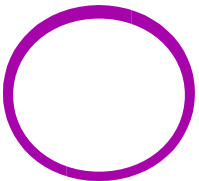# TeraGrid Communication and Computation

**Tal Lavian**
**tlavian@cs.berkeley.edu**

**Brainstorm and concepts**
**Many slides and most of the graphics are taken from other slides**

# Agenda

**Introduction**

**Some applications**
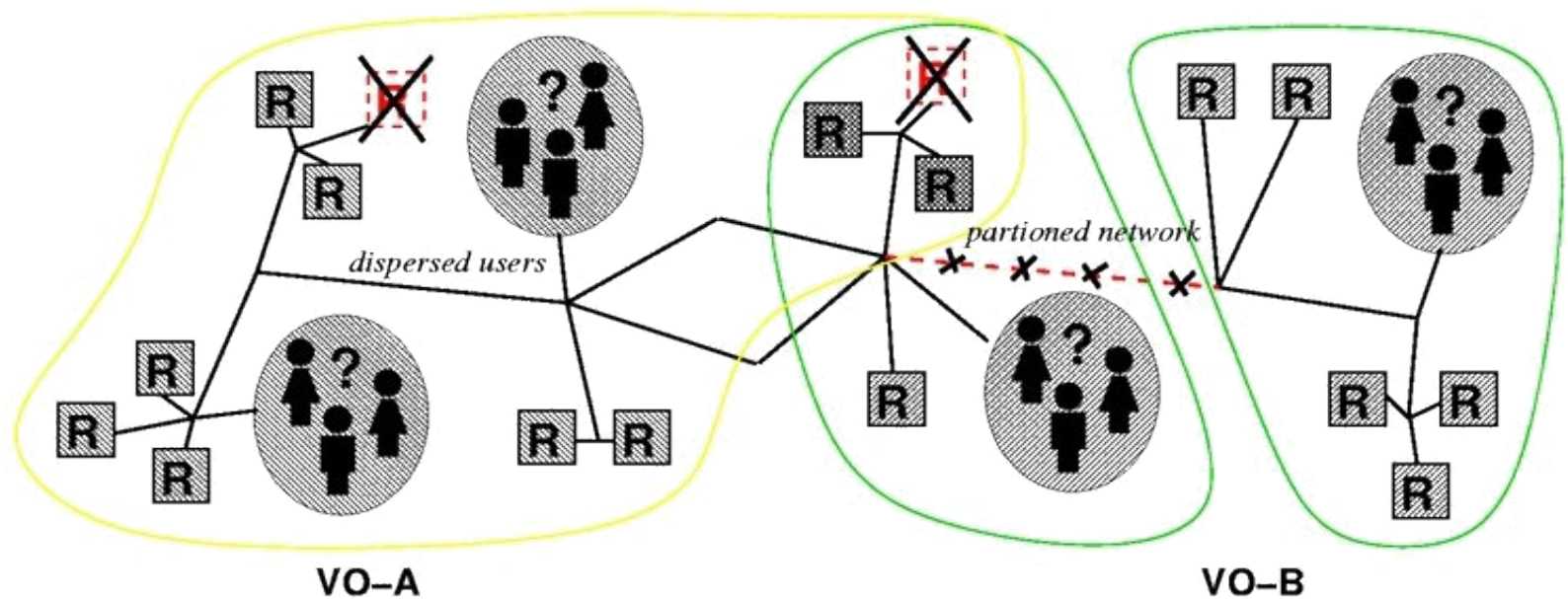
**TeraGrid Architecture**

**Globus toolkit**

**Future comm direction**

**Summary**

# The Grid Problem

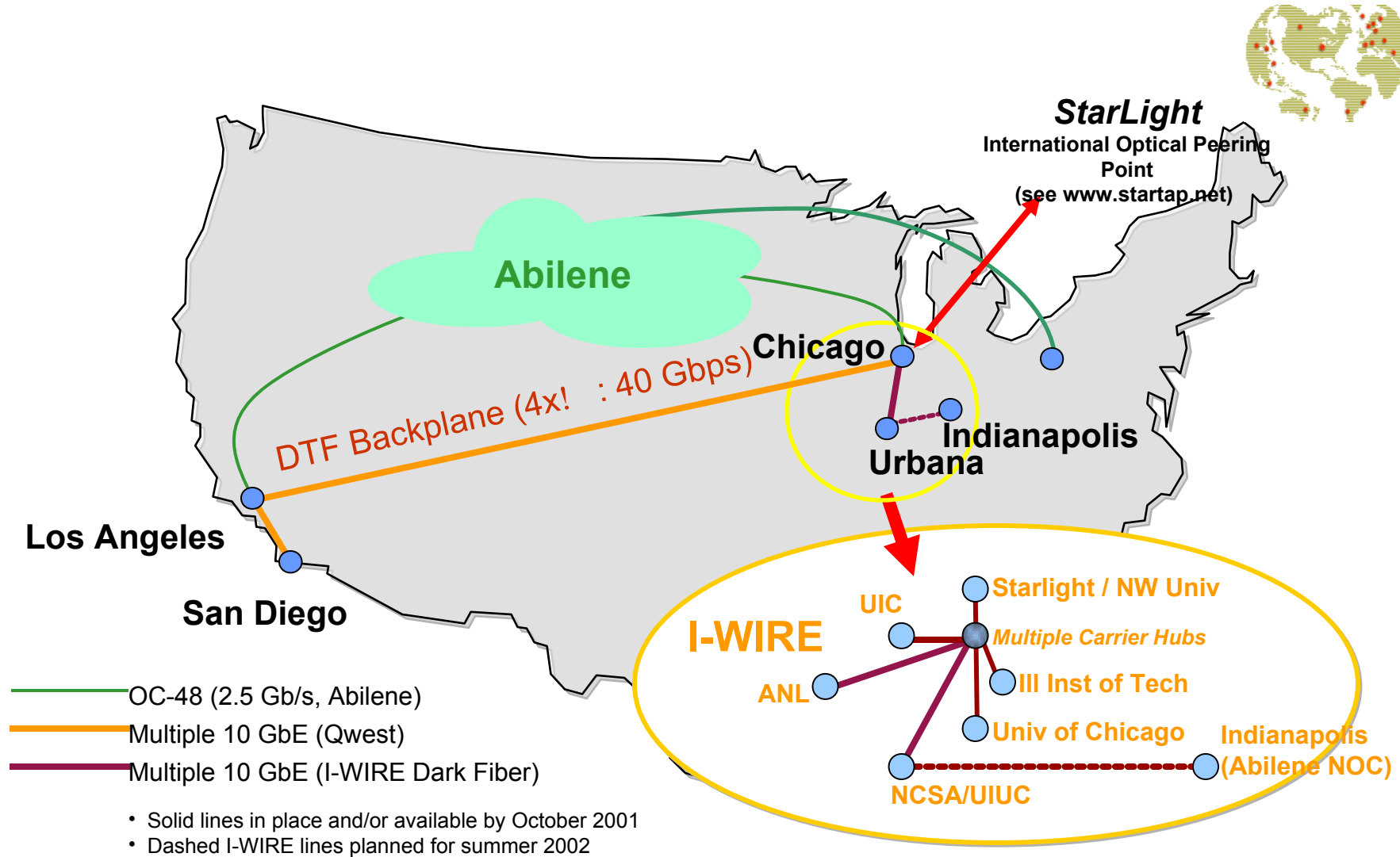**Resource sharing & coordinated problem solving in dynamic, multi-institutional virtual organizations**
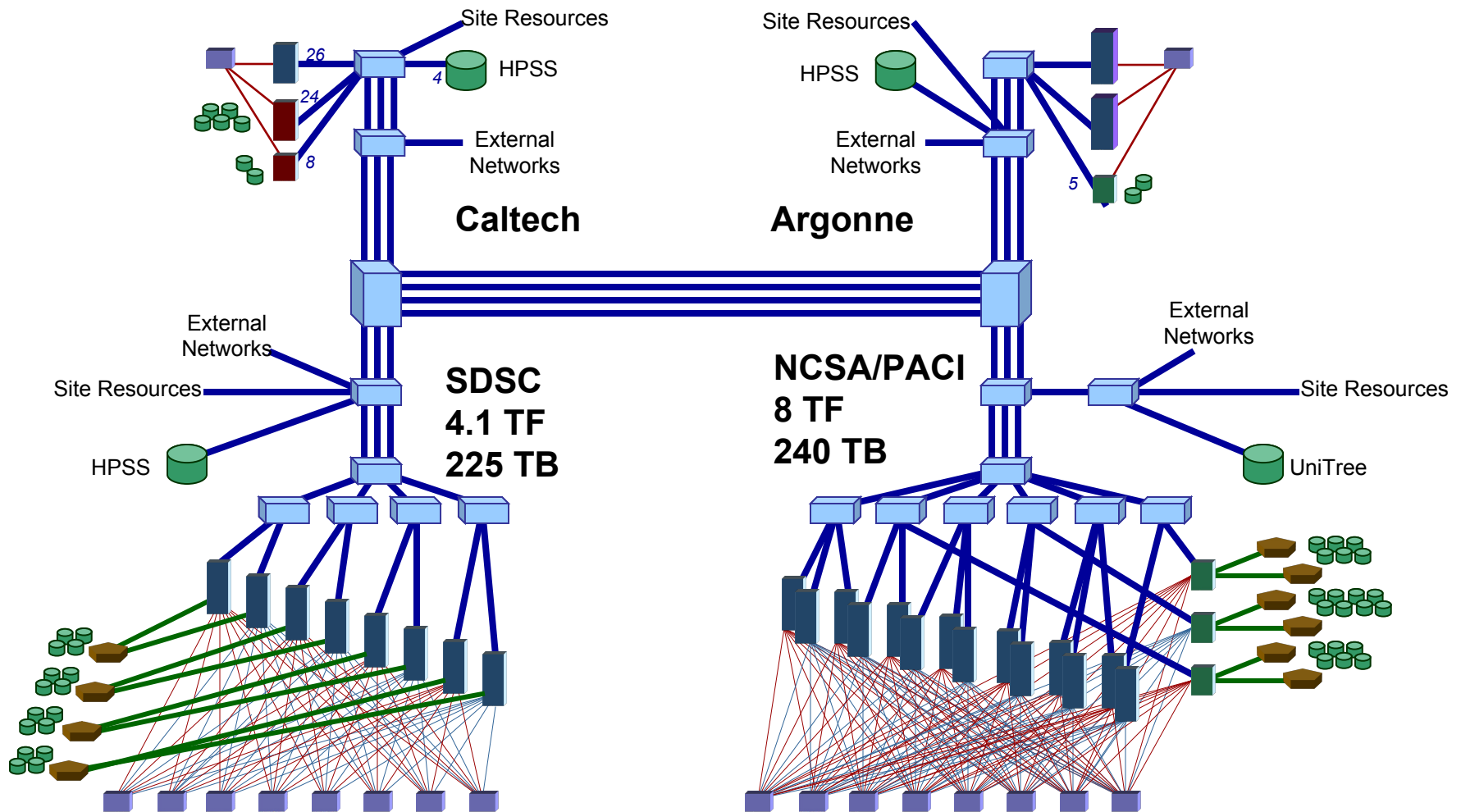


