

On the total external length of the Kingman coalescent

Svante Janson*

Götz Kersting[†]

Abstract

In this paper we prove asymptotic normality of the total length of external branches in Kingman's coalescent. The proof uses an embedded Markov chain, which can be described as follows: Take an urn with n *black* balls. Empty it in n steps according to the rule: In each step remove a randomly chosen pair of balls and replace it by one *red* ball. Finally remove the last remaining ball. Then the numbers U_k , $0 \leq k \leq n$, of red balls after k steps exhibit an unexpected property: (U_0, \dots, U_n) and (U_n, \dots, U_0) are equal in distribution.

Key words: coalescent, external branch, reversibility, urn model.

AMS 2010 Subject Classification: Primary 60K35; Secondary: 60F05, 60J10.

Submitted to EJP on February 3, 2011, final version accepted October 9, 2011.

*Department of Mathematics, Uppsala University, PO Box 480, SE-751 06 Uppsala, Sweden.
svante.janson@math.uu.se

[†]Fachbereich Informatik und Mathematik, Goethe Universität, Fach 187, D-60054 Frankfurt am Main, Germany.
kersting@math.uni-frankfurt.de

1 Introduction and results

Our main result in this paper is that the total length L_n of all external branches in Kingman's coalescent with n external branches is asymptotically normal as $n \rightarrow \infty$.

Kingman's coalescent (1982) consists of two components. First there are the coalescent times $T_1 > T_2 > \dots > T_n = 0$. They are such that

$$\binom{k}{2} (T_{k-1} - T_k), \quad k = 2, \dots, n$$

are independent, exponential random variables with expectation 1. Second there are partitions $\pi_1 = \{\{1, \dots, n\}\}, \pi_2, \dots, \pi_n = \{\{1\}, \dots, \{n\}\}$ of the set $\{1, \dots, n\}$, where the set π_k contains k disjoint subsets of $\{1, \dots, n\}$ and π_{k-1} evolves from π_k by merging two randomly chosen elements of π_k . Moreover, (T_n, \dots, T_1) and (π_n, \dots, π_1) are independent. For convenience we put $\pi_0 := \emptyset$.

As is customary the coalescent can be represented by a tree with n leaves labelled from 1 to n . Each of these leaves corresponds to an external branch of the tree. The other node of the branch with label i is located at level

$$\rho(i) := \max\{k \geq 1 : \{i\} \notin \pi_k\}$$

within the coalescent. The length of this branch is $T_{\rho(i)}$. The total external length of the coalescent is given by

$$L_n := \sum_{i=1}^n T_{\rho(i)}.$$

This quantity is of a certain statistical interest. Coalescent trees have been introduced by Kingman (1982) as a model for the genealogy of n individuals, down to their most recent common ancestor. Mutations can be located everywhere on the branches. Then mutations on external branches affect only single individuals. This fact was used by Fu and Li (1993) in designing their D -statistic and providing a test whether or not data fit to Kingman's coalescent.

Elsewhere the total external length of coalescents has been studied by Möhle (2010). He obtained results on the asymptotic distribution for a class of coalescents, which differ substantially from Kingman's coalescent. It includes so-called Beta($2 - \alpha, \alpha$)-coalescents with $0 < \alpha < 1$. For $1 < \alpha < 2$ Berestycki et al (2006) proved a law of large numbers (see the quantity $M_1(n)$ in their Theorem 9); a more general result is contained in Berestycki et al (2011). Otherwise single external branches have been investigated in the literature. The asymptotic distribution of $T_{\rho(i)}$ has been obtained by Caliebe et al (2007), using a representation of its Laplace transform due to Blum and François (2005). Freund and Möhle (2009) studied the Bolthausen-Sznitman coalescent, and Gnedin et al (2008) the general Λ -coalescent.

Here is our main result.

Theorem 1. *As $n \rightarrow \infty$,*

$$\frac{1}{2} \sqrt{\frac{n}{\log n}} (L_n - 2) \xrightarrow{d} N(0, 1).$$

Here \xrightarrow{d} denotes convergence in distribution. The proof will show that the limiting normal distribution originates from the random partitions and not from the exponential waiting times.

A second glance on this result reveals a peculiarity: The normalization of L_n is carried out using its expectation, but only half of its variance. These two terms have been determined by Fu and Li (1993) (with a correction given by Durrett (2002)). They obtained

$$\mathbf{E}(L_n) = 2, \quad \mathbf{Var}(L_n) = \frac{8nh_n - 16n + 8}{(n-1)(n-2)} \sim \frac{8 \log n}{n}$$

with $h_n := 1 + \frac{1}{2} + \dots + \frac{1}{n}$, the n -th harmonic number. Below we derive a more general result.

To uncover this peculiarity we shall study the external lengths in more detail. First we look at the point processes η_n on $(0, \infty)$, given by $\eta_n = \sum_{i=1}^n \delta_{\sqrt{n}T_{\rho(i)}}$, i.e.

$$\eta_n(B) := \#\{i : \sqrt{n}T_{\rho(i)} \in B\} \tag{1}$$

for Borel sets $B \subseteq (0, \infty)$.

Theorem 2. *As $n \rightarrow \infty$ the point process η_n converges in distribution, as point processes on $(0, \infty]$, to a Poisson point process η on $(0, \infty)$ with intensity measure $\lambda(dx) = 8x^{-3} dx$.*

We use $(0, \infty]$ in the statement of Theorem 2 instead of $(0, \infty)$ since it is stronger, including for example $\eta_n(a, \infty) \xrightarrow{d} \eta(a, \infty)$ for every $a > 0$. The significance is that, as $n \rightarrow \infty$, there will be points clustering at 0 but not at ∞ . (Below in the proof we recall the definition of convergence in distribution of point processes.) It is not evident, whether there exists a connection to the Poisson point processes introduced by Pitman (1999) for the construction of coalescent processes.

Theorem 2 permits a first orientation. Since $\sqrt{n}L_n = \int x \eta_n(dx)$, one is tempted to resort to infinitely divisible distributions. However, the intensity measure $\lambda(dx)$ is slightly outside the range of the Lévy-Chintchin formula. Shortly speaking this means that small points of η_n have a dominant influence on the distribution of L_n and we are within the domain of the normal distribution.

