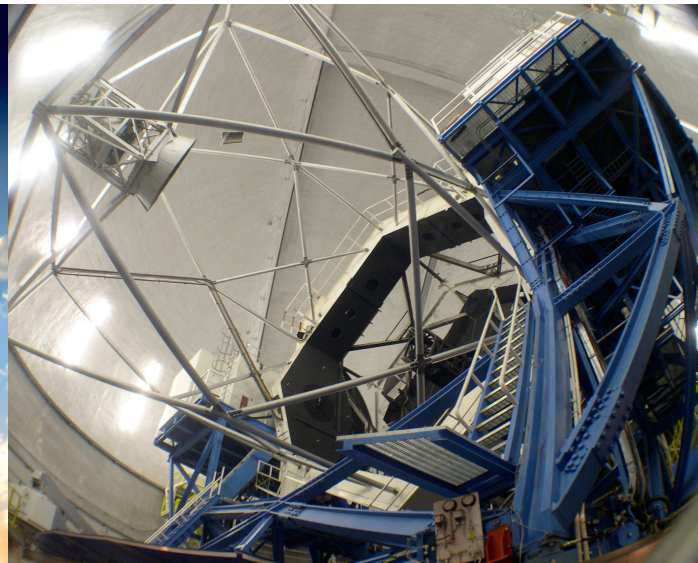


Scaling Astronomy to the Next Decade



Dark Energy, Dark Matter and Baryons

- **Nature of the Universe**

- **Dark energy**

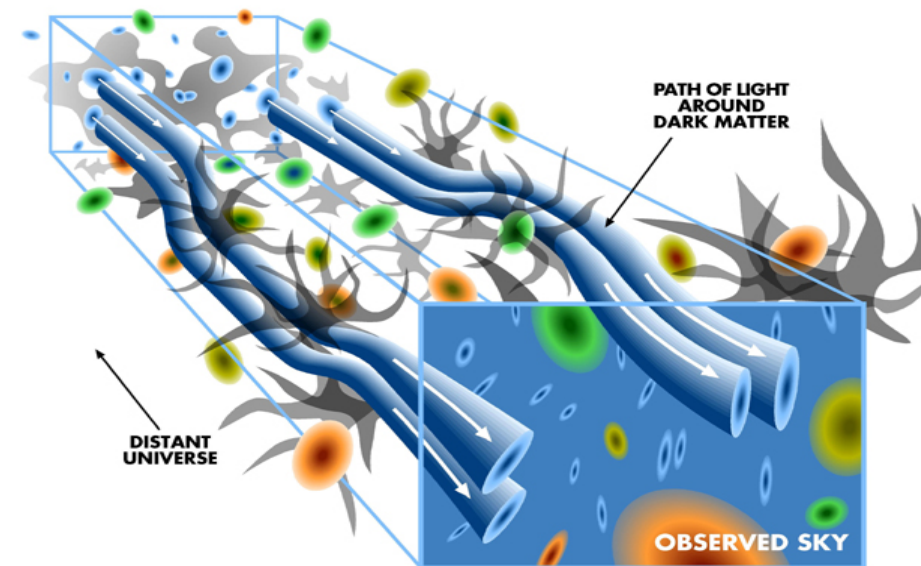
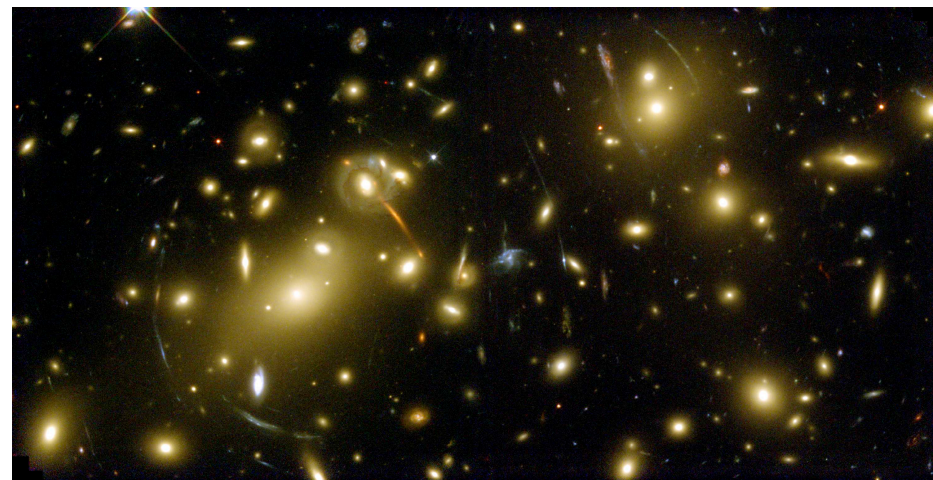
- 73% of energy density
 - Drives acceleration
 - Physics unknown

