



cnrs

Centre de Calcul
de l'Institut National de Physique Nucléaire
et de Physique des Particules

Rubin-LSST data processing at CC-IN2P3

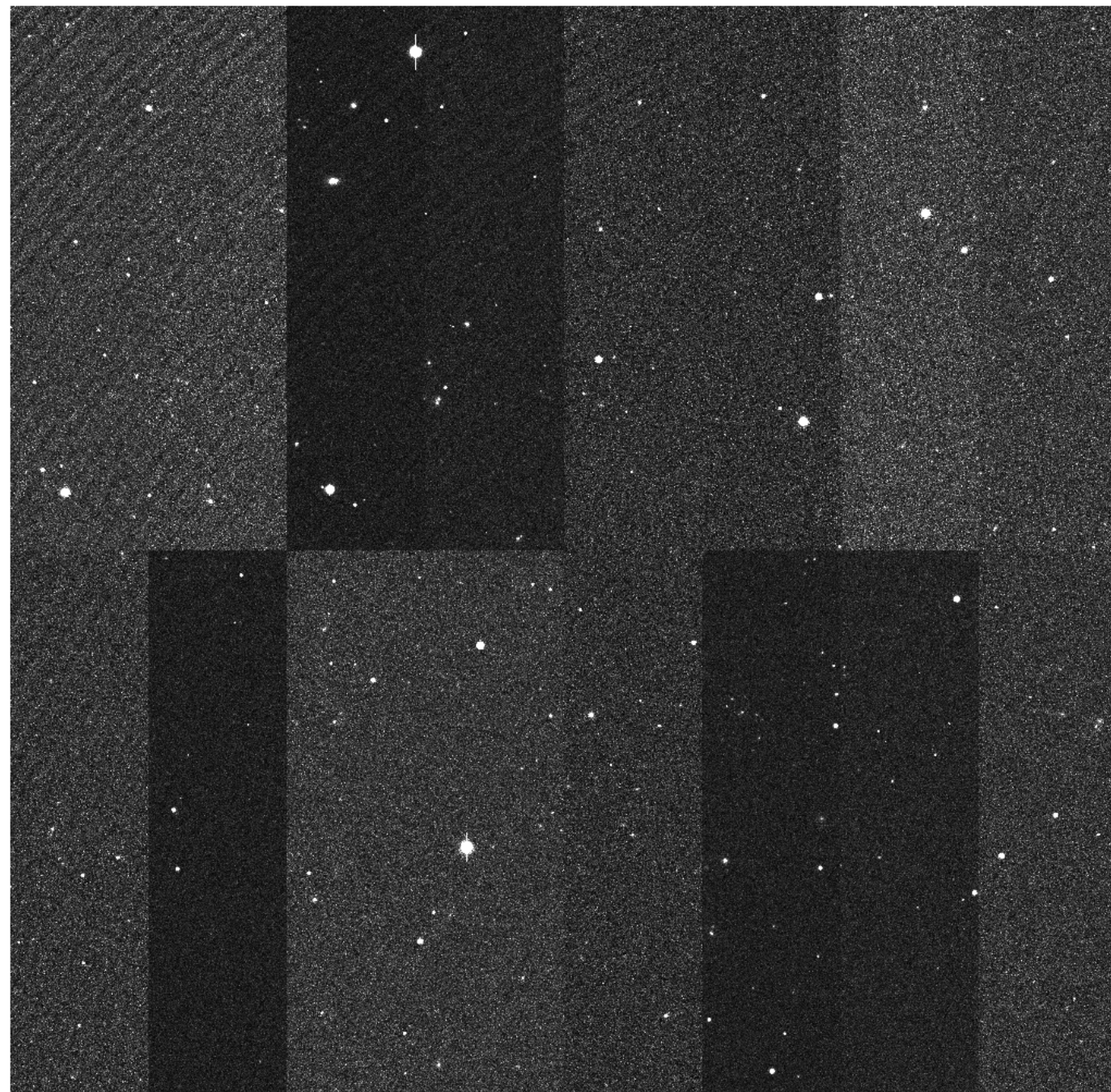
Quentin Le Boulc'h



Data Release Production pipelines

DRP in a nutshell:

