

ESnet Planning, Status, and Future Issues

ASCAC, August 2008

William E. Johnston,
ESnet Department Head and Senior Scientist

Joe Burrescia, General Manager

Mike Collins, Chin Guok, and Eli Dart, Engineering

Brian Tierney, Advanced Development

Jim Gagliardi, Operations and Deployment

Stan Kluz, Infrastructure and ECS

Mike Helm, Federated Trust

Dan Peterson, Security Officer

Gizella Kapus, Business Manager

and the rest of the ESnet Team

Energy Sciences Network

Lawrence Berkeley National Laboratory

wej@es.net, this talk is available at www.es.net

Networking for the Future of Science



DOE Office of Science and ESnet – the ESnet Mission

- ESnet is an Office of Science (“SC”) facility in the Office of Advanced Scientific Computing Research (“ASCR”)
- **ESnet’s primary mission is to enable the large-scale science that is the mission of the Office of Science (SC) and that depends on:**
 - Sharing of massive amounts of data
 - Thousands of collaborators world-wide
 - Distributed data processing
 - Distributed data management
 - Distributed simulation, visualization, and computational steering
- In order to accomplish its mission ESnet provides high-speed networking and various collaboration services to Office of Science laboratories
 - As well as to many other DOE programs on a cost recovery basis

ESnet Stakeholders and their Role in ESnet

- SC/ASCR Oversight of ESnet
 - High-level oversight through the budgeting process
 - Near term input is provided by weekly teleconferences between ASCR ESnet Program Manager and ESnet
 - Indirect long term input is through the process of ESnet observing and projecting network utilization of its large-scale users
 - Direct long term input is through the SC Program Offices Requirements Workshops (more later)
- Site input to ESnet
 - Short term input through many daily (mostly) email interactions
 - Long term input through bi-annual ESCC (ESnet Coordinating Committee – all of the Lab network principals) meetings
- SC science collaborators input
 - Through numerous meeting, primarily with the networks that serve the science collaborators – mostly US and European R&E networks

Talk Outline

- I.** How are SC program requirements communicated to ESnet and what are they
- II.** ESnet response to SC requirements
- III.** Re-evaluating the ESnet strategy and identifying issues for the future
- IV.** Research and development needed to secure the future

I. SC Science Program Requirements

- Requirements are determined by
 - 1) Exploring the plans and processes of the major stakeholders:
 - 1a) Data characteristics of instruments and facilities
 - What data will be generated by instruments coming on-line over the next 5-10 years (including supercomputers)?
 - 1b) Examining the future process of science
 - How and where will the new data be analyzed and used – that is, how will the process of doing science change over 5-10 years?
 - 2) Observing current and historical network traffic patterns
 - What do the trends in network patterns predict for future network needs?

(1) Exploring the plans of the major stakeholders

- Primary mechanism is SC network Requirements Workshops
- Workshop agendas and invitees are determined by the SC science Program Offices
- Two workshops per year
- Workshop schedule
 - BES (2007 – published)
 - BER (2007 – published)
 - FES (2008 – published)
 - NP (2008 – published)
 - IPCC (Intergovernmental Panel on Climate Change) special requirements (BER) (August, 2008)
 - ASCR (Spring 2009)
 - HEP (Summer 2009)
- Future workshops - ongoing cycle
 - BES, BER – 2010
 - FES, NP – 2011
 - ASCR, HEP – 2012
 - (and so on...)
- Workshop reports: <http://www.es.net/hypertext/requirements.html>

Major Facilities Examined

- Some of these are done outside (in addition to) the Requirements Workshops
- Advanced Scientific Computing Research (ASCR)
 - NERSC (supercomputer center) (LBNL)
 - NLCF (supercomputer center) (ORNL)
 - ACLF (supercomputer center) (ANL)
- Basic Energy Sciences (BES)
 - Advanced Light Sources
 - Macromolecular Crystallography
 - Chemistry/Combustion
 - Spallation Neutron Source (ORNL)
- Biological and Environmental (BER)
 - Bioinformatics/Genomics
 - Climate Science
 - IPCC
- Fusion Energy Sciences (FE)
 - Magnetic Fusion Energy/ITER
- High Energy Physics (HEP)
 - LHC (Large Hadron Collider, CERN), Tevatron (FNAL)
- Nuclear Physics (NP)
 - RHIC (Relativistic Heavy Ion Collider) (BNL)
- These are representative of the data generating ‘hardware infrastructure’ of DOE science

Requirements from Instruments and Facilities

- ***Bandwidth***
 - Adequate network capacity to ensure timely movement of data produced by the facilities
- ***Connectivity***
 - Geographic reach sufficient to connect users and analysis systems to SC facilities
- ***Services***
 - Guaranteed bandwidth, traffic isolation, end-to-end monitoring
 - Network ***service delivery architecture***
 - SOA / Grid / “Systems of Systems”

Requirements from Instruments and Facilities - Services

- Fairly consistent requirements are found across the large-scale sciences
- **Large-scale science uses distributed systems** in order to:
 - Couple existing pockets of code, data, and expertise into “systems of systems”
 - Break up the task of massive data analysis into elements that are physically located where the data, compute, and storage resources are located
- Such systems
 - are **data intensive and high-performance**, typically moving terabytes a day for months at a time
 - are **high duty-cycle**, operating most of the day for months at a time in order to meet the requirements for data movement
 - **are widely distributed** – typically spread over continental or inter-continental distances
 - **depend on network performance and availability**, but these characteristics cannot be taken for granted, even in well run networks, when the multi-domain network path is considered
- The system elements must be able to get **guarantees** from the network that there is adequate bandwidth to accomplish the task at hand
- The systems must be able to **get information from the network** that allows graceful failure and auto-recovery and adaptation to unexpected network conditions that are short of outright failure

See, e.g., [ICFA SCI@]

The International Collaborators of DOE's Office of Science, Drives ESnet Design for International Connectivity



Most of ESnet's traffic (>85%) goes to and comes from outside of ESnet. This reflects the highly collaborative nature of large-scale science (which is one of the main focuses of DOE's Office of Science).

◆ = the R&E source or destination of ESnet's top 100 sites (all R&E)
(the DOE Lab destination or source of each flow is not shown)

Aside



- At present, ESnet traffic is dominated by data flows from large instruments – LHC, RHIC, Tevatron, etc.
- Supercomputer traffic is a small part of ESnet's total traffic, though it has the potential to increase dramatically
 - However not until appropriate system architectures are in place to allow high-speed communication among supercomputers

Other Requirements

- Assistance and services are needed for smaller user communities that have significant difficulties using the network for bulk data transfer
 - Part of the problem here is that WAN network environments (such as the combined US and European R&E networks) are large, complex systems like supercomputers and you cannot expect to get high performance when using this “system” in a “trivial” way – this is especially true for transferring a lot of data over distances $> 1000\text{km}$

Science Network Requirements Aggregation Summary

Science Drivers Science Areas / Facilities	End2End Reliability	Near Term End2End Band width	5 years End2End Band width	Traffic Characteristics	Network Services
ASCR: ALCF	-	10Gbps	30Gbps	<ul style="list-style-type: none"> • Bulk data • Remote control • Remote file system sharing 	<ul style="list-style-type: none"> • Guaranteed bandwidth • Deadline scheduling • PKI / Grid
ASCR: NERSC	-	10Gbps	20 to 40 Gbps	<ul style="list-style-type: none"> • Bulk data • Remote control • Remote file system sharing 	<ul style="list-style-type: none"> • Guaranteed bandwidth • Deadline scheduling • PKI / Grid
ASCR: NLCF	-	Backbone Bandwidth Parity	Backbone Bandwidth Parity	<ul style="list-style-type: none"> • Bulk data • Remote control • Remote file system sharing 	<ul style="list-style-type: none"> • Guaranteed bandwidth • Deadline scheduling • PKI / Grid
BER: Climate	-	3Gbps	10 to 20Gbps	<ul style="list-style-type: none"> • Bulk data • Rapid movement of GB sized files • Remote Visualization 	<ul style="list-style-type: none"> • Collaboration services • Guaranteed bandwidth • PKI / Grid
BER: EMSL/Bio	-	10Gbps	50-100Gbps	<ul style="list-style-type: none"> • Bulk data • Real-time video • Remote control 	<ul style="list-style-type: none"> • Collaborative services • Guaranteed bandwidth
BER: JGI/Genomics	-	1Gbps	2-5Gbps	<ul style="list-style-type: none"> • Bulk data 	<ul style="list-style-type: none"> • Dedicated virtual circuits • Guaranteed bandwidth

