

**MUTATIONS IN TIME:**

SOME BASICS OF  
POPULATION GENETICS

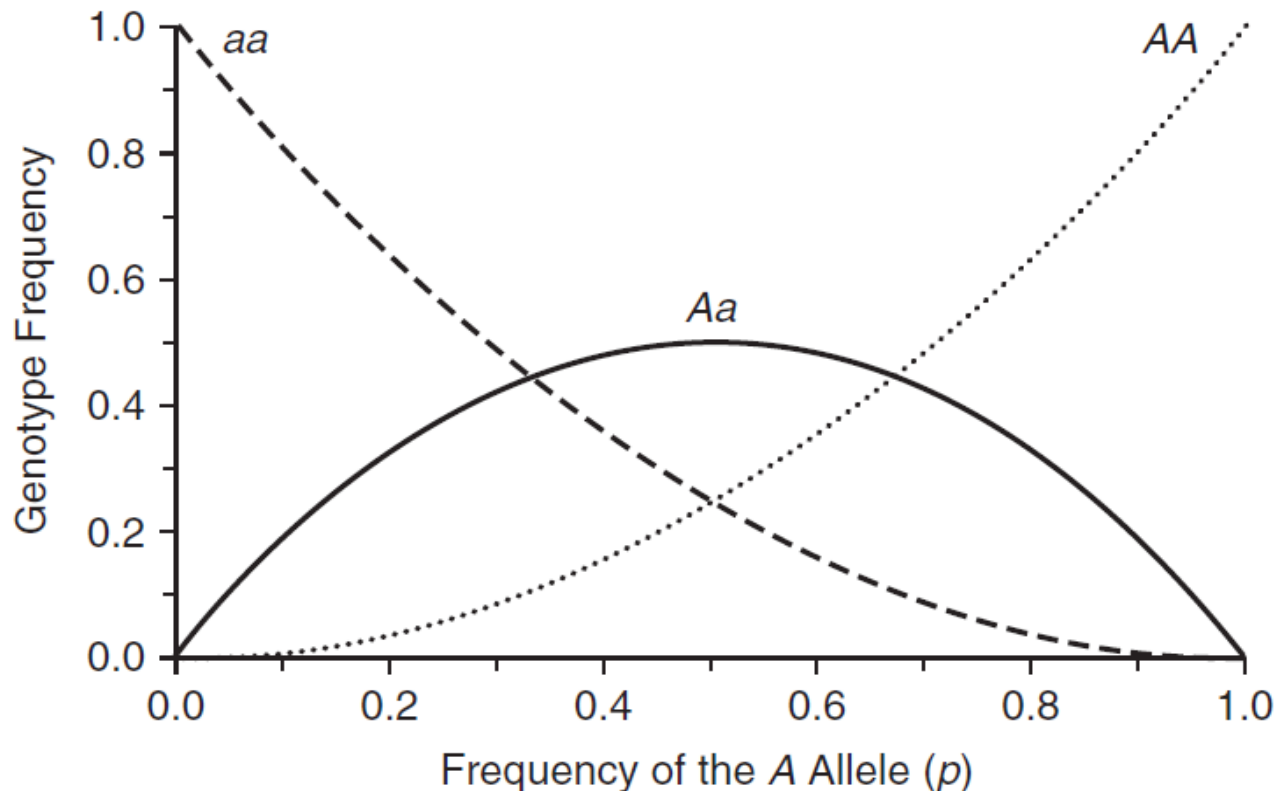
# Lecture plan

- Hardy-Weinberg equilibrium
- Random genetic drift without mutations
- Effective population size
- Random genetic drift and mutations
- The coalescent theory
- Natural selection. Mutation-selection balance
- Random genetic drift, positive selection
- Selection coefficients, deleterious alleles
- Non-random mating, population subdivision, gene flow, admixture, adaptation

# Hardy-Weinberg equilibrium (1908)

Generation  $N$  :  $f_A = p$ ,  $f_a = q$ ,  $p + q = 1$

Generation  $N + 1$  :  $F_{AA} = p^2$ ,  $F_{Aa} = 2pq$ ,  $F_{aa} = q^2$



# Hardy-Weinberg equilibrium

Generation  $N$  :  $f_A = p$ ,  $f_a = q$ ,  $p + q = 1$

Generation  $N + 1$  :  $F_{AA} = p^2$ ,  $F_{Aa} = 2pq$ ,  $F_{aa} = q^2$

Implications:

1. The allele frequencies does not change:

$$p' = f'_A = F'_{AA} + F'_{Aa}/2 = p^2 + pq = p$$

*Exercise: derive this*

2. HWE frequencies are attained in one generation



# Hardy-Weinberg equilibrium

## Assumptions:

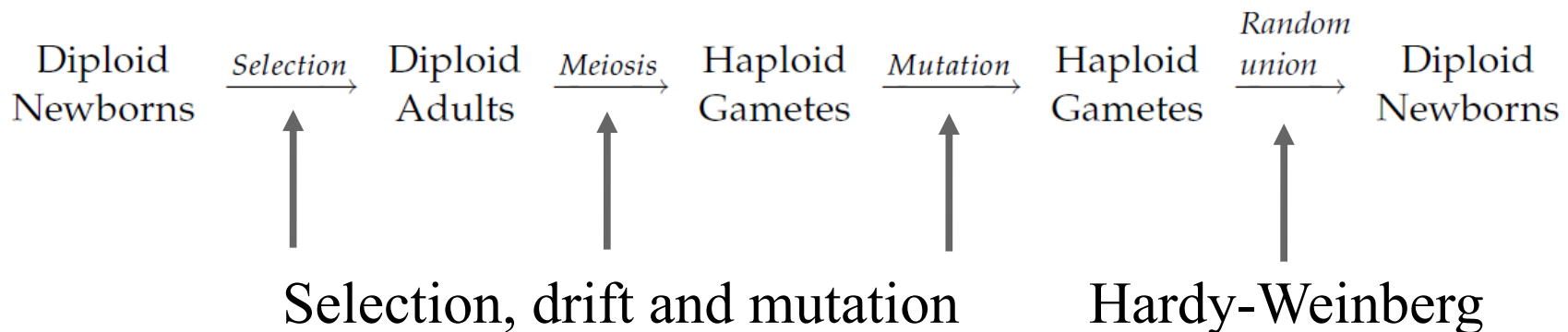
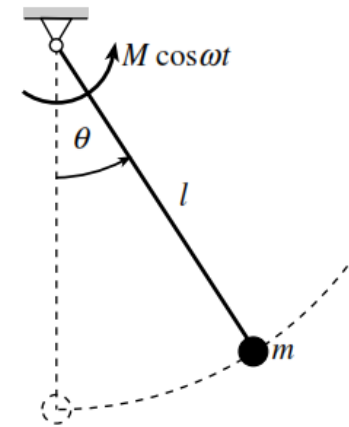
- Diploid species with sexual reproduction and random (not assortative) mating
- Same allele frequencies in males and females
- Non-overlapping generations
- Biallelic (autosomal) locus
- Population size is infinite
- No change in allele frequencies by migration, natural selection or mutation
- No genotyping errors

# Hardy-Weinberg equilibrium

Does it still make sense with so many assumptions? Yes:

1. A baseline for more realistic models

2. The H-W model splits life history into two intervals: gametes  $\rightarrow$  zygotes and zygotes  $\rightarrow$  adults



# Hardy-Weinberg equilibrium

## Testing for HWE:

$df = n - k - 1$ , where  $n = 3$  is the number of classes and  $k = 1$  is the number of independent parameters

| Genotype | Observed Number (O) | Expected Number (E) | (O - E) | (O - E) <sup>2</sup> | (O - E) <sup>2</sup> /E |
|----------|---------------------|---------------------|---------|----------------------|-------------------------|
| AA       | 90                  | 83.2                | 6.8     | 46.24                | 0.5558                  |
| Aa       | 28                  | 41.6                | -13.6   | 184.96               | 4.4462                  |
| aa       | 12                  | 5.2                 | 6.8     | 46.24                | 8.8923                  |

After performing the calculations in this table, we get a chi-square ( $\chi^2$ ) statistic of

$$\chi^2 = 0.5558 + 4.4462 + 8.8923 = 13.8943$$

This value is *much* larger than the critical value of 3.841, so we reject the hypothesis of Hardy-Weinberg equilibrium.

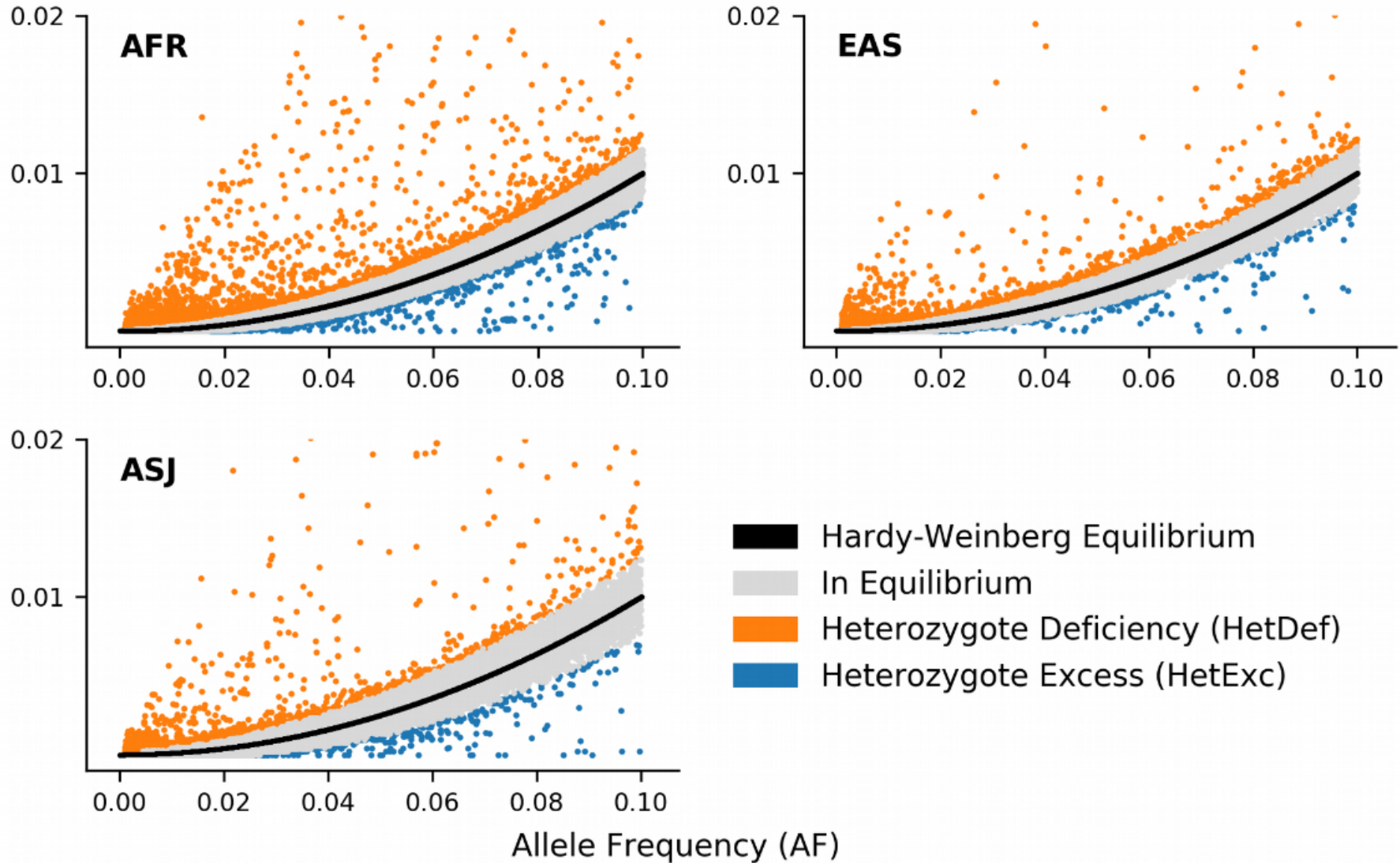
$$\chi^2 = \sum \frac{(O - E)^2}{E}$$

*Exercise: do it yourself*

# Hardy-Weinberg Equilibrium in the Large Scale Genomic Sequencing Era

 Nikita Abramovs,  Andrew Brass,  May Tassabehji

doi: <https://doi.org/10.1101/859462>



**gnomAD:** 137,842 predominantly healthy individuals from 7 major ethnic populations

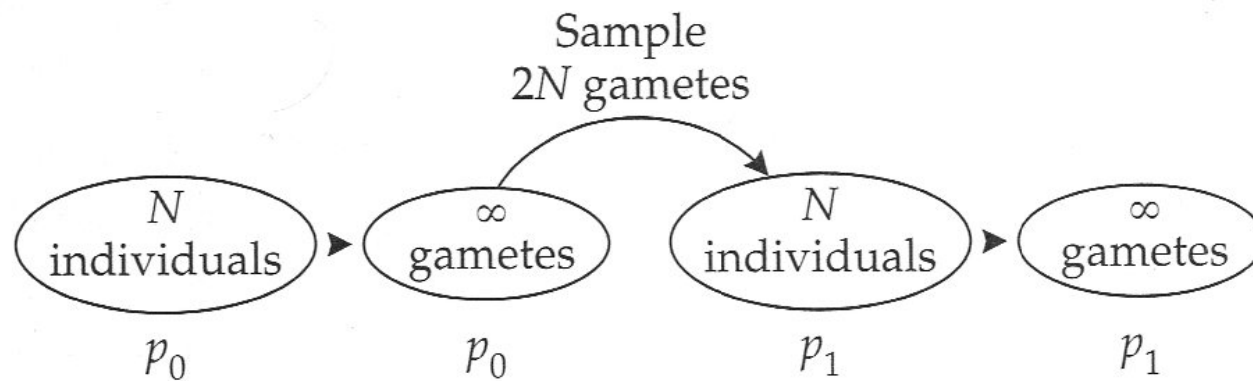
# Random genetic drift (Wright-Fisher, 1930)

## Assumptions:

- Diploid species with sexual reproduction and random (not assortative) mating
- Same allele frequencies in males and females
- Non-overlapping generations
- Biallelic (autosomal) locus
- ~~Population size is infinite~~
- No change in allele frequencies by migration, natural selection or mutation
- No genotyping errors

# Random genetic drift

Finite population  $\Rightarrow$  Sampling variation  $\Rightarrow$   
Allele frequency fluctuations  $\Rightarrow$  Random genetic drift



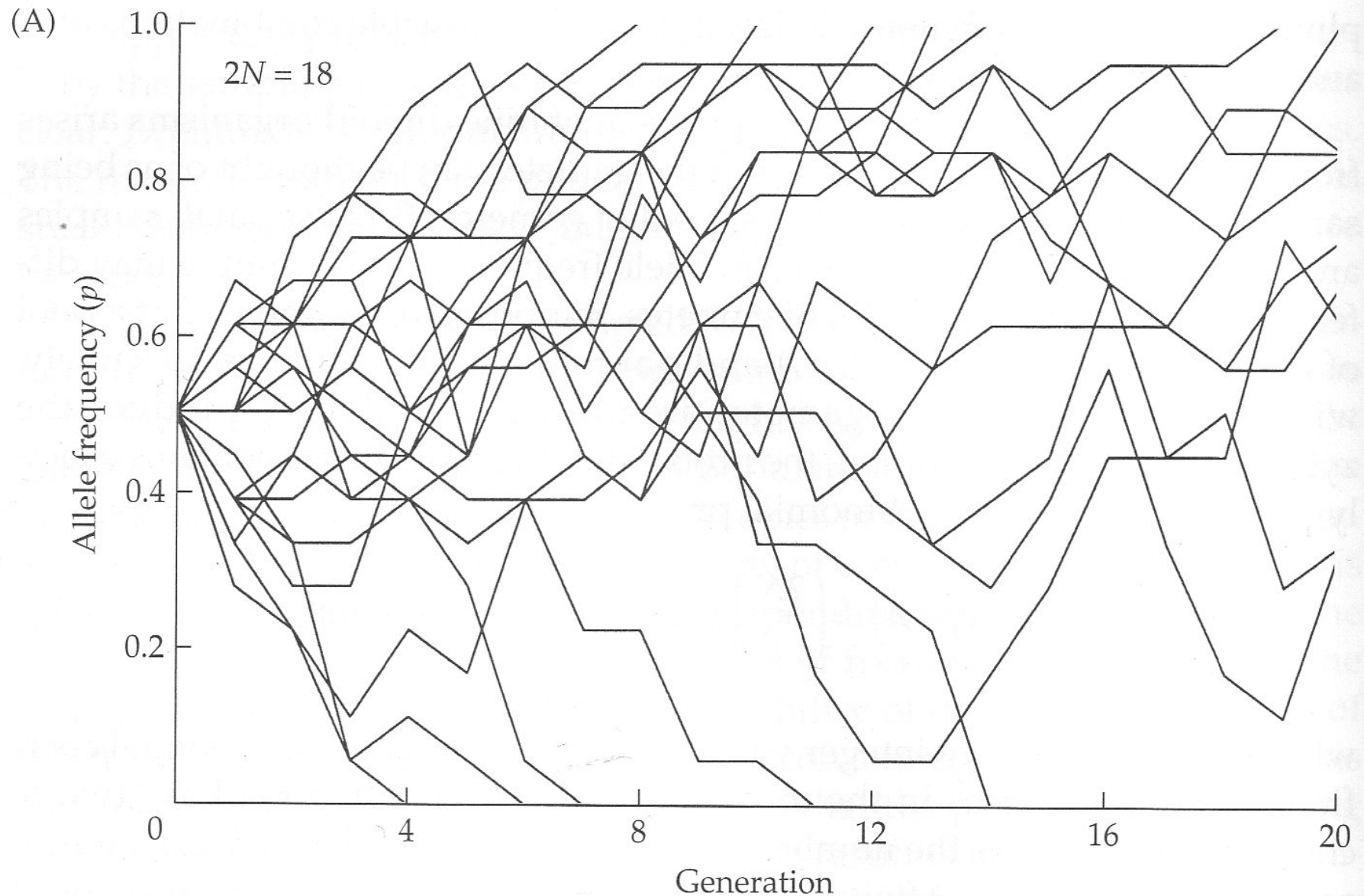
$$P(k) = \binom{2N}{k} p_0^k (1 - p_0)^{2N-k}$$

