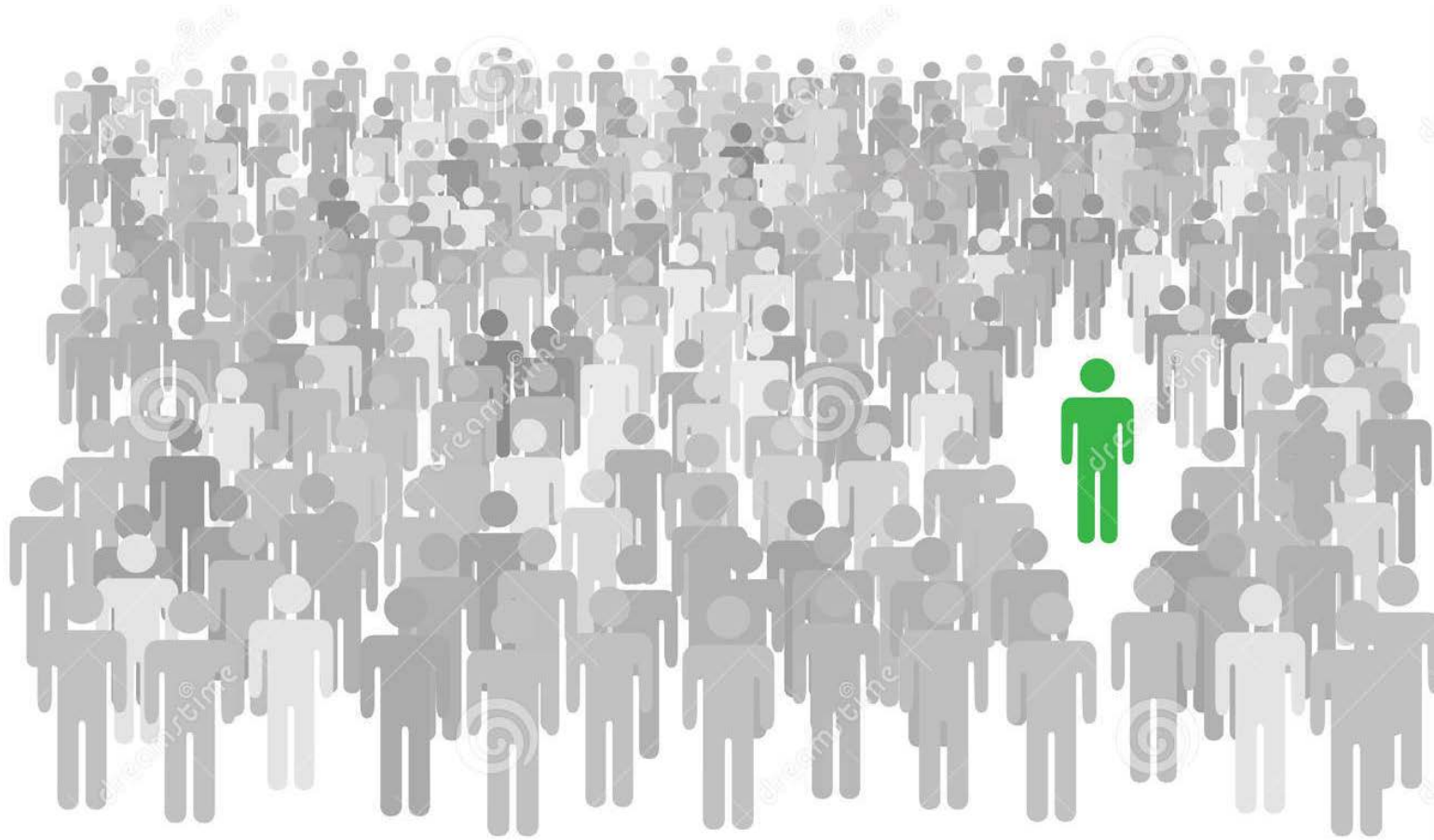


# Population genomics



John Wang  
April 29, 2020  
Intro NGS



# Population genetics:

The study of polymorphism (within species) and divergence (between species). [Hartl 2003]



Aims to understand evolution

Changes in frequency of alleles within populations over space and time

- natural selection
- genetic drift
- gene flow
- mutation
- recombination

# Population genomics:

The field of population genomics surveys patterns in the genome within and among populations to make inferences about evolution and the genome. (Nosil and Buerkle 2010)



# Population genomics:



The field of population genomics surveys patterns in the genome within and among populations to make inferences about evolution and the genome. (Nosil and Buerkle 2010)

Population genomics can be broadly defined as the **simultaneous study of numerous loci or genome regions** to better understand the roles of evolutionary processes (such as mutation, random genetic drift, gene flow and natural selection) that influence variation across genomes and populations. (Luikart et al 2003)

# Population genomics:



The field of population genomics surveys patterns in the genome within and among populations to make inferences about evolution and the genome. (Nosil and Buerkle 2010)

Population genomics can be broadly defined as the **simultaneous study of numerous loci or genome regions** to better understand the roles of evolutionary processes (such as mutation, random genetic drift, gene flow and natural selection) that influence variation across genomes and populations. (Luikart et al 2003)

Population genomics is the use of genome-wide sampling to identify and to **separate locus-specific effects** (such as selection, mutation, assortative mating and recombination) from **genome-wide effects** (such as drift or bottlenecks, gene flow and inbreeding) to improve our understanding of microevolution. (Black et al 2001)

# Population genomics:

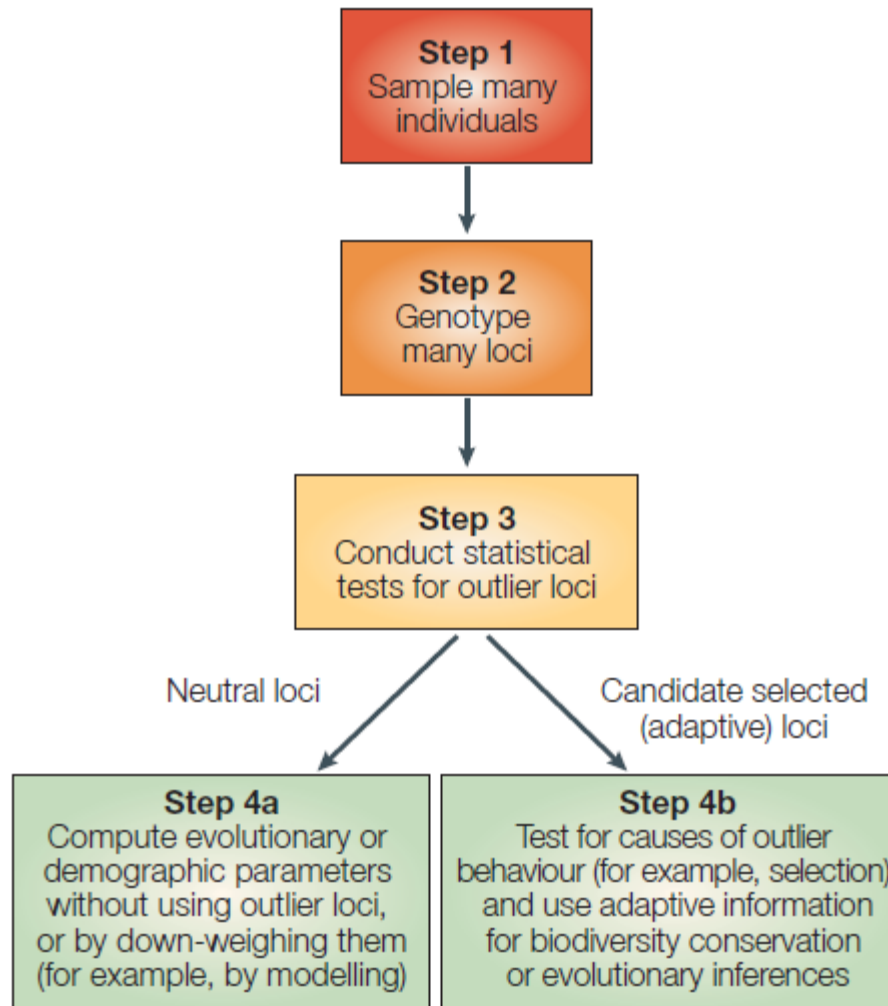


Figure 1 | **Flow chart of the four main steps in the population-genomic approach.** The approach summarized

# Population genomics:

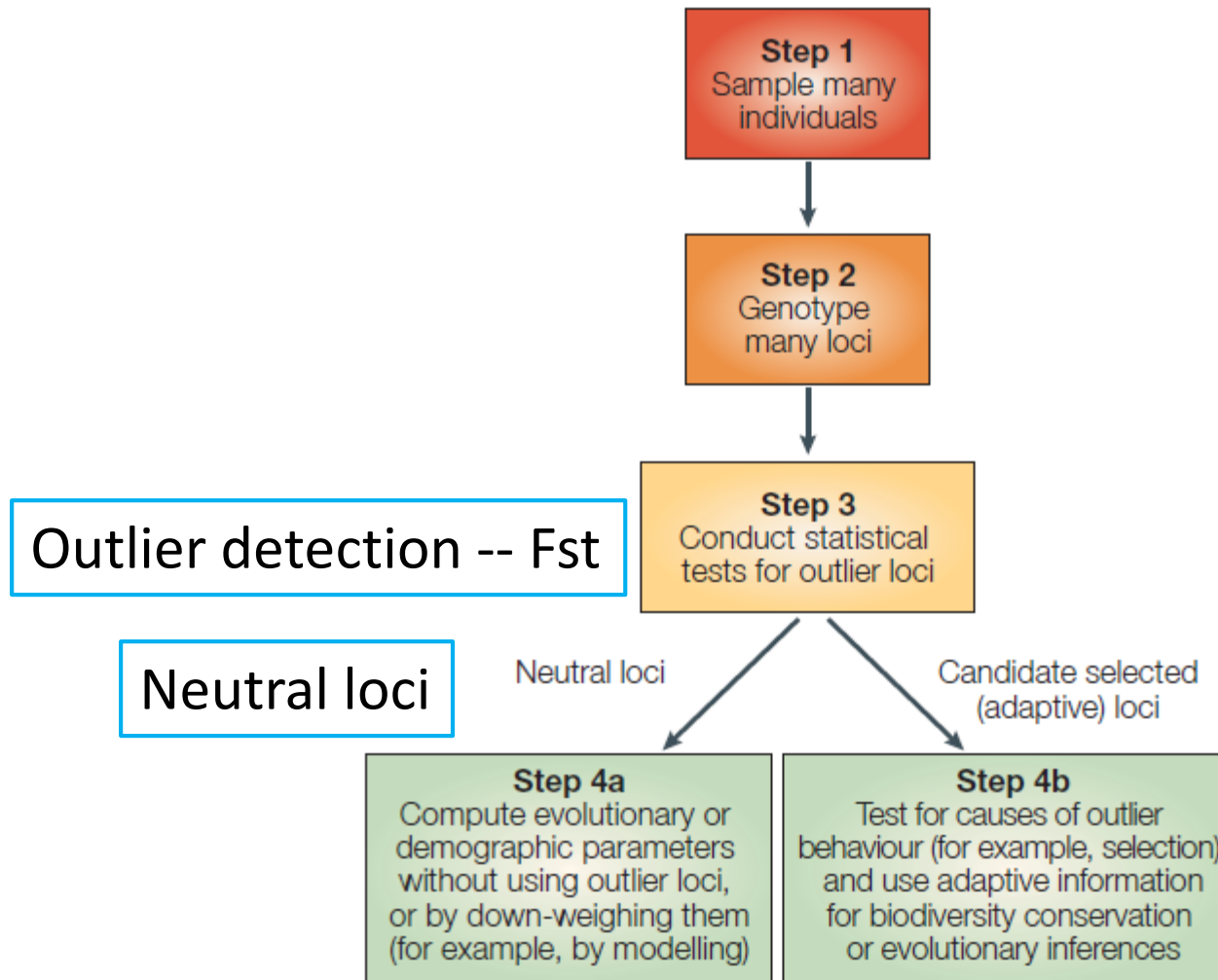


Figure 1 | **Flow chart of the four main steps in the population-genomic approach.** The approach summarized

# Population genomics:

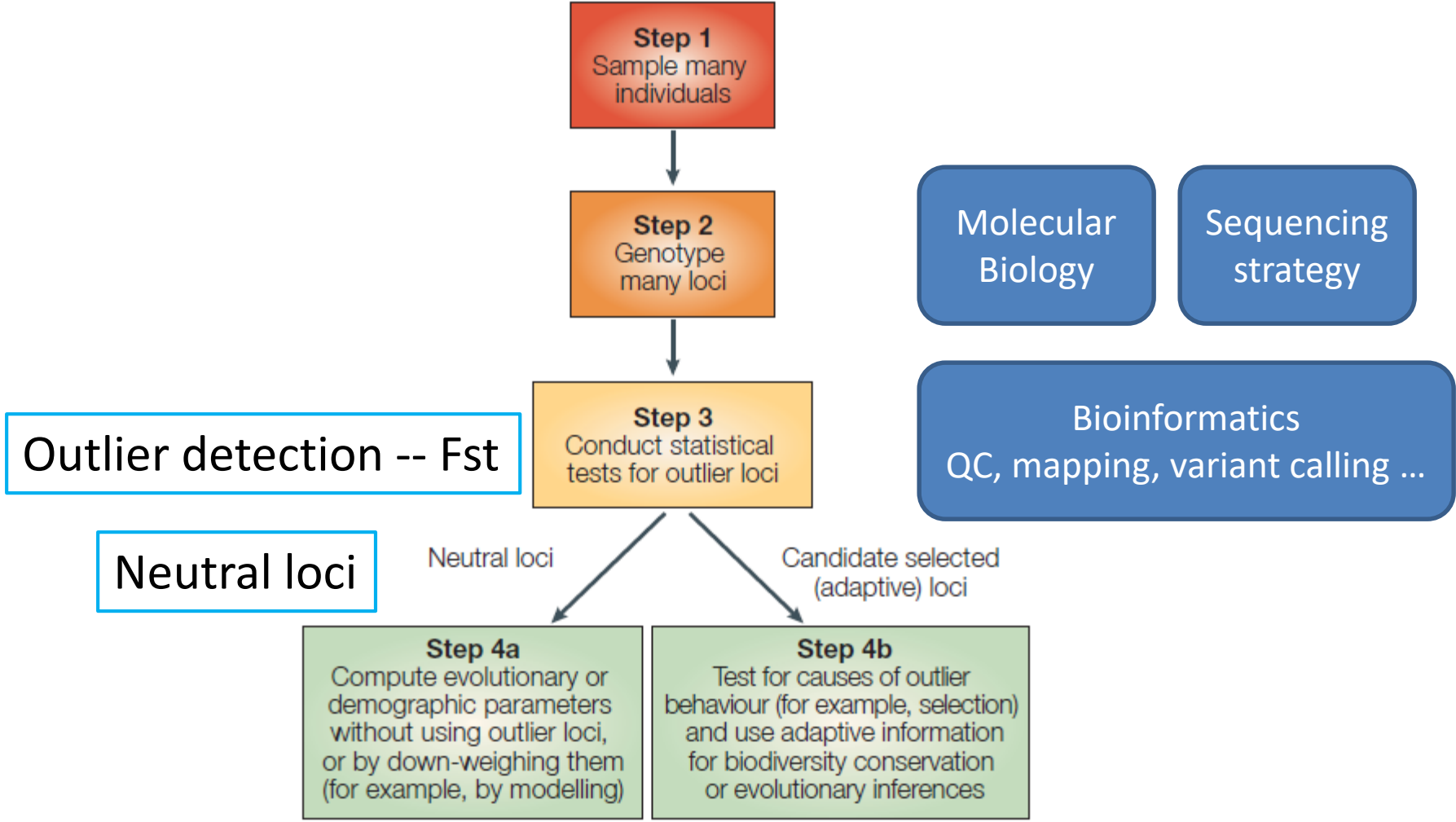
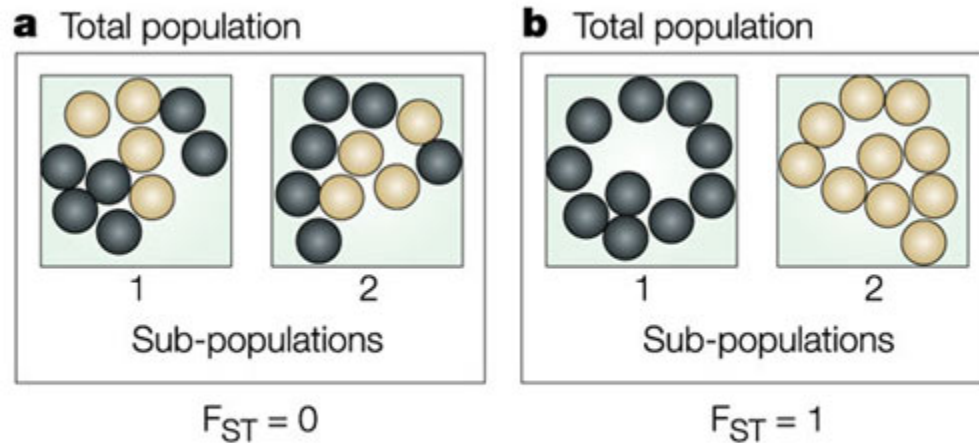


Figure 1 | **Flow chart of the four main steps in the population-genomic approach.** The approach summarized



# Fixation index – $F_{ST}$



Nature Reviews | **Genetics**

- Measure of population differentiation.
- Most widely used index of genetic divergence b/w populations.
- High  $F_{ST}$  → two populations are fixed for allelic differences.  
Drift or selection
- $F_{ST}$  scans across genome (locus-by-locus → e.g., manhattan plots)

# Some population genomics/genetics questions

# Some population genomics/genetics questions

What genes are under selection? [speciation genes, domestication]

# Some population genomics/genetics questions

What genes are under selection? [speciation genes, domestication]

Has selection been acting or can drift be sufficient explanation? [allele frequency clines]

# Some population genomics/genetics questions

What genes are under selection? [speciation genes, domestication]

Has selection been acting or can drift be sufficient explanation? [allele frequency clines]

What gene contributes to a trait? [agriculturally important genes]

# Some population genomics/genetics questions

What genes are under selection? [speciation genes, domestication]

Has selection been acting or can drift be sufficient explanation? [allele frequency clines]

What gene contributes to a trait? [agriculturally important genes]

What was past demographic history? [out of Africa]

# Some population genomics/genetics questions

What genes are under selection? [speciation genes, domestication]

Has selection been acting or can drift be sufficient explanation? [allele frequency clines]

What gene contributes to a trait? [agriculturally important genes]

What was past demographic history? [out of Africa]

Is this one big population or many small ones? [climate change]

# Some population genomics/genetics questions

What genes are under selection? [speciation genes, domestication]

Has selection been acting or can drift be sufficient explanation? [allele frequency clines]

What gene contributes to a trait? [agriculturally important genes]

What was past demographic history? [out of Africa]

Is this one big population or many small ones? [climate change]

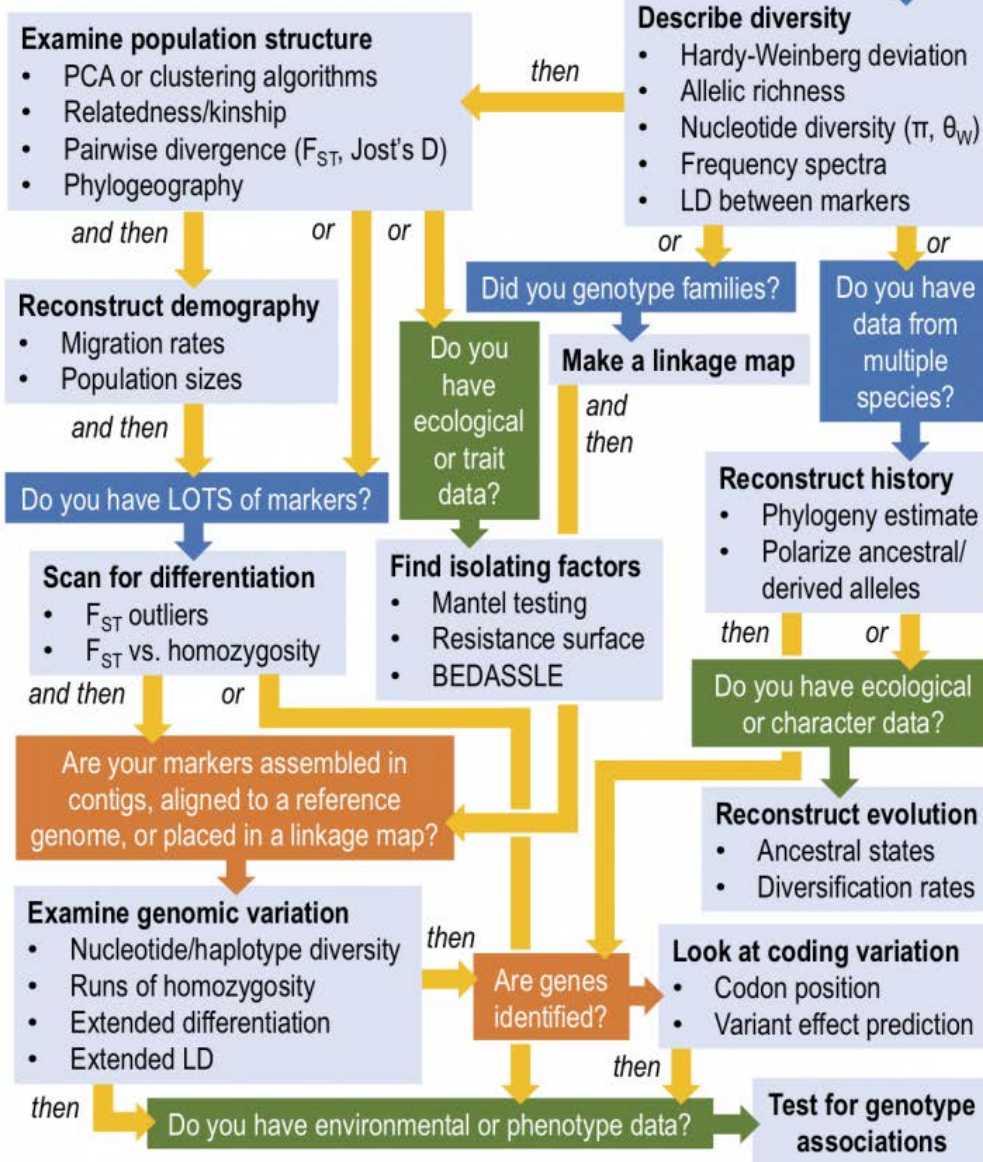
“Just so” stories of selection are easy to invoke...

... but there might be alternative scenarios, e.g. genetic surfing



# What can I do with this genetic data?

If you have allozymes, AFLPs, sequences, microsatellites, or SNPs, you can



# Approaches

- Reduced representation sequencing (various)
- Low coverage sequencing → Pool-seq
- Whole genome re-sequencing (WGS)

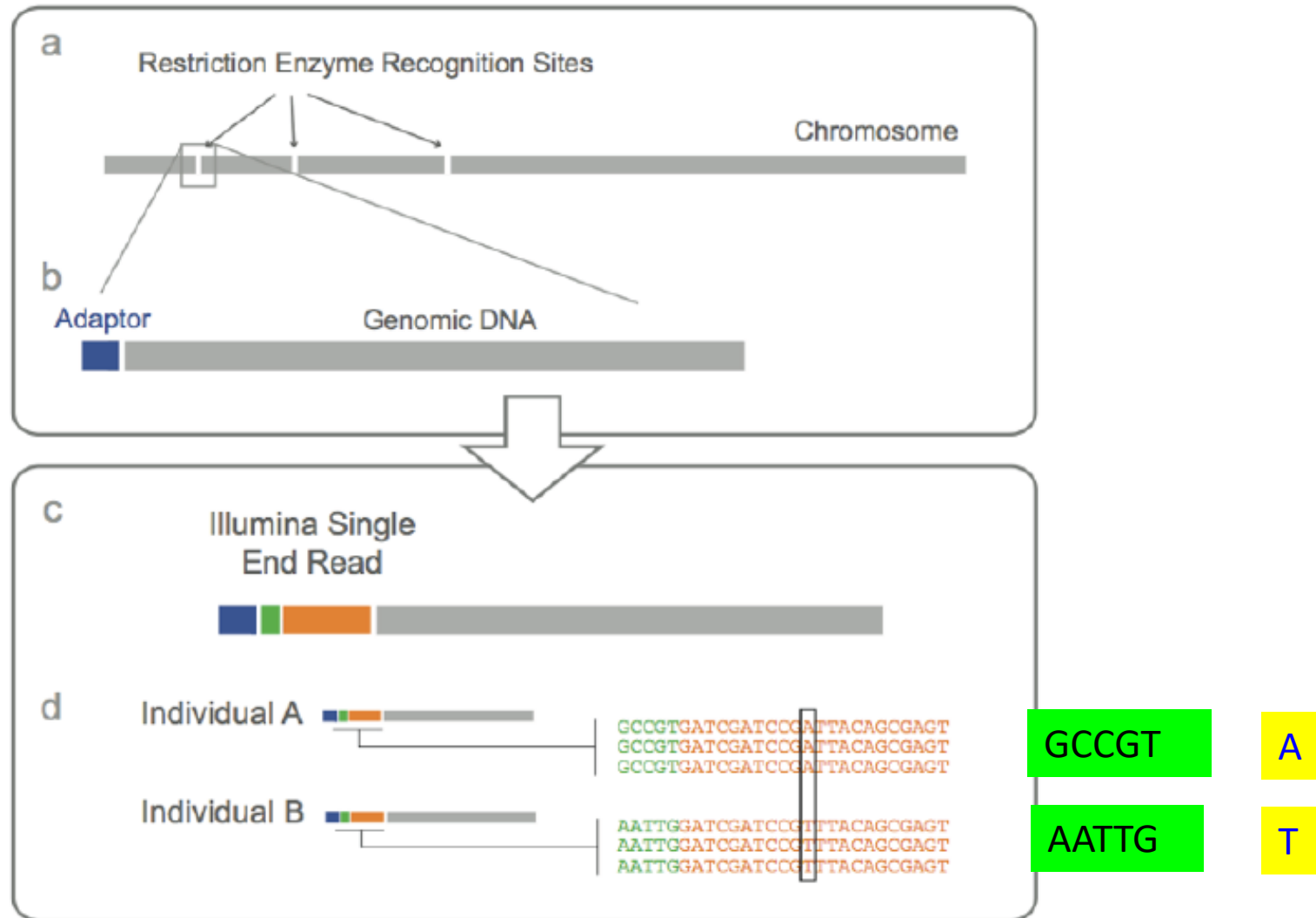
# Reduced representation sequencing approaches

- RAD-seq [genotyping-by-synthesis]
- Exon capture (targeted sequencing)
- Transcriptome (RNA-seq)
- Restriction enzyme size selection

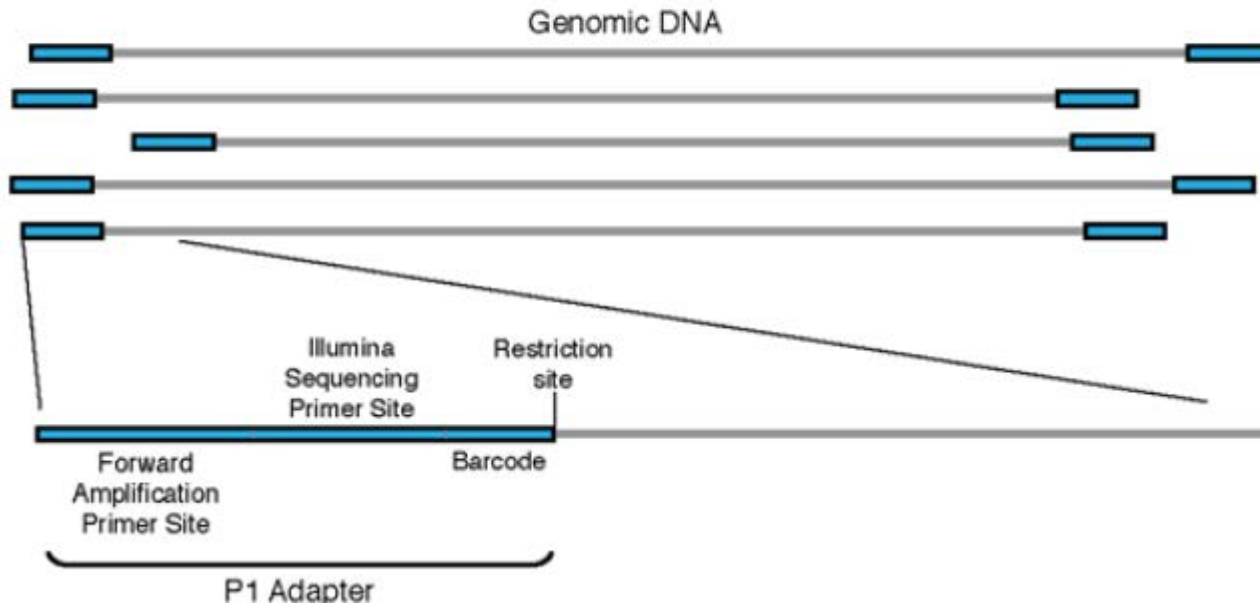
# Genotyping-by-sequencing

e.g., RADseq

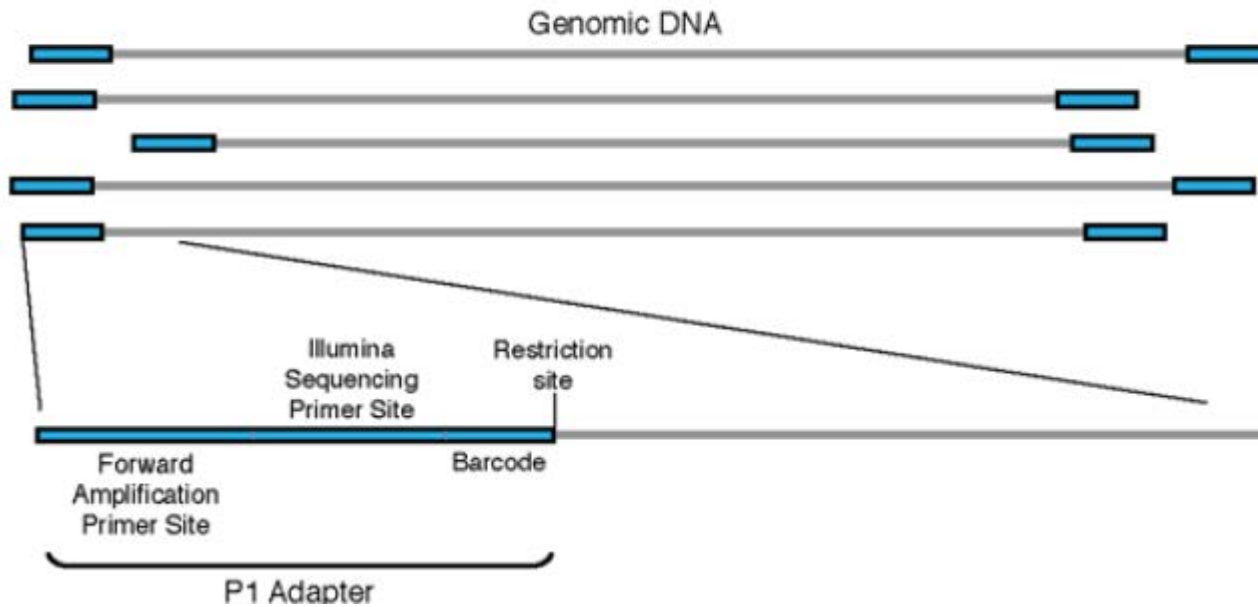
# Restriction site associated DNA (RAD) marker sequencing with barcoding



**A** *Ligate P1 Adapter to digested genomic DNA*



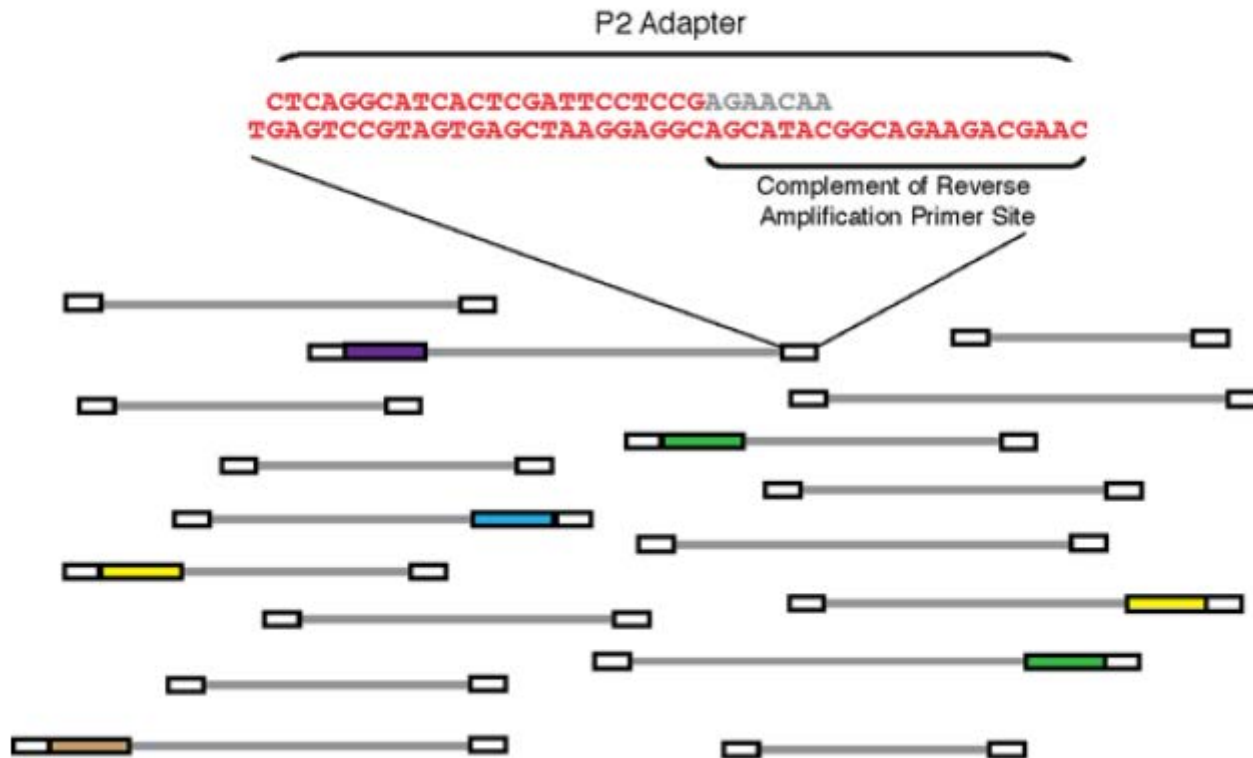
**A** *Ligate P1 Adapter to digested genomic DNA*



**B** *Pool barcoded samples and shear*



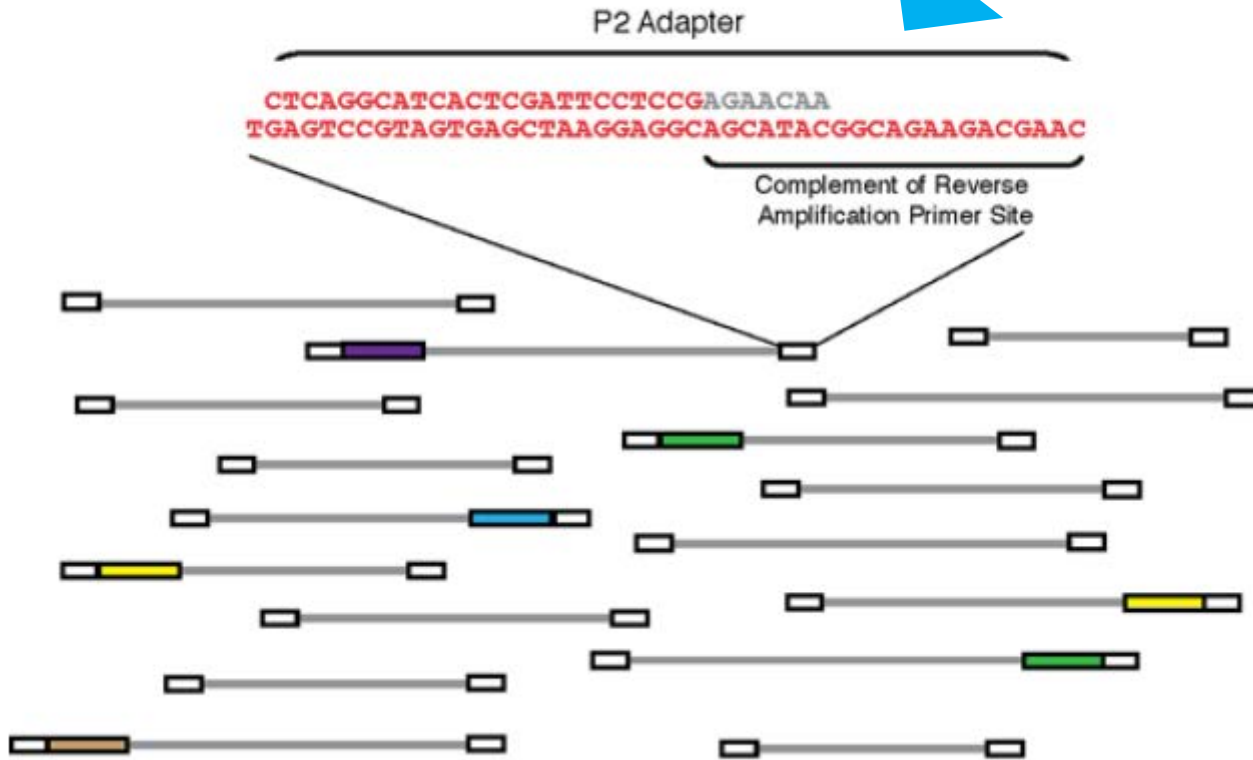
**C** *Ligate P2 Adapter to sheared fragments*



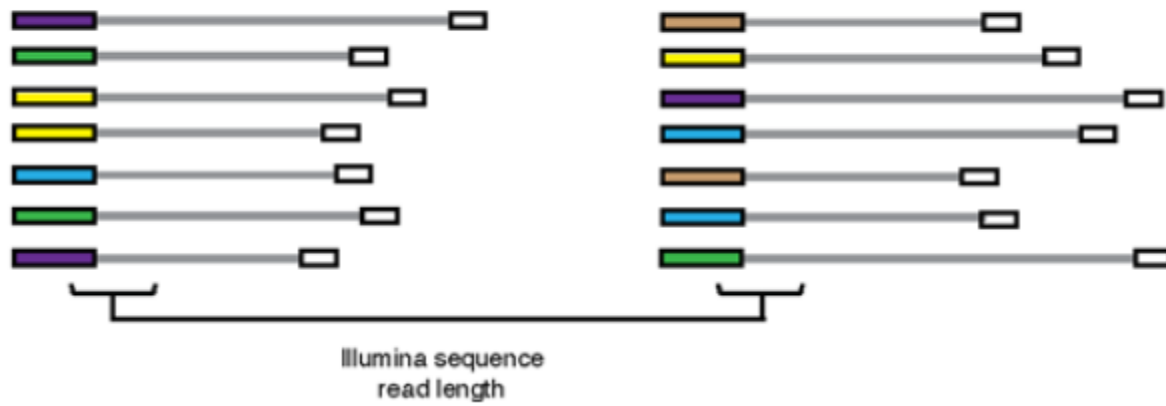


**C** *Ligate P2 Adapter to sheared fragments*

This is really clever:  
Only PCR's what you want



**D** *Selectively amplify RAD tags*



# RAD markers sequence a subset of the genome

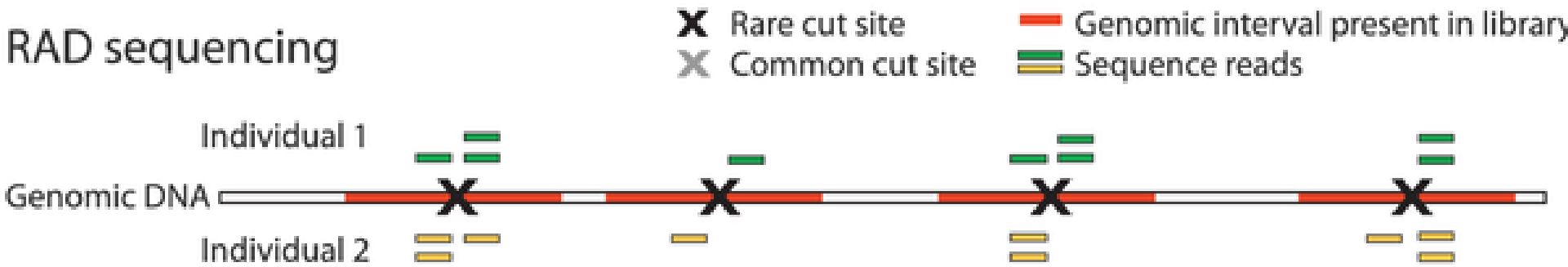


# There are several variants of RAD-seq

ddRAD  
ezRAD  
2bRAD  
...

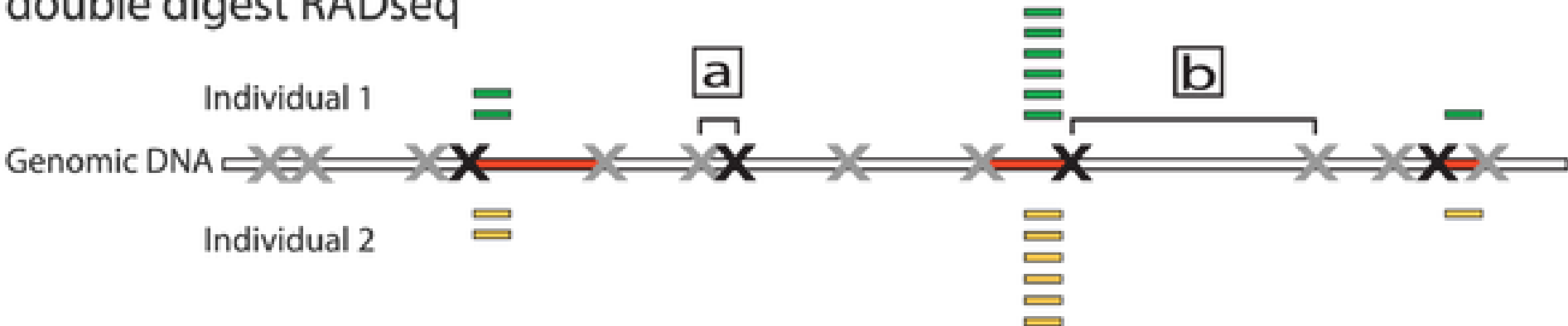
A

## RAD sequencing



B

## double digest RADseq



# ddRAD-seq

Easier tunability; better reproducibility

But PCR duplicates unknowable

Fraction of genome

Sanger sequencing

Whole genome re-sequencing

ddRAD

RADtag (Baird 2008)

Phylogeny

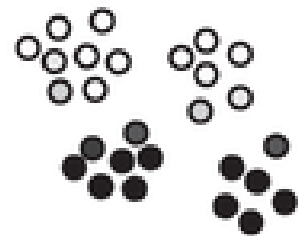
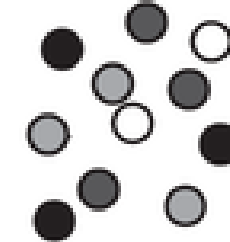
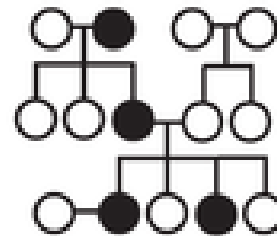
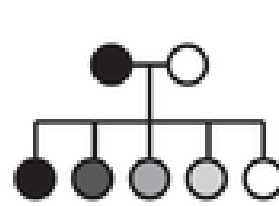
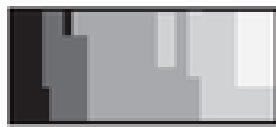
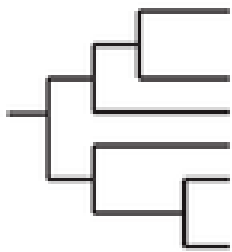
Population Structure

QTL Mapping

Pedigree Mapping

Association Mapping

Population Genomic Scans



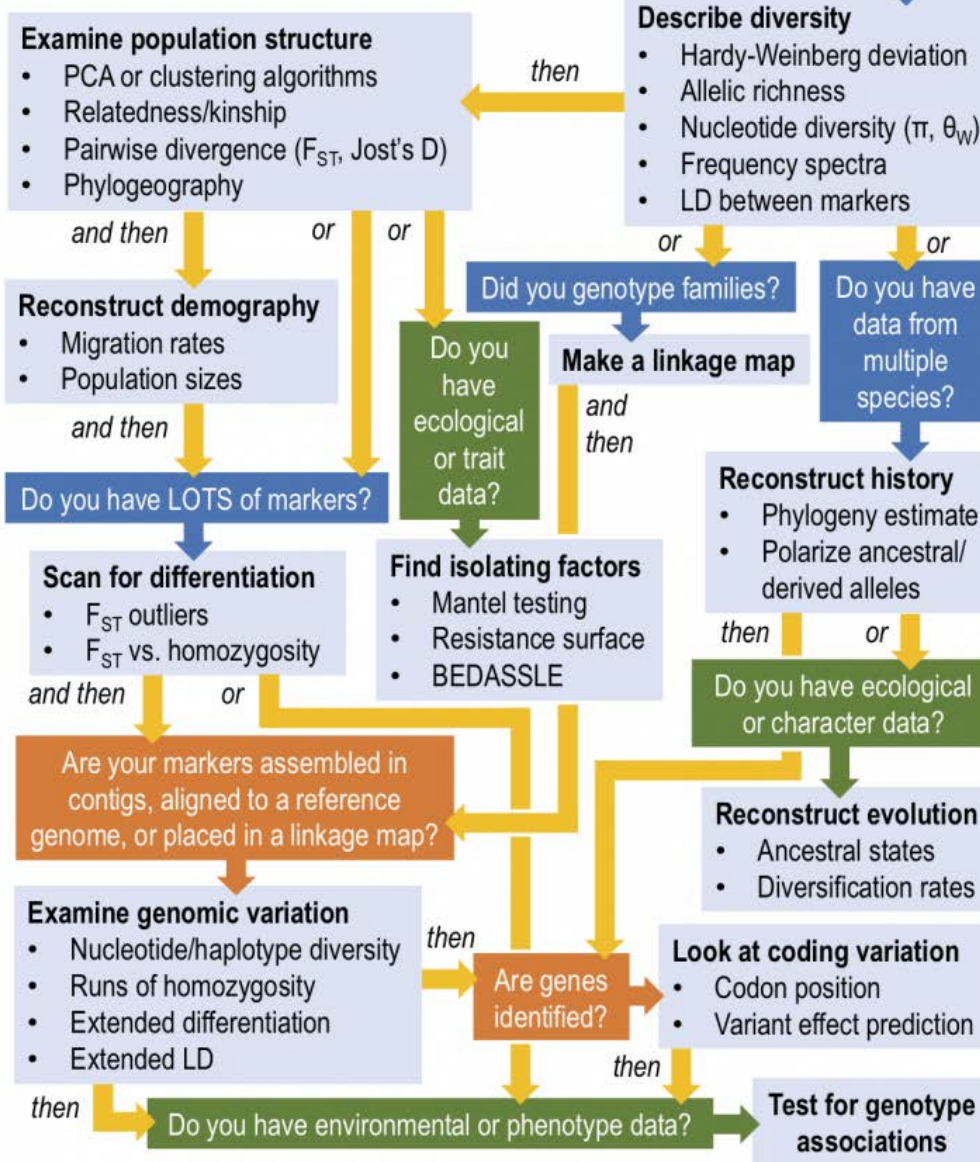
Divergence limited

Recombination limited

Linkage Diseq. limited

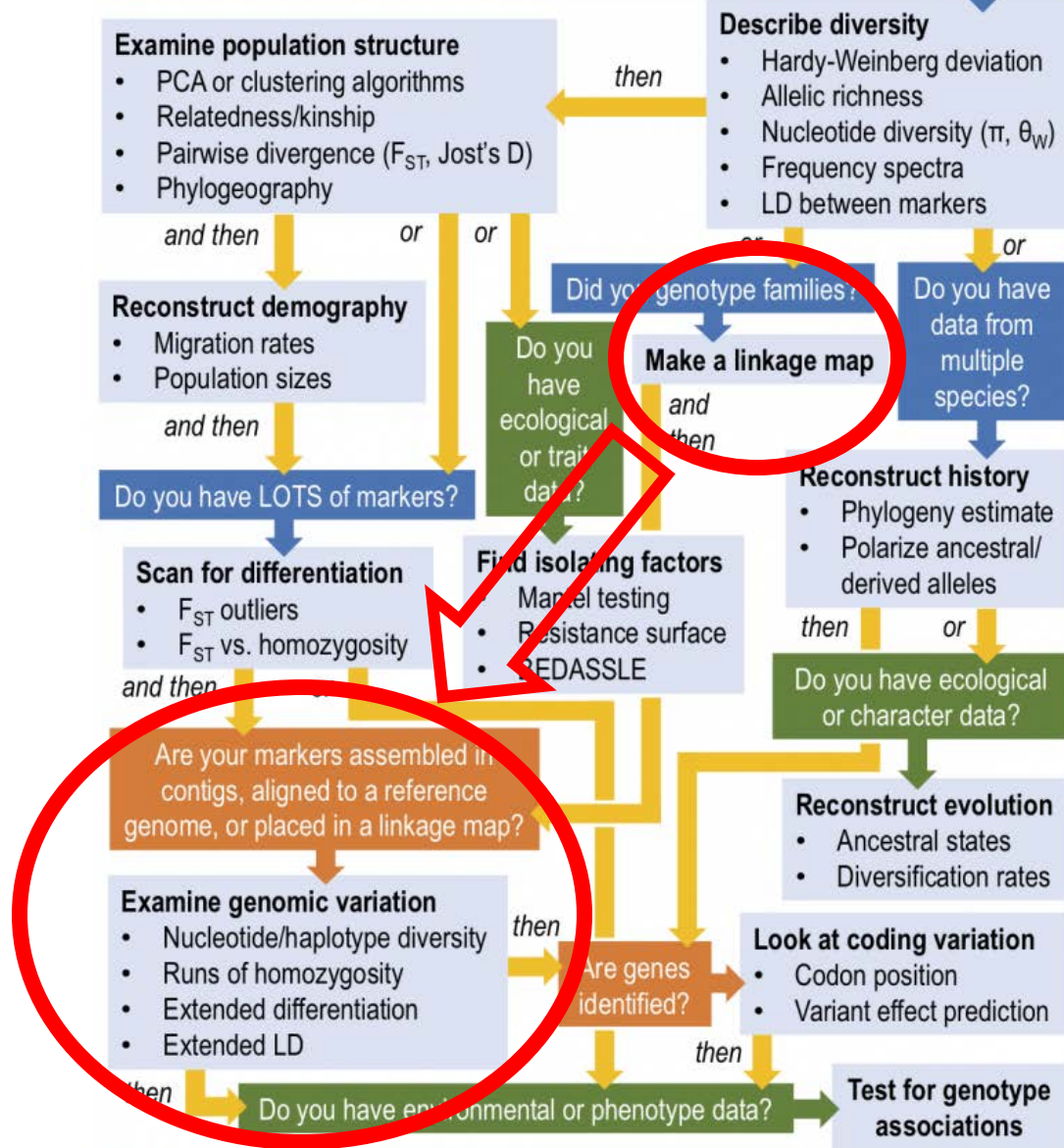
# What can I do with this genetic data?

