# The NSF TeraGrid:
# A Pre-Production Update

**2nd Large Scale Cluster Computing Workshop**

**FNAL**

**21 Oct 2002**

**Rémy Evard, evard@mcs.anl.gov**

**TeraGrid Site Lead**

**Argonne National Laboratory**

# A Pre-Production Introspective

- **Overview of The TeraGrid**
  - For more information:
    - www.teragrid.org
      - See particularly the "TeraGrid primer".
  - Funded by the National Science Foundation
  - Participants:
    - NCSA
    - SDSC
    - ANL
    - Caltech
    - PSC, starting in October 2002

- **Grid Project Pondering**
  - Issues encountered while trying to build a complex, production grid.

# Motivation for TeraGrid

- **The Changing Face of Science**
  - Technology Drivers
  - Discipline Drivers
  - Need for Distributed Infrastructure

- **The NSF's Cyberinfrastructure**
  - "provide an integrated, high-end system of computing, data facilities, connectivity, software, services, and sensors that …"
  - "enables all scientists and engineers to work on advanced research problems that would not otherwise be solvable"
    - Peter Freeman, NSF

- **Thus the Terascale program**

- **A key point for this workshop:**
  - TeraGrid is meant to be an infrastructure supporting <u>many</u> scientific disciplines and applications.
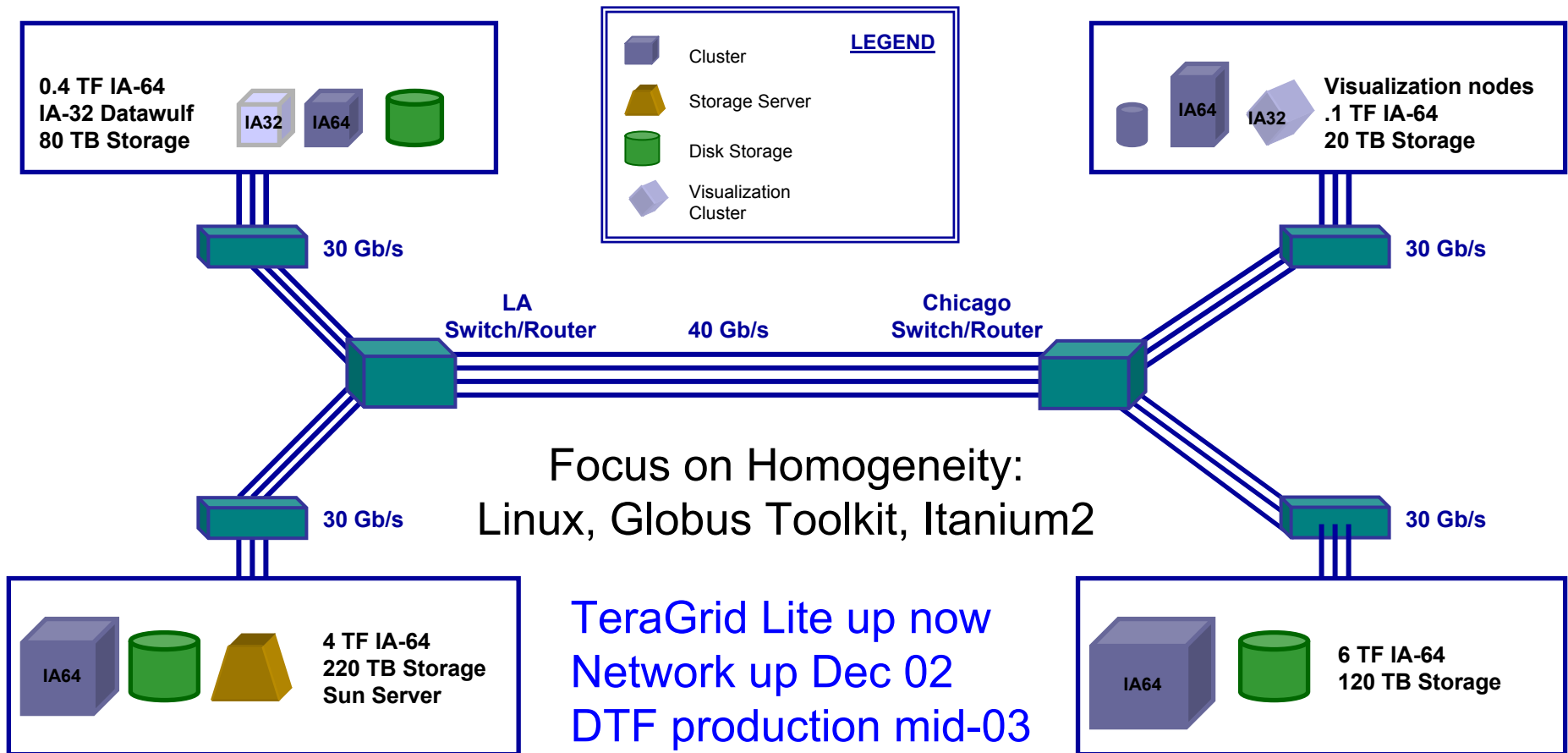
# Historical Context

- **Terascale funding arrived in FY00**
- **Three competitions so far:**
  - FY00 – Terascale Computing System
    - Funded PSC's EV68 6TF Alpha Cluster
  - FY01 – Distributed Terascale Facility (DTF)
    - Initial TeraGrid Project
  - FY02 – Extensible Terascale Facility  (ETF)
    - Expansion of the TeraGrid
- **An additional competition is now underway for community participation in ETF**

