



# **“Genomic resource development and analysis in Apples & Pears”**

Dr. Alan E. Yocca for the position of Genomics Data Scientist

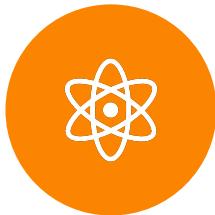
# Career Background



The Pennsylvania State University:  
2014 – 2016  
Undergraduate Research  
Assistant - dePamphilis lab



Michigan State University:  
2017 – 2022  
Graduate Assistant - Edger lab  
Dissertation title: “Leveraging  
Angiosperm Pangenomics to  
Understand Genome Evolution”



HudsonAlpha Institute for  
Biotechnology: 2022 – present  
Postdoctoral Research Associate -  
Harkess lab



USDA:  
2023 – present  
Postdoctoral Research Associate -  
Honaas lab

# Evolution of Conserved *Nothogardia thaliana* Origin and evolution genome

# Uncovering a miltiradiene biosynthetic gene cluster in the Lamiaceae reveals a dynamic evolutionary trajectory

Patrick P. Edger<sup>1,2\*</sup>, Thomas J. Poorten<sup>3</sup>, R Michael R. McKain<sup>4</sup>, Ronald D. Smith<sup>6</sup>, Elizabeth I. Alger<sup>1</sup>, Kevin A. Bird<sup>1,2</sup>, Avital Brodt<sup>8</sup>, Kobi Baruch<sup>8</sup>, Thomas Jeffrey P. Mower<sup>10</sup>, Kevin L. Childs<sup>11,12</sup> and Steven J. Knapp<sup>3\*</sup>

Received: 8 April 2022

Accepted: 4 January 2023

Abigail E. Bryson<sup>1,7</sup>, Emily R. Lanier<sup>1,7</sup>, Kin H. Lau<sup>2,3</sup>, John P. Hamilton<sup>2,4</sup>, Brieanne Vaillancourt<sup>1,2,4</sup>, Davis Mathieu<sup>1</sup>, Alan E. Yocca<sup>1,2,5</sup>, Garret P. Miller<sup>1,6</sup>

## Genome of the North American wild apple species *Malus angustifolia*

Ben N. Mansfeld, Shujun Ou, Erik Burchard, Alan Yocca, Alex Harkess, Ben Gutierrez, Steve van Nocker, Lisa Tang, Christopher Gottschalk

## Machine learning approach to genes in pangenomes

## Disease Resistance Genetics and Genomics in Octoploid Strawberry

Christopher R. Barbey,<sup>\*,†,1</sup> Seonghee Lee,<sup>‡</sup> Sujeet Verma,<sup>‡</sup> Kevin A. Bird,<sup>§,\*\*</sup> Alan E. Yocca,<sup>††</sup> Daniel C. Folta,<sup>‡,††</sup> Michael J. Feltner,<sup>‡,††</sup> Michael J. Folta,<sup>‡,††</sup> and Kevin M. Folta<sup>\*,†</sup>

## Exceptional subgenome stability and functional divergence in the allotetraploid Ethiopian cereal teff

Robert VanBuren<sup>1,2,9</sup>, Ching Man Wai<sup>1,2,9</sup>, Xuewen Wang<sup>3,9</sup>, Jeremy Pardo<sup>1,2,4</sup>, Alan E. Yocca<sup>1,4</sup>, Hao Wang<sup>3</sup>, Srinivasa R. Chaluvadi<sup>3</sup>, Guomin Han<sup>1,3</sup>, Douglas Bryant<sup>5</sup>, Patrick P. Edger<sup>1,1</sup>, Joachim Messing<sup>1,6</sup>, Mark E. Sorrells<sup>7</sup>, Todd C. Mockler<sup>1,5</sup>, Jeffrey L. Bennetzen<sup>1,3</sup> & Todd P. Michael<sup>1,8</sup>

## In-genome dynamics of

Alan E. Yocca<sup>1,2,5,\*,†,§,||</sup>, Adrian Platts<sup>1,2,5,†,§,||</sup>, Elizabeth Alger<sup>1,2,5,†,§,||</sup>, Qin Qiao<sup>1,2,5,†,§,||</sup>, Patrick P. Edger<sup>1,2,5,†,§,||</sup>, Li Xue<sup>1,2,5,†,§,||</sup>, Li Qiong<sup>1,2,5,†,§,||</sup>, Jie Lu<sup>3,8</sup>, Yichen Zhang<sup>3</sup>, Qiang Cao<sup>3</sup>, Alan E. Yocca<sup>1,2,5,†,§,||</sup>, Michal Babinski<sup>1,1</sup>, Maria Magallanes-Lundback<sup>1</sup>, Philip Adrian E. Platts<sup>1,2,5,†,§,||</sup>, Steven J. Knapp<sup>1,2,5,†,§,||</sup>, Marc Van Montagu<sup>9,10,11,12</sup>, Yves Van de Peer<sup>9,10,11,12</sup>, Jiajun Lei<sup>10,12</sup>, and Ticao Zhang<sup>1,2,5,†,§,||</sup>, Nick W. Albert<sup>11,12</sup>, Sara Montanari<sup>12</sup>, Nicholi Vorsa<sup>13</sup>, Massimo Iorizzo<sup>4,14</sup> and Patrick P. Edger<sup>1,2,5,†,§,||</sup>

# Postdoc Projects:

- Chromosome scale haplotype phased genome assemblies with PacBio HiFi
- *Malus* crop wild relatives assembly and annotation
- Detecting recombination events using  $k$ -mers
- When is the perfect time to pick an apple or pear?
- Current progress in graph-based pangenomes



# Apples & Pears

- *Malus x domestica* (Apple) and *Pyrus communis* (European Pear)
- Pome fruits in Rosaceae
- Diverged ~8 MYA
- Share ancient WGD ~40 MYA
- Genome size ~600-750Mb
- 17 chromosomes
- >50% of genes found in duplicate from last WGD



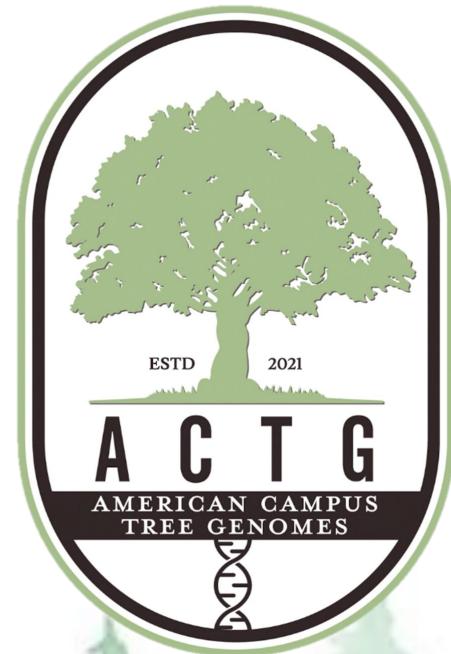
# Postdoc Projects:

- **Chromosome scale haplotype phased genome assemblies with PacBio HiFi**
- *Malus* crop wild relatives assembly and annotation
- Detecting recombination events using  $k$ -mers
- When is the perfect time to pick an apple or pear?
- Current progress in graph-based pangenomes



# American Campus Tree Genomes

- “American Campus Tree Genomes Project pushes equality in genomic science”
- Students sequence, assembly, annotate, and analyze a chromosome-scale genome
- Final project is a figure on a manuscript





Genes | Genomes | Genetics



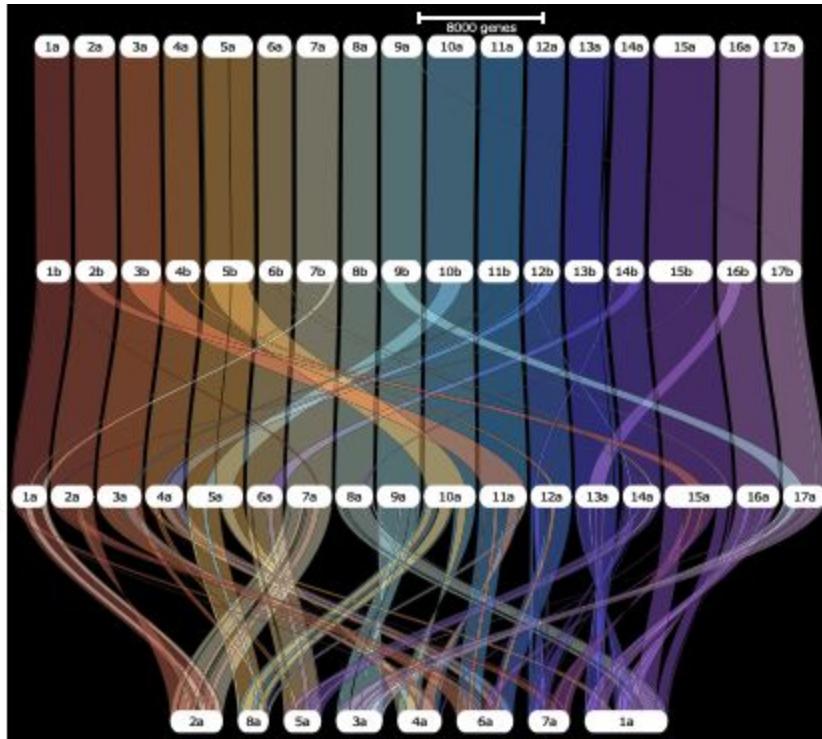
MARCH 2024  
VOLUME 14 • ISSUE 3  
[academic.oup.com/g3journal](http://academic.oup.com/g3journal)

