

Research Data Storage Platform

The Royal Marsden Hospital

Founded in 1851 as the Free Cancer Hospital to Classify Tumours, research the causes and find new treatments.

ICR Became a separate entity in 1909.



Institute of Cancer Research - at a glance



Top 4 global cancer research organisation



Top-ranked UK academic institution (REF)



20 drug candidates discovered since 2005



More than 1,000 staff



£161.9m income
£110.0m expenditure



Awarded Athena SWAN Silver



More than 900 scientific papers



Partnerships with 163 different companies



Top UK university for invention income



141 research students
143 MSc students

Making the discoveries

Our strategy to defeat cancer

The ICR and The Royal Marsden delivered a joint strategy covering the next five years

Our vision

We will overcome the challenges posed by cancer's complexity, adaptability and evolution through scientific and clinical excellence, innovation and partnership

First pillar: Unravelling cancer's complexity

Comprehend the full complexity of cancer by **harnessing the power of new technologies and Big data**



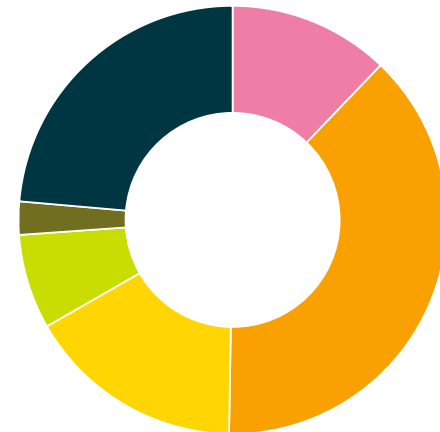
Substantial and diverse funding

Total incoming resources 2016

2016

- Total income £162m
- HEFCE 12% based on research excellence
- Grant income 38%
- Legacies and donations 7%
- Invention income from our discoveries 16%

-> *Services not centrally funded*



12%
Higher Education
Funding
Council for England

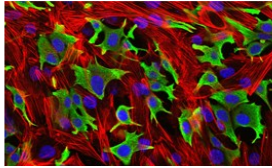
38%
Research grants

16%
Royalty income

7%
Legacies and donations
3%
Investment and tuition fees

24%
Sale of part of our future
royalty stream

ICR Research Divisions



Breast
Cancer
Research



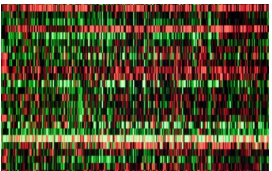
Cancer
Biology



Cancer
Therapeutics



Clinical
Studies



Genetics and
Epidemiology



Molecular
Pathology



Radiotherapy
and Imaging



Structural
Biology

Scientific Computing Service

- HPC
 - ~2000 cores, 2PB Lustre storage
- Storage
 - RDS Platform
 - Future: Archive, Data Transfer Service
- Scientific Software
- Edge / User support
 - Instrument connectivity, Staging areas etc

History of Data Storage

- 2006 – Pandorica
 - SGI CXFS 3 Tier system based on 2 sites
- 2011 – Gallifrey
 - Quantum StorNext 2 Tier (Disk + tape) system
 - Issues caused by retrieving small files from tape
- 2016 – RDS...

RDS Service - User Requirements

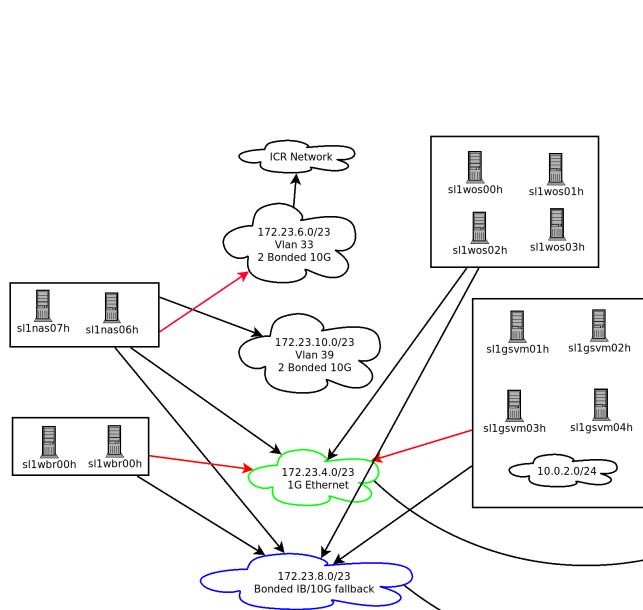
Collected from researcher workshops:

