# Security Challenges for High Throughput Data Transfers

Dr Tim Chown, Network Development Manager, Jisc

WISE Meeting – The Cosener's House – 28 Feb 2018

» There is growing interest in moving research data around the network
  › Data capture or generation to computing facility, and perhaps back
  › Remote visualisation
  › Data replication / distributed storage / backups
  › To / from cloud

» Data set volumes are increasing
  › 100 TB is no longer 'large'
  › Moving 100 TB takes 10Gbps of throughput for 24 hours

» How do we do this securely, AND with the necessary performance?
  › This deck has some thoughts on this topic...

# What are we not talking about?

» Security embraces many perspectives and methods

» WISE has no doubt covered many topics at this meeting

» We're not talking here about the more classic security challenges
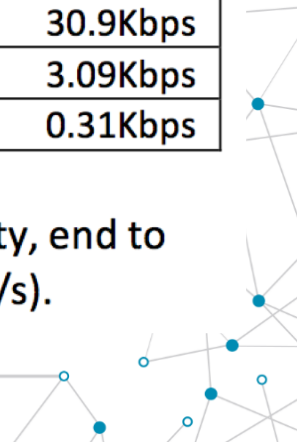  › Authentication, authorisation, certificates, etc.
  › Or how to do a GDPR audit!

» The question to consider here is how do we achieve high throughput on the end-to-end data transfer path, while applying appropriate security measures to the traffic in flight
  › The path needs to be performant, else the research suffers
  › Can argue that necessary performance is a **requirement**

The following table, taken from a publication by ESnet[3], shows the *theoretical* throughput required to transfer a given size of data set in a range of example time periods.

|        | 1 Min | 5 Mins | 20 Mins | 1 Hour | 8 Hours | 1 Day | 7 Day | 30 Days |
|--------|-------|--------|---------|--------|---------|-------|-------|---------|
| 10 PB  | 1,333Tbps | 266.7Tbps | 66.7Tbps | 22.2Tbps | 2.78Tbps | 926Gbps | 132Gbps | 30.9Gbps |
| 1 PB   | 133.3Tbps | 26.7Tbps | 6.67Tbps | 2.2Tbps | 278Gbps | 92.6Gbps | 13.2Gbps | 3.09Gbps |
| 100 TB | 13.3Tbps | 2.67Tbps | 667Gbps | 222Gbps | 27.8Gbps | 9.26Gbps | 1.32Gbps | 309Mbps |
| 10 TB  | 1.33Tbps | 266.7Gbps | 66.7Gbps | 22.2Gbps | 2.78Gbps | 926Mbps | 132Mbps | 30.9Mbps |
| 1 TB   | 133.3Gbps | 26.67Gbps | 6.67Gbps | 2.22Gbps | 278Mbps | 92.6Mbps | 13.2Mbps | 3.09Mbps |
| 100 GB | 13.3Gbps | 2.67Gbps | 667Mbps | 222Mbps | 27.8Mbps | 9.26Mbps | 1.32Mbps | 309Kbps |
| 10 GB  | 1.33Gbps | 266.7Mbps | 66.7Mbps | 22.2Mbps | 2.78Mbps | 926Kbps | 132Kbps | 30.9Kbps |
| 1 GB   | 133.3Mbps | 26.7Mbps | 6.67Mbps | 2.22Mbps | 278Kbps | 92.6Kbps | 13.2Kbps | 3.09Kbps |
| 100 MB | 13.3Mbps | 2.67Mbps | 667Kbps | 222Kbps | 27.8Kbps | 9.26Kbps | 1.32Kbps | 0.31Kbps |

Thus, in principle, if you need to move 100GB in 20 minutes, you will need at least a 1Gbit/s capacity, end to end. Or, if you have a 10Gbit/s link, you can in principle move 100TB in a day (at a rate of 9.26Gbit/s).

# Understanding the factors affecting E2E

» Achieving optimal end-to-end performance is a multi-faceted problem.

