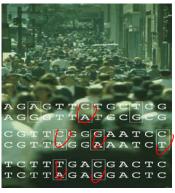


Genetic Demography







- Are there genetic prototypes among them?
- What are they?
- How many ? (how many ancestors do we have ?)

Multi-population Genetic Demography





• Inference done separately, or jointly?





The coalescent

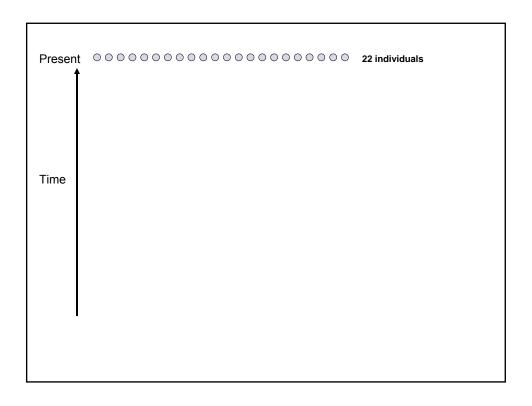


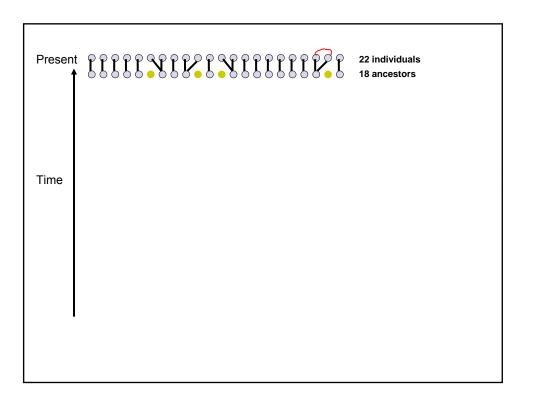
Sir John Kingman, Head of the Isaac Newton Institute of Mathematical Sciences

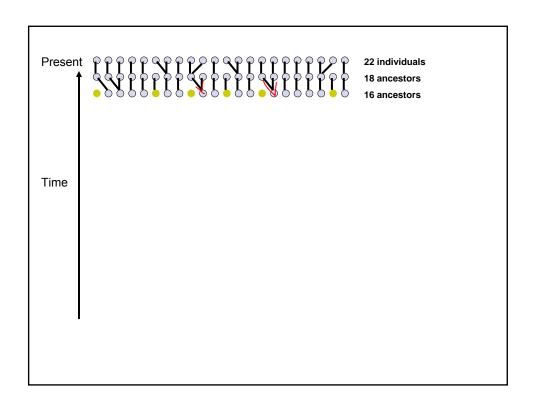
Coalescent Theory

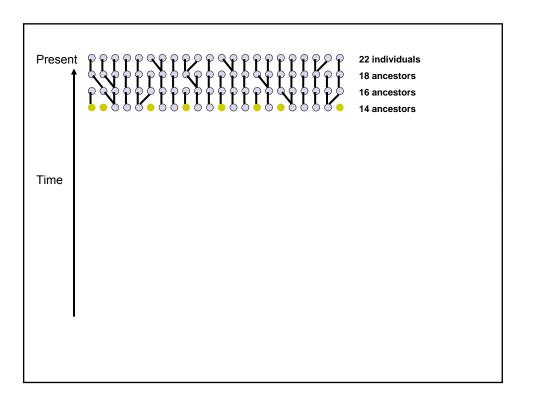


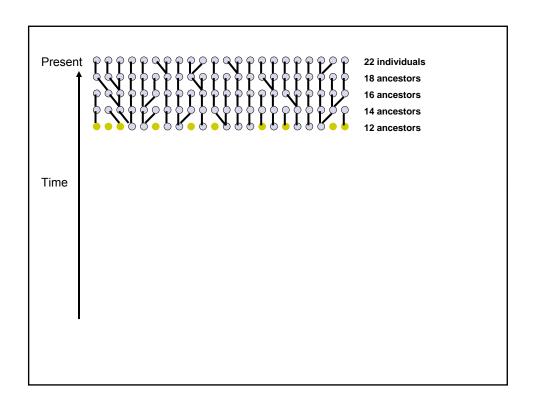
- how we can build up a genealogical tree to relate a sample of <u>n</u> haploid individuals, collected in the present day?
 - The following series of slides shows how you can build up a
 genealogical tree to relate a sample of 22 individuals, collected in the
 present day, at a single haplotype locus (e.g. the non-recombining Y
 chromosome).
 - Because (for the Y chromosome) one son has only one father, but one
 father can have more than one son, coalescent events occur in the
 genealogy which inevitably result in a reduction of ancestors. Eventually,
 one ancestor remains the Most Recent Common Ancestor (MRCA).

