 seconds of observation.
- A catalog of orbits for  $\sim 6$  million bodies in the Solar System.
- A catalog of  $\sim 37$  billion objects (20B galaxies, 17B stars),  $\sim 7$  trillion single-epoch detections (“sources”), and  $\sim 30$  trillion forced sources, produced annually, accessible through online databases.
- Deep co-added images.
- Services and computing resources at the Data Access Centers to enable end-user analysis and generation of more added-value data products.

Prompt

Data Rel.

User Gen.

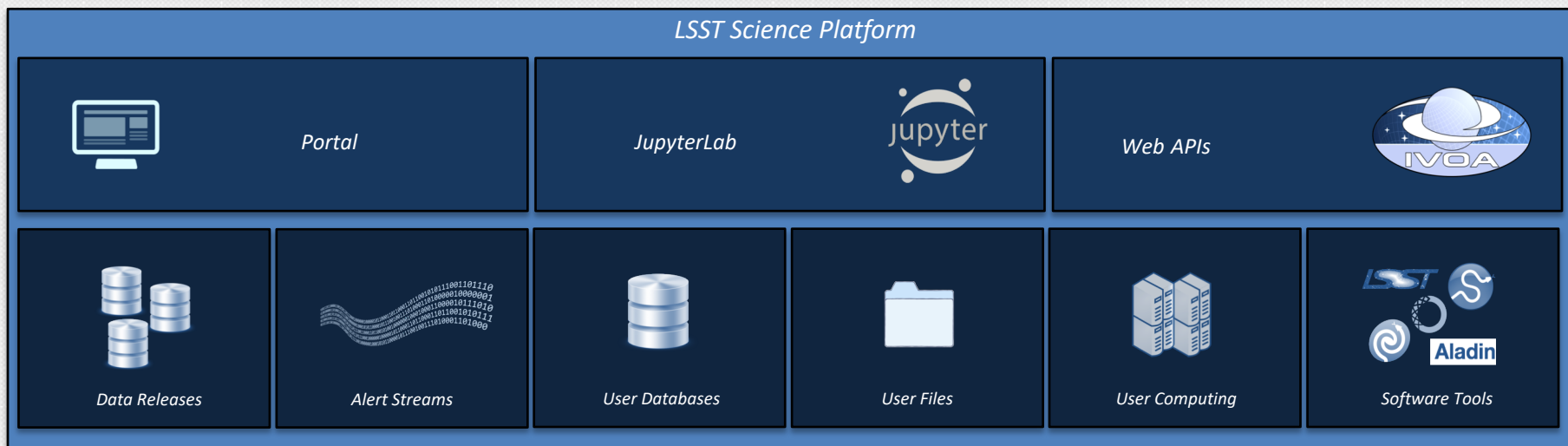
# The LSST Science Platform: Accessing LSST Data and Enabling LSST Science



LSST Users

Internet

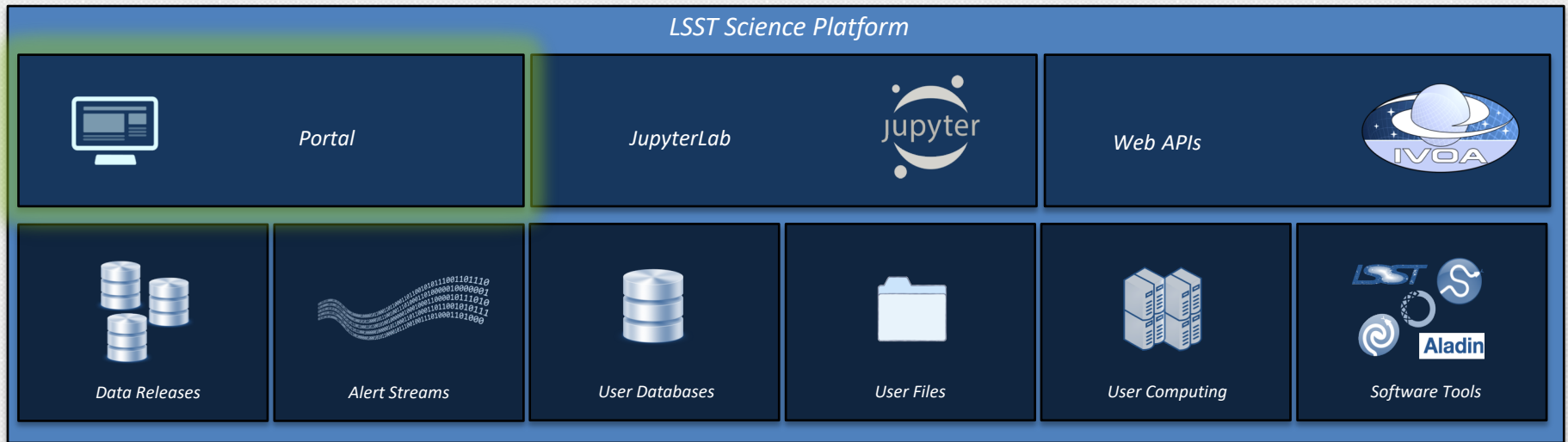
LSST Science Platform



The **LSST Science Platform** is a set of integrated web applications and services deployed at the LSST Data Access Centers (DACs) through which the scientific community will access, visualize, subset and perform next-to-the-data analysis of the data.



# LSST Portal: The Web Window into the LSST Archive

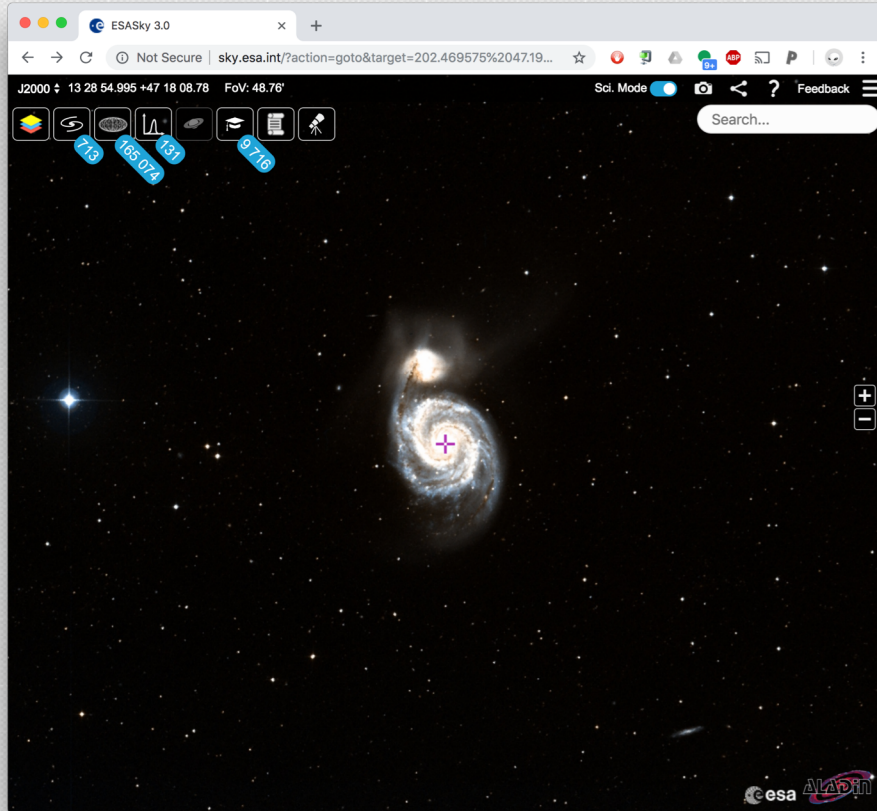


The Web Portal to the archive will enable browsing and visualization of the available datasets in ways the users are accustomed to at archives such as IRSA, MAST, or the SDSS archive, with an added level of interactivity.

