

Population Genomics - Where are we going? (in 60 minutes....)



Andrew Clark
Cornell University

EMBO short course
Napoli 21 May 2018

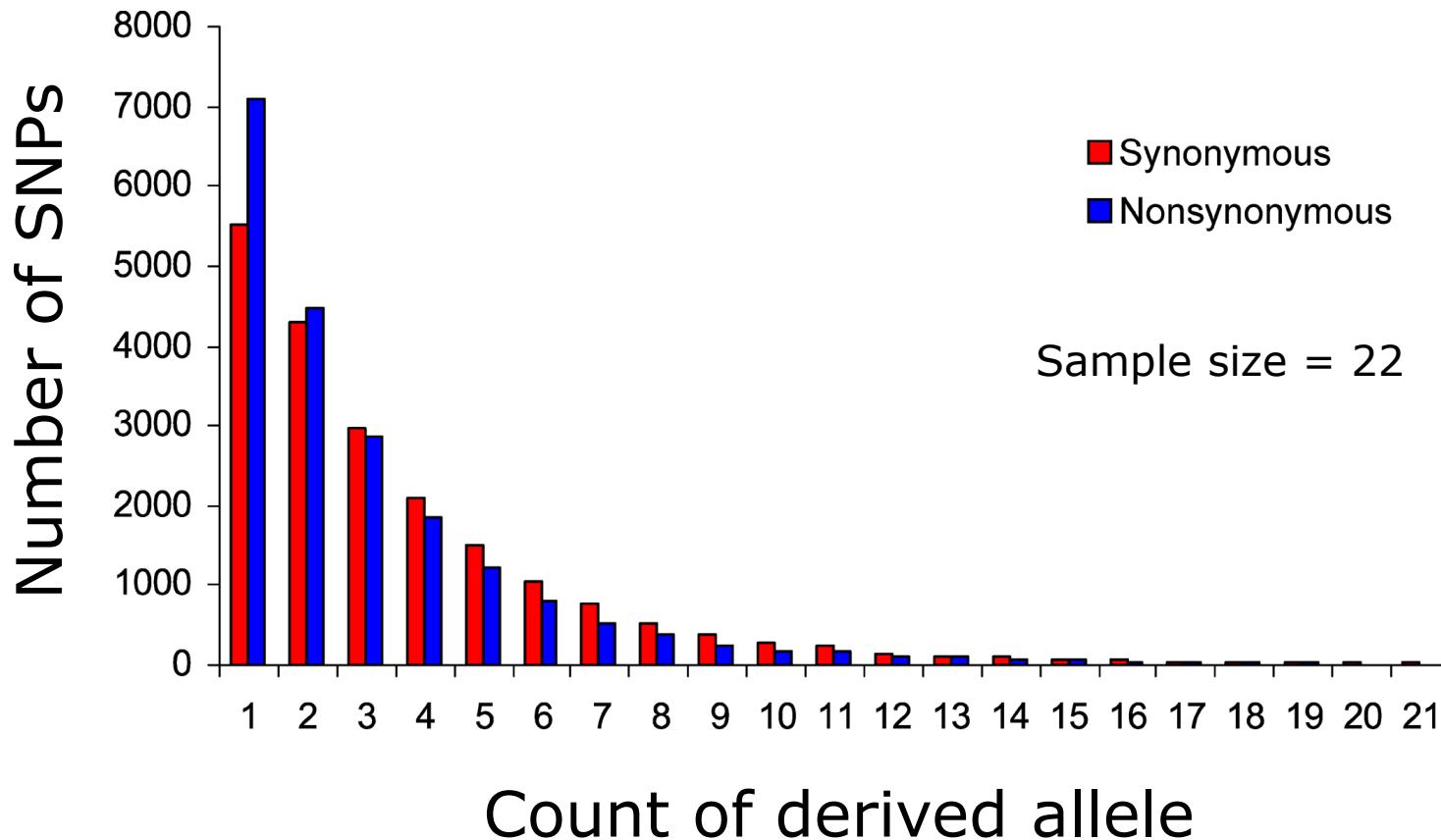
Outline

- Demographic inference
- Population structure and history
- Admixture / Introgression
- Random genetic drift
- Natural selection
- Mutation spectrum
- Disease association
- Genome function

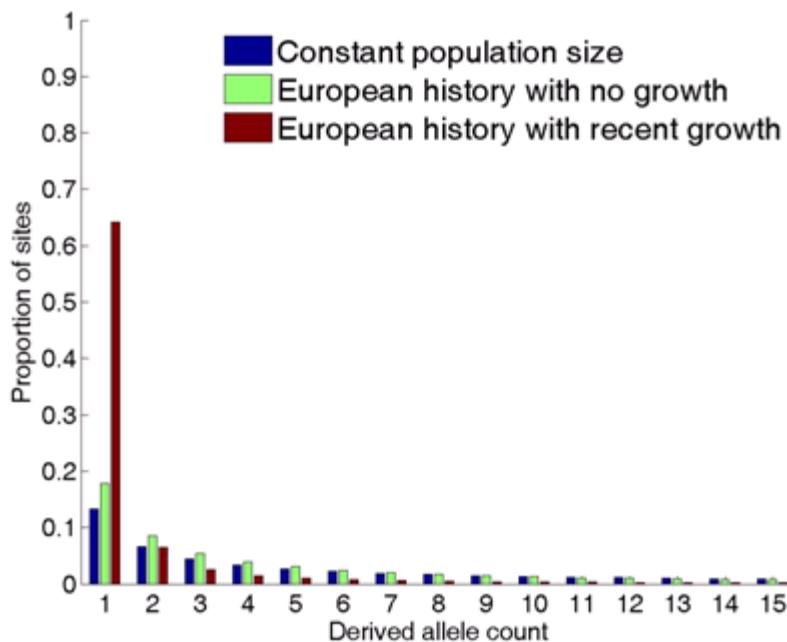
DEMOGRAPHY

How can we infer past changes in population size, bottlenecks, etc.?

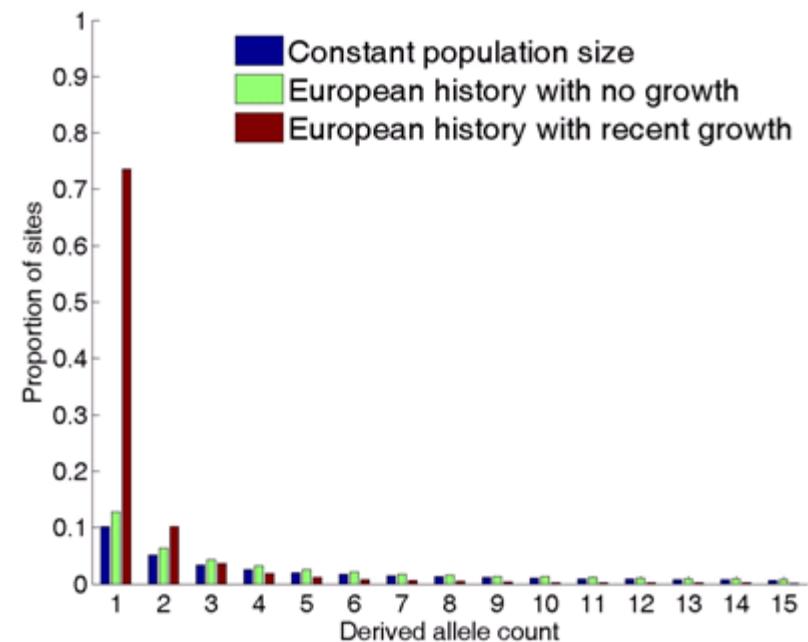
The Site Frequency Spectrum



Large samples are needed to see extent of skew to SFS

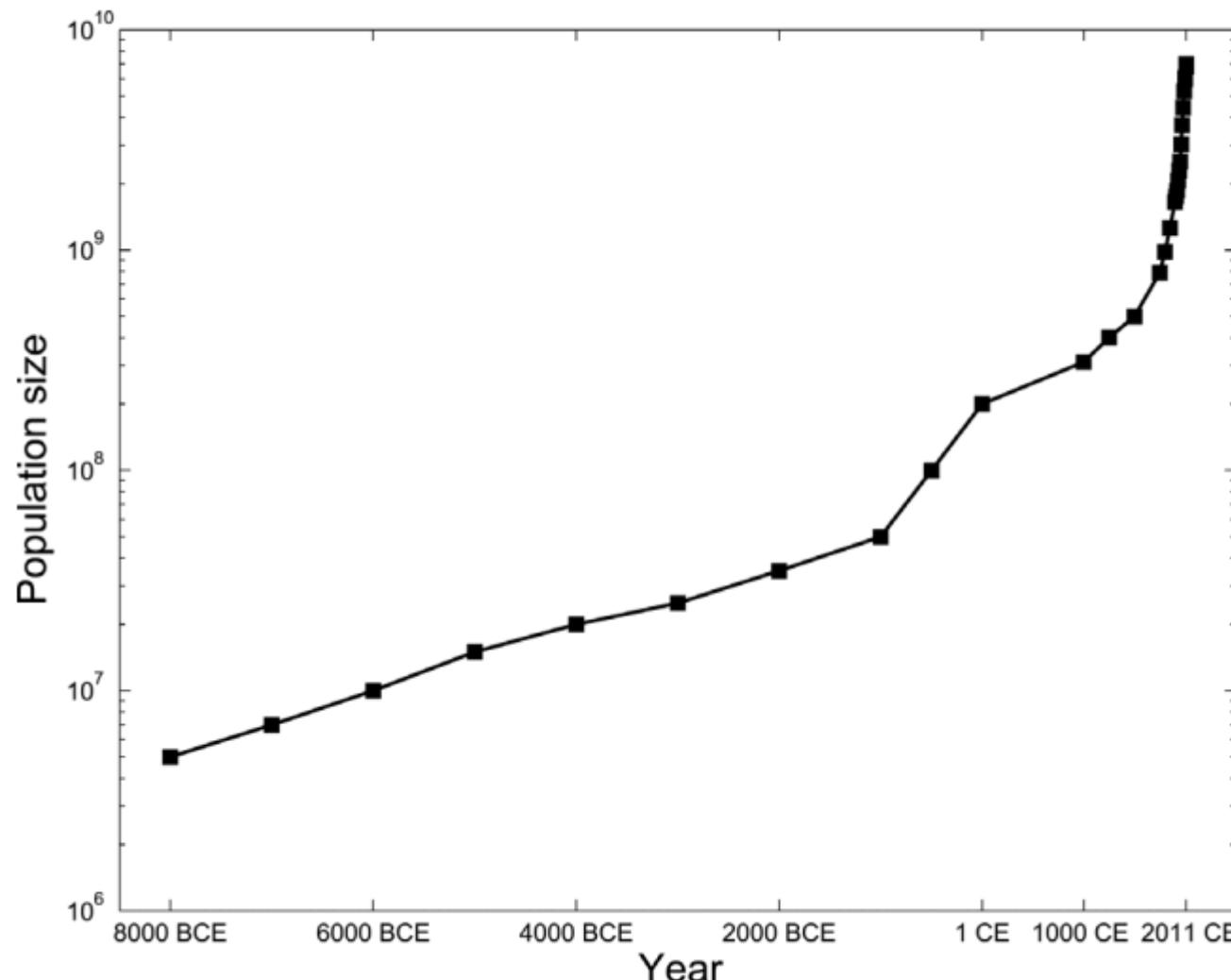


$n = 1000$



$n = 10,000$

The human population has grown super-exponentially



An Abundance of Rare Functional Variants in 202 Drug Target Genes Sequenced in 14,002 People

Matthew R. Nelson,^{1,*†} Daniel Wegmann,^{2*} Ma...
Pamela St. Jean,¹ Claudio Verzilli,³ Judo...
Dana Fraser,¹ Liling Warren,¹ Jennifer...
Yong Zhang,⁴ Jun Li,⁷ Yun Li,⁵ ...
Keith Nangle,¹ Jun Wan,³ ...
John C. Whittaker,³ ...

International weekly journal of science

nature LETTER

Evolutionary analysis of 6,515 exomes reveals the recent origin of most human protein-coding variants

Jacob A. Tennessen,^{1,*} Al... Timothy D. O'Connor,^{1*} Wenqing Fu,¹
Eimear E. Kenny,³ Simon... McGee,¹ Ron Do,^{4,5} Xiaoming Liu,⁶ Goo Jun,⁷
Hyun Min Kang,⁷ Daniel Jo... Suzanne M. Leal,⁹ Stacey Gabriel,⁴ Mark J. Rieder,¹
Goncalo Abecasis,⁷ David Altshuler,⁴ Deborah A. Nickerson,¹ Eric Boerwinkle,^{6,10}
Shamil Sunyaev,^{4,8} Carlos D. Bustamante,³ Michael J. Bamshad,^{1,2,‡} Joshua M. Akey,^{1,‡}
Broad GO, Seattle GO, on behalf of the NHLBI Exome Sequencing Project

doi:10.1038/nature11690

Deep

Implications of recent growth for complex disease studies

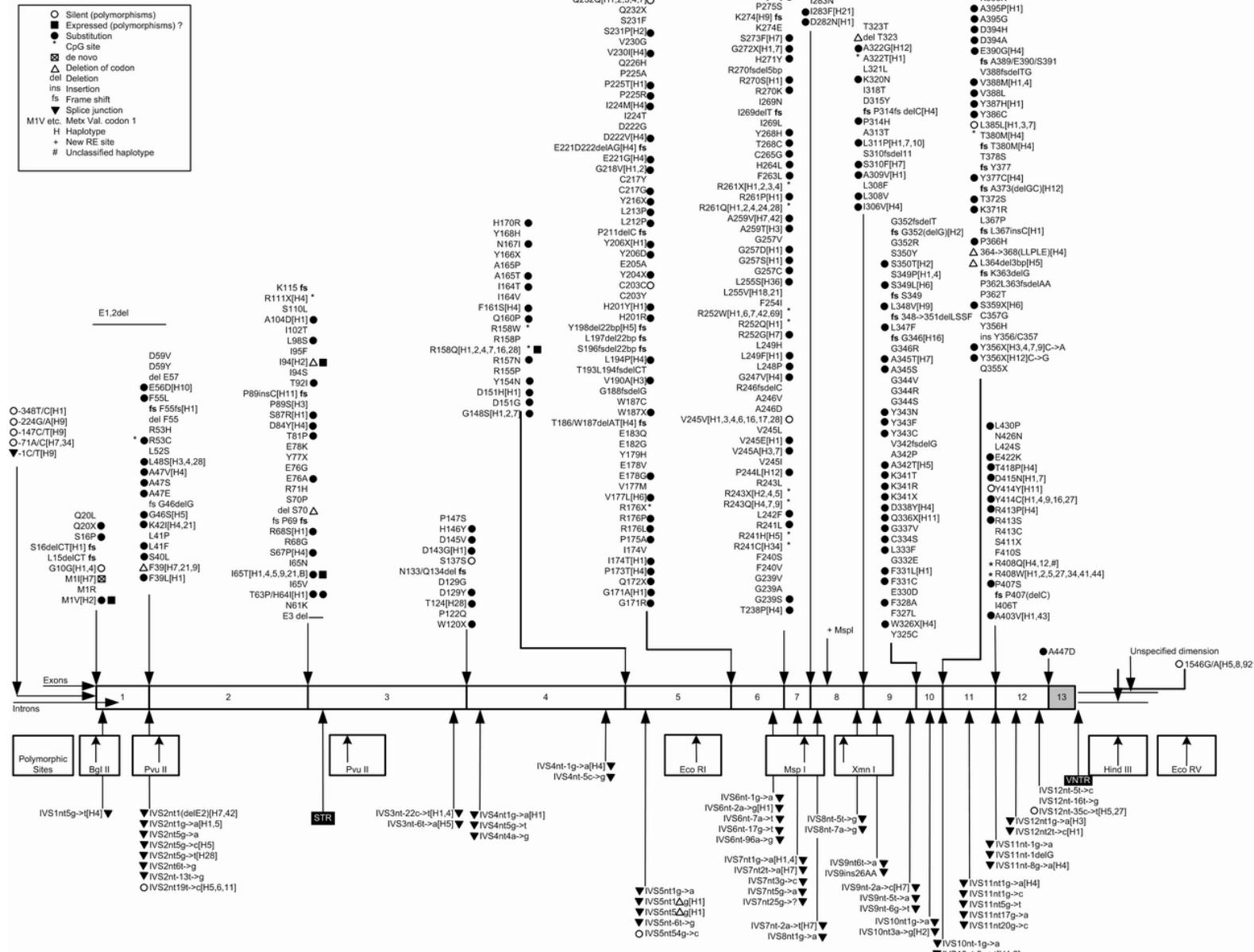
- Many more variants of all classes.
- Massive excess of rare variants.
- Rare variants are more likely to have deleterious effects.

Implications of recent growth for complex disease studies

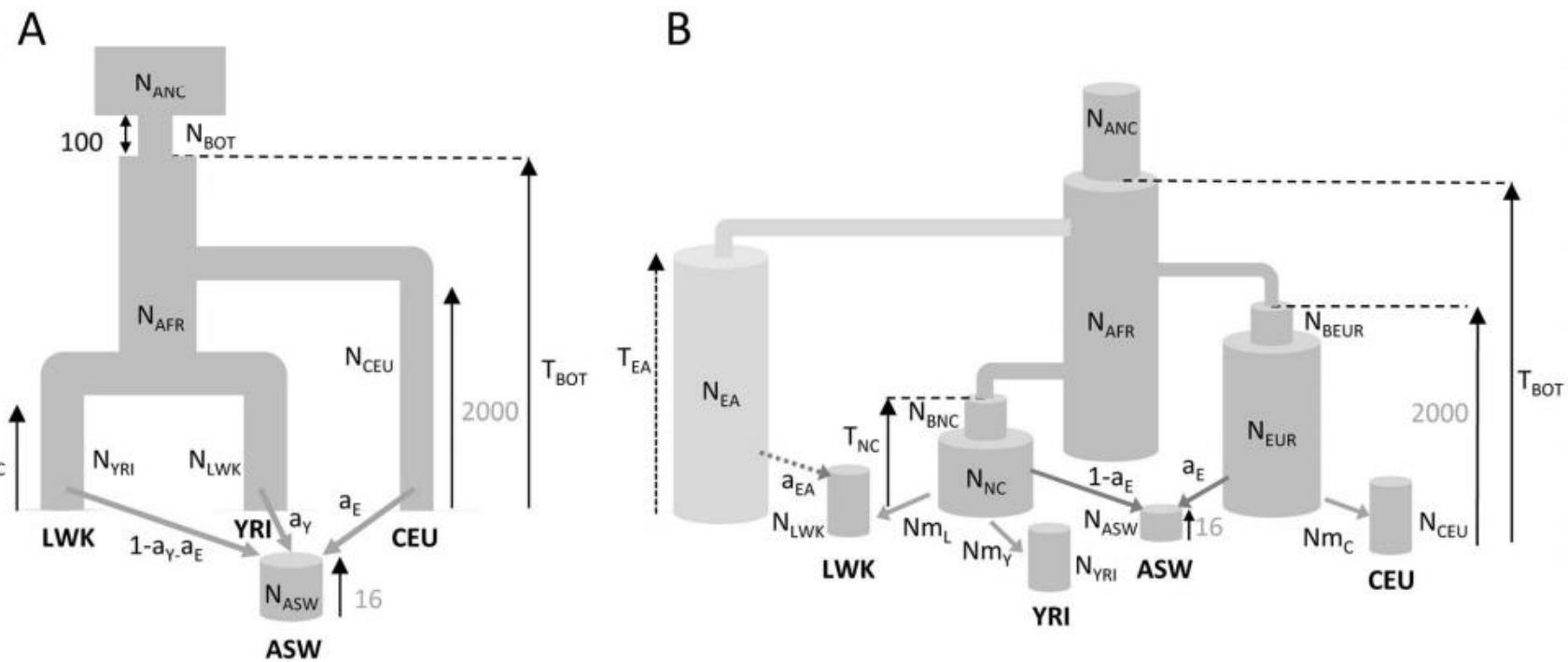
- Many more variants of classes.
- Frequency shifted to other variants.
- Rare variants more likely to have deleterious effects.