Science Network Requirements Aggregation Summary

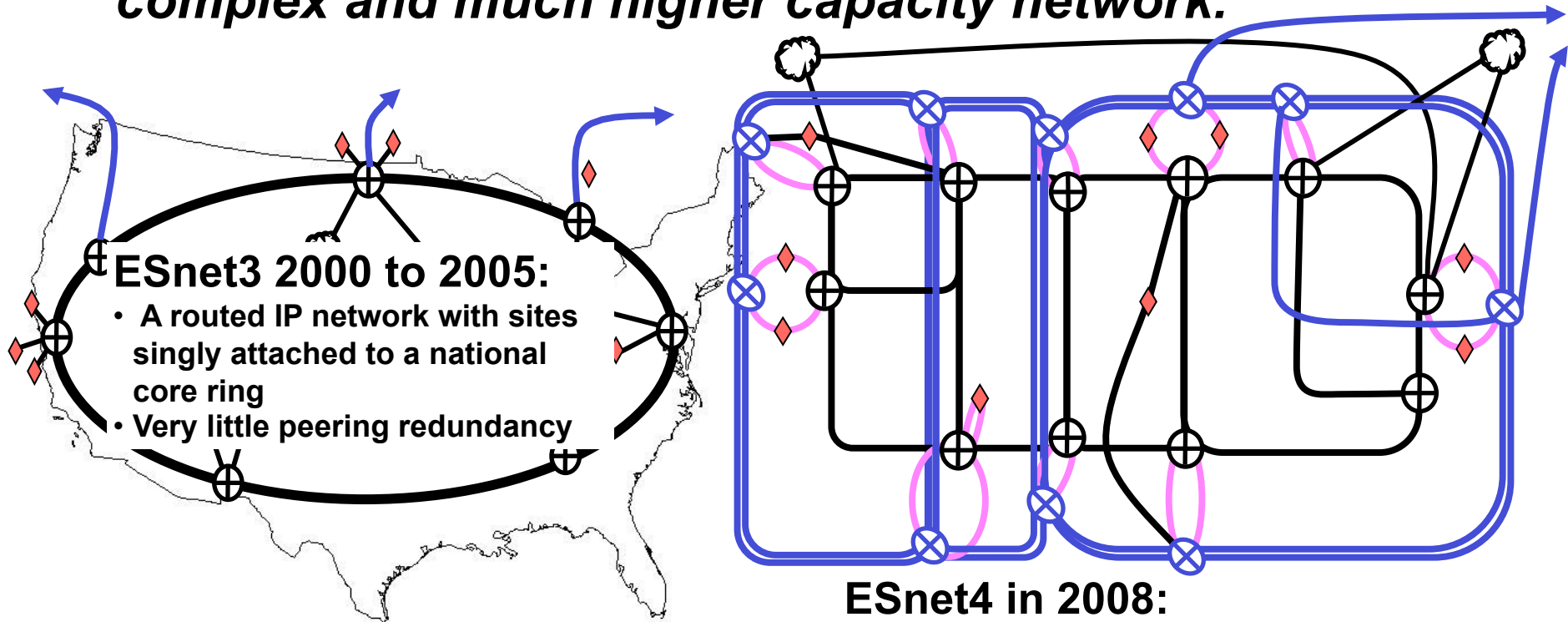
Science Drivers Science Areas / Facilities	End2End Reliability	Near Term End2End Band width	5 years End2End Band width	Traffic Characteristics	Network Services
BES: Chemistry and Combustion	-	5-10Gbps	30Gbps	<ul style="list-style-type: none"> • Bulk data • Real time data streaming 	<ul style="list-style-type: none"> • Data movement middleware
BES: Light Sources	-	15Gbps	40-60Gbps	<ul style="list-style-type: none"> • Bulk data • Coupled simulation and experiment 	<ul style="list-style-type: none"> • Collaboration services • Data transfer facilities • Grid / PKI • Guaranteed bandwidth
BES: Nanoscience Centers	-	3-5Gbps	30Gbps	<ul style="list-style-type: none"> • Bulk data • Real time data streaming • Remote control 	<ul style="list-style-type: none"> • Collaboration services • Grid / PKI
FES: International Collaborations	-	100Mbps	1Gbps	<ul style="list-style-type: none"> • Bulk data 	<ul style="list-style-type: none"> • Enhanced collaboration services • Grid / PKI • Monitoring / test tools
FES: Instruments and Facilities	-	3Gbps	20Gbps	<ul style="list-style-type: none"> • Bulk data • Coupled simulation and experiment • Remote control 	<ul style="list-style-type: none"> • Enhanced collaboration service • Grid / PKI
FES: Simulation	-	10Gbps	88Gbps	<ul style="list-style-type: none"> • Bulk data • Coupled simulation and experiment • Remote control 	<ul style="list-style-type: none"> • Easy movement of large checkpoint files • Guaranteed bandwidth • Reliable data transfer

Science Network Requirements Aggregation Summary

Science Drivers Science Areas / Facilities	End2End Reliability	Near Term End2End Band width	5 years End2End Band width	Traffic Characteristics	Network Services
Immediate Requirements and Drivers for ESnet4					
HEP: LHC (CMS and Atlas)	99.95+% (Less than 4 hours per year)	73Gbps	225-265Gbps	<ul style="list-style-type: none"> • Bulk data • Coupled analysis workflows 	<ul style="list-style-type: none"> • Collaboration services • Grid / PKI • Guaranteed bandwidth • Monitoring / test tools
NP: CMS Heavy Ion	-	10Gbps (2009)	20Gbps	<ul style="list-style-type: none"> • Bulk data 	<ul style="list-style-type: none"> • Collaboration services • Deadline scheduling • Grid / PKI
NP: CEBF (JLAB)	-	10Gbps	10Gbps	<ul style="list-style-type: none"> • Bulk data 	<ul style="list-style-type: none"> • Collaboration services • Grid / PKI
NP: RHIC	Limited outage duration to avoid analysis pipeline stalls	6Gbps	20Gbps	<ul style="list-style-type: none"> • Bulk data 	<ul style="list-style-type: none"> • Collaboration services • Grid / PKI • Guaranteed bandwidth • Monitoring / test tools

II. ESnet Response to the Requirements

- ESnet4 was *built to address specific Office of Science program requirements. The result is a much more complex and much higher capacity network.*



- The new Science Data Network (blue) is a switched network providing guaranteed bandwidth for large data movement
- All large science sites are dually connected on metro area rings or dually connected directly to core ring for reliability
- Rich topology increases the reliability of the network

New ESnet Services

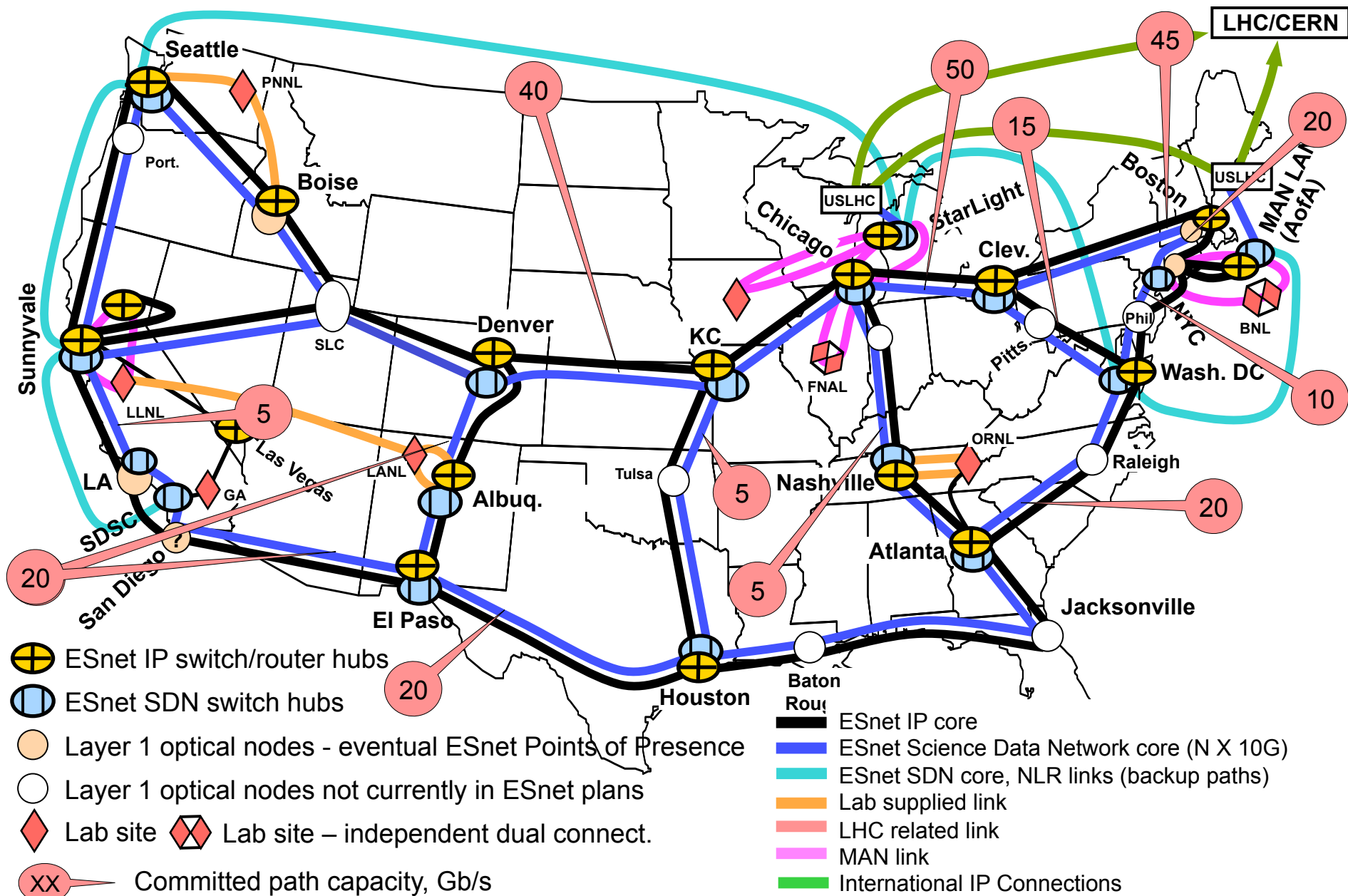
- Virtual circuit service providing schedulable bandwidth guarantees, traffic isolation, etc
 - ESnet OSCARS service
 - <http://www.es.net/OSCARS/index.html>
 - Successfully deployed in early production today
 - Additional R&D is needed in many areas of this service
- Assistance for smaller communities in using the network for bulk data transfer
 - fasterdata.es.net – web site devoted to information on bulk data transfer, host tuning, etc. established
 - Other potential approaches
 - Various latency insensitive forwarding devices in the network (R&D)

