

LSST Data Management: Building the Data System for the Era of Petascale Optical Astronomy

Mario Juric

LSST Data Management Project Scientist

WRF Data Science Chair in Astronomy, University of Washington

and the LSST Data Management Team.



ADASS XXV

Sydney, Australia, October 26th, 2015

- Large Survey and Why They're Different

Hipparchus of Rhodes (180-125 BC)

Discovered the precession of the equinoxes.

Measured the length of the year to ~6 minutes.

In 129 BC, constructed one of the first star catalogs, containing about 850 stars.



n.b.: also the one to blame for the magnitude system ...

Galileo Galilei (1564-1642)

Researched a variety of topics in physics, but called out here for the introduction of the *Galilean telescope*.

Galileo's telescope allowed us for the first time to *zoom in* on the cosmos, and study the individual objects in great detail.



The Astrophysics Two-Step

- Surveys
 - Construct catalogs and maps of objects in the sky. Focus on coarse classification and discovering targets for further follow-up.
- Large telescopes
 - Acquire detailed observations of a few representative objects. Understand the details of astrophysical processes that govern them, and extrapolate that understanding to the entire class.

Analogy: Google Search

A screenshot of a Google search for "two step". The search bar shows "two step" and the search button is a magnifying glass. Below the search bar, there are tabs for "Web", "Videos", "Images", "Maps", "Shopping", and "More". The search results show "About 263,000,000 results (0.52 seconds)". The first result is "Two-step - Wikipedia, the free encyclopedia" with a link to "en.wikipedia.org/wiki/Two-step". The second result is "Two-step (dance move) - Wikipedia, the free encyclopedia" with a link to "en.wikipedia.org/wiki/Two-step_(dance_move)". The third result is "How to dance the Two-Step. Free 2-Step Dancing Lessons ..." with a video thumbnail and a link to "www.youtube.com/watch?v=XRYOqjDkcc". The fourth result is "Texas Two Step - Texas Lottery" with a link to "www.tlottery.org/export/sites/lottery/.../Texas_Two_Step/". The fifth result is "Two Step Restaurant and Cantina San Antonio" with a link to "www.twosteprestaurant.com".

A screenshot of the Wikipedia article for "Two-step (dance move)". The article title is "Two-step (dance move)" and it is from Wikipedia, the free encyclopedia. There is a warning box that says "This article does not cite any references or sources. Please help improve this article by adding citations to reliable sources. Unsourced material may be challenged and removed." Below the warning box, the article text says "The **two-step** is a step found in various dances, including many folk dances. It seems to take its name from the 19th century dance related to the Polka." and "A two-step consists of two steps in approximately the same direction onto the same foot, separated by a closing step with the other foot. For example, a right two-step forward is a forward step onto the right foot, a closing step with the left foot, and a forward step onto the right foot. The closing step may be done directly beside the other foot, or obliquely beside, or even crossed, as long as the closing foot does not go past the other foot." and "Some types of two-step, or related steps, are named "lock step".

A screenshot of a YouTube video titled "How to dance the Two-Step. Free 2-Step Dancing Lessons w/Shawn Trautman". The video shows a man and a woman dancing in a studio. The man is wearing a plaid shirt and jeans, and the woman is wearing a red top and a dark skirt. The video player shows a progress bar at 0:10 / 7:46. The video has 852,553 views and a "Subscribe" button for Shawn Trautman.

Google's index is a catalog of the Web. We use it to "zoom in" on individual entries to find out more.

weather sydney

About 105,000,000 results (0.35 seconds)

Sydney NSW, Australia

Monday 7:00 AM
Mostly Sunny

19°C | 64°F

Precipitation: 0%
Humidity: 78%
Wind: 10 km/h

Temperature Precipitation Wind

19 26 29 28 23 19 18 17

8 AM 11 AM 2 PM 5 PM 8 PM 11 PM 2 AM 5 AM

Mon	Tue	Wed	Thu	Fri	Sat	Sun	Mon
28° 16°	18° 14°	21° 13°	22° 14°	22° 16°	23° 16°	24° 18°	26° 19°

More on weather.com Feedback

Display a menu

But, it's more than just a catalog of pointers – more and more, Google itself collects, processes, indexes, visualizes, and serves the actual information we need.

More and more often, our “research” begins and ends with Google!

Entering the Era of Massive Sky Surveys

- There's a close parallel with large surveys in astronomy, in scale, quality, and richness of the collected information
 - Scale: We're entering the era when we can image and catalog the entire sky
 - Quality: The measurements can be as precise as those taken with "pointed" observations (used to be ~5-10x worse)
 - Richness: Those catalogs contain not only positions and magnitudes, but also shapes, profiles, and temporal behavior of the objects.
- Quite often, the research can begin and end with the survey.
- This is what makes large surveys of today **not just bigger, but better, more information rich, and therefore different.** **Big data is now about complexity and optimal knowledge extraction, not PBs (IMHO).**

Sloan Digital Sky Survey

2.5m telescope

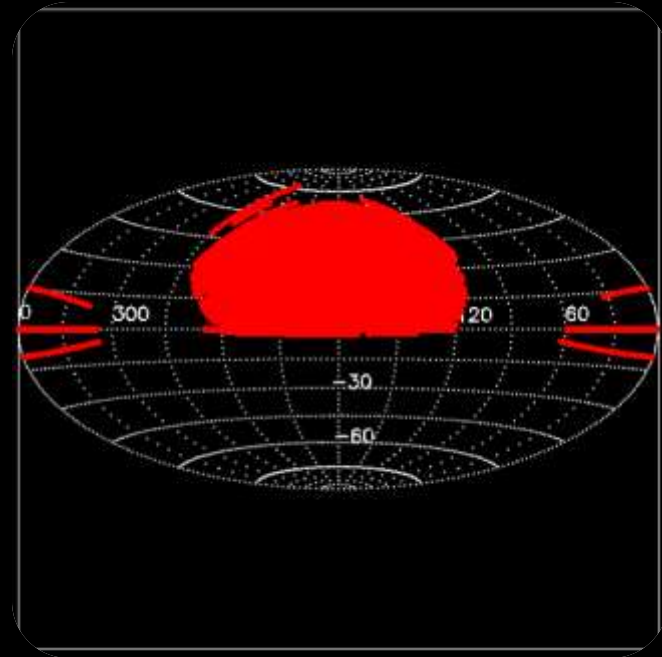
>10000 deg²

0.1" astrometry

r<22.5 flux limit

5 band, 2%, photometry for >50M stars

>300k R=2000 stellar spectra



Panoramic Survey Telescope and Rapid Response System

1.8m telescope

30000 deg²

50mas astrometry

r<23 flux limit

5 band, better than 1% photometry (goal)