Genetic heterogeneity

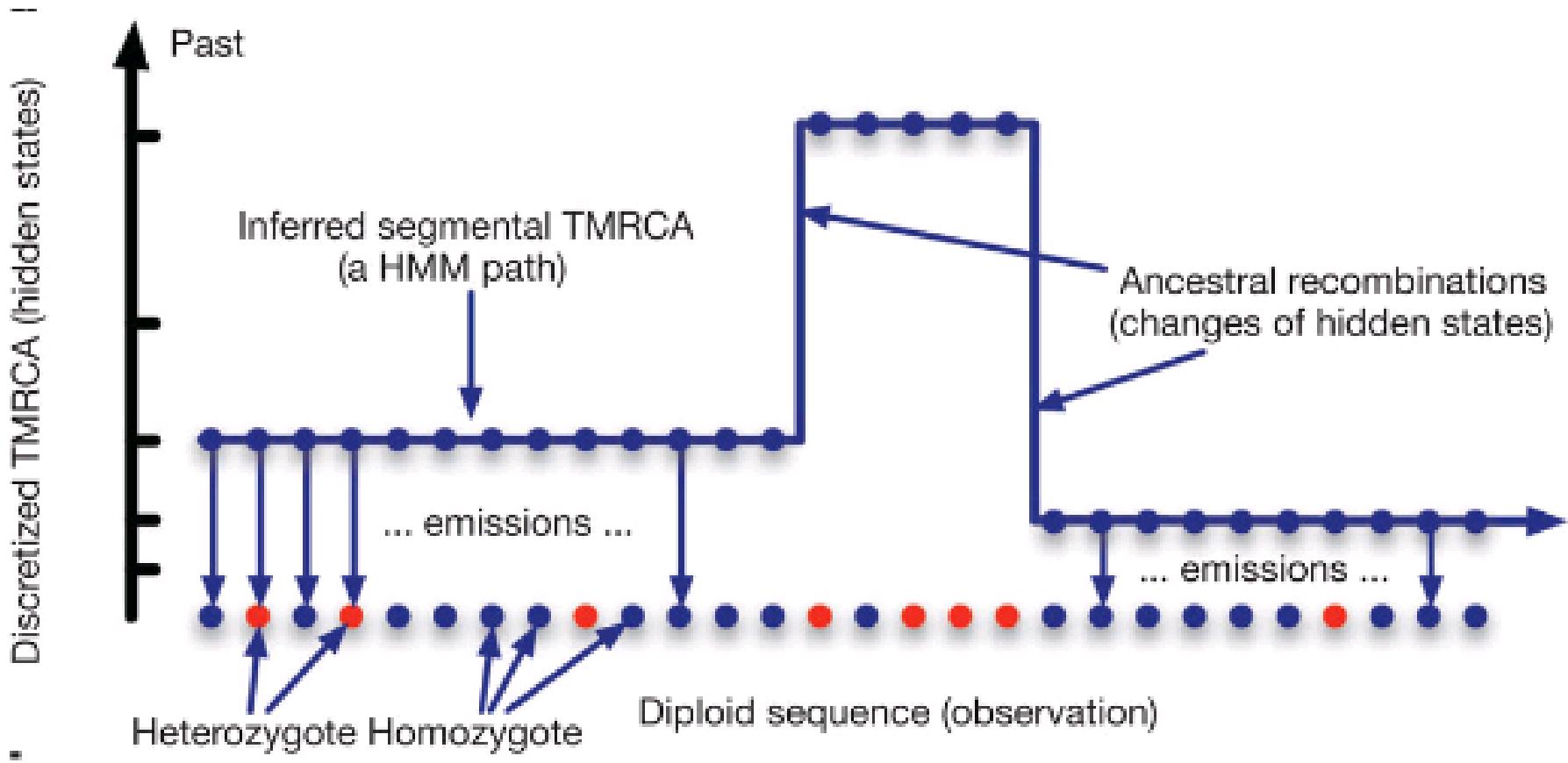
Allelic heterogeneity of PAH



Many methods for inference of demography: IM, dadi, fastsimcoal

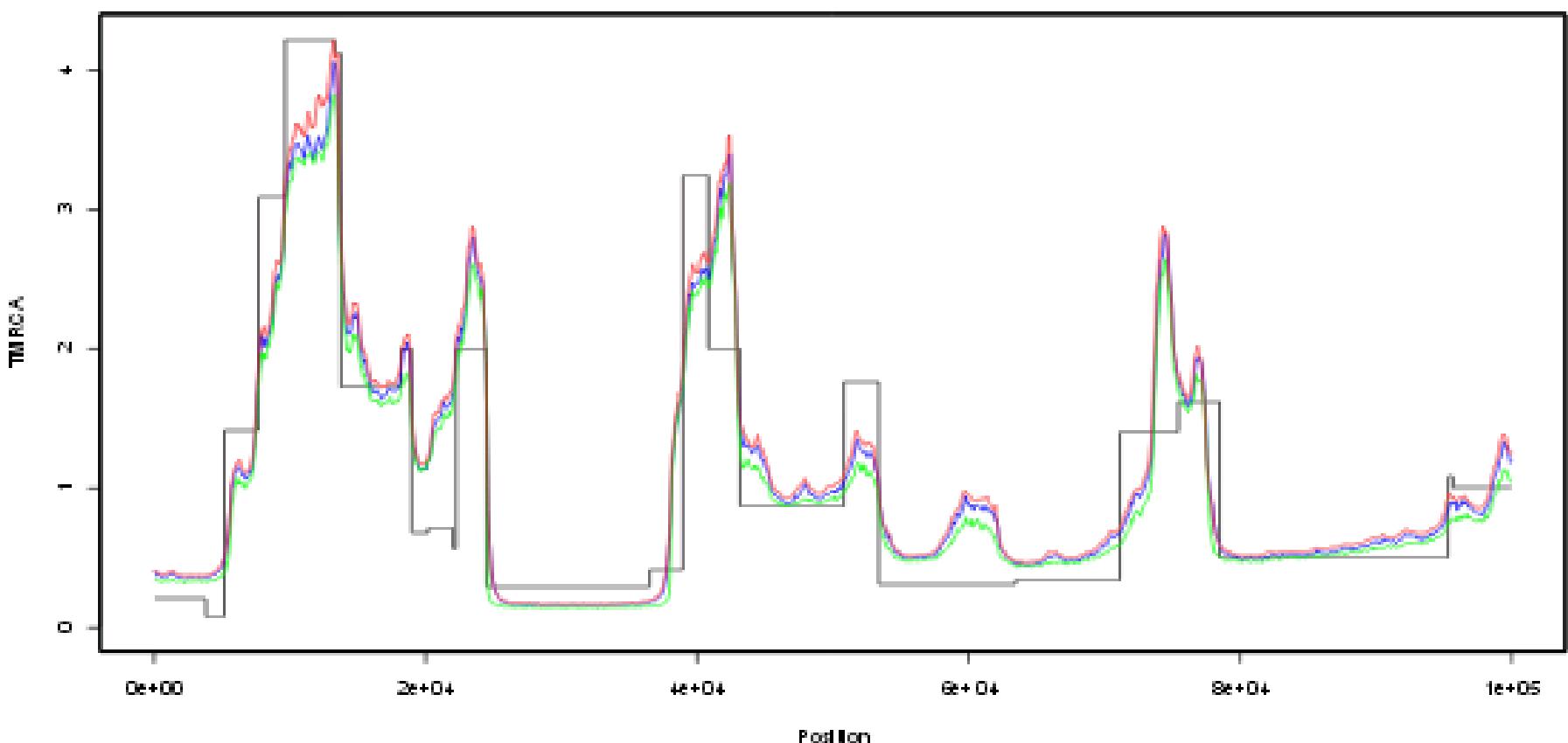


Inferring demography from a single individual: Pairwise Sequentially Markovian Coalescent

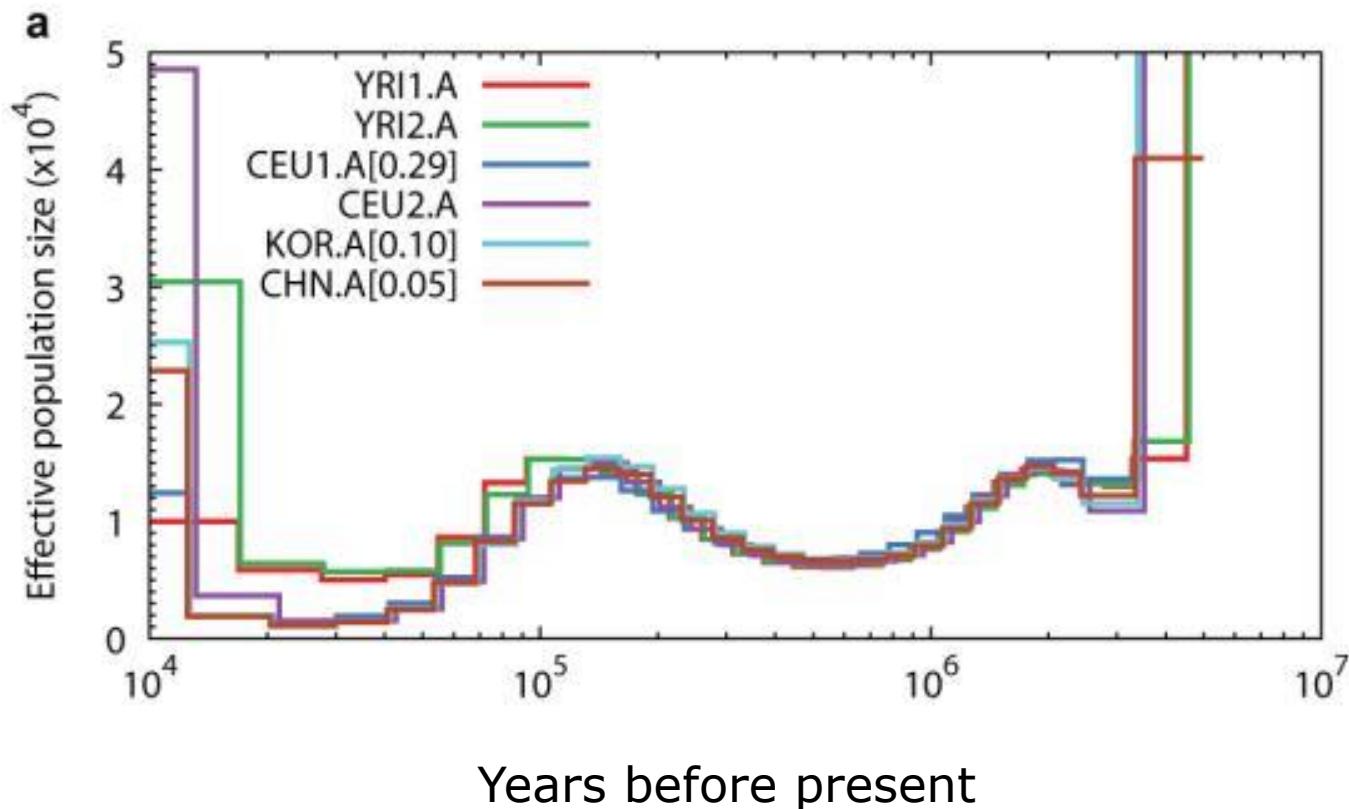


Simulations of past population demography shows reasonably good accuracy of PSMC, with an important caveat

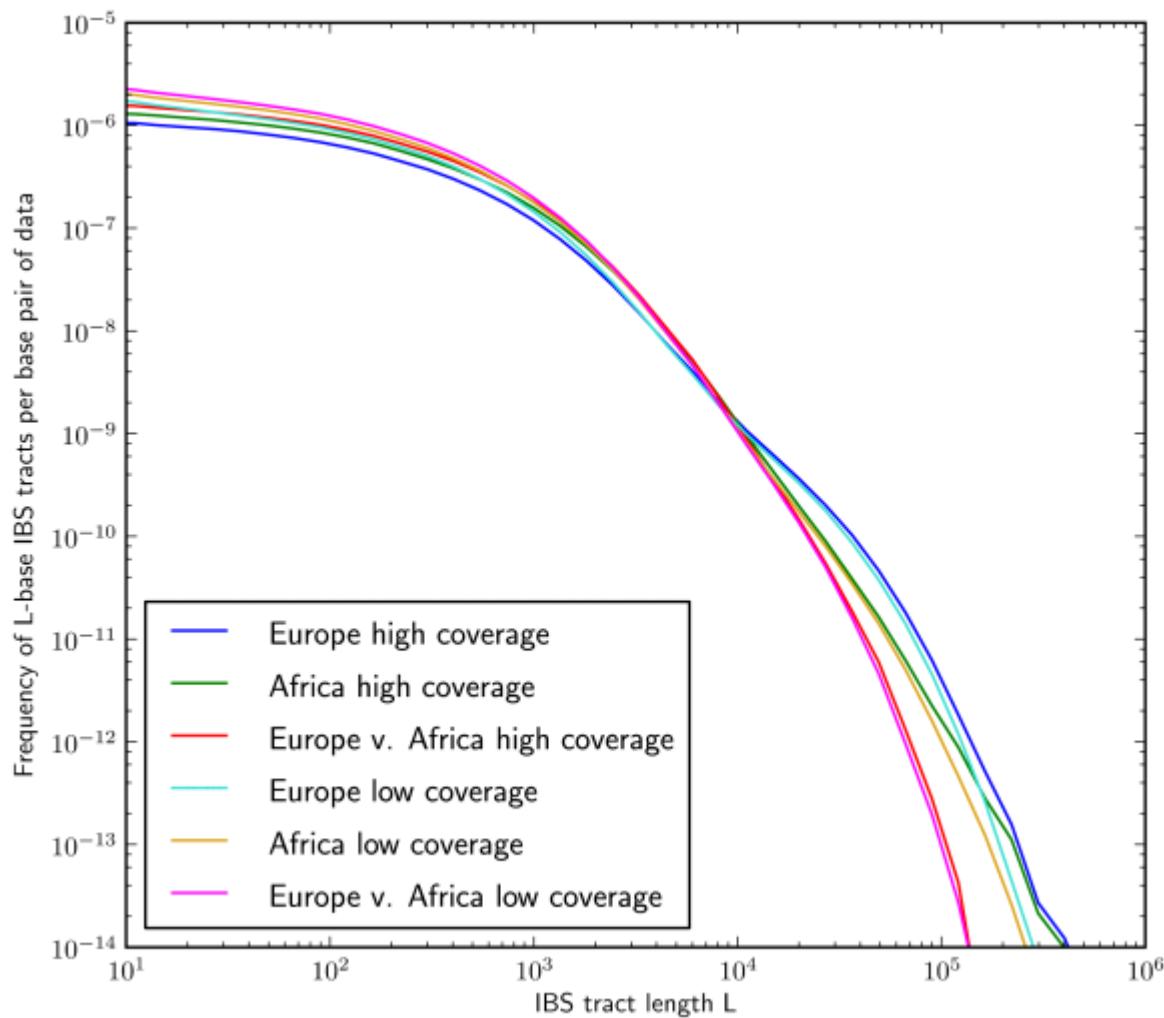
Black line = True TMRCA, Red line = psmc meanTMRCA, Blue line = our meanTMRCA, Green line = our medianTMRCA



If the population has deme structure, and we sample from just one deme, we can spuriously infer a bottleneck !



Identity-by-Descent tracts for inference of demography



DEMOGRAPHY

(Population collapse)

What happens to genetic variation in genomes of populations that are crashing?

Florida Scrub-Jay

(*Aphelocoma coerulescens*)

Cooperative breeder

Federally Threatened

Non-migratory

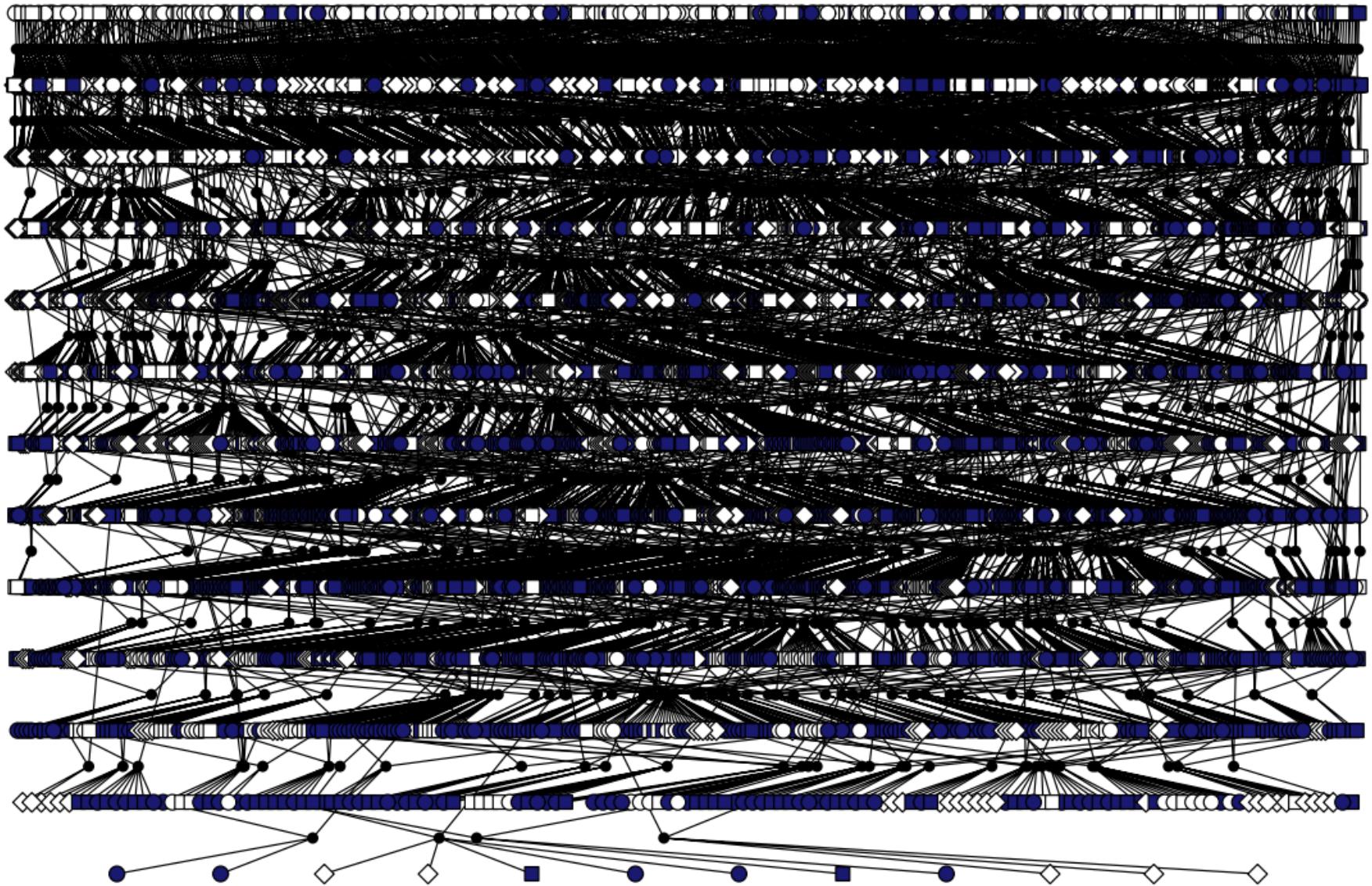
Highly territorial & philopatric

Socially & (mostly) genetically
monogamous

**All individuals banded and
tracked throughout life**

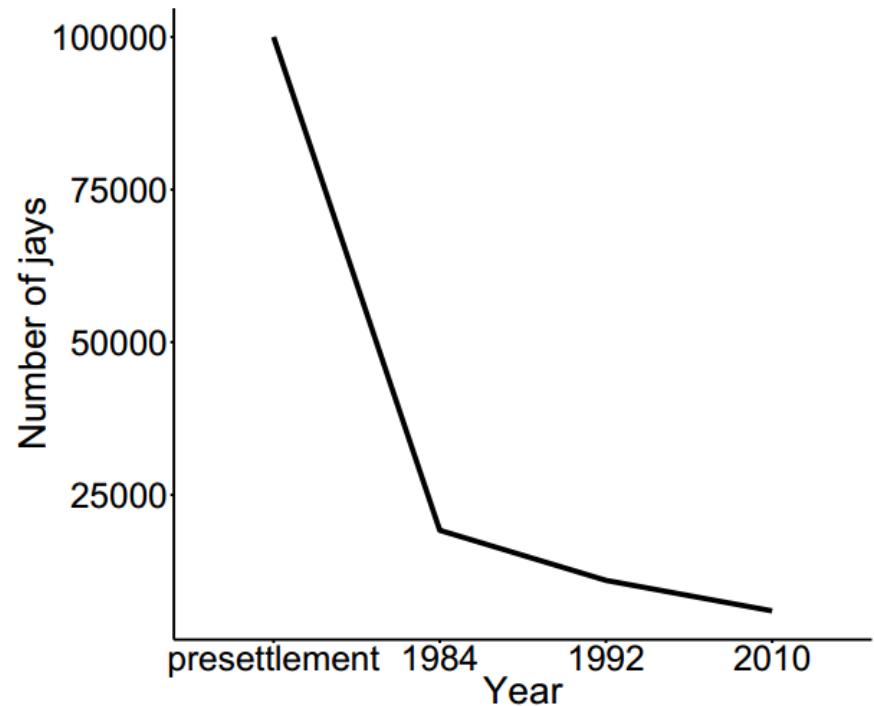
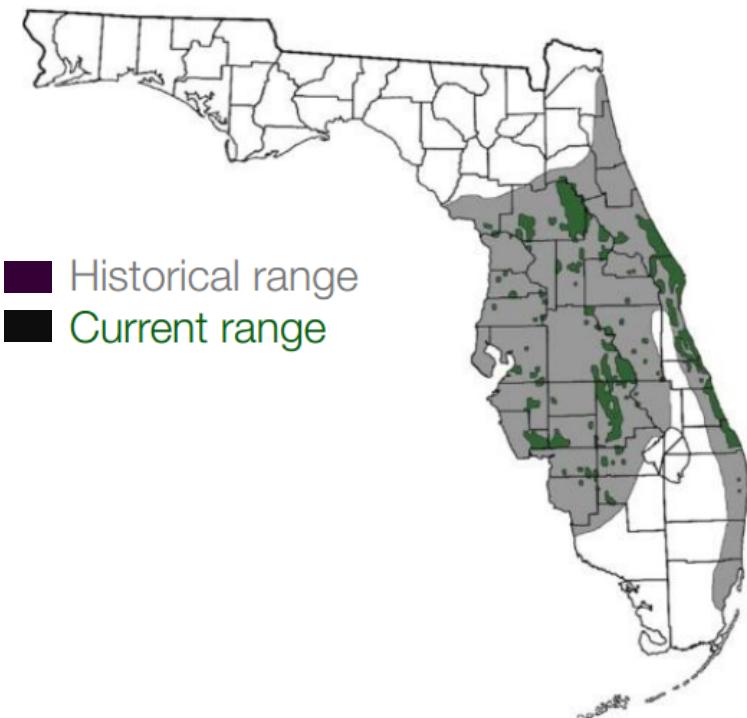
Nancy Chen





We have blood samples as well as life history & morphological data for >4,000 pedigreed individuals

Florida Scrub-Jay populations have drastically declined due to habitat loss



97% decline in past century. 50% decline in past 20 years

Florida Scrub-Jay genomic resources



Genome-wide SNPs

Genome assembly

Transcriptome assembly

Genome annotation

Linkage map construction

Florida Scrub-Jay genomic resources



Genome-wide SNPs

Genome assembly

Transcriptome assembly

Genome annotation

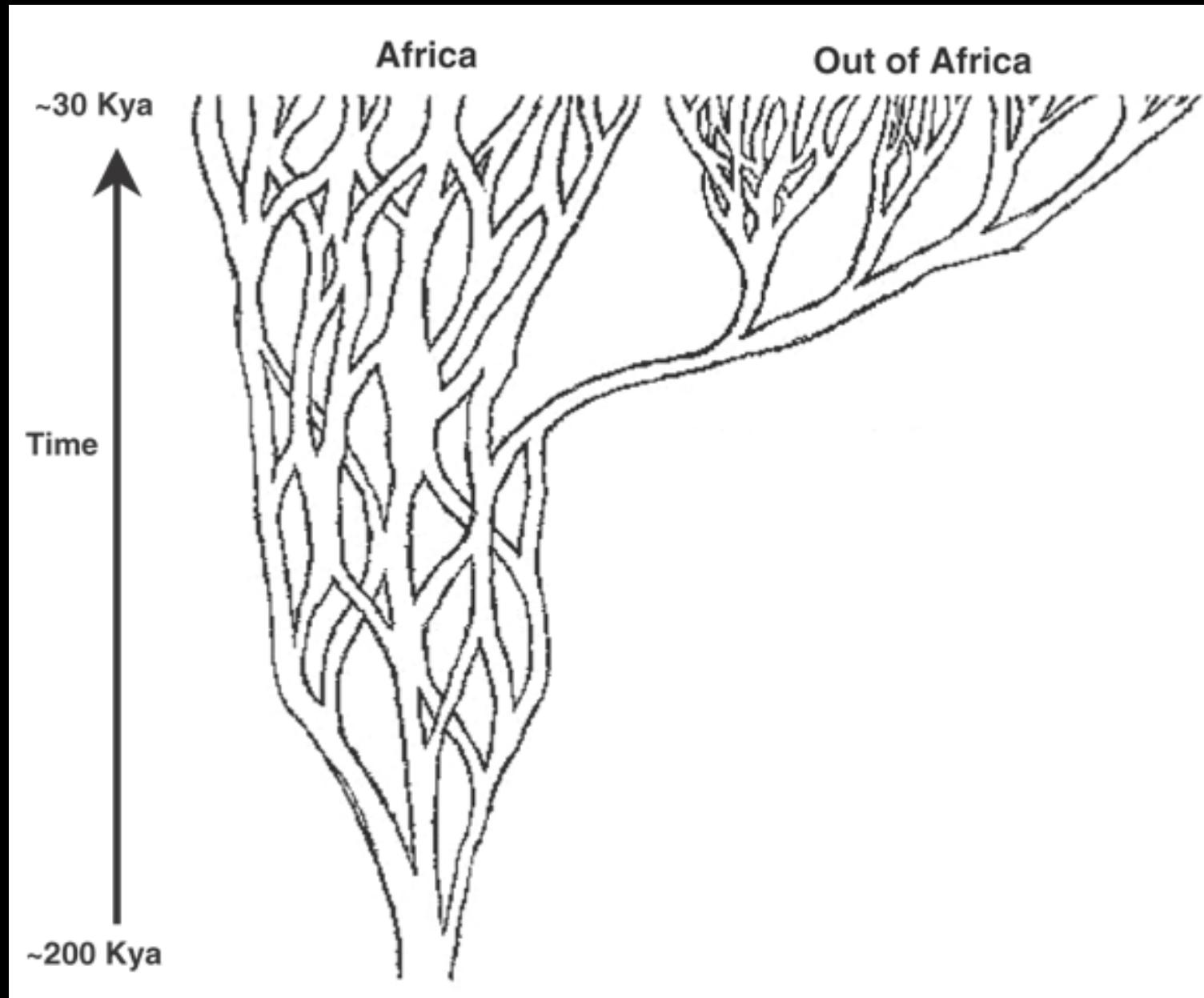
Linkage map construction

How does the recent population crash manifest itself in the structure of genetic variation?