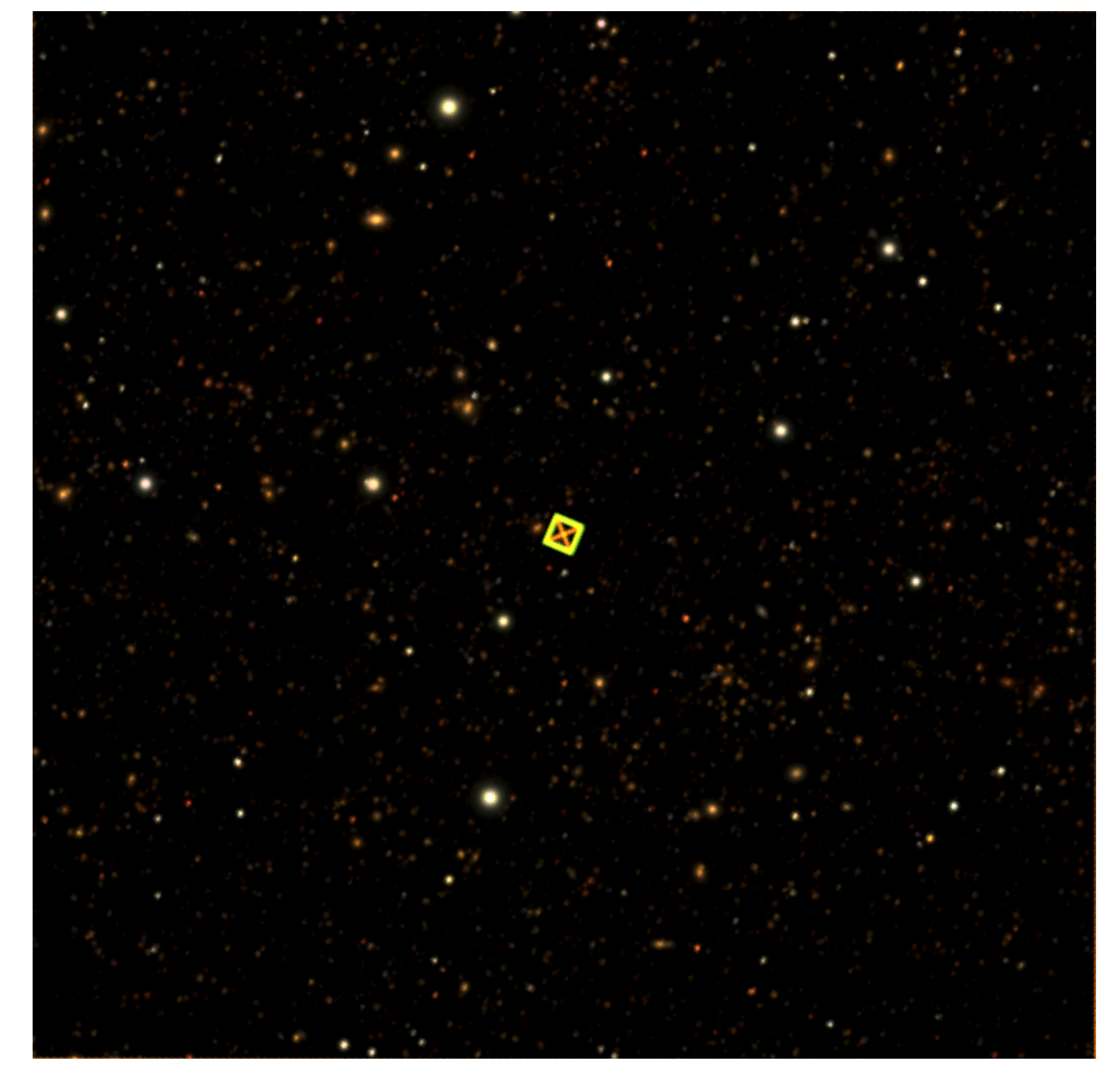
- **Process raw images** and calibrate them to remove any artifact from instrument
- Merge images together ("**coaddition**") to improve signal/noise
- **Measure differences** between images to detect changes
- Detect objects, measure their properties and populate **catalogs**



07/04/2023



Rubin-LSST data processing at CC-IN2P3



Data Preview 0.2 (DP0.2)

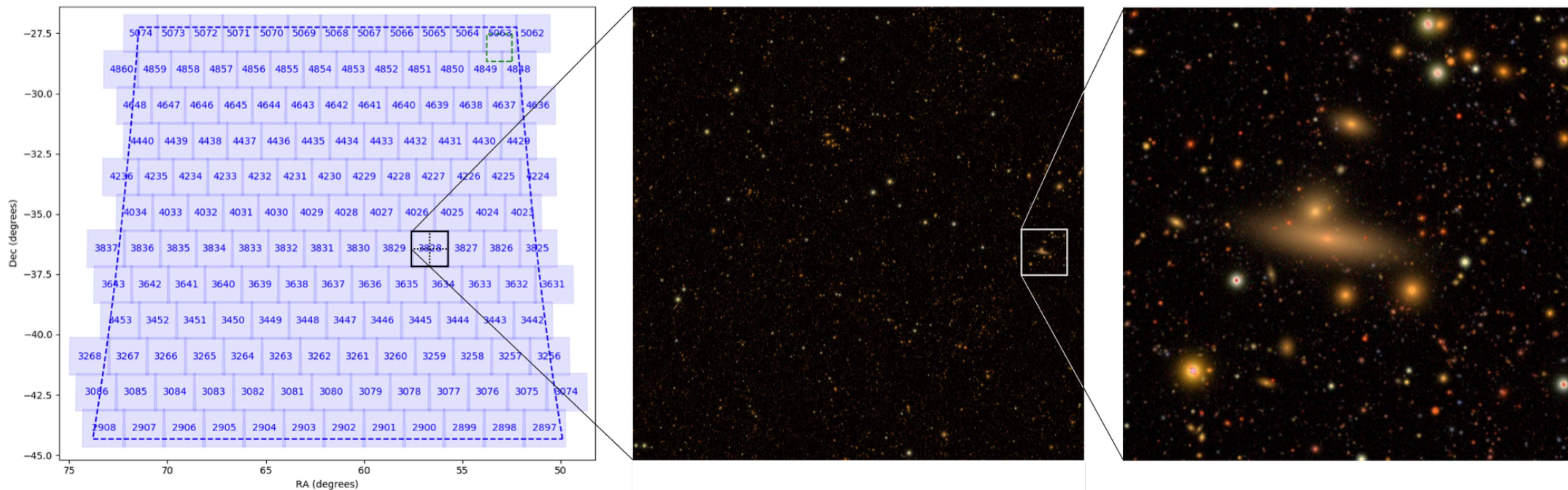
- Data Preview 0:
 - Early integration test of pipelines and Rubin Science Platform
 - Provide data for preparation of science analysis
- DP0.2:
 - Reprocessing of simulated data using latest pipelines
 - Gen3 butler (data management)
 - Workflow automation