LSST: A Deep, Wide, Fast, Optical Sky Survey



8.4m telescope

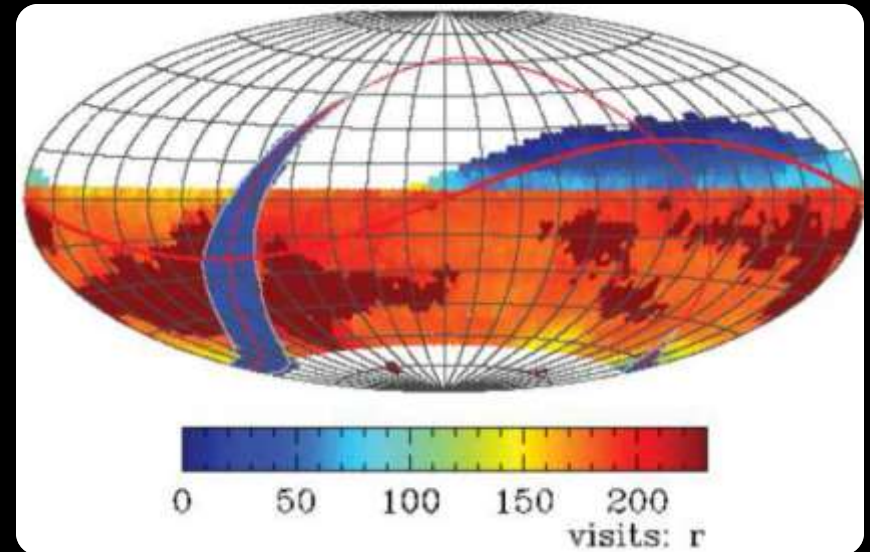
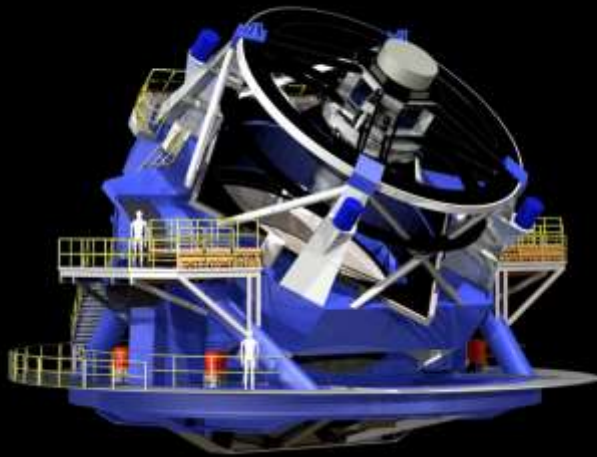
18000+ deg²

10mas astrom.

r<24.5 (<27.5@10yr)

ugrizy

0.5-1% photometry



3.2Gpix camera

30sec exp/4sec rd

15TB/night

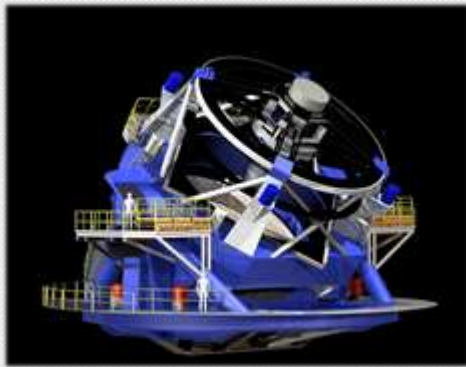
37 B objects

Imaging the visible sky, once every 3 days, for 10 years (825 revisits)

LSST: Turning the Sky into a Database



- A wide (half the sky), deep (24.5/27.5 mag), fast (image the sky once every 3 days) survey telescope. Beginning in 2022, it will repeatedly image the sky for 10 years.
- The LSST will be an automated survey system. In nighttime, the observatory and the data system will operate with minimum human intervention.
- **The ultimate deliverable of LSST is not the telescope, nor the instruments; it is the fully reduced data.**
 - All science will be come from survey catalogs and images



Telescope



Images



Table 4: Level 2 Catalog Object Table

Name	Type	Unit	Description
psRadecl	double	time	Point source model: Time at which the object was at position radecl.
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 ⁵⁸ .
psCov	float[66]	various	Point-source model covariance matrix ⁵⁹ .
psLnL	float		Natural log likelihood of the observed data given the point source model.
bdRadecl	double[2]	degrees	B+D model ⁶⁰ : (α, δ) position of the object at time radecl. In each band...

Catalogs

Location: Cerro Pachon, Chile



Cerro Pachón – Future site of the LSST



Leveling of El Peñón (the summit of Cerro Pachón)

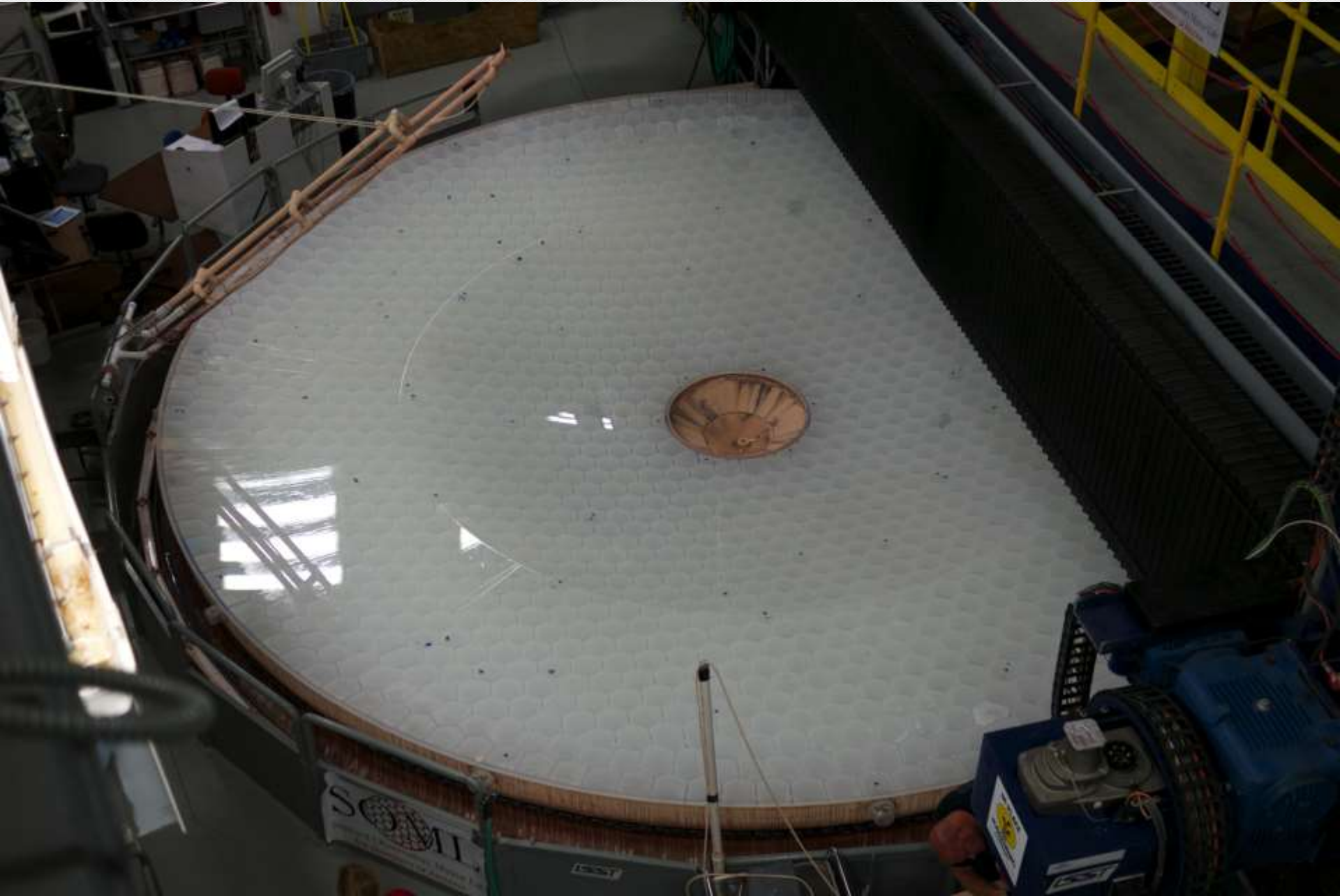




LSST Observatory (cca. late ~2018)