- **Dark matter**

- 25% of energy density
 - Drives growth of structure
 - Particle unknown

- **Small effects**

- Signals are small
 - Systematics can be large
 - Image distortions measured to



How did the universe at 300,000 years

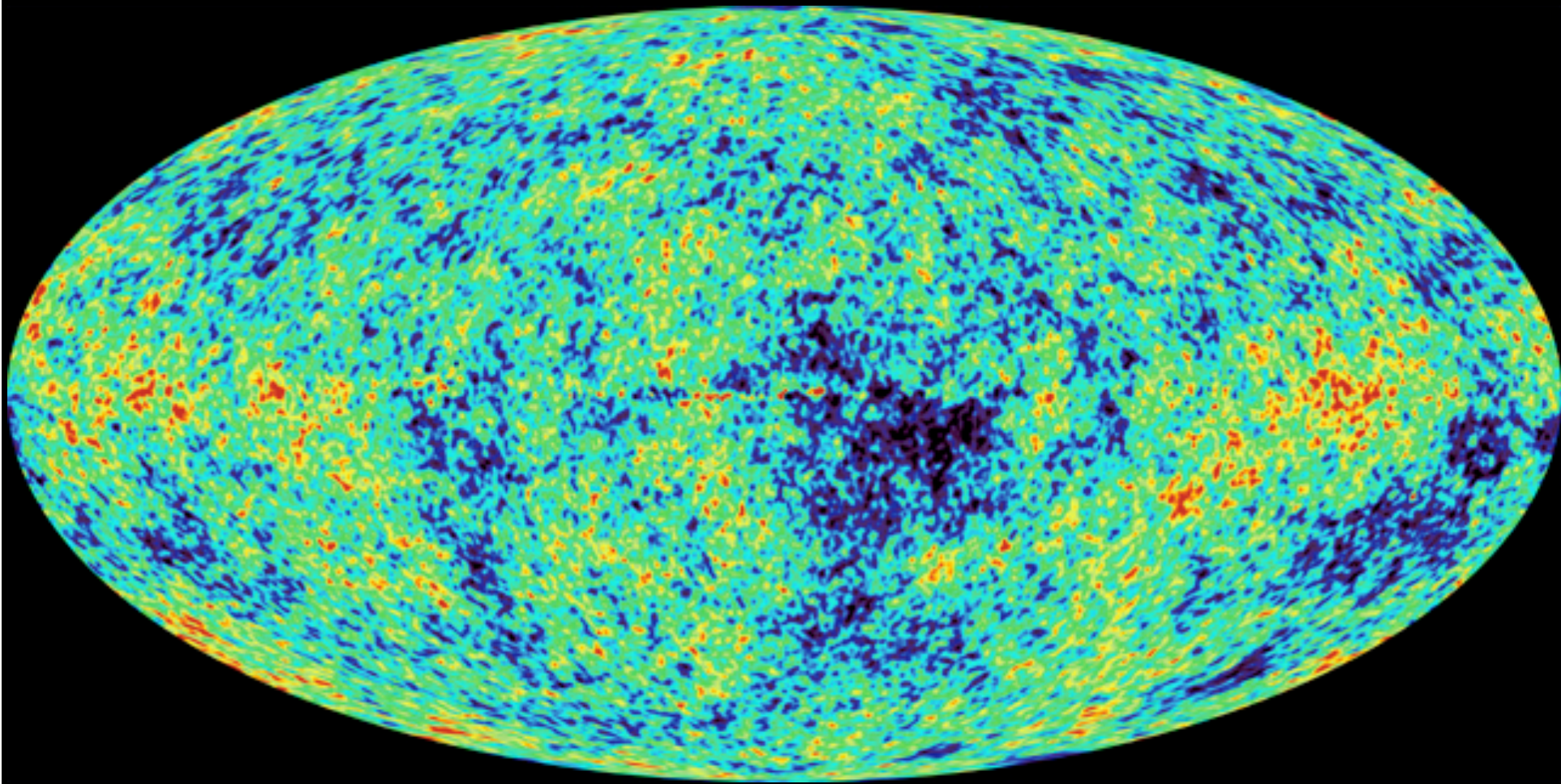
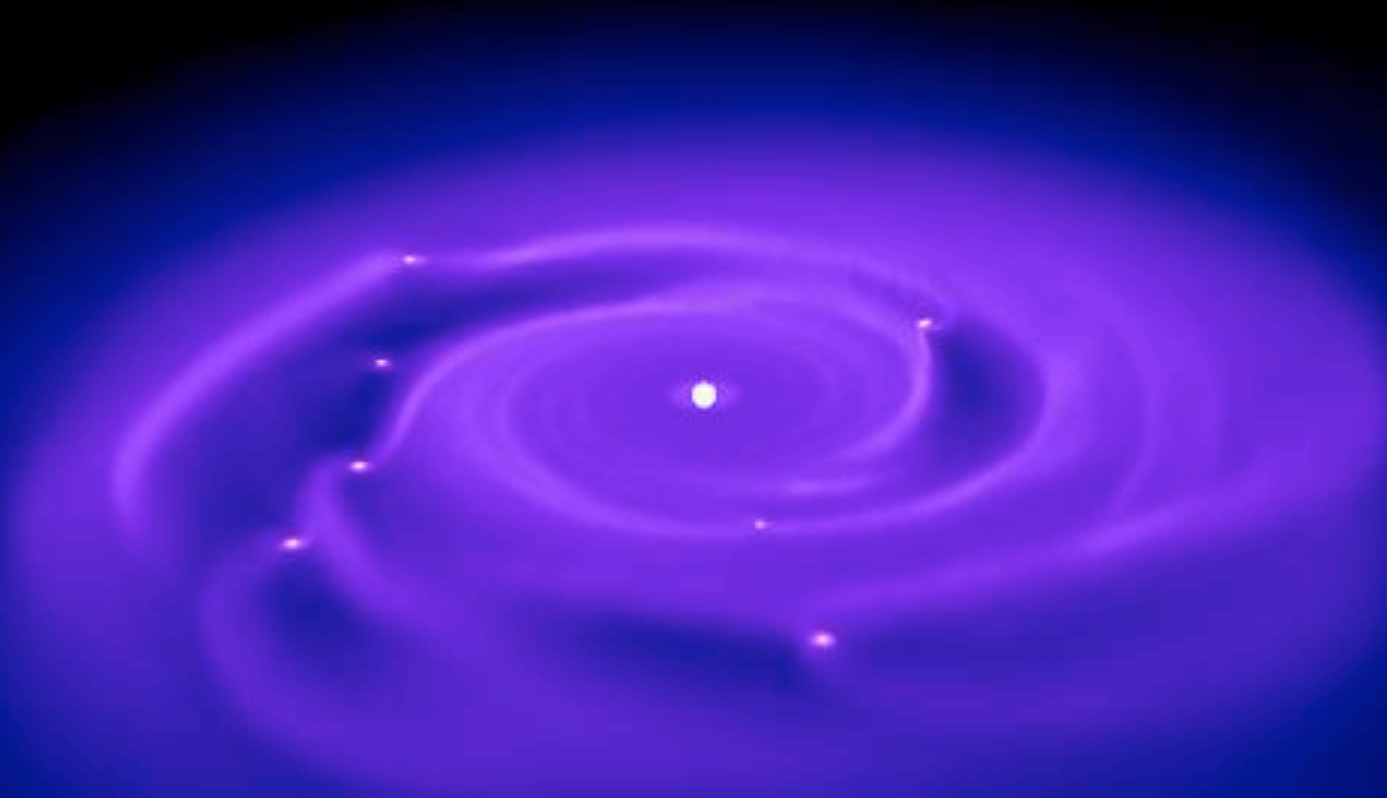


