# A reservation system in the Nectar Research Cloud for GPU and large memory instances
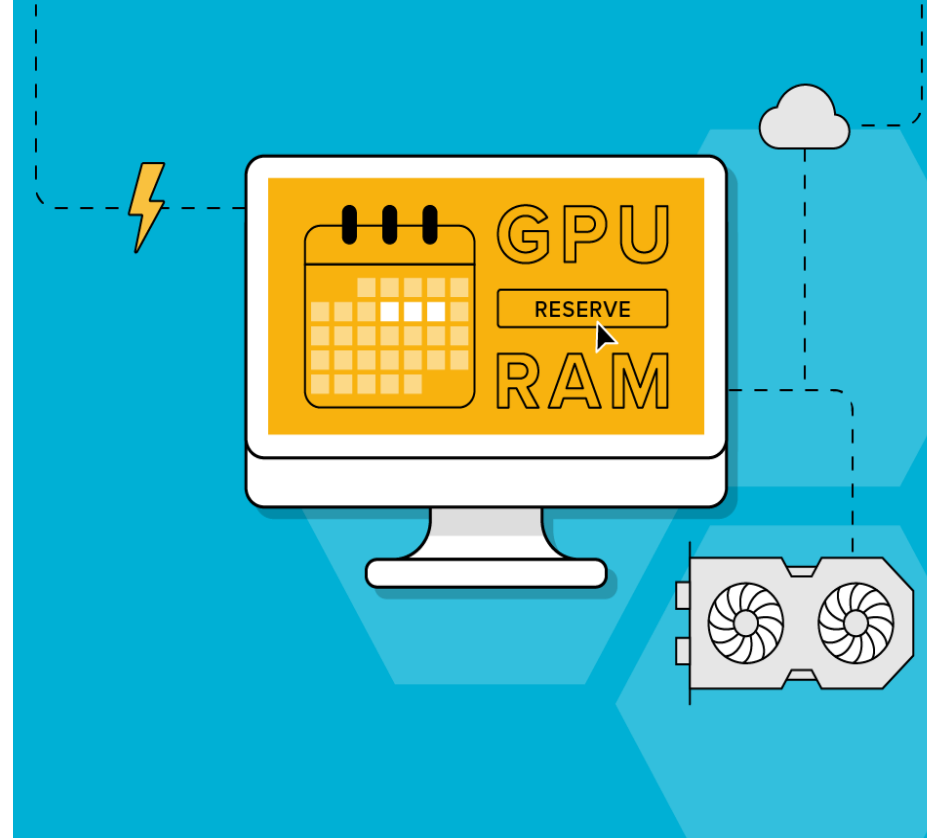
**20 October 2022**
**14.05 - 14.25**

**PRESENTED BY**
Paul Coddington

Australian Research Data Commons

**ARDC** Australian Research Data Commons

**NCRIS** National Research Infrastructure for Australia — An Australian Government Initiative
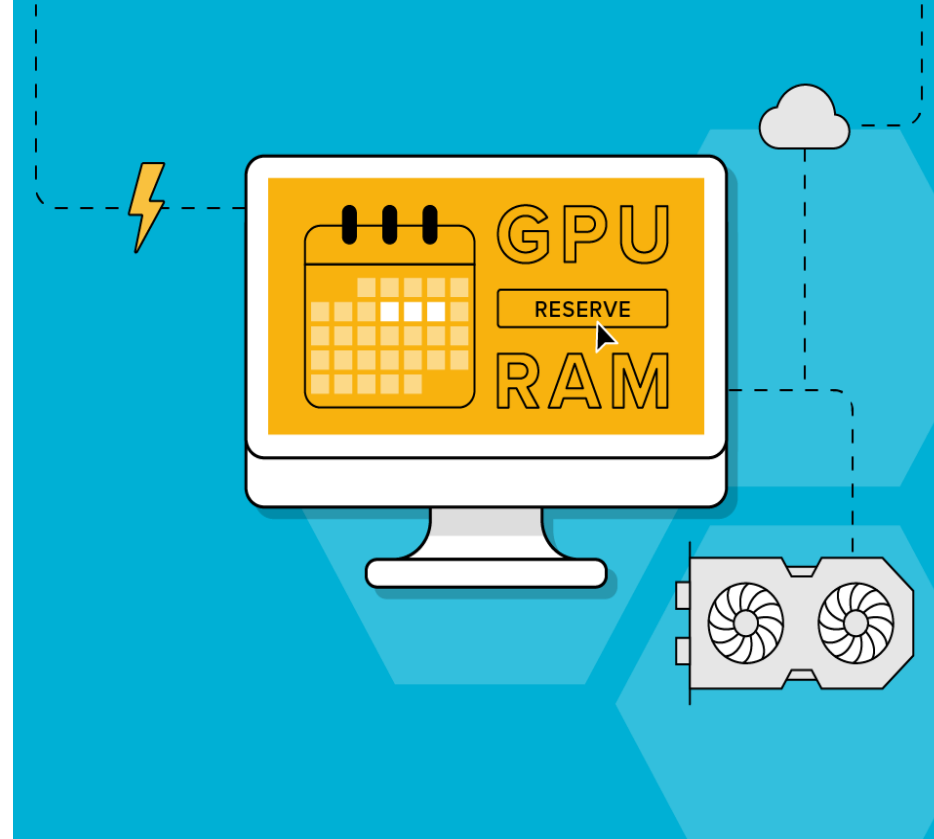
ARDC is enabled by NCRIS

