

# THE CRITICAL BETA-SPLITTING RANDOM TREE I: HEIGHTS AND RELATED RESULTS

BY DAVID ALDOUS<sup>1,a</sup> AND BORIS PITTEL<sup>2,b</sup>

<sup>1</sup>*Department of Statistics, University of California, Berkeley, <sup>a</sup>[aldous@stat.berkeley.edu](mailto:aldous@stat.berkeley.edu)*

<sup>2</sup>*Department of Mathematics, The Ohio State University, <sup>b</sup>[pittel.1@osu.edu](mailto:pittel.1@osu.edu)*

In the critical beta-splitting model of a random  $n$ -leaf binary tree, leaf-sets are recursively split into subsets, and a set of  $m$  leaves is split into subsets containing  $i$  and  $m - i$  leaves with probabilities proportional to  $1/i(m - i)$ . We study the continuous-time model in which the holding time before that split is exponential with rate  $h_{m-1}$ , the harmonic number. We (sharply) evaluate the first two moments of the time-height  $D_n$  and of the edge-height  $L_n$  of a uniform random leaf (i.e., the length of the path from the root to the leaf), and prove the corresponding CLTs. We study the correlation between the heights of two random leaves of the same tree realization, and analyze the expected number of splits necessary for a set of  $t$  leaves to partially or completely break away from each other. We give tail bounds for the time-height and the edge-height of the *tree*, that is, the maximal leaf heights. We show that there is a limit distribution for the size of a uniform random subtree, and derive the asymptotics of the mean size. Our proofs are based on asymptotic analysis of the attendant (sum-type) recurrences. The essential idea is to replace such a recursive equality by a pair of recursive inequalities for which matching asymptotic solutions can be found, allowing one to bound, both ways, the elusive explicit solution of the recursive equality. This reliance on recursive inequalities necessitates usage of Laplace transforms rather than Fourier characteristic functions.

**1. Introduction.** The topic of this paper is a certain random tree model, described below, introduced and discussed briefly in 1996 in [3] but not subsequently studied. There is a slight “applied” motivation as a toy model of phylogenetic trees (see Section 1.3), but our purpose here is to commence a theoretical study of the model. A key point is that it has qualitatively different properties from those of the well-studied random tree models that can be found in (for instance) [7, 12].

One fundamental question about a random tree concerns the height of leaves. It turns out that this question, and many extensions, can be answered extremely sharply via an analysis of recursions, and this paper gives a thorough and detailed account of the range of questions that can be answered by this methodology.

In addressing the Applied Probability community, let us observe that there are also other questions about the model that can be studied via a wide range of general modern probability techniques. Some such work in progress is briefly described in a final Section 3, and a current overview can be found in the preprint [4].

1.1. *The model.* For  $m \geq 2$ , consider the distribution  $(q(m, i), 1 \leq i \leq m - 1)$  constructed to be proportional to  $\frac{1}{i(m-i)}$ . Explicitly (by writing  $\frac{1}{i(m-i)} = (\frac{1}{i} + \frac{1}{m-i})/m$ )

$$(1) \quad q(m, i) = \frac{m}{2h_{m-1}} \cdot \frac{1}{i(m-i)}, \quad 1 \leq i \leq m - 1,$$

Received February 2023; revised April 2024.

*MSC2020 subject classifications.* Primary 60C05; secondary 05C05, 92B10.

*Key words and phrases.* Markov chain, phylogenetic tree, random tree, recurrence.

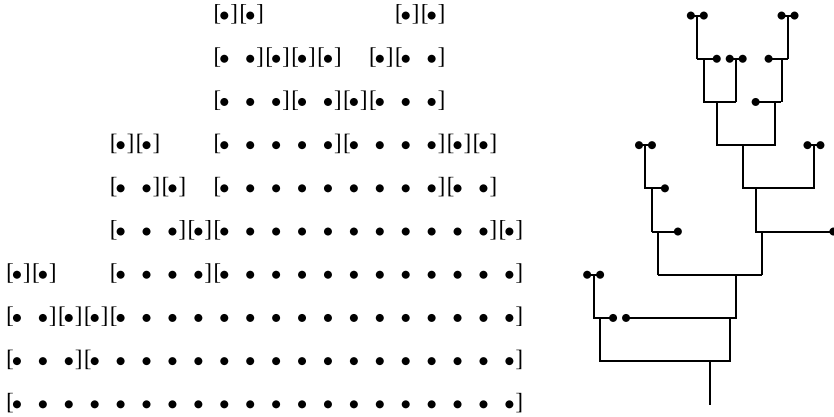


FIG. 1. The discrete time construction for  $n = 20$ . In the tree, by edges we mean the  $n - 1$  vertical edges. The leaves have edge-heights from 2 to 9.

where  $h_{m-1}$  is the harmonic sum  $\sum_{i=1}^{m-1} 1/i$ . Now fix  $n \geq 2$ . Consider the process of constructing a random tree by recursively splitting the integer interval  $[n] = \{1, 2, \dots, n\}$  of “leaves” as follows. First, specify that there is a left edge and a right edge at the root, leading to a left subtree which will have the  $L_n$  leaves  $\{1, \dots, L_n\}$  and a right subtree which will have the  $R_n = n - L_n$  leaves  $\{L_n + 1, \dots, n\}$ , where  $L_n$  (and also  $R_n$ , by symmetry) has distribution  $q(n, \cdot)$ . Recursively, a subinterval with  $m \geq 2$  leaves is split into two subintervals of random size from the distribution  $q(m, \cdot)$ . Continue until reaching intervals of size 1, which are the leaves. This process has a natural tree structure, illustrated schematically<sup>1</sup> in Figure 1. In this discrete-time construction we regard the edges of the tree as having length 1. It turns out (see Section 1.3) to be convenient to consider the continuous-time construction in which a size- $m$  interval is split at rate  $h_{m-1}$ , that is, after an exponential( $h_{m-1}$ ) holding time. Once constructed, it is natural to identify “time” with “distance”: a leaf that appears at time  $t$  has *time-height*  $t$ . Of course the discrete-time model is implicit within the continuous-time model, and a leaf which appears after  $\ell$  splits has *edge-height*  $\ell$ .

We call the continuous-time model the *critical beta-splitting random tree*, but must emphasize that the word *critical* does not have its usual meaning within branching processes. Instead, among the one-parameter family of splitting probabilities with

$$(2) \quad q(m, i) \propto i^\beta (m - i)^\beta, \quad -2 < \beta < \infty$$

our parameter value  $\beta = -1$  is *critical* in the sense that leaf-heights change from order  $n^{-\beta-1}$  to order  $\log n$  at that value, as noted many years ago when this family was introduced [3].

Finally, our results do not use the leaf-labels  $\{1, 2, \dots, n\}$  in the interval-splitting construction. Instead they involve a uniform *random* leaf. Equivalently, one could take a uniform random permutation of labels and then talk about the leaf with some arbitrary label.

1.2. *Outline of results.* Our main focus is on two related random variables associated with the continuous-time random tree on  $n$  leaves:

- $D_n$  = time-height of a uniform random leaf;
- $L_n$  = edge-height of a uniform random leaf.

We start with sharp asymptotic formulas for the moments of  $D_n$  and  $L_n$ . They are of considerable interest in their own right, and also because the techniques are then extended for

<sup>1</sup>Actual simulations appear in [4].

analysis of the limiting distributions, with the moments estimates enabling us to guess what those distributions should be.

Write  $\zeta(\cdot)$  for the Riemann zeta-function,  $\zeta(r) := \sum_{j=1}^{\infty} \frac{1}{j^r}$ , ( $r > 1$ ). Note that  $\zeta(2) = \pi^2/6$  and that  $\zeta^{-1}(2)$  below means  $1/\zeta(2)$ , not the inverse function. Write  $\gamma$  for the Euler–Masceroni constant, which will appear frequently in our analysis:  $\sum_{j=1}^n \frac{1}{j} = \log n + \gamma + O(n^{-1})$ . Asymptotics are as  $n \rightarrow \infty$ .

THEOREM 1.1.

$$\begin{aligned}\mathbb{E}[D_n] &= \zeta^{-1}(2) \log n + O(1), \\ \text{var}(D_n) &= (1 + o(1)) \frac{2\zeta(3)}{\zeta^3(2)} \log n,\end{aligned}$$

and, contingent on a numerically supported “h-ansatz” (see Section 2.2),

$$\mathbb{E}[D_n] = \zeta^{-1}(2) \log n + c_0 - \frac{1}{2\zeta(2)} n^{-1} + O(n^{-2})$$

for a constant  $c_0$  estimated numerically, and

$$\text{var}(D_n) = \frac{2\zeta(3)}{\zeta^3(2)} \log n + O(1).$$

THEOREM 1.2.

$$\begin{aligned}\mathbb{E}[L_n] &= \frac{1}{2\zeta(2)} \log^2 n + \frac{\gamma\zeta(2) + \zeta(3)}{\zeta^2(2)} \log n + O(1), \\ \text{var}(L_n) &= \frac{2\zeta(3)}{3\zeta^3(2)} \log^3 n + O(\log^2 n).\end{aligned}$$

The various parts of Theorem 1.1 are proved in Sections 2.1–2.4 and 2.7, and Theorem 1.2 is proved in Section 2.8. These theorems immediately yield the WLLNs (weak laws of large numbers) for  $D_n$  and  $L_n$ , with rates, as follows.

COROLLARY 1.3. *In probability*

$$\mathbb{P}\left(\left|\frac{D_n}{\mathbb{E}[D_n]} - 1\right| \geq \varepsilon\right), \mathbb{P}\left(\left|\frac{L_n}{\mathbb{E}[L_n]} - 1\right| \geq \varepsilon\right) = O(\varepsilon^{-2} \log^{-1} n).$$

Consider next the time-height  $\mathcal{D}_n$  and the edge-height  $\mathcal{L}_n$  of the random tree itself, that is, the largest time length and the largest edge length of a path from the root to a leaf. By upper-bounding the Laplace transforms of  $D_n$  and  $L_n$ , we prove in Sections 2.6 and 2.9

THEOREM 1.4. *There exists  $\rho > 0$  such that for all  $\varepsilon \in (0, 1)$  we have*

$$\mathbb{P}(\mathcal{D}_n \geq (2 + \varepsilon) \log n) \leq \frac{1}{n^{\rho\varepsilon}},$$

THEOREM 1.5. *Let  $\beta = \min_{\alpha > 1/\log 2} [\alpha + \frac{4\alpha^2\zeta(3)}{\alpha \log 2 - 1}] \approx 42.9$ . For  $\varepsilon \in (0, 1)$ ,*

$$\mathbb{P}(\mathcal{L}_n \geq (1 + \varepsilon)\beta \log^2 n) \leq \exp(-\Theta(\varepsilon \log n)).$$

We conjecture that both  $\frac{D_n}{\log n}$  and  $\frac{L_n}{\log^2 n}$  converge, in probability, to constants.

The definitions of  $D_n$  and  $L_n$  involve two levels of randomness, the random tree and the random leaf within the tree. To study the interaction between levels, it is natural to consider the correlation between the heights of two leaves within the same realization of the random tree. Write  $D_n^{(1)}$  and  $D_n^{(2)}$  for the time-heights of two distinct leaves chosen uniformly from all pairs of leaves. We study the correlation coefficient defined by

$$r_n = \frac{\mathbb{E}[D_n^{(1)} D_n^{(2)}] - \mathbb{E}^2[D_n]}{\text{Var}(D_n)},$$

and prove in Section 2.5

**THEOREM 1.6.** *Contingent on the  $h$ -ansatz,  $r_n = O(\log^{-1} n)$ , that is, asymptotically  $D_n^{(1)}$  and  $D_n^{(2)}$  are uncorrelated.*

We conjecture that a similar result holds for the correlation coefficient of  $L_n^{(1)}$  and  $L_n^{(2)}$ , the edge-heights of two distinct, uniformly random leaves, *independently* of the  $h$ -ansatz.

Returning to properties of  $D_n$  and  $L_n$ , in Sections 2.7 and 2.10 we will prove the CLTs corresponding to the means and variances in Theorems 1.1 and 1.2.

**THEOREM 1.7.** *In distribution, and with all their moments,*

$$\frac{D_n - \zeta^{-1}(2) \log n}{\sqrt{\frac{2\zeta(3)}{\zeta^3(2)} \log n}}, \frac{L_n - (2\zeta(2))^{-1} \log^2 n}{\sqrt{\frac{2\zeta(3)}{3\zeta^3(2)} \log^3 n}} \implies \text{Normal}(0, 1).$$

The sharp asymptotic estimates of the moments of  $D_n$  and  $L_n$ , and the ample numeric evidence in the case of  $D_n$ , provided a compelling evidence that both  $D_n$  and  $L_n$  must be asymptotically normal. However, the proof of Theorem 1.7 does not use these estimates, providing instead an alternative verification of the *leading* terms in those estimates, without relying on the  $h$ -ansatz.

After posting the original preprint version of this article, alternative proofs of these CLTs have appeared in preprints. Via a martingale CLT [4] (continuous model); via the contraction method [13] (discrete model); and via the theory of regenerative composition structures [11] (discrete model). Presumably these methods can also be applied to the alternate model. Of course, in Theorem 1.7 there is presumably *joint convergence* to a bivariate Gaussian limit. It would be interesting to see which method would be best for proving such joint convergence.

Like Theorems 1.4 and 1.5, the proof of Theorem 1.7 is based on showing convergence of the Laplace transform for the (properly centered and scaled) leaf height to that of  $\text{Normal}(0, 1)$ . Why Laplace, but not Fourier? Because, even though there is enough independence to optimistically expect asymptotic normality, our variables are too far from being the sums of essentially independent terms. So, the best we could do is to use recurrences to bound the (real-valued) Laplace transforms recursively both ways, by those of the Normals, whose parameters we choose to satisfy, asymptotically, the respective recursive inequalities. The added feature here is that we get convergence of the moments as well.

Leaving Laplace versus Fourier issue aside, there are many cases when a limited moment information and the recursive nature of the process can be used to establish asymptotic normality, but the standard techniques hardly apply; see [8, 15, 17–20]. The concrete details vary substantially, of course. For instance, in [19] it was shown that the total number of linear extensions of the random, tree-induced, partial order is lognormal, by showing convergence

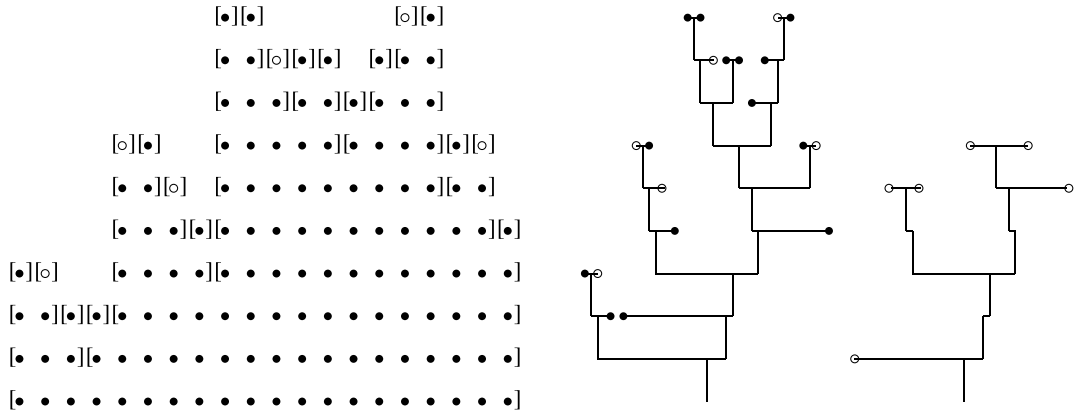


FIG. 2. The left and center diagrams show  $t = 6$  leaves  $\circ$  in the  $n = 20$ -leaf tree in Figure 1. The right diagram is the pruned spanning tree on those leaves, with 8 edges.

of all semi-invariants, rather than of the Laplace transforms. In [20], for the proof of a two-dimensional CLT for the number of vertices and arcs in the giant strong component of the random digraph, boundedness of the Fourier transform made it indispensable. The unifying feature of these diverse arguments is the recurrence equation for the chosen transform.

