# Janet End-to-End Performance Initiative

Dr Tim Chown, Network Development Manager, Jisc

HEPSYSMAN – RAL – 14 June 2017

» Overall objective: help sites make optimal use of their Janet connectivity

» Seeing a growth in data-intensive science applications
  › Includes established areas, most notably GridPP (many of you here ☺)
  › As well as new areas like cryo-electron microscopy

» Seeing an increasing number of remote computation scenarios
  › e.g., scientific networked equipment, no local compute - might require 10Gbit/s to return computation results on a 100GB data set to a researcher for timely visualisation

» Starting to see more 100Gbit/s connectivity requests
  › Likely to have challenging data transfer requirements behind them, e.g. SKA

# Janet End-to-End Performance Initiative

» The goals include:

› Engaging with existing and data-intensive research communities and identifying emerging communities

› Creating dialogue between Jisc, computing service groups, and research communities

› Holding workshops, facilitating discussion on e-mail lists, etc.

› Helping researchers manage expectations

› Establishing and sharing best practices in identifying and rectifying causes of poor performance

› Promoting good practices in campus network engineering, esp. 'Science DMZ'

› Promoting deployment of performance measurement tools, esp. perfSONAR

» More information:

› https://www.jisc.ac.uk/rd/projects/janet-end-to-end-performance-initiative

# Understanding the factors affecting E2E

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

» It includes:

› Appropriate Janet provisioning between the end sites

› Properties of the local campus network (at each end), including capacity of the Janet connectivity, internal LAN design, the performance of firewalls and the configuration of other devices on the path

› End system configuration and tuning; network stack buffer sizes, disk I/O, memory management, etc.

› The choice of tools used to transfer data, and the underlying network protocols, e.g. xrootd, Globus, rsync, Aspera, or...?

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

» There will inevitably be a bottleneck somewhere

› Our initial focus has been network elements, but includes the end systems

# Getting a handle on requirements

» Is there good dialogue at your site between data-intensive researchers and the IT staff?
 › If so, is this happening on a regular basis?
 › Ideally want to be able to plan ahead, rather than adapt on the fly

» Do you conduct networking 'future looks'?
 › The LHC community is quite good at this; experiments are well-defined in advance
 › Any step changes – certainly at 10Gbps - might dwarf your site's organic growth
 › The capacity issue ought to have attention at CIO or PVC Research level

» Do you know what your application elephants are?
 › What's the breakdown of your site's network traffic?
 › How are you monitoring network flows?

# Researcher expectations

» How can we help set and manage researcher expectations?

» One aspect is helping them articulating their network requirements
  › At a simple level, volume of data in time X => data rate required
» It's also about understanding practical limitations
  › Can determine theoretical network throughput
  › e.g., in principle, you can transfer 100TB over a 10Gbit/s link in 1 day
  › But in practice, many factors may prevent this, as mentioned earlier
» We're encouraging researchers to speak to their computing service
  › Not to try, give up, and start sending their data on hard disks!

» Initial guidance document (feedback welcomed!):
  › https://community.jisc.ac.uk/system/files/1862/Network-Expectations-v1-0.pdf

# Example from the doc - theoretical throughput

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).
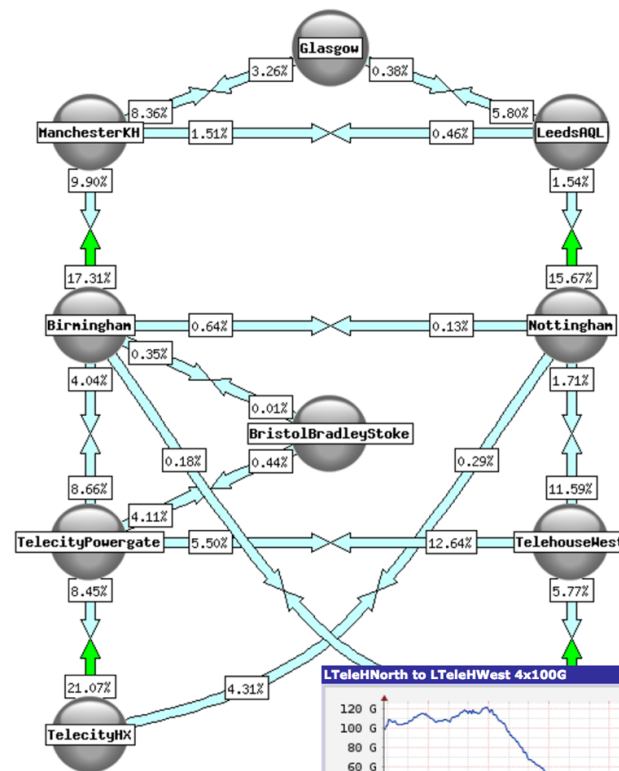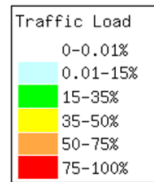
# Janet network engineering – capacity planning

