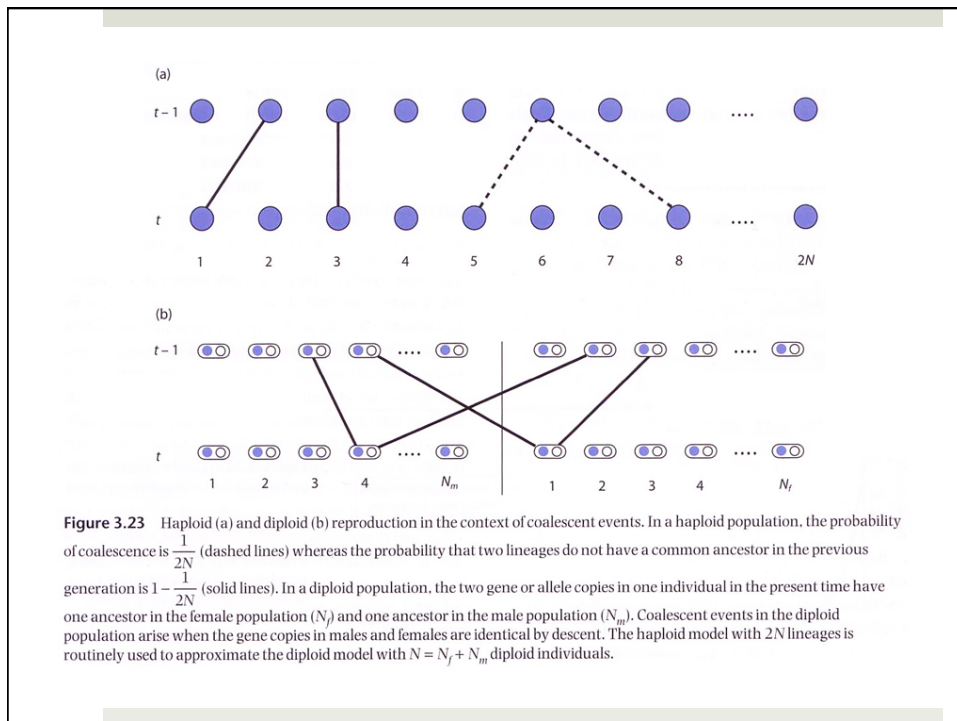
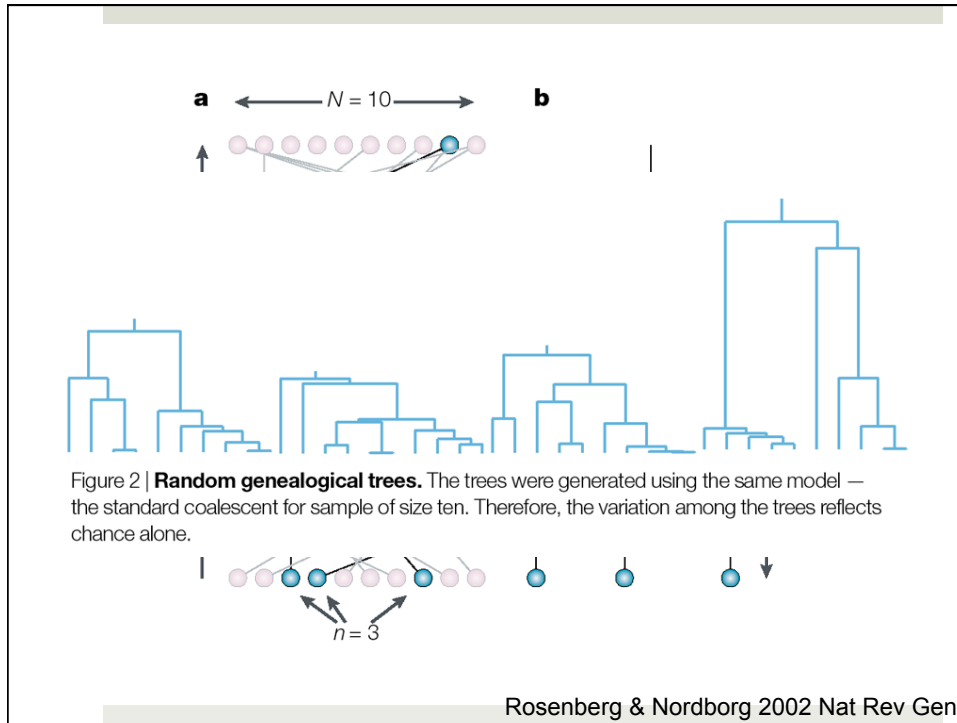