The structure theory studied in [4] involves the notion of *pruned spanning tree*, illustrated in Figure 2, and here we study its edge-length. Given a set  $T$  of  $t := |T| < n$  leaves of the tree on  $n$  leaves, there is spanning tree on those leaves and the root; the edges of the spanning tree are the union of the edges on the paths to these leaves. Now we can “prune” this spanning tree by cutting the end segment of each path back to the internal vertex  $v$  where it branches from the other paths; the spanning tree on those branchpoints  $v$  forms the pruned spanning tree. Equivalently, the edges of the pruned spanning tree are the edges in the paths from the root to vertices  $v$  such that *each* of the two subtrees rooted at  $v$ ’s children has at least one leaf from  $T$ . Write  $S_{n,t}^*$  for the number of such edges, when  $T$  is a uniform random choice of  $t < n$  leaves. In Section 2.11 we prove the following theorem.

THEOREM 1.8. With  $B(t_1, t_2) = \frac{\Gamma(t_1)\Gamma(t_2)}{\Gamma(t)}$ , we have

$$\mathbb{E}[S_{n,t}^*] = \alpha(t) \log n + O(1), \quad \alpha(t) = \left( h_{t-1} - \sum_{t_1+t_2=t} B(t_1, t_2) \right)^{-1},$$

along with a related result (Proposition 2.14) for the edge-height of the first branch-point in the pruned tree. At long last, the Riemann zeta-function has suddenly loosened its grip, and appropriately the beta-function has taken the stage.

Finally in Section 2.12 we prove

THEOREM 1.9. Let  $\mathbf{u}_n := \{u_n(t)\}_{t \leq n}$  be the distribution of the number of leaves in a subtree rooted at a uniform random vertex, that is, one of the  $2n - 1$  leaves or branchpoints. The sequence  $\{\mathbf{u}_n\}_{n \geq 1}$  converges to a proper distribution  $\mathbf{u}$ . However  $\sum_{t \geq 1} t u_n(t) \sim \frac{3}{2\pi^2} \log^2 n$ .

1.3. *Motivation and background.* The one-parameter family at (2) was introduced in [3] in 1996 as a toy model for phylogenetic trees. It was observed in [5] that, in splits  $m \rightarrow (i, m - i)$  in real-world phylogenetic trees, the median size of the smaller subtree scaled roughly as  $m^{1/2}$ . This data is not consistent with more classical random tree models, where

the median size would be  $O(\log m)$  or  $\Theta(m)$ . However, in the “critical” case  $\beta = -1$  of the model studied in this paper, that median size is indeed order  $m^{1/2}$ , because

$$2 \sum_{i < m^{1/2}} q(m, i) \rightarrow \frac{1}{2} \quad \text{as } m \rightarrow \infty.$$

Of course the model is not a biologically meaningful “forwards in time” model, but is a mathematically basic model that could be used for comparison with more realistic models [14] that are deliberately constructed in the mathematical biology literature to reproduce features of real phylogenetic trees. The distributional results of this model might therefore be useful for some future “applied” project of that kind.

**2. The proofs.** Let  $\tau_\nu$  be the holding time before a split of a subset of size  $\nu$ . So  $\tau_\nu$  has exponential distribution with rate  $h_{\nu-1}$ . By the definition of the splitting process, for  $\nu \geq 2$  we have, with  $q(\nu, i) = \frac{\nu}{2h_{\nu-1}} \frac{1}{i(\nu-i)}$  as at (1),

$$D_\nu = \begin{cases} \tau_\nu + D_i, & \text{with probability } q(\nu, i) \frac{i}{\nu}, i = 1, \dots, \nu - 1, \\ \tau_\nu + D_{\nu-i}, & \text{with probability } q(\nu, i) \frac{\nu-i}{\nu}, i = 1, \dots, \nu - 1. \end{cases}$$

Introduce  $\phi_\nu(u) = \mathbb{E}[e^{u D_\nu}]$ , the Laplace transform of the distribution of  $D_\nu$ ; so,  $\phi_1(u) = 1$ . The equation above implies that for  $\nu \geq 2$ ,

$$\begin{aligned} (3) \quad \phi_\nu(u) &= \sum_{k=1}^{\nu-1} q(\nu, k) \left( \frac{k}{\nu} \mathbb{E}[\exp(u(\tau_\nu + D_k))] + \frac{\nu-k}{\nu} \mathbb{E}[\exp(u(\tau_\nu + D_{\nu-k}))] \right) \\ &= 2 \mathbb{E}[\exp(u\tau_\nu)] \sum_{k=1}^{\nu-1} \frac{k}{\nu} \phi_k(u) q_{\nu,k} = \frac{1}{h_{\nu-1} - u} \sum_{k=1}^{\nu-1} \frac{\phi_k(u)}{\nu - k}. \end{aligned}$$

Furthermore, introduce  $f_\nu(u) = \mathbb{E}[e^{u L_\nu}]$ , the Laplace transform of the distribution of  $L_\nu$ ; so  $f_1(u) = 1$ . In this case we have, for  $\nu \geq 2$ ,

$$L_\nu = \begin{cases} 1 + L_i, & \text{with probability } q(\nu, i) \frac{i}{\nu}, i = 1, \dots, \nu - 1, \\ 1 + L_{\nu-i}, & \text{with probability } q(\nu, i) \frac{\nu-i}{\nu}, i = 1, \dots, \nu - 1. \end{cases}$$

Therefore,

$$\begin{aligned} (4) \quad f_\nu(u) &= \sum_{k=1}^{\nu-1} q(\nu, k) \left( \frac{k}{\nu} \mathbb{E}[\exp(u(1 + L_k))] + \frac{\nu-k}{\nu} \mathbb{E}[\exp(u(1 + L_{\nu-k}))] \right) \\ &= 2e^u \sum_{k=1}^{\nu-1} \frac{k}{\nu} f_k(u) q_{\nu,k} = \frac{e^u}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{f_k(u)}{\nu - k}. \end{aligned}$$

In particular, we make extensive use of the following fundamental recurrence for  $\mathbb{E}[D_\nu]$ :

$$(5) \quad \mathbb{E}[D_\nu] = \frac{1}{h_{\nu-1}} \left( 1 + \sum_{k=1}^{\nu-1} \frac{\mathbb{E}[D_k]}{\nu - k} \right).$$

This follows directly from the hold-jump construction of the random tree, or by differentiating both sides of (3) at  $u = 0$ .

2.1. *The moments of  $D_n$ .* Our first result includes one part of Theorem 1.1.

PROPOSITION 2.1.

$$\zeta^{-1}(2) \log n \leq \mathbb{E}[D_n] \leq \max\{0, 1 + \log(n - 1)\}, \quad n \geq 2,$$

$$\mathbb{E}[D_n] = \zeta^{-1}(2) \log n + O(1).$$

PROOF. The proof has three steps.

(i) Let us prove that  $\mathbb{E}[D_n] \geq \frac{6}{\pi^2} \log n$ . Introduce  $\theta_n = A \log n$ . Then  $\mathbb{E}[D_1] = 0 = \theta_1$ . If we find  $A$  such that

$$(6) \quad \theta_n \leq \frac{1}{h_{n-1}} \left( 1 + \sum_{k=1}^{n-1} \frac{\theta_k}{n-k} \right), \quad n \geq 2,$$

then, by induction on  $n$ ,  $\mathbb{E}[D_n] \geq \theta_n$  for all  $n \geq 1$ . We compute

$$\begin{aligned} \frac{1}{h_{n-1}} \left( 1 + \sum_{k=1}^{n-1} \frac{\theta_k}{n-k} \right) &= \frac{1}{h_{n-1}} \left( 1 + \sum_{k=1}^{n-1} \frac{A \log k}{n-k} \right) \\ &= \frac{1}{h_{n-1}} \left( 1 + A(\log n)h_{n-1} + A \sum_{k=1}^{n-1} \frac{\log(k/n)}{n(1-k/n)} \right) \\ &= \theta_n + \frac{1}{h_{n-1}} \left( 1 + A \sum_{k=1}^{n-1} \frac{\log(k/n)}{n(1-k/n)} \right) \\ &\geq \theta_n + \frac{1}{h_{n-1}} \left( 1 - A \int_0^1 \frac{\log(1/x)}{1-x} dx \right). \end{aligned}$$

The inequality holds since the integrand is positive and decreasing. Since

$$\int_0^1 \frac{\log(1/x)}{1-x} dx = \sum_{j \geq 0} \int_0^1 x^j \log(1/x) dx = \sum_{j \geq 0} \frac{1}{(j+1)^2} = \zeta(2) = \frac{\pi^2}{6},$$

we deduce that (6) holds if we select  $A = \frac{6}{\pi^2} = \zeta^{-1}(2)$ .

NOTE. The proof above is the harbinger of things to come, including the next part. The seemingly naive idea is to replace a recurrence equality by a recurrence *inequality* for which an exact solution can be found and then to use it to upper bound the otherwise-unattainable solution of the recurrence *equality*. Needless to say, it is critically important to have a good guess as to how that “hidden” solution behaves asymptotically.

(ii) Let us prove that  $\mathbb{E}[D_n] \leq f(n) := \max\{0, 1 + \log(n - 1)\}$  for  $n \geq 2$ . This is true for  $n = 1, 2$  since  $\mathbb{E}[D_1] = 0$ ,  $\mathbb{E}[D_2] = 1$ . Notice that  $1 + \log(x - 1) \leq x - 1$  for  $x \in (1, 2]$ . So  $f(x) \leq g(x)$ ,  $\forall x > 1$ , where  $g(x) = x - 1$  for  $x \in [1, 2]$ ,  $g(x) = 1 + \log(x - 1)$  for  $x \geq 2$ , and  $g(x)$  is concave for  $x \geq 1$ . So, similar to (6), it is enough to show that  $g(n)$  satisfies

$$(7) \quad g(n) \geq \frac{1}{h_{n-1}} \left( 1 + \sum_{i=1}^{n-1} \frac{g(i)}{n-i} \right), \quad n \geq 2.$$

By concavity of  $g(x)$  for  $x \geq 1$ , we have

$$\begin{aligned} \frac{1}{h_{n-1}} \left( 1 + \sum_{i=1}^{n-1} \frac{g(i)}{n-i} \right) &\leq \frac{1}{h_{n-1}} + g\left(\sum_{i=1}^{n-1} \frac{i}{n-i}\right) \\ &= \frac{1}{h_{n-1}} + g\left(n - \frac{n-1}{h_{n-1}}\right) \leq \frac{1}{h_{n-1}} + g(n) - g'(n)\left(\frac{n-1}{h_{n-1}}\right), \end{aligned}$$

which is exactly  $g(n)$ , since  $g'(n) = \frac{1}{n-1}$  for  $n > 1$ .

(iii) Write  $\mathbb{E}[D_n] = \frac{6}{\pi^2} \log n + u_n$ , so that  $u_n \geq 0$  and  $u_1 = 0$ . Let us prove that  $u_n = O(1)$ . Using (5) we have

$$\begin{aligned}
 (8) \quad u_n &= \frac{1}{h_{n-1}} \left( 1 + \sum_{k=1}^{n-1} \frac{u_k}{n-k} \right) + \frac{6}{\pi^2} \left( \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{\log k}{n-k} - \log n \right) \\
 &= \frac{1}{h_{n-1}} \left( 1 + \sum_{k=1}^{n-1} \frac{u_k}{n-k} \right) + \frac{6}{\pi^2 h_{n-1}} \sum_{k=1}^{n-1} \frac{\log(k/n)}{n-k}.
 \end{aligned}$$

The proof of (iii) depends on the following rather sharp asymptotic formula for the last sum, which we believe to be new. We defer the proof of the lemma.

LEMMA 2.2.

$$\sum_{k=1}^{n-1} \frac{\log(k/n)}{n-k} = -\zeta(2) + \frac{\log(2\pi e)}{2n} + \frac{\log n}{12n^2} + O(n^{-2}).$$

Granted this estimate, the recurrence (8) becomes

$$\begin{aligned}
 (9) \quad u_n &= \frac{\zeta^{-1}(2)}{h_{n-1}} \left( \frac{\log(2\pi en)}{2n} + \frac{\log n}{12n^2} + O(n^{-2}) \right) \\
 &\quad + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{u_k}{n-k}, \quad n \geq 2, u_1 = 0.
 \end{aligned}$$

It is easy to check that the sequence  $x_n := \frac{n-1}{n}$  satisfies the recurrence

$$x_n = \frac{1}{n} + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{x_k}{n-k}, \quad n \geq 2, x_1 = 0.$$

As the explicit term on the RHS of (9) is asymptotic to  $\frac{\zeta^{-1}(2)}{2n}$ , we can deduce that  $u_n = O(1)$ , establishing (iii). Indeed, by the triangle inequality, the equation (9) implies that

$$|u_n| \leq \frac{c}{n} + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{|u_k|}{n-k}.$$

By induction on  $n$ , this inequality coupled with the recurrence for  $x_n$  imply that  $|u_n| \leq 2cx_n \leq 2c$ .  $\square$

*Proof of Lemma 2.2.* First, we have, for  $n \geq 2$

$$\begin{aligned}
 (10) \quad \sum_{k=1}^{n-1} \frac{\log(k/n)}{n-k} &= \sum_{k=1}^{n-1} \left( \frac{\log(k/n)}{n} + \frac{k \log(k/n)}{n(n-k)} \right) \\
 &= \frac{1}{n} \log \frac{(n-1)!}{n^{n-1}} + \sum_{k=1}^{n-1} \frac{(k/n) \log(k/n)}{n-k}.
 \end{aligned}$$

By Euler's summation formula (Graham, Knuth, and Patashnik [10], (9.78)), if  $f(x)$  is a smooth differentiable function for  $x \in [a, b]$  such that the even derivatives are all of the same



sign, then for every  $m \geq 1$ ,

$$(11) \quad \sum_{a \leq k < b} f(k) = \int_a^b f(x) dx - \frac{1}{2} f(x) \Big|_a^b + \sum_{\ell=1}^m \frac{B_{2\ell}}{(2\ell)!} f^{(2\ell-1)}(x) \Big|_a^b + \theta_m \frac{B_{2m+2}}{(2m+2)!} f^{(2m+1)}(x) \Big|_a^b.$$

