# Establishing Gen3 to enable better human genome data sharing in Australia

Welcome! The webinar will commence at 1pm AEDT /12pm AEST/ 11:30am ACDT/ 10am AWST

Australian
**BioCommons**

Actively supporting Australian life sciences research through bioinformatics and bioscience data infrastructure

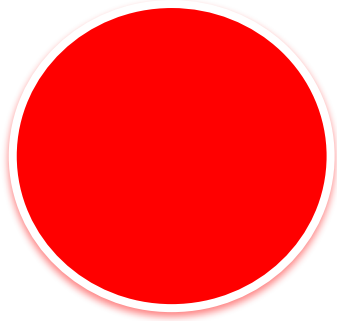biocommons.org.au     AustralianBioCommons     @AusBiocommons

# Acknowledgement of Country

We acknowledge the Traditional Owners and their custodianship of the lands on which we meet today.

We pay our respects to their Ancestors and their descendants, who continue cultural and spiritual connections to Country.

We recognise their valuable contributions to Australian and global society.

Australian
**BioCommons**

# Housekeeping

Session is recorded

Autogenerated captions available

Questions via Q&A function

Australian
**BioCommons**

# Establishing Gen3 to enable better human genome data sharing in Australia

SPEAKERS

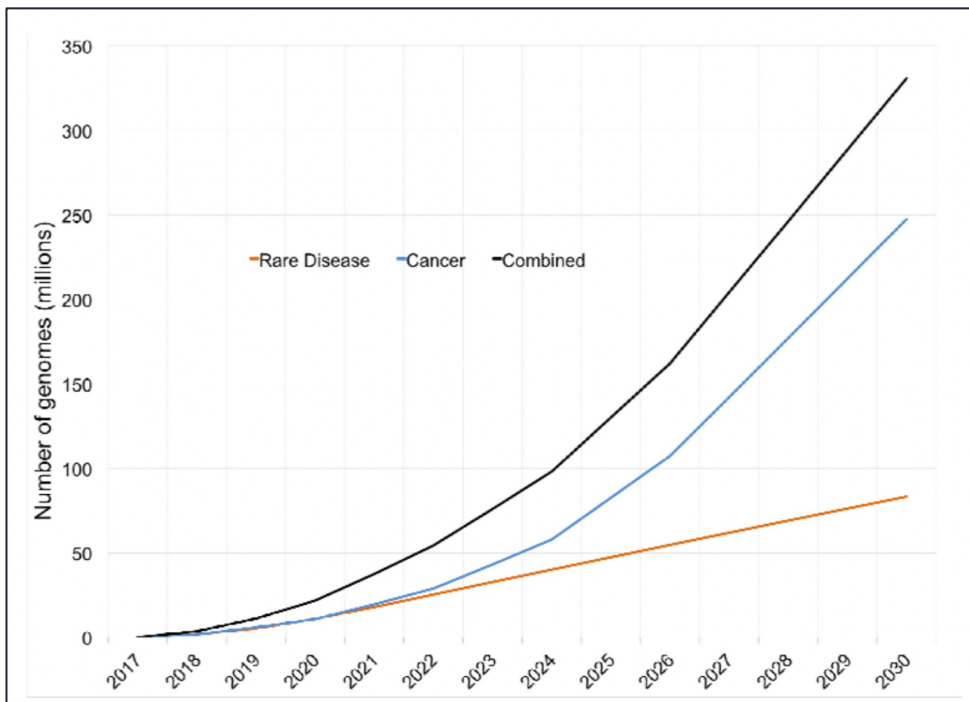Associate Professor Bernie Pope, Australian BioCommons / Melbourne Bioinformatics
Professor Oliver Hofmann, University of Melbourne Center for Cancer Research
Mr Kamile Taouk, Children's Cancer Institute
Dr Marie Wong-Erasmus, Children's Cancer Institute

