# Reads2Map

**Developing best practices for genotyping-by-sequencing analysis using linkage maps as benchmarks**

Cristiane Taniguti

chtaniguti@tamu.edu

TEXAS A&M UNIVERSITY

NC STATE UNIVERSITY

STATISTICAL GENETICS LAB · ESALQ / USP

TOOLS FOR POLYPLOIDS

# Motivation: GBS data



Issues while building linkage maps

► Computational intensive
► Time consuming
► Wrong grouping
► Wrong ordering
► Inflated linkage maps

OneMap

- Since 2007 - 52k downloads
- Diploid species
- Bi-parental populations
- Backcross, RILs, F2 and outcrossing
- Biallelic and Multiallelic markers

Maintainer since 2017
Updates in version 3.0

STATISTICAL GENETICS ·LAB·
ESALQ / USP

Augusto Garcia
Marcelo Mollinari
Gabriel Margarido

TOOLS FOR POLYPLOIDS

# Motivation: GBS data

Issues while building linkage maps

- ► Computational intensive
- ► Time consuming
- ► Wrong grouping
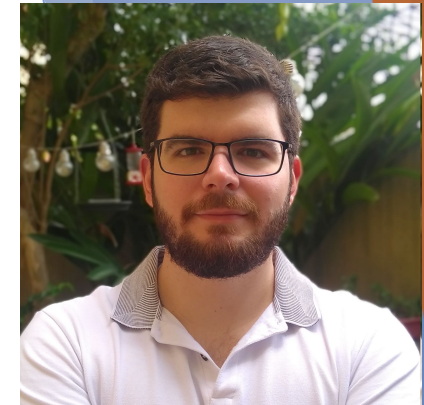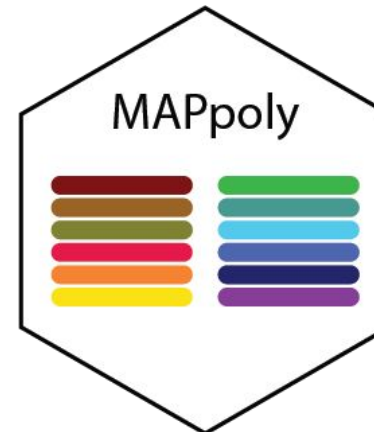- ► Wrong ordering
- ► Inflated linkage maps

Tips and Tricks:



RGC11 - Poster



Augusto Garcia



Gabriel Gesteira



MAPpoly

- ● Since 2018
- ● Diploid and polyploid species
- ● Bi-parental populations
- ● Outcrossing
- ● All dosages markers
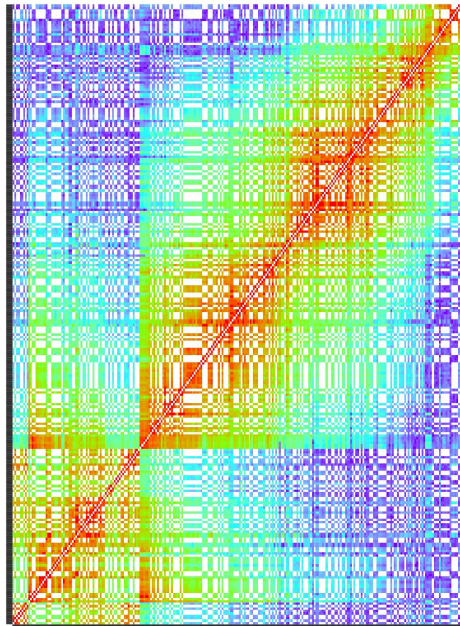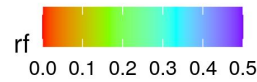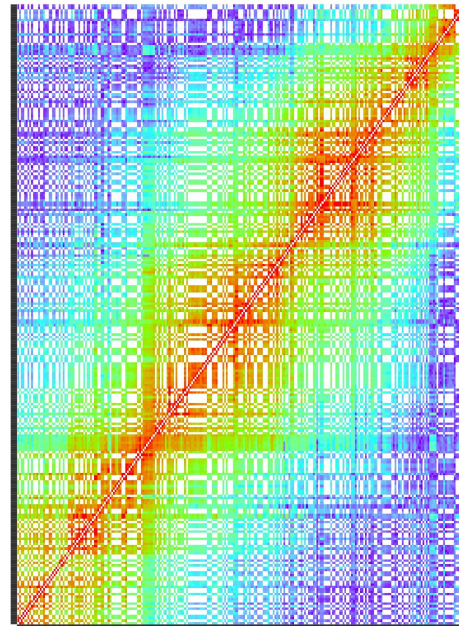- ● Updates by Marcelo and Gabriel



