

MOVING DATA: Getting up to speed with Globus and Science DMZ

Brian Flaherty
Richard Tumaliuan



NeSI

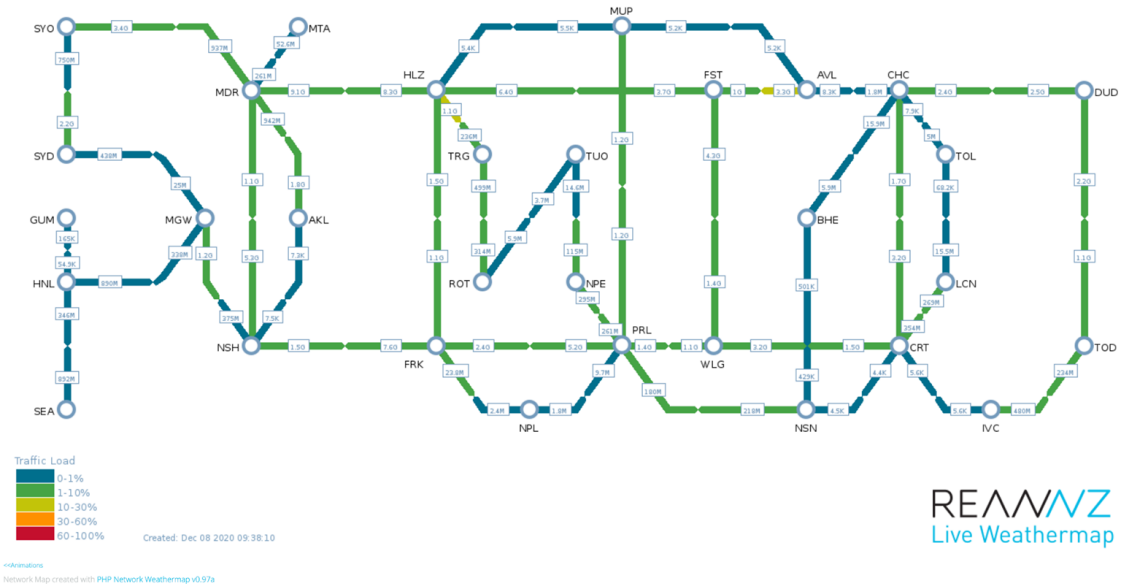
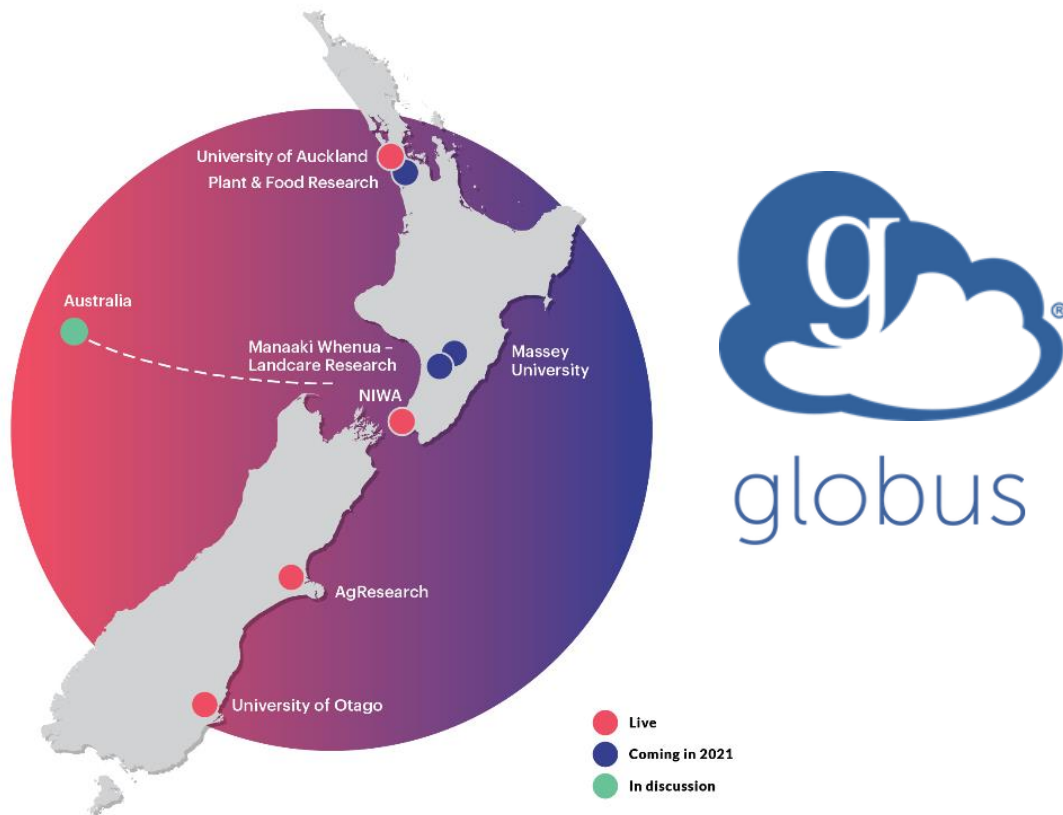
New Zealand eScience
Infrastructure

REANZ



Overview of the tools that support and enable fast, secure, and reliable transfers with Globus and Science DMZ

How these tools fit together to improve performance across the whole data transfer journey



REANNZ
Live Weathermap

Science DMZ

A lightweight and high performing on-ramp to the REANNZ international research and education network.

An opportunity for discussion groups to delve deeper into the user experience of Globus and the technical aspects of Science DMZ

Globus break out group - Show Globus users and potential Globus users how they can get value from the service, what is the difference between managed vs. personal endpoints? How users can engage with more advance transfer options and the steps for getting started.



Science DMZ break out group – technical aspects of the Science DMZ, transfer node reference design, and the REANNZ managed network edge and Science DMZ. Examples and technical requirements.

Science DMZ

A lightweight and high performing on-ramp to the REANNZ international research and education network.

