

# IPv6 Only at Imperial

David Stockdale  
ICT Networks Group  
[david@imperial.ac.uk](mailto:david@imperial.ac.uk)

---

## Facts and figures

- Over 65,000 unique hosts on wired network
  - Over 60,000 unique hosts on wireless network
  - Over 26,000 concurrent wireless clients at peak time
  - 2x100G to Janet
  - Most hosts within VRFs (MPLS L3VPNs)
  - Firewalls between VRFs
  - No NAT(44)
  - Firewalls capable of NAT64
-

## Our current position

- ~35% of our Internet traffic IPv6 (nearly 50% on BYOD)
- Dual stack on production, guest & BYOD (including wireless)
- AAAAs on most load-balanced services
- Other services enabled:
  - Home directories (>95% IPv6!)
  - 10PB research data storage (~~IPv6 only~~)
  - Mail, DNS, HEP systems
- SLAAC rather than DHCPv6
- Feature parity mandated in tenders

## HPC refresh

- Multi-year programme to replace HPC estate
  - Year 1: 7 racks, 30 servers in each
  - By the end: 40 racks, 1PFLOPS
  - Existing servers 1/10G Ethernet plus Infiniband
  - New servers 100G Ethernet with RoCE
  - Speaks to IPv6 enabled research data storage
-

## HPC refresh

- An opportunity to go IPv6 only! How hard can it be!?
- Pesky IPv4 only PBS for starters
- 2x spine switches and leaf per rack
- EBGP, ASN per switch
- /64 IPv6 and /24 IPv4 per leaf, MTU 8000
- MP-BGP between switches, IPv6 sessions only
- IPv6 only... to rest of College network
- /32 IPv4 route on servers via local gateway for PBS
- 1G management network truly IPv6 only

## HPC refresh

- So far, so good!
- How do we boot them?
- ... DHCPv6 ... and UEFI
- SLAAC, RDNSS
- Plan A:
  - Stateless DHCPv6 server on switch returning PXE options
  - Server: “I only support DUID-UUID”
  - Switch: “Unsupported DUID type 4”
  - Us: :-(

## HPC refresh

- Plan B:
  - Stateless DHCPv6 relay to Kea returning PXE options
  - Server: “Send me a Bootfile URL in option 59”
  - Kea: “Sure thing”
  - Server: “And reflect Vendor Class option 16 back at me”
  - Kea: “What?”

## HPC refresh

- Plan C:
  - Stateless DHCPv6 relay to ISC returning PXE options
  - Server: “I see your RA. SLAAC address, done”
  - Server: “Other Configuration flag set? Will do”
  - Server: “Managed Address Configuration flag unset...”
  - Server: “... Request! Request! Request!”



## HPC refresh

- Plan D:
  - Stateful\* DHCPv6 relay to ISC
  - Success!
- Sort of... hello iPXE
- iPXE: “Information-Request!”
- Switch: *drops packet*
- iPXE: “RA said I needed that reply. I’ll wait... forever...”

## HPC refresh

- iPXE recompiled to not send Information-Requests
  - iPXE: “I don’t care, the NIC wasn’t initialised anyway”
  - iPXE recompiled with bodgetastic initial sleep
  - Us, 2 days in: “OMG it’s actually booting!!!”
-

## HPC refresh

- Machines booted!
- How to talk to rest of the world?
- ... NAT64/DNS64
  
- Presenting software exhibit A
- Exhibit A: “Where’s my licence server?”
- DNS: “Here A or here AAAA”
- Exhibit A: “I’ll go with the A then”
- One AAAA only hostname later... Fixed!

## HPC refresh

- Presenting software exhibit B
- Exhibit B: “Where’s my licence server?”
- DNS: “Here A or here AAAA”
- Exhibit B: “What’s a AAAA? What’s IPv6?”
- yum install clatd
- Exhibit B: *springs into life*
- Everyone: “Let’s party like it’s 1980 again!”