Marcelo Mollinari

TOOLS FOR
POLYPLOIDS

# Recombination fraction matrix as benchmarks

## Inversion

- Aspen chromosome 12



rf  0.0  0.1  0.2  0.3  0.4  0.5

Genomic order

Order fixed
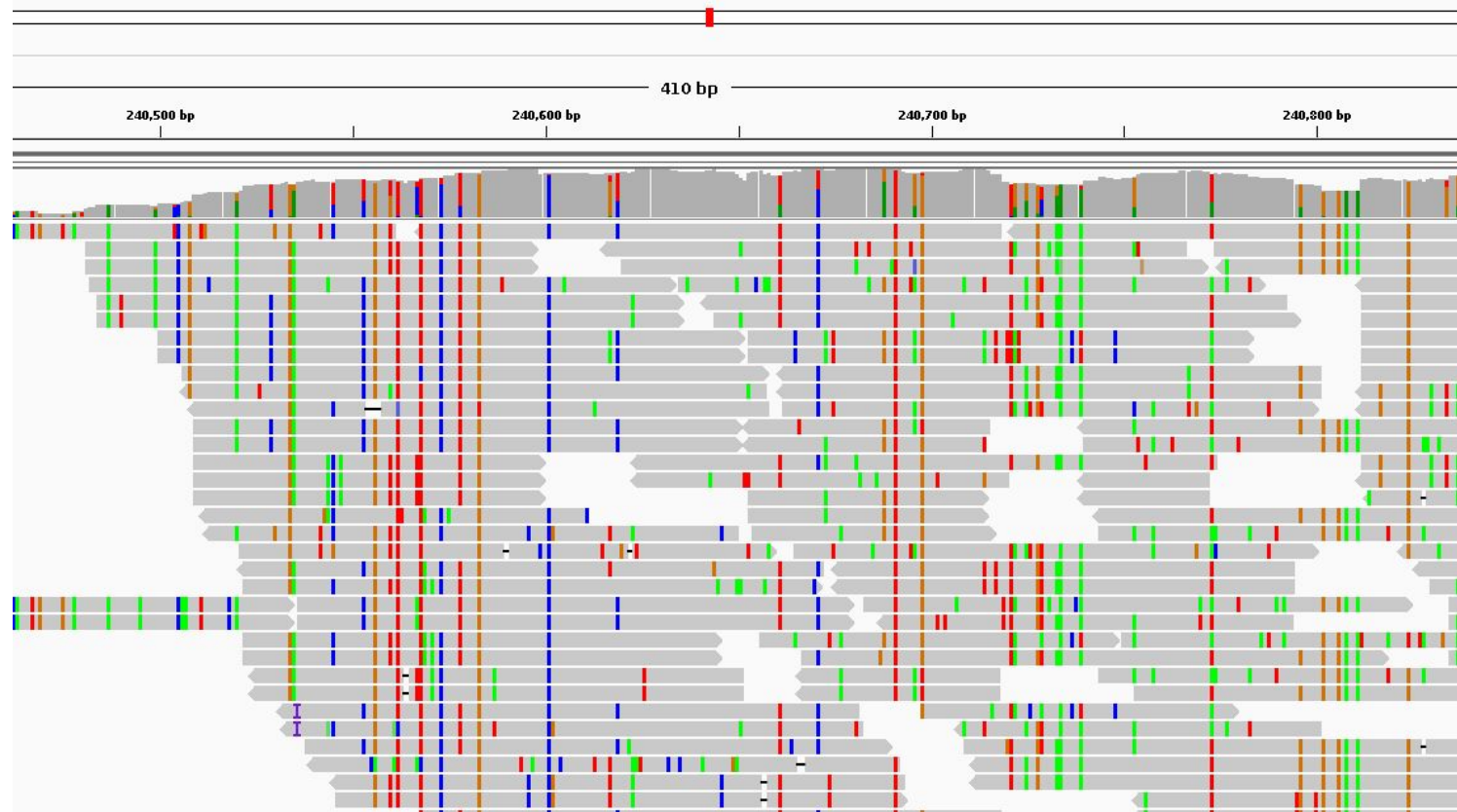
Most likely bad data

TOOLS FOR
POLYPLOIDS

# GBS Overview
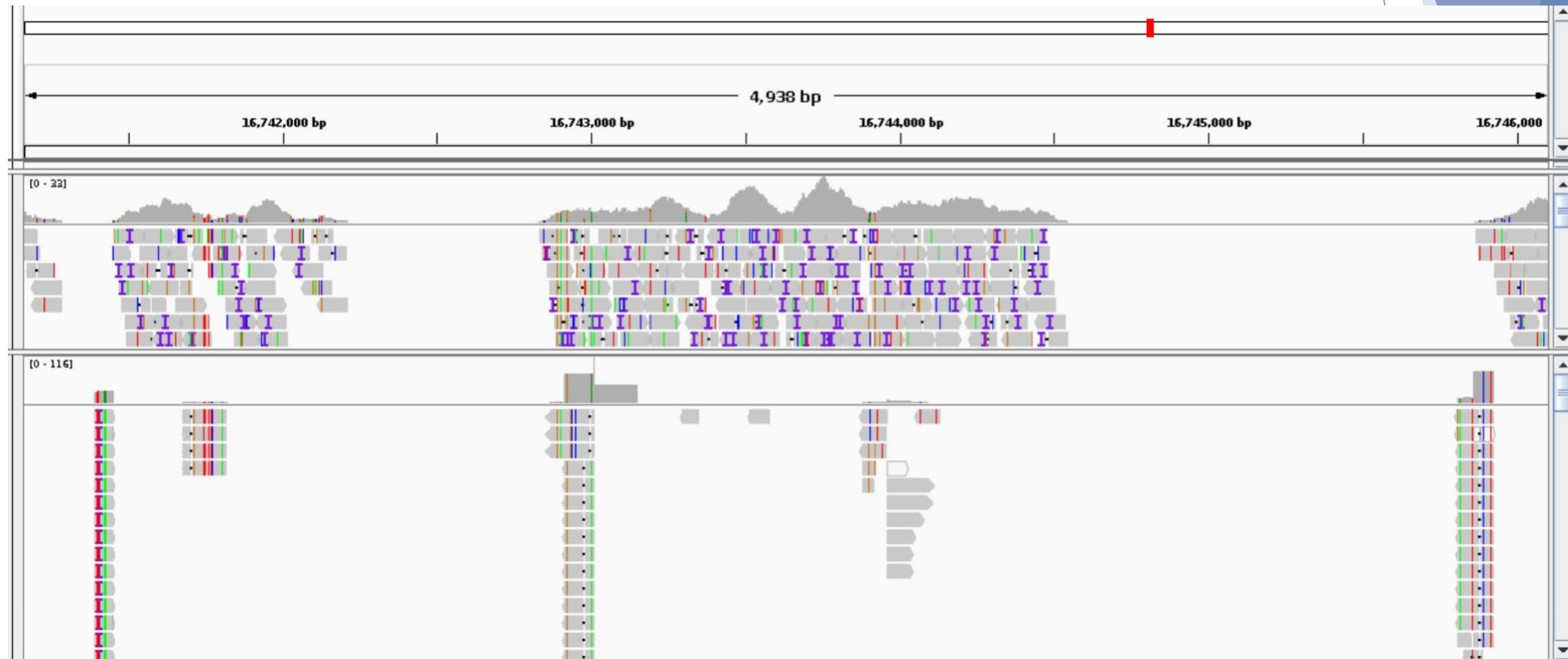
# SNP Calling

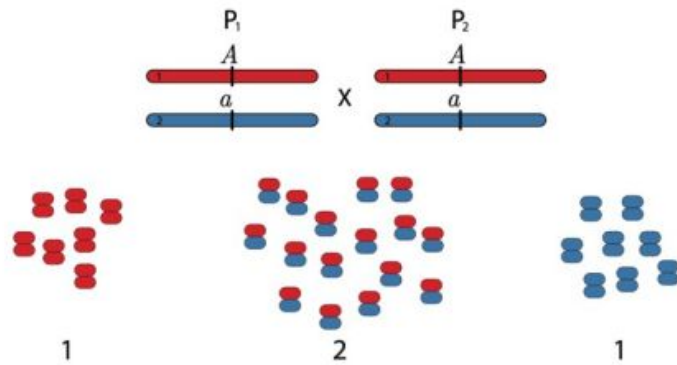► Whole Genome Sequencing (WGS)

