

Using Parsl for Image Processing of Rubin Observatory Data

Jim Chiang
SLAC
ParSl & funcX Fest 2022
2022-09-13

Thanks to the Parsl & funcX Fest organizers and the Parsl developers!



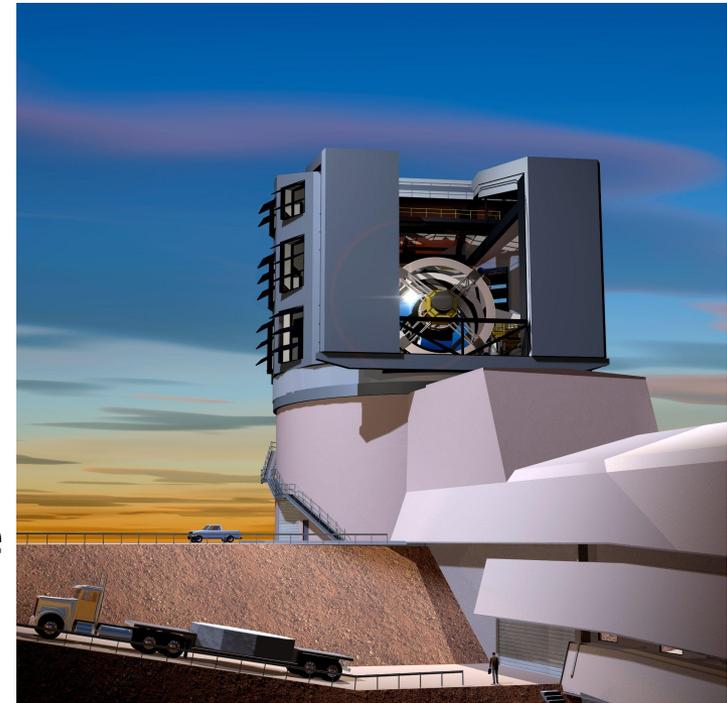
Vera C. Rubin Observatory



Starting in ~ 2024 , [Vera C. Rubin Observatory](#) will conduct the 10-year Legacy Survey of Space and Time (LSST). LSST is designed to address four science areas:

- Probing dark energy and dark matter
- Taking an inventory of the Solar System
- Exploring the transient optical sky
- Mapping the Milky Way

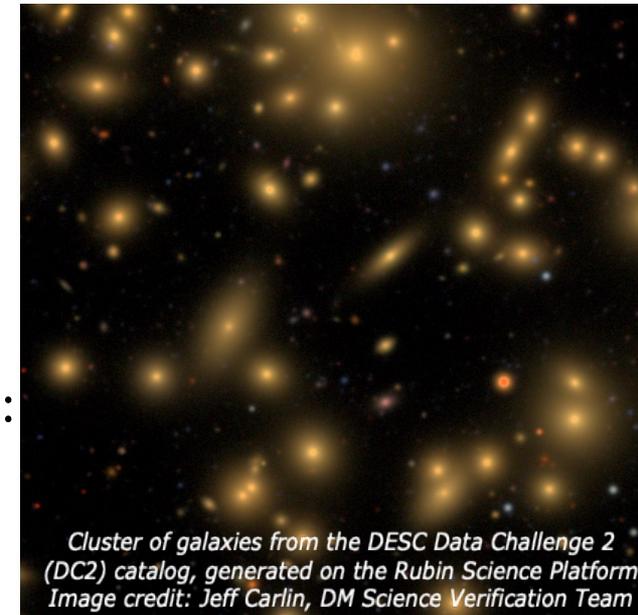
Sited atop Cerro Pachón, Chile, LSST will survey the Southern sky every ~ 4 -5 nights in six bands, *ugrizy*. The LSST survey will cover ~ 20 k sq degrees with a ~ 10 sq degree FOV. Each location will be observed ≥ 800 times with 30 s exposures.