Thus let us look in more detail on the external lengths and focus on

$$L_n^{\alpha, \beta} := \sum_{n^\alpha \leq \rho(i) < n^\beta} T_{\rho(i)}, \quad 0 \leq \alpha < \beta \leq 1,$$

which is the total length of those external branches having their internal nodes between level $\lceil n^\alpha \rceil$ and $\lceil n^\beta \rceil$ within the coalescent. Obviously $L_n = L_n^{0,1}$.

Proposition 3. *For $0 \leq \alpha < \beta \leq 1$*

$$\mathbf{E}(L_n^{\alpha, \beta}) = \frac{2}{n(n-1)} (\lceil n^\beta \rceil - \lceil n^\alpha \rceil) (2n + 1 - \lceil n^\beta \rceil - \lceil n^\alpha \rceil)$$

and

$$\mathbf{Var}(L_n^{\alpha, \beta}) \sim 8(\beta - \alpha) \frac{\log n}{n},$$

as $n \rightarrow \infty$.

In particular $\mathbf{E}(L_n^{1-\varepsilon, 1}) \sim \mathbf{E}(L_n^{0,1})$, whereas $\mathbf{Var}(L_n^{1-\varepsilon, 1}) \sim \varepsilon \mathbf{Var}(L_n^{0,1})$. Thus the proposition indicates that the systematic part of L_n and its fluctuations arise in different regions of the coalescent tree, the former close to the leaves and the latter closer to the root.

However, this proposition gives an inadequate impression.

Theorem 4. For $0 \leq \alpha < \beta < 1/2$

$$\mathbf{P}(L_n^{\alpha,\beta} = 0) \rightarrow 1$$

as $n \rightarrow \infty$. Moreover

$$\sqrt{n}L_n^{0,\frac{1}{2}} \xrightarrow{d} \int_2^\infty x \eta(dx)$$

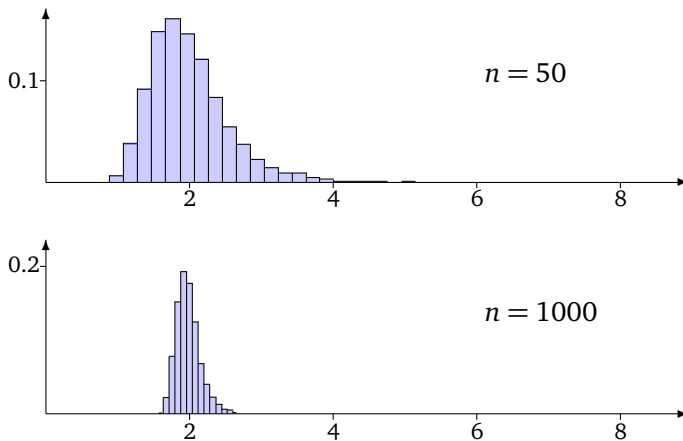
and for $1/2 \leq \alpha < \beta \leq 1$

$$\frac{L_n^{\alpha,\beta} - \mathbf{E}(L_n^{\alpha,\beta})}{\sqrt{\mathbf{Var}(L_n^{\alpha,\beta})}} \xrightarrow{d} N(0, 1).$$

In addition $L_n^{\alpha,\beta}$ and $L_n^{\gamma,\delta}$ are asymptotically independent for $\alpha < \beta \leq \gamma < \delta$.

This result implies Theorem 1: In $L_n = L_n^{0,\frac{1}{2}} + L_n^{\frac{1}{2},1}$ the summands are of order $\sqrt{1/n}$ and $\sqrt{\log n/n}$, such that in the limit the second, asymptotically normal component dominates. To this end, however, n has to become exponentially large, otherwise the few long branches, which make up $L_n^{0,\frac{1}{2}}$, cannot be neglected and may produce extraordinary large values of L_n . Thus the normal approximation for the distribution of L_n seems little useful for practical purposes. One expects a fat right tail compared to the normal distribution. Indeed $\int_2^\infty x \eta(dx)$ has finite mean but infinite variance.

This is illustrated by the following two histograms from 10000 values of L_n , where the length of the horizontal axis to the right indicates the range of the values.



The heavy tails to the right are clearly visible. Also very large outliers appear: For $n = 50$ the simulated values of L_n range from 0.685 to 8.38, and for $n = 1000$ from 1.57 to 7.87.

Also it turns out that the approximation of the variance in Proposition 3 is good only for very large n . This can be seen already from the formula of Fu and Li. To get an exact formula for the variance we look at a somewhat different quantity, namely

$$\hat{L}_n^{\alpha,\beta} := \sum_{i=1}^n (T_{\rho(i)} \wedge T_{\lfloor n^\alpha \rfloor} - T_{\rho(i)} \wedge T_{\lfloor n^\beta \rfloor})$$

with $0 \leq \alpha < \beta \leq 1$, which is the portion of the external length between level $\lfloor n^\alpha \rfloor$ and $\lfloor n^\beta \rfloor$ within the coalescent.

Proposition 5. For $0 \leq \alpha \leq 1$ with $m := \lfloor n^\alpha \rfloor$

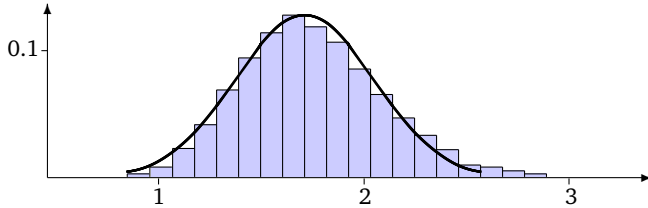
$$\mathbf{E}(\hat{L}_n^{\alpha,1}) = 2 \frac{n-m}{n-1}$$

and

$$\mathbf{Var}(\hat{L}_n^{\alpha,1}) = \frac{8(h_{n-1} - h_{m-1})(n+2m-2)}{(n-1)(n-2)} - \frac{4(n-m)(4n+m-5)}{(n-1)^2(n-2)}.$$

For $\alpha = 0$ we recover the formula of Fu and Li. A similar expression holds for $\hat{L}_n^{\alpha,\beta}$.