Hap 1

Hap 2

*Pirus*  
*vesticata*

*P. pyrifera*  
*Yasudae*



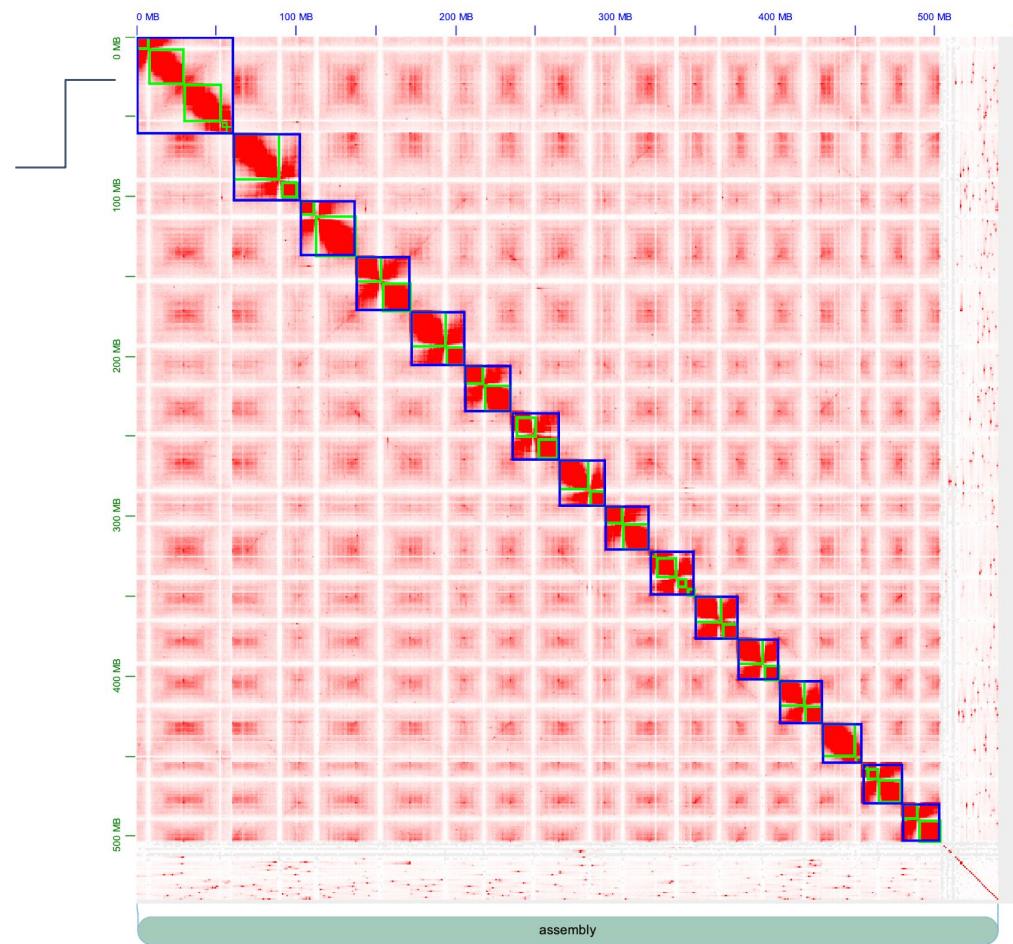
# HiFi genome assembly

- ~40x HiFi coverage per haplotype
- ~80x Omni-C coverage per haplotype
- Assembled with hifiasm
- Quality Control:
  - Omni-C contact map
  - Dot plots to existing assemblies and to itself
  - Remap HiFi data to the assembly



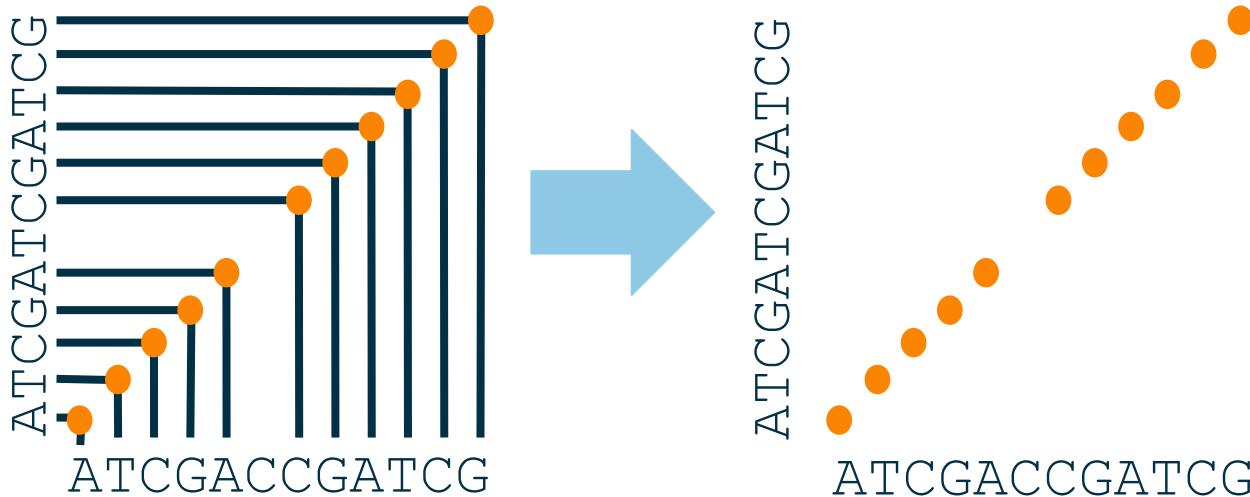
# Omni-C contact map

Chromosome  
fusion, not real!



# Dot plots

- Sequence alignment in 2-dimensions



chr3A vs. chr3B

Zoom: 59448 : 1

Word length: 20

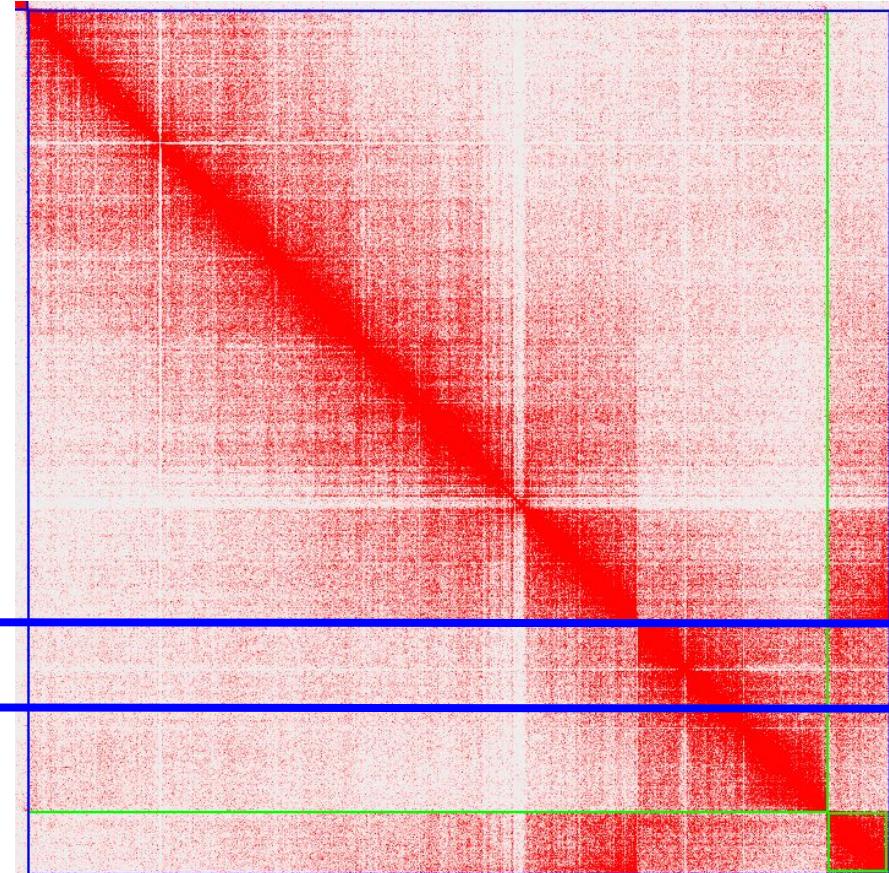
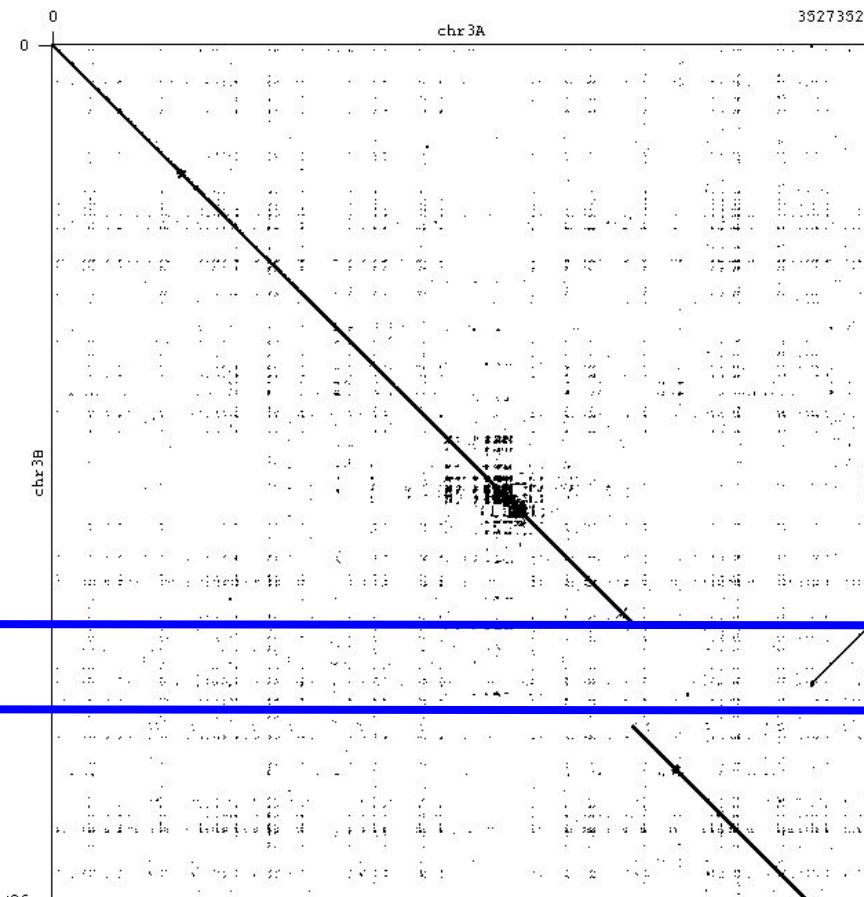
GC ratio seq1: 0.3789