**Some relation to Sahara**

**Service composition: computation, servers, storage, disk, network…**
**Sharing, cooperating, peering, brokering**…

# TeraGrid Wide Area Network - NCSA, ANL, SDSC, Caltech



*StarLight*
**International Optical Peering Point**
**(see www.startap.net)**

**Abilene**

**Chicago**

DTF Backplane (4x! : 40 Gbps)

**Indianapolis**
**Urbana**

**Los Angeles**

**San Diego**

**I-WIRE**

**UIC**

**Starlight / NW Univ**

*Multiple Carrier Hubs*

**ANL**

**Ill Inst of Tech**

**Univ of Chicago**

**Indianapolis (Abilene NOC)**

**NCSA/UIUC**

— OC-48 (2.5 Gb/s, Abilene)
— Multiple 10 GbE (Qwest)
— Multiple 10 GbE (I-WIRE Dark Fiber)

- Solid lines in place and/or available by October 2001
- Dashed I-WIRE lines planned for summer 2002

**TeraGrid Comm & Comp**

# The 13.6 TF TeraGrid: Computing at 40 Gb/s



Site Resources

HPSS

External Networks

**Caltech**

External Networks

Site Resources

HPSS

**SDSC**
**4.1 TF**
**225 TB**

Site Resources

HPSS

Site Resources

External Networks

**Argonne**

**NCSA/PACI**
**8 TF**
**240 TB**

External Networks

Site Resources

UniTree

TeraGrid/DTF: NCSA, SDSC, Caltech, Argonne          www.teragrid.org

# 4 TeraGrid Sites Have Focal Points

## SDSC – The Data Place

Large-scale and high-performance data analysis/handling

Every Cluster Node is Directly Attached to SAN

## NCSA – The Compute Place

Large-scale, Large Flops computation

## Argonne – The Viz place

Scalable Viz walls

## Caltech – The Applications place

Data and flops for applications – Especially some of the GriPhyN Apps

## Specific machine configurations reflect this

# TeraGrid building blocks

**Distributed, multisite facility**

single site and "Grid enabled" capabilities
- ➢ **uniform compute node selection and interconnect networks at 4 sites**
- ➢ **central "Grid Operations Center"**

at least one 5+ teraflop site and newer generation processors
- ➢ **SDSC at 4+ TF, NCSA at 6.1-8 TF with McKinley processors**

at least one additional site coupled with the first
- ➢ **four core sites: SDSC,  NCSA,  ANL, and Caltech**

**Ultra high-speed networks (Static configured)**

multiple gigabits/second
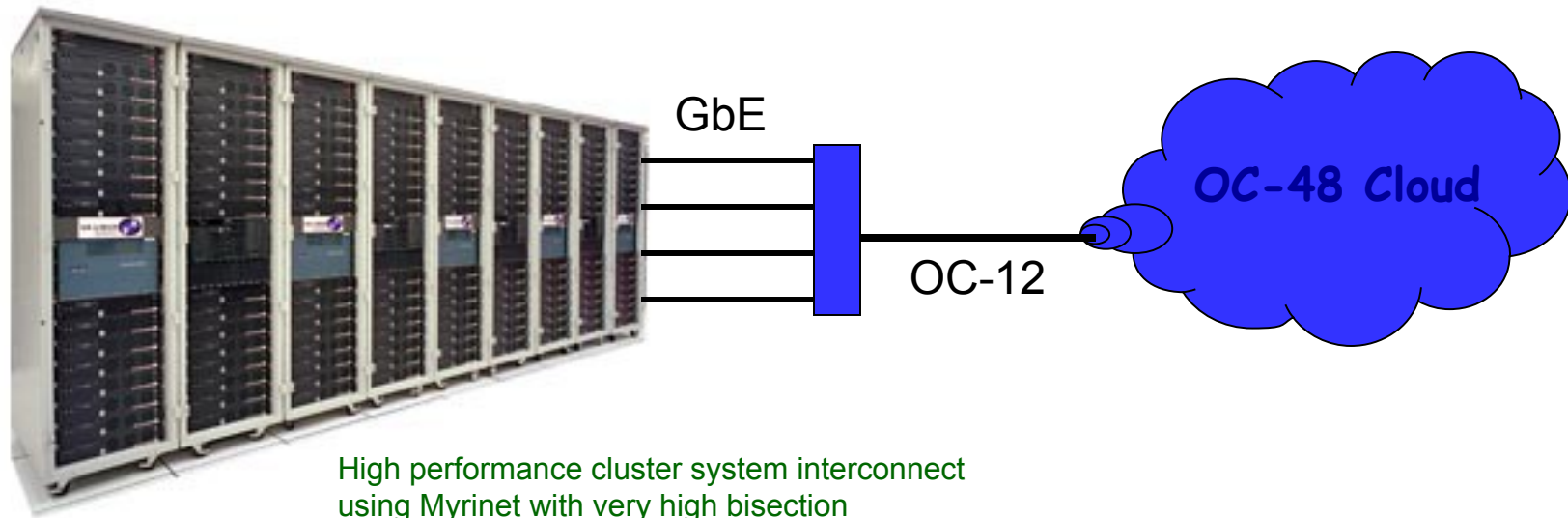- ➢ **modular 40 Gb/s backbone (4 x 10 GbE)**

**Remote visualization**

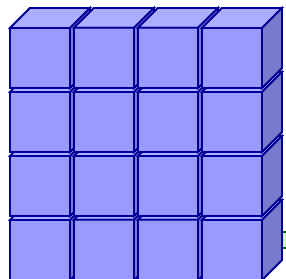data from one site visualized at another
- ➢ **high-performance commodity rendering and visualization system**
- ➢ **Argonne hardware visualization support**
- ➢ **data serving facilities and visualization displays**

**NSF - $53M award in August 2001**

# Traditional Cluster Network Access

GbE

OC-48 Cloud

OC-12

High performance cluster system interconnect using Myrinet with very high bisection bandwidth (hundreds of GB/s) with external connection of n x GbE, n is small integer.

(Time to move entire contents of memory)

2000 s (33 min)
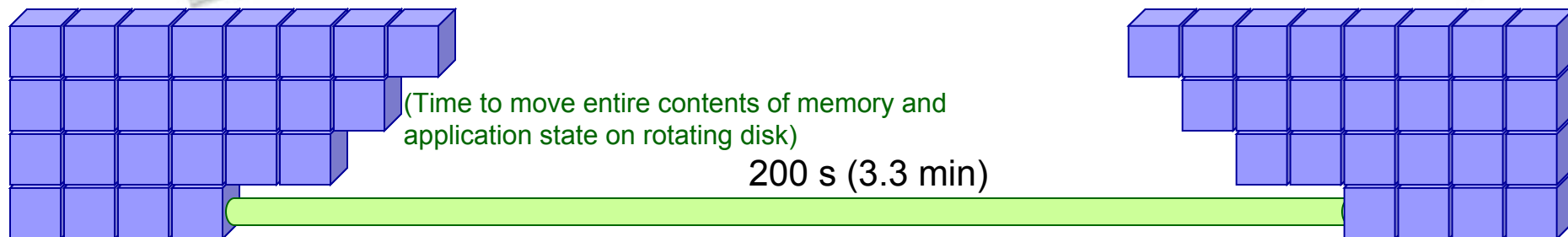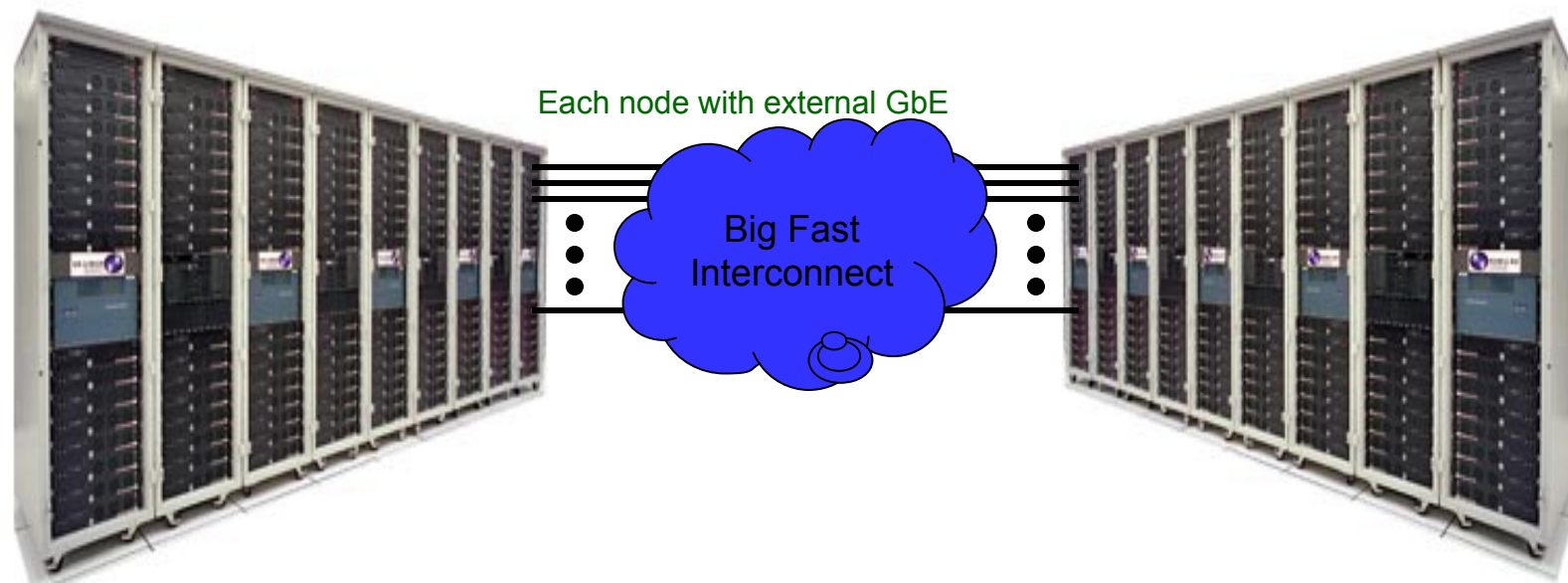
13k s (3.6h)

