

BRIDGEcereal streamlines unsupervised learning to graph indel-based haplotype from pan-genome*

Xianran Li

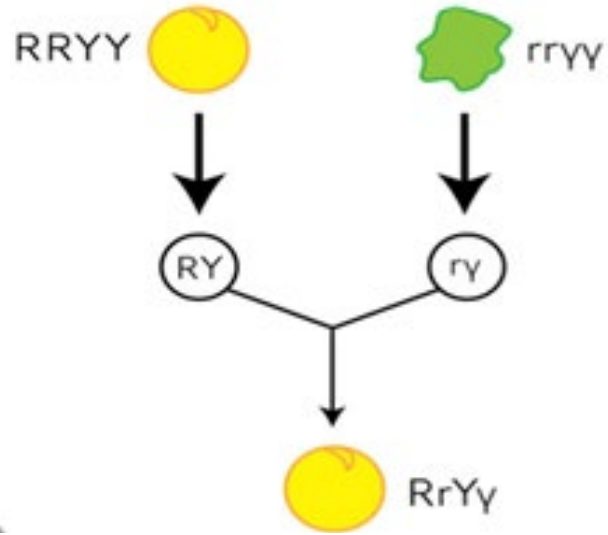
USDA-ARS Wheat Health, Genetics, and Quality

Pullman, WA

<https://compbiolab.org/>

*How to find large indel polymorphisms for your favorite genes from available pan-genome

Genetics: phenotype & casual DNA polymorphisms



MENDEL'S LAW OF INHERITANCE



Seed		Flower		Pod		Stem	
Form	Cotyledons	Color	Form	Color	Place	Size	
Round	Yellow	White	Full	Yellow	Axial pods, Flowers along	Long (6-7ft)	1
Wrinkled	Green	Violet	Constricted	Green	Terminal pods, Flowers top	Short (¾-1ft)	2

Cell

Supports open access

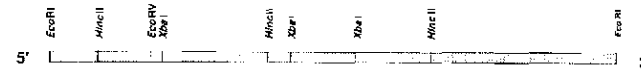
The first cloned Mendel's gene segregates a large insertion (0.8-kb)!

ARTICLE | VOLUME 60, ISSUE 1, P115-122, JANUARY 12, 1990 [Download Full Issue](#)

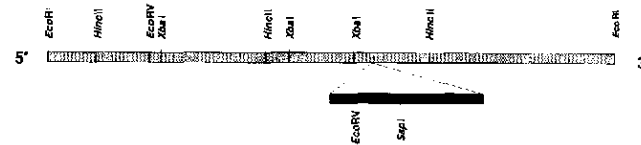
The wrinkled-seed character of pea described by Mendel is caused by a transposon-like insertion in a gene encoding starch-branching enzyme

Madan K. Bhattacharyya • Alison M. Smith • T.H.Noel Ellis • Cliff Hedley • Cathie Martin

A

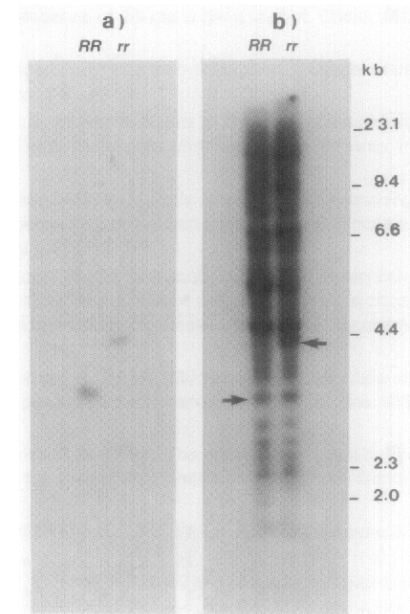


pJSBE102

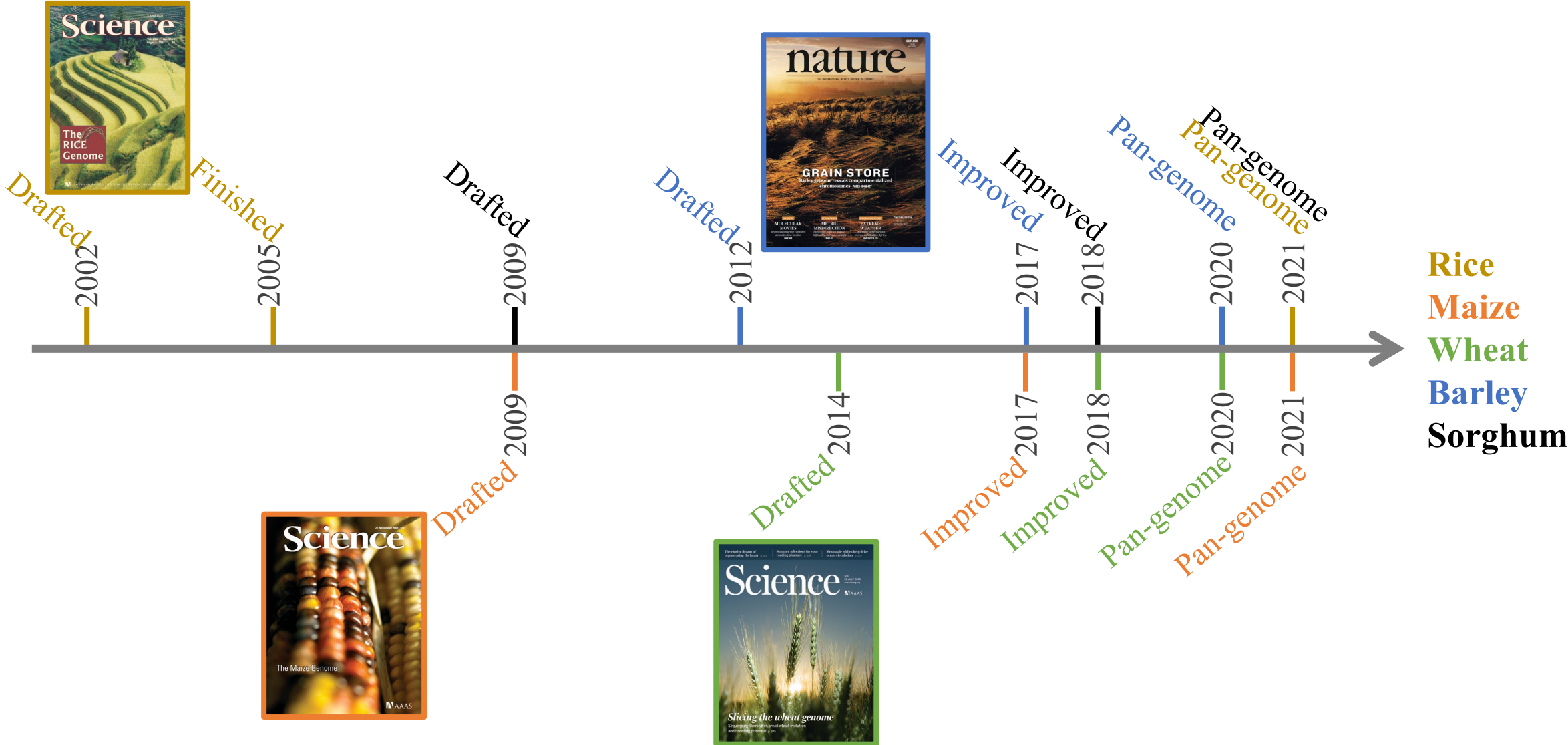


pJSBE206

B



Two decades of cereal genome sequencing: from single reference to pan-genome

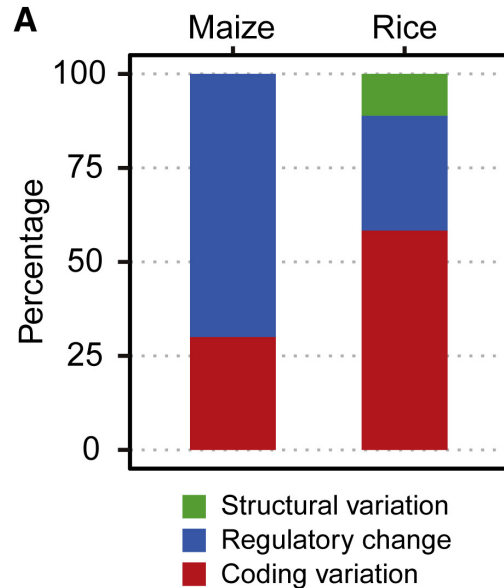


REVIEW

How the pan-genome is changing crop genomics and improvement!!!

