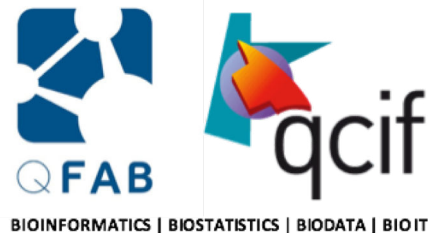# Galaxy Australia – the Bring Your Own Data platform enabling multi-omics analysis for biology researchers

**Dr Gareth Price, Service Manager of Galaxy Australia (Queensland Facility for Advanced Bioinformatics)**

# Biosciences: the nature of the Australian research community

**30,000** health/biosciences researchers
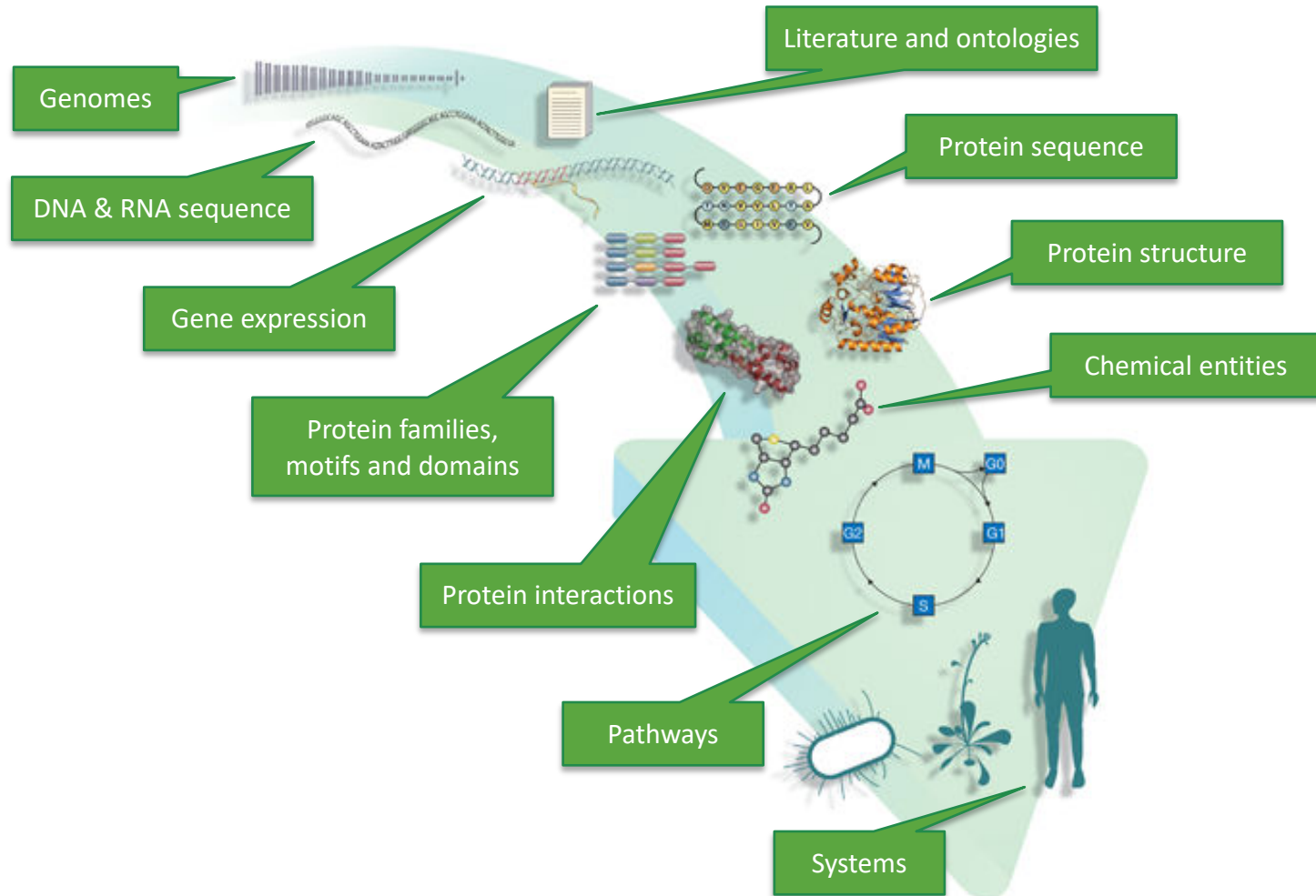
**18,000** health/biosciences HDR students

