**Rubin Observatory**

Near term workflow for pre-operations with PanDA

**William O'Mullane and Richard Dubois and Hsin Fang Chiang**

2020-12-18

# 1   Introduction

We need a workflow system and production tools which can process DESC DC2 for DP0.2. Nominally the processing starts in June 2021.  We have a milestone L3-MW-0050 in March (see Table 1) for Batch system installation and configurations on IDF and L3-MW-0060 for the DP0.2 production.  The preferable way to do this would be with BPS in front of PanDA but there are potentially other solutions (see Section 5).

For PREOPS we should focus on PanDA in the near term get all the hooks in place and make it work for DP0.2.  Getting this in place requires some leadership and decision making.  We need a product owner and manager (see Section 3).  This is separate from the construction side's HSC reprocessing at NCSA for development needs. The constructions team at NCSA can continue to use Condor-based BPS for the biweekly HSC reprocessing at NCSA.

How effective these tools are will determine how effort-intensive (and successful) the large-scale processing campaigns will be.

## 1.1   Milestones

General DP0 information is RTN-001.  For simplicity some milestones are copied here in Table 1. Jira is the source of truth for dates on these though some may need revising.

Table 1: FY21 Middleware Milestones

| Milestone | Jira ID | Rubin ID | Due Date | Level | Team |
|---|---|---|---|---|---|
| Read only Gen3 butler for DP0 at IDF | PREOPS-143 | L3-MW-0030 | 2021-03-31 | 3 | Science Users Middleware |
| Qserv installation on IDF | PREOPS-142 | L3-MW-0010 | 2021-03-31 | 3 | Science Users Middleware |
| PanDA based workflow system in place | PREOPS-154 | L3-MW-0050 | 2021-03-31 | 3 | Science Users Middleware |
| DP0.1 data loaded into Qserv on IDF | PREOPS-144 | L3-MW-0020 | 2021-04-30 | 3 | Science Users Middleware |
| Gen3 butler and pipeline task ready for DP0 production | PREOPS-156 | L3-MW-0070 | 2021-06-10 | 3 | Science Users Middleware |
| PanDA based workflow system with tooling (e.g. restart) added. | PREOPS-155 | L3-MW-0060 | 2021-06-30 | 3 | Science Users Middleware |
| Evaluate Batch Production System | PREOPS-153 | L3-MW-0040 | 2021-07-31 | 3 | Science Users Middleware |

**Rubin** Observatory

## 2   Requirements and priories

LDM-636 forms the formal requirements baseline.

Concisely we need the execution team to be able to run DP0.2 with minimum hand holding. Hence the top priorities for the near term would be:

1. Documentation: preferably on lsst.io, enough for the execution team to kick off pipelines, monitor and to first order troubleshoot them.

2. Workflow monitoring - some sort of web page which gives status (perhaps slightly customized)

3. Restart: Can resume an unfinished workflow. Can automatically retry jobs killed by preemption, DB connection, or other transient issues.

4. Logstash: On IDF this will be Google Logging. Any logging should end up in the same central logging system.

5. Troubleshooting failed jobs: Features to help understand non-transient failures, such as error messages aggregation and ways to reproduce failures. This kind of error usually is caused by pipeline failures and needs follow-up investigation.

Longer term (which may not be for DP0.2)we need

- Installation at SLAC

- Multi site execution with France and eventually UK.

- Campaign execution monitoring

### 2.1   Timeline

We have a milestone L3-MW-0050 in March (see Table 1) for Batch system installation and configurations on IDF and L3-MW-0060 for the DP0.2 production. We should track these two milestones: L3-MW-0050 for an initial system and L3-MW-0060 to have the system to run DP0.2.

**Rubin Observatory**

## 2.2 Evaluation

L3-MW-0060 will see the commencement of the processing run - we assume there may be some hiccups at that point. But at L3-MW-0060 + one month we should decide if this is the long term approach for Rubin Operations with DOE buy in. Hence L3-MW-0040 is approximately the evaluation date.