Building the Network as Opposed to Planning the Budget

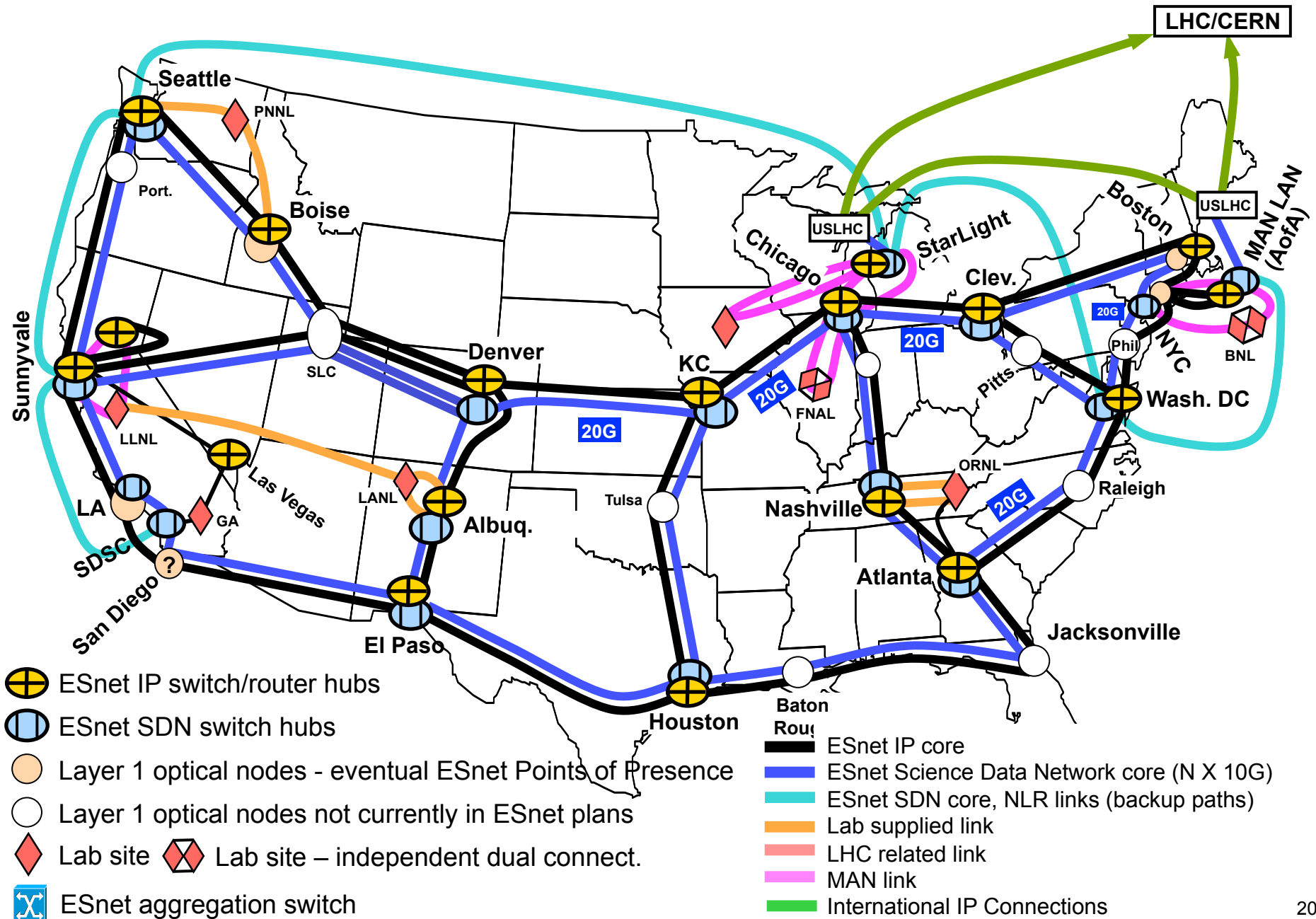
- Aggregate capacity requirements like those above indicate how to budget for a network but do not tell you how to build a network
- To actually build a network you have to look at where the traffic originates and ends up and how much traffic is expected on specific paths
- So far we have specific bandwidth and path (collaborator location) information for
 - LHC (CMS, CMS Heavy Ion, Atlas)
 - SC Supercomputers
 - CEBF/JLab
 - RHIC/BNL

this specific information has lead to the current and planned configuration of the network for the next several years

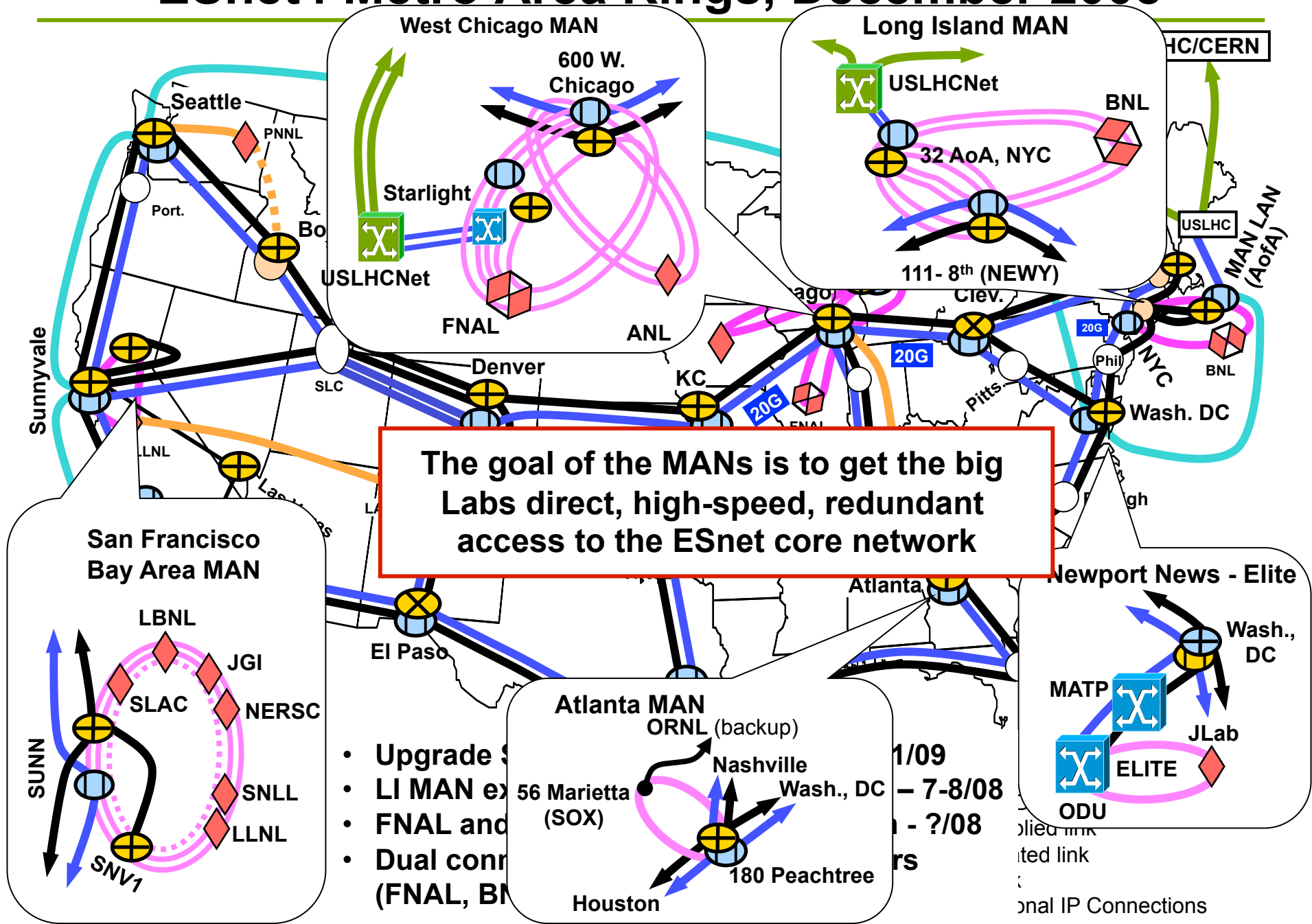
How do the Bandwidth – Path Requirements Map to the Network? (Core Network Planning - 2010)