Population Structure

How can we infer past patterns of migration from genome sequence?

Human population history



Principal Components Analysis for population structure

OPEN  ACCESS Freely available online

PLOS GENETICS

Population Structure and Eigenanalysis

Nick Patterson^{1*}, Alkes L. Price^{1,2}, David Reich^{1,2}

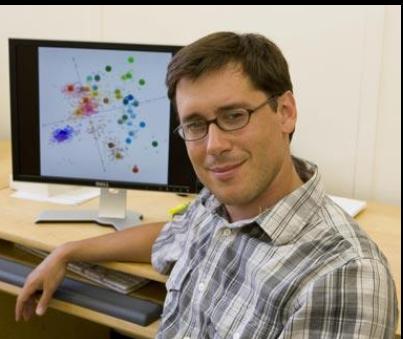
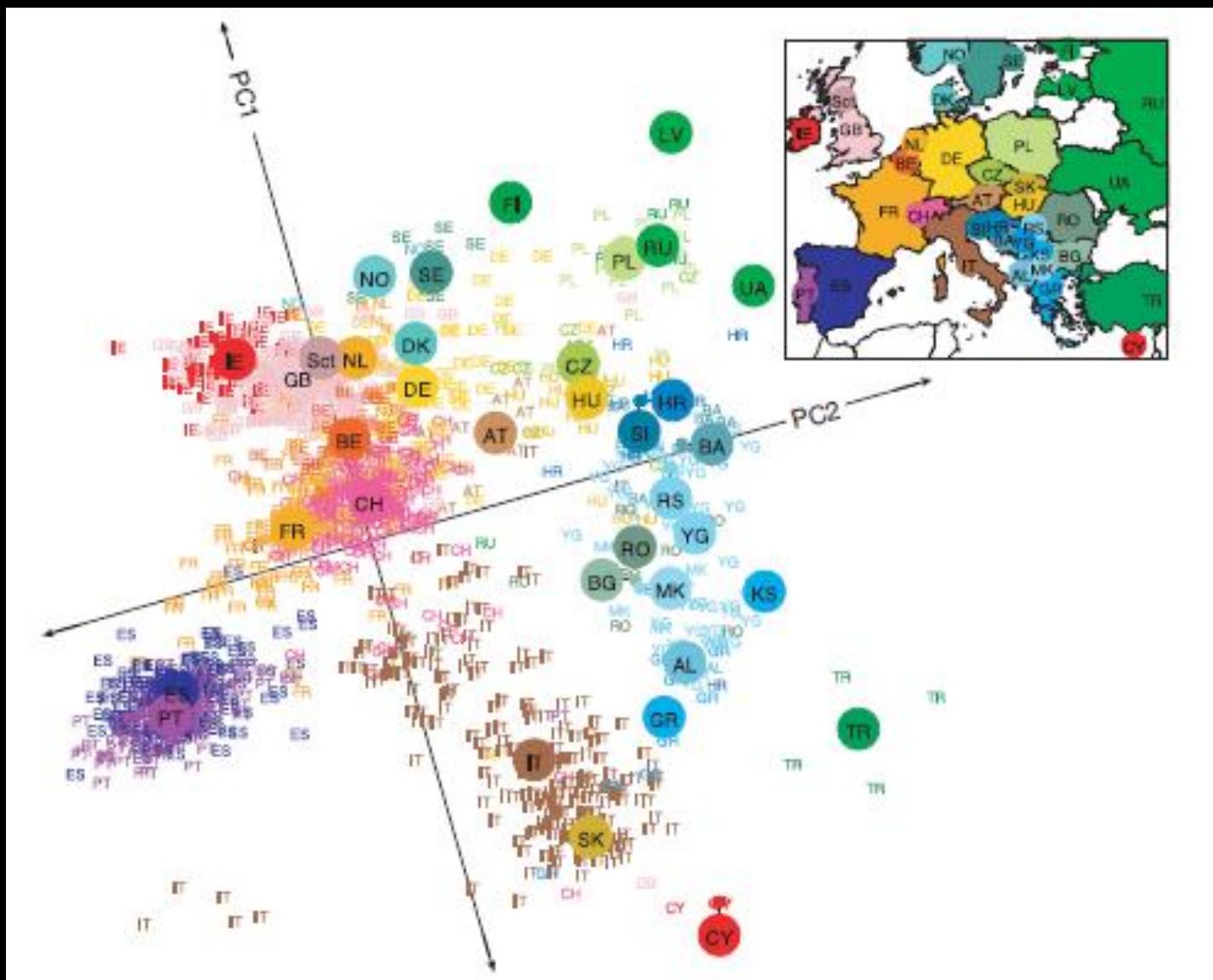
1 Broad Institute of Harvard and MIT, Cambridge, Massachusetts, United States of America, 2 Department of Genetics, Harvard Medical School, Boston, Massachusetts, United States of America

Current methods for inferring population structure from genetic data do not provide formal significance tests for population differentiation. We discuss an approach to studying population structure (principal components analysis) that was first applied to genetic data by Cavalli-Sforza and colleagues. We place the method on a solid statistical footing, using results from modern statistics to develop formal significance tests. We also uncover a general “phase change” phenomenon about the ability to detect structure in genetic data, which emerges from the statistical theory we use, and has an important implication for the ability to discover structure in genetic data: for a fixed but large dataset size, divergence between two populations (as measured, for example, by a statistic like F_{ST}) below a threshold is essentially undetectable, but a little above threshold, detection will be easy. This means that we can predict the dataset size needed to detect structure.

Citation: Patterson N, Price AL, Reich D (2006) Population structure and eigenanalysis. PLoS Genet 2(12): e190. doi:10.1371/journal.pgen.0020190

Patterson N, Price AL, Reich D (2006) *PLoS Genet* 2(12): e190.

European genetic differentiation using PCA

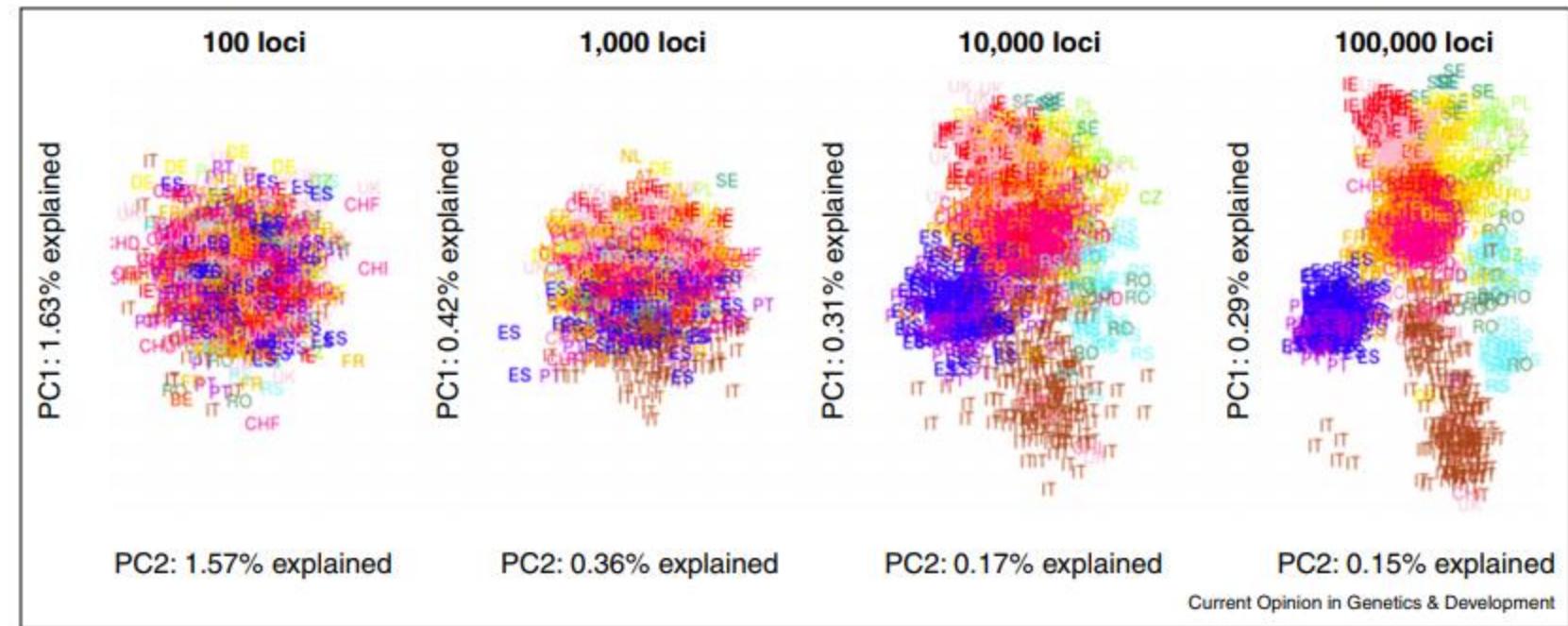


John Novembre

Novembre et al. 2008 *Nature* 456:274.

Accurate fine structure inference from PCA requires a huge number of SNPs

Figure 2



Rosenberg NA, Pritchard JK, Weber JL, Cann HM,
Kidd KK, Zhivotovsky LA, Feldman MW. 2002
Genetic structure of human populations.
Science. 298:2381–2385.

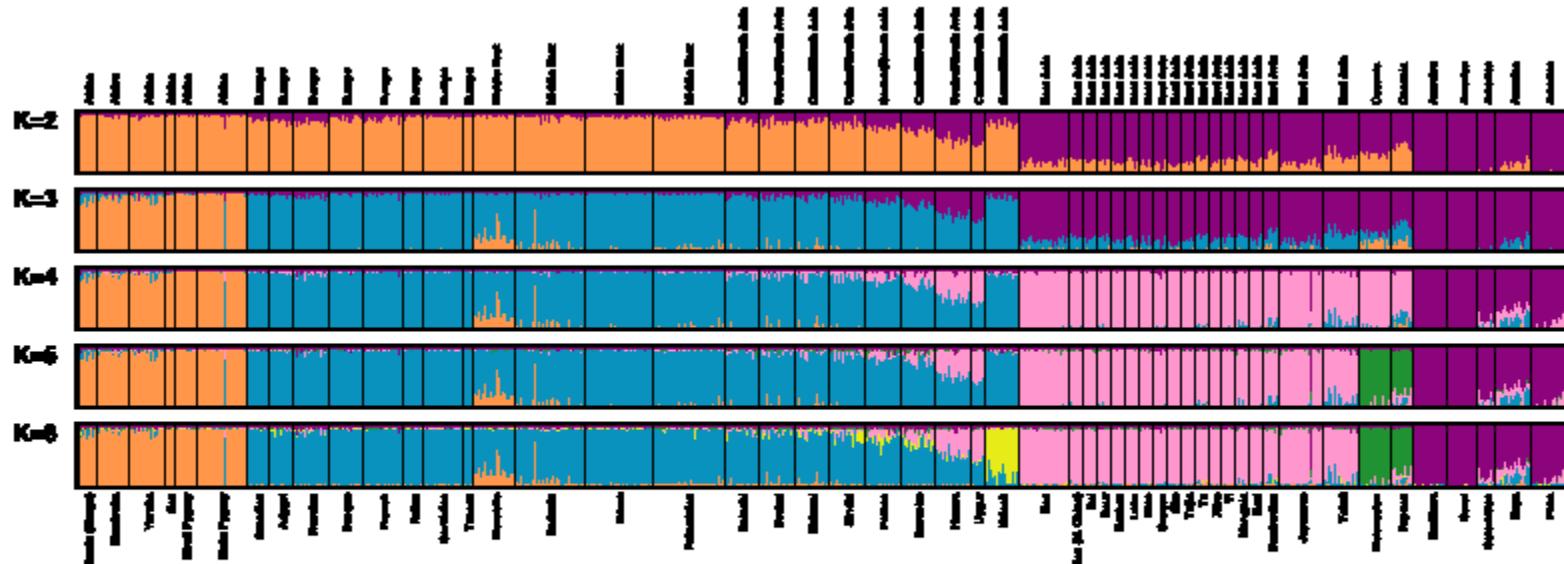
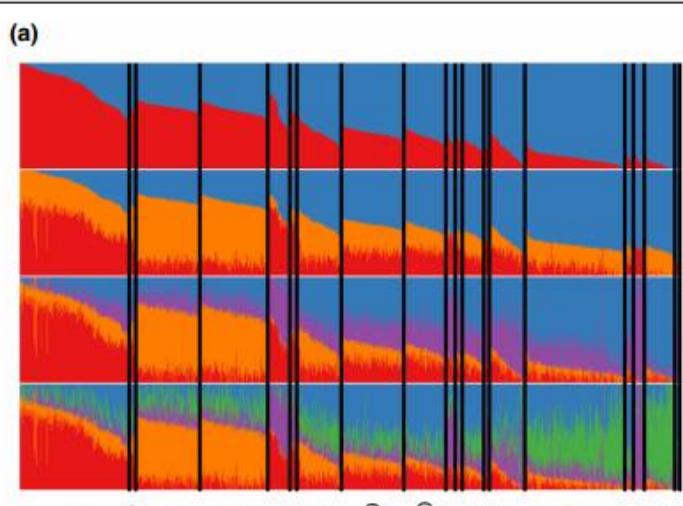
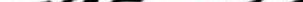


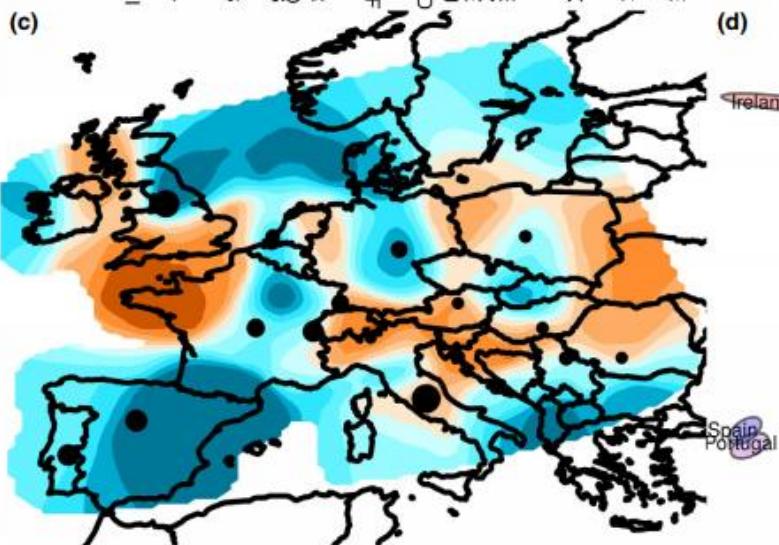
Fig. 1. Estimated population structure. Each individual is represented by a thin vertical line, which is partitioned into K colored segments that represent the individual's estimated membership fractions in K clusters. Black lines separate individuals of different populations. Populations are labeled below the figure, with their regional affiliations above it. Ten structure runs at each

K produced nearly identical individual membership coefficients, having pairwise similarity coefficients above 0.97, with the exceptions of comparisons involving four runs at $K = 3$ that separated East Asia instead of Eurasia, and one run at $K = 6$ that separated Karitiana instead of Kalash. The figure shown for a given K is based on the highest probability run at that K .

ADMIXTURE

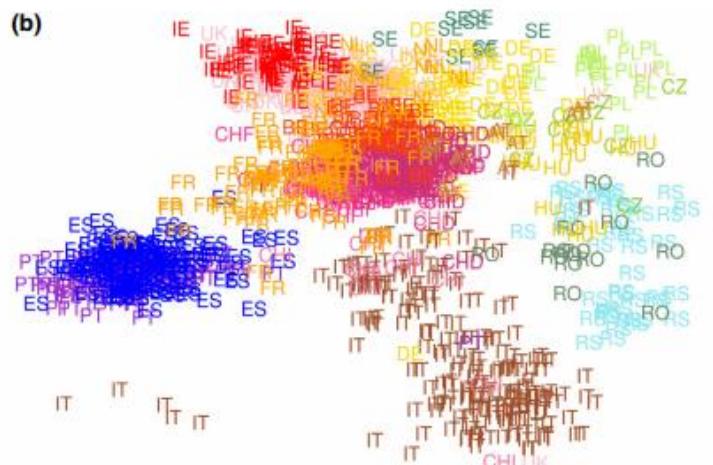


(c) 

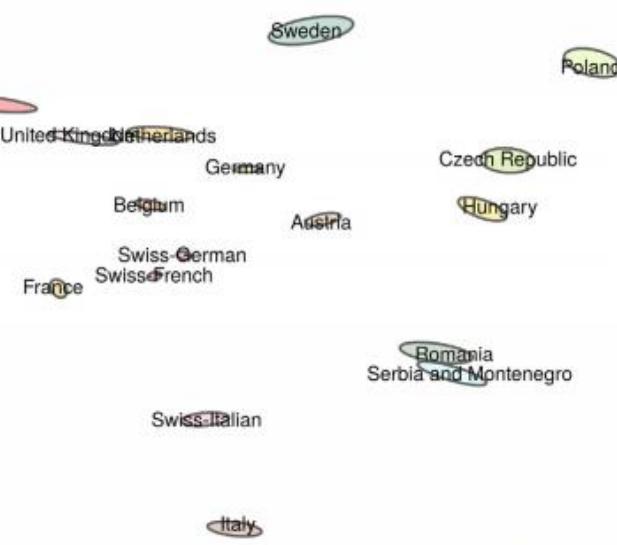


EEMS

PCA



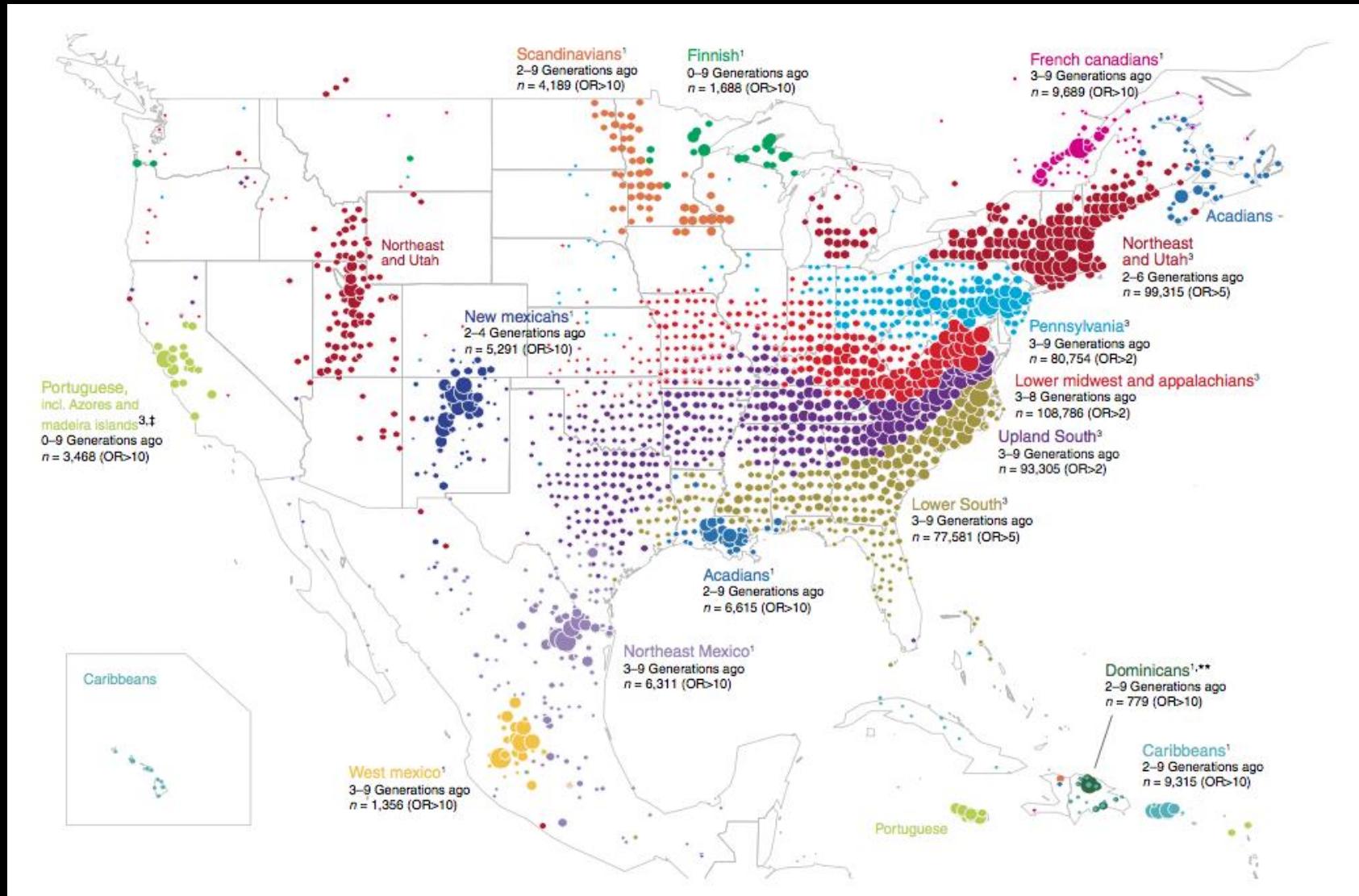
(d)