$$E(\Delta p | p) = E(k/2N - p | p) = 0$$

*Exercise: derive*

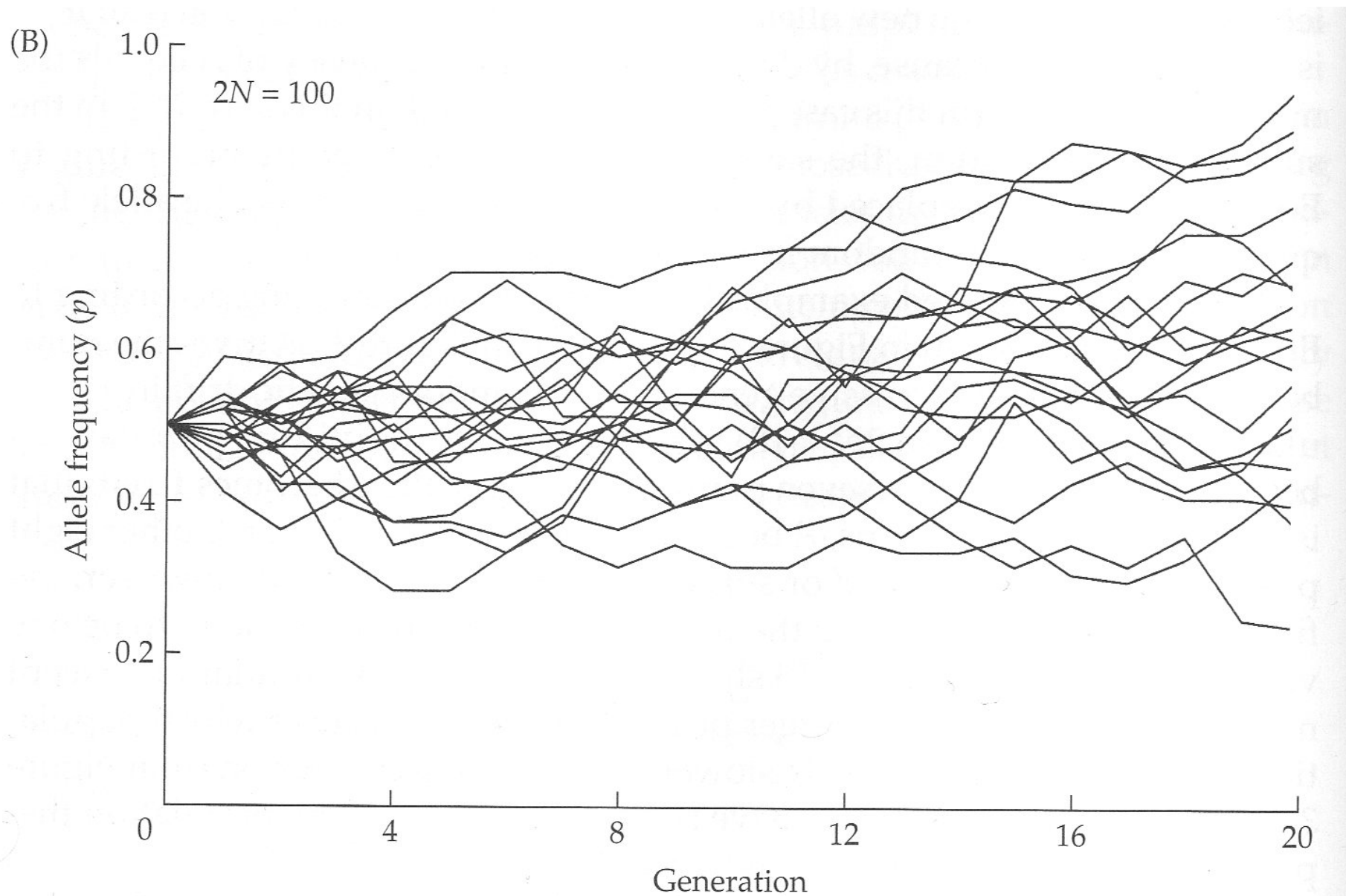
$$Var(\Delta p | p) = Var(k/2N - p | p) = p(1 - p)/2N$$

# Random genetic drift





# Random genetic drift





# Random genetic drift

The endpoint is allele fixation or loss:  $P(F|p) = p$

Mean time to fixation, if fixed:  $\bar{t}_F(p) = -4N \left( \frac{1-p}{p} \right) \ln(1-p)$

Mean time to loss, if lost:  $\bar{t}_L(p) = -4N \left( \frac{p}{1-p} \right) \ln(p)$

Mean persistence time:  $\bar{t}(p) = p\bar{t}_F(p) + (1-p)\bar{t}_L(p) =$   
 $= -4N[(1-p)\ln(1-p) + p \cdot \ln(p)]$

*Exercise:* at which  $p$  persistence time is maximal and what is it?

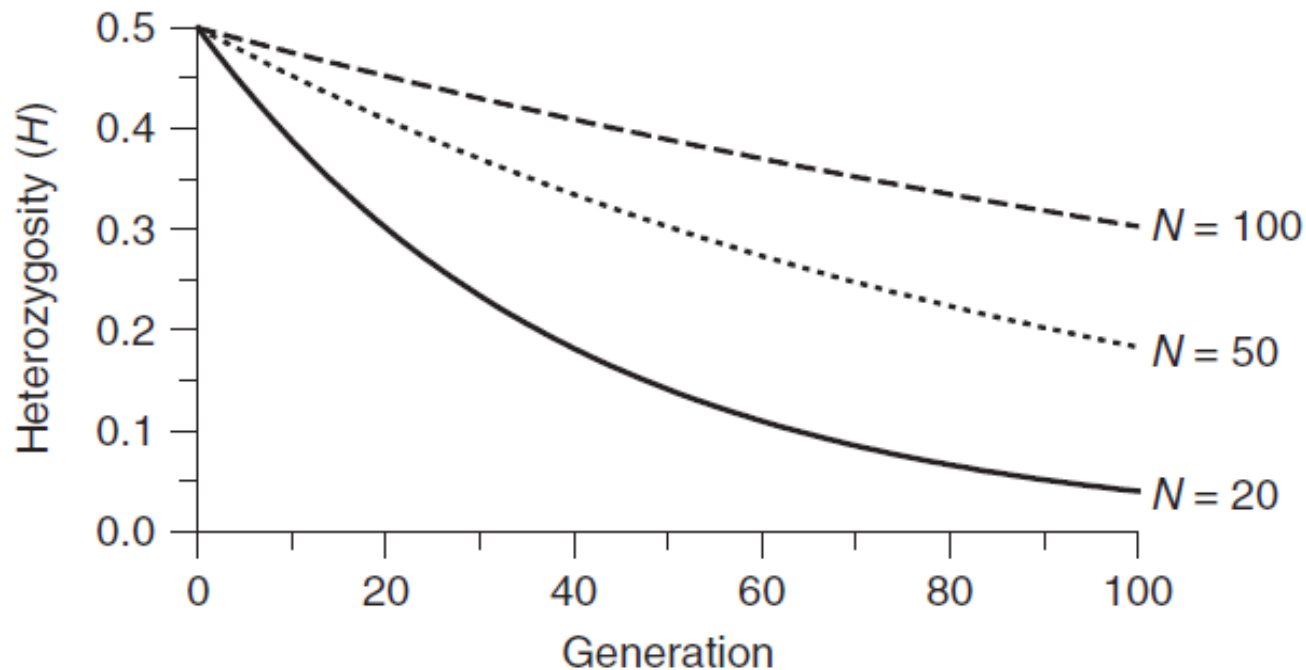
*Exercise:* estimate  $t_F(p)$  when  $p \rightarrow 0$

# Random genetic drift and genetic variation

**Heterozygosity:** probability that an individual is heterozygous at a locus:  $H = 2pq$

Heterozygosity decay due to drift:  $H_t = H_0(1 - 1/2N)^t$

Decay is slow:  $H_t = H_0/2 : t \approx 2N \ln(2)$  for  $N \gg 1$



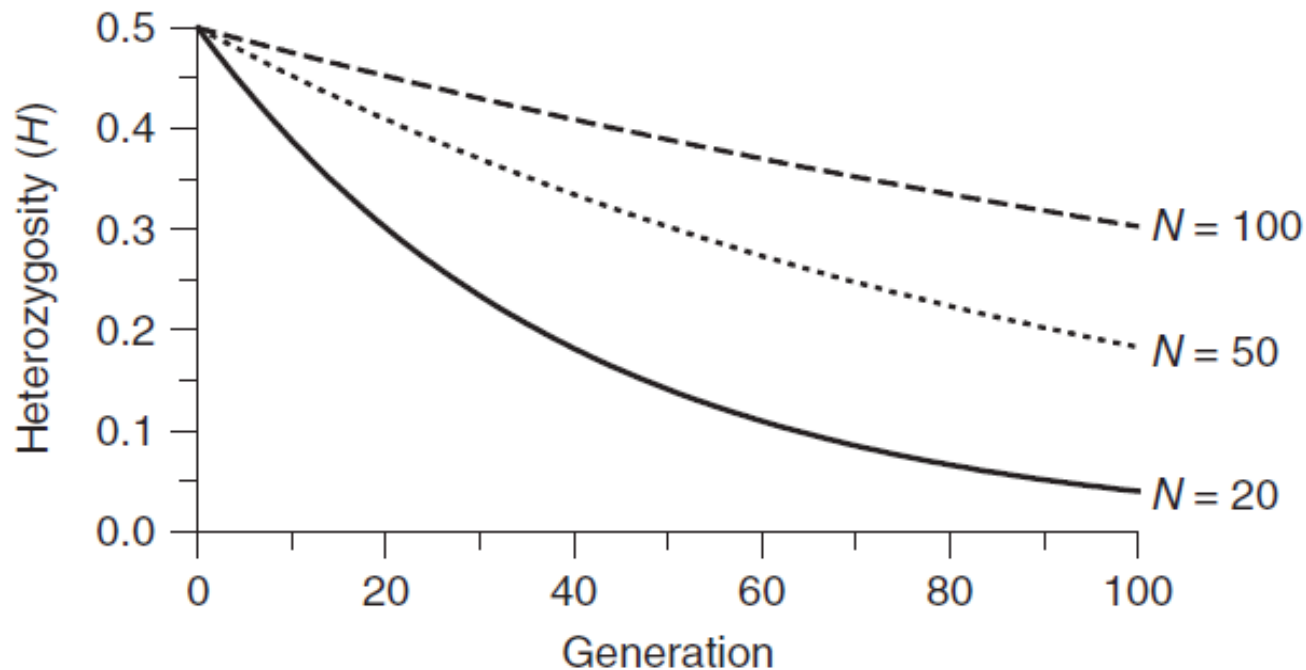
# Random genetic drift and genetic variation

**Heterozygosity:** probability that an individual is heterozygous at a locus:  $H = 2pq$

**Drift strength is  $\approx 1/2N$**

Heterozygosity decay due to drift:  $H_t = H_0(1 - 1/2N)^t$

Decay is slow:  $H_t = H_0/2 : t \approx 2N \ln(2)$  for  $N \gg 1$



# Effective population size

**Effective population size** of an actual population is the number of individuals in a theoretically ideal population having the same magnitude of genetic drift as the actual population (Hartl & Clark, *Principles of population genetics*)

• Fluctuation in population size  $\frac{1}{N_e} = \frac{1}{t} \left( \frac{1}{N_0} + \frac{1}{N_1} + \dots + \frac{1}{N_{t-1}} \right)$

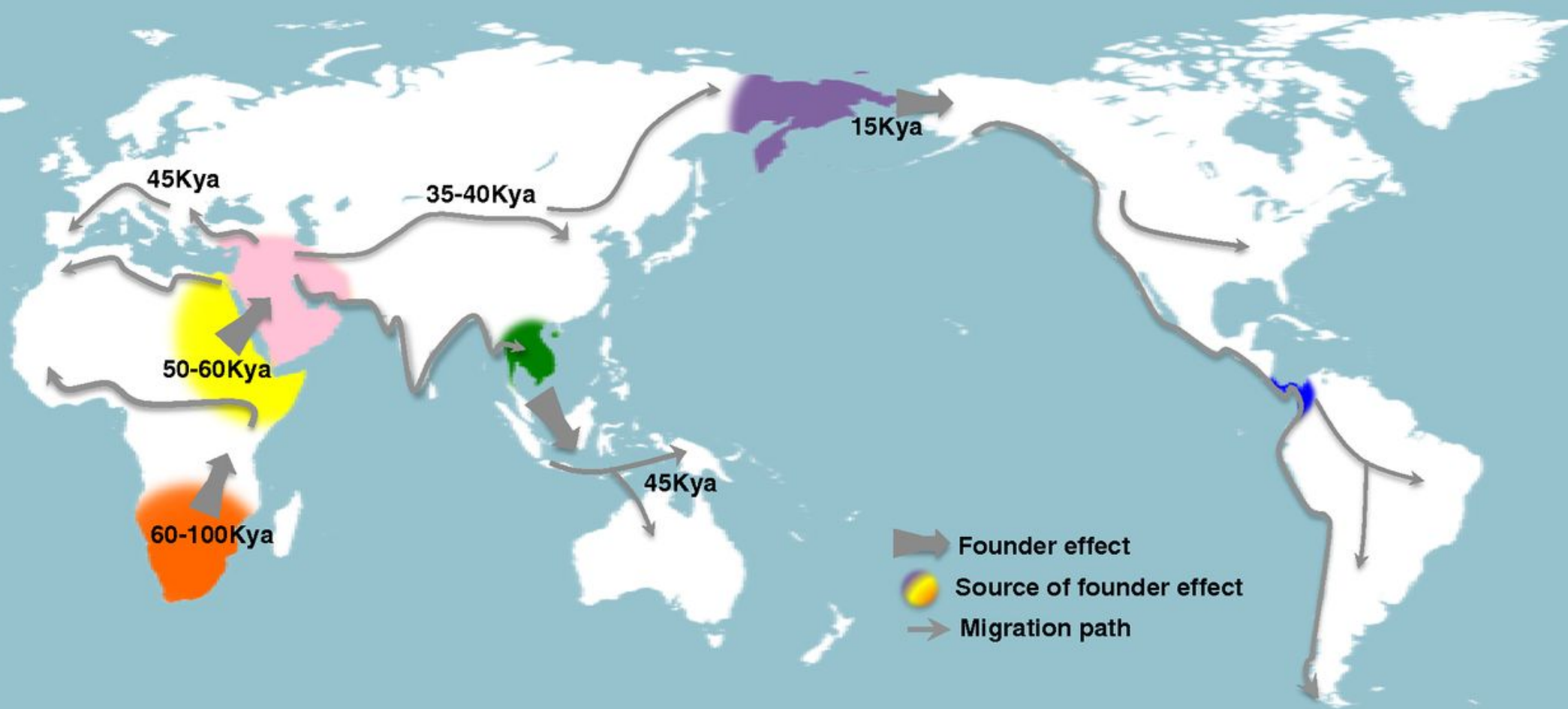
• Unequal sex ratio:  $N_e = \frac{4N_m N_f}{N_m + N_f}$

*Exercise: bottleneck consequences for  $N_e$*

• Variance in offspring number:  
 $\sigma, \xi$  – offspring mean and variance  $N_e = \frac{N - 1}{(\sigma^2/\xi) + (\xi - 1)}$

• Subdivided population:  
 $d$  sub-populations of size  $N$ ;  $m$ , migration  $N_e = Nd \left( 1 + \frac{1}{4Nm} \right)$

