



Vera C. Rubin Observatory
Rubin Observatory Operations

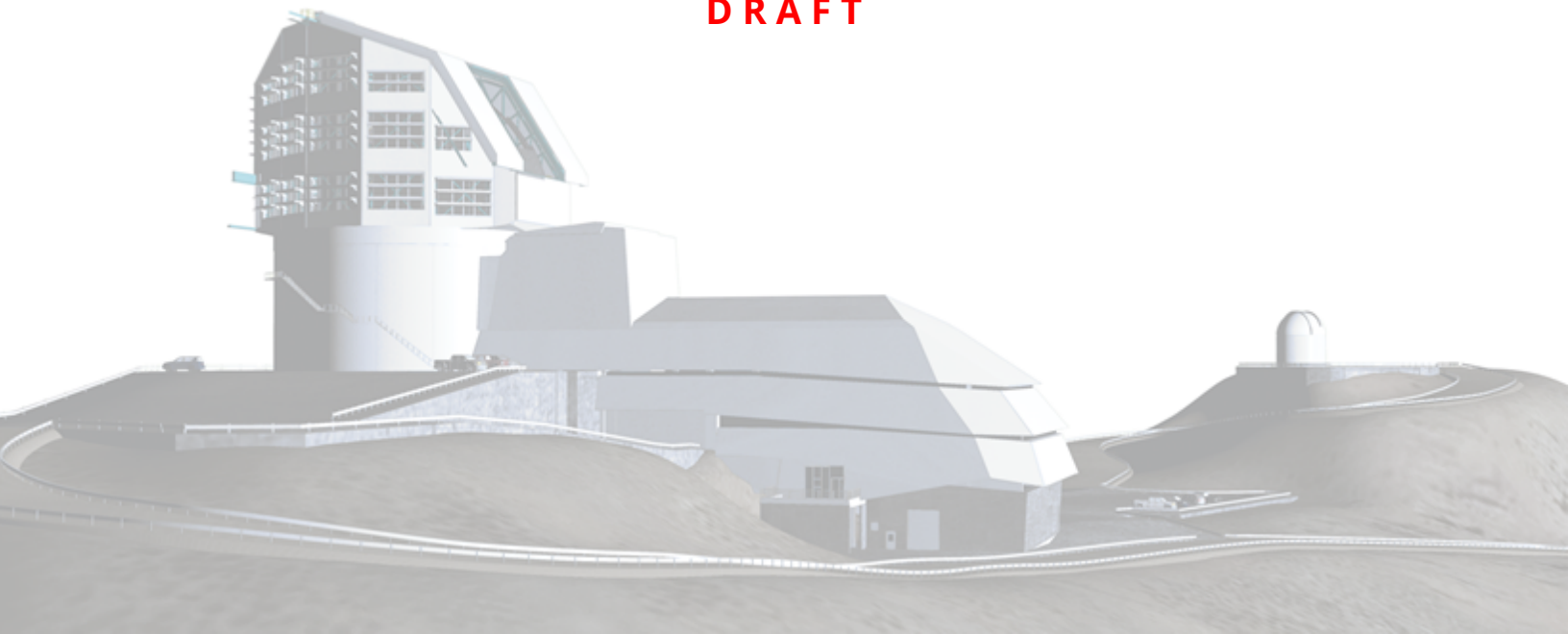
Management and Execution plan for Data Management Operations.

William O'Mullane

RTN-046

Latest Revision:

DRAFT



Abstract

This is the management plan for operations of Data Management - this includes software products and data products.

Draft

Change Record

Version	Date	Description	Owner name
1	YYYY-MM-DD	Unreleased.	William O'Mullane

Document source location: <https://github.com/lstt/rtn-046>

Draft

Contents

1 Introduction	1
1.1 Purpose	1
1.2 Mission Statement	1
1.3 Goals and Objectives	1
2 Architecture	2
3 Functionality based teams and organisation	2
3.1 Data services	4
3.2 Making Data	5
3.3 Data Abstraction	5
3.3.1 Data Engineering	5
3.3.2 Middleware Assumptions	6
3.4 Infrastructure	6
4 Chile DevOps	6
5 Data facilities and access centers	6
5.1 USDF	6
5.2 FrDF	6
5.3 UKDF	7
5.4 US DAC	7
5.5 Chile DAC	7
6 Products	7
7 Roles	8
A References	8
B Acronyms	10

Management and Execution plan for Data Management Operations.

