

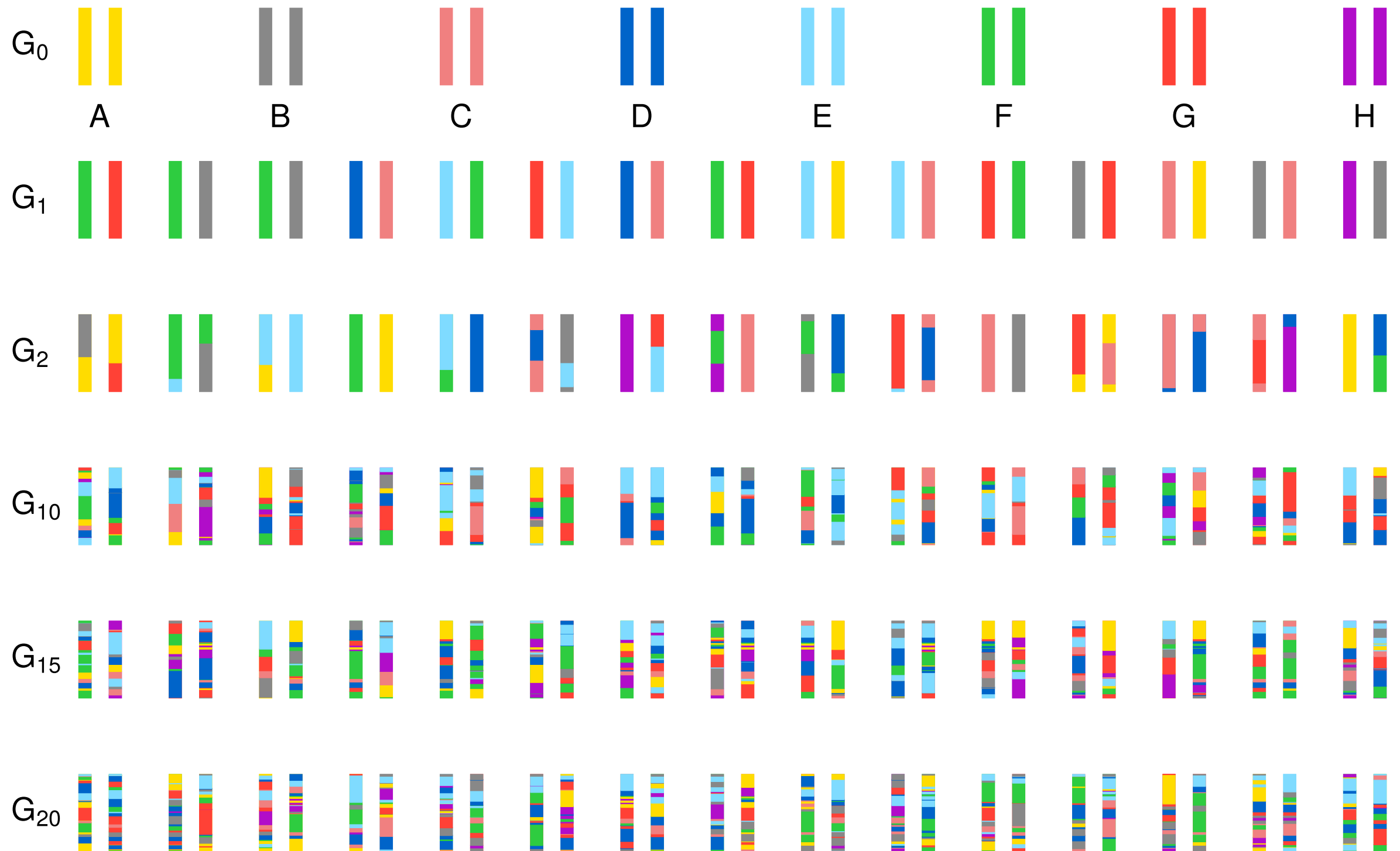
Cleaning genotype data for diversity outbred mice

Karl Broman

Biostatistics & Medical Informatics
University of Wisconsin–Madison

kbroman.org
github.com/kbroman
[@kwbroman](https://twitter.com/kwbroman)
Slides: bit.ly/jax18

Multi-parent advanced intercross



Diversity outbred mouse data

- 500 DO mice
- GigaMUGA SNP arrays (114k SNPs)
- RNA-seq data on pancreatic islets
- Microbiome data (16S and shotgun sequencing)
- protein and lipid measurements by mass spec
- Collaboration with Alan Attie, Gary Churchill, Brian Yandell, Josh Coon, Federico Rey, and many others

Principles

What might have gone wrong?

How could it be revealed?

Principles

What might have gone wrong?

How could it be revealed?

Also, just make a bunch of graphs.

Principles

What might have gone wrong?

How could it be revealed?

Also, just make a bunch of graphs.

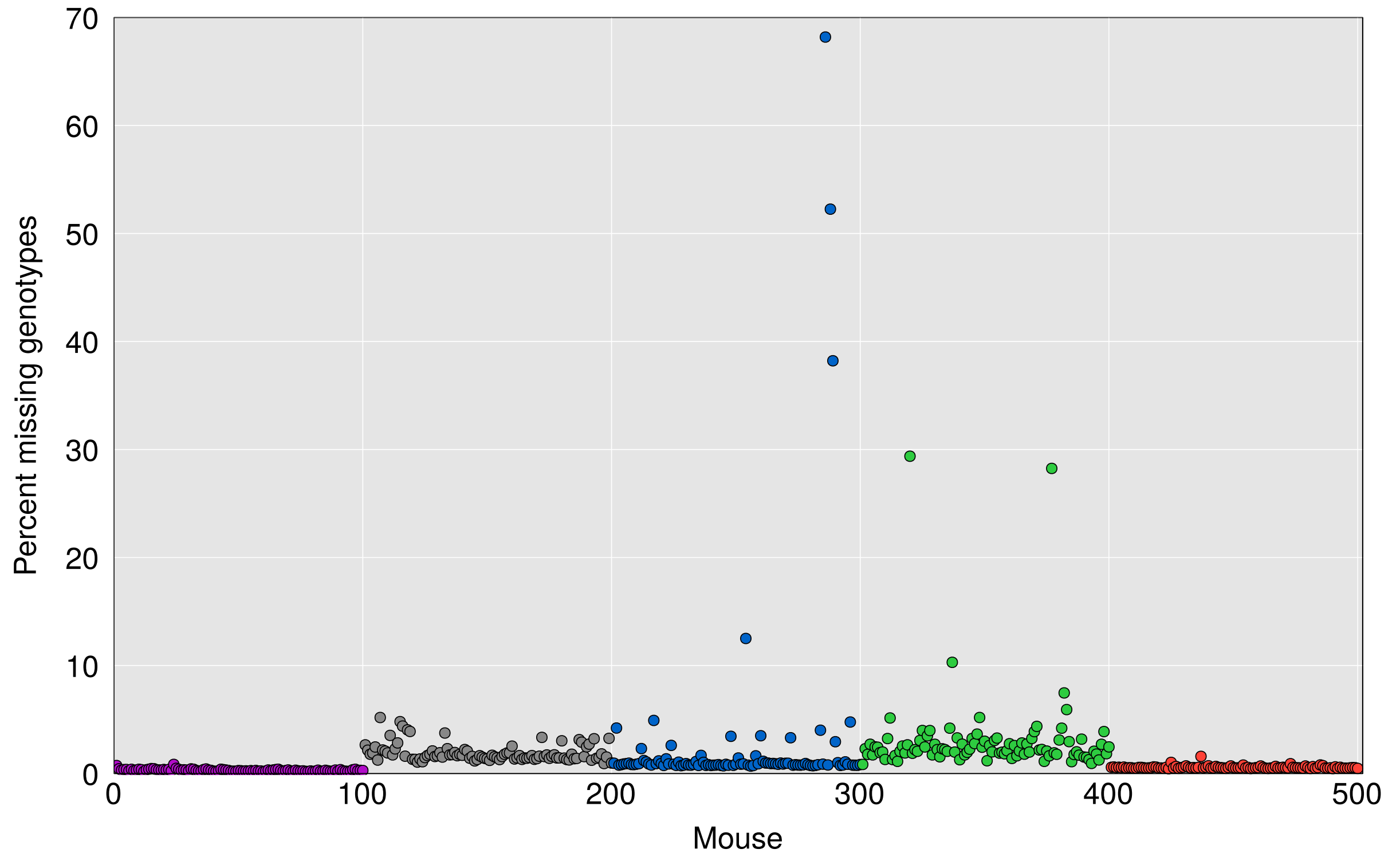
If you see something weird, try to figure it out.

Possible problems

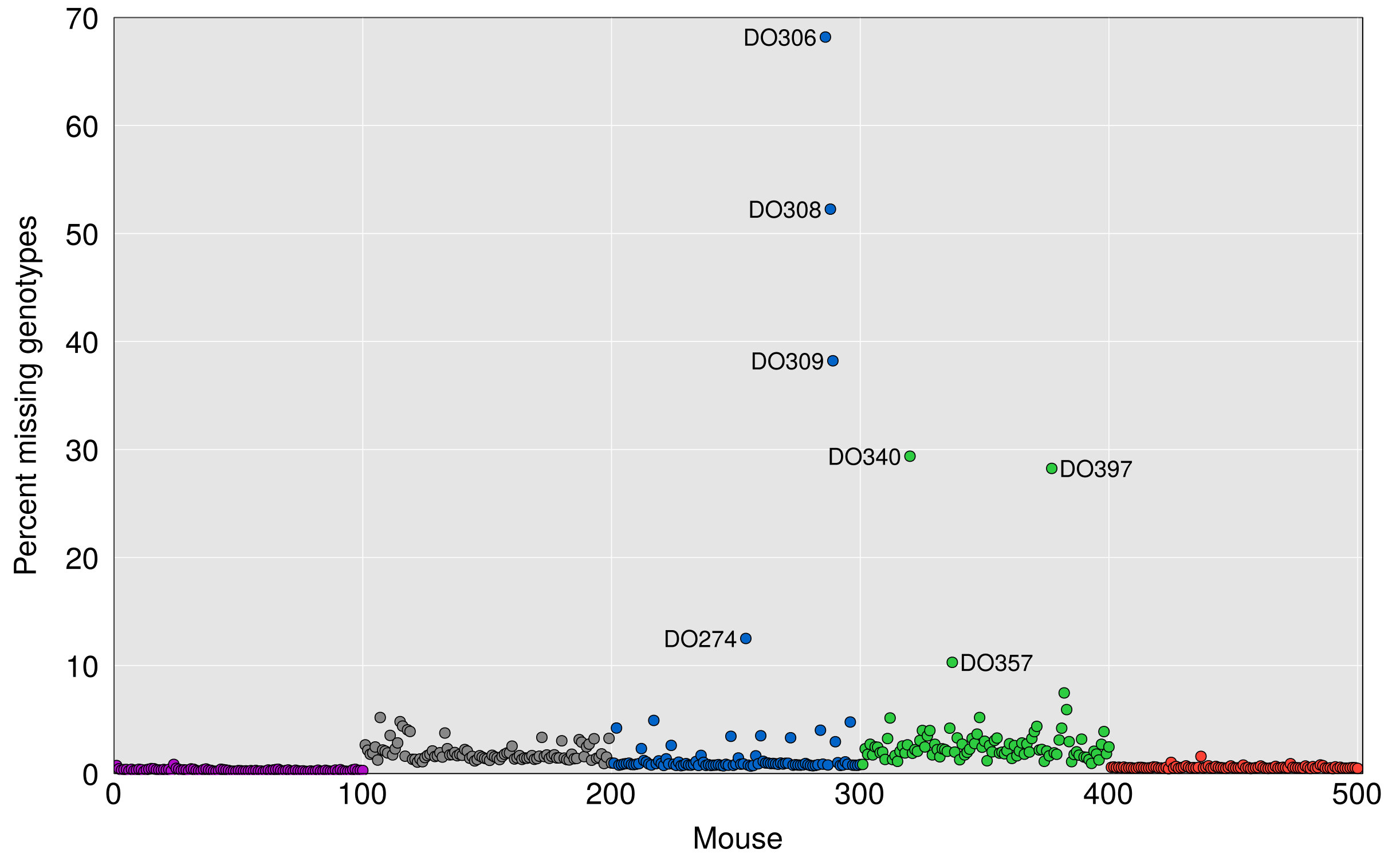
- Sample duplicates
- Sample mix-ups
- Bad samples
- Bad markers
- Genotyping errors in founders

What to look at first?