1 TB

0.5 GB/s

78 MB/s

64 GB

1024 MB

Traditionally, high-performance computers have been islands of capability separated by wide area networks that provide a fraction of a percent of the internal cluster network bandwidth.

# To Build a Distributed Terascale Cluster...

Each node with external GbE

Big Fast Interconnect

(Time to move entire contents of memory and application state on rotating disk)

200 s (3.3 min)

10 TB

5 GB/s

10 TB

**5 GB/s = 200 nodes x 25 MB/s (=20% of GbE per node)**

4096 GB

64 GB

TeraGrid is building a "machine room" network across the country while increasing external cluster bandwidth to many GbE. Requires edge systems that handle n x 10 GbE and hubs that handle minimum 10 x 10 GbE.

# Agenda

**Introduction**

**Some applications**

**TeraGrid Architecture**

**Globus toolkit**

**Future comm direction**

**Summary**

# What applications are being targeted for Grid-enabled computing? Traditional

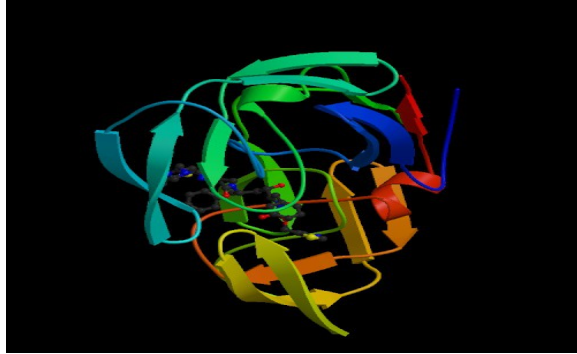**Quantum Chromodynamics**

**Biomolecular Dynamics**

**Weather Forecasting**

**Cosmological Dark Matter**

**Biomolecular Electrostatics**

**Electric and Magnetic Molecular Properties**

# Beginning of the Digital Millennium:
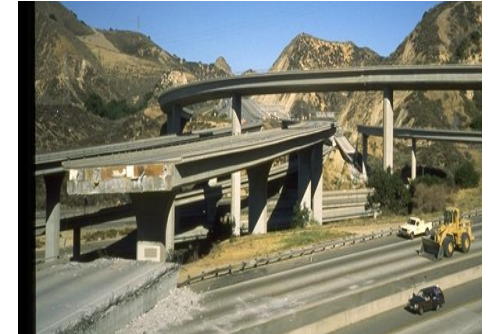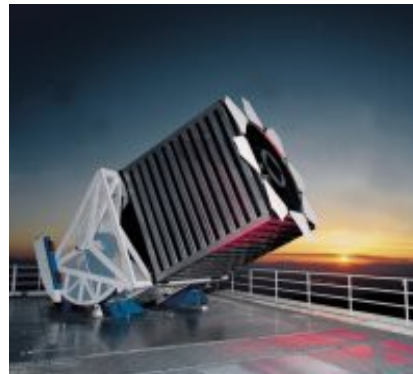# The Data Decade!

**Sensors**

**Genomics**

**Disaster response**

**Physics**

**Digital Libraries**

**Astronomy**

# Multi-disciplinary Simulations: Aviation Safety



**Wing Models**
- Lift Capabilities
- Drag Capabilities

**Stabilizer Models**
- Deflection capabilities
- Responsiveness

**Airframe Models**

Crew Capabilities
- accuracy
- perception
- stamina
- re-action times

**Human Models**

**Engine Models**
- Thrust performance
- Reverse Thrust performance
- Responsiveness
- Fuel Consumption

**Landing Gear Models**
- Braking performance
- Steering capabilities
- Traction
- Dampening capabilities

**Source NASA**

**Whole system simulations are produced by coupling all of the sub-system simulations**

# New Results Possible on TeraGrid

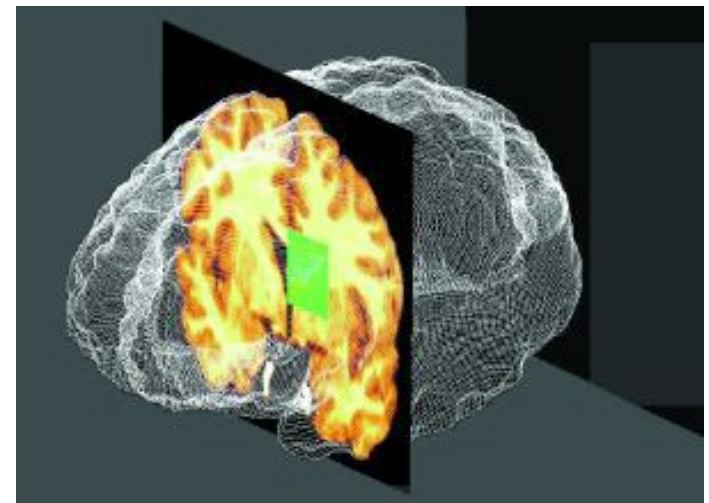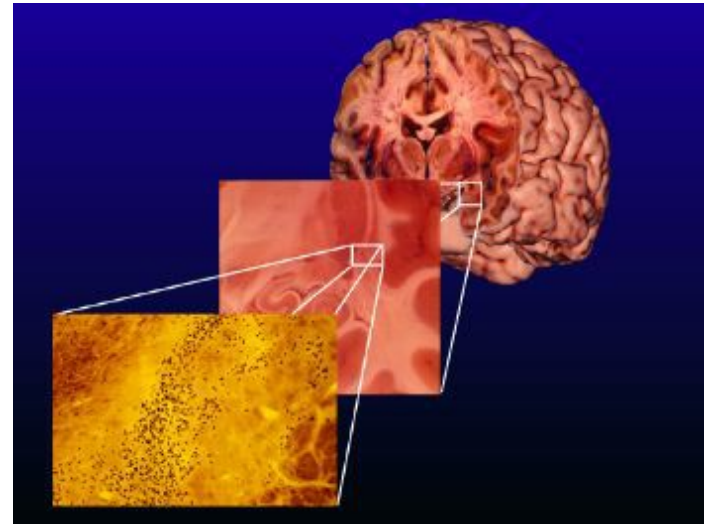**Biomedical Informatics Research Network (National Inst. Of Health):**

Evolving reference set of brains provides essential data for developing therapies for neurological disorders (Multiple Sclerosis, Alzheimer's, etc.).
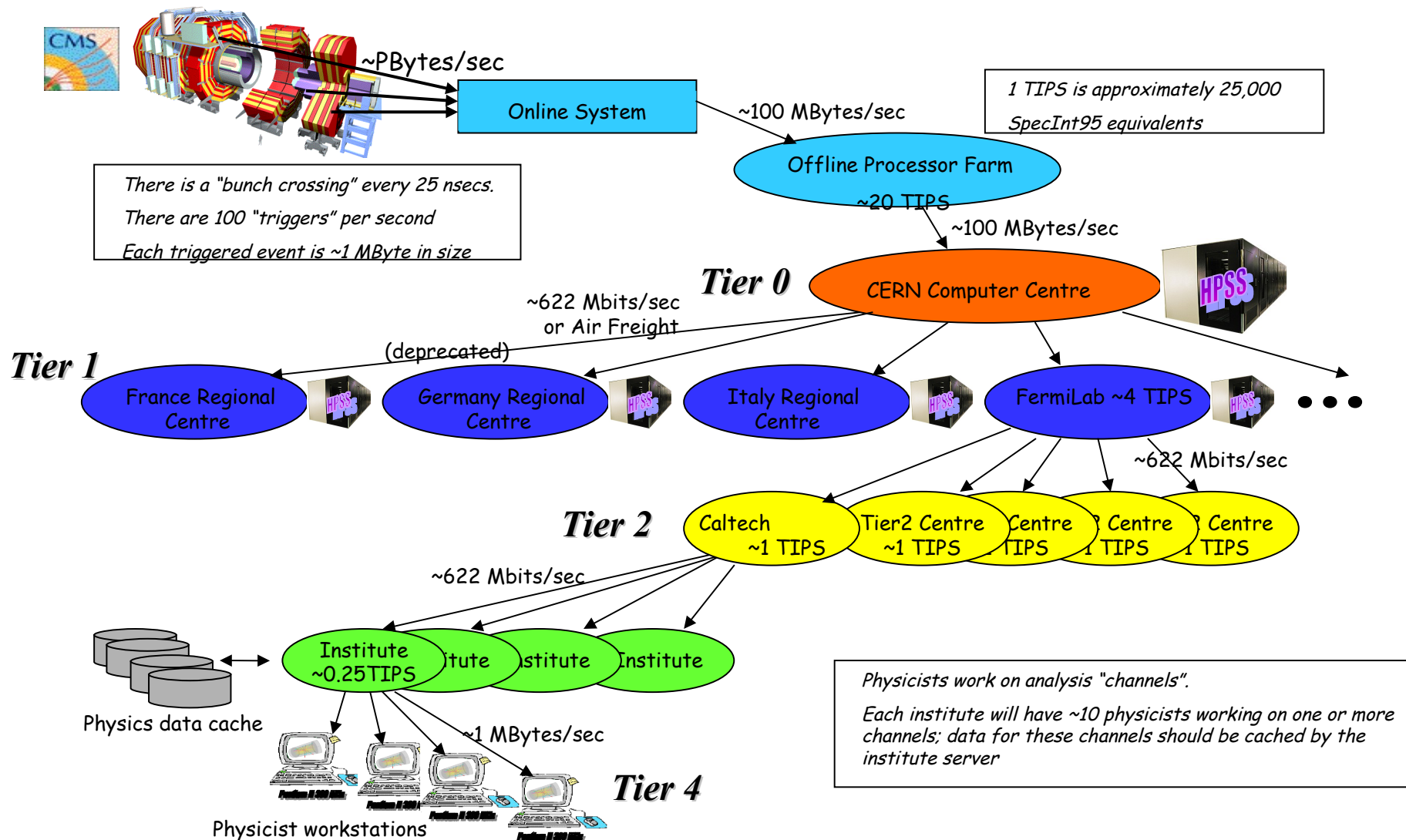
**Pre-TeraGrid:**

