

Introgression & Deriving Patterson's D

BIOL 1435

March 7, 2023

Overview

1. Motivation

2. Preliminaries

3. Derivations

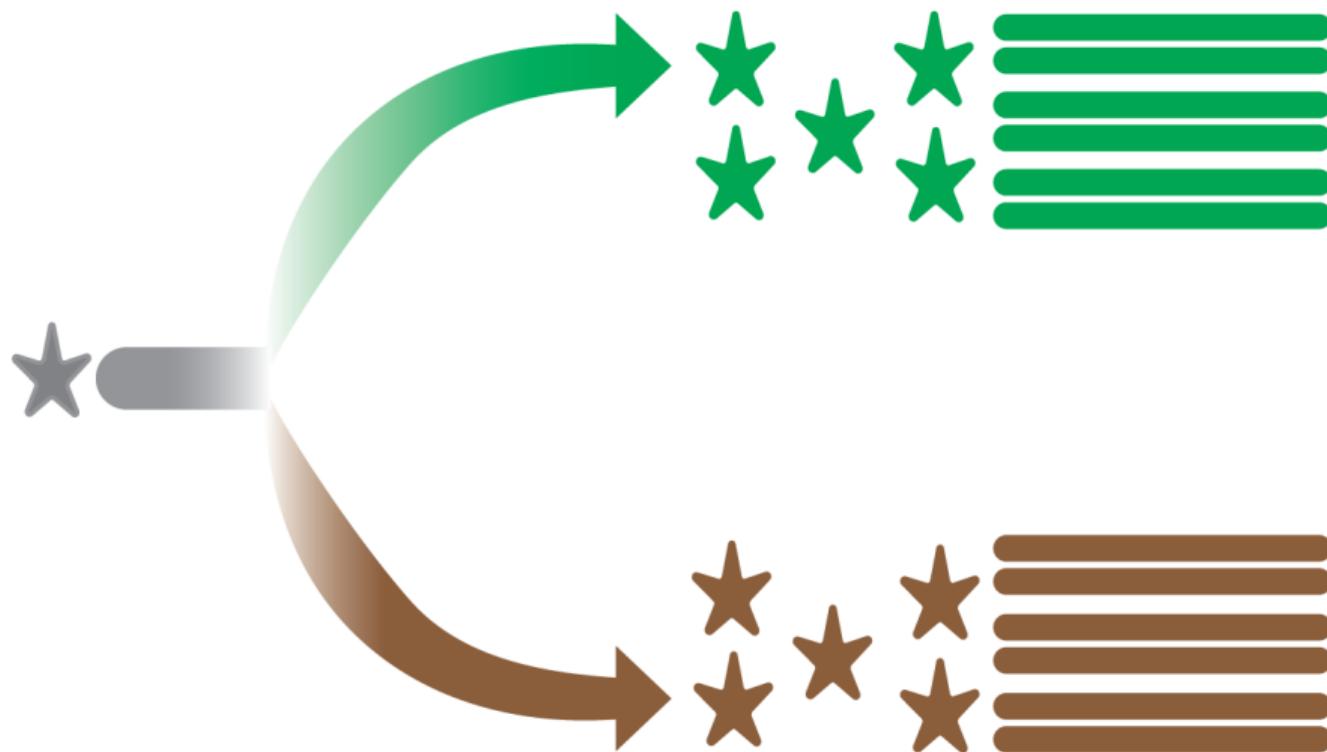
Overview

1. Motivation

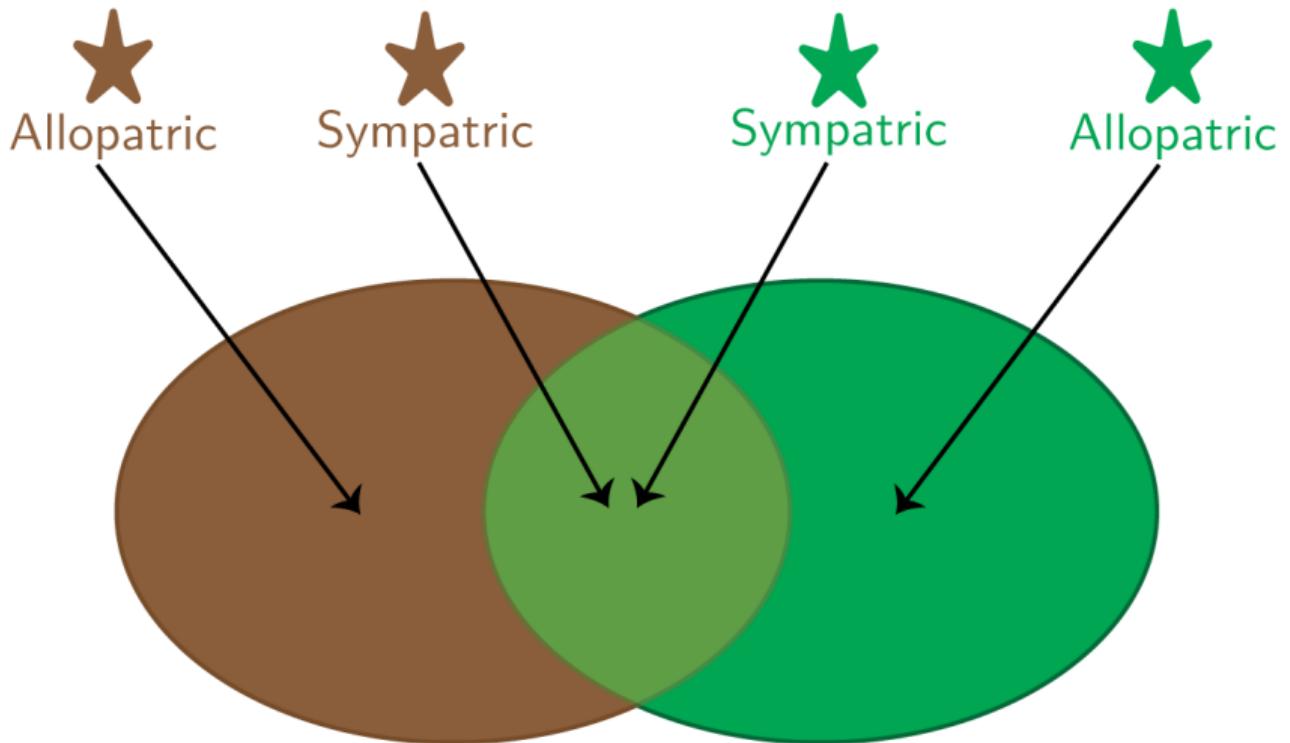
2. Preliminaries

3. Derivations

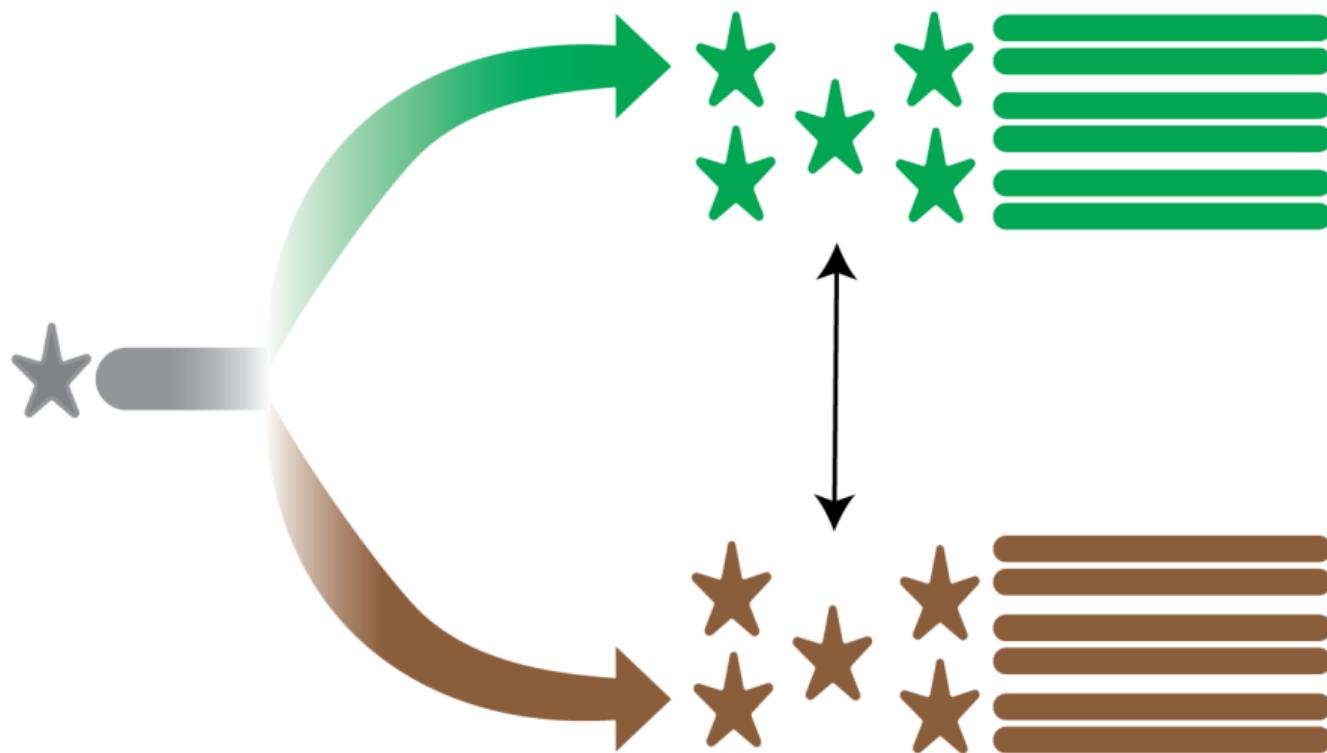
Evolution as a bifurcating process



Sympatry is necessary for gene flow



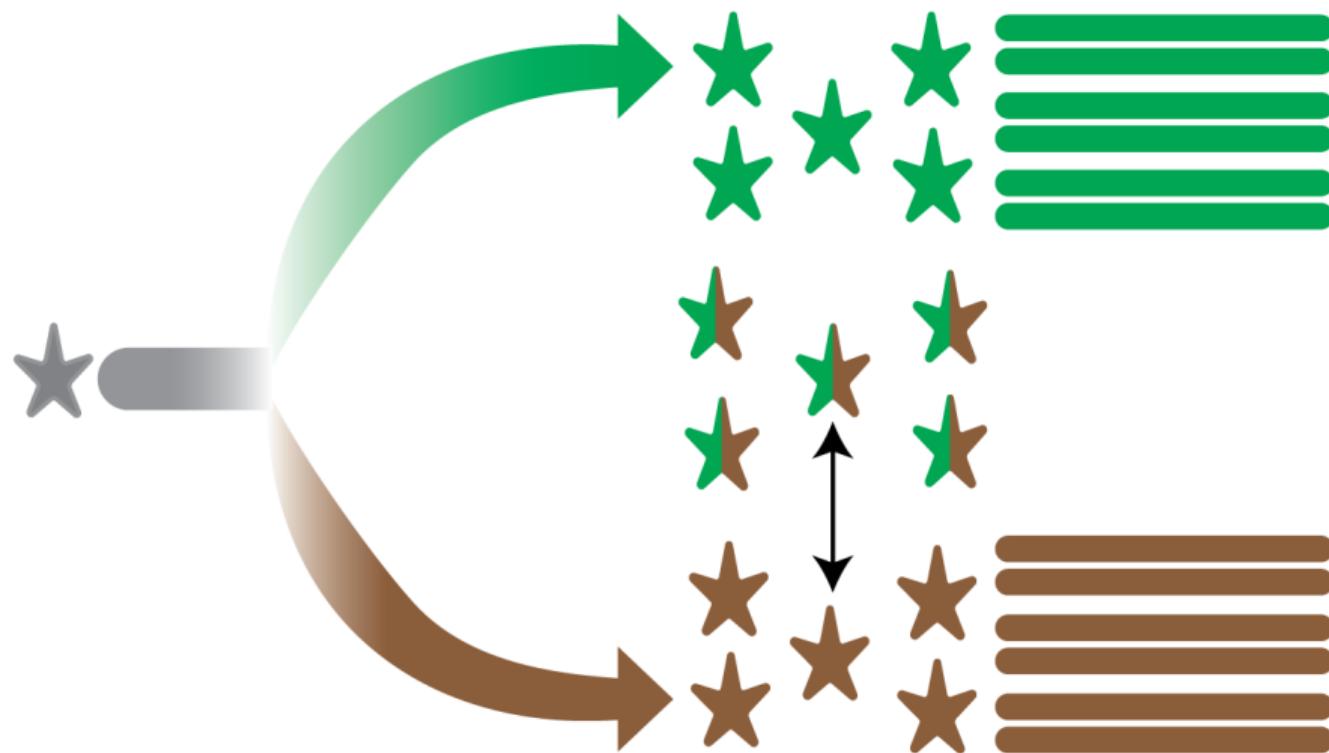