Rafael Della Coletta¹, Yinjie Qiu¹, Shujun Ou², Matthew B. Hufford^{2*} and Candice N. Hirsch^{1*}

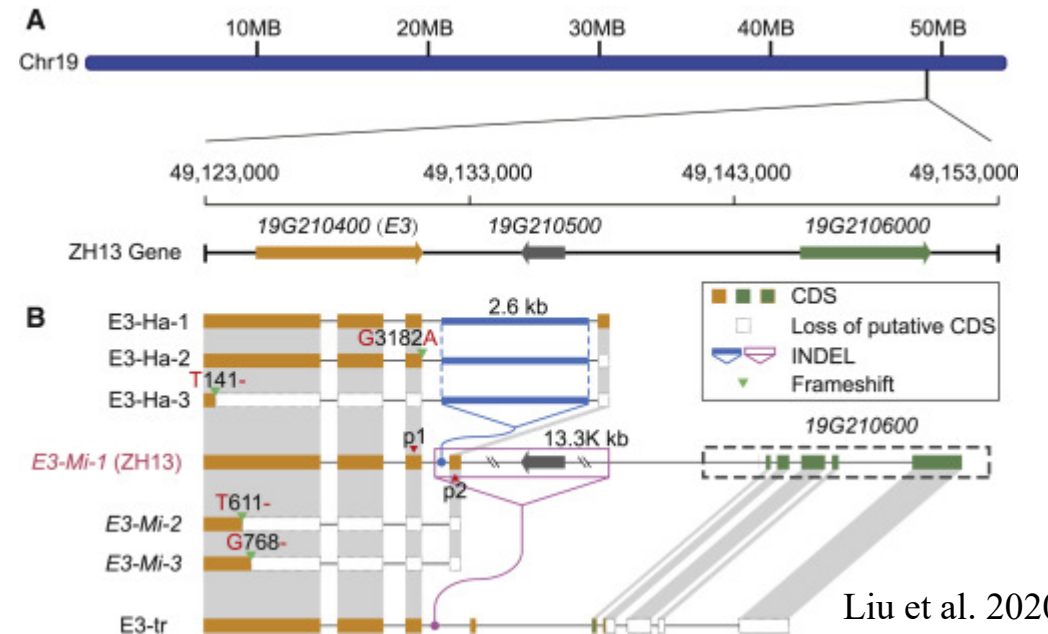
Half of identified casual polymorphisms are large indels.



Chen et al. 2021

Advancing QTL mapping and GWAS using pan-genomes

Haplotypes of the *E3* gene because of indels



Liu et al. 2020

How the pan-genome ^{is} changing crop genomics and improvement?

How to leverage pan-genomes for your genes?



ChatGPT

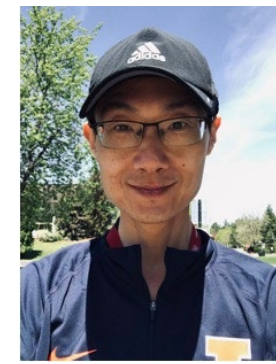
Li how to identify potential indels for a gene from pan-genomes

To identify potential indels for a gene from pan-genomes, you can follow these general steps:

1. Obtain the pan-genome data: This can be obtained from various sources, such as public databases or through sequencing and assembling the genomes of multiple individuals or strains of the same species.
2. Identify the gene of interest: Use bioinformatics tools to identify the gene of interest in the pan-genome data. This can be done by comparing the gene sequence to the pan-genome data using alignment tools such as BLAST or MUMmer.
3. Compare the gene sequence across the pan-genomes: Compare the gene sequence across the pan-genomes to identify any insertions or deletions (indels). You can use alignment tools such as MAFFT or Clustal Omega to align the gene sequences across the pan-genomes.
4. Identify potential indels: Look for gaps or differences in the alignment that indicate potential indels. These can be large or small and can be insertions or deletions relative to the reference sequence.
5. Validate potential indels: Once you have identified potential indels, you can validate them using various methods such as PCR, Sanger sequencing, or high-throughput sequencing techniques.

It's important to note that the analysis of pan-genomes can be complex and requires specialized expertise in genomics and bioinformatics. Additionally, the accuracy of the

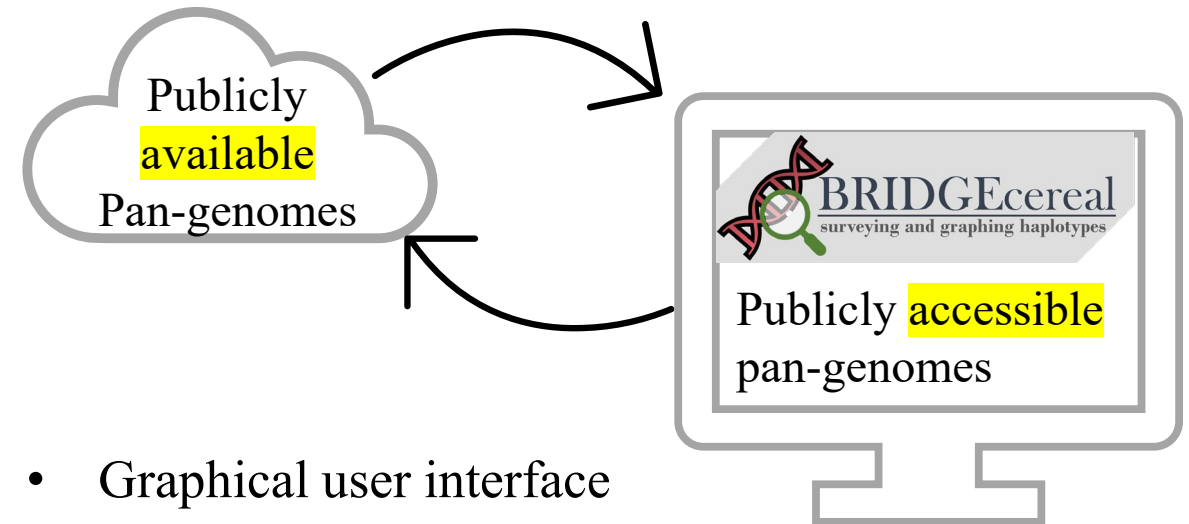
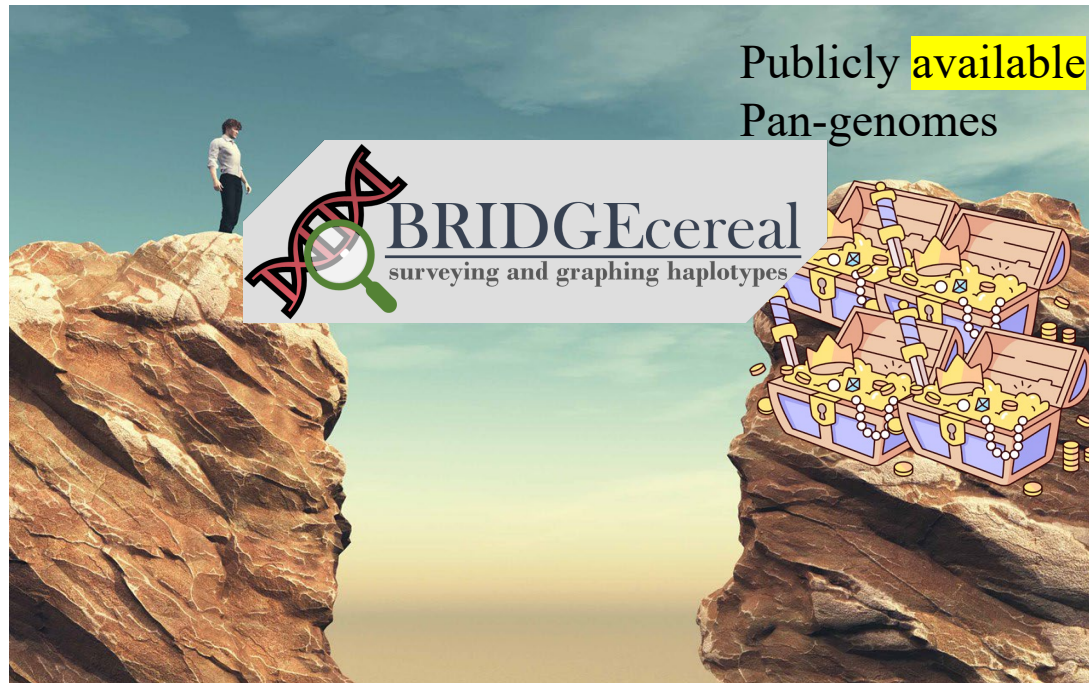
BRIDGEcereal: new bioinformatic bridge for identifying large indels