If you have allozymes, AFLPs, sequences, microsatellites, or SNPs, you can



# What can I do with this genetic data?

If you have allozymes, AFLPs, sequences, microsatellites, or SNPs, you can





# Our RADseq example



Monogyne



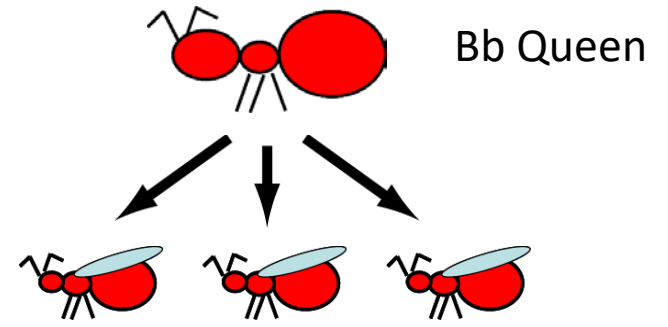
Polygyne

Fire ant *Solenopsis invicta*

Goals:

- 1) Build genetic map
- 2) test for the presence of a supergene



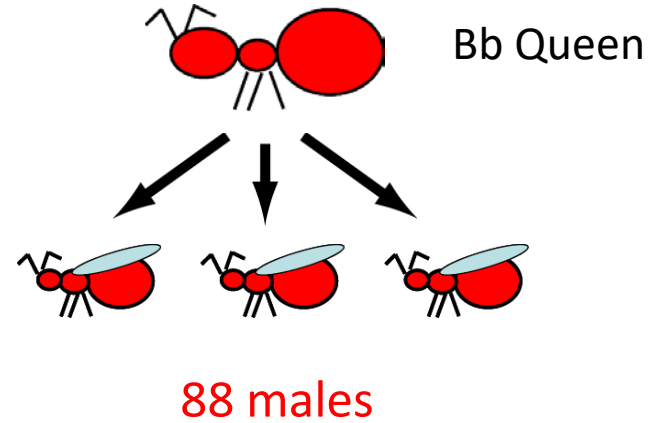


88 males

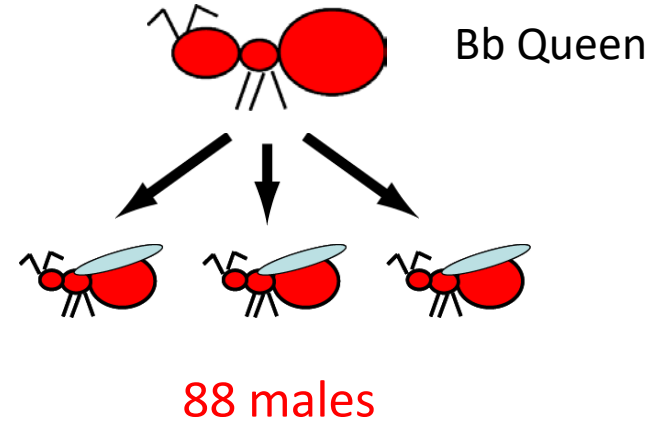
Isolate DNA, barcode each  
Sequence RAD tags with Illumina



192,236 RAD tags in genome  
(96,118 EcoRI sites)



Isolate DNA, barcode each  
Sequence RAD tags with Illumina



192,236 RAD tags in genome  
(96,118 EcoRI sites)

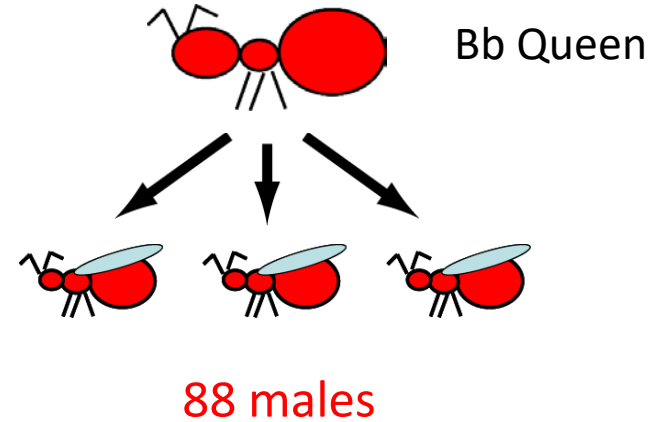
Isolate DNA, barcode each  
Sequence RAD tags with Illumina

Filter

12,684 RAD tags with multiple alleles

4,232 w/ unambiguous genotype info

2,724 w/ <25% missing data



192,236 RAD tags in genome  
(96,118 EcoRI sites)

Isolate DNA, barcode each  
Sequence RAD tags with Illumina

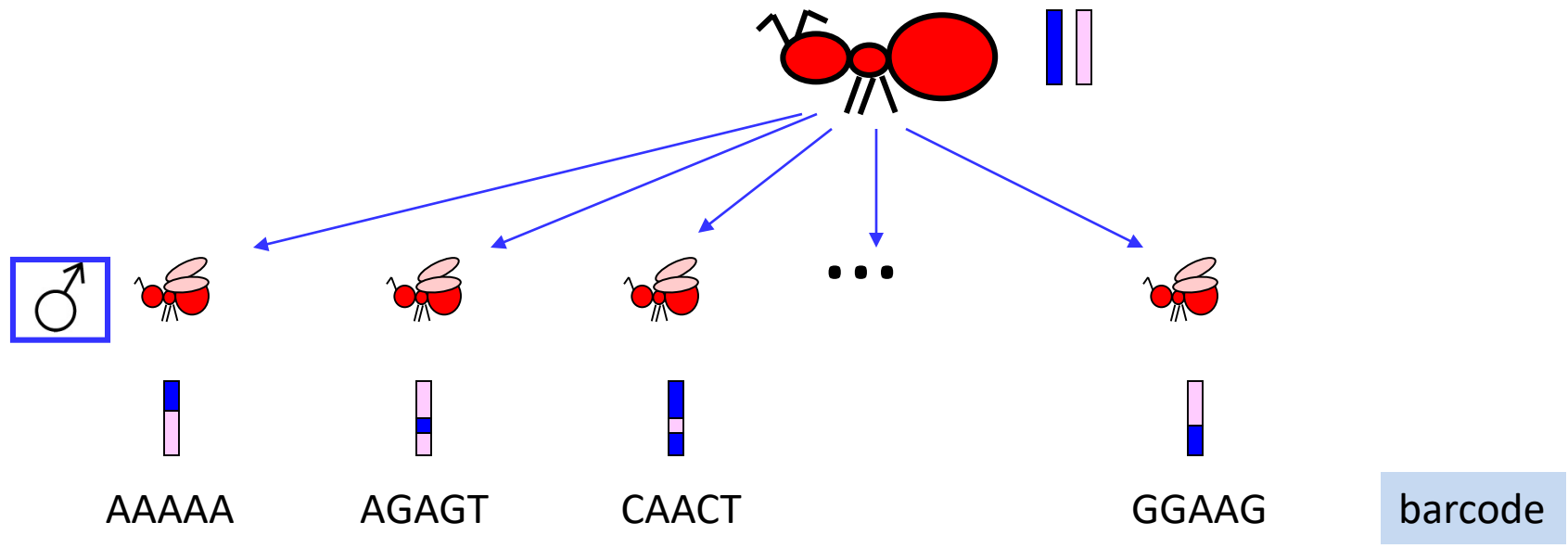
Filter

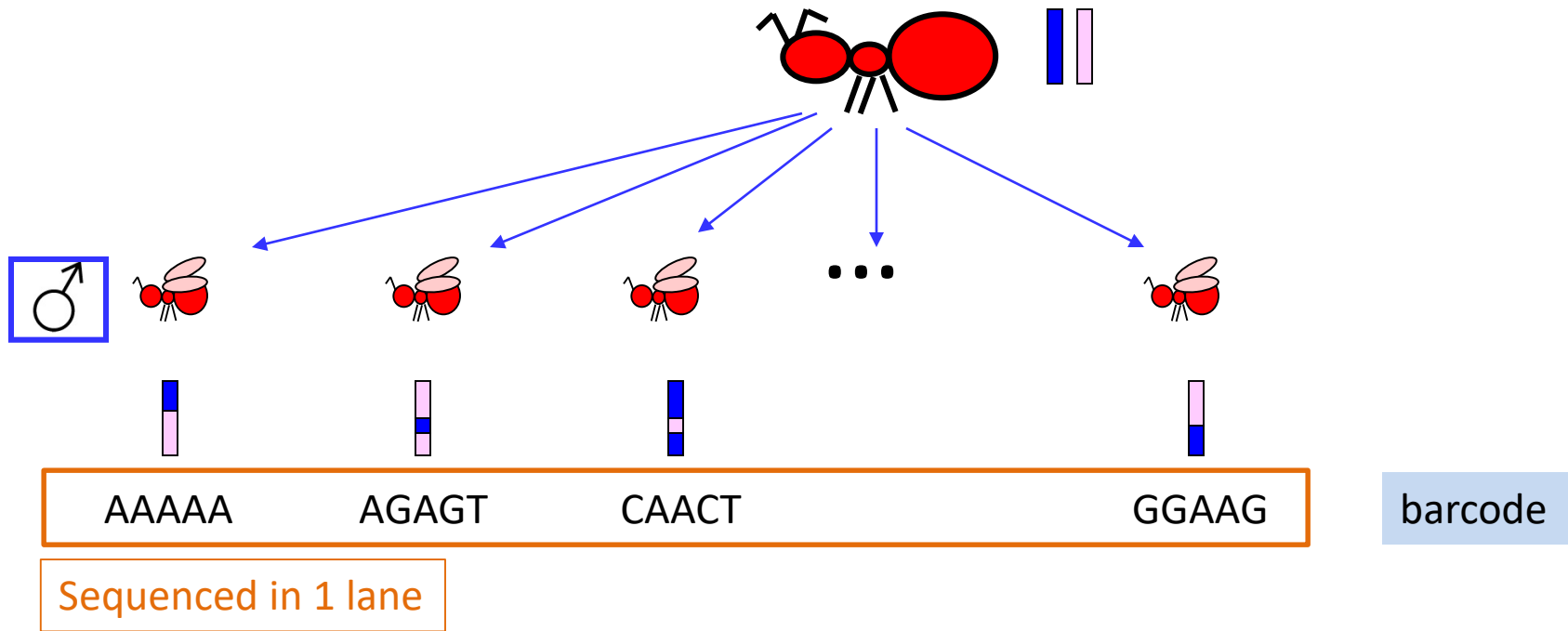
12,684 RAD tags with multiple alleles

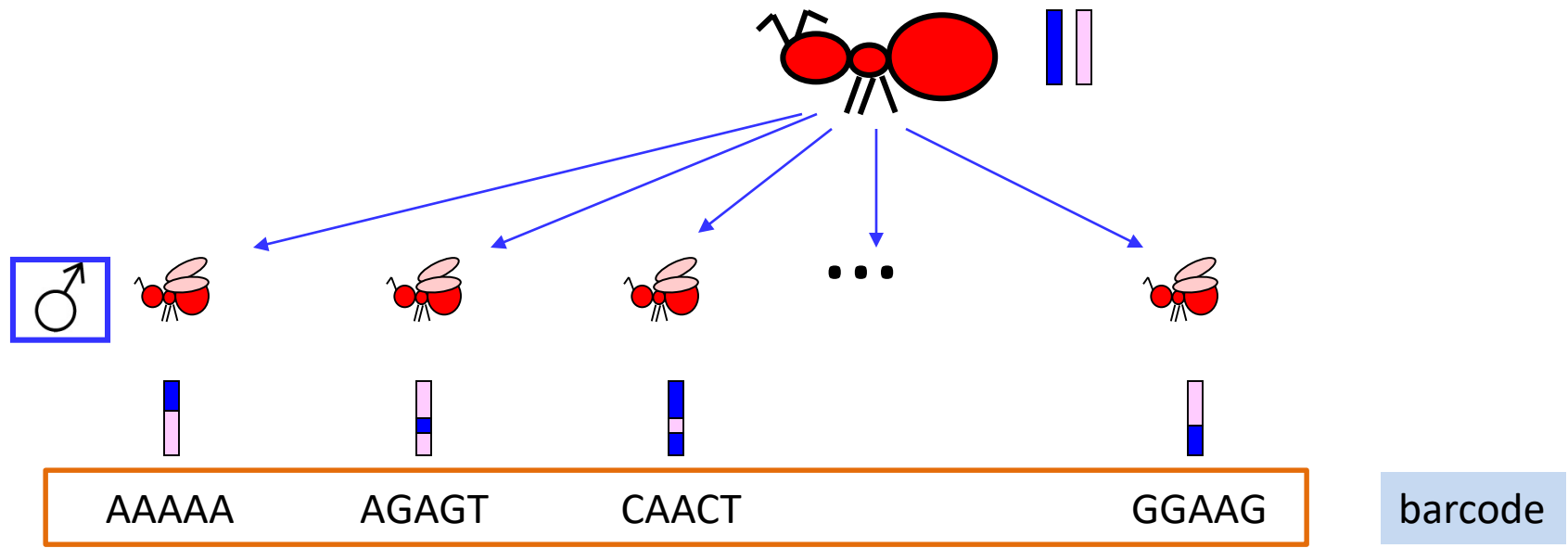
4,232 w/ unambiguous genotype info

2,724 w/ <25% missing data

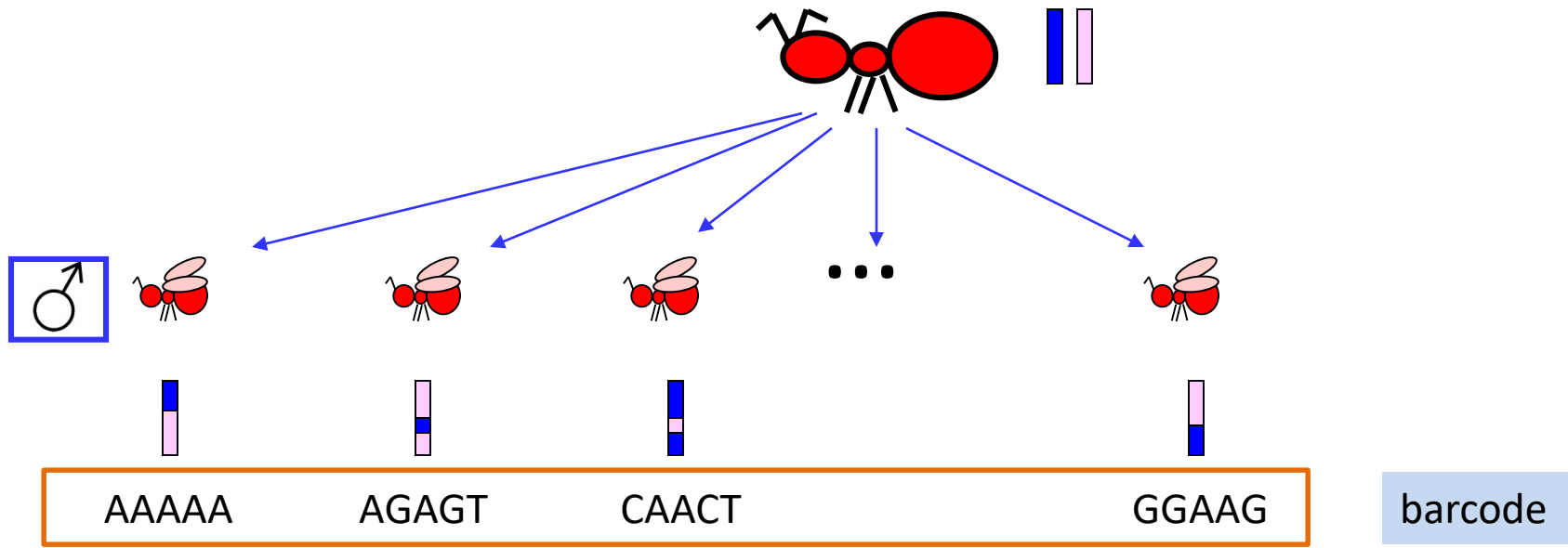
16 linkage groups  
(+1 locus outlier)







- 1) Split by barcode
- 2) Trim off barcode
- 3) Quality trimmer
- 4) Chop to 50 bp
- 5) Merge identical
- 6) Make fake FASTQ file
- 7) Run MAQ or bwa

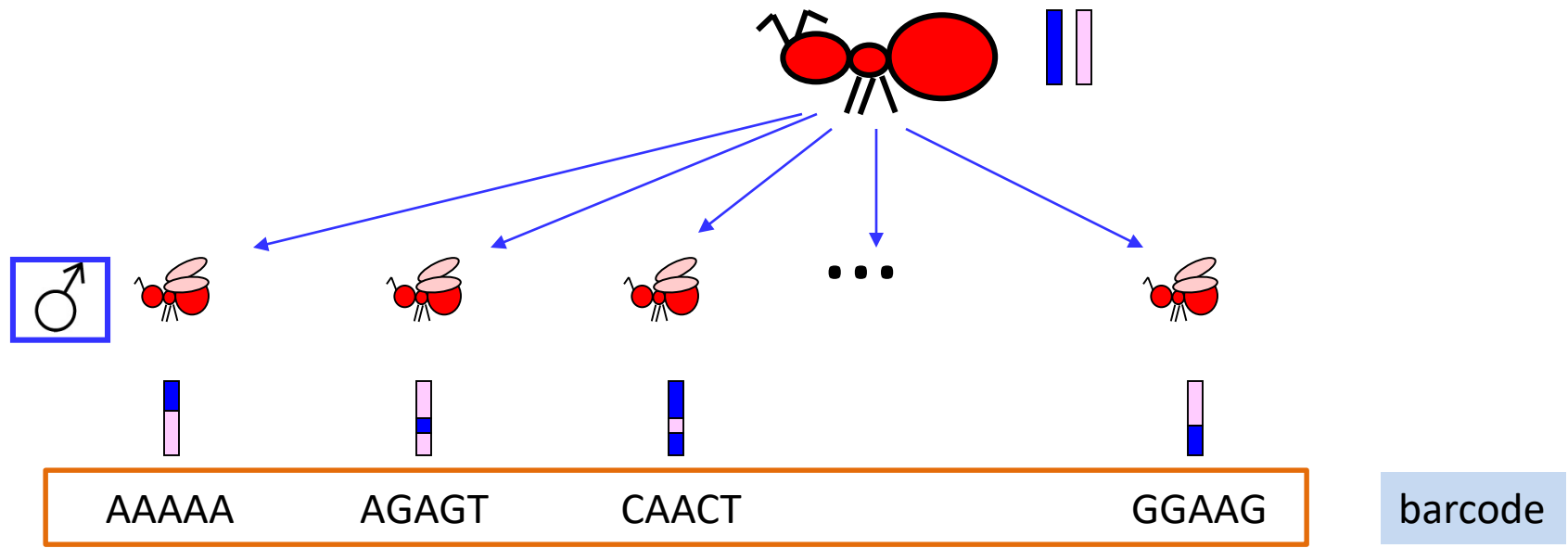


Sequenced in 1 lane

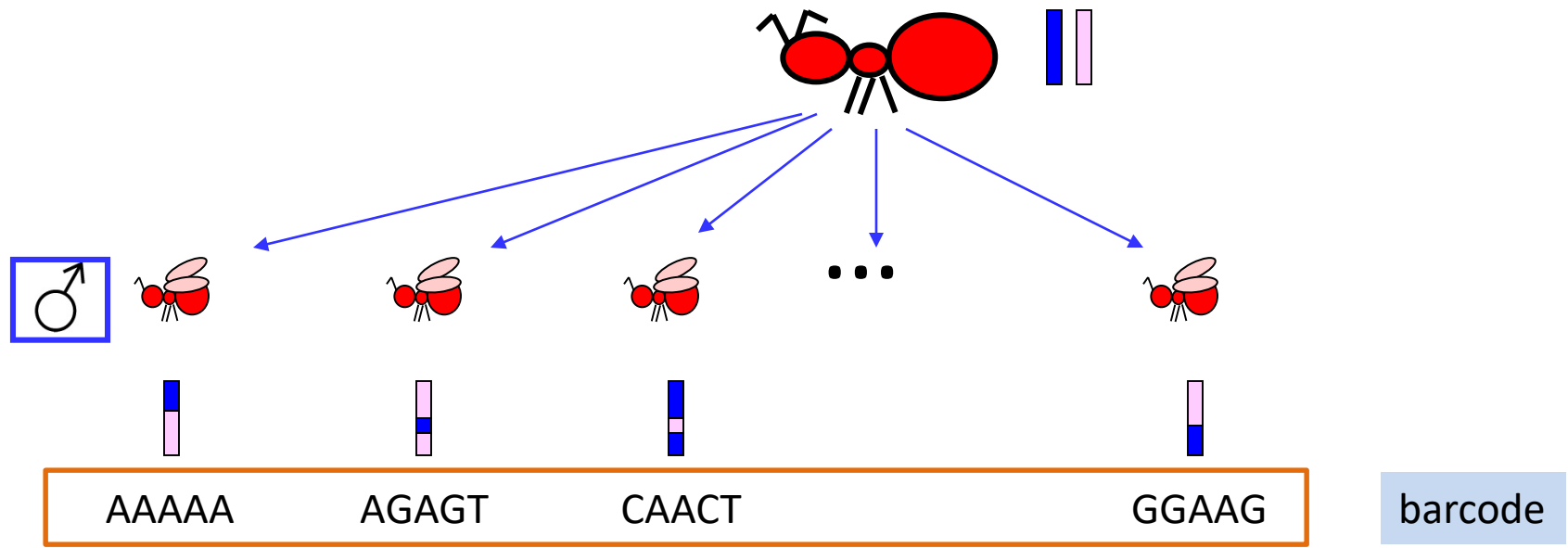


- 1) Split by barcode      GGAAGaattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 2) Trim off barcode
- 3) Quality trimmer
- 4) Chop to 50 bp
- 5) Merge identical
- 6) Make fake FASTQ file
- 7) Run MAQ or bwa





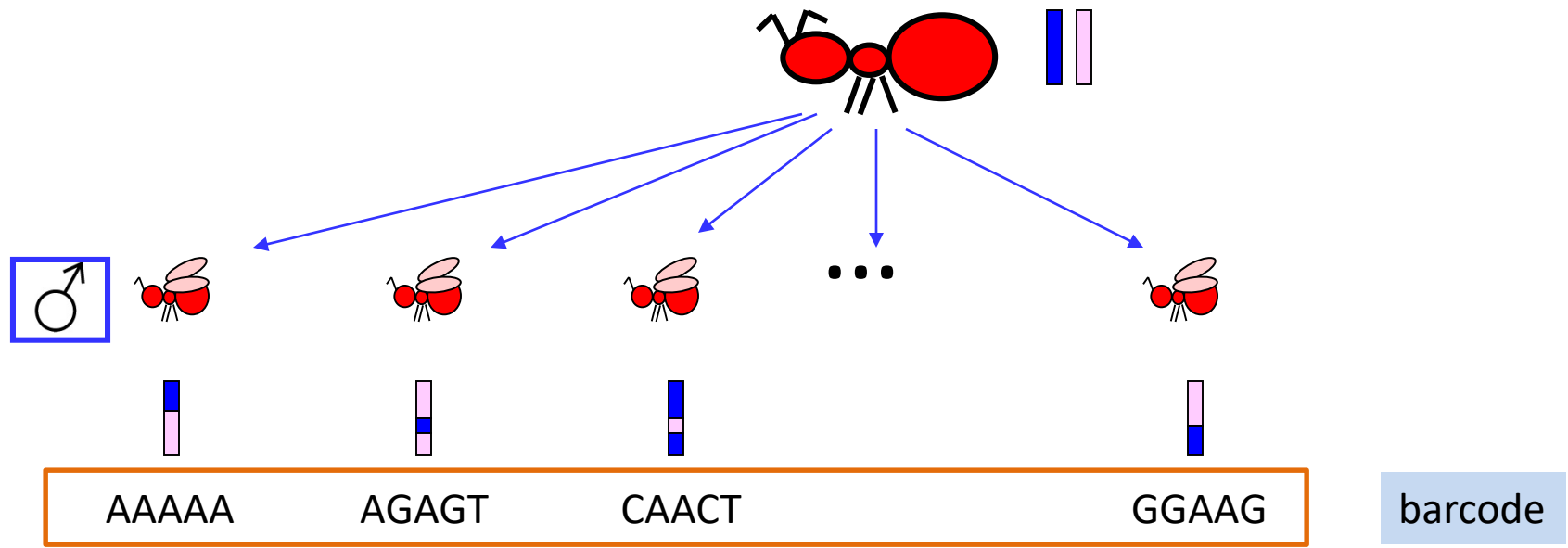
- 1) Split by barcode      GGAAGaattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 2) Trim off barcode      \_\_\_\_\_aattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 3) Quality trimmer
- 4) Chop to 50 bp
- 5) Merge identical
- 6) Make fake FASTQ file
- 7) Run MAQ or bwa



Sequenced in 1 lane



- 1) Split by barcode      GGAAGaattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 2) Trim off barcode      \_\_\_\_\_aattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 3) Quality trimmer      \_\_\_\_\_aattcacgtac ... gtacgt ... acgtacgt(63)
- 4) Chop to 50 bp
- 5) Merge identical
- 6) Make fake FASTQ file
- 7) Run MAQ or bwa



Sequenced in 1 lane

barcode

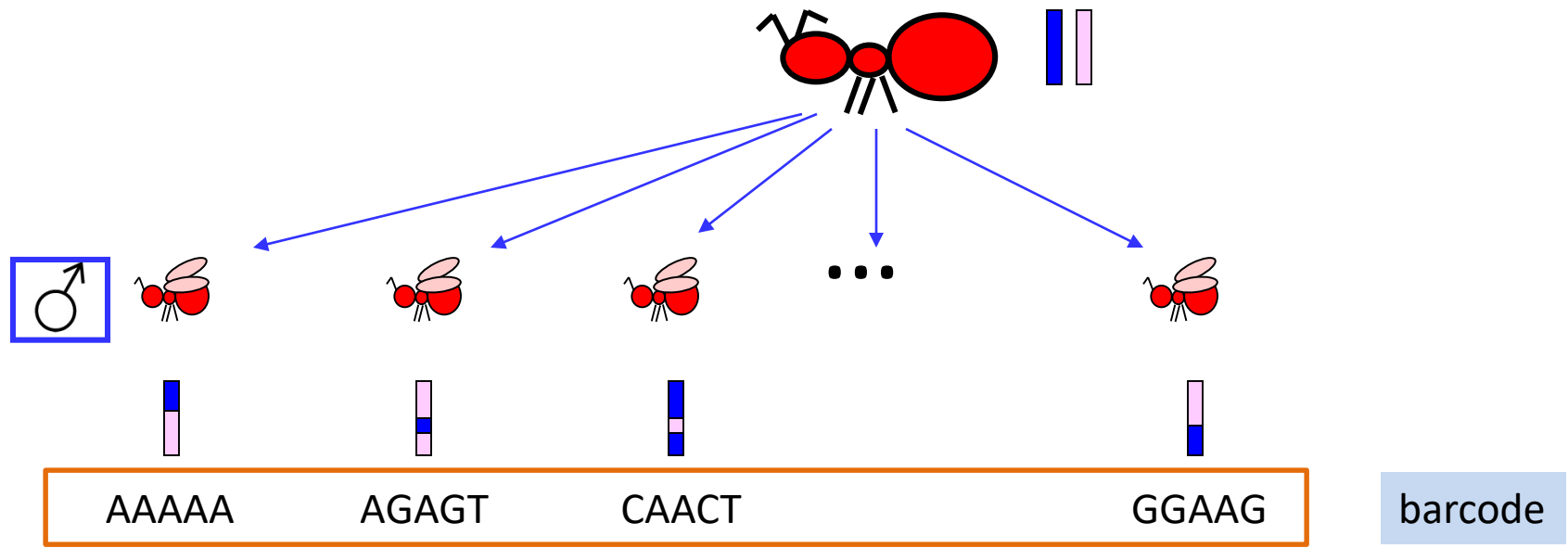


- 1) Split by barcode
- 2) Trim off barcode
- 3) Quality trimmer
- 4) Chop to 50 bp
- 5) Merge identical
- 6) Make fake FASTQ file
- 7) Run MAQ or bwa

```

GGAAGaattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
_____aattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
_____aattcacgtac ... gtacgt ... acgtacgt(63)
_____aattcacgtac ... gtacgt(50)

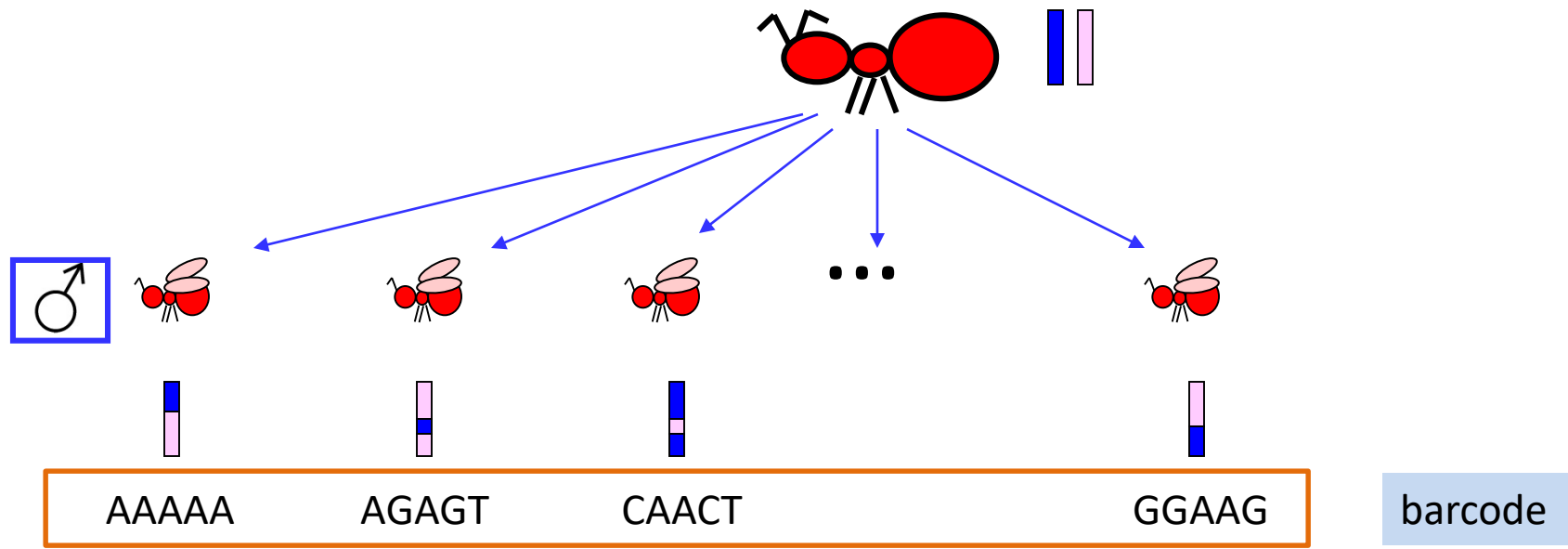
```



Sequenced in 1 lane



- 1) Split by barcode      GGAAGaattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 2) Trim off barcode      \_\_\_\_\_aattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 3) Quality trimmer      \_\_\_\_\_aattcacgtac ... gtacgt ... acgtacgt(63)
- 4) Chop to 50 bp      \_\_\_\_\_aattcacgtac ... gtacgt(50)
- 5) Merge identical      \_\_\_\_\_aattcacgtac ... gtacgt(50)
- 6) Make fake FASTQ file      \_\_\_\_\_aattcacgtac ... gtacgt(50)
- 7) Run MAQ or bwa



Sequenced in 1 lane



- 1) Split by barcode      GGAAGaattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 2) Trim off barcode      \_\_\_\_\_aattcacgtac ... gtacgt ... acgtacgt ... acgt(76)
- 3) Quality trimmer      \_\_\_\_\_aattcacgtac ... gtacgt ... acgtacgt(63)
- 4) Chop to 50 bp      \_\_\_\_\_aattcacgtac ... gtacgt(50)
- 5) Merge identical      \_\_\_\_\_aattcacgtac ... gtacgt(50)
- 6) Make fake FASTQ file      \_\_\_\_\_aattcacgtac ... gtacgt(50)
- 7) Run MAQ or bwa
  
- 8) Insert data into MySQL database (persistence)
- 9) New automatic scripts for filtering and QC of SNPs
- 10) MSTMap to determine genetic linkage



# Stacks

Stacks is a software pipeline for building loci from short-read sequences, such as those generated on the Illumina platform. Stacks was developed to work with restriction enzyme-based data, such as RAD-seq, for the purpose of building genetic maps and conducting population genomics and phylogeography.



**Download Stacks**

Version 1.19

[Recent Changes \[updated April 23, 2014\]](#)

## Stacks Pipeline

---



### Genetic Maps

Stacks can be used to generate mappable markers from RAD-seq data. Thousands of markers can be generated from a single generation, F1 map as well as markers for traditional F2 and backcross designs. Stacks can export data to JoinMap, OneMap, or R/qtl. These data can be used for examining

### Getting started with Stacks

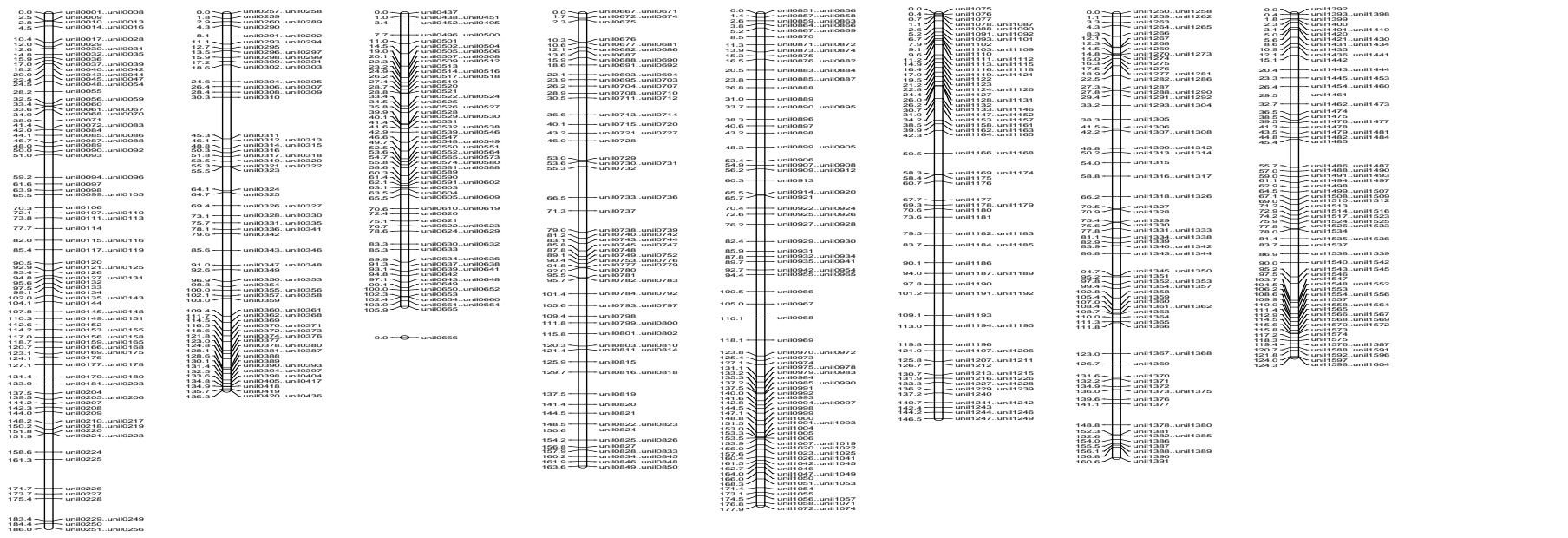
[Stacks Manual](#)

### Frequently Asked Questions

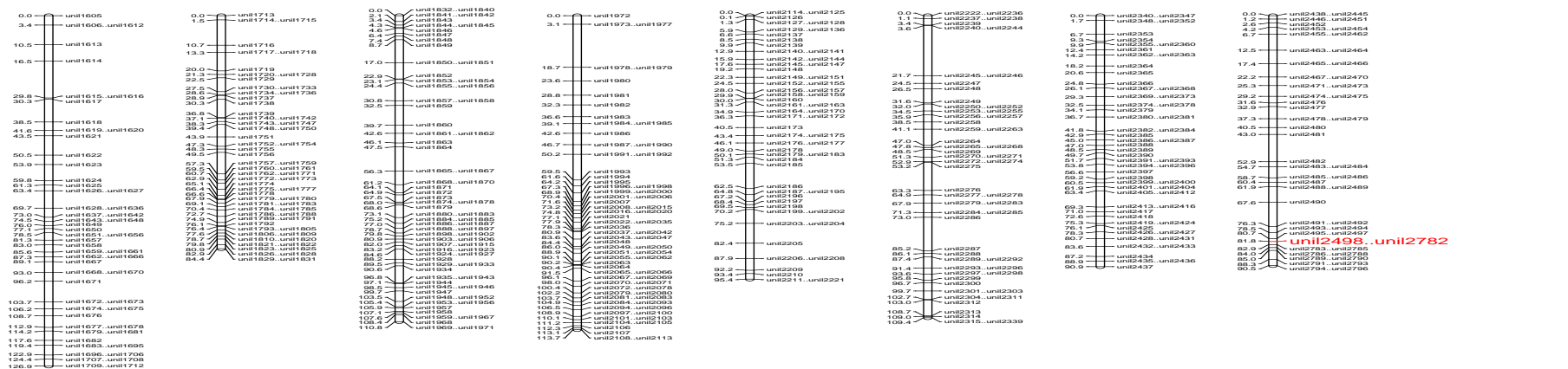
# Fire ant genetic map v1



## LG1 LG2 LG3 LG4 LG5 LG6 LG7 LG8



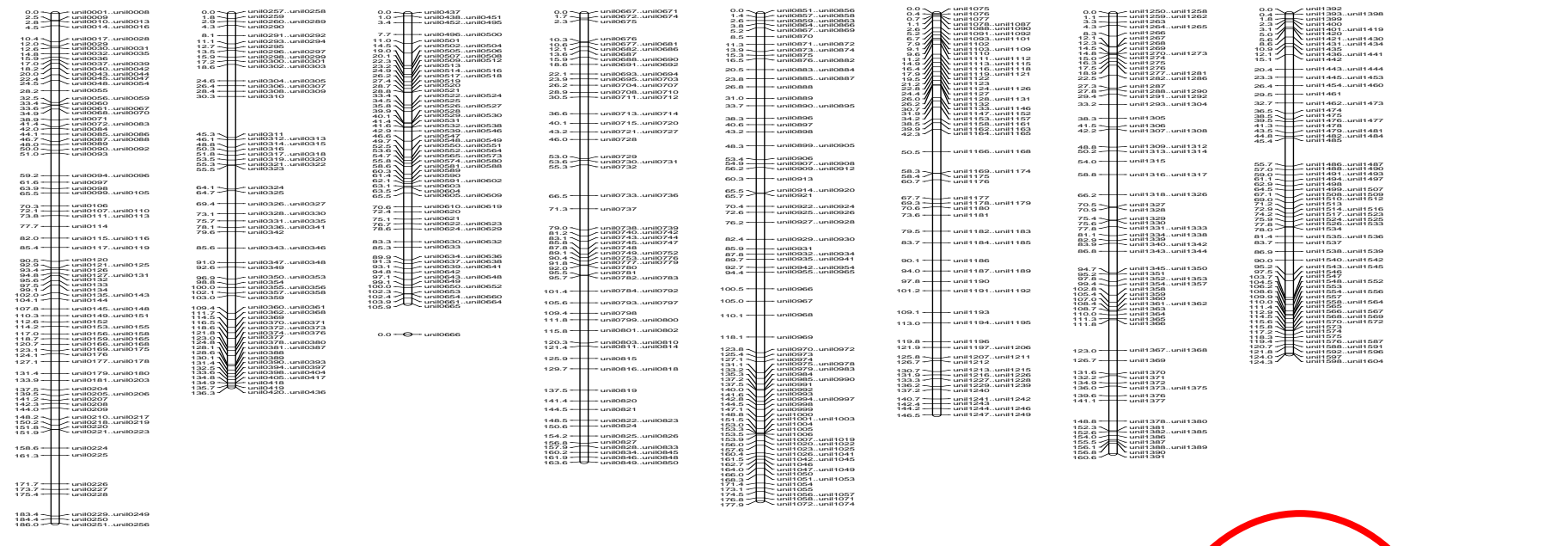
## LG9 LG10 LG11 LG12 LG13 LG14 LG15 LG 16 or S



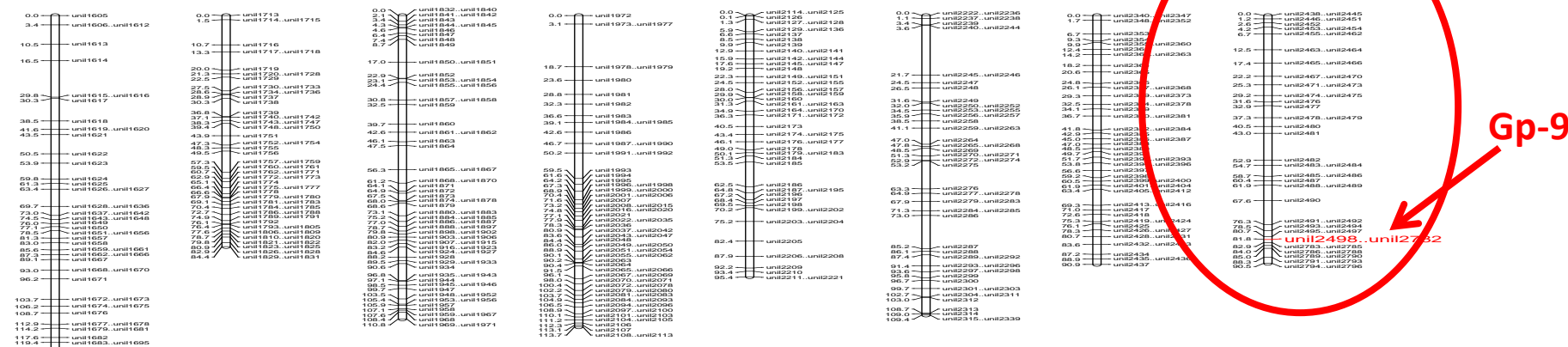
# Fire ant genetic map v1



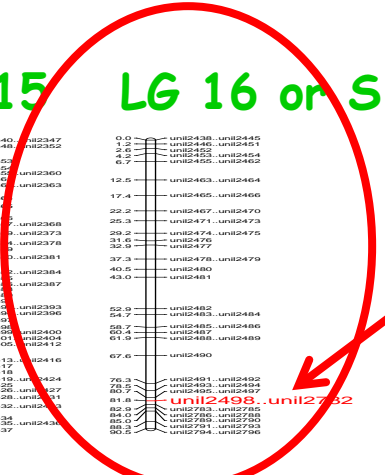
## LG1 LG2 LG3 LG4 LG5 LG6 LG7 LG8



## LG9 LG10 LG11 LG12 LG13 LG14 LG15 LG 16 or S



Gp-9



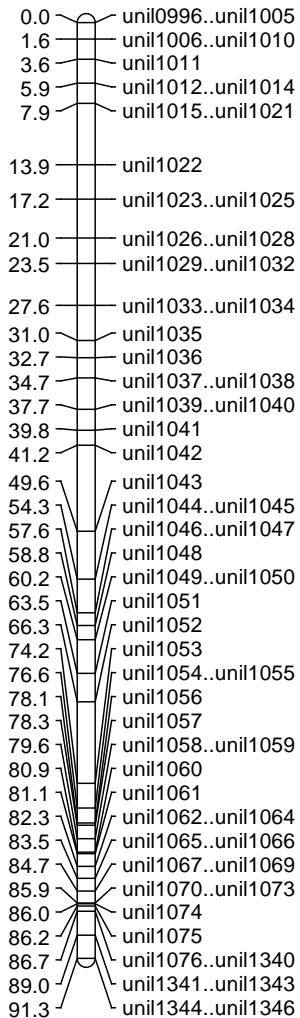


# Large non-combining region around Gp-9



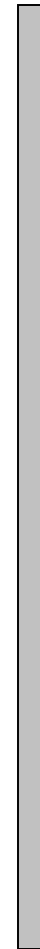
LG 16 = SB or Sb

LG 16 (~23 Mbp)

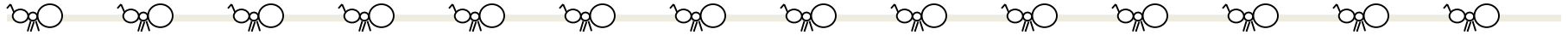


*Genetic map*

*Physical map*

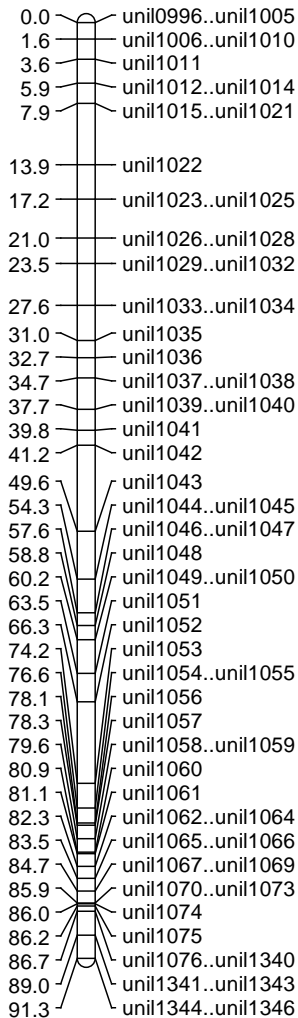


# Large non-combining region around Gp-9



LG 16 = SB or Sb

LG 16 (~23 Mbp)



*Genetic map*

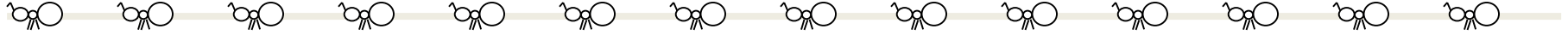
*Physical map*



- Gp-9**
- 265 markers non-recombining
  - 44 scaffolds
  - ~12.7 Mbp

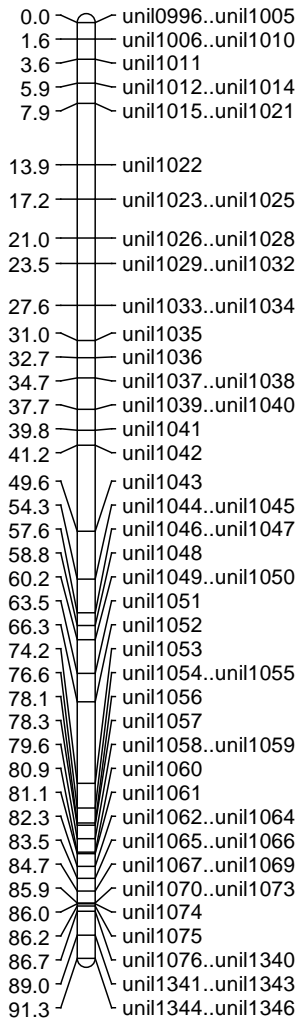


# Large non-combining region around Gp-9



LG 16 = SB or Sb

LG 16 (~23 Mbp)



Genetic map

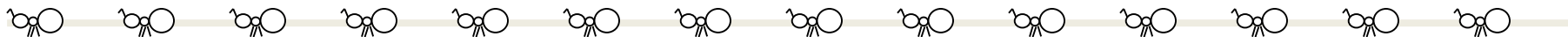
Physical map

- Gp-9
- 265 markers non-recombining
  - 44 scaffolds
  - ~12.7 Mbp



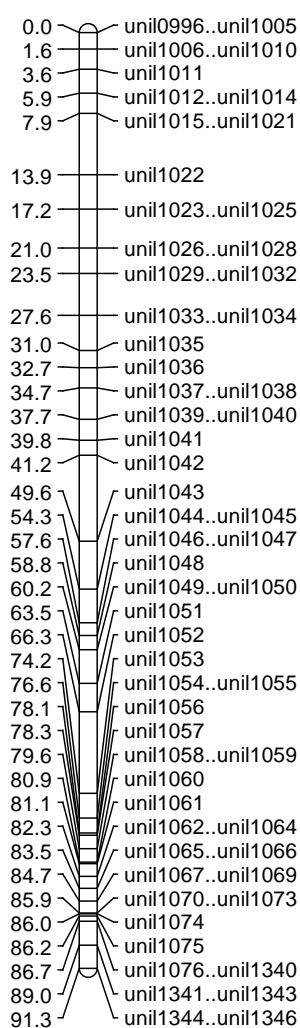
12.7 Mbp  
(55%)

# Large non-combining region around Gp-9



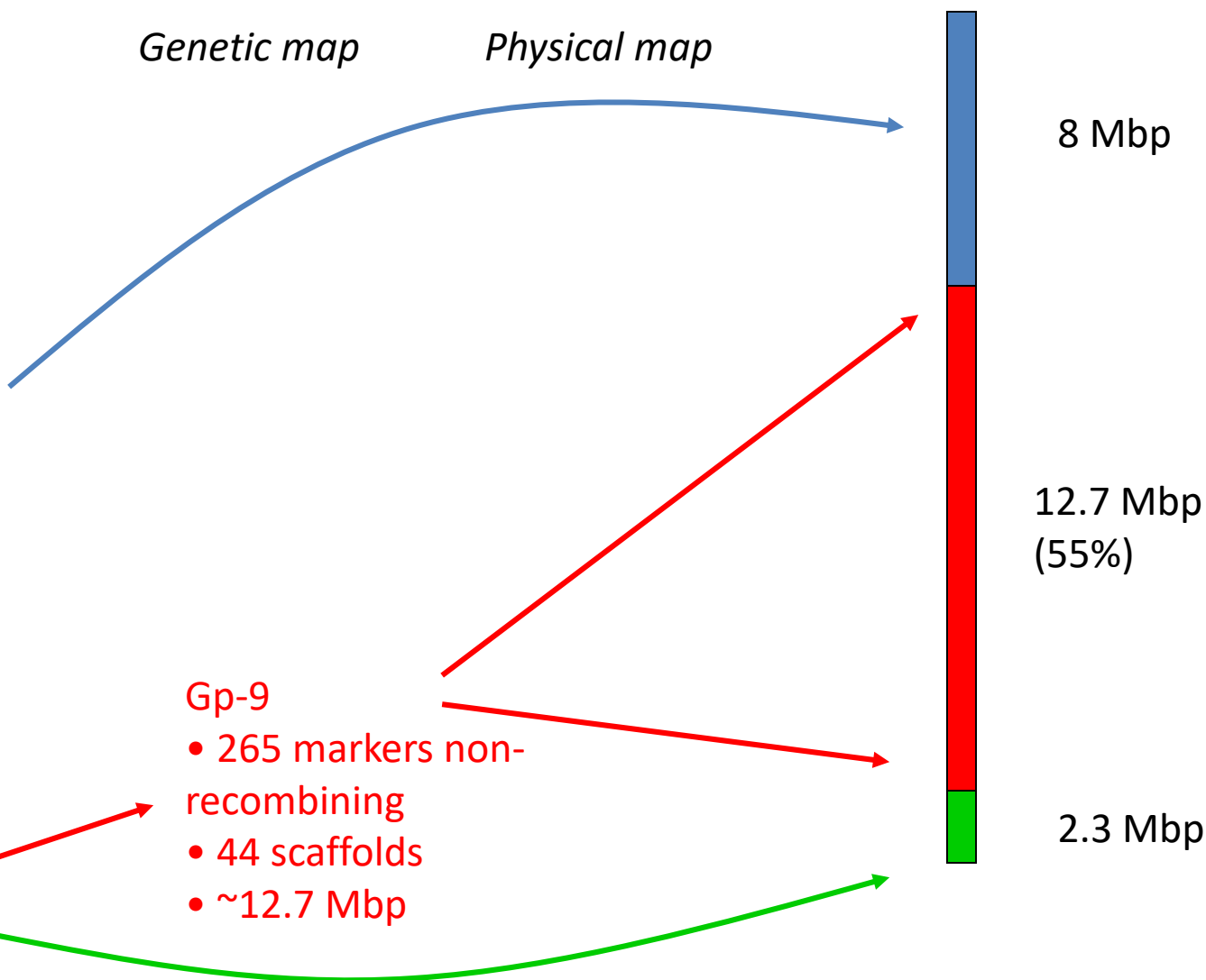
LG 16 = SB or Sb

LG 16 (~23 Mbp)



Genetic map

Physical map



# 2<sup>nd</sup> ant also has supergene



Fire ant case just RAD-seq  
but not really population genomics



# 2<sup>nd</sup> ant also has supergene



Fire ant case just RAD-seq  
but not really population genomics

*Formica selysi*

Also monogyne and polygyne social forms

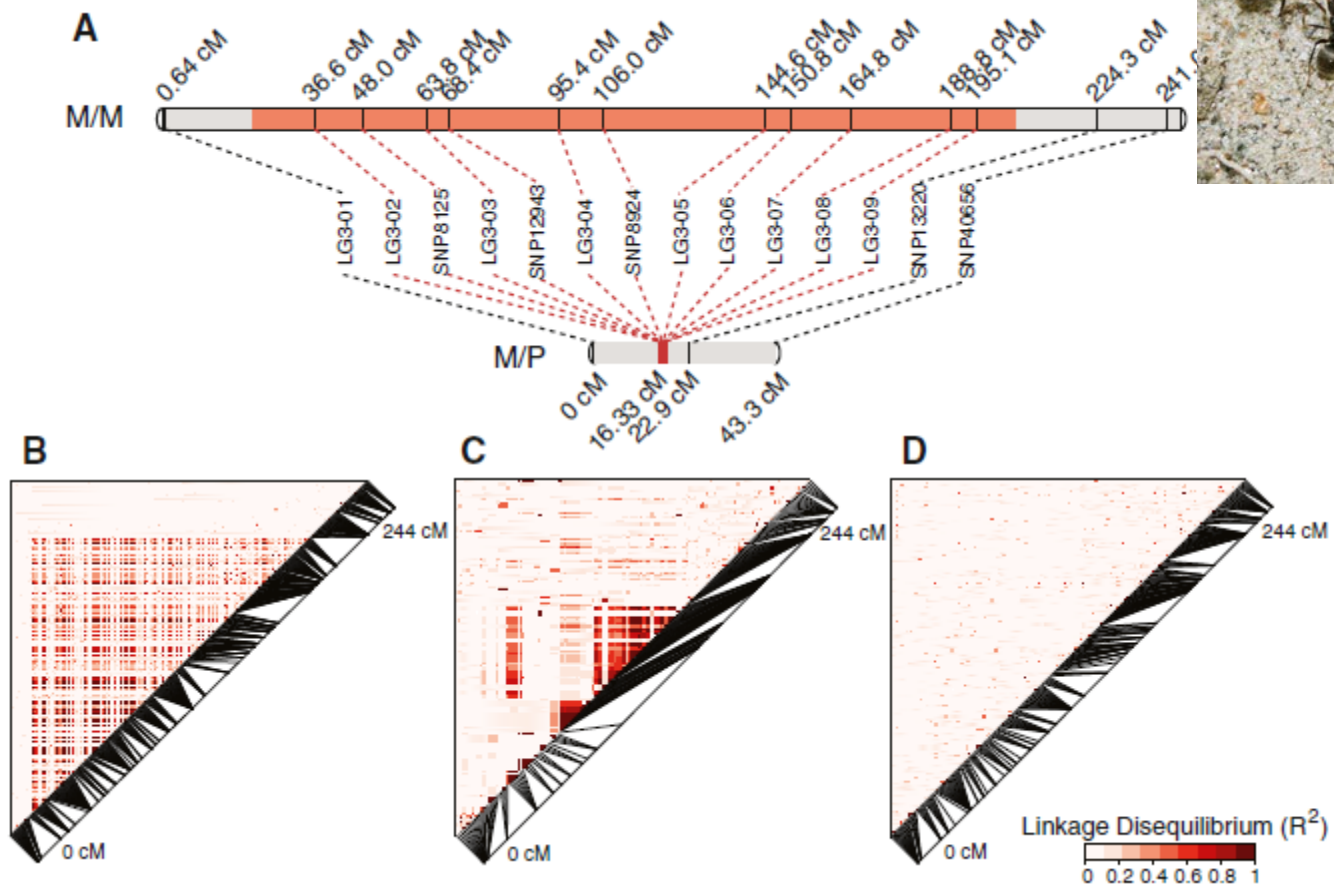
Purcell et al: GBS  
Linkage map

Then population LD verified supergene





# Linkage disequilibrium → supergene



# RAD-seq in phylogeography

*Wyeomyia smithii* mosquito

larvae live in

*Sarracenia purpurea* carnivorous  
pitcher plant



photo: Leonora bug level



# RAD-seq in phylogeography

*Wyeomyia smithii* mosquito

larvae live in

*Sarracenia purpurea* carnivorous  
pitcher plant



photo: Leonora bug level

What is their population structure in the Eastern USA?

Is there an expansion pattern associated with glacial retreat?

# What can I do with this genetic data?

If you have allozymes, AFLPs, sequences, microsatellites, or SNPs, you can

## Examine population structure

- PCA or clustering algorithms
- Relatedness/kinship
- Pairwise divergence ( $F_{ST}$ , Jost's D)
- Phylogeography

## Describe diversity

- Hardy-Weinberg deviation
- Allelic richness
- Nucleotide diversity ( $\pi$ ,  $\theta_W$ )
- Frequency spectra
- LD between markers

and then

or

or

or

or

## Reconstruct demography

- Migration rates
- Population sizes

and then

Did you genotype families?

Do you have data from multiple species?

Do you have ecological or trait data?

## Make a linkage map

## Reconstruct history

- Phylogeny estimate
- Polarize ancestral/derived alleles

Do you have LOTS of markers?

## Scan for differentiation

- $F_{ST}$  outliers
- $F_{ST}$  vs. homozygosity

and then

or

## Find isolating factors

- Mantel testing
- Resistance surface
- BEDASSLE

Do you have ecological or character data?

Are your markers assembled in contigs, aligned to a reference genome, or placed in a linkage map?

## Reconstruct evolution

- Ancestral states
- Diversification rates

## Examine genomic variation

- Nucleotide/haplotype diversity
- Runs of homozygosity
- Extended differentiation
- Extended LD

then

Are genes identified?

## Look at coding variation

- Codon position
- Variant effect prediction

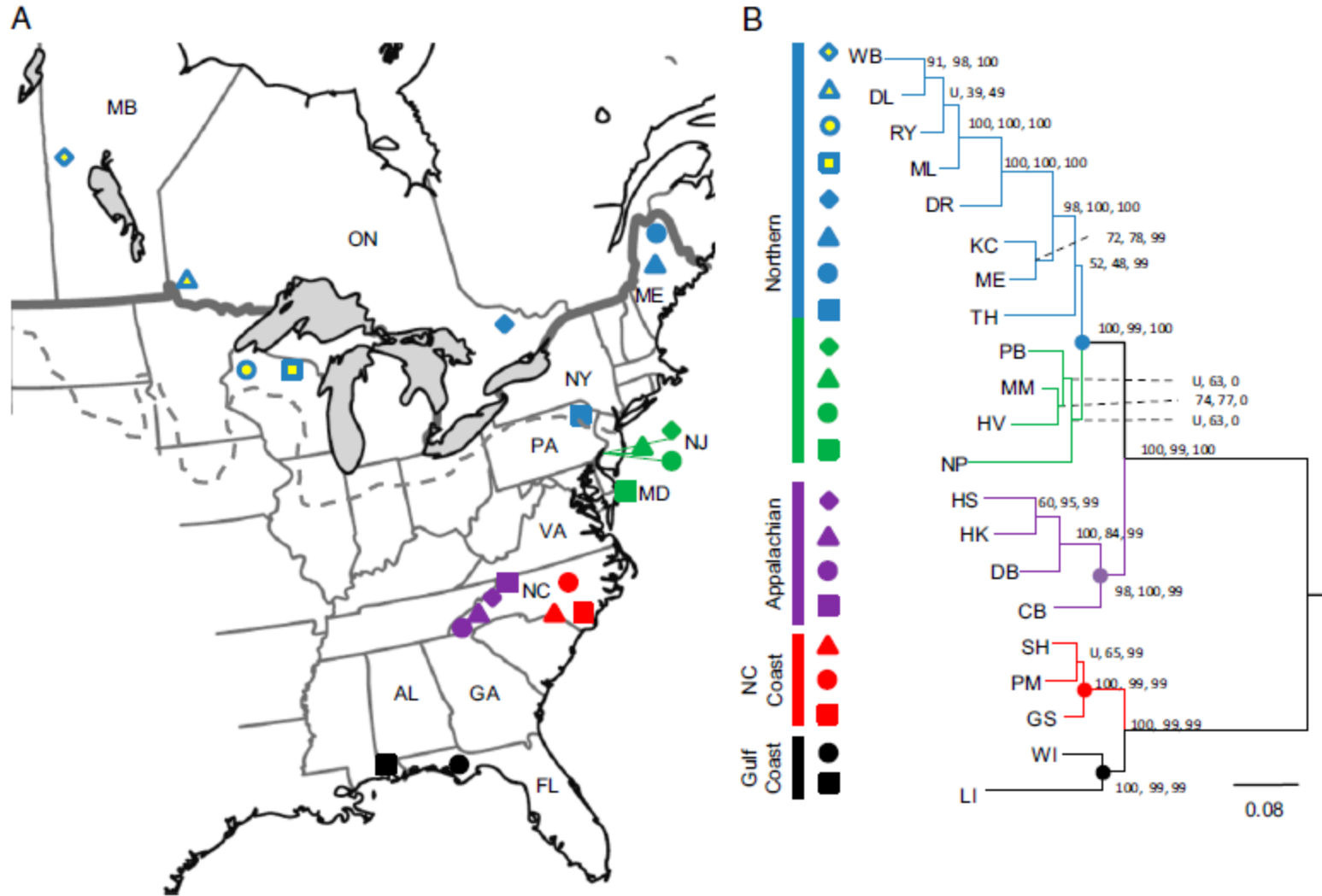
then

Do you have environmental or phenotype data?