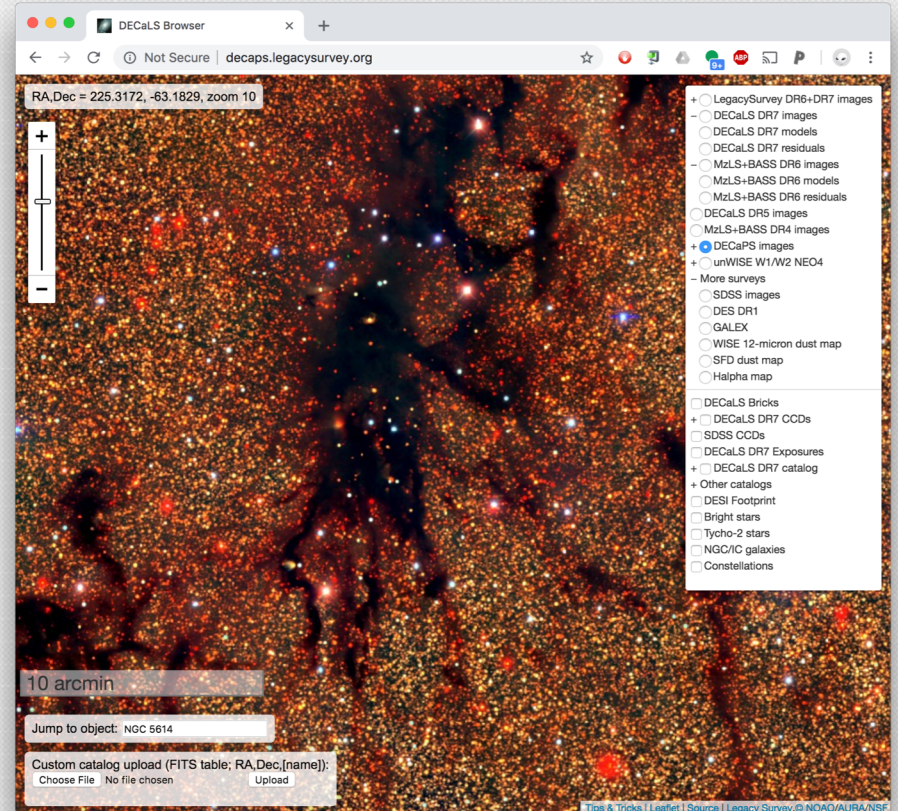
Through the Portal, the users will be able to view the LSST images, request subsets of data (via simple forms or SQL queries), construct simple plots, and generally explore the LSST dataset in a way that allows them to identify and access (subsets of) data required by their science case.

This will all be backed by a petascale-capable RDBMS.

# What to expect

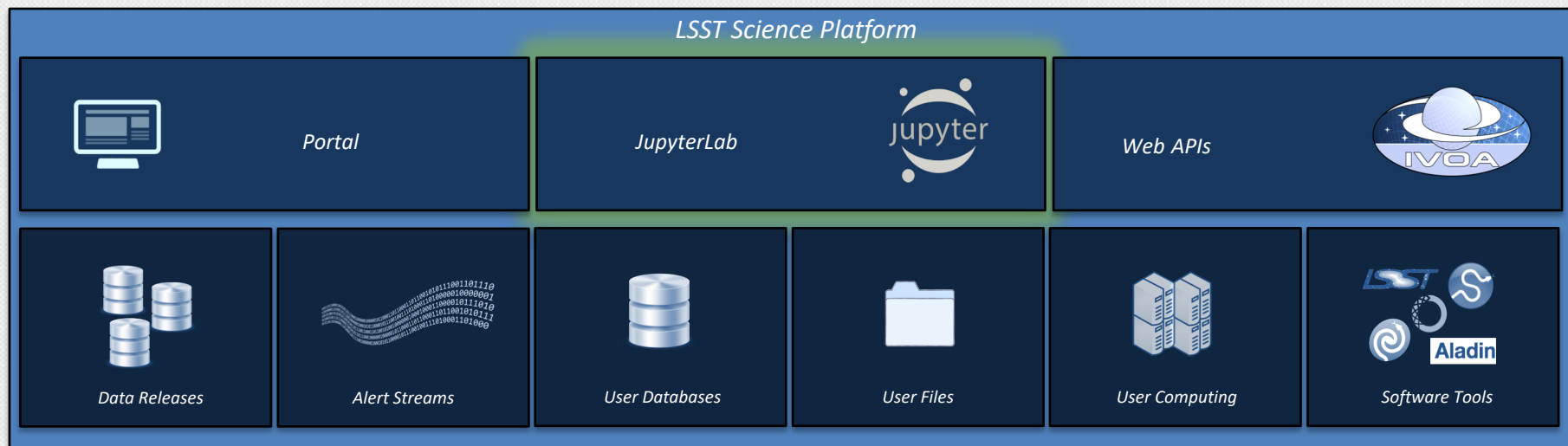


<http://sky.esa.int/>



<http://decaps.legacysurvey.org/>

# JupyterLab: Next-to-the-data Analysis



The tools exposed through the Web Portal will permit simple exploration, subsetting, and visualization of LSST data. They may not, however, be suitable for more complex data selection or analysis tasks.

To enable that next level of next-to-the-data work, we plan to enable the users to launch their own Jupyter notebooks at our computing resources at the DAC. These will have fast access to the LSST database and files. They will come with commonly used and useful tools preinstalled (e.g., AstroPy, LSST data processing software stack).

This service is similar in nature to efforts such as SciServer at JHU.

# JupyterLab: Next-to-the-data Analysis

The image displays a JupyterLab environment. On the left, a file browser shows a directory structure with files like 'analysis', 'singlechip\_sample', and 'test.ipynb'. The main area is a code editor with a file named 'test.ipynb' open. The code in the editor includes a table of file statistics, a list of file paths, and two code cells. The first cell defines a Butler object and retrieves data. The second cell imports 'lsst.afw.geom' and performs geometric operations. Below the code, a scatter plot shows astronomical data points on a 2D plane, with axes ranging from 0 to 500. A mouse cursor is visible over the plot. On the right, a browser window shows the JupyterLab interface, including a 'Logout' button and a 'Control Panel' link. The browser address bar shows the URL: 'https://epyc.astro.washington.edu/jupyter/user/mjuric/tree/ep...'. The browser window also displays a file browser with a list of folders and files, including 'axes-documentation', 'hack\_day', 'lsd', 'lsd-archive', 'lsd2', 'mops-iod', 'plasticc-ml', 'sssc', 'thor', 'trilegal', 'ztf-alerts', 'ztf-conda', 'ztf-jupyter', 'ztf\_experiments', and 'DEADJOE'.

YouTube demo of the LSST JupyterLab Aspect Demo: <http://ls.st/bgt>

Why Jupyter is data scientists' x +  
https://www.nature.com/articles/d41586-018-07196-1  
nature > toolbox > article  
a nature research journal

TOOLBOX · 30 OCTOBER 2018

# Why Jupyter is data scientists' computational notebook of choice

*An improved architecture and enthusiastic user base are driving uptake of the open-source web tool.*

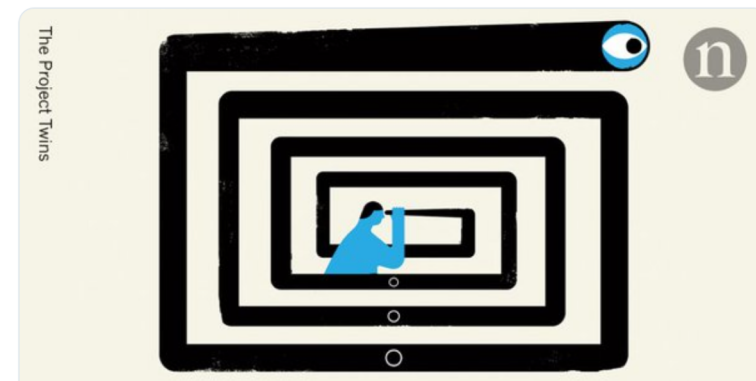
Jeffrey M. Perkel



**Nature Careers**  
@NatureCareers

Follow

"I've never seen any migration this fast. It's just amazing." -- @mjuric on the rise of @ProjectJupyter in data science