Bosen Zhang

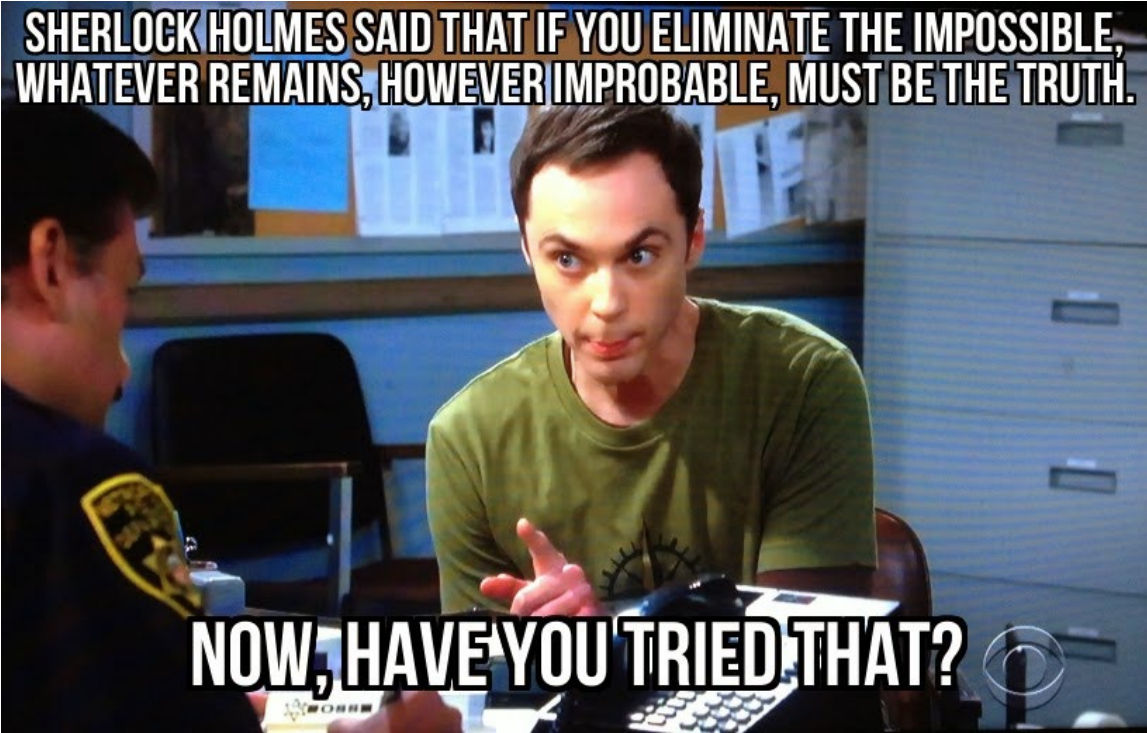


Laura Tibbs-Cortes

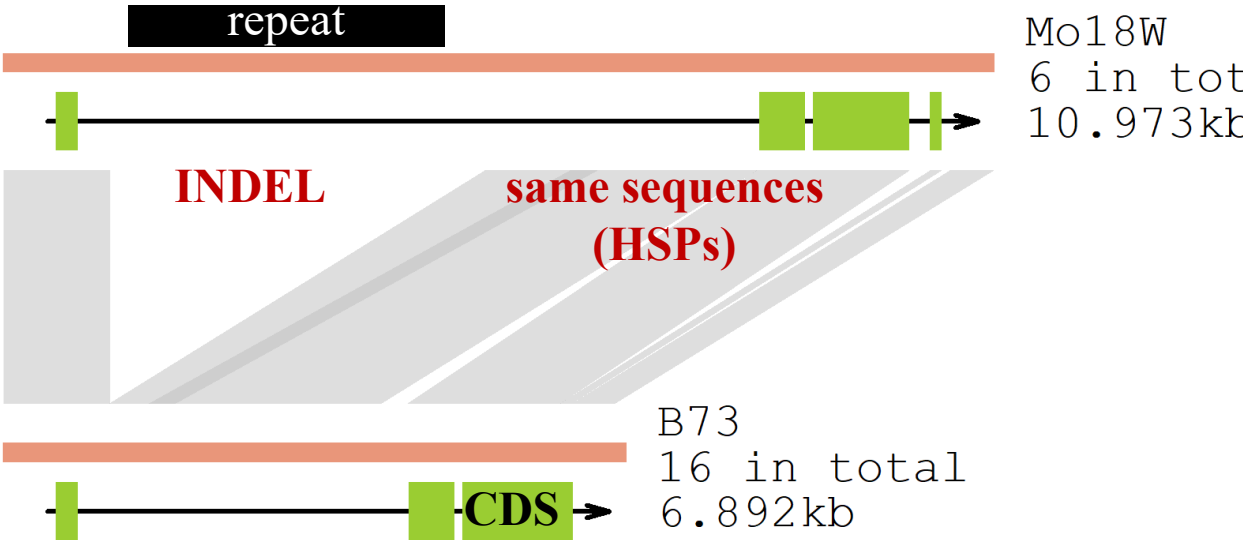


- Graphical user interface
- Minimal input for users
- Adjustable parameters

Blueprint of BRIDGEcereal



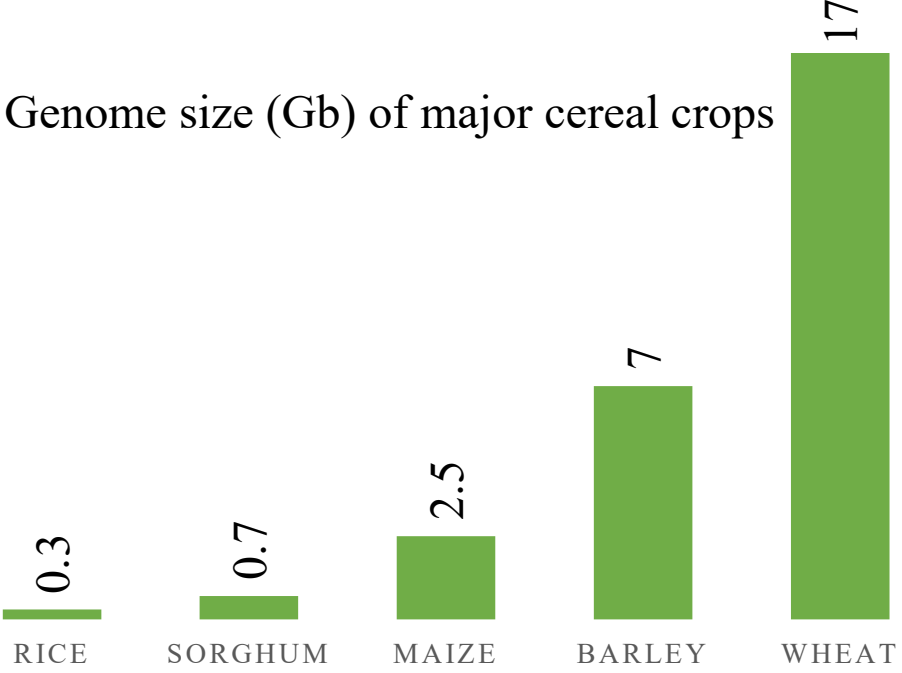
“If you eliminate **same sequences**, whatever remains, however improbable, must be **INDELs**”