Data Product	DP0.1	DP0.2	DP1	DP2	DR1
	DC2 Simulated Sky Survey	Reprocessed DC2 Survey	ComCam On-Sky Data	LSSTCam On-Sky Data	LSST First 6 Months Data
Raw images	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
DRP Processed Visit Images and Visit Catalogs	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
DRP Coadded Images	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
DRP Object and ForcedSource Catalogs	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
DRP Difference Images and DIASources	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
DRP ForcedSource Catalogs including DIA outputs	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
PP Processed Visit Images	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
PP Difference Images	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
PP Catalogs (DIASources, DIAObjects, DIAForcedSources)	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
PP Alerts (Canned)	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
PP Alerts (Live, Brokered)	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
PP SSP Catalogs	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>	<input checked="" type="checkbox"/>
DRP SSP Catalogs	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input type="checkbox"/>	<input checked="" type="checkbox"/>

Data Preview 0.2 (DP0.2)

Simulated Rubin images generated by DESC (Dark Energy Science Collaboration):

- 300 deg² (full survey ~ 20 000 deg²)
- 5 years (full survey: 10 years)
 - approximately **0.5% of the full survey**, or 10 nights of data gathering



Data Preview 0.2 (DP0.2)

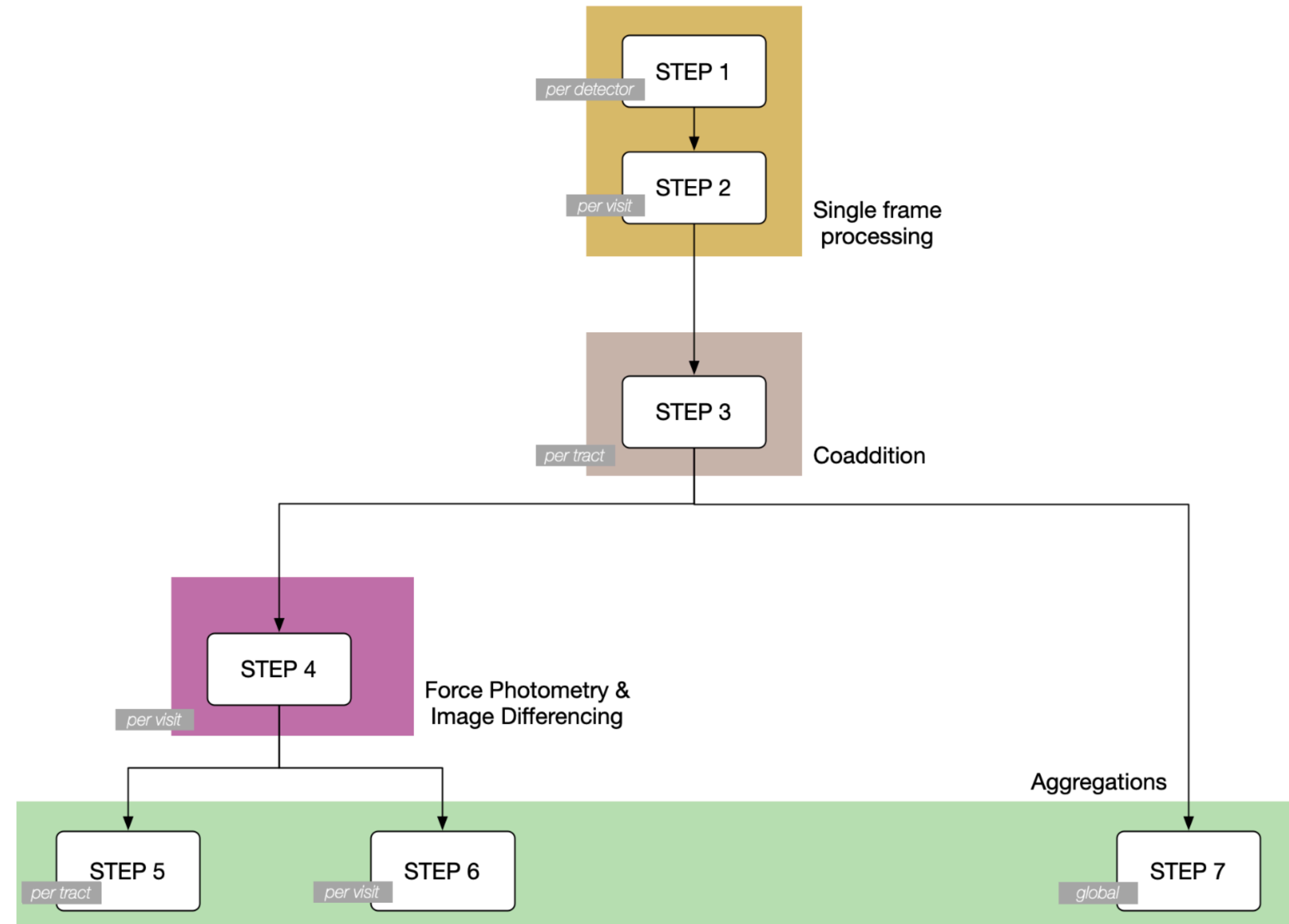
Simulated Rubin images generated by DESC (Dark Energy Science Collaboration):