# The great human expansion

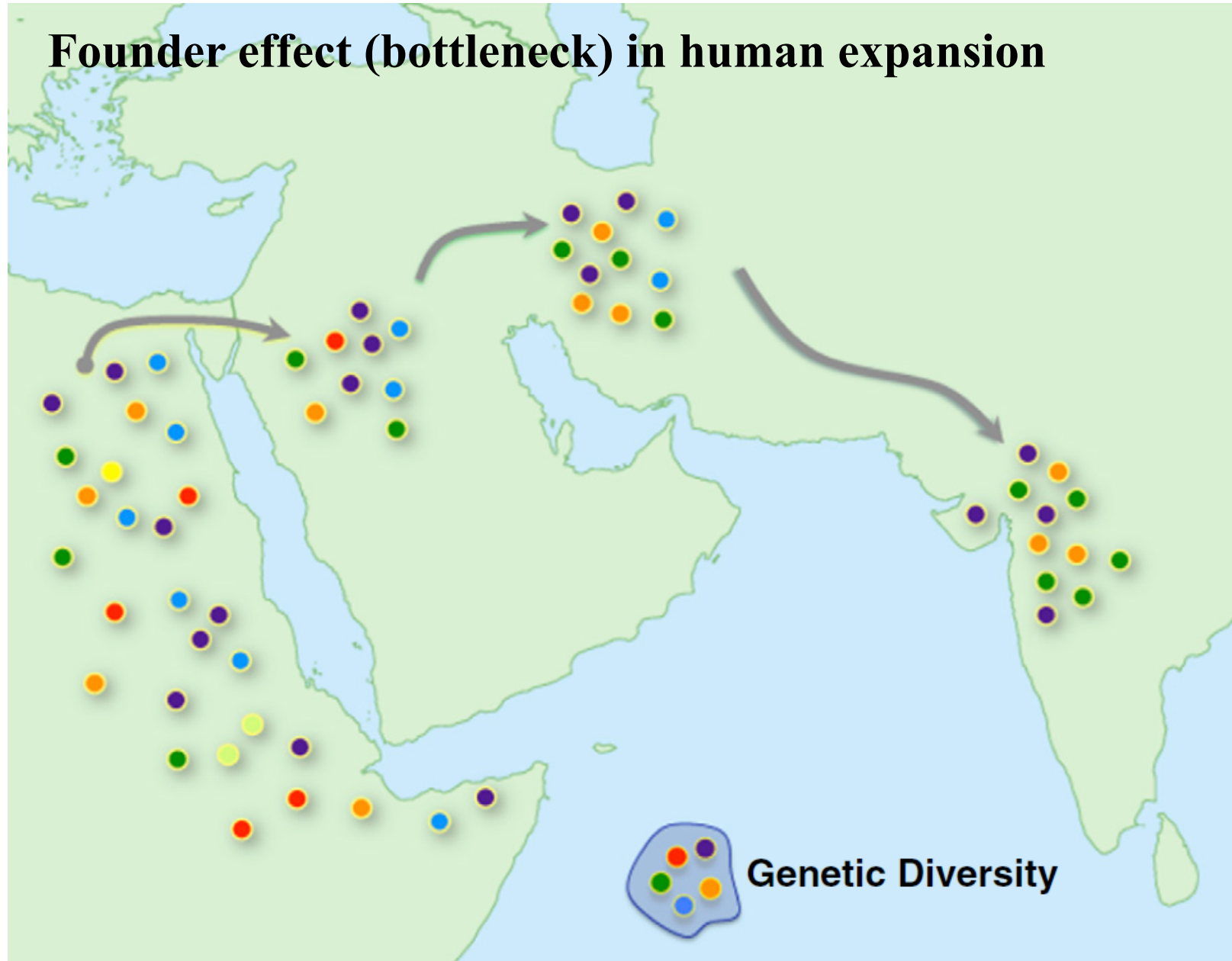


Resequencing studies have estimated the ancestral effective population size at 12,800 to 14,400, with a 5- to 10-fold bottleneck beginning approximately 65,000 to 50,000 y ago (although see ref. 15 for a bottleneck to only 450 individuals).

Henn *et al* (2012) *PNAS*

# The great human expansion

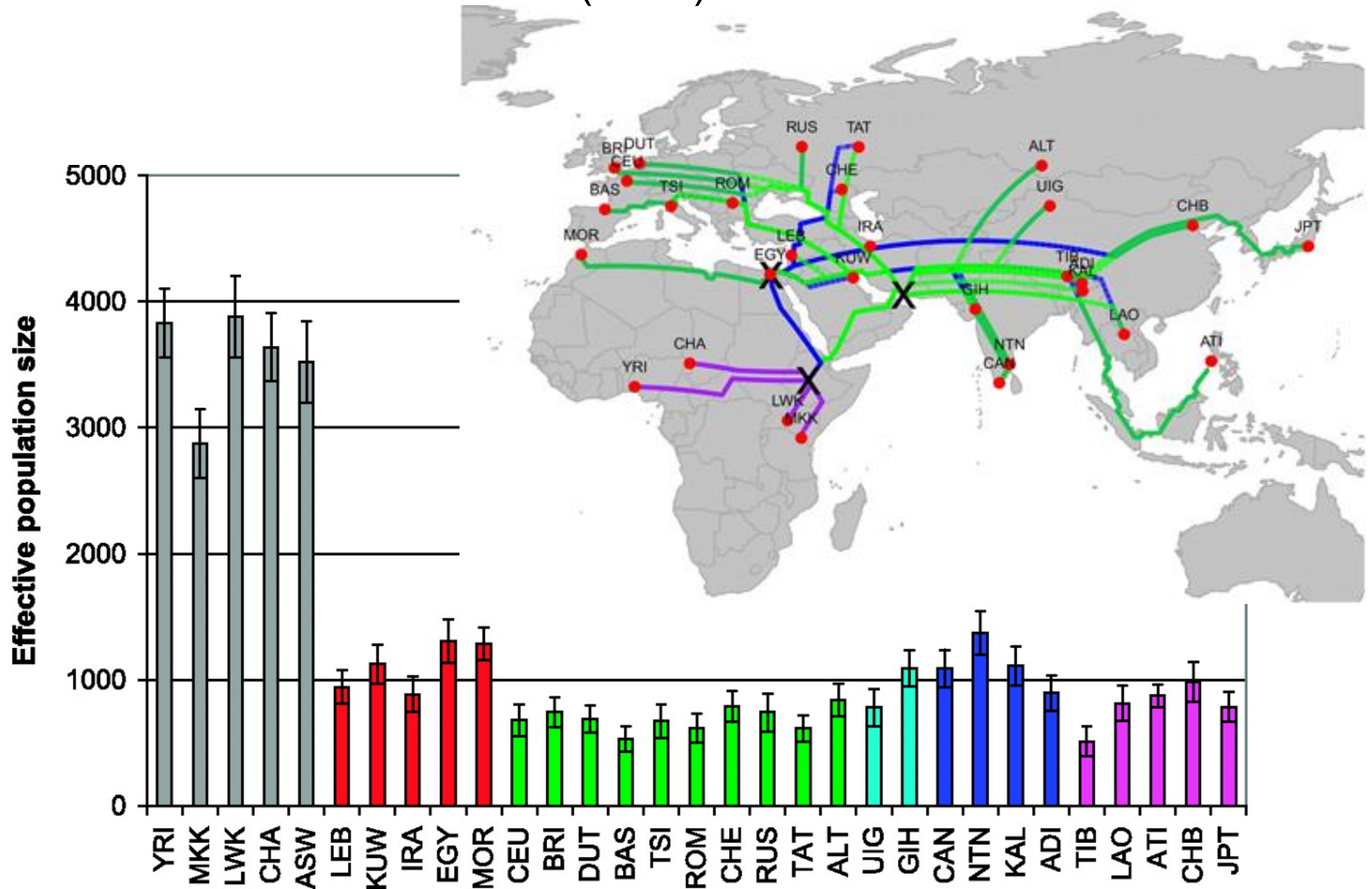
Founder effect (bottleneck) in human expansion





# Recombination Gives a New Insight in the Effective Population Size and the History of the Old World Human Populations

Mele *et al* (2011) *Mol Biol Evol*







# Random genetic drift and mutations

**The neutral theory:** most mutations are selectively neutral with allele frequency determined by random genetic drift (Kimura 1968)

$2N$  gametes  $\Rightarrow 2N\mu$  mutations in each generation, where  $\mu$  = mutations per gamete per generation

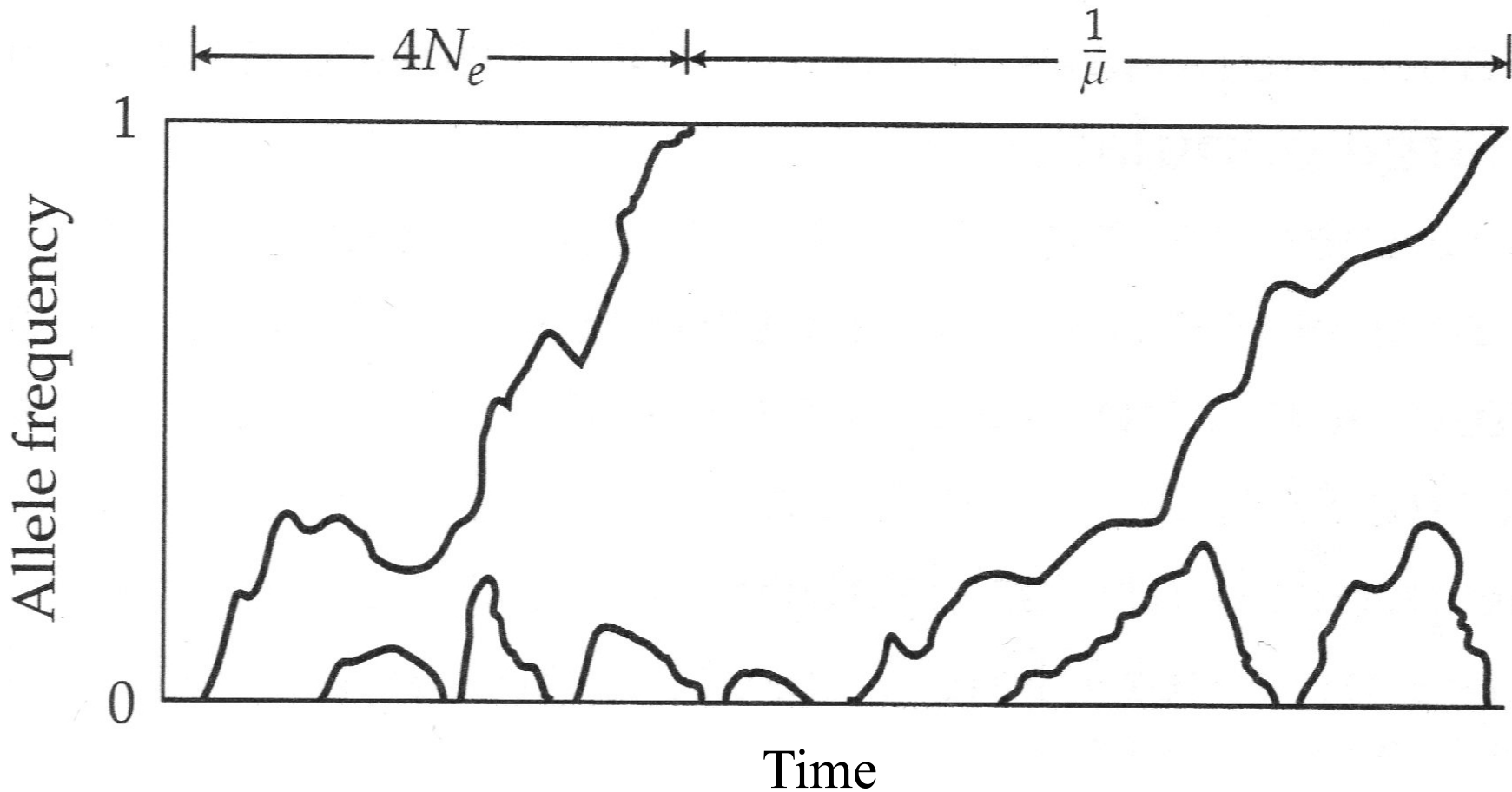
Each mutation  $p_0 = 1/2N \Rightarrow P_{\text{Fix}} = 1/2N$

The steady-state rate at which neutral mutations are fixed in a population:  $k = 2N\mu P_{\text{Fix}} = \mu$

Q: What is the average time between fixation events?

Mean time to fixation, if fixed:  $t_{\text{F}}(p) = 4N_e$  for  $p \approx 0$

# Random genetic drift and mutations

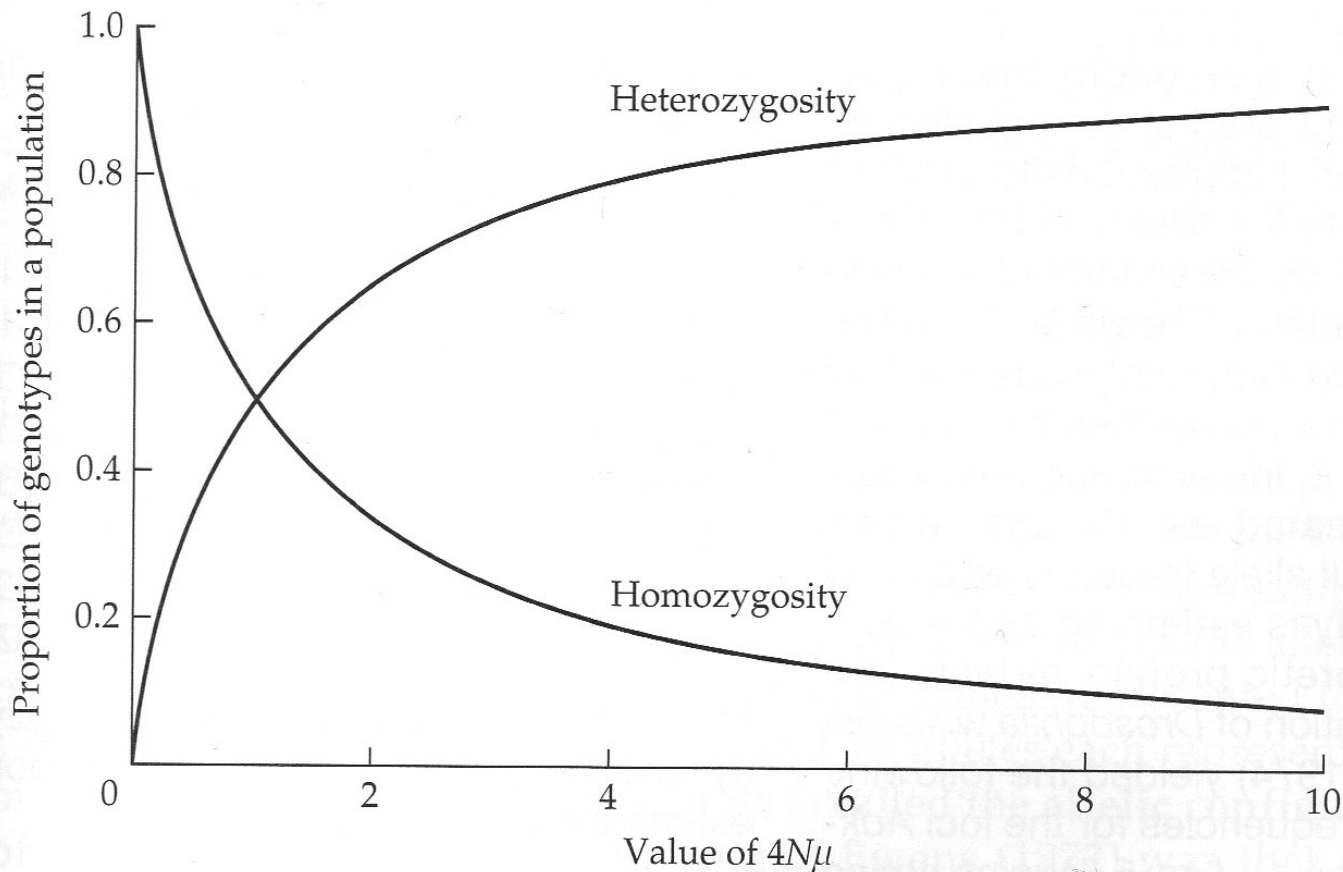


*Exercise:* estimate fixation time for a new neutral allele

# Random genetic drift and mutations

**The infinite-alleles model:** each mutation creates a new allele in the population

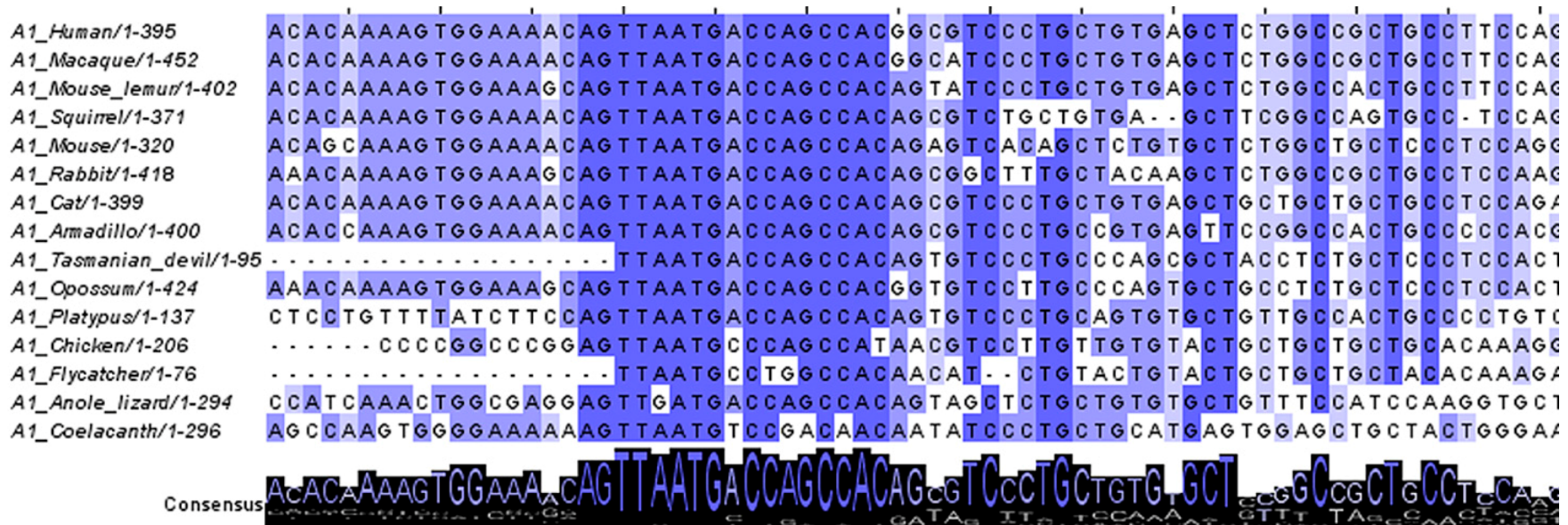
$$\text{Heterozygosity } H = \frac{\theta}{1 + \theta}, \text{ where } \theta = 4N_e\mu$$