Same sequences = HSP in BLAST

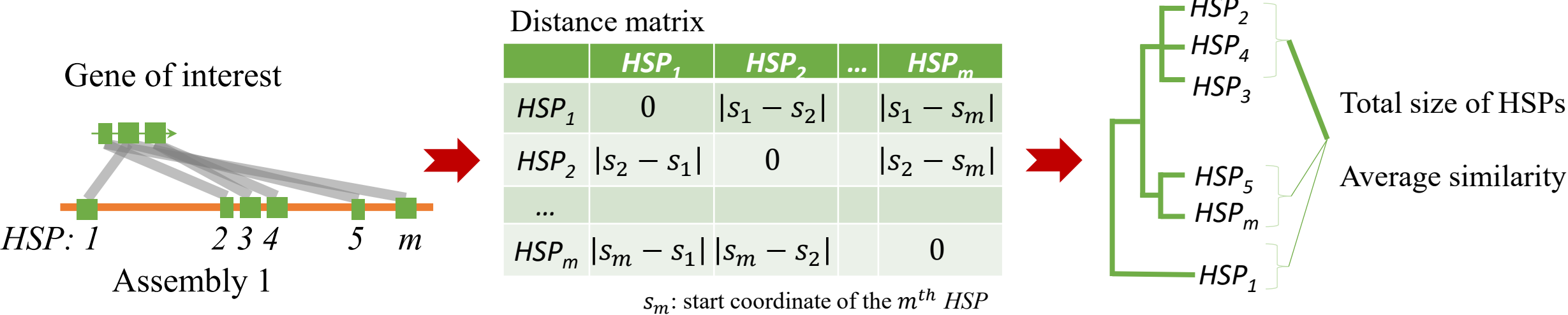
HSP: High-scoring Segment Pairs

Challenge 1: where is the homolog of a gene of interest in each assembly?



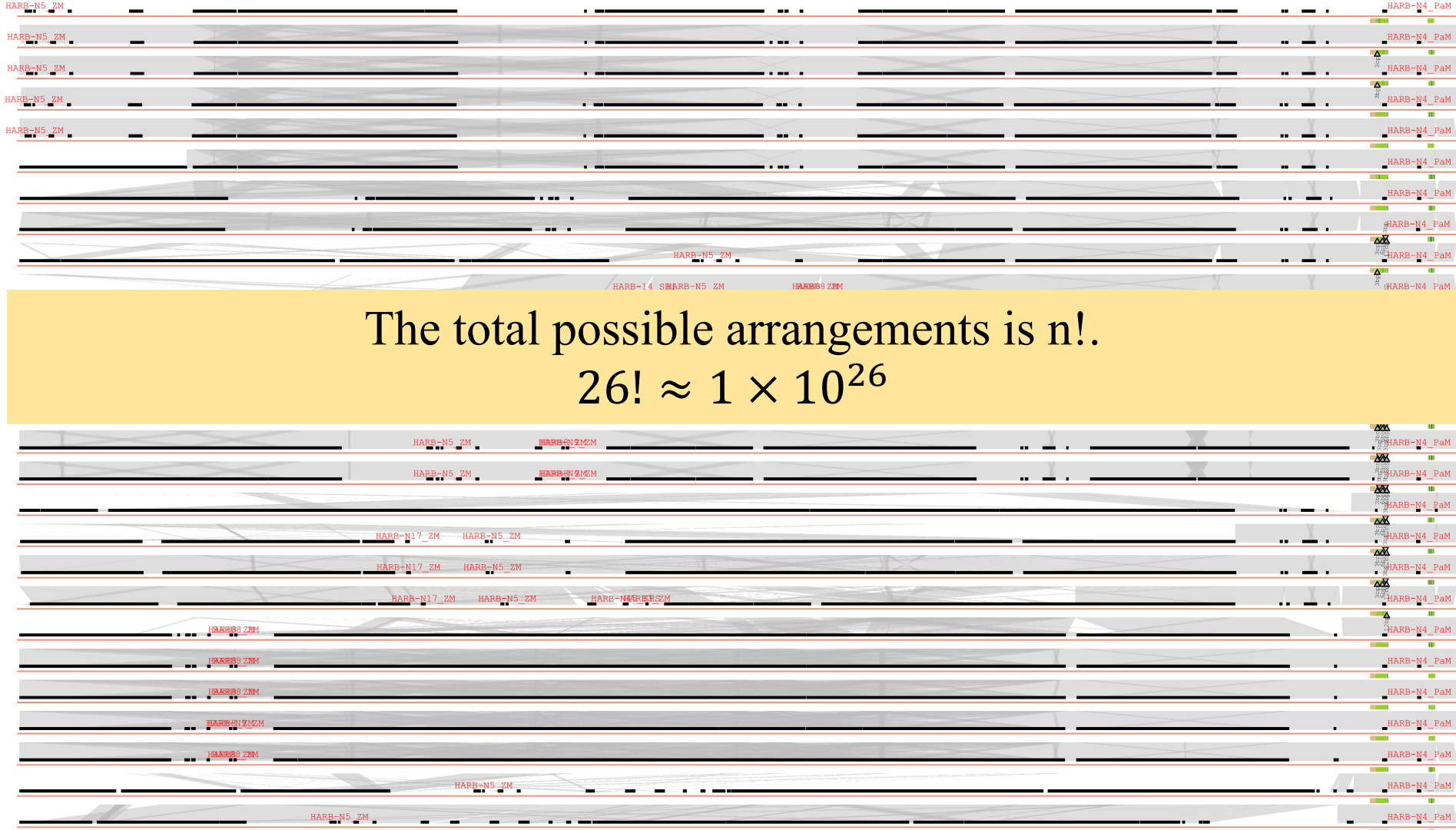
Gene annotation is typically only available for the first reference.

CHOICE to identify and extract segments harboring the ortholog from each assembly



Clustering HSPs for Ortholog Identification
via Coordinates and Equivalence

Challenge 2: Complex and indiscernible pattern from a large number assemblies

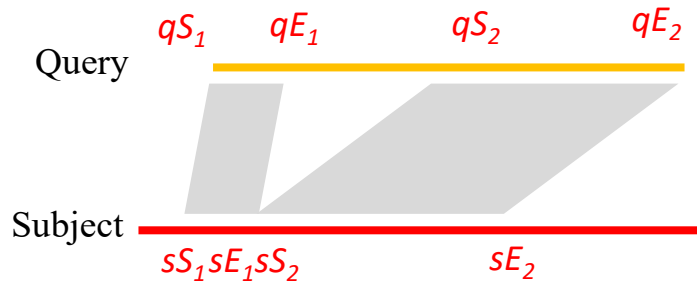


- Kill (63.85kb)
- NC358 (63.88kb)
- Tzi8 (63.88kb)
- CML333 (63.87kb)
- Mo18W (63.85kb)
- CML228 (63.89kb)
- CML322 (63.89kb)
- CML277 (63.85kb)
- CML69 (63.88kb)
- Ky21 (63.86kb)
- B97 (63.86kb)
- HP301 (63.86kb)
- Oh7B (63.88kb)
- Tx303 (63.88kb)
- I114H (63.88kb)
- P39 (63.88kb)
- Ki3 (63.85kb)
- M162W (63.85kb)
- M37W (63.85kb)
- CML247 (63.88kb)
- Oh43 (63.88kb)
- Ms71 (63.89kb)
- Mo17 (63.89kb)
- B73 (63.89kb)
- NC350 (63.89kb)
- CML103 (63.88kb)
- CML52 (63.88kb)