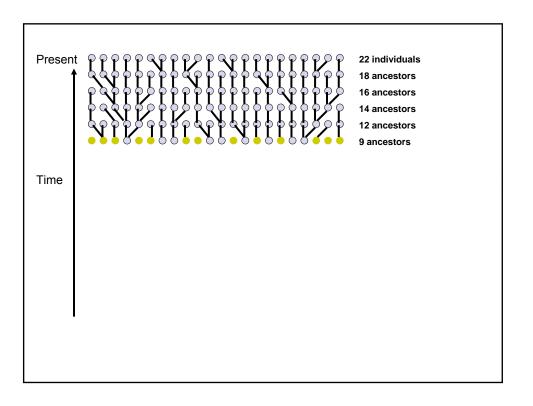


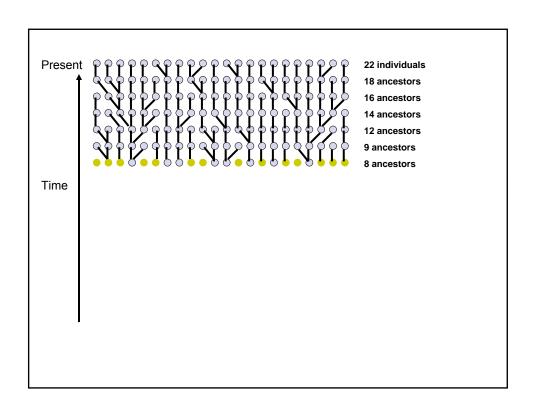


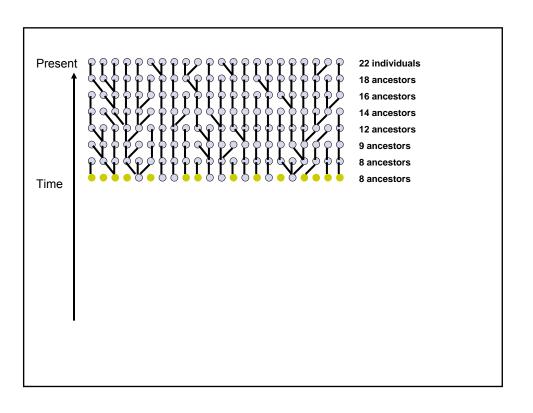


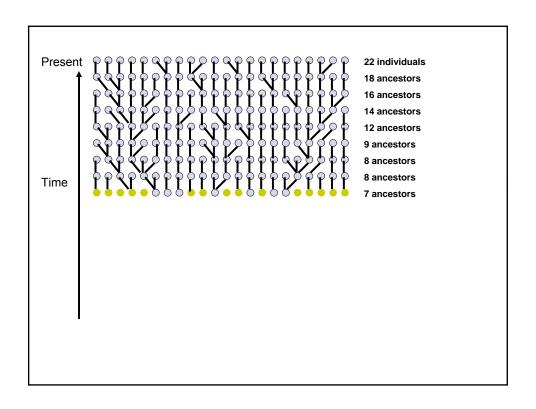


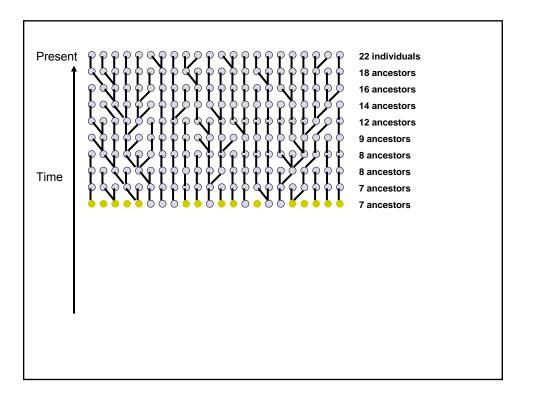


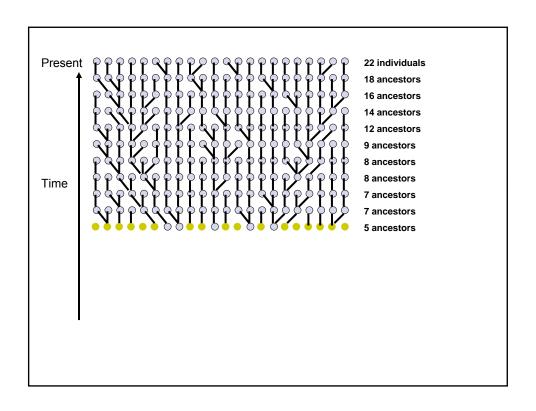


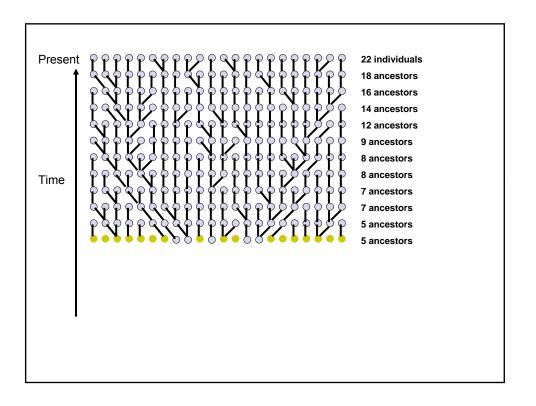


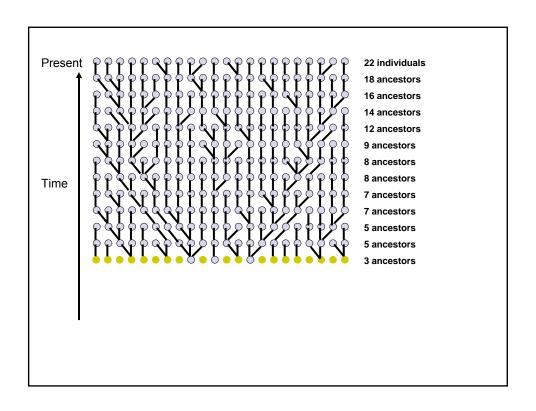


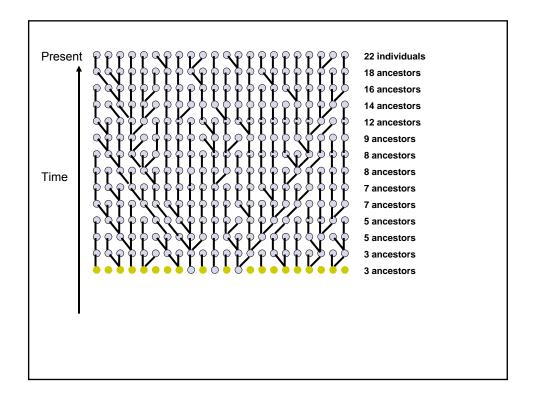


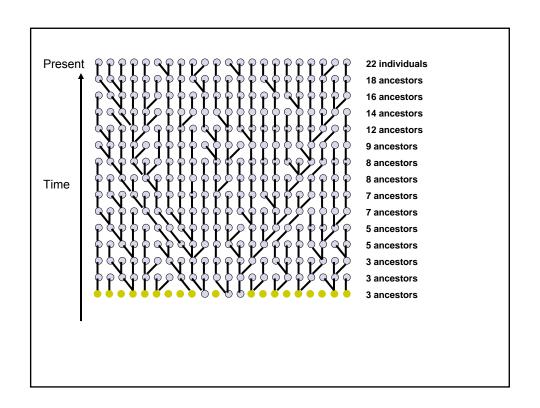


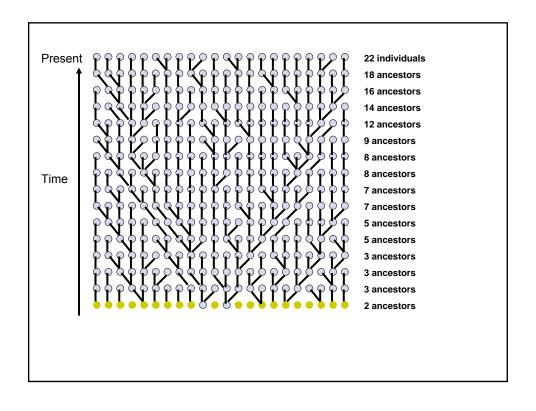


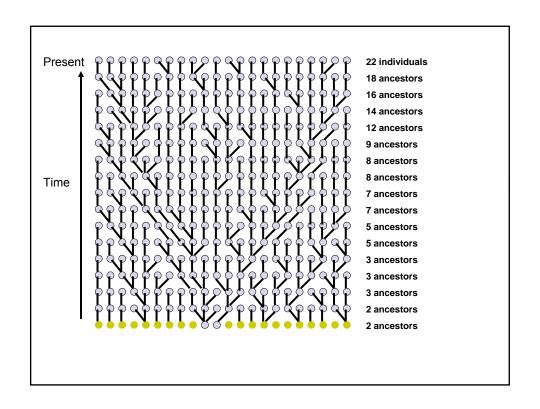


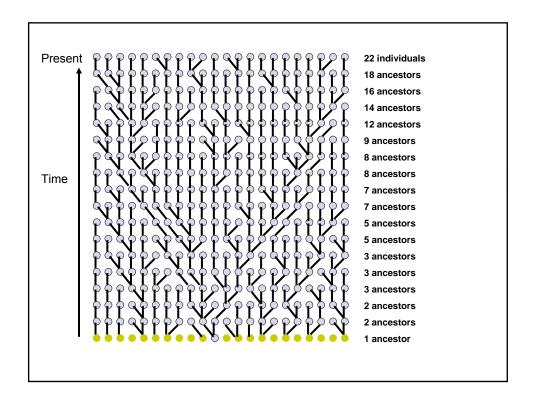


