

# GENOTYPING, LINKAGE ANALYSIS AND HAPLOTYPE PHASING IN AUTOPOLYPLOIDS

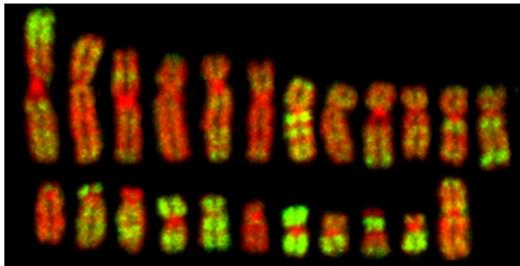
Marcelo Mollinari  
[mmollin@ncsu.edu](mailto:mmollin@ncsu.edu)

Bioinformatics Research Center - Hicks Hall 368  
Department of Statistics  
North Carolina State University

September 12 and 17, 2019 - GN 757 001 - North Carolina State University

## PLOIDY LEVEL

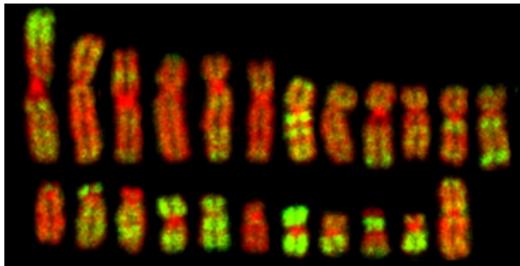
- ▶ **Basic chromosome number:** the number of different chromosomes that make up a single complete set. Humans,  $x = 23$



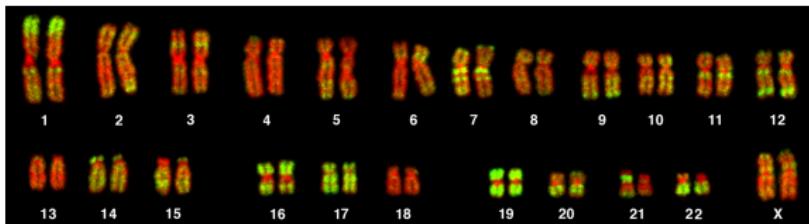
- ▶ **Ploidy level:** Number of basic chromosome sets. Humans  $p = 2$

## PLOIDY LEVEL

- ▶ **Basic chromosome number:** the number of different chromosomes that make up a single complete set. Humans,  $x = 23$

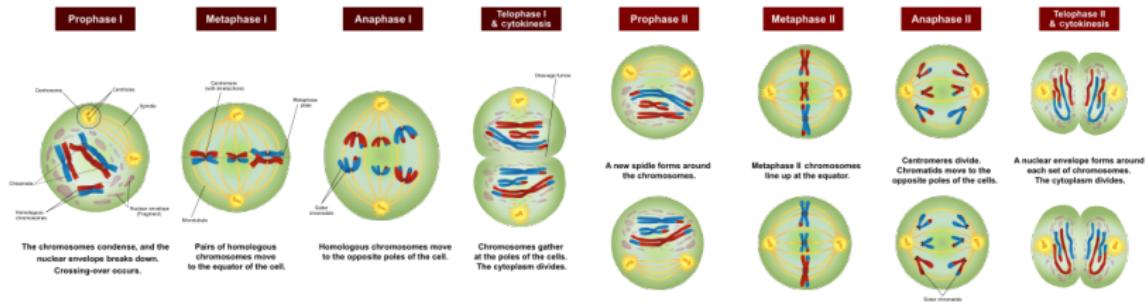


- ▶ **Ploidy level:** Number of basic chromosome sets. Humans  $p = 2$



# MEIOTIC PAIRING AND GAMETE FORMATION IN DIPLOIDS - REVIEW

- ▶ Meiotic process is quite stable
- ▶ It results in four products with  $n = x = \text{basic number}$



# MEIOTIC PAIRING AND GAMETE FORMATION IN DIPLOIDS

Animation

# POLYPLOID SPECIES

Organisms that have multiple sets of chromosomes

Diploid



Diploid



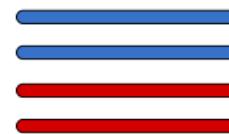
# POLYPLOID SPECIES

Organisms that have multiple sets of chromosomes

Autotetraploid



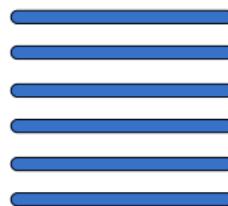
Allotetraploid



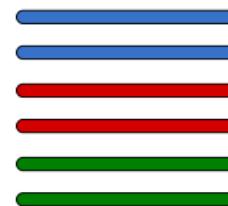
# POLYPLOID SPECIES

Organisms that have multiple sets of chromosomes

Autohexaploid



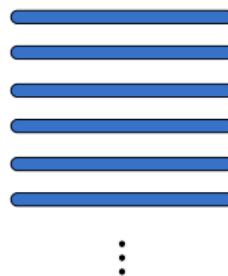
Allohexaploid



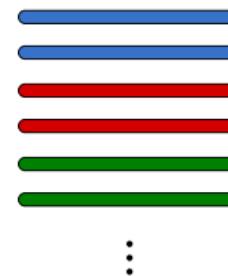
# POLYPLOID SPECIES

Organisms that have multiple sets of chromosomes

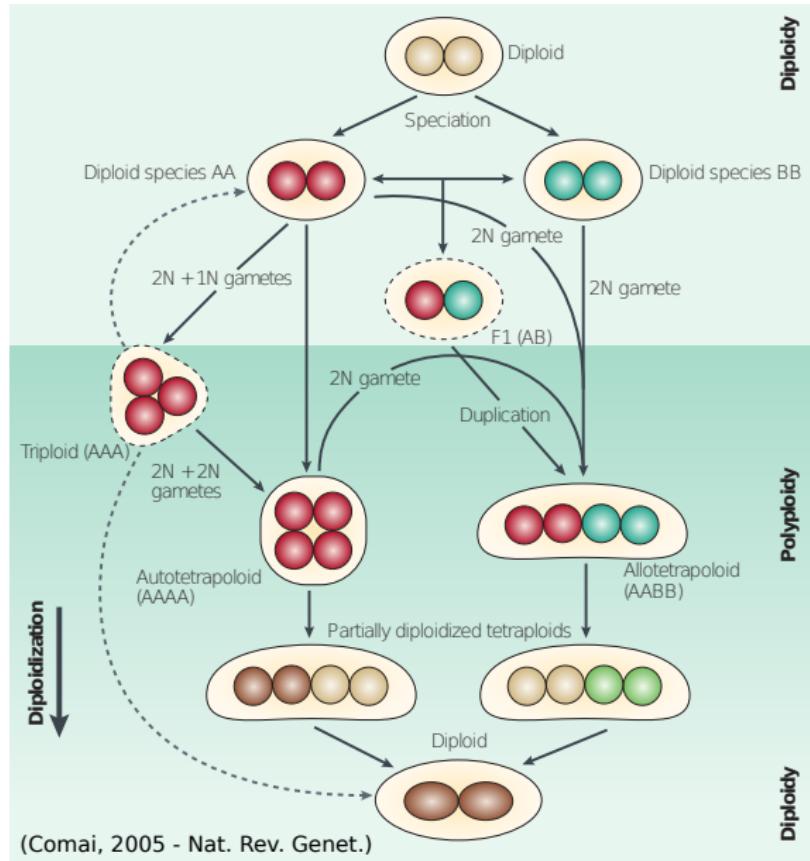
Autohexaploid



Allohexaploid

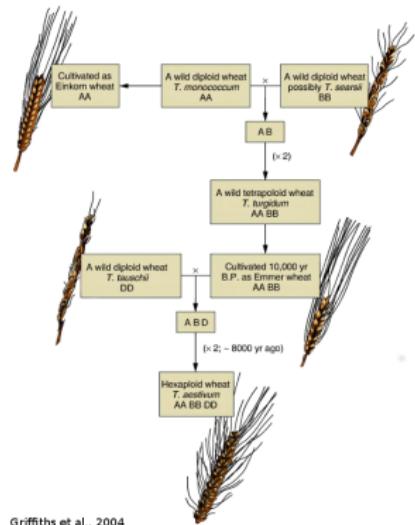


# HOW ARE POLYPLOIDS FORMED?

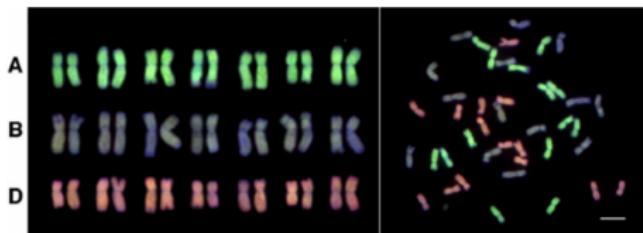


# ALLOPOLYPLOID EXAMPLE - WHEAT

WHEAT - *Triticum aestivum*

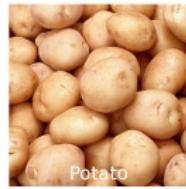


Griffiths et al., 2004



Zhang H et al. PNAS 2013;110:3447-3452

# AUTOPOLYPLOID EXAMPLE



Chrysanthemum

Sweet potato

Potato

Rose

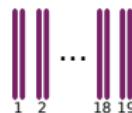
Forage crops

Sugarcane

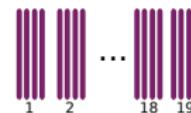


Copyright Leonard Lessing/Peter Arnold Inc.

Diploid

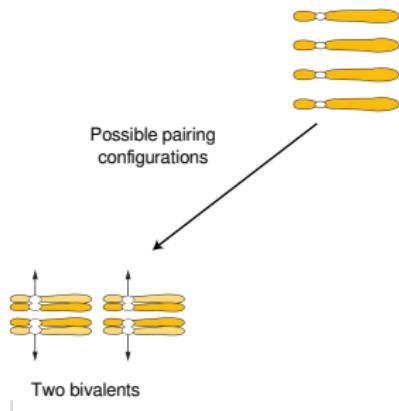


Autotetraploid



# MEIOTIC PAIRING AND GAMETE FORMATION IN AUTOPOLYPLOIDS

## Tetraploid example

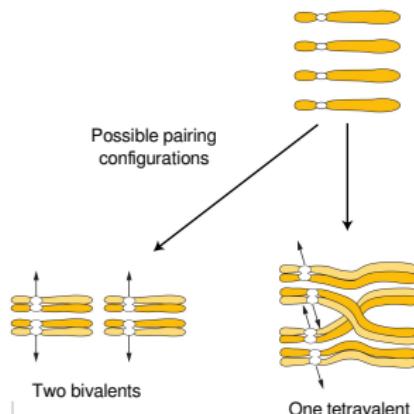


Griffiths et al. (2004)

Usually, for high ploidy levels, we assume **bivalent** pairing

# MEIOTIC PAIRING AND GAMETE FORMATION IN AUTOPOLYPLOIDS

## Tetraploid example

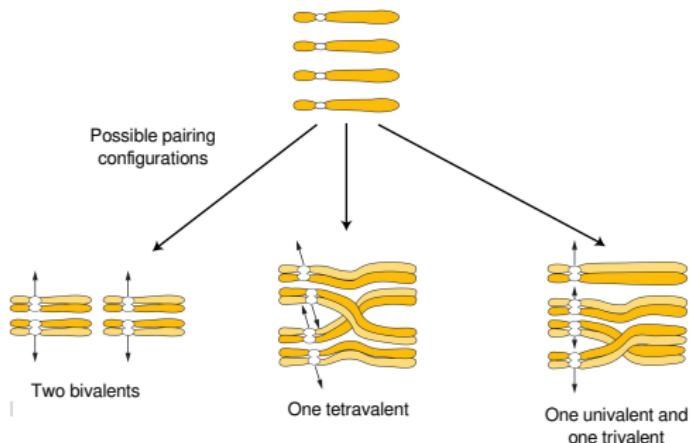


**Griffiths et al. (2004)**

Usually, for high ploidy levels, we assume **bivalent** pairing

# MEIOTIC PAIRING AND GAMETE FORMATION IN AUTOPOLYPLOIDS

## Tetraploid example

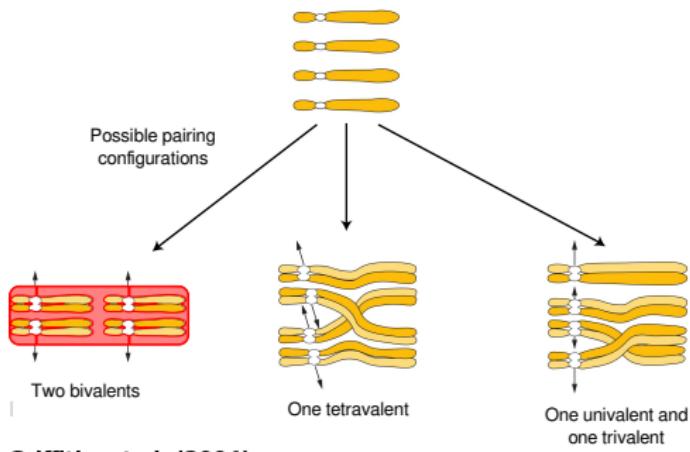


Griffiths et al. (2004)

Usually, for high ploidy levels, we assume **bivalent** pairing

# MEIOTIC PAIRING AND GAMETE FORMATION IN AUTOPOLYPLOIDS

## Tetraploid example

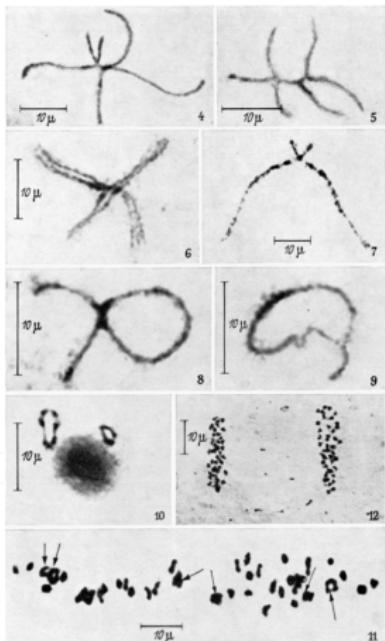


Griffiths et al. (2004)

Usually, for high ploidy levels, we assume **bivalent** pairing

# MEIOTIC PAIRING IN AUTOPOLYPLOIDS

Sweetpotato



Hexavalents, quadrivalents and bivalents in sweetpotato (Magoon *et al.* 1970)

Sugarcane

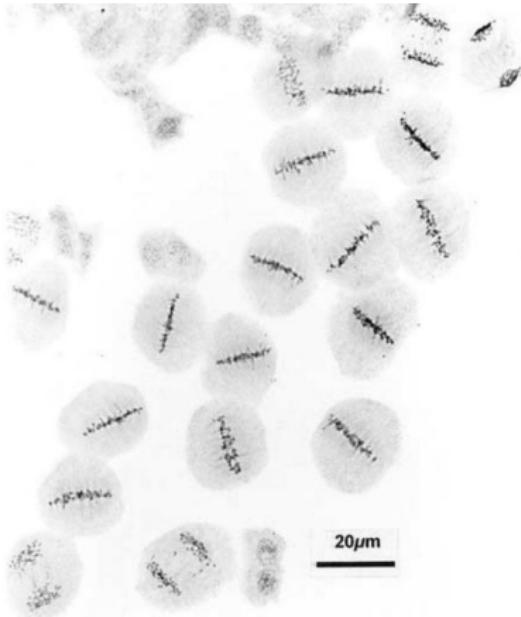
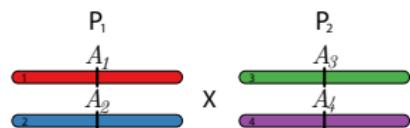


Figure 1. PMCs of *Saccharum* spp. hybrid clone 79N9059 at meiosis. As was the case in other clones, pairing was regular, bivalents generally formed. (Bielig *et al.* 2003)

# ACCESSING THE ALLELIC VARIATION - DIPLOIDS

Multiallelic



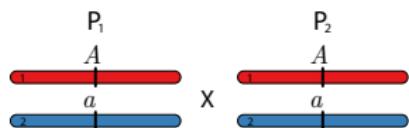
Gametes

	$P_1$	$A_1$	$A_2$
$P_2$		Red	Blue
$A_3$		Green	Green
$A_4$		Red	Blue

4 possible genotypes

1:1:1:1

Biallelic



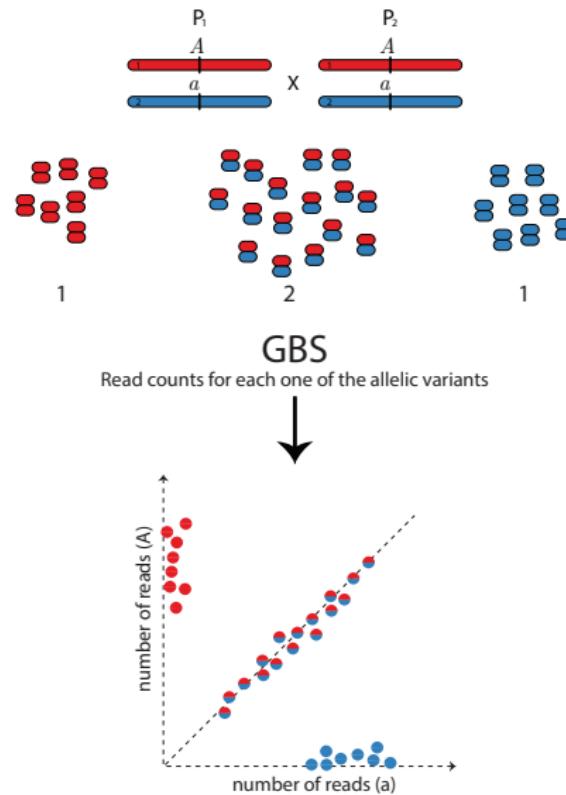
Gametes

	$P_1$	$A$	$a$
$P_2$		Red	Blue
$A$		Red	Red
$a$		Blue	Blue

3 possible genotypes

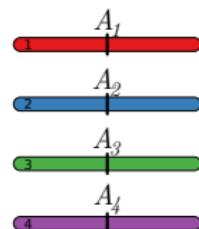
1:2:1

# ACCESSING THE ALLELIC VARIATION - DIPLOIDS



# BIALLELIC VS. MULTIALLELIC - POLYPLOIDS

Codominant  
Multiallelic



$A_1A_2A_3A_4$

$A_1A_2$

$A_1A_3$

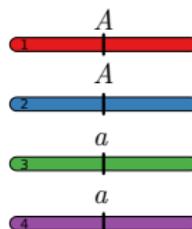
$A_1A_4$

$A_2A_3$

$A_2A_4$

$A_3A_4$

Codominant  
Biallelic



$AAaa$   
(2 doses)

$AA$

$Aa$

$Aa$

$Aa$

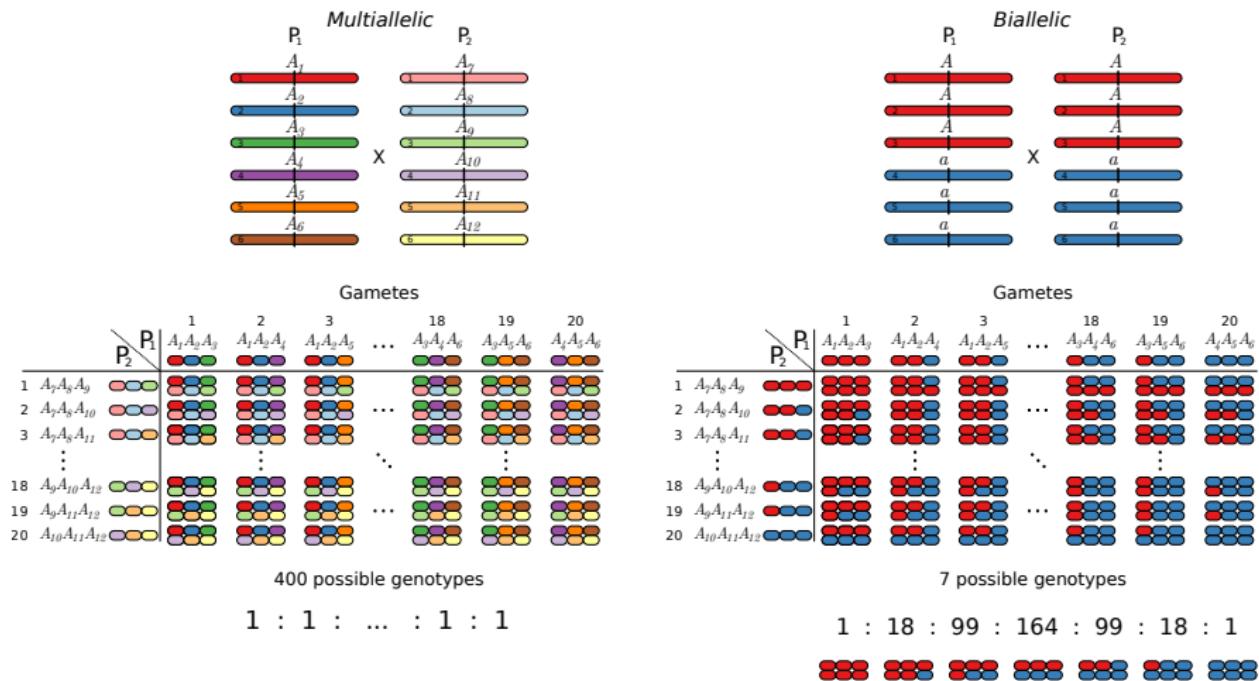
$Aa$

$aa$

Number of possible gametes considering  
one locus in one and two parents

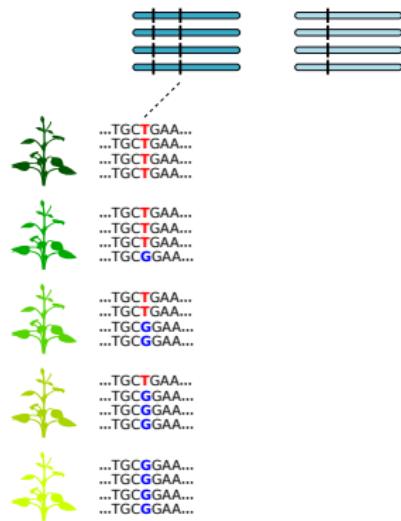
Ploidy	$\left(\frac{p}{2}\right)$	$\left(\frac{p}{2}\right)^2$
4	6	36
6	20	400
8	70	4900
10	252	63504
12	924	853776
14	3432	11778624
16	12870	165636900