- 6PiB storage, expandable to growing research need; minimum 20PiB capacity
- High level redundancy ensuring robust solution
- Cost effective and competitively priced solution (-> vendors offered two Tiers)
 - Rapid access to data held in Tiers
 - Ability for researchers to manage data transfers between Tiers
 - Single namespace
 - Direct access to data in Tier 2
- Ability to protect against accidental loss of datasets

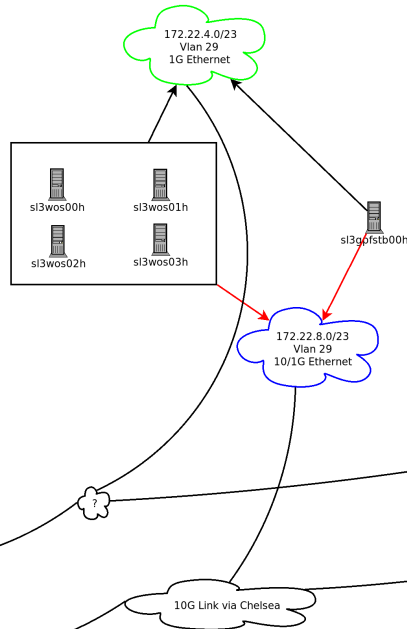
RDS Architecture

- DDN GS7Ks at Sutton and Slough providing 2PiB of mirrored storage
- DDN WOS using 3 site GOA policy (Sutton and Slough DH1 / DH3)
- GS Bridge to manage data between tiers
 - DMAPI handler enables transparent reads of migrated files
 - Mover uses policy engine rules to move data between tiers
- Total of 8 MediaScaler servers
 - DDN's package of CTDB / Samba / NFS
- Networking within each site (mostly) FDR InfiniBand
- 10G links between sites
- Daily snapshots (but no backup)

