

# Genotyping – By – Sequencing (GBS) in Crop Plants

Rajneesh Paliwal



RESEARCH  
PROGRAM ON  
DrylandCereals

## Topics presented

**GBS OVERVIEW AND SNP CALLING**

**GENOME WIDE ASSOCIATION STUDY (GWAS)**

**APPLICATION OF GBS-SNP INCLUDING GWAS**

# Genotyping by sequencing (GBS) in any large genome species requires reduction of genome complexity

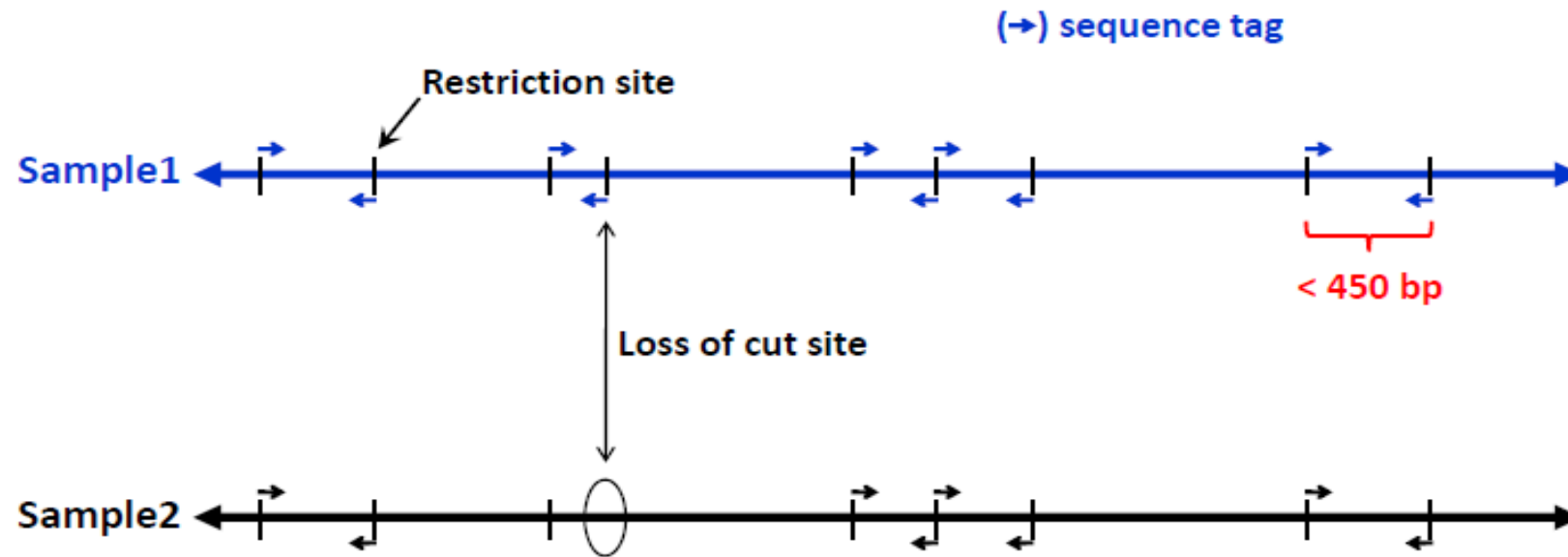
## Target enrichment

- Long range PCR of specific genes or genomic subsets
- Molecular inversion probes
- Sequence capture approaches hybridization-based (microarrays)

## Restriction Enzymes (REs)

- **\*Technically less challenging\***
- Methylation sensitive Res filter out repetitive genomic fraction

# Overview of Genotyping by Sequencing (GBS)



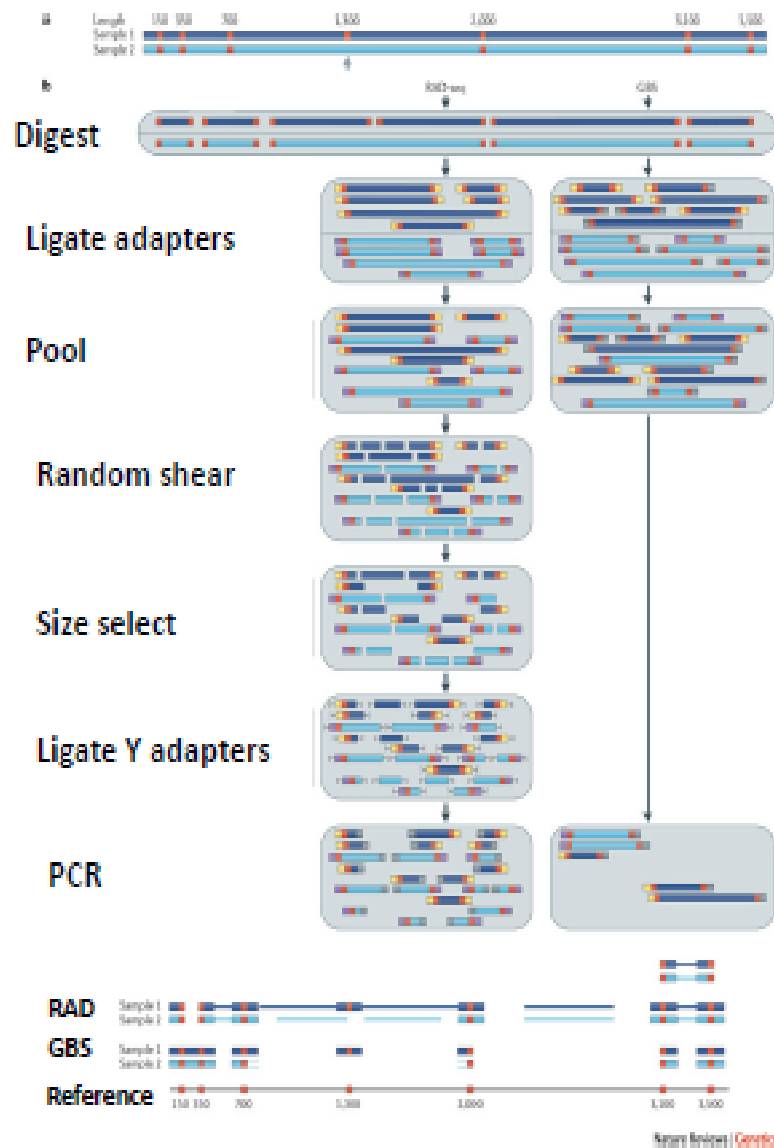
- Focuses NextGen sequencing power to ends of restriction fragments in both reference and non-reference genome plants
- Both SNPs and presence/absence markers can be scored
- Small Indels are identified but are not scored



# **GBS is a simple, highly multiplexed system for constructing libraries for NGS**

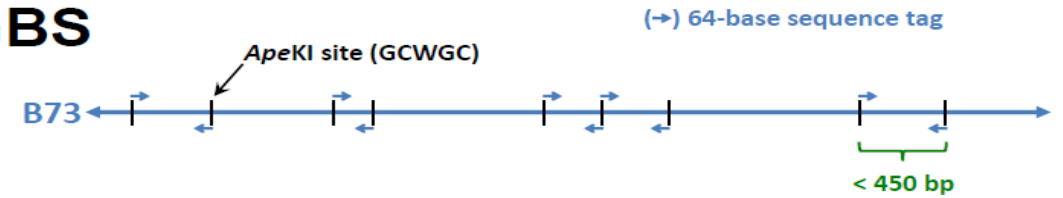
- Reduced Sample handling**
- Few PCR & purification steps**
- No DNA size fractionation**
- Efficient barcoding system**
- Simultaneous marker discovery and genotyping**
- Scales very well**

# RADs vs GBS SNPs

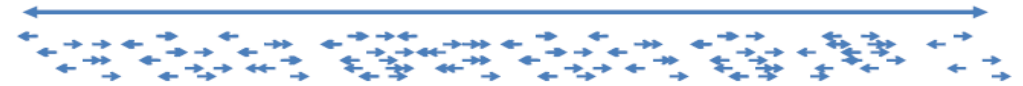


## Reduced Genome Representation through GBS

**GBS**



**WGS**



Davey et al. 2011

# Most frequently asked question for new species

**“How many SNP will I get?”**

**Answer: “It depends.....”**

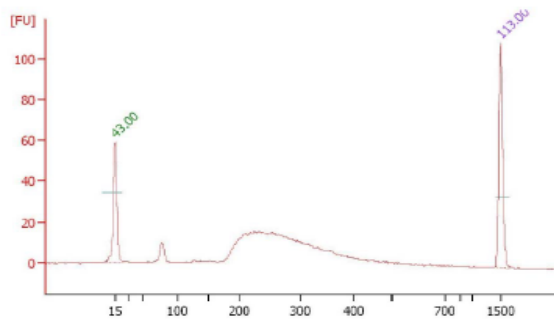
- Genome size and expected heterozygosity affects size of fragment pool for desired amount of sequence coverage (enzyme choice and multiplex level).
- Amount of extent diversity and how well your samples capture that diversity
- Reference genome sequence? 3-4x more SNPs attained by aligning to a reference sequence

Cont...

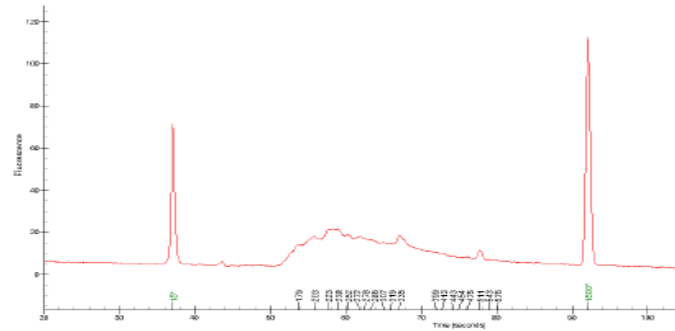
Species	Genome Size (Mb)	Enzyme	Sample Size	No. SNPs
Maize	2,600	<i>ApeKI</i>	33,000	2,200K
Rice	400	<i>ApeKI</i>	850	60K
Grape	500	<i>ApeKI</i>	8,000	1,200K
Willow*	460	<i>ApeKI</i>	459	23K
Pine*	16,000	<i>ApeKI</i>	12	63K
Vole*	3,400	<i>PstI</i>	283	53K
Fox*	2,400	<i>EcoT22I</i>	48	16K
Cow	3,000	<i>PstI</i>	48	64K
<i>Neurospora</i> (fungus isolates)	40	<i>ApeKI</i>	384	100K

\*No reference genome. UNEAK analysis pipeline used for analysis. To avoid homology/paralogy issues this pipeline calls SNPs very conservatively.

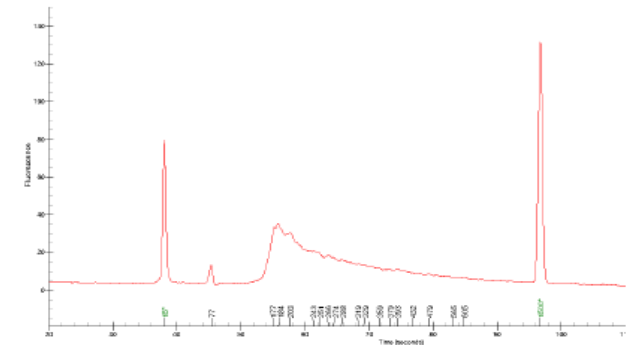
*PstI* works well for mammals and birds.



*EcoT22I* works well for fish, amphibians and invertebrates.



*ApeKI* works well for grasses.



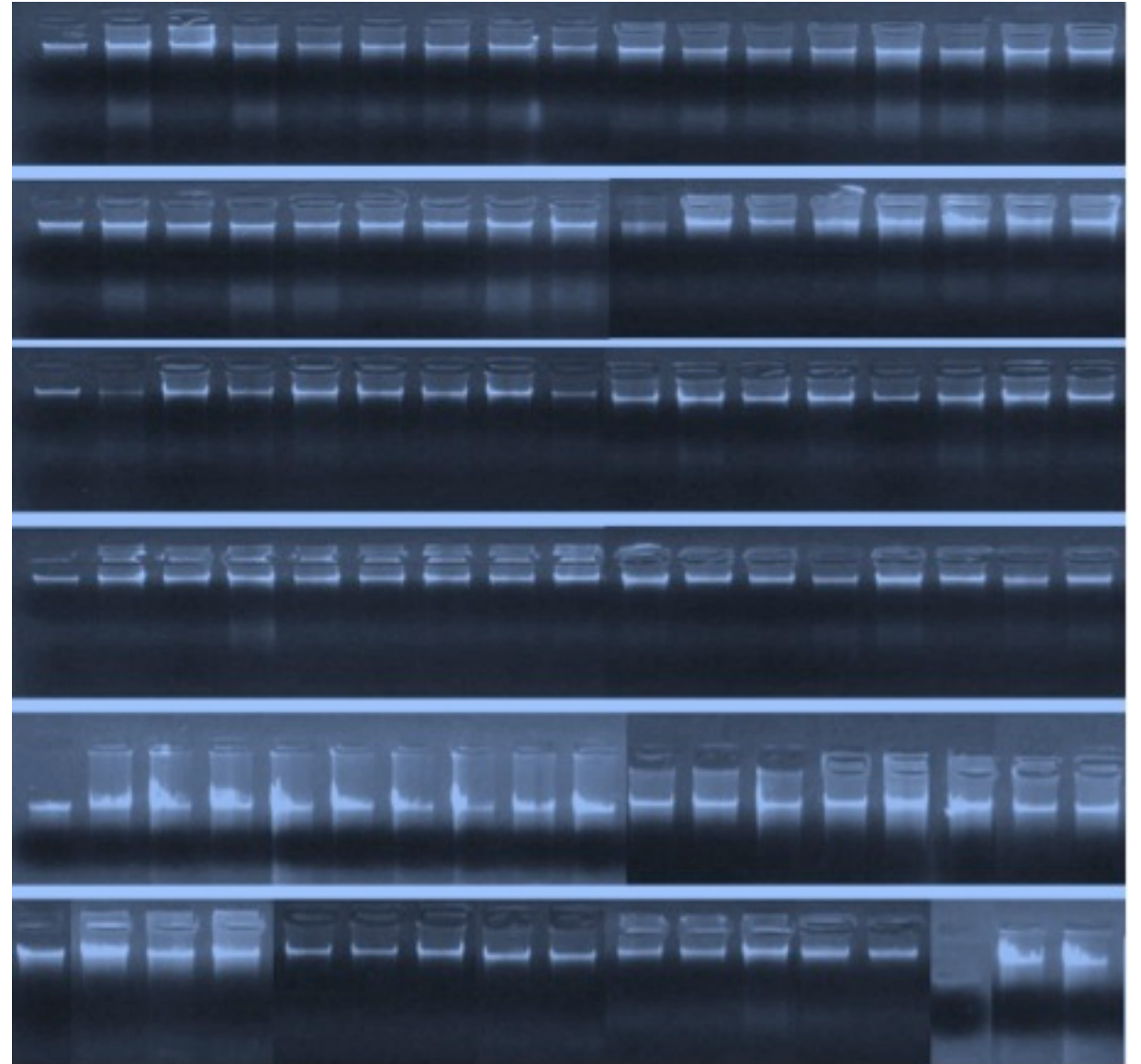
# Most significant GBS technical issues?