# Distributed Terascale Facility (DTF) TeraGrid

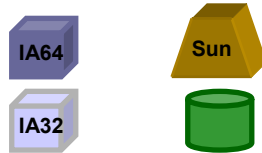**Caltech:** Data collection analysis

**ANL:** Visualization

0.4 TF IA-64
IA-32 Datawulf
80 TB Storage

IA32  IA64

**LEGEND**

Cluster

Storage Server

Disk Storage

Visualization
Cluster

IA64  IA32  Visualization nodes
.1 TF IA-64
20 TB Storage

30 Gb/s

30 Gb/s

LA
Switch/Router          40 Gb/s          Chicago
Switch/Router

30 Gb/s

Focus on Homogeneity:
Linux, Globus Toolkit, Itanium2

30 Gb/s

IA64  4 TF IA-64
220 TB Storage
Sun Server

TeraGrid Lite up now
Network up Dec 02
DTF production mid-03

IA64  6 TF IA-64
120 TB Storage

**SDSC**: Data Intensive

**NCSA**: Compute Intensive

# Extensible TeraGrid Facility

**Caltech**: Data collection analysis

IA64
Sun
IA32

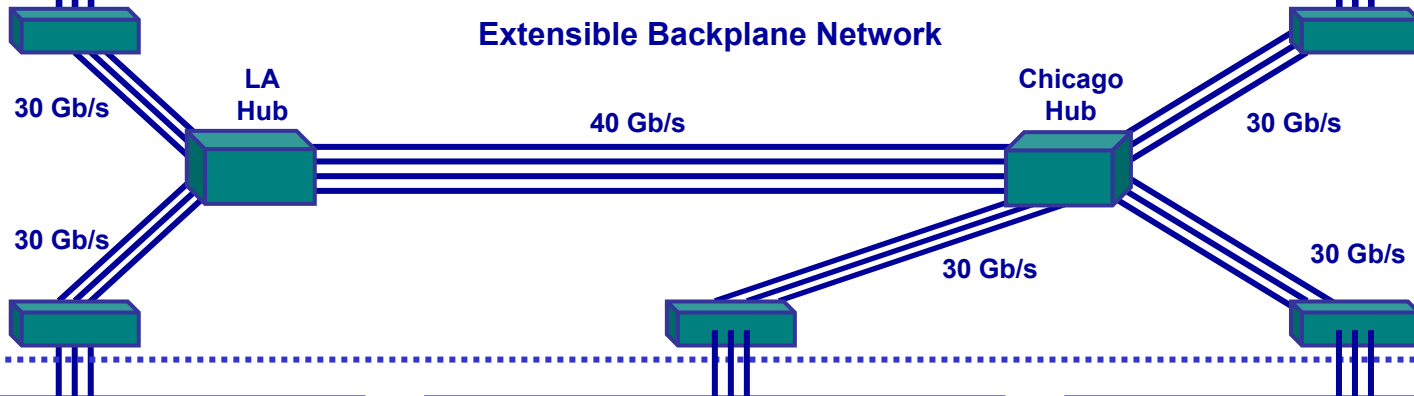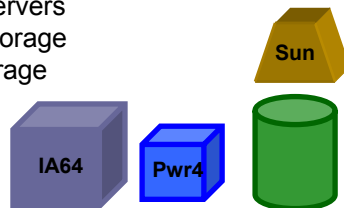0.4 TF IA-64
IA32 Datawulf
80 TB Storage

**LEGEND**

| | | | |
|---|---|---|---|
| Cluster | | Visualization Cluster | |
| Storage Server | | Shared Memory | |
| Disk Storage | | Backplane Router | |

**ANL**: Visualization

IA64
IA32

.5 TF IA-64
96 Visualization nodes
20 TB Storage

**Extensible Backplane Network**

LA Hub

Chicago Hub

30 Gb/s

40 Gb/s
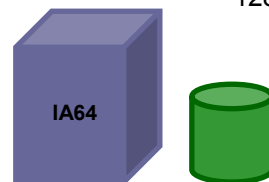
30 Gb/s

30 Gb/s

30 Gb/s

30 Gb/s

4 TF IA-64
DB2, Oracle Servers
500 TB Disk Storage
6 PB Tape Storage
1.1 TF Power4
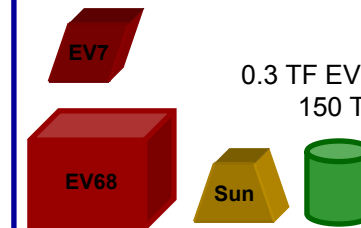
IA64
Pwr4
Sun

