



U.S. National
Science Foundation



Office of Science

Rubin US Data Facility Overview and Update

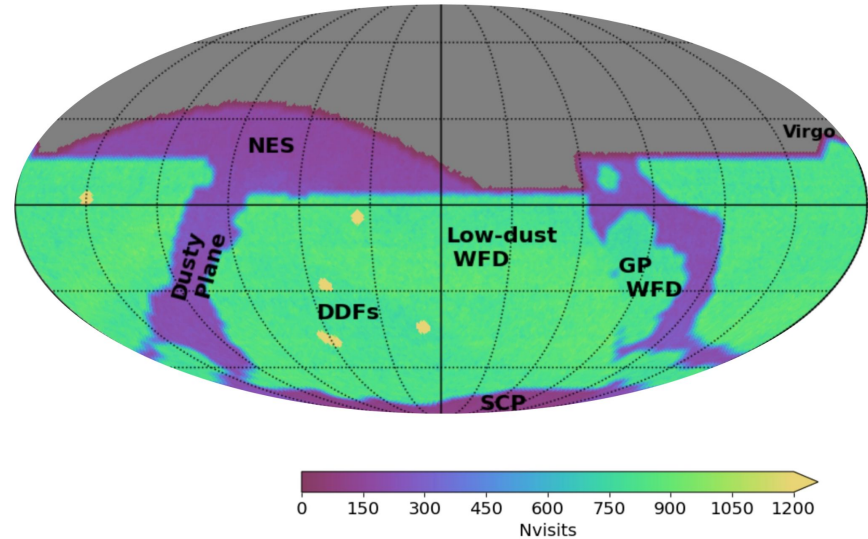
Adam S. Bolton, SLAC National Lab

w/ K-T Lim, Wei Yang, Fabio Hernandez

SA3CC Meeting, 07 May 2025



Rubin Observatory / Legacy Survey of Space and Time

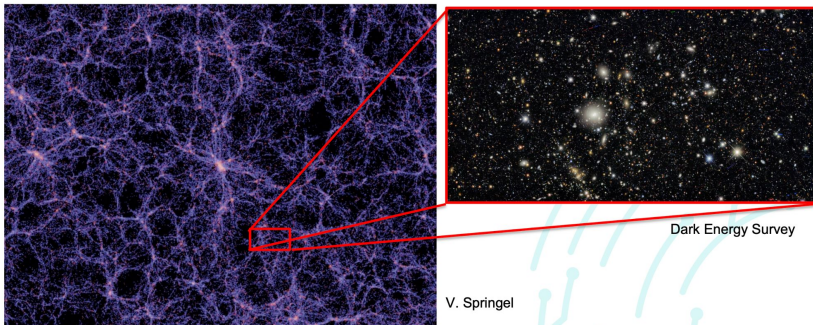


Imaging the entire visible sky every 3-4 nights for 10 years with the world's largest digital camera coupled to the world's highest-extended telephoto "lens" at an excellent astronomical site in the Chilean Andes to build a time-resolved 6-color database of 20 billion stars and 20 billion galaxies which we will release to scientists worldwide