- DNA Quality
- DNA quantification

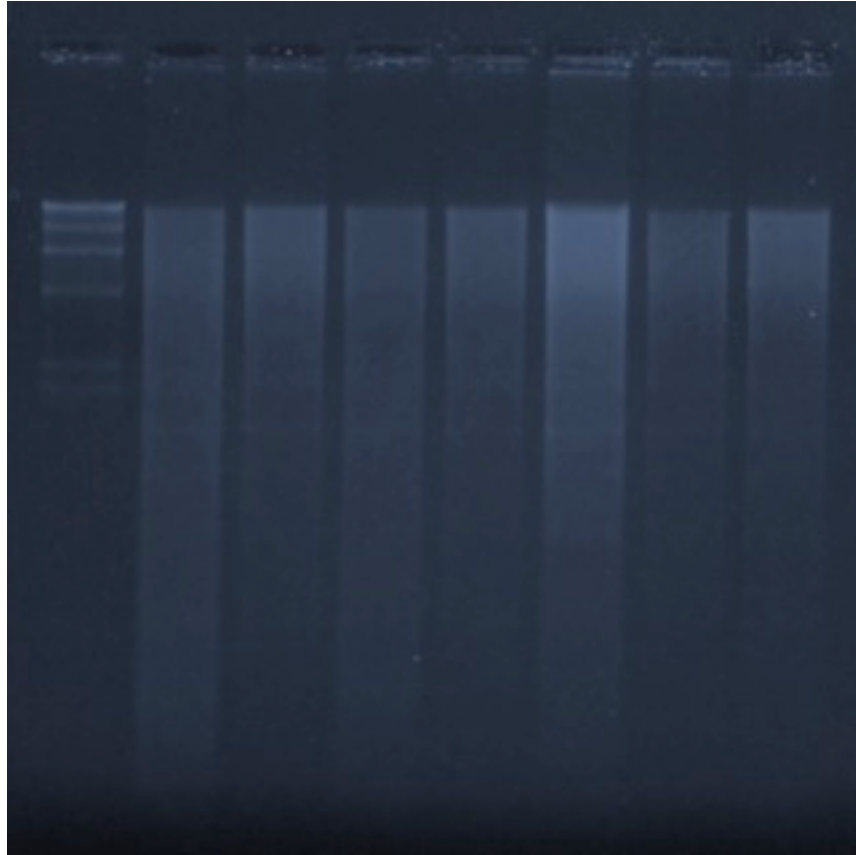
# SAMPLE PROCESSING

## DNA – Next Generation Sequencing

- High molecular weight
- Good quality

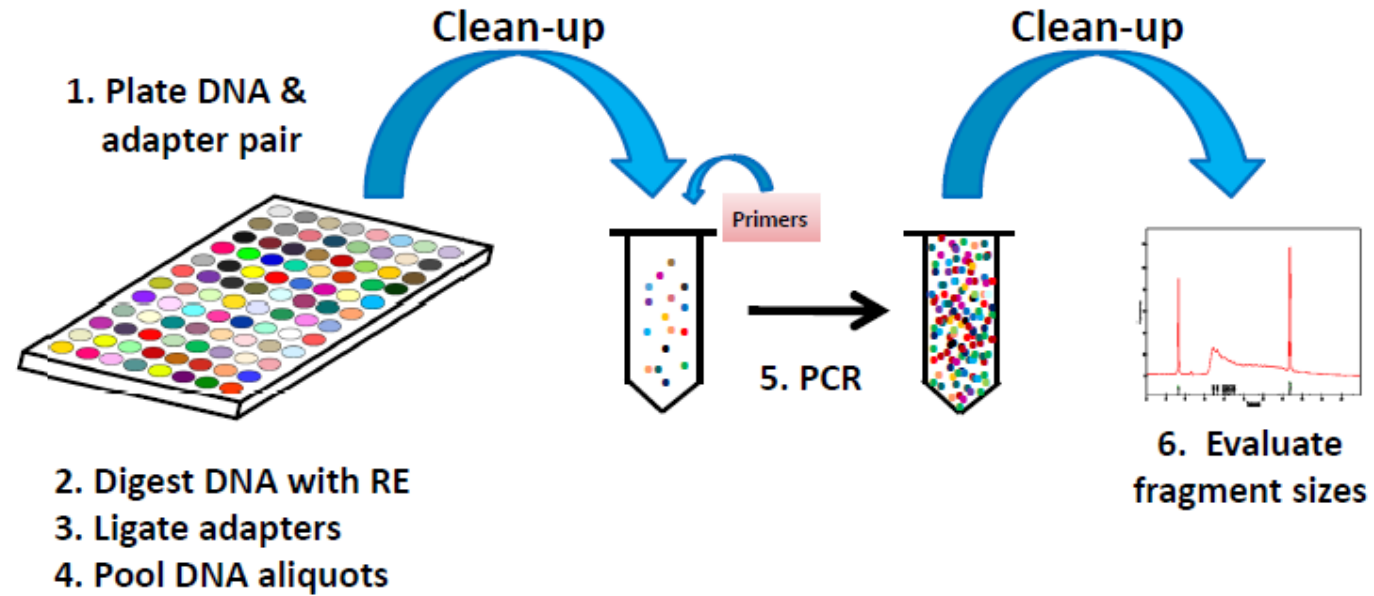


# Restriction Digestion



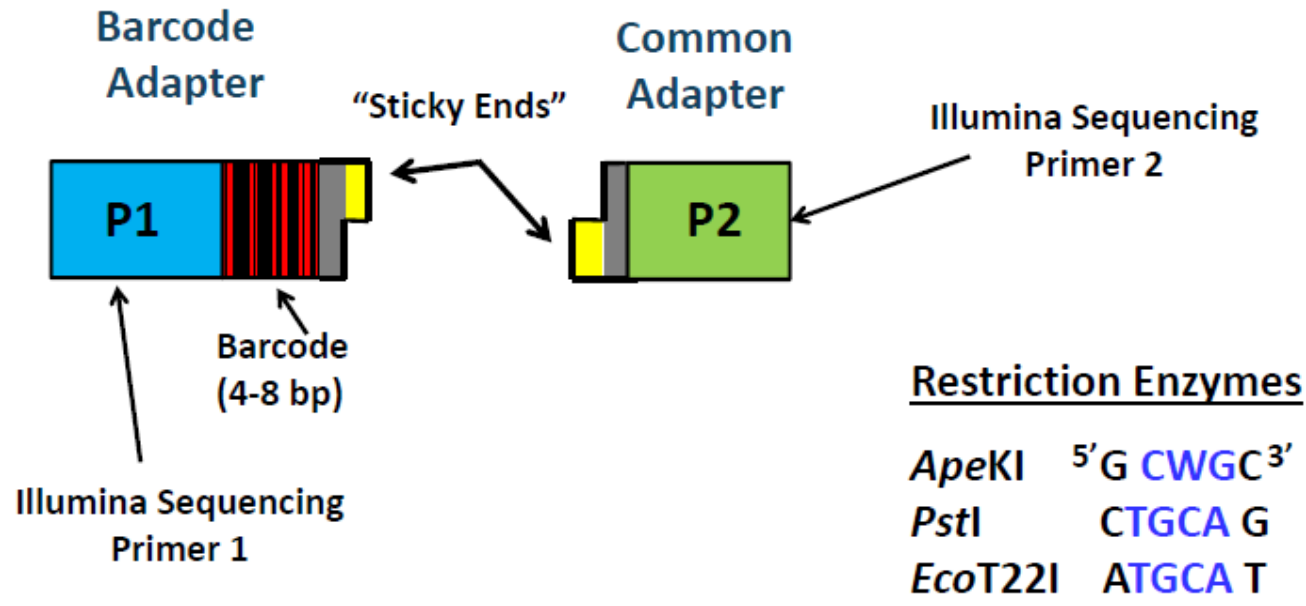
# GBS 96- or 384-plex Protocol

(<http://www.maizegenetics.net/gbs-overview>)



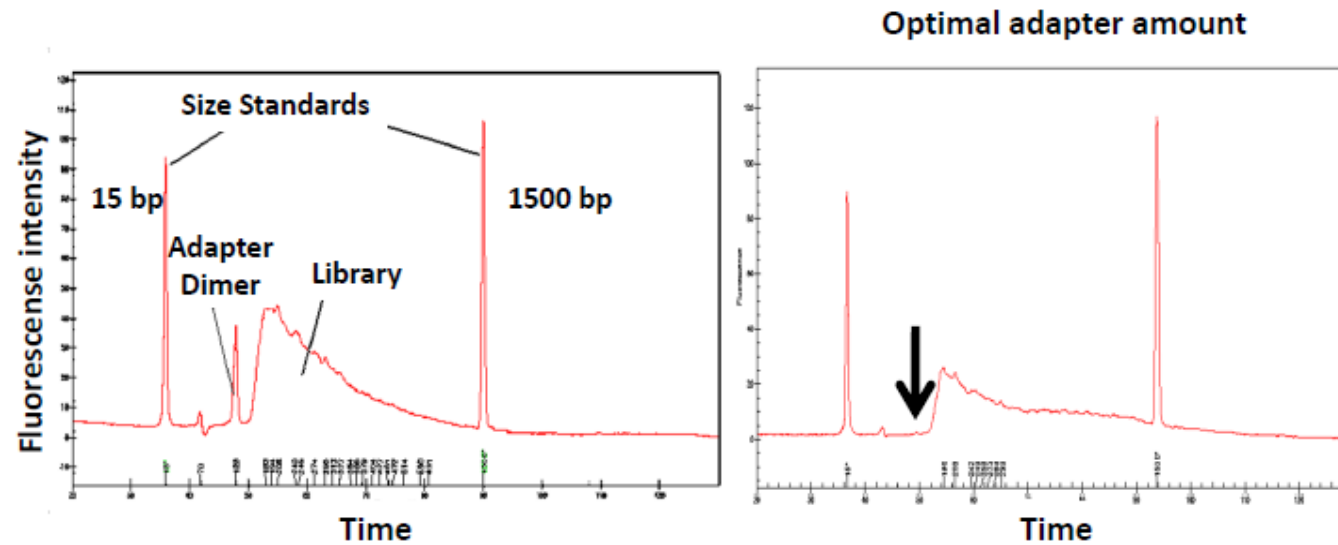


# GBS Adapters and Enzymes

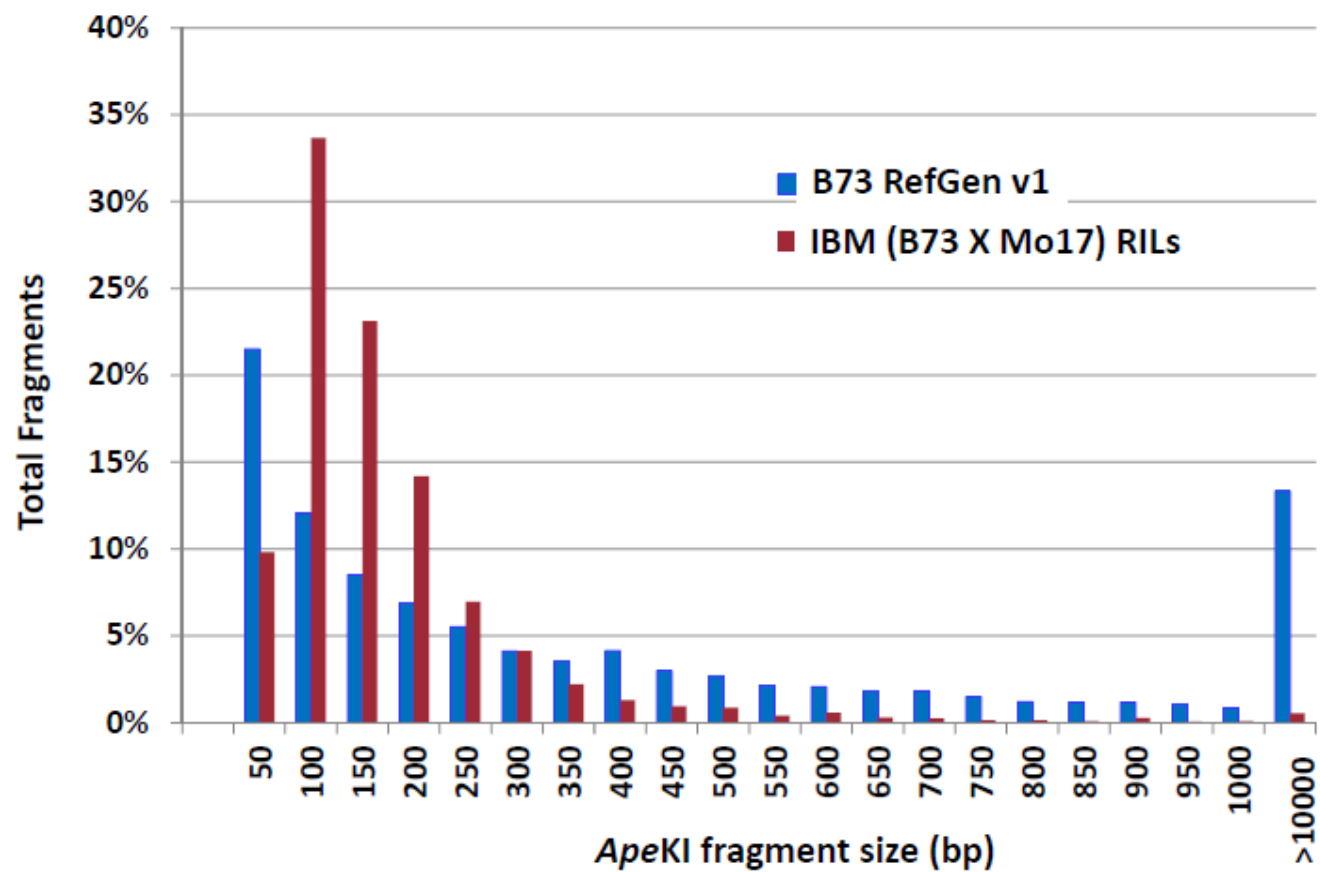


## Perform Titration to Minimize Adapter Dimers Before Sequencing

NOTE: Done once with a small number of samples.  
Adapter dimers constitute only 0.05% of raw sequence reads

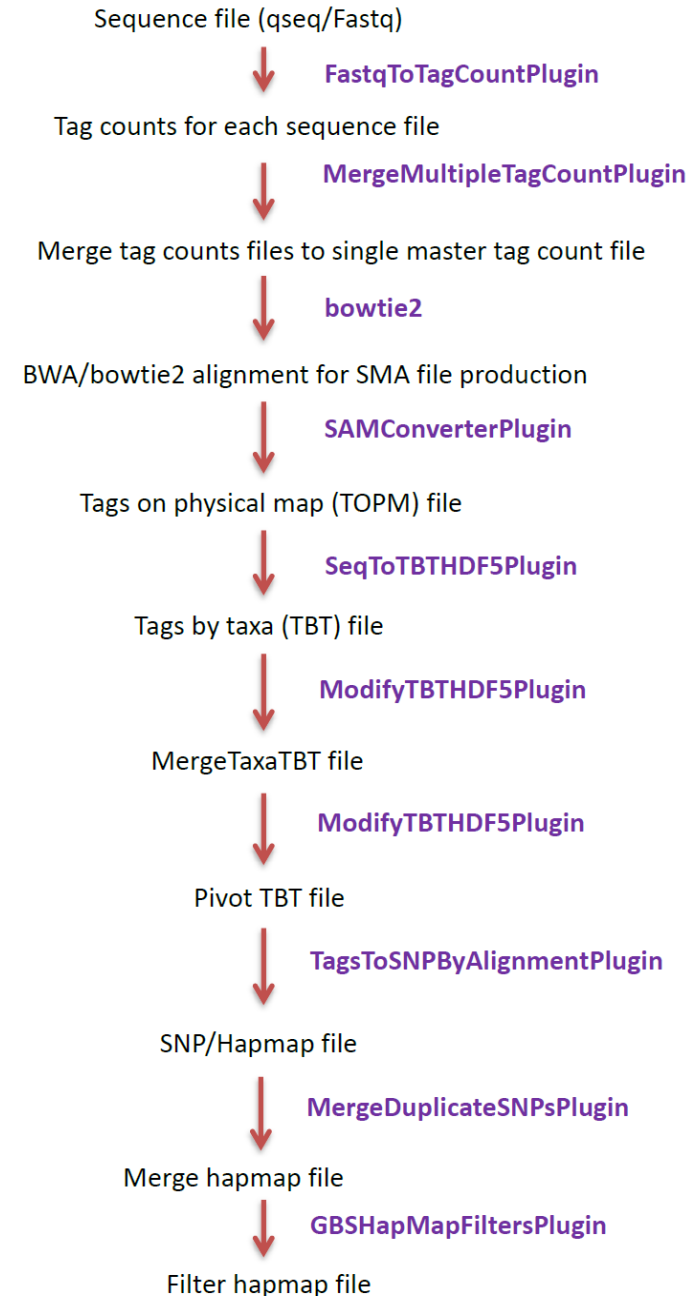


## Small Fragments are Enriched in GBS Libraries



## GBS Work flow:

# Bioinformatics (GBS Pipelines)



# GBS pipelines

TASSEL – tassel4.0standalone

[http://www.maizegenetics.net/index.php?option=com\\_content&task=view&id=89&Itemid=119](http://www.maizegenetics.net/index.php?option=com_content&task=view&id=89&Itemid=119)

User manual: Reference genome

<http://www.maizegenetics.net/tassel/docs/TasselPipelineGBS.pdf> (Elshire et al. 2011)

User manual: without reference genome

[http://www.maizegenetics.net/images/stories/bioinformatics/TASSEL/uneak\\_pipeline\\_documentation.pdf](http://www.maizegenetics.net/images/stories/bioinformatics/TASSEL/uneak_pipeline_documentation.pdf) (Lu et al. 2013)

For more information on GBS

<http://www.maizegenetics.net/gbs-bioinformatics>

**Computationally intensive**

**JAVA environment**

**PERL + JAVA**

# GBS Vocabulary

Taxa	Individual sample
Key file	Text file containing <ul style="list-style-type: none"><li>- Sample information</li><li>- Barcode</li><li>- Flowcell and lane number</li><li>- Sample prepID</li></ul>
Barcode	Unique DNA sequence associated with each taxa
Sequence file	Text file containing DNA sequences information from Illumina <ul style="list-style-type: none"><li>- Qseq or Fastq file</li></ul>
Read	DNA Sequence produced in sequencing..
GBS Tag	DNA sequence starts with cut site remnant and having additional sequence without
barcode	
TagsByTaxa (TBT)	Matrix of GBS tags (row) with taxa (columns)

# GBS pipelines

## Types

### Discovery pipeline

- ❖ Requires a reference genome
- ❖ Multiple steps to process the data into genotypes

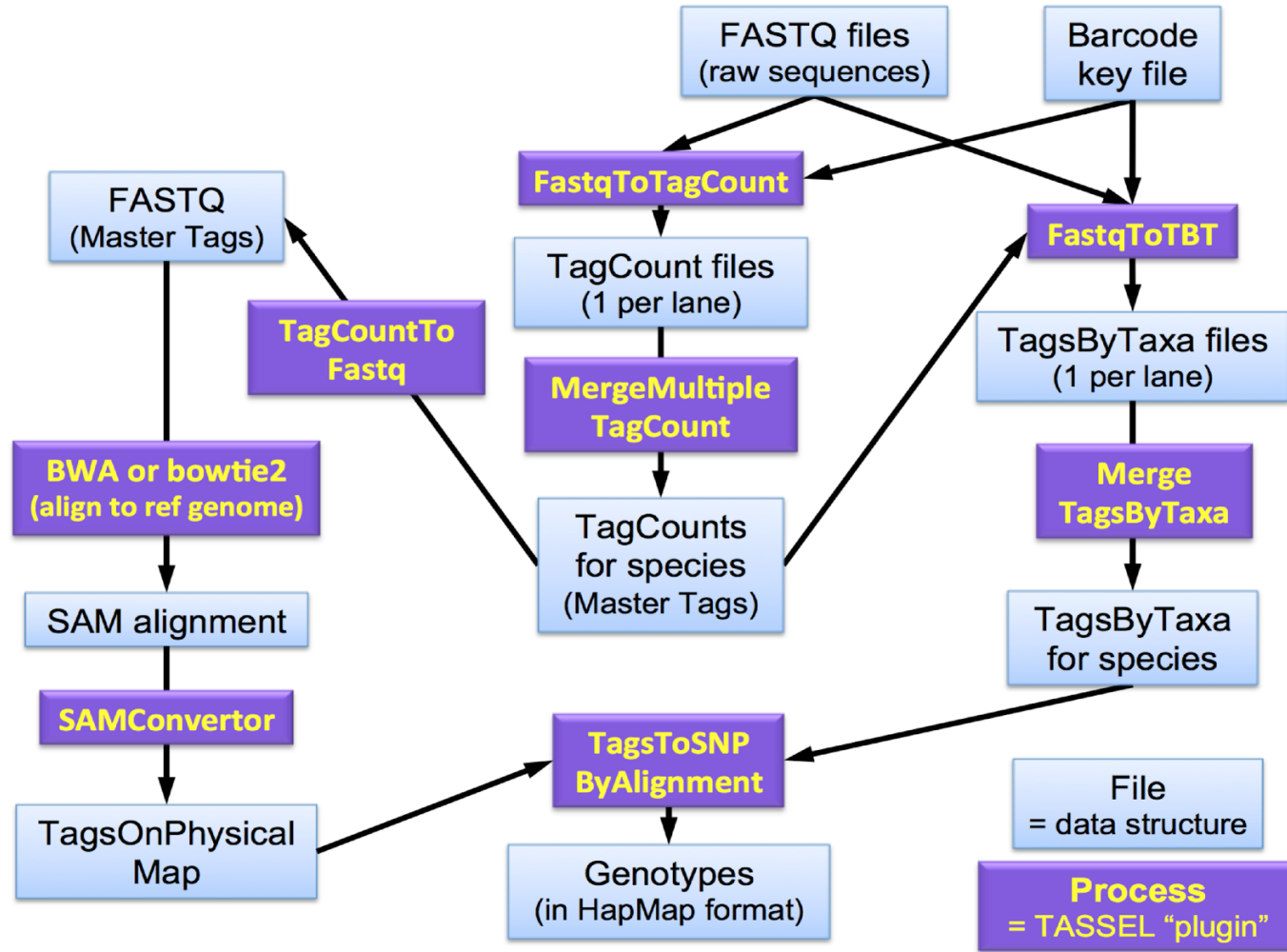
### Production pipeline

- ❖ Uses info from discovery pipeline
- ❖ Only one step from sequence to genotypes

