Queen Mary
University of London

# Lessons from Optimising data transfers at QMUL

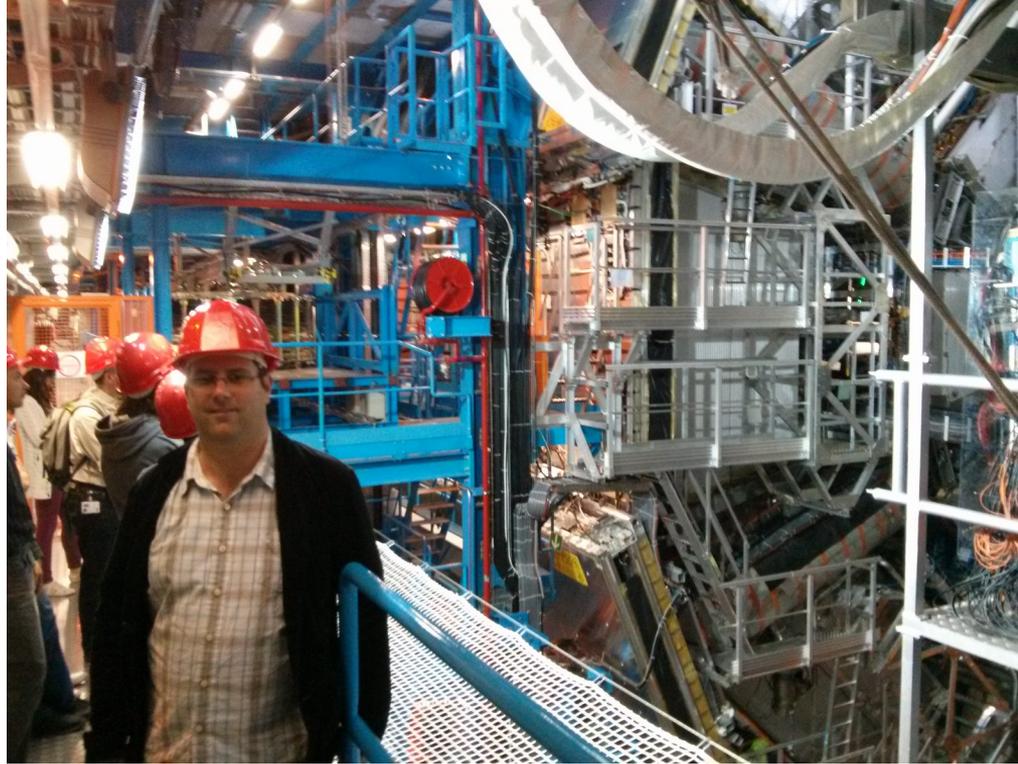Christopher J. Walker
<C.J.Walker@qmul.ac.uk>

# Overview

- What can be achieved
- TCP/IP explained
- Bottlenecks found
- Conclusions

# Motivation (LHC@CERN)

- Collisions  25ns
  - 100 PB/year
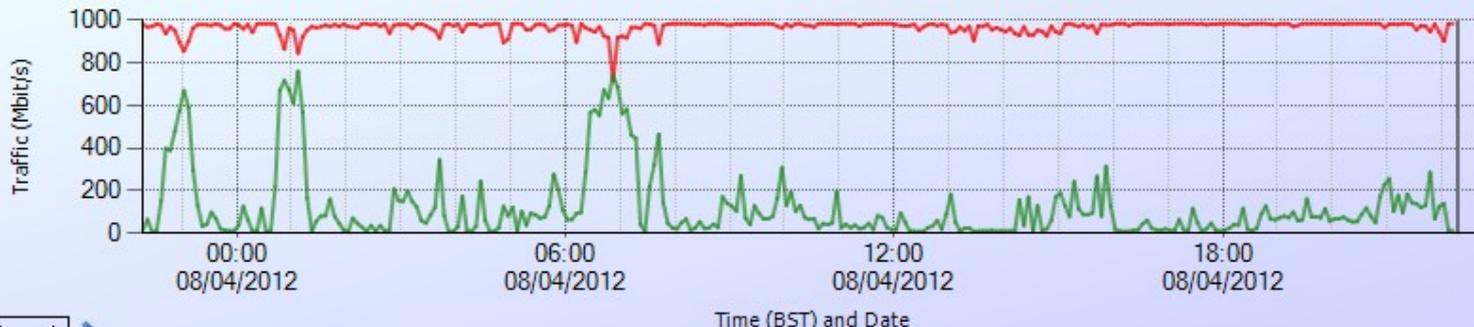- QMUL
  - Small fraction