# Random genetic drift and mutations

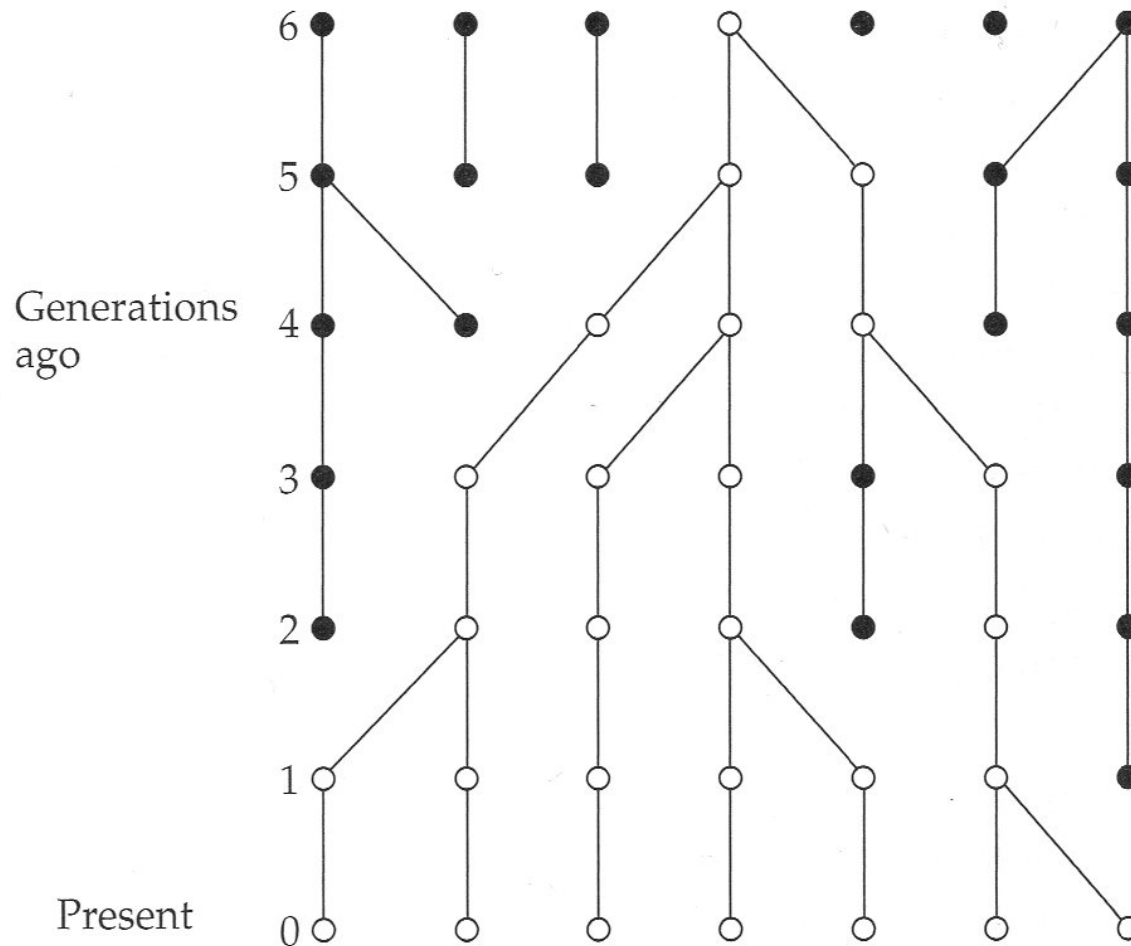
**The neutral (Motoo Kimura) and nearly neutral (Tomoko Ohta) theory of molecular evolution (1960-70):**

- Random genetic drift of [nearly] neutral alleles is the source of polymorphism, not balancing selection.
- Most substitutions (fixations) are due to random drift of neutral mutants, not advantageous mutations
- Missing substitutions are then evolutionary forbidden



# The coalescent theory

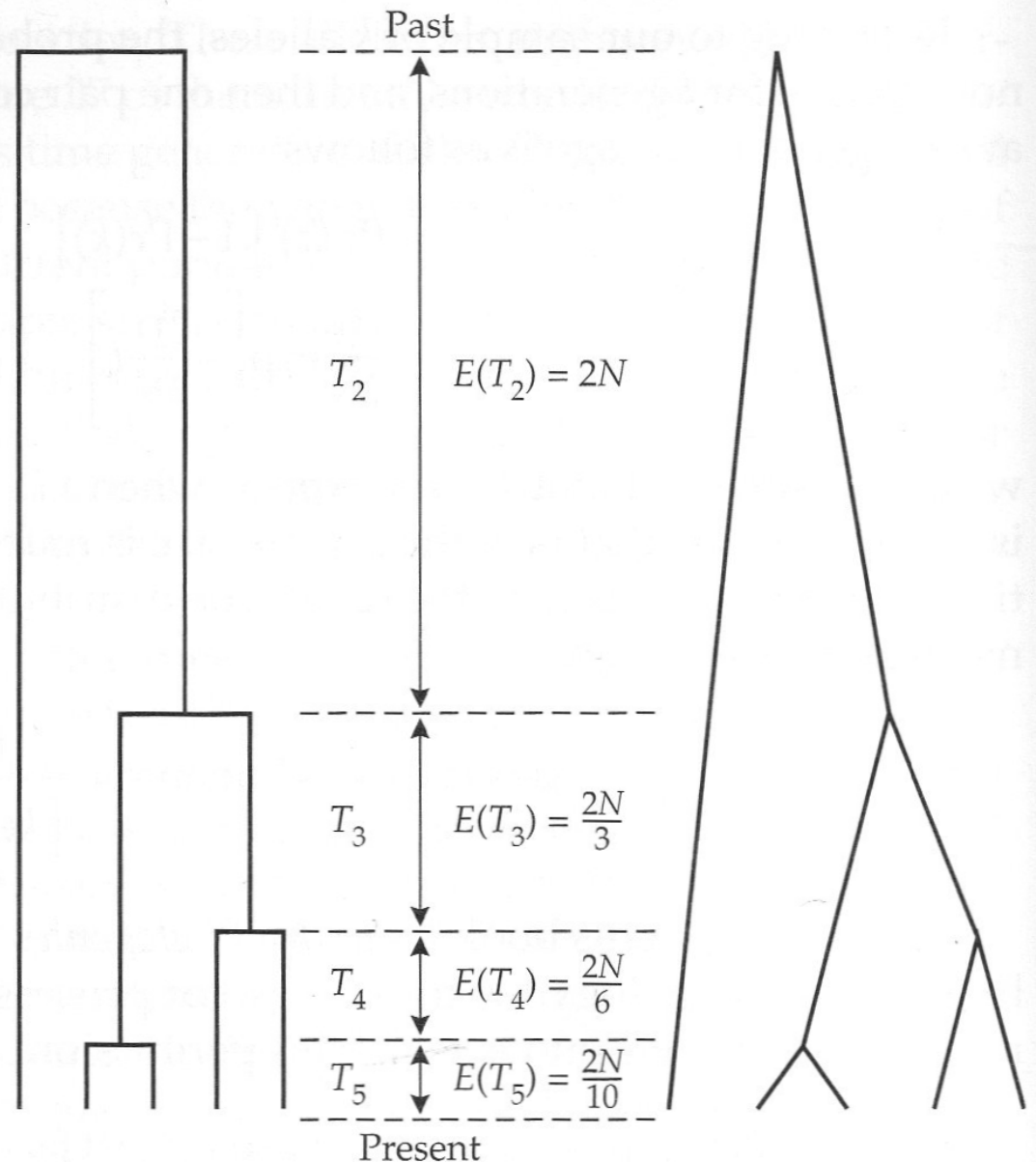
**Coalescent theory** looks back in time and merges sequences originating from a common ancestor





# The coalescent theory

**FIGURE 3.15** Two completely equivalent ways of illustrating the coalescences in a gene tree. On the left, the coalescent events are represented as horizontal lines, on the right they are represented as nodes. In any each generation, if there are  $k$  alleles present, the expected time back to the next coalescence is given by  $4N/[k(k-1)]$ . For example, starting with five alleles, the expected time back to the first coalescence is  $4N/[(5)(4)] = 2N/10$ . Note that the successive times get longer. When there are only two alleles, the time back to the final coalescence is  $2N$  generations.



# The coalescent theory: application

**The infinite-sites model:** each mutation alters a new site in a [very long] nucleotide sequence

|          |          |   |   |          |          |   |   |          |          |   |   |          |          |   |   |
|----------|----------|---|---|----------|----------|---|---|----------|----------|---|---|----------|----------|---|---|
| A        | A        | A | A | T        | T        | T | T | G        | G        | G | G | C        | C        | C | C |
| A        | A        | A | A | T        | T        | T | T | G        | G        | G | G | C        | C        | C | C |
| <b>G</b> | A        | A | A | <b>C</b> | T        | T | T | <b>A</b> | G        | G | G | <b>T</b> | C        | C | C |
| A        | <b>G</b> | A | A | T        | <b>C</b> | T | T | G        | <b>A</b> | G | G | C        | <b>T</b> | C | C |
| 1        | 2        | 3 | 4 | 5        | 6        | 7 | 8 | 9        | 0        | 1 | 2 | 3        | 4        | 5 | 6 |

Sequences:  $n = 4$

Segregating sites:  $S = 8$

Sequence length:  $L = 16$

Average mismatches:  $\Pi = 24/6 = 4$

Nucleotide diversity:  $\pi = H = \Pi/L$

$$E(S) = \theta_s L \sum_{k=1}^{n-1} \frac{1}{k}, \quad \text{where } \theta_s = 4N_e\mu_s$$

Mutation per site per generation:  $\mu_s$

$$E(\Pi) = \theta_s L$$

$$E(\pi) = \theta_s$$

*Exercise:* sample size and variant discovery

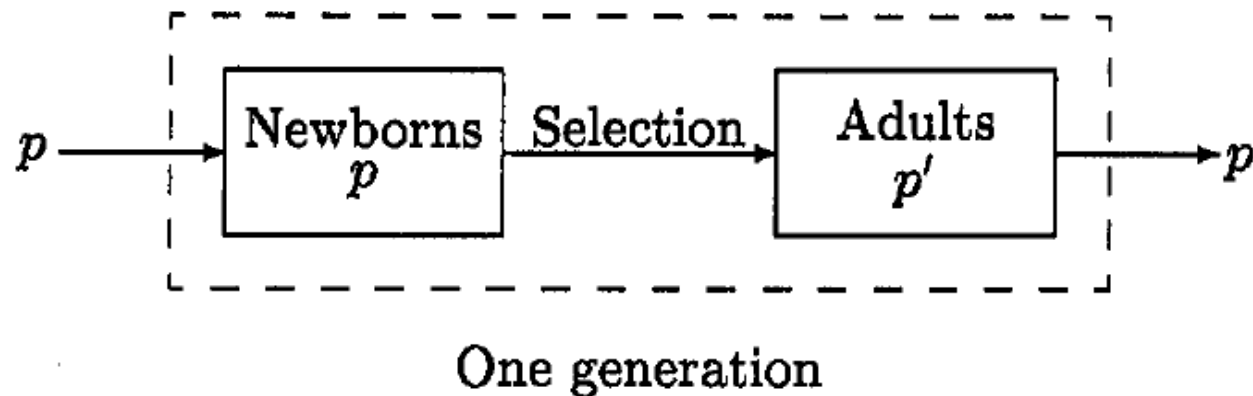




# Natural selection

**Natural selection** is the evolutionary force most responsible for the adaptation to the environment.

Natural selection changes allele frequencies:



**Fitness**  $\approx$  viability [+fertility+developmental time+mating, ...]

Alleles that reduce fitness are **deleterious**.

# Natural selection

**TABLE 5.2** Diploid Selection for Survivorship (Viability)

|                              | Genotype                    |   |                             | Total   |
|------------------------------|-----------------------------|---|-----------------------------|---|
| Generation $t - 1$           | $AA$                        | $Aa$  | $aa$                        |   |
| Frequency before selection   | $p^2$                       | $2pq$                                       | $q^2$                       | $1 = p^2 + 2pq + q^2$                         |
| Relative fitness (viability) | $w_{11}$                    | $w_{12}$                                    | $w_{22}$                    |   |
| After selection              | $p^2w_{11}$                 | $2pqw_{12}$                                 | $q^2w_{22}$                 | $\bar{w} = p^2w_{11} + 2pqw_{12} + q^2w_{22}$ |
| Normalized                   | $\frac{p^2w_{11}}{\bar{w}}$ | $\frac{2pqw_{12}}{\bar{w}}$                 | $\frac{q^2w_{22}}{\bar{w}}$ |   |
| Generation $t$               |                             | $p' = \frac{p^2w_{11} + pqw_{12}}{\bar{w}}$ |                             |   |
|                              |                             | $q' = \frac{pqw_{12} + q^2w_{22}}{\bar{w}}$ |                             |   |

$$\Delta p = \frac{pq[p(w_{11} - w_{12}) + q(w_{12} - w_{22})]}{\bar{w}}$$

# Natural selection

---

|                     |          |                 |                 |
|---------------------|----------|-----------------|-----------------|
| Genotype            | $A_1A_1$ | $A_1A_2$        | $A_2A_2$        |
| Viability (fitness) | $w_{11}$ | $w_{12}$        | $w_{22}$        |
| Relative fitness    | 1        | $w_{12}/w_{11}$ | $w_{22}/w_{11}$ |

|                  |   |        |       |
|------------------|---|--------|-------|
| Relative fitness | 1 | $1-hs$ | $1-s$ |
|------------------|---|--------|-------|

where  $0 \leq s \leq 1$  is the **selection coefficient**,

$h$  is the **heterozygous effect** and measures **dominance**

---

$h = 0$        $A_1$  dominant,  $A_2$  recessive

$h = 1$        $A_1$  recessive,  $A_2$  dominant

$0 < h < 1$       incomplete dominance ( $h = 1/2$  additive)

$h < 0$       overdominance

$h > 1$       underdominance

---

# Natural selection

$$\Delta p = \frac{pq[p(w_{11} - w_{12}) + q(w_{12} - w_{22})]}{\bar{w}}$$

Switch to relative fitness:  $w_{12}/w_{11} = 1 - hs$ ,  $w_{22}/w_{11} = 1 - s$

$$\Delta p = \frac{pq s [ph + q(1 - h)]}{\tilde{w}}$$

$$\tilde{w} = 1 - 2pqhs - q^2s$$

*Exercise: derive*

# Natural selection

## 1. Directional (positive, negative, purifying) selection

Recessive allele:  $w_{11}=1$ ,  $w_{12}=1$ ,  $w_{22}=1-s$

Dominant allele:  $w_{11}=1$ ,  $w_{12}=1-s$ ,  $w_{22}=1-s$

Incomplete dominance:  $w_{11}=1$ ,  $w_{12}=1-hs$ ,  $w_{22}=1-s$ ,  $0 < h < 1$

## 2. Balancing selection

Overdominance:  $w_{11}=1$ ,  $w_{12}=1-hs$ ,  $w_{22}=1-s$ ,  $h < 0$

## 3. Disruptive selection

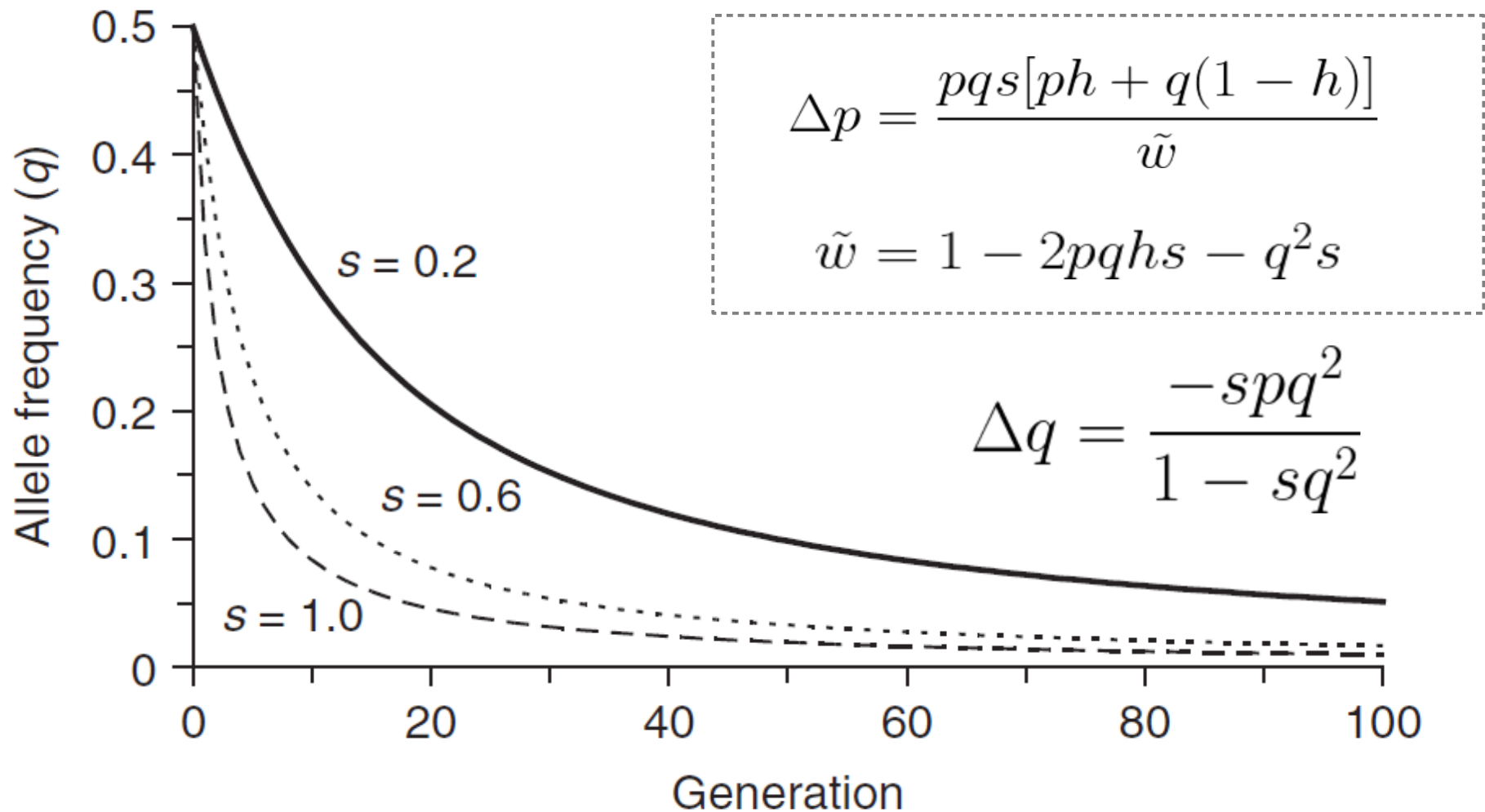
Underdominance:  $w_{11}=1$ ,  $w_{12}=1-hs$ ,  $w_{22}=1-s$ ,  $h > 1$

