# Discover. Innovate. Grow.™
## Hybrid Genome Assembly: A Practical Guide

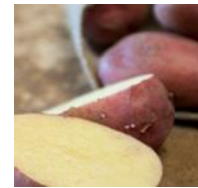**Cecilia Deng,** Senior Bioinformatician
邓泓

THE SCIENCE OF PREMIUM™

The New Zealand Institute for Plant & Food Research Limited

# Overview

- Genome Assembly Project
- Sequencing Platforms
- Sequencing Considerations & Experimental Design
- Hybrid Genome Assembly Strategies
- A Complete Genome Assembly Workflow
- Conclusions
- Acknowledgements

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI
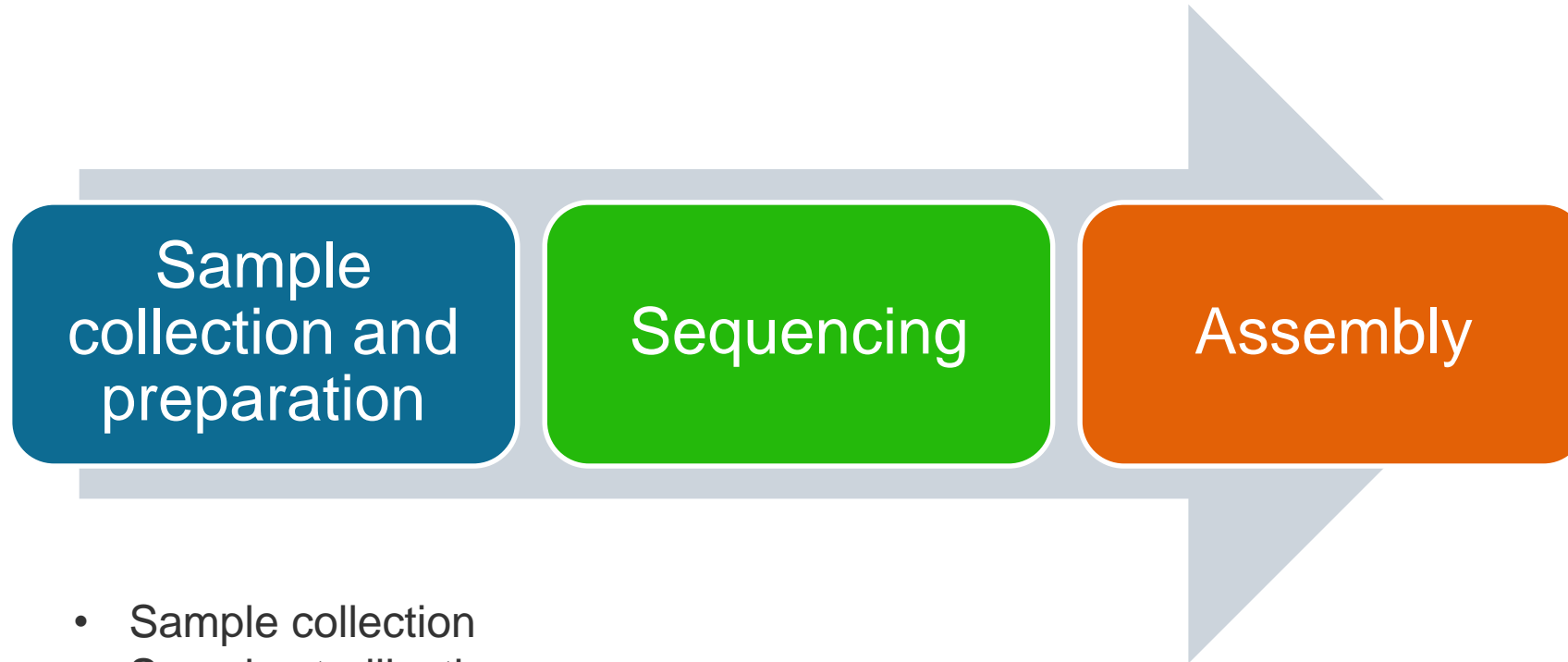
# Plant & Food Research (PFR)



**NEW ZEALAND**

- Kerikeri
- Aucklan
- Pukekohe
- Ruakura
- Te Puke
- Hawke's Bay
- Palmerston North
- Nelson
- Motueka
- Wellington
- Blenheim
- Lincoln
- Clyde
- Dunedin
- Gore

**AUSTRALIA**

- Brisbane
- Adelaide
- Albury

**USA**

- Davis, California

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

3. 测序:

乙方保证采用 ███████ 测序平台,共测序 40 Gb 的数据.