NATIONAL DATA TRANSFER PLATFORM

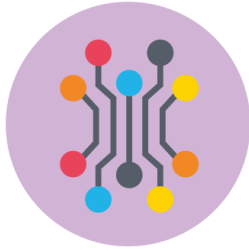
Globus – the world wide network of DTNs



NeSI REANZ

New Zealand eScience
Infrastructure





High performance computing (HPC) and analytics

- Fit-for-purpose national HPC platform including data analytics



Data transfer and share

- High speed, secure data transfer with end-to-end integration
- Hosting of large actively used research datasets, repositories, and archives



Training and researcher skill development

- In-person and online training to grow capabilities in NZ research sector
- Partnership with The Carpentries (global programme to teach foundational coding and data science skills to researchers)



Consultancy

- Computational science experts available to lift the computational capabilities of research teams, as well as optimise tools & workflows

National data transfer platform activities in 2020:

870 TB 138% increase on 2019	5,107 28% increase	182 42% increase
Amount of data transferred	Number of transfers made	Number of users

Globus benefits:

High-speed data transfer

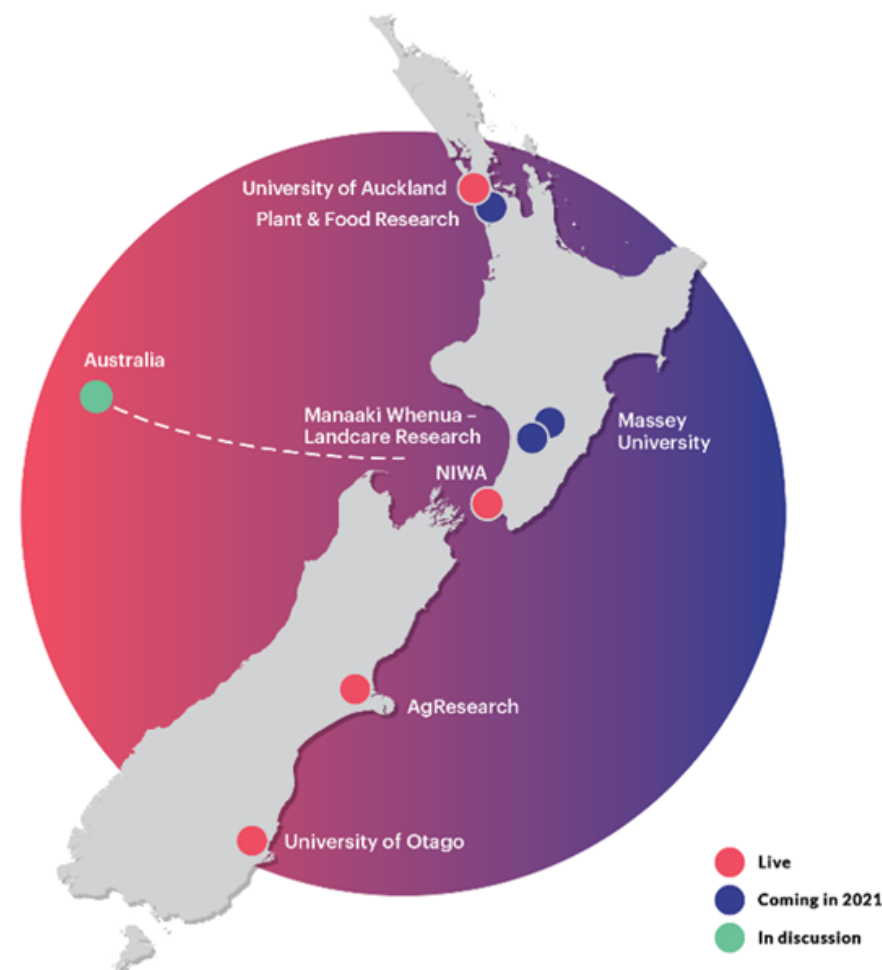
- Move gigabits of data on a network 1,000 times faster than broadband Internet. Powered by REANNZ (NZ's national advanced network provider), research data transfers can be done at 10Gbps.

Secure and easy data sharing

- Share your research data with collaborators locally, nationally, and around the world. Control who has access using group management tools.

Research data delivery network

- Take advantage of our partnering institutions across NZ who have existing data delivery network nodes. Transfer to and from any laptop or server using a Globus connect endpoint.