### UNEAK pipeline

- ❖ Species without reference genome

# Discovery GBS Pipeline





# Discovery pipeline

## Key file format

Flowcell	Lane	Barcode	DNA Sample	LibraryPlate	Row	Col	Library PrepID	Library PlateID	Enzyme	FullSampleName
628NVAAXX	5	TGCA	IS2787	ICRISAT Plate 1A	A	2	250006891	450013296	ApeKI	IS2787:628NVAAXX:5:250006891
628NVAAXX	5	ACTA	IS3957	ICRISAT Plate 1A	A	3	250006893	450013296	ApeKI	IS3957:628NVAAXX:5:250006893
628NVAAXX	5	GTCT	IS6193	ICRISAT Plate 1A	A	4	250006895	450013296	ApeKI	IS6193:628NVAAXX:5:250006895
628NVAAXX	5	GAAT	IS9303	ICRISAT Plate 1A	A	5	250006897	450013296	ApeKI	IS9303:628NVAAXX:5:250006897
628NVAAXX	5	GCGT	IS11119	ICRISAT Plate 1A	A	6	250006899	450013296	ApeKI	IS11119:628NVAAXX:5:250006899
628NVAAXX	5	TGGC	IS13845	ICRISAT Plate 1A	A	7	250006901	450013296	ApeKI	IS13845:628NVAAXX:5:250006901
628NVAAXX	5	CGAT	IS15752	ICRISAT Plate 1A	A	8	250006903	450013296	ApeKI	IS15752:628NVAAXX:5:250006903
628NVAAXX	5	CTTGA	IS19132	ICRISAT Plate 1A	A	9	250006905	450013296	ApeKI	IS19132:628NVAAXX:5:250006905
628NVAAXX	5	TCACC	IS20351	ICRISAT Plate 1A	A	10	250006907	450013296	ApeKI	IS20351:628NVAAXX:5:250006907
628NVAAXX	5	CTAGC	IS22330	ICRISAT Plate 1A	A	11	250006909	450013296	ApeKI	IS22330:628NVAAXX:5:250006909
628NVAAXX	5	ACAAA	IS23669	ICRISAT Plate 1A	A	12	250006911	450013296	ApeKI	IS23669:628NVAAXX:5:250006911
628NVAAXX	5	TTCTC	IS303	ICRISAT Plate 1A	B	1	250006913	450013296	ApeKI	IS303:628NVAAXX:5:250006913
628NVAAXX	5	AGCCC	IS2807	ICRISAT Plate 1A	B	2	250006915	450013296	ApeKI	IS2807:628NVAAXX:5:250006915
628NVAAXX	5	GTATT	IS3971	ICRISAT Plate 1A	B	3	250006917	450013296	ApeKI	IS3971:628NVAAXX:5:250006917
628NVAAXX	5	CTGTA	IS6351	ICRISAT Plate 1A	B	4	250006919	450013296	ApeKI	IS6351:628NVAAXX:5:250006919
628NVAAXX	5	AGCAT	IS9468	ICRISAT Plate 1A	B	5	250006921	450013296	ApeKI	IS9468:628NVAAXX:5:250006921
628NVAAXX	5	ACTAT	IS12169	ICRISAT Plate 1A	B	6	250006923	450013296	ApeKI	IS12169:628NVAAXX:5:250006923
628NVAAXX	5	GAGAAT	IS13926	ICRISAT Plate 1A	B	7	250006925	450013296	ApeKI	IS13926:628NVAAXX:5:250006925
628NVAAXX	5	CCAGCT	IS16044	ICRISAT Plate 1A	B	8	250006927	450013296	ApeKI	IS16044:628NVAAXX:5:250006927

# GBS Restriction Fragment Structure



## Accepted read



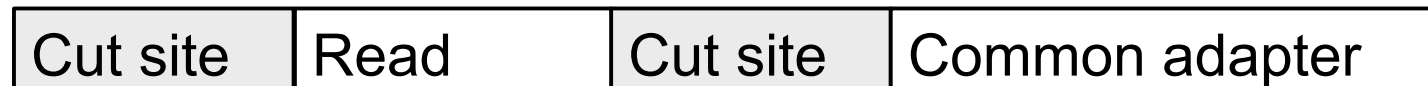
ATGCCTCTGCCTCGTTTAAGTCCGATTGCTGATAGTTATCTTTGTTTCTTACCAAATCCACATGGCTG

## Rejected or Trimmed reads

Potential chimeric sequence



Short sequence



Adapter dimer



# Raw sequence data

Read

ATGCCTCTGCCCTCGTTTAAGTCCGATTGCTGATAGTTATCTTTGTTTCTTACCAAATCCACATGGCTGAGATCGGAAGAGCGGT



GBS pipelines



ATGCCTCTGCCCTCGTTTAAGTCCGATTGCTGATAGTTATCTTTGTTTCTTACCAAATCCACATGGCTGAGATCGGAAGAGCGGT  
CGATCAGCAGTTGACTGGACATCTAGGGGCAAAGCACTGTTTCGGTGCGGGCTGAGATCGGAAGAGCGGTTCCAGCAGGAATGCCGA  
CCGGATATCAGCAGGCCGTGGTACATGTAATGGAGCATGGATTGAAGGTGGATGCCTTCATGTCTGGACGGCGATCGAGCTCGC



TagsCounts

CTGCCCTCGTTTAAGTCCGATTGCTGATAGTTATCTTTGTTTCTTACCAAATCCACATGGCTG  
CAGCAGTTGACTGGACATCTAGGGGCAAAGCACTGTTTCGGTGCGGGCTGAAAAAAAAA  
CAGCAGGCCGTGGTACATGTAATGGAGCATGGATTGAAGGTGGATGCCTGCAAAAAAAAA

Tag count > 5/10

MergeMultipleTagCountPlugin



Merge Tags counts

# MasterTagCounts

Total # Tags			Tags size
1453006	2		
	CTGCCCTCGTTTAAGTCCGATTGCTGATAG...GTTTCTTACCAAATCCACATGGCTG	64	409
	CAGCAAAAAAAAAAAGGACATGGGTATCCGGTAACTGCAGAAAATTTGAGTAAAAGGCGGTA	64	314
	CAGCAAAAAAAAAAATCTCAGAAAATGCAGTGCAGAGTGATATTTTGCTTGATCCTGGCTCAC	64	235
	CAGCAAAAAAAAAAACAAAAACAAAAACAAAAGGATTTCTTTTAGGGAAAAAACATAAGGTGCG	64	338
	CAGCAAAAAAAAAACAAATTGTCAATTATTACAATCCCAAAGGCGAATGACGAACTAAAATCT	64	
	CAGCAAAAAAAAAACAAGACAGGAAGAACGGTGGTGGAGTTGAGTAGAGAGAGGCCAGCAAAGC	64	
	CAGCAAAAAAAAAACCCTCTGCTCTGAAGGTGGAAGGTAGGATAATGCACCTGGATATCAGAAA	64	281
	CAGCTTACCCTTGTTTGGTTGCCCGCATAACGCTTTATGCAAGTTTGCACAGTGCTTAAAGCAG	64	171
	CTGCAACAACAAGCTCGTCGGCGCCAAGTTCTTCGGCCTGGGGTACGAGGCCGCGCACGGCGGG	64	290
	CTGCAACAACGTGTCCTTCGCGCCATATGGCGGGAAGTGGCGCCGGGGCAAGAAGATCGCGGTG	64	382
	CTGCAACAAGACAGAATTTGTTAAAATGTAATATGAAACAGACACATCAAGTTAAATTCTGGTA	64	239
	CTGCAACAAGCAGAAGGTCCGGCGGGGCTGTGGTCGCCGGAGGAGGACGAGAAGCTCATCAAG	64	197
	CTGCAACAAGGAAATTGCAGGCGCGCGGATAACAATGTACCAATACAACGAAGATGCTCTCTGC	64	125
	CTGCAACAAGTTTTTTGCGTGGGCAAGACTTCCATGCTTCTAGAGCCTGATCTCTCAAATCTC	64	141

Tags

Reads

# MergeTagsByTaxa

## Tags

```
CTGCCCTCGTTTAAGTCCGATTGCTGATAGTTATCTTTGTTTCTTACCAAATCCACATGGCTG  
CAGCAGTTGACTGGACATCTAGGGGCAAAGCACTGTTTCGGTGCGGGCTGAGATCGGAAGAGCG  
CAGCAGGCCGTGGTACATGTAATGGAGCATGGATTGAAGGTGGATGCCTTCATGTCCTGGACG
```



```
ATGCCTCTGCCCTCGTTTAAGTCCGATTGCTGATAGTTATCTTTGTTTCTTACCAAATCCACATGGCTGAGATCGGAAGAGCGGT  
GCGTCTGCCCTCGTTTAAGTCCGATTGCTGATAGTTATCTTTGTTTCTTACCAAATCCACATGGCTGAGATCGGAAGAGCGGT  
CGATCAGCAGTTGACTGGACATCTAGGGGCAAAGCACTGTTTCGGTGCGGGCTGAGATCGGAAGAGCGGTTTCAGCAGGAATGCCGA  
CCGGATATCAGCAGGCCGTGGTACATGTAATGGAGCATGGATTGAAGGTGGATGCCTTCATGTCCTGGACGCGGATCGAGCTCGC
```

SeqToTBTHDF5Plugin

Qseq file

Tags file

key file

**TagsByTaxa (TBT)**

(Presence/Absence matrix)

ATGCCT	BTx623
CGAT	N13
CCGGATAT	E36-1
CAGA	ICSL1118
AACT	ICSL1126
GCGT	ICSL1134

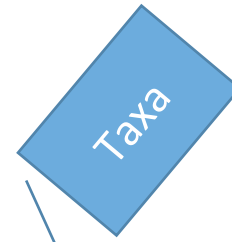
MergeTBTFiles

Merge TBT files

# TBT file

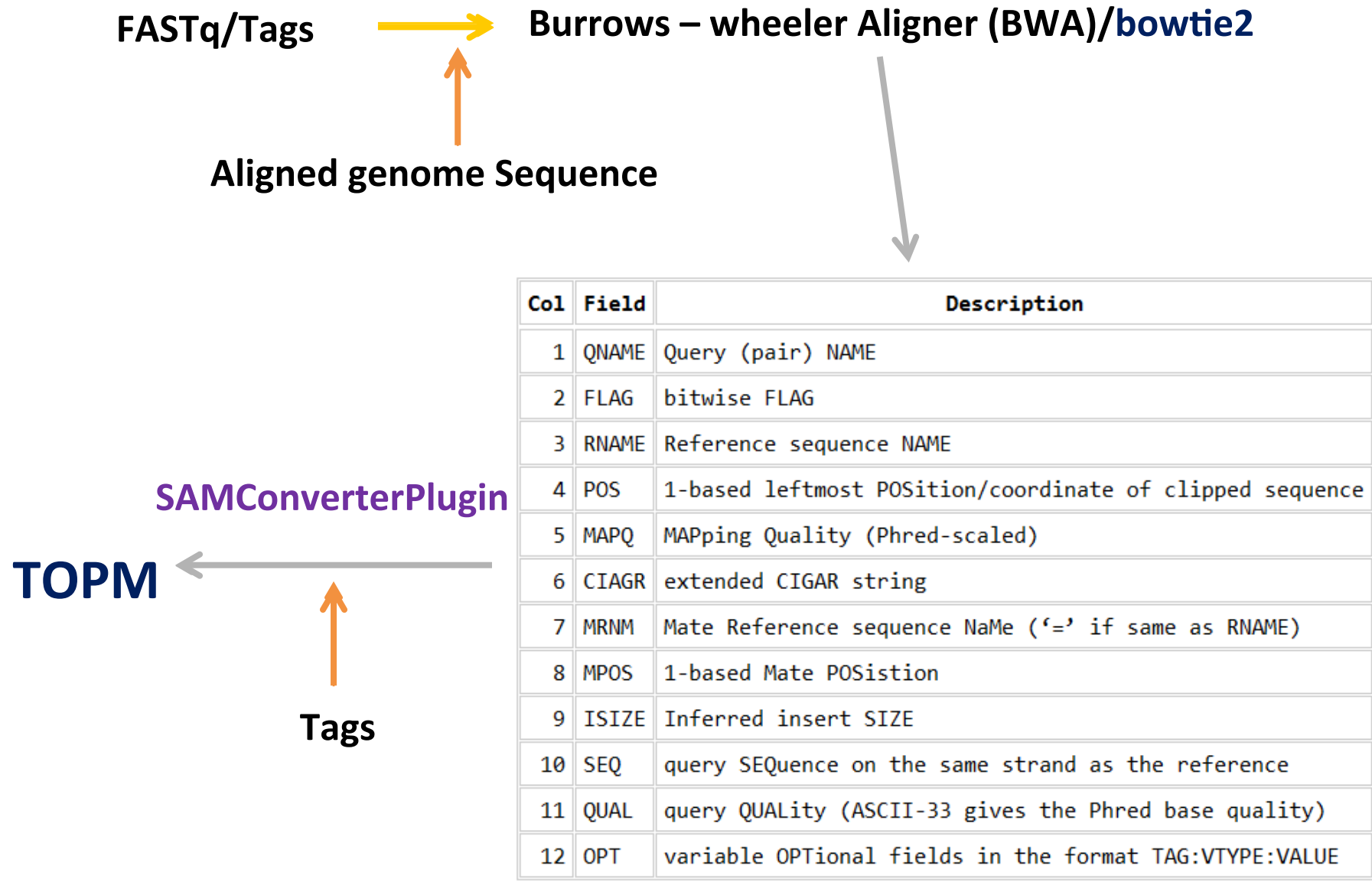
Tags

CTGCCCTCGTTTAAGTCCGATTGCTGATAGTTATCTTTGTTTCTTACCAAAATCCACATGGCTG  
CAGCAAAAAAAAAAAGGACATGGGTATCCGGTAACACTGCAGAAAATTTGAGTAAAAGGCGGTA  
CAGCAAAAAAAAAAATCTCAGAAAATGCAGTGCAGAGTGTATTTTTGCTTGATCCTGGCTCAC  
CAGCAAAAAAAAAACAAAAACAAAAACAAAAGGATTTCTTTAGGGAAAAAACATAAGGTGCG  
CAGCAAAAAAAAAACAAATTGTCAATTATTACAATCCCAAAGGCGAATGACGAAACTAAAATCT  
CAGCAAAAAAAAAACAAGACAGGAAGAACGGTGGTGGAGTTGAGTAGAGAGAGGCCAGCAAAGC  
CAGCAAAAAAAAAACCCTCTGCTCTGAAGGTGGAAGGTAGGATAATGCACCTGGATATCAGAAA  
CAGCTTACCCTTGTGGTTGCCGCATACGCCTTTATGCAAGTTTGCACAGTGCTTAAAGCAG  
CTGCAACAACAAGCTCGTCGGCGCCAAGTTCTTCGGCCTGGGGTACGAGGCCGCGCACGGCGGG  
CTGCAACAACGTGTCCTTCGCGCCATATGGCGGGAAGTGGCGCCGGGGCAAGAAGATCGCGGTG  
CTGCAACAAGACAGAATTTGTTAAATGTAATATGAAACAGACACATCAAGTTAAATTCTGGTA  
CTGCAACAAGCAGAAGGTCCGGCGGGGCTGTGGTCCGGGAGGAGGACGAGAAGCTCATCAAG  
CTGCAACAAGGAAATTGCAGGCGCGCGGATAACAATGTACCAATACAACGAAGATGCTCTCTGC  
CTGCAACAAGTTTTTGCCTGGGCAAGACTTCCATGCTTCTAGAGCCTGATCTCTCAAAATCTC



	BTx623	N13	E36-1	ICSL1118	ICSL1126	ICSL1134
64	1	0	0	0	0	1
64	0	0	0	0	0	0
64	0	0	1	0	0	0
64	1	1	0	1	0	0
64	0	0	0	0	0	0
64	1	0	0	0	0	1
64	0	0	0	0	0	0
64	0	0	0	1	0	0
64	0	0	0	0	0	1
64	0	0	0	0	0	0
64	0	0	0	0	0	0
64	1	1	1	0	0	1
64	0	0	0	0	0	0
64	0	0	0	0	0	0

# TagsOnPhysicalMap (TOPM)



# TOPM file

Tag	length	chr	start	end
CACCCCCGCAAGTGTTTCTGTGTTGAGGAGTCAACAAGAACCAATTGCAGGATGTAAAAGCAG	64	1	12982	13046
CTGCTGGTTAATACGTACTCCAGTTGGTTACCATGCACCAAGCAGGGCAAGGCACTCCTGCTCC	64	1	2132	2196
CTGCTGGTTAATACGTACTCCAGTTGGTTACTATGCACCAAGCAGGGCAAGGCACTCCTGCTCC	64	1	2132	2196
CAGCCCTACGGCCGCAACTGCAAAAATCAGATGTGAGCATAGTTCTCTAATGTATCATAAAAAGT	64	2	2266316	2266380
CTGCTACCTCTGCCCTGCGTGCGTGACCCATACATACGGTGGGGCCGGGGTGTAATAATCTTGTG	64	2	2926106	2926170
CAGCGGCTGAACACTGACCAGTTGACGGCGTCGCGCGCGGGCAGGGCGACGATGTCGGCGCAG	64	3	2129422	2129486
CTGCGGCGGGCTCCTACTCCCTCAACGGAAGAGGCGTTACGACGGGTGGAGGGCAGTTTCTCT	64	3	2140833	2140897
CTGCGGCGGGCTCCTACTCCCTCAACGGAAGAGGCGTTACGACGGGGGGAGGGCAGTTTCTCT	64	3	2140833	2140897
CTGCGGCGGGCTCCTACTCCCTCAACGGAAGAGGCGTTACGACGGGGGGAGGGCAGTTTCTCT	64	3	2140833	2140897
CTGCGGCGGGCTCCTACTCCCTCAACGGAAGAGGCGTTACGACGGGGGGAGTGGCAGTTTCTCT	64	3	2140833	2140897
CTGCGGCGGGCTCCTACTCCCTCAACGGAAGAGGCGTTACGACGGGGGGGGGGCAGTTTCTCT	64	3	2140833	2140897
CTGCGGCGGGCTCCTACTCCCTCAACGGAAGAGGCGTTACGACGGGGGGGGTGGCAGTTTCTCT	64	3	2140833	2140897
CTGCGGCGGGCTCCTACTCCCTCAACGGAAGAGGCGTTACGACGGGTGGAGGGCAGTTTCTCT	64	3	2140833	2140897
CTGCGGCGGGCTCCTACTCCCTCAACGGAAGAGGCGTTACGACGGGTGGAGTGGCAGTTTCTCT	64	3	2140833	2140897
CAGCTGTCCCTTGTAACCTGCGGACATGACCGAACATCCAAGTACTCGAGAGCCTTGAAGGCTT	64	4	56916960	56917024
CTGCTTGTGTTGTTGATTGCTCGTTGGCCGATCGGTGGTCGTTGTCGTCGTTGAGCTGGTAGC	64	4	56942951	56943015
CAGCGCCCAAACCTTCTTTCCCGCTCTTGTTGATTCTCAGTCTCTCAGTTCTTGAACATCG	64	5	54508281	54508345
CTGCCTCAGCACAGCCGGTGCTCTTGGATCCTCGAGCACCTGCACTCTTTCAGGAGGCCAAGGA	64	5	13747204	13747268
CAGCACGCGTATCCTTGGCCCTTATTGTTTGAAAATGACCTGTGACGGTGCCTTCTGCGTAGT	64	6	52338128	52338192
CTGCGGCGGGCATGACCGACGGGCGGCGTGGGTGCCGAGTGCCGGCAACTTCCGCCGTTCT	64	6	58276481	58276545
CTGCCACAGACGGAGTTACGACAACCTTGCAAACTCCGCCGGCTTAAGTTAATTTAATCAGG	64	7	18492751	18492815
CAGCCAAAAAATGGCAAAAGGTGTAGGACTTGGGTACGTCAAAATTCATTCATGTGAGAATGA	64	7	19393217	19393281



# Hapmap file

Key file  
Qseq file  
Tags file  
TOPM file  
TBT file

TagsToSNPByAlignmentPlugin

Hapmap file  
(Sequence alignment file)

```
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGGGGAGTGGCAGTTT
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGGGGGGTGGCAGGTT
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGGGGGGGGCAGGTT
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGGGGAGTGGCAGGTT
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGTGGAGTGGCAGTTT
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGTGGAGTGGCAGGTT
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGGGGAGGGGCAGGTT
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGTGGAGGGGCAGTTT
CTGCGGC GGGCTCC TACTCCC TCAACGGAAGAGGCGTTACGACGGGGGGAGGGGCAGTTT
*****
```



