

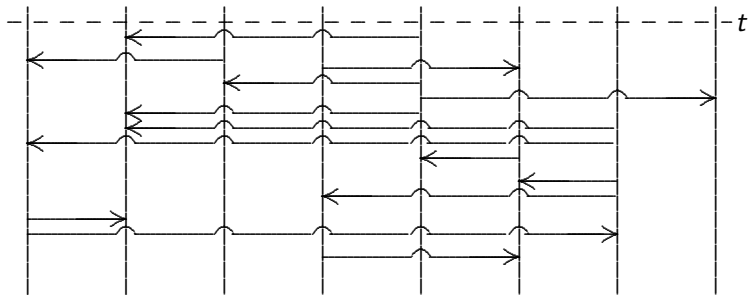
# The evolution of genealogies

Peter Pfaffelhuber  
(joint with Anton Wakolbinger, Heinz Weisshaupt)

Luminy, May 2009

## The Moran model

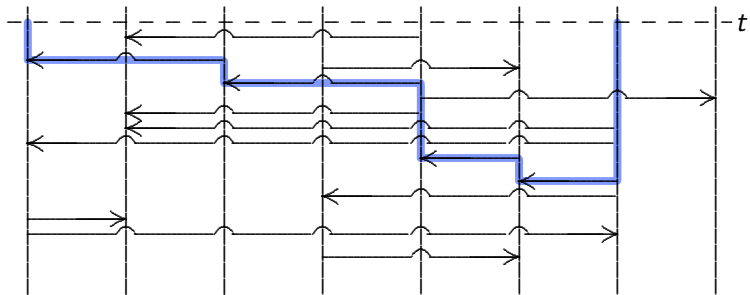
- ▶ A population consists of  $N$  individuals
- ▶ Each pair of individuals **resamples** at rate 1
- ▶ Resampling means: one individual **dies**, the other **reproduces**



## The Moran model

- ▶ A population consists of  $N$  individuals
- ▶ Each pair of individuals **resamples** at rate 1
- ▶ Resampling means: one individual **dies**, the other **reproduces**

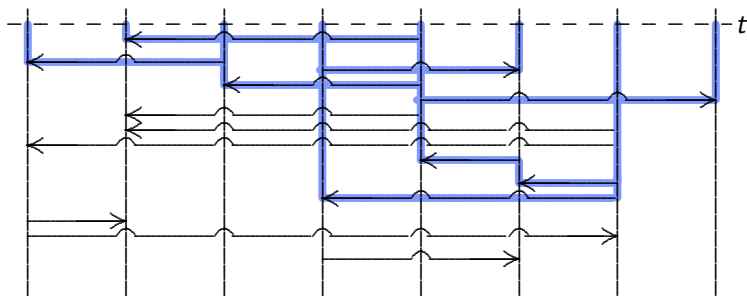
Ancestral lineages coalesce



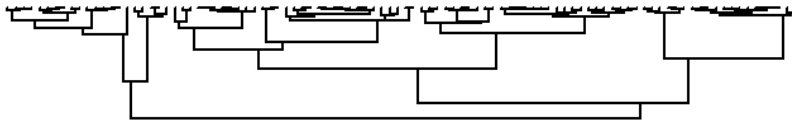
## The Moran model

- ▶ A population consists of  $N$  individuals
- ▶ Each pair of individuals **resamples** at rate 1
- ▶ Resampling means: one individual **dies**, the other **reproduces**

Ancestral lineages coalesce



## Kingman's coalescent



- ▶ **Genealogies** are given by **Kingman's N-coalescent**
- ▶ Coalescence rate is  $\binom{k}{2}$ .
- ▶ What are properties of the **tree length**?

