# Building a healthier world together



**Precision medicine**

… in Australia



**Find Disease Genes**

… using Apache Spark



**Share Disease Insight**

… using serverless architecture

# CSIRO: Top 1% of global research agencies

- Invented **WiFi**, used in five billion devices globally.
- Developed the vaccine for the **Hendra Virus**.
- Developed the **Total Wellbeing** & Low-Carb Diets.

Credit https://toolstotal.com/

# CSIRO's vision for the #FutureOfHealth



The health system will shift...

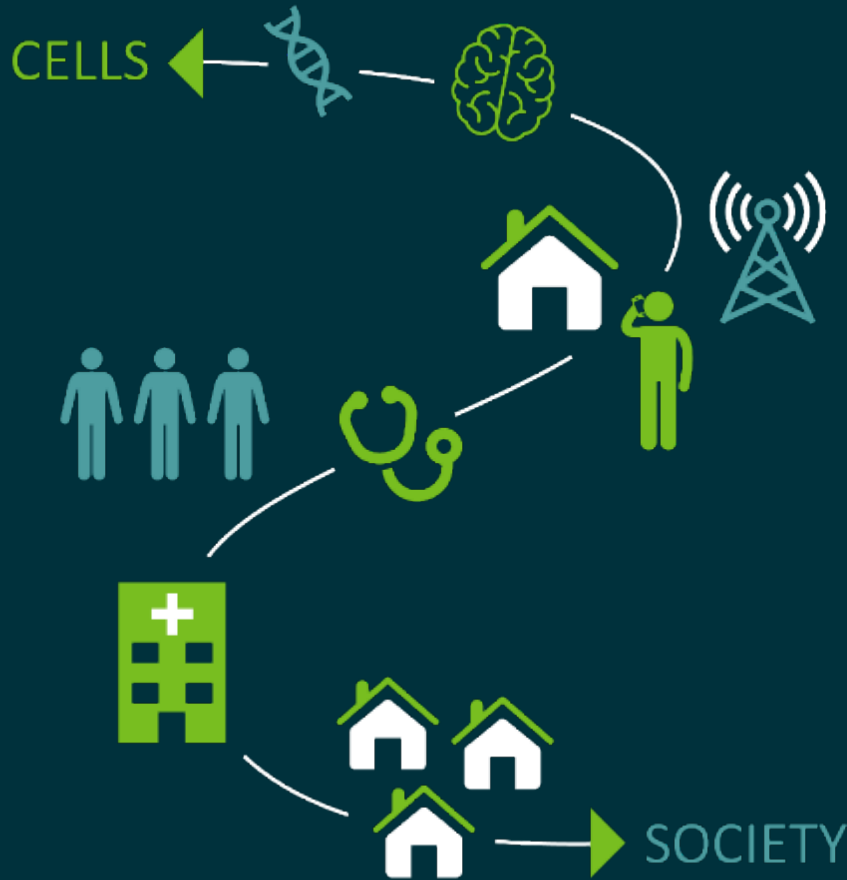...from treating patient illness to managing consumer health and wellbeing

...from accepting one-size-fits-all to precision health solutions

...from a reactive system to a holistic and predictive approach

...from extending life to improving quality of life over a lifetime
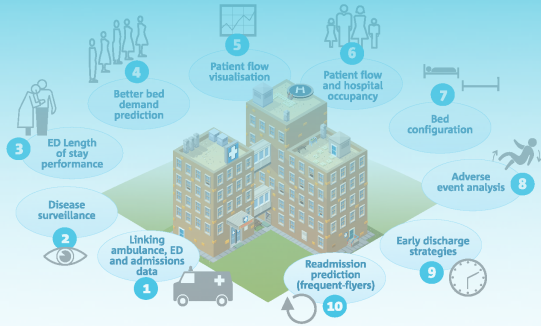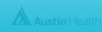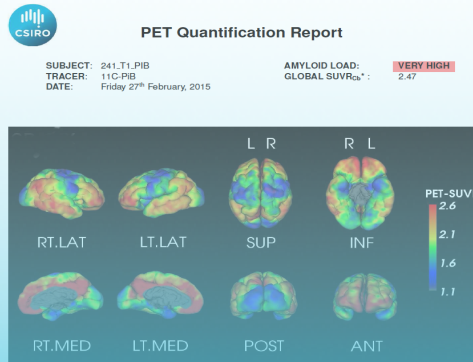
# Precision medicine is at the heart of our research
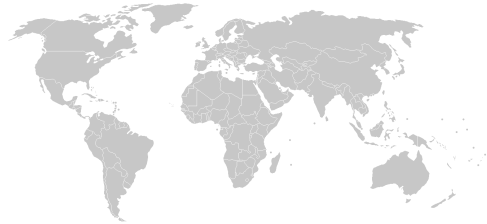
## Hospital forecasts