Proposition 3 and Theorem 4 carry over to $\hat{L}_n^{\alpha,\beta}$, up to a change in expectation and with the limit $\sqrt{n}\hat{L}_n^{0,\frac{1}{2}} \xrightarrow{d} \int_2^\infty (x-2)\eta(dx)$. The following histogram from a random sample of length 10000 shows that already for $n = 50$ the distribution of $\hat{L}_n^{\frac{1}{2},1}$ fits well to the normal distribution when using the values for expectation and variance, given in Proposition 5.



Our main tool for the proofs is a representation of L_n by means of an imbedded Markov chain U_0, U_1, \dots, U_n , which is of interest of its own. We shall introduce it as an urn model. The relevant fact is that this model possesses an unexpected hidden symmetry, namely it is reversible in time. This is our second main result. For the proof we use another urn model, which allows reversal of time in a simple manner.

The urn models are introduced and studied in Section 2. Proposition 3 is proven in Section 3, Theorems 2 and 4 are derived in Section 4 and Proposition 5 in Section 5.

2 The urn models

Take an urn with n black balls. Empty it in n steps according to the rule: In each step remove a randomly chosen pair of balls and replace it by one red ball. In the last step remove the last remaining ball. Let

$$U_k := \text{number of red balls in the urn after } k \text{ steps}.$$

Obviously $U_0 = U_n = 0$, $U_1 = U_{n-1} = 1$ and $1 \leq U_k \leq \min(k, n-k)$ for $2 \leq k \leq n-2$. U_0, \dots, U_n is a time-inhomogeneous Markov chain with transition probabilities

$$\mathbf{P}(U_{k+1} = u' \mid U_k = u) = \begin{cases} \binom{u}{2} / \binom{n-k}{2}, & \text{if } u' = u - 1, \\ u(n-k-u) / \binom{n-k}{2}, & \text{if } u' = u, \\ \binom{n-k-u}{2} / \binom{n-k}{2}, & \text{if } u' = u + 1. \end{cases}$$

We begin our study of the model by calculating expectations and covariances.

Proposition 6. For $0 \leq k \leq l \leq n$

$$\mathbf{E}(U_k) = \frac{k(n-k)}{n-1}, \quad \mathbf{Cov}(U_k, U_l) = \frac{k(k-1)(n-l)(n-l-1)}{(n-1)^2(n-2)}.$$

Proof. Imagine that the black balls are numbered from 1 to n . Let Z_{ik} be the indicator variable of the event that the black ball with number i is not yet removed after k steps. Then $U_k = n - k - \sum_{i=1}^n Z_{ik}$ and consequently

$$\mathbf{E}(U_k) = n - k - n\mathbf{E}(Z_{1k})$$

and for $k \leq l$ in view of $Z_{1l} \leq Z_{1k}$

$$\begin{aligned} \mathbf{Cov}(U_k, U_l) &= \sum_{i=1}^n \sum_{j=1}^n \mathbf{Cov}(Z_{ik}, Z_{jl}) \\ &= n(n-1)\mathbf{E}(Z_{1k}Z_{2l}) + n\mathbf{E}(Z_{1l}) - n^2\mathbf{E}(Z_{1k})\mathbf{E}(Z_{1l}). \end{aligned}$$

Also

$$\mathbf{E}(Z_{1k}) = \mathbf{P}(Z_{1k} = 1) = \frac{\binom{n-1}{2}}{\binom{n}{2}} \cdots \frac{\binom{n-k}{2}}{\binom{n-k+1}{2}} = \frac{(n-k)(n-k-1)}{n(n-1)}$$

and for $k \leq l$

$$\begin{aligned} \mathbf{E}(Z_{1k}Z_{2l}) &= \mathbf{P}(Z_{1k} = Z_{2l} = 1) = \frac{\binom{n-2}{2}}{\binom{n}{2}} \cdots \frac{\binom{n-k-1}{2}}{\binom{n-k+1}{2}} \cdot \frac{\binom{n-k-1}{2}}{\binom{n-k}{2}} \cdots \frac{\binom{n-l}{2}}{\binom{n-l+1}{2}} \\ &= \frac{(n-k-1)(n-k-2)(n-l)(n-l-1)}{n(n-1)^2(n-2)}. \end{aligned}$$

Our claim now follows by careful calculation. □

Note that these expressions for expectations and covariances are invariant under the transformation $k \mapsto n-k, l \mapsto n-l$. This is not by coincidence:

Theorem 7. (U_0, U_1, \dots, U_n) and $(U_n, U_{n-1}, \dots, U_0)$ are equal in distribution.

Proof. Leaving aside $U_0 = U_n = 0$ we have $U_k \geq 1$ a.s. for the other values of k . Instead we shall look at $U'_k = U_k - 1$ for $1 \leq k \leq n-1$. It turns out that for this process one can specify different dynamics, which are more lucid and amenable to reversing time.

Consider the following alternative box scheme: There are two boxes A and B . At the beginning A contains $n-1$ black balls whereas B is empty. The balls are converted in $2n-2$ steps into $n-1$ red balls lying in B . Namely, in steps number $1, 3, \dots, 2n-3$ a randomly drawn ball from A is shifted to B and in steps number $2, 4, \dots, 2n-2$ a randomly chosen black ball (whether from A or B) is recolored to a red ball. These $2n-2$ operations are carried out independently.

For $1 \leq k \leq n-1$ let

$$U'_k := \text{number of red balls in box } A \text{ after } 2k-1 \text{ steps,}$$

that is at the moment after the k th move and before the k th recoloring. Obviously the sequence is a Markov chain, also $U'_1 = 0$.

