

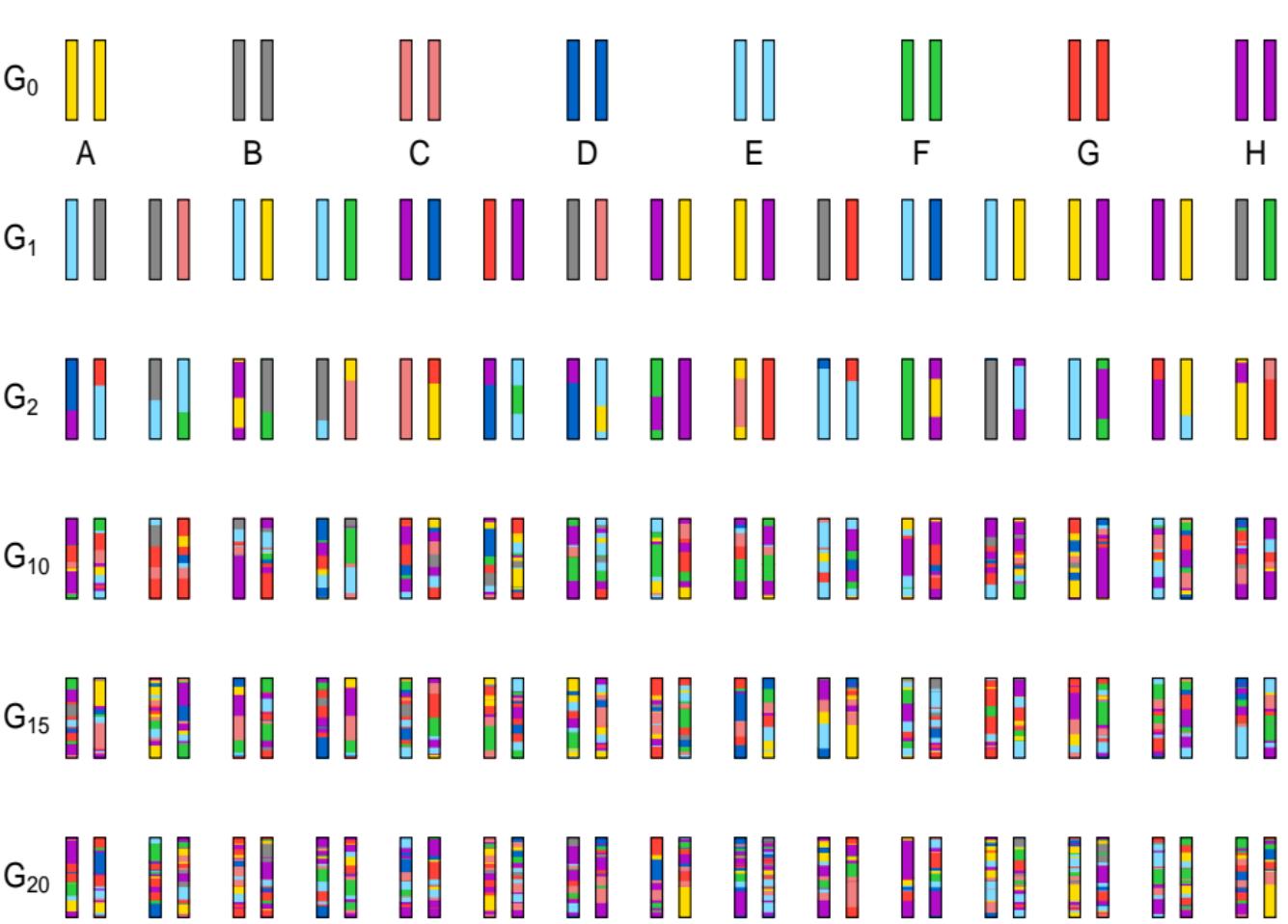
# Genotype Reconstruction for Diversity Outbred Mice