Image: IGV

# SNP Calling

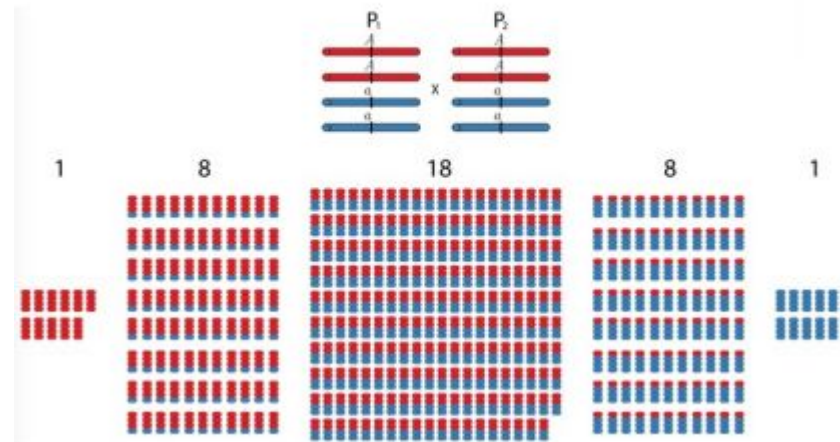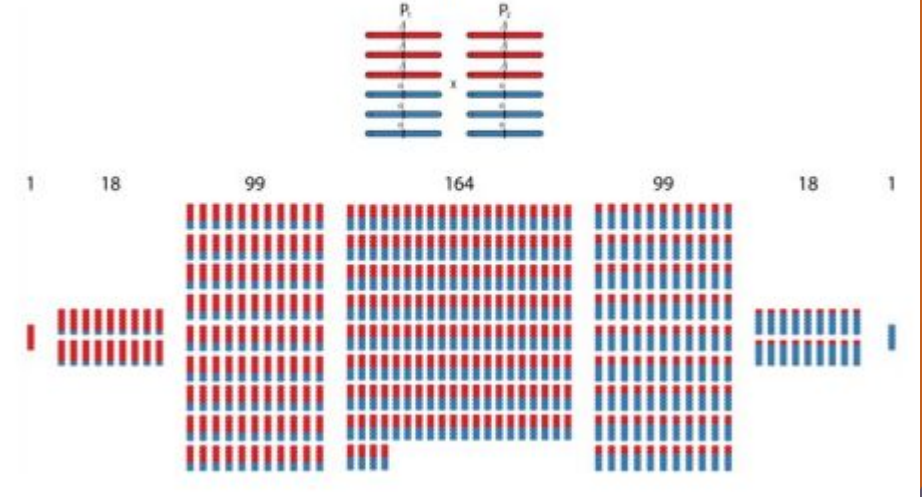- ► Exome sequencing (top) and Genotyping-by-Sequencing (bottom)