Current Opinion in Genetics & Development

SpaceMix

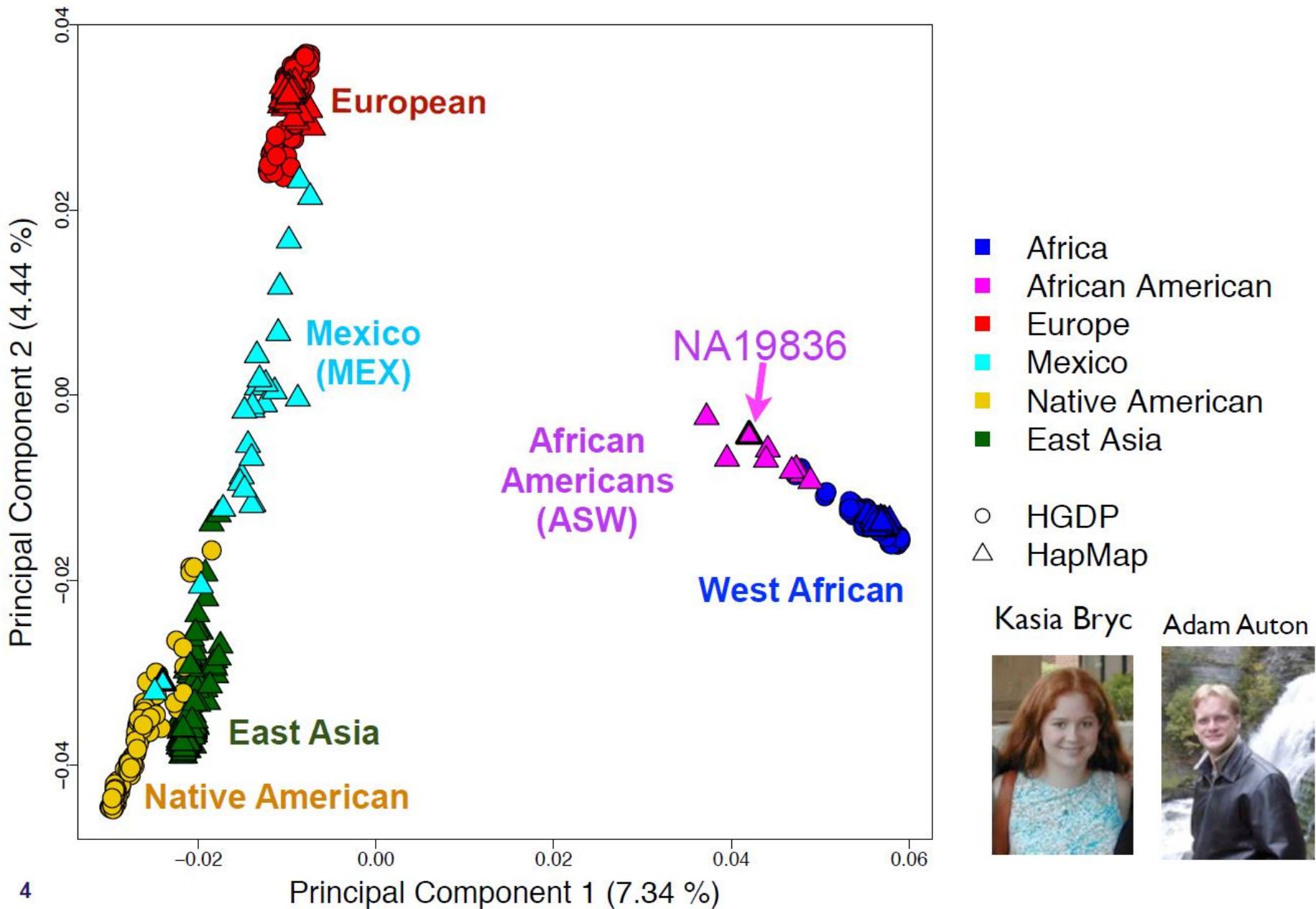
Inference of population fine-structure from IBD



ADMIXTURE

How can we infer consanguinity and
admixture from genetic data?

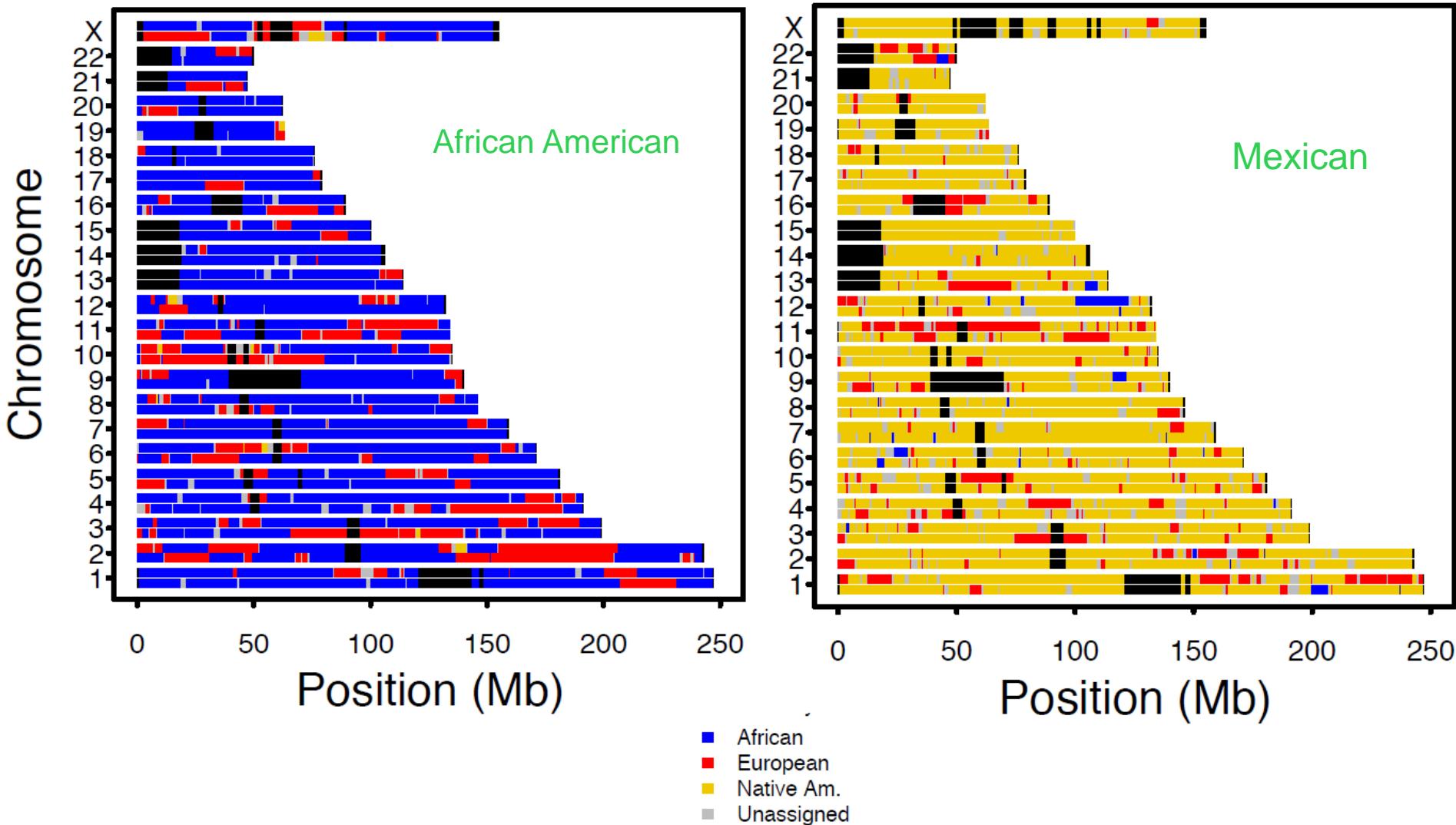
Using PCA to infer admixed individuals



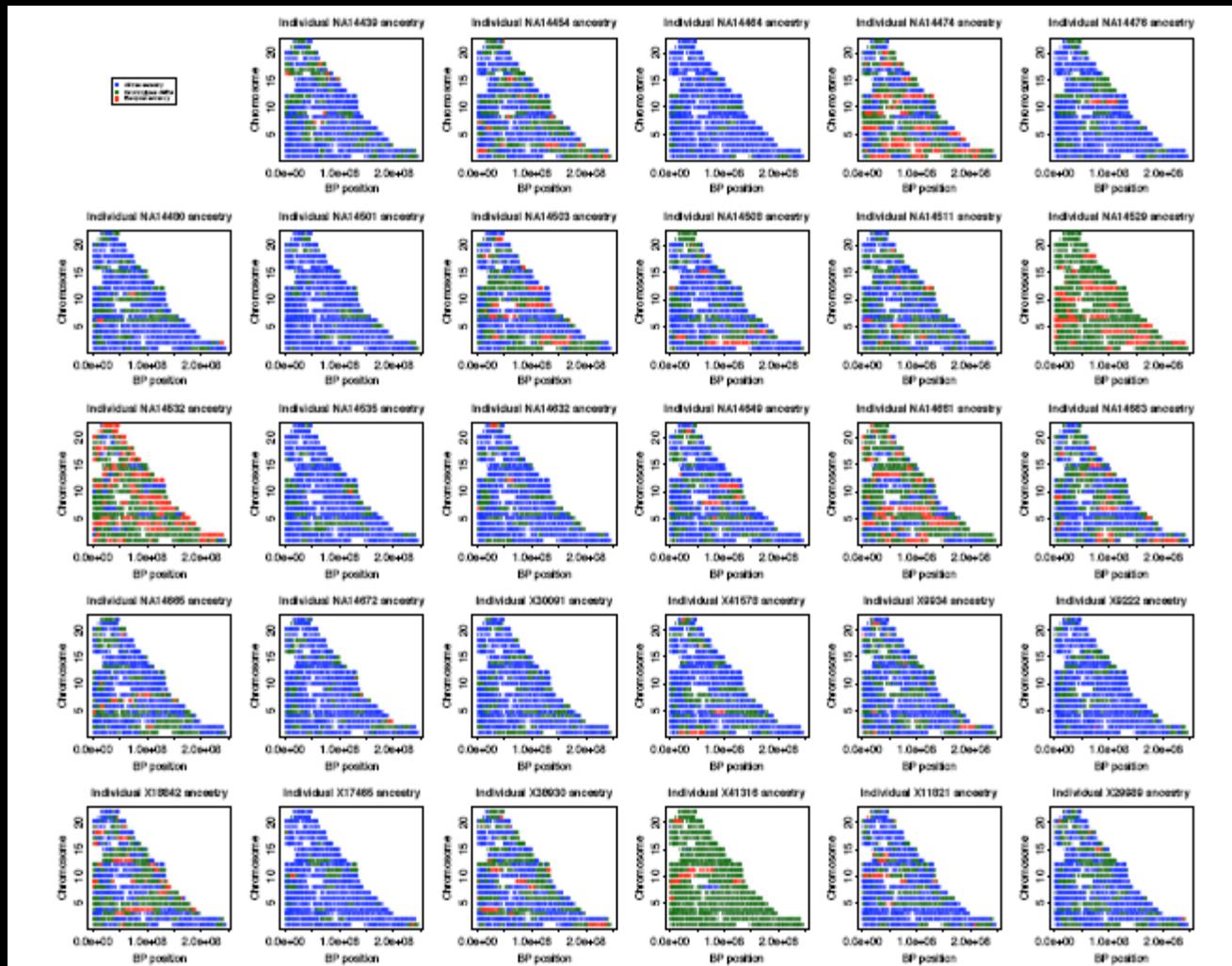
Local ancestry inference

- Run PCA on African, European, and African American samples.
- From PC loadings of 15 SNP windows, infer 0, 1, or 2 copies of African ancestry.
- Slide along the genome.
- RESULT: Call of 0, 1, 2 copies of African ancestry for each chunk of the genome in each individual.

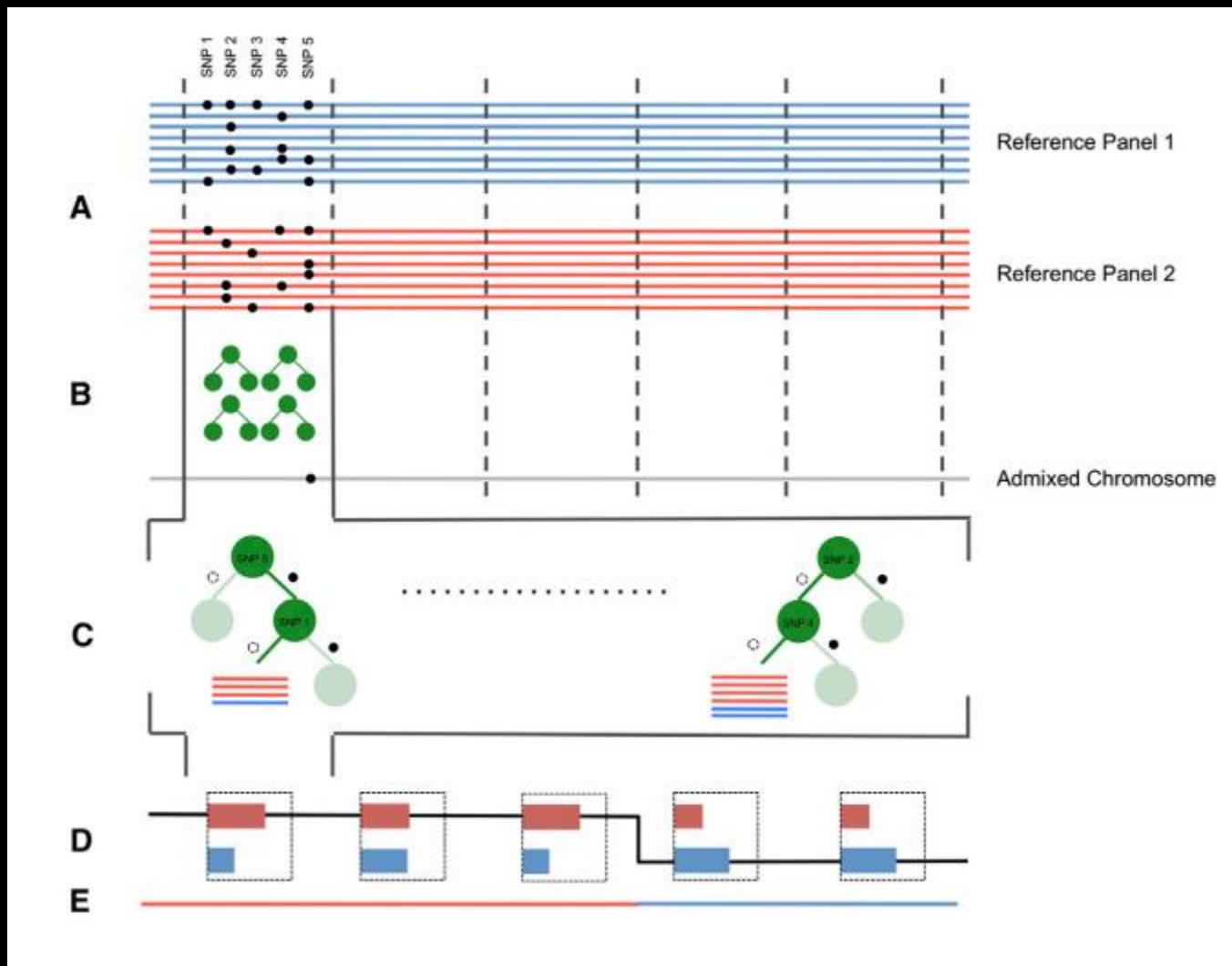
PCAdmix can identify the population-of-origin of segments of the genome



High variability among individuals in admixture patterns



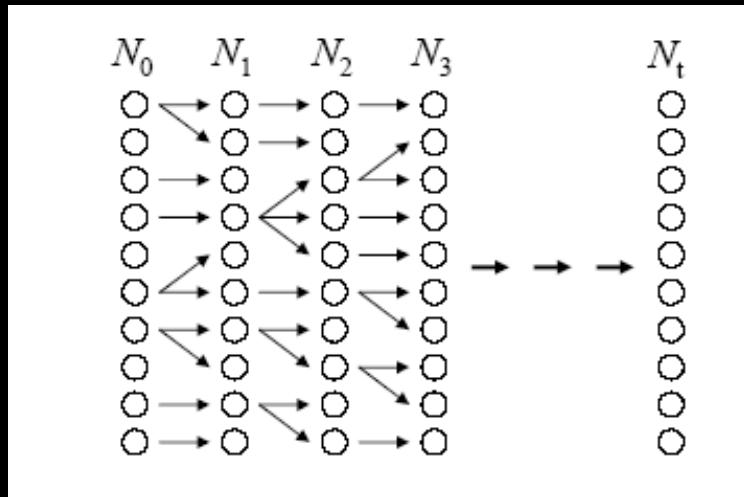
RFMix uses sequence from ancestral populations and a window-based random forest approach to map ancestry blocks



Random Genetic Drift

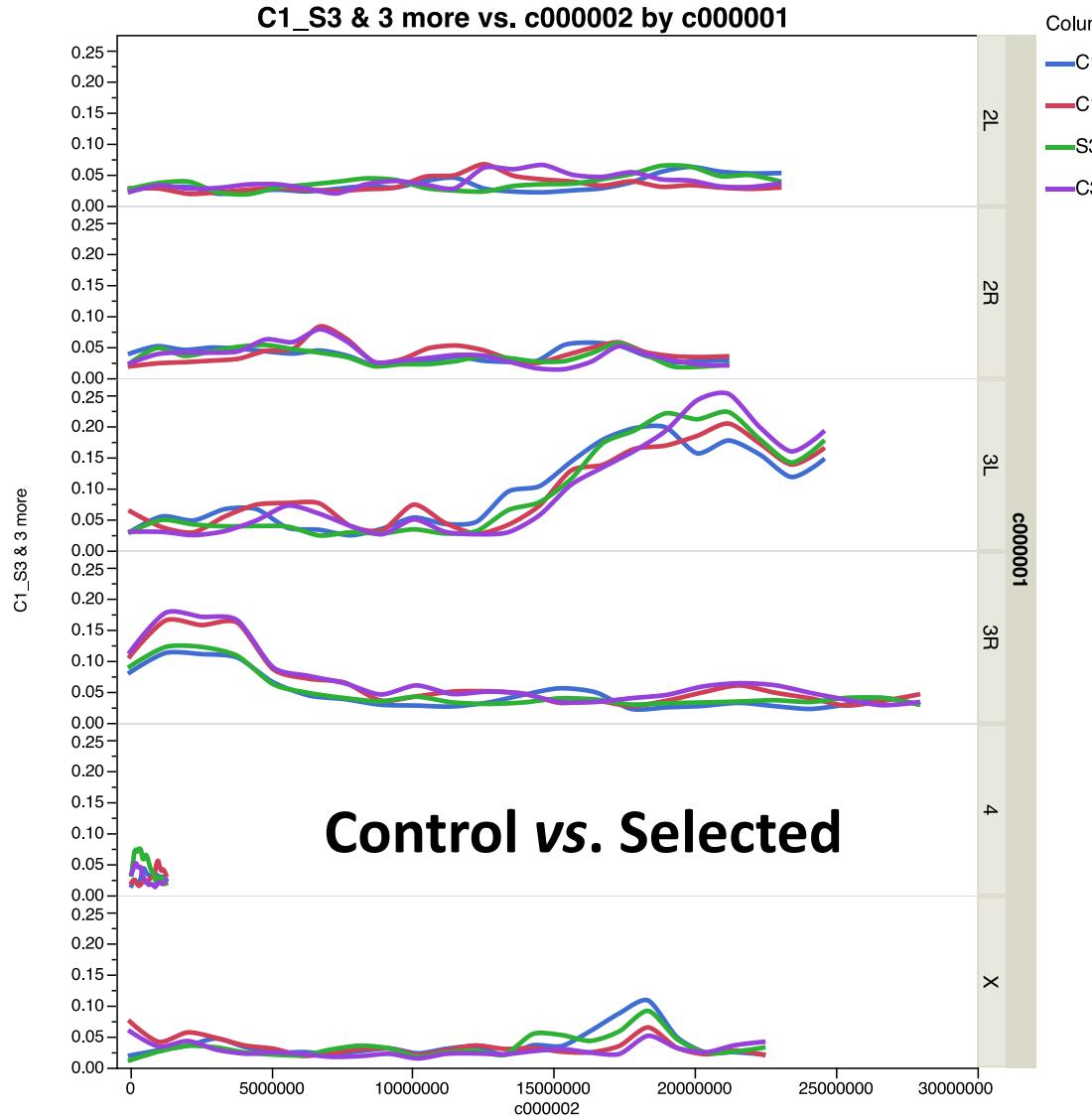
What can we infer from allele frequency dynamics of every nucleotide in the genome?