Missing data per sample

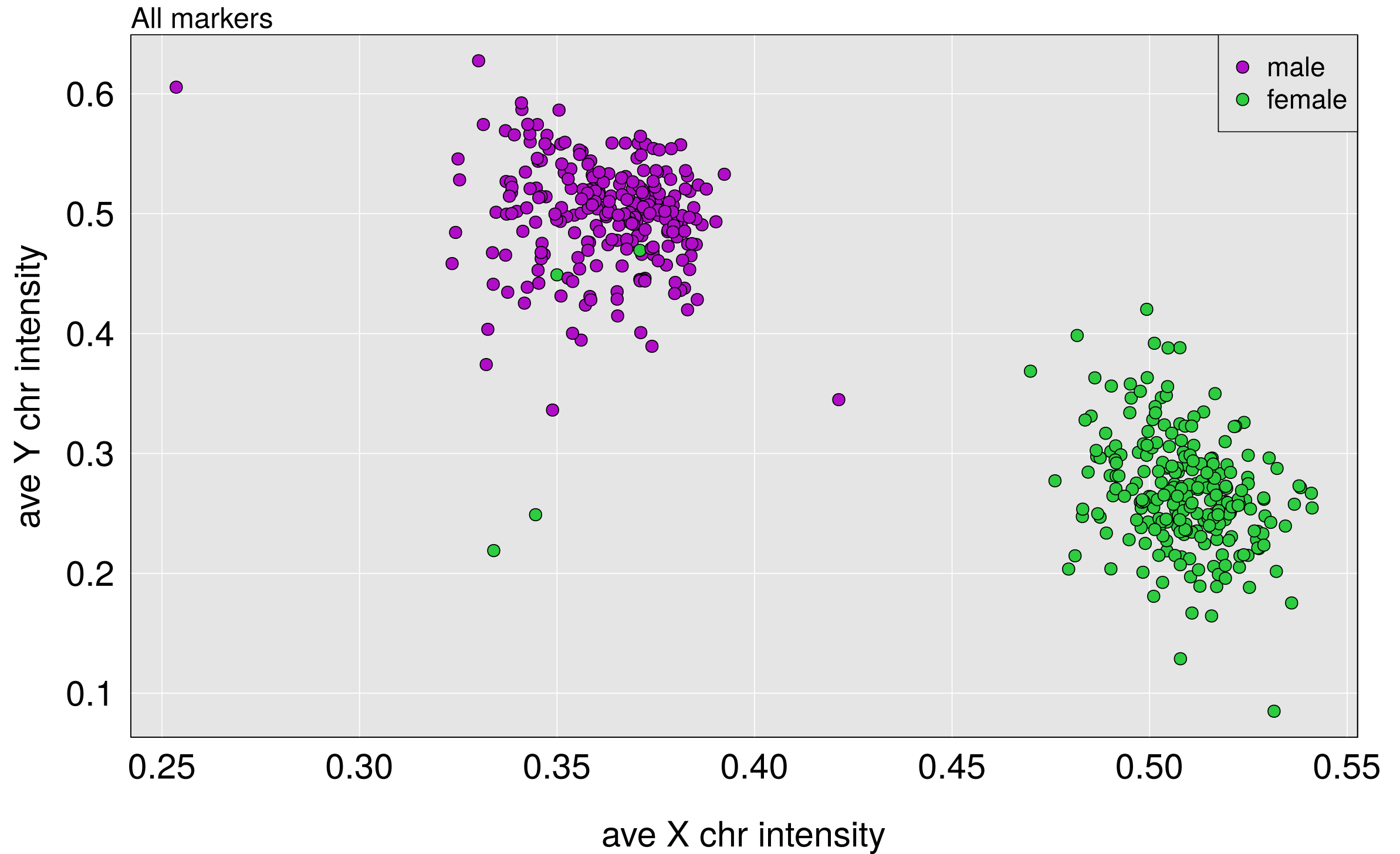


Missing data per sample

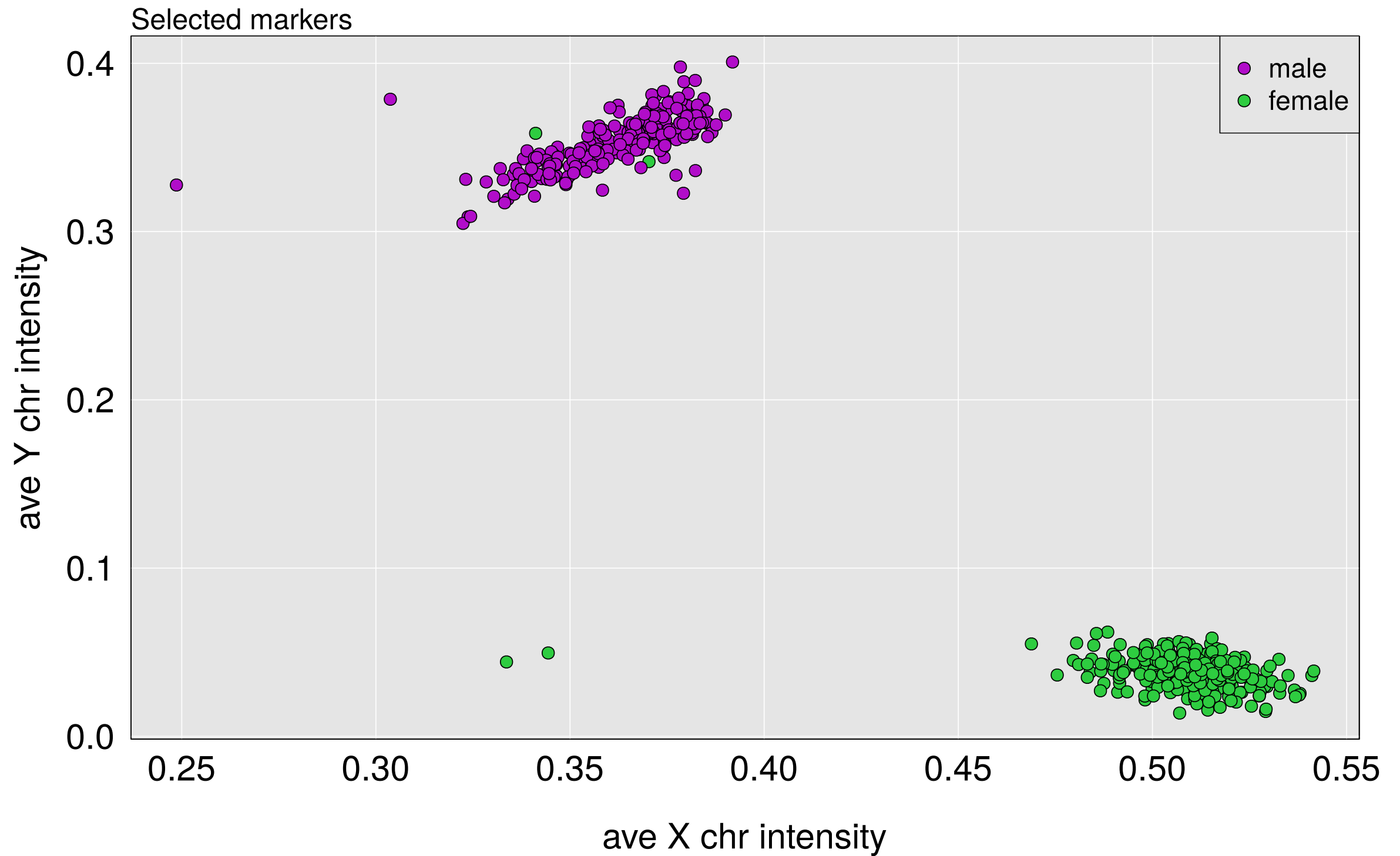


Swapped sex labels

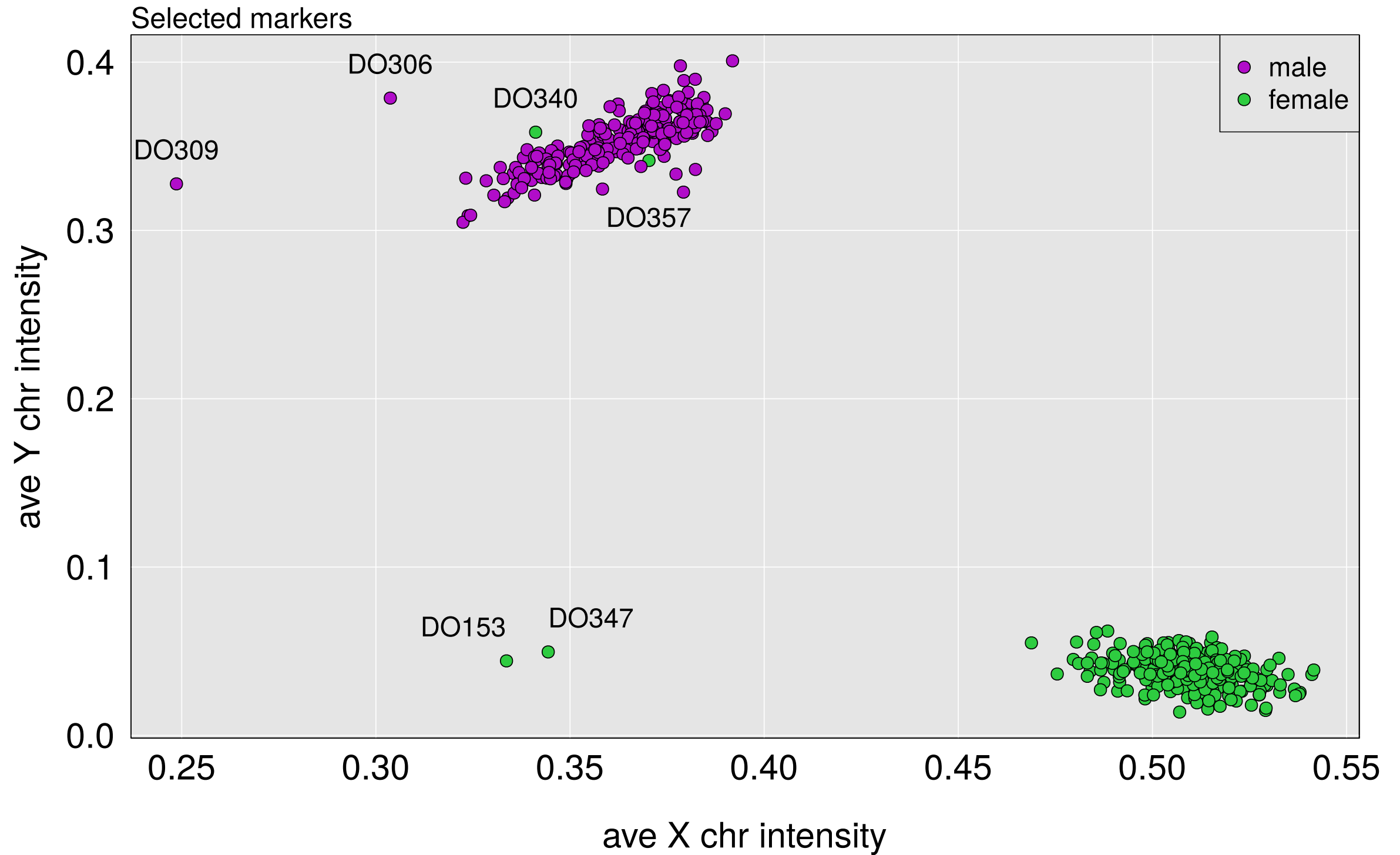
Average SNP intensity on X and Y chr



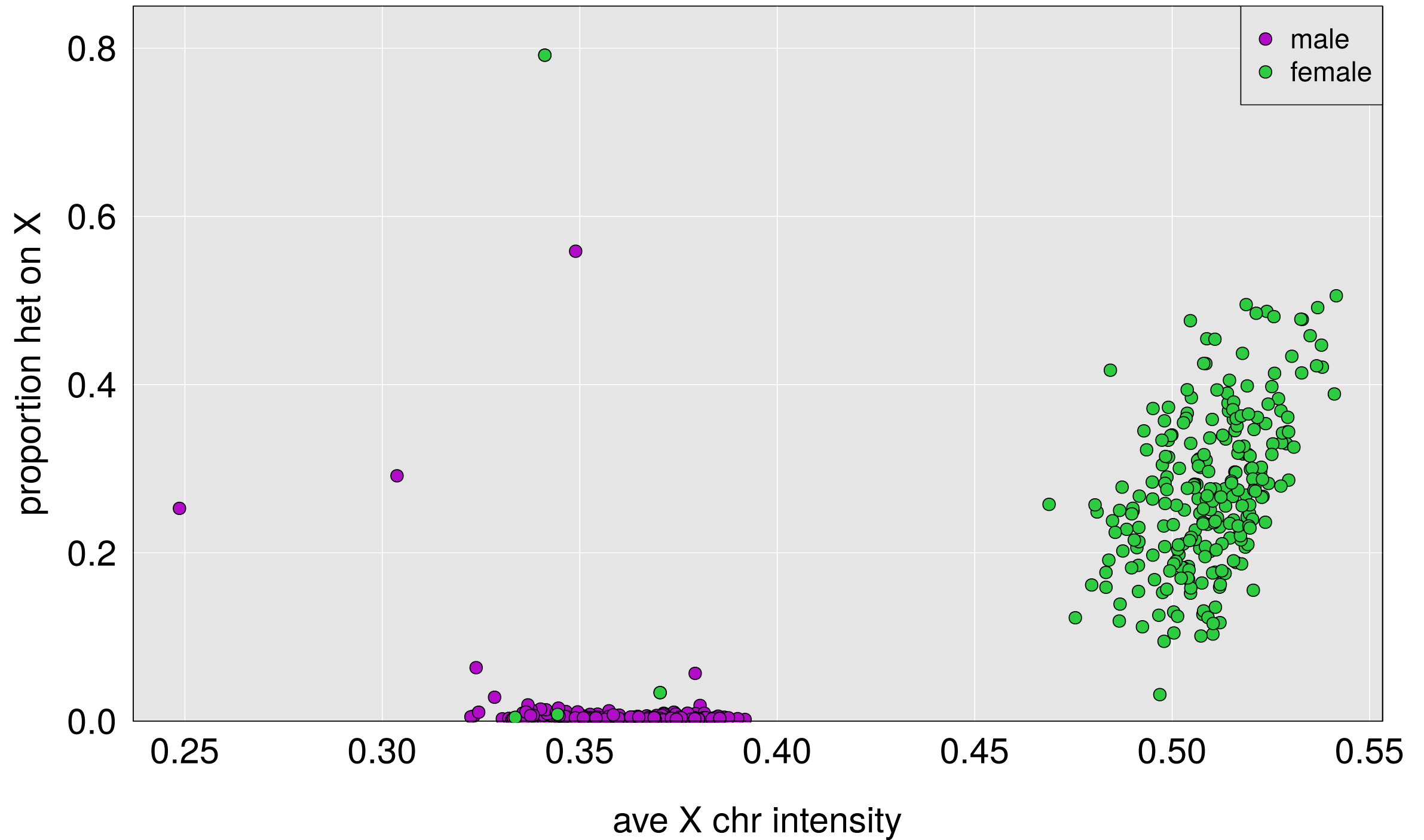
Average SNP intensity on X and Y chr



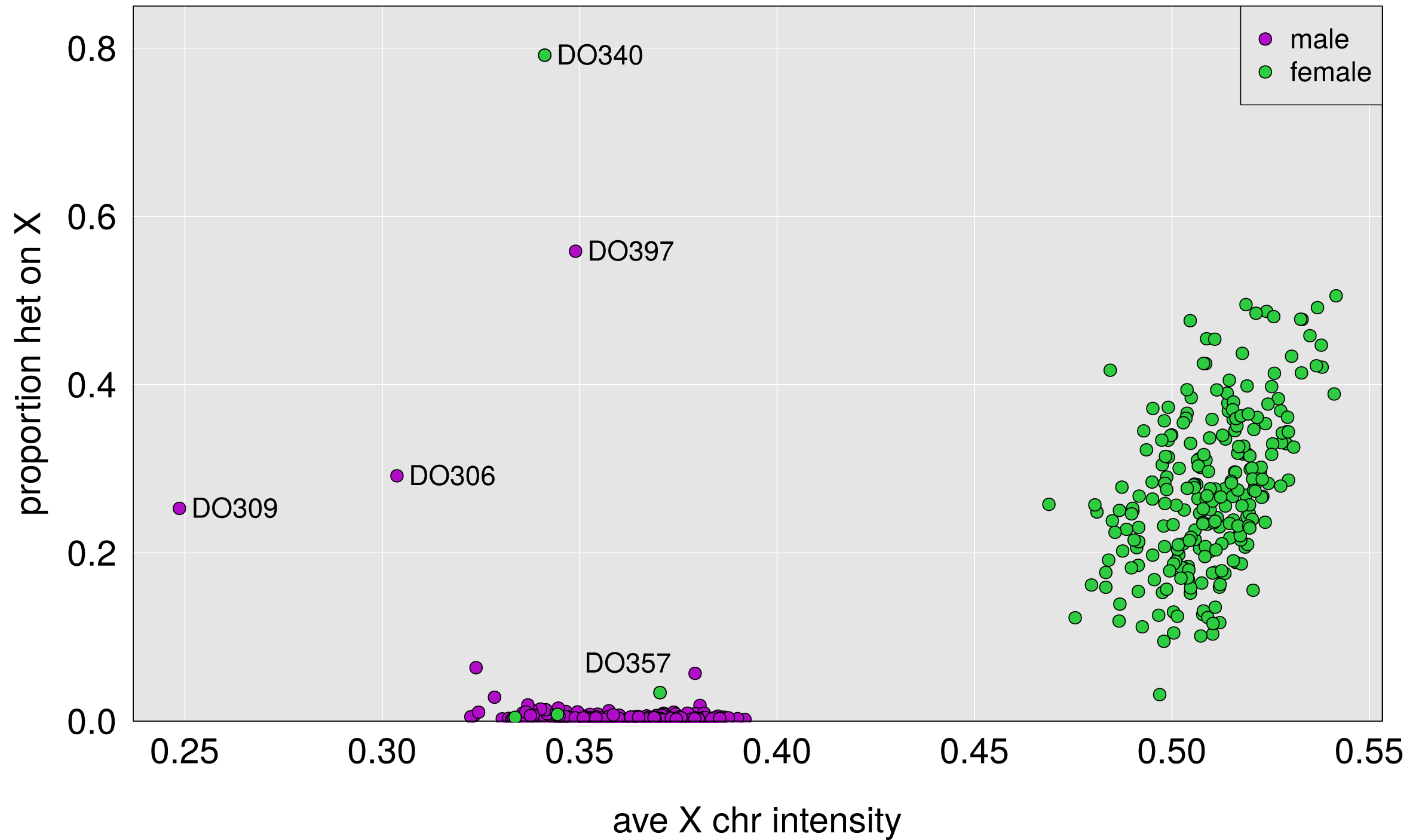
Average SNP intensity on X and Y chr



Heterozygosity vs SNP intensity on X chr

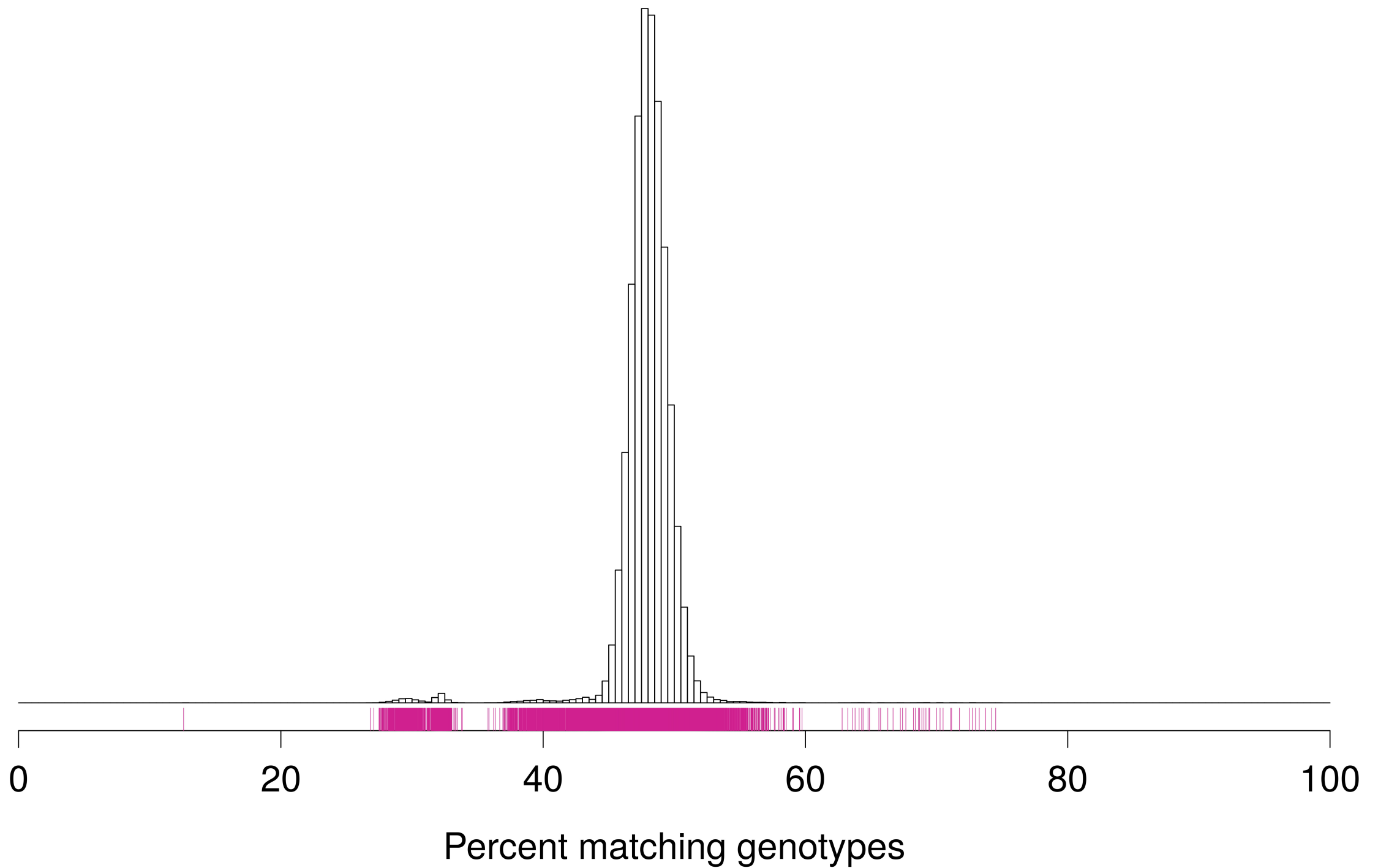


Heterozygosity vs SNP intensity on X chr

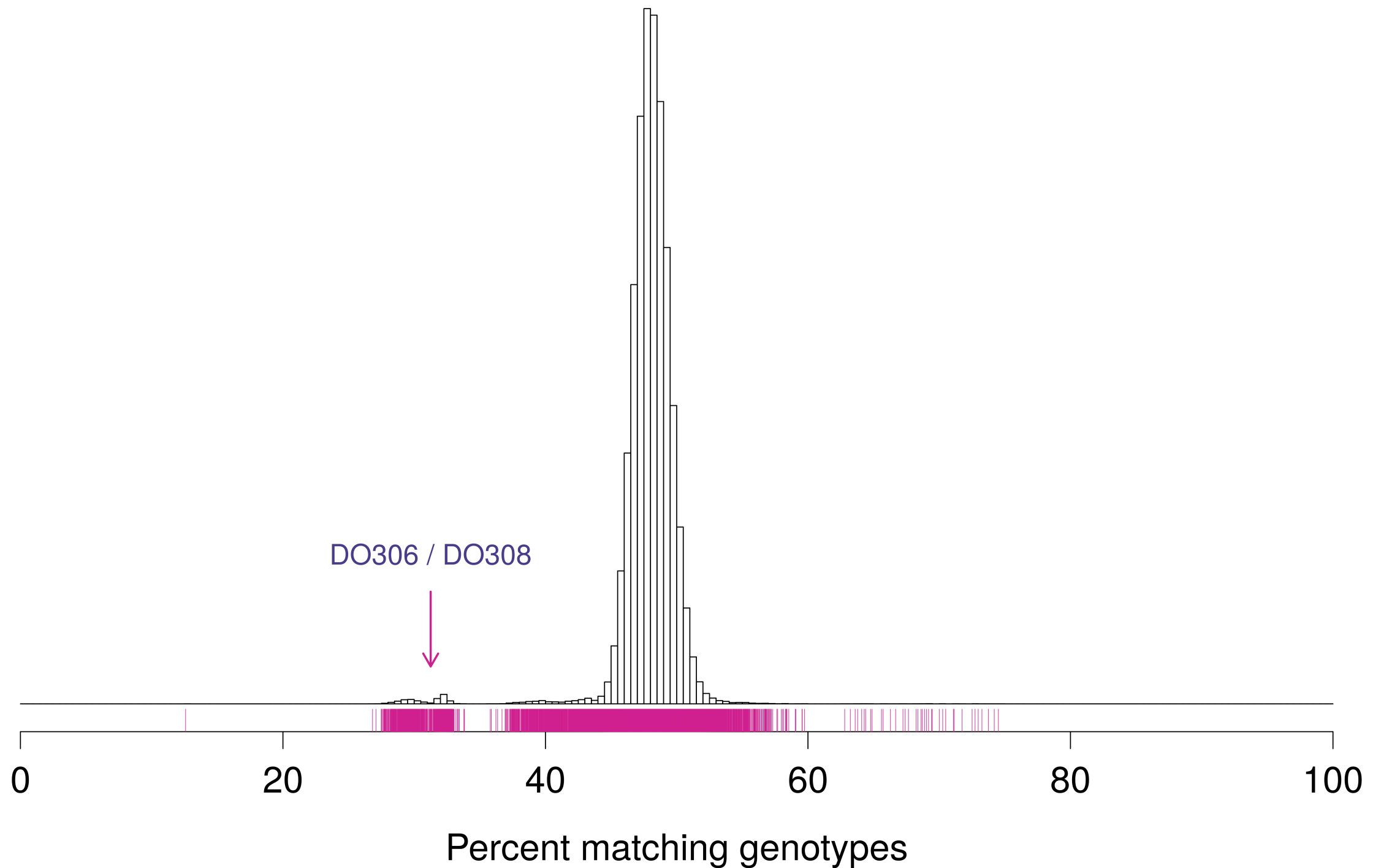


Sample duplicates

Percent matching genotypes between pairs

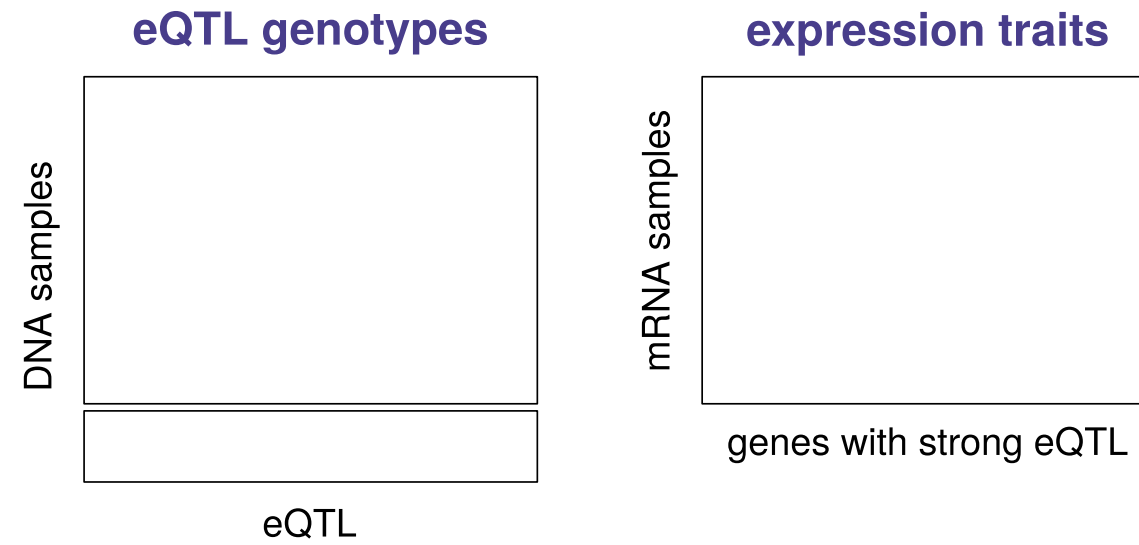


Percent matching genotypes between pairs

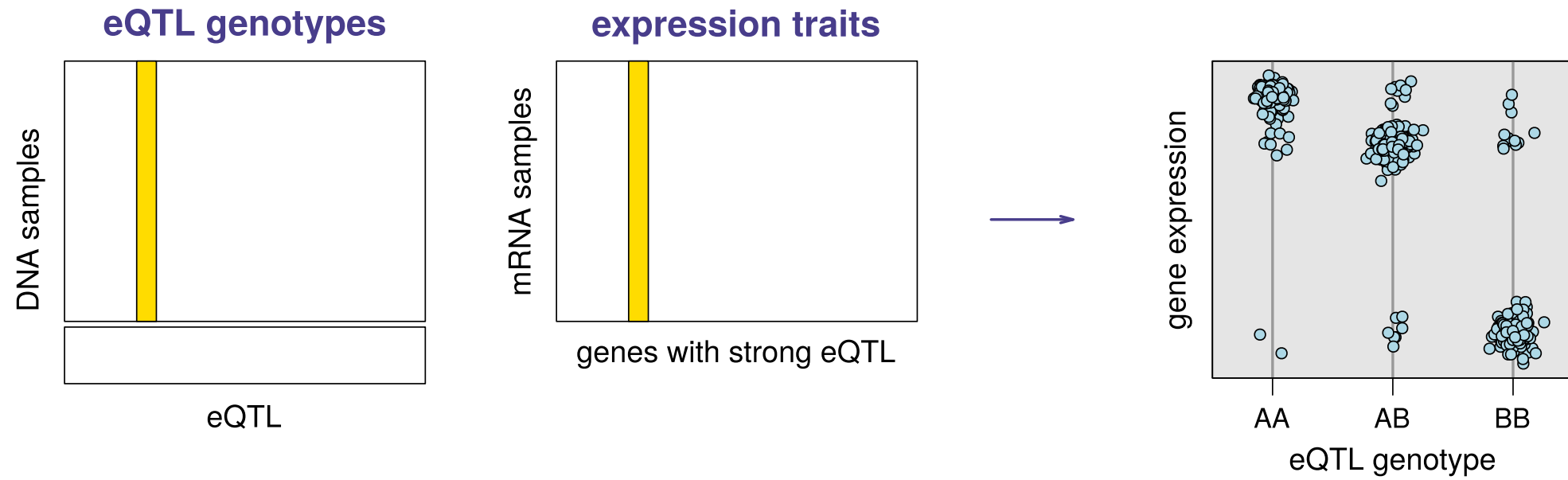


Sample mix-ups

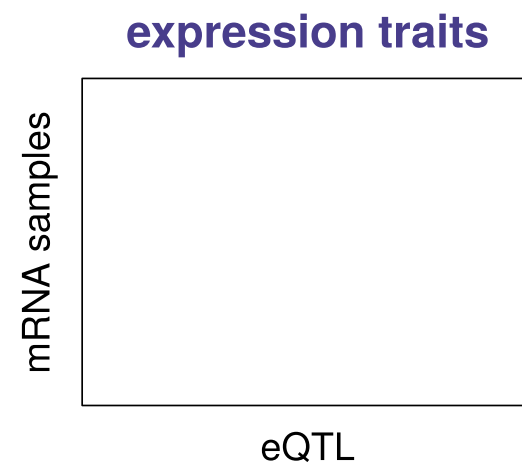
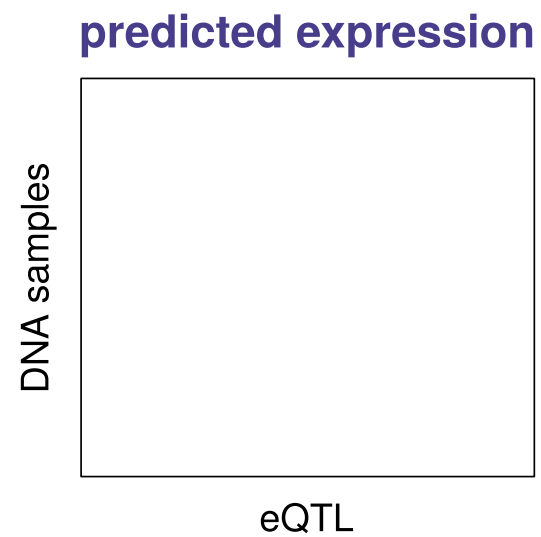
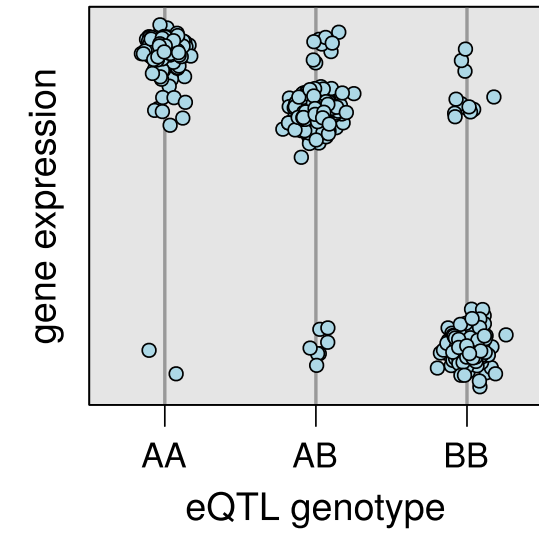
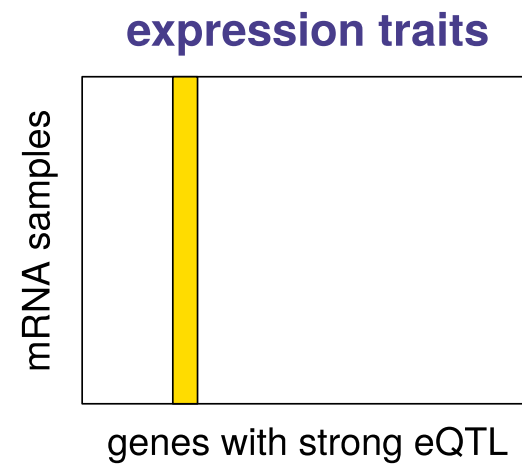
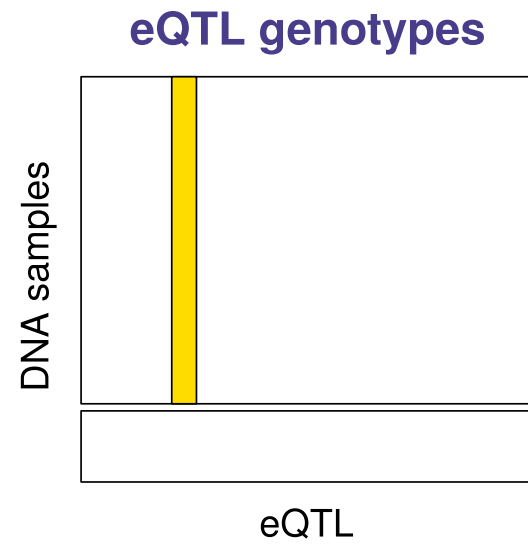
Sample mix-ups



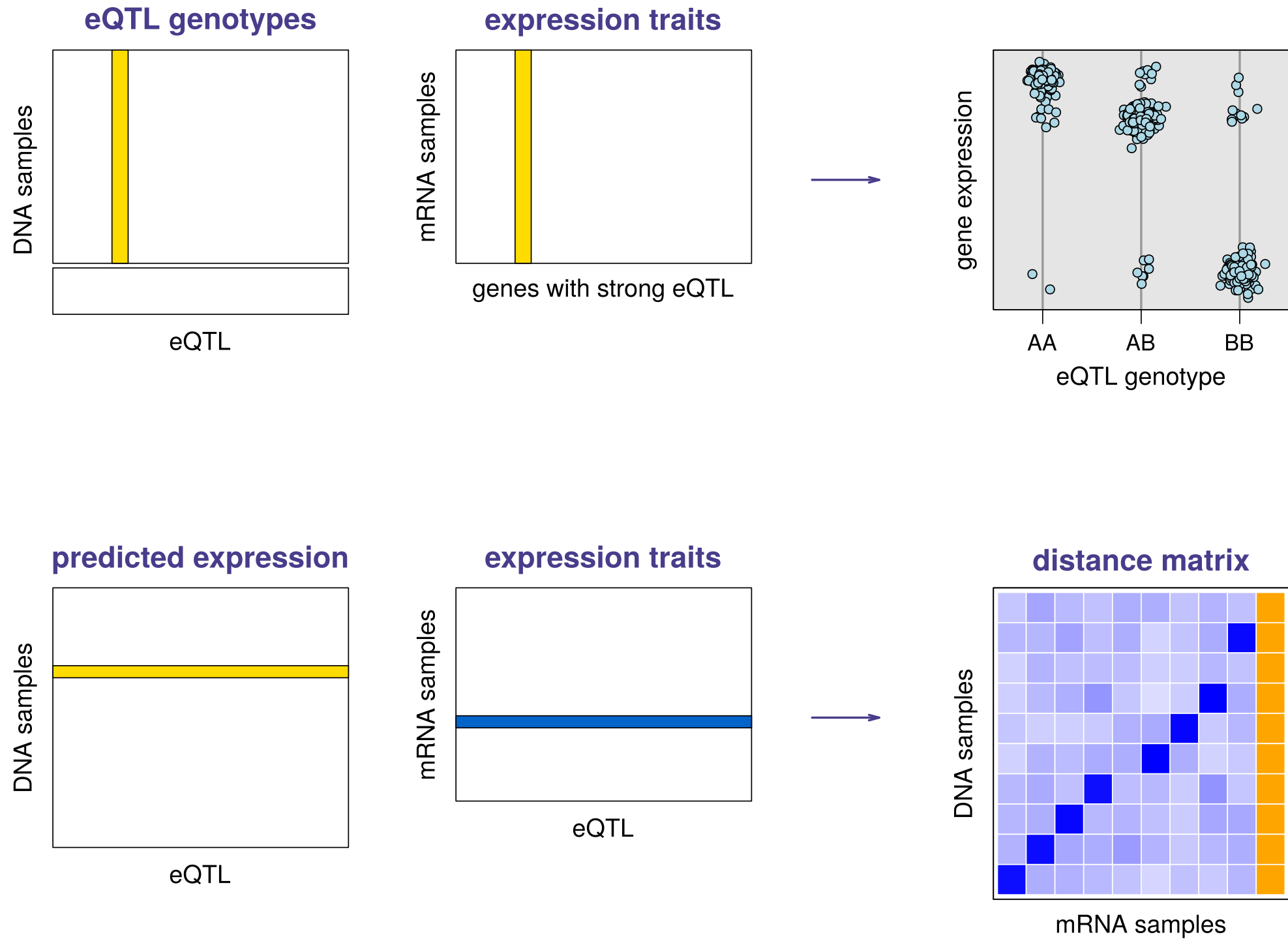
Sample mix-ups



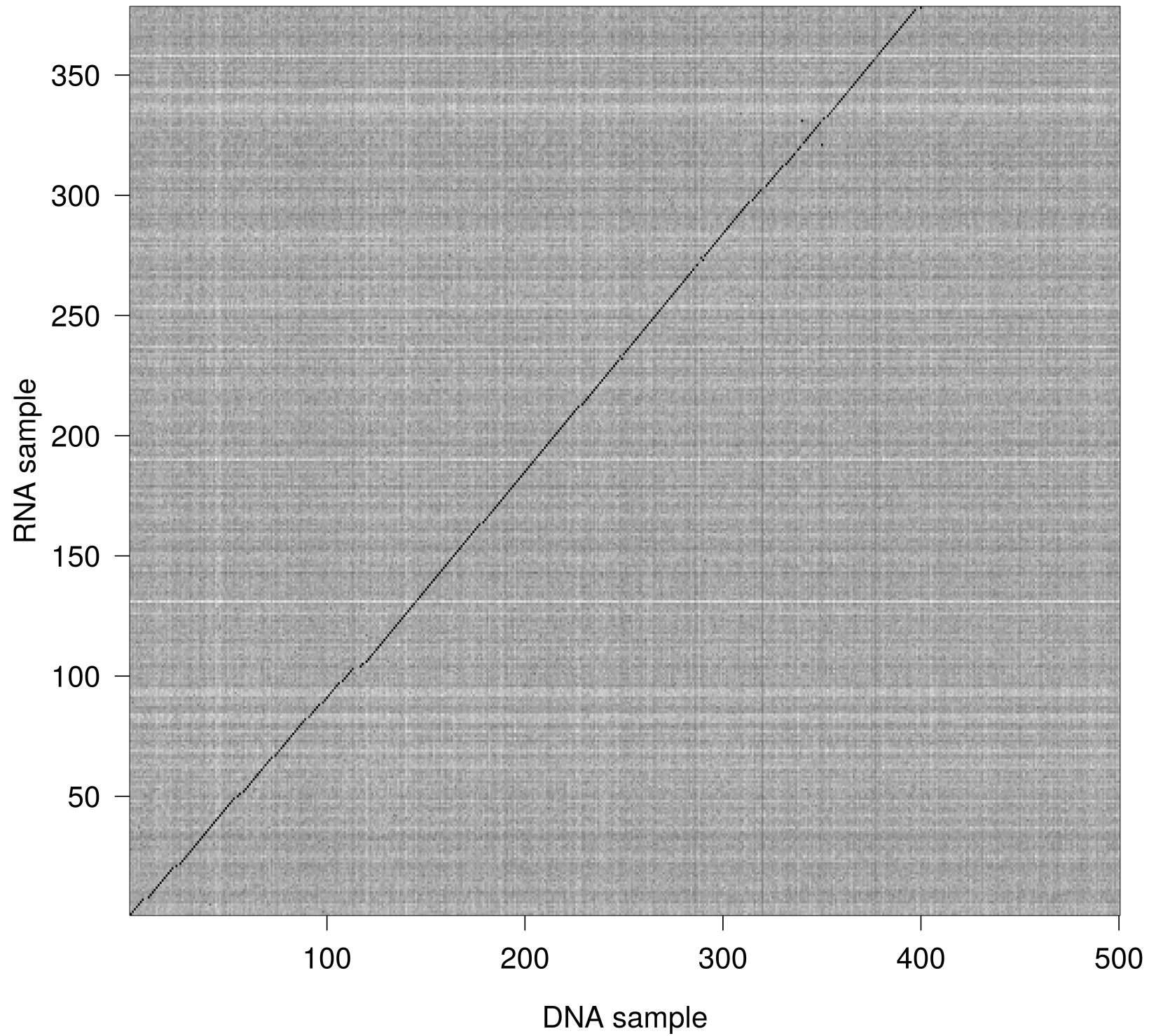
Sample mix-ups



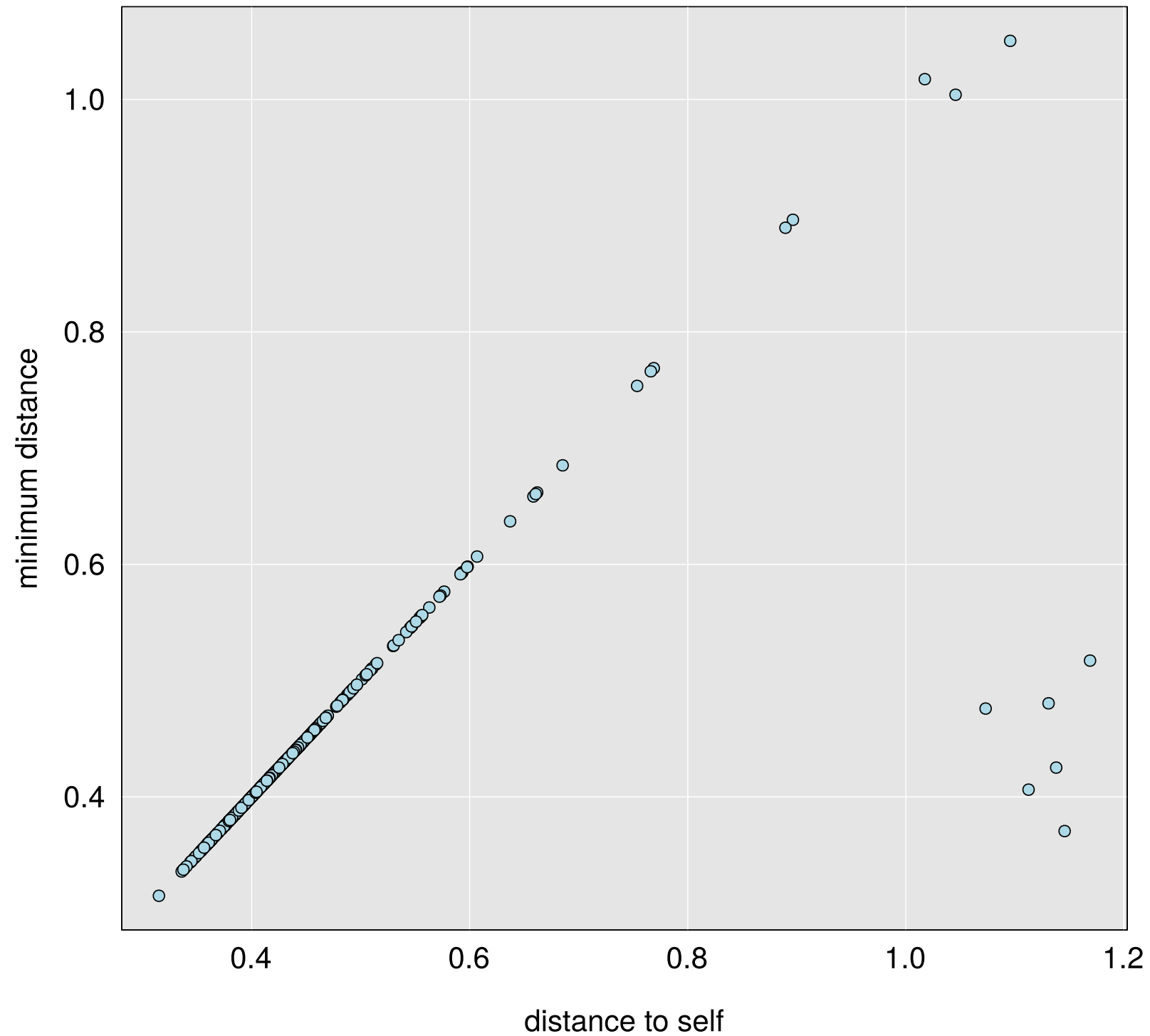
Sample mix-ups



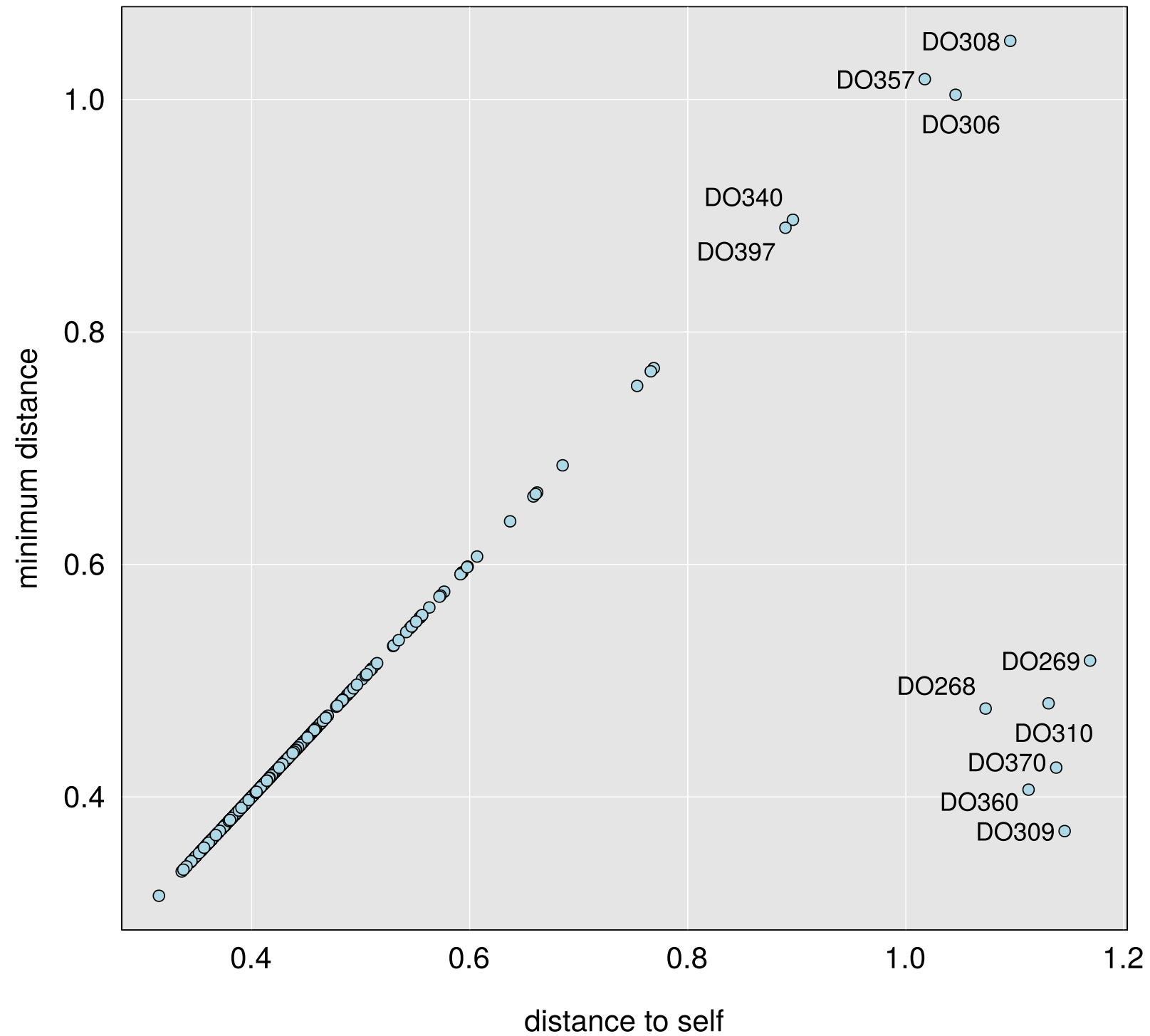
RNA-seq sample mix-ups: distance matrix



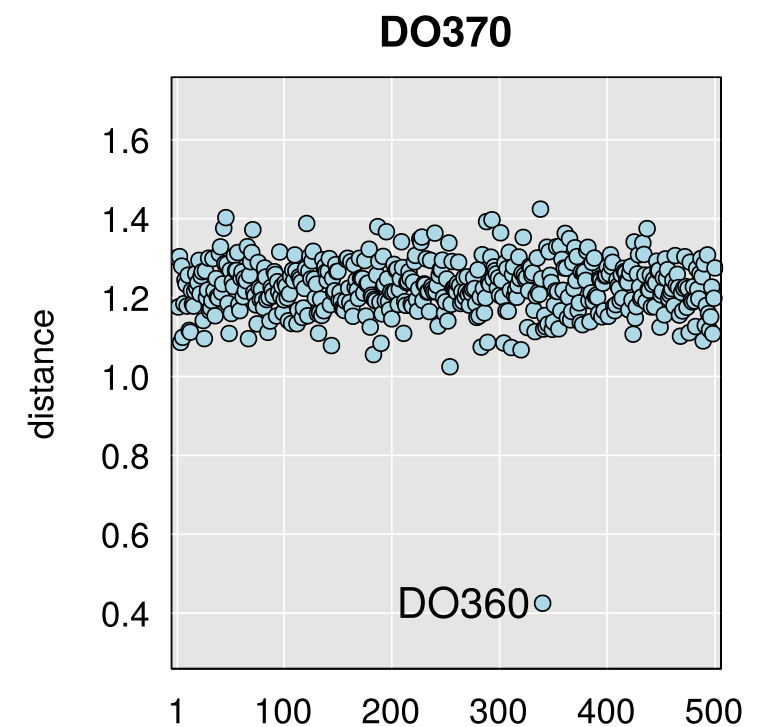
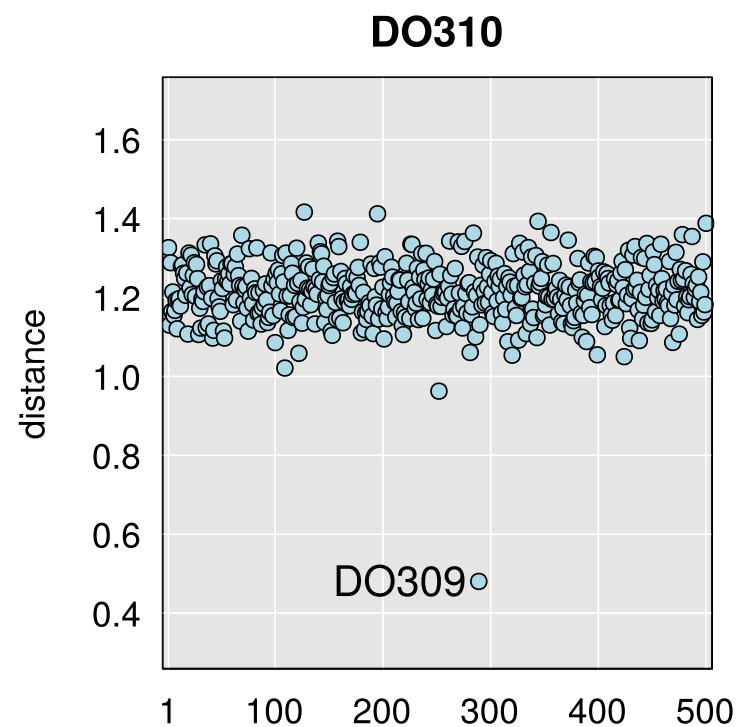
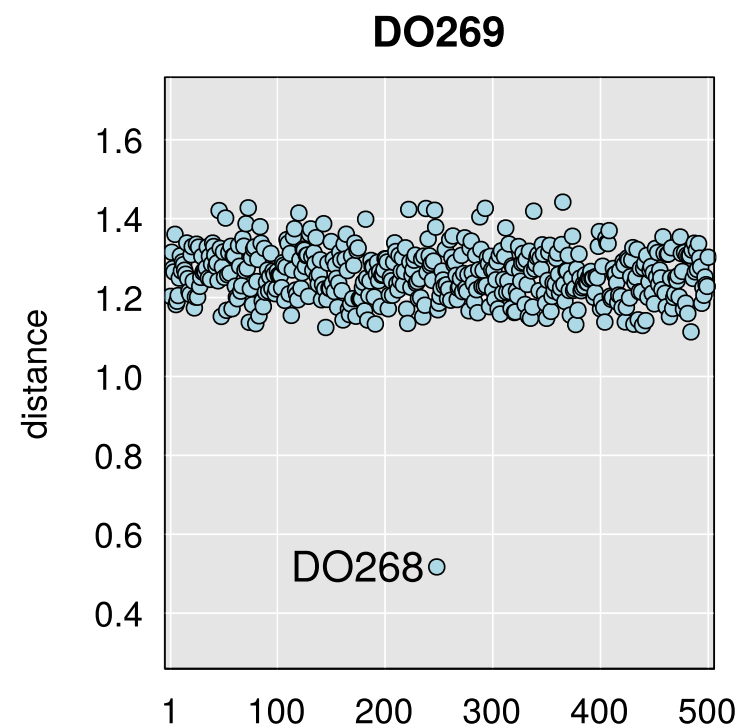
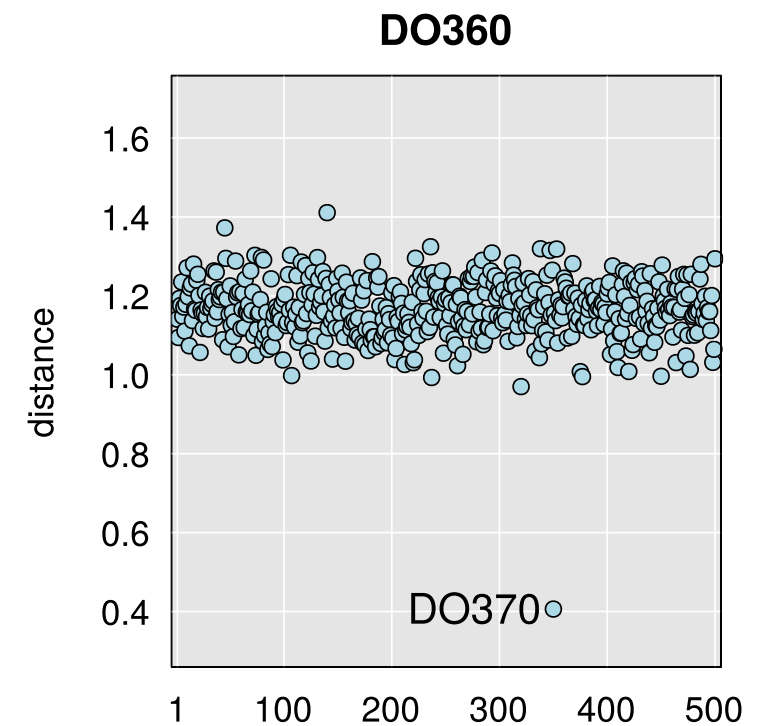
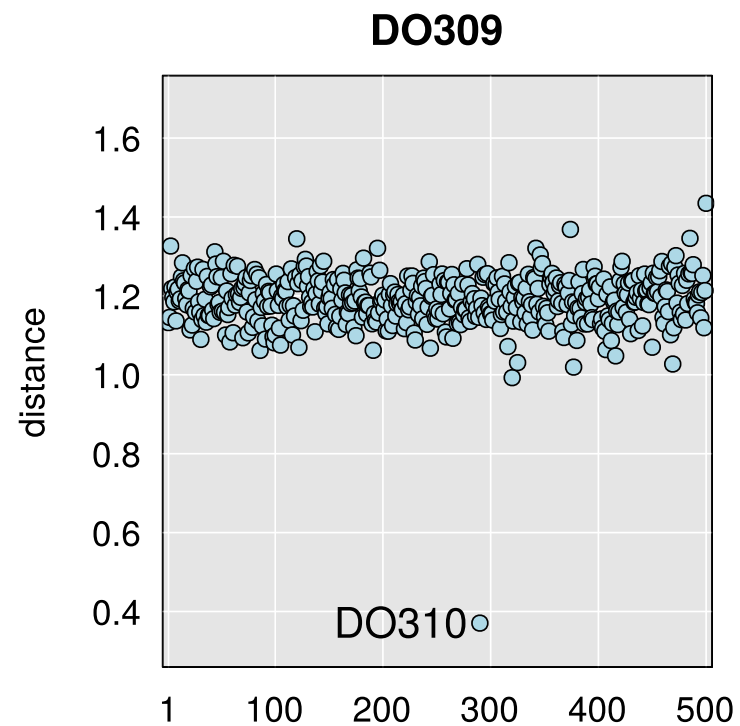
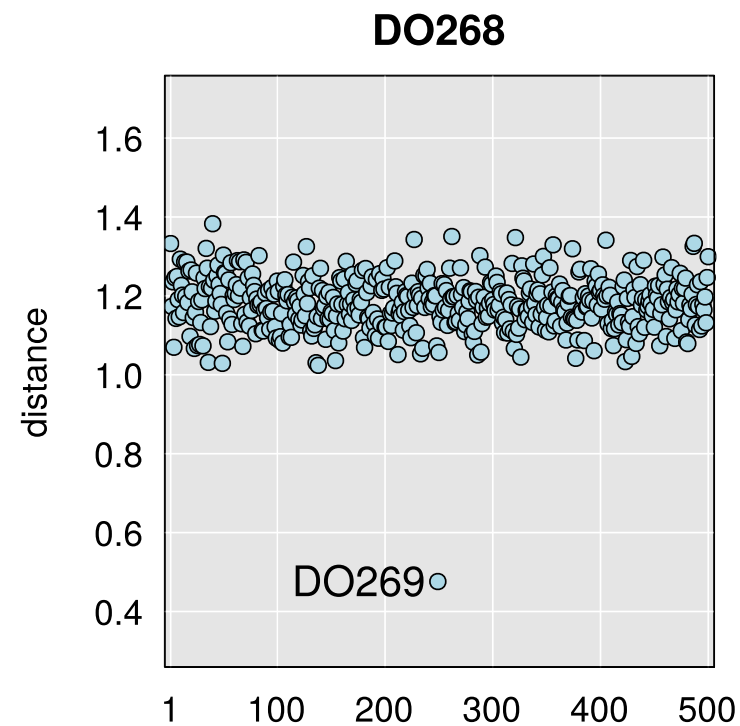
RNA-seq sample mix-ups: min vs self distance



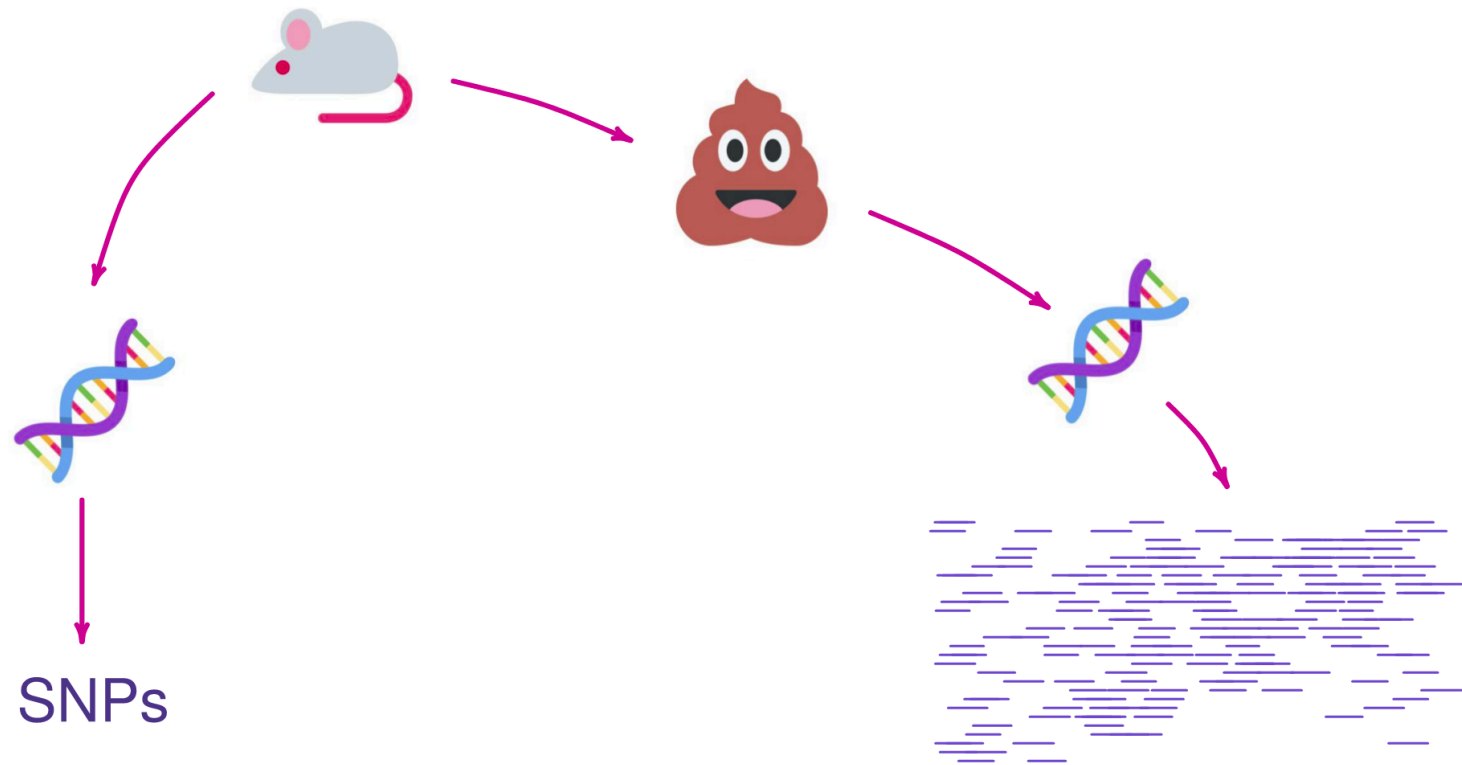
RNA-seq sample mix-ups: min vs self distance



RNA-seq sample mix-ups: detail



Microbiome data



Sample mix-ups: Microbiome data

- Impute genotypes at all SNPs in DNA samples
- Map microbiome reads to mouse genome; find reads overlapping a SNP
- For each pair of samples (DNA + microbiome):
 - Focus on reads that overlap a SNP where that DNA sample is homozygous
 - Distance = proportion of reads where SNP allele doesn't match DNA sample's genotype

