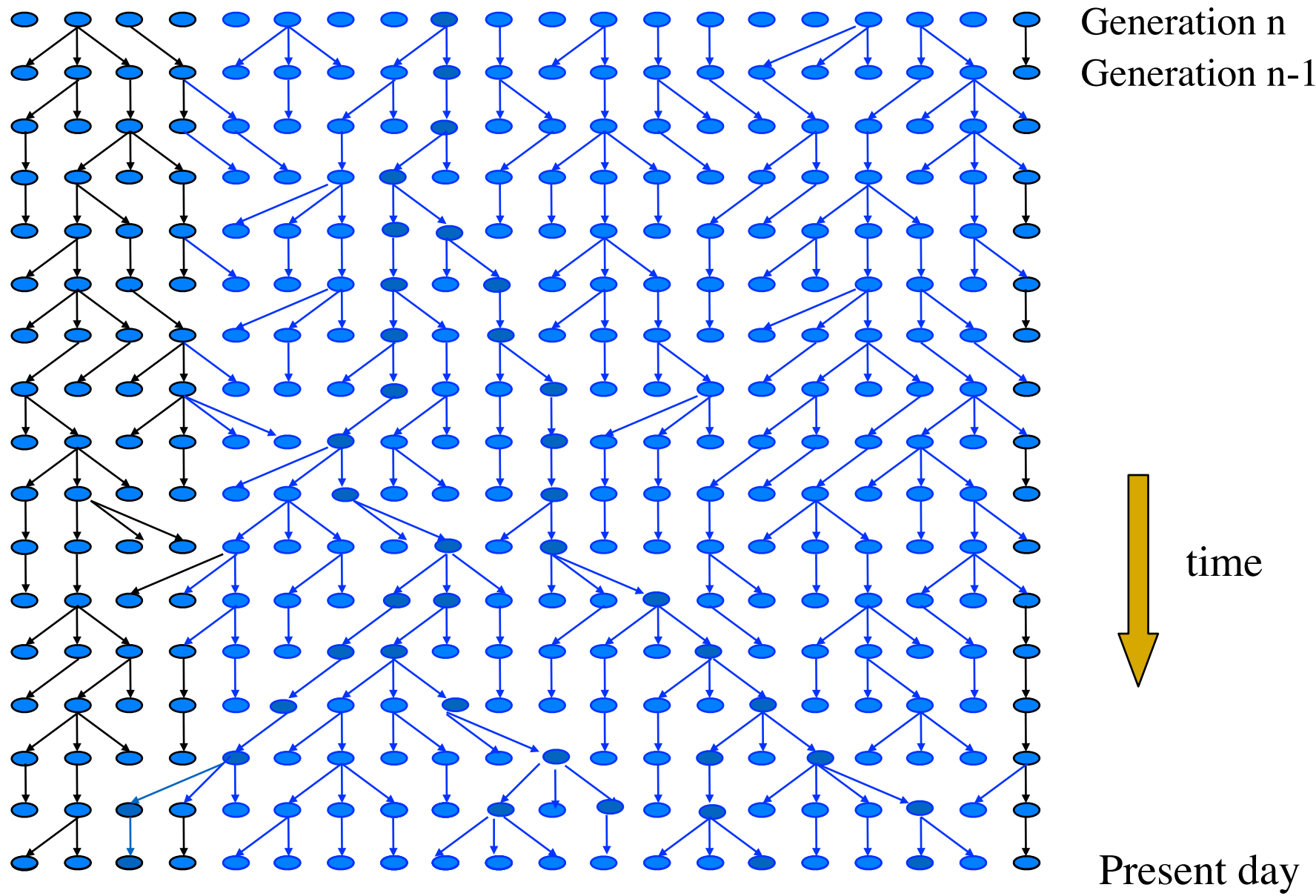


Population Genetics, Coalescent Trees and Urns

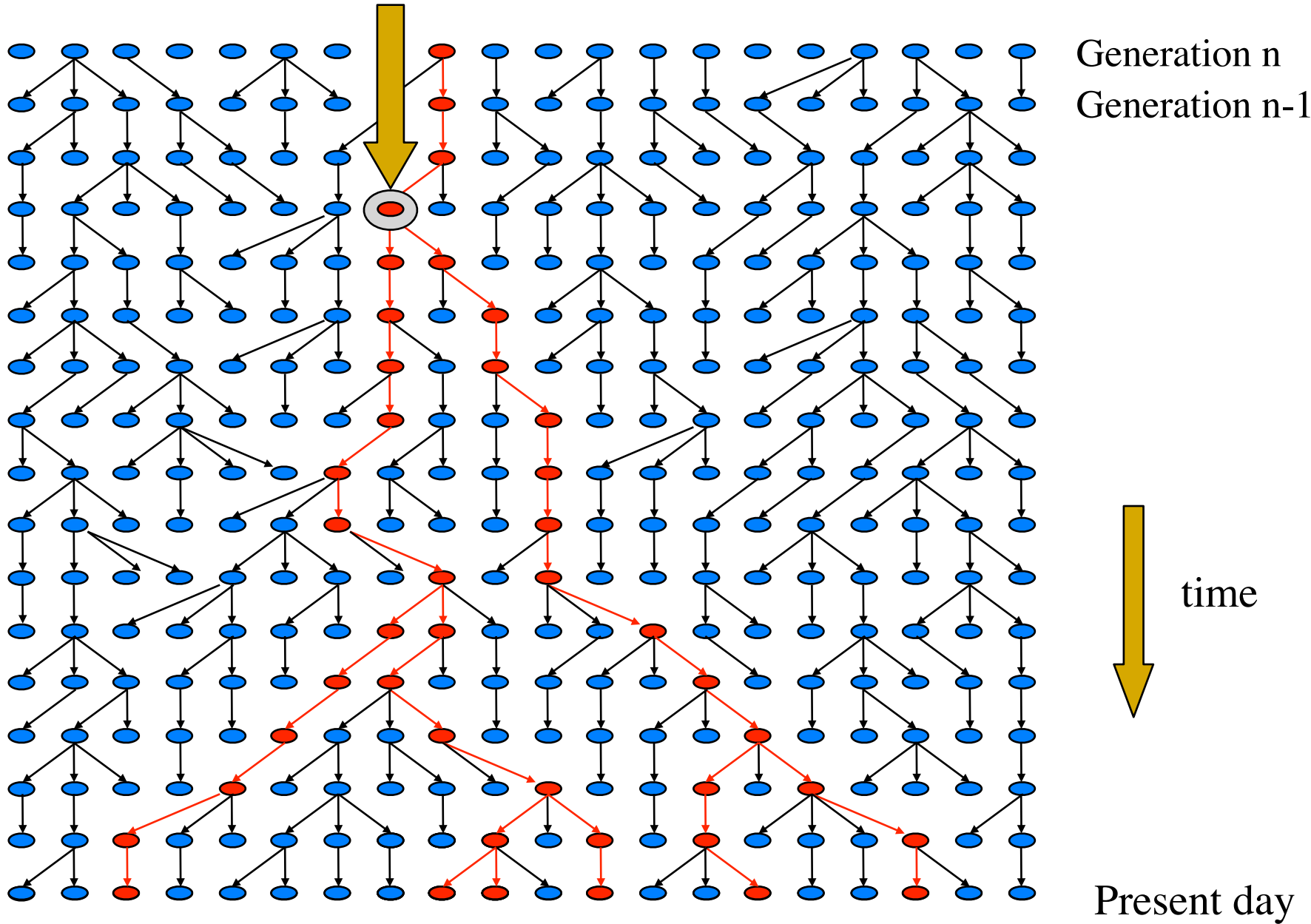


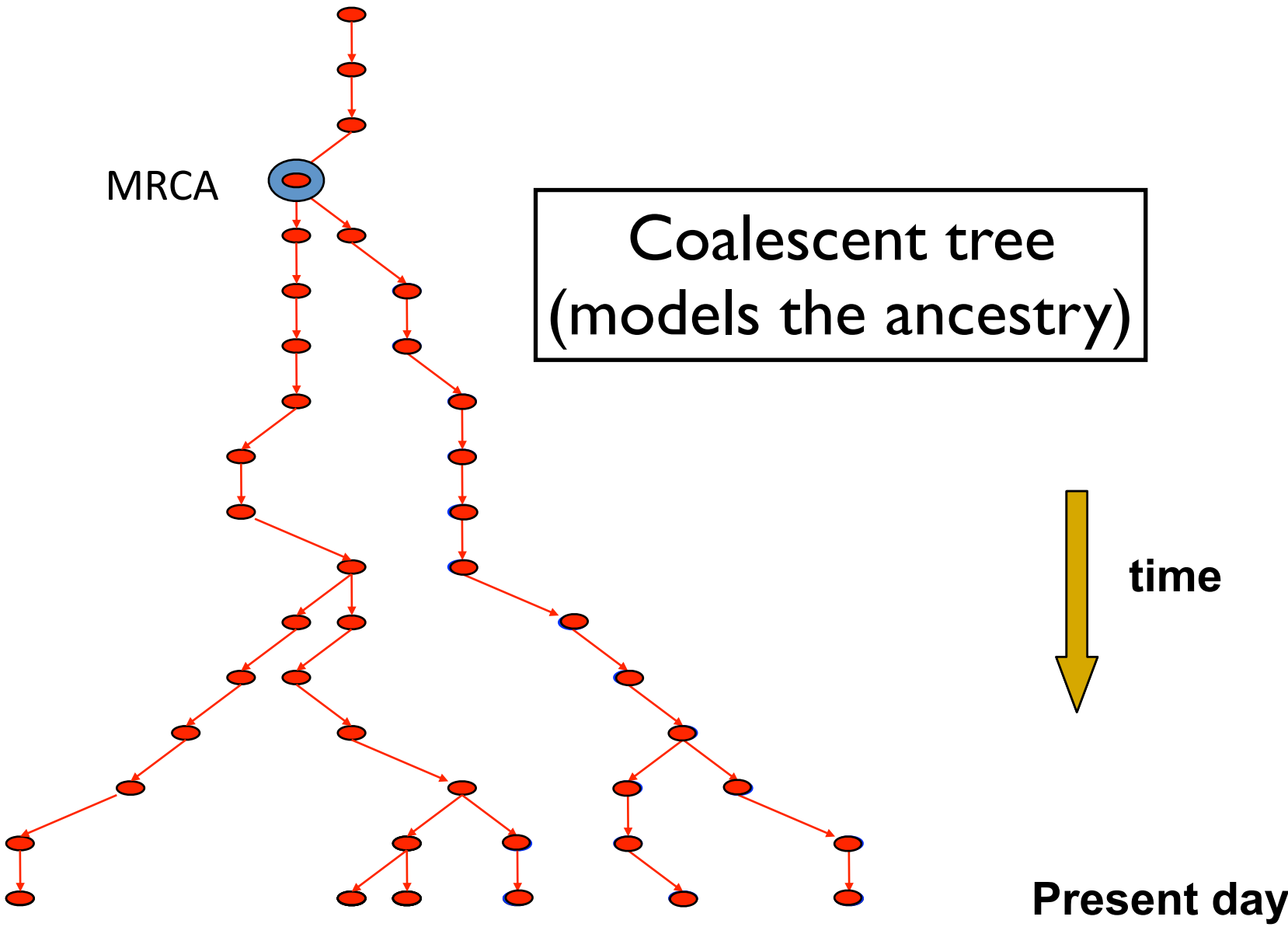
Caspar Friedrich - "Dolmen in the snow" (1807)

Schematic of evolving population



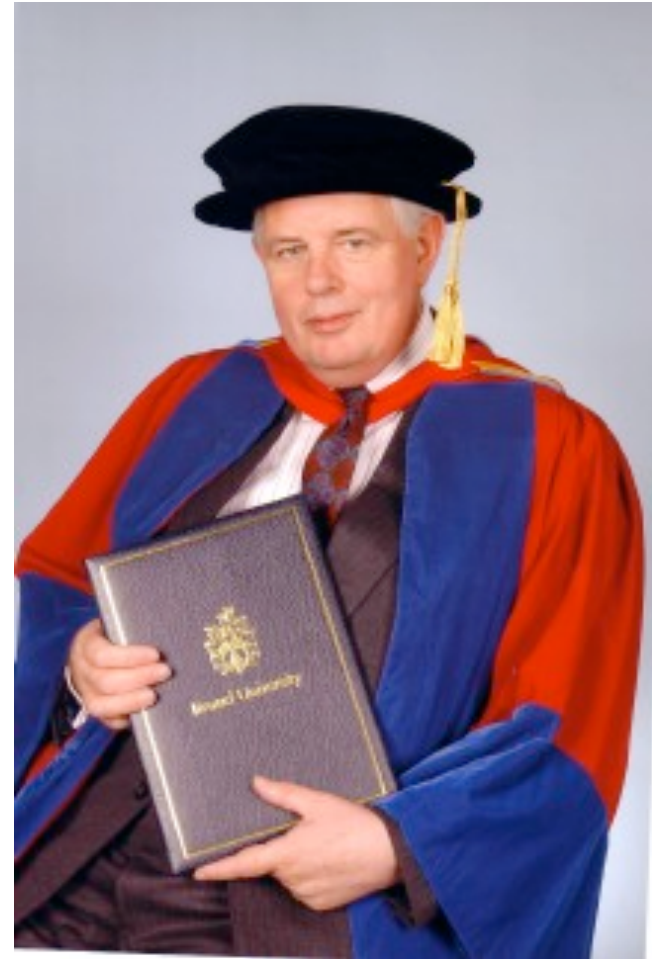
Most recent common ancestor (MRCA)





Ancestral methods

- The coalescent was introduced in Kingman (1982) as a mathematical description of the genealogy that underlies the evolution of a population.
- It has become a standard tool for the analysis of molecular population data.
- It allows for **efficient** modeling/analysis of random samples drawn from a population.



John Kingman, 1939-

Coalescent trees - algorithm

- Time runs in reverse! (We start at the bottom and work our way up.)
- Start with a sample of n individuals (set $k=n$)
- Time to next event $\sim \exp(k(k-1)/2)$
 - Choose two lines of ancestry, uniformly at random, to coalesce. (So keep a list of the all the lines of ancestry that still exist.)
 - Label the new node as $n+(n-k+1)$. Add it to the list of lines of ancestry that exist. Remove the two lines that coalesced.
 - Set $k=k-1$
- Keep going until you get down to one line (i.e., $k=1$)