**SDSC**: Data-Intensive

10 TF IA-64
128 large memory nodes
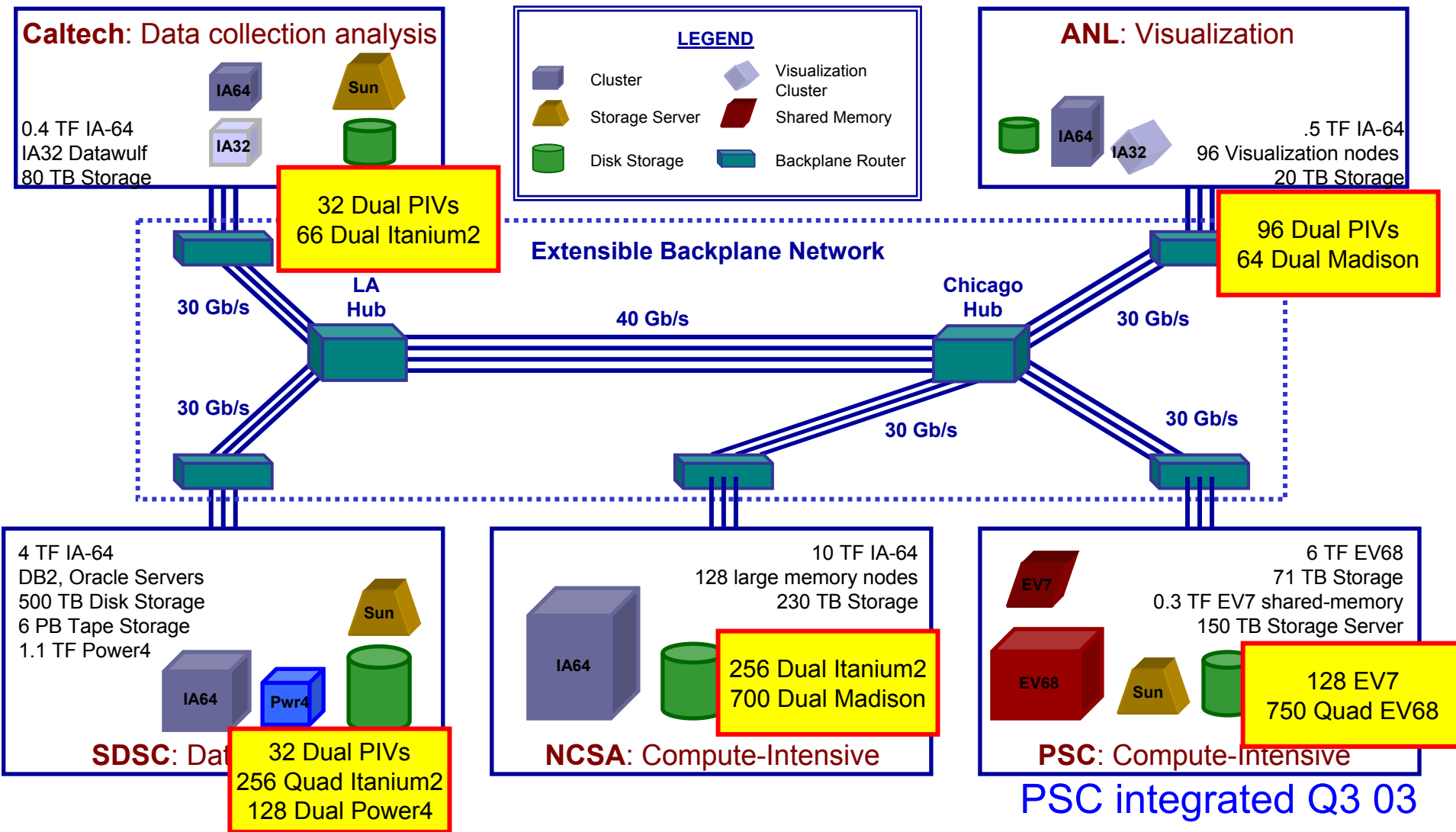230 TB Storage

IA64

**NCSA**: Compute-Intensive

6 TF EV68
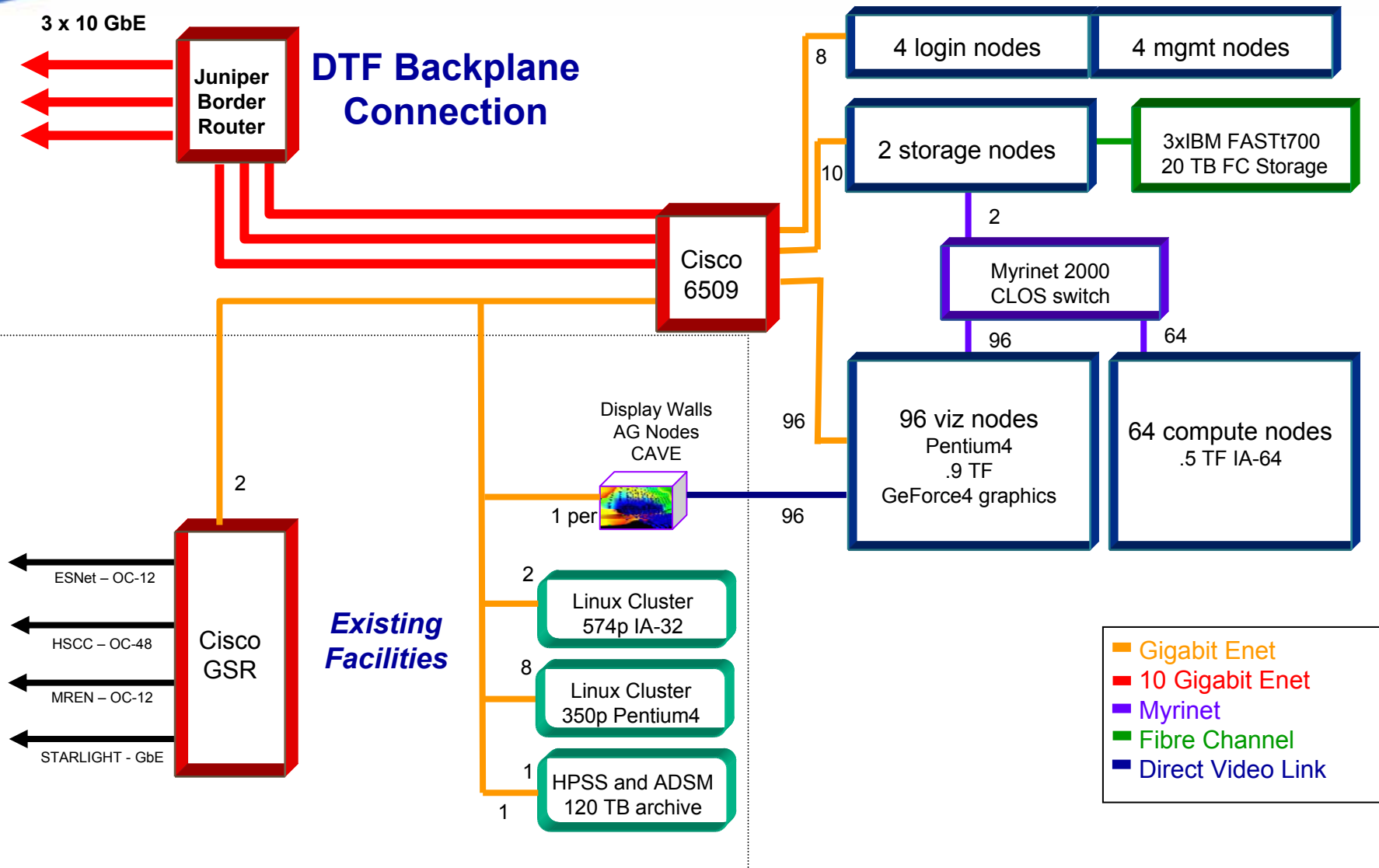71 TB Storage
0.3 TF EV7 shared-memory
150 TB Storage Server

EV7
EV68
Sun

**PSC**: Compute-Intensive

PSC integrated Q3 03

# Extensible TeraGrid Facility

# Argonne ETF Cluster Schematic



**3 x 10 GbE**

Juniper Border Router

**DTF Backplane Connection**

4 login nodes — 4 mgmt nodes

8

2 storage nodes — 3xIBM FASTt700 20 TB FC Storage

10

2

Cisco 6509

Myrinet 2000 CLOS switch

96

64

96

96 viz nodes
Pentium4
.9 TF
GeForce4 graphics

64 compute nodes
.5 TF IA-64

Display Walls
AG Nodes
CAVE

1 per

96

*Existing Facilities*

Cisco GSR

ESNet – OC-12

HSCC – OC-48

MREN – OC-12

STARLIGHT - GbE

2

2

Linux Cluster
574p IA-32

8

Linux Cluster
350p Pentium4

1

HPSS and ADSM
120 TB archive

1

**Legend:**
- Gigabit Enet
- 10 Gigabit Enet
- Myrinet
- Fibre Channel
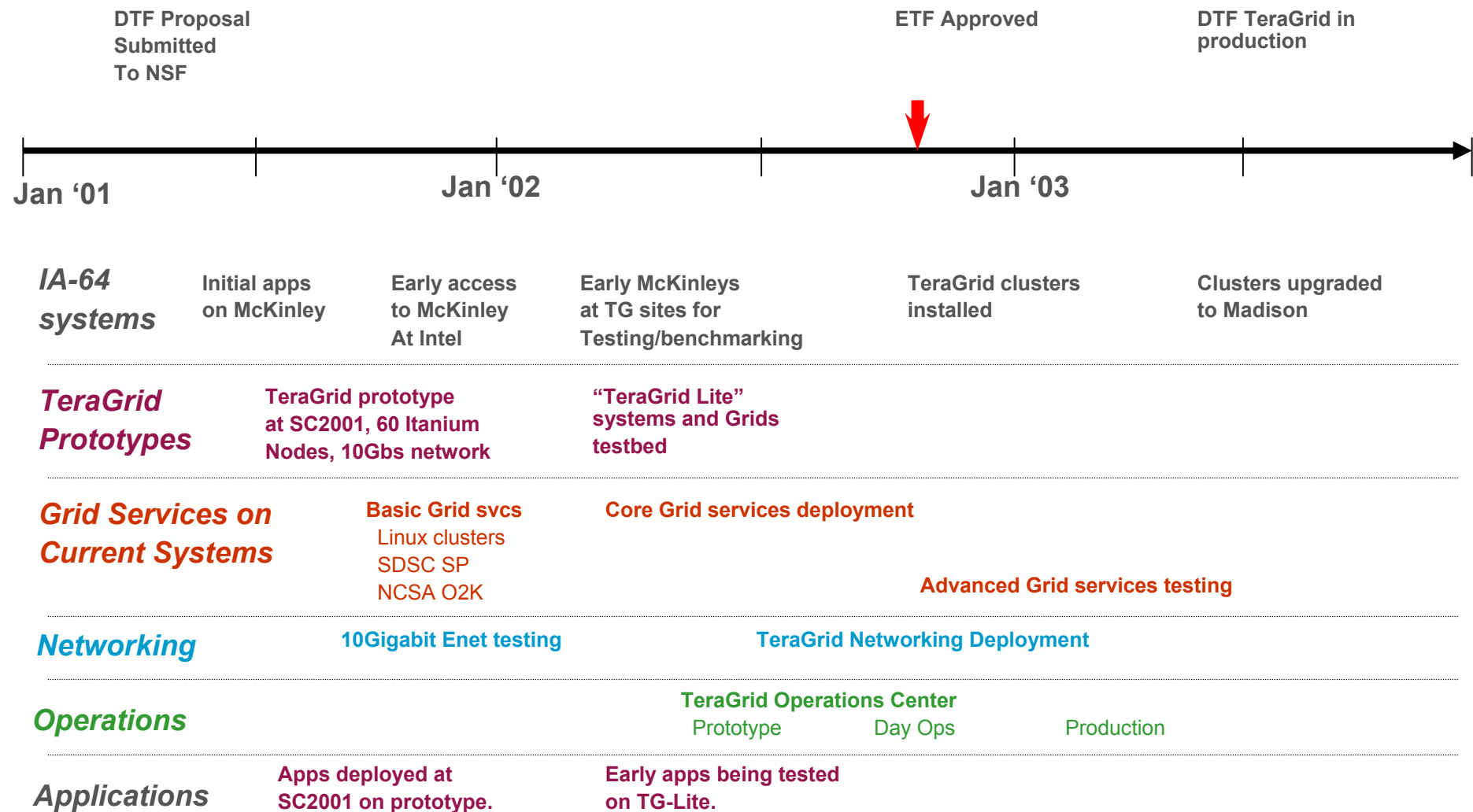- Direct Video Link

Charlie Catlett

# TeraGrid Objectives

- **Create significant enhancement in capability**
  - Beyond capacity, provide basis for exploring new application capabilities
- **Deploy a balanced, distributed system**
  - Not a "distributed computer" but rather
  - A distributed "system" using Grid technologies
    - Computing and data management
    - Visualization and scientific application analysis
- **Define an open and extensible infrastructure**
  - An "enabling cyberinfrastructure" for scientific research
  - Extensible beyond the original four sites

# Where We Are



**Timeline:** Jan '01 — Jan '02 — Jan '03

- DTF Proposal Submitted To NSF
- ETF Approved