# Network Expectations

- 1 Terabyte can be transferred in:
  - 100 Mbps network : 30 hrs
  - 1 Gbps network : 3 hrs
  - 10 Gbps network : 20 minutes
- Takes work to achieve this in practice
  - TCP tuning
  - Find and eliminate bottlenecks
  - Reduce packet loss
- Fasterdata.es.net
  - Excellent source of information

# WAN Performance
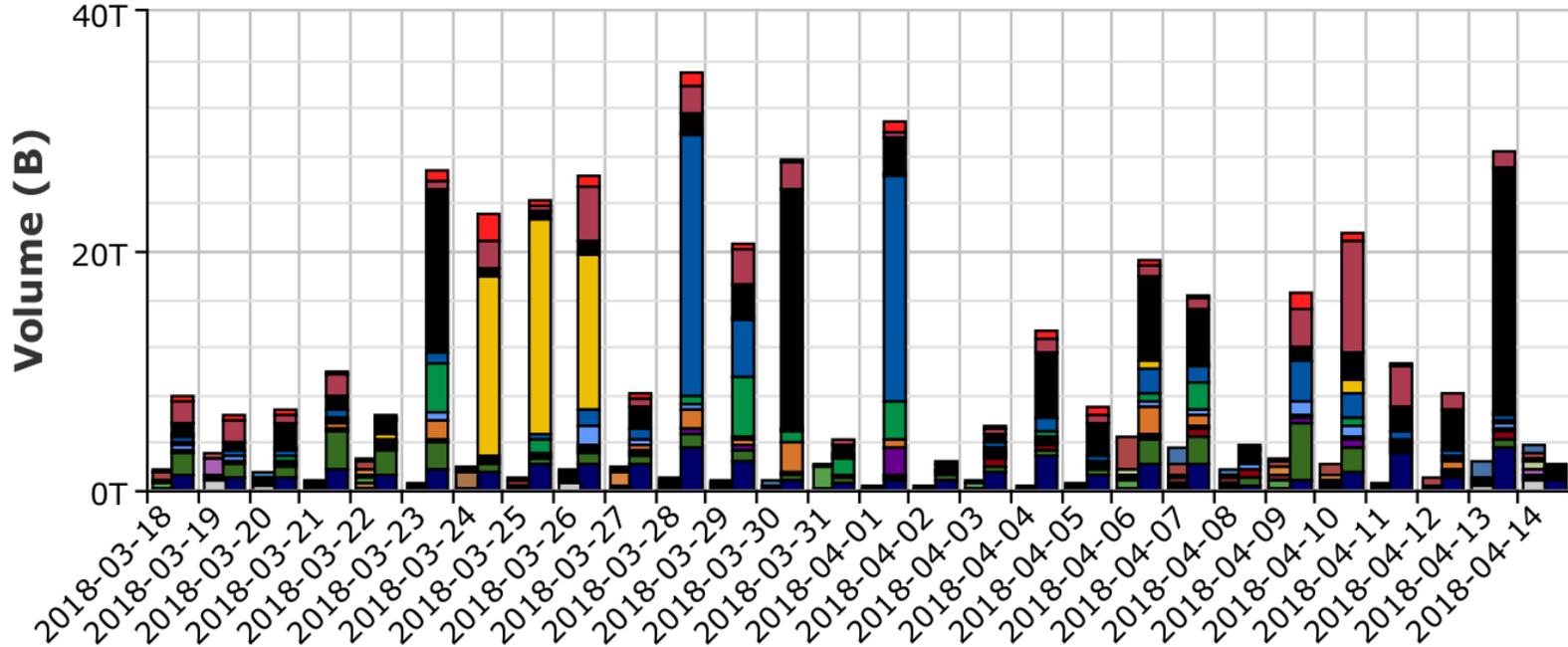


- April 2012:
  - 1 Gbit

- Sept 2013
  - 10Gbit

WLCG Sites

# Transfer Volume In/Out
## 2018-03-18 00:00 to 2018-04-15 00:00 UTC

Sources / Destinations

- UK RAL-LCG2
- UK RAL-LCG2-ECHO
- UK UKI-LT2-RHUL
- UK UKI-NORTHGRID-MAN-HE
- UK UKI-NORTHGRID-SHEF-HEP
- UK UKI-SCOTGRID-ECD
- UK UKI-SCOTGRID-GLASGOW