ESnet 4 Core Network – December 2008



ESnet4 Metro Area Rings, December 2008

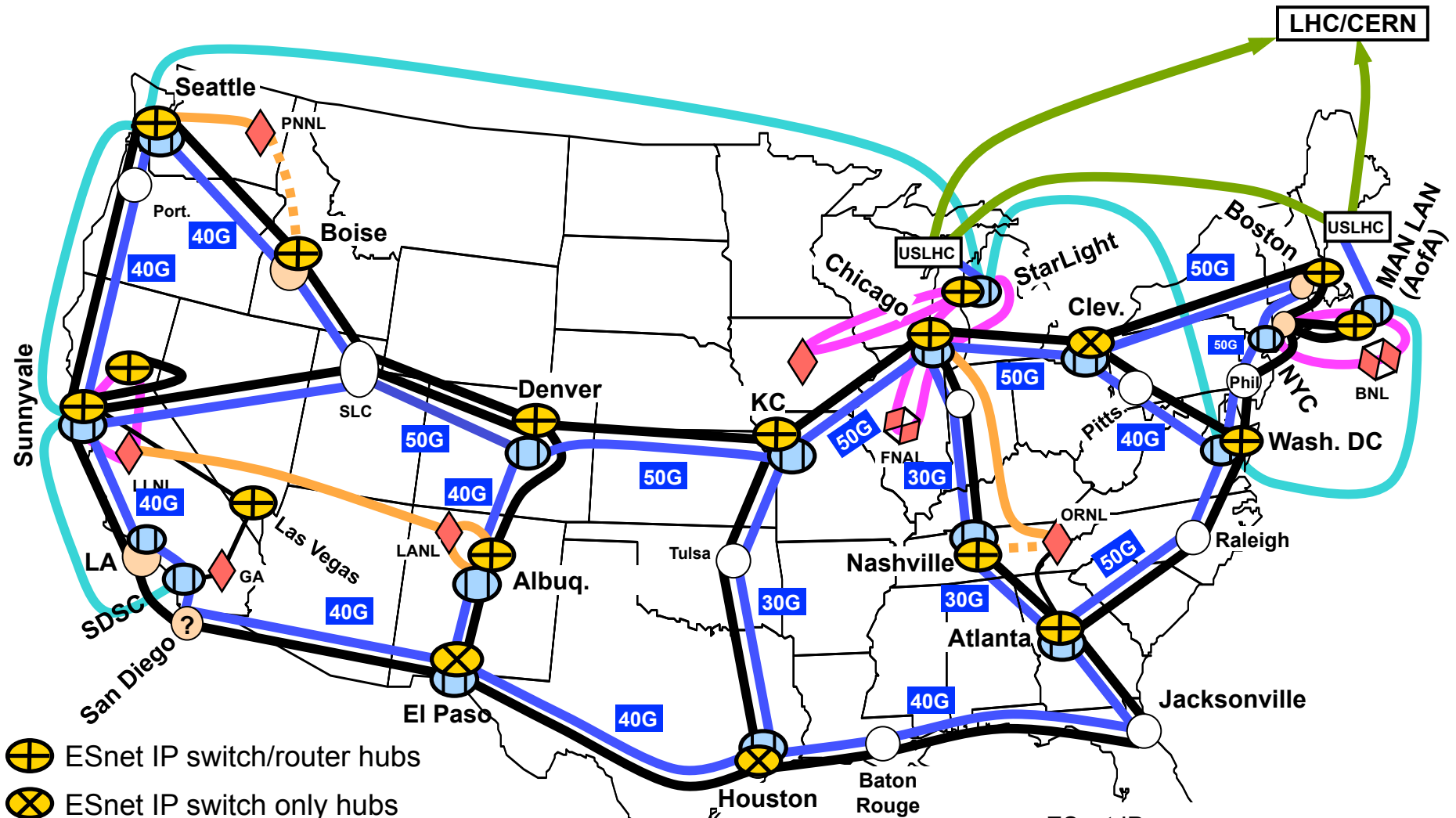


The goal of the MANs is to get the big Labs direct, high-speed, redundant access to the ESnet core network

- Upgrade S
 - LI MAN ex 56 Marietta (SOX)
 - FNAL and
 - Dual conn (FNAL, BI
- 1/09
- 7-8/08
- ?/08
rs

onal IP Connections

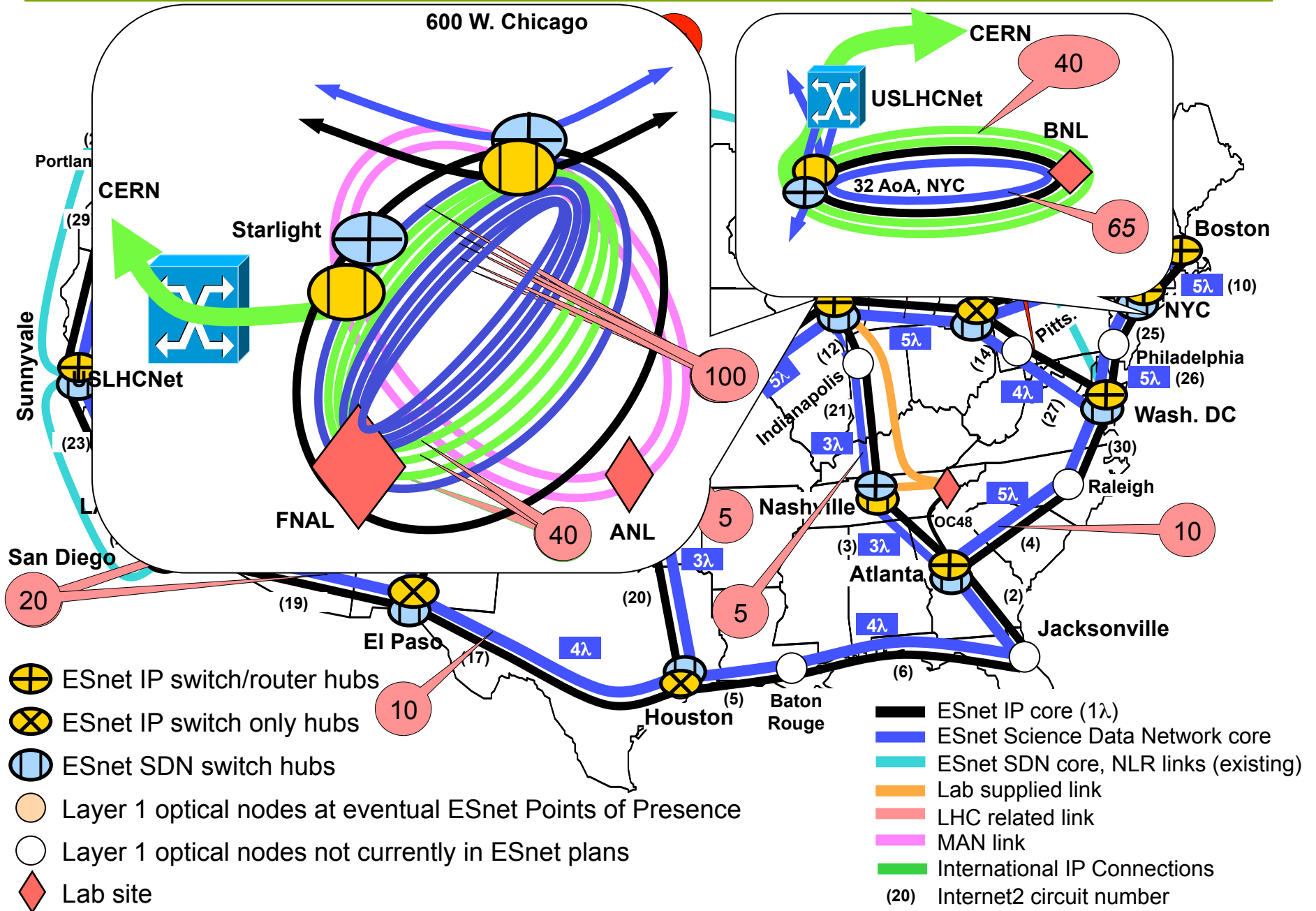
ESnet 4 As Planned for 2010



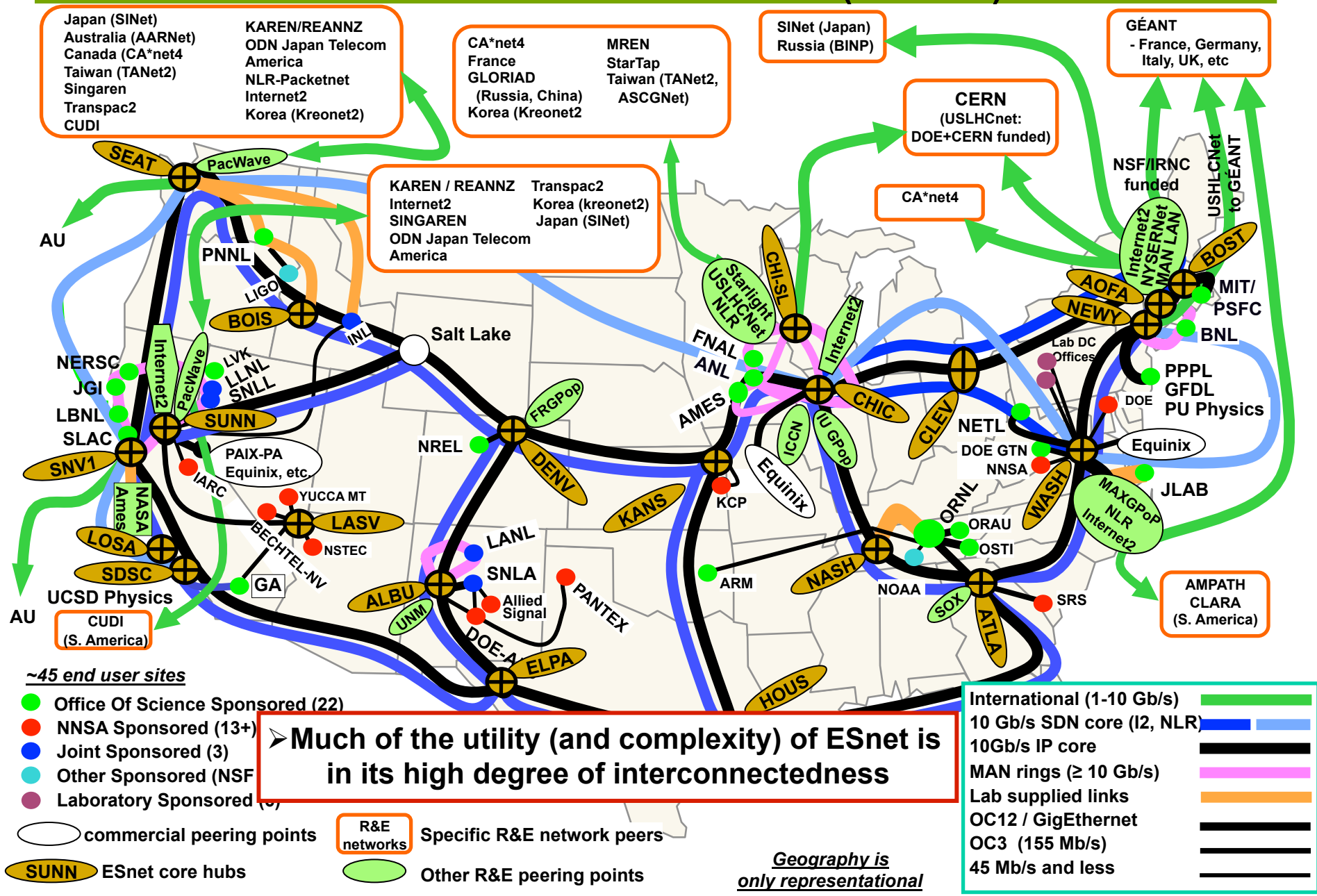
This growth in the network capacity is based on the current ESnet 5 yr. budget as submitted by SC/ASCR to OMB