- 300 deg² (full survey ~ 20 000 deg²)
- 5 years (full survey: 10 years)
 - approximately **0.5% of the full survey**, or 10 nights of data gathering

- 300 deg² = about 30 pointings
- 5 year = hundreds of visits per pointing
 - approximately **20,000 Rubin camera simulated exposures**, each with up to 189 detectors
 - **3 M files**, 18 MB each, **50 TB** in total

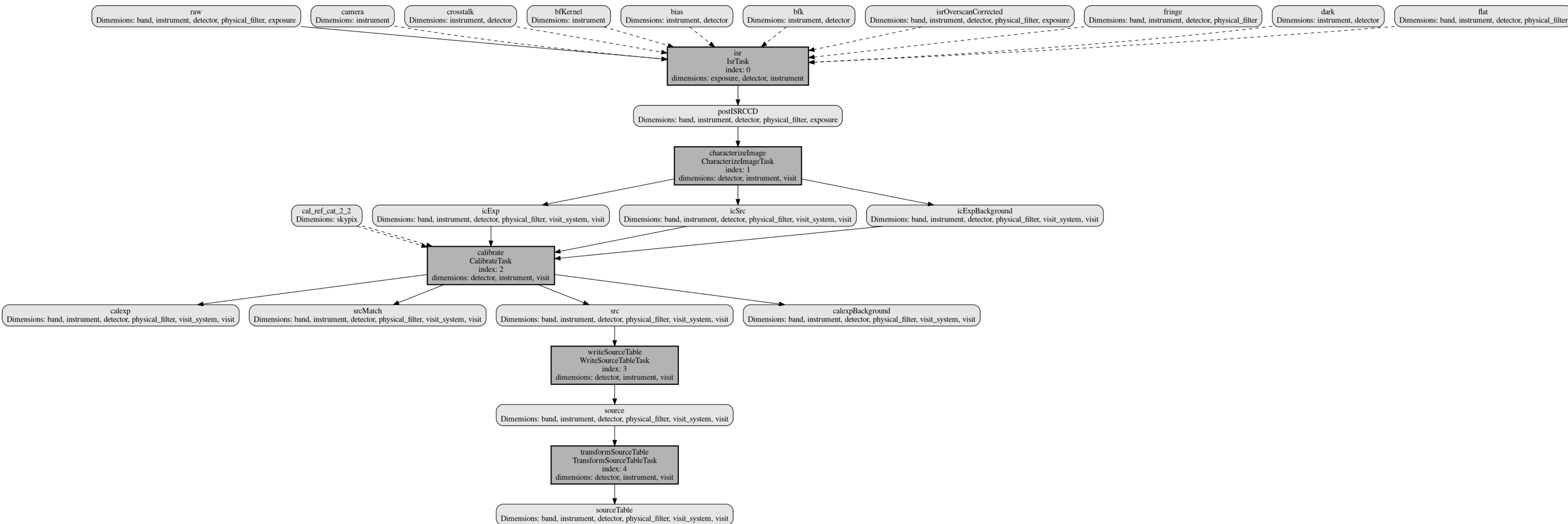
Workflows of pipeline tasks

- Set of ~ 80 pipeline tasks, grouped in 7 steps
- Each step process data at different levels: single CCD, full visit, coadded images, catalogs...



