



Summary Results From NASA High End Computing (HEC) WAN File Accessing Experiments/Demonstrations At SC10

Pat Gary/GSFC
Pat.Gary@nasa.gov

Bill Fink/GSFC
William.E.Fink@nasa.gov

Paul Lang/ADNET
Paul.A.Lang@nasa.gov

Computational and Information Sciences and Technology Office (CISTO), Code 606
NASA Goddard Space Flight Center
January 10, 2011

Presentation for 13Jan11 Meeting of GSFC's
Network Evolution and Architecture Transformation Working Group (NEATWG)



01/10/11
GODDARD SPACE FLIGHT CENTER

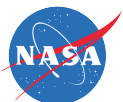
J. P. Gary



Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Topics

- Background Context
 - HECN Team's fundamental objectives
 - SC10's SCinet Research Sandbox opportunity
- Accomplishments
 - Joint plans and implementation
 - Sample results
- Significance
 - Continuing and new partnerships
 - Usefulness of HECN net-test workstations
 - Preparing future plans

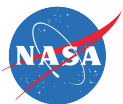




Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Problem Statement

- GSFC's NASA Center for Computational Sciences (NCCS) is increasing its data production/analysis/storage capabilities and its accessing of other large remote data
- Higher bandwidth networks can be deployed
- But bottlenecks in the combination of our file copying applications, disk I/O subsystems, server/workstation configurations, protocol stack tuning and/or NICs is preventing full use of our higher bandwidth networks
- Need to determine server/workstation configurations, data transfer utilities and protocols that enable higher throughput, especially given the emergence of 40- to 100-Gbps networks





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Notional Milestone Schedule

	CY09	CY10	CY11	CY12
	JFMAMJJASON	DJFMAMJJASON	DJFMAMJJASON	DJFMAMJJASON
<u>Phase 0 10-Gbps Testbeds</u>				
O LAN & Region/MAN	-----**			
O WAN	-----*			
<u>Phase 1 20-Gbps Testbeds</u>				
O LAN & Region/MAN		-----***		
O WAN		-----***		
<u>Phase 2 40-Gbps Testbeds</u>				
O LAN & Region/MAN			-----***	
O WAN			-----***	
<u>Phase 3 100-Gbps Testbeds</u>				
O LAN & Region/MAN				-----***
O WAN				-----***

Legend for Milestone Schedule

----- Planning and acquisition subphase

***** I&T plus demo subphase





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Key Considerations (1 of 2)

- Use of Layer 1+2 DCN-enabled VLANs versus Layer 1+2+3 full IP routed networks in both the Regional/MAN and WAN testbeds is critical
 - Sufficient to enable more effort to be focused on the primary subjects of this effort which are the processor interfaces and LAN infrastructure needed on the ends of the intervening links
 - Core IP routing issues (while otherwise interesting with many R&D challenges remaining) are not the primary subject of this effort
 - Costs of 40 and 100 Gbps Layer 3 router interfaces are likely to be two or more orders of magnitude greater than 40 and 100 Gbps Layer 2 Ethernet switch interfaces which are sufficient to enable the needed VLANs

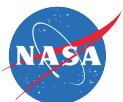




Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Key Considerations (2 of 2)

- “Suitable” near-typical server/workstations used by the science community are essential to the acceptance of the testbed findings
- HECN Team has chosen to iteratively specify and then assemble net-test server/workstations to achieve the disk-to-disk throughput performance goals of the respective Phases
 - Lowest cost approach
 - Update specs and assemble new server/workstations to overcome newly discovered bottlenecks
 - Memory-to-memory throughput tests significantly help to calibrate the infrastructure’s wire-speed





10-Gbps Disk-to-Disk File Copies Achieved Via Workstations Costing Less Than \$7,000

- As part of plans to assess the throughput performance of wide-area file transfer applications, GSFC's High End Computer Network Team specified and assembled workstations that individually costs less than \$7,000 and are capable of over 10 gigabits per second (Gbps) disk-to-disk file copying.
- Each workstation consists of a 3.2-GHz single-processor (quad core) Intel Core i7 (Nehalem) with two HighPoint RocketRaid 4320 RAID disk controllers and a Myricom 10 Gigabit Ethernet network interface card in the PCIe Gen2 slots of a Asus P6T6 WS Revolution motherboard. Each RAID controller hosts eight Western Digital WD5001AALS 500-gigabyte disks.
- Over 10-Gbps disk-to-disk file-copying throughput between two of the workstations was measured using the nuttscp (www.nuttcp.net) file copying tool.
- Demonstrations of these workstations supporting network-performance testing, wide-area file systems, and file transfer applications ranging from traditional to experimental are planned in the NASA research exhibit at the SC09 conference, Portland, OR, November 16–19 .



Figure: Two Core i7 workstations interconnected via 10 Gigabit Ethernet in test configuration prior to shipping to SC09.

POC: Bill Fink, William.E.Fink@nasa.gov, (301) 286-7924, GSFC Computational and Information Sciences and Technology Office





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Phase 1 Network-Test Workstations: Nominal “B” System

- Chassis: Supermicro 836TQ-R800B (3u 16bay 7slot 800W RPS)
- Motherboard: Asus P6T6 WS Revolution (5 PCIe V2 x8)
- Processors: one Intel i7 965 (3.2GHz quad-core Nehalem)
- Memory: Kingston KHX16000D3ULT1K3 (6GB 2000MHz DDR3 CL8)
- System disks: one Western Digital WD2500BEKT (2.5” 250GB)
- NICs: two Myricom 10G-PCIE2-8B2-2S+E (Dual 10GE SFP+)
- Raid controllers: two HighPoint RocketRaid 4320 (internal, 8 disks each)
- User disks: 16 Western Digital WD5001AALS (500GB)
- IB HCA: one Qlogic QLE7280 (DDR, 8x)
- For more detail, contact Paul.Lang@nasa.gov





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Some More Phase 1 Network-Test Workstations

- “C” system: Like B, but with twice the number of data disks and disk controllers
 - E.g., i7test3 and i7test4
- “A” system: Like B, but without any data disks
 - Useful for memory-to-memory wire-speed assessments
 - Also useful as WAN emulators and firewalls
- “A+” system: Like A, but with extra 10-GE NICs
- “X” system: Like A, but Xeon-based
- “X++” system: Like X, but with extra 10-GE NICs
- “XSSD++” system: Like X++, but with SSD data disks
- For more detail, contact Paul.Lang@nasa.gov





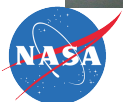
17.8-Gbps Disk-to-Disk File Copies Achieved Via Workstations Costing Less Than \$9,000

- As part of plans to assess the throughput performance of wide-area file transfer applications, GSFC's High End Computer Network Team specified and assembled workstations that individually costs less than \$9,000 and are capable of over 17.8 gigabits per second (Gbps) disk-to-disk file copying.
- Each workstation consists of a 3.2-GHz single-processor (quad core) Intel Core i7 (Nehalem) with four HighPoint RocketRaid 4320 RAID disk controllers and a Myricom 2-port 10 Gigabit Ethernet network interface card in the PCIe Gen2 slots of a Asus P6T6 WS Revolution motherboard. Each RAID controller hosts eight Western Digital WD5001AALS 500-gigabyte disks.
- Over 17.8-Gbps disk-to-disk file-copying throughput between two of the workstations was measured using the nuttscp (www.nuttcp.net) file copying tool.
- While SSD technology is next to be investigated, parallelism of multiple cores and multiple streams is likely to be key to going to 40-Gbps and beyond, since individual cores are not getting significantly faster.



Figure: Right case houses Core i7 cores, DDR3 memory, NIC, two “internal” controllers each with eight disks and two “external” controllers; left case houses sixteen SAS-connected disks.

POC: Bill Fink, William.E.Fink@nasa.gov, (301) 286-7924, GSFC Computational and Information Sciences and Technology Office





Single Link Results To Date [i.e. 5Oct10*]

Source: Hoot Thompson/PTP (GSFC/NCCS)

- i7test3 to i7test4 disk-to-disk transfers [not same as Fink used]
- Single 10 Gigabit link, 0 msec RTT [not same as Fink used]
- Single file, size as noted [Fink used 64GB files]
- Source file generated to be random data
- Destination directory empty, file does not exist
- Performance in MB/sec
- Theoretical maximum of 1212 MB/sec

App	4GB	8GB	16GB	32GB	64GB	128GB
gridFTP	890	942	1057	1018	1040	1047
Aspera	██████	██████	██████	██████	██████	██████
xdd	730	852	992	1094	1156	1186
FDT	455	630	862	993	1074	1120
dsync	418	609	548	795	772	694

***These tests are being re-run with newer application versions**





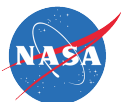
Dual Link Results To Date [i.e. 5Oct10*]

Source: Hoot Thompson/PTP (GSFC/NCCS)

- i7test3 to i7test4 disk-to-disk transfers [same as Fink used]
- Dual 10 Gigabit links, 0 msec RTT [same as Fink used]
- Single file, size as noted [Fink used 64GB files]
- Source file generated to be random data
- Destination directory empty, file does not exist
- Performance in MB/sec
- Theoretical maximum of 2424 MB/sec [Fink got 2225 using nuttscp]
- Applications a work in progress, especially Aspera and xdd

App	4GB	8GB	16GB	32GB	64GB	128GB
gridFTP	1500	1633	1546	1539	1533	1528
Aspera	██████	██████	██████	██████	██████	██████
xdd	975	1235	1361	1672	1947	2104

***These tests are being re-run with newer application versions**





100 Gigabits per Second Transmissions Achieved Via A Single Workstation

- As part of plans to assess the throughput performance of wide-area file transfer applications, GSFC's High End Computer Network (HECN) Team specified and assembled a workstation that costs less than \$11,000 and is capable of over 100 gigabits per second (Gbps) data transmission – 10 times the transmission speed of most high end computers.
- The workstation consists of a 3.2-GHz dual-processor (quad core) Intel Xeon W5580 (Nehalem) with six Myricom dual-port 10-Gigabit Ethernet network interface cards in the PCIe Gen2 slots of a Supermicro X8DAH+-F motherboard.
- Over 100-Gbps aggregate-throughput transmissions from the Xeon-workstation to two Intel Core i7 workstations (also specified and assembled by the HECN Team) were measured using the nuttcp (www.nuttcp.net) network-performance testing tool.
- Demonstrations of these workstations supporting network-performance testing, wide-area file systems, and file transfer applications ranging from traditional to experimental are planned in the NASA research exhibit at the SC09 conference, Portland, OR, Nov. 16–19 .

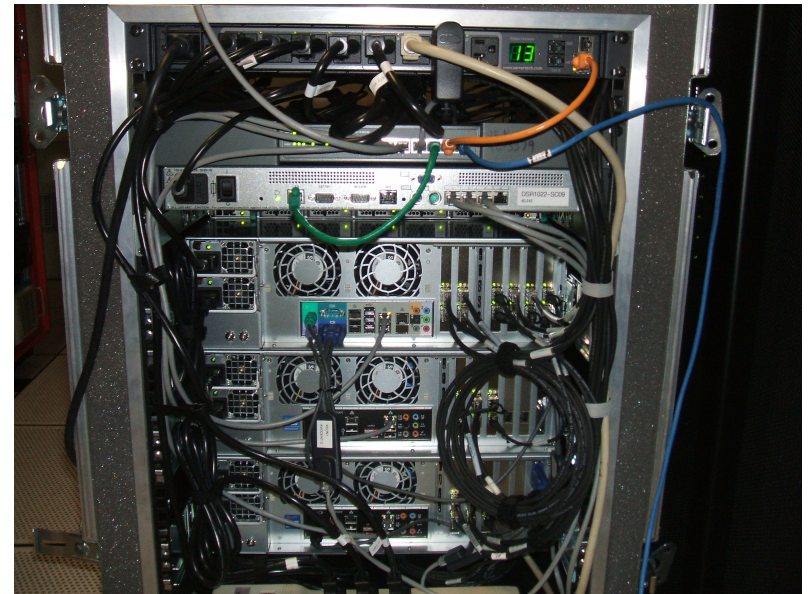
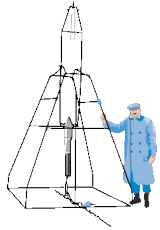


Figure: Xeon and two Core i7 workstations (bottom) interconnected with 10 Gigabit Ethernet switch and management units (top) in a rack for shipping to SC09.

POC: Bill Fink, William.E.Fink@nasa.gov,
(301) 286-7924, GSFC Computational and
Information Sciences and Technology Office

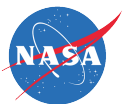




Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

“X++” Server Approximate Costs *(With components acquired via SEWP IV in lot-sizes of 3 - 15, and self assembly. Source: Paul Lang (ADNET))*

• Supermicro 836TQ-R800B 3u 16bay 800W RPS Chassis	\$850
• Supermicro X8DAH+-F motherborad	\$508
• Intel W5580 XEON 3.2GHz processor \$1669 x 2	\$3338
• Kingston KHX2000C8D3T1K3 6GB DDR3 2000 CL8 memory x 2	\$500
• CBL-0084 front pannel cable	\$3
• 12" 3pin fan extension cable	\$1
• ArkTech slim IDE DVD to SATA adapter	\$10
• Myri 10G-PCIE2-8B2-2S+E Dual SFP+ NIC \$950 x 6	\$5700
• Dynatron G666 CPU cooler	\$35
• Western Digital WD2500BEKT 250GB 2.5" system disk	\$73
• Red Greatland 18" Slimline SATA adapter	\$6
• Supermicro MCP-220-83601-0B FDD tray for 2.5" disk	\$8
• eVGA GeForce 8400GS video card	\$40
• 8" 8pin power extension cable	\$8
	<hr/>
	\$11080

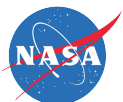




Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

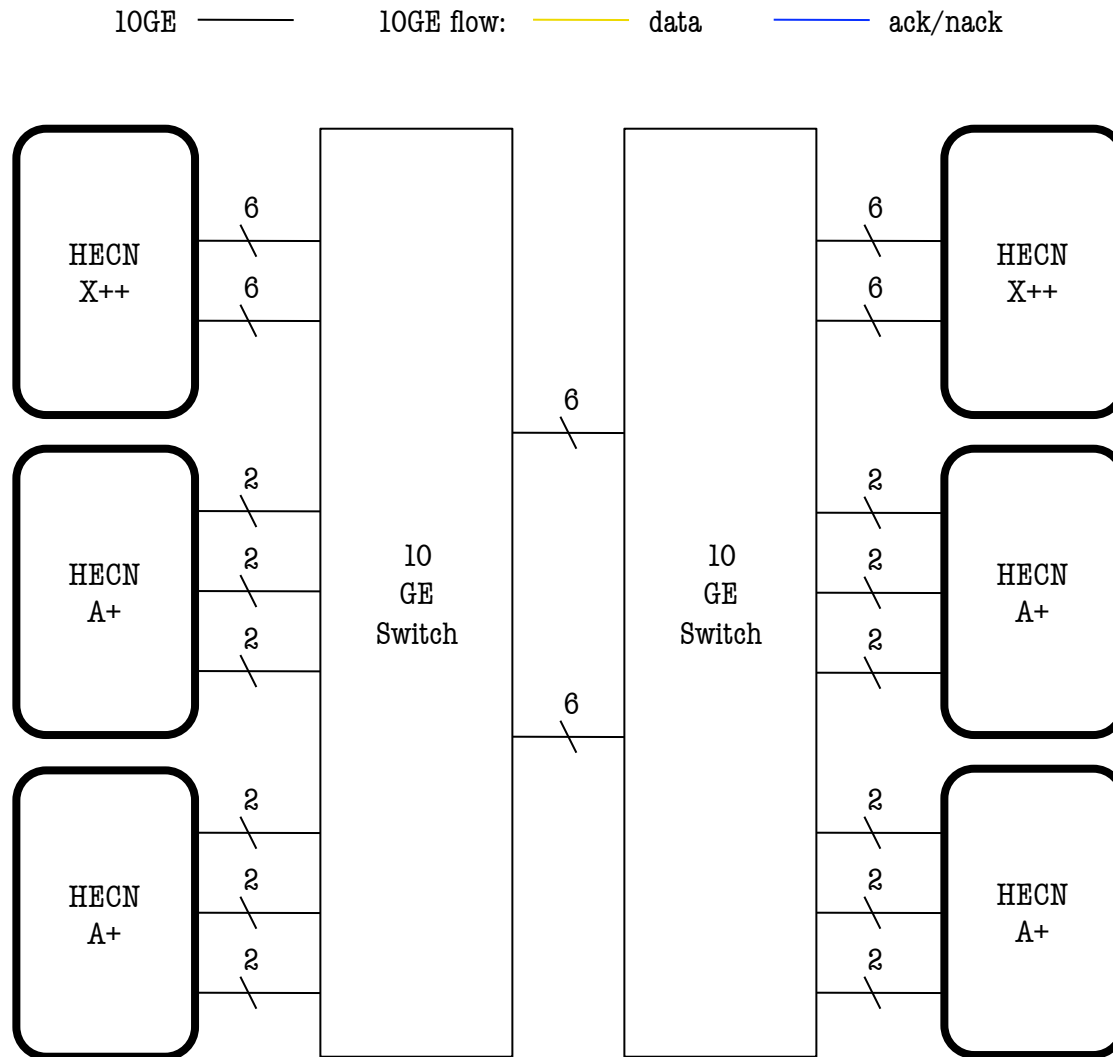
Nuttcp >100 Gbps Uni-Directional Memory-to-Memory Flow Paths

- Receive limitations in HECN's Xeon-based workstations cause HECN to use two Core i7-based workstations to receive >100-Gbps uni-directional memory-to-memory data transmissions from one HECN Xeon-based transmitter
- Bi-directionally filling an intervening network with 100-Gbps data flows requires (only) two sets, each consisting of one HECN Xeon-based transmitter and two HECN Core i7-based receivers





Nuttcp >100 Gbps Uni-Directional Memory-to-Memory Flow Paths (1 of 3)

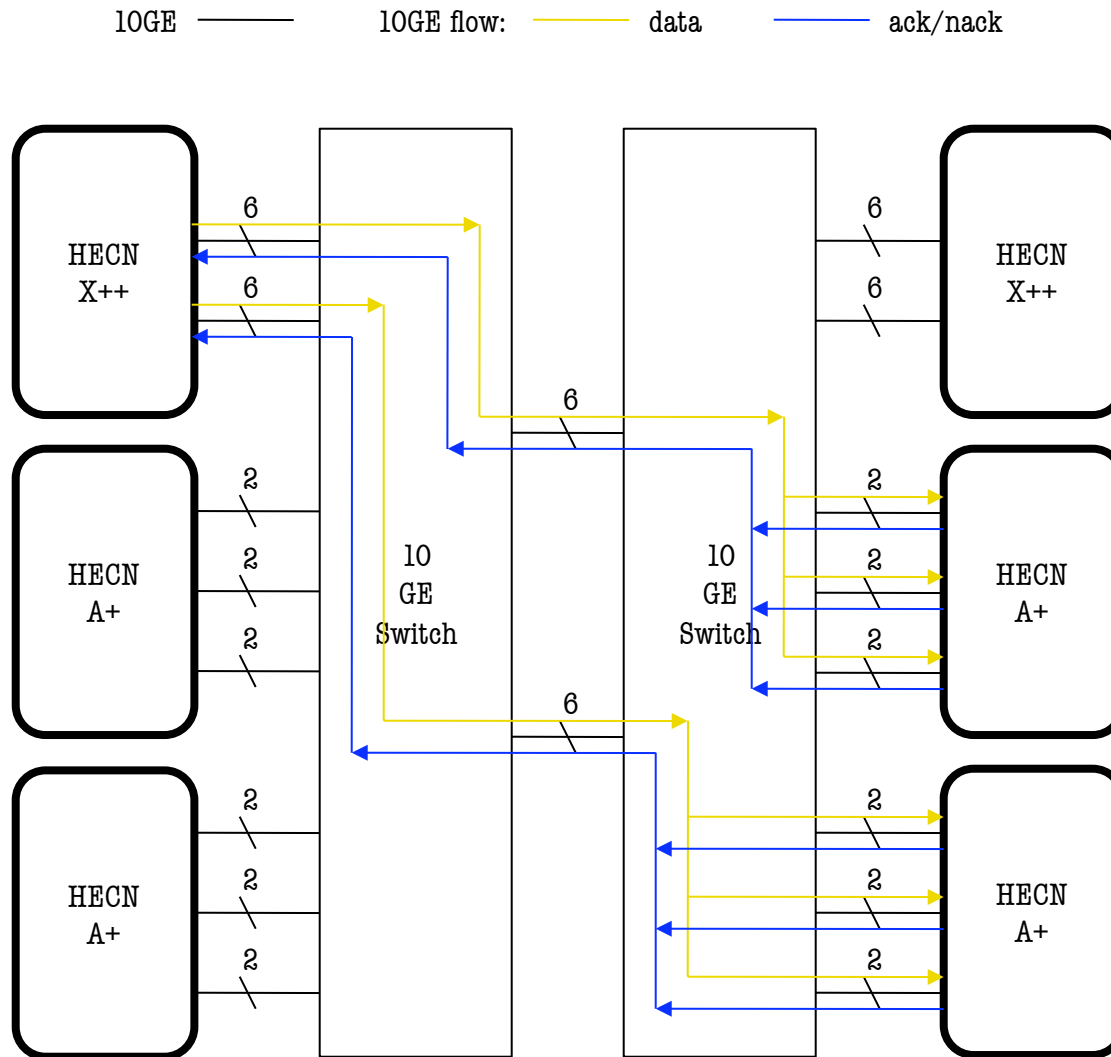


J. P. Gary 01/09/11

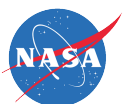




Nuttcp >100 Gbps Uni-Directional Memory-to-Memory Flow Paths (2 of 3)

