

Data Preview 0: Definition and planning.

William O'Mullane

2020-06-02

1 Introduction

Table 1 shows the FY21 milestones for the Vera C. Rubin Observatory, many of which concern, or relate to, data previews. Section 2 defines what Data Preview 0 is about and covers possible risks and mitigations to that definition. Section 6 Sets out the planning for achieving DP0.

Table 1: Milestones for Rubin Observatory FY21

Milestone	Label	Year	Q	Type	Team
Read only Gen3 butler for DP0 at IDF	DP-MW-M-01	FY21	Q1	Code Release	Science Users Middleware
IDF DP0-Ready: Complete IDF installation and IDF staff preparations for DP0.	DP-IN-01	FY21	Q1	Event	Infrastructure and Support
Submit FY20 POP Annual Progress Report to NSF (via OIR Lab POPPR)		FY21	Q1	Reporting	Rubin Observatory Directorate
Open up the Help Desk for LSST Operations staff to use.		FY21	Q1	Process Definition	Community Engagement
Establish new media presence on at least one new channel for Rubin operations.		FY21	Q1	Process Definition	Outreach
Produce FY22 POP input for DOE budget briefing		FY21	Q2	Reporting	Rubin Observatory Directorate
DP0.1 Early Access: Provide access to processed images and visit level catalogs from the IDF	DP-SP-01	FY21	Q2	Data Release	Science Platform and Reliability Engineering
Announce Initial Survey Strategy		FY21	Q2	Event	Survey Scheduling
Deliver Q1 Report to NSF on POP21 status (via OIR Lab Q report)		FY21	Q2	Reporting	Rubin Observatory Directorate
USDF Decision: obtain confirmed location of US Data Facility		FY21	Q3	Event	Rubin Observatory Directorate
Deliver Q2 Report to NSF on POP21 status (via OIR Lab Q report)		FY21	Q3	Reporting	Rubin Observatory Directorate
Identify Observatory Operations Team Leads (Observatory Software and Summit and Engineering Operation) or launch external searches.		FY21	Q3	Hiring	Observatory Operations Management
Submit FY23 DOE FWP(s)		FY21	Q3	Reporting	Rubin Observatory Directorate
Gen3 butler backed by S3 for processing DP0	DP-MW-M-02	FY21	Q3	Code Release	Science Users Middleware
DP0.2 Reprocessing Start: Begin early DRP-like re-processing of DP0 simulated image data, at the IDF.	DP-SP-02	FY21	Q3	Event	Execution
Begin operations of IDF to support shared-risk simulated data distribution to community		FY21	Q3	Event	Community Engagement
Demonstrate EPO interface with DP0	DP-SP-M-01	FY21	Q3	Process Definition	Science Platform and Reliability Engineering
Stand up Users Committee so that it is active for DP0.1 feedback from Science Collaborations.	DP-PM-01	FY21	Q3	Process Definition	Community Engagement
DP0.1 Data Release: science-ready catalogs released from the IDF		FY21	Q3	Data Release	Verification and Validation
USDF Transition Plan: work with selected USDF team to plan start-up of USDF.	DP-SP-03	FY21	Q3	Process Definition	Data Production Management
Deliver Q3 Report to NSF on POP21 status (via OIR Lab Q report)		FY21	Q4	Reporting	Rubin Observatory Directorate
DP0.2 Early Access: Provide access to reprocessed images and visit level catalogs from the IDF	DP-SP-04	FY21	Q4	Data Release	Science Platform and Reliability Engineering
Deploy early instantiation of service desk providing second-tier technical support for community	DP-SP-M-02	FY21	Q4	Event	Science Platform and Reliability Engineering
Submit FY21 Management Report for AURA Operations of LSST to NSF		FY21	Q4	Reporting	Rubin Observatory Directorate
Submit FY22 POP to NSF (via OIR Lab POP)		FY21	Q4	Reporting	Rubin Observatory Directorate
Incorporate ComCam and/or simulated data into 2 EPO formal education investigations		FY21	Q4	Data Release	Education

Establish Communications Strategy for Operations		FY21	Q4	Process Definition	EPO Management
Establish Communications Strategy for Operations		FY21	Q4	Process Definition	Communications

2 Data Preview 0

In LSO-011 we outlined a number of scenarios for early releases of Rubin Observatory data. The purpose of these releases are not only to prepare the community for LSST data, but also to serve as an early integration test of existing elements of the Data Management systems and to familiarize the community with our access mechanisms.

Two major new developments have occurred since LSO-011 was drafted:

- There have since been delays in construction such that we are now planning on making Data Previews with Rubin Observatory simulated data or on-sky data from other observatories (see Section 3.1) which would still allow us to meet some of the goals of the early releases.
- We are planning on carrying these activities at the Intermediate Data Facility, which is dedicated to Pre-Ops activities infrastructure needs such as serving data and training operations staff (commissioning activities will continue at NCSA and Chile).

In this document we outline notable elements of DP0, the first of these planned data previews, from the Data Management and Pre-Operations perspective.