## 四、报酬及其支付方式

（一）项目报酬

对于本协议包括的乙方需要完成的所有技术服务工作，甲方需要按服务内容向乙方支付的技术服务报酬为 ████████（大写：人民币 ████████ 圆整）。

# Genome Assembly Project

**Sample collection and preparation** → **Sequencing** → **Assembly**

- Sample collection
- Sample sterilization
- gDNA extraction
- gDNA QC
- Transportaion
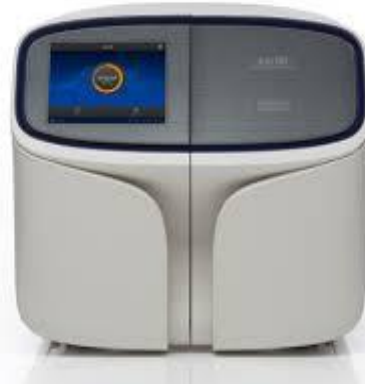- Sequencing library construction
- Library QC

# Sequencing Platforms

|  | NovaSeq 6000 | HiSeq X Ten | HiSeq 4000 |
|---|---|---|---|
| Output Range | 167–6000 Gb | 900–1800 Gb | 125–1500 Gb |
| Run Time | 19–40 hr | <3 days | 3.5 days |
| Reads per Run | 1.4–20 billion | 3–6 billion | 2.5 - 5 billion |
| Maximum Read Length | 2 × 150 bp | 2 × 150 bp | 2 × 150 bp |
| Samples per Run | 4–48 | 8–16 | 6 - 12 |
| Relative Price per Sample | Higher Cost | Lower Cost | Mid Cost |

Short read:
- High accuracy
- Deep coverage
- Cheap

- Illumina SLR

X Ten is NOT 10X !

# Sequencing Platforms

|  | Sequel | RS II |
|---|---|---|
| **Average read length** | 10 - 15 Kb | 10 Kb |
| **Throughput per cell** | ~5 - 10 Gb | 500 Mb ~ 1 Gb |
| **SMRT Cells per run** | 1 - 16 | 1 - 16 |
| **Movie lengths per SMRT Cell** | 30 mins - 6 hrs | 30 mins - 6 hrs |


PACBIO®

- Single Molecular Real Time
- Long read length
- No PCR, less bias
- Higher error rate
- More expensive
- Read length distribution

MinION

SmidgION

**ONT**


Oxford NANOPORE Technologies

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

1304KHS-0092/CKVAopDA4_2.fastq.gz

D1B7JACXX_SD1_29_3_GCCAAT_L002_R1_001.fastq.bz2

C3BC5ACXX_NoIndex_L004_R1.fastq.gz

**PE?**
**MP?**
**Hi-C?**
**10X?**
**RNA?**
**GBS?**

**?**

**Genome Size?**

**Ploidy?**

**Heterozygosity?**

```
@HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
TTAATTGGTAAATAGTGCTTAGATNTTACCTTNNNNNNNNNTAGTTTCTTGAGATTTGTTGGGGGAGACATTTTTGTGATTGCCTTGAT
+HWI-EAS209_0006_FC706VJ:5:58:5894:21141#ATCACG/1
efcfffffcfeefffffdfeceeffffff|fc]_[YBBBBBBBBBBRTT\]][]dddd`ddd^dddadd^BBBBBBBBBBBBBBB_BB_BfB
```

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

3. 测序:

乙方保证采用 ▮▮▮▮▮▮ 测序平台, 共测序 40 Gb 的数据.

## 四、报酬及其支付方式

（一）项目报酬

对于本协议包括的乙方需要完成的所有技术服务工作，甲方需要按服务内容向乙方支付的技术服务报酬为 ▮▮▮▮（大写：人民币 ▮▮▮▮▮▮ 圆整）。|

# Sequencing Considerations & Experimental Design



**Illumina**
- PE 300bp? PE 600bp?
- MP 8KB? MP 12KB?
- Sequencing depth?

**PacBio or ONT?**
- Bluepinning?
- Sequencing depth?

**Illumina SLR or 10X?**

**Bionano mapping?**

**Hi-C on Illumina?**

**Genetic maps available?**

Sample collection and preparation → Sequencing → Assembly

Plant & Food RESEARCH
RANGAHAU AHUMĀRA KAI

https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4231593/

Plant & Food
**RESEARCH**
RANGAHAU AHUMĀRA KAI

- OLC: Overlap-Layout-Consensus
  - Suitable for long reads
  - Newbler, Celera Assembler, PCAP, etc.



Amplified DNA

Shear DNA

Sequenced reads

Overlaps

Layout

Consensus → "Contigs"

