

The New Zealand Institute for Plant & Food Research Limited

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI



FLOSSing in the Lab

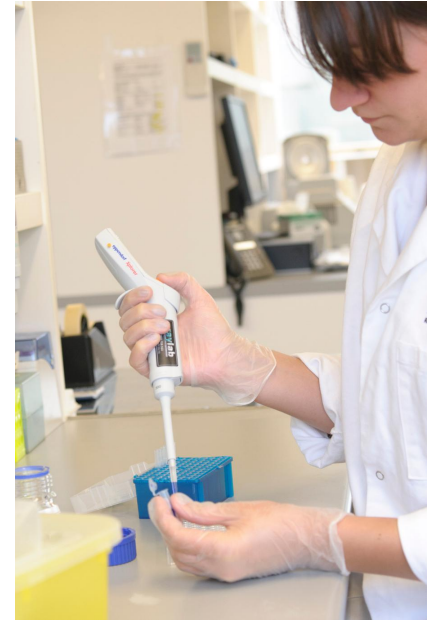
Plant and Food's use of Free/Libre Open Source technologies

Zane Gilmore, Ben Warren, (Eric Burgueno, Roy Storey)

FLOSSing in the Lab

What you are in for:

- Who is Plant & Food?
- What do we do?
- Why do we need software?
- Why we use OSS
- Some examples
- Genetic science
- Genetics and FLOSS



Crown Research Institutes

- AgResearch
- ESR
- Scion
- GNS
- Landcare Research
- NIWA
- **Plant & Food Research**



Manaaki Whenua
Landcare Research



Who we are

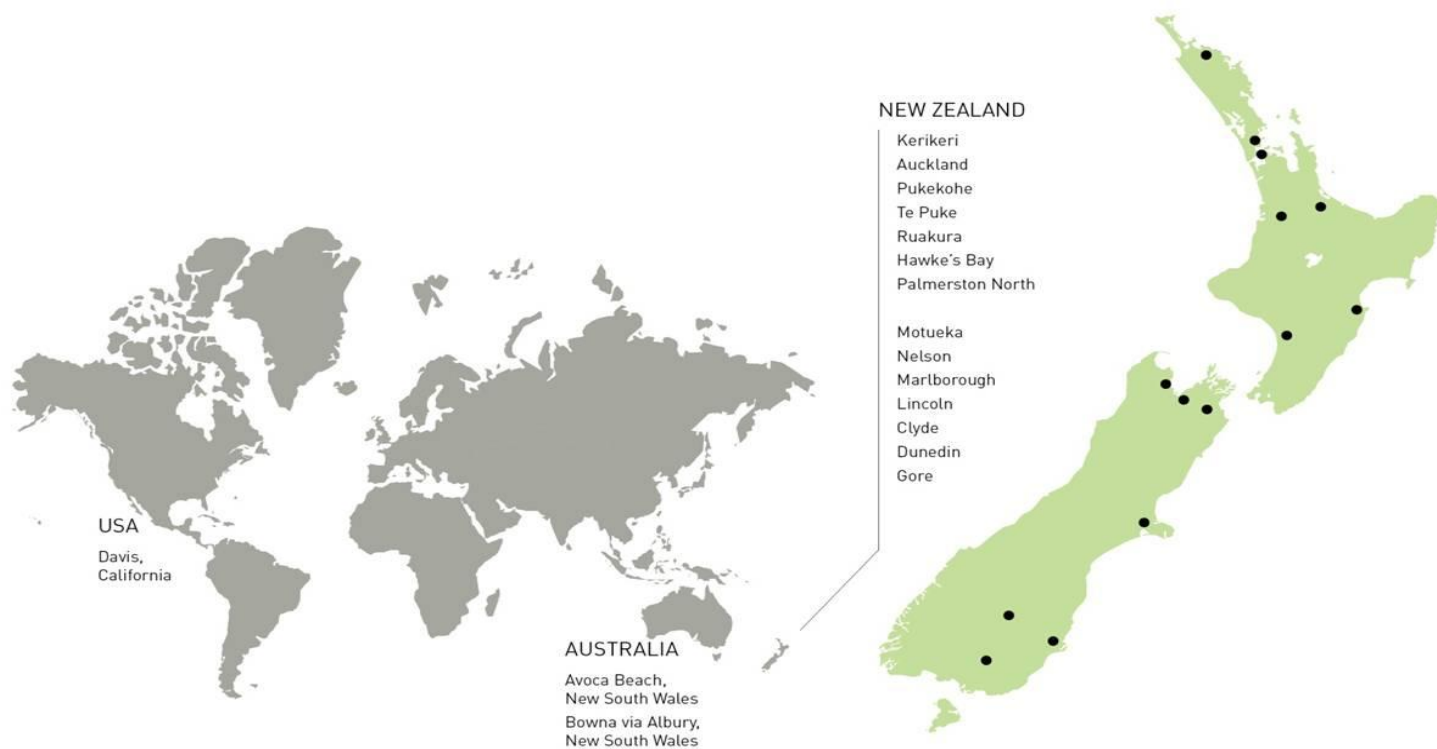
- >> Based in New Zealand
- >> Government-owned Crown Research Institute
- >> Revenue NZ\$119.6 million (2013/14)