Dark Energy Science Collaboration



- **Our scientific aim:** Exploring the physics of the Dark Universe
 - Dark energy, dark matter, neutrinos and signatures of inflation
- **Our objectives:** (see lsstdesc.org for more details)
 - ⚙️ Accurate cosmology
 - 🎓 A vibrant & inclusive scientific community
 - 🏭 Meeting LSST's big data challenge
 - 🗣️ Learning continuously from each other
- **Our approach:** Combining five cosmological probes:
 - **Clusters of galaxies, large-scale structure, supernovae, weak and strong lensing**
- **Our challenges:** Systematics and more systematics
- **Our requirements:** [DESC Science Requirements Document](#)
- **Our plan:** Prepare with simulated and precursor data



We have built DESC as a vibrant, large (>1200 members), international collaboration over the last 10 years to explore the physics of the Dark Universe with Rubin data!

Rubin LSST Observations and DESC Science



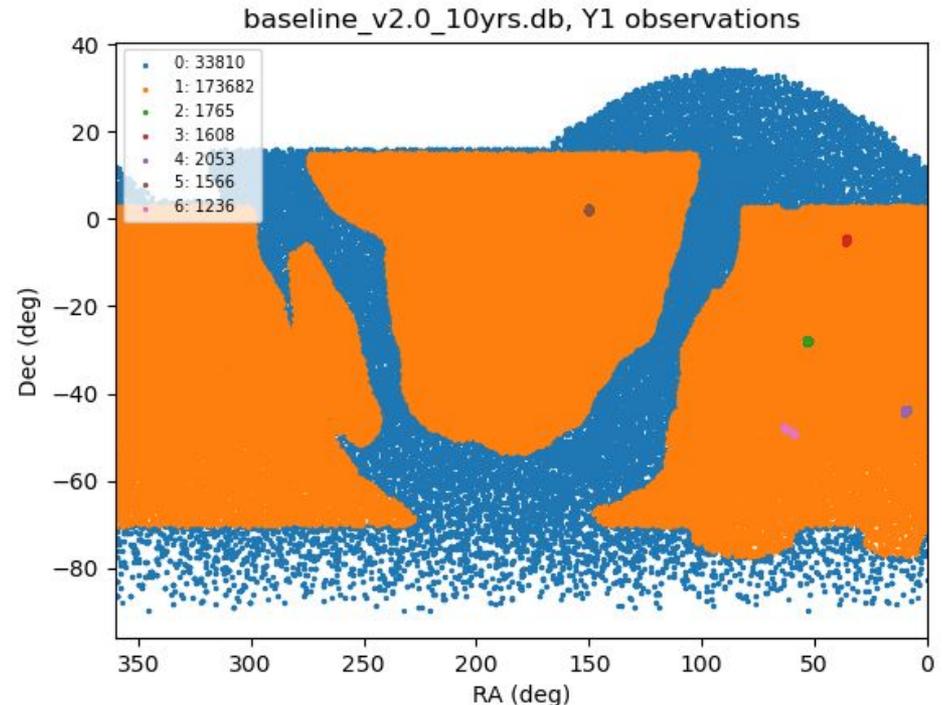
DESC Science is enabled by Rubin's observing sequence: hundreds of short "visits" (~ 30 s) for each location. For each location, the data are co-added to produce deep exposures from which the properties of faint galaxies can be measured.

To characterize systematic uncertainties, DESC will reprocess the LSST data, applying alternative analysis algorithms and data selections.

For Year 1 observations, there are

- 216k visits
- **41M CCD-visit combinations**
- 464k "patches" on the sky
- **3M co-added images** (6 bands)

The Rubin focal plane consists of 189 CCDs, and each CCD-visit or co-add corresponds to a "task instance" (or job) for particular task types.