» It includes:

> Appropriate network capacity provisioning between the end sites

> Properties of the local campus network (at each end), including capacity of the external connectivity, internal LAN design, the performance of firewall / IDS devices, and the configuration of other devices on the path

> End system configuration and tuning; network stack buffer sizes, disk I/O, …

> The choice of tools used to transfer data, e.g. scp, Globus, rsync, Aspera, …

» To optimise end-to-end performance, you need to address each aspect

» There will inevitably be a bottleneck somewhere

# Campus network engineering

» Question: how to design the local campus network for optimal end-to-end inter-site data transfer performance?

» Problem: An application using TCP will see its performance degrade if packets are lost, with more degradation the higher the path's RTT

› Very small loss can have a surprisingly significant impact

› Therefore we need to engineer towards zero packet loss

» Zero loss implies both sufficient capacity and performant network elements

» The challenge is that many campus security appliances, esp. corporate firewall/IDS, are designed for 1000's of small flows, not tens of very large flows, and they can thus drop packets

» Answer? The Science DMZ

# TCP with a small amount of packet loss...



Throughput vs. Increasing Latency with .0046% Packet Loss

With loss, high performance beyond metro distances is essentially impossible

Local (LAN)
Metro Area
Regional
Continental
International

Measured (TCP Reno)  Measured (HTCP)  Theoretical (TCP Reno)  Measured (no loss)

# Site network engineering – the Science DMZ

» ESnet published the Science DMZ 'design pattern' in 2012/13
  › https://www.es.net/assets/pubs_presos/sc13sciDMZ-final.pdf
» Three key elements:
  › Network architecture improvements; avoiding local bottlenecks
  › Network performance measurement
  › Data transfer node (DTN) design and configuration
» Also termed a "high speed on-ramp" to the campus storage
  › Splits the internal and external latency domains

» The NSF Cyberinfrastructure (CC*) Program funded this model in over 100 US universities:
  › See https://www.nsf.gov/funding/pgm_summ.jsp?pims_id=504748

# Science DMZ Design Pattern (Abstract)



**Border Router**

perfS●NAR

**Enterprise Border Router/Firewall**

**WAN**

10G

10GE

*Clean, High-bandwidth WAN path*

*Site / Campus access to Science DMZ resources*

10GE

perfS●NAR

10GE

**Site / Campus LAN**

10GE

**Science DMZ Switch/Router**

*Per-service security policy control points*

perfS●NAR

10GE

**High performance Data Transfer Node with high-speed storage**

# Local and wide area data flows



**Border Router**

**Enterprise Border Router/Firewall**

perfSONAR

perfSONAR

**WAN**

10G

10GE

*Clean, High-bandwidth WAN path*

*Site / Campus access to Science DMZ resources*

10GE

perfSONAR

**Site / Campus LAN**

10GE

**Science DMZ Switch/Router**

perfSONAR

*Per-service security policy control points*

10GE

**High performance Data Transfer Node with high-speed storage**

*High Latency WAN Path*

*Low Latency LAN Path*

# Examples of Science DMZ in use at UK sites

» There are several examples of sites in the UK that have some form of Science DMZ deployment

» In many cases the deployments were made without knowledge of the Science DMZ model!

» Science DMZ is just a set of good principles to follow, so it's not surprising that some Janet sites were already doing it, especially the GridPP sites

» Examples in the UK:
  › Diamond Light Source
  › JASMIN/CEDA Data Transfer Zone
  › Imperial College GridPP; supports up to 35Gbit/s of IPv4/IPv6
  › BUT, to realise the end-to-end benefit, both ends need to apply the principles

# Science DMZ as a Security Architecture

» Rationale?

» It allows for better segmentation of risks, and more granular application of controls to those segmented risks.