# ACCESSING THE ALLELIC VARIATION - HEXAPLOID EXAMPLE



# GENOTYPING IN POLYPLOIDS

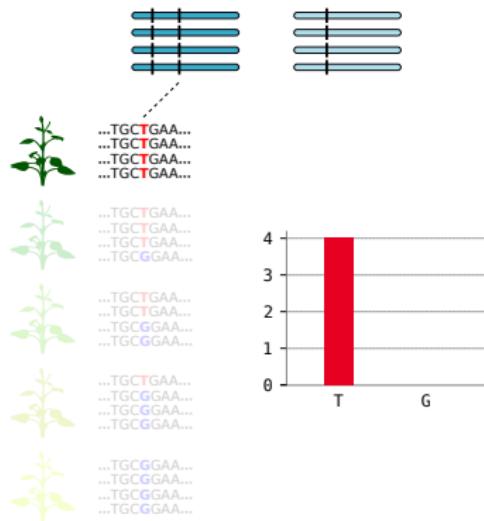
## BIALLELIC CO-DOMINANT MARKER



**Quantitative genotyping:** access the abundance of each one of the allelic forms (in a biallelic marker)

# GENOTYPING IN POLYPLOIDS

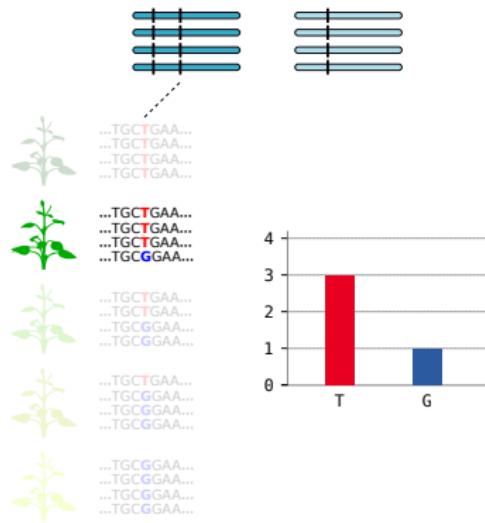
## BIALLELIC CO-DOMINANT MARKER



**Quantitative genotyping:** access the abundance of each one of the allelic forms (in a biallelic marker)

# GENOTYPING IN POLYPLOIDS

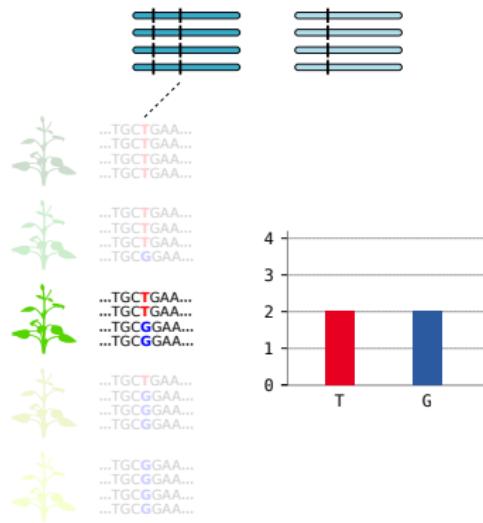
## BIALLELIC CO-DOMINANT MARKER



**Quantitative genotyping:** access the abundance of each one of the allelic forms (in a biallelic marker)

# GENOTYPING IN POLYPLOIDS

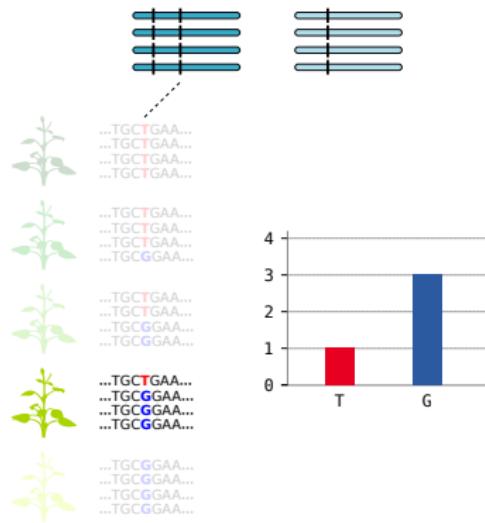
## BIALLELIC CO-DOMINANT MARKER



**Quantitative genotyping:** access the abundance of each one of the allelic forms (in a biallelic marker)

# GENOTYPING IN POLYPLOIDS

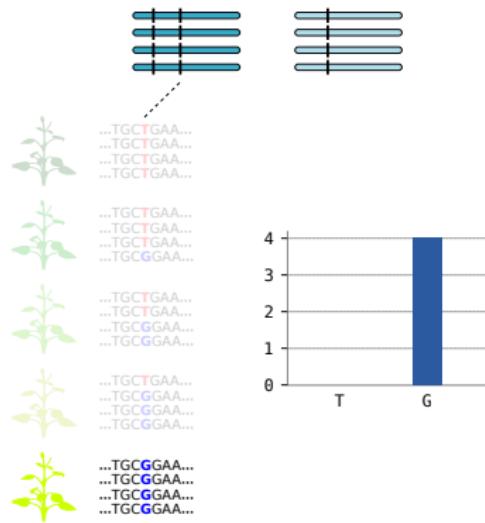
## BIALLELIC CO-DOMINANT MARKER



**Quantitative genotyping:** access the abundance of each one of the allelic forms (in a biallelic marker)

# GENOTYPING IN POLYPLOIDS

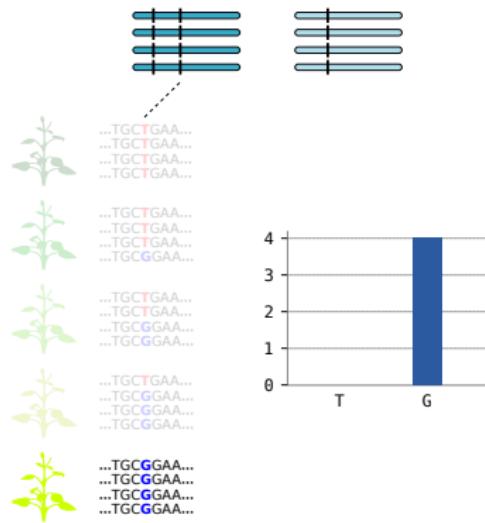
## BIALLELIC CO-DOMINANT MARKER



**Quantitative genotyping:** access the abundance of each one of the allelic forms (in a biallelic marker)

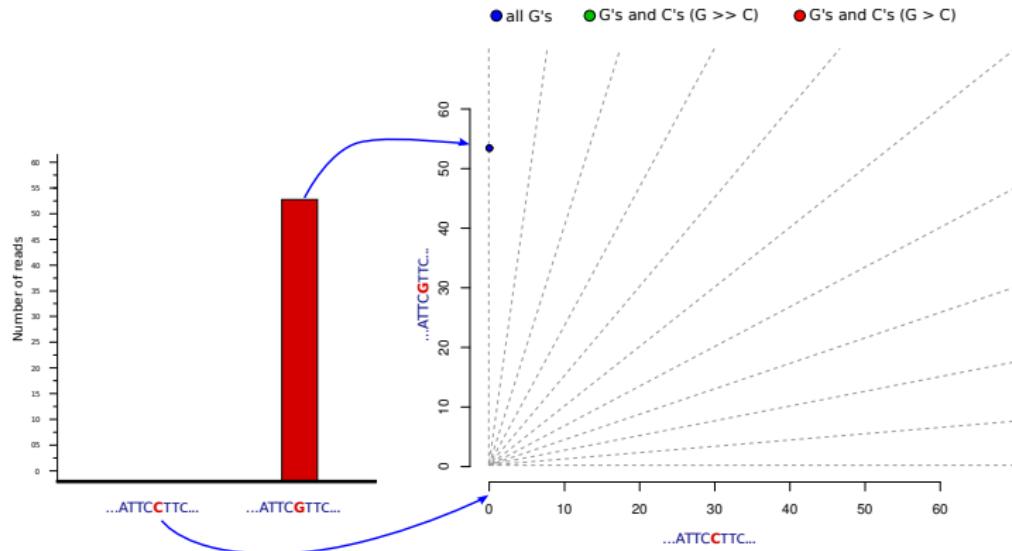
# GENOTYPING IN POLYPLOIDS

## BIALLELIC CO-DOMINANT MARKER



**Quantitative genotyping:** access the abundance of each one of the allelic forms (in a biallelic marker)

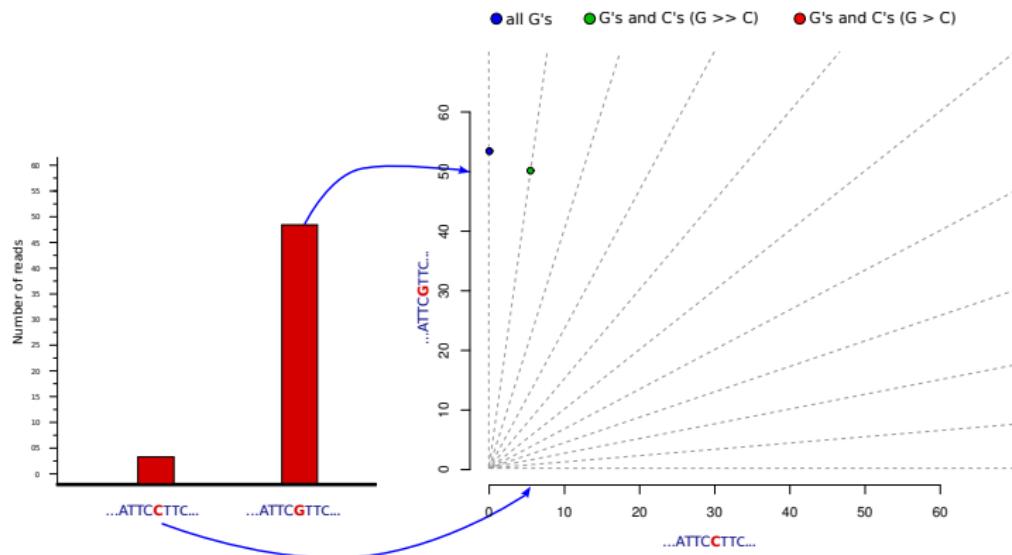
# GBS SCATTER PLOT



IMPORTANT

- ▶ Ratios of peaks or areas (angle): dosage

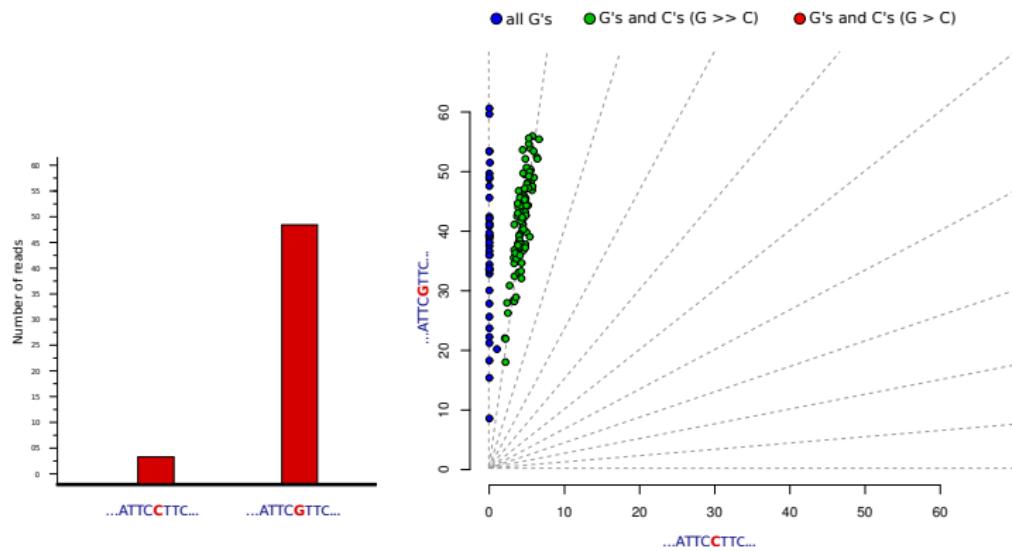
# GBS SCATTER PLOT



IMPORTANT

- ▶ Ratios of peaks or areas (angle): dosage

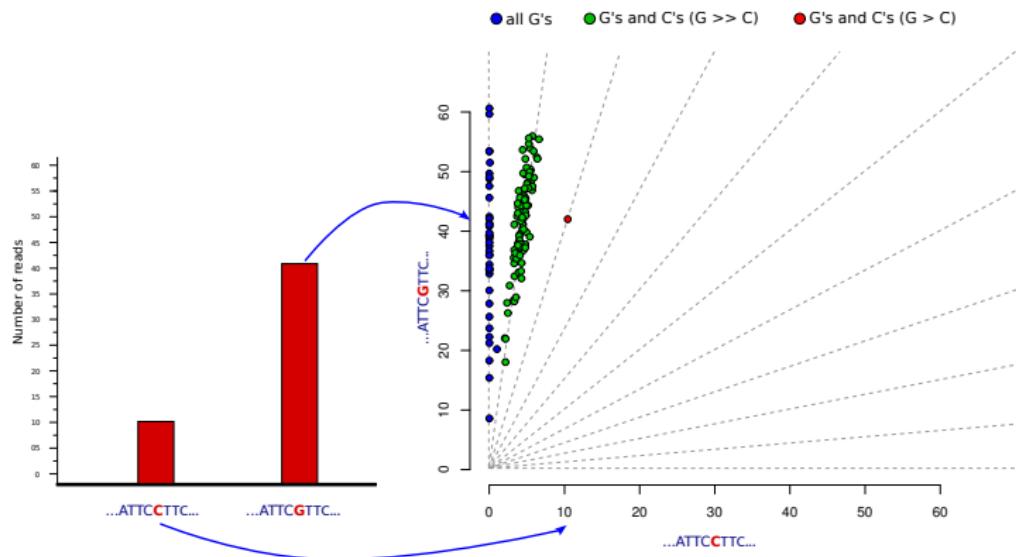
# GBS SCATTER PLOT



IMPORTANT

- ▶ Ratios of peaks or areas (angle): dosage

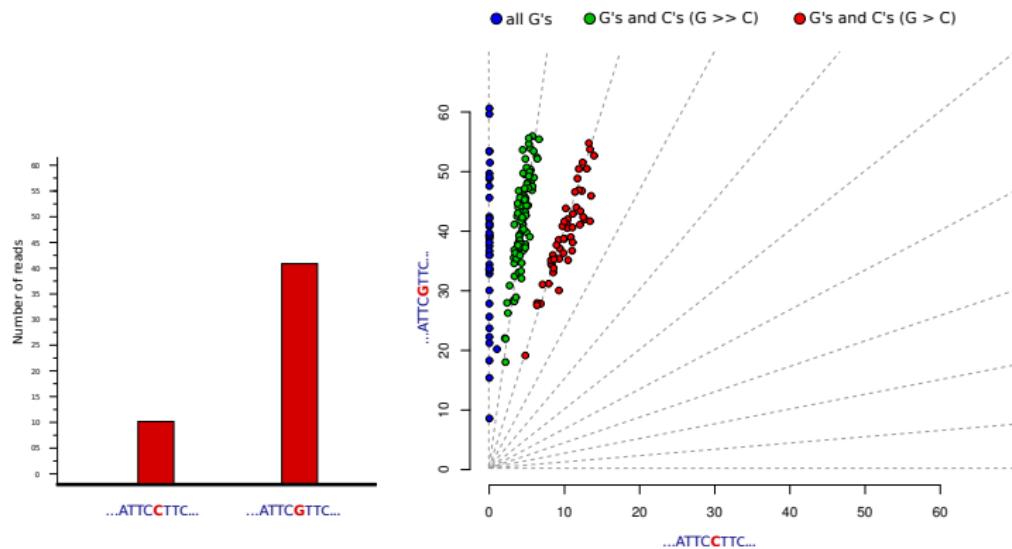
# GBS SCATTER PLOT



IMPORTANT

- ▶ Ratios of peaks or areas (angle): dosage

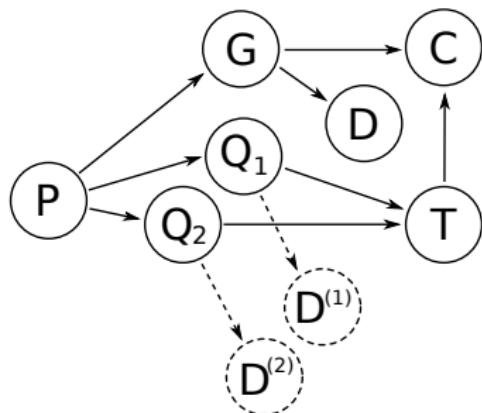
# GBS SCATTER PLOT



## IMPORTANT

- ▶ Ratios of peaks or areas (angle): dosage

# PROBABILISTIC GRAPHICAL MODEL - F<sub>1</sub>



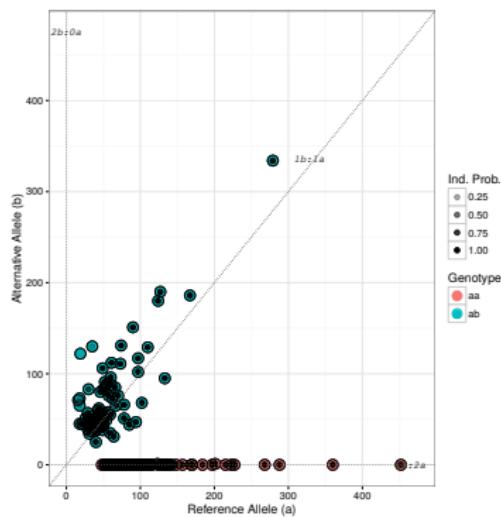
Serang et al. (2012)  
Molinari and Serang (2015)

$$\Pr(P, G, D, T, C, Q_1, Q_2, D_1, D_2)$$

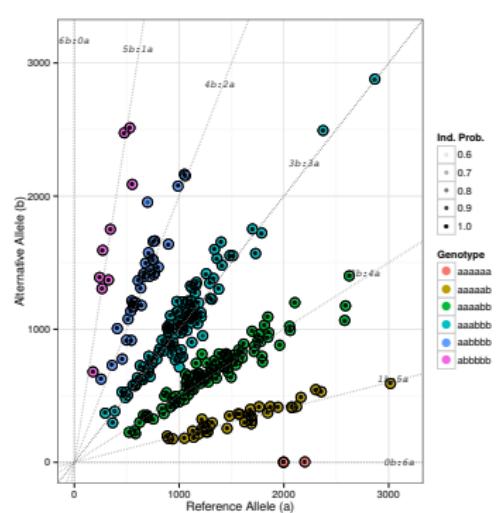
- ▶  $P$ : ploidy
- ▶  $G$ : genotype of all individuals
- ▶  $D$ : observed data
- ▶  $T$ : theoretical distribution of genotypes
- ▶  $C$ : histogram of genotypes
- ▶  $Q_1$  and  $Q_2$ : parent genotypes, with data  $D_1$  and  $D_2$  (if available)

# EXAMPLE: SNPs IN SWEET POTATO

Diploid



Hexaploid



## GENOTYPE CALLING IN POLYPLOIDS

- ▶ fitTetra (tetraploids):

<https://www.wur.nl/en/show/Software-fitTetra.htm>

- ▶ SuperMASSA (any ploidy level):

<https://bitbucket.org/orserang/supermassa>

- ▶ updog (any ploidy level, allows preferential pairing):

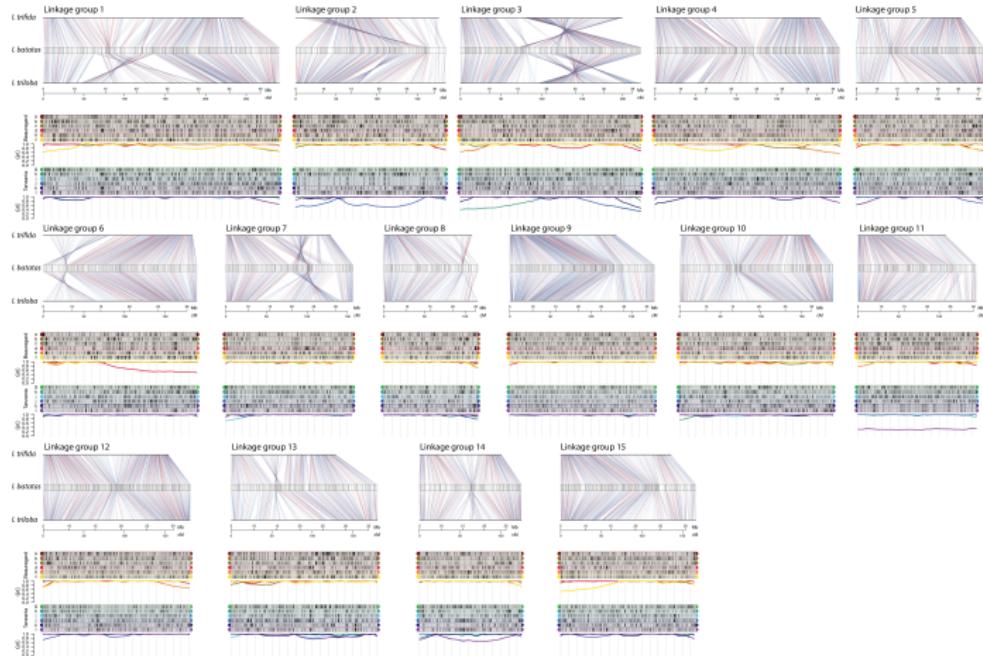
<https://github.com/dcgerard/updog>

- ▶ polyRAD (any ploidy level, reads VCF, BAM, etc):

<https://github.com/lvclark/polyRAD>

# LINKAGE ANALYSIS

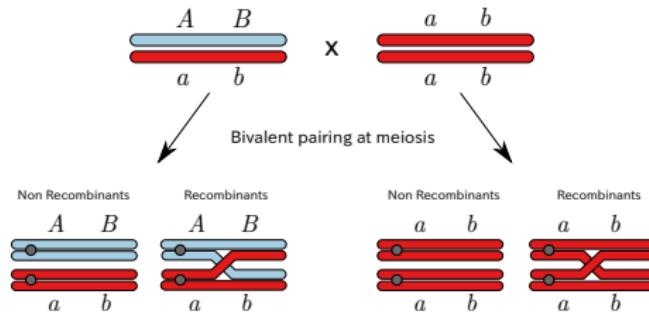
Statement about the patterns of co-segregation of alleles or chromosomal segments. In outcrossing species, we do not know how the alleles are positioned in the parental chromosomes, thus we also need to infer the linkage phase of these alleles.