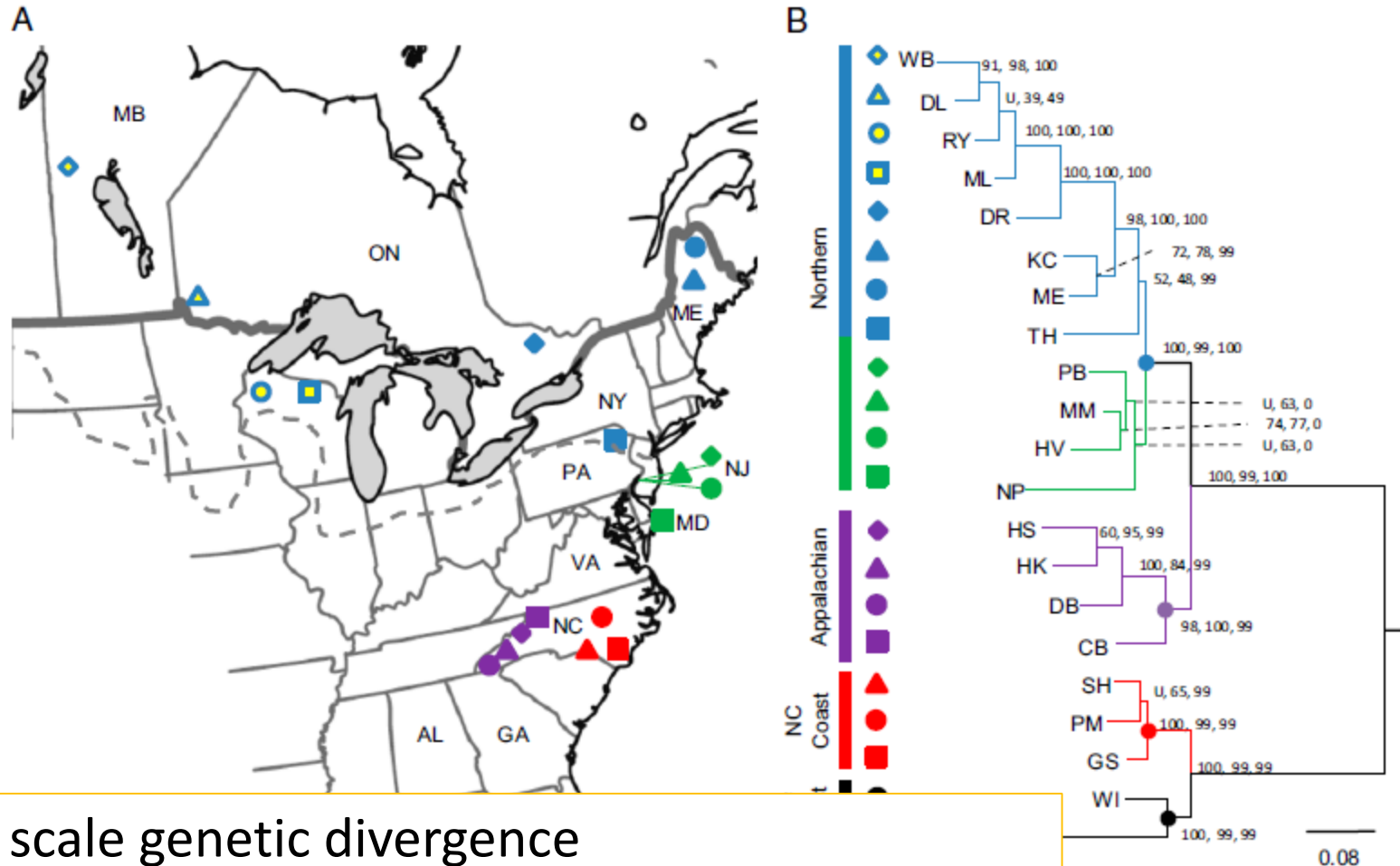
## Test for genotype associations



# *W. smithii* samples and phylogenetic tree from RAD

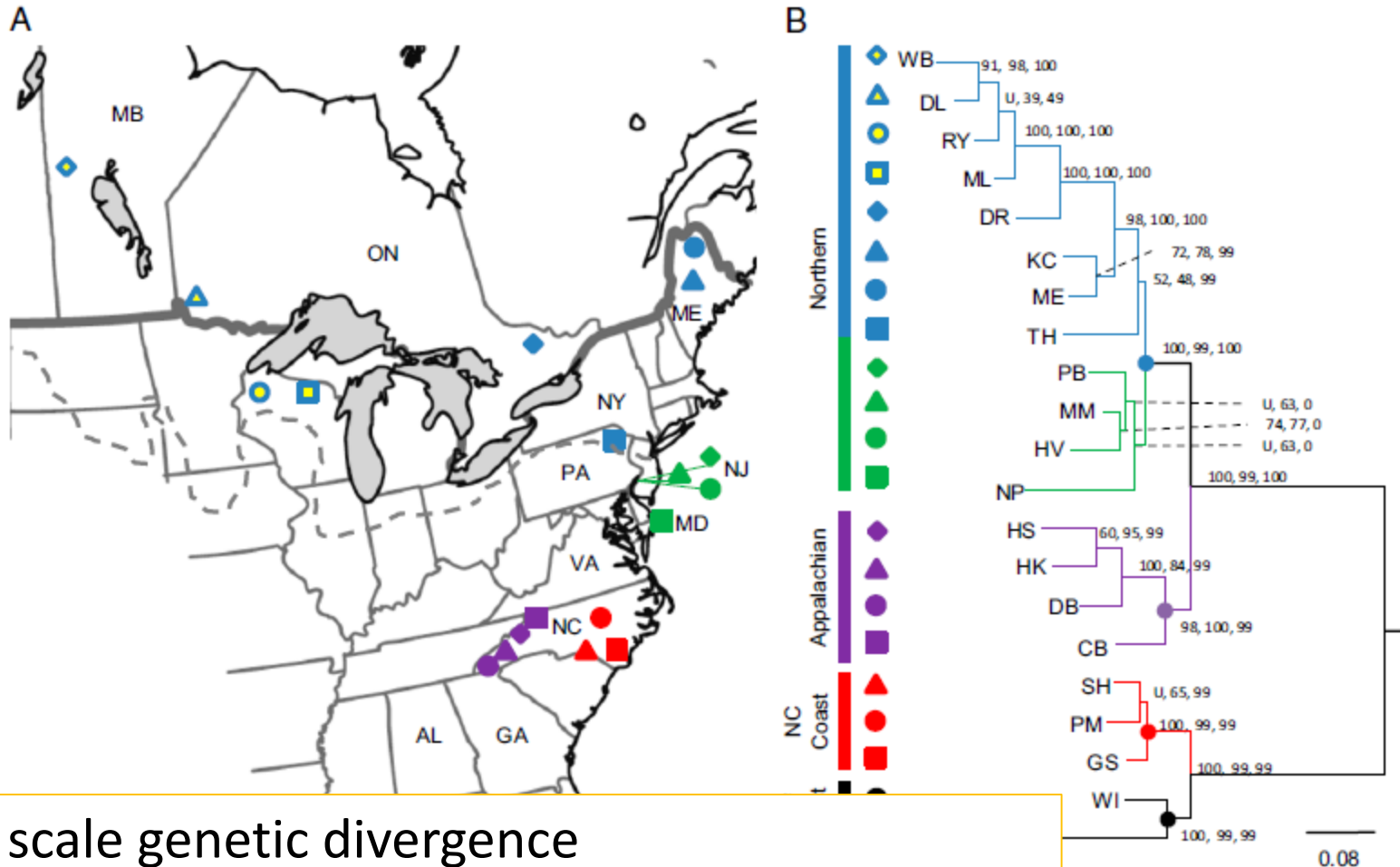


# *W. smithii* samples and phylogenetic tree from RAD



Fine scale genetic divergence  
Evidence for glacial retreat associated expansion

# *W. smithii* samples and phylogenetic tree from RAD



Fine scale genetic divergence  
Evidence for glacial retreat associated expansion

Future – local photoperiodic response adaptations

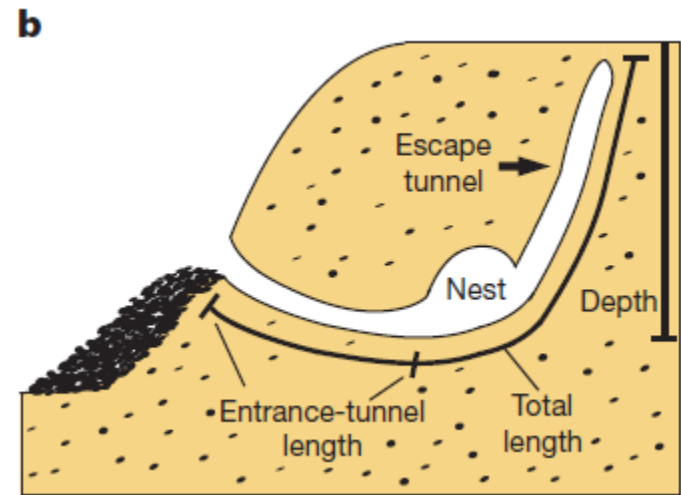


# Discrete genetic modules are responsible for complex burrow evolution in *Peromyscus* mice

Jesse N. Weber<sup>1,†</sup>, Brant K. Peterson<sup>1,2</sup> & Hopi E. Hoekstra<sup>1,2</sup>

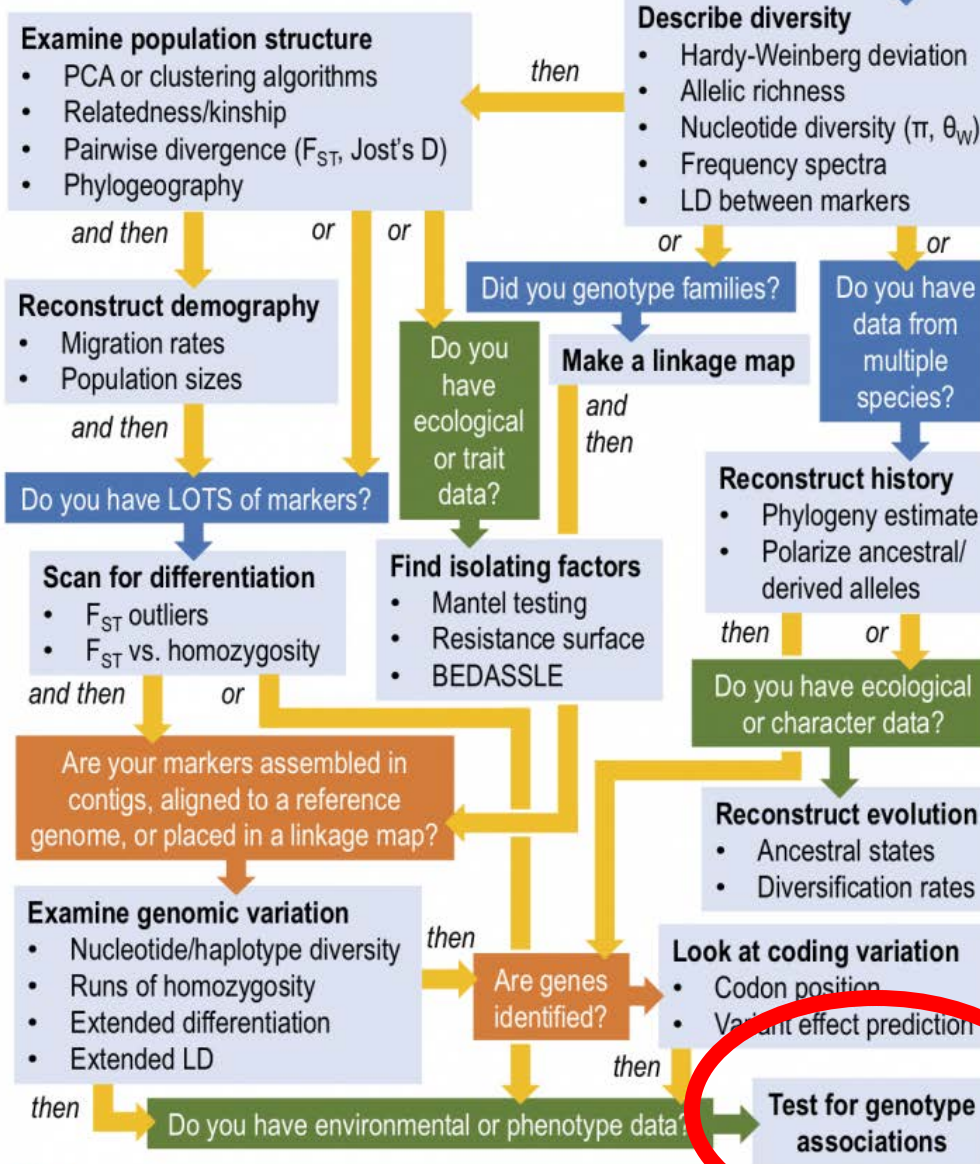


Vera Domingues/Hopi Hoekstra



# What can I do with this genetic data?

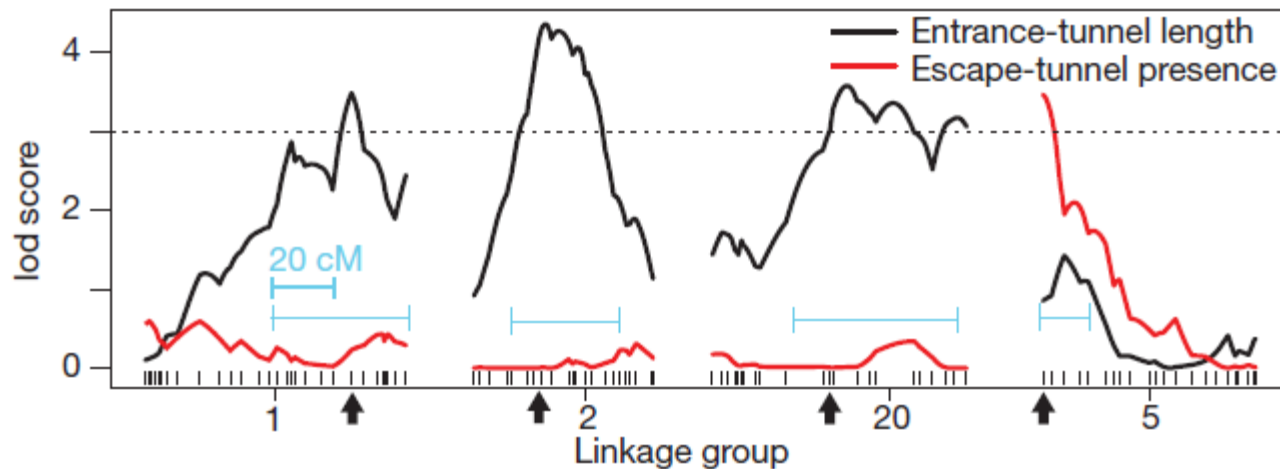
If you have allozymes, AFLPs, sequences, microsatellites, or SNPs, you can



# Tunnel building at least 4 loci

Crossed 2 species of mice:  
oldfield mice (escape tunnels) x deer mice (simple nests)

QTL mapping with ddRAD-seq

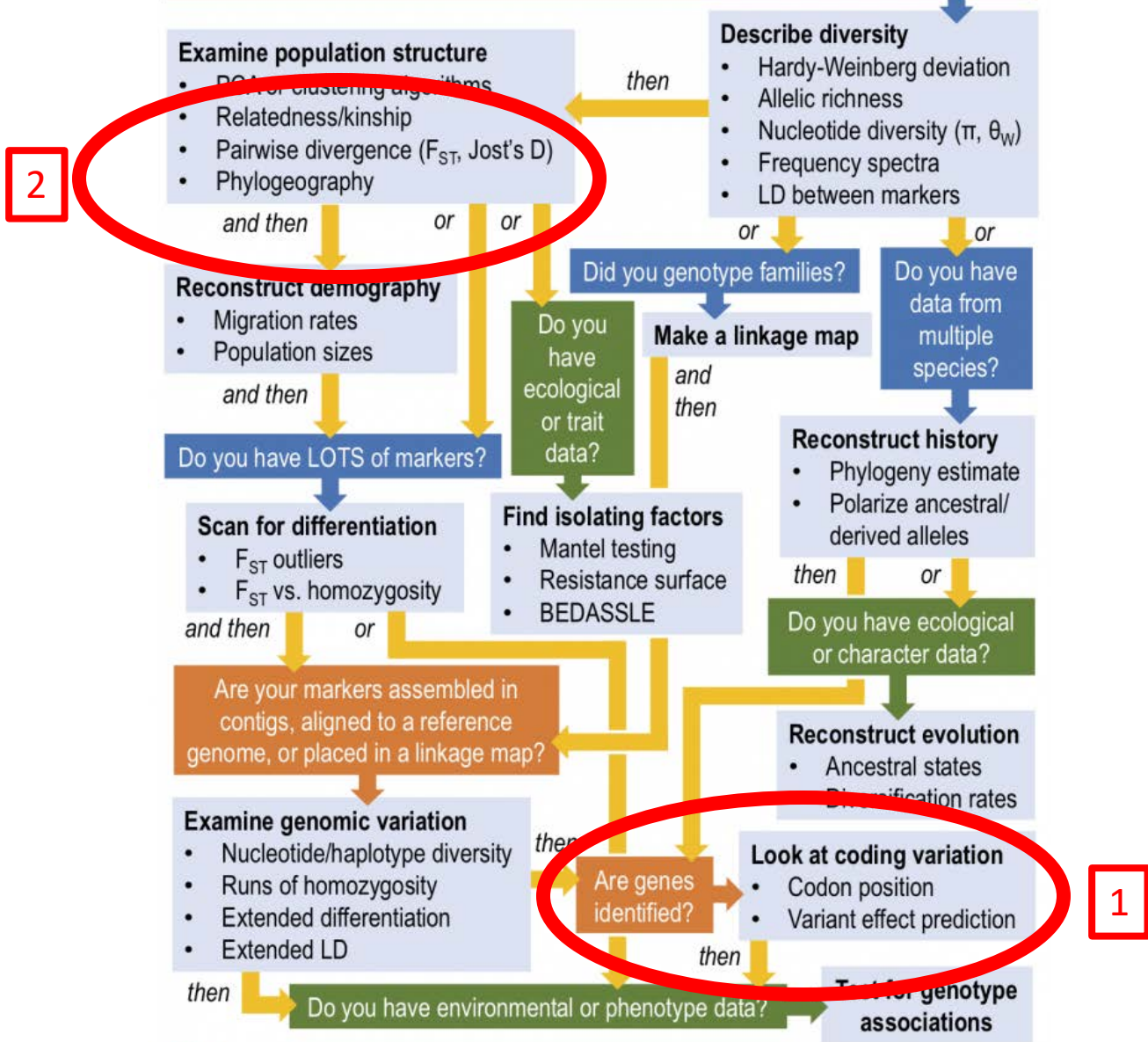




# Exon sequencing

# What can I do with this genetic data?

If you have allozymes, AFLPs, sequences, microsatellites, or SNPs, you can



Targeted (re)sequencing or **exon capture**

# Targeted (re)sequencing or **exon capture**

- Target for specific regions (e.g., exons)

# Targeted (re)sequencing or **exon capture**

- Target for specific regions (e.g., exons)
- Subset of the genome can answer your question

# Targeted (re)sequencing or **exon capture**

- Target for specific regions (e.g., exons)
- Subset of the genome can answer your question
- Reduce effort and cost
  - Reduced data storage (sometimes never analyzed)

# Targeted (re)sequencing or **exon capture**

- Target for specific regions (e.g., exons)
- Subset of the genome can answer your question
- Reduce effort and cost
  - Reduced data storage (sometimes never analyzed)
- Increase sample size

# Targeted (re)sequencing or **exon capture**

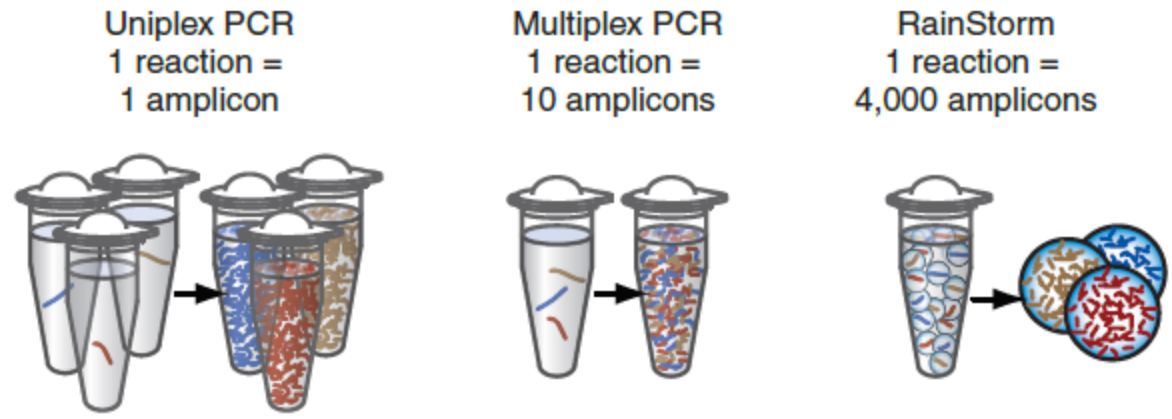
- Target for specific regions (e.g., exons)
- Subset of the genome can answer your question
- Reduce effort and cost
  - Reduced data storage (sometimes never analyzed)
- Increase sample size
- Something of a stop gap until WGS is too cheap
  - Genome too big



# Exon capture

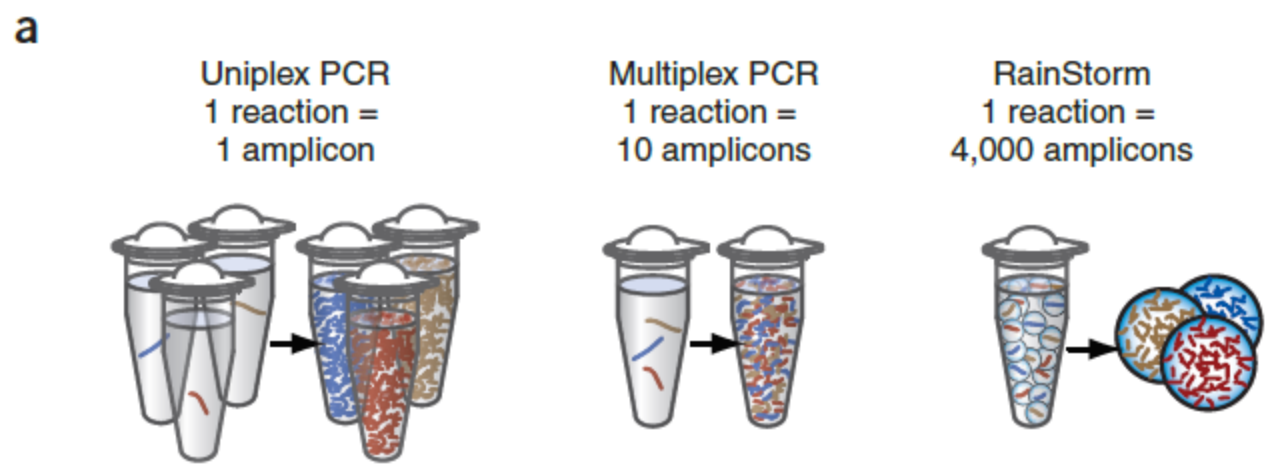
1,000's

a

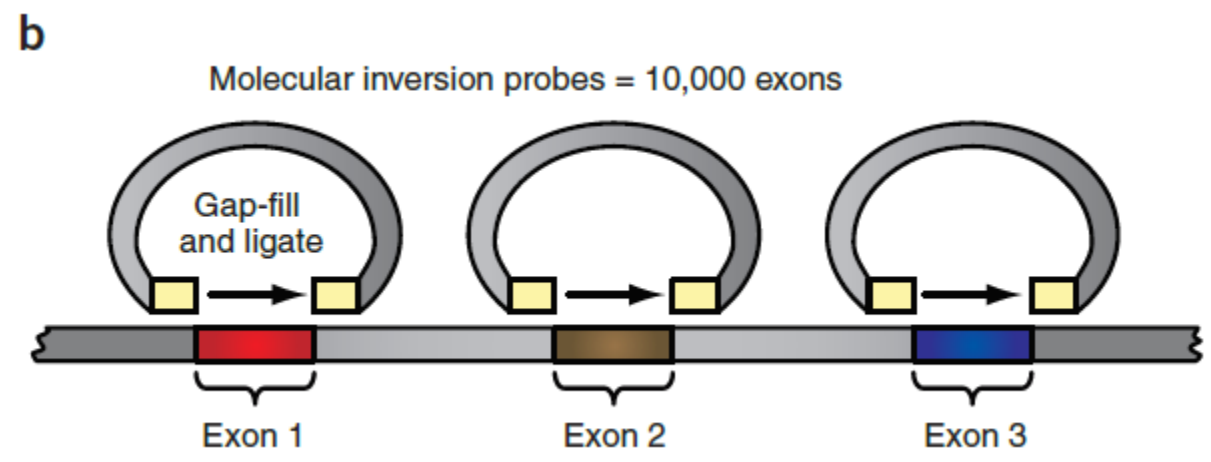


# Exon capture

1,000's



MIP  
10,000's



# Exon capture

c

Hybrid capture > 100,000 exons



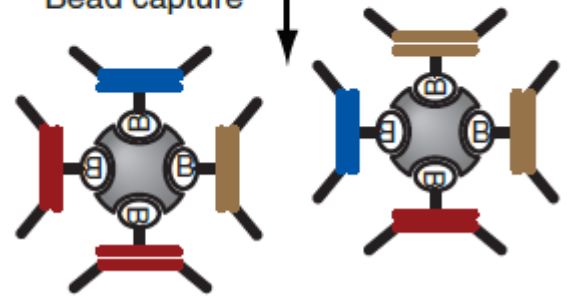
Adapter-modified  
shotgun library

Array capture

Solution  
hybridization

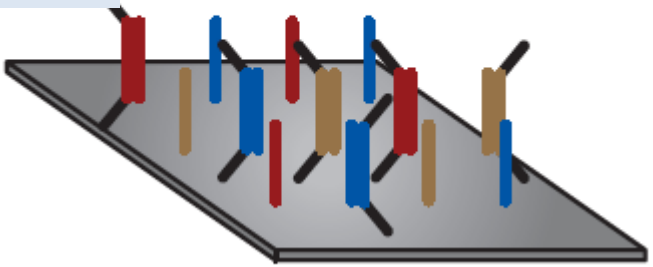


Bead capture



100,000's

Agilent  
Nimblegen  
Illumina



# Exon capture – applications

# Exon capture – applications

- Mendelian disease discovery  
Since many are protein coding mutations

Strategies:

Family based

Extreme phenotypes (e.g., cardiovascular health)

# Exon capture – applications

- Mendelian disease discovery

Since many are protein coding mutations

Strategies:

Family based

Extreme phenotypes (e.g., cardiovascular health)

- Non-human organisms

Anything you can imagine with many (protein coding) loci

Cloning homozygous lethals (via heterozygotes)

Climate change driven subdivision

# Exon capture – applications

- Mendelian disease discovery

Since many are protein coding mutations

Strategies:

Family based

Extreme phenotypes (e.g., cardiovascular health)

- Non-human organisms

Anything you can imagine with many (protein coding) loci

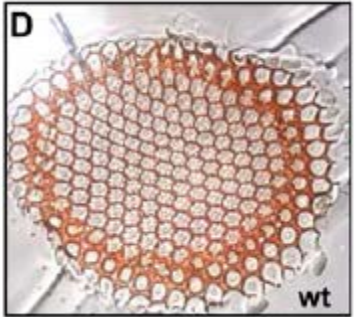
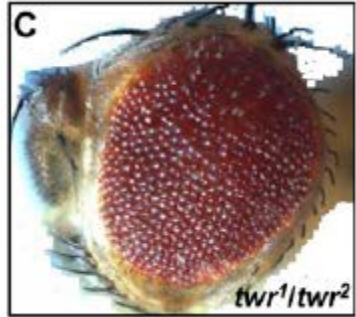
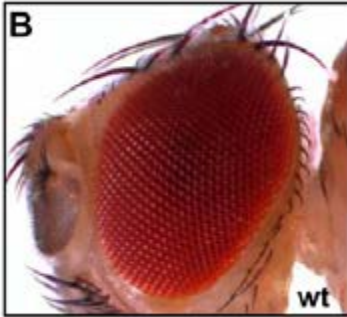
Cloning homozygous lethals (via heterozygotes)

Climate change driven subdivision

- Caveat: nontrivial (although not huge) cost of setup

# Exon capture – homozygous lethal case study find in heterozygote

*twisted bristles roughened eye (twr)* mutant



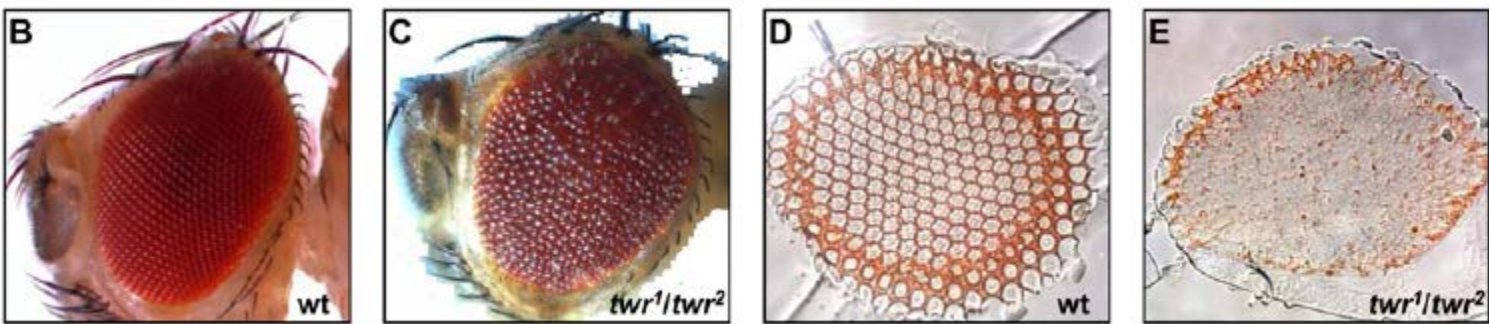
*twr/twr* → lethal

*twr<sup>1</sup>/twr<sup>2</sup>* → viable

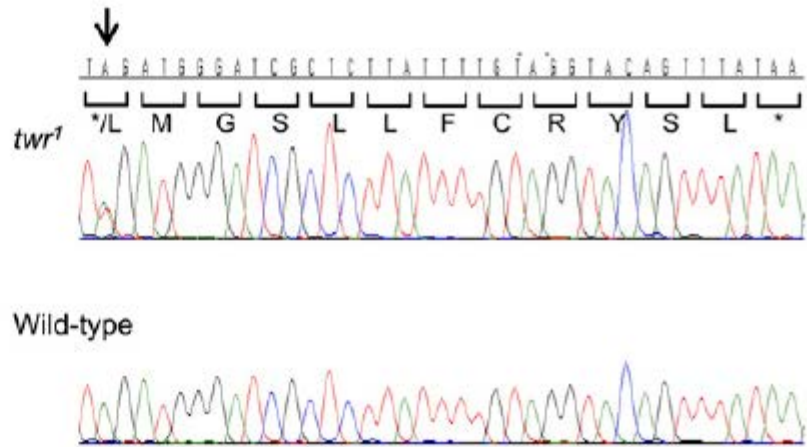
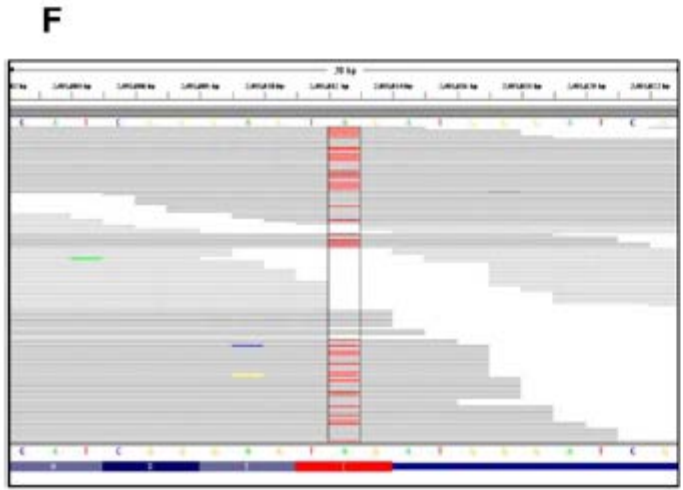


# Exon capture – homozygous lethal case study find in heterozygote

*twisted bristles roughened eye (twr)* mutant



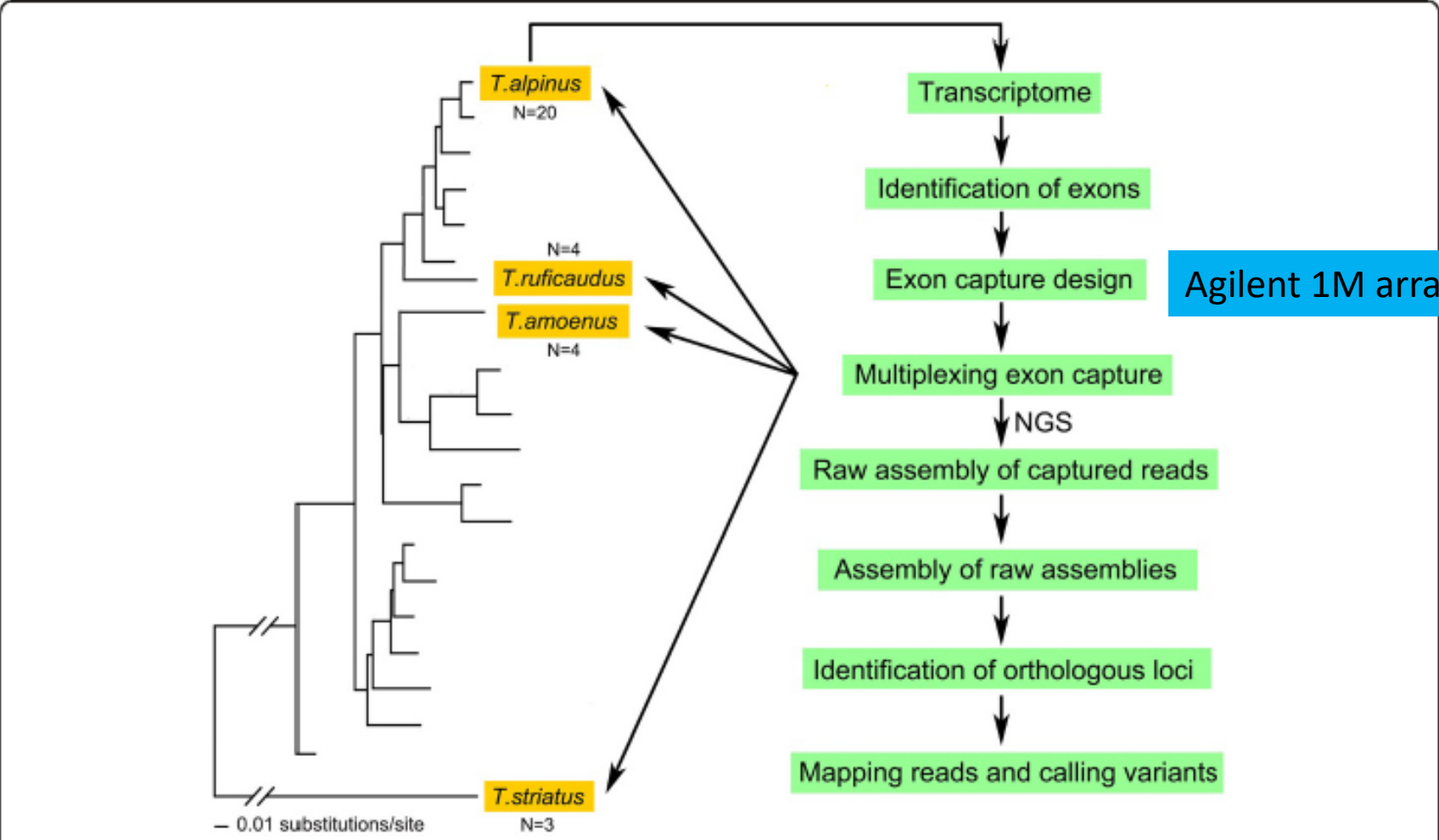
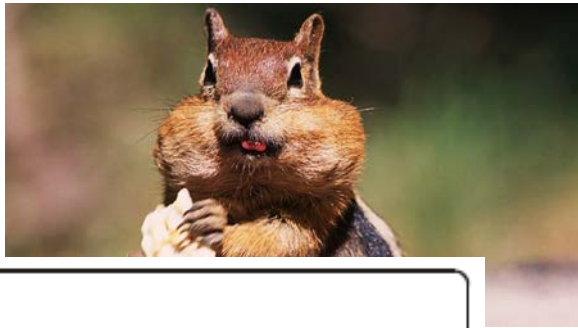
twr/twr → lethal  
twr1/twr2 → viable



# Exon capture – chipmunk with no reference genome

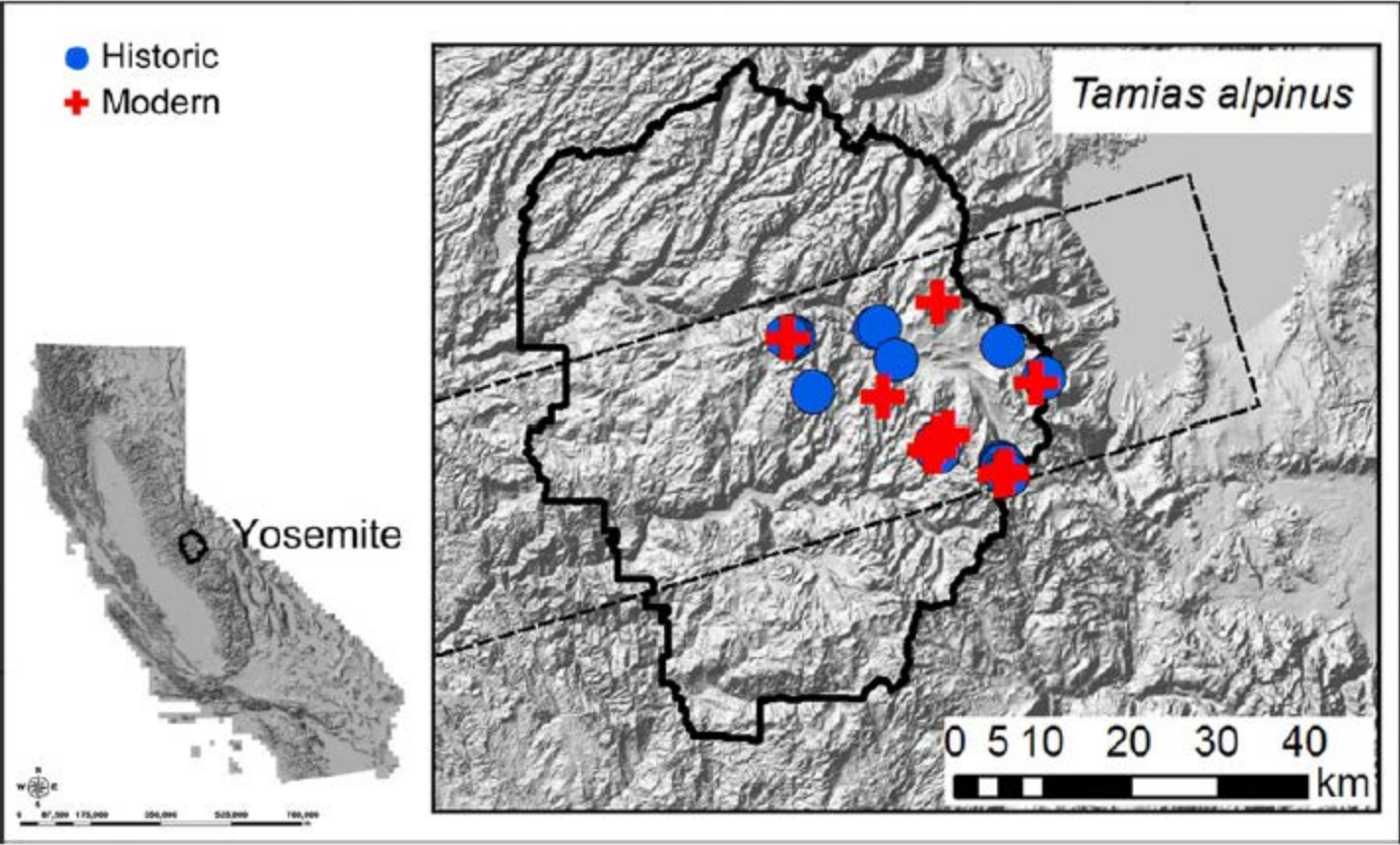


# Exon capture – chipmunk with no reference genome

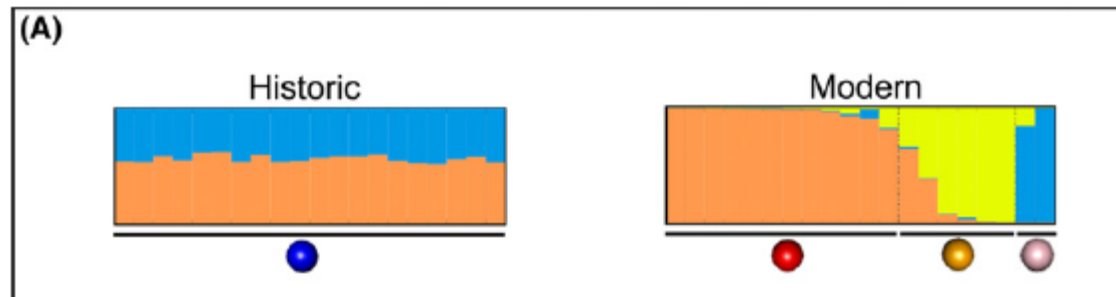


**Figure 1** An overall work flow of this study. The *Tamias* phylogenetic tree is modified from [13] by replacing the outgroup species with *T. striatus*. The *Tamias* species that were not under investigation in the present study are not shown.

# Exon capture – chipmunk → museum samples and range contraction



# Exon capture – chipmunk → museum samples and range contraction



STRUCTURE analysis reveals more modern subdivision  
Consistent with reduced gene flow associated with  
climate change

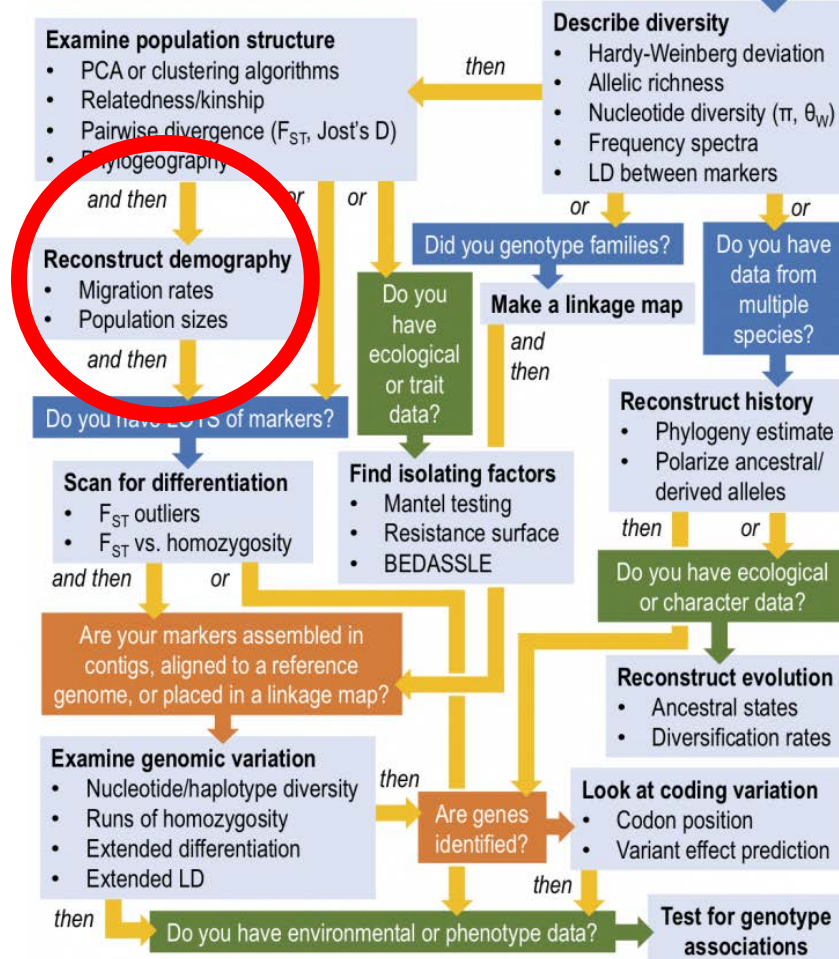
# RE delimited sequencing



# RE delimited sequencing

## What can I do with this genetic data?

If you have allozymes, AFLPs, sequences, microsatellites, or SNPs, you can



# A reduced representation approach to population genetic analyses and applications to human evolution

2011

Francesca Luca,<sup>1</sup> Richard R. Hudson,<sup>1,2,3</sup> David B. Witonsky,<sup>1</sup> and Anna Di Rienzo<sup>1,3</sup>

Reduced representation simply by digestion (Van Tassell et al 2008)

Population genomics goals:

- Test serial founder model of human global dispersal

- Estimate timing of out of Africa

- First estimate of the colonization of Australia



# A reduced representation approach to population genetic analyses and applications to human evolution

2011

Francesca Luca,<sup>1</sup> Richard R. Hudson,<sup>1,2,3</sup> David B. Witonsky,<sup>1</sup> and Anna Di Rienzo<sup>1,3</sup>

Reduced representation simply by digestion (Van Tassell et al 2008)

Population genomics goals:

Test serial founder model of human global dispersal

Estimate timing of out of Africa

First estimate of the colonization of Australia

Nucleotide diversity

Heterozygosity

Allele frequency spectrum

...

# Samples: 19 individuals (18 pops) around the world

A



N = 1-2 per population

Use power of many loci to estimate population parameters

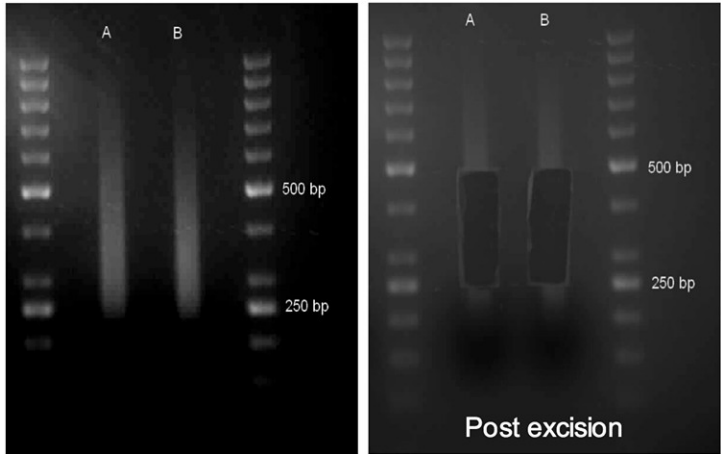
- nucleotide diversity

# Molecular protocol – relatively simple

- 1) Cut with RsaI restriction enzyme  
~9 million sites total  
less near CpG islands  
not near repeats



- 2) Select 70-75 bp fragments  
(really 50-100 bp)  
~51,000 sites in human genome



Amador, D. M., Wang, C., Holland, K. H. and Mou, Z. (2013)

- 3) Illumina single end sequencing  
36 bp (?)  
1 individual/1-2 lanes  
wanted ~50x coverage  
(got <20x)



# Bioinformatics – mapping, filtering

## Strategy:

Map to reduced genome

Limited to 36 bp flanking RsaI sites

Removed redundant loci with  $\leq 4$  mismatch

# Bioinformatics – mapping, filtering

## Strategy:

- Map to reduced genome
- Limited to 36 bp flanking RsaI sites
- Removed redundant loci with  $\leq 4$  mismatch



## Filtering: QC or clean up data

- Remove 2 nts from RE <GTAC>
- 3<sup>rd</sup> nt also biased
- keep phred score  $\geq 20$
- if >2 alleles (keep best 2)

# Bioinformatics – mapping, filtering

## Strategy:

- Map to reduced genome
- Limited to 36 bp flanking RsaI sites
- Removed redundant loci with  $\leq 4$  mismatch



## Filtering: QC or clean up data

- Remove 2 nts from RE <GTAC>
- 3<sup>rd</sup> nt also biased
- keep phred score  $\geq 20$
- if >2 alleles (keep best 2)

## Caveats:

- Cannot know pcr duplicates
- Loss due to RE polymorphism

# Bioinformatics – mapping, filtering

## Strategy:

- Map to reduced genome
- Limited to 36 bp flanking RsaI sites
- Removed redundant loci with  $\leq 4$  mismatch



## Filtering: QC or clean up data

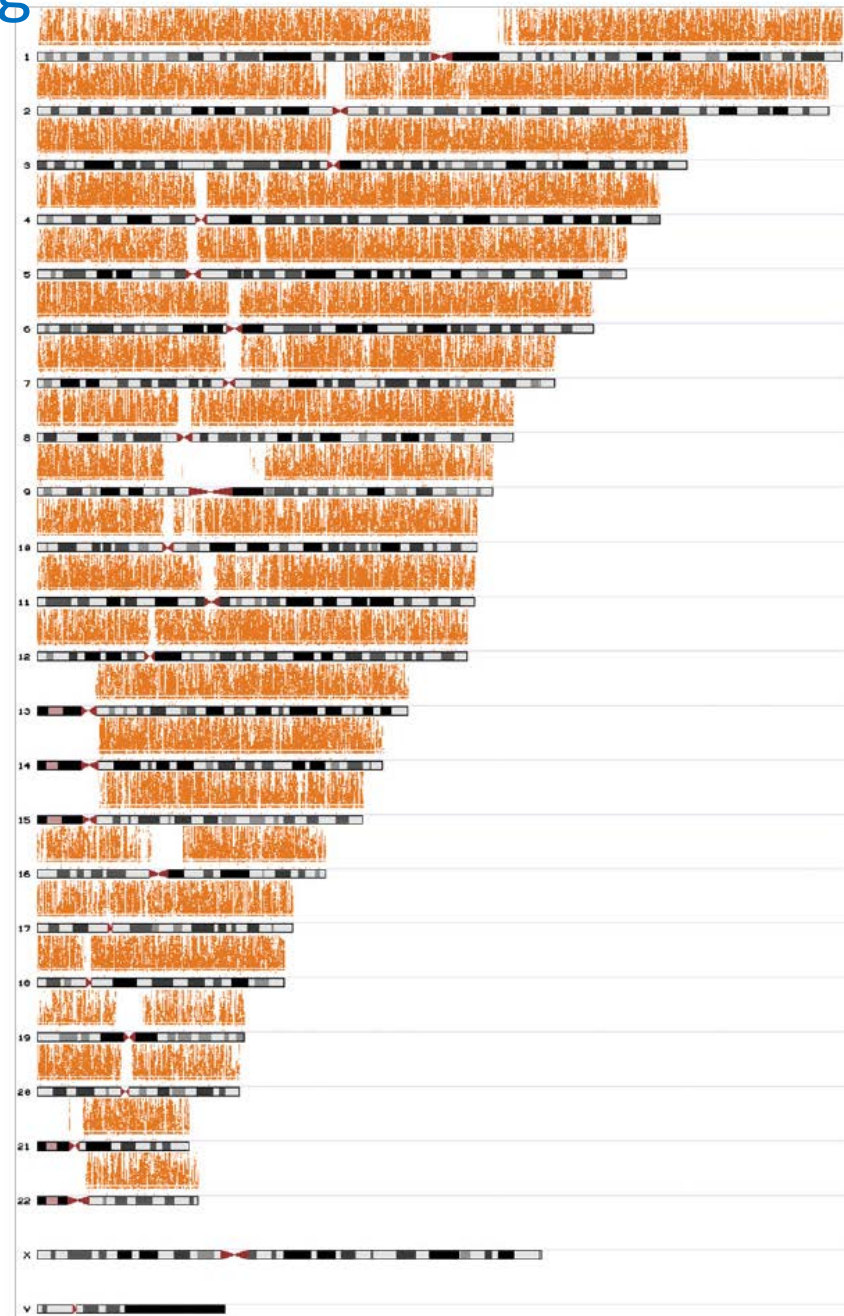
- Remove 2 nts from RE <GTAC>
- 3<sup>rd</sup> nt also biased
- keep phred score  $\geq 20$
- if >2 alleles (keep best 2)

## Caveats:

- Cannot know pcr duplicates
- Loss due to RE polymorphism

## RESULT:

- 61.7% uniquely mapping
- Even/random coverage



# Data analysis #1

This method: 1000's of sites per genome

Nucleotide diversity ( $\pi$ ):

Calculate: Proportion of heterozygous sites within each genome is estimate of its  $\pi$  population

0.61-1.08/1000 <lower than previous estimates>

Test serial founder effect:

Expect negative correlation

$\pi \sim$  geographic distance



# Data analysis #1

This method: 1000's of sites per genome

Nucleotide diversity ( $\pi$ ):

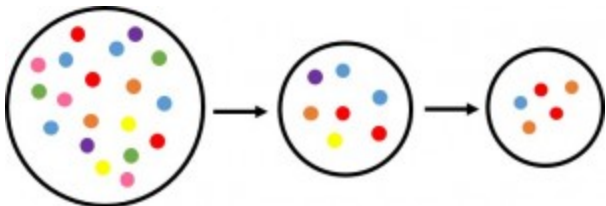
Calculate: Proportion of heterozygous sites within each genome is estimate of its  $\pi$  population

0.61-1.08/1000 <lower than previous estimates>

Test serial founder effect:

Expect negative correlation

$\pi \sim$  geographic distance



# Data analysis #1

This method: 1000's of sites per genome

Nucleotide diversity ( $\pi$ ):

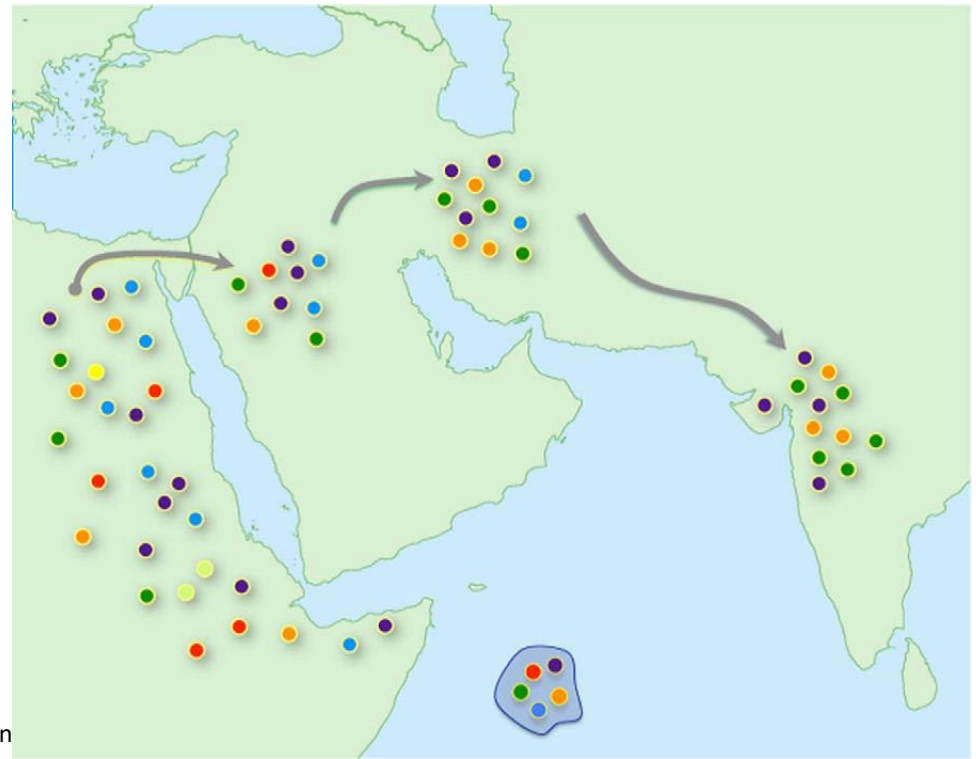
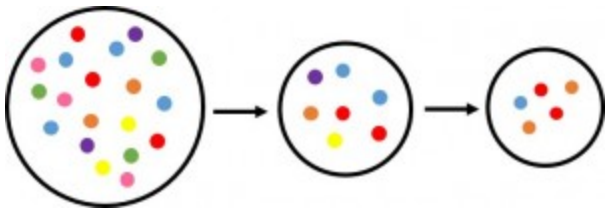
Calculate: Proportion of heterozygous sites within each genome is estimate of its  $\pi$  population