Image: IGV

# Dosage calling



Diploid

Tetraploid

Hexaploid

# Dosage Calling
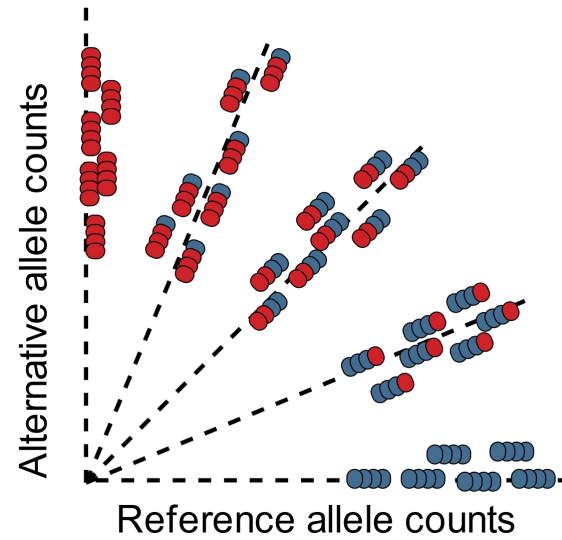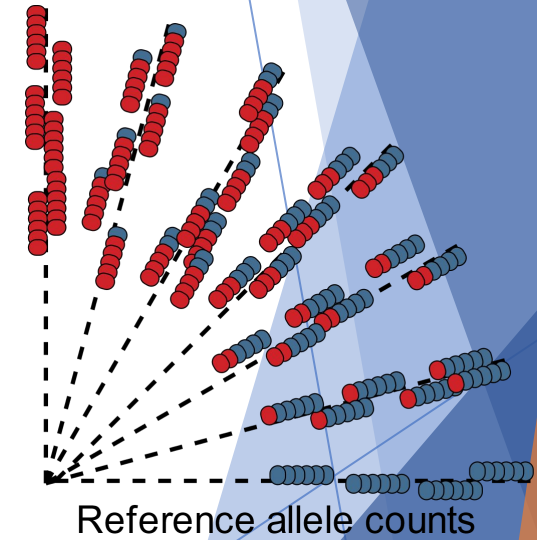
► The theory

### Diploid
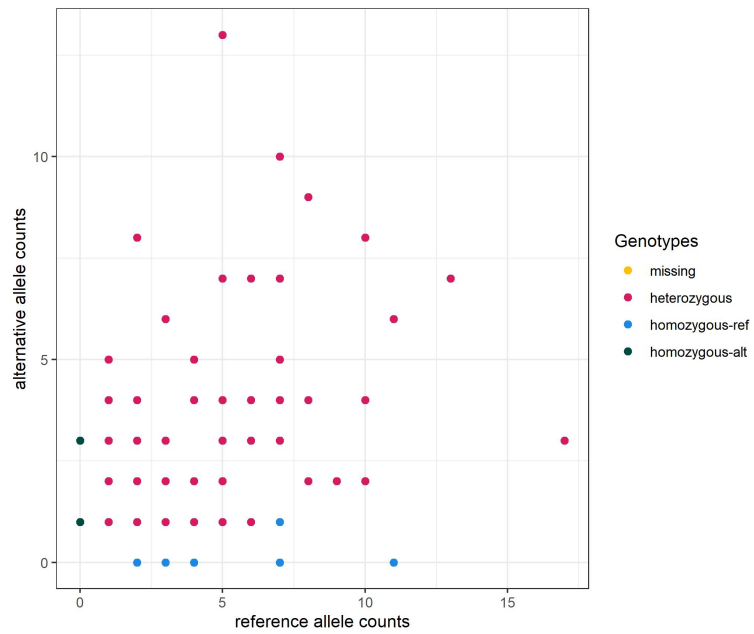


### Tetraploid



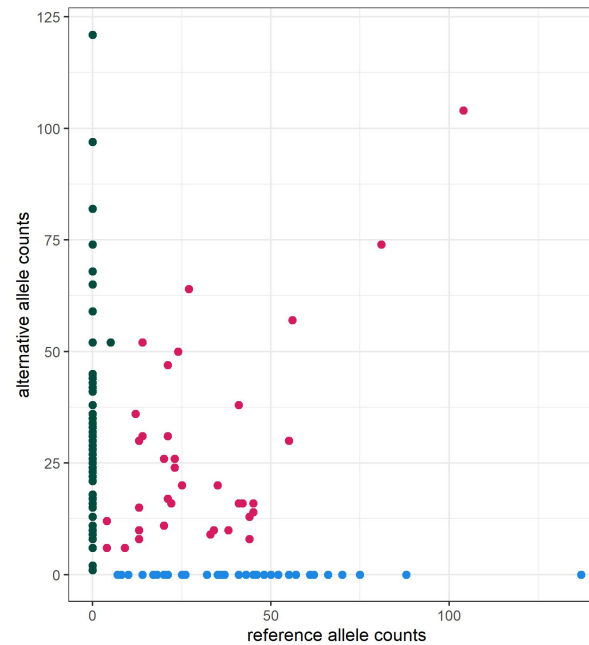### Hexaploid



TOOLS FOR POLYPLOIDS

# Dosage Calling

► The reality
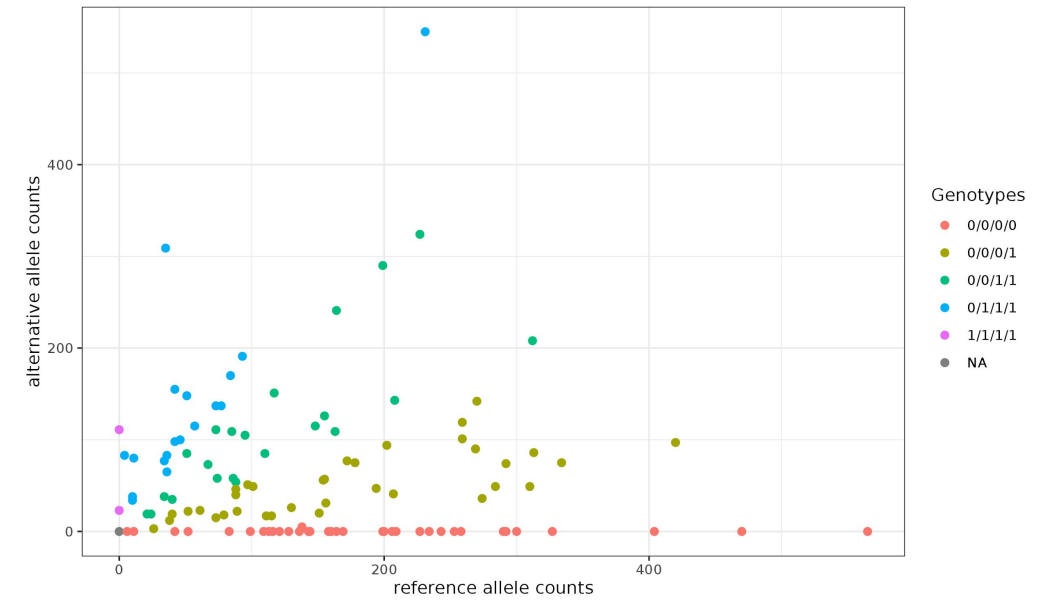
Diploid (mean depth 6)
N = 200
Aa x Aa

Diploid (mean depth 96)
N = 138
Aa x Aa

Tetraploid (mean depth 83)
N = 114
AAaa x AAaa



TOOLS FOR
POLYPL⚪IDS

# Sequencing Data - Technical Difficulties

- ► Large files

- ► Many software
- ► Many programming languages
- ► Different Operational Systems
- ► Updates

**Quality Control**
cutadapt
STACKS  FASTQC
MULTIQC
samtools
Others

**Alignment**
Bowtie
BWA  HISAT2
STAR
Others

**SNP Calling**
freebayes
GATK  samtools
TASSEL  STACKS
VarScan
Others

**Dosage Calling**
updog
SuperMASSA
GATK  polyRAD
freebayes
Others

**Files Manipulation**
picard
bedtools  tabix
samtools  GATK
bcftools
Others

TOOLS FOR
POLYPLOIDS

# Sequencing Data - Technical Difficulties

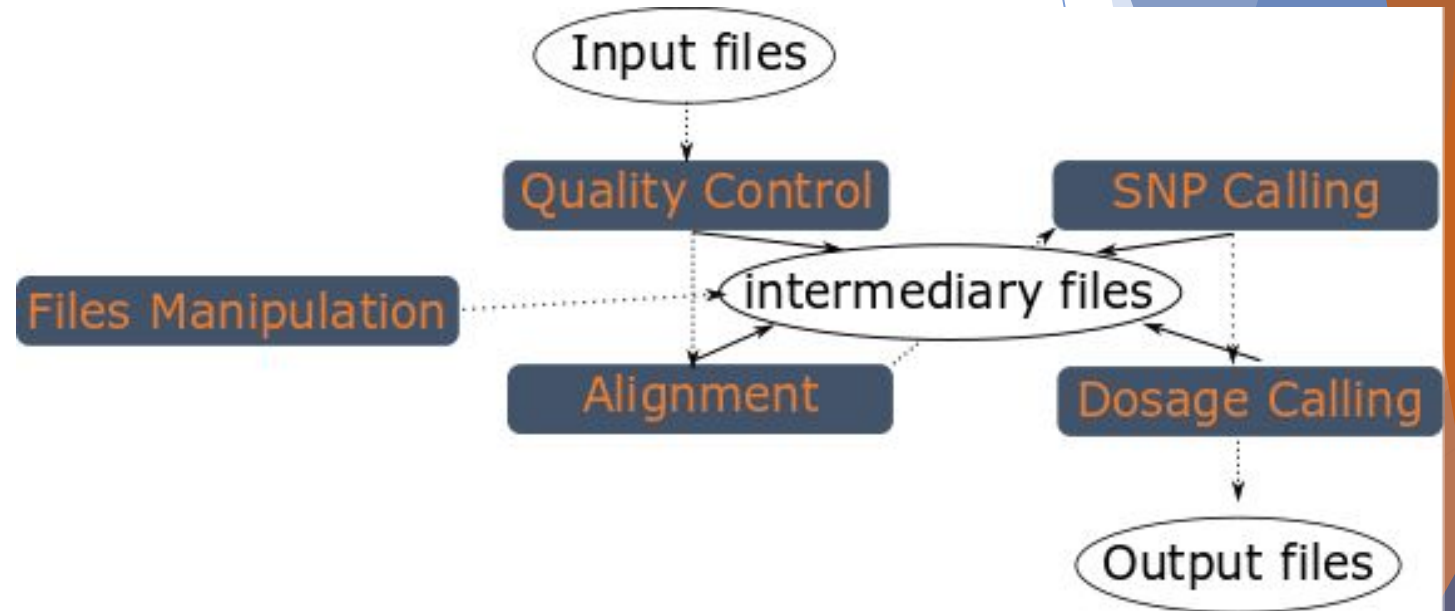- ► Large files
  - ► High Performance Computing (HPC)
  - ► Management systems (SLURM, SGE)
  - ► Cloud (Google, Amazon)