## 3 Team

SLAC obviously have long term interest in this working and on a single track so it would be good to have some SLAC oversight on the topic. A product owner to shepherd requirements and priorities as well as a manager to guide resources must be identified. Currently (all at partial fractions) the team consists of:

- Brian Yanny and team at FERMILAB for execution

- Monica Adamow - Execution NCSA

- Michelle Gower, Mikolaj Kowalik - BPS and deployment

- Sergey Padowski and Shuwei Ye (starting in January)- PanDA

## 4 PanDA

The PanDA ("Processing and Data Analysis") system was created by ATLAS at LHC to manage its massive processing efforts. In that capacity it handles several hundred thousand processing jobs per day across heterogeneous systems, supporting multiple parallel campaigns. Its main services (PanDA, Harvester, iDDS) are driven from a central database. The system can ingest DAGs, handle the workflow and then the workload management. Currently PanDA cannot rerun parts of workflow, but the feature is being actively considered for addition.

PanDA satisfies a number of criteria:

- Multi-site authentication

**Rubin** Observatory

- Multi-site processing - Harvester can be used to mitigate network traffic between sites and central workflow db; also handles site-specific submission properties allowing a range of different kinds of resources

- Manages workflow (via iDDS) as well as workload

- Good monitoring tools for the submitted workflow. Can be customized.

While support would be dependent on BNL expertise, several installations of PanDA have been undertaken outside of ATLAS, so there is experience in doing installs and of ongoing maintenance for other organizations.

In order to demonstrate the viability and customizability of PanDA for Rubin, BNL has set a target of doing processing with PanDA in the IDF by the March 2021 time frame. As a part of that demonstration, they will provide documentation of the PanDA system.

It would be additionally instructive to set up multi-site processing to include the French Data Facility and US Data Facility during 2021.

However, campaign management is outside PanDA's scope, so a layer on ctrl_bps would be needed to chunk up and keep track of elements of a campaign. Ctrl_bps would likely also need to handle resubmissions.

## 5   Potential solutions

Conceptually this is done in two steps: (a) workflow generation and (b) job execution. In step (a) the workflow generation defines executable jobs and job interdependency as a graph. In step (b) job execution includes workflow status monitoring, pausing/resuming/killing work-flows, debugging/retrying failed jobs, resource usage monitoring, and relevant toolkits to facilitate execution management on a large scale.

1. ctrl_bps workflow generation + PanDA-plugin execution tools developed by BNL

2. ctrl_bps workflow generation + Condor-plugin execution tools developed by NCSA

3. ctrl_bps workflow generation + Pegasus as the execution tools

**Rubin** Observatory

4. ctrl_bps workflow generation tools can't work on IDF, one can use customized scripts to generate workflow for any execution tools.

# 6 Risks and worries

1. Lack of documentation for PanDA: it is a complex system and will be the heart of processing.. operating for 12 years in this mode is unwise.

2. Dependence on an institution or individual because of 1, also suggests the need to spread the expertise more broadly across the team.

3. Having a LOT of scripting to make a production run of any size

4. dependence on Oracle: is an open source project would you not like to depend on commercial products furthermore some of us have had bad experience with Oracle.

# A   References

**[LDM-636]**, Kowalik, M., Gower, M., Kooper, R., 2019,  *Batch Production Service Requirements*, LDM-636, URL `https://ls.st/LDM-636`

**[RTN-001]**, O'Mullane, W., 2020,  *Data Preview 0: Definition and planning.*,  RTN-001, URL `http://RTN-001.lsst.io`

# B   Acronyms

| Acronym | Description |
| --- | --- |
| ATLAS | A Toroidal LHC Apparatus |
| BNL | Brookhaven National Laboratory |
| BPS | Batch Production Service |
| DB | DataBase |
| DC2 | Data Challenge 2 (DESC) |

**Rubin Observatory**

| | |
|---|---|
| DESC | Dark Energy Science Collaboration |
| DM | Data Management |
| DOE | Department of Energy |
| DP0 | Data Preview 0 |
| FY21 | Financial Year 21 |
| HSC | Hyper Suprime-Cam |
| IDF | Interim Data Facility |
| L3 | Lens 3 |
| LDM | LSST Data Management (Document Handle) |
| LHC | Large Hadron Collider (at CERN) |
| NCSA | National Center for Supercomputing Applications |
| PanDA | Production ANd Distributed Analysis system |
| RTN | Rubin Technical Note |
| SLAC | SLAC National Accelerator Laboratory |
| UK | United Kingdom |
| US | United States |