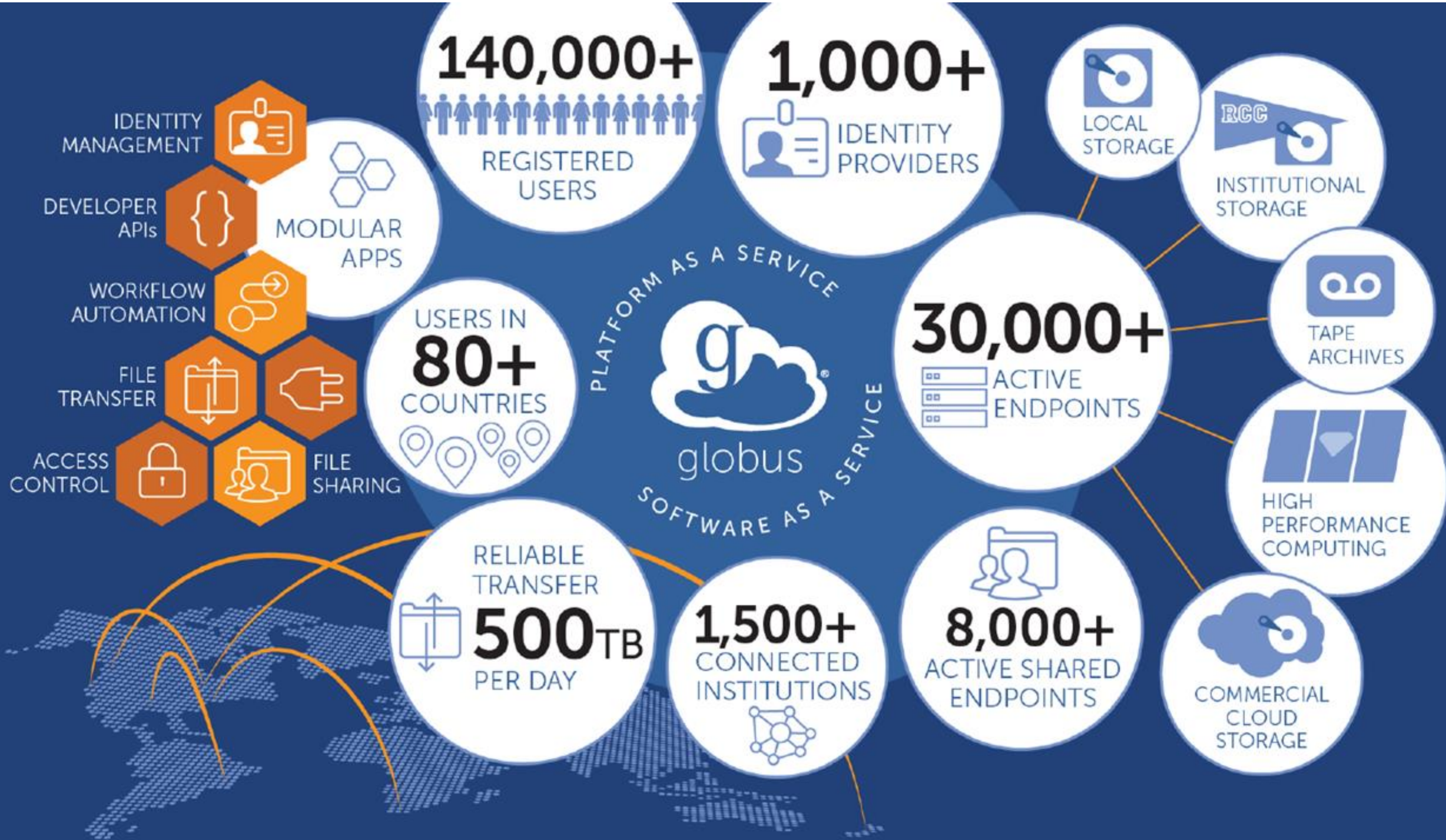


Globus is ...