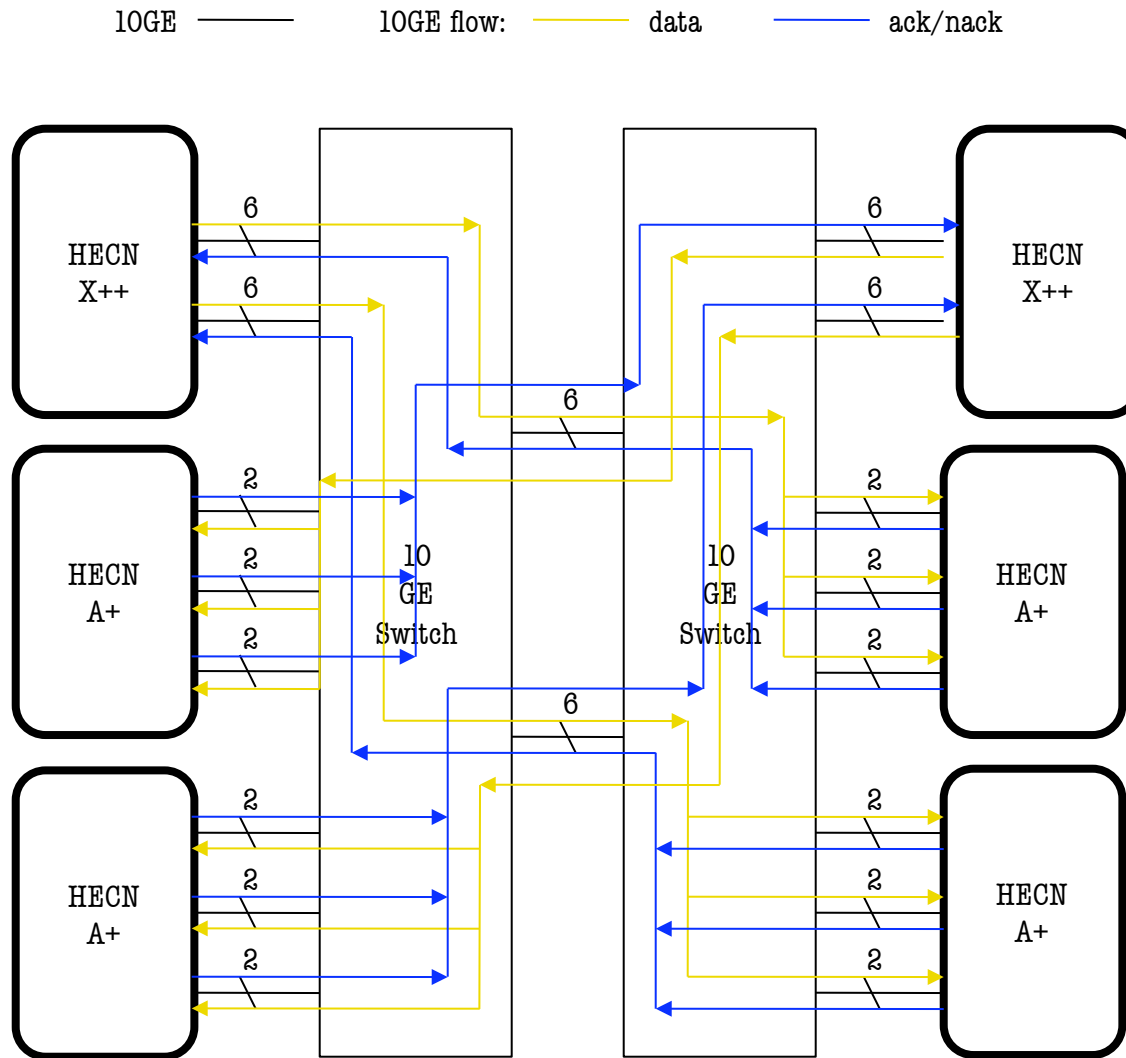


J. P. Gary 01/09/11





Nuttcp >100 Gbps Uni-Directional Memory-to-Memory Flow Paths (3 of 3)



J. P. Gary 01/09/11





Aggregate 55+ Gigabits per Second (Gbps) Transmits, 52+ Gbps Receives and 75+ Gbps Bi-Directional Transmissions Achieved Via A Single Workstation With a Single 6x10-Gigabit Ethernet Network Interface Card

- As part of plans to assess the throughput performance of wide-area file transfer applications, GSFC's High End Computer Network (HECN) Team tested a HotLava six-port 10-Gigabit Ethernet network interface card (NIC) in a HECN Team-assembled workstation that costs less than \$ 6,800 with the NIC and achieved aggregate 55+ Gbps transmits, 52+ Gbps receives and 75+ Gbps bi-directional memory-to-memory data transmissions.
- The workstation consists of a 3.2-GHz single-processor (quad core) Intel Core i7 (Nehalem) with one HotLava NIC in one PCIe Gen2 x16 slot of an Asus P6T6 WS Revolution motherboard.
- Transmissions between the above workstation and two other HECN Team-assembled Intel Core i7 workstations with other NICs were measured using the nuttcp (www.nuttcp.net) network-performance testing tool.
- Demonstrations of similar workstations supporting 100 Gbps network testing and near-40 Gbps file transfer applications are planned in the NASA research exhibit at the SC10 conference, New Orleans, LA, Nov. 15–18.

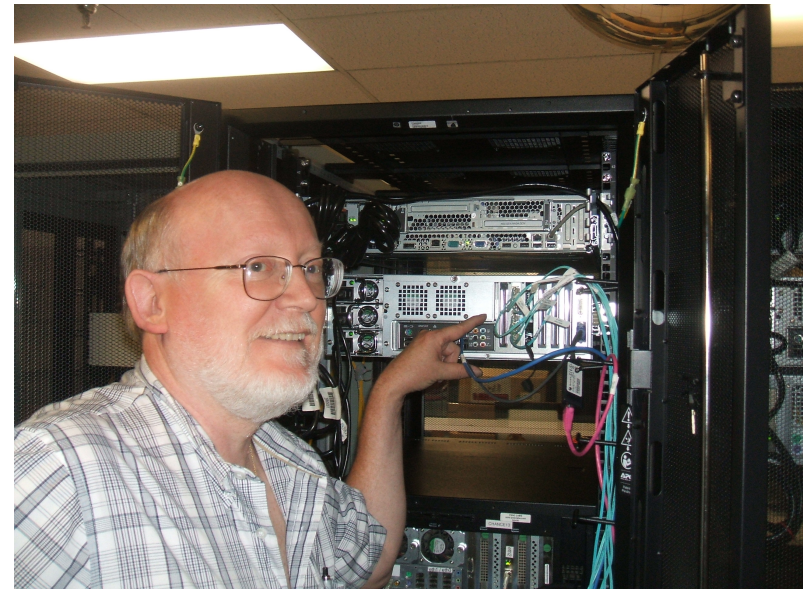
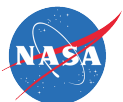


Figure: Bill Fink, author of nuttcp and the throughput performance tests, pointing to the 6x10GE HotLava NIC in the HECN Team's Intel Core i7 based workstation.

POC: Bill Fink, Bill.Fink@nasa.gov,
GSFC Computational and Information
Sciences and Technology Office





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

“XSSD1++” Server Approximate Costs (With components acquired via SEWP
IV in lot-sizes of 3 - 15, and self assembly. Source: Paul Lang (ADNET)) (1 of 2)

- Supermicro 836TQ-R800B 3u 16bay 800W RPS chassis \$828
- Supermicro X8DAH+-F motherboard \$506
- Intel X5590 4-core 3.3GHz Xeon proc. \$1,612 X 2 = \$3,224*
- 3X2GB 1600 MHz DDR3 memory \$227 X 2 = \$454
- CBL-0084 front pannel cable \$3
- 12" 3pin fan extension cable \$1
- ArkTech slim IDE DVD to SATA adapter \$10
- HotLava Tanbora 6xSFP+ NIC 6ST2A30A1F1 \$1,401 X 2 = \$2,802*
- Dynatron G666 CPU cooler \$35 X 2 = \$70
- 2.5" SATA system disk (WD2500BEKT 250GB) \$60
- Red Greatland 18" Slimline SATA adapter \$6
- Supermicro MCP-220-83601-0B FDD tray for 2.5" disk \$8
- PCIe Video card eVGA GeForce 8400GS \$40
- 8" 8pin power extension cable \$8



01/10/11

GODDARD SPACE FLIGHT CENTER

J. P. Gary



Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

“XSSD1++” Server Approximate Costs (With components acquired via SEWP
IV in lot-sizes of 3 - 15, and self assembly. Source: Paul Lang (ADNET)) (2 of 2)

• LSI MegaRAID 9261-8i Raid Controllers	\$460*
• LSI MegaRAID 9280-8e Raid Controllers	\$660 x 3 = \$1,920*
• Supermicro 216-R900LPB chassis 2u, 24x2.5"bay	\$891*
• OCZ Vertex 2 EX 50GB SLC SSD	\$837 x 32 = \$26,784*
• 2.5" to 3.5" adapter (IcyDock MB882SP-1S-2B)	\$12 x 8 = \$96*
• Dual SFF-8087/SFF-8088 (CoolDrives 36Hx2-26TX2)	
	\$49 x 3 = \$147*
• SAS to 4SATA cable (3ware CBL-SFF8087OCF-06M)	\$16*
• SAS-8888-05m .5m SFF-8088 SAS cable	\$42 x 3 = \$126*
	<hr/>
	\$38,456*

*Importantly different from the X++ Server

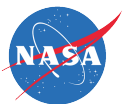




Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

SC10's SCinet Research Sandbox Opportunity

- SC10
 - 23rd annual international Supercomputing 2010 (SC10) conference on high-performance computing, networking, storage and analysis, New Orleans, Nov. 13-18, 2010
- SCinet
 - Provisioned each year for the short duration of the conference
 - One of the most powerful and advanced networks in the world: ~300-Gbps LAN capacity
- SCinet Research Sandbox
 - For network researchers: network monitoring, performance optimization, power / thermal research, network security...
 - <http://sc10.supercomputing.org/?pg=scinetsandboxprojects.html>

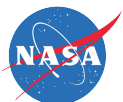




Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

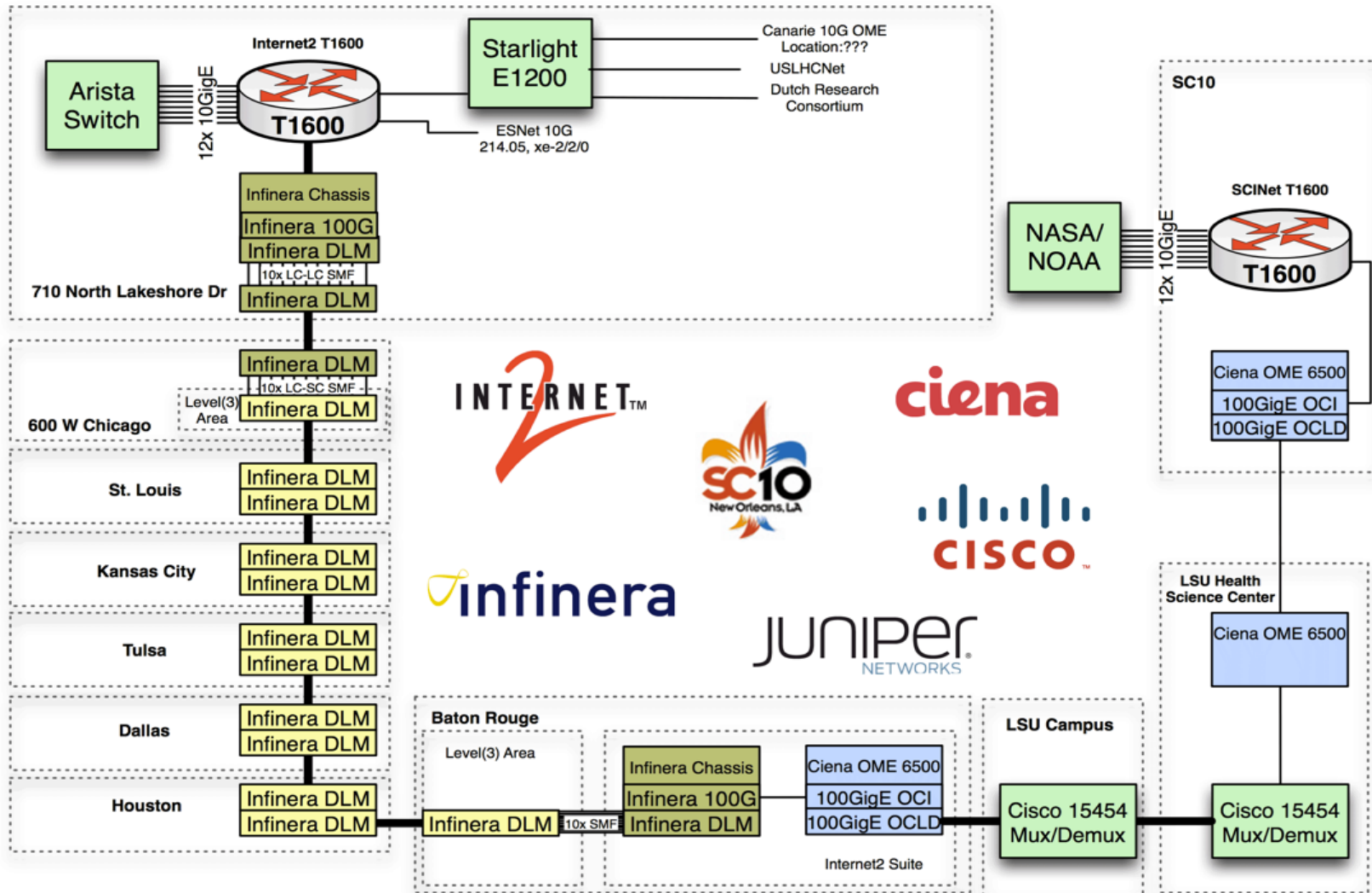
WAN Connectivity for SC10's SCinet

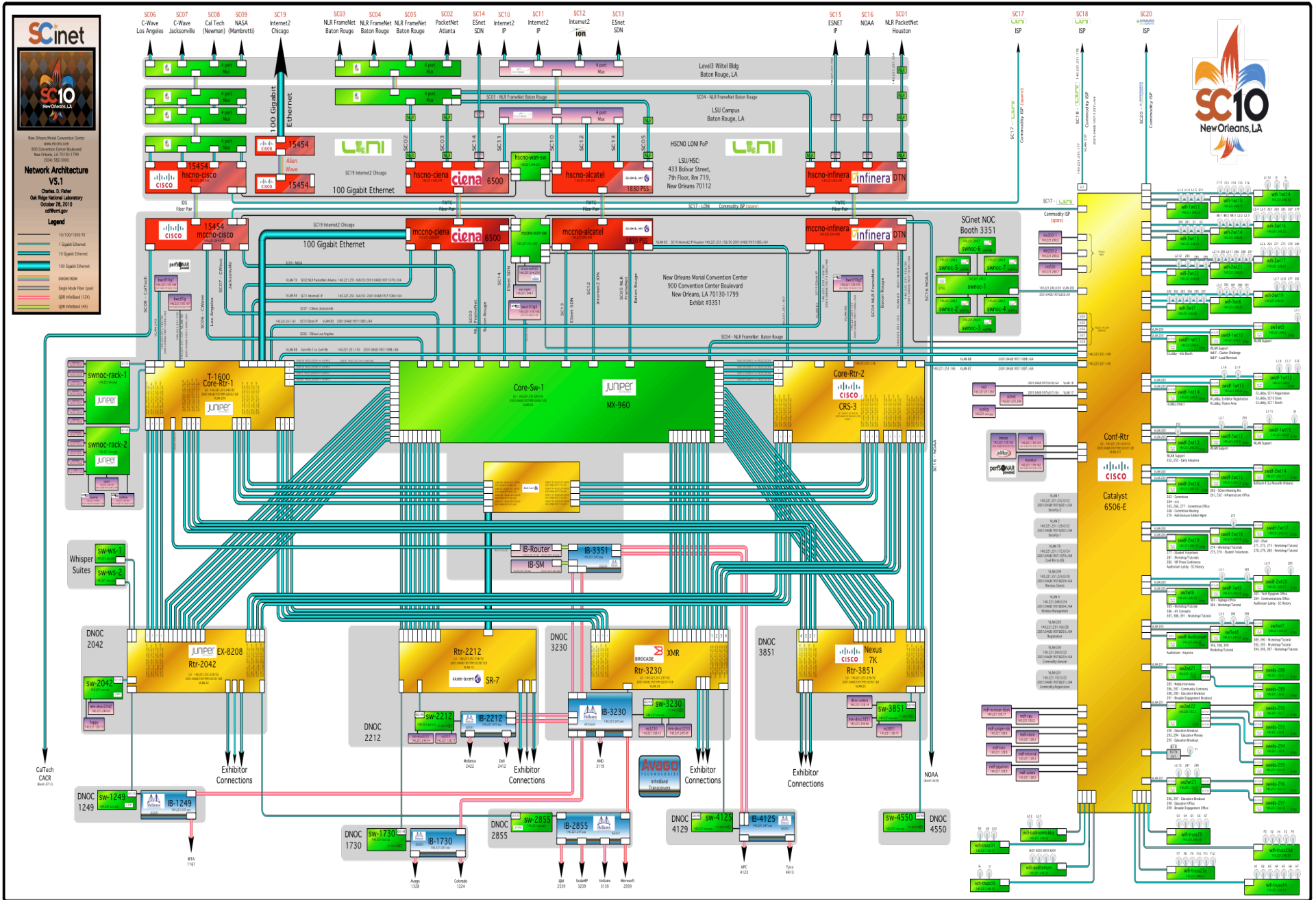
- National LambdaRail (NLR): Eight 10-Gbps links, from StarLight:
 - L2 FrameNet Wave 1_Western Route
 - L2 FrameNet Wave 2_Eastern Route
 - L2 FrameNet_STAR EXPRESS (STAR to NEW ORLEANS)
 - L3 PacketNet Wave 1_Western Route (HOUS)
 - L3 PacketNet Wave 2_Eastern Route (ATLA)
 - Cisco CWAVE West (LOSA)
 - Cisco CWAVE East (JACK)
 - L1 WaveNet_STAR EXPRESS (STAR to NEW ORLEANS) [for NASA]
- ESnet: Four 10-Gbps links, including one each supporting three projects using the SCinet Research Sandbox
 - <http://www.lbl.gov/cs/SC10/demos.html>
- Internet2: Four “normal” 10-Gbps links (two to the IP Network, one to the ION Infrastructure & one for NOAA) plus a special 1x100G pathway for their Multi-Vendor 100GigE Demo Between Chicago and SC10



Internet2's SC2010 Multi-Vendor 100G Demonstration

Source: Chris Robb/Internet2





01/10/11

J. P. Gary

25



01/10/11
GODDARD SPACE FLIGHT CENTER

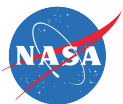
J. P. Gary



Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Objectives of NASA HEC WAN File Accessing Experiments

- Determine optimal ‘tuning parameter’ settings to obtain maximum user throughput performance with several traditional and new (or emerging) disk-to-disk file-copying utilities when operating over multi-10Gbps WANs using new state-of-the-art high performance workstations and servers
- Inter-compare throughput findings from traditional versus new file-copying utilities
- As a baseline, determine maximum memory-to-memory throughput performance among the workstations and servers using nuttcp (<http://www.nuttcp.org/>)
- Are an integral part of GSFC/HEC’s 20, 40 & 100 Gbps Network Testbed Plan (http://science.gsfc.nasa.gov/606.1/docs/HECN_10G_Testbeds_082210.pdf)





Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

NASA HEC WAN File Accessing Team

- GSFC NASA Center for Climate Simulation (NCCS)
 - Dan Duffy/GSFC
 - Hoot Thompson/PTP
 - Kirk Hunter/PTP
- GSFC/NCCS HEC Network (HECN) Team
 - Pat Gary/GSFC
 - Paul Lang/ADNET
 - Jeff Martz/ADNET
 - Aruna Muppalla/ADNET
 - Mike Stefanelli/ADNET
- ARC/CIO Network Team
 - Kevin Jones/ARC
 - Dave Hartzel/CSC
 - Hugh LaMaster/ARC
 - Mark Foster/CSC
 - Matt Mountz/CSC





Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

NASA Partners in “Using 100G Network Technology in Support of Petascale Science” Special Demos (1 of 2)

- International Center for Advanced Internet Research (iCAIR), PI: Dr. Joe Mambretti/Northwestern University
- Laboratory for Advanced Computing (LAC), PI: Dr. Bob Grossman/UIC
- Mid-Atlantic Crossroads (MAX), PM: Peter O’Neil/UMCP
- National LambdaRail (NLR), POC: Bonnie Hurst/NLR
- National Oceanic and Atmospheric Administration (NOAA), POC: Jerry Janssen
- SCinet Research Sandbox (SRS), Chair: Rodney Wilson/Ciena
- Vendors who loaned equipment – see following charts

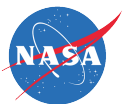




Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Acknowledgement of Vendor Equipment On Loan (1 of 2)

- Arista: Two 7148SX 48-port 10-GE switches
- Ciena: Two Optical Multiservice Edge 6500 units each with 100-Gbps transport and 10x10Gbps-to-1x100Gbps muxponder interfaces
- Cisco: Two CRS-3 routers each with one 100-GE and 14x10-GE interfaces, plus use of a third CRS-3 with two 100-GE interfaces
- ColorChip: Two DragonFly 40G-LR (up to 10km) QSFP transceivers (beta)
- cPacket: Two cVu 320G 32-port 10-GE traffic monitoring switches
- Extreme Networks: Two VIM3-40G4X 4-port 40-GE modules (beta, for Summit X650 10-GE switches)





Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Acknowledgement of Vendor Equipment On Loan (2 of 2)

- Fusion-io: Two Octal cards (SSDs on PCIe Gen2 x16)
- HP: Two ProLiant DL580 G7 servers with two 2x10-GE NICs and one QDR IB HCA
- Panduit: Two CN1 Net-Access Switch Cabinets with accessories

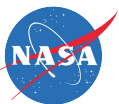




Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Acknowledgement of Partner Contributions (Partial)

- iCAIR: Many 10-GE connections in and through the StarLight@Chicago
- LAC: Exhibit booth use at SC10
- MAX: 4x10Gbps lambdas between GSFC@Greenbelt & NLR@McLean/DC
- NLR: 4x10Gbps lambdas between DC & StarLight and other 8x10Gbps pathways (two are really Cisco C-Waves; one is dedicated for NASA) between StarLight & Baton Rouge (plus coordination across the Louisiana Optical Network Initiative (LONI) regional optical network (RON) to SC10@NewOrleans)
- NOAA: Exhibit booth use at SC10
- SRS: Cost-discounted fiber-pair-bundles among the exhibit booths of NASA, NCDM-LAC/iCAIR, NOAA & SCinet NOC





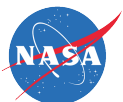
Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

NASA Partners in “Using 100G Network Technology in Support of Petascale Science” Special Demos (2 of 2)

- Internet2, POC: Chris Robb/Internet2

Acknowledgement of Partner Contributions (Partial)

- Internet2: Use of their 1x100G pathway between StarLight@Chicago & SC10@NewOrleans for their Multi-Vendor 100GigE Demo Between Chicago and SC10; on the ends:
 - Juniper: Two T1600 routers each with one 100-GE and 12x10-GE interfaces





Sample Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Special SC10 Demonstration/Evaluation Experiments

- Use a set of the NASA/HECN Team's network-testing-workstations deployed into both the NCDM-LAC/iCAIR and NASA Exhibit Booths, capable of:
 - >100-Gbps uni-directional memory-to-memory data flows
 - >80-Gbps aggregate-bidirectional memory-to-memory data flows
 - ~40-Gbps uni-directional disk-to-disk file copies (using SSDs)
- Demonstrate/evaluate different vendor-provided 40G/100G network technology solutions with full-duplex 40G and 100G LAN data flows across SCinet Research Sandbox inter-booth fiber
- Use existing 4x10G dedicated pathway across NLR and MAX/DRAGON between GSFC and StarLight, plus a mix of 8 other 10G NLR+Cisco-provisioned pathways and a 1x100G Internet2-provisioned pathway between StarLight and SC10, to conduct science-oriented WAN data flow demonstrations