One lab

Small patient base

4 TB collection

**Post-TeraGrid:**

Tens of collaborating labs

Larger population sample

400 TB data collection:  more brains, higher resolution

Multiple scale data integration and analysis

# Grid Communities & Applications:
# Data Grids for High Energy Physics



~PBytes/sec

Online System

~100 MBytes/sec

*1 TIPS is approximately 25,000 SpecInt95 equivalents*

*There is a "bunch crossing" every 25 nsecs.*

*There are 100 "triggers" per second*

*Each triggered event is ~1 MByte in size*

Offline Processor Farm

~20 TIPS

~100 MBytes/sec

**Tier 0**   CERN Computer Centre   HPSS

~622 Mbits/sec
or Air Freight

(deprecated)

**Tier 1**   France Regional Centre   HPSS   Germany Regional Centre   HPSS   Italy Regional Centre   HPSS   FermiLab ~4 TIPS   HPSS   • • •

~622 Mbits/sec

**Tier 2**   Caltech ~1 TIPS   Tier2 Centre ~1 TIPS   Centre ~1 TIPS   Centre ~1 TIPS   Centre ~1 TIPS

~622 Mbits/sec

Institute ~0.25TIPS   Itute   stitute   Institute

Physics data cache

*Physicists work on analysis "channels".*

*Each institute will have ~10 physicists working on one or more channels; data for these channels should be cached by the institute server*

~1 MBytes/sec

**Tier 4**

Physicist workstations

# Agenda

**Introduction**

**Some applications**

**TeraGrid Architecture**

**Globus toolkit**

**Future comm direction**

**Summary**

# Grid Computing Concept

**New applications enabled by the coordinated use of geographically distributed resources**

E.g., distributed collaboration, data access and analysis, distributed computing

**Persistent infrastructure for Grid computing**

E.g., certificate authorities and policies, protocols for resource discovery/access

**Original motivation, and support, from high-end science and engineering; but has wide-ranging applicability**

# Globus Hourglass

## Focus on architecture issues

Propose set of core services as basic infrastructure

Use to construct high-level, domain-specific solutions

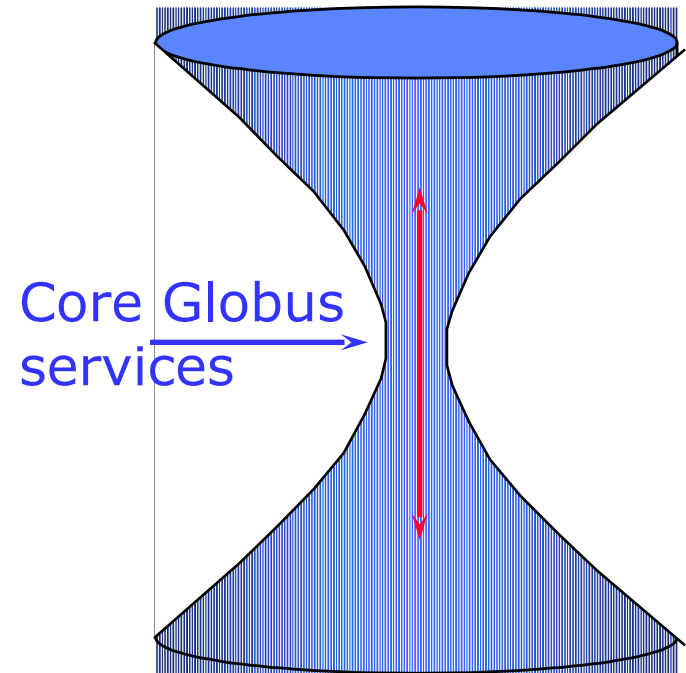## Design principles

Keep participation cost low

Enable local control

Support for adaptation

"IP hourglass" model



**Applications**

Diverse global services

Core Globus services

Local OS

# Elements of the Problem

## Resource sharing

Computers, storage, sensors, networks, …

Sharing always conditional: issues of trust, policy, negotiation, payment, …

## Coordinated problem solving

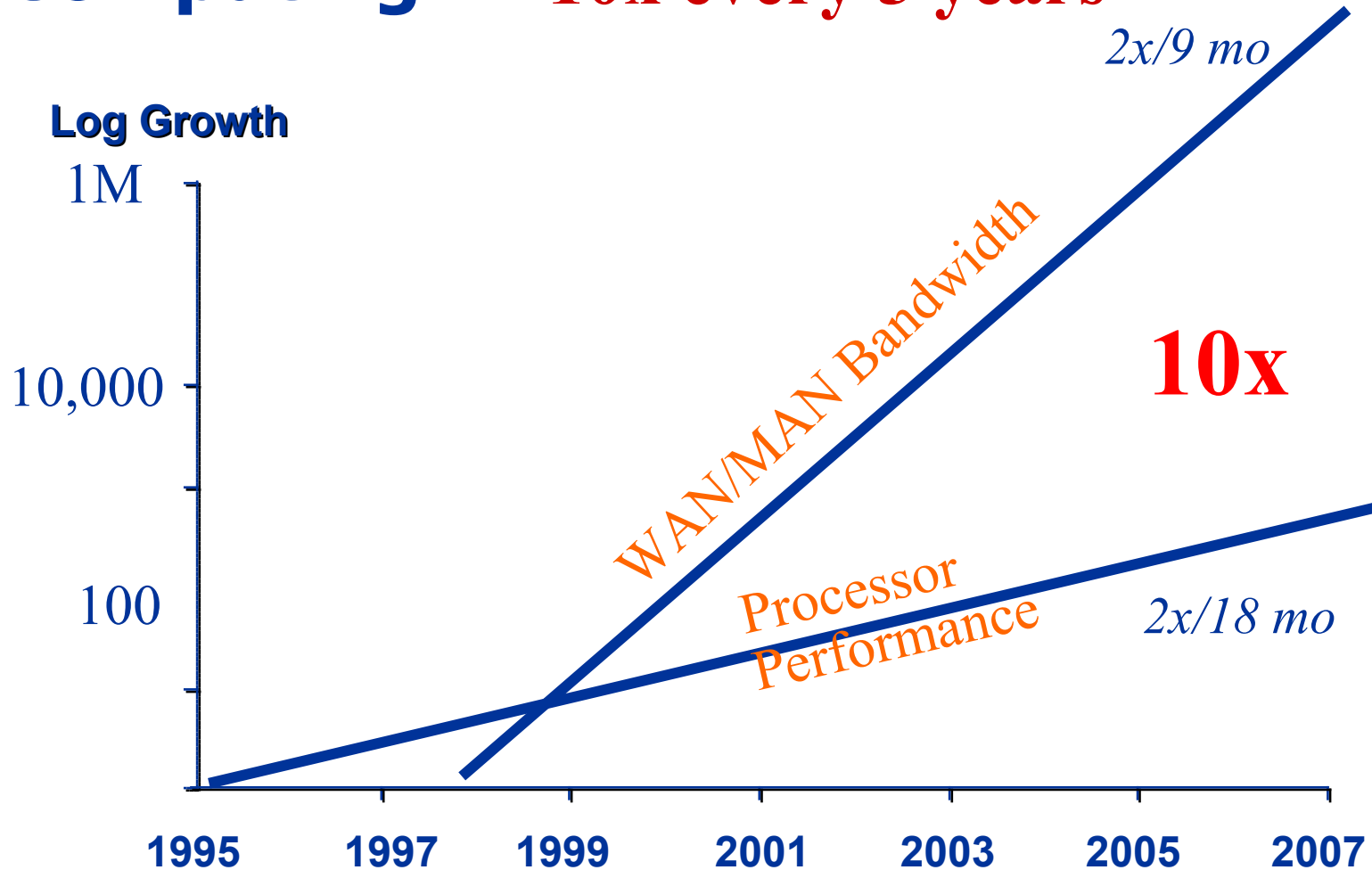Beyond client-server: distributed data analysis, computation, collaboration, …