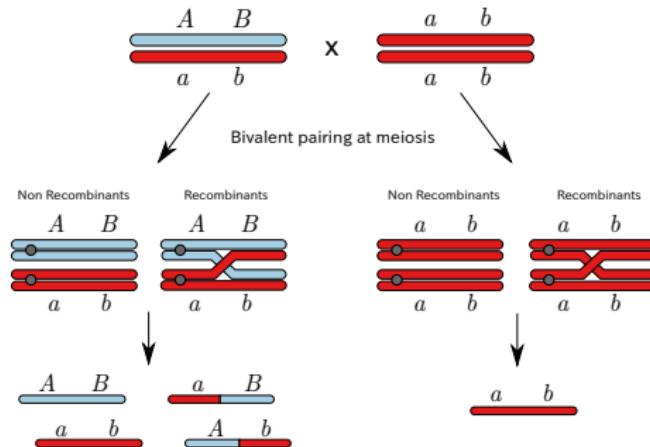
# RECOMBINATION FRACTION IN DIPLOIDS



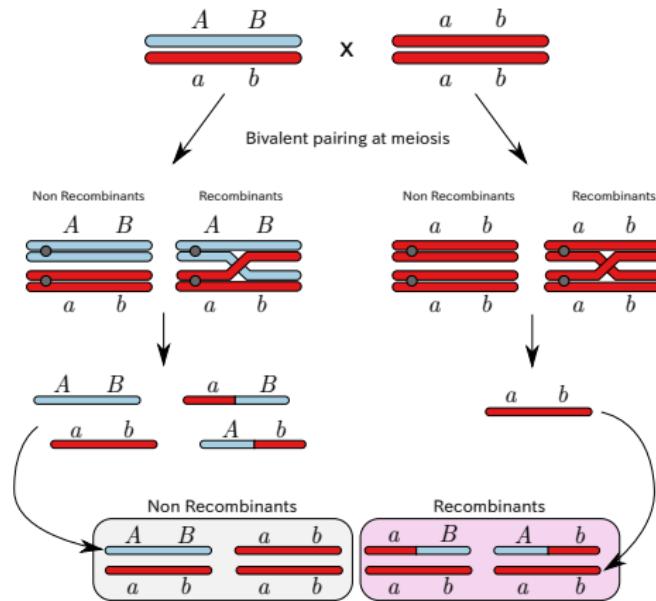
# RECOMBINATION FRACTION IN DIPLOIDS



# RECOMBINATION FRACTION IN DIPLOIDS



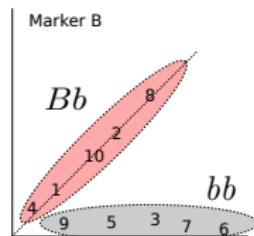
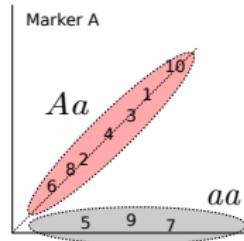
# RECOMBINATION FRACTION IN DIPLOIDS



$$\hat{r} = \frac{\#\text{recombinants}}{\#\text{total}}$$

# RECOMBINATION FRACTION IN DIPLOIDS

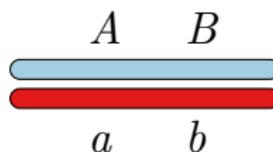
## EXAMPLE



Individual	Obs. Gen.
1	(A ,B)
2	(A ,B)
3	(A ,b )
4	(A ,B)
5	(a ,b )
6	(A ,b )
7	(a ,b )
8	(A ,B)
9	(a ,b )
10	(A ,B)

$$\hat{r} = 2/10 = 0.2$$

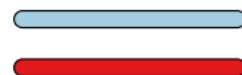
# EXPECTED GAMETIC FREQUENCY



$r$   
Recombinants



$1-r$   
Non Recombinants



$$\mathbf{V} = \begin{bmatrix} \Pr(A, B) & \Pr(a, B) \\ \Pr(A, b) & \Pr(a, b) \end{bmatrix} = \begin{bmatrix} \frac{1-r}{2} & \frac{r}{2} \\ \frac{r}{2} & \frac{1-r}{2} \end{bmatrix}$$

# RECOMBINATION FRACTION IN DIPLOIDS

## EXAMPLE (CONT.)

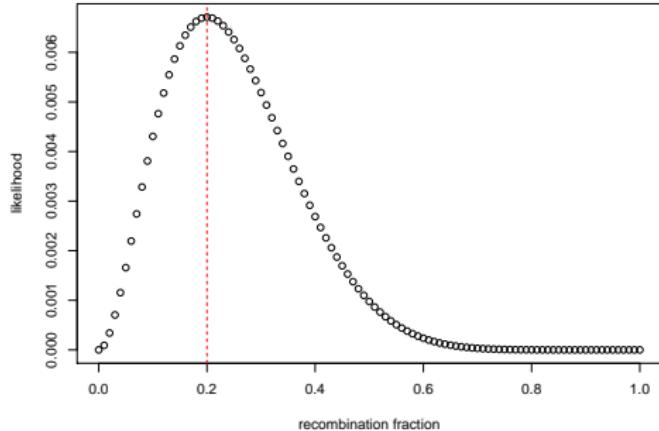


$$L = \prod_n \Pr(\text{loc}_B, \text{loc}_A \mid \text{Data})$$

Individual	Obs. Gen.	$\Pr(\text{loc}_B, \text{loc}_A)$
1	(A ,B)	$\frac{1}{2}(1-r)$
2	(A ,B)	$\frac{1}{2}(1-r)$
3	(A ,b )	$\frac{1}{2} (r)$
4	(A ,B)	$\frac{1}{2}(1-r)$
5	(a ,b )	$\frac{1}{2}(1-r)$
6	(A ,b )	$\frac{1}{2} (r)$
7	(a ,b )	$\frac{1}{2}(1-r)$
8	(A ,B)	$\frac{1}{2}(1-r)$
9	(a ,b )	$\frac{1}{2}(1-r)$
10	(A ,B)	$\frac{1}{2}(1-r)$

$$L = \left(\frac{r}{2}\right)^2 \left(\frac{1-r}{2}\right)^8$$

# LIKELIHOOD PROFILE

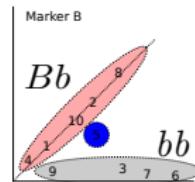
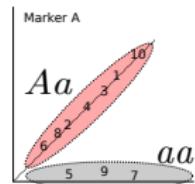
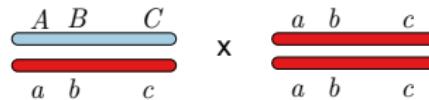


$$L = \prod_n \Pr(G_A, G_B) = \left(\frac{r}{2}\right)^2 \left(\frac{1-r}{2}\right)^8$$

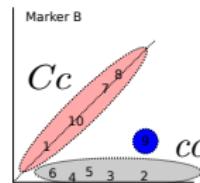
The MLE (maximum likelihood estimate) of  $r$  is  $\hat{r} = 0.2$

# RECOMBINATION FRACTION IN DIPLOIDS

THREE MARKERS, ELIMINATING NOISY DATA POINTS



Individual	Obs. Gen.
1	(A ,B,C )
2	(A ,B,c )
3	(A ,b ,c )
4	(A ,B,c )
5	(a ,?,c )
6	(A ,b ,c )
7	(a ,b ,C )
8	(A ,B,C )
9	(a ,b ,?)
10	(A ,B,C )



$$\hat{r}_{AB} = 2/9 = 0.222$$

$$\hat{r}_{BC} = 3/8 = 0.375$$

# HIDDEN MARKOV MODEL

## INTUITION

DO YUO FNID TIHS SMILPE TO RAED?

Bceause of the phaonmneal  
pweor of the hmuani mind,  
msot pleope do. Aoccdrnig  
to rscheearch at Cmabrigde  
Uinervtisy, it dseno't mtaetr  
what oerdr the ltteres are  
in. The olny iproamtnt thing  
is that the frsit and last  
ltteer be in the rghit pclae.  
The rset can be a taotl mses  
and you can still raed it  
whotuit a pboerlm. This is  
bcuseae the huamn mnid  
deos not raed ervey lteter  
by istlef, but the word as a  
wlohe. Takl abuot cool.

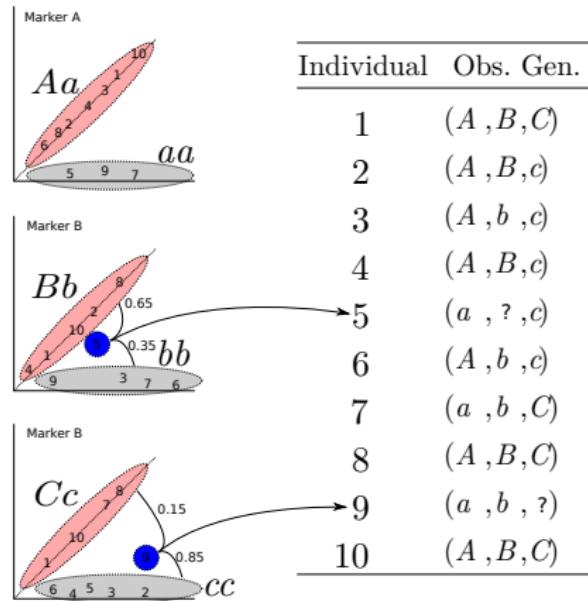
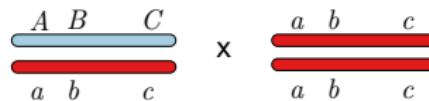
Find out why you see the world the  
way you do. You won't believe your mind.

**BRAIN GAMES**



# RECOMBINATION FRACTION IN DIPLOIDS

THREE MARKERS, HMM TO HANDLE NOISY DATA POINTS



# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



Aa

aa

# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



Aa

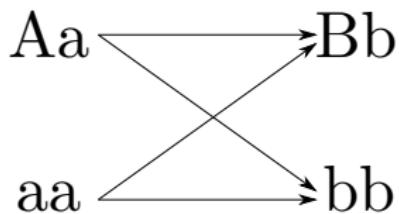
Bb

aa

bb

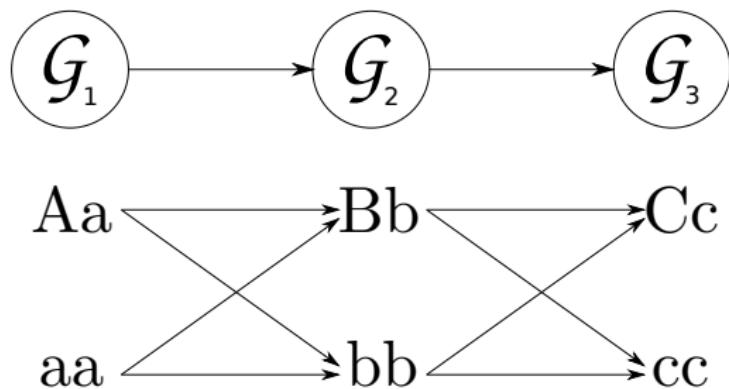
# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



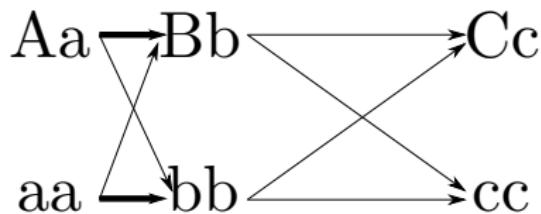
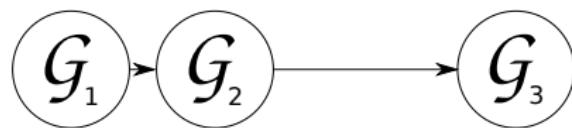
# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



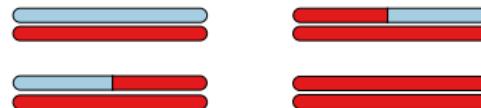
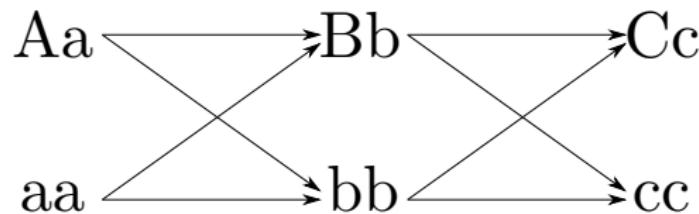
# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



# HIDDEN MARKOV MODELS (HMM)

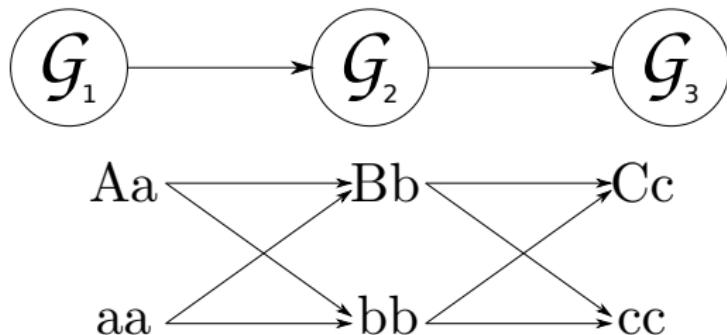
## TRANSITION PROBABILITIES



$$\mathbf{V} = \begin{bmatrix} \Pr(\mathcal{G}_2=Bb | \mathcal{G}_1=Aa) & \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=Aa) \\ \Pr(\mathcal{G}_2=Bb | \mathcal{G}_1=aa) & \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=aa) \end{bmatrix} = \begin{bmatrix} 1-r & r \\ r & 1-r \end{bmatrix}$$

# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES

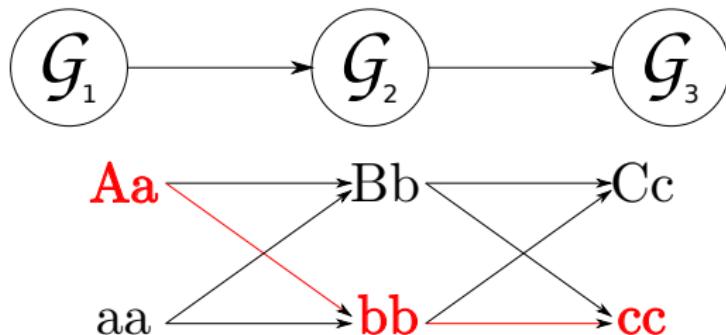


$$\begin{bmatrix} \Pr(\mathcal{G}_2=Bb | \mathcal{G}_1=Aa) & \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=Aa) \\ \Pr(\mathcal{G}_2=Bb | \mathcal{G}_1=aa) & \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=aa) \end{bmatrix} = \begin{bmatrix} 1-r_{AB} & r_{AB} \\ r_{AB} & 1-r_{AB} \end{bmatrix}$$

$$\begin{bmatrix} \Pr(\mathcal{G}_3=Cc | \mathcal{G}_2=Bb) & \Pr(\mathcal{G}_3=cc | \mathcal{G}_2=Bb) \\ \Pr(\mathcal{G}_3=Cc | \mathcal{G}_2=bb) & \Pr(\mathcal{G}_3=cc | \mathcal{G}_2=bb) \end{bmatrix} = \begin{bmatrix} 1-r_{BC} & r_{BC} \\ r_{BC} & 1-r_{BC} \end{bmatrix}$$

# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



$$\begin{bmatrix} \Pr(\mathcal{G}_2=Bb | \mathcal{G}_1=Aa) & \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=Aa) \\ \Pr(\mathcal{G}_2=Bb | \mathcal{G}_1=aa) & \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=aa) \end{bmatrix} = \begin{bmatrix} 1-r_{AB} & r_{AB} \\ r_{AB} & 1-r_{AB} \end{bmatrix}$$

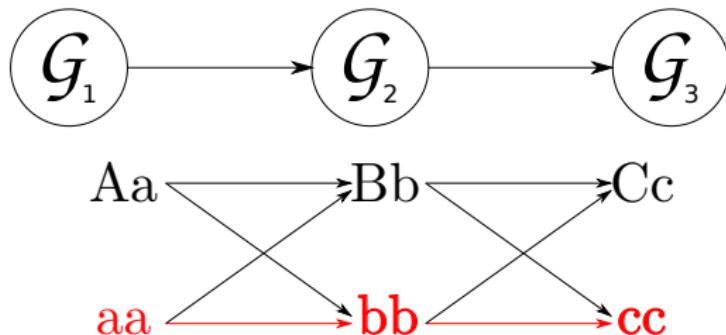
$$\begin{bmatrix} \Pr(\mathcal{G}_3=Cc | \mathcal{G}_2=Bb) & \Pr(\mathcal{G}_3=cc | \mathcal{G}_2=Bb) \\ \Pr(\mathcal{G}_3=Cc | \mathcal{G}_2=bb) & \Pr(\mathcal{G}_3=cc | \mathcal{G}_2=bb) \end{bmatrix} = \begin{bmatrix} 1-r_{BC} & r_{BC} \\ r_{BC} & 1-r_{BC} \end{bmatrix}$$

Observed Genotype:    Aa        bb        cc

$$\Pr(Aa, bb, cc) = \Pr(\mathcal{G}_1=Aa) \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=Aa) \Pr(\mathcal{G}_3=cc | \mathcal{G}_2=bb)$$

# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



$$\begin{bmatrix} \Pr(\mathcal{G}_2=Bb|\mathcal{G}_1=Aa) & \Pr(\mathcal{G}_2=bb|\mathcal{G}_1=Aa) \\ \Pr(\mathcal{G}_2=Bb|\mathcal{G}_1=aa) & \Pr(\mathcal{G}_2=bb|\mathcal{G}_1=aa) \end{bmatrix} = \begin{bmatrix} 1-r_{AB} & r_{AB} \\ r_{AB} & 1-r_{AB} \end{bmatrix}$$

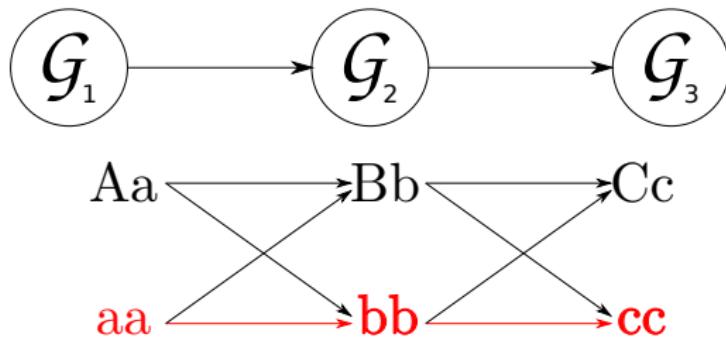
$$\begin{bmatrix} \Pr(\mathcal{G}_3=Cc|\mathcal{G}_2=Bb) & \Pr(\mathcal{G}_3=cc|\mathcal{G}_2=Bb) \\ \Pr(\mathcal{G}_3=Cc|\mathcal{G}_2=bb) & \Pr(\mathcal{G}_3=cc|\mathcal{G}_2=bb) \end{bmatrix} = \begin{bmatrix} 1-r_{BC} & r_{BC} \\ r_{BC} & 1-r_{BC} \end{bmatrix}$$

Observed Genotype: aa      bb      cc

$$\Pr(aa, bb, cc) = \Pr(\mathcal{G}_1=aa) \Pr(\mathcal{G}_2=bb|\mathcal{G}_1=aa) \Pr(\mathcal{G}_3=cc|\mathcal{G}_2=bb)$$