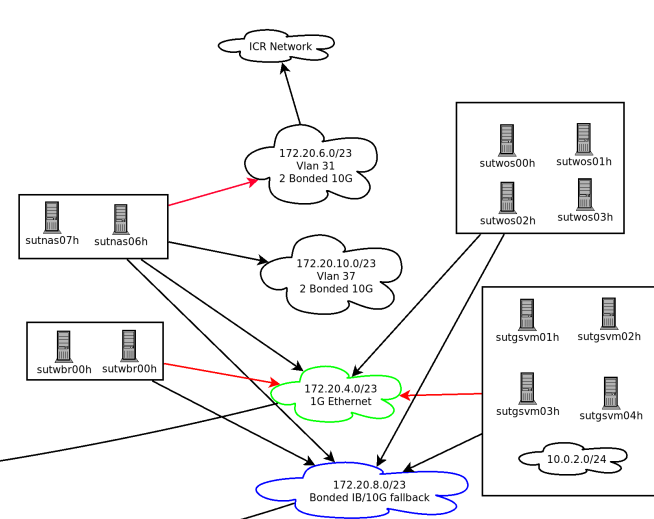
Slough DH1

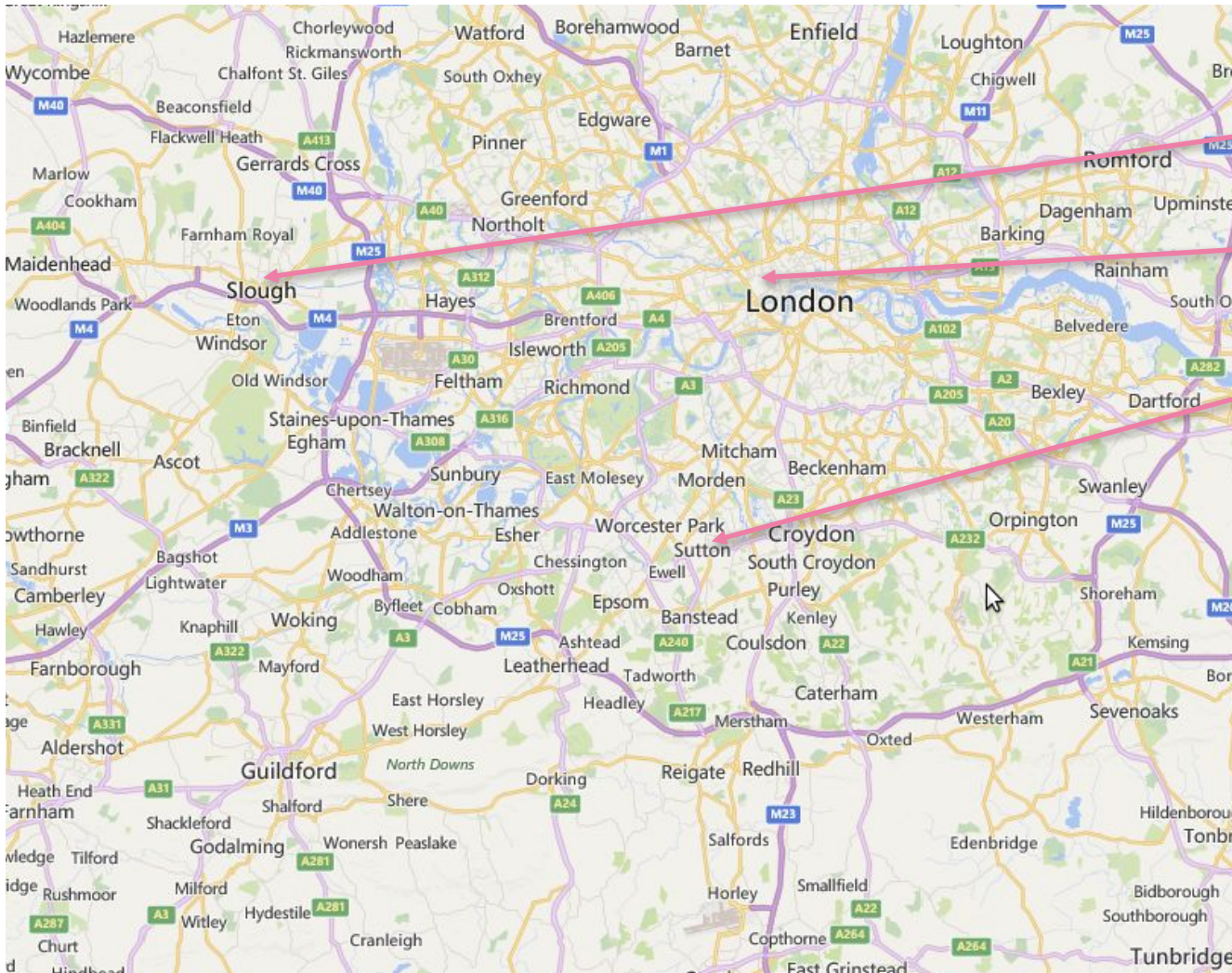


Slough DH3



Sutton





Slough

Chelsea

Sutton

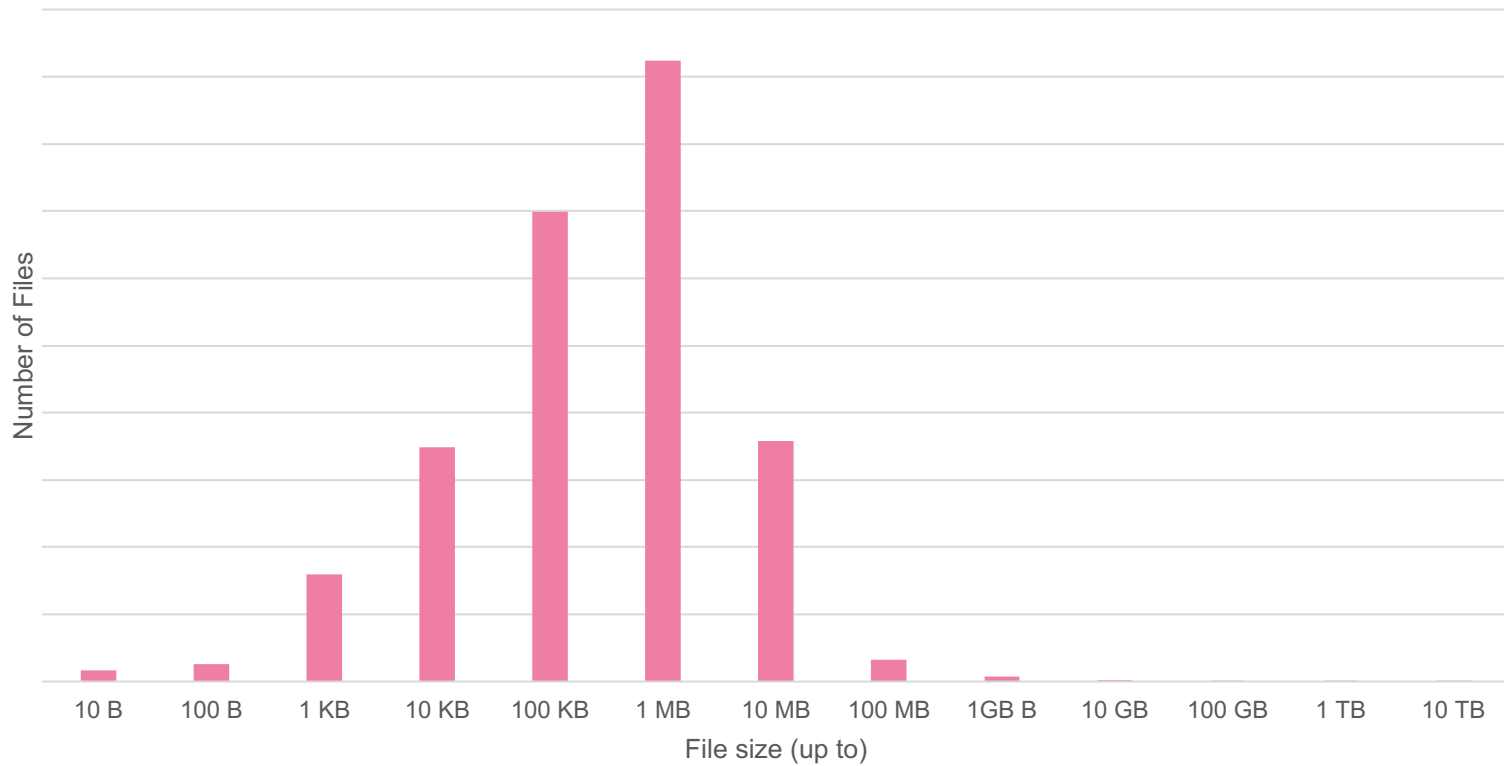
SpectrumScale Performance

- For large files - write performance limited by 10Gb/s inter-site links
 - (but need some network tuning first!)
- Can do better with read performance by setting readReplicaPolicy=fastest
- Small file performance less good (particularly via SMB)
- Mirroring working well
 - We have lost inter site link on a few occasions but recovery has been pretty smooth
 - Takes us ~2 hours for policy scan (100m files) + time needed to resync data
- “Double Size” files causes some confusion.