**48,000** health/biosciences PG course work students

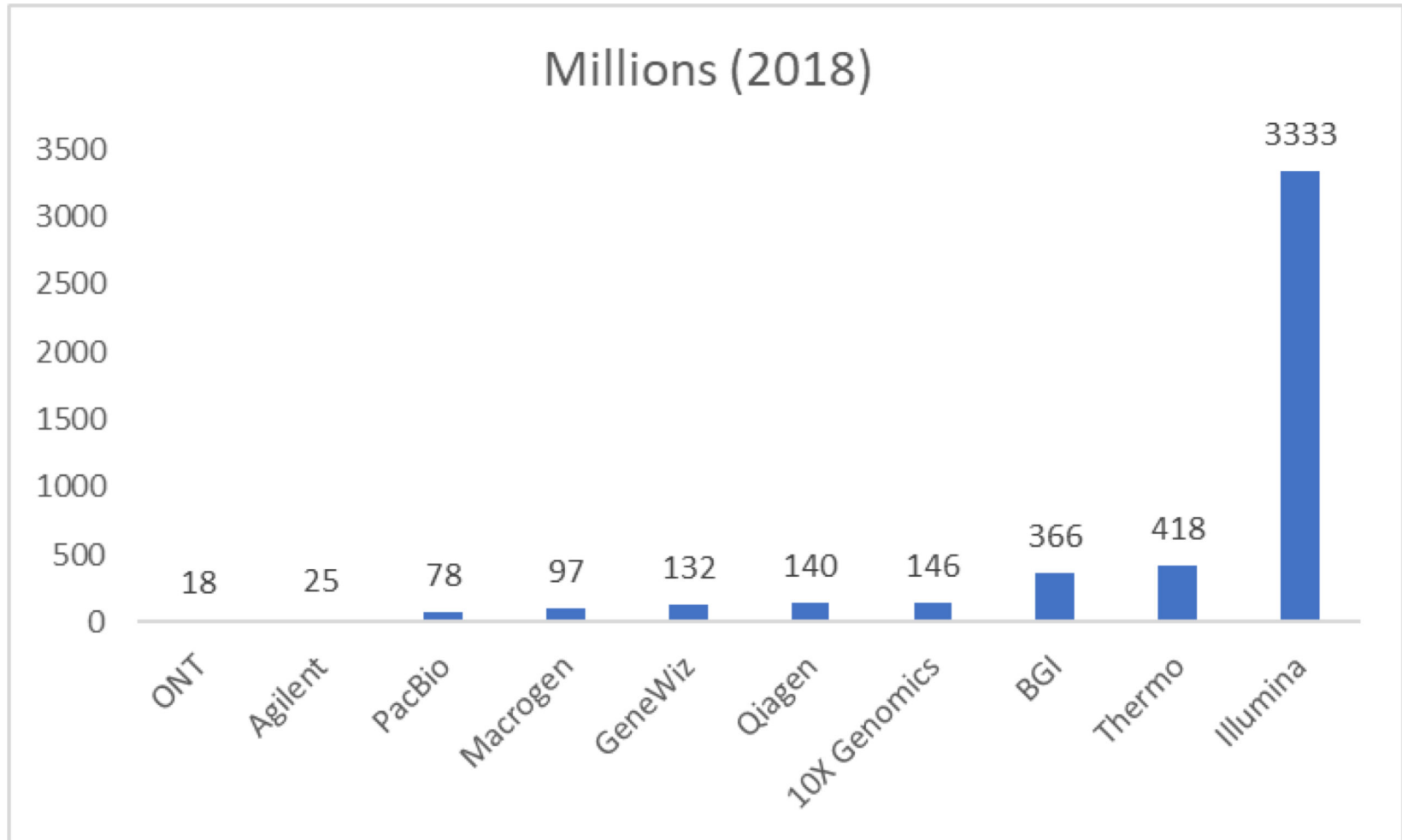**(163,000 + 40,000 =) 200,000 health/biosciences UG students**

*1,000 to 1,500 bioinformatician/computational biologists*

*=> the IT smarts need to be packaged up and delivered cleanly*

# From Genomes to Systems



Genomes

DNA & RNA sequence

Gene expression

Protein families, motifs and domains

Protein interactions

Literature and ontologies

Protein sequence

Protein structure

Chemical entities

Pathways

Systems

EMBL Australia

**Bioinformatics Resource**

# The Production landscape of genomics

Millions (2018)

| Company | Value |
|---|---|
| ONT | 18 |
| Agilent | 25 |
| PacBio | 78 |
| Macrogen | 97 |
| GeneWiz | 132 |
| Qiagen | 140 |
| 10X Genomics | 146 |
| BGI | 366 |
| Thermo | 418 |
| Illumina | 3333 |