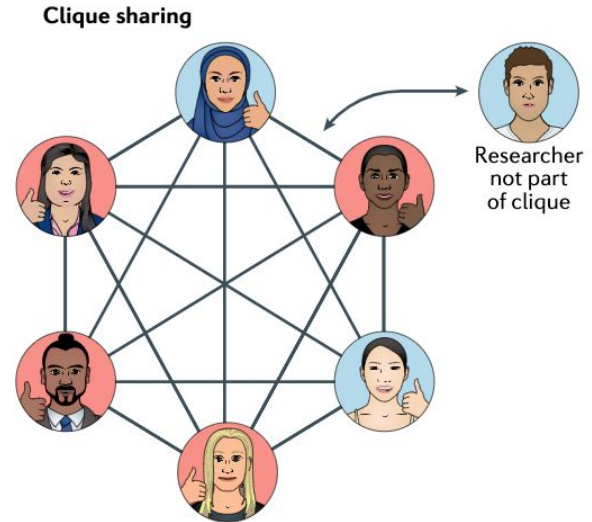# Predicted global growth of healthcare funded sequenced human genomes



Global storage requirements in 2025 to be exabytes to low zettabytes.

Birney, E., Vamathevan, J., and Goodhand, P. (2017). Genomics in healthcare: GA4GH looks to 2022. bioRxiv

# Siloed data

- Human genomics data has often been siloed.

- This limits reuse and reanalysis.

- Public benefit is increased when data is shared.

- Sharing is frequently necessary in human health, especially in rare disease and cancer.

- Large cohorts are needed for statistical power.

- National and international data sharing is highly beneficial but requires considerable collaboration and coordination.



Byrd, J. B., Greene, A. C., Prasad, D. V., Jiang, X. & Greene, C. S. Responsible, practical genomic data sharing that accelerates research. Nat. Rev. Genet. 21, 615–629 (2020).

Australian BioCommons

# Towards global standards for data sharing

- The Global Alliance for Genomics and Health (GA4GH) is a policy-framing and technical standards-setting organization, seeking to enable responsible genomic data sharing within a human rights framework.

- Australian Genomics is a driver project of GA4GH.

- A key outcome is the specification for standard APIs for data sharing technology.

- Recognition that the data life cycle in human genomics is complex and data storage and analysis are parts of a bigger ecosystem.



**GA4GH Standards in the Data Life Cycle**

Rehm, H. L. *et al.* GA4GH: International policies and standards for data sharing across genomic research and healthcare. *Cell Genomics* 1, (2021).

Australian BioCommons

# Infrastructure ecosystem



Example solutions:

- IAM: CILogon, GA4GH passports

- Data commons: Gen3

- DAC approval: REMS

- Analytics: national infrastructure, institutional infrastructure, commercial cloud

- Integrated: Broad Terra + DUOS

Marion Shadbolt

# Establishing Gen3 to enable better Human Genome Data sharing in Australia

- Gen3 was identified as a leading candidate for building a human genomics data commons.

- In Q3 2021 we established a pilot project to assess the use of Gen3 as the foundation for a human genomics data commons.

- That project is now complete, and today we provide an overview of the motivations, process, and findings.

# Zero Childhood Cancer - Australia's national paediatric cancer precision medicine program

ZERO2: by 2023
all children in Australia



> 650 high-risk patients recruited

# Scale up and share

The **Zero Childhood Cancer** program involves all paediatric oncology units across Australia, these hospitals will work with key medical research institutes, both nationally and internationally.



Curie Institute, France

Princess Maxima Centre, Netherlands

DKFZ (German Cancer Research Center), Germany

University of Cambridge (Cancer Research UK Cambridge Institute), UK

Institute for Molecular Medicine Finland, University of Helsinki, Finland

Children's Hospital of Philadelphia, USA

St. Jude's Children's Research Hospital, USA

Terry Fox PROFYLE, The Hospital for Sick Children (SickKids), Canada

Queensland Children's Hospital, Brisbane

The University of Queensland Diamantina Institute

Perth Children's Hospital, Perth

Centre for Childhood Cancer Research, Telethon Kids Institute

Women and Children's Hospital, Adelaide

South Australian Health and Medical Research Institute

Centre for Cancer Biology, Adelaide

The Children's Hospital at Westmead

Kid's Research Institute

Children's Medical Research Institute

John Hunter Children's Hospital, Newcastle

**The Kid's Cancer Centre at Sydney Children's Hospital, Randwick**

**Children's Cancer Institute**

Hereditary Cancer Centre and NSW Health Pathology, Prince of Wales Hospital

Kinghorn Centre for Clinical Genomics, Garvan Institute of Medical Research

Centre for Economic Impacts of Genomic Medicine, Macquarie University

Australian Institute of Health Innovation, Macquarie University

Royal Children's Hospital, Melbourne

Peter MacCallum Cancer Centre

Walter and Eliza Hall Institute of Medical Research

Monash Children's Hospital, Clayton

Murdoch Children's Research Institute

Hudson Institute of Medical Research

- Children's hospitals participating in Zero Childhood Cancer National Clinical Trial, sponsored by ANZCHOG
- Research Institutes participating in Zero Childhood Cancer
- Zero Childhood Cancer is a joint initiative of Children's Cancer Institute, the Kids Cancer Centre at Sydney Children's Hospital, Randwick and The Sydney Children's Hospitals Network