- ► Many software
- ► Many programming languages
- ► Different Operational Systems
- ► Updates
  - ► Containers
    - ► Docker
    - ► Singularity (usually available in HPC)
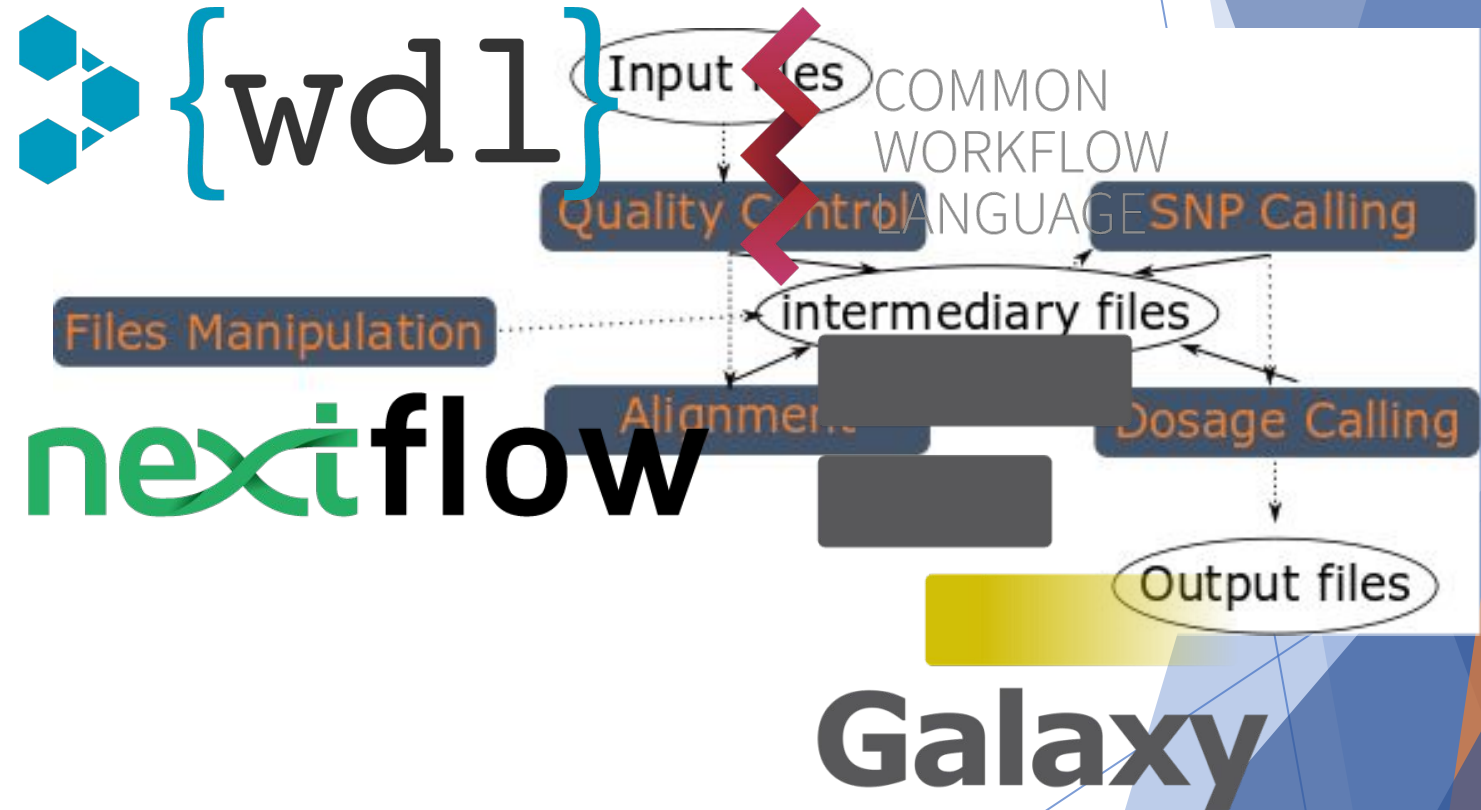    - ► BioContainers



TOOLS FOR
POLYPLOIDS

# Sequencing Data - Technical Difficulties

► Many steps
► Many file formats

# Sequencing Data - Technical Difficulties

- ► Many steps
- ► Many file formats
  - ► Workflows systems
    - ► Galaxy
    - ► Nextflow
    - ► Snakemake
    - ► CWL
    - ► WDL
  - ► Workflows repositories
    - ► Dockerstore
    - ► WorkflowHub
  - ► Run workflows on Cloud
    - ► Galaxy
    - ► DNAnexus
    - ► Terra
    - ► AnVIL
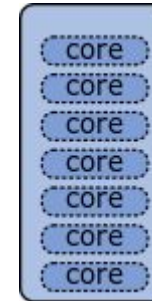    - ► SevenBridges



TOOLS FOR
POLYPLOIDS

# Sequencing Data - Technical Difficulties

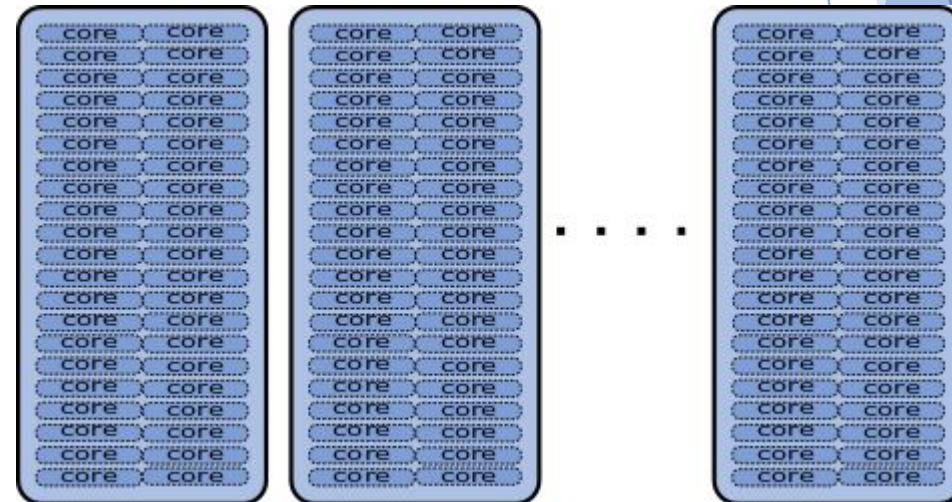- ► Resources optimization
  - ► Time
  - ► Cores
  - ► Nodes
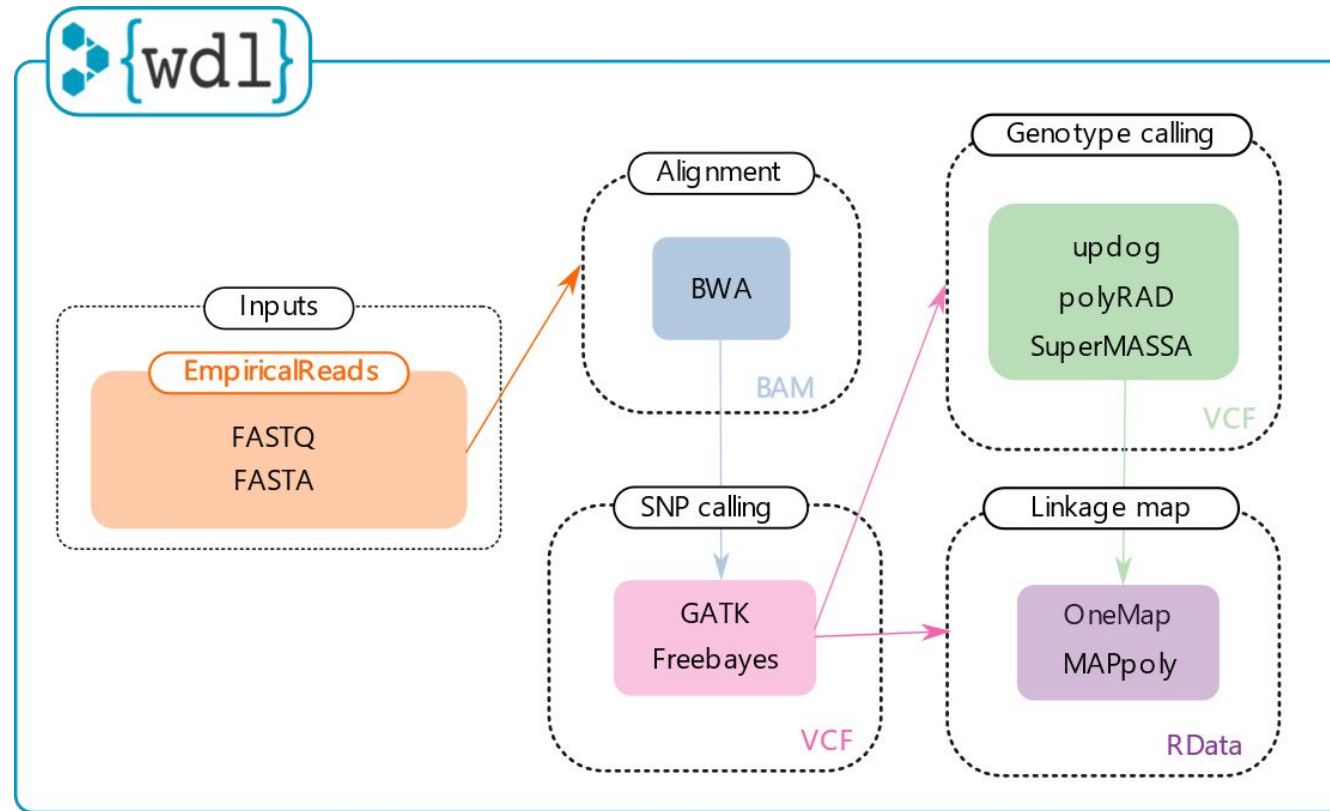  - ► RAM memory