As to the transition probabilities note that after $2k - 1$ steps there are $n - k$ black balls in all and $n - k - 1$ balls in A . Thus given $U'_k = r$ there are r red and $n - k - r - 1$ black balls in A , and the remaining $r + 1$ black balls belong to B . Then $U'_{k+1} = r + 1$ occurs only, if in the next step the ball recolored from black to red belongs to A and subsequently the ball shifted from A to B is black. Thus

$$\mathbf{P}(U'_{k+1} = r + 1 \mid U'_k = r) = \frac{n-k-r-1}{n-k} \cdot \frac{n-k-r-2}{n-k-1} = \binom{n-k-r-1}{2} / \binom{n-k}{2}.$$

Similarly $U'_{k+1} = r - 1$ occurs, if the recolored ball belongs to B and next the ball shifted from A to B is red. The corresponding probability is

$$\mathbf{P}(U'_{k+1} = r - 1 \mid U'_k = r) = \frac{r+1}{n-k} \cdot \frac{r}{n-k-1} = \binom{r+1}{2} / \binom{n-k}{2}.$$

Since $U_1 = 1 = U'_1 + 1$ and in view of the transition probabilities of (U_k) and (U'_k) we see that (U_1, \dots, U_{n-1}) and $(U'_1 + 1, \dots, U'_{n-1} + 1)$ indeed coincide in distribution.

Next note that $U'_{n-1} = 0$. Therefore U'_k can be considered as a function not only of the first $2k - 1$ but also of the last $2n - 2k - 1$ shifting and recoloring steps. Since the steps are independent, the process backwards is equally easy to handle. Taking into account that backwards the order of moving and recoloring balls is interchanged, one may just repeat the calculations above to obtain reversibility.

But this repetition can be avoided as well. Let us put our model more formally: Label the balls from 1 to $n - 1$ and write the state space as

$$S := \{((L_1, c_1), \dots, (L_{n-1}, c_{n-1})) \mid L_i \in \{A, B\}, c_i \in \{b, r\}\},$$

where L_i is the location of ball i and c_i its color. Then in our model the first and second coordinate are changed in turn from A to B and from b to r . This is done completely at random, starting within the first coordinates. Clearly we may interchange the role of the first and second coordinate. Thus our box model is equivalent to the following version:

Again initially A contains $n - 1$ black balls whereas B is empty. Now in the steps number $1, 3, \dots, 2n - 3$ a randomly chosen black ball is recolored to a red ball and in the steps number $2, 4, \dots, 2n - 2$ a randomly drawn ball from A is shifted to B . Again these $2n - 2$ operations are carried out independently. Here we consider

$$U''_k := \text{number of black balls in box } B \text{ after } 2k - 1 \text{ steps.}$$

Then from the observed symmetry it is clear that the quantities (U'_1, \dots, U'_{n-1}) and $(U''_1, \dots, U''_{n-1})$ are equal in distribution.

If we finally interchange both colors and boxes as well, then we arrive at the dynamics of the backward process. This finishes the proof. \square

There is a variant of our proof, which makes the reversibility of (U'_k) manifest in a different manner. Let again the balls be labelled from 1 to $n - 1$. Denote

$$\begin{aligned} \nu_m &:= \text{instance between 1 and } n - 1, \text{ when ball } m \text{ is colored to red,} \\ \sigma_m &:= \text{instance between 1 and } n - 1, \text{ when ball } m \text{ is shifted to box } B. \end{aligned}$$

Then from our construction it is clear that $\nu = (\nu_m)$ and $\sigma = (\sigma_m)$ are two independent random permutations of the numbers $\{1, \dots, n - 1\}$. Moreover, at instance k (i.e. after $2k - 1$ steps) ball

number m is red and belongs to box A , if it was colored before and shifted afterwards, i.e. $v_m < k < \sigma_m$. Thus we obtain the formula

$$U'_k = \#\{1 \leq m \leq n-1 : v_m < k < \sigma_m\} \quad (2)$$

and we may conclude the following result.

Corollary 8. *Let v and σ be two independent random permutations of $\{1, \dots, n-1\}$. Then (U_1, \dots, U_{n-1}) is equal in distribution to the process*

$$(\#\{1 \leq m \leq n-1 : v_m < k < \sigma_m\} + 1)_{1 \leq k \leq n-1}.$$

Certainly this representation implies Theorem 7 again. Also it contains additional information. For example, it is immediate that $U_k - 1$ has a hypergeometric distribution with parameters $n-1, k-1, n-k-1$.

One might think to apply similar diminishing urn schemes to other coalescent processes. However, reversibility will hardly be preserved. For related urn models compare the sock-sorting process studied in Steinsaltz (1999) and Janson (2009), Section 8.

We conclude this section by imbedding our urn model into the coalescent. Let

$$V_k := k - \#\{i : \rho(i) < k\}, \quad (3)$$

and $U_k := V_{n-k}$, $0 \leq k \leq n$. Thus V_k is the number of internal branches among the k branches after the $(n-k)$ -th coalescing event and U_k is the number of internal branches among the $n-k$ branches after the k -th coalescing event. The coalescing mechanism takes two random branches and combines them into one internal branch. If we code the external branches by black balls and the internal branches by red, this completely conforms to our urn model; thus (U_0, \dots, U_n) is as above. By Theorem 7, (V_0, \dots, V_n) has the same distribution as (U_0, \dots, U_n) . In the next sections we make use of the Markov chain V_0, \dots, V_n and its properties.

3 Proof of Proposition 3

We use the representation

$$L_n^{\alpha, \beta} = \sum_{n^\alpha \leq k < n^\beta} T_k X_k,$$

where

$$X_k := \#\{i : \rho(i) = k\},$$

$1 \leq k < n$. In view of the coalescing procedure X_k takes only the values 0, 1, 2, and from the definition (3) of V_k

$$X_k = 1 + V_k - V_{k+1}. \quad (4)$$

From (4), $V_k = U_{n-k}$ and Proposition 6 we obtain after simple calculations

$$\mathbf{E}(X_k) = \frac{2k}{n-1}, \quad \mathbf{Var}(X_k) = \frac{2k(n-k-1)(n-3)}{(n-1)^2(n-2)} \quad (5)$$

and for $k < l$

$$\mathbf{Cov}(X_k, X_l) = -\frac{4k(n-l-1)}{(n-1)^2(n-2)}. \quad (6)$$

Also from $T_k = \sum_{j=k+1}^n (T_{j-1} - T_j)$ we have $\mathbf{E}(T_k) = 2 \sum_{j=k+1}^n \frac{1}{(j-1)j}$ and $\mathbf{Var}(T_k) = 4 \sum_{j=k+1}^n \frac{1}{(j-1)^2 j^2}$; thus

$$\mathbf{E}(T_k) = 2\left(\frac{1}{k} - \frac{1}{n}\right), \quad \mathbf{Var}(T_k) \leq \frac{c}{k^3} \quad (7)$$

for a suitable $c > 0$, independent of n .