# Why production pipeline?

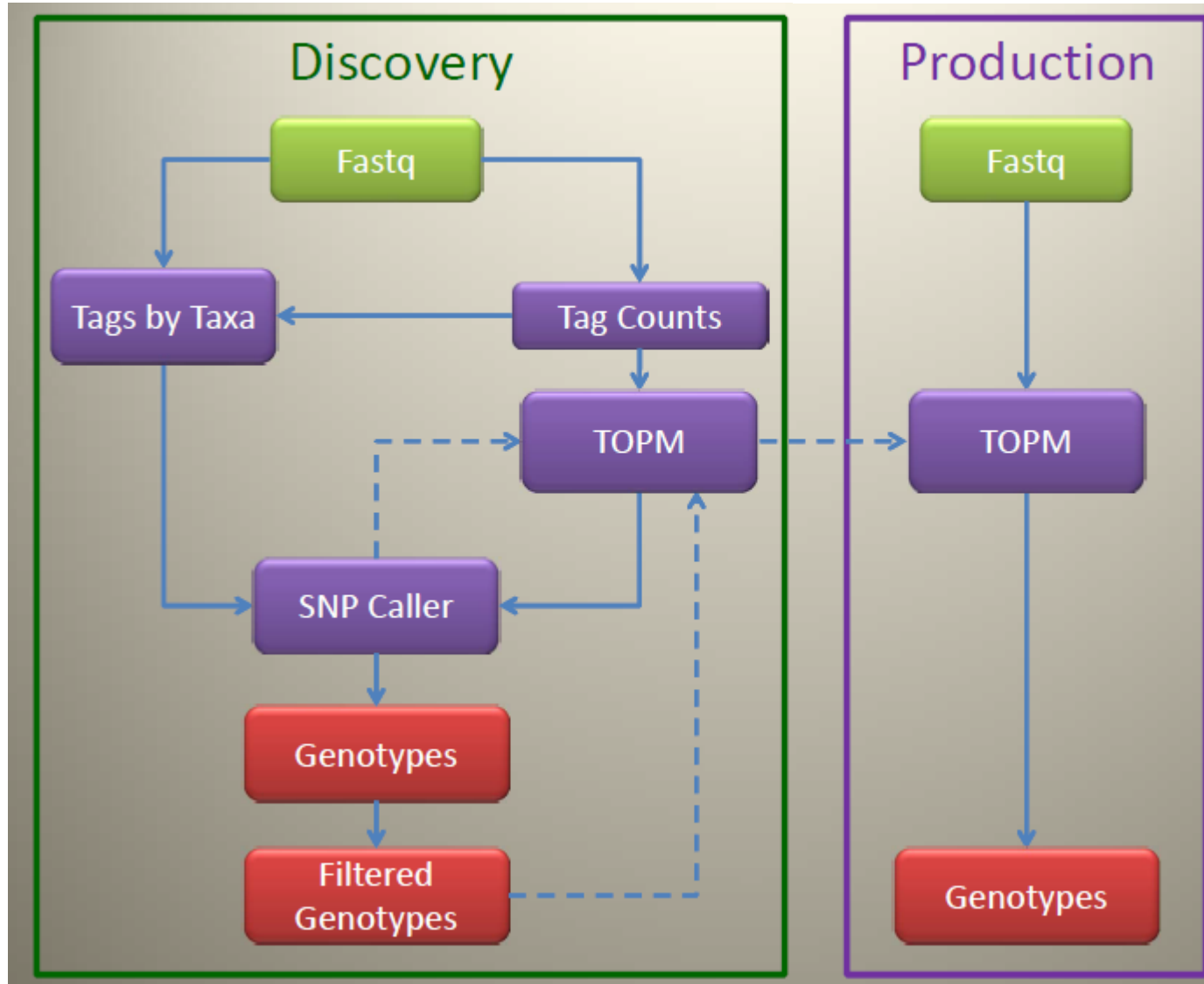
- ❖ **Discovery pipeline takes more time**

  - Eg: Last maize build took > 3 months**

- ❖ **Most common alleles have been identified in discovery pipelines**

- ❖ **Use of available information from discovery pipelines to call SNPs  
in new run**

# GBS Production pipeline



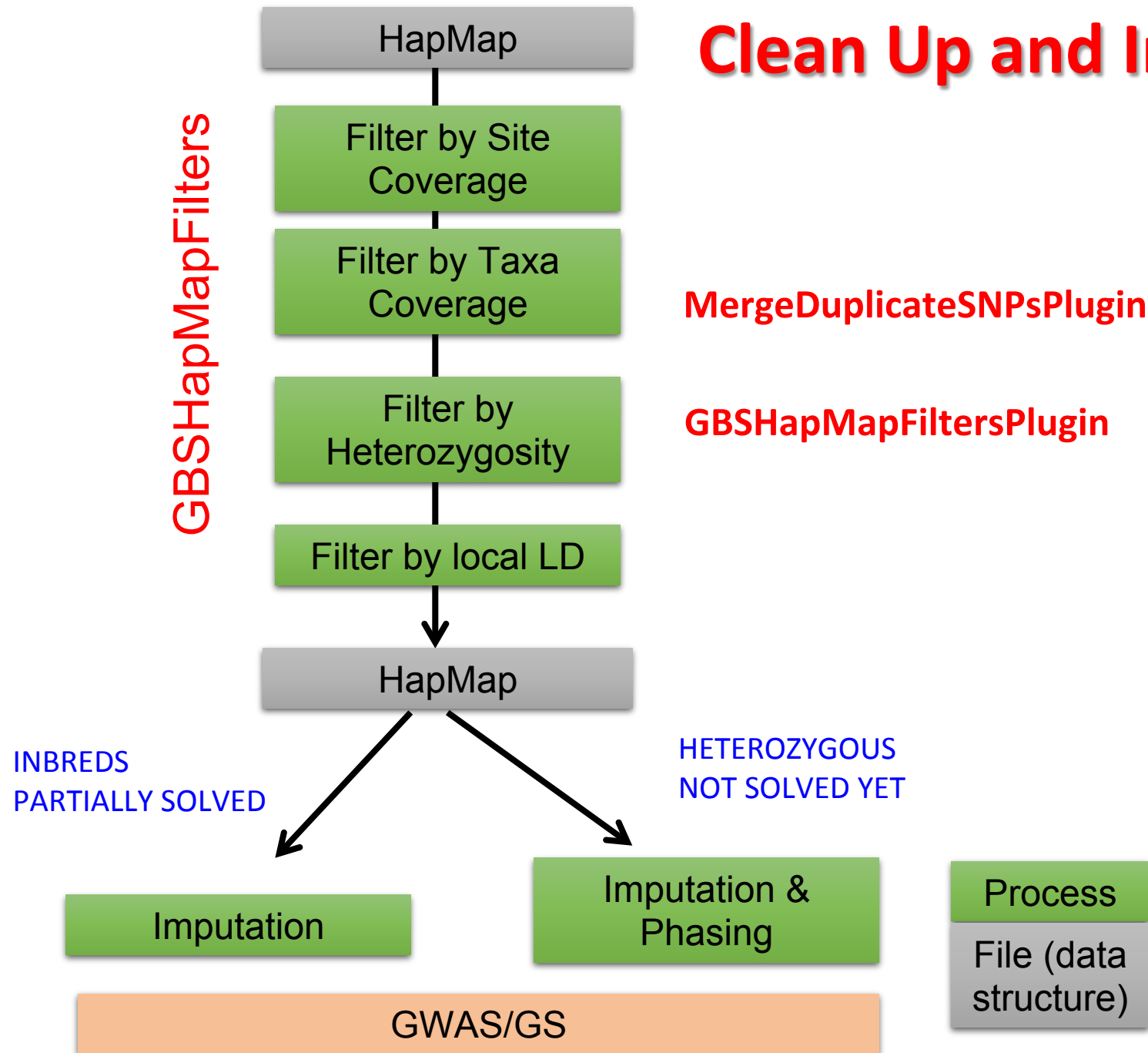
# UNEAK pipeline

(Universal Network Enabled Analysis Kit)  
A reference free SNP calling pipeline

**Designed for species that.....**

- lack a reference genome
- diploid or polyploid
- inbreeders or outcrosses
- having limited genetic or genomic resources











# Clean Up and Imputation



# Hapmap file

rs#	alleles	chrom	pos	ICSB617	ICSB212	ICSB536	ICSB79	ICSB538	ICSB639	ICSB417	ICSB748	ICSB563
S9_20491	G/A	9	20491	G	N	N	N	N	N	N	N	N
S9_23432	C/Y	9	23432	N	N	Y	C	N	N	C	C	N
S9_23452	T/K	9	23452	N	N	T	T	N	N	T	T	N
S9_26213	G/A	9	26213	N	N	N	W	N	N	N	G	N
S9_27701	T/C	9	27701	T	N	T	N	N	Y	C	T	N
S9_27720	T/K	9	27720	T	N	T	N	N	T	T	T	N
S9_30044	C/T	9	30044	C	N	N	N	N	N	N	N	N
S9_36986	A/G	9	36986	N	N	A	N	N	N	N	N	N
S9_39902	C/T	9	39902	N	N	N	N	T	N	T	N	T
S9_39912	A/G	9	39912	N	N	N	N	G	N	G	N	G
S9_40450	A/C	9	40450	N	K	N	N	N	N	N	N	N
S9_40452	G/T	9	40452	N	N	N	N	N	N	N	N	N
S9_40459	A/G	9	40459	N	N	N	N	N	N	N	N	N
S9_43958	T/C	9	43958	N	N	N	Y	N	M	N	N	N
S9_43982	C/T	9	43982	N	N	N	N	N	N	N	N	N
S9_44191	G/C	9	44191	N	M	N	N	N	N	N	N	N
S9_45394	T/A	9	45394	N	N	N	N	N	N	N	N	N
S9_56733	C/M	9	56733	N	N	N	N	N	N	C	N	N
S9_56734	C/S	9	56734	N	N	N	N	N	N	C	N	N
S9_56735	T/W	9	56735	N	N	N	N	N	N	T	N	N
S9_69001	C/Y	9	69001	C	N	N	N	N	C	N	C	N
S9_72163	C/T	9	72163	N	N	N	N	N	N	T	N	N
S9_74270	C/T	9	74270	C	N	T	N	N	N	C	N	N
S9_76335	A/G	9	76335	N	A	N	N	N	N	R	A	N
S9_76341	T/G	9	76341	N	T	N	N	N	N	K	T	N
S9_78501	A/C	9	78501	N	A	N	N	N	N	N	N	N

# Hapmap Files

 sorg_Jan2013_Filter2_c10.hmp.txt	1001,347,276	Text Doc...
 sorg_Jan2013_Filter2_c9.hmp.txt	972,503,003	Text Doc...
 sorg_Jan2013_Filter2_c8.hmp.txt	667,513,134	Text Doc...
 sorg_Jan2013_Filter2_c7.hmp.txt	616,327,520	Text Doc...
 sorg_Jan2013_Filter2_c6.hmp.txt	1048,787,660	Text Doc...
 sorg_Jan2013_Filter2_c5.hmp.txt	1047,007,421	Text Doc...
 sorg_Jan2013_Filter2_c4.hmp.txt	1284,135,657	Text Doc...
 sorg_Jan2013_Filter2_c3.hmp.txt	1505,358,845	Text Doc...
 sorg_Jan2013_Filter2_c2.hmp.txt	1493,688,785	Text Doc...
 sorg_Jan2013_Filter2_c1.hmp.txt	1675,855,191	Text Doc...

Chromosome	Unfilter	Filter2 (mnMAF 0.01, mnTCov 0.01, mnSCov 0.01, mnF 0.8)
1	1288506	83460
2	999490	74192
3	1077534	74976
4	868854	64199
5	528132	52401
6	706426	51633
7	569982	28153
8	462068	30221
9	618291	48143
10	659650	49595
<b>Total</b>	<b>7778933</b>	<b>556973</b>



# Optimizing GBS in New Species





# A Robust, Simple Genotyping-by-Sequencing (GBS) Approach for High Diversity Species

Robert J. Elshire<sup>1</sup>, Jeffrey C. Glaubitz<sup>1</sup>, Qi Sun<sup>2</sup>, Jesse A. Poland<sup>3</sup>, Ken Kawamoto<sup>1</sup>, Edward S. Buckler<sup>1,4</sup>, Sharon E. Mitchell<sup>1\*</sup>

<sup>1</sup> Institute for Genomic Diversity, Cornell University, Ithaca, New York, United States of America, <sup>2</sup> Computational Biology Service Unit, Cornell University, Ithaca, New York, United States of America, <sup>3</sup> Hard Winter Wheat Genetics Research Unit, United States Department of Agriculture/Agricultural Research Service, Manhattan, Kansas, United States of America, <sup>4</sup> Plant, Soil and Nutrition Research Unit, United States Department of Agriculture/Agricultural Research Service, Ithaca, New York, United States of America

# Switchgrass Genomic Diversity, Ploidy, and Evolution: Novel Insights from a Network-Based SNP Discovery Protocol

Fei Lu<sup>1</sup>, Alexander E. Lipka<sup>1,2</sup>, Jeff Glaubitz<sup>1</sup>, Rob Elshire<sup>1</sup>, Jerome H. Cherney<sup>3</sup>, Michael D. Casler<sup>4,5</sup>, Edward S. Buckler<sup>1,2</sup>, Denise E. Costich<sup>1,2\*</sup>

<sup>1</sup> Institute for Genomic Diversity, Cornell University, Ithaca, New York, United States of America, <sup>2</sup> Agricultural Research Service, United States Department of Agriculture, Ithaca, New York, United States of America, <sup>3</sup> Department of Crop and Soil Sciences, Cornell University, Ithaca, New York, United States of America, <sup>4</sup> Agricultural Research Service, United States Department of Agriculture, Madison, Wisconsin, United States of America, <sup>5</sup> Department of Agronomy, University of Wisconsin–Madison, Madison, Wisconsin, United States of America

# Population genomic and genome-wide association studies of agroclimatic traits in sorghum

Geoffrey P. Morris<sup>a,1,2</sup>, Punna Ramu<sup>b,1</sup>, Santosh P. Deshpande<sup>b</sup>, C. Thomas Hash<sup>c</sup>, Trushar Shah<sup>b</sup>, Hari D. Upadhyaya<sup>b</sup>, Oscar Riera-Lizarazu<sup>b</sup>, Patrick J. Brown<sup>d</sup>, Charlotte B. Acharya<sup>e</sup>, Sharon E. Mitchell<sup>e</sup>, James Harriman<sup>e</sup>, Jeffrey C. Glaubitz<sup>e</sup>, Edward S. Buckler<sup>e,f,g</sup>, and Stephen Kresovich<sup>a</sup>

<sup>a</sup>Department of Biological Sciences, University of South Carolina, Columbia, SC 29208; <sup>b</sup>International Crops Research Institute for the Semi-Arid Tropics (ICRISAT), Hyderabad 502 324, Andhra Pradesh, India; <sup>c</sup>ICRISAT-Sadoré, BP 12404 Niamey, Niger; <sup>d</sup>Department of Crop Sciences, University of Illinois, Urbana, IL 61801; <sup>e</sup>Institute for Genomic Diversity and <sup>f</sup>Department of Plant Breeding and Genetics, Cornell University, Ithaca, NY 14853; and <sup>g</sup>Agricultural Research Service, Department of Agriculture, Ithaca, NY 14853

# **Genome Wide Association Study (GWAS)**

**Association Mapping (AM)**

(or)

**Linkage Disequilibrium  
mapping (LD mapping)**

(or)

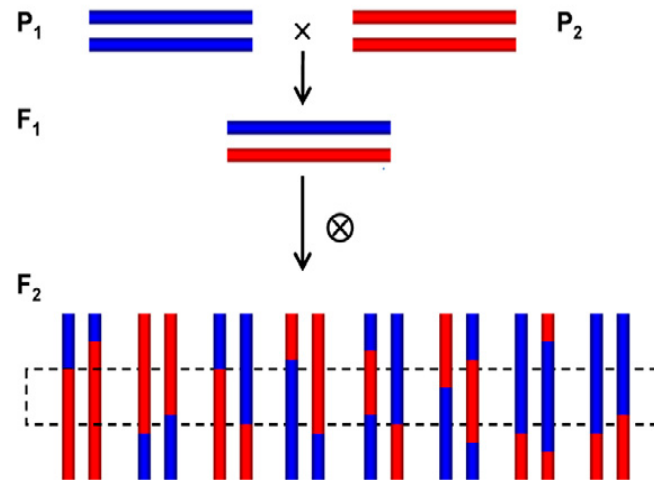
**Population mapping**

# Linkage mapping (or) Family mapping

Generation of mapping population (RILs, NILs, DH, BC, F2)

Genotyping – polymorphic markers

Phenotyping – trait of interest



## Limitations

Resolution power is low (10–30 cM)

Small population size

Modest degree of recombination within the population

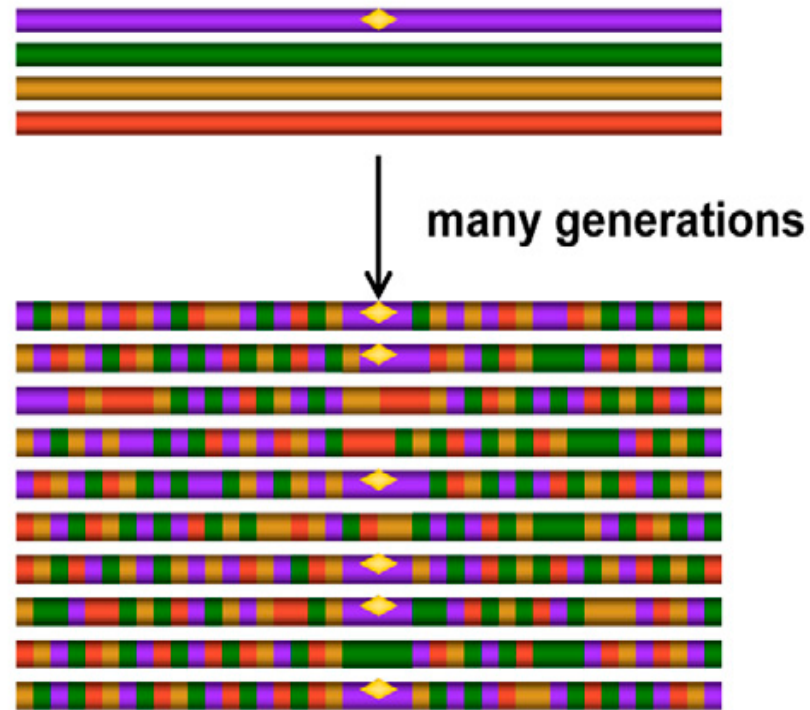
**Linkage mapping – limited to sampling only two alleles at a given locus in any given bi-parental population**