# Kingman's coalescent

- $\mathcal{L}_t^N$ : tree length of N-coalescent at time  $t$

$$\mathbb{E}[\mathcal{L}_t^N] = \sum_{k=2}^N k \frac{1}{\binom{k}{2}} \stackrel{N \rightarrow \infty}{\approx} 2 \log(N)$$

$$\mathbb{V}[\mathcal{L}_t^N] = \sum_{k=2}^N k^2 \frac{1}{\binom{k}{2}^2} \stackrel{N \rightarrow \infty}{\approx} 4 \frac{\pi^2}{6}$$

## Kingman's coalescent

- ▶  $\mathcal{L}_t^N$ : **tree length** of N-coalescent at time  $t$
- ▶  $\mathcal{E}(\cdot)$ : independent exponential distributions

$$\frac{1}{2}\mathcal{L}_t^N \stackrel{d}{=} \frac{1}{2} \sum_{k=2}^N k \cdot \mathcal{E}\left(\binom{k}{2}\right) \stackrel{d}{=} \sum_{k=1}^{N-1} \mathcal{E}(k) \stackrel{d}{=} \max_{1 \leq k \leq N-1} \mathcal{E}(1)$$

$$\mathbb{P}\left[\frac{1}{2}(\mathcal{L}_t^N - 2 \log(N)) \leq t\right] = (1 - e^{-\log(N)-t})^{N-1} \approx e^{-e^t}$$

- ▶  $\Rightarrow \frac{1}{2}(\mathcal{L}_t^N - 2 \log(N)) \xrightarrow{N \rightarrow \infty} \text{Gumbel}$

## The Gumbel variable in the coalescent

- ▶ Are there **stronger versions** of tree length convergence on a coalescent?
- ▶ Consider subtrees with  $N$  leaves of a full Kingman coalescent
- ▶  $X_i \stackrel{d}{=} \mathcal{E}(\frac{1}{2})$ : time full coalescent stays with  $i$  lines

$$L_N^1 = \sum_{i=2}^N iX_i \quad \text{Temporal coupling}$$

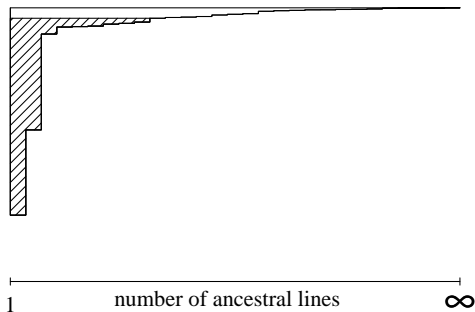
- ▶  $K_i^N$ : # lines in  $N$ -tree while full tree has  $i$  lines

$$L_N^2 = \sum_{i=2}^{\infty} K_i X_i \quad \text{Natural coupling}$$



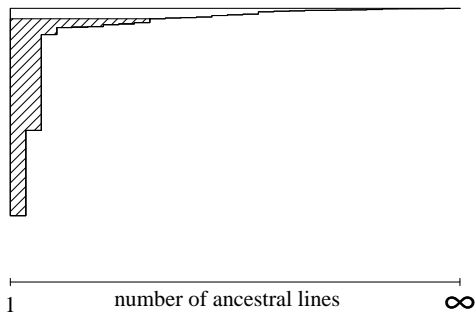


## The Gumbel variable in the coalescent



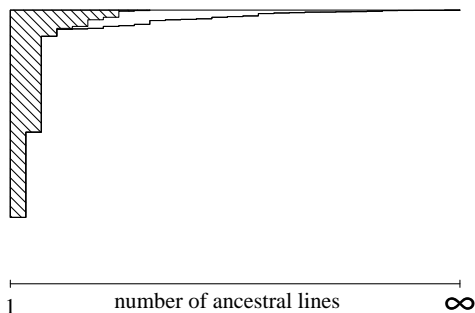
- ▶  $\mathcal{L}_t^{N,1} = \sum_{i=2}^N iX_i$
- ▶  $X_i$ : time full coalescent stays with  $i$  lines

## The Gumbel variable in the coalescent



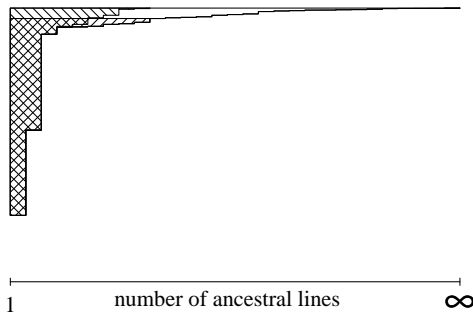
- ▶  $\frac{1}{2}(\mathcal{L}_t^{N,1} - 2 \log N) \xrightarrow{N \rightarrow \infty} \mathcal{L}_t$  almost surely and in  $L^2$
- ▶  $\mathcal{L}_t$ : Gumbel distributed