A mix of private contracts and royalties,
and NZ Government contracts

Over 900 employees

- >> 650 research staff
- >> 2 dedicated programmers
- >> 15 sites in New Zealand
- >> Representatives in USA, Australia

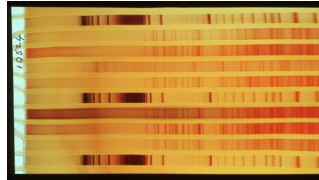
Our Locations



What PFR does

>> Plants

- » Breed new cultivars
- » Cultivation
- » Diseases
- » Insect pests



>> Food

- » Nutritional health
- » Nutrient analysis
- » Food manufacturing

>> Seafood and fishing

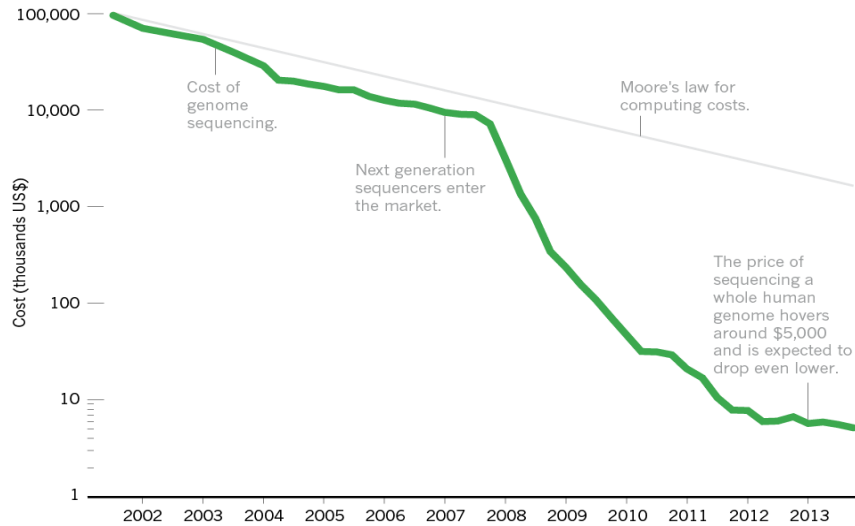
- >> Other stuff but mainly in the service of, or related to the above e.g. soil science and electro-spinning



Computing problems we face

Falling fast

In the first few years after the end of the Human Genome Project, the cost of genome sequencing roughly followed Moore's law, which predicts exponential declines in computing costs. After 2007, sequencing costs dropped precipitously.