*Exercise: valid range for  $h$  ?*

# Natural selection

Directional selection against a recessive allele:

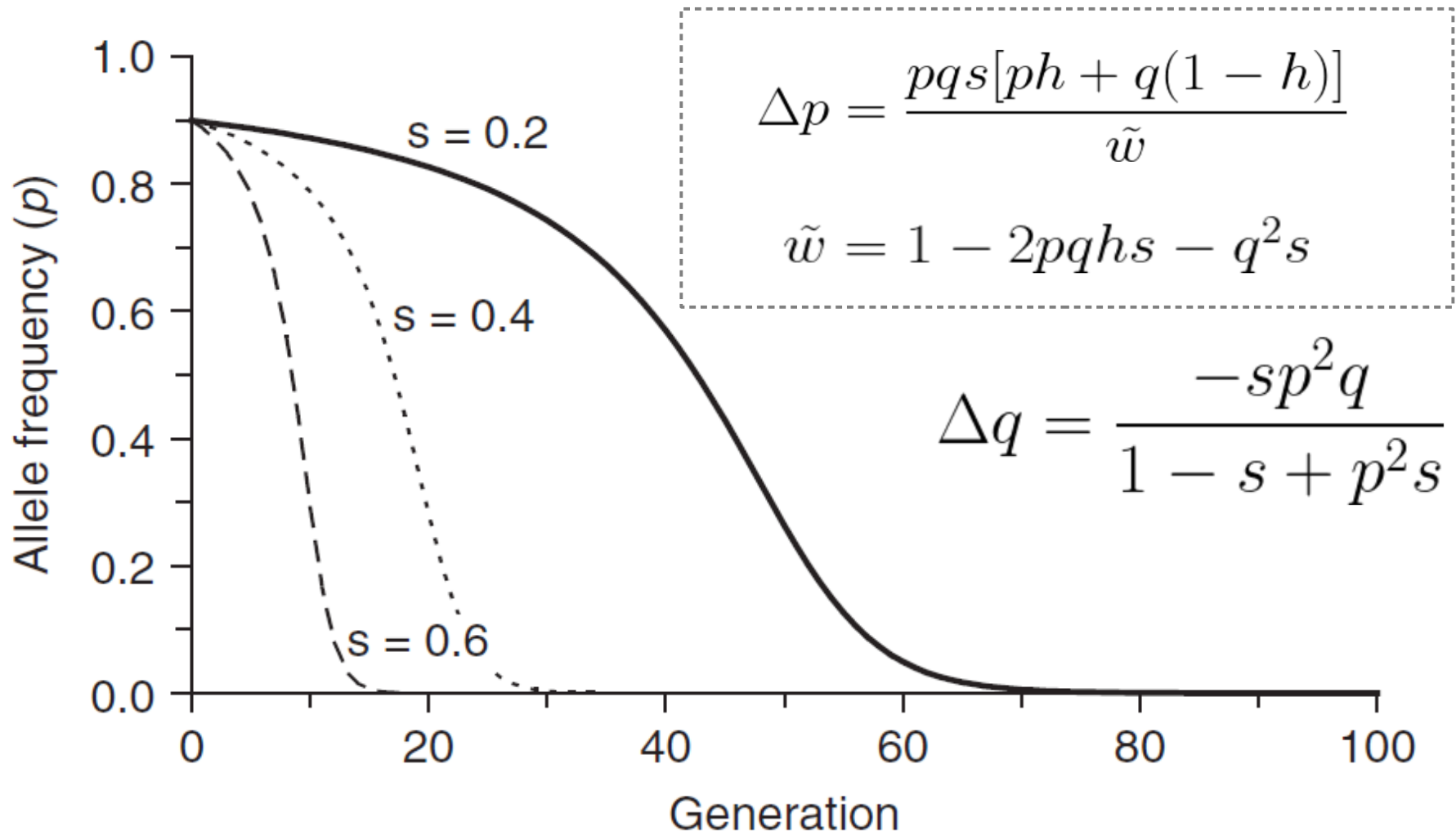
$$w_{11} = w_{12} = 1, \quad w_{22} = 1 - s$$



# Natural selection

Directional selection against a dominant allele:

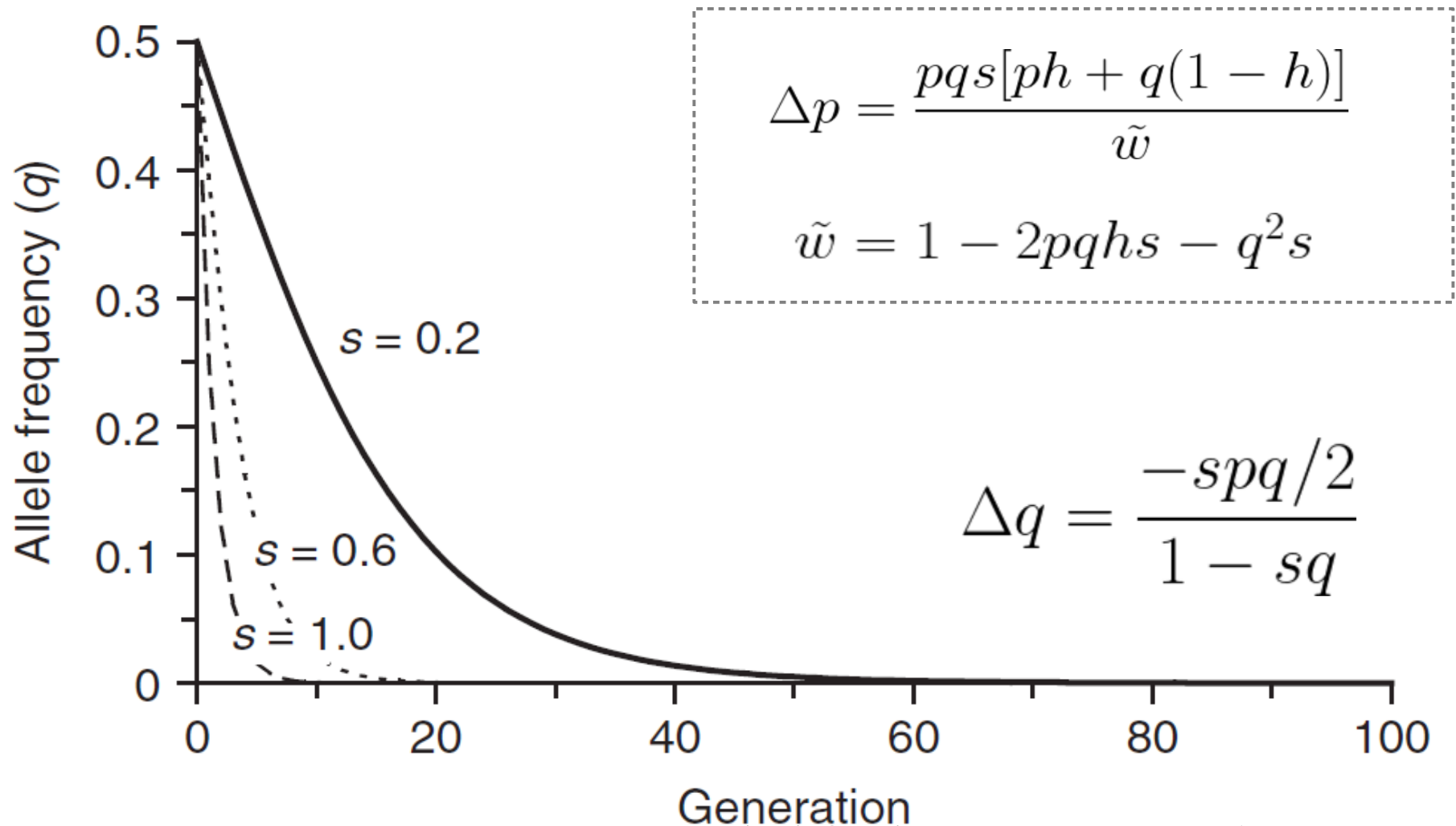
$$w_{11} = 1, w_{12} = w_{22} = 1 - s$$



# Natural selection

Directional selection against a codominant additive allele:

$w_{11} = 1$ ,  $w_{12} = 1 - s/2$ ,  $w_{22} = 1 - s$  // incomplete dominance

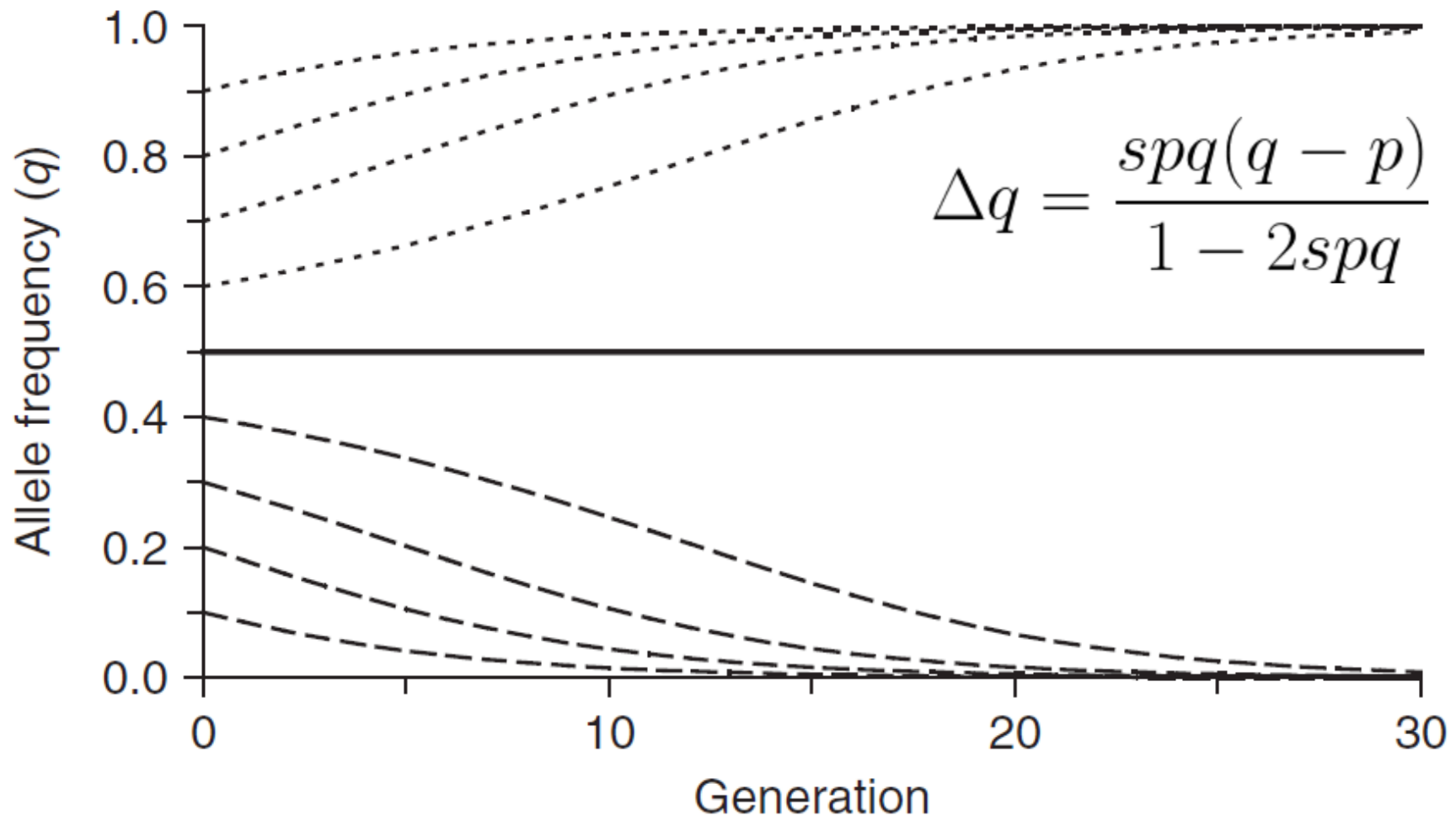




# Natural selection

Disruptive selection against a heterozygote:

$w_{11} = 1, w_{12} = 1 - s, w_{22} = 1$  // underdominance



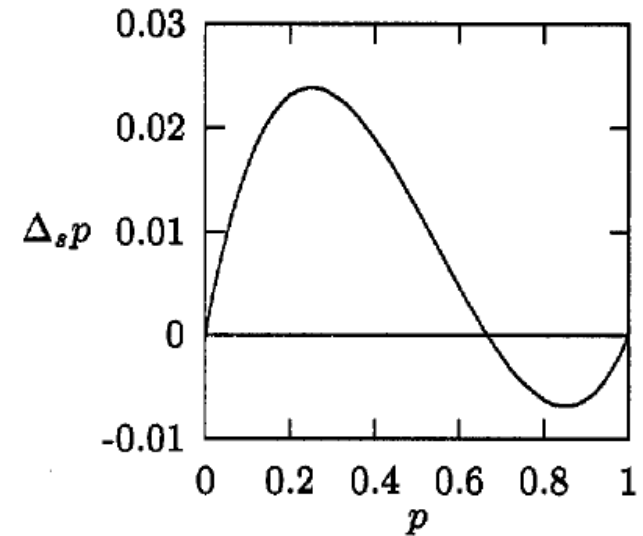
# Natural selection

Balancing selection for a heterozygote:

$w_{11}=1$ ,  $w_{12} = 1 - hs$ ,  $w_{22} = 1 - s$ ,  $h < 0$  // overdominance

$$\Delta p = \frac{pqs[ph + q(1 - h)]}{\tilde{w}}$$

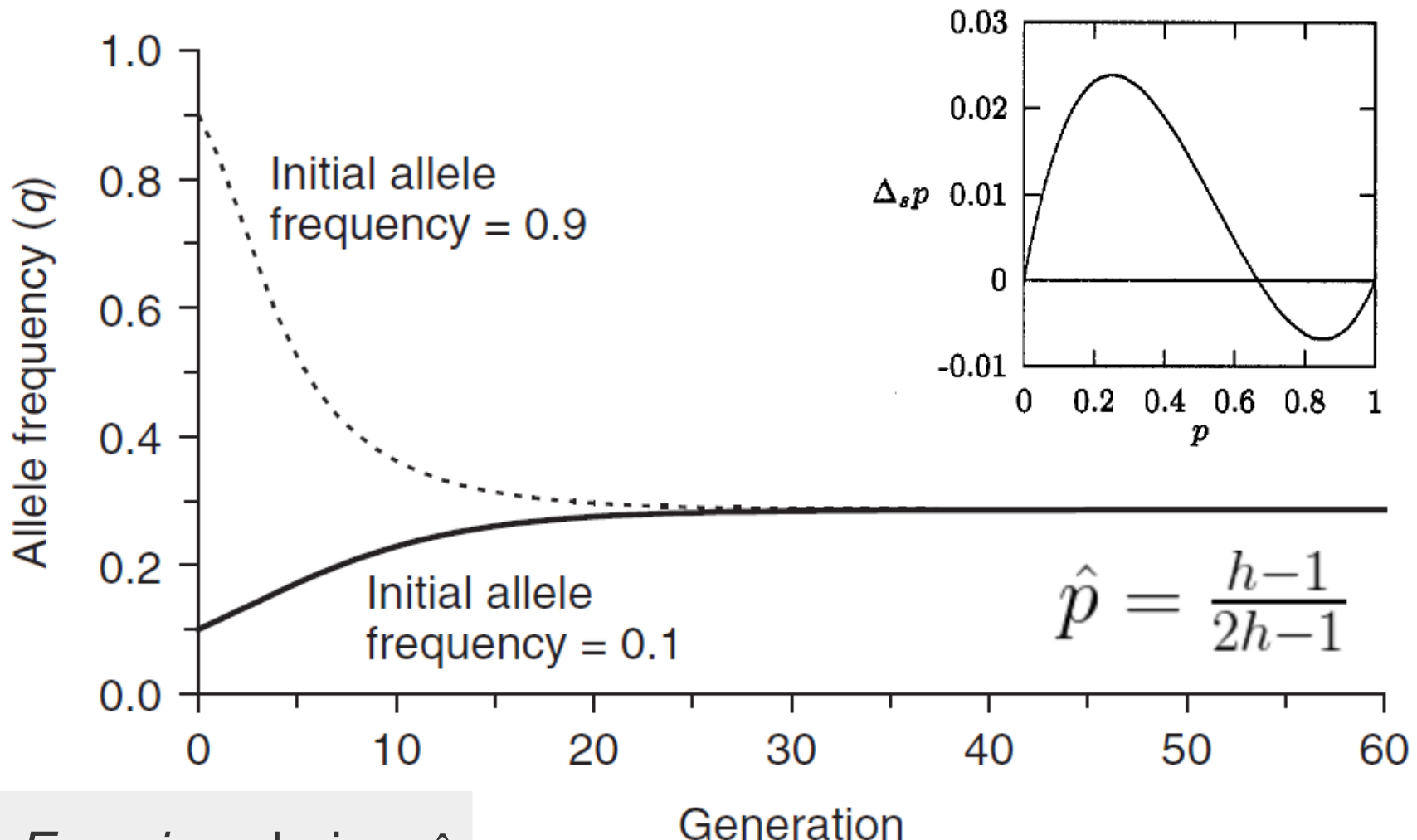
$$\tilde{w} = 1 - 2pqhs - q^2s$$



# Natural selection

Balancing selection for a heterozygote:

$$w_{11}=1, w_{12} = 1 - hs, w_{22}= 1 - s, h < 0 \quad // \text{ overdominance}$$



Exercise: derive  $\hat{p}$

# Balancing selection: the case of CF

## BOX 3.7 SELECTION IN FAVOR OF HETEROZYGOTES FOR CYSTIC FIBROSIS

For CF, the disease frequency in Denmark is about one in 2000 births.

|              |            |       |                |
|--------------|------------|-------|----------------|
| Phenotypes:  | Unaffected |       | Affected       |
| Genotypes:   | AA         | Aa    | aa             |
| Frequencies: | $p^2$      | $2pq$ | $q^2 = 1/2000$ |

$q^2$  is  $5 \times 10^{-4}$ ; therefore  $q = 0.022$  and  $p = 1 - q = 0.978$ .