Incomplete RI leads to hybridization



F1 hybrids = uniform mixture of parental ancestry



Hybridization is NOT sufficient for introgression



Incorporation of heterospecific loci via backcrossing



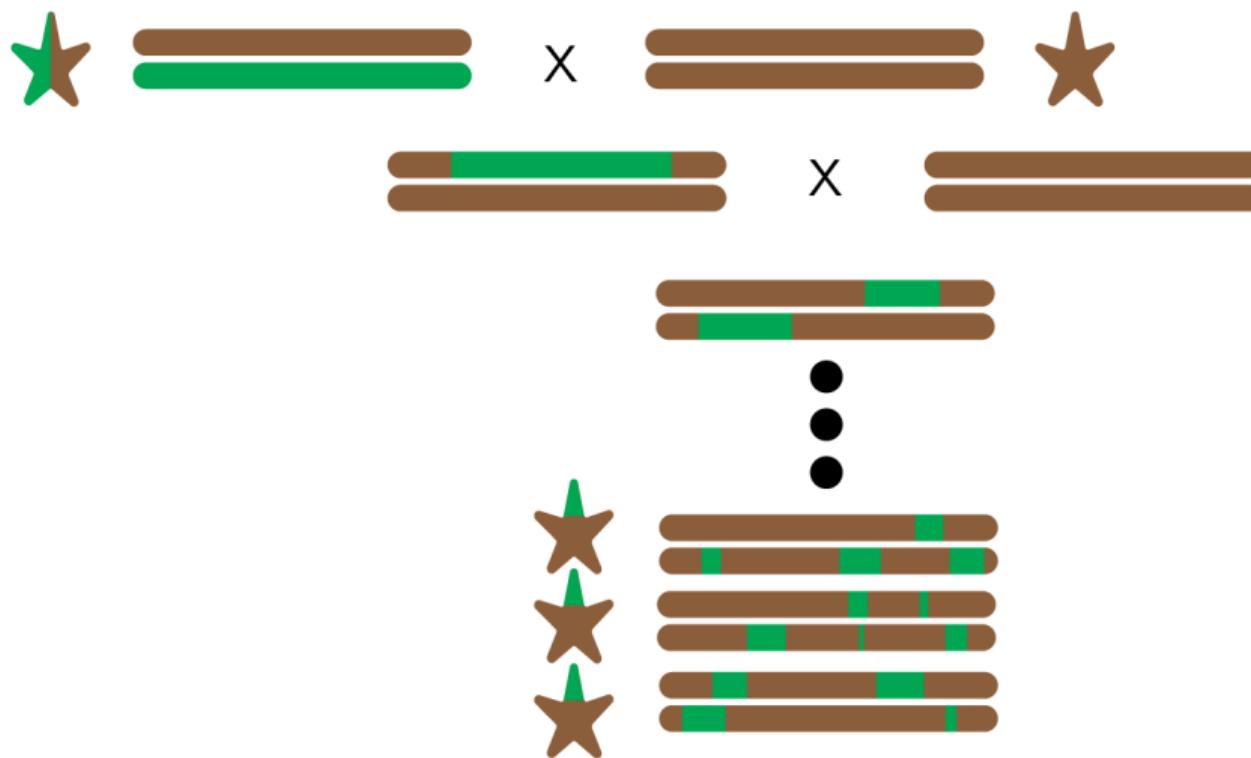
Incorporation of heterospecific loci via backcrossing



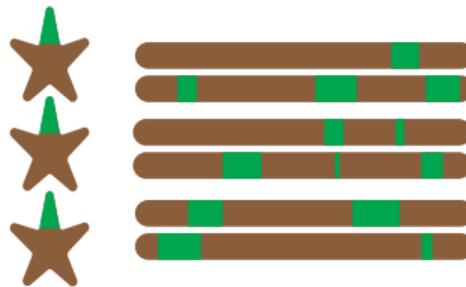
Incorporation of heterospecific loci via backcrossing