Genomic DO361 vs Microbiome DO361

	A	B
AA	939,918	2,998
BB	1,044	125,962

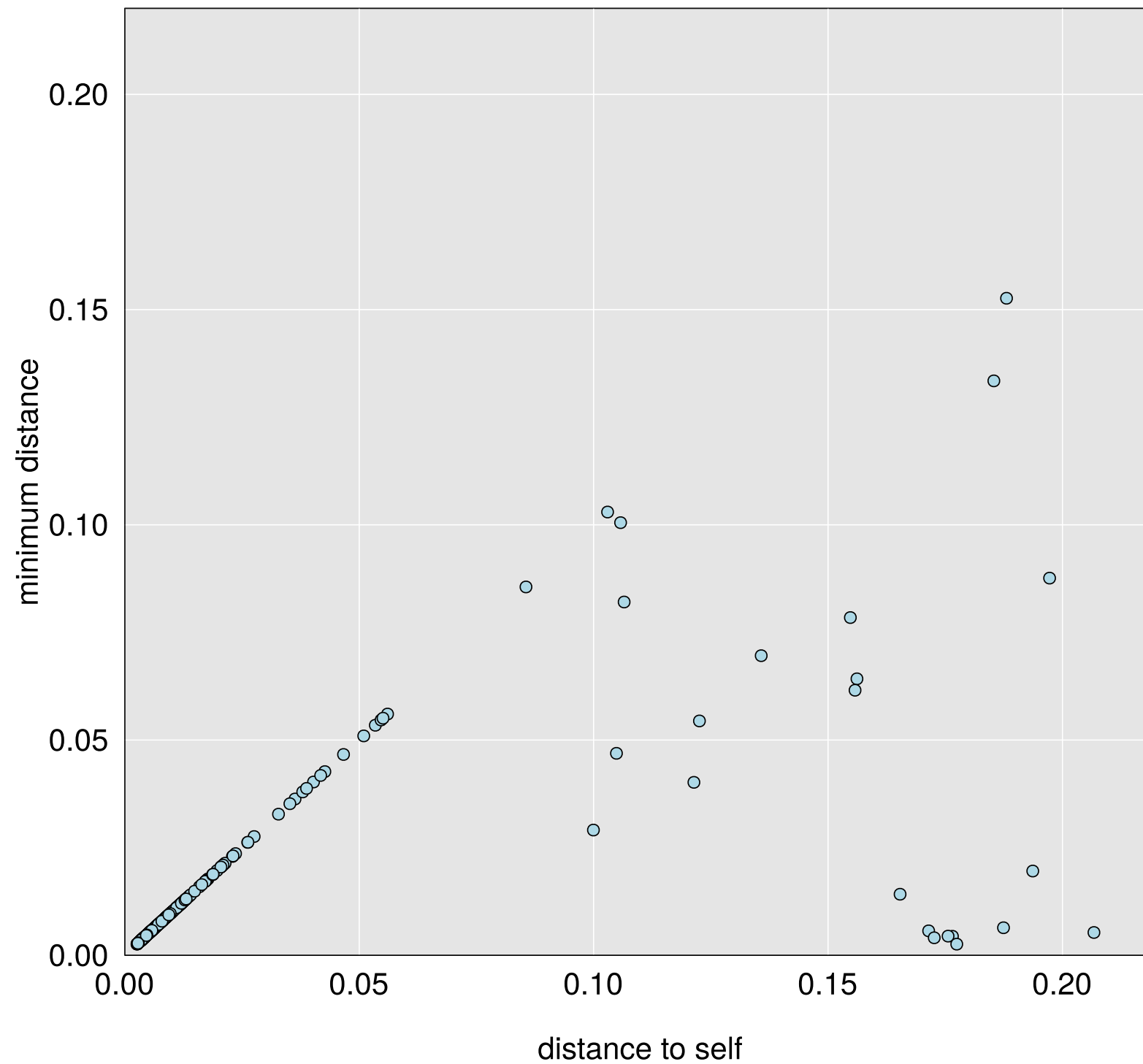
Genomic DO360 vs Microbiome DO360

	A	B
AA	2,661,645	427,685
BB	190,188	202,335

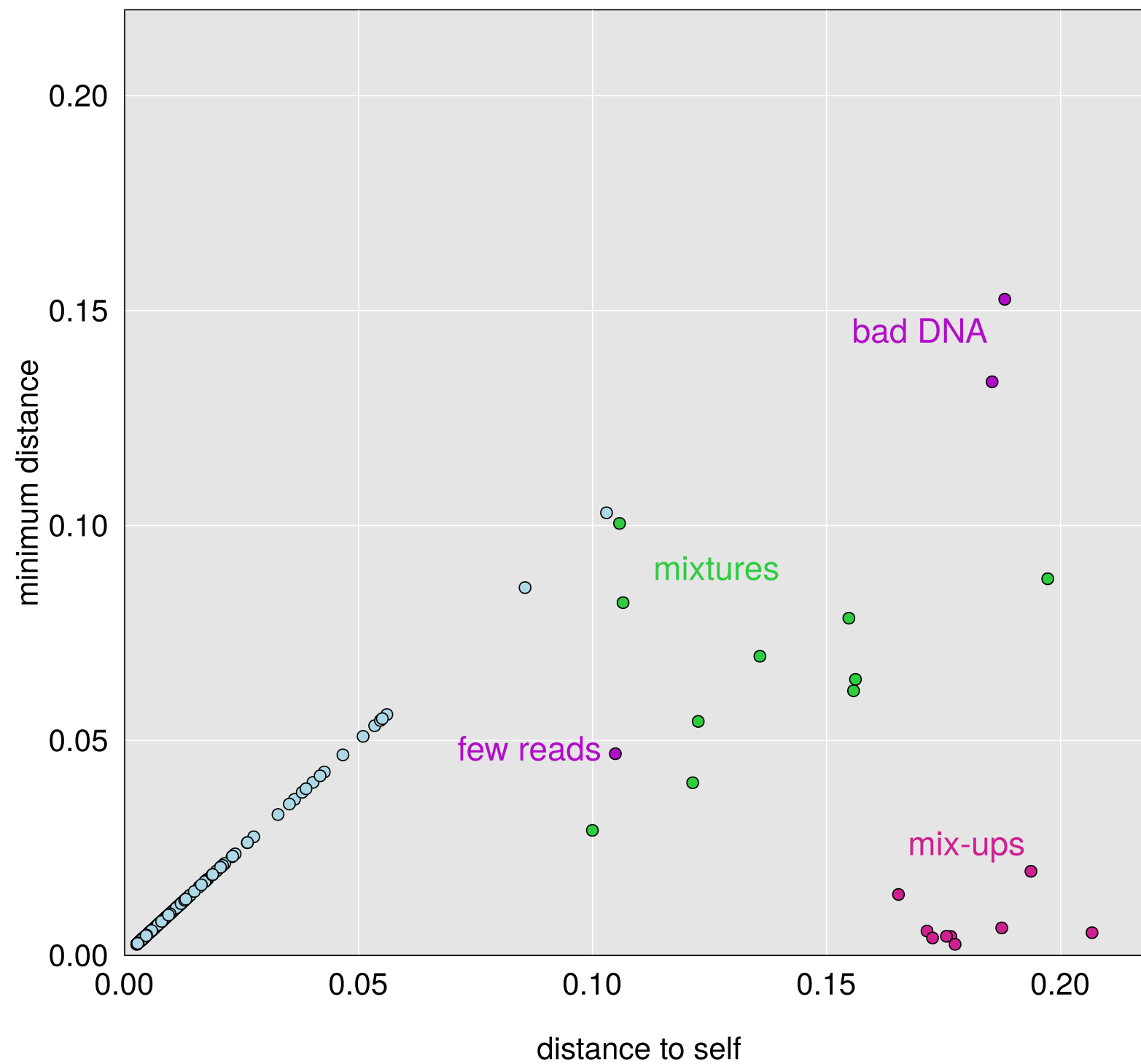
Genomic DO370 vs Microbiome DO360

	A	B
AA	3,137,751	7,461
BB	1,475	310,369

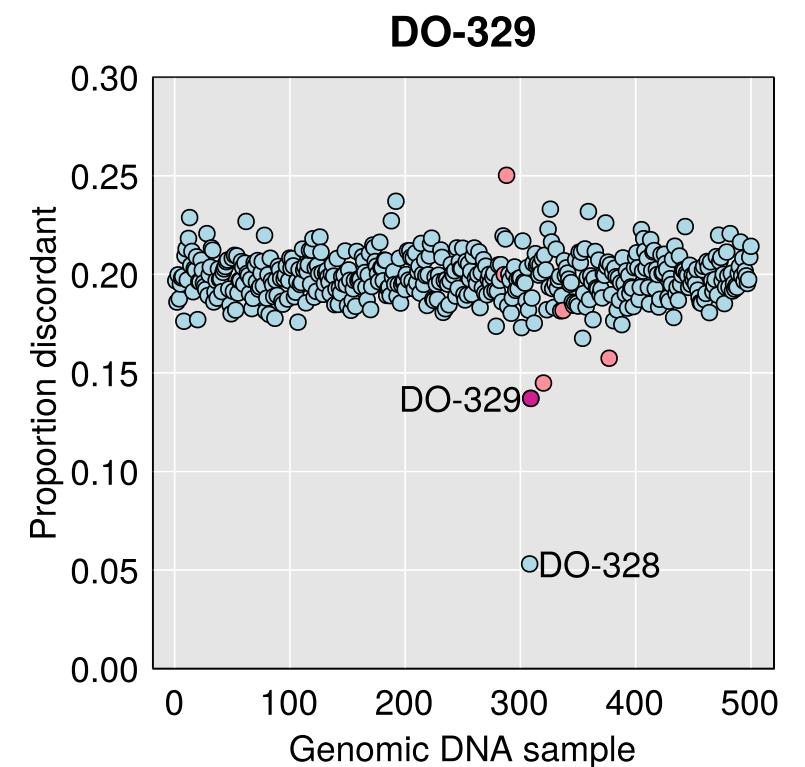
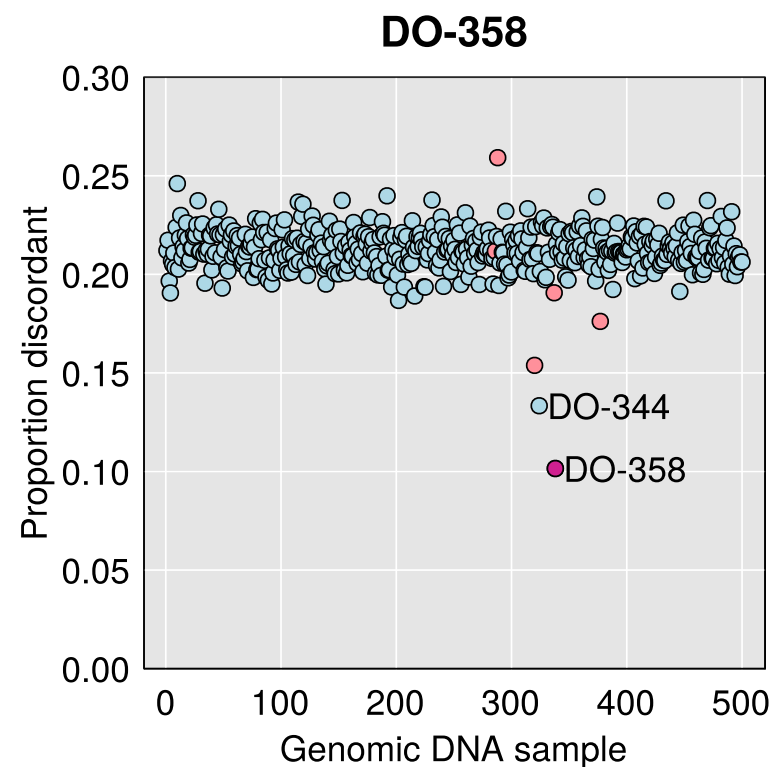
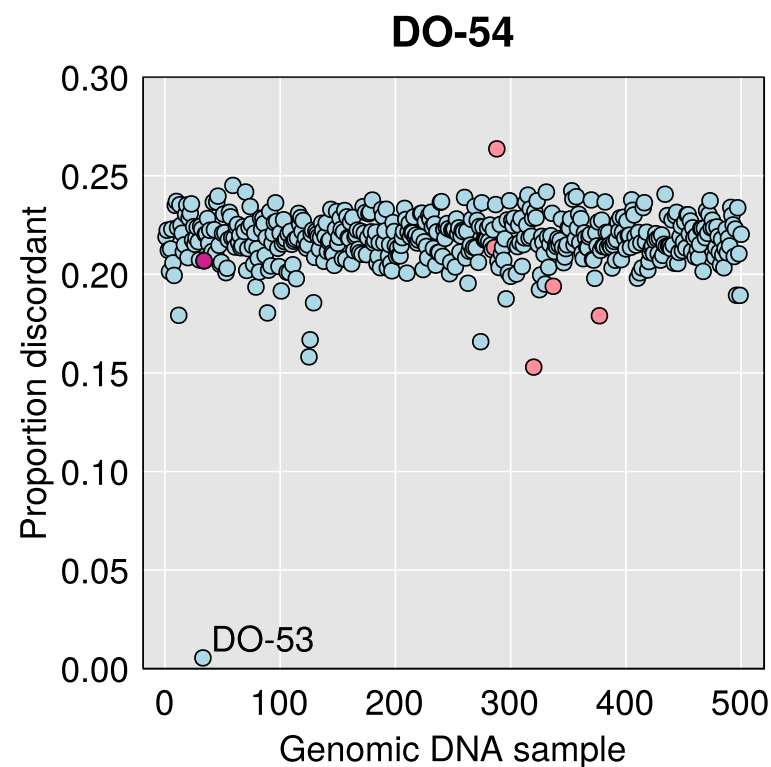
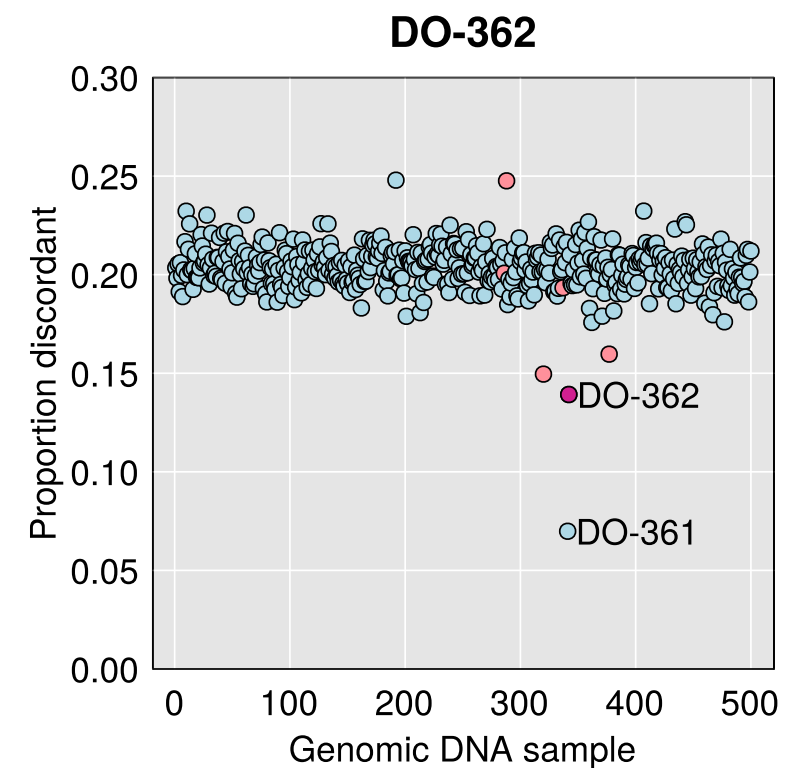
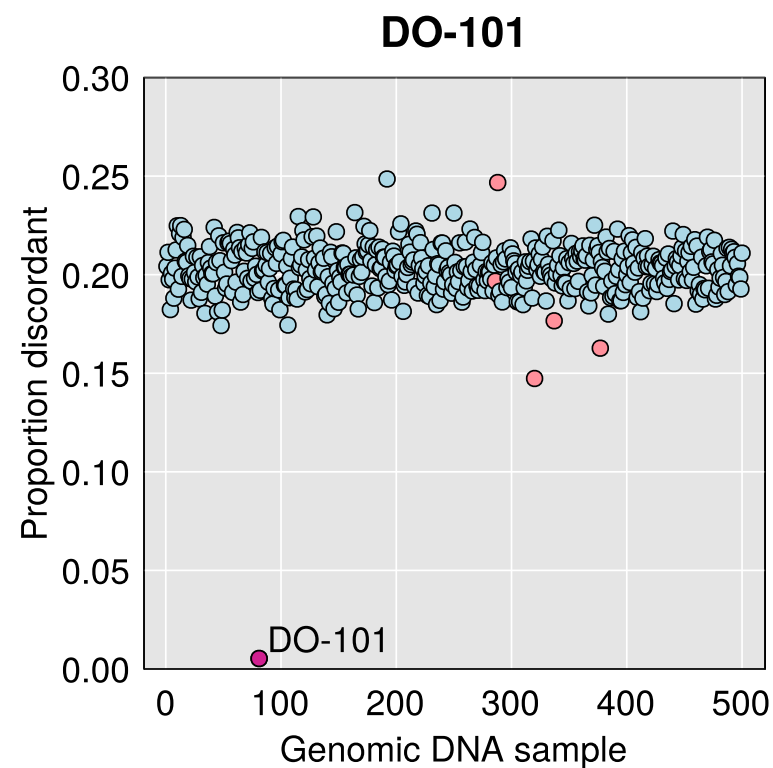
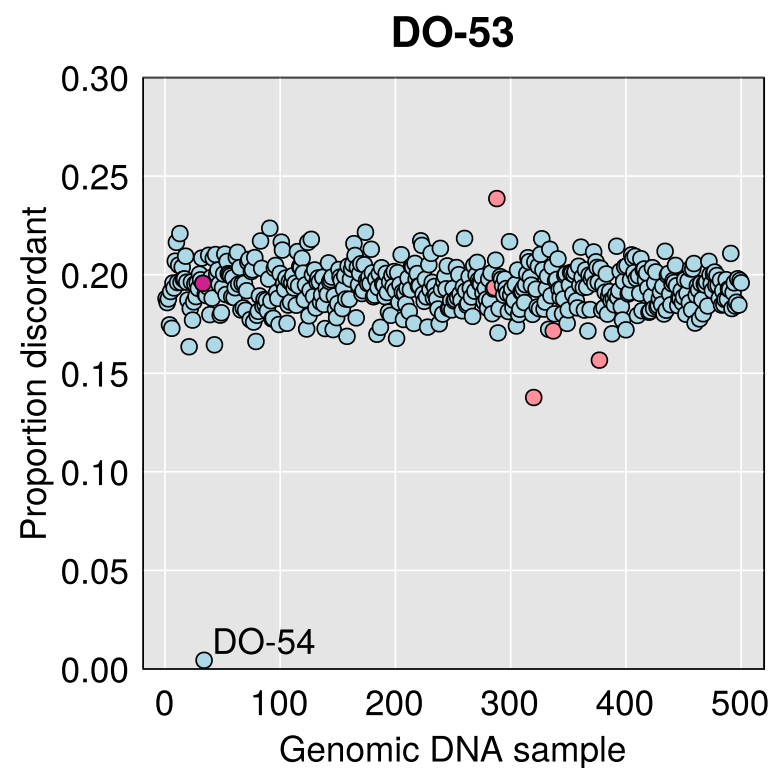
Microbiome mix-ups: min vs self distance



Microbiome mix-ups: min vs self distance



Microbiome mix-ups and mixtures: detail



Gen DO101 & DO102 vs Mic DO101

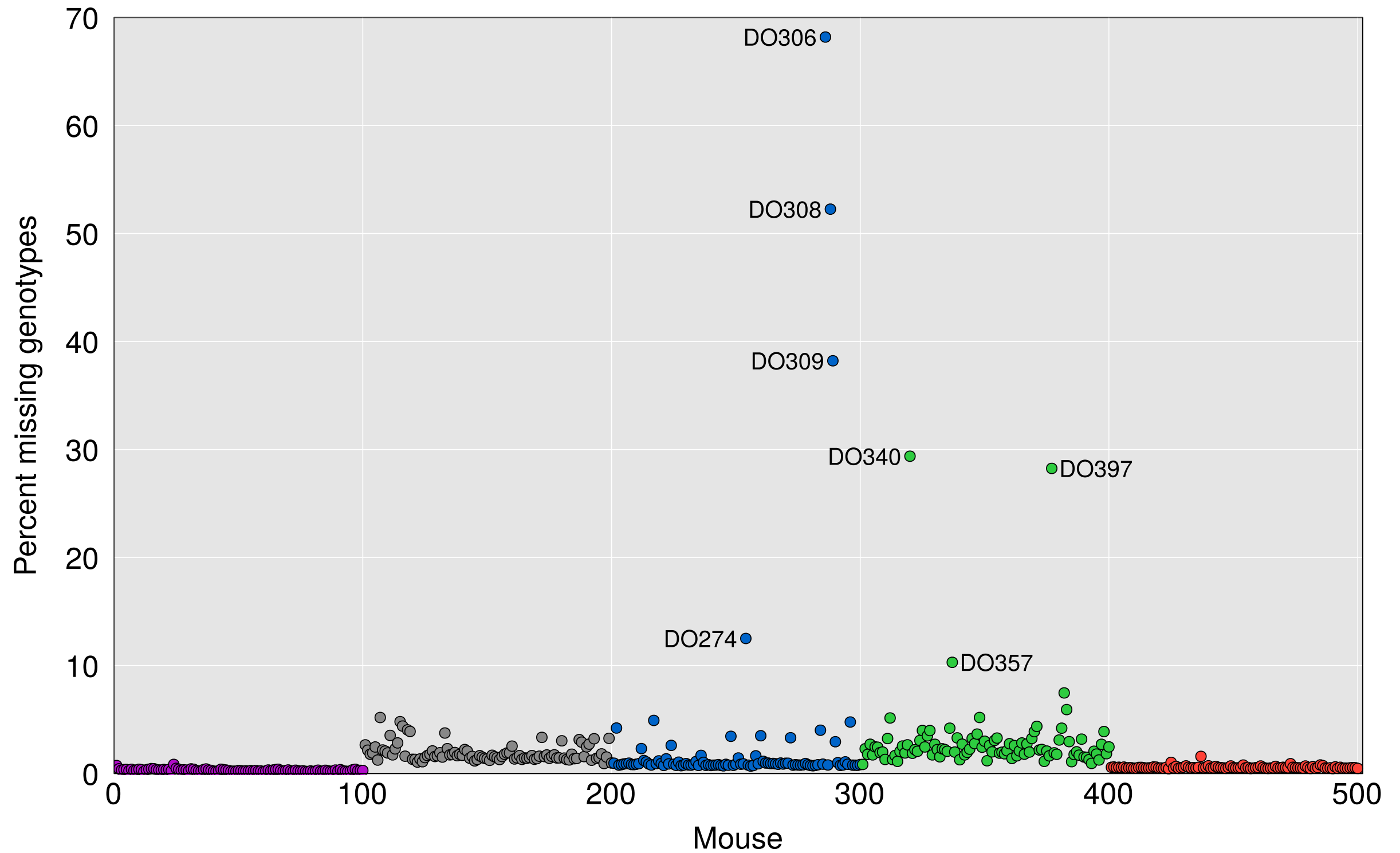
DO-101 genotype	DO-102 genotype	allele in DO-102 microbiome			
		A (%)		B (%)	
AA	AA	3,664,076	(99.6%)	14,305	(0.4%)
AA	AB	1,161,383	(99.5%)	6,187	(0.5%)
AA	BB	153,501	(99.3%)	1,067	(0.7%)
AB	AA	651,287	(52.0%)	600,434	(48.0%)
AB	AB	378,800	(51.8%)	352,828	(48.2%)
AB	BB	155,967	(51.2%)	148,703	(48.8%)
BB	AA	3,088	(1.6%)	185,825	(98.4%)
BB	AB	3,210	(1.3%)	240,712	(98.7%)
BB	BB	2,162	(1.0%)	217,882	(99.0%)

Gen DO358 & DO344 vs Mic DO358

DO-358 genotype	DO-344 genotype	allele in DO-358 microbiome			
		A (%)		B (%)	
AA	AA	2,394,215	(99.7%)	6,050	(0.3%)
AA	AB	869,613	(79.5%)	224,483	(20.5%)
AA	BB	103,036	(59.1%)	71,332	(40.9%)
AB	AA	686,970	(71.8%)	269,447	(28.2%)
AB	AB	297,500	(51.4%)	280,958	(48.6%)
AB	BB	55,982	(29.9%)	131,111	(70.1%)
BB	AA	73,727	(42.9%)	98,257	(57.1%)
BB	AB	47,000	(21.9%)	167,513	(78.1%)
BB	BB	542	(0.5%)	117,802	(99.5%)