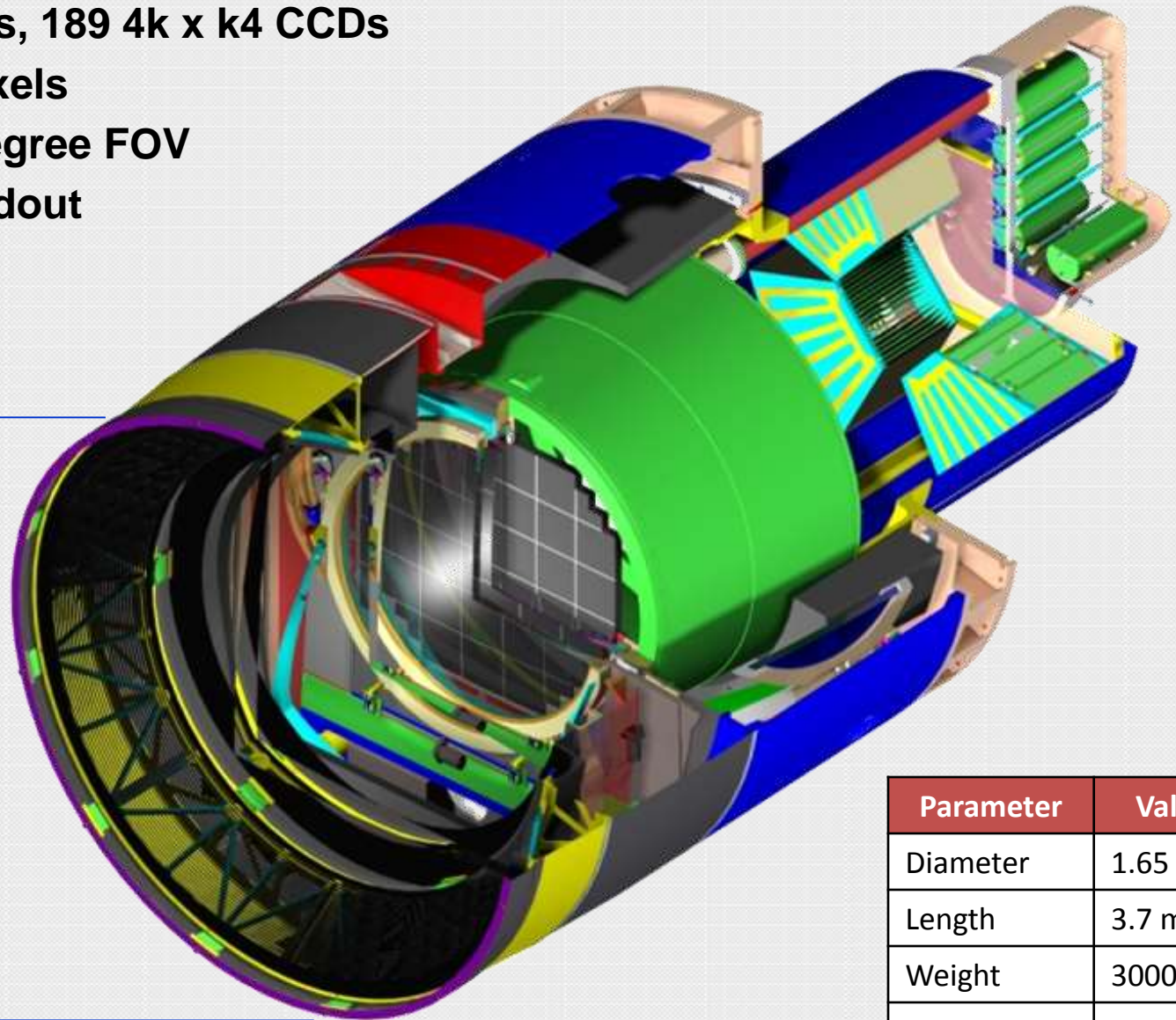
Done!



LSST Camera



- 3.2 Gigapixels, 189 4k x 4k CCDs
- 0.2 arcsec pixels
- 9.6 square degree FOV
- 2 second readout
- 6 filters



1.65 m
5'-5"

Parameter	Value
Diameter	1.65 m
Length	3.7 m
Weight	3000 kg
F.P. Diam	634 mm



LSST Imaging: ~5 PB/yr

(~5 expensive, information rich PB/yr)

LSST Operations: Sites and Data Flows



Satellite Processing Center

(CC-IN2P3, Lyon, France)

Data Release Production (50%)
French DAC



Archive Site

Archive Center

Alert Production
Data Release Production (50%)
EPO Infrastructure
Long-term Storage (copy 2)

Data Access Center

Data Access and User Services

Summit and Base Sites

Telescope and Camera
Data Acquisition
Crosstalk Correction
Long-term storage (copy 1)
Chilean Data Access Center



HQ Site

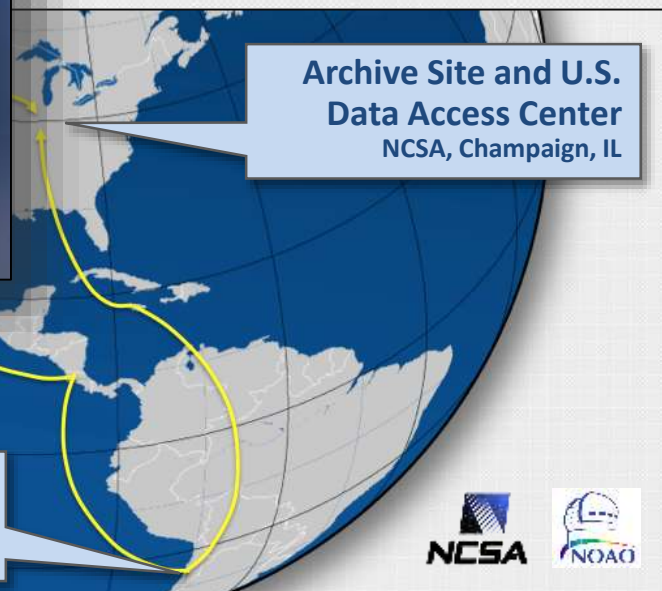


Science Operations
Observatory Management
Education and Public Outreach



*The computing cluster at the **LSST Archive** (at NCSA) will run the processing pipelines.*

- *Single-user, single-application, dedicated data center*
- *Process images in real-time to detect changes in the sky*
- *Produce annual data releases*