0.61-1.08/1000 <lower than previous estimates>

Test serial founder effect:

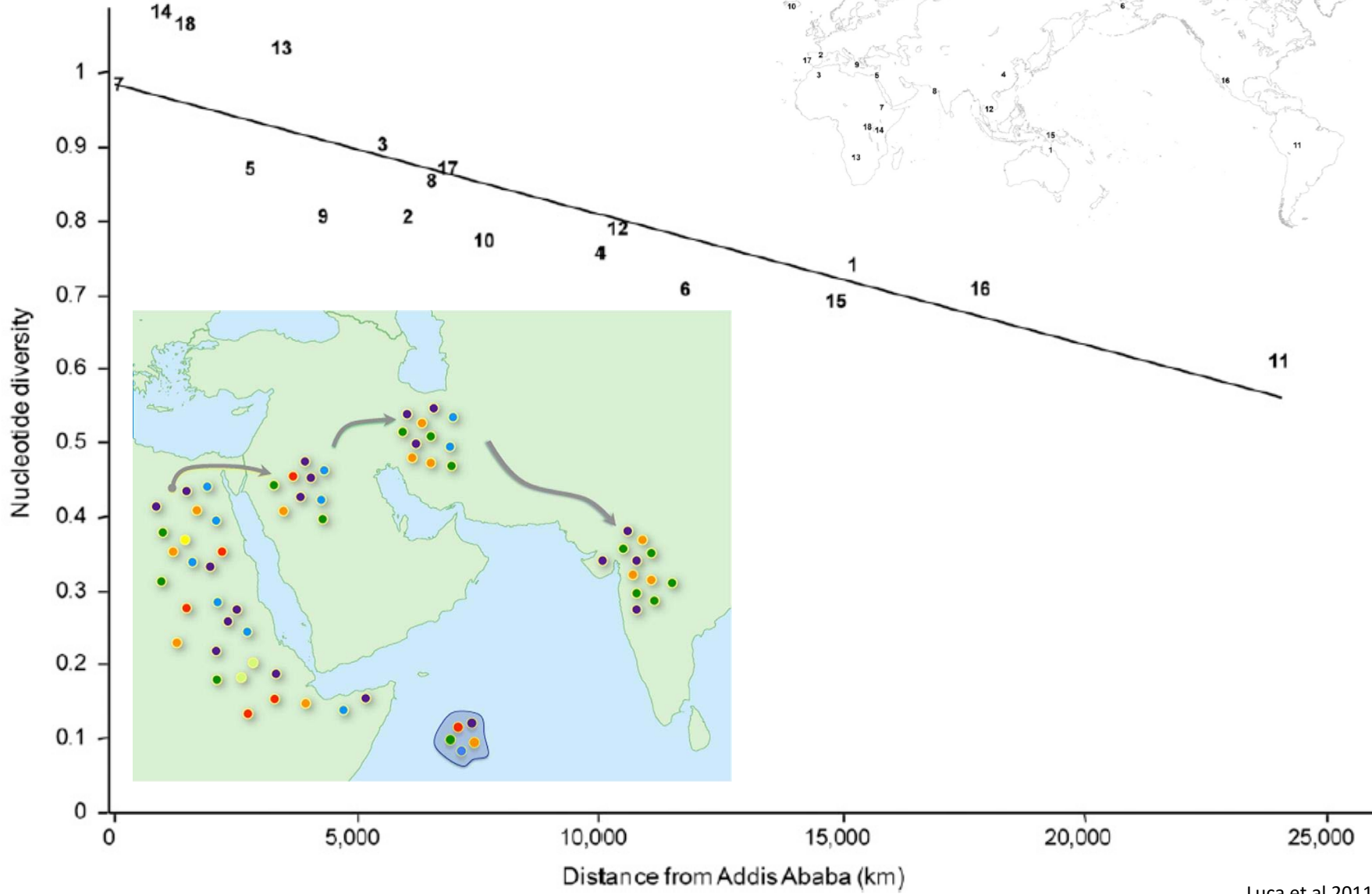
Expect negative correlation

$\pi \sim$  geographic distance



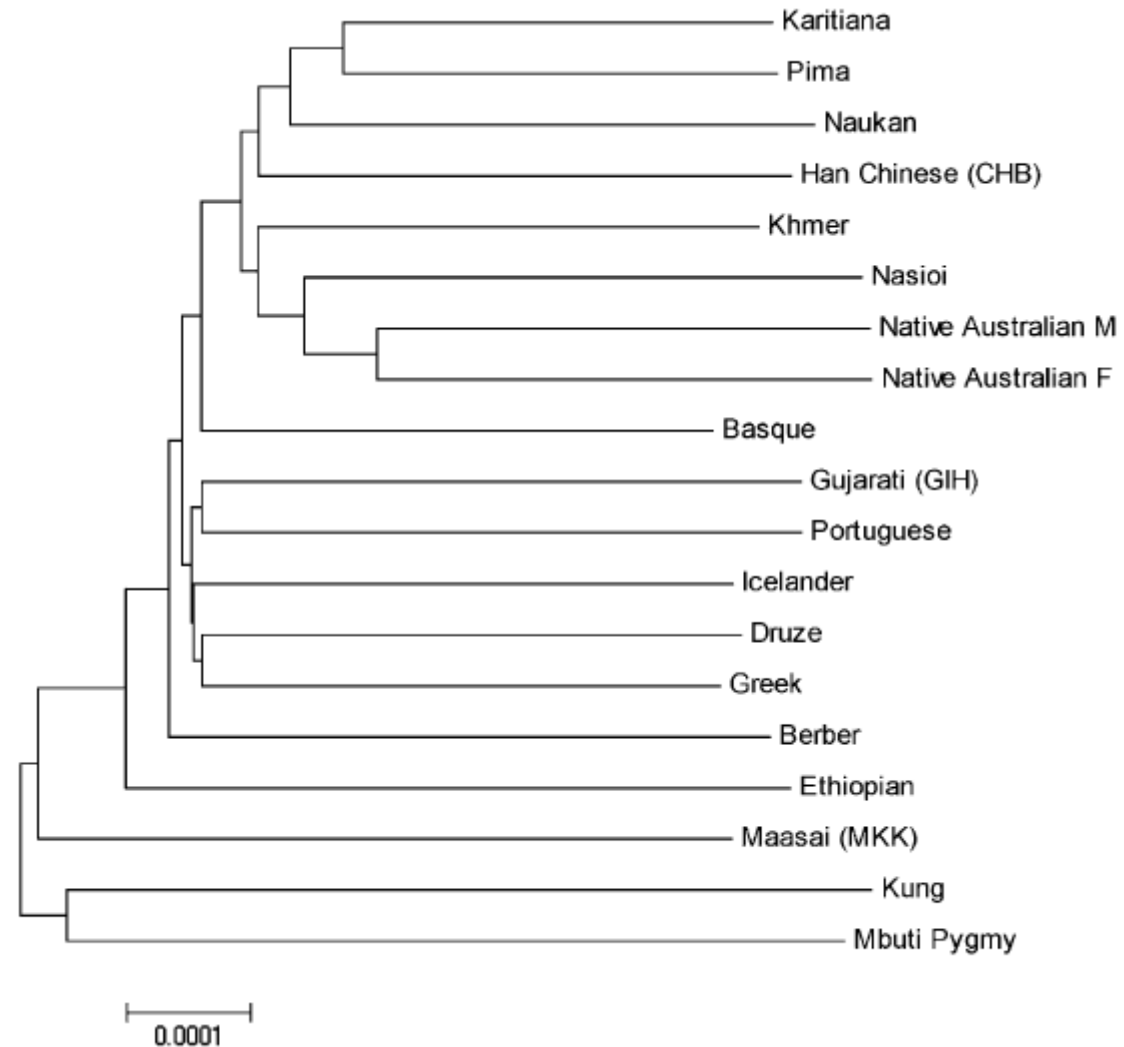
# Data consistent with serial founder effect

B



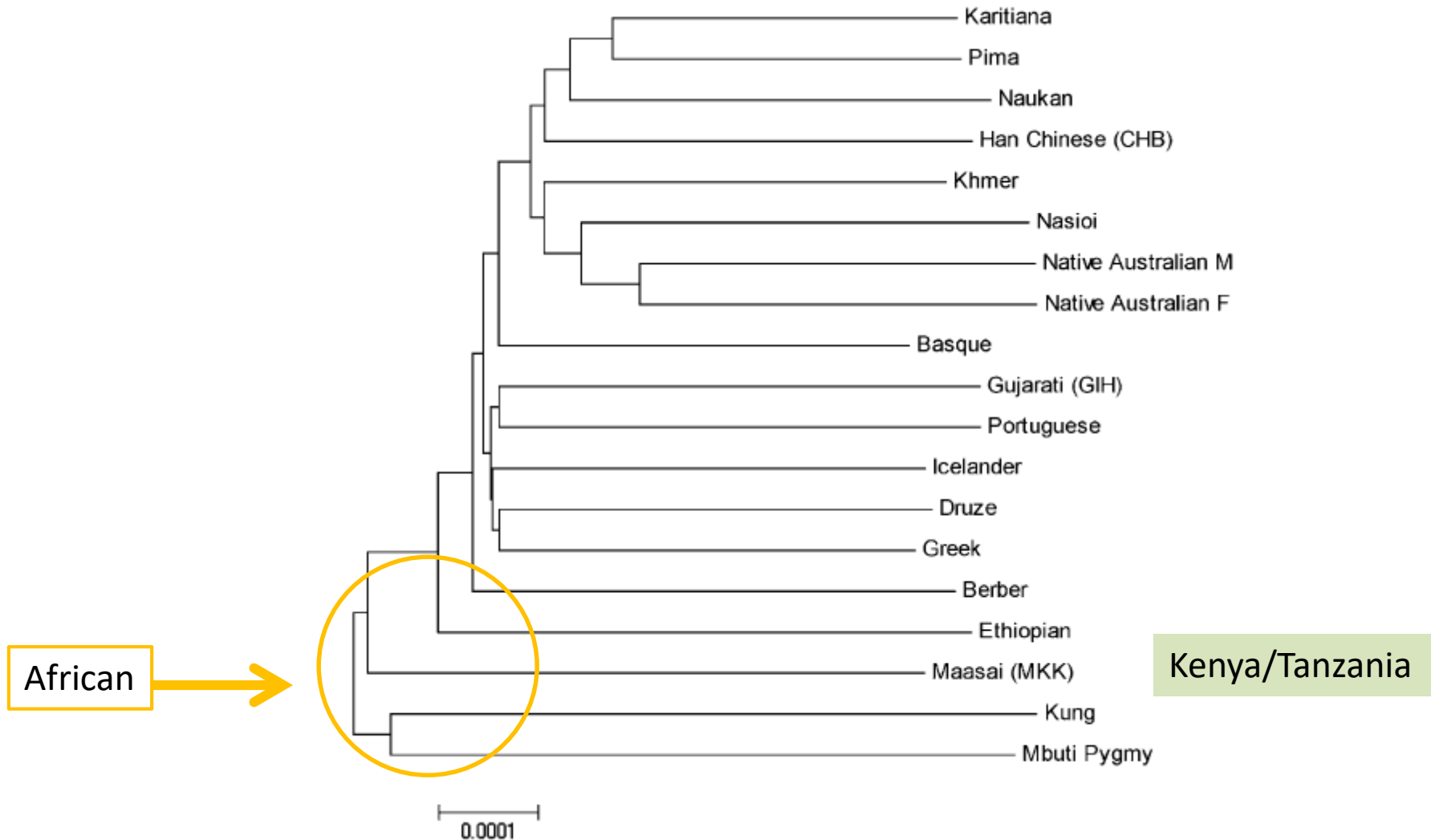
# Data analysis #2: Population divergence and split times

Pairwise differences then phylogenetic tree (neighbor joining)



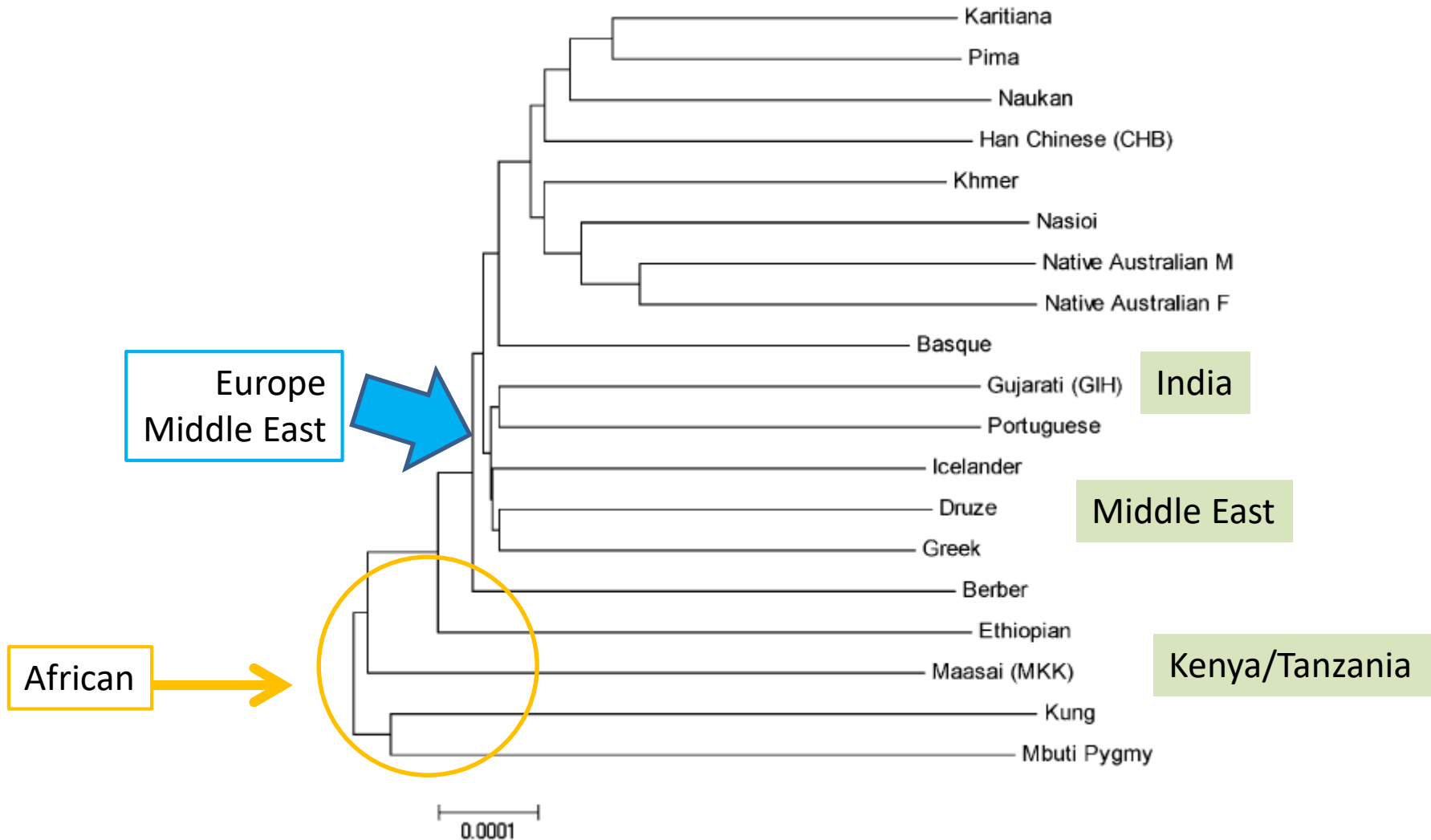
# Data analysis #2: Population divergence and split times

Pairwise differences then phylogenetic tree (neighbor joining)



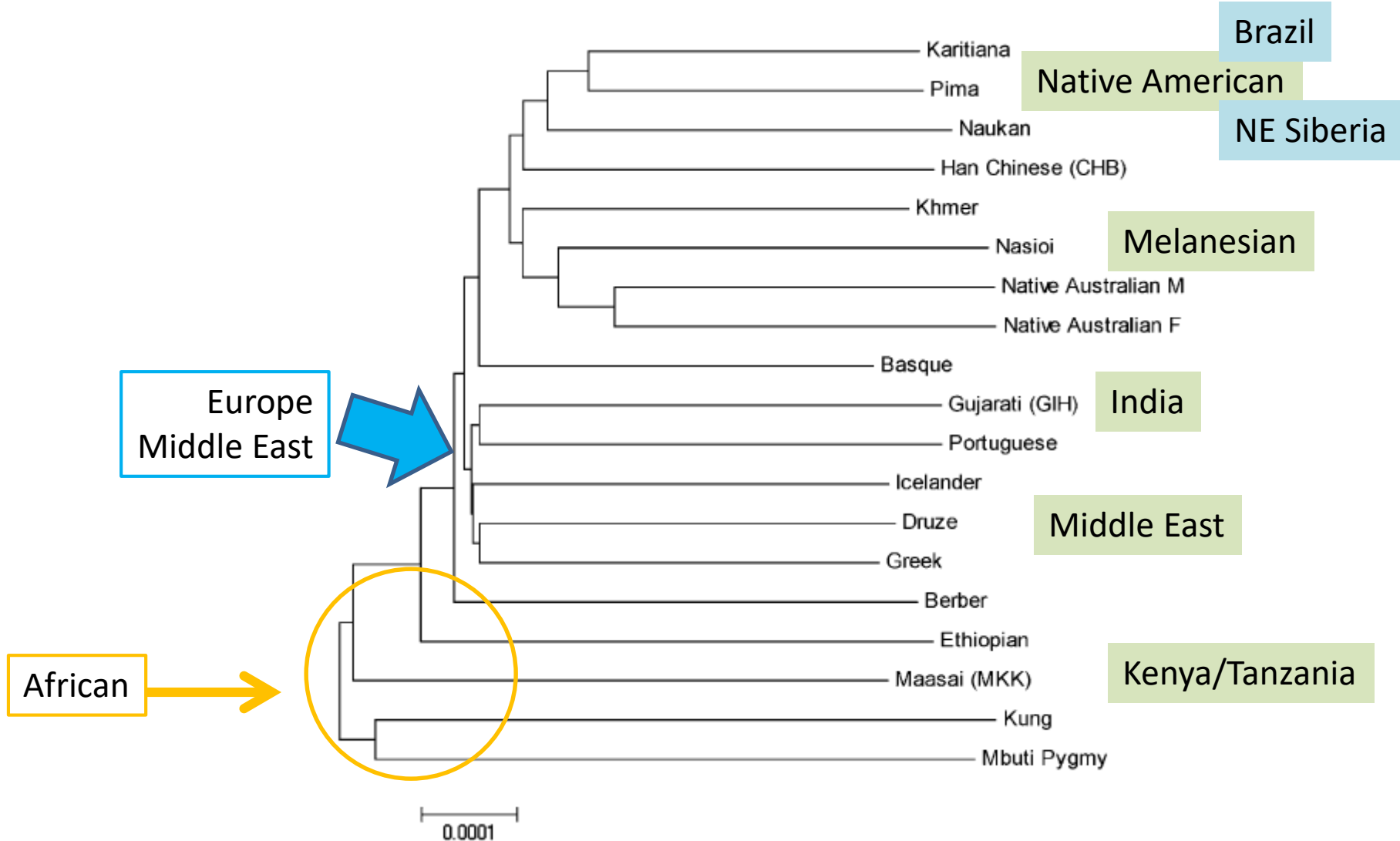
# Data analysis #2: Population divergence and split times

Pairwise differences then phylogenetic tree (neighbor joining)



# Data analysis #2: Population divergence and split times

Pairwise differences then phylogenetic tree (neighbor joining)



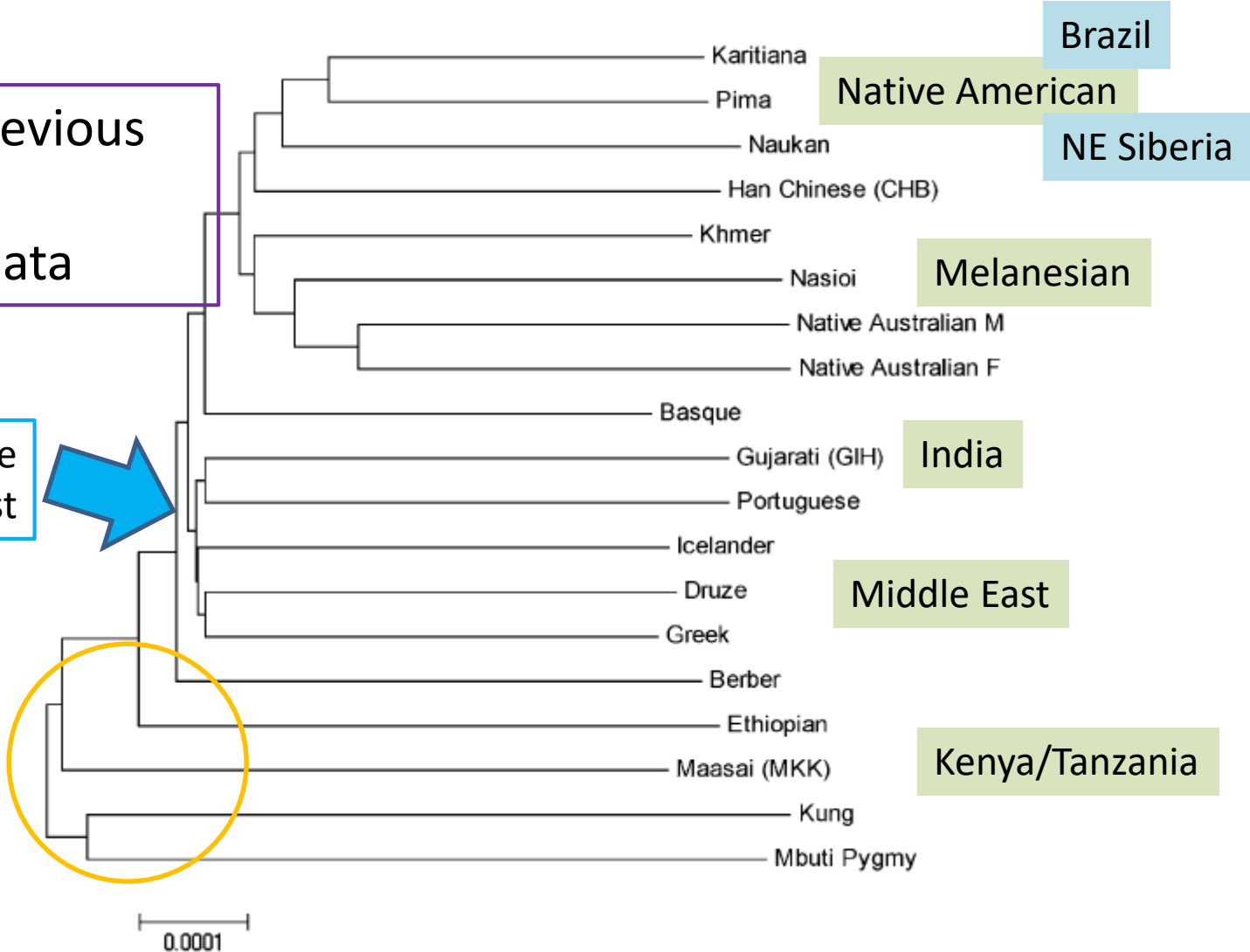
# Data analysis #2: Population divergence and split times

Pairwise differences then phylogenetic tree (neighbor joining)

Recapitulates previous genetic and archaeological data

Europe  
Middle East

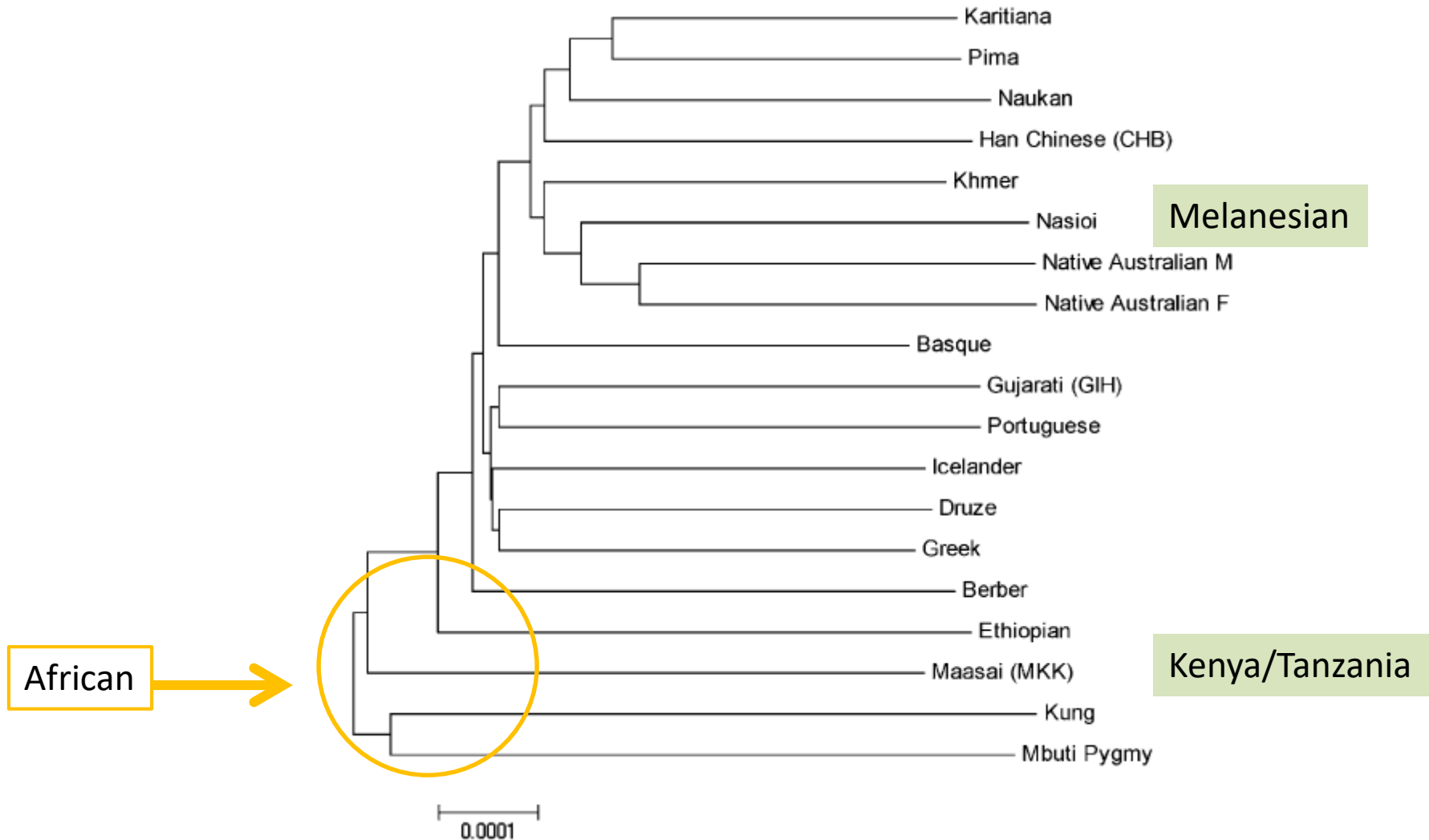
African





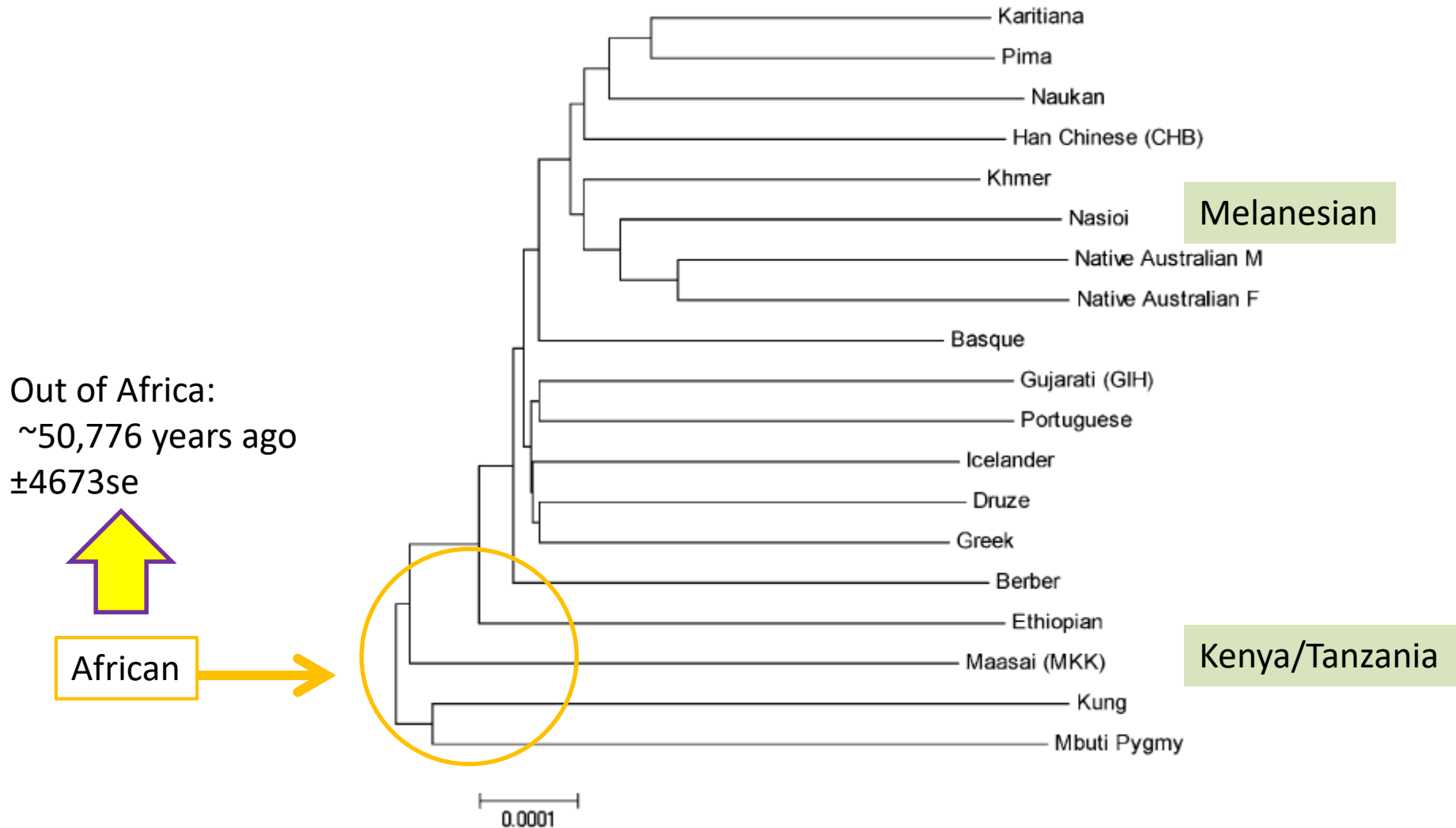
# Data analysis #2: Population divergence and split times

Pairwise differences then phylogenetic tree (neighbor joining)



# Data analysis #2: Population divergence and split times

Pairwise differences then phylogenetic tree (neighbor joining)





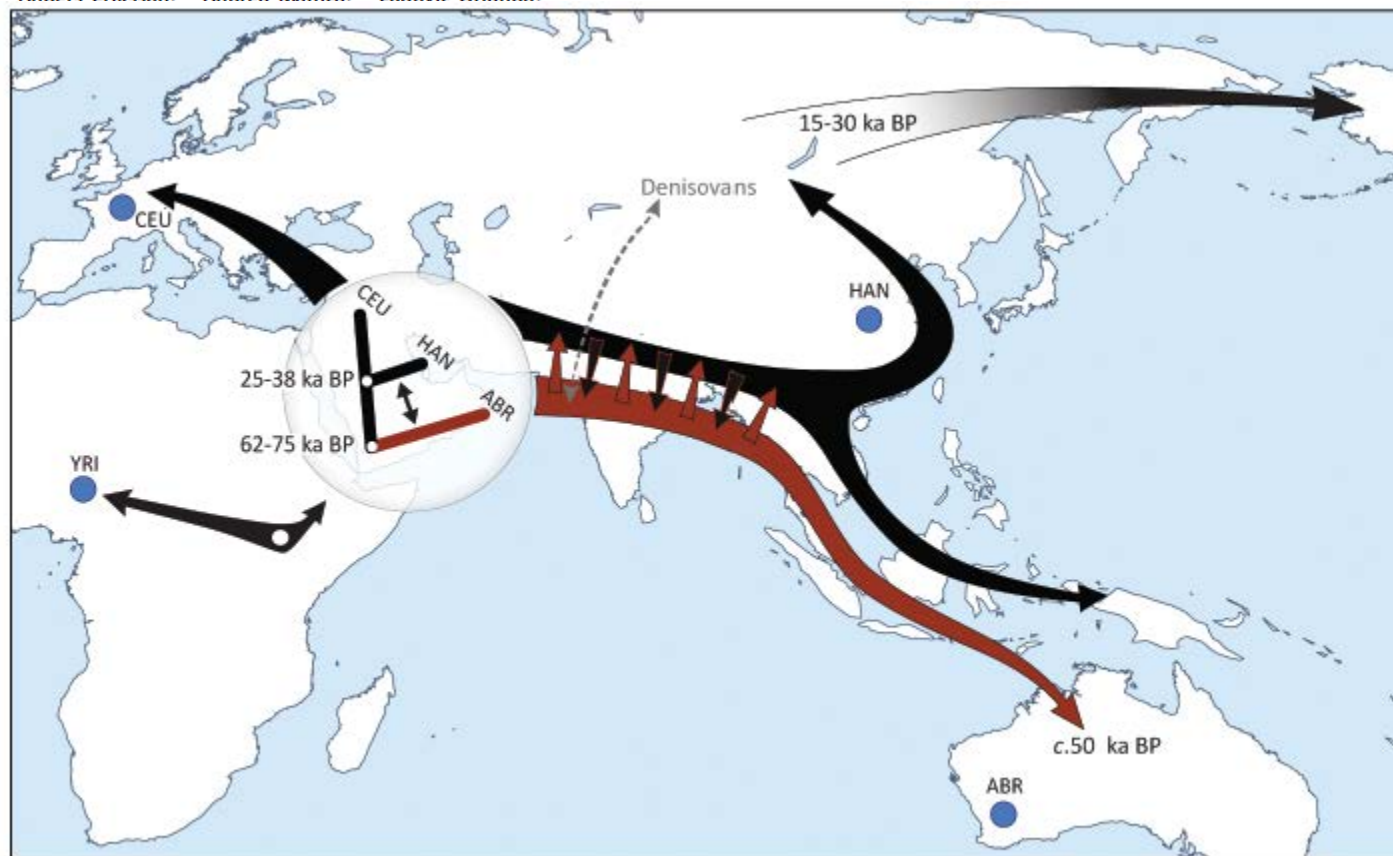
# An Aboriginal Australian Genome Reveals Separate Human Dispersals into Asia

Morten Rasmussen,<sup>1,2\*</sup> Xiaosen Guo,<sup>2,3\*</sup> Yong Wang,<sup>4\*</sup> Kirk E. Lohmueller,<sup>4\*</sup> Simon Rasmussen,<sup>5</sup> Anders Albrechtsen,<sup>6</sup> Line Skotte,<sup>6</sup> Stinus Lindgreen,<sup>1,6</sup> Mait Metspalu,<sup>7</sup> Thibaut Jombart,<sup>8</sup> Toomas Kivisild,<sup>9</sup> Weiwei Zhai,<sup>10</sup> Anders Eriksson,<sup>11</sup> Andrea Manica,<sup>11</sup> Ludovic Orlando,<sup>1</sup> Francisco M. De La Vega,<sup>12</sup> Silvana Tridico,<sup>13</sup> Ene Metspalu,<sup>7</sup> Kasper Nielsen,<sup>5</sup> María C. Ávila-Arcos,<sup>1</sup> J. Víctor Moreno-Mayar,<sup>1,14</sup> Craig Muller,<sup>15</sup> Joe Dortch,<sup>16</sup> M. Thomas P. Gilbert,<sup>1,2</sup> Ole Lund,<sup>5</sup> Agata Wesolowska,<sup>5</sup> Monika Karmin,<sup>7</sup> Lucy A. Weinert,<sup>8</sup> Bo Wang,<sup>3</sup> Jun Li,<sup>3</sup> Shuaishuai Tai,<sup>3</sup> Fei Xiao,<sup>3</sup> Tsunehiko Hanihara,<sup>17</sup> George van Driem,<sup>18</sup> Aashish R. Jha,<sup>19</sup> François-Xavier Ricaut,<sup>20</sup> Peter de Knijff,<sup>21</sup> Andrea B. Migliano,<sup>9,22</sup> Irene Gallego Romero,<sup>19</sup> Karsten Kristiansen,<sup>2,3,6</sup> David M. Lambert,<sup>23</sup> Søren Brunak,<sup>5,24</sup> Peter Forster,<sup>25,26</sup> Bernd Brinkmann,<sup>26</sup> Olaf Nehlich,<sup>27</sup> Michael Bunce,<sup>13</sup> Michael Richards,<sup>27,28</sup> Ramneek Gupta,<sup>5</sup> Carlos D. Bustamante,<sup>12</sup> Anders Krogh,<sup>1,6</sup> Robert A. Foley,<sup>9</sup> Marta M. Lahr,<sup>9</sup> François Balloux,<sup>8</sup> Thomas Sicheritz-Pontén,<sup>5,29</sup> Richard Villems,<sup>7,30</sup> Rasmus Nielsen,<sup>4,6†</sup> Jun Wang,<sup>2,3,6,31†</sup> Eske Willerslev,<sup>1,2†</sup>

# An Aboriginal Australian Genome Reveals Separate Human Dispersals into Asia

Morten Rasmussen,<sup>1,2\*</sup> Xiaosen Guo,<sup>2,3\*</sup> Yong Wang,<sup>4\*</sup> Kirk E. Lohmueller,<sup>4\*</sup> Simon Rasmussen,<sup>5</sup> Anders Albrechtsen,<sup>6</sup> Line Skotte,<sup>6</sup> Stinus Lindgreen,<sup>1,6</sup> Mait Metspalu,<sup>7</sup> Thibaut Jombart,<sup>8</sup> Toomas Kivisild,<sup>9</sup> Weiwei Zhai,<sup>10</sup> Anders Eriksson,<sup>11</sup> Andrea Manica,<sup>11</sup> Ludovic Orlando,<sup>1</sup>

Francisco M. De La Vega,<sup>12</sup> J. Víctor Moreno-Mayar,<sup>1,14</sup> Agata Wesolowska,<sup>5</sup> Monik Fei Xiao,<sup>3</sup> Tsunehiko Hanih Peter de Knijff,<sup>21</sup> Andrea B David M. Lambert,<sup>23</sup> Søren Michael Bunce,<sup>13</sup> Michael I Anders Krogh,<sup>1,6</sup> Robert A. Richard Villemis,<sup>7,30</sup> Rasmu



**Fig. 2.** Reconstruction of early spread of modern humans outside Africa. The tree shows the divergence of

Low coverage approaches

# PoolSeq

Eponymous:

pool samples then sequence

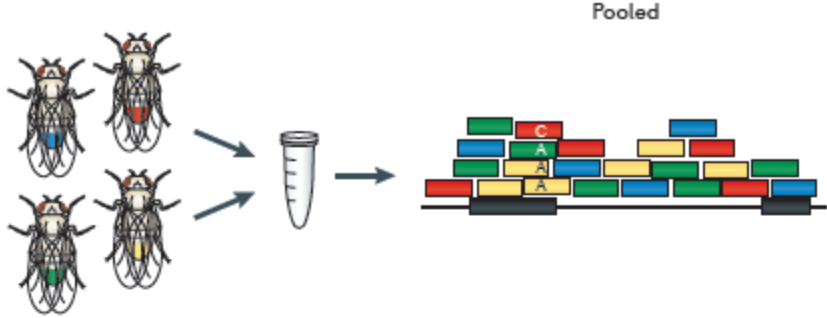
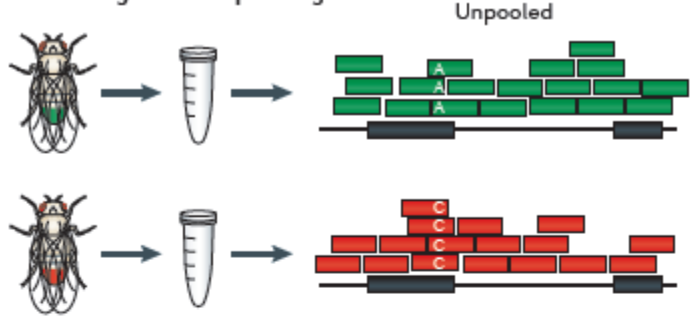
often ~1x depth for each sample

# PoolSeq

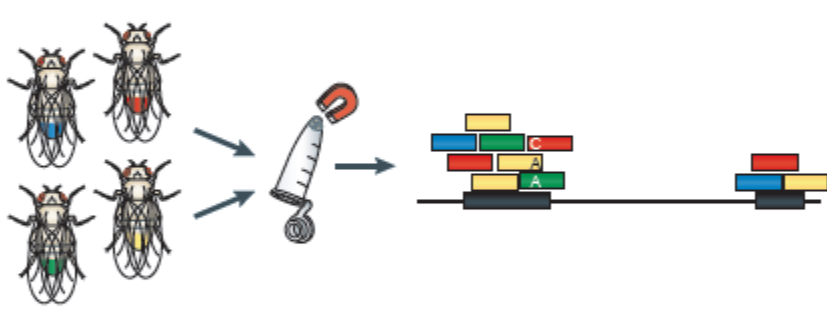
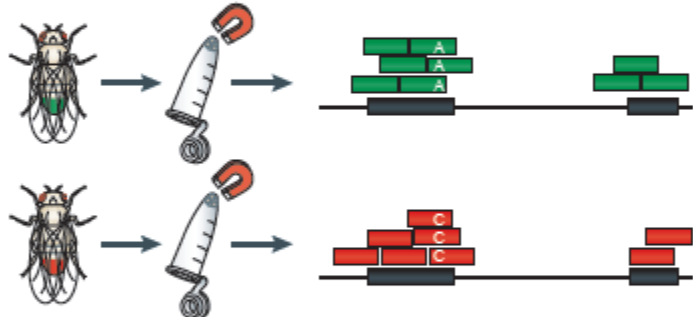
UNPOOLED

POOLED

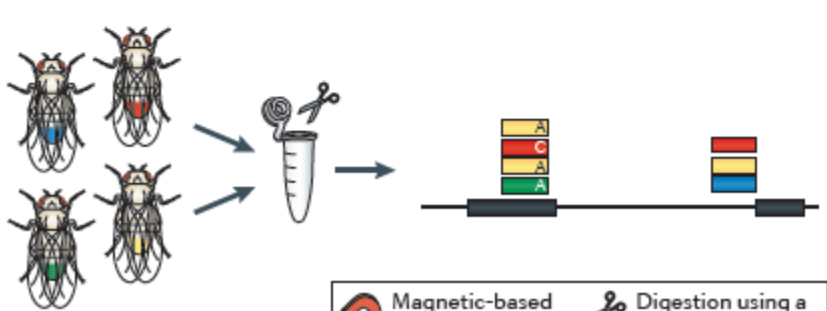
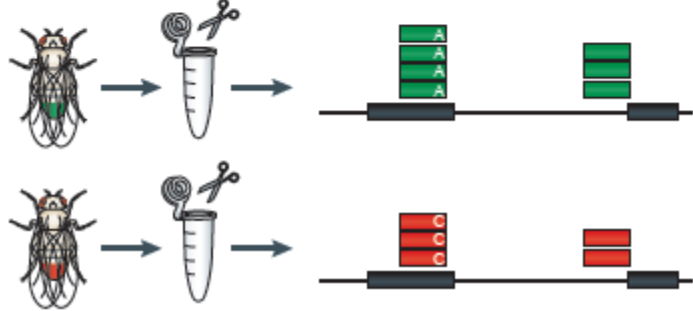
**a** Whole-genome sequencing



**b** Exome sequencing



**c** RAD-seq



 Magnetic-based exome capture       Digestion using a restriction enzyme



# PoolSeq – applications

- Identification of functionally diverged alleles in mapping experiments
- Characterization of alleles involved in the adaptation of populations to their local environment
- Analysis of temporal allele frequency trajectories

# PoolSeq – advantages

# PoolSeq – advantages

Better estimate of allele frequencies (polymorphism data)

Individual → often redundant coverage (e.g., 10x)

Pool → not every individual, BUT basically only 1x

# PoolSeq – advantages

Better estimate of allele frequencies (polymorphism data)

Individual → often redundant coverage (e.g., 10x)

Pool → not every individual, BUT basically only 1x

Protocol easier

[ No individualization ]

# PoolSeq – advantages

Better estimate of allele frequencies (polymorphism data)

Individual → often redundant coverage (e.g., 10x)

Pool → not every individual, BUT basically only 1x

Protocol easier

[ No individualization ]

Cheaper

[ Barcoding/library construction entails a big cost ]

# PoolSeq – advantages

Better estimate of allele frequencies (polymorphism data)

Individual → often redundant coverage (e.g., 10x)

Pool → not every individual, BUT basically only 1x

Protocol easier

[ No individualization ]

Cheaper

[ Barcoding/library construction entails a big cost ]

Ideally many individuals

>100

>50 okay...

# PoolSeq – advantages

Better estimate of allele frequencies (polymorphism data)

Individual → often redundant coverage (e.g., 10x)

Pool → not every individual, BUT basically only 1x

Protocol easier

[ No individualization ]

Cheaper

[ Barcoding/library construction entails a big cost ]

Ideally many individuals

>100

>50 okay...

Unbiased (genome wide)

Compared to RADseq or exome seq

# PoolSeq – concerns



# PoolSeq – concerns

SNP call versus Illumina error rate

Need enough depth, but assaying superficially

# PoolSeq – concerns

SNP call versus Illumina error rate

Need enough depth, but assaying superficially

Loss of low frequency alleles

But such alleles are informative for many pop genetics analyses:

hitchhiking

quantification of positive and purifying selection

demographic parameters

# PoolSeq – concerns

SNP call versus Illumina error rate

Need enough depth, but assaying superficially

Loss of low frequency alleles

But such alleles are informative for many pop genetics analyses:

hitchhiking

quantification of positive and purifying selection

demographic parameters

Bad for:

Heterozygosity, CNVs, inversions and transposable elements,

# PoolSeq – concerns

SNP call versus Illumina error rate

Need enough depth, but assaying superficially

Loss of low frequency alleles

But such alleles are informative for many pop genetics analyses:

hitchhiking

quantification of positive and purifying selection

demographic parameters

Bad for:

Heterozygosity, CNVs, inversions and transposable elements,

No haplotype info, no linkage disequilibrium (LD) information

e.g., evidence for African lactose tolerance

# PoolSeq – concerns

SNP call versus Illumina error rate

Need enough depth, but assaying superficially

Loss of low frequency alleles

But such alleles are informative for many pop genetics analyses:

hitchhiking

quantification of positive and purifying selection

demographic parameters

Bad for:

Heterozygosity, CNVs, inversions and transposable elements,

No haplotype info, no linkage disequilibrium (LD) information

e.g., evidence for African lactose tolerance

Variance in pooling (Less of a problem with greater pooling)

# PoolSeq – some software

<i>Population genetics</i>		
<a href="#">PoPoolation</a>	Estimates variation within populations	39
<a href="#">PoPoolation2</a>	Estimates differentiation between multiple populations	132
<a href="#">Pool-HMM</a>	Detects selective sweeps from the allele frequency spectrum using a hidden Markov model	133
<a href="#">npstat</a>	Computes a wide range of population genetic estimators; may be used in conjunction with an external SNP caller; every contig needs to be analysed separately	134
<a href="#">Stacks</a>	Developed for population genomics with RAD-seq; may also be used with pooled RAD-seq data	135
<a href="#">Bayenv2</a>	Estimates differentiation between populations	79
<a href="#">SelEstim</a>	Detects and measures selection	136
<a href="#">KimTree</a>	Infers population histories	137
<i>Haplotype information</i>		
<a href="#">harp</a>	Estimates frequencies of known haplotypes using read counts; supports a sliding window approach	107
<a href="#">PoolHap</a>	Estimates frequencies of known haplotypes using a regression on allele frequencies	106
<a href="#">eALPS</a>	Estimates the abundance of individuals in pools given the genotypes of at least some individuals	109
<a href="#">LDx</a>	Estimates linkage disequilibrium between pairs of SNPs spanned by single-end or paired-end reads	138
<i>Forward genetic screens</i>		
<a href="#">SHOREmap</a>	Identifies causative recessive variants from a large pool of recombinants that have the recessive genotype	44
<a href="#">CloudMap</a>	A cloud-based pipeline for localizing mutations	139
<a href="#">MULTIPOOL</a>	Identifies candidate loci from bulk segregant analysis in which progeny are grouped by phenotype and sequenced as pools	140
<a href="#">Fishyskeleton</a>	Detects mutations in zebrafish from a pool of mutant F <sub>2</sub> fish	141
<a href="#">NGM</a>	A web-based tool for localizing mutations from a small pool of F <sub>2</sub> population	142
<a href="#">SNPtrack</a>	A web-based tool for localizing mutations using a hidden Markov model	143

... there are others

# PoolSeq – some applications

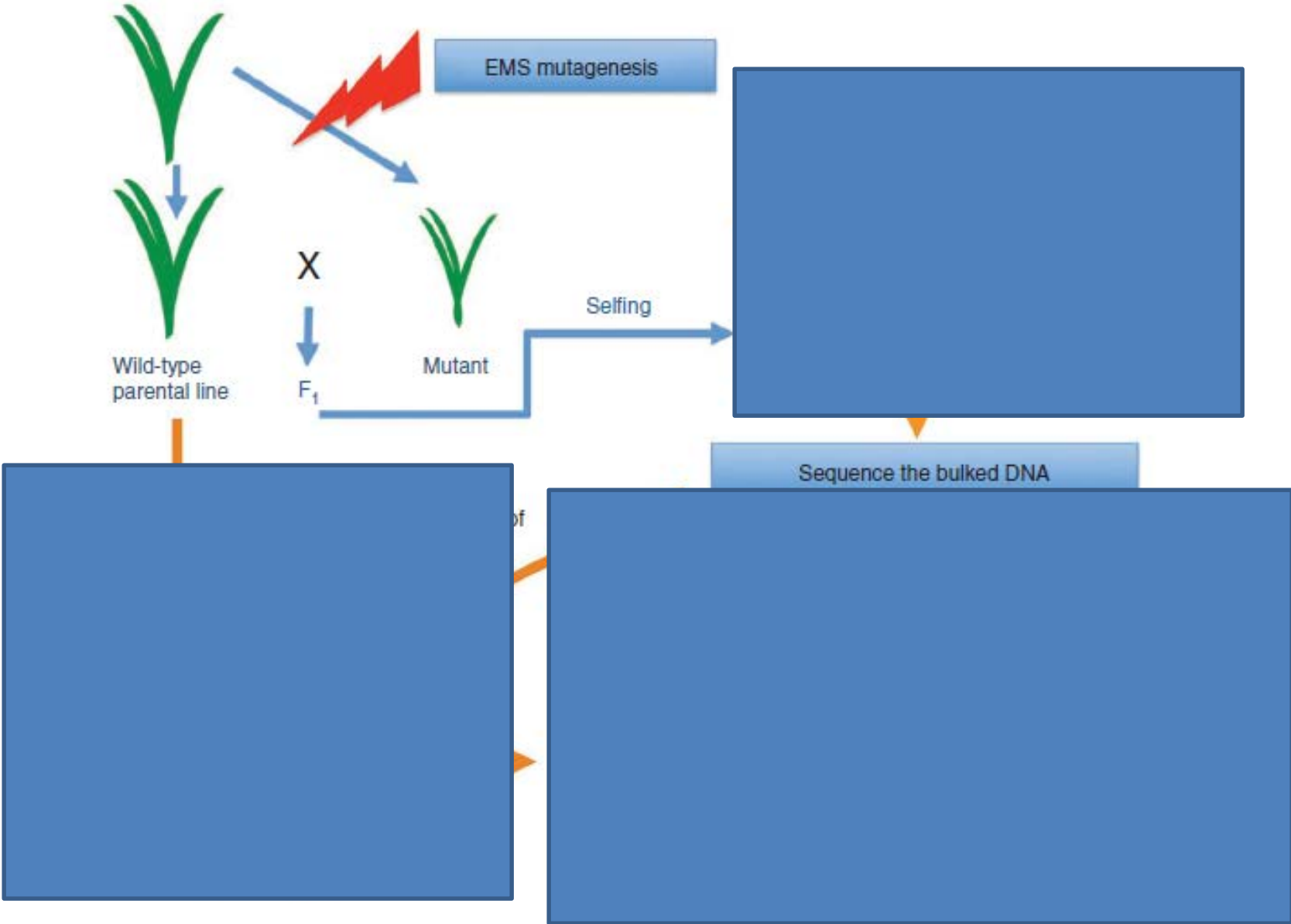
- Induced mutation
- QTL
- GWAS
- Evolve and re-sequence
- Domestication gene identification
- Recombination rate analysis
- ...





