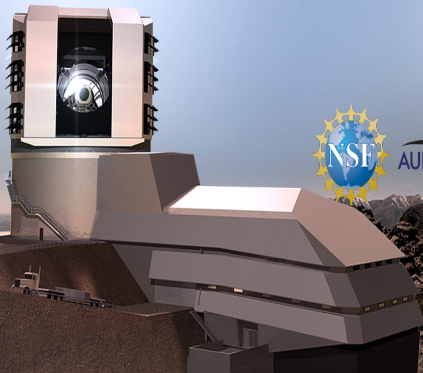**Rubin**
Observatory

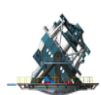# Optimizing the Rubin DRP
# CC–IN2P3, Lyon

**J. Bregeon, D. Boutigny, F. Hernandez, Q. Leboulc'h, D. Parello, C**

**December 14, 2023**

NSF  AURA  U.S. DEPARTMENT OF ENERGY  SLAC  CHARLES AND LISA SIMONYI FUND  LSST CORPORATION
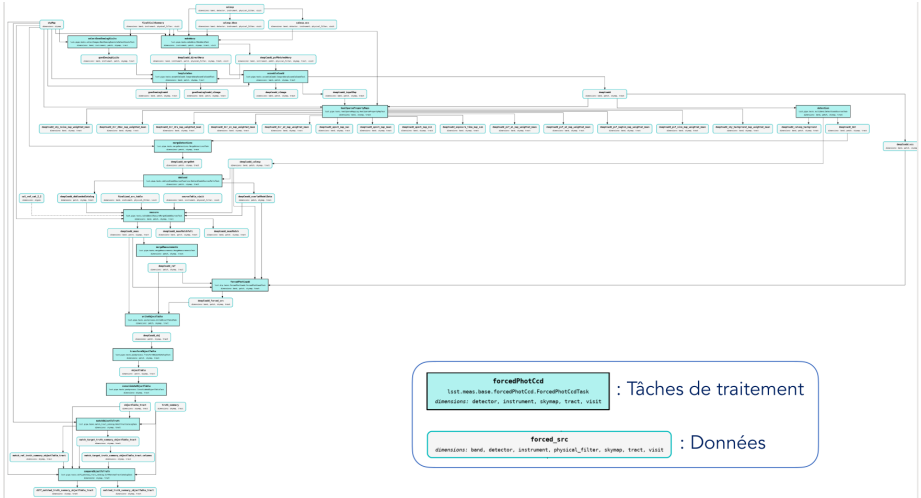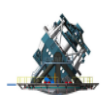
**Rubin** Observatory

- CC-IN2P3 is one of the 3 main Rubin LSST data center (35% US, 25% UK)
  - IN2P3 signed to process 40% of LSST data
  - Rubin **DRP**
- A lot of data
  - 6.4 GB per exposure, i.e. 16 TB per night
  - 15 PB of "catalogues" after 10 years

**ENERGY** Office of Science

Rubin
Observatory



- 80 tasks in 7 main steps
- each task is a workflow of its own

U.S. DEPARTMENT OF **ENERGY** | Office of Science

# really complex



forcedPhotCcd
lsst.meas.base.forcedPhotCcd.ForcedPhotCcdTask
dimensions: detector, instrument, skymap, tract, visit

: Tâches de traitement

forced_src
dimensions: band, detector, instrument, physical_filter, skymap, tract, visit

: Données

4

– Interdisciplinary Team
  - Quentin Leboulc'h @ CC-IN2P3
  - Camille Parisel @ APC
  - David Parello @ Univ. Perpignan & LIRMM (Researcher in Computer Science)
  - and Fabio, Dominique and myself. . .

– Goal: Optimize the DRP
  - CPU consumption
  - RAM requirements
  - disk space needs

– Expected results
  - reduce costs and carbon footprint of DRP
  - free some resources for cosmology analysis

Elapsed and CPU time spent by pipetask kind

Rubin Observatory – Processing for Data Preview 0.2 at FrDF (v23.0.1)

# Task level memory profiling



Memory used by the most compute-intensive pipetasks

# py-spy and cprofile on ISR

# Memray on Isr

– Optimizing the code is good on all aspects

– Profiling first

– The team has made a lot of progress

⇒ let's keep our momentum