# Big Data: Storing, Moving, Mining, Visualizing

Things HiPACC Could be Leveraging

# Topics

- Moving: CENIC @100G
- Storing: Data Oasis
- Mining: Gordon, Comet
- Visualizing: yt
- Sharing: SeedMe.org

- 3,800+ miles of optical fiber
- **Members in all 58 counties connect via fiber-optic cable or leased circuits from telecom carriers**
  - **Nearly *10,000* sites connect to CENIC**
    - *10,000,000+* Californians use CENIC each day
    - Governed by members on the segmental level

Corning
Sacramento
Oakland
San Francisco
Palo Alto
San Jose
Merced
Fresno
Soledad
San Luis Obispo
Bakersfield
Los Angeles
Riverside
Palm Desert
Tustin
San Diego
El Centro

cenic

it²

# CENIC, Internet2, ESNet now connected at 100G

# UCSD will be connected to CENIC and ESNet at 100G by end of 2014
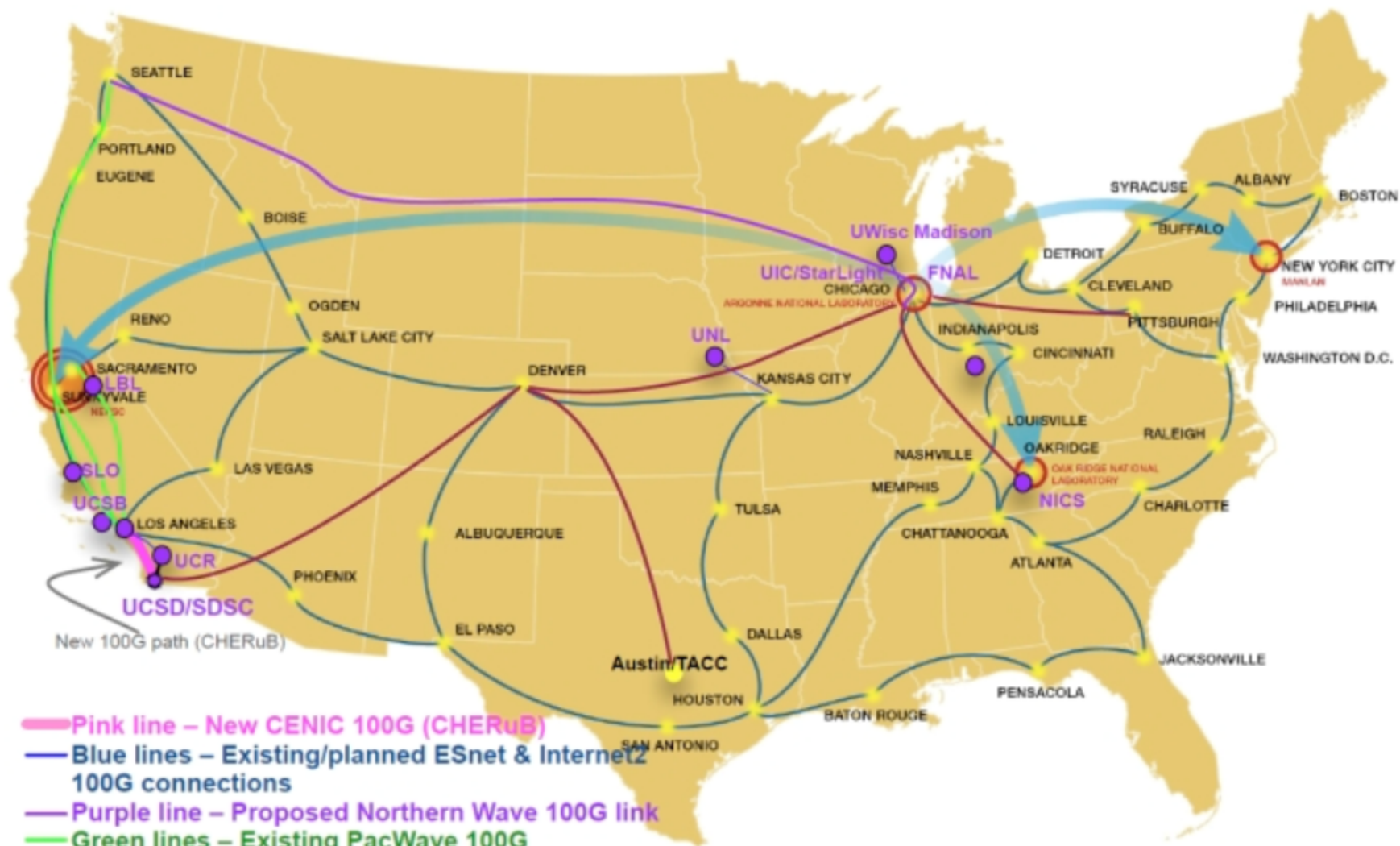


NSF CC-NIE grant
PI: M. Norman

**CHERuB**

Configurable, High-speed, Extensible Research Bandwidth

CHERuB is bringing 100Gbps (Gigabit per second) data network connectivity to UCSD to accommodate big data and network research.
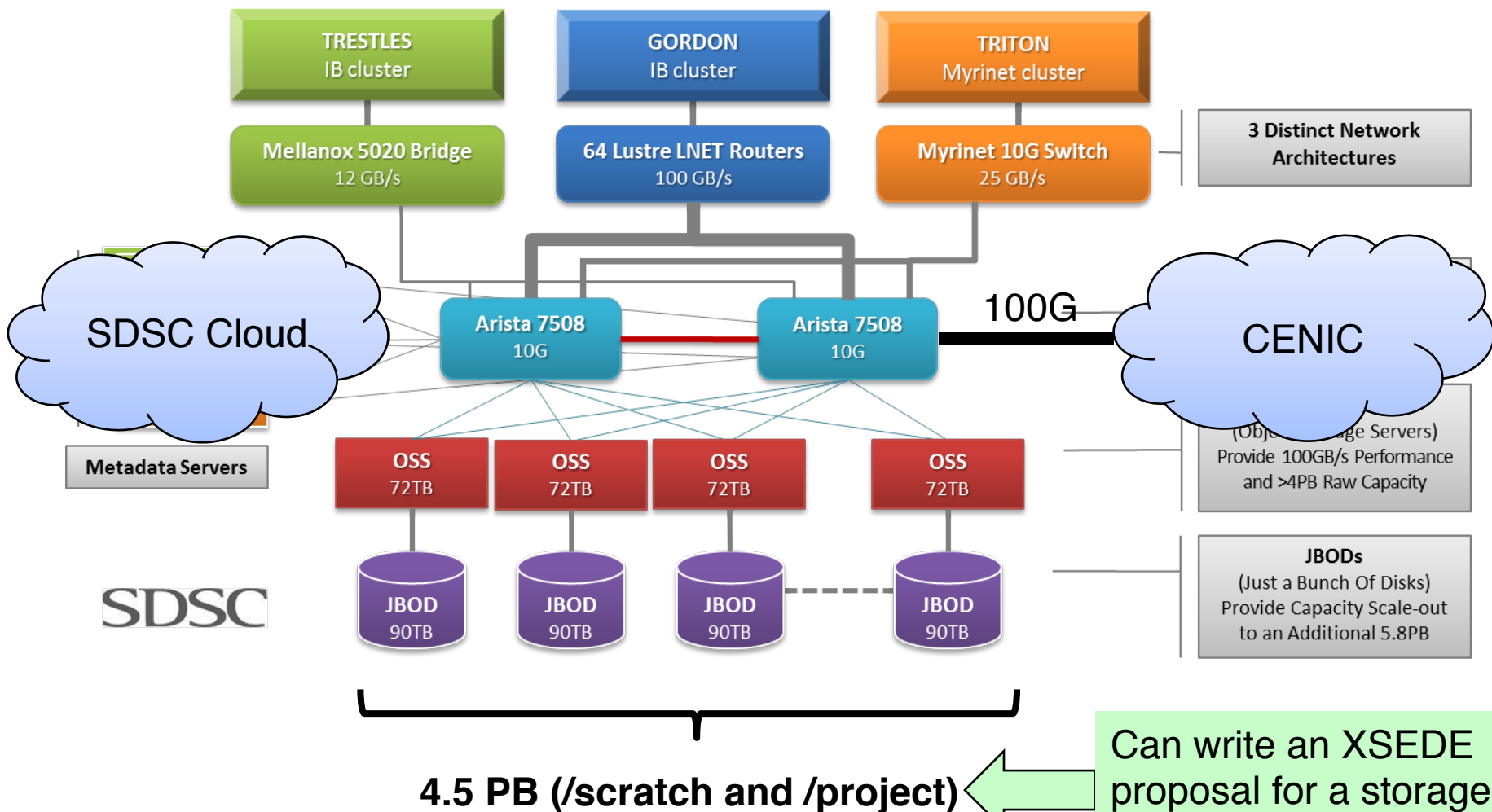
# Other UC campuses with 100G connectivity

| Campus | Connectivity | Date |
|--------|--------------|------|
| UCLA | 100G | Now |
| UCSC | 100G | 2014 |
| UCD | 100G | 2015 |
| UCB | 100G | 2015 |