Here  $\theta_m \in (0, 1)$  and  $\{B_{2\ell}\}$  are even Bernoulli numbers, defined by  $\frac{z}{e^z-1} = \sum_{\mu \geq 0} B_\mu \frac{z^\mu}{\mu!}$ . The equation (11) was used in [10] to show that

$$\sum_{1 \leq k < n} \log k = n \log n - n + \frac{1}{2} \log \frac{2\pi}{n} + \sum_{\ell=1}^m \frac{B_{2\ell}}{2\ell(2\ell-1)n^{2\ell-1}} + \theta_{m,n} \frac{B_{2m+2}}{(2m+2)(2m+1)n^{2m+1}},$$

$\theta_{m,n} \in (0, 1)$ . Here  $f(x) = \log x$ , so that  $f^{(2\ell)}(x) < 0$  for  $x \geq 1$  and  $\ell \geq 1$ . Using this estimate for  $m = 1$ , we obtain a sharp version of Stirling’s formula:

$$(12) \quad \frac{1}{n} \log \frac{(n-1)!}{n^{n-1}} = -1 + \frac{\log(2\pi n)}{2n} + O(n^{-2}).$$

Consider the sum in the bottom RHS of (10). This time, take  $f(x) = \frac{(x/n)\log(x/n)}{n-x}$ ,  $x \in [1, n-1]$ , and  $f(n) := -\frac{1}{n}$ . Let us show that  $f^{(2\ell)}(x) > 0$  for  $x \in (0, n)$ , or equivalently that  $g^{(2\ell)}(y) > 0$  for  $y \in (0, 1)$ , where  $g(y) := \frac{y \log y}{1-y}$ . We have

$$\begin{aligned} g(y) &= -\log y + \frac{\log y}{1-y} = -\log y - \sum_{j \geq 1} \frac{(1-y)^{j-1}}{j} \\ &= -\log(1-z) - \sum_{j \geq 1} \frac{z^{j-1}}{j}, \quad z := 1-y. \end{aligned}$$

So, we need to show that

$$(-\log(1-z))^{(2\ell)} \geq \left( \sum_{j \geq 1} \frac{z^{j-1}}{j} \right)^{(2\ell)},$$

or equivalently that

$$\frac{(2\ell-1)!}{(1-z)^{2\ell}} \geq \sum_{j > 2\ell} \frac{(j-1)_{2\ell} z^{j-1-2\ell}}{j}.$$

This inequality will follow if we prove a stronger inequality,<sup>2</sup> namely that, for every  $\nu \geq 0$ ,

$$[z^\nu] \frac{(2\ell-1)!}{(1-z)^{2\ell}} \geq [z^\nu] \sum_{j > 2\ell} \frac{(j-1)_{2\ell} z^{j-1-2\ell}}{j}.$$

But this is equivalent to

$$(2\ell + \nu - 1)! \geq \frac{(2\ell + \nu)!}{2\ell + \nu + 1},$$

<sup>2</sup> $[z^\nu]$  denotes the coefficient of  $z^\nu$ .

which is obviously true. Therefore, applying (11), we have, with  $\theta'_{m,n} \in (0, 1)$ ,

$$(13) \quad \begin{aligned} \sum_{k=1}^{n-1} \frac{(k/n) \log(k/n)}{n-k} &= \int_{1/n}^1 g(y) dy - \frac{1}{2n} g(y) \Big|_{1/n}^1 \\ &+ \sum_{\ell=1}^m \frac{B_{2\ell}}{n^{2\ell} (2\ell)!} g^{(2\ell-1)}(y) \Big|_{1/n}^1 \\ &+ \theta'_{m,n} \frac{B_{2m+2}}{n^{2m+2} (2m+2)!} g^{(2m+1)}(y) \Big|_{1/n}^1. \end{aligned}$$

For the first terms in (13)

$$\begin{aligned} \int_{1/n}^1 g(y) dy &= \int_0^1 \frac{y \log y}{1-y} dy - \sum_{j \geq 1} \int_0^{1/n} y^j \log y dy \\ &= -\zeta(2) + 1 + (\log n) \sum_{j \geq 2} n^{-j} j^{-1} + \sum_{j \geq 2} n^{-j} j^{-2}; \\ g(y) \Big|_{1/n}^1 &= -1 + \frac{\log n}{n-1}. \end{aligned}$$

The integrals were evaluated using the more general identities (23) and (24) later.

For the next term in (13) we need  $g^{(2\ell-1)}(y) \Big|_{1/n}^1$ . We use the Newton–Leibniz formula and evaluate  $g^{(2\ell-1)}(1/n)$  and  $g^{(2\ell-1)}(1)$  using respectively

$$\begin{aligned} g^{(2\ell-1)}(y) &= \sum_{j=0}^{2\ell-1} \binom{2\ell-1}{j} (\log y)^{(j)} \left(\frac{y}{1-y}\right)^{(2\ell-1-j)}, \\ g^{(2\ell-1)}(y) &= \sum_{j=0}^{2\ell-1} \binom{2\ell-1}{j} y^{(j)} \left(\frac{\log y}{1-y}\right)^{(2\ell-1-j)}. \end{aligned}$$

In the second sum there are only two nonzero terms, for  $j = 0$  and  $j = 1$ , and using  $\frac{\log y}{1-y} = -\sum_{j \geq 1} \frac{(1-y)^{j-1}}{j}$  we obtain, with some work, that

$$g^{(2\ell-1)}(1) = -\frac{(2\ell-2)!}{2\ell}.$$

For  $g^{(2\ell-1)}(1/n)$ , we use  $\left(\frac{y}{1-y}\right)^{(\mu)} = \left(\frac{1}{1-y}\right)^{(\mu)}$  for  $\mu > 0$ , and after some more protracted work we obtain

$$g^{(2\ell-1)}(1/n) = -(\log n) \frac{(2\ell-1)!}{(1-n^{-1})^{2\ell}} + \sum_{j=1}^{2\ell-2} n^j \cdot \frac{(2\ell-1)_j}{j(1-n^{-1})^{2\ell-j}} + n^{2\ell-2} \cdot \frac{(2\ell-2)!}{1-1/n}.$$

Therefore,

$$\begin{aligned} g^{(2\ell-1)}(y) \Big|_{1/n}^1 &= -\frac{(2\ell-2)!}{2\ell} + (\log n) \frac{(2\ell-1)!}{(1-n^{-1})^{2\ell}} \\ &- \sum_{j=1}^{2\ell-2} n^j \cdot \frac{(2\ell-1)_j}{j(1-n^{-1})^{2\ell-j}} - n^{2\ell-2} \cdot \frac{(2\ell-2)!}{1-1/n}. \end{aligned}$$

This term enters the RHS of (13) with the factor  $n^{-2\ell}$ , making the product of order  $n^{-2}$  regardless of  $m \geq 1$ . And the remainder term in (13) is of order  $n^{-2}$ , again independently of  $m \geq 1$ . So we choose the simplest  $m = 1$ . Collecting all the pieces we transform (13) into

$$(14) \quad \sum_{k=1}^{n-1} \frac{(k/n) \log(k/n)}{n-k} = -\zeta(2) + 1 + \frac{1}{2n} + \frac{\log n}{12n^2} + O(n^{-2}).$$

So, combining (10), (12), and (14), we have

$$\sum_{k=1}^{n-1} \frac{\log(k/n)}{n-k} = -\zeta(2) + \frac{\log(2\pi e)}{2n} + \frac{\log n}{12n^2} + O(n^{-2})$$

which is the assertion of Lemma 2.2.

This completes the proof of Proposition 2.1.

2.2. *An ansatz for sharper results.* Knowing that  $\mathbb{E}[D_n] = \zeta^{-1}(2) \log n + O(1)$ , it seems natural to seek more refined estimates by imagining that

$$\mathbb{E}[D_n] = \zeta^{-1}(2) \log n + \sum_{j \geq 0} c_j n^{-j}$$

almost satisfies the recurrence, and then calculating  $c_j$ . Let us call this the *h-ansatz*, being analogous to a known expansion for  $h_n$ . So to use this ansatz we write

$$w_n := \sum_{j \geq 0} c_j n^{-j}$$

and seek to identify the  $c_j$  from the recurrence (9), which we rewrite as follows:

$$(15) \quad \begin{aligned} w_n &= \frac{d_1 \log n}{nh_{n-1}} + \frac{d_2}{nh_{n-1}} + \frac{1}{h_{n-1}} \sum_{k=2}^{n-1} \frac{w_k}{n-k}, \quad n \geq 2, \\ d_1 &= \frac{\zeta^{-1}(2)}{2}, \quad d_2 = \frac{\zeta^{-1}(2)}{2} \log(2\pi e). \end{aligned}$$

Here

$$\frac{\log n}{h_{n-1}} = 1 - \frac{\gamma}{\log n} + O(\log^{-2} n),$$

where

$$(16) \quad \gamma := 1 - \sum_{j=2}^{\infty} \frac{\zeta(j) - 1}{j} \approx 0.5772156649,$$

is the Euler–Masceroni constant coming from  $h_\nu = \log \nu + \gamma + O(\nu^{-1})$ , [10]. For  $n \geq 3$ , using  $\frac{1}{k(n-k)} = n^{-1}(\frac{1}{k} + \frac{1}{n-k})$ , we have

$$\begin{aligned} \sum_{k=2}^{n-1} \frac{w_k}{n-k} &= \sum_{j \geq 0} c_j \sum_{k=2}^{n-1} \frac{1}{k^j(n-k)} \\ &= c_0 \left( h_{n-1} - \frac{1}{n-1} \right) + c_1 n^{-1} \left( 2h_{n-1} - \frac{n}{n-1} \right) \\ &\quad + n^{-1} \sum_{j \geq 2} c_j \sum_{k=2}^{n-1} \left( \frac{1}{k^j} + \frac{1}{k^{j-1}(n-k)} \right) \end{aligned}$$

$$\begin{aligned}
 &= c_0 \left( h_{n-1} - \frac{1}{n-1} \right) + c_1 n^{-1} \left( 2h_{n-1} - \frac{n}{n-1} \right) \\
 &\quad + n^{-1} \sum_{j \geq 2} c_j (\zeta(j) - 1) + O(n^{-2} \log n).
 \end{aligned}$$

Therefore,

$$\begin{aligned}
 &\frac{d_1 \log n}{nh_{n-1}} + \frac{d_2}{nh_{n-1}} + \frac{1}{h_{n-1}} \sum_{k=2}^{n-1} \frac{w_k}{n-k} - w_n \\
 &= \frac{d_1 + c_1}{n} + \frac{1}{nh_{n-1}} \left( -d_1 \gamma + d_2 - c_0 - c_1 + \sum_{j \geq 2} c_j (\zeta(j) - 1) \right) + O(n^{-2}).
 \end{aligned}$$

So, selecting

$$(17) \quad c_1 = -d_1 = -\frac{3}{\pi^2}, \quad c_0 = d_2 + \sum_{j \geq 2} \left( c_j + \frac{d_1}{j} \right) (\zeta(j) - 1),$$

(as suggested by (16)) we have

$$\frac{d_1 \log n}{nh_{n-1}} + \frac{d_2}{nh_{n-1}} + \frac{1}{h_{n-1}} \sum_{k=2}^{n-1} \frac{w_k}{n-k} - w_n = O(n^{-2}).$$

Therefore,  $w_n = \sum_{j \geq 0} c_j n^{-j}$  satisfies (9) within the additive error  $O(n^{-2})$ , provided that  $\{c_j\}_{j \geq 0}$  satisfies (17). It is worth noticing that  $c_0$  is well defined for every  $\{c_j\}_{j \geq 2}$  provided that the series in (17) converges. The constant  $c_0$  can be viewed as a counterpart of the Euler–Masceroni constant  $\gamma$ . Strikingly,  $c_0$  depends on *all*  $c_j$ ,  $j \geq 2$ , while  $c_1$  is determined uniquely from the requirement that  $w_n$  satisfies (15) within  $O(n^{-2})$  error.

So the conclusion is as follows.

PROPOSITION 2.3. *Assuming the h-ansatz, there exists a constant  $c_0$  such that*

$$(18) \quad \mathbb{E}[D_n] = \frac{6}{\pi^2} \log n + c_0 - \frac{3}{\pi^2} n^{-1} + O(n^{-2}).$$

This is another part of Theorem 1.1. One can calculate  $\mathbb{E}[D_n]$  numerically via the basic recurrence, and doing so up to  $n = 400,000$  gives a good fit<sup>3</sup> to (18) with  $c_0 = 0.7951556604\dots$ . We do not have a conjecture for the explicit value of  $c_0$ .

In what follows, we will use only a weak corollary of (18), namely

$$(19) \quad \mathbb{E}[D_n] = \frac{6}{\pi^2} \log n + c_0 + O(n^{-1}).$$

Paradoxically, the actual value of  $c_0$  will be immaterial as well.

2.3. *The recursion for variance.* Parallel to the recursion (5) for expectations, here is the recursion for variance.

LEMMA 2.4. *Setting  $v_n := \text{var}(D_n)$ , we have*

$$(20) \quad v_n = \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{v_k + (\mathbb{E}[D_n] - \mathbb{E}[D_k])^2}{n-k}.$$

<sup>3</sup>Taking the coefficient of  $n^{-1}$  as unknown, the fit to this data is 0.30408, compared to  $\frac{3}{\pi^2} = 0.30396$ .

PROOF. Differentiating twice both sides of (3) at  $u = 0$ , we get

$$\begin{aligned} \mathbb{E}[D_n^2] &= \frac{2}{h_{n-1}^3} \cdot h_{n-1} + \frac{2}{h_{n-1}^2} \sum_{k=1}^{n-1} \frac{\mathbb{E}[D_k]}{n-k} + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{\mathbb{E}[D_k^2]}{n-k} \\ &= \frac{2}{h_{n-1}^2} \left( 1 + \sum_{k=1}^{n-1} \frac{\mathbb{E}[D_k]}{n-k} \right) + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{\mathbb{E}[D_k^2]}{n-k} \\ &= \frac{2\mathbb{E}[D_n]}{h_{n-1}} + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{\mathbb{E}[D_k^2]}{n-k}. \end{aligned}$$

Since  $v_n = \mathbb{E}[D_n^2] - \mathbb{E}^2[D_n]$ , the equation above becomes

$$v_n = \frac{2\mathbb{E}[D_n]}{h_{n-1}} + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{v_k + \mathbb{E}^2[D_k]}{n-k} - \mathbb{E}^2[D_n].$$

The identity (20) holds because, by (5),

$$\frac{2\mathbb{E}[D_n]}{h_{n-1}} + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{\mathbb{E}^2[D_k]}{n-k} - \mathbb{E}^2[D_n] = \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{(\mathbb{E}[D_n] - \mathbb{E}[D_k])^2}{n-k}. \quad \square$$

NOTE. The equation (46) could be obtained by using the ‘‘law of total variance’’. We preferred the above derivation as more direct in the present context, inconceivable without Laplace transform. Besides, the similar argument will be used later to derive a recurrence for variance of the edge length of the random path. It will be almost the ‘‘same’’ as (20), but with an unexpected, if not shocking, additive term  $-1$  on the RHS.

2.4. *Sharp estimates of  $\text{var}(D_n)$ .* Assuming the h-ansatz, and using (20), we are able to obtain the following sharp estimate, asserted as part of Theorem 1.1.