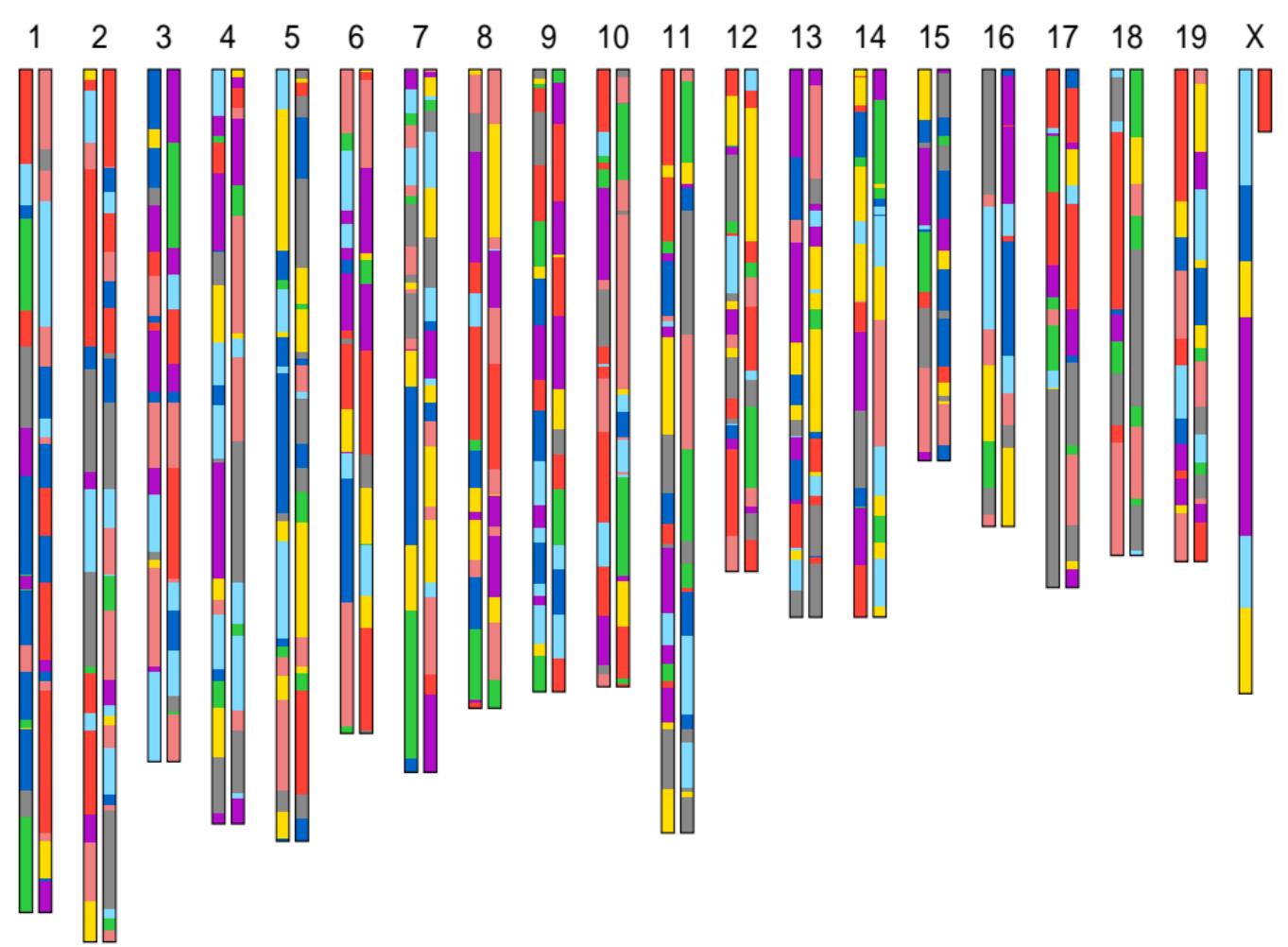
A comparison of r/qt12 and DOQTL

John Spaw<sup>1</sup>  
Karl Broman<sup>2</sup>

[github.com/JohnPSpaw<sup>1</sup>](https://github.com/JohnPSpaw)  
[github.com/kbroman<sup>2</sup>](https://github.com/kbroman)







# Hidden Markov Model

# Hidden Markov Model

Probabilities are generated from:

# Hidden Markov Model

Probabilities are generated from:

1. Transition Model

# Hidden Markov Model

Probabilities are generated from:

1. Transition Model
2. **Emission Model**

# Hidden Markov Model

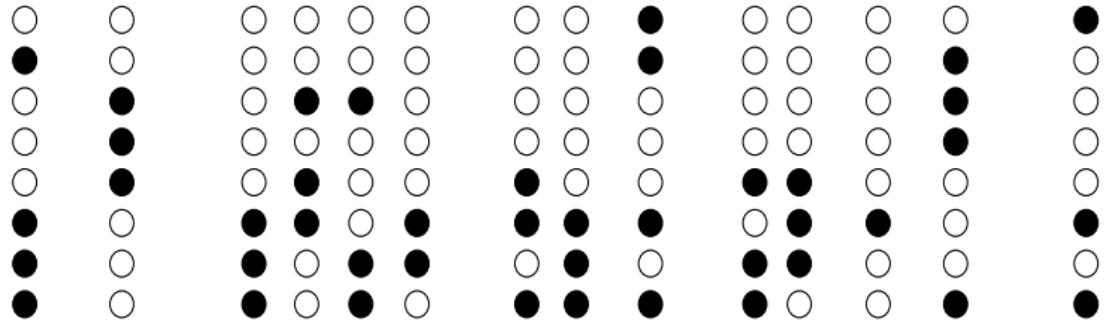
Probabilities are generated from:

1. Transition Model
2. **Emission Model**

Conditional probability of observed data given underlying diplotype state

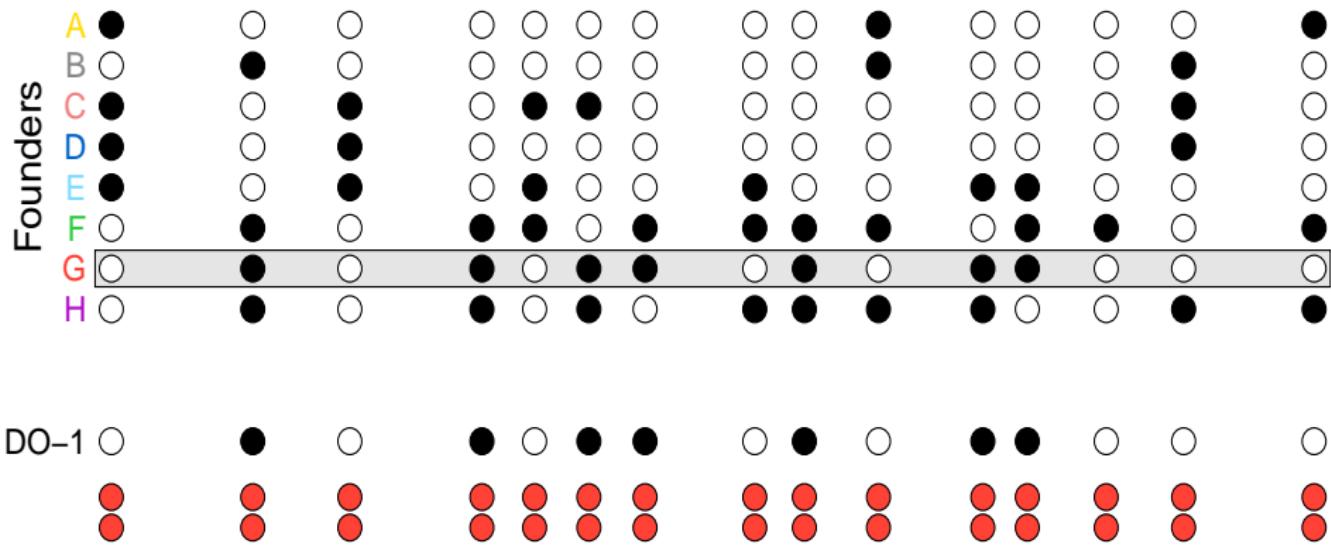
Founders

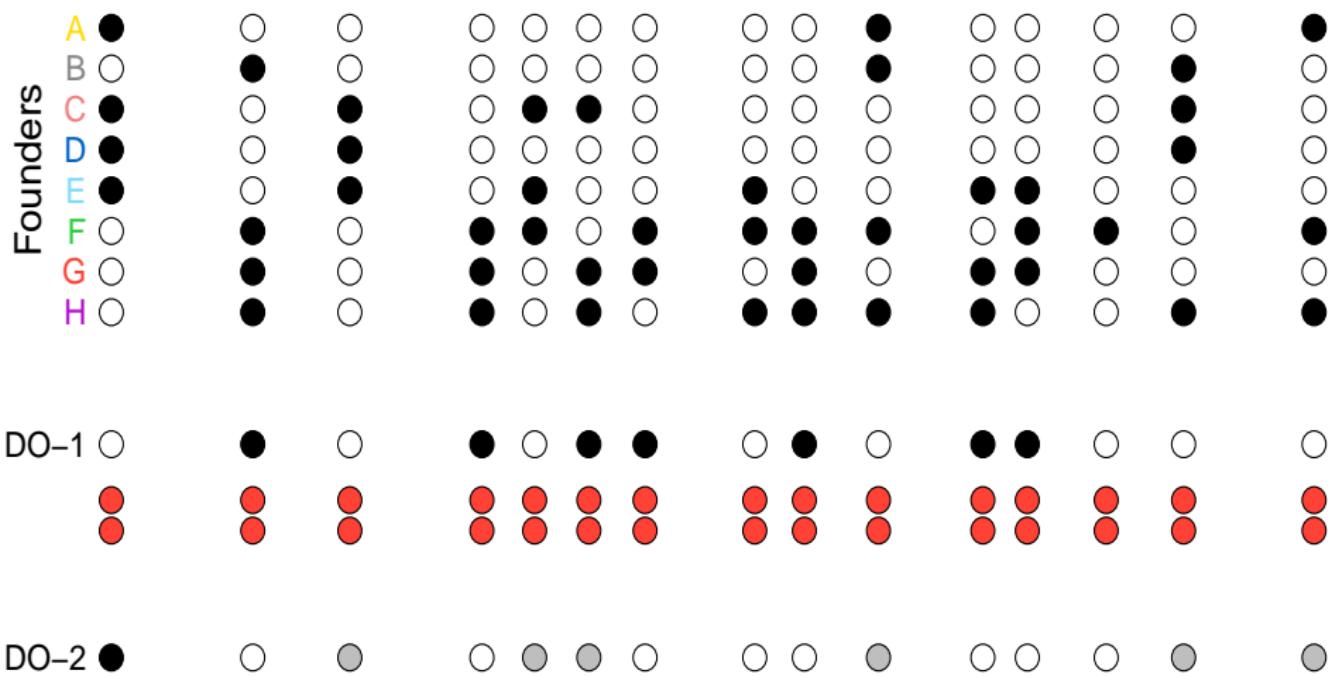
A ●  
B ○  
C ●●  
D ●●●  
E ●●  
F ○  
G ○  
H ○