# NEW SERVICE - why?

- Growing demand for high-end compute infrastructure in the research sector
  - Seen in ARDC Platforms
  - Seen across institutions with demand for GPUs for ML, image processing, simulation
  - Requirement for large memory machines to handle big data sets and large scale analysis and simulation
- Not limited to one type of research discipline

# NEW SERVICE - why?

- GPU and large memory servers in the Nectar Research Cloud were
  - dedicated to Platforms projects or
  - for local Node use
- Want to provide a national IaaS for projects meeting national merit criteria
- Nectar infrastructure is provided at no cost to researchers - but GPU and large memory servers are very expensive
- Need to ensure high utilisation of these expensive and limited resources

- **Solution - virtualisation and a reservation system**

# Steps required to move to new service

1. Move from limit-based allocation quotas to usage based quotas

2. Design and develop a reservation system

3. Define and develop standard flavors for GPU and large memory virtual machines

4. Provide new GPU and large memory hardware at Nodes to underpin the service

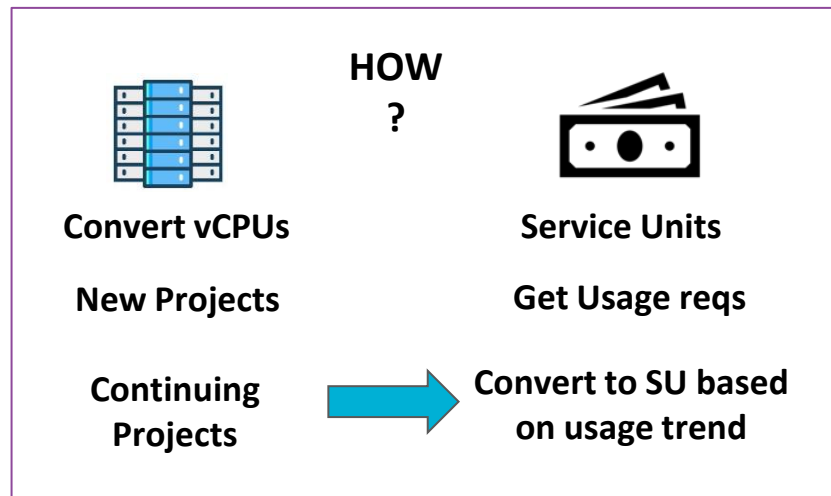5. Design and test with a Pilot Phase

6. Launch service

# 1. Move to Usage Based Allocations - Service Units

## OVERVIEW

➢ Move from maximum capacity limit to credit/budget of service units (SUs) for the period of the project allocation
➢ This will then account for actual usage of resources where each resource has a specific cost
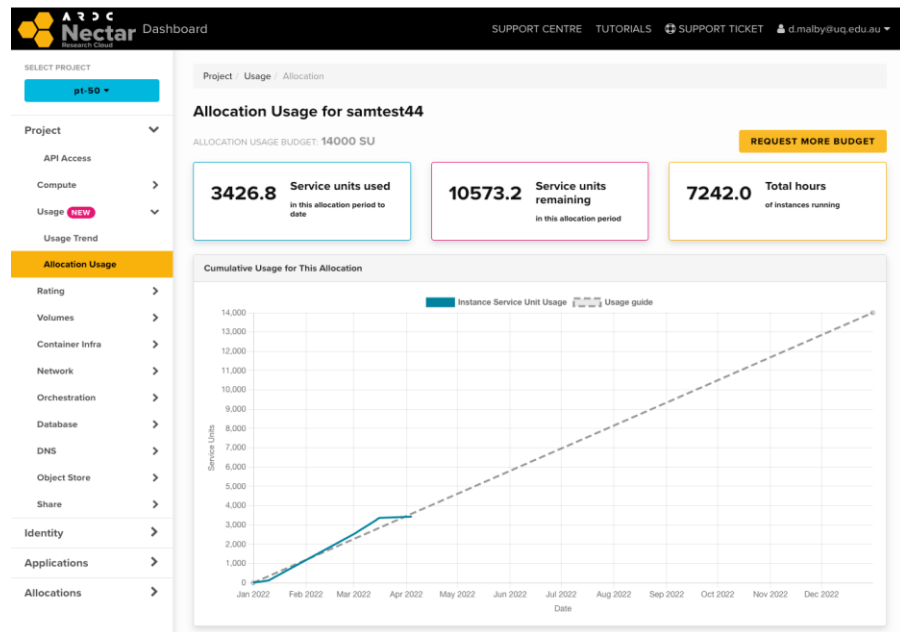➢ Service developed for generic cloud allocations (all users)

## WHY?

• *Maximum capacity allocation limits flexibility*
• *Adapt to bursty data analysis workflows*
• *Service designed for diversity of compute workflows*
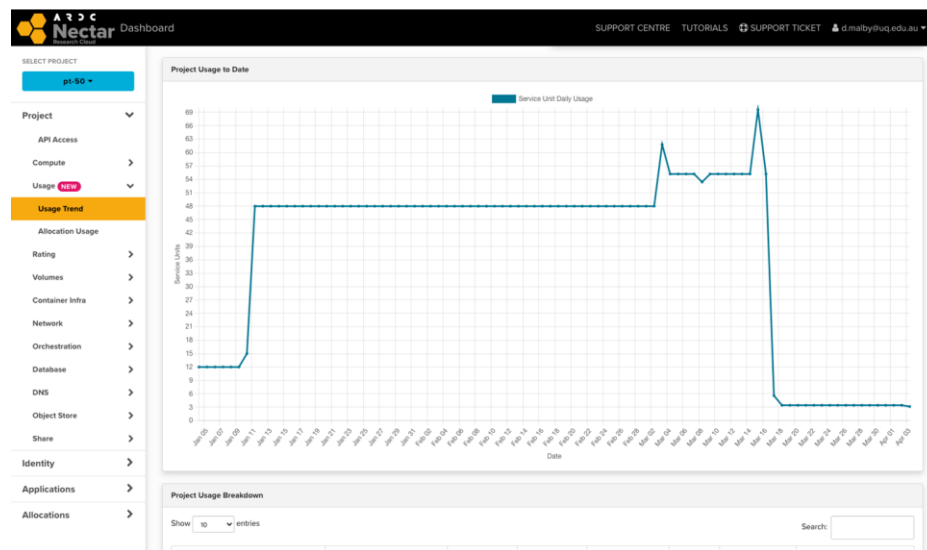• *More efficient use of available capacity*