**https://www.genengnews.com/a-lists/top-10-sequencing-companies/**

# usegalaxy.org.au

# Efficiency through Galaxy controlled scheduling

# Taking the IT out of bioinformatics

- 946 tools supporting analysis applying more than 200 reference datasets
- simple data upload with optional rule based data management
- retains "histories" of analyses
- a large number of cutting edge tools
- an app store (over 7,000 tools)
- All aspect peer reviewable and transparent
  - Workflows
  - Citations / origins
  - Histories
  - Data

# A Brief History in Time

- Genomics Virtual Lab for public and private VLs, with Galaxy as a managed public service
  - Galaxy-Melb, Galaxy-Qld and Galaxy-Tut (Melbourne)
    - Tool version alignment / tutorial consequences
    - Duplicate staff roles
    - Tut was re-imaged over and over, no longevity to user accounts
- Rebrand Gal-Qld to Galaxy Australia and move to org.au
  - Staged migration of Melb to Aus
- Bring online 2 x Pulsar (repurpose Galaxy-Melb allocation and new allocation at NCI Canberra)
- Hybrid training events: 4 run in Galaxy with local facilitators and a lead trainer.
  - 11 of sites, 1064 of participants
- Policy, policy and policy – have a lot to thank Galaxy Europe for on this front

# Galaxy Australia – distributed architecture

# Community Impact



Jobs/Month and Total Jobs Run on Galaxy Australia.

"We have sequenced well over 500 transcriptomes and genomes, and routinely use Galaxy Australia for many bioinformatics processes.

It is easy to use, has high computational power, a sophisticated support structure and enables global collaboration through straightforward data sharing.

We greatly appreciate the service."

Dr Fabio Cortesi & Prof Justin Marshall, Queensland Brain Institute

# Growing our Australian Community

- Ongoing funding (for 4years) with a view to establish an enduring national research infrastructure

- adding to service functionality by adding metabolomics and phylogenetics

- greater security through institutional authentication, linking this to higher resourcing for authenticated users

- providing data sharing and data movement options through AARNET Cloudstor

- anticipating needs of other additional communities – esp. single cell and genome assembly and ultra long read technologies

- adding new resources: Galaxy slave servers (aka Pulsars) around the country

# Growing our Australian Community – Metabolomics Australia

- Provide vendor agnostic tools and visualisations

- Leverage an increasingly mature global open source community

# Growing our Australian Community – Oz Mammal Genomes and Target Capture Panel pipelines

- Build on 2018 experience with BPA Data Portal Australian Microbiome

- User interact at BPA Data Portal (and potentially never leave!) to trigger pipeline on Galaxy Australia

- Analysed data returned to Data Portal and made publicly visible

# Growing our Australian Community – authentication and storage links

- Galaxy Project is an open collaborative global project
- Users on public Galaxy services rarely are required to authenticate
- But with authentication comes opportunity for tailored or seamless servicing

- Link Galaxy to AAF
  - Link to institutional HPC destinations

- Link Galaxy to AARNET Cloudstor
  - Data Input
  - Data Export
  - Service provider mediated Input

![Galaxy Australia logo]

# usegalaxy.* - a global platform and support network

**Distributed reference data between .* servers**

- reduced System Administration per locale
- Australian contribution to global efforts
- users are not restricted to "local" content

**Intergalactic Data Commission**

- formed in 2018
- regulation, automation and documentation of the CVMFS reference collections
- Australian representation on the IDC

**Galaxy Project Executive Steering Committee**

- formed in 2019
- Australian representation on the Committee

# usegalaxy.* - a global platform and support network

**Galaxy Training Network**
- Australian content contributions
- Multi-language options
- Simple and comprehensive options
- Many peer reviewed best practices
- Synergised tool set

**More usegalaxy.* services**
- Growing number of countries and regions are forming their own usegalaxy.* service
- Increasing content development
- Distributing SysAdmin activity
- Thursday's keynote address: **The Development of ASEAN Federated Identity and Login Management, and Galaxy ASEAN Community**

# Computational Power

- Per patient sequencing data size is not changing, but the scope of analysis is/will:
  - WES and WGS
  - Short read for SNV
  - Long read of CNV and SV
  - Multi-omics (epigenome and transcriptome)
- Number of requests growing

- Choices are:
  - More computational power
  - Better algorithms
  - Appropriate use of computational architecture

# Solutions for Data Analysis

- **Freeware**
    - Galaxy
    - R Studio
    - Command Line / HPC

**Most equivocal solution**

Galaxy Australia user numbers (as of Sept 2018 - 2268) as CLC-Bio users is a difference of 750k funds (+750K in kind) vs approx. $11million annual licence fees

