



# How DM uses precursor datasets and stories from HS

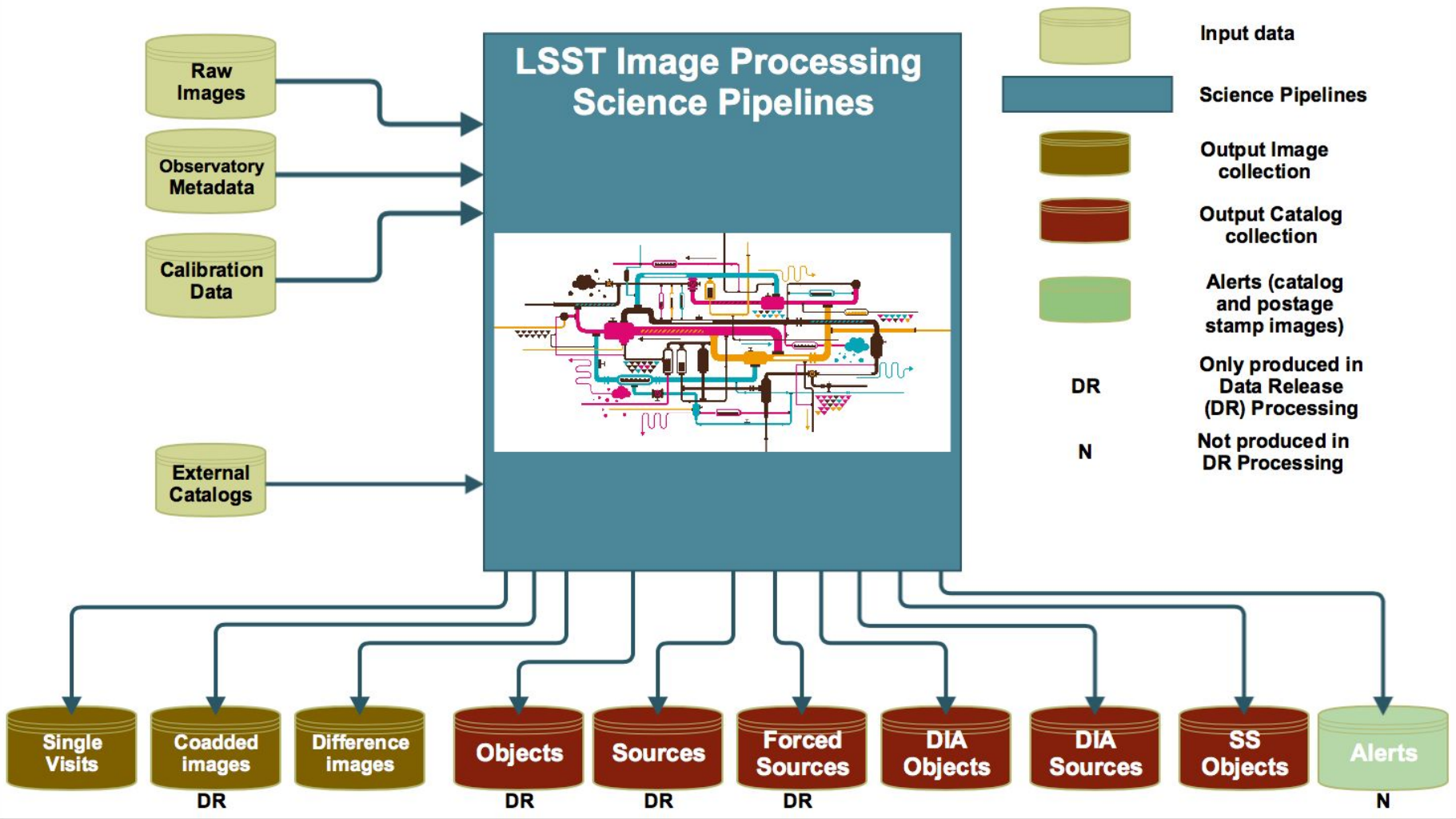
Yusra AlSayyad

Project and Community Workshop  
August 10 2023



U.S. DEPARTMENT OF  
**ENERGY**



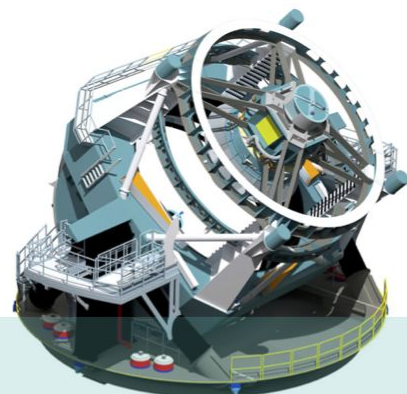


# You will get data products fast and slow

## Raw Data: 20TB/night



Sequential 30s images covering the entire visible sky every few days



Access to proprietary data and the Science Platform require Rubin data rights



## Prompt Data Products

Alerts: up to 10 million per night



via nightly alert streams

Raw & Processed Visit Images, Difference Images, Templates

Transient and variable sources from Difference Image Analysis

Solar System Objects: ~ 6 million



via Prompt Products DB

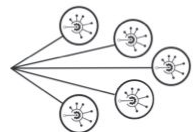
## Data Release Data Products

Final 10yr Data Release:

- Images: 5.5 million x 3.2 Gpixels
- Catalog: 15PB, 37 billion objects



via Data Releases



Community  
Brokers

Rubin Data Access Centres  
(DACs)

USA (USDF)  
Chile (CLDF)  
France (FRDF)  
United Kingdom (UKDF)

Independent Data Access  
Centers (IDACs)

## Rubin Science Platform

Provides access to LSST Data Products and services for all science users and project staff.

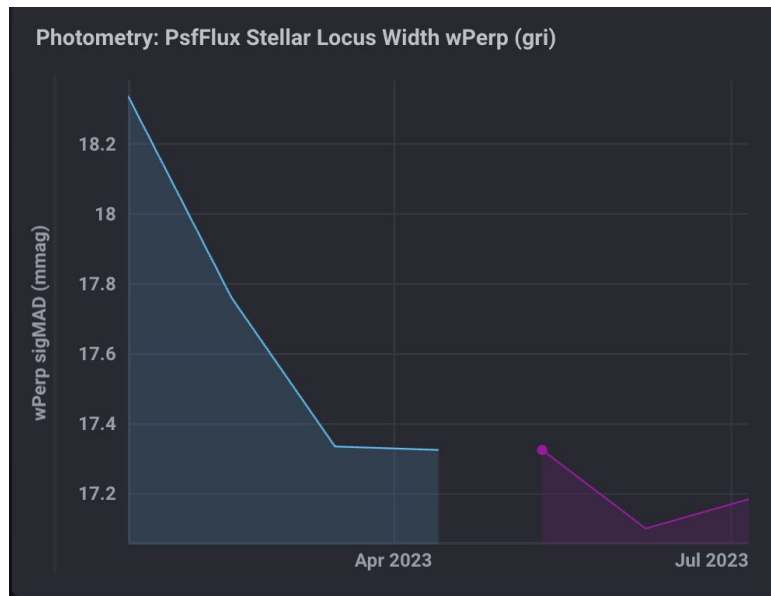
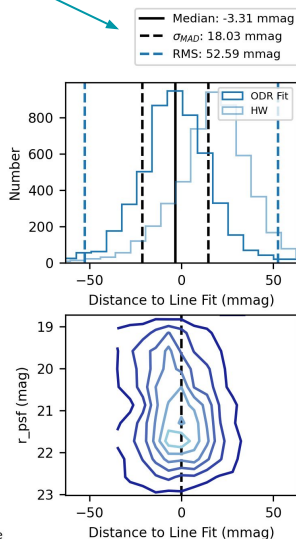
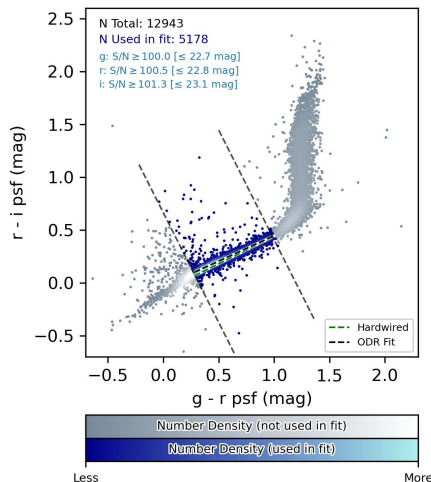