## Dynamic, multi-institutional virtual orgs

Community overlays on classic org structures

Large or small, static or dynamic

# Gilder vs. Moore – Impact on the Future of Computing

**10x every 5 years**

*2x/9 mo*

**Log Growth**

1M

10,000

100

**10x**

WAN/MAN Bandwidth

Processor Performance

*2x/18 mo*

1995   1997   1999   2001   2003   2005   2007

# Improvements in Large-Area Networks

**Network vs. computer performance**

Computer speed doubles every 18 months

Network speed doubles every 9 months

Difference = order of magnitude per 5 years
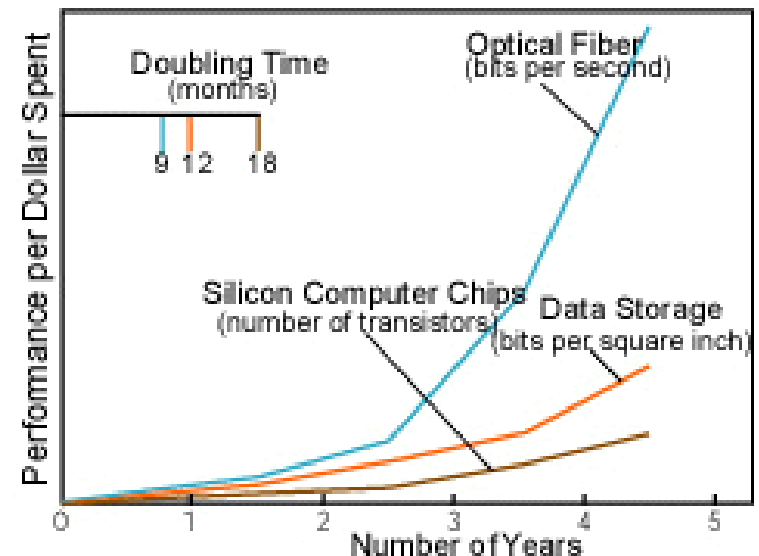
**1986 to 2000**

Computers: x 500
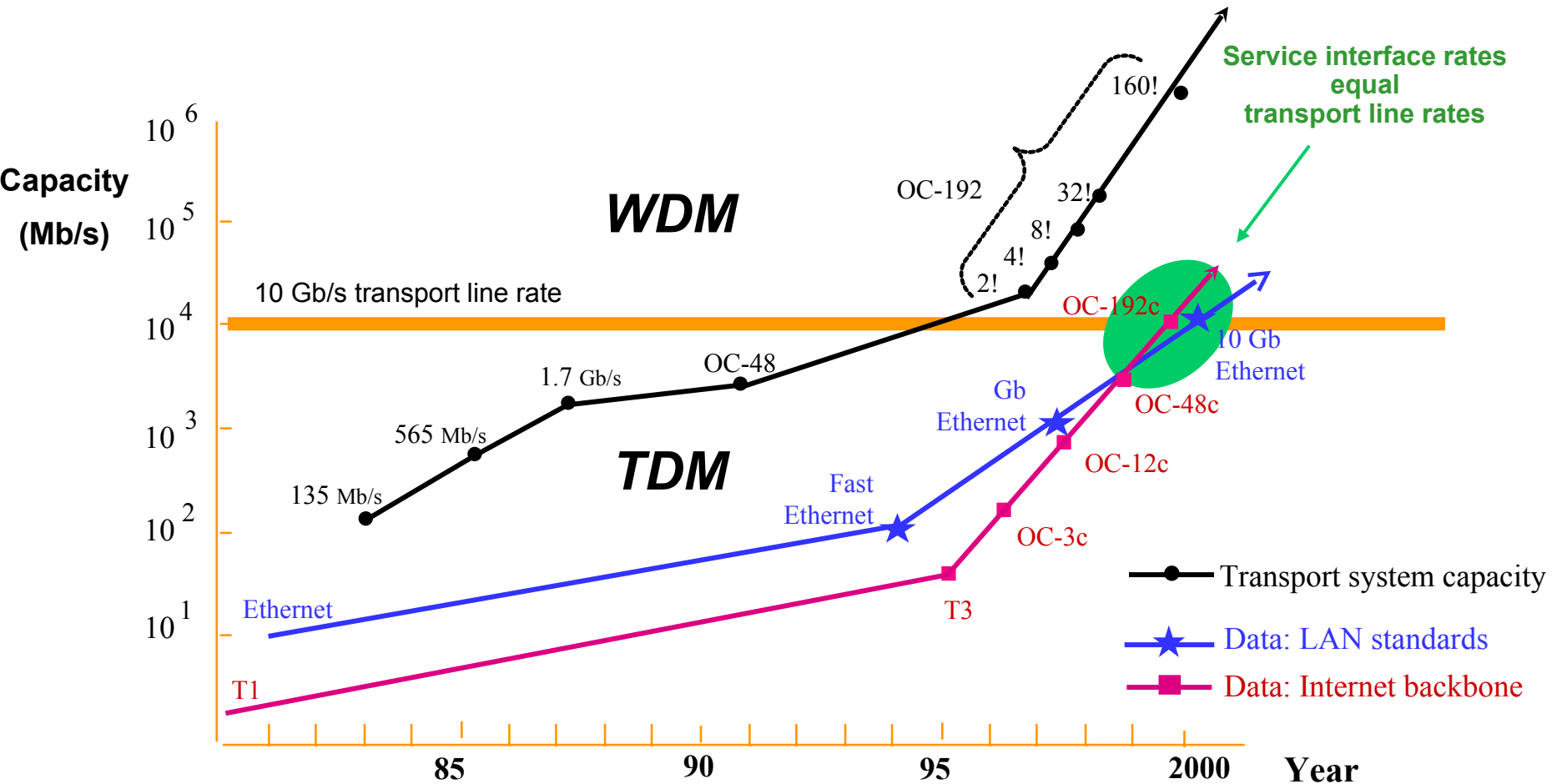
Networks: x 340,000

**2001 to 2010**

Computers: x 60

Networks: x 4000



**Moore's Law vs. storage improvements vs. optical improvements.** Graph from **Scientific American** (Jan-2001) by Cleo Vilett, source Vined Khoslan, Kleiner, Caufield and Perkins.

# Evolving Role of Optical Layer



Source: IBM WDM research

# Scientific Software Infrastructure
## *One of the Major Software Challenges*

Peak Performance is skyrocketing (more than Moore's Law)

but ...

**Efficiency has declined from 40-50% on the vector supercomputers of 1990s to as little as 5-10% on parallel supercomputers of today and may decrease further on future machines**

Research challenge is software

**Scientific codes to model and simulate physical processes and systems**

**Computing and mathematics software to enable use of advanced computers for scientific applications**

**Continuing challenge as computer architectures undergo fundamental changes:** *Algorithms that scale to thousands-millions processors*

# Agenda

**Introduction**

**Some applications**

**TeraGrid Architecture**

**Globus toolkit**

**Future comm direction**

**Summary**

# Globus Approach

**A toolkit and collection of services addressing key technical problems**

   Modular "bag of services" model

   Not a vertically integrated solution

   General infrastructure tools (aka middleware) that can be applied to many application domains

**Inter-domain issues, rather than clustering**

   Integration of intra-domain solutions

**Distinguish between local and global services**

# Globus Technical Focus & Approach

**Enable incremental development of grid-enabled tools and applications**

*Model neutral*: Support many programming models, languages, tools, and applications

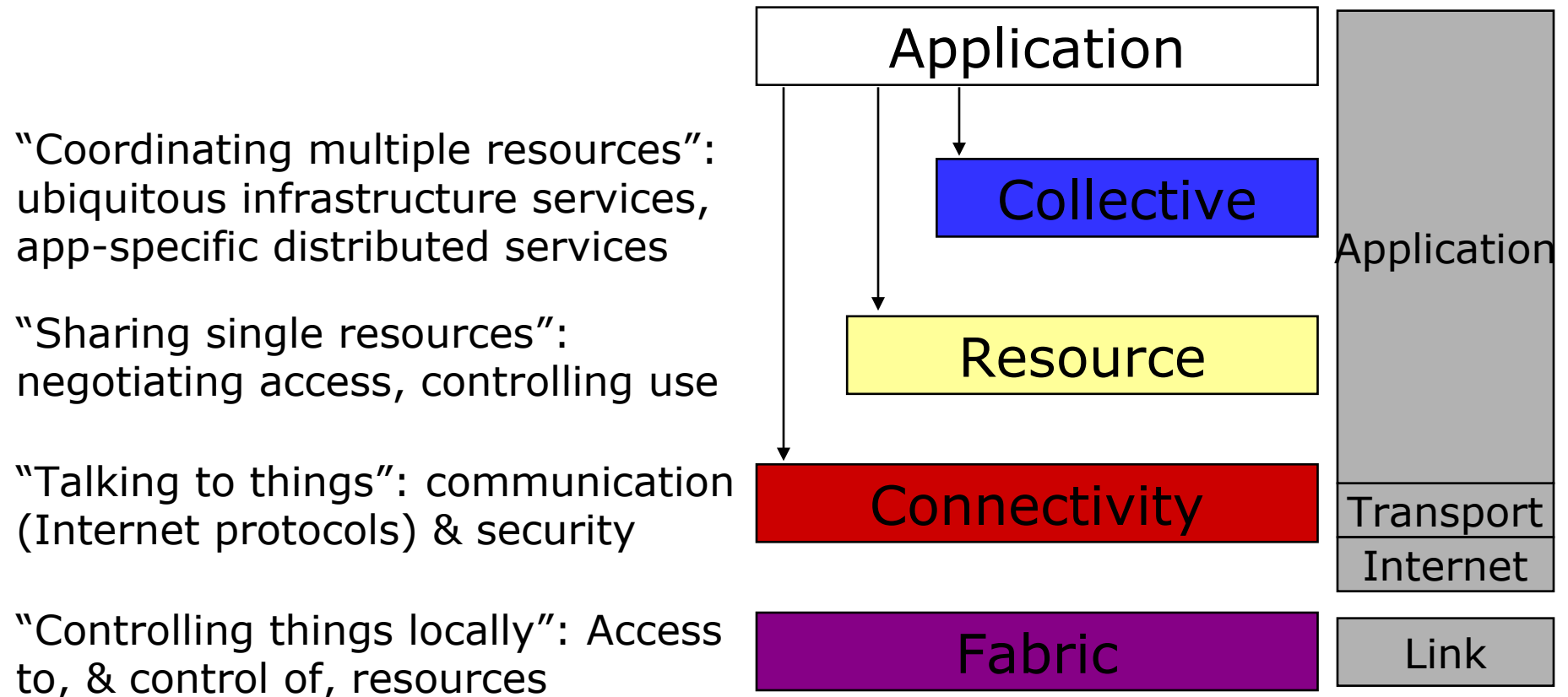Evolve in response to user requirements

**Deploy toolkit on international-scale production grids and testbeds**

Large-scale application development & testing

**Information-rich environment**

Basis for configuration and adaptation

# Layered Grid Architecture
# (By Analogy to Internet Architecture)

"Coordinating multiple resources": ubiquitous infrastructure services, app-specific distributed services