The Rubin pointing directions for the baseline cadence for Y1 observations. Orange points are the WFD survey, blue points are Galactic plane survey, and remaining locations are "deep drilling fields"

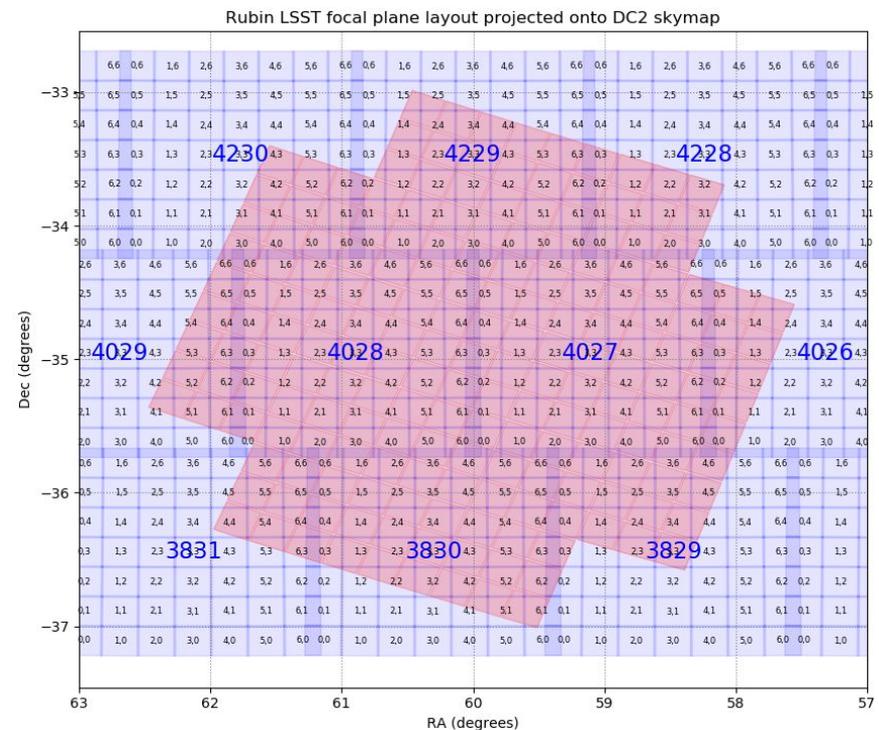
Properties of Rubin data and how they affect workflow complexity



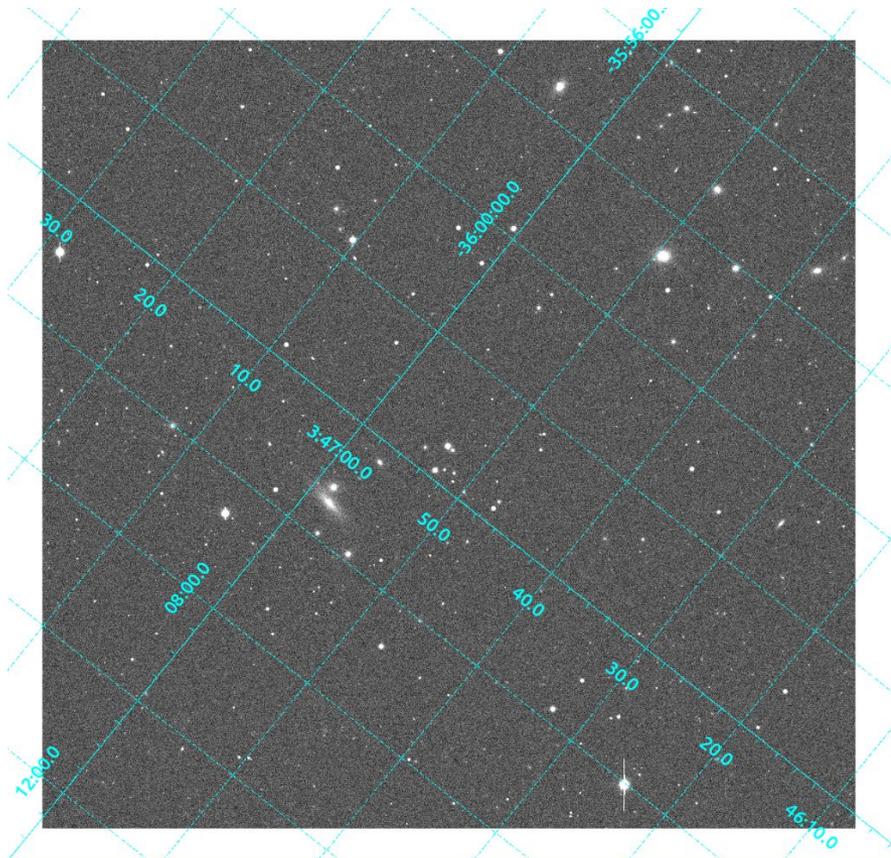
To mitigate systematics arising from the layout of the focal plane, features of the CCDs, etc., each **visit** is randomly dithered in angle and position to average out these effects. Ultimately, this dithering contributes substantially to the complexity of the image processing workflow.