Using 100G Network Technology in Support of Petascale Science

A Collaborative Initiative Among NASA, NLR, NOAA, Northwestern/iCAIR, SCinet & UIC/LAC
 Also Using Internet2's Multi-Vendor 100GigE Infrastructure Between StarLight and SC10

Demo Summary



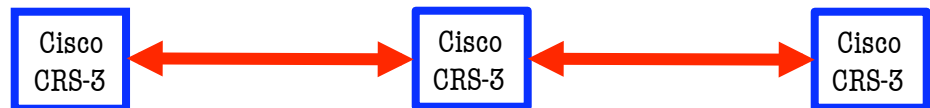
1x40Gbps full-duplex across:



1x100Gbps full-duplex across:



1x100Gbps full-duplex across:



1x100Gbps full-duplex across:



8x10Gbps full-duplex across:



40Gbps disk-to-disk between:



40Gbps disk-to-disk across:

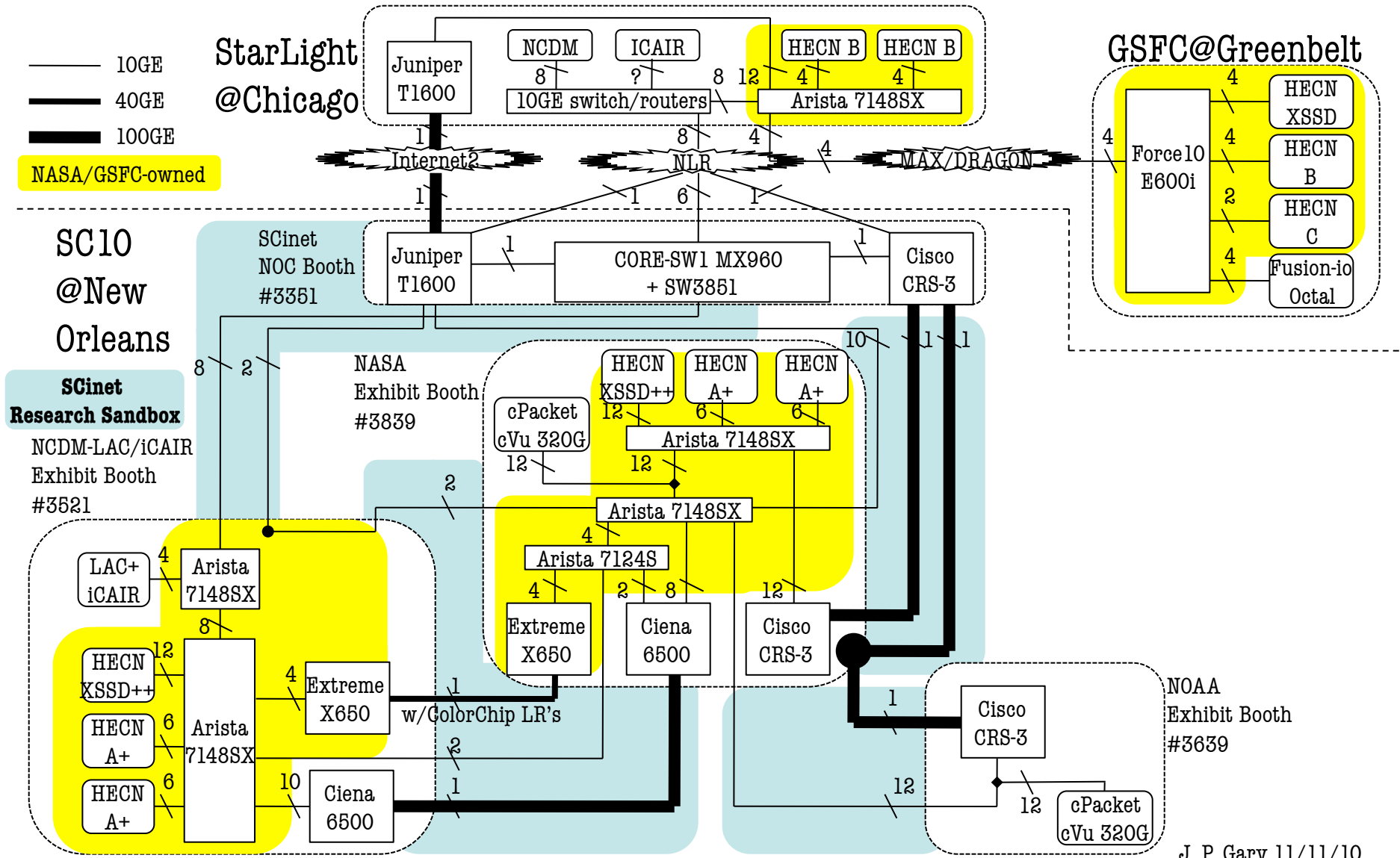


***bi-directionally**

Using 100G Network Technology in Support of Petascale Science

A Collaborative Initiative Among NASA, NLR, NOAA, Northwestern/iCAIR, SCinet & UIC/LAC

Also Using Internet2's Multi-Vendor 100GigE Infrastructure Between StarLight and SC10



01/10/11

J. P. Gary

J. P. Gary 11/11/10



Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Demo Configuration Factoids

- Connections to the respective booths from the other booths:
 - NASA: equivalent to 48x10-GE
 - NCDM: equivalent to 22x10-GE
 - NOAA: equivalent to 22x10-GE
 - SCinet: equivalent to 40x10-GE
 - **Total: equivalent to 132x10-GE**, but (after some barrel connections) only **needed 36 fiber-pairs** (not 66) as the four fiber-pairs carrying 40G and/or 100G handle the equivalent of 34x10-GE
- Intra-booth 10-GE ports used:
 - NASA: 154 (note: cPacket really uses 24)
 - NCDM: 108
 - NOAA: 36 (note: cPacket really uses 24)
 - SCinet: 20
 - **Total: 318**



Snapshots During Setup of NASA Exhibit Booth for SC10

NASA and Partners Demonstrate 40- and 100-Gigabit Network Technologies

http://science.gsfc.nasa.gov/606.1/HECN-highlights/HECN_SC10_Net-Demo_announce_110210.html



NASA and partners network-demo racks in NCDM Exhibit Booth.



NASA and partners network-demo rack in NASA Exhibit Booth.



Network-demo rack with NASA network-demo posters displayed in NOAA Exhibit Booth.



From left Paul Lang and Bill Fink.



From left Jeff Martz and Matt Mountz.



From left Pat Gary and Dave Hartzel.

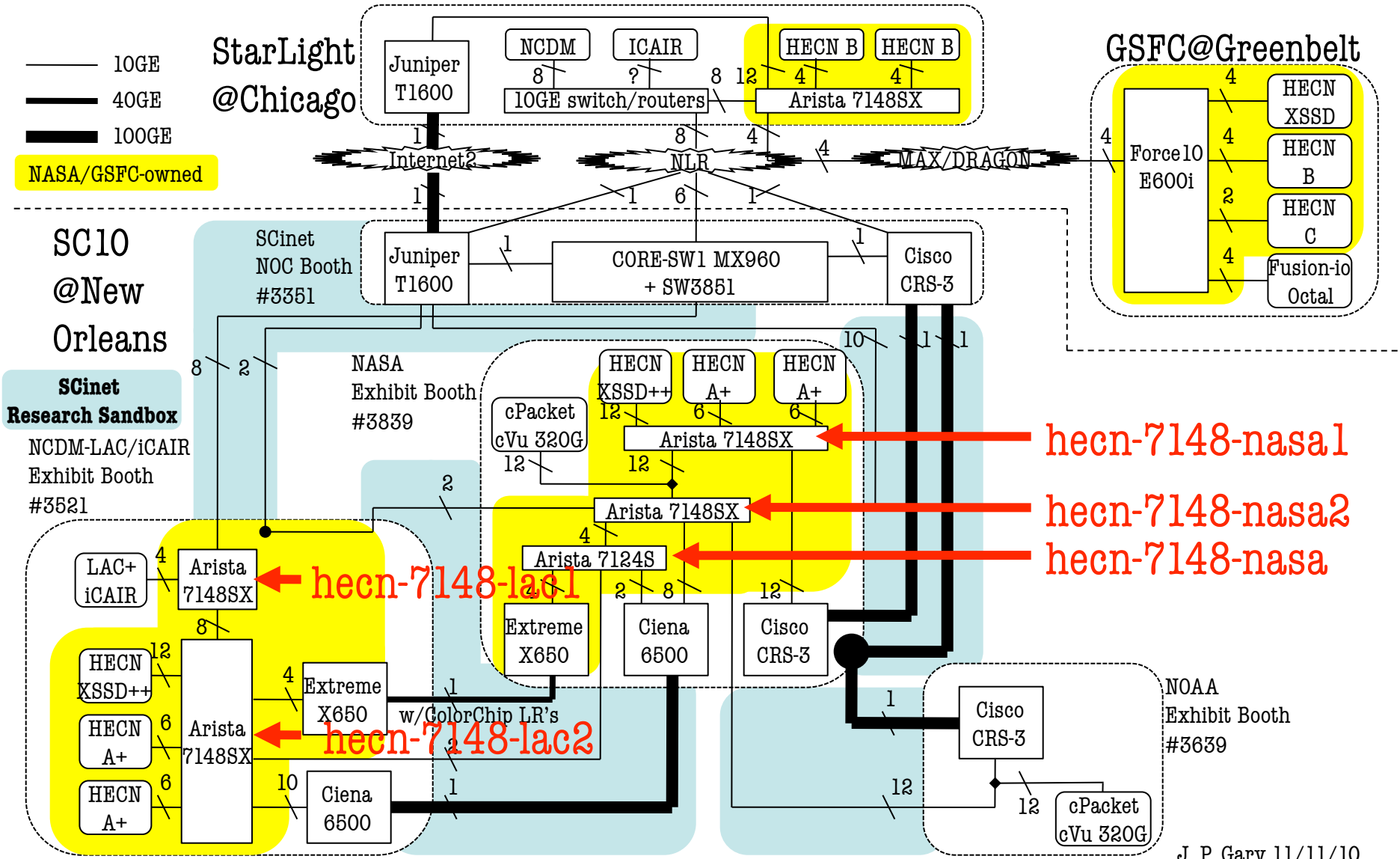
01/10/11

J. P. Gary

J. P. Gary 11/30/10
38

Using 100G Network Technology in Support of Petascale Science

A Collaborative Initiative Among NASA, NLR, NOAA, Northwestern/iCAIR, SCinet & UIC/LAC
 Also Using Internet2's Multi-Vendor 100GigE Infrastructure Between StarLight and SC10



01/10/11

J. P. Gary

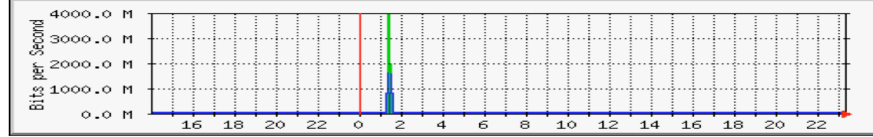
J. P. Gary 11/11/10

1 of 4

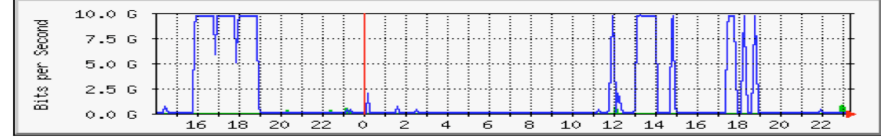
MRTG Index Page for hecn-7148-lac2

Nov 17, 2010, 10:28 PM CT

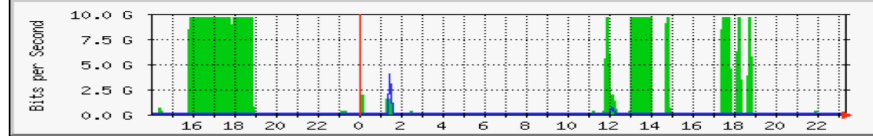
1. Traffic Analysis for Ethernet1 -- hecn-7148-lac2



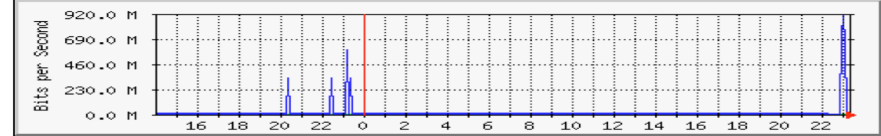
2. Traffic Analysis for xeontest1



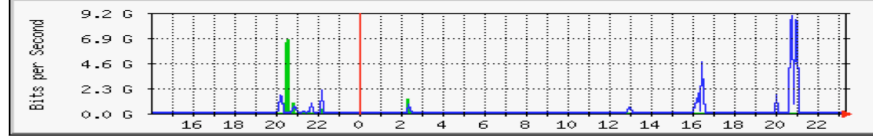
3. Traffic Analysis for xssd1



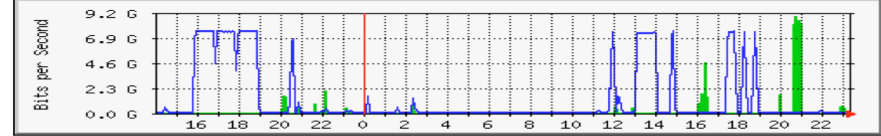
4. Traffic Analysis for i7test1



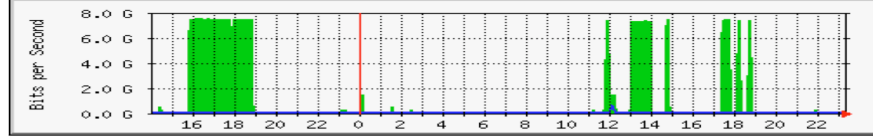
5. Traffic Analysis for Ethernet5 -- hecn-7148-lac2



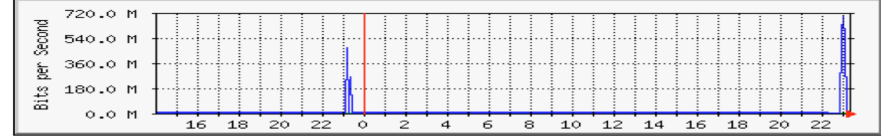
6. Traffic Analysis for xeontest1



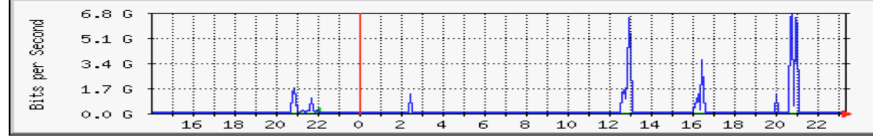
7. Traffic Analysis for xssd1



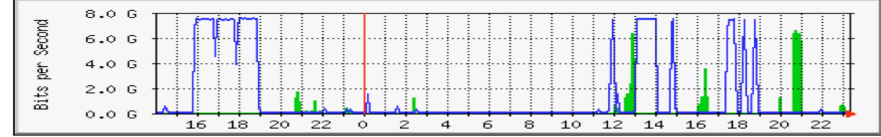
8. Traffic Analysis for i7test1



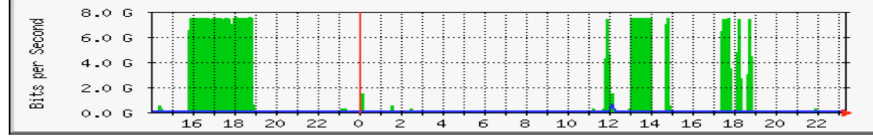
9. Traffic Analysis for xeontest1



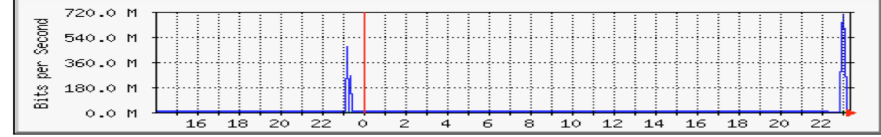
10. Traffic Analysis for xssd1



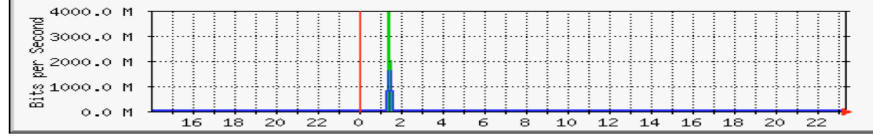
11. Traffic Analysis for xssd1



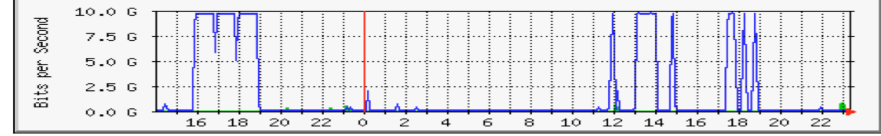
12. Traffic Analysis for i7test1



13. Traffic Analysis for Ethernet13 -- hecn-7148-lac2



14. Traffic Analysis for xeontest1

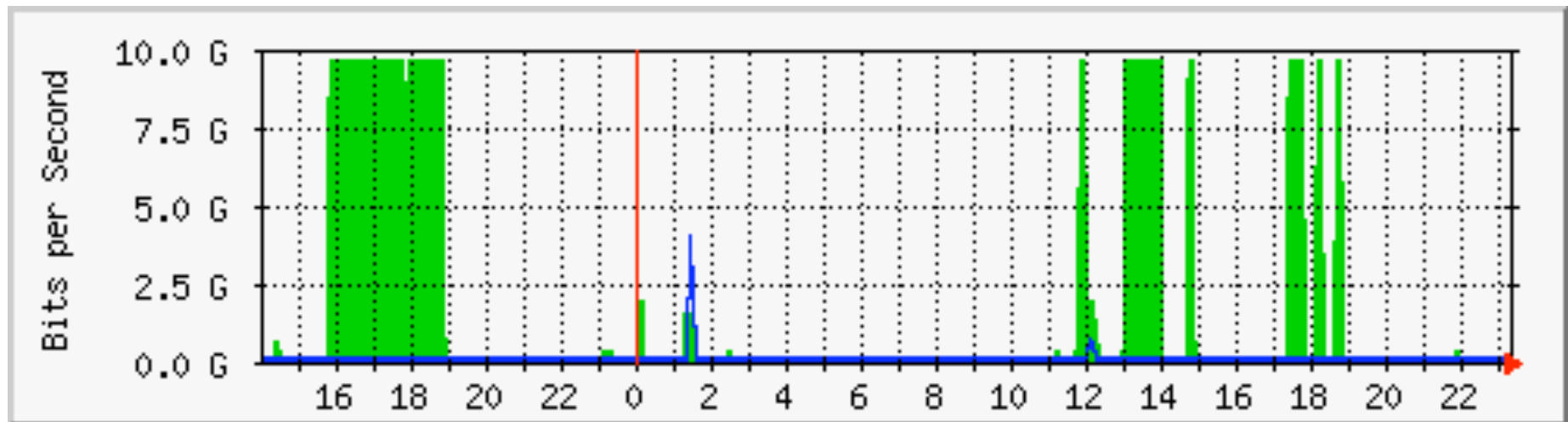




A Sample MRTG-Generated Traffic Analysis From NASA's hecn-7148-lac2 10-GE Switch In NCDM's Exhibit Booth During SC10 Of NASA Workstation XSSD1's 10-GE Network Interface #1

The statistics were last updated **Wednesday, 17 November 2010 at 22:28**

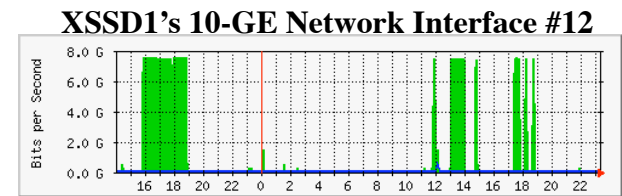
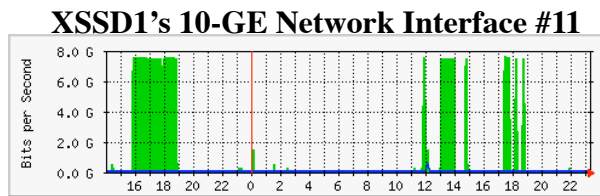
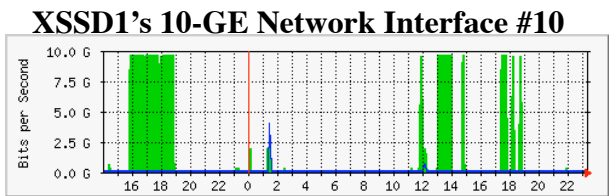
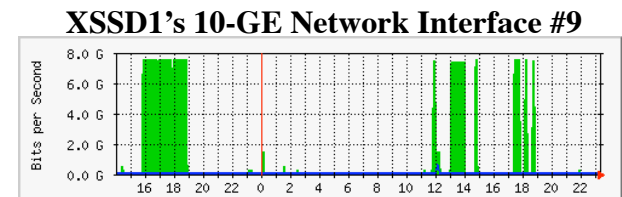
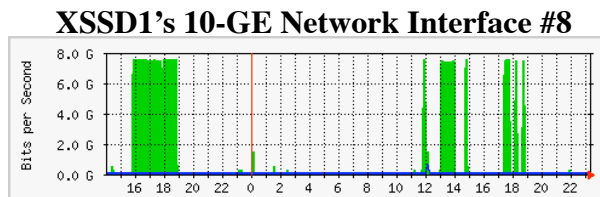
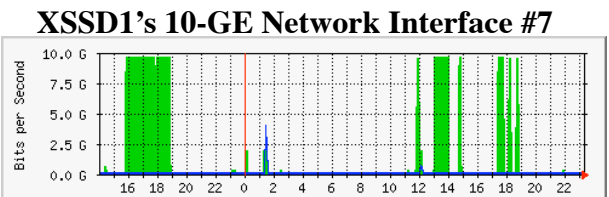
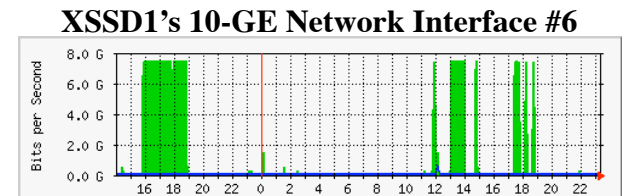
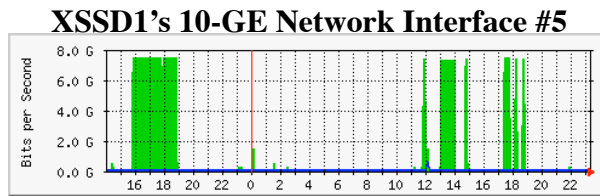
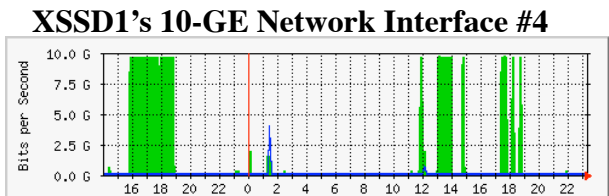
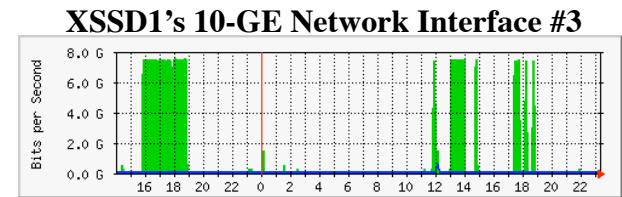
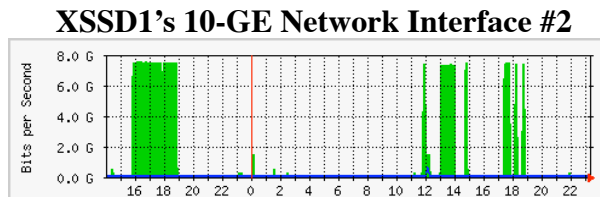
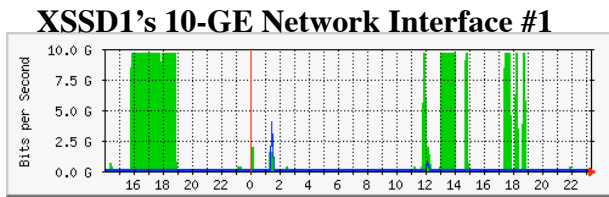
`Daily' Graph (5 Minute Average)