ATGTTCCGATTAGGAAACCTATCTGTAACTGTTTCATTCAGTAAAAGGAGGAAATATAA

Plant & Food
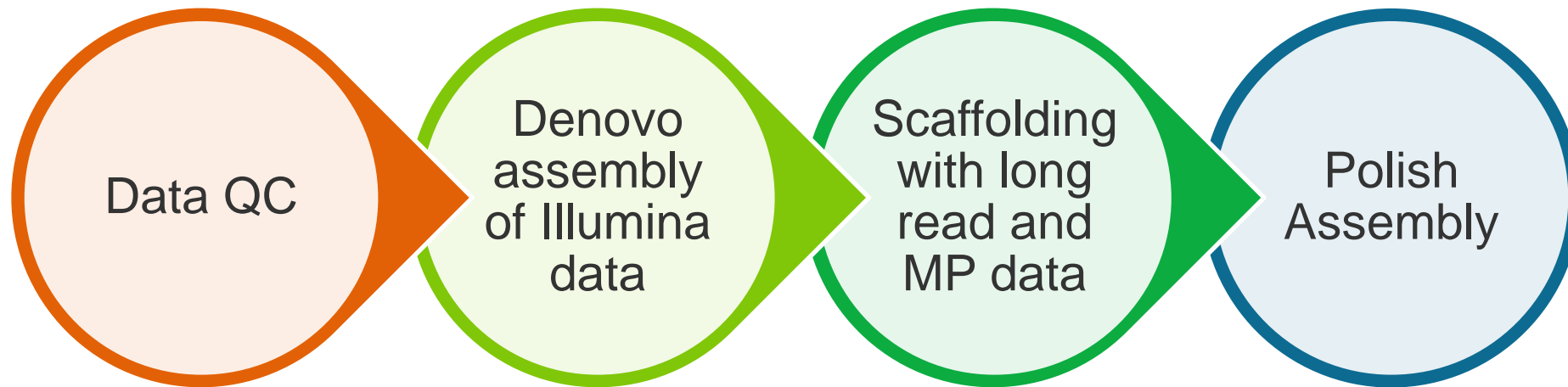RESEARCH
RANGAHAU AHUMĀRA KAI

Limited

- DBG: De Bruijn Graph
  - An *n*-dimensional **De Bruijn graph** of *m* symbols is a <u>directed graph</u> representing overlaps between sequences of symbols
  - Each read is broken into fixed-size *k*-mers. A graph is directly constructed where each vertex is a *k*-mer and each edge indicates two adjacent *k*-mers overlapping by *k* – 1 letters.
  - Suitable for short reads
  - Velvet, AllPath-LG, ABySS, etc.

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

- SG: String Graph
  - In graph theory, a **string graph** is an intersection graph of curves in the plane; each curve is called a "string". Given a graph *G*, *G* is a string graph if and only if there exists a set of curves, or strings, drawn in the plane such that no three strings intersect at a single point and such that the graph having a vertex for each curve and an edge for each intersecting pair of curves is isomorphic to *G*.
  - Falcon

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

# Hybrid Genome Assembly Strategies: Option 1

**Data QC** → **Denovo assembly of Illumina data** → **Scaffolding with long read and MP data** → **Polish Assembly**
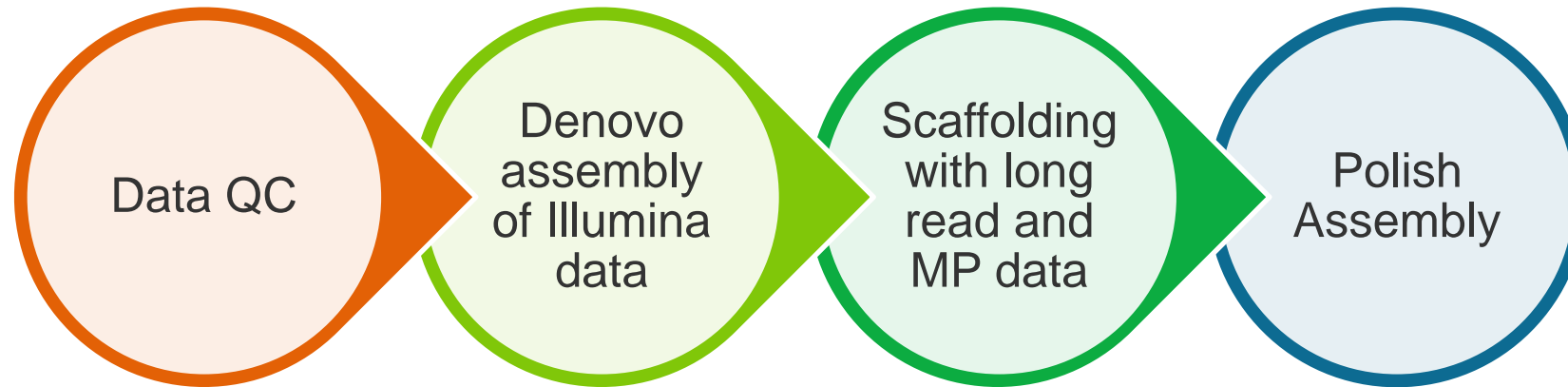
- Data integrity (md5sum)
- FastQC and MultiQC
- Fastp
- Trimmomatic
- Trim Galore!
- ErrorCorrection

- Velvet
- ALLPATHS-LG
- ABySS
- SOAPdenovo2
- MIRA
- SGA
- And many more

- SSPACE
- SOAPdenovo2
- SOPRA

- PBJelly
- OPERA-LG

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

Data QC → Denovo assembly of Illumina data → Scaffolding with long read and MP data → Polish Assembly