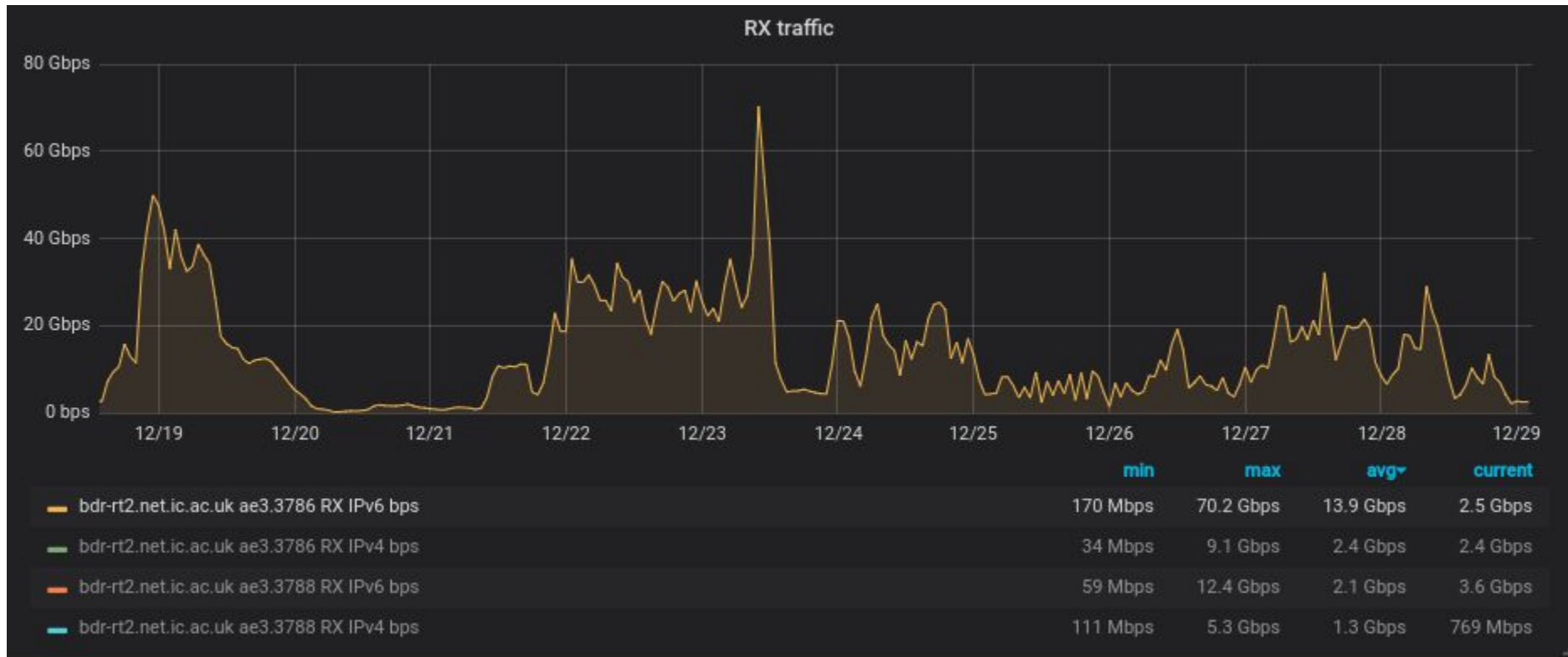
## HPC refresh

- We have a fully functional system
- One last thing... RoCE
- Oh and remember those jumbo frames?
  - >1440 MSS, IPv6 src -> NAT64 -> IPv4 dst
  - ... >1480 byte IPv4 responses from Internet to NAT64
  - ... >1500 bytes when translated to IPv6; too big for College network
  - Broken PMTUD to some remote hosts. Clamped MSS on firewalls.
- All 7 racks now in full production use
- Looking at options for PBS

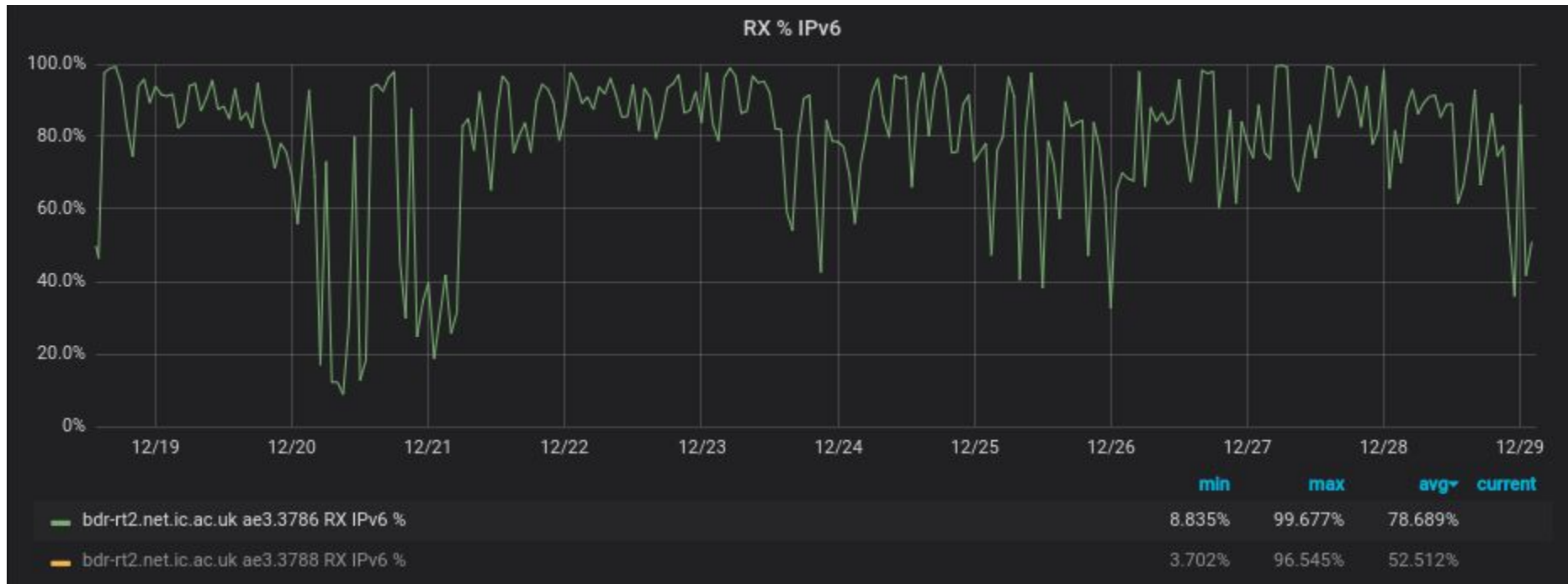
## What next?