» From Jisc's perspective, it's important to ensure there is sufficient capacity in the network for its connected member sites

» We perform an ongoing review of network utilisation
  › Provision the backbone network and regional links
  › Provision external connectivity, to other NRENs and networks
  › Observe the utilisation, model growth, predict ahead
  › Step changes have bigger impact at the local rather than backbone scale

» Janet has no differential queueing for regular IP traffic
  › The Janet Netpath service exists for dedicated / overlay links
  › In general, Jisc plans regular network upgrades with a view to ensuring that there is sufficient latent capacity in the network

# Janet backbone, Jun 2017 (exam period!)

**Major External Links**

# 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; avoiding local bottlenecks
  › Network performance measurement
  › Data transfer node (DTN) design and configuration
» Important to apply your security policy without impacting performance
  › Note some of Eric's comments yesterday on applying stateless filters / ACLs

» 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
» No current funding equivalent in the UK
  › But this can and should be part of your network architecture evolution

# Good news – Science DMZ happening already

» There are several examples of sites in the UK that have a 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 are already doing it, especially GridPP sites ☺

» Examples in the UK:

› Diamond Light Source

› JASMIN/CEDA Data Transfer Zone

› Imperial College GridPP; supports up to 40Gbit/s of IPv4/IPv6

› To realise the end-to-end benefit, both ends need to apply the principles

# Other examples of campus network engineering

» In principle you can just use your Janet IP service for data-intensive science

» Many sites split their connectivity, e.g. 10G campus, 10G GridPP, 20G resilience

› And then apply Science DMZ principles to the GridPP path

» Where specific point-to-point guarantees are required, there's Janet Netpath Plus

› See https://www.jisc.ac.uk/netpath  for details

› But then you'll not be able to exceed / burst beyond that capacity


» The WLCG has used physical / virtual overlays

› LHCOPN (private optical network) / LHCONE (virtual network)

» But how should campuses cater for multiple data-intensive science disciplines?

› Would one new overlay network per research community scale?


» At least one site is exploring SDN, using Cisco Application Centric Infrastructure (ACI)

# Choice of data transfer tool?

» Researchers will find a wide range of data transfer tools available
  › Choice may be dictated by the facility they're accessing, or made by word of mouth
» There's the simpler old friends
  › FTP, scp
  › But these are likely to give a bad initial impression of what your network can do
» There's TCP-based tools designed to mitigate the impact of packet loss
  › GridFTP typically uses four parallel TCP streams
  › Globus Connect has become popular - https://www.globus.org/globus-connect
» There's specific tools to support management of transfers
  › FTS – see http://astro.dur.ac.uk/~dphoelh/documentation/transfer-data-to-ral-v1.4.pdf
» There's also a commercial UDP-based tool, Aspera
  › See http://www.asperasoft.com/
» It would be good to establish more benchmarking of these tools at Janet campuses

» We (Duncan Rand and I) have visited a dozen or so sites so far
  › Met with a variety of networking staff and researchers
  › And spoken to many others via email
» Really interesting to hear what sites are doing
  › Some good practice evident, especially in local network engineering
  › e.g. routing those elephants around main campus firewalls
  › Campus firewalls typically not designed for single high throughput flows