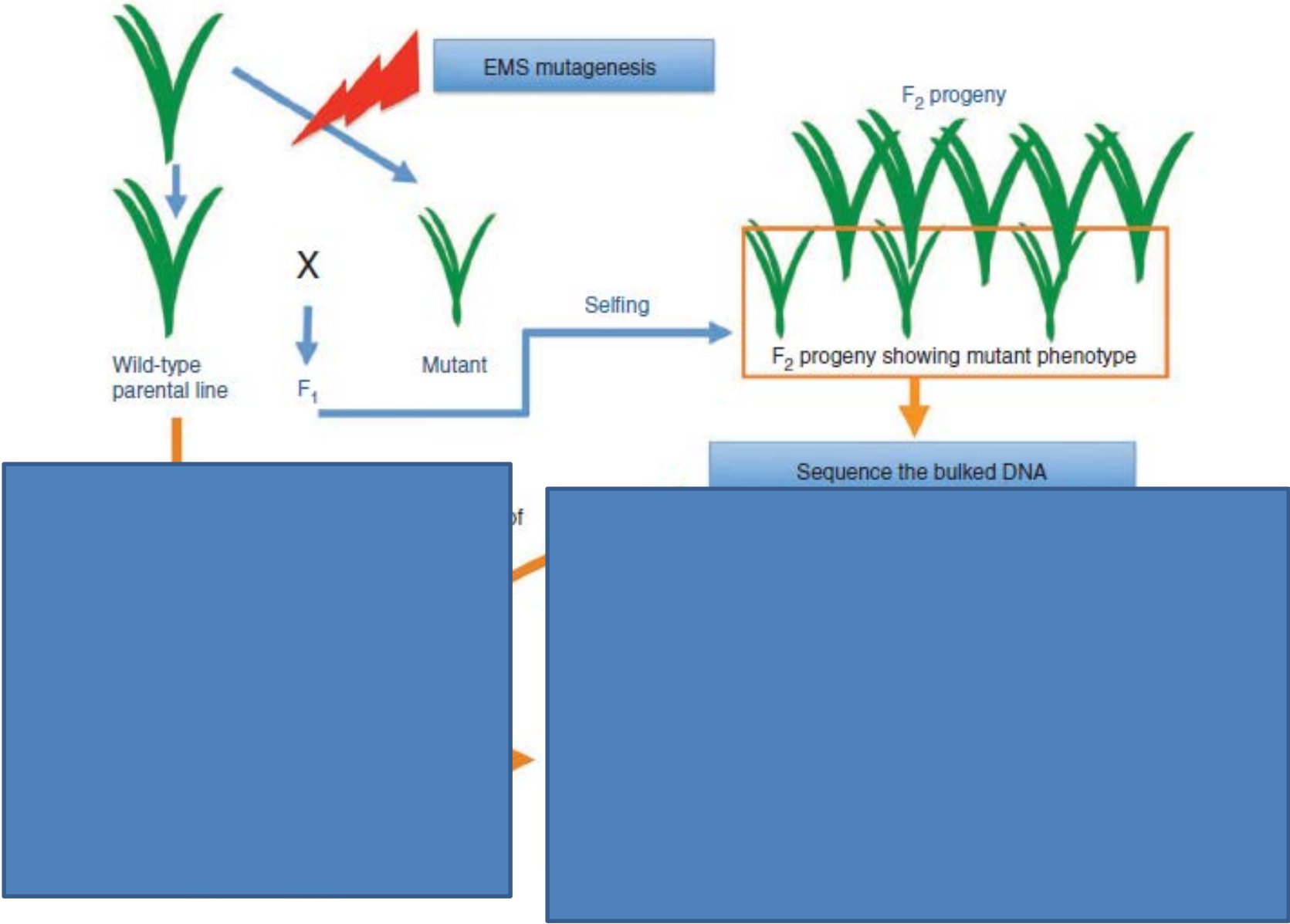
# PoolSeq – MutMap for induced mutation

Induced mutation



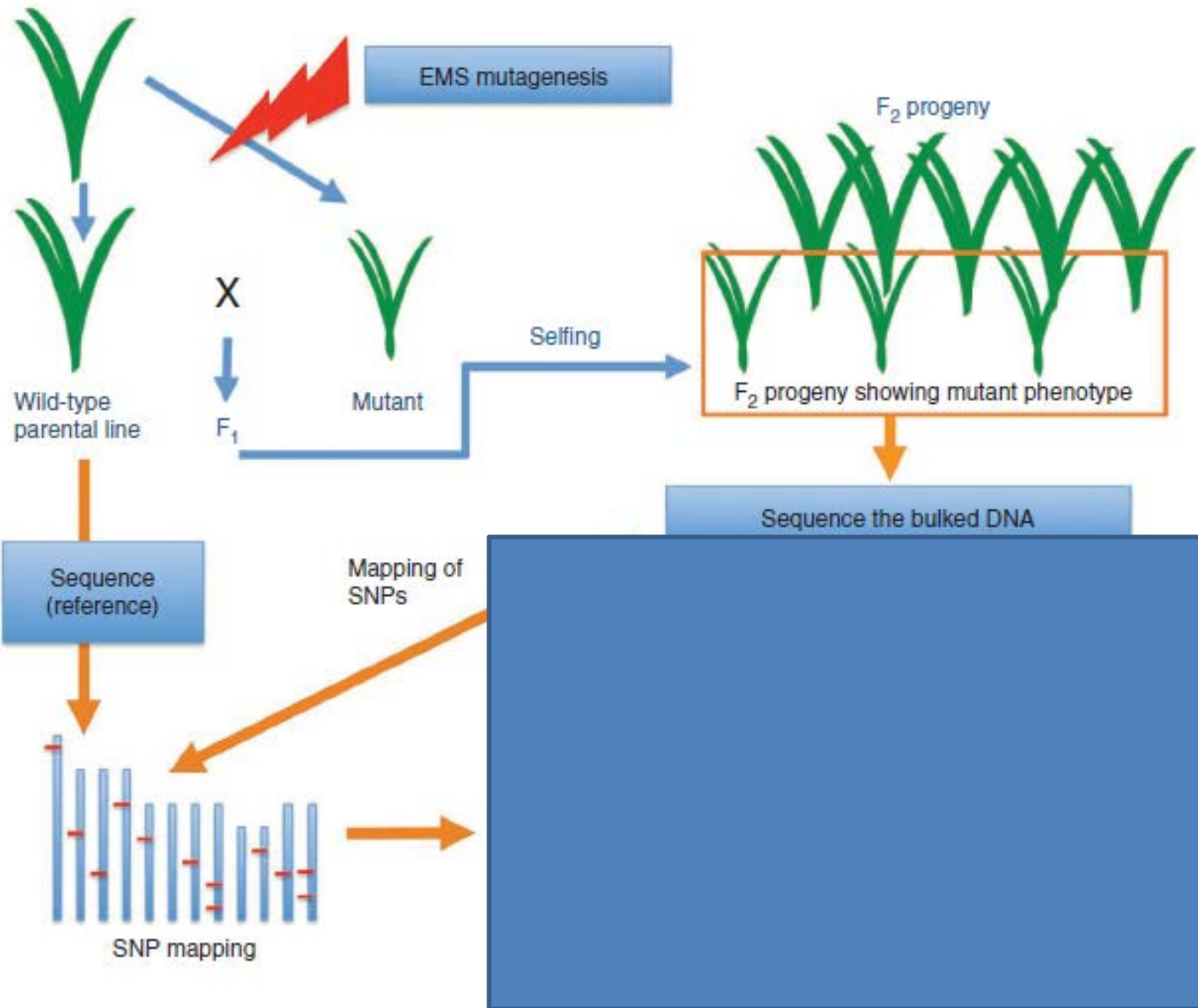
# PoolSeq – MutMap for induced mutation

Induced mutation



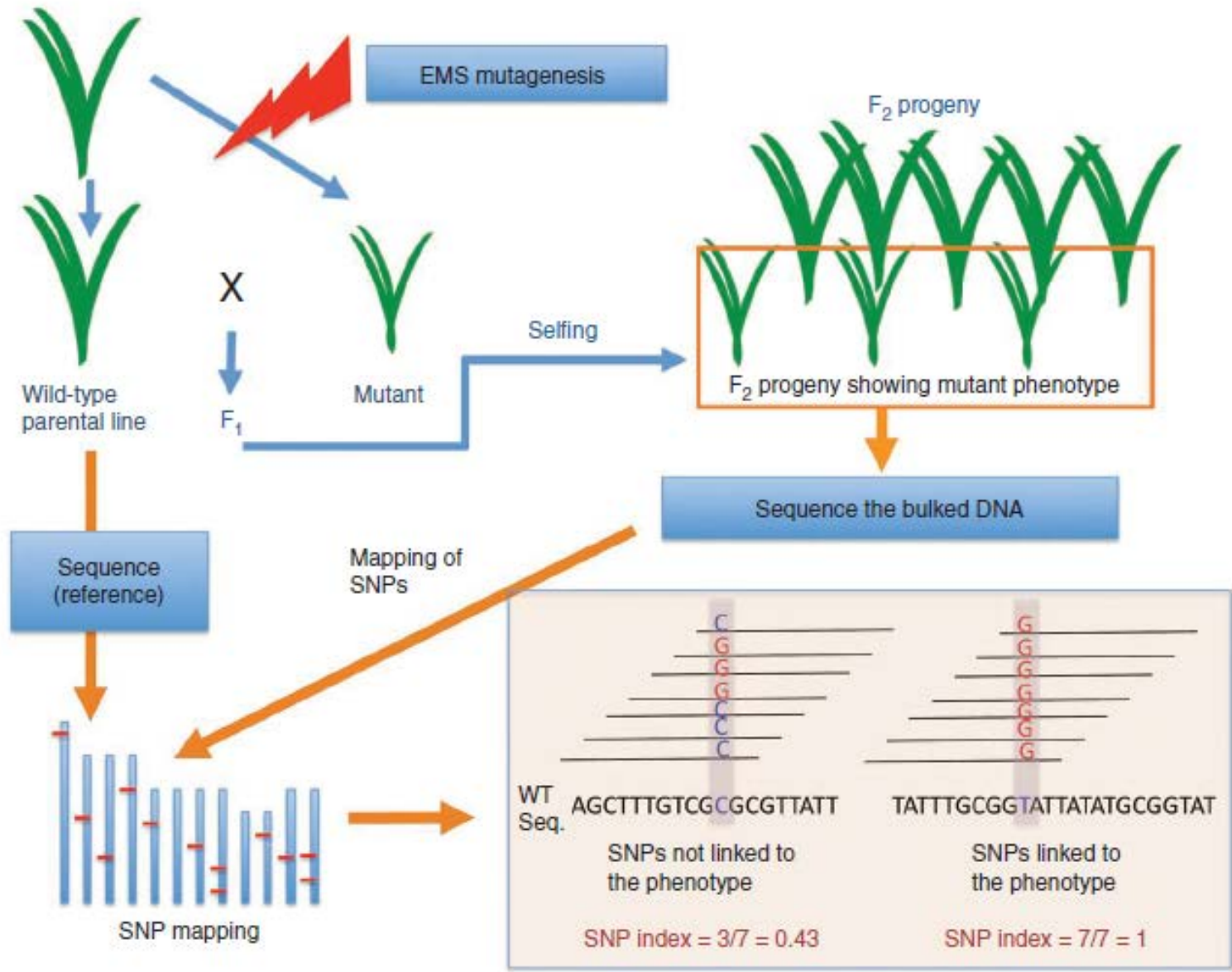
# PoolSeq – MutMap for induced mutation

Induced mutation



# PoolSeq – MutMap for induced mutation

Induced mutation



A **recessive** mutant trait due to a  
**single** gene mutation

A **recessive** mutant trait due to a  
**single** gene mutation

Mapping mutation within same strain, so less noise

**SNP index:** The ratio of the number of reads containing a mutant SNP to the total number of reads containing the site of the SNP.

SNP index:

100% for a SNP tightly linked to the causal gene

50% for a SNP unlinked to the mutation

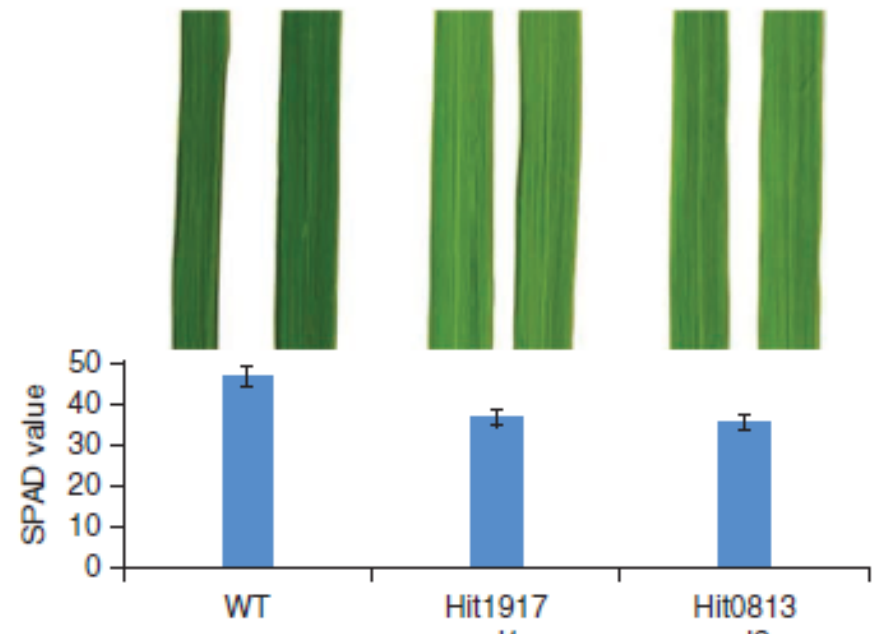
>50% for a mutant SNP loosely linked to the causal mutation

<50% for the wild type allele

# PoolSeq – MutMap for induced mutation

Induced mutation

a

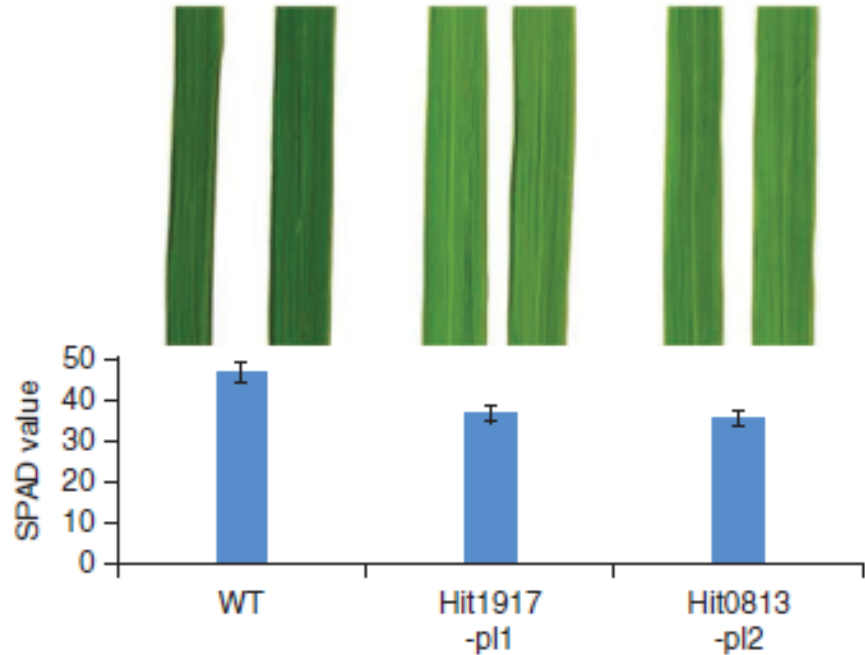




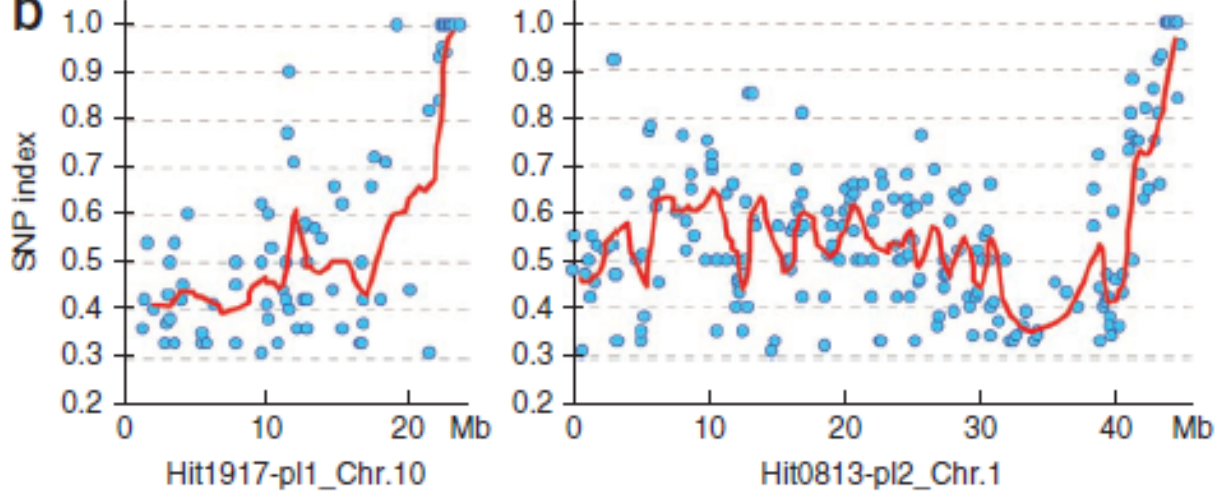
# PoolSeq – MutMap for induced mutation

Induced mutation

**a**

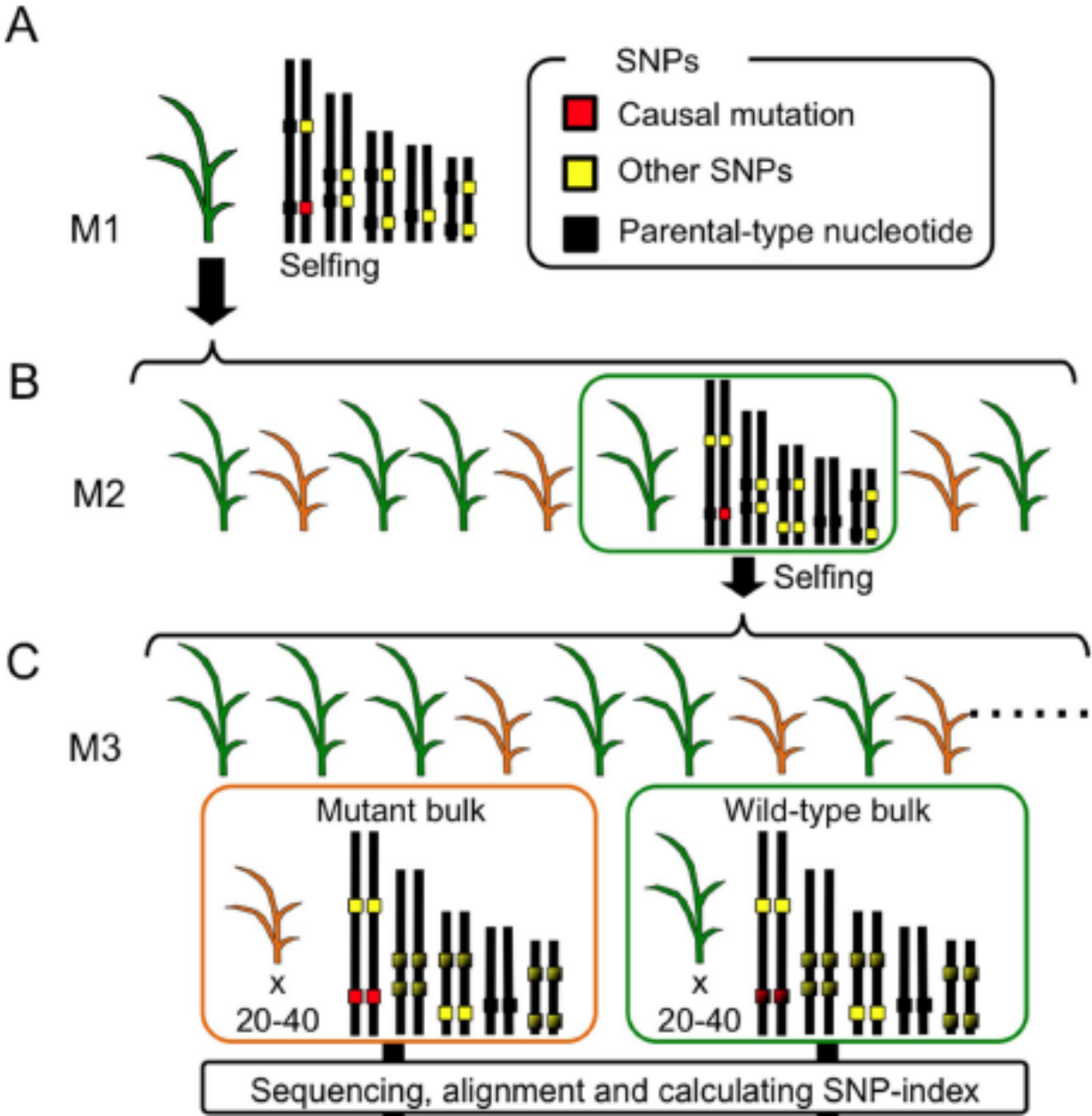


**b**



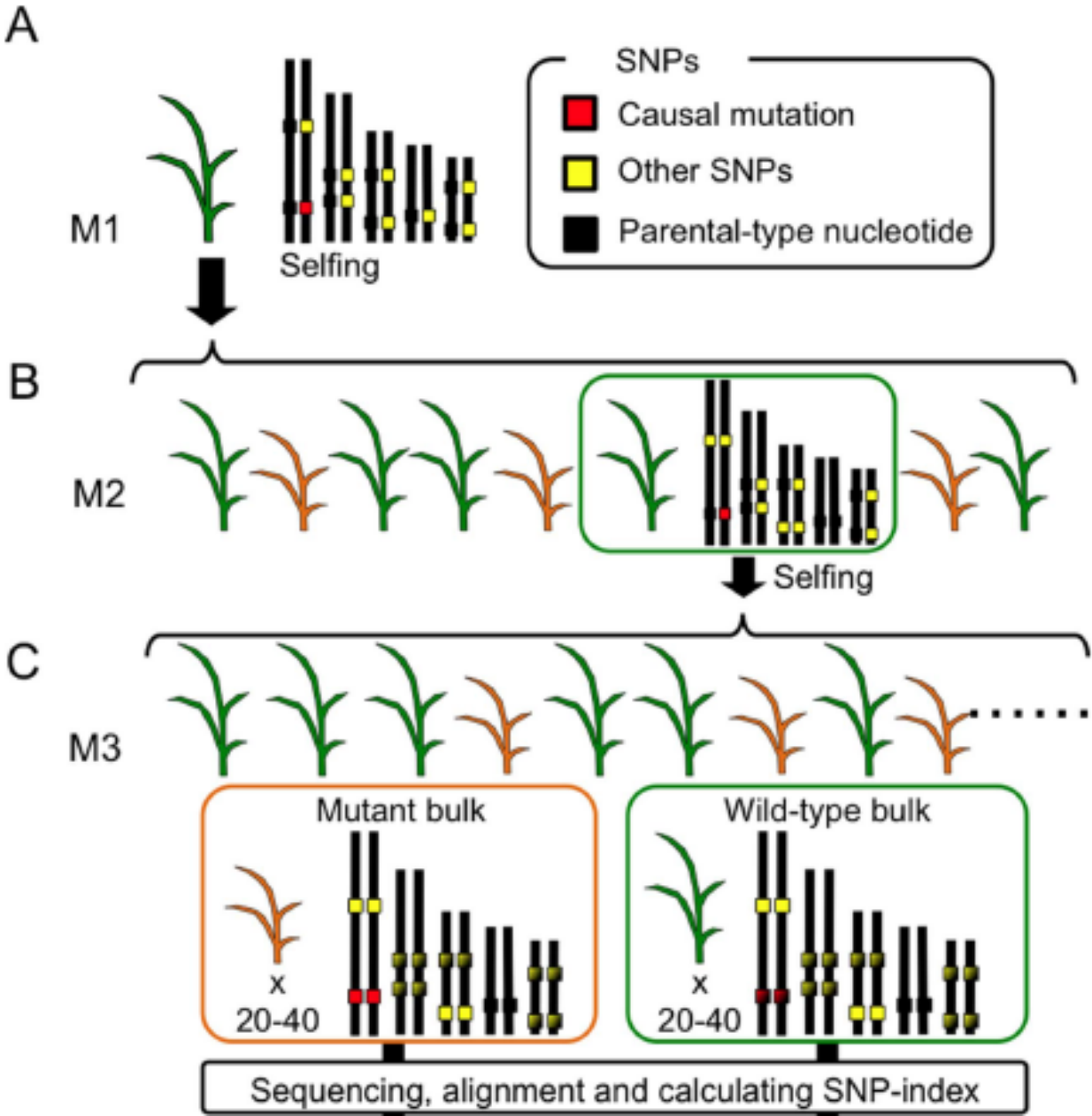
# Variant: MutMap+ → No need to cross

Induced mutation

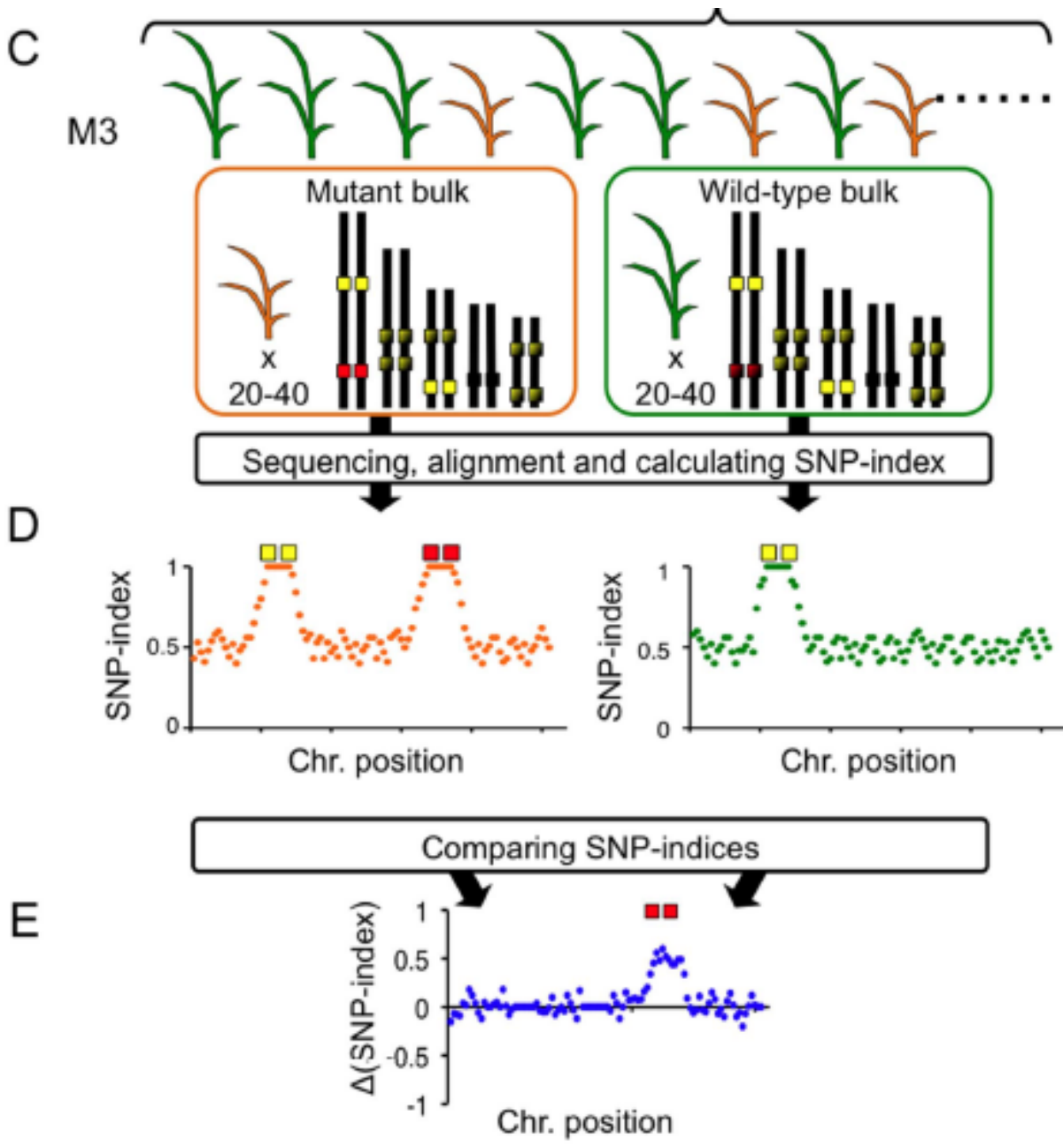


# Variant: MutMap+ → No need to cross

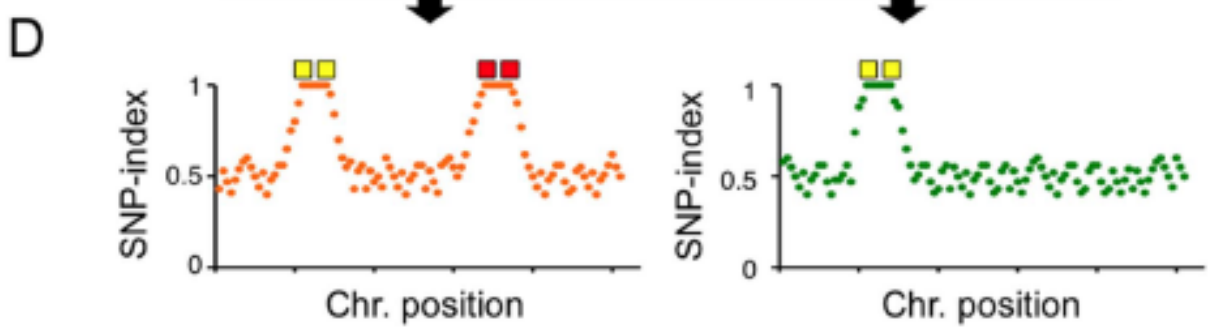
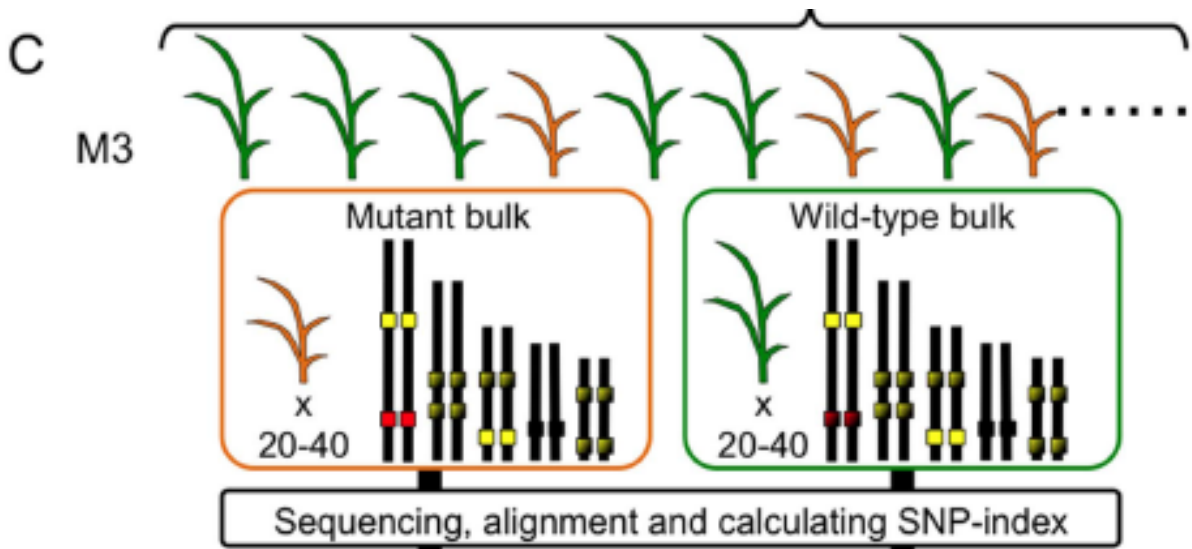
Induced mutation



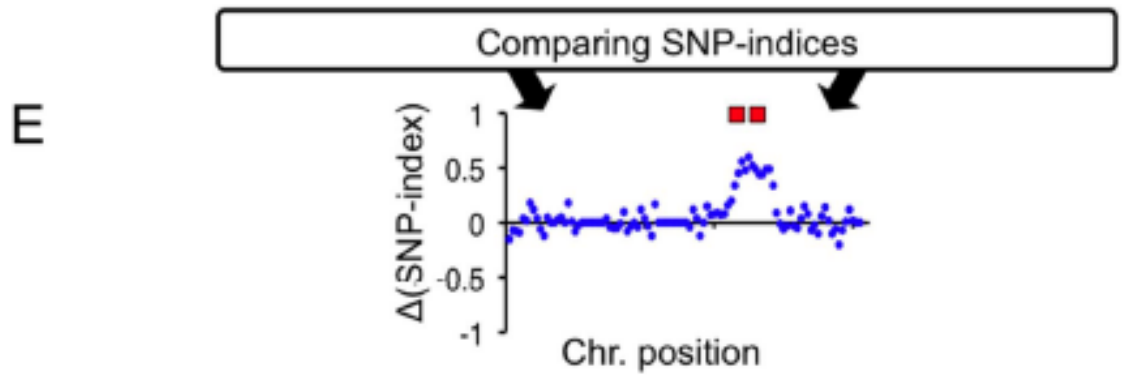
Some wt → all wt ignore  
wt → wt + mut keep

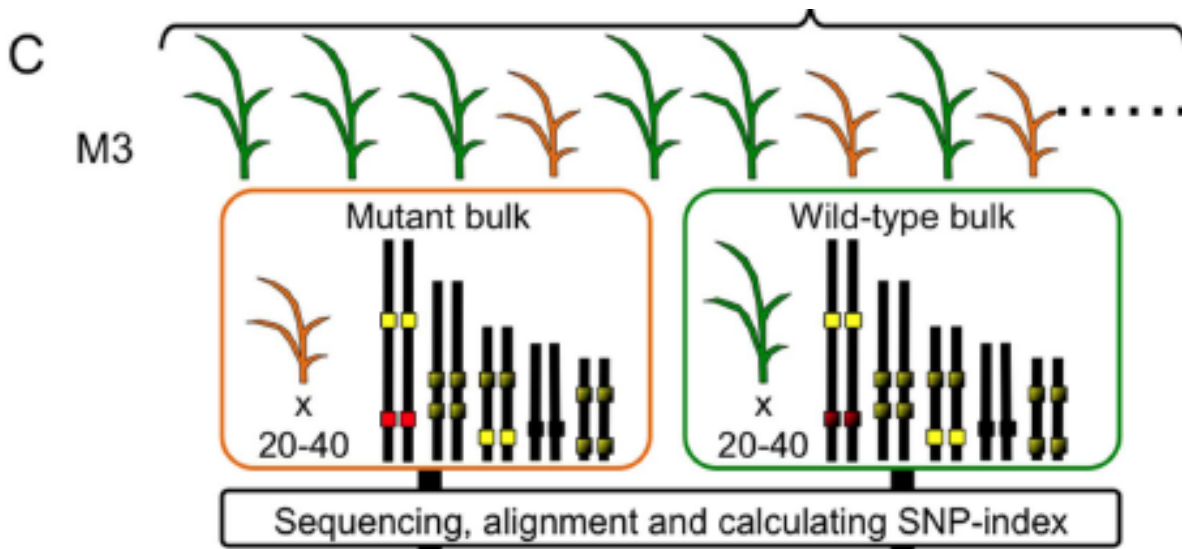


Induced mutation

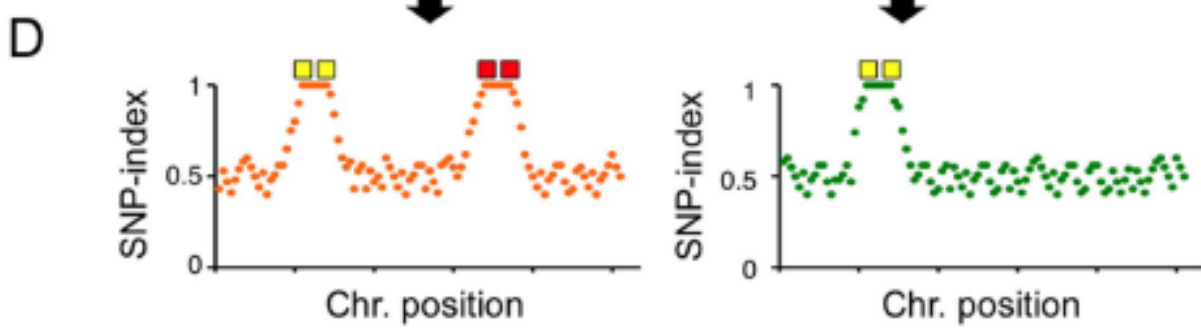


Some SNPs fixed due to chance or drift



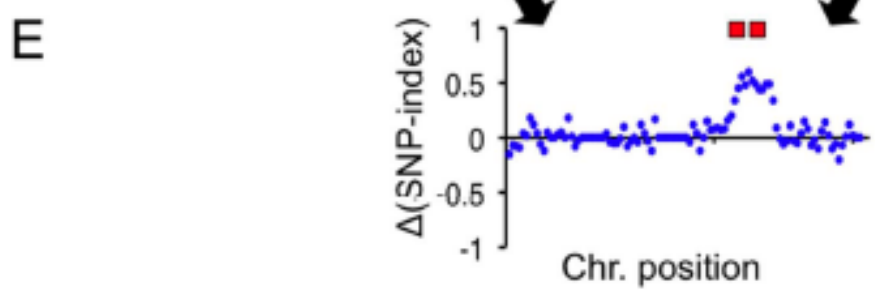


Induced mutation



Some SNPs fixed due to chance or drift

Comparing SNP-indices



$\Delta$ SNP index

# NGS and Poolseq

# NGS and Poolseq

OLD:

many graduate students and many years



# NGS and Poolseq

OLD:

many graduate students and many years

NEW:

single graduate student and fewer years

# NGS and Poolseq

OLD:

many graduate students and many years

NEW:

single graduate student and fewer years

PI's and peons ostensibly happier

Poolseq

xQTL mapping

QTL-seq

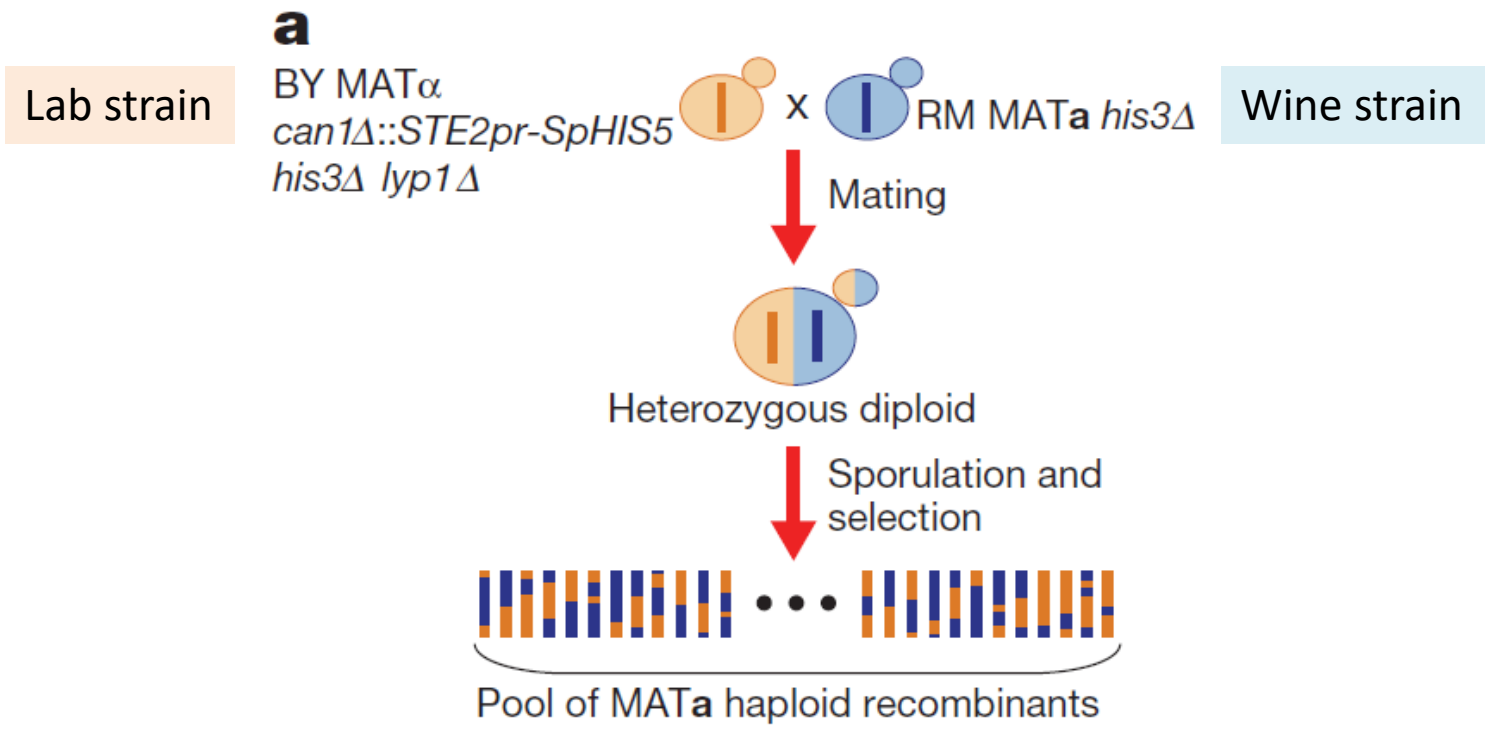
Highly complex traits → Multiple loci  
Often each with small effect

Need large populations for QTL mapping

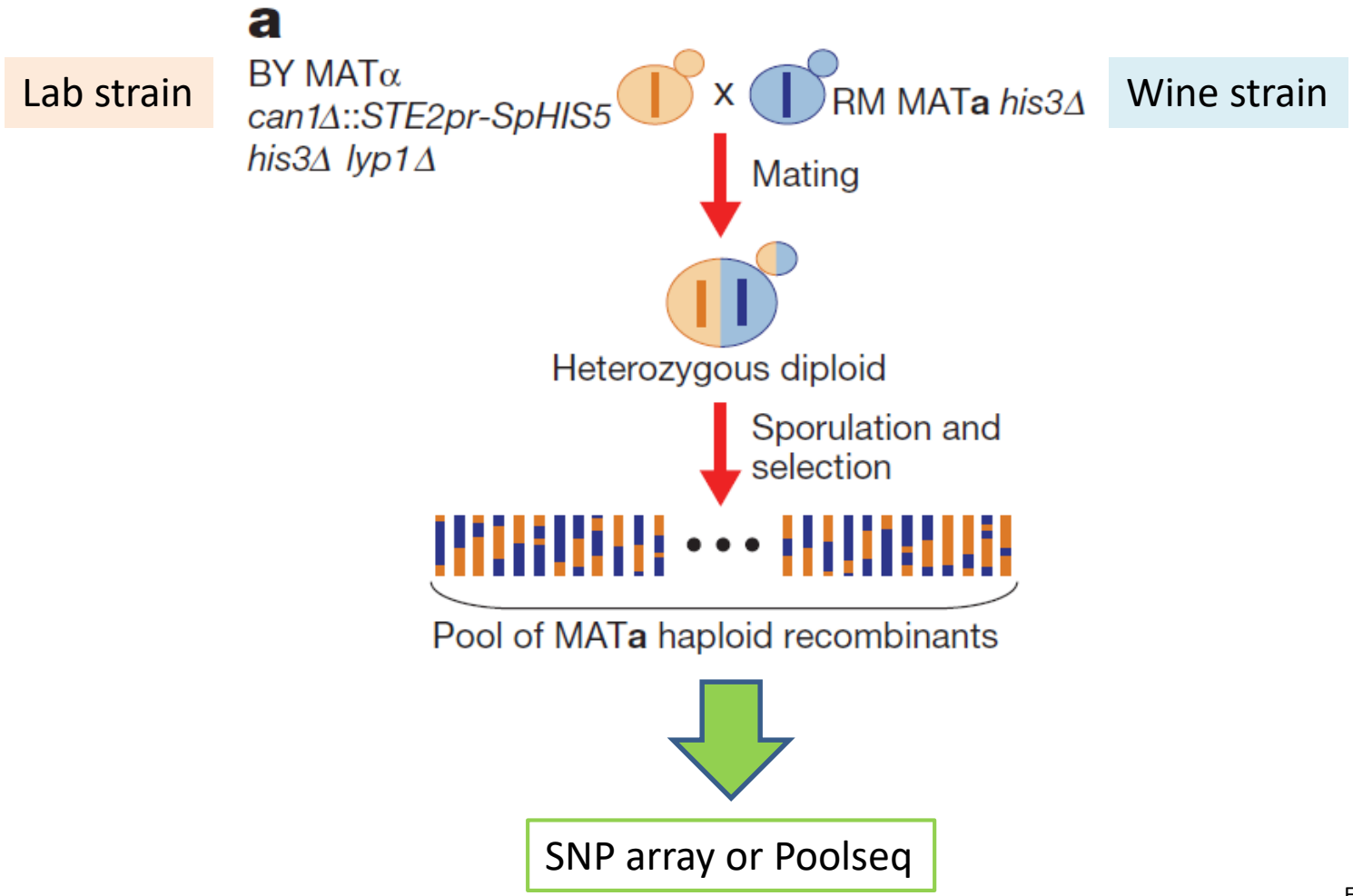
## **Dissection of genetically complex traits with extremely large pools of yeast segregants**

Ian M. Ehrenreich<sup>1,2,3</sup>, Noorossadat Torabi<sup>1,4</sup>, Yue Jia<sup>1,3</sup>, Jonathan Kent<sup>1</sup>, Stephen Martis<sup>1</sup>, Joshua A. Shapiro<sup>1,2,3</sup>, David Gresham<sup>1†</sup>, Amy A. Caudy<sup>1</sup> & Leonid Kruglyak<sup>1,2,3</sup>

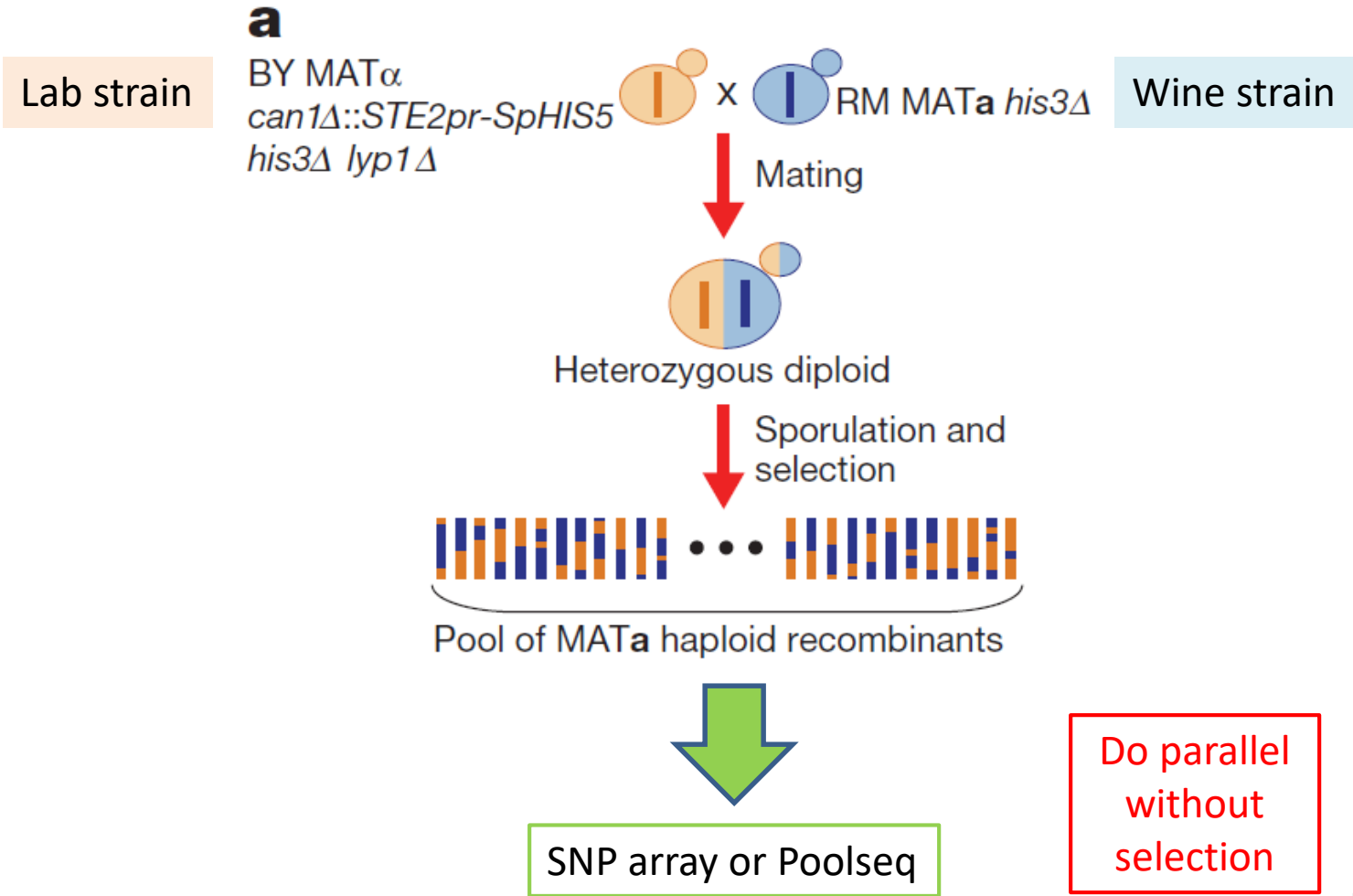
- step 1. generate segregating population of very large size
- step 2. selection based phenotyping (drug resistance or cell sorter)
- step 3. measure pooled allele frequencies (NGS)



- step 1. generate segregating population of very large size
- step 2. selection based phenotyping (drug resistance or cell sorter)
- step 3. measure pooled allele frequencies (NGS)



- step 1. generate segregating population of very large size
- step 2. selection based phenotyping (drug resistance or cell sorter)
- step 3. measure pooled allele frequencies (NGS)



# 4-NQO resistance case study

4-nitroquinoline

DNA damaging agent

Previously:

Conventional QTL

RAD5 (DNA repair gene)

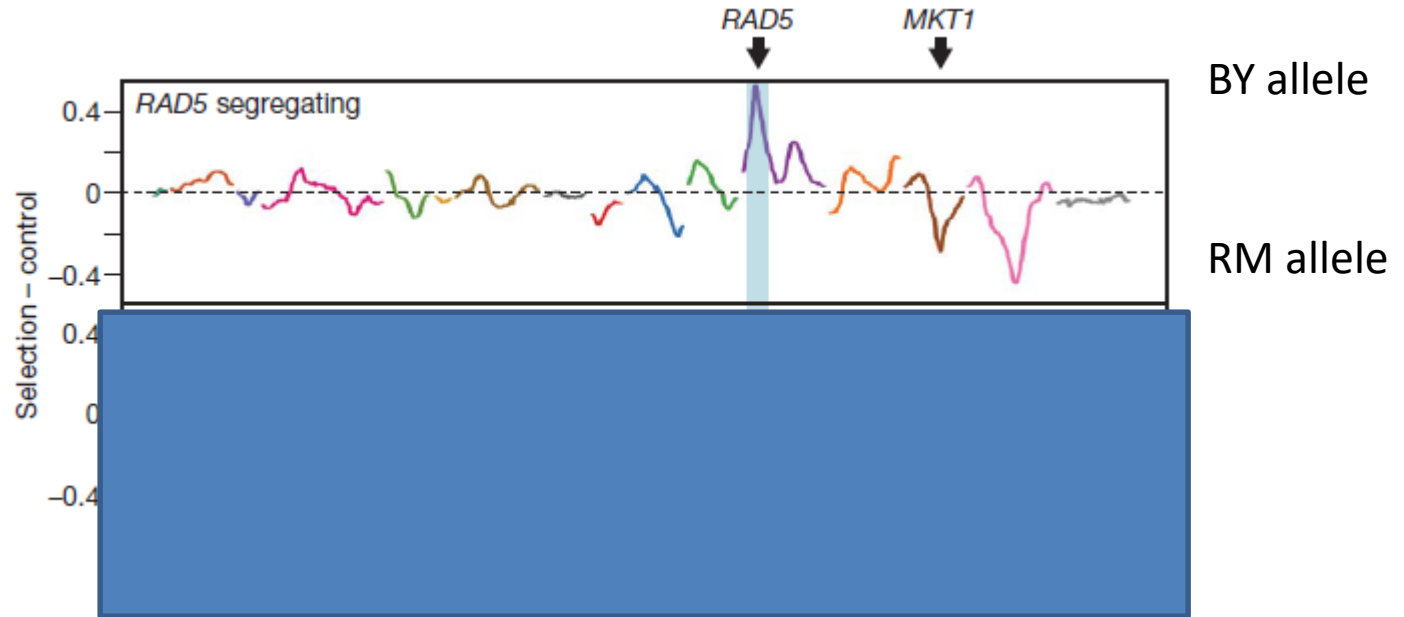
Follow up backcrossing revealed a 2<sup>nd</sup> gene:

MKT1



# 4-NQO resistance case study

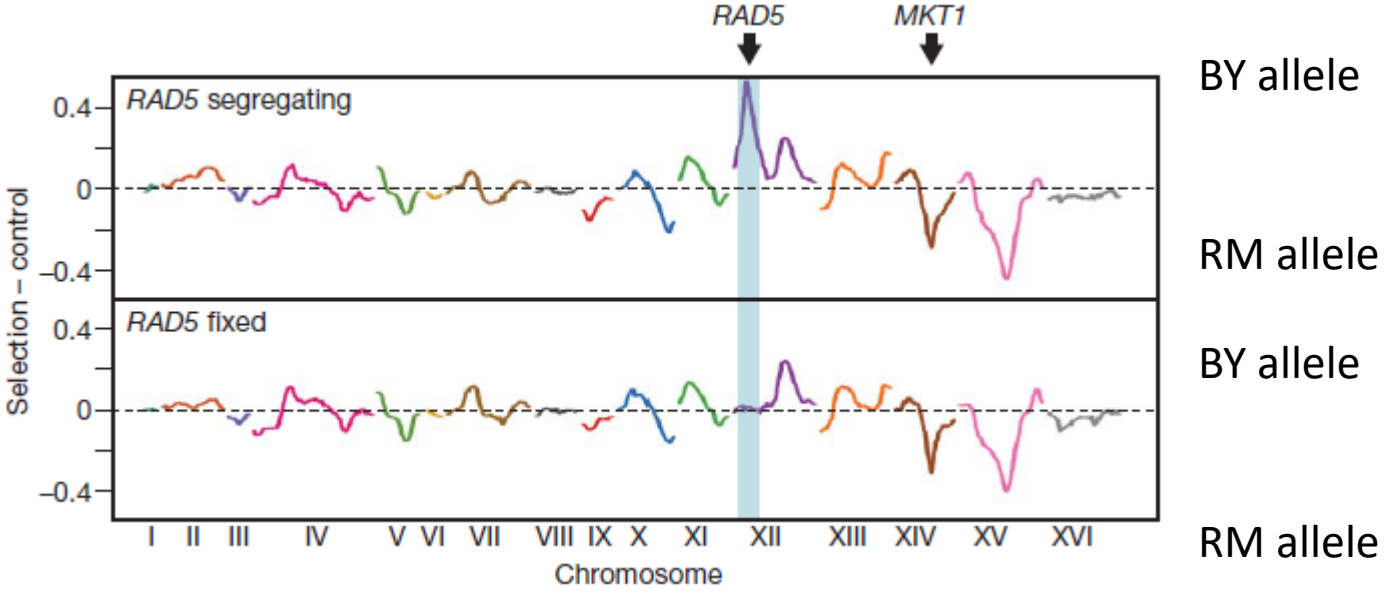
14 loci 70%



# 4-NQO resistance case study

14 loci 70%

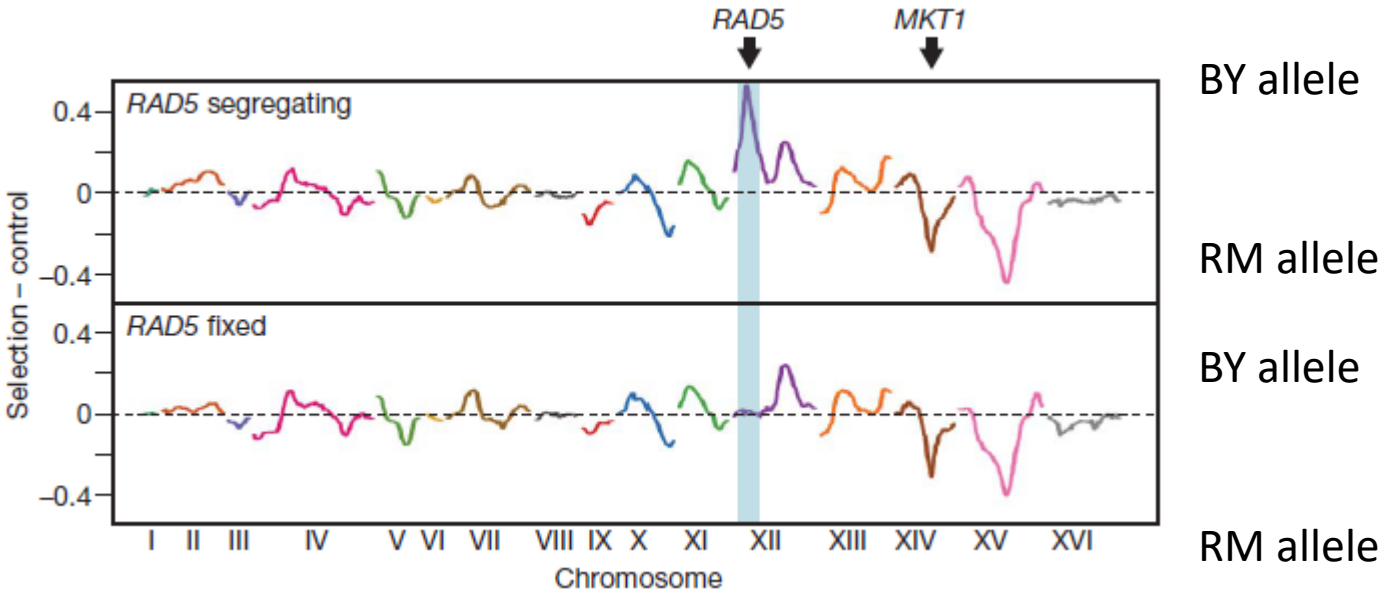
**FIX RAD5**



# 4-NQO resistance case study

14 loci 70%

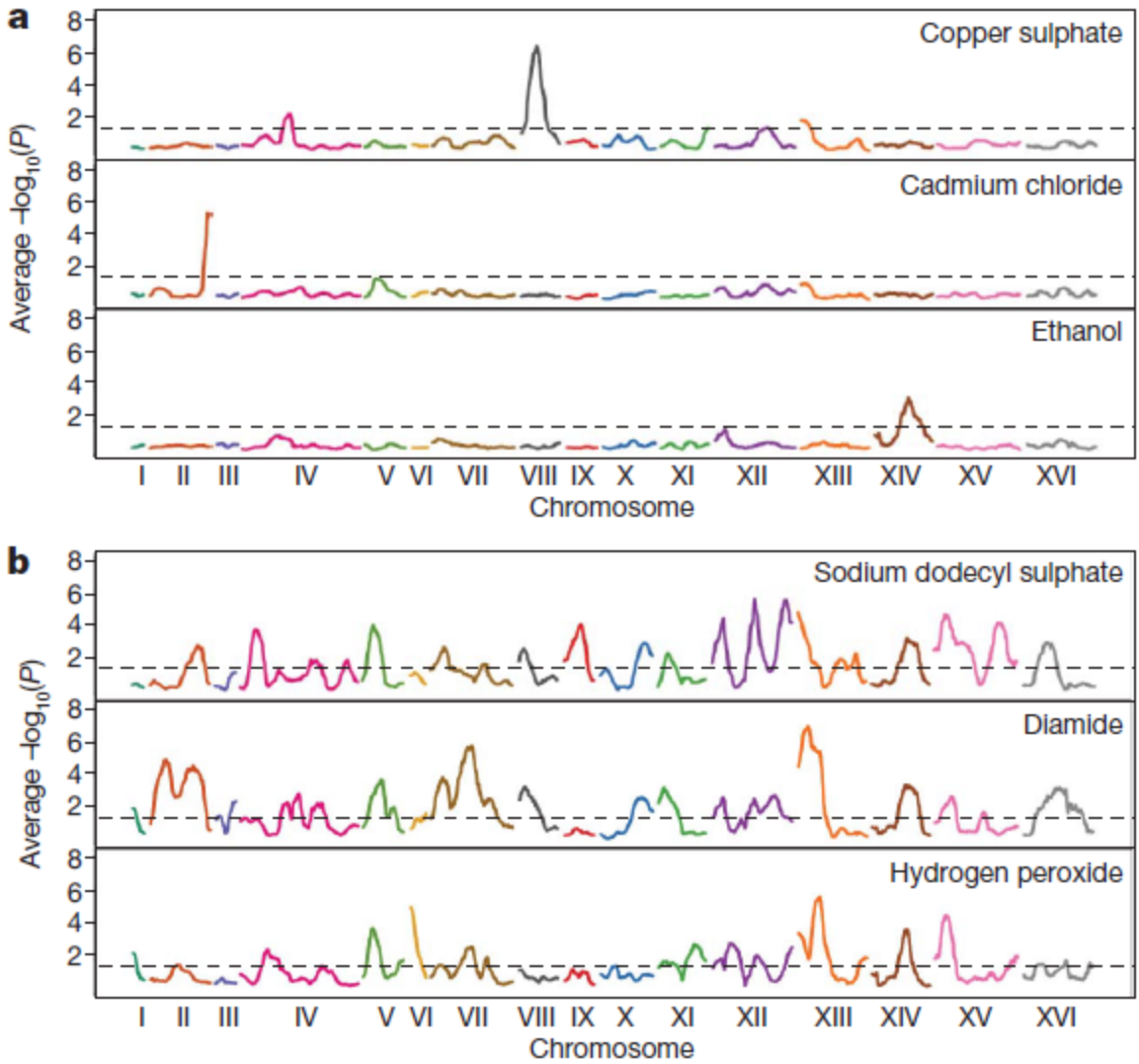
**FIX RAD5**



Very powerful  
Basically, never uncover this with single gene approach

# Other examples of simple and complex traits

x-QTL

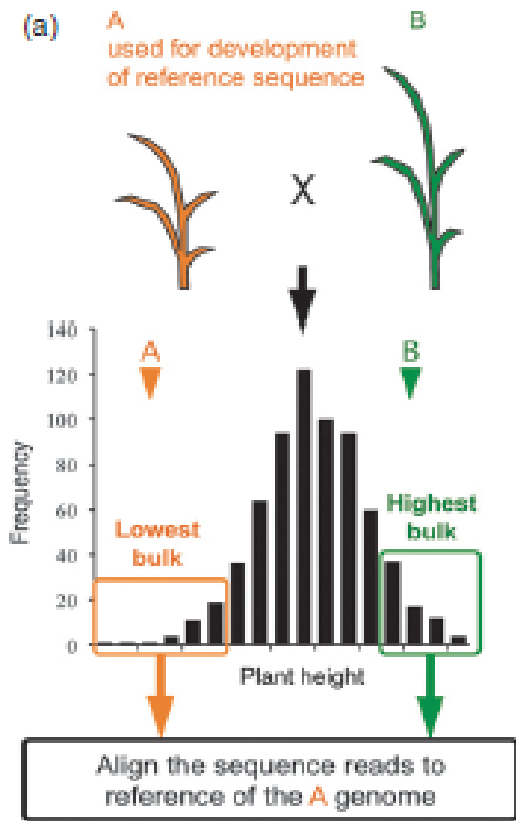


## xQTL mapping

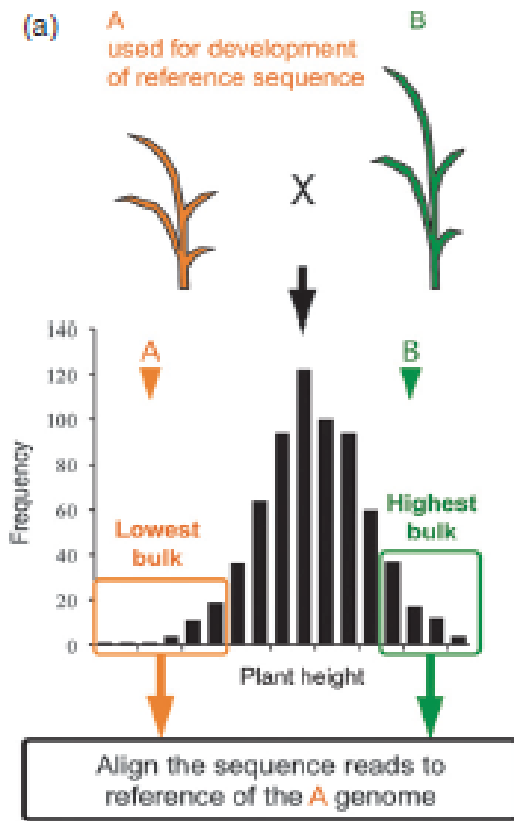
is advertised to require big N  
(more recombinants better resolution)

## QTL-seq

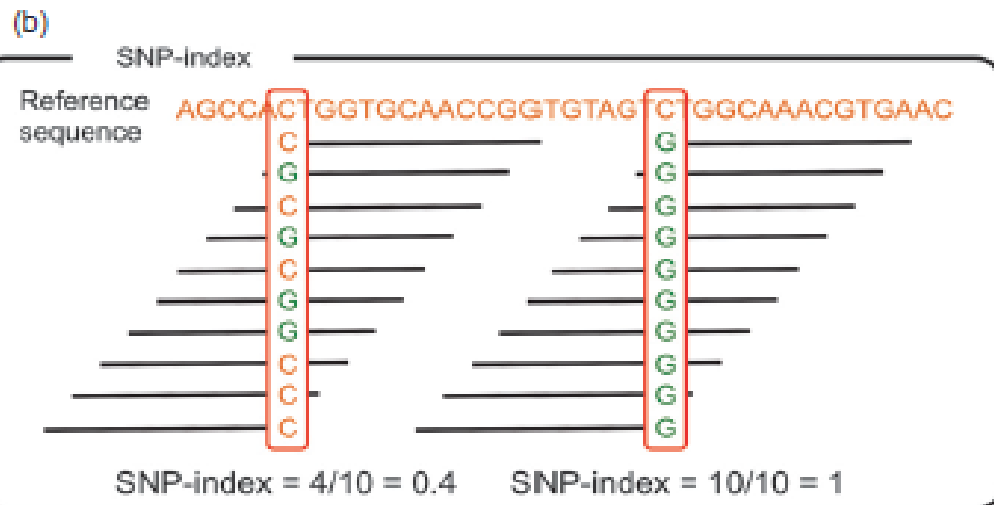
results suggest only 20-50 can be informative  
(although finding the exact gene is unclear)

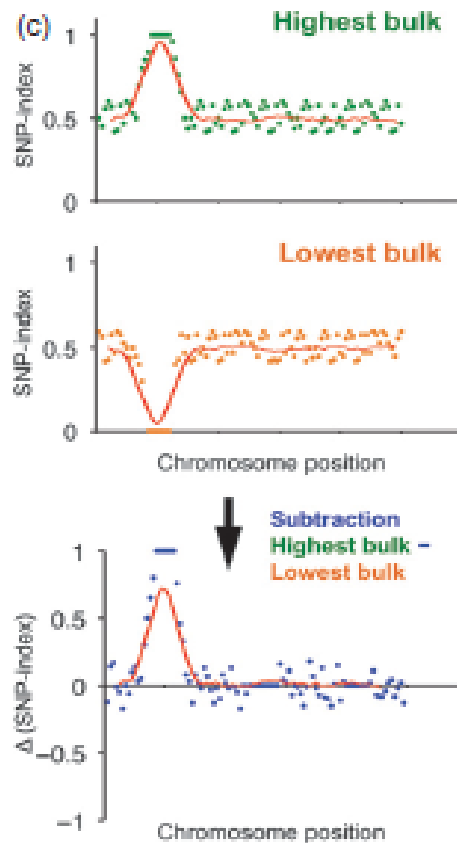
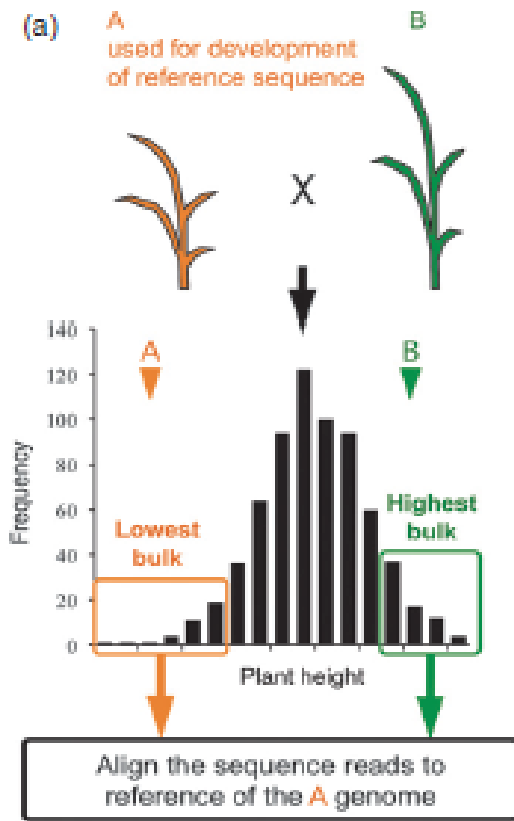


- N = 20-50
- Larger genome
- Works on crops
- ~MutMap with SNP index

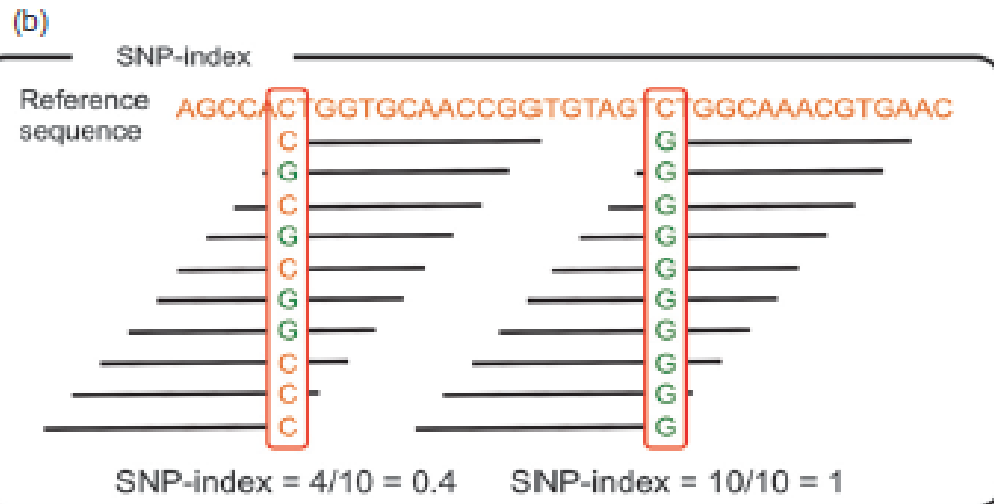


- N = 20-50
- Larger genome
- Works on crops
- ~MutMap with SNP index





- N = 20-50
- Larger genome
- Works on crops
- ~MutMap with SNP index





# Case study: seed vigor

(a)

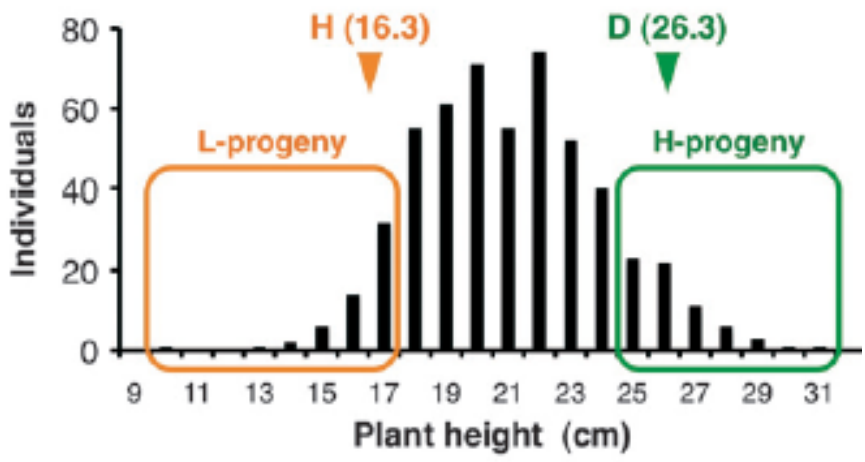


# Case study: seed vigor

(a)



(b)

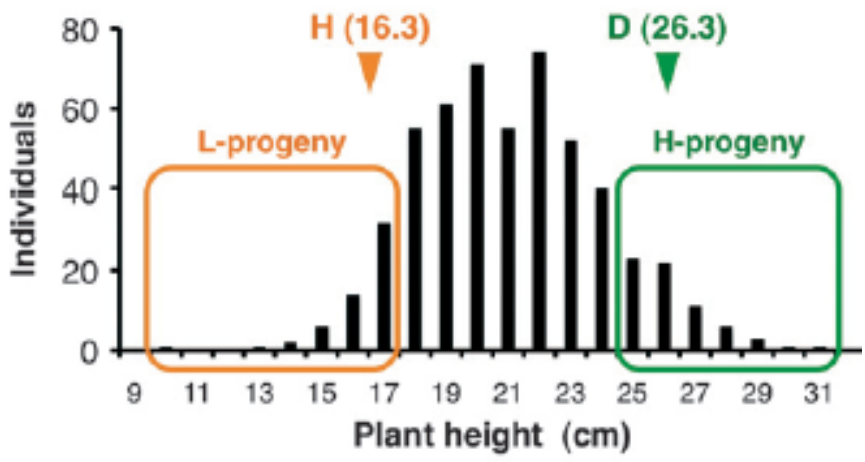


# Case study: seed vigor

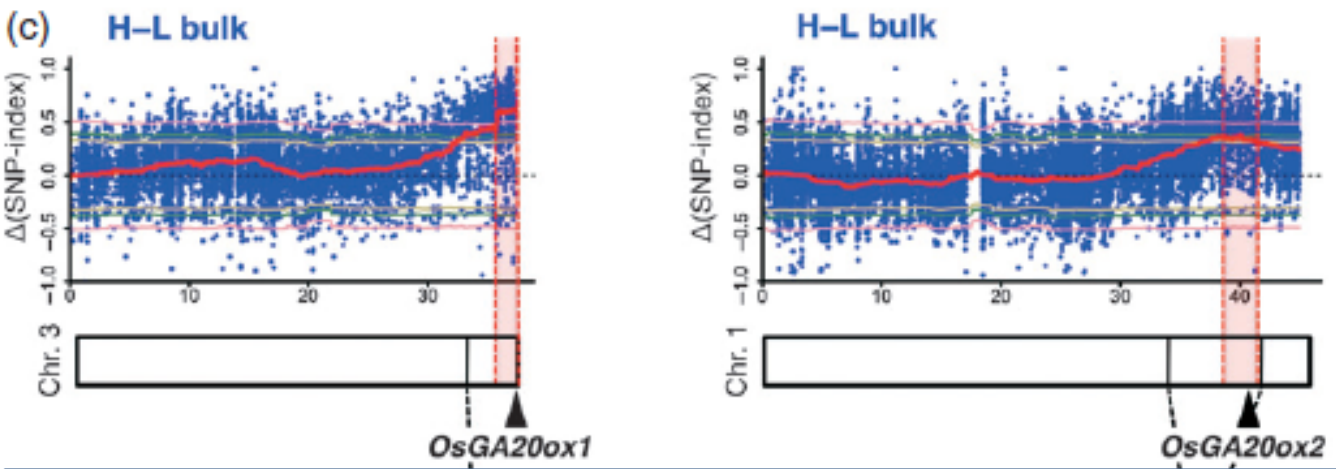
(a)



(b)



(c)

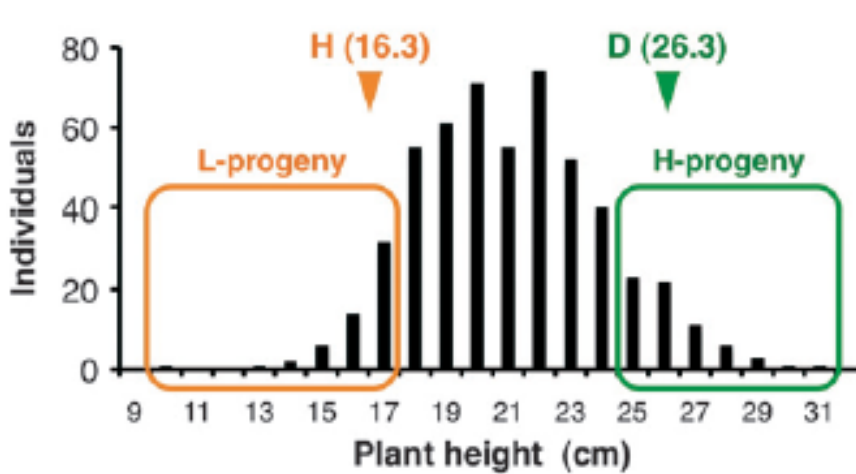


# Case study: seed vigor

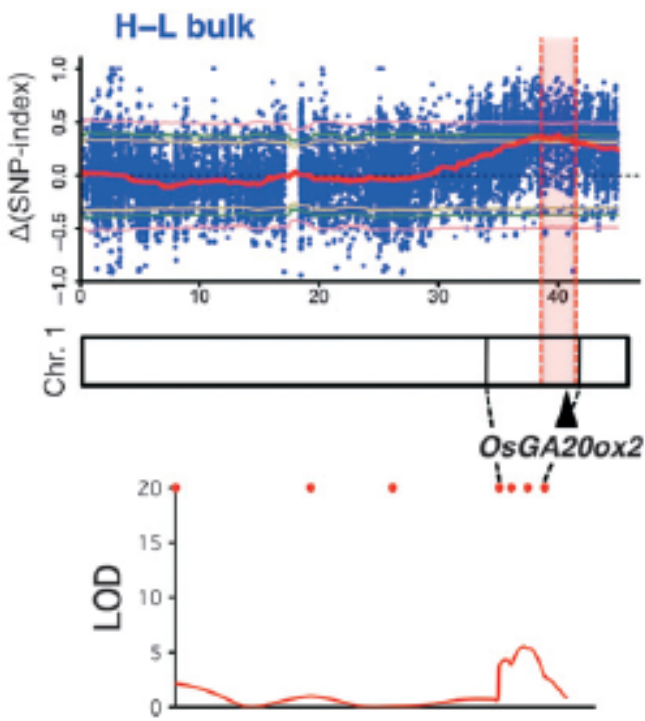
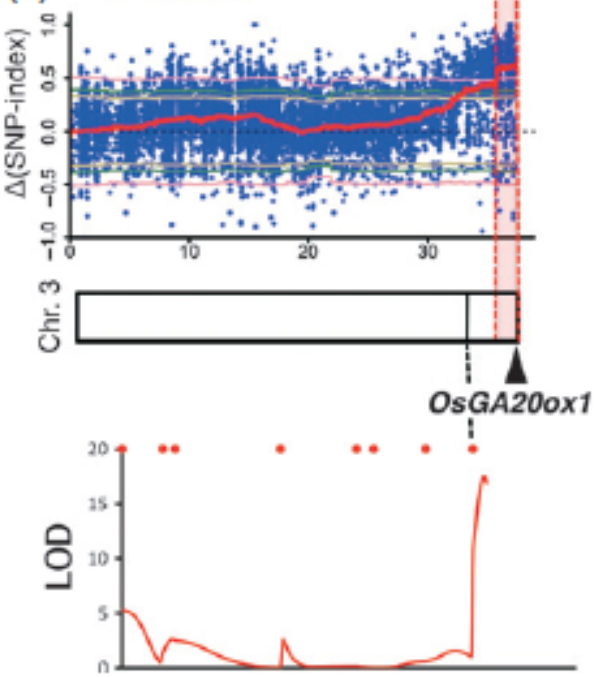
(a)



(b)



(c) H-L bulk



# Poolseq

## GWAS

taking advantage of historical recombination

# Poolseq

## GWAS

taking advantage of historical recombination

Basically, the more recombinants, the better map resolution.

F2 crosses or Recombinant Inbred Lines have limited recombination events.

Wild samples have history to help

# Poolseq

## GWAS

taking advantage of historical recombination

Basically, the more recombinants, the better map resolution.

F2 crosses or Recombinant Inbred Lines have limited recombination events.

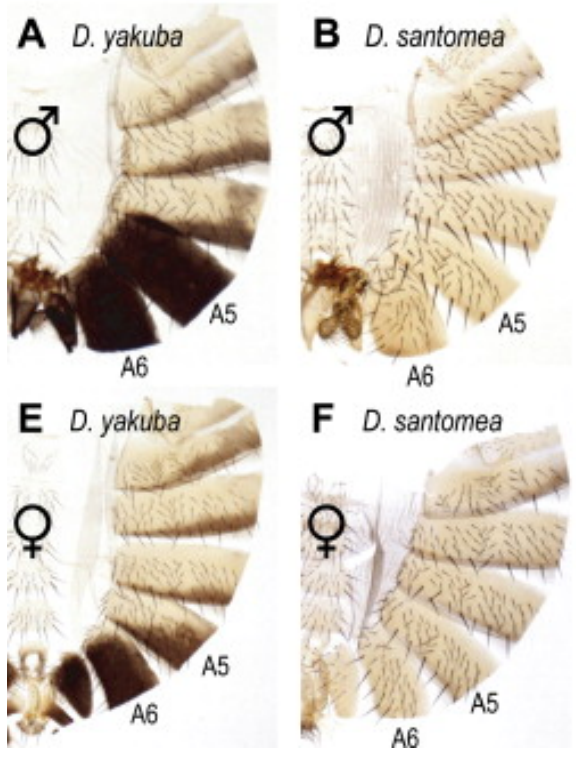
Wild samples have history to help

*Drosophila* body pigmentation example

# *Drosophila* pigmentation background

## *D. yakuba* x *D. santomea*

- There is variation in abdominal pigmentation
- Probably sexual selection



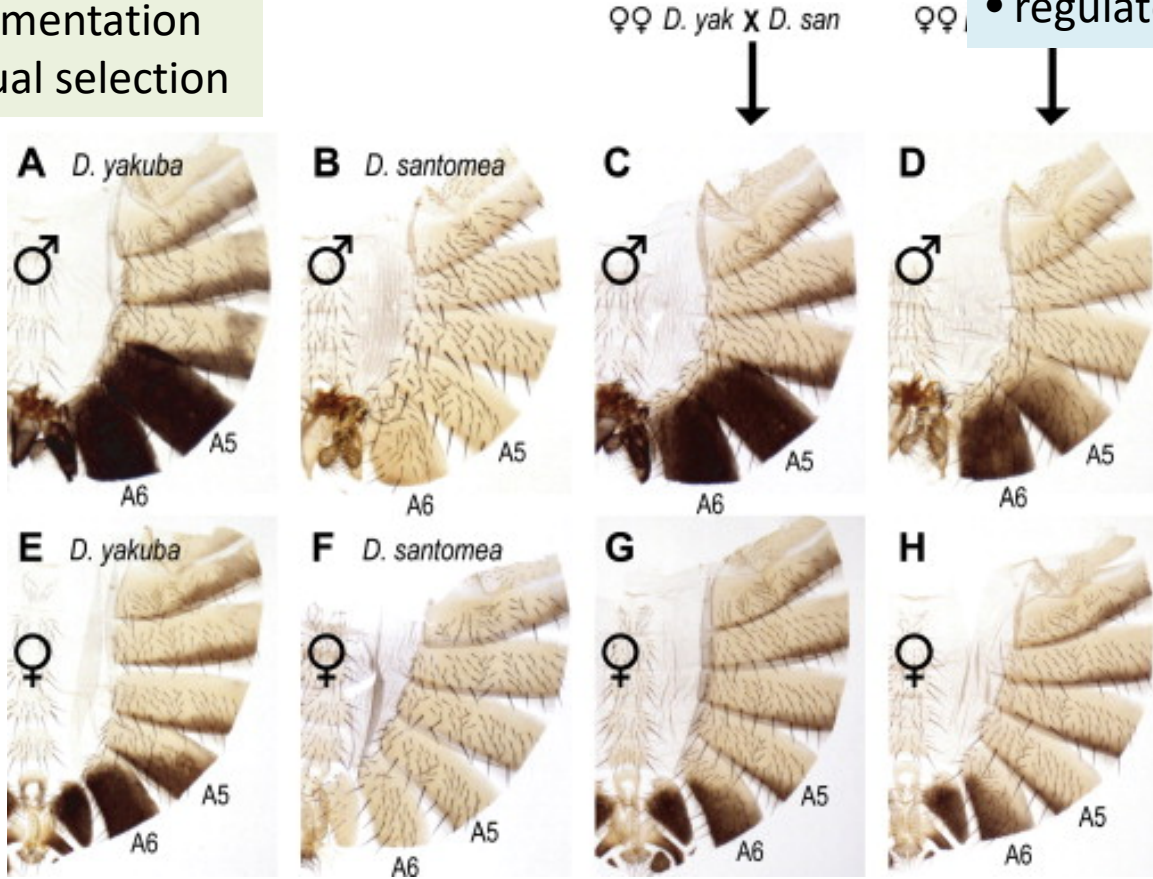


# Drosophila pigmentation background

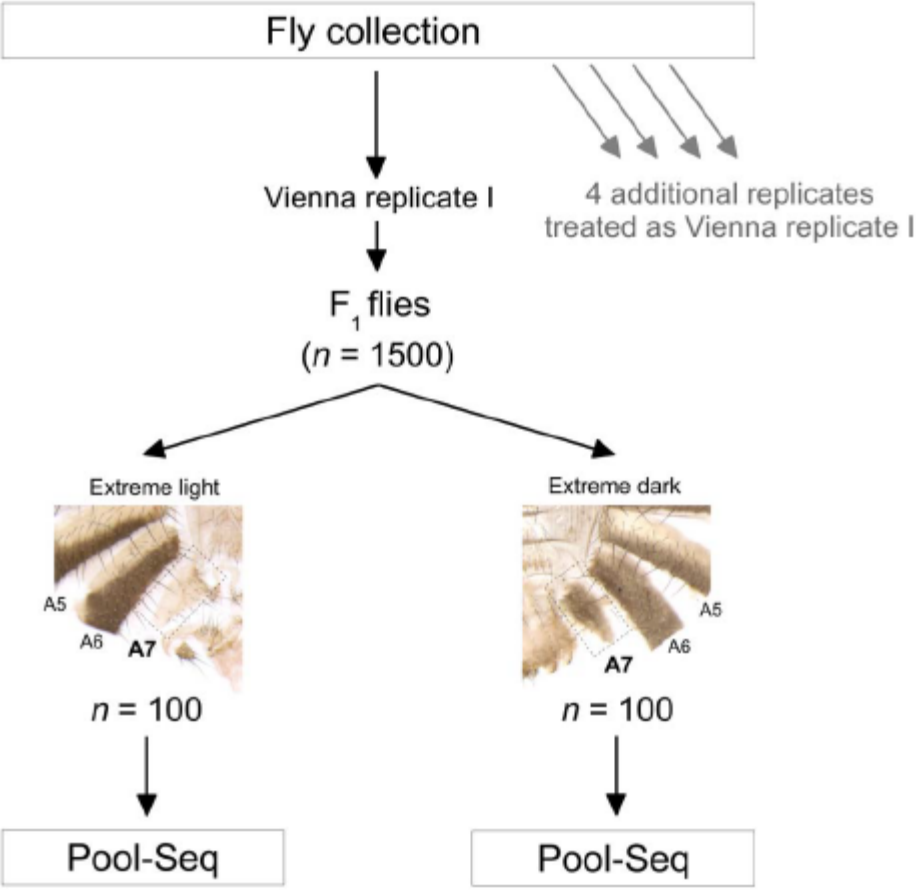
## D. yakuba x D. santomea

- There is variation in abdominal pigmentation
- Probably sexual selection

- Hybridization experiments identified: *tan* and *yellow*
- regulatory differences

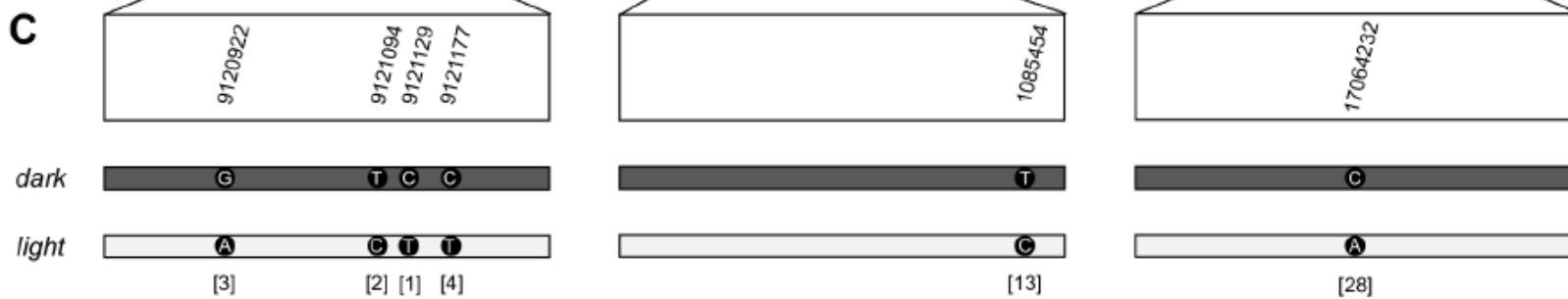
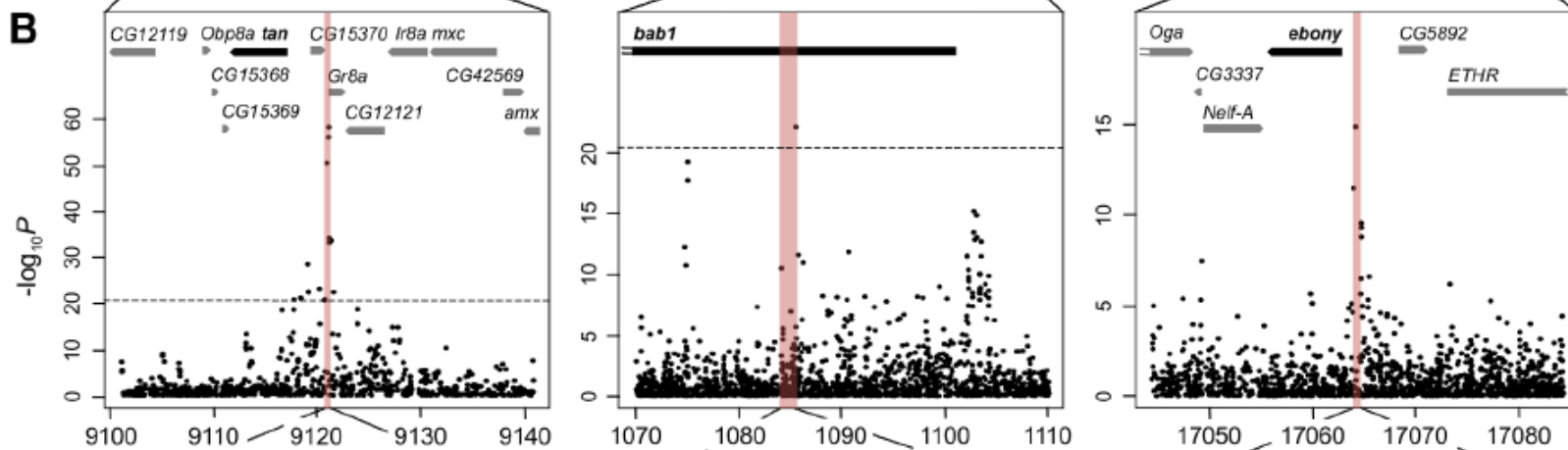
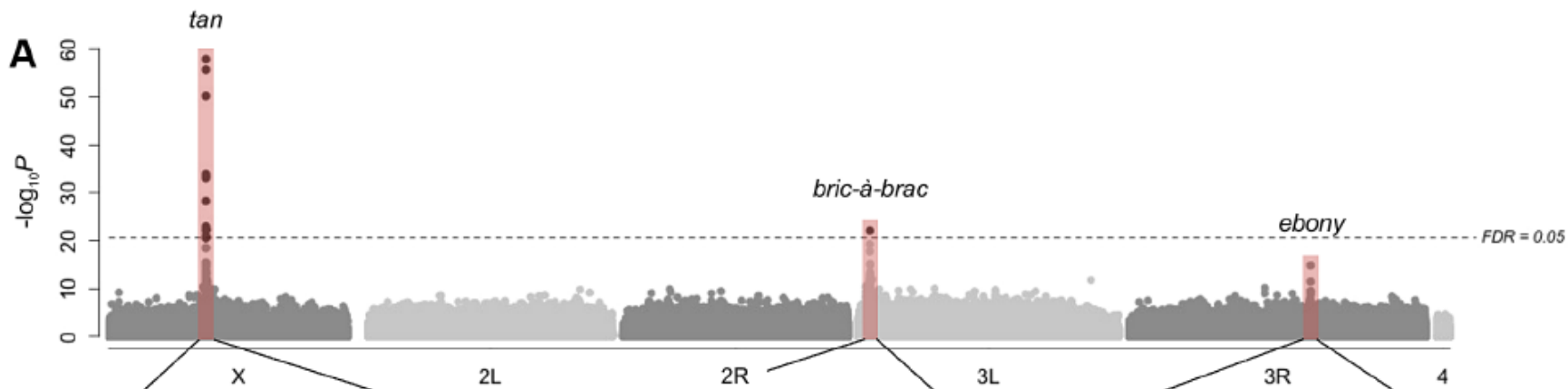


# *D. melanogaster* → natural variation in pigmentation

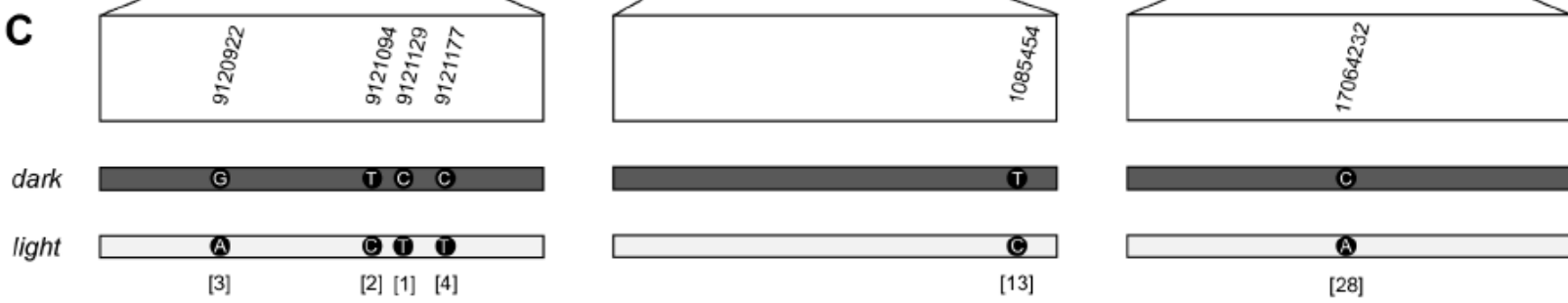
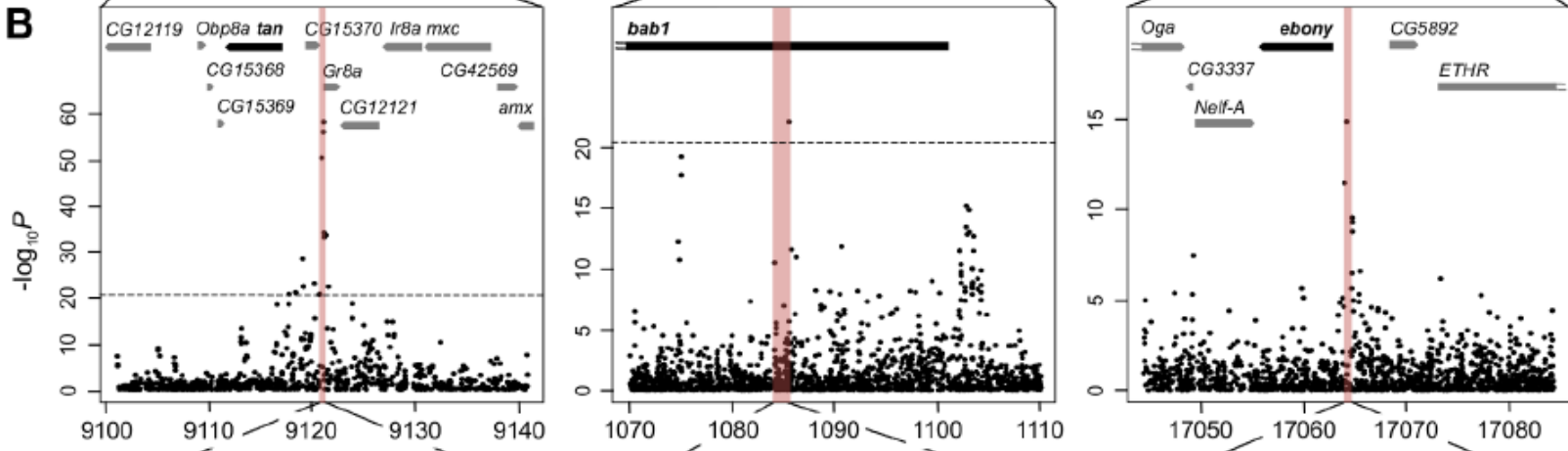
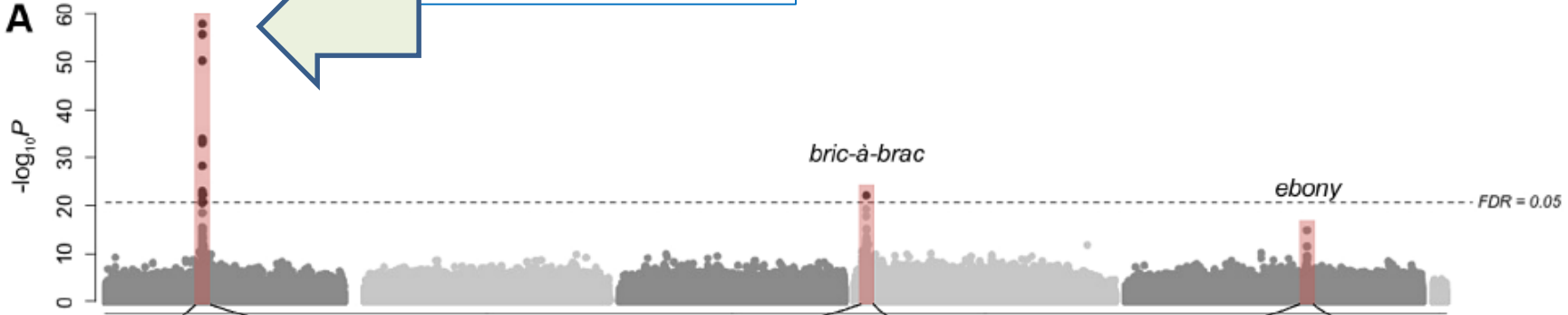


2 pops from Vienna, Austria  
3 pops from Bolzano, Italy

Grow in lab controlled conditions 1 gen  
Sort by color  
Take extremes  
Pool-seq with replicates

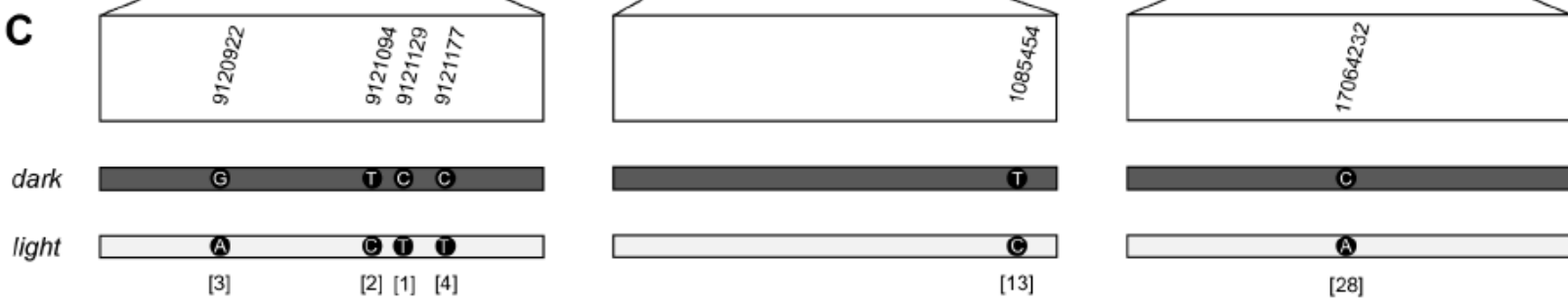
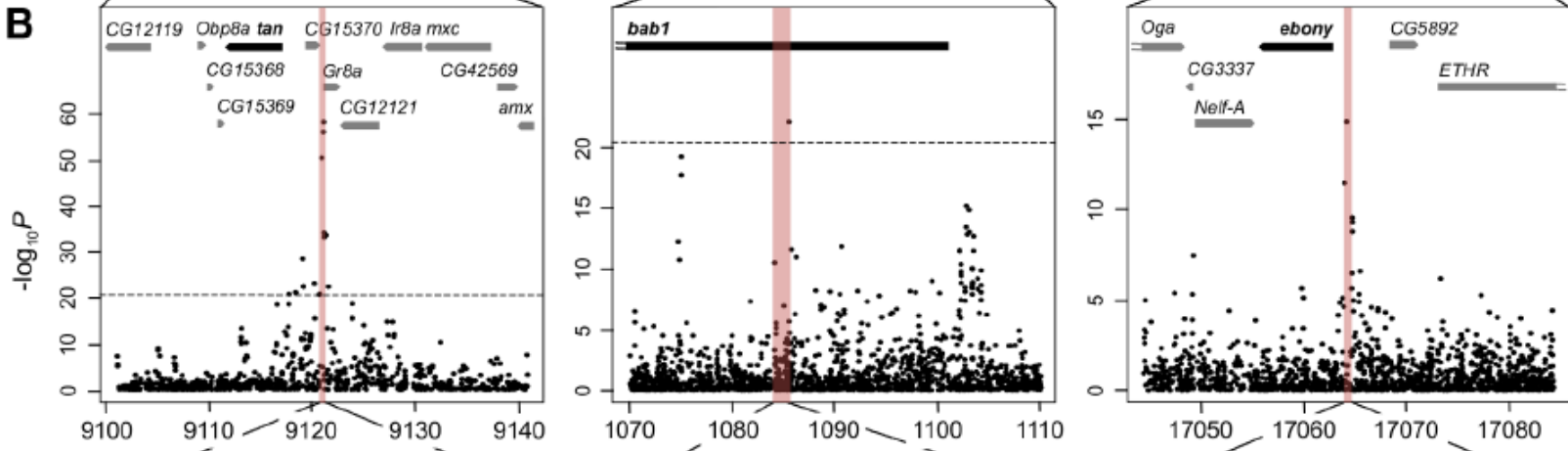
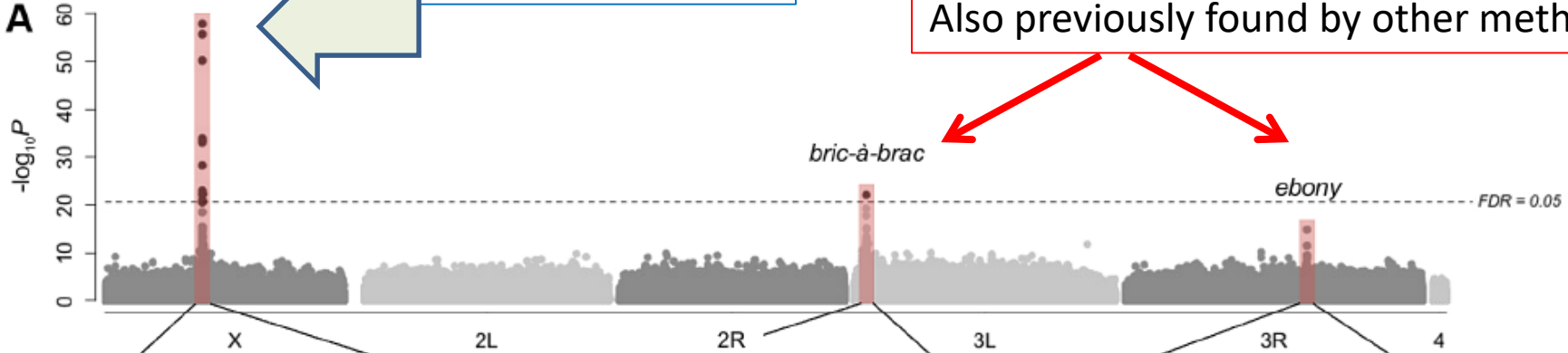


Matches *D. yakuba/ santomea* result



Matches *D. yakuba/santomea* result

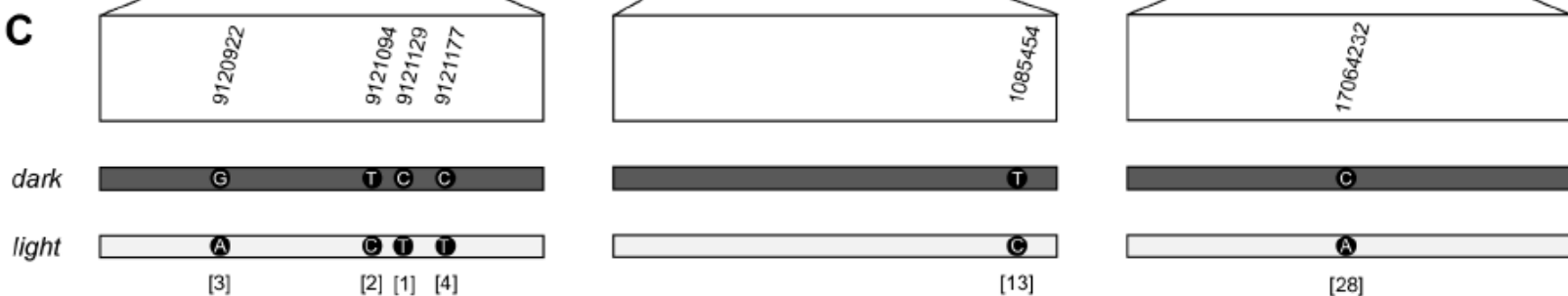
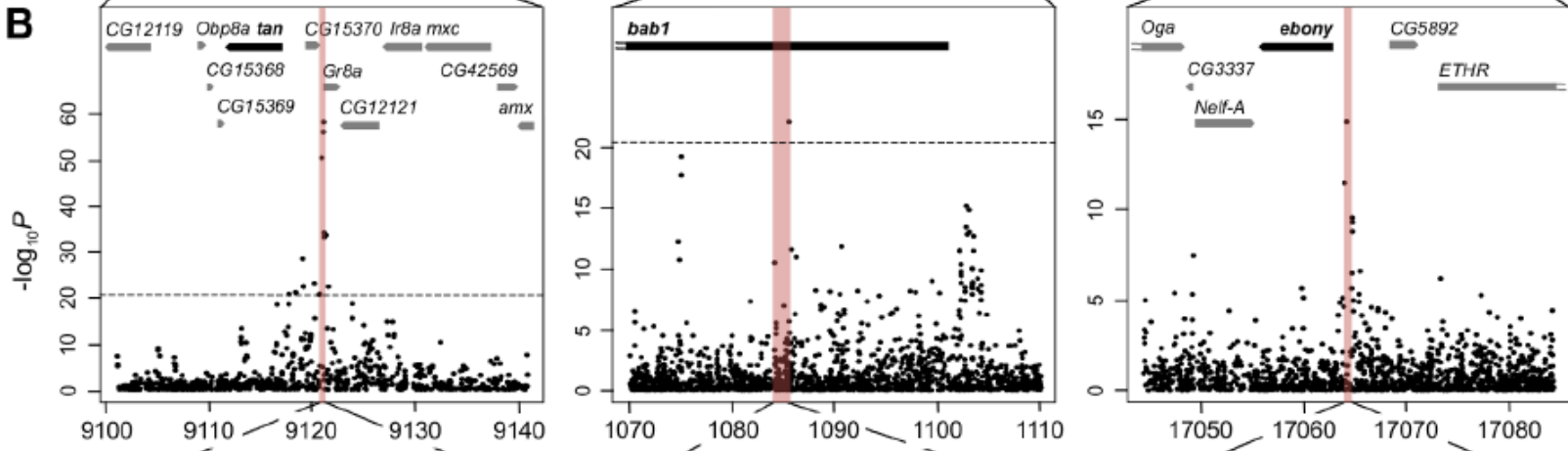
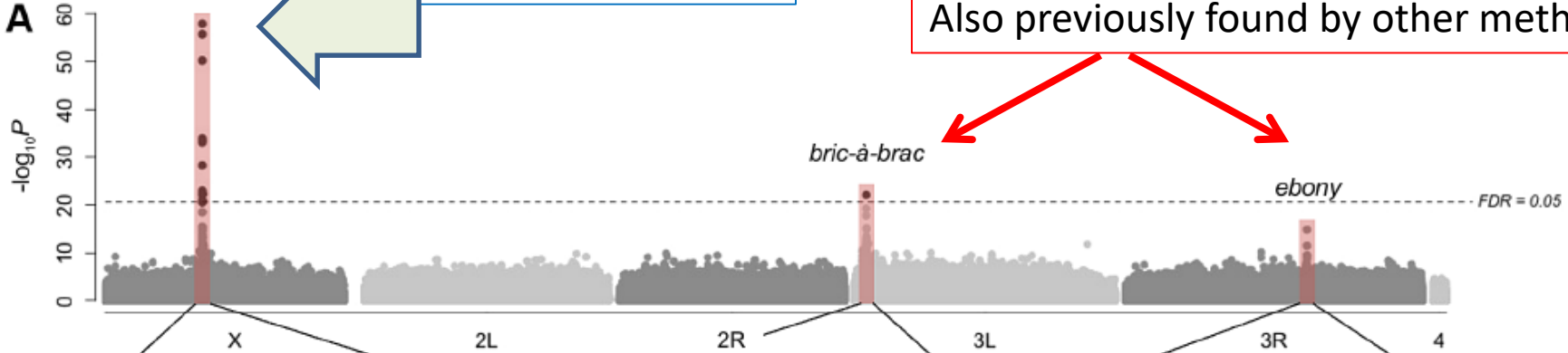
Also previously found by other methods



All regulatory

Matches *D. yakuba/santomea* result

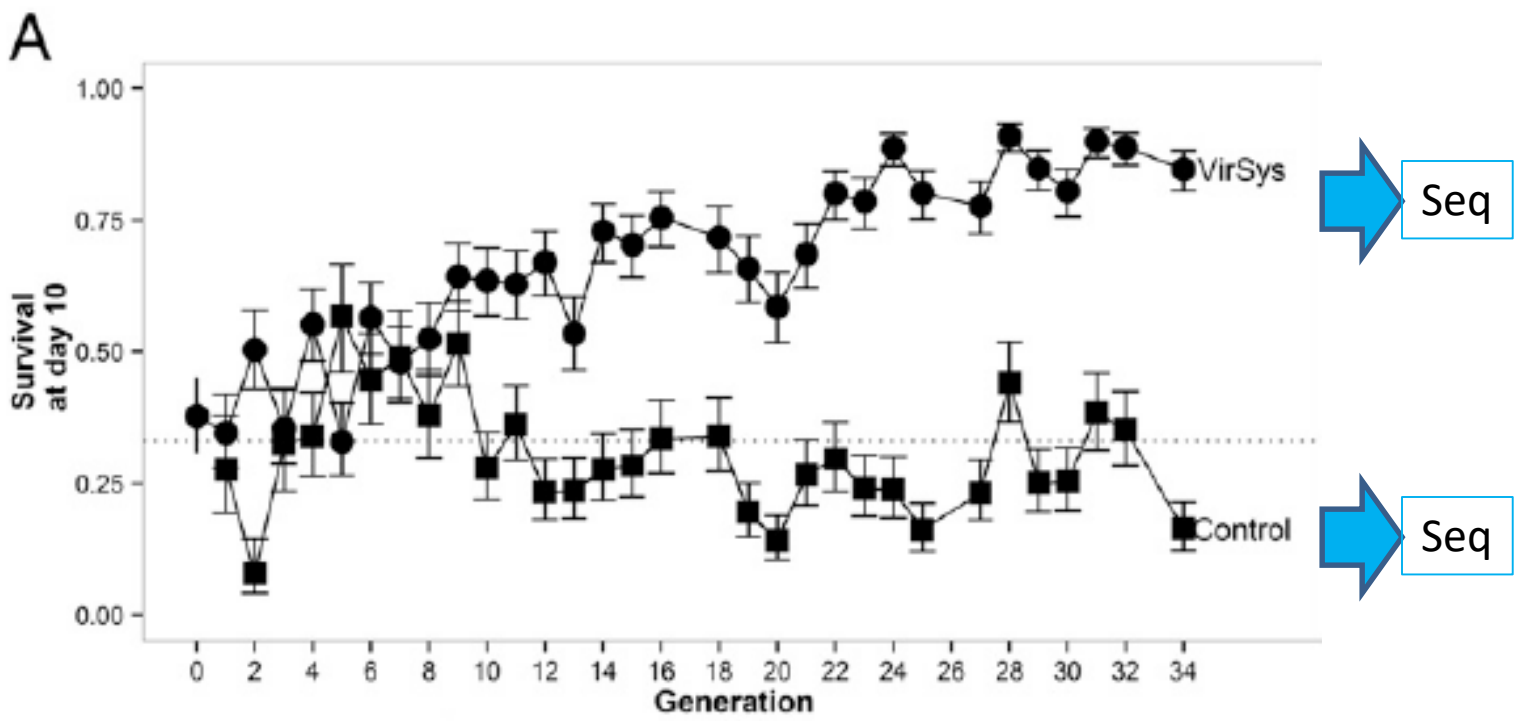
Also previously found by other methods



Evolve and re-sequence

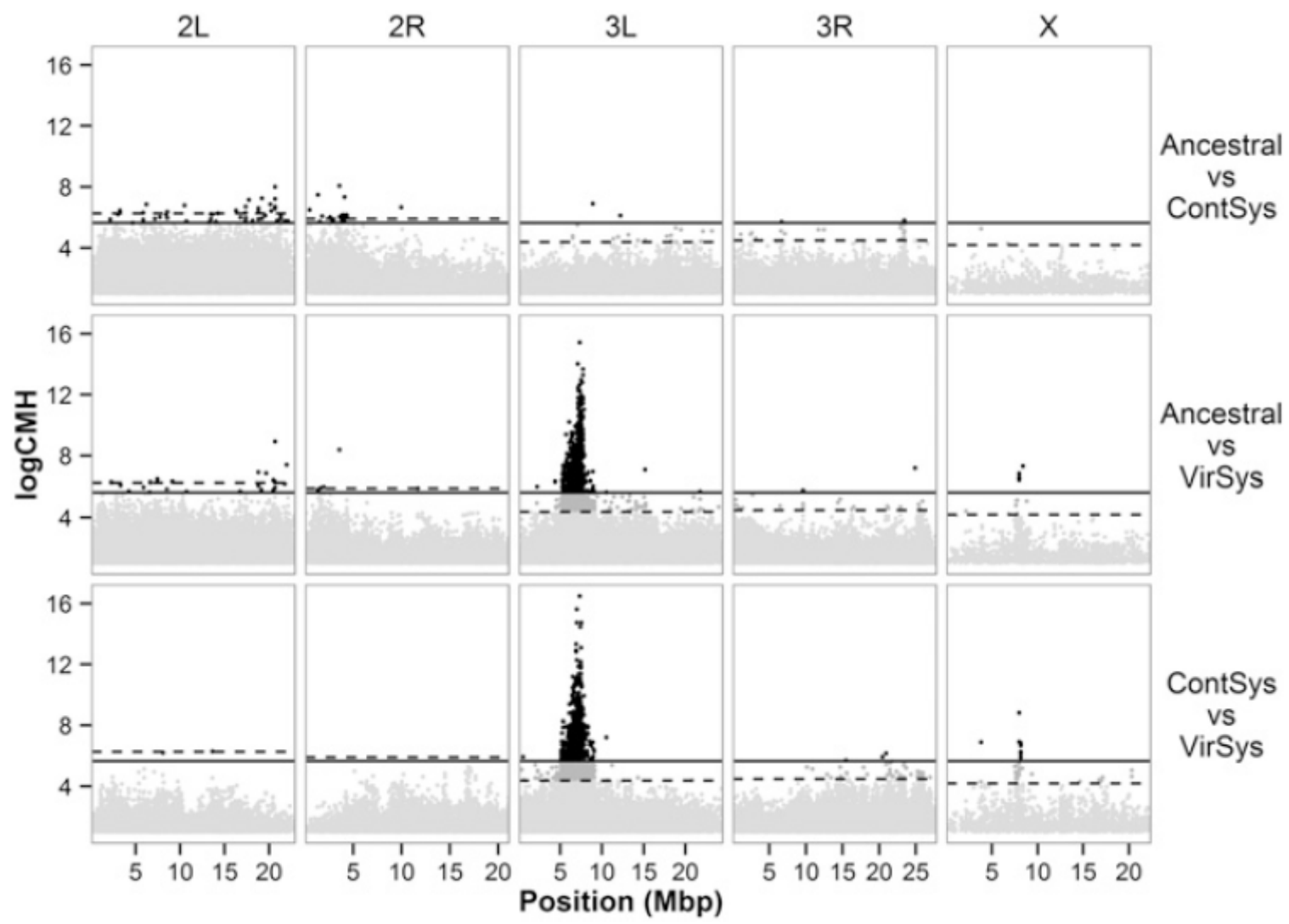
# Evolve and re-sequence

Evolve *Drosophila* lines to be more resistant to C virus  
4x replicates [virus, control1, control2]