This figure shows a projection of the Rubin focal plane onto the sky. The 189 individual CCDs are plotted in **red**, and the “tracts” and “patches” that divide the sky into manageable units are plotted in **blue**.

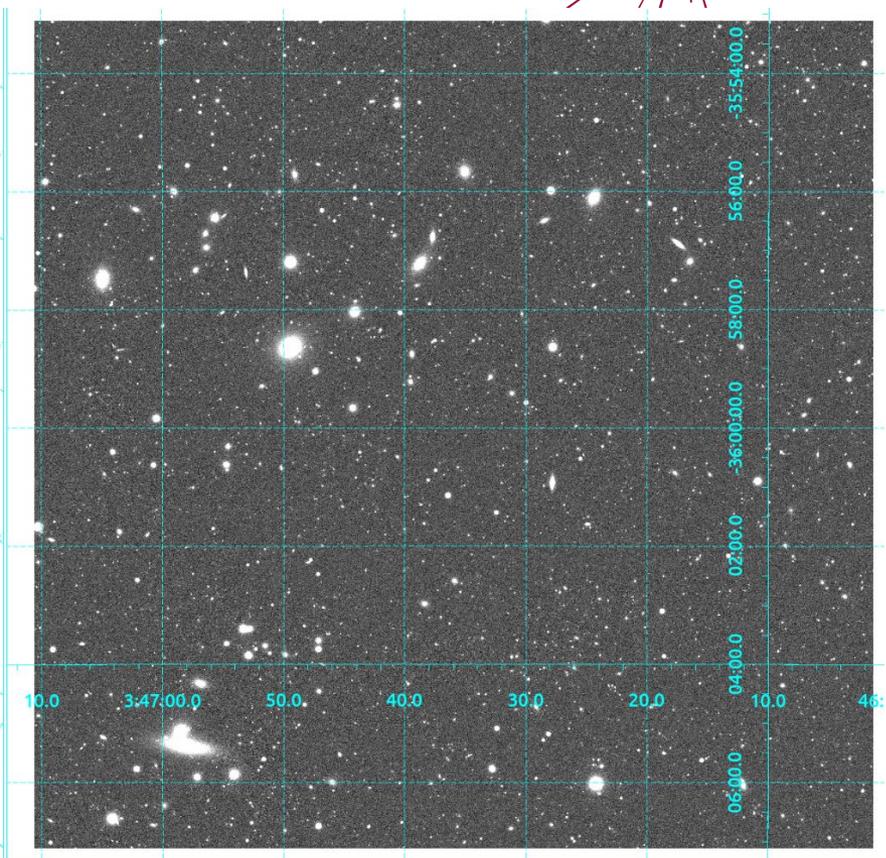
The random placements of the projected focal plane mean that a many different CCD-visits will contribute to the co-added image for a given patch.



Single CCD-visit vs co-added image



Raw CCD image for a single DC2 i-band 30s visit. Visit-level processing removes instrumental effects and performs the initial astrometric and photometric calibration.