*This does not include the cost of computers to support CLC-Bio installations*

- Commercial
    - **Office**
        - Excel
    - **Agilent**
        - Cartagenia Bench Lab for Molecular Pathology
    - **Illumina**
        - BaseSpace
    - **Qiagen**
        - CLC-Bio Suite of Analysis Products
    - **ThermoFisher**
        - Ion Reporter

RESEARCH ARTICLE

# Comprehensively benchmarking applications for detecting copy number variation

Le Zhang [1,2,3] *, Wanyu Bai[1], Na Yuan[4], Zhenglin Du[4] *

# Computational Power - more

- Cloud Life Sciences (formerly Google Genomics)
  - *Broad Institute replaced its in-house genome sequence analysis computers and storage with Google Cloud Platform, which delivers greater speed, scalability, and data security.*
- AWS Genomics
- Galaxy Australia

# AnVIL



## [https://anvilproject.org/](https://anvilproject.org/)

- [Terra](#) is an analysis platform that allows users to access data, run analysis tools, and collaborate, powered by Google Cloud Platform.

- [Gen3](#) is a cloud-based software platform for managing, analyzing, harmonizing, and sharing large datasets.

- [Dockstore](#) is an open platform used by the GA4GH for sharing Docker-based tools described with the Common Workflow Language (CWL), the Workflow Description Language (WDL), or Nextflow (NFL).

# Computational Power – better algorithms / architecture

- Analytical code can be deployed in a systematic way (Ansible, Docker, Singularity, Windows Updates!) but this makes assumption about the computer architecture

- Solution: Graphics processing unit (GPU) card or Field Programmable Gate Array (FPGA) card
  - Parabricks
  - Illumina / Dragen
  - ONT – NVIDIA

# Parabricks



**Processing Time in Minutes for 40x Whole Genome**

- 32 vCPU, Open Source Solution
- Parabricks Software on 8 GPU Server

| | BWA + Preprocessing | Haplotypecaller | Mutect2 | GenotypeGVCF | DeepVariant |
|---|---|---|---|---|---|
| 32 vCPU, Open Source Solution | 1090 | 1250 | 1086 | 332 | 580 |
| Parabricks Software on 8 GPU Server | 45 | 15 | 14 | 28 | 25 |

**Accuracy**
BAM : **100%** Match, VCF : **99.99%** Match

# Illumina - Dragen

- The DRAGEN Platform can process NGS data for an entire human genome at 30× coverage in about 25 minutes on premise vs. > 15 hours with a traditional CPU-based system. It set two world speed records for genomic data analysis.

- The DRAGEN Platform can reduce on-premise investments in server clusters and utilization of cloud computing resources.

- Available on premises and on the cloud (Basespace, at AWS Sydney)

# Oxford Nanopore Technologies - NVIDIA

- With MinIT on NVIDIA AGX, they're approaching a 10x performance improvement over previous versions to help unlock real-time human and plant genomics. Its benchtop PromethION product is powered by NVIDIA Volta GPUs and can crank out a human genome for under $800.

- *Stay tuned for Illumina to do the same with Dragen (GATK)*

# You know you've made it when..

Genome **Medicine**

**METHOD** ~~19.5~~ **Open Access**

CrossMark

# A 26-hour system of highly sensitive whole genome sequencing for emergency management of genetic diseases

Neil A. Miller[1†], Emily G. Farrow[1,2,3,4†], Margaret Gibson[1], Laurel K. Willig[1,2,4], Greyson Twist[1], Byunggil Yoo[1], Tyler Marrs[1], Shane Corder[1], Lisa Krivohlavek[1], Adam Walter[1], Josh E. Petrikin[1,2,4], Carol J. Saunders[1,2,3,4], Isabelle Thiffault[1,3], Sarah E. Soden[1,2,4], Laurie D. Smith[1,2,3,4], Darrell L. Dinwiddie[5], Suzanne Herd[1], Julie A. Cakici[1], Severine Catreux[6], Mike Ruehle[6] and Stephen F. Kingsmore[1,2,3,4,7*]

- *Dr. Kingsmore receives the GUINNESS WORLD RECORDS™ certificate for the fastest genetic diagnosis.*
- **San Diego—Feb. 12, 2018**
- https://www.rchsd.org/about-us/newsroom/press-releases/new-guinness-world-records-title-set-for-fastest-genetic-diagnosis/

# Galaxy Australia Team Members

Gareth Price            Nick Rhodes

Igor Makunin            Simon Gladman

Thom Cuddihy            Nuwan Goonasekera

Special Thanks:

Sarah Richmond, Ecoscience Research Cloud (ecocloud)

Derek Benson, CSIRO

Anna Syme, Royal Botanical Gardens, AU

Grahame Bowland, QCIF