

# Flow Data at 10 GigE and Beyond

## What can (or should) we do ?

Scott Pinkerton

[pinkerton@anl.gov](mailto:pinkerton@anl.gov)

Argonne National Laboratory

[www.anl.gov](http://www.anl.gov)



## About me ....

- Involved in network design, network operation & network security for the last 10 years
- Flow data practitioner
- Campus perspective
- Our flow data uses typically include:
  - Real-time anomaly detection
  - Forensic analysis



# User facilities



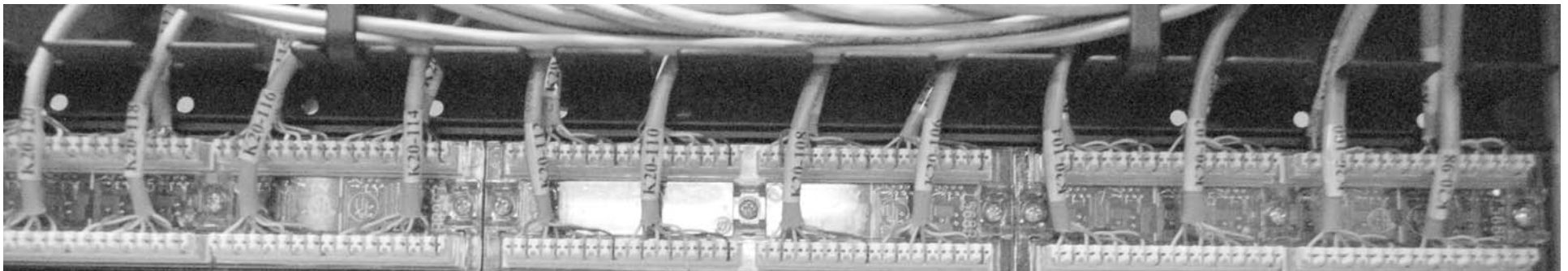
# Using Flow Data in a campus environment

- In ~2000 started collecting Netflow data from all of the core campus network devices using the OSU Flowtools package
- By 2004, we were collecting Netflow data down in the distribution and access layers of the campus network
- Today, still consider flow data to be a critical part of our anomaly detection systems. Goals are to:
  - Protect the Laboratory computers from the Internet
  - Protect the Internet from the Laboratory computers
  - Have visibility into “lateral movement” of compromised hosts
- Campus environments can be large ....



# Texas A&M Campus Network

- Wired Network
  - 10 Gbps backbone
  - 50,000 computers
  - 90,000 wired ports
- Gateway to regional and national networks
- Wireless Network
  - 11 million square ft. of wireless access
  - 340+ buildings with wireless access across 5200 acres



---

# U of M Twin Cities Campus Network

- 23 Cisco 6509s
- 4,323 Cisco 3750s
- 1,133 Switch Stacks
- 74,414 Switchports
- Redundant 10-Gigabit Backbone
- Topology: 18 layer-2 switched domains interconnected by a layer-3 MPLS-VPN backbone




# A “big science” Perspective - driving speeds & feeds

- Data networks continue to evolve in support of the scientific mission
- Key drivers include:
  - Large Hadron Collider (LHC), CERN
    - CERN to US Tier1 data rates: 10 Gbps by 2007, 30-40 Gbps by 2010/11
  - Leadership Computing Facilities (LCF), ANL and ORNL
  - Relativistic Heavy Ion Collider (RHIC), BNL
  - Large-scale Fusion (ITER), France
  - Climate Science
    - Significant data set growth is likely in the next 5 years, with corresponding increase in network bandwidth requirement for data movement (current data volume is ~200TB, 1.5PB/year expected rate by 2010)