- Forwards in time, this is equivalent to a model in which the offspring distribution is Multinomial($N, 1/N, 1/N, \dots, 1/N$)

Markov Chain

A *Markov chain* is a sequence of random variables X_1, X_2, X_3, \dots with the Markov property, namely that, given the present state, the future and past states are independent (Wikipedia).

Generic example: a random walk: $X(i+1) = X(i) + R$, where R is some random variable.



Born: 14 June 1856 in Ryazan, Russia

Died: 20 July 1922 in Petrograd (now St Petersburg), Russia

<http://www-history.mcs.st-and.ac.uk/Biographies/Markov.html>

Theoretical results

- Let T_i be the time for the coalescence from i to $i-1$ lines.
- Let T be the overall height of the tree. Then:

$$\begin{aligned} E(T) &= E\left(\sum_{i=2}^n T_i\right) \\ &= \sum_{i=2}^n E(T_i) \\ &= \sum_{i=2}^n 2/i(i-1) \\ &= \frac{2}{n(n-1)} + \frac{2}{(n-1)(n-2)} + \cdots + \frac{2}{3 \times 2} + \frac{2}{2 \times 1} \\ &= \dots \end{aligned}$$

Theoretical results

- Claim $E(\text{Tree height}) = 2(1 - 1/n)$
- Can think of time here as being measured in units of N generations, where N is the population size (Kingman, 1982).
- Note that Expected tree height increases as sample size n increases.
- But also note that $E(\text{Tree height})$ never exceeds (or even reaches) 2 (or $2N$ generations).