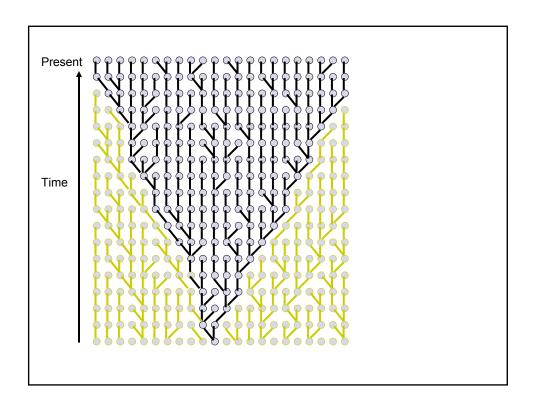


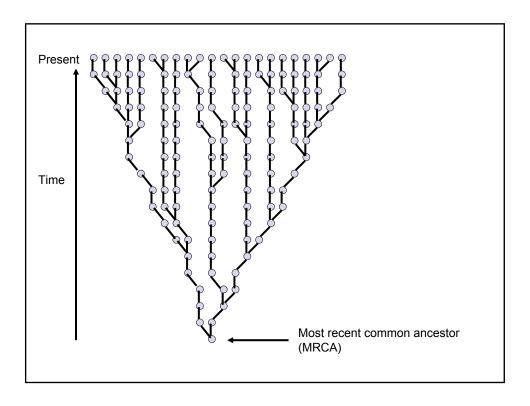






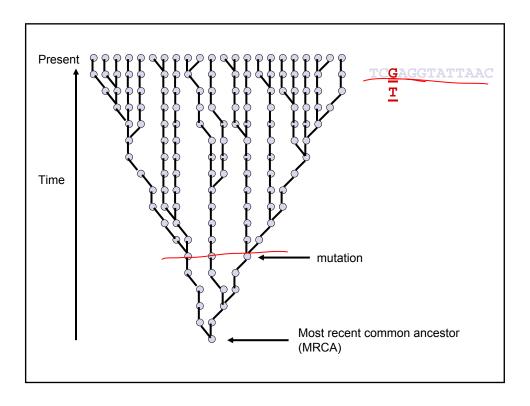


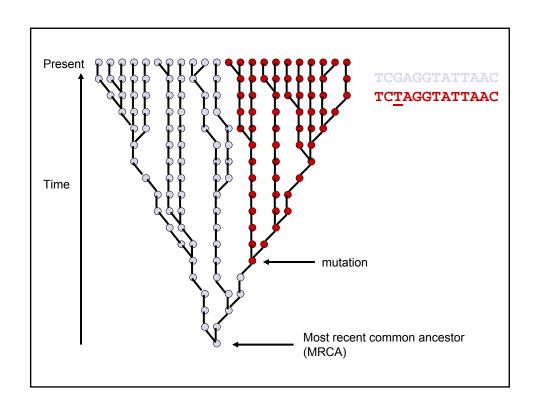


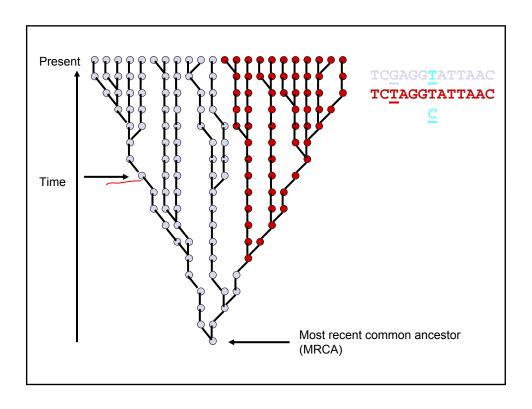


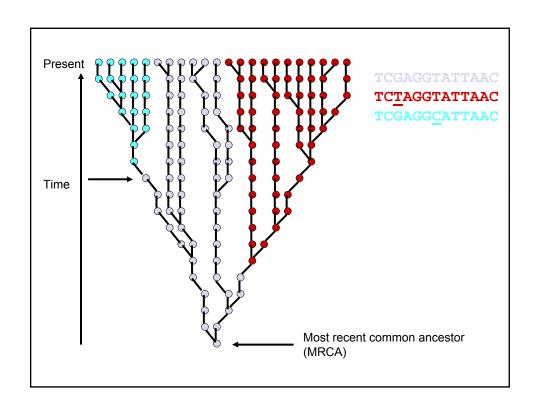


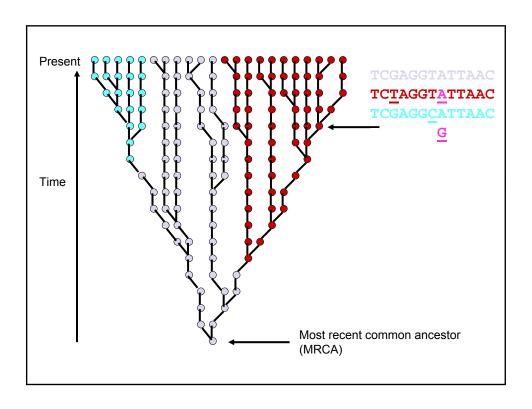
 Mutational events can now be added to the genealogical tree, resulting in polymorphic sites. If these sites are typed in the modern sample, they can be used to split the sample into sub-clades (represented by different colors)