# ASSOCIATION MAPPING

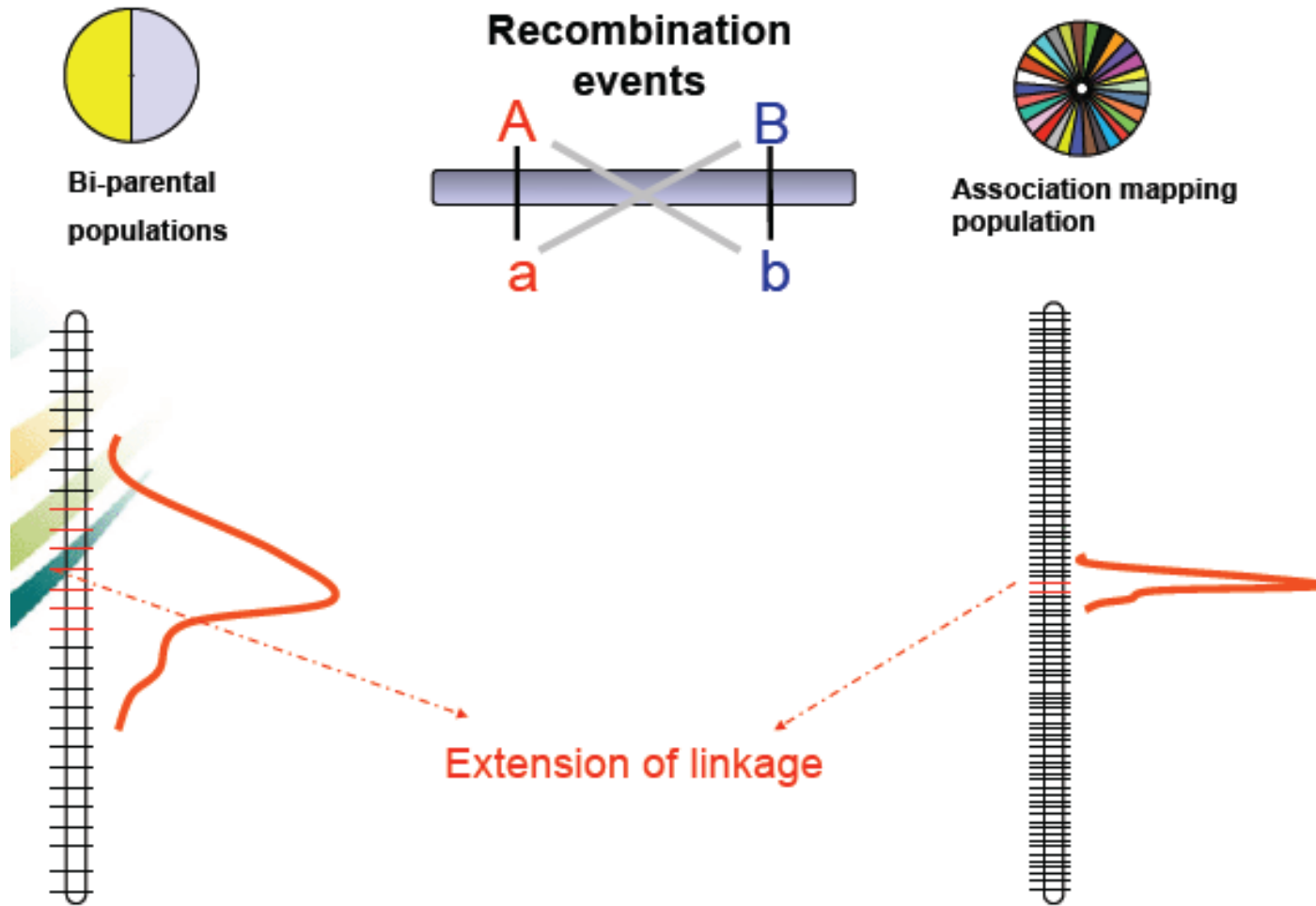
- ❖ Currently existing natural populations are used Vs generating a population via a biparental cross
- ❖ No need to develop mapping population
- ❖ A potentially large number of alleles per locus – as opposed to only two – can be surveyed simultaneously
- ❖ Resolution can be dramatically increased (e.g. 2000 bp in diverse maize inbred lines) - - - Fine mapping.
- ❖ Reduces time
- ❖ Considering recombination of history/evolution

## AM is a multi-disciplinary field

- Genomics
- Genetics
- Molecular Biology
- Statistical Genetics
- Bioinformatics



# Recombination: key of Genetic variation or Success of Breeding



# GWAS Types

Success of either methods depends on population size and degree of LD

## 1. Genome wide scanning or AM

Markers spanned across the genome

Moderate to extensive LD

## 2. Candidate gene scanning or AM

Sequencing only candidate gene

Low LD

# AM in human and plants

## HUMANS

AM is responsible for identification and cloning of

- ❑ Cystic fibrosis gene
- ❑ Diastrophic displasia gene
- ❑ Alzheimer's disease

## PLANTS

First AM – Candidate gene analysis

**Flowering time and dwarf8 (d8) gene in Maize**

First AM – Genome scan is in

**Sea beet (*Beta vulgaris* ssp. *maritima*)**



# GENOME-WIDE ASSOCIATION MAPPING (GWA)

Sps – Self-fertile : Arabidopsis, rice

Clonally propagated : Switch grass, grape

If LD is high, GWA is useful with low resolution mapping

Number of markers to screen determined by

- sample size,
- Extent of LD

E.g.:

Human – 70,000 markers

Arabidopsis – 2,000 markers

Diverse Maize Landraces – 750,000 markers

Elite Maize lines – 50,000 markers

Sorghum – 556,000 markers

# CANDIDATE GENE APPROACH

Multi-disciplinary  
approach

Mutagenesis

Biochemical analysis

Expression profiling

Comparative genome  
mapping

Bioinformatics

Linkage mapping

**Positional  
candidates**

**or**

**Candidate  
genes**

***How to start?***

# Basic resources for AM

- **Germplasm choice**
- **Trait evaluation**
- **Identification of candidate polymorphism**
- **Estimation of population structure**
- **Statistical analysis**

# GERMPLASM CHOICE

- Encompass as much **phenotypic variation** as possible, and perhaps represent the breeding pool of a crop species
- **Genetic or phenotypic surveys** can be used to identify genotypically diverse subsets of the available germplasm in order to maximize the range of alleles sampled in the population
- In some species, **core sets of germplasm** have already been defined and characterized, and can be used to initiate preliminary association studies
- Germplasm should have maximum diversity of gene pool with more extensive recombination events

## Community resources

<b>Maize</b>	300 maize inbred lines (Flint-Garcia et al., 2005) NAM panel consists of 5000 RILs (Yu et al. 2008)
<b>Sorghum</b>	377 Association mapping panel (Casa et al. 2008) 107 Sorghum Diversity Research Set (SDRS) (Shehzad et al. 2009) 384 Reference set of sorghum (Billot and Ramu et al. 2013,)
<b>Barley</b>	3840 Diverse germplasm set

- ✓ Model based study found that more power is achieved by increasing the no of individuals rather than increasing the no of markers

# Linkage Disequilibrium (LD)

- ❑ Non-random association between alleles at different loci
- ❑ LD extends to a much longer distance in self-pollinated crops than in cross-pollinated species
- ❑ Genome-wide LD determines the mapping resolution and marker density for a genome scan

**Plant populations amenable for association studies can be classifiable into one of five groups**

- (i) ideal sample with subtle population structure and familial relatedness
- (ii) multi-family sample
- (iii) sample with population structure
- (iv) **sample with both population structure and familial relationships, and**
- (v) sample with severe population structure and familial relationships.

Due to local adaptation, selection, and breeding history in many plant species, many populations for association mapping would fall into category four.

# Factors affecting LD

## Mutation

– Base for producing polymorphism

## Recombination

– weakens the intra-chromosomal LD

inter-chromosomal LD broken by **independent assortment**

## Population size

– Small pop – effect of genetic drift result in constant loss of rare allele combinations, thus LD increases

## Population mating system – Self vs Out cross

LD decay more – out cross sps – because recombination leads to heterozygosity

## Admixture

– introduction of chromosomes of different ancestry and allele frequency

Gene flow between genetically distant populations

## Selection

– Bottle neck effect very few allele

- combinations are passed to next generations
- Locus specific alleles

# LD Decay

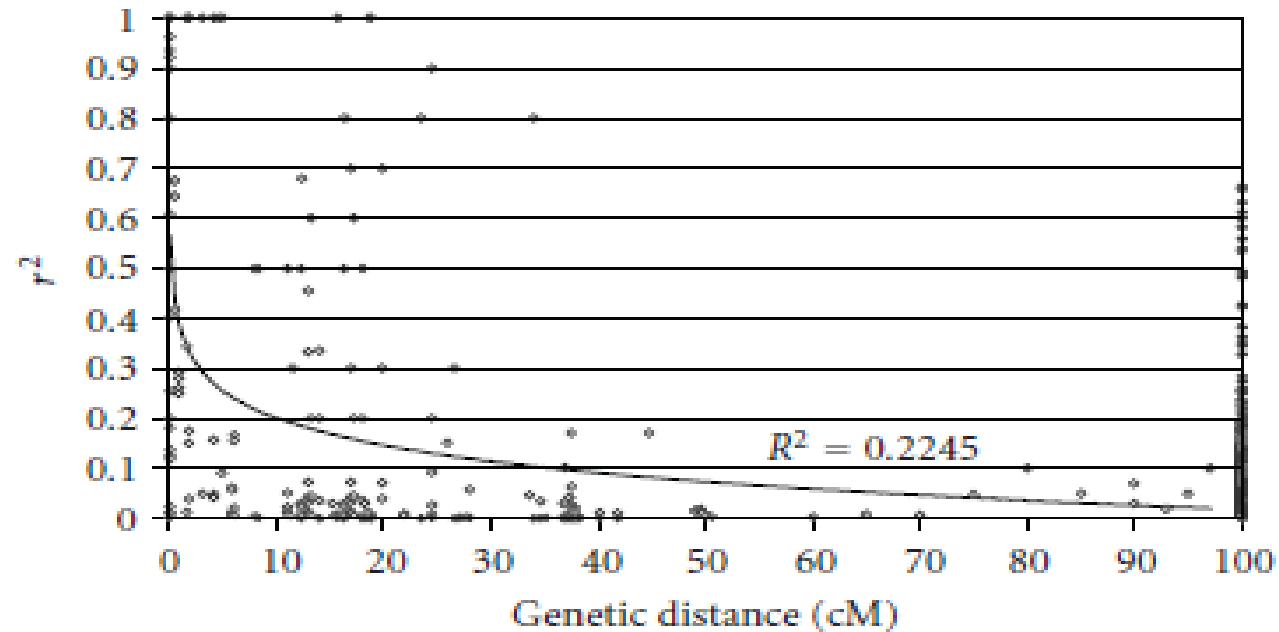


FIGURE 3: Linkage disequilibrium (LD) decay plot depicted from the LD values of a hypothetical marker data to demonstrate a measure of an average genome-wide LD block sizes. A pairwise LD values ( $r^2$ ) are plotted against a genetic distance. Inner fitted trend line is a nonlinear logarithmic regression curve of  $r^2$  on genetic distance. LD decay is considered below  $r^2 = 0.1$  threshold and based on trend line it is around 38–40 cM in above plot. A pairwise LD between unlinked marker loci is assigned to 100 cM distance point. Note: this is for demonstration purposes only and does not have any real impact or correspond to any genomic fragment of an organism.



# 2. Genotyping

## CANDIDATE GENE MARKERS

- ❑ May be from related or unrelated
- ❑ Comparative genomics and Bioinformatics
- ❑ SNPs are preferred markers system in candidate gene approach
- ❑ Promoters, Introns, Exons and 5' / 3' UTRs – targets for identification of candidate gene SNPs
- ❑ Rate of LD decay – denotes the no of SNPs per unit length
- ❑ It is not essential to screen all SNPs. SNPs that cause phenotypic variation due to alter in protein function (Coding SNPs) or gene expression (regulatory SNPs) should be a top priority
- ❑ If there are block of several SNPs in significant LD, there is no need to do all SNPs. An alternatively, select and score a fraction of SNPs (Tag SNPs) that capture most of the haplotype blocks in a candidate gene region – more cost effective

## WHOLE GENOME SCANNING

The extent of LD - dictates the number of markers

SNP marker – large number.....resolution is more

SSR markers – small number .....resolution is low

# 3. PHENOTYPING DATA

**AM** –

a relatively large number of diverse accessions phenotypic data collection with adequate replications across multiple years multiple locations is challenging

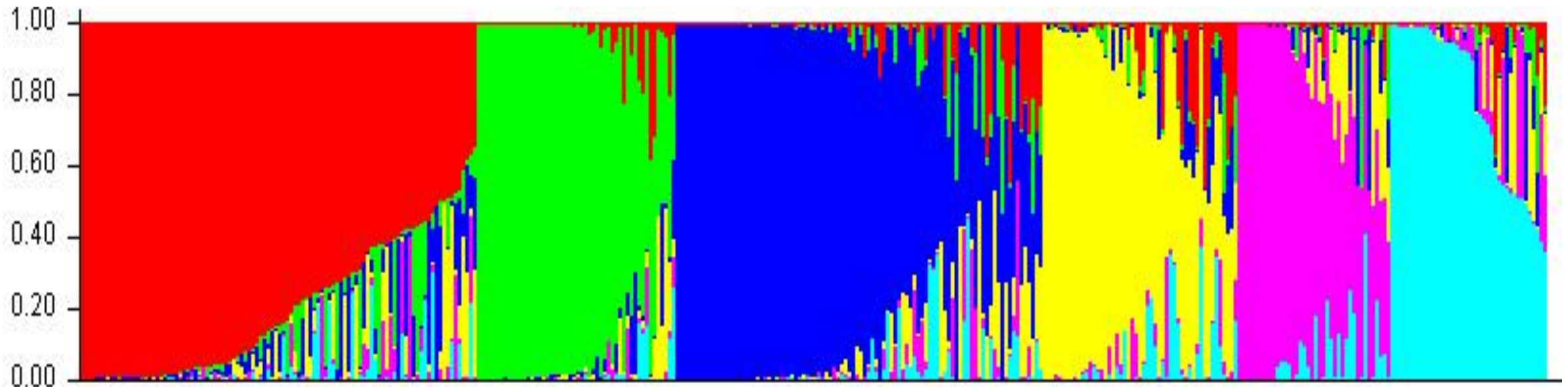
- Efficient field design with incomplete block design (e.g.,  $\alpha$ -lattice)
- Appropriate statistical methods (e.g., nearest neighbor analysis and spatial models)
- Consideration of QTL  $\times$  environmental interaction should be explored to increase the mapping power, particularly if the field conditions are not homogenous (Eskridge, 2003).

**Data collection – repository (Database)**

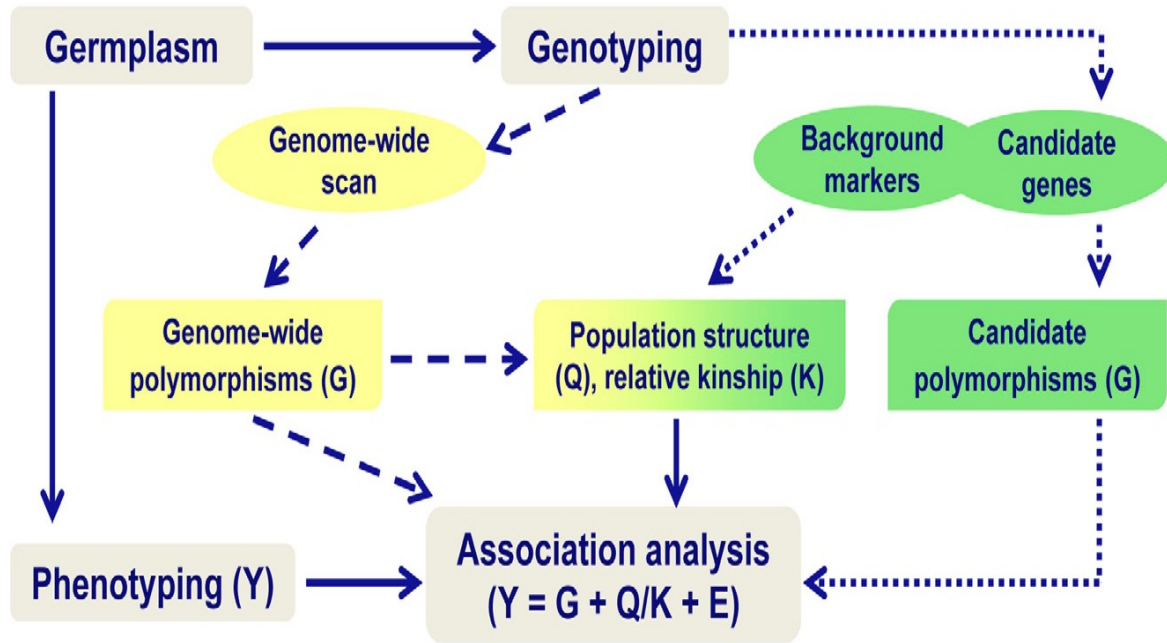
# 4. POPULATION STRUCTURE

## Statistical methods for calculating population structure

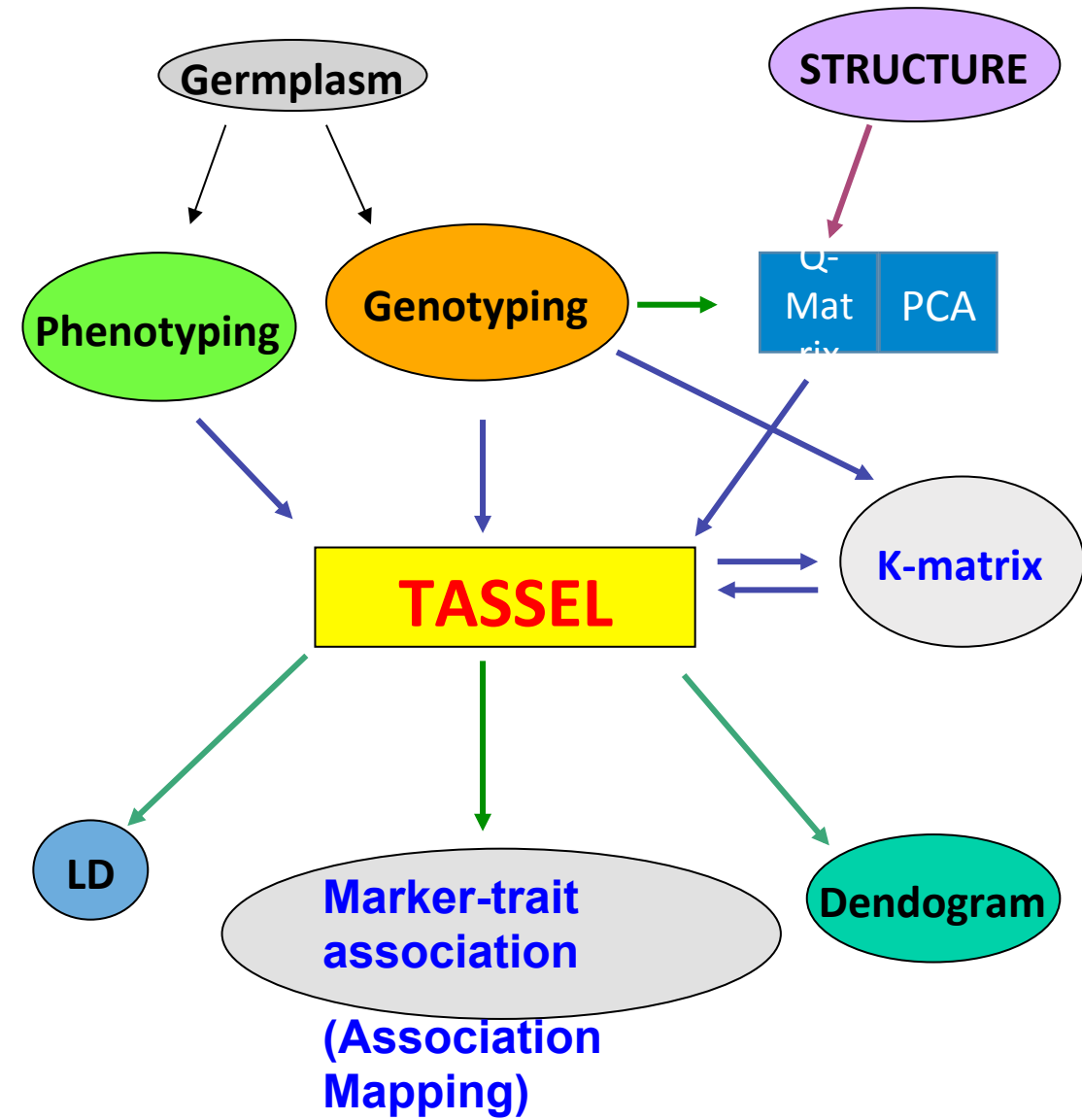
- **Structured associations (SA)** - uses a set of random markers to estimate population structure (Q) and then incorporates this estimate into further statistical analysis
- **Mixed model approach** - random markers are used to estimate Q and a relative kinship matrix (K), which are then fit into a mixed-model framework to test for marker-trait associations
- **Principal component analysis (PCA)** - summarizes variation observed across all markers into a smaller number of underlying component variables



# 5. Statistical Analysis



Genome-wide association mapping	Candidate-gene association mapping
<p>It is a comprehensive approach to systematically search the genome for causal genetic variation. A large number of markers are tested for association with various complex traits, and prior information regarding candidate genes is not required. It works best for a research consortium with complementary expertise and adequate funding.</p>	<p>Candidate genes are selected based on prior knowledge from mutational analysis, biochemical pathway, or linkage analysis of the trait of interest. An independent set of random markers needs to be scored to infer genetic relationships. It is a low cost, hypothesis-driven, and trait-specific approach but will miss other unknown loci.</p>



**TASSEL** = Trait **A**nalysis by **a**SSociation, **E**volution, **L**inkage

# Points for Analysis

**GLM** – Genotypic data + Phenotypic data + Q-matrix/PCA

**MLM** – Genotypic data + Phenotypic data + Q-matrix/PCA + K-matrix

- ❑ Association mapping without consideration of population structure would result in a high rate of false positive **Type I** errors
- ❑ Mixed-model approach to account for multiple levels of relatedness simultaneously, as detected by use of genetic markers, has improved control of both **type I** and **type II** error rates
- ❑ **False negative** : the declaration of an outcome as statistically non-significant, when the effect is actually genuine.
- ❑ **False positive** : the declaration of an outcome as statistically significant, when there is no true effect.