Thus from independence

$$\mathbf{E}(L_n^{\alpha, \beta}) = \sum_{n^\alpha \leq k < n^\beta} 2\left(\frac{1}{k} - \frac{1}{n}\right) \frac{2k}{n-1}.$$

Now the first claim follows by simple computation.

Further from independence

$$\mathbf{Var}\left(\sum_{n^\alpha \leq k < n^\beta} (T_k - \mathbf{E}(T_k))X_k\right) = \sum_{n^\alpha \leq k, l < n^\beta} \mathbf{Cov}(T_k, T_l)\mathbf{E}(X_k X_l). \quad (8)$$

Using (5)–(7) we have for $k < l$,

$$\mathbf{Cov}(T_k, T_l)\mathbf{E}(X_k X_l) = \mathbf{Var}(T_l)\mathbf{E}(X_k X_l) \leq \mathbf{Var}(T_l)\mathbf{E}(X_k)\mathbf{E}(X_l) \leq \frac{c}{l^3} \cdot \frac{4kl}{(n-1)^2},$$

and it follows that

$$\begin{aligned} 0 \leq \sum_{n^\alpha \leq k < l < n^\beta} \mathbf{Cov}(T_k, T_l)\mathbf{E}(X_k X_l) &\leq \sum_{n^\alpha \leq k < l < n^\beta} \frac{4ck}{l^2}(n-1)^{-2} \\ &\leq \sum_{n^\alpha \leq k < n^\beta} 4c(n-1)^{-2} = O(n^{-1}). \end{aligned}$$

Consequently, (8) yields, using again (5)–(7),

$$\begin{aligned} \mathbf{Var}\left(\sum_{n^\alpha \leq k < n^\beta} (T_k - \mathbf{E}(T_k))X_k\right) &= \sum_{n^\alpha \leq k < n^\beta} \mathbf{Var}(T_k)\mathbf{E}(X_k^2) + O(n^{-1}) \\ &\leq c \sum_{n^\alpha \leq k < n^\beta} \frac{1}{k^3} \left(\frac{2k}{n-1} + \frac{4k^2}{(n-1)^2}\right) + O(n^{-1}) \\ &\leq \frac{6c}{n-1} \sum_{n^\alpha \leq k < n^\beta} \frac{1}{k^2} + O(n^{-1}) = O(n^{-1}). \end{aligned} \quad (9)$$

It remains to show that

$$\mathbf{Var}\left(\sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k)X_k\right) \sim 8(\beta - \alpha) \frac{\log n}{n}.$$

Now

$$\begin{aligned} & \left| \sum_{n^\alpha \leq k < l < n^\beta} \mathbf{E}(T_k) \mathbf{E}(T_l) \mathbf{Cov}(X_k, X_l) \right| \\ & \leq \sum_{n^\alpha \leq k < l < n^\beta} \frac{2}{k} \cdot \frac{2}{l} \cdot \frac{4k}{(n-1)^2} = 16 \sum_{n^\alpha < l < n^\beta} \frac{l - \lceil n^\alpha \rceil}{l(n-1)^2} = O(n^{-1}) \end{aligned}$$

and consequently

$$\begin{aligned} & \mathbf{Var} \left(\sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k) X_k \right) \\ & = \sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k)^2 \mathbf{Var}(X_k) + O(n^{-1}) \\ & = \sum_{n^\alpha \leq k < n^\beta} \frac{4}{k^2} \cdot \frac{2k}{n} \left(1 + O\left(\frac{k}{n}\right) \right) + O(n^{-1}) = 8(\beta - \alpha) \frac{\log n}{n} + O(n^{-1}). \end{aligned}$$

This gives our claim.

4 Proof of Theorems 2 and 4

In this section we use Theorem 7. Namely, V_0, \dots, V_n is a Markov chain with transition probabilities, which can be expressed by means of X_1, \dots, X_{n-1} as follows:

$$\mathbf{P}(X_k = x \mid V_k = v) = \begin{cases} \binom{n-k-v}{2} / \binom{n-k}{2}, & \text{if } x = 0, \\ v(n-k-v) / \binom{n-k}{2}, & \text{if } x = 1, \\ \binom{v}{2} / \binom{n-k}{2}, & \text{if } x = 2. \end{cases}$$

We would like to couple these random variables with suitable independent random variables taking values 0 or 1. Note that V_k takes only values $v \leq k$, thus for $k \leq n/3$

$$\binom{n-k-v}{2} / \binom{n-k}{2} \geq \binom{n-2k}{2} / \binom{n-k}{2} \geq \frac{n-3k}{n-k}.$$

Therefore we may enlarge our model by means of random variables Y_k , $k \leq n/3$, such that

$$\begin{aligned} & \mathbf{P}(X_k = x, Y_k = y \mid V_k = v, V_{k-1}, \dots, V_0, Y_{k-1}, \dots, Y_1) \\ & = \begin{cases} \frac{n-3k}{n-k}, & \text{if } x = 0, y = 0, \\ \binom{n-k-v}{2} / \binom{n-k}{2} - \frac{n-3k}{n-k}, & \text{if } x = 0, y = 1, \\ v(n-k-v) / \binom{n-k}{2}, & \text{if } x = 1, y = 1, \\ \binom{v}{2} / \binom{n-k}{2}, & \text{if } x = 2, y = 1. \end{cases} \end{aligned}$$

For $\mathbf{P}(X_k = x \mid V_k = v)$ this gives the above formula, whereas

$$\mathbf{P}(Y_k = y \mid V_k = v, V_{k-1}, \dots, V_0, Y_{k-1}, \dots, Y_1) = \begin{cases} \frac{n-3k}{n-k}, & \text{if } y = 0, \\ \frac{2k}{n-k}, & \text{if } y = 1. \end{cases}$$