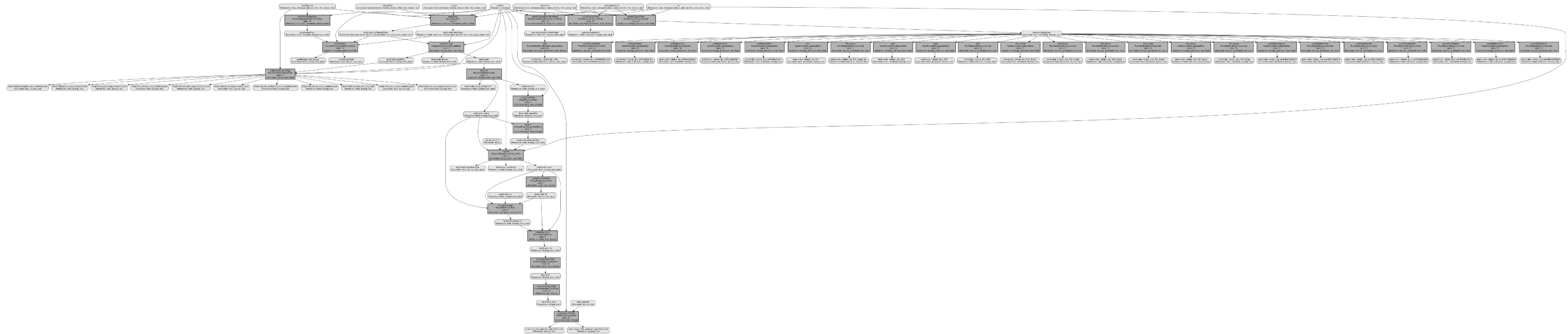
Workflows of pipeline tasks

- Pipeline tasks depend on each other through data production / consumption
- Some of the steps have rather straightforward workflows:



Workflows of pipeline tasks

- Pipeline tasks depend on each other through data production / consumption
- Others are more complex (step3):



Workflows of pipeline tasks

- Each task must be executed between 1 and 3 millions times
 - Tasks execution time goes from few seconds to 24+ hours
 - Maximum memory usage for each task varies from 1 GB to 200+ GB
 - For a given task, memory usage depends on input data
- **Automation!**

Batch Processing Service (BPS): LSST framework for distributed pipeline execution

- Generate a workflow and submit it to an external **Workflow Management System (WMS)**
- A plugin is needed to interface BPS to the WMS
- Explore existing WMS that could be plugged into BPS to run the processing workflows on our computing platform

Workflows of pipeline tasks

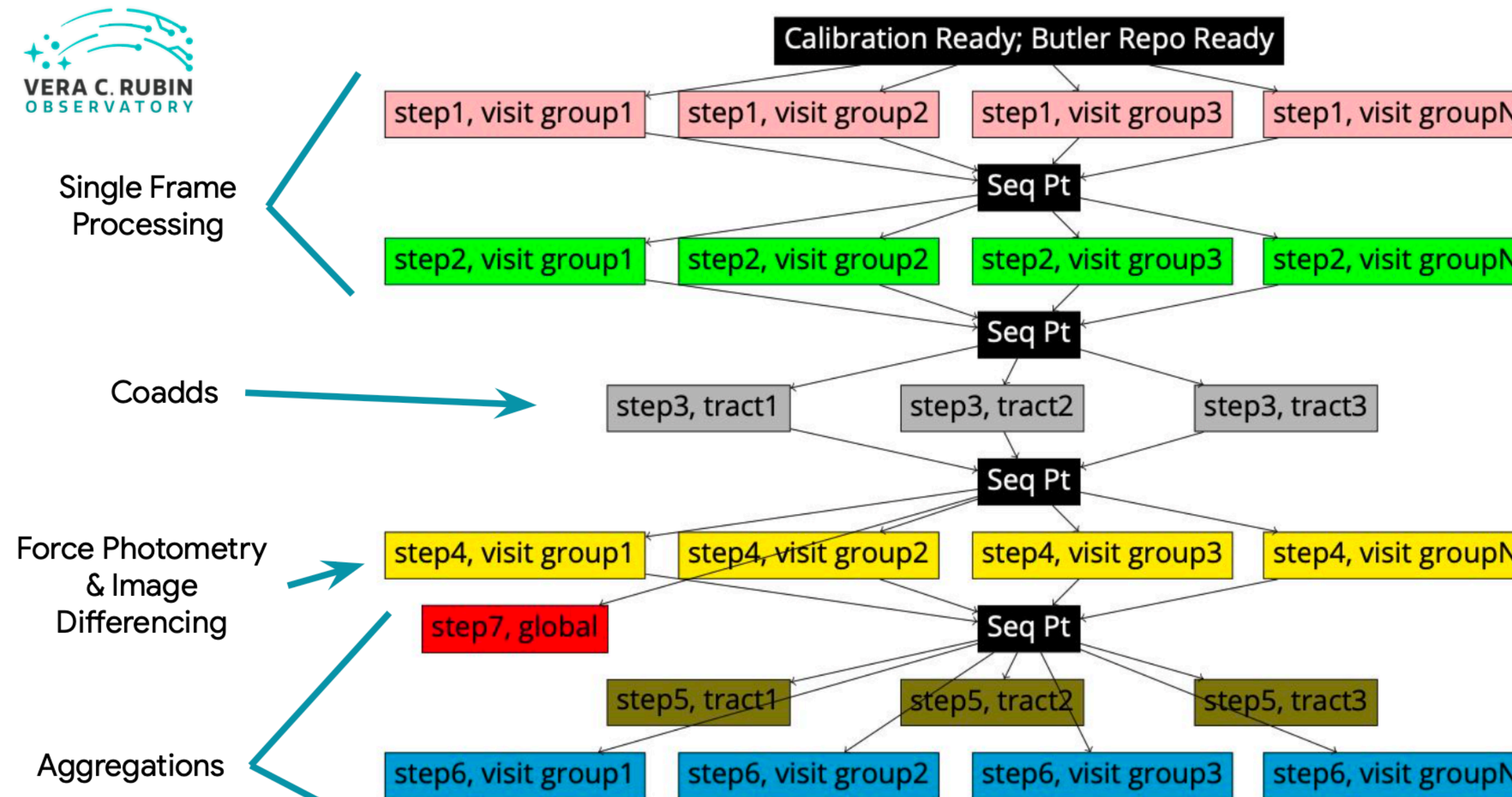
- HTCondor: only available through grid interfaces at CC-IN2P3
- We needed a WMS that:
 - Can talk to **Grid Engine and Slurm**
 - Use **pilots jobs**
 - Is **scalable**
 - Does not require too much development to interface it with BPS
- **Parsl**: "extends parallelism in Python beyond a single computer"
 - Can submit jobs to Grid Engine and Slurm
 - Provides pilot jobs through its HighThroughputExecutor feature
 - Claimed to be scalable
 - Already existing plugin developed by J. Chiang (SLAC)
- All processing done using Parsl + plugin + configuration/customization