The Wright-Fisher drift model: the Null model for Evolve-and-Resequence



- Selfing allowed
- Random mating
- Non-overlapping generations
- Constant population size
- No migration
- No selection

Genomic regions show consistent changes in allele frequency



Bioinformatics, 31(11), 2015, 1762–1770

doi: 10.1093/bioinformatics/btv014

Advance Access Publication Date: 21 January 2015

Original paper

OXFORD

Genetics and population analysis

Gaussian process test for high-throughput sequencing time series: application to experimental evolution

**Hande Topa^{1,*†}, Ágnes Jónás^{2,3,*†}, Robert Kofler², Carolin Kosiol^{2,*}
Antti Honkela^{4,*}**

¹Helsinki Institute for Information Technology (HIIT), Department of Information and Computer Science, Aalto University, Espoo, Finland, ²Institut für Populationsgenetik, Vetmeduni Vienna, 1210 Wien, Austria, ³Vienna Graduate School of Population Genetics, Wien, Austria and ⁴Helsinki Institute for Information Technology (HIIT), Department of Computer Science, University of Helsinki, Helsinki, Finland

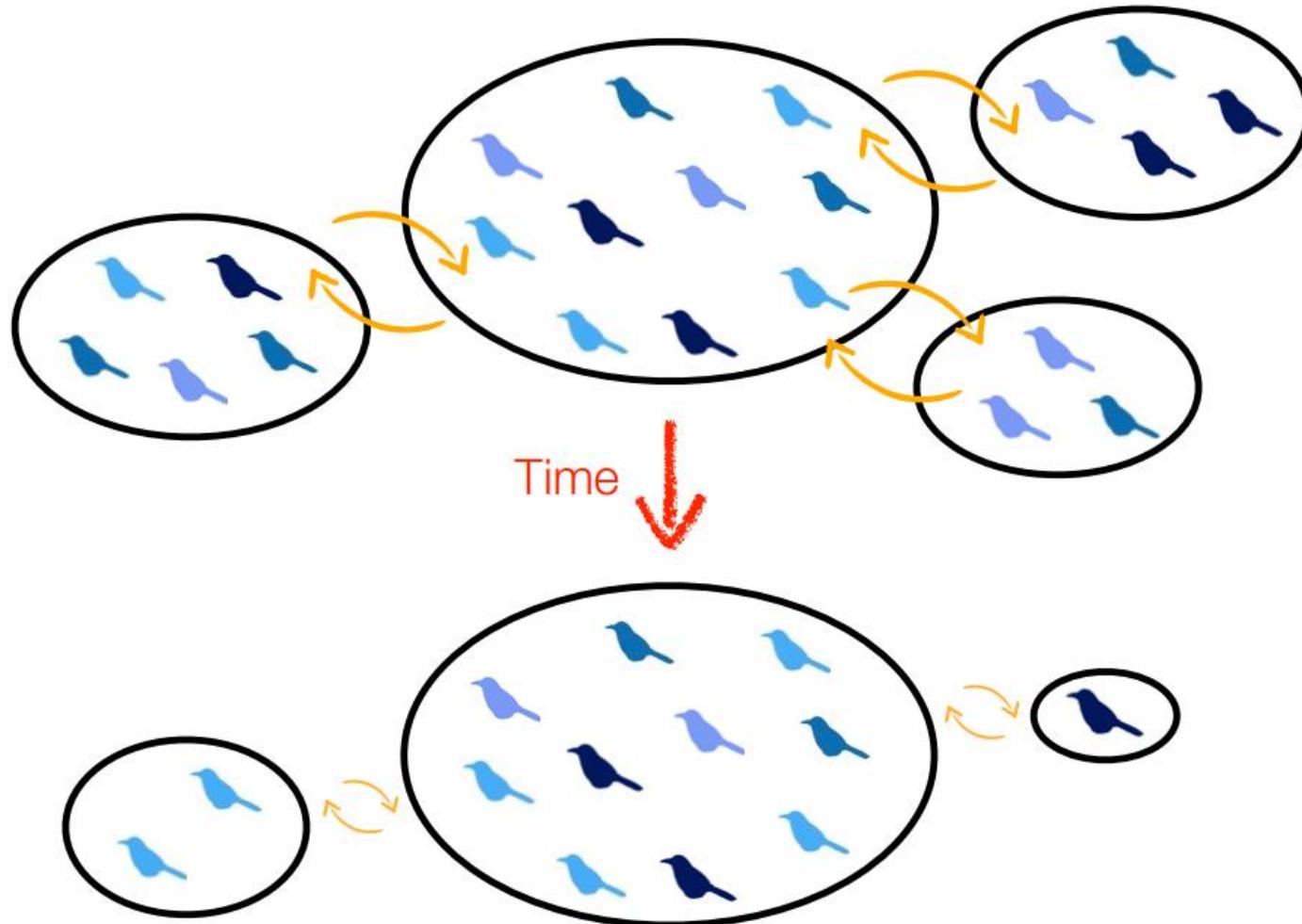
Fitting pool-seq data by Gaussian process

- Sampling from one generation to the next each generation is Binomial(p_i, n_j)
- The p_i are drawn from a Beta distribution.
- The sampling is correlated along the genome, with covariance depending on linkage and LD.
- This model is fitted by Gaussian Process regression (least square).

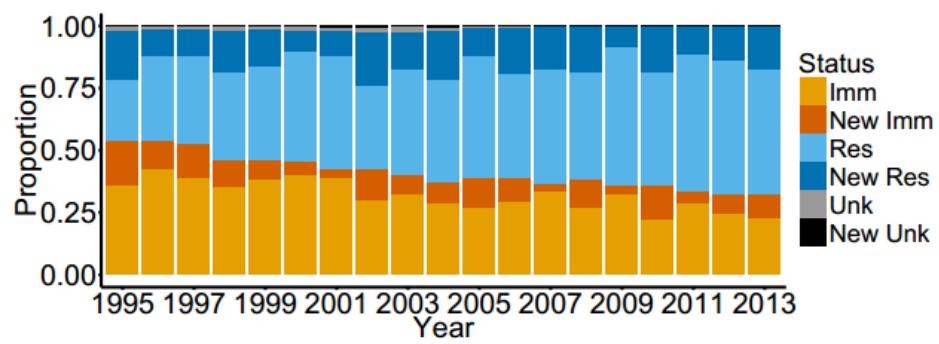
Migration

How does spatial movement of individuals interact with other forces?

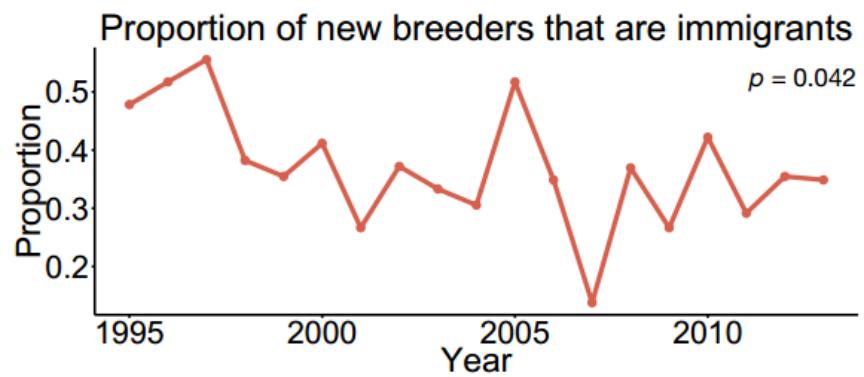
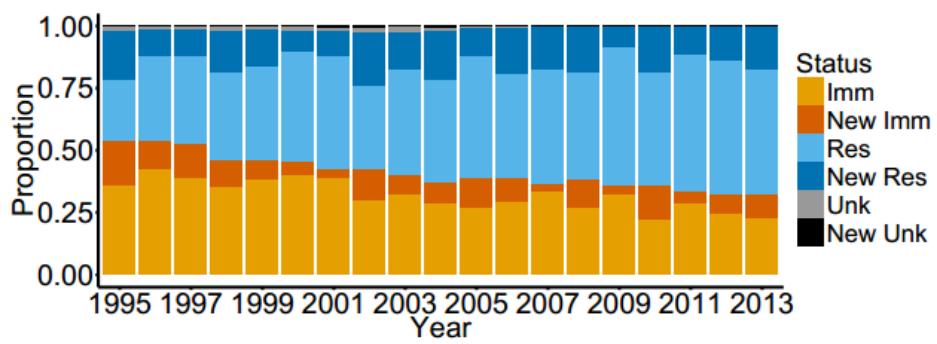
Consequences of habitat fragmentation



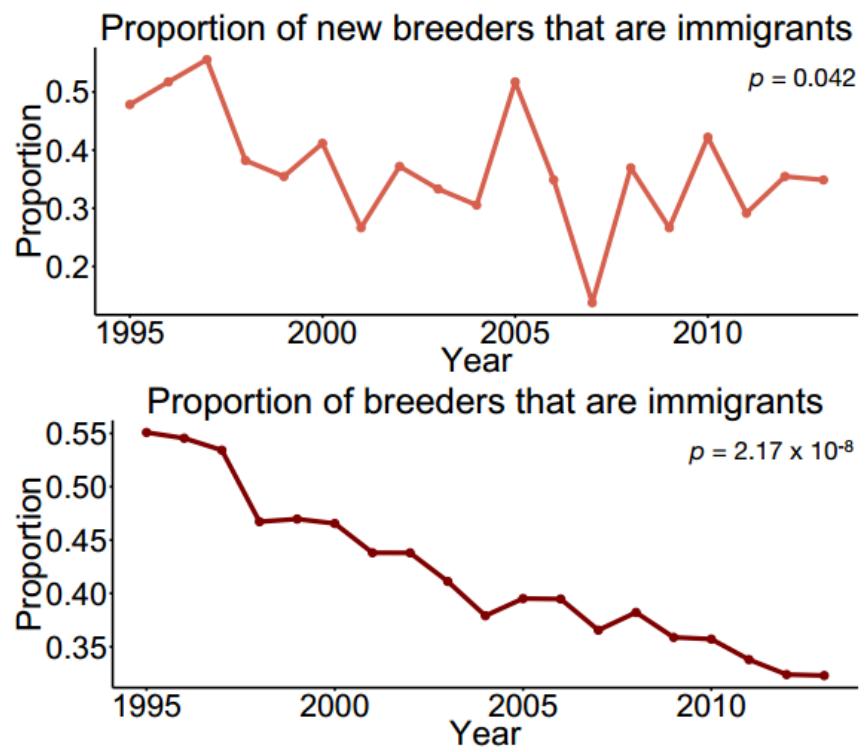
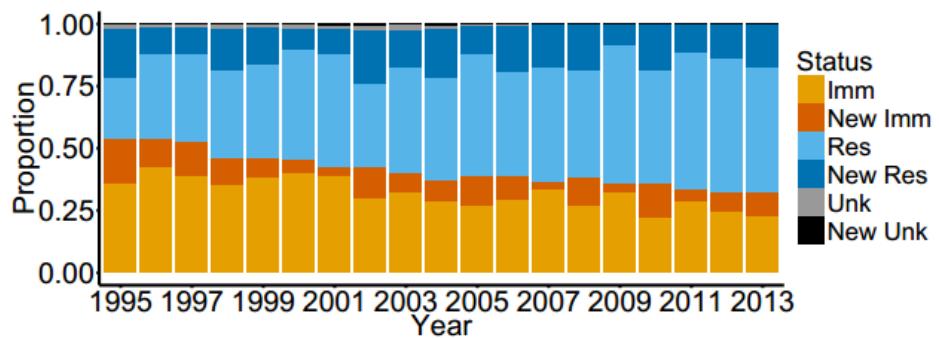
Immigration rate decreased over time



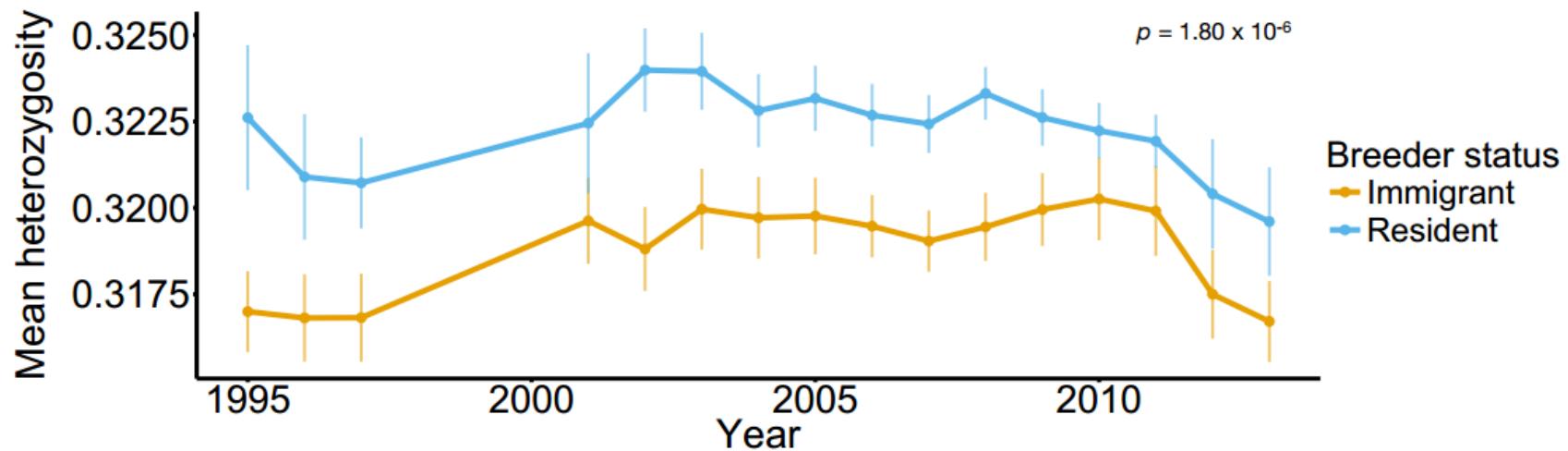
Immigration rate decreased over time



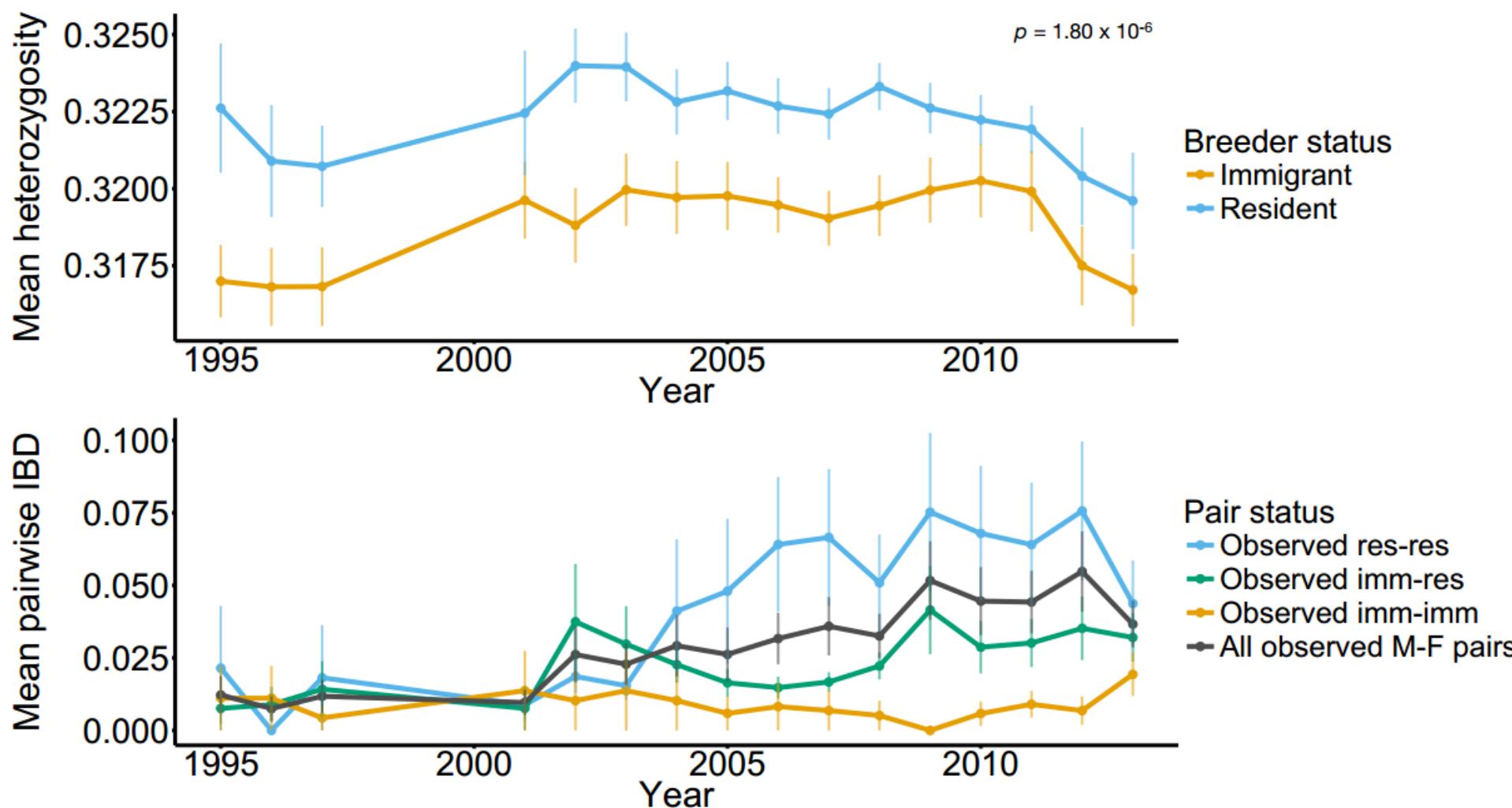
Immigration rate decreased over time



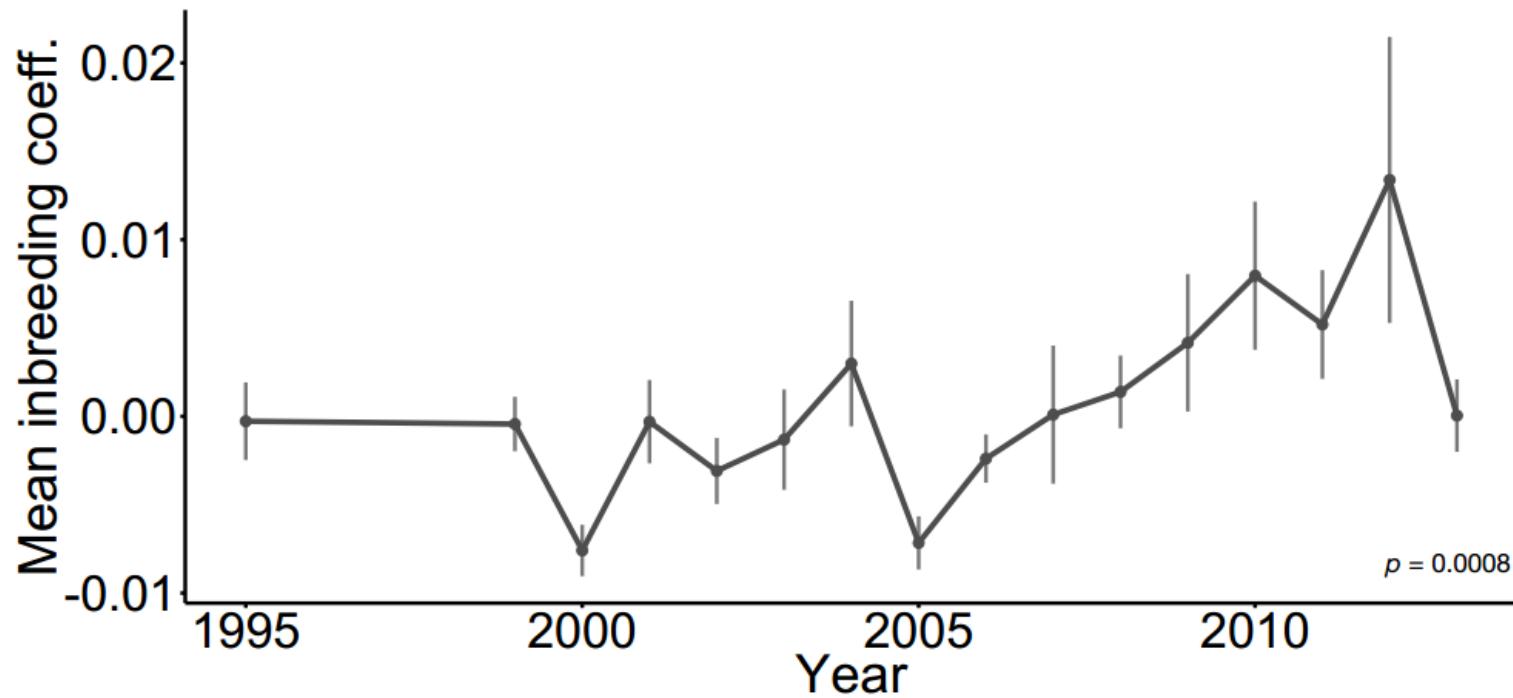
Immigrants were less heterozygous than residents



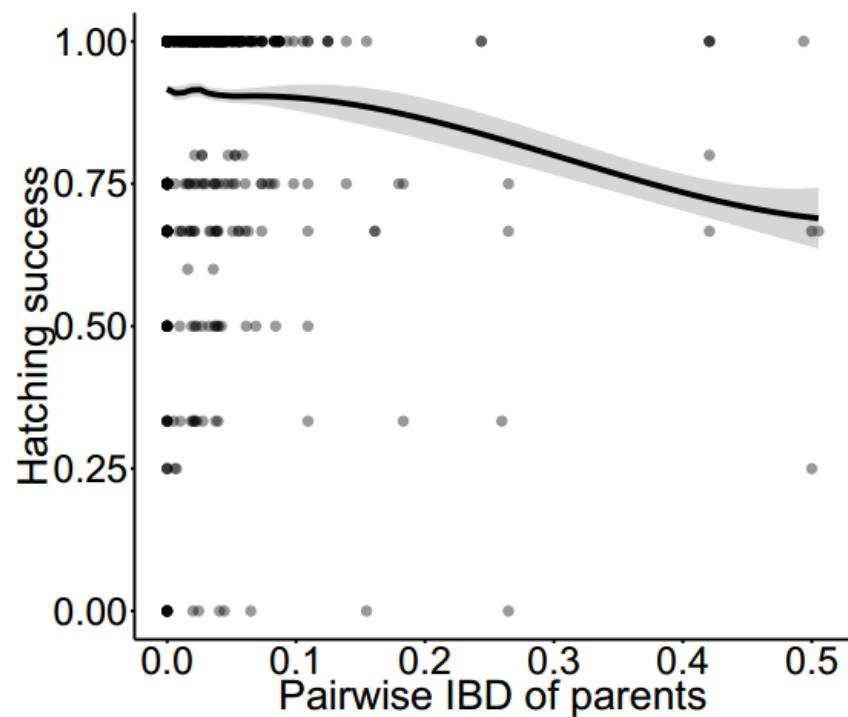
Immigrants were less heterozygous than residents
but still contributed genetic variation