- UK UKI-SOUTHGRID-OX-HEP
- CA
- CERN
- DE
- ES
- FR
- IT
- ND
- NL
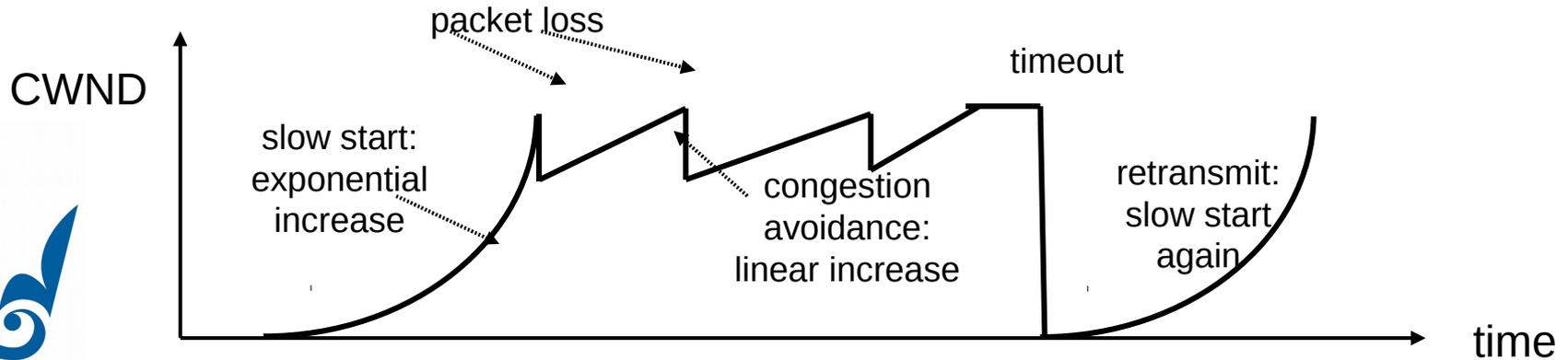- RU
- TW
- UK
- US
- 10 OTHERS

# scp or ssh+rsync: good enough?

- 100MB testfile copy to QMUL HPC in Slough

  - 4G mobile: 32s

  - Southampton eduroam: 13s

  - 100Mbit desktop: 10s

  - QMUL physics (1Gbit): 2.5s

# TCP: A short overview

- Congestion window (CWND) = the number of packets the sender is allowed to send
  - The larger the window size, the higher the throughput
  - Throughput = Window size / Round-trip Time
- TCP Slow start
  - exponentially increase the congestion window size until a packet is lost
  - this gets a rough estimate of the optimal congestion window size