- ◆ Lab site ◆⊗ Lab site – independent dual connect.
- MAN link
- International IP Connections

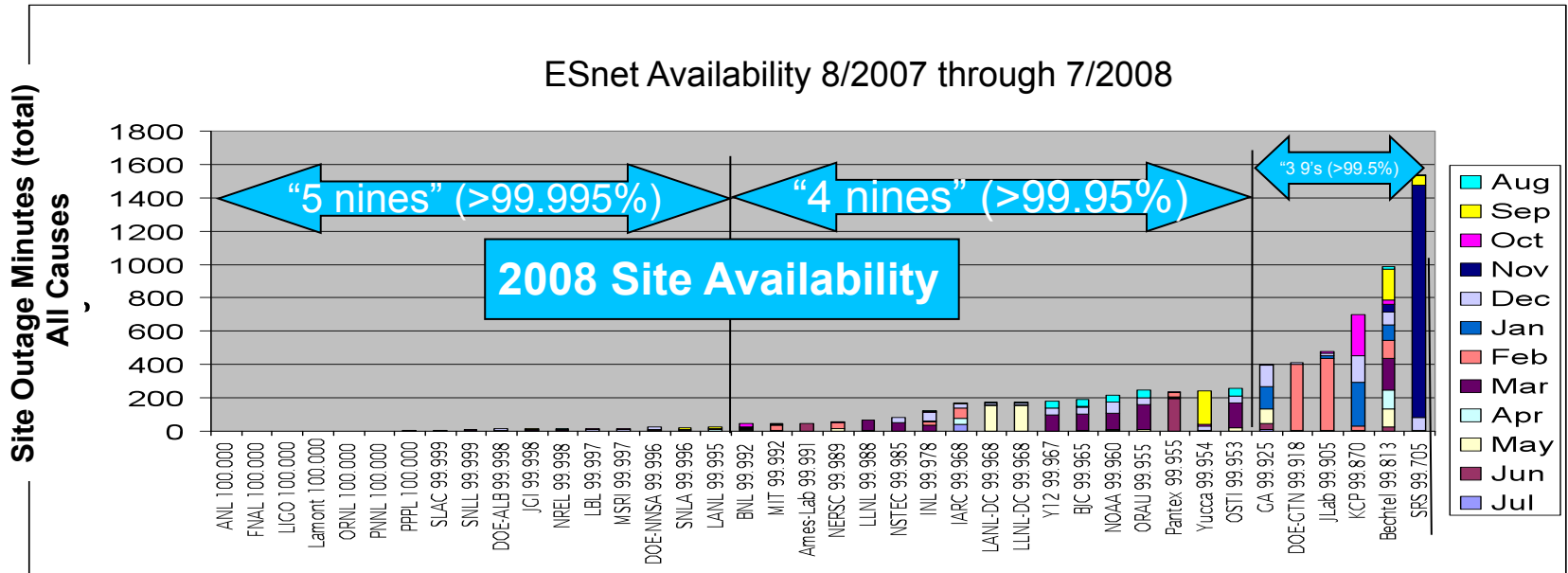
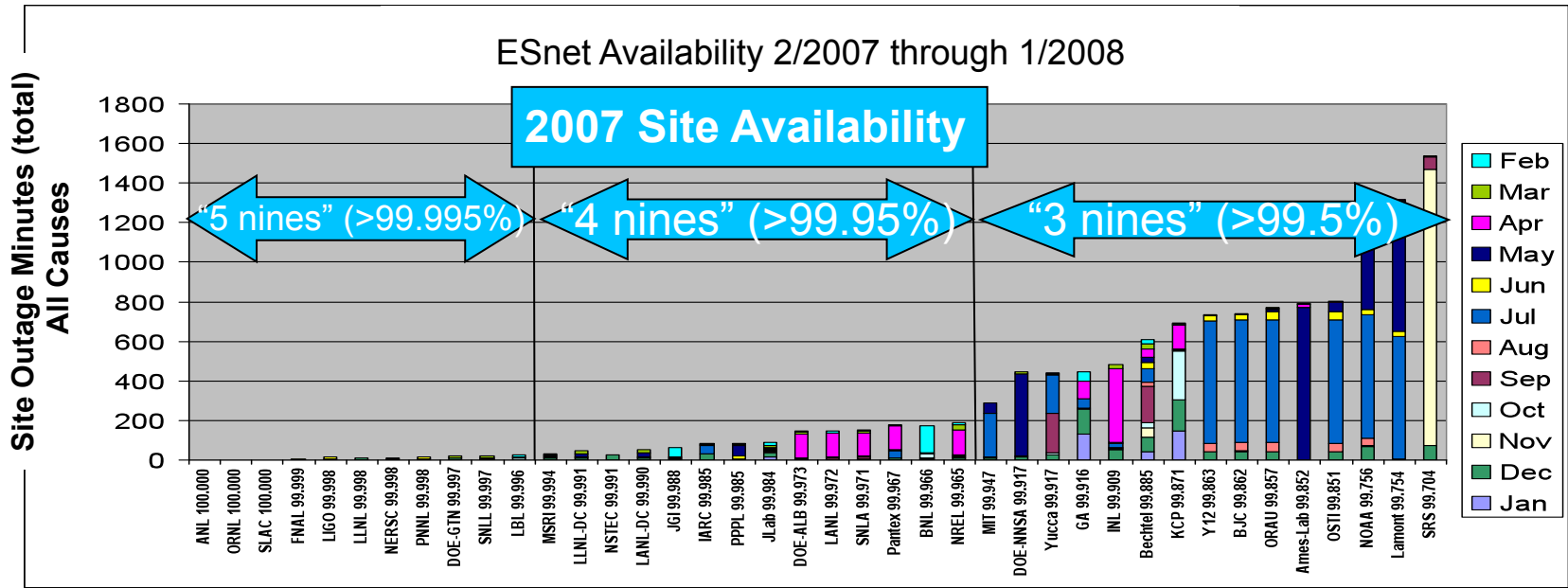
MAN Capacity Planning - 2010



ESnet Provides Global High-Speed Internet Connectivity for DOE Facilities and Collaborators (12/2008)



One Consequence of ESnet's New Architecture is that Site Availability is Increasing



III. Re-evaluating the Strategy and Identifying Issues

- The current strategy (that lead to the ESnet4, 2012 plans) was developed primarily as a result of the information gathered in the 2003 and 2003 network workshops, and their updates in 2005-6 (including LHC, climate, RHIC, SNS, Fusion, the supercomputers, and a few others) [workshops]
- So far the more formal requirements workshops have largely reaffirmed the ESnet4 strategy developed earlier
- ***However – is this the whole story?***

Where Are We Now?

How do the science program identified requirements compare to the network capacity planning?

- The current network is built to accommodate the known, path-specific needs of the programs
- However this is not the whole picture: The **core path capacity planning** (see map above) **so far only accounts for 405 Gb/s out of 789 Gb/s identified** aggregate requirements provided by the science programs

Synopsis of "Science Network Requirements Aggregation Summary," 6/2008								
	5 year requirements			Accounted for in current ESnet path planning			Unacc'ted for	
Requirements (aggregate Gb/s)	789			405			384	
ESnet Planned Aggregate Capacity (Gb/s) Based on 5 yr. Budget								
	2006	2007	2008	2009	2010	2011	2012	2013
ESnet "aggregate"	57.50	192	192	842	1442	1442	1442	2042