With only **~1000 new cases** of high-risk paediatric cancer per year, it is imperative to **aggregate** Australian data with global data to develop strategies to **effectively treat high-risk childhood cancer**

**We need to share, analyse, integrate data more easily**

**image credit:** https://www.zerochildhoodcancer.org.au/about/research---clinical-partners

Australian
BioCommons

# Paediatric data sources

# CCI - where's our data?



**NetApp**

StorageGRID: smart, fast, and future proof object storage

Data commons?

**GEN3** DATA COMMONS

**SevenBridges** Announces International Collaboration Focused on Personalized Treatment for Kids with Cancer

**Expansion of the CAVATICA Platform to Australia Enables Harmonized Analysis of Geographically Separated and Jurisdictionally Protected Data Resources**

BOSTON, June 2, 2020—Seven Bridges, the industry-leading bioinformatics ecosystem provider, today announced a collaborative partnership between The Gabriella Miller Kids First Data Resource Center (Kids First DRC), ZERO Childhood Cancer (ZERO), the Children's Brain Tumor Tissue Consortium (CBTTC), the Australian BioCommons and the Australian Research Data Commons (ARDC). The multinational genomic

CAVATICA

Australian
**BioCommons**

# How to share?

We Need:
- a way to SEARCH different catalogues of paediatric data
- a way to easily gain and grant ACCESS to the data
- a way to ANALYSE this data in place if possible

# University of Melbourne Centre for Cancer Research

Precision Oncology Program

Sean Grimmond

# Precision Oncology Program

Recalcitrant Cancers, Rare Cancers, Cancers of the Unknown Primary

# UMCCR Genomics Platform Group

**Workflow Development**

Technology Assessment

Standards Development & Implementation

Workflow Development: Rapid WGTS

Supporting Precision Oncology

**IT Infrastructure**

Primary Processing → Post-process → Storing Data → Curation & Reporting / Data sharing

Workflow Development: **Data Flow**

Supporting Precision Oncology

**IT Infrastructure**

Primary Processing → Post-process → Storing Data → Curation & Reporting / Data sharing

Workflow Development: **Primary Analysis**

Supporting Precision Oncology

Illumina-University of Melbourne Partnership

Partnering to provide the infrastructure, expertise, systems and analysis to translate and implement genomics into routine clinical care.

Core Platform: DRAGEN/ICA

# IT Infrastructure

Primary Processing → Post-process → Storing Data → Curation & Reporting / Data sharing

Workflow Development: Post-process / Reporting

Supporting Precision Oncology

Somatic SNV
Germline SNV
Mutation consequence
Somatic CNV
Ploidy
Purity
Clonality
Mutational burden
Mutational signature

Somatic SV
SV- gene consequence
CNV-gene consequence
HRD Detection
MSI sensing
Viral integration
Pathogen detection

Genomic stability
Clinical prioritization

# Workflow Development: Post-processing

umccrise

Workflow Development: **Reporting**

MultiQC, PCGR/CPSR, ...

Workflow Development: Portal

# Sharing data saves lives
## THE GLOBAL ALLIANCE FOR GENOMICS & HEALTH

### The need:

- Data from **millions of samples** is needed to address questions in rare disease, complex disease and cancer.

### The challenge:

- Data in silos.
- Lack of standard analysis methods.
- Different approaches to regulation, consent and data sharing.

Global Alliance
for Genomics & Health

Australian Genomics
Health Alliance
Slide courtesy of Tiffany Boughtwood, Australian Genomics

(Aggregate) Data Sharing: cBio Portal

# UMCCR Genomics Platform Group

Workflow Development

**Technology Assessment**

Standards Development & Implementation

# GEN3
DATA COMMONS

# Gen3 is how data commons are made.

A data commons is a cloud-based software platform for managing, analyzing, harmonizing, and sharing large datasets. Gen3 is an open source platform for developing data commons. Data commons accelerate and democratize the process of scientific discovery, especially over large or complex datasets.