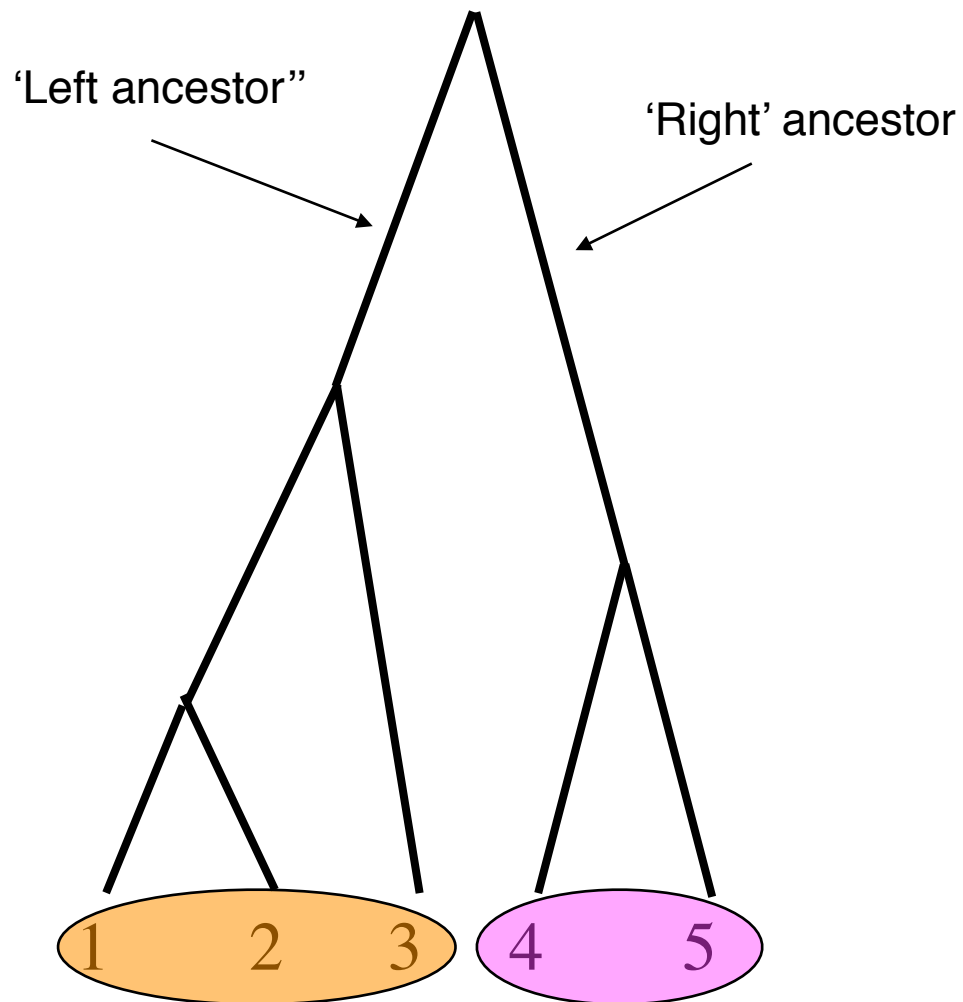
Theoretical results

- Let L_i be the tree length for coalescence from i to $i-1$ lines.
- Let L be the overall height of the tree. Then:

-

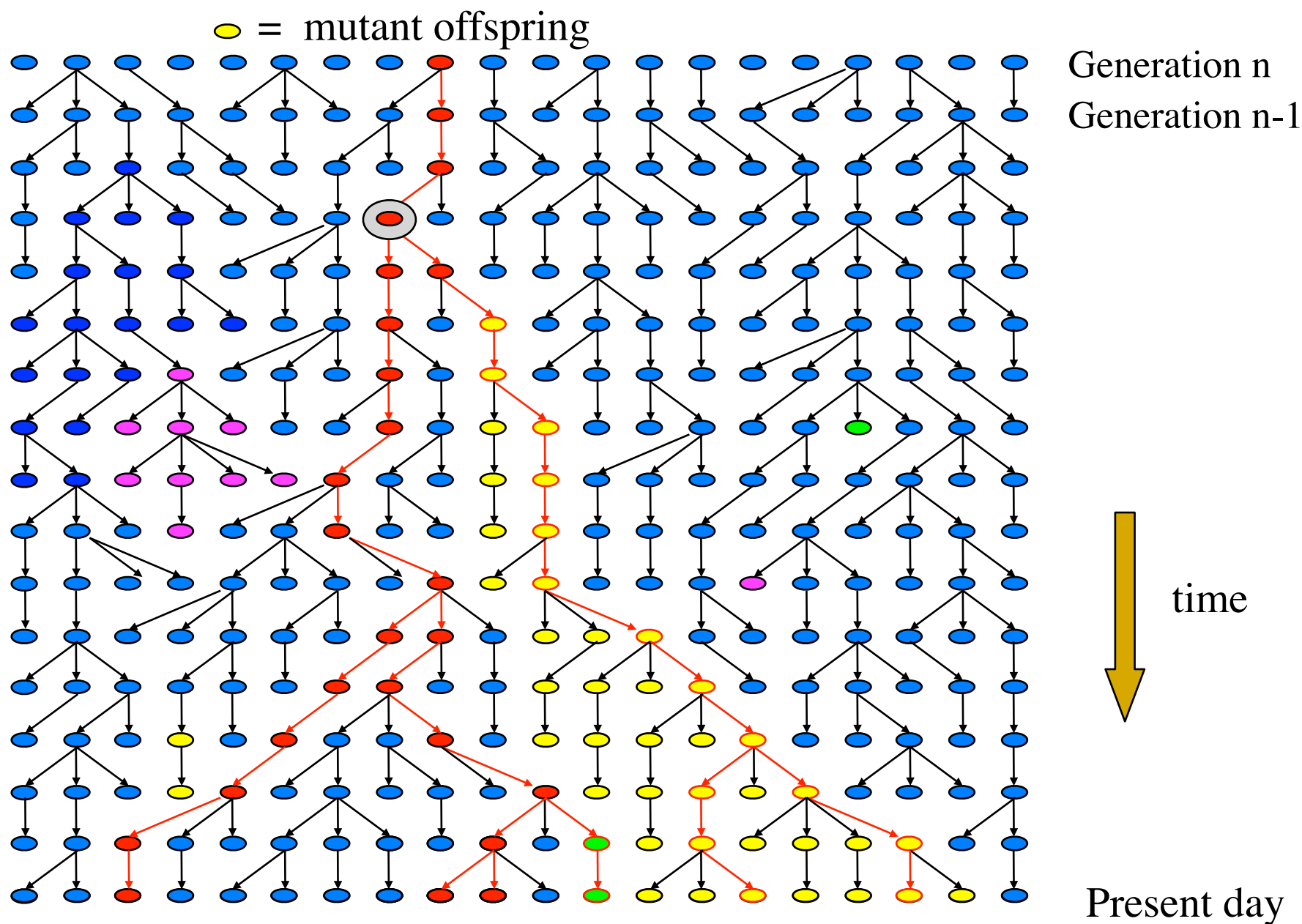
$$\begin{aligned}
 E(L) &= E\left(\sum_{i=2}^n L_i\right) = E\left(\sum_{i=2}^n iT_i\right) \\
 &= \sum_{i=2}^n E(iT_i) \\
 &= \sum_{i=2}^n 2i/i(i-1) \\
 &= \frac{2}{(n-1)} + \frac{2}{(n-2)} + \cdots + \frac{2}{2} + \frac{2}{1} \\
 &= 2 \sum_{i=2}^n \frac{1}{n-1} \\
 &= 2 \sum_{i=1}^{n-1} \frac{1}{n} \\
 &\longrightarrow 2 \log(n-1), \quad \text{as } n \longrightarrow \infty
 \end{aligned}$$