# Plots and Metrics are computed during pipeline execution alongside the algorithms

The **metrics** printed on the **plot**, **match** those written to the butler and dispatched to Sasquatch

colorColorFitPlot\_wFit\_PSF

HSC/runs/RC2/w\_2023\_23/DM-39610/step3/group0/w01\_000  
PhotoCalib: None, Astrometry: None  
Table: objectTable\_tract, Tract: 9615, Bands: r,i,g, S/N > 100.0 (psfFlux)



Sasquatch is the Rubin Observatory service for recording, displaying, and alerting on telemetry data and metrics  
See [sasquatch.lsst.io](https://sasquatch.lsst.io)

# We welcome your pull requests to analysis tools

Last year, the metric and plotting code was refactored and redesigned into a framework/package called `analysis_tools`.

`analysis_tools` was intentionally designed in a modular way such that members of the science community can contribute analysis code to be automatically run as part of science pipelines processing and generate science performance diagnostic plots and metrics

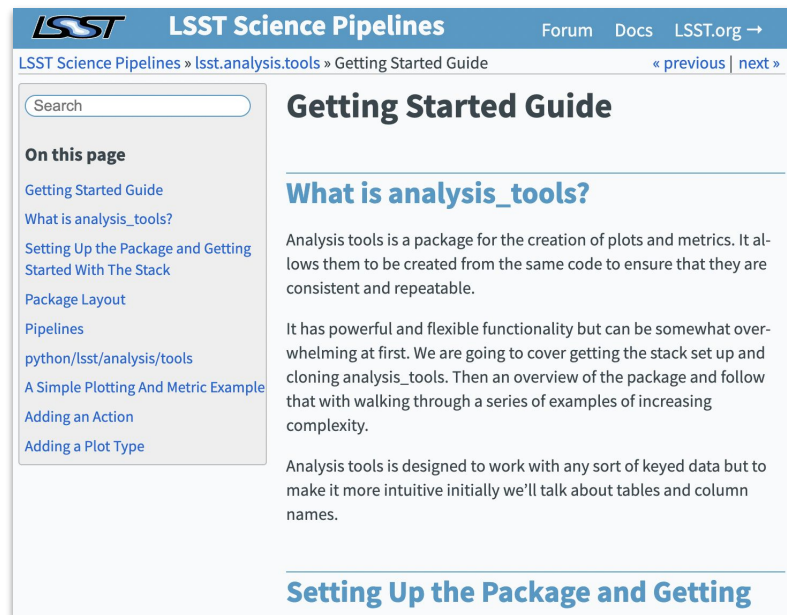
See tutorials from the [May 2023 Commissioning Science Validation Bootcamp](#) and the new [getting started guide](#)



Nate Lust



Sophie Reed



The screenshot shows the 'Getting Started Guide' page for the LSST Science Pipelines. The page has a blue header with the LSST logo and navigation links for 'Forum', 'Docs', and 'LSST.org'. Below the header, the breadcrumb trail reads 'LSST Science Pipelines » lsst.analysis.tools » Getting Started Guide'. A search bar is present. On the left, a sidebar titled 'On this page' lists links: 'Getting Started Guide', 'What is analysis\_tools?', 'Setting Up the Package and Getting Started With The Stack', 'Package Layout', 'Pipelines', 'python/lsst/analysis/tools', 'A Simple Plotting And Metric Example', 'Adding an Action', and 'Adding a Plot Type'. The main content area is titled 'Getting Started Guide' and contains a section 'What is analysis\_tools?' which explains that the package is for creating plots and metrics from consistent code. It also mentions that the package is powerful but can be overwhelming at first. Below this, there is a section 'Setting Up the Package and Getting'.



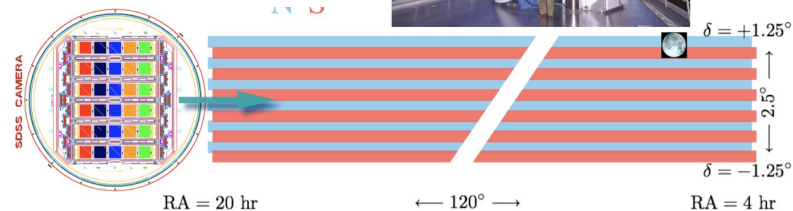
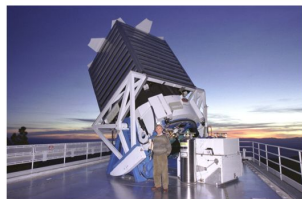
# We've been running the pipelines on 100s of sq deg precursor dataset since the beginning

In the pre-construction era we called them “data challenges”

In 2012 we were validating the coaddition algorithms and forced photometry **on SDSS Stripe 82** ([dmtn-034.lsst.io](https://dmtn-034.lsst.io))

In 2013 we did a joint reprocessing with the FrDF, with some improvements (e.g. background matching)

300 sq deg  
ugriz  
80 epochs over  
10 yr baseline



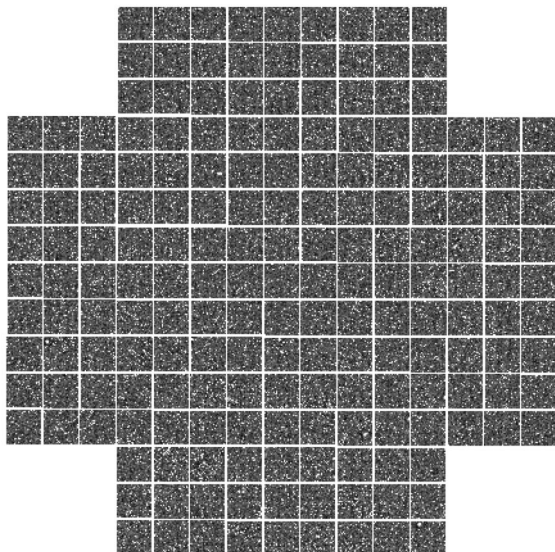
# Since then we regularly process precursor datasets from 4 cameras (1 simulated) with more in common with LSSTCam