**Experience Demo**    **Get Started**

...data commons. Data commons accelerate and democratize the process of scientific discovery, especially over large or complex datasets.

Experience Demo          Get Started

OpenCGA

OpenCB    IVA v2.0.0-dev    Variant Browser ▾    Variant Analysis ▾    Clinical Analysis ▾    Catalog Metadata ▾    GA4GH ▾

Studies ▾    About    demouser ▾

Projects / family / corpasome

## Variant Browser

▶ RUN

Filters    Aggregation

**STUDY AND COHORTS**

Studies Filter

Corpus Family

In any of [OR]    In all (AND)

**GENOMIC**

Chromosomal Location

3:444-55555,1:1-100000

Feature IDs (gene, SNPs, ...)

Search for Gene Symbols

BRCA2,ENSG00000139618,ENST0000 0544455,rs28897700

Select SO terms

10 items selected

Disease Panels

---

Kids First Data Resource Center

🏠 Dashboard    🔍 Explore Data BETA    🗄 File Repository

Resources ▾    Adam ▾

### Explore Data

▶ NEW    ▶ OPEN    ▶ SAVE    ▶ DELETE    ▶ SHARE

Search all filters    Quick Filters    Study    Demographic    Clinical    Biospecimens    Available Data

Combine Queries:    and    or    CLEAR ALL

☐ #1    Available Data Types is any of    Aligned Reads, gV...    ✕

☐ #2    #1 ✕    AND ▾    Kf Id is    set_id:97eede81-ee...    ✕

⊕ START NEW QUERY    ⎘ DUPLICATE QUERY

**Cohort Results** for Query 2    ⚡ **6,257 Participants** with **29,441 Files**

▦ Summary View    Table View

**Available Data Files** 29,445

| Data Type | Experimental Strategy | Files |
|---|---|---|
| Aligned Reads | WGS | 12,130 |
| Aligned Reads | RNA-Seq | 2,243 |
| Aligned Reads | WXS | 498 |
| Aligned Reads | miRNA-Seq | 246 |
| Aligned Reads | – | 55 |
| gVCF | WGS | 5,738 |

**Studies** 8

■ Probands  ■ Other Participants

Congenital Heart Defects
Orofacial Cleft: European...
Ewing Sarcoma: Genetic Risk
Pediatric Brain Tumors: C...
Disorders of Sex Development
TARGET: Neuroblastoma
TARGET: Acute Myeloid Leu...

0    500    1,000    1,500    2,000
# Participants

**Most Frequent Diagnoses**

■ Probands  ■ Other Participants

Ewing Sarcoma
Cleft Lip Palate
Low Grade Glioma
Neuroblastoma
Acute Myeloid Leukemia
Ventricular Septal Defect...
Adolescent Idiopathic Sco...
Tetralogy Of Fallot
Atrial Septal Defect Osti...
Medulloblastoma

0    100    200
# Participants

---

data c...
demo...
over la...