Window size: 0

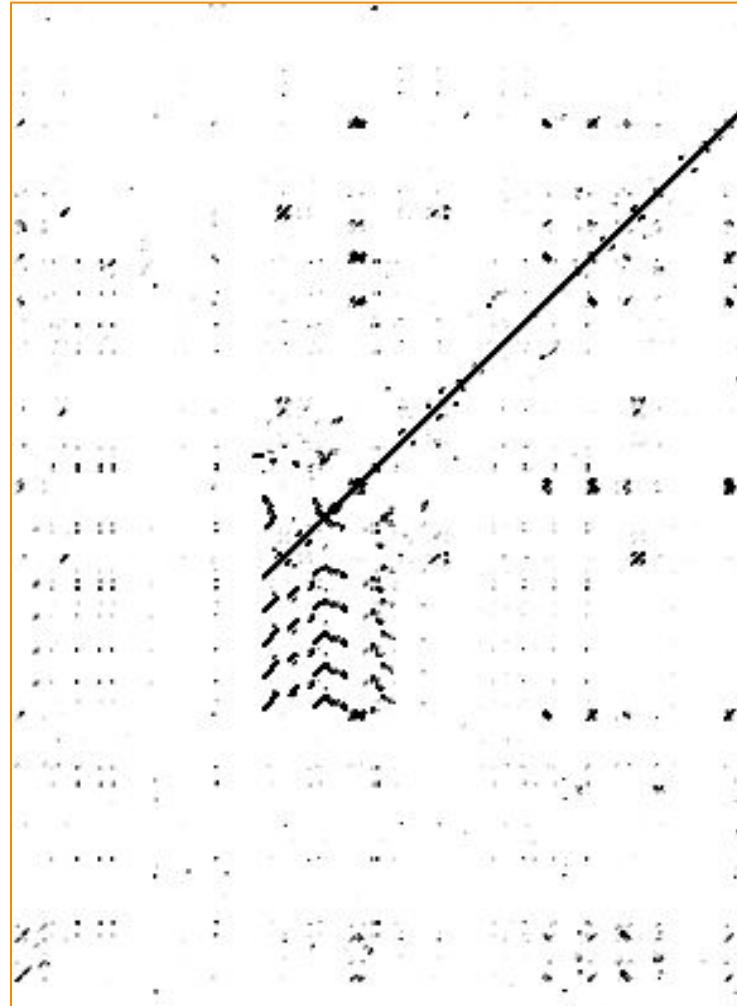
GC ratio seq2: 0.3792

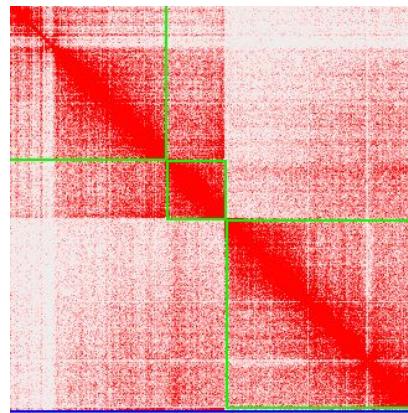
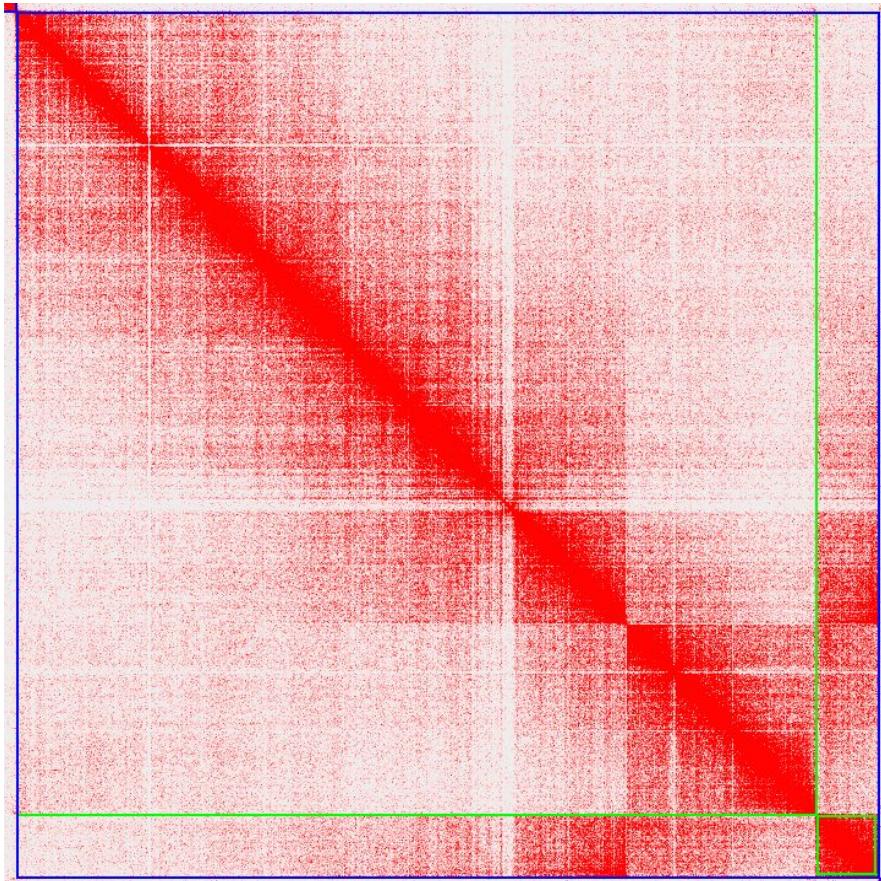
Matrix: edna.mat

Program: Gepard (2.0)



Zoomed in on  
breakpoints of  
the misplaced  
contig





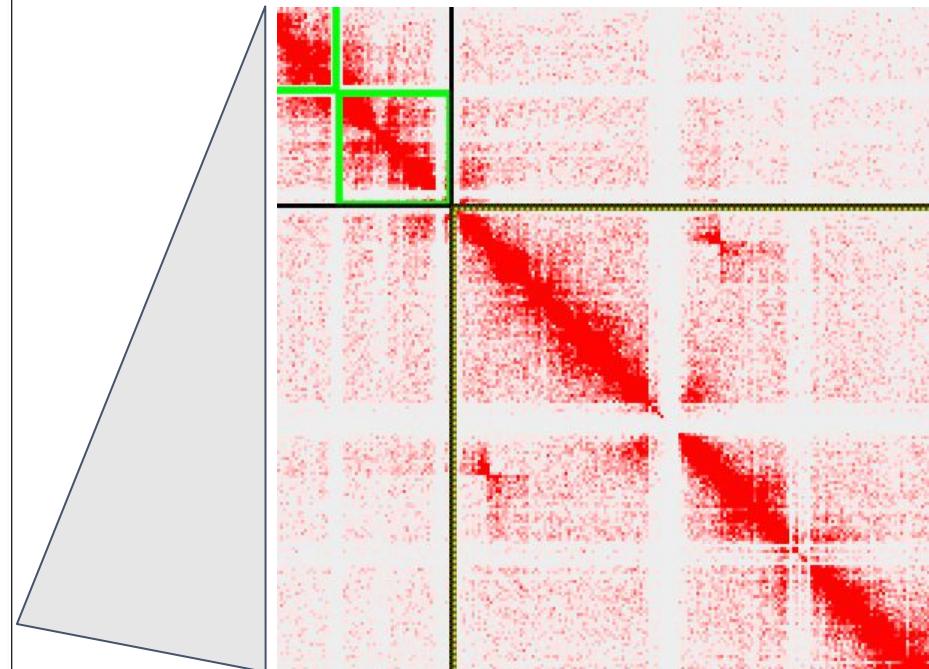
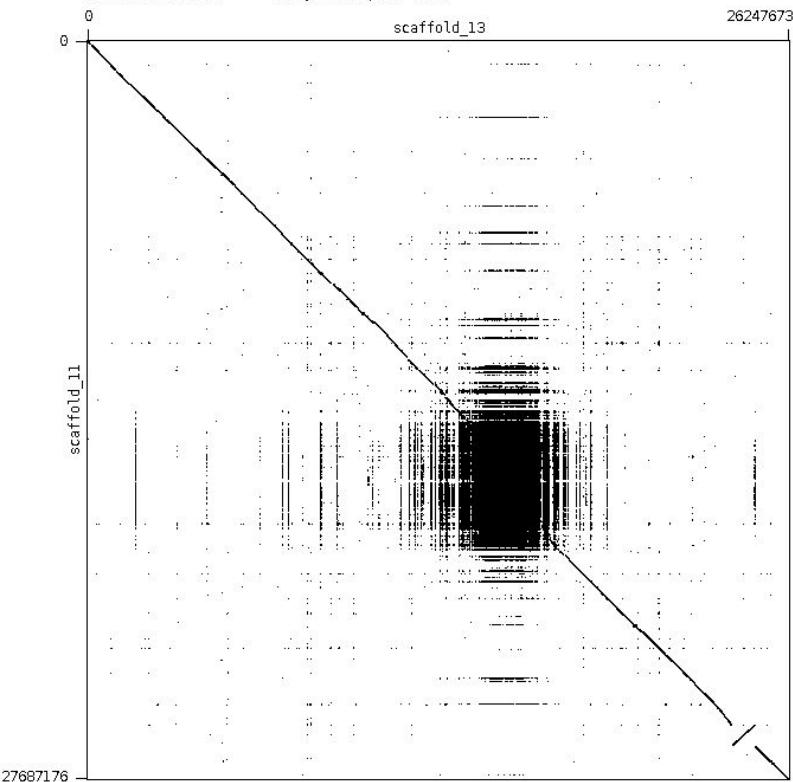
scaffold\_13 vs. scaffold\_11