WOS Bridge Performance

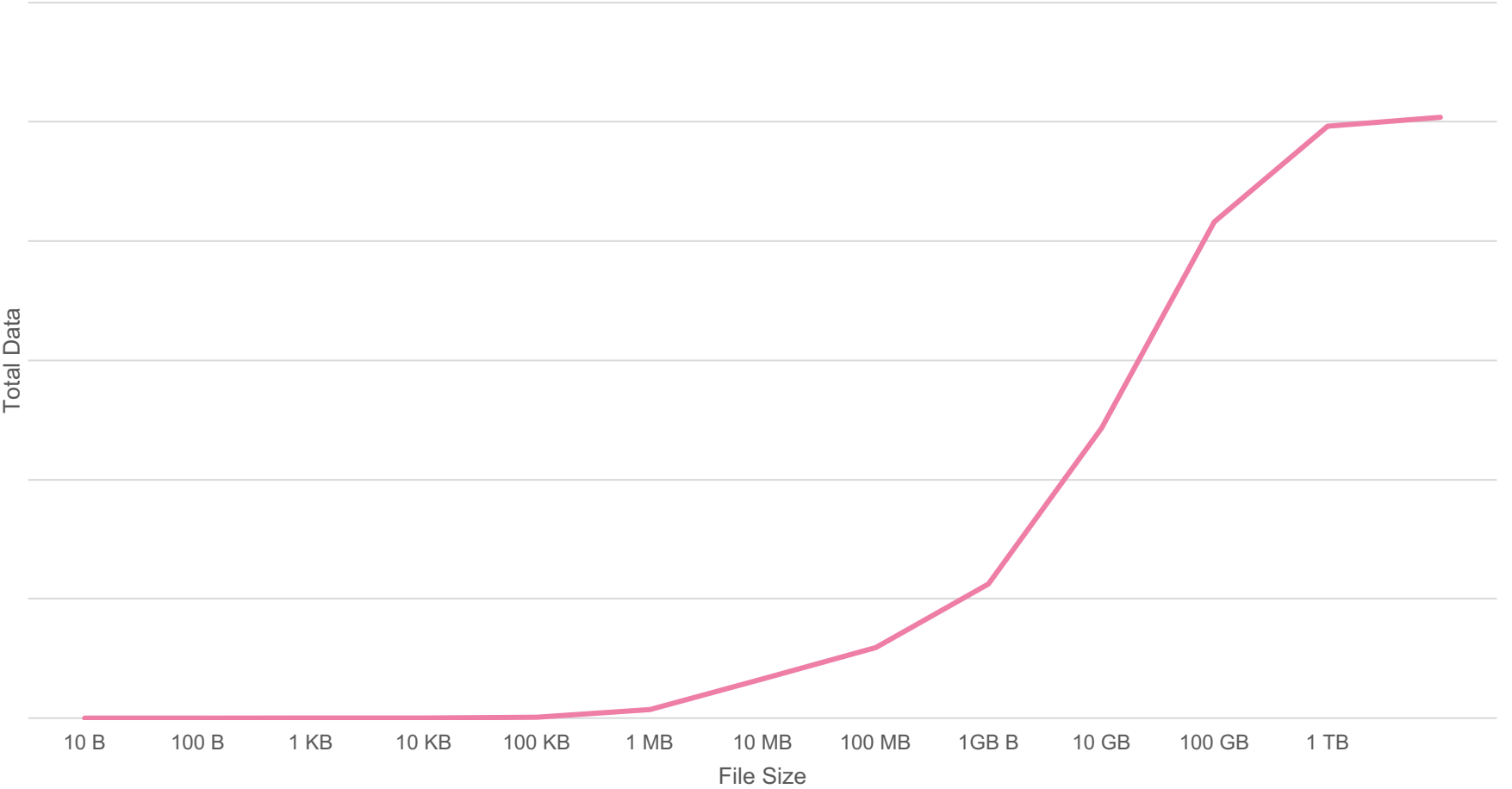
- Significantly slower than GPFS data - ~100MB/s
 - ...but can scale fairly well (in the right circumstances)
- Sometimes stops working at all
 - DDN have identified some issues with WOS but problems currently ongoing
 - Unfortunate interaction with NFS causes all MediaScaler nodes to fall over!
- Difficult to work out which data should be migrated just based on access times etc
 - Given performance difference we really want all active large files on GPFS