# AM results

Trait	Marker	marker_p	markerR2
FT	S6_51258751	0.002526	0.118313
FT	S6_49476889	0.006728	0.104095
FT	S6_49476896	0.006728	0.104095
FT	S6_51074660	0.007553	0.06921
FT	S6_51074666	0.007553	0.06921
FT	S6_57125381	0.008782	0.073529
FT	S6_57125382	0.008782	0.073529
FT	S6_50090655	0.010302	0.09183
FT	S6_48617413	0.011146	0.092174
FT	S6_50083373	0.011984	0.060825
FT	S6_49946596	0.01502	0.06651
FT	S6_48622968	0.019185	0.085937
FT	S6_48825076	0.021259	0.063186
FT	S6_51074670	0.021644	0.074513

# Association Mapping vs Linkage Mapping

## *WHICH ONE IS BEST?*

AM is not a replacement of LM rather is a **complementary** to LM

- AM – species specific and population specific
- Sps with low genetic diversity LM is superior to AM

## **Ideal method**

### **Combination of AM and LM**

- where strength of each method are used to conduct
- high-resolution power and
- high-resolution tests

# JOINT FAMILY POPULATION MAPPING

**Low detection power** – low allele frq and QTLs with small effects  
**More false positives** – genotype – phenotype covariance



Manipulate allele frequency and population structure by controlled crosses and using family mapping to enhance the power



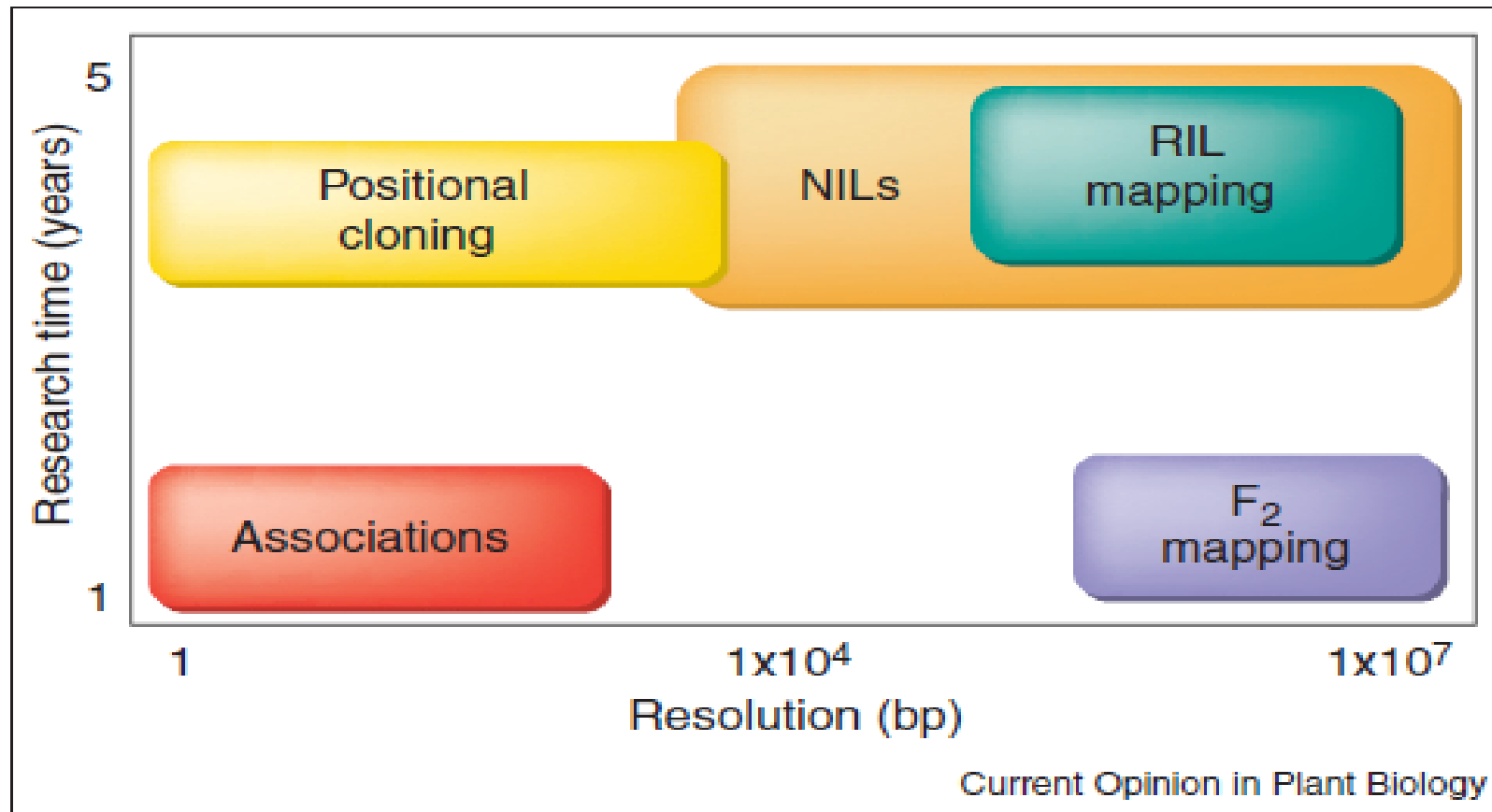
**Nested Association Mapping (NAM)**

**(Joint linkage association mapping)**

Relevant question is not what the nature of the genotyping data will be, but rather how to select germplasm to maximize the allelic diversity and the power to detect complex traits



# Resolution vs Research Time



Comparison of resolution and research time for various approaches to dissect quantitative variation. The research times assume the target species has only two generations per year. NIL, near-isogenic line; RIL, recombinant inbred line.

# **GBS-SNP APPLICATIONS IN CROP PLANTS**

## **SORGHUM & MAIZE**

- Marker Discovery**
- Phylogeny/Kinship**
- Fine-Mapping**
- Genomic Selection**
- NAM-GWAS**
- GWAS**
- Improving reference genome assembly**

# Marker Discovery

## MAIZE

Total Accessions : > 30000

Total SNPs : > 1.5M

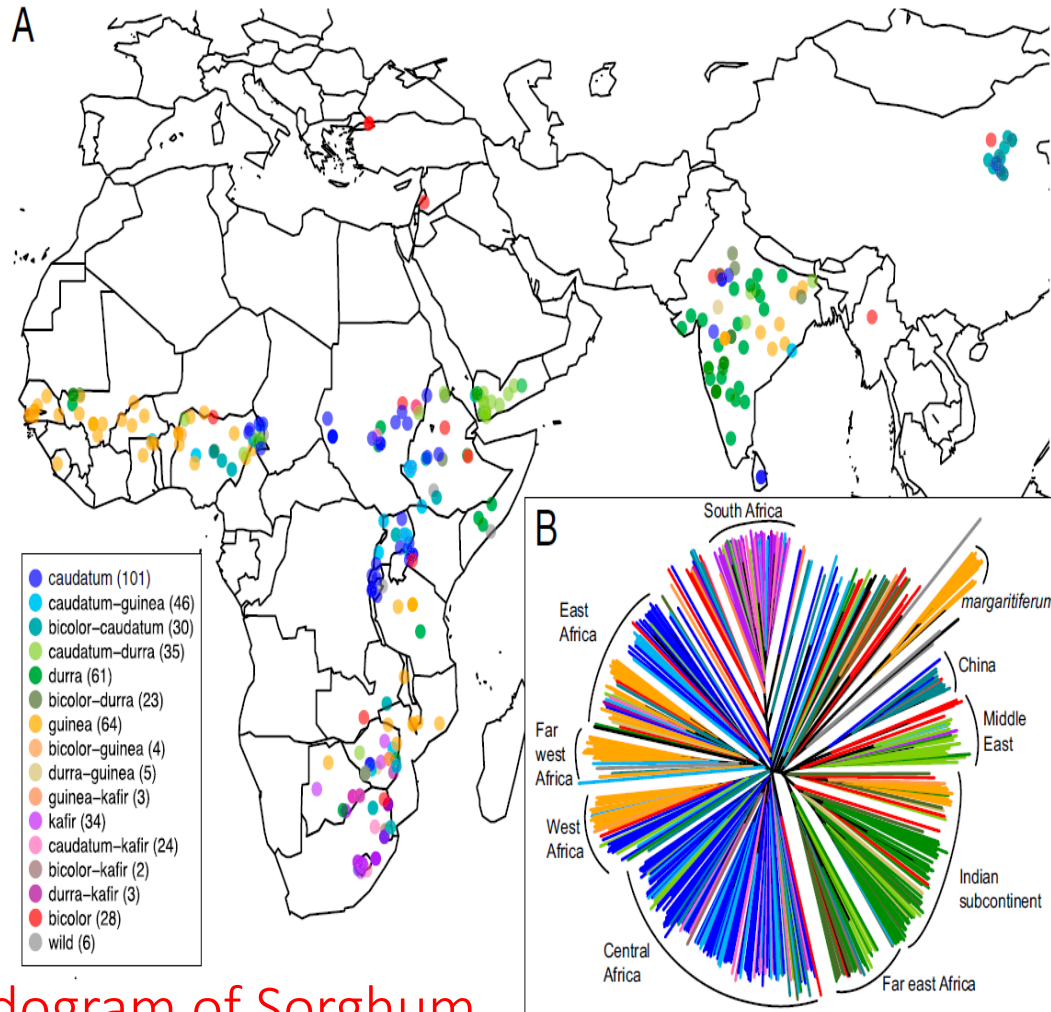
## SORGHUM

Total Accessions : 7173

Total SNPs : 556,969

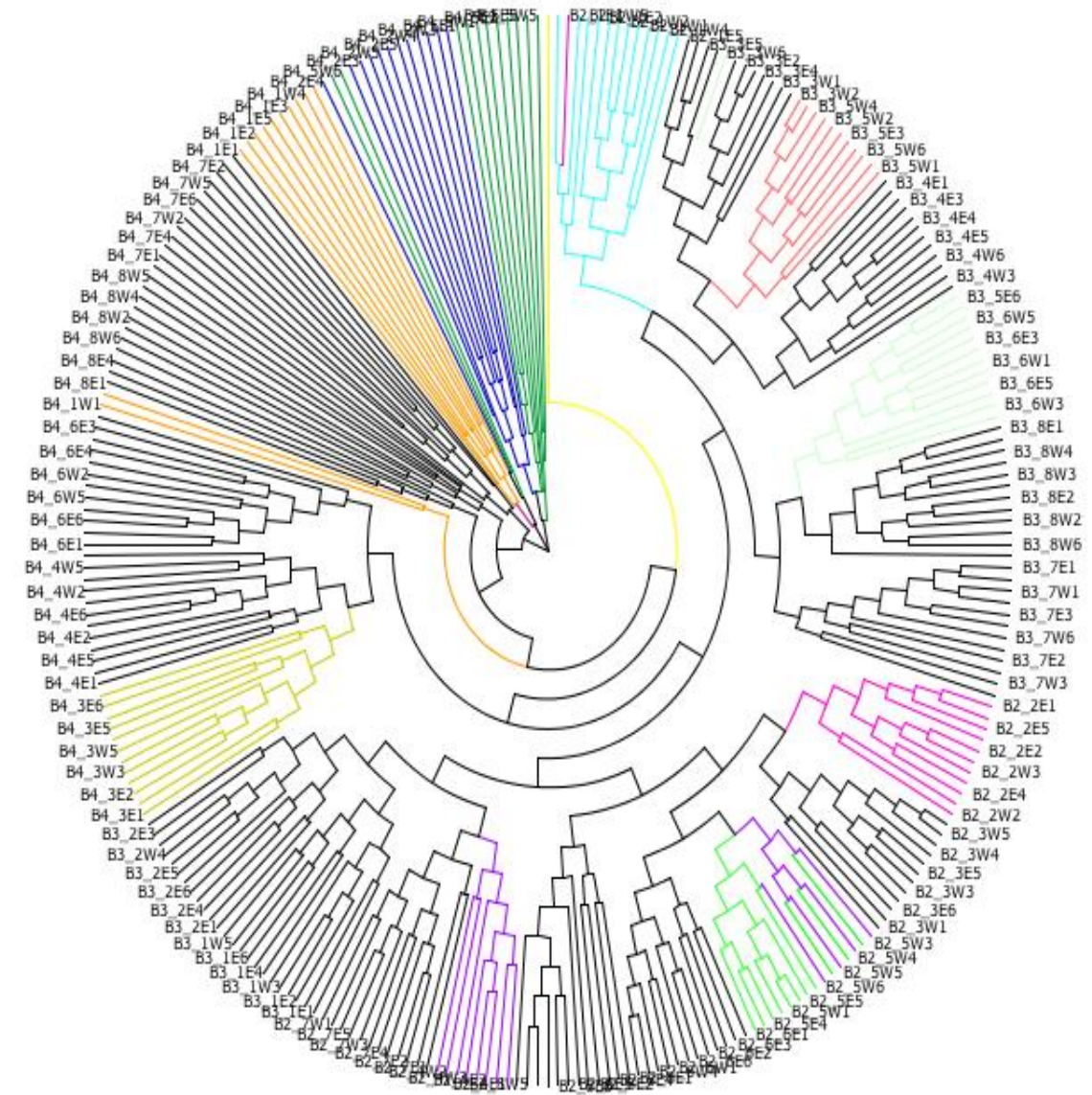
Total Data points : 3,995,138,637

# Genetic Diversity/phylogeny



Cladogram of Sorghum

Source: Dr. Ram Sharma



Cladogram of Maize Hybrid and Landraces

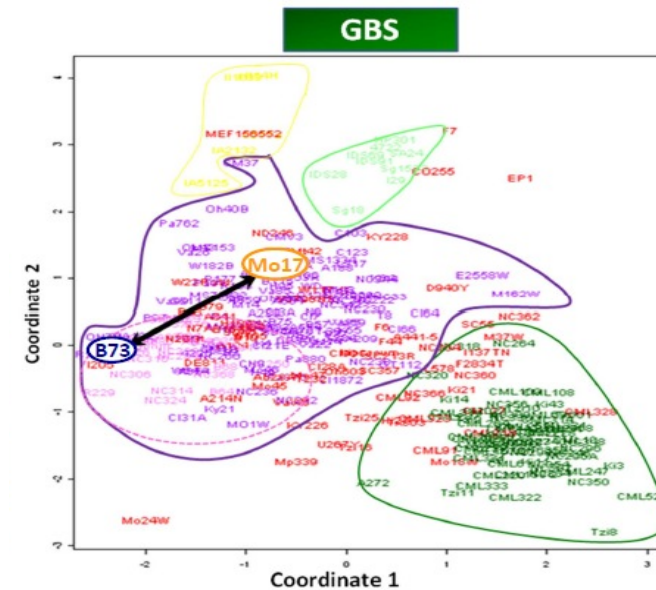
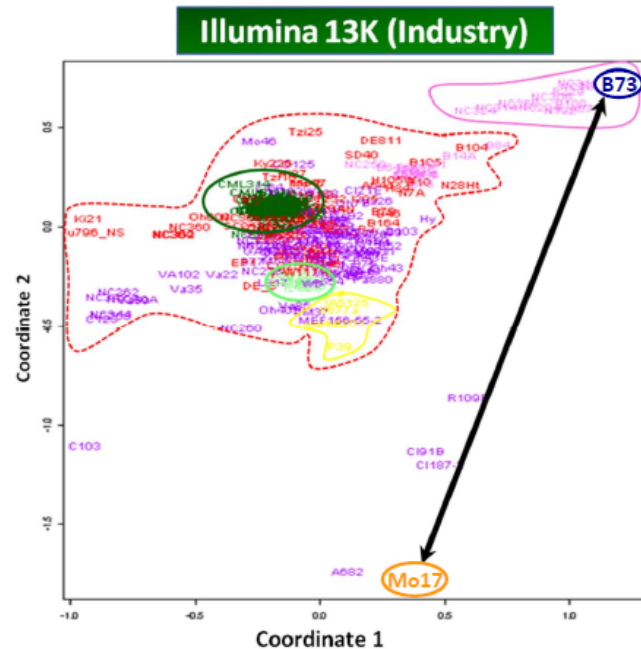
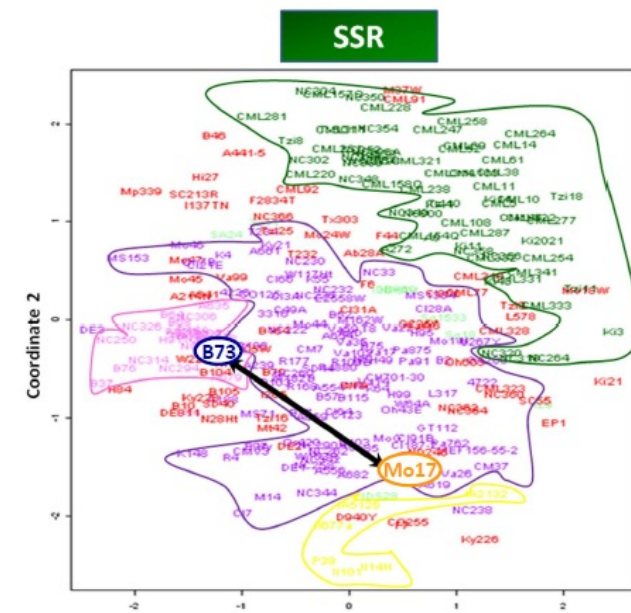
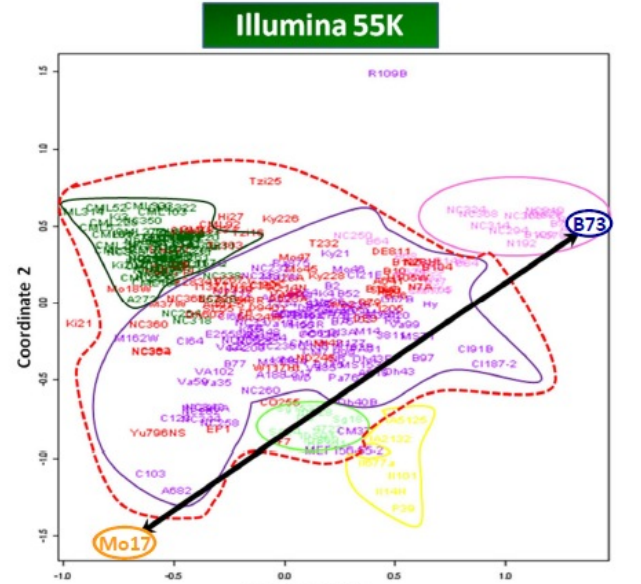