› Limit risk profile for high-performance data transfer applications

› Apply specific controls to data transfer nodes (DTNs)

› Avoid including unnecessary risks, unnecessary controls

» Remove degrees of freedom – focus only on what is necessary

› Easier to secure

› Easier to achieve performance

› Easier to troubleshoot

» Performance is a key requirement; e.g., use efficient ACLs

› See https://www.slideshare.net/JISC/science-dmz-security (Kate Mace)

# Other examples of campus network engineering

» Many Janet sites split their external connectivity (their choice…)
  › e.g., 40G total; 1x10G campus, 1x10G research science data, 2x10G resilience
  › And then apply Science DMZ principles to the dedicated research data path
  › Or employ Science DMZ in their data centre

» The Worldwide LHC Computing Grid  (WLCG) has used physical / virtual overlays
  › LHCOPN (private optical network) / LHCONE (virtual network)
  › LHCONE implicitly becomes a 'trusted' network
» But how should campuses cater for multiple data-intensive science disciplines?
  › Would one new overlay network per research community scale?

» Some sites are exploring SDN, to direct traffic dynamically and efficiently on campus
  › The classic on-ramp 'Science DMZ' is a rather static architecture

# Aside 1: One person's (D)DoS...

» ... is another person's research data transfer!

» What might you see in/out of a campus?
  › High volume UDP data transfer flows, e.g., Aspera
  › New protocols, e.g., QUIC (UDP/HTTP2)
  › 'Smarter' TCP algorithms, e.g., TCP-BBR
  › Highly parallelised flows, e.g., Globus / Grid FTP

» Applications behaving this way might seem to be out of profile, and thus potential (D)DoS
  › Need to keep abreast of application protocol developments
  › May white-list certain applications / address space

» If you need to encrypt sent data, this might be implemented by
› Pre-encrypting the data
› Encryption on the fly as you transmit
› Encryption between gateways on the path

» All these have potential performance issues or limitations

» Some NRENs offer solutions in this space
› e.g., Jisc has Safe Share
› https://www.jisc.ac.uk/safe-share
› But by default only up to 1Gbit/s; higher throughput costs £££
› Genomics project data sets can easily be 100-200TB.

# Measuring network characteristics

» Important to have telemetry on your network

» The Science DMZ model recommends perfSONAR for this
» Collects telemetry over time
  › Throughput, loss, latency, path
  › Allows retrospective viewing of data
  › Uses proven tools under the hood such as iperf
» Can run tests between two perfSONAR systems, or build a mesh

» Helps you assess the impact of changes to your network or systems
» It can highlight poor performance, but doesn't troubleshoot per se
» May indicate impact of security appliances on performance

# Janet - London perfSONAR node

Added a Jisc certificate

Dual-stack

Possible to set up tests manually, but better to set up a mesh…
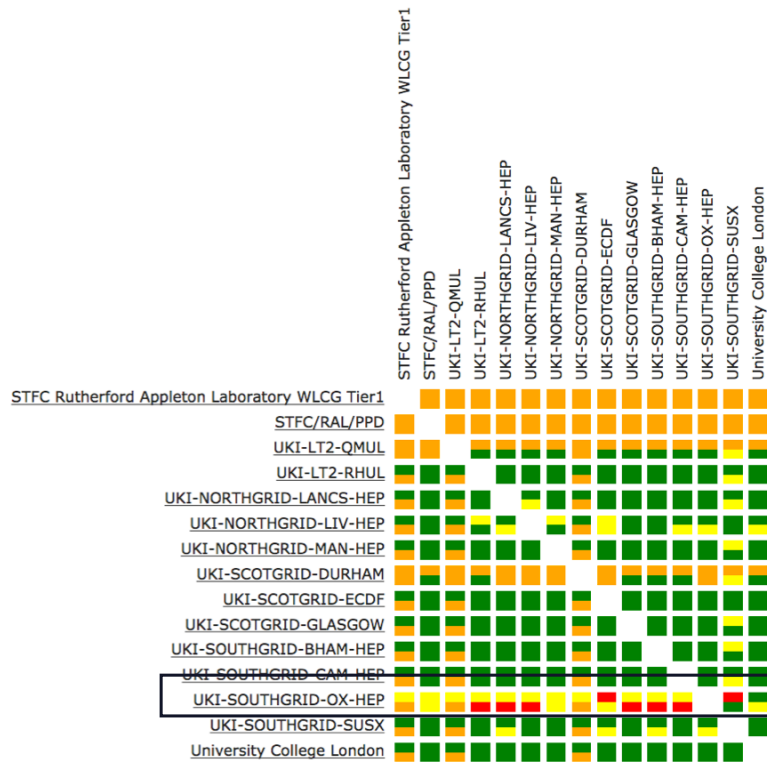
# Example perfSONAR mesh – UK GridPP