Incorporation of heterospecific loci via backcrossing

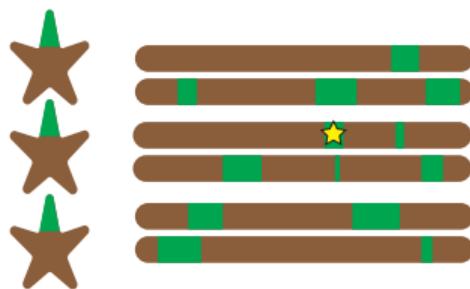


Adaptive introgression



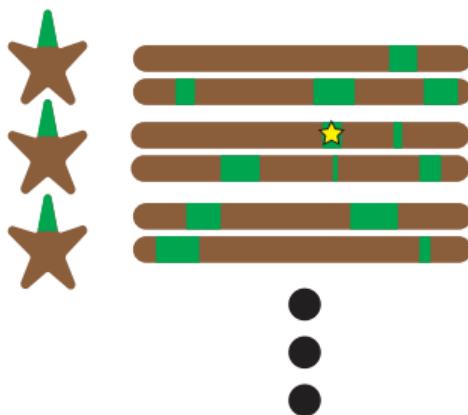
Before selection

Introgression may introduce beneficial mutations



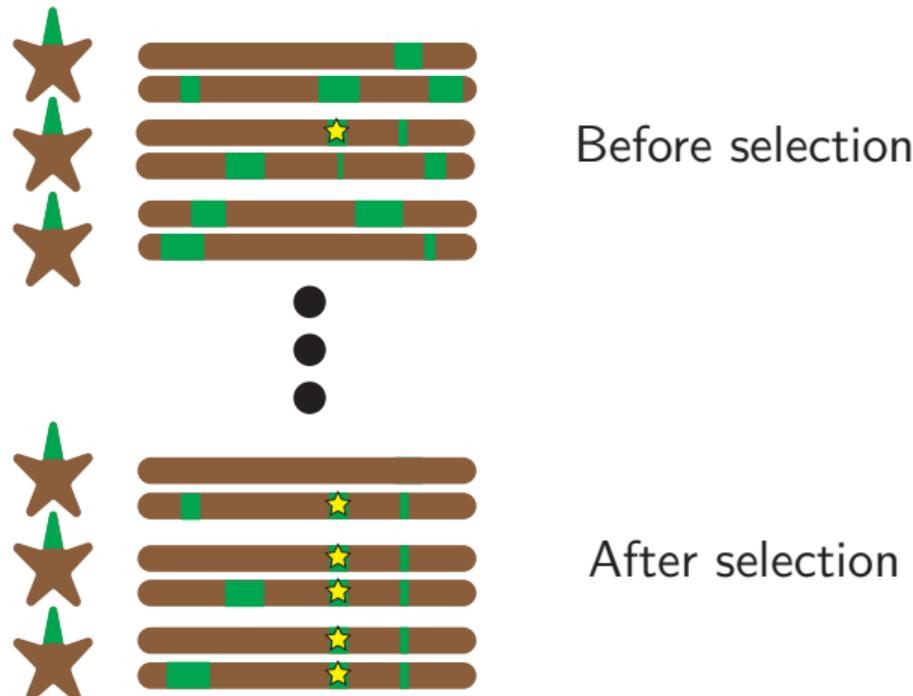
Before selection

The beneficial mutation rises in frequency



Before selection

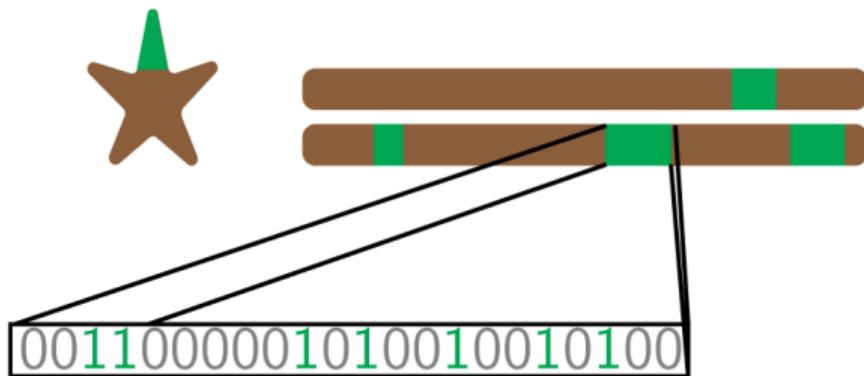
The beneficial mutation is eventually fixed or nearly fixed



Genetic composition of an introgressed segment



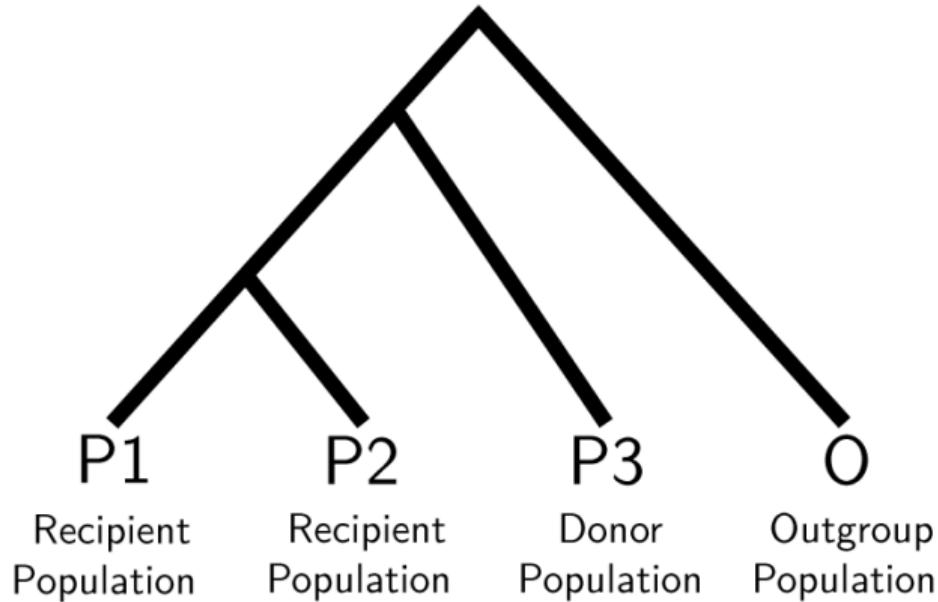
Introgression segments leaves a genomic footprint