Sample MRTG-Generated Traffic Analyses From NASA's hecn-7148-lac2 10-GE Switch In NCDM's Exhibit Booth During SC10 Of NASA Workstation XSSD1's 10-GE Network Interfaces

The statistics were last updated Wednesday, 17 November 2010 at 22:28
'Daily' Graph (5 Minute Average)

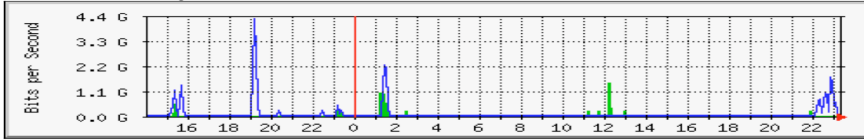


1 of 4

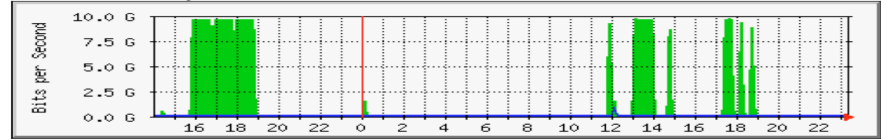
MRTG Index Page hecn-7148-nasa1

Nov 17, 2010, 10:23 PM CT

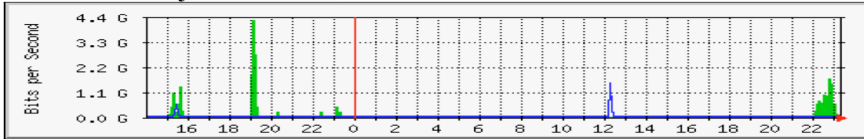
1. Traffic Analysis for hecn-7148-nasa2



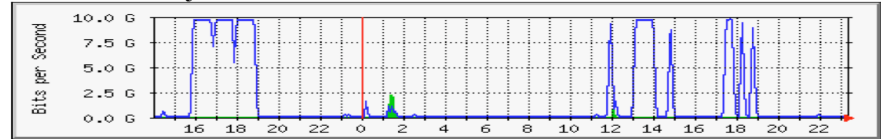
2. Traffic Analysis for cisco-crs3-nasa



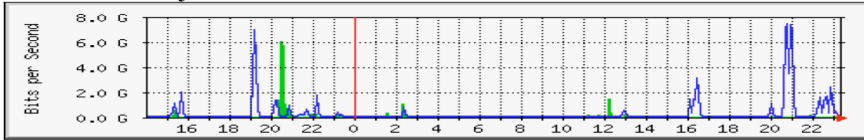
3. Traffic Analysis for xssd2



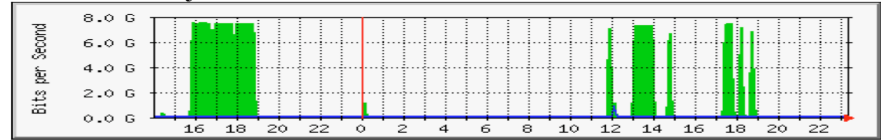
4. Traffic Analysis for i7test17



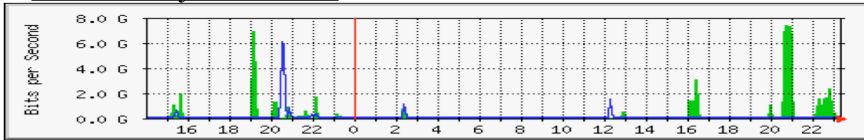
5. Traffic Analysis for hecn-7148-nasa2



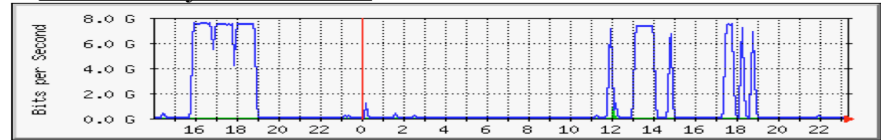
6. Traffic Analysis for cisco-crs3-nasa



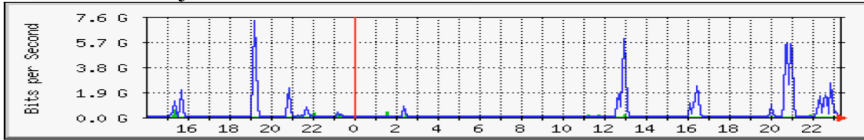
7. Traffic Analysis for xssd2



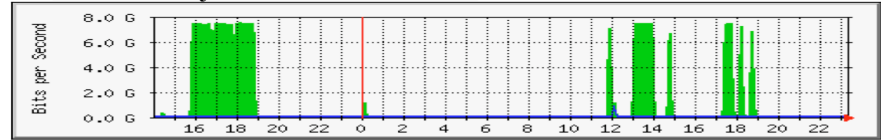
8. Traffic Analysis for i7test17



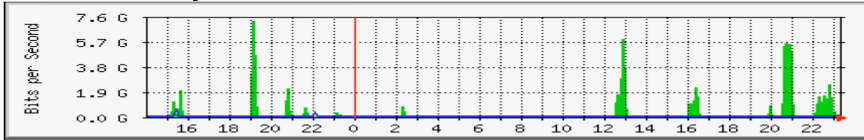
9. Traffic Analysis for hecn-7148-nasa2



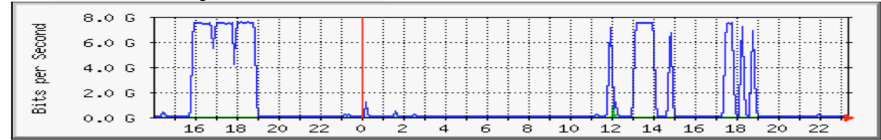
10. Traffic Analysis for cisco-crs3-nasa



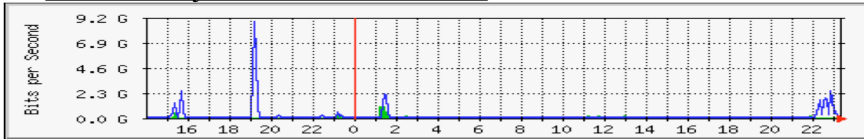
11. Traffic Analysis for xssd2



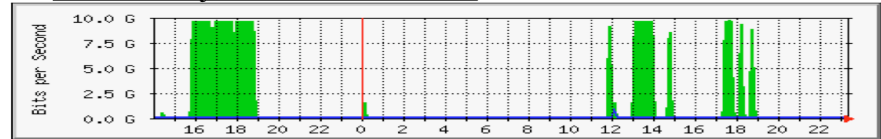
12. Traffic Analysis for i7test17



13. Traffic Analysis for hecn-7148-nasa2



14. Traffic Analysis for cisco-crs3-nasa

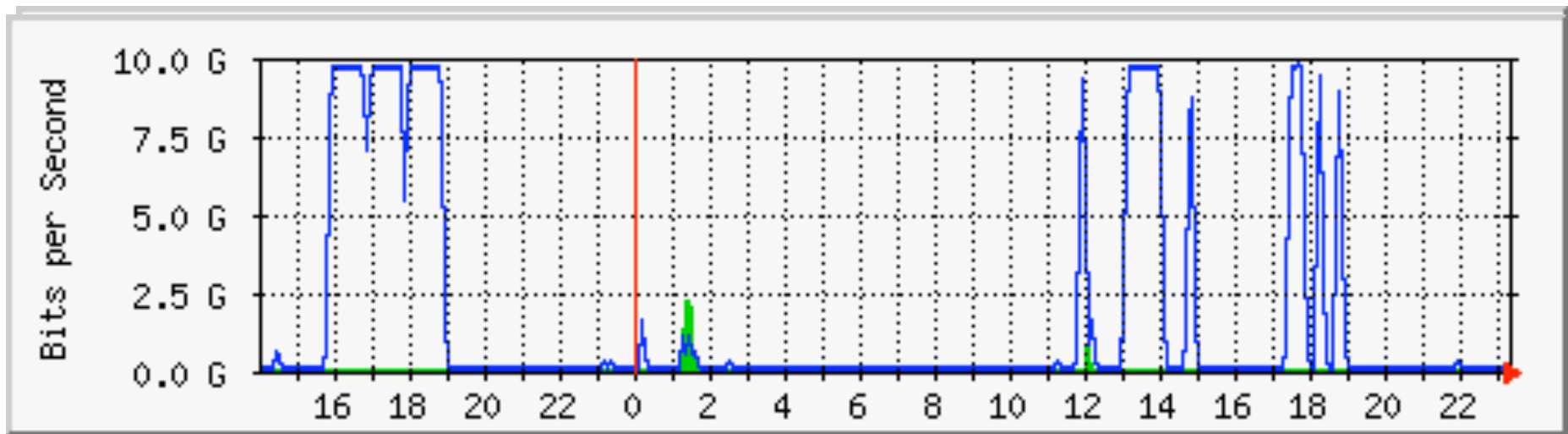




A Sample MRTG-Generated Traffic Analysis From NASA's hecn-7148-nasa1 10-GE Switch In NASA's Exhibit Booth During SC10 Of NASA Workstation i7test17's 10-GE Network Interface #1

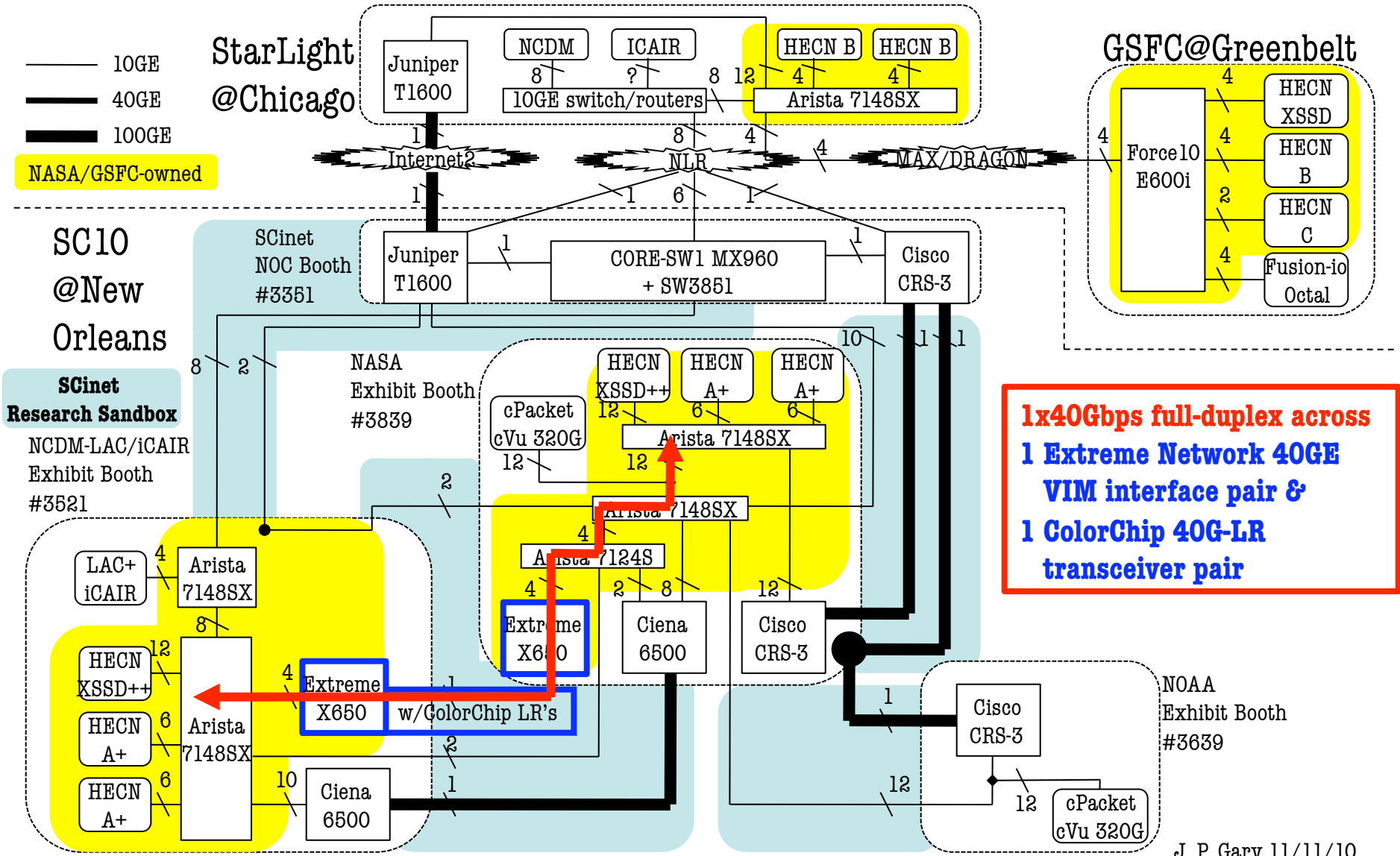
The statistics were last updated **Wednesday, 17 November 2010 at 22:23**

'Daily' Graph (5 Minute Average)



Using 100G Network Technology in Support of Petascale Science

A Collaborative Initiative Among NASA, NLR, NOAA, Northwestern/iCAIR, SCinet & UIC/LAC
 Also Using Internet2's Multi-Vendor 100GigE Infrastructure Between StarLight and SC10



01/10/11

J. P. Gary

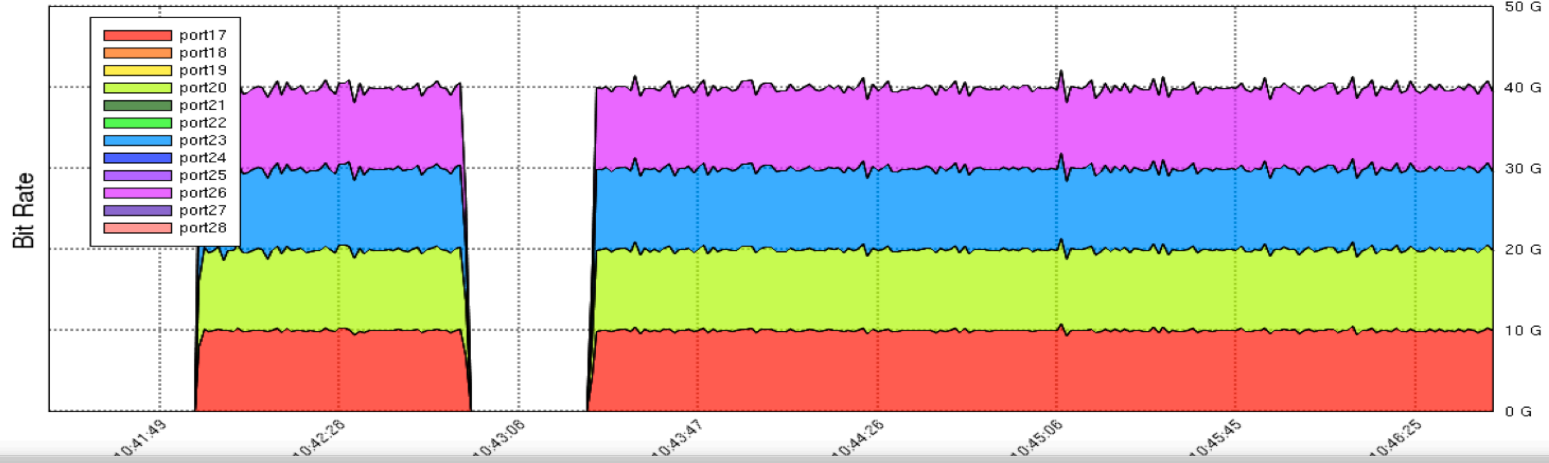
J. P. Gary 11/11/10



SC10 100G DEMO



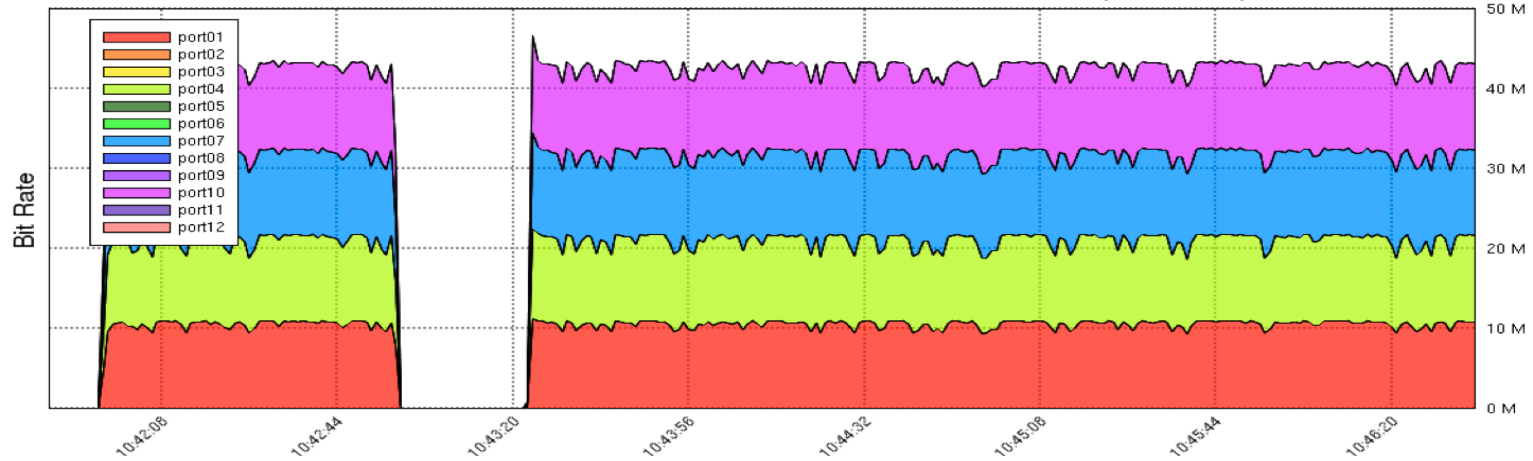
NASA-NOAA 100G Demo 18/Nov/2010 : 10:41:24 - 10:46:42 (GMT-6:00)



SC10 100G DEMO

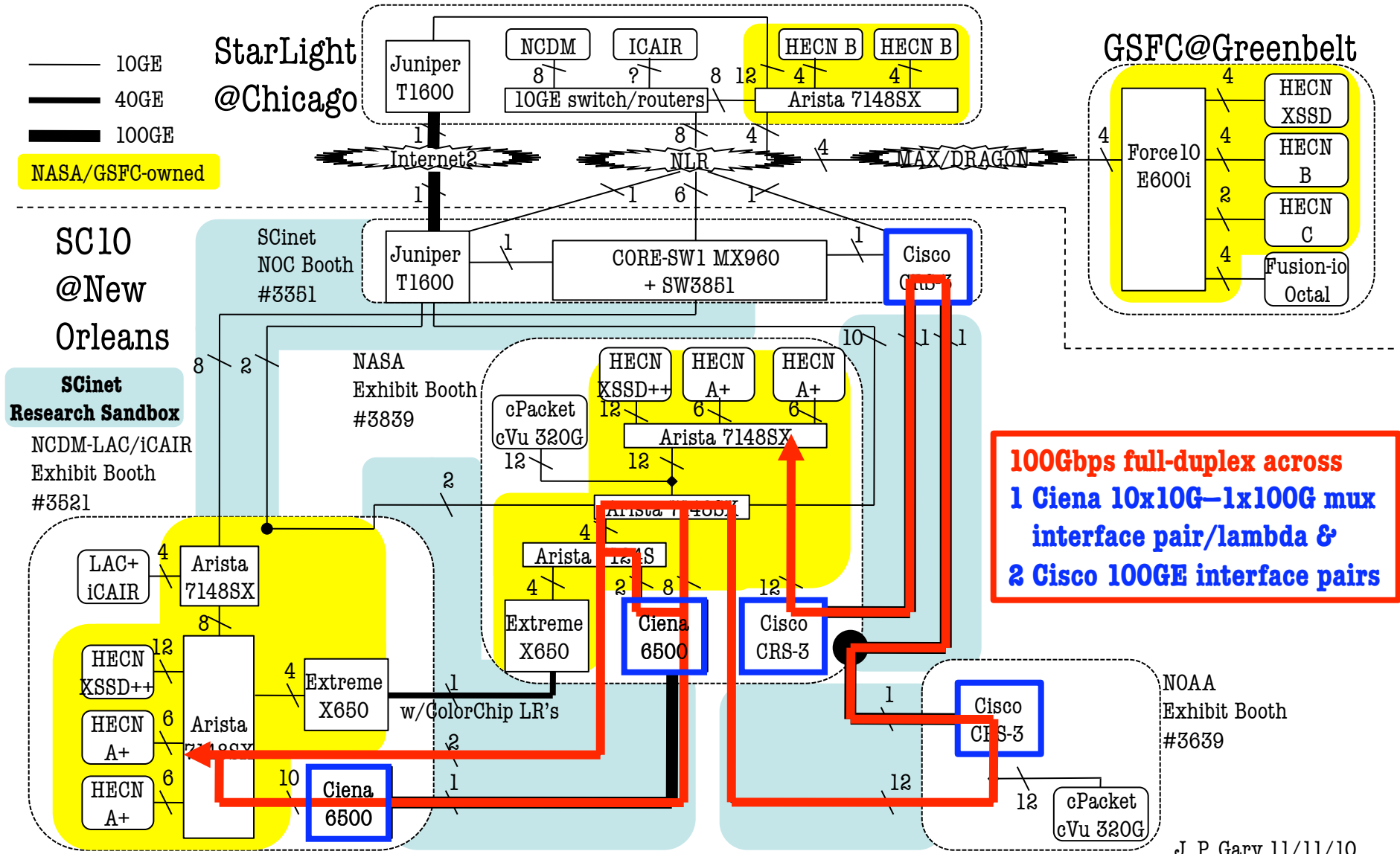


NASA-NOAA 100G Demo 18/Nov/2010 : 10:41:45 - 10:46:37 (GMT-6:00)