**HOW ?**

Convert vCPUs          Service Units

New Projects          Get Usage reqs

Continuing Projects → Convert to SU based on usage trend

# Change to Dashboard User Interface (UI)

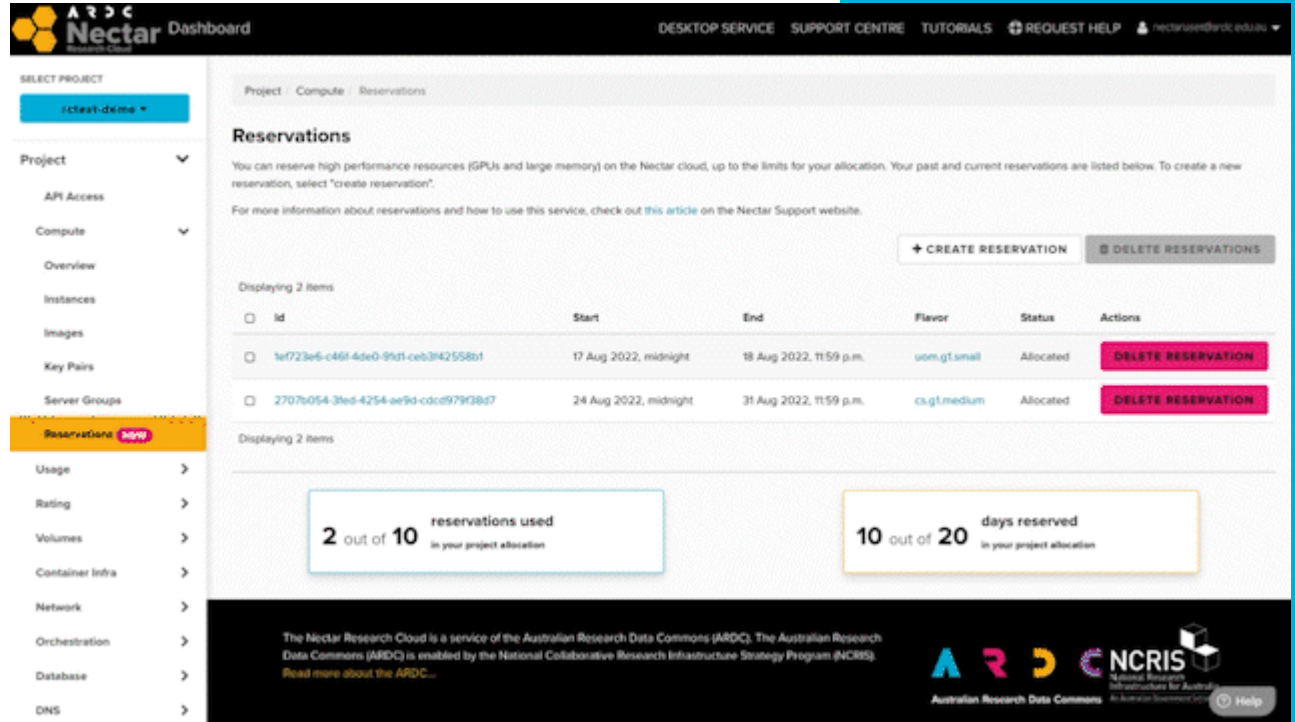## Allocation Usage

## Usage Trend

# 2. Design & Implement a Reservation System

Built on the OpenStack Blazar reservation service.

Users can reserve access to specialised high-end computing power in the Dashboard.

\* Allocations must first be approved for reservations.

# SERVICE UNITS - in the Reservation System



☑ **Enough SU budget is available.**

☒ **Not enough SU budget is available.**

# 3. Defining Flavors

The following standard flavor classes are offered on the Nectar Research Cloud:

- Tiny (t3)
- Balanced (m3)
- RAM Optimised (r3)
- CPU Optimised (c3)
- Preemptible (p3)
- **GPU Visualisation (g1) - A40**
- **GPU Compute (g2) - A100-80 sliced up 1/2 to 1/10**
- **Huge RAM (h4) - up to 128 vcpus and 960GB RAM**
- Huge RAM (h3)

Detailed list and recommended uses of Flavors available here:
[https://support.ehelp.edu.au/support/solutions/articles/6000205341-nectar-flavors](https://support.ehelp.edu.au/support/solutions/articles/6000205341-nectar-flavors)

| Name | VCPUS | RAM | Root Disk | Ephemeral Disk | Public | SU/hour | |
|---|---|---|---|---|---|---|---|
| > t3.xsmall | 1 | 1 GB | 10 GB | 0 GB | Yes | 0.014 | ↑ |
| > ⚠ p3.xsmall | 1 | 2 GB | 30 GB | 0 GB | Yes | 0.007 | ↑ |
| > t3.small | 2 | 2 GB | 10 GB | 0 GB | Yes | 0.029 | ↑ |
| > c3.xsmall | 1 | 2 GB | 30 GB | 0 GB | No | 0.043 | ↑ |
| > m3.xsmall | 1 | | | | | | |
| > r3.xsmall | 1 | | | | | | |
| > m3.small | 2 | | | | | | |
| > t3.medium | 4 | | | | | | |
| > ⚠ p3.small | 2 | | | | | | |
| > c3.small | 2 | | | | | | |