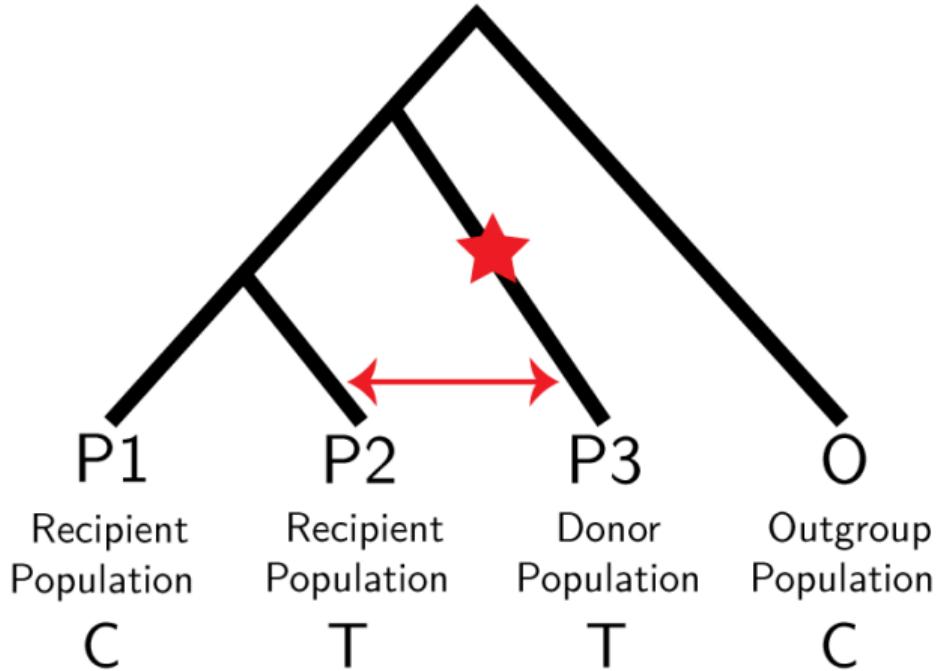
0 = reintroduced ancestral allele

1 = newly introduced derived allele

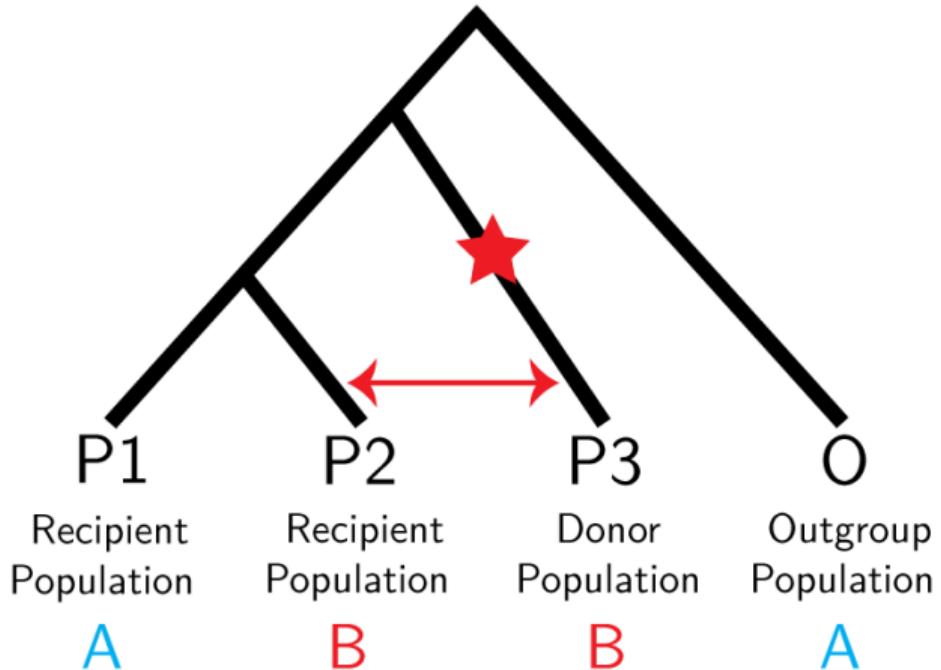
Site pattern tests of introgression



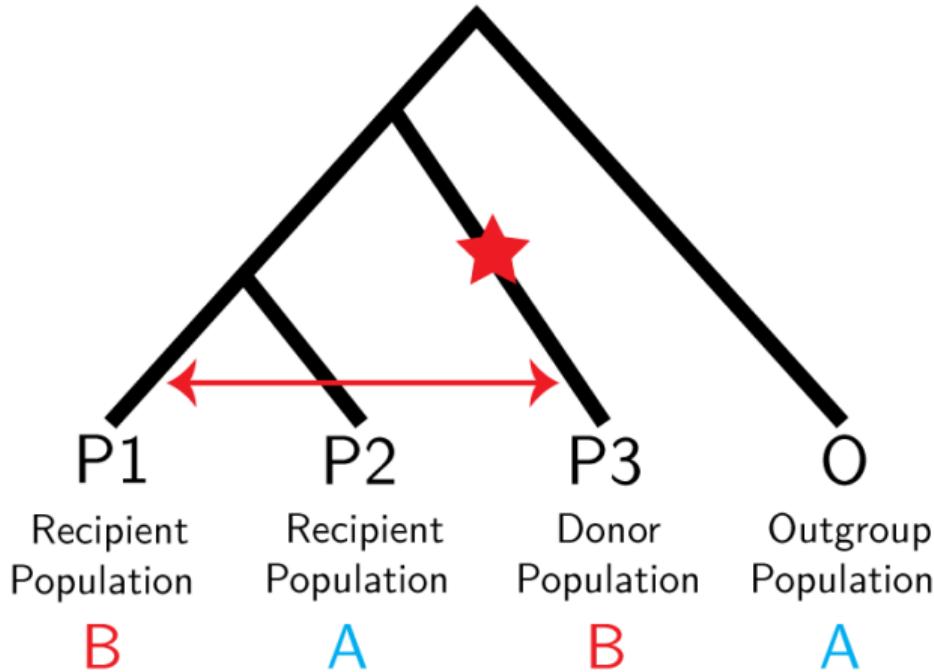
Derived allele sharing ($P_2 \longleftrightarrow P_3$)



Derived allele sharing ($P_2 \longleftrightarrow P_3$)



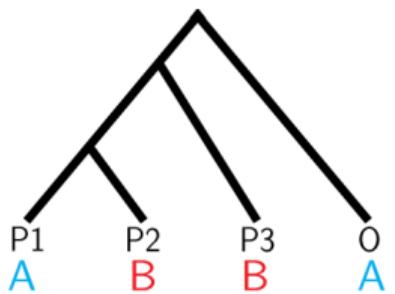
Derived allele sharing ($P_1 \longleftrightarrow P_3$)



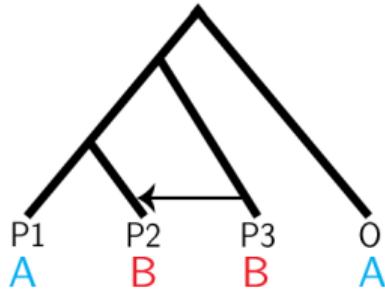
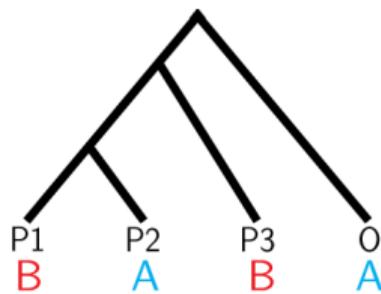
Patterson's D



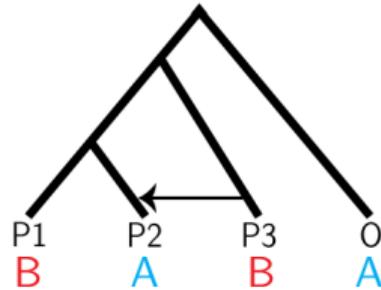
Patterson's D



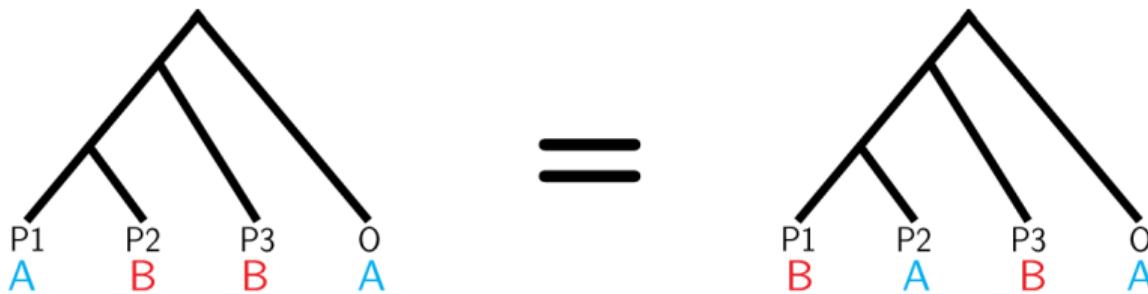
=



>

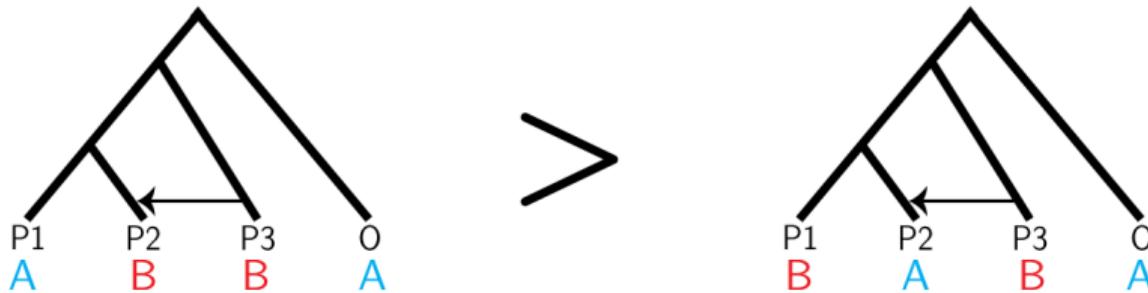


Patterson's D

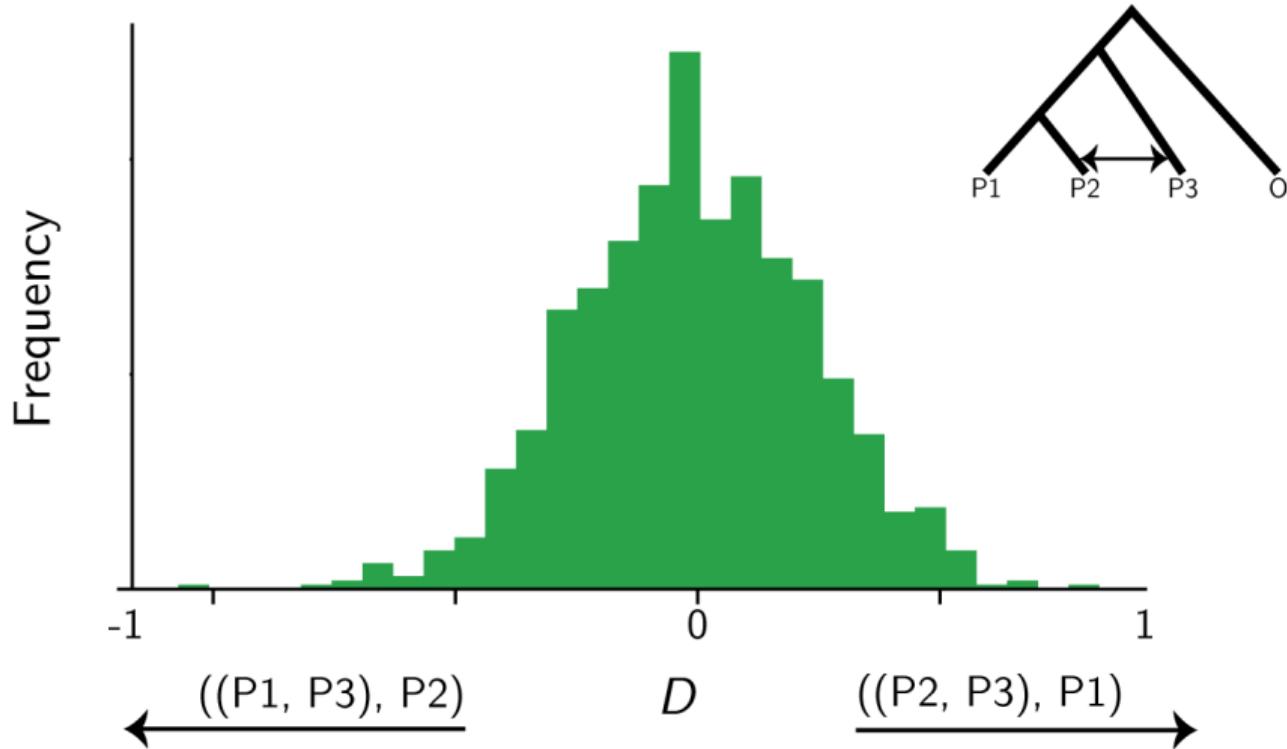


$$\text{Patterson's } D = \frac{\sum_{i=1}^n ABBA_i - \sum_{i=1}^n BABA_i}{\sum_{i=1}^n ABBA_i + \sum_{i=1}^n BABA_i}$$