Primary science drivers



Dark matter and dark energy



Dark Energy Survey

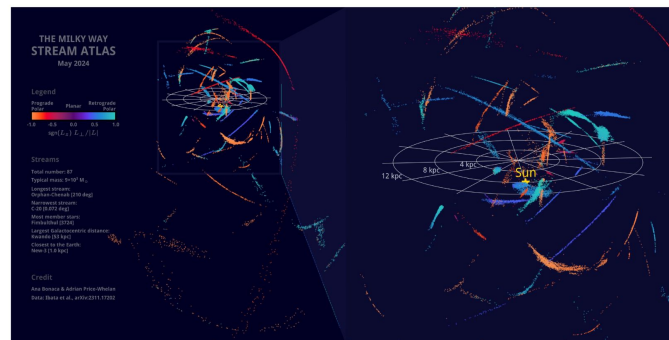
V. Springel

A. Bolton | Monterey Data Conference | 20 August 2024

3



Structure & formation of our Milky Way



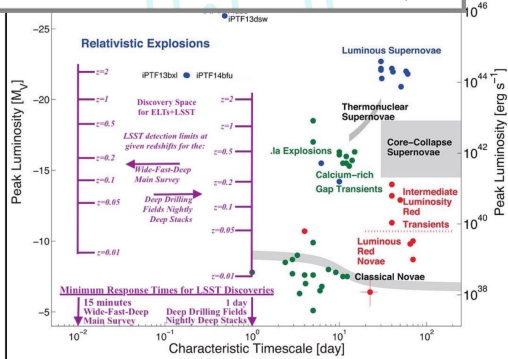
Bonaca & Price-Whelan arXiv:2405.19410

A. Bolton | Monterey Data Conference | 20 August 2024

6

Exploring the transient Universe

M. Graham et al.
arXiv:1904.05957

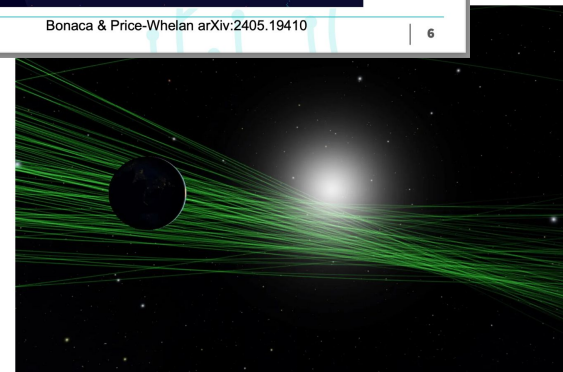


A. Bolton | Monterey Data Conference | 20 August 2024

4

Cataloging our Solar System

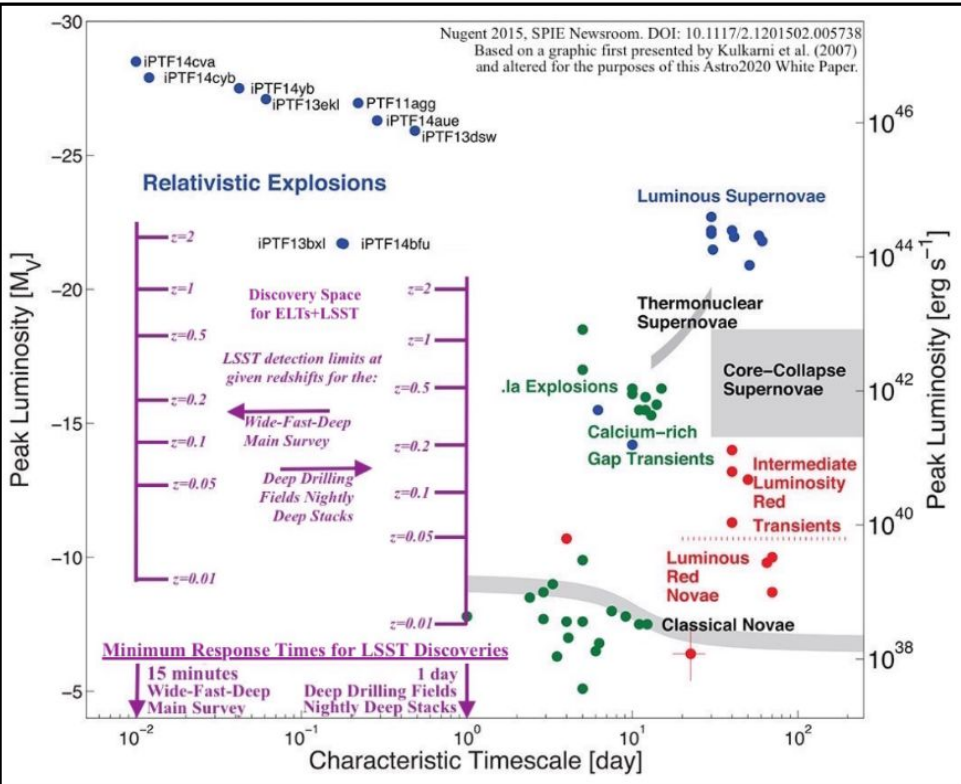
B612 Asteroid Institute /
U. of Washington DIRAC Institute /
OpenSpace Project