$p/q = 0.978/0.022 = 43.72 = s_2/s_1$ .

If  $s_2 = 1$  (affected homozygotes never reproduce),  $s_1 = 0.023$ .

The present CF gene frequency will be maintained, even without fresh mutations, if Aa heterozygotes have on average 2.3% more surviving children than AA homozygotes.

*Exercise:* express heterozygous advantage  $h$  as a function of  $p^{\wedge}$ , verify estimate above



# Balancing selection: the case of $\beta$ -hemoglobin

The most thoroughly studied example of overdominance is the sickle-cell hemoglobin polymorphism found in many human populations in Africa. Hemoglobin, the oxygen-carrying red protein found in red blood cells, is a tetramer composed of two alpha chains and two beta chains. In native West and Central African populations, the  $S$  allele of beta hemoglobin reaches a frequency as high as 0.3 in some areas. The more common  $A$  allele is found at very high frequency in most other areas of the world. The two alleles differ only in that the  $S$  allele has a glutamic acid at its sixth amino position while the  $A$  allele has a valine. The glutamic acid causes the hemoglobin to form crystal aggregates under low partial pressures of oxygen, as occur, for example, in the capillaries. As a result,  $SS$  homozygotes suffer from sickle-cell anemia, a disease that is often fatal.

The  $S$  allele could not have reached a frequency of 0.3 unless  $AS$  heterozygotes are more fit than  $AA$  homozygotes. This is precisely the case in regions where malaria is endemic, for there the heterozygotes are somewhat resistant to severe forms of malaria. The resistance is due to the sickling phenomena, which makes red blood cells less suitable for *Plasmodium falciparum*. In an old study from 1961, it was shown that the viability of  $AS$  relative to  $AA$  is 1.176 in regions with malaria. Assuming that the fitness of  $SS$  is zero ( $s = 1$ ),  $h = -0.176$ . Plugging this into Equation 3.4 gives  $\hat{p} = 0.87$  or  $\hat{q} = 0.13$  for the  $S$  allele, which is nestled right in the middle of allele frequencies in regions with endemic malaria.



# Natural selection

$$\Delta p = \frac{pq s [ph + q(1 - h)]}{\tilde{w}} \quad \tilde{w} = 1 - 2pqhs - q^2 s$$

Sewall Wright:

$$\Delta p = \frac{pq}{2\tilde{w}(p)} \frac{d\tilde{w}(p)}{dp}$$

“Natural selection always increases the mean fitness and does so at a rate that is proportional to the genetic variation”





# Mutation-selection balance

- Many new alleles are deleterious and incompletely dominant.
- They enter the population by mutation and are removed by negative selection.

$$A_1(p \approx 1) \xrightarrow{\mu} A_2(q \approx 0)$$

- Balance: the rate of introduction of mutations equals rate of loss due to selection

$$\Delta_{mut} p = -\mu p \approx -\mu$$

$$\Delta_{sel} p = \frac{pqs[ph + q(1 - h)]}{1 - 2pqhs - q^2s} \approx qhs$$

$$\Delta_{mut} p + \Delta_{sel} p = 0$$

$$\hat{q} \approx \frac{\mu}{hs}$$



# Mutation-selection balance

- Many new alleles are deleterious and incompletely dominant.
- They enter the population by mutation and are removed by negative selection.

$$A_1(p \approx 1) \xrightarrow{\mu} A_2(q \approx 0)$$

- Balance: the rate of introduction of mutations equals rate of loss due to selection

$$\Delta_{mut} p = -\mu p \approx -\mu$$

$$\Delta_{sel} p = \frac{pqs[ph + q(1 - h)]}{1 - 2pqhs - q^2s} \approx qhs$$

$$\Delta_{mut} p + \Delta_{sel} p = 0$$

$$\hat{q} \approx \frac{\mu}{hs}$$

**Large effect →  
Low frequency**

*Exercise:* derive  $\hat{q}$   
for a recessive allele

# Random drift and advantageous allele

Selection in finite population is very weak for *de novo* alleles:  
New allele:  $\Delta p \approx (1+s)p - p = sp = s/2N \ll 1/2N$  (drift),  
unless  $s \approx 1$

$$P_F(p) = \frac{1 - e^{-2Nsp}}{1 - e^{-2Ns}}, \text{ if } h = 1/2$$

$$P_F(1/2N) = \frac{1 - e^{-s}}{1 - e^{-2Ns}} \quad \mathbf{P_F \approx s \text{ if } s \approx 0 \text{ and } 2Ns \gg 1}$$

# Random drift and advantageous allele

Selection in finite population is very weak for *de novo* alleles:  
New allele:  $\Delta p \approx (1+s)p - p = sp = s/2N \ll 1/2N$  (drift),  
unless  $s \approx 1$

$$P_F(p) = \frac{1 - e^{-2Nsp}}{1 - e^{-2Ns}}, \text{ if } h = 1/2$$

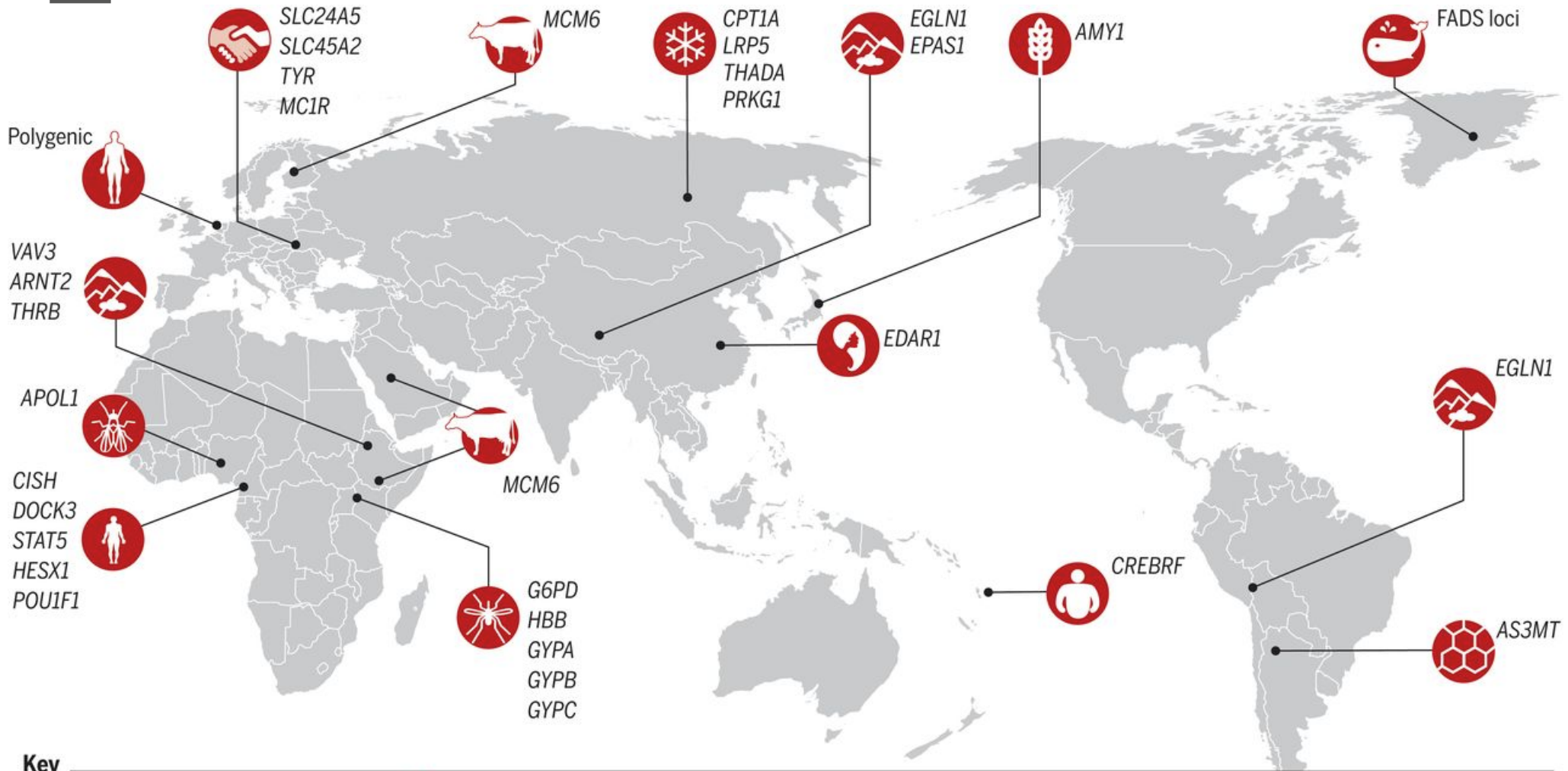
$$P_F(1/2N) = \frac{1 - e^{-s}}{1 - e^{-2Ns}} \quad P_F \approx s \text{ if } s \approx 0 \text{ and } 2Ns \gg 1$$

- Most advantageous alleles are lost.
- Adaptive evolution is random













Exercise:  $P_F$  for  $s, 2Ns \approx 0$



# Examples of human local adaptations



## Key

- |  |  |   |  |  |  |
|--|--|---|--|--|--|
| <br>Lactase persistence | <br>Height        | <br>Arctic environment     | <br>High-fat diet | <br>Thick hair                      | <br>Starchy food  |
| <br>Skin pigmentation   | <br>High altitude | <br>Trypanosome resistance | <br>Malaria       | <br>Toxic arsenic-rich environments | <br>Increased BMI |

# Random drift and deleterious allele

Can a deleterious allele fix in a finite population?

$$P_F(q) = 1 - P_F(1 - q) = \frac{e^{2Nsq} - 1}{e^{2Ns} - 1}$$

$$P_F(1/2N) \approx \frac{s}{e^{2Ns} - 1} \quad \mathbf{P_F \approx 0 \text{ if } 2Ns \gg 1}$$

# Random drift and deleterious allele

Can a deleterious allele fix in a finite population?

$$P_F(q) = 1 - P_F(1 - q) = \frac{e^{2Nsq} - 1}{e^{2Ns} - 1}$$

$$P_F(1/2N) \approx \frac{s}{e^{2Ns} - 1} \quad \mathbf{P_F \approx 0 \text{ if } 2Ns \gg 1}$$

Fixation rate for deleterious alleles:

$$k = 2N\mu P_F(1/2N) = \frac{2N\mu s}{e^{2Ns} - 1}$$

Exercise:  $P_F$  for  $s \rightarrow 0$

Exercise:  $k$  for  $s \rightarrow 0$  ?

# Mildly deleterious vs neutral mutations

Mutations can be placed in three main categories:

- those that are selected (either positively or negatively);
- those that are neutral (i.e. have no effect on fitness) and
- those that have low selection coefficients, and thus behave as neutral in small populations (where the effects of drift dominate) or are selected in large populations, where the deterministic effects of selection prevail

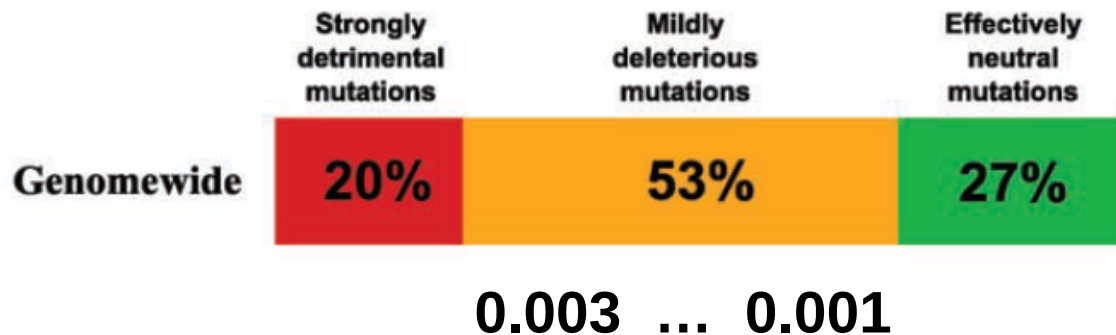
Meyer, Diogo; and, Harris, Eugene E (March 2008) Selection Operating on Protein-coding Genes in the Human Genome. In: Encyclopedia of Life Sciences (ELS). John Wiley & Sons, Ltd: Chichester.  
DOI: 10.1002/9780470015902.a0020791

# Mildly deleterious vs neutral mutations

## Most Rare Missense Alleles Are Deleterious in Humans: Implications for Complex Disease and Association Studies

Gregory V. Kryukov, Len A. Pennacchio, and Shamil R. Sunyaev

The American Journal of Human Genetics Volume 80 April 2007



We combined analysis of mutations causing human Mendelian diseases, of human-chimpanzee divergence, and of systematic data on human genetic variation and ... estimated that  $>50\%$  of *de novo* missense mutations in an average human gene and  $70\%$  of missense SNPs detected only once among 1,500 chromosomes are mildly deleterious. Such **mildly deleterious mutations are associated with selection coefficients within a surprisingly narrow range of 0.001–0.003**