Image courtesy NASA/WMAP

... turn into this?



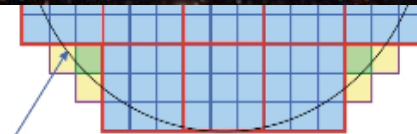
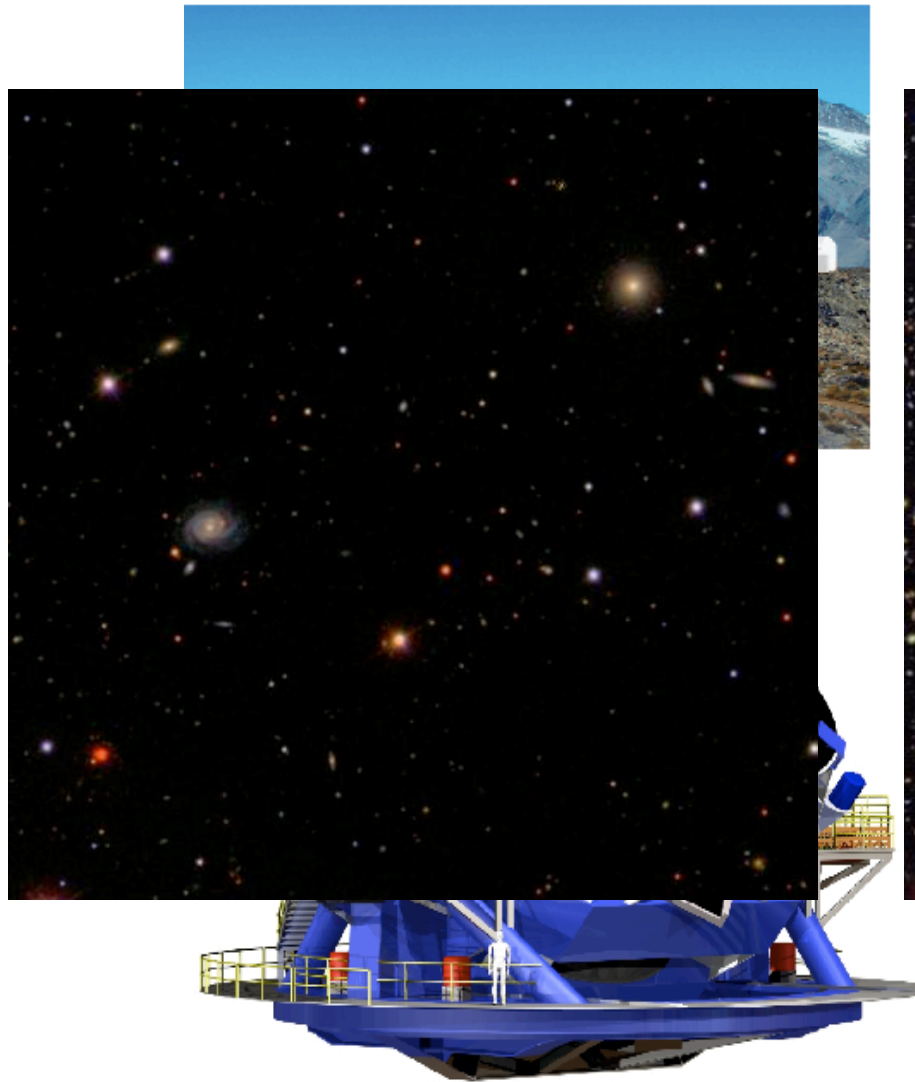
...and this....



...and this?



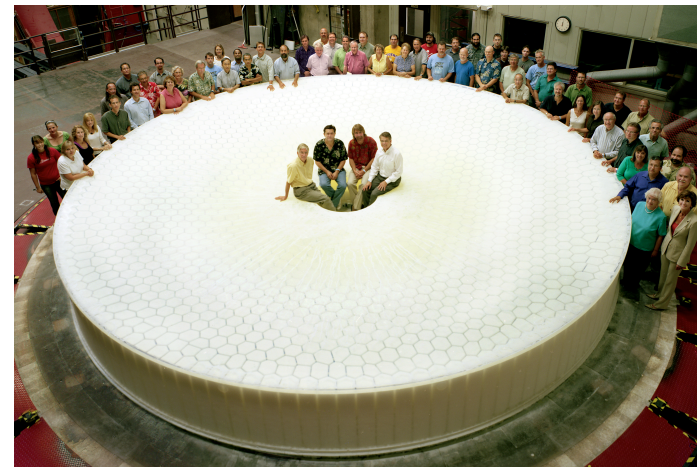
The Large Synoptic Survey Telescope



3.5 degree Field
of View (634 mm diameter)

Challenges from new astronomy

- **Sloan Digital Sky Survey (SDSS)**
 - 7 years of imaging
 - 8000 sq degrees of the sky (1/5th)
 - 200 million stars and galaxies
 - 80 TB raw images
- **LSST data flow**
 - 20,000 sq degrees every 3 nights
 - 40 TB of imaging per night
 - 10^8 sources a night (10^3 “events”)
 - 1000 repeat observations over 10 years
 - 10 Petabytes of catalogs (10 years)
 - 100 PBs of images
 - 5 months to watch 1 year of data (HDTV)
 - Data public as soon as taken



That all seems easy enough....

- **Data Warehouses**
 - **Skyserver** (<http://www.sdss.org>)
 - 4TB SQL server database
 - Flat schema for sources
 - Extended to 10 TB for other applications
 - User definable tables (Casjobs)
 - **LSST**
 - 1 PB (catalogs), 6 PB (images) in 2015
 - Time sampled data (with variable attribute quality)
 - Streaming over the network (or a single disk) not feasible
 - Need to run applications on the data not move the data to the user
- **What techniques should we running on the data**

Science in the coming decade: the temporal sky

- **Finding the unusual**
 - Billion sources a night
 - Nova, supernova, GRBs
 - Instantaneous discovery
- **Finding moving sources**
 - Asteroids and comets
 - Proper motions of stars
- **Mapping the Milky Way**
 - Tidal streams
 - Galactic structure
- **Dark energy and dark matter**
 - Gravitational lensing
 - Slight distortion in shape
 - Trace the nature of dark energy



Science in the coming decade: the temporal sky

- **Finding the unusual**
 - Anomaly detection
 - Dimensionality reduction
 - Weak classifiers
- **Finding moving sources**
 - Tracking algorithms
 - Kalman filters
- **Mapping the Milky Way**
 - Density estimation
 - Clustering (n-tuples)
- **Dark energy and dark matter**
 - Computer vision
 - Model fitting
 - Non-parametric estimation