This means that the 0/1-valued random variables Y_k , $k \leq n/3$, are independent. For convenience we put $Y_k = 0$ for $k > n/3$. A straightforward computation gives

$$\mathbf{E}(Y_k - X_k \mid V_k = \nu) = \frac{2(k - \nu)}{n - k}, \quad (10)$$

$$\begin{aligned} \mathbf{E}((Y_k - X_k)^2 \mid V_k = \nu) &= \frac{2(k - \nu)}{n - k} + \frac{2\nu(\nu - 1)}{(n - k)(n - k - 1)} \\ &\leq \frac{2(k - \nu)}{n - k} + \frac{2k(k - 1)}{(n - k)(n - k - 1)} \end{aligned} \quad (11)$$

for $k \leq n/3$. Since $k - \mathbf{E}(V_k) = k(k - 1)/(n - 1)$ from Proposition 6, it follows

$$\mathbf{E}((Y_k - X_k)^2) \leq \frac{4k(k - 1)}{(n - k)(n - k - 1)}. \quad (12)$$

Proof of Theorem 2. Recall that, by (1) and (4),

$$\eta_n = \sum_{i=1}^n \delta_{\sqrt{n}T_{\rho(i)}} = \sum_{k=1}^{n-1} X_k \delta_{\sqrt{n}T_k}. \quad (13)$$

Recall also that $\eta_n \xrightarrow{d} \eta$ as point processes on the interval $(0, \infty]$ means that $\int f d\eta_n \xrightarrow{d} \int f d\eta$ for every continuous f with compact support in $(0, \infty]$, or equivalently $\eta_n(B) \xrightarrow{d} \eta(B)$ for every relatively compact Borel subset B of $(0, \infty]$ such that $\eta(\partial B) = 0$ a.s. (Here B is relatively compact, if $B \subseteq [\delta, \infty]$ for some $\delta > 0$.) See, for example, the Appendix in Janson and Spencer (2007) and Chapter 16 (in particular Theorem 16.16) in Kallenberg (2002).

Let us first look at the point process

$$\eta'_n := \sum_{k=1}^{n-1} Y_k \delta_{2\sqrt{n}/k}. \quad (14)$$

For $0 < a < b \leq \infty$

$$\eta'_n([a, b]) = \sum_{\frac{2\sqrt{n}}{b} < k \leq \frac{2\sqrt{n}}{a}} Y_k$$

and

$$\mathbf{E}(\eta'_n([a, b])) = \sum_{\frac{2\sqrt{n}}{b} < k \leq \frac{2\sqrt{n}}{a}} \frac{2k}{n - k} \rightarrow 4(a^{-2} - b^{-2}) = 8 \int_a^b \frac{dx}{x^3},$$

thus we obtain from standard results on sums of independent 0/1-valued random variables that $\eta'_n([a, b])$ has asymptotically a Poisson distribution. Also $\eta'_n(B_1), \dots, \eta'_n(B_i)$ are independent for disjoint B_1, \dots, B_i . Therefore we obtain from standard results on point processes (for example Kallenberg (2002), Proposition 16.17) weak convergence of η'_n to the Poisson point process η on $(0, \infty]$ with intensity $8x^{-3} dx$.

Next we prove that for all $0 < a < b \leq \infty$

$$\eta_n([a, b]) - \eta'_n([a, b]) \rightarrow 0$$

in probability. To this end note that from (12)

$$\mathbf{E} \left[\sum_{k \leq \frac{2\sqrt{n}}{a}} (Y_k - X_k)^2 \right] = O(n^{-1/2}),$$

which implies that $\mathbf{P}(X_k = Y_k \text{ for all } k \leq \frac{2\sqrt{n}}{a}) \rightarrow 1$. Therefore we may well replace Y_k by X_k in $\eta'_n([a, b])$.

Also, by (7), $\sqrt{n}T_k - 2\sqrt{n}/k = \sqrt{n}T_k - \sqrt{n}\mathbf{E}(T_k) - 2/\sqrt{n}$. From (7) and Doob's inequality for any $\varepsilon > 0$

$$\mathbf{P} \left(\max_{k \geq n^{2/5}} \sqrt{n}|T_k - \mathbf{E}(T_k)| \geq \varepsilon \right) \leq \frac{n}{\varepsilon^2} \mathbf{Var}(T_{\lceil n^{2/5} \rceil}) = O(n^{-1/5}).$$

Since $\mathbf{P}(Y_k = 0 \text{ for all } k < n^{2/5}) \rightarrow 1$, we may as well also replace $2\sqrt{n}/k$ by $\sqrt{n}T_k$ in η'_n , which yields η_n by (13) and (14) (use for example Kallenberg (2002), Theorem 16.16). Thus the proof of Theorem 2 is complete. \square

Proof of Theorem 4. As to the first claim of Theorem 4 observe that the events $\{L_n^{0,\beta} = 0\} = \{X_k = 0 \text{ for all } k < n^\beta\}$ and $\{V_{\lceil n^\beta \rceil} = \lceil n^\beta \rceil\}$ are equal. Thus

$$\begin{aligned} \mathbf{P}(L_n^{\alpha,\beta} > 0) &\leq \mathbf{P}(L_n^{0,\beta} > 0) = \mathbf{P}(\lceil n^\beta \rceil - V_{\lceil n^\beta \rceil} \geq 1) \\ &\leq \mathbf{E}(\lceil n^\beta \rceil - V_{\lceil n^\beta \rceil}) = \frac{\lceil n^\beta \rceil (\lceil n^\beta \rceil - 1)}{n - 1}. \end{aligned}$$

For $\beta < 1/2$ this quantity converges to zero, which gives the first claim of the theorem.

For the next claim we use that because of (7) $\sqrt{n}T_{\lceil n^{1/2} \rceil}$ has expectation $2 + O(n^{-1/2})$ and variance of order $n^{-1/2}$. Thus $\mathbf{P}(2 - \varepsilon < \sqrt{n}T_{\lceil n^{1/2} \rceil} < 2 + \varepsilon) \rightarrow 1$ for all $\varepsilon > 0$. This implies that the probability of the event

$$\begin{aligned} \int_{[2+\varepsilon, \infty)} x \eta_n(dx) &= \sqrt{n} \sum_{k=1}^n T_k X_k I_{\{\sqrt{n}T_k \geq 2+\varepsilon\}} \\ &\leq \sqrt{n} \sum_{k < \sqrt{n}} T_k X_k = \sqrt{n} L_n^{0, \frac{1}{2}} \\ &\leq \sqrt{n} \sum_{k=1}^n T_k X_k I_{\{\sqrt{n}T_k \geq 2-\varepsilon\}} = \int_{[2-\varepsilon, \infty)} x \eta_n(dx) \end{aligned}$$

goes to 1. Also for $a > 0$ from Theorem 2 $\int_a^\infty x \eta_n(dx) \rightarrow \int_a^\infty x \eta(dx)$ in distribution. Altogether we obtain, letting $\varepsilon \rightarrow 0$,

$$\sqrt{n} L_n^{0, \frac{1}{2}} \rightarrow \int_2^\infty x \eta(dx),$$

which is our second claim.