**Why Jupyter is data scientists' computational notebook of choice**

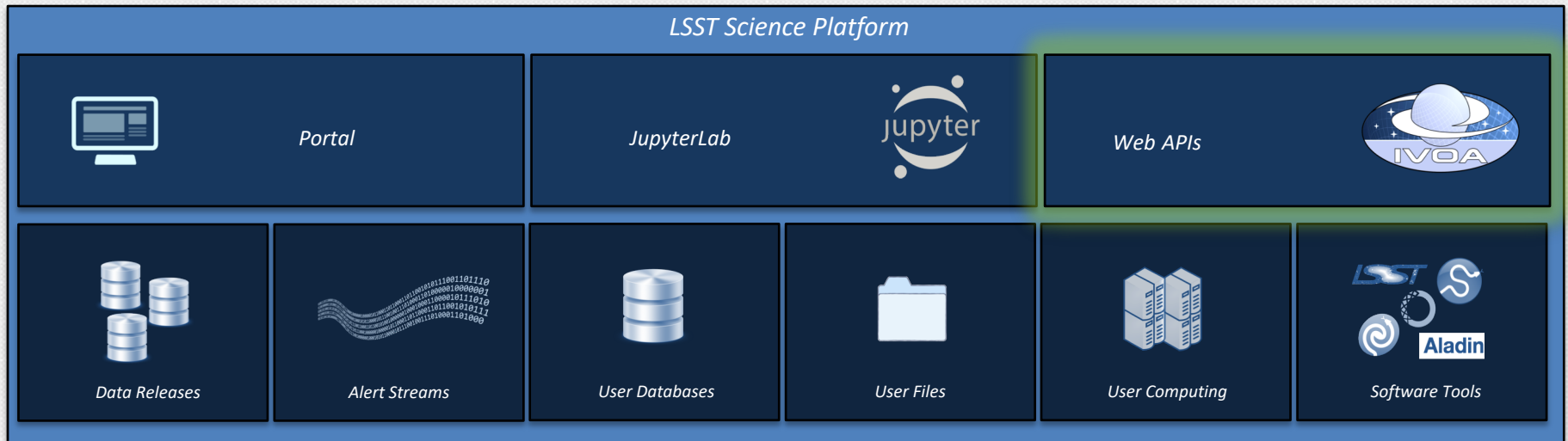
An improved architecture and enthusiastic user base are driving uptake of the open-source web tool.

[nature.com](https://www.nature.com)

10:30 PM - 5 Nov 2018



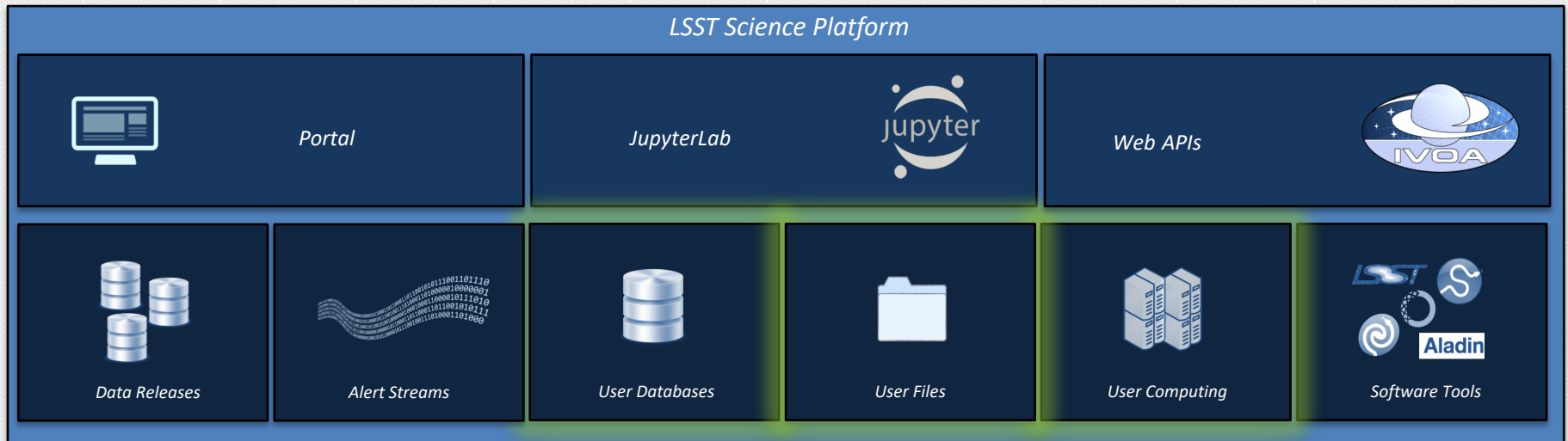
# Web APIs: Integrating With Existing Tools



Backend Platform services – such as access to databases, images, and other files – will be exposed through machine-accessible web APIs.

We have a preference for industry standard and/or VO APIs (e.g., WebDAV, TAP, SIA, etc.) – the goal is to support what’s broadly accepted within the community. This will allow the discoverability of LSST data products from within the Virtual Observatory, federation of the LSST data set to other archives, and the use of widely utilized tools (eg., TOPCAT or others).

# Computing, Storage, and Database Resources



Computing, file storage, and personal databases (the “*user workspace*”) will be made available to support the work via the Portal and within the Notebooks.

An important feature is that no matter how the user accesses the DAC (Portal, Notebook, or VO APIs) they always “see” the same workspace.

# How big is the “LSST Science Cloud” (@ DR2)?

## – Computing:

- ~2,400 cores
- ~18 TFLOPs

**This is shared by all users.** We’re estimating the number of potential DAC users not to exceed 7500 (relevant for file and database storage).

Not all users will be accessing the computing cluster concurrently. **We are estimating on order of a ~100.**

## – File storage:

- ~4 PB

Though this is a relatively small cluster by 2020-era standards, it will be **sufficient to enable preliminary end-user science analyses** (working on catalogs, smaller number of images) and creation of some added-value (Level 3) data products.

## – Database storage

- ~3 PB

*Think of this as having your own server with a few TB of disk and database storage, right next to the LSST data, with a chance to use tens to hundreds of cores for analysis.*

This kind of approach will become increasingly common for *all* big data archives.



# Already Here for Gaia: GAVIP

The screenshot shows a web browser window with the GAVIP Portal. The browser address bar shows <https://gavip.esac.esa.int>. The page header includes "GAVIP v1.1.2" and "UTC: 2018/11/8 5:52:34".

**Sign in or Register**

Buttons: Sign in, Register

Or use a temporary (limited) account

Continue as anonymous user.

## WELCOME TO GAVIP

**Gaia Added Value Interface Platform**  
A science-exploitation platform for the Gaia data archive

### About Gaia

Gaia is a major ESA mission designed to survey 1,000,000,000 (1 billion) stars in the Milky Way in order to make a 3D map of unprecedented accuracy. This will allow us to answer fundamental questions about the origin and evolution of our galaxy.

However, the data-products of Gaia are expected to surpass 1 PetaByte in size, making it difficult for the scientific community to engage and analyse the data.

The Gaia Added Value Interface Platform (GAVIP) is designed to address this issue by providing an innovative platform for scientists to re-use and deploy, close to the data, their own existing code, packaged as "Added Value Interfaces" (AVIs).

**Tweets by @ESAGaia**

**ESA Gaia** @ESAGaia  
More #GaiaScience: secrets of the "Wild Duck Cluster" revealed! Read through our #IoW story: [cosmos.esa.int/web/gaia/iow\\_2...](https://cosmos.esa.int/web/gaia/iow_2...) and discover

### Platform Features

- Jupyter**  
Use hosted Jupyter notebooks to run code next to the archive, then evolve that into an AVI for others.
- AVIs**  
Turn your analysis into a reusable tool for the world to use! Create arbitrarily complex pipelines and interfaces using the AVI framework.
- Isolation**  
GAVIP uses [Docker](#) to isolate and run AVIs on demand within the platform. It's an AVI just for you.
- AVI framework**  
Resource management and asynchronous pipeline execution are **automatically managed**. Interfaces to GACS, SAMP, and more are provided for your pipeline.