Rubin AuxTel  
LATISS



1.2; 6.7 arcmin diam  
1 real LSST CCD

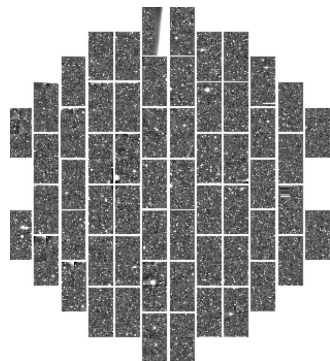
SIMULATIONS **LSST ImSim**  
DESC's DC2 Run2.2i



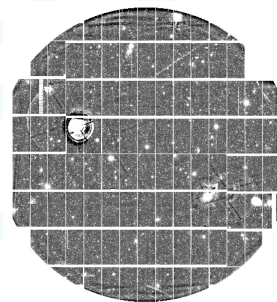
6.5m; 3.5 deg diam  
189 4k X 4k CCDs, ugrizy

Meredith and Lee  
speaking to this next

Dark Energy Camera (DECam) Hyper Suprime-Cam (HSC)  
Subaru Strategic Program



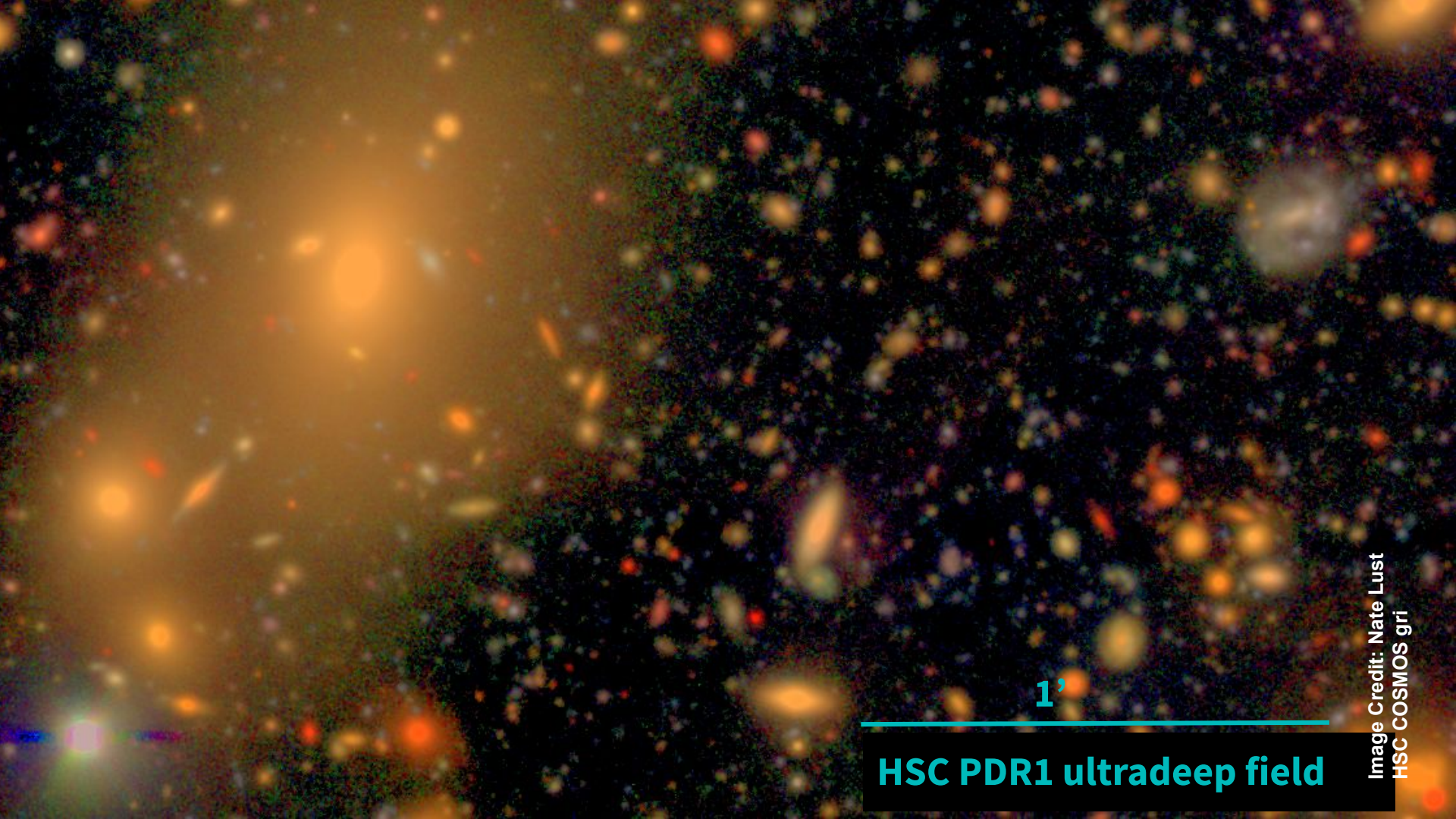
3.9m; 2.2 deg diam  
62 2k X 4k CCDs



8.2m; 1.5 deg diam  
103 2k X 4k CCDs  
grizy

Figures: Pipeline-processed visit images (PVIs) aka calexps





1'

**HSC PDR1 ultradeep field**

Image Credit: Nate Lust  
HSC COSMOS gri



## And the LSST Pipelines are the Hyper Suprime-Cam (HSC SSP) Pipelines

Survey Comparison	LSST	HSC (Subaru Strategic Program)
Effective Aperture	6.5m	8.2m
Filters	ugrizy	grizy + narrow
Exp time per visit	~30s	~240s
Field of View	10 deg <sup>2</sup> 3.5 deg diam	1.8 deg <sup>2</sup> 1.5 deg diam
Num CCDs	189 (4k x 4k)	103 (4k x 2k)

# We run pipelines on precursor data in two modes now

## Fix images vary pipelines vs. Fix pipelines vary images

- 1) As always, we analyze pipeline performance on **fixed datasets** of 3 sizes and cadences:
  - Small areas ( $< 1 \text{ deg}^2$ ) on a **nightly** cadence
  - Medium areas ( $\sim 10 \text{ of deg}^2$ ) on a **monthly** cadence
  - Large areas (100s of  $\text{deg}^2$ ) on **annual** cadence
  
- 2) And now also run pipelines routinely on **real-time AuxTel observations** via 3 campaigns:
  - **Rapid Analysis**
  - **10am DRP Processing**
  - **Nightly Prompt Processing**