As to the last claim of Theorem 4 we note that from (9)

$$L_n^{\alpha,\beta} = \sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k) X_k + O(n^{-1/2}) \quad (15)$$

meaning that the remainder term is of order $O(n^{-1/2})$ in the L^1 -norm. In this representation, we would like to replace X_k by Y_k . We assume first $\beta < 1$. Note that for $\beta < 1$ in view of (7) and (12)

$$\begin{aligned} & \mathbf{Var}\left(\sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k)(Y_k - X_k - \mathbf{E}(Y_k - X_k | V_k))\right) \\ & \leq \sum_{n^\alpha \leq k < n^\beta} \frac{4}{k^2} \mathbf{E}((Y_k - X_k)^2) = O(n^{\beta-2}) \end{aligned}$$

and from (10), (7) and Proposition 6

$$\begin{aligned} \mathbf{Var}\left(\sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k)\mathbf{E}(Y_k - X_k | V_k)\right) &= \mathbf{Var}\left(\sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k)\frac{2V_k}{n-k}\right) \\ &\leq 2 \sum_{n^\alpha \leq k \leq l < n^\beta} 4 \frac{\mathbf{E}(T_k)\mathbf{E}(T_l)}{(n-k)(n-l)} \mathbf{Cov}(V_k, V_l) \\ &\leq 32 \sum_{n^\alpha \leq k \leq l < n^\beta} \frac{k}{l} \cdot \frac{(n-l)}{(n-k)(n-1)^2(n-2)} = O(n^{2\beta-3}). \end{aligned}$$

Thus $\sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k)((Y_k - X_k) - \mathbf{E}(Y_k - X_k)) = O_p(n^{-1/2})$ and (15) yields

$$L_n^{\alpha, \beta} - \mathbf{E}(L_n^{\alpha, \beta}) = \sum_{n^\alpha \leq k < n^\beta} \mathbf{E}(T_k)(Y_k - \mathbf{E}(Y_k)) + O_p(n^{-1/2}).$$

Also $\mathbf{Var}(\frac{1}{n} \sum_{n^\alpha \leq k < n^\beta} Y_k) \leq n^{-2} \sum_{n^\alpha \leq k < n^\beta} 2k/(n-k) = O(n^{-1})$, and because of (7) we end up with

$$L_n^{\alpha, \beta} - \mathbf{E}(L_n^{\alpha, \beta}) = 2 \sum_{n^\alpha \leq k < n^\beta} \frac{Y_k - \mathbf{E}(Y_k)}{k} + O_p(n^{-1/2}). \quad (16)$$

This is a representation of the external length by a sum of independent random variables.

Now $\mathbf{Var}(Y_k) = \frac{2k}{n-k} - \frac{4k^2}{(n-k)^2}$, thus for $\beta < 1$

$$\begin{aligned} \mathbf{Var}\left(2 \sum_{n^\alpha \leq k < n^\beta} \frac{Y_k - \mathbf{E}(Y_k)}{k}\right) &= 4 \sum_{n^\alpha \leq k < n^\beta} \left(\frac{2}{k(n-k)} - \frac{4}{(n-k)^2}\right) \\ &\sim 8(\beta - \alpha) \frac{\log n}{n}. \end{aligned}$$

Moreover for $\delta > 0$ we have $\mathbf{E}(|Y_k - \mathbf{E}(Y_k)|^{2+\delta}) \leq \frac{2k}{n-k} + (\frac{2k}{n-k})^{2+\delta} \leq \frac{4k}{n-k}$ for $k \leq n/3$, thus

$$\sum_{n^\alpha \leq k < n^\beta} \frac{1}{k^{2+\delta}} \mathbf{E}(|Y_k - \mathbf{E}(Y_k)|^{2+\delta}) \leq 4 \sum_{n^\alpha \leq k < n^\beta} \frac{1}{k^{1+\delta}(n-k)} \leq \frac{8}{\delta n} \frac{1}{(n^\alpha - 1)^\delta}.$$

Thus for $\alpha \geq 1/2$ we get

$$\sum_{n^\alpha \leq k < n^\beta} \frac{1}{k^{2+\delta}} \mathbf{E}(|Y_k - \mathbf{E}(Y_k)|^{2+\delta}) = o\left(\frac{(\log n)^{1+\delta/2}}{n^{1+\delta/2}}\right),$$

and we may use Lyapunov's criterion for the central limit theorem. Consequently, (16) implies

$$\frac{L_n^{\alpha,\beta} - \mathbf{E}(L_n^{\alpha,\beta})}{\sqrt{8(\beta - \alpha)\log n/n}} \xrightarrow{d} N(0, 1).$$

This finishes the proof in the case $\beta < 1$, using Proposition 3.

The case $\beta = 1$ then follows from $L_n^{\alpha,1} = L_n^{\alpha,1-\varepsilon} + L_n^{1-\varepsilon,1}$ using Proposition 3.