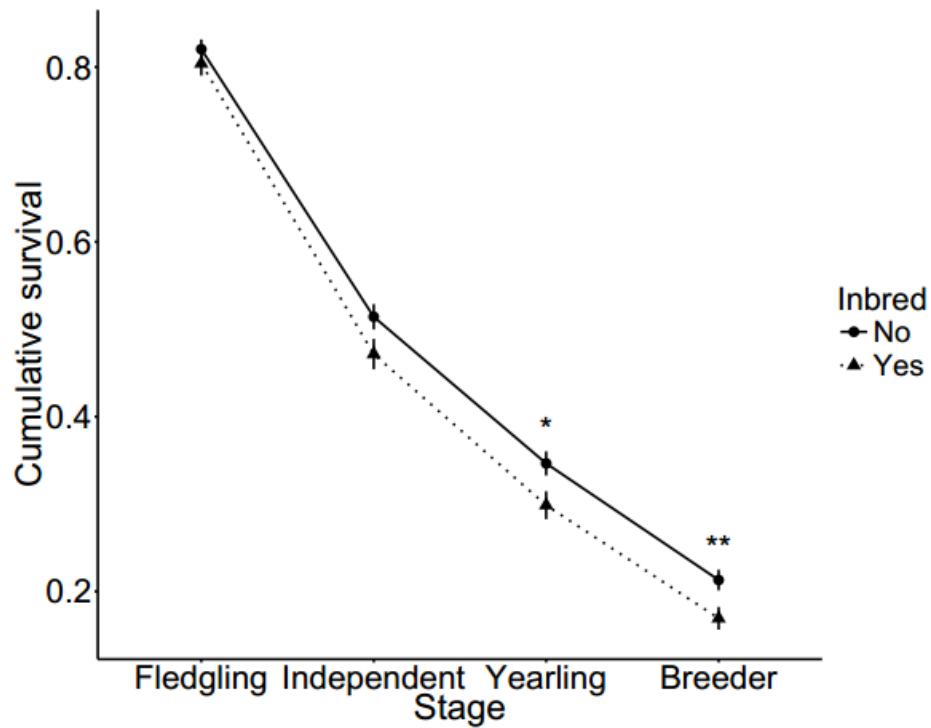
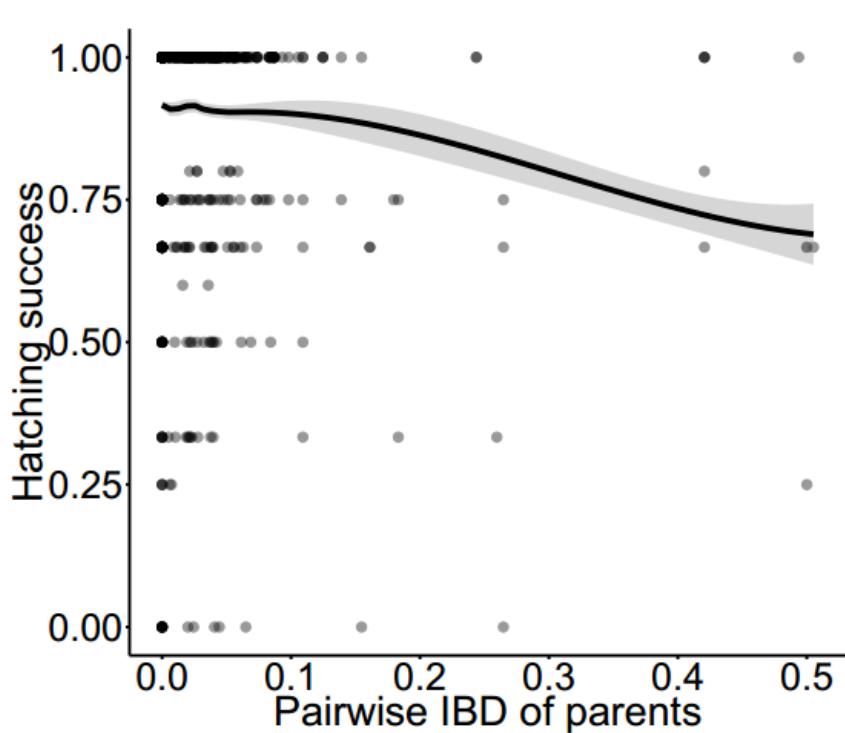
Mean inbreeding coefficient of the birth cohort increased over time



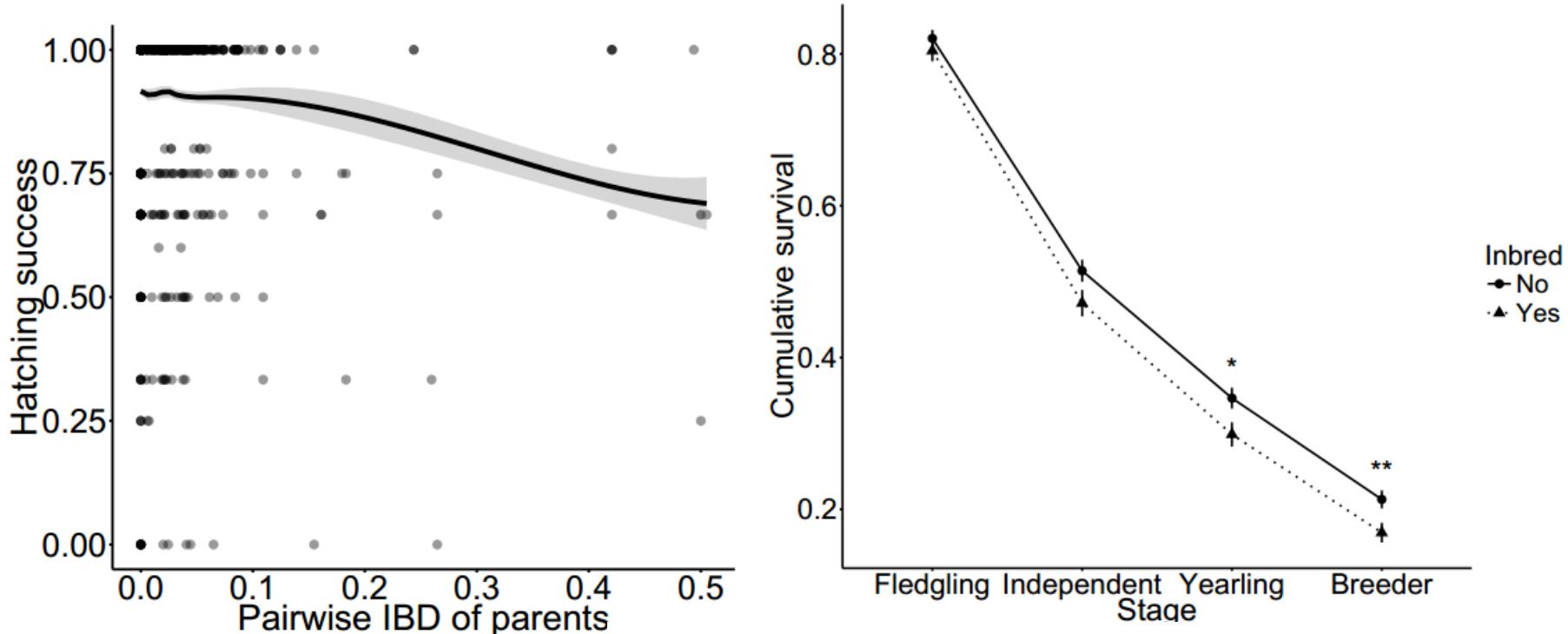
Inbreeding depression in multiple life-history stages



Inbreeding depression in multiple life-history stages

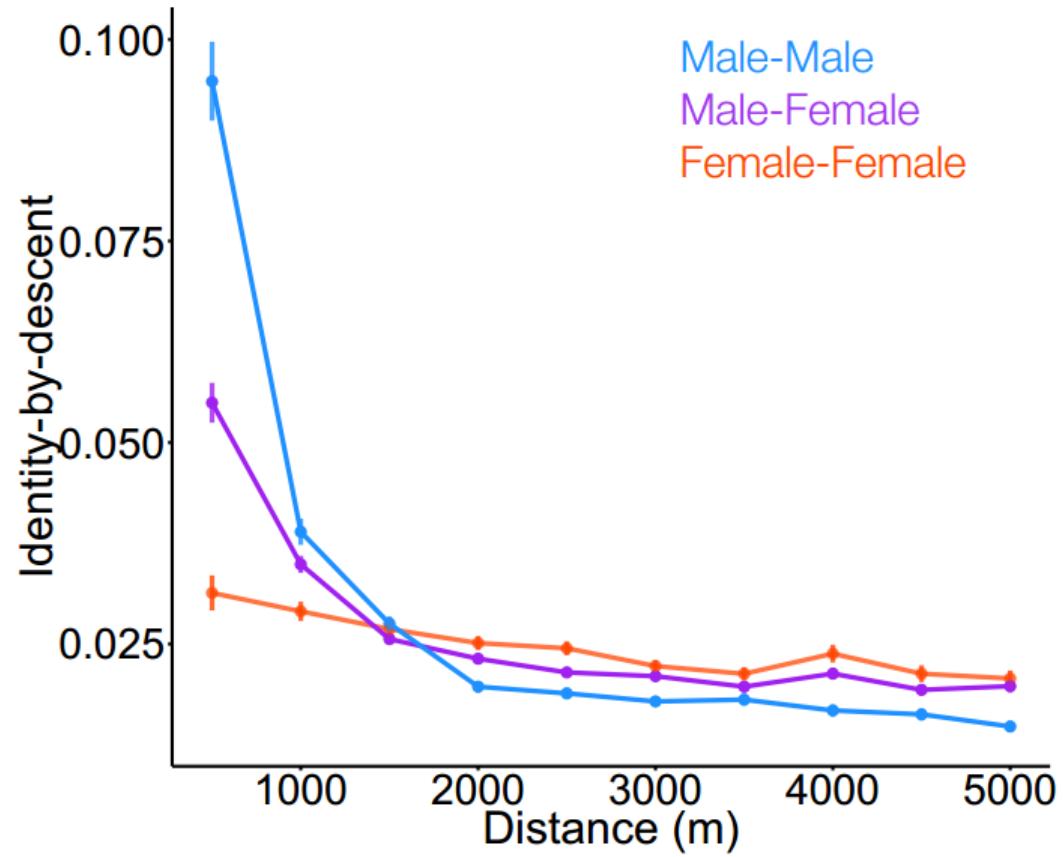
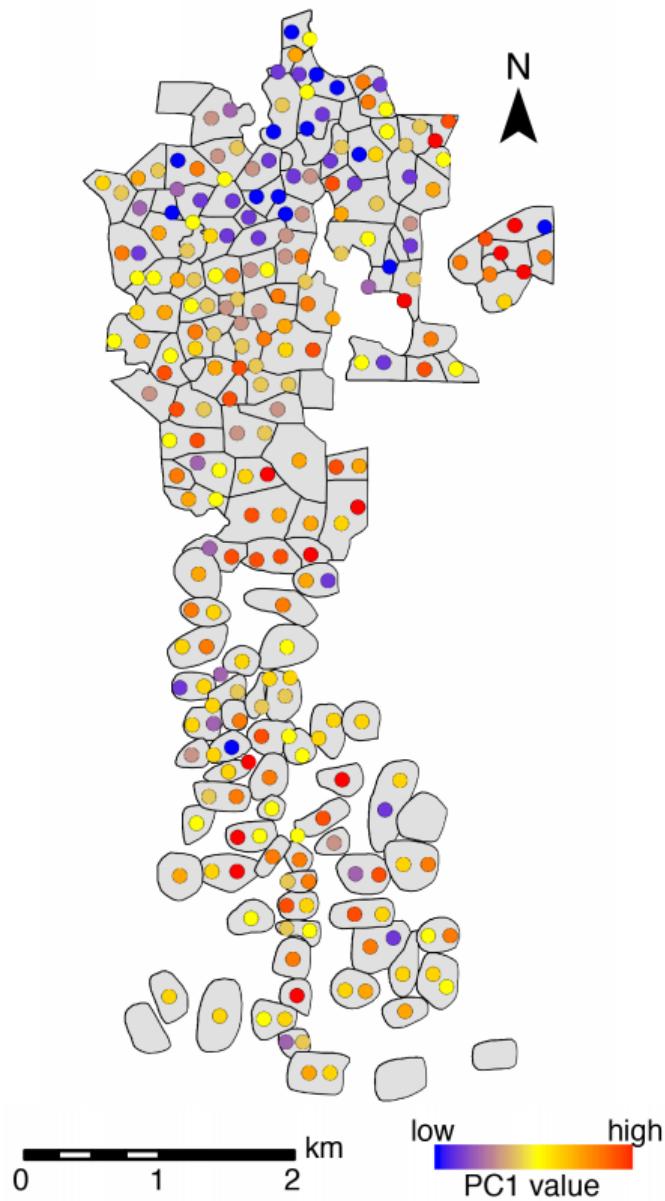


Inbreeding depression in multiple life-history stages

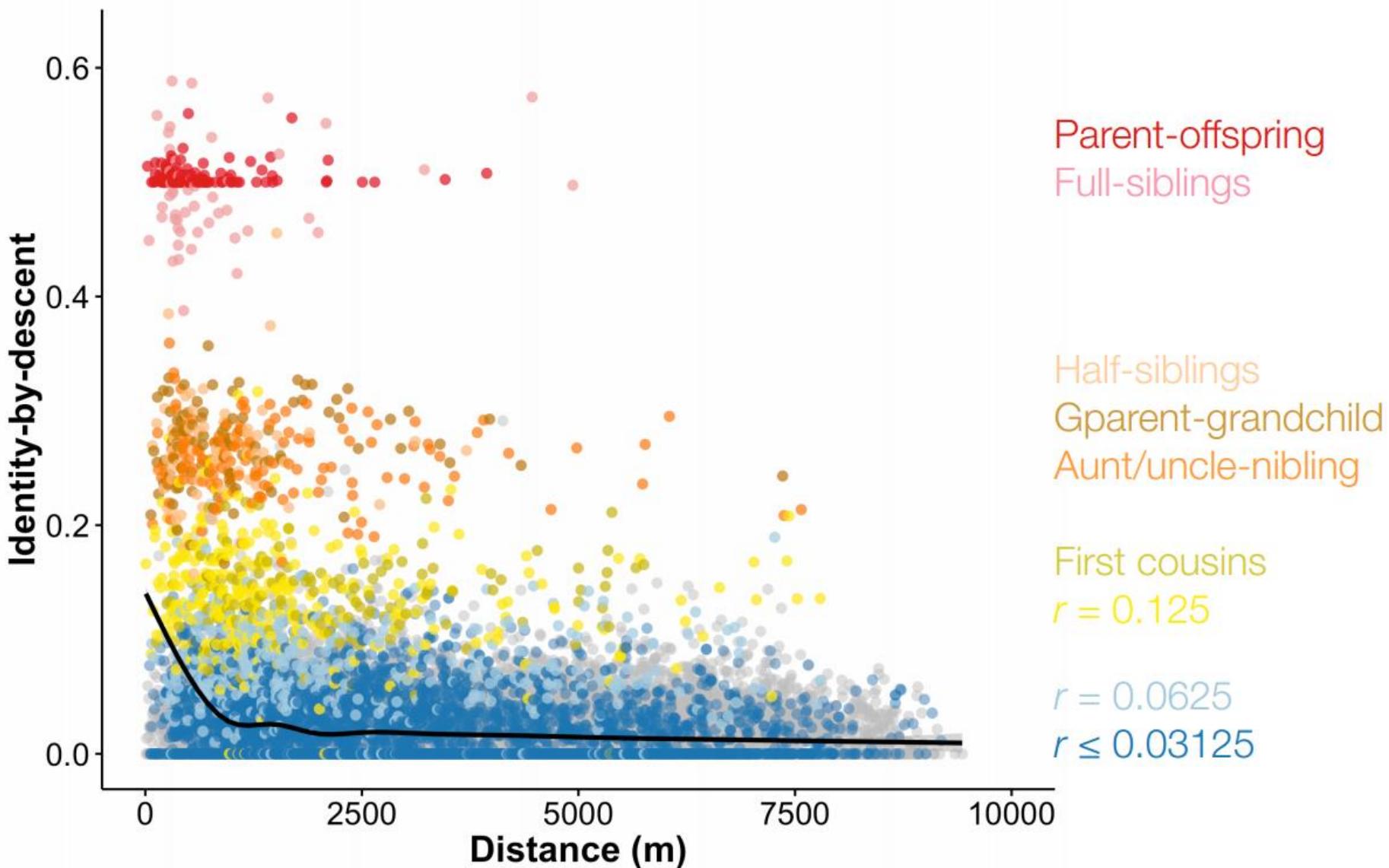


- ✓ Hatching success
- ✓ Nestling weight
- ✓ Juvenile survival
- ✓ Breeder lifespan
- ✓ Lifetime reproductive success

Limited dispersal leads to isolation-by-distance



A closer look at the isolation-by-distance pattern

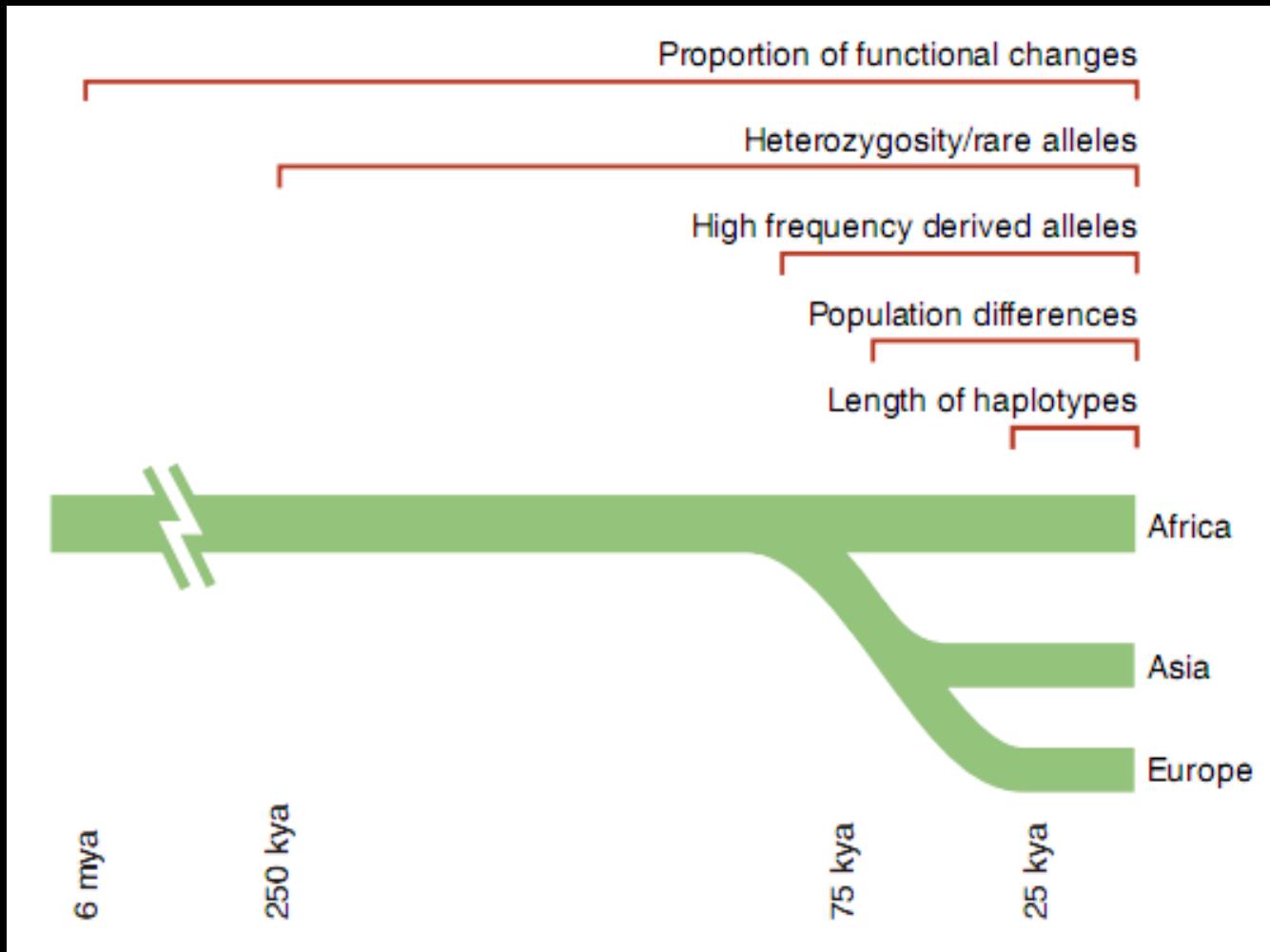


Natural Selection

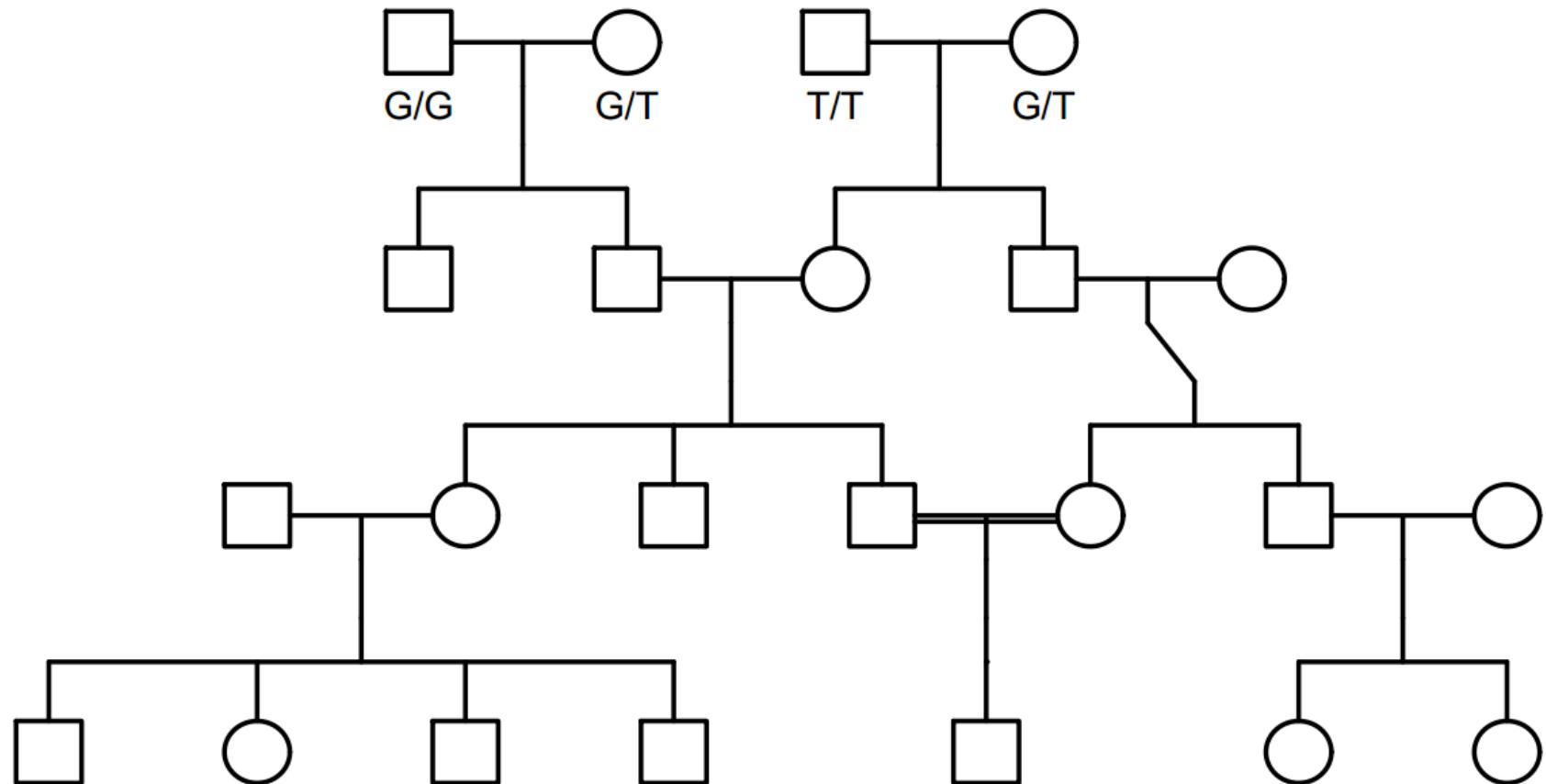
How best to do inference of selection
genome-wide?