Pink line – New CENIC 100G (CHERuB)
Blue lines – Existing/planned ESnet & Internet2 100G connections
Purple line – Proposed Northern Wave 100G link
Green lines – Existing PacWave 100G
Maroon lines – XSEDE 10G network
Black lines – Other existing 10-40G

# Data Oasis Heterogeneous Architecture



**4.5 PB (/scratch and /project)**

Can write an XSEDE proposal for a storage allocation; e.g., /oasis/proj/hipacc

# SDSC Cloud
## http://cloud.sdsc.edu

### OpenStack components

- **SWIFT object store**
- **NOVA compute**



### Functionality

- **Easily locating research data (web accessible)**
- **Storage facilities that can hold huge datasets**
- **New web, compute servers with few clicks**
- **Publishing and sharing data**

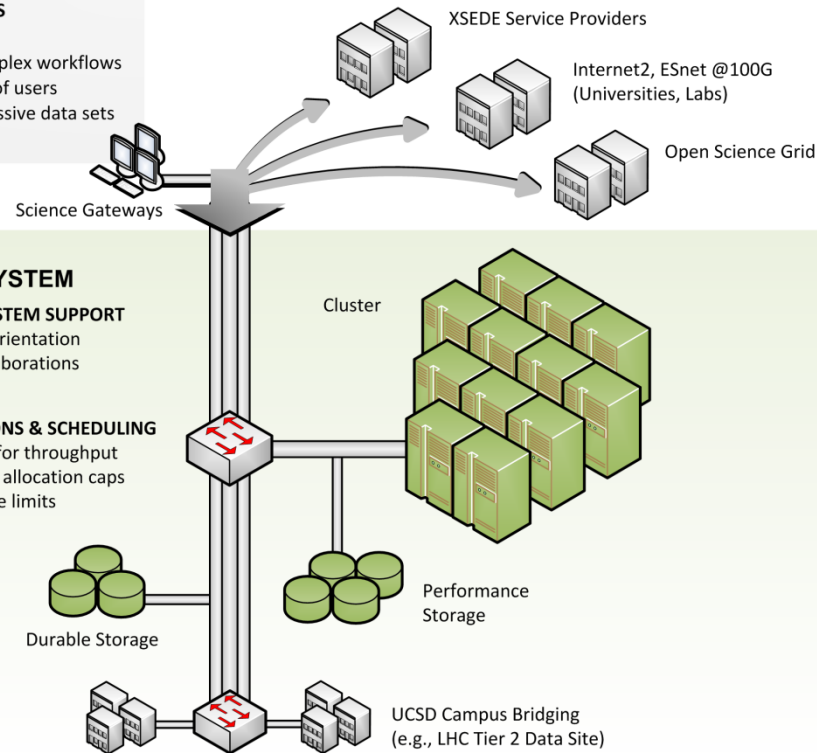# *Data Mining: Gordon Data-Intensive Supercomputer (XSEDE resource)*

- **First HPC system with flash SSD (300 TB)**

- **All I/O channels "maxed out" to accommodate Big Data movement**

- **Can "spin up" a SSD based Hadoop cluster from the batch queue individual users**

- **Fore-runner to Catalyst system at LLNL**



**SDSC** SAN DIEGO SUPERCOMPUTER CENTER

# *Comet supercomputer (2015) will support Big Data workflows for all fields of scholarship*



**CHALLENGES OUR PROPOSAL ADDRESSES**
- ✓ Attract new users and communities
- ✓ Support diverse applications with complex workflows
- ✓ Ensure responsiveness for thousands of users
- ✓ Transfer, store, analyze, and share massive data sets
- ✓ Integrate with XSEDE

Science Gateways

XSEDE Service Providers

Internet2, ESnet @100G
(Universities, Labs)

Open Science Grid

**COMET COMPUTE SYSTEM**

**Cluster architecture**
Fast standard nodes
Large-memory nodes
GPU-accelerated nodes
FDR InfiniBand

**Storage architecture**
Performance Storage
Durable Storage

**Software**
Science Gateways
Rich base of installed apps
Virtualization

**USER & SYSTEM SUPPORT**
New user orientation
XSEDE collaborations
FutureGrid

**ALLOCATIONS & SCHEDULING**
Optimized for throughput
Per-project allocation caps
Per-job core limits

Cluster

Performance Storage

Durable Storage

UCSD Campus Bridging
(e.g., LHC Tier 2 Data Site)

2 PF compute
7 PB Lustre PFS
4.5 PB 2$^{nd}$ copy PFS
(old Data Oasis)

COMET
SDSC

# *Sharing: SeedMe.org*



- A new SaaS web service under development at SDSC

- Think of it as Flickr+YouTube for scientists