## Quantitative Imaging



## Disease risk prediction

# Precision medicine is enabled by genomics



project MinE — Make it yours — $70 Million

Global Alliance for Genomics & Health

Queensland Genomics Health Alliance — $25 Million

NSW GOVERNMENT — Genomics Strategy Health — $25 Million

$25 Million — Australian Genomics Health Alliance

Melbourne Genomics Health Alliance — $25 Million

CSIRO

# By **2025** it is estimated that **50%** of the world population will have been sequenced.

Frost&Sullivan

Data acquisition of BigData disciplines in 2025

YouTube

Genomics

Astronomy

Twitter

20 EB Storage / year

Stephens *et al.* BigData: Astronomical or Genomical (2015)

CSIRO

# Overview

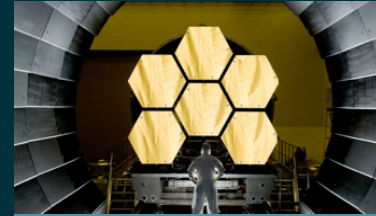**Precision medicine**

… in Australia

CSIRO

**Find Disease Genes**

… using Apache Spark

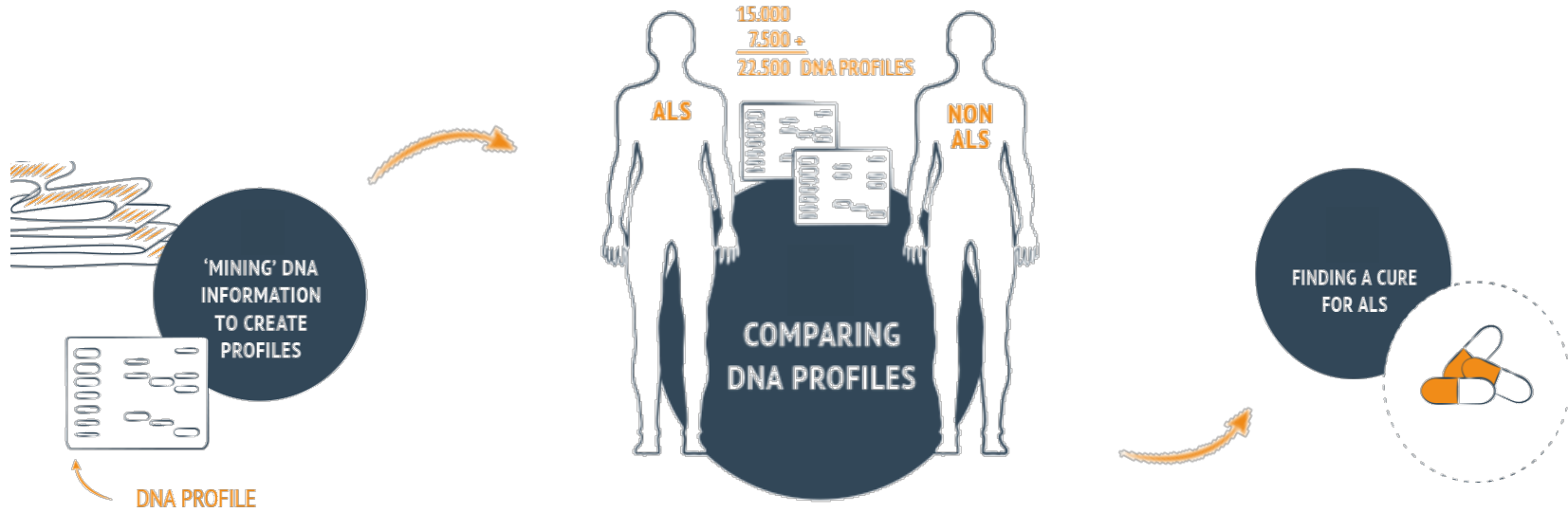VariantSpark
Machine Learning for
Genomic Variants