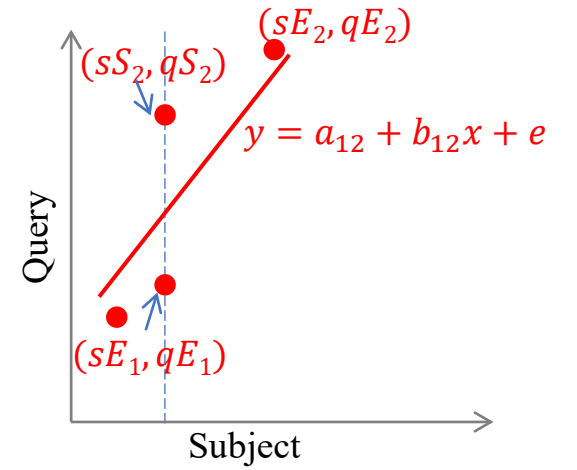
26 assemblies



CLIPS: Clustering via Large Indel Permuted Slopes

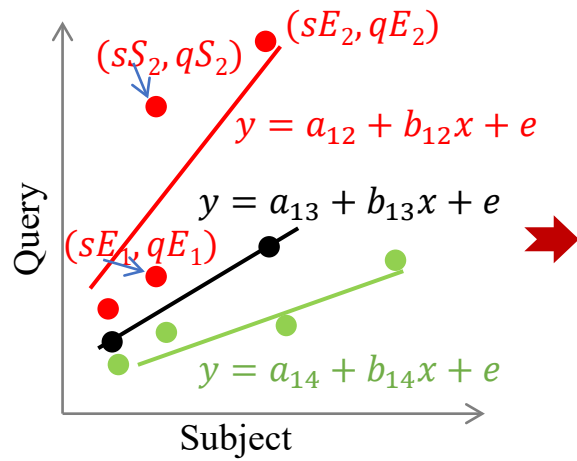
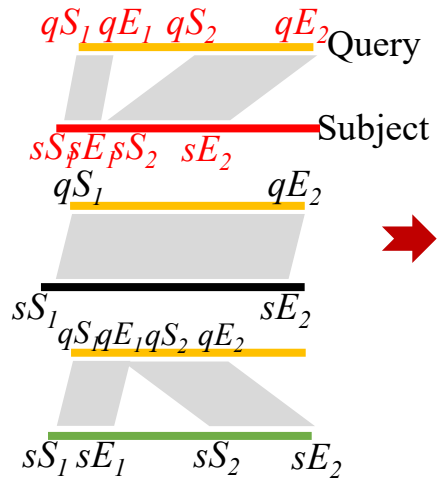


Query	$y = (qS_1, qE_1, qS_2, qE_2 \dots qS_n, qE_n)$
Subject	$x = (sS_1, sE_1, sS_2, sE_2 \dots sS_n, sE_n)$



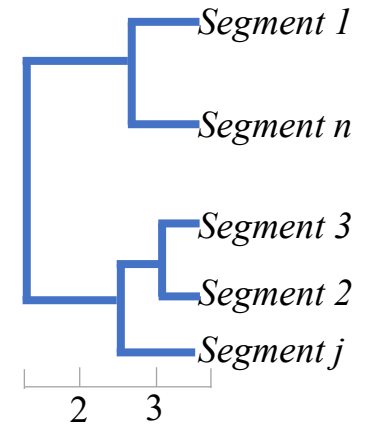


CLIPS: Clustering via Large Indel Permuted Slopes

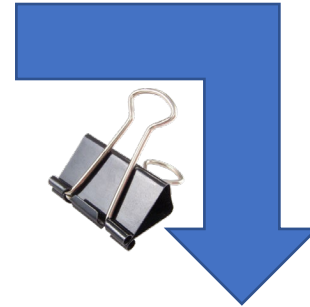
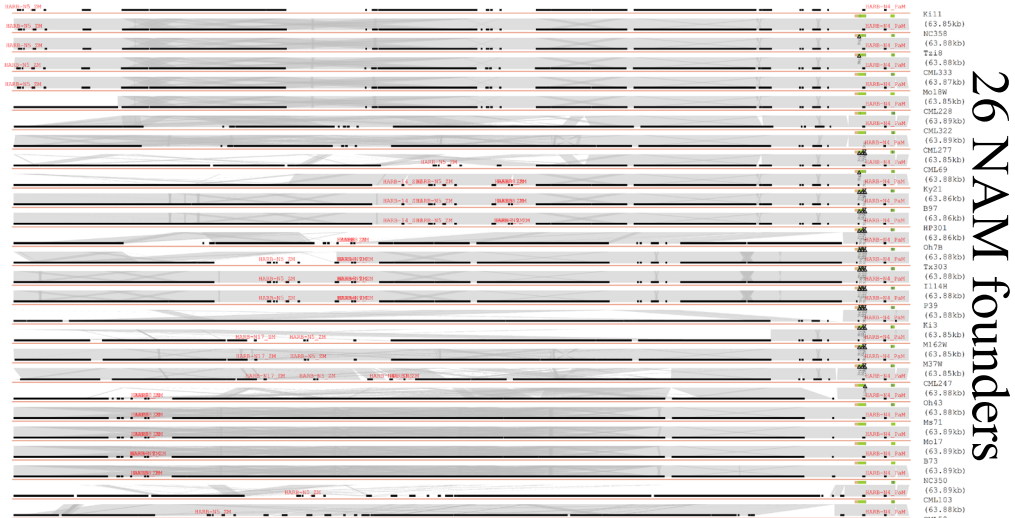


Slope matrix

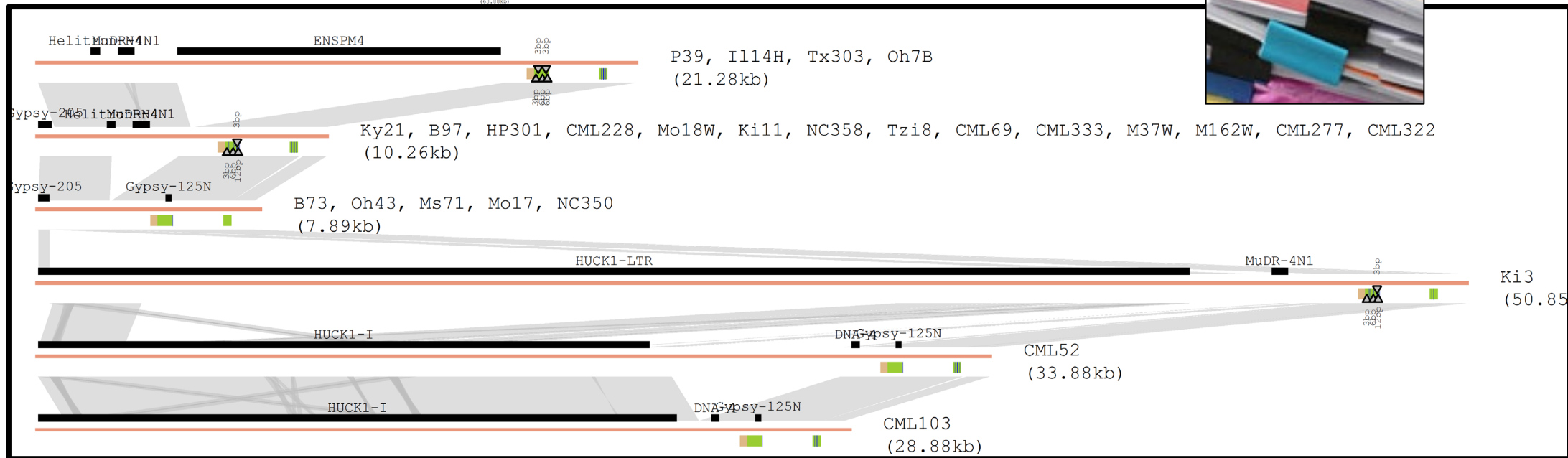
	Segment 1	Segment 2	...	Segment n
Segment 1	1	b_{12}		b_{1n}
Segment 2	b_{21}	1		b_{2n}
...				
Segment n	b_{n1}	b_{n2}		1