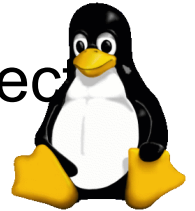
<http://www.nature.com/news/technology-the-1-000-genome-1.14901>

Reproducible research



FLOSS issues

- >> Biologists often aren't at home in the world of computing
- >> Managers (who are often biologists) don't understand FLOSS concepts
- >> CRI funding model
- >> Geneticists ARE good informaticians
- >> Battle is not futile as scientists are clever and respect data



Food Composition (FCDB)

- > 2600 Foods
- > 300 Nutrients/Components/Attributes
- > 400 recipes
- Produce Food Files for Ministry of Health
- Present system is old and creaky
- Data has high “coolness coefficient”
- www.foodcomposition.co.nz
- We are going to rebuild it



More FCDB

- »» Attribute calculator
- »» Recipe calculator
- »» Recipes of Recipes
- »» Meat pie example
 - »» Recipe for pastry
 - »» Recipe for meat stew filling



Kea



- » Plant breeding needs to be done faster
- » We use genetic and chemical analysis for breeding decisions
- » Thousands of plants
- » Kea sample tracking (in-house then with help from Encode)
- » Linux-Django-Postgres stack with Elastic search
- » Just produced alternative provenance system
- » Working on getting it Open Sourced



Other stuff

- >> Data loggers: Lysimeters, rain-shelters
- >> Chemistry databases
- >> Continuous requests



Next Guy

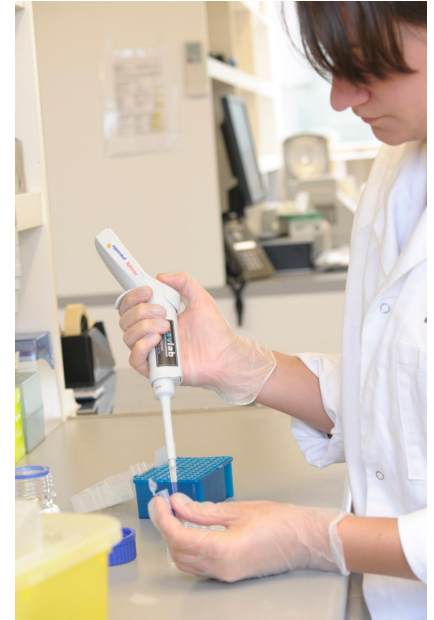
Time for Ben



FLOSSing in the Lab

What you are in for:

- >> Who is Plant and Food?
- >> What do we do?
- >> Why do we need software?
- >> Why we use OSS
- >> Some examples
- >> Genetic science
- >> Genetics and FLOSS