# Science Network Requirements Aggregation Summary

Science Drivers	End2End Reliability	Connectivity	2006 End2End Band width	2010 End2End Band width	Traffic Characteristics	Network Services
Science Areas / Facilities						
<b>Advanced Light Source</b>	-	<ul style="list-style-type: none"> <li>• DOE sites</li> <li>• US Universities</li> <li>• Industry</li> </ul>	1 TB/day 300 Mbps	5 TB/day 1.5 Gbps	<ul style="list-style-type: none"> <li>• Bulk data</li> <li>• Remote control</li> </ul>	<ul style="list-style-type: none"> <li>• Guaranteed bandwidth</li> <li>• PKI / Grid</li> </ul>
<b>Bioinformatics</b>	-	<ul style="list-style-type: none"> <li>• DOE sites</li> <li>• US Universities</li> </ul>	625 Mbps 12.5 Gbps in two years	250 Gbps	<ul style="list-style-type: none"> <li>• Bulk data</li> <li>• Remote control</li> <li>• Point-to-multipoint</li> </ul>	<ul style="list-style-type: none"> <li>• Guaranteed bandwidth</li> <li>• High-speed multicast</li> </ul>
<b>Chemistry / Combustion</b>	-	<ul style="list-style-type: none"> <li>• DOE sites</li> <li>• US Universities</li> <li>• Industry</li> </ul>	-	10s of Gigabits per second	<ul style="list-style-type: none"> <li>• Bulk data</li> </ul>	<ul style="list-style-type: none"> <li>• Guaranteed bandwidth</li> <li>• PKI / Grid</li> </ul>
<b>Climate Science</b>	-	<ul style="list-style-type: none"> <li>• DOE sites</li> <li>• US Universities</li> <li>• International</li> </ul>	-	5 PB per year 5 Gbps	<ul style="list-style-type: none"> <li>• Bulk data</li> <li>• Remote control</li> </ul>	<ul style="list-style-type: none"> <li>• Guaranteed bandwidth</li> <li>• PKI / Grid</li> </ul>
<b>High Energy Physics (LHC)</b> 	99.95+% (Less than 4 hrs/year)	<ul style="list-style-type: none"> <li>• US Tier1 (DOE)</li> <li>• US Tier2 (Universities)</li> <li>• International (Europe, Canada)</li> </ul>	10 Gbps	60 to 80 Gbps (30-40 Gbps per US Tier1)	<ul style="list-style-type: none"> <li>• Bulk data</li> <li>• Remote control</li> </ul>	<ul style="list-style-type: none"> <li>• Guaranteed bandwidth</li> <li>• Traffic isolation</li> <li>• PKI / Grid</li> </ul>