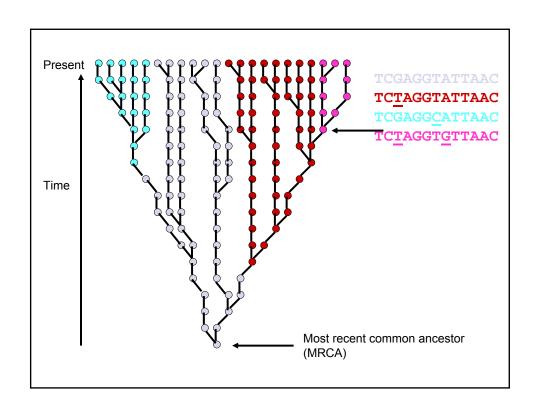


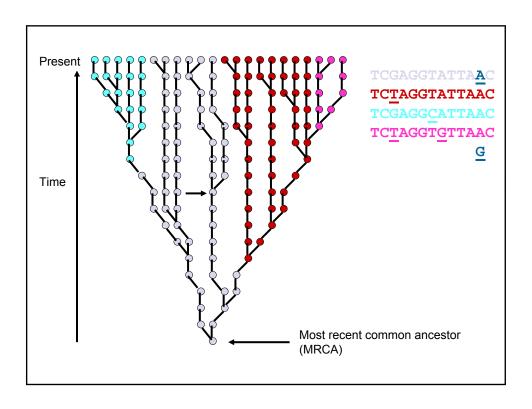


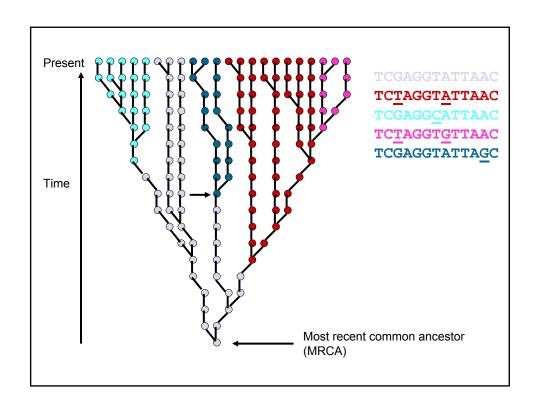


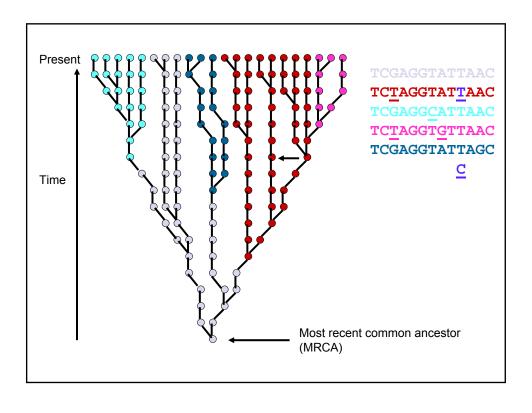


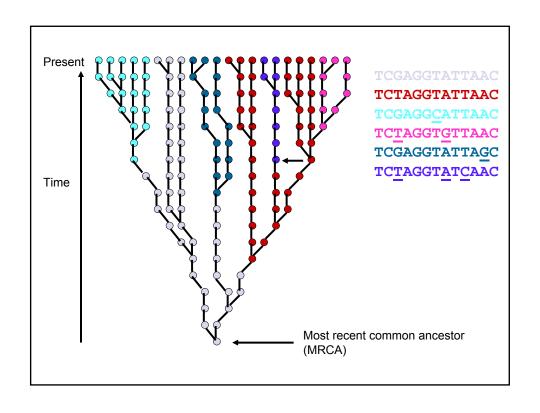


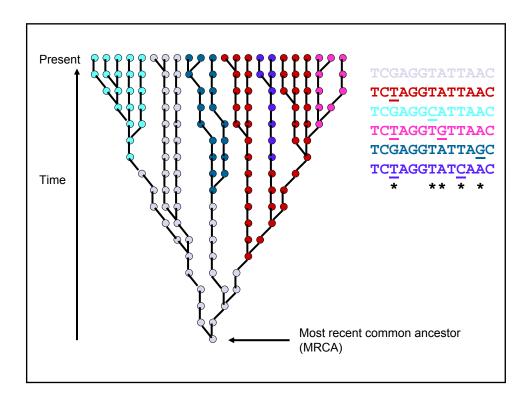












The Statistical Models



- To move beyond mere description, and to attempt such things as estimating the TMRCA (Time to Most Recent Common Ancestor) of the tree, it is necessary to adopt certain modeling assumptions.
- For now lets forget about mutations, but just concern ourselves with the coalescence