The last claim on asymptotic independence follows from (16), too. □

5 Proof of Proposition 5

Let $0 \leq \alpha \leq 1$ and $m = \lfloor n^\alpha \rfloor$. Since $k - V_k = \#\{i : \rho(i) < k\}$ is the number of external branches, which are found between level $k - 1$ and k ,

$$\hat{L}_n^{\alpha,1} = \sum_{m < k \leq n} (T_{k-1} - T_k)(k - V_k).$$

From independence

$$\mathbf{E}(\hat{L}_n^{\alpha,1}) = \sum_{m < k \leq n} \frac{2}{k(k-1)} \cdot \frac{k(k-1)}{n-1}.$$

This gives the first claim. Next, letting

$$E_n := \mathbf{E}(\hat{L}_n^{\alpha,1} | V_0, \dots, V_n) = \sum_{m < k \leq n} \frac{k - V_k}{\binom{k}{2}},$$

we have

$$\mathbf{Var}(\hat{L}_n^{\alpha,1}) = \mathbf{Var}(\hat{L}_n^{\alpha,1} - E_n) + \mathbf{Var}(E_n).$$

Now, using Proposition 6,

$$\begin{aligned} \mathbf{Var}(\hat{L}_n^{\alpha,1} - E_n) &= \sum_{m < k \leq n} \mathbf{E} \left(\left(T_{k-1} - T_k - \frac{1}{\binom{k}{2}} \right)^2 \right) \mathbf{E}((k - V_k)^2) \\ &= \sum_{m < k \leq n} \frac{1}{\binom{k}{2}^2} \left(\frac{k^2(k-1)^2}{(n-1)^2} + \frac{k(k-1)(n-k)(n-k-1)}{(n-1)^2(n-2)} \right) \\ &= 4 \frac{n-m}{(n-1)^2} + 4 \sum_{m < k \leq n} \frac{(n-k)(n-k-1)}{k(k-1)(n-1)^2(n-2)} \end{aligned}$$

and

$$\begin{aligned} \mathbf{Var}(E_n) &= \sum_{m < k, l \leq n} \frac{1}{\binom{k}{2} \binom{l}{2}} \mathbf{Cov}(V_k, V_l) \\ &= 4 \sum_{m < k \leq n} \frac{(n-k)(n-k-1)}{k(k-1)(n-1)^2(n-2)} + 8 \sum_{m < k < l \leq n} \frac{(n-l)(n-l-1)}{l(l-1)(n-1)^2(n-2)} \\ &= 4 \sum_{m < k \leq n} \frac{(n-k)(n-k-1)}{k(k-1)(n-1)^2(n-2)} + 8 \sum_{m < l \leq n} \frac{(l-m-1)(n-l)(n-l-1)}{l(l-1)(n-1)^2(n-2)}. \end{aligned}$$

Thus

$$\mathbf{Var}(\hat{L}_n^{\alpha,1}) = 4 \frac{n-m}{(n-1)^2} + 8 \sum_{m < k \leq n} \frac{(k-m)(n-k)(n-k-1)}{k(k-1)(n-1)^2(n-2)}.$$

Now

$$\begin{aligned} & (k-m)(n-k)(n-k-1) \\ &= (k-1-(m-1))(k(k-1)-2(n-1)k+n(n-1)) \\ &= k(k-1)^2 - (2n+m-3)k(k-1) \\ & \quad + (n+2m-2)(n-1)k - mn(n-1), \end{aligned}$$

thus

$$\begin{aligned} & \frac{1}{2}(n-m)(n-2) + \sum_{m < k \leq n} \frac{(k-m)(n-k)(n-k-1)}{k(k-1)} \\ &= \frac{1}{2}(n-m)(n-2) + \frac{1}{2}(n-m)(n+m-1) - (n-m)(2n+m-3) \\ & \quad + (h_{n-1} - h_{m-1})(n+2m-2)(n-1) - \left(\frac{1}{m} - \frac{1}{n}\right)mn(n-1) \\ &= (h_{n-1} - h_{m-1})(n+2m-2)(n-1) - \frac{1}{2}(n-m)(4n+m-5). \end{aligned}$$

Combining our formulas the result follows. □

References

- [1] Berestycki, J. Berestycki, N. and Schweinsberg, J. (2007) Beta-coalescents and continuous stable random trees. *Ann. Probab.* **35**, 1835–1887. MR2349577
- [2] Berestycki, J. Berestycki, N. and Limic, V. (2011) Asymptotic sampling formulae and particle system representations for Λ -coalescents. arXiv:1101.1875
- [3] Blum, M.G.B. and François, O. (2005) Minimal clade size and external branch length under the neutral coalescent. *Adv. Appl. Prob.* **37**, 647–662. MR2156553
- [4] Caliebe, A., Neininger, R., Krawczak, M. and Rösler, U. (2007) On the length distribution of external branches in coalescent trees: Genetic diversity within species. *Theor. Population Biology* **72**, 245–252.
- [5] Durrett, R. (2002) Probability models for DNA sequence evolution. Probability and its Applications (New York). *Springer-Verlag, New York*. MR1903526
- [6] Freund F and Möhle, M. (2009) On the time back to the most recent common ancestor and the external branch length of the Bolthausen-Sznitman coalescent. *Markov Proc. Rel. Fields* **15**, 387–416. MR2554368
- [7] Fu, Y.X. and Li, W.H. (1993) Statistical tests of neutrality of mutations. *Genetics* **133**, 693–709.

- [8] Gnedin, A., Iksanov, A. and Möhle, M. (2008) On asymptotics of exchangeable coalescents with multiple collisions. *J. Appl. Probab.* **45**, 1186–1195. MR2484170
- [9] Janson, S. (2009) Sorting using complete subintervals and the maximum number of runs in a randomly evolving sequence. *Ann. Comb.* **12**, 417–447. MR2496126
- [10] Janson S. and Spencer J. (2007) A point process describing the component sizes in the critical window of the random graph evolution. *Combin. Probab. Comput.* **16**, 631–658. MR2334588
- [11] Kallenberg, O. (2002) *Foundations of Modern Probability*. 2nd ed., Springer, New York. MR1876169
- [12] Kingman, J.F.C. (1982) The coalescent. *Stoch. Process. Appl.* **13**, 235–248. MR0671034
- [13] Möhle, M. (2010) Asymptotic results for coalescent processes without proper frequencies and applications to the two-parameter Poisson-Dirichlet coalescent. *Stoch. Process. Appl.* **120**, 2159–2173. MR2684740
- [14] Pitman, J. (1999). Coalescents with multiple collisions. *Ann. Probab.* **27**, 1870–1902. MR1742892
- [15] Steinsaltz, D. (1999) Random time changes for sock-sorting and other stochastic process limit theorems. *Electron. J. Probab.* **4**, no. 14, 25 pp. MR1692672