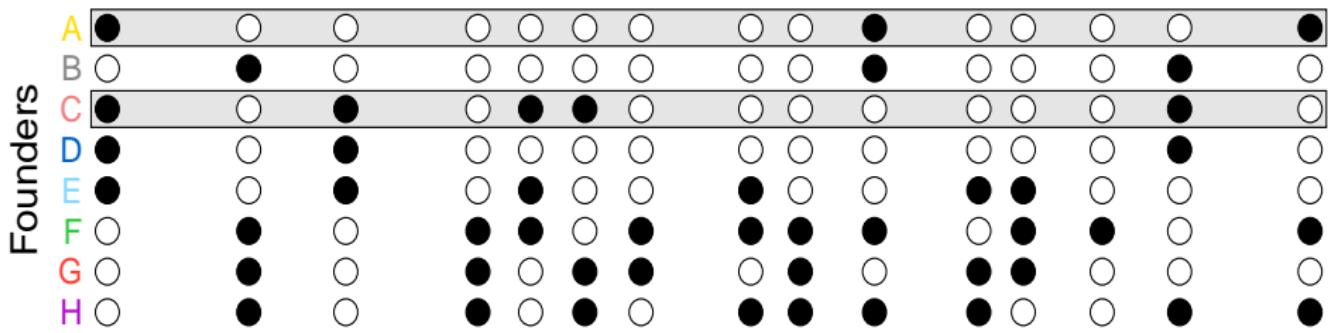


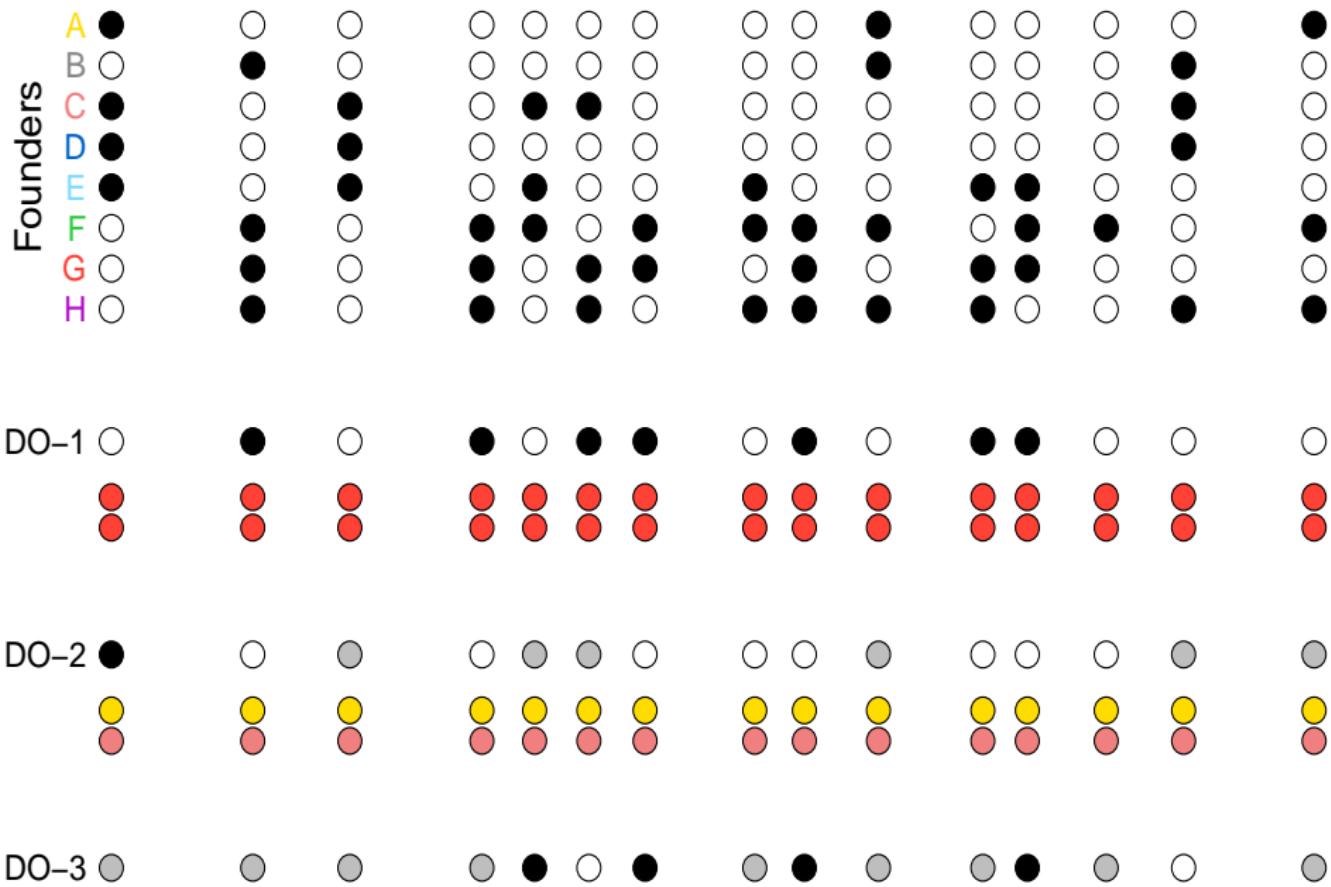
DO-1 ○

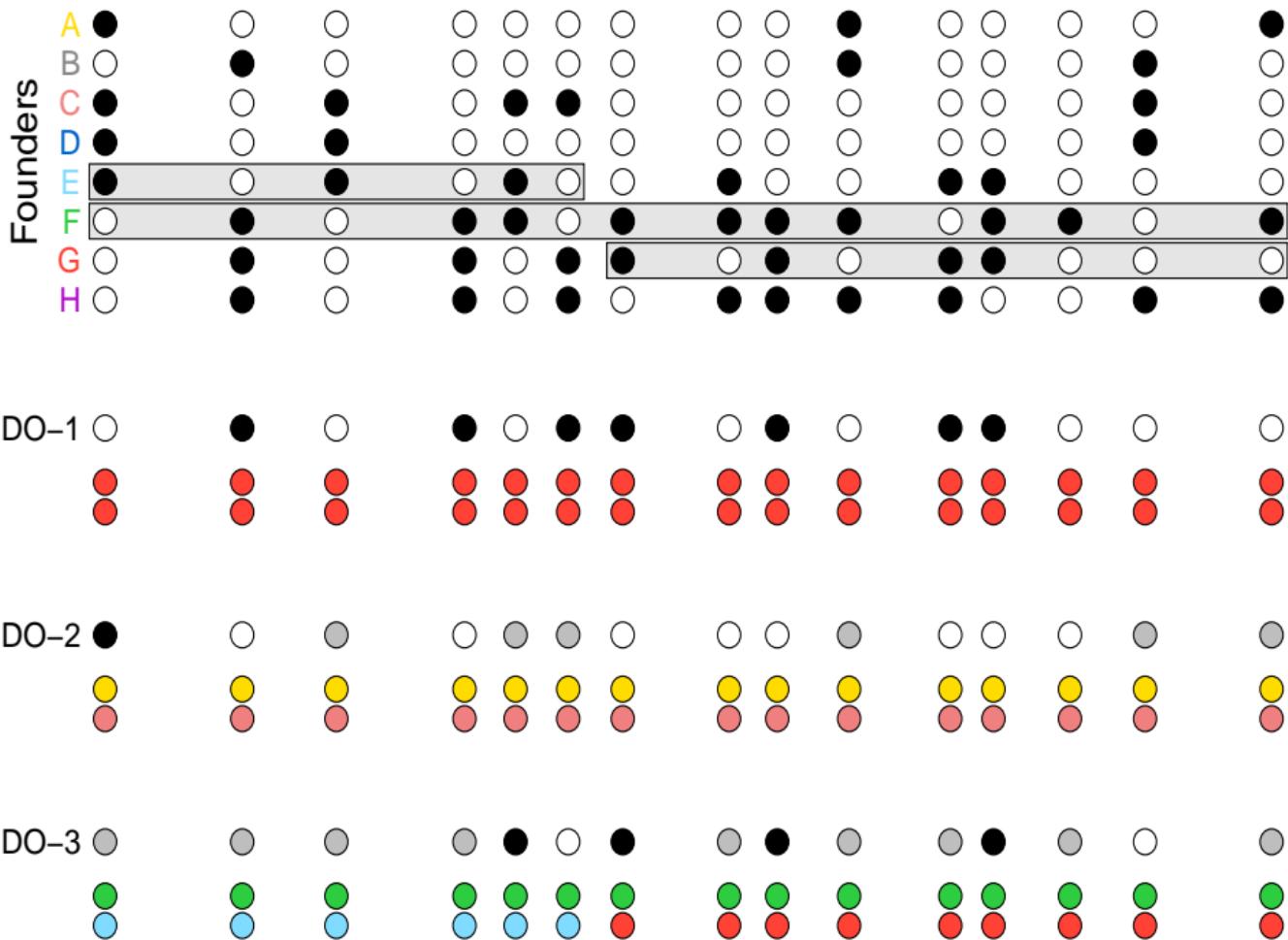
● ○ ● ○ ● ● ○ ● ● ○ ○ ○ ○



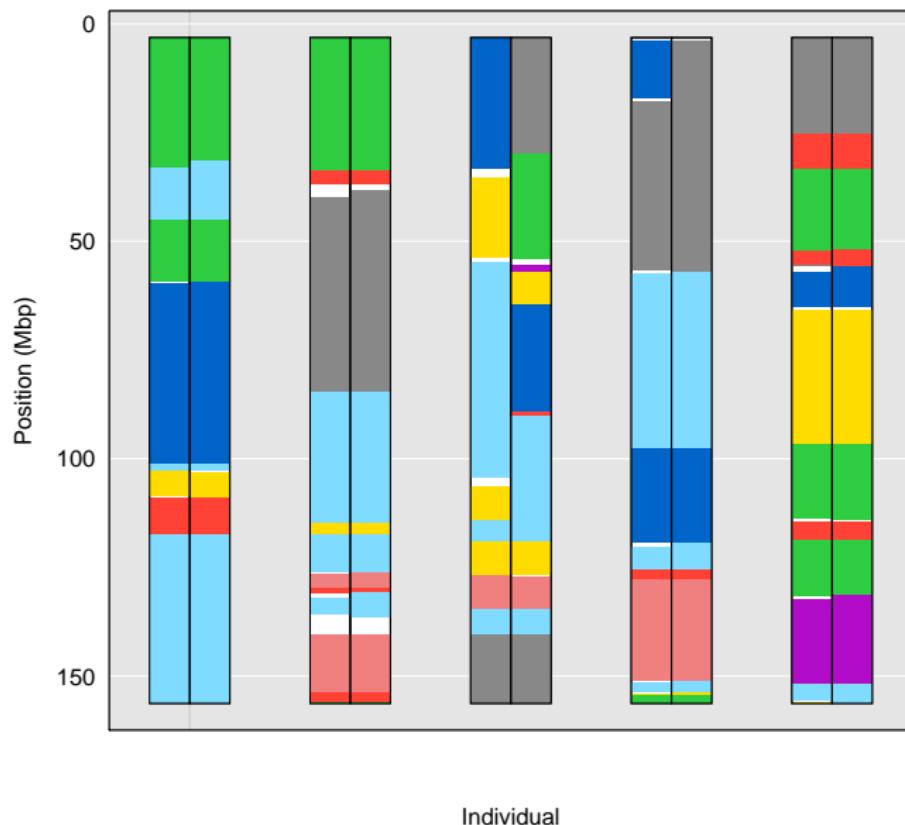








# Inferred Haplotypes



# Data

Two *large* 3D arrays of emission probabilities

- ▶ `r\qt12` (Broman)
- ▶ `DOQTL` (Gatti)

# Data

Two *large* 3D arrays of emission probabilities

- ▶ r\qt12 (Broman)
- ▶ DOQTL (Gatti)

$500 \times 120,000 \times 36$  (Individual  $\times$  Markers  $\times$  Diplotype)