**Select Flavor Availabilty**

Flavor Class [ ALL ] [ HUGE RAM ] [ GPU ]     Availability Zone [ All ]

**NOTE:** All times are displayed in **UTC** time.

**October/2022**

01 02 03 04 05 06 07 08 09 10 11 12 **13** 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 01 02 03 04 05 06 07

- mon.g2.small ▾
- qld.h4.large ▾
- qld.h4.small ▾
- tas.g2.xlarge ▾
- tas.h4.large ▾
- tas.h4.xlarge ▾

# 4. Provision the hardware and licenses

- RFP to Nectar nodes asking for proposals for hardware based on their researcher requirements
- ARDC provided capex, Nodes provide opex as co-investment
- 16 GPU servers and 7 large memory servers from ARDC investment
- Additional servers from Node investment for use by Node members
- GPUs are a mix of
  - A100-80 GPUs (mainly for compute)
  - A40 GPUs (mainly for image processing and visualisation)
- All have large NVMe drives for fast local disk storage
- A pool of Nvidia licenses for virtualisation - VCS for A100 and VWS for A40

# 5. Pilot and Testing Snapshot

**Design and Development**

**Pilot (limited user groups)**

**Post production monitoring**

**Final Testing for Pilot**

**Prod Release 1**

Pilot environment for user testing released on 2nd Aug 2022

24 users across ~14 projects

2 GPU and 4 large RAM servers running with different Flavors

44 GPUs in total to be deployed

~248 vGPUs available for researchers

# NATIONAL GPU SERVICE: **Benefits Realised**

| **Resource Allocation** | **Virtualisation** | **Unique flavors** |
|---|---|---|

- Fair allocation of resources that are limited & expensive
- Ensures that the limited resources are reserved and released when required
- Reserved access for training courses

- Virtualisation enables improved GPU utilisation and more users
- First virtualised GPUs service on a national scale

- A variety of flavors designed for research
- Provides access to large GPU flavours not yet available on any cloud platform in Australia
- GPU servers can be reconfigured to adapt to usage trends for different sized flavors

# By end of 2023

**248**
vGPUs

## Participating Nodes

**University of Tasmania**
Hobart

**Monash University**
Melbourne

**QCIF**
Brisbane

**Intersect**
Sydney

**Swinburne University of Technology**
Melbourne

# 6. Launch and expand

- Production environment released as 'BETA' on 13th September 2022
- Uplift in capacity as more infrastructure comes on line
- Available to all national merit research projects
- Capacity will be added as infrastructure comes on line
- Review and revise the mix of flavors, limits, etc

Learn more!
Webinar October 25th
2-3pm AEST

ARDC
Australian Research Data Commons

NCRIS
National Research Infrastructure for Australia
An Australian Government Initiative

ARDC is enabled by NCRIS

# ACKNOWLEDGEMENTS

- Shubhra Dargar - project management
- Sam Morrison - reservation system
- Darcelle Maltby - web dashboard interface
- Sengor Kusturica, Rocky Yan, Andy Botting, Dylan McCulloch - GPU virtualisation, licensing, standard flavors
- Jo Morris and Sonia Ramza - user guides, user support, promotion
- ARDC comms team - communications and promotion
- And technical assistance from many Node operations staff

# HOW CAN THE ARDC ACCELERATE YOUR RESEARCH?

*Visit us at*
**eResearch Australasia - Stand 14**

# THANK YOU

🌐 ardc.edu.au

✉️ contact@ardc.edu.au

📞 +61 3 9902 0585

🐦 @ARDC_AU

in Australian-Research-Data-Commons

Subscribe to the
**ARDC CONNECT**
newsletter

**About the ARDC**