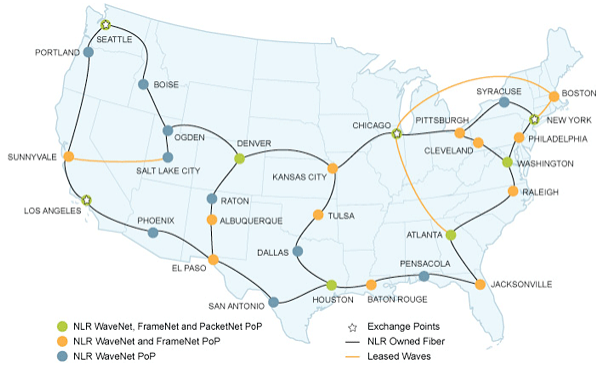


DM Centers optimized for specific functions

- **Base Center**
 - Real-time Processing and Alert Generation
 - Long-term storage (copy 1)
- **Archive Center**
 - Nightly Reprocessing
 - Data Release Processing
 - Long-term Storage (copy 2)
- **Co-located Data Access Centers**
 - Data Access and User Services (load sharing)
 - Provided via the Community Services Subsystem
 - Shares Infrastructure with the Base/Archive Center also at Site
- **System Operations Center**
 - Monitors/manages activities across centers
- **Education & Public Outreach Center**
 - Specialized data access and services for outreach applications
 - Not part of DM System, but several data interfaces



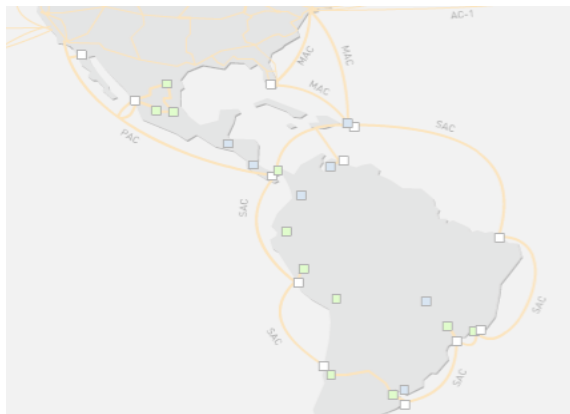
Existing long-haul fiber optic networks, known providers