packet loss

timeout

CWND

slow start:
exponential
increase

congestion
avoidance:
linear increase
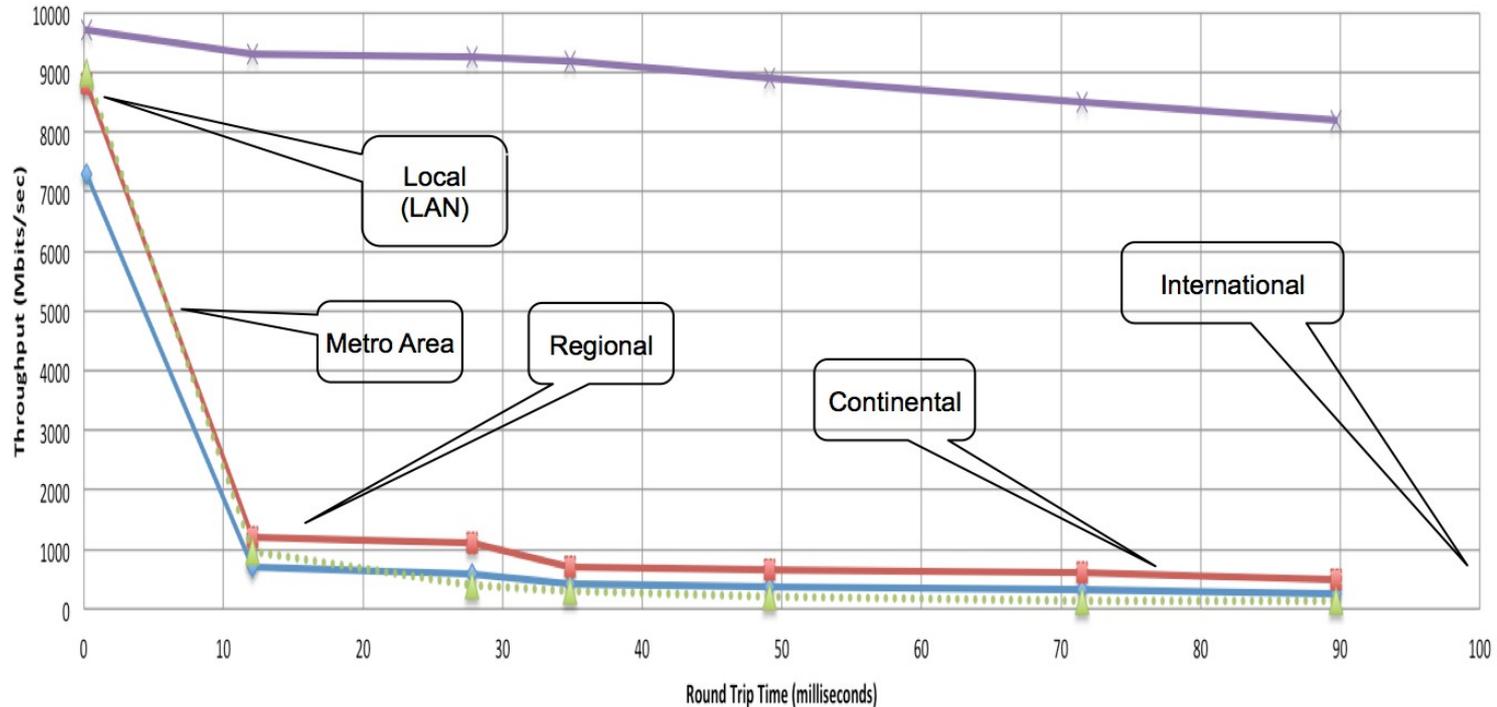
retransmit:
slow start
again

time

# TCP Tuning

- Latency: time to send 1 packet from the source to the destination
- RTT: Round-trip time
- Bandwidth*Delay  = Bandwidth Delay Product
  - The number of bytes in flight to fill the entire path
  - Example:  10 Gbps path; ping shows a 90 ms RTT (QMUL->BNL)
    - BDP = 10 * 0.090  = 0.9 Gbits (112 MBytes)
  - QMUL ->Taiwan 273ms RTT (at 10Gbps path)
    - BDP = 10*0.273 = 2.73 Gbits (340 MBytes)

# Effect of Packet loss with distance



**Throughput vs. increasing latency on a 10Gb/s link with _0.0046%_ packet loss**

Local (LAN)