Zoom: 44300 : 1

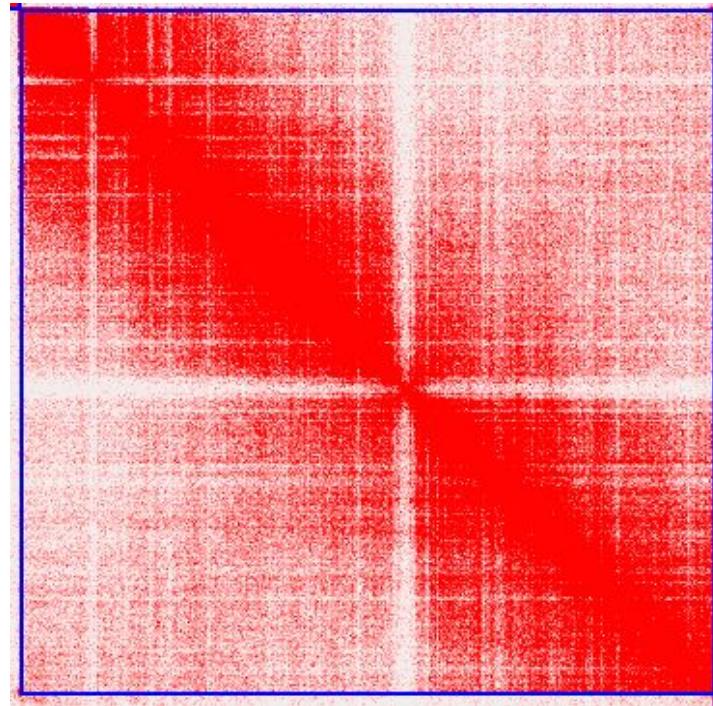
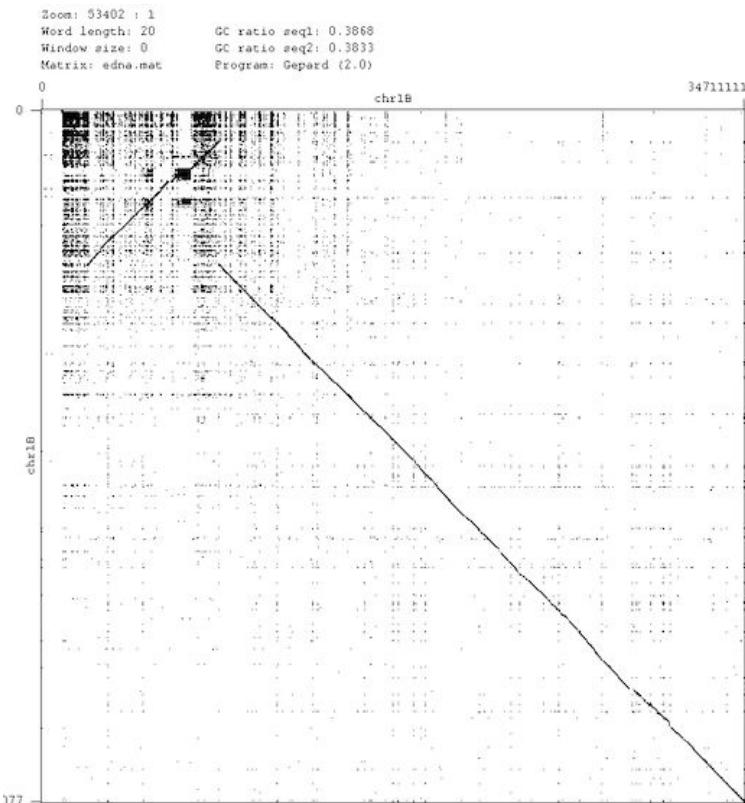
Word length: 20 GC ratio seq1: 0.3734

Window size: 0 GC ratio seq2: 0.3738

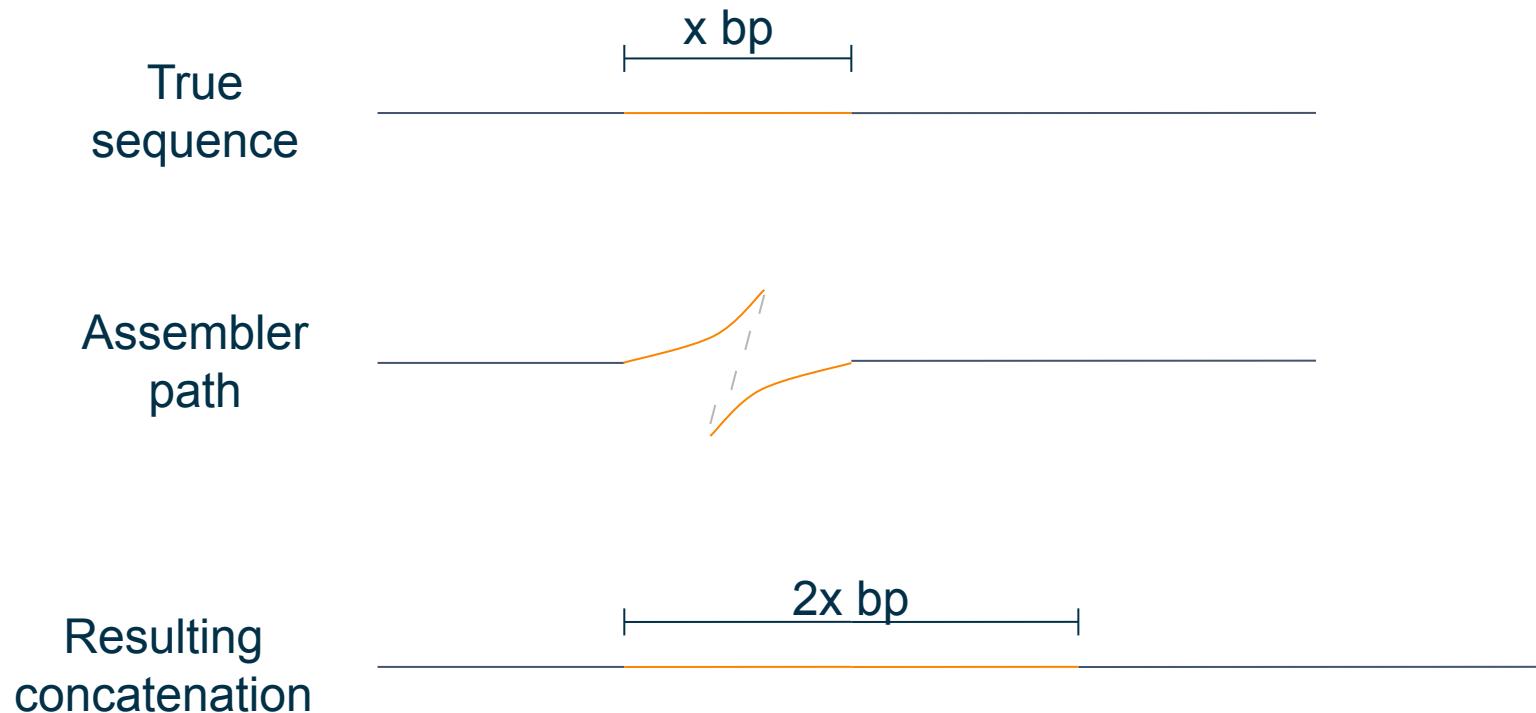
Matrix: edna.mat Program: Gepard (2.0)



# Some assembly issues are not as straightforward

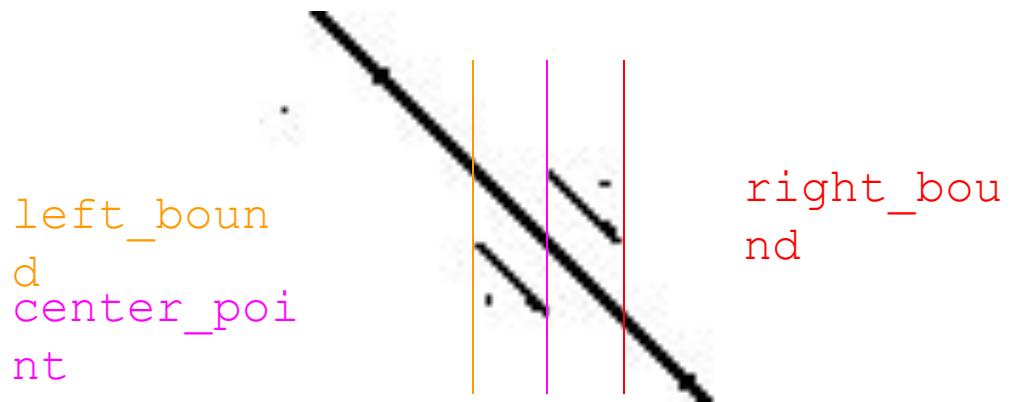


# Alternate assembled haplotype

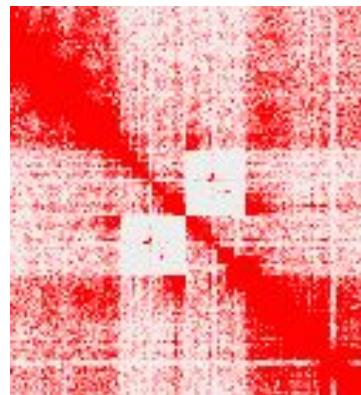


# Alternate assembled haplotype

Dot plot “self vs. self”



Omni-C map



# Postdoc Projects:

- Chromosome scale haplotype phased genome assemblies with PacBio HiFi
- ***Malus* crop wild relatives assembly and annotation**
- Detecting recombination events using  $k$ -mers
- When is the perfect time to pick an apple or pear?
- Current progress in graph-based pangenomes