## Coalescent Theory

- ❖ the Wright-Fisher model considers changes in the ideal population as time moves forward
- ❖ coalescent theory (~1980+) looks backwards in time
- ❖ how long does it take for  $k$  alleles to coalesce to  $k - 1$  alleles, then  $k - 2$ ,  $k - 3$ , ..., and finally a single ancestral allele?

## Stochastic elements of the coalescent

- ❖ alleles randomly sample their parents in the previous generation
  - ❖ results in binomial variation in offspring number
- ❖ sample of loci from the genome
  - ❖ different loci have different genealogical histories
- ❖ sample of alleles from the population
  - ❖ different samples of the same locus may have different coalescent trees
- ❖ distribution of mutations on the genealogy
  - ❖ mutations allow estimates of coalescence times



## Coalescent probabilities 1

- ❖ the present is time 0 (zero)
- ❖ probability that *two* alleles had a common ancestor in generation 1

$$= \frac{1}{2N}$$

- ❖ probability that *two* alleles **did not have** a common ancestor in generation 1

$$= \left(1 - \frac{1}{2N}\right)$$

## Coalescent probabilities 2

- ❖ probability that *two* alleles have still not coalesced by generation  $t$

$$= \left[1 - \left(\frac{1}{2N}\right)\right]^t$$

## Coalescent probabilities 3

- ❖ probability that two alleles had a common ancestor in generation  $t+1$

$$= \frac{1}{2N} \left[ 1 - \left( \frac{1}{2N} \right) \right]^t$$

probability of  
coalescence in  
generation  $t + 1$

probability of NO  
coalescence in  
generations 1  