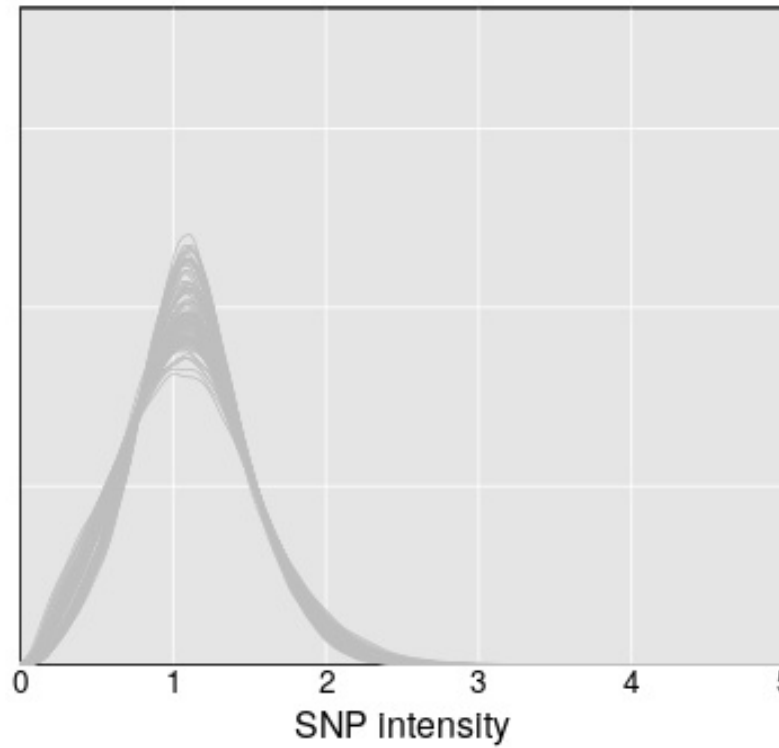
Sample quality

Missing data per sample

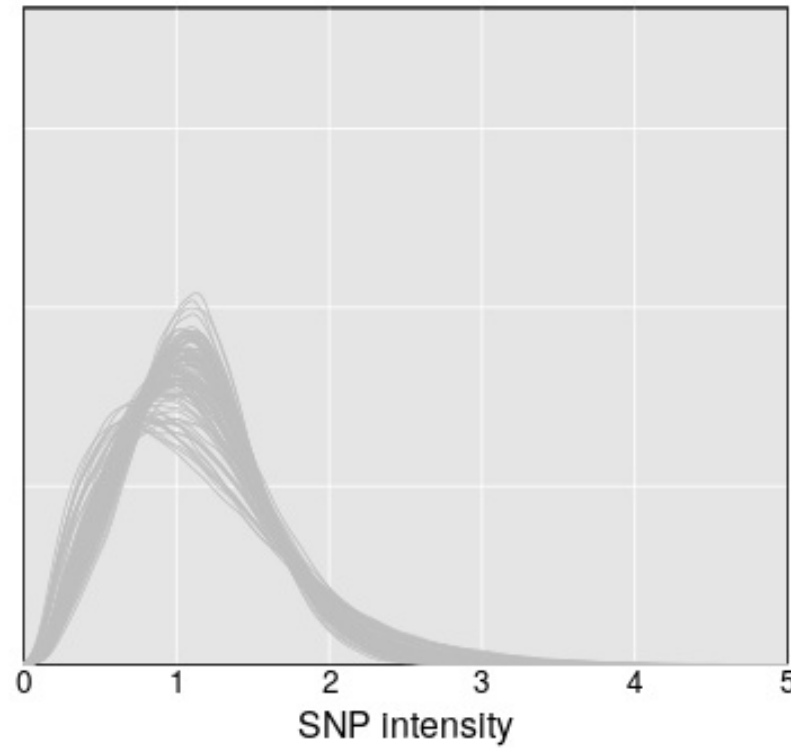


Array intensities

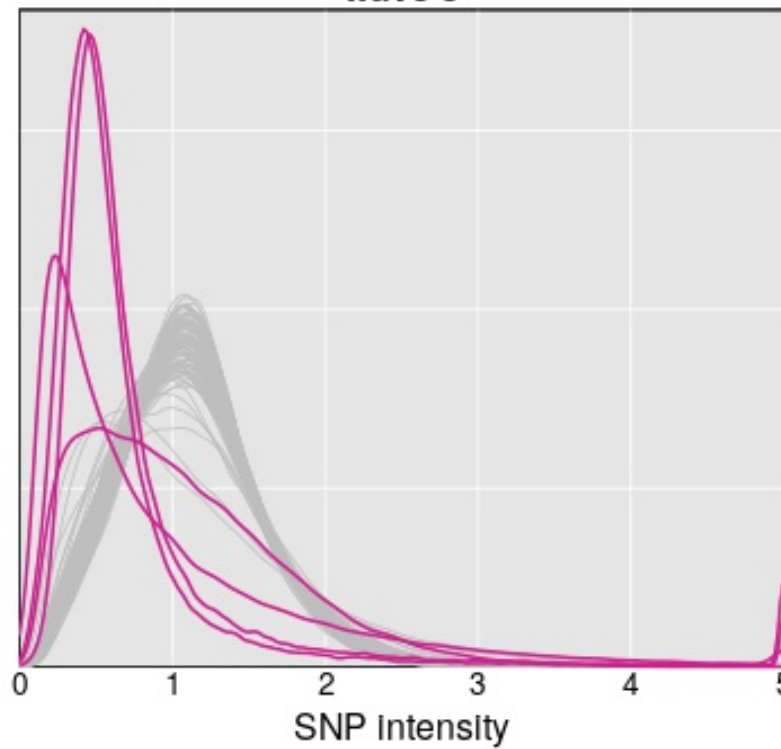
wave 1



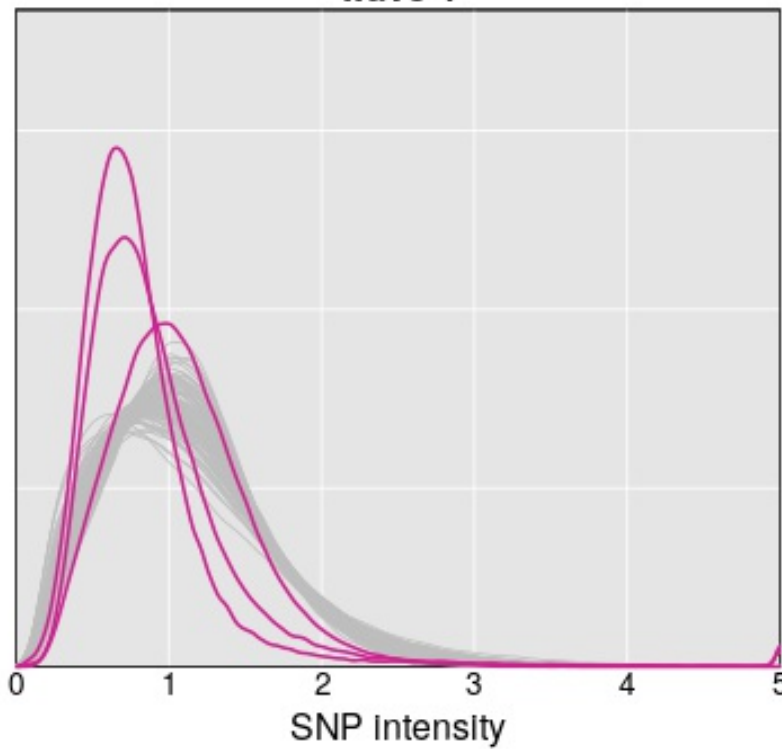
wave 2



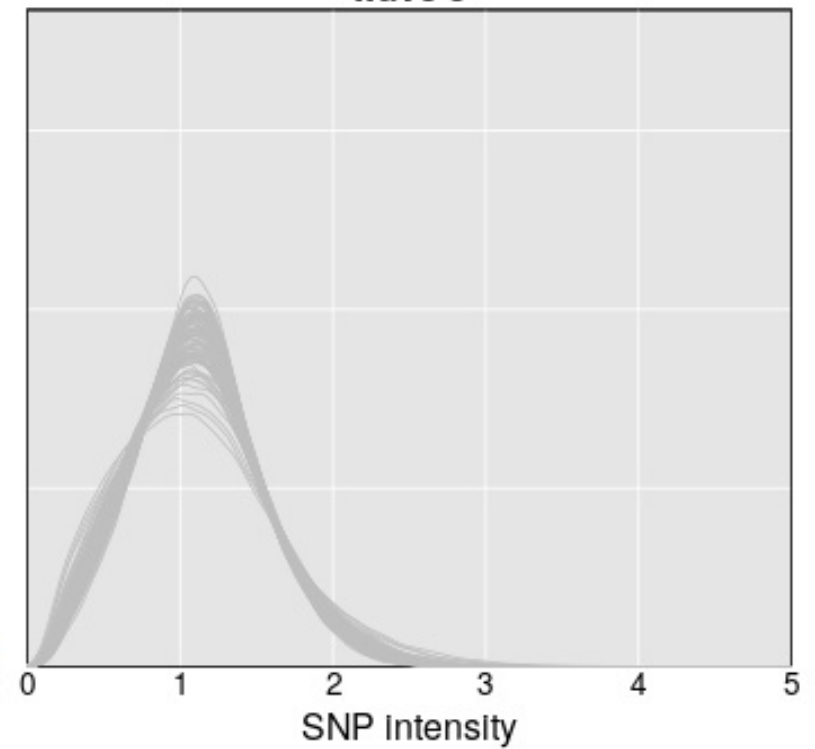
wave 3



wave 4

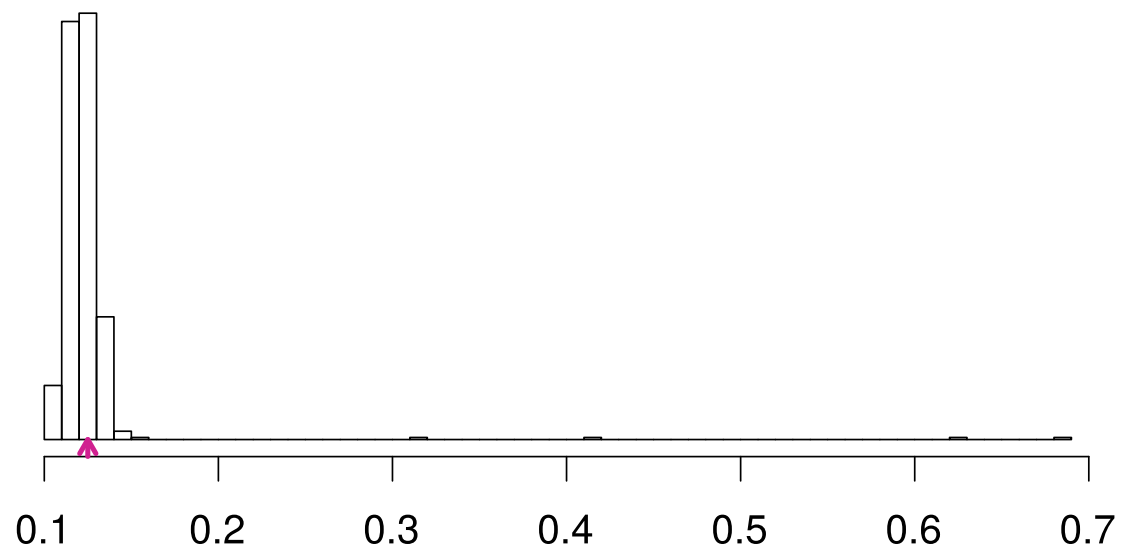


wave 5



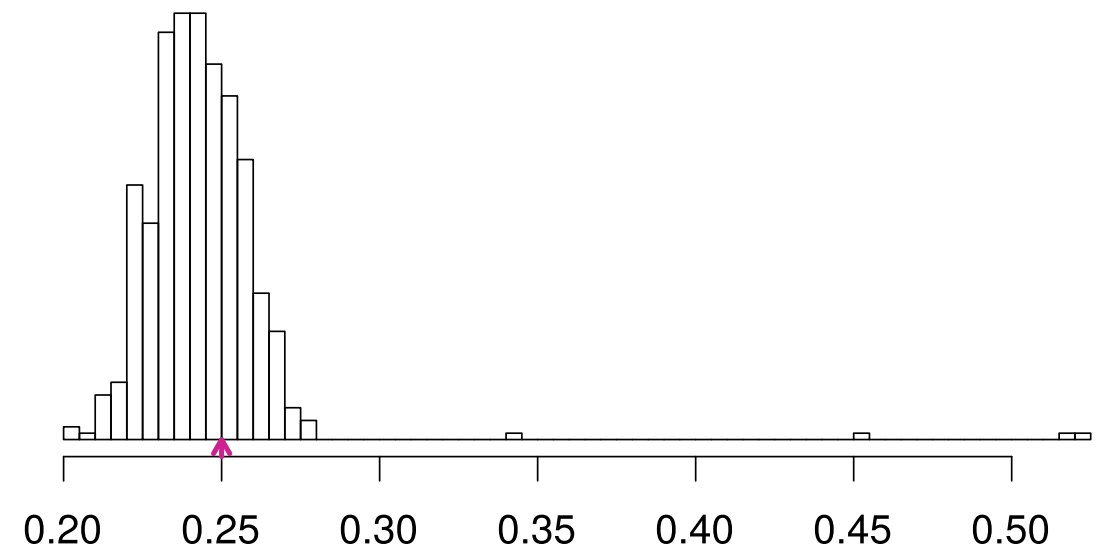
Allele frequencies, by individual

founder MAF = 1/8



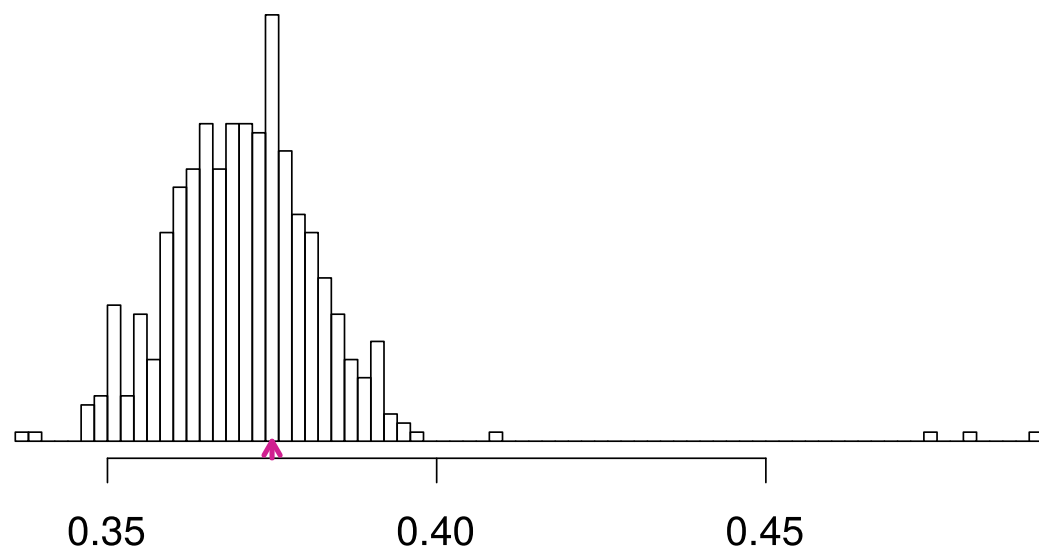
Frequency of minor allele

founder MAF = 2/8



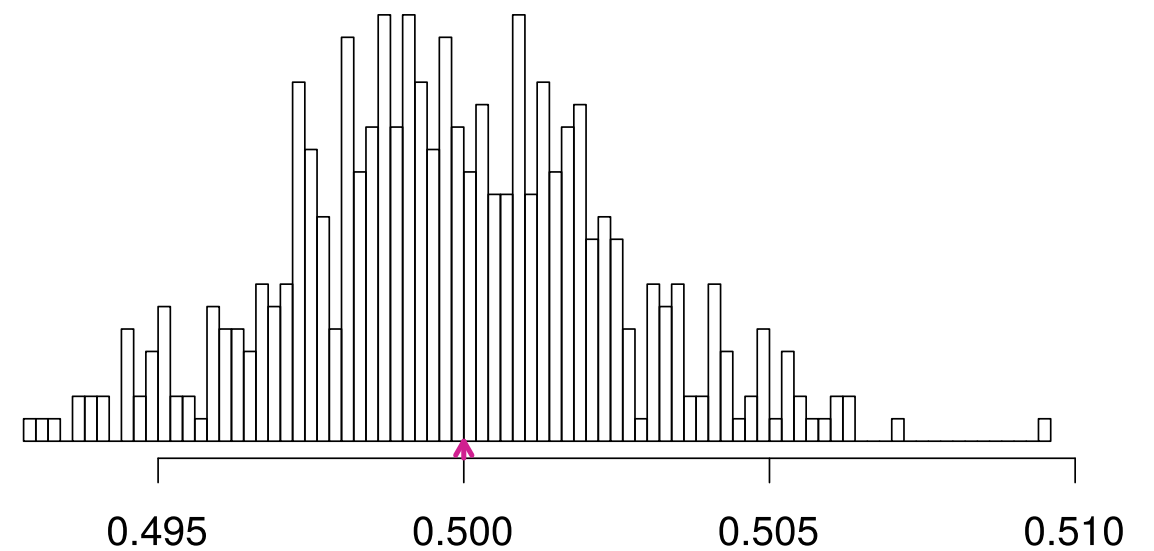
Frequency of minor allele

founder MAF = 3/8



Frequency of minor allele

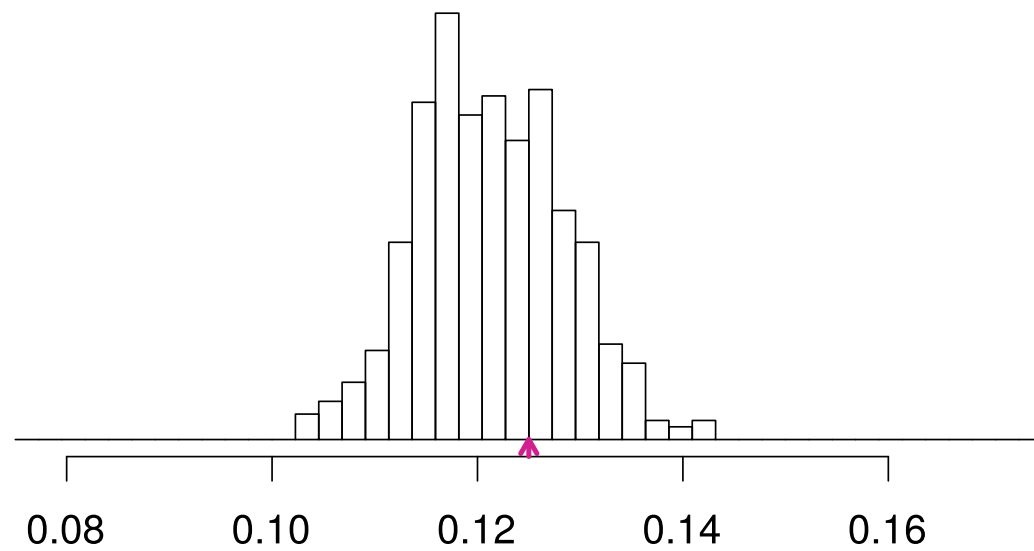
founder MAF = 4/8



Frequency of minor allele

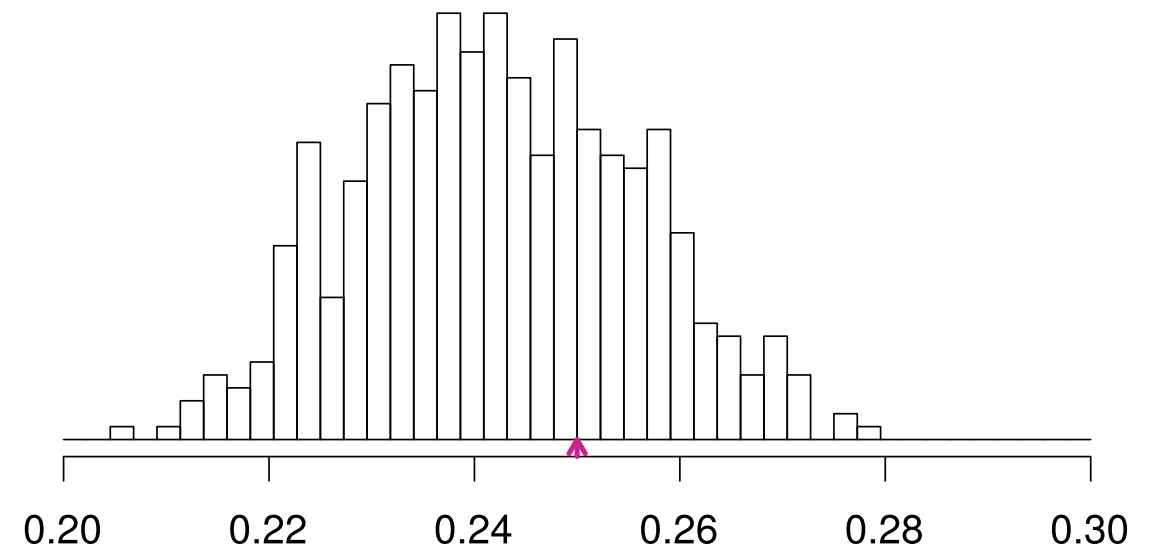
Allele frequencies, by individual

founder MAF = 1/8



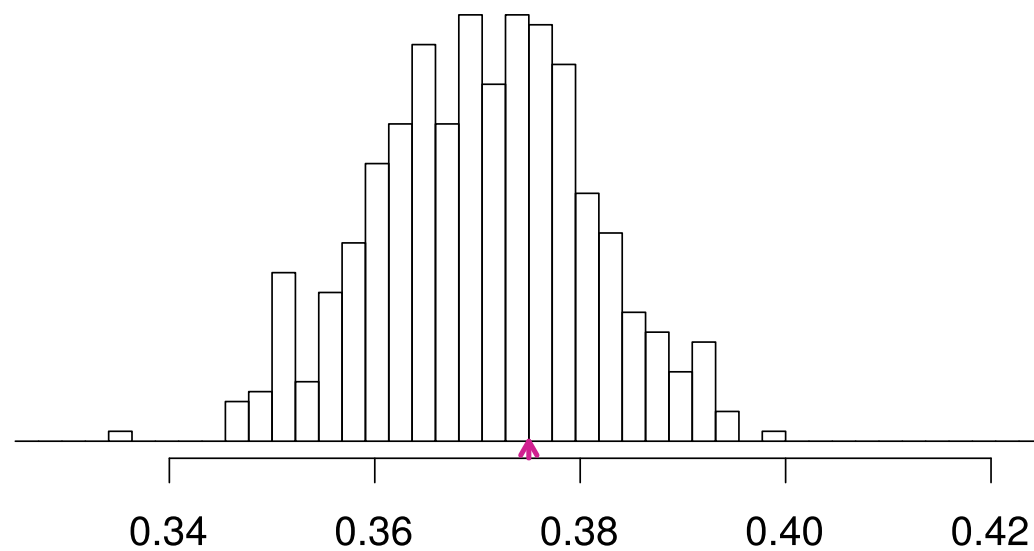
Frequency of minor allele

founder MAF = 2/8



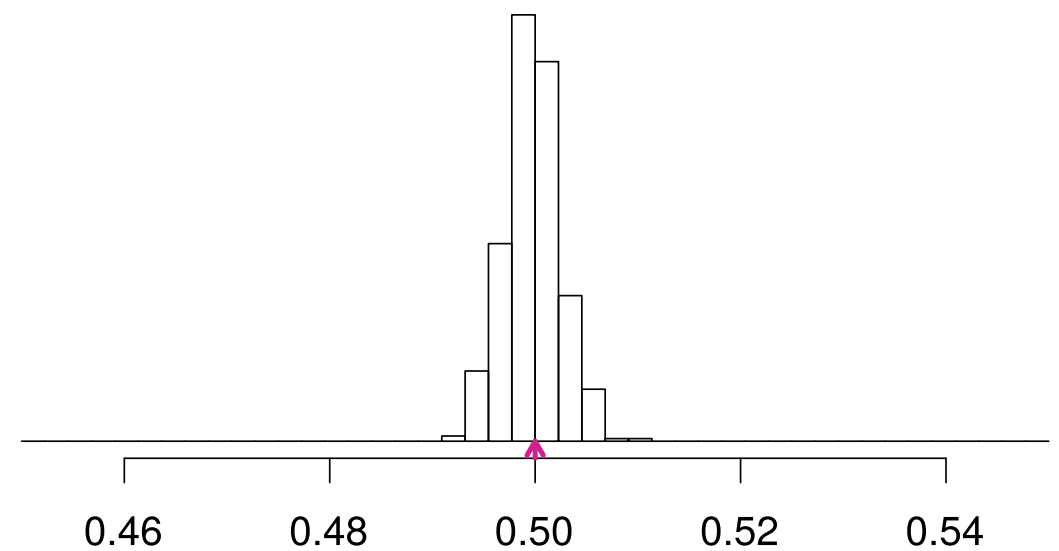
Frequency of minor allele

founder MAF = 3/8



Frequency of minor allele

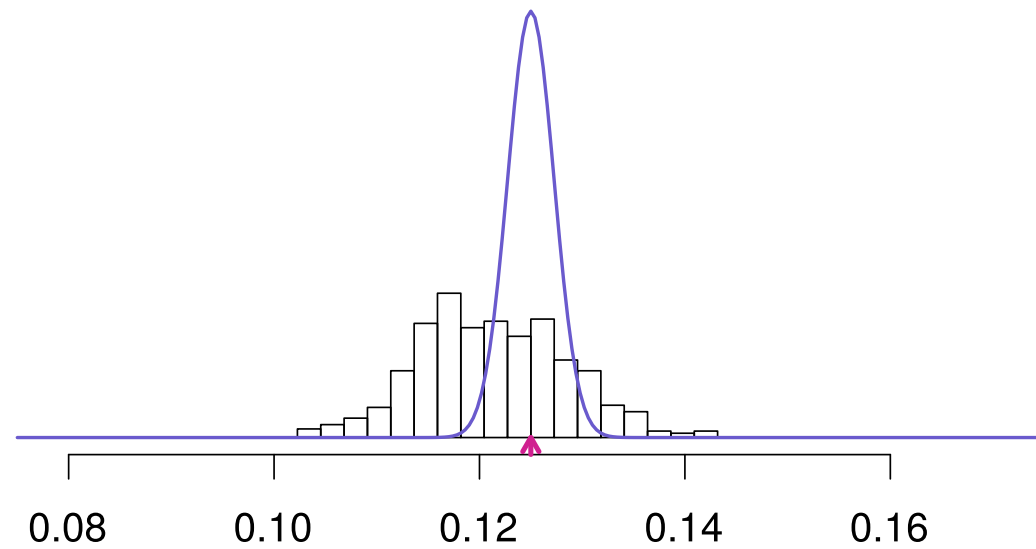
founder MAF = 4/8



Frequency of minor allele

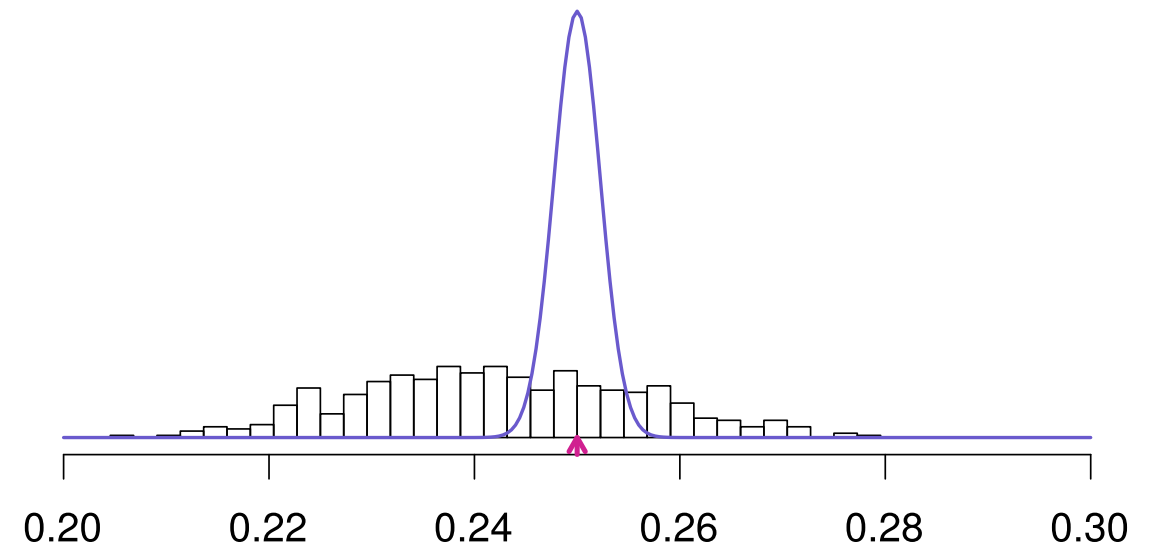
Allele frequencies, by individual

founder MAF = 1/8



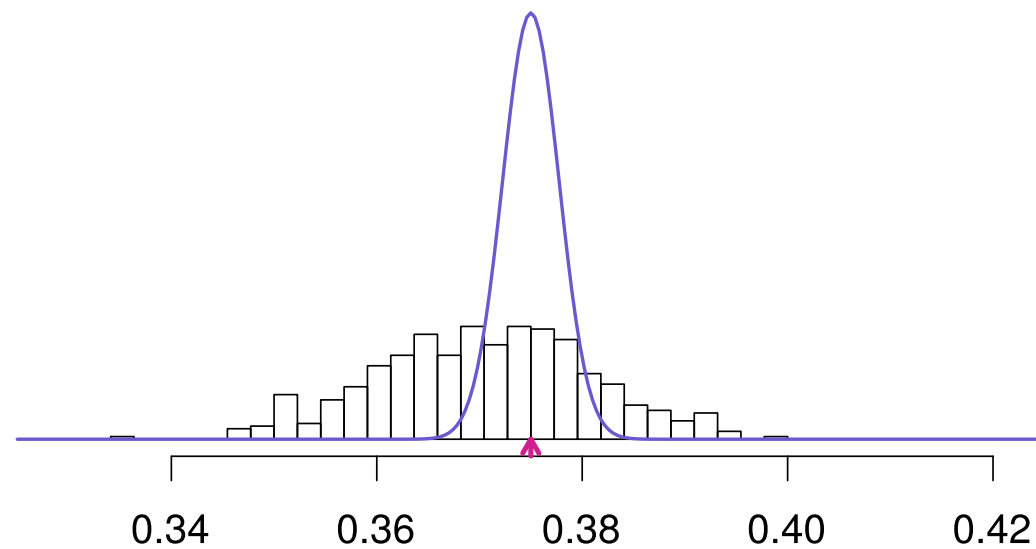
Frequency of minor allele

founder MAF = 2/8



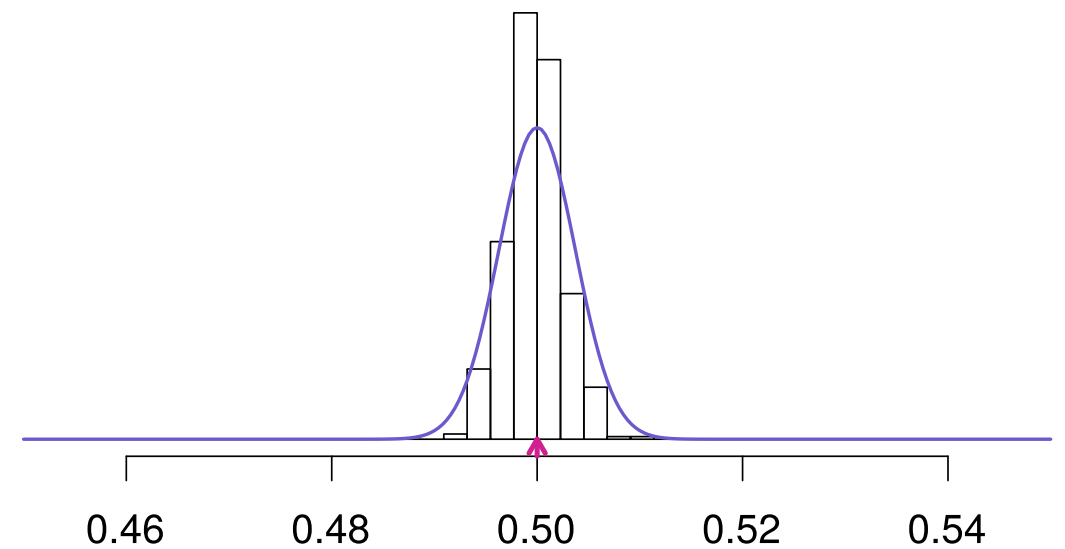
Frequency of minor allele

founder MAF = 3/8



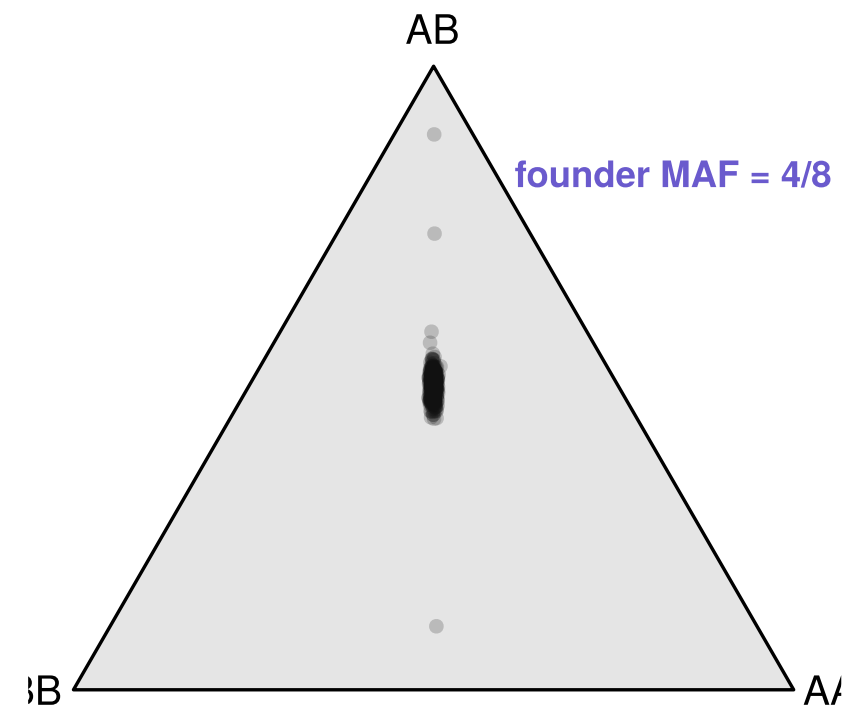
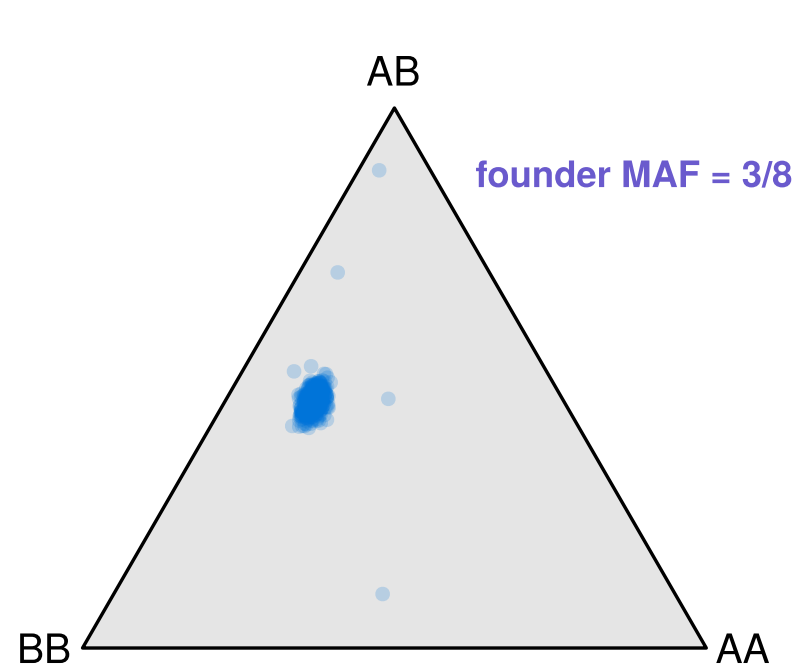
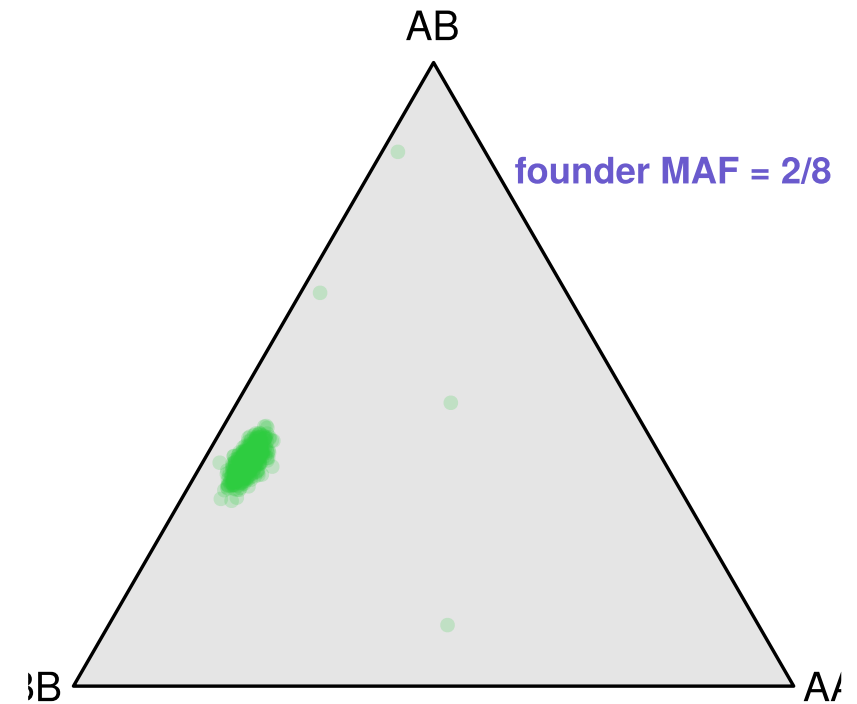
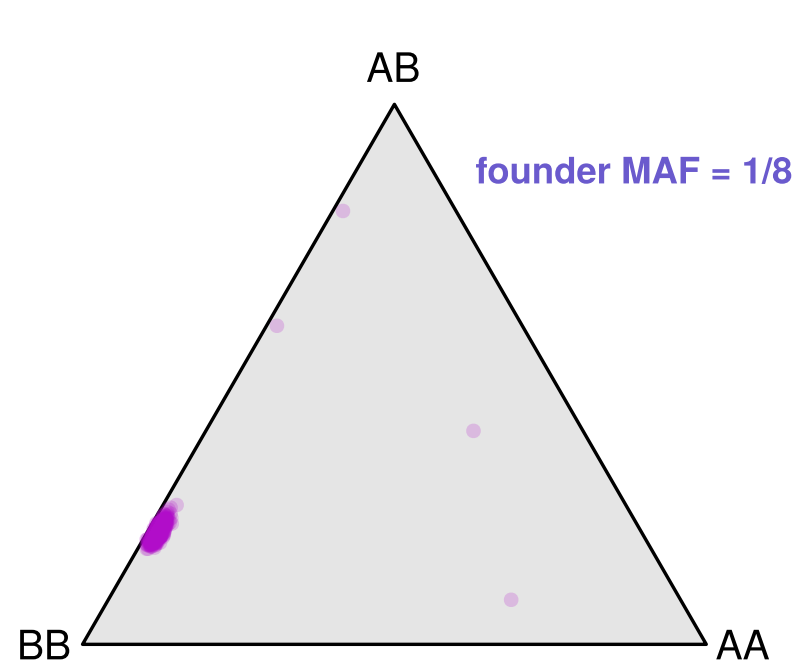
Frequency of minor allele

founder MAF = 4/8

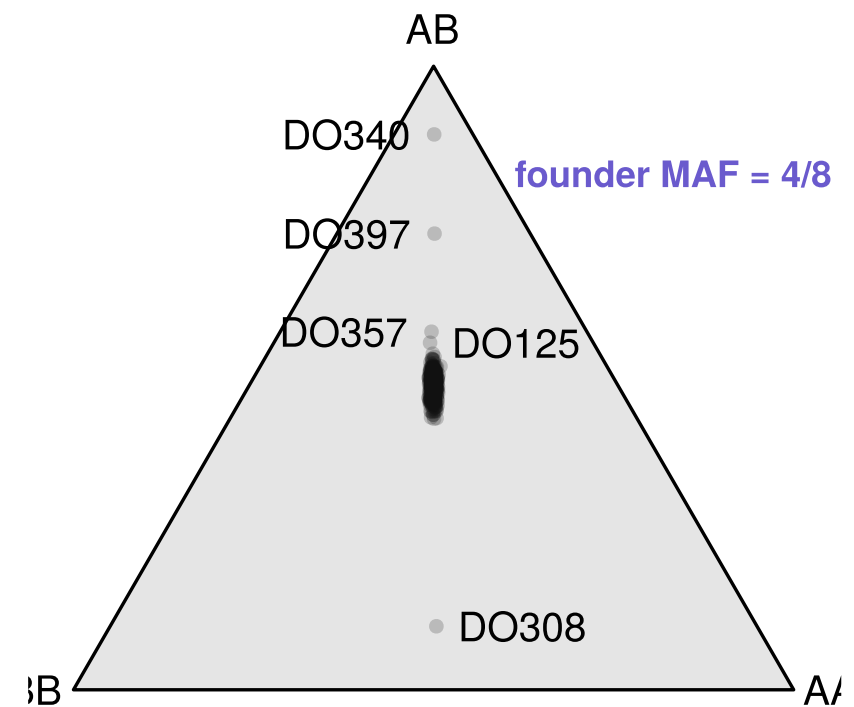
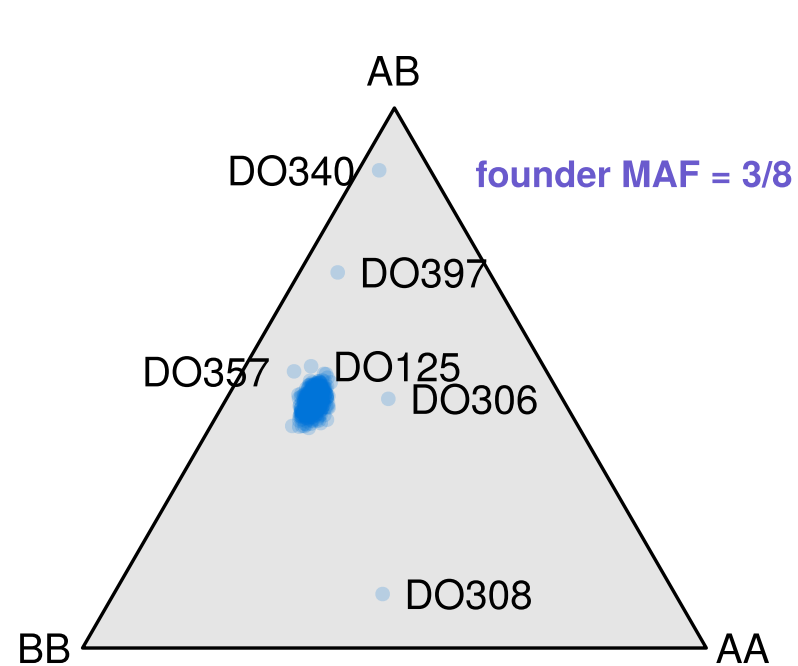
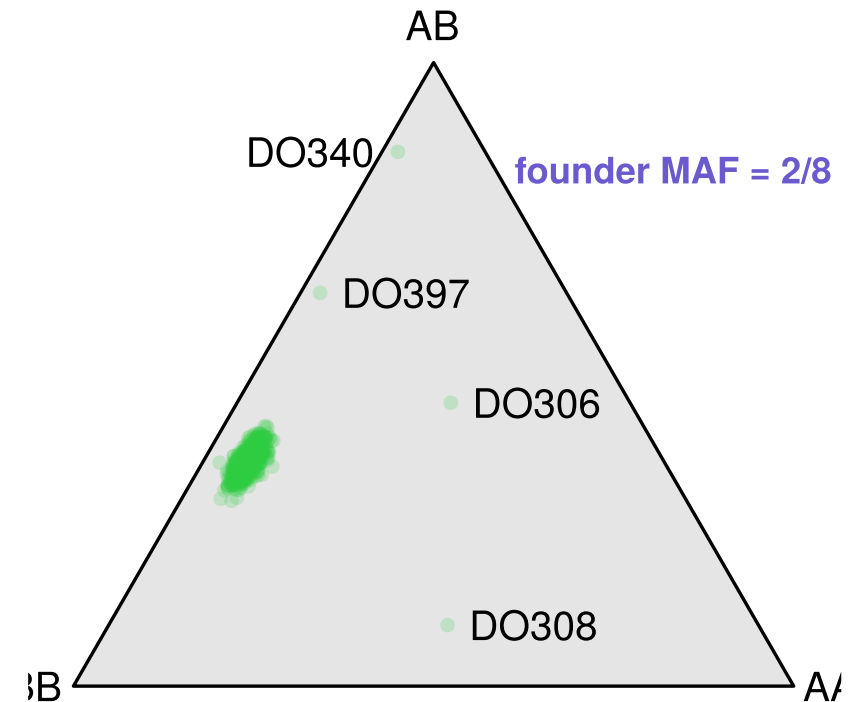
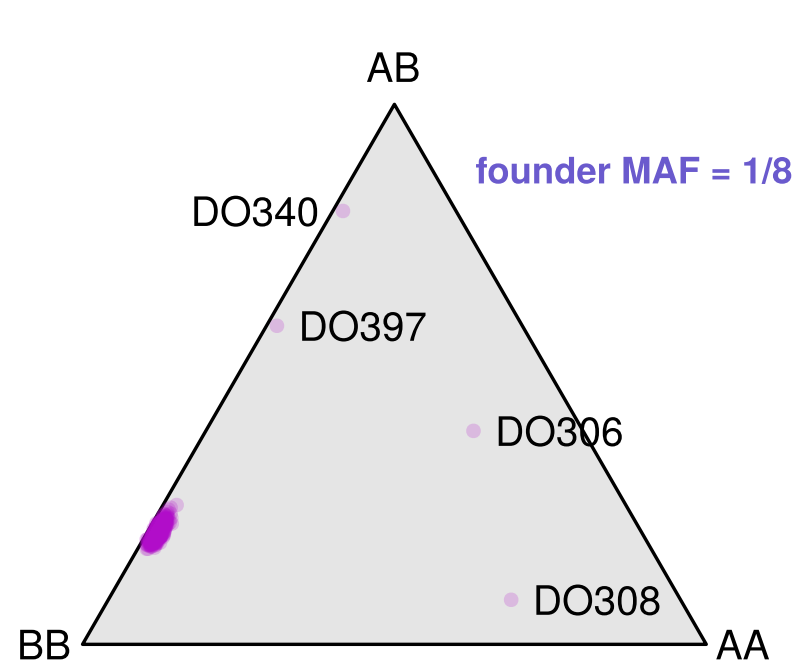