## Challenges (part 1)

Better Together

(joining datasets is powerful)

I Want it All

(science demands whole dataset operations)



# 3D Dust Mapping

## with Pan-STARRS 1

Query Map



Usage Notes



Read Papers



## The Map

Interstellar dust attenuates ultraviolet, optical and near-infrared light. Because the extent of this attenuation is wavelength-dependent, dust both dims and reddens the light of stars and galaxies before it can reach our telescopes. In many areas of astrophysics, an accurate correction for the effects of interstellar extinction and reddening is critical. Historically, the most widely used maps of dust have been two-dimensional, tracing integrated dust reddening out to infinite distance. Here, we describe three-dimensional maps of interstellar dust reddening, which trace dust reddening both as a function of angular position on the sky and distance. These dust maps are based on [Pan-STARRS 1](#) photometry of 800 million stars, along with [2MASS](#) photometry of 200 million stars.

To read about how to download the map, or how to query it remotely, read our [usage notes](#). To explore our map in the browser, see our [interactive query page](#). To read in detail about our map, read our [published papers](#).

# Whole Dataset Operations

- Galactic structure: density/proper motion maps of the Galaxy
  - => forall stars, compute distance, bin, create 5D map
- Galactic structure: dust distribution
  - => forall stars, compute g-r color, bin, find blue tip edge, infer dust distribution
- Near-field cosmology: MW satellite searches
  - => forall stars, compute colors, convolve with spatial filters, report any satellite-like peaks
- Variability: Bayesian classification of transients and discovery of variables
  - => forall stars, get light curves, compute likelihoods, alert if interesting
- ...

## Challenges (part 2)

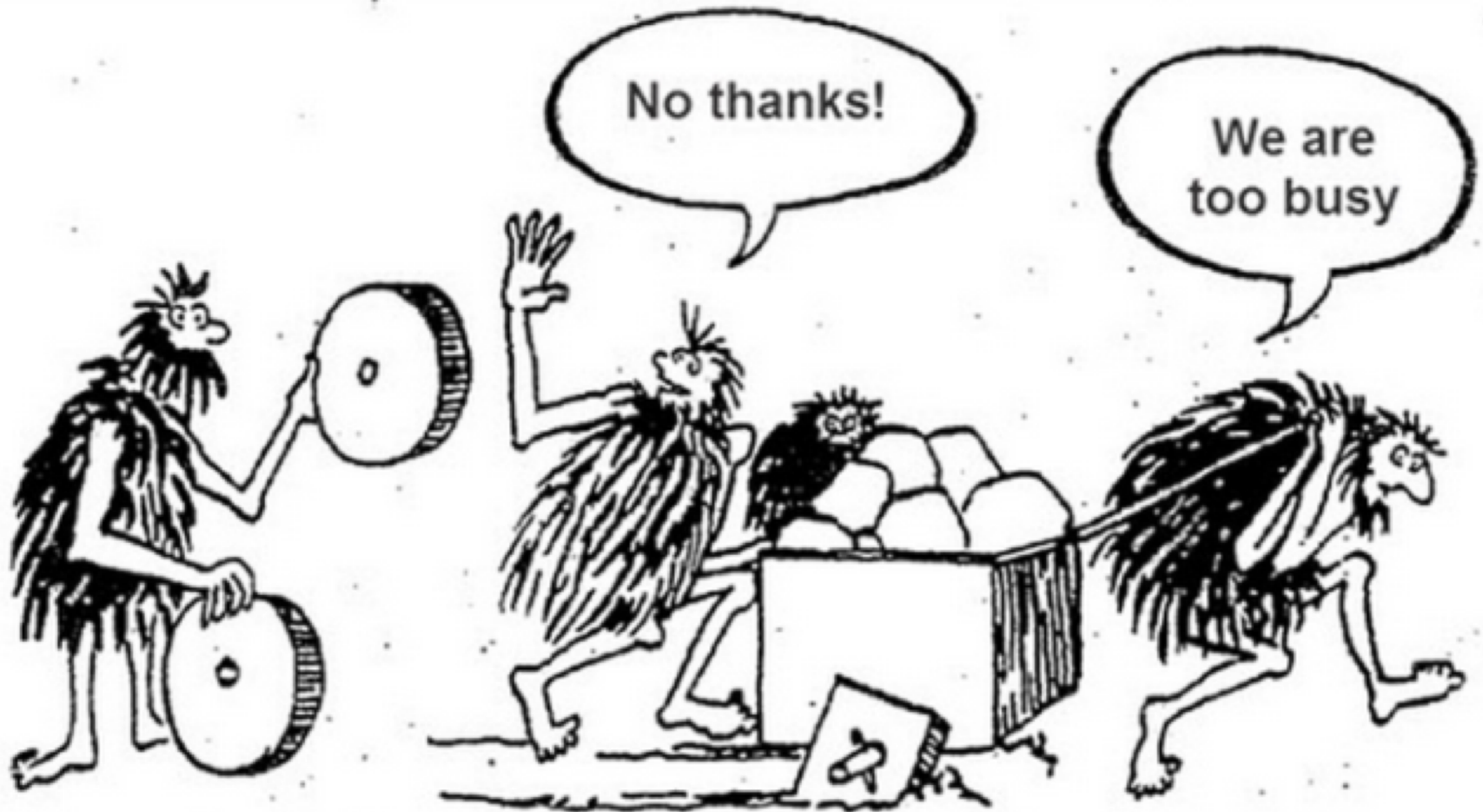
Scalability

(how do I write an analysis code that will scale to petabytes of data?)

Resources

(where are the resources to run this code?)

## Remember Yesterday...



# Writing Scalable Applications: MapReduce and Apache Spark

The screenshot shows the Apache Spark website homepage. At the top, there's a navigation bar with links for Download, Libraries, Documentation, Examples, Community, Developers, and Apache Software Foundation. Below this is a main heading: "Apache Spark™ is a unified analytics engine for large-scale data processing." To the left, under "Speed", it says "Run workloads 100x faster." and includes a bar chart comparing Hadoop (110s) and Spark (0.9s) for logistic regression. Below that, under "Ease of Use", it says "Write applications quickly in Java, Scala, Python, R, and SQL." and includes a code snippet for reading JSON files. On the right, there's a "Latest News" section with recent updates and an "APACHECON North America" event announcement for September 24-27, 2018 in Montréal, Canada, with a "Download Spark" button.

Apache Spark™ is a unified analytics engine for large-scale data processing.

### Speed

Run workloads 100x faster.

Apache Spark achieves high performance for both batch and streaming data, using a state-of-the-art DAG scheduler, a query optimizer, and a physical execution engine.

Framework	Running time (s)
Hadoop	110
Spark	0.9

Logistic regression in Hadoop and Spark

### Ease of Use

Write applications quickly in Java, Scala, Python, R, and SQL.

```
df = spark.read.json("logs.json")
df.where("age > 21")
  .select("name.first").show()
```

Spark's Python DataFrame API  
Read JSON files with automatic schema inference

### Latest News

- Spark 2.3.2 released (Sep 24, 2018)
- Spark+AI Summit (October 2-4th, 2018, London) agenda posted (Jul 24, 2018)
- Spark 2.2.2 released (Jul 02, 2018)
- Spark 2.1.3 released (Jun 29, 2018)

[Archive](#)

### APACHECON North America

September 24-27, 2018  
Montréal, Canada

[Download Spark](#)

### Built-in Libraries:

- [SQL and DataFrames](#)
- [Spark Streaming](#)
- [MLlib \(machine learning\)](#)
- [GraphX \(graph\)](#)

[Third-Party Projects](#)

Apache Spark is an open-source distributed general-purpose cluster-computing framework.

**Spark provides an interface for programming entire clusters with implicit data parallelism and fault tolerance.**

-- Wikipedia

## Examples

# Pi Estimation

Spark can also be used for compute-intensive tasks. This code estimates  $\pi$  by "throwing darts" at a circle. We pick random points in the unit square ((0, 0) to (1,1)) and see how many fall in the unit circle. The fraction should be  $\pi / 4$ , so we use this to get our estimate.