Personal Computer:

4GB RAM; 8 cores; 1 node

High Performance Computing (Texas A&M):

384GB; 48 cores per node; 900 nodes

TOOLS FOR
POLYPLOIDS

# How we solved it: Reads2Map
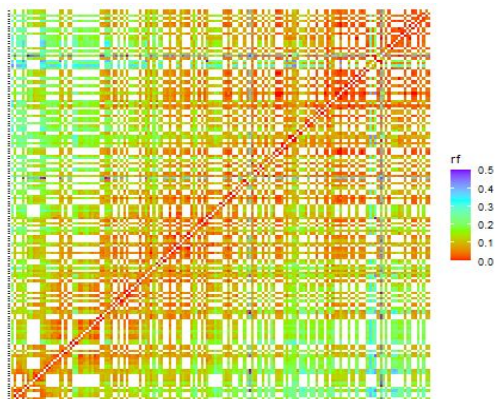


Available in Github, Dockerstore and WorkflowHub
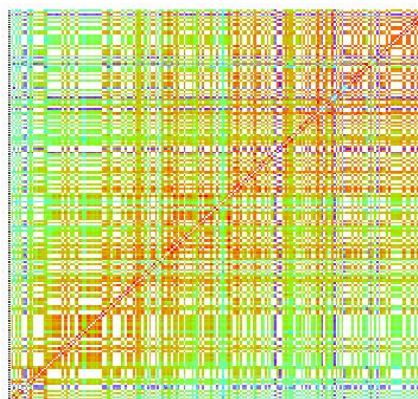
# Reads2Map results - Diploid roses
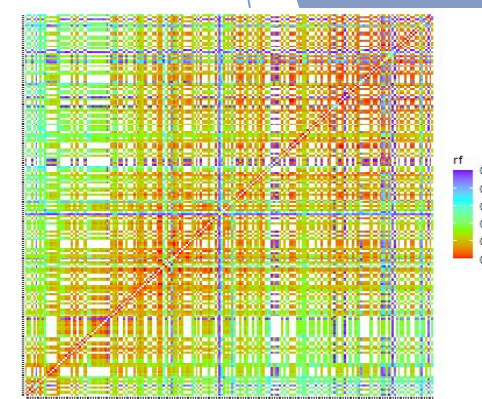
- 37% of chromosome 1
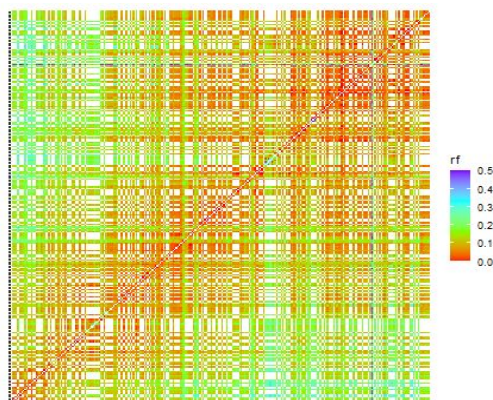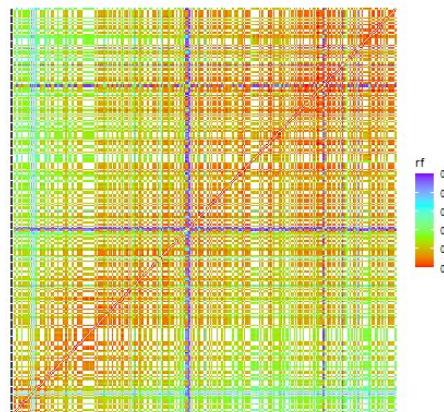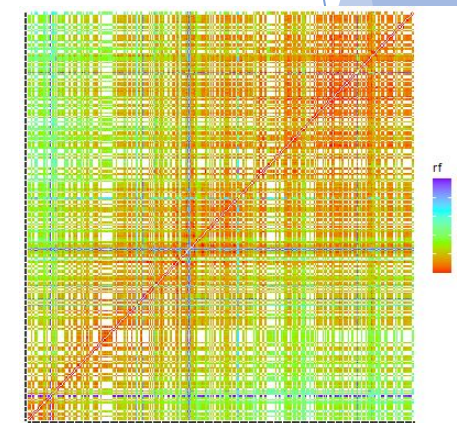
- ~ 38 cM

- Sequencing depth ~ 94X

OneMap

GATK

GATK + polyRAD

GATK + updog

freebayes

freebayes + polyRAD

freebayes + updog

TOOLS FOR POLYPLOIDS

GATK + updog | GATK + polyRAD | freebayes + polyRAD | freebayes + updog

Simulation studies: 10,880 maps
Empirical studies: 816 maps

Filters:
- genotype probabilities
- non-informative markers
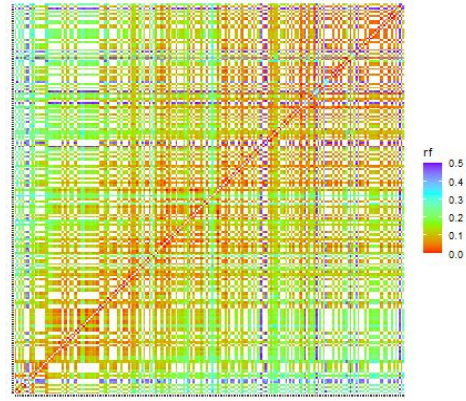- replace AD by missing when GT is missing

GATK + updog | GATK + polyRAD | freebayes + polyRAD | freebayes + updog

TOOLS FOR
POLYPLOIDS

This project is funded by USDA NIFA Specialty Crop Research Initiative
Award # 2020-51181-32156 (09/01/20 - 08/31/24)

# Map size
## 37% of chromosome 1 ~38cM

# Map size

- Hidden Markov Model Emission Function
  - global error rate
  - genotype probabilities (PL)



changed = 0.25%

imputed = 1.1%

unchanged = 98.65%

Figure by Jeekin Lau

TOOLS FOR
POLYPLOIDS

# Map size
## 37% of chromosome 1 ~38cM

- Hidden Markov Model Emission Function
  - global error rate
  - genotype probabilities (PL)