**Experience Demo**    **Get Started**

Overture

# Gen3 is how data commons are made.

A data commons is a cloud-based software platform for managing, analyzing, harmonizing, and sharing large datasets. Gen3 is an open source platform for developing data commons. Data commons accelerate and democratize the process of scientific discovery, especially over large or complex datasets.
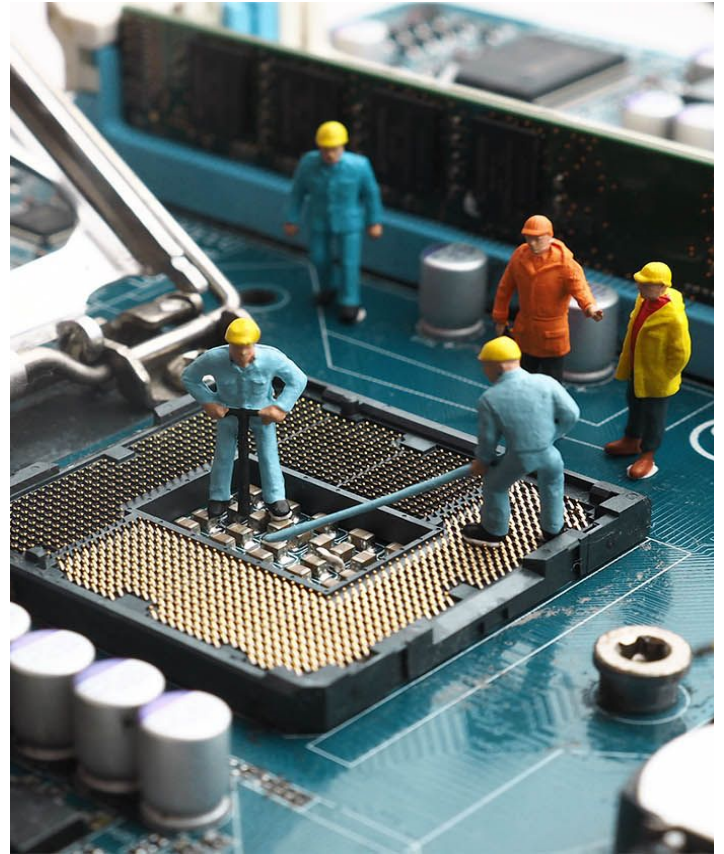
**Experience Demo**　　　**Get Started**

# Actively Maintained Code

*"A piece of software is being sustained if people are using it, fixing it, and improving it rather than replacing it."*

software-carpentry.org/blog/2014/08/sustainability.html

Active development on Github

github.com/uc-cdis

Lively Slack channel

Gen3: Graph Data Model

Relationships between subjects, clinical, biological and molecular data

Gen3: Data Dictionaries

Rules for structured data using external terminology standards / ontologies

Gen3: Microservices

Modular components with defined interfaces

Gen3: Sheepdog

Data ingestion and validation service (UI and API)

Gen3: Windmill

Web portal for data submission, query, exploration, and analysis

Gen3: IndexD

ID management, checksum and size catalogue

Gen3: IndexD

Supports multiple URLs for stored objects

Gen3: Fence

Authentication and Authorisation – OpenID Connect with support for Google, eRA Commons, eduGain, …

Gen3: **Peregrine**

Graph-based metadata queries

# Gen3: Lightweight Workspaces

Basic support for Jupyter notebooks for analysis and visualization in R, Python

Setting up Gen3

master    compose-services / docs / setup.md    Go to file   ...

uwwint Highlighted the note to update docker config to 6GB. I walked straigh... ...  ✕    Latest commit ba1a880 on Nov 2, 2021   History

3 contributors

114 lines (84 sloc)   10 KB    &lt;&gt;    Raw   Blame

# Setup

## Dependencies

- OpenSSL
- Docker and Docker Compose

## Docker and Docker Compose Setup

If you've never used Docker before, it may be helpful to read some of the Docker documentation to familiarize yourself with containers. You can also read an overview of what Docker Compose is here if you want some extra background information.

The official *Docker* installation page can be found here. The official *Docker Compose* installation page can be found here. For Windows and Mac, Docker Compose is included into Docker Desktop. If you are using Linux, then the official Docker installation does not come with Docker Compose; you will need to install Docker Engine before installing Docker Compose. Go through the steps of installing Docker Compose for your platform, then proceed to set up credentials. Note, that Docker Desktop is set to use 2 GB runtime memory by default.

# Quick: Compose-Services

github.com/uc-cdis/compose-services

How to: POC setup on AWS?

- We discuss high level architecture about POC setup on AWS using Gen3 `compose-services`.
- Centre who interested to adopt Gen3 will do trial POC run before committing to a more production oriented Cloud Automation (Kubernetes cluster) setup.
- Additionally, `compose-services` stack also gives a quick dive into Gen3 foundation services and, it is a good perk for your centre data dictionary development.

Idea:

- We choose AWS EC2 instance `m5.2xlarge` with Hibernation support.
- Gen3 `compose-services` stack simply run on this EC2 instance.
- We hibernate this instance when not in use (over weekend, doing other priority tasks, etc).
- This EC2 instance is front-ed by ALB -- Application Load Balancer. Hence, it shows graceful `504 Gateway Time-out` when EC2 is in hibernation.