1 Introduction

1.1 Purpose

This document defines the mission, goals and objectives, organization and responsibilities of Vera C. Rubin Observatory Data Management Operations.

1.2 Mission Statement

Maintain, improve and operate a suite of Vera C. Rubin data management services to produce and serve to the community high-quality data products from the Legacy Survey of Space and Time.

1.3 Goals and Objectives

These are similar to our construction goals outlined in LDM-294. Rubin Data Management Operations will:

- Produce the data products as outlined in LSE-61
- Maintain and improve data production mechanisms.
- Maintain and improve data access mechanisms.
- Maintain and improve data abstraction mechanisms.
- Assess current and operations-era technologies for use in providing engineered solutions for Vera C. Rubin Observatory .
- Maintain appropriate cybersecurity measures throughout Vera C. Rubin Observatory and especially on external facing services.
- Document the operational procedures associated with using and maintaining DM capabilities.
- Evaluate, select, recruit, hire/contract and direct permanent staff, contract, and in-kind resources in Rubin and from partner organizations participating in DM initiatives.

The goals in selecting and, where necessary, developing Rubin software solutions are:

- We prefer to acquire and configure existing, off-the-shelf, solutions. Where no satisfactory off-the-shelf solutions are available, we develop the software and hardware systems necessary to meet our objectives. This extends into maintenance where we will continue to probe choices and may replace custom systems with off-the-shelf solutions where appropriate.
- The software architecture is actively managed at the subsystem level. A well engineered and cleanly designed codebase is less buggy, more maintainable, and makes developers who work on it more productive. We continue to follow and maintain the developer guide¹.
- Other than when prohibited by licensing, security, or other similar considerations, all newly developed source code, and in particular that pertaining to scientific algorithms, is public. Our primary goals in publicizing the code are to simplify reproducibility of LSST data products and to provide insight into algorithms used. Achieving these goals requires that the software must be properly documented.
- Background decision material on choices made will be documented in technical notes with "DMTN", "RTN" or similar series handles. (see `lsst.io`)

2 Architecture

The construction era DM architecture is defined in LDM-148.

As stated in the introduction our operational goals now include production of the data products. In broad terms we may think of two prongs in data management: Data Production and Data Serving. This is depicted in Figure 3.

We also now have three operational data facilities for data release production and a Cloud Facility on Google for the science users. This is all depicted in Figure 1.

Details about the build up the data facilities is given in RTN-021.

3 Functionality based teams and organisation

While Figure 2 Shows the reporting structure Figure 3 puts this more in a operations concept. We consider the main functions to be Data Production and Data Serving,

¹`developer.lsst.io`

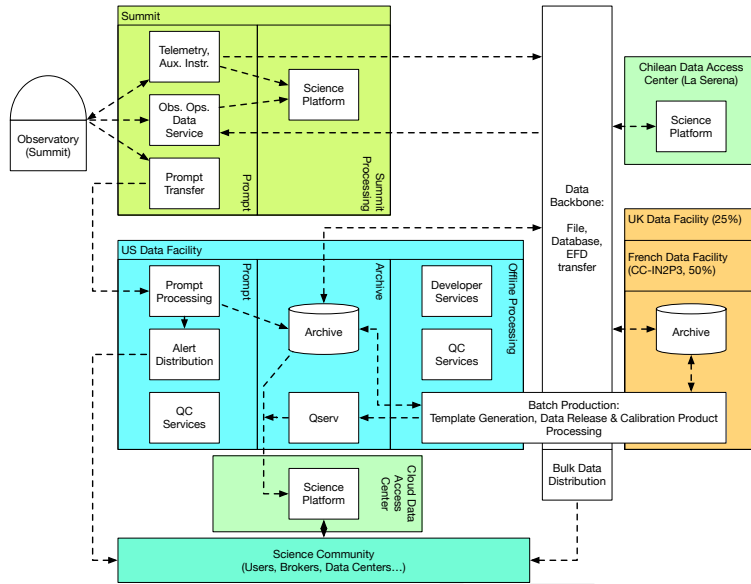


FIGURE 1: Simplified operations architecture for Data Management.