Python

Scala

Java

```
def inside(p):  
    x, y = random.random(), random.random()  
    return x*x + y*y < 1  
  
count = sc.parallelize(xrange(0, NUM_SAMPLES)) \  
    .filter(inside).count()  
print "Pi is roughly %f" % (4.0 * count / NUM_SAMPLES)
```

<https://spark.apache.org/examples.html>



Map

$\{x_i\} \text{ ---map--> } \{y_i=f(x_i)\}$

Apply a function  $f$  to every element of dataset  $X$ , producing dataset  $Y$

Reduce

$\{(k_i, v_{ij})\} \rightarrow \{y_i=(k_i, f(\{v_{ij}\}))\}$  Apply a function  $f$  to all values with a common key

Example:

$\{("dog", 2), ("dog", 1), ("cat", 3), ("dog", 2), ("cat", 1)\}$

-> reduce w.  $sum()$  ->

$\{("dog", 5), ("cat", 4)\}$

## Examples

```
{ ("dog", 2), ("dog", 1), ("cat", 3), ("dog", 2), ("cat", 1) }
```

-> reduce w. *sum()* ->

```
{ ("dog", 5), ("cat", 4) }
```

## Word Count

In this example, we use a few transformations to build a dataset of (String, Int) pairs called counts and then save it to a file.

Python

Scala

Java

```
text_file = sc.textFile("hdfs://...")
counts = text_file.flatMap(lambda line: line.split(" ")) \
    .map(lambda word: (word, 1)) \
    .reduceByKey(lambda a, b: a + b)
counts.saveAsTextFile("hdfs://...")
```

<https://spark.apache.org/examples.html>

# Astronomy Example: Compute Light Curve Features

*This works on arbitrarily large datasets!*

```
In [10]: from pyspark.sql.types import ArrayType, FloatType, DoubleType
from pyspark.sql.functions import col, pandas_udf, explode
import pandas as pd

import cesium
from cesium.time_series import TimeSeries
from cesium.featurize import featurize_single_ts, featurize_time_series

#####

features_to_use = ["amplitude", "percent_beyond_1_std", "maximum", "max_slope",
                  "median", "median_absolute_deviation", "percent_close_to_median",
                  "minimum", "skew", "std", "weighted_average"]

ls_features = ["freq1_amplitude1", "freq1_amplitude2", "freq1_amplitude3",
              "freq1_amplitude4", "freq1_freq", "freq1_lambda", "freq1_rel_phase2",
              "freq1_rel_phase3", "freq1_rel_phase4", "freq1_signif", "freq2_amplitude1",
              "freq2_amplitude2", "freq2_amplitude3", "freq2_amplitude4", "freq2_freq",
              "freq2_rel_phase2", "freq2_rel_phase3", "freq2_rel_phase4"]

def featurize_udf(mjd, psfflux):
    feat_outs = []
    for row_mjd, row_psfflux in zip(mjd, psfflux):
        feat_out = featurize_time_series(np.array(row_mjd), np.array(row_psfflux),
                                       features_to_use=features_to_use + ls_features)
        feat_outs.append(feat_out.values.flatten())
    return pd.Series(feat_outs)

#####

feat_udf = pandas_udf(featurize_udf, returnType = ArrayType(DoubleType()))
spark_session.udf.register("FEATURIZE", feat_udf)

pdf = ztf.where("SIZE(mjd)>50").selectExpr("FEATURIZE(mjd, psfflux)").toPandas()
```



# Want to try it out?

```
conda install -c conda-forge pyspark
```



# Scaling with Spark

Today, Spark is being adopted by major players like Amazon, eBay, and Yahoo! Many organizations run Spark on clusters with thousands of nodes. According to the Spark FAQ, the largest known cluster has over 8000 nodes. Indeed, Spark is a technology well worth taking note of and learning about.



<https://www.toptal.com/spark/introduction-to-apache-spark>



## Cloud services

- Essentially, companies who rent computers (or a few million of them)
  - The same for storage
- Pay only for what you use (by the second/minute/hour)
- Scalable: ask for 1000 machines, get a 1000 machines
- Becoming cost effective (TCO)
  - Especially “spot” pricing



# Meeting the Challenges

Dataset Storage

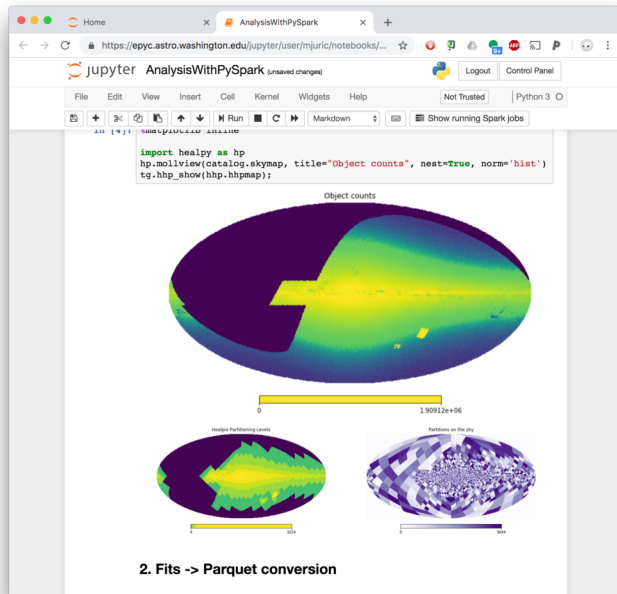
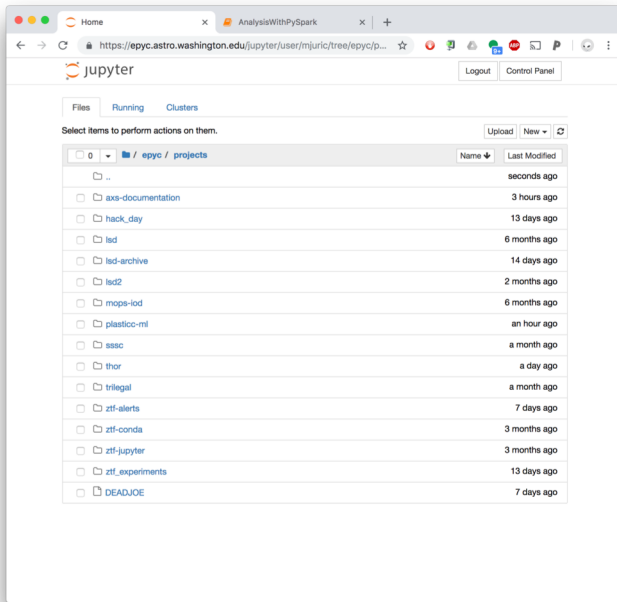
Resources

Scalable Analysis  
Code

Interface



# “Analysis 2025”



# A Number of Projects are Working to Make this Happen

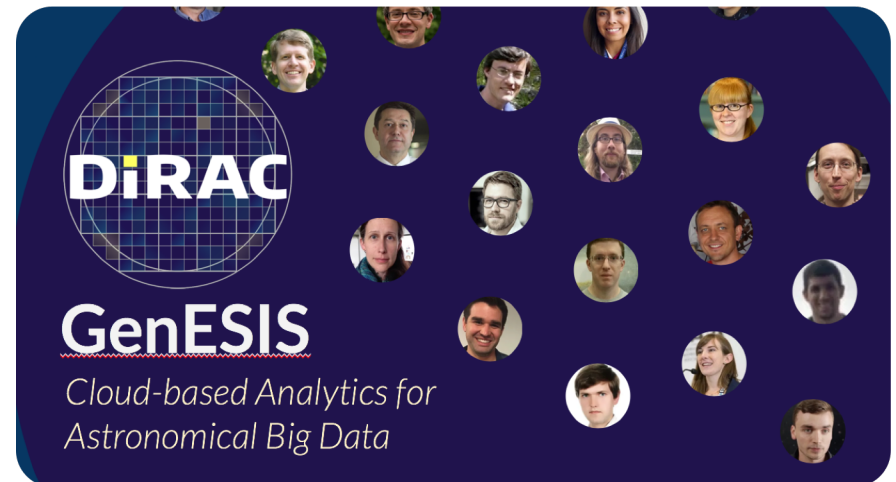


**PANGEO**

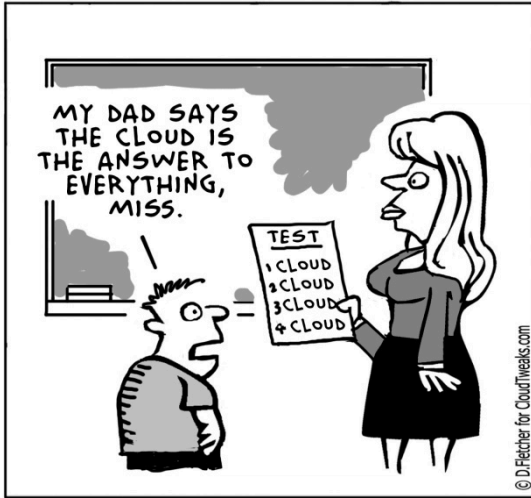
A community platform for Big Data geoscience

<http://pangeo.io/>

*Coming soon w. ZTF !*



# Some Words of Caution



Just like with machine learning / A.I., there's no need to throw cloud at everything.