Using 100G Network Technology in Support of Petascale Science

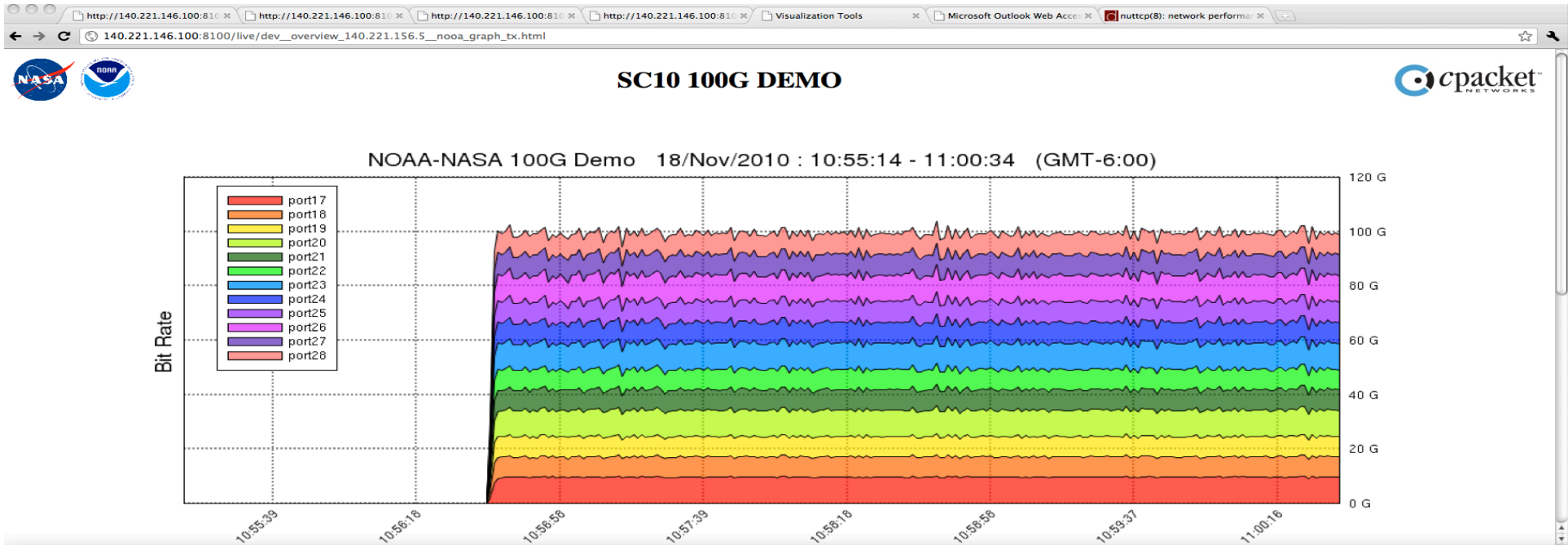
A Collaborative Initiative Among NASA, NLR, NOAA, Northwestern/iCAIR, SCinet & UIC/LAC
 Also Using Internet2's Multi-Vendor 100GigE Infrastructure Between StarLight and SC10



01/10/11

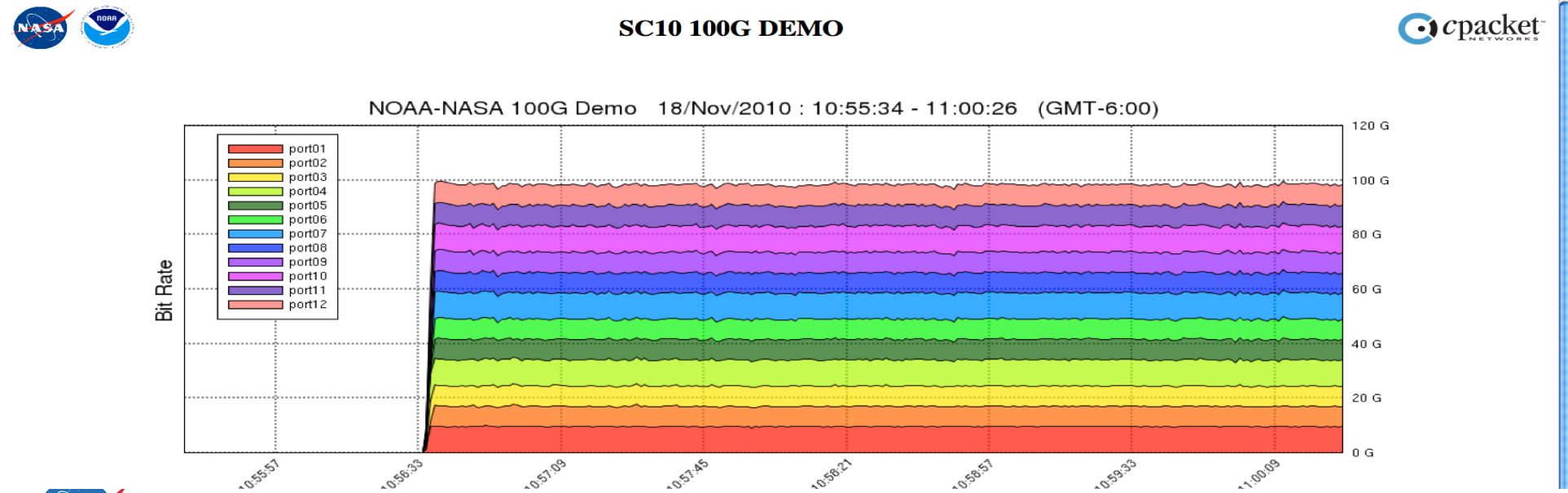
J. P. Gary

J. P. Gary 11/11/10



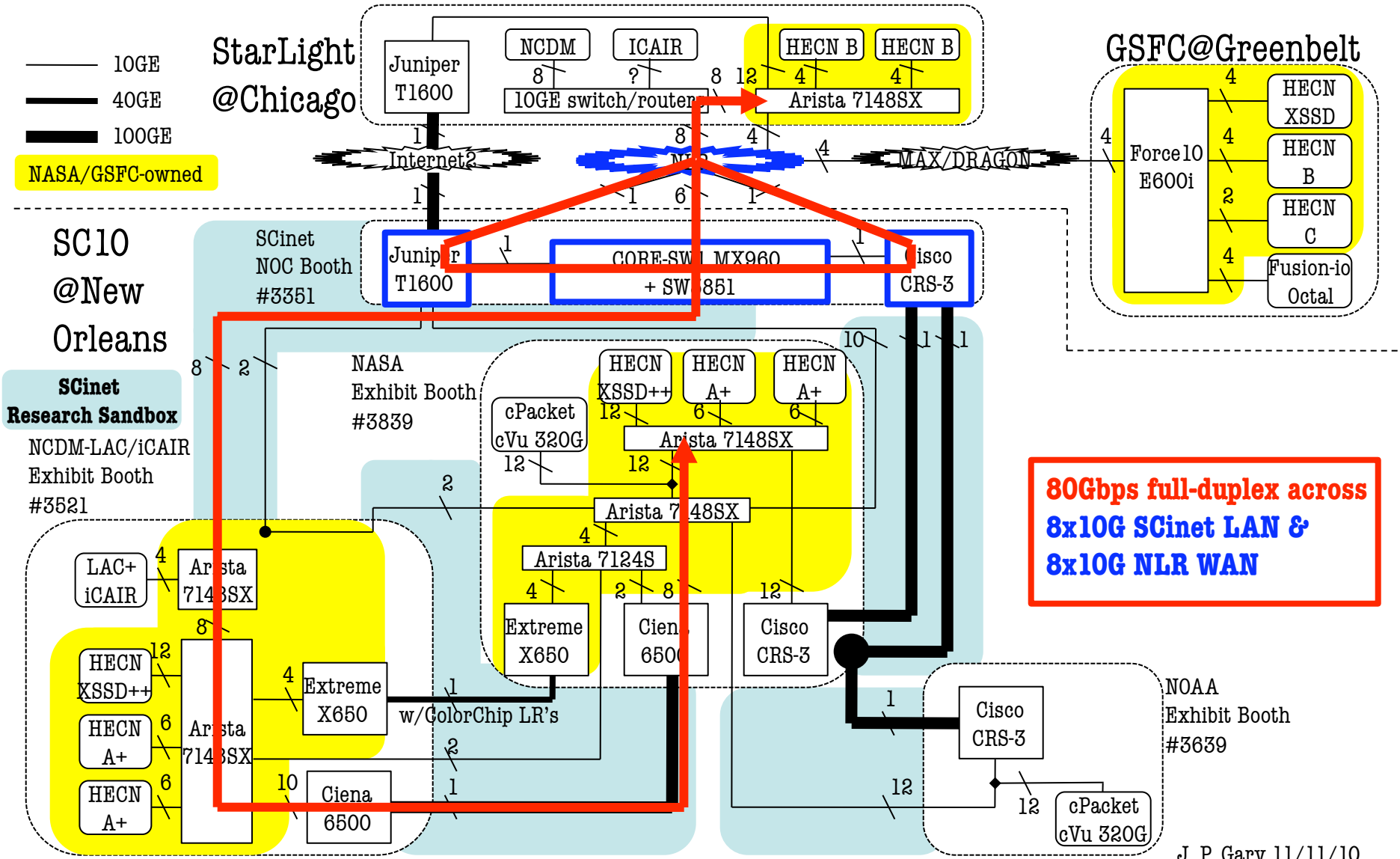
Visualization Tools

140.221.146.100:8100/live/dev_140.221.156.5.csv__noaa_graph.html



Using 100G Network Technology in Support of Petascale Science

A Collaborative Initiative Among NASA, NLR, NOAA, Northwestern/iCAIR, SCinet & UIC/LAC
 Also Using Internet2's Multi-Vendor 100GigE Infrastructure Between StarLight and SC10

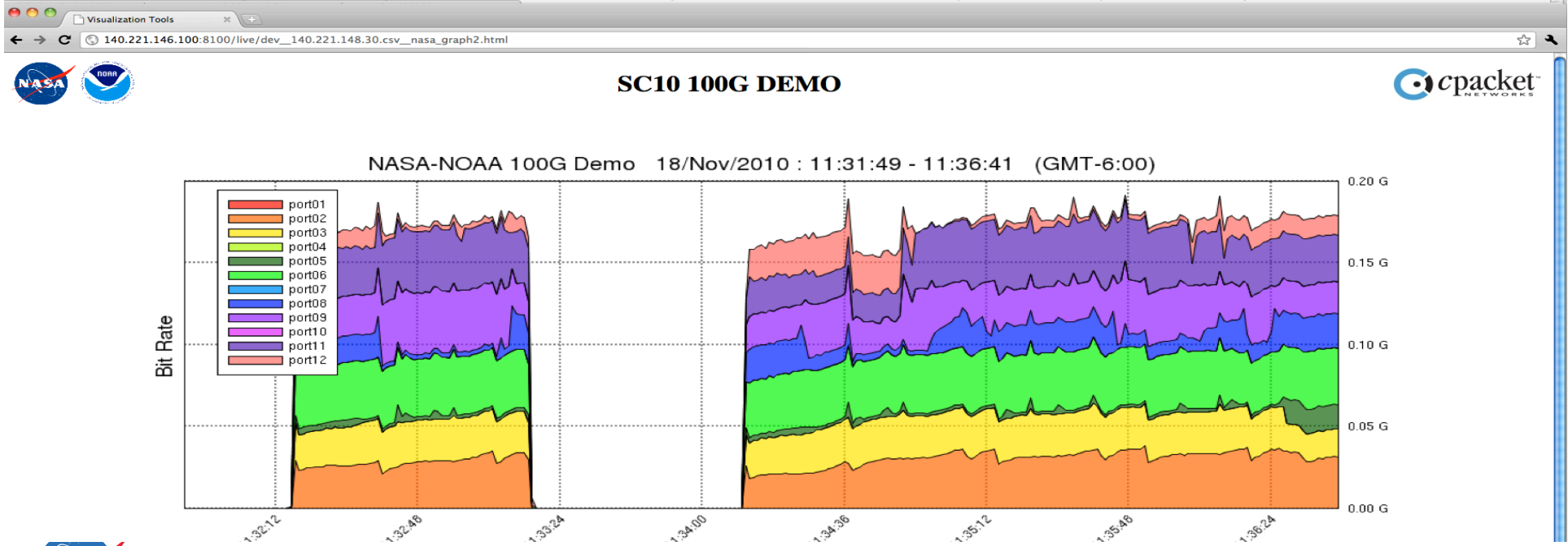
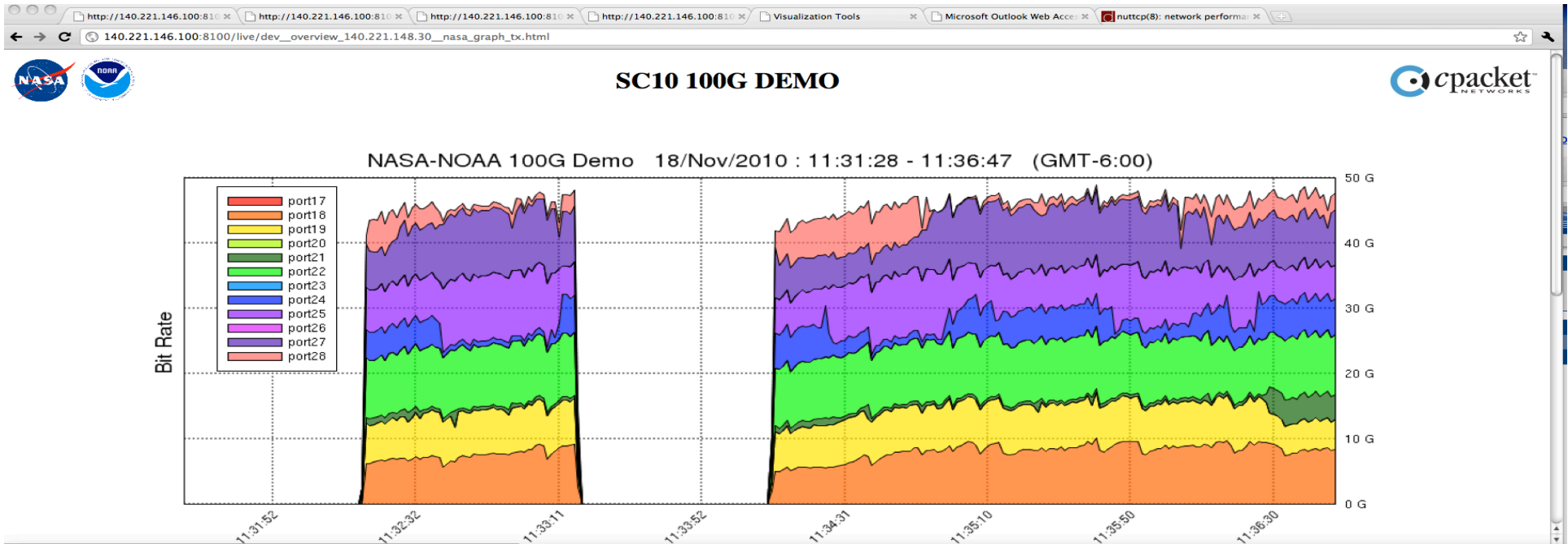


**80Gbps full-duplex across
 8x10G SCinet LAN &
 8x10G NLR WAN**

01/10/11

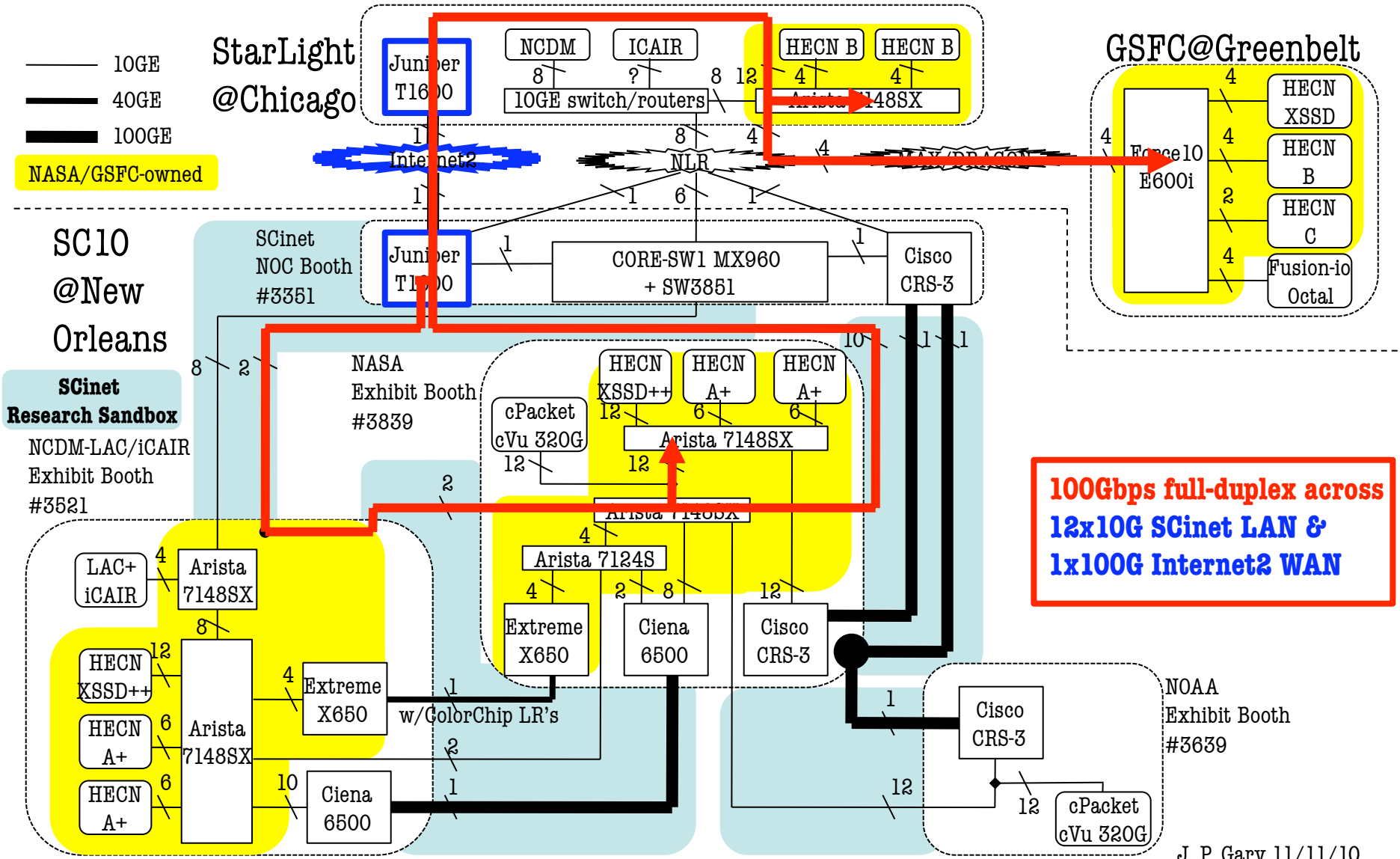
J. P. Gary

J. P. Gary 11/11/10



Using 100G Network Technology in Support of Petascale Science

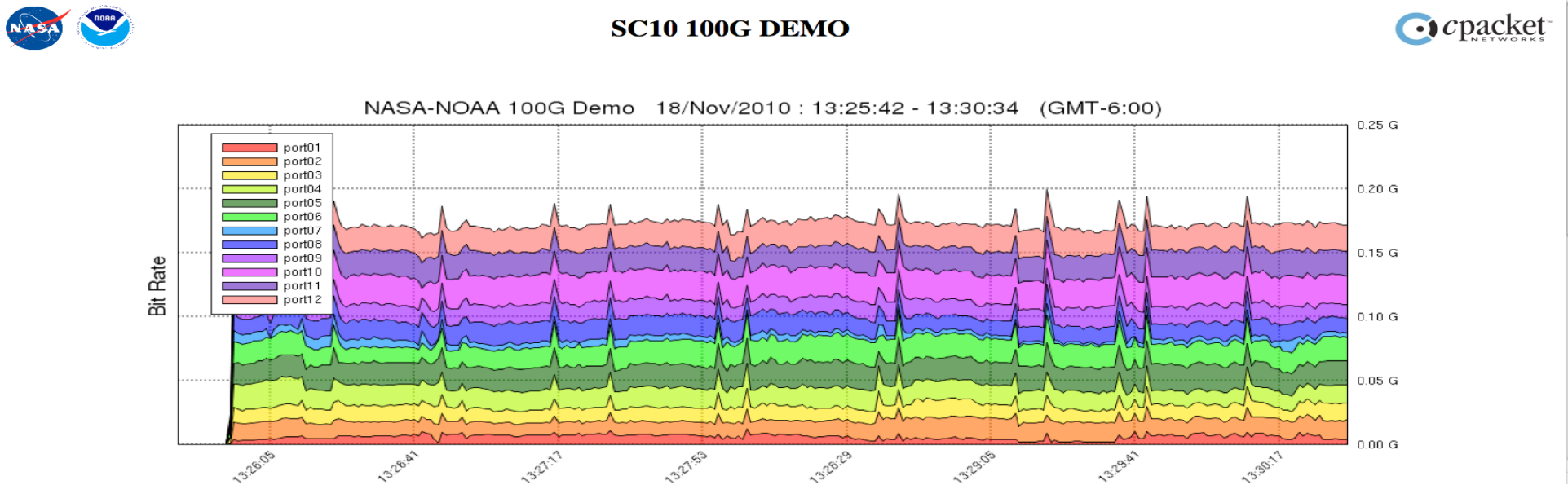
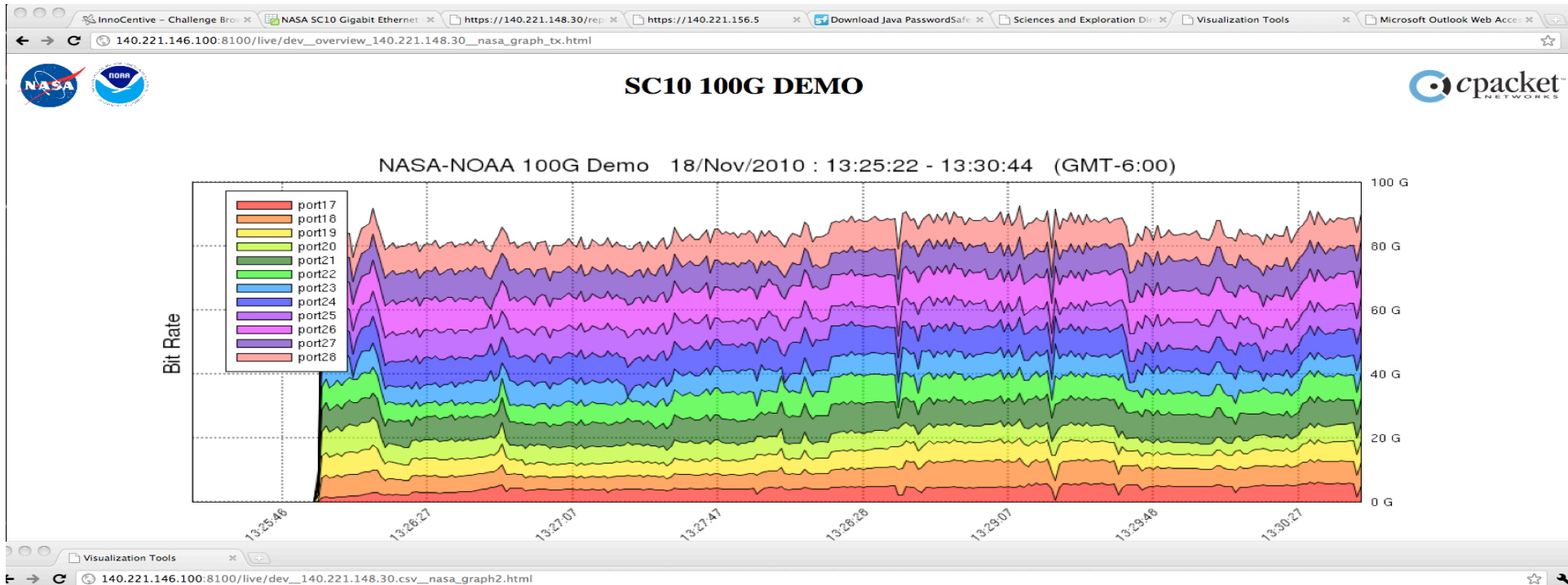
A Collaborative Initiative Among NASA, NLR, NOAA, Northwestern/iCAIR, SCinet & UIC/LAC
 Also Using Internet2's Multi-Vendor 100GigE Infrastructure Between StarLight and SC10