# Data

Two *large* 3D arrays of emission probabilities

- ▶ r\qt12 (Broman)
- ▶ DOQTL (Gatti)

$500 \times 120,000 \times 36$  (Individual  $\times$  Markers  $\times$  Diplotype)

**How are they different?**

# Data

Two *large* 3D arrays of emission probabilities

- ▶ r\qt12 (Broman)
- ▶ DOQTL (Gatti)

$500 \times 120,000 \times 36$  (Individual  $\times$  Markers  $\times$  Diplotype)

**How are they different?**

**How can we visualize this?**

# Measure of distance

# Measure of distance

For each individual at each marker:

Compute *sum of absolute differences*

$$\sum_{i=1}^{36} |p_{1,i} - p_{2,i}|$$

## Measure of distance

For each individual at each marker:

Compute *sum of absolute differences*

$$\sum_{i=1}^{36} |p_{1,i} - p_{2,i}|$$

Each entry represents distance between methods

## Measure of distance

For each individual at each marker:

Compute *sum of absolute differences*

$$\sum_{i=1}^{36} |p_{1,i} - p_{2,i}|$$

Each entry represents distance between methods

**Reduces problems to two-dimensions:**

## Measure of distance

For each individual at each marker:

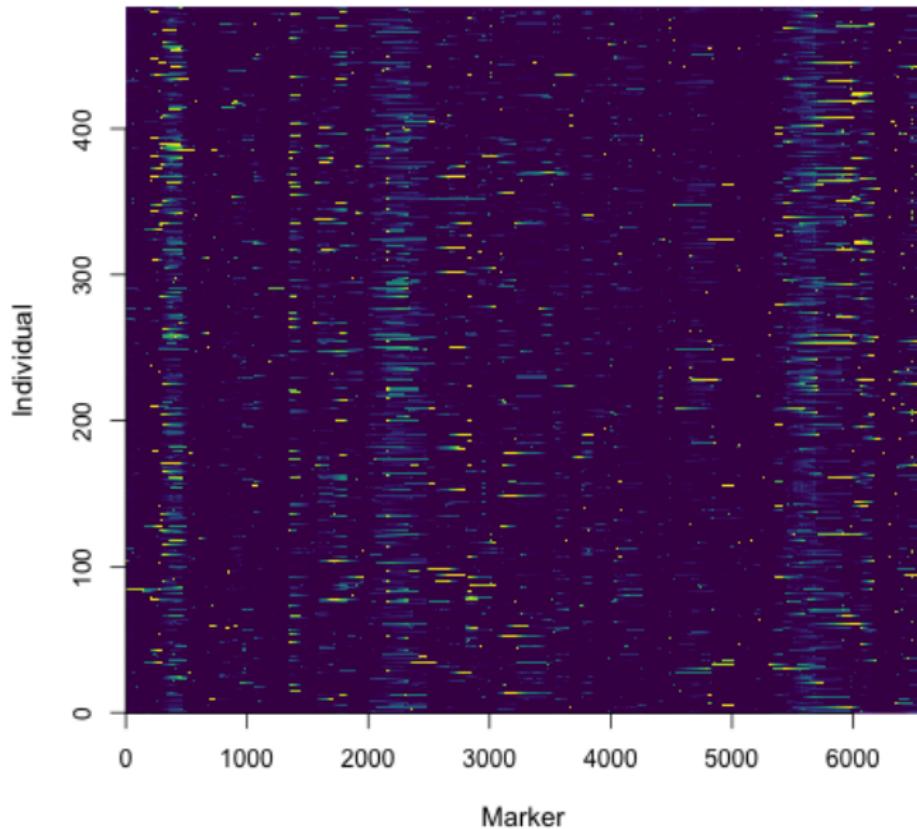
Compute *sum of absolute differences*

$$\sum_{i=1}^{36} |p_{1,i} - p_{2,i}|$$

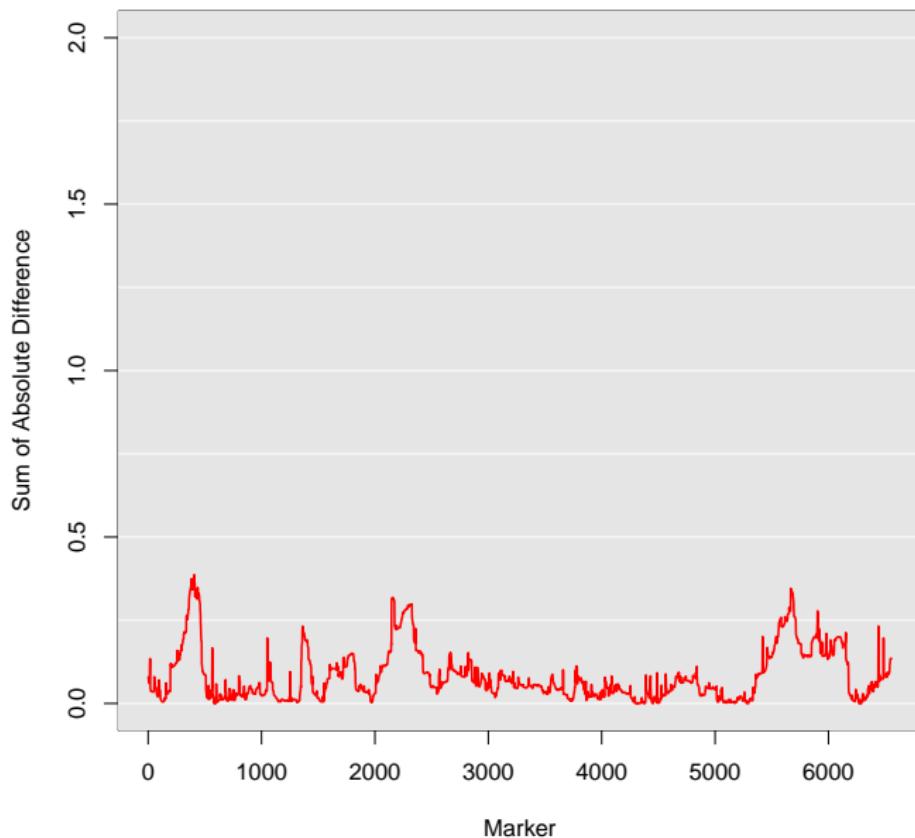
Each entry represents distance between methods