Archive Site and U.S. Data Access Center
NCSA, Champaign, IL

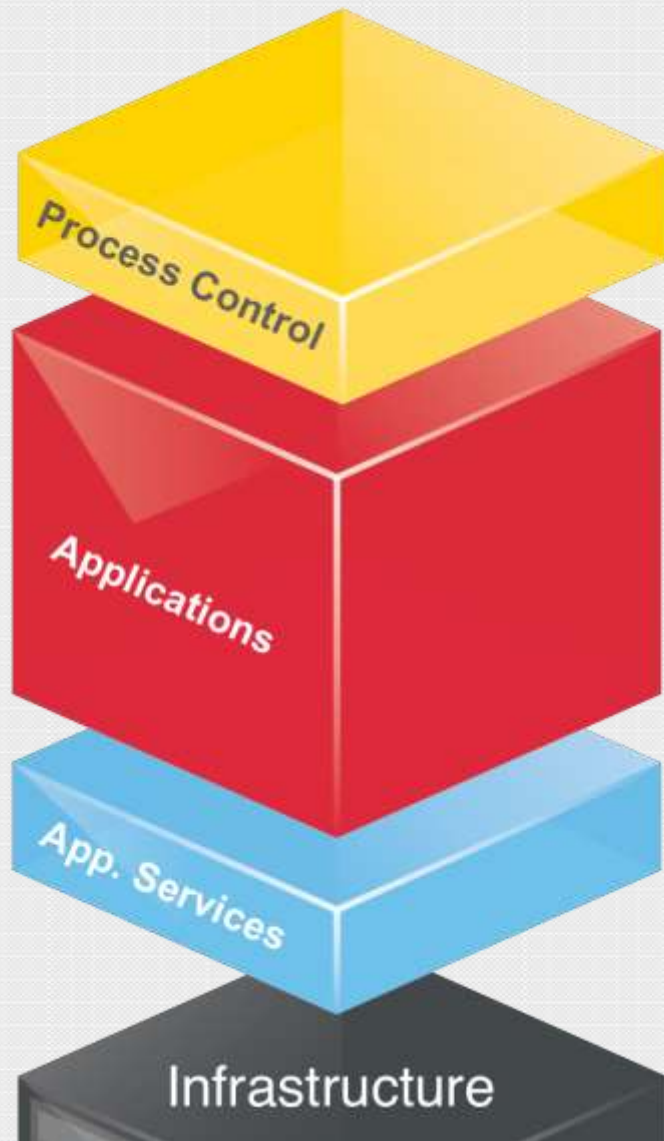
Base Site and Chilean Data Access Center
La Serena, Chile

Long Haul Networks to transport data from Chile to the U.S.

- *200 Gbps from Summit to La Serena (new fiber)*
- *2x40 Gbit (minimum) for La Serena to Champaign, IL (protected, existing fiber)*



“Applications”: Scientific Core of LSST DM



- *Applications* carry core scientific algorithms that process or analyze raw LSST data to generate output Data Products
- Variety of processing
 - Image processing
 - Measurement of source properties
 - Associating sources across space and time, e.g. for tracking solar system objects
- *Applications framework* layer (*afw*; not shown) allows them to be written in a high-level language



Middleware Layer: Isolating Hardware, Orchestrating Software



Enabling execution of science pipelines on hundreds of thousands of cores.

- Frameworks to construct pipelines out of basic algorithmic components
- Orchestration of execution on thousands of cores
- Control and monitoring of the whole DM System

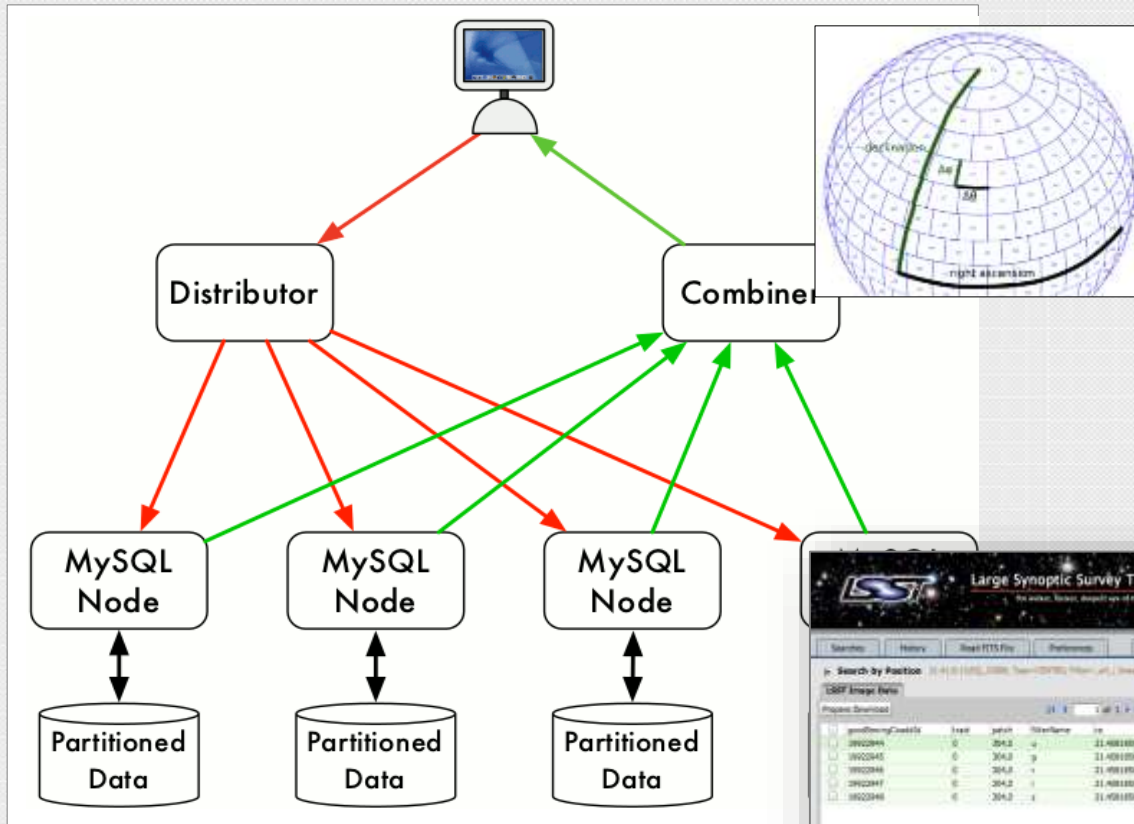


Isolating the science pipelines from details of underlying hardware

- Services used by applications to access/produce data and communicate
- "Common denominator" interfaces handle changing underlying technologies



Database and Science UI: Delivering to Users



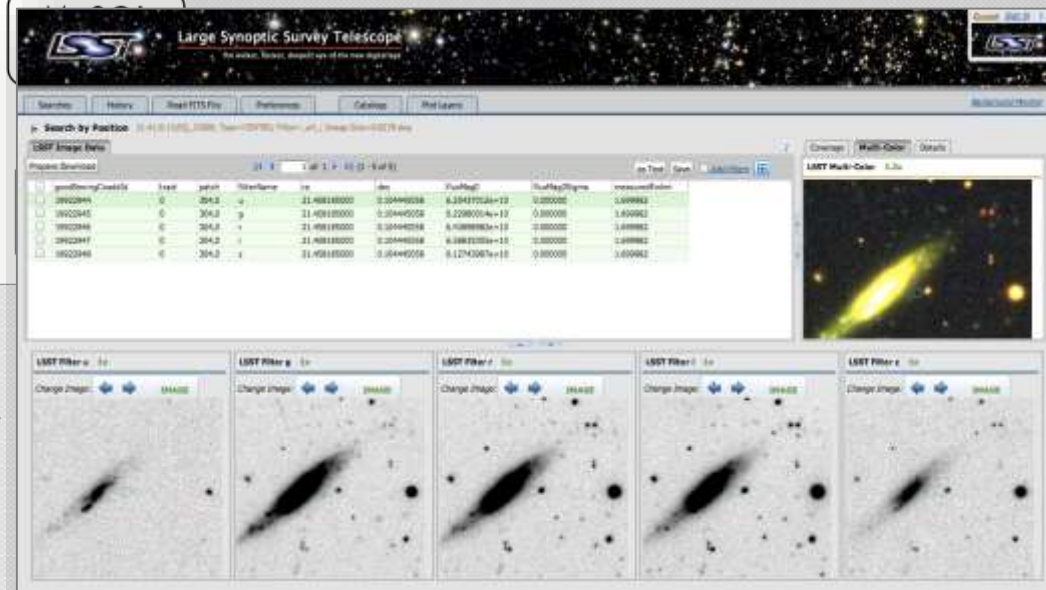
*Massively parallel,
distributed, fault-tolerant
relational database.*