# Science Network Requirements Aggregation Summary

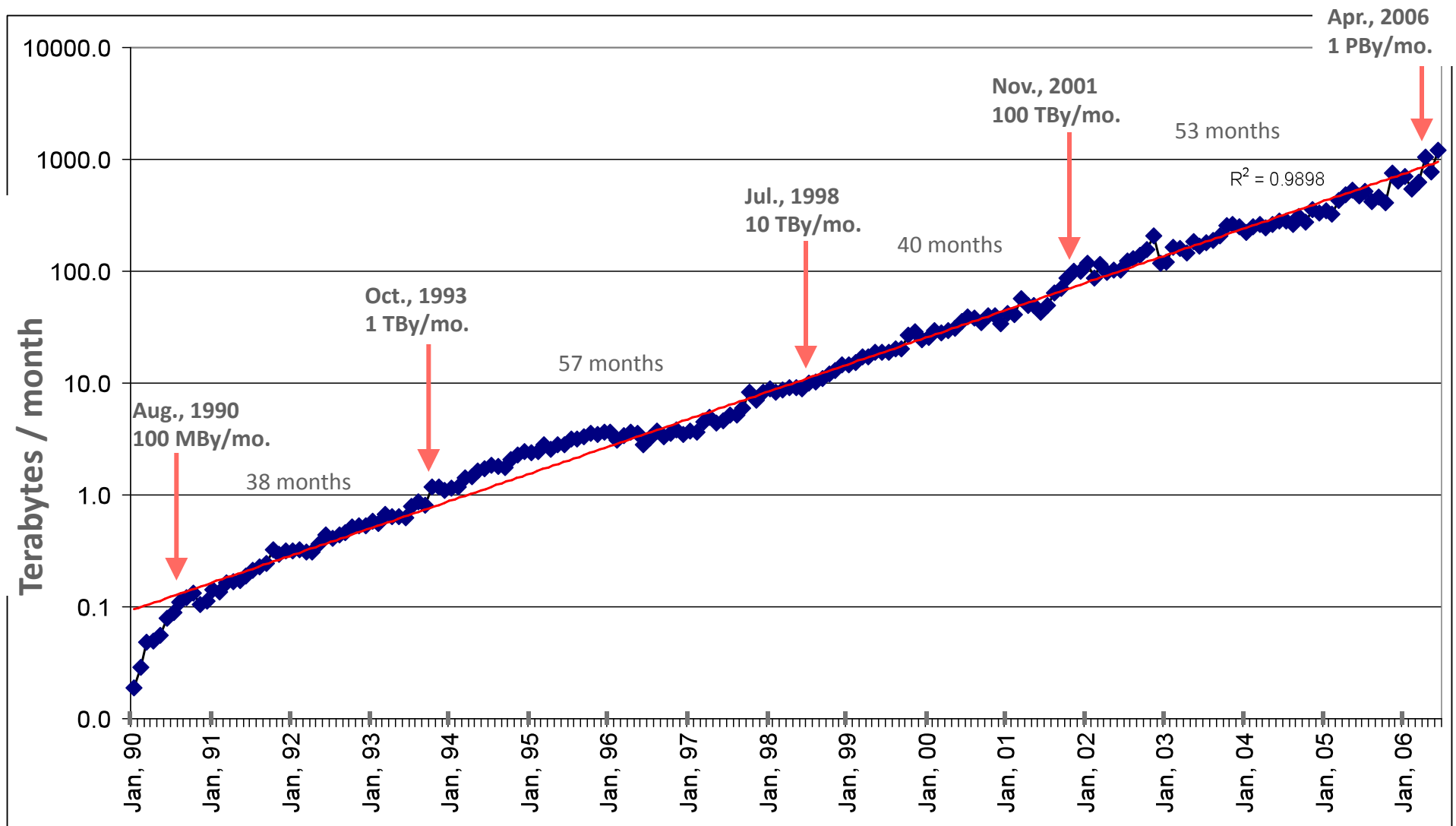
Science Drivers	End2End Reliability	Connectivity	2006 End2End Band width	2010 End2End Band width	Traffic Characteristics	Network Services
Science Areas / Facilities						
<b>Magnetic Fusion Energy</b>	99.999% (Impossible without full redundancy)	<ul style="list-style-type: none"> <li>• DOE sites</li> <li>• US Universities</li> <li>• Industry</li> </ul>	200+ Mbps	1 Gbps	<ul style="list-style-type: none"> <li>• Bulk data</li> <li>• Remote control</li> </ul>	<ul style="list-style-type: none"> <li>• Guaranteed bandwidth</li> <li>• Guaranteed QoS</li> <li>• Deadline scheduling</li> </ul>
<b>NERSC</b>	-	<ul style="list-style-type: none"> <li>• DOE sites</li> <li>• US Universities</li> <li>• Industry</li> <li>• International</li> </ul>	10 Gbps	20 to 40 Gbps	<ul style="list-style-type: none"> <li>• Bulk data</li> <li>• Remote control</li> </ul>	<ul style="list-style-type: none"> <li>• Guaranteed bandwidth</li> <li>• Guaranteed QoS</li> <li>• Deadline Scheduling</li> <li>• PKI / Grid</li> </ul>
<b>NLCF</b>	-	<ul style="list-style-type: none"> <li>• DOE sites</li> <li>• US Universities</li> <li>• Industry</li> <li>• International</li> </ul>	Backbone Band width parity	Backbone band width parity	<ul style="list-style-type: none"> <li>• Bulk data</li> </ul>	
<b>Nuclear Physics (RHIC)</b>	-	<ul style="list-style-type: none"> <li>• DOE sites</li> <li>• US Universities</li> <li>• International</li> </ul>	12 Gbps	70 Gbps	<ul style="list-style-type: none"> <li>• Bulk data</li> </ul>	<ul style="list-style-type: none"> <li>• Guaranteed bandwidth</li> <li>• PKI / Grid</li> </ul>
<b>Spallation Neutron Source</b> 	High (24x7 operation)	<ul style="list-style-type: none"> <li>• DOE sites</li> </ul>	640 Mbps	2 Gbps	<ul style="list-style-type: none"> <li>• Bulk data</li> </ul>	