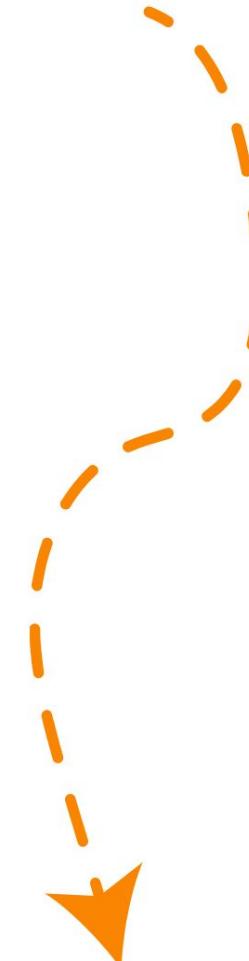
# *Malus* crop wild relatives

- 4 native North American *Malus* species
- Readily hybridize with domesticated Apple
- Native species exhibit abiotic and biotic resistance
  - Fire blight resistance in *M. fusca*



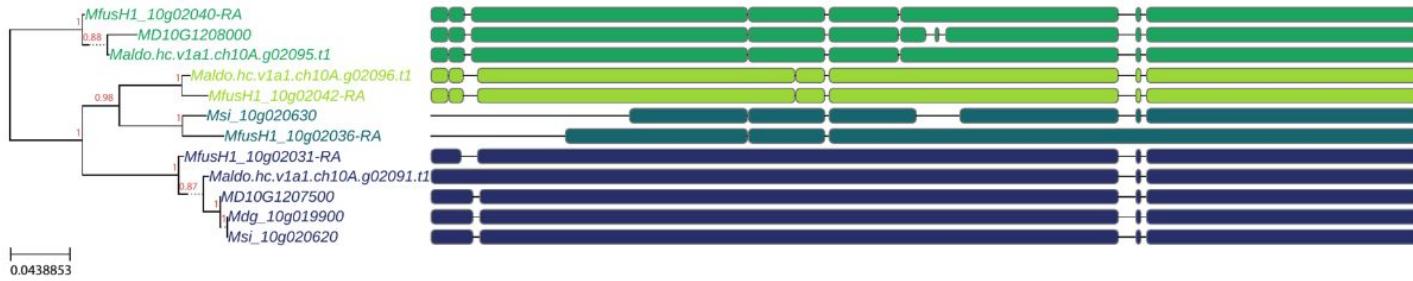
# Annotation workflow

- MAKER pipeline
- Evidence-based annotation
  - Arabidopsis proteins
  - “Honeycrisp” apple proteins
  - Available genotype-specific RNA-seq of diverse tissue types
- Two-rounds of *ab-initio* gene prediction with AUGUSTUS and SNAP
- Quality control
  - Gene length distributions
  - BUSCO
  - Exon count distributions
  - Genome-browser

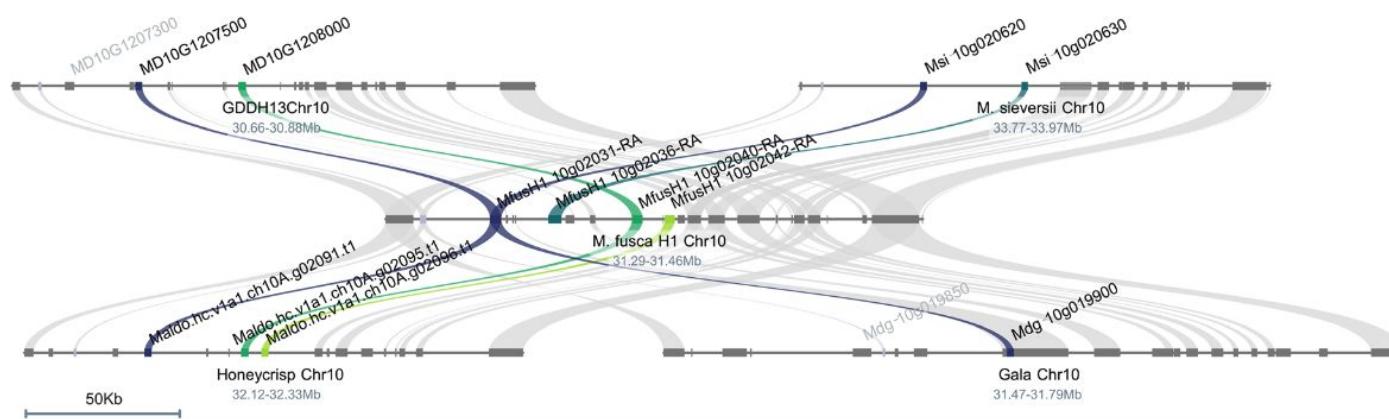


# *Malus fusca* Fire Blight resistance locus

(c)

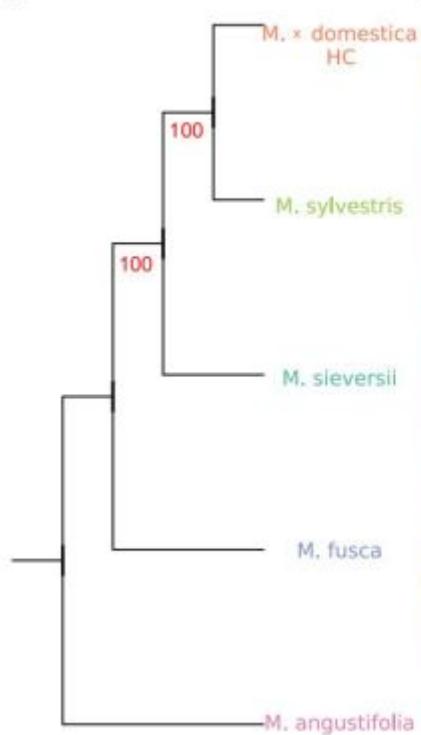


(d)

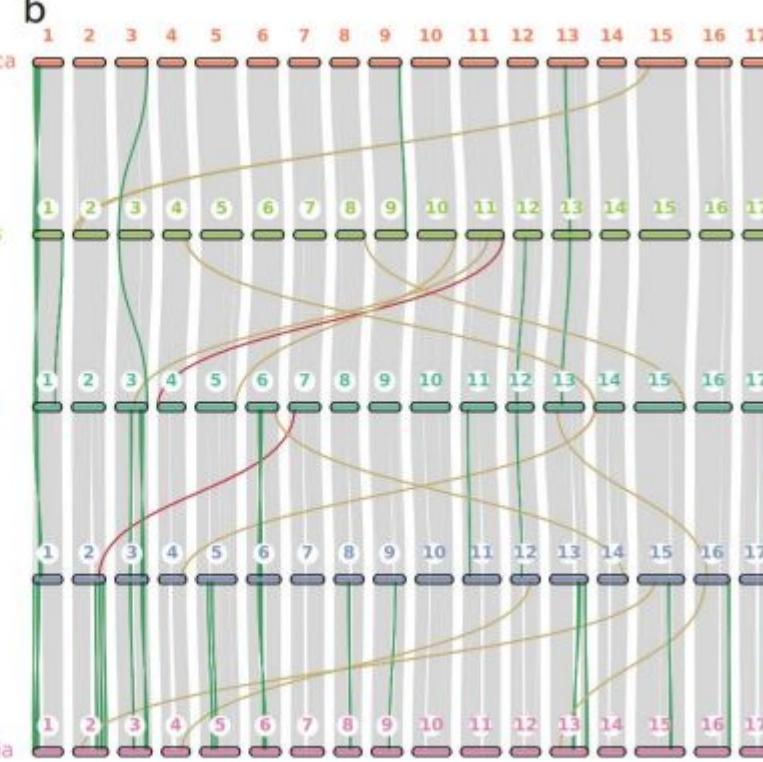


# *Malus angustifolia*

a



b



# Postdoc Projects:

- Chromosome scale haplotype phased genome assemblies with PacBio HiFi
- *Malus* crop wild relatives assembly and annotation
- **Detecting recombination events using  $k$ -mers**
- When is the perfect time to pick an apple or pear?
- Current progress in graph-based pangenomes

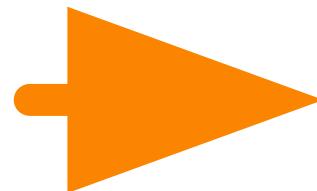


# Experimental System

“Honeycrisp”



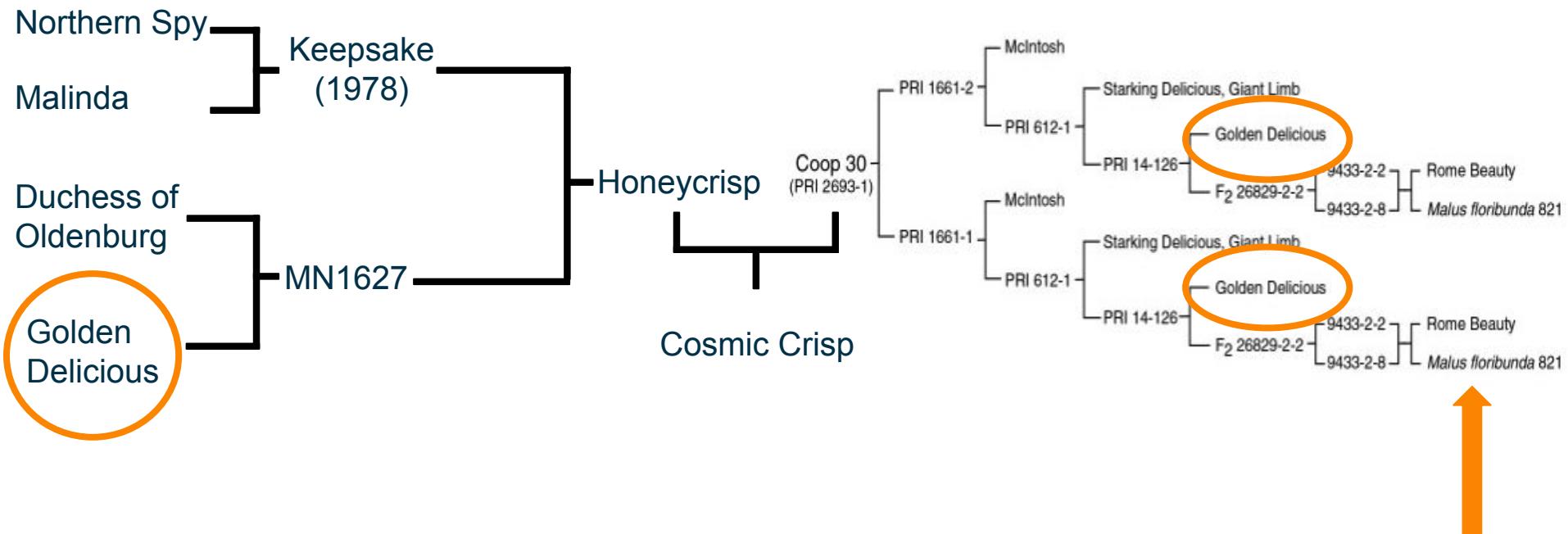
“Enterprise”