**Reduces problems to two-dimensions:**  $500 \times 120,000$

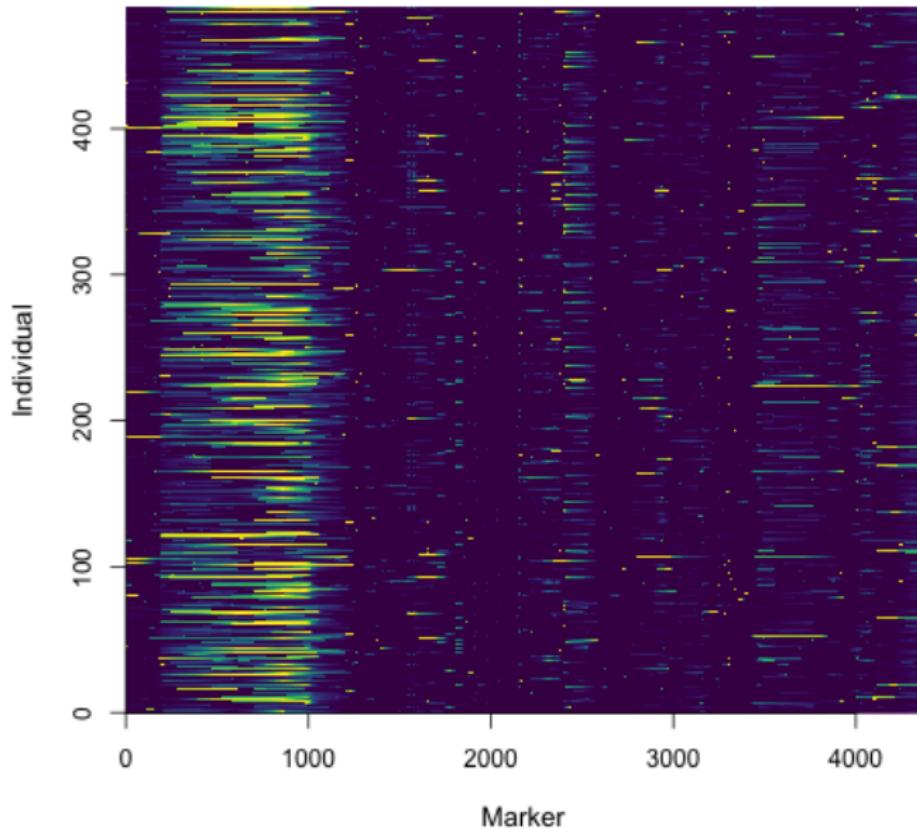
## Chromosome 5



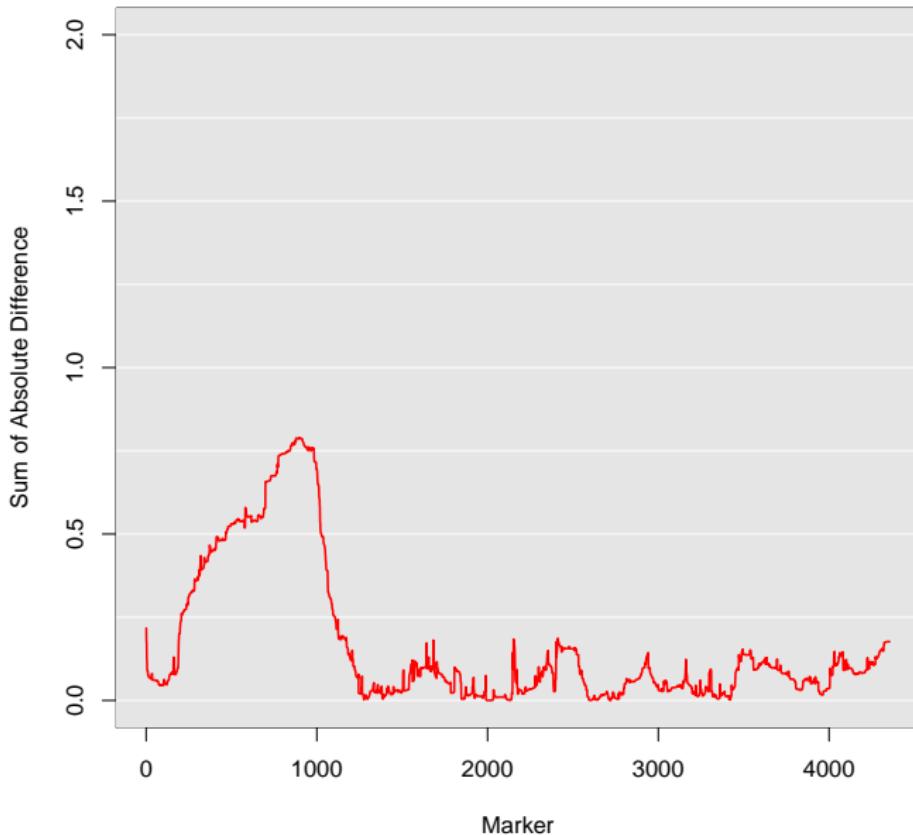
## Chromosome 5

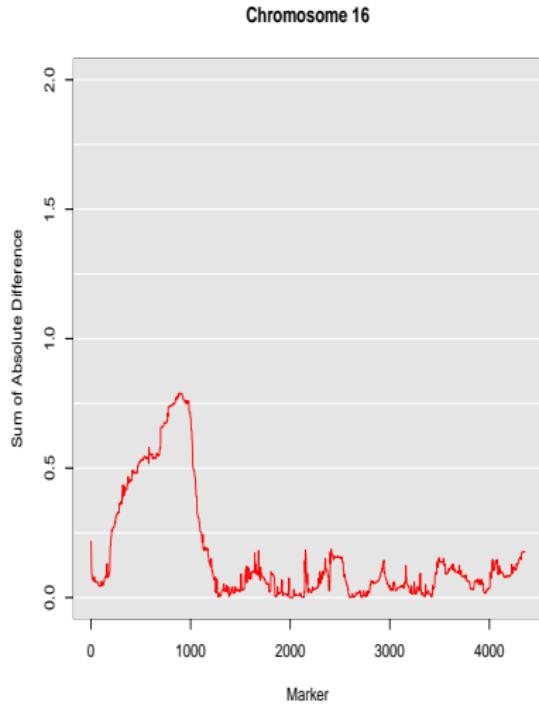
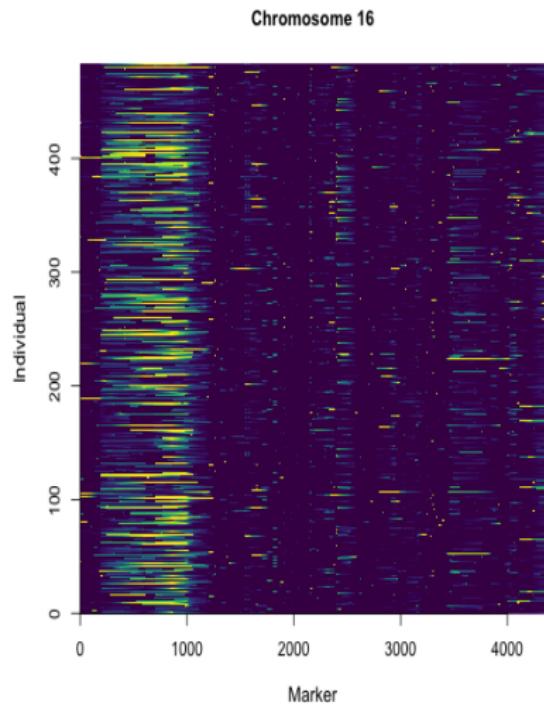