Metro Area

Regional

Continental

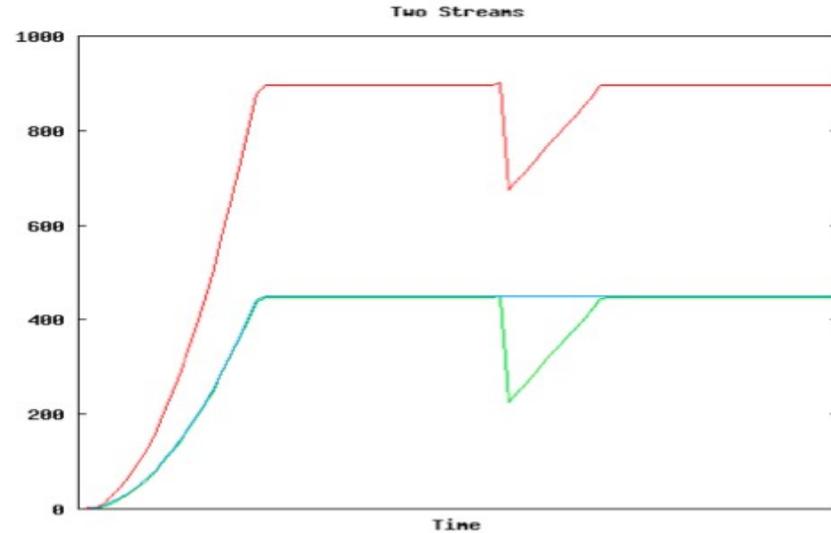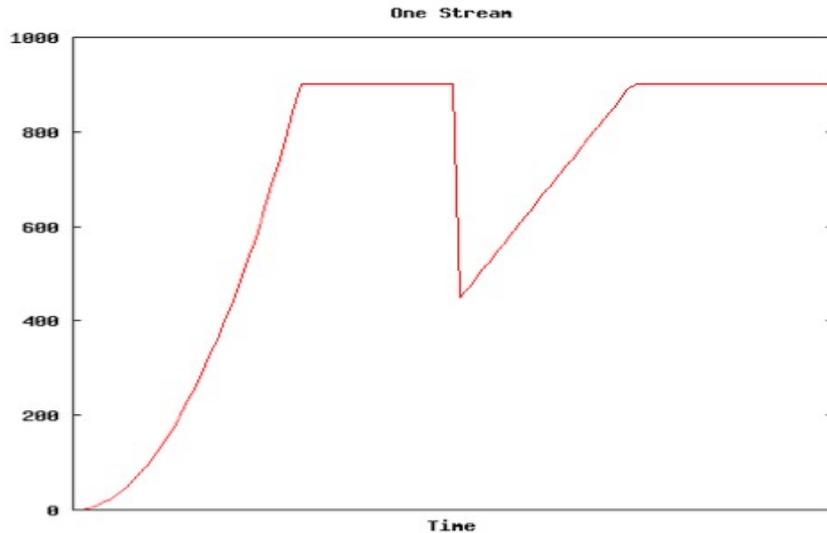International

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

# Multiple streams

- Parallel streams can help
  - Potentially unfair on other users

# TCP lessons

- Increase TCP buffers for distant transfers
    - Fasterdata.es.net has good recommendations
- Packet loss needs eliminating
- Application
    - large buffers (not scp)
    - Multiple streams
    - GridFTP has these
- Aspera uses UDP (and GridFTP can)
- Fasterdata.es.net has excellent recommendations

# Slough backup – case study

- 10Gbit (*2) link QMUL ↔ Slough
  - 1.5ms RTT
  - 1GB/s expected transfer rate
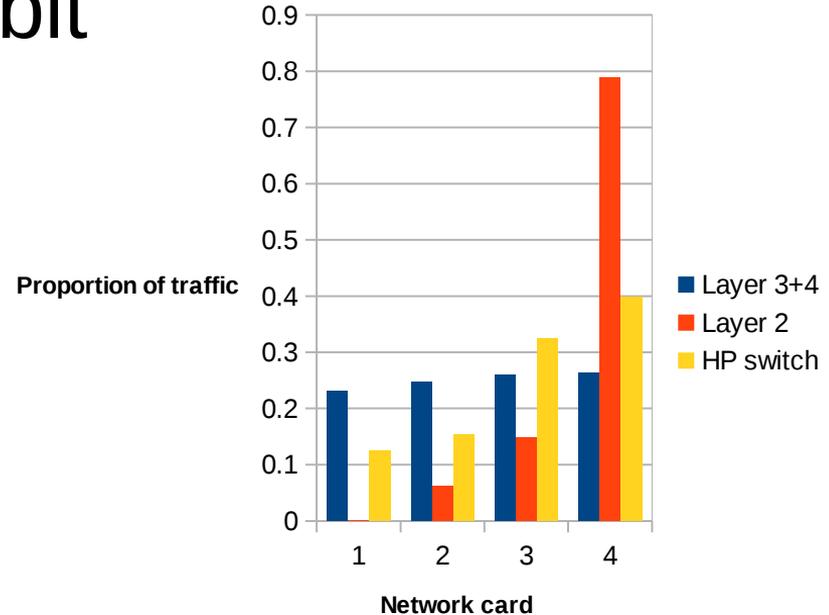    - 40MB/s achieved initially