Execution of DP0.2 campaign

- Full production from step 1 to step 7
 - From April to October '22
 - Using up to **3,000 simultaneous Slurm jobs**
 - Executed 57,903,740 pipetasks which consumed 2,347,306 elapsed hours
- Butler repository details:
 - Input dataset: 51 TiB, 2.9 M files
 - Products: **3 PiB**, 54 M directories, **201 M files** (including intermediates)
 - Registry database: 314 GiB

Execution of DP0.2 campaign

- Some of the steps can not be executed in a single workflow:
 - Workflow is too large and its generation used too much memory and time
 - Execution too long
 - Scalability issues
- ➔ Divided into sub-workflows of "reasonable" size (up to 32 workflows for step3)

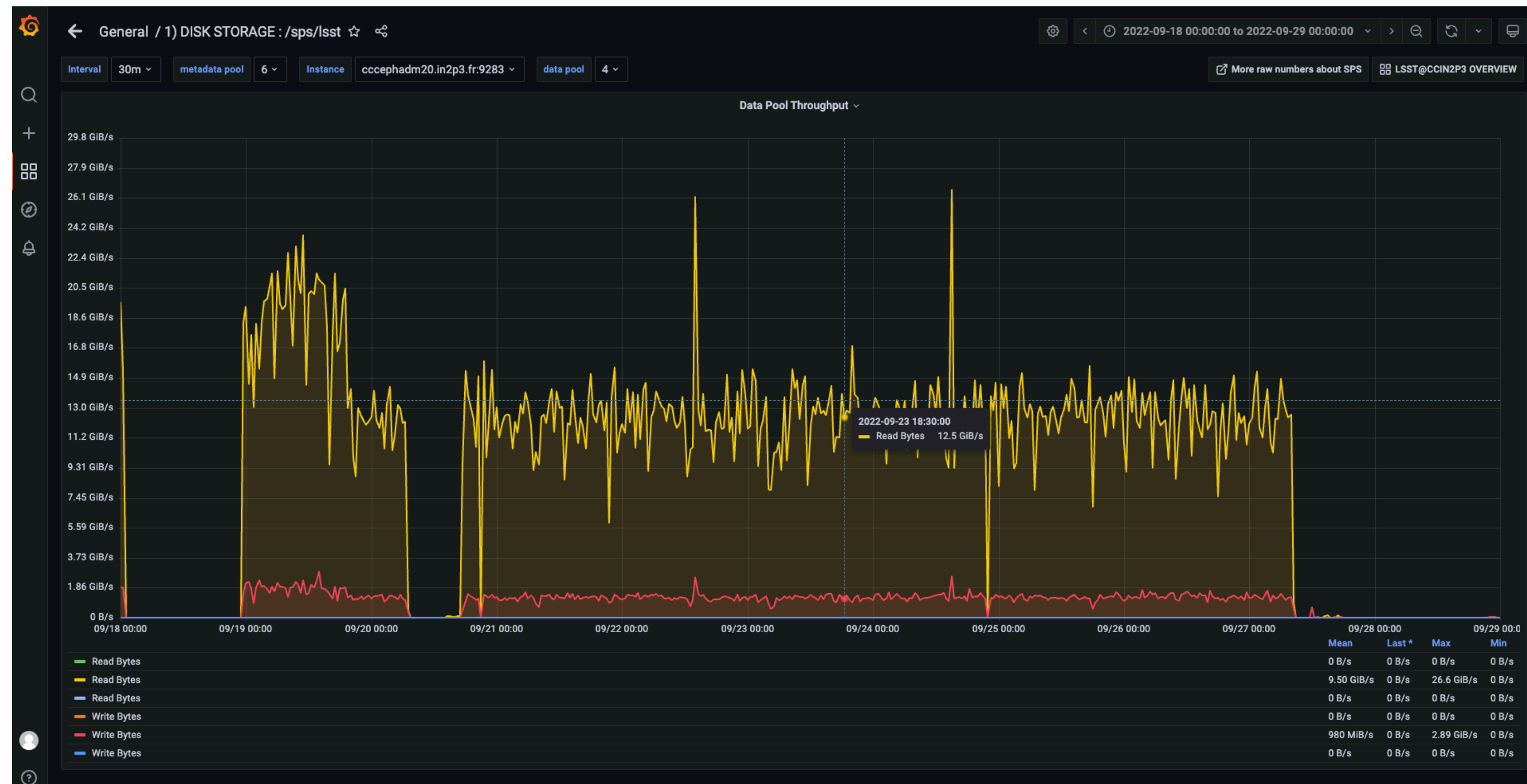


Execution of DP0.2 campaign

- Number of tasks in the workflows can be reduced thanks to the **clustering** feature of BPS:
 - Workflow graph split into subgraph where tasks are grouped together ("cluster")
 - A cluster is seen as a single task to be executed
 - Better scalability, but makes resources tuning and debugging more complex

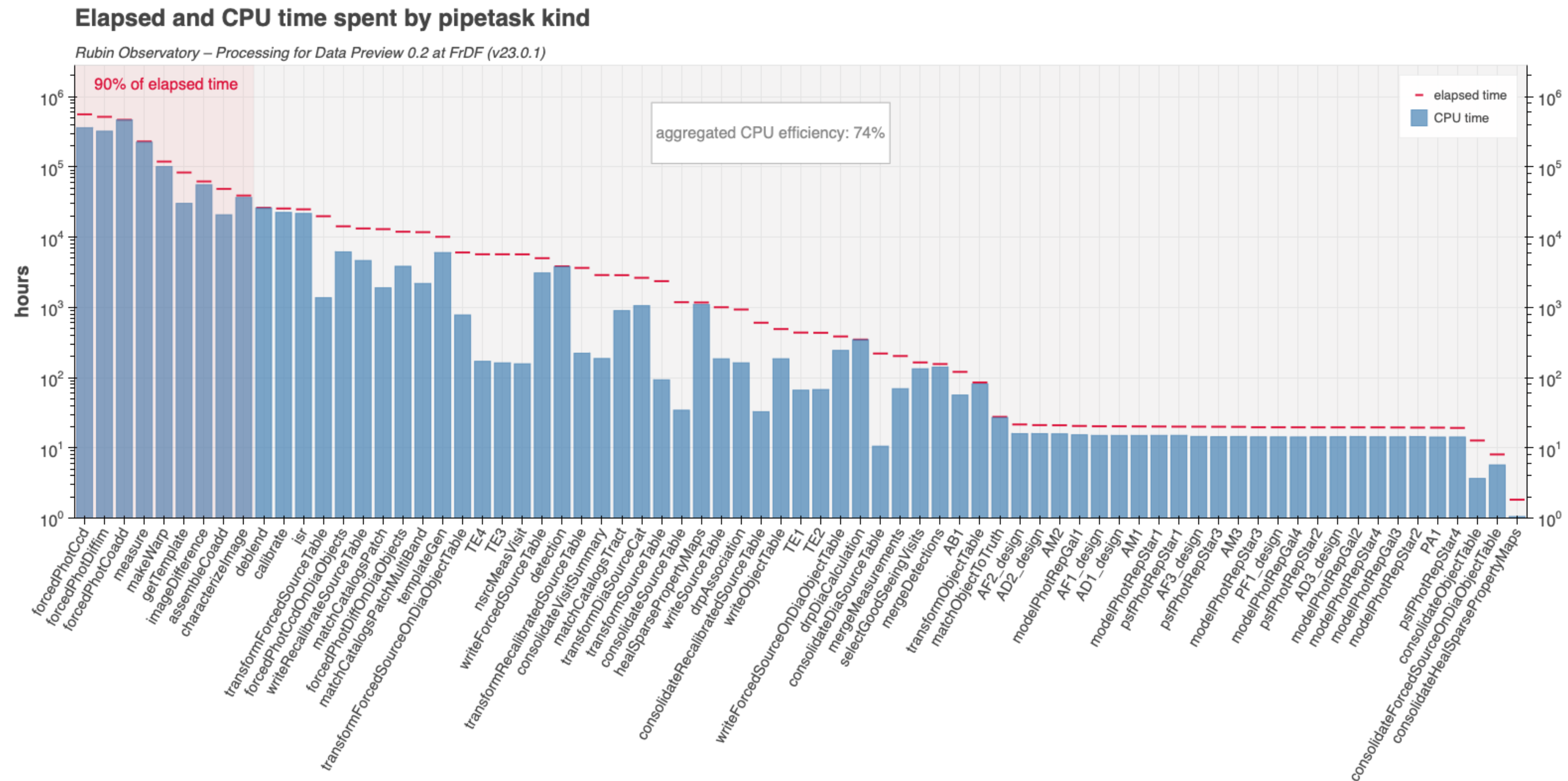
Execution of DP0.2 campaign

- Storage system: /sps/lst (CephFS)
 - **Very good behaviour** overall
 - Performance issue for some tasks that required to copy the files to the local disk of the compute node: significant improvement (80% CPU efficiency, induced I/O throughput reduced to 12 GiB/s)