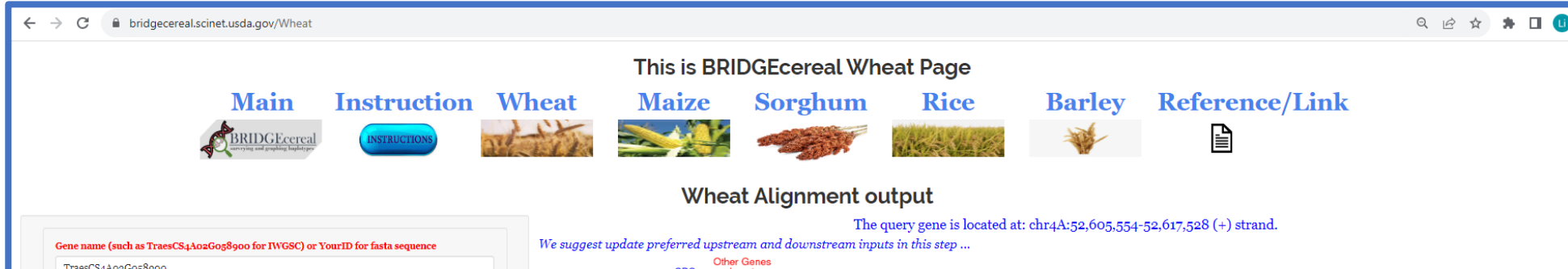


CLIPS clamps a complex gene pile into a concise and legible haplotype graph



6 haplotypes





This is BRIDGEcereal Wheat Page

Main Instruction Wheat Maize Sorghum Rice Barley Reference/Link

Wheat Alignment output

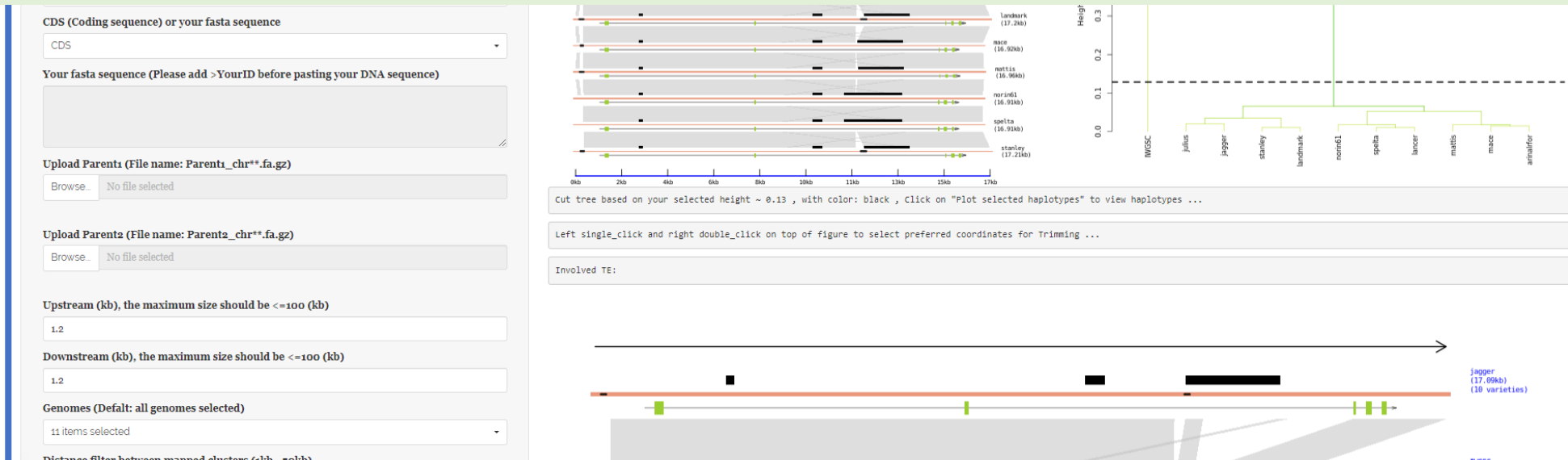
The query gene is located at: chr4A:52,605,534-52,617,528 (+) strand.

We suggest update preferred upstream and downstream inputs in this step ...

Gene name (such as TraesCS4A02G058900 for IWGSC) or YourID for fasta sequence

TraesCS4A02G058900

- One input (A gene model or CDS sequence) + 4 Clicks + 30 seconds
- Five major cereal crops (Wheat, Barley, Maize, Sorghum, Rice) with 120 assemblies



CDS (Coding sequence) or your fasta sequence

CDS

Your fasta sequence (Please add >YourID before pasting your DNA sequence)

Upload Parents (File name: Parent1_chr**.fa.gz)

Browse... No file selected

Upload Parents (File name: Parent2_chr**.fa.gz)

Browse... No file selected

Upstream (kb), the maximum size should be <=100 (kb)

1.2

Downstream (kb), the maximum size should be <=100 (kb)

1.2

Genomes (Default: all genomes selected)

11 items selected

Distance filter between mapped clusters (kb - 50kb)

landmark (17.28kb)

maize (16.92kb)

mattis (16.96kb)

norin61 (16.93kb)

spelta (16.93kb)

stanley (17.21kb)

Height

0.0 0.1 0.2 0.3

landmark norin61 spelta lance mattis maize annafor

cut tree based on your selected height ~ 0.13 , with color: black , Click on "Plot selected haplotypes" to view haplotypes ...

Left single_click and right double_click on top of figure to select preferred coordinates for Trimming ...

Involved TE:

jagger (17.09kb) (10 varieties)

Demonstration on BRIDGEcereal facilitating gene identification and characterization

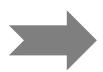
Population	Gene	Prior information Gene	Causal site	BRIDGEcereal results
Chinese Spring × Paragon RIL	<i>Rc-D1</i>	✓	✓	New recessive alleles
	<i>B1</i>	✓	✗	A 13-kb deletion 6-kb upstream
	<i>Hooded</i>	✗	✗	3-kb deletion removing 3 exons as a potential causal site

Wheat Chinese Spring × Paragon RIL population:

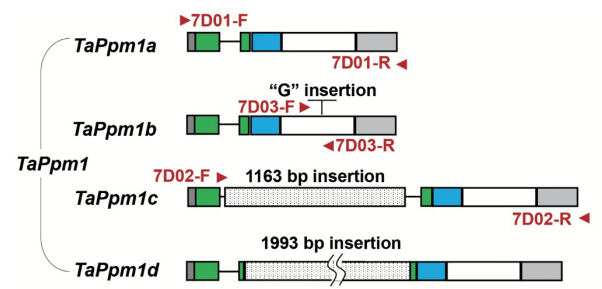
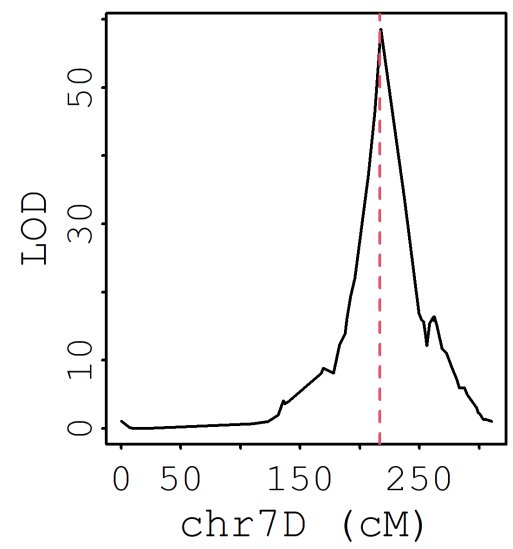
- Chinese Spring is the reference genome.
- Paragon was assembled at the scaffold level.