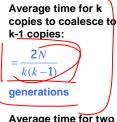
Kingman's coalescent process



- Random collision of lineages as go back in time
- Collision is faster the smaller the effective population size
 - In a haplotype population of effective population size

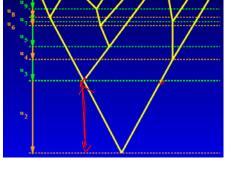
Average time for proper to coalesce: $= 2N\left(1 - \frac{1}{n}\right)$

generations



Average time for two copies to coalesce:

Support of two coalesce:



Derivation? ---- Hw!

Hint of the derivation





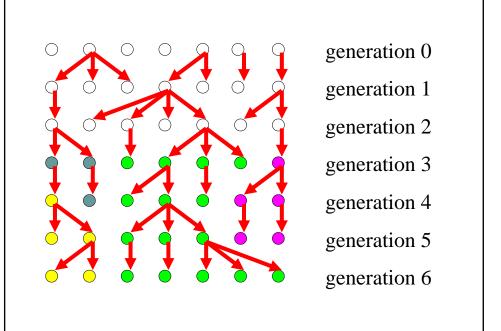
The Wright-Fisher (WF) model



- The coalescent is descriptive, but not generative!
- A classic generative model is the <u>Wright-Fisher model</u>. This
 is the canonical model of genetic drift in populations. It was
 invented in 1932 and 1930 by Sewall Wright and R. A. Fisher.
- It starts with the following assumptions:
 - <u>random mating</u> and a <u>random number of offspring</u> (strictly, following a Poisson distribution)
 - no recombination (i.e. a single locus),
 - constant population size,
 - no selection,

The Wright-Fisher (WF) model

- It is a forwards-in-time model of a neutral locus in a constant-size, random-mating, haploid population evolving in discrete generations.
- Each individual in generation t has a random number (possibly 0) of offspring in generation t+1. Each is:
 - identical to the parent with probability 1-μ;
 - otherwise a mutation occurs.
- With WF, one can attempt such things as estimating the TMRCA (Time to Most Recent Common Ancestor) of the tree, etc.

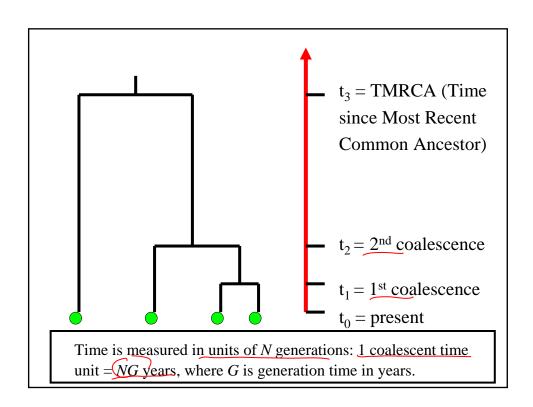


Coalescent theory



When we consider the same set of assumptions but now simulate going "backwards in time", we arrive at the standard coalescent model with *infinite-allele-mutations*.

- A coalescent is the backwards-in-time "cousin" of the WF model: similar assumptions, but traces the ancestry of n observed alleles.
- Ancestry is represented via a genealogical tree: leaves are observed alleles, root is the most recent common ancestor (MRCA).



Time back to the next coalescence when there are k lineages has the exponential distribution with mean and standard deviation both 2/k(k-1);

e.g. k=4:

mean = sd = 1/3

mean = sd = 1/6

The TMRCA under the coalescent

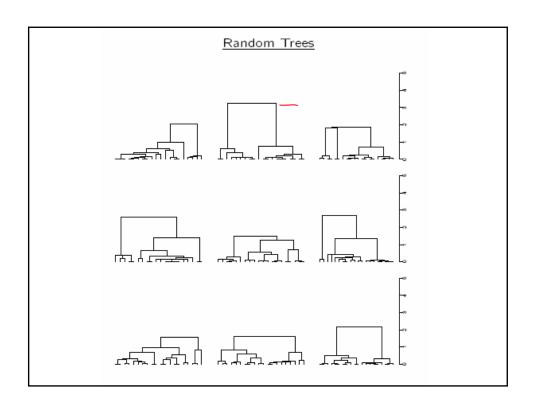
Height of tree:

Total branch length:

mean = 3/2 sd = 1.07

mean = 11/3 sd = 2.33

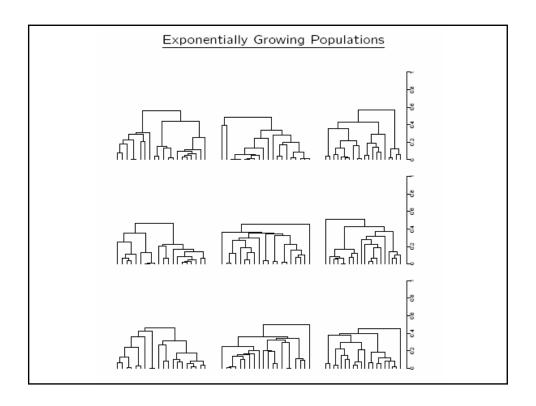
- The TMRCA (height of the genealogical tree) is on average 2(n-1)/n; the average time in which there are just two ancestral lineages is 1.
 - the number of ancestors of a sample drops rapidly (backwards in time);
 - for more than half its history, on average, a sample has only two ancestors:
 - data often clustered.
- When we simulate from the standard coalescent, we find that there
 is considerable variation in the TMRCA from one simulation to the
 next.
 - Most coalescent event occur in the recent past (at the tips of the tree)