- DTF TeraGrid in production

| Category | | | | |
|---|---|---|---|---|
| **IA-64 systems** | Initial apps on McKinley | Early access to McKinley At Intel | Early McKinleys at TG sites for Testing/benchmarking | TeraGrid clusters installed |
| | | | | Clusters upgraded to Madison |
| **TeraGrid Prototypes** | TeraGrid prototype at SC2001, 60 Itanium Nodes, 10Gbs network | | "TeraGrid Lite" systems and Grids testbed | |
| **Grid Services on Current Systems** | Basic Grid svcs Linux clusters SDSC SP NCSA O2K | | Core Grid services deployment | |
| | | | Advanced Grid services testing | |
| **Networking** | 10Gigabit Enet testing | | TeraGrid Networking Deployment | |
| **Operations** | | | TeraGrid Operations Center Prototype    Day Ops    Production | |
| **Applications** | Apps deployed at SC2001 on prototype. | | Early apps being tested on TG-Lite. | |

# Challenges and Issues

- **Technology and Infrastructure**
  - Networking
  - Computing and Grids
  - Others (not covered in this talk):
    - Data
    - Visualization
    - Operation
    - …

- **Social Dynamics**

- **To Be Clear…**
  - While the following slides discuss problems and issues in the spirit of this workshop, the TG project is making appropriate progress and is on target for achieving milestones.

# Networking Goals

- **Support high bandwidth between sites**
  - Remote access to large data stores
  - Large data transfers
  - Inter-cluster communication
- **Support extensibility to N sites**
  - 4 <= N <= 20 (?)
- **Operate in production, but support network experiments.**
- **Isolate the clusters from network faults and vice versa.**
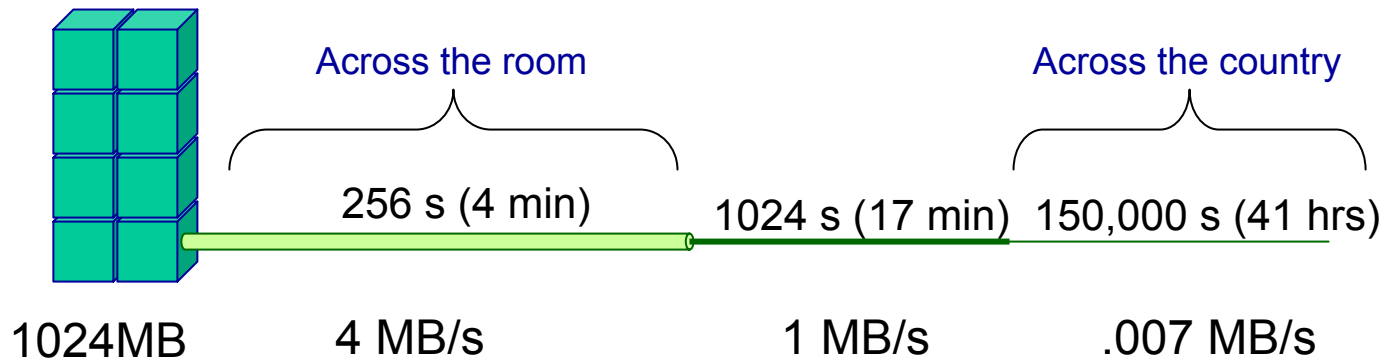
# NSFNET 56 Kb/s Site Architecture