Scans for selection

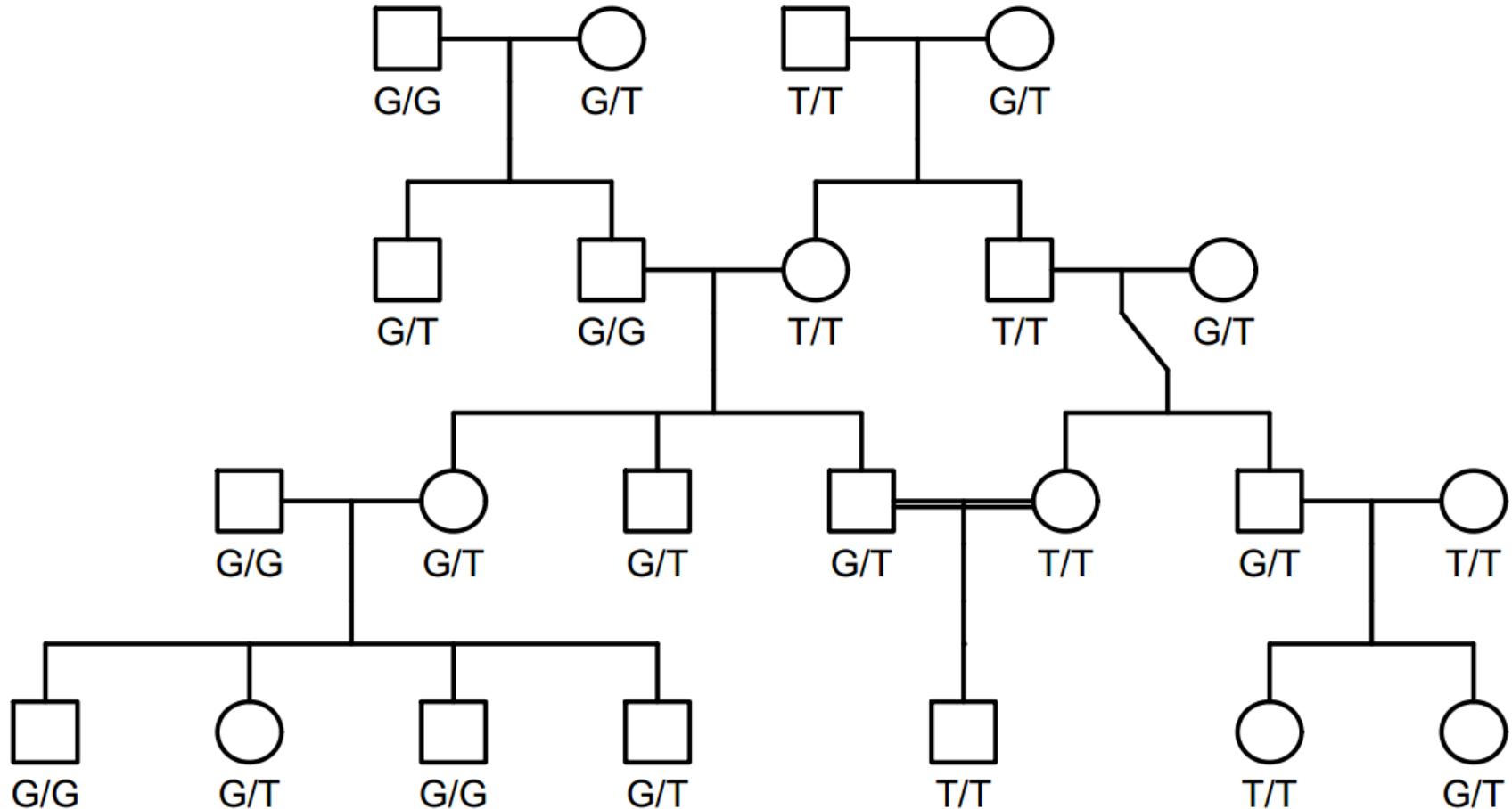
Different signals at different time depths



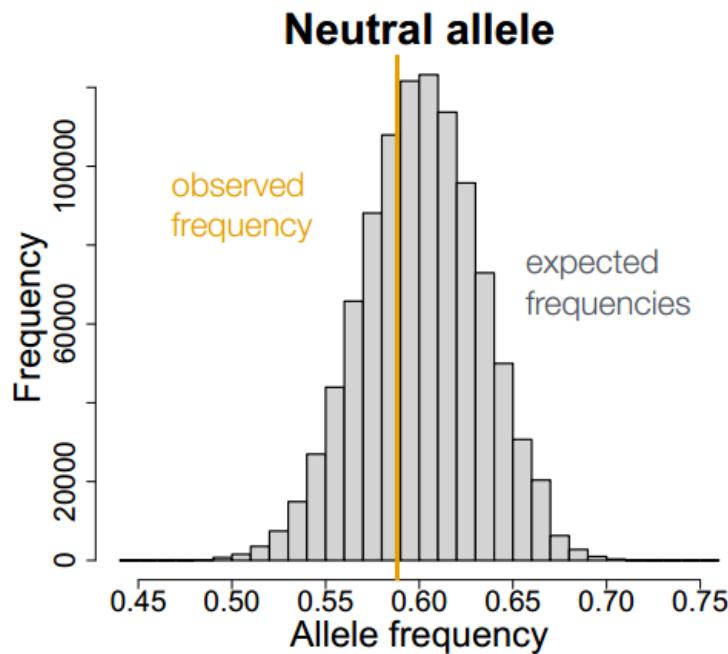
Testing for selection using pedigree information



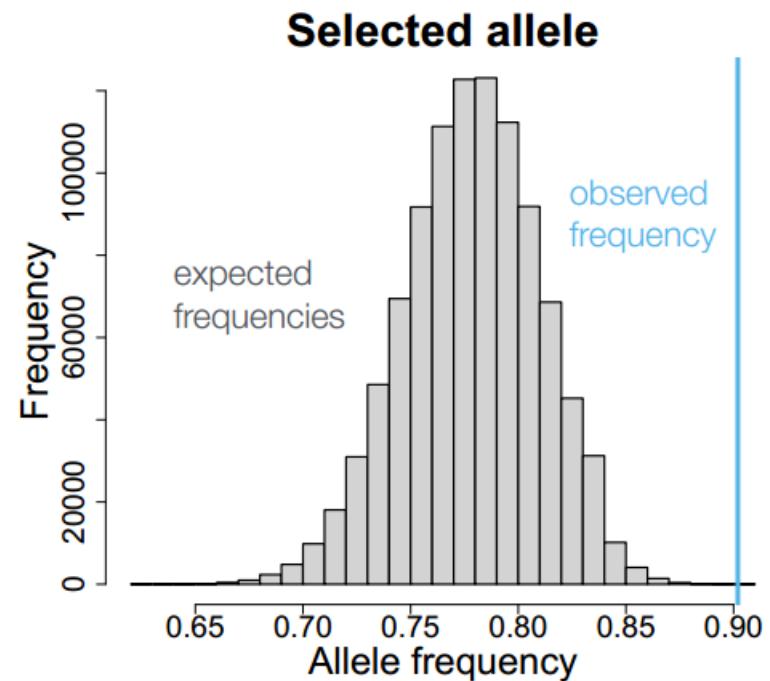
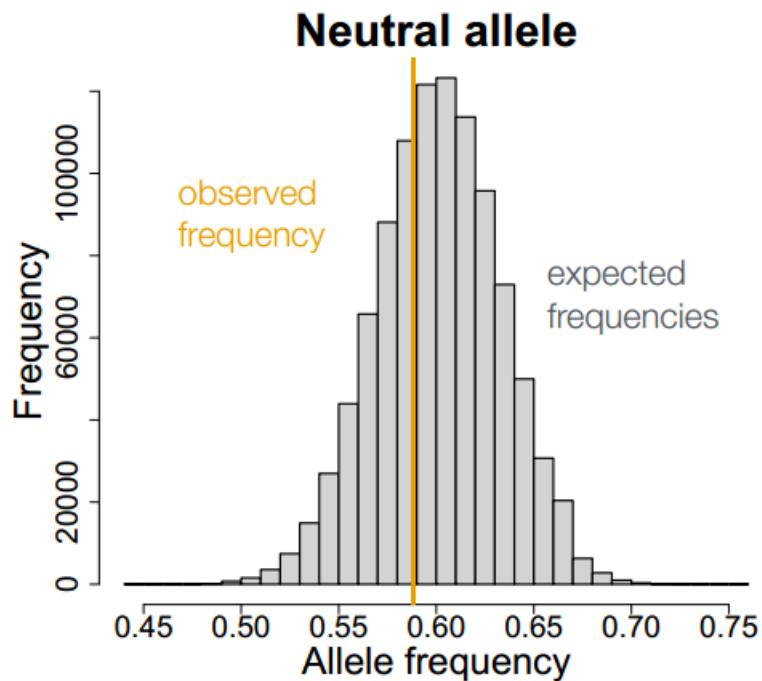
Gene-dropping to generate neutral expectations



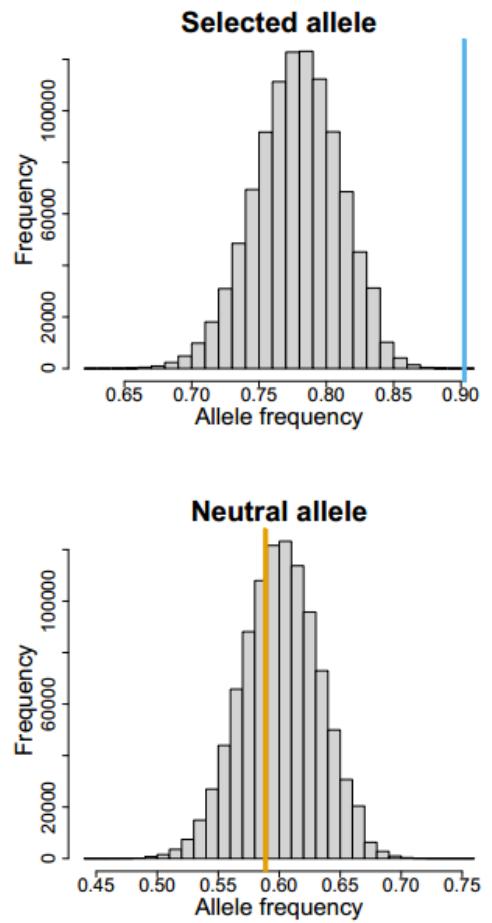
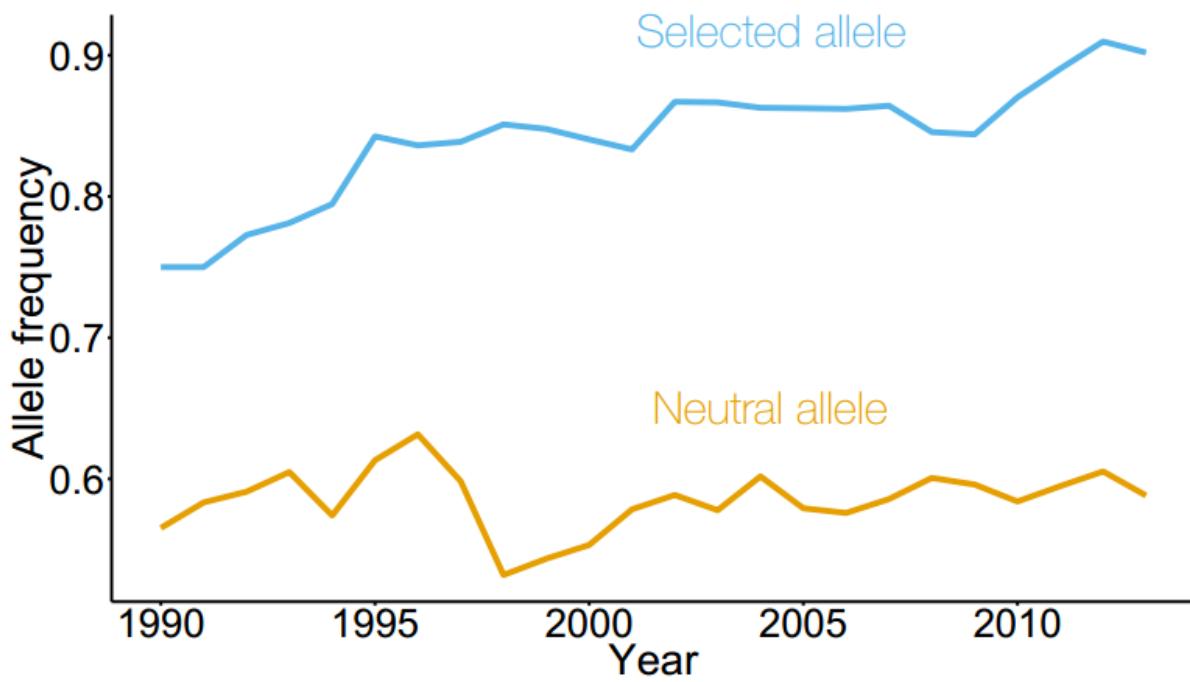
Most SNPs conform to the null hypothesis



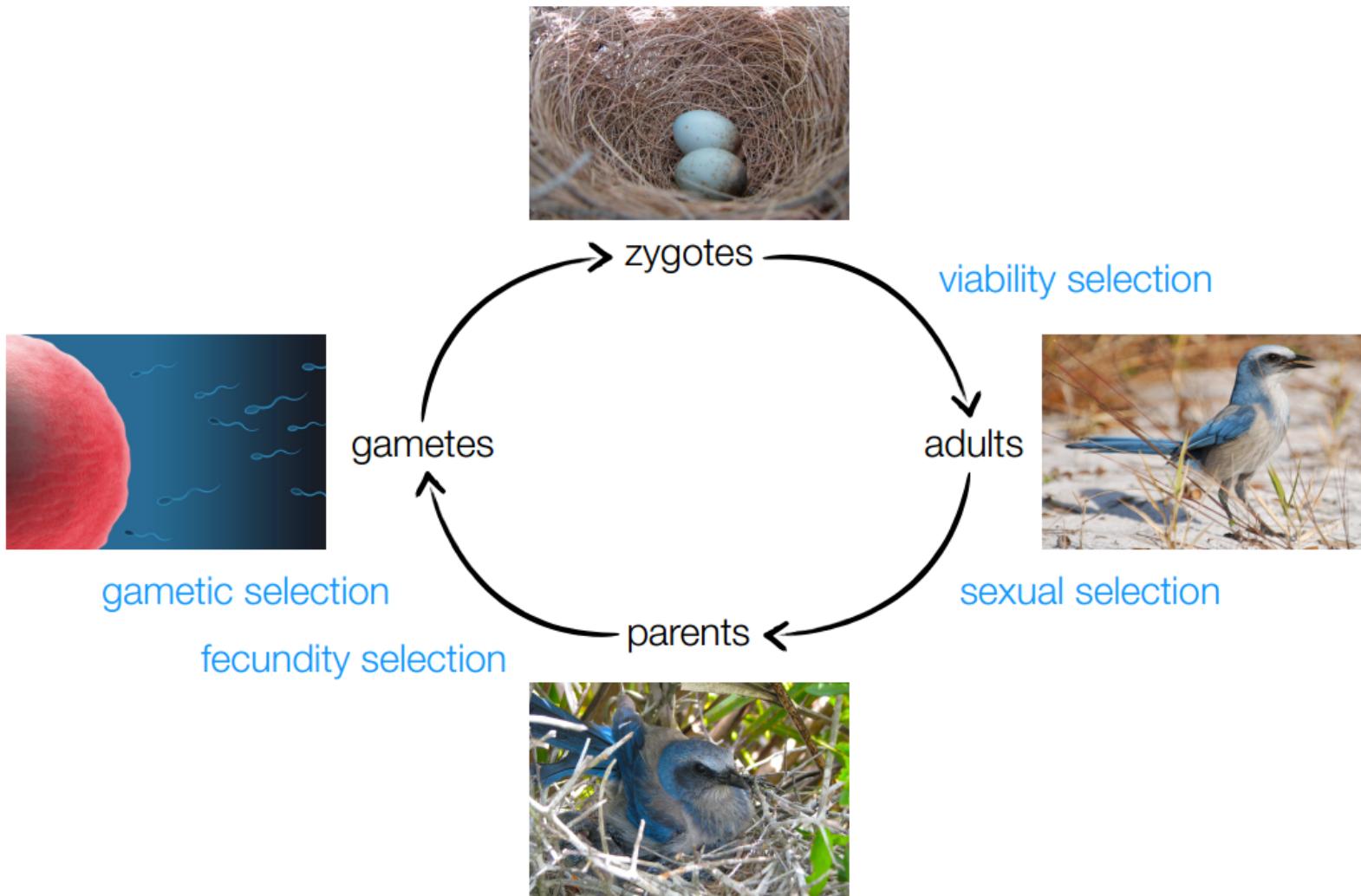
... but 67 SNPs display a significant departure.



Evidence of selection at 67 SNPs



Selection can act at different stages of the life cycle



Disease Association

How can variation in genome sequence be used to infer disease risk?

Population genetics and disease-causing variants

GWAS is fundamentally a population genetic idea – common variants that are causal to a phenotype are likely to show direct or indirect statistical association

Early onset disease-risk enhancing alleles ought to have lower allele frequency.

Rare variants can also be found by association, but require huge samples.

Variance explained



Peter Visscher

Genetic variance estimation with imputed variants finds negligible missing heritability for human height and body mass index

Jian Yang^{1,2,24}, Andrew Bakshi¹, Zhihong Zhu¹, Gibran Hemani^{1,3}, Anna A E Vinkhuyzen¹, Sang Hong Lee^{1,4}, Matthew R Robinson¹, John R B Perry⁵, Ilja M Nolte⁶, Jana V van Vliet-Ostaptchouk^{6,7}, Harold Snieder⁶, The LifeLines Cohort Study⁸, Tonu Esko⁹⁻¹², Lili Milani⁹, Reedik Mägi⁹, Andres Metspalu^{9,13}, Anders Hamsten¹⁴, Patrik K E Magnusson¹⁵, Nancy L Pedersen¹⁵, Erik Ingelsson^{16,17}, Nicole Soranzo^{18,19}, Matthew C Keller^{20,21}, Naomi R Wray¹, Michael E Goddard^{22,23} & Peter M Visscher^{1,2,24}

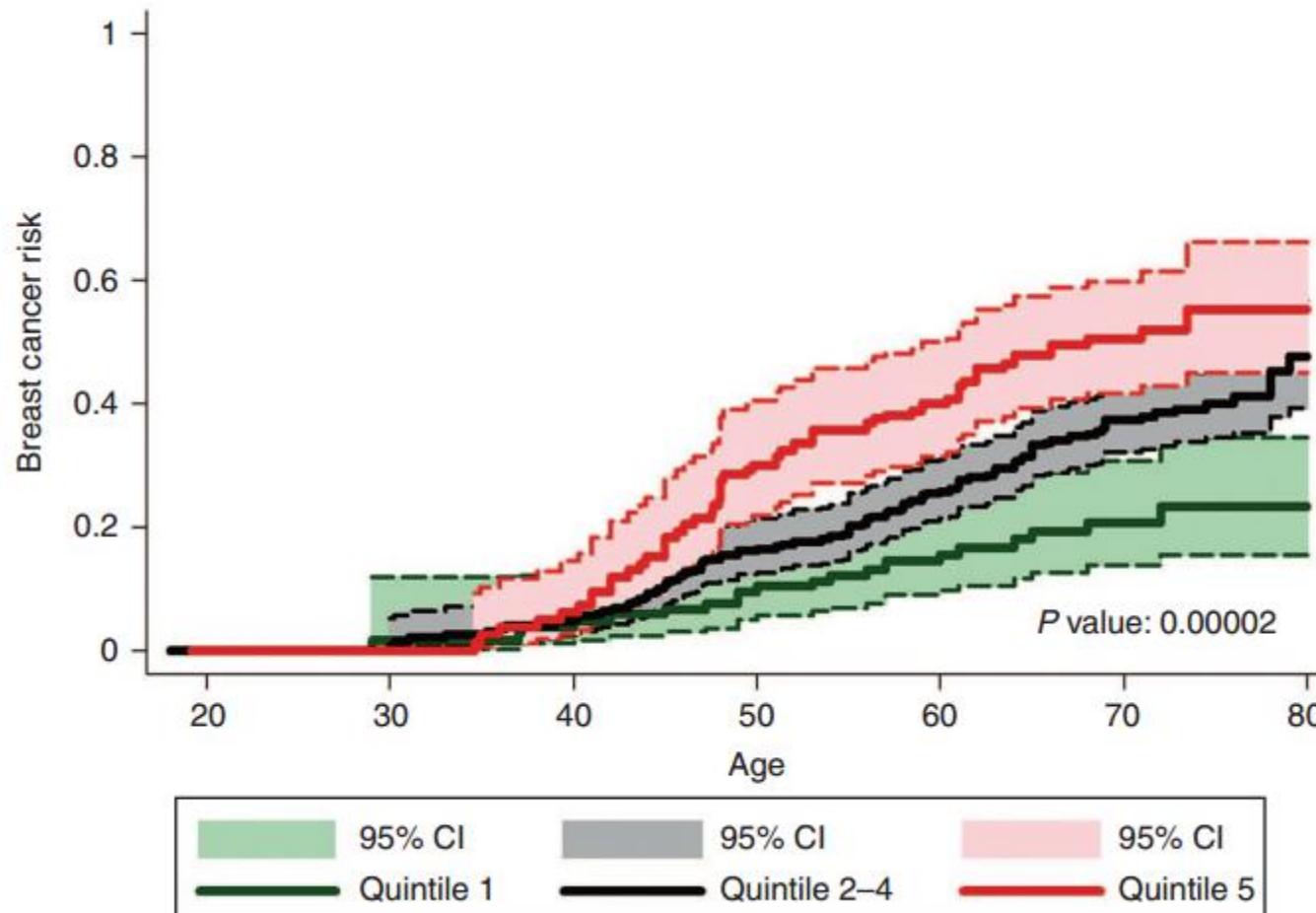
GCTA

The sum of infinitesimal effects of 17 M imputed SNPs yields 56% of variance.