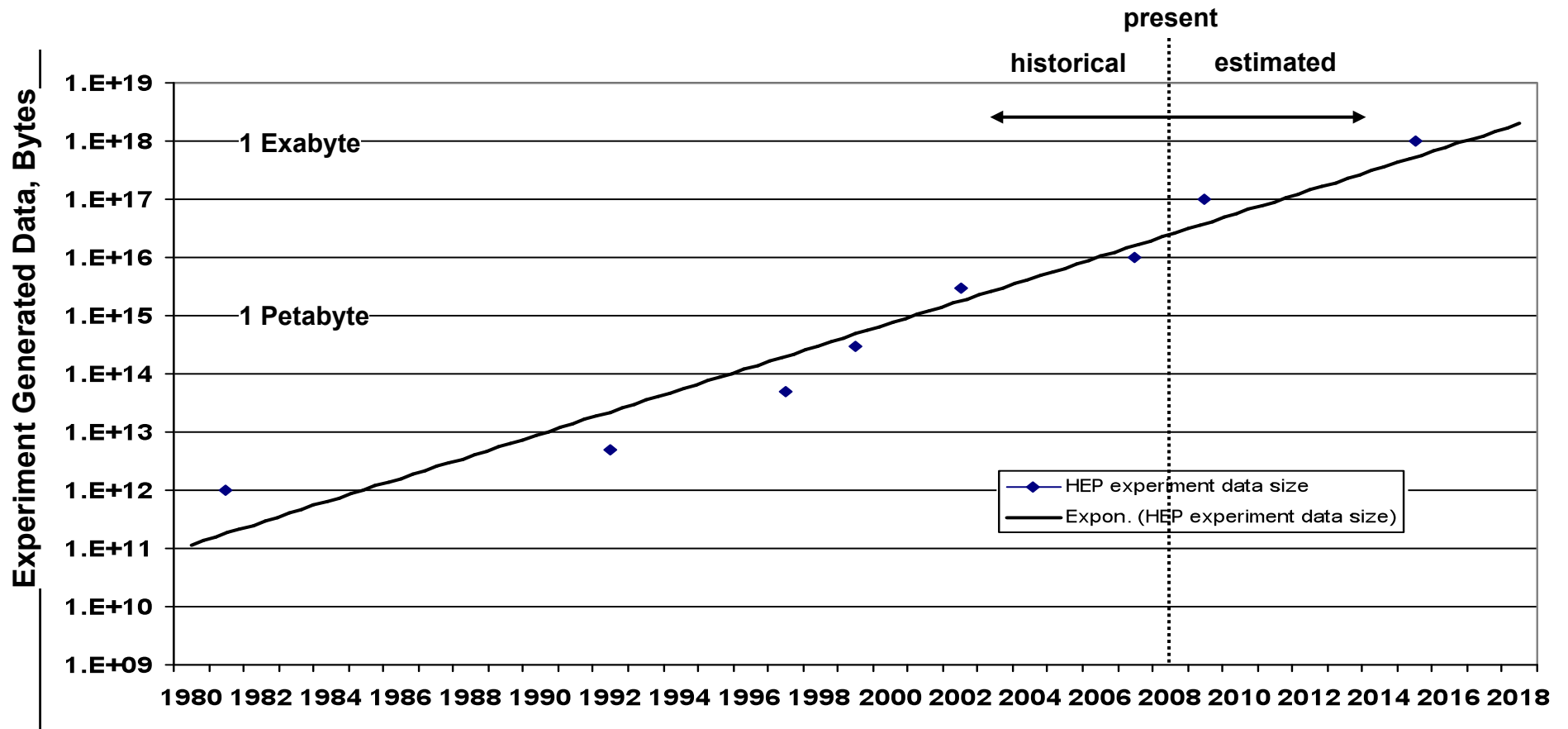
- The **planned aggregate capacity growth of ESnet matches the known requirements**
 - The **"extra" capacity** indicated above **is needed to account for the fact that there is much less than complete flexibility in mapping specific path requirements to the aggregate capacity planned network** and we won't know specific paths until several years into building the network
- Whether this approach works is TBD, but indications are that it probably will

Is ESnet Planned Capacity Adequate? E.g. for LHC? (Maybe So, Maybe Not)

- Several Tier2 centers (mostly at Universities) are capable of 10Gbps now
 - Many Tier2 sites are building their local infrastructure to handle 10Gbps
 - We won't know for sure what the "real" load will look like until the testing stops and the production analysis begins
 - ***Scientific productivity will follow high-bandwidth access to large data volumes ⇒ incentive for others to upgrade***
- Many Tier3 sites are also building 10Gbps-capable analysis infrastructures – ***this was not in LHC plans a year ago***
 - Most Tier3 sites do not yet have 10Gbps of network capacity
 - It is likely that this will cause a "second onslaught" in 2009 as the Tier3 sites all upgrade their network capacity to handle 10Gbps of LHC traffic
 - ***It is possible that the USA installed base of LHC analysis hardware will consume significantly more network bandwidth than was originally estimated***
 - N.B. Harvey Newman (HEP, Caltech) predicted this eventuality years ago

Reexamining the Strategy: The Exponential Growth of HEP Data is “Constant”

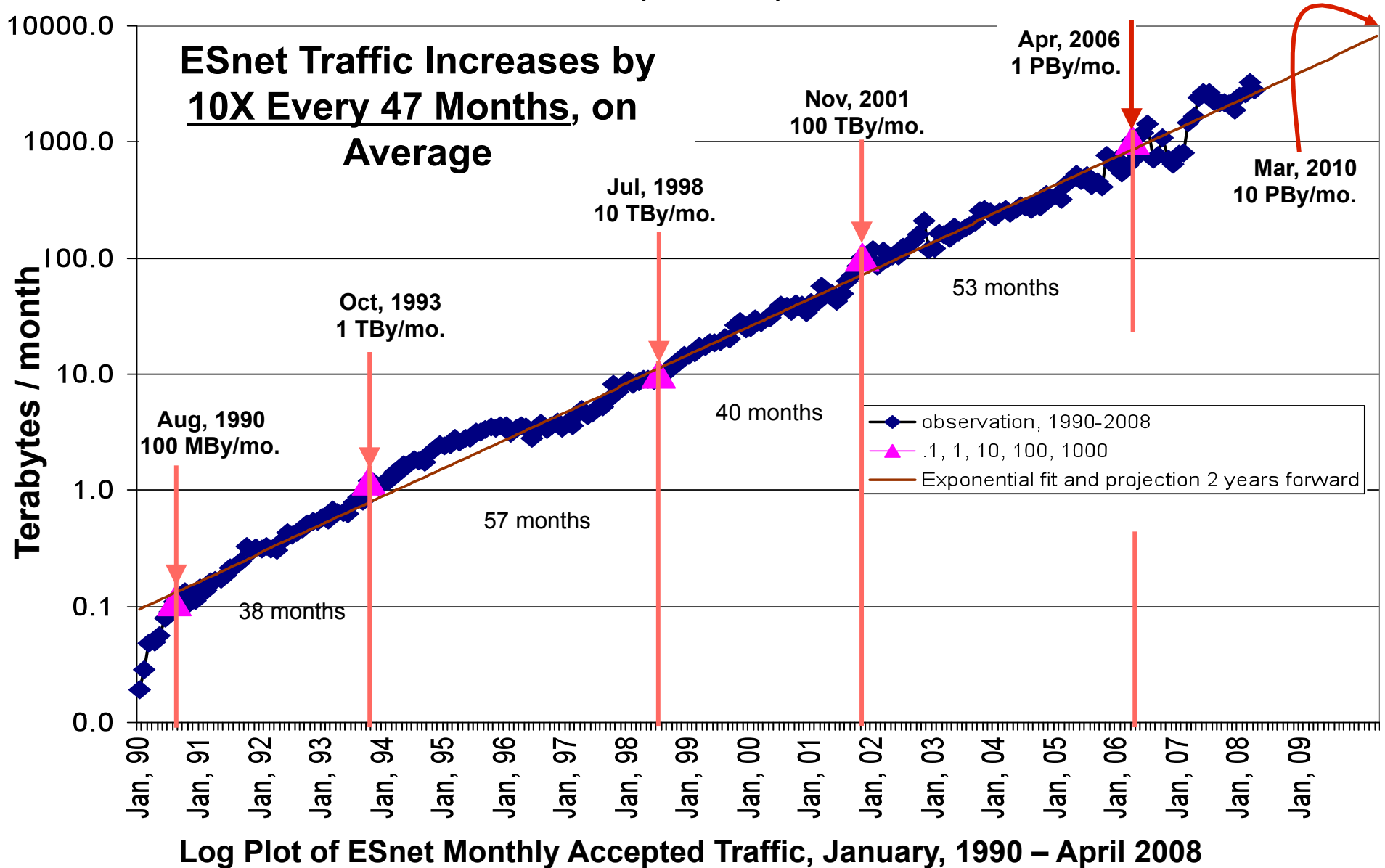
For a point of “ground truth” consider the historical growth of the size of HEP data sets – The trends as typified by the FNAL traffic will continue.



Data courtesy of Harvey Newman, Caltech,
and Richard Mount, SLAC

Reexamining the Strategy

- Consider network traffic patterns – “ground truth”
 - What do the trends in network patterns predict for future network needs