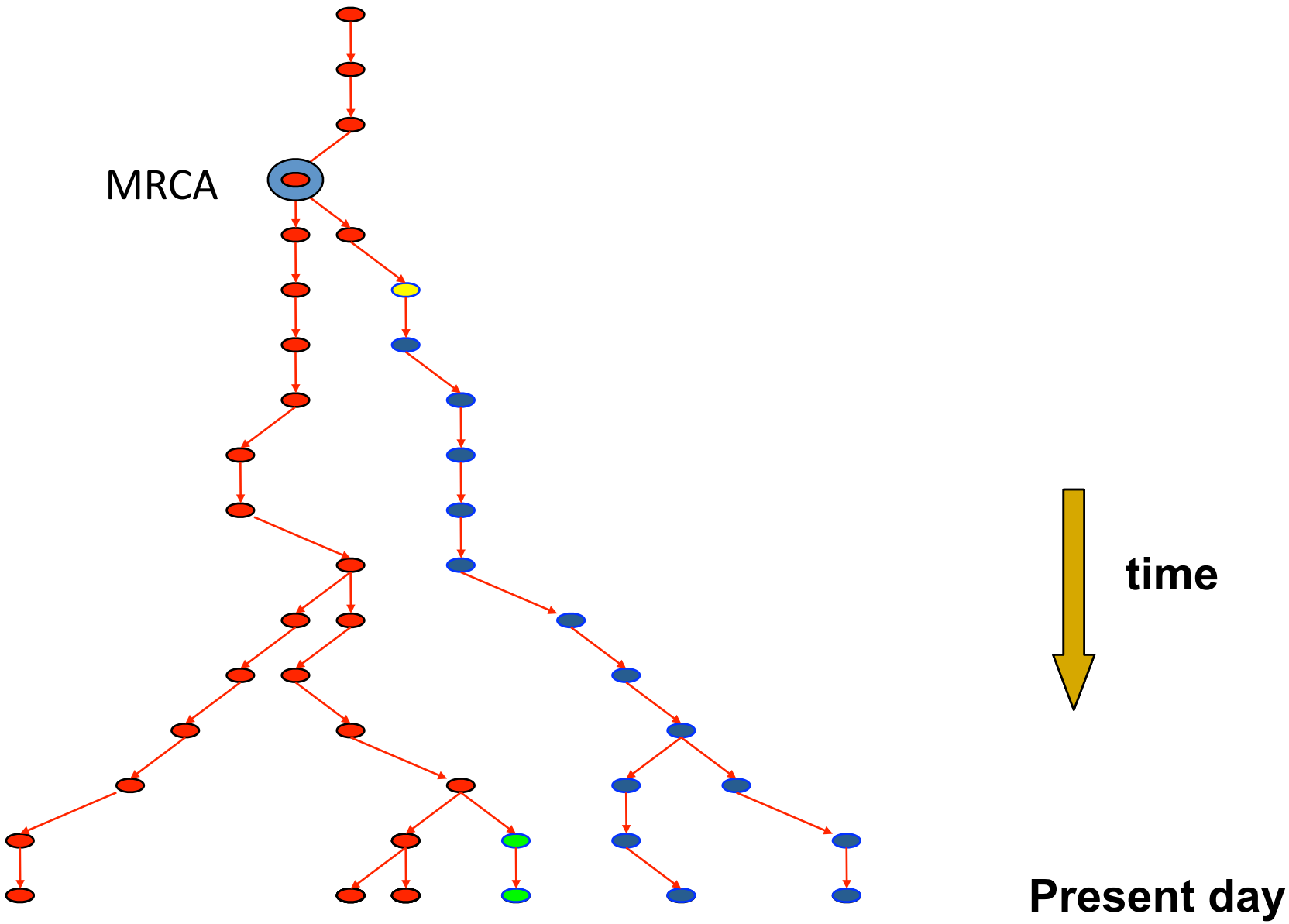
Optional Coalescent - Number of descendants of last two lines

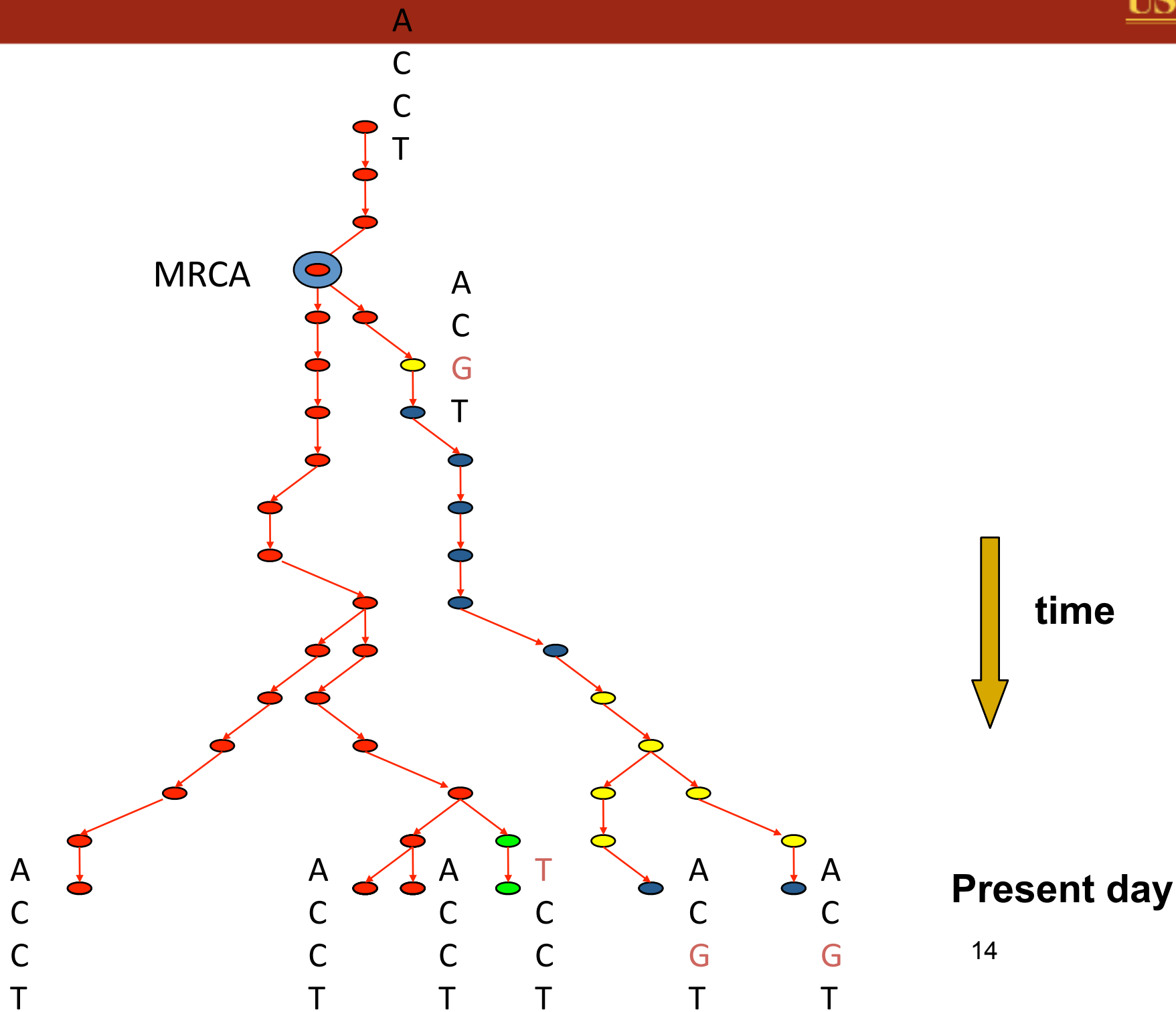


Suppose we have a sample of size 50. What is the distribution of the number of descendants of the 'left (or right) ancestor'?

Now we add mutation







- Mutations may now appear on lines of ancestry.
- Recall, N = population size
- Suppose we have k lines of ancestry, indexed $1, 2, \dots, k$
- In a discrete generation model:
 - $P(\text{lines } i \text{ and } j \text{ coalesce}) = 1/N$
 - k lines: $P(\text{some pair coalesce}) = k(k-1)/2N$
 - $P(\text{line } i \text{ mutates}) = u$
 - k lines: $P(\text{some line mutates}) \sim ku$
 - Define $\theta = 2Nu$: $P(\text{some line mutates}) \sim k\theta/2N$

- $P(\text{coalesce}) = k(k-1)/2N$
- $P(\text{mutation}) = k\theta/2N$
- The **Jump chain** (the chain observed only at times when the state changes) is as follows: one of two things will happen
 - $P(\text{coalesce}) = [k(k-1)/2N] / [k(k-1)/2N + k\theta/2N]$
 $= (k-1)/(k-1+\theta)$
 - $P(\text{mutation}) = \theta/(k-1+\theta)$
- Again, the jump chain is a **Markov chain**. (“What you do next depends only upon where you are, not how you got there.”)

