# Workflows for the Dark Energy Science Collaboration and the Rubin Observatory

Jim Chiang
SLAC
HEP-CCE All Hands, 2023-12-18

# Rubin Workflows

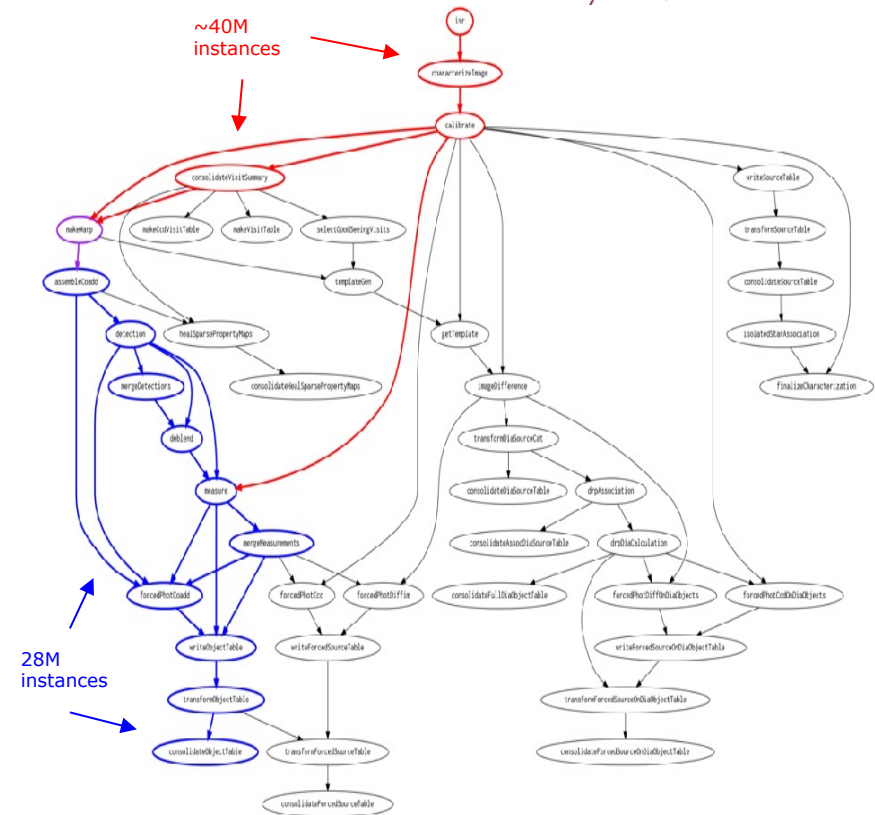**Large image processing pipelines at scale:**

- Complicated DAGs
- Large numbers of small, interdependent tasks, with widely varying resource needs.

**Orchestrating processing campaigns at multiple sites:**

- Moving large numbers of files and distributing payloads in an optimal fashion

**Long campaigns:**

- We expect Data Release Processing campaigns to last ~200 days/year.



Graph showing dependencies between task types for Rubin image processing.

# Rubin Workflows:
# Current solutions/technologies

**Rubin LSST middleware provides:**

- DAG-generation
- Plugin-based batch processing framework that works with various workflow management systems: PanDA, Parsl, DAGMan/HTCondor.  *HEP-CCE*
- Data "Butler" as an abstraction layer between data products and processing tasks allowing for different kinds of backend storage, e.g., posix, S3, etc..

**For Rubin production work, PanDA chosen for multi-site capabilities, and Rucio for data management between sites.**

- Both are extensively used and supported in the HEP community.
- Significant friction because of differences between Rubin processing needs and more typical HEP workloads.  *HEP-CCE*
- Development needed to integrate Rubin Butler with Rucio.
- Monitoring the finely grained Rubin processing across the three different data facilities–USDF, FrDF, UKDF–is challenging.  *HEP-CCE*

# DESC Workflows

**DESC image processing pipelines:**

- Running subsets of Rubin pipelines with alternative data selections (including injecting simulated objects), algorithms, and configurations to understand systematics.
- Joint pixel analysis of Rubin data with data from other observatories.