- To be built on existing, robust, well-understood, technologies (MySQL and xrootd)
- Commodity hardware, open source
- Advanced prototype in existence (qserv)



Science User Interface to enable the access to and analysis of LSST data

- Web and machine interfaces to LSST databases
- Visualization and analysis capabilities



LSST DM: A Distributed Development Team



UI

Database

Core Algorithms (“Apps”)

Middleware

Infrastructure

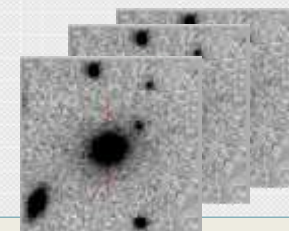
LSST
Mgmt, I&T, and Science QA



LSST's #1 Challenge:

*General purpose processing while
minimizing information loss.*

From Data to Knowledge



And metadata!

Computationally (and cognitively) expensive, science-case specific

Scientists **Model** ← *Scientists* *inference* – **Data** *Project*

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



- 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 LSST 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

- 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, a Database, and a (small) Cloud



- A stream of ~10 million time-domain events per night, detected and transmitted to event distribution networks within 60 seconds of observation.
- A catalog of orbits for ~6 million bodies in the Solar System.
- A catalog of ~37 billion objects (20B galaxies, 17B stars), ~7 trillion single-epoch detections (“sources”), and ~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 user-specified custom processing and analysis.
- Software and APIs enabling development of analysis codes.

Level 1

Level 2

Level 3



- 02C.01.02.01/02. **Data Quality Assessment Pipelines** *(slides by Juric)*
- 02C.01.[02.01.04,04.01,04.02] **Calibration Pipelines** *(slides by Axelrod, Yoachim)*
- 02C.03.01. **Single-Frame Processing Pipeline** *(slides by Krughoff, Lupton)*
- 02C.03.02. **Association pipeline** *(slides by Lupton)*
- 02C.03.03. **Alert Generation Pipeline** *(slides by Becker)*
- 02C.03.04. **Image Differencing Pipeline** *(slides by Becker)*
- 02C.03.06. **Moving Object Pipeline** *(slides by Jones)*
- 02C.04.03. **PSF Estimation Pipeline** *(slides by Lupton)*
- 02C.04.04. **Image Coaddition Pipeline** *(slides by AlSayyad)*
- 02C.04.05. **Deep Detection Pipeline** *(slides by Lupton)*
- 02C.04.06. **Object Characterization Pipeline** *(slides by Lupton, Bosch)*
- 02C.01.02.03. **Science Pipeline Toolkit**
(slides by Dubois-Felsmann)
- 02C.03.05/04.07 **Application Framework**
(slides by Lupton)



Level 1

Level 2

L3



- **Difficulty adapting existing public codes to LSST requirements (AstroMatic suite, PHOTO, Elixir, IRAF-based pipelines, etc.)**
 - Need to run efficiently at scale
 - Need to be flexible (plugging/unplugging of algorithms at runtime)
 - Need to have it developed by a large team (50+ scientists and programmers)
 - Need to be maintainable over ~25 years of R&D, Construction, and Survey Operations
 - Need to run on a variety of hardware and software platforms
 - Need to have logging and provenance built into the design
 - Need to be able to travel back in time, find Sarah Connor, stop Terminator 3, Salvation, and Genisys from ever being made.
- LSST software stack is largely being written from scratch (transferring the algorithmic knowledge!), in **Python 2.7 (w. Python 3 compatible syntax)**, unless computational demands require the use of **C++**

Example: Toy SExtractor with lsst primitives (1/2)



```
exposure = afwImage.ExposureF(fileName)
mi = exposure.getMaskedImage()
im = mi.getImage()

#
# Subtract background
#
back_size = 64
bctrl = afwMath.BackgroundControl(im.getWidth()//back_size + 1,
                                   im.getHeight()//back_size + 1)
backobj = afwMath.makeBackground(im, bctrl)

im -= backobj.getImageF("LINEAR")

#
# Smooth image with a
#   1 2 1
#   2 4 2
#   1 2 1
# filter
oneD = afwMath.PolynomialFunction1D([2, 0, -1])
kernel = afwMath.SeparableKernel(3, 3, oneD, oneD)

smoothedIm = im.Factory(im.getDimensions())
afwMath.convolve(smoothedIm, im, kernel)

threshold = afwDetection.Threshold(threshold)
fs = afwDetection.FootprintSet(smoothedIm, threshold, npixMin)

if grow > 0:
    isotropic = False
    fs = afwDetection.FootprintSet(fs, grow, isotropic)
```

Example: Toy SExtractor with lsst primitives (2/2)



```
fs.setMask(mi.getMask(), "DETECTED")

# Define the measurements we want to make
ctrlCentroid = measAlg.SdssCentroidControl()
ctrlAperture = measAlg.SincFluxControl()
ctrlAperture.radius2 = apRad

schema = afwTable.SourceTable.makeMinimalSchema()
algorithms = [
    measAlg.MeasureSourcesBuilder().addAlgorithm(ctrlCentroid).build(schema),
    measAlg.MeasureSourcesBuilder().addAlgorithm(ctrlAperture).build(schema)]

cat = afwTable.SourceCatalog(schema)

table = cat.table
table.defineCentroid("centroid.sdss")
table.defineApFlux("flux.sinc")

# Measure sources
fs.makeSources(cat)
print "Measuring %d objects" % (len(cat))

for source in cat:
    for alg in algorithms:
        alg.apply(source, exposure)

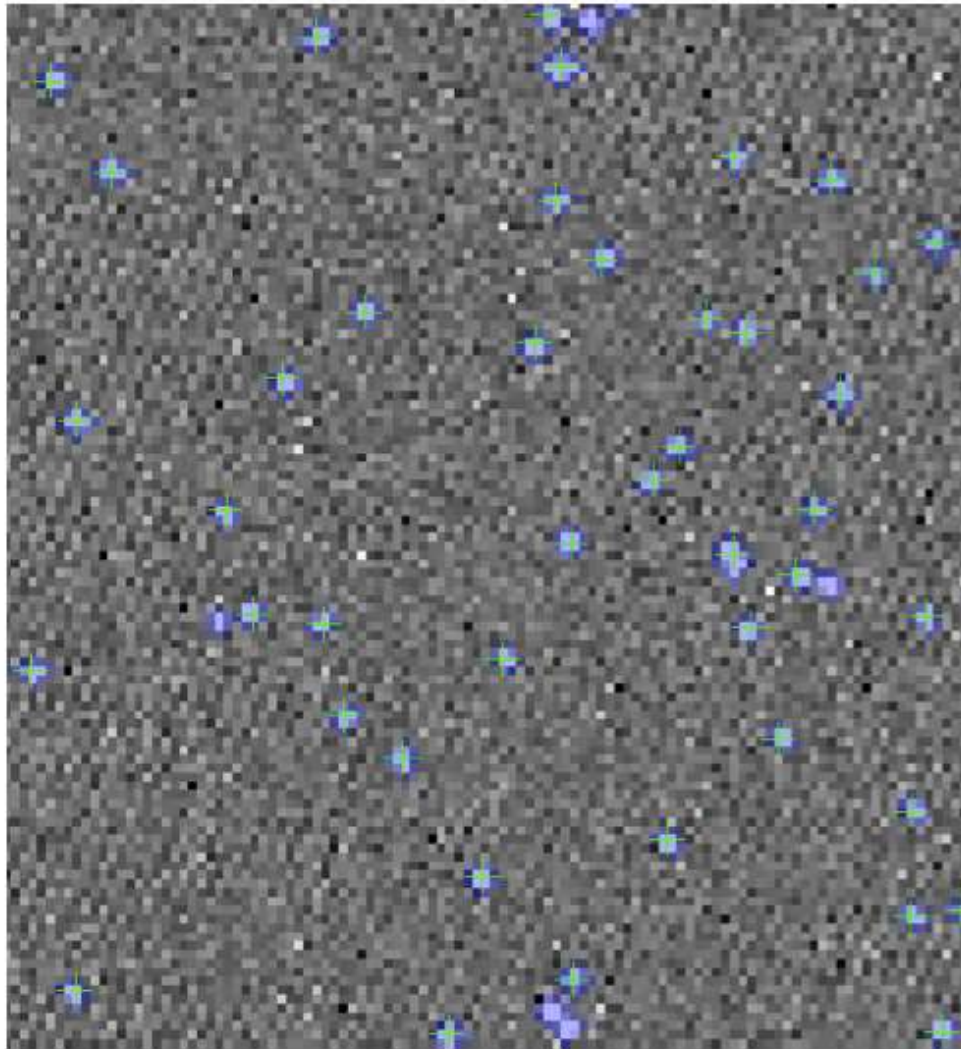
if display:
    ds9.mtv(mi, frame=1, title="Detection")
    with ds9.Buffering():
        for source in cat:
            ds9.dot("+", *source.getCentroid(), frame=1)

return mi, cat
```