# Biasness in coordinate analysis in SSR, Chip-SNPs and GBS-SNPs in maize

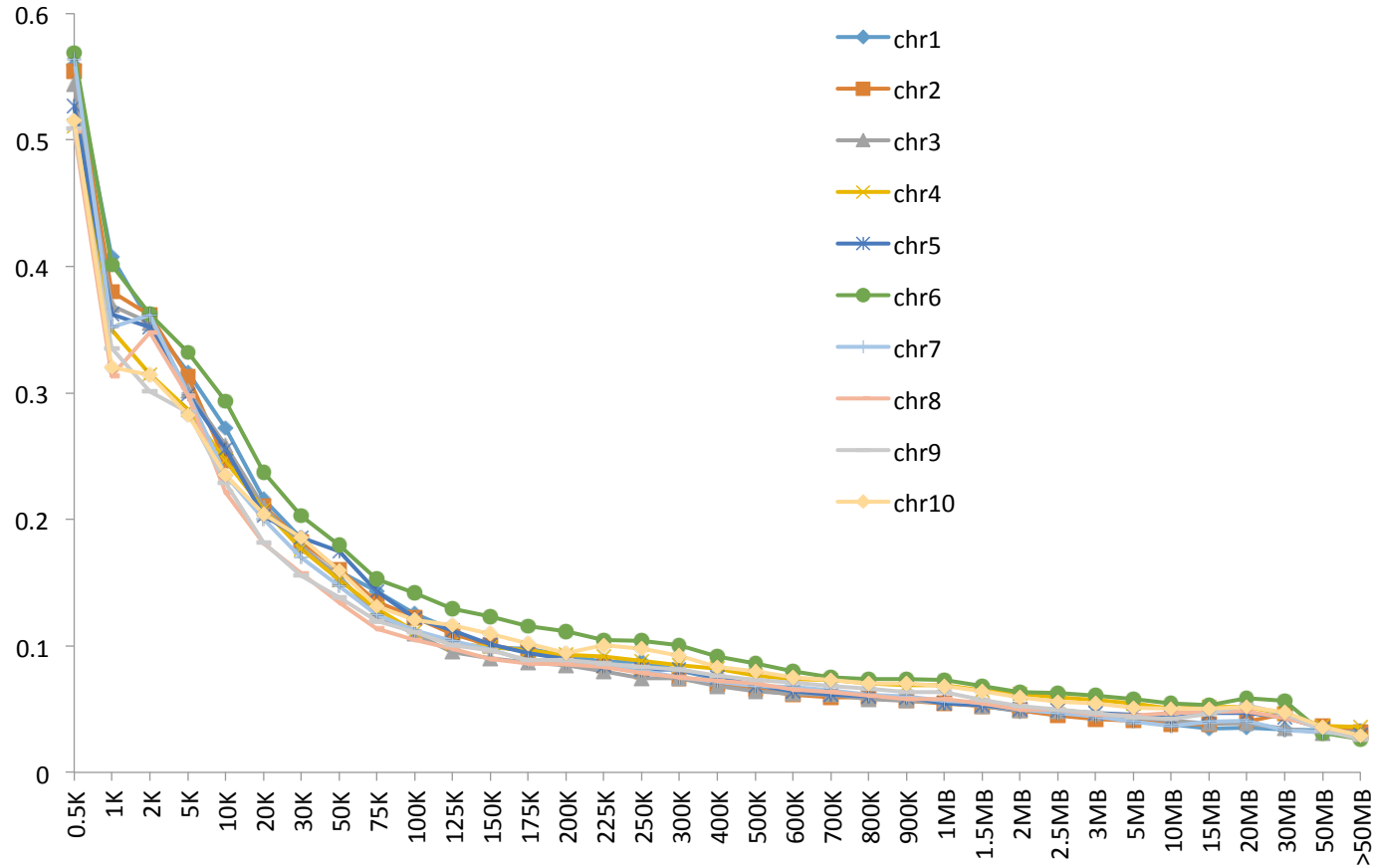
- Non stiff stalk, nss (106)
- Stiff stalk, ss (28)
- Tropical/subtropical (66)
- Popcorn (9)
- Sweet corn (6)
- Unclassified (67)
- B73 (ss)
- Mo17 (nss)



Source: Dr. Ram Sharma



# LD decay in sorghum



Chr	LD decay (kb)
1	125-150
2	150-175
3	100-125
4	125-150
5	150-175
6	300-400
7	125-150
8	100-125
9	125-150
10	225-250

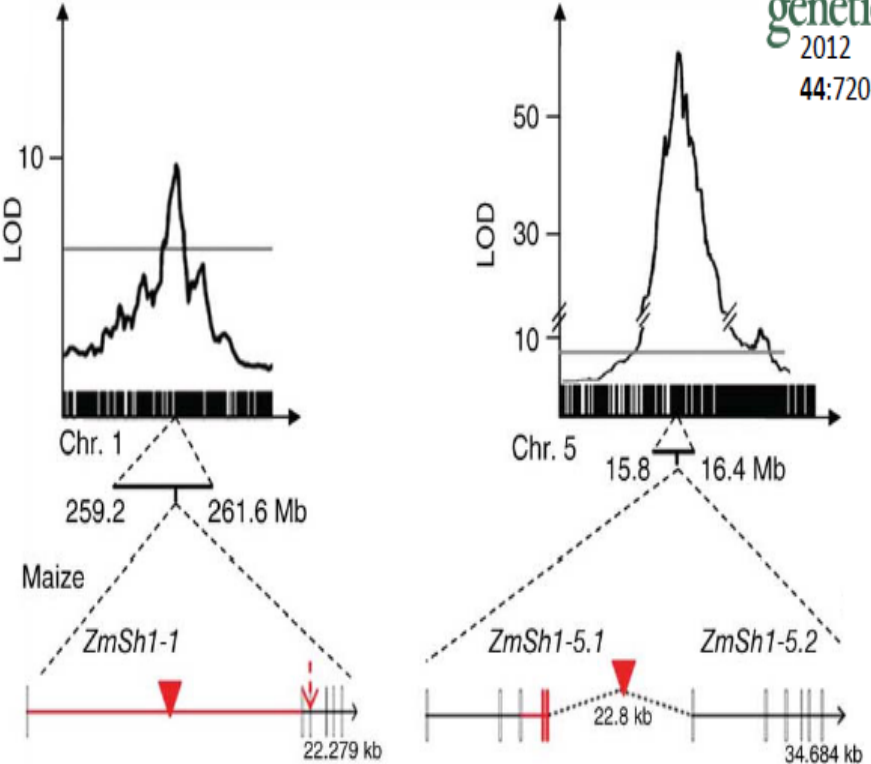
Source: Dr. Punna Ramu

# Bi-Parental Linkage mapping/Fine Mapping

## Parallel domestication of the *Shattering1* genes in cereals

Zhongwei Lin<sup>1</sup>, Xianran Li<sup>1</sup>, Laura M Shannon<sup>2</sup>, Cheng-Ting Yeh<sup>3,4</sup>, Ming L Wang<sup>5</sup>, Guihua Bai<sup>1,6</sup>, Zhao Peng<sup>7</sup>, Jiarui Li<sup>7</sup>, Harold N Trick<sup>7</sup>, Thomas E Clemente<sup>8</sup>, John Doebley<sup>2</sup>, Patrick S Schnable<sup>3,4</sup>, Mitchell R Tuinstra<sup>9</sup>, Tesfaye T Tesso<sup>1</sup>, Frank White<sup>7</sup> & Jianming Yu<sup>1</sup>

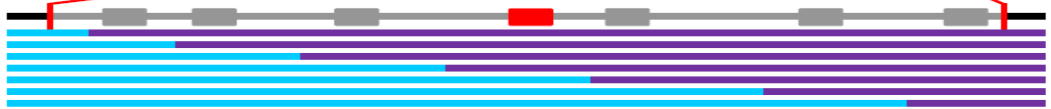
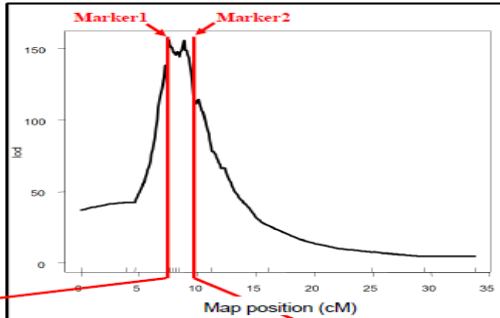
nature  
genetics  
2012  
44:720-725



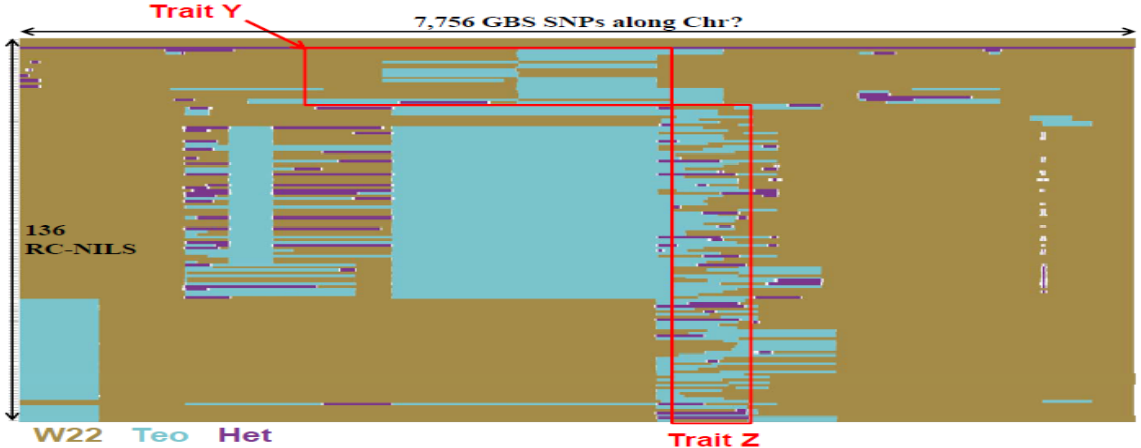
Maize *Sh1* orthologs are located at seed shattering QTLs

## Fine mapping QTL

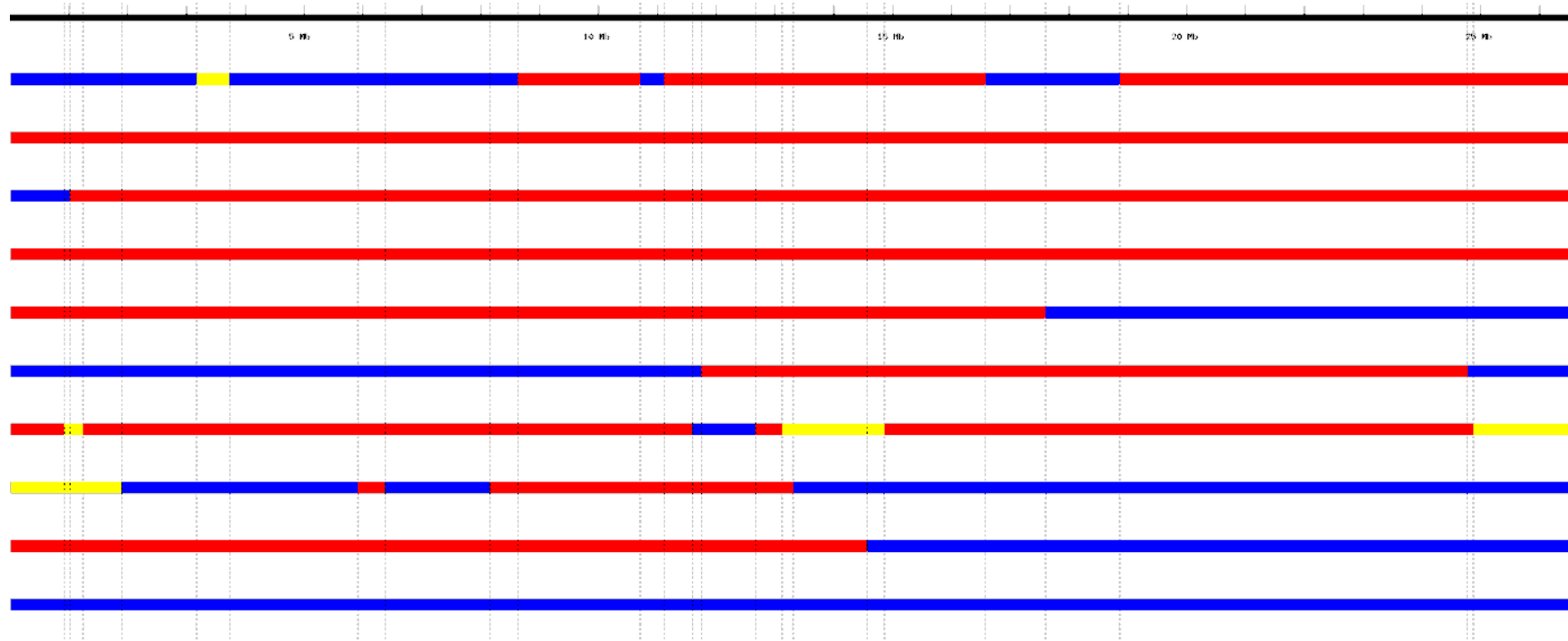
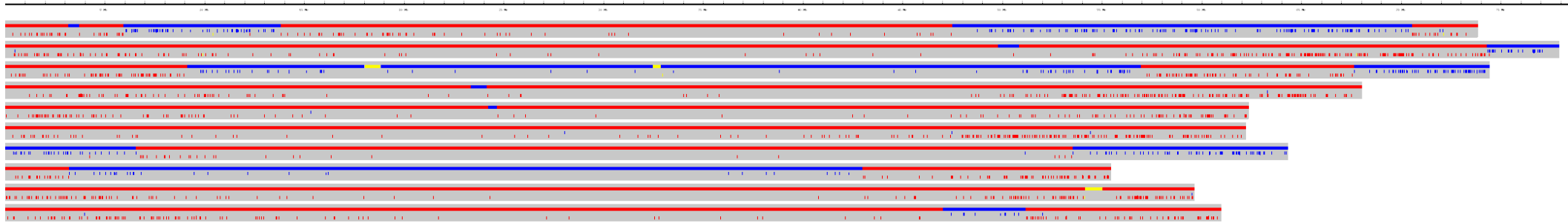
- Need to saturate interval containing QTL with markers
- GBS a good source of markers
- Also need to collect recombinants in the interval
- Near-isogenic lines (NILs) helpful (Mendelize)
- Good reference genome



## Fine Mapping of Domestication QTL in Maize



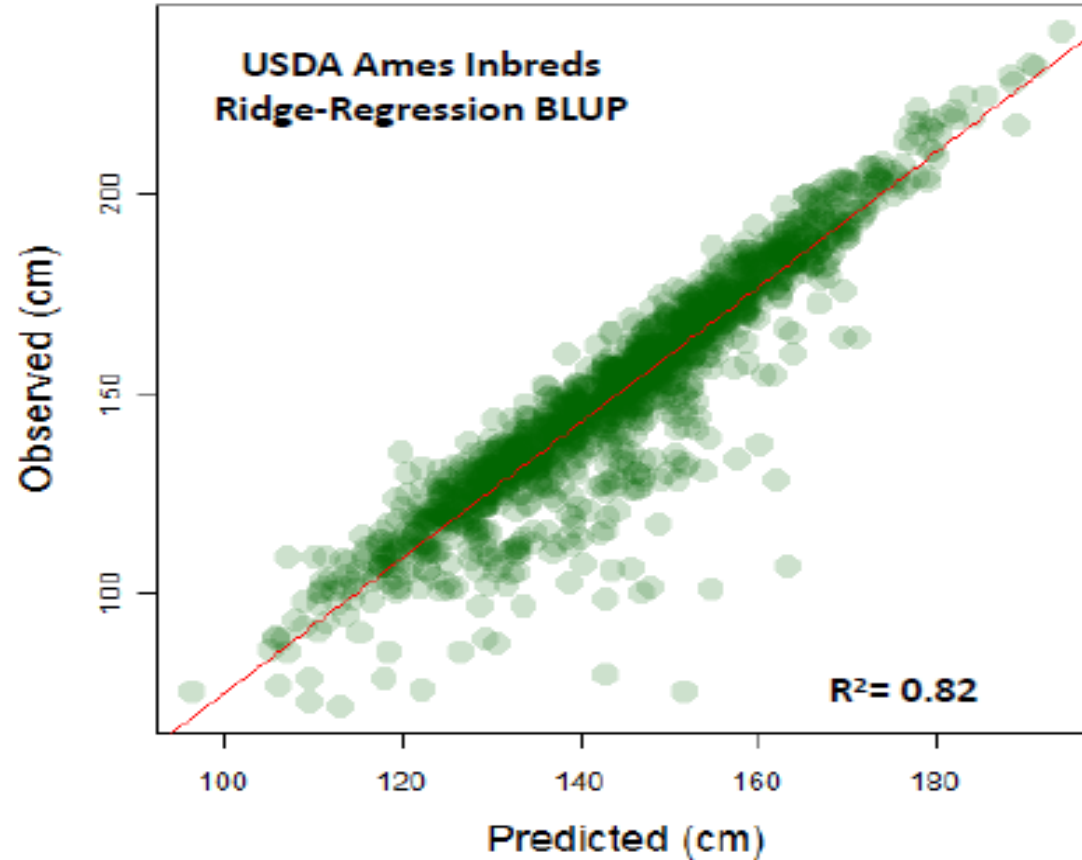
# Recombination break points





# Genomic prediction in maize

2800 diverse maize breeding lines



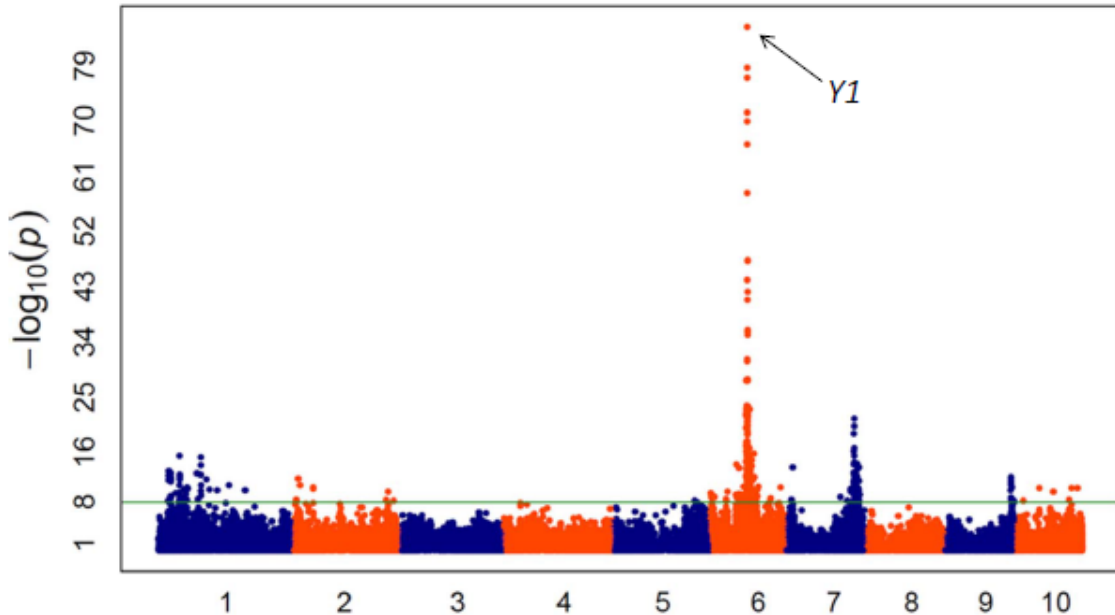
Source : Jason Peiffer

Accurate genomic prediction of height based on GBS data in maize

# GWAS examples in plants

GWAS directly hits known Mendelian traits

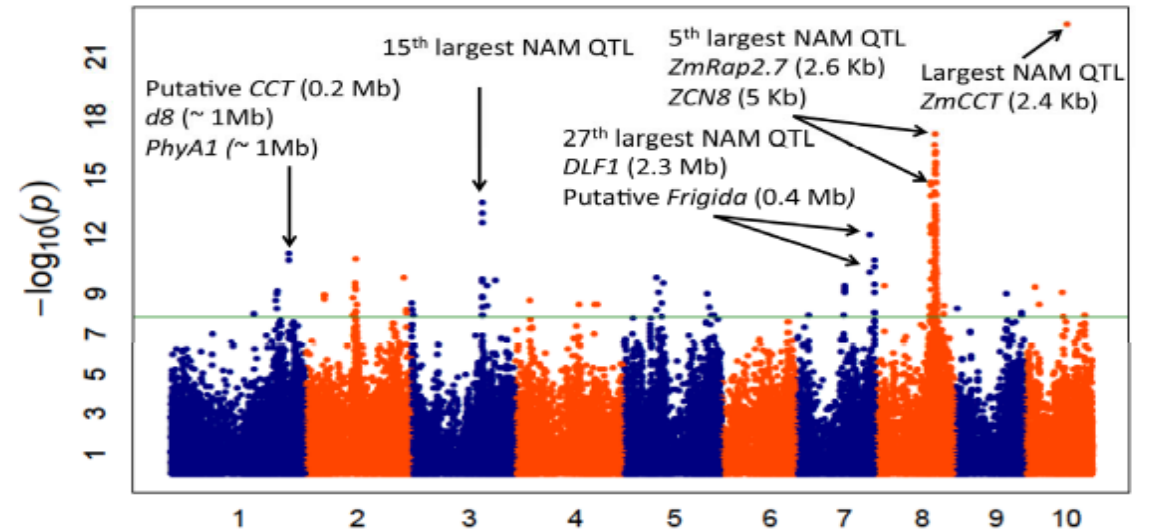
GWAS of white vs. yellow kernels in 1,595 USDA Ames inbreds