PROPOSITION 2.5. *Contingent on the h-ansatz,*

$$v_n = \frac{2\zeta(3)}{\zeta^3(2)} \log n + O(1). \quad n \geq 2.$$

NOTE. It is the term  $(\mathbb{E}[D_n] - \mathbb{E}[D_k])^2$  in (20) that necessitates our reliance on the h-ansatz. Comfortingly, the first-order result  $\text{var}(D_n) \sim \frac{2\zeta(3)}{\zeta^3(2)} \log n$  follows from the CLT proof in Section 2.7, independently of the h-ansatz.

PROOF. By (19), we have

$$\begin{aligned} (\mathbb{E}[D_n] - \mathbb{E}[D_k])^2 &= \zeta^{-2}(2)(\log(n/k) + O(k^{-1}))^2 \\ (21) \qquad \qquad \qquad &= \zeta^{-2}(2)(\log^2(n/k) + O(k^{-1} \log(n/k)) + O(k^{-2})). \end{aligned}$$

We need the estimates

$$\begin{aligned} \sum_{k=1}^{n-1} \frac{\log(n/k)}{k(n-k)} &= n^{-1} \sum_{k=1}^{n-1} (k^{-1} + (n-k)^{-1}) \log(n/k) = O(n^{-1} \log^2 n), \\ \sum_{k=1}^{n-1} \frac{1}{k^2(n-k)} &= n^{-1} \sum_{k=1}^{n-1} (k^{-2} + n^{-1}(k^{-1} + (n-k)^{-1})) = O(n^{-1}). \end{aligned}$$

Consider the dominant term in (21). Observe that the function  $\frac{\log^2(n/x)}{n-x}$  is convex. So, using (11) for  $m = 0$ , we obtain

$$\begin{aligned}
 \sum_{k=1}^{n-1} \frac{\log^2(n/k)}{n-k} &= \int_1^n \frac{\log^2(n/x)}{n-x} dx + O(n^{-1} \log^2 n) \\
 (22) \qquad \qquad \qquad &= \int_0^1 \frac{\log^2(1/x)}{1-x} dx + O(n^{-1} \log^2 n) \\
 &= 2\zeta(3) + O(n^{-1} \log^2 n).
 \end{aligned}$$

To explain the final equality, by induction on  $r$  and integrating by parts, we obtain

$$(23) \qquad \int_0^1 z^j \log^r z dz = (-1)^r \frac{r!}{(j+1)^{r+1}}.$$

Consequently

$$(24) \qquad \int_0^1 \frac{\log^r z}{1-z} dz = \int_0^1 (\log^r z) \sum_{j \geq 0} z^j dz = (-1)^r r! \zeta(r+1), \quad r \geq 1$$

used for  $r = 2$  at (22). Now the recursion in Lemma 2.4 becomes

$$v_n = \frac{1}{h_{n-1}} \left( \frac{2\zeta(3)}{\zeta^2(2)} + O(n^{-1} \log^2 n) + \sum_{k=1}^{n-1} \frac{v_k}{n-k} \right).$$

Recalling that

$$\mathbb{E}[D_n] = \frac{1}{h_{n-1}} \left( 1 + \sum_{k=1}^{n-1} \frac{\mathbb{E}[D_k]}{n-k} \right),$$

it follows that  $w_n := |v_n - \frac{2\zeta(3)}{\zeta^2(2)} \mathbb{E}[D_n]|$  satisfies

$$(25) \qquad w_n \leq \frac{1}{h_{n-1}} \left( cn^{-1} \log^2 n + \sum_{k=1}^{n-1} \frac{w_k}{n-k} \right), \quad n \geq 2, w_1 = 0,$$

for some constant  $c > 0$ . Let us prove that the sequence

$$z_n := c \left( \log^2(14) - \frac{\log^2(14n)}{n} \right)$$

satisfies

$$(26) \qquad z_n \geq \frac{1}{h_{n-1}} \left( cn^{-1} \log^2 n + \sum_{k=1}^{n-1} \frac{z_k}{n-k} \right), \quad n \geq 2.$$

Because  $z_1 = 0 = w_1$ , we will get then, predictably by induction using (25), that  $w_n \leq z_n$ .

Let us prove (26). For  $g(x) := -\frac{\log^2(14x)}{x}$ , we have

$$\begin{aligned}
 g'(x) &= x^{-2} (\log^2(14x) - 2 \log(14x)), \\
 g''(x) &= -\frac{2}{x^3} [\log^2(14x) - 3 \log(14x) + 1] < 0, \quad x \geq 1,
 \end{aligned}$$

because  $\log(14) > 2.63 > \frac{3+\sqrt{5}}{2}$ , the larger of two roots of  $x^2 - 3x + 1$ . Therefore,  $g(x)$  is concave on  $[1, \infty)$ . So,

$$\begin{aligned} \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{g(k)}{n-k} &\leq g\left(\frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{k}{n-k}\right) = g\left(n - \frac{n-1}{h_{n-1}}\right) \\ &\leq g(n) - g'(n) \frac{n-1}{h_{n-1}} \\ &= g(n) - n^{-2}(\log^2(14n) - 2\log(14n)) \frac{n-1}{h_{n-1}}. \end{aligned}$$

Since  $z_k = c(\log^2(14) + g(k))$ , we obtain then

$$\begin{aligned} \frac{1}{h_{n-1}} \left( \frac{c \log^2 n}{n} + \sum_{k=1}^{n-1} \frac{z_k}{n-k} \right) \\ \leq z_n + \frac{c}{h_{n-1}} \left[ \frac{\log^2 n}{n} - n^{-2}(n-1)(\log^2(14n) - 2\log(14n)) \right] < z_n, \end{aligned}$$

because the expression within square brackets is easily shown to be negative for  $n \geq 2$ . This establishes (26).  $\square$

*2.5. How correlated are leaf-heights?* Recall the statement of Theorem 1.6, copied below as Theorem 2.6. To study the interaction between the two levels of randomness, it is natural to consider the correlation between leaf heights. Write  $D_n^{(1)}$  and  $D_n^{(2)}$  for the time-heights, within the same realization of the random tree, of two distinct leaves chosen uniformly over pairs of leaves. Both time-heights individually are distributed as  $D_n$ , the time height of the uniformly random leaf. We study the correlation coefficient defined by

$$r_n = \frac{\mathbb{E}[D_n^{(1)} D_n^{(2)}] - \mathbb{E}^2[D_n]}{\text{Var}(D_n)}.$$

**THEOREM 2.6.** *Contingent on the h-ansatz,  $r_n = O(\log^{-1} n)$ .*

**PROOF.** Recall the splitting distribution  $n \rightarrow (L_n, R_n)$  at (1):

$$(27) \quad \mathbb{P}(L_n = i) = q(n, i) = \frac{n}{2h_{n-1}} \frac{1}{i(n-i)} = q(n, n-i), \quad 1 \leq i \leq n-1.$$

There is a natural recursion for  $Z_v := D_v^{(1)} \cdot D_v^{(2)}$ , as follows:

$$(28) \quad Z_v \stackrel{d}{=} \begin{cases} (\tau_v + D_i^{(1)})(\tau_v + D_i^{(2)}), & \text{with probability } q(v, i) \cdot \frac{\binom{i}{2}}{\binom{v}{2}}, \\ (\tau_v + D_{v-i}^{(1)})(\tau_v + D_{v-i}^{(2)}), & \text{with probability } q(v, i) \cdot \frac{\binom{v-i}{2}}{\binom{v}{2}}, \\ (\tau_v + D_i^{(1)})(\tau_v + D_{v-i}^{(2)}), & \text{with probability } q(v, i) \cdot \frac{i \binom{v-i}{2}}{\binom{v}{2}}, \\ (\tau_v + D_i^{(2)})(\tau_v + D_{v-i}^{(1)}), & \text{with probability } q(v, i) \cdot \frac{i \binom{v-i}{2}}{\binom{v}{2}}. \end{cases}$$

Here  $\tau_v$  is the exponential( $h_{v-1}$ ) hold time. The first two cases correspond to the two leaves being in the same subtree, so their heights are dependent, whereas the last two cases correspond to the two leaves being in the different subtrees, so their heights are (conditionally) independent.

Consequently

$$\begin{aligned} \mathbb{E}[Z_\nu | L_\nu = i] &= \left( \frac{2}{h_{\nu-1}^2} + \frac{2}{h_{\nu-1}} \mathbb{E}[D_i] + \mathbb{E}[Z_i] \right) \frac{(i)_2}{(\nu)_2} \\ &\quad + \left( \frac{2}{h_{\nu-1}^2} + \frac{2}{h_{\nu-1}} \mathbb{E}[D_{\nu-i}] + \mathbb{E}[Z_{\nu-i}] \right) \frac{(\nu-i)_2}{(\nu)_2} \\ &\quad + 2 \left( \frac{2}{h_{\nu-1}^2} + \frac{1}{h_{\nu-1}} (\mathbb{E}[D_i] + \mathbb{E}[D_{\nu-i}]) + \mathbb{E}[D_i] \cdot \mathbb{E}[D_{\nu-i}] \right) \frac{i(\nu-i)}{(\nu)_2}, \end{aligned}$$

or, with a bit of algebra,

$$\begin{aligned} \mathbb{E}[Z_\nu | L_\nu = i] &= \frac{2}{h_{\nu-1}^2} + \frac{2i \mathbb{E}[D_i]}{\nu h_{\nu-1}} + \frac{2(\nu-i) \mathbb{E}[D_{\nu-i}]}{\nu h_{\nu-1}} \\ &\quad + \frac{1}{(\nu)_2} ((i)_2 \mathbb{E}[Z_i] + (\nu-i)_2 \mathbb{E}[Z_{\nu-i}] + 2i(\nu-i) \mathbb{E}[D_i] \mathbb{E}[D_{\nu-i}]). \end{aligned}$$

Using (27) we obtain then

$$\begin{aligned} \mathbb{E}[Z_\nu] &= \sum_{i=1}^{\nu-1} q(\nu, i) \mathbb{E}[Z_\nu | L_\nu = i] \\ &= \frac{2}{h_{\nu-1}^2} + \frac{2}{h_{\nu-1}^2} \sum_{i=1}^{\nu-1} \frac{\mathbb{E}[D_i]}{\nu-i} \\ &\quad + \frac{1}{(\nu-1)h_{\nu-1}} \sum_{i=1}^{\nu-1} \mathbb{E}[D_i] \mathbb{E}[D_{\nu-i}] + \frac{1}{(\nu-1)h_{\nu-1}} \sum_{i=1}^{\nu-1} \frac{(i-1) \mathbb{E}[Z_i]}{\nu-i}. \end{aligned}$$

So, using  $\mathbb{E}[D_\nu] = \frac{1}{h_{\nu-1}} (1 + \sum_{i=1}^{\nu-1} \frac{\mathbb{E}[D_i]}{\nu-i})$ , we arrive at

$$\begin{aligned} \mathbb{E}[Z_\nu] &= \frac{1}{(\nu-1)h_{\nu-1}} \sum_{i=1}^{\nu-1} \frac{(i-1) \mathbb{E}[Z_i]}{\nu-i} \\ (29) \quad &\quad + \frac{2 \mathbb{E}[D_\nu]}{h_{\nu-1}} + \frac{1}{(\nu-1)h_{\nu-1}} \sum_{i=1}^{\nu-1} \mathbb{E}[D_i] \mathbb{E}[D_{\nu-i}]. \end{aligned}$$

We use (29) to sharply estimate  $\mathbb{E}[Z_\nu]$  and then estimate  $r_n = \frac{\mathbb{E}[Z_\nu] - \mathbb{E}^2[D_n]}{\text{Var}(D_n)}$ . To start,

$$\frac{2 \mathbb{E}[D_\nu]}{h_{\nu-1}} = 2\zeta^{-1}(2) + O(\log^{-1} \nu).$$

Second,

$$\begin{aligned} \mathbb{E}[D_i] \mathbb{E}[D_{\nu-i}] &= [\zeta^{-1}(2) \log i + c_0 + O(i^{-1})] \\ &\quad \times [\zeta^{-1}(2) \log(\nu-i) + c_0 + O((\nu-i)^{-1})]. \end{aligned}$$

The leading contribution to  $\sum_i \mathbb{E}[D_i] \mathbb{E}[D_{\nu-i}]$  comes from

$$\begin{aligned} &\zeta^{-2}(2) \sum_{i=1}^{\nu-1} \log i \cdot \log(\nu-i) \\ &= \zeta^{-2}(2)(\nu-1) \log^2 \nu + 2\zeta^{-2}(2) \log \nu \sum_{i=1}^{\nu-1} \log(i/\nu) \end{aligned}$$



$$\begin{aligned}
 &+ \zeta^{-2}(2) \sum_{i=1}^{v-1} \log(i/v) \log((v-i)/v) \\
 &= \zeta^{-2}(2)v \log^2 v + 2\zeta^{-2}(2)v \log v \int_0^1 \log x \, dx + O(v) \\
 &= \zeta^{-2}(2)(v \log^2 v - 2v \log v) + O(v).
 \end{aligned}$$

The secondary contribution to  $\sum_i \mathbb{E}[D_i] \mathbb{E}[D_{v-i}]$  comes from  $c_0 \zeta^{-1}(2)(\log i + \log(v-i))$ , and it equals  $2c_0 \zeta^{-1}(2)v \log v + O(v)$ . The terms  $c_0, O(i^{-1}), O((v-i)^{-1})$  contribute jointly  $O(v)$ . Altogether,

$$\sum_{i=1}^{v-1} \mathbb{E}[D_i] \mathbb{E}[D_{v-i}] = \zeta^{-2}(2)(v \log^2 v - 2v \log v) + 2c_0 \zeta^{-1}(2)v \log v + O(v).$$

Therefore, the equation (29) becomes

$$\begin{aligned}
 (30) \quad \mathbb{E}[Z_v] &= \frac{1}{(v-1)h_{v-1}} \sum_{i=1}^{v-1} \frac{(i-1) \mathbb{E}[Z_i]}{v-i} + 2\zeta^{-1}(2) \\
 &+ \zeta^{-2}(2)(\log v - 2 - \gamma) + 2c_0 \zeta^{-1}(2) + O(\log^{-1} v).
 \end{aligned}$$

Let us look at an approximate solution  $\tilde{E}(v) := A \log^2 v + B \log v$ . The RHS of the above equation is

$$\begin{aligned}
 &\frac{1}{(v-1)h_{v-1}} \sum_{i=1}^{v-1} \frac{(i-1)(A \log^2 i + B \log i)}{v-i} + 2\zeta^{-1}(2) \\
 &+ \zeta^{-2}(2)(\log v - 2 - \gamma) + 2c_0 \zeta^{-1}(2) + O(\log^{-1} v).
 \end{aligned}$$

Here, since  $\sum_i \frac{i-1}{v-i} = (v-1)(h_{v-1} - 1)$ , we have

$$\begin{aligned}
 \frac{1}{(v-1)h_{v-1}} \sum_{i=1}^{v-1} \frac{(i-1) \log^2 i}{v-i} &= \frac{1}{(v-1)h_{v-1}} \sum_{i=1}^{v-1} \frac{(i-1)(\log(i/v) + \log v)^2}{v-i} \\
 &= \frac{h_{v-1} - 1}{h_{v-1}} \log^2 v + \frac{2 \log v}{(v-1)h_{v-1}} \sum_{i=1}^{v-1} \frac{(i-1) \log(i/v)}{v-i} \\
 &+ \frac{1}{(v-1)h_{v-1}} \sum_{i=1}^{v-1} \frac{(i-1) \log^2(i/v)}{v-i} \\
 &= \log^2 v - \log v + \gamma + 2 \int_0^1 \frac{x \log x}{1-x} \, dx + O(\log^{-1} v) \\
 &= \log^2 v - \log v + \gamma + 2(1 - \zeta(2)) + O(\log^{-1} v),
 \end{aligned}$$

and

$$\frac{1}{(v-1)h_{v-1}} \sum_{i=1}^{v-1} \frac{(i-1) \log i}{v-i} = \log v - 1 + O(\log^{-1} v).$$