"Sharing single resources": negotiating access, controlling use

"Talking to things": communication (Internet protocols) & security

"Controlling things locally": Access to, & control of, resources

| Application |
| Collective |
| Resource |
| Connectivity |
| Fabric |

| Application |
| Transport |
| Internet |
| Link |

For more info: www.globus.org/research/papers/anatomy.pdf

# Globus Architecture?

**No "official" standards exist**

**But:**

Globus Toolkit has emerged as the de facto standard for several important Connectivity, Resource, and Collective protocols

Technical specifications are being developed for architecture elements: e.g., security, data, resource management, information
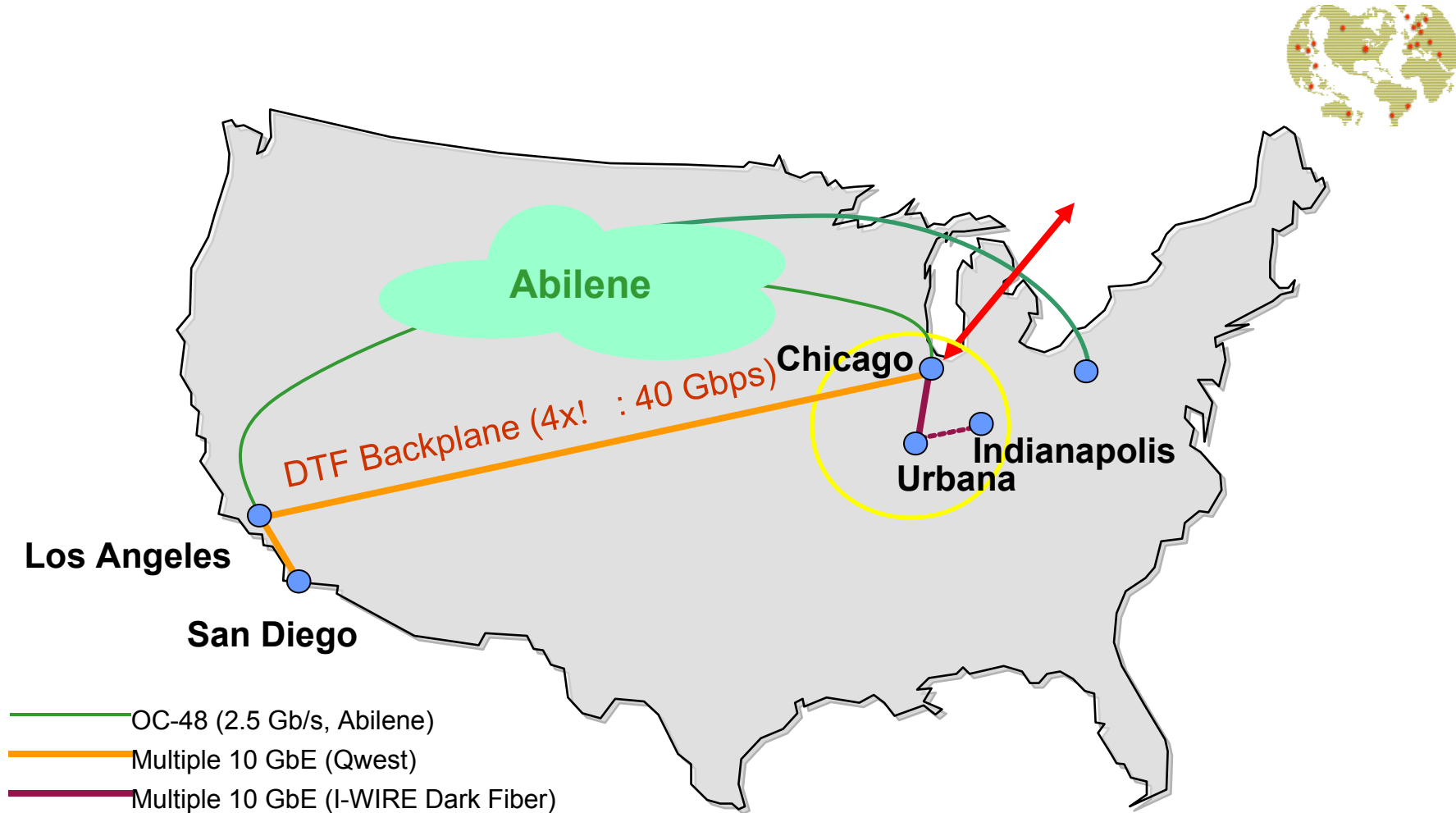
# Agenda

**Introduction**

**Some applications**

**TeraGrid Architecture**

**Globus toolkit**

**Future comm direction**

**Summary**

# Static lightpath setting NCSA, ANL, SDSC, Caltech

**Abilene**

**Chicago**

DTF Backplane (4x! : 40 Gbps)

**Indianapolis**

**Urbana**

**Los Angeles**

**San Diego**

—— OC-48 (2.5 Gb/s, Abilene)

—— Multiple 10 GbE (Qwest)

—— Multiple 10 GbE (I-WIRE Dark Fiber)

- Solid lines in place and/or available by October 2001
- Dashed I-WIRE lines planned for summer 2002

**Source: Charlie Catlett, Argonne**

# Lightpath for OVPN

**Lightpath setup**
- One or two-way
- Rates: OC48, OC192 and OC768
- QoS constraints
- On demand

**Aggregation of BW**
- OVPN
- Video
- HDTV

ASON

ASON

**Optical Ring**

ASON

*Mirror Server*

OVPN
video
HDTV

Optical fiber and channels

**TeraGrid Comm & Comp**

# Dynamic Lightpath setting

**Resource optimization (route 2)**
Alternative lightpath

**Route to mirror sites (route 3)**
Lightpath setup failed
Load balancing
Long response time
  ➢ **Congestion**
  ➢ **Fault**

ASON

Route 3

Route 2

Optical Ring

ASON

ASON

ASON

Route 1

*Mirror Server*

*main Server*

**Multiple Architectural Considerations**

Apps

Clusters

Dynamically Allocated Lightpaths

Switch Fabrics

Physical Monitoring

CONTROL PLANE

# Agenda

**Introduction**

**Some applications**

**TeraGrid Architecture**

**Globus toolkit**

**Future comm direction**

**Summary**

# Summary

**The Grid problem: Resource sharing & coordinated problem solving in dynamic, multi-institutional virtual organizations**

**Grid architecture: Emphasize protocol and service definition to enable interoperability and resource sharing**

**Globus Toolkit a source of protocol and API definitions, reference implementations**

**Current static communication. Next wave dynamic optical VPN**

**Some relation to Sahara**

> **Service composition: computation, servers, storage, disk, network…**

> **Sharing, cooperating, peering, brokering…**

# References

globus.org

 griphyn.org

gridforum.org

grids-center.org

nsf-middleware.org

# Backup

# Wavelengths and the Future

**Wavelength services are causing a network revolution:**

Core long distance SONET Rings will be replaced by meshed networks using wavelength cross-connects
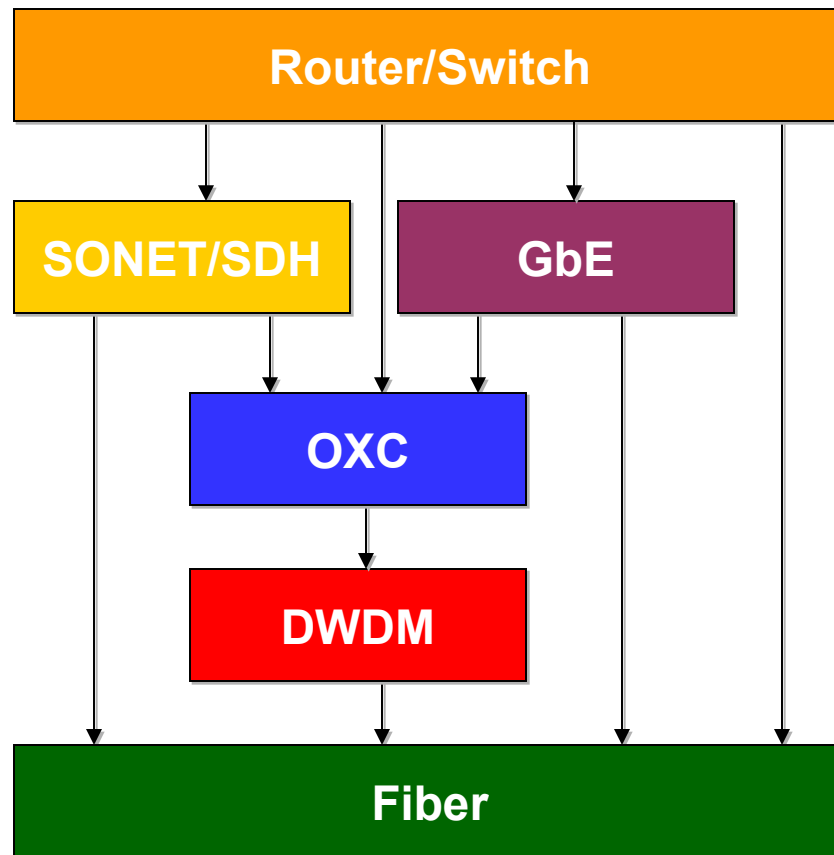
Re-invention of pre-SONET network architecture