“Cosmic Crisp”

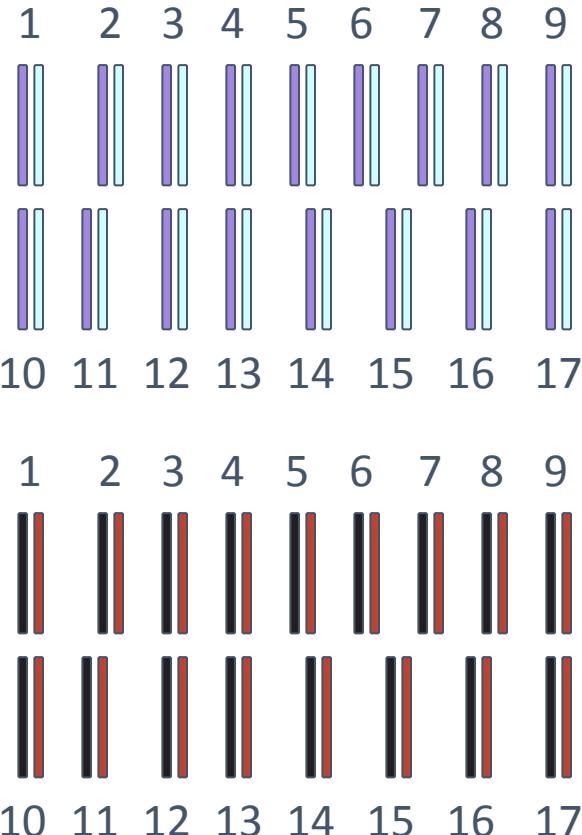


# Experimental System



# Experimental System

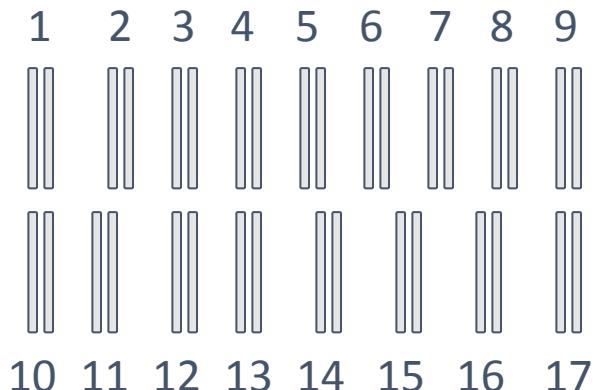
“Honeycrisp”



“Enterprise”



“Cosmic Crisp”



??

# Which “Cosmic Crisp” chromosome was inherited from which parent?

- We have a phased “Cosmic crisp” and a phased “Honeycrisp” (HC from here on) which one parent of “Cosmic crisp” (CC from here on)
- Lets use kmers to find which copy in CC came from HC and which from the other parent “Enterprise”



# Workflow

- Use meryl to find kmers specific to one CC haplotype not found in the other CC haplotype
- Intersect these with kmers found in both HC haplotypes
- The chromosome that came from HC should have a higher abundance of kmers shared than the other haplotype



# Results based on $k$ -mer abundance

Chromosome	HapA-specific shared $k$ -mers	HapB-specific shared $k$ -mers	Difference	Ratio HapA/HapB
1	2051610	5471251	-3419641	0.3749800548
2	3077972	9635778	-6557806	0.3194316017
3	7550411	3455526	4094885	2.185025087
4	6170561	1835473	4335088	3.361836976
5	2910988	6149547	-3238559	0.4733662496
6	6958204	4067368	2890836	1.710738738
7	1436727	8390383	-6953656	0.1712349722
8	7203946	3592405	3611541	2.005326794
9	1265578	3941492	-2675914	0.3210910995

Negative value reflect more CC hapB kmers matching HC than hapA kmers

# Results based on kmer abundance (cont)



Chromosome	HapA-specific shared kmers	HapB-specific shared kmers	Difference	Ratio HapA/HapB
10	7976761	2412989	5563772	3.305759371
11	2105156	6412403	-4307247	0.3282944007
12	5573204	1887142	3686062	2.953251001
13	2326268	9489397	-7163129	0.2451439222
14	5953204	1372498	4580706	4.337495574
15	9694659	4465069	5229590	2.171222662
16	10090592	3721436	6369156	2.711478042
17	5902270	2390097	3512173	2.469468812

7 chromosome where hapB matches better than hapA: 1,2,5,7,9,11,13

# k-mers to detect recombination events

- Now we know which “Cosmic Crisp” haplotype came from which parent
- Can we detect recombination events between parental haplotypes?



# Workflow

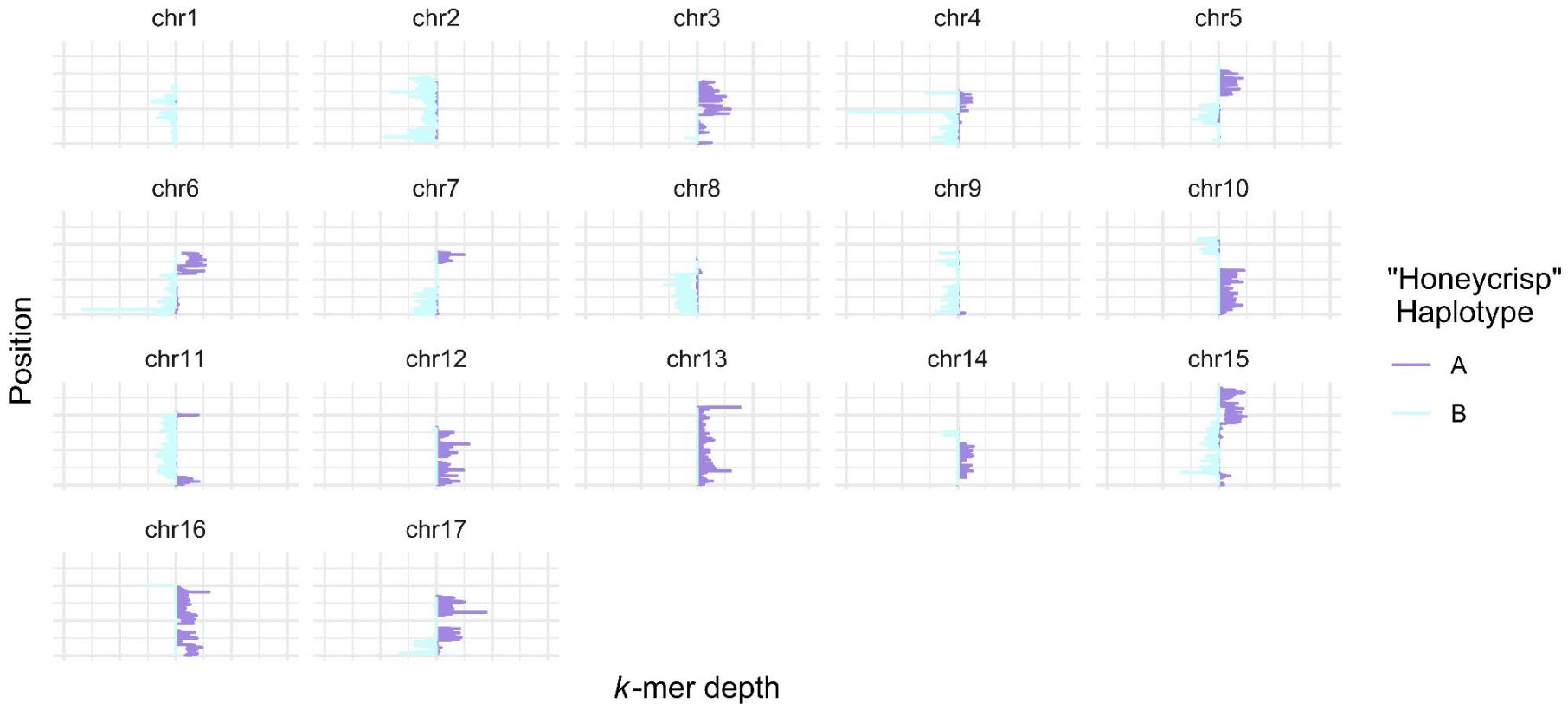
```
PROCESSING TREE #1 using 10 threads.
```

```
opDifference
  opDifference
    opDifference
      /cluster/home/ayocca/01_Rosaceae_Evo/Malus/01_cosmic_phasing//03_meryl/Mdom_CC_hapA_chr8A.meryl/
      /cluster/home/ayocca/01_Rosaceae_Evo/Malus/01_cosmic_phasing//03_meryl/Mdom_CC_hapB_chr8B.meryl/
    opIntersect
      /cluster/home/ayocca/01_Rosaceae_Evo/Malus/01_cosmic_phasing//03_meryl/Mdom_HC_hapA_chr8A.meryl/
      /cluster/home/ayocca/01_Rosaceae_Evo/Malus/01_cosmic_phasing//03_meryl/Mdom_HC_hapB_chr8B.meryl/
/cluster/home/ayocca/01_Rosaceae_Evo/Malus/01_cosmic_phasing//03_meryl/Mdom_EP_hapBoth_chr8Both.meryl/
output to /cluster/home/ayocca/01_Rosaceae_Evo/Malus/01_cosmic_phasing//03_meryl/Mdom_CC_hapA_chr8A.meryl.diff_hapB.diff_EP/
```

