

Genetic diversity stats

segregating sites, S , equals # mutations in the sample's history (infinite sites approximation) :

$$E[S] = 4N\mu \sum_{i=1}^{n-1} \frac{1}{i}$$

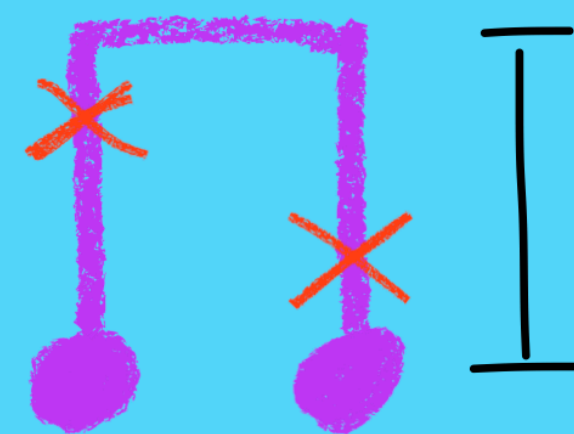
Pairwise divergence, π , # mutations in the history of two sampled haplotypes :

$$E[\pi] = 4N\mu$$

Tajima's D : null hypothesis, standard neutral coalescent (constant N)

$$\frac{E[S]}{\sum_{i=1}^{n-1} \frac{1}{i}} - E[\pi] = 0$$

Derivation:

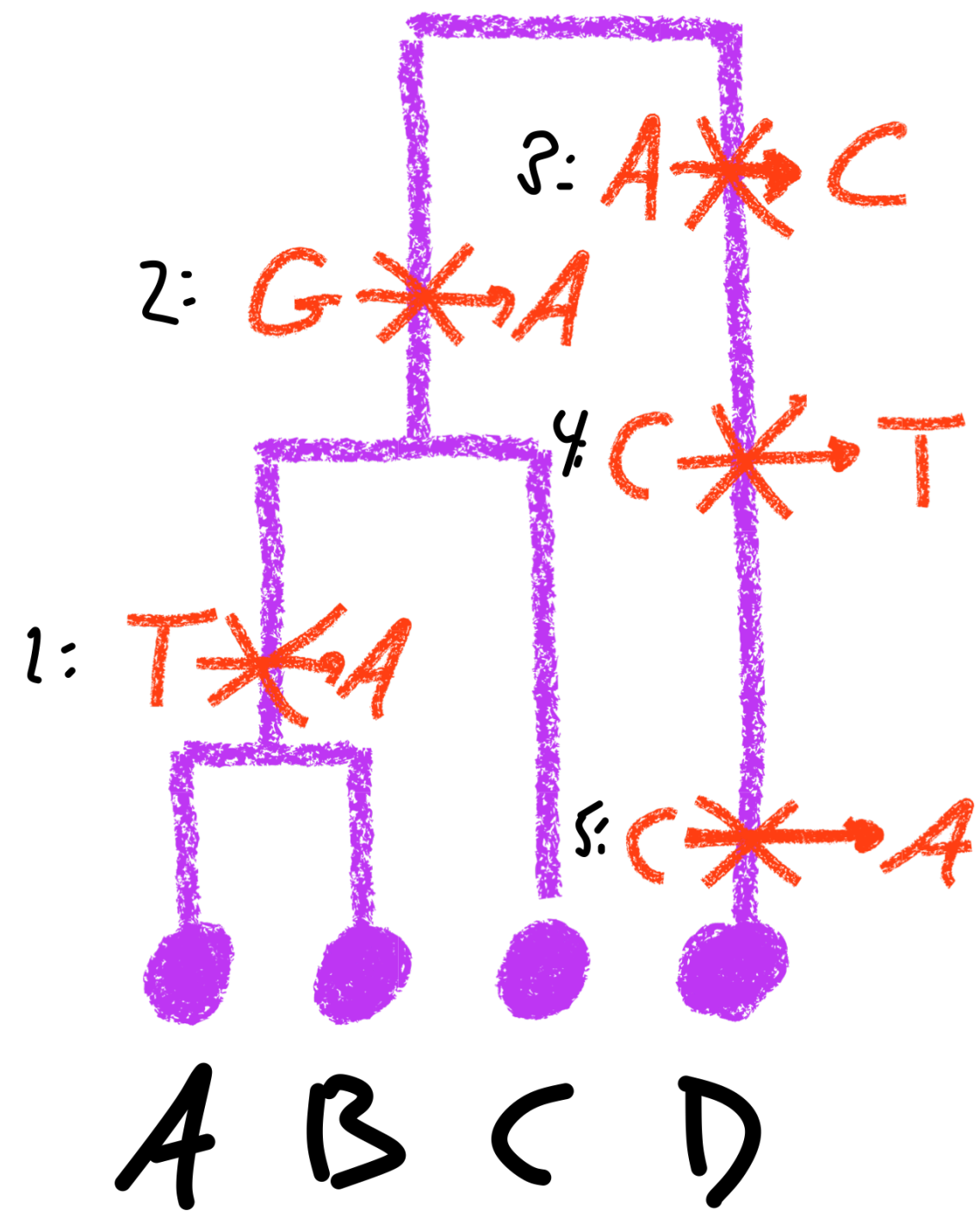


$$E[T_2] = \frac{1}{1/2N} = 2N$$

$$\rightarrow E[\pi] = 4N\mu$$

Question: why 4?

Example



| | 1 | 2 | 3 | 4 | 5 |
|---|---|---|---|---|---|
| A | A | A | A | C | C |
| B | A | A | A | C | C |
| C | T | A | A | C | C |
| D | T | G | C | T | A |

$$S = 5$$

$$\pi = \frac{1}{\binom{4}{2}} (0 + 1 + 5 + 1 + 5 + 4) = \frac{8}{3}$$

$$D = \frac{5}{1 + \frac{1}{2} + \frac{2}{3}} - \frac{8}{3} = \dots \frac{2}{33}$$