THE ASSEMBLATHON

Genome Assembly Gold-Standard Evaluations

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

Data QC

Denovo assembly of long reads

Error correction with short read PE data

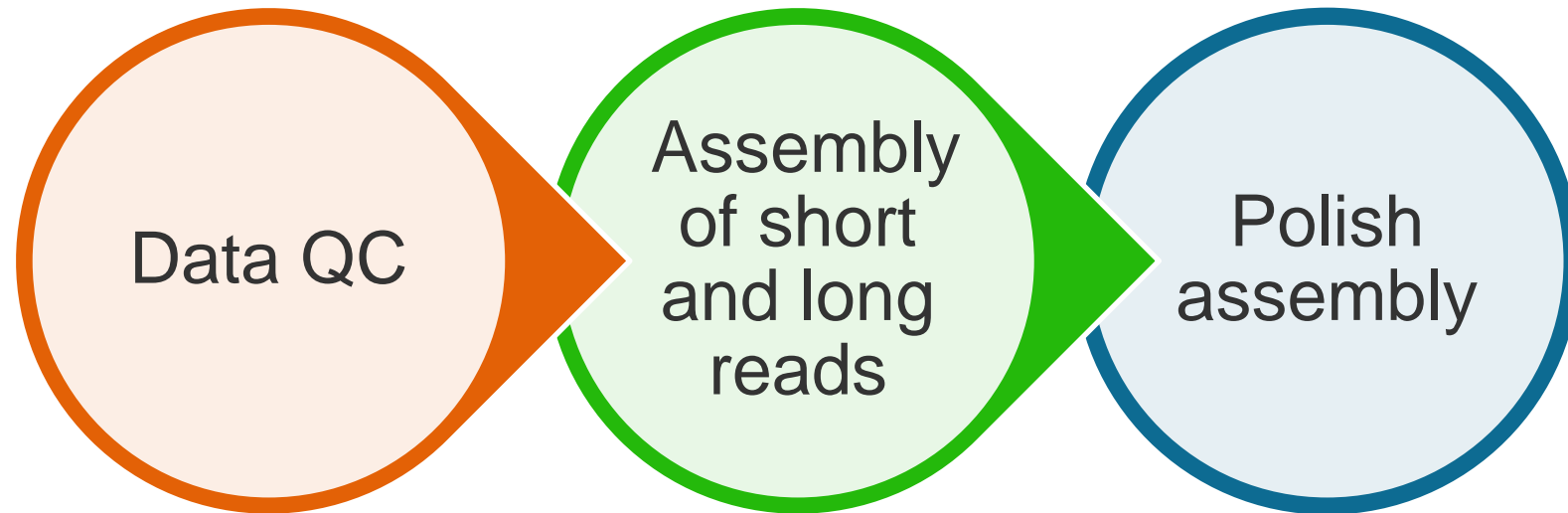Polish assembly

- Canu
- Falcon and Falcon-Unzip
- HGAP4
- Celera Assembler
- mugqic/genpipes/ pacbio_assembly