A. Bolton | Monterey Data Conference | 20 August 2024

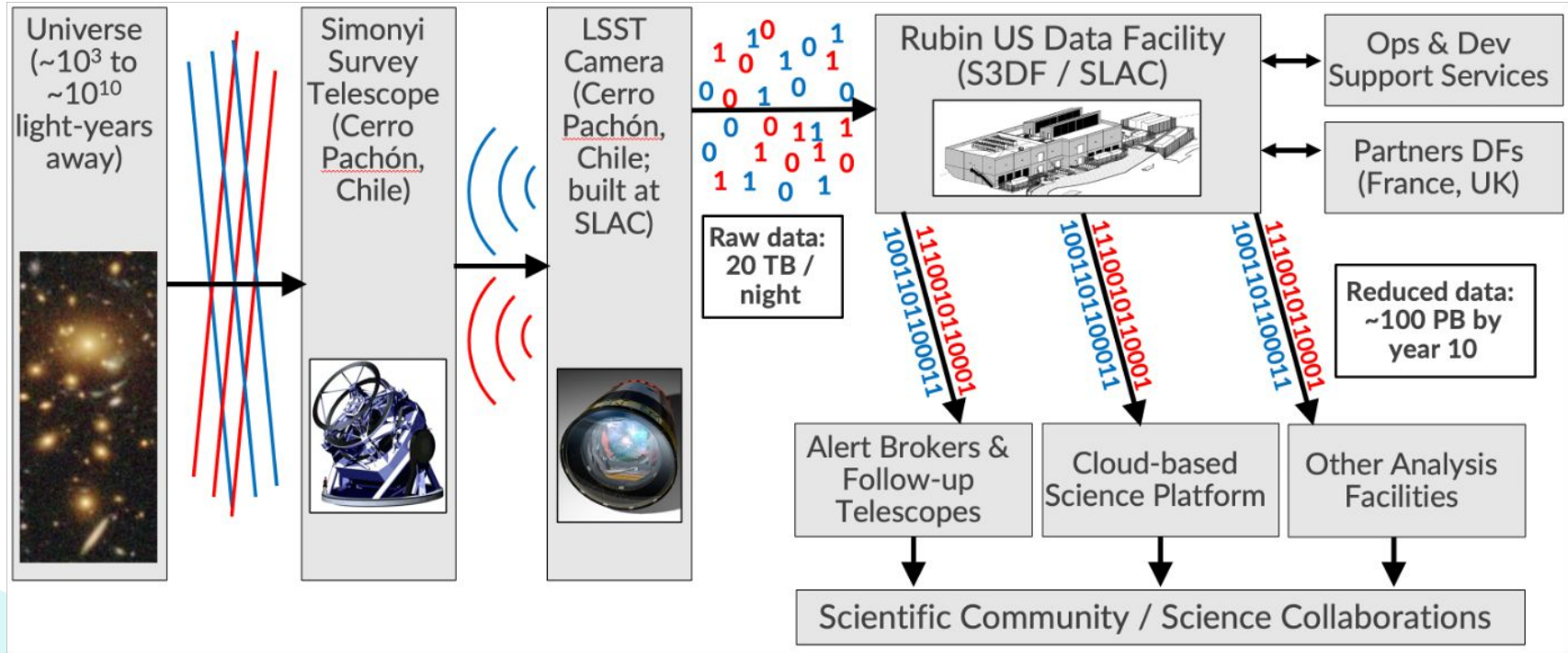
5

Exploring the transient Universe in real time



M. Graham et al. arXiv:1904.05957

Telescope + Camera + Data System





Sites & Data Flow



US Data Facility
SLAC National Accelerator Laboratory
Menlo Park, CA Processing Center
Processing Center
Alert Production
Data Release Production (25%)
Calibration Products Production
Long-term Storage (copy 1)

HQ Site
Tucson, AZ
Science Operations
Observatory Management
Education & Public Outreach

Base Site
La Serena, Chile
Base Center
Data Access Center
Data Access & User Services

UK Data Facility
IRIS Network, UK
Data Release Production (25%)

France Data Facility
CC-IN2P3, Lyon, France
Data Release Production (50%)
Long-term Storage (copy 2)

US
Data Access Center
Data Access and User Services
EPO Infrastructure

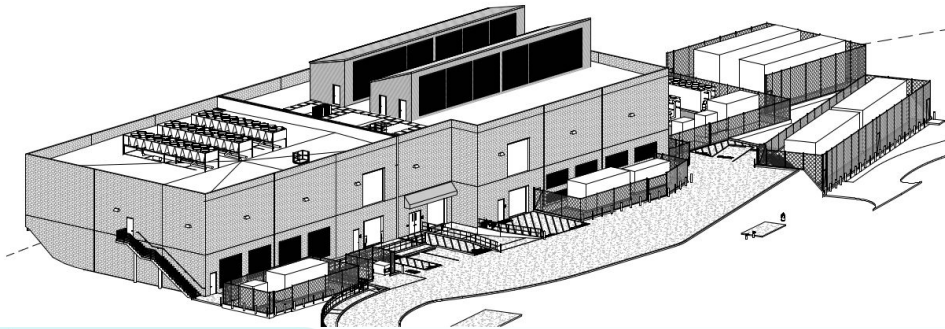
Summit Site
Cerro Pachón, Chile
Telescope & Camera
Data Acquisition
Crosstalk Correction

USDF roles within Rubin Observatory

- **Stewardship and curation** of Rubin Observatory's raw data and generated data products, including backup and disaster recovery systems
- **Computing infrastructure and environment** for prompt processing, alert generation, and data release processing, including framework for distributed multi-site processing
- **Movement and mirroring of data** between USDF and other processing and analysis sites, including FrDF, UKDF, and IDACs
- **Platform to support diverse set of Rubin scientific applications, services, and databases**, many of which are developed and operated by other Rubin teams
- **Operation of the hybrid US DAC**, in collaboration with SQuaRE, Qserv, Middleware, others
- **Coordination with the Summit facility** to ensure reliable and performant data transfer
- **Implementation of USDF aspects of Rubin Data Security policy** (DMTN-199)
- **Liaison and support for Data Facility stakeholders** within Rubin and broader community