# Evolve and re-sequence

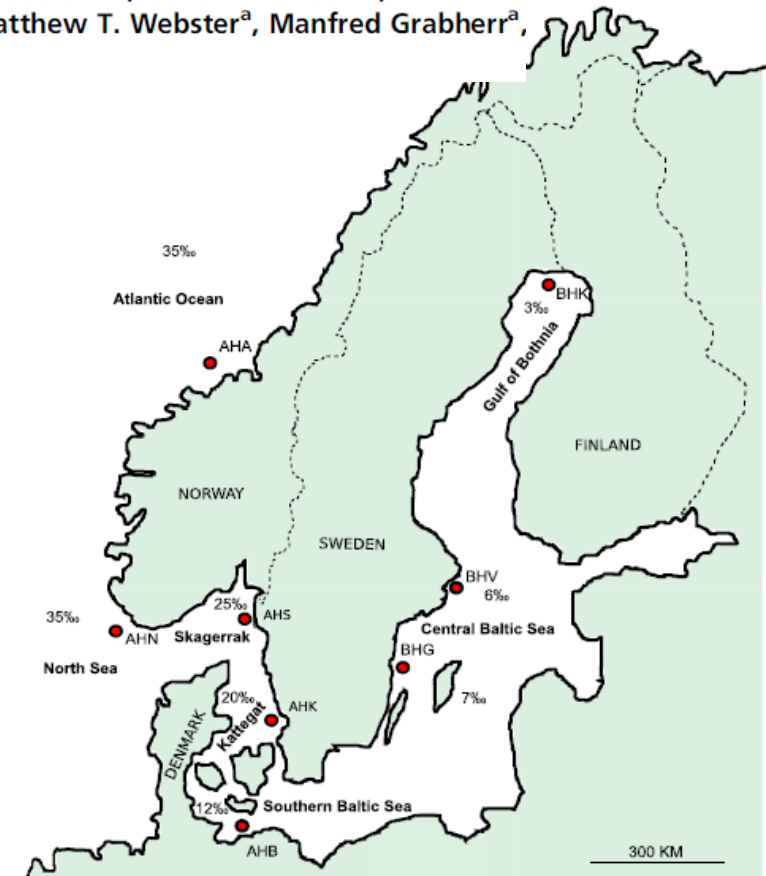
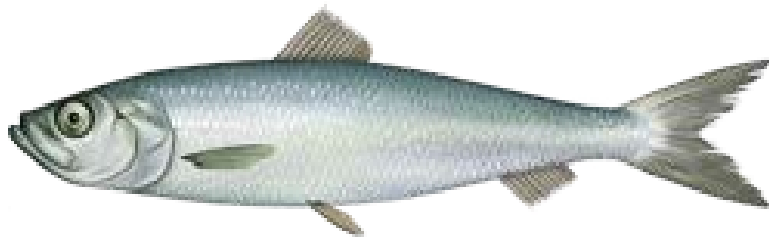


*pastrel* gene (and lesser extent *Ubc-E2H*)

# Ecology (reverse ecology) and local adaptation

## Population-scale sequencing reveals genetic differentiation due to local adaptation in Atlantic herring

Sangeet Lamichhane<sup>a,1</sup>, Alvaro Martinez Barrio<sup>a,1</sup>, Nima Rafati<sup>a,1</sup>, Görel Sundström<sup>a,1</sup>, Carl-Johan Rubin<sup>a</sup>, Elizabeth R. Gilbert<sup>a,2</sup>, Jonas Berglund<sup>a</sup>, Anna Wetterbom<sup>b</sup>, Linda Laikre<sup>c</sup>, Matthew T. Webster<sup>a</sup>, Manfred Grabherr<sup>a</sup>, Nils Ryman<sup>c</sup>, and Leif Andersson<sup>a,d,3</sup>



# Ecology (reverse ecology) and local adaptation



8 Pool-seq populations

Each with 50 individuals

30x coverage

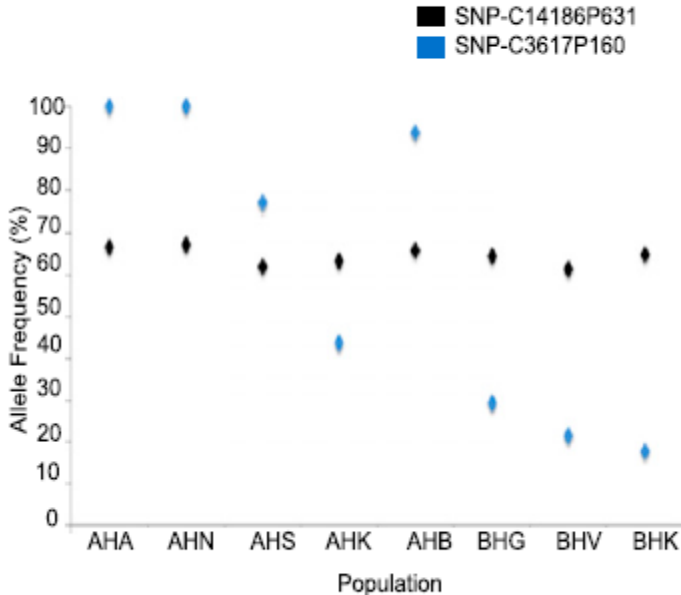
Map against transcriptome

# Ecology (reverse ecology) and local adaptation



8 Pool-seq populations  
Each with 50 individuals  
30x coverage  
Map against transcriptome

A



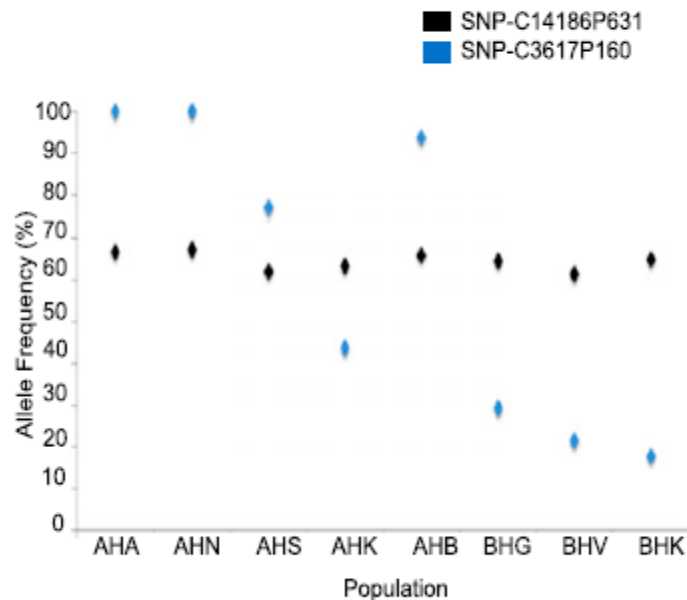
2 SNPs with different allele frequency patterns

# Ecology (reverse ecology) and local adaptation



8 Pool-seq populations  
Each with 50 individuals  
30x coverage  
Map against transcriptome

A



2 SNPs with different  
allele frequency patterns

Most undifferentiated

~440,000

Few highly differentiated <1%

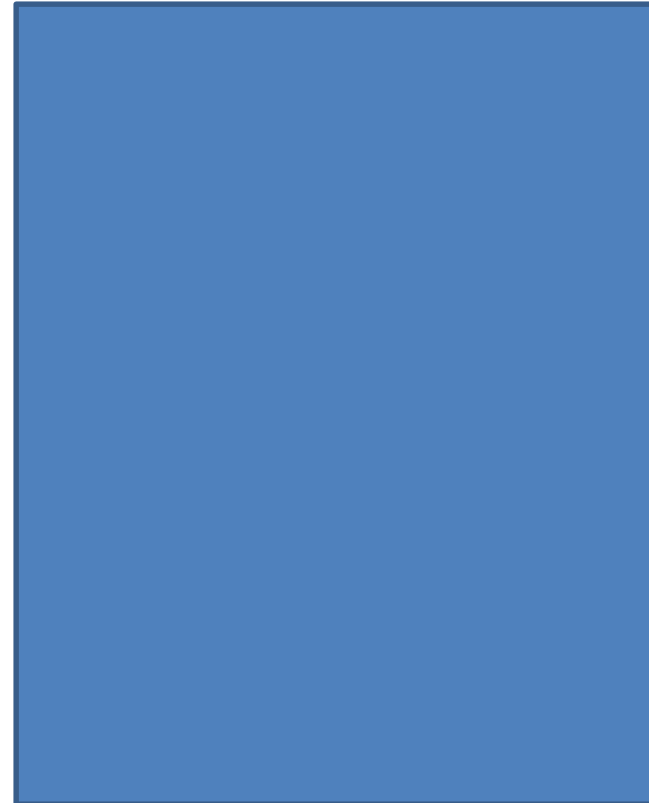
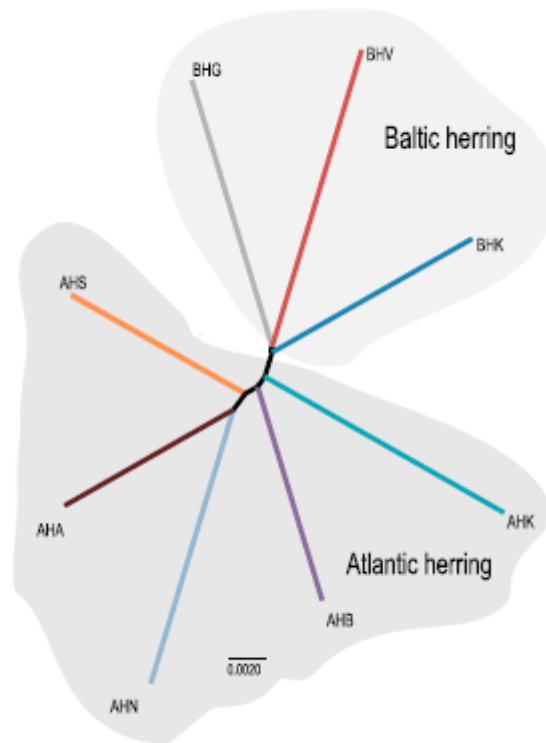
~3800

# Little difference between Baltic and Atlantic if using all SNPs...



All SNPs →  
(star phylogeny)

B

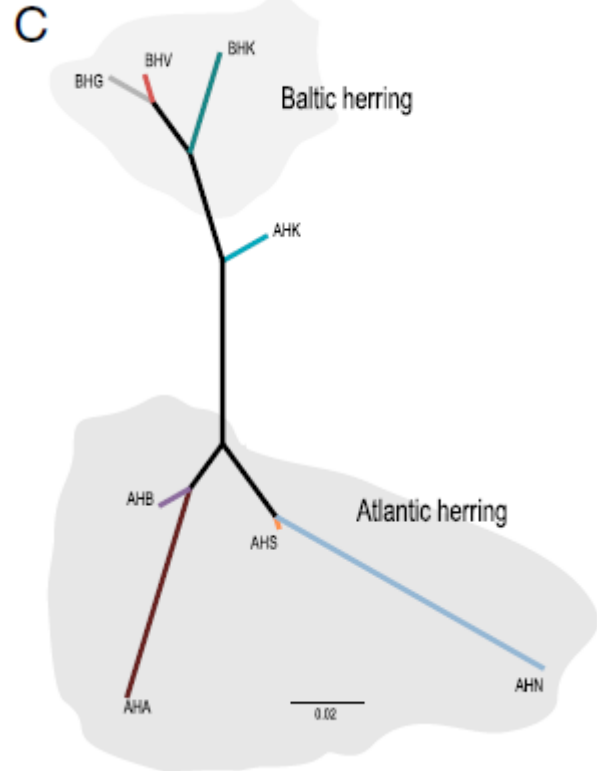
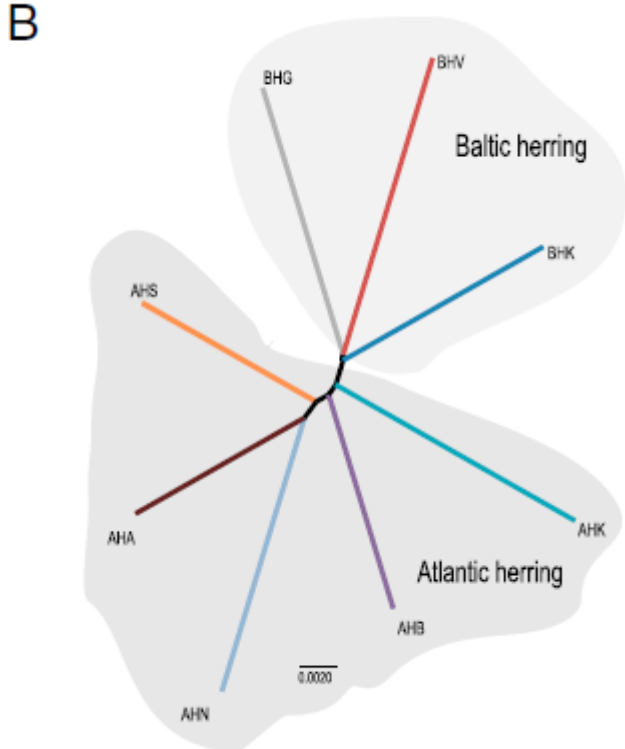


# Little difference between Baltic and Atlantic if using all SNPs...



All SNPs →  
(star phylogeny)

High Fst SNPs →  
separate populations



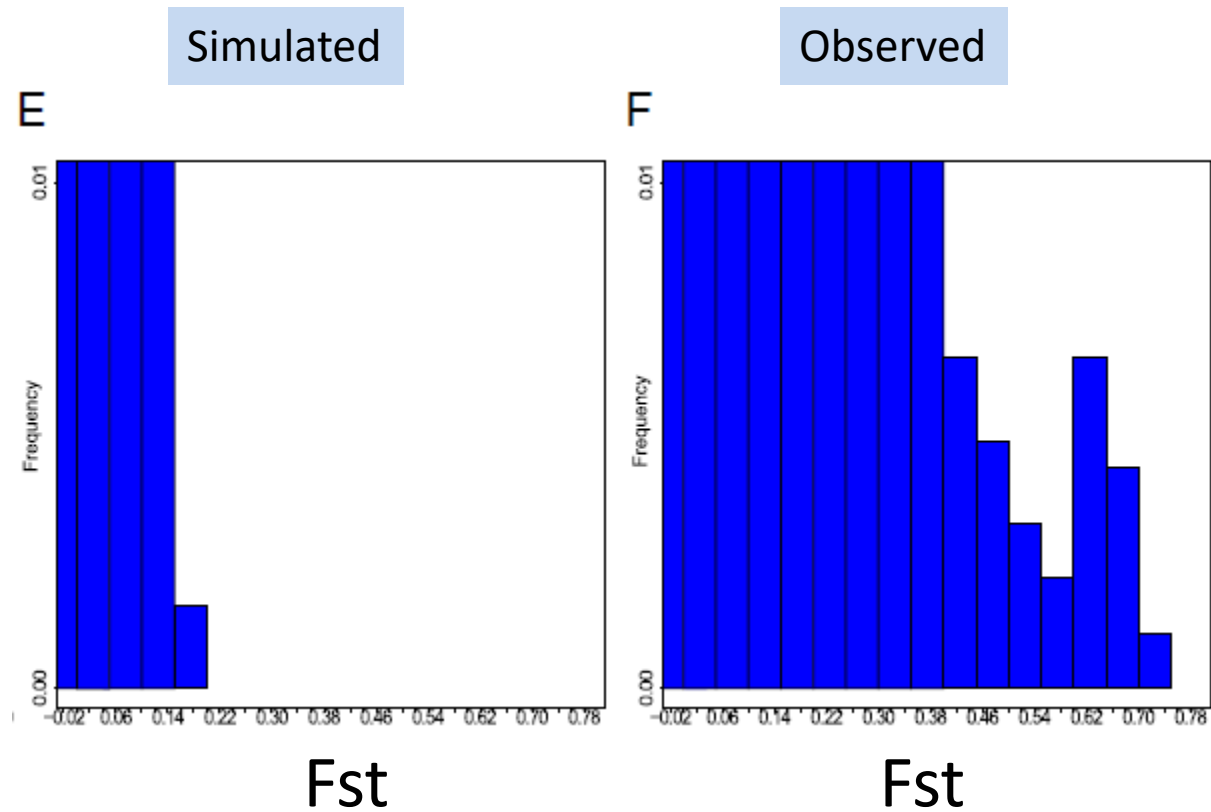
# Are high $F_{st}$ SNPs from drift or selection?

Simulation analysis to test



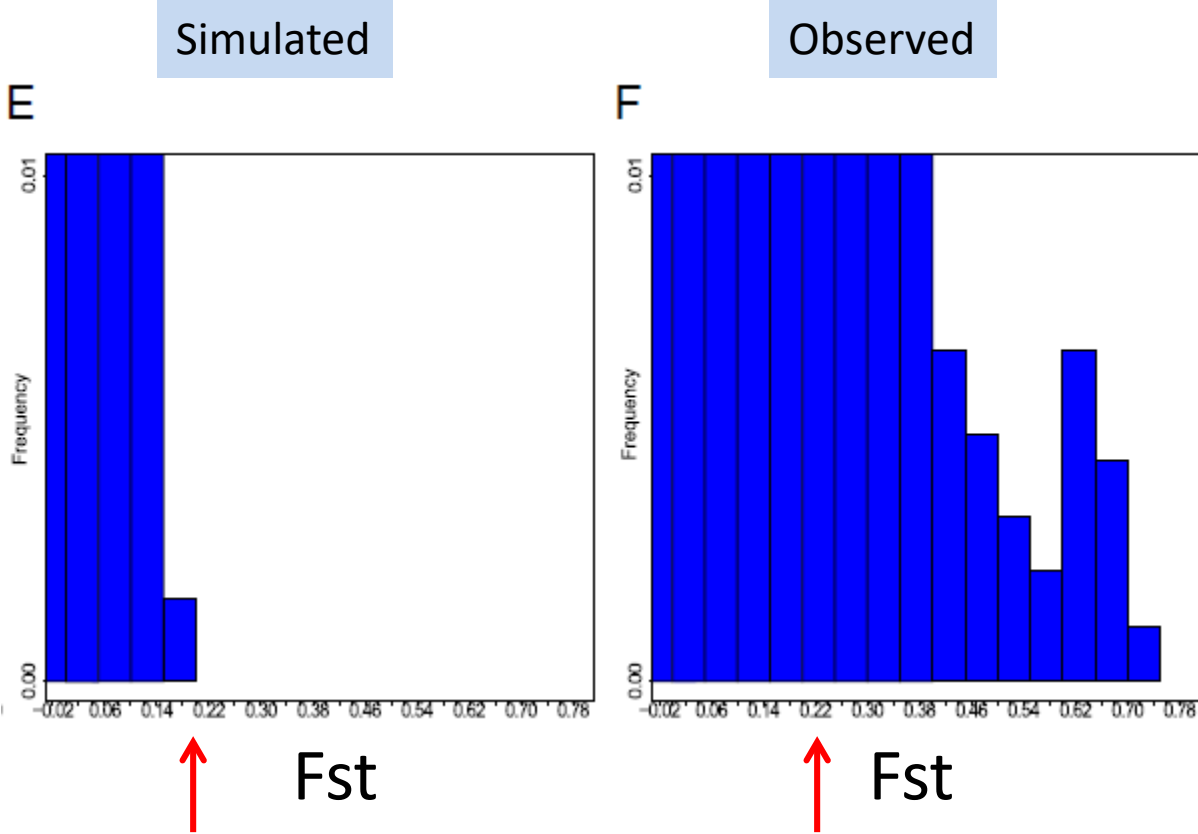
# Are high $F_{st}$ SNPs from drift or selection?

Simulation analysis to test



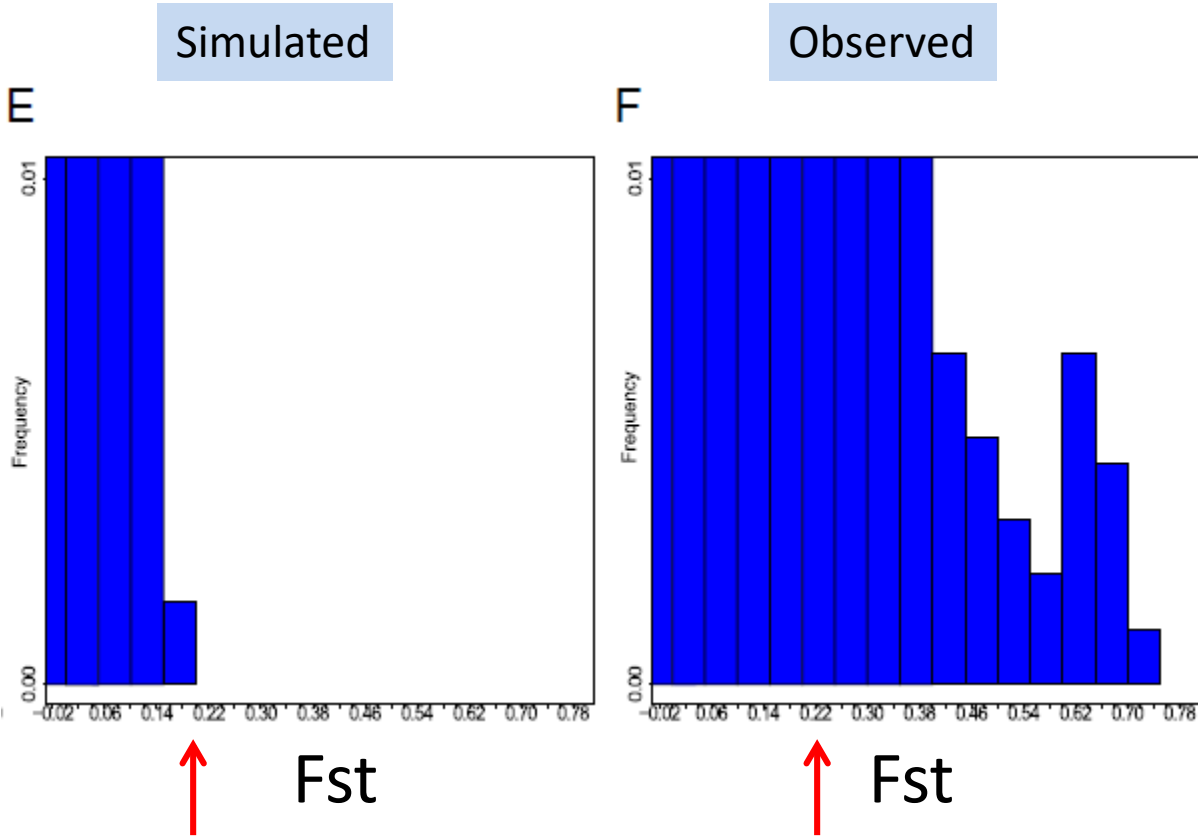
# Are high Fst SNPs from drift or selection?

Simulation analysis to test



# Are high Fst SNPs from drift or selection?

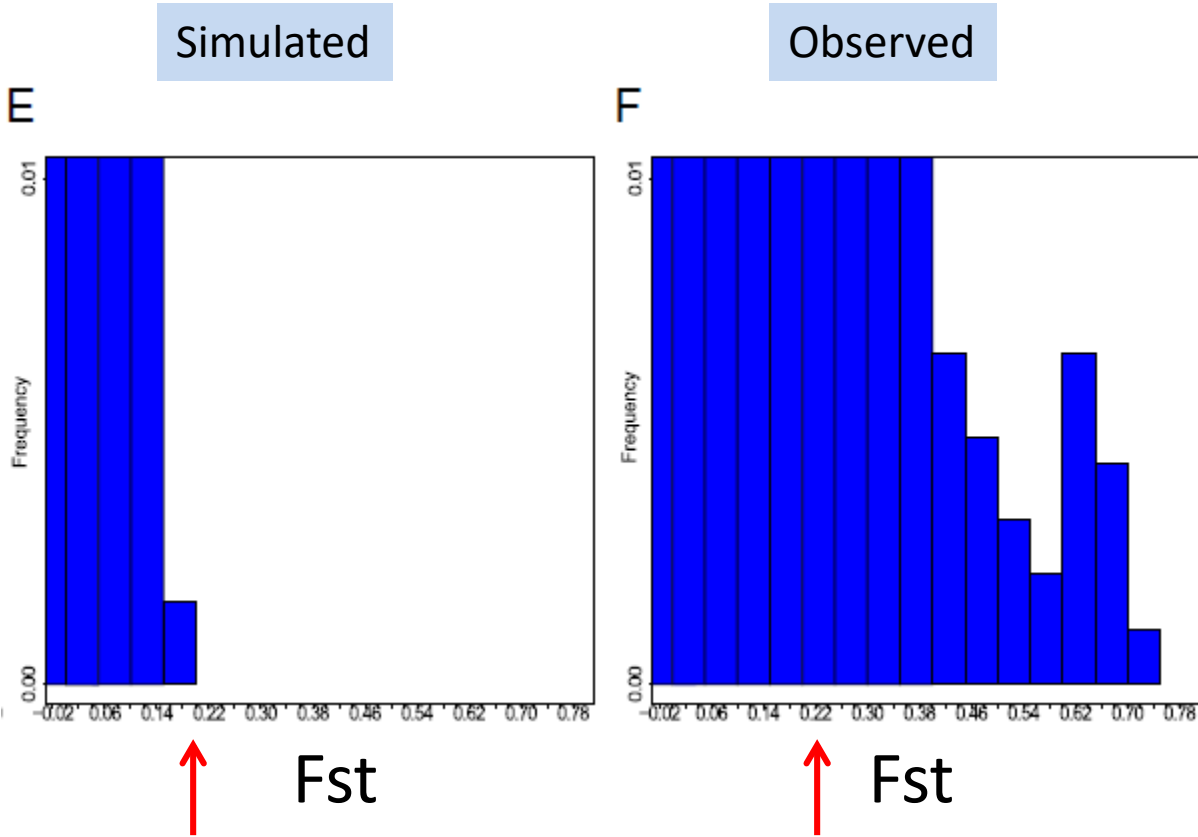
Simulation analysis to test



Unlikely drift → therefore selection!  
local adaptation

# Are high Fst SNPs from drift or selection?

Simulation analysis to test



Unlikely drift → therefore selection!  
local adaptation

Now need to find causative genes and mechanism

# Domestication – artificial selection → Wolf vs Dog



By Retron - self-made now, Public Domain,  
<https://commons.wikimedia.org/w/index.php?curid=3865957>



Charles M. Schulz

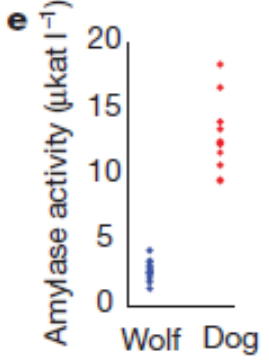
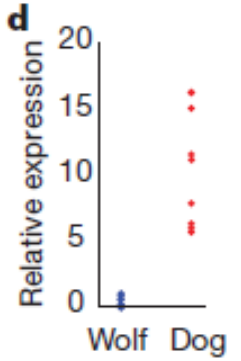
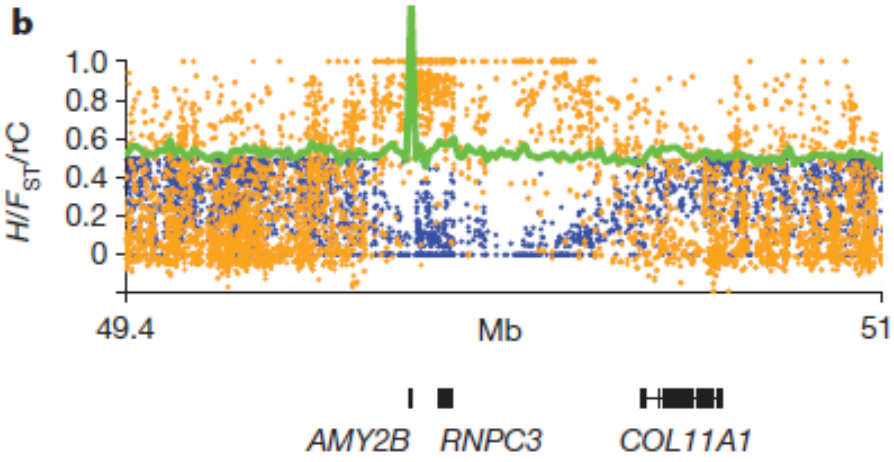
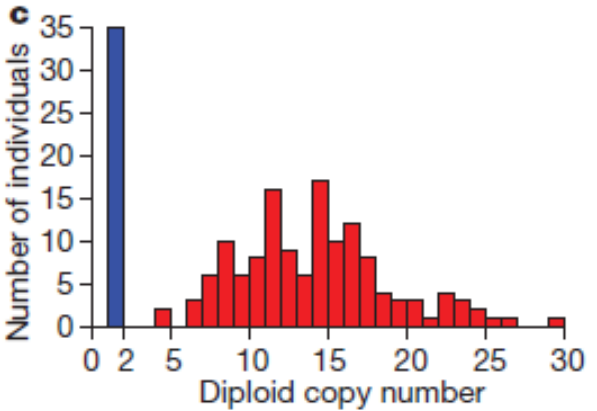
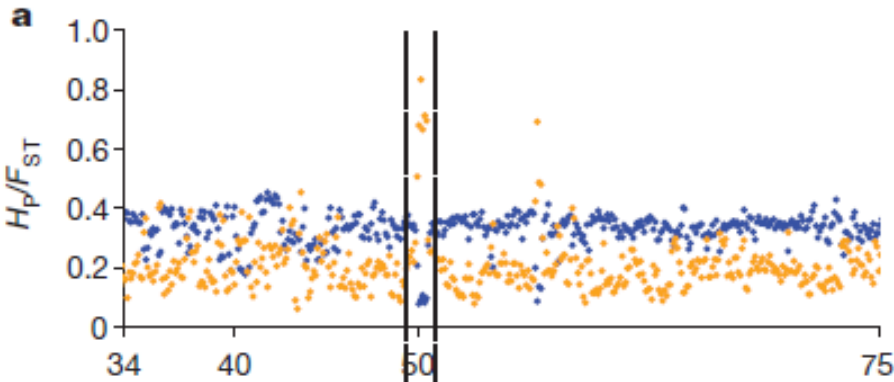




# Dog adaptation to starch diet

36 candidate domestication loci here → Amylase

Fst



also maltase activity



# A basic science question a very population genetics question

## Recombination analysis using Pool-seq



Photo James Gaither

*Mimulus guttatus*  
(Common yellow monkeyflower)

# Recombination analysis using Pool-seq



Earlier, mentioned that Pool-seq not so good for recombination/LD studies because **too short**

However, monkeyflower has high nt diversity (2.9%)  
i.e., within 1 Illumina read, there can be  $\geq 2$  nt differences



# Samples and analysis sketch



98 samples

4 localities

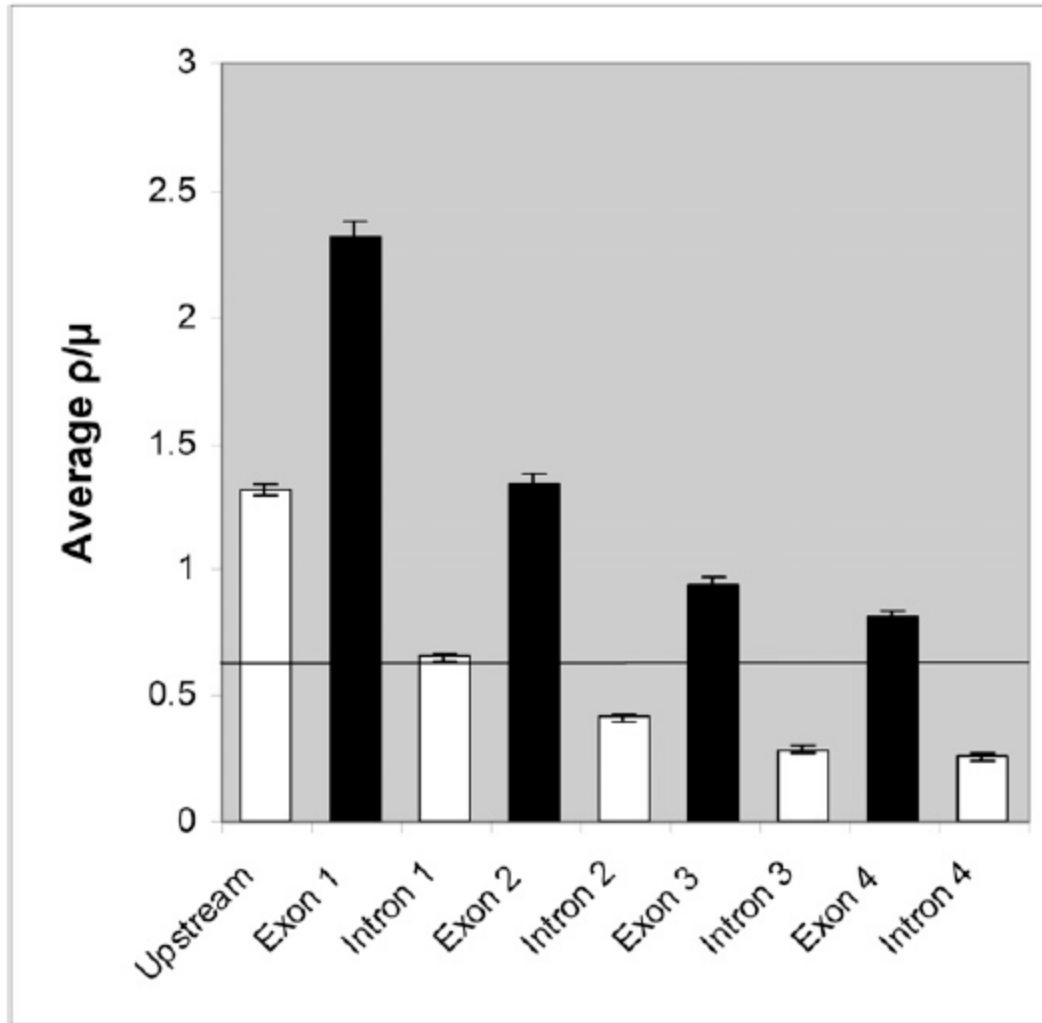
4 Pool-seq runs

PE 75 bp runs (of ~200 bp fragments)

After QC → 255x coverage

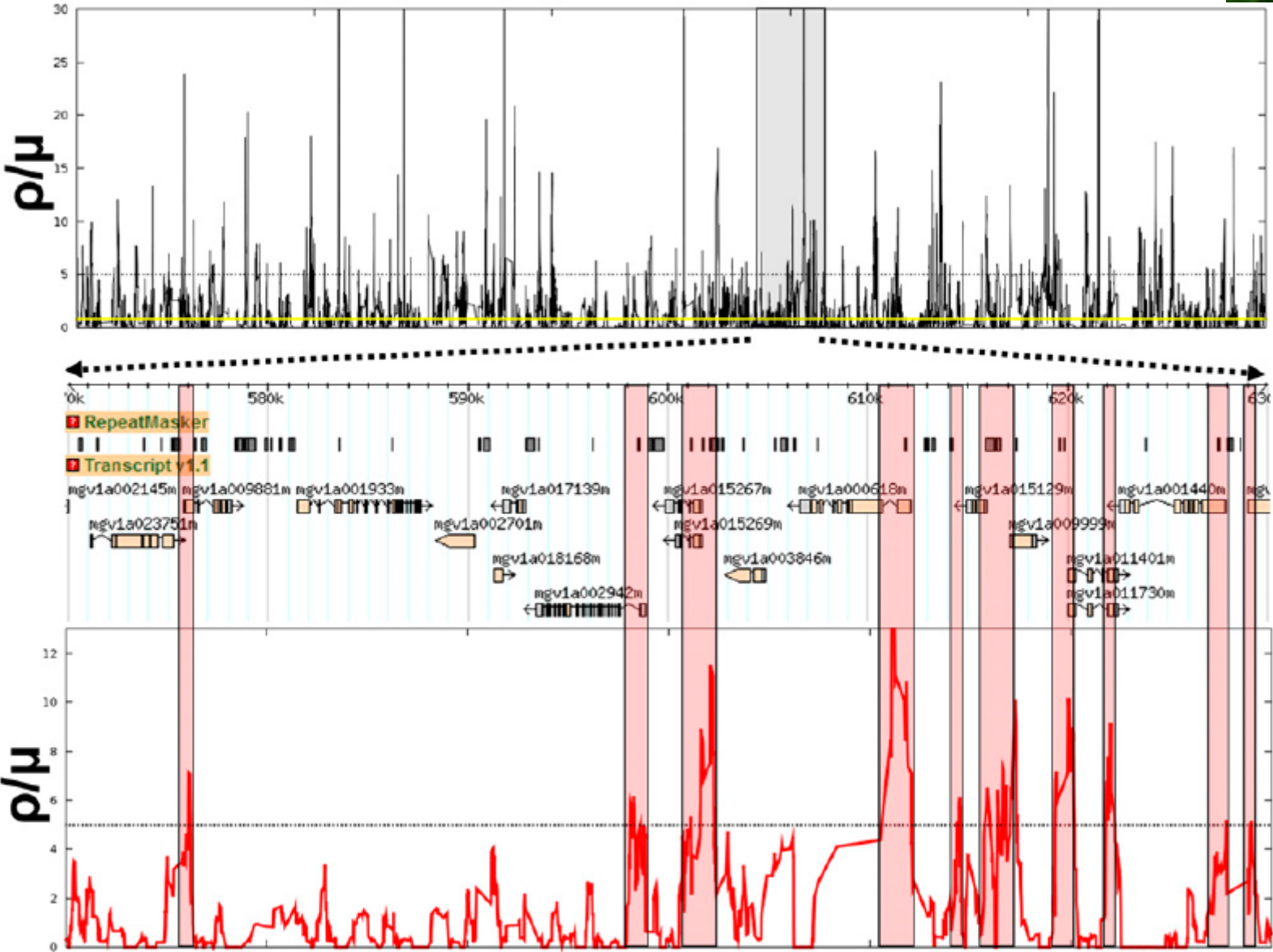
Examined pairs of SNPs within 50 bp

# Interesting recombination patterns



- Exons > introns
- Polarity  
5' > 3' recomb. rate

# Recombination landscape example



# Whole genome (re)sequencing

Can do lots of stuff  
Here, demographic history,  
Scan for regions of diversity,  
Identify gene for phenotype

# What can I do with this genetic data?

If you have allozymes, AFLPs, sequences, microsatellites, or SNPs, you can

## Examine population structure

- PCA or clustering algorithms
- Relatedness/kinship
- Pairwise divergence ( $F_{ST}$ , Jost's D)
- Phylogeography

## Describe diversity

- Hardy-Weinberg deviation
- Allelic richness
- Nucleotide diversity ( $\pi$ ,  $\theta_W$ )
- Frequency spectra
- LD between markers

## Reconstruct demography

- Migration rates
- Population sizes

Did you genotype families?

Do you have data from multiple species?

## Make a linkage map

Do you have ecological or trait data?

## Reconstruct history

- Phylogeny estimate
- Polarize ancestral/derived alleles

## Scan for differentiation

- $F_{ST}$  outliers
- $F_{ST}$  vs. homozygosity

## Find isolating factors

- Mantel testing
- Resistance surface
- BEDASSLE

## Reconstruct evolution

- Ancestral states
- Diversification rates

## Examine genomic variation

- Nucleotide/haplotype diversity
- Runs of homozygosity
- Extended differentiation
- Extended LD

Are genes identified?

## Look at coding variation

- Codon position
- Variant effect prediction

## Test for genotype associations



# Panda demographic history

34 wild panda + ref seq

Pairwise sequentially Markovian coalescent (PSMC)

Can infer past history from 1 or few individuals





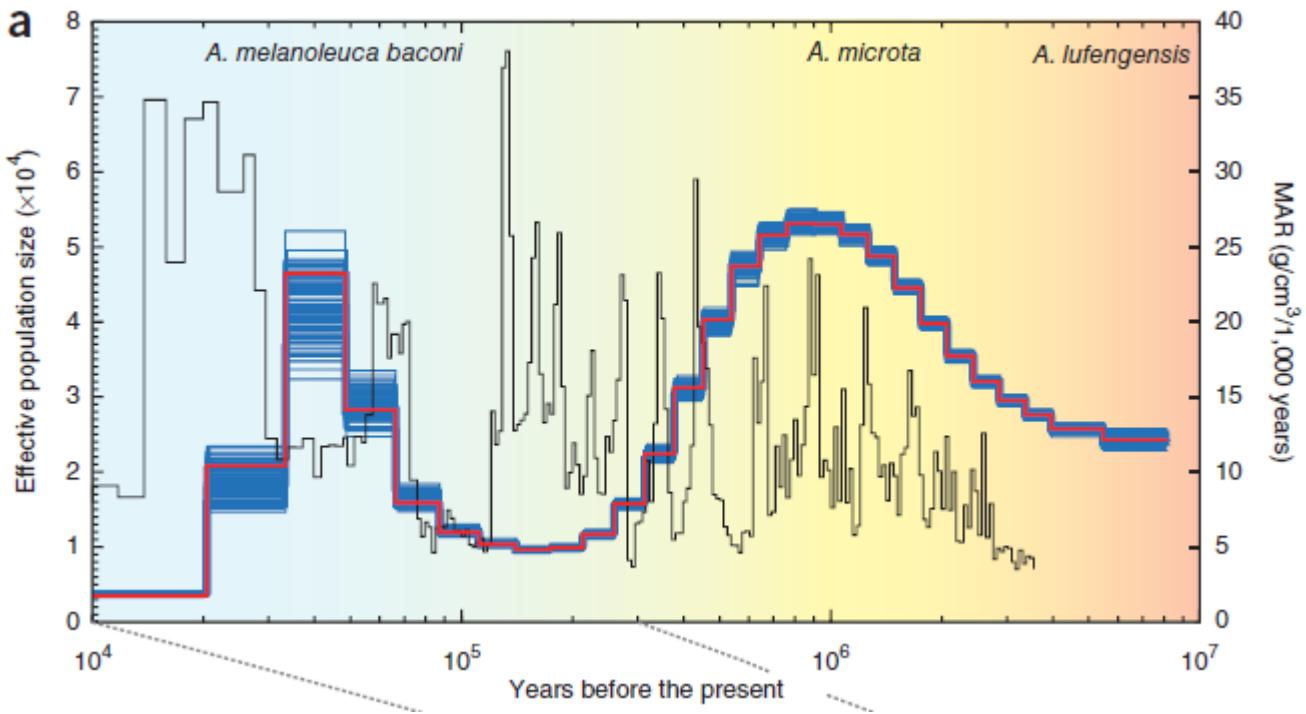
# Panda demographic history



34 wild panda + ref seq

Pairwise sequentially Markovian coalescent (PSMC)

Can infer past history from 1 or few individuals

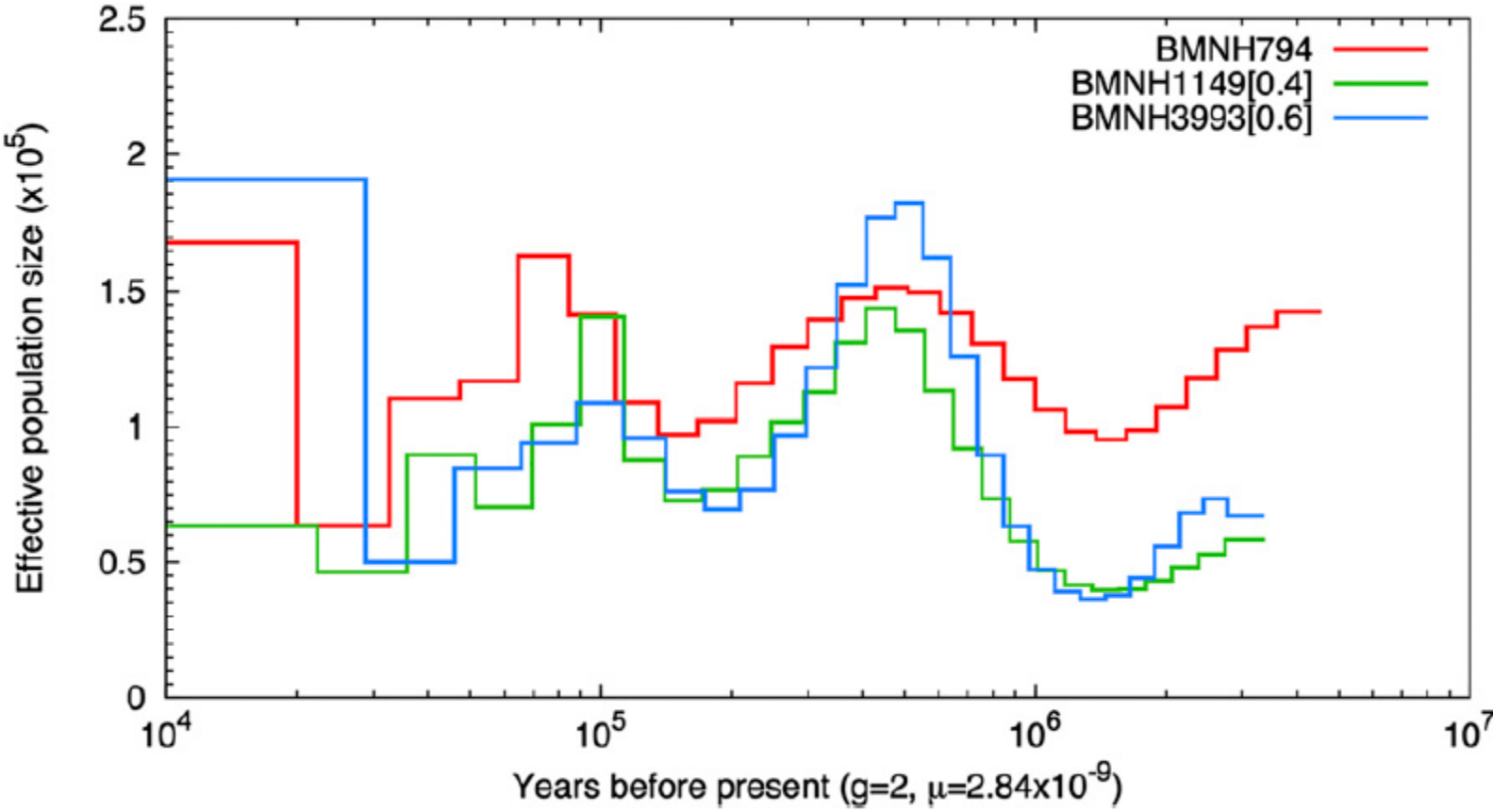


# Passenger pigeon demographic history



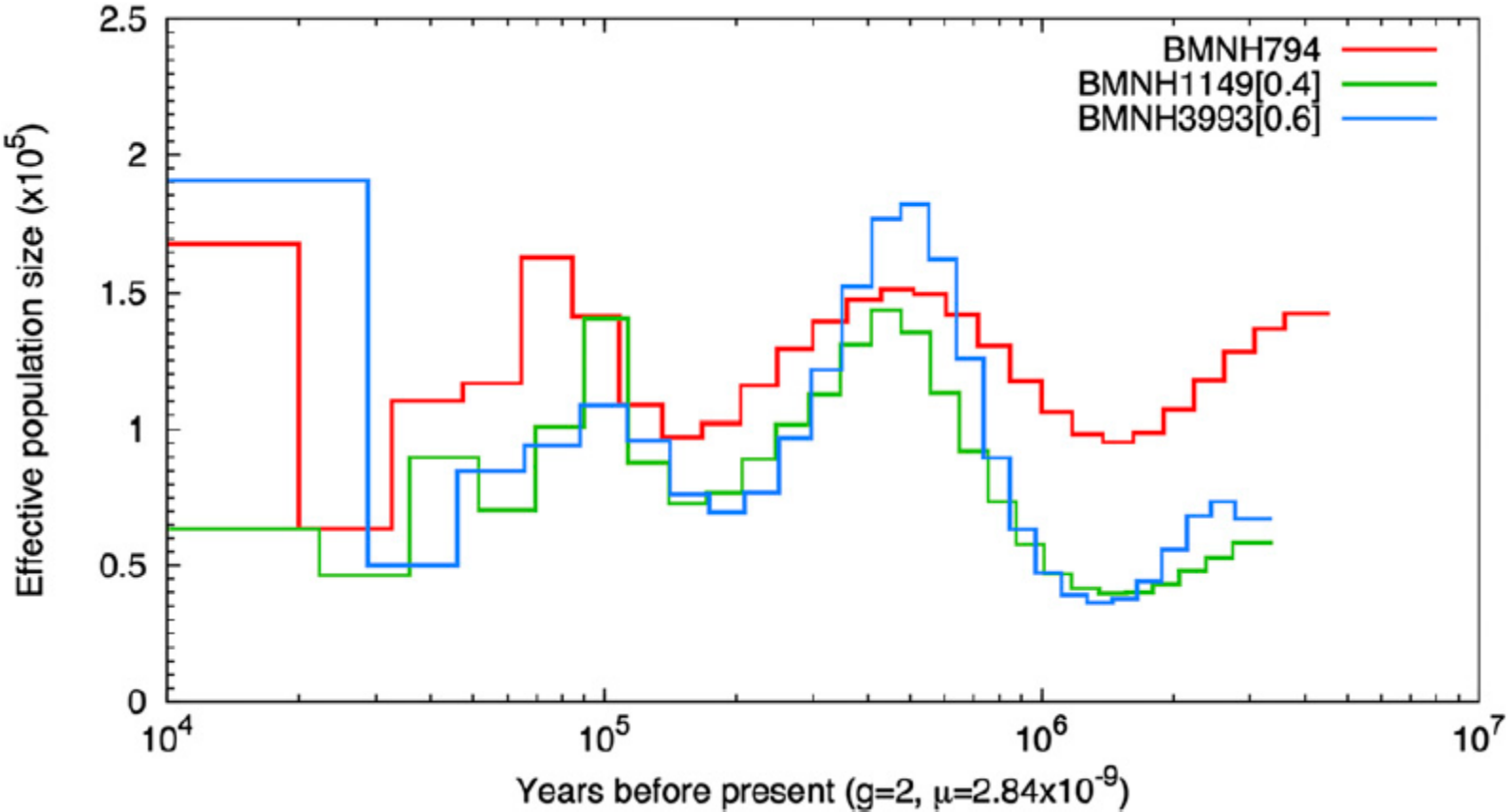
# Passenger pigeon demographic history (+ museum DNA)

N=3 from museum, PSMC



# Passenger pigeon demographic history (+ museum DNA)

N=3 from museum, PSMC



Ne (effective population size; 100,000's) << census N (billions)  
Past fluctuations → climatic, food-resource [acorn], and other ecological variations  
Increased extinction risk