**Improved transport infrastructure will exist for IP/packet services**

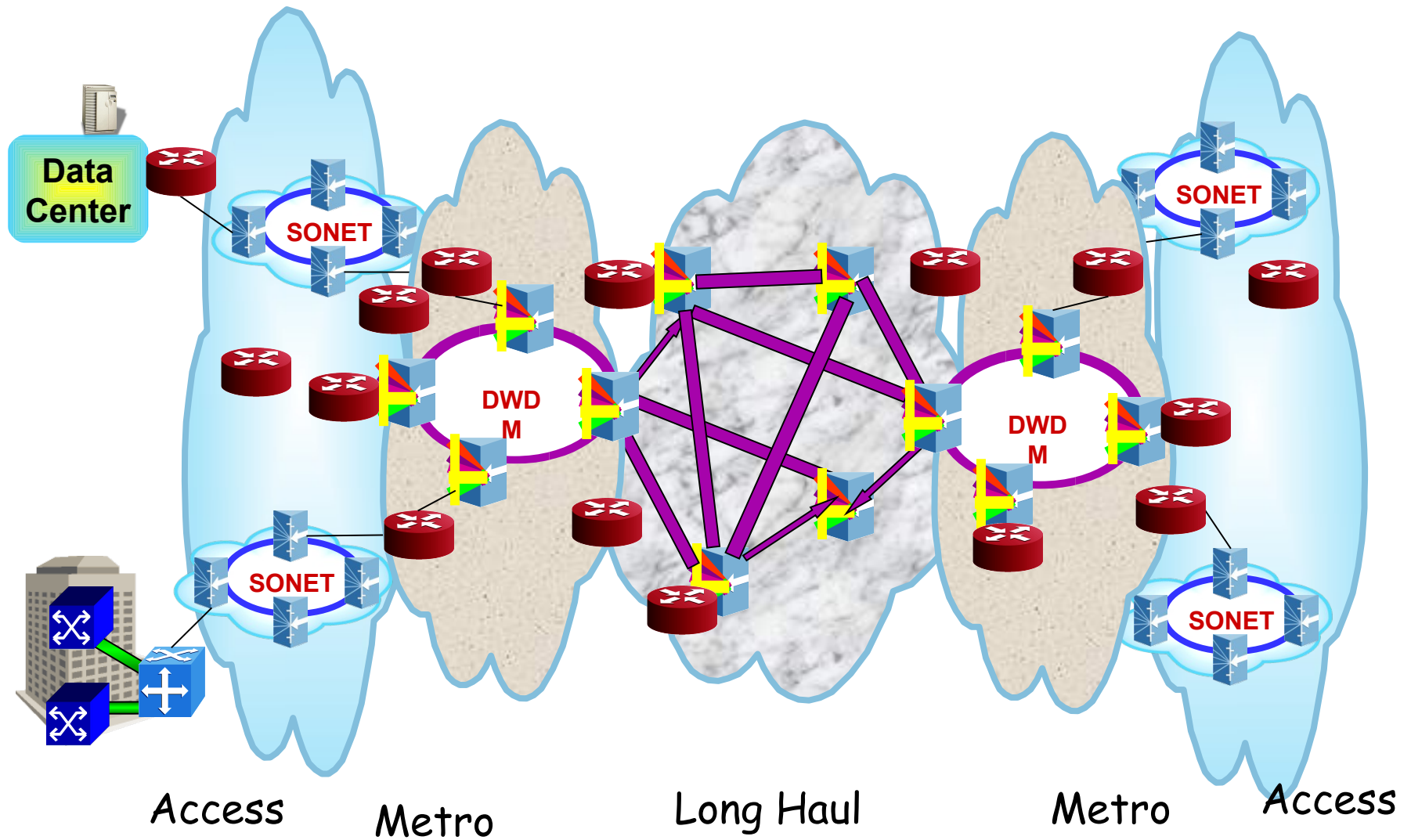**Electrical/Optical grooming switches will emerge at edges**

**Automated Restoration (algorithm/GMPLS driven) becomes technically feasible.**
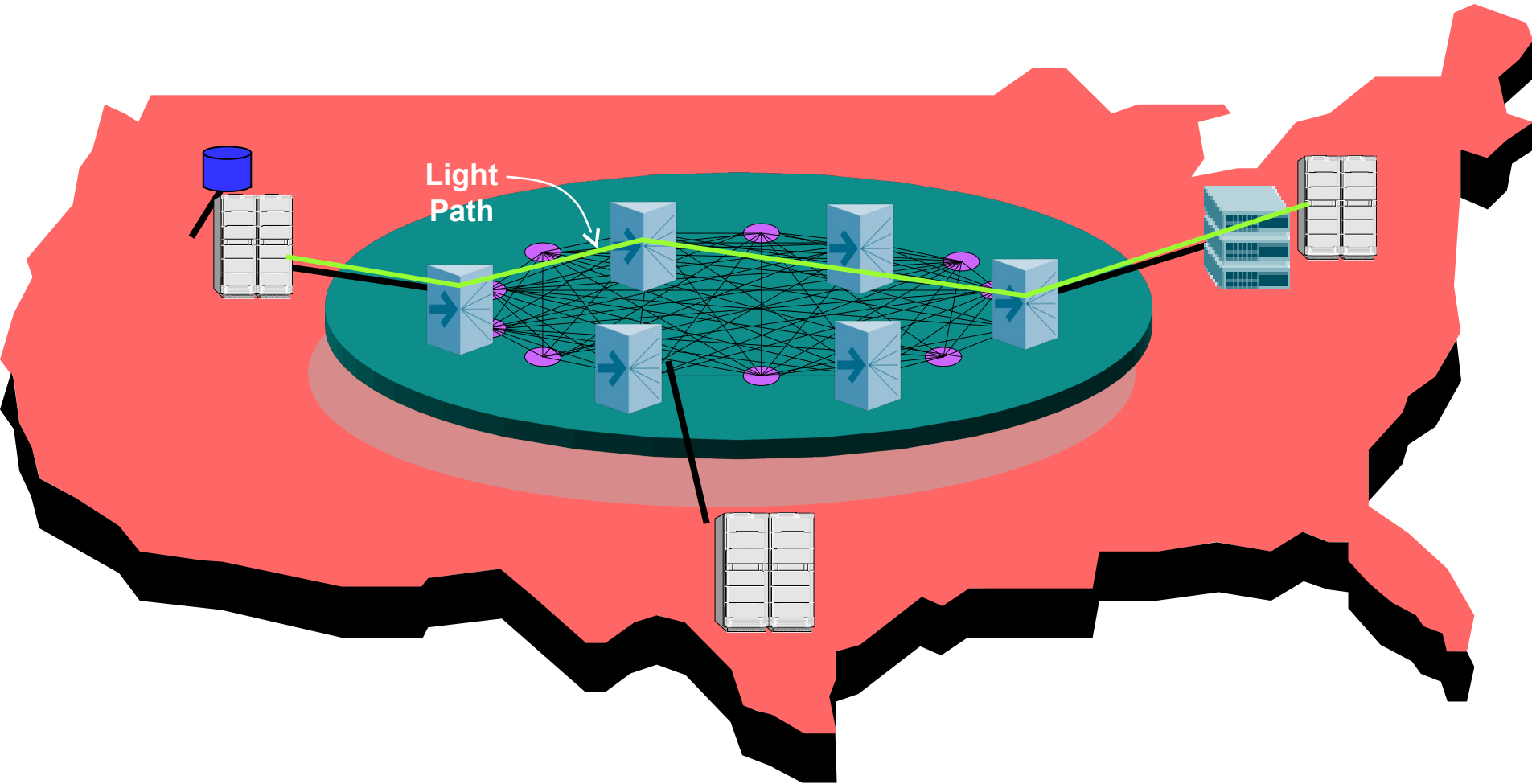
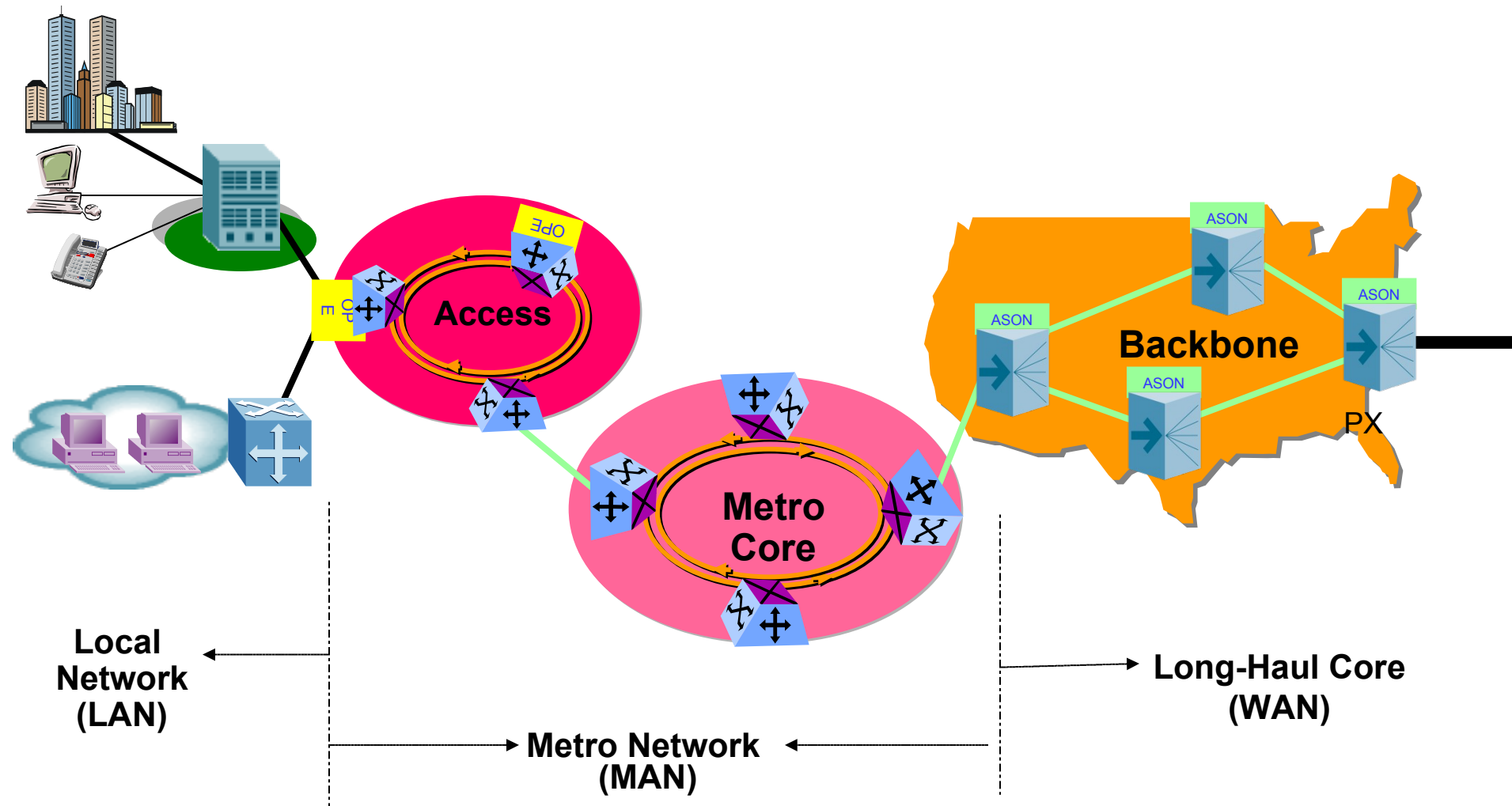**Operational implementation will take some time**