Number of files by size



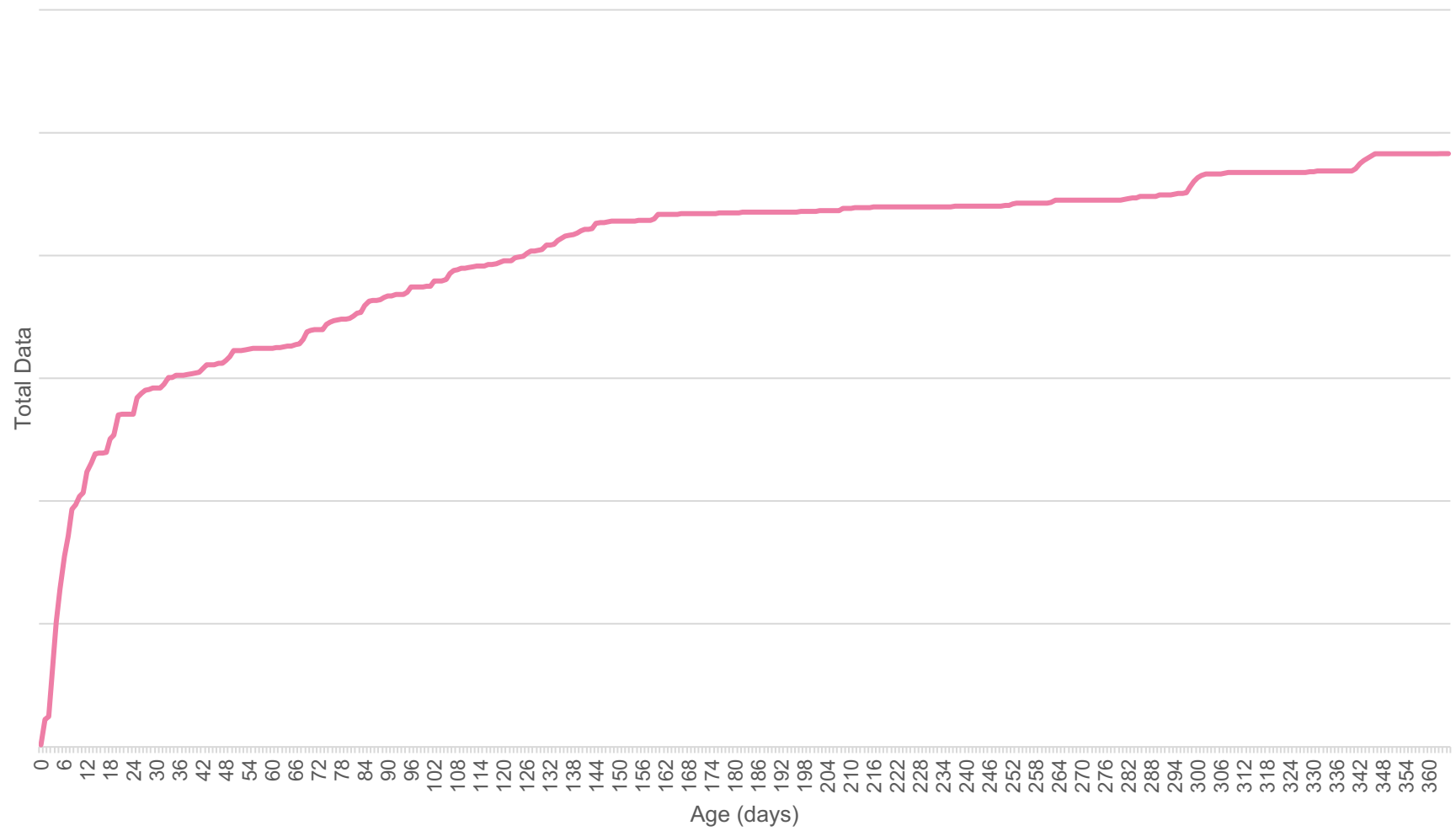


Total data by file size





Total data by access time

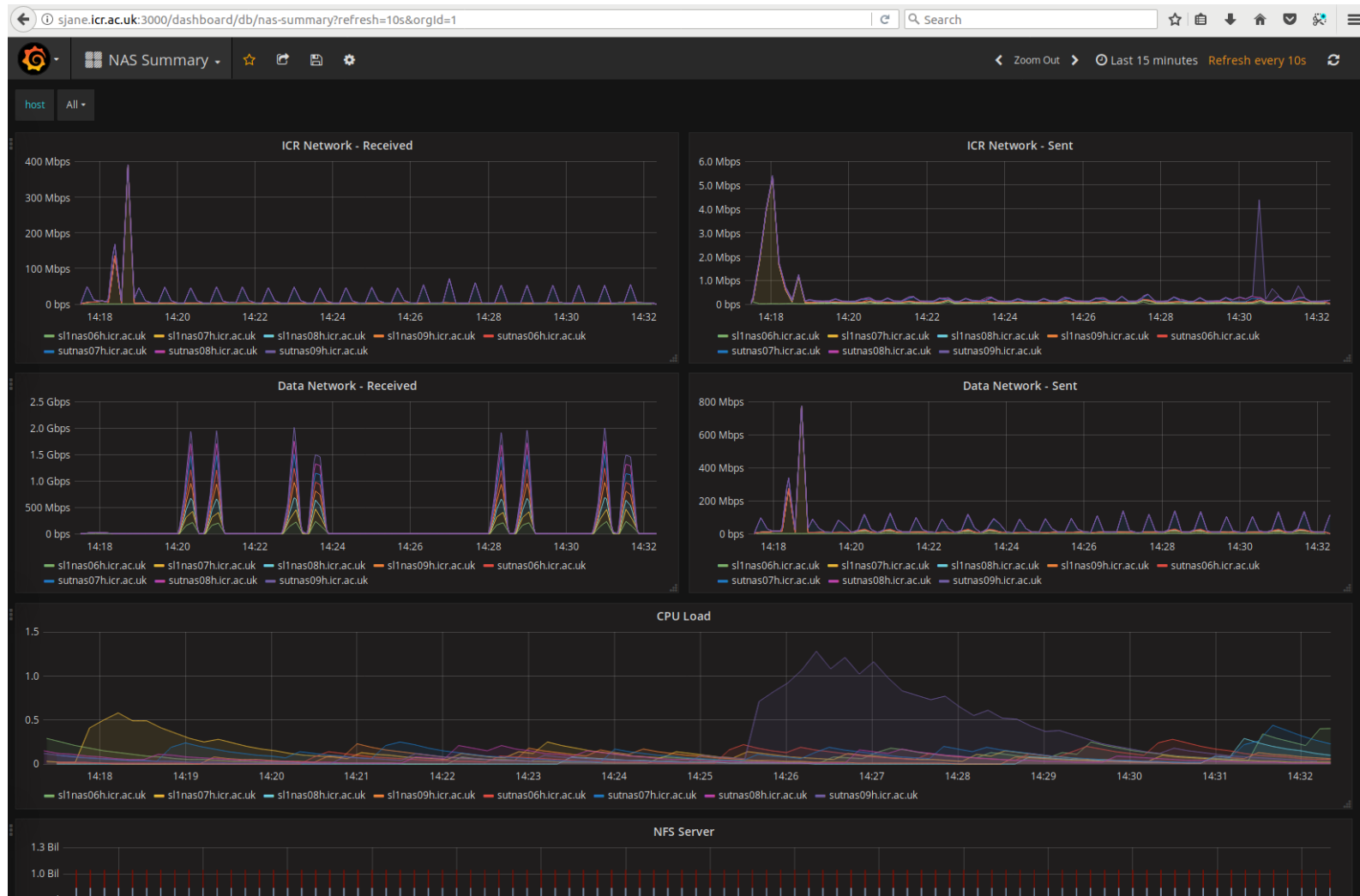


Monitoring

Currently using:

- Icinga2 to check health of system
 - Spectrum Scale Plugins from IBM work well
 - Monitoring NFS / Samba / CTDB more tricky
- Telegraf + InfluxDB + Grafana for metric storage / analysis
 - Haven't got any Spectrum Scale specific metrics in (yet)