Frequency of minor allele

Genotype frequencies, by individual

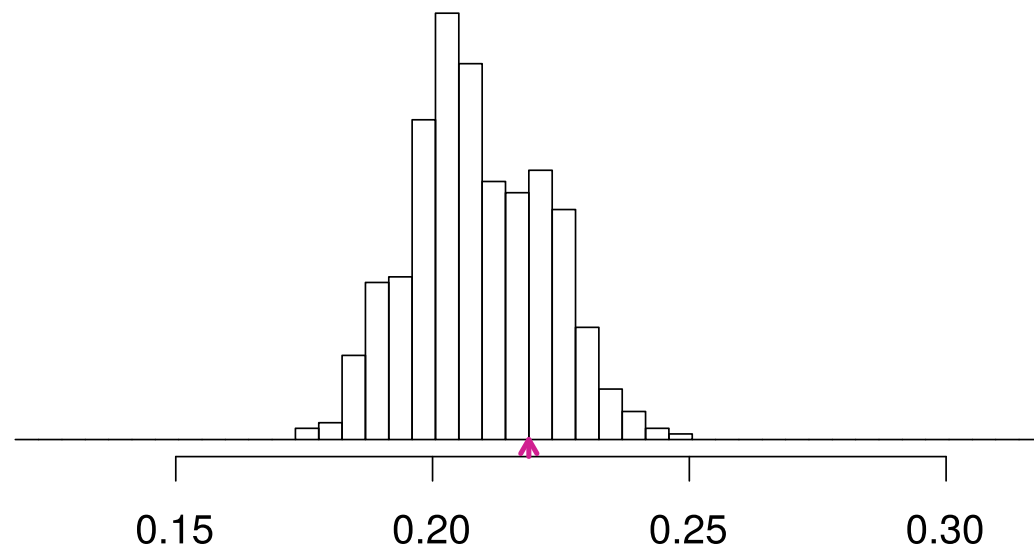


Genotype frequencies, by individual



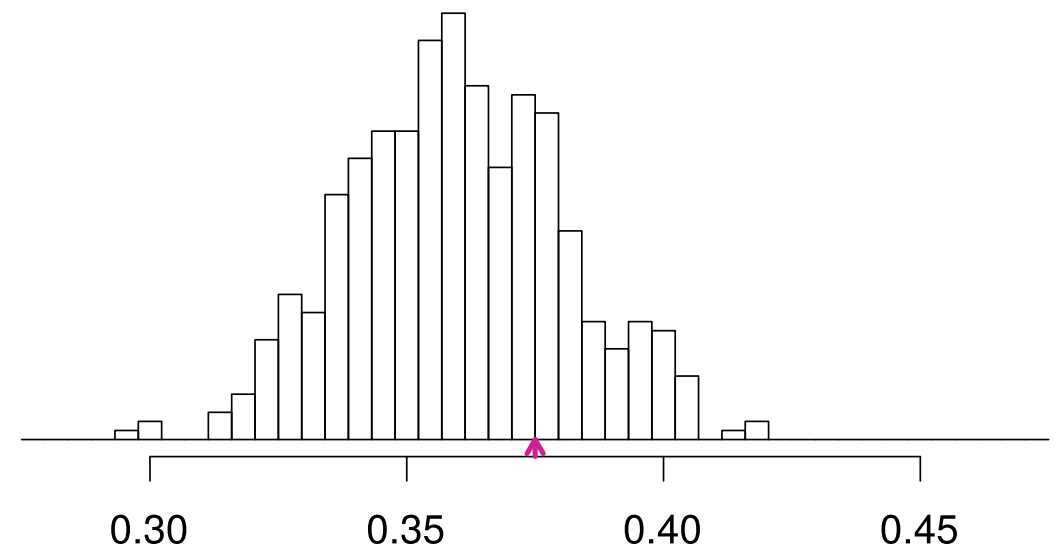
Heterozygosities, by individual

founder MAF = 1/8



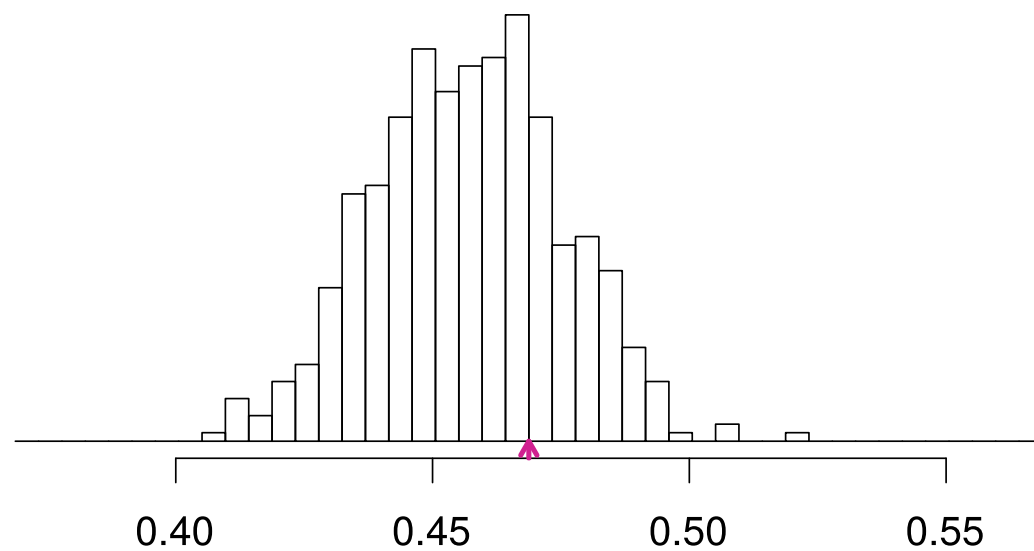
Frequency of minor allele

founder MAF = 2/8



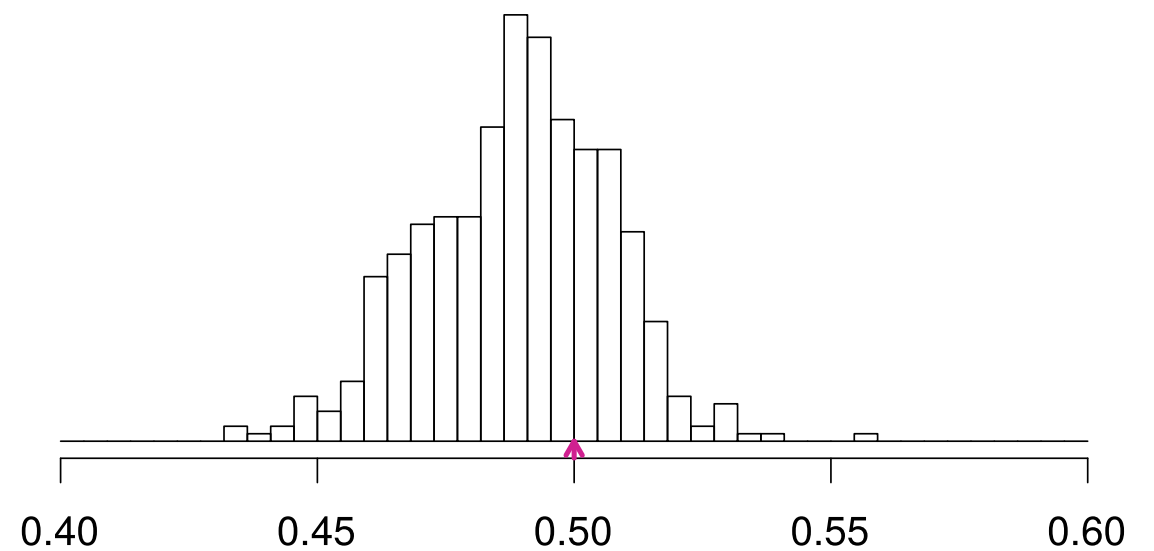
Frequency of minor allele

founder MAF = 3/8



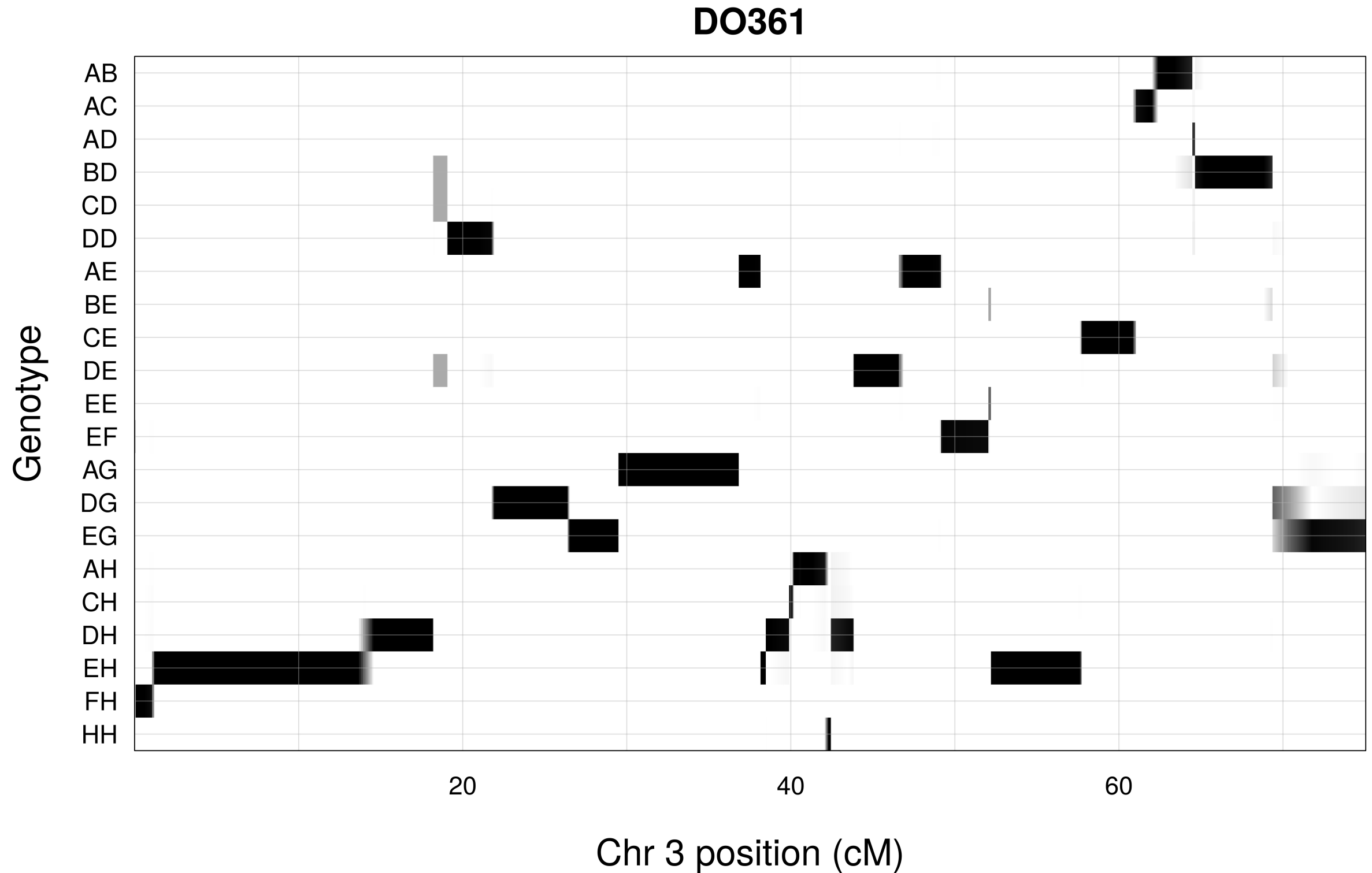
Frequency of minor allele

founder MAF = 4/8



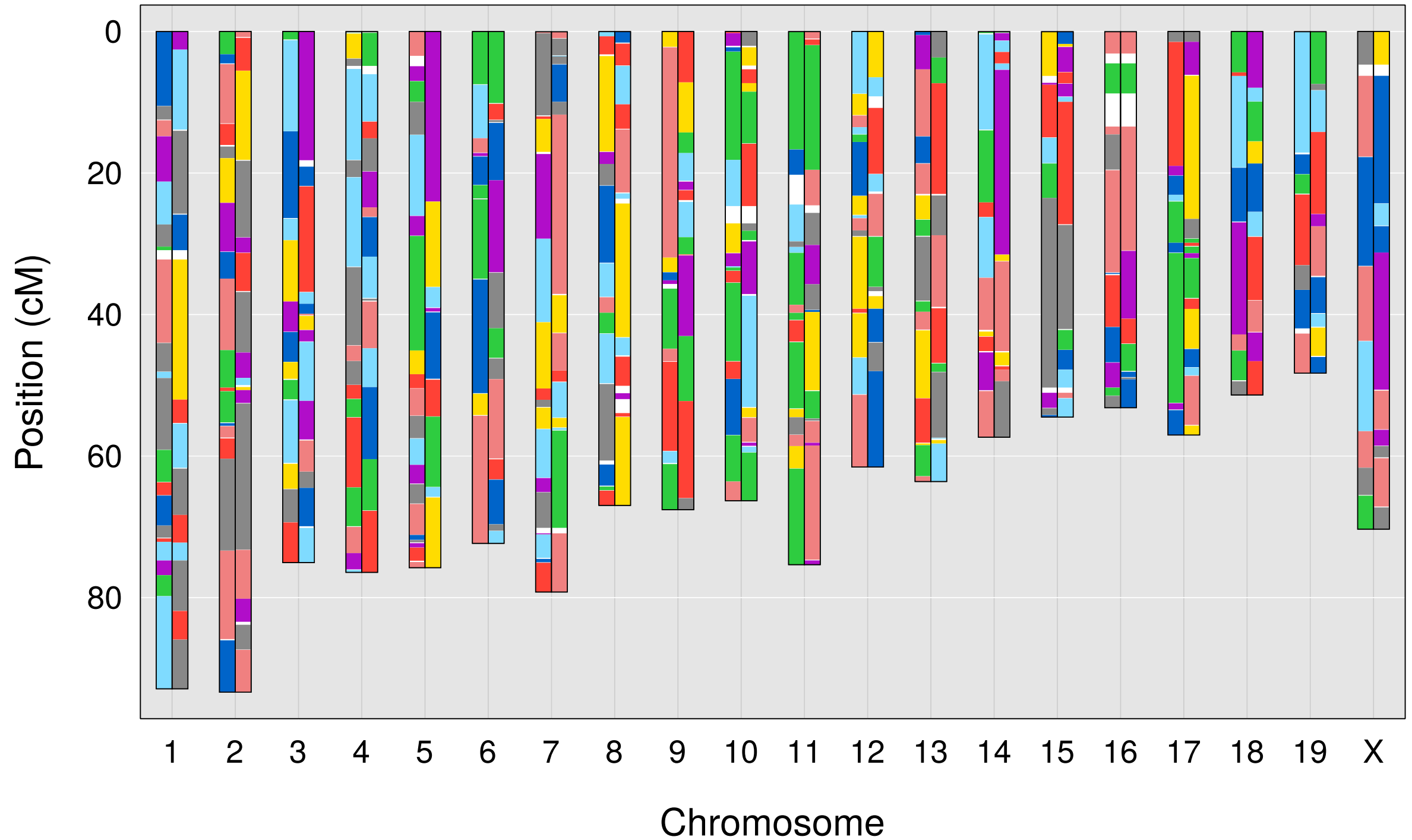
Frequency of minor allele

Genotype probabilities (one mouse on one chr)

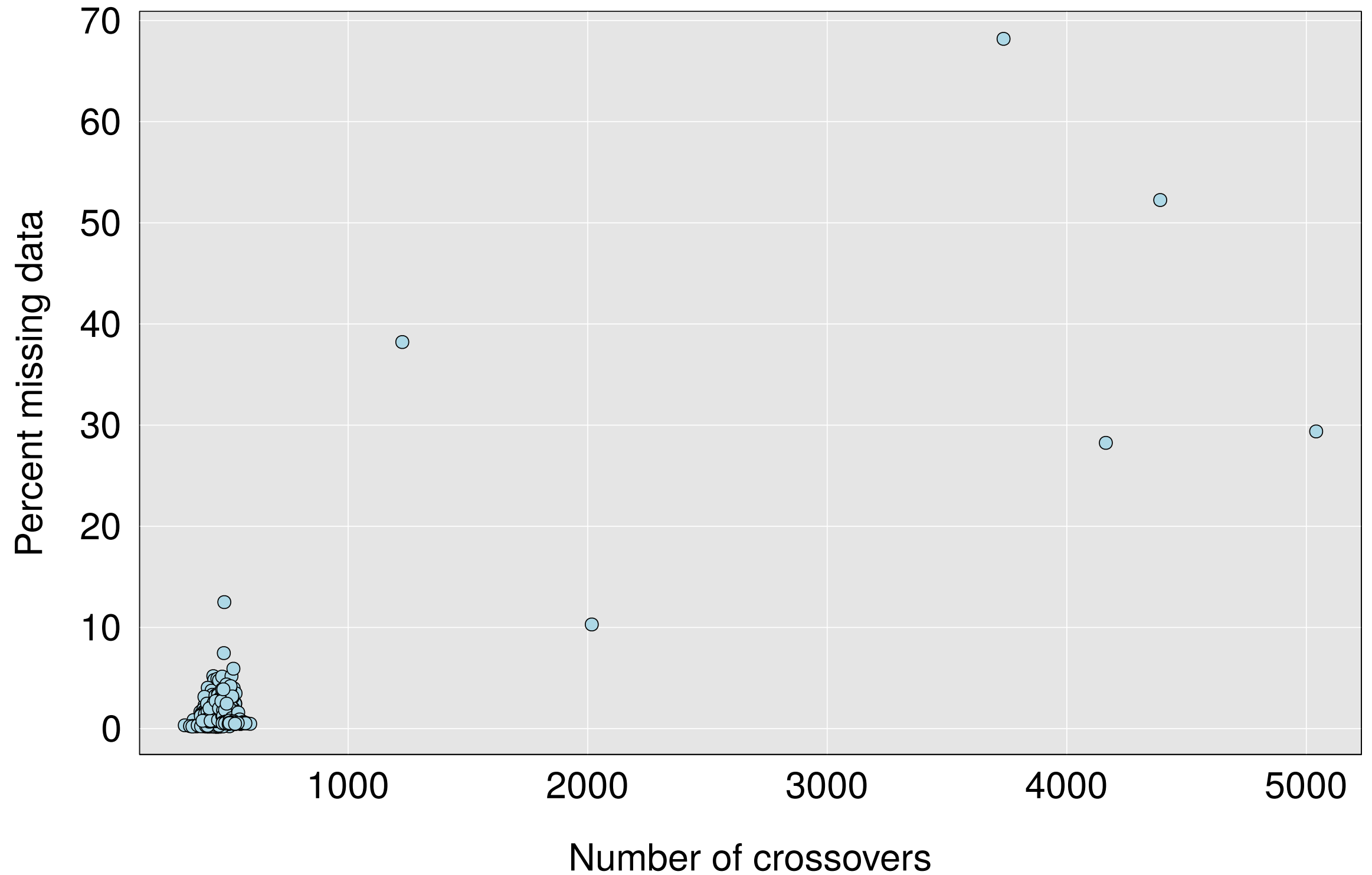


Genome reconstruction (one mouse)