- Arrow/Quiver
- Pilon

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

Data QC

Assembly of short and long reads

Polish assembly

- MaSuRCA
- Spades
- MIRA
- CABOG

Plant & Food
**RESEARCH**
RANGAHAU AHUMĀRA KAI

# Hybrid Genome Assembly

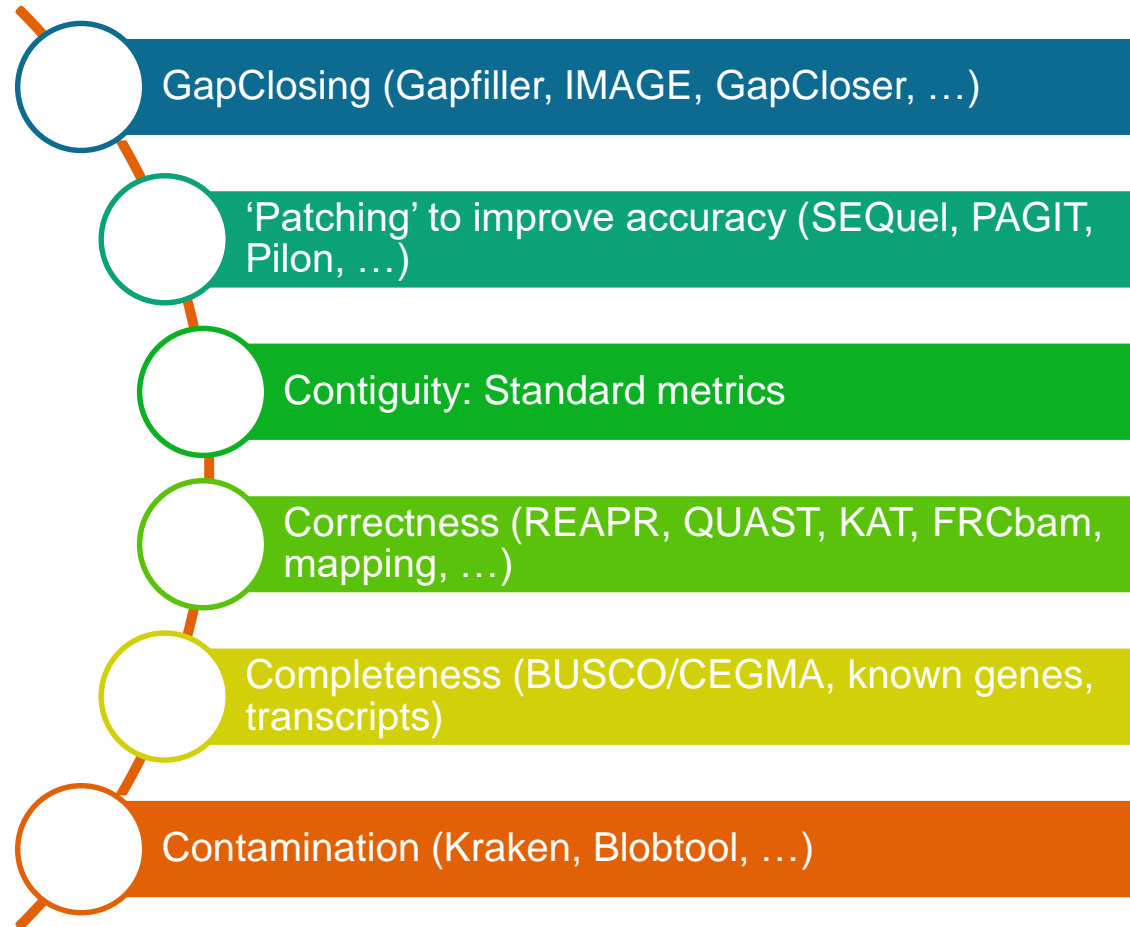| Assembler | Algorithm | Input |
|-----------|-----------|-------|
| Arachne | OLC | Sanger |
| CAP3 | OLC | Sanger |
| TIGR | Greedy | Sanger |
| Newbler | OLC | 454/Roche |
| Edena | OLC | Illumina |
| SGA | OLC | Illumina |
| MaSuRCA | De Bruijn/OLC | Illumina/PacBio |
| MIRA | De Bruijn/OLC | Illumina/PacBio/454/Sanger |
| Velvet | De Bruijn | Illumina |
| ALLPATHS | De Bruijn | Illumina/PacBio |
| ABySS | De Bruijn | Illumina |
| SOAPdenovo | De Bruijn | Illumina |
| Spades | Paired De Bruijn | Illumina/PacBio |
| CLC | De Bruijn | Illumina/454 |
| CABOG | OLC | Hybrid |
| Falcon | String graph | PacBio |
| StriDe | String graph + De Brujin | Illumina |

- Every species has it's own surprises and characters
- Every sequencing chemistry has it's strengths and weaknesses
- Every assembler has it's own set of heuristics.

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

# An Example Workflow To Assemble A Bacteria Genome



http://sepsis-omics.github.io/tutorials/modules/cmdline_assembly/

# Assembly Assessment and Improvement At Scaffolds Level

Polish assembly

GapClosing (Gapfiller, IMAGE, GapCloser, …)

'Patching' to improve accuracy (SEQuel, PAGIT, Pilon, …)

Contiguity: Standard metrics

Correctness (REAPR, QUAST, KAT, FRCbam, mapping, …)

Completeness (BUSCO/CEGMA, known genes, transcripts)

Contamination (Kraken, Blobtool, …)

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI
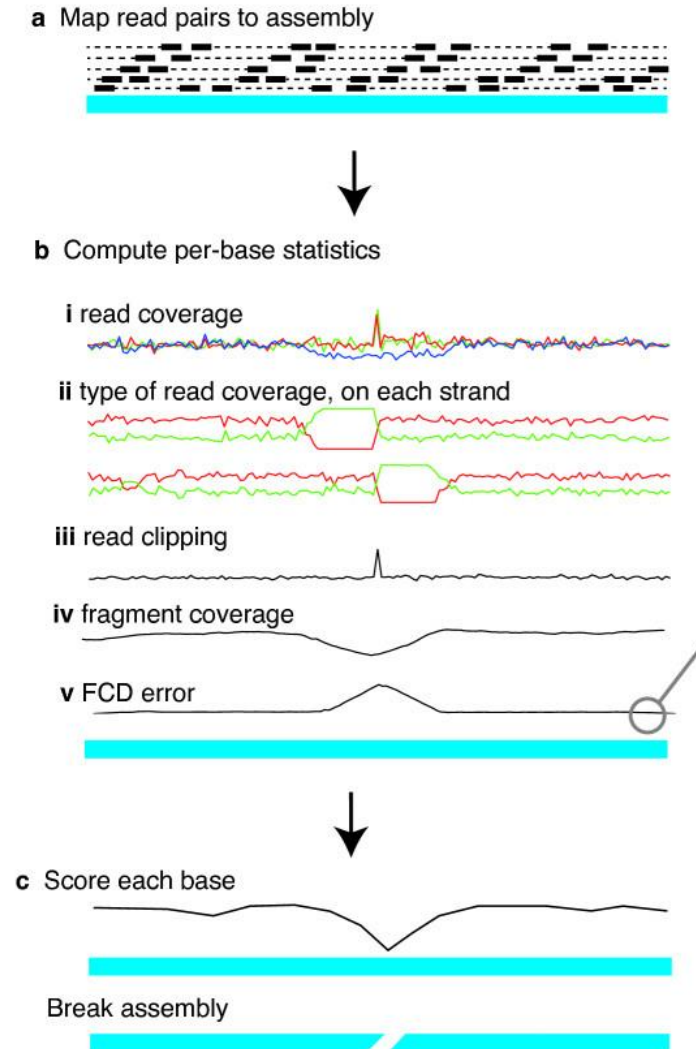
- Standard metrics
  - Assembled size, # of contigs, # of scaffolds, N50, size of the longests contig, size of the longest scaffold, etc.
- N50: The  length of the longest sequence such that the sum of sequences longer than it reaches half of the assembled size
- NG50: The  length of the longest sequence such that the sum of sequences longer than it reaches half of the genome size