» GridPP = UK Particle Physics computing collaboration
› UK component of the WLCG (Worldwide LHC Computing Grid)

» Nineteen sites forming one Tier-1 and four distributed Tier-2 sites

» Most sites have perfSONAR nodes next to their storage servers

» They are running a dual-stack mesh
› Measure IPv4 and IPv6 performance independently

» Provides an insight into network performance across the sites

» Live version:
› http://ps-dash.dev.ja.net/maddash-webui/index.cgi?dashboard=UK%20Mesh%20Config

# GridPP mesh: Traceroute

# GridPP mesh: Latency and Loss



**UK Mesh Config - IPv4 Latency Tests**
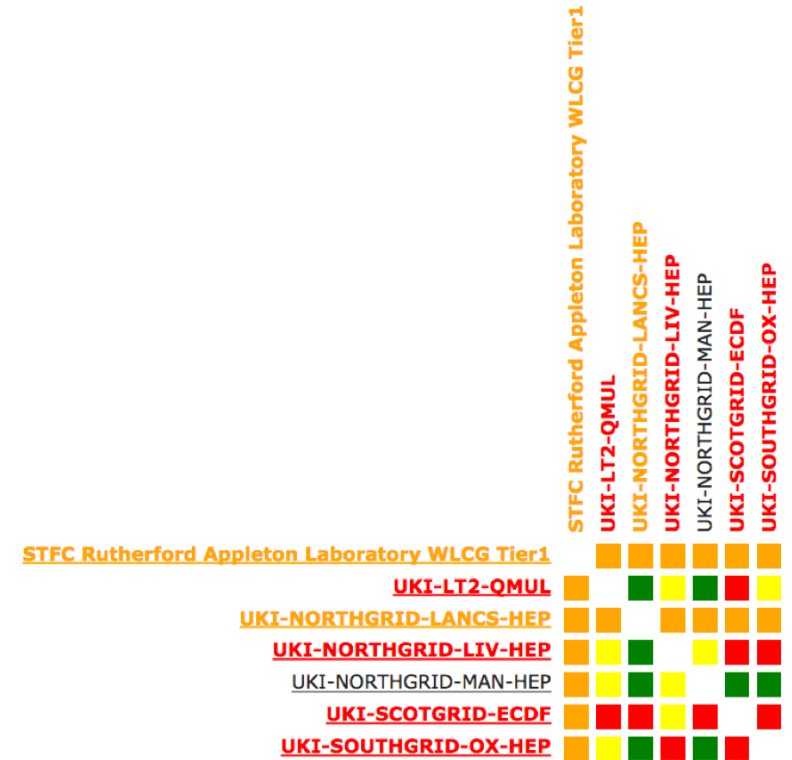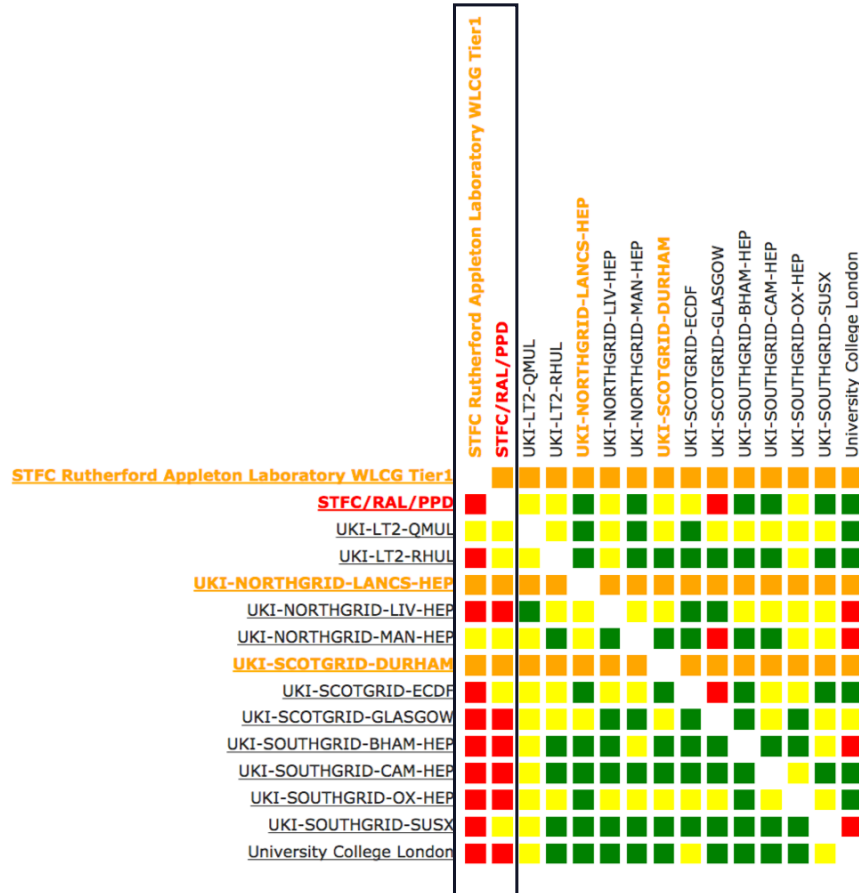
Loss rate is <= 0 | Loss rate is >= 0 | Loss rate is >= 0.01 | Unable to retrieve data | Check has