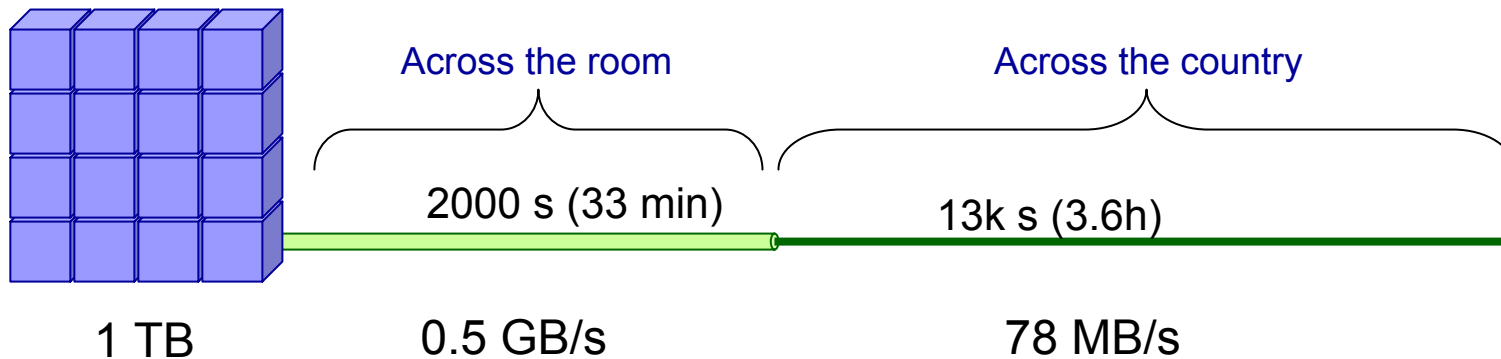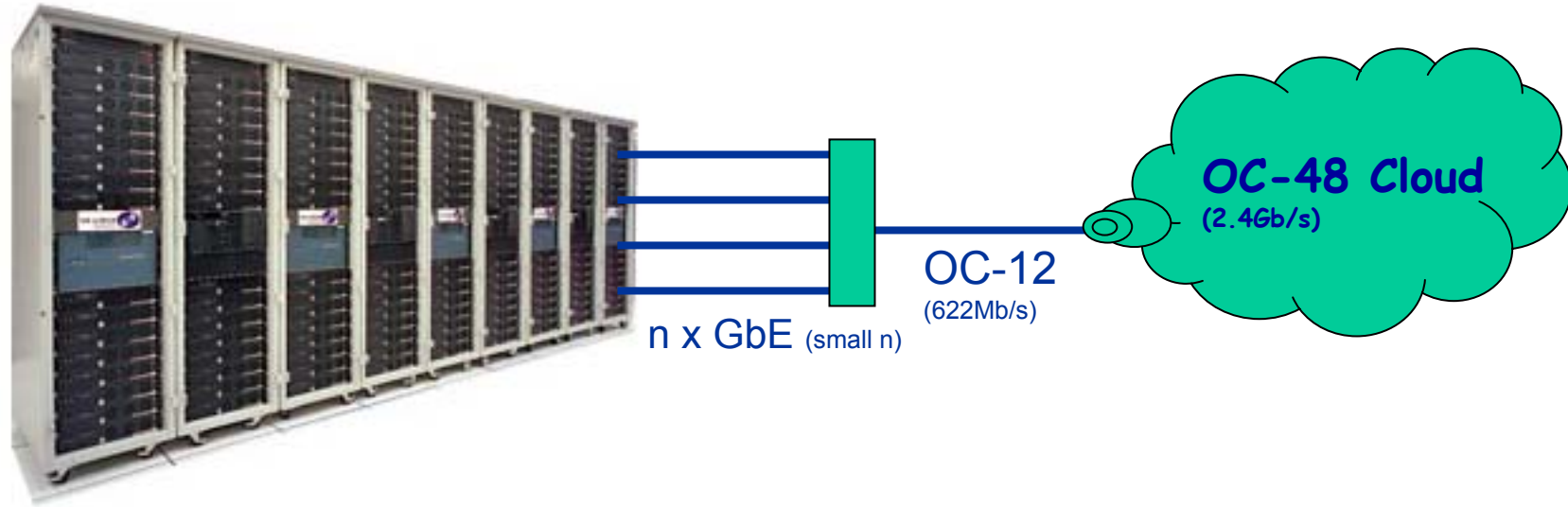Bandwidth in terms of **burst** data transfer and user **wait time**.
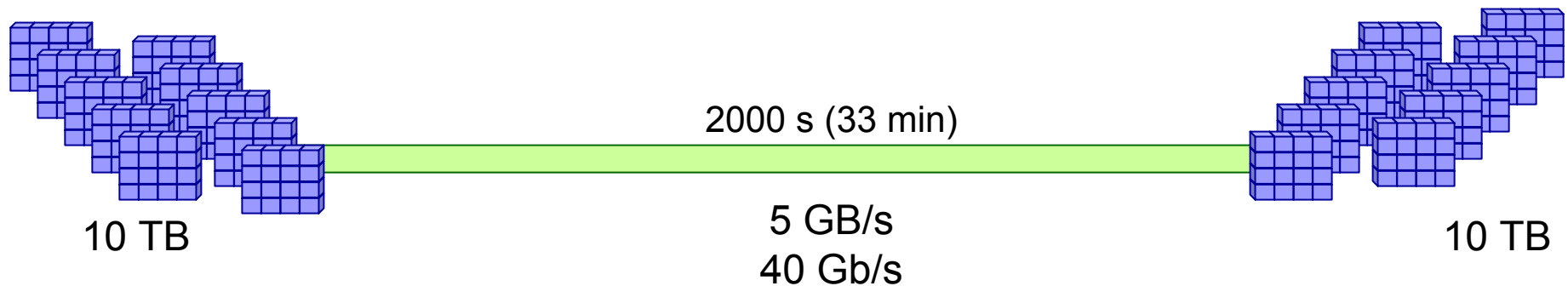
**VAX**