# Optical components

# Internet Reality



Access     Metro     Long Haul     Metro     Access

# OVPN on Optical Network



Light Path

# Three networks in The Internet



Access

Metro Core

Backbone

ASON

PX

Local Network (LAN)

Metro Network (MAN)

Long-Haul Core (WAN)

# Data Transport Connectivity

## Packet Switch

**data-optimized**

    Ethernet

    TCP/IP

**Network use**

    LAN

**Advantages**

    Efficient

    Simple

    Low cost

**Disadvantages**

    Unreliable

## Circuit Switch

**Voice-oriented**

    SONET

    ATM

**Network uses**

    Metro and Core

**Advantages**
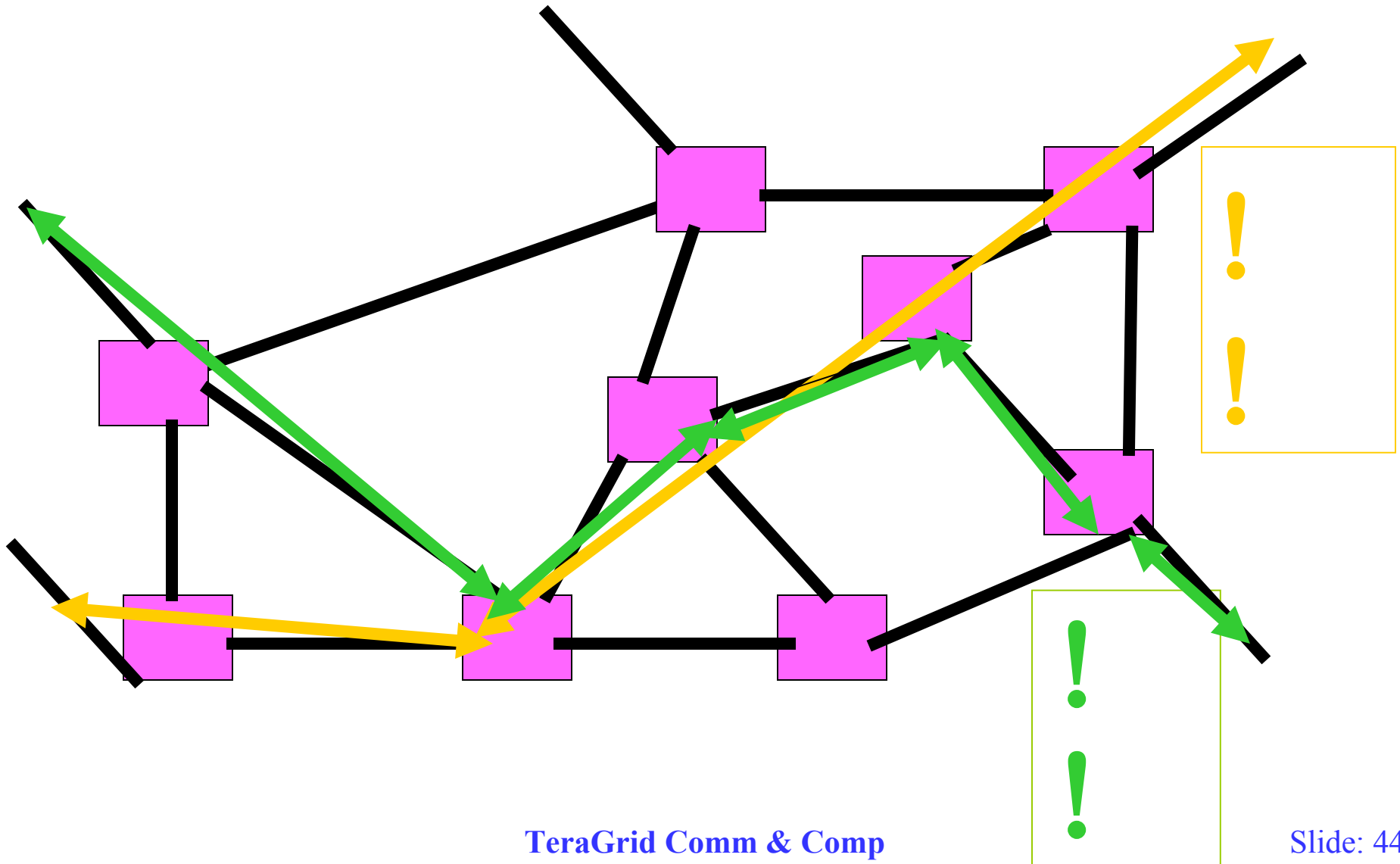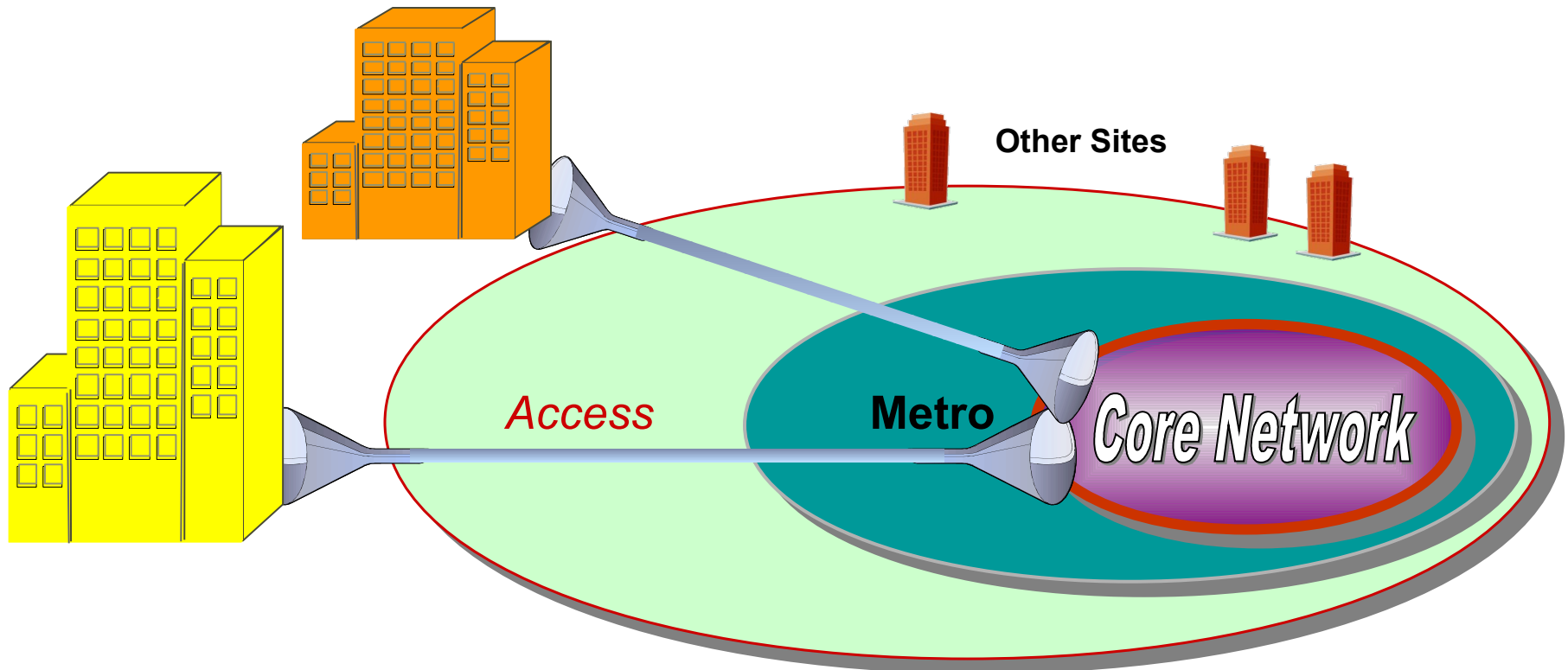
    Reliable

**Disadvantages**

    Complicate

    High cost

**Efficiency ? Reliability**

# Global Lambda Grid - Photonic Switched Network

# The Metro Bottleneck



| **End User** | **Access** | **Metro** | **Core** |
|---|---|---|---|
| Ethernet LAN | DS1<br>DS3 | OC-12<br>OC-48<br>OC-192 | OC-192<br>DWDM n x<br>! |
| IP/DATA<br>1GigE | LL/FR/ATM<br>1-40Meg | 10G | 40G+ |

**TeraGrid Comm & Comp**