**Table 1: ALCF requirements summary**

Feature	Key Science Drivers		Anticipated Network Requirements	
	Science Instruments and Facilities	Process of Science	Local Area Network Bandwidth and Services	Wide Area Network Bandwidth and Services
Near-term (0-2 years)	<ul style="list-style-type: none"> <li>ALCF production resources (intrepid)</li> </ul>	<ul style="list-style-type: none"> <li>Large file transfers. Other labs and computing centers are common targets, but it can be any institution based on INCITE users needs.</li> <li>Some real-time video, computational steering, real-time control apps possible.</li> </ul>	<ul style="list-style-type: none"> <li>Node to node is handled by proprietary vendor interconnect. 425 MB/s per link.</li> <li>Node to storage is approx. 1,000 ports of 10 gigabit.</li> <li>Other local needs are primarily admin-related and are trivial.</li> </ul>	<ul style="list-style-type: none"> <li>10s of TB/day</li> <li>10-30 Gbps</li> </ul>
2-5 years	<ul style="list-style-type: none"> <li>Next major machine upgrade</li> </ul>	<ul style="list-style-type: none"> <li>Large file transfers. Other labs and computing centers are common targets, but it can be any institution based on INCITE users needs.</li> <li>Real-time video, computational steering, real-time control apps more common, but still relatively small in comparison to file transfers.</li> </ul>	<ul style="list-style-type: none"> <li>Node to node is handled by proprietary vendor interconnect. 1-5 GB/s per link.</li> <li>Node to storage is likely InfiniBand-based and on the order of 3K-5K ports</li> <li>Other local needs are primarily admin-related and are trivial.</li> </ul>	<ul style="list-style-type: none"> <li>100s of TB/day</li> <li>100-300 Gbps</li> </ul>
5+ years	<ul style="list-style-type: none"> <li>Push towards exascale computing</li> </ul>	<ul style="list-style-type: none"> <li>Massive data sets are common. File transfers still dominate, but WAN file systems, distributed databases use grows.</li> <li>Machines are sufficiently powerful that computational steering, real time simulations are used regularly</li> <li>Use of collaboration tools continues to grow</li> </ul>	<ul style="list-style-type: none"> <li>Node to node is probably still handled by proprietary vendor interconnect, but could be standards based, such as InfiniBand.</li> <li>Node to storage is likely InfiniBand or other standards-based interconnect.</li> <li>Other local needs are primarily admin-related and are trivial.</li> </ul>	<ul style="list-style-type: none"> <li>Petabytes per day</li> <li>Terabit networks</li> </ul>



# ESnet Traffic has Increased by 10X Every 47 Months, on Average, Since 1990



Log Plot of ESnet Monthly Accepted Traffic, January, 1990 – June, 2006

# Key Take Aways

- Building networks for the future – takes a lot of planning
- Or, maybe more importantly it takes a lot of predicting (future requirements)
- Without the planning (and the predicting) how can the vendors gear up to provide the necessary capabilities ?
- Are we doing a good job communicating future requirements for flow data ?