a non-profit service
developed and operated by



THE UNIVERSITY OF
CHICAGO



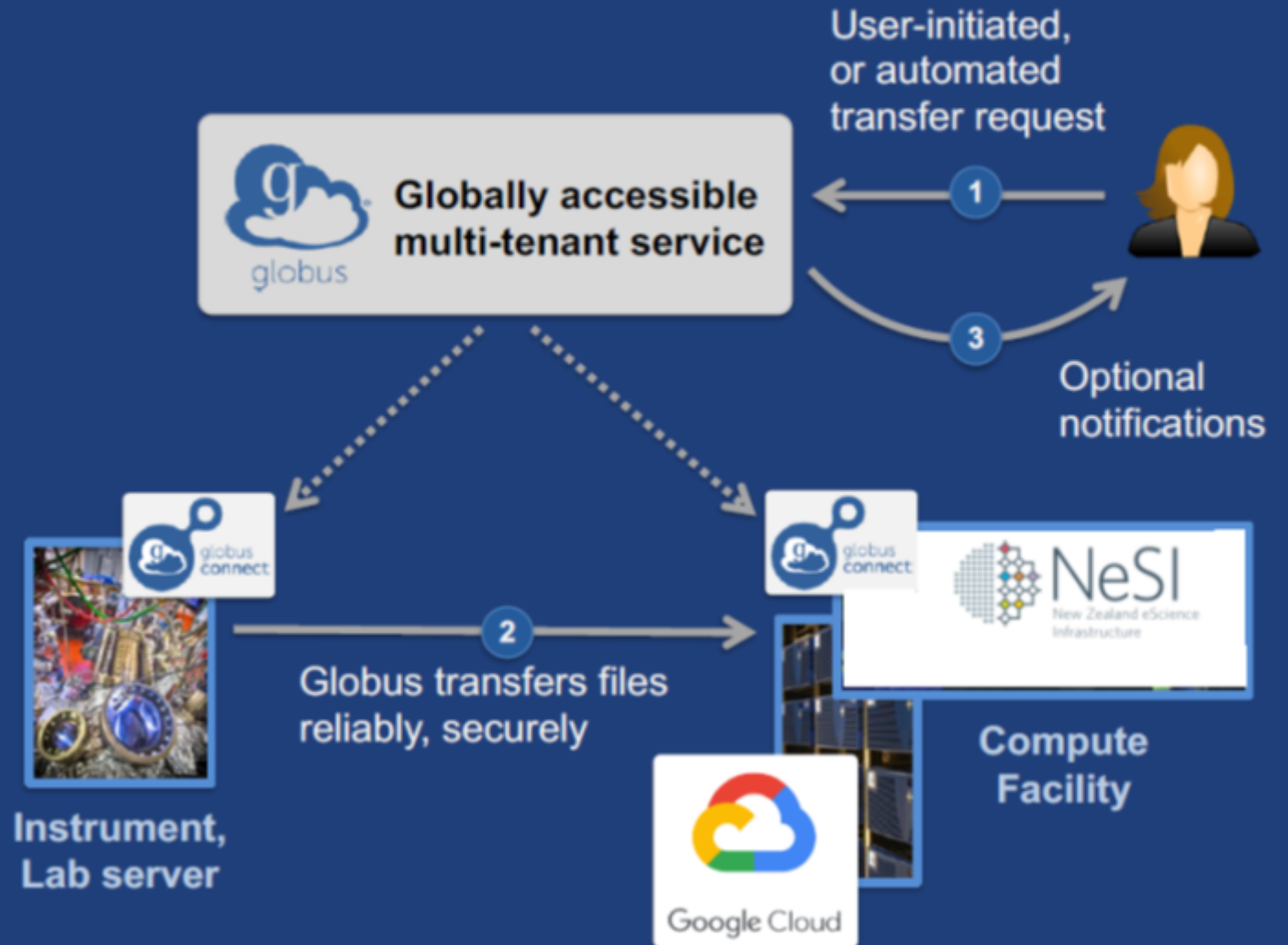


1,000,000,000,000,000,000
1 EXABYTE – A QUINTILLION BYTES TRANSFERRED BY GLOBUS



Fast, reliable file transfer ...from any to any system

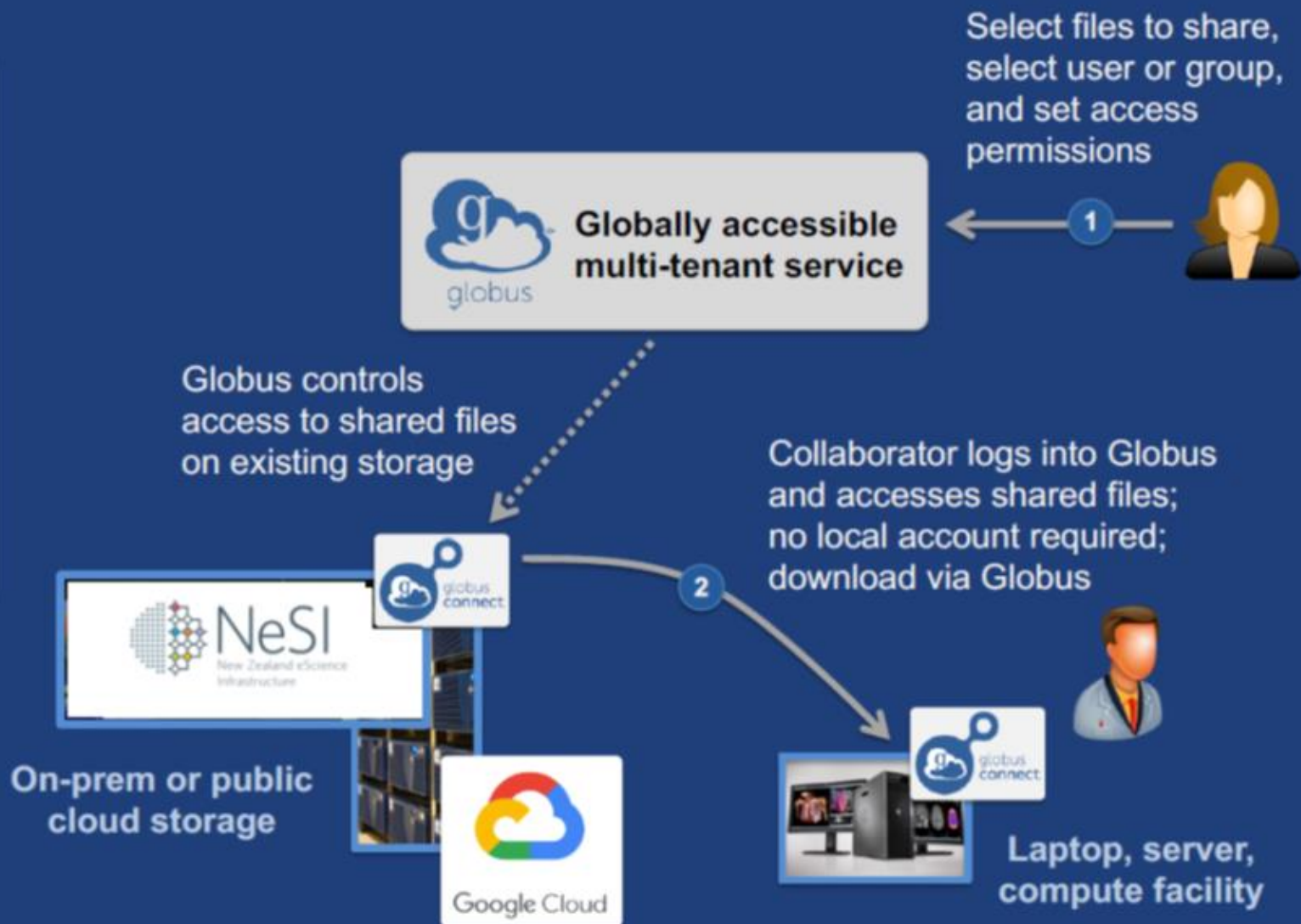
- Fire-and-forget transfers
- Optimized speed
- Assured reliability
- Unified view of storage
- Browser, REST API, CLI





Secure data sharing ...from any storage

- Fine-grained access control “overlay” on storage system
- Share with any identity, email, group
- No need to stage data just for sharing





Globus Connectors



Google Cloud



Caringo Swarm™

Quantum

ActiveScale
Object
Storage



IBM Spectrum Scale



Google Drive



SCALITY

IBM Cloud
Object Storage



ceph

HPSS



wasabi™
hot cloud storage

lustre™

Planned



OneDrive

Microsoft Azure
Blob Storage



Dropbox

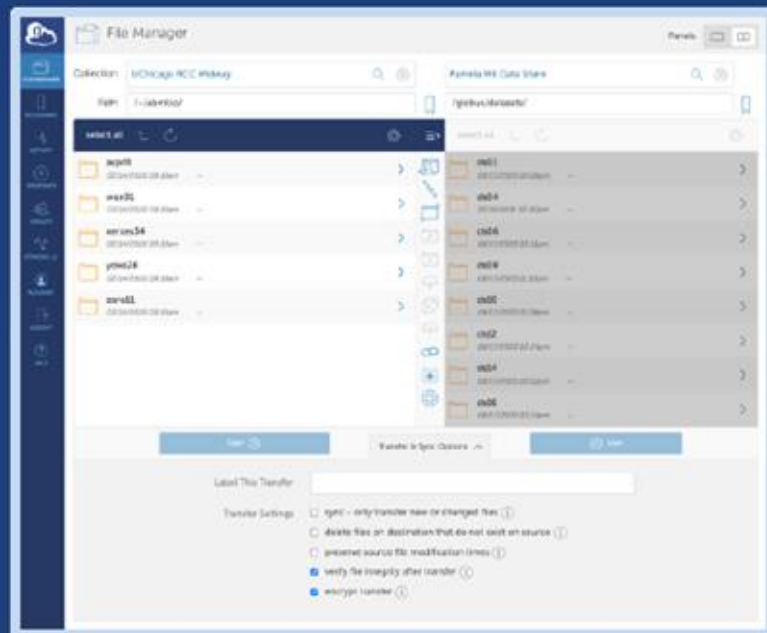


Use(r)-appropriate interfaces

Globus
service



Web



Platform
(RESTful APIs)

```
GET /endpoint/go%23ep1
PUT /endpoint/demodoc#my_endpt
200 OK
X-Transfer-API-Version: 0.10
Content-Type: application/json
...
```