- IPv6 enable remaining services
  - NAT64/DNS64 trials on wireless and wired
  - IPv6 only internal services
  - External services fronted by load-balancers
  - DHCPv6 in DC, PXE
  - Retire IPv4
  - Free up IPv4 address space - \$\$\$!
-

# HEP Internet traffic

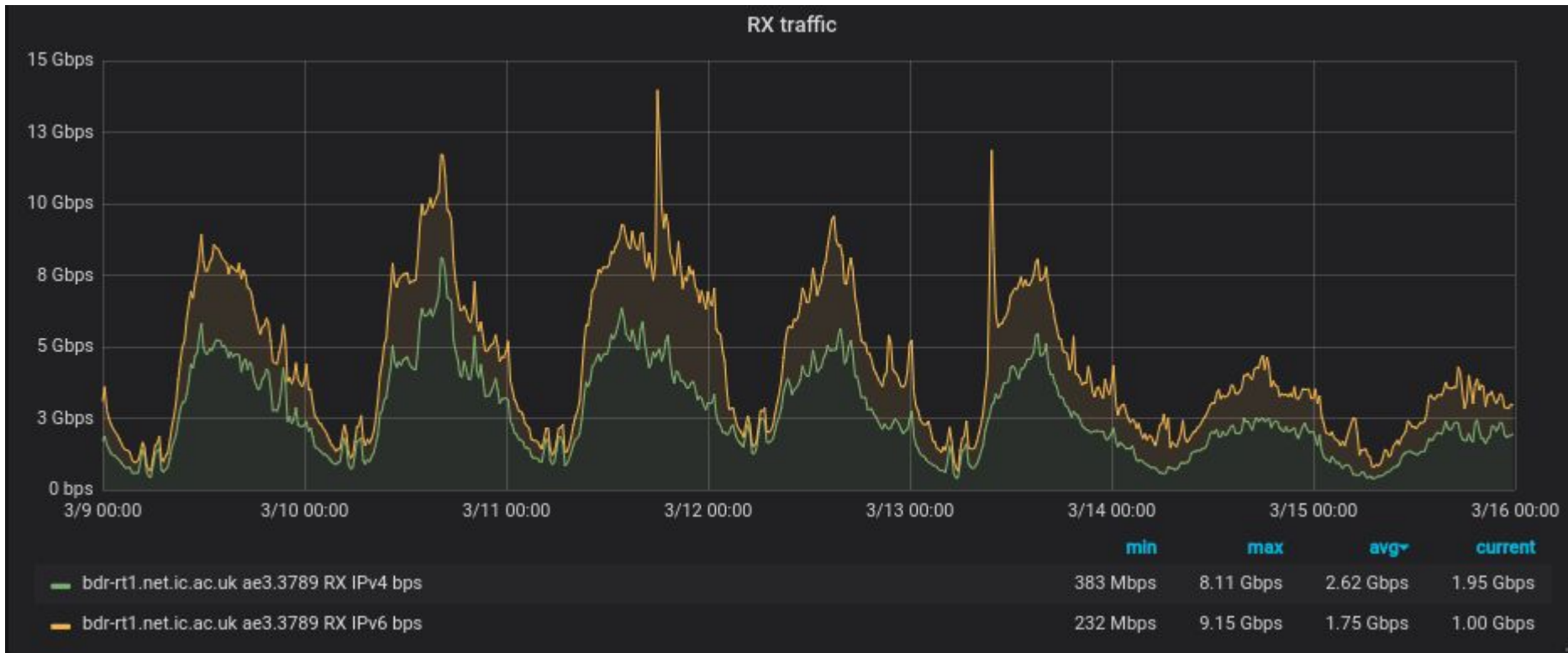


# HEP Internet traffic





# College Internet traffic



## College Internet traffic