SLAC Shared Science Data Facility

- SLAC's new shared high-throughput experimental computing infrastructure
 - Consolidating historically siloed scientific computing at SLAC
- Hosted within Stanford Research Computing Facility (SRCF)
 - Split across two adjoining facilities (SRCF-I and SRCF-II)
- Key mission area: critical, data-heavy, scientific computing workflows
 - Supports other large experiments in addition to Rubin (LCLS, UED, CryoEM, SSRL)
- Operated by Scientific Computing (SCS) Division of SLAC's Technology and Innovation Directorate (TID)



- 6MW Facility with air cooling (SLAC can use up to 2.5 MW)
- Flywheel + Generator allows for resilient power
- SLAC has over 100 racks
- 400 Gbps Networking to SLAC backbone

Slide content courtesy Jay Srinivasan

Rubin data flow and interfaces

Colors encode different levels of data security

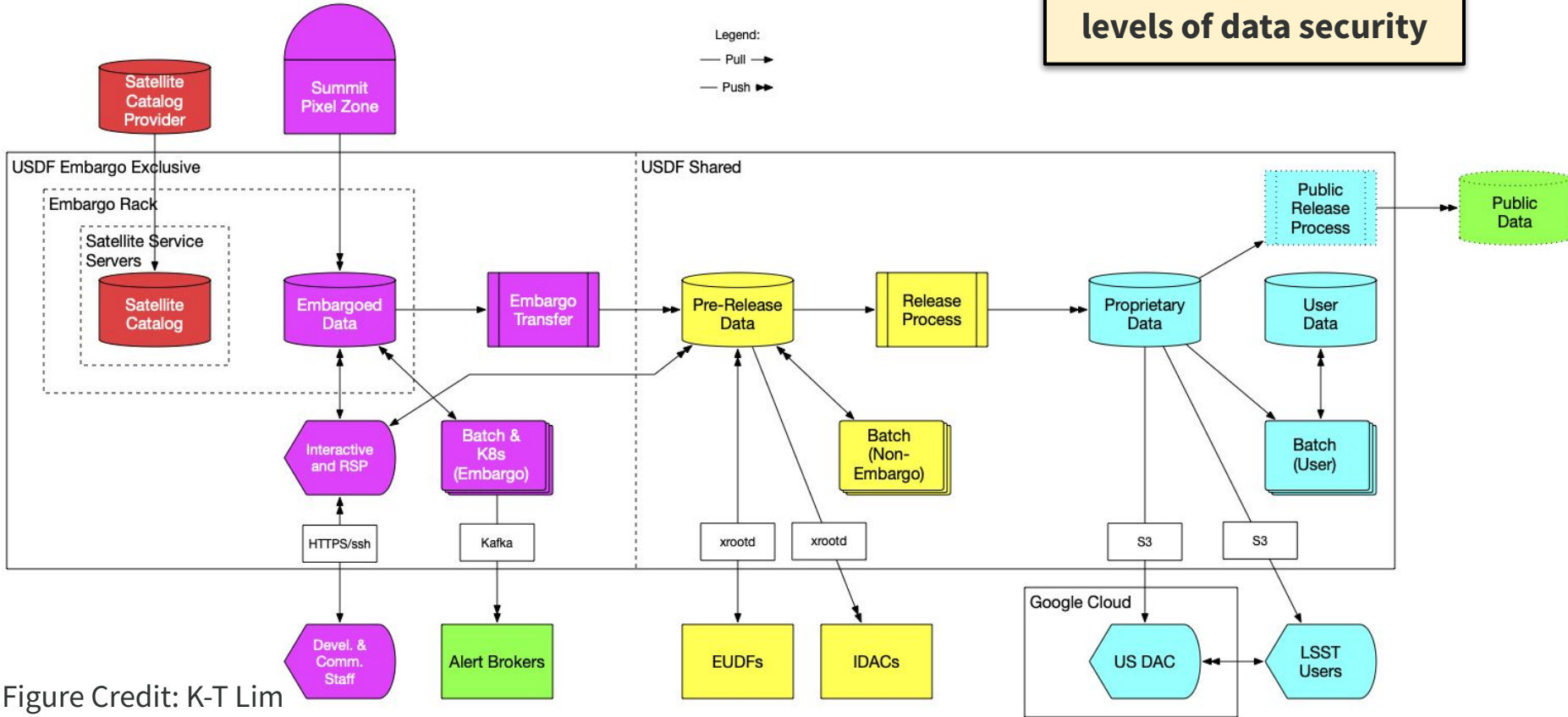


Figure Credit: K-T Lim

Time-domain alert science