CLI

```
Usage: globus [OPTIONS] COMMAND [ARGS]...

Options:
  -v, --verbose          Control level of output
  -h, --help             Show this message and exit.
  -F, --format [unix|json|text] Output format for stdout. Defaults to text
  --jmespath, --jq TEXT  A JMESPath expression to apply to json
                        output. Takes precedence over any specified '
                        --format' and forces the format to be json
                        processed by this expression
  --map-http-status TEXT Map HTTP statuses to any of these exit codes:
                        0,1,50-99. e.g. "404=50,403=51"

Commands:
  bookmark      Manage endpoint bookmarks
  config        Manage your Globus config file. (Advanced Users)
  delete        Submit a delete task (asynchronous)
  endpoint      Manage Globus endpoint definitions
  get-identities Lookup Globus Auth Identities
  list-commands List all CLI Commands
  login         Log into Globus to get credentials for the Globus CLI
  logout        Logout of the Globus CLI
  ls            List endpoint directory contents
  mkdir         Make a directory on an endpoint
  rename        Rename a file or directory on an endpoint
  rm            Delete a single path; wait for it to complete
  session       Manage your CLI auth session
  task          Manage asynchronous tasks
  transfer      Submit a transfer task (asynchronous)
  update        Update the Globus CLI to its latest version
  version       Show the version and exit
  whoami        Show the currently logged-in primary identity.
```

transfer.nesi.org.nz

The screenshot displays the NeSI File Manager interface. The left sidebar contains the NeSI logo, a 'File Manager' section with a 'RECENTLY USED' list (Globus Usage Reports, Flaherty Laptop, NeSI Wellington DTN, Genomics Aotearoa D...), a 'BOOKMARKS' section (Flaherty Laptop, Genomics Aotearoa D..., Globus Usage Reports, /nesi/nobackup/nesi99...), and a navigation menu (Activity, Endpoints, Groups, Console, Account: brianflaherty@globusid..., Logout, Help, NeSI Home).

The main area is divided into two panels. The left panel, titled 'File Manager', shows the 'Collection' as 'NeSI Wellington DTN' and the 'Path' as '/nesi/share/ga/kakapo/'. It lists several folders and files:

Item	Size	Modified
analysis	-	10/06/2019 03...
bioinf_report_LaTeX	-	10/06/2019 03...
config	51 B	10/04/2019 01...
credentials	119 B	10/04/2019 01...
other_genomes	-	10/06/2019 03...
raw	-	10/06/2019 02...
reference	-	10/06/2019 02...
scripts	-	10/06/2019 02...

The right panel, titled 'Genomics Aotearoa Data Repository Otago v1', shows the 'Path' as '/projects/'. It lists three folders:

Item	Size	Modified
Chrysophrys_auratus	-	04/05/2019 1...
MetagenomicsBenchmarkData	-	07/26/2019 0...
SnapperRNAseq	-	04/04/2019 0...

At the bottom, there are 'Start' buttons for both panels and a 'Transfer & Sync Options' dropdown menu.



Automation Examples

- **Syncing a directory**
 - bash script; calls the Globus CLI
 - Python module; run as script or import as module
- **Staging data for distribution**
 - bash and Python variants
- **Removing directories after files are transferred**
 - Python script

github.com/globus/automation-examples



Globus Automate

A platform service for defining, applying, and sharing distributed research automation **flows**

- **Triggers** start flows based on subscribed events
- Flows call **Action Providers** to perform tasks