[DMTN-091](#): Test Datasets for  
Scientific Performance Monitoring



Colin, Merlin, and AuxTel

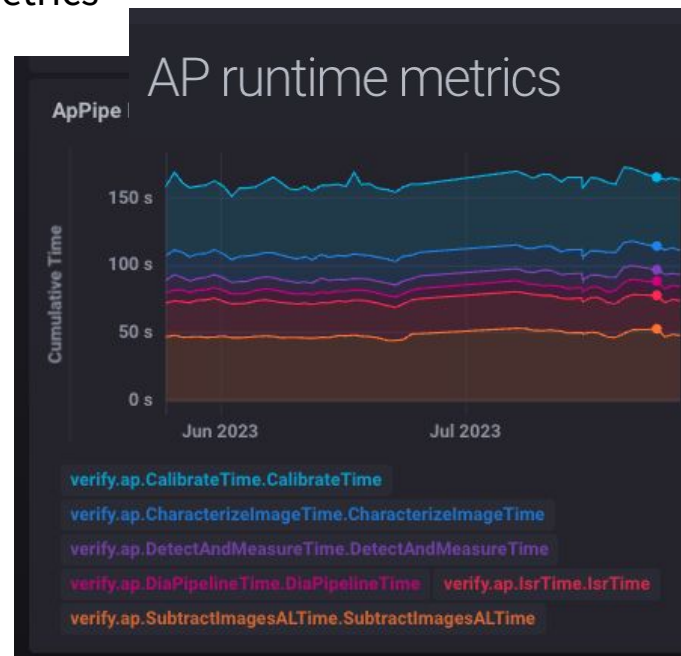
# Daily small ( $<1 \text{ deg}^2$ ) reruns in continuous integration keep pipelines healthy

Two builds are launched during **each nightly software release** that **produce metrics** dispatched to Sasquatch on **Chronograf**.

- **rc2\_subset** tests the **DRP** and is a one patch subset of **HSC RC2** (next slide) that is the same one as the getting started guide on [pipelines.lsst.io](https://pipelines.lsst.io)
- **ap\_verify** runs nightly on a few
  - HSC ccds [ap\\_verify ci cosmos\\_pdr2](#)
  - DECam ccds [ap\\_verify ci hits2015](#)
  - ImSim ccds [ap\\_verify ci dc2](#)

Tracks science quality metrics and comp perf metrics

Runtime (s) on an ImSim ccd



Date of nightly release

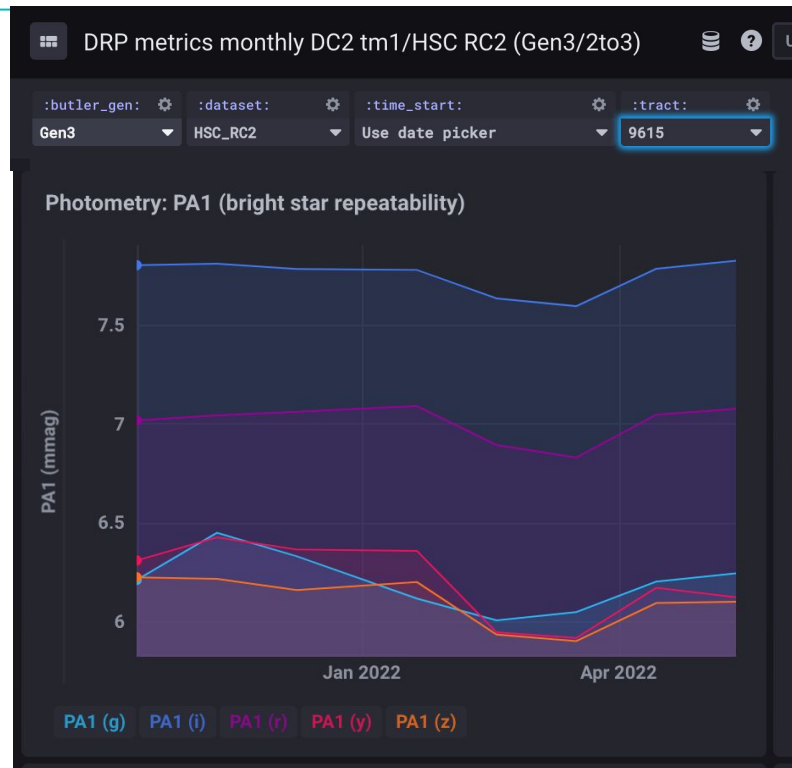
# Monthly medium ( $<10 \text{ deg}^2$ ) reruns provide a testbed for new algorithms

Monthly continuous integration on datasets of  $\sim 5 \text{ deg}^2$

- Tract-sized datasets are the **minimum** to test the speed, robustness, and performance of any algorithmic change. and **plots** as a function of (i.e., hold data fixed and change pipeline)
- AP on HSC COSMOS and DECam HiTS ([Förster+16](#))
- DRP on **HSC RC2** (3 tracts) and ImSim **DC2 2.2i test-med-1** (2 tracts at 1.5 yr depth). All bands.

Next talk!

Photometric  
Repeatability (mmag)

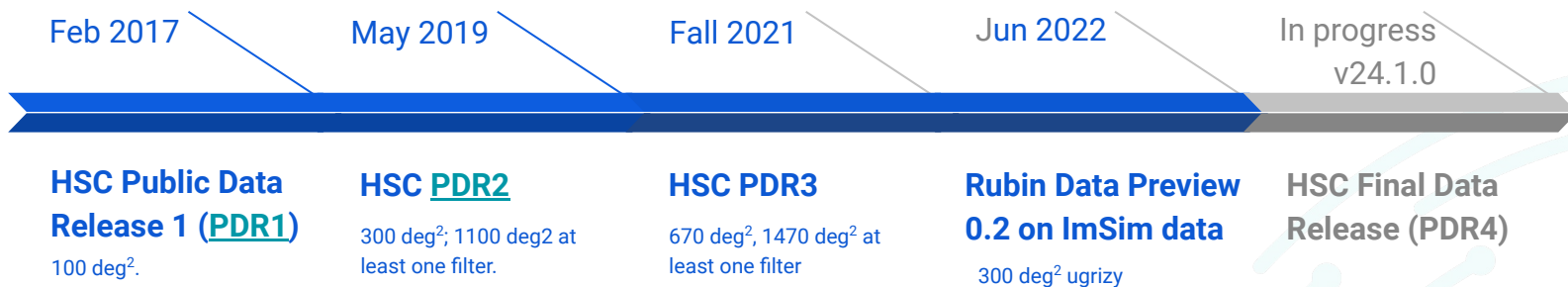


Date of weekly release

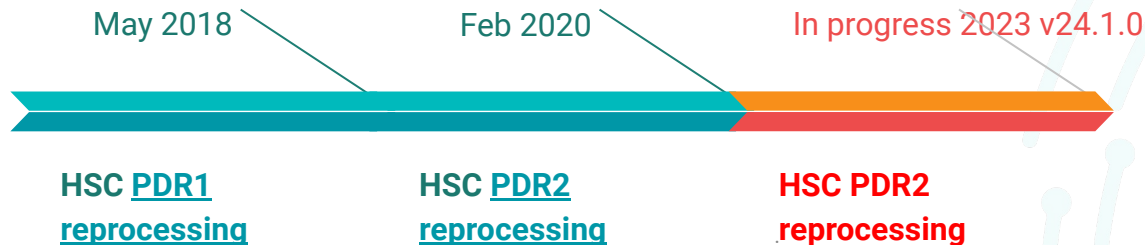


# Large (100s deg<sup>2</sup>) processing campaigns test whole of data management and rare edge cases

Pipelines are used in **external data releases**; the best QA is astronomers publishing papers with your data products:

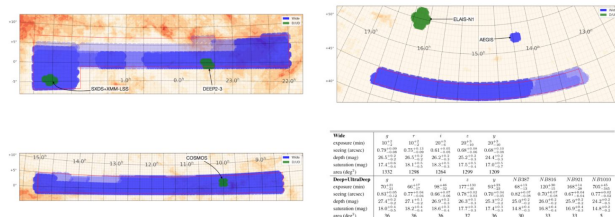


**Internal reprocessing campaigns** integrate all of DM, provide data for characterization reports, and test algorithms for robustness:



- <https://hsc-release.mtk.nao.ac.jp/>
- **Benefits:**
  - Real data of comparable depth
  - Same algorithms and flags
- **Caveats:**
  - The column names are different
  - No Science Platform available
  - Does not include DIA Data Products

HSC-SSP PDR3 includes over 600 square degrees of multi-band data at the nominal survey depth. See the figures below for the survey footprints. The blue and green areas show the Wide and Deep+UltraDeep layers, respectively. The darker blue regions are covered in more filters (max. 5).



The table gives a quick overview of the quality of our data. The depths are given as 5 sigma limiting magnitudes for point sources. Area is the area covered in at least 1 exposure in each filter.

## Data Retrieval

The data can be retrieved in multiple ways. The simplest way to retrieve catalog data is to use the database. We have online/offline SQL tools. For image data, most users will find hscMap, an online image browser, very useful. For binary files, we have a data search tool as well as image cutout tool. All these tools are summarized in the [Data Access page](#). In order to access the data, you first have to [sign up for an account](#). Before you use our data products, we strongly recommend you to go over the [data release paper](#) and the [Known Problems page](#). If you use the HSC data in your publication, please [acknowledge us](#). This site serves only the processed data. Raw data can be retrieved from [SMOKA](#).

## Data Quality

We have performed a number of validation tests for our data products. A complete set of the plots can be found [here](#).

**Hyper Suprime-Cam Subaru Strategic Program**

Public Data Release

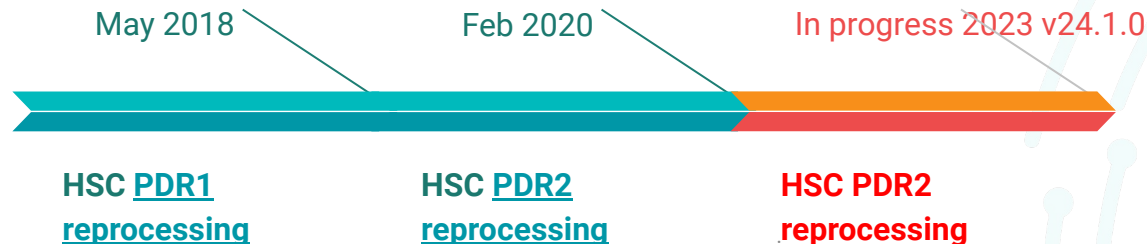
[Home](#)   [Survey](#)   [Processing](#)   [Database](#)   [Available Data](#)   [Data Access](#)   [FAQ](#)

# Large (100s deg<sup>2</sup>) processing campaigns test whole of data management and rare edge cases

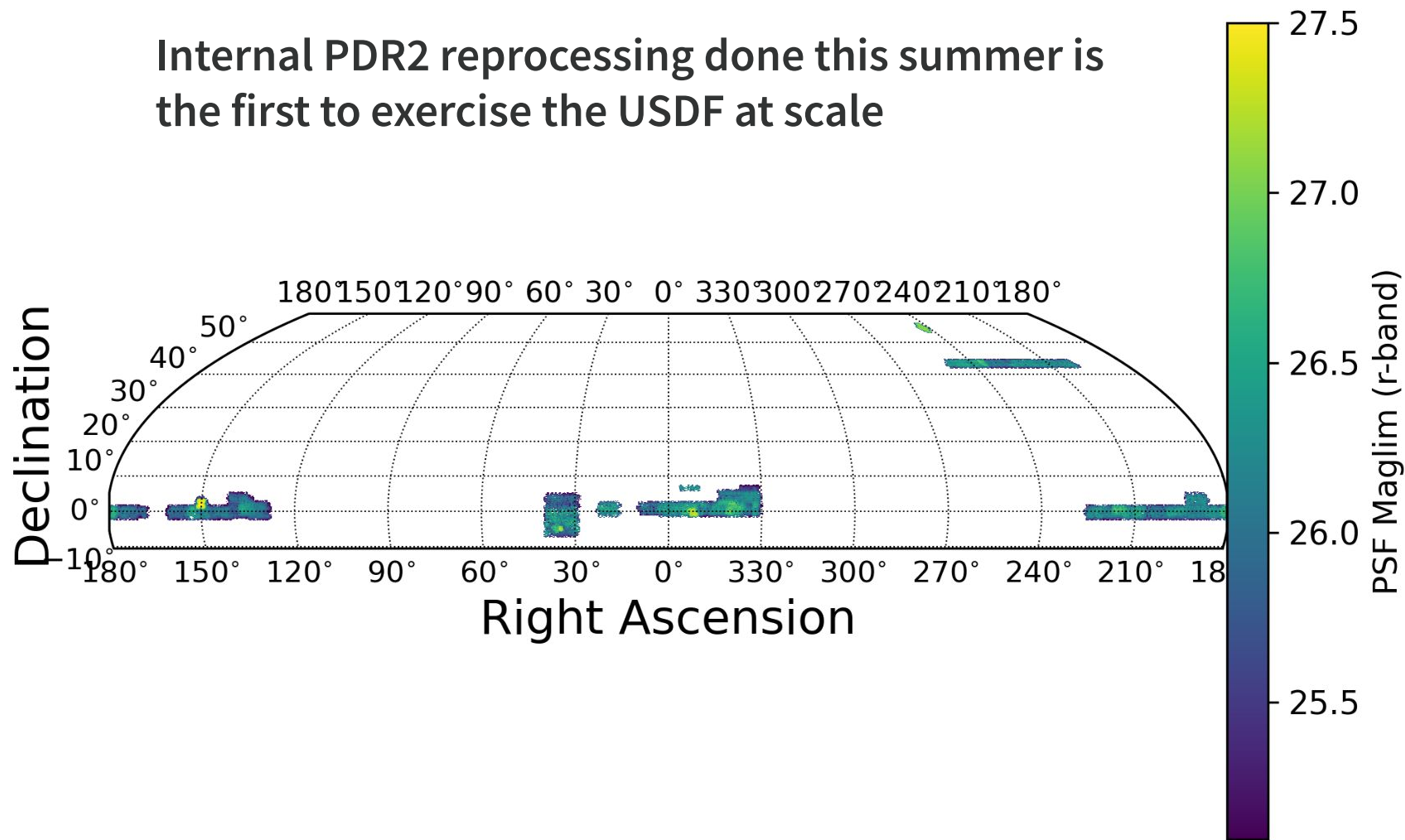
Pipelines are used in **external data releases**; the best QA is astronomers publishing papers with your data products:



**Internal reprocessing campaigns** integrate all of DM, provide data for characterization reports, and test algorithms for robustness:



Internal PDR2 reprocessing done this summer is  
the first to exercise the USDF at scale









# Stories from HSC

(if there's time)



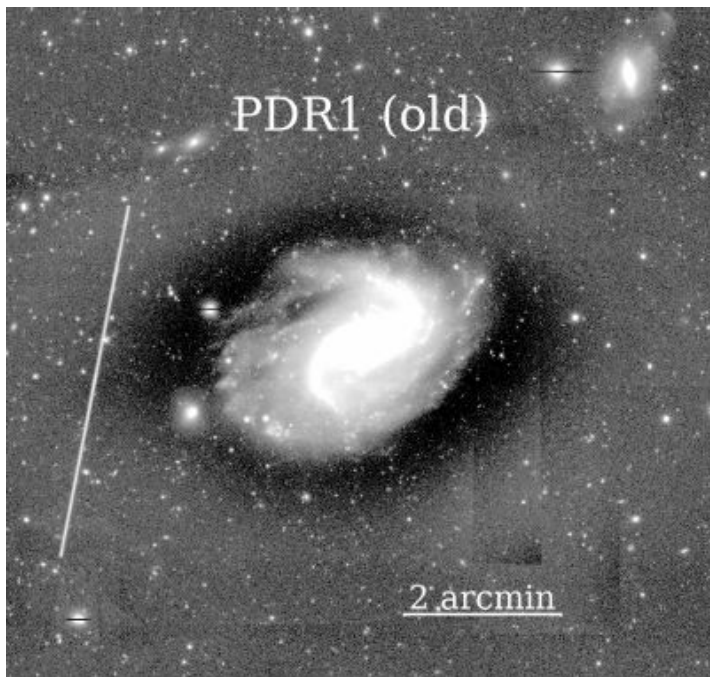
U.S. DEPARTMENT OF  
**ENERGY**

**SLAC**

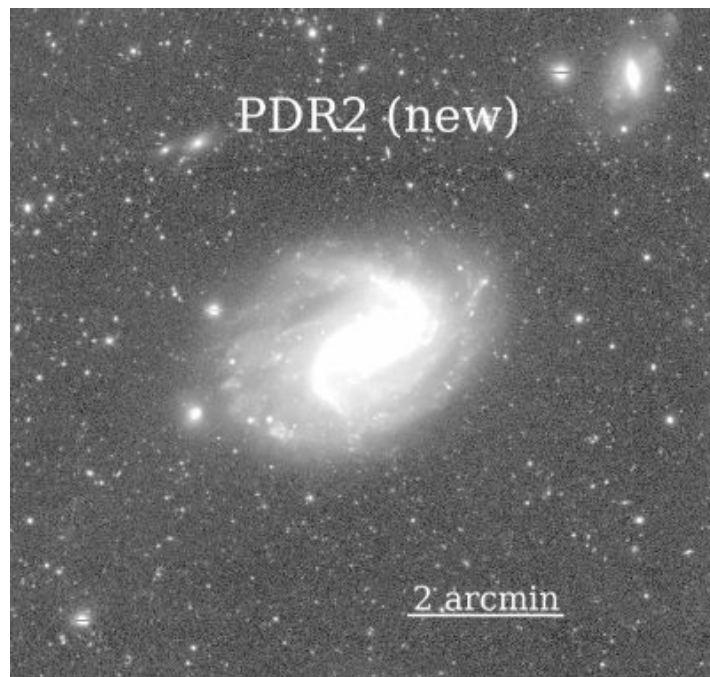
CHARLES AND LISA SIMONYI FUND  
... FOR ARTS AND SCIENCES ...

**LSST**  
CORPORATION

## Low Surface Brightness community was happy with PDR2 full focal plane “SkyCorrection”



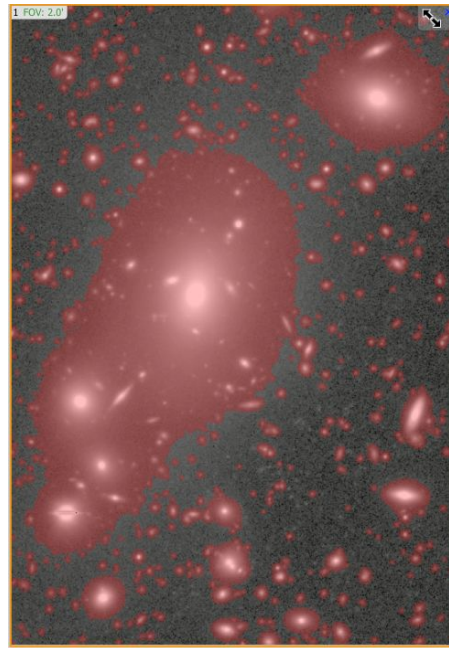
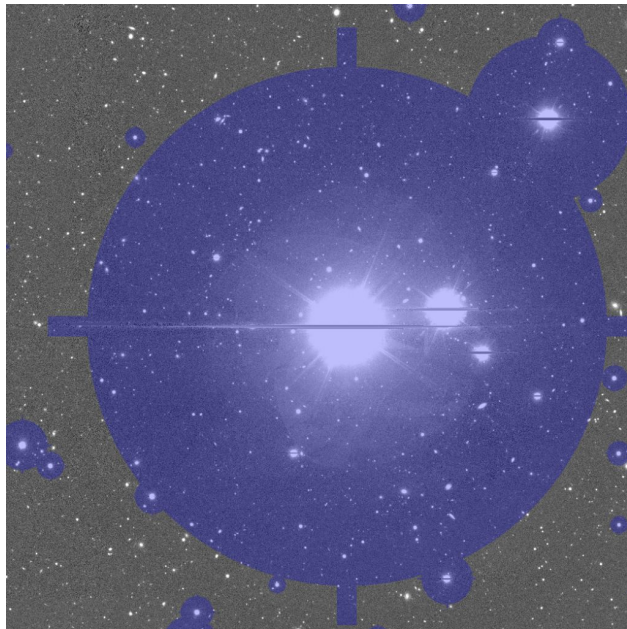
Coadd with  
PDR1 Local Background subtraction