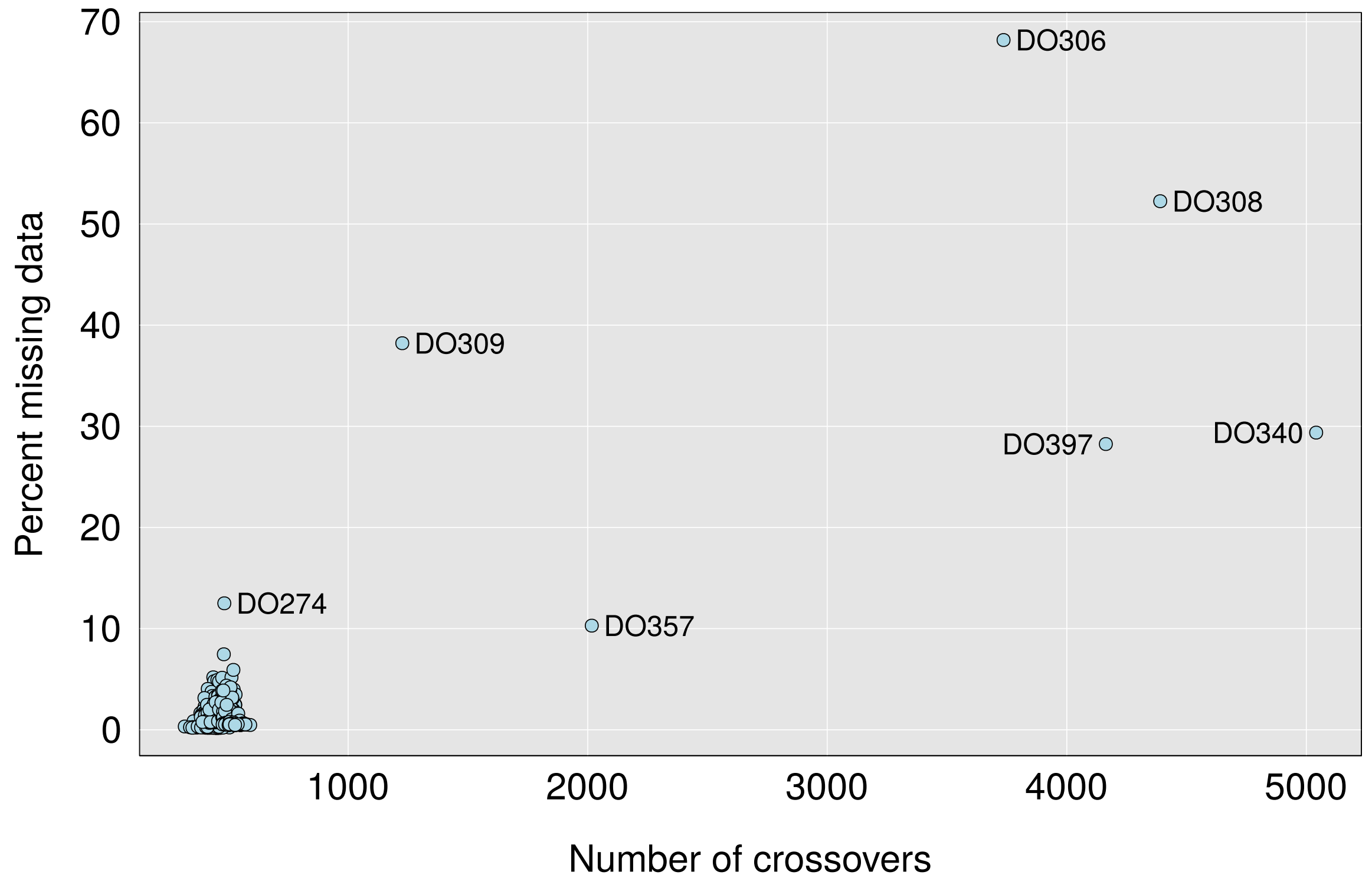
DO361



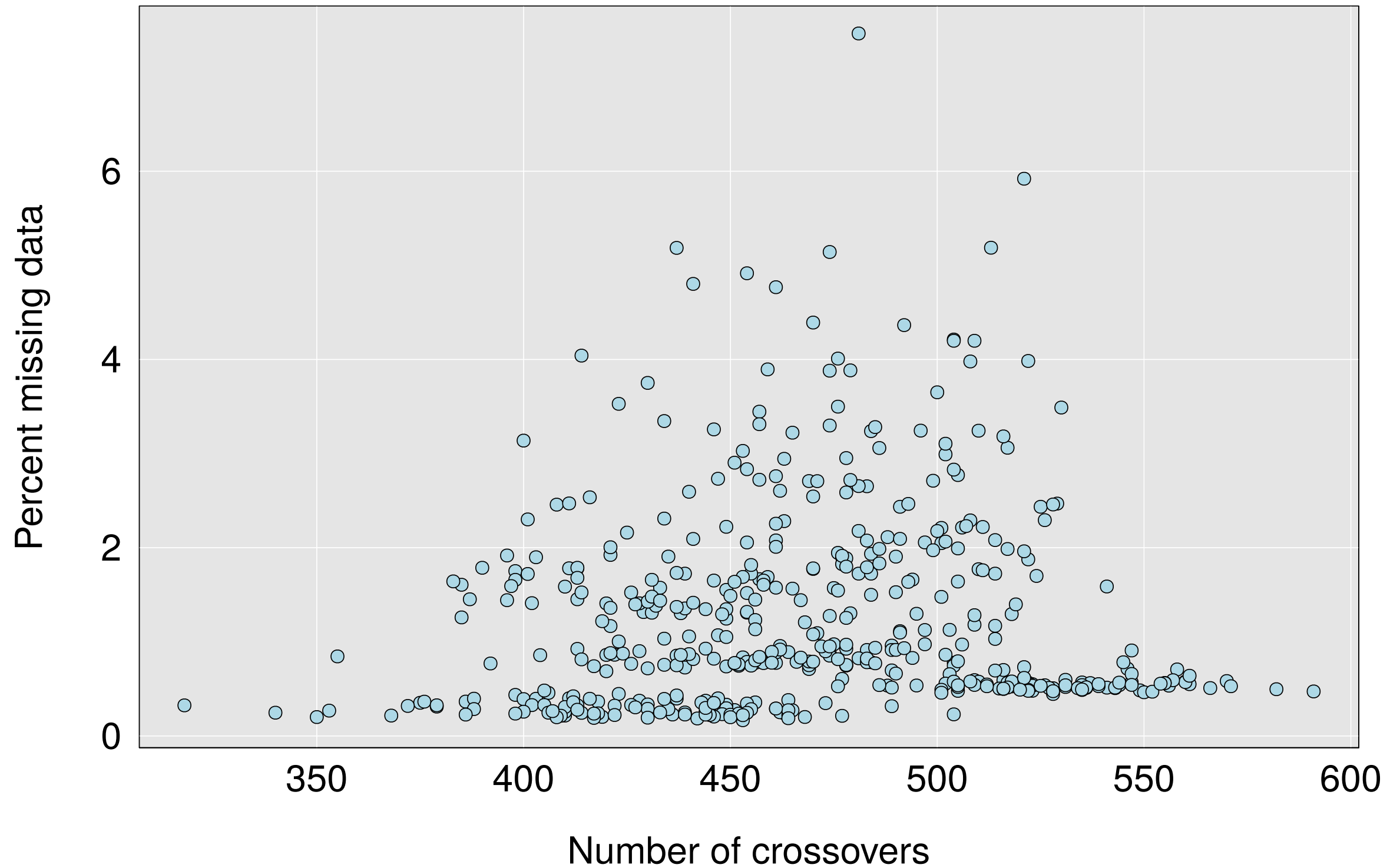
Percent missing vs number of crossovers



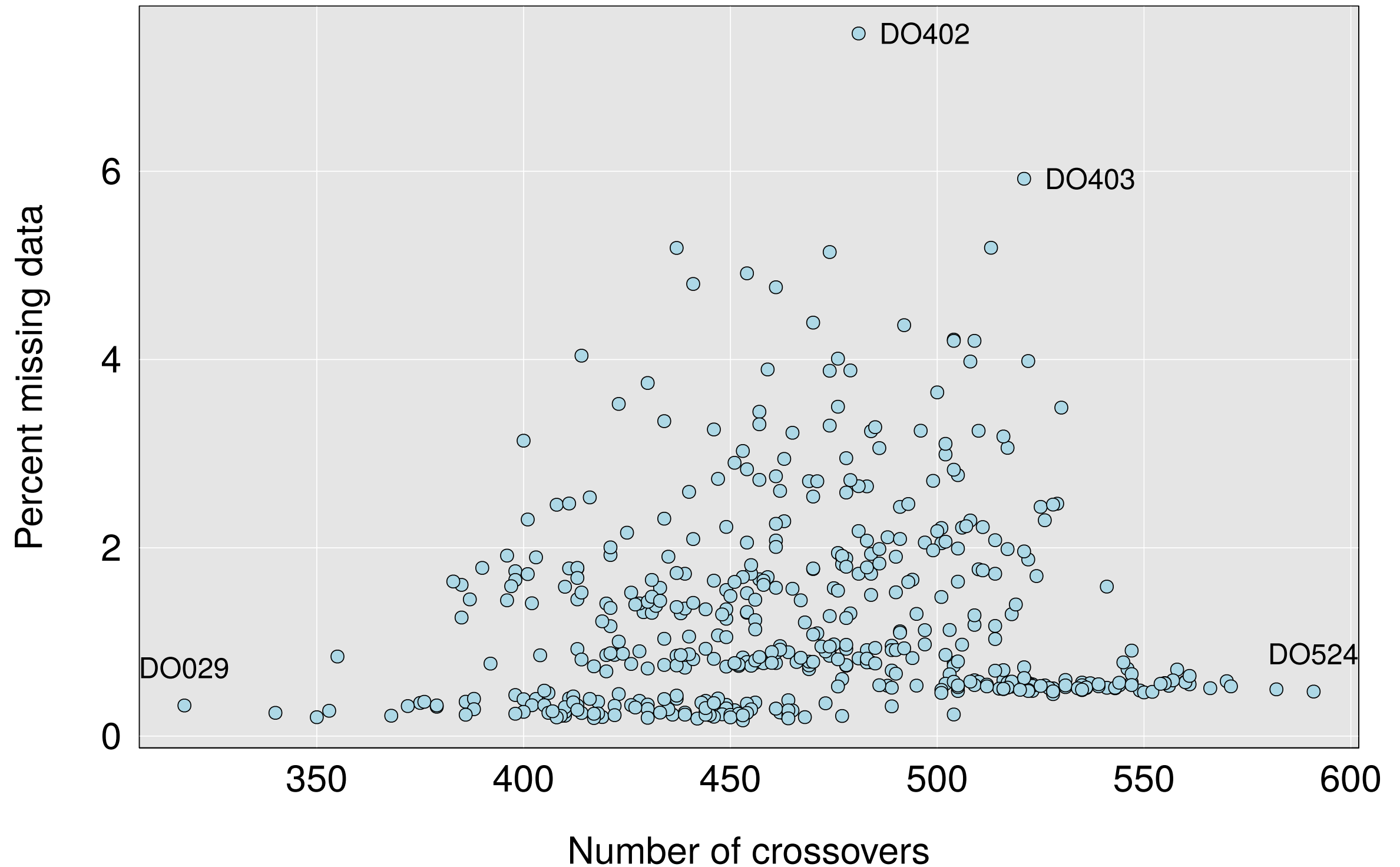
Percent missing vs number of crossovers



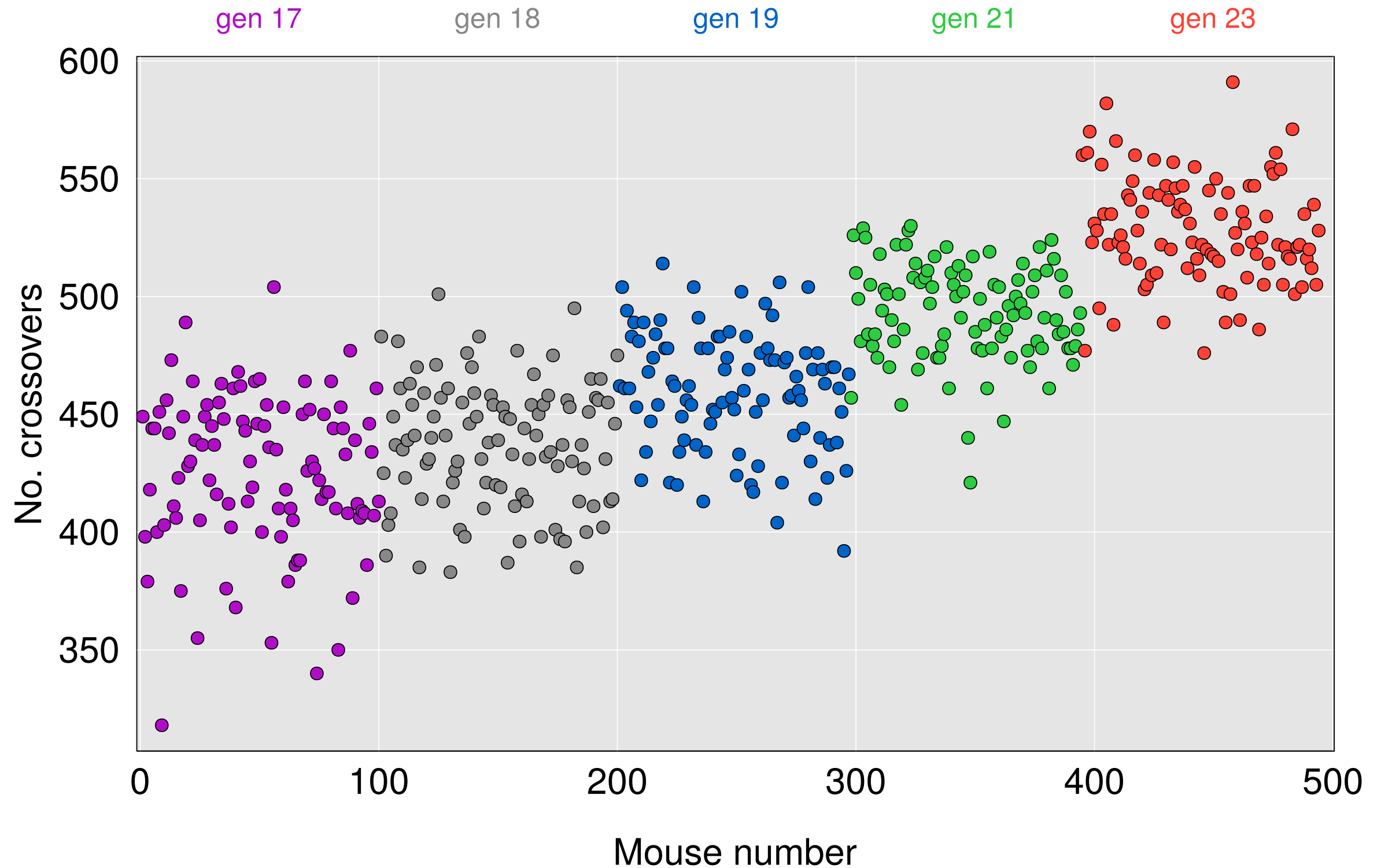
Percent missing vs number of crossovers



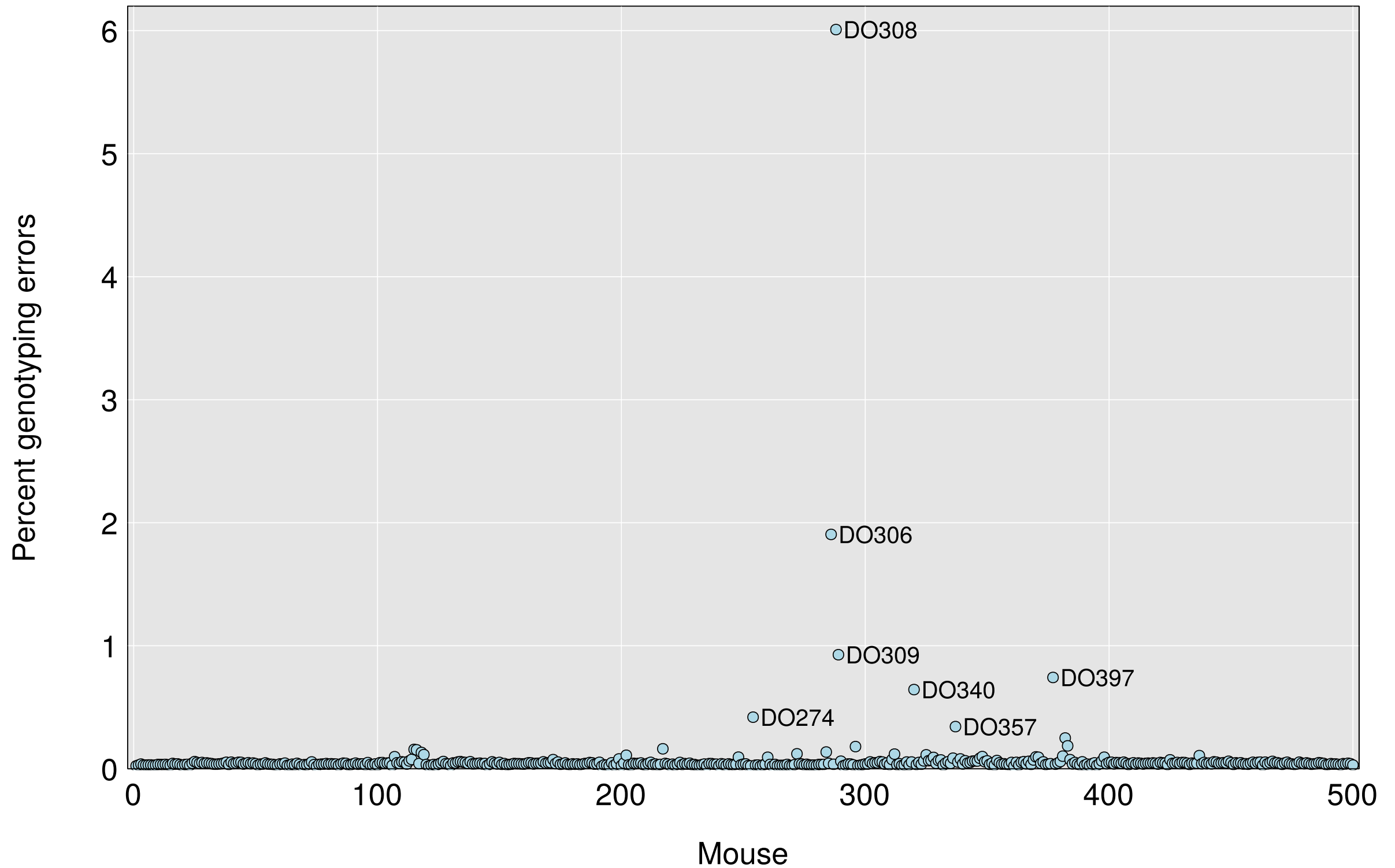
Percent missing vs number of crossovers



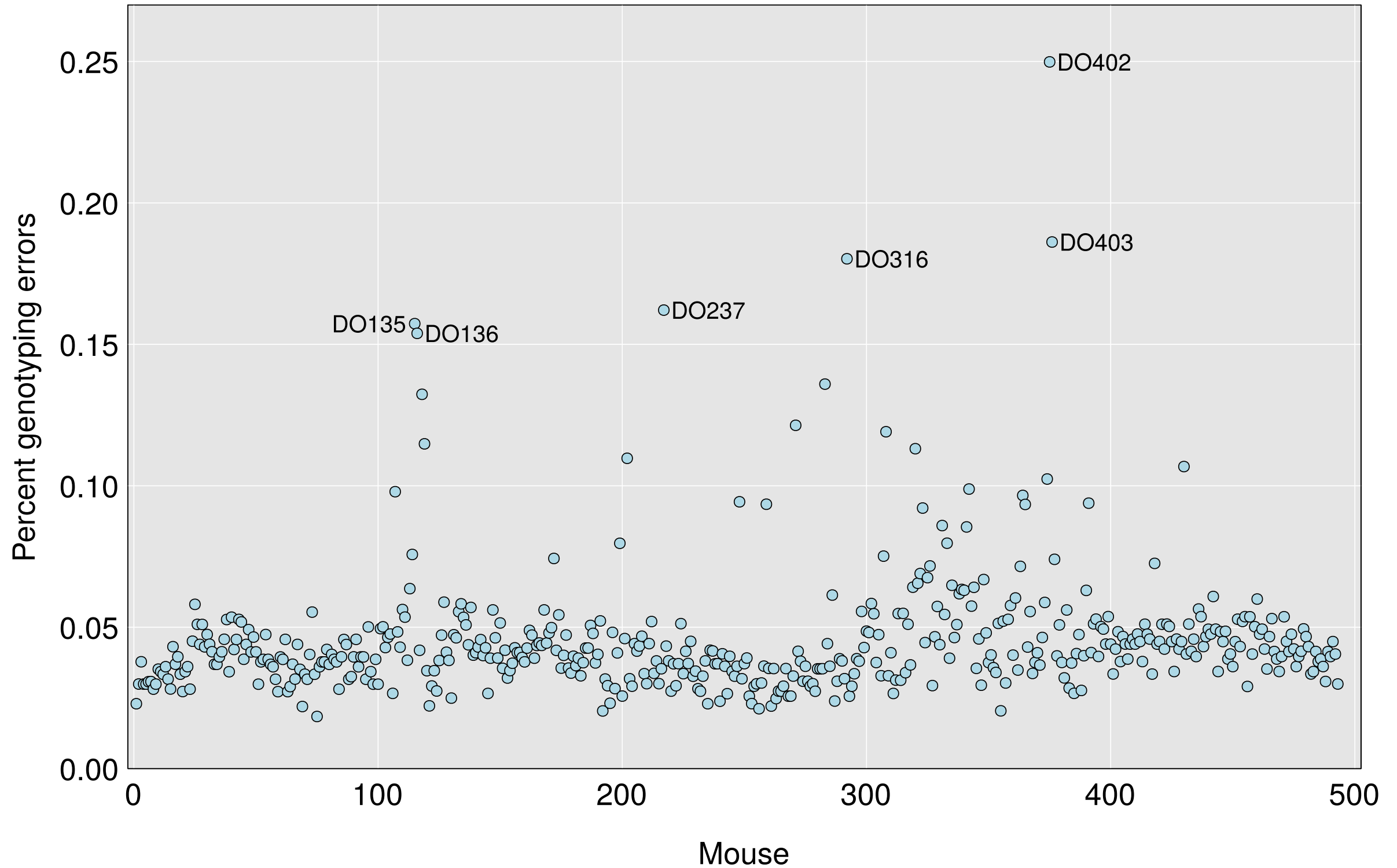
No. crossovers by generation



Estimated percent of genotyping errors



Estimated percent of genotyping errors



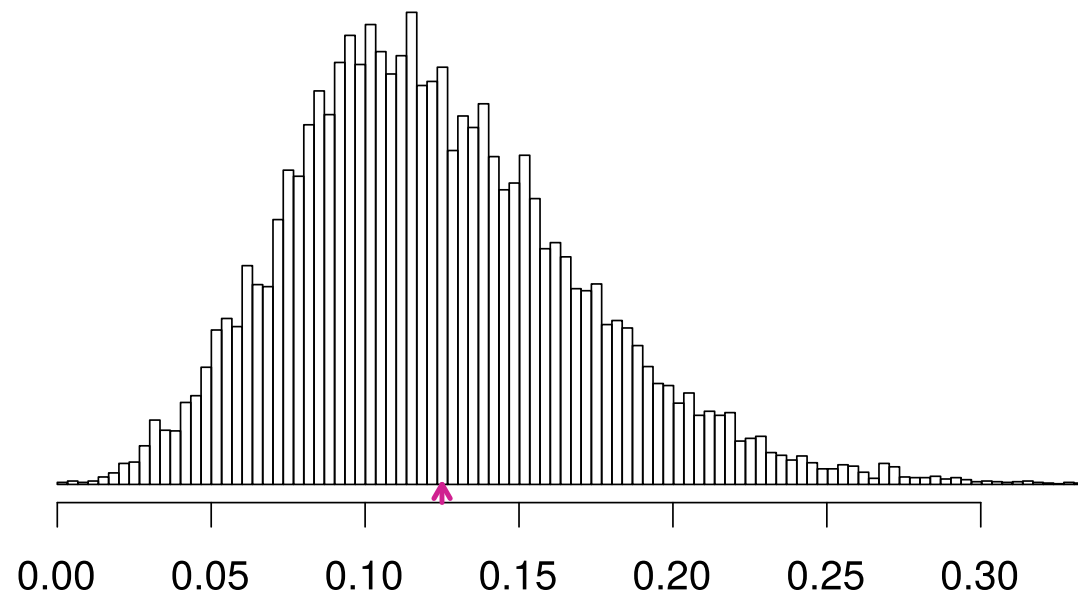
Marker quality

Proportion missing data



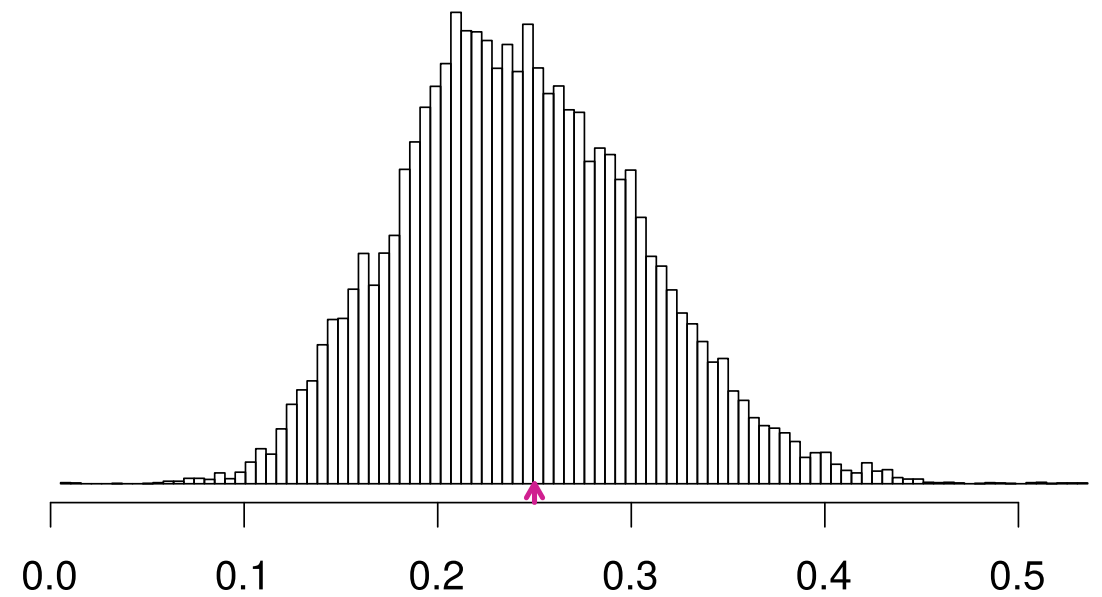
Allele frequencies, by marker

founder MAF = 1/8



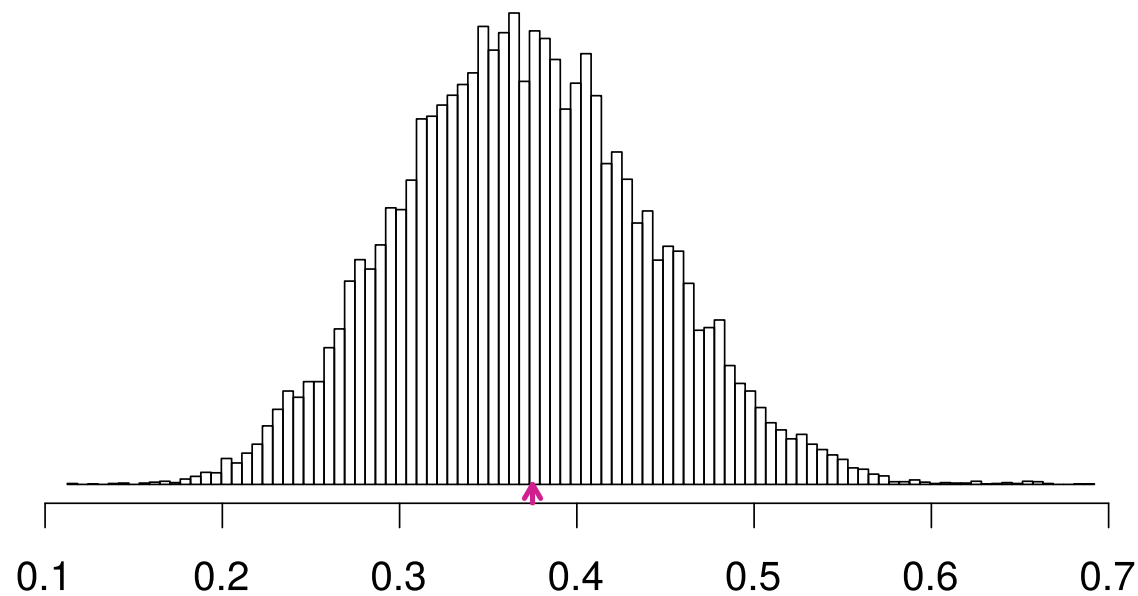
Frequency of minor allele

founder MAF = 2/8



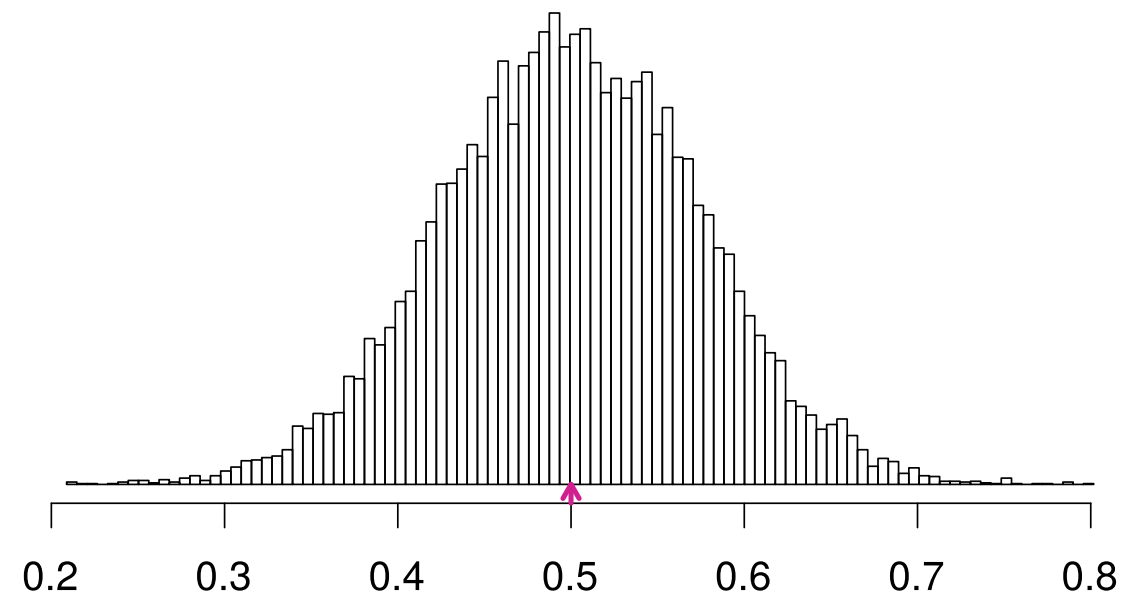
Frequency of minor allele

founder MAF = 3/8



Frequency of minor allele

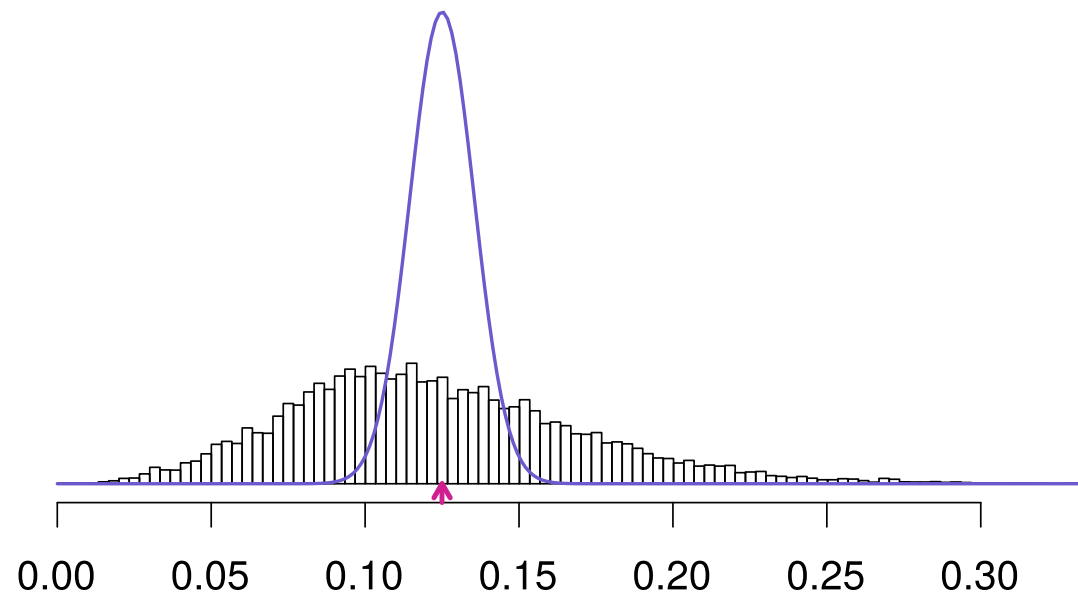
founder MAF = 4/8



Frequency of minor allele

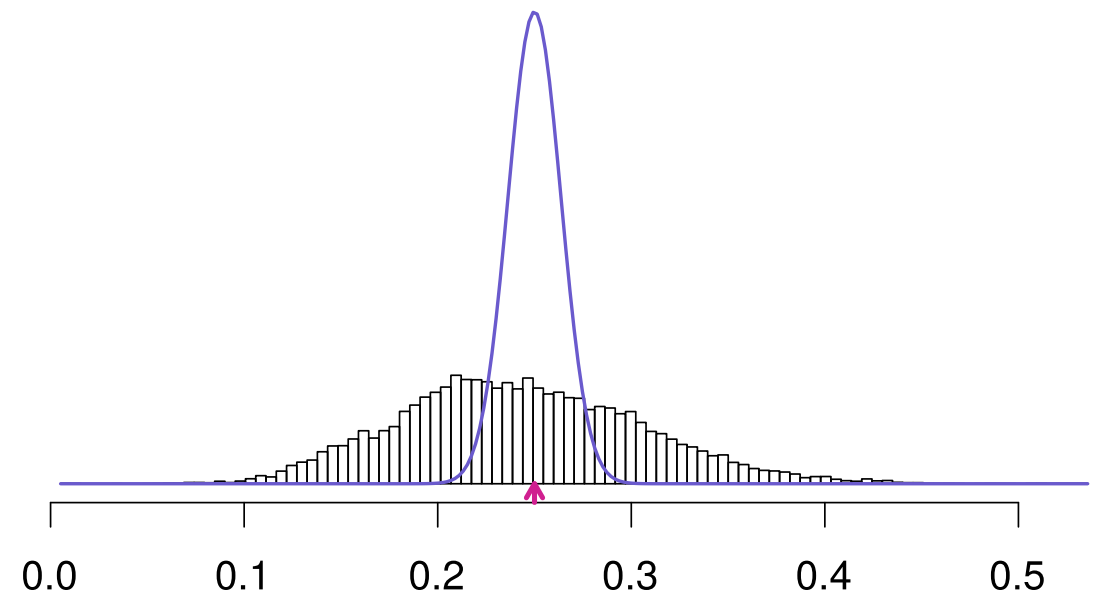
Allele frequencies, by marker

founder MAF = $1/8$



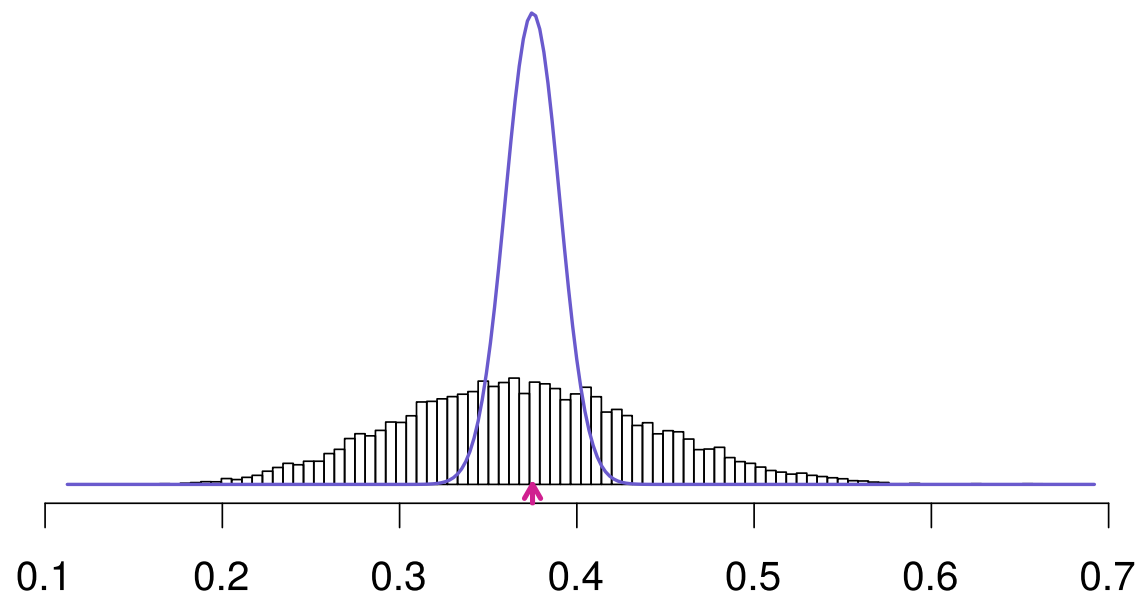
Frequency of minor allele

founder MAF = $2/8$



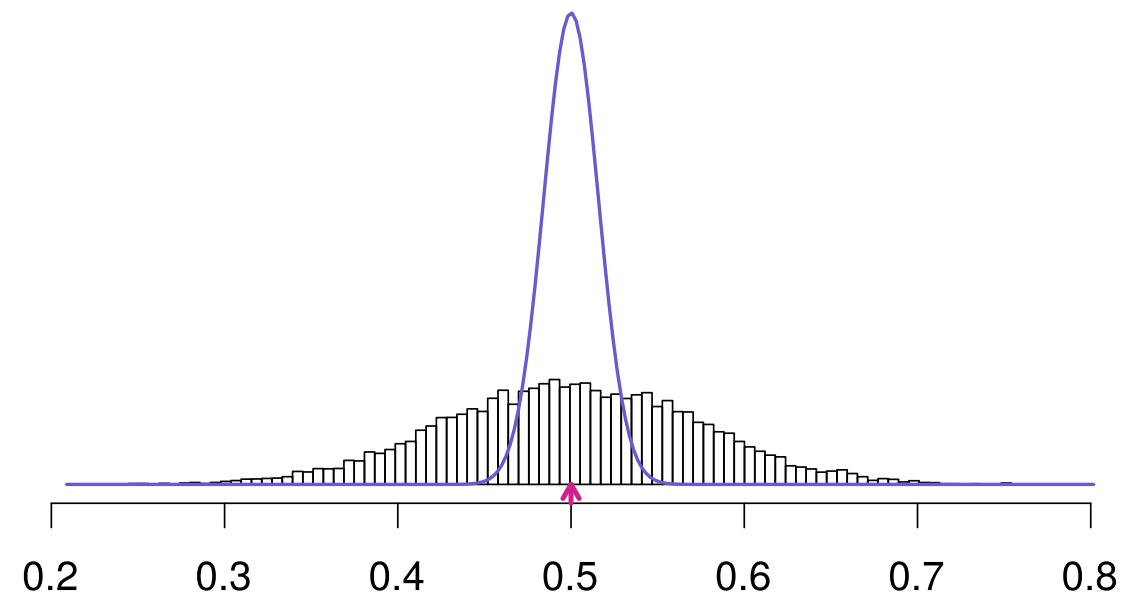
Frequency of minor allele

founder MAF = $3/8$



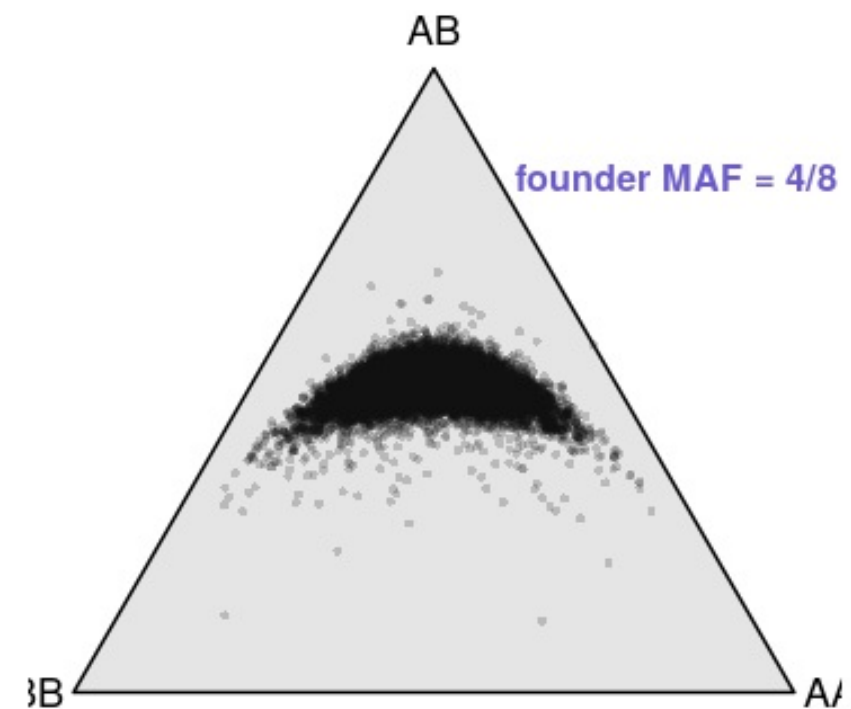
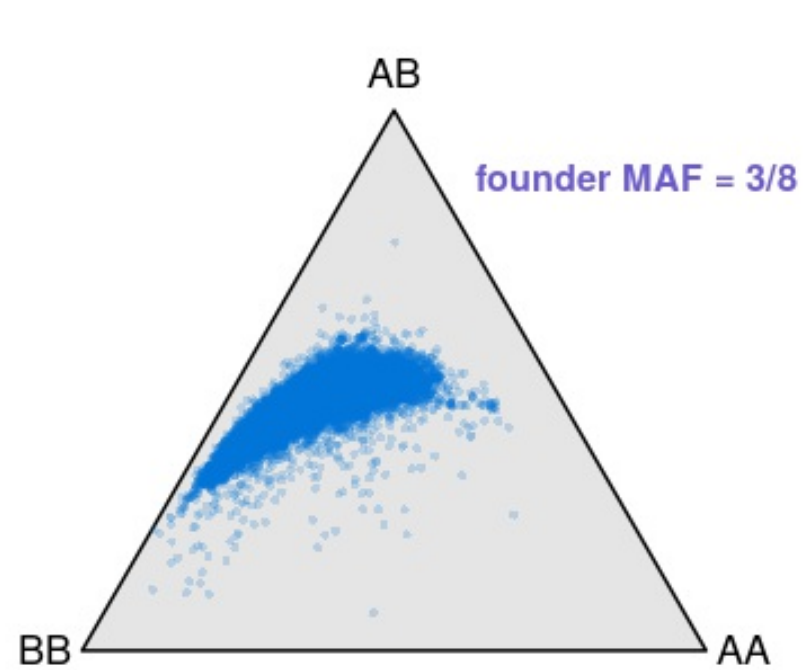
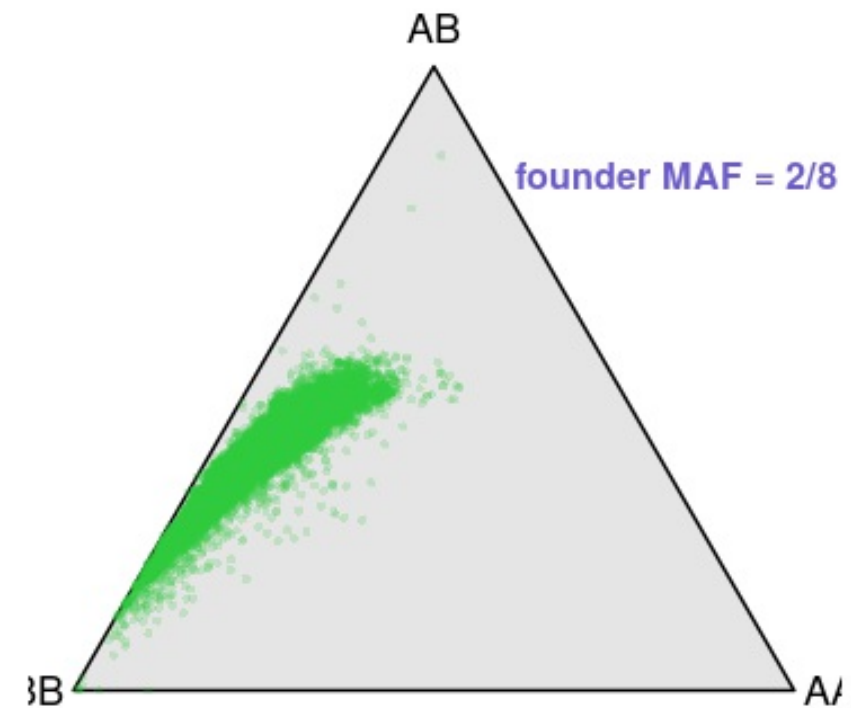
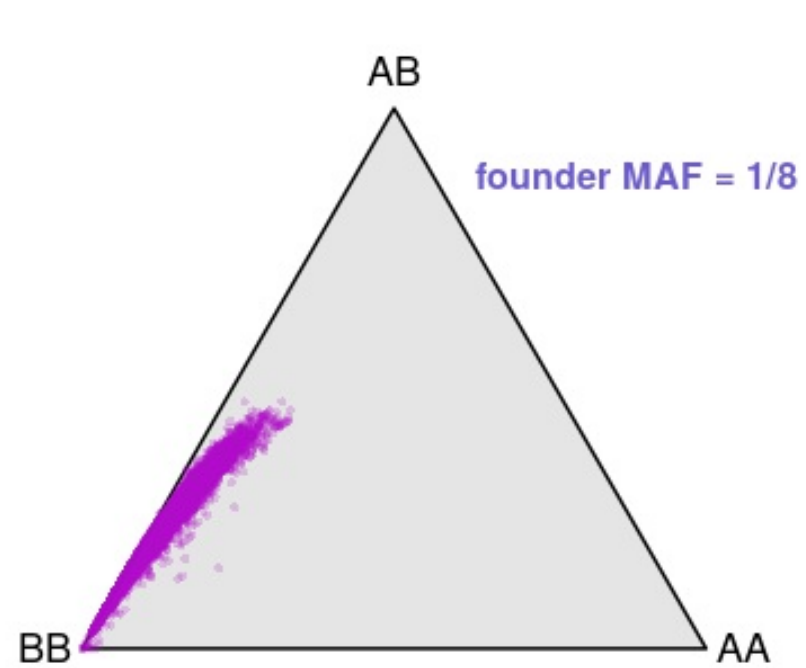
Frequency of minor allele