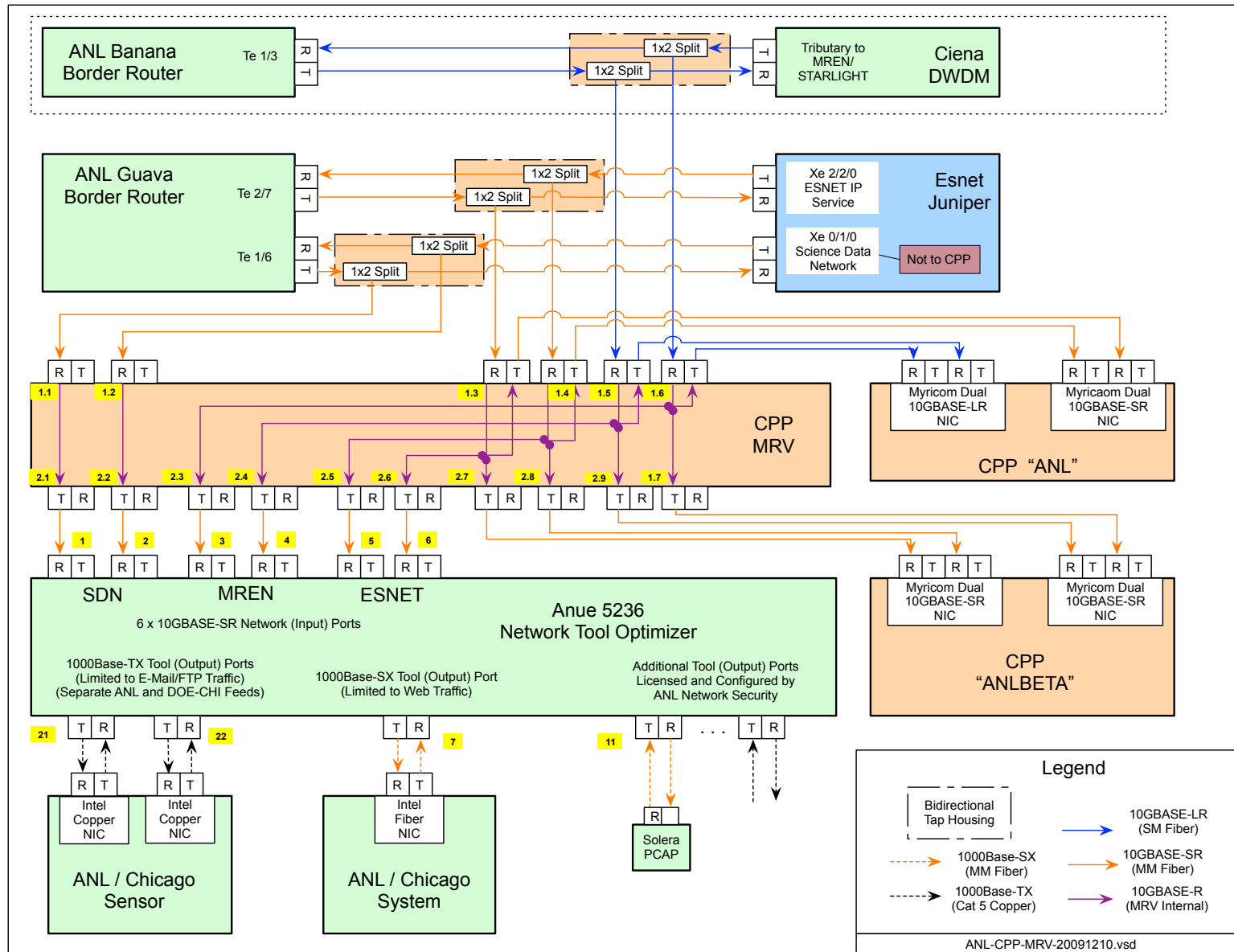


# Future of non-sampled Flow data seems bleak (IMHO)

- Speeds and feeds increasing to keep pace with scientific demand
- Many/most vendors are struggling to provide non-sampled flow data directly from the switches or routers just @ 10 Gbps (much less at 40 or 100 Gbps)
- Can optical taps really scale up to provide the needed number of monitor points ?
  - For me, I think the answer is no



# Leveraging taps to create monitor points





# What Can We Do - Process Perspective ?

- Identify our needs/requirements
- Write it down
- Communicate it to the vendors



# What are our needs/requirements/drivers ?

- Strong support for “existence” analysis
- Scalable
  - From a campus perspective (monitoring at the border & internally)
  - From a “big science” perspective (speeds/feeds & large file txfer)
  - Equals – built in to the switches & routers (IMHO)
- Non-sampled data
  - Sampled data has its place (traffic engineering & other)
  - Painful to perform forensic analysis with sampled data





# What if ideas ?

- Leveraging “cores” internal to the switches & routers for custom applications
  - Bloom filter ?
  - What info/data types would we want available to an internal app ?
- Adapting to the Very Very Large data txfers
  - Do we need a new scale for active timeout ? Was 30 seconds, now 30 minutes or 3 hours ?
- Notifications at start of a new flow – first packet ?





## As a community - what can we do ?

- Should we try and develop future requirements ?
- Do we have enough energy/motivation to do it ?
- Can we agree on requirements ?
- Can we influence the networking equipment purchase decisions ?
- Your thoughts ??