## The Gumbel variable in the coalescent



- ▶  $K_i^N$ : # lines in N-tree while full tree has  $i$  lines
- ▶  $\mathcal{L}_t^{N,2} = \sum_{i=2}^{\infty} K_i X_i$

# The Gumbel variable in the coalescent



►  $\mathcal{L}_t^{N,1} - \mathcal{L}_t^{N,2} \xrightarrow{N \rightarrow \infty} \mathbf{0}$  in  $L^2$

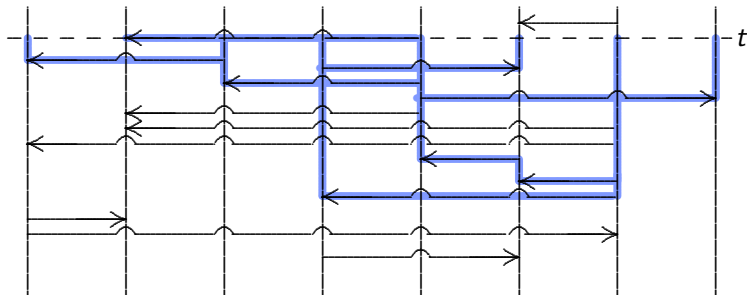
## Sample path

- ▶ Genealogies evolve together with the population
- ▶ Show movie
- ▶ Rest of the talk:

**What does the evolution of tree lengths  $\mathcal{L}_t^N$  look like?**

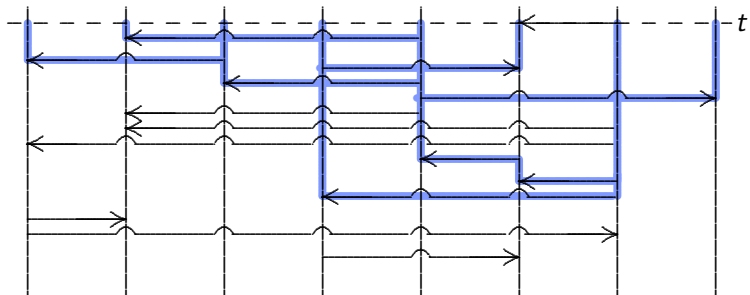
## The Moran model

- ▶ Genealogies evolve as time proceeds
- ▶ Tree growth at speed  $Ndt$  between resampling events



## The Moran model

- ▶ Genealogies evolve as time proceeds
- ▶ Tree growth at speed  $Ndt$  between resampling events

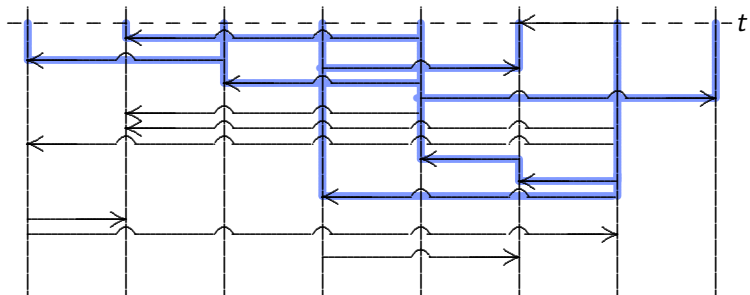




## The Moran model

- ▶ Genealogies evolve as time proceeds
- ▶ Tree growth at speed  $Ndt$  between resampling events
- ▶ At resampling times the tree length changes

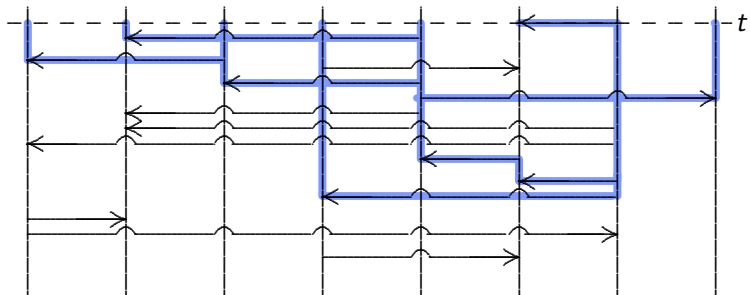
External branches break off



## The Moran model

- ▶ Genealogies evolve as time proceeds
- ▶ Tree growth at speed  $Ndt$  between resampling events
- ▶ At resampling times the tree length changes