Therefore, with  $\tilde{E}(\cdot)$  instead of  $\mathbb{E}[Z_\cdot]$ , the RHS of the equation (30) becomes

$$\begin{aligned}
 &A(\log^2 v - \log v + \gamma + 2(1 - \zeta(2))) + B(\log v - 1) \\
 &+ \zeta^{-2}(2)(\log v - 2 - \gamma) + 2(c_0 + 1)\zeta^{-1}(2) + O(\log^{-1} v).
 \end{aligned}$$

And we need this to be equal to  $\tilde{E}(v) := A \log^2 v + B \log v$  within an additive error  $O(\log^{-1} v)$ , meaning that

$$-A + B + \zeta^{-2}(2) = B,$$

$$A[\gamma + 2(1 - \zeta(2))] - B - (2 + \gamma)\zeta^{-2}(2) + 2(c_0 + 1)\zeta^{-1}(2) = 0,$$

or explicitly

$$(31) \quad A = \zeta^{-2}(2), \quad B = 2c_0\zeta^{-1}(2).$$

With these  $A$  and  $B$ , our approximation  $\tilde{E}(v)$  satisfies the same equation (30) as  $\mathbb{E}[Z_v]$ , excluding an exact value of the remainder term  $O(\log^{-1} v)$ , of course. Consequently,  $\Delta(v) := |\mathbb{E}[Z_v] - \tilde{E}(v)|$  satisfies

$$(32) \quad \Delta(v) \leq \frac{1}{(v-1)h_{v-1}} \sum_{i=1}^{v-1} \frac{(i-1)\Delta(i)}{v-i} + O(\log^{-1} v), \quad \Delta(1) = 0.$$

With  $\mathcal{U}_v := (v-1)\Delta(v)$ , the resulting equation is a special case of the later equation (65) with the remainder term  $O(v^{t-1} \log^{-1} v)$ , when  $t = 2$ . Applying the bound for the solution proved there, we obtain that  $\mathcal{U}_v = O(v)$ , or that  $\Delta(v) = O(1)$ . Thus

$$\mathbb{E}[Z_v] = A \log^2 v + B \log v + O(1).$$

Combining this formula with (31),  $r_n = \frac{\mathbb{E}[Z_n] - \mathbb{E}^2[D_n]}{\text{Var}(D_n)}$  and  $\mathbb{E}[D_n] = \zeta^{-1}(2) \log n + c_0 + O(n^{-1})$ , we compute

$$r_n = \frac{\zeta^{-2}(2) \log^2 n + 2c_0\zeta^{-1}(2) \log n - (\zeta^{-1}(2) \log n + c_0)^2 + O(1)}{\frac{2\zeta(3)}{\zeta^3(2)} \log n + O(1)} = O(\log^{-1} n). \quad \square$$

NOTE. We do not need the h-ansatz in the rest of the paper.

2.6. *Bounding the time-height of the random tree.* Consider now the time-height  $\mathcal{D}_n$  of the random tree itself, that is, the maximum leaf time-height. We re-state Theorem 1.4, together with a tail bound on  $D_n$ .

PROPOSITION 2.7. (i) For some  $\rho > 0$  and all  $\varepsilon \in (0, \pi^2/6 - 1)$ ,

$$\mathbb{P}\left(D_n \geq \frac{6}{\pi^2}(1 + \varepsilon) \log n\right) = O(n^{-\rho\varepsilon}).$$

(ii) For some  $\rho'$  and all  $\varepsilon \in (0, 1)$ ,

$$\mathbb{P}(\mathcal{D}_n \geq 2(1 + \varepsilon) \log n) = O(n^{-\rho'\varepsilon}).$$

PROOF. (i) Since the tree with  $v$  leaves has  $v - 1$  nonleaf vertices, rather crudely  $D_v$  is stochastically dominated by the sum of  $v - 1$  independent exponentials with rate 1. Therefore, for  $u < 1$ , the Laplace transform  $\phi_v(u) := \mathbb{E}[e^{-uD_v}]$  is bounded above by  $(1 - u)^{-v}$ . Recall (3):

$$\phi_v(u) = \frac{1}{h_{v-1} - u} \sum_{k=1}^{v-1} \frac{\phi_k(u)}{v - k}, \quad v \geq 2.$$

Pick  $\varepsilon' < \varepsilon$  and introduce  $\alpha = \frac{6}{\pi^2}(1 + \varepsilon')$  (so  $\alpha < 1$ ) and  $\psi_v(u) = \exp(u\alpha \log v)$ . Let us prove that

$$(33) \quad \psi_v(u) \geq \frac{1}{h_{v-1} - u} \sum_{k=1}^{v-1} \frac{\psi_k(u)}{v - k},$$

if  $u \in (0, 1)$  is sufficiently small, and  $v > 1$  sufficiently large.

First note that

$$\psi_k(u) = \psi_v(u) \exp(u\alpha \log(k/v)), \quad k \leq v.$$

Therefore,

$$\begin{aligned} \frac{1}{\psi_v(u)(h_{v-1} - u)} \sum_{k=1}^{v-1} \frac{\psi_k(u)}{v-k} &= \frac{1}{h_{v-1} - u} \sum_{k=1}^{v-1} \frac{\exp(u\alpha \log(k/v))}{v-k} \\ &= \left(1 - \frac{u}{h_{v-1}}\right)^{-1} \cdot \left(1 + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\exp(u\alpha \log(k/v)) - 1}{v-k}\right) \\ &= \left(1 + \frac{u}{h_{v-1}} + O\left(\frac{u^2}{h_{v-1}^2}\right)\right) \\ &\quad \times \left[1 + \frac{u}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\alpha \log(k/v)}{v-k} + O\left(\frac{u^2}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\log^2(k/v)}{v-k}\right)\right]; \end{aligned}$$

(where we used  $|e^x - 1 - x| \leq x^2/2$ , for  $x \leq 0$ ). So, since  $\alpha = \zeta^{-1}(2)(1 + \varepsilon')$ ,

$$\begin{aligned} \frac{1}{\psi_v(u)(h_{v-1} - u)} \sum_{k=1}^{v-1} \frac{\psi_k(u)}{v-k} &= 1 + \frac{u}{h_{v-1}} \left(1 + \alpha \sum_{k=1}^{v-1} \frac{\log(k/v)}{v-k}\right) + O\left(\frac{u^2}{h_{v-1}}\right) \\ (34) \quad &\leq 1 + \frac{u}{h_{v-1}} \left(1 + \alpha \left(-\zeta(2) + \frac{\log(v e)}{v-1}\right) + O\left(\frac{u^2}{h_{v-1}}\right)\right) \\ &= 1 - \frac{u}{h_{v-1}} \left(\varepsilon' - \zeta^{-1}(2)(1 + \varepsilon) \frac{\log(v e)}{v-1}\right) + O\left(\frac{u^2}{h_{v-1}}\right). \end{aligned}$$

To justify the inequality above:  $\frac{\log x}{1-x}$  increases for  $x \leq 1$ , so that

$$\begin{aligned} \sum_{k=1}^{v-1} \frac{\log(k/v)}{v-k} &\leq \int_0^1 \frac{\log x}{1-x} dx - \int_0^{1/v} \frac{\log x}{1-x} dx \\ &\leq -\zeta(2) + \frac{v}{v-1} \int_0^{1/v} \log(1/x) dx = -\zeta(2) + \frac{\log(v e)}{v-1}. \end{aligned}$$

The big-O term is uniform over all  $u \in (0, 1)$  and  $v > 1$ . It follows then from (34) that there exist  $u(\varepsilon') \in (0, 1)$  and  $v(\varepsilon') > 1$  such that (33) holds for  $u \in (0, u(\varepsilon'))$  and  $v \geq v(\varepsilon')$ . Furthermore, for  $u \in (0, u(\varepsilon'))$  and  $v \leq v(\varepsilon')$ ,

$$\frac{\phi_v(u)}{\psi_v(u)} \leq A(\varepsilon') := \frac{(1 - u(\varepsilon'))^{-v(\varepsilon')}}{\exp(u(\varepsilon')\alpha \log(v(\varepsilon')))},$$

since, for  $\alpha \leq 1$ ,  $\frac{(1-u)^{-v}}{\exp(u\alpha \log v)}$  attains its maximum on  $[0, v(\varepsilon')]$  at  $v(\varepsilon')$ . Combining this inequality with (33), by induction on  $v$  we obtain that  $\phi_v(u) \leq A(\varepsilon')\psi_v(u)$  for all  $v > 1$  and  $u \leq u' := u(\varepsilon')$ . The rest is easy:

$$\begin{aligned} \mathbb{P}\left(D_n \geq \frac{6}{\pi^2}(1 + \varepsilon) \log n\right) &\leq \frac{\mathbb{E}[\exp(u' D_n)]}{\exp(u' \frac{6}{\pi^2}(1 + \varepsilon) \log n)} \leq \frac{A(\varepsilon')\psi_v(u')}{\exp(u' \frac{6}{\pi^2}(1 + \varepsilon) \log n)} \\ &\leq A(\varepsilon') \exp\left[u' \left(\alpha - \frac{6}{\pi^2}(1 + \varepsilon)\right) \log n\right] = \frac{A(\varepsilon')}{n^{\frac{6u'}{\pi^2}(\varepsilon - \varepsilon')}}. \end{aligned}$$

(ii) Predictably, we will use the union bound, which makes it necessary to upper-bound  $\mathbb{P}(D_n \geq 2(1 + \varepsilon) \log n)$ . To this end, we use a cruder version of the argument in the part (i). Set  $\alpha = 1 + \varepsilon/2$  and choose  $u = \frac{1}{\alpha}$ . Denoting  $z_\nu = u/h_{\nu-1}$  we bound

$$\begin{aligned} \frac{1}{\psi_\nu(u)(h_{\nu-1} - u)} \sum_{k=1}^{\nu-1} \frac{\psi_k(u)}{\nu - k} &= \frac{1}{h_{\nu-1} - u} \sum_{k=1}^{\nu-1} \frac{\exp(u\alpha \log(k/\nu))}{\nu - k} \\ &= \frac{h_{\nu-1}}{h_{\nu-1} - u} \cdot \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{k/\nu}{\nu - k} = \frac{h_{\nu-1}}{h_{\nu-1} - u} \cdot \left(1 - \frac{\nu - 1}{\nu h_{\nu-1}}\right)^{u\alpha} \\ &\leq \exp\left(-\log(1 - z_\nu) - z_\nu \frac{\alpha(\nu - 1)}{\nu}\right). \end{aligned}$$

Since  $z_\nu \rightarrow 0$ , the last expression is below 1 for  $\nu \in [\nu(\alpha), n]$ . Therefore, arguing closely to the part (i), we see that  $\phi_n(u) = O(\psi_n(u))$ . Consequently

$$\mathbb{P}(D_n \geq 2(1 + \varepsilon) \log n) = O\left(\frac{\psi_n(u)}{\exp(2u(1 + \varepsilon) \log n)}\right) = O(n^{-\frac{2(1+\varepsilon)}{1+\varepsilon/2} + 1}),$$

implying, by the union bound, that

$$\mathbb{P}(D_n \geq 2(1 + \varepsilon) \log n) \leq n \mathbb{P}(D_n \geq 2(1 + \varepsilon) \log n) = O(n^{-\frac{2(1+\varepsilon)}{1+\varepsilon/2} + 2}) = O(n^{-\frac{\varepsilon}{1+\varepsilon/2}}). \quad \square$$

2.7. *Asymptotic normality of  $D_n$ .* Here is one part of Theorem 1.7.

PROPOSITION 2.8. *In distribution, and with all of its moments,*

$$\frac{D_n - \zeta^{-1}(2) \log n}{\sqrt{\frac{2\zeta(3)}{\zeta^3(2)} \log n}} \implies \text{Normal}(0, 1).$$

In particular, this provides a proof of the first-order result

$$\text{var}(D_n) \sim \frac{2\zeta(3)}{\zeta^3(2)} \log n,$$

without having to rely on the h-ansatz, as stated in Theorem 1.1.

PROOF. By a general theorem due to Curtis [6], it suffices to show that for  $|u| = \Theta(\log^{-1/2} n)$  and properly chosen  $\alpha_1, \alpha_2 > 0$ , the Laplace transform  $\phi_n(u) = \mathbb{E}[e^{uD_n}]$  satisfies

$$(35) \quad \phi_n(u) = (1 + o(1)) \exp[(u\alpha_1 + u^2\alpha_2) \log n].$$

Recall from (3) that

$$(36) \quad \phi_\nu(u) = \frac{1}{h_{\nu-1} - u} \sum_{k=1}^{\nu-1} \frac{\phi_k(u)}{\nu - k}, \quad \nu \geq 2.$$

Define a function

$$\Psi_\nu(u) = \exp[(u\alpha_1 + u^2\alpha_2) \log \nu], \quad \nu \in [1, n];$$

obviously  $\Psi_1(u) = 1 = \phi_1(u)$ . We will use induction on  $\nu$  to prove a stronger result, namely that there exist  $\alpha_1$  and  $\alpha_2$  such that for  $|u| = \Theta(\log^{-1/2} n)$ , the ratio  $\frac{\phi_\nu(u)}{\Psi_\nu(u)}$  converges to 1,

uniformly over  $n \geq \nu \rightarrow \infty$ , sufficiently fast. Pick  $\delta \in (0, 1/6)$ , and set  $\nu_n = \lceil \exp(\log^\delta n) \rceil$ , so in particular  $u \log \nu_n \rightarrow 0$ . Introduce  $\Psi_\nu^*(u) := 1 + u\alpha \log \nu$ . For  $u > 0$ , we have

$$\begin{aligned}
 (37) \quad & \frac{1}{(h_{\nu-1} - u)\Psi_\nu^*(u)} \sum_{k=1}^{\nu-1} \frac{\Psi_k^*(u)}{\nu - k} \\
 &= \frac{1}{h_{\nu-1}} \sum_{j_1, j_2 \geq 0} u^{j_1 + j_2} \left(\frac{1}{h_{\nu-1}}\right)^{j_1} (-\alpha \log \nu)^{j_2} \\
 &\quad \times \left( (1 + u\alpha \log \nu)h_{\nu-1} + u\alpha \sum_{k=1}^{\nu-1} \frac{\log(k/\nu)}{\nu - k} \right) \\
 &= \left( 1 + u \left( \frac{1}{h_{\nu-1}} - \alpha \log \nu \right) + O(\alpha^2 u^2 \log^2 \nu) \right) \cdot \left( 1 + u\alpha \log \nu - \frac{u}{h_{\nu-1}} \Theta(\alpha) \right) \\
 &= 1 + \frac{u}{h_{\nu-1}} (1 - \Theta(\alpha)) + O((\alpha^2 + 1)u^2 \log^2 \nu) \\
 &\begin{cases} > 1, & \text{if } \nu \leq \nu_n, \alpha > 0 \text{ and small,} \\ < 1, & \text{if } \nu \leq \nu_n, \alpha > 0 \text{ and large.} \end{cases}
 \end{aligned}$$