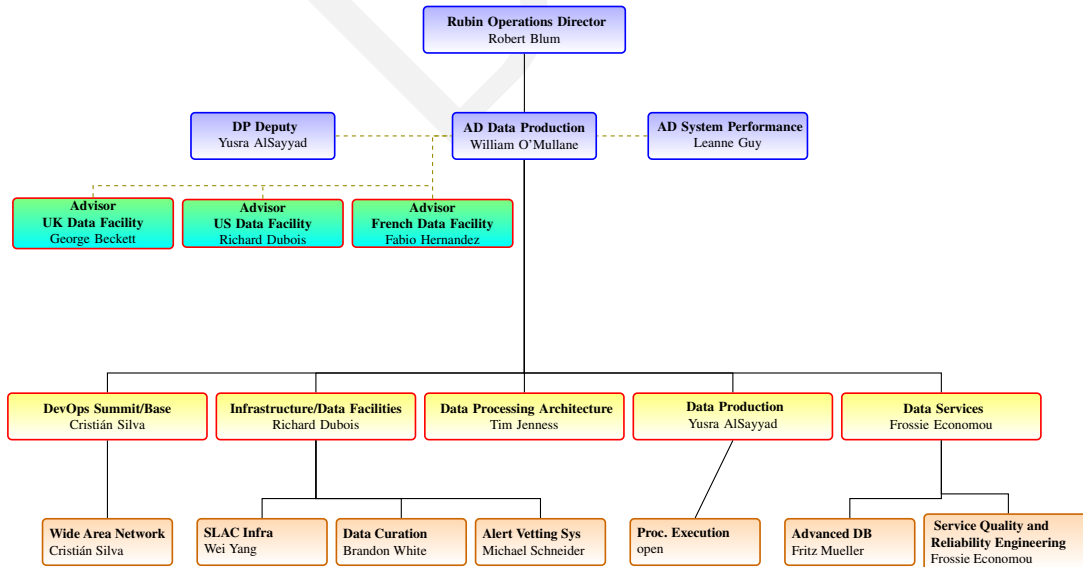


FIGURE 2: Reporting lines in Data Management Operations.

these are supported by the data abstraction team and the data facilities.