Small datasets?  
Large-ish datasets?

But the *programming model* works across all scales.



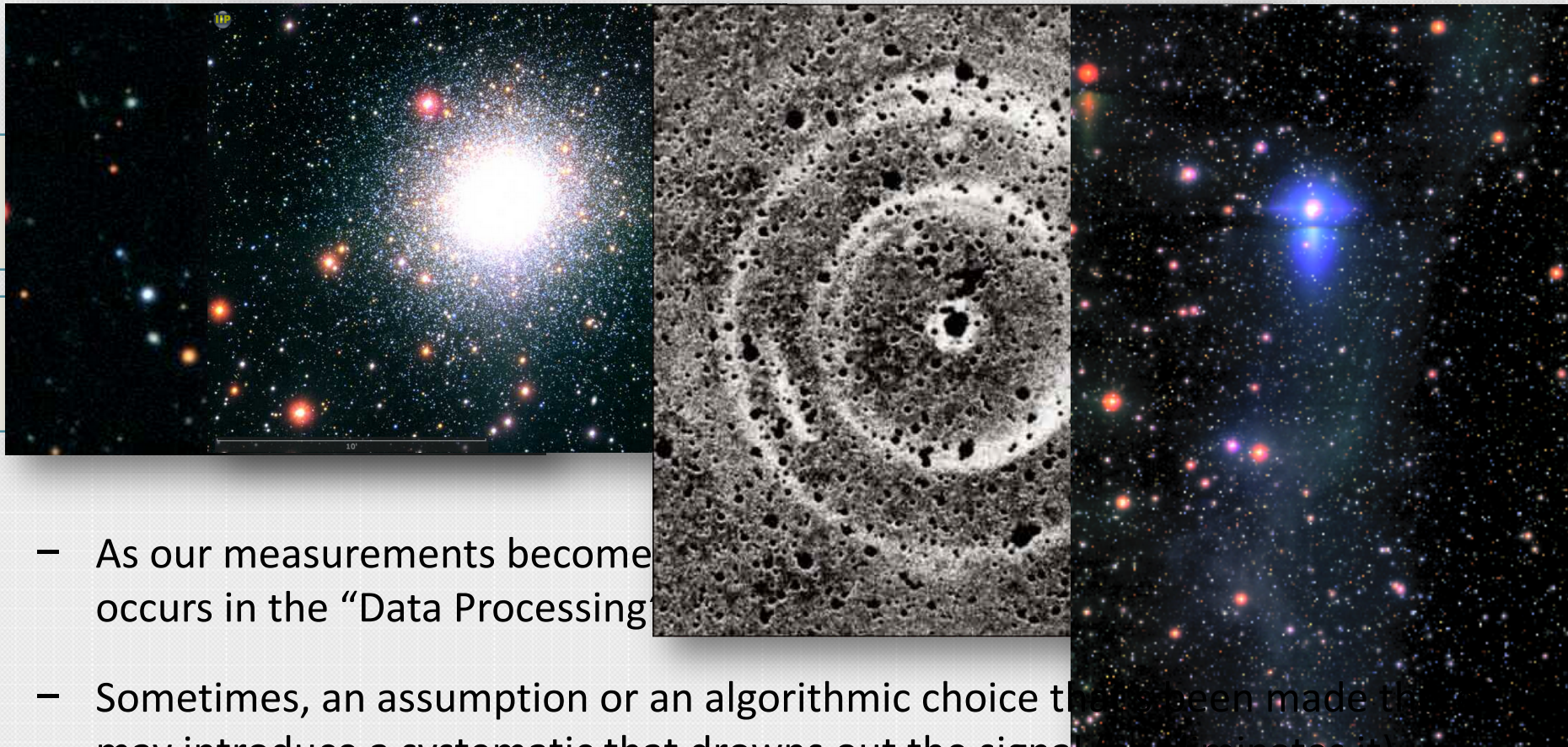
The implementation of these technologies is still in its infancy. They change incredibly quickly.

Expect you may need to shift from framework to framework (e.g., Spark → Dask).

That said, the *programming models* change on a much longer timescale (e.g., MR 2004 → ).

## ***Looking Ahead: Leave no Information Behind***

*(or why software and services are even more important than we think)*



- As our measurements become more complex, the amount of data that occurs in the “Data Processing” stage increases.
- Sometimes, an assumption or an algorithmic choice that has been made throughout the pipeline may introduce a systematic that drowns out the signal (or eliminates it).
- For optimal inference, one wants to design measurements that directly probe the relevant aspects of the *original (imaging data)*, and not the (lossy-compressed) catalog.
  - Or derive more appropriate catalogs/feature sets/etc.

# Pushing the Boundaries of Optimal Inference

**Model** ← *inference* – **Data**

**Model** ← *inference* – **Catalog** ← Data Processing – **Data**

## – Reasons we don't do this today:

1. Computationally (and I/O) intensive
2. Sociologically difficult
  - Expertise in statistics, applied math, and software engineering is often not there
  - Catalogs are too often taken as “God given”, fundamental, result of a survey

## – Things are changing

- Big data problems are becoming computationally tractable (see prev. discussion)
- Average astronomer in the 2020s will grow up with an expectation of being well versed in Stats, SE, Appl. Math.
- A concerted effort is under way, primarily driven by people in large survey and telescope projects, to create the necessary software to make this possible.

## Astronomy 2025: “Personalized Medicine”

- In the next decade, it may be possible for any one of you to re-reduce large datasets for optimally your science case.
- You will be able to do this because the software building blocks (AstroPy, LSST stack, etc.) will be there, with frameworks and cloud resources for large-scale computation.
- Right now, we see the data releases as the key product of a survey. By the end of the next decade, I wouldn't be surprised if we saw **the software as the key product**, with hundreds specialized (and likely ephemeral) catalogs being generated by it.
- The official “data releases” will just be some of those catalogs, designed to be more broadly useful than others, and retained for a longer period of time.