through  $t$

Error on pg. 93

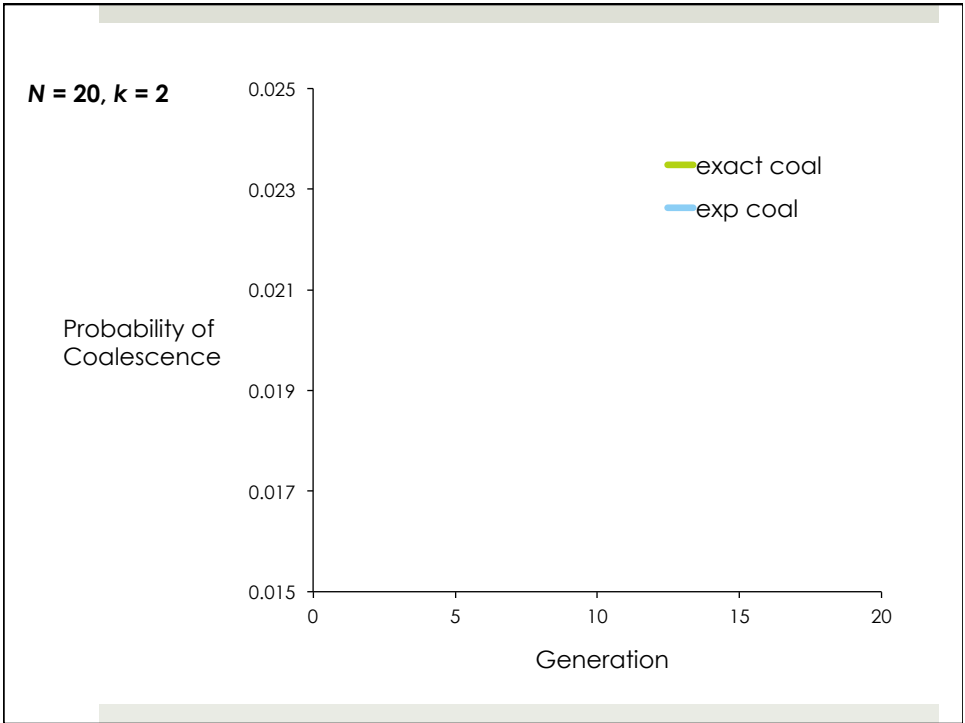
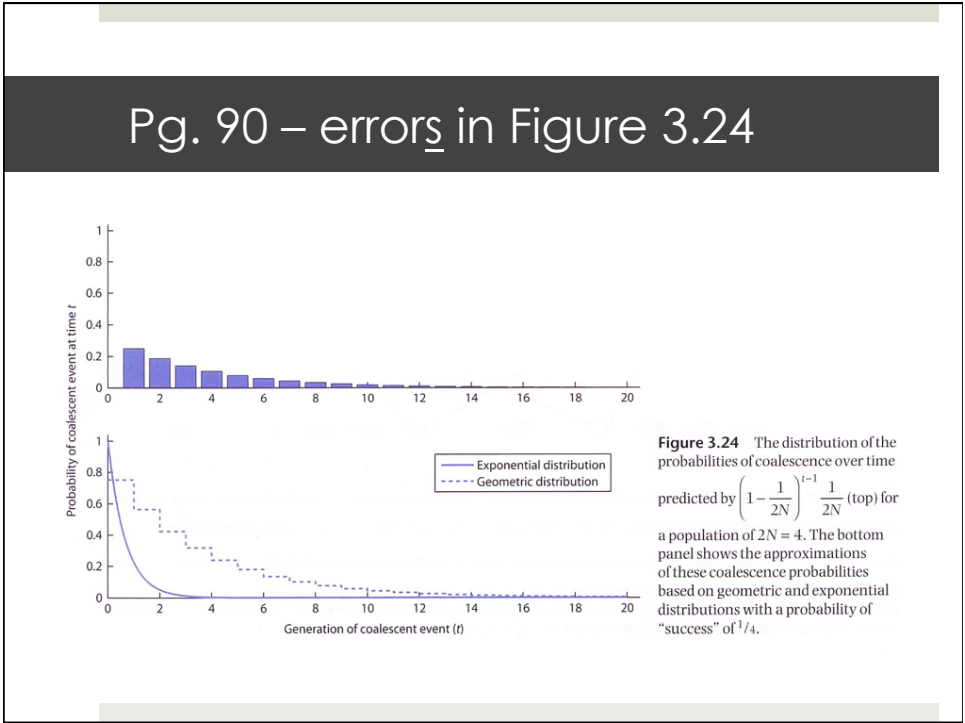
the exact probability of coalescence approximately equal to the exponential gives an expression

$$\left( 1 - \frac{1}{2N} \right)^{t-1} \frac{1}{2N} \approx \frac{1}{2N} e^{-\frac{t}{2N}} \quad (3.72)$$

that can be simplified by canceling the  $\frac{1}{2N}$  constant term on both sides

$$\left( 1 - \frac{1}{2N} \right)^{t-1} \approx e^{-\frac{t}{2N}} \quad (3.73)$$

Therefore, the exponential distribution approximates the probability of non-coalescence at each time  $t$ .



## Coalescent probabilities 4

- ❖ can we randomly choose a coalescence time from the exponential distribution?
- ❖ need to solve for  $t$  as a function of a random variable from 0 to 1

$$P_{NC} \approx e^{-t/(2N)}$$

$$\ln(P_{NC}) \approx -t/(2N)$$

$$\ln(P_{NC}) \times 2N \approx -t$$

$$t \approx -\ln(P_{NC}) \times 2N$$

## Coalescent probabilities 5

- ❖ what if we consider  $k$  alleles and not just 2?
- ❖ probability that  $k$  alleles had  $k$  distinct parental alleles the previous generation

$$\Pr(k) = \prod_{i=0}^{k-1} \left(1 - \frac{i}{2N}\right) \approx \left(1 - \frac{\binom{k}{2}}{2N}\right)$$

## Coalescent probabilities 6

- ❖ probability that  $k$  alleles do not coalesce for  $t$  generations

$$P_{NC} = \left(1 - \frac{\binom{k}{2}}{2N}\right)^t \approx e^{\left[-\frac{\binom{k}{2}}{2N}t\right]}$$

$$t \approx -\ln(P_{NC}) \times \frac{2N}{\binom{k}{2}} = -\ln(P_{NC}) \times \frac{4N}{k(k-1)}$$