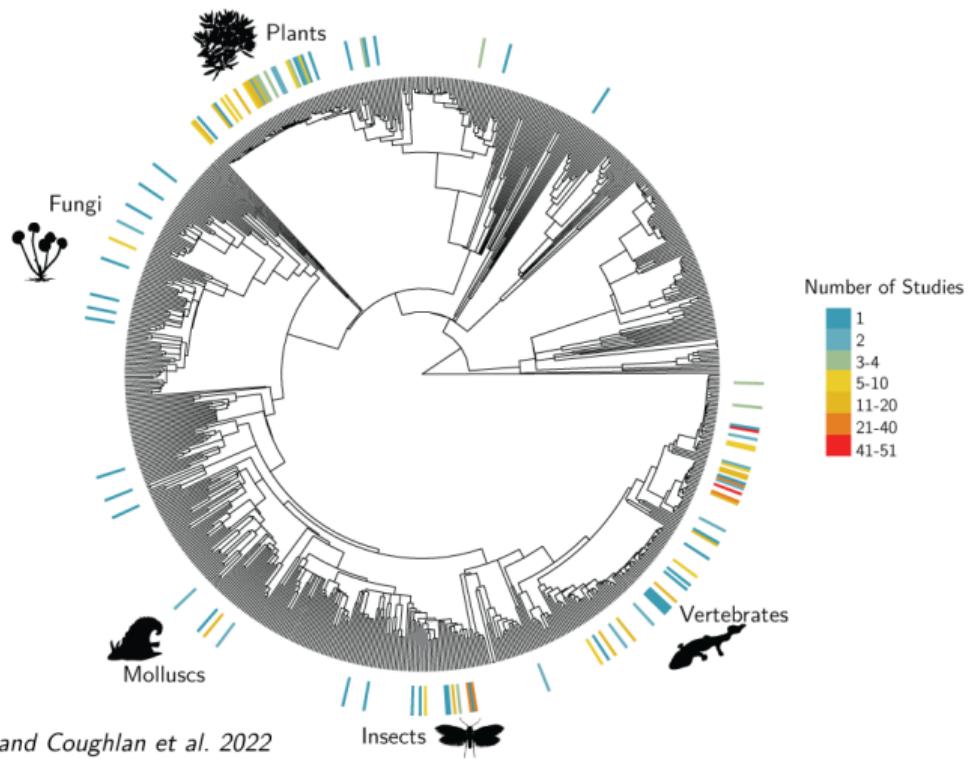
Green et al. 2010 & Durand et al. 2011



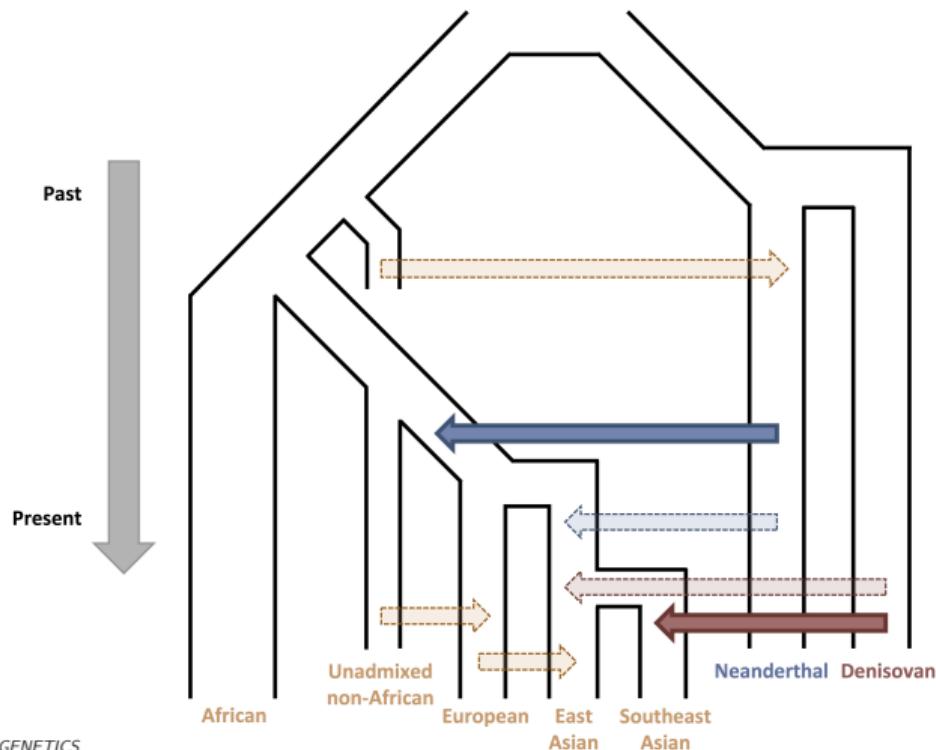
Patterson's D



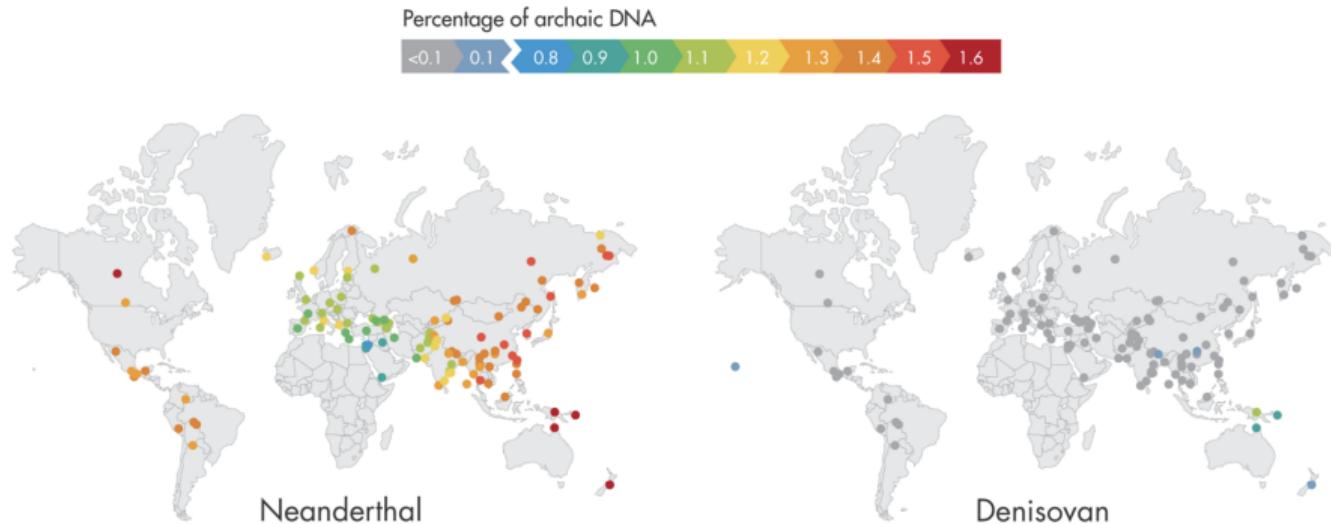
How common is introgression in eukaryotes?