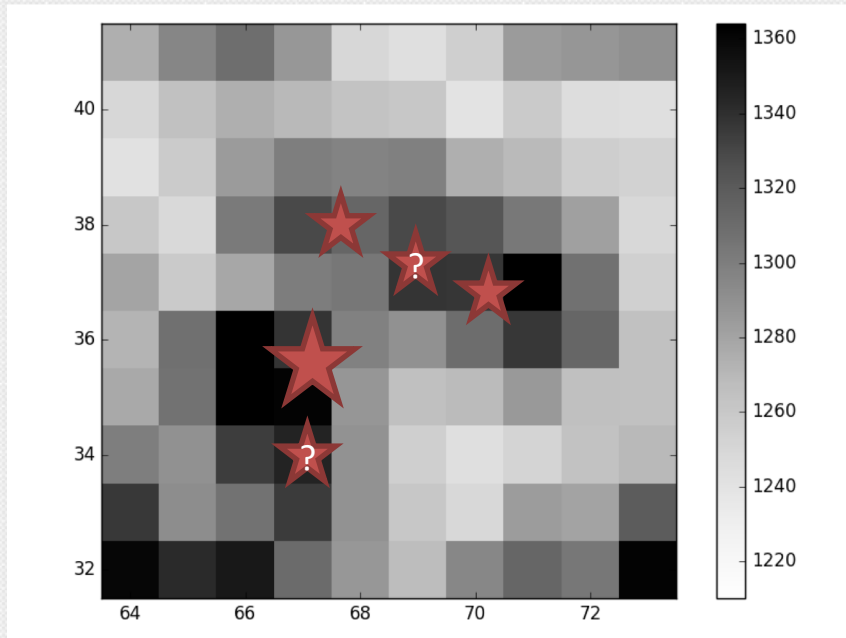
*An Example*

***Probabilistic cataloging for crowded  
fields***

*(Stephen Portillo et al.)*

# Crowded Field Cataloguing

- Neighbouring sources are covariant
- Deblending can be difficult or even ambiguous
- The inferred properties depend on how the image is deblended

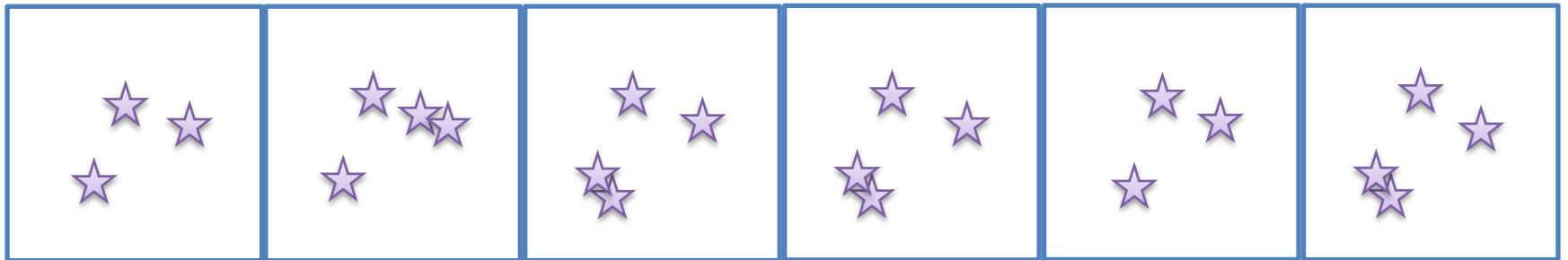


SDSS DR14



## Probabilistic Cataloguing

- Instead of having one catalog, produce an ***ensemble of catalogs*** (each with an associated probability of occurrence)
- Naturally handles deblending ambiguities and source-source covariance in crowded fields

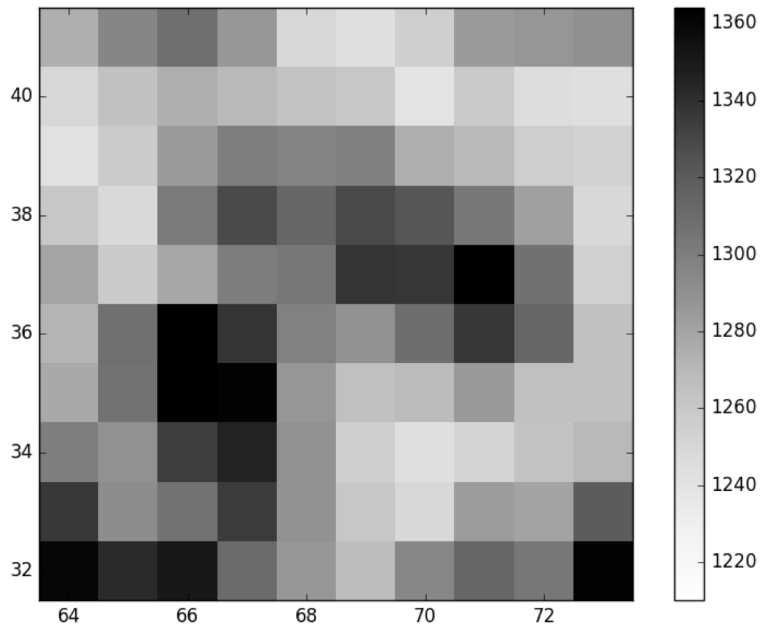


- Space of possible catalogues is ***transdimensional***

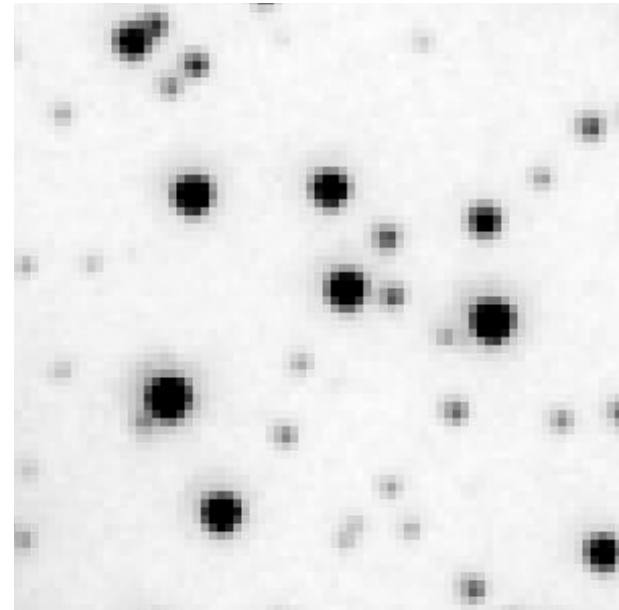
Brewer, Foreman-Mackey, and Hogg (2013)

# Application: Deblending the Cluster M2

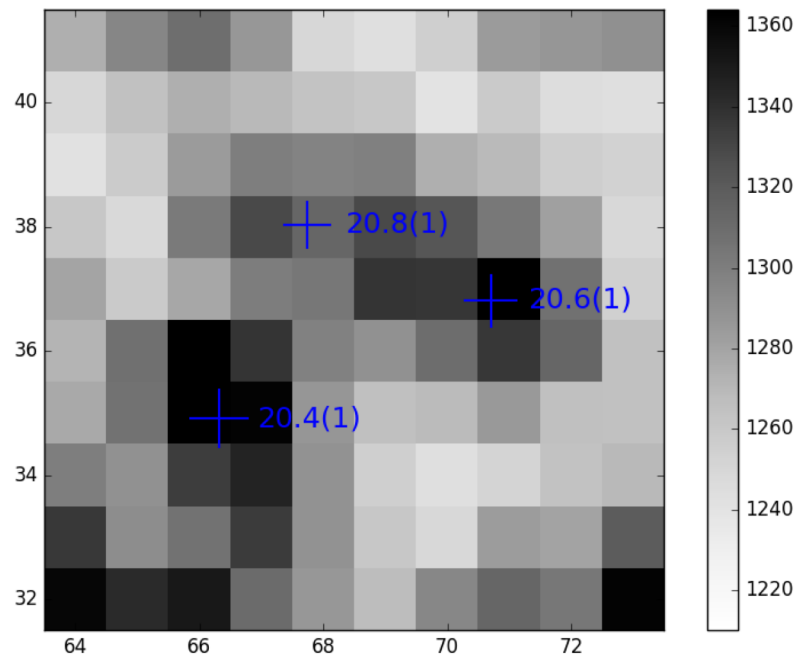
- Sloan Digital Sky Survey



- Hubble Space Telescope



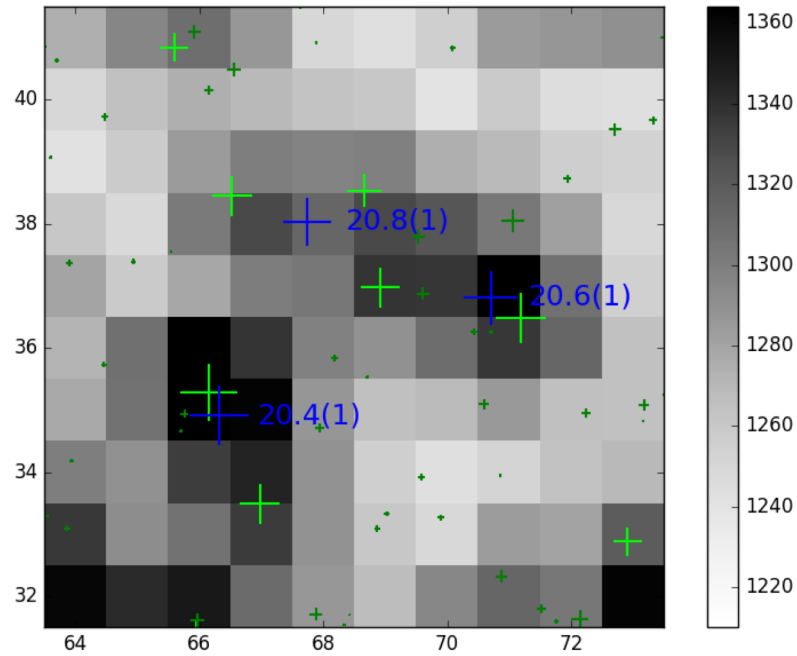
# Traditional Catalogue



DAOPhot

SDSS DR14  
An et al. (2008)