3 Elements of Data Preview 0

In this section we discuss the following key topics:

- Dataset choice considerations
- Data products offered
- Services offered
- Audience considerations

3.1 Dataset choice considerations

The Construction Project has been working for some time now with a number of pre-cursor datasets and simulated data. There are two leading candidates for forming the basis of DP0:

- The Subaru Hyper Suprime-Cam PDR2 dataset, provided permission can be secured from our HSC colleagues. As real (on-sky) data it is likely that users will interact with it in more realistic ways. It is a well understood dataset, and it is regularly re-processed with software that shares a common codebase with the LSST Science Pipelines.
- The simulated precursor to LSST data produced by the Dark Energy Survey, DESC DC2, provided permission can be secured. This is a very large dataset and putting DC2 catalogs in Qserv would be an excellent demonstration of its abilities.

Data Management is currently in transition between its 2nd and 3rd generation data abstraction layer (aka “Butler”). For DP0 to fulfill its aim as an early deployment/integration exercise, Gen 3 Butler must be used, preferably (stretch goal) using an S3 compliant Object Store as is the intent in production. This has bearing on the choice of dataset. HSC PDR2 can either be converted from Gen 2 to Gen 3 or (stretch goal but ideally?) reprocessed naively with Gen3. Hence this is preferred choice from an engineering point of view.

DESC2 is available through Gen2 Butler and as we do not process that data with the Science Pipelines, the only option is conversion to Gen3 but estimates are that this is such a time-consuming process that it cannot be done in time for DC2. Therefore if DC2 is to be involved, a significantly smaller subset would have to be selected.

Questions:

- Which dataset has the broader scientific interest
- If DC2 and we take a subset to avoid the Gen2-Gen3 conversion issues, will that reduce the usefulness of picking DC2 in the first place?
- Will we be able to do Butler over S3 at production grade by DP0 ?

3.2 Data Products Offered

We will offer access to images and catalogs, though in more limited ways that will be available in Operations.

Images will be stored in read-only Butler Gen3 repo.

Catalogues will be stored in Qserv.

Questions:

- Are we offering parquet files?
- We should presumably explicitly rule out bulk download — YES
- When does ingest into Qserv has to start to be ready by DP0?

3.3 Services Offered

Although DP0 as a milestone described LSO-011 can be fulfilled with simple data distribution, we intend to offer limited Science Platform functionality as part of DP0. This includes:

- Provided the data is stored in Qserv or a Postgres database, catalogue access through TAP
- Access to the Science Platform's notebook-based analysis environment (Nublado); images can be accessed pragmatically via the Butler.
- Federated Authentication

Shell access (except through Nublado) will not be offered.

Questions:

- Is it understood that portal is not included? Not necessarily ..

3.4 Audience Considerations

Care should be taken to limit the target audience for the data previews; it is most critical that this is done for DP0.

- We have limited capacity to divert resources to support users.
- We will not have performed scaling tests on the Science Platform services by that point; current Science Platform usage is under 100 users, and any intent to exceed that should be communicated well in advance
- We will not yet have the ability to throttle excessive IDF usage

Authorization will be provided in an all-in basis (users will have the same level of access as project members currently have) since finer access control mechanisms will not be available by DP0; care should be taken in selecting them.

Questions:

- What is the authorization constraints for this data? For example, are DC2 data products only available to DESC science collaboration members? If so, if DC2 is chosen, does only DESC participate in DP0? **When agreed DC2 would be available to all data rights holders.**
- How do we handle access? First come first served? Do we need a sign-up process?
- How do we intend to do support? Slack? JIRA? CLO?

4 DP0.2 - processing

The Milestone DP-SP-02 includes re processing on IDF of the data set previously served as part of DP-SP-01. This requires a workflow system and associated tools to preferable make this quite automated. Demonstrating a portable set of cloud enabled tools based on Butler Gen 3 and HTCondor would help to allay the main risk of moving to a new Data Facility in operations.

5 Risks and mitigation

The biggest schedule risk is not getting an interim data facility in place in time. This would delay the entire schedule and there is not much mitigation.

In the long run costs may be higher than expected in a cloud based IDF. This will be due to storage. An mitigation to this would be to store data on our own systems (NCSA or Chile) and expose it through S3. NCSA already have this in place and we should consider testing this for lesser used data sets.

6 Planning and team(s) fro DP0

6.1 Teams

The main departments involved in this are Data Production and System Performance. With in those departments various people will be involved from the underlying teams but in small numbers. It makes most sense to approach DP0 with a task force approach. This might best be seen as two teams:

- Data production - with a focus on middleware and execution (Section 6.2);
- System Performance - with a focus on quality assurance and community support (Section 6.3).

6.2 DP Middleware and Execution

For DP0 on IDF Hsin-Fang Chiang would coordinate Data Production activities and be the point of contact for the IDF provider.

6.3 SP Quality and Community Support

Leanne ..

A References

References

[**LSO-011**], William O'Mullane, L.G., Phil Marshall, 2019, *Release Scenarios for LSST Data*, LSO-011, URL <https://lso-011.lsst.io>

B Acronyms

Draft