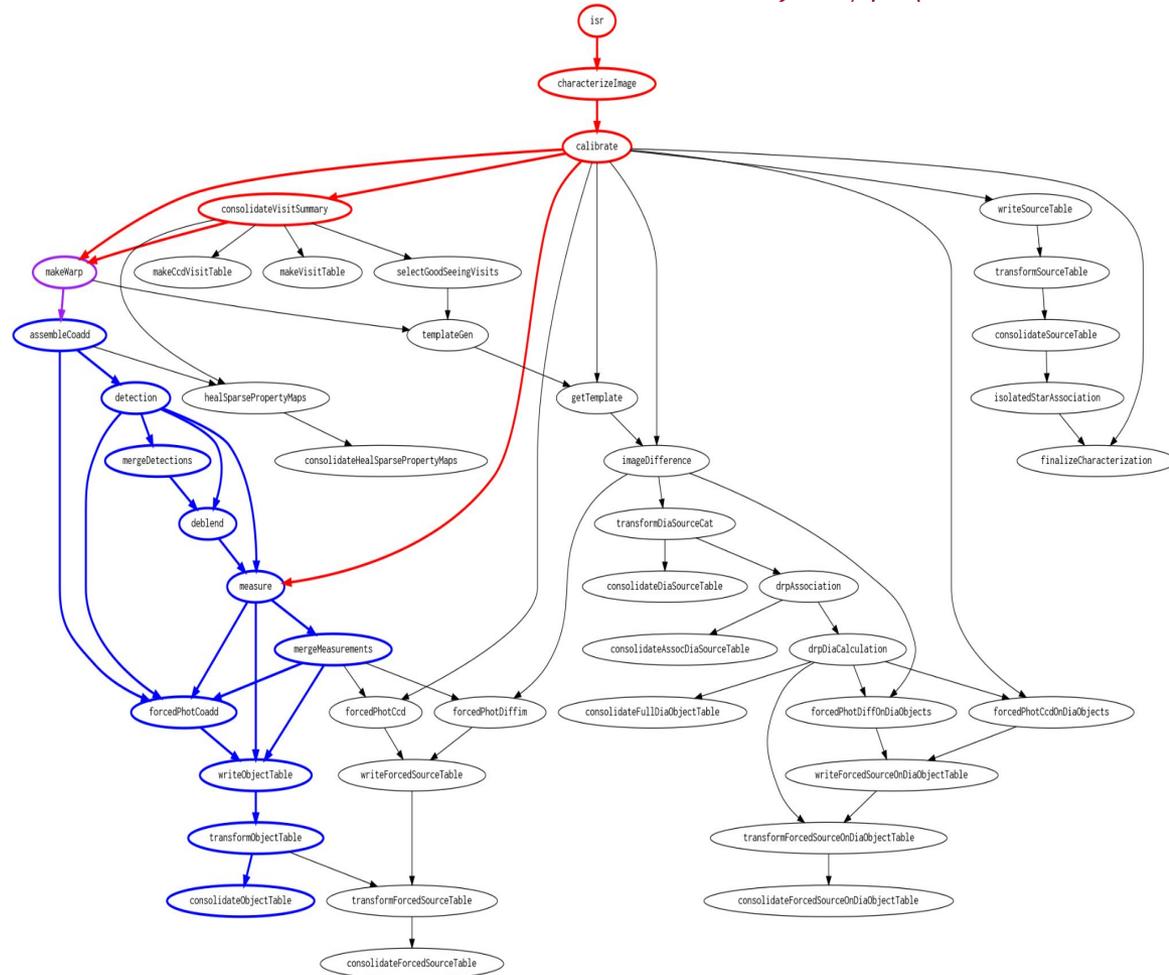
Co-added DC2 image for a "patch" on the sky covering the ~same region. This co-add combines ~18 i-band visits acquired during the Y1 survey.

Rubin LSST image processing steps

The bolded parts correspond to the processing needed by DESC “static” science, i.e., cosmological inference based on measurements of galaxies.

The tasks in **red** operate on CCD-visits, **blue** tasks on co-added images/patches, and **purple** on both.

All of these tasks are implemented as single-threaded jobs, and for a given task type, they are several millions of instances for Y1 data.