The best hit for kernel color lies **within Y1**

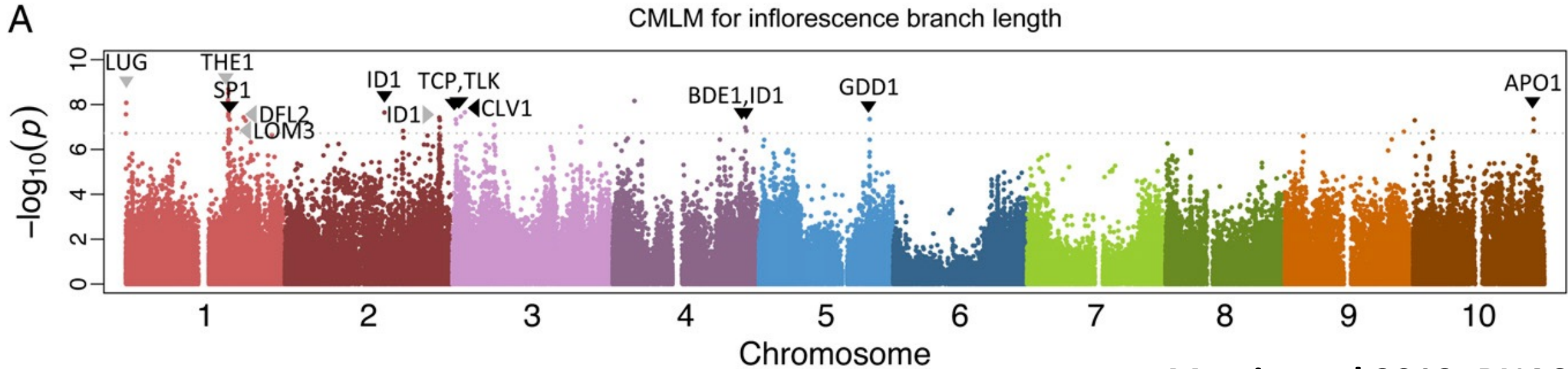
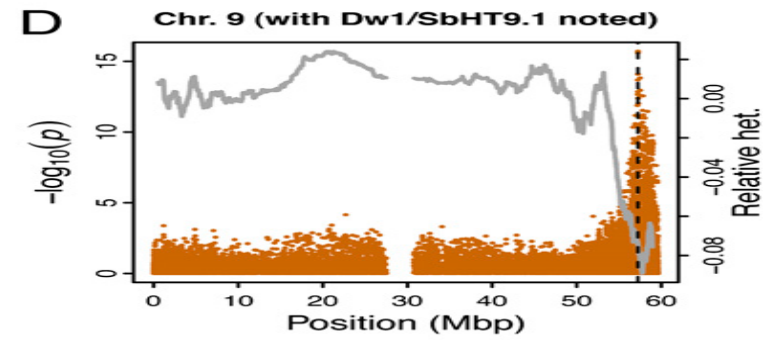
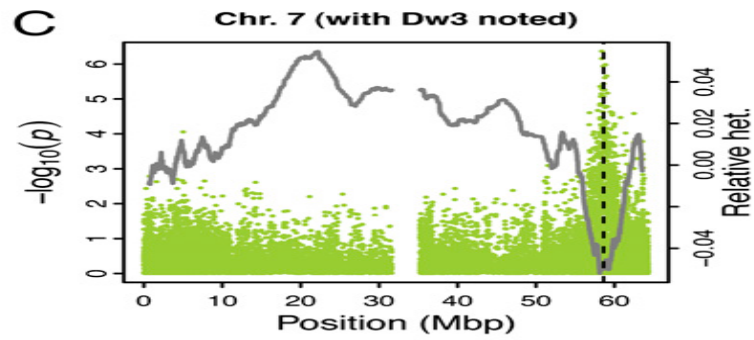
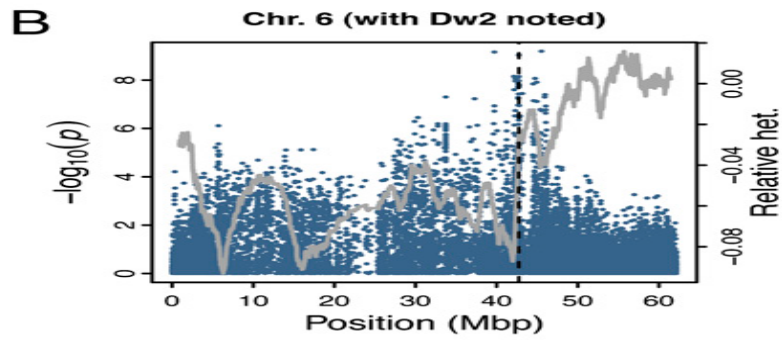
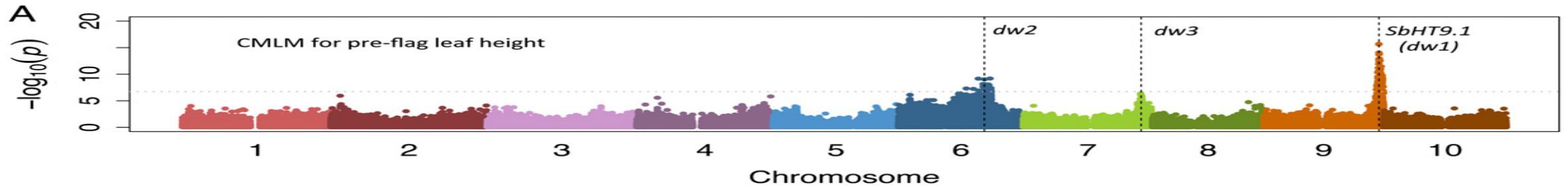
GWAS of a more complex trait directly hits known flowering time

GWAS of growing degree days to silking in 2,279 inbreds



Even with ~660K SNPs we almost missed ZmCCT (only one significant SNP)

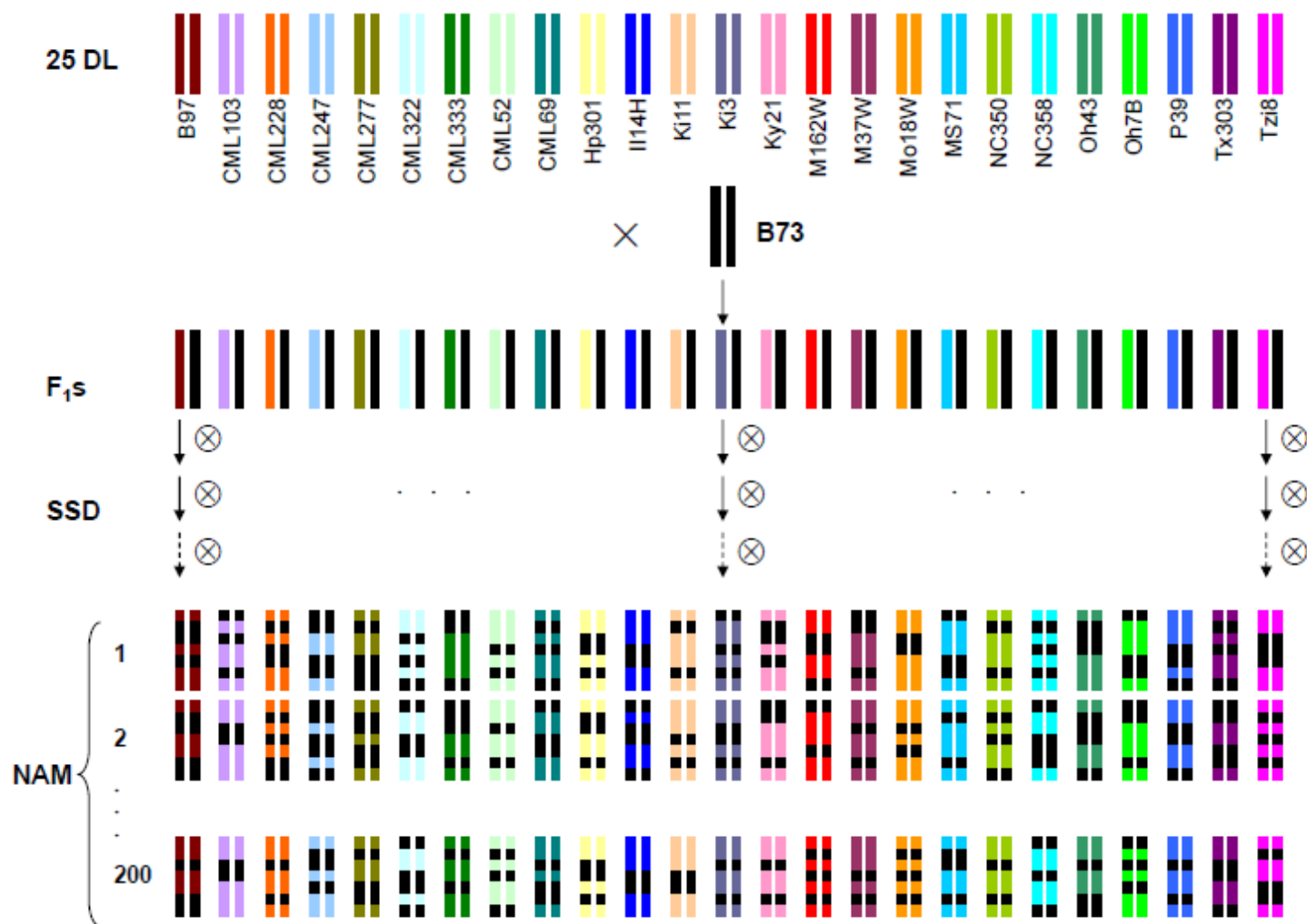
Cont....



# NAM (nested association mapping)-GWAS in maize

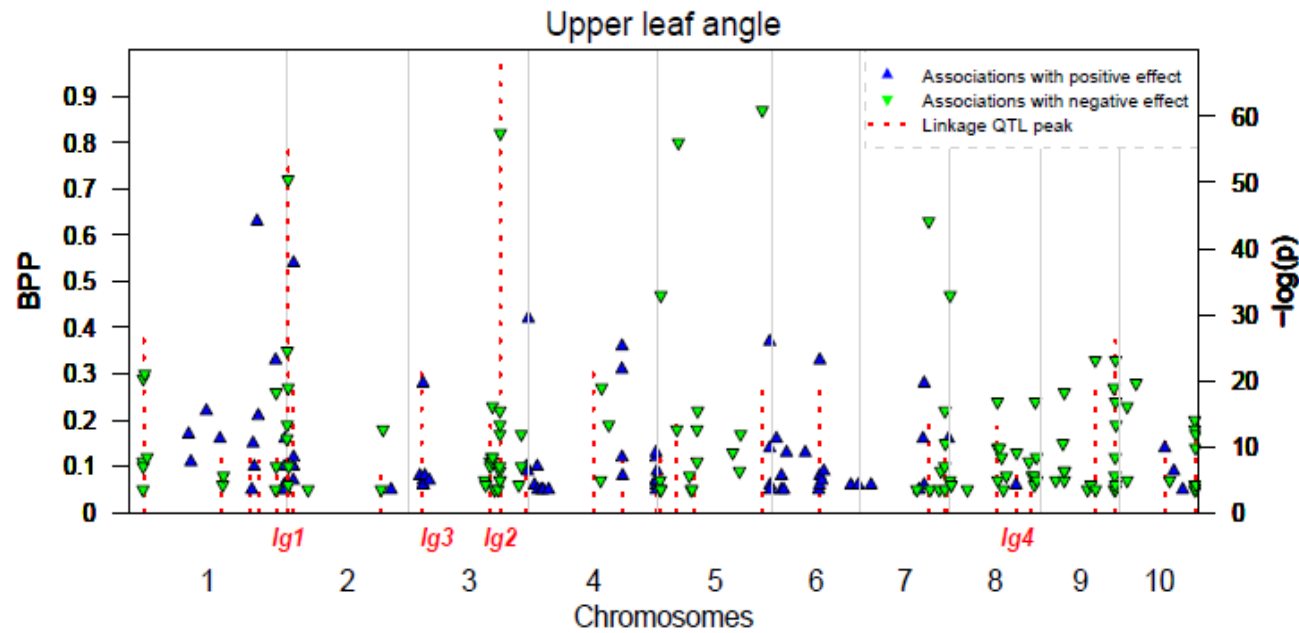
Development of NAM mapping population in maize

The maize NAM population was built for NAM-GWAS



Cont....

## *liguleless1* and *liguleless2* explain the two “biggest” leaf angle QTL



Tian, Bradbury, et al 2011 Nature Genetics

GBS data improves the resolution of joint linkage analysis in the NAM population

Trait: Days to silk (flowering time)

Markers	Median QTL support interval
1,106 array SNPs	6.2 cM
171,479 GBS SNPs	2.6 cM

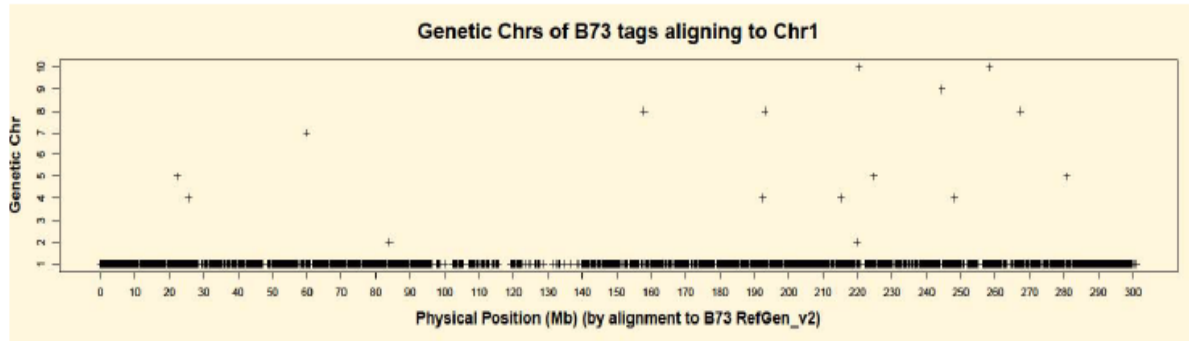
# The Maize B73 reference genome: room for improvement?

- ❑ The B73 reference genome accurate for B73 but less so for other maize lines like M017
- ❑ Even for B73, some regions of the genome are in the wrong location
- ❑ Some large (multiple BAC) contigs could not be anchored
  - assigned to 'chromosome 0'
  - 30 chr0 contigs in B73 RefGenV1
  - 30 chr0 contigs in B73 RefGenV2
- ❑ Some regions of the genome are missing
  - ~ 5% of the B73 sequence is not in the B73 reference genome



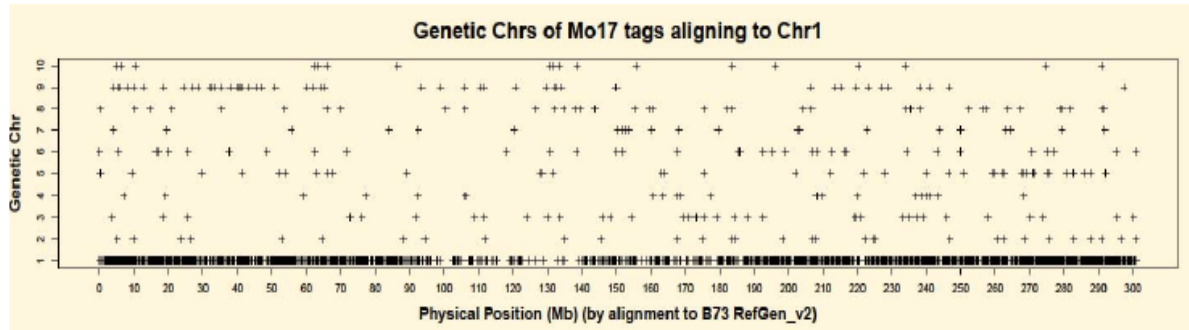
# Cont....

## B73 reference genome highly accurate for B73...



- 0.4% of B73 tags genetically map to different chromosome than they align to

## ...but far less so for other maize lines



- 9.3% of Mo17 tags genetically map to different chromosome than they align to

## Some chunks of the B73 reference genome are in the wrong place

Physical Chr	Start (Mb)	End (Mb)	Genetic Chr	Approx. Genetic Location (Mb)	# Tags
10	139.3	139.8	2	16.5–16.8	49
9	102.5	106.9	9	15–32	49
7	150.1	161.8	5	192–214	13
10	0.2	0.4	4	83–151	12
8	48.4	50	2	61–127	12
10	0.07	0.2	7	47–100	9
2	231.2	231.2	7	18–26	8
3	228.1	230.5	5	194–212	6

# Acknowledgements



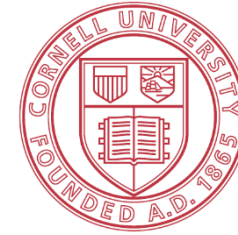
Odeny Damaris

Stefania Grando

Santosh Deshpande

Trushar Shah

Punna Ramu



Ed Buckler

Jeff Glaubitz

James Harriman

Terry C

Rob Elshire

Sharon E. Mitchell

Alex Lipka

Fei Lu



# Thank You



for your kind attention