# A FAIRer future for genomics research in Australia: creating an Australian Reference Genome Atlas

K. Hall, J. Christiansen, N. dos Remedios, S. Richmond, N. Ward & H. Holewa
*eResearch Australasia Conference*, 19 October 2022

***Myrmecia nigrocincta***
Jumping Jack Ant

**Myrmecia nigrocincta**
Jumping Jack Ant

"More than a century of research has led to the identification of some key navigational strategies, such as compass navigation, path integration, and route following. Ants have been shown to rely on visual, olfactory, and idiothetic cues for navigational guidance."

Freas, C.A. & Schultheiss, P. (2018) How to Navigate in Different Environments and Situations: Lessons From Ants. *Frontiers in Psychology*, **9**: 841, https://doi.org/10.3389/fpsyg.2018.00841.

# Genomic sequencing and analysis have been identified by the Australian Government as being a critical technology for our national prosperity and security.

**Key sectors**
- Healthcare & Medical
- Agriculture
- Environment
- Defence & Defence Industry
- Energy

**Estimated impact on national interest**
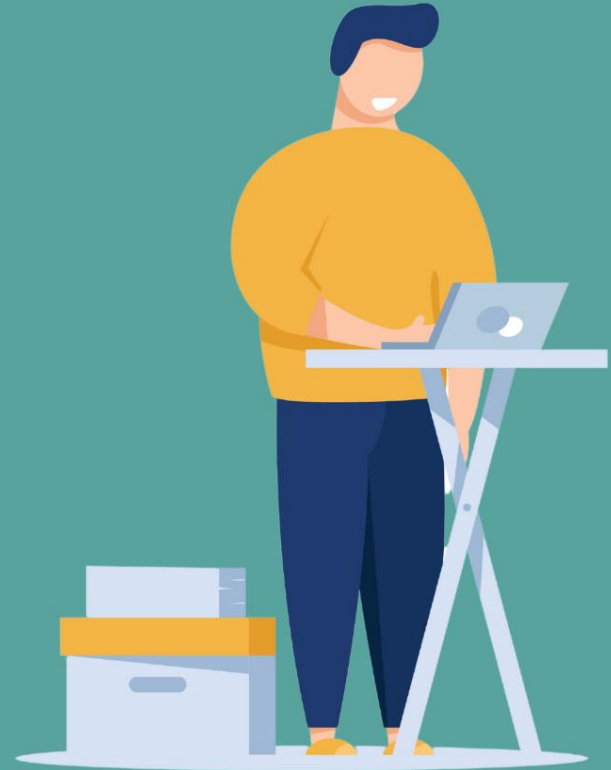- Economic Prosperity - High
- National Security - Med

Policy responses to bushfires (and another environmental catastrophe) responses can be proactive, not reactive, when driven by these data.

The Australian BioCommons estimates that there are currently 15,000 life science researchers in Australia currently using genomics data to answer research questions.

## Genomics data sources are:

- **nested**
- **hidden**
- **variable**
- **complex**
- **different**
- **scattered**
- **embedded**
- **distributed**
- **disconnected**

# GenBank (NCBI)

- created in 1982

- currently holds:

  traditional GenBank records:
  - 240,539,282 sequences
  - 1,562,963,366,851 bases

  set-based (WGS/TSA/TLS) records
  - 2,857,043,692 sequences
  - 18,787,298,109,534 bases

- acceptance standards hinge on genetic annotation, not taxonomic or other isolate metadata accuracy
- not peer-reviewed



**GenBank and EMBL database 1986/1987**
© David Landsman, Bethesda, Maryland (CC-BY)
https://twitter.com/bffo/status/289529886484348928

# GenBank growth since 2014



Total Base Pairs and Entries in GenBank (log scale)



Total Basepairs in GenBank (billions)



GenBank Entries (millions)

Data source: NCBI-GenBank Flat File Release 252.0, Genetic Sequence Data Bank, October 15 2022, https://ftp.ncbi.nih.gov/genbank/gbrel.txt

100% of researchers have told us that they do not trust data they download from GenBank*



**Turning to idiothetic cues**

* may not actually be true

# Current researcher strategies to find data

- GenBank
- BioPlatforms data portal
- Barcode of Life Database
- EMBL-ENA
- Project repositories and websites (*e.g.* Apollo, GoaT, Darwin Tree of Life, ReefGenomics)
- Taxon or community databases (*e.g.* WormBase, ANEMONE)
- Read papers and track back to repositories (Dryad, Zenodo, OSF)
- Word-of-mouth
- Own library generation

# Current researcher strategies to trust data

- metadata for voucher in recognised and curated collection (*i.e.* a museum or herbarium)
- reputation of researcher generating and depositing data
- metadata about methods

# ARGA concept development

Born out of researcher frustration at having to search across multiple data sources to find and access genomics data, ARGA aims to aggregate data from a number of reputed domestic and international sources in a single location.

STEP 1: NEED IDENTIFIED
## 2021

PHASE 2: PRODUCT TESTING
## 2022

PHASE 3: PORTAL RELEASE
## 2023

# Acknowledging ARGA partnerships

# ARGA objectives and vision

The Australian Reference Genome Atlas is an indexing service for discovering, filtering and accessing complex life science data.

For plants, animals, microbiota and other species endemic or relevant to Australia, ARGA will build a platform to locate and aggregate genomic data, including:

- reference genome assemblies
- genome annotations
- population and variant sets
- DNA barcodes
- coding and non-coding DNA sequences

# Genomics data cycle

Data that are newly generated by research projects can be consumed by those researchers and also made available for consumption to others.

Genomic data from specimens can also be enriched by intersecting it with other observations, using metadata and processing pipelines.

Enriched data can then be consumed to answer novel questions.

Data enrichment can seed the generation of new data by identifying targets.

Consume

Generate

Enrich

## Prototype of ARGA Portal
Source code:
https://github.com/ARGA-Genomes

- Data from ingested and processed via GBIF pipelines using Darwin Core Archives (DwC-A) metadata standards.
- Working prototype interface built using React.js implementing Solr searching.
- Traits facet filters built to slice results, including: vernacular groups, biomes, data types.

## Data sources

Indexed:

- NCBI-GenBank (RefSeq, Genome)
- Bioplatforms Australia
- BOLD systems (Barcode of Life)

Next to be indexed:

- NCBI-GenBank (nucleotides, assembly, SRA)
- The European Nucleotide Archive (ENA)
- Genomes on a Tree (GoaT)
  - Genome Size
  - Plant DNA C-values Database

Data exploration site: **https://nectar-arga-dev-1.ala.org.au**

# A FAIRer future for genomics research

**Findable**
- data from multiple reputed repositories indexed using recognised metadata standards (Darwin Core and Darwin Core extensions)
- enriched by intersecting with specimen metadata from the Atlas of Living Australia to increase trust and quality assessment
- PIDs, citations and original sources linked

**Accessible and Interoperable**
- multiple data formats, user-selectable
- new and flexible search strategies by taxon, phenotype and genomic annotations

**Reusable**
- DOIs of searches
- integration with BioCommons platforms

# ARGA Development Team

| | | |
|---|---|---|
| Nick dos Remedios | Atlas of Living Australia | Lead Systems Engineer |
| Christopher Mangion | Australian BioCommons | Data Engineer |
| Matt Andrews | Atlas of Living Australia | Systems Support |
| Mok | Australian BioCommons | UX/UI Designer |
| Keeva Connolly | Australian BioCommons | Scientific Business Analyst |
| Kathryn Hall | Atlas of Living Australia | Project Manager |