⚠ Found a total of 4 problems involving 4 hosts in the grid

**UK Mesh Config - IPv6 Latency Tests**

Loss rate is <= 0 | Loss rate is >= 0 | Loss rate is >= 0.01 | Unable to ret
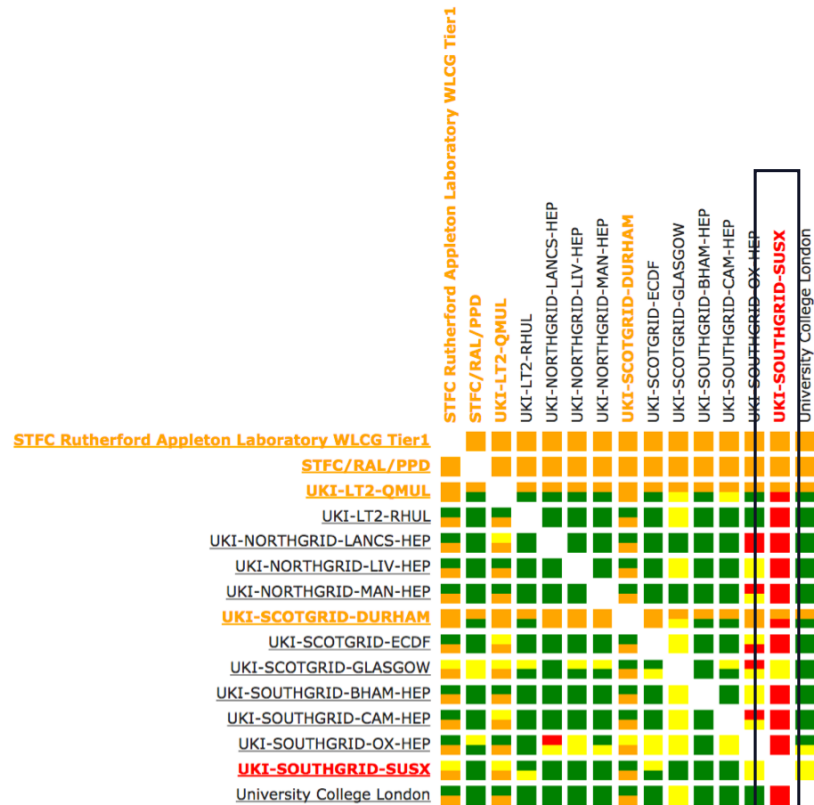
⚠ Found a total of 8 problems involving 6 hosts in the grid