Execution of DP0.2 campaign

- 90% of total computing time used by 9 tasks: to be optimized in priority
- These tasks have very good CPU efficiency
- Overall CPU efficiency is good: 74%

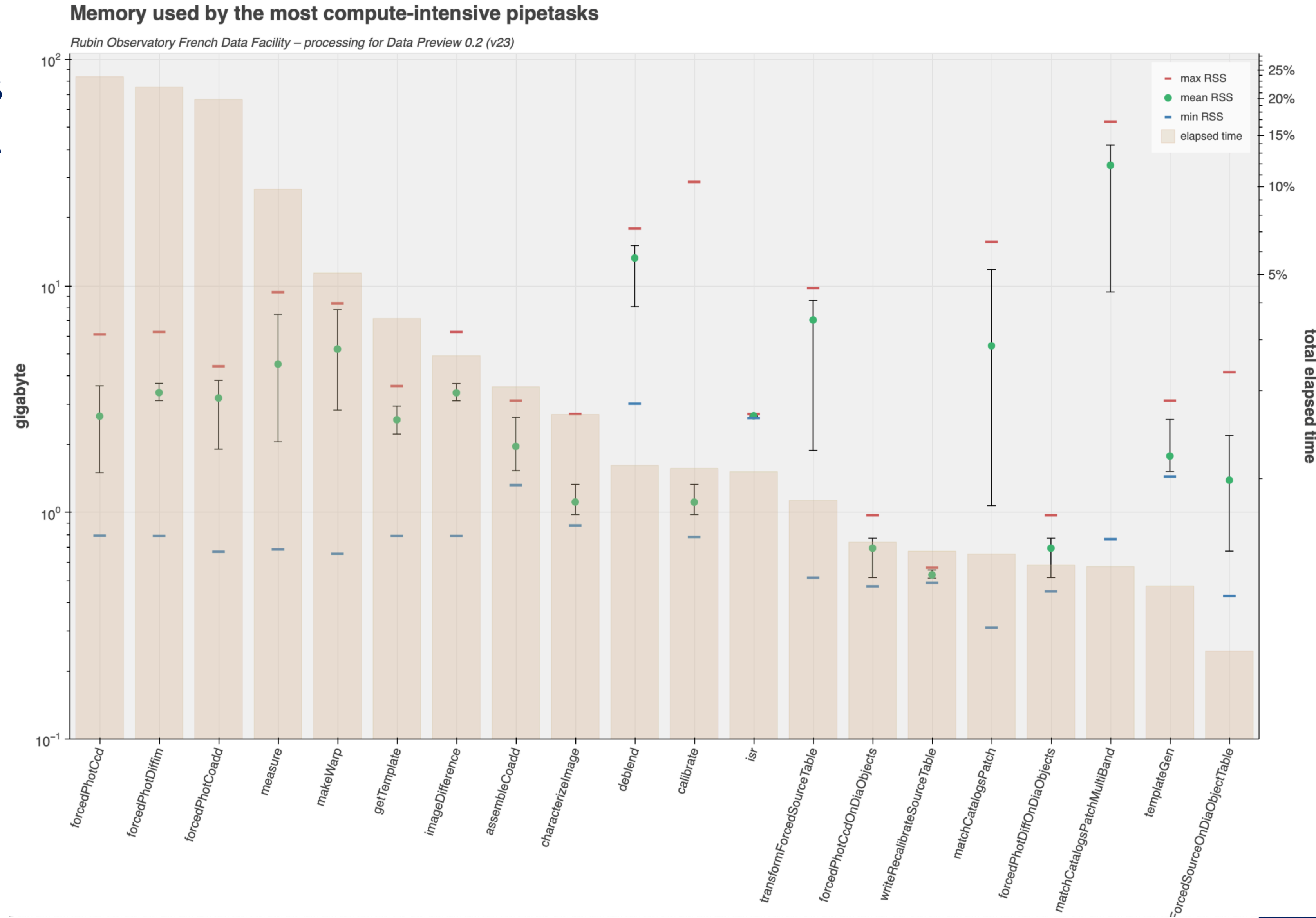


Execution of DP0.2 campaign

- Tasks with highest CPU usage have memory needs between 1 and 10 GB
- Tasks with moderate CPU usage have higher memory needs, up to 50 GB

Several options for improvement:

- Jobs tuning
- Pipeline optimisation
- Batch system configuration
- Hardware



Current sources of informations:

- Metrics provided directly by the LSST processing framework
- Informations provided by the job scheduler (Slurm)
- Internal system metrics on nodes, storage system, network...
- Profiling?