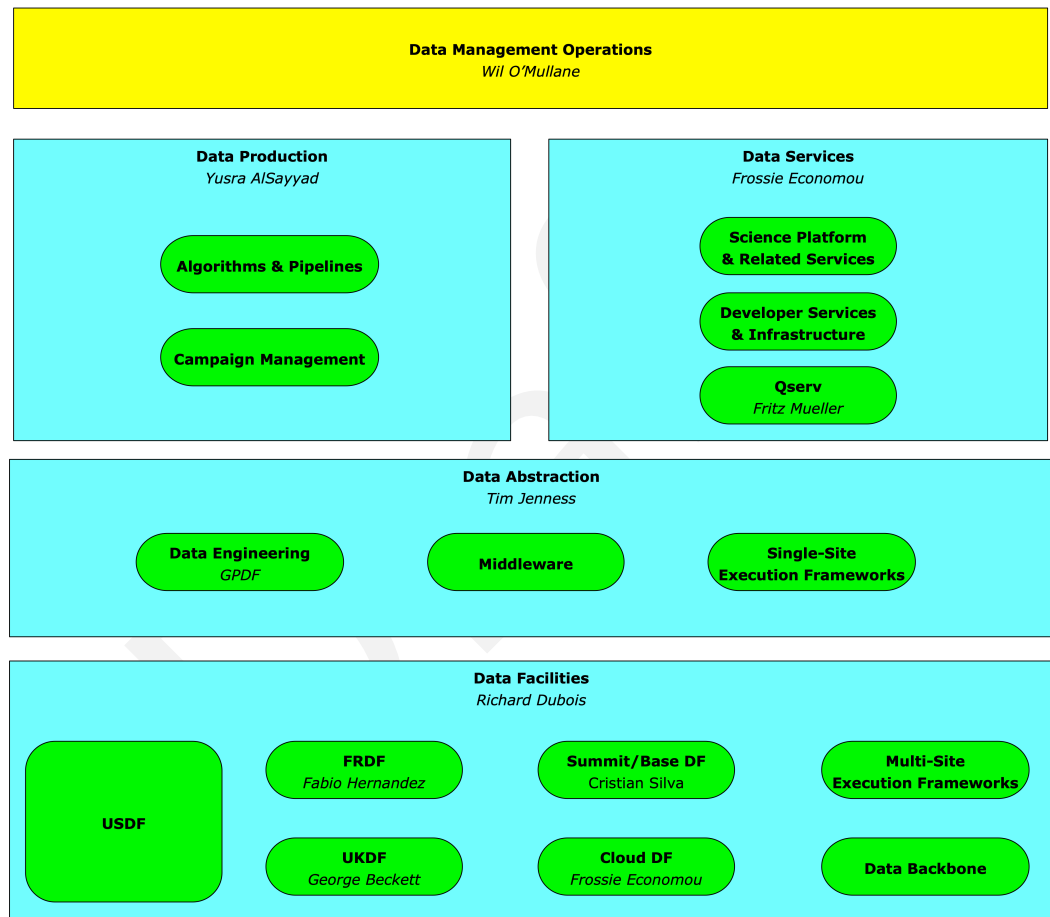


FIGURE 3: Functions in operations of Rubin Data Management.

3.1 Data services

All services associated with data serving are in this group. As depicted in Figure 3 this includes:

- The Science Platform
- Developer Services Infrastructure
- Qserv advanced Database
- The Engineering Facilities Database

A more complete list of items under may be found in the Section 6

3.2 Making Data

All services associated with data making are in this group. As depicted in Figure 3 this includes:

- The Science Pipelines code
- Execution of science pipelines to produce data products
 - Alert production
 - Data Release Production

A more complete list of items under may be found in the Section 6

3.3 Data Abstraction

Underpinning Data Making and Data Serving is out abstraction of data and services. This includes middleware such as butler and batch production systems etc. But also Prompt Processing execution and Data engineering. It is crucial for our system portability to maintain the abstraction layer.

A comprehensive list is given in Section 6.

Some of these require a little more discussion here.

3.3.1 Data Engineering

- Support the metadata translation infrastructure (astro_metadata_translator) and monitor correctness of FITS headers.
- Advise on file formats and file metadata for all systems writing files that are to be archived.
- Support the Felis system for specifying schemas.
- Define the data models for everything in the consolidated database. (“global data model” schema?)
- Write and support code that populates the consolidated database (for example, code that analyzes the EFD and creates the “exposure” and “visit” tables).
- Gregory Dubois-Felsmann is product owner (“Data Scientist”)
- Staffing: tiny in construction. 0.5 FTE in ops + fractional GPDF.

3.3.2 Middleware Assumptions

As can be seen in the product list there are a lot of elements to Middleware. A few assumptions are made.

- Assumes butler/Rucio integration is entirely handled by the infrastructure team.
- Assumes database administration is done by the infrastructure team.
- Staffing: 0.5 Andy S; 0.5 new NateP; NateL 0.25; Matthias 0.5; MichelleG 0.5; Kowalik 0.25; at least 0.25 PanDA ongoing person for ctrl_bps_panda.

3.4 Infrastructure

text here

4 Chile DevOps

The Chile DevOps team delivers and maintains the base and summit facilities.

The team provides networking and machines plus a kubernetes layer ready for deployment of services from data management as well as telescope and site software.

On the summit there are also some bare metal machines which are setup usually with puppet.

5 Data facilities and access centers

Hardware underpins all of our operations. This is arranged in three data facilities in US, UK and France as outlined below. We also have two on project Data Access Centers to provide services to the scientific users.

The plan for building up the data facilities is in RTN-021. A more complete list of items under may be found in the Section 6

5.1 USDF

The USDF will be the main archive of Rubin data. It performs the daily processing of data including alert generation. It performs 25% of th DRP processing. There is a full description in DMTN-189. User batch will run at the USDF [DMTN-223].

5.2 FrDF

The French Data facility will hold a copy of the Raw data. The FRDF will run 50% of the DRP processing.

5.3 UKDF

The FRDF will run 25% of the DRP processing.

5.4 US DAC

The USDAC is hosted on Google Cloud. Most image data remains at USDF but some catalogs and possibly coadds will be kept on Google. All User files spaces and the RSP will be on google [DMTN-209].

5.5 Chile DAC

The Chilean Data Access Center will be built after operations commences. Some discussions are still pending on its exact shape see LDM-572.