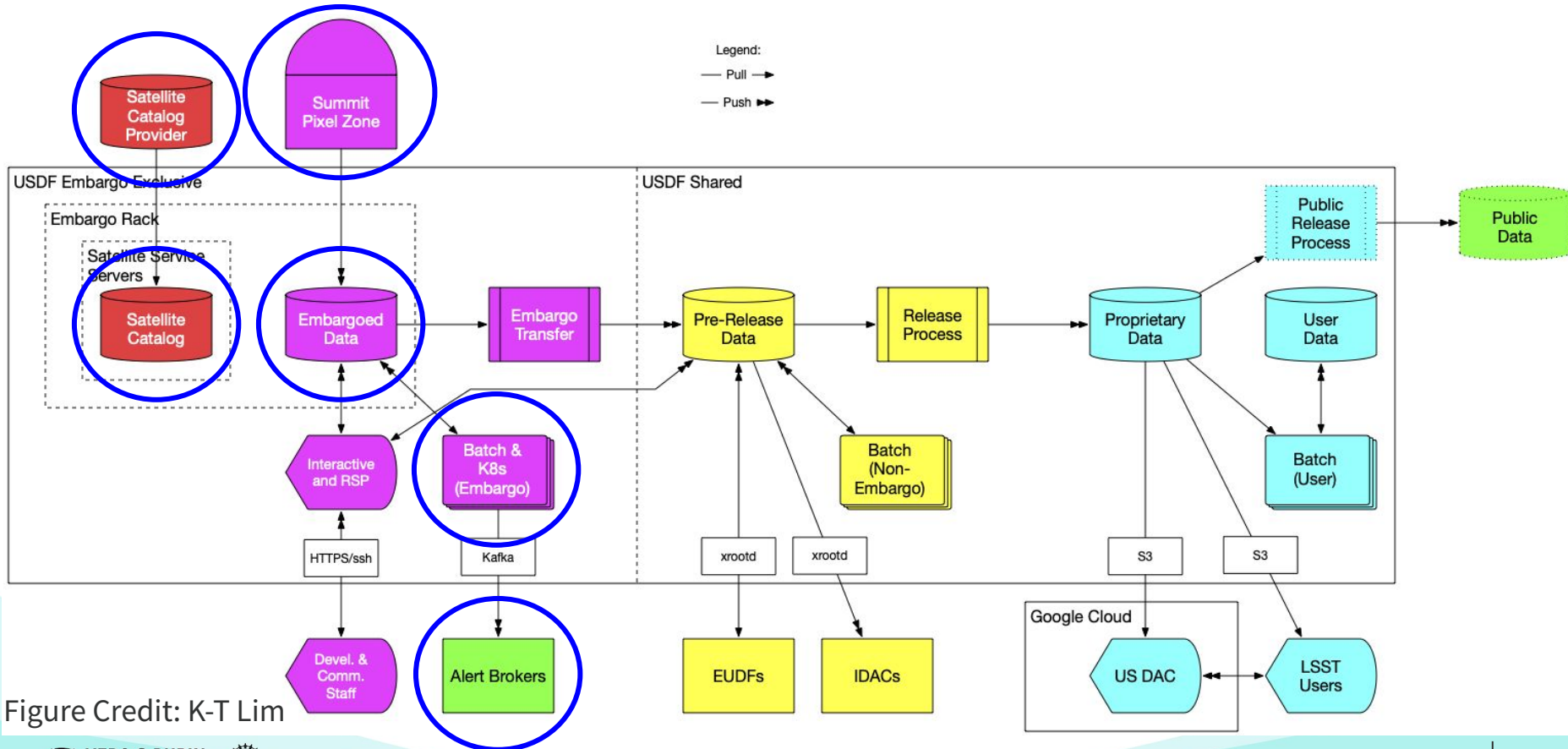


Figure Credit: K-T Lim

Data-release processing

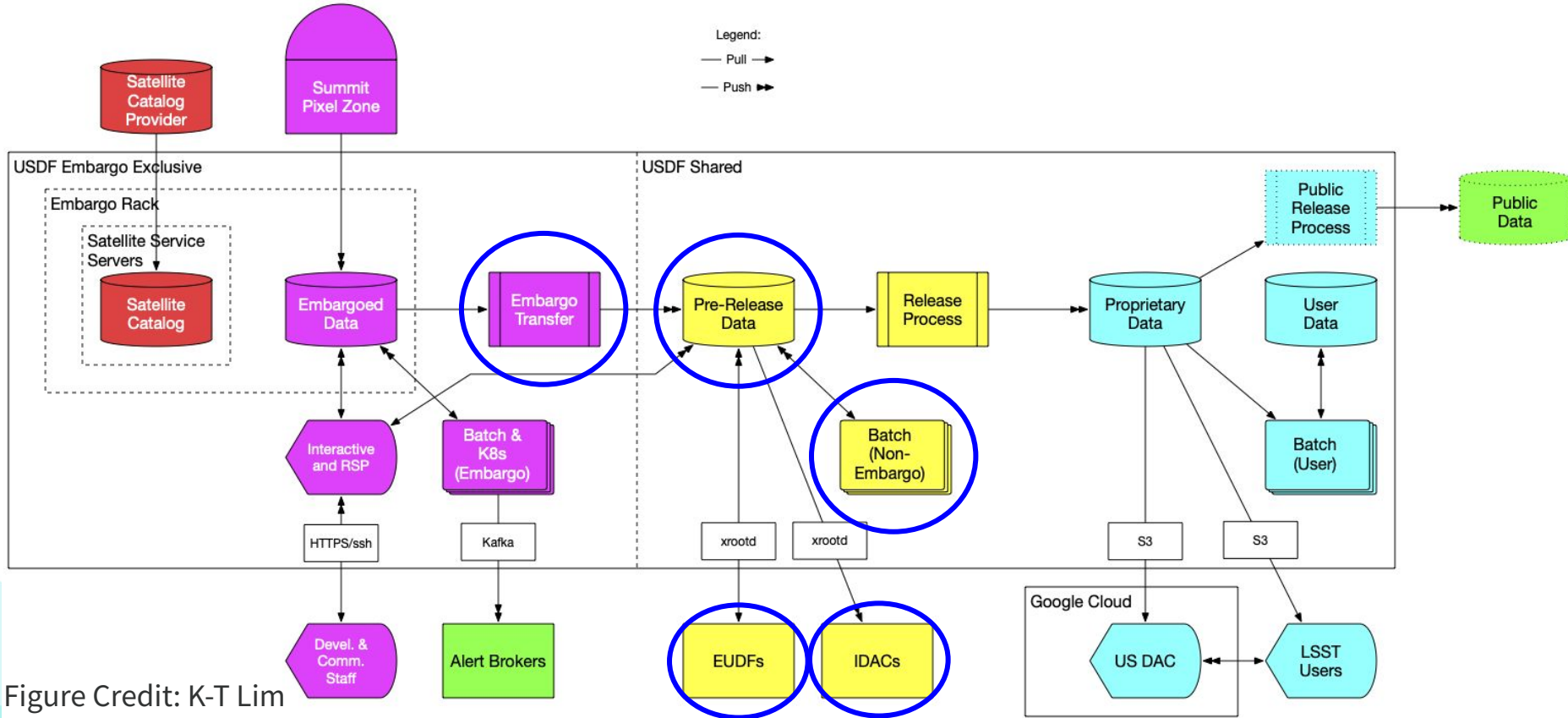


Figure Credit: K-T Lim

Scientific community data access

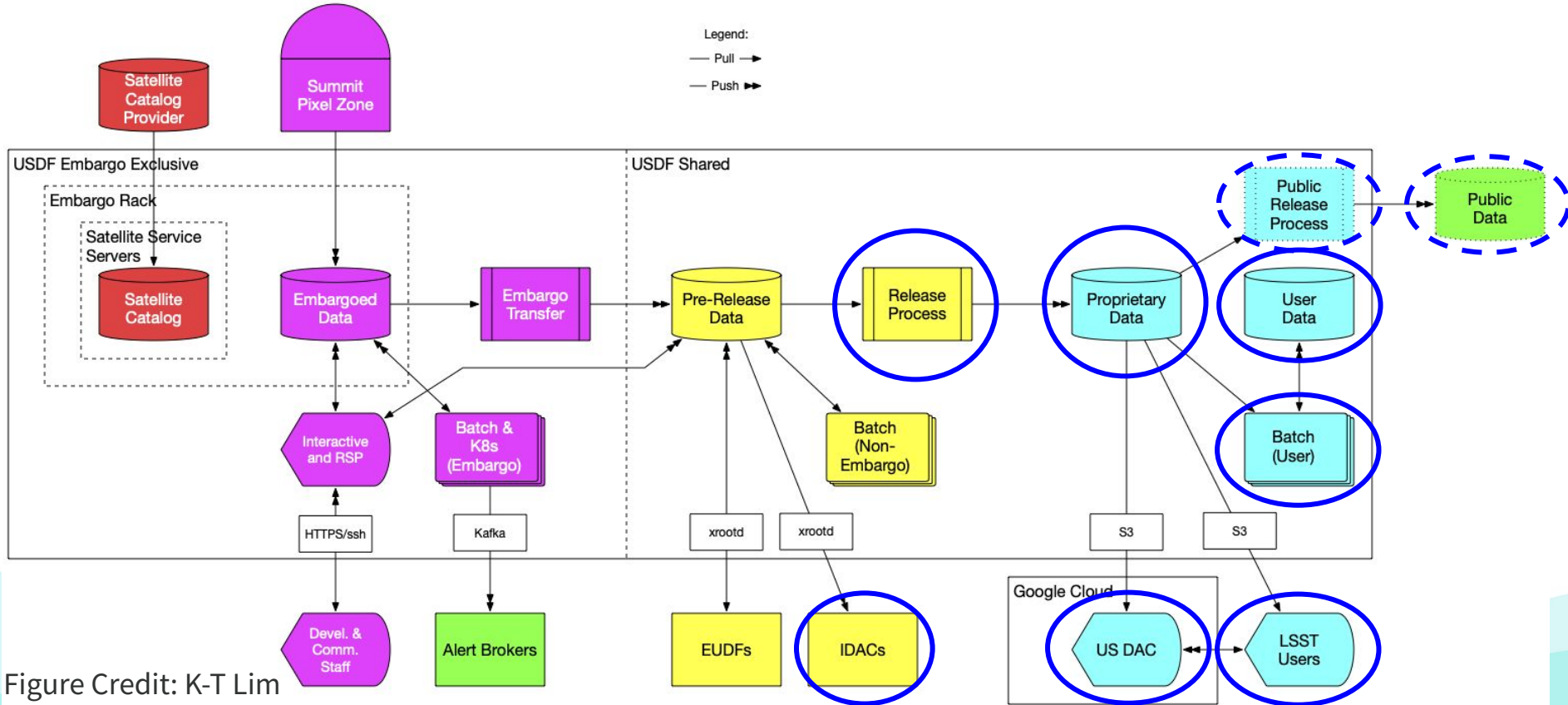


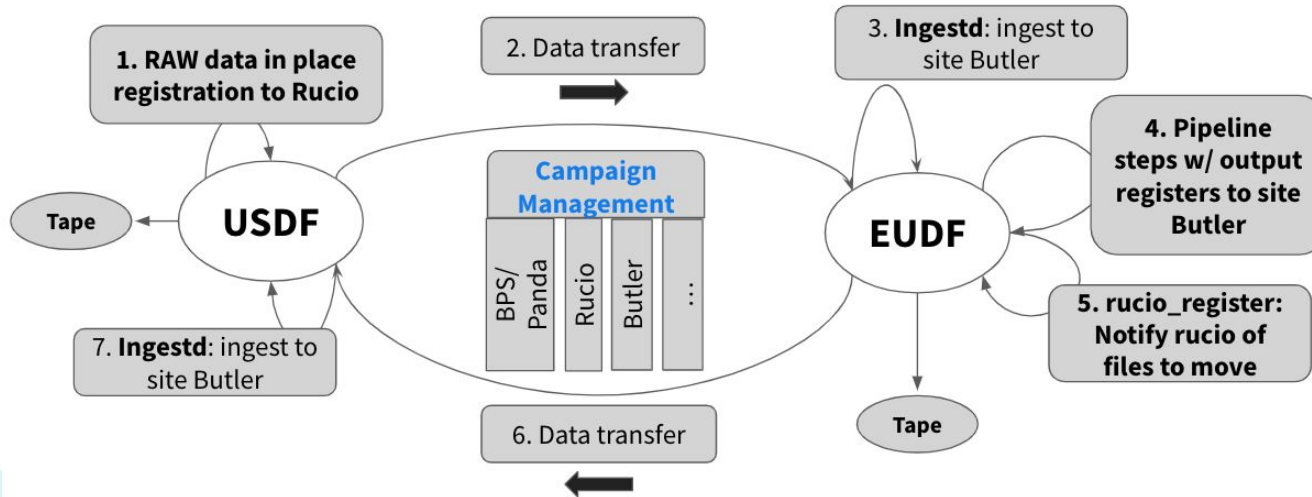
Figure Credit: K-T Lim

Update on multi-site exercise

For more detail see:
[R. Dubois 2024 presentation](#)
[F. Hernandez et al. presentation](#)

Full chain multi-site exercise involving Campaign Management Team orchestrated:

- Panda-driven pipeline software execution
- Rucio/FTS driven data movement of input data (USDF → European facilities) and data products (European facilities → USDF)



Update on Rucio-based inter-facility data movement

Rucio-driven input data for Data Preview 1 (LSSTComCam) transferred to FrDF:

- Completed successfully
- ctrl_ingstd, the software that connects Rucio and Bulter is involved
- Using FTS service at RAL. Working on setting up a FTS service at USDF.