Automation Action Providers

Transfer



Delete



ACLs



funcX



DLHub



Identifier



User Form



Notification



Xtract



Web Form



Ingest



Expression
Evaluation



Search



Describe



Globus action
providers

Custom action
providers

NATIONAL RESEARCH AND EDUCATION NETWORK

Research Education Advanced Network NZ



NeSI REANZ

New Zealand eScience
Infrastructure



REANNZ

Research and Education,
Advanced Network New
Zealand

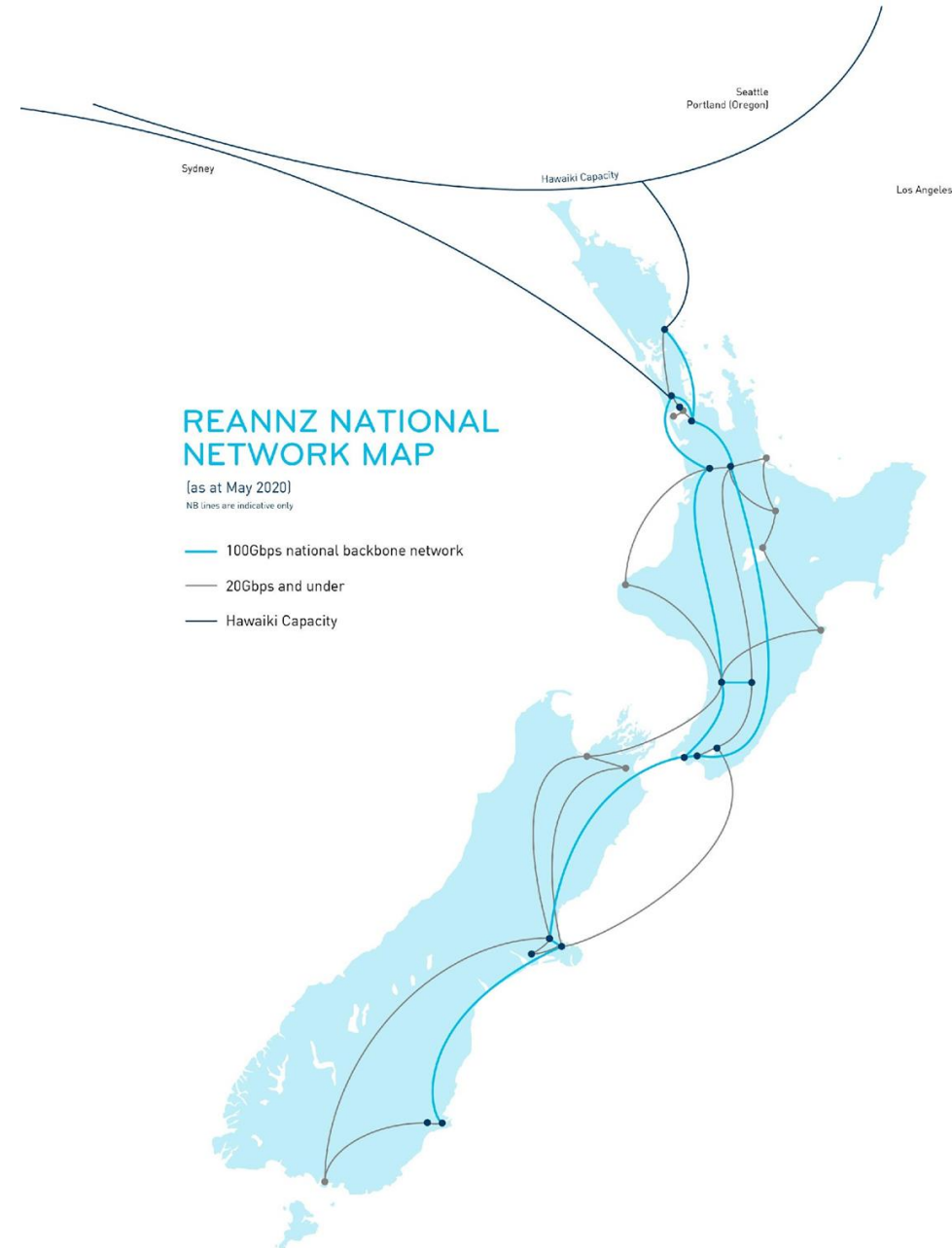
NREN

National, Research and
Education Network

120 global NRENs

Supporting and enabling research outcomes

Crown Owned entity, not for profit



Enabling research in NZ

Operating the network and underlying infrastructure

Core services

domestic network (26 POPs)

International network (5 POPs)

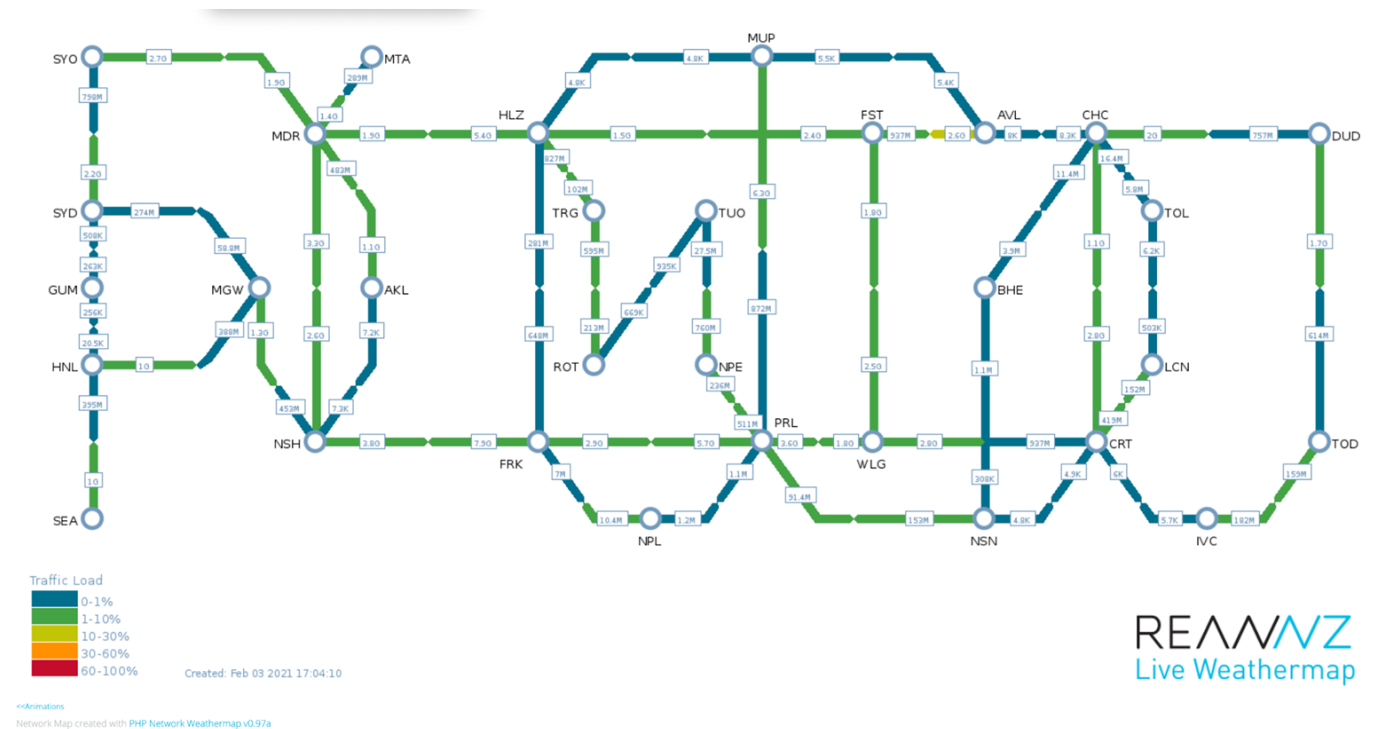
Connectivity solutions

(managed access/edge, manage firewall, science DMZ, cloud connect, eduroam, tuakiri (trust ident), professional services)