National Lambda Rail



Internet2 Abilene



Global Crossing

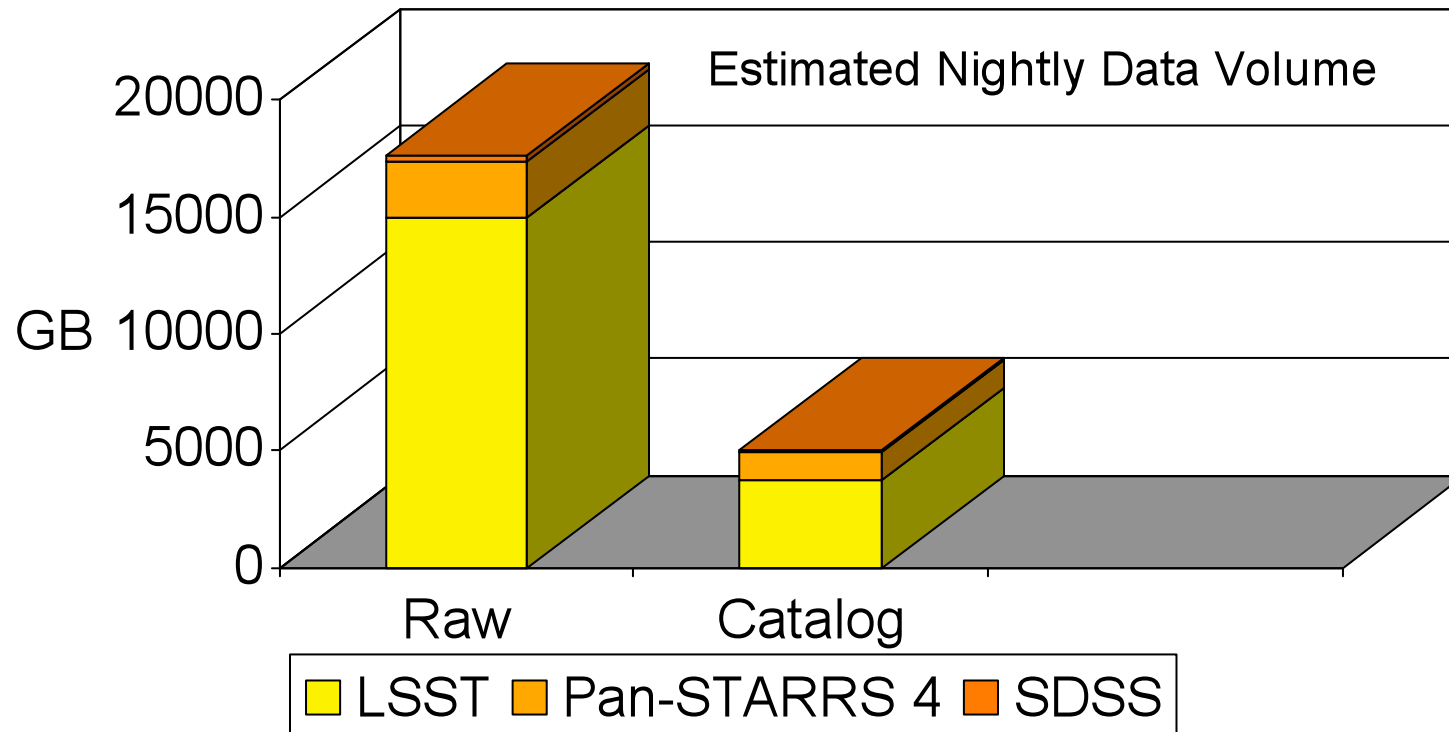


Ampath/WHREN-LILA



telefonicas

DM data volumes/rates are unprecedented in astronomy

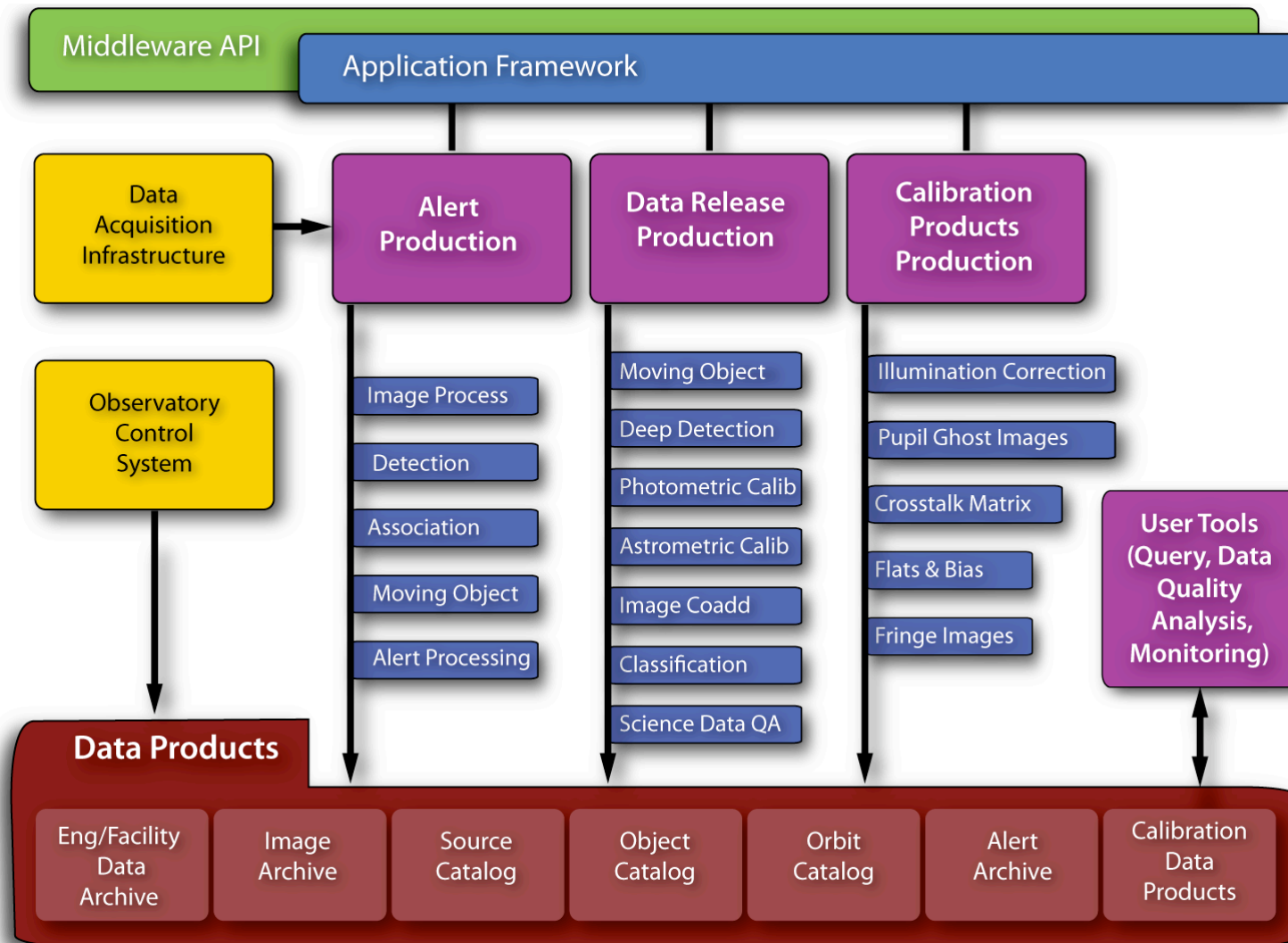


- **LSST will make tens of trillions photometric observations over tens of billions objects**