a Map read pairs to assembly

b Compute per-base statistics

i read coverage

ii type of read coverage, on each strand

iii read clipping

iv fragment coverage

v FCD error

c Score each base

Break assembly

- Uses the same principle of feature response curve (FRC)
- Captures trade-off between quality and contiguity
- Identifies erroneous positions
- Breaks sequences at suspicious positions

Plant & Food
**RESEARCH**
RANGAHAU AHUMĀRA KAI

```
# BUSCO was run in mode: genome

        C:89.3%[S:68.6%,D:20.7%],F:3.0%,M:7.7%,n:1440

        1286      Complete BUSCOs (C)
        988       Complete and single-copy BUSCOs (S)
        298       Complete and duplicated BUSCOs (D)
        43        Fragmented BUSCOs (F)
        111       Missing BUSCOs (M)
        1440      Total BUSCO groups searched
```
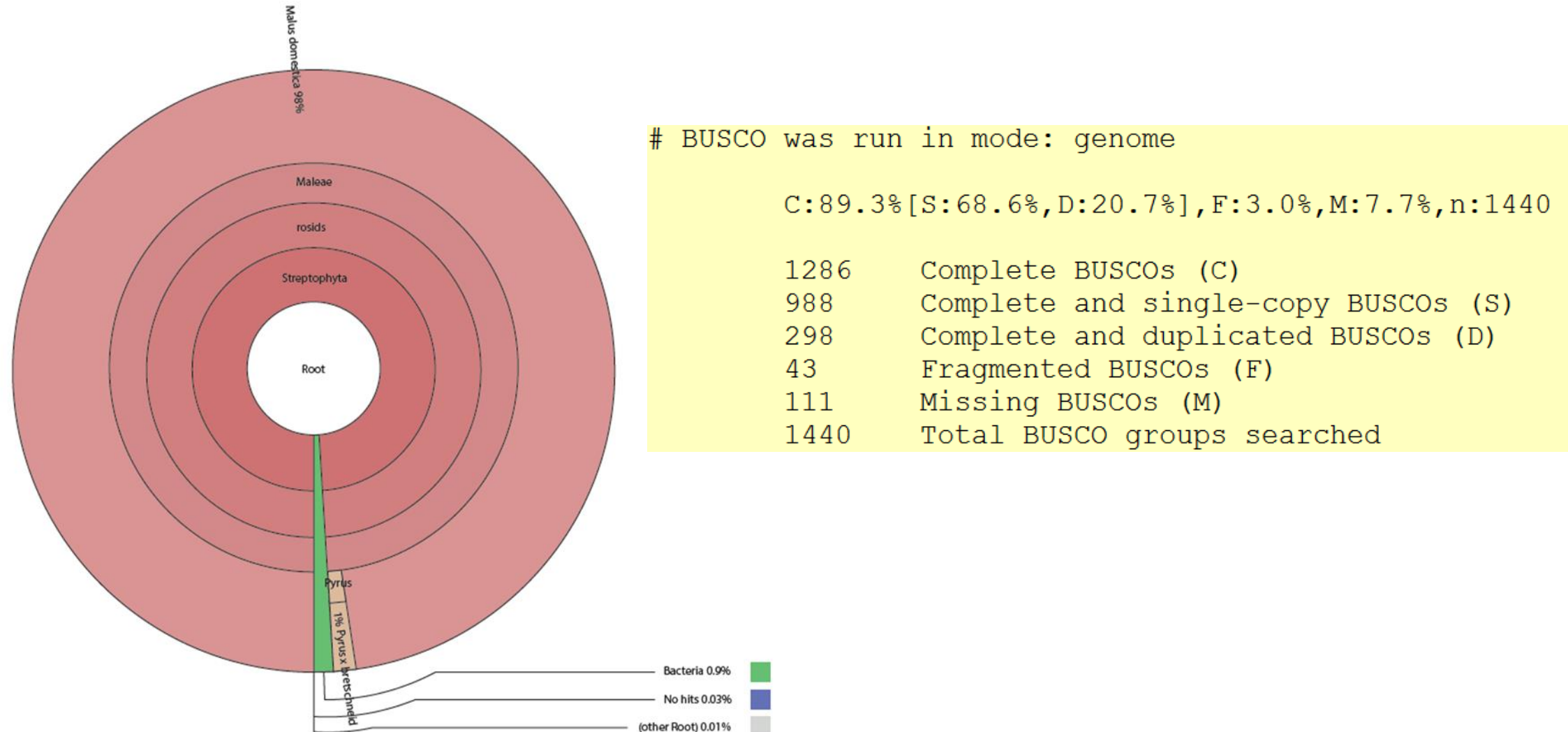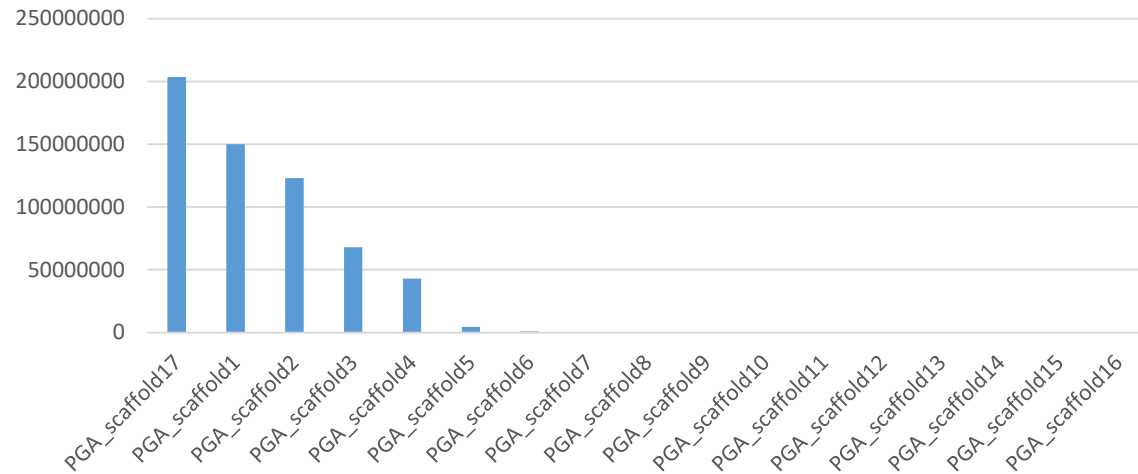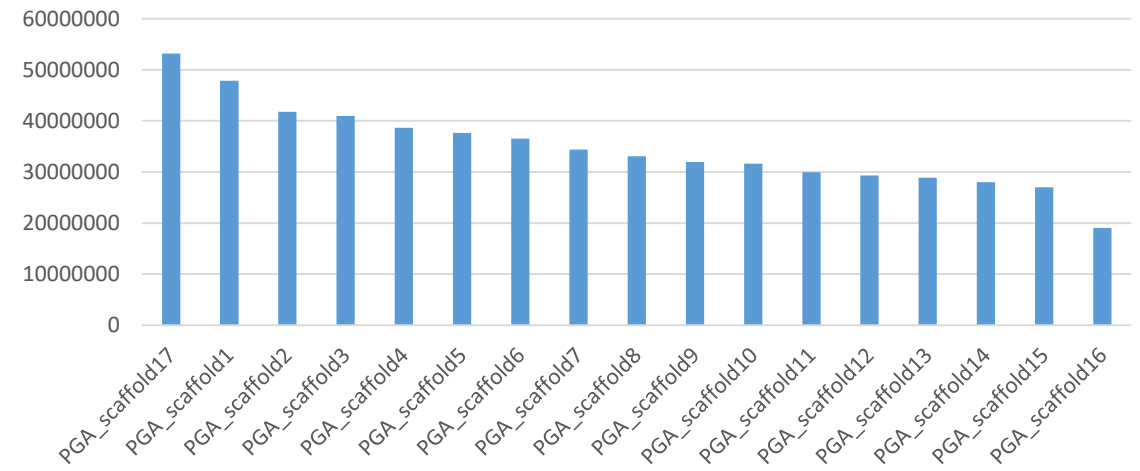
**Figure 4:** 'Royal Gala' assembly C3: Scaffold classification and contamination check.