» Seeing varying use of site links
  › As mentioned before, GridPP sites tend to split out their GridPP traffic
  › Some sites are using their 'resilient' link for bulk data transfers (be careful here!)
» Some rate limiting of researcher traffic
  › We'd encourage sites to talk to us (Jisc) about capacity rather than rate limit

# Example case: Diamond – University X

» Spoke with Diamond Light Source staff

   › Become evident that some of their university customers are using hard disk transfers

   › Science DMZ deployed at Diamond, but not at many of the customer sites

» Example: University X

   › Running ~6 experiments at DLS per year, producing 10-40TB of data per run

   › Takes 3 weeks from start to finish to get data locally for further processing

   › Remember that moving 100TB requires approx 10Gbps throughput for 24 hours

» Arranged a staged plan with X:

   › Benchmark current performance, with existing 1Gbps connectivity to data store

   › Repeat using Globus tools

   › Upgrade data store node to 10Gbps   *<- at this stage now; encouraging results*

   › Then repeat with 10G DTN at campus edge (Science DMZ)

» Has shown the value of getting researchers and network admins sat together ☺
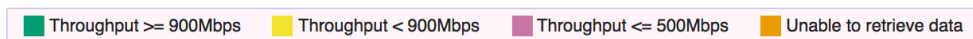
# Measuring network characteristics

» Important to have telemetry on your network

» The Science DMZ model recommends perfSONAR for this
  › Current version now 4.0; would encourage upgrades; but do check h/w requirements first

» Collects telemetry over time; throughput, loss, latency, path, allows retrospective viewing of data
  › Uses proven tools under the hood such as iperf
» Can run 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
» See our 'case for perfSONAR deployment' document (again, feedback welcomed!)
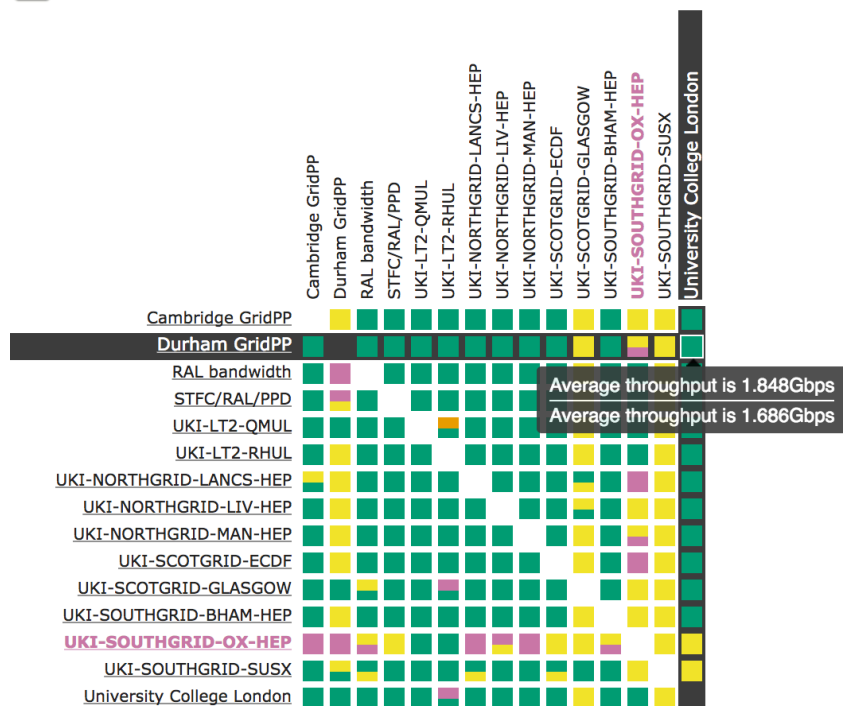» https://community.jisc.ac.uk/system/files/1862/The%20case%20for%20perfSONAR%20deployment-v1-0.pdf

UK Mesh Config - IPv4 Bandwidth Tests

UK Mesh Config - IPv4 Latency Tests

# Janet perfSONAR / DTN test node(s)

» We've installed a 10G perfSONAR node at a Janet PoP in London
  › Lets you test your site's throughput to/from Janet backbone
  › Useful if you want to run data-intensive applications in the near future, but don't have perfSONAR at the far end, or to benchmark your site's connectivity
  › Ask us if you're interested in using it