Coadd with  
PDR2 Focal Plane Background subtraction

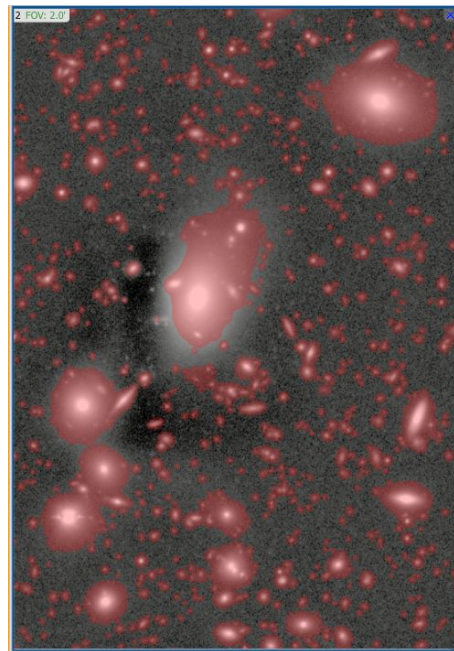
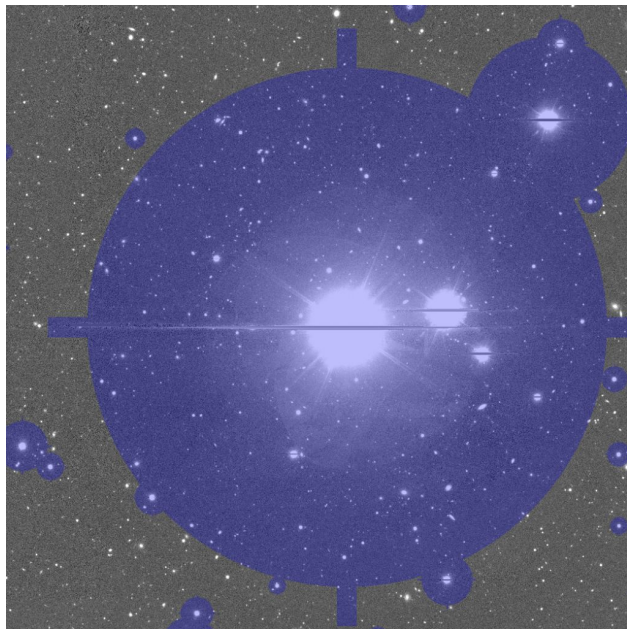
Aihara+19 (PDR2 release  
paper)

## But everyone else was unhappy





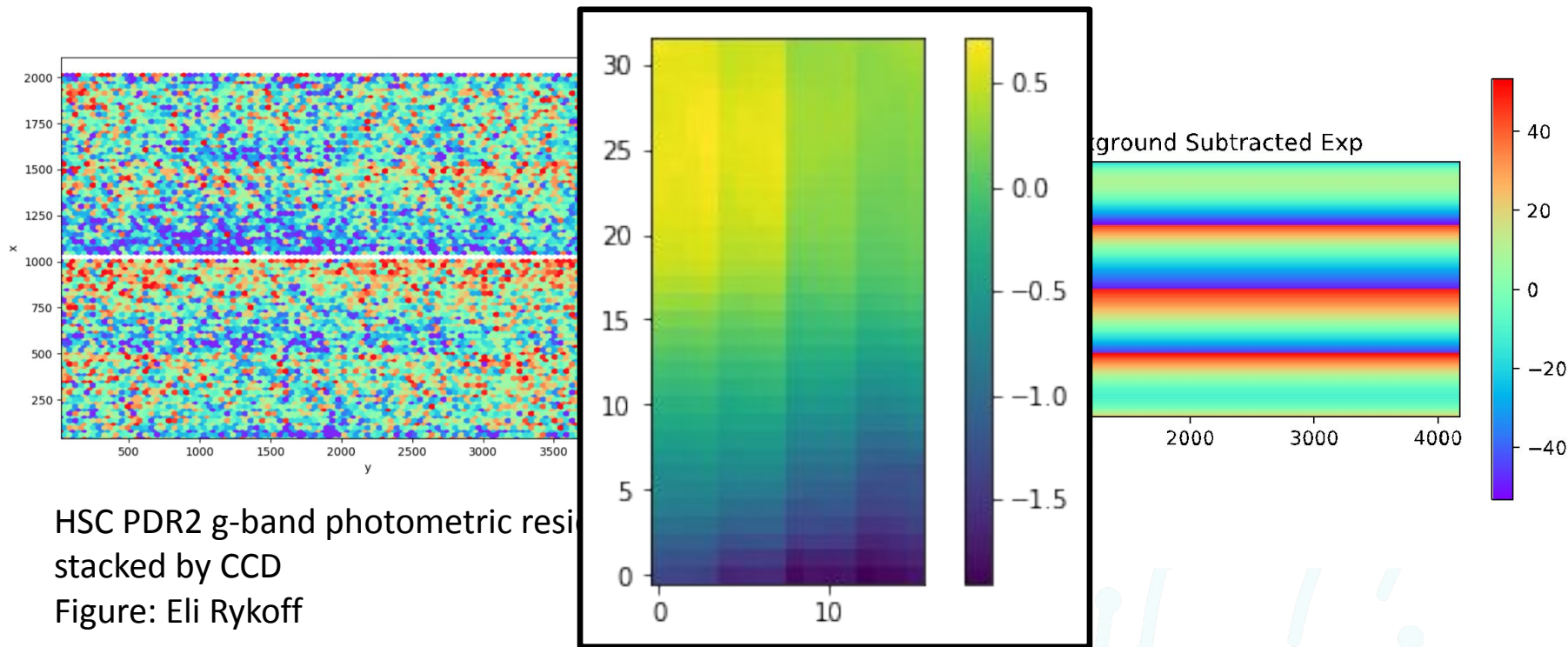
## But everyone else was unhappy



PDR3 adds a very aggressive 128x128 binned spline background subtraction on the deepCoadd\_calexps

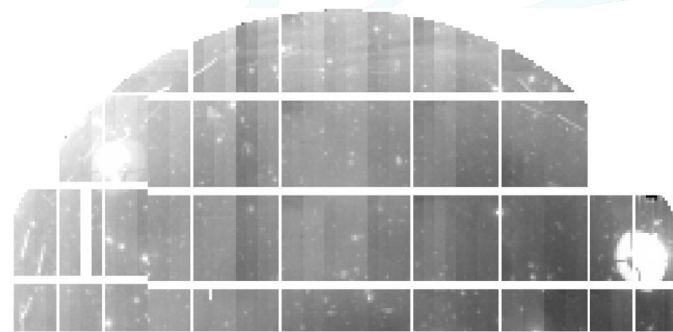
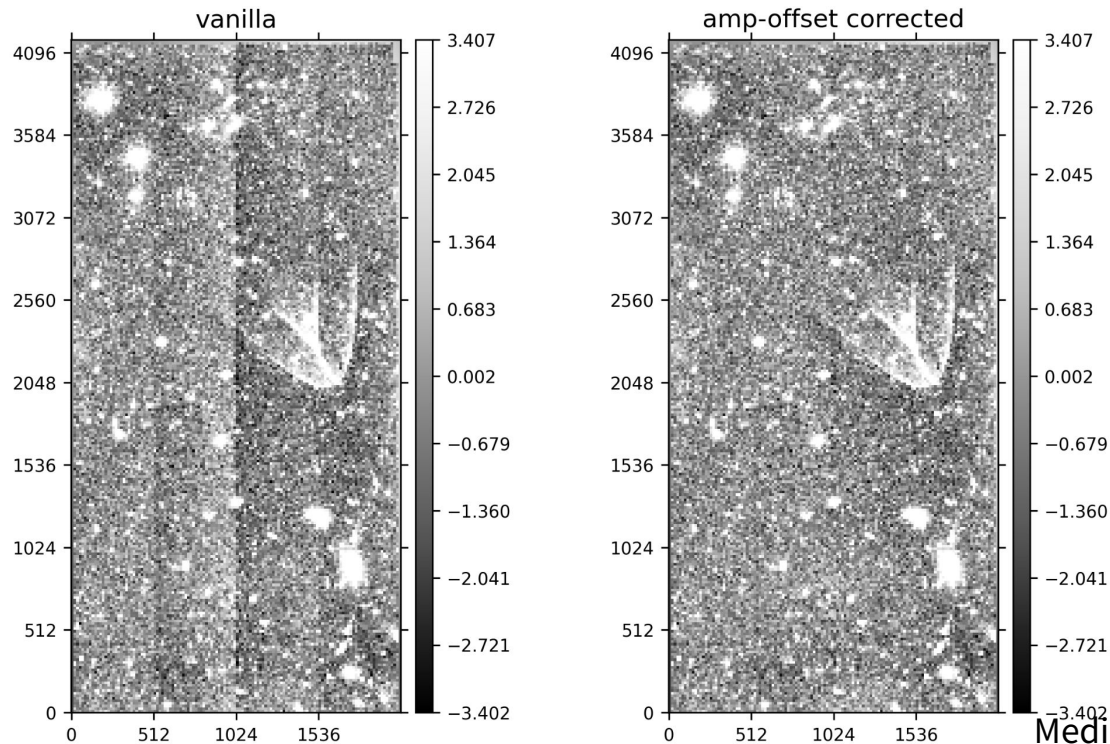
deepCoadd still available for LSB detection