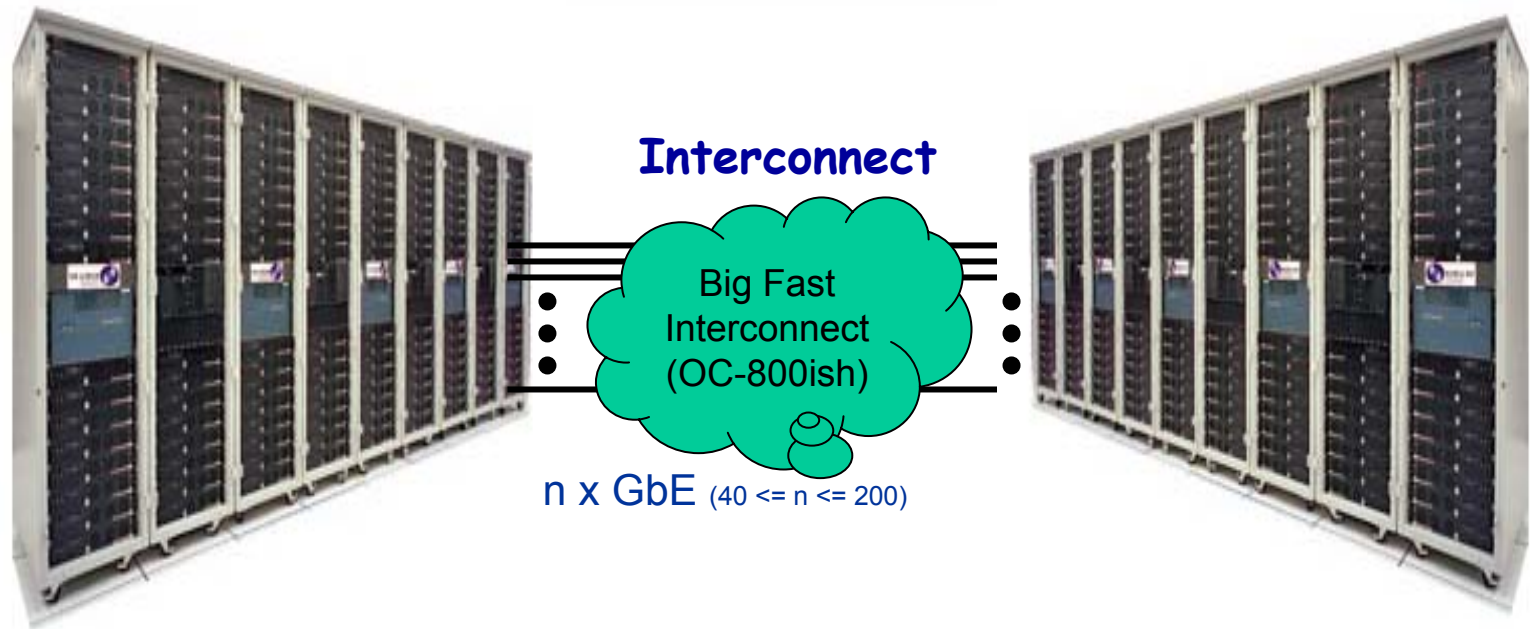**Fuzzball**

Across the room

Across the country

256 s (4 min)

1024 s (17 min)  150,000 s (41 hrs)

1024MB

4 MB/s

1 MB/s

.007 MB/s

# 2002 Cluster-WAN Architecture

**OC-48 Cloud**
(2.4Gb/s)

OC-12
(622Mb/s)

n x GbE (small n)

1 TB

Across the room

2000 s (33 min)

0.5 GB/s

Across the country

13k s (3.6h)

78 MB/s

# To Build a Distributed Terascale Cluster…

## Interconnect

Big Fast
Interconnect
(OC-800ish)

n x GbE (40 <= n <= 200)
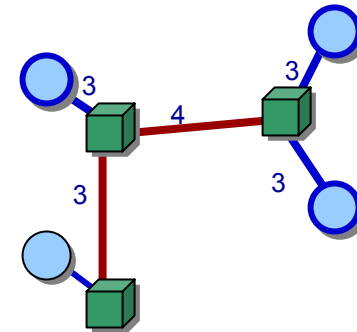
2000 s (33 min)

10 TB

5 GB/s
40 Gb/s

10 TB

# TeraGrid Interconnect: Qwest Partnership

**Phase 0 (June 2002)**

**Phase 1 (November 2002)**

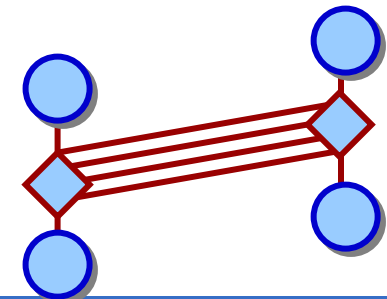**Physical**

# denotes λ count

Pasadena

Argonne

LA

1

Chicago

2

1

2

Urbana

La Jolla

San Diego

3

4

3

3

3

**Light Paths (Logical)**

Caltech/JPL

ANL

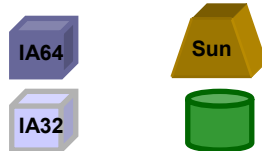SDSC/UCSD

NCSA/UIUC
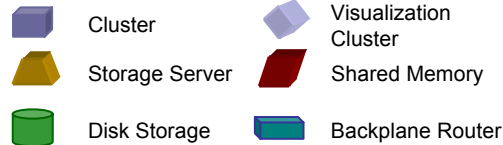
Original: Lambda Mesh

Extensible: Central Hubs

# Extensible TeraGrid Facility



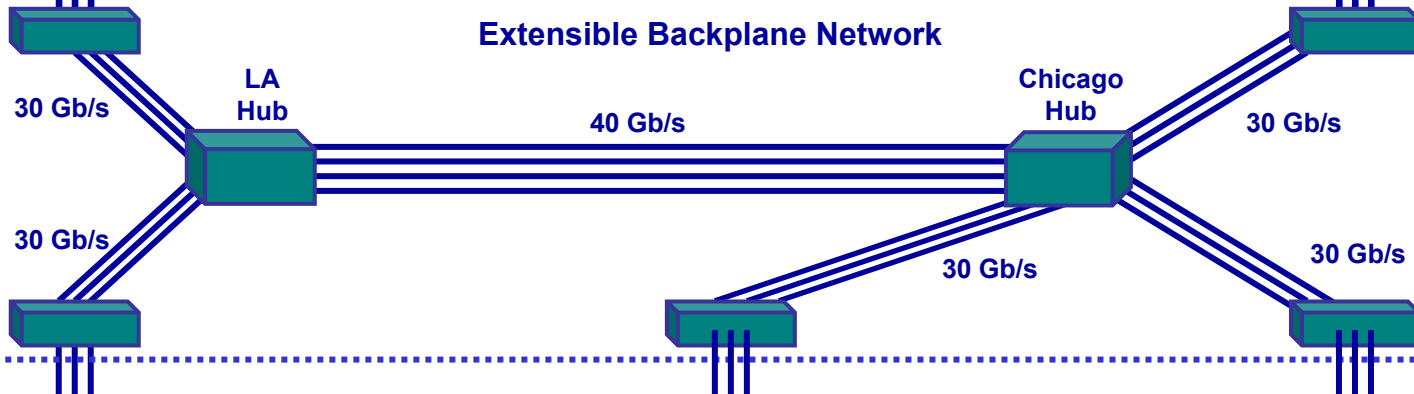**Caltech**: Data collection analysis

0.4 TF IA-64
IA32 Datawulf
80 TB Storage

**LEGEND**

- Cluster
- Visualization Cluster
- Storage Server
- Shared Memory
- Disk Storage
- Backplane Router