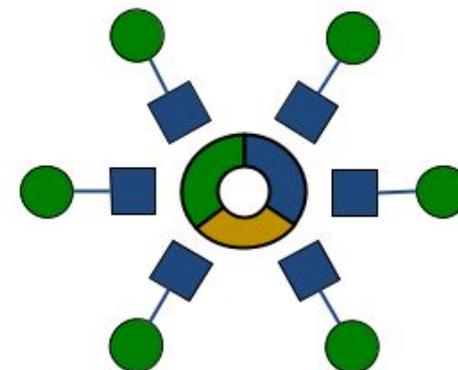
Graph showing dependencies between task types for Rubin image processing.

Running the Rubin Pipelines with Parsl



NERSC is DESC's main computing resource, so we will use HPC resources for this processing.

- Individual jobs for a given pipeline task are embarrassingly parallelizable and are wrapped as Parsl bash apps.
- The Rubin code computes the DAG for the 10s of thousands of jobs in a typical payload, i.e., a set of pipeline tasks + data selection.
- The resources used by the jobs vary significantly both in memory required ($\lesssim 1$ GB to ~ 10 s GB) and CPU time (\sim mins to 10s hours), so Parsl's **WorkQueueExecutor** is a natural candidate for managing execution on large (~ 1000 nodes) HPC allocations.

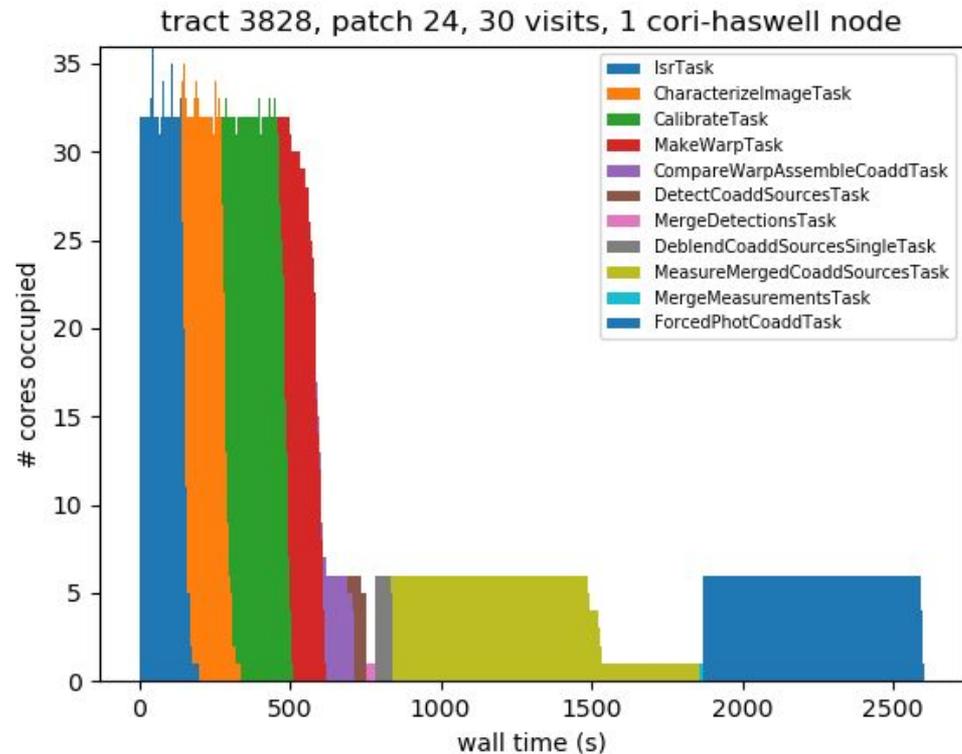


Higher-level orchestration of the Rubin pipelines



However, we need an additional orchestration layer to manage the overall execution of the Rubin image processing workflow:

- Transition from processing in CCD coordinates to sky coordinates is complicated by the many-to-many mapping between CCD-visits and patches.
- Structural bottlenecks prevent effective backfilling of resources.
- These issues suggest dividing into sub-pipelines:
 - Visit-level processing
 - Warp generation
 - Co-add assembly
 - Deblending
 - Object measurement.



Rubin static pipeline run on a small test data set for single patch.