# HIDDEN MARKOV MODELS (HMM)

## TRANSITION PROBABILITIES



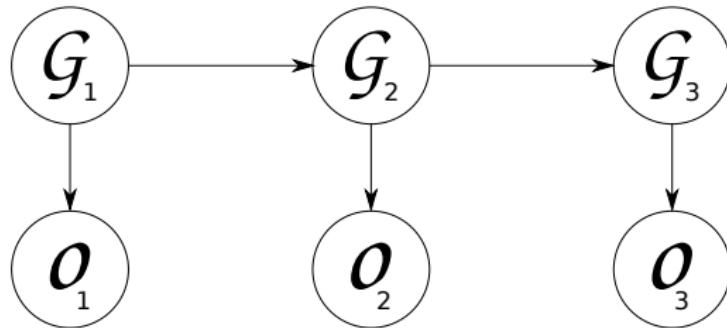
$$\begin{bmatrix} \Pr(\mathcal{G}_2=Bb | \mathcal{G}_1=Aa) & \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=Aa) \\ \Pr(\mathcal{G}_2=Bb | \mathcal{G}_1=aa) & \Pr(\mathcal{G}_2=bb | \mathcal{G}_1=aa) \end{bmatrix} = \begin{bmatrix} 1-r_{AB} & r_{AB} \\ r_{AB} & 1-\cancel{r}_{AB} \end{bmatrix}$$

$$\begin{bmatrix} \Pr(\mathcal{G}_3=Cc | \mathcal{G}_2=Bb) & \Pr(\mathcal{G}_3=cc | \mathcal{G}_2=Bb) \\ \Pr(\mathcal{G}_3=Cc | \mathcal{G}_2=bb) & \Pr(\mathcal{G}_3=cc | \mathcal{G}_2=bb) \end{bmatrix} = \begin{bmatrix} 1-r_{BC} & \cancel{r}_{BC} \\ r_{BC} & 1-\cancel{r}_{BC} \end{bmatrix}$$

$$L = \prod_n \Pr(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3 | Data)$$

# HIDDEN MARKOV MODELS (HMM)

## EMISSION PROBABILITIES

 $\Pr(O_1|G_1)$ 

		$G_1$	
		Aa	aa
Aa	1	0	
	—	1/2 or p <sub>1</sub>	
aa	0	1	
	—	1/2 or 1-p <sub>1</sub>	

 $\Pr(O_2|G_2)$ 

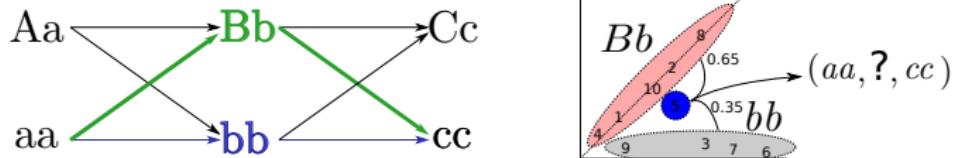
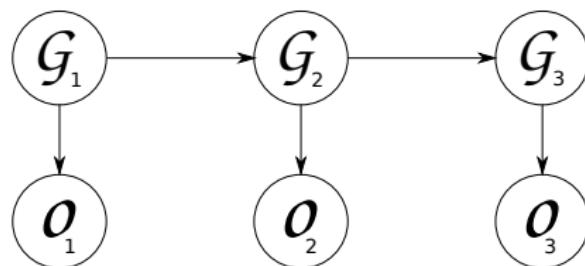
		$G_2$	
		Bb	bb
Bb	1	0	
	—	1/2 or p <sub>2</sub>	
bb	0	1	
	—	1/2 or 1-p <sub>2</sub>	

 $\Pr(O_3|G_3)$ 

		$G_3$	
		Cc	cc
Cc	1	0	
	—	1/2 or p <sub>3</sub>	
cc	0	1	
	—	1/2 or 1-p <sub>3</sub>	

# HIDDEN MARKOV MODELS (HMM)

## EMISSION PROBABILITIES

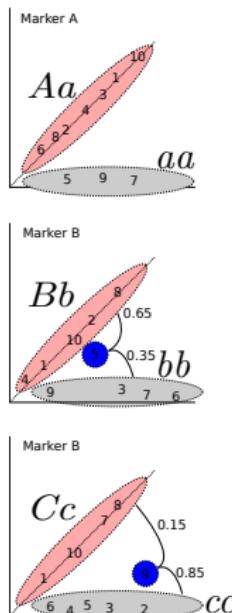


$$\Pr(\text{aa})\Pr(\text{Bb}|\text{aa})\Pr(\text{O}|\text{Bb})\Pr(\text{cc}|\text{Bb}) = 1/2 \cdot (r_{ab}) \cdot (1-r_{bc}) \cdot 0.65 \times$$

$$\Pr(\text{aa})\Pr(\text{Bb}|\text{aa})\Pr(\text{O}|\text{bb})\Pr(\text{cc}|\text{bb}) = 1/2 \cdot (1-r_{ab}) \cdot (1-r_{bc}) \cdot 0.35$$

# HIDDEN MARKOV MODELS (HMM)

## LIKELIHOOD



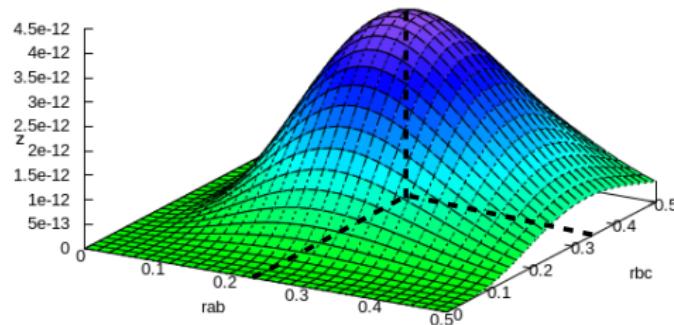
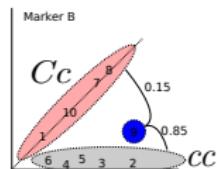
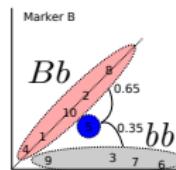
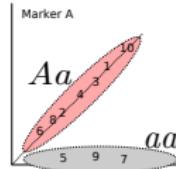
Ind.	Obs.	Gen.	Pr(G <sub>1</sub> , G <sub>2</sub> , G <sub>3</sub>   Data)
1	(A ,B,C)	Pr(Aa)Pr(Bb Aa)Pr(Cc Bb) = 1/2 . (1-r <sub>ab</sub> ) .(1-r <sub>bcc</sub> )	
2	(A ,B,c)	Pr(Aa)Pr(Bb Aa)Pr(cc Bb) = 1/2 . (1-r <sub>ab</sub> ) .(r <sub>bcc</sub> )	
3	(A ,b ,c)	Pr(Aa)Pr(bb Aa)Pr(cc bb) = 1/2 . (r <sub>ab</sub> ) .(1-r <sub>bcc</sub> )	
4	(A ,B,c)	Pr(Aa)Pr(Bb Aa)Pr(cc Bb) = 1/2 . (1-r <sub>ab</sub> ) .(r <sub>bcc</sub> )	
5	(a ,? ,c)	Pr(aa)Pr(Bb aa).Pr(O Bb).Pr(cc bb) = 1/2 . (r <sub>ab</sub> ) .(1-r <sub>bcc</sub> ).0.65 Pr(aa)Pr(bb aa).Pr(O bb).Pr(cc bb) = 1/2 . (1-r <sub>ab</sub> ) .(1-r <sub>bcc</sub> ).0.35	
6	(A ,b ,c)	Pr(Aa)Pr(bb Aa)Pr(cc bb) = 1/2 . (r <sub>ab</sub> ) .(1-r <sub>bcc</sub> )	
7	(a ,b ,C)	Pr(aa)Pr(bb aa)Pr(Cc bb) = 1/2 . (1-r <sub>ab</sub> ) .(r <sub>bcc</sub> )	
8	(A ,B,C)	Pr(Aa)Pr(Bb Aa)Pr(Cc Bb) = 1/2 . (1-r <sub>ab</sub> ) .(1-r <sub>bcc</sub> )	
9	(a ,b ,?)	Pr(aa)Pr(bb aa)Pr(cc bb).Pr(O cc) = 1/2 . (1-r <sub>ab</sub> ) .(1-r <sub>bcc</sub> ).0.85 Pr(aa)Pr(bb aa)Pr(Cc bb).Pr(O Cc) = 1/2 . (1-r <sub>ab</sub> ) .(r <sub>bcc</sub> ).0.15	
10	(A ,B,C)	Pr(Aa)Pr(Bb Aa)Pr(Cc Bb) = 1/2 . (1-r <sub>ab</sub> ) .(1-r <sub>bcc</sub> )	

$$L = \prod_n \Pr(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3 | Data) = 7.1 \cdot 10^{-6} (1-r_{ab})^9 r_{ab}^3 (1-r_{bcc})^8 r_{bcc}^4$$

# HIDDEN MARKOV MODELS (HMM)

## MAXIMUM LIKELIHOOD ESTIMATE

$$L = \prod_n \Pr(\mathcal{G}_1, \mathcal{G}_2, \mathcal{G}_3 | Data) = 7.1 \cdot 10^{-6} (1-rab)^9 rab^3 (1-rbc)^8 rbc^4$$



HMM

$$\hat{r}_{AB} = 1/4 = 0.250$$

$$\hat{r}_{BC} = 1/3 = 0.333$$

two-point

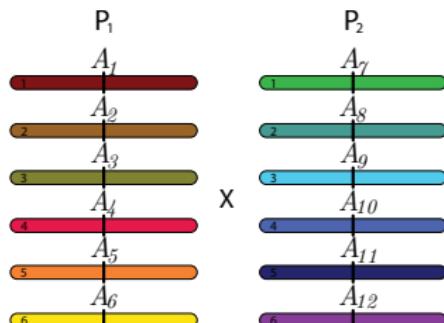
$$\hat{r}_{AB} = 2/9 = 0.222$$

$$\hat{r}_{BC} = 3/8 = 0.375$$

## Linkage analysis in autopolyploids

- ▶ Linkage Analysis and Haplotype Phasing in Experimental Autopolyploid Populations with High Ploidy Level Using Hidden Markov Models (<https://doi.org/10.1534/g3.119.400378>)
- ▶ Unraveling the hexaploid sweetpotato inheritance using ultra-dense multilocus mapping (<https://doi.org/10.1101/689638>)

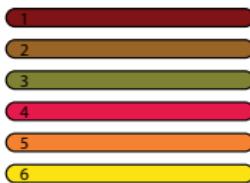
# HEXAPLOID BI-PARENTAL CROSS



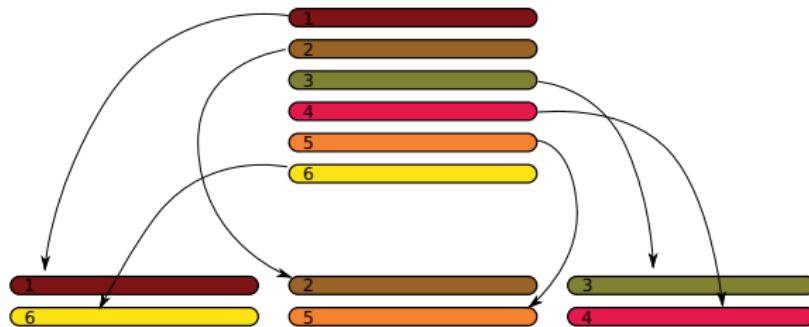
		Gametes						
		1 $A_1A_2A_3$	2 $A_1A_2A_4$	3 $A_1A_2A_5$	...	18 $A_3A_4A_6$	19 $A_3A_5A_6$	20 $A_4A_5A_6$
$P_2 \setminus P_1$					...			
1	$A_7A_8A_9$				...			
2	$A_7A_8A_{10}$				...			
3	$A_7A_8A_{11}$				...			
...	...	...	...	...	...	...	...	...
18	$A_gA_{10}A_{12}$				...			
19	$A_gA_{11}A_{12}$				...			
20	$A_{10}A_{11}A_{12}$				...			

400 possible genotypes

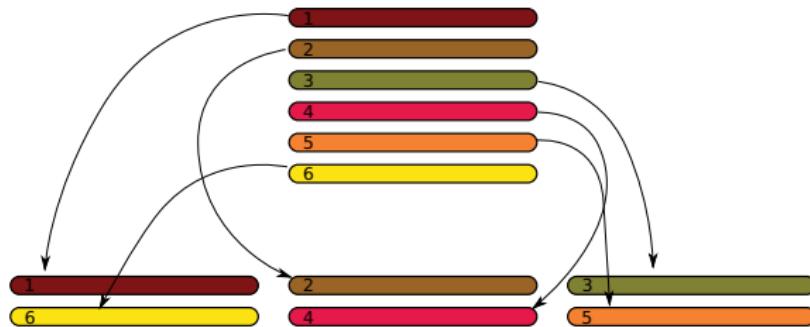
# GAMETE FORMATION - RANDOM PAIRING



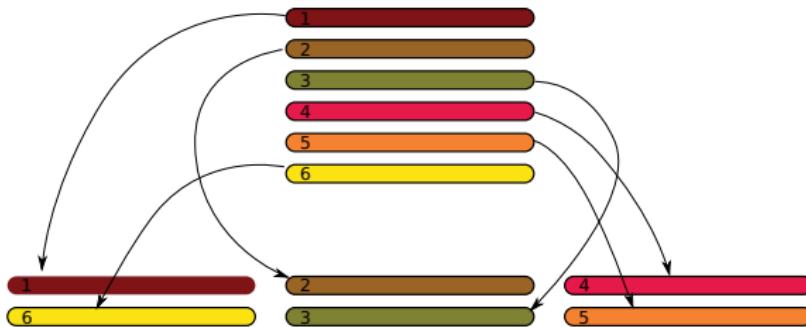
# GAMETE FORMATION - RANDOM PAIRING

 $\psi_1$

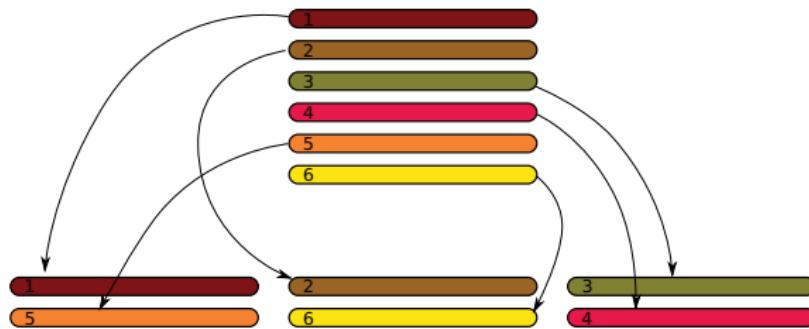
# GAMETE FORMATION - RANDOM PAIRING


$$\psi_2$$

# GAMETE FORMATION - RANDOM PAIRING

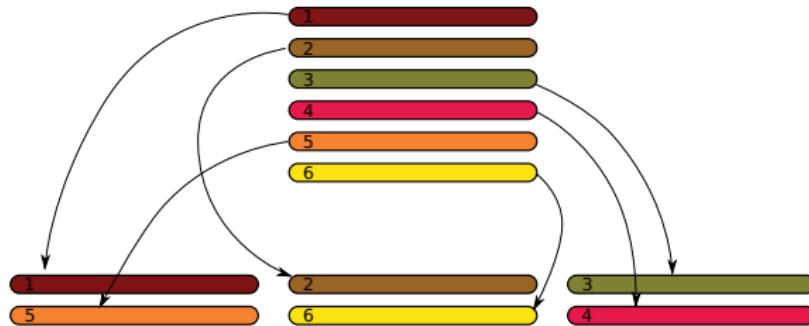
 $\psi_3$

# GAMETE FORMATION - RANDOM PAIRING



$$\psi_4$$

## GAMETE FORMATION - RANDOM PAIRING

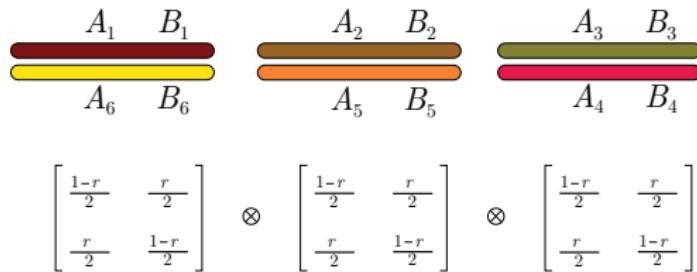


$$\psi_4$$

In this case: 15 possible configurations. For any ploidy level  $p$

$$\frac{1}{\frac{p}{2}!} \prod_{i=1}^{\frac{p}{2}} \binom{2i}{2}$$

# EXPECTED GAMETIC FREQUENCY FOR A BIVALENT CONFIGURATION



- ▶ In general:

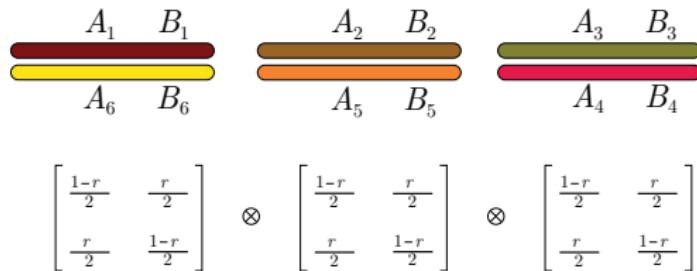
$$\mathbf{V}_1 \otimes \mathbf{V}_2 \otimes \cdots \otimes \mathbf{V}_{\frac{p}{2}}$$

- ▶ All elements of this product are of the form

$$\Pr(G_A, G_B \mid \psi_j) = \frac{(1-r)^{\frac{p}{2}-i}(r)^i}{2^{\frac{p}{2}}}$$

*i*: known number of recombinant bivalents between loci A and B

# EXPECTED GAMETIC FREQUENCY FOR A BIVALENT CONFIGURATION



- ▶ In general:

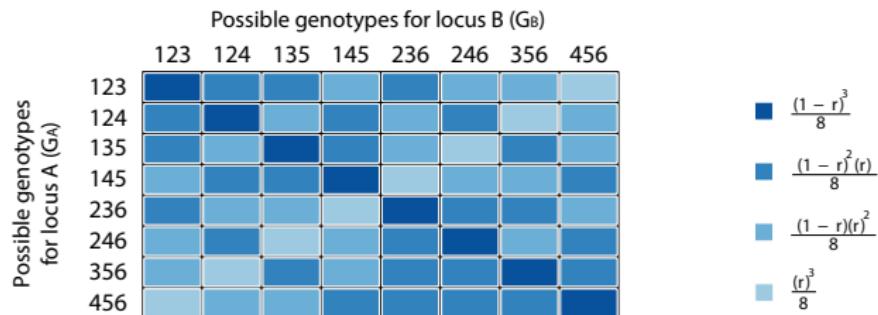
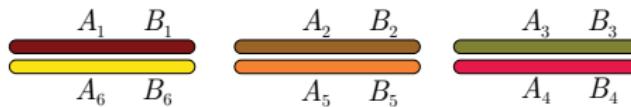
$$\mathbf{V}_1 \otimes \mathbf{V}_2 \otimes \cdots \otimes \mathbf{V}_{\frac{p}{2}}$$

- ▶ All elements of this product are of the form

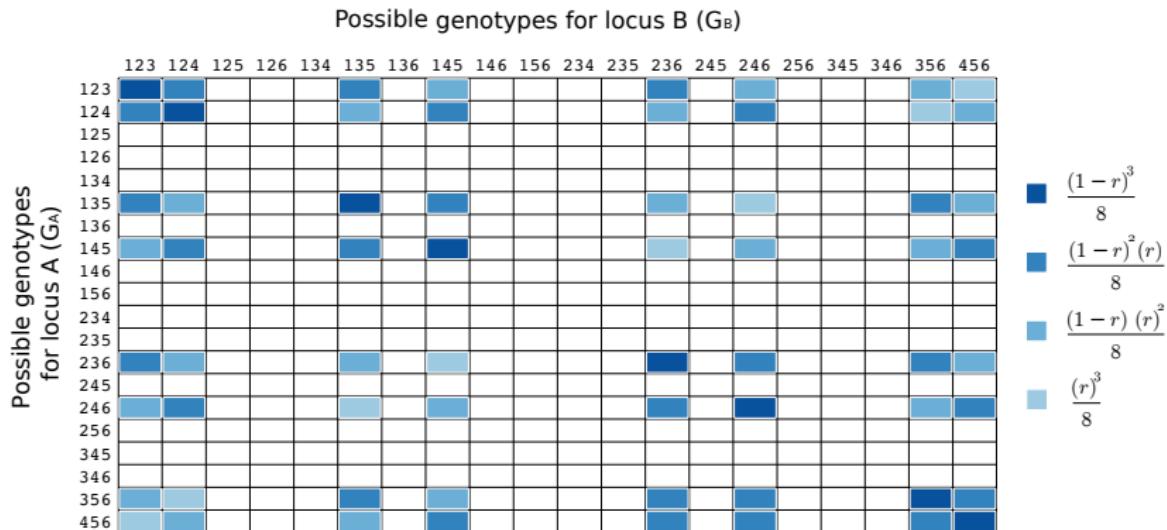
$$\Pr(G_A, G_B \mid \psi_j) = \frac{(1-r)^{\frac{p}{2}-i}(r)^i}{2^{\frac{p}{2}}}$$