# Map size
## 37% of chromosome 1 ~38cM

TOOLS FOR POLYPLOIDS

# Simulation study

- ## Match recombinations breakpoints
  - Large maps - always bad
  - Small maps - not always good

- ## Other tested scenarios with:
  - Segregation distortion
  - Contaminants samples
  - Multiallelic markers



Red square: no inflated size (1 or less Euclidean distance) but have from 10 to 100 wrong recombination breakpoints

# Preprint

# Diploid Aspen

- 37% of chromosome 10

- Sequencing depth ~ 6X



OneMap

GATK

GATK + polyRAD

GATK + updog

freebayes

freebayes + polyRAD

freebayes + updog

TOOLS FOR POLYPLOIDS

TEXAS A&M UNIVERSITY

# Tetraploid rose

- Chromosome 2

- Sequencing depth ~ 50X

MAPpoly

GATK

GATK + polyRAD

GATK + updog

freebayes

freebayes + polyRAD

freebayes + updog

# Hands-on Workshop!!

https://github.com/Cristianetaniguti/Reads2Map

Tutorial:

# bit.ly/Reads2Map2023

TOOLS FOR
POLYPLOIDS

# Reads2Map

# Reads2Map

- Cloud environments
  - terra.bio
- HPC
  - Cromwell
  - MiniWDL
  - dxWDL

inputs.json

```json
{
"SNPCalling.max_cores": 2,
"SNPCalling.ploidy": 4,
"SNPCalling.rm_dupli": false,
"SNPCalling.replaceAD": false,
"SNPCalling.run_gatk": true,
"SNPCalling.run_freebayes": true,
"SNPCalling.hardfilters": true,
"SNPCalling.n_chrom": 1,
"SNPCalling.chunk_size": 2,
"SNPCalling.samples_info": "tests/data/polyploid/fastq/samples_info.txt",
"SNPCalling.gatk_mchap": false,
"SNPCalling.references": {
 "ref_fasta": "tests/data/polyploid/RchinensisV1.0/Chr04_sub.fasta",
 "ref_dict": "tests/data/polyploid/RchinensisV1.0/Chr04_sub.dict",
 "ref_ann": "tests/data/polyploid/RchinensisV1.0/Chr04_sub.fasta.ann",
 "ref_sa": "tests/data/polyploid/RchinensisV1.0/Chr04_sub.fasta.sa",
 "ref_amb": "tests/data/polyploid/RchinensisV1.0/Chr04_sub.fasta.amb",
 "ref_pac": "tests/data/polyploid/RchinensisV1.0/Chr04_sub.fasta.pac",
 "ref_bwt": "tests/data/polyploid/RchinensisV1.0/Chr04_sub.fasta.bwt",
 "ref_fasta_index": "tests/data/polyploid/RchinensisV1.0/Chr04_sub.fasta.fai"
}
}
}
```
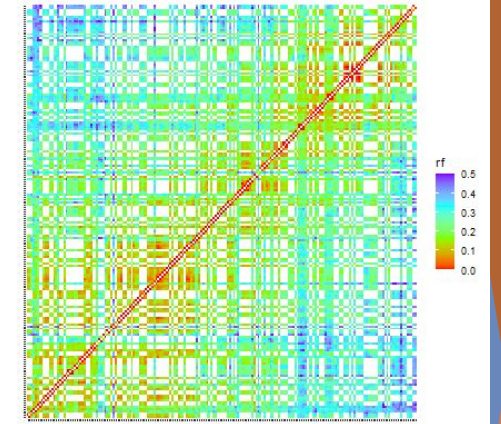
```
$ java -jar /path/to/cromwell.jar run -i EmpiricalSNPCalling/inputs.json EmpiricalSNPCalling .wdl
```

TOOLS FOR
POLYPLOIDS

# task

Cristianetaniguti joint wf

..

BWA.wdl

JointReports.wdl

bcftools.wdl

chunk_lists.wdl

cutadapt.wdl

freebayes.wdl

gatk.wdl

gusmap.wdl

mchap.wdl

pedigree_simulator.wdl

gusmap.wdl

mchap.wdl

pedigree_simulator.wdl

pedigree_simulator_utils.wdl

pirs.wdl

radinitio.wdl

simuscop.wdl

stacks.wdl

utils.wdl

utilsR.wdl

vcf2diploid.wdl

TEXAS A&M UNIVERSITY

# task
## example: freebayes.wdl

```
task RunFreebayes {

    input { ... }

    Int disk_size = ceil(size(reference, "GiB") + size(bam, "GiB") +  50)
    Int memory_size = ceil(size(bam, "MiB") * 25 * max_cores + 10000)

    command <<< ...
    >>>

    runtime { ... }

    meta { ... }

    output { ... }

}
```

```
input {
    File reference
    File reference_idx
    File bam
    File bai
    Int max_cores
    Int ploidy
}

Int disk_size = ceil(size(reference, "GiB") + size(bam, "GiB") +  50)
Int memory_size = ceil(size(bam, "MiB") * 25 * max_cores + 10000)
```

inputs.json

```
{
    "reference": "tests/data/PtrichocarpaV3.0/Chr10.2M.fa",
    "reference_idx": "tests/data/PtrichocarpaV3.0/Chr10.2M.fa.fai",
    "bam": "tests/data/Ptremula_PRJNA395596_subset/merged.bam",
    "bai": "tests/data/Ptremula_PRJNA395596_subset/merged.bam.bai",
    "max_cores": 2,
    "ploidy": 4
}
```

TOOLS FOR POLYPLOIDS

TEXAS A&M UNIVERSITY

# task
## example: freebayes.wdl



```
task RunFreebayes {

>   input { ⋯
    }

    Int disk_size = ceil(size(reference, "GiB") + size(bam, "Gi
    Int memory_size = ceil(size(bam, "MiB") * 25 * max_cores +

>   command <<< ⋯
    >>>

>   runtime { ⋯
    }

>   meta { ⋯
    }

>   output { ⋯
    }
}
```

```
command <<<

    ln -s ~{bam} .
    ln -s ~{bai} .

    freebayes-parallel <(fasta_generate_regions.py ~{reference_idx} 100000) ~{max_cores} \
    --genotype-qualities --ploidy ~{ploidy} -f ~{reference} *bam > "freebayes.vcf"

>>>
```

# task
## example: freebayes.wdl

```
task RunFreebayes {

>   input { ...
    }

    Int disk_size = ceil(size(reference, "GiB") + size(bam, "GiB") +  50)
    Int memory_size = ceil(size(bam, "MiB") * 25 * max_cores + 10000)

>   command <<< ...
    >>>

>   runtime { ...
    }

>   meta { ...
    }

>   output { ...
    }
}
```

```
runtime {
    docker: "cristaniguti/freebayes:0.0.1"
    cpu: max_cores
    # Cloud
    memory:"~{memory_size} MiB"
    disks:"local-disk " + disk_size + " HDD"
    # Slurm
    job_name: "RunFreebayes"
    mem:"~{memory_size}M"
    time:"48:00:00"
}
```