» We're planning to add a second perfSONAR test node in our Slough DC
  › Also planning to install a 10G reference DTN there
  › Will be a FIONA appliance - see https://fasterdata.es.net/science-dmz/DTN/
  › Will allow disk-to-disk tests, using a variety of transfer tools

» We can also run a perfSONAR mesh for you, using MaDDash on a VM
» We may also deploy an experimental perfSONAR node for TCP-BBR tests

# Small node perfSONAR

» In some cases, just an indicative perfSONAR test is useful

» i.e., run loss/latency tests as normal, but limit throughout tests to 1Gbit/s

» For this scenario, you can build a small node perfSONAR system for under £250

» Jisc took part in the GEANT small node pilot project, using Gigabyte Brix:
  › IPv4 and IPv6 test mesh at http://perfsonar-smallnodes.geant.org/maddash-webui/

» We now have a device build that we can offer to communities for testing
  › Aim to make them as 'plug and play' as possible
  › And a stepping stone to a full perfSONAR node

» Further information and TNC2016 meeting slide deck:
  › https://lists.geant.org/sympa/d_read/perfsonar-smallnodes/

# perfSONAR small node test mesh

# Aside: Google's TCP-BBR

» Traditional TCP performs poorly even with just a fraction of 1% loss rate
» Google have been developing a new version of TCP
  › TCP-BBR was open-sourced last year
  › Requires just sender-side deployment
  › Seeks high throughput with a small queue
  › Good performance at up to 15% loss
  › Google using it in production today
» Would be good to explore this further
  › Understand impact on other TCP variants
  › And when used for parallelised TCP applications like GridFTP
» See the presentation from the March 2017 IETF meeting:
  › https://www.ietf.org/proceedings/98/slides/slides-98-iccrg-an-update-on-bbr-congestion-control-00.pdf



**Optimal operating point**

RTT

Delivery rate

Optimal: max BW and min RTT (Kleinrock)

BDP    amount in flight    BDP + BufSize

# E2E performance to cloud compute

» We're seeing interest in the use of commercial cloud compute
  › e.g. to provide remote CPU for scientific equipment, or for cloudbursting
» Complements compute available at ESPRC Tier-2 HPC facilities
  › https://www.epsrc.ac.uk/research/facilities/hpc/tier2/
» Anecdotal reports of 2-3Gbit/s into AWS
  › e.g. by researchers at the Institute of Cancer Research
  › See presentations at RCUK Cloud Workshop - https://cloud.ac.uk/
» Bandwidth fo AWS depends on the VM size
  › See https://aws.amazon.com/ec2/instance-types/

» We're exploring measurement of cloud compute connectivity performance further
  › Includes AWS, MS ExpressRoute, GCP, …
  › And scaling Janet connectivity to these services as appropriate

# Building on Science DMZ?

» We should seek to establish good principles and practices at all campuses
  › And the research organisations they work with
  › There's already a good foundation at many GridPP sites
» The Janet backbone is heading towards 1Tbps capacity and beyond
» Jisc can help to seed further communities of good practice on this foundation
  › e.g. DiRAC HPC community, SES consortium, …
» And grow a Research Data Transfer Zone (RDTZ) within and between campuses
  › Towards a UK RDTZ
  › Inspired by the US Pacific Research Platform (PRP) model of multi-site, multi-discipline research-driven collaboration built on the NSF's Science DMZ investment
» Many potential benefits, such as enabling new types of workflow
  › e.g. streaming data to CPUs without the need to store locally
  › And hopefully no more hard disk transfers!

# Useful links

» Janet E2EPI project page
  › https://www.jisc.ac.uk/rd/projects/janet-end-to-end-performance-initiative
» E2EPI Jisc community page
  › https://community.jisc.ac.uk/groups/janet-end-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
» Fasterdata knowledge base
  › http://fasterdata.es.net/
» eduPERT knowledge base
  › http://kb.pert.geant.net/PERTKB/WebHome

# Please feel free to get in touch!

**Dr Tim Chown**

Network Development Manager

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

**jisc.ac.uk**