# TSM's TCPWindowSize

- TCPWindowSize 63 #TSM default
  - 39,389.44 KB/sec
- TCPWindowsize 10 # Reduced for testing
  - 22,591.03 KB/sec
- TCPWindowSize 0 # Use Linux default
  - 129,721.81 KB/sec

# Packet loss

- Iperf measurements from TSM to:
  - Backup1 (10.x.y.101): 9.6GBit/s
  - Backup2 (10.x.y.102): 0.6 Gbit/s
    - Adjacent IPs
      - Packets take different legs of bond – packet loss on one of them
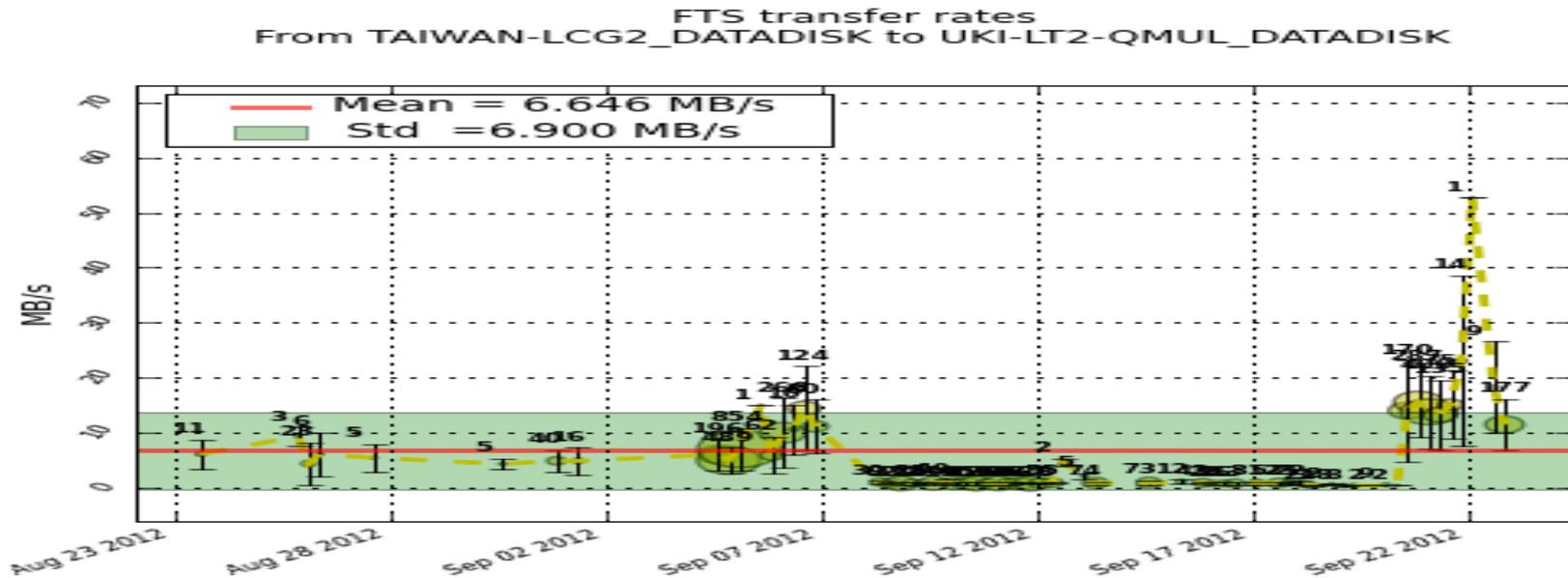  - Packet loss: 0.015%
    - http://fasterdata.es.net/network-tuning/tcp-issues-explained/packet-loss/