6 Products

Product	Manager	Owner	Notes
DM Ops	Wil O'Mullane		Data Management (Ops)
Data Abstraction	Tim Jenness		Data Abstraction
Data Engineering	Gregory Dubois Felsmann		Data Engineering
Felis	TBD		Felis
Metadata	TBD		Metadata
Middleware	Tim Jenness		Middleware
BPS	TBD		BPS
Butler	Tim Jenness	Y (Jlm B)	Butler
Control Interface	TBD		Control Interface
ctrl_bps	TBD		ctrl_bps
ctrl_mpexec	TBD		ctrl_mpexec
user batch envelope	TBD		user batch envelope
Pipeline interfaces	TBD		Pipeline interfaces
pex_config	TBD		pex_config
pipe_base	TBD		pipe_base
Single-site Exec	TBD		Single-site Exec
OCPS	KT Lim ?		Observatory Controlled Processing System
Prompt f/ w	TBD		Prompt forwarder
Data Facilities	Richard Dubois		Data Facilities
Data Curation	TBD		Data Curation
Data Backbone	TBD		Data Backbone
Backups	TBD		Backups
Bulk Download	TBD		Bulk Download
Consolidated DB	TBD		Consolidated DB
Butler repos	TBD		Butler repos
Rucio	TBD		Rucio
Infrastructures	TBD		Infrastructures
CDF	Frossie Economou		CDF

FrDF	Fabio Hernandez		FrDF
TDF	Cristián Silva		TDF
UKDF	George Beckett		UKDF
USDF	TBD		USDF
Multi-site & User Exec	TBD		Multi-site & User Exec
PanDA	TBD		PanDA
User Batch	TBD		User Batch
Data Production	Yusra AlSaiyyad		Data Production
Campaign Management	?	Y (N/ A)	Campaign Management
Algorithms & Pipelines	Yusra AlSaiyyad	? (Jim Bosch)	Algorithms & Pipelines: In ops our construction POs (JimB+EricB) become our group leads, so PO prob not necessary.
Data Services	Frossie Economou		Data Services
Complex. DB	Fritz Mueller	Y (Colin)	Complex Database Support: Should this be complex Databases or somethign ? Qserv under
Big Databases	Fritz Mueller		Big Databases
PromptDV	TBD		Prompt Products DB
Qserv	Fritz Mueller		Qserv
User Databases	TBD		User Databases
RSP Portal	TBD		RSP Portal
SQuaRE	TBD		SQuaRE
Doc Services	TBD		Doc Services
Documentation standards	TBD		Documentation standards
LtD	TBD		LtD
Templating	TBD		Templating
Phalanx	TBD		Phalanx
Authorisation	TBD		Authorisation
Reliability Engineering	TBD		Reliability Engineering
Secrets	TBD		Secrets
RSP	TBD		RSP
APIs		Y (GPDF)	APIs: IVOA and non-VO Apis
data.lsst.cloud	TBD		data.lsst.cloud
Authentication			Authentication: and security engineering
Notebook		Y (KSK)	Notebook
Portal		Y (GPDF)	Portal
Square One	TBD		Square One
User Support			User Support: clo service and helpdesk
Sasquatch	TBD		Sasquatch
EFD	TBD		EFD
Metrics	TBD		Metrics
Telemetry Gateway	TBD		Telemetry Gateway