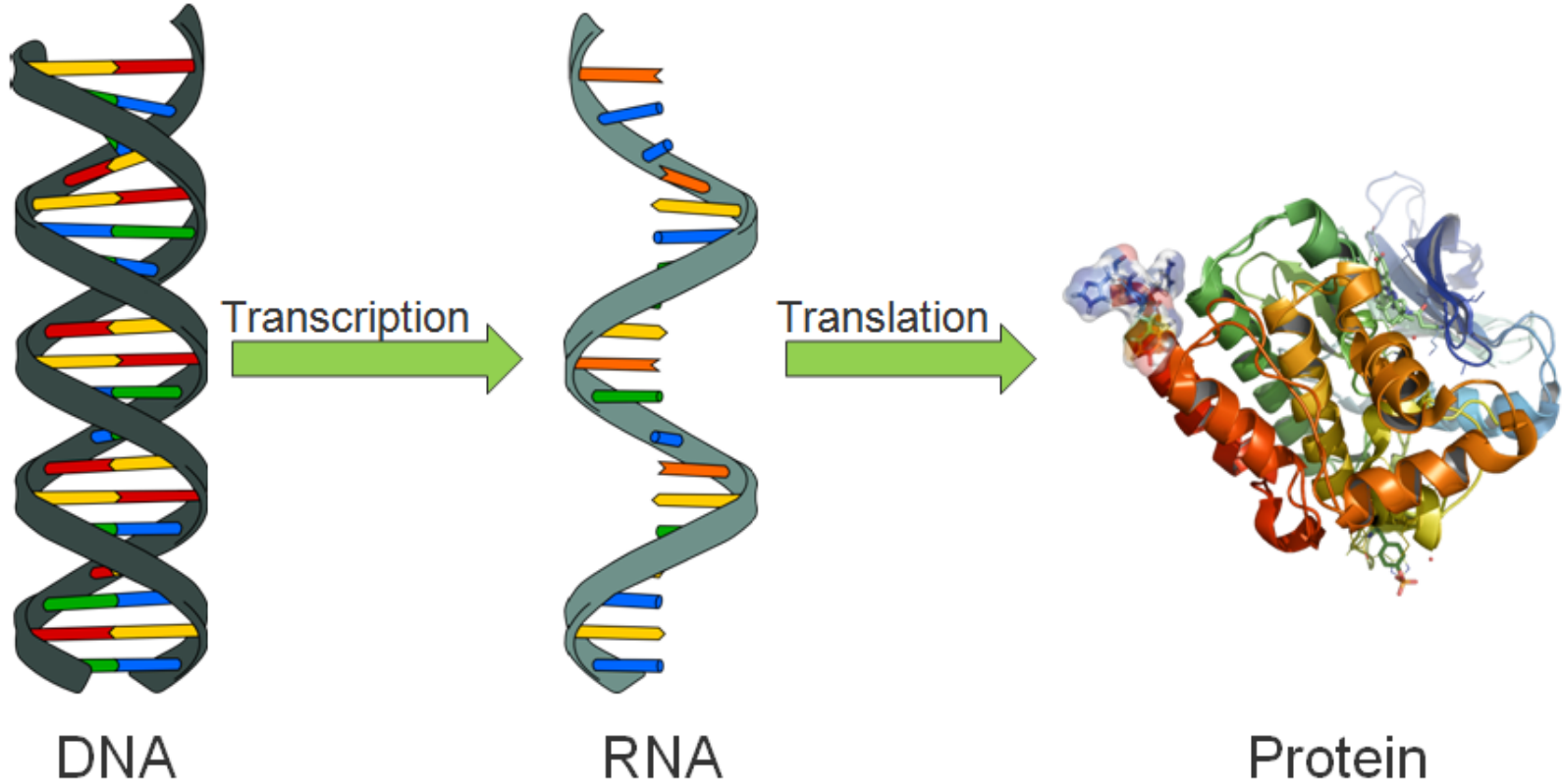
We Do *omics

What is an *omics?

There are many species of *omics.

In the bioinformatics department at PFR we mainly do **genomics** and **transcriptomics**. This is the study of the genome (DNA) and the transcriptome(RNA) respectively.

The Central Dogma



Genome Assembly - A Computational Problem

The assembly problem:



[Mike Haw](#) / [CC-BY-SA-3.0](#)

Given **N** of the same textbooks (possibly differing editions) cut into strips and put in a pile, reconstruct the **N** original texts.

We Need Software for Computation

Assembly and other *omics tasks often require large computations.

- openLava¹ - Job scheduler Software
 - Assign jobs to appropriate nodes
 - Priority queues
- powerPlant - Compute cluster
 - Shared data store (~1PB)
 - Virtual compute nodes
 - Physical compute nodes (e.g. 2TB of memory)



We Need Software for Visualisation

Visual representations of data enhance understanding and spark new ideas about data.

Ensembl² allows us to visualise genomic data.