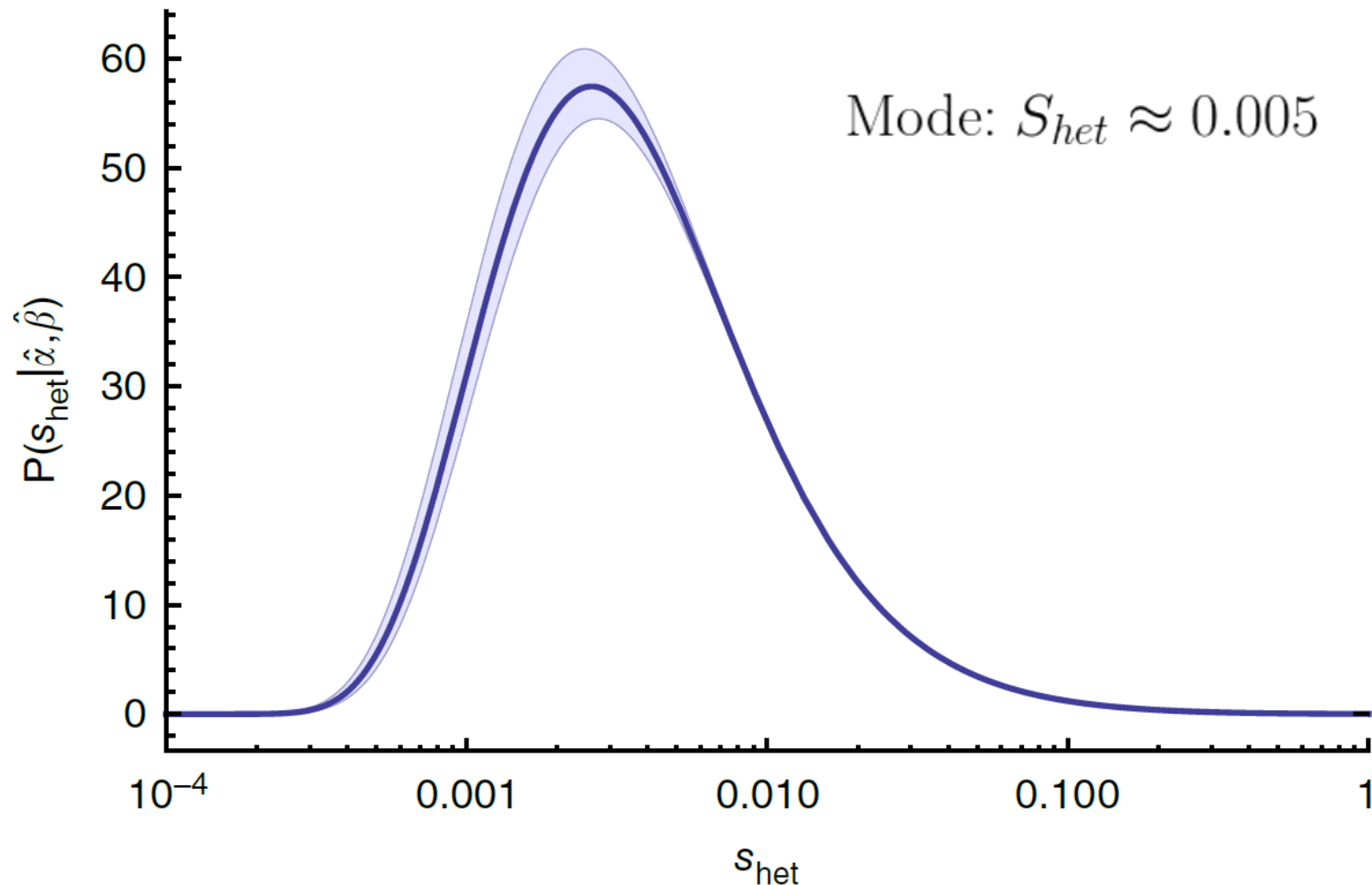
Kryukov (2007) *Am J Hum Genet*



# Estimating the selective effects of heterozygous protein-truncating variants from human exome data

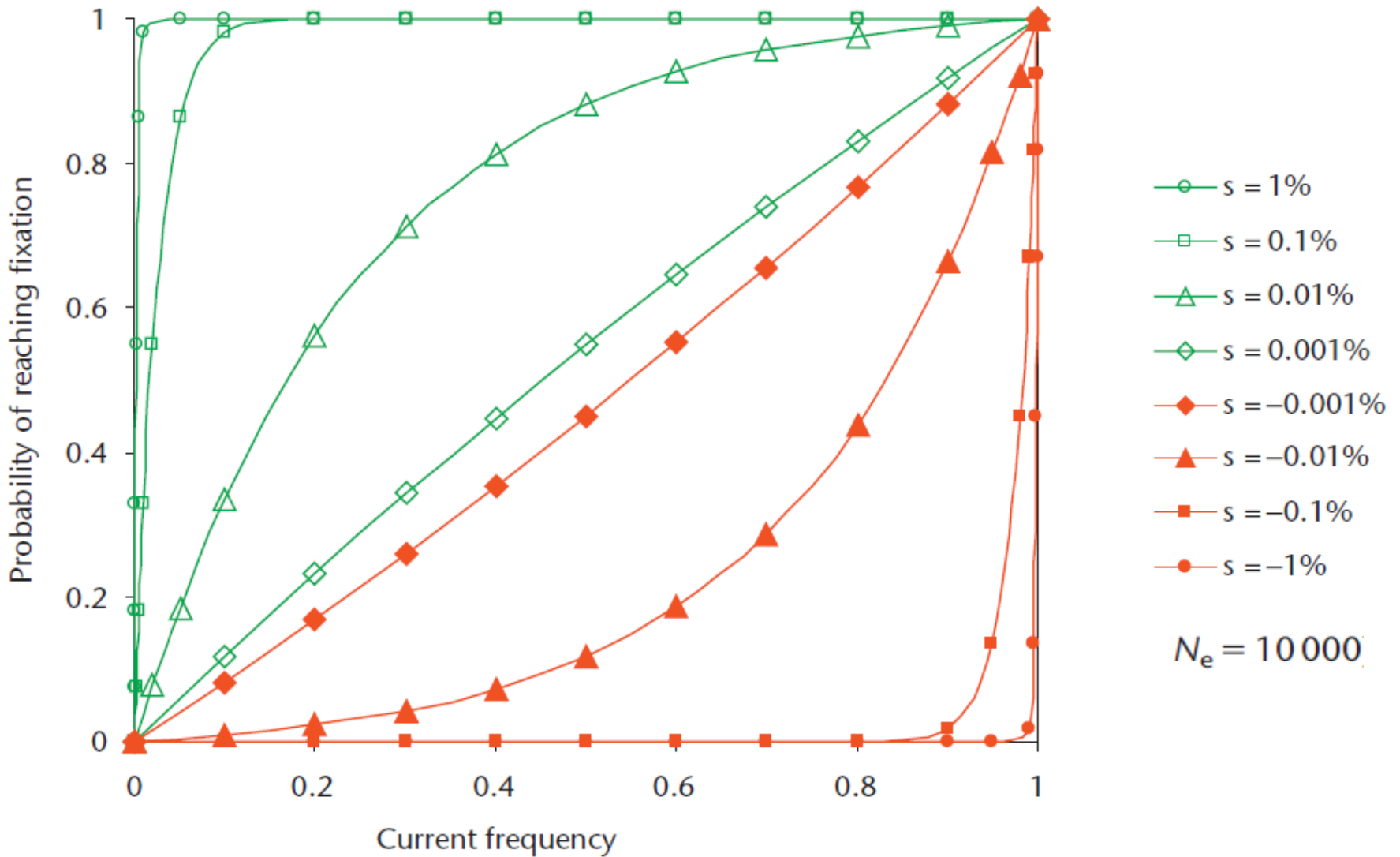
Christopher A Cassa<sup>1,2,9</sup>, Donate Weghorn<sup>1,9</sup>, Daniel J Balick<sup>1,9</sup>, Daniel M Jordan<sup>3,9</sup>, David Nusinow<sup>1</sup>, Kaitlin E Samocha<sup>4,5</sup>, Anne O'Donnell-Luria<sup>4,6</sup>, Daniel G MacArthur<sup>2,4</sup>, Mark J Daly<sup>2,4</sup>, David R Beier<sup>7,8</sup> & Shamil R Sunyaev<sup>1,2</sup>

VOLUME 49 | NUMBER 5 | MAY 2017 **NATURE GENETICS**



Cassa (2017) *Nat Genet*

# Fixation probabilities for all alleles



Thomas, Paul D (July 2008) Single Nucleotide Polymorphisms in Human Disease and Evolution: Phylogenies and Genealogies. In: Encyclopedia of Life Sciences (ELS). John Wiley & Sons, Ltd: Chichester.  
DOI: 10.1002/9780470015902.a0020763

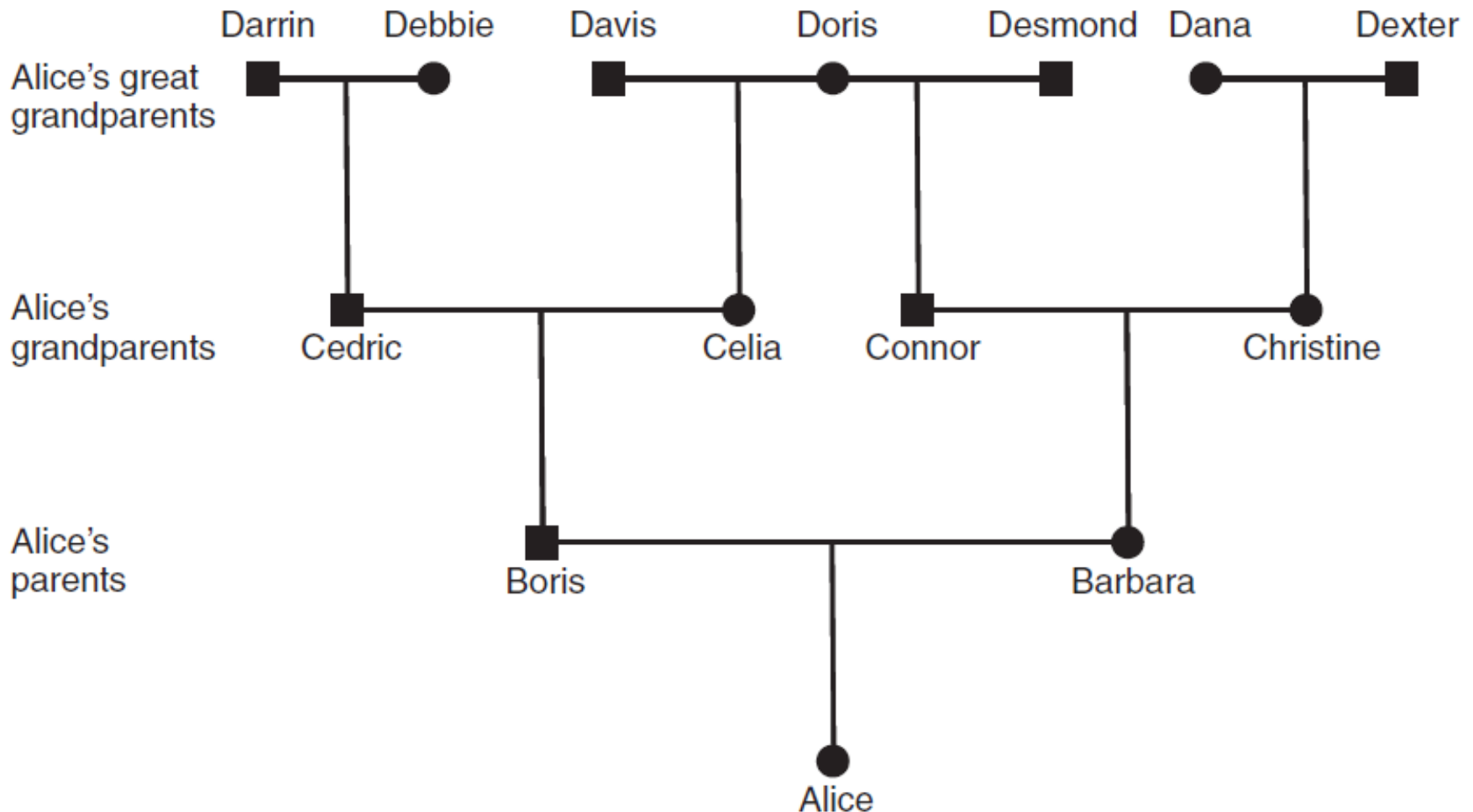




# Non-random mating

**Inbreeding:** mating with relatives

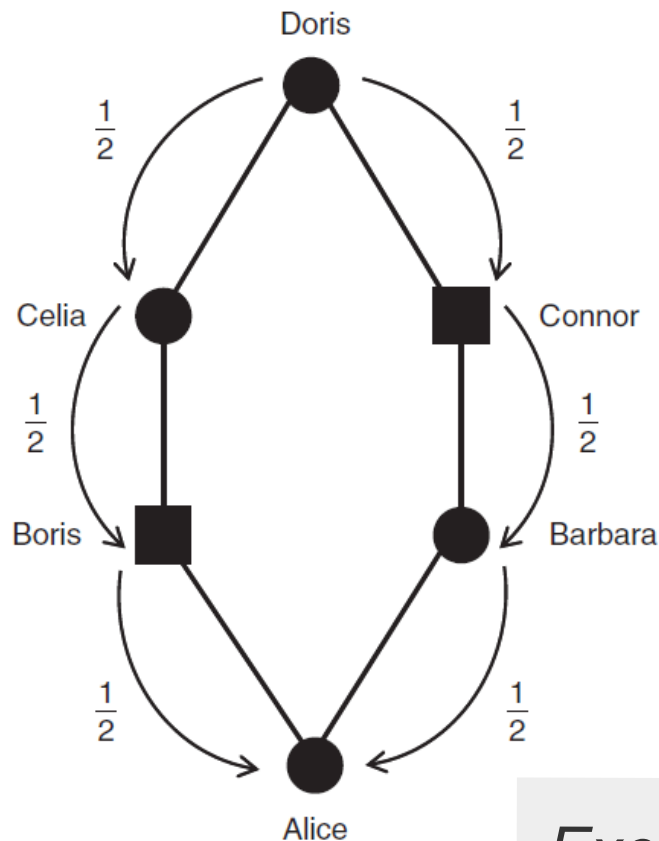
Boris and Barbara are *half first cousins*



# Non-random mating

**Identity by descent (IBD):** two identical alleles are inherited from a common ancestor. **Identity by state (IBS):** identical alleles that are *not* from a common ancestor.

**Inbreeding coefficient  $F$ :** the probability of being IBD at a locus



Generalized Hardy-Weinberg principle:

|                 |            |                 |
|-----------------|------------|-----------------|
| $A_1A_1$        | $A_1A_2$   | $A_2A_2$        |
| $p^2(1-F) + pF$ | $2pq(1-F)$ | $q^2(1-F) + qF$ |

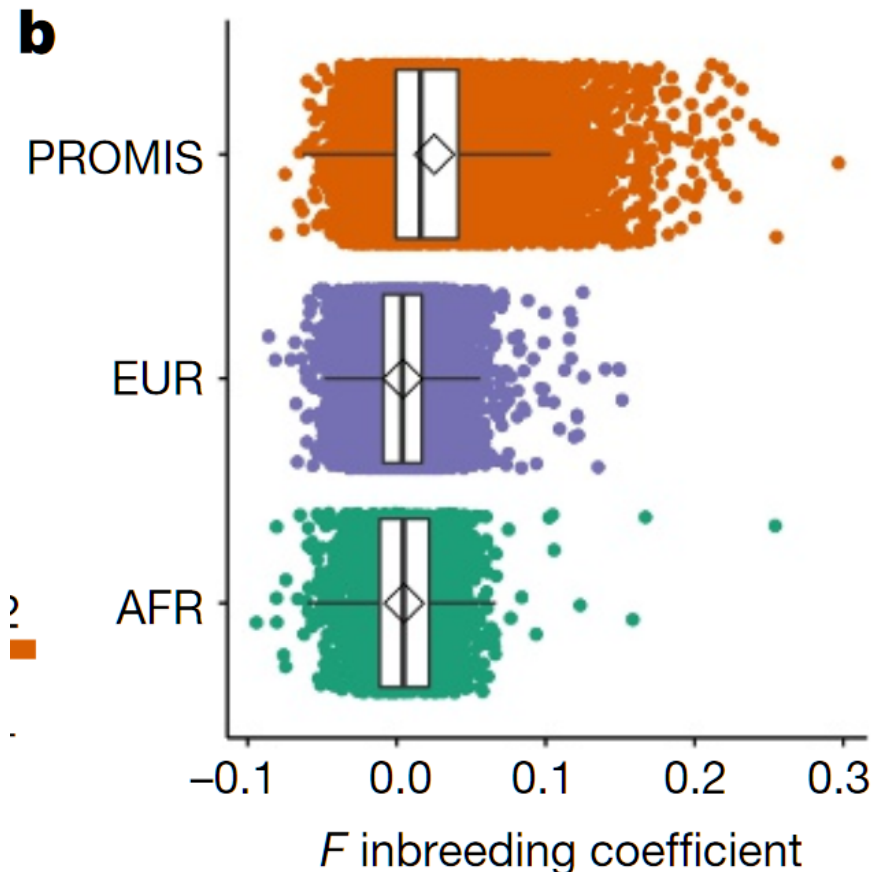
*Exercise:* explain what  $\frac{1}{2}$  means, and why  $F = 1/32$  for Alice

*Exercise:* calculate GT frequencies for  $p = 0.4$

# Human knockouts and phenotypic analysis in a cohort with a high rate of consanguinity

Danish Saleheen<sup>1,2\*</sup>, Pradeep Natarajan<sup>3,4\*</sup>, Irina M. Armean<sup>4,5</sup>, Wei Zhao<sup>1</sup>, Asif Rasheed<sup>2</sup>, Sumeet A. Khetarpal<sup>6</sup>, Hong-Hee Won<sup>7</sup>, Konrad J. Karczewski<sup>4,5</sup>, Anne H. O'Donnell-Luria<sup>4,5,8</sup>, Kaitlin E. Samocha<sup>4,5</sup>, Benjamin Weisburd<sup>4,5</sup>, Namrata Gupta<sup>4</sup>, Mozzam Zaidi<sup>2</sup>, Maria Samuel<sup>2</sup>, Atif Imran<sup>2</sup>, Shahid Abbas<sup>9</sup>, Faisal Majeed<sup>2</sup>, Madiha Ishaq<sup>2</sup>, Saba Akhtar<sup>2</sup>, Kevin Trindade<sup>6</sup>, Megan Mucksavage<sup>6</sup>, Nadeem Qamar<sup>10</sup>, Khan Shah Zaman<sup>10</sup>, Zia Yaqoob<sup>10</sup>, Tahir Saghir<sup>10</sup>, Syed Nadeem Hasan Rizvi<sup>10</sup>, Anis Memon<sup>10</sup>, Nadeem Hayyat Mallick<sup>11</sup>, Mohammad Ishaq<sup>12</sup>, Syed Zahed Rasheed<sup>12</sup>, Fazal-ur-Rehman Memon<sup>13</sup>, Khalid Mahmood<sup>14</sup>, Naveeduddin Ahmed<sup>15</sup>, Ron Do<sup>16,17</sup>, Ronald M. Krauss<sup>18</sup>, Daniel G. MacArthur<sup>4,5</sup>, Stacey Gabriel<sup>4</sup>, Eric S. Lander<sup>4</sup>, Mark J. Daly<sup>4,5</sup>, Philippe Frossard<sup>2§</sup>, John Danesh<sup>19,20§</sup>, Daniel J. Rader<sup>6,21§</sup> & Sekar Kathiresan<sup>3,4§</sup>