For more, see the poster by Tim Jenness et al (P056):

“The LSST Data Processing Software Stack: Summer 2015 Release”


```
$ main.py C0_20090717-214738-141.fits.gz --display
Measuring 1502 objects
Elapsed time = 1.39s
```

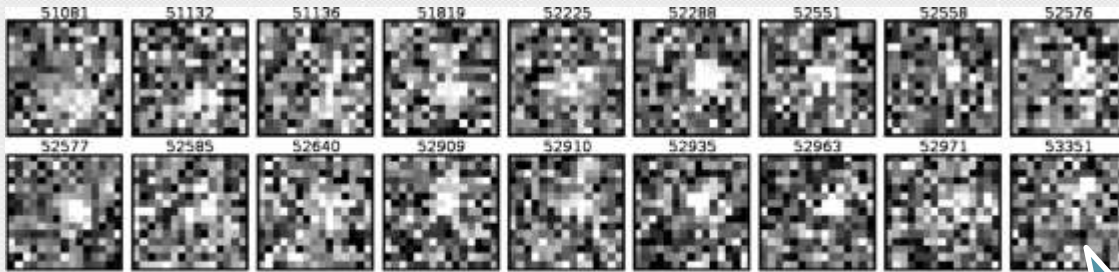


Right: Detection of ^{56}Fe hits (lab characterization of LSST CCDs)

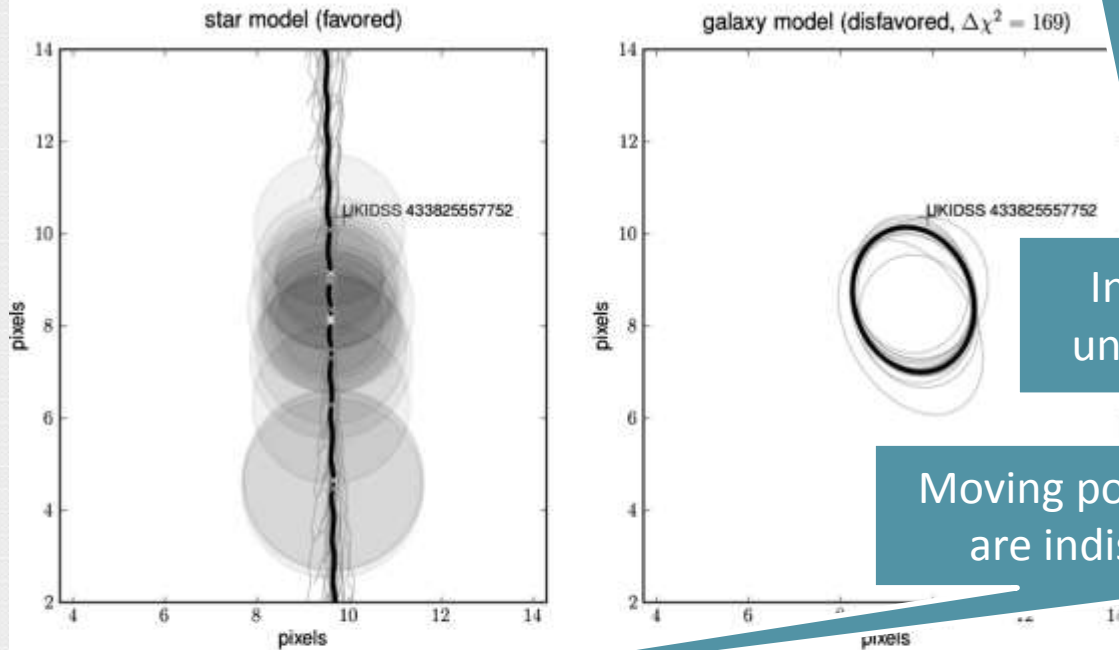
For more, see the poster by Tim Jenness et al (P056):

“The LSST Data Processing Software Stack: Summer 2015 Release”

Detecting and Estimating Proper Motions Below the Single-Epoch Flux Limit

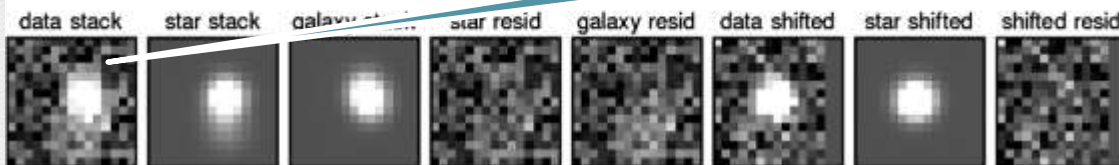


Optimal measurement of properties of objects imaged in multiple epoch. Left: extraction of a moving point source (Lang 2009).