Random nice results:

$$\begin{aligned} E(\# \text{ mutations}) &= \sum_{k=2}^n E\left(\exp\left(\frac{k}{2}\right) k\theta/2\right) \\ &= \sum_{k=2}^n \frac{\theta}{k-1} = \theta \sum_{k=1}^{n-1} \frac{1}{k} \end{aligned}$$

- Watterson's estimator of mutation rate:

$$\hat{\theta} = S / \sum_{k=1}^{n-1} \frac{1}{k}$$

where S is the observed number of segregating sites
(estimate)

A sample is enough

- $P(\text{sample of } n \text{ has same MRCA as population of } N) =$

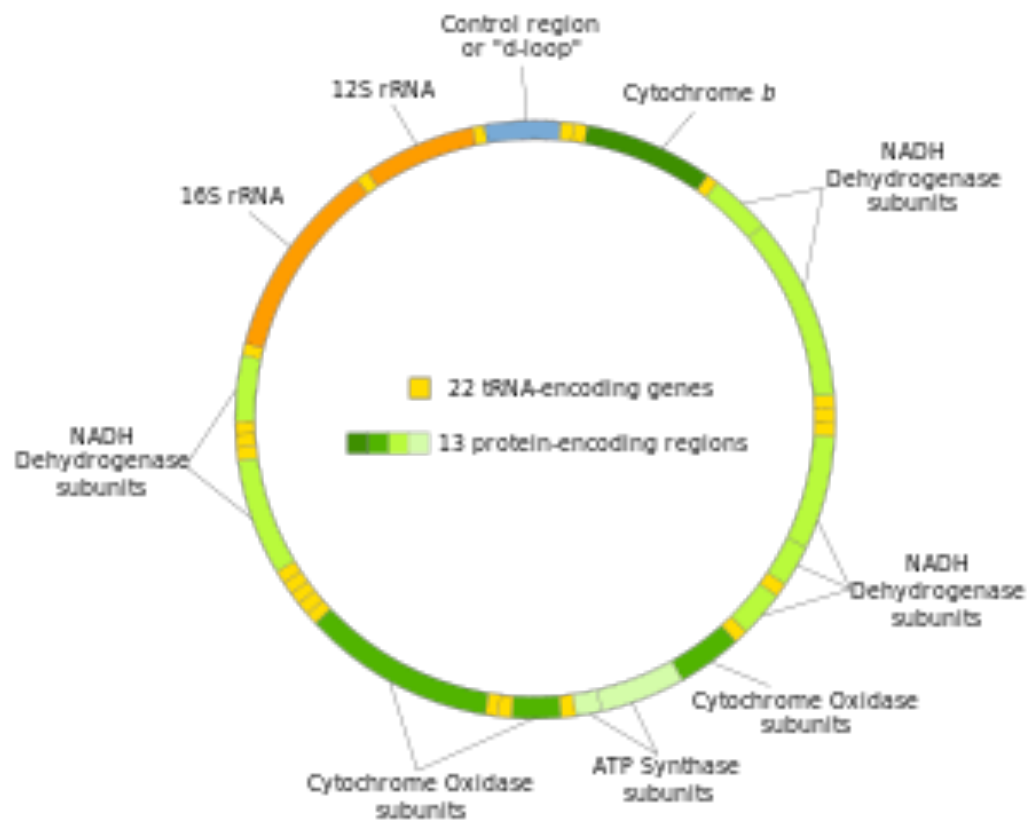
$$\frac{n-1}{n+1} \frac{N+1}{N-1}$$

(Saunders, I.W., Tavaré, S., Watterson, G.A.: On the genealogy of nested subsamples from a haploid population. Adv. Appl. Prob. 16, 471–491 (1984))

Mutation models

- There are many different mutation models:
 - **Infinite Alleles** - each mutation creates a previously unseen allele [type]
 - **Finite Sites** - there are a finite number of sites at which mutation can occur (e.g. mtDNA)
 - **Infinite Sites** - assume there are an infinite number of sites that could mutate. Thus, each mutation will occur at a unique site. (Sounds a bit like Infinite alleles, but actually carries more information about ancestry.).
 - What model for the human genome?

mtDNA



- Just 16000 base pairs
- Loop structure
- Codes for ~40 genes
- Energy factories
- Up to ~1000/cell
- Originally from bacteria engulfed by early eukaryote ancestors

http://en.wikipedia.org/wiki/Mitochondrial_DNA

```

      1 1 1 1 1 1 1 2 2 2 2 2 2 2 2 2 3 3 3 3 3 3
6 8 9 0 2 4 6 6 9 9 0 1 3 4 5 5 6 7 7 9 0 0 0 1 3 4
9 8 1 6 4 9 2 6 0 4 0 9 3 7 1 5 7 1 5 6 1 2 4 9 9 4

```

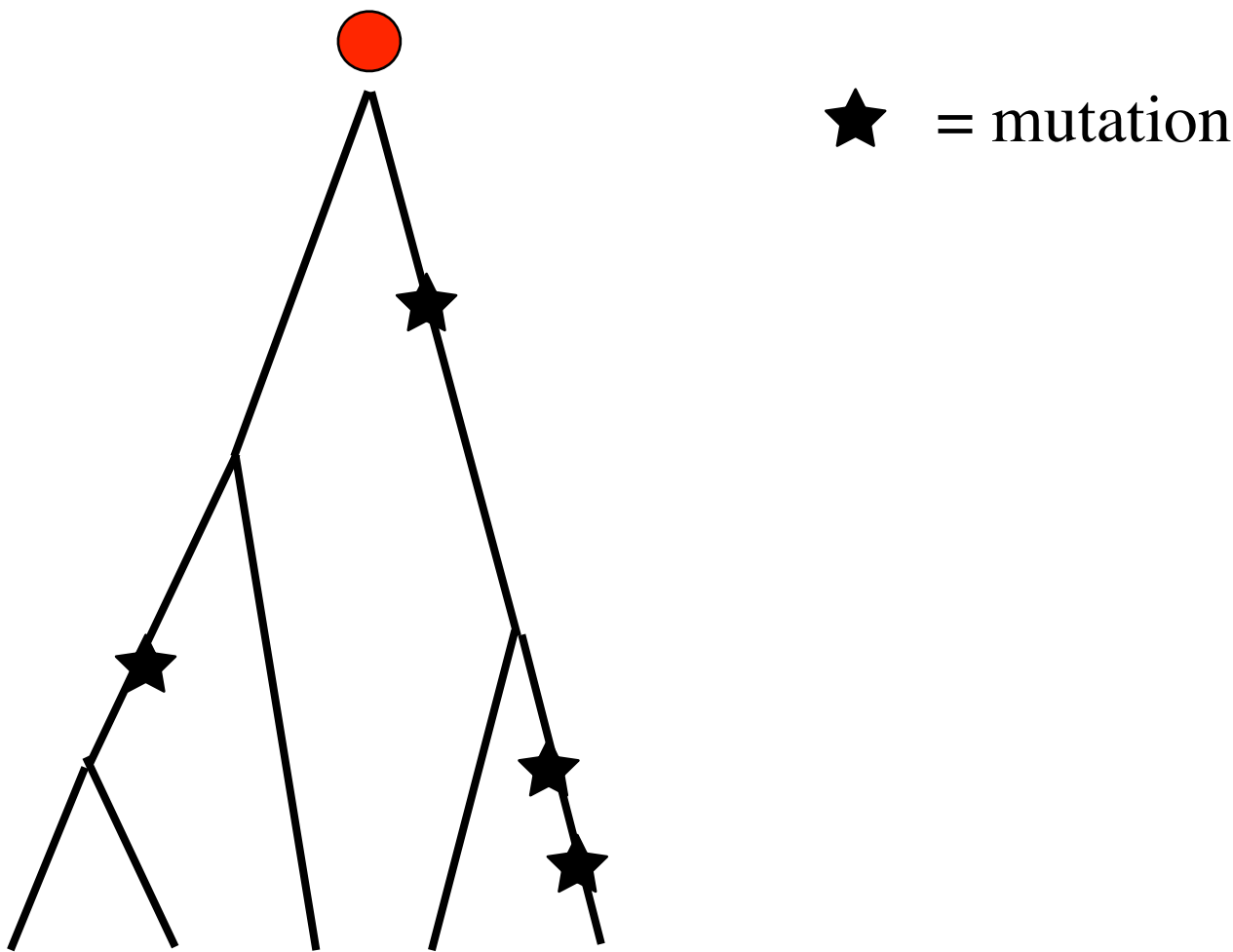
T C C G C T C T G T C C C G C C C T G T T C T T A

