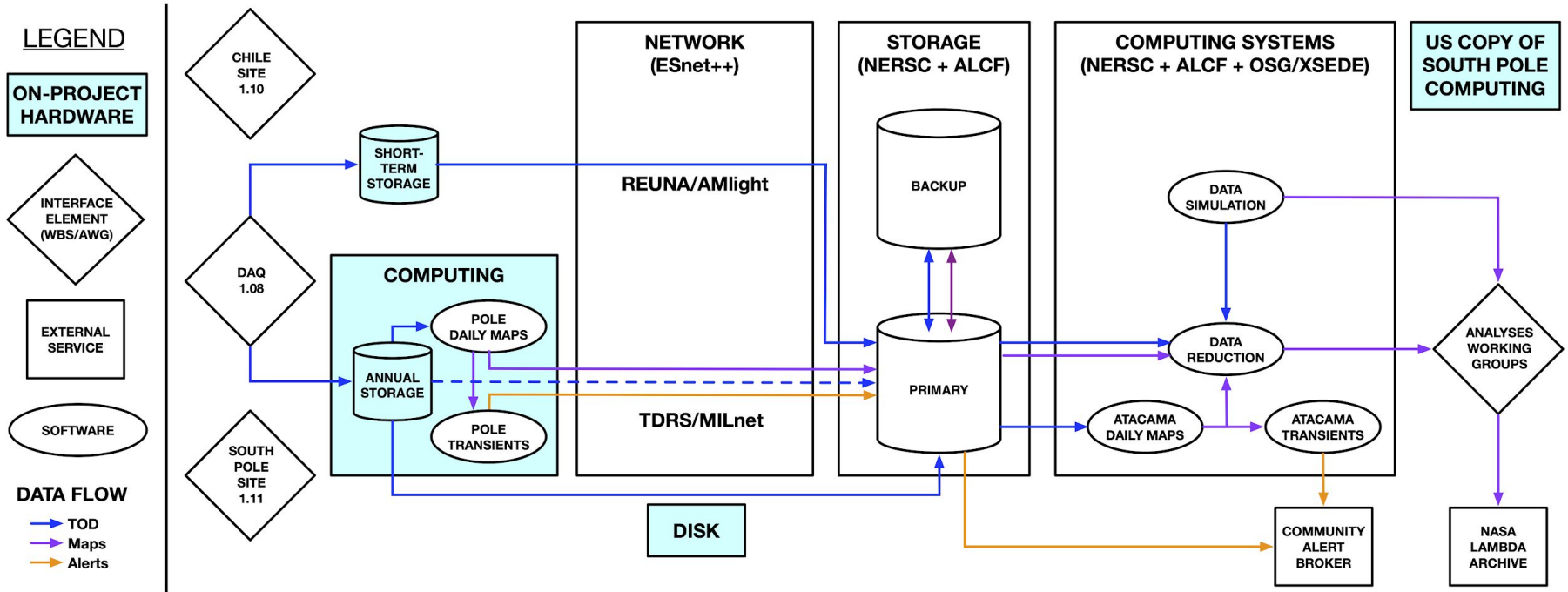# Data Movement

## Sasha Rahlin and Eli Dart
## Data Movement L3

# Data Registration

- Data products added to the file catalog
  - Registry of files - where they are, when/where they were created, etc.
  - File catalog instances at South Pole, Chile, Primary and secondary data centers
- Data placement occurs according to policy
  - Workflow uses file catalog to determine what needs to move
  - File catalog updated as data is moved to data centers
- Data registration drives the rest of the workflow
  - Transmission, Archiving, Distribution, etc. can all be done based on catalog
  - APIs are key
- Data transfer tools should be robust and reliable
- Tools exist in the community today
  - Librarian (https://github.com/HERA-Team/librarian)
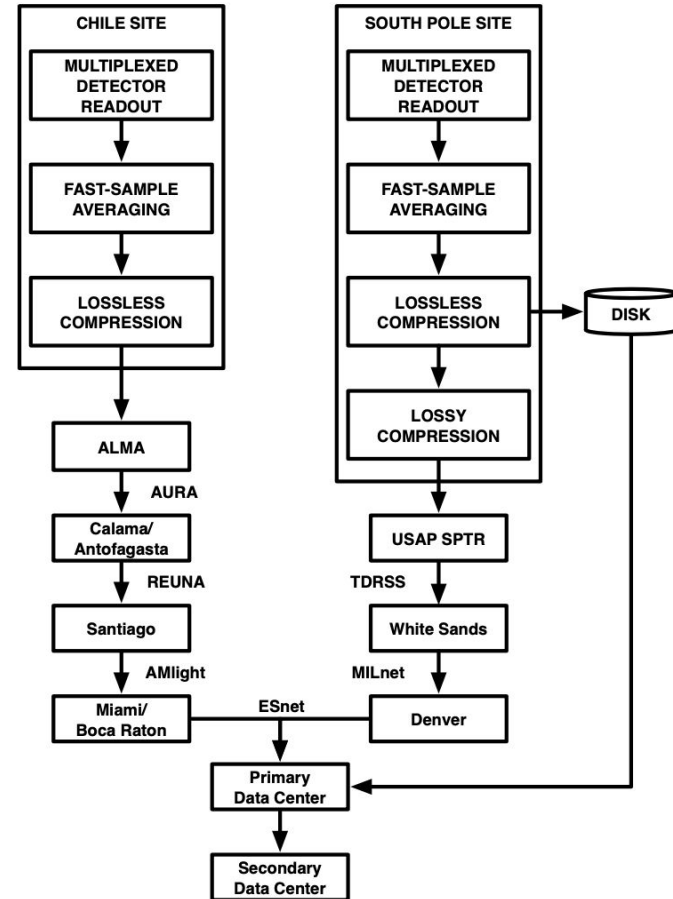  - Rucio (https://rucio.cern.ch)

CMB-S4

# Data Flow



*Allocated computational resources are planned, not confirmed.*

# Data Transmission

- Network paths for South Pole and Chile are very different
- South Pole:
  - Reduced data products via TDRS and MILnet
  - Data stored onsite, annual shipments of disk
  - Lightweight copy of registry, as needed
- Chile:
  - Prompt transfers via fiber-optic networks
  - One month disk buffer to handle outages
  - Up to date registry
- Data Challenges
  - Previous work with Simons Observatory has shown that data transfers are workable
  - This model has been adopted by other experiments and collaborations, and works well



CMB-S4

4

# Data Archiving

- Primary and Secondary data centers host data archives
- File catalog indicates data objects that have been archived
- Archiving process driven by regular automated checks of file catalog