**DESC cosmology analysis pipelines:**

- These will be run at various sites, mostly HPC, but also at local university clusters.
- They consist of many heterogeneous pipelines using different kinds of compute: MPI-based, AI/ML, MCMC.
- We expect challenges in managing data products exchanged between these pipelines because of distributed nature of the processing.
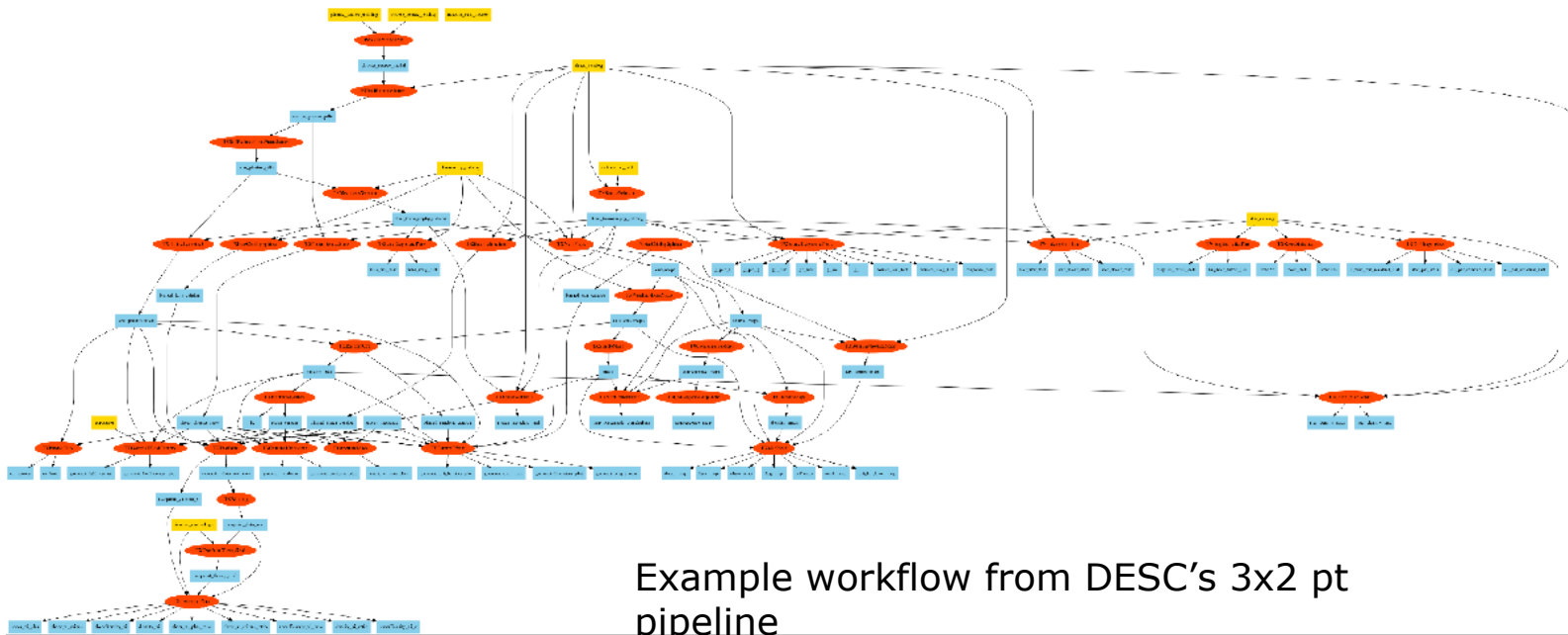
# DESC Workflows:
# Current solutions/technologies

- Custom pipeline framework for analysis pipelines (Ceci)
- Parsl for workflow management at HPC facilities and local clusters
  *HEP-CCE*
- Custom solution for data management  *HEP-CCE*



Example workflow from DESC's 3x2 pt pipeline