Coalescent with variable population size



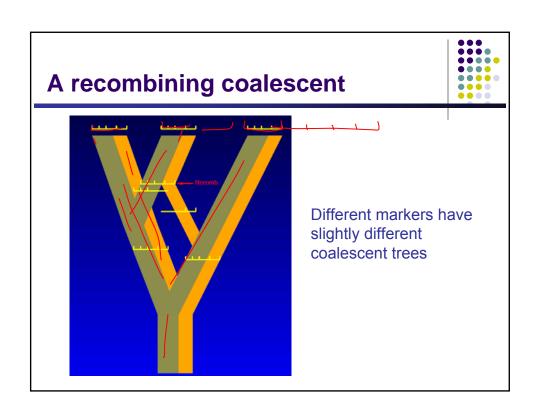
- The situation changes if we expand the coalescent model to incorporate a factor of **exponential population growth**.
- Now there is less variation in the TMRCA between simulations, and more coalescent events occur in the more distant past (near the root of the tree).

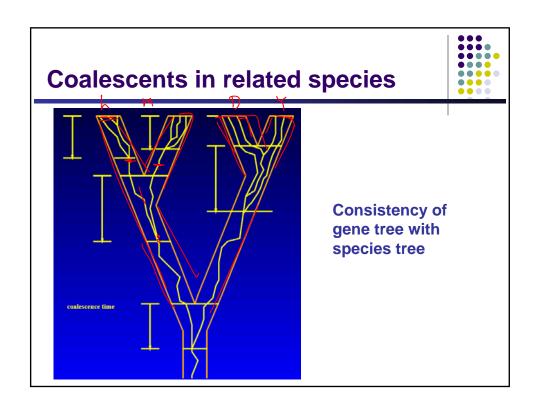


Generalisations of the standard coalescent model



- Variable population size: coalescences occur more rapidly when the population size is small.
- Population subdivision with migration.
- Some forms of selection.
- Recombination: the ancestral recombination graph (ARG)

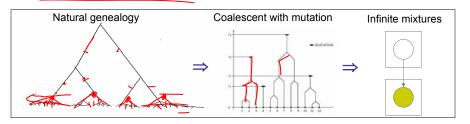


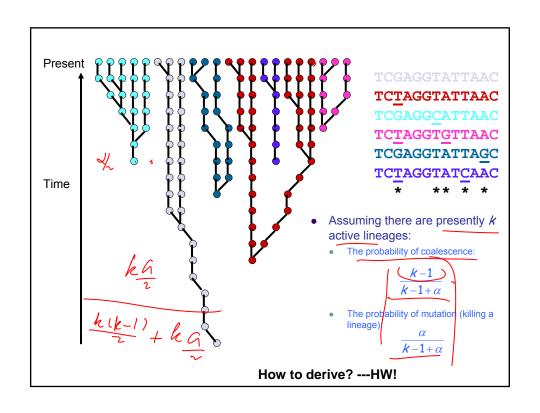


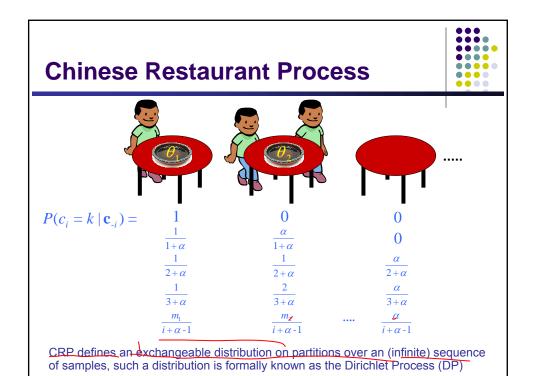
How to approximate a coalescent?



- Kingman coalescent process with binary lineage merging
- New population haplotype alleles emerge along all branches of the coalescence tree at rate a/2 per unit length
- This can be approximated by an infinite mixture model (aks, Dirichlet process mixture)







The DP Mixture of Ancestral Haplotypes



- The customers around a table form a cluster
 - associate a mixture component (i.e., a population haplotype) with a table
 - sample $\{a, \theta\}$ at each table from a base measure G_0 to obtain the population haplotype and nucleotide substitution frequency for that component









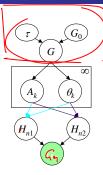




• With $p(h/\{A, \theta\})$ and $p(g/h_p h_p)$, the CRP yields a posterior distribution on the number of population haplotypes (and on the haplotype configurations and the nucleotide substitution frequencies)

A Hierarchical Bayesian Infinite Allele model





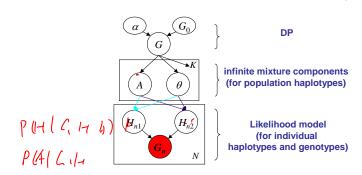
• Assume an individual haplotype h is stochastically derived from a population haplotype a_k with nucleotide-substitution frequency θ_k :

$$h \sim p(h|\{a,\theta\}_k).$$

- Not knowing the correspondences between individual and population haplotypes, each individual haplotype is a mixture of population haplotypes.
- The number and identity of the population haplotypes are unknown
 - use a Dirichlet Process to construct a prior distribution G on $\mathcal{H} \times \mathcal{R}'$.

DP-haplotyper





- Inference: Markov Chain Monte Carlo (MCMC)
 - Gibbs sampling
 - Metropolis Hasting

Model components



· Choice of base measure:

$$G_0 \sim \operatorname{Unif}(\boldsymbol{a}) \cdot \prod_j \operatorname{Beta}(\theta_j)$$

• Nucleotide-substitution model:

$$p(h_{i} | \{a, \theta\}_{k}) = \prod_{j} p(h_{i,j} | a_{k,j}, \theta_{k,j})$$
where
$$p(h_{i,j} | a_{k,j}, \theta_{k,j}) = \begin{cases} \theta_{k,j} & \text{if } h_{i,j} = a_{k,j} \\ 1 - \theta_{k,j} & \text{if } h_{i,j} = a_{k,j} \end{cases}$$

• Noisy genotyping model:

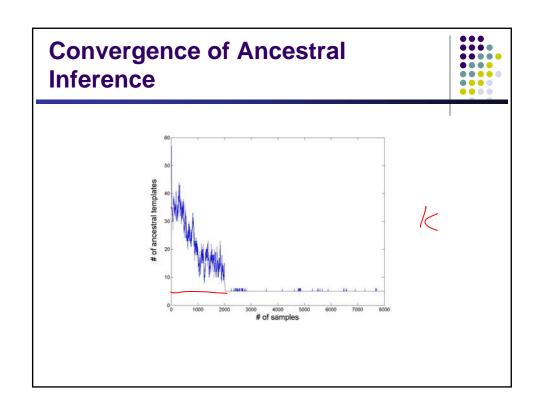
$$p(g_{i} | h_{i_{1}}, h_{i_{2}}) = \prod_{j} p(g_{i,j} | h_{i_{1},j}, h_{i_{2},j})$$
where
$$p(g_{i,j} | h_{i_{1},j}, h_{i_{2},j}) = \begin{cases} \gamma & \text{if } h_{i_{1},j} \oplus h_{i_{2},j} = g_{i,j} \\ \frac{1-\gamma}{2} & \text{if } h_{i_{1},j} \oplus h_{i_{2},j} \neq g_{i,j} \end{cases}$$

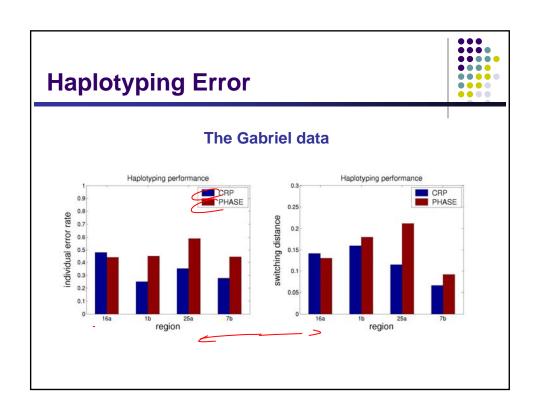
Gibbs sampling



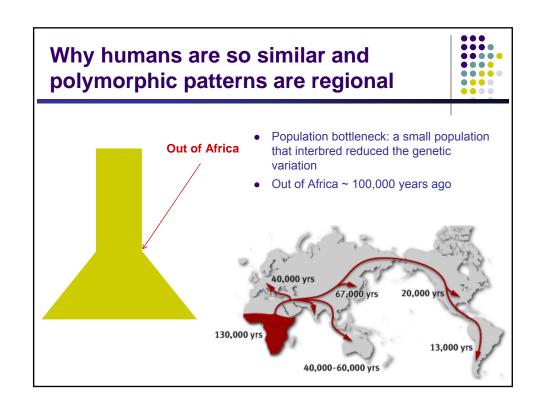
Starting from some initial haplotype reconstruction $H^{(0)}$, pick a first table with an arbitrary $a_I^{(0)}$, and form initial population-hap pool $\mathbf{A}^{(0)} = \{a_I^{(0)}\}$:

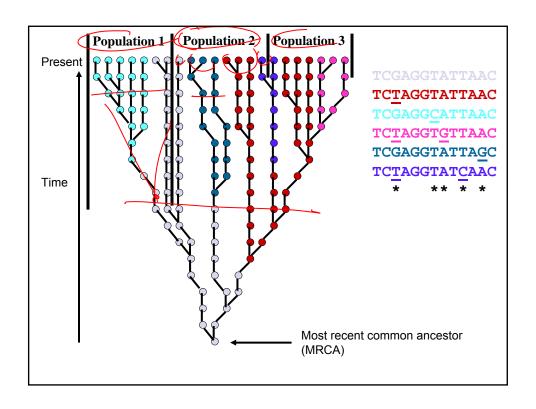
- i) Choose an individual *i* and one of his/her two haplytopes *t*, uniformly and at random, from all ambiguous individuals;
- ii) Sample $c_{i_t}^{(t+1)}$ from $p(c_{i_t}^{(t+1)} \mid c_{-i_t}^{(t)}, H^{(t)}, \mathbf{A}^{(t)})$, update $c^{(t+1)}$;
- iii) Sample $a_k^{(t+1)}$, where $k = c_{i_t}^{(t+1)}$, from $p(a_k^{(t+1)} \mid \forall h_{-i_{i_t}}^{(t)} \text{ s.t. } c_{i_{i_t}}^{(t+1)} = k)$; update $\mathbf{A}^{(t+1)}$;
- iii) Sample $h_{i_t}^{(t+1)}$ from $p(h_{i_t}^{(t+1)} \mid c_{i_t}^{(t+1)}, H_{-i_t}^{(t)}, \mathbf{A}^{(t+1)})$, update $H^{(t+1)}$.





Multi-population Genetic Gemography Pool everything together and solve 1 hap problem? - rignore population structures Solve 4 hap problems separately? - rdata fragmentation Co-clustering ... solve 4 coupled hap problems jointly



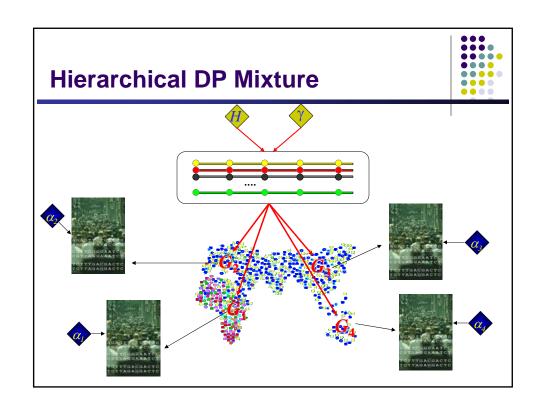


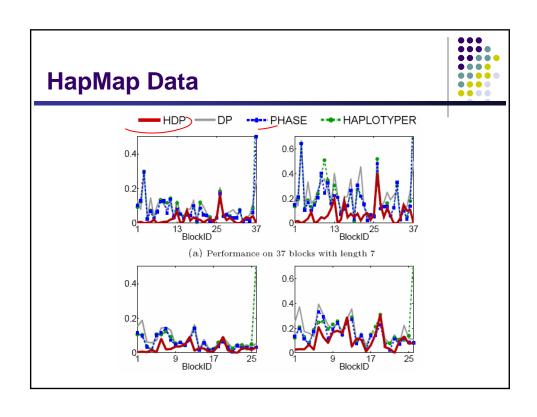
Population Specific DPs



- Each population can be associated with a unique DP capturing population-specific genetic demography
- Different population may have unique haplotypes
- Different population may share common haplotypes
- Thus Population specific DPs are marginally dependent