01/10/11

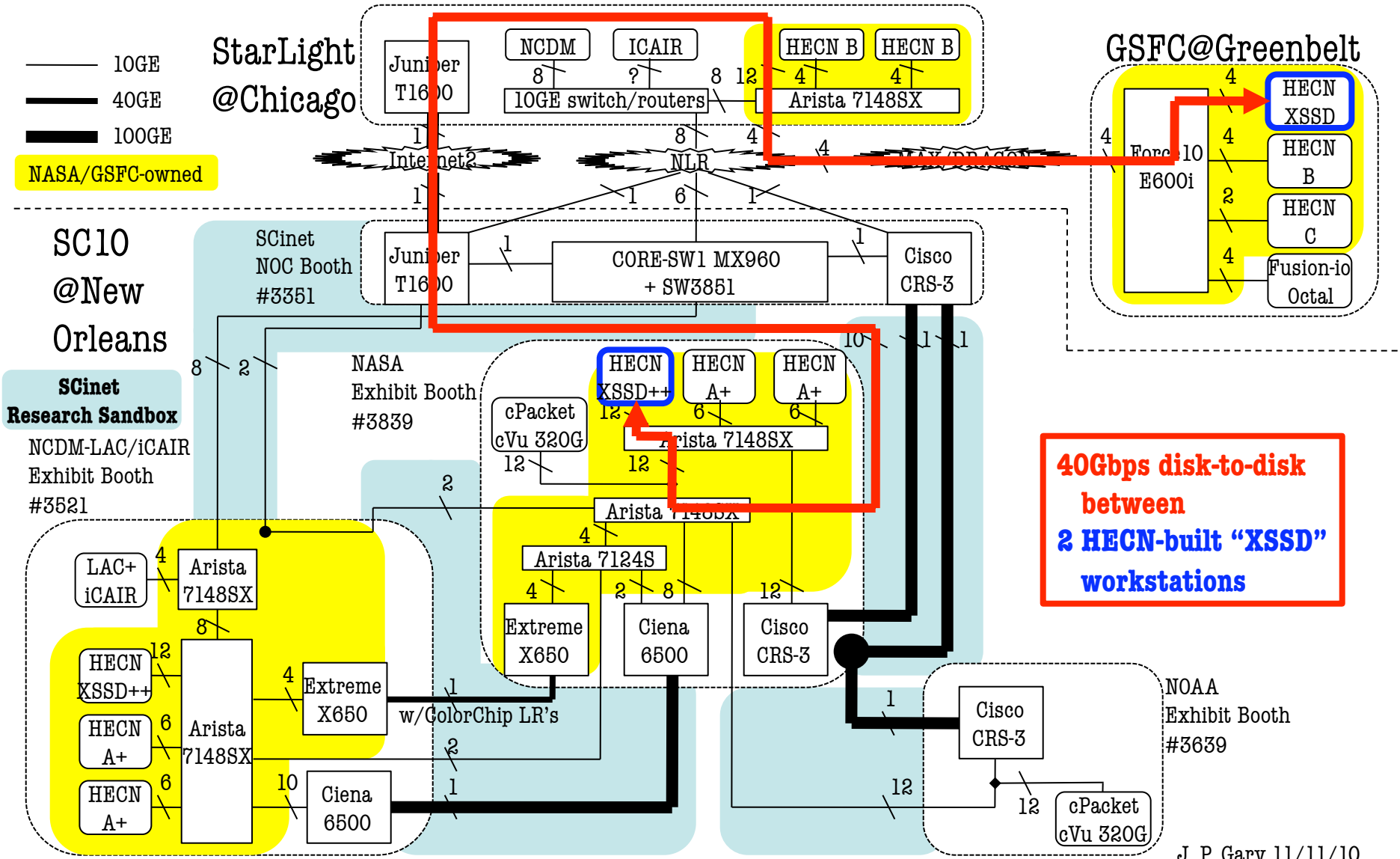
J. P. Gary

J. P. Gary 11/11/10



Using 100G Network Technology in Support of Petascale Science

A Collaborative Initiative Among NASA, NLR, NOAA, Northwestern/iCAIR, SCinet & UIC/LAC
 Also Using Internet2's Multi-Vendor 100GigE Infrastructure Between StarLight and SC10

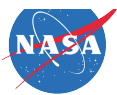
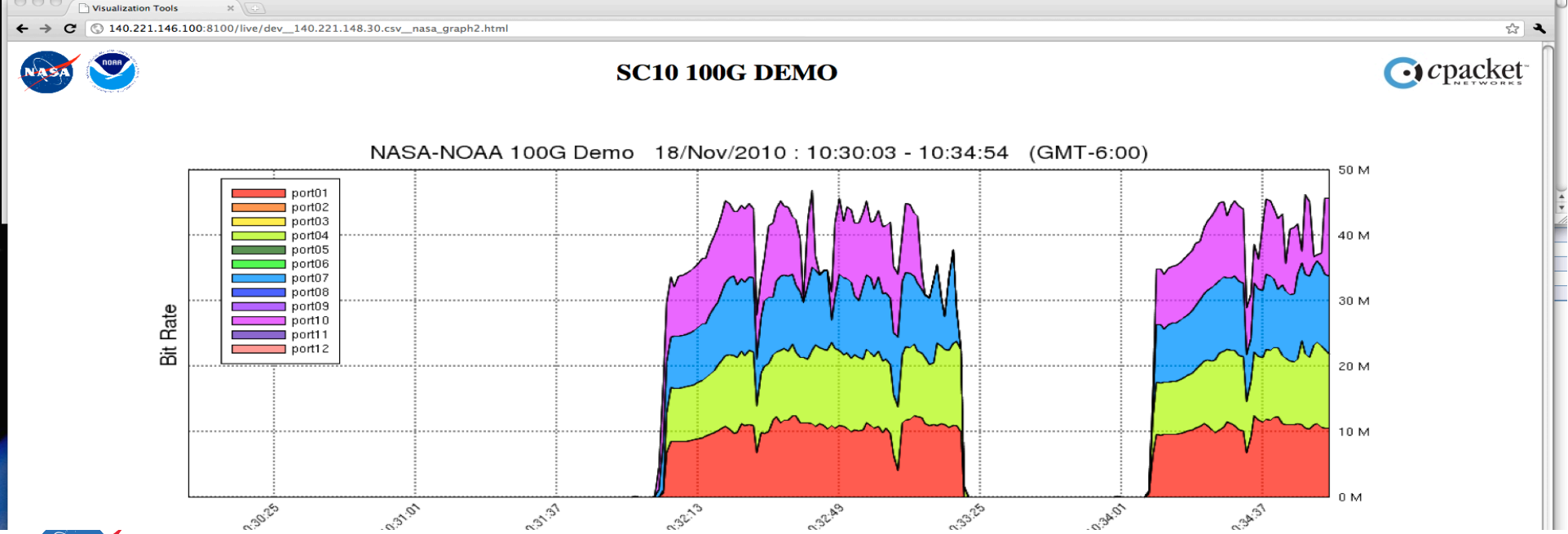
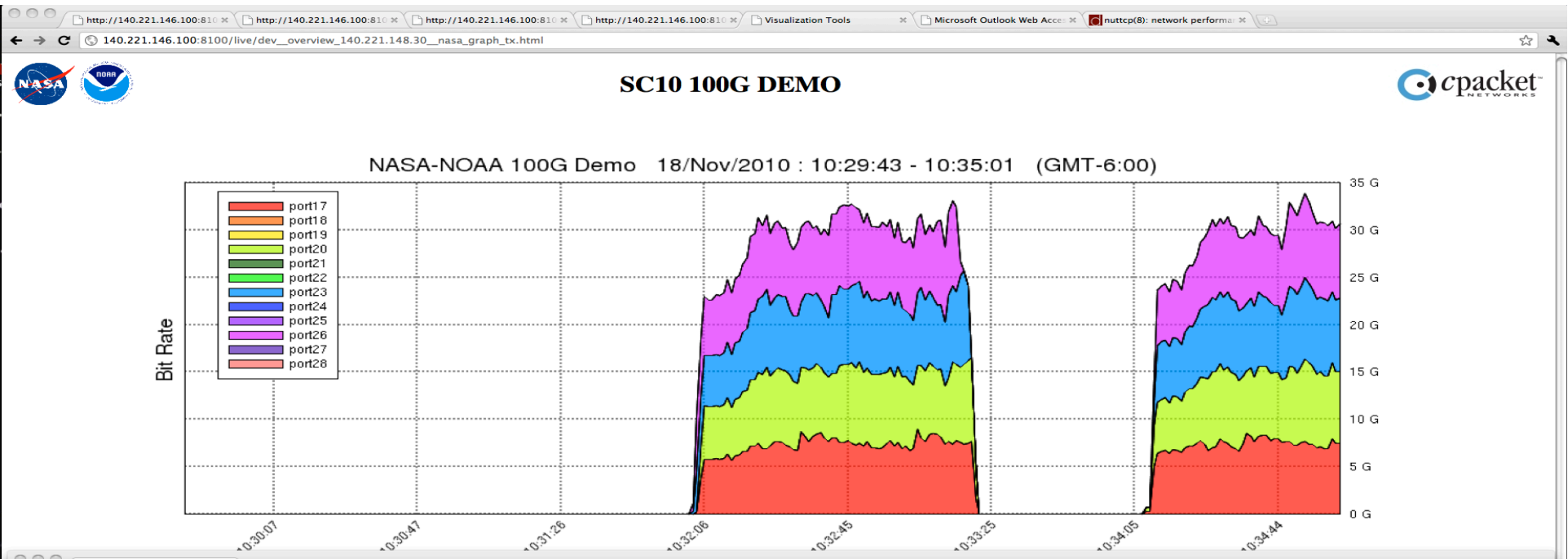


**40Gbps disk-to-disk
 between
 2 HECN-built "XSSD"
 workstations**

01/10/11

J. P. Gary

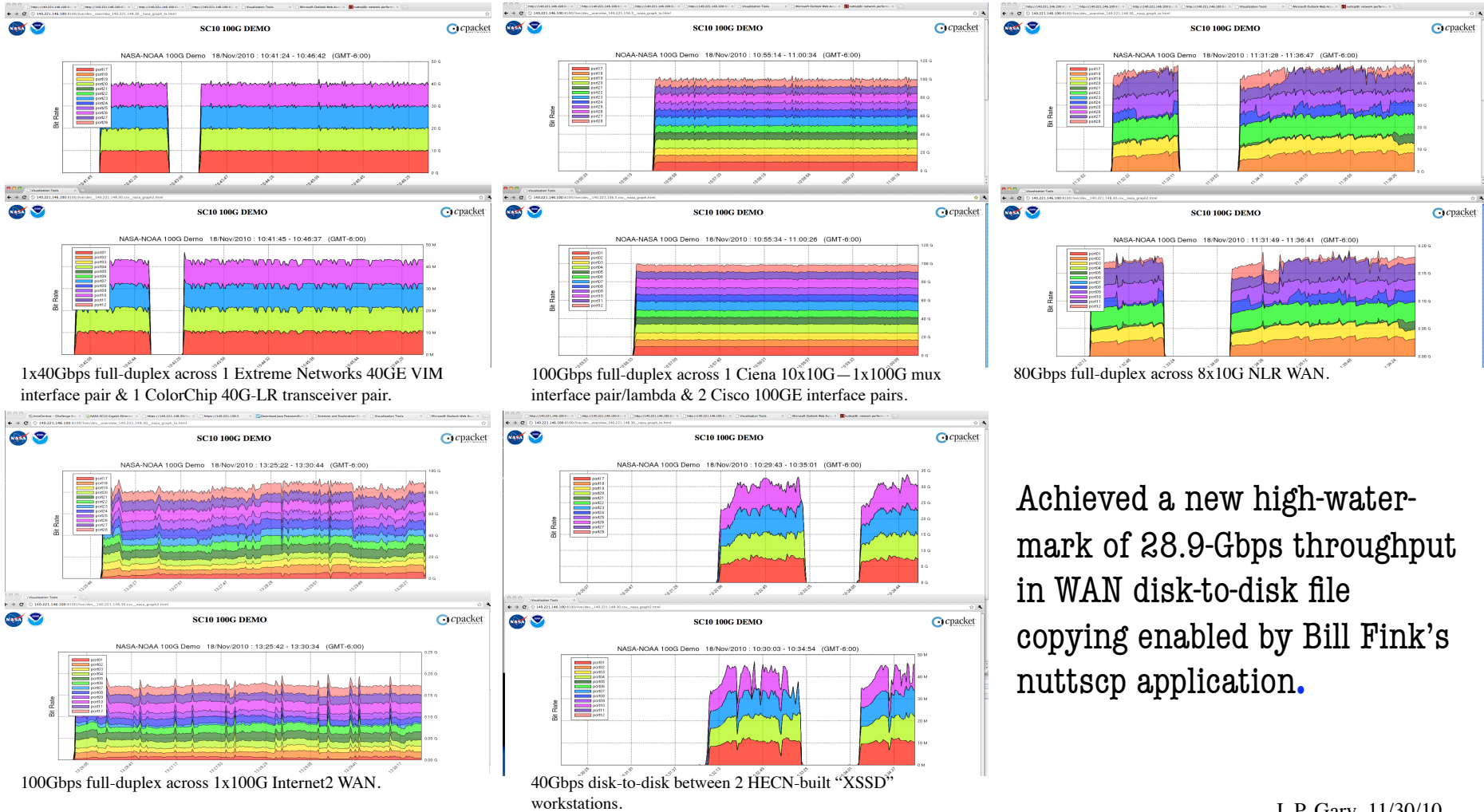
J. P. Gary 11/11/10



Snapshots of Test Results at NASA Exhibit Booth During SC10

NASA and Partners Demonstrate 40- and 100-Gigabit Network Technologies

http://science.gsfc.nasa.gov/606.1/HECN-highlights/HECN_SC10_Net-Demo_announce_110210.html



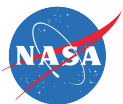
Achieved a new high-watermark of 28.9-Gbps throughput in WAN disk-to-disk file copying enabled by Bill Fink's nuttsep application.



Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Publications Referencing the NASA Experiments/ Demonstrations At SC10

- Our announcement
 - <http://www.hpcwire.com/offthewire/NASA-to-Demo-40100-Gigabit-Networking-at-SC10-106691913.html>
 - <http://supercomputingonline.com/latest/nasa-and-partners-to-demonstrate-40-and-100-gigabit-network-technologies-at-sc10>
- GCN:
 - http://gcn.com/articles/2010/11/15/update-1-nasa-40g-networking.aspx?sc_lang=en
- Vendor Press Releases
 - Ciena: <http://www.ciena.com/corporate/news-events/press-releases/Ciena-Puts-High-Performance-Networking-into-Action-at-SC10.html>
 - Extreme Networks: <http://investor.extremenetworks.com/releasedetail.cfm?ReleaseID=530684>
 - Internet2: <https://lists.internet2.edu/sympa/arc/i2-news/2010-11/msg00004.html>
 - NLR: <http://www.nlr.net/release.php?id=75>



01/10/11

GODDARD SPACE FLIGHT CENTER

J. P. Gary



Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Reference URL Summary

- Our announcement: "NASA and Partners to Demonstrate 40- and 100-Gigabit Network Technologies at SC10"
 - http://science.gsfc.nasa.gov/606.1/HECN-highlights/HECN_SC10_Net-Demo_announce_110210.html
- Our last updated detailed plan
 - http://science.gsfc.nasa.gov/606.1/docs/SC10_HECN-demos_111110.pdf
- “Collage” slides mostly of the NASA Exhibit Booth and network-demo rack setup and take down during SC10
 - http://science.gsfc.nasa.gov/606.1/docs/SC10-rack_stand-up_collage_113010.pdf
- Our current sample results
 - http://science.gsfc.nasa.gov/606.1/HECN-highlights/HECN_Net-Demo-SC10-Results_hilite_112910.html





Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Significance (Partial)

- Continuing and new partnerships
 - Key Non-NASA Collaborators for SC10: Internet2, NLR, NOAA, Northwestern University's iCAIR, SCinet Research Sandbox, University of Illinois at Chicago's LAC, University of Maryland College Park's MAX
 - Vendors who loaned equipment for SC10: Arista, Ciena, Cisco, ColorChip, cPacket, Extreme Networks, Fusion-io, HP, Panduit
- Usefulness of HECN net-test workstations
 - Isolating problems with, and then stress-testing, the data flow capabilities of leading-edge optical transport, Ethernet switch and Internet Protocol router equipment supporting 40- and 100-Gbps wire-speed rates - using Bill Fink's nuttcp-enabled memory-to-memory data transfers
 - Determining throughput limits in WAN disk-to-disk file copying – using Bill Fink's nuttscp application
- Preparing future plans

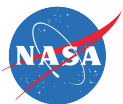




Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

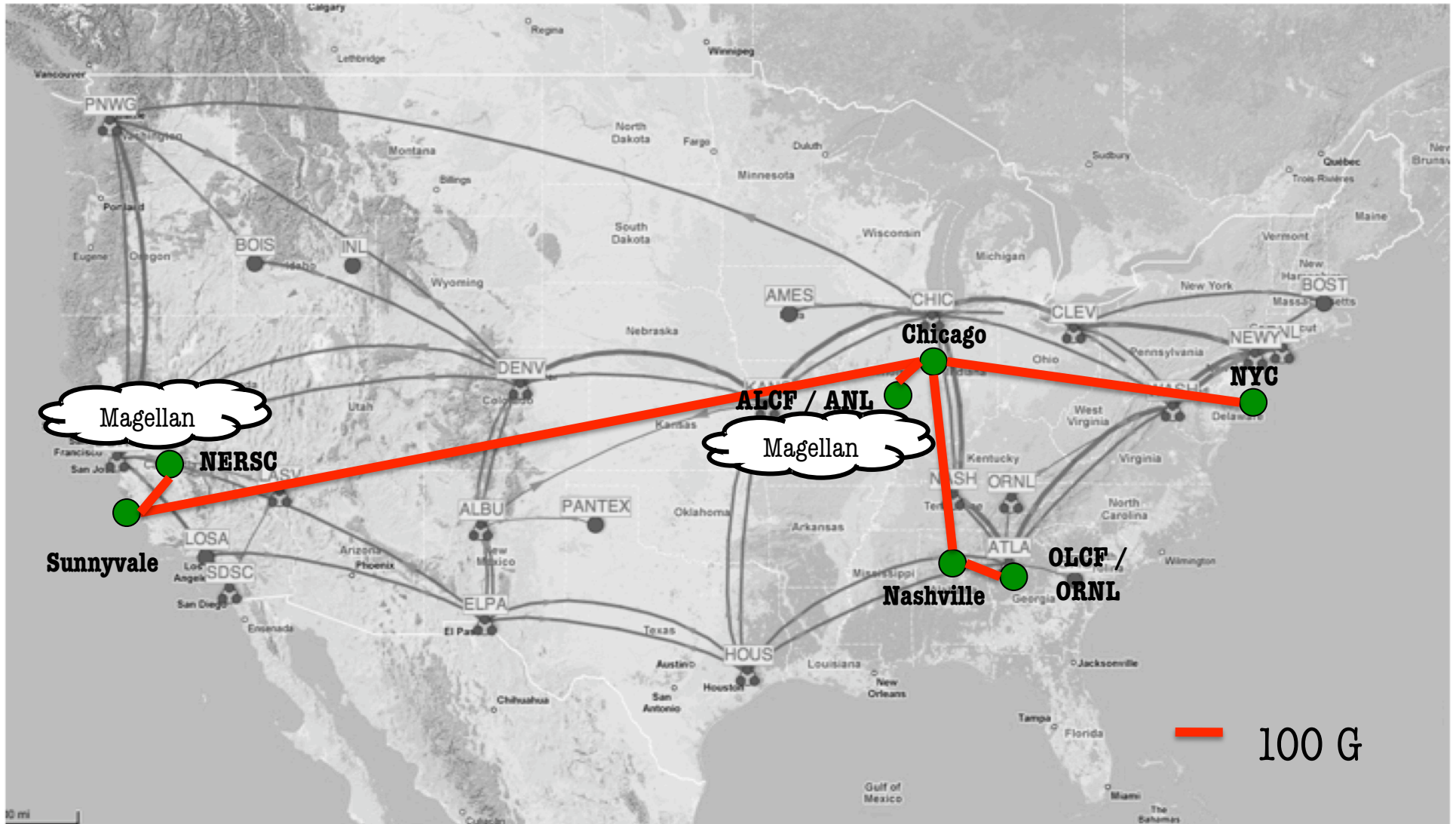
Preparing Future Plans

- Need to share and apply knowledge gained in achieving greater than the Phase 1 Goal of GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds initiative: 20-Gbps WAN disk-to-disk user-throughput
 - NASA: NCCS climate simulation data flows for IPCC AR5
 - NOAA: Extreme data flows between GFDL and ORNL
 - DoE: Assist in assessing throughput performance of DoE Advanced Network Initiative (ANI)'s 100G Prototype Network
- Need to refine approaches to achieving the respective Phase 2 & 3 Goals of GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds initiative: 40- & 100-Gbps WAN disk-to-disk user-throughput
 - File copying applications
 - Workstations with NICs and disk controllers using PCIe Gen 3
 - WAN test infrastructure



Source: Brian Tierney (LBL/ESnet)