## Large area also alerted us to very small amplitude systematics



# Which led to the addition of the the Pan-STARRS “pattern continuity” (a.k.a amp-to-amp matching)

{'visit': 1252, 'ccd': 68} (HSC-I, 30.0s)



Lee Kelvin

Median binned (128x128) HSC-I visit=1252 (short exp)



# Some answers to the discussion questions



U.S. DEPARTMENT OF  
**ENERGY**

**SLAC**

CHARLES AND LISA SIMONYI FUND  
... FOR ARTS AND SCIENCES ...

**LSST**  
CORPORATION



Do HSC data releases use a recent version of the LSST pipelines, or rather very similar underlying algorithmic code that was “forked off” some time ago?

- For HSC PDR1 we **forked** the pipelines hard, meaning that there was development on the fork that was not mirrored upstream. We had a bad time upstreaming changes back to main. Not recommended!
- For HSC PDR2/3 we did a **soft fork**: see [hscPipe](#). Bugfixes were applied to the LSST Science Pipelines and backported to the fork.
- For HSC PDR4, we are maintaining v24.0.x as the **release branch** on the LSST Science Pipelines. Any bugfixes we know will be needed for PDR4 are backported to v24.
- For DP0.2. We maintained the v23.0.x release branch. Any bugfixes needed for DP0.2 were backported. See process outlined here: <https://developer.lsst.io/work/backports.html>

When processing precursor data sets, how do you balance keeping up with recent LSST pipelines releases versus maintaining some “stable” choice of LSST pipelines version?

Depends on how long your processing will take:

For >100s of sq deg data release:

- We maintain a **release branch** on the LSST Science Pipelines. Bugfixes are backported to that release branch.
- For DP0.2. We maintained the v23.0.x release branch. Any bugfixes needed for DP0.2 were backported. See process outlined here: <https://developer.lsst.io/work/backports.html>

*For the monthly reprocessings we're a little more loose, use w\_2023\_XX and use ticket branches if we need to patch anything.*

*For the nightlies, we just use the daily release. If it fails, too bad, try again the next night.*

What types of computational resources are your large LSST pipelines processing of precursor data sets deployed on? Cloud? Academic HPC?

*Data releases:*

- *DP0.2 was produced in google cloud*
  - *ctrl\_bps\_panda*
  - *storage: s3*
- *HSC public data releases at the NOAJ Mitaka cluster*
  - *ctrl\_bps\_htcondor*
  - *Condor over PBS*
  - *Storage: gpfs*

*Monthlies:*

- *USDF, ctrl\_bps\_panda and ctrl\_bps\_htcondor over slurm.*

Do you have any tips for deploying the LSST pipelines at scale in an efficient way, for instance database optimizations (or similar) that make a big difference?

For the tiny rc2\_subset example in the getting started guide uses sqlite as the registry and pipetask -j. For larger datasets you'll need a real DBMS (USDF uses postgres for it now) and BPS: <https://pipelines.lsst.io/modules/lsst.ctrl.bps>

When processing precursor data sets, how do you manage large output data volumes from the LSST pipelines? Do you delete any intermediate products?

- *Yes. Process all the data through singleframe processing, then all the data through calibration, all the data through coaddition etc... rather than region by region all the way to the end. When you get to the end of a stage and validate it, you can delete the intermediates. The **per-visit intermediate exposures** in particular take space. These include:*

DatasetType	Name	producedInStage	Can be deleted...
icExp	1	After single-frame processing is validated	
postISRCCD	1	After single-frame processing is validated	
deepCoadd_directWarp	3	After coadds are validated	
deepCoadd_psfMatchedWarp	3	After coadds are validated	
goodSeeingDiff_templateExp	4	After DIA outputs are validated	
goodSeeingDiff_matchedExp	4	After DIA outputs are validated	
<u>goodSeeingDiff_differenceTempExp</u>	<u>4</u>	<u>After DIA outputs are validated</u>	



Are there any upcoming LSST pipelines developments that those of us working on precursor data sets like DECam and HSC should particularly be on the lookout for?

- *There's a memory leak in the PiffPsfs that cause trouble when you coadds 100s of visits deep.*
- *QuantumBackedButler (QBB) has landed (w\_2023\_31) and is being tested now.*

Have you (or any other groups you know of) done multi-instrument forced photometry on precursor data sets using the LSST pipelines?

- *Talk to Raphael Shirley about his joint processing of HSC and VISTA data with the pipelines. He treated the two surveys as different bands of the same survey, and achieved good results.*
- *If you just have a list of ra/decls, check out `ForcedPhotCcdFromDataFrameTask`*

What resources do you recommend for someone learning to run the LSST pipelines on precursor data?

- *The starting point is the getting started guide on <https://pipelines.lsst.io/v/weekly> It's a tiny amount of data, just so you can get from start to finish in a day and know where you're going with lots of data.*
- *If you're ready to ingest your data into a Gen3 butler take a look at **Lee's guide**.*
- *Then it's off to [community.lsst.org](https://community.lsst.org) for help*

*I highly recommend watching Jim's Tour of Middleware talk at the 2022 Pipelines Bootcamp:  
<https://confluence.lsstcorp.org/display/DM/DM+Pipelines+Bootcamp+2022>*

How actively are DECam calibration procedures/algorithms within the LSST pipelines being developed/upgraded/changed?

What is the status of the LSST pipelines as far as run-time optimization? Is that currently a point of emphasis, or are the LSST pipelines already viewed as meeting their formal requirements in this regard?

Same question but for memory usage rather than run-time.

- *Early on we optimized the most expensive kernels on then-hardware (warping and convolution especially). We've kept an eye on scalability (e.g. swapped jointcal out for GBDES. Working on cell-based coadds). Now hardware/DF We've started profiling the higher level Task code and pipelines.*

Possibly likely question from the audience: have the LSST pipelines been adapted for [fill in the blank] facility that I use/manage/operate yet?