Level 1, 2, and 3 Data Products

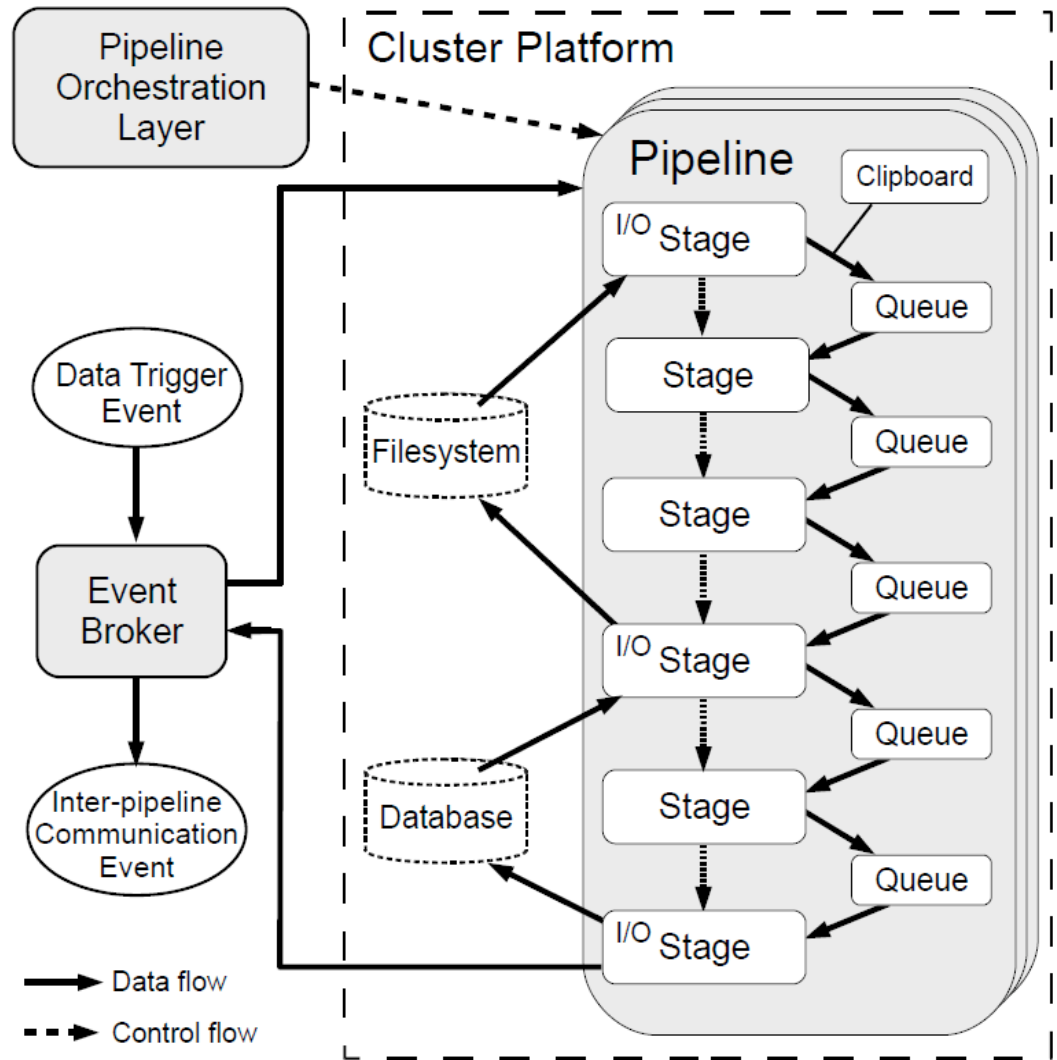
- **Level 1 is generated during nightly processing of the data stream from the camera**
 - Science exposures
 - Transient alerts
 - Updates to Science Database
- **Level 2 is generated as part of a Data Release**
 - Science Database
 - Co-added exposures
- **Level 3 is generated outside the Data Management System, e.g. by science collaborations**
 - Usually federated with L2 science database
 - Supported with tools supplied by LSST DM