External branches break off



## Typical jumps

- ▶  $F$ : jump time of  $N$ -coalescent
- ▶  $J^N$ : length of a random external branch

$$\mathcal{L}_F^N - \mathcal{L}_{F-}^N \stackrel{d}{=} J^N$$

- ▶ Fu, Li; Durrett; Caliebe, Neiniger, Krawczak, Rösler

$$N \cdot J^N \xrightarrow{N \rightarrow \infty} J, \quad \mathbb{E}[J] = 2, \quad \mathbb{V}[J] = \infty$$

- ▶  $\Rightarrow \binom{N}{2}$  jumps of size  $2/N$  per time unit.

## Main result

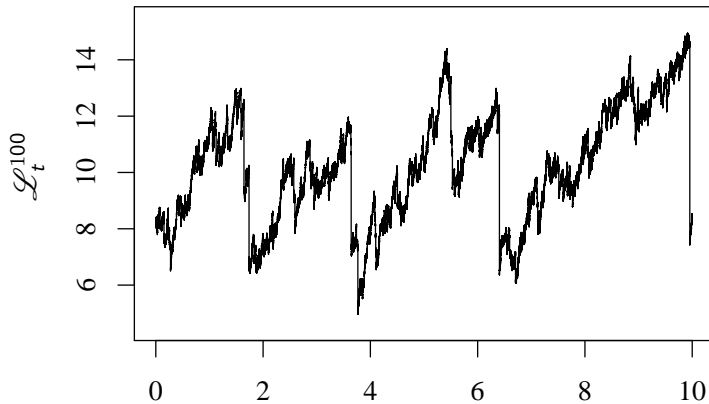
- **There is a** process  $\mathcal{L} = (\mathcal{L}_t)_{t \in \mathbb{R}}$  with càdlàg paths such that

$$\mathcal{L}^N - 2 \log(N) \Longrightarrow \mathcal{L} \text{ as } N \rightarrow \infty.$$

The process  $\mathcal{L}$  has **infinite quadratic variation**; in particular,

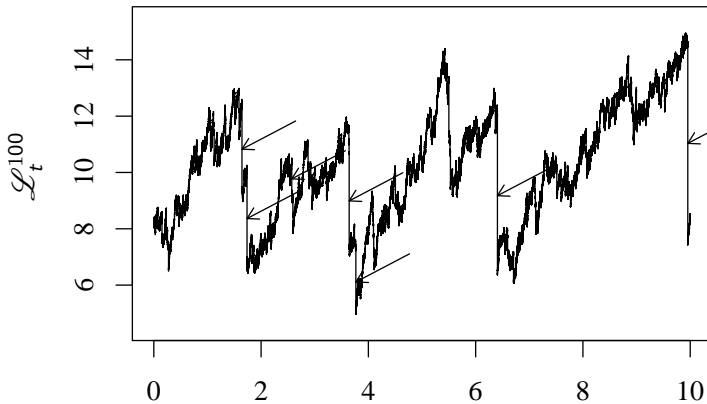
$$\frac{1}{t |\log t|} \mathbb{E}[(\mathcal{L}_t - \mathcal{L}_0)^2] \stackrel{t \rightarrow 0}{\sim} 2.$$