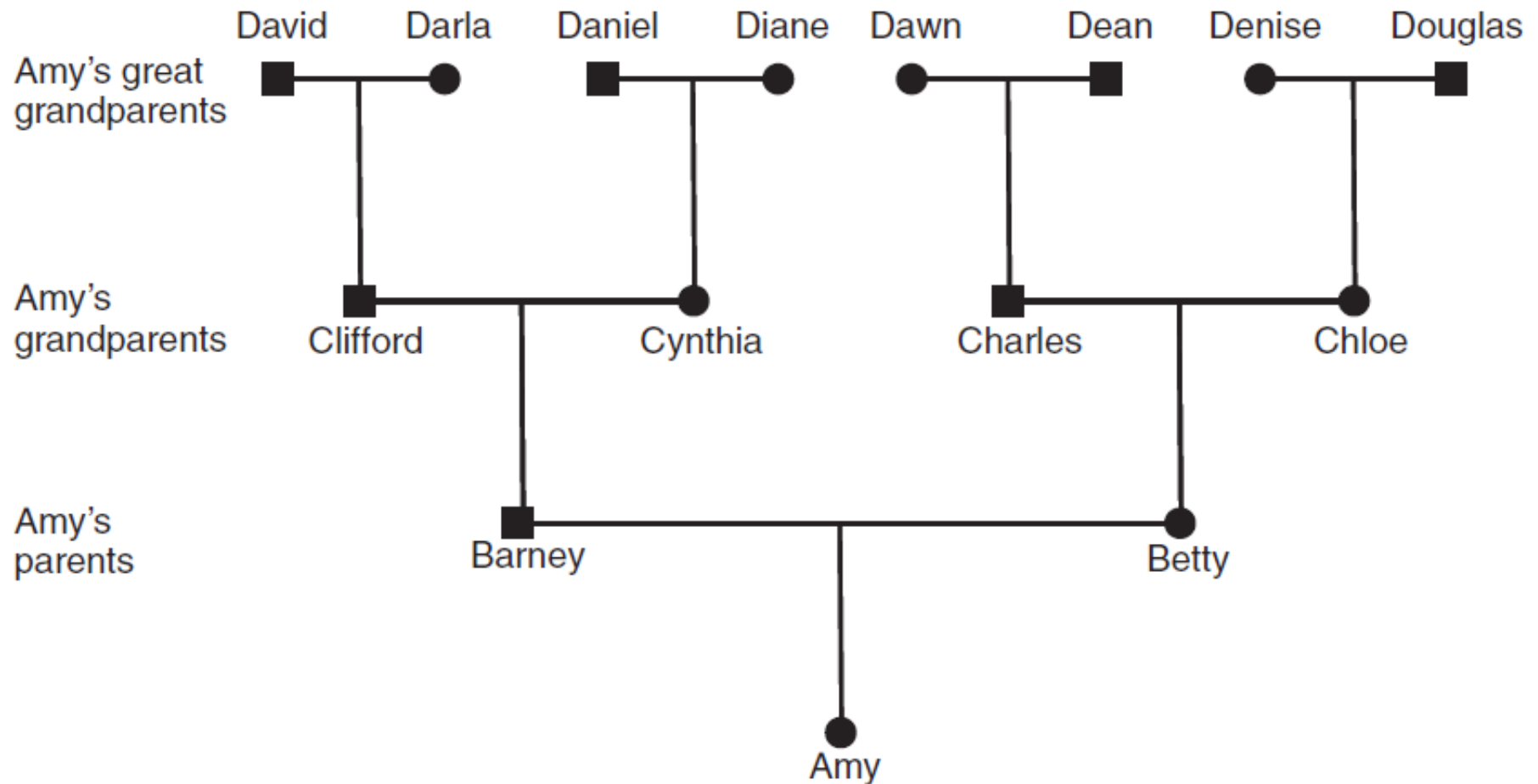
A major goal of biomedicine is to understand the function of every gene in the human genome<sup>1</sup>. Loss-of-function mutations can disrupt both copies of a given gene in humans and phenotypic analysis of such 'human knockouts' can provide insight into gene function. Consanguineous unions are more likely to result in offspring carrying homozygous loss-of-function mutations. In Pakistan, consanguinity rates are notably high<sup>2</sup>. Here we sequence the protein-coding regions of 10,503 adult participants in the Pakistan Risk of Myocardial Infarction Study (PROMIS), designed to understand the determinants of cardiometabolic diseases in individuals from South Asia<sup>3</sup>. We identified individuals carrying homozygous predicted loss-of-function (pLoF) mutations, and performed phenotypic analysis involving more than 200 biochemical and disease traits. We enumerated 49,138 rare (<1% minor allele frequency) pLoF mutations. These pLoF mutations are estimated to knock out 1,317 genes, each in at least one participant.



# Non-random mating

Non-unique  $2^n$  ancestors: everyone is inbred

Not an evolutionary force: affects genotype frequencies, but not allele frequencies. Genotypes are subject to selection, though.





ARTICLE

<https://doi.org/10.1038/s41467-019-12424-x>

OPEN

# Genetic evidence for assortative mating on alcohol consumption in the UK Biobank

Laurence J. Howe<sup>1,2\*</sup>, Daniel J. Lawson<sup>1</sup>, Neil M. Davies<sup>1</sup>, Beate St. Pourcain<sup>1,3,4</sup>, Sarah J. Lewis<sup>1</sup>, George Davey Smith<sup>1,5</sup> & Gibran Hemani<sup>1,5</sup>

Alcohol use is correlated within spouse-pairs, but it is difficult to disentangle effects of alcohol consumption on mate-selection from social factors or the shared spousal environment. We hypothesised that genetic variants related to alcohol consumption may, via their effect on alcohol behaviour, influence mate selection. Here, we find strong evidence that an individual's self-reported alcohol consumption and their genotype at rs1229984, a missense variant in *ADH1B*, are associated with their partner's self-reported alcohol use. Applying Mendelian randomization, we estimate that a unit increase in an individual's weekly alcohol consumption increases partner's alcohol consumption by 0.26 units (95% C.I. 0.15, 0.38;  $P = 8.20 \times 10^{-6}$ ). Furthermore, we find evidence of spousal genotypic concordance for rs1229984, suggesting that spousal concordance for alcohol consumption existed prior to cohabitation. Although the SNP is strongly associated with ancestry, our results suggest some concordance independent of population stratification. Our findings suggest that alcohol behaviour directly influences mate selection.



# Population subdivision, gene flow, admixture

## Population subdivision: restricted migration

Let  $p_i$  be the frequency of the  $A_1$  allele in the  $i$ th subpopulation. Let the relative contribution of this subpopulation to the species or sample be  $c_i$ ,  $\sum c_i = 1$ . Let  $p$  be the average frequency of the  $A_1$  allele across patches,  $p = \sum c_i p_i$ , and let  $q = 1 - p$ . As with the example, the frequencies of genotypes are

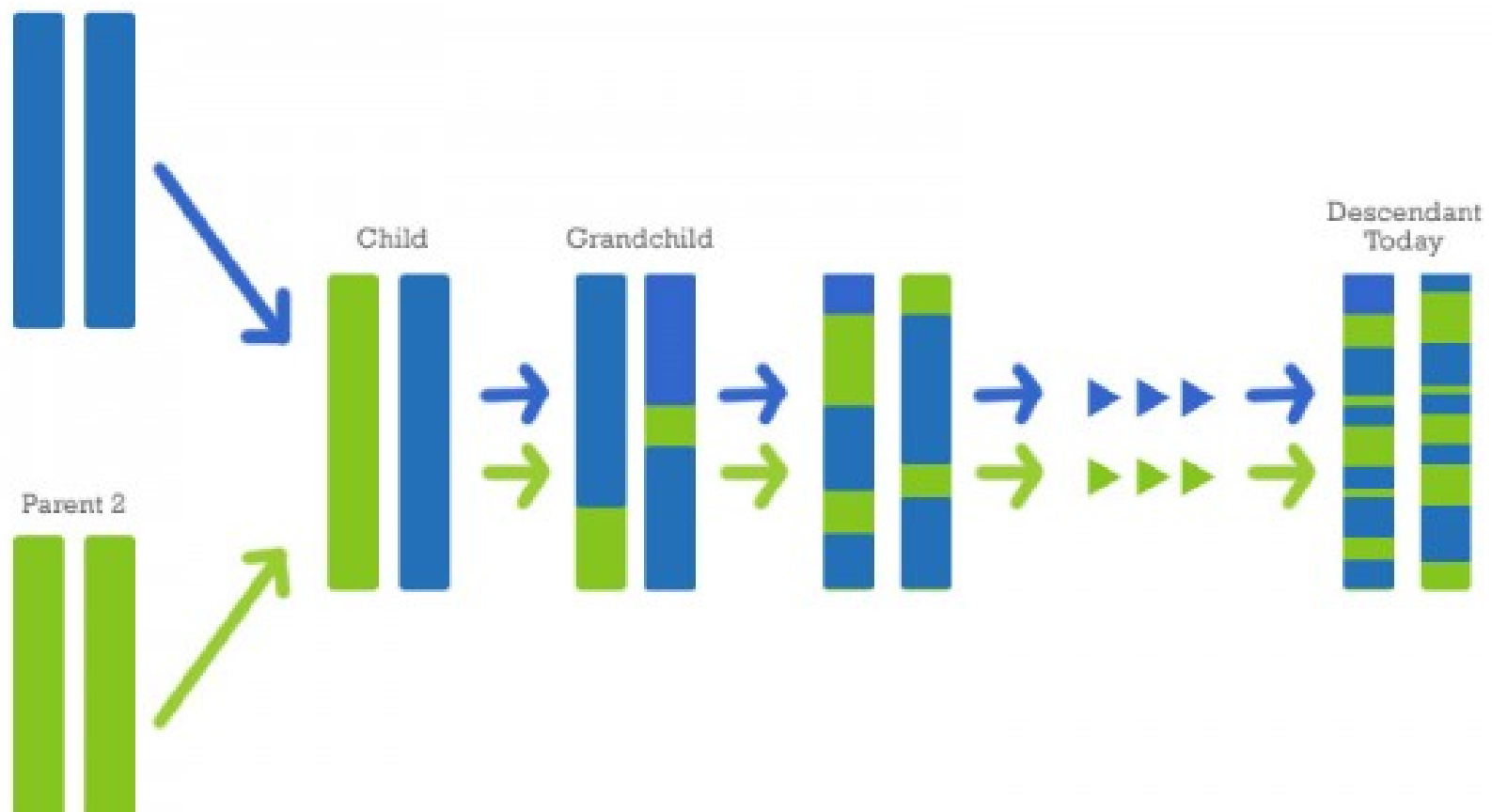
| Genotype:       | $A_1 A_1$                   | $A_1 A_2$           | $A_2 A_2$                   |
|-----------------|-----------------------------|---------------------|-----------------------------|
| In $i$ th patch | $p_i^2$                     | $2p_i q_i$          | $q_i^2$                     |
| In species:     | $\sum c_i p_i^2$            | $\sum c_i 2p_i q_i$ | $\sum c_i q_i^2$            |
| In species:     | $p^2(1 - F_{ST}) + pF_{ST}$ | $2pq(1 - F_{ST})$   | $q^2(1 - F_{ST}) + qF_{ST}$ |

**Gene flow** (gene migration or allele flow): the transfer of genetic variation from one population to another, creates population **admixture**

- Introduces new alleles into population
- Reduces genetic differences between populations, in particular, caused by genetic drift

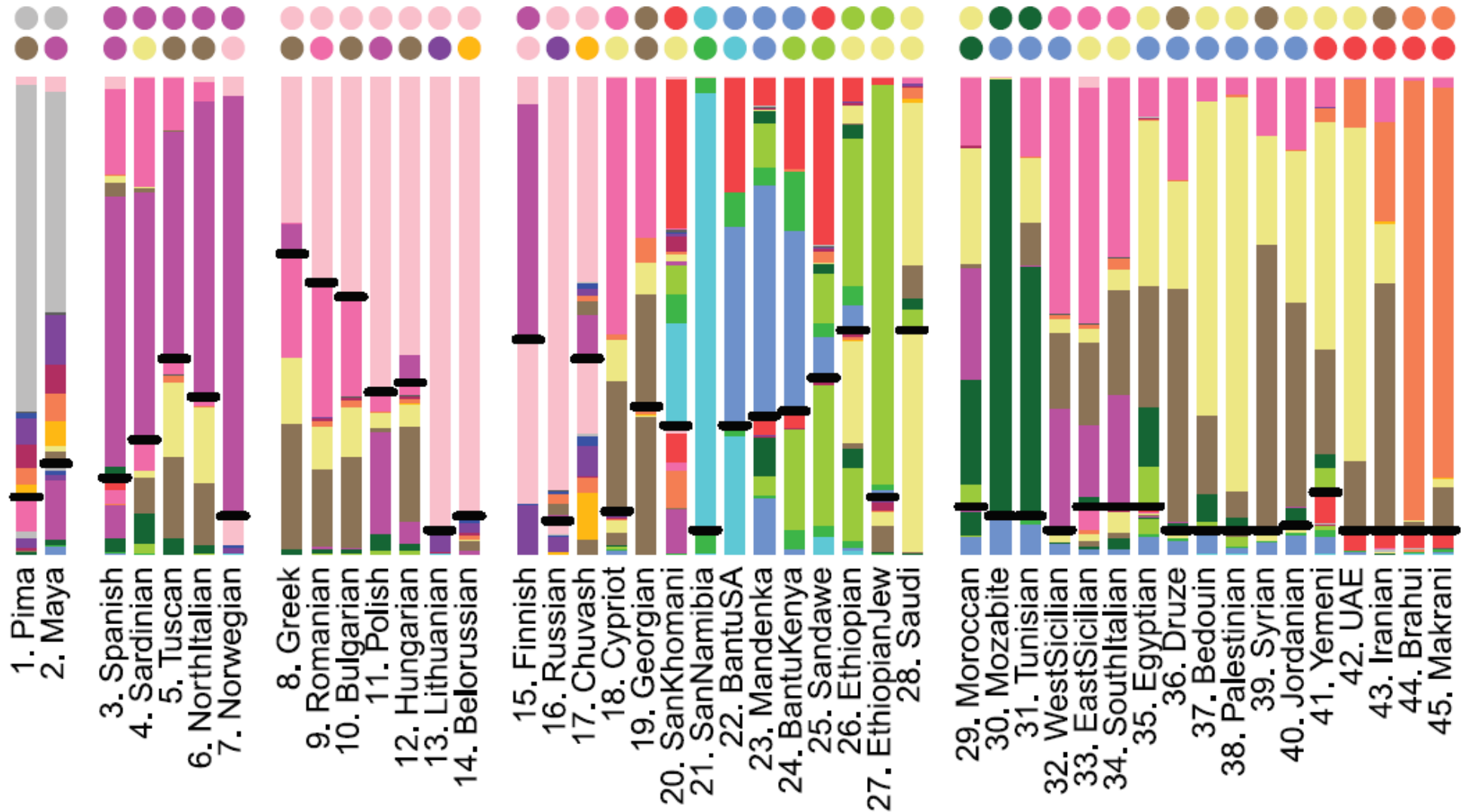
# Population subdivision, gene flow, admixture

**Admixture:** gene flow between previously separated (or partially separated) populations



# A Genetic Atlas of Human Admixture History

Garrett Hellenthal<sup>1</sup>, George B. J. Busby<sup>2</sup>, Gavin Band<sup>3</sup>, James F. Wilson<sup>4</sup>, Cristian Capelli<sup>2</sup>, Daniel Falush<sup>5,\*</sup>, Simon Myers...



# Summary

What changes allele/genotype frequencies?

- **Mutation:** introduction of new alleles into a population
- **Genetic drift:** sampling variation of transmitted alleles
- **Selection:** different probabilities of survival/reproduction depending on genotypes
- **Gene flow:** movement of alleles due to migration
- **Non-random mating** of individuals in a population

# Summary

- Hardy-Weinberg equilibrium describes how zygotes originate from gametes
- Random genetic drift drives alleles to loss or fixation and reduces heterozygosity
- Neutral theory postulates that most inter- and intra-species changes are due to negative selection and random drift
- A coalescent is the lineage of alleles in a sample traced backward in time to their common ancestor allele
- Natural selection changes allele frequencies. It always increases the mean fitness and does so at a rate that is proportional to the genetic variation
- Most new alleles are deleterious and incompletely dominant. They appear by mutation and are subject to negative selection (mutation-selection balance).
- In a finite population, a new advantageous mutation is usually lost because of random drift. On the other hand, a deleterious allele can fix.

# Further reading

- Meyer, D., Harris, E. (2008) Selection Operating on Protein-coding Genes in the Human Genome. In: *Encyclopedia of Life Sciences* (ELS). doi:10.1002/9780470015902.a0020791
- Nei, M., Suzuki, Y., and Nozawa, M. (2010). The neutral theory of molecular evolution in the genomic era. *Annu Rev Genomics Hum Genet* 11, 265–289
- Hurst, L.D. (2009). Genetics and the understanding of selection. *Nature Reviews Genetics* 10, 83–93.
- Fan, S., Hansen, M.E.B., Lo, Y., and Tishkoff, S.A. (2016). Going global by adapting local: A review of recent human adaptation. *Science* 354, 54–59.
- John H. Gillespie – Population Genetics. A concise guide
- John H. Relethford – Human population genetics