[illegible]

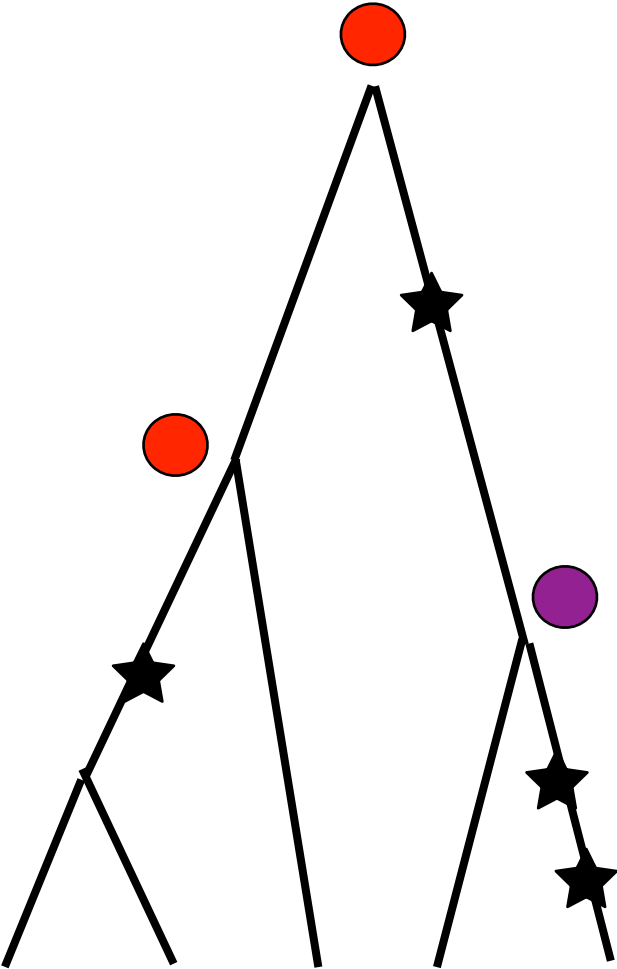
infinite alleles model

- Think of types as colors
 - The type of the bottom of a branch is the same as its type at the top if there are no mutations on the branch
 - Otherwise, it is of some unique, new type