$i$ : known number of recombinant bivalents between loci A and B

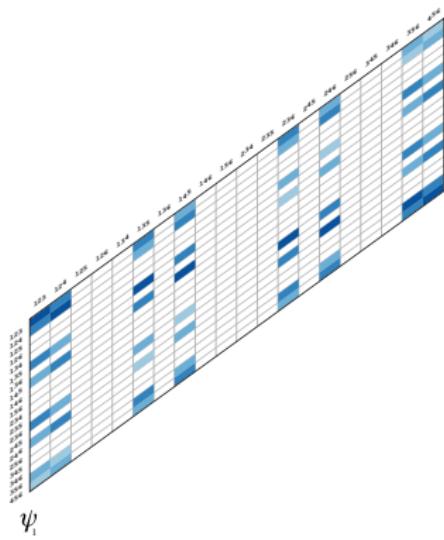
# GAMETIC PROBABILITY FOR $\psi_1$



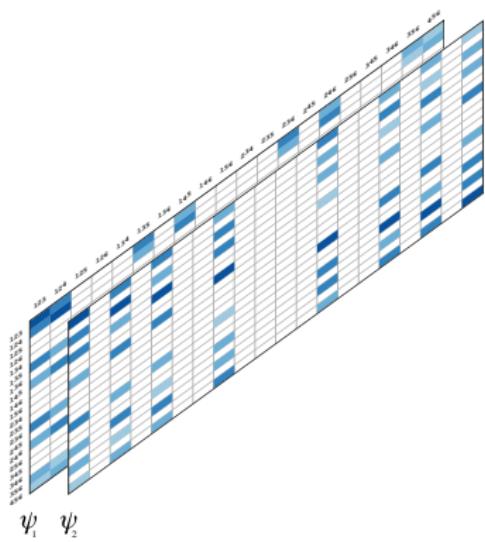
## GAMETIC PROBABILITY FOR $\psi_1$



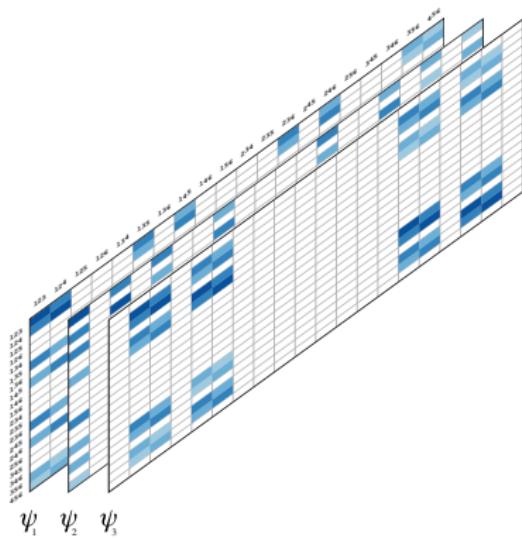
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



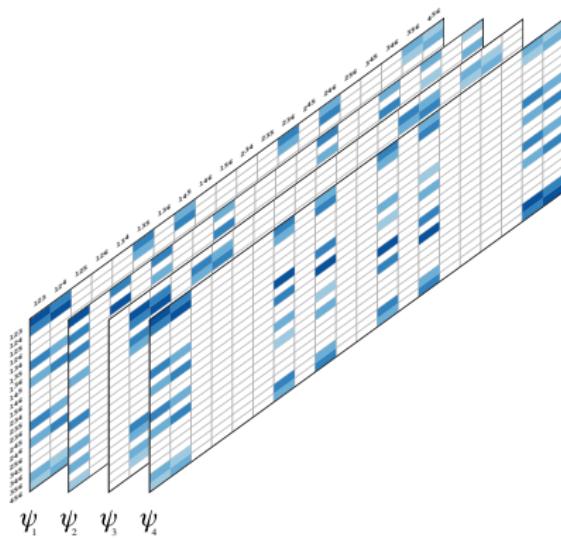
## GAMETIC PROBABILITY FOR ALL $\psi$ 's



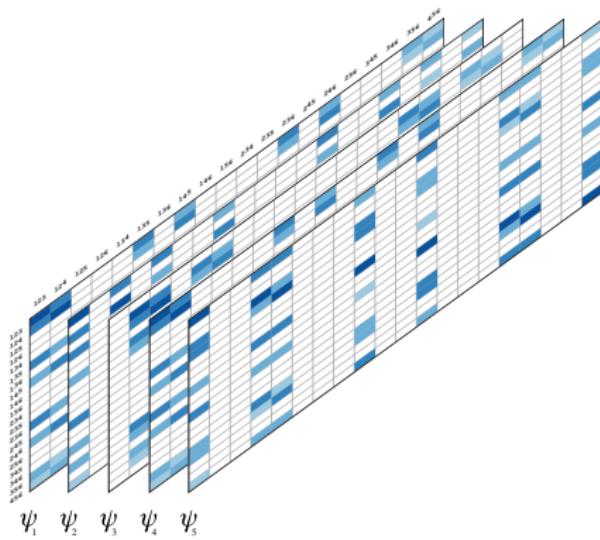
# GAMETIC PROBABILITY FOR ALL $\psi'$ s



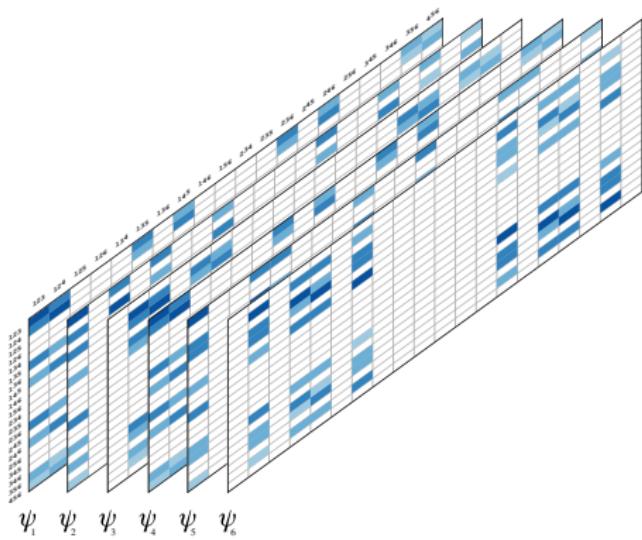
# GAMETIC PROBABILITY FOR ALL $\psi'$ 's



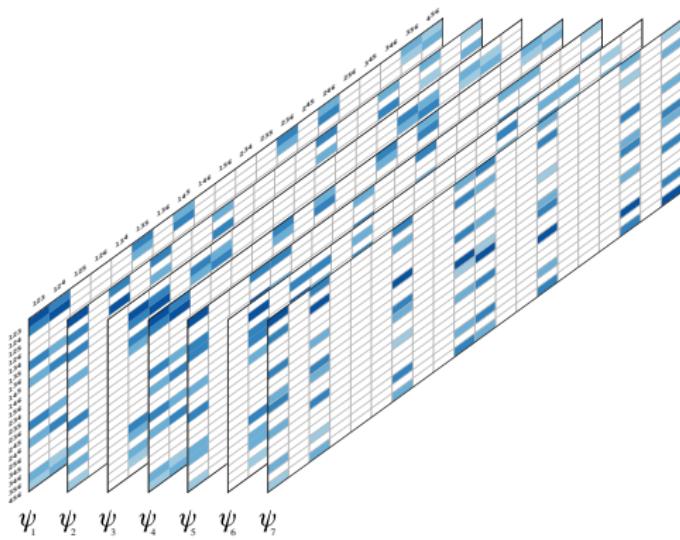
## GAMETIC PROBABILITY FOR ALL $\psi$ 's



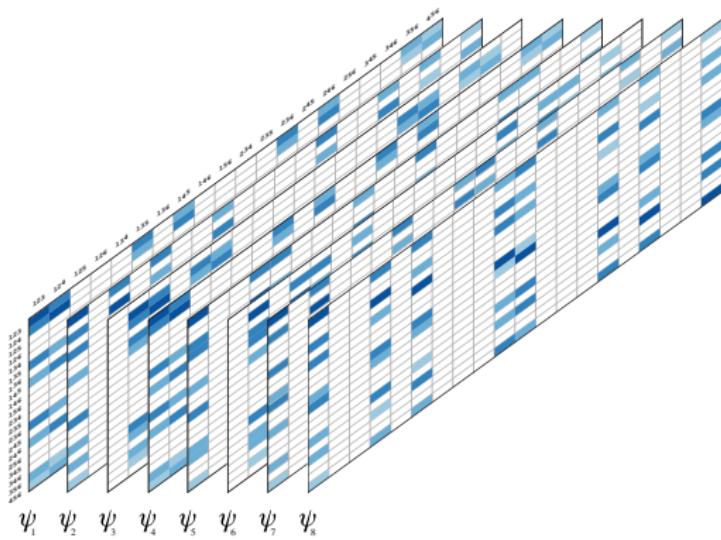
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



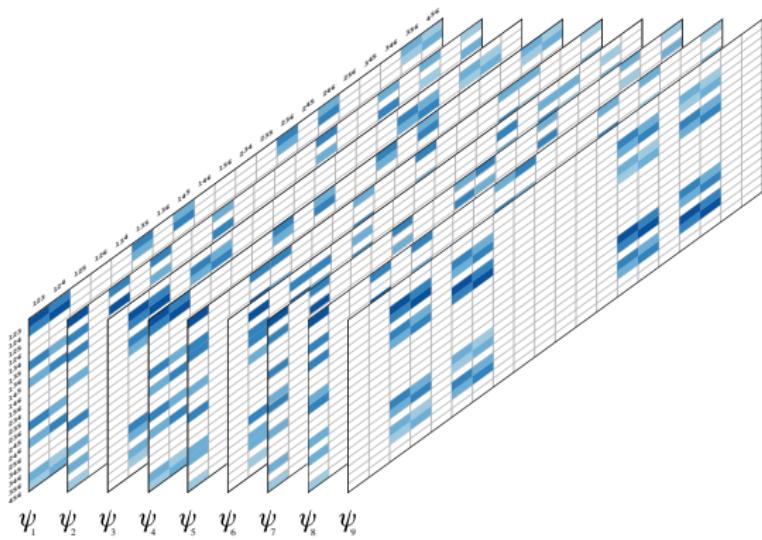
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



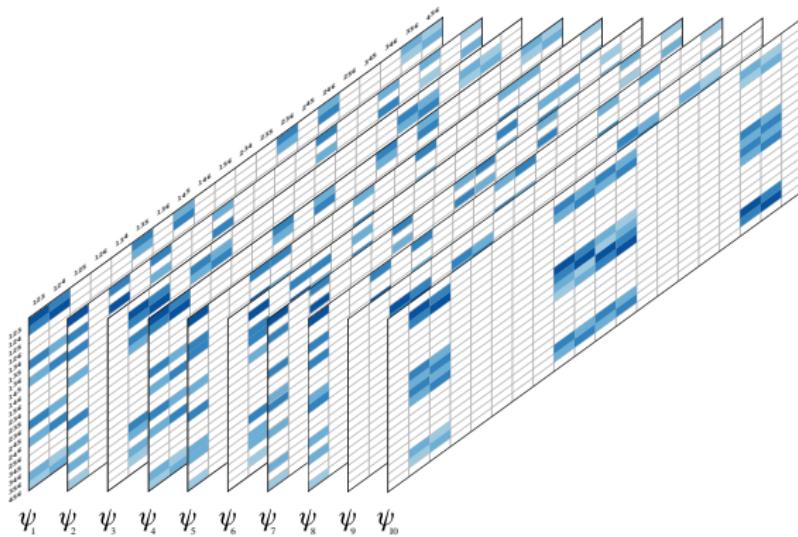
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



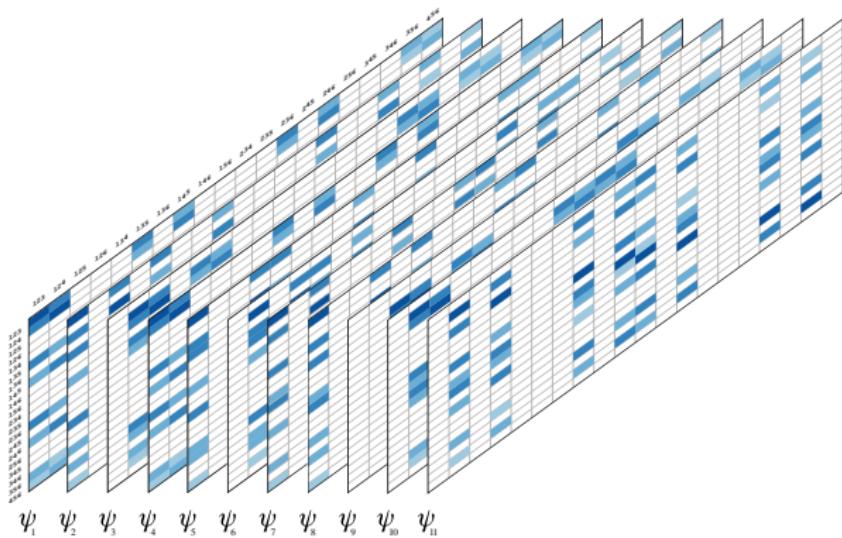
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



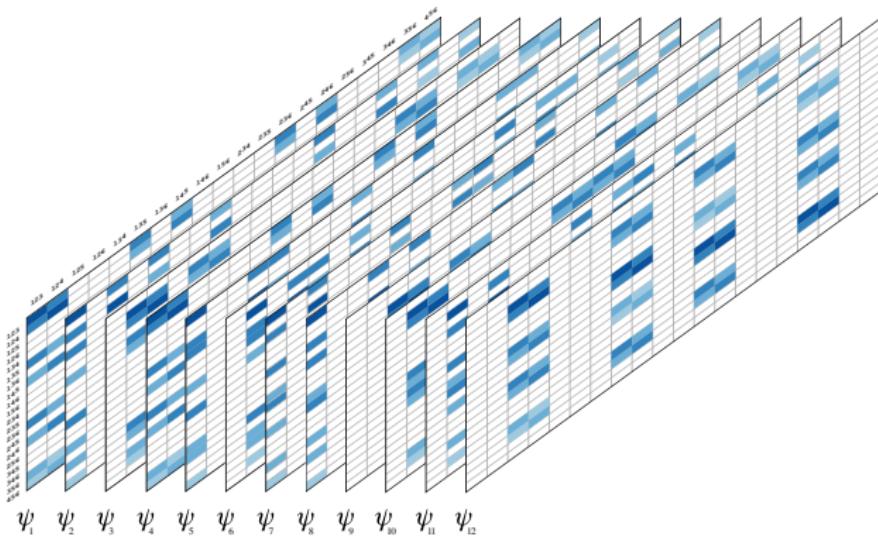
# GAMETIC PROBABILITY FOR ALL $\psi'$ 's



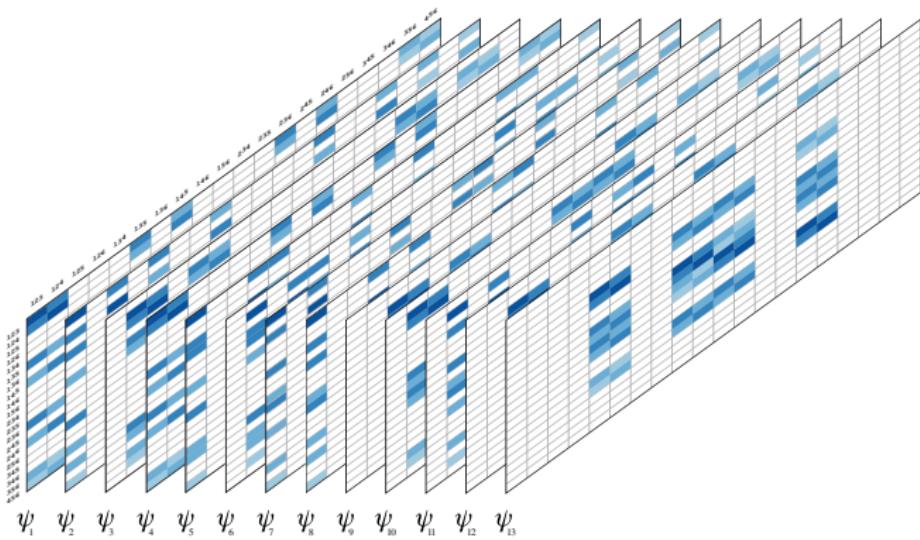
# GAMETIC PROBABILITY FOR ALL $\psi'$ s



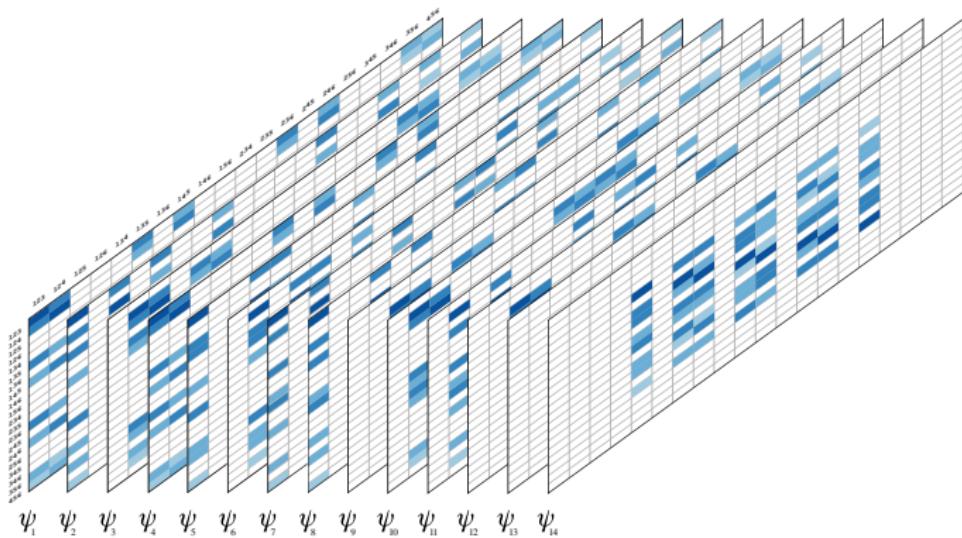
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



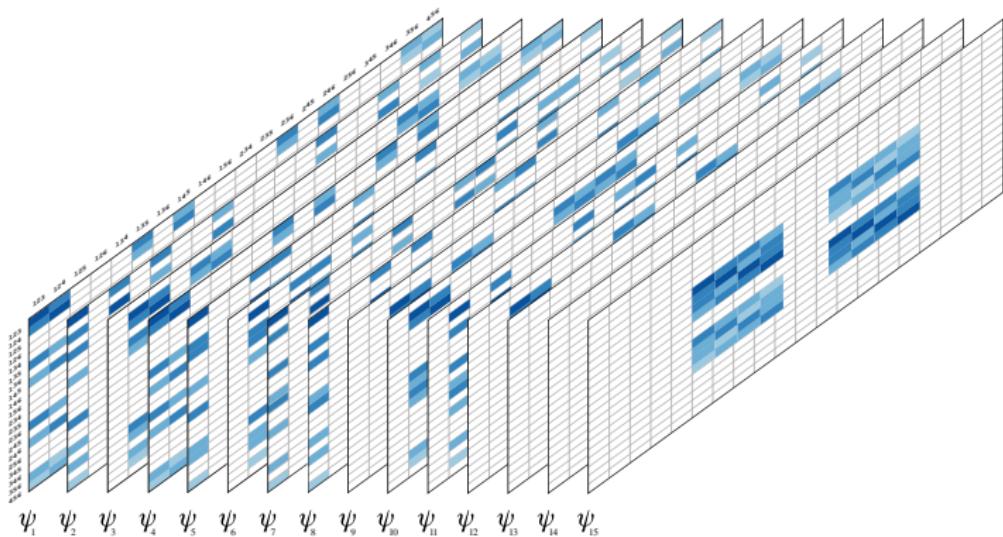
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



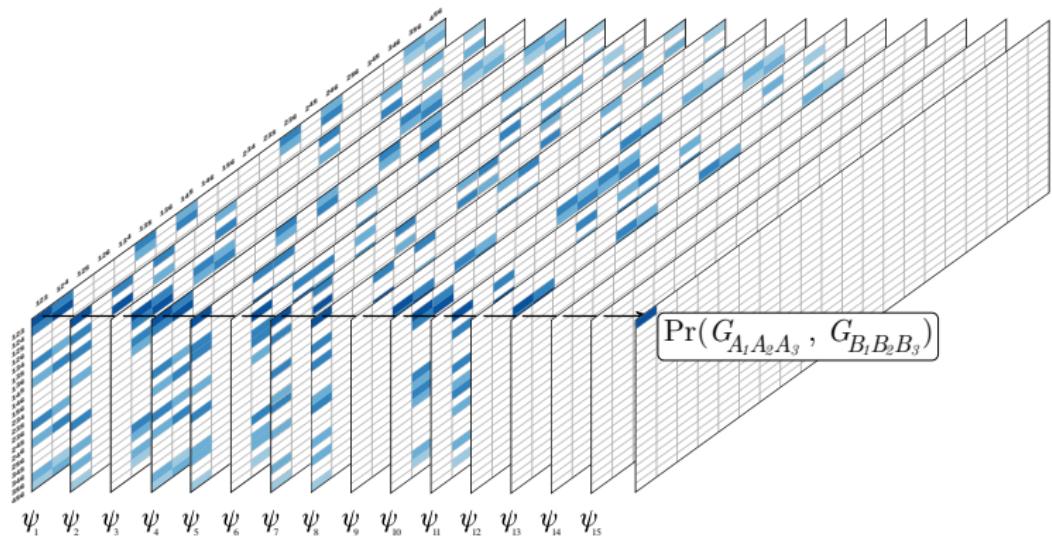
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