(For the bottom part we used  $\delta < \frac{1}{6}$ .) And the inequalities are interchanged if  $u < 0$ . Combining this with (36), we conclude that  $\phi_\nu(u) = 1 + O(|u| \log \nu) = \exp(O(|u| \log \nu))$ , uniformly for  $\nu \leq \nu_n$ . So, for bounded  $\alpha_1, \alpha_2$ ,

$$(38) \quad \lim_{n \rightarrow \infty} \max_{\nu \leq \nu_n} \left| \frac{\phi_\nu(u)}{\Psi_\nu(u)} - 1 \right| = 0.$$

Thus, we need to prove existence of  $\alpha_1, \alpha_2$  such that the property above holds for  $\nu \geq \nu_n$ , as well. To this end, let us determine  $\alpha_1$  and  $\alpha_2$  from the condition that  $\Psi_\nu(u)$ , ( $\nu \in [\nu_n, n]$ ), satisfies the recursive inequality

$$(39) \quad \Psi_\nu(u) \geq (\leq) \frac{1}{h_{\nu-1} - u} \left( \sum_{k=1}^{\nu-1} \frac{\Psi_k(u)}{\nu - k} \right), \quad \nu \in [\nu_n, n].$$

First of all, we have

$$\Psi_k(u) = \Psi_\nu(u) \exp[(u\alpha_1 + u^2\alpha_2) \log(k/\nu)], \quad k \leq \nu.$$

Therefore,

$$\begin{aligned}
 (40) \quad & \frac{1}{\Psi_\nu(u)(h_{\nu-1} - u)} \sum_{k=1}^{\nu-1} \frac{\Psi_k(u)}{\nu - k} \\
 &= \frac{1}{h_{\nu-1} - u} \sum_{k=1}^{\nu-1} \frac{\exp[(u\alpha_1 + u^2\alpha_2) \log(k/\nu)]}{\nu - k} \\
 &= \left( 1 - \frac{u}{h_{\nu-1}} \right)^{-1} \cdot \left( 1 + \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\exp[(u\alpha_1 + u^2\alpha_2) \log(k/\nu)] - 1}{\nu - k} \right) \\
 &= \left( 1 - \frac{u}{h_{\nu-1}} \right)^{-1} \\
 &\quad \cdot \left( 1 + \frac{1}{h_{\nu-1}} \int_0^1 \frac{\exp[(u\alpha_1 + u^2\alpha_2) \log x] - 1}{1 - x} dx + O\left(\frac{|u| \log \nu_n}{\nu_n}\right) \right).
 \end{aligned}$$

In the final line, the bottom integral does not depend on  $\nu$ . Let us first justify the remainder term. Define  $f(k/\nu)$  as the  $k$ th term in the previous sum, ( $k < \nu$ ), and, for continuity, set  $f(\nu/\nu) = -\nu^{-1}(u\alpha_1 + u^2\alpha_2)$ . It can be checked that  $f'_k(k/\nu)$  does not change its sign on  $[1, \nu]$ . So, replacing the sum with the integral for  $k$  varying continuously from 1 to  $\nu$ , we introduce the error on the order of the sum of absolute values of

$$f(k/\nu)|_1^\nu \quad \text{and} \quad f'_k(k/\nu)|_1^\nu.$$

The dominant contribution to each of these terms comes from  $k = 1$ . Since for  $\sigma \in (0, 1)$  the function  $\frac{z^\sigma - 1}{z}$  decreases for  $z \geq \sigma^{-1} \log \frac{1}{1-\sigma}$ , we bound

$$|f(1/\nu)| \leq \frac{\exp(|u\alpha_1 + u^2\alpha_2| \log \nu_n) - 1}{\nu_n} = O(\nu_n^{-1} |u| \log \nu_n).$$

And the bound for  $|f'_k(1/\nu)|$  is even better. So the sum in question is of order  $O(\frac{|u| \log \nu_n}{\nu_n})$  uniformly for  $\nu \geq \nu_n$ . Extending the resulting integral to the full  $[0, \nu]$ , we introduce the second error on the order of

$$(41) \quad \int_0^1 \frac{\exp[(u\alpha_1 + u^2\alpha_2) \log(k/\nu)] - 1}{\nu - k} dk = O\left(\frac{|u| \log \nu_n}{\nu_n}\right).$$

The sum of the two error terms is  $O(\nu^{-1} |u| \log \nu)$ , and dividing it by  $h_{\nu-1}$  we get  $O(\frac{|u|}{\nu_n})$ .

Let us sharply estimate the bottom integral in (40). By (41), the contribution to this integral coming from  $x \in (0, 1/\nu_n)$  is  $O(\nu_n^{-1} |u| \log \nu_n)$ . And for  $x \in [1/\nu_n, 1]$ , we have  $|u| \log(1/x) \leq |u| \log \nu_n \rightarrow 0$ , that is, we can use the Taylor expansion

$$\begin{aligned} \frac{\exp[(u\alpha_1 + u^2\alpha_2) \log x] - 1}{1 - x} &= \frac{(u\alpha_1 + u^2\alpha_2) \log x}{1 - x} + \frac{(u\alpha_1 + u^2\alpha_2)^2 \log^2 x}{2(1 - x)} \\ &\quad + O\left(\frac{|u|^3 \log^3(1/x)}{1 - x}\right). \end{aligned}$$

This means that, at the price of the error term of the order  $\nu_n^{-1} |u| \log \nu_n + |u|^3 \int_0^1 \frac{\log^3(1/x)}{1-x} dx$ , we can use the expansion above for all  $x \in (0, 1]$ .

So, using (24), we obtain

$$\begin{aligned} &\frac{1}{h_{\nu-1}} \int_0^1 \frac{\exp[(u\alpha_1 + u^2\alpha_2) \log x] - 1}{1 - x} dx \\ &= -\frac{\alpha_1 \zeta(2)u}{h_{\nu-1}} + \frac{u^2}{h_{\nu-1}} [\alpha_1^2 \zeta(3) - \alpha_2 \zeta(2)] + O\left(\frac{|u|^3}{h_{\nu-1}} + \nu_n^{-1} |u|\right). \end{aligned}$$

Consequently, for  $\nu \geq \nu_n (= \lceil \exp(\log^\delta n) \rceil)$ ,

$$(42) \quad \begin{aligned} \frac{1}{\Psi_\nu(u)(h_{\nu-1} - u)} \sum_{k=1}^{\nu-1} \frac{\Psi_k(u)}{\nu - k} &= 1 + \frac{u}{h_{\nu-1}} (1 - \alpha_1 \zeta(2)) \\ &\quad + \frac{u^2}{h_{\nu-1}} [\alpha_1^2 \zeta(3) - \alpha_2 \zeta(2)] + O\left(\frac{|u|^3}{h_{\nu-1}} + \frac{|u|}{\nu_n}\right) \\ &= 1 + \frac{u}{h_{\nu-1}} (1 - \alpha_1 \zeta(2)) + O\left(\frac{|u|^3}{h_{\nu-1}}\right), \end{aligned}$$

if we select  $\alpha_2 = \frac{\alpha_1^2 \zeta(3)}{\zeta(2)}$ , which we certainly do. Suppose  $u > 0$ ; set  $\alpha_1 = \zeta^{-1}(2) + u^b$ ,  $b \in (1, 2)$ . Then, uniformly for  $\nu \in [\nu_n, n]$ , we have

$$1 + \frac{u}{h_{\nu-1}} (1 - \alpha_1 \zeta(2)) + O\left(\frac{|u|^3}{h_{\nu-1}}\right) = 1 - \frac{\zeta^{-1}(2)u^{b+1}}{h_{\nu-1}} (1 + O(u^{2-b})) < 1.$$

So, (42) becomes

$$\frac{1}{h_{v-1} - u} \sum_{k=1}^{v-1} \frac{\Psi_k(u)}{v - k} \leq \Psi_v(u).$$

This equation and the equation (36) together imply, by induction on  $v \in [v_n, n]$ , that  $\limsup_{n \rightarrow \infty} \max_{v \in [v_n, n]} \frac{\phi_v(u)}{\Psi_v(u)} \leq 1$ . Now,

$$\begin{aligned} \Psi_v(u) &= \exp[(u\alpha_1 + u^2\alpha_2) \log v] \\ &= \exp\left[\left(u\zeta^{-1}(2) + u^2 \frac{\zeta(3)}{\zeta^3(2)}\right) \log v + O(u^{b+1} \log v)\right] \\ &\sim \exp\left[\left(u\zeta^{-1}(2) + u^2 \frac{\zeta(3)}{\zeta^3(2)}\right) \log v\right], \end{aligned}$$

since  $u^{b+1} \log n = O(\log^{-\frac{b-1}{2}} n)$  and  $b > 1$ . Therefore,

$$\limsup_{n \rightarrow \infty} \max_{v \in [v_n, n]} \frac{\phi_v(u)}{\Psi_v(u)} \leq 1.$$

Analogously, setting  $\alpha_1 = \zeta^{-1}(2) - u^b$ , we have

$$\liminf_{n \rightarrow \infty} \min_{v \in [v_n, n]} \frac{\phi_v(u)}{\Psi_v(u)} \geq 1.$$

So, for  $u = \Theta(\log^{-1/2} n) > 0$  we have

$$\lim_{n \rightarrow \infty} \frac{\phi_n(u)}{\exp[(u\zeta^{-1}(2) + u^2 \frac{\zeta(3)}{\zeta^3(2)}) \log n]} = 1.$$

The case  $u < 0$  is completely similar, so that the last equation holds for  $u = -\Theta(\log^{-1/2} n) < 0$  as well.  $\square$

2.8. *The moments of edge-heights of the leaves.* Recall that  $L_n$  denotes the edge-height of a uniform random leaf. In this section we prove Theorem 1.2 via the two propositions below.

PROPOSITION 2.9.

$$(43) \quad \mathbb{E}[L_n] = \frac{1}{2\zeta(2)} \log^2 n + \frac{\gamma\zeta(2) + \zeta(3)}{\zeta^2(2)} \log n + O(1).$$

PROOF. The straightforward recurrence for  $\mathbb{E}[L_v]$  is

$$(44) \quad \mathbb{E}[L_v] = 1 + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\mathbb{E}[L_k]}{v - k}.$$

Write  $\mathbb{E}[L_v] = A \log^2 v + B \log v + u_v$ , so that  $u_1 = 0$ . We need to show that  $u_v = O(1)$ , if we select  $A$  and  $B$  appropriately. (Sure enough, these will be the constants in the claim.) Using (44), we have

$$(45) \quad \begin{aligned} u_v &= 1 + \frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{u_k}{v - k} + A \left( \frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{\log^2 k}{v - k} - \log^2 v \right) \\ &\quad + B \left( \frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{\log k}{v - k} - \log v \right). \end{aligned}$$

Here, by (14),

$$\frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{\log k}{v-k} - \log v = \frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{\log(k/v)}{v-k} = -\frac{\zeta(2)}{h_{v-1}} + \frac{\log(2\pi e)}{vh_{v-1}} + O(v^{-2}),$$

and, combining the equation above with (22), we also have

$$\begin{aligned} & \frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{\log^2 k}{v-k} - \log^2 v \\ &= \frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{\log(k/v) \cdot (\log(k/v) + 2 \log v)}{v-k} \\ &= \frac{2\zeta(3)}{h_{v-1}} + O(v^{-1} \log v) + 2 \left( -\frac{\zeta(2) \log v}{h_{v-1}} + \frac{\log(2\pi e) \log v}{vh_{v-1}} + O(v^{-2} \log v) \right). \end{aligned}$$

Plugging the estimates above into (45) and using  $\log v = h_{v-1} - \gamma + O(v^{-1})$ , we get

$$\begin{aligned} u_v &= \frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{u_k}{v-k} + (1 - 2A\zeta(2)) + \frac{1}{h_{v-1}} [2A(\gamma\zeta(2) + \zeta(3)) - B\zeta(2)] \\ &+ O(v^{-1} \log v). \end{aligned}$$

So, selecting  $A$  and  $B$  such that the  $(A, B)$ -dependent coefficients are both zeros, that is,  $A = \frac{1}{2\zeta(2)}$ ,  $B = \frac{\gamma\zeta(2) + \zeta(3)}{\zeta(2)}$ , we arrive at

$$u_v = \frac{1}{h_{v-1}} \left( \sum_{k \in [v-1]} \frac{u_k}{v-k} + O(v^{-1} \log^2 v) \right).$$

From the proof of Proposition 2.5 (starting with (25)), it follows that  $u_v = O(1)$ .  $\square$

PROPOSITION 2.10.  $\text{var}(L_n) = \frac{2\zeta(3)}{3\zeta^3(2)} \log^3 n + O(\log^2 n)$ .

PROOF. (i) The key is

LEMMA 2.11. *Setting  $\bar{v}_n := \text{var}(L_n)$ , we have*

$$(46) \quad \bar{v}_n = -1 + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{\bar{v}_k + (\mathbb{E}[L_n] - \mathbb{E}[L_k])^2}{n-k}.$$

NOTE. In particular,  $\bar{v}_2 = 0$  as it should be, since  $L_2 \equiv 1$ , unlike  $D_2$  which is distributed exponentially with rate 1.