An oversimplified view of human evolution



Global distribution of archaic DNA



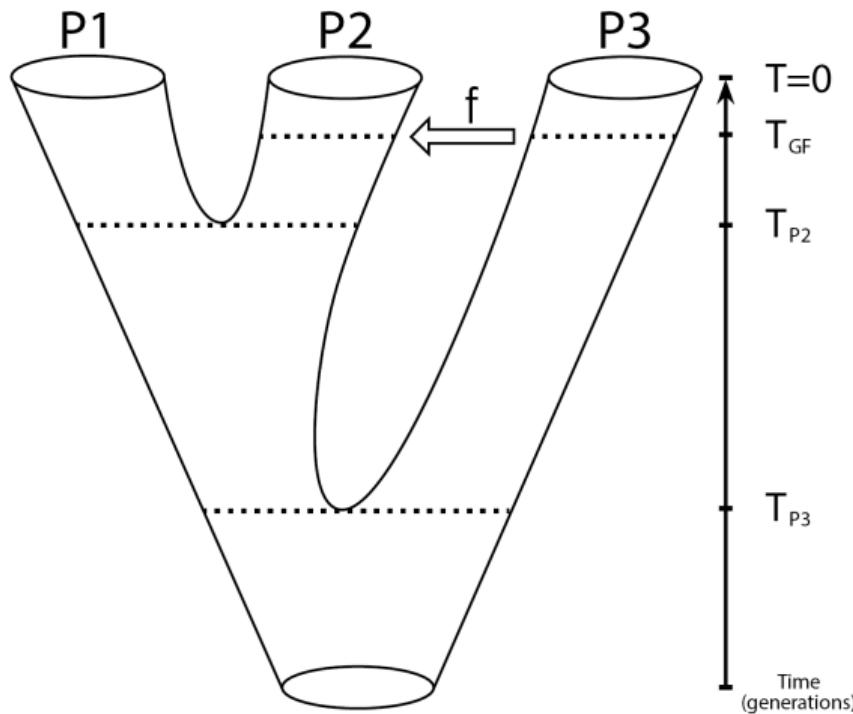
Overview

1. Motivation

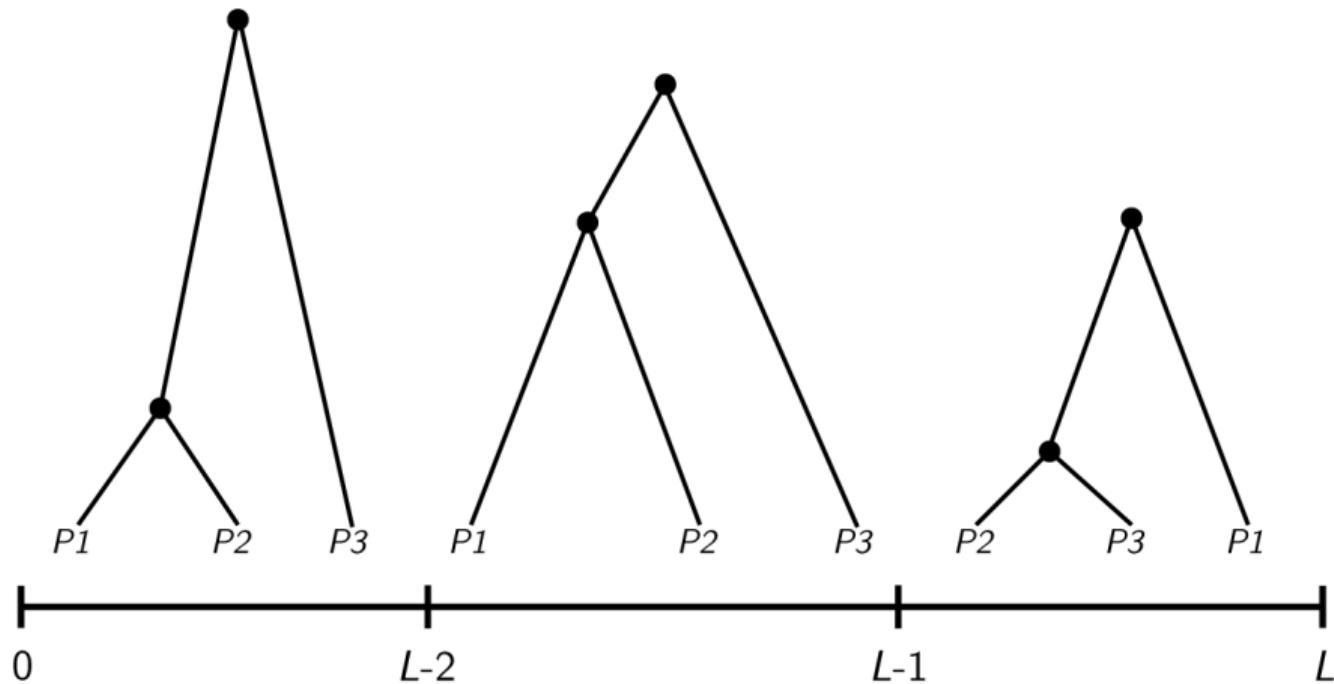
2. Preliminaries

3. Derivations

Instantaneous unidirectional admixture model



Each locus has its own coalescent history



Probability of gene flow

Equation

$$Pr(\text{gene flow}) = f \quad (1)$$

$$Pr(\text{no gene flow}) = (1 - f) \quad (2)$$

Where f represents the admixture proportion—the probability that any lineage from $P3$ migrates to $P2$.

Probability of no coalescence during time interval t

Probability of no coalescence during time interval t

Equation

$$Pr(\text{no coalescences}) = \left(1 - \frac{1}{2N}\right)^t \quad (3)$$

Where t denotes the time interval where coalescence can occur.

Probability of coalescence during time interval t

Probability of coalescence during time interval t

Equation

$$Pr(\text{coalescences}) = 1 - \left(1 - \frac{1}{2N}\right)^t \quad (4)$$

Where t denotes the time interval where coalescence can occur.

Expected time of coalescence T_2

Expected time of coalescence T_2

Equation

$$T_2 = 2N \tag{5}$$

Expected time of coalescence | coalescence during t

Expected time of coalescence | coalescence during t

Equation

$$i \sim T_{geo}(t|p, c) \quad (6)$$

Expected time of coalescence | coalescence during t

Equation

$$i \sim Tgeo(t|p, c) \quad (6)$$

$$f(i) = \frac{\frac{1}{2N} \left(1 - \frac{1}{2N}\right)^{i-1}}{1 - \left(1 - \frac{1}{2N}\right)^c} \quad (7)$$

Expected time of coalescence | coalescence during t

Equation

$$i \sim Tgeo(t|p, c) \quad (6)$$

$$f(i) = \frac{\frac{1}{2N} \left(1 - \frac{1}{2N}\right)^{i-1}}{1 - \left(1 - \frac{1}{2N}\right)^c} \quad (7)$$

$$\bar{t} = \mathbb{E}(i) = \sum_{i=1}^c i \frac{\frac{1}{2N} \left(1 - \frac{1}{2N}\right)^{i-1}}{1 - \left(1 - \frac{1}{2N}\right)^c} = \frac{2N - \left(\left(1 - \frac{1}{2N}\right)^c (c + 2N)\right)}{1 - \left(1 - \frac{1}{2N}\right)^c} \quad (8)$$

Where c denotes the time interval where coalescence must occur.

Overview

1. Motivation

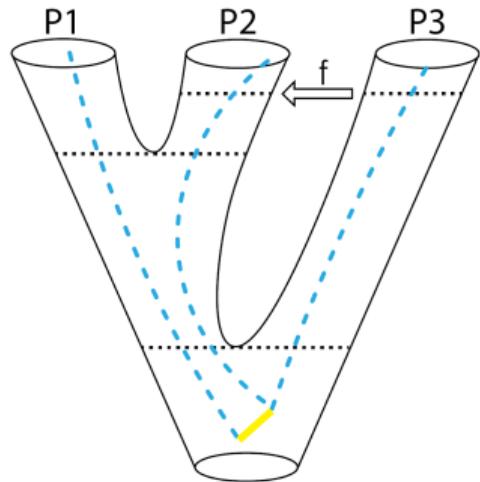
2. Preliminaries

3. Derivations

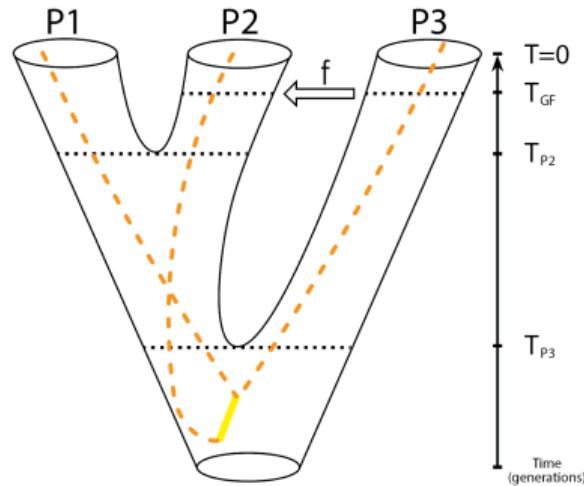
Coalescent history 1

No gene flow from $P_3 \rightarrow P_2$, P_1 & P_2 don't coalesce between T_{P_2} & T_{P_3} , and P_1/P_2 & P_3 coalesce first.

ABBA



BABA



Coalescent history 1

No gene flow from $P_3 \rightarrow P_2$, P_1 & P_2 don't coalesce between T_{P_2} & T_{P_3} , and P_1 / P_2 & P_3 coalesce first.

Coalescent history 1

No gene flow from $P_3 \rightarrow P_2$, P_1 & P_2 don't coalesce between T_{P_2} & T_{P_3} , and P_1 / P_2 & P_3 coalesce first.

Derivation

$$Pr(\text{no gene flow}) = (1 - f) \quad (9)$$

Coalescent history 1

No gene flow from $P_3 \rightarrow P_2$, P_1 & P_2 don't coalesce between T_{P_2} & T_{P_3} , and P_1 / P_2 & P_3 coalesce first.

Derivation

$$Pr(\text{no gene flow}) = (1 - f) \quad (9)$$

$$Pr(P_1 \& P_2 \text{ don't coalesce between } T_{P_2} \& T_{P_3}) = \left(1 - \frac{1}{2N}\right)^{T_{P_3} - T_{P_2}} \quad (10)$$

Coalescent history 1

No gene flow from $P_3 \rightarrow P_2$, P_1 & P_2 don't coalesce between T_{P_2} & T_{P_3} , and P_1 / P_2 & P_3 coalesce first.