**Share Disease Insight**

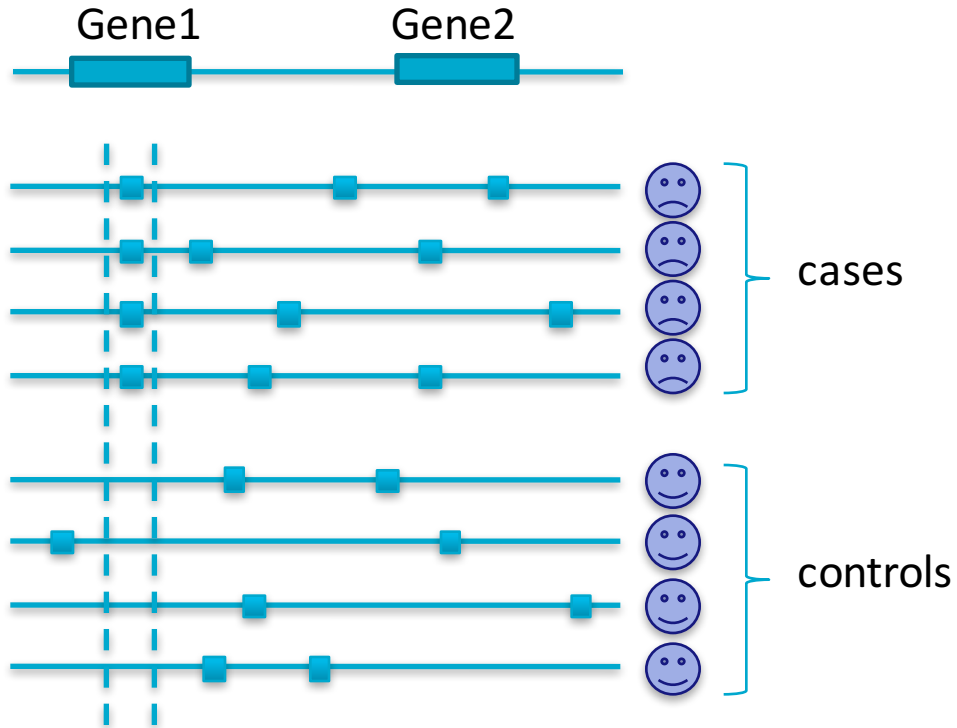… using serverless architecture

GT-Scan2
Computationally Guiding
Genome Engineering

# Finding the cure for ALS

# Finding the disease gene(s)

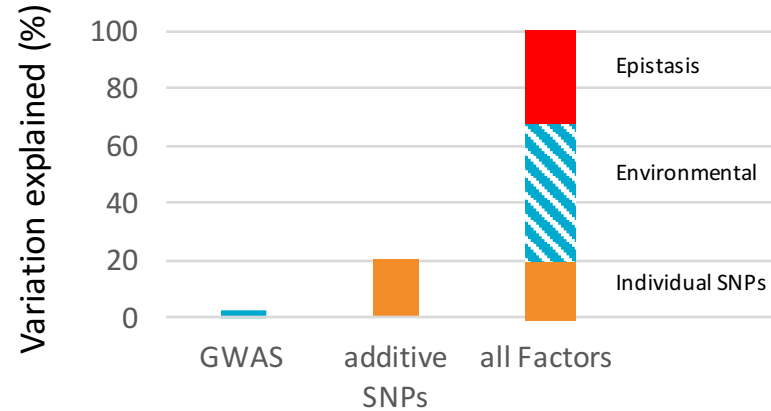# Complex diseases have polygenic risk factors



Polygene analyses suggest that SNPs with *P*-values **well below GWS add significantly** to [obesity] variance explained.

Locke *et al.* Nature 2015

CSIRO

# Complex diseases are driven by multiple genes

Complex diseases are driven by
# multiple interacting genes with variable contribution