Grafana

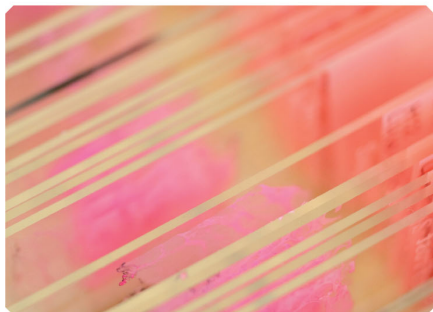
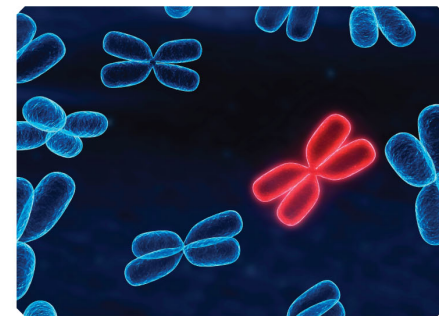


Filesystem Analysis Tool(s)

- Interested in tools for
 - Helping users see how much data they have
 - Understanding size / age profiles etc
 - Managing quotas
- Looked at
 - Starfish
 - ClarityNow
 - Perl
- In the long term, need to move to being more of a data management service

Unrivalled
track record

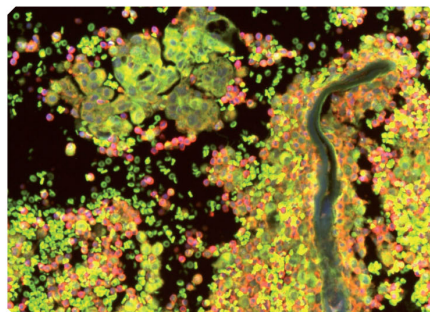
ICR The Institute of
Cancer Research



Making the
discoveries that
defeat cancer



ICR



One of the world's
most influential
cancer research
institutes

[Title slide with Mission Statement]

Our mission
is to make the
discoveries that
defeat cancer.

The Institute of Cancer Research is one of the world's most influential cancer research institutes.

Making discoveries

Scientists and clinicians are working every day in our labs to make a real impact on cancer patients' lives.

Translational research

Our unique partnership with The Royal Marsden and bench-bedside approach means we create and deliver results in a way that other institutions cannot.

Our unrivalled track record

Whether in cancer biology, genetics, personalised therapies or new drug discoveries, our achievements speak for themselves.

[Introduction] Heading

Activate second <List Level>
for lighter grey subtitle style

1

Use [Introduction] from the 'Slide Layouts' to highlight three specific areas

2

3

[Title and Content]

Activate second <List Level>
for lighter grey subtitle style

The first level is 19pt Arial and does not have bullets

- <Increase List Level> to activate the second style which does have a bullet style

[Biographies x1]

Activate second <List Level>
for lighter grey subtitle style

Image size:

H=4cm

W=2.8cm

The first level is 19pt Arial and does not have bullets

- <Increase List Level> to activate the second style which does have a bullet style

[Biographies x2]

Activate second <List Level>
for lighter grey subtitle style

Image size:

H=4cm

W=2.8cm

The first level is 19pt Arial and does not have bullets

- <Increase List Level> to activate the second style which does have a bullet style

[Biographies x3]

Activate second <List Level>
for lighter grey subtitle style

Image size:

H=4cm

W=2.8cm

The first level is 19pt Arial and does not have bullets

- <Increase List Level> to activate the second style which does have a bullet style

[V1 Image with caption]

Second <List Level> for subtitle style

The first level is 19pt Arial and does not have bullets

- <Increase List Level> to activate the second style which does have a bullet style

**Delete this positional and
insert required picture here**

Ensure that the Image size is:

H=13.9cm

W=16.8cm

[V2 Image with caption]

Second <List Level> for subtitle style

The first level is 19pt Arial and does not have bullets

- <Increase List Level> to activate the second style which does have a bullet style

**Delete this positional and
insert required picture here**

**Ensure that the Image size is:
H=13.9cm
W=12cm**

[Image] Heading

Second <List Level> for subtitle style

**Delete this positional and
insert required picture here**

Ensure that the Image size is:

H=13.9cm

W=24cm



[Title only] Heading

Second <List Level> for subtitle style

[Author] 19pt Arial

Activate second <List Level> for Title

Department

Location

name@icr.ac.uk



ICR