# Sample path

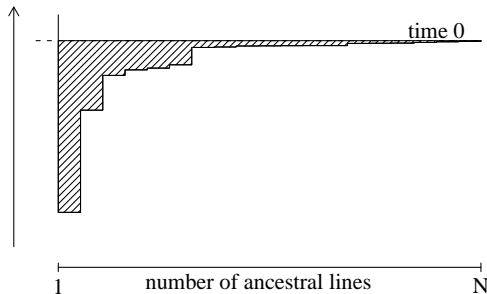


# Sample path

When **MRCA jumps**, tree lengths jump as well

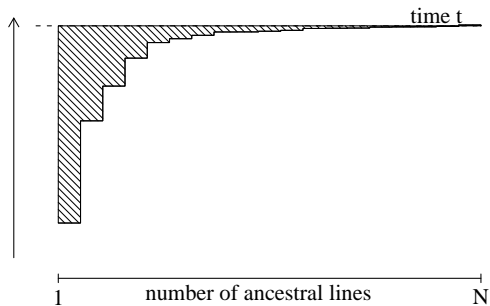


## Ideas for tightness



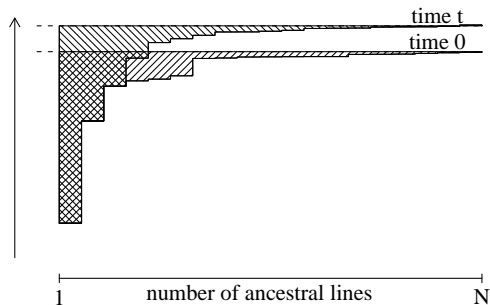
- ▶  $S_s^N := \#$  ancestors at time  $-s$  ( $:=0$  before MRCA)
- ▶  $\mathcal{L}_0^N = \int_0^\infty S_s^N ds$

## Ideas for tightness





## Ideas for tightness



- ▶ Trees at times 0,  $t$  **overlap**



## Ideas for tightness

- ▶ **Gain in tree length:**  $A_{0,t}^N \stackrel{d}{=} \int_0^t S_s^N ds$
- ▶ Coalescent comes down from  $\infty \Rightarrow S_s^N \xrightarrow{N \rightarrow \infty} S_s$
- ▶ Aldous (1999):

$$\frac{S_s - 2/s}{\sqrt{2/(3s)}} \xrightarrow{s \rightarrow 0} N(0, 1)$$

- ▶ Extension:  $r \leq s \Rightarrow \text{COV}[S_r, S_s] \stackrel{s \rightarrow 0}{\approx} \frac{2}{3} \frac{r}{s^2}$

$$\lim_{N \rightarrow \infty} \mathbb{V}[A_{0,t}^N] = 2 \int_0^t \int_0^s \text{COV}[S_r, S_s] dr ds \stackrel{t \rightarrow 0}{\approx} \frac{2}{3} t$$

## Ideas for tightness

▶ **Loss in tree length**

$$B_{0,t}^N \stackrel{d}{=} \sum_{i=2}^N (i - K_i^{N, S_t^N}) \mathcal{E} \left( \binom{i}{2} \right)$$

- ▶  $K_i^{N,K} := \#$  lines in  $K$ -tree while the  $N$ -tree has  $i$  lines.
- ▶  $S_t^N \approx 2/t$ ,  $(K_i^{N,K})_{i=N, N-1, \dots}$  is Markov Chain

$$\xrightarrow{\text{some calculations}} \lim_{N \rightarrow \infty} \mathbb{V}[B_{0,t}^N] \stackrel{t \rightarrow 0}{\approx} 2t |\log t|$$

## Ideas for tightness

► **Collecting** terms

$$\begin{aligned} \lim_{N \rightarrow \infty} \mathbb{E}[(\mathcal{L}_t^N - \mathcal{L}_0^N)^2] &= \lim_{N \rightarrow \infty} \mathbb{V}[A_{0,t}^N - B_{0,t}^N] \\ &\stackrel{t \rightarrow 0}{\approx} \lim_{N \rightarrow \infty} \mathbb{V}[B_{0,t}^N] \\ &\stackrel{t \rightarrow 0}{\approx} 2t |\log(t)| \end{aligned}$$

# Outlook

- ▶ General theory shows:  $\mathcal{L}^N \xrightarrow{N \rightarrow \infty} \mathcal{L}$   
**in probability** on the Lookdown probability space  
Does **almost sure convergence** hold as well?
- ▶ What is the joint evolution of  $(\mathcal{D}_t, \mathcal{L}_t)$  of the tree?  
( $\mathcal{D}_t =$  **depth of the tree** at time  $t$ )
- ▶ Take a Cannings model with **finite offspring variance**. Does

$$\mathcal{L}^{\text{Cannings}, N} \xrightarrow{N \rightarrow \infty} \mathcal{L}?$$

- ▶ What are **other limits** of  $\mathcal{L}^{\text{Cannings}, N}$  for Cannings models with infinite offspring variance?

# Summary

- ▶ Tree lengths in Kingman's coalescent are **Gumbel** distributed
- ▶ Evolution of tree lengths gives a **càdlàg** process  $\mathcal{L}$
- ▶  $\mathcal{L}$  has **infinite quadratic variation**