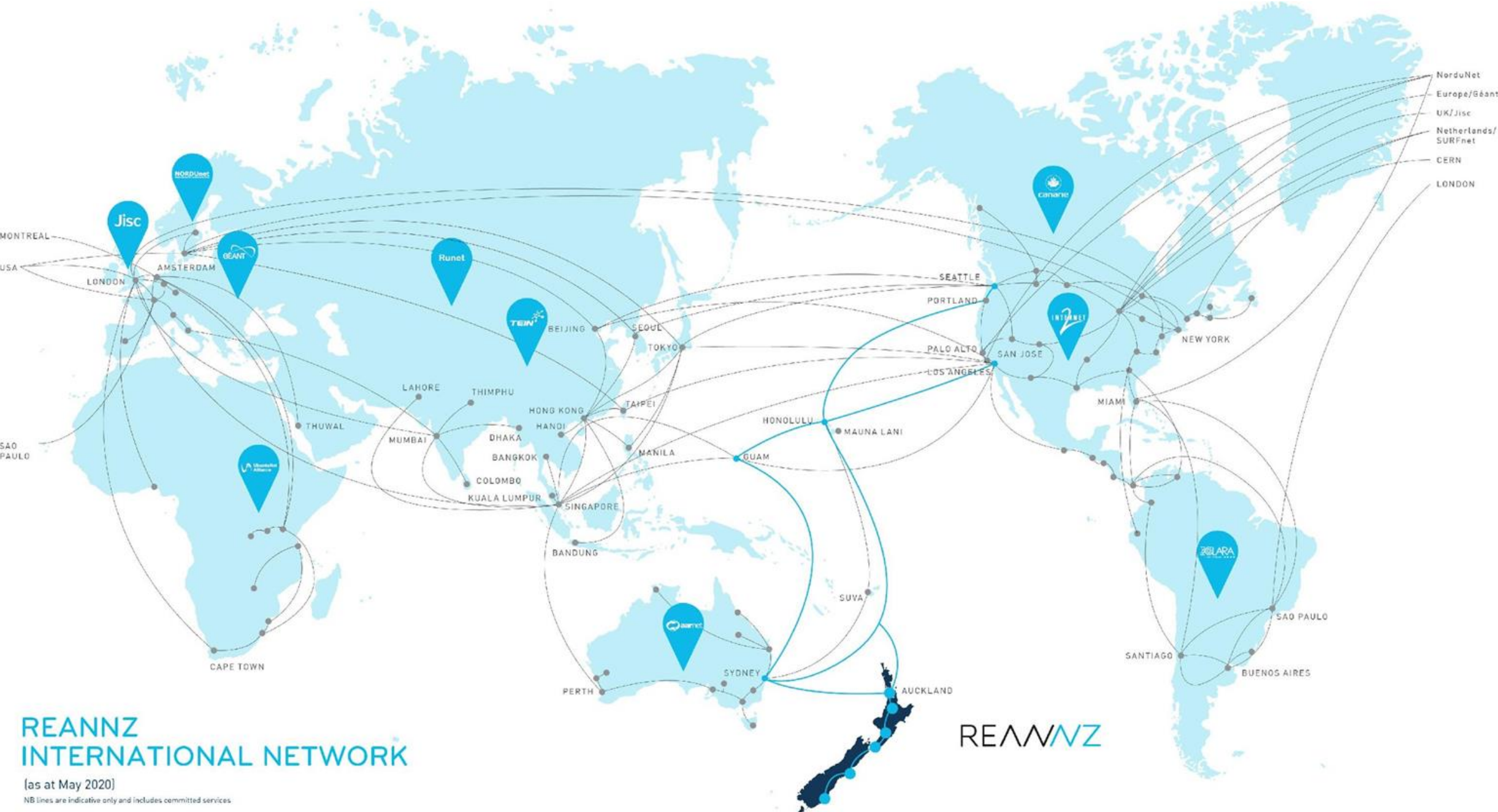
On call 24/7

end to end visibility

Partnership with member technology teams



Global National Research and Education Networks



Global National Research and Education Networks



NRENs operate nationally, but connect globally
Seamless connectivity and tailored services
Sharing information and collaboration, driving research

National R&E networks leverage the global community to support the bespoke and demanding needs of the research and education sector.

By sharing information and developing collaborations of best practice, the global R&E network community creates efficiencies and avoids reinventing the wheel.

In many countries, R&E networks provide connectivity for universities and institutions and organisations with a research and education mission at a cost and capacity that commercial networks cannot.

Network Challenges for a Data Transfer Service

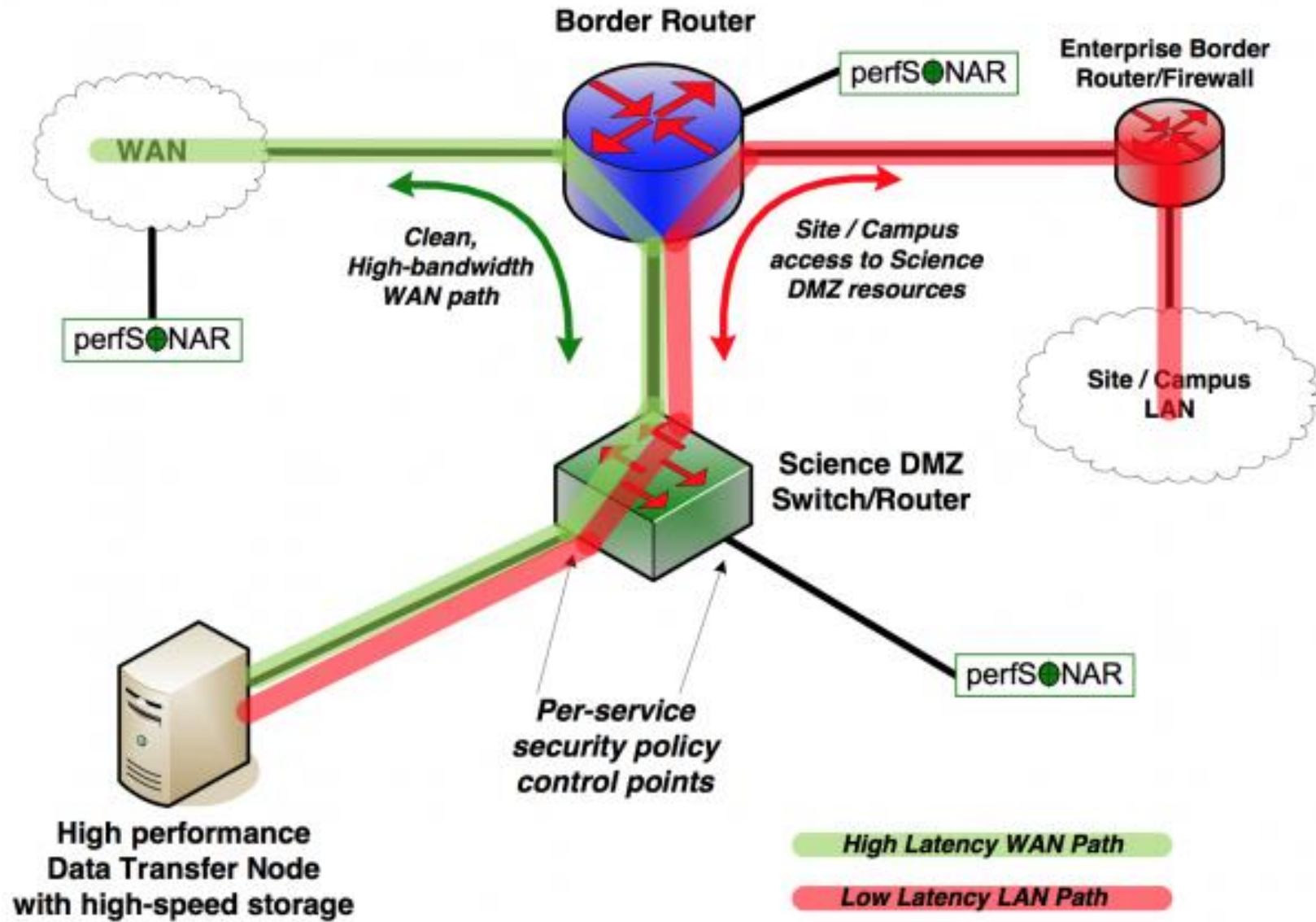
- Think of a network 'pipe' – how to get the most data flowing through it per second, every second.
- Make sure your network 'pipes' are large – everywhere and at every junction.
- Make sure you have enough data ready to 'pour' into the network pipe – so you need fast multiplexed data storage at each end.
- How do you get data into the network pipe fast enough – use multiple high speed data 'pumps' - each uses many large 'data buckets'. They all pump at once.
- Keep the network data pipes below capacity and minimise the data flow controls.