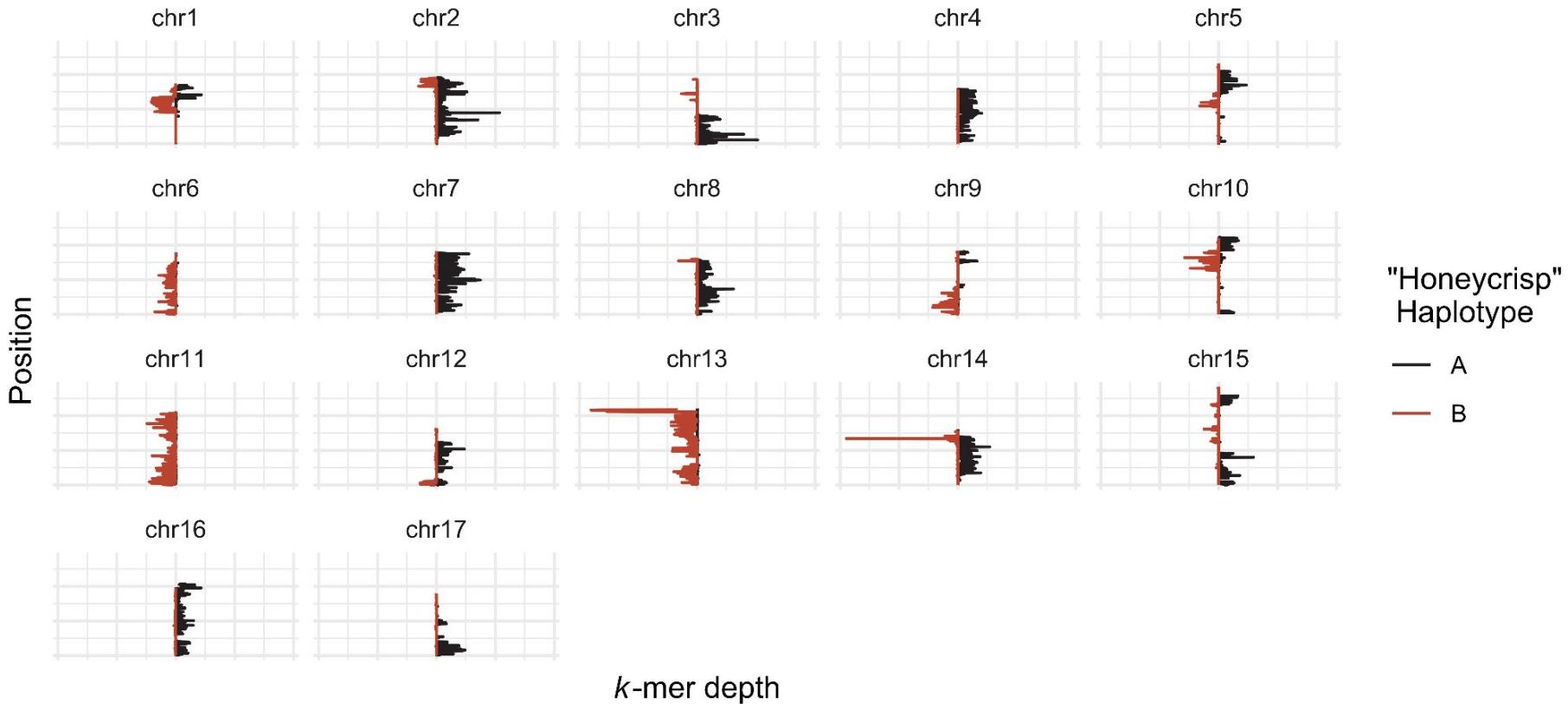
- k-mers specific to one haplotype
- Remove any k-mers present in the other parent
- Remove any k-mers not present in the parent of interest
- Map k-mers back to the parent of interest for only exact matches



## "Cosmic Crisp" Haplotype A $k$ -mer depth bias



## "Cosmic Crisp" Haplotype B $k$ -mer depth bias



# Motivating Questions:

- How good is PacBio HiFi?
- Can we generate chromosome scale haplotype phased genome assemblies?
- What is the genetic diversity of crop wild relatives?
- Can we detect recombination events using  $k$ -mers?
- **When is the perfect time to pick an apple or pear?**
- Current progress in graph-based pangenomes



# Apple Harvest

- Apples require cold storage after harvest to reach optimum maturity
- Response to cold storage and chilling hour requirements depend on maturity at harvest
- Current tests for maturity are unreliable especially across years and locations
- Is there a stage of physiological maturity measured by gene expression where post-harvest physiology and response to cold storage are predictable?



# Experimental Design

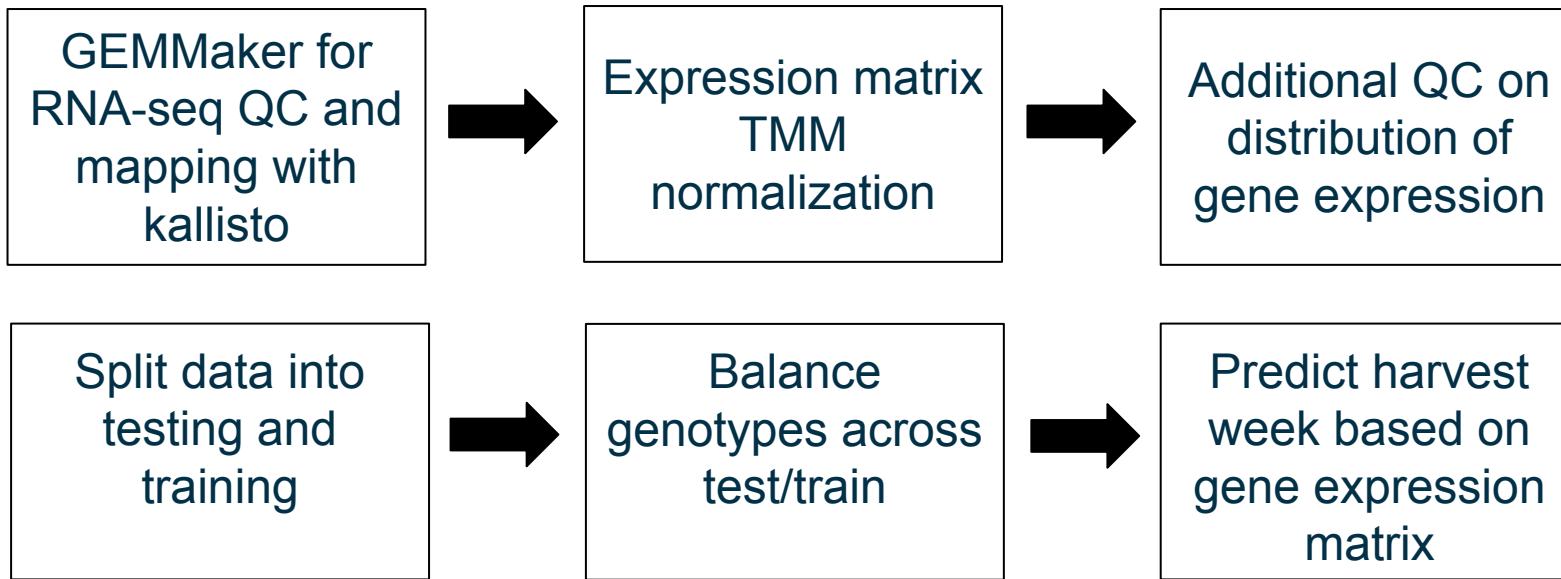
- 14 genotypes, 11 orchards, 3 years, and 6 time points relative to commercial harvest
- Fruit RNA-seq 3 bio replicates ~15 million read pairs
- 276 total libraries