Derivation

$$Pr(\text{no gene flow}) = (1 - f) \quad (9)$$

$$Pr(P_1 \& P_2 \text{ don't coalesce between } T_{P_2} \& T_{P_3}) = \left(1 - \frac{1}{2N}\right)^{T_{P_3} - T_{P_2}} \quad (10)$$

$$Pr(P_1/P_2 \& P_3 \text{ coalesce first}) = \frac{1}{3} \quad (11)$$

Coalescent history 1

No gene flow from $P_3 \rightarrow P_2$, P_1 & P_2 don't coalesce between T_{P_2} & T_{P_3} , and P_1 / P_2 & P_3 coalesce first.

Derivation

$$Pr(\text{no gene flow}) = (1 - f) \quad (9)$$

$$Pr(P_1 \& P_2 \text{ don't coalesce between } T_{P_2} \& T_{P_3}) = \left(1 - \frac{1}{2N}\right)^{T_{P_3} - T_{P_2}} \quad (10)$$

$$Pr(P_1/P_2 \& P_3 \text{ coalesce first}) = \frac{1}{3} \quad (11)$$

$$\mathbb{E} \left(\text{Branch length between the } 1^{\text{st}} \& 2^{\text{nd}} \text{ coalescent event} \right) = 2N \quad (12)$$

Coalescent history 1

No gene flow from $P_3 \rightarrow P_2$, P_1 & P_2 don't coalesce between T_{P_2} & T_{P_3} , and P_1/P_2 & P_3 coalesce first.

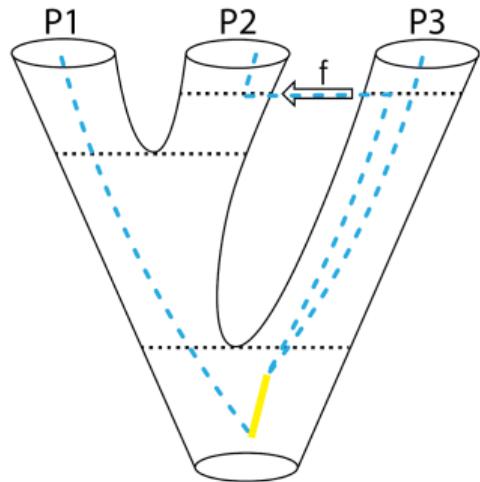
Derivation

$$\mathbb{E}(C_{ABBA_1}) = \mathbb{E}(C_{BABA_1}) = (1-f) \left(1 - \frac{1}{2N}\right)^{T_{P_3} - T_{P_2}} \frac{2N}{3} \quad (13)$$

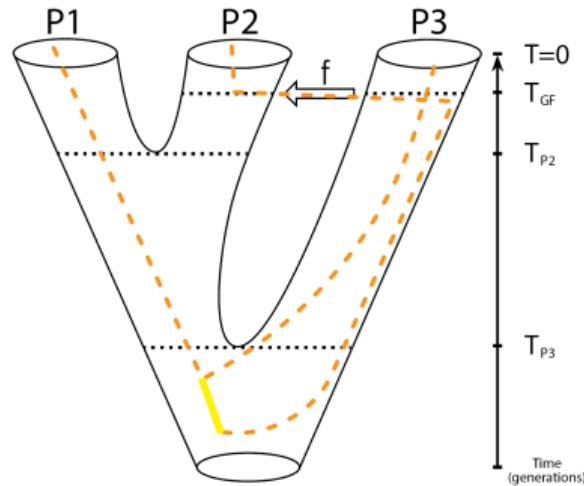
Coalescent history 2

Gene flow from $P_3 \rightarrow P_2$, P_2 & P_3 don't coalesce between T_{GF} & T_{P_3} , and P_1 / P_2 & P_3 coalesce first.

ABBA



BABA



Coalescent history 2

Gene flow from $P_3 \rightarrow P_2$, P_2 & P_3 don't coalesce between T_{GF} & T_{P3} , and P_1 / P_2 & P_3 coalesce first.

Coalescent history 2

Gene flow from $P_3 \rightarrow P_2$, P_2 & P_3 don't coalesce between T_{GF} & T_{P3} , and P_1 / P_2 & P_3 coalesce first.

Derivation

$$Pr(\text{gene flow}) = f \quad (14)$$

Coalescent history 2

Gene flow from $P3 \rightarrow P2$, $P2$ & $P3$ don't coalesce between T_{GF} & T_{P3} , and $P1 / P2$ & $P3$ coalesce first.

Derivation

$$Pr(\text{gene flow}) = f \quad (14)$$

$$Pr(P2 \& P3 \text{ don't coalesce between } T_{GF} \& T_{P3}) = \left(1 - \frac{1}{2N}\right)^{T_{P3} - T_{GF}} \quad (15)$$

Coalescent history 2

Gene flow from $P3 \rightarrow P2$, $P2$ & $P3$ don't coalesce between T_{GF} & T_{P3} , and $P1 / P2$ & $P3$ coalesce first.

Derivation

$$Pr(\text{gene flow}) = f \quad (14)$$

$$Pr(P2 \& P3 \text{ don't coalesce between } T_{GF} \& T_{P3}) = \left(1 - \frac{1}{2N}\right)^{T_{P3} - T_{GF}} \quad (15)$$

$$Pr(P1/P2 \& P3 \text{ coalesce first}) = \frac{1}{3} \quad (16)$$

Coalescent history 2

Gene flow from $P3 \rightarrow P2$, $P2$ & $P3$ don't coalesce between T_{GF} & T_{P3} , and $P1 / P2$ & $P3$ coalesce first.

Derivation

$$Pr(\text{gene flow}) = f \quad (14)$$

$$Pr(P2 \& P3 \text{ don't coalesce between } T_{GF} \& T_{P3}) = \left(1 - \frac{1}{2N}\right)^{T_{P3} - T_{GF}} \quad (15)$$

$$Pr(P1/P2 \& P3 \text{ coalesce first}) = \frac{1}{3} \quad (16)$$

$$\mathbb{E} \left(\text{Branch length between the 1}^{st} \& 2^{nd} \text{ coalescent event} \right) = 2N \quad (17)$$

Coalescent history 2

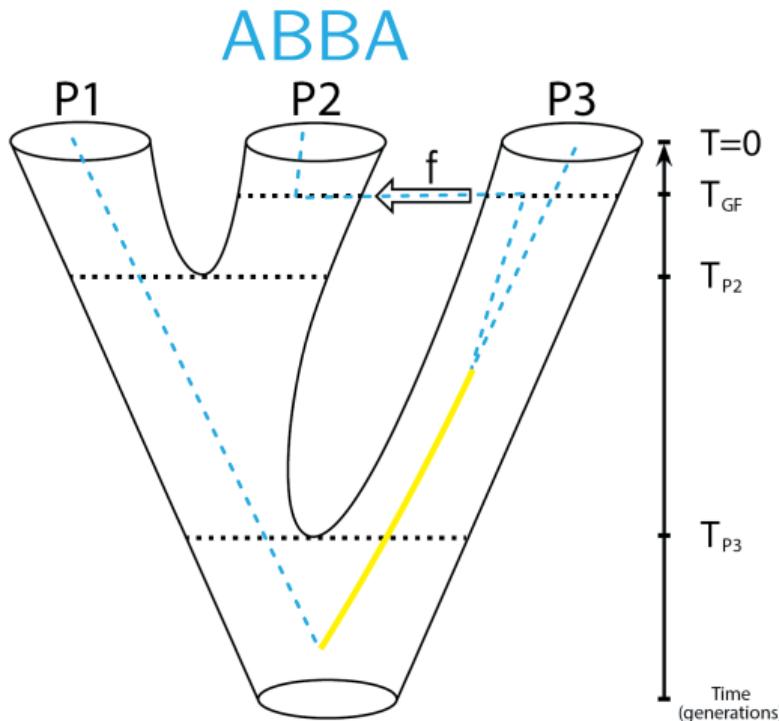
Gene flow from $P_3 \rightarrow P_2$, P_2 & P_3 don't coalesce between T_{GF} & T_{P_3} , and P_1 / P_2 & P_3 coalesce first.

Derivation

$$\mathbb{E}(C_{ABBA_2}) = \mathbb{E}(C_{BABA_2}) = f \left(1 - \frac{1}{2N}\right)^{T_{P_3} - T_{GF}} \frac{2N}{3} \quad (18)$$

Coalescent history 3

Gene flow from $P_3 \rightarrow P_2$ and $P_2 \& P_3$ coalesce between T_{GF} & T_{P_3} .



Coalescent history 3

Gene flow from $P_3 \rightarrow P_2$ and $P_2 \& P_3$ coalesce between T_{GF} & T_{P3} .

Coalescent history 3

Gene flow from $P_3 \rightarrow P_2$ and $P_2 \& P_3$ coalesce between T_{GF} & T_{P3} .

Derivation

$$Pr(\text{gene flow}) = f \quad (19)$$

Coalescent history 3

Gene flow from $P_3 \rightarrow P_2$ and $P_2 \& P_3$ coalesce between T_{GF} & T_{P3} .