founder MAF = $4/8$



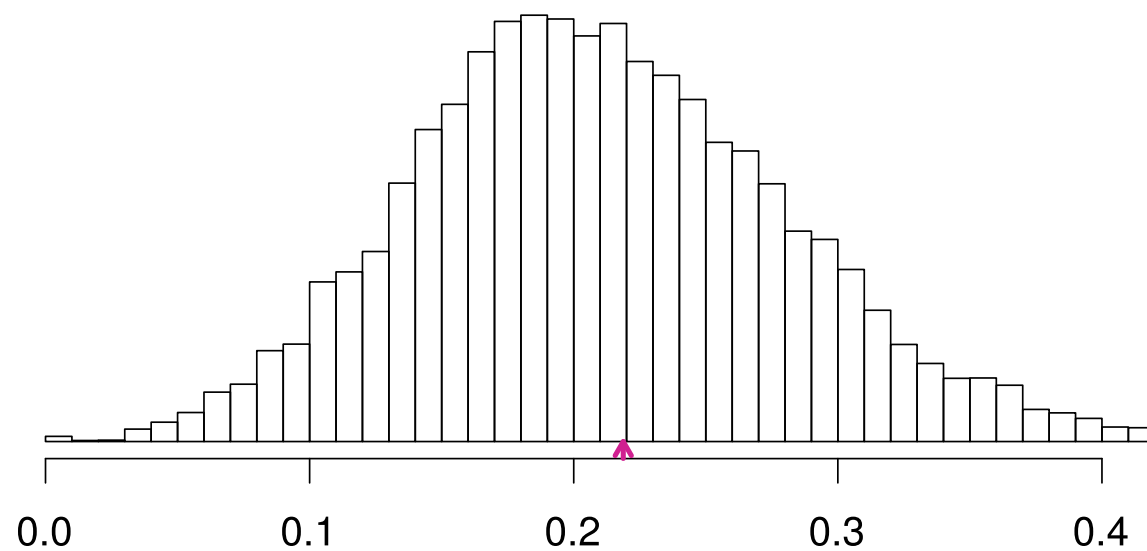
Frequency of minor allele

Genotype frequencies, by marker



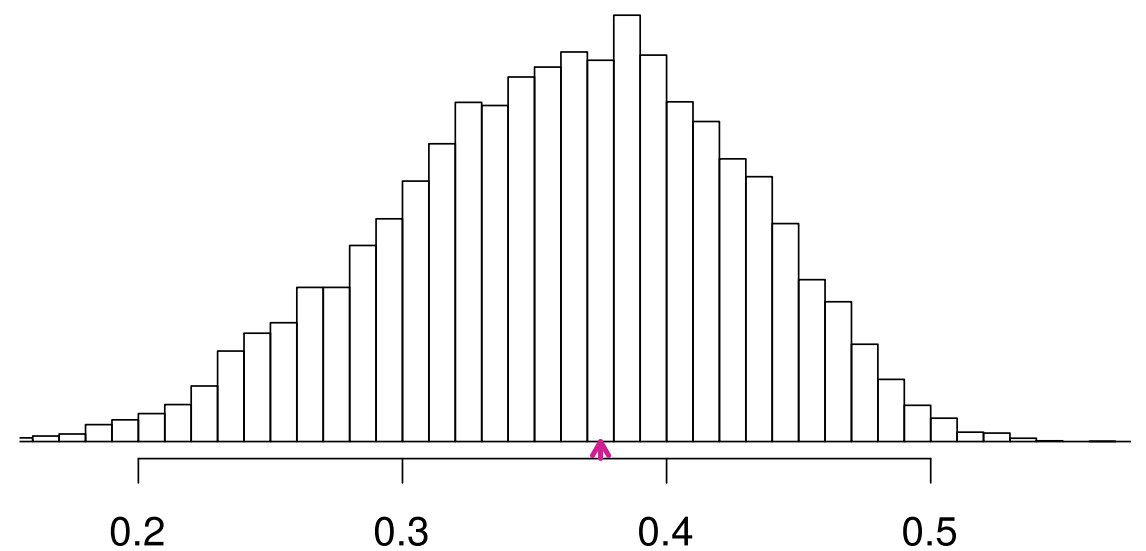
Heterozygosities, by marker

founder MAF = 1/8



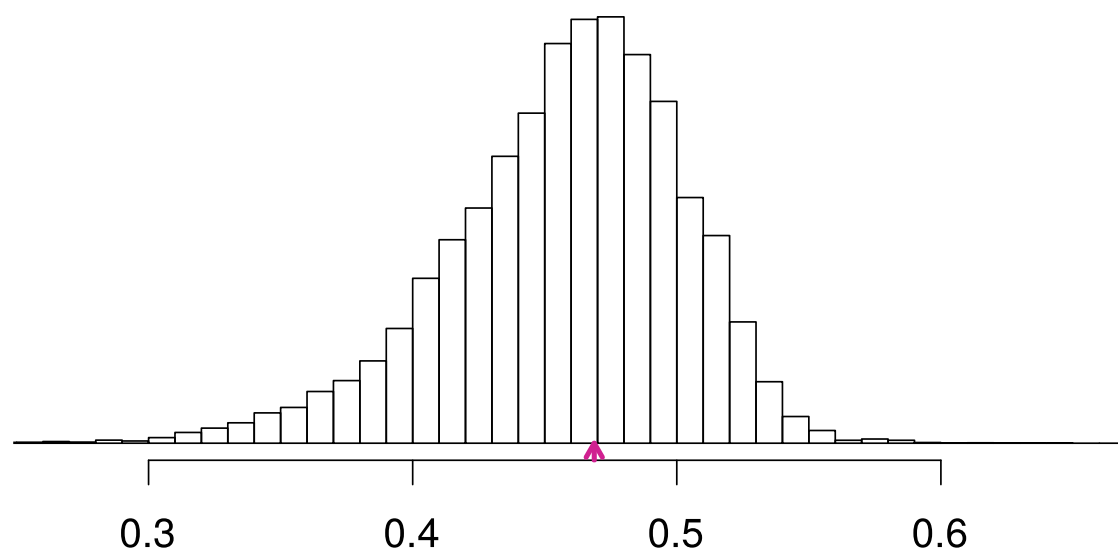
Frequency of minor allele

founder MAF = 2/8



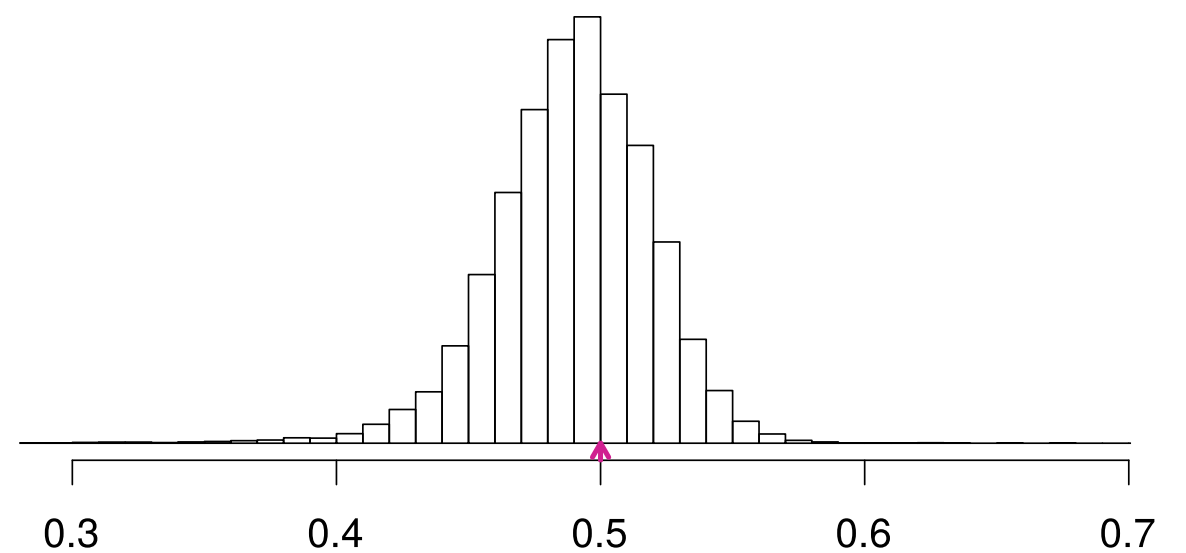
Frequency of minor allele

founder MAF = 3/8



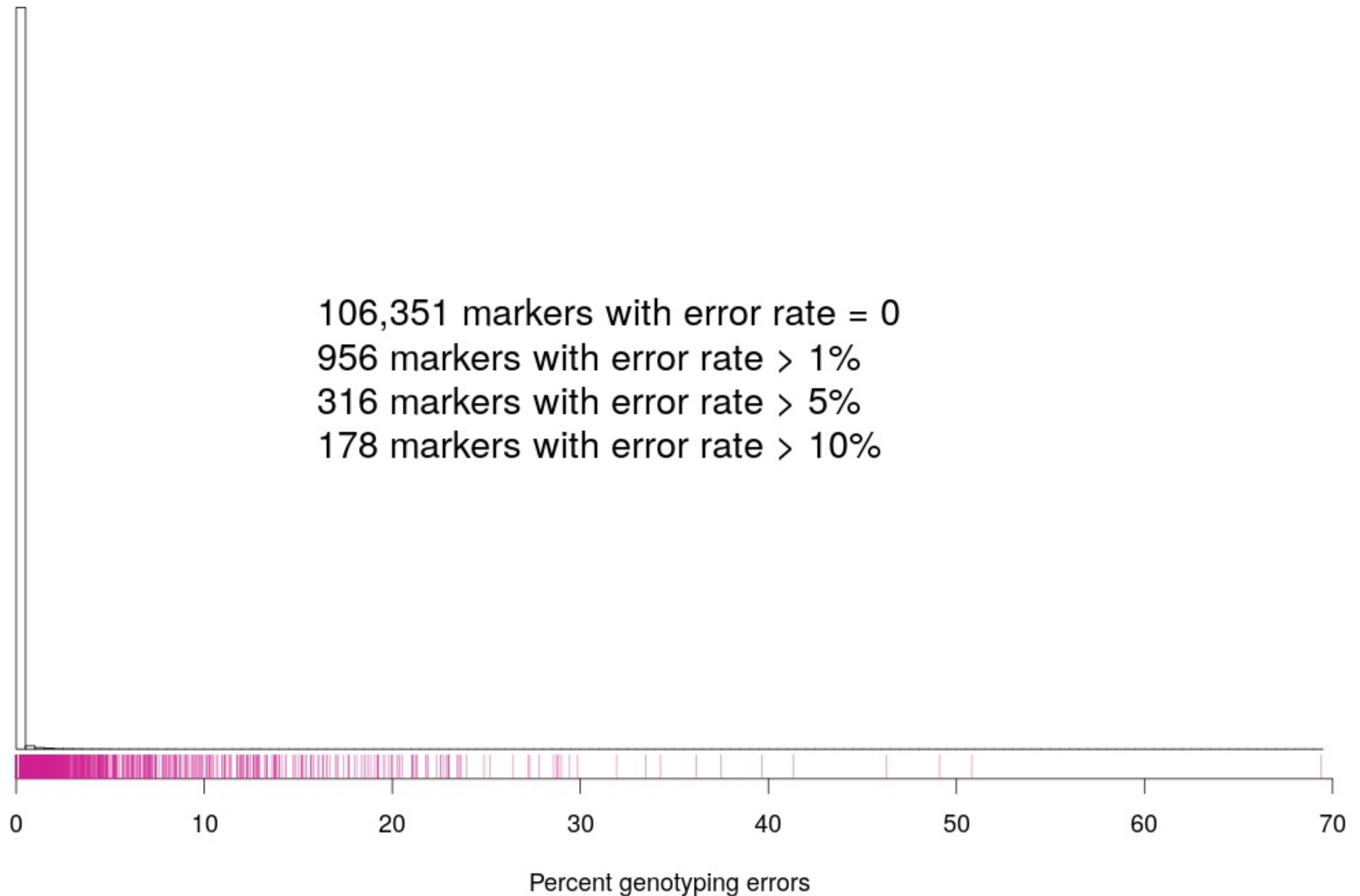
Frequency of minor allele

founder MAF = 4/8

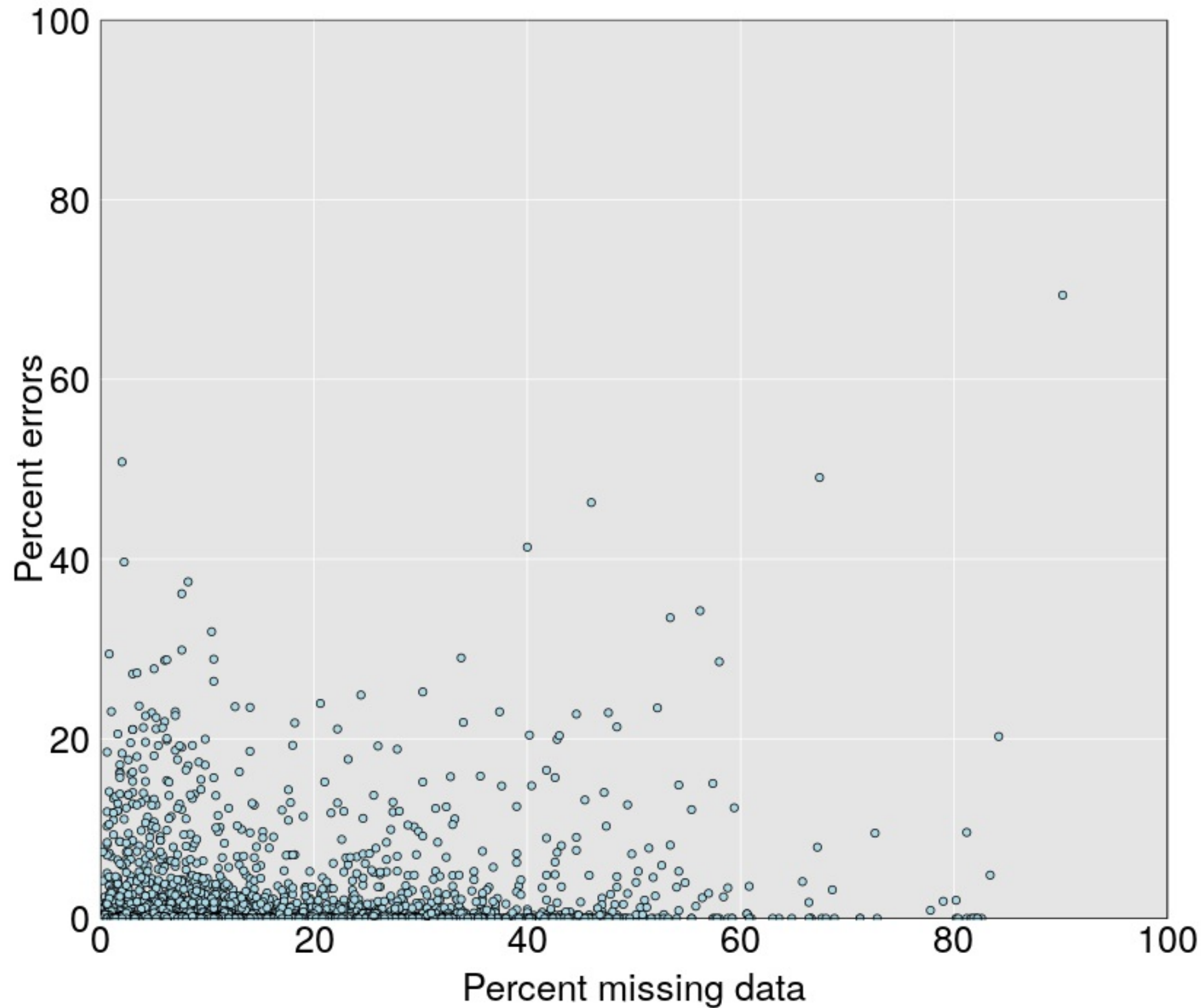


Frequency of minor allele

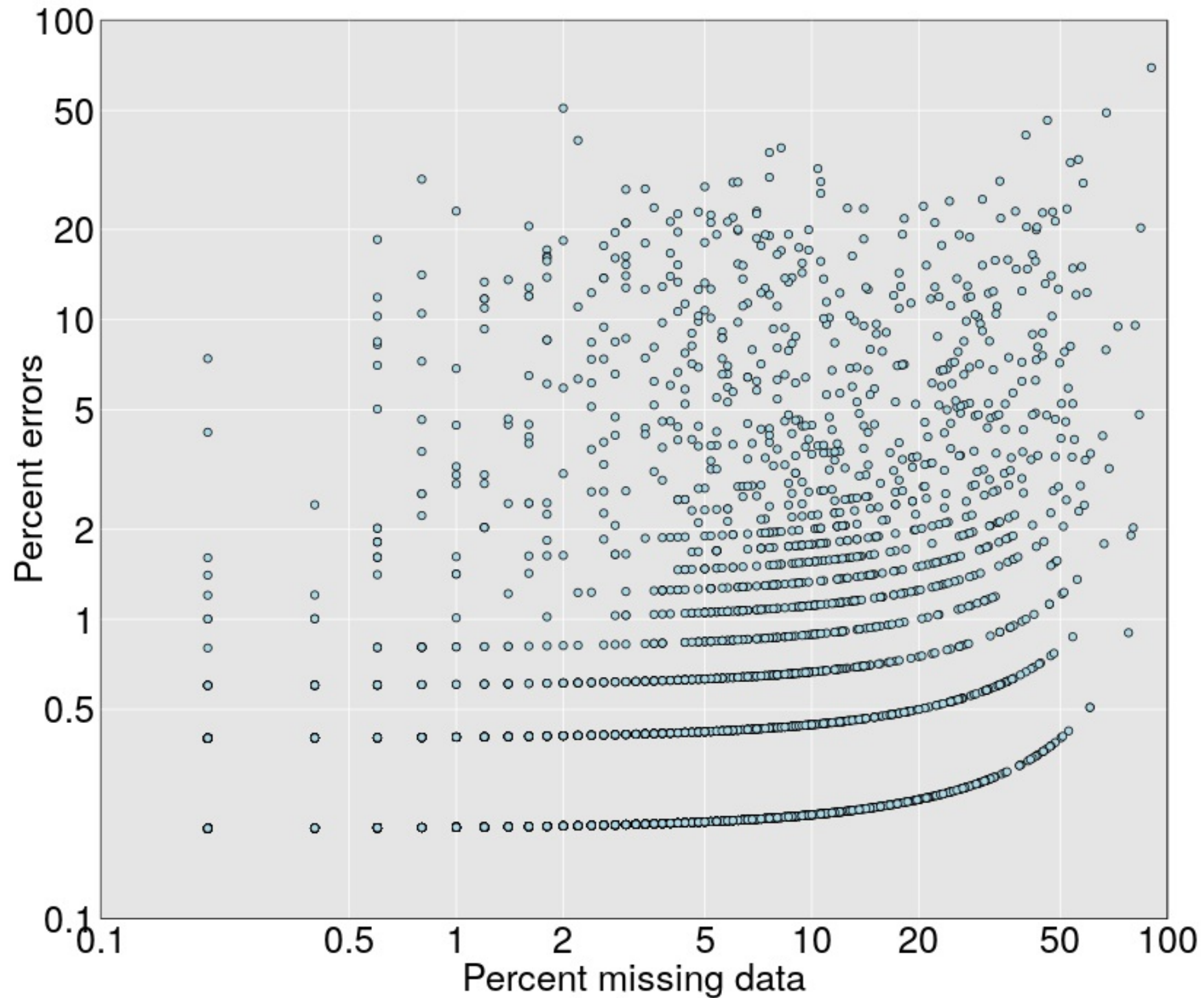
Genotyping error rates



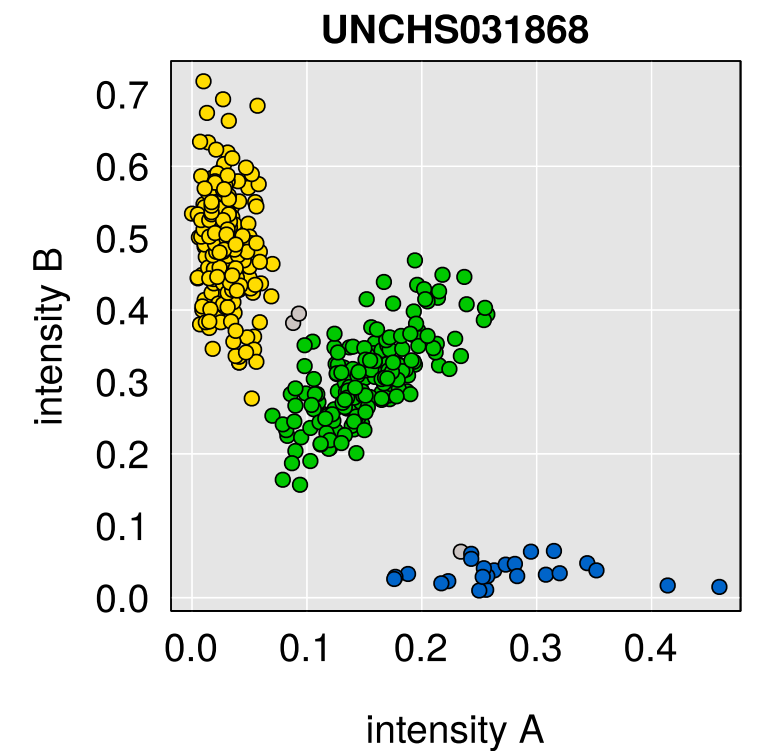
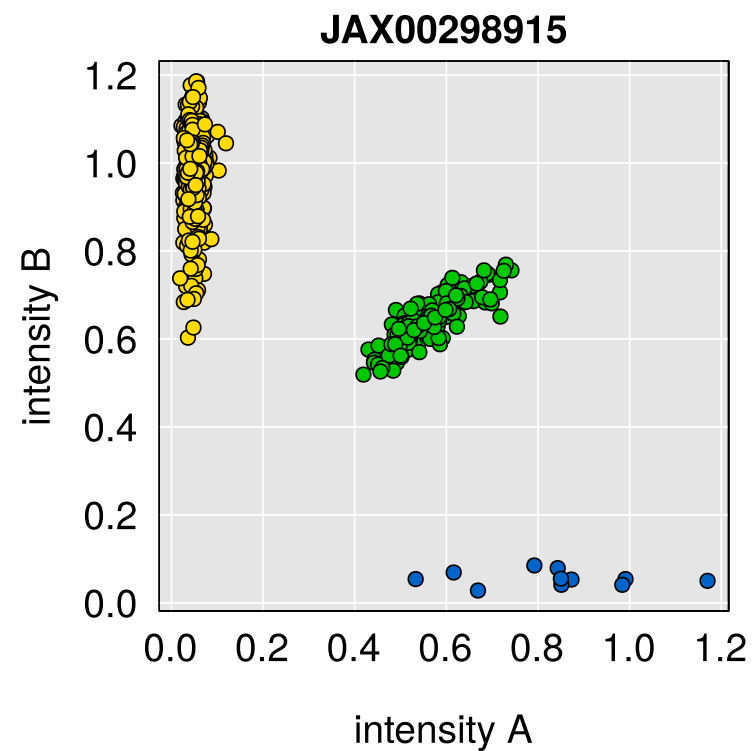
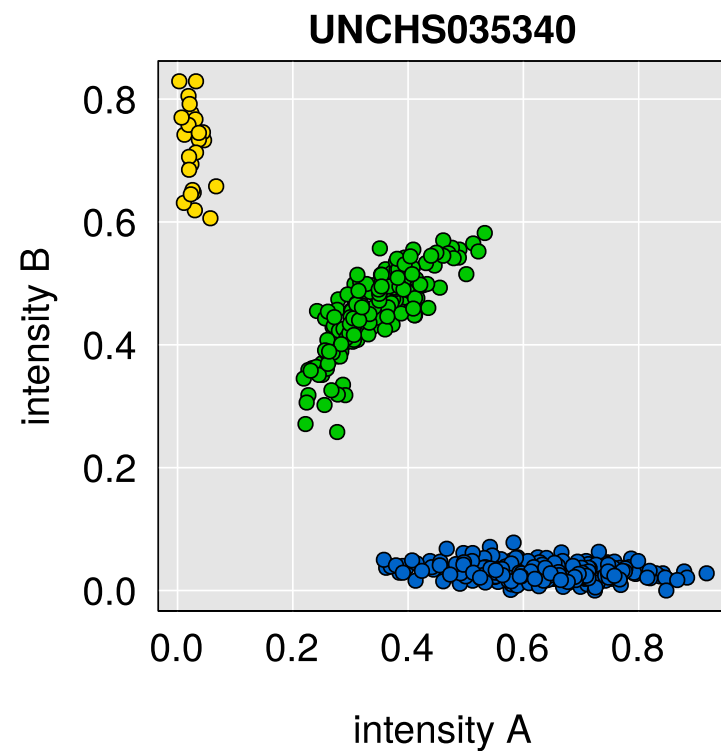
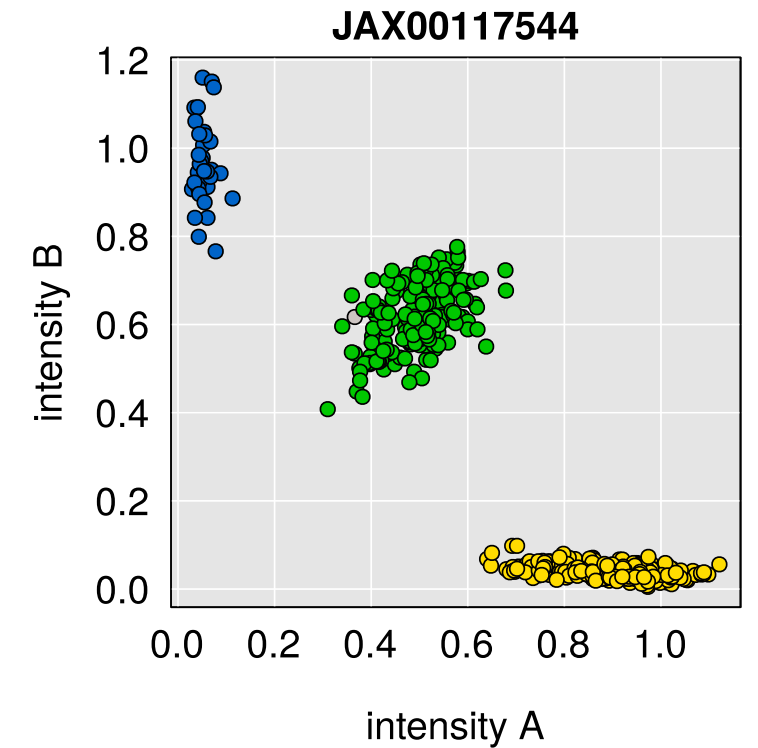
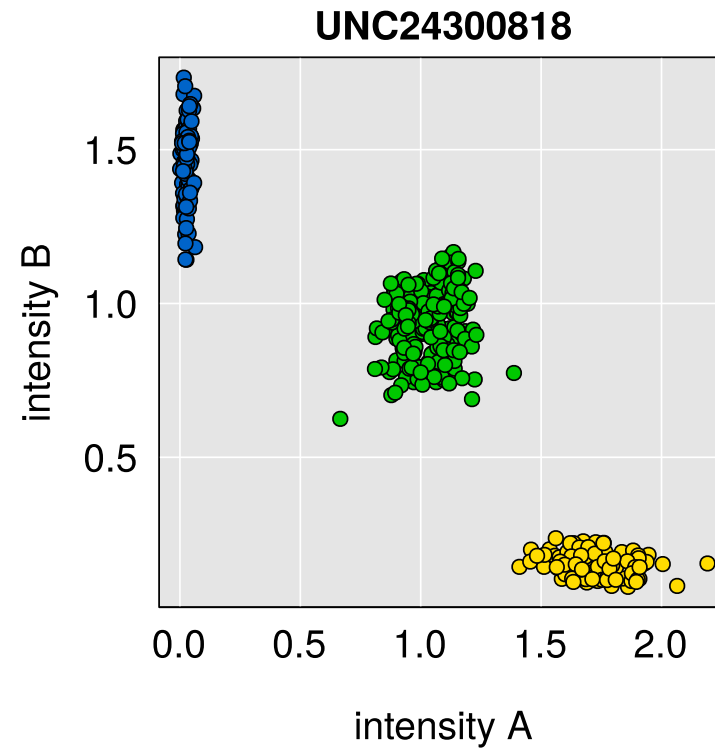
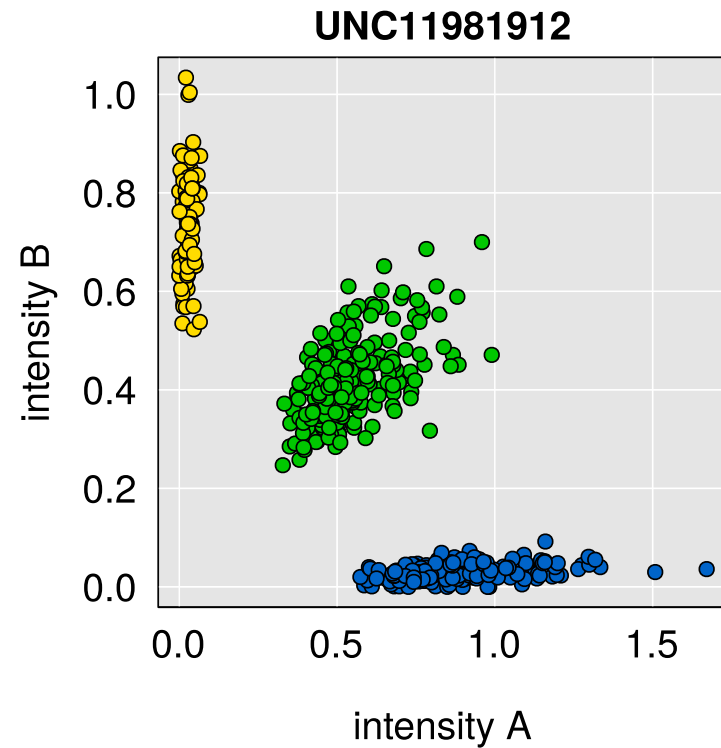
Genotyping error rate vs percent missing



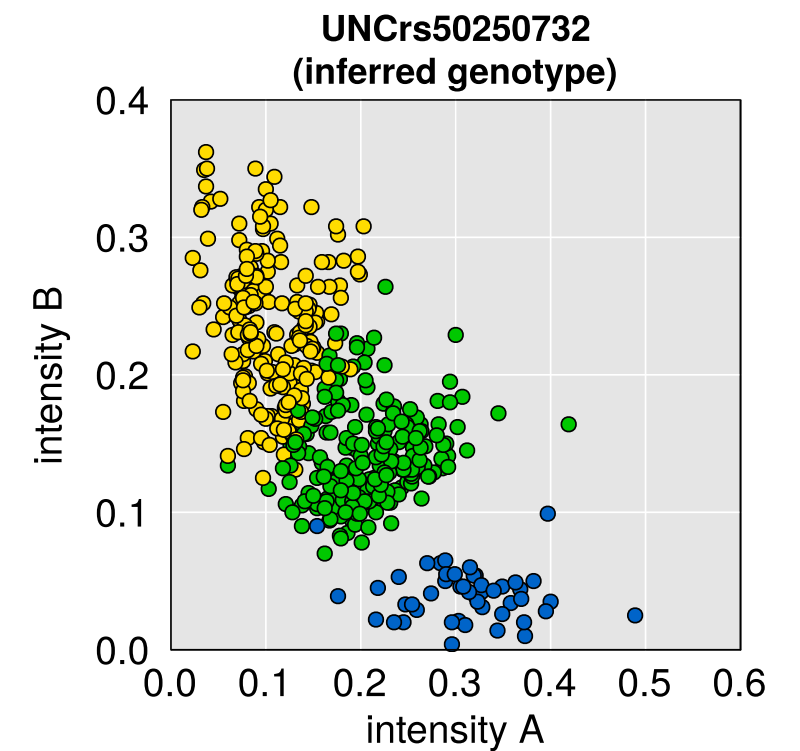
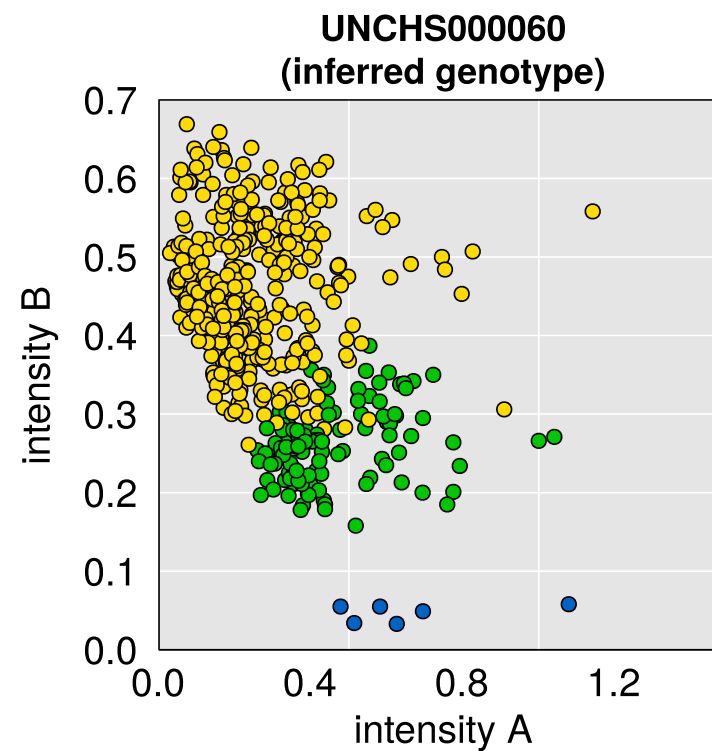
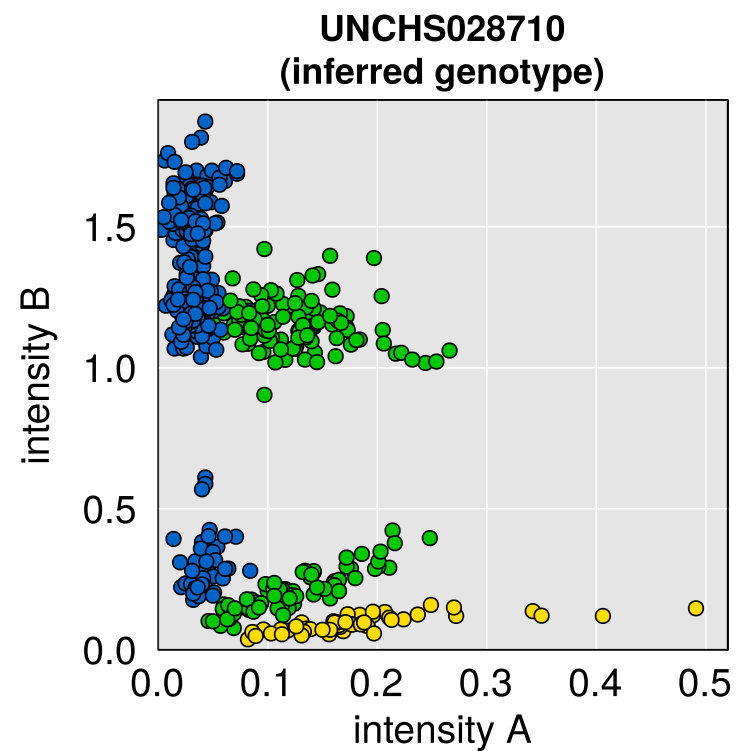
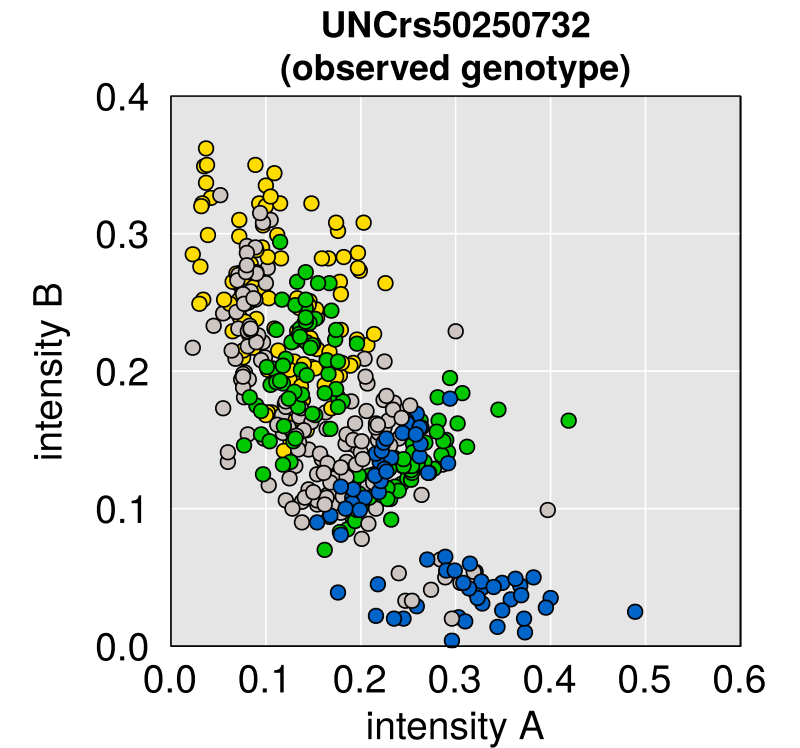
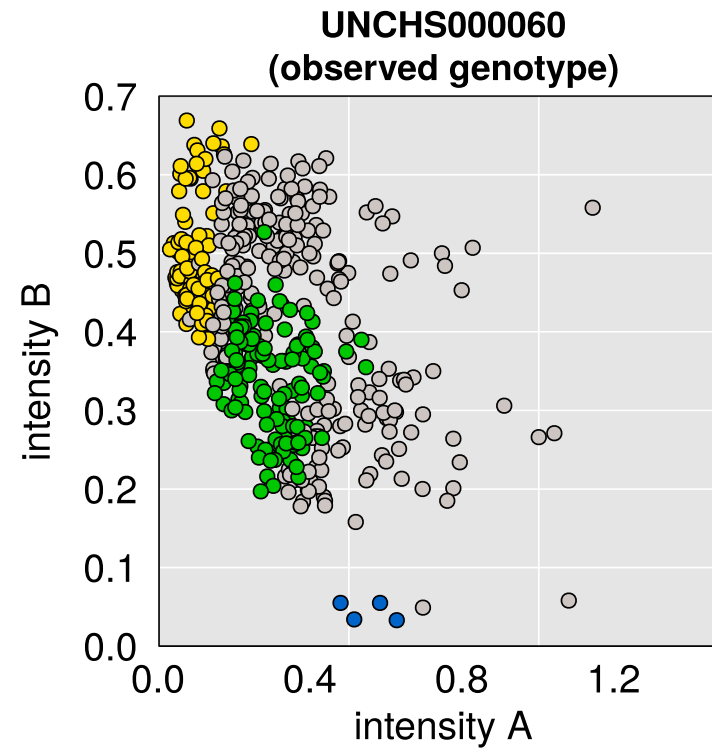
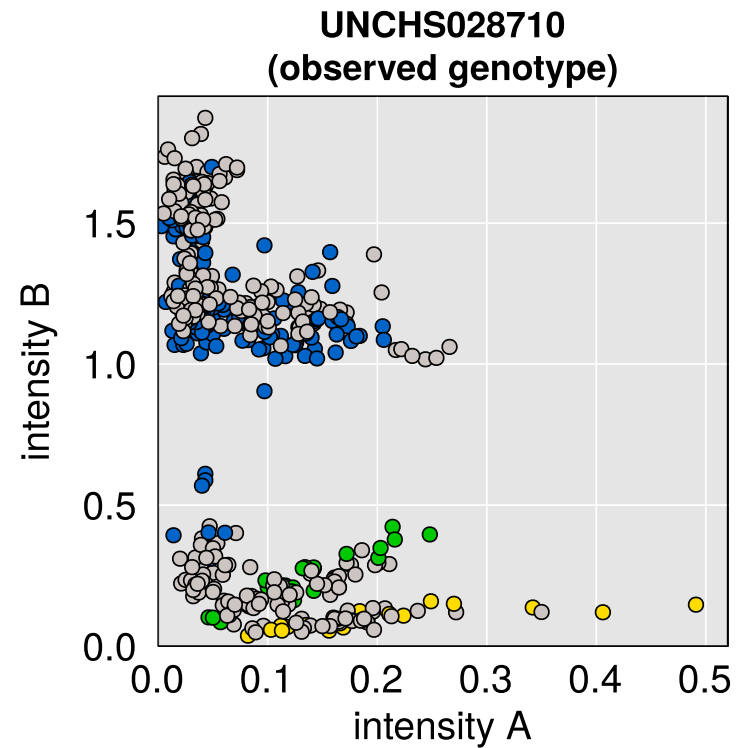
Genotyping error rate vs percent missing



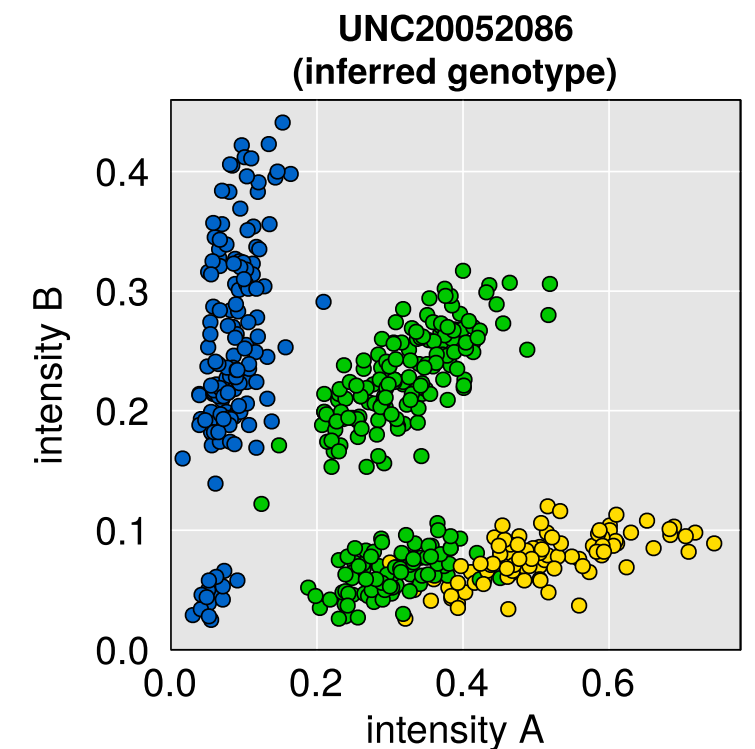
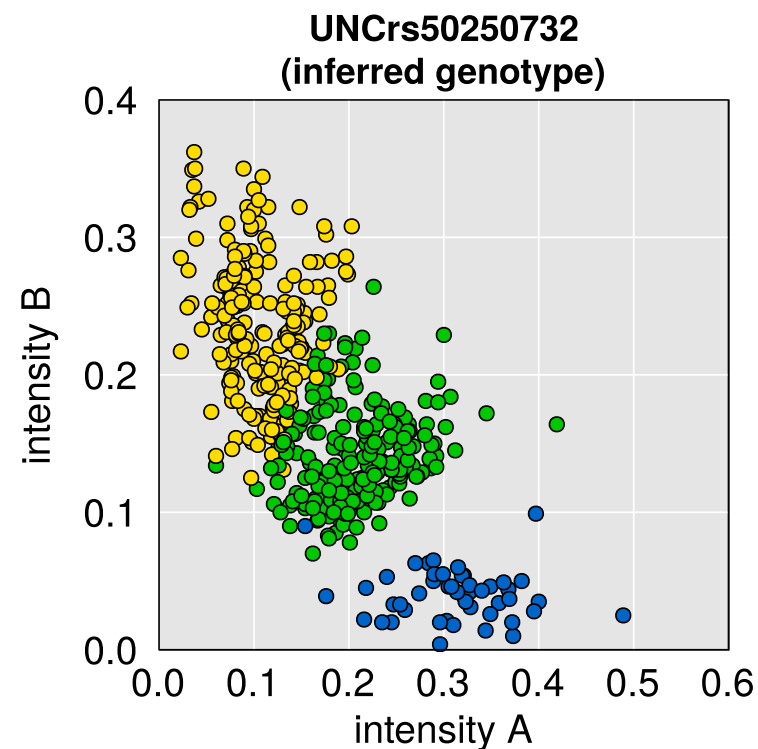
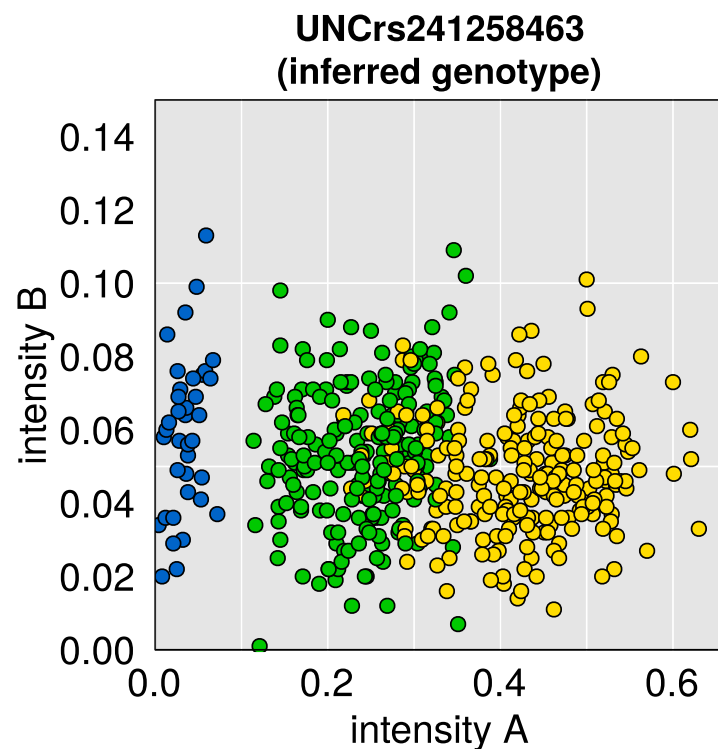
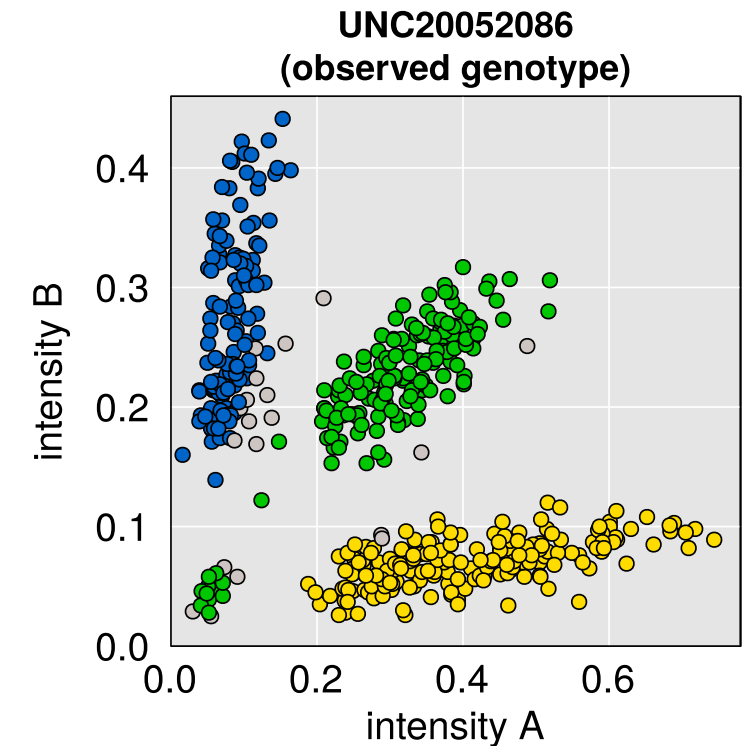
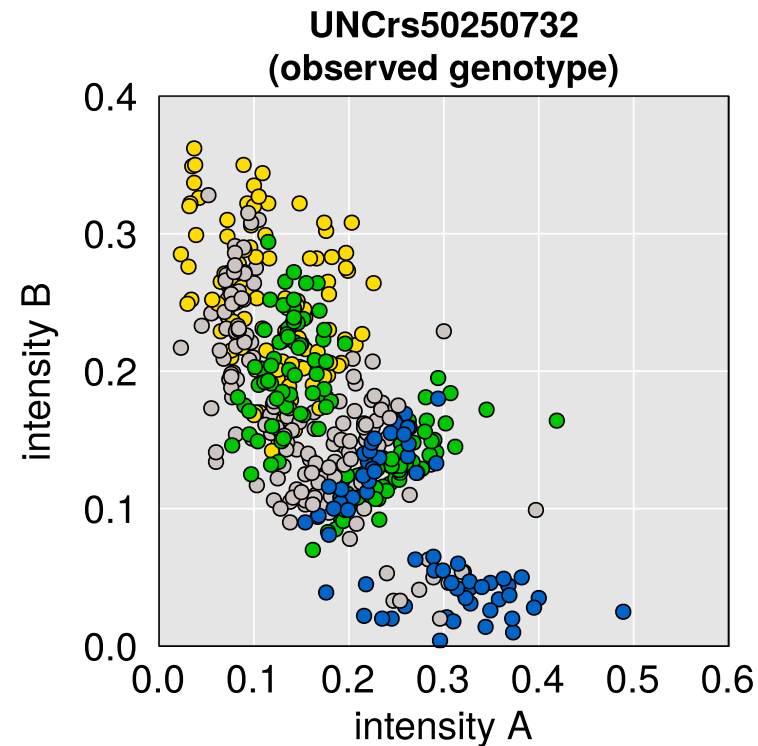
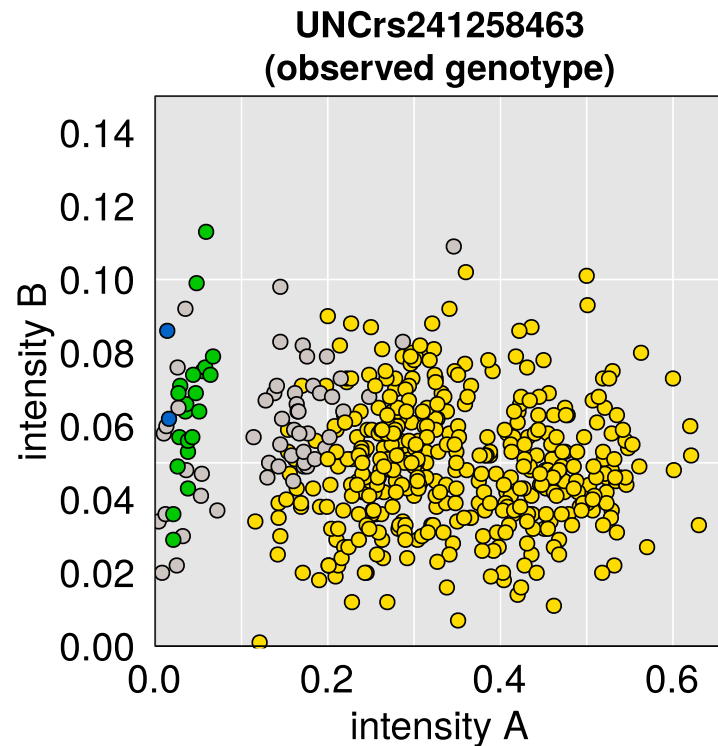
Nice markers