PROOF. Differentiating twice both sides of (4) at  $u = 0$ , we get

$$\begin{aligned} \mathbb{E}[L_v^2] &= 1 + \frac{1}{h_{v-1}} \sum_{k \in [v-1]} \frac{\mathbb{E}[L_k^2]}{v-k} + \frac{2}{h_{v-1}} \sum_{k \in [v-1]} \frac{\mathbb{E}[L_k]}{v-k} \\ &= 2\mathbb{E}[L_v] - 1 + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{\mathbb{E}[L_k^2]}{n-k} \\ &= 2\mathbb{E}[L_v] - 1 + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\mathbb{E}^2[L_k]}{v-k} + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\bar{v}_k}{v-k}. \end{aligned}$$



Since  $\bar{v}_\nu = \mathbb{E}[L_\nu^2] - \mathbb{E}^2[L_\nu]$ , the equation above becomes

$$\bar{v}_\nu = 2E[L_\nu] - 1 + \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\bar{v}_k + \mathbb{E}^2[L_k]}{\nu - k} - \mathbb{E}^2[L_\nu],$$

and it is easy to check that this equation is equivalent to the claim.  $\square$

(ii) Using Proposition 2.9, we compute, for  $A = \frac{1}{2\zeta(2)}$ ,  $B = \frac{\nu\zeta(2)+\zeta(3)}{\zeta(2)}$ ,

$$\begin{aligned} (\mathbb{E}[L_\nu] - \mathbb{E}[L_k])^2 &= (A(\log^2 \nu - \log^2 k) + B(\log \nu - \log k) + O(1))^2 \\ &= [2A(\log(k/\nu)) \log \nu]^2 + O[\mathcal{P}(\log(\nu/k)) \log \nu], \end{aligned}$$

where  $\mathcal{P}(\eta)$  is a fourth-degree polynomial. Therefore, invoking (22), we have

$$\begin{aligned} \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{(\mathbb{E}[L_\nu] - \mathbb{E}[L_k])^2}{\nu - k} &= \frac{4A^2 \log^2 \nu}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\log^2(k/\nu)}{\nu - k} + O(1) \\ &= \frac{8A^2 \zeta(3) \log^2 \nu}{h_{\nu-1}} + O(1) = 8A^2 \zeta(3) \log \nu + O(1). \end{aligned}$$

So, since  $A = \frac{1}{2\zeta(2)}$ , the equation (46) becomes

$$(47) \quad \bar{v}_\nu = \frac{2\zeta(3)}{\zeta(2)^2} \log \nu + O(1) + \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\bar{v}_k}{\nu - k}.$$

Let us use this recurrence to show that, for appropriately chosen  $A^*$ ,

$$\bar{v}_\nu = \mathcal{V}_\nu + O(\log^2 \nu), \quad \mathcal{V}_\nu := A^* \log^3 \nu.$$

Here  $O(\log^2 \nu)$  is uniform over all  $\nu \geq 2$ . We compute

$$\begin{aligned} \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\log^3 k}{\nu - k} &= \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{(\log \nu + \log(k/\nu))^3}{\nu - k} \\ &= \frac{1}{h_{\nu-1}} \left( \log^3 \nu h_{\nu-1} + 3 \log^2 \nu \sum_{k=1}^{\nu-1} \frac{\log(k/\nu)}{\nu - k} \right. \\ &\quad \left. + 3 \log \nu \sum_{k=1}^{\nu-1} \frac{\log^2(k/\nu)}{\nu - k} + \sum_{k=1}^{\nu-1} \frac{\log^3(k/\nu)}{\nu - k} \right) \\ &= \log^3 \nu + \frac{3 \log^2 \nu}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\log(k/\nu)}{\nu - k} + O(1) \\ &= \log^3 \nu - 3\zeta(2) \log \nu + O(1). \end{aligned}$$

It follows that

$$\begin{aligned} \frac{2\zeta(3)}{\zeta(2)^2} \log \nu + \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\mathcal{V}_k}{\nu - k} \\ = \mathcal{V}_\nu + \left( \frac{2\zeta(3)}{\zeta(2)^2} - 3A^* \zeta(2) \right) \log \nu + O(1) = \mathcal{V}_\nu + O(1), \end{aligned}$$

if we select  $A^* = \frac{2\zeta(3)}{3\zeta^3(2)}$ . Combining this equation with (47), we obtain that  $W_v := \bar{v}_v - \mathcal{V}_v$  satisfies

$$W_v = O(1) + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{W_k}{v-k}.$$

Comparing with (44), we see that  $W_v = O(E[L_v]) = O(\log^2 v)$ .  $\square$

2.9. *Bounding the edge-height of the random tree.* As with the time-height in Section 2.6, we use a tail bound on the edge-height of a random leaf to obtain a tail bound for the edge-height of the tree itself.

PROPOSITION 2.12. *Let  $L_n$  denote the edge-height of the uniformly random leaf, and let  $\mathcal{L}_n$  denote the largest edge-height of a leaf.*

(1) For  $\varepsilon > 0$ ,

$$\mathbb{P}\left(L_n \geq \frac{3}{\pi^2} (1 + \varepsilon) \log^2 n\right) = O(n^{-\Theta(\varepsilon)}).$$

(2) Let  $\beta = \min_{\alpha > 1/\log 2} [\alpha + \frac{4\alpha^2\zeta(3)}{\alpha \log 2 - 1}] \approx 42.9$ . For  $\varepsilon \in (0, 1)$ ,

$$\mathbb{P}(\mathcal{L}_n \geq (1 + \varepsilon)\beta \log^2 n) \leq \exp(-\Theta(\varepsilon \log n)).$$

PROOF. (1) First of all,  $f_v(u) := \mathbb{E}[e^{uL_v}] \leq e^{u(v-1)}$  for  $u \geq 0$ . By (4), we have

$$f_v(u) := \mathbb{E}[e^{uL_v}] = \frac{e^u}{h_{v-1}} \sum_{k=1}^{v-1} \frac{f_k(u)}{v-k}, \quad v \in [2, n].$$

Consider  $u = \frac{v}{\log n}$ , and introduce  $g_k(u) = \exp(u\alpha \log^2 k)$ ,  $\alpha > 0$  yet to be determined. Let us prove that there exists  $v = v(\alpha)$  sufficiently small, and  $v(\alpha)$  sufficiently large such that

$$(48) \quad g_v(u) \geq \frac{e^u}{h_{v-1}} \sum_{k=1}^{v-1} \frac{g_k(u)}{v-k}, \quad \forall v \in [v(\alpha), n].$$

We compute

$$(49) \quad \begin{aligned} \frac{1}{g_v(u)h_{v-1}} \sum_{k=1}^{v-1} \frac{g_k(u)}{v-k} &= \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\exp[u\alpha(\log^2 k - \log^2 v)]}{v-k} \\ &= 1 + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\exp[u\alpha(\log^2 k - \log^2 v)] - 1}{v-k} \\ &\leq 1 + \frac{u\alpha}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\log^2 k - \log^2 v}{v-k} + O\left(u^2 \log v \sum_{k=1}^{v-1} \frac{\log^2(k/v)}{v-k}\right). \end{aligned}$$

The big-Oh term is  $O(u^2 \log v)$ , and

$$\sum_{k=1}^{v-1} \frac{\log^2 k - \log^2 v}{v-k} = \sum_{k=1}^{v-1} \frac{\log^2(k/v)}{v-k} + 2 \log v \sum_{k=1}^{v-1} \frac{\log(k/v)}{v-k} = -\frac{\pi^2}{3} \log v + O(1).$$

Therefore,

$$\frac{1}{g_v(u)h_{v-1}} \sum_{k=1}^{v-1} \frac{g_k(u)}{v-k} \leq 1 - \frac{\alpha u \pi^2}{3} + O(u \log^{-1} v + u^2 \log v),$$

implying that

$$(50) \quad \frac{e^u}{g_\nu(u)h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{g_k(u)}{\nu-k} \leq 1 - u \left( \frac{\alpha\pi^2}{3} - 1 \right) + O(u \log^{-1} \nu + u^2 \log \nu).$$

Recalling that  $u = \frac{\nu}{\log n}$ , we obtain that for  $\alpha > \frac{3}{\pi^2}$  there exists a sufficiently small  $\nu(\alpha) > 0$ , and a sufficiently large  $\nu(\alpha)$ , such that for  $\nu \leq \nu(\alpha)$  and  $n \geq \nu \geq \nu(\alpha)$ , we have

$$(51) \quad \frac{e^u}{g_\nu(u)h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{g_k(u)}{\nu-k} \leq 1.$$

Furthermore, for  $u \leq \frac{\nu(\alpha)}{\log n}$  and  $\nu \leq \nu(\alpha)$ ,

$$(52) \quad \frac{f_\nu(u)}{g_\nu(u)} \leq \frac{\exp(u\nu)}{\exp(u\alpha \log^2 \nu)} \leq \exp\left[\frac{\nu(\alpha)\nu(\alpha)}{\log n}\right].$$

By induction on  $\nu \in [\nu(a), n]$ , it follows that for those  $\nu$ 's

$$f_\nu(u) \leq \exp\left[\frac{\nu(\alpha)\nu(\alpha)}{\log n}\right] g_\nu(u) \implies f_n(u) \leq \exp\left[\frac{\nu(\alpha)\nu(\alpha)}{\log n}\right] g_n(u),$$

provided that  $\alpha > \frac{3}{\pi^2}$ , and  $u \leq \frac{\nu(\alpha)}{\log n}$ . So, given  $\varepsilon > 0$ , we set  $\alpha = \frac{3}{\pi^2}(1 + \varepsilon/2)$ ,  $\alpha' = \frac{3}{\pi^2}(1 + \varepsilon)$ , and bound

$$\mathbb{P}\left(L_n \geq \frac{3}{\pi^2}(1 + \varepsilon) \log^2 n\right) = O\left(\frac{\exp(u\alpha \log^2 n)}{\exp(u\alpha' \log^2 n)}\right) \Big|_{u=\frac{\nu(\alpha)}{\log n}} = O(n^{-\Theta(\varepsilon)}).$$

This is the assertion of (1).

(2) We need a more explicit, but cruder, version of (50), again for  $u = O(\log^{-1} n)$ . Instead of (49), we bound

$$\frac{1}{g_\nu(u)h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{g_k(u)}{\nu-k} \leq 1 + \frac{u\alpha}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\log^2 k - \log^2 \nu}{\nu-k} + \frac{2\alpha^2 u^2 \log^2 \nu}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\log^2(k/\nu)}{\nu-k}.$$

Here

$$\begin{aligned} \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\log^2 k - \log^2 \nu}{\nu-k} &= \frac{1}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\log^2(k/\nu)}{\nu-k} + \frac{2 \log \nu}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\log(k/\nu)}{\nu-k} \\ &\leq \frac{2\zeta(3)}{h_{\nu-1}} + 2 \log \nu \cdot \log\left(1 - \frac{\nu-1}{\nu h_{\nu-1}}\right) \\ &\leq \frac{\zeta(3)}{\log \nu} - 2 \frac{(\nu-1) \log \nu}{h_{\nu-1} \nu} \leq \frac{2\zeta(3)}{\log \nu} - \log 2, \end{aligned}$$

since  $\frac{\log^2 x}{1-x}$  is increasing,  $\log x$  is concave, and  $\log(1+z) \leq z$ , ( $z > -1$ ). So, we replace (50) with

$$(53) \quad \frac{e^u}{g_\nu(u)h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{g_k(u)}{\nu-k} \leq 1 - u \left( \alpha \left( \log 2 - \frac{2\zeta(3)}{\log \nu} \right) - 1 \right) + u^2(1 + 4\alpha^2 \zeta(3) \log \nu).$$

Given  $\alpha > 0$ , the coefficient by  $u$  can be made arbitrarily close to  $\alpha \log 2 - 1$  for  $\nu$  sufficiently large, thus positive if  $\alpha > \frac{1}{\log 2}$ . Assuming the latter, for those large  $\nu$ 's, still below  $n$ , the RHS

expression in (53) is below 1 if  $0 < u \leq \frac{\alpha \log 2 - 1}{(4 + \delta)\alpha^2 \zeta(3) \log n}$ , ( $\delta > 0$ ), and  $n \geq n(\delta)$ . It follows that, as  $n \rightarrow \infty$ , we have  $f_n(u) = O(g_n(u))$ . Consequently, given  $\alpha' > \alpha > \frac{1}{\log 2}$ ,

$$(54) \quad \begin{aligned} \mathbb{P}(\mathcal{L}_n \geq \alpha' \log^2 n) &\leq n \mathbb{P}(L_n \geq \alpha' \log^2 n) = O\left(n \frac{e^{\alpha u \log^2 n}}{e^{\alpha' u \log^2 n}}\right) \\ &= O(\exp[\log n - u(\alpha' - \alpha) \log^2 n]) \rightarrow 0, \end{aligned}$$

if

$$u = \frac{\alpha \log 2 - 1}{(4 + \delta)\alpha^2 \zeta(3) \log n}, \quad \alpha' > \alpha + \frac{(4 + \delta)\alpha^2 \zeta(3)}{\alpha \log 2 - 1}.$$

Set

$$\beta = \min_{\alpha > 1/\log 2} \left[ \alpha + \frac{4\alpha^2 \zeta(3)}{\alpha \log 2 - 1} \right],$$

and let  $\hat{\alpha}$  stand for the point where the minimum is attained. Given  $\varepsilon > 0$ , there exists  $\delta = \delta(\varepsilon) = \Theta(\varepsilon)$  such that

$$\alpha' := (1 + \varepsilon)\beta > \hat{\alpha} + \frac{(4 + \delta)\hat{\alpha}^2 \zeta(3)}{\hat{\alpha} \log 2 - 1},$$

implying, by (54) with  $\alpha = \hat{\alpha}$ , that  $\mathbb{P}(\mathcal{L}_n \geq \alpha' \log^2 n) = O(\exp(-\Theta(\varepsilon) \log n))$ .  $\square$

2.10. *Asymptotic normality of  $L_n$ .* Here is the second part of Theorem 1.7.

PROPOSITION 2.13. *In distribution, and with all of its moments,*

$$\frac{L_n - (2\zeta(2))^{-1} \log^2 n}{\sqrt{\frac{2\zeta(3)}{3\zeta^3(2)} \log^3 n}} \implies \text{Normal}(0, 1).$$

PROOF. Analogously to the proof of Proposition 2.8, it would seem natural to show that for  $|u| = \Theta(\log^{-3/2} n)$  and properly chosen  $\alpha_1 > 0$ ,  $\alpha_2 > 0$ , the Laplace transform  $f_\nu(u) = \mathbb{E}[e^{uL_\nu}]$  satisfies

$$(55) \quad f_\nu(u) = (1 + o(1))g_\nu(u), \quad g_\nu(u) := \exp(u\alpha_1 \log^2 \nu + u^2 \alpha_2 \log^3 \nu),$$

uniformly for  $\nu \leq n$ . But we could only prove (55) for  $0 < u \leq \nu \log^{-3/2} n$  and a fixed  $\nu > 0$ . Fortunately, that is all we need, thanks to a relatively recent extension of the Curtis lemma: it suffices to prove convergence of the sequence of Laplace transforms for the parameter  $\nu$  confined to a fixed interval  $(0, \sigma]$ ; see [16] or [21].