## Coalescent probabilities 7

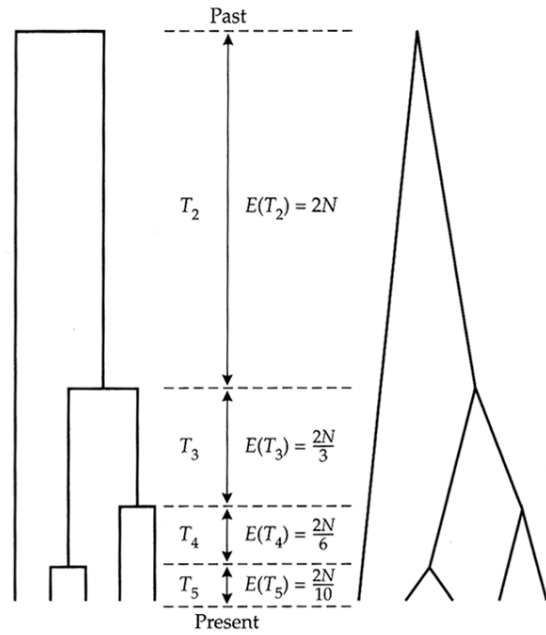
- ❖ probability that  $k$  alleles do not coalesce for  $t$  generations, and then one pair coalesces to give  $k - 1$  alleles at  $t + 1$  generations

$$= \Pr(k)^t [1 - \Pr(k)] \approx \frac{\binom{k}{2}}{2N} e^{\left[-\frac{\binom{k}{2}}{2N}t\right]}$$

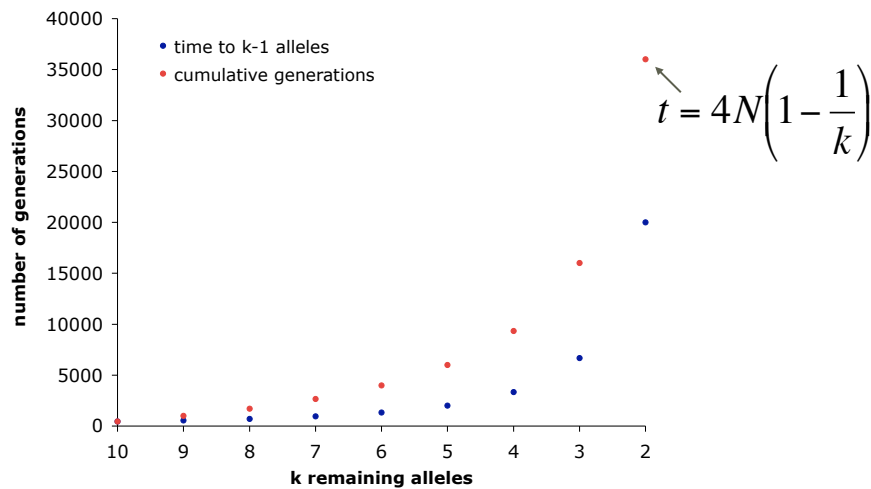
- ❖ distribution has mean and variance:

$$\text{Mean} = \frac{4N}{k(k-1)} \text{generations} \quad \text{Var} = \frac{16N^2}{[k(k-1)]^2} \text{generations}^2$$