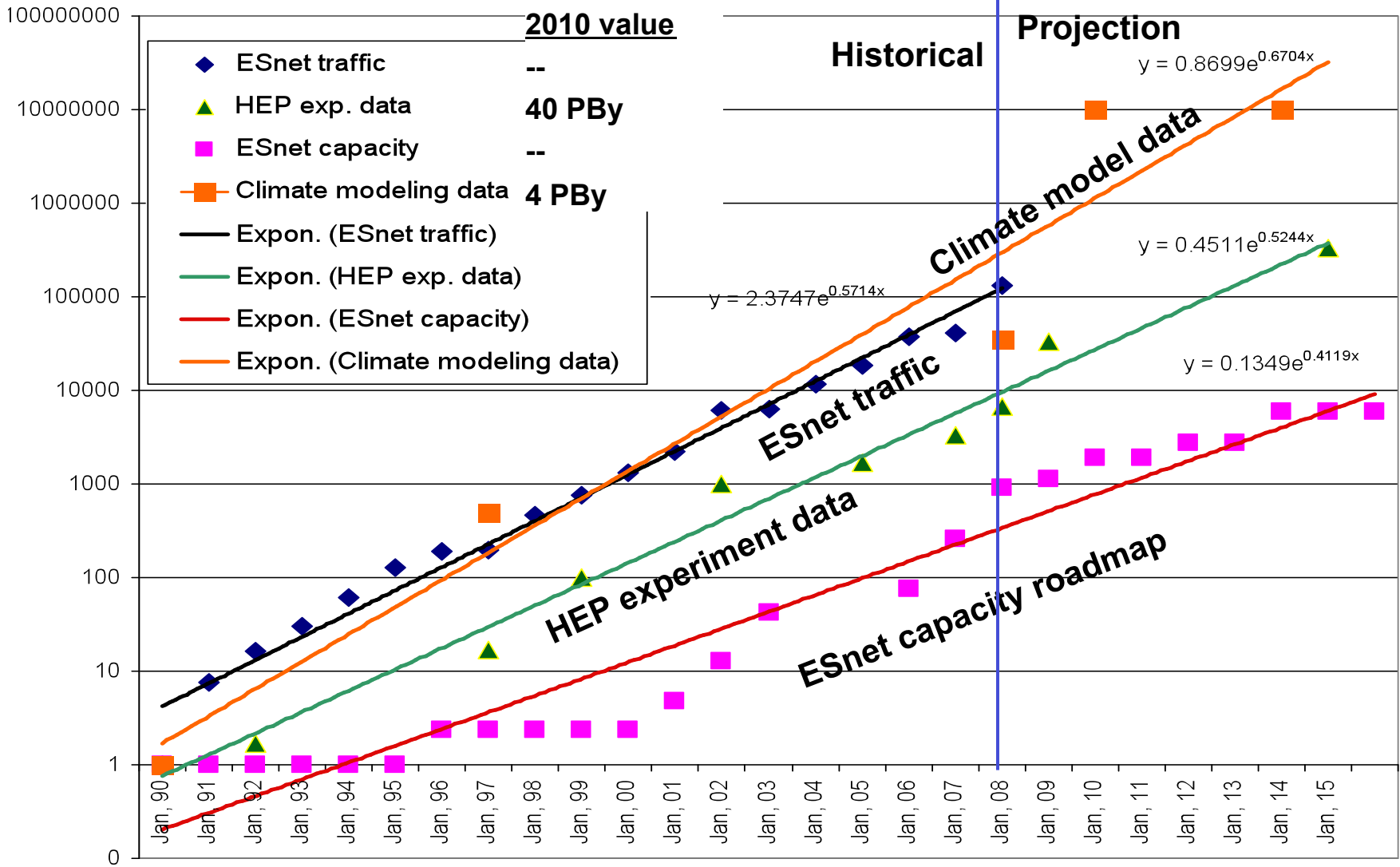
Where Will the Capacity Increases Come From?

- ESnet4 planning assumes technology advances will provide 100Gb/s optical waves (they are 10 Gb/s now) which gives a potential 5000 Gb/s core network by 2012
 - The ESnet4 SDN switching/routing platform is designed to support ***new 100Gb/s network interfaces***
 - With ***capacity planning based on the ESnet 2010 wave count, together with some considerable reservations about the affordability of 100 Gb/s network interfaces***, we can probably assume some fraction of the 5000 Gb/s of potential core network capacity by 2012 depending on the cost of the equipment – perhaps 20% – about 1000-2000 Gb/s of ***aggregate capacity***
- **Is this adequate to meet future needs?**
- Not Necessarily!**

Network Traffic, Physics Data, and Network Capacity

Ignore the units of the quantities being graphed they are normalized to 1 in 1990, just look at the long-term trends: **All of the "ground truth" measures are growing significantly faster than ESnet projected capacity**

All Three Data Series are Normalized to "1" at Jan. 1990



➤ Issues for the Future Network

- The current estimates from the LHC experiments and the supercomputer centers **have the currently planned ESnet 2011 wave configuration operating at capacity** and there are several other major sources that will be generating significant data in that time frame (e.g. Climate)
- The significantly **higher exponential growth of traffic** (total accepted bytes) **vs. total capacity** (aggregate core bandwidth) means **traffic will eventually overwhelm the capacity** – “when” cannot be directly deduced from aggregate observations, but if you add this fact
 - Nominal average load on busiest backbone paths in June 2006 was ~1.5 Gb/s - In 2010 average load will be ~15 Gbps based on current trends and 150 Gb/s in 2014

My (wej) guess is that capacity problems will start to occur by 2015-16 without new technology approaches

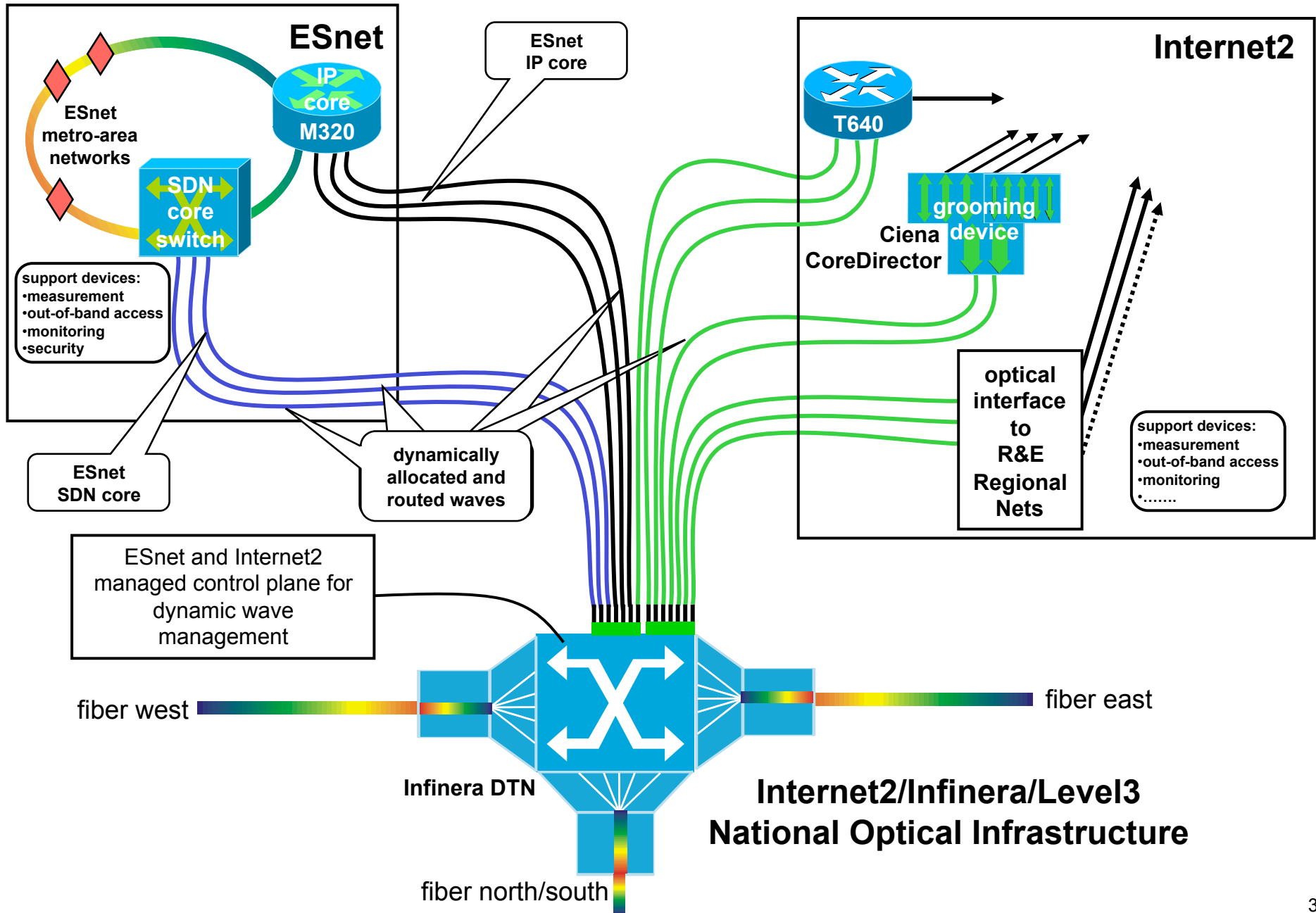
Issues for the Future Network

- The “casual” increases in overall network capacity based on straightforward commercial channel capacity that have sufficed in the past are less likely to easily meet future needs due to the (potential) un-affordability of the hardware
 - the few existing examples of >10G/s interfaces are ~10x more expensive than the 10G interfaces (~\$500K each – not practical)

Where Do We Go From Here?

- The Internet2-ESnet partnership optical network is build on dedicated fiber and optical equipment
 - The current optical network is configured with 10 × 10G waves / fiber path and more waves will be added in groups of 10 up to 80 waves
- The **current wave transport topology is essentially static** or only manually configured - our current network infrastructure of routers and switches assumes this
- With completely flexible traffic management extending down to the optical transport level we **should be able to extend the life of the current infrastructure** by moving significant parts of the capacity to the specific routes where it is needed
- We must integrate the optical transport with the “network” and provide for dynamism / route flexibility at the optical level in order to make optimum use of the available capacity

Internet2 and ESnet Optical Node in the Future



IV. Research and Development Needed to Secure the Future

- In order for “R&D” to be useful to ESnet it must be “directed R&D” – that is, R&D that has ESnet as a partner ***so that the result is deployable in the production network where there are many constraints arising out of operational requirements***
 - Typical undirected R&D either produces interesting results that are un-deployable in a production network or that have to be reimplemented in order to be deployable

Research and Development Needed to Secure the Future: Approach to R&D

- Partnership R&D is a successful “directed R&D” approach that is used with ESnet’s OSCARS virtual circuit system that provides bandwidth reservations and integrated layer 2/3 network management
 - OSCARS is a partnership between ESnet, Internet2 (university network), USC/ISI, and several European network organizations – because of this it has been successfully deployed in several large R&E networks

Research and Development Needed to Secure the Future: Approach to R&D

- OSCARS ...
 - DOE has recently informed ESnet that funding for OSCARS R&D will end with this year – presumably because their assessment is that research is “done” and the R&D program will not fund development
 - This is a persistent problem in the DOE R&D programs and is clearly described in the **ASCAC networking report** [ASCAC, Stechel and Wing]

“In particular, ASCR needs to establish processes to review networking research results, as well as to select and fund promising capabilities for further development, with the express intent to accelerate the availability of new capabilities for the science community.