New recessive alleles of the well characterized wheat *Rc-D1* gene

Prior information	
Gene	Causal site
✓	✓



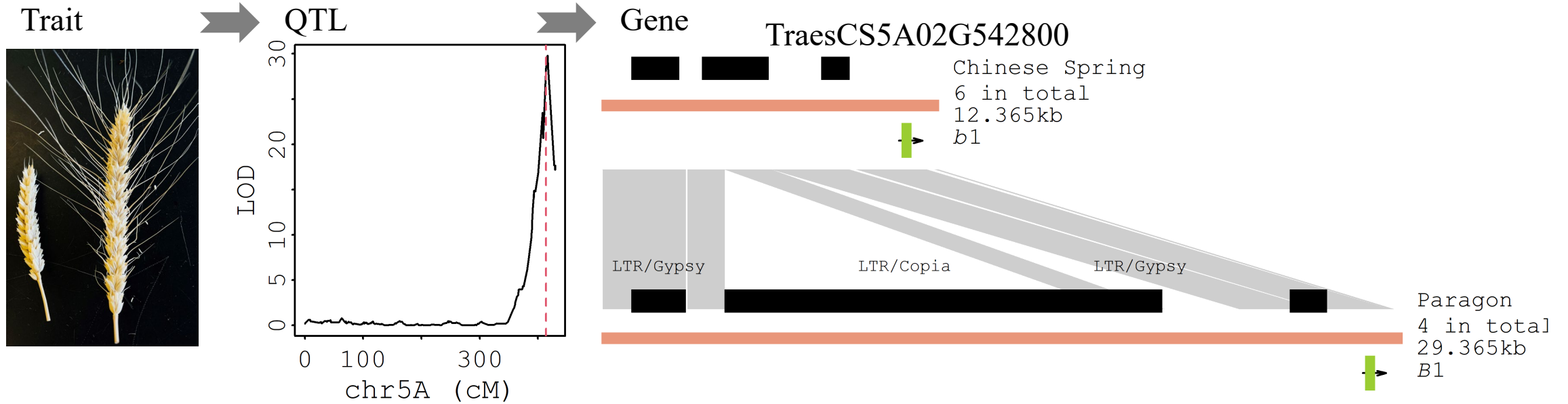
QTL



Jiang et al. 2018

Two large indels (13-kb + 4-kb) in 6-kb upstream as potential causal site for *B1*

Prior information	
Gene	Causal site
✓	✗

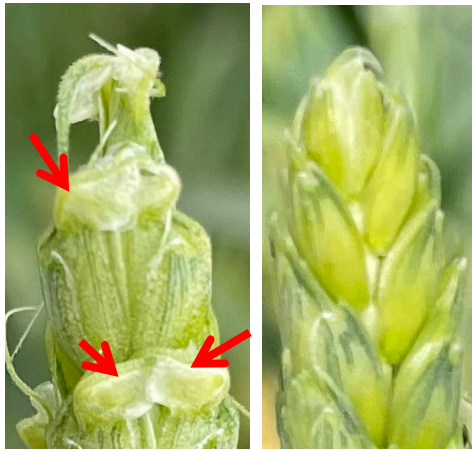


- In 2020, three papers reported TraesCS5A02G542800 as the underlying gene with the consensus of an unknown casual polymorphism outside of gene altering its expression.

TaDL is less likely as candidate for the classic wheat Hooded QTL

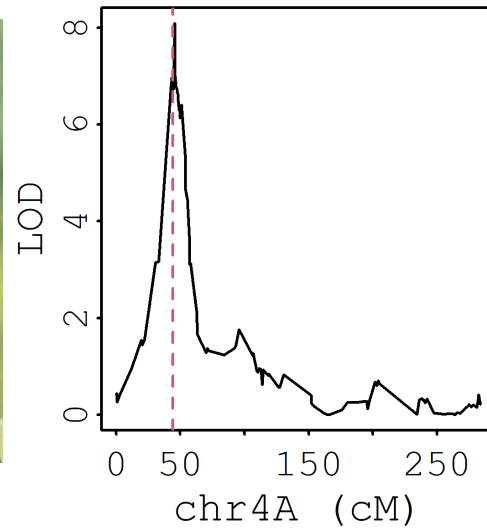
Prior information	
Gene	Causal site
×	×

Trait



Chinese Spring × Paragon

→ QTL



TraesCS4A02G058800 (*TaDL*)

```

TraesPAR_scaffold_035072_01G00      MQSMDLVSPSEHLCYVRCTYCNTVLA VGVPC KRLMDT VTVKCGHCNNLSF
TraesCS4A02G058800                  MQSMDLVSPSEHLCYVRCTYCNTVLA VGVPC KRLMDT VTVKCGHCNNLSF
*****

TraesPAR_scaffold_035072_01G00      LSPRPPPMVQPLSPNDHHHPMGPFQGWTD CRRNQPLPPLASPTSSDASPR
TraesCS4A02G058800                  LSPRPPPMVQPLSPNDHHHPMGPFQGWTD CRRNQPLPPLASPTSSDASPR
*****

TraesPAR_scaffold_035072_01G00      APFVVKPPEKKHRLPSAYNRFMREEIQRIKAAKPDIPHREAFSMAAKNWA
TraesCS4A02G058800                  APFVVKPPEKKHRLPSAYNRFMREEIQRIKAAKPDIPHREAFSMAAKNWA
*****

TraesPAR_scaffold_035072_01G00      KCDPRCSTTVSASNSAPEPRIVVPGPQLQERATEQVVESFDIFKQMERSA
TraesCS4A02G058800                  KCDPRCSTTVSASNSAPEPRIVVPGPQLQERATEQVVESFDIFKQMERSA
*****:*****
    
```



Two parental lines have not obvious causal functional polymorphisms

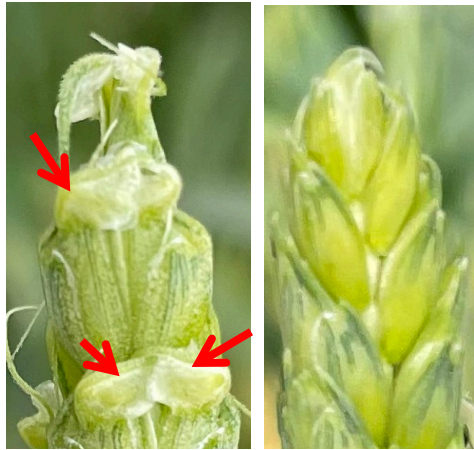
A *MADS*-box gene more likely underlying the classic wheat *Hooded* QTL