**ANL**: Visualization

.5 TF IA-64
96 Visualization nodes
20 TB Storage

**Extensible Backplane Network**

LA Hub

Chicago Hub

30 Gb/s

40 Gb/s

30 Gb/s

30 Gb/s

30 Gb/s

30 Gb/s
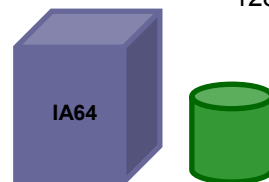
4 TF IA-64
DB2, Oracle Servers
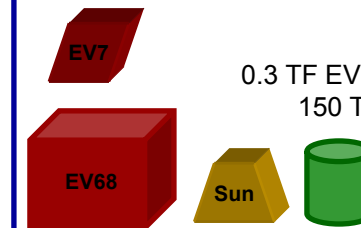500 TB Disk Storage
6 PB Tape Storage
1.1 TF Power4

**SDSC**: Data-Intensive

10 TF IA-64
128 large memory nodes
230 TB Storage

**NCSA**: Compute-Intensive

6 TF EV68
71 TB Storage
0.3 TF EV7 shared-memory
150 TB Storage Server

**PSC**: Compute-Intensive

PSC integrated Q3 03

# Teragrid Logical Network Diagram

# I-WIRE Geography

**ANL** — **Starlight / NU-C**

**UChicago** Gleacher Ctr

**UChicago** Main

**UIC** **IIT** **ICN** **UIUC/NCSA**

○ Commercial Fiber Hub

- Status:
  Done: ANL, NCSA, Starlight
  Laterals in process: UC, UIC, IIT
- Investigating extensions to
  Northwestern Evanston, Fermi,
  O'Hare, Northern Illinois Univ,
  DePaul, etc.

**Northwestern Univ-Chicago "Starlight"**

I-290

**UI-Chicago**

I-294

I-55

**Illinois Inst. Tech**

Dan Ryan Expwy (I-90/94)

**U of Chicago**

**Argonne Nat'l Lab**
(approx 25 miles SW)

**UIUC/NCSA**
Urbana (approx 140 miles South)

# State of Illinois I-WIRE

- **I-WIRE timeline**
  - 1994: Governor interest
    - schools and networks
  - 1997: task force formed
  - 1999: I-WIRE funding approved
  - 2002: fiber in place and operational

- **Features**
  - fiber providers
    - Qwest, Level(3)
    - McLeodUSA, 360Networks
  - 10 segments
  - 190 route miles and 816 fiber miles
    - longest segment is 140 miles
  - 4 strands minimum to each site



Argonne

Starlight
(NU-Chicago)

18

4

UC Gleacher Ctr
450 N. Cityfront

Qwest
455 N. Cityfront

4    10    4

UIC

4

McLeodUSA
151/155 N. Michigan
Doral Plaza

UIUC/NCSA

12

Level(3)
111 N. Canal

12

State/City Complex
James R. Thompson Ctr
City Hall
State of IL Bldg

NU-Evanston

2    2

FNAL

2

JIT

UChicago

Numbers indicate fiber count (strands)

# I-Wire Transport

**TeraGrid Linear**
3x OC192
1x OC48
First light: 6/02

**Starlight Linear**
4x OC192
4x OC48 (→8x GbE)
Operational

**Metro Ring**
1x OC48 per site
First light: 8/02

**Argonne**

**Starlight (NU-Chicago)**

UC Gleacher Ctr
450 N. Cityfront

Qwest
455 N. Cityfront

**UIC**

**UIUC/NCSA**

State/City Complex
James R. Thompson Ctr
City Hall
State of IL Bldg

**IIT**

**UChicago**

- Each of these three ONI DWDM systems have capacity of up to 66 channels, up to 10 Gb/s per channel
- Protection available in Metro Ring on a per-site basis

# Network Policy Decisions

- **The TG backplane is a closed network, internal to the TG sites.**
  - Open question: what is a TG site?

- **The TG network gear is run by the TG network team.**
  - I.e. not as individual site resources.

# Network Challenges

- **Basic Design and Architecture**
  - We think we've got this right.
- **Construction**
  - Proceeding well.
- **Operation**
  - We'll see.

# Computing and Grid Challenges

- **Hardware configuration and purchase**
  - I'm *still* not 100% sure what we'll be installing.
  - The proposal was written in early 2001.
  - The hardware is being installed in late 2002.
  - The IA-64 line of processors is young.
  - Several vendors, all defining new products, are involved.

  - Recommendations:
    - Try to avoid this kind of long-wait, multi vendor situation.
    - Have frequent communication with all vendors about schedule, expectations, configurations, etc.

# Computing and Grid Challenges

- **Understanding application requirements and getting people started before the hardware arrives.**

- **Approach: TG-Lite**
  - a small PIII testbed
  - 4 nodes at each site
  - Internet/Abilene connectivity
  - For early users and sysadmins to test configurations.

# Computing and Grid Challenges

- **Multiple sites, one environment:**
  - Sites desire different configurations.
  - Distributed administration.
  - Need a coherent environment for applications.
    - Ideal: binary compatibility

- **Approach: service definitions.**

# NSF TeraGrid: 14 TFLOPS, 750 TB

# Defining and Adopting Standard Services

Finite set of TeraGrid services-applications see *standard services* rather than *particular implementations*…

Grid Applications

…but sites also provide additional services that can be discovered and exploited.

- IA-64 Linux TeraGrid Cluster Runtime
- File-based Data Service
- IA-64 Linux Cluster Interactive Development
- Interactive Collection-Analysis Service
- Volume-Render Service
- Collection-based Data Service

# Strategy: Define Standard Services

- **Finite Number of TeraGrid Services**
  - Defined as specifications, protocols, API's
  - Separate from implementation (magic software optional)
- **Extending TeraGrid**
  - Adoption of TeraGrid specifications, protocols, API's
    - What protocols does it speak, what data formats are expected, what features can I expect (how does it behave)
    - Service Level Agreements (SLA)
  - Extension and expansion via:
    - Additional services not initially defined in TeraGrid
      - e.g. Alpha Cluster Runtime service
    - Additional instantiations of TeraGrid services
      - e.g. IA-64 runtime service implemented on cluster at a new site
- **Example: File-based Data Service**
  - API/Protocol: Supports *FTP* and *GridFTP, GSI* authentication
  - SLA
    - All TeraGrid users have access to $N$ TB storage
    - available 24/7 with $M$% availability
    - >= $R$ Gb/s read, >= $W$ Gb/s write performance

# Standards → Cyberinfrastructure

**If** done _openly_ and well…

- other IA-64 cluster sites would adopt TeraGrid service specifications, increasing users' leverage in writing to the specification
- others would adopt the framework for developing similar services on different architectures