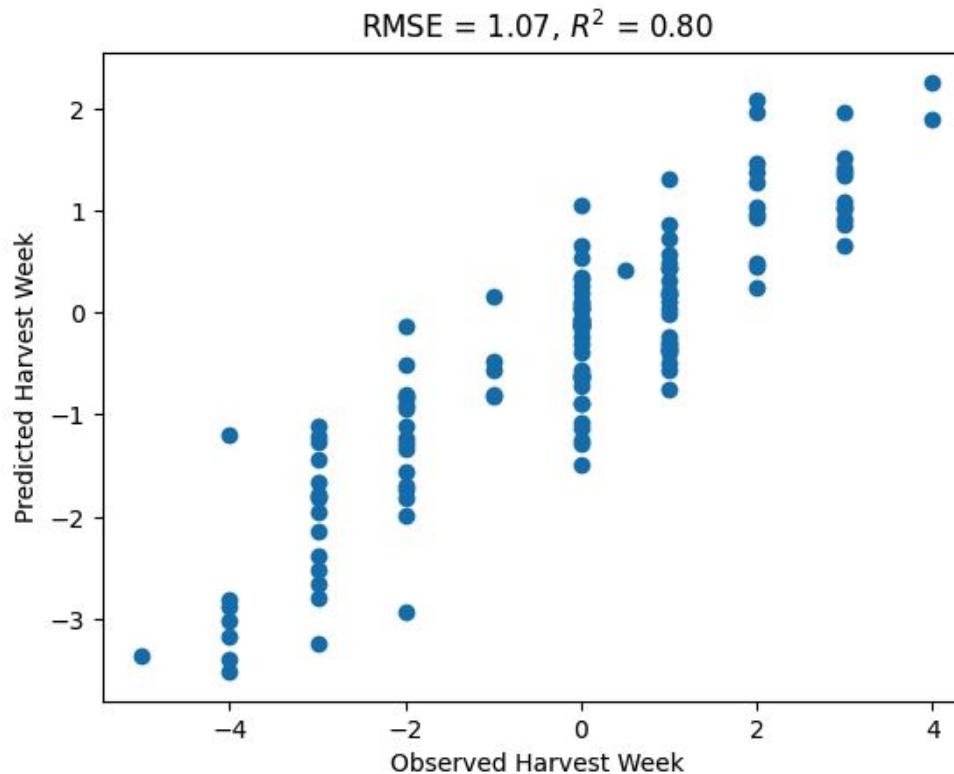


# GEMmaker

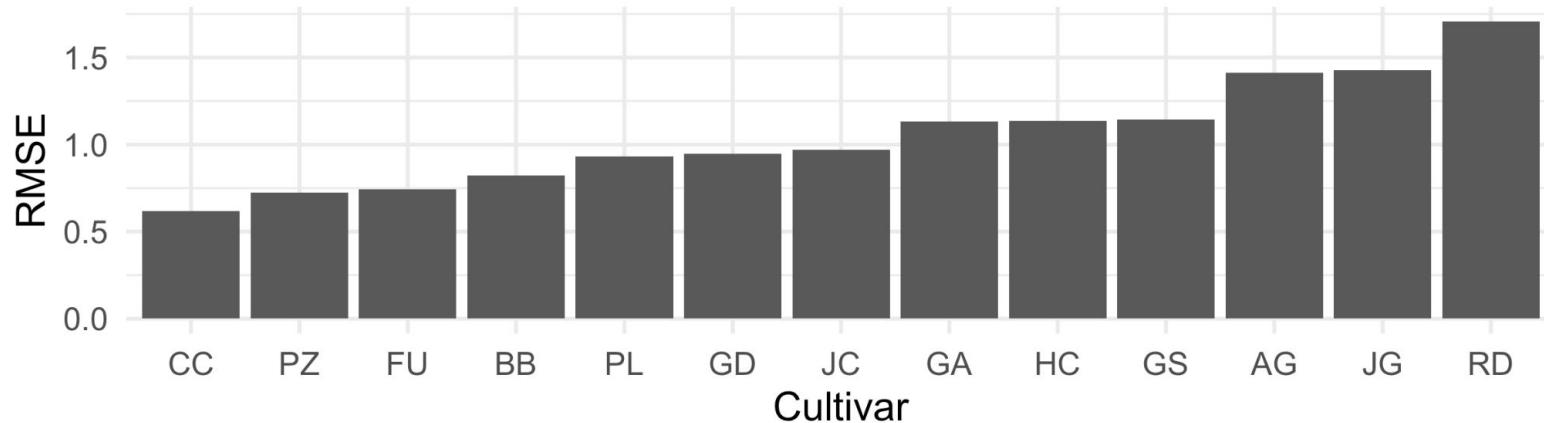
## Workflow



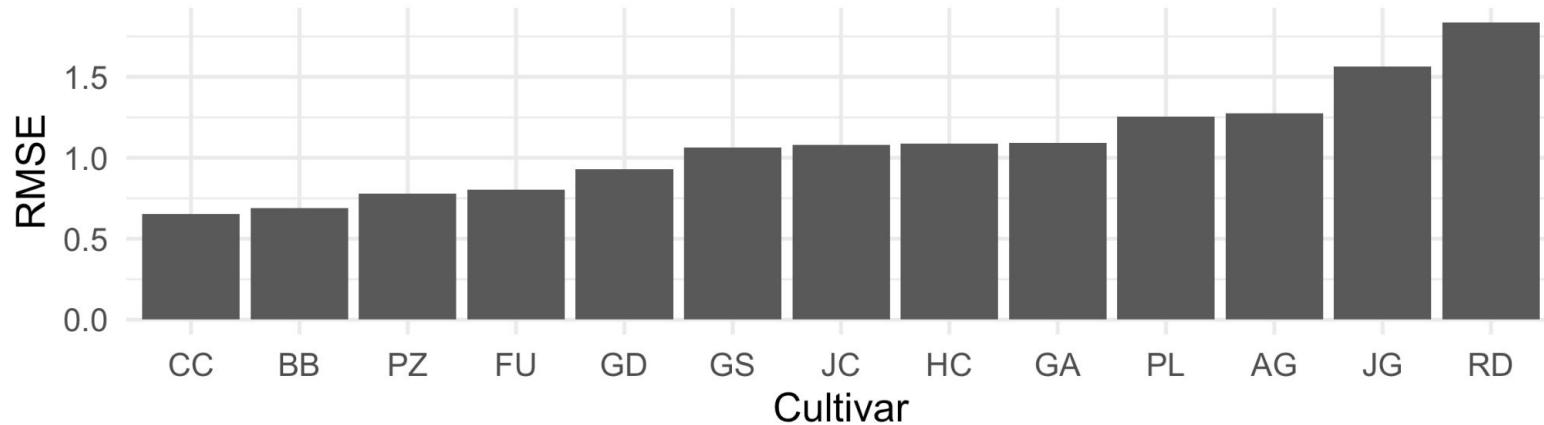
# We can predict relative harvest week!



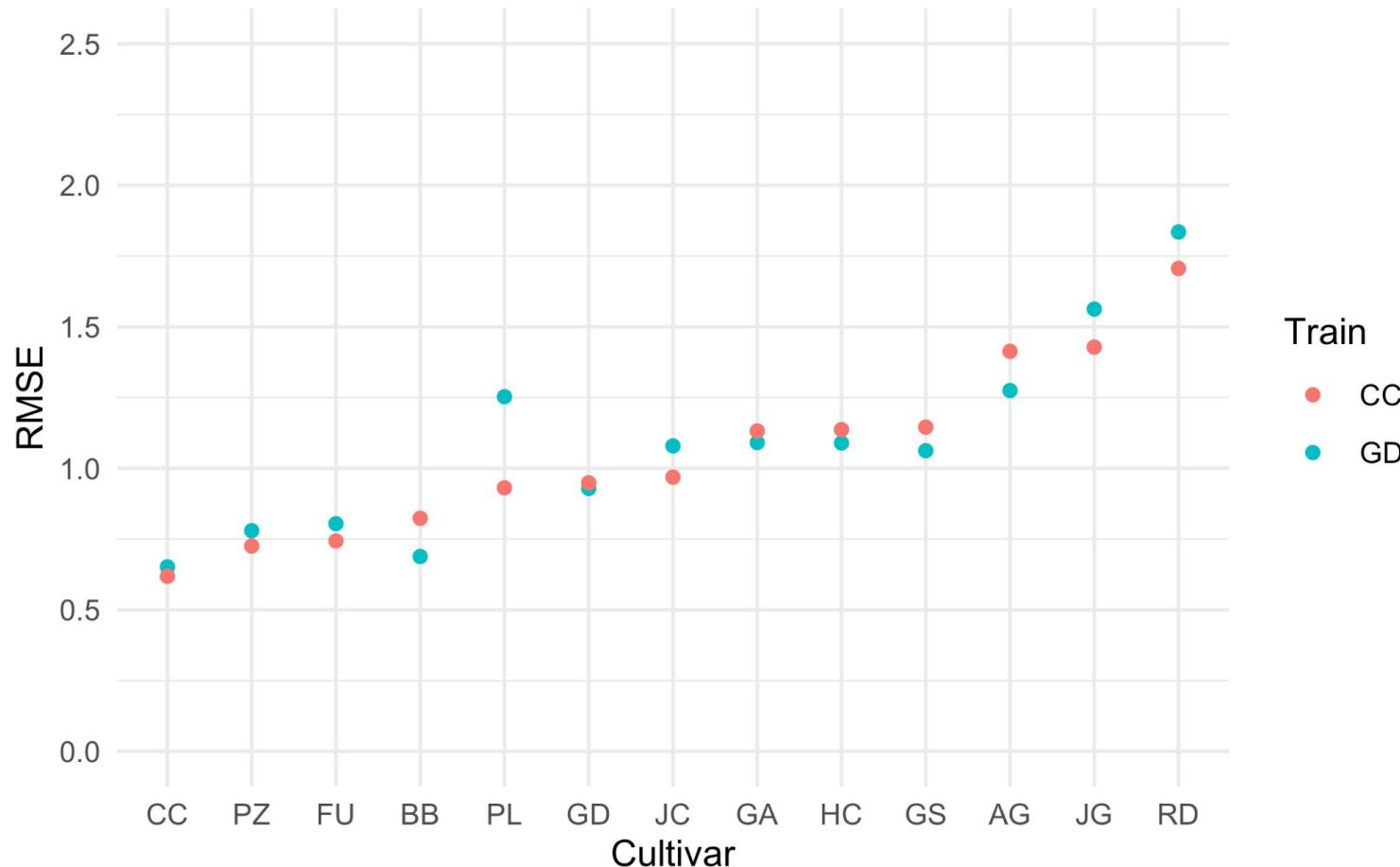
### RF RMSE on test data for model trained on CC



### RF RMSE on test data for model trained on GD



## RF RMSE on test data for model trained on CC



# Motivating Questions:

- How good is PacBio HiFi?
- Can we generate chromosome scale haplotype phased genome assemblies?
- What is the genetic diversity of crop wild relatives?
- Can we detect recombination events using  $k$ -mers?
- When is the perfect time to pick an apple or pear?
- **Current progress in graph-based pangenomes**



# Graph-based pangenome

- Build haplotype graph between “Cosmic Crisp” haplotype A and haplotype B
- Visualize in sequenceTubeMap
- Identify alleles between haplotypes
- Map RNA-seq to the graph



# Identifying alleles between haplotypes

- MCScanX
  - Blast, identify syntenic blocks
  - 56% of genes between hapA/B “orthologs”
- Cactus-pangenome
  - Align haplotypes with minimap2
  - Intersect genome annotations on the pangenome graph
  - 73.3% of genes between hapA/B “orthologs”

# Mapping RNA-seq... a work in progress



**PROPOSED DURATION:** 2 Years

**Project Title:** Risk Assessment for Loss of Firmness During Storage for Gala

**PI:** Dr. Alan Yocca

**Organization:** USDA-ARS Tree Fruit Research Lab

**Telephone:** 5096642280

**Email:** alan.yocca@usda.gov; ayocca@hudsonalpha.org

**Co-PI:** Dr. Alex Harkess

**Organization:** HudsonAlpha Institute of Biotechnology

**Telephone:** 2563270475

**Email:** aharkess@hudsonalpha.org

**Cooperators:** \_\_\_\_\_

**Year 1:** \$ 115,000.00

**Year 2:** \$ 90,000.00

**Year 3:** \_\_\_\_\_

**Total Budget:** \$ 205,000.00