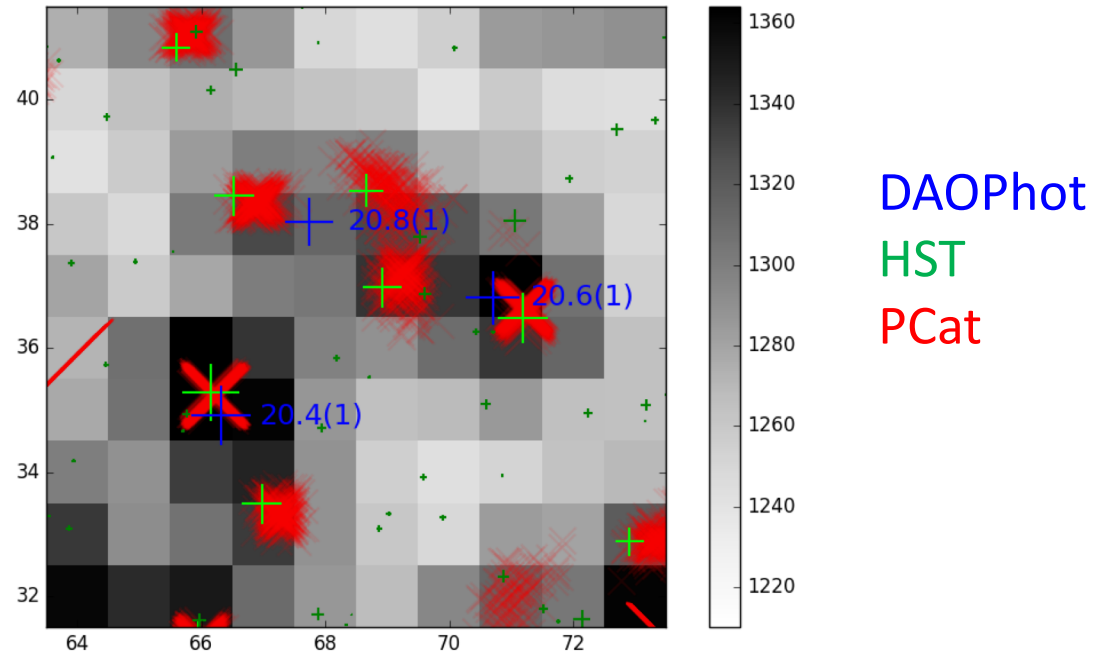
# Compared to Hubble



DAOPhot  
HST

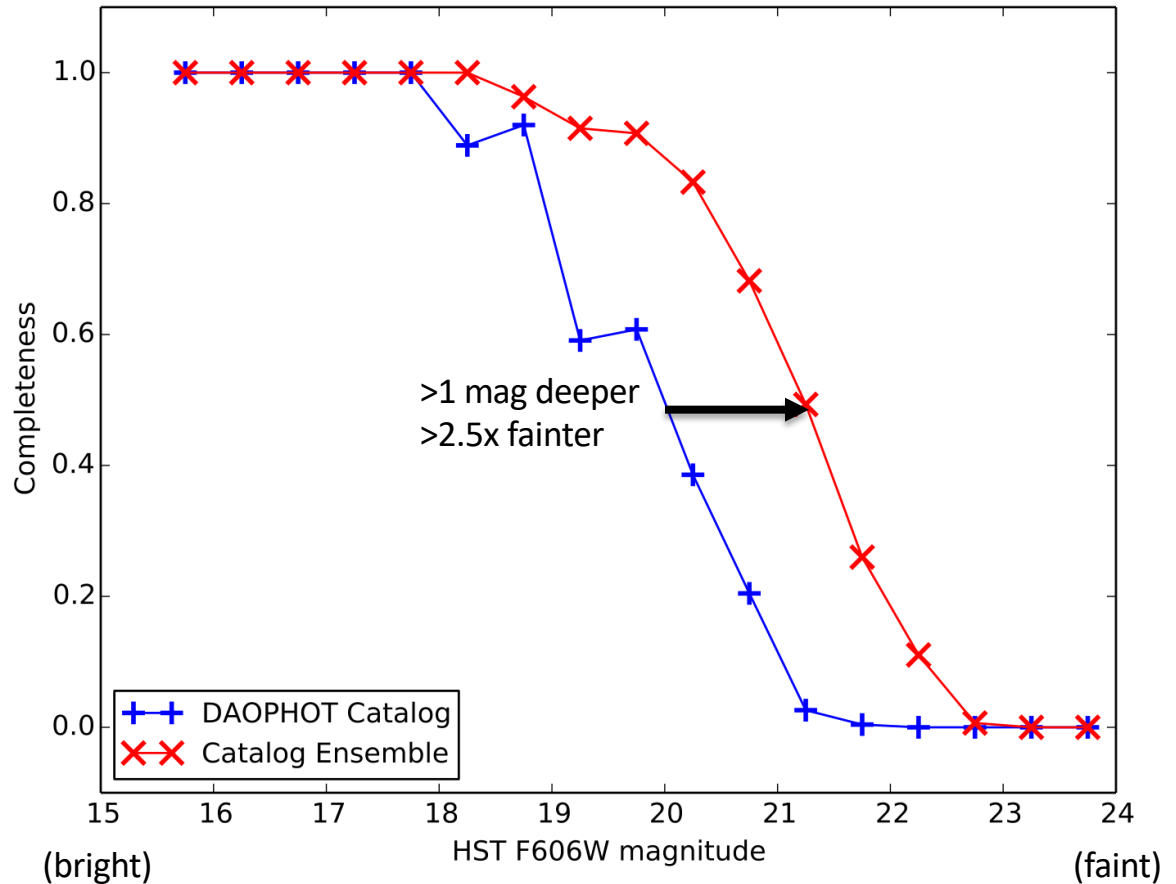
SDSS DR14  
An et al. (2008)  
Sarajedini et al. (2007)

# Stacked Catalogue Ensemble



# Completeness

$$\frac{\text{true positives}}{\text{all real sources}}$$





## Where is all this going

- The data is big, but not unmanageable. But technologies exist (in the industry) to meet the challenges.
- Two changes in paradigms:
  - New programming models (and frameworks): MapReduce (Spark)
  - Analysis on cloud services, rather than on local machines
- This is an ***opportunity***: we'll soon be able to take data analysis one level closer to the *images* (and therefore extract more data). Or devise custom, complex, analyses over entire datasets.
  - E.g., crowded field codes.
- Adding ML/AI in the mix, sky's the limit....

Machine learning	Statistics
network, graphs	model
weights	parameters
learning	fitting
generalization	test set performance
supervised learning	regression/classification
unsupervised learning	density estimation, clustering
large grant = \$1,000,000	large grant = \$50,000
nice place to have a meeting: Snowbird, Utah, French Alps, <b>Tenerife</b>	nice place to have a meeting: Las Vegas in August

Some important differences between machine learning and statistics.