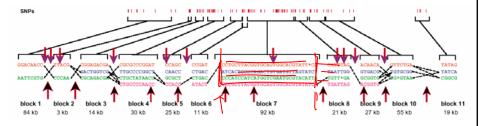




The block structure of haplotypes



- The Daly et al (2001) data set
 - This consists of 103 common SNPs (>5% minor allele frequency) in a 500 kb region implicated in Crohn disease, genotyped in 129 trios (mom, pop, kid) from a European derived population, giving 258 transmitted and 258 untransmitted chromosomes.

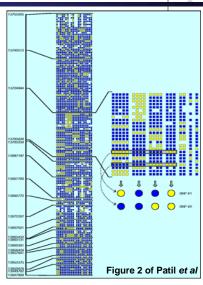


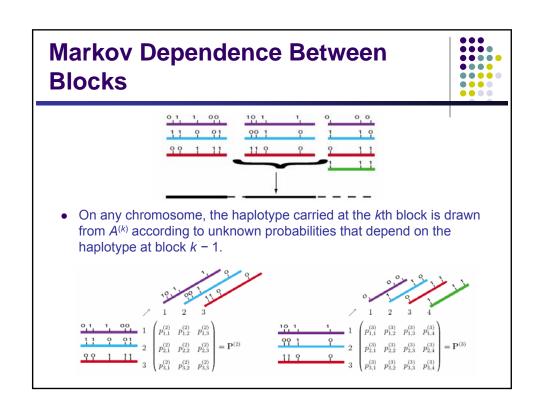
The haplotype blocks span up to 100kb and contain 5 or more common SNPs.
 For example, one 84 kb block of 8 SNPs shows just two distinct haplotypes accounting for 95% of the observed chromosomes.

Another study: the Patil et al data



- The haplotype patterns for 20 independent globally diverse chromosomes defined by 147 common human chr 21 SNPs spanning 106 kb of genomic sequence.
 - Each row represents an SNP.
 - Blue box = major, yellow = minor allele.
 - Each column represents a single chromosome.
- The 147 SNPs are divided into 18 blocks defined by black lines.
 - The expanded box on the right is an SNP block of 26 SNPs over 19kb of genomic DNA.
 - The 4 most common of 7 different haplotypes include 80% of the chromosomes, and can be distinguished with 2 SNPs.

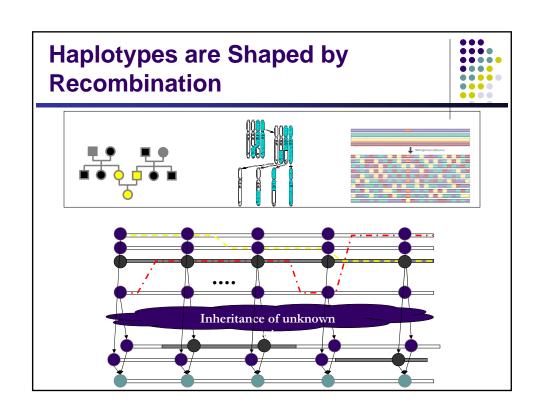




Block Partitioning Algorithms

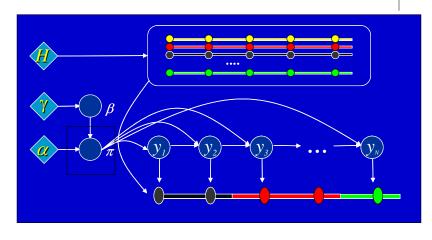


- Greedy algorithm (Patil et al 2001).
 - Begin by considering all possible blocks of ≥1 consecutive SNPs.
 - Next, exclude all blocks in which < 80% of the chromosomes in the data are defined by haplotypes represented more than once in the block (80% coverage).
 - Considering the remaining overlapping blocks simultaneously, select the one which maximizes the ratio of total SNPs in the block to the number required to uniquely discriminate haplotypes represented more than once in the block. Any of the remaining blocks that physically overlap with the selected block are discarded, and the process repeated until we have selected a set of contiguous, non-overlapping blocks that cover the 32.4 Mb of chr 21 2ith no gaps and with every SNP assigned to a block.
- Hidden Markov Models
 - Maximum a posterior inference (i.e., viterbi) (Daly et al. 2001)
 - Minimum description length (Anderson et al.2003)
- Dynamic programming (Sun et al. 2002)



Hidden Markov DP for Recombination





Reference



- E.P. Xing, R. Sharan and M.I Jordan, Bayesian Haplotype Inference via the Dirichlet Process. Proceedings of the 21st International Conference on Machine Learning (ICML2004),
- E. P. Xing and K. Sohn, Hidden Markov Dirichlet Process: Modeling Genetic Recombination in Open Ancestral Space, Journal of Bayesian Analysis, 2007
- E.P. Xing, K. Sohn, M.I. Jordan and Y.W. Teh, Bayesian Multi-Population Haplotype Inference via a Hierarchical Dirichlet Process Mixture, Proceedings of the 23st International Conference on Machine Learning (ICML 2006).
- N Patil et al. Blocks of limited haplotype diversity revealed by high-resolution scanning of human chromosome 21 Science 294 2001:1719-1723.
- M J Daly et al. High-resolution haplotype structure in the human genome Nat. Genet. 29 2001: 229-232
- Anderson, E.C., Novembre, J. (2003) "Finding haplotype block boundaries using the minimum description length principle." American Journal of Human Genetics 73(2):336-354.