Certificate Authority

TeraGrid Certificate Authority

Grid Info Svces

**Grid Applications**

Alpha Clusters

IA-64 Linux Clusters

IA-32 Linux Clusters

Visualization Services

_Data/Information_

File-based Data Service

Collection-based Data Service

Relational dBase Data Service

_Compute_

Interactive Development

Runtime

_Analysis_

Interactive Collection-Analysis Service

Visualization Services

# Computing and Grid Challenges

- **Architecture**
  - Individual clusters architectures are fairly solid.
  - Aggregate architecture is a bigger question.
    - Being defined in terms of services.
- **Construction and Deployment**
  - We'll see, starting in December.
- **Operation**
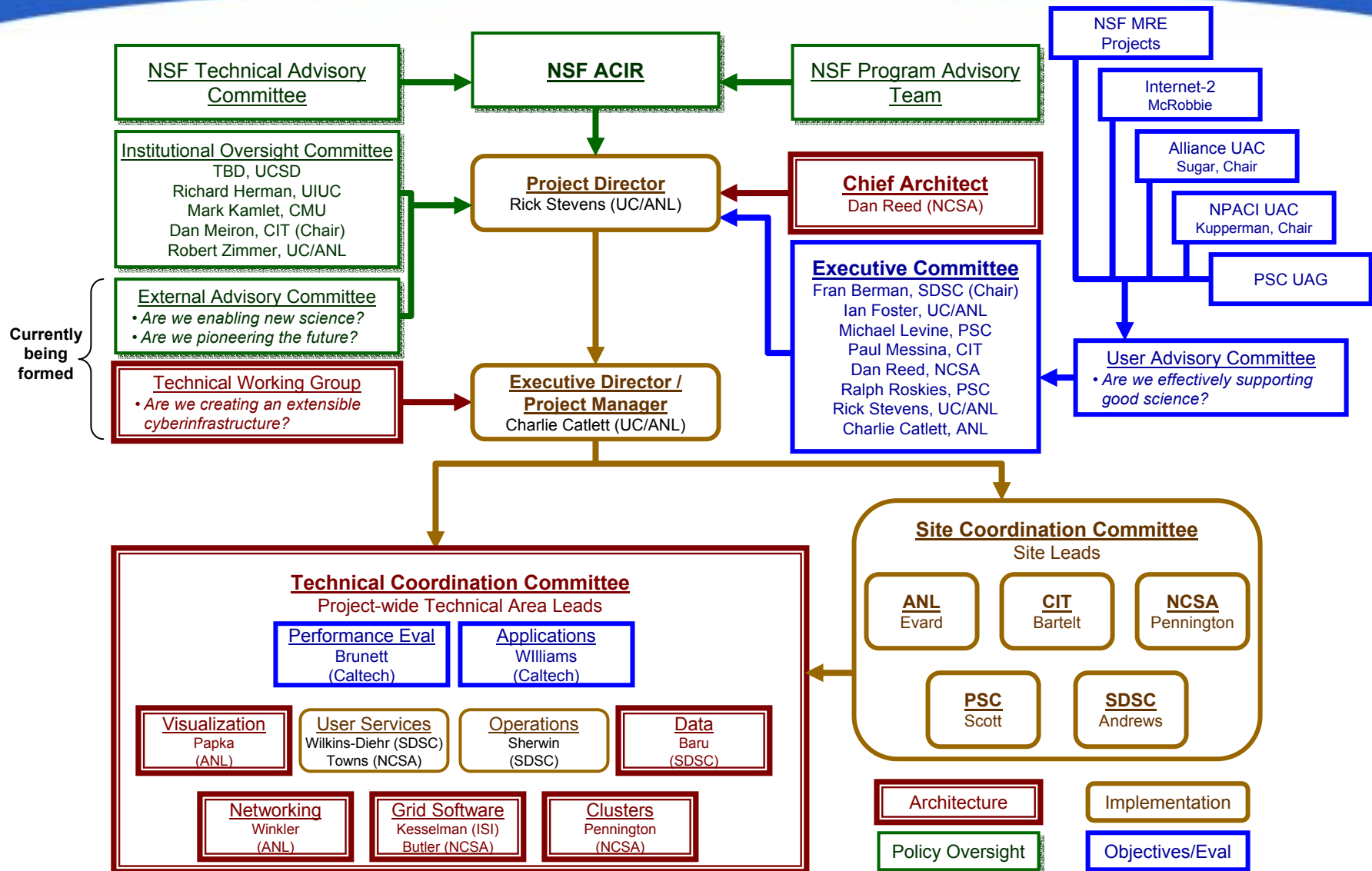  - We'll see.  Production by June 2003.

# Social Issues: Direction

- **4 sites tend to have 4 directions.**
  - NCSA and SDSC have been competitors for over a decade.
    - This has created surprising cultural barriers that must be recognized and overcome.
    - Including PSC, a 3$^{rd}$ historical competitor, will complicate this.
  - ANL and Caltech are smaller sites with fewer resources but specific expertise. And opinions.

# Social Issues: Organization

- **Organization is a big deal.**
  - Equal/fair participation among sites.
    - To the extreme credit of the large sites, this project has been approached as 4 peers, not 2 tiers. This has been extremely beneficial.
  - Project directions and decisions affect all sites.
    - How best to distribute responsibilities but make coordinated decisions?
  - Changing the org chart is a heavyweight operation, best to be avoided…

# The ETF Organizational Chart

# Social Issues: Working Groups

- **Mixed effectiveness of working groups**
  - The networking working group has turned into a team.
  - The cluster working group is less cohesive.
  - Others range from teams to just email lists.
- **Why?**
  - Not personality issues, not organizational issues.
- **What makes the networking group tick:**
  - Networking people already work together:
    - The individuals have a history of working together on other projects.
    - They see each other at other events.
    - They're expected to travel.
    - They held meetings to decide how to build the network before the proposal was completed.
  - The infrastructure is better understood:
    - Networks somewhat like this have been built before.
    - They are building one network, not four clusters.
    - There is no separation between design, administration, and operation.
- **Lessons:**
  - Leverage past collaborations that worked.
  - Clearly define goals and responsibilities.

# Social Issues: Observations

- **There will nearly always be four opinions on every issue.**
  - Reaching a common viewpoint takes a lot of communication.
  - Not every issue can actually be resolved.
  - Making project-wide decisions can be tough.
- **Thus far in the project, the social issues have been just as complex as the technical.**
  - … but the technology is just starting to arrive…
- **It's possible we should have focused more on this in the early proposal stage, or allocated more resources to helping with these.**
  - We have, just appointed a new "Director of Engineering" to help guide technical decisions and maintain coherency.

# Conclusion

- **Challenges abound!  Early ones include:**
  - Network design and deployment.
  - Cluster design and deployment.
  - Building the right distributed system architecture into the grid.
  - Getting along and having fun.
  - Expansion.

- **The hardware arrives in December, production is in mid-2003.**
- **Check back in a year to see how things are going…**