# Bottlenecks found

- Gbit connected at 100Mbit
  - GridFTP node
  - Dept
  - College
  - NOC ↔ Archer
- Routing:
  - QMUL ↔ CERN via USA

# Routing problems with 10Gbit/s upgrade

- 8th September 10Gbit/s WAN upgrade
- UK sites – increased rates
- ASGC (Taiwan) decrease
  - Route not advertised via GEANT.



FTS transfer rates
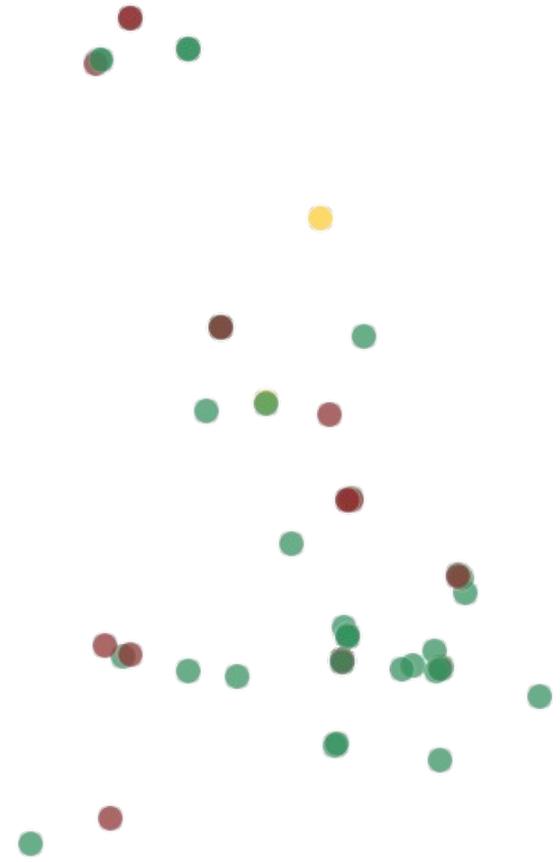From TAIWAN-LCG2_DATADISK to UKI-LT2-QMUL_DATADISK

# Firewalls

- ICMP
  - Often blocked
    - timeout rather than failure
- IPv6
  - Tracepath6 blocked  (ICMP blocking )
- Barclays bank blocked
  - Deep packet inspection and rewriting of http packets (but not https)
- Scp failing half way through transfer
- GridFTP Slow performance
  - 1 MB/s through firewall, 50MB/s avoiding firewall
- GridFTP control connection forgotten

# QMUL plans

- Slough Firewall upgrade
  - 10Gig module
- HPC Data transfer node
- Improve monitoring
  - Perfsonar
  - RIPE atlas

/QMUL     @QMUL

Queen Mary University of London

# Conclusions

- Large transfers routine
  - Need to transfer multiple files at once and deal with failures
  - Takes work  (GridPP sites have this experience)
- Many "small" transfers in the UK can work with rsync/ssh
- Monitoring vital
  - Transfers
  - Network
- Network
  - Need Low packet loss (most problems in the last mile)
  - Good relationship with network team useful
- Information
  - Fasterdata.es.net

# Acknowledgements

- Fasterdata.es.net (Brian Tierney)
  - Many thanks for the TCP tuning slides

# Backup slides

/QMUL    @QMUL

Queen Mary University of London

# IPv6

- Routes
  - May be different to IPv4
    - Geneva ->QMUL via New York (fixed)
- Software (IPv6) / ASIC (IPv4)
  - Older routers may give poor performance
- Preferred over IPv4
  - If IPv6 address (AAAA record) in DNS, it will be used by machines that think they are IPv6 connected.
- Blocked differently by firewalls