Illustration

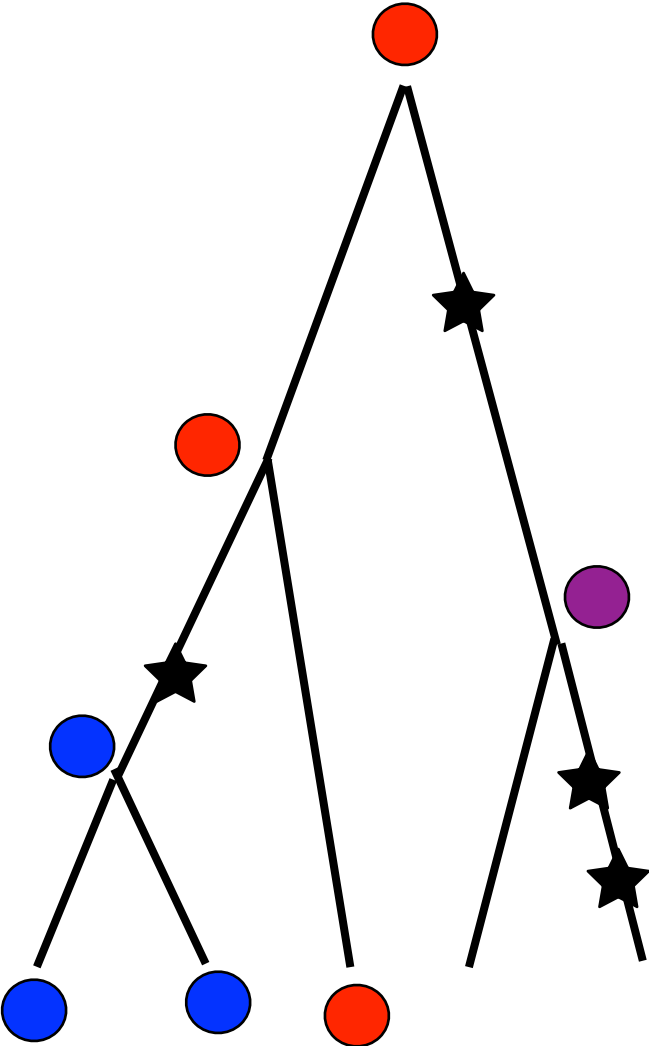


Illustration



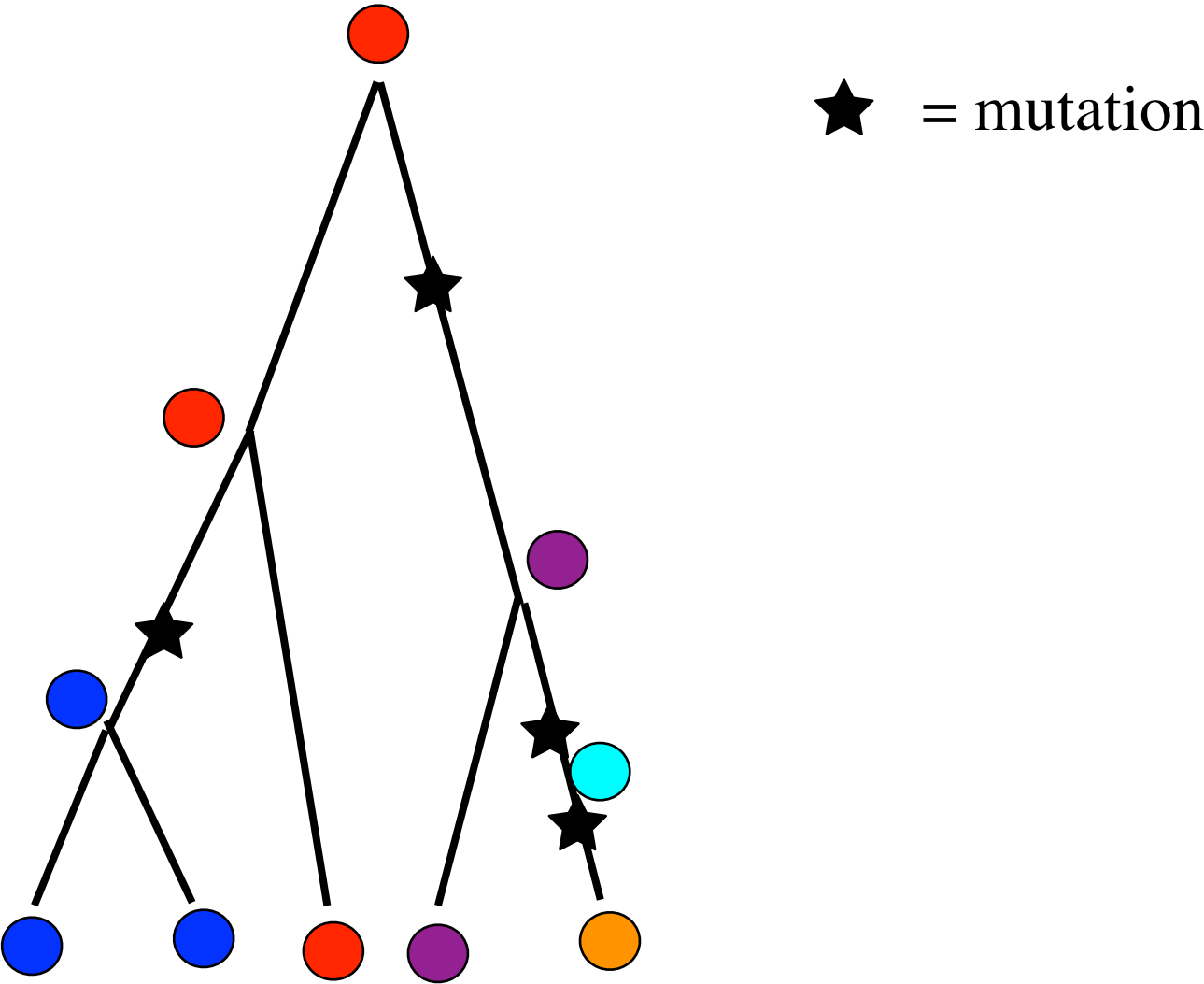
★ = mutation

Illustration

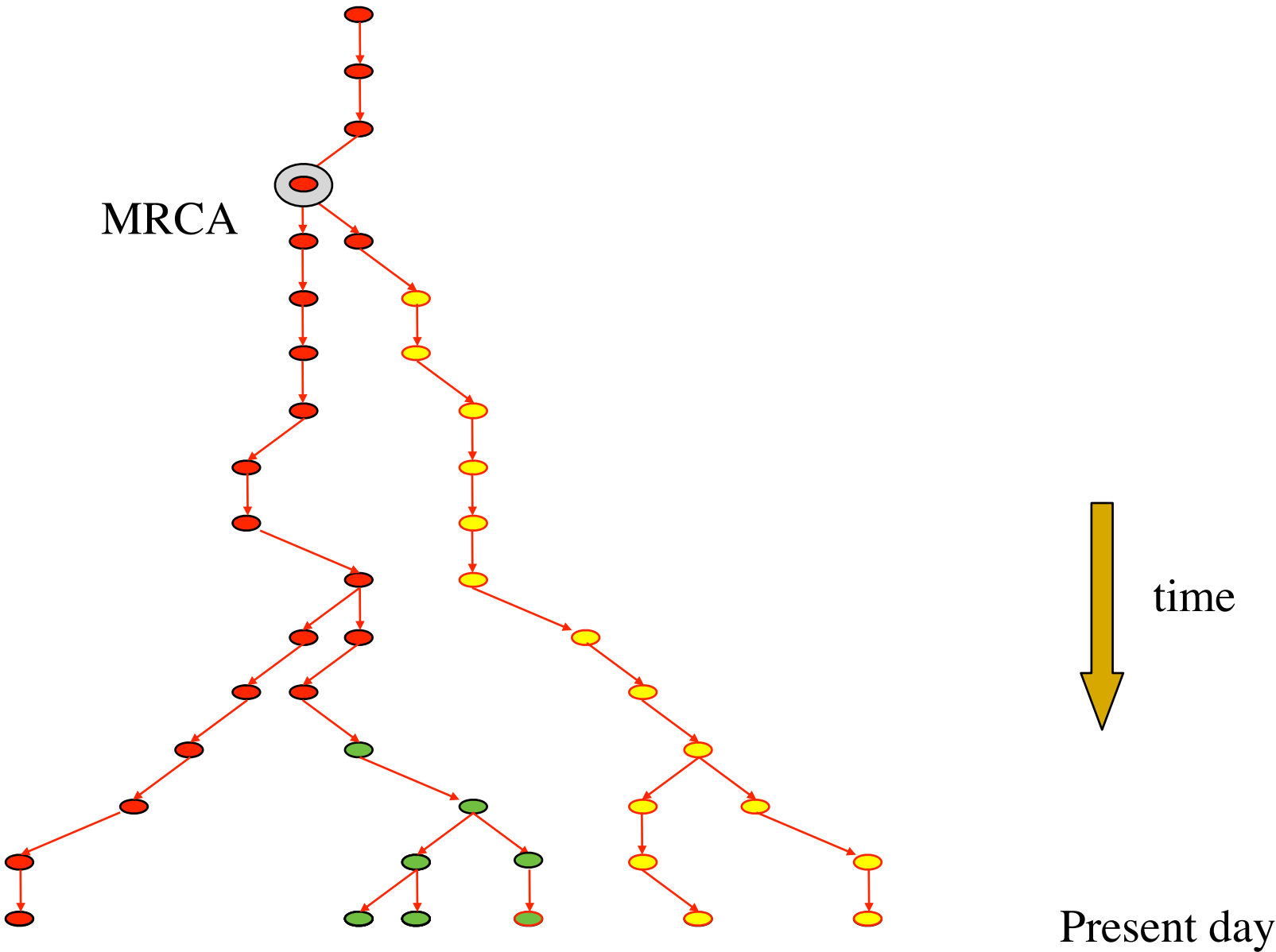


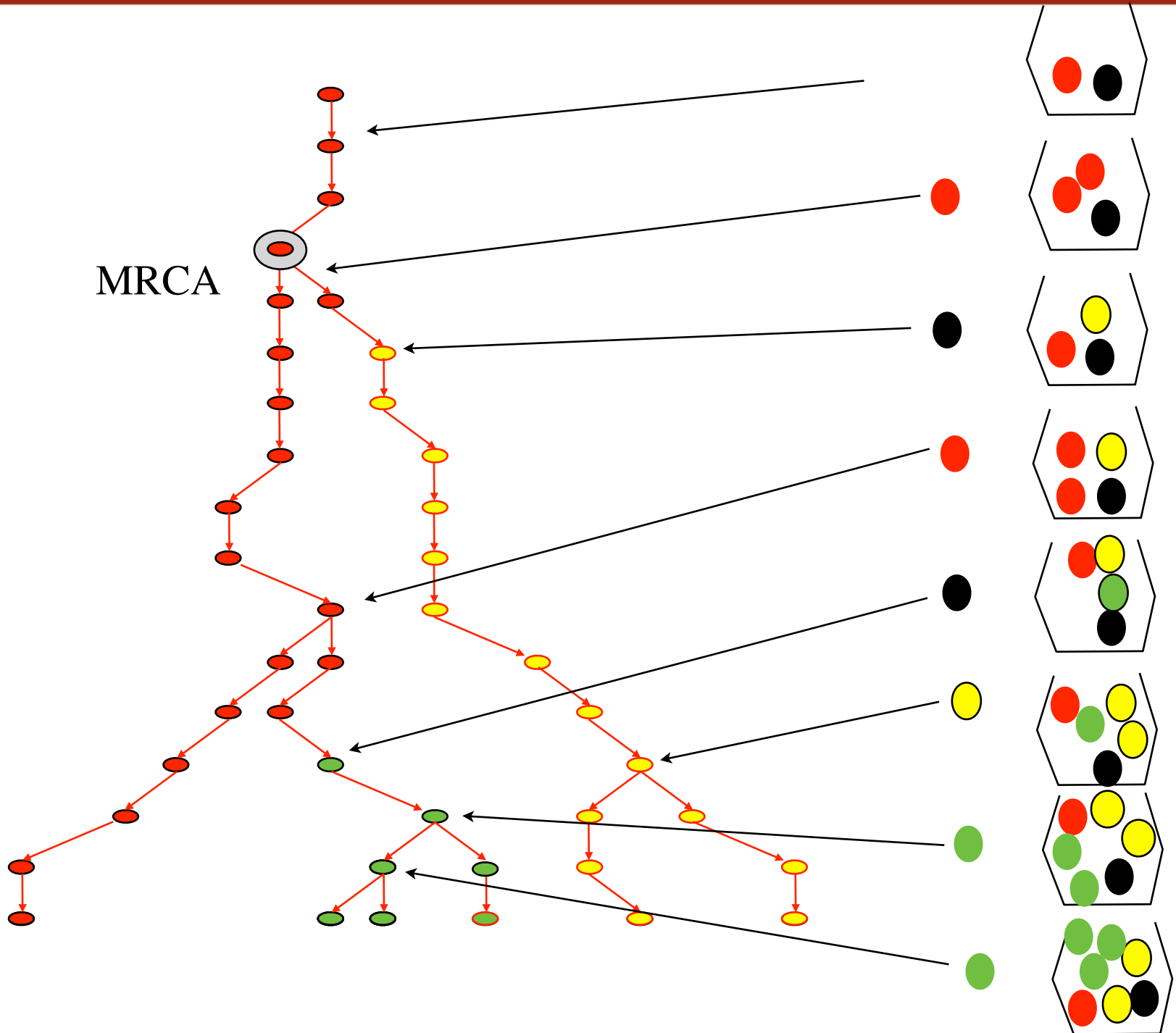
★ = mutation

Illustration



Comparing the coalescent to the Urn model





Real-life is
complicated....

The coalescent with
recombination.