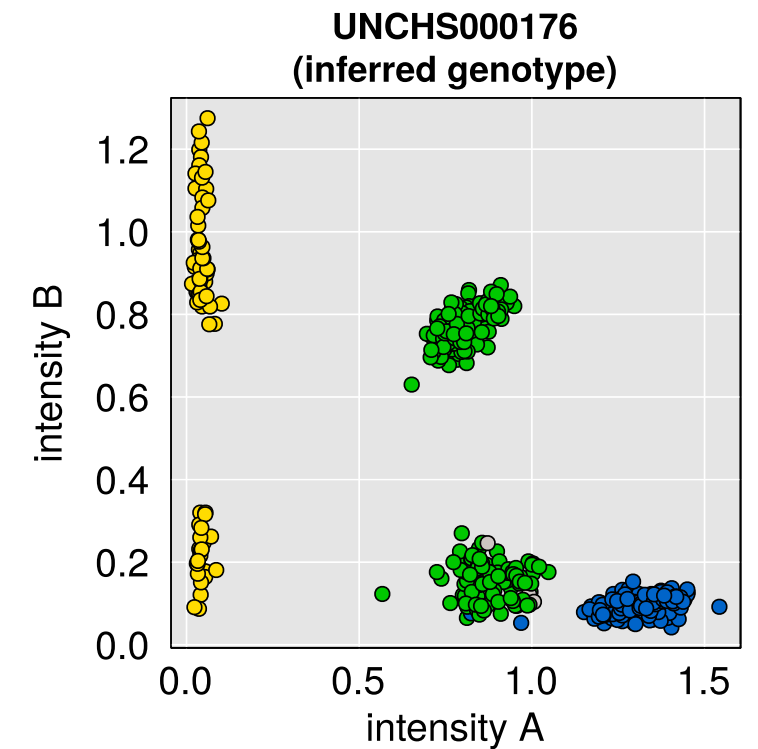
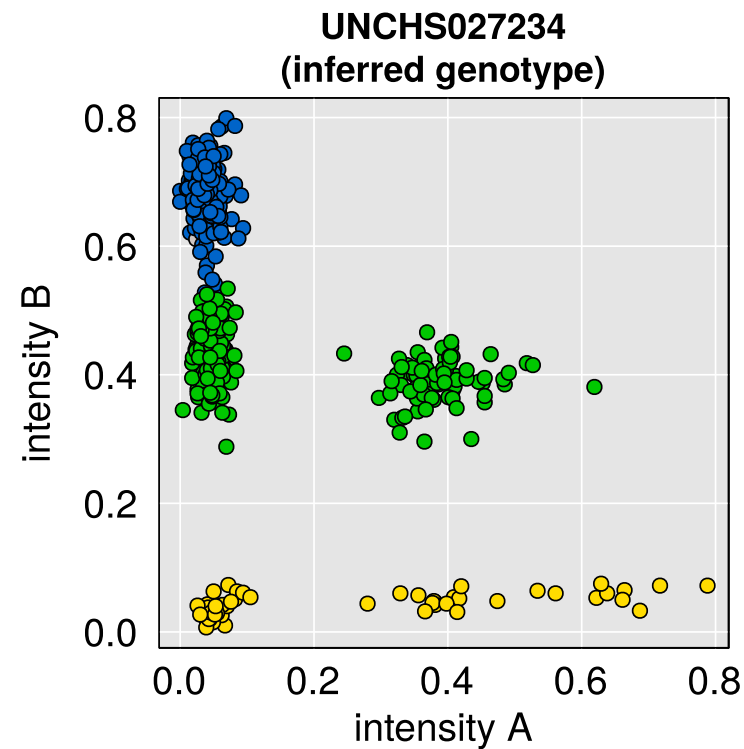
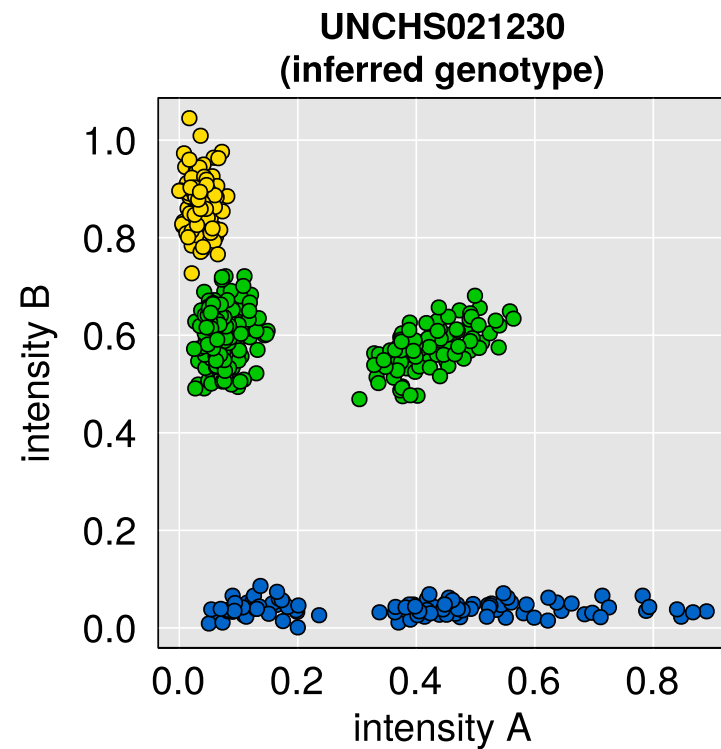
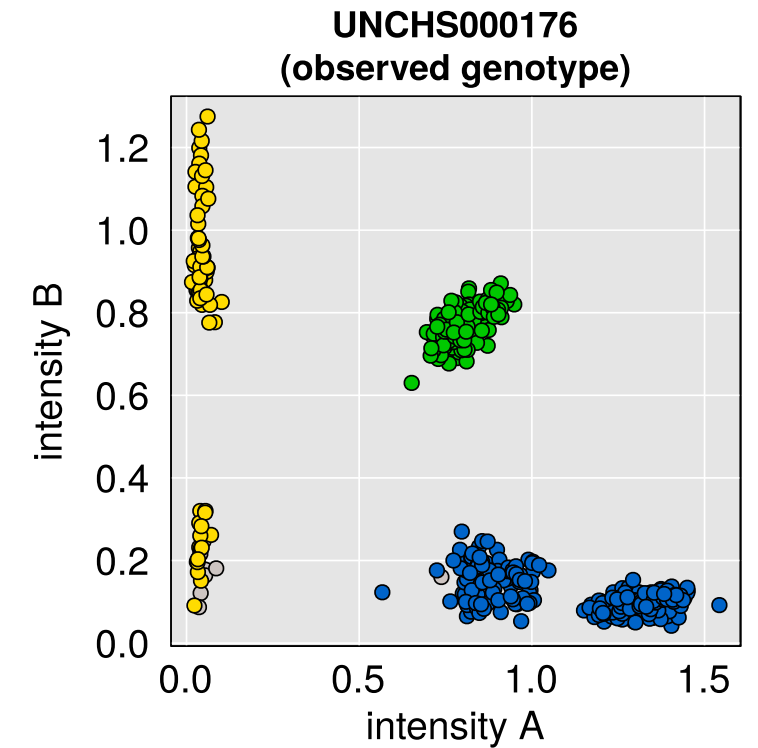
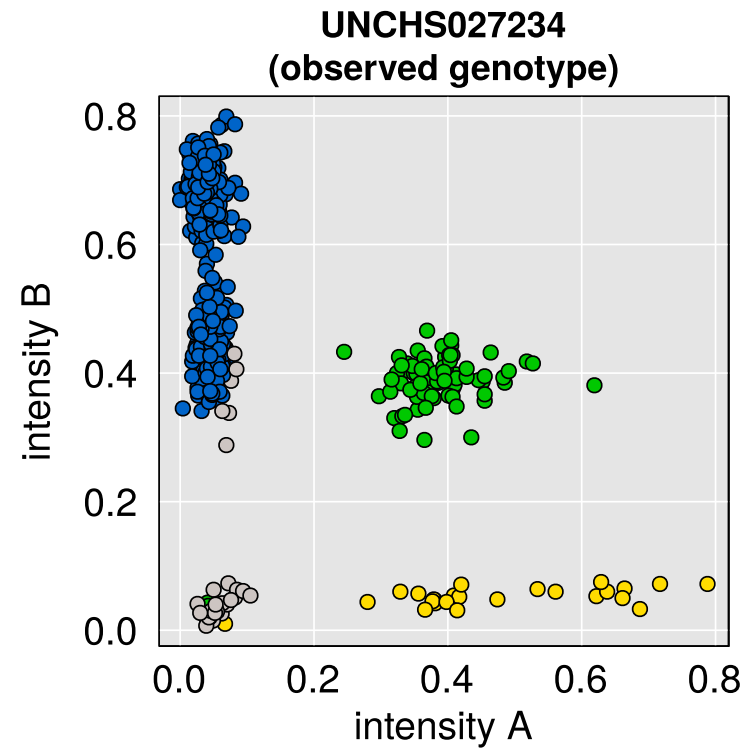
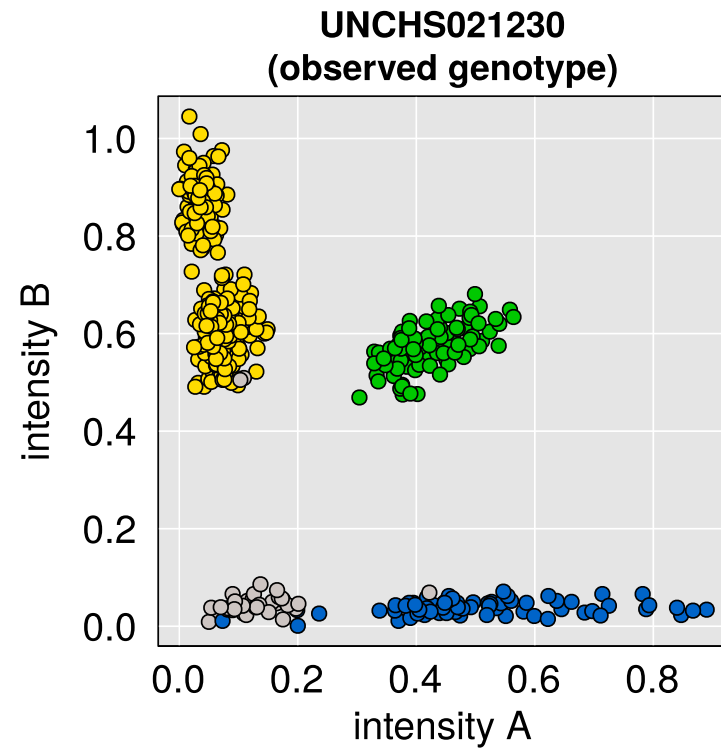
Crap markers



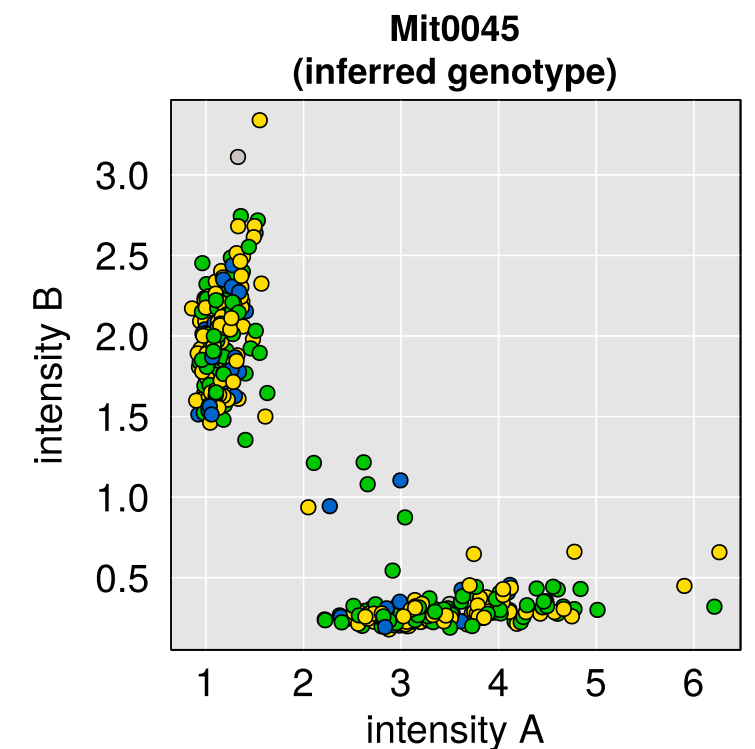
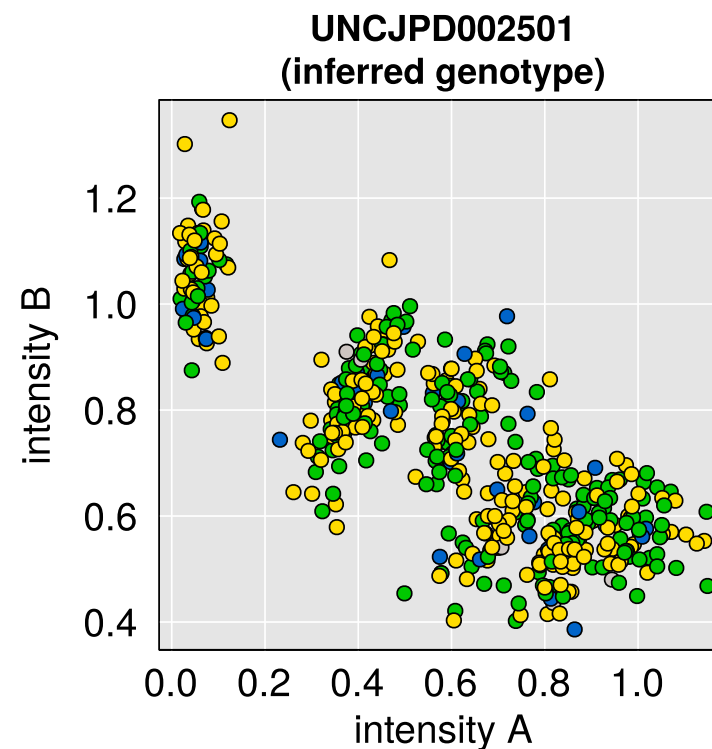
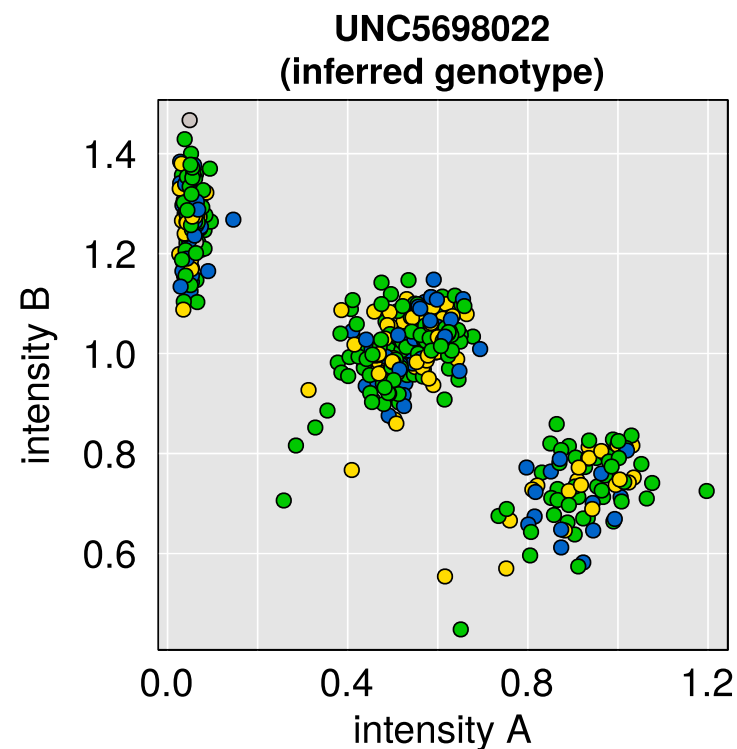
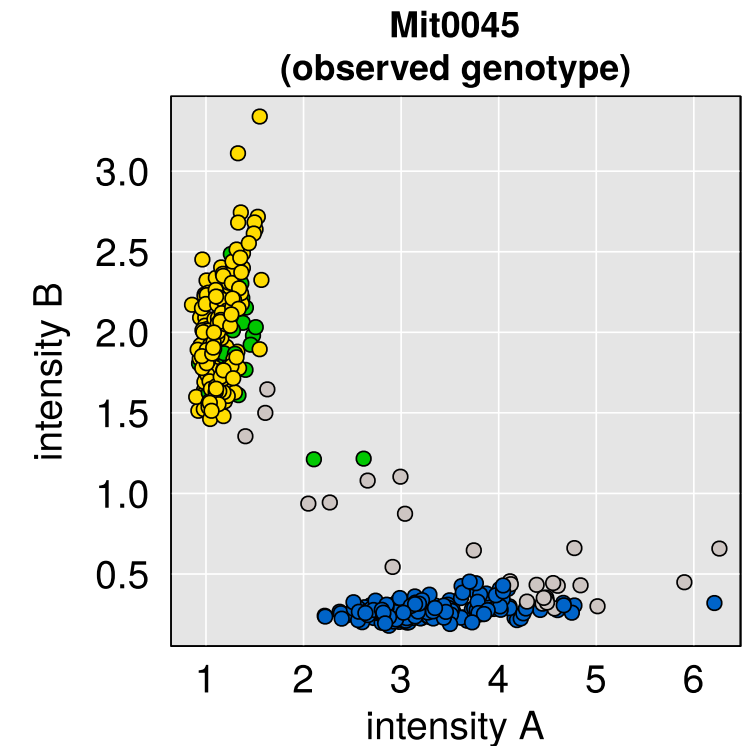
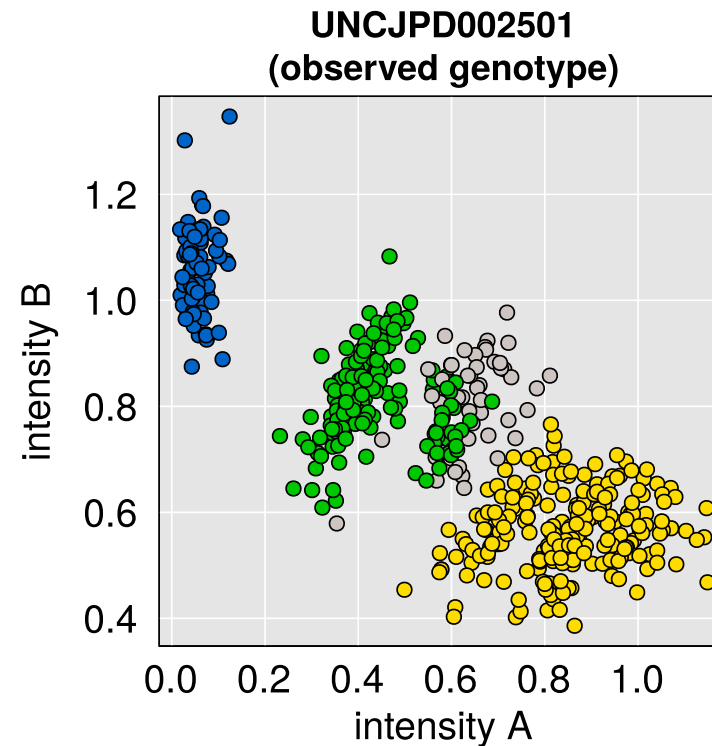
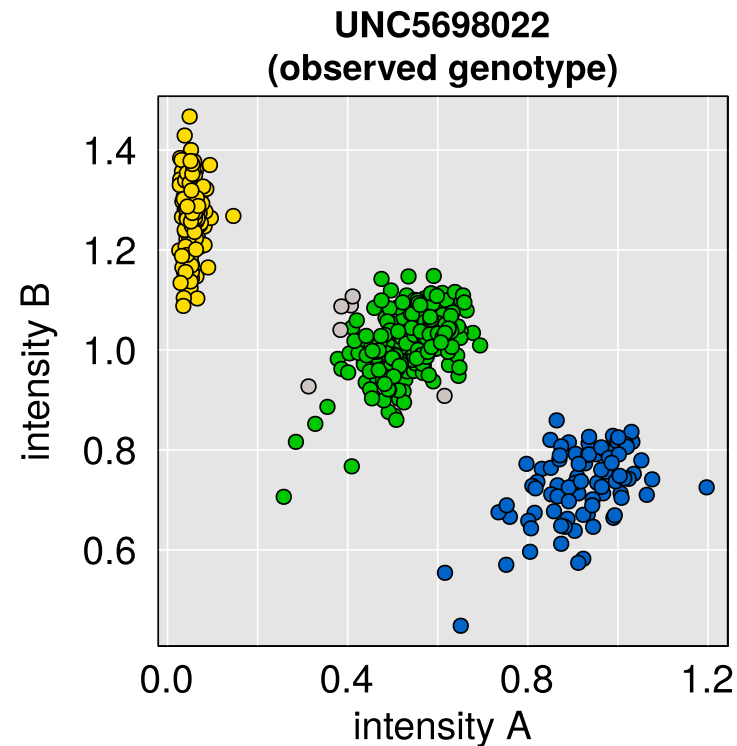
More crap markers



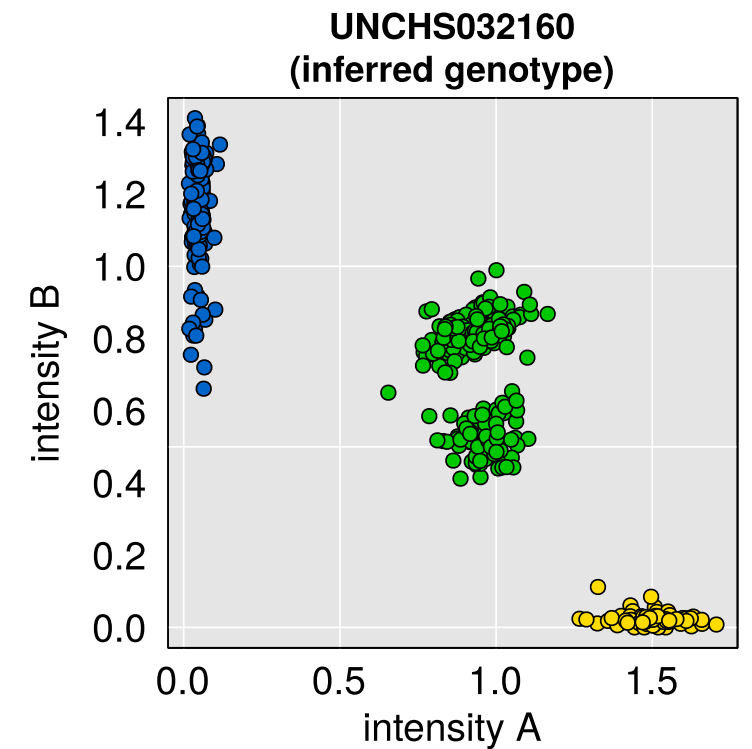
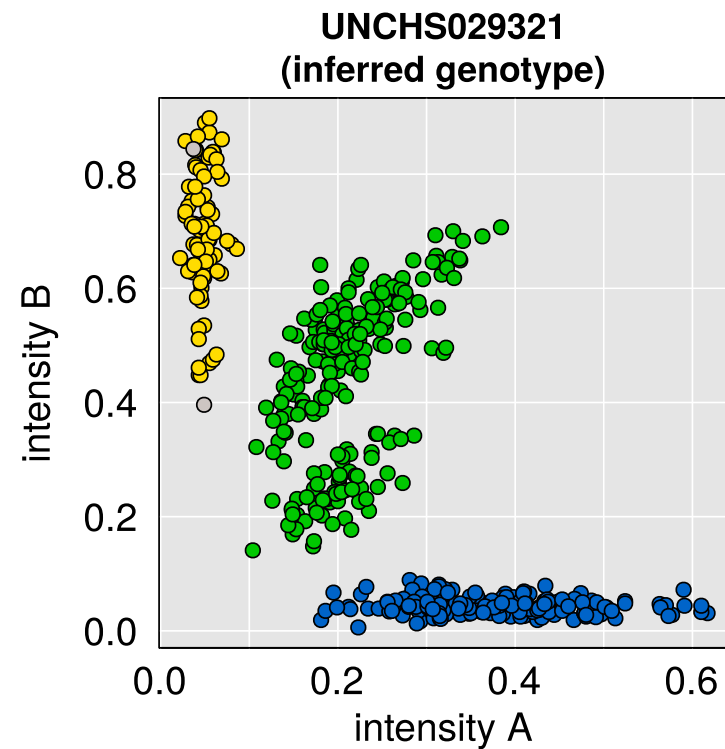
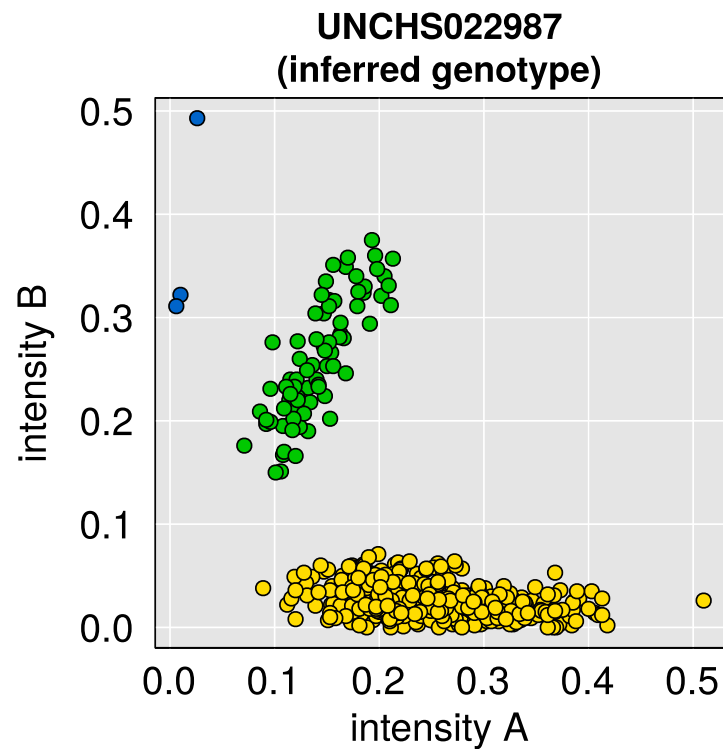
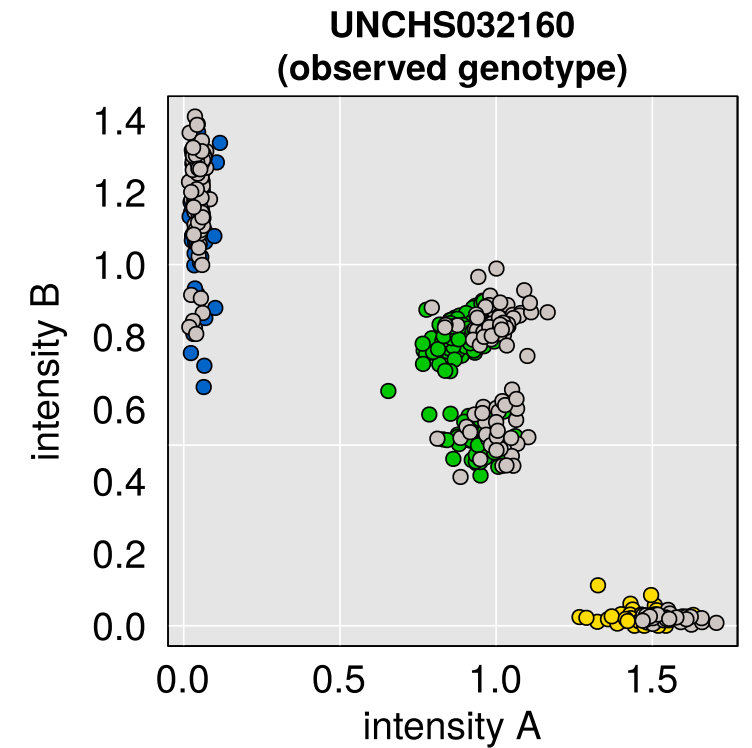
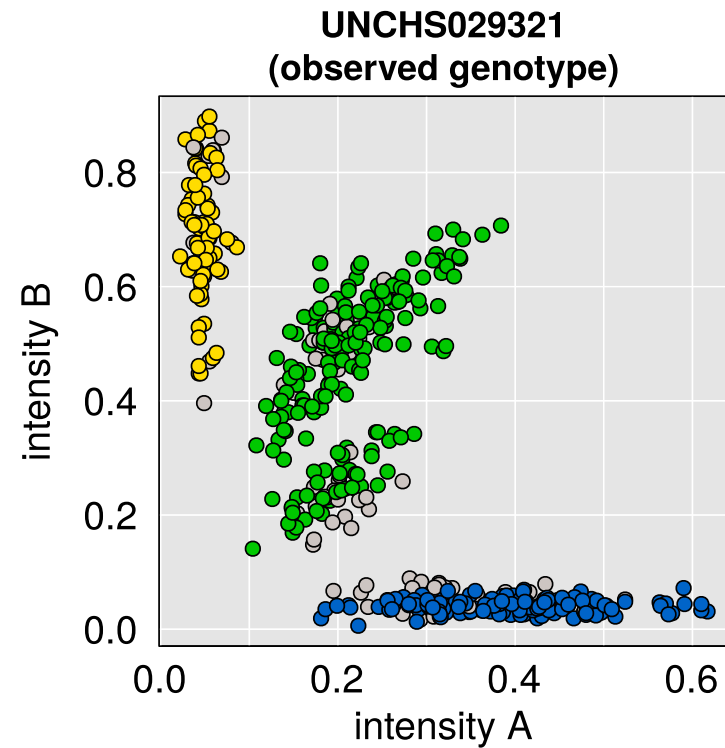
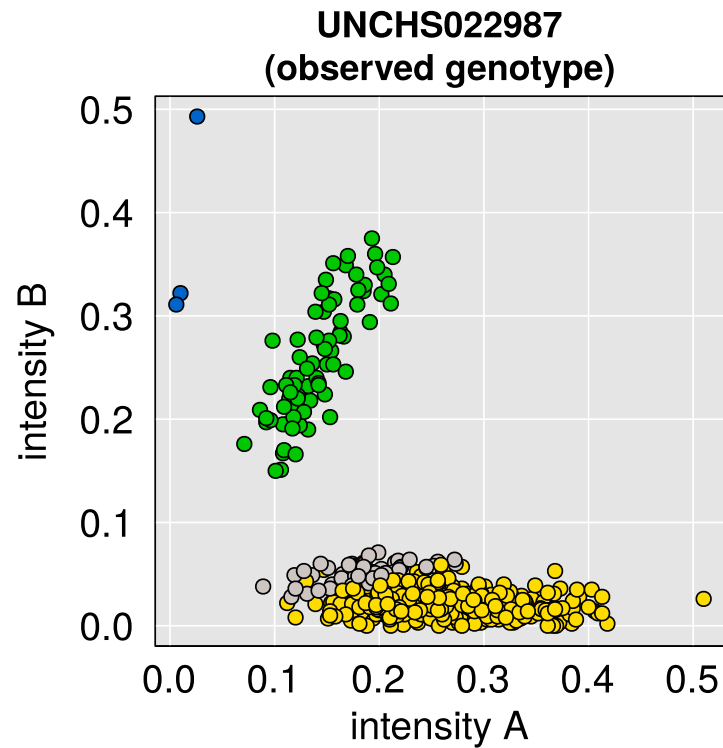
One bad blob



Wrong genomic coordinates

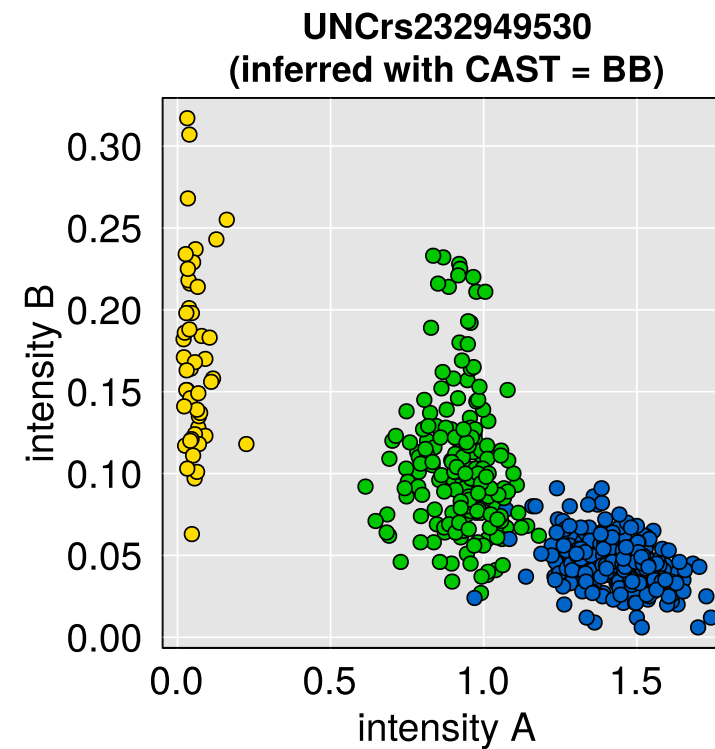
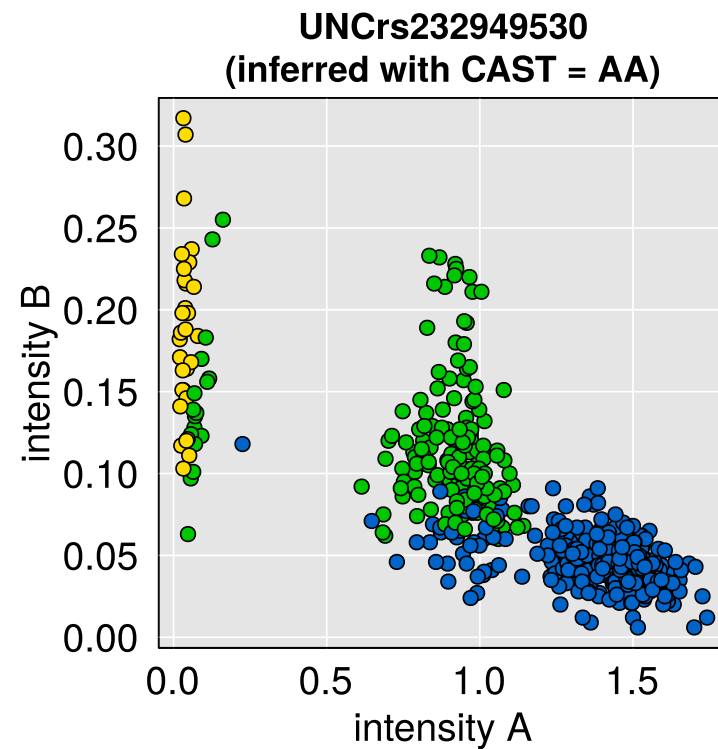
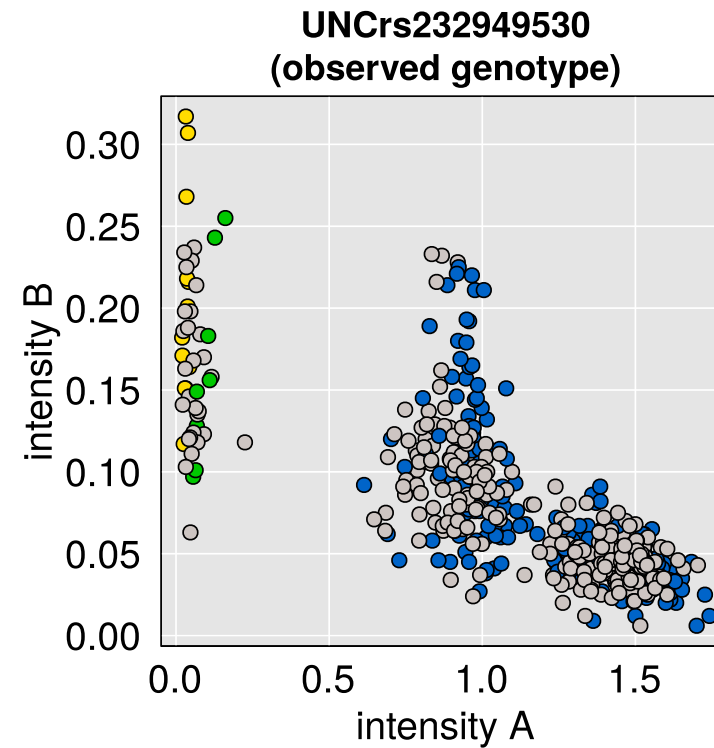


Puzzling no calls

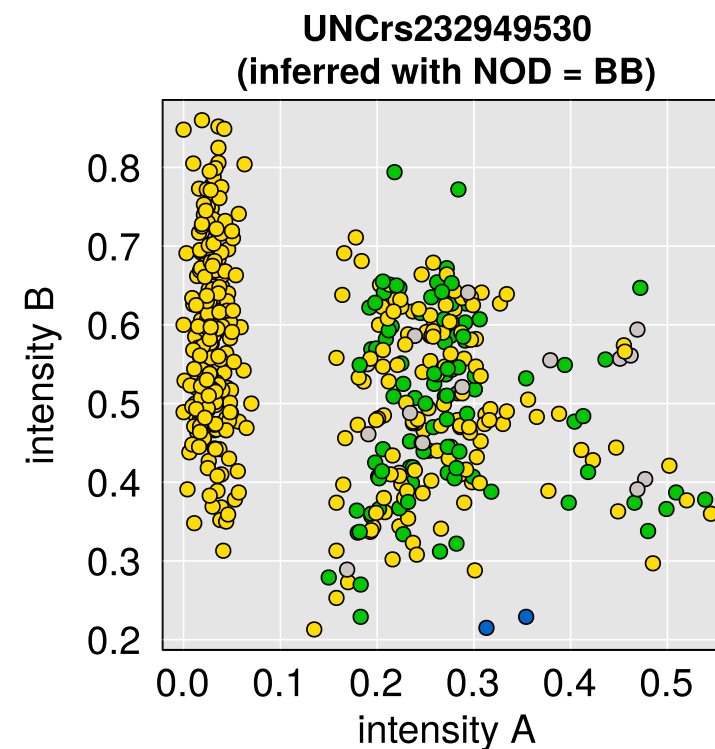
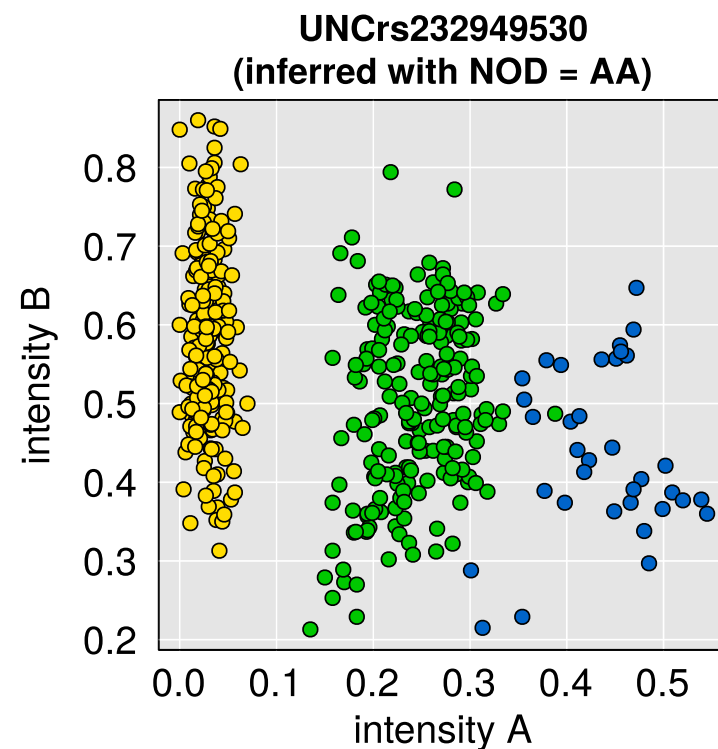
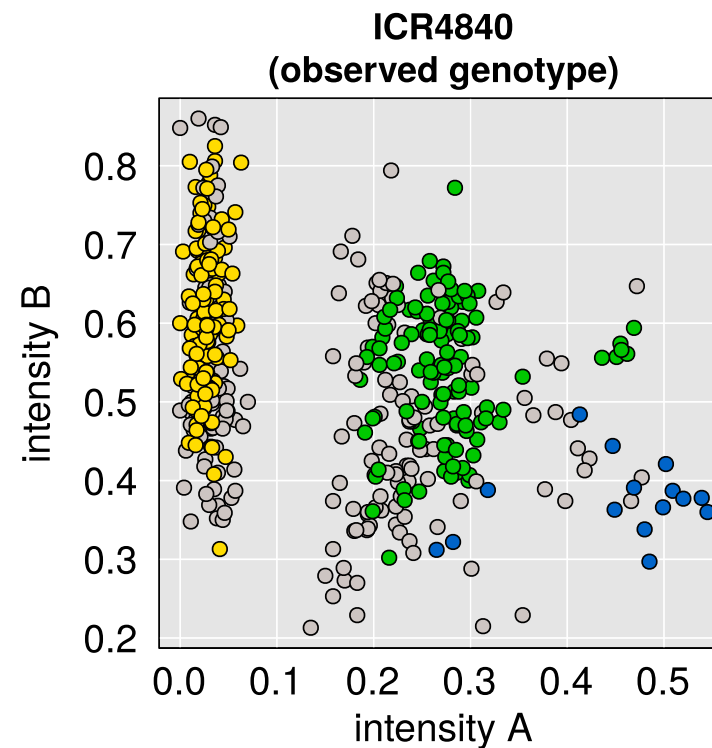


Founder genotyping errors

One founder missing



Another case with one founder missing



Summary

- Quality of results depends on quality of data
- Think about what might have gone wrong, and how it might be revealed
- Pulling out the bad samples is the most important thing
- Sex swaps: look at array intensities
- Look for sample duplicates, and if possible sample mix-ups
- Samples: missing data, array intensities, crossovers, errors
- Markers: lots of reasons for the bad ones

Acknowledgments

Alan Attie
Gary Churchill
Dan Gatti
Alexandra Lobo
Federico Rey
Śaunak Sen
Lindsay Traeger
Brian Yandell

NIH/NIGMS, NIH/NIDDK

Slides: bit.ly/jax18



kbroman.org

github.com/kbroman

[@kwbroman](#)