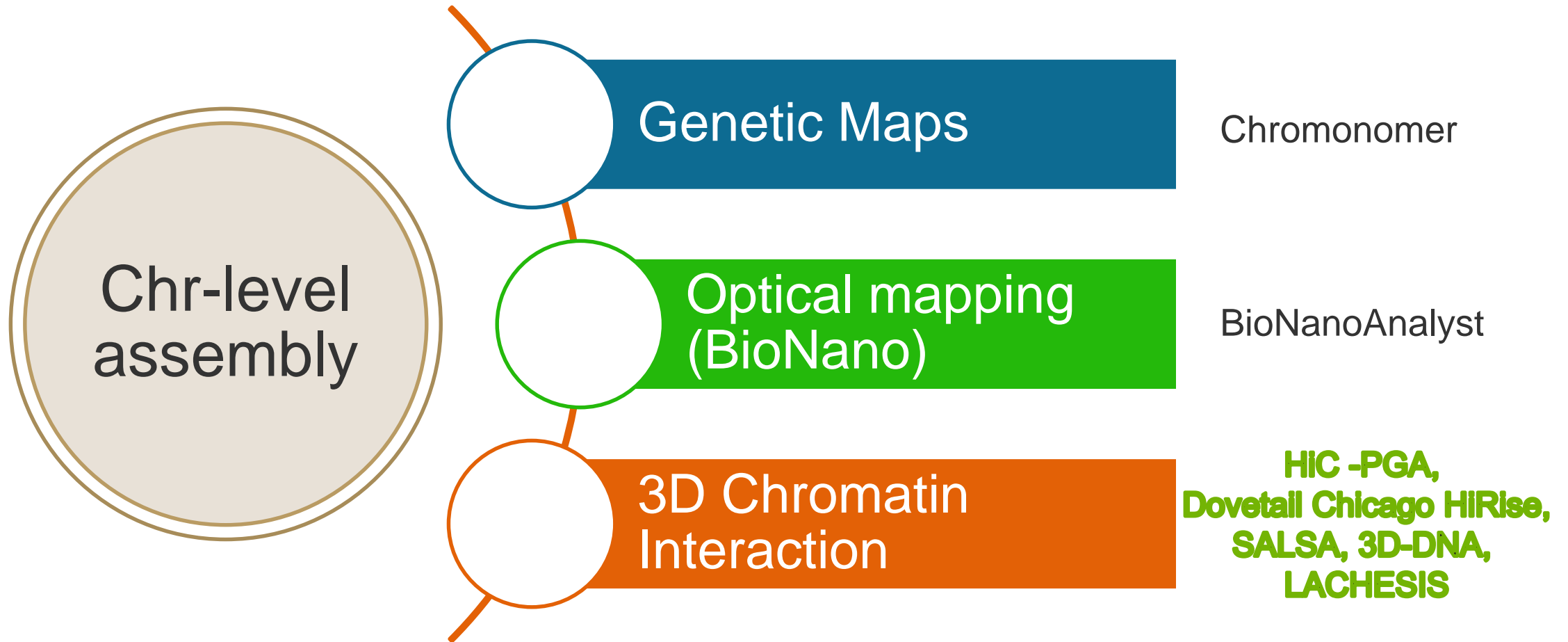Assembly Assessment: Completeness and Contamination Check

- Example 2: Puccinia co.
- Example 3: Puccinia tr.

# Genome Assembly Post-Processing



Repeat Detection and Genome Masking

Gene Prediction and Functional Annotation

Visualization and Release

Publish to Public Domain

Genome Annotation

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

# A Complete Genome Project Workflow



Experimental Design

Sample collection and preparation

Sequencing

Assembly

Applications:
- Pan genome construction
- Comparative genomics
- Gene mining
- Genotyping and marker discovery
- Population genetics
- Evolution
- GWAS
- GS

Validation and Annotation

Visualization and Release

Publish to public domain

Plant & Food
RESEARCH
RANGAHAU AHUMĀRA KAI

# Conclusions

» Genome assembly can be complicated

» Experimental design is critical

» 4Cs: Contiguity/Correctness/Completeness/Contamination

» Assemblies are not perfect

  » Species specific difficulties (repeat, polymorphism, ploidy)

  » Sequencing chemistry

  » Regions not clone/sequence/assemble well

  » Software heuristics

» Know when to stop!

**Time frame**

**Funds**        **Quality**

## Bioinformatics landscape changes fast

Plant & Food
**RESEARCH**
RANGAHAU AHUMĀRA KAI

# Acknowledgements

» Pipfruit Breeding Team
» Mapping and Markers Team
» Molecular Biology

» Bioinformatics Team
» Flavour Team

» Kiwifruit Breeding Team

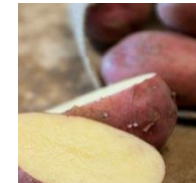**Zespri** Kiwifruit

**Prevar** PREMIUM APPLE & PEAR PRODUCTS

**SlipStream Automation**

Plant & Food **RESEARCH** RANGAHAU AHUMĀRA KAI

The New Zealand Institute for Plant & Food Research Limited

# Plant & Food Research (PFR)



NEW ZEALAND
- Kerikeri
- Auckland
- Pukekohe
- Te Puke
- Ruakura
- Palmerston North
- Hawke's Bay
- Nelson
- Wellington
- Motueka
- Blenheim
- Lincoln
- Clyde
- Dunedin
- Gore

AUSTRALIA
- Brisbane
- Adelaide
- Albury

USA
- Davis, California

plantandfood.co.nz

Cecilia.Deng@plantandfood.co.nz

OUR SCIENCE IS GROWING FUTURES™

The New Zealand Institute for Plant & Food Research Limited