- >> Can incorporate user data easily
- >> Extendable and customisable



Ensembl - Wine Grape Genome

Chromosome 7: 381,035-383,836



We Need Software for Reproducible Research

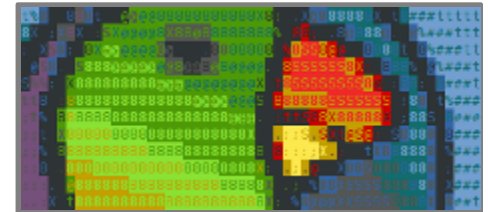
- A workflow is a recipe describing how to get from input data to results
- A well-documented workflow allows the process to be reproduced exactly
- This is necessary for;
 - transparency
 - verification
 - sanity

We Need Software for Reproducible Research

Moa⁵ provides extendable templates based on common workflows.

“Moa hopes to make meticulous organization of a command line project much less of a burden - leaving you to focus on the fun parts.” - Mark Fiers, <http://moa.readthedocs.org/en/latest/>

- Integration with Git
- Integration with openLava



We Need Software for Reproducible Research

We can use Git³ to store workflows, allowing reproduction of the workflow at any version.

- Branches can store specific instances of a workflow
- Github⁴ allows easy workflow sharing and collaboration on development



We Need Software for Scientists

Galaxy⁶ delivers:

- A GUI to command line tools
- History of processes
- Construction of workflows
- Running workflows
- Integration with job schedulers
- Per-user management
- Extendable tool suites



Galaxy Example

Galaxy Analyze Data Workflow Shared Data Visualization Admin Help User Using 33%

Tools
search tools

- Get Data
- Send Data
- ENCODE Tools
- Lift-Over
- Text Manipulation
- Filter and Sort
- Join, Subtract and Group
- Convert Formats
- Extract Features
- Fetch Sequences
- Fetch Alignments
- Get Genomic Scores
- Operate on Genomic Intervals
- Statistics
- Wavelet Analysis
- Graph/Display Data
- Regional Variation
- Multiple regression
- Multivariate Analysis
- Evolution
- Motif Tools
- Multiple Alignments
- Meta-genomic analyses
- FASTA manipulation
- NGS: QC and manipulation
- NGS: Mapping
- NGS: Indel Analysis
- NGS: RNA Analysis
- NGS: SAM Tools
- NGS: GATK Tools (beta)
- NGS: Peak Calling
- NGS: Simulation
- SNP/WGA: Data: Filters
- SNP/WGA: QC: LD: Plots
- SNP/WGA: Statistical Models
- Phenotype Association
- VCF Tools
- BED Tools

Filtered Sequences 0
Sequence length 66
%GC 54

Per base sequence quality

Quality scores across all bases (Sanger / Illumina 1.9 encoding)

Position in read (bp)

Per sequence quality scores

Quality score distribution over all sequences

History
data 16: unmatched reads (L)
20: Bowtie2 on data 12 and data 17: aligned reads
19: Bowtie2 on data 12 and data 17: unaligned reads (L)
18: virus: PVS_TuMV_Macluravirus.fa
17: non-virus_contaminants.fa
16: blastn on data 15
15: Trinity on data 12: Assembled Transcripts
14: Trinity on data 12: log
13: FastQC on data 12
12: FASTQ Trimmer on data 1
11: FastQC on data 10
10: FASTQ Trimmer on data 1
9: FastQC on data 3
8: FastQC on data 2
7: FastQC on data 1
6: unmatched_R2.fq
5: unmatched_R1.fq
4: BC48_R2.fq
3: BC48_R1.fq
2: BC1_R2.fq
1: BC1_R1.fq

Why FLOSS?

- **Open:** Similar philosophy to scientific research
- **Current:** Keeps up with the scientific community
- **Community:** Collaboration, knowledge sharing
- **Flexible:** Adaptation to related problems
- **Trust:** Scientists do not trust what they cannot read/understand

References

1. www.openlava.org
2. www.ensembl.org
3. git-scm.com
4. github.com
5. <https://github.com/mfiers/Moa>
6. galaxyproject.org