**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 left 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.



expected coalescence times for  $k = 10$  and  $N = 10,000$



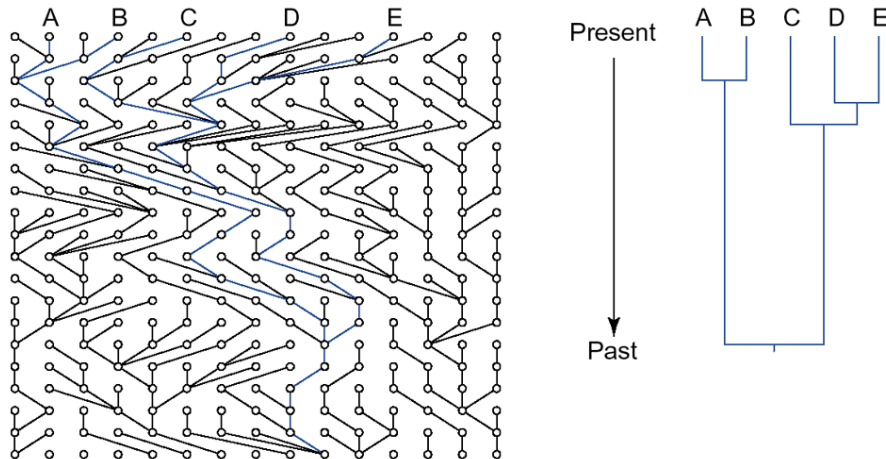


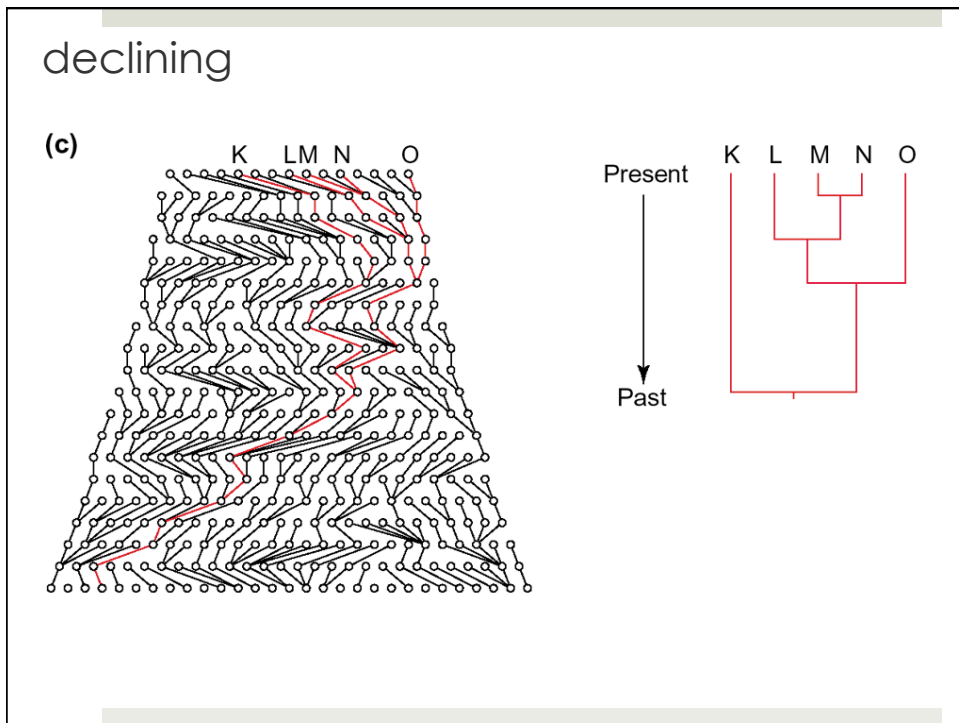
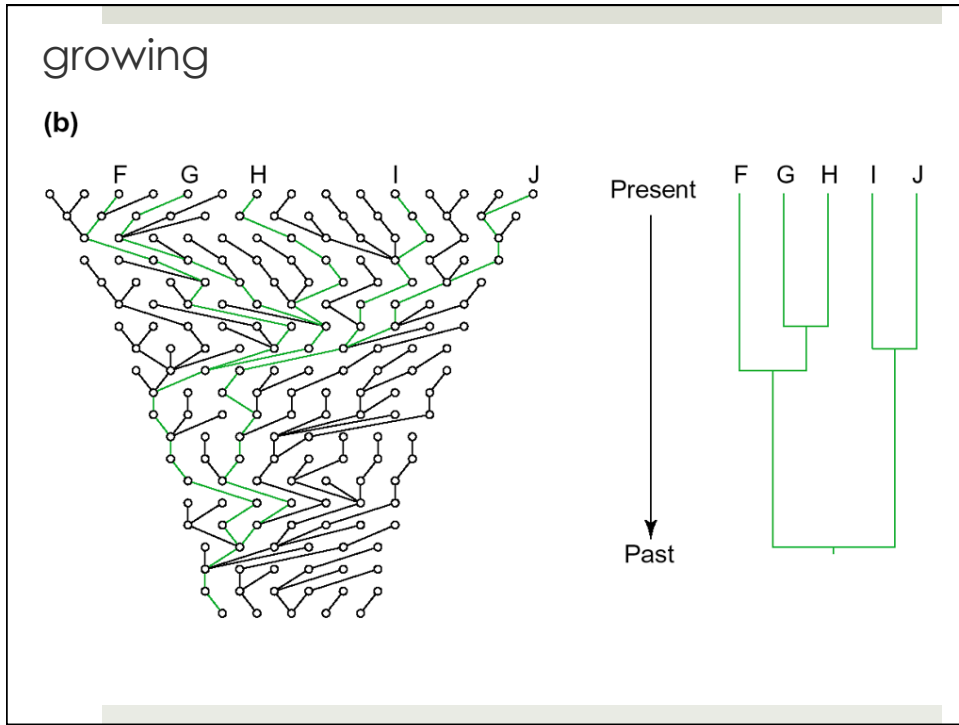
## The coalescent with population growth