# GridPP mesh: Throughput

**UK Mesh Config - IPv4 Bandwidth Tests**

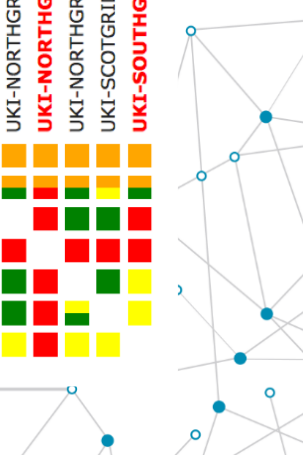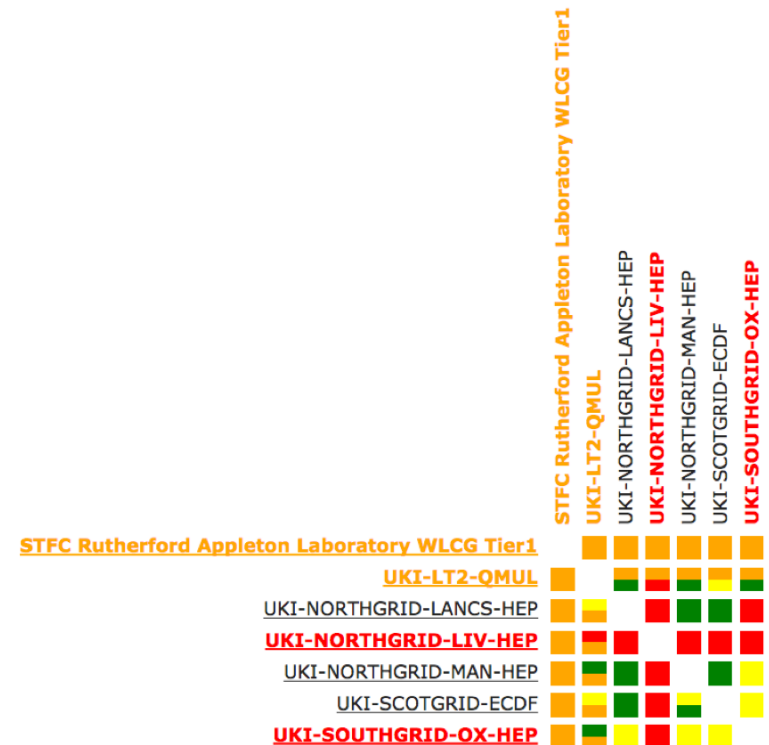Throughput >= 900Mbps | Throughput < 900Mbps | Throughput <= 500Mbps | Unable to retrieve da