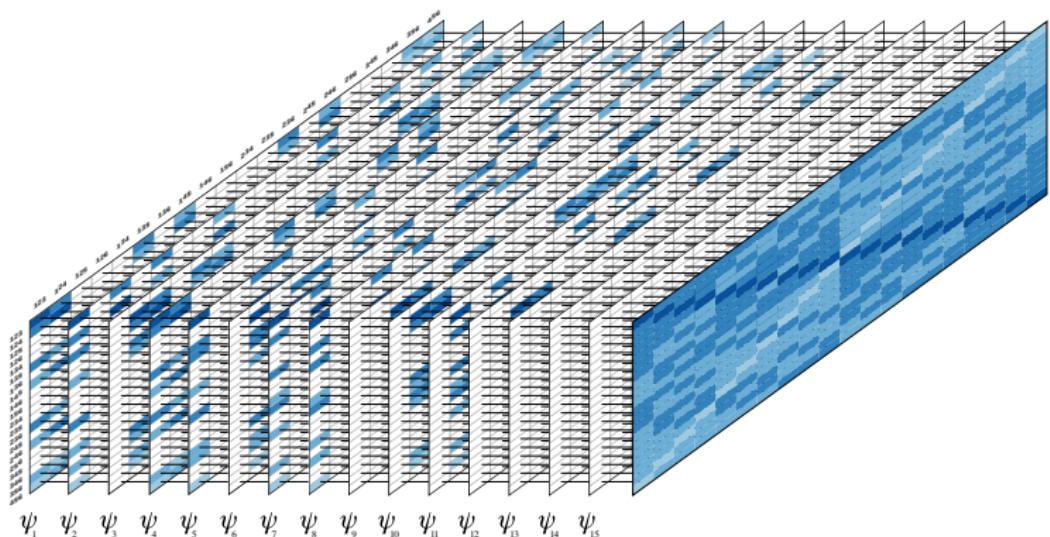
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



# GAMETIC PROBABILITY FOR ALL $\psi$ 's

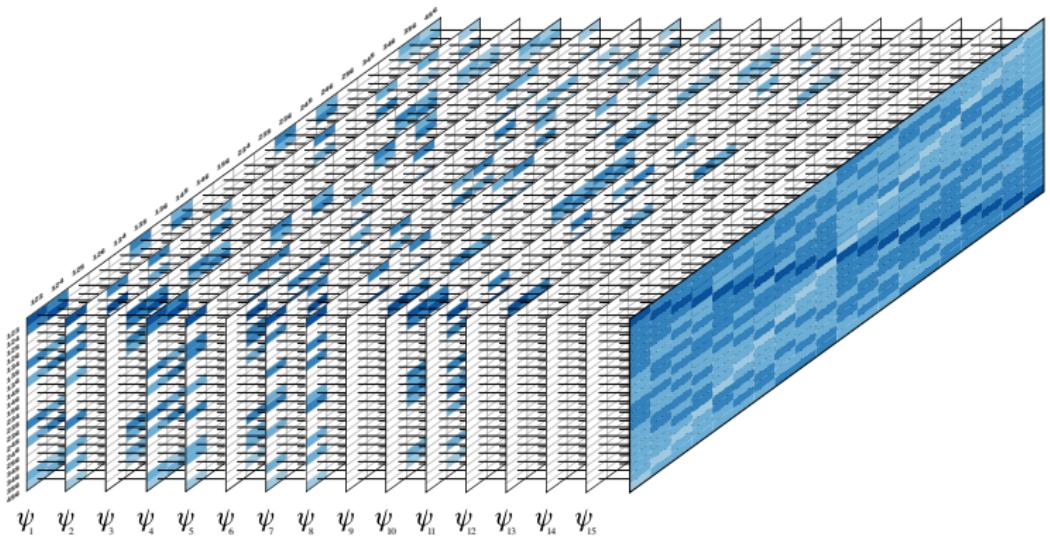


# GAMETIC PROBABILITY FOR ALL $\psi$ 's



$$\Pr(G_A, G_B) = \sum_j \Pr(G_A, G_B \mid \psi_j) \Pr(\psi_j) \quad (1)$$

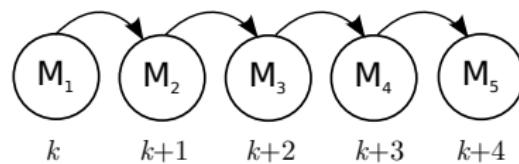
# GAMETIC PROBABILITY FOR ALL $\psi$ 's



$$\begin{aligned}
 \Pr(G_A, G_B) &= \sum_j \Pr(G_A, G_B \mid \psi_j) \Pr(\psi_j) \\
 &= \frac{l! \left(\frac{m}{2} - l\right)! (1 - r_k)^{\frac{m}{2} - l} (r_k)^l}{w_m 2^{\frac{m}{2}}}. \tag{1}
 \end{aligned}$$

# TRANSITION PROBABILITIES AND THE HIDDEN MARKOV MODEL (HMM)

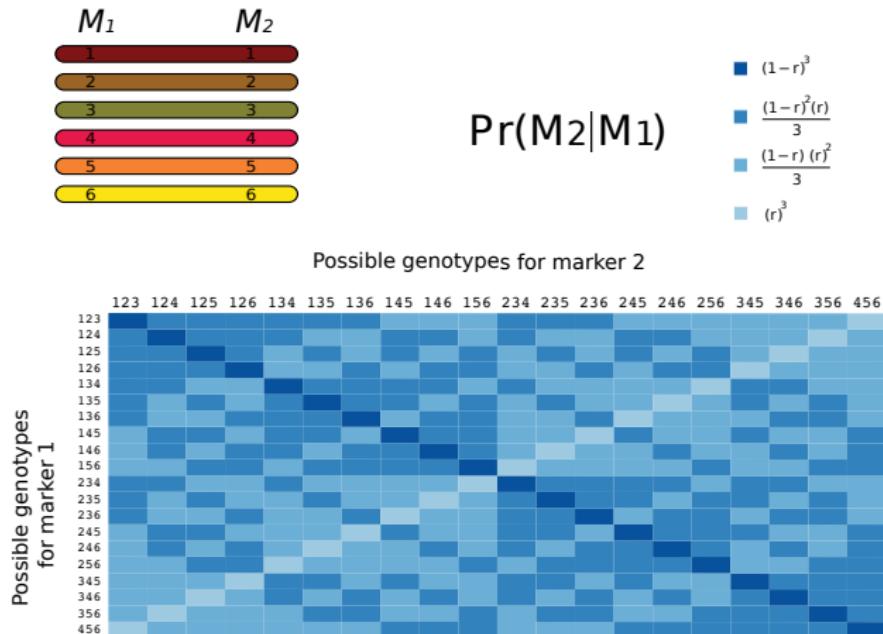
Markov Model: conditional independence



$$\Pr(G_{k+1}|G_k) = \frac{(1 - r_k)^{\frac{p}{2} - l} (r_k)^l}{\binom{\frac{p}{2}}{l}}$$

where  $r_k$  is the recombination fraction between loci  $k$  and  $k + 1$ ,  $p$  is the ploidy level and  $l$  is the number of recombinant events between  $k$  and  $k + 1$ .

# FINAL TRANSITION MATRIX FOR HEXAPLOIDS (ONE PARENT)



## THE MARKOV MODEL

	$M_1$	$M_2$	$M_3$	$M_4$
Individual 1	$\{A_1, A_2, A_3\}$	$\{B_1, B_3, B_6\}$	$\{C_1, C_5, C_6\}$	$\{D_3, D_4, D_6\}$

123 ○  
124 ○  
125 ○  
126 ○  
134 ○  
135 ○  
136 ○  
145 ○  
146 ○  
156 ○  
234 ○  
235 ○  
236 ○  
245 ○  
246 ○  
256 ○  
345 ○  
346 ○  
356 ○  
456 ○

# THE MARKOV MODEL

Individual 1       $M_1$        $M_2$        $M_3$        $M_4$

$\{A_1, A_2, A_3\}$      $\{B_1, B_3, B_6\}$      $\{C_1, C_5, C_6\}$      $\{D_3, D_4, D_6\}$

123    ○  
124    ○  
125    ○  
126    ○  
134    ○  
135    ○  
136    ○  
145    ○  
146    ○  
156    ○  
234    ○  
235    ○  
236    ○  
245    ○  
246    ○  
256    ○  
345    ○  
346    ○  
356    ○  
456    ○

○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○

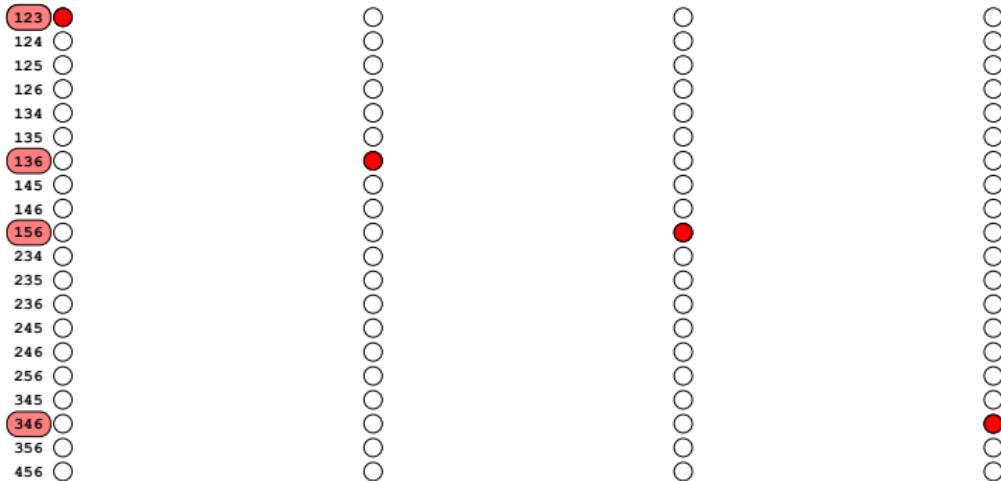
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○

○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○  
○

$\Pr(\text{ Sequence}_1 | \text{Model}) =$

# THE MARKOV MODEL

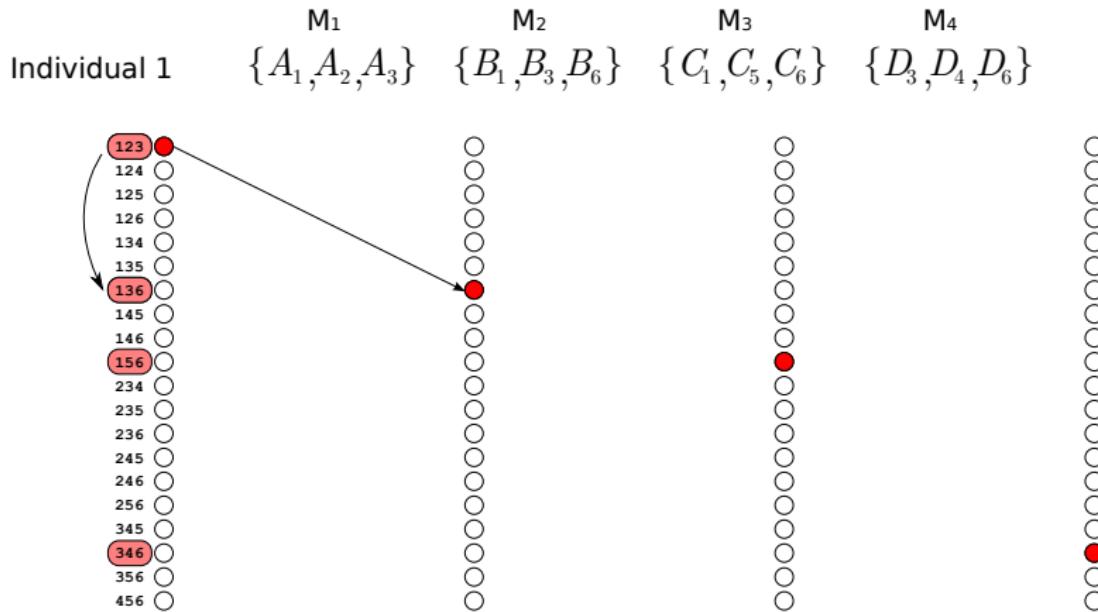
Individual 1       $M_1$        $M_2$        $M_3$        $M_4$   
 $\{A_1, A_2, A_3\}$      $\{B_1, B_3, B_6\}$      $\{C_1, C_5, C_6\}$      $\{D_3, D_4, D_6\}$



$\Pr(\text{ Sequence}_1 | \text{Model}) =$

$\Pr(\{A_1, A_2, A_3\})$

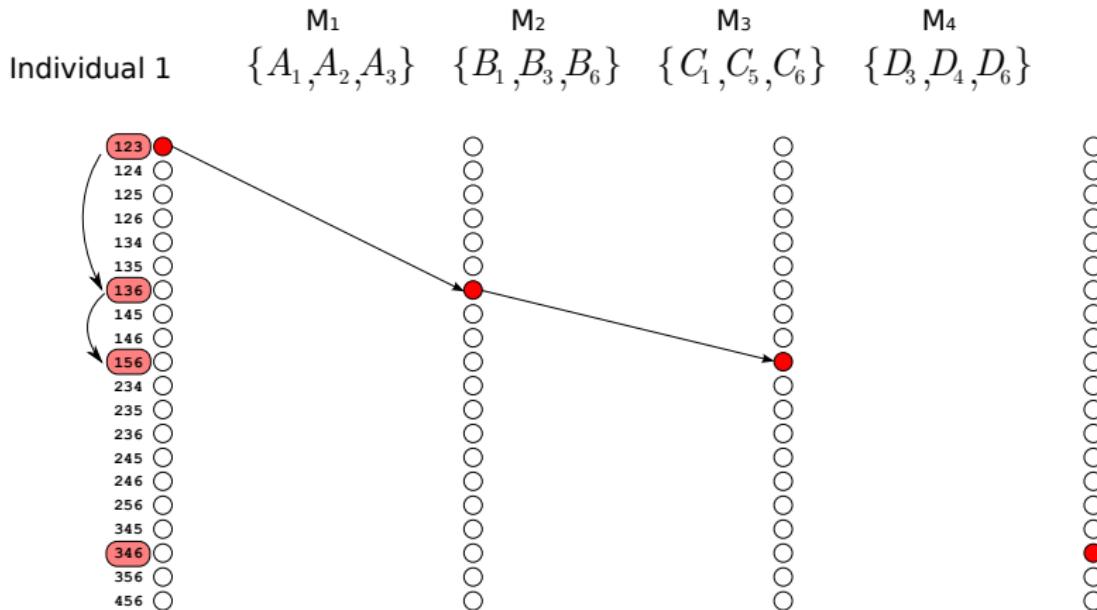
# THE MARKOV MODEL



$\Pr(\text{ Sequence}_1 | \text{Model}) =$

$$\Pr(\{A_1 A_2 A_3\}) \times \Pr(\{B_1 B_3 B_6\} | \{A_1 A_2 A_3\})$$

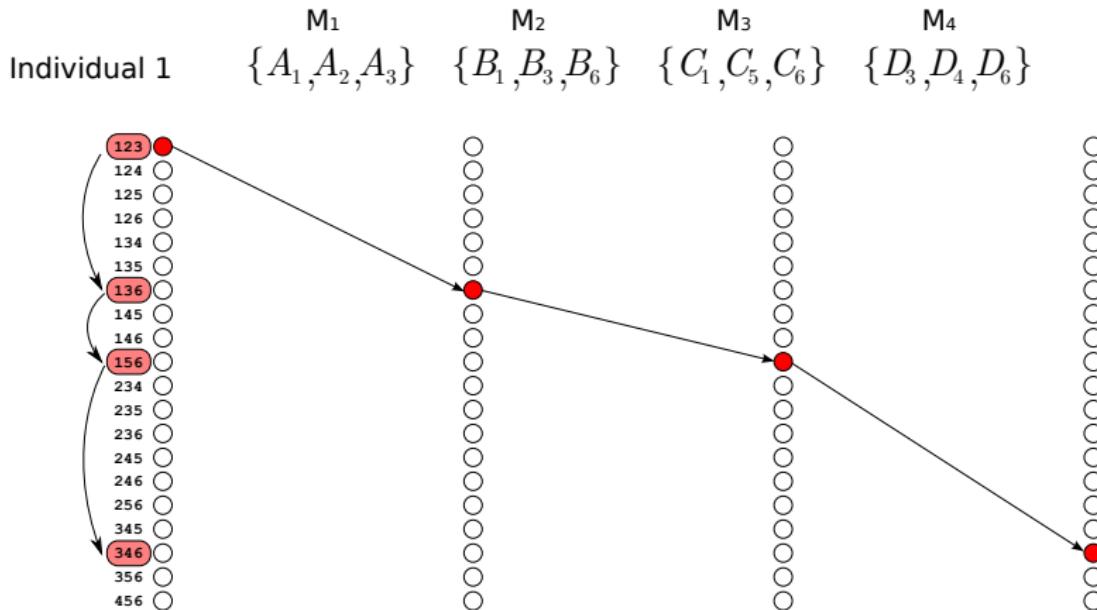
# THE MARKOV MODEL



$\Pr(\text{ Sequence}_1 | \text{Model}) =$

$$\Pr(\{A_1 A_2 A_3\}) \times \Pr(\{B_1 B_3 B_6\} | \{A_1 A_2 A_3\}) \times \Pr(\{C_1 C_5 C_6\} | \{B_1 B_3 B_6\})$$

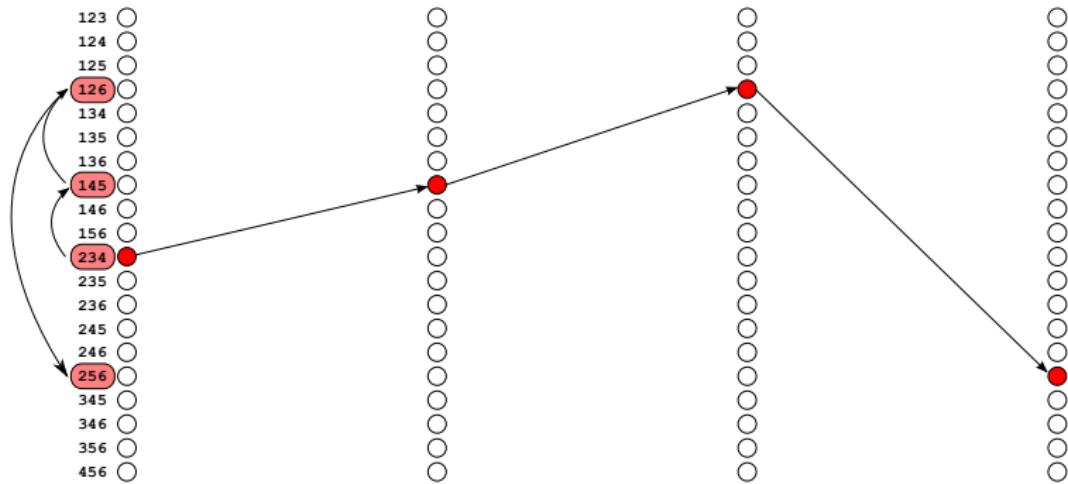
# THE MARKOV MODEL



$$\Pr(\{A_1 A_2 A_3\}) \times \Pr(\{B_1 B_3 B_6\} | \{A_1 A_2 A_3\}) \times \Pr(\{C_1 C_5 C_6\} | \{B_1 B_3 B_6\}) \times \Pr(\{D_3 D_4 D_6\} | \{C_1 C_5 C_6\})$$

# THE MARKOV MODEL

Individual 2       $M_1 \quad \{A_2 A_3 A_4\}$      $M_2 \quad \{B_1 B_4 B_5\}$      $M_3 \quad \{C_1 C_2 C_6\}$      $M_4 \quad \{D_2 D_5 D_6\}$

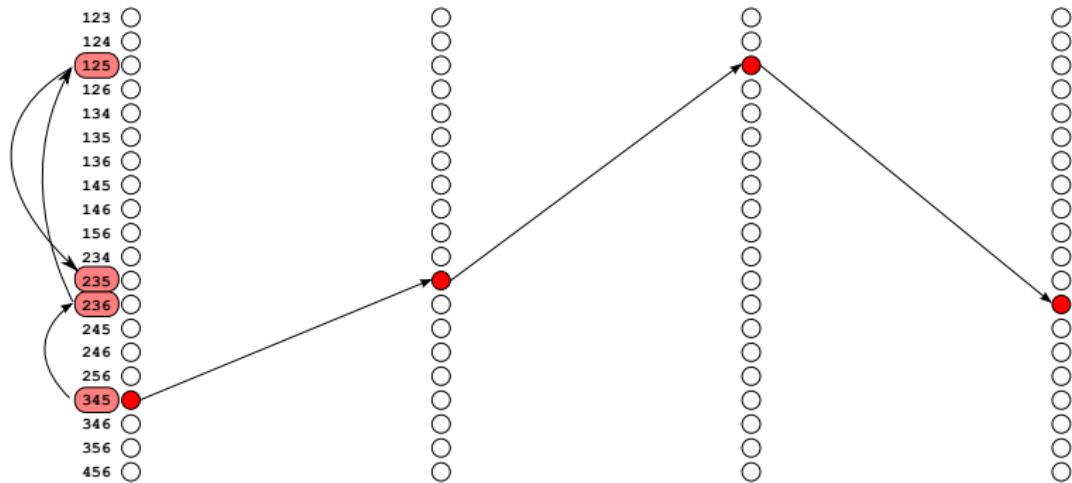


$\Pr(\text{ Sequence}_2 | \text{Model}) =$

$$\Pr(\{A_2 A_3 A_4\}) \times \Pr(\{B_1 B_4 B_5\} | \{A_2 A_3 A_4\}) \times \Pr(\{C_1 C_2 C_6\} | \{B_1 B_4 B_5\}) \times \Pr(\{D_2 D_5 D_6\} | \{C_1 C_2 C_6\})$$

# THE MARKOV MODEL

Individual 3       $M_1 \quad \{A_3 \ A_4 \ A_5\}$      $M_2 \quad \{B_2 \ B_3 \ B_6\}$      $M_3 \quad \{C_1 \ C_2 \ C_5\}$      $M_4 \quad \{D_2 \ D_3 \ D_5\}$



$\Pr(\text{ Sequence}_3 | \text{Model}) =$

$$\Pr(\{A_3 \ A_4 \ A_5\}) \times \Pr(\{B_2 \ B_3 \ B_6\} | \{A_3 \ A_4 \ A_5\}) \times \Pr(\{C_1 \ C_2 \ C_5\} | \{B_2 \ B_3 \ B_6\}) \times \Pr(\{D_2 \ D_3 \ D_5\} | \{C_1 \ C_2 \ C_5\})$$

## LIKELIHOOD

$$L = \prod_{\text{all individuals}} \Pr(\text{Sequence}_i | \text{Model})$$

Maximizing the likelihood function, it is possible to obtain the estimates of

$$[r_{AB}, r_{BC}, r_{CD}]$$

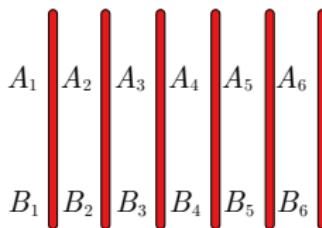
# LIKELIHOOD

$$L = \prod_{\text{all individuals}} \Pr(\text{Sequence}_i | \text{Model})$$

Maximizing the likelihood function, it is possible to obtain the estimates of

$$[r_{AB}, r_{BC}, r_{CD}]$$

However, we access the genotypes using biallelic SPNs.