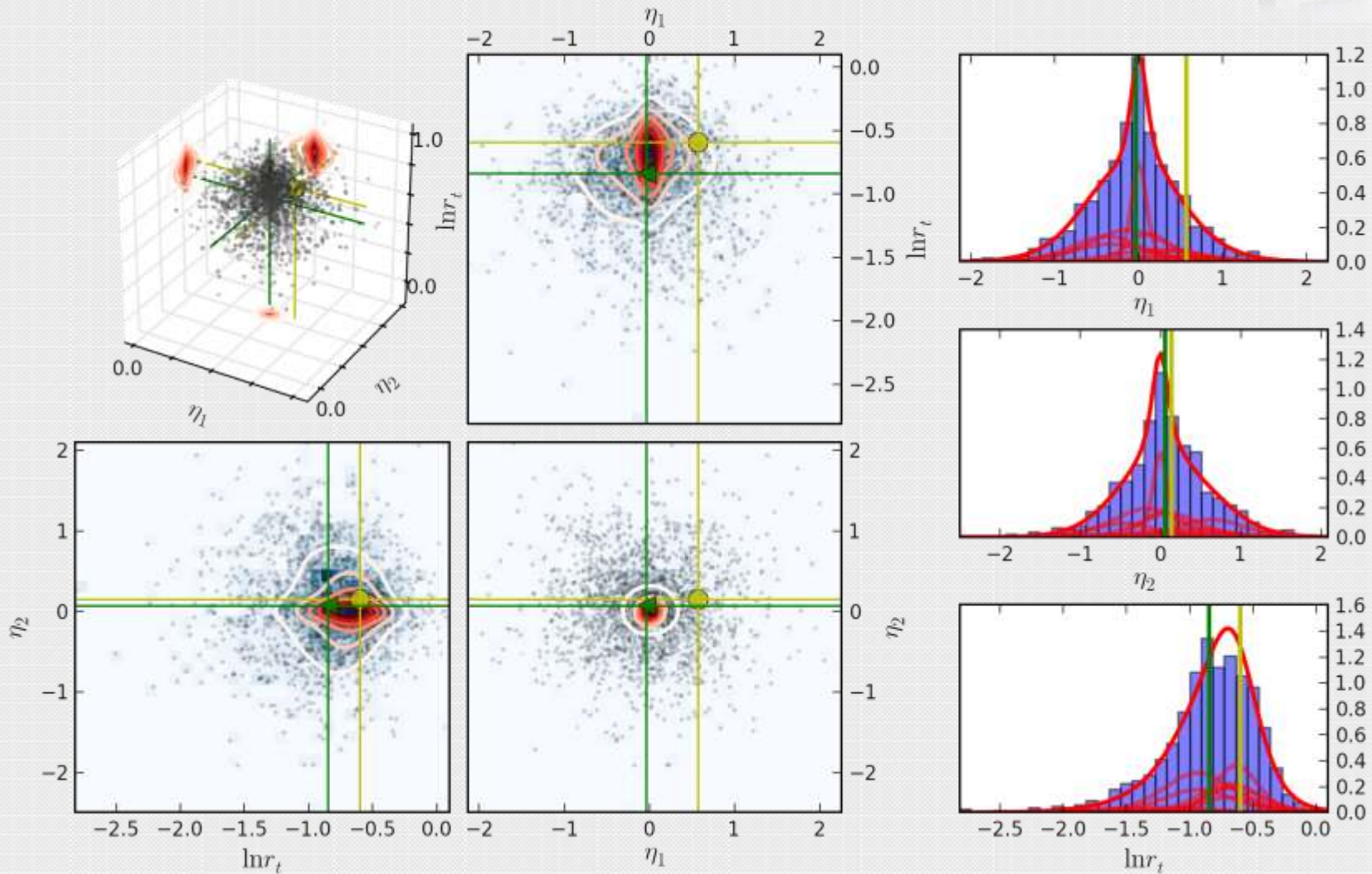


Individual exposures: objects are undetected or marginally detected

Moving point-source and galaxy models are indistinguishable on the coadd



Example: Sampling and retaining the Likelihoods



Perform importance sampling from a proposal distribution determined on the coadd. Plan to characterize (and keep!) the full posterior for each object. **(Unexplored) possibilities for compression.**

Improved Algorithms: Background-matched co-add of SDSS Stripe 82 in the vicinity of M2.

Background matching preserves diffuse structures.

Generated with LSST pipeline prototypes.



Figure:
5 sq. deg.
background-matched
coadd composite

(g,r,i)
~55 epochs

Region: Aqr
Galactic lat = -35.0

Slide: Yusra AlSayyad

Dec (J2000)

00'
-0°30'
-1°00'

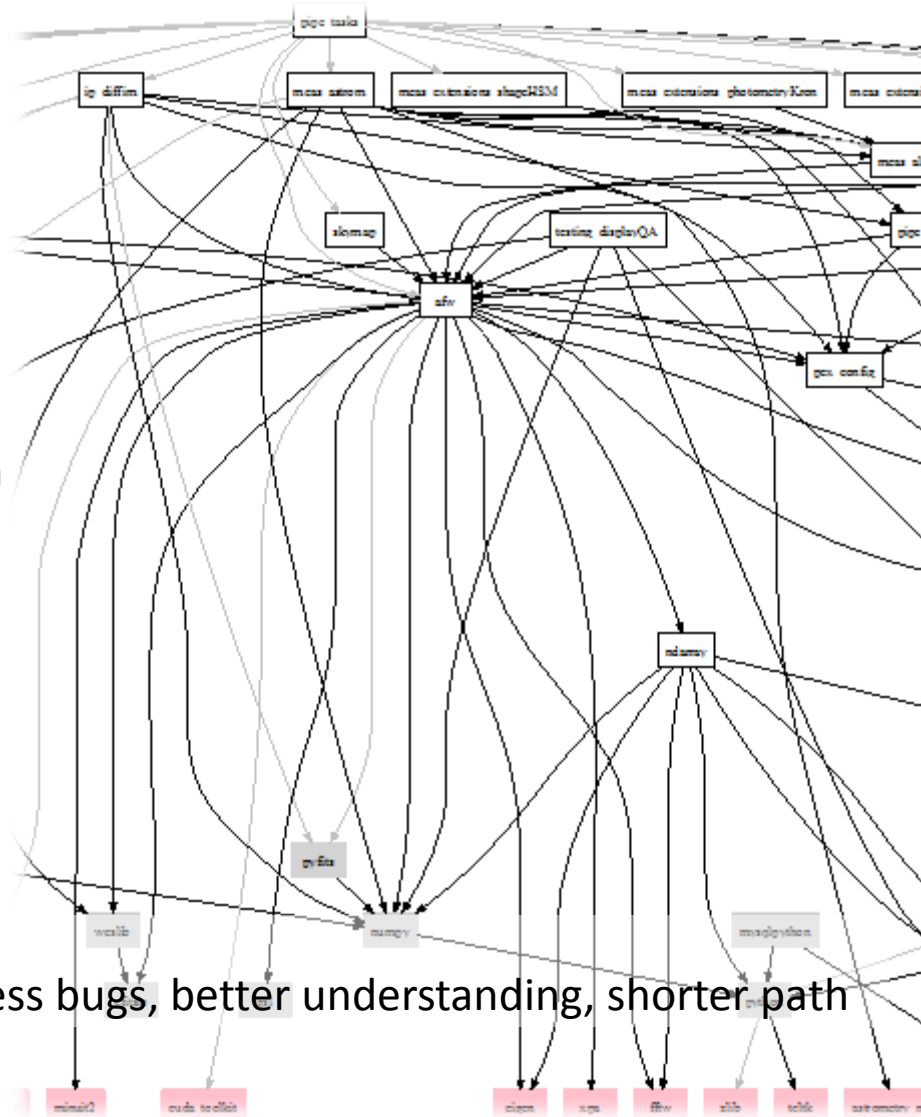
38m 36m 34m
RA (J2000)

<http://moe.astro.washington.edu/sdss/>



LSST Vision: Reusable and Open Code and Development

- LSST software is designed to be ***general purpose*** and highly ***reusable***.
 - Necessary to enable “Level 3” use cases (user-driven processing)
 - Necessary to deal with real-world hardware
 - Necessary to be able to process precursor data
- **Opportunities for using LSST-derived code on other data sets**
 - **Used to reprocess SDSS Stripe 82 data**
 - **Used in production for HSC Survey**
 - CFHT and DECam coming this year
 - Possibilities: PS1, HSC, DES, WFIRST, Euclid, ...
 - Good basis for analysis frameworks (LSST DESC)
- **The benefits feed back to LSST:** more users, less bugs, better understanding, shorter path to science.