# Speciation islands? (divergence islands)



collared



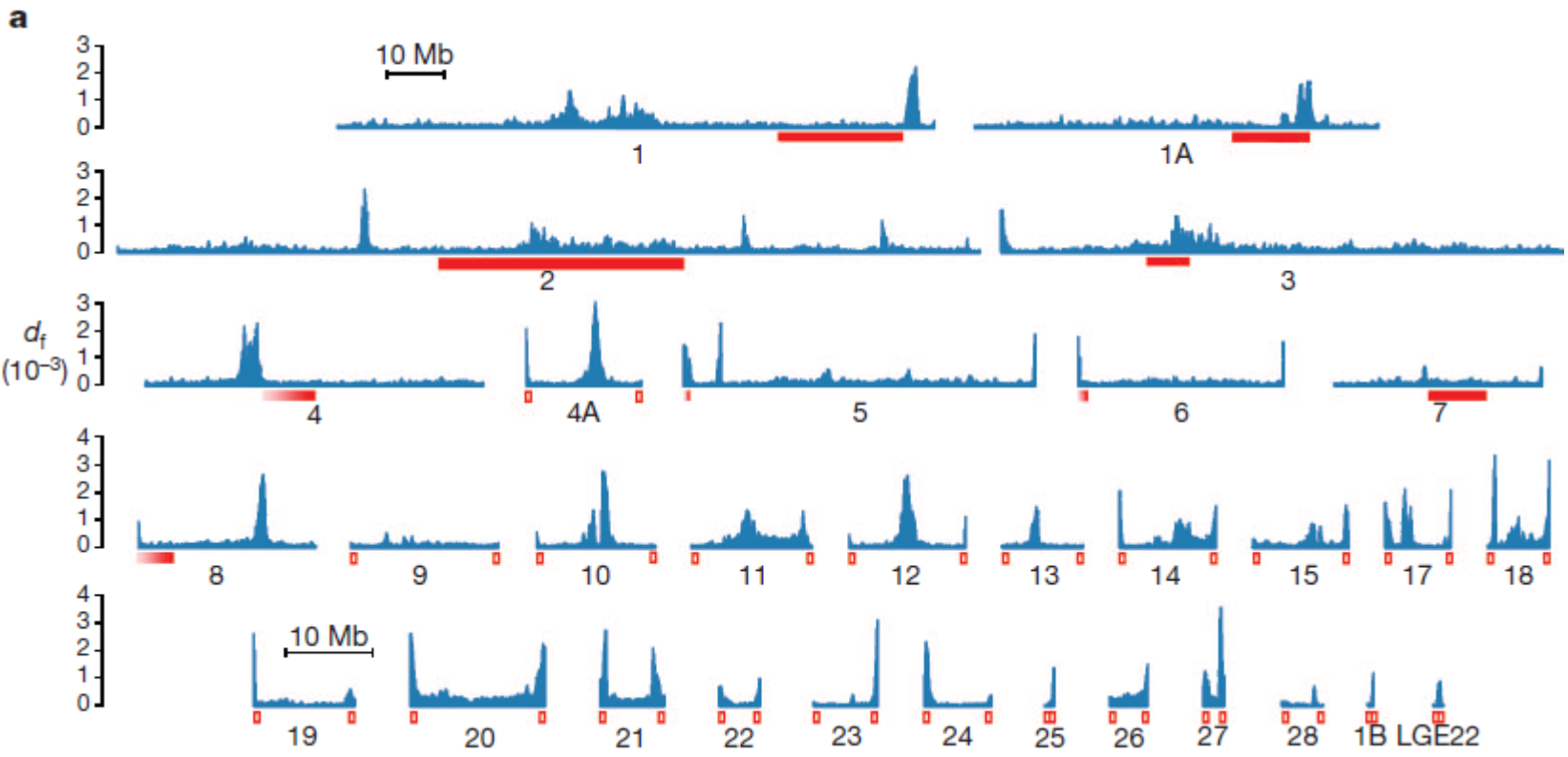
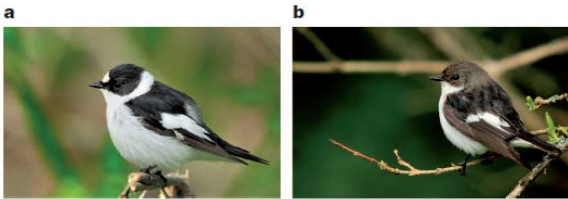
pied

flycatchers  
diverged < 2Mya  
separated, 2° contact  
hybridize

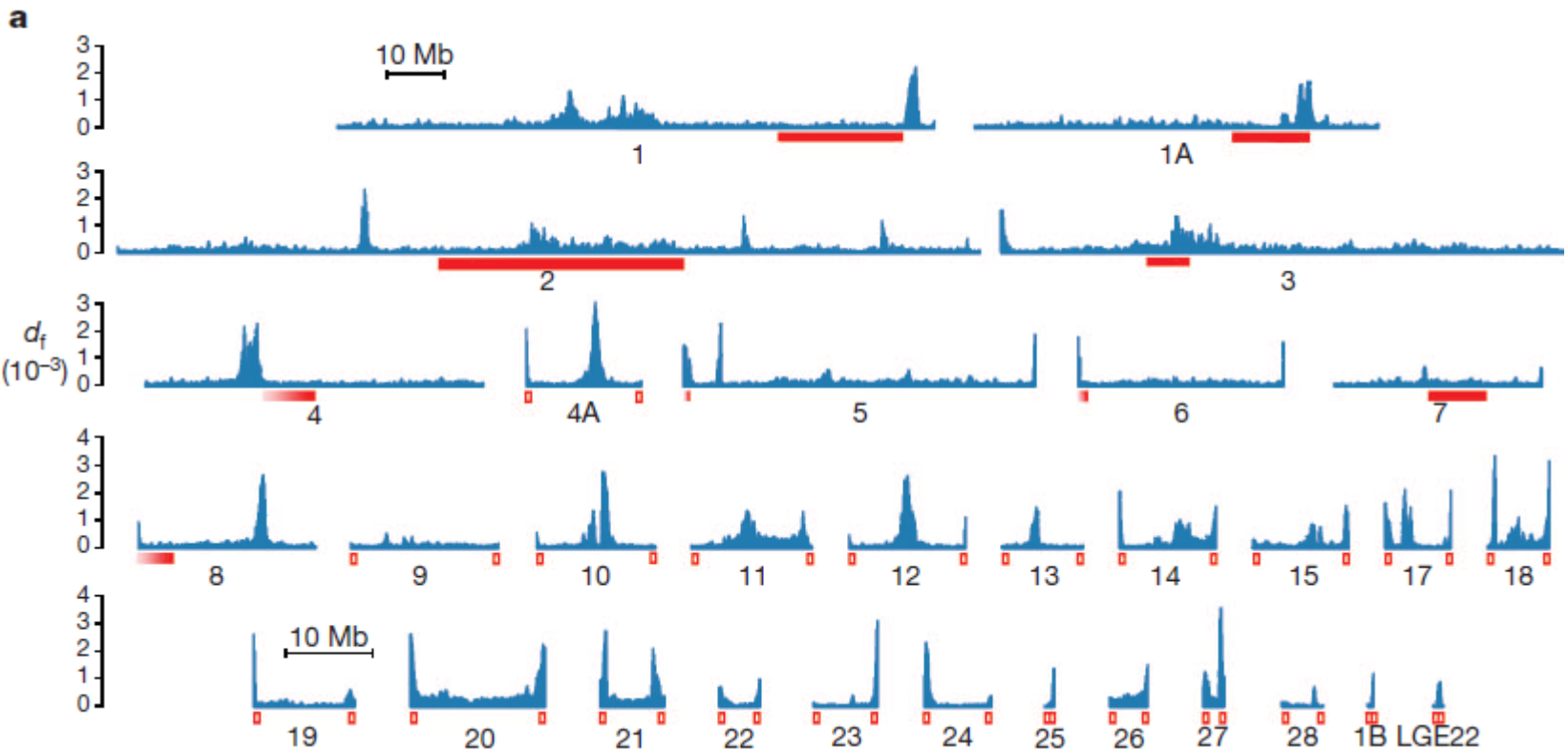
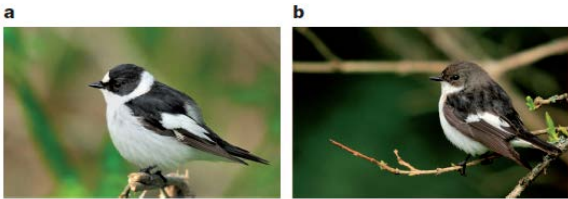
10 re-sequenced individuals + 1 ref seqs (collared)

Compared genomewide divergence to each other

# Divergence islands often near centromeres and telomeres



# Divergence islands often near centromeres and telomeres



Authors postulate meiotic drive

a



b

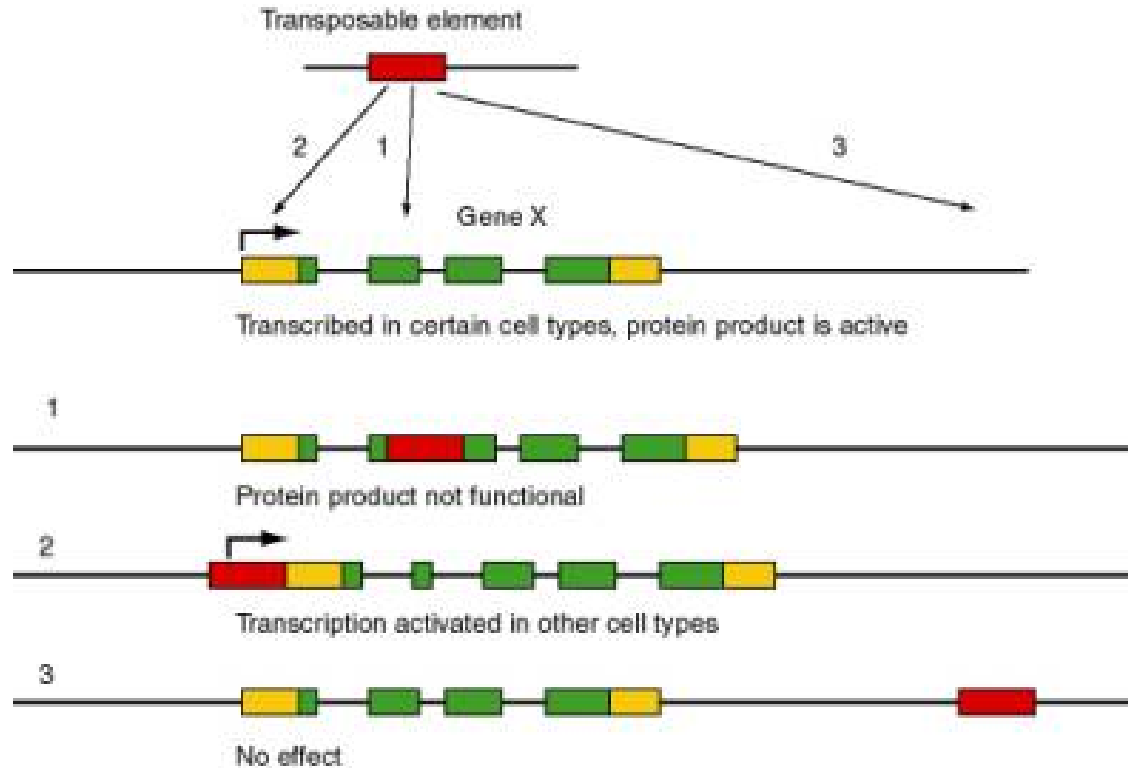


Inspect for transposable elements



# Transposable elements (TEs) are selfish “jumping genes”

Neutral marker or evolutionarily important



Also:

Alternative splicing, methylation, ectopic recombination, etc

# TEs are an important source of genomic variation

## Humans

Alu element → ~800 insertion differences between 2 specific individuals

Estimated new insertion rates:

Alu → 1 in 21 people

L1 → 1 in 108 to 212

## Maize

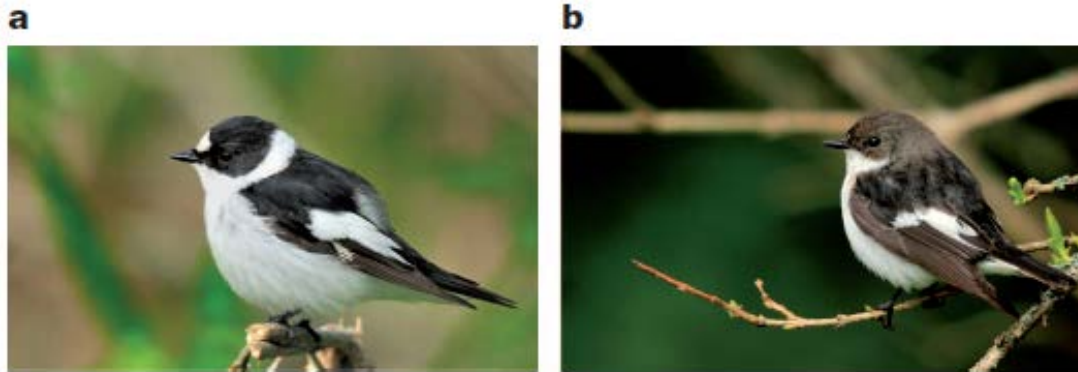
>177,000 TEs in maize (B73 reference strain)

B73 versus W22 strain → >22,000 TEs at different locations

In fact, >10,000 **strain-specific** subfamilies of TEs

# Birds have been suggested to have low TE densities

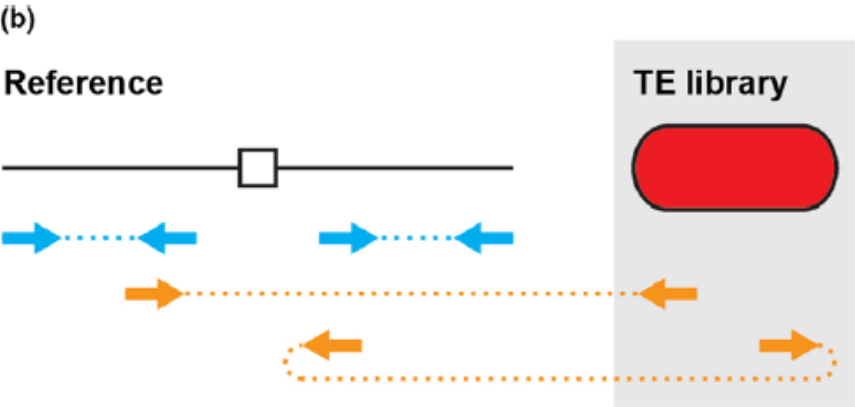
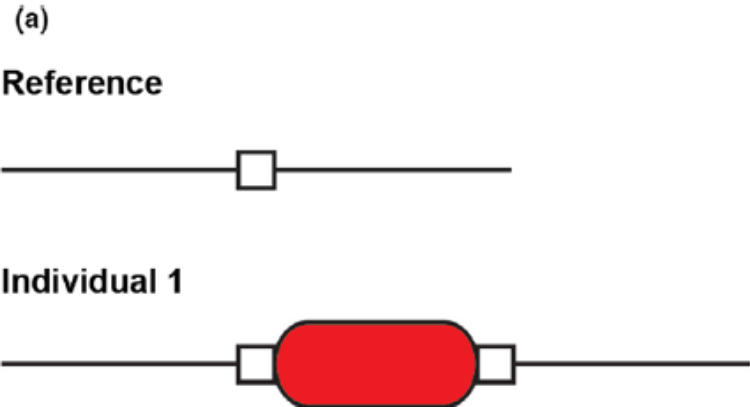
Implying that TEs have not transposed much in bird history



Suh et al looked at:

- TE repertoire in flycatcher genome
- Then polymorphism in 200 resequenced genomes

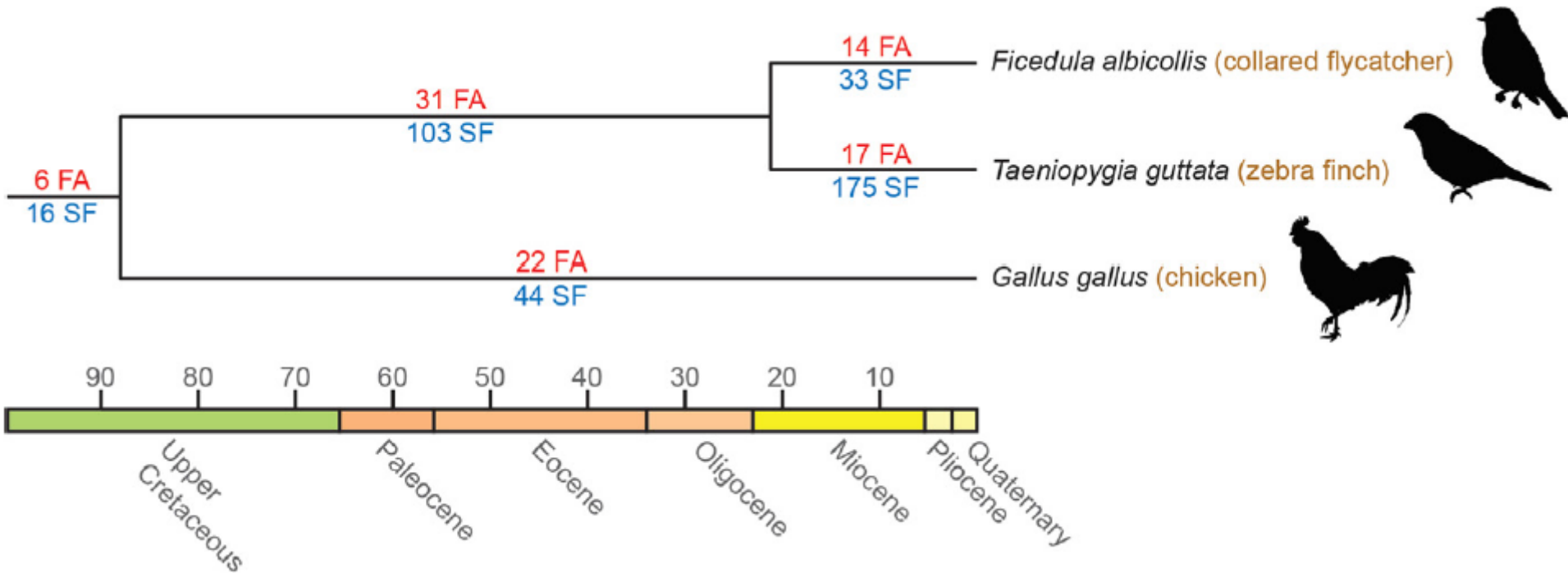
# NGS to locate TEs



*(There's also a "reduced" method using primers from the TEs)*

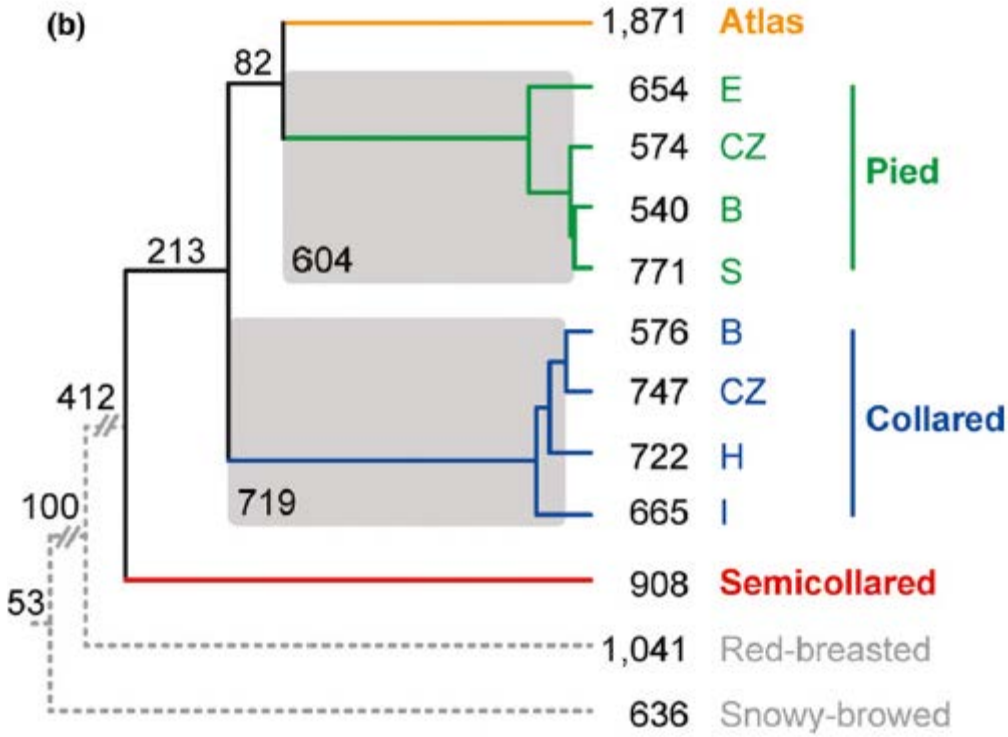
# Novel transposon LTR families in 3 birds

i.e., These are NEW to science



# TE variation in flycatchers

Retrotransposition events per lineage



# Flycatcher case study conclusions



# Flycatcher case study conclusions

Flycatchers (songbirds) actually have a lot of transposition events





# Flycatcher case study conclusions

Flycatchers (songbirds) actually have a lot of transposition events

New TE families and current polymorphism



# Flycatcher case study conclusions

Flycatchers (songbirds) actually have a lot of transposition events

New TE families and current polymorphism

These may have phenotypic consequences for their diversification



# Flycatcher case study conclusions

Flycatchers (songbirds) actually have a lot of transposition events

New TE families and current polymorphism

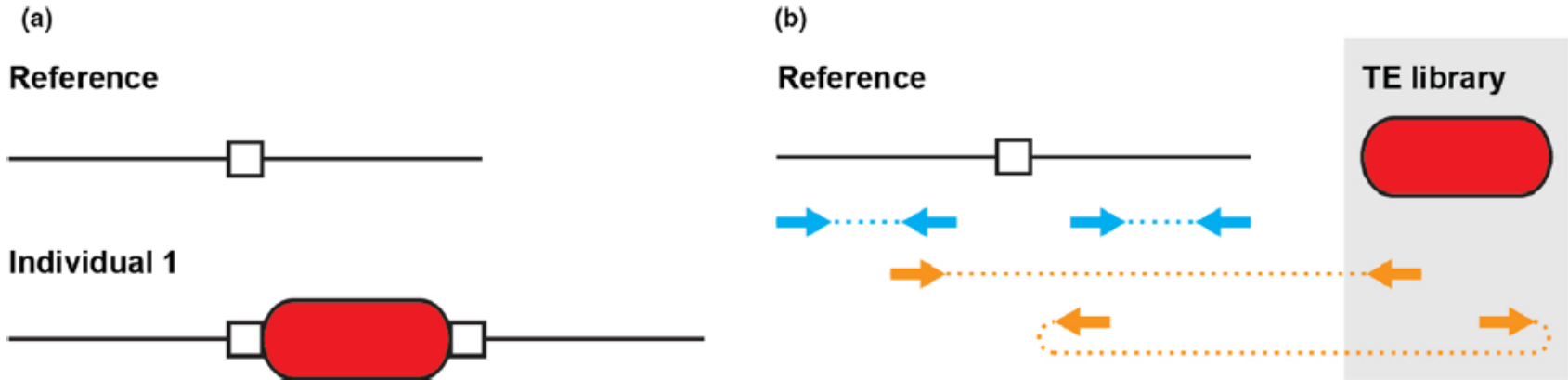
These may have phenotypic consequences for their diversification



In general TEs seem to be underappreciated for their impact

# Paired-end method has drawbacks

TEs can be 1000's bp long



Short reads have higher **False Discovery Rates**

Need lots of depth

Somewhat inferred (e.g., don't know the internal sequence)

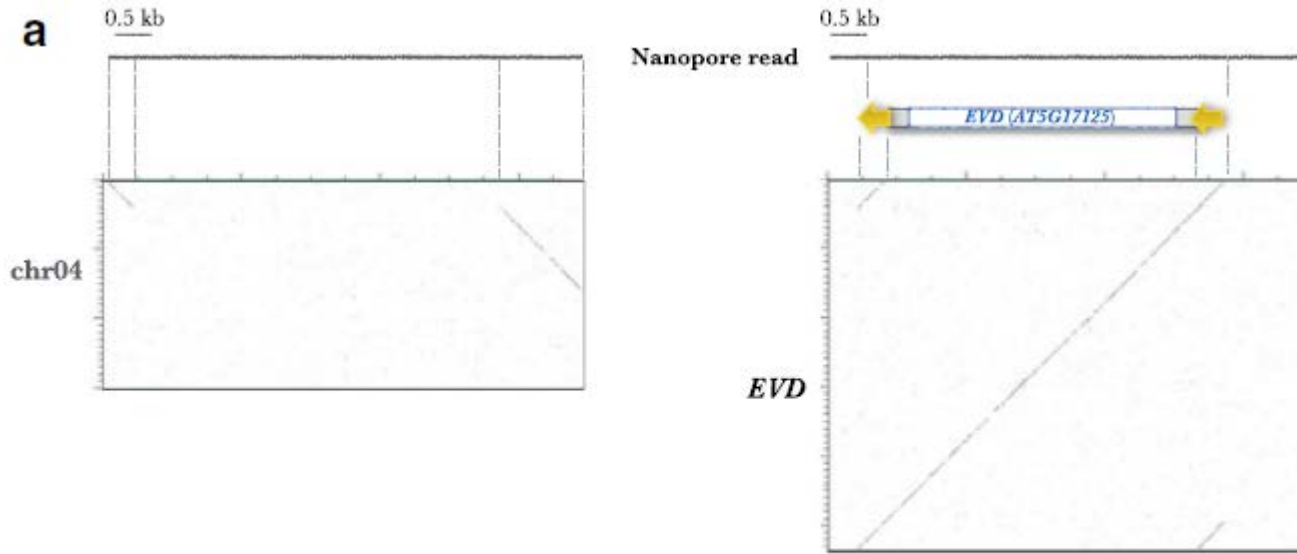
etc

# Long reads (PacBio, ONT) can just read across TE insertions

TEs can be 1000's bp long

ONT analysis in Arabidopsis

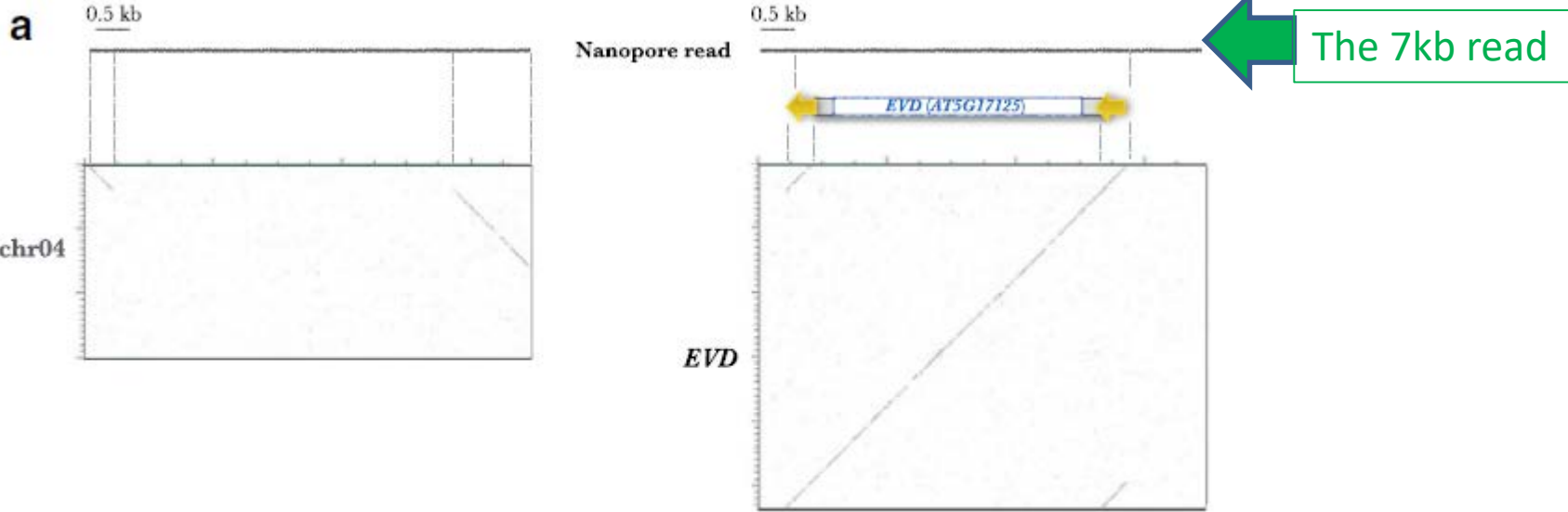
Example of a ~7kb that crosses 1 TE



# Long reads (PacBio, ONT) can just read across TE insertions

TEs can be 1000's bp long

ONT analysis in Arabidopsis  
Example of a ~7kb that crosses 1 TE

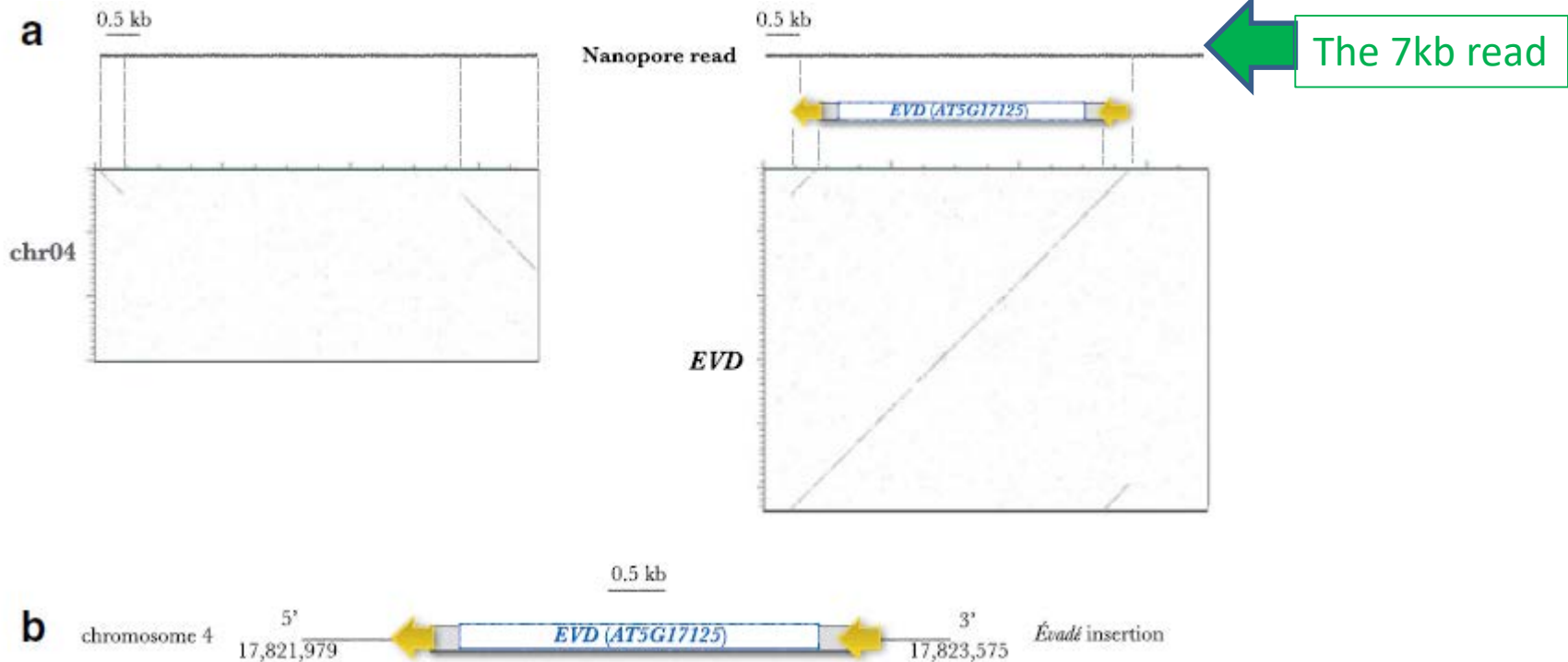


# Long reads (PacBio, ONT) can just read across TE insertions

TEs can be 1000's bp long

ONT analysis in Arabidopsis

Example of a ~7kb that crosses 1 TE

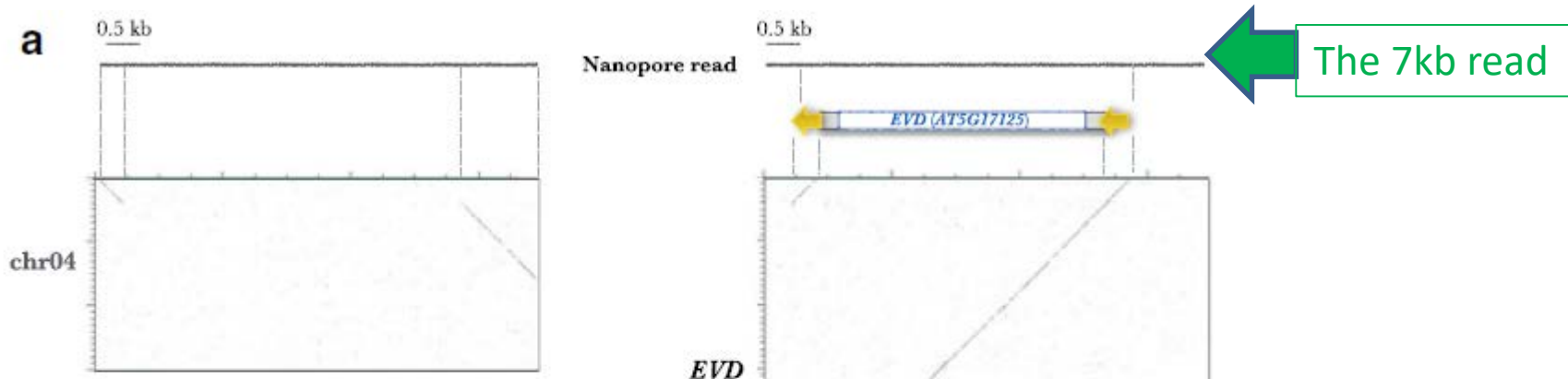


# Long reads (PacBio, ONT) can just read across TE insertions

TEs can be 1000's bp long

ONT analysis in Arabidopsis

Example of a ~7kb that crosses 1 TE



Less depth needed, sort of

Does tell you if “whole/intact” or “partial” copy

Goes from algorithm to just “inspecting”

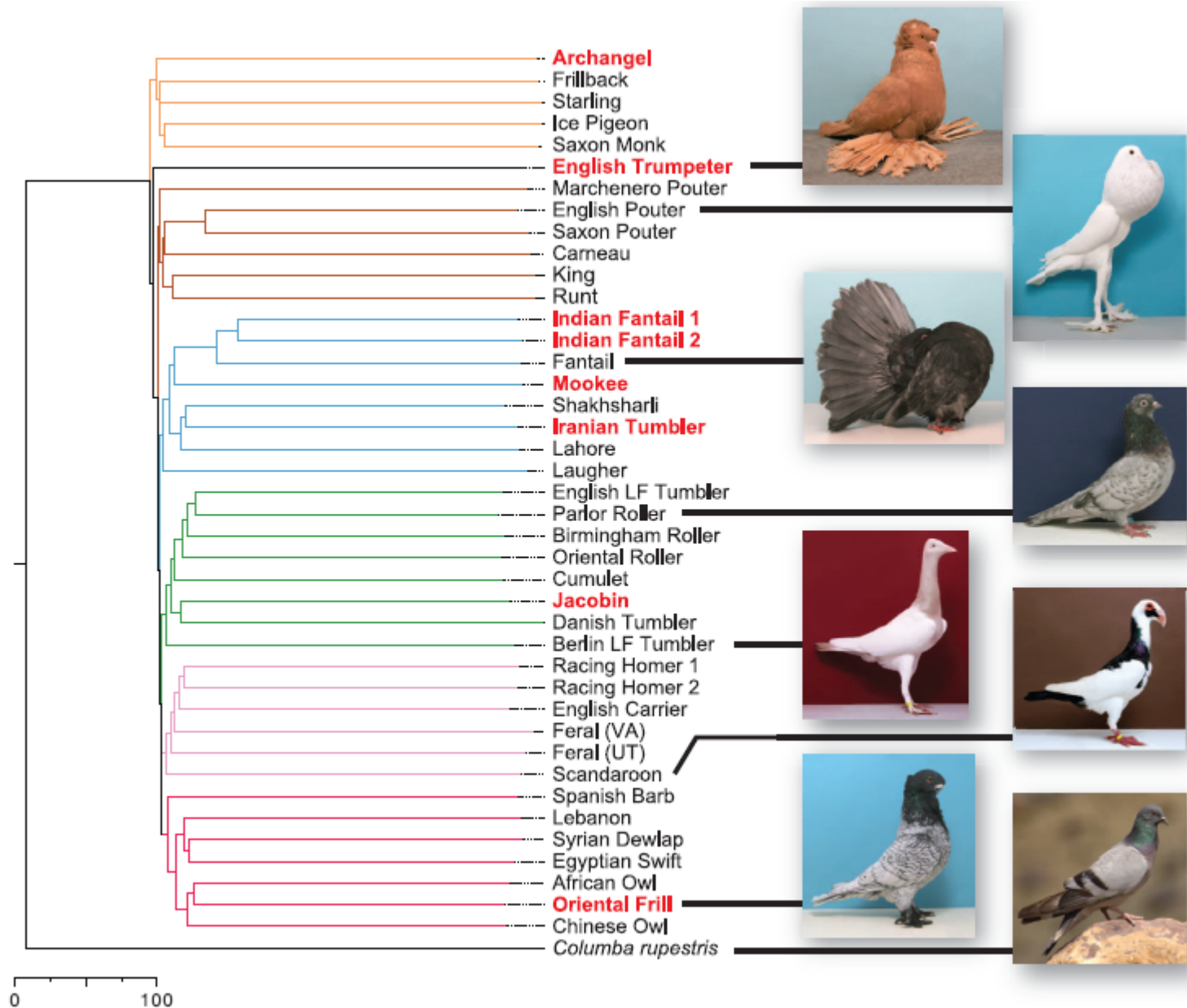


# Genomic Diversity and Evolution of the Head Crest in the Rock Pigeon

Michael D. Shapiro,<sup>1\*</sup> Zev Kronenberg,<sup>2</sup> Cai Li,<sup>3,4</sup> Eric T. Domyan,<sup>1</sup> Hailin Pan,<sup>3</sup>  
Michael Campbell,<sup>2</sup> Hao Tan,<sup>3</sup> Chad D. Huff,<sup>2,5</sup> Haofu Hu,<sup>3</sup> Anna I. Vickrey,<sup>1</sup>  
Sandra C. A. Nielsen,<sup>4</sup> Sydney A. Stringham,<sup>1</sup> Hao Hu,<sup>5</sup> Eske Willerslev,<sup>4</sup>  
M. Thomas P. Gilbert,<sup>4,6</sup> Mark Yandell,<sup>2</sup> Guojie Zhang,<sup>3</sup> Jun Wang<sup>3,7,8\*</sup>

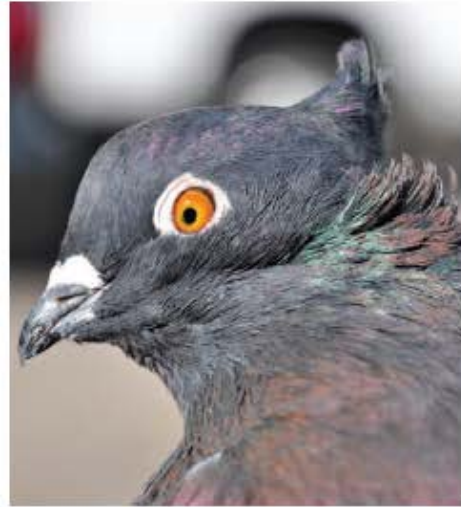
1 ref + 40 additional re-sequenced pigeon genomes

# 1 ref + 40 additional re-sequenced pigeon genomes



# What is the genetic basis for head crests?

Recessive mendelian trait



Peak crest



Shell crest



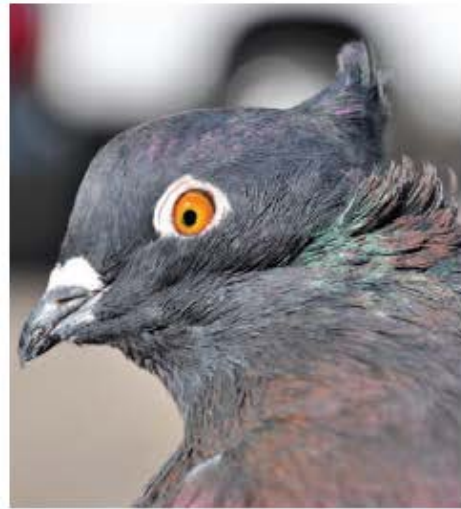
Mane



Hood

# What is the genetic basis for head crests?

Recessive mendelian trait



Peak crest



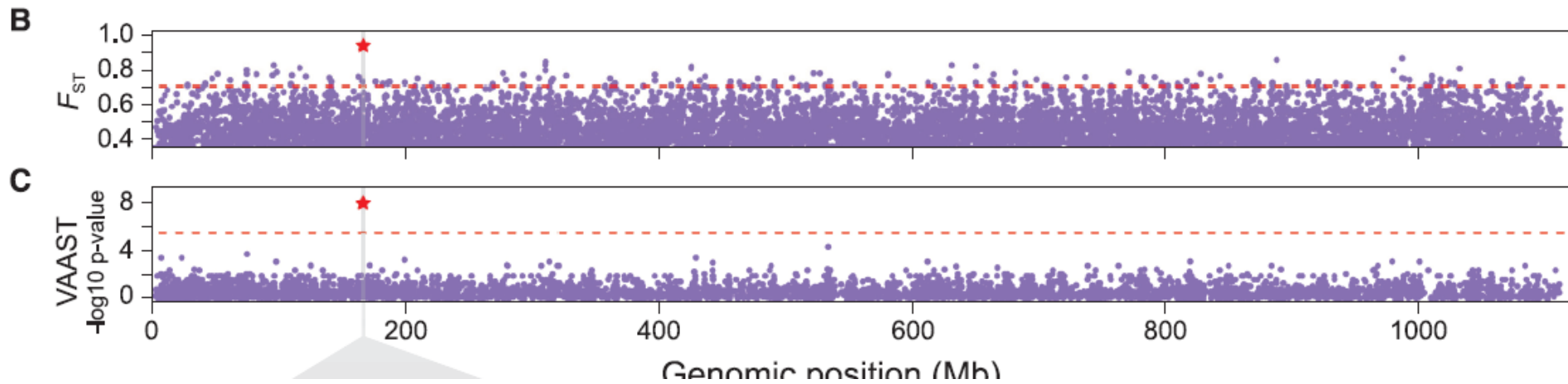
Shell crest



Mane



Hood





# What is the genetic basis for head crests?

Recessive mendelian trait



Peak crest



Shell crest

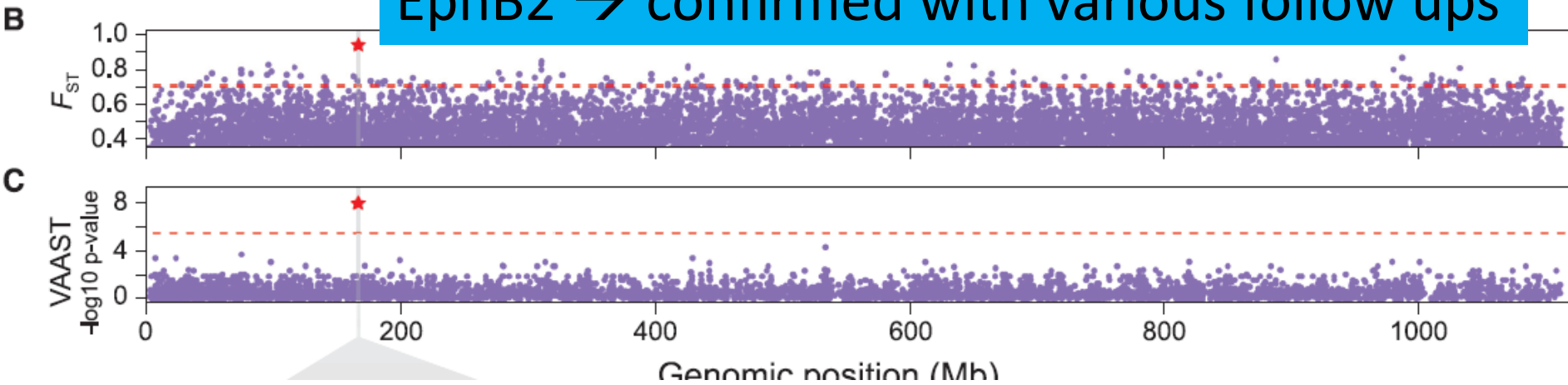


Mane



Hood

EphB2 → confirmed with various follow ups





# Human genetics and health

GWAS  
using SNP arrays

# The UK Biobank resource with deep phenotyping and genomic data

Clare Bycroft<sup>1,13</sup>, Colin Freeman<sup>1,13</sup>, Desislava Petkova<sup>1,12,13</sup>, Gavin Band<sup>1</sup>, Lloyd T. Elliott<sup>2</sup>, Kevin Sharp<sup>2</sup>, Allan Motyer<sup>3</sup>, Damjan Vukcevic<sup>3,4</sup>, Olivier Delaneau<sup>5,6,7</sup>, Jared O'Connell<sup>8</sup>, Adrian Cortes<sup>1,9</sup>, Samantha Welsh<sup>10</sup>, Alan Young<sup>11</sup>, Mark Effingham<sup>10</sup>, Gil McVean<sup>1,11</sup>, Stephen Leslie<sup>3,4</sup>, Naomi Allen<sup>11</sup>, Peter Donnelly<sup>1,2,14</sup> & Jonathan Marchini<sup>1,2,14\*</sup>



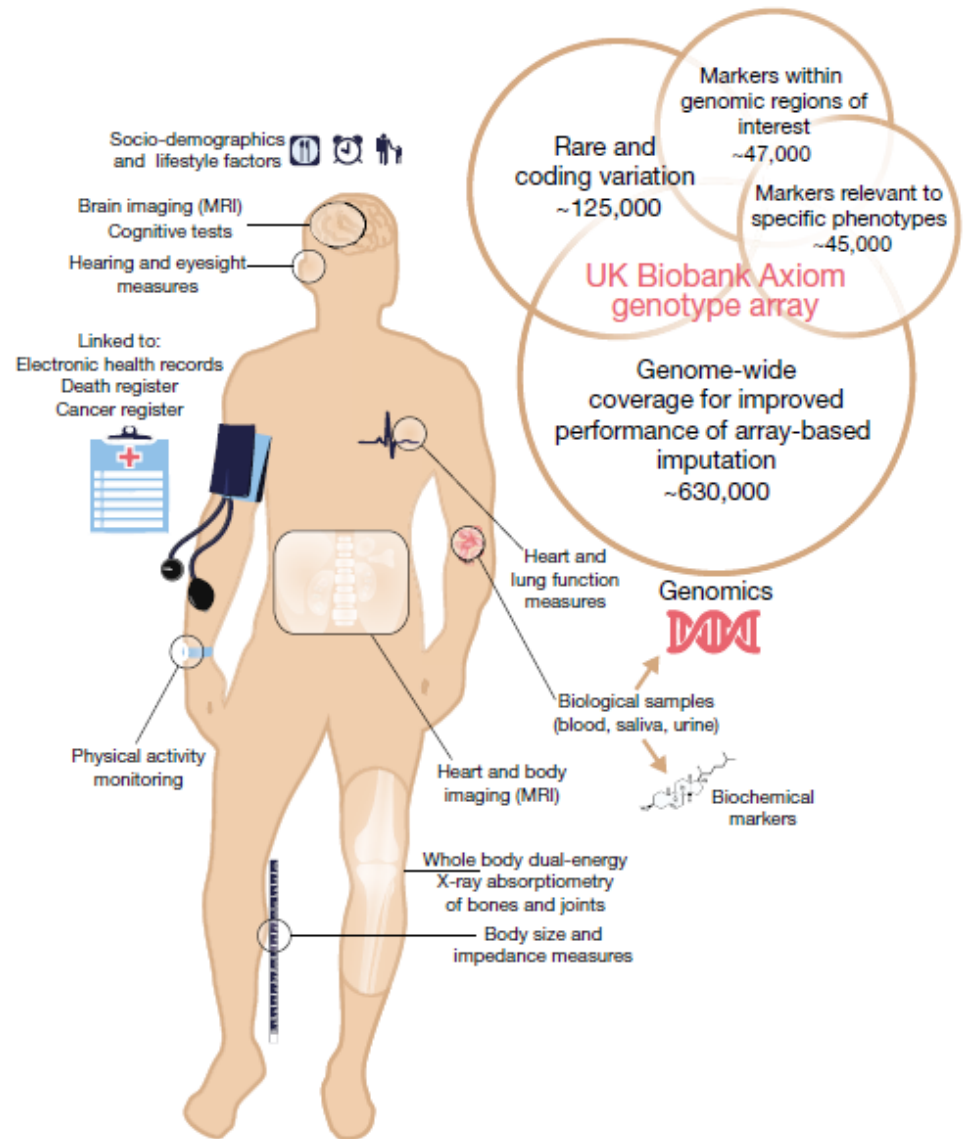
# UK Biobank: overview

~500,000 people

Genotyping

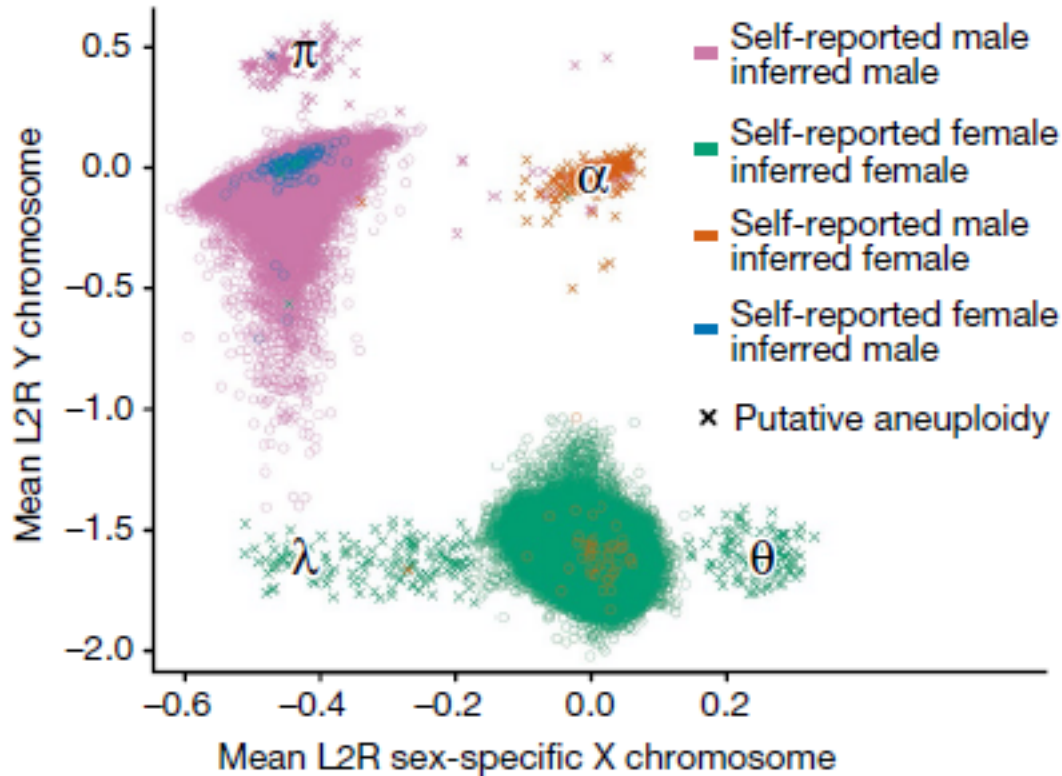
Phenotyping

Prospective cohort study

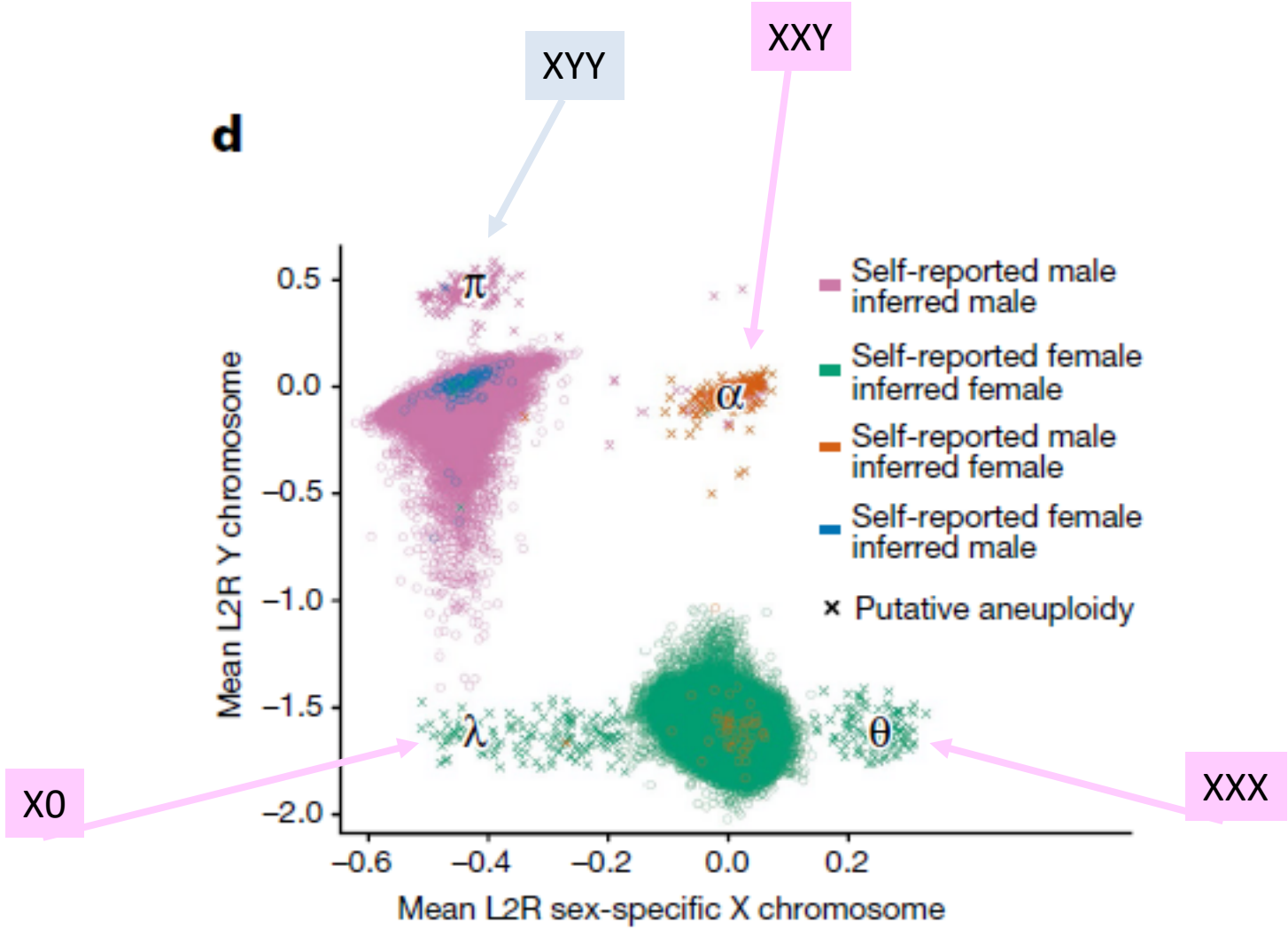


# UK Biobank: non-standard XX/XY sex genotypes

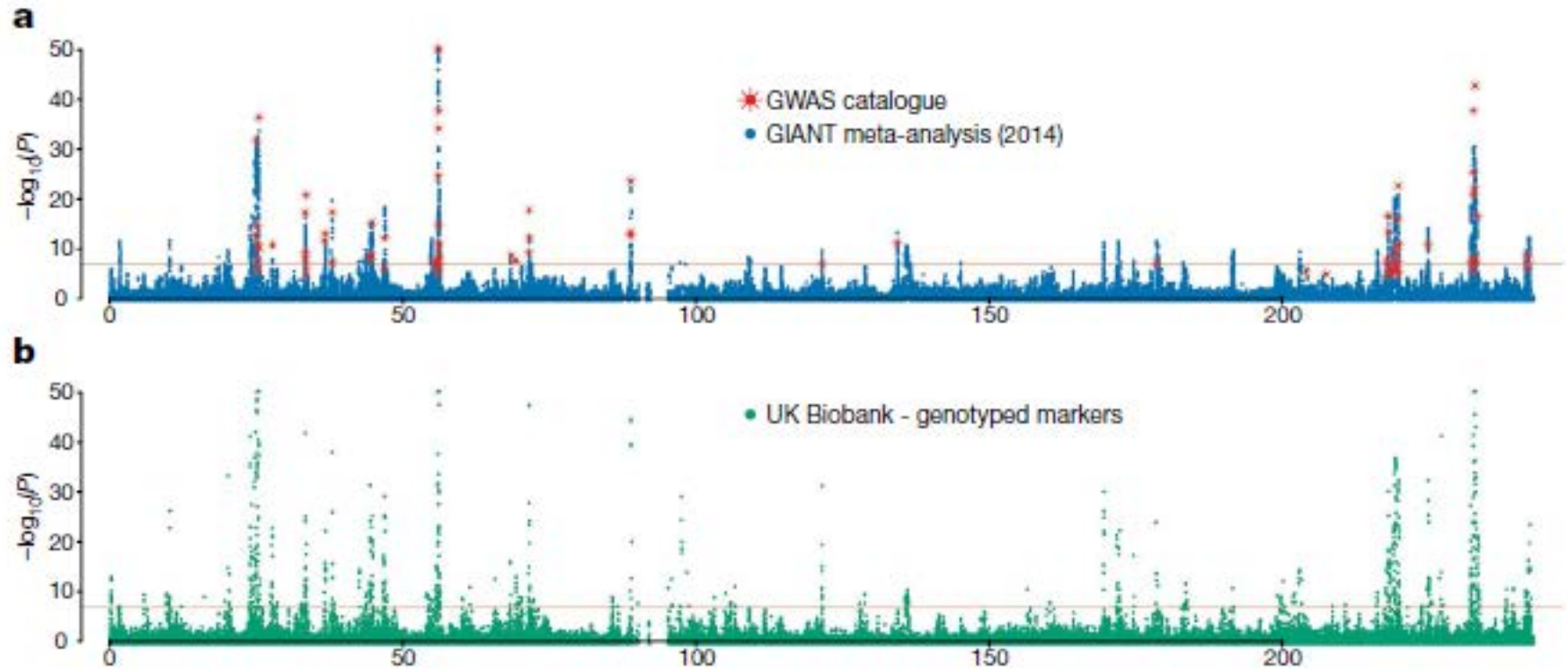
**d**



# UK Biobank: non-standard XX/XY sex genotypes



# UK Biobank: height GWAS, more detection power



# Concluding remarks 1

- Population genomics can be very powerful
- Find the right question (like all of science)
  - Better to know some population genetics
- Studies often stronger when combined with mechanism

## Concluding remarks 2

- There are economical and expensive ways
- Getting cheaper
  - Trend to whole genomes
  - More samples
- No plug and play; no established pipelines
- Finally, hope you get some ideas for your projects



THANKS FOR YOUR ATTENTION



## Further reading

Sequencing pools of individuals - mining genome-wide polymorphism data without big funding.  
Schlötterer C, Tobler R, Kofler R, Nolte V.  
Nat Rev Genet. 2014 Nov;15(11):749-63. doi: 10.1038/nrg3803. Epub 2014 Sep 23. Review.