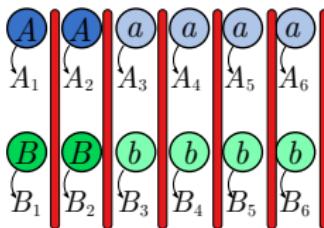
# LIKELIHOOD

$$L = \prod_{\text{all individuals}} \Pr(\text{Sequence}_i | \text{Model})$$

Maximizing the likelihood function, it is possible to obtain the estimates of

$$[r_{AB}, r_{BC}, r_{CD}]$$

However, we access the genotypes using biallelic SPNs.



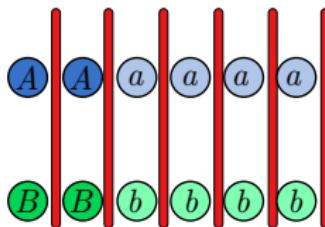
# LIKELIHOOD

$$L = \prod_{\text{all individuals}} \Pr(\text{Sequence}_i | \text{Model})$$

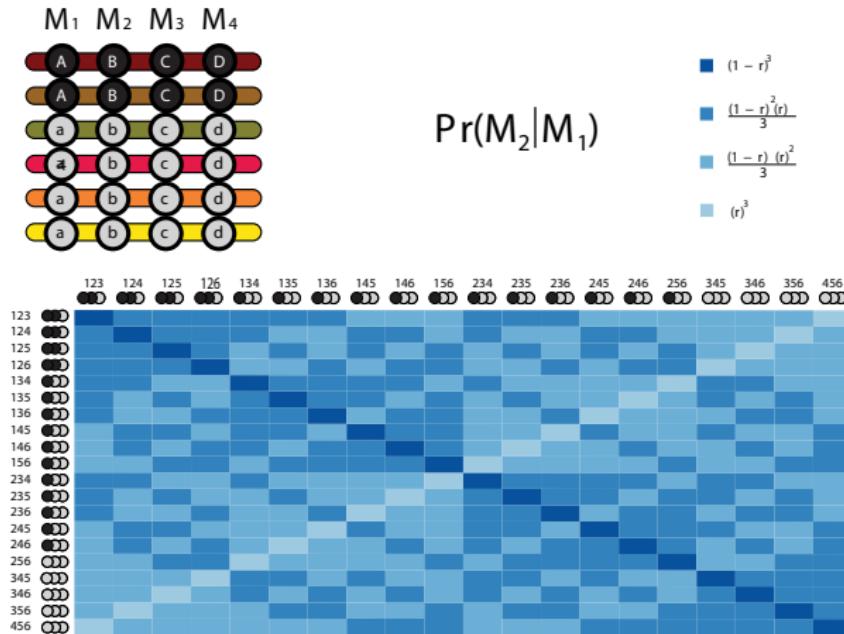
Maximizing the likelihood function, it is possible to obtain the estimates of

$$[r_{AB}, r_{BC}, r_{CD}]$$

However, we access the genotypes using biallelic SPNs.

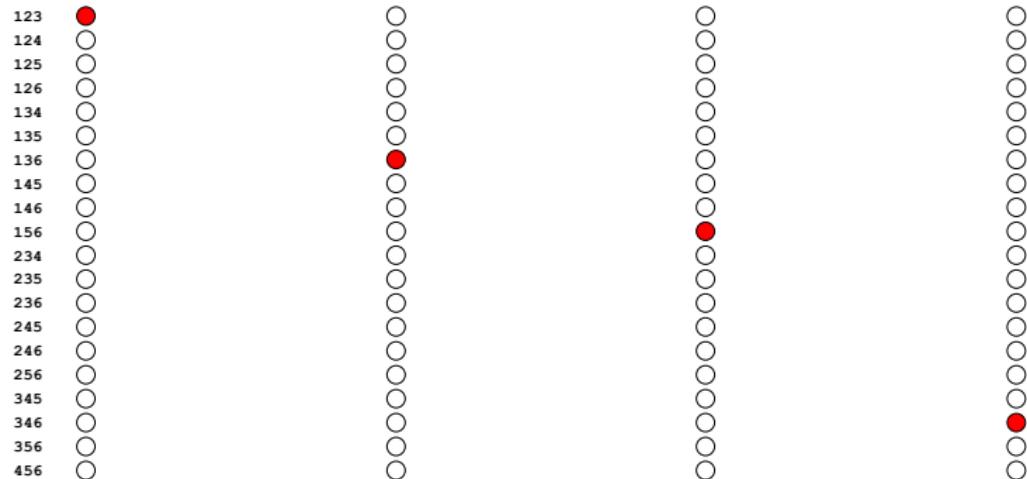


## FINAL TRANSITION MATRIX FOR HEXAPLOIDS (ONE PARENT)

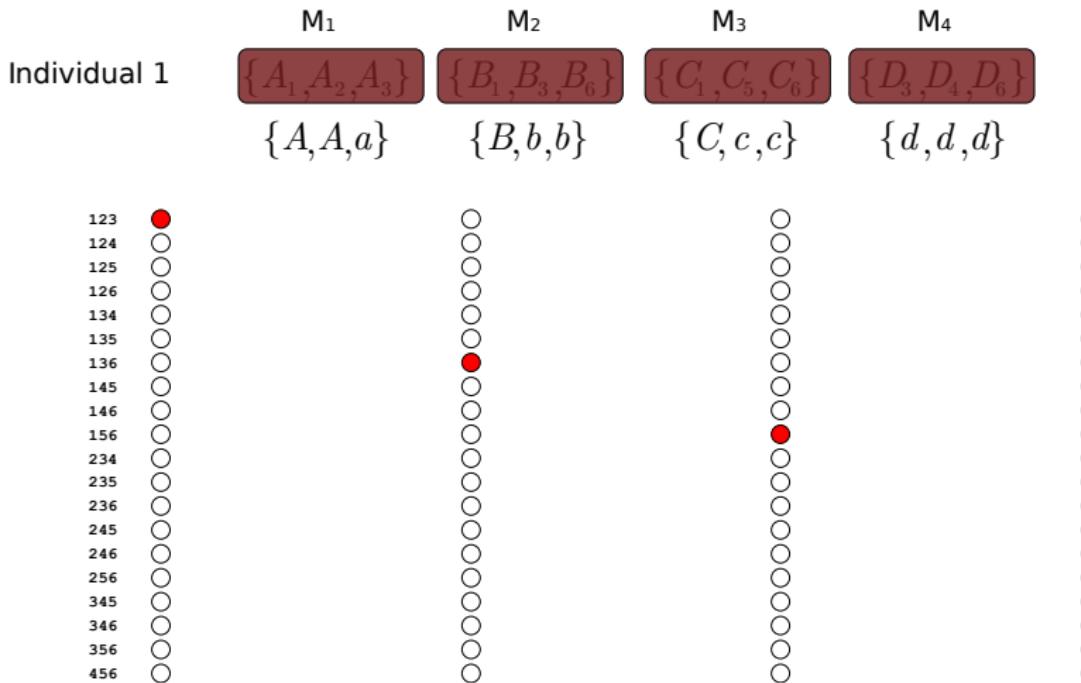


# HIDDEN MARKOV MODEL (HMM)

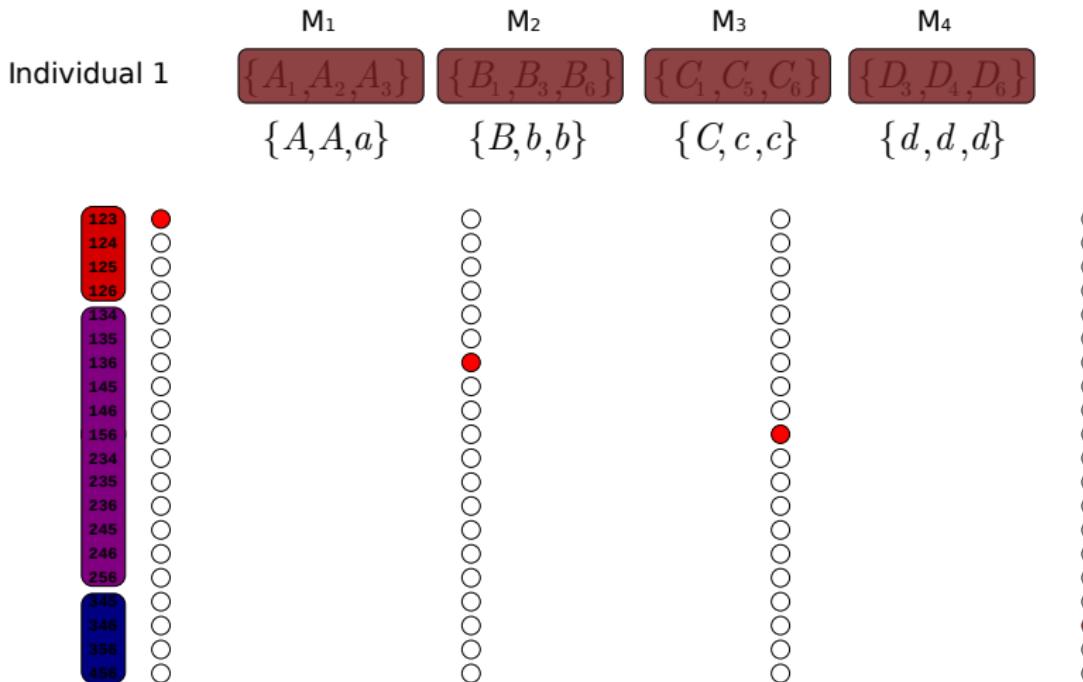
	M <sub>1</sub>	M <sub>2</sub>	M <sub>3</sub>	M <sub>4</sub>
Individual 1	{P <sub>1</sub> <sup>1</sup> , P <sub>1</sub> <sup>2</sup> , P <sub>1</sub> <sup>3</sup> }	{P <sub>2</sub> <sup>1</sup> , P <sub>2</sub> <sup>3</sup> , P <sub>2</sub> <sup>6</sup> }	{P <sub>3</sub> <sup>1</sup> , P <sub>3</sub> <sup>5</sup> , P <sub>3</sub> <sup>6</sup> }	{P <sub>4</sub> <sup>3</sup> , P <sub>4</sub> <sup>4</sup> , P <sub>4</sub> <sup>6</sup> }



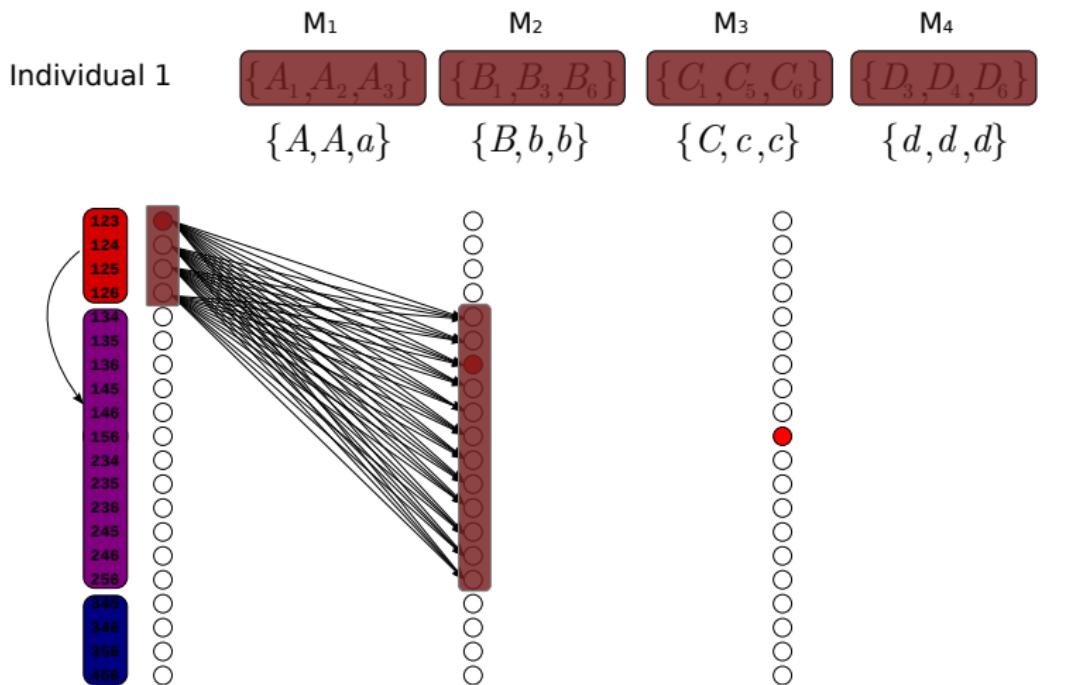
# HIDDEN MARKOV MODEL (HMM)



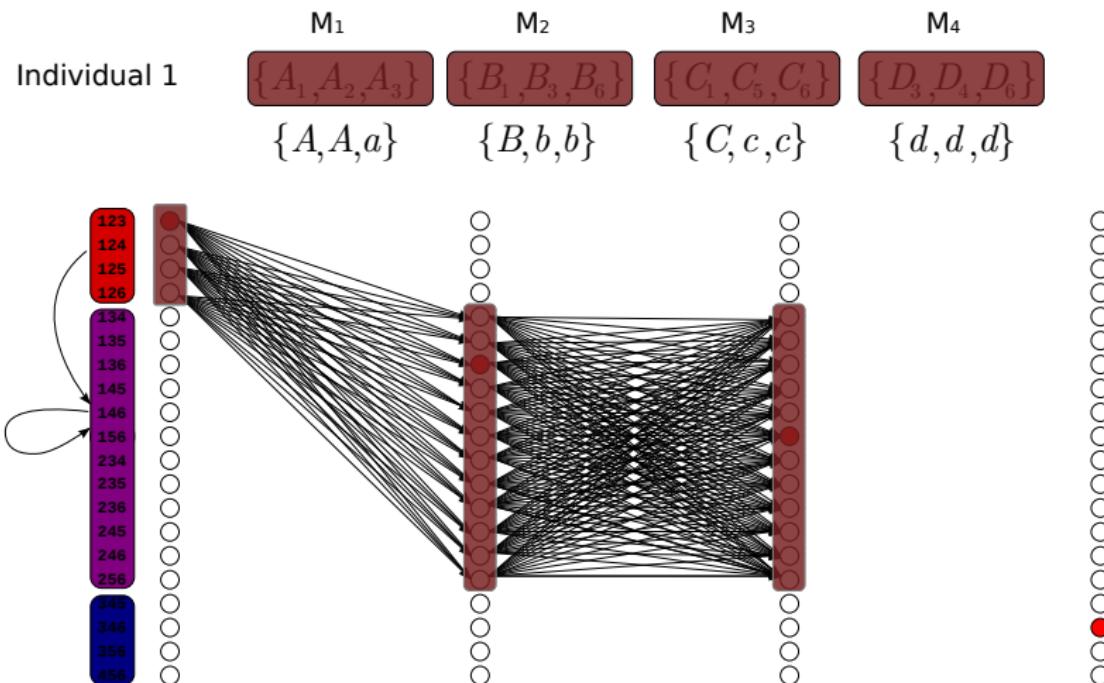
# HIDDEN MARKOV MODEL (HMM)



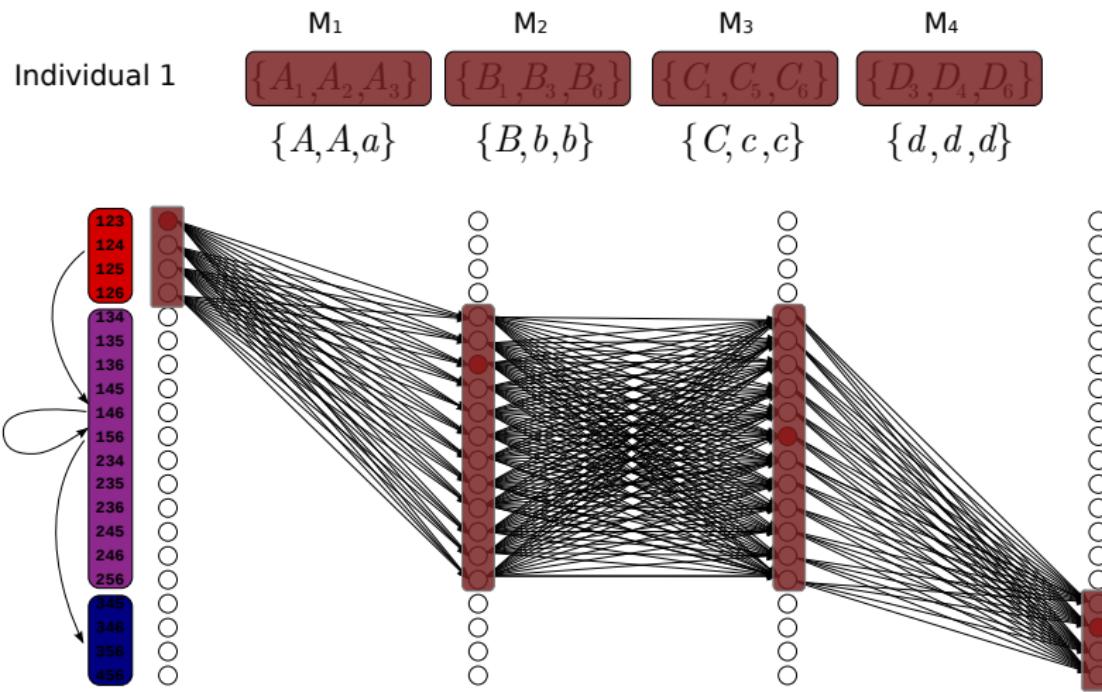
# HIDDEN MARKOV MODEL (HMM)



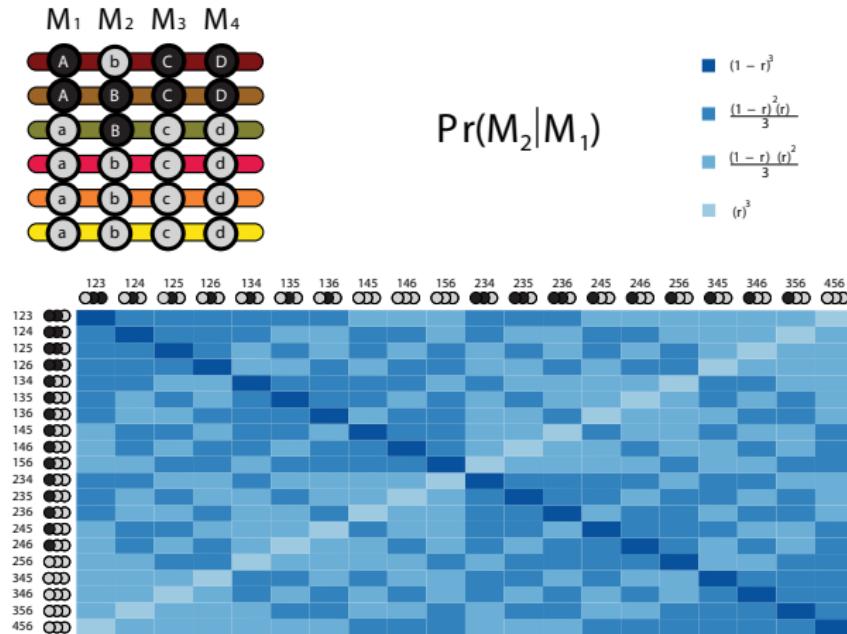
# HIDDEN MARKOV MODEL (HMM)