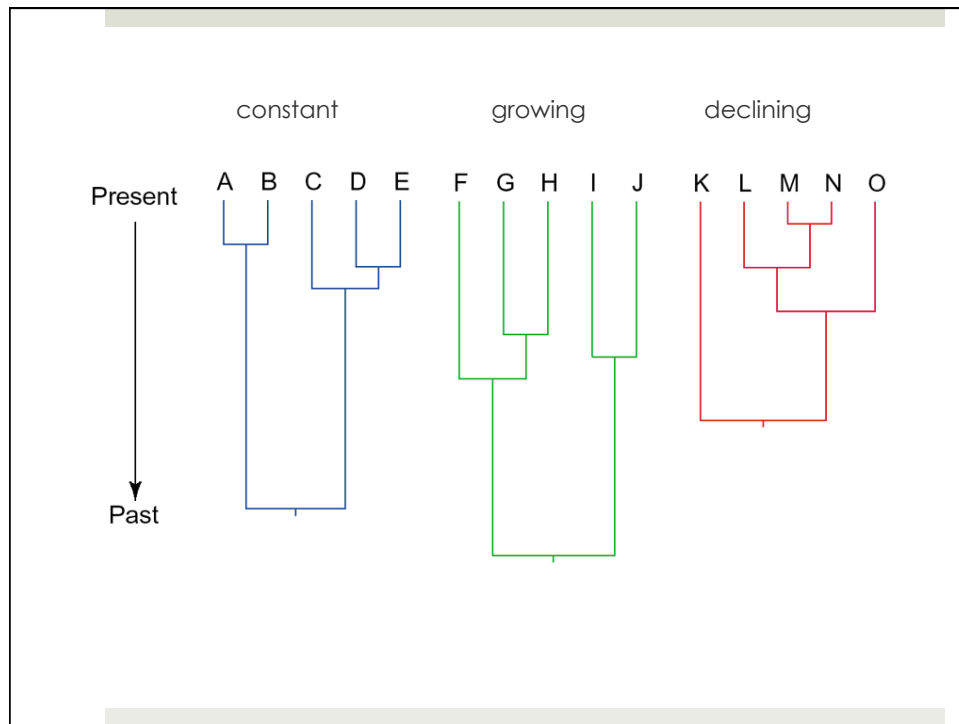
- ❖ coalescent trees are expected to be sparse (few lineages) near the root for populations of constant size
- ❖ in a growing or shrinking population, the distribution of coalescence times differs from expectations for the ideal population
- ❖ expanding populations have more nodes closer to the root of the tree
  - ❖ takes longer for alleles to “find each other” in a growing population

constant

(a)







## So, what's the point?

- ❖ coalescent modeling
  - ❖ simulate genealogies under a given set of population parameters
  - ❖ “hang” random mutations on those trees in equal number (or at the same rate) as in the observed data
  - ❖ evaluate whether the observed data could have been produced by a random coalescent process (the null hypothesis)