LSST Software is Freely Available



The screenshot shows a web browser window with the URL `dm.lsst.org/#code`. The page title is "LSST :: Data Management" and the navigation menu includes "Getting the Code", "Getting Involved", and "About LSST". The main heading is "Getting the Code".

The text under "Getting the Code" states: "The LSST data processing codes are being developed in an iterative, agile, fashion. Though engineering first light is still six years away, prototype versions of a number of LSST codes are already being tested on simulations and being applied to existing data (e.g., [reprocessing SDSS Stripe 82](#), or [processing HSC Survey data](#)).

While already state-of-the-art in many areas, LSST software is still in its infancy when it comes to end-user friendliness, documentation, and API stability. There is no binary distribution yet — builds must be done from source. Knowledge of Python (and willingness to write some Python code) are necessary to work with the current code base.

Warning At this stage, the LSST software will be of greatest interest to the LSST Science Collaborations, large survey builders (or those reprocessing large survey data sets), and astronomical image processing enthusiasts. If you're just looking to reduce a few observations with a ready-to-use tool, it may be better to look into one of the more polished and/or established packages such as [AstroPy](#) or the [AstrOmatic](#) suite.

Installing

Assuming you have the prerequisites and are running `bash`:

```
curl -O https://sw.lsstcorp.org/eupspkg/newinstall.sh
bash newinstall.sh

source loadLSST.bash

eups distrib install -t v9_2 lsst_distrib
```

This will download and build from source a specific release of the LSST Stack (v9.2, in the example above). For complete instructions, see the [documentation](#).

Once you've installed the stack, see [here](#) for examples of what you can do with it.

Cloning the sources

All LSST DM code is visible on [GitHub](#), spread across 100+ repositories. You may find the [LSST software build tool](#) helpful for cloning and (re)building from git. Feel free to subscribe to the [dm-devel](#) mailing list for help.

<http://dm.lsst.org>

and

<http://github.com/LSST>

#1 Outstanding Issue:
Incomplete documentation & insufficient end-user friendliness

Planning to address that over the next ~year.

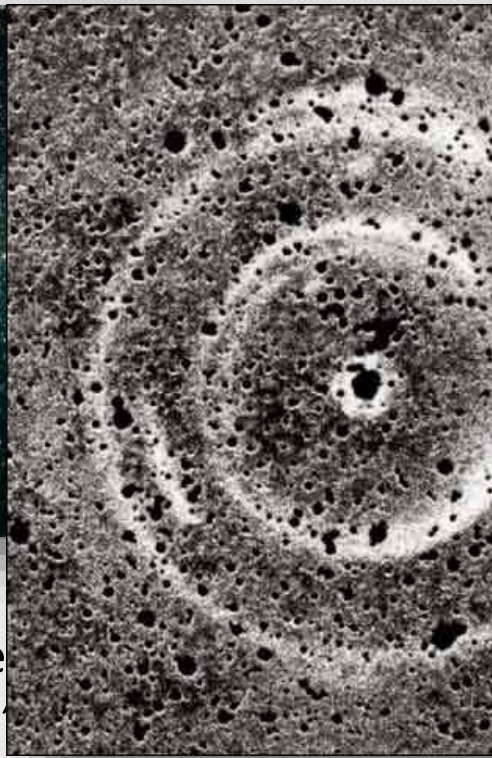
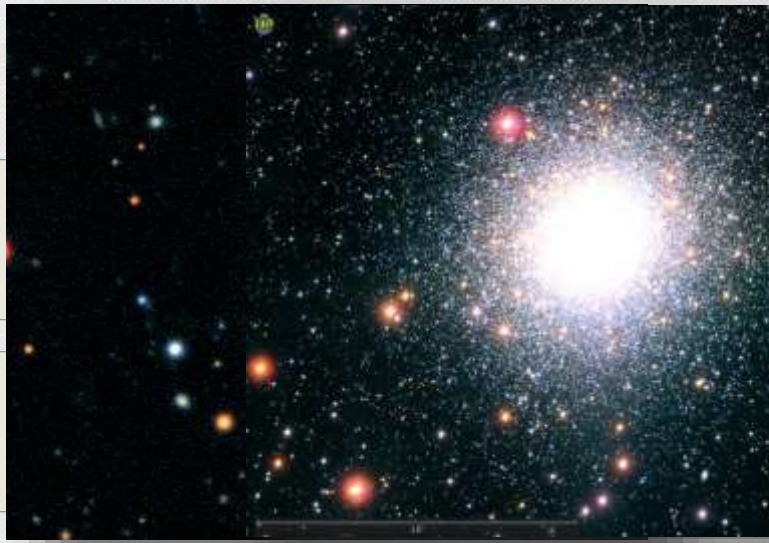
WARNING:

Still under HEAVY development; don't use unless you have a real need (e.g., want to reprocess an existing or are building a new survey).



Looking Ahead: Astrophysical Inference in the 2020s

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



- As our measurements become more precise, the amount of data that occurs in the “Data Processing” stage increases significantly.
- Sometimes, an assumption or an algorithmic choice that’s 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.



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 increasingly computationally tractable
- 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.



- LSST “Level 3” concept our first step in that direction: **Enabling the community to create new products using LSST’s software, services (APIs), and/or computing resources.** This means:
 - **Providing the software primitives to construct custom measurement/inference codes**
 - **Enabling the users to run those codes at the LSST data center, leveraging the investment in I/O (piggyback onto LSST’s data trains).**
- Looking ahead: Right now, we see the data releases as the key product of a survey. By the end of LSST, 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.
- LSST “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.
- **LSST software software and hardware is being engineered to make this possible.**



- Traditionally, astronomy was a data-starved science. Our methods and approach to research were shaped by this environment. Surveys are altering it; data is becoming abundant, well characterized, and rich.
- LSST is a poster-child of this transformation: it will deliver the positions, magnitudes and variability information for virtually *everything* in the southern sky to 24th-27th magnitude.
- To enable science, we're building a new data processing system, with software written in Python and C++. Our pipelines are designed to be maintainable and reusable for other cameras.
- As the costs of new surveys and telescopes are entering ~Bn\$1 range, the importance of software and optimal data processing is growing as well. **Reusability at every level will be crucial.** Collaboratively developed, open, extensible toolkits (e.g., AstroPy, LSST) and “Data Center in a Box” concepts (like some aspects of LSST Level 3) are beginning to address this.



BoF Session 2

4:15pm –
5:45pm

Sarah Brough
BoF2: LSST and Australia
Ballroom 1

LSST IS HIRING



WE'RE SEEKING TOP TALENT TO WORK IN A TEAM ENVIRONMENT THAT INSPIRES EXCELLENCE.



JOIN US IN:

LSST HEADQUARTERS
TUCSON, AZ

SLAC/STANFORD
MENLO PARK, CA

PRINCETON UNIVERSITY
PRINCETON, NJ

NCSA / UIUC
URBANA-CHAMPAIGN, IL

UNIVERSITY OF WASHINGTON
SEATTLE, WA

LSST OBSERVATORY SITE
CERRO PACHÓN, CHILE



LSST.ORG/HIRING



Office of
Science

CHARLES AND LISA SIMONYIFUND

• • • FOR ARTS AND SCIENCES • • •