⚠ Found a total of 5 problems involving 5 hosts in the grid

**UK Mesh Config - IPv6 Bandwidth Tests**

Throughput >= 900Mbps | Throughput < 900Mbps | Throughput <= 500Mbps

⚠ Found a total of 5 problems involving 4 hosts in the grid

# perfSONAR – performance visualization over time
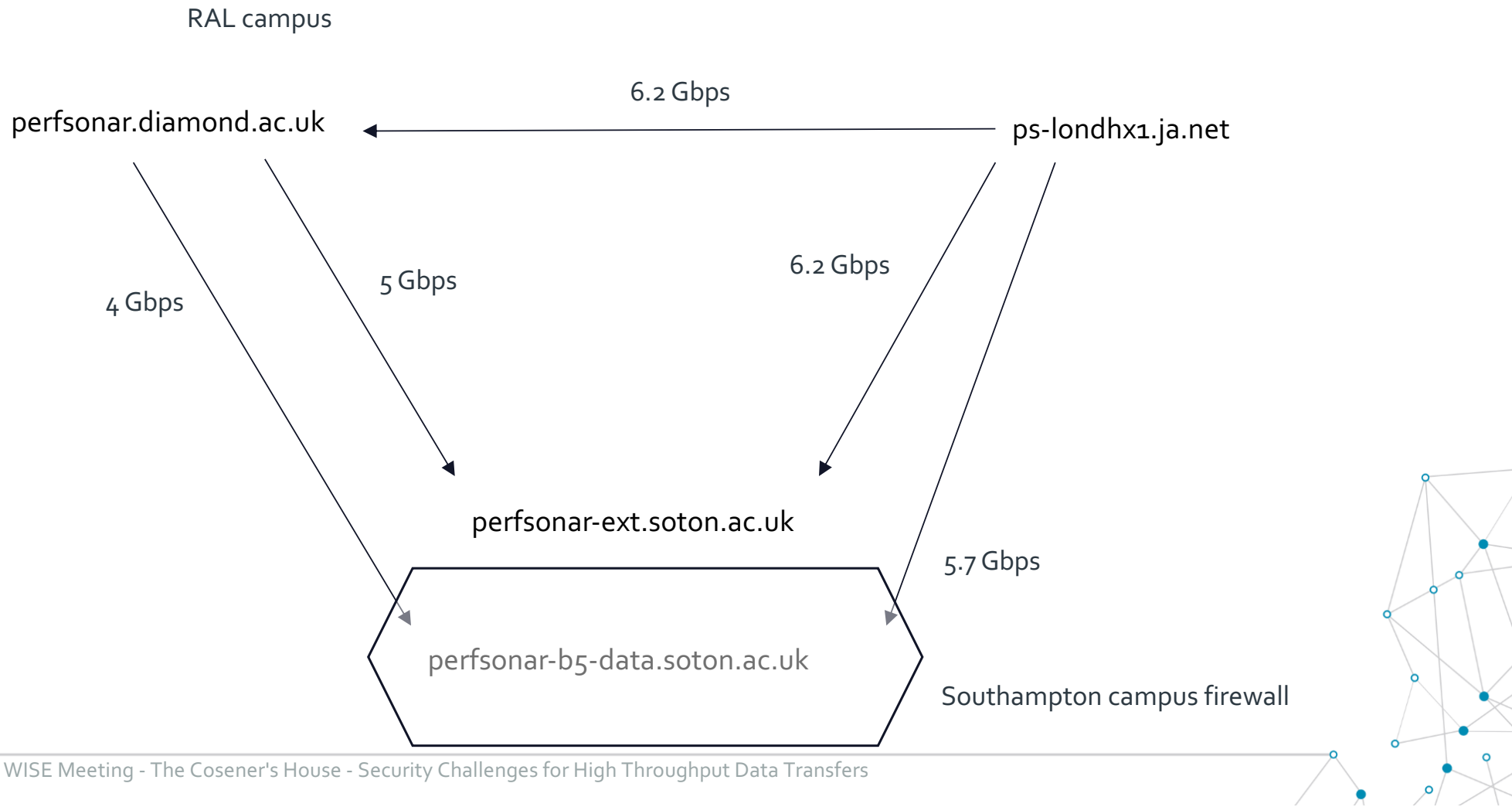
# Example: Diamond LS <-> Southampton

» Southampton researchers go to DLS to run experiments and collect data which they wish to transport back home