Nationwide 100G Prototype Network





Overall Timeline for the HEC 20, 40 & 100 Gbps Network Testbeds

| FY10 | FY11 | FY12 | FY13 | FY14 | FY15

4x10G

A. LAN+MAN

Testing tttttttttttttttttt

B. WAN betw GSFC+StarLight

Testing TTTTTTTTTTTTTTTTTT

1x40G

A. LAN+MAN

Testing ttttttttoooooooooooooooo

B. WAN betw GSFC+StarLight

Testing TTTTTTTTTTTTTTTTTT

C. WAN betw ARC+GSFC

Operations OOOOOOOOOOOOOOOOOO

1x100G

A. LAN+MAN

Testing ttttttttoooooooooooooooo

B. WAN betw GSFC+StarLight

Testing TTTTTTTTTTTTTTTTTT

C. WAN betw ARC+GSFC

Operations OOOOOOOOOO

Wherein:

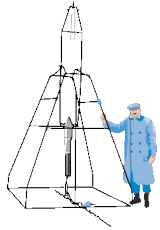
ttt...ttt implies only pre-operational testing use

ooo...ooo implies operational use in LAN+MAN

TTT...TTT implies only pre-operational testing use in WAN

OOO...OOO implies operational use in WAN

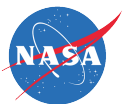




Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Preparing Future Plans – Near Term (1 of 2)

- Replace HECN's two "B" net-test workstations at StarLight
 - Upgrade to two XSSDs: each >100G memory-to-memory, ~30G disk-to-disk
 - Assist DoE with checkout of their ANI 100G Prototype Network
 - Assist StarLight with checkout of their NSF-awarded 100G-capable upgrades in part of StarLight's infrastructure
- Upgrade DRAGON's lambda pathways between GSFC and McLean in cooperation with the MAX
 - Replace ADVA-based 10-Gbps DWDM
 - Assist MAX with checkout of their NSF-awarded 100G-capable upgrades in part of MAX's infrastructure
- Leverage Internet2's and/or NLR's 100G pathways between McLean and StarLight
 - <https://lists.internet2.edu/sympa/arc/i2-news/2010-11/msg00003.html>
 - <http://www.nlr.net/release.php?id=62>

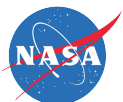




Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Preparing Future Plans – Near Term (2 of 2)

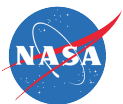
- HECN also will participate in DoE's February 2011 by-invitation-only Terabit Networks for Extreme Scale Science Workshop
 - To identify the major research challenges in developing and deploying end-to-end high-capacity networks to support distributed extreme scale science
 - In the next decade, the Office of Science at the US Department of Energy anticipates deploying exascale computers; storing, distributing, and analyzing massive data sets estimated in the zetabyte-scale generated by complex simulations and large instruments; and engaging in extreme scale scientific activities involving thousands of distributed researchers
 - DOE envisions the development of a new generation of networks with unprecedented end-to-end performance. **An early implementation of this vision calls for Energy Sciences Network (ESnet) to deploy a terabit network by 2014**





Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Q & A





Summary Results From NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

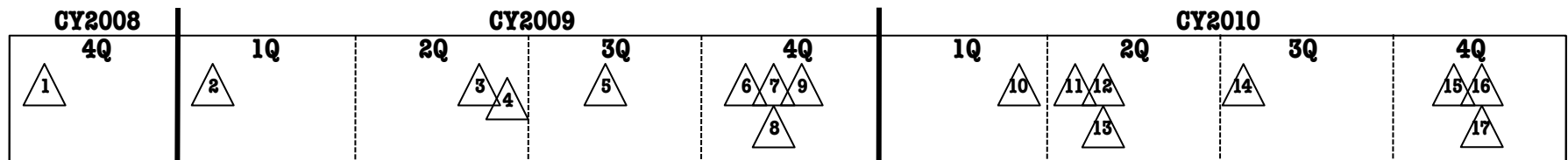
Backup Slides





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Significant >10G Accomplishments of HECN Team & Partners (1 of 7)



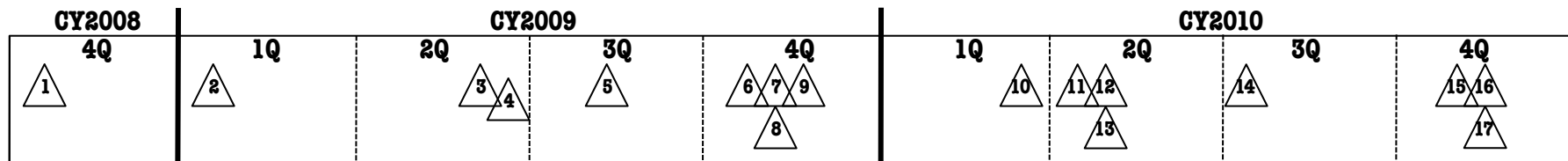
- 1: 10/14/08 In collaboration with MAX across MAX-provisioned link between College Park and McLean, bi-directionally tested maximum throughput of Fujitsu 40G optical transponders & Juniper OC-768c interfaces on T1600 routers
- 2: 01/29/09 Created initial drafts of HECN 20, 40 & 100G Network Testbed plans
- 3: 06/12/09 Bench-tested HECN's first >10G net-test-workstations (i7-based) measuring nuttcp-enabled memory-to-memory throughput flows unidirectionally at >69.2G and bi-directionally at >77.2G aggregate
- 4: 06/24/09 Presented "Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds" at MAX Spring 2009 All Hands Meeting



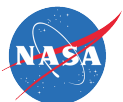


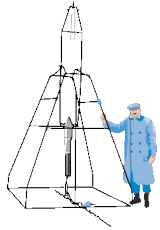
Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Significant >10G Accomplishments of HECN Team & Partners (2 of 7)



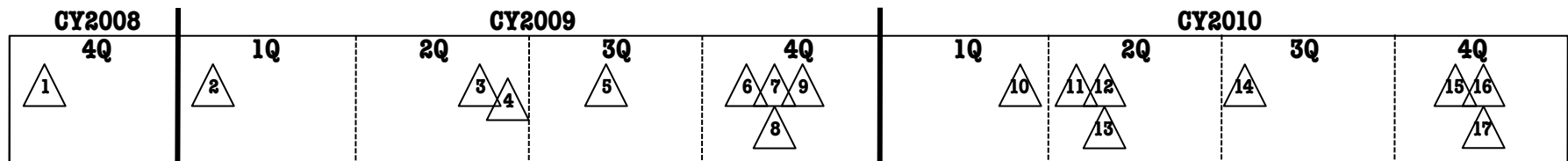
- 5: 08/06/09 Bench-tested a new HECN net-test-workstation (Xeon-based) measuring nuttcp-enabled memory-to-memory throughput flows unidirectionally (transmit) at >100.4G
- 6: 10/30/09 Two HECN net-test-workstations (i7-based) each with 4x10G NIC interfaces, deployed at ARC, used by NREN in prep for SC09 to fully check out new 4x10G links on single fiber pair between ARC and Sunnyvale using ADVA FSP3000 dwdm mux/demuxes
- 7: 11/02/09 Using two HECN net-test-workstations (i7-based) each with two RAID5 disk controllers nested as RAID50, measured nuttscp-enabled disk-to-disk data throughput unidirectionally at >9.8G in a 0.1ms RTT testbed; on 11/10/09 measured >9.5G in a 80.1ms RTT testbed





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Significant >10G Accomplishments of HECN Team & Partners (3 of 7)



8: 11/02/09 Using two HECN net-test-workstations (i7-based) in nuttcp-enabled memory-to-memory throughput flows, with kernel bonding/standard link-aggregation among each workstation's 4x10G NIC interfaces measured unidirectionally >31.6G and with nuttcp-enabled "application bonding" measured unidirectionally >39.5G

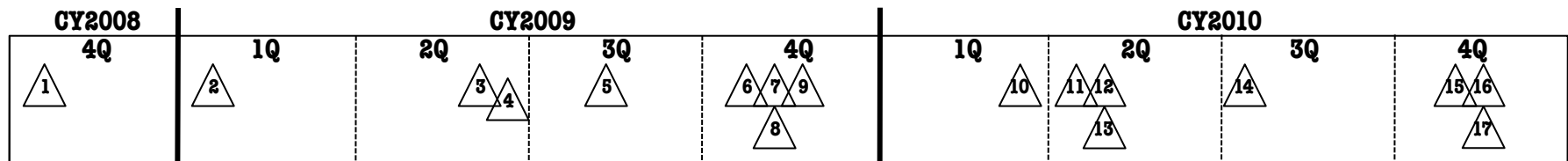
9: 11/16/09 In the NASA research exhibit at the SC09 conference, Portland, OR, demoed: >100G uni-directional memory-to-memory data throughput between in-booth HECN servers; 40G bi-directional memory-to-memory data throughput between HECN servers in-booth and at ARC across 4x10G NLR/C-Wave links; and 10G disk-to-disk data throughput between in-booth HECN servers, between HECN servers in-booth and at ARC, and between HECN servers in-booth and at GSFC



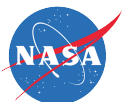


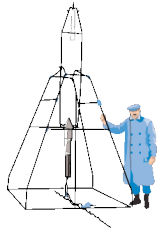
Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Significant >10G Accomplishments of HECN Team & Partners (4 of 7)



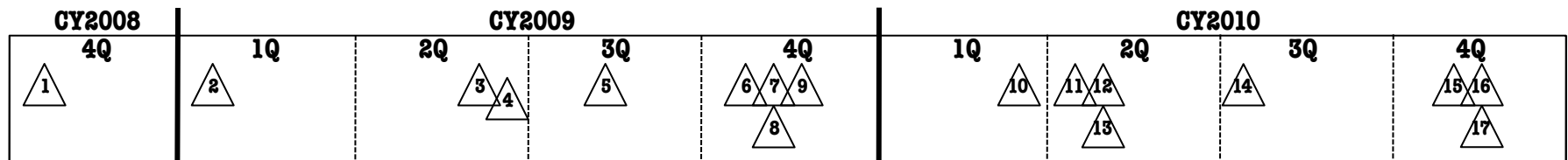
- 10: 03/24/10 In collaboration with MAX, added 4x10G links across MAX's DRAGON regional optical network between GSFC and the NASA rack in the Level3 PoP at McLean, and then through MAX's ADVA dwdm mux/demuxes to NLR's Cisco 15454 optical transponder at that PoP
- 11: 04/07/10 In collaboration with iCAIR, deployed two HECN net-test-workstations (i7-based) at StarLight for 4x10G connections with NLR and 8x10G connections with other R&D networks
- 12: 04/14/10 In collaboration with NREN, used two HECN net-test-workstations (i7-based) at ARC to measure unidirectional and bidirectional throughput of a Fortinet FortiGate-3810A security gateway with four 10GE ports





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Significant >10G Accomplishments of HECN Team & Partners (5 of 7)



13: 05/09/10 NLR turned up new 4x10G links between McLean and StarLight for HECN maximum throughput testing; but due to errors that ultimately were isolated to and fixed on HECN equipment at MCLN, HECN did not “accept” the NLR links as active until 08/03/10. In between HECN conducted limited 40G bi-directional memory-to-memory and 10G unidirectional disk-to-disk data throughput between HECN servers at GSFC and StarLight

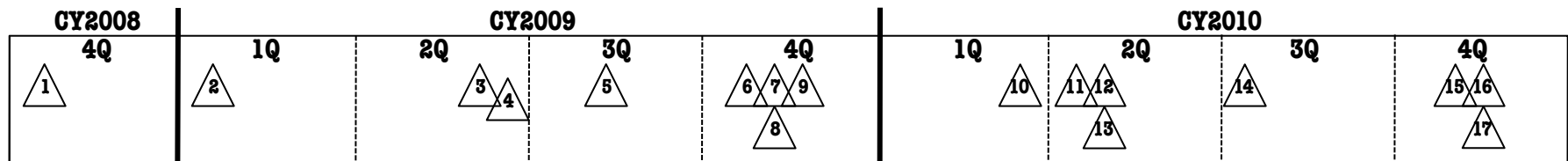
14: 08/14/10 Using two HECN net-test-workstations (i7-based) each with four RAID5 disk controllers nested as two separate RAID50s, with each disk controller hosting eight rotating disks, measured nuttscp-enabled aggregate disk-to-disk data throughput unidirectionally at >17.8G (90% of maximum 19.8G) in a 0.1ms RTT testbed



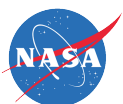


Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Significant >10G Accomplishments of HECN Team & Partners (6 of 7)



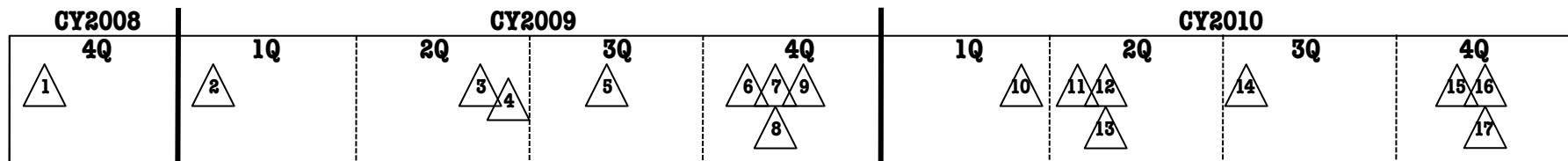
- 15: 11/01/10 Using two HECN XSSD net-test-workstations (Xeon-based) each with four RAID5 disk controllers nested as two separate RAID50s, with each disk controller hosting eight SSD, measured nuttscp-enabled aggregate disk-to-disk data throughput unidirectionally at >30.4G in a 0.1ms RTT testbed
- 16: 11/18/10 Using two HECN XSSD net-test-workstations (Xeon-based) each with four RAID5 disk controllers nested as two separate RAID50s, with each disk controller hosting eight SSD, measured nuttscp-enabled aggregate disk-to-disk data throughput unidirectionally at >28.9G in a ~45.0ms RTT pathway between GSFC and New Orleans during SC10





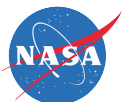
Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Significant >10G Accomplishments of HECN Team & Partners (7 of 7)



17: 11/18/10 During SC10 HECN net-test-workstations were used to evaluate the maximum data flow capabilities of:

- Ciena: 100-Gbps transport and 10x10Gbps-to-1x100Gbps muxponder interfaces of Optical Multiservice Edge 6500 units
- Cisco: 100-GE interfaces of CRS-3 switch/routers
- ColorChip: 40-Gbps DragonFly 40G-LR (up to 10km) QSFP transceivers (beta)
- Extreme Networks: 40-GE VIM3-40G4X modules (beta) for Summit X650 10-GE switches
- Juniper: 100-GE interfaces of T1600 routers in Internet2 pathway





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Phase 1 Network-Test Workstations: Nominal “C” System

- Nominal “B” (Baseline) System
- Minus:
 - NICs: One Myricom 10G-PCIE2-8B2-2S+E (Dual 10GE SFP+)
 - IB HCA: one Voltaire (DDR, 8x)
- Plus:
 - Raid controllers: two HighPoint RocketRaid 4322 (external, 8 disks each)
- Plus via SAS-connection:
 - Chassis: one Supermicro 836TQ-R800B (3u 16bay 7slot 800W RPS) with SAS converter/adaptor and cables
 - User disks: 16 Western Digital WD5001AALS (500GB)
- For more detail, contact Paul.Lang@nasa.gov

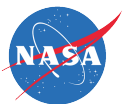




Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Phase 1 Network-Test Workstations: Nominal “A” System

- Nominal “B” (Baseline) System
- Minus:
 - Raid controllers: two HighPoint RocketRaid 4320 (internal, 8 disks each)
 - User disks: 16 Western Digital WD5001AALS (500GB)
 - IB HCA: one Voltaire (DDR, 8x)
- For more detail, contact Paul.Lang@nasa.gov





Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Phase 1 Network-Test Workstations: “A+” System

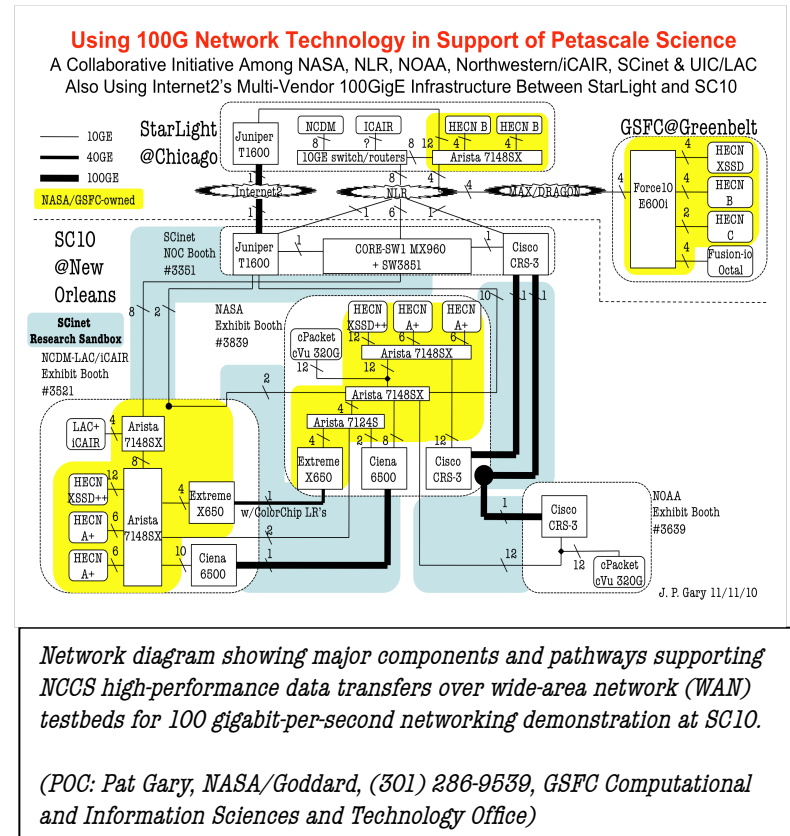
- Nominal “A” System
- Plus:
 - NICs: Two Myricom 10G-PCIE2-8B2-2S+E (Dual 10GE SFP+)
- For more detail, contact Paul.Lang@nasa.gov





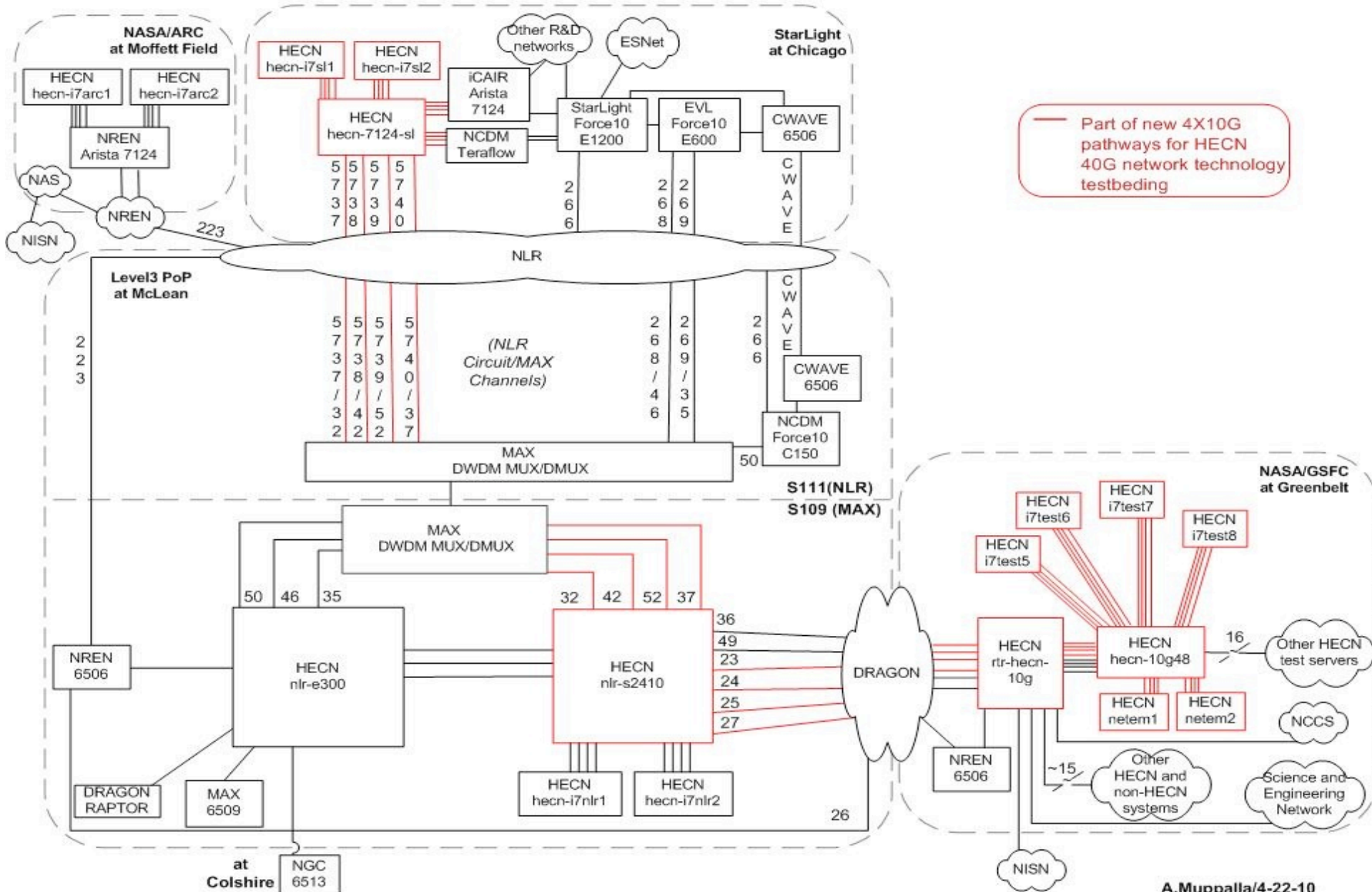
NASA and Partners Demonstrated 40- and 100-Gigabit Per Second Network Technologies at SC10

- Workstations assembled by GSFC's High End Computer Network (HECN) Team (Code 606.1) stress-tested the data flow capabilities of several new leading-edge optical transport, Ethernet switch and Internet Protocol router equipment just beginning to support 40 and 100 gigabit per second (Gbps) wire-speed rates.
- Staff used HECN workstations to isolate problems with, and then saturate, various 40-, 80- and 100-Gbps links with Bill Fink's nuttcp-enabled memory-to-memory data transfers, and also achieved a new 28.9-Gbps high-water-mark in wide area network (WAN) disk-to-disk file transfers.
- SC10 stress-tests were conducted with network equipment
 - Located at SC10 among the Exhibit Booths of NASA, the National Center for Data Mining, NOAA and the SCinet Network Operations Center across 36 inter-booth single-mode fiber-pairs.
 - Across special SC10-only 80- and 100-Gbps WAN links between the NASA Exhibit Booth in New Orleans and the StarLight Communications Exchange facility in Chicago (provisioned respectively by the National LambdaRail and Internet2), then on to GSFC across a 40-Gbps link provisioned by the Mid-Atlantic Crossroads (MAX).
- Success of the SC10 demonstrations is strengthening the NASA's HECN Team's relationships with its existing and new partners, e.g., NOAA representatives now have requested HECN Team's assistance in NOAA's extreme data flows between the Geophysical Fluid Dynamics Laboratory and NOAA's HPC resources to be located at Oak Ridge National Laboratory.