DM Applications are framework-based Productions/Pipelines, Data Products, Tools



DM System relies on large-scale computational parallelism

- **With few exceptions, LSST pipeline processing is “embarrassingly parallel”**
 - 3024 parallel image readouts
 - $O(10^8)$ sky tiles
 - $O(10^9)$ objects
- **Computational clusters are well matched to the available parallelism**
 - 5000 cores at Base
 - 12000 (yr1) – 33000 (yr10) cores at Archive
- **Middleware implements flexible pipeline/production model of parallelism**



Data Challenge 4 – Database performance and Data Access (complete December 2011)

- **Goals**

- **Demonstrate database scale-up**
- **Alert generation and distribution**
- **Demonstrate fault tolerance**
- **Alert Production throughput 20% of full operational DMS**
- **Tools for external data access and visualization**

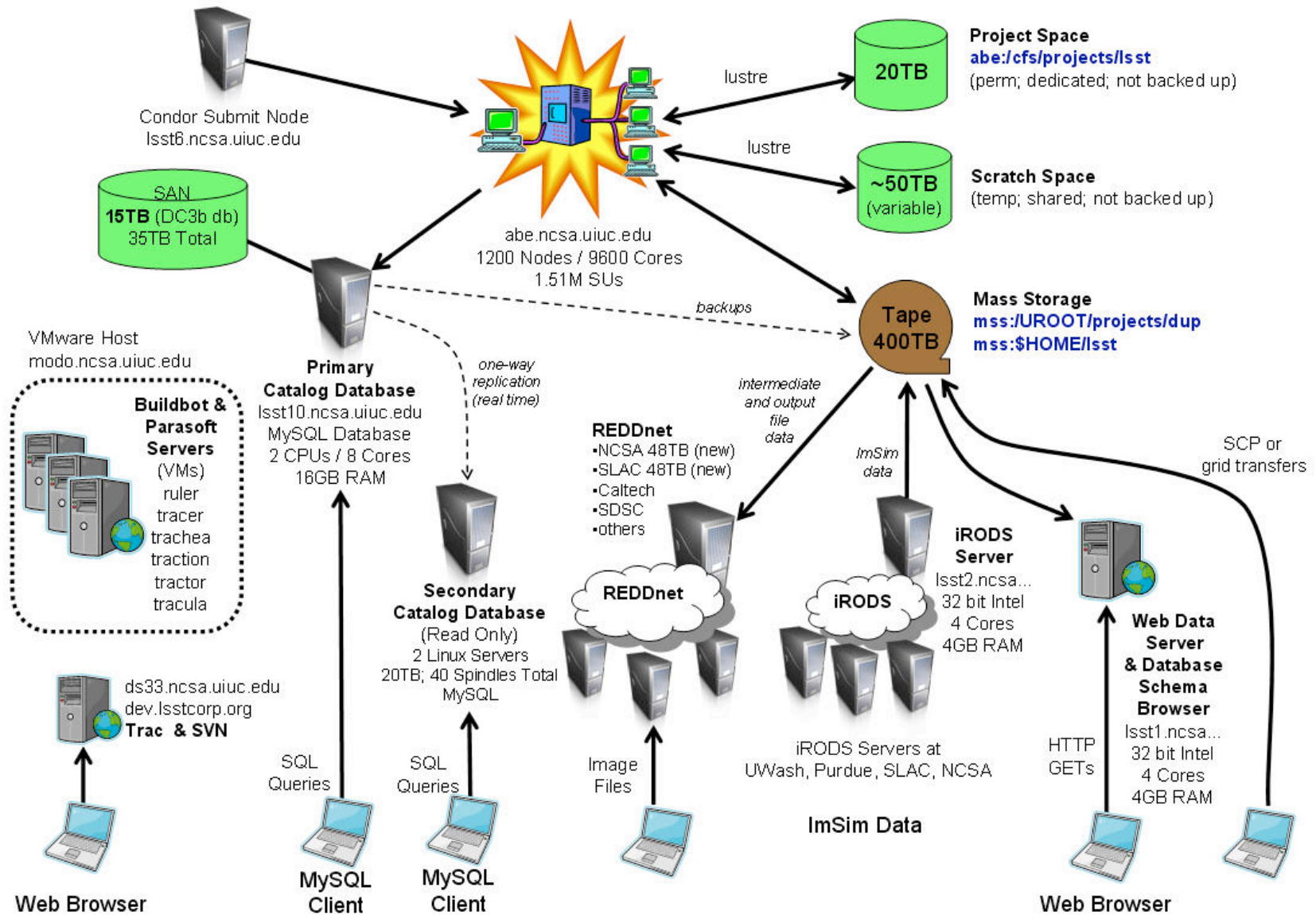
- **Execution**

- **2 Phases to scale up performance: stepping up in amount of data processed**

- **Need significant infrastructure allocations**

- **~5 TFLOPS Computing resources**
- **~700 TB Storage**





DM System is widely distributed