Need a more sophisticated ML approach, such as **Random Forest**

# Machine learning on 1.7 Trillion datapoints

Genomic profile

80 Million features

Disease status

Individuals
22,500 samples

Disease genes

project MinE
Make it yours

CSIRO

# Population-scale genomic data analysis requires BigData solutions

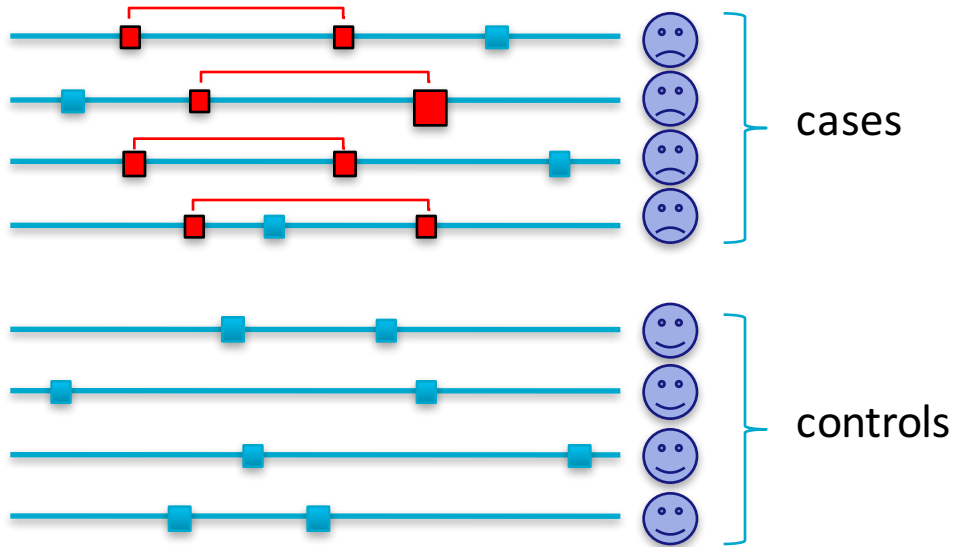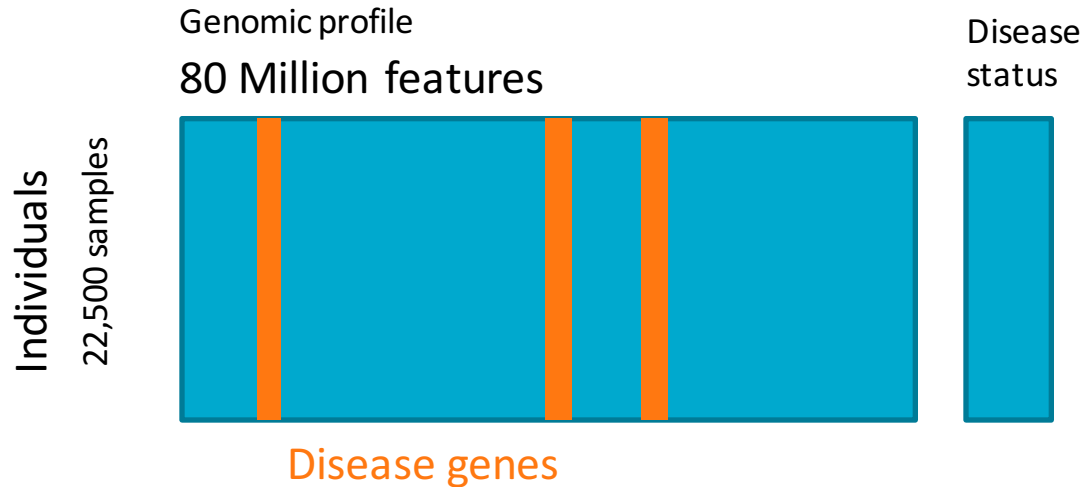| | Desktop compute | High-performance compute cluster | Hadoop/Spark compute cluster |
| --- | --- | --- | --- |
| Focus | small data | Compute-intensive | Data-intensive |
| Node-bound | Yes | Yes | No |
| Parallelization | 10 CPU | 100+ CPU | 1000+ CPU |
| Parallelization procedure | bespoke | bespoke | standardized |

CSIRO solution

VariantSpark
Machine Learning for
Genomic Variants

# VariantSpark: Machine Learning to find markers for complex diseases



**Faster**

"Analyzes 3000 individuals with 80M features in 30 minutes"

**Smarter**

"Requires 80% fewer samples to detect statistical significant signal"
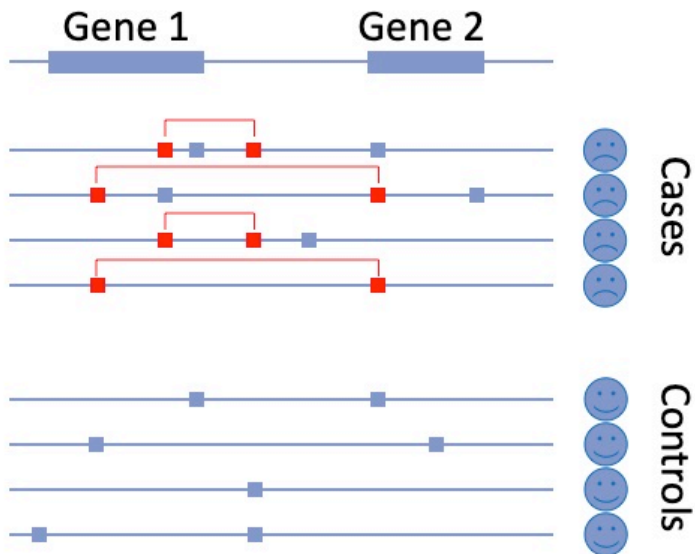
Used by