What is a Science DMZ?

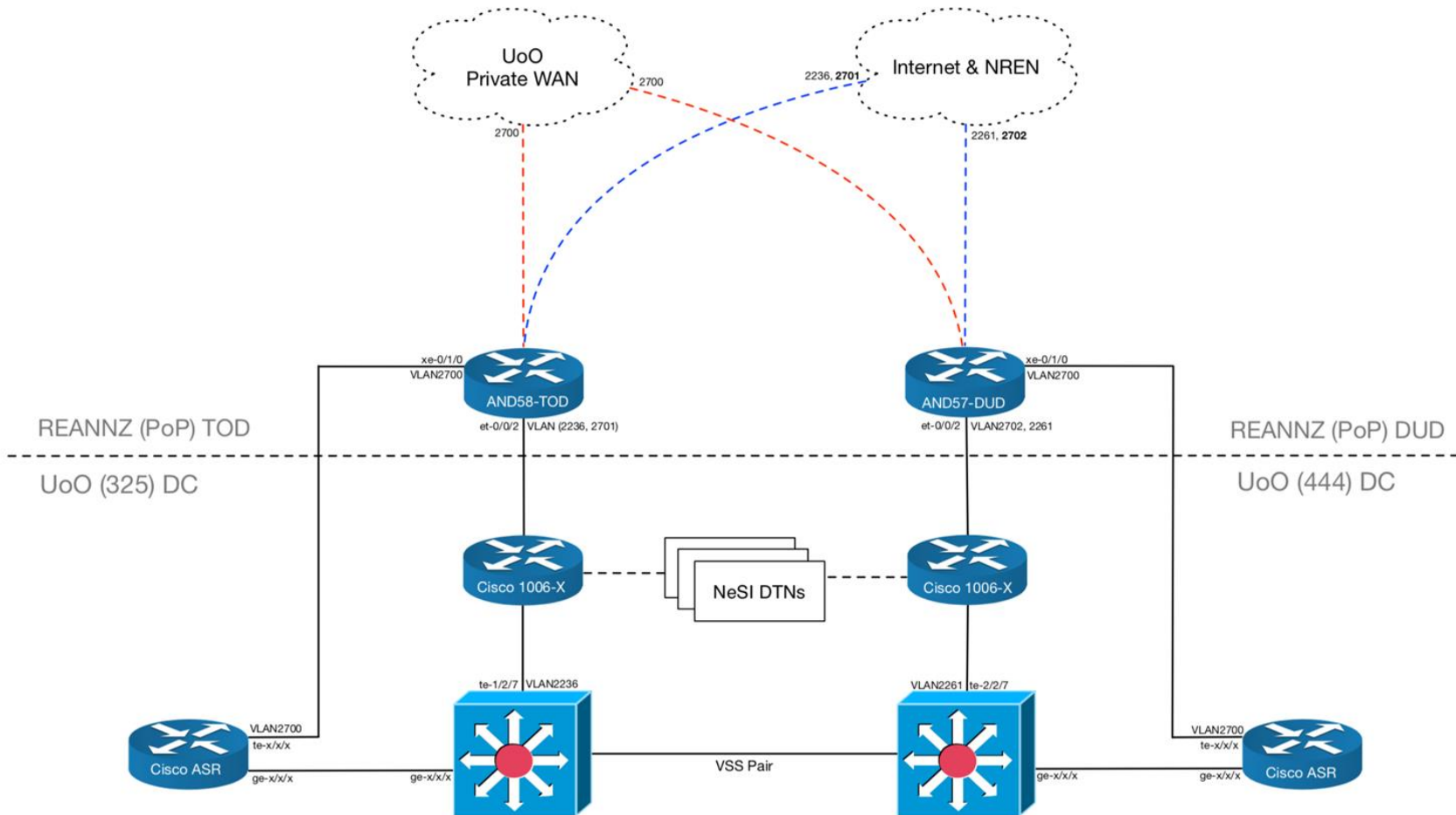
- A Science DMZ is a high-performing on-ramp at or near the site network perimeter, dedicated to supporting data-intensive science resources.
- The Science DMZ is a portion of the network, built at or near the campus or laboratory's local network perimeter. It is designed so that the equipment, configuration, and security policies are optimised for high-performance scientific applications, rather than for general-purpose business systems or “enterprise” traffic.
- Can be used as dedicated systems for data transfer (DTNs)
- Integrated performance management and security
- Built with high performance components
- Equipment, configuration and policies are optimized for high performance scientific applications

How does it work?

- Addresses common network performance problems encountered at research institution
- Science applications will go through clean high bandwidth path
- Equipment, configuration and policies are optimized for high performance scientific applications
- Managed edge switch, so that everyday internet traffic could be directed down a normal path protected by enterprise firewalls etc. while research data would go down a different, faster path, which bypassed all of the gates it didn't need to go through.
- Creating a clear, high-speed research data path to and from the REANNZ network.
- REANNZ can manage the switch on behalf of the member, so if changes or upgrades need to be made in the future REANNZ can make them, meaning less work for the member.
- Talk about the following example of a Science DMZ topology



University of Otago – NeSI Science DMZ



Thank you!

Questions?

Brian Flaherty

Richard Tumaliuan



NeSI REANZ

New Zealand eScience
Infrastructure