  - Workflow should not require a human in the loop
- Mechanisms for retrieving files from the archive if needed
  - Registry/file catalog provides data location information
  - APIs for retrieval
- Multiple data centers provide archive robustness

CMB-S4

# Internal Data Distribution

- Multiple data products
    - Raw and calibrated time-ordered data
    - Intermediate data products (single-frequency maps) - DM deliverable to AWGs
    - Derived data products - AWG products
- Primary data center is the main source for internal data distribution
- Data volumes dictate use cases
    - Bulk time-domain data must be processed at the data centers
    - Small-scale time-domain data and map-domain data may be processed at the data centers or transferred to other resources
- APIs for data location, registration of data products, etc.

# External Data Distribution

- Public data releases at regular intervals
- Many different data types archived at NASA's LAMBDA archive
  - Maps
  - Catalogs
  - Power spectra
  - Cosmological parameter likelihood codes
- Other data products served directly from NERSC
  - Time domain data
  - Monte Carlo suites
- Data management software stack (with documentation) will also be available

CMB-S4

# Path Forward - Data Challenges

- Plan for data challenges on the path to commissioning
- Progressive testing of features, capabilities, performance
  - Registration, Archive, etc (end2end)
- Preliminary system in place for first data challenge
- Finalize registry system with APIs prior to commissioning on site
- Data challenges will include a replica of South Pole computing to allow testing without using South Pole network bandwidth

# Questions?

**Thanks!**