GLM and prediction – polygenic risk score

- Purcell *et al.* (*Nature* 2009)
 - Polygenic Risk Score (schizophrenia $R^2 = 3\%$)
- Chen et al. (*Genet. Epidemiol.* 2015)
 - Big improvement with ancestry correction
 - hair color $R^2 = 4\text{-}7\%$, tanning ability $R^2 = 1\text{-}3\%$, basal cell carcinoma $R^2 = 1\text{-}2\%$
- Vilhalmsson *et al.* (2015 *AJHG*)
 - Direct modeling LD increases accuracy of PRS

Polygenic Risk Score and Breast Cancer Risk

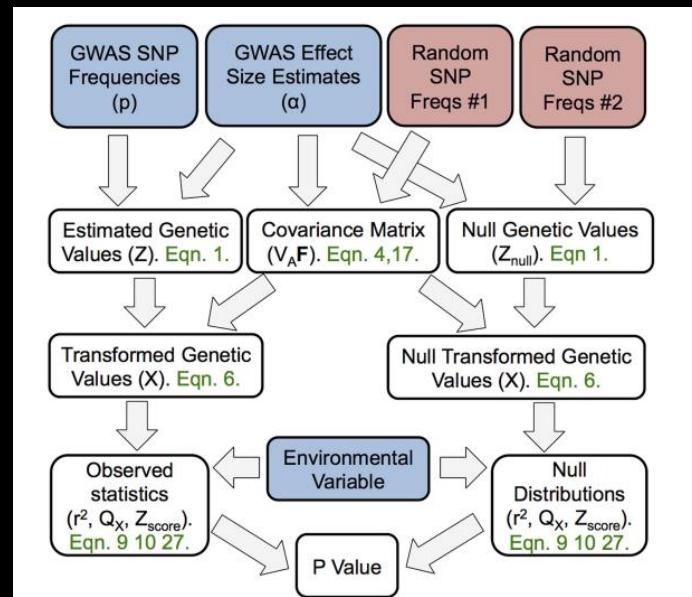


Polygenic selection

Estimate effect size from GWAS

Contrast population from different places or times

Ask if allele frequencies change in a correlated way



Evidence for selection on height in humans

GWAS on height finds many SNPs.

Many of these SNPs show a consistent cline from North to South in Europe.

There is a significant tendency of “short” alleles to be found in the south.

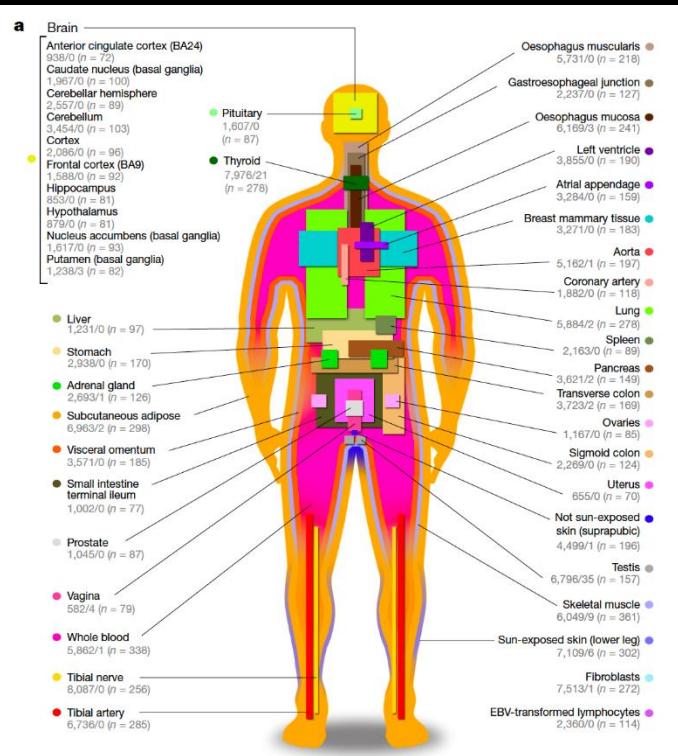
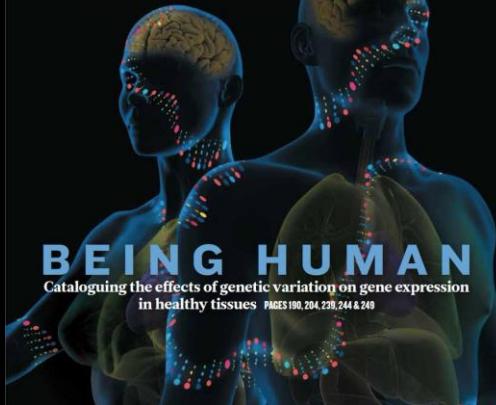
Caveats about Polygenic Selection

1. Population stratification in GWAS – errors in effect sizes.
2. Assumes constant effect sizes over space and time.
3. SNPs that impact fitness have inflated ascertainment, producing a bias

Genome functional variation

How can analysis of variation at the sequence and gene-expression levels inform us about genome function?

GTEX (Genotype-Tissue-Expression) Project



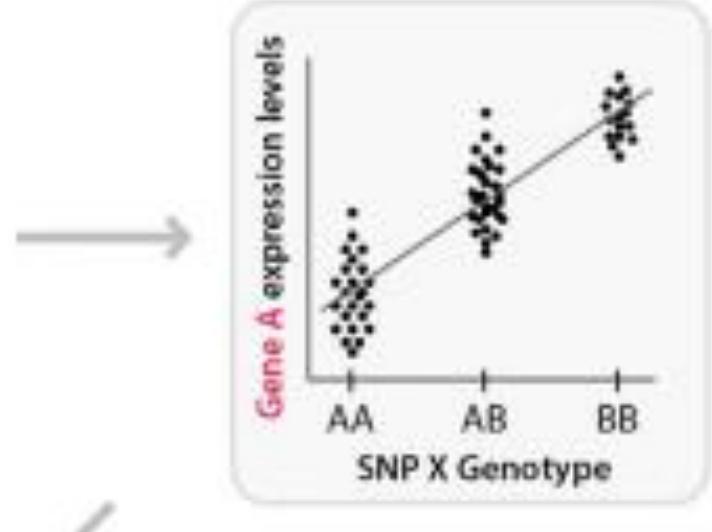
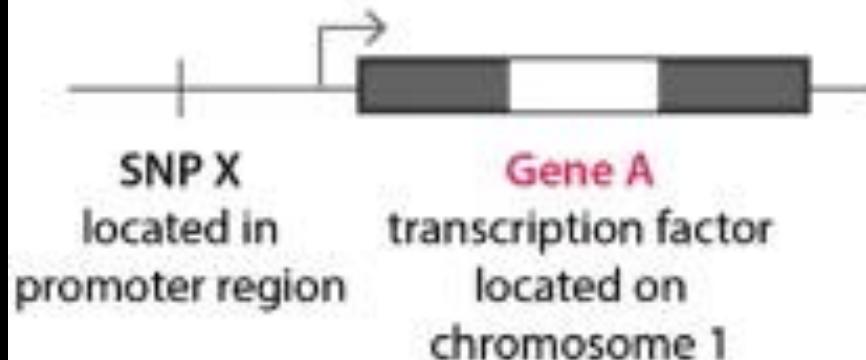
- Pilot (2015): 1658 samples from 175 donors
- Mid-phase (2017) : 7,051 samples from 450 donors
- Final (ongoing analysis): 17,382 from 838 donors
- Transcriptome: mRNA-seq
- Genotyping: < 2017: arrays, WES, WGS
- Banked biospecimens & procurement protocols
- Raw data (dbGap) + summary data (GTEX Portal) + browsers (Portal, UCSC, Ensembl)

Courtesy of Tuuli Lappalainen

Expression QTL have revealed regulatory variation

Cis-eQTL

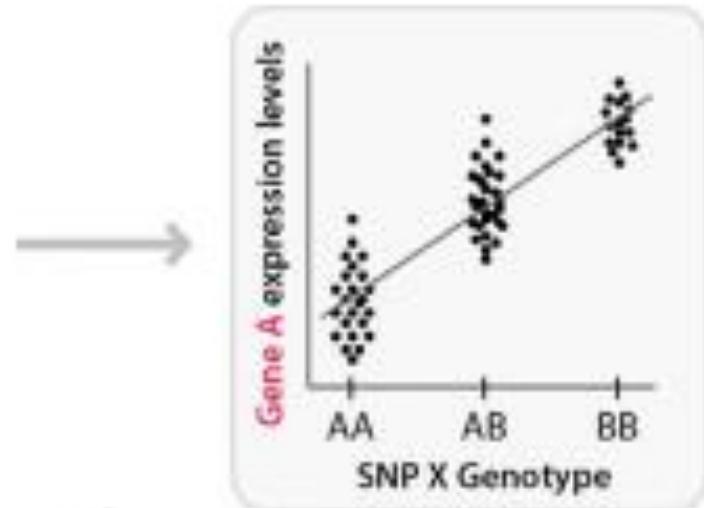
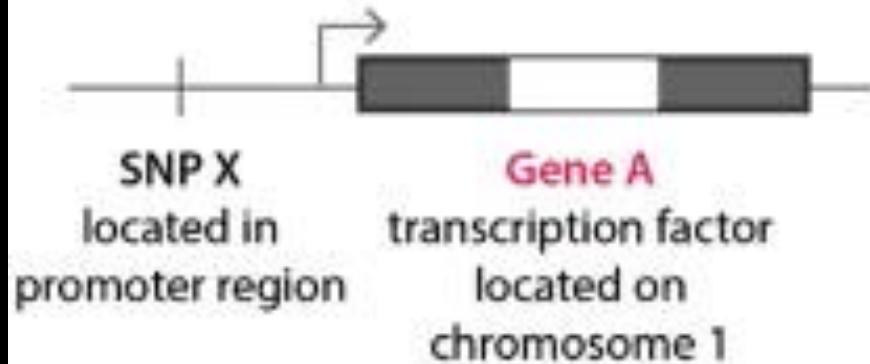
SNP X has an effect on local Gene A



Expression QTL have revealed regulatory variation

Cis-eQTL

SNP X has an effect on local Gene A

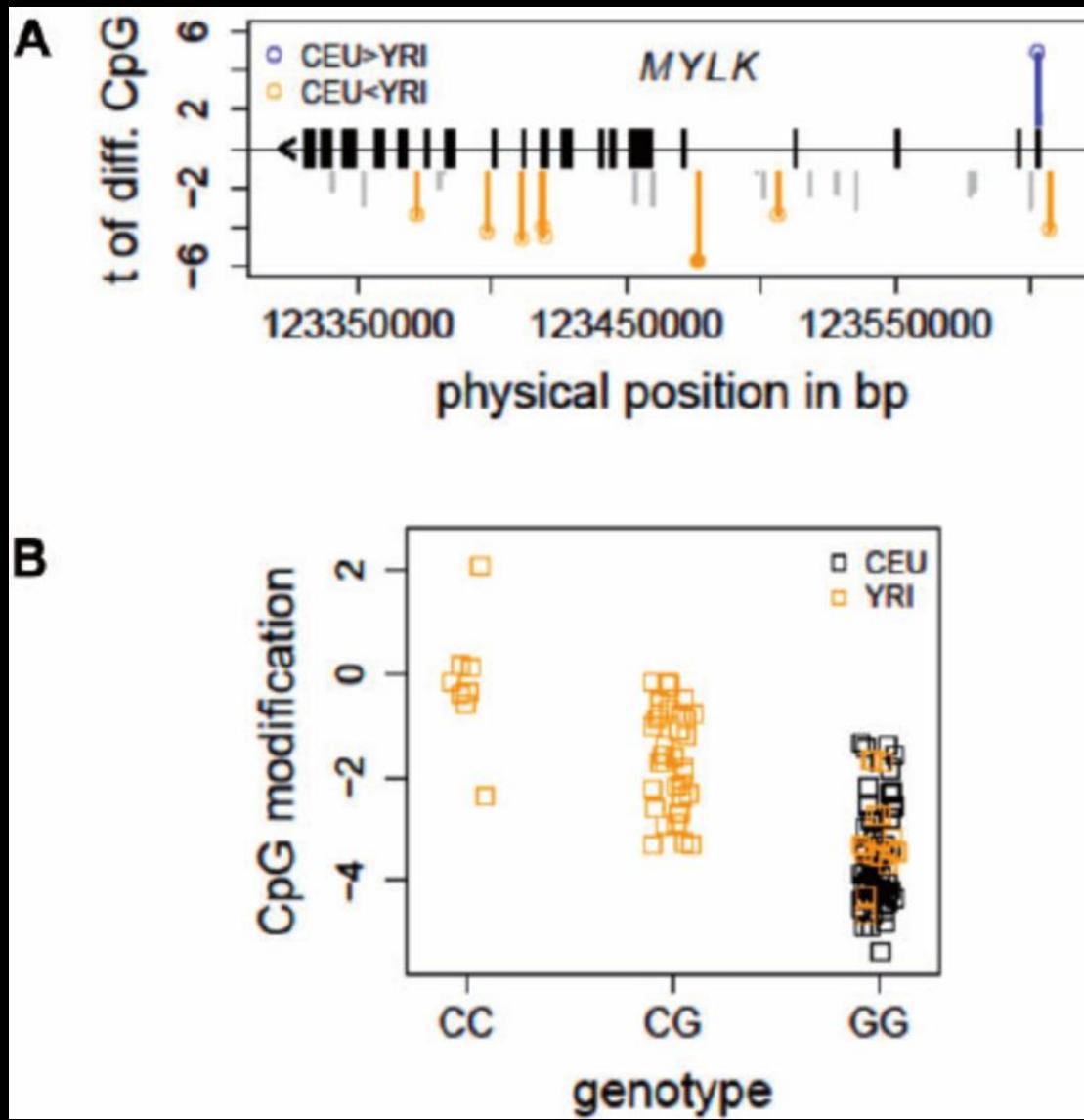


Altered Protein A levels,
effect on the binding to
the transcription factor
binding sites of
downstream genes

Trans-eQTL

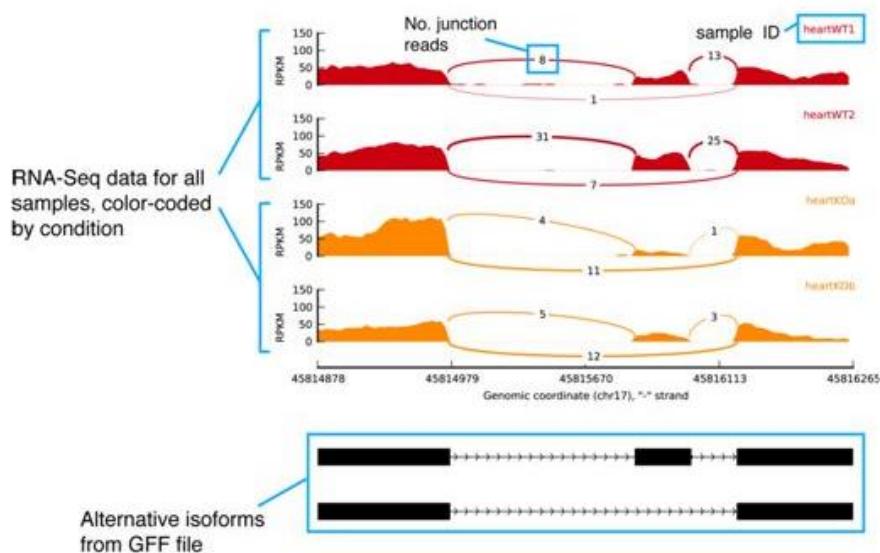
SNP X has an effect on distant Gene B through an intermediary factor (such as a transcription factor)

Methylation QTL (mQTL)

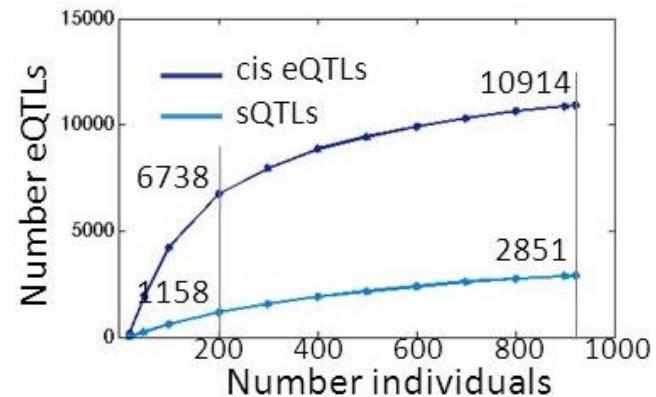


Splicing QTL – polymorphism for splicing

Can investigate relative transcript ratios or reads across junctions.



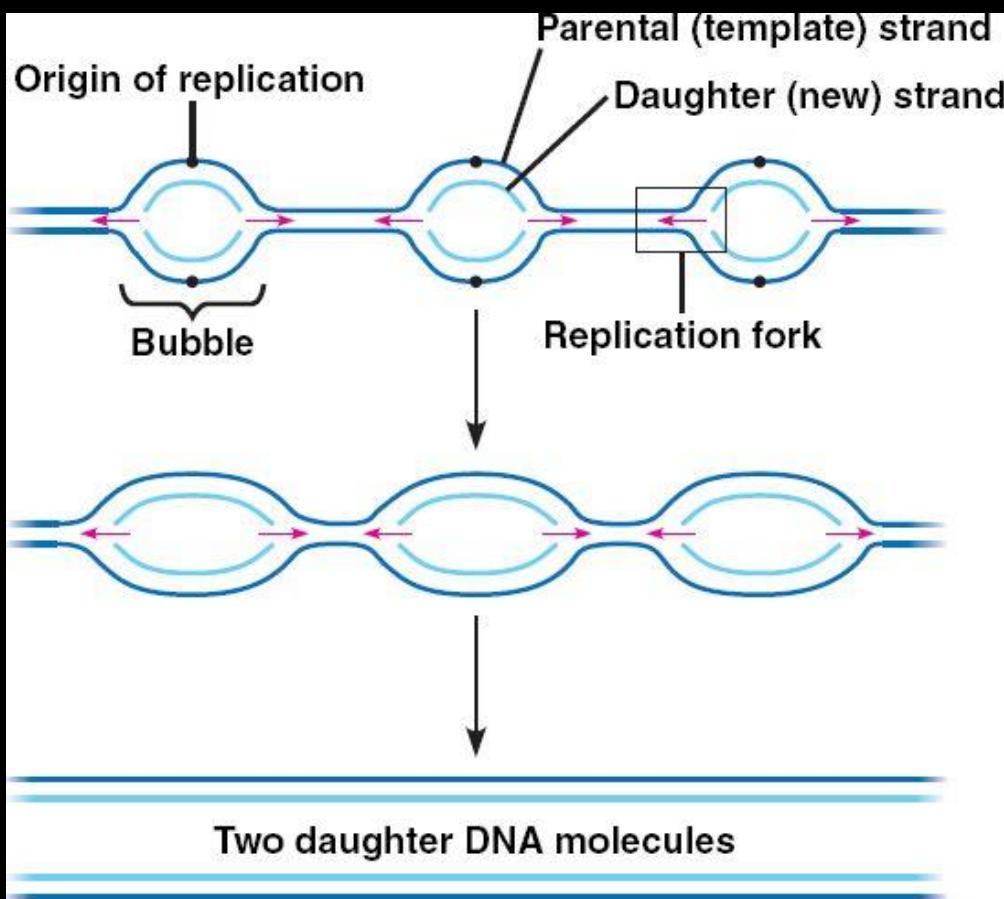
- Splicing also affected for many genes



Katz et al, Nature Methods, 2010

Battle et al, Genome Research, 2014

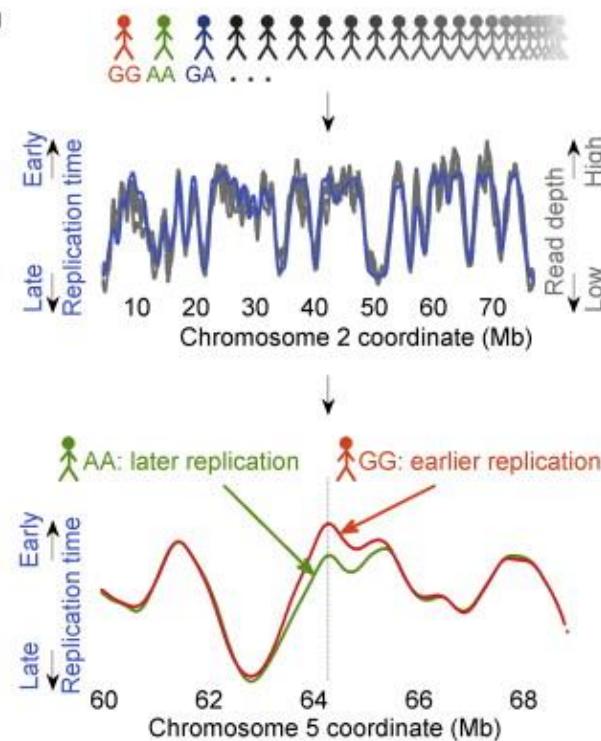
Replication timing QTL – variation in DNA copying!



Population sequencing
using proliferating
cell cultures

Read depth
along chromosomes
~
DNA replication
timing

Replication timing
quantitative trait loci
(rtQTLs)



Outline

- Demographic inference
- Population structure and history
- Admixture/ Introgression
- Random genetic drift
- Natural selection
- Mutation spectrum
- Disease association
- Genome function