## Chromosome 16

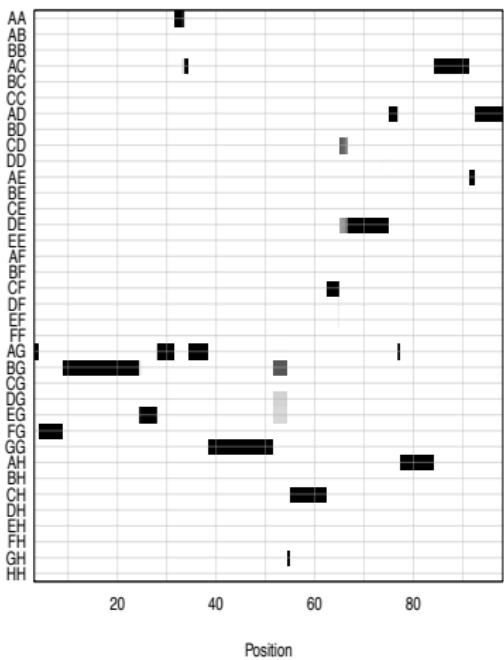


## Chromosome 16

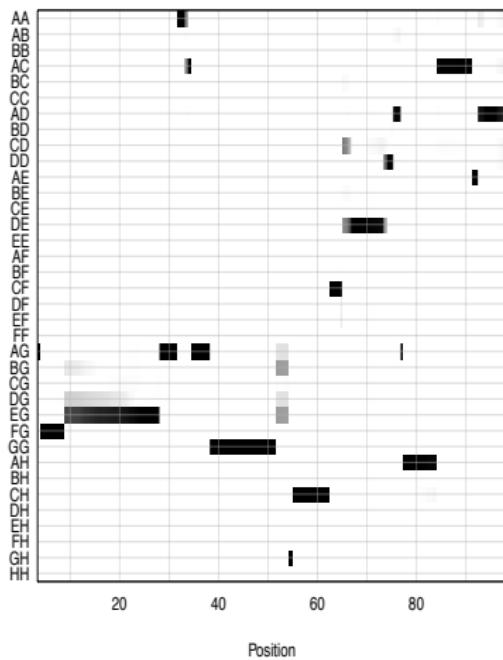




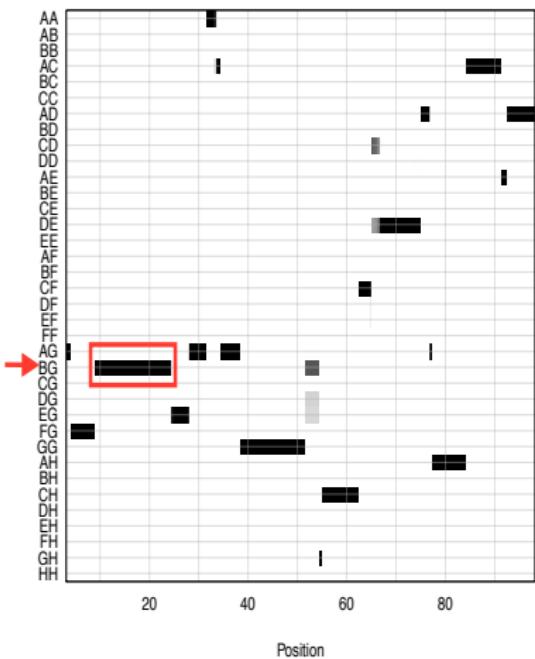
DOQTL DO-171



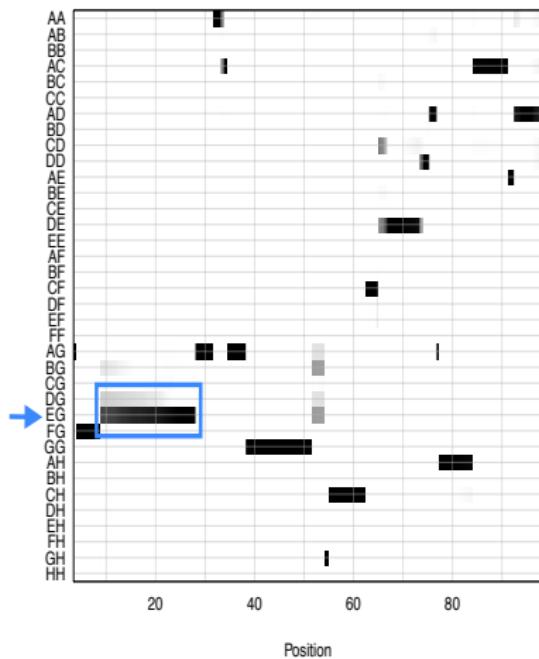
r/qtl2 DO-171



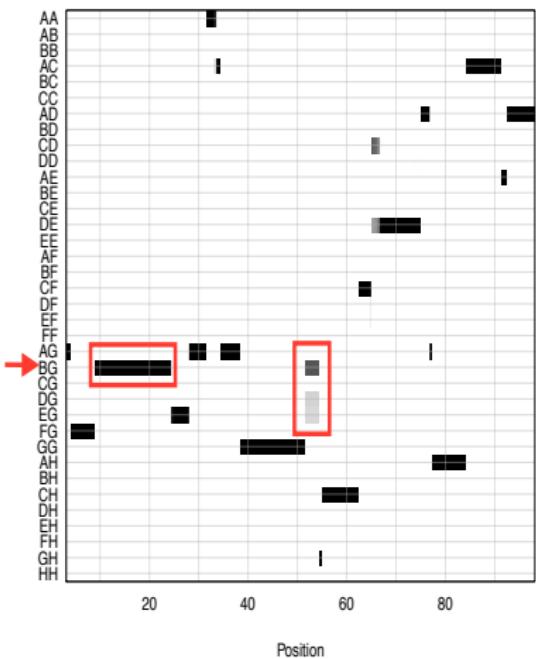
DOQTL DO-171



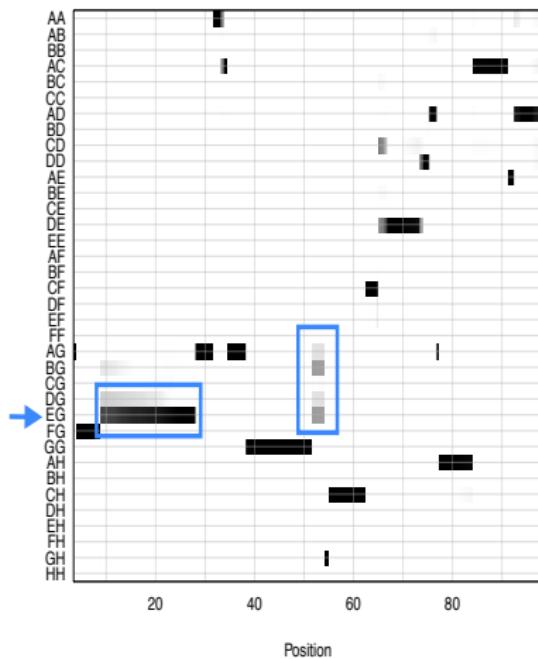
r/qtl2 DO-171

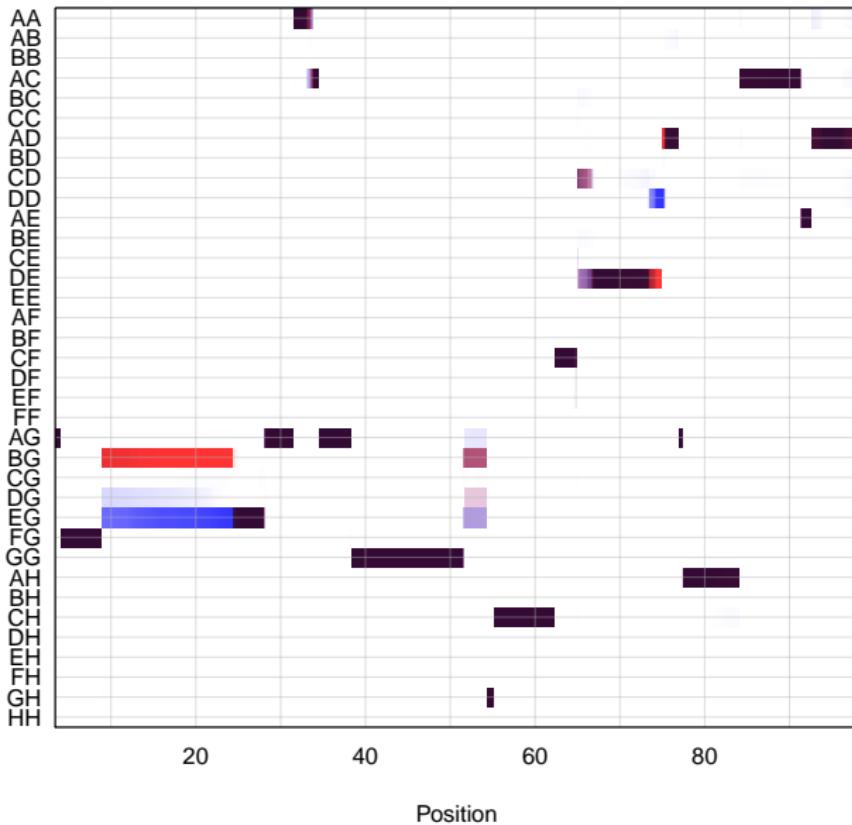


DOQTL DO-171



r/qtl2 DO-171





Which answer is 'correct'?

## Which answer is 'correct'?

- ▶ Pick an interval with differing, confident answers

## Which answer is 'correct'?

- ▶ Pick an interval with differing, confident answers
- ▶ Compare SNP genotypes with expected genotypes

## Which answer is 'correct'?

- ▶ Pick an interval with differing, confident answers
- ▶ Compare SNP genotypes with expected genotypes
- ▶ Compare SNP genotypes between discordant founders

Which answer is 'correct'?

Which answer is 'correct'?

**Individual:** DO-171      **Chromosome:** 16

# Which answer is 'correct'?

**Individual:** DO-171

**Chromosome:** 16

	<b>Prediction</b>	<b>Probability</b>	<b>Match Proportion</b>
r\qtl2	EG	0.892	0.9958848
doqtl	BG	0.999	0.9917695

# Which answer is 'correct'?

**Individual:** DO-171

**Chromosome:** 16

	<b>Prediction</b>	<b>Probability</b>	<b>Match Proportion</b>
r\qtl2	EG	0.892	0.9958848
doqtl	BG	0.999	0.9917695

**Best match:** EG

# Which answer is 'correct'?

**Individual:** DO-171

**Chromosome:** 16

	<b>Prediction</b>	<b>Probability</b>	<b>Match Proportion</b>
r\qtl2	EG	0.892	0.9958848
doqtl	BG	0.999	0.9917695

**Best match:** EG

Disagreement between **E** and **B**

# Which answer is 'correct'?

**Individual:** DO-171

**Chromosome:** 16

	<b>Prediction</b>	<b>Probability</b>	<b>Match Proportion</b>
<b>r\qtl2</b>	EG	0.892	0.9958848
<b>doqtl</b>	BG	0.999	0.9917695

**Best match:** EG

Disagreement between **E** and **B** → Proportion match = 0.997

# Which answer is 'correct'?

**Individual:** DO-171

**Chromosome:** 16

	<b>Prediction</b>	<b>Probability</b>	<b>Match Proportion</b>
<b>r\qtl2</b>	EG	0.892	0.9958848
<b>doqtl</b>	BG	0.999	0.9917695

**Best match:** EG

Disagreement between **E** and **B** → Proportion match = 0.997