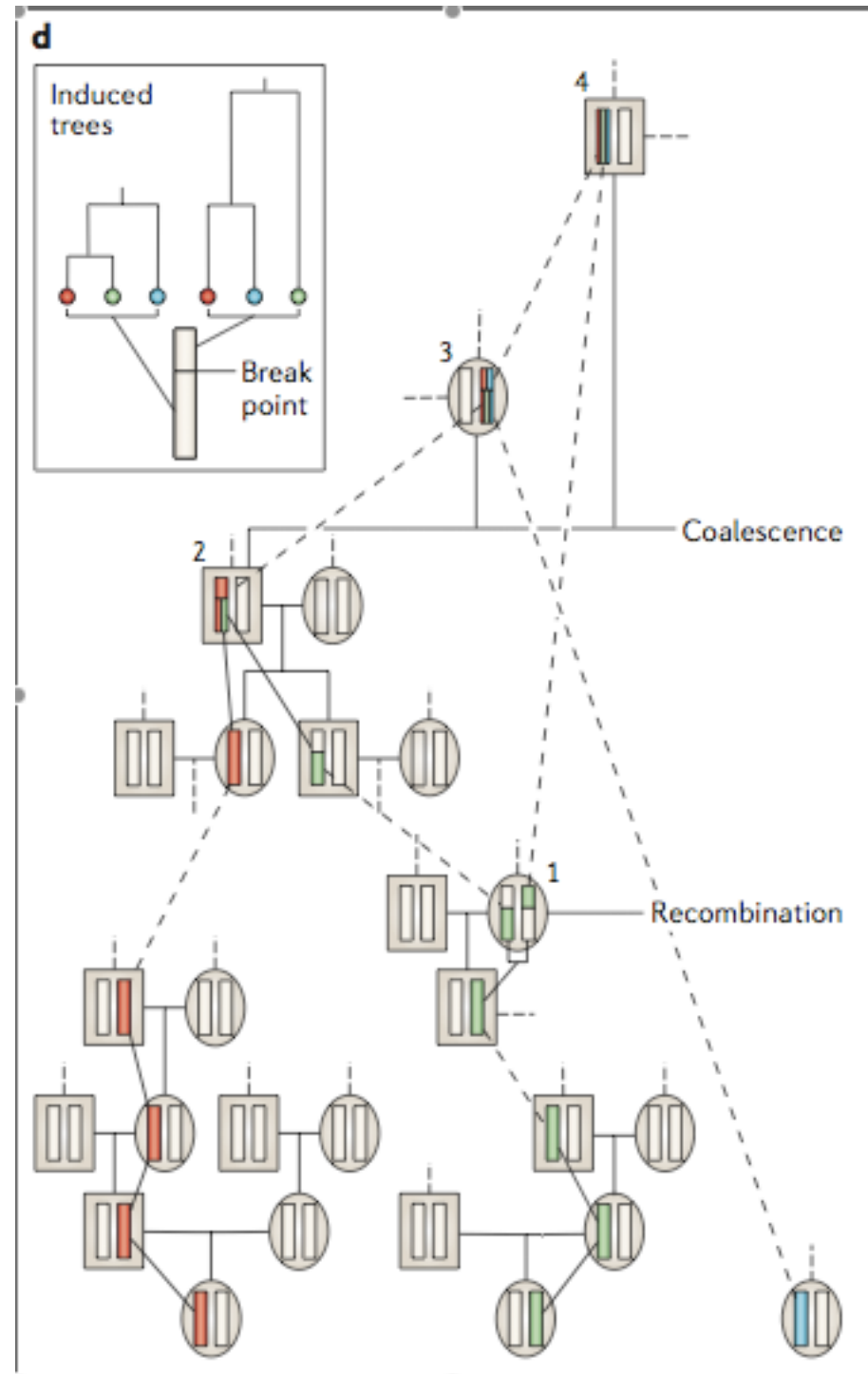
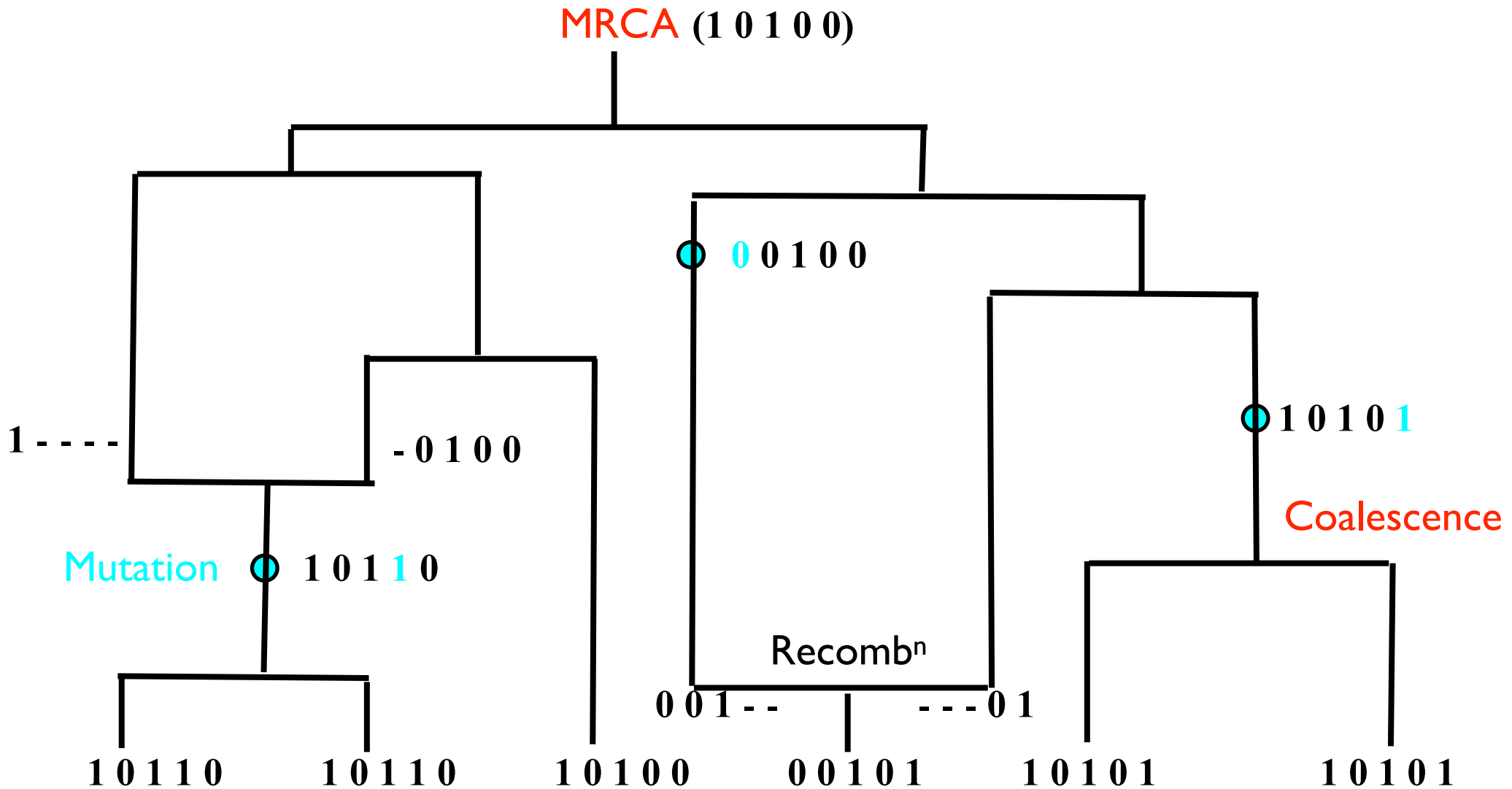


Figure 5: Representation of an ancestry for markers subject to recombination



We trace the ancestry of a sample of 6 marker sequences, until we reach the MRCA. Mutational events are marked in blue. (Markers not ancestral to the sample are marked '-')

Not all mutations on this **graph** will appear in the final sample.

References

Coalescent Theory, M. Nordborg:

<https://onlinelibrary.wiley.com/doi/abs/10.1002/0470022620.bbc21>

Coalescent Theory: An Introduction, J. Wakeley (2009).

https://www.amazon.com/Coalescent-Theory-Introduction-John-Wakeley/dp/0974707759/ref=sr_1_1?crid=AK4KGA9B24JO&keywords=coalescent+theory+wakeley&qid=1582304288&sprefix=coalescent+theory%2Caps%2C207&sr=8-1

Partition structures, Polya urns, the Ewens sampling formula, and the ages of alleles. P. Donnelly (1986).

<https://www.sciencedirect.com/science/article/pii/0040580986900377>

END