7 Roles

A References

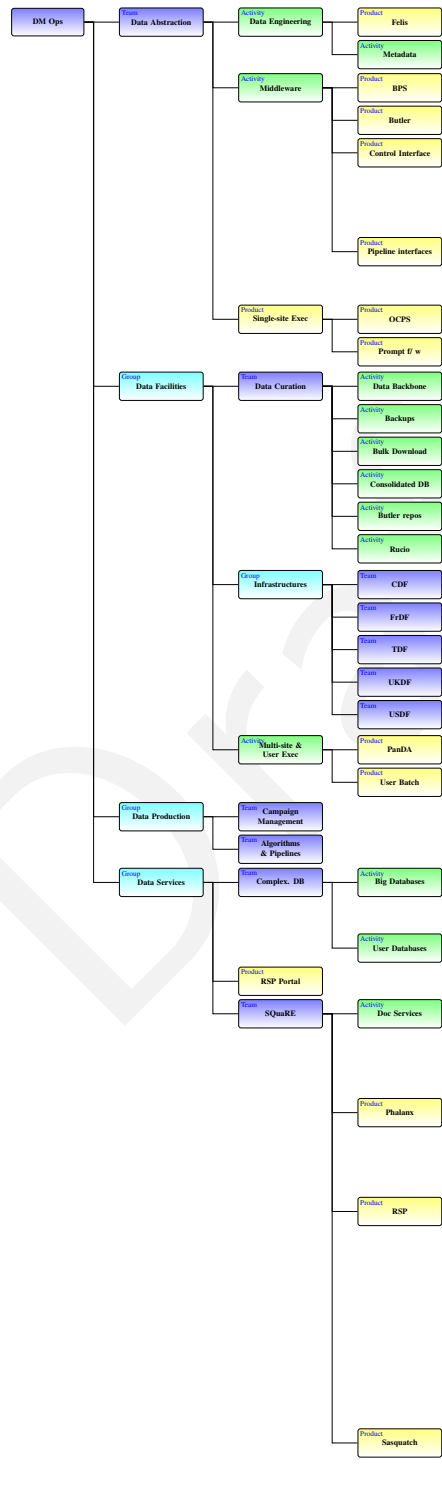


FIGURE 4: DM operations product tree

- [RTN-021]**, Dubois, R., O'Mullane, W., 2022, *Data Facilities Transition Plan*, RTN-021, URL <https://rtn-021.lsst.io/>,
Vera C. Rubin Observatory Technical Note
- [LSE-61]**, Dubois-Felsmann, G., Jenness, T., 2019, *Data Management System (DMS) Requirements*, LSE-61, URL <https://lse-61.lsst.io/>,
Vera C. Rubin Observatory
- [DMTN-189]**, Lim, K.T., 2021, *Data Facility Specifications*, DMTN-189, URL <https://dmtn-189.lsst.io/>,
Vera C. Rubin Observatory Data Management Technical Note
- [LDM-148]**, Lim, K.T., Bosch, J., Dubois-Felsmann, G., et al., 2018, *Data Management System Design*, LDM-148, URL <https://ls.st/LDM-148>
- [LDM-572]**, O'Mullane, W., 2021, *Chilean Data Access Center*, LDM-572, URL <https://ldm-572.lsst.io/>,
Vera C. Rubin Observatory Data Management Controlled Document
- [DMTN-223]**, O'Mullane, W., 2022, *User batch - possibilities and plans.*, DMTN-223, URL <https://dmtn-223.lsst.io/>,
Vera C. Rubin Observatory Data Management Technical Note
- [DMTN-209]**, O'Mullane, W., Economou, F., Huang, F., et al., 2021, *Rubin Science Platform on Google: the story so far.*, DMTN-209, URL <https://dmtn-209.lsst.io/>,
Vera C. Rubin Observatory Data Management Technical Note
- [LDM-294]**, O'Mullane, W., Swinbank, J., Juric, M., Guy, L., DMLT, 2022, *Data Management Organization and Management*, LDM-294, URL <https://ldm-294.lsst.io/>,
Vera C. Rubin Observatory Data Management Controlled Document

B Acronyms

Acronym	Description
B	Byte (8 bit)
BPS	Batch Production Service
CDF	Cumulative Distribution Function
DAC	Data Access Center

DB	DataBase
DF	Data Facility
DM	Data Management
DMTN	DM Technical Note
DP	Data Production
DRP	Data Release Production
EFD	Engineering and Facility Database
FITS	Flexible Image Transport System
FTE	Full-Time Equivalent
FrDF	French Data Facility
IVOA	International Virtual-Observatory Alliance
LDM	LSST Data Management (Document Handle)
LSE	LSST Systems Engineering (Document Handle)
LSST	Legacy Survey of Space and Time (formerly Large Synoptic Survey Telescope)
OCPS	OCS Controlled Pipeline System
OPS	Operations
PO	Program Operations
PanDA	Production ANd Distributed Analysis system
RSP	Rubin Science Platform
RTN	Rubin Technical Note
SQuaRE	Science Quality and Reliability Engineering
TBD	To Be Defined (Determined)
UK	United Kingdom
UKDF	United Kingdom Data Facility
US	United States
USDF	United States Data Facility
VO	Virtual Observatory
VRO	(not to be used)Vera C. Rubin Observatory
bps	bit(s) per second