  › Rather than carrying disks, they ought to be able to use the Janet network

» Initial work at Southampton involved adopting Globus Connect to transfer files from DLS

» Achieved a significant improvement up to the then available 1 Gbps

» Also installed a perfSONAR host (*perfsonar-b5-data.soton.ac.uk*) on campus next to the data storage

» Network to storage then upgraded to 10 Gbps

» Later a perfSONAR host (*perfsonar-ext .soton.ac.uk*) was installed at the Southampton border, outside the firewall
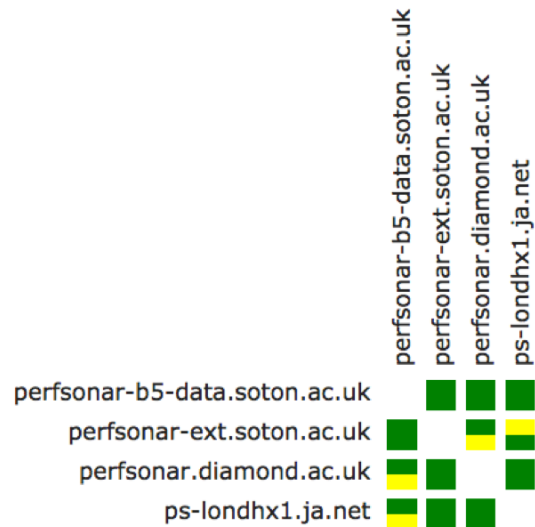
» See http://ps-dash.dev.ja.net/maddash-webui/index.cgi?dashboard=SES

# Current observed average throughput

RAL campus

perfsonar.diamond.ac.uk ← 6.2 Gbps ← ps-londhx1.ja.net

4 Gbps

5 Gbps

6.2 Gbps

perfsonar-ext.soton.ac.uk

5.7 Gbps

perfsonar-b5-data.soton.ac.uk

Southampton campus firewall

# Last 12 month view – London -> B5 network

» University wanting to backup data to RAL
  › Obtaining 300-400Mbit/s
  › Introduced Science DMZ principles
  › Now sustaining 3-4 Gbit/s

» University on WLCG; IDS imposed
  › Was obtaining up to 18Gbit/s
  › IDS throughput maximum 8 Gbit/s
  › An example of appropriate application of policy

» University with 'bug' on firewall
  › Capacity reduced to 1Gbit/s or less on any flow
  › Normal campus users did not report the issue; perfSONAR detected it

» Consider how to apply the necessary policy efficiently

» Is this an area that interests the WISE community?

» Design in appropriate network engineering

» The classic 'Science DMZ' model has value; many were doing it anyway
  › Well-tuned DTNs with host-based security
  › SDN may provide a more agile 'DMZ' architecture

» Consider emerging data-intensive application network protocols; QUIC, etc.

» Measure performance over time (perfSONAR)

» Don't disrupt research; performance is a requirement
  › Just like the confidentiality, integrity and availability (CIA) principles

# Some useful links

» Janet E2EPI project page
  › https://www.jisc.ac.uk/rd/projects/janet-end-to-end-performance-initiative
» JiscMail E2EPI list (approx 100 subscribers)
  › https://www.jiscmail.ac.uk/cgi-bin/webadmin?A0=E2EPI
» Camus Network Engineering for Data-Intensive Science workshop slides
  › https://www.jisc.ac.uk/events/campus-network-engineering-for-data-intensive-science-workshop-19-oct-2016
  › https://www.slideshare.net/JISC/science-dmz-security (Kate Mace, ESnet)
» Fasterdata knowledge base
  › http://fasterdata.es.net/
» eduPERT knowledge base
  › http://kb.pert.geant.net/PERTKB/WebHome

**Jisc**

# Please feel free to get in touch!

**Dr Tim Chown**

Network Development Manager

**tim.chown@jisc.ac.uk**

**jisc.ac.uk**