We see this consistently across the marker region

# Overall Results

# Overall Results

**Goal:** Compare large 3D arrays

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes
  - Heat map of individuals × markers

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes
  - Heat map of individuals × markers
  - Plot average distance by marker position

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes
  - Heat map of individuals × markers
  - Plot average distance by marker position
- ▶ Visually compare probabilities *at individual level*

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes
  - Heat map of individuals × markers
  - Plot average distance by marker position
- ▶ Visually compare probabilities *at individual level*
  - Univariate probability heatmap

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes
  - Heat map of individuals × markers
  - Plot average distance by marker position
- ▶ Visually compare probabilities *at individual level*
  - Univariate probability heatmap
  - Bivariate probability heatmap

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes
  - Heat map of individuals × markers
  - Plot average distance by marker position
- ▶ Visually compare probabilities *at individual level*
  - Univariate probability heatmap
  - Bivariate probability heatmap
- ▶ Compare observed SNP genotypes to expected

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes
  - Heat map of individuals × markers
  - Plot average distance by marker position
- ▶ Visually compare probabilities *at individual level*
  - Univariate probability heatmap
  - Bivariate probability heatmap
- ▶ Compare observed SNP genotypes to expected
  - r\qtl2 consistently gives better answers

# Overall Results

**Goal:** Compare large 3D arrays

- ▶ Measure of distance to identify problem regions
  - Sum of absolute differences across haplotypes
  - Heat map of individuals × markers
  - Plot average distance by marker position
- ▶ Visually compare probabilities *at individual level*
  - Univariate probability heatmap
  - Bivariate probability heatmap
- ▶ Compare observed SNP genotypes to expected
  - r\qtl2 consistently gives better answers
  - Disagreement often related to similarity of founders