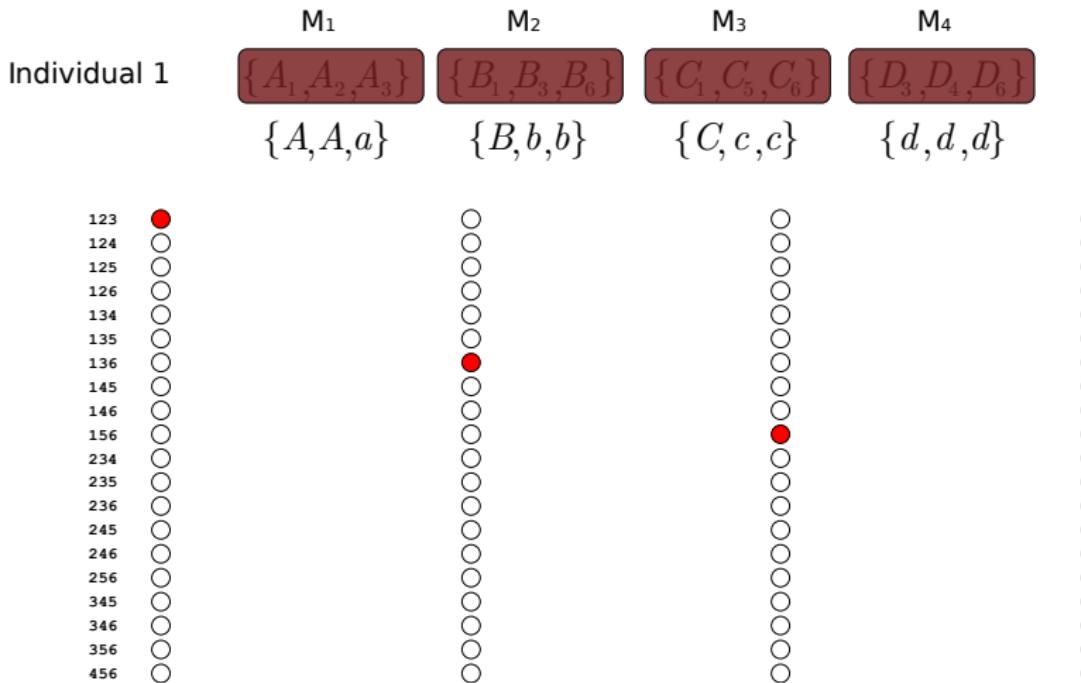
# HIDDEN MARKOV MODEL (HMM)



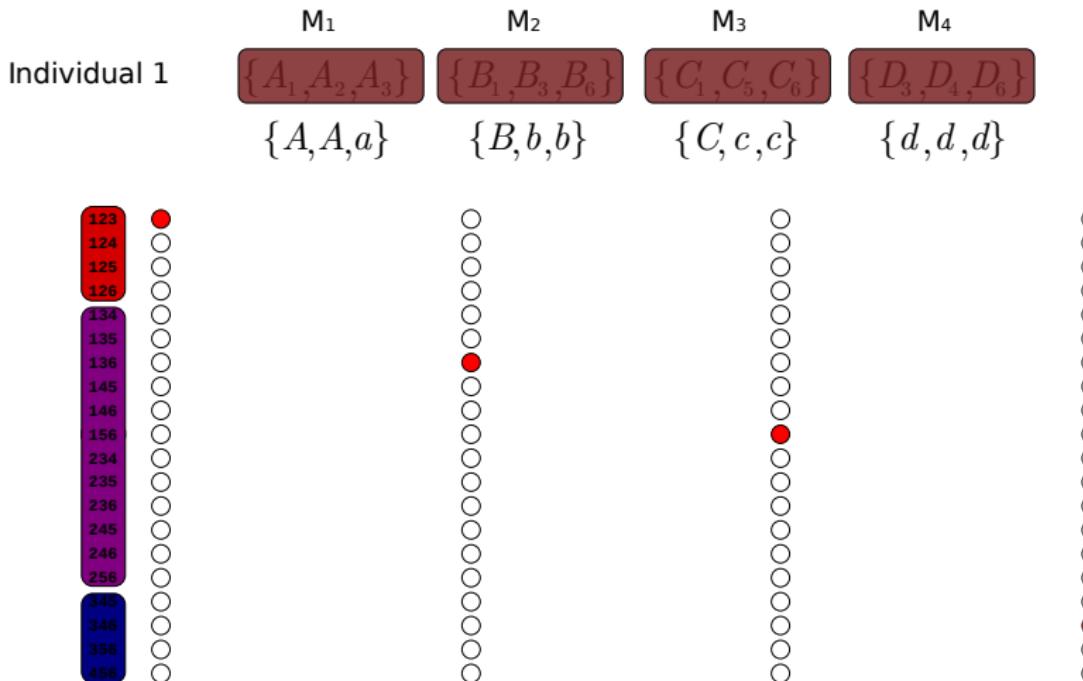
## FINAL TRANSITION MATRIX FOR HEXAPLOIDS (ONE PARENT)



# HIDDEN MARKOV MODEL (HMM)

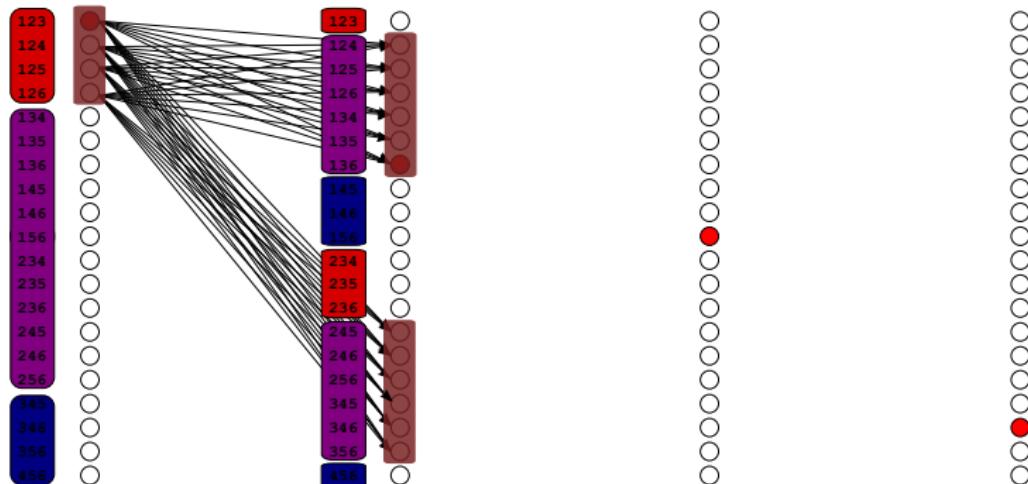


# HIDDEN MARKOV MODEL (HMM)



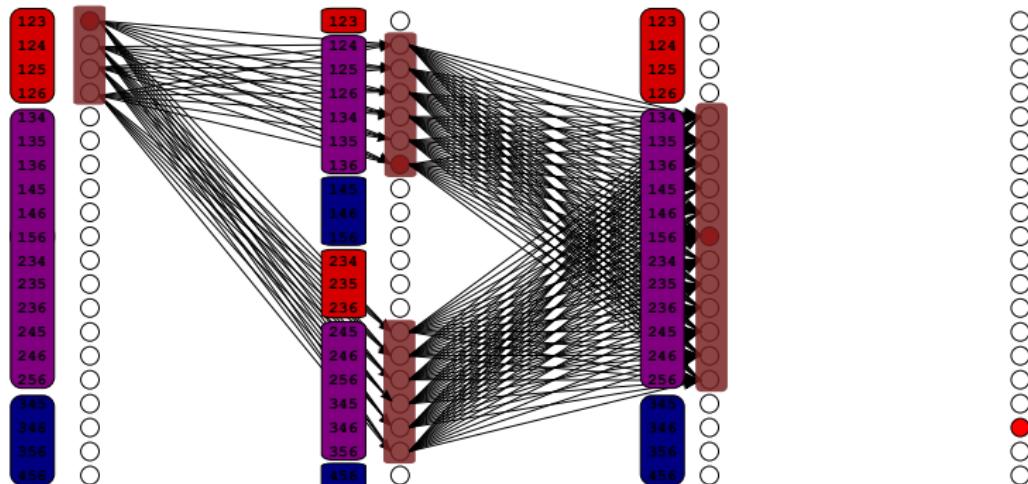
# HIDDEN MARKOV MODEL (HMM)

	$M_1$	$M_2$	$M_3$	$M_4$
Individual 1	$\{A_1, A_2, A_3\}$	$\{B_1, B_3, B_6\}$	$\{C_1, C_5, C_6\}$	$\{D_3, D_4, D_6\}$
	$\{A, A \ a\}$	$\{B, b \ b\}$	$\{C, c \ c\}$	$\{d, d, d\}$

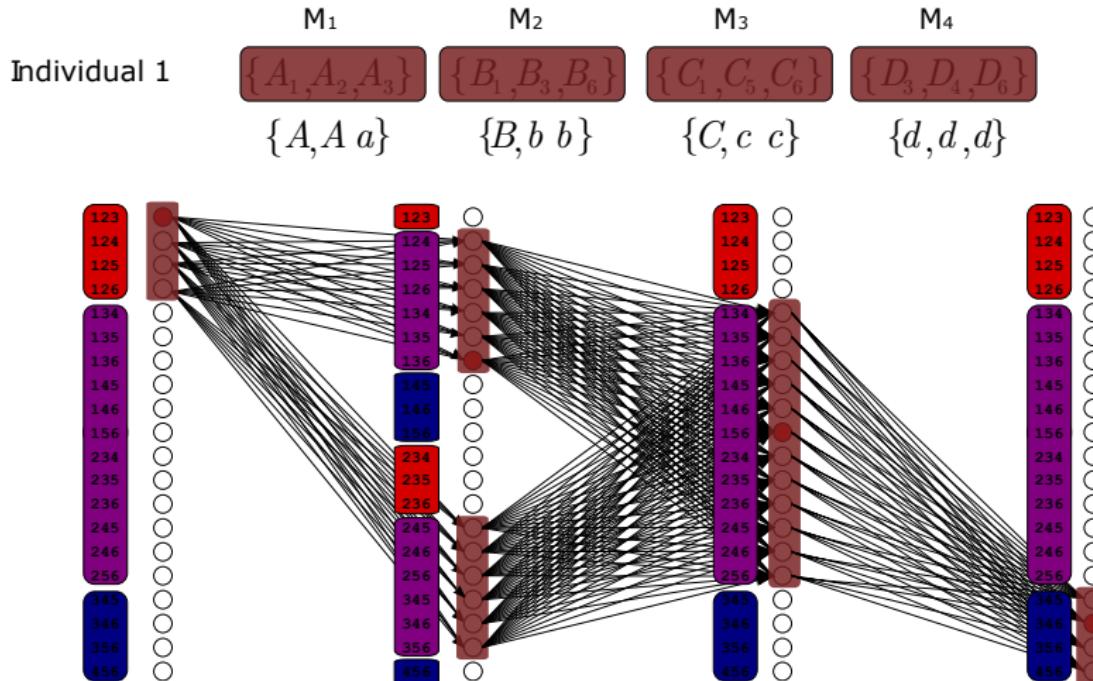


# HIDDEN MARKOV MODEL (HMM)

	$M_1$	$M_2$	$M_3$	$M_4$
Individual 1	$\{A_1, A_2, A_3\}$	$\{B_1, B_3, B_6\}$	$\{C_1, C_5, C_6\}$	$\{D_3, D_4, D_6\}$
	$\{A, A \ a\}$	$\{B, b \ b\}$	$\{C, c \ c\}$	$\{d, d, d\}$



# HIDDEN MARKOV MODEL (HMM)



# HMM - LIKELIHOOD

## LIKELIHOOD - *Forward* PROCEDURE

1. Initialization:

$$\alpha_1(j) = \pi_j b_j(o_1), j = 1, \dots, g_m$$

2. Induction:

$$\alpha_{k+1}(j') = \left[ \sum_j^{g_m} \alpha_k(j) t_k(j, j') \right] b_{j'}(o_{k+1})$$

where  $k = 1, \dots, z-1$  and  $j = 1, \dots, g_m$

3. Termination:

$$\Pr(o_1, \dots, o_z | \mathbf{r}, \Omega_P, \Omega_Q) = \sum_{j=1}^{g_m} \alpha_z(j)$$

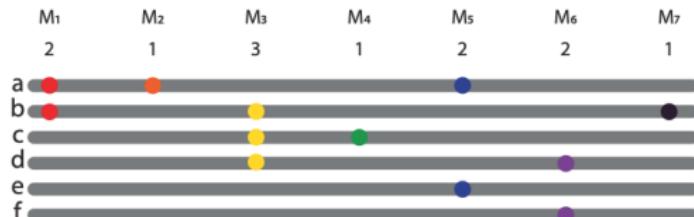
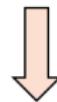
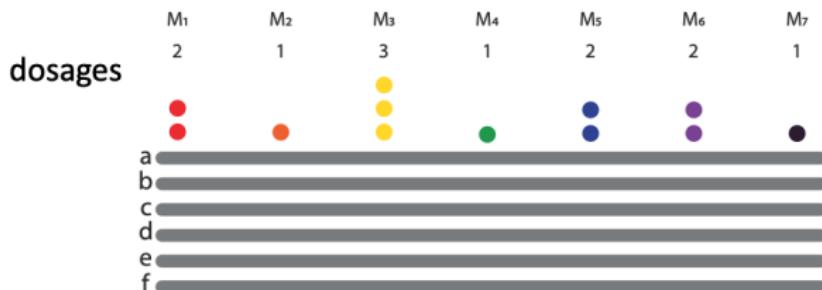
Then, the likelihood of the model can be defined as

$$L = \prod_{i=1}^n \Pr(o_1, \dots, o_z | \mathbf{r}, \Omega_P, \Omega_Q)$$

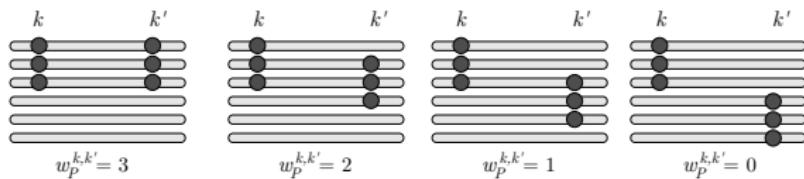
where  $n$  is the number of individuals on the full sib family.

# HAPLOTYPING OF PHASING

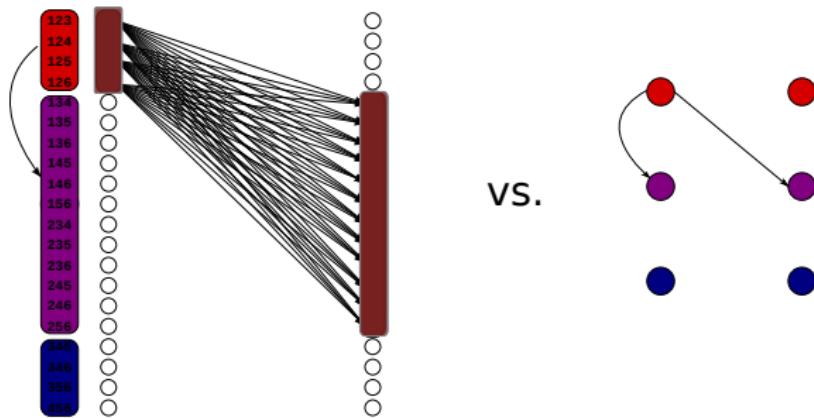
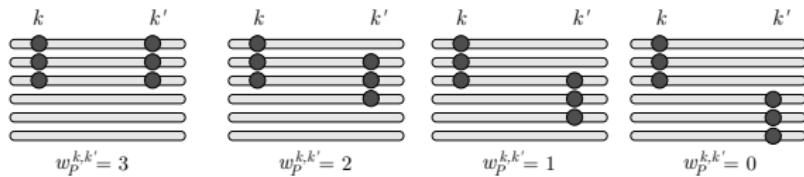
- Disposition of allelic variants in the homologs in a homology group



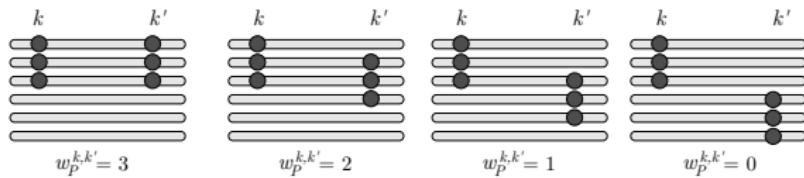
# REDUCTION OF COMPLEXITY - THE TWO-POINT APPROACH



# REDUCTION OF COMPLEXITY - THE TWO-POINT APPROACH

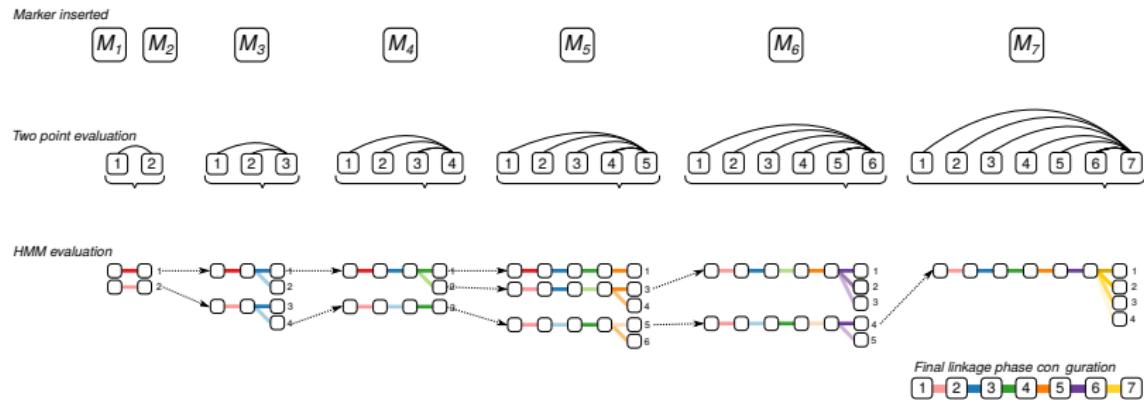


# REDUCTION OF COMPLEXITY - THE TWO-POINT APPROACH

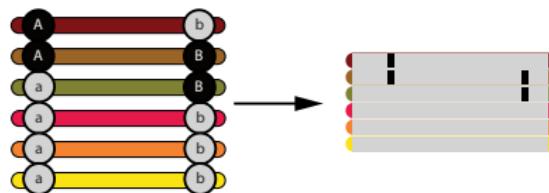
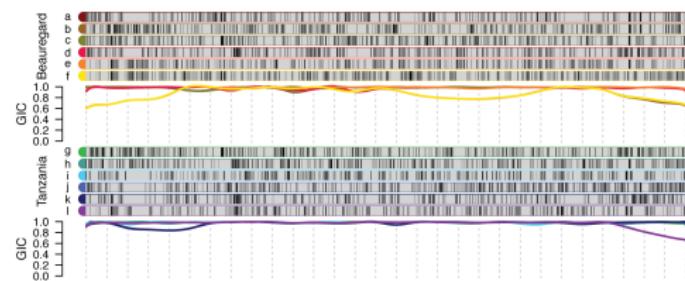
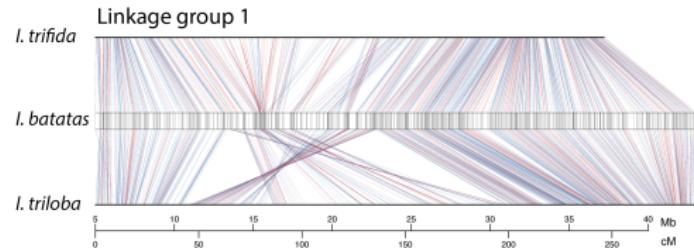


	LOD_ph	rf	LOD_rf
1	0.000000	0.07239439	6.7773070858
2	2.302351	0.09233424	4.8678664544
3	2.855213	0.36392079	4.0093145291
0	6.864653	0.49994804	0.0001254893

# HYBRID PHASING - TWO-POINT AND HMM APPROACHES



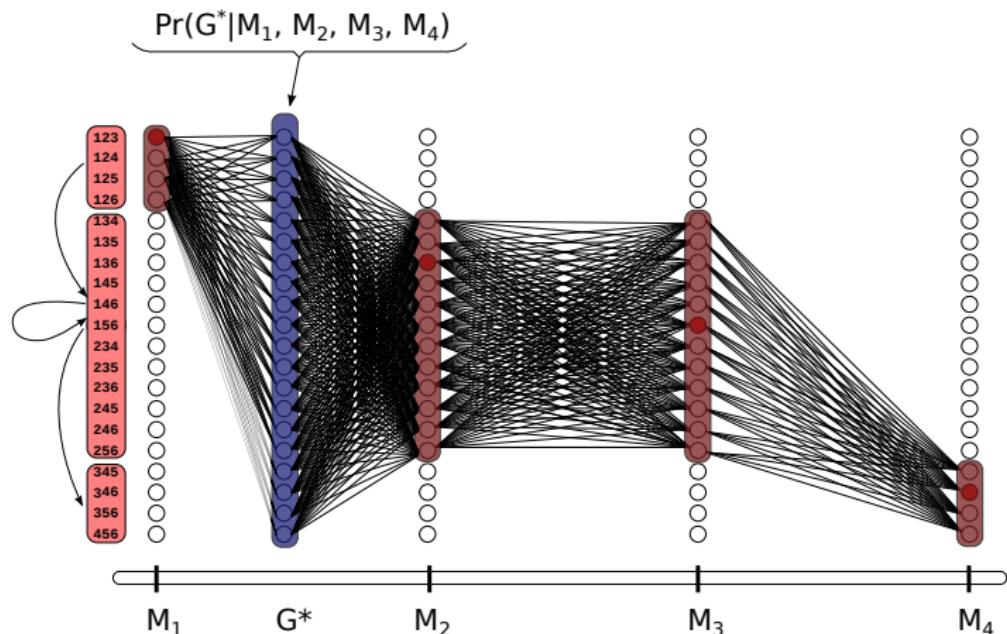
# HOMOLOGOUS GROUP 1 - SWEETPOTATO



# HMM - ESTIMATION PROCEDURE AND IMPLEMENTATION

- ▶ Maximum likelihood estimation: Baum-Welch algorithm
- ▶ Implemented in C++ and R
- ▶ Available under GPL.
- ▶ GitHub: <https://github.com/mmollina/>
- ▶ R package - MAPpoly

# QTL CONDITIONAL PROBABILITIES AND HMM



QTL CONDITIONAL PROBABILITIES AND HMM