Intermediate: Compose-Services on AWS

github.com/umccr/gen3-doc/blob/main/poc/AWS.md

Scalable: Cloud

github.com/uc-cdis/cloud-automation and github.com/umccr/gen3-doc

**AWS EC2 instances:**

- 4x Worker nodes ( `t3.xlarge` )
- 1x Admin VM ( `t2.micro` )
- 1x Forward Proxy VM ( `t2.medium` )

**AWS RDS Databases:**

- 3x RDS PostgreSQL instances ( `db.t2.small` )

**AWS Elasticsearch:**

- 1x Elasticsearch ( `t3.small.elasticsearch` )

**AWS Elastic Kubernetes Service (EKS):**

- 1x Kubernetes cluster

**Others:**

- 1x Virtual Private Cloud (VPC)
- 1x NAT Gateway
- 1x Elastic Load Balancer (ELB)

# Scalable: Cloud

github.com/umccr/gen3-doc/blob/main/cloud/AWS.md

AWS Services

Terraform

Kubernetes, Docker, Linux

ElasticSearch

PostgreSQL

GraphQL, Graph and DAG, ETL process

ReactJS SPA

Identity Provider (IdP), Federated AuthN/Z,

Single-SignOn (SSO) setup

…

Good troubleshooting skills

# Scalable: Cloud

github.com/umccr/gen3-doc/blob/main/cloud/AWS.md

# Quickstart: Sample Data Models

Define a data model

Generate a commons with a Gen3 API

Load data into the commons

Start exploring

# Limitations and Difficulties

- Complicated infrastructure

- Data models are complex; one per instance

- Lack of granular control over permissions and data access

# A national approach to genomics information management (NAGIM)

- The vision for human genomics data sharing in Australia requires considerable coordination and collaboration.

- The NAGIM Blueprint sets out a series of principles to guide decision-making on the responsible collection, storage, use and management of genomic data.

- Australian Genomics is developing recommendations for implementing NAGIM.

- In 2021 Australian Genomics led an implementation prototyping phase in response to NAGIM.

- A panel of external assessors are evaluating prototype submissions presently.



Figure 1. Components of a national genomics infrastructure

A National Approach to Genomic Information Management, Australian Genomics Implementation Recommendations Progress Report, November 2021

HUMAN GENOMES PLATFORM PROJECT

Documentation & training

Streamlined submission to data archives

Virtual cohort querying via federated data portals

Automated Data Access Committee submission & approval

Identity and access management

Project Partners

AUSTRALIAN ACCESS FEDERATION

Australian BioCommons

ARDC
Australian Research Data Commons

BIOPLATFORMS AUSTRALIA

Garvan Institute of Medical Research

NCI AUSTRALIA

QIMR Berghofer Medical Research Institute

THE UNIVERSITY OF MELBOURNE

zero CHILDHOOD CANCER

Jess Holliday and Marion Shadbolt

# Supporting Australian Cardiovascular Disease Research

- We have begun working with partners from the Australian Cardiovascular Alliance to establish systems to support identification of biomarkers of increased risk of heart attack.

- We're currently mid way through an 8 month project to establish a new Gen3 instance and populate with 3 coronary artery disease cohorts.

- Data harmonisation across the cohorts is underway.

- We've populated the instance with synthetic data to allow functionality testing.

# Acknowledgements

Australian BioCommons &

Melbourne Bioinformatics

- Jeff Christiansen
- Lisa Phippard
- Jess Holliday
- Marion Shadbolt
- Steven Manos
- Uwe Winter
- Andrew Lonie
- Nuwan Goonasekera

UMCCR

- Victor San-Kho Lin
- Florian Reisinger
- Andrew Patterson
- (Grant Lee)
- (Lavinia Gordon)

Australian Genomics

- Tiffany Boughtwood
- Marie-Jo Brion
- Sarah Casauria

Children's Cancer Institute, Sydney
*Zero Childhood Cancer program*
- Marie Wong-Erasmus
- Kamile Taouk
- Mark Cowley
- Vanessa Tyrrell

Children's Hospital of Philadelphia
- Allison Heath
- Adam Resnick
- Miguel Brown
  Yuankun Zhu
- Bailey Farrow

# Questions?

# NEXT …

*Conservation genomics in the age of extinction*

Dr Carolyn Hogg, University of Sydney

8 March 2022

biocommons.org.au/events

Australian
**BioCommons**

# Tell us what you thought …

Feedback survey

# Thanks for joining us!

## The Australian BioCommons is enabled by NCRIS via Bioplatforms Australia funding