# task
## example: freebayes.wdl



```
task RunFreebayes {

    input { ⋯
    }

    Int disk_size = ceil(size(reference, "GiB") + size(bam, "GiB") +  50)
    Int memory_size = ceil(size(bam, "MiB") * 25 * max_cores + 10000)

    command <<< ⋯
    >>>

    runtime { ⋯
    }

    meta { ⋯
    }

    output { ⋯
    }
}
```

```
meta {
    author: "Cristiane Taniguti"
    email: "chtaniguti@tamu.edu"
    description: "Split genomic regions and runs [freebayes](https://github.com/freebayes/freebayes) parallelized."
}
```

# task
## example: freebayes.wdl

```
task RunFreebayes {

>   input {···
    }

    Int disk_size = ceil(size(reference, "GiB") + size(bam, "GiB") + 50)
    Int memory_size = ceil(size(bam, "MiB") * 25 * max_cores + 10000)

>   command <<< ···
    >>>

>   runtime {···
    }

>   meta {···
    }

>   output {···
    }
}
```

**inputs.json**

```
{
    "reference": "tests/data/PtrichocarpaV3.0/Chr10.2M.fa",
    "reference_idx": "tests/data/PtrichocarpaV3.0/Chr10.2M.fa.fai",
    "bam": "tests/data/Ptremula_PRJNA395596_subset/merged.bam",
    "bai": "tests/data/Ptremula_PRJNA395596_subset/merged.bam.bai",
    "max_cores": 2,
    "ploidy": 4
}
```

```
output {
    File vcf = "freebayes.vcf"
}
```

```
miniwdl run --task RunFreebayes -i tests/tasks/freebayes/inputs.json tasks/freebayes.wdl
```

# subworkflow
## example

```
workflow FreebayesGenotyping {
>    input {···
    }

>    call chunk_lists.CreateChunksBamByChr {···
    }

    scatter (chunk in zip(CreateChunksBamByChr.bams_chunks, CreateChunksBamByChr.bais_chunks)) {

>       call freebayes.RunFreebayes {···
       }
    }

>    call utils.mergeVCFs {···
    }

>    call norm_filt.Normalization {···
    }

    Map[String, Array[File]] map_bams = {"bam": CreateChunksBamByChr.bams_chunks, "bai": CreateChunk

>    if(replaceAD){···
    }

    Array[File] freebayes_vcfs = select_all([Normalization.vcf_norm, ReplaceAD.bam_vcf])
    Array[String] freebayes_software = select_all([Normalization.software, ReplaceAD.software])
    Array[String] freebayes_counts_source = select_all([Normalization.source, ReplaceAD.source])

>    output {···
    }
}
```

```
$ java -jar /path/to/cromwell.jar  run -i freebayes_genotyping/inputs.json
freebayes_genotyping.wdl
```

TOOLS FOR
POLYPLOIDS

This project is funded by USDA NIFA Specialty Crop Research Initiative
Award # 2020-51181-32156 (09/01/20 - 08/31/24)

ĀTM̄ | TEXAS A&M
UNIVERSITY.

# pipeline
## example

```
workflow SNPCalling {

    input { ...
    }

    call fam.CreateAlignmentFromFamilies { ...
    }

    if(run_gatk){
        call gatk.GatkGenotyping { ...
        }
    }

    if(run_freebayes){
        call freebayes.FreebayesGenotyping { ...
        }
    }

    Array[Array[File]] vcfs_sele = select_all([GatkGenotyping.vcfs, FreebayesGenotyping.vcfs])
    Array[Array[String]] software_sele = select_all([GatkGenotyping.vcfs_software, FreebayesGenotyping.vcfs_software
    Array[Array[String]] source_sele = select_all([GatkGenotyping.vcfs_counts_source, FreebayesGenotyping.vcfs_count

    output { ...
    }
}
```

```
$ java -jar /path/to/cromwell.jar  run -i EmpiricalSNPCalling/inputs.json  EmpiricalSNPCalling.wdl
```

# Thank you



Oscar Riera-Lizarazu
David Byrne
Jeekin Lau
Tessa Hochhaus

Augusto Garcia
Rodrigo Amadeu
Getulio Caixeta

Marcelo Mollinari
Gabriel Gesteira

Lucas Taniguti

Guilherme Pereira

Thiago Oliveira

# Thank you!!

- ► Susan Thomson
- ► Cecilia Deng
- ► Ben Warren

Plant & Food™
Research
Rangahau Ahumāra Kai

TOOLS FOR
POLYPLOIDS

# Project Members

# Other Collaborators

# Links for Tutorials

- [polyRAD tutorial](#)

- [updog tutorial](#)

- [fitPoly tutorial](#)

- [(TASSEL) Variant and Genotype Calling in Highly Duplicated Genomes (Lindsay Clark)](#)

- [Using Reads2Map workflows for SNP and dosage calling in polyploid sequencing data](#)

More cool tools related to Workflows systems:

- [cromwell-cli](#) (by my brother)

- [cromwell server](#)

- [womtools](#)

TOOLS FOR
POLYPL⬡IDS

# References

- Catchen, J., Hohenlohe, P. A., Bassham, S., Amores, A.; Cresko, W. A. (2013). Stacks: an analysis tool set for population genomics. Molecular Ecology, 22(11), 3124–3140. https://doi.org/10.1111/mec.12354

- Glaubitz, J. C., Casstevens, T. M., Lu, F., Harriman, J., Elshire, R. J., Sun, Q.; Buckler, E. S. (2014). TASSEL-GBS: a high capacity genotyping by sequencing analysis pipeline. PLoS ONE, 9(2), 1–11. https://doi.org/10.1371/journal.pone.0090346

- Garrison, E.; Marth, G. (2012). Haplotype-based variant detection from short-read sequencing. ArXiv E-Prints, 9. https://doi.org/1207.3907

- McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M.; DePristo, M. A. (2010). The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. Genome Research, 20(9), 1297–1303. https://doi.org/10.1101/gr.107524.110

TOOLS FOR
POLYPLOIDS

# References

► Gerard, D., Ferrão, L. F. V., Garcia, A. A. F.,& Stephens, M. (2018). Genotyping Polyploids from Messy Sequencing Data. Genetics, 210(3), 789-807. doi: 10.1534/genetics.118.301468.

► Wadl, P. A., Olukolu, B. A., Branham, S. E., Jarret, R. L., Yencho, G. C.; Jackson, D. M. (2018). Genetic Diversity and Population Structure of the USDA Sweetpotato (Ipomoea batatas) Germplasm Collections Using GBSpoly. Frontiers in Plant Science, 9, 1166. https://doi.org/10.3389/fpls.2018.01166

► Serang, O., Mollinari, M.; Garcia, A. A. F. (2012). Efficient exact maximum a posteriori computation for bayesian SNP genotyping in polyploids. PLoS ONE, 7(2), 1–13. https://doi.org/10.1371/journal.pone.0030906

► Clark, L. v., Lipka, A. E.; Sacks, E. J. (2019). polyRAD: Genotype Calling with Uncertainty from Sequencing Data in Polyploids and Diploids. G3: Genes|Genomes|Genetics, 9(March), g3.200913.2018. https://doi.org/10.1534/g3.118.200913