GSFC/High End Computer Network (HECN) and Partners 10GE and 10G Lambda Connections Through McLean

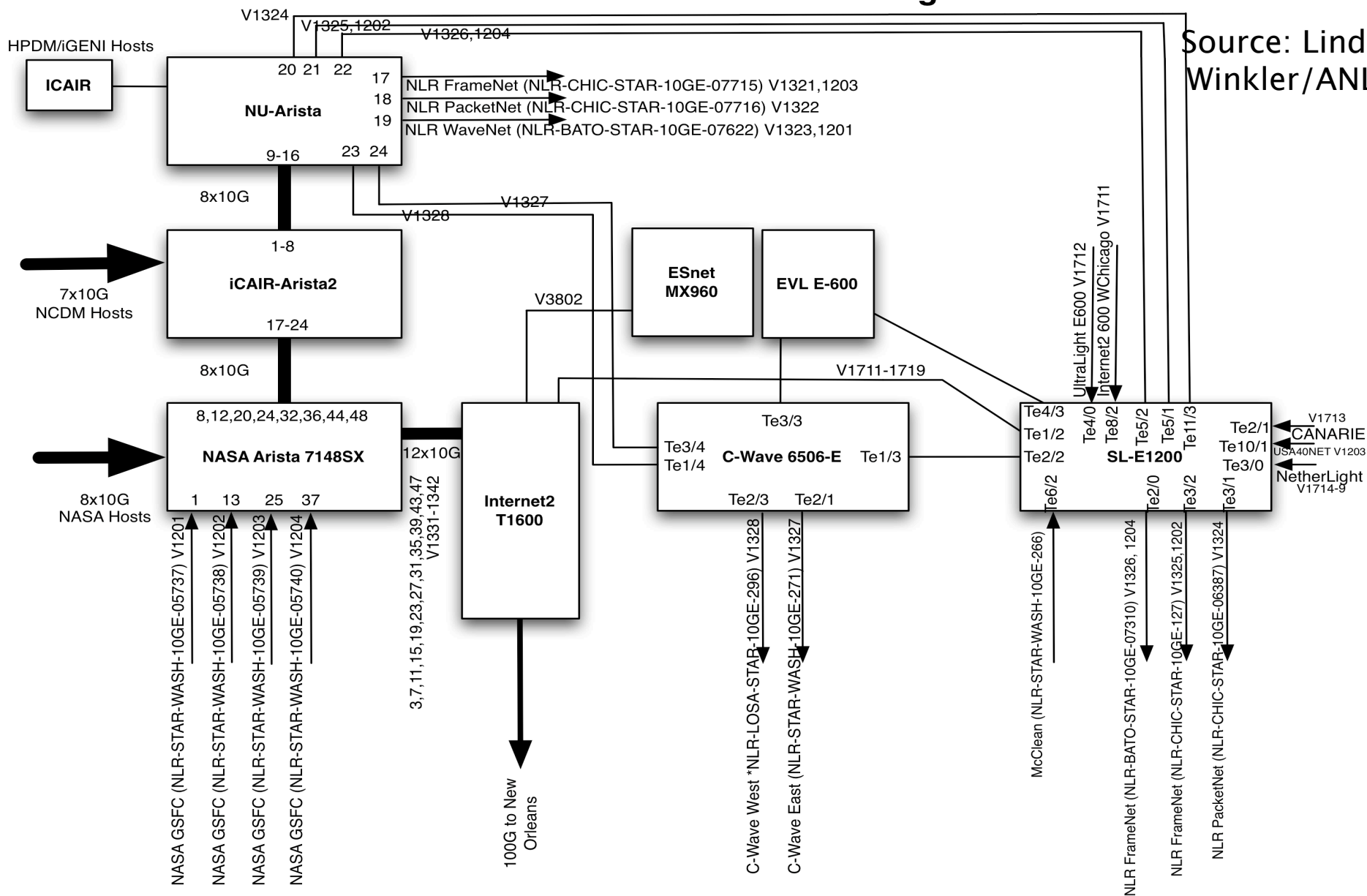
Note: The non-GSFC/HECN systems shown typically have other connections that are not shown in this diagram, as the focus is primarily GSFC/HECN connections



ICAIR/NCDM/NASA Resources @ StarLight for SC10

L. Winkler 11/10/2010

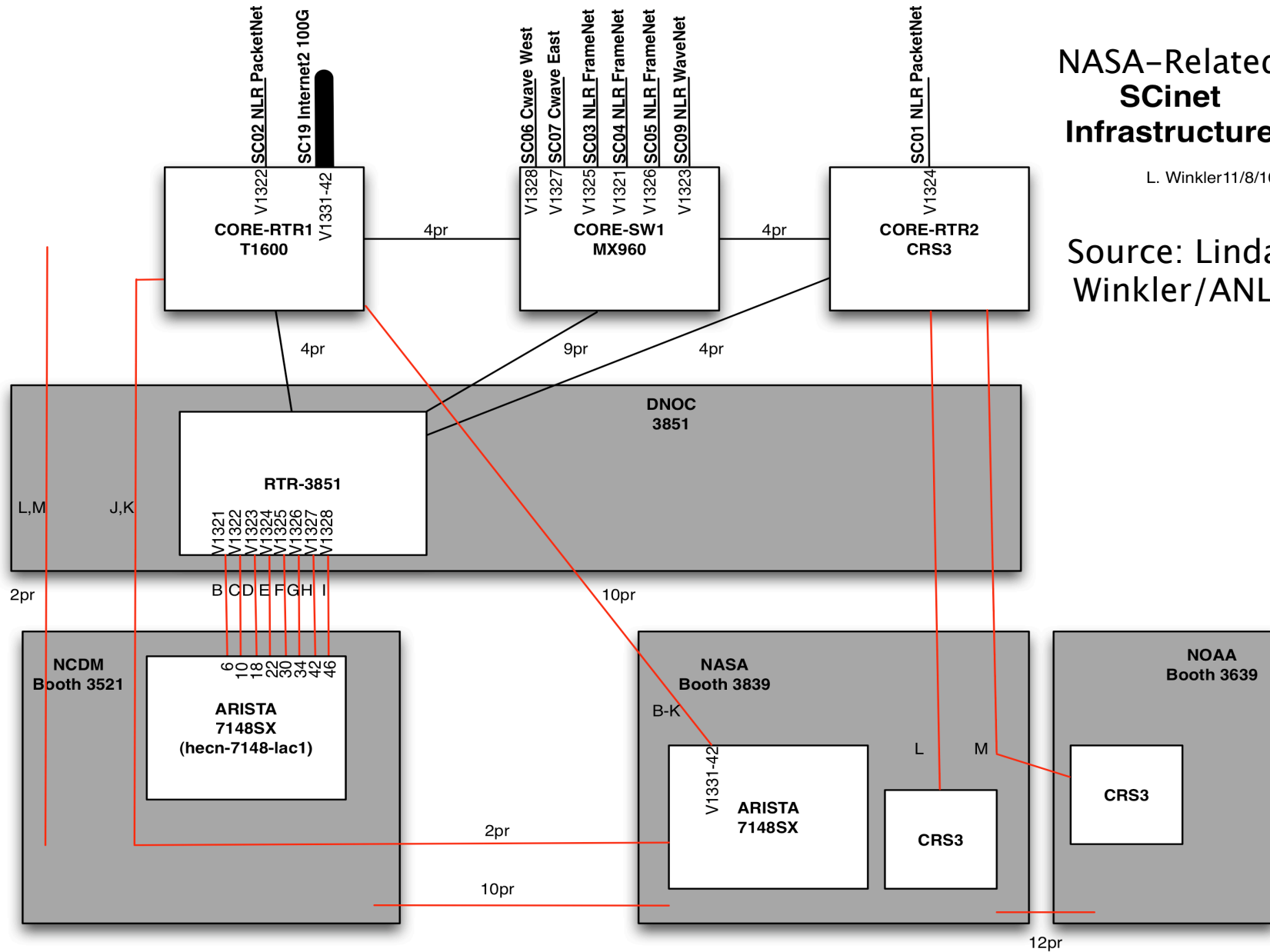
Source: Linda Winkler/ANL

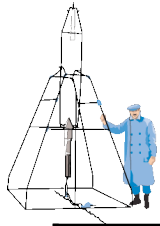


NASA-Related SCinet Infrastructure

L. Winkler11/8/10

Source: Linda Winkler/ANL





Test Matrix for Optimizing Wide-Area File Transfer

Tool	Type
aspera fasp**	Disk ↔ Disk
nuttscp	Disk ↔ Disk
GridFTP	Disk ↔ Disk
RocketStream**	Disk ↔ Disk
iRODS	Disk ↔ Disk
bbFTP	Disk ↔ Disk
FDT	Disk ↔ Disk
HPN-SCP	Disk ↔ Disk
dsync*	Disk ↔ Disk
xdd	Disk ↔ Disk
rsync	Disk ↔ Disk
FTP	Disk ↔ Disk
gpfs	Mem/Disk ↔ Disk
pNFS	Mem/Disk ↔ Disk
NFS	Mem/Disk ↔ Disk
NFS-RDMA*	Mem/Disk ↔ Disk
nuttcp	Mem ↔ Mem

*RDMA options – IB/Obsidian, iWARP, RoCE

**Commercial product

- Desired measurements: disk-to-disk file-copying-throughput performance (in Gbps), plotted against different file-sizes and different conditions
- Key single-file-sizes in GBs: 16, 32, 64, 128
- Primary different conditions:
 - File-copying-applications, e.g., GridFTP, bbFTP, nuttscp, ...
 - Both well-established and experimental/emerging ones
 - Key round trip times (RTTs) in milliseconds: 0, 15, 90, 180
 - Corresponding very roughly to LAN, large MAN/ROn, trans-USA/WAN, trans-Atlantic
- Secondary different conditions, when time permits:
 - Several-file-sizes yet with a constant 256 GB total volume: 16@16, 8@32, 4@64, 2@128
 - Many-file-sizes yet with a constant total volume
 - Cases with packet loss, corruption, etc
 - Real tests
 - @10Gbps: GSFC-GSFC (RTT=0); GSFC-StarLight (RTT= ~ 17); GSFC-ARC (RTT= ~ 87)
 - @20Gbps: GSFC-GSFC; GSFC-StarLight
 - @40Gbps: GSFC-GSFC; GSFC-SC10 (RTT= ~ 35)

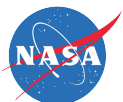




Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

Elaboration of Test Matrix for Optimizing Large File Transfers Over Wide Areas (1 of 2)

- Desired measurements (y-axis): disk-to-disk file-copying-throughput performance (in Gbps), plotted against different file-sizes and different conditions
- Key single-file-sizes in GBs: 16, 32, 64, 128
- Primary different conditions:
 - File-copying-applications, e.g., GridFTP, bbFTP, nuttscp, ...
 - Both well-established and experimental/emerging ones
 - Key round trip times (RTTs) in milliseconds: 0, 15, 90, 180
 - Corresponding very roughly to LAN, large MAN/RON, trans-USA/WAN, trans-Atlantic



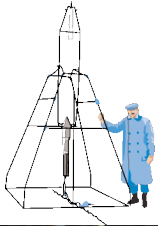


Introduction To NASA HEC WAN File Accessing Experiments/Demonstrations At SC10

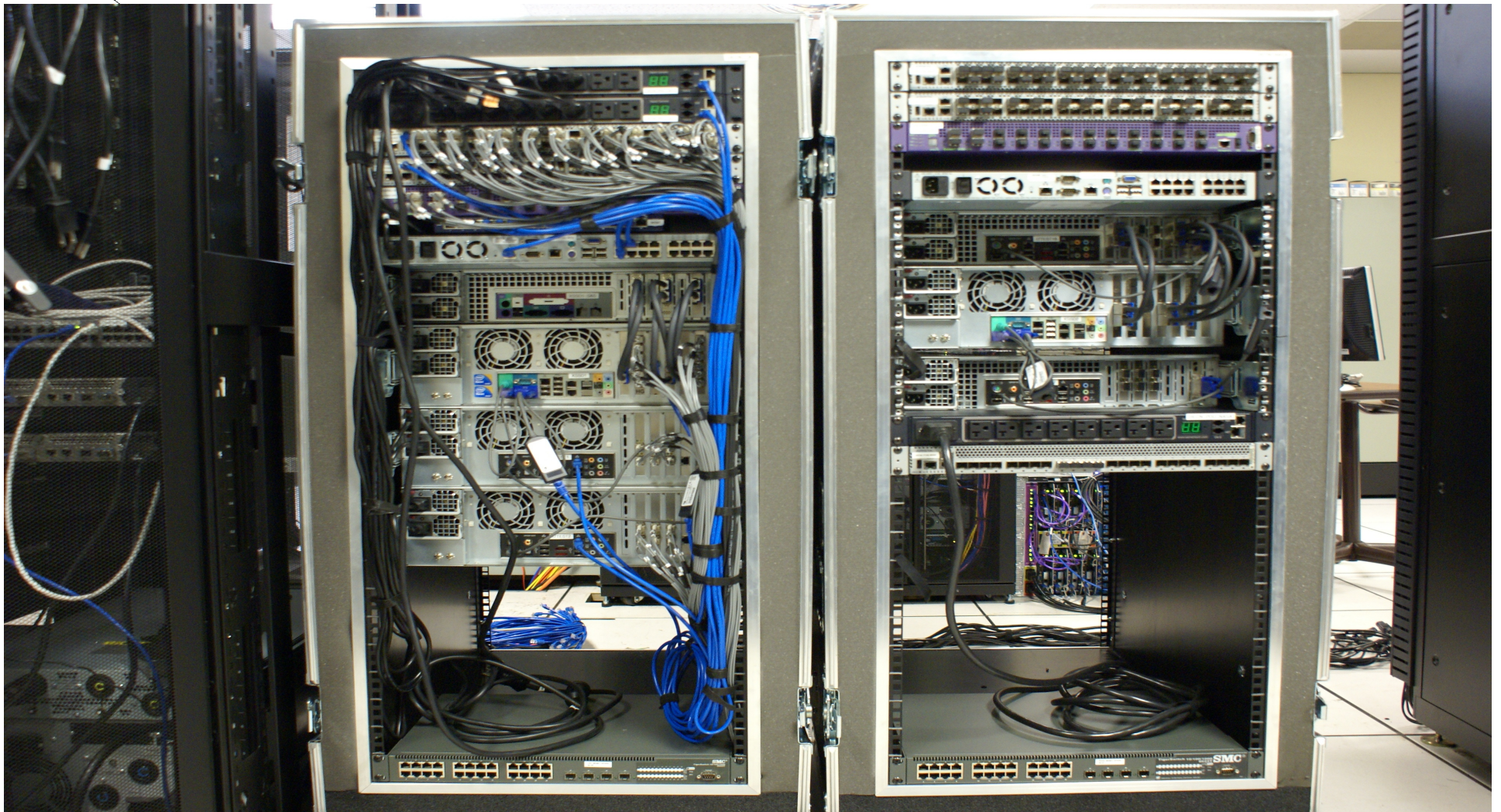
Elaboration of Test Matrix for Optimizing Large File Transfers Over Wide Areas (2 of 2)

- Secondary different conditions, when time permits:
 - Several-file-sizes yet with a constant 256 GB total volume: 16@16, 8@32, 4@64, 2@128
 - Many-file-sizes yet with a constant total volume, e.g., a single directory of small-to-medium-to-partly-large size files
 - Cases with packet loss, corruption, etc, particularly if they can be controlled by netem and “be realistic”
 - Real tests near-term
 - @10Gbps: GSFC-GSFC (RTT=0); GSFC-StarLight (RTT=~17); GSFC-ARC (RTT=~87)
 - @20Gbps: GSFC-GSFC; GSFC-StarLight
 - Real tests very soon
 - @40Gbps: GSFC-GSFC; GSFC-SC10 (RTT=~35)





NASA/GSFC High End Computer Network Team's Equipment Ready for Shipment to SC10



01/10/11
GODDARD SPACE FLIGHT CENTER

J. P. Gary

83

Example Showing NASA's Daily Scheduled Use of NLR During SC10

Source: Bonnie Hurst/NLR

All Times are Central	NLR L1 WaveNet_STAR EXPRESS (STAR to NEW ORLEANS)	NLR L2 FrameNet Wave 1_Western Route	NLR L2 FrameNet Wave 2_Eastern Route	NLR L2 FrameNet_STAR EXPRESS (STAR to NEW ORLEANS)	Cisco CWAVE West (LOSA)	Cisco CWAVE East (JACK)	NLR L3_PacketNet Wave 1_Western Route (HOUS)	NLR L3 PacketNet Wave 2_Eastern Route (ATLA)		
Tuesday, November 16th										
8:00 AM	NASA/NOAA/ICAIR/NCDM-									
8:30 AM	LAC/DICE; see details in Scinet									
9:00 AM										
9:30 AM										
10:00 AM	NASA/NOAA/ICAIR/NCDM-									
10:30 AM	NASA/NOAA/ICAIR/NCDM-									
11:00 AM	NASA/NOAA/ICAIR/NCDM-									
11:30 AM	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	
12:00 PM	NASA/NOAA/ICAIR/NCDM-									
12:30 PM	NASA/NOAA/ICAIR/NCDM-									
1:00 PM	NASA/NOAA/ICAIR/NCDM-									
1:30 PM	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	
2:00 PM	NASA/NOAA/ICAIR/NCDM-									
2:30 PM	NASA/NOAA/ICAIR/NCDM-									
3:00 PM	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	
3:30 PM	NASA/NOAA/ICAIR/NCDM-									
4:00 PM	NASA/NOAA/ICAIR/NCDM-									
4:30 PM	NASA/NOAA/ICAIR/NCDM-									
5:00 PM	NASA/NOAA/ICAIR/NCDM-									
5:30 PM	NASA/NOAA/ICAIR/NCDM-									
6:00 PM	End of Show Hours									
6:00 PM - 8:00 PM	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NC	NASA/NOAA/ICAIR/NC	NASA/NOAA/ICAIR/NC	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NCDM-	NASA/NOAA/ICAIR/NC	NASA/NOAA/ICAIR/NC		
8:00 PM - 9:00 PM	NASA/NOAA/ICAIR/NCDM-									
9:00 PM - 10:00 PM	NASA/NOAA/ICAIR/NCDM-									

Example Showing NASA's Daily Scheduled Use of Internet2 During SC10

Source: Chris Robb/Internet2

All Times are Central	10% of 100G	20% of 100G	30% of 100G	40% of 100G	50% of 100G	60% of 100G	70% of 100G	80% of 100G	90% of 100G	100% of 100G
Tuesday, November 16th										
8:00 AM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
8:30 AM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
9:00 AM	LAC/DICE; see details in Scinet	LAC/DICE; see details in Scinet	LAC/DICE; see details in Scinet	LAC/DICE; see details in Scinet						
9:30 AM										
10:00 AM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-
10:30 AM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
11:00 AM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
11:30 AM										
12:00 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
12:30 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
1:00 PM	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see
1:30 PM										
2:00 PM	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see						
2:30 PM	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see	NASA/NOAA/ICA IR/NCDM- LAC/DICE; see						
3:00 PM										
3:30 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
4:00 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-
4:30 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
5:00 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
5:30 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
6:00 PM	End of Show Hours									
6:00 PM - 8:00 PM										
8:00 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-						
9:00 PM										
9:00 PM - 10:00 PM	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-	NASA/NOAA/ICA IR/NCDM-

Example Showing NASA's Daily Scheduled Demonstrations During SC10

All Times are Central	Type of Demonstration/Experiment	Comments
Tuesday, November 16th		
8:00 AM	NASA open testing period	
8:30 AM		
9:00 AM		
9:30 AM		
10:00 AM	1x100G Internet2 Use	
10:30 AM	Inter-Booth: 1x40G Extreme Networks+ColorChip, 1x100G Ciena, 1x100G Cisco	DICE using 4x10G NLR links with GSFC
11:00 AM	NASA open testing period	
11:30 AM	8x10G NLR Use	
12:00 PM	NASA open testing period	
12:30 PM	Inter-Booth: 1x40G Extreme Networks+ColorChip, 1x100G Ciena, 1x100G Cisco	DICE using 4x10G NLR links with GSFC
1:00 PM	1x100G Internet2 Use	
1:30 PM	8x10G NLR Use	
2:00 PM	NASA open testing period	
2:30 PM	NASA open testing period	
3:00 PM	8x10G NLR Use	
3:30 PM	NASA open testing period	
4:00 PM	1x100G Internet2 Use	
4:30 PM	Inter-Booth: 1x40G Extreme Networks+ColorChip, 1x100G Ciena, 1x100G Cisco	DICE using 4x10G NLR links with GSFC
5:00 PM	NASA open testing period	
5:30 PM	NASA open testing period	
6:00 PM	End of Show Hours	
6:00 PM -	8x10G NLR Use	
8:00 PM		
8:00 PM -	NASA open testing period	
9:00 PM		
9:00 PM -	1x100G Internet2 Use	
10:00 PM		

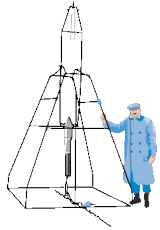


Introduction To GSFC High End Computing 20, 40 & 100 Gbps Network Testbeds

Summary of HEC Test WAN Requirements

- Key SLA Requirements of HEC's Test WAN Network
 - Availability (percent): 80.0
 - Restoral Time: <48 hours
 - Coverage Period: 24x7
 - RTT between sites in CONUS: <100ms
 - Packet Loss: <1E-7
 - Jumbo Frames: Transport of 9000-byte IP MTU jumbo frames without fragmentation
 - Bandwidth between NAS@ARC and NCCS@GSFC:
 - 1Jul11 through 30Jun13: 40Gbps (i.e., 1x40G, not 4x10G or other LAG approaches)
 - 1Jul13 through 30Jun15: 100Gbps (i.e., 1x100G, not 10x10G or other LAG approaches)
- The above SLA requirements apply only to the “TTT...TTT” links
- The SLA requirements of “OOO...OOO” links are similar to “TTT...TTT” links except that “OOO...OOO” links have their Availability (percent) parameters at 99.50 and their Restoral Time parameters at 4 hours; and therefore they likely need to be provisioned from a different supplier than “TTT...TTT” links





GSFC High End Computer Network (HECN) Team



From left: Pat Gary, Bill Fink, Paul Lang, Aruna Muppalla, Jeff Martz, Mike Stefanelli

Science and Engineering Network (SEN)

Typically enabling 1-10 gigabit per second (Gbps) user connections, the SEN is a non-mission-dedicated high-end computer network at GSFC serving GSFC projects/users who have computer network performance requirements greater than those baselined for GSFC's general-use campus-wide Center Network Environment (CNE)

<http://science.gsfc.nasa.gov/606.1/SENuser.html>

20, 40 & 100 Gbps Network Testbed Plans & Accomplishments to Date

http://science.gsfc.nasa.gov/606.1/docs/HECN_10G_Testbeds_082210.pdf





~~NETWORK BOTTLENECKS~~