Prior information
Gene Causal site

×

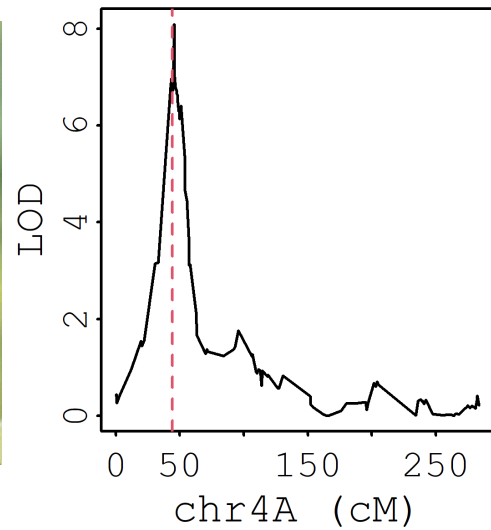
×

Trait

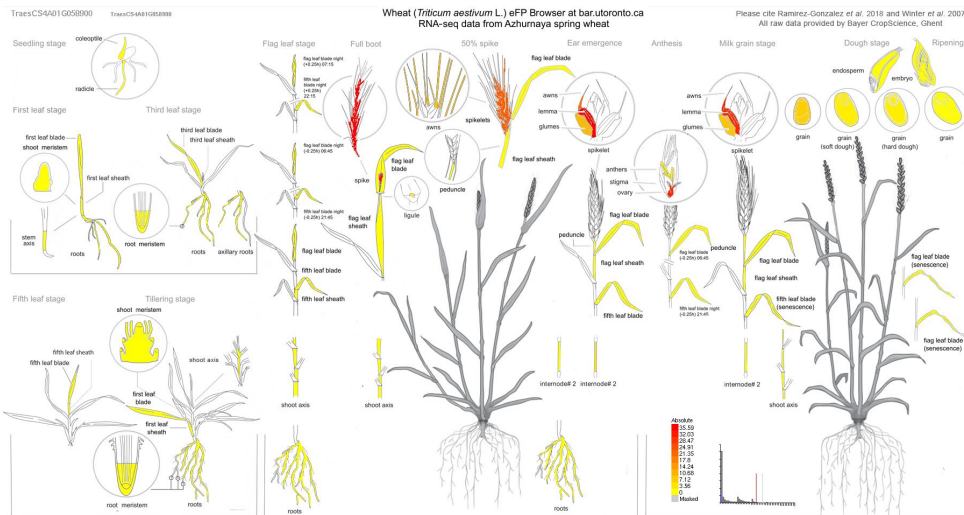
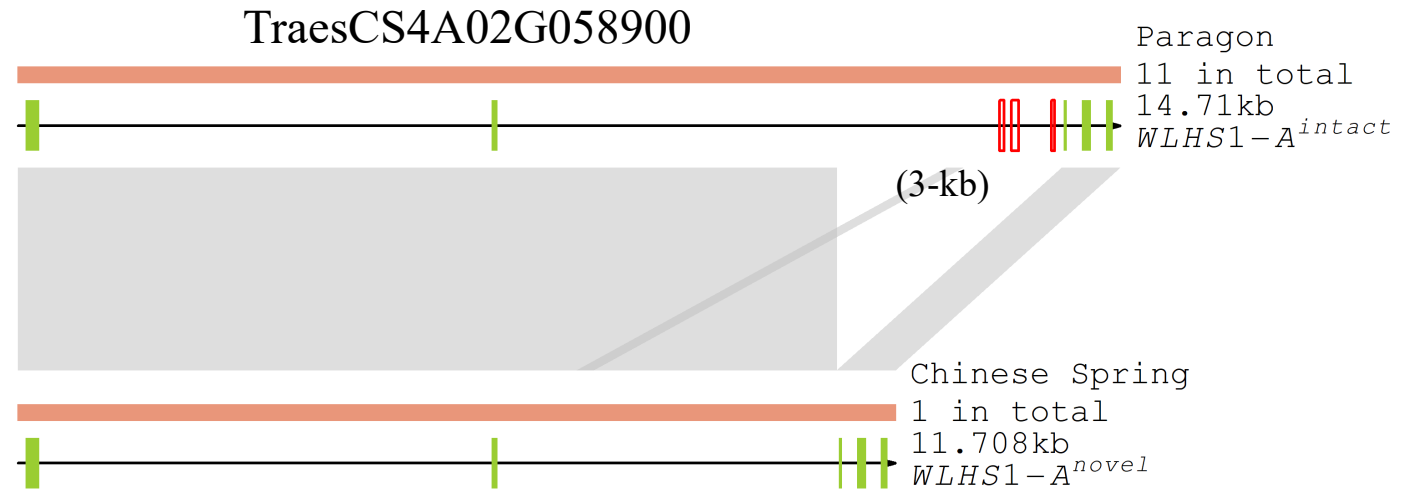


Chinese Spring × Paragon

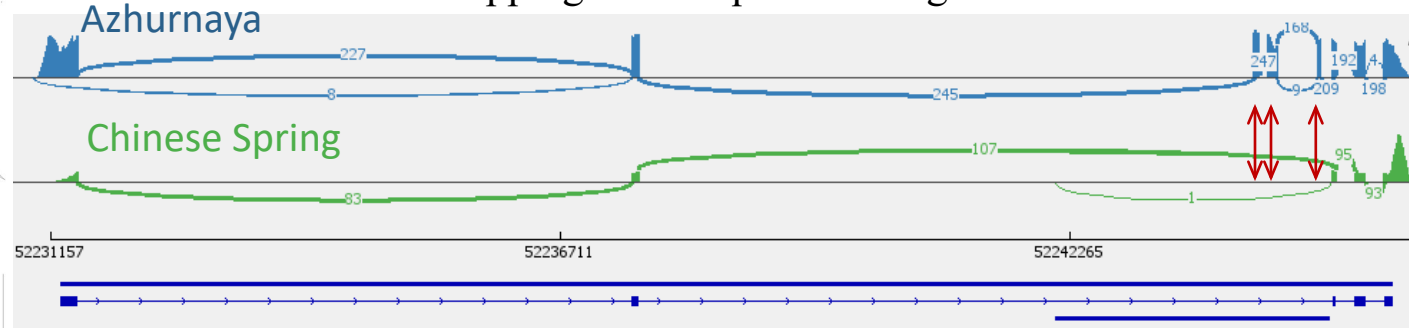
QTL

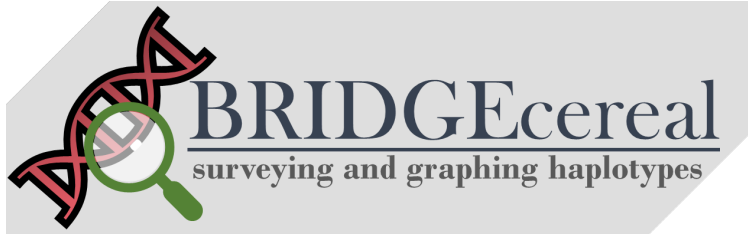


Gene



Mapping RNAseq reads to a genome without the deletion





<https://bridgecereal.scinet.usda.gov/>

①

Search “BRIDGEcereal” :

Google

Bing

yahoo!

②

Molecular Plant *Supports open access*

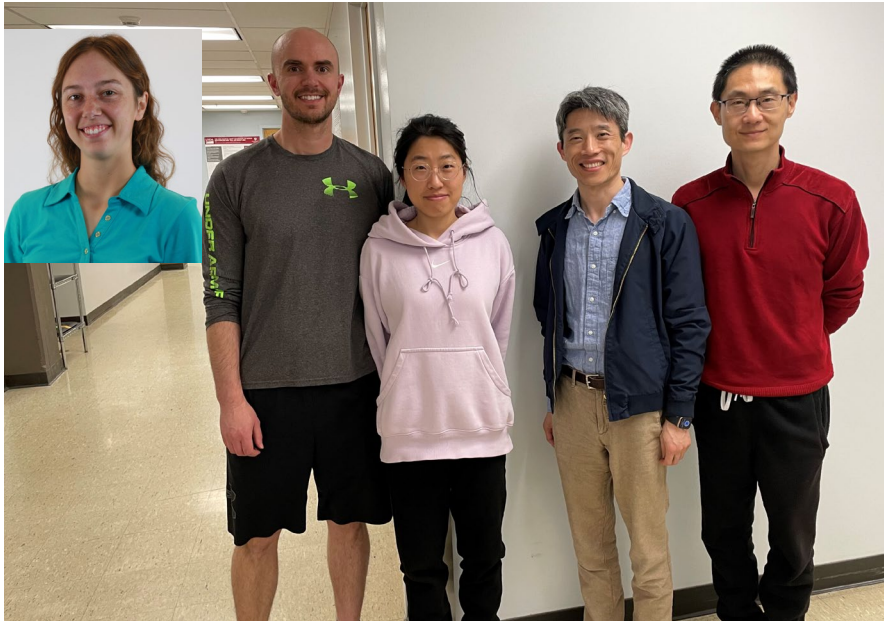
CORRESPONDENCE | VOLUME 16, ISSUE 6, P975-978, JUNE 05, 2023 [Download Full Issue](#)

Streamline unsupervised machine learning to survey and graph indel-based haplotypes from pan-genomes

Bosen Zhang • Haiyan Huang • Laura E. Tibbs-Cortes • Adam Vanous • Zhiwu Zhang • Karen Sanguinet • Kimberly A. Garland-Campbell • Jianming Yu • Xianran Li [✉](#) • [Show less](#)

[Open Access](#) • Published: May 17, 2023 • DOI: <https://doi.org/10.1016/j.molp.2023.05.005>

Acknowledgements



Compbio Lab

Bosen Zhang

Laura Tibbs-Cortes

Ryan Benke

Linqian Han

Haiyan Huang

Adam Vanous

Kimberly Garland-Campbell

IOWA STATE UNIVERSITY
Department of Agronomy

Jianming Yu



Zhiwu Zhang

Karen Sanguinet

Germplasm Resources Unit

..... a national capability supported by the BBSRC at the John Innes Centre

Chinese Spring × Paragon RILs