Research and Development Needed to Secure the Future: *Example Needed R&D*

- To best utilize the total available capacity we must integrate the optical (L1) transport with the “network” (L2 and L3) and provide for dynamism / route flexibility at all layers
 - The L1 control plane manager approach currently being considered is based on an extended version of the OSCARS dynamic circuit manager – but a good deal of R&D is needed for the integrated L1/2/3 dynamic route management
 - For this – or any such new approach to routing – to be successfully (and safely) introduced into the production network it will first have to be developed and extensively tested in a testbed that has characteristics (e.g. topology and hardware) very similar to the production network

Research and Development Needed to Secure the Future:

Example Needed R&D

- It is becoming apparent that another aspect of the most effective utilization of the network requires the ability to transparently direct routed IP traffic onto SDN
 - There are only ideas in this area at the moment

Research and Development Needed to Secure the Future

- End-to-end monitoring as a service: Provide useful, comprehensive, and meaningful ***information on the state of end-to-end paths***, or potential paths, to the user –
 - perfSONAR, and associated tools, provide real time information in a form that is useful to the user (via appropriate abstractions) and that is delivered through standard interfaces that can be incorporated in to SOA type applications (See [E2EMON] and [TrViz].)
 - Techniques need to be developed to:
 - 1) Use “standardized” network topology from all of the networks involved in a path to give the user an appropriate view of the path
 - 2) Monitoring for virtual circuits based on the different VC approaches of the various R&E nets
 - e.g. MPLS in ESnet, VLANs, TDM/grooming devices (e.g. Ciena Core Directors), etc.,and then integrate this into a perfSONAR framework

Research and Development Needed to Secure the Future: Data Transfer Issues Other than HEP and an Approach

- Assistance and services are needed for smaller user communities that have significant difficulties using the network for bulk data transfer
- This issue cuts across several SC Science Offices
- These problems **MUST** be solved if scientists are to effectively analyze the data sets produced by petascale machines
- Consider some case studies

Data Transfer Problems – Light Source Case Study

- Light sources (ALS, APS, NSLS, etc) serve many thousands of users
 - Typical user is one scientist plus a few grad students
 - 2-3 days of beam time per year
 - Take data, then go home and analyze data
 - Data set size up to 1TB, typically 0.5TB
- Widespread frustration with network-based data transfer among light source users
 - WAN transfer tools not installed
 - Systems not tuned
 - Lack of available expertise for fixing these problems
 - Network problems at the “other end” – typically a small part of a university network
- Users copy data to portable hard drives or burn stacks of DVDs today, but data set sizes will probably exceed hard disk sizes in the near future

Data Transfer Problems – Combustion Case Study

- Combustion simulations generate large data sets
- User awarded INCITE allocation at NERSC, 10TB data set generated
- INCITE allocation awarded at ORNL → need to move data set from NERSC to ORNL
- Persistent data transfer problems
 - Lack of common toolset
 - Unreliable transfers, low performance
 - Data moved, but it took almost two weeks of babysitting the transfer

Data Transfer Problems – Fusion Case Study

- Large-scale fusion simulations (e.g. GTC) are run at both NERSC and ORNL
- Users wish to move data sets between supercomputer centers
- Data transfer performance is low, workflow software unavailable or unreliable
- Data must be moved between systems at both NERSC and ORNL
 - Move data from storage to WAN transfer resource
 - Transfer data to other supercomputer center
 - Move data to storage or onto computational platform

Proper Configuration of End Systems is Essential

- Persistent performance problems exist throughout the DOE Office of Science
 - Existing tools and technologies (e.g. TCP tuning, GridFTP) are not deployed on end systems or are inconsistently deployed across major resources
 - Performance problems impede productivity
 - Unreliable data transfers soak up scientists' time (must babysit transfers)
- Default system configuration is inadequate
 - Most system administrators don't know how to properly configure a computer for WAN data transfer
 - System administrators typically don't know where to look for the right information
 - Scientists and system administrators typically don't know that WAN data transfer can be high performance, so they don't ask for help
 - WAN transfer performance is often not a system administration priority

High Performance WAN Data Transfer is Possible

- Tools and technologies for high performance WAN data transfer exist today
 - TCP tuning documentation exists
 - Tools such as GridFTP are available and are used by sophisticated users
 - DOE has made significant contribution to these tools over the years
- Sophisticated users and programs are able to get high performance
 - User groups with the size and resources to “do it themselves” get good performance (e.g. HEP, NP)
 - Smaller groups do not have the internal staff and expertise to manage their own data transfer infrastructures, and so get low performance
- The WAN is the same in the high and low performance cases but the end system configurations are different

Data Transfer Issues Other than HEP and an Approach

- DOE/SC should task one entity with development, support and advocacy for WAN data transfer software
 - Support (at the moment GridFTP has no long-term funding)
 - Port to new architectures – we need these tools to work on petascale machines and next-generation data transfer hosts
 - Usability – scientific productivity must be the goal of these tools, so they must be made user-friendly so scientists can be scientists instead of working on data transfers
 - Consistent deployment – all major DOE facilities must deploy a common, interoperable, reliable data transfer toolkit (NERSC, ORNL, light sources, nanocenters, etc)
 - Workflow engines, GridFTP and other file movers, test infrastructure
- These problems **MUST** be solved if scientists are to effectively analyze the data sets produced by petascale machines

Research and Development Needed to Secure the Future

- Artificial (network device based) reduction of end-to-end latency as seen by the user application is needed in order to allow small, unspecialized systems (e.g. a Windows laptop) do “large” data transfers with good throughput over national and international distances
 - There are several approaches possible here and R&D is needed to determine the “right” direction
 - The answer to this may be dominated by deployment issues that are sort of outside ESnet’s realm – for example deploying data movement “accelerator” systems at user facilities such as the Light Sources and Nanotechnology Centers

New in ESnet – Advanced Technologies Group / Coordinator

- Up to this point individual ESnet engineers have worked in their “spare” time to do the R&D, or to evaluate R&D done by others, and coordinate the implementation and/or introduction of the new services into the production network environment – and they will continue to do so
- In addition to this – looking to the future – ESnet has implemented a more formal approach to investigating and coordinating the R&D for the new services needed by science
 - An ESnet Advanced Technologies Group / Coordinator has been established with a twofold purpose:
 - 1) To provide a unified view to the world of the several engineering development projects that are on-going in ESnet in order to publicize a coherent catalogue of advanced development work going on in ESnet.
 - 2) To develop a portfolio of exploratory new projects, some involving technology developed by others, and some of which will be developed within the context of ESnet.
- A highly qualified Advanced Technologies lead – Brian Tierney – has been hired and funded from current ESnet operational funding, and by next year a second staff person will be added. Beyond this, growth of the effort will be driven by new funding obtained specifically for that purpose.

Needed in ESnet – Science User Advocate

- A position within ESnet to act as a direct advocate for the needs and capabilities of the major SC science users of ESnet
 - At the moment ESnet receives new service requests and requirements in a timely way, but no one acts as an active advocate to represent the user's point of view once ESnet gets the requests
 - Also, the User Advocate can suggest changes and enhancements to services that the Advocate sees are needed to assist the science community even if the community does not make this connection on their own

➤ Summary

- Transition to ESnet4 is going smoothly
 - New network services to support large-scale science are progressing
 - OSCARS virtual circuit service is being used, and the service functionality is adapting to unforeseen user needs
 - Measurement infrastructure is rapidly becoming widely enough deployed to be very useful
- Revaluation of the 5 yr strategy indicates that the future will not be qualitatively the same as the past – and this must be addressed
 - R&D, testbeds, planning, new strategy, etc.
- New ESC hardware and service contract are working well
 - Plans to deploy replicate service are delayed to early CY 2009
- Federated trust - PKI policy and Certification Authorities
 - Service continues to pick up users at a pretty steady rate
 - Maturing of service - and PKI use in the science community generally

References

[OSCARS]

For more information contact Chin Guok (chin@es.net). Also see <http://www.es.net/oscars>

[Workshops]

see <http://www.es.net/hypertext/requirements.html>

[LHC/CMS]

<http://cmsdoc.cern.ch/cms/aprom/phedex/prod/Activity::RatePlots?view=global>

[ICFA SCIC] “Networking for High Energy Physics.” International Committee for Future Accelerators (ICFA), Standing Committee on Inter-Regional Connectivity (SCIC), Professor Harvey Newman, Caltech, Chairperson.

<http://monalisa.caltech.edu:8080/Slides/ICFASCIC2007/>

[E2EMON] Geant2 E2E Monitoring System –developed and operated by JRA4/WI3, with implementation done at DFN

http://cnmdev.lrz-muenchen.de/e2e/html/G2_E2E_index.html

http://cnmdev.lrz-muenchen.de/e2e/lhc/G2_E2E_index.html

[TrViz] ESnet PerfSONAR Traceroute Visualizer

<https://performance.es.net/cgi-bin/level0/perfsonar-trace.cgi>

[ASCAC] “Data Communications Needs: Advancing the Frontiers of Science Through Advanced Networks and Networking Research” An ASCAC Report: Ellen Stechel, Chair, Bill Wing, Co-Chair, February 2008