Daily export of unembargoed in-dome calibration exposures from USDF to FrDF ongoing since mid-April:

- Currently exploring solutions to fix a transfer error rate (5 to 10%) which seems linked to limitations in the network configuration of Rubin's Rucio instance

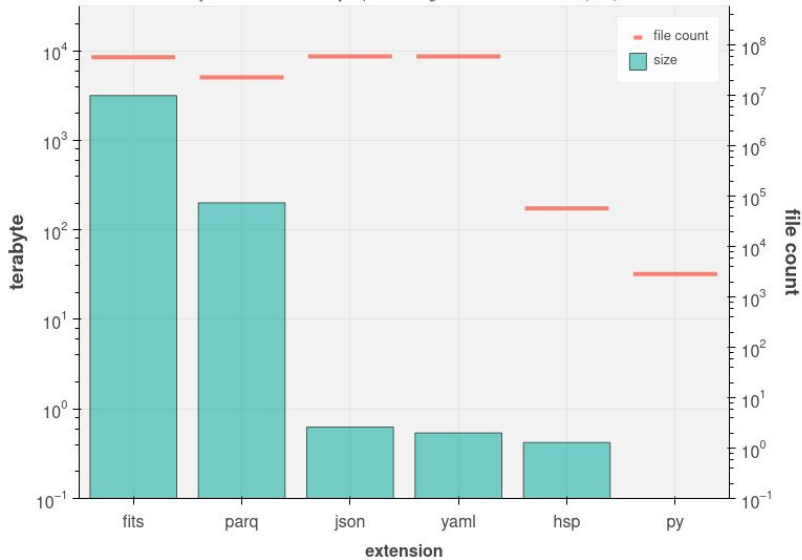
Expanding Rucio usage:

- Setup and performed real calibration data transfer to Summit and Base
- Successfully completed a pilot test involving two IDACs (Canada and Poland)
- Now planning to add a few more IDACs

File zipping

DP0.2 products: file count and aggregated size

Rubin Observatory French Data Facility – processing for Data Preview 0.2 (v23)



DP0.2 exercise at FrDF indicated that Rubin would create billions of small file objects a year.

- This is a significant stress to Butler/Rucio/FTS and storage systems.
- Proposed zipping small files and we now see the benefits:
 - Unembargoed raw images of a focal plane are put in a single zip, which includes all CCDs (20 MB to 4 GB per file)
 - Middleware supports zipping of any data type.
 - Reduce load on Rucio/FTS and storage systems
 - Butler won't benefit though
 - Will make larger zip (of Rucio datasets) when backing up to USDF HPSS system

Current USDF focus areas

- Support for on-sky engineering and commissioning
- S3DF/USDF facility performance, reliability, and observability
- Data embargo / data security requirements
- Preparation for Data Preview 1
- Multi-site processing (see above)
- Standardizing application deployment and operations support model
- End-to-end dashboard development for real-time science data flow
- ConOps / R2A2's for “data wrangling”
- Hardware roadmap & sizing model evolution
- Navigating cybersecurity landscape

Major 2025 Rubin milestones

- March 2025: LSSTCam on telescope
- April 2025: “First Photon” (*LSSTCam on-sky engineering starts*)
- June: Rubin First Look event (*Public splash of early Rubin images*)
- Also June: Data Preview 1
 - *First release of Rubin (ComCam) data to scientific community*
- July(-ish): start Science Validation Surveys, through ~September
- Early October: 10-year Legacy Survey of Space and Time starts

end / thanks!

List of applications & services at USDF

Consolidated Database (ConsDB)
Butler Databases (multiple)
Alert Production Database (APDB)
Prompt Products Database (PPDB)*
Exposure/narrativelog replica
DP0.3 Postgres
Minor Planet Survey replica
InfluxDB Enterprise for the EFD
Qserv (Docker)
Qserv (k8s operator)
Postgres for CM-service
Postgres for PanDA
Postgres for Rucio
Alert database
Embargo Transfer

(* to GCP?)

PanDA
Rucio + FTS3
Rucio-Butler integration
Jenkins Control
Embargo and LFA ingest
Plot navigator
Scheduler Pre-Night
RubinTV
Rapid Analysis
Alert Distribution
CM-service
Large File Annex (LFA) replication
OpenSearch
HTCondor workflow system
BPS (Batch Production Service)
Prompt Processing

Rubin Science Platform at USDF
Sasquatch
Shared stack maintenance
Shared datasets maintenance
ArgoCD
Phalanx
Sasquatch Kafka
Shared Butler repos
Sat db
Metrics Analysis Framework
s3proxy
Data Transfer Summit-USDF
sattle deployment (pending)
Exposure Checker
Scheduler snapshot dashboard

(+ a few recent additions not listed!)