# **Compute sovereignty** becomes of growing importance

- Retain domestic HPC capability but enable **cloud-like flexibility** (e.g. multi-tenancy)
- **Connect international** datasets without lag
- **-> SuperCloud**

Technical details:  March 13<sup>th</sup> 1630 – 1700
SuperCloud – The evolution of HPC to Software Defined Computing
Mr. Jacob Anders, CSIRO and
Mr. Garry Swan, CSIRO

# Deploying cloud-like workflows on HPC



| Access to bare metal instances | Create Docker for VariantSpark | Enable instances to run docker | Provision Spark-cluster | Connect Jupyter to trigger runs and capture output |

# Demo



Genetic
**Hipster-Index**

Powered by
**VariantSpark**
Machine Learning for
Genomic Variants

CSIRO

and512@razor:~

File   Edit   View   Search   Terminal   Help

[and512@razor ~(keystone_sca19)]$

# Overview



**Precision medicine**

… in Australia



**Find Disease Genes**

… using Apache Spark



**Share Disease Insight**

… using serverless architecture

# Beacon:
## sharing genomic data

- Today, **70 Beacon are lit** to share information about rare genetic diseases.
- **Reduce cost** to enable more sharing.
- Using **serverless (FaaS)** technology.

CSIRO

# Recruiting instantaneous appropriately powered compute

| | Desktop compute | High-performance compute | Hadoop/Spark | Serverless |
|---|---|---|---|---|
| Focus | small data | Compute-intensive | Data-intensive | Agility |
| Node-bound | Yes | Yes | No | No |
| Parallelization | 10 CPU | 100+ CPU | 1000+ CPU | 1-1000+ CPU |
| Parallelization procedure | bespoke | bespoke | standardized | standardized |
| Overhead in the cloud | NA | spin-up lag | spin-up lag | instantaneously |

CSIRO solution

VariantSpark
Machine Learning for
Genomic Variants

Beacon
Serverless lookup of geno-
types and frequencey

# Serverless-Beacon to scale up discovery across continents



**Powerful**

*"Scaling up to large volumes of distributed variant data."*

**Cheaper**

*"Only pay for the resources consumed – zero downtime cost."*

# Three things to remember

- Complex multigenic diseases should be studied using **'wide' ML** (VariantSpark).

- **Serverless architecture** makes even data-intensive web-apps affordable  (Serverless Beacon).

- SuperCloud offers an exciting and potentially cheaper supplement to public cloud providers: **let's build a healthier future together!**

CSIRO

# Let's build a healthier world together

**Team**



Denis Bauer, PhD

Oscar Luo, PhD

Laurence Wilson, PhD

Aidan O'Brien

Natalie Twine, PhD

Arash Bayat

Brendan Hosking

We are hiring…
You?
…email Denis

DATA 61 · CSIRO

Rob Dunne, PhD

Piotr Szul

**Collaborators**

MACQUARIE University
SYDNEY·AUSTRALIA

Lynn Langit

ANU THE AUSTRALIAN NATIONAL UNIVERSITY

Epsagon

Microsoft Azure

aws

Cloudera

databricks

Alibaba Cloud

DiUS

**Software**

GT-Scan2
Computationally Guiding Genome Engineering

NGSANE
Production Informatics for High Throughput Data

VariantSpark
Machine Learning for Genomic Variants

GenPhen-Insight
Genome-Phenome Discovery Framework

Tribes
Detecting distantly related individuals

Beacon
Serverless lookup of geno-types and frequencey

**News**

Top 10 Australian IT stories of 2017

ComputerWeekly.com   Top 10 IT stories of 2017   CW+ Content

#ODSC INDIA

**OPEN DATA SCIENCE CONFERENCE**
Bengaluru, India • Aug 30th - Sept 2nd 2018
#ODSC_india
https://india.odsc.com/

KEYNOTE
PUBLIC SECTOR SUMMIT 2018
aws

CSIRO