Derivation

$$Pr(\text{gene flow}) = f \quad (19)$$

$$Pr(P_2 \& P_3 \text{ coalesce between } T_{GF} \& T_{P3}) = 1 - \left(1 - \frac{1}{2N}\right)^{T_{P3} - T_{GF}} \quad (20)$$

Coalescent history 3

Gene flow from $P_3 \rightarrow P_2$ and $P_2 \& P_3$ coalesce between T_{GF} & T_{P3} .

Derivation

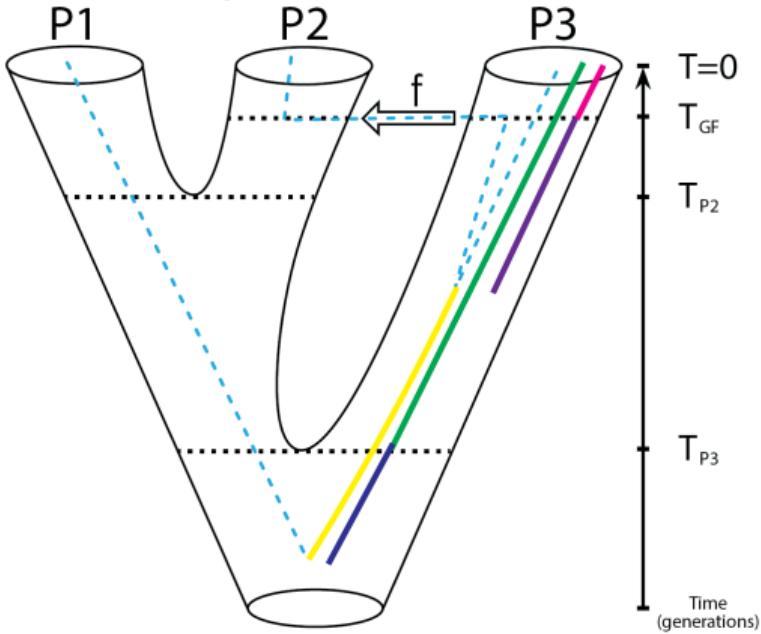
$$Pr(\text{gene flow}) = f \quad (19)$$

$$Pr(P_2 \& P_3 \text{ coalesce between } T_{GF} \& T_{P3}) = 1 - \left(1 - \frac{1}{2N}\right)^{T_{P3} - T_{GF}} \quad (20)$$

$$\mathbb{E} \left(\text{Branch length between the 1}^{\text{st}} \& 2^{\text{nd}} \text{ coalescent event} \right) = (T_{P3} + 2N) - (T_{GF} + \bar{t}) \quad (21)$$

Coalescent history 3 (\bar{t})

$$IBL = T_{P_3} + 2N - (T_{GF} + \bar{t})$$



Coalescent history 3

Gene flow from $P_3 \rightarrow P_2$ and $P_2 \& P_3$ coalesce between T_{GF} & T_{P3} .

Derivation

$$\mathbb{E}(C_{ABBA_3}) = f \left(1 - \left(1 - \frac{1}{2N} \right)^{T_{P3} - T_{GF}} \right) ((T_{P3} + 2N) - (T_{GF} + \bar{t})) \quad (22)$$

Coalescent history 3

Gene flow from $P_3 \rightarrow P_2$ and $P_2 \& P_3$ coalesce between T_{GF} & T_{P3} .

Derivation

$$\mathbb{E}(C_{ABBA_3}) = f \left(1 - \left(1 - \frac{1}{2N} \right)^{T_{P3} - T_{GF}} \right) ((T_{P3} + 2N) - (T_{GF} + \bar{t})) \quad (22)$$

$$\mathbb{E}(C_{BABA_3}) = 0 \quad (23)$$

$$\mathbb{E}(ABBA)$$

Derivation

$$\mathbb{E}(\tau_{ABBA}) = C_{ABBA_1} + C_{ABBA_2} + C_{ABBA_3} \quad (24)$$

$$\mathbb{E}(ABBA)$$

Derivation

$$\mathbb{E}(\tau_{ABBA}) = C_{ABBA_1} + C_{ABBA_2} + C_{ABBA_3} \quad (24)$$

$$\begin{aligned} \mathbb{E}(\tau_{ABBA}) &= (1-f) \left(\frac{2N}{3} \left(1 - \frac{1}{2N}\right)^{T_{P3}-T_{P2}} \right) \\ &\quad + (f) \left(\left(\frac{2N}{3} \left(1 - \frac{1}{2N}\right)^{T_{P3}-T_{GF}} \right) + (T_{P3} - T_{GF}) \right) \end{aligned} \quad (25)$$

$$\mathbb{E}(ABBA)$$

Derivation

$$\mathbb{E}(\tau_{ABBA}) = C_{ABBA_1} + C_{ABBA_2} + C_{ABBA_3} \quad (24)$$

$$\begin{aligned} \mathbb{E}(\tau_{ABBA}) &= (1-f) \left(\frac{2N}{3} \left(1 - \frac{1}{2N}\right)^{T_{P3}-T_{P2}} \right) \\ &\quad + (f) \left(\left(\frac{2N}{3} \left(1 - \frac{1}{2N}\right)^{T_{P3}-T_{GF}} \right) + (T_{P3} - T_{GF}) \right) \end{aligned} \quad (25)$$

$$\mathbb{E}(ABBA_{sites}) = \mathbb{E}(\tau_{ABBA}) \times \mu \times L \quad (26)$$

Where μ represents the mutation rate and L represents the sequence length.

$$\mathbb{E}(BABA)$$

Derivation

$$\mathbb{E}(\tau_{BABA}) = C_{BABA_1} + C_{BABA_2} \quad (27)$$

$$\mathbb{E}(BABA)$$

Derivation

$$\mathbb{E}(\tau_{BABA}) = C_{BABA_1} + C_{BABA_2} \quad (27)$$

$$\begin{aligned} \mathbb{E}(\tau_{BABA}) &= (1-f) \left(\frac{2N}{3} \left(1 - \frac{1}{2N}\right)^{T_{P3}-T_{P2}} \right) \\ &\quad + (f) \left(\frac{2N}{3} \left(1 - \frac{1}{2N}\right)^{T_{P3}-T_{GF}} \right) \end{aligned} \quad (28)$$

$\mathbb{E}(BABA)$

Derivation

$$\mathbb{E}(\tau_{BABA}) = C_{BABA_1} + C_{BABA_2} \quad (27)$$

$$\begin{aligned} \mathbb{E}(\tau_{BABA}) &= (1 - f) \left(\frac{2N}{3} \left(1 - \frac{1}{2N}\right)^{T_{P3} - T_{P2}} \right) \\ &\quad + (f) \left(\frac{2N}{3} \left(1 - \frac{1}{2N}\right)^{T_{P3} - T_{GF}} \right) \end{aligned} \quad (28)$$

$$\mathbb{E}(BABA_{sites}) = \mathbb{E}(\tau_{BABA}) \times \mu \times L \quad (29)$$

Where μ represents the mutation rate and L represents the sequence length.

Patterson's D

Derivation

$$\mathbb{E}(D) = \frac{\mathbb{E}(\tau_{ABBA}) - \mathbb{E}(\tau_{BABA})}{\mathbb{E}(\tau_{ABBA}) + \mathbb{E}(\tau_{BABA})} \quad (30)$$

Patterson's D

Derivation

$$\mathbb{E}(D) = \frac{\mathbb{E}(\tau_{ABBA}) - \mathbb{E}(\tau_{BABA})}{\mathbb{E}(\tau_{ABBA}) + \mathbb{E}(\tau_{BABA})} \quad (30)$$

$$\mathbb{E}(D) = \frac{(f)(T_{P3} - T_{GF})}{(1-f) \left[{}^{4N/3} (1 - {}^{1/2}N)^{T_{P3} - T_{P2}} \right] + (f) \left[\left({}^{4N/3} (1 - {}^{1/2}N)^{T_{P3} - T_{GF}} \right) + (T_{P3} - T_{GF}) \right]} \quad (31)$$

Patterson's D

Derivation

$$\mathbb{E}(D) = \frac{\mathbb{E}(\tau_{ABBA}) - \mathbb{E}(\tau_{BABA})}{\mathbb{E}(\tau_{ABBA}) + \mathbb{E}(\tau_{BABA})} \quad (30)$$

$$\mathbb{E}(D) = \frac{(f)(T_{P3} - T_{GF})}{(1-f) \left[\frac{4N}{3} (1 - \frac{1}{2N})^{T_{P3} - T_{P2}} \right] + (f) \left[\left(\frac{4N}{3} (1 - \frac{1}{2N})^{T_{P3} - T_{GF}} \right) + (T_{P3} - T_{GF}) \right]} \quad (31)$$

$$\mathbb{E}(D) = \frac{\sum_{i=1}^L (1 - p_{i1}) p_{i2} p_{i3} (1 - p_{iO}) - p_{i1} (1 - p_{i2}) p_{i3} (1 - p_{iO})}{\sum_{i=1}^L (1 - p_{i1}) p_{i2} p_{i3} (1 - p_{iO}) + p_{i1} (1 - p_{i2}) p_{i3} (1 - p_{iO})} \quad (32)$$

Where $p_{i\#}$ represents the derived allele frequency at site i and L represents the sequence length.