Recall

$$(56) \quad f_\nu(u) = \frac{e^u}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{f_k(u)}{\nu - k}, \quad \nu \geq 2.$$

Pick  $\delta \in (0, 3/4)$  and set  $\nu_n = \lceil \exp(\log^\delta n) \rceil$ . For a constant  $\alpha$ , introduce  $g_\nu^*(u) := \exp(u\alpha \log^2 \nu)$ . For  $u > 0$ , we have

$$\begin{aligned} &\frac{e^u}{h_{\nu-1} g_\nu^*(u)} \sum_{k=1}^{\nu-1} \frac{g_k^*(u)}{\nu - k} \\ &= \frac{e^u}{h_{\nu-1}} \sum_{k=1}^{\nu-1} \frac{\exp(u\alpha \log(k\nu) \log(k/\nu))}{\nu - k} \end{aligned}$$

$$\begin{aligned}
 &= \frac{e^u}{h_{v-1}} \sum_{k=1}^{v-1} [1 + u\alpha \log(kv) \log(k/v) + O(u^2\alpha^2 \log^v \log^2(k/v))] \\
 (57) \quad &= e^u \left[ 1 + \frac{u\alpha\Theta(\log v)}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\log(k/v)}{v-k} + O\left(u^2\alpha^2 \log v \sum_{k=1}^{v-1} \frac{\log^2(k/v)}{v-k}\right) \right] \\
 &= e^u [1 - u\Theta(\alpha) + O(u^2\alpha^2 \log v)] = 1 + u(1 - \Theta(\alpha)) + O(u^2(1 + \alpha^2 \log v)) \\
 &\begin{cases} > 1, & \text{if } v \leq v_n, \alpha > 0 \text{ and small,} \\ < 1, & \text{if } v \leq v_n, \alpha > 0 \text{ and large.} \end{cases}
 \end{aligned}$$

(For the second line we used  $e^z = 1 + z + O(z^2/2)$ , uniformly for  $z < 0$ . For the bottom line we used  $\delta < \frac{3}{4}$ .) Combining this with (56), we conclude that  $f_v(u) = \exp(O(u \log^2 v))$ , uniformly for  $v \leq v_n$ . So, for bounded  $\alpha_1, \alpha_2$ ,

$$(58) \quad \lim_{n \rightarrow \infty} \max_{v \leq v_n} \left| \frac{f_v(u)}{g_v(u)} - 1 \right| = 0.$$

Thus, we need to prove existence of  $\alpha_1, \alpha_2$  such that the analogous relation holds uniformly for all  $v \in [v_n, n]$ . Predictably, we select  $\alpha_1$  and  $\alpha_2$ , requiring that  $g_v(u)$  is the asymptotically best fit for the recurrence (56) for all  $v \in [v_n, n]$ . To begin,

$$\begin{aligned}
 &g_k(u) = g_v(u) \exp[u\alpha_1 G_1(k/v, v) + u^2\alpha_2 G_2(k/v, v)], \\
 (59) \quad &G_1(k/v, v) := \log\left(\frac{k}{v}\right) \log(kv) \leq 0, \\
 &G_2(k/v, v) := \log\left(\frac{k}{v}\right) \left[ 3 \log k \log v + \log^2\left(\frac{k}{v}\right) \right] \leq 0.
 \end{aligned}$$

And, since  $G_j(k/v, v)$  are *nonpositive*, the Taylor expansion of the exponential function holds for  $u > 0$ , even though  $|uG_1(1/n, n)| = O(\sqrt{\log n})$ . (Notice that  $u^2|G_2(1/v, v)| = O(1)$ .) So, proceeding analogously to (40),

$$\begin{aligned}
 &\frac{e^u}{g_v(u)h_{v-1}} \sum_{k=1}^{v-1} \frac{g_k(u)}{v-k} \\
 &= \frac{e^u}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\exp[u\alpha_1 G_1(k/v, v) + u^2\alpha_2 G_2(k/v, v)]}{v-k} \\
 (60) \quad &= e^u \cdot \left( 1 + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{\exp[u\alpha_1 G_1(k/v, v) + u^2\alpha_2 G_2(k/v, v)] - 1}{v-k} \right) \\
 &= e^u \cdot \left( 1 + \frac{1}{h_{v-1}} \int_0^1 \frac{\exp[u\alpha_1 G_1(x, v) + u^2\alpha_2 G_2(x, v)] - 1}{1-x} dx \right. \\
 &\quad \left. + O\left(\frac{\exp(-u \log^2 v)}{v \log v}\right) \right),
 \end{aligned}$$

where *uniformly* for  $x \in (0, 1]$

$$\begin{aligned}
 \frac{\exp[u\alpha_1 G_1(x, v) + u^2\alpha_2 G_2(x, v)] - 1}{1-x} &= \frac{u\alpha_1 G_1(x, v) + u^2\alpha_2 G_2(x, v)}{1-x} \\
 &\quad + \frac{(u\alpha_1 G_1(x, v) + u^2\alpha_2 G_2(x, v))^2}{2(1-x)}
 \end{aligned}$$

$$+ O\left(\frac{u^3 \log^3(1/x) \log^3 v}{1-x}\right).$$

Using (59), and (24), we have then

$$\begin{aligned} & \int_0^1 \frac{\exp[u\alpha_1 G_1(x, v) + u^2\alpha_2 G_2(x, v)] - 1}{1-x} dx \\ &= \alpha_1 u(-2\zeta(2) \log v + 2\zeta(3)) + u^2 \left( \frac{\alpha_1^2}{2} (8\zeta(3) \log^2 v - 24\zeta(4) \log v) \right. \\ & \quad \left. + \alpha_2(-3\zeta(2) \log^2 v + 6\zeta(3) \log v - 6\zeta(4)) \right) + O(u^3 \log^3 v). \end{aligned}$$

Upon expansion  $e^u = 1 + u + u^2/2 + O(u^3)$ , the bottom RHS in (60) then becomes

$$\begin{aligned} & 1 + u \left( 1 + \alpha_1 \frac{-2\zeta(2) \log v + 2\zeta(3)}{h_{v-1}} \right) + u^2 \left( \frac{\frac{\alpha_1^2}{2} (8\zeta(3) \log^2 v - 24\zeta(4) \log v)}{h_{v-1}} \right. \\ (61) \quad & \left. + \frac{\alpha_2(-3\zeta(2) \log^2 v + 6\zeta(3) \log v - 6\zeta(4))}{h_{v-1}} + \alpha_1 \frac{-2\zeta(2) \log v + 2\zeta(3)}{h_{v-1}} \right) \\ & + O(u^3 \log^2 v) \\ & = 1 + u \left( 1 + \alpha_1 \frac{-2\zeta(2) \log v + 2\zeta(3)}{h_{v-1}} \right) + O(u^3 \log^2 v), \end{aligned}$$

if, leaving  $\alpha_1 = \alpha_1(v) > 0$  to be determined shortly, we select  $\alpha_2 = \alpha_2(v)$  to make the coefficient by  $u^2$  equal to zero. Looking closer at the coefficient by  $u^2$ , we see that

$$\alpha_2 = \frac{4\alpha_1^2 \zeta(3)}{3\zeta(2)} + O(\log^{-1} v).$$

The rest is short. Pick  $\alpha_1 = (2\zeta(2))^{-1}(1 + u^b)$ ,  $b < 2\delta/3$ . Since  $u = \Theta(\log^{-3/2} n)$ , the bottom expression in (61) becomes

$$\begin{aligned} & 1 + u(-u^b + O(\log^{-1} v)) + O(u^3 \log^2 v) \\ & = 1 - u^{b+1}(1 + O(u^{-b} \log^{-1} v) + O(u^{2-b} \log^2 n)) \\ & = 1 - u^{b+1}[1 + O((\log n)^{-\delta+3b/2}) + O((\log n)^{-1+3b/2})] < 1. \end{aligned}$$

So, it follows from (60) that

$$\frac{e^u}{h_{v-1}} \sum_{k=1}^{v-1} \frac{g_k(u)}{v-k} < g_v(u), \quad v \in [v_n, n].$$

Combining this recursive inequality with (58), we conclude that

$$\limsup_{n \rightarrow \infty} \max_{v \in [v_n, n]} \frac{f_v(u)}{g_v(u)} \leq 1.$$

Now,

$$\begin{aligned} g_v(u) &= \exp(u\alpha_1 \log^2 v + u^2\alpha_2 \log^3 v) \\ &= \exp\left[ u((2\zeta(2))^{-1}(1 + u^b)) \log^2 v + u^2 \left( \frac{\zeta(3)}{3\zeta^3(2)} + o(1) \right) \log^3 v \right] \\ &= \exp\left[ u(2\zeta(2))^{-1} \log^2 v + u^2 \frac{\zeta(3)}{3\zeta^3(2)} \log^3 v + o(1) + O(u^{b+1} \log^2 v) \right] \end{aligned}$$

$$= (1 + o(1)) \exp \left[ u(2\zeta(2))^{-1} \log^2 v + u^2 \frac{\zeta(3)}{3\zeta^3(2)} \log^3 v \right],$$

if we select  $b > 1/3$ . Since  $b < 2\delta/3$ , a desired  $b$  exists provided that  $\delta > 1/2$ , the constraint compatible with the initial restriction  $\delta < 3/4$ . We conclude that for  $\delta \in (1/2, 3/4)$

$$\limsup_{n \rightarrow \infty} \max_{v \in [v_n, n]} f_v(u) \exp \left[ -u(2\zeta(2))^{-1} \log^2 v - u^2 \frac{\zeta(3)}{3\zeta^3(2)} \log^3 v \right] \leq 1.$$

Likewise, picking  $\alpha_1 = (2\zeta(2))^{-1}(1 - u^b)$ ,  $b < 2\delta/3$ , we obtain

$$\liminf_{n \rightarrow \infty} \min_{v \in [v_n, n]} f_v(u) \exp \left[ -u(2\zeta(2))^{-1} \log^2 v - u^2 \frac{\zeta(3)}{3\zeta^3(2)} \log^3 v \right] \geq 1.$$

This verifies (55), as required.  $\square$

2.11. *How soon do the species part their ways?* Recall from Section 1.2 the notion of *pruned spanning tree* on  $t$  random leaves within the tree model on  $n$  leaves. Write  $S_{n,t}$  for the edge height of the first branchpoint in the pruned tree. In other words, the number of edges from the root to the vertex after which the  $t$  sampled leaves are first split into some  $(k, t - k)$  leaf subsets. Conditioned on the size  $k$  of the left subtree at the root of the tree with  $n$  leaves, the probability that the  $t$  sampled leaves are all in this left subtree is  $\frac{\binom{k}{t}}{\binom{n}{t}}$ . Therefore, since  $q(n, k) = \frac{n}{2h_{n-1}k(n-k)}$ , we obtain the recursion

$$(62) \quad \mathbb{E}[S_{n,t}] = 1 + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{(n/k) \mathbb{E}[S_{k,t}]}{n-k} \frac{\binom{k}{t}}{\binom{n}{t}}, \quad n \geq t \geq 2,$$

( $\mathbb{E}[S_{k,1}] = 0$ ), or, introducing  $\Phi_{n,t} = (n - 1)_{t-1} \mathbb{E}[S_{n,t}]$ ,

$$(63) \quad \Phi_{n,t} = (n - 1)_{t-1} + \frac{1}{h_{n-1}} \sum_{k=1}^{n-1} \frac{\Phi_{k,t}}{n-k}.$$

PROPOSITION 2.14.

$$\mathbb{E}[S_{n,t}] = \frac{\log n}{h_{t-1}} + O(1) \quad \text{as } n \rightarrow \infty.$$

PROOF. Given  $\alpha > 0$ , define

$$U_{v,t} = \Phi_{v,t} - \alpha v^{t-1} \log v.$$

Then, by (63), we have

$$(64) \quad U_{v,t} = (v - 1)_{t-1} + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{U_{k,t}}{v-k} + \alpha \left( \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{k^{t-1} \log k}{v-k} - v^{t-1} \log v \right),$$

and the coefficient by  $\alpha$  equals

$$\begin{aligned} & \frac{v^{t-1}}{h_{v-1}} \sum_{k=1}^{v-1} \frac{(k/v)^{t-1} [\log v + \log(k/v)]}{v-k} - v^{t-1} \log v \\ &= \frac{v^{t-1}}{h_{v-1}} \left( \log v \sum_{k=1}^{v-1} \frac{(k/v)^{t-1} - 1}{v-k} + \log v \sum_{k=1}^{v-1} \frac{1}{v-k} + \sum_{k=1}^{v-1} \frac{(k/v)^{t-1} \log(k/v)}{v-k} \right) \\ & \quad - v^{t-1} \log v \end{aligned}$$

$$\begin{aligned}
 &= \frac{v^{t-1}}{h_{v-1}} \left( \log v \int_0^1 \frac{x^{t-1} - 1}{1-x} dx + h_{v-1} \log v + O(1) \right) - v^{t-1} \log v \\
 &= -\frac{v^{t-1} \log v}{h_{v-1}} h_{t-1} + O(v^{t-1} \log^{-1} v).
 \end{aligned}$$

So, the equation (64) becomes

$$\begin{aligned}
 (65) \quad U_{v,t} &= (v-1)_{t-1} + \alpha \left( -\frac{v^{t-1} \log v}{h_{v-1}} h_{t-1} + O(v^{t-1} \log^{-1} v) \right) + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{U_{k,t}}{v-k} \\
 &= O(v^{t-1} \log^{-1} v) + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{U_{k,t}}{v-k}
 \end{aligned}$$

if we choose  $\alpha = \frac{1}{h_{t-1}}$ . Consequently, for some constant  $\beta$ ,

$$|U_{v,t}| \leq \beta v^{t-1} \log^{-1} v + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{|U_{k,t}|}{v-k}.$$

For a constant  $B$ , to be chosen shortly, we have

$$\begin{aligned}
 &\beta v^{t-1} \log^{-1} v + \frac{1}{h_{v-1}} \sum_{k=1}^{v-1} \frac{Bk^{t-1}}{v-k} \\
 &= \beta v^{t-1} \log^{-1} v + \frac{Bv^{t-1}}{h_{v-1}} \sum_{k=1}^{v-1} \frac{(k/v)^{t-1}}{v-k} \\
 &= \beta v^{t-1} \log^{-1} v + \frac{Bv^{t-1}}{h_{v-1}} \left( h_{v-1} + \int_0^1 \frac{x^{t-1} - 1}{1-x} dx + O(v^{-1}) \right) \\
 &= \beta v^{t-1} \log^{-1} v + \frac{Bv^{t-1}}{h_{v-1}} (h_{v-1} - h_{t-1} + O(v^{-1})) < Bv^{t-1},
 \end{aligned}$$

provided that

$$\beta \log^{-1} v - B \left( \frac{h_{t-1}}{h_{v-1}} + O(v^{-1}) \right) < 0.$$

And this inequality holds for all  $v \geq 2$ , if we choose  $B$  sufficiently large. It follows, by induction on  $v$ , that  $|U_{v,t}| \leq Bv^{t-1}$ . Consequently

$$\Phi_{v,t} = \alpha v^{t-1} \log v + O(v^{t-1}),$$

so that

$$\mathbb{E}[S_{v,t}] = \frac{\Phi_{v,t}}{(v-1)_{t-1}} = \alpha \log v + O(1), \quad \alpha = \frac{1}{h_{t-1}}. \quad \square$$

Within the same notion of *pruned spanning tree* on  $t$  random leaves within the tree model on  $n$  leaves, a more complicated statistic is the edge-length of the pruned tree, which we denote as  $S_{n,t}^*$ . To derive the counterpart of (62), notice that the total number of ways to partition the set  $[n] \setminus [t]$  into two trees, the left one of cardinality  $k$ , with  $t_1 \leq t$  vertices from  $[t]$  and the right one of cardinality  $n - k$ , with  $t_2 = t - t_1$  remaining vertices from  $[t]$ , equals



$\binom{n-t}{k-t_1}$ . Defining  $S_{n,0}^* = 0, S_{n,1}^* = 0, \forall n \geq 0$ , we have the recursion: for  $n \geq t \geq 2$ ,

$$\begin{aligned} \mathbb{E}[S_{n,t}^*] &= 1 + \sum_{k=1}^{n-1} \frac{n}{2h_{n-1}k(n-k)} \cdot \binom{n}{k}^{-1} \\ &\quad \times \sum_{t_1 \leq t} \binom{n-t}{k-t_1} (\mathbb{E}[S_{k,t_1}^*] + \mathbb{E}[S_{n-k,t_2}^*]) \\ &= 1 + \sum_{k=1}^{n-1} \frac{n}{2h_{n-1}k(n-k)} \sum_{t_1 \leq t} \frac{(k)_{t_1} (n-k)_{t_2}}{(n)_t} (\mathbb{E}[S_{k,t_1}^*] + \mathbb{E}[S_{n-k,t_2}^*]) \\ &= 1 + \frac{1}{h_{n-1}} \sum_{k=2}^{n-1} \sum_{t_1=2}^t \frac{(k-1)_{t_1-1} (n-k)_{t_2}}{(n-1)_{t-1} (n-k)} \mathbb{E}[S_{k,t_1}^*]. \end{aligned}$$

Therefore, with  $\Psi_{n,t} := (n-1)_{t-1} \mathbb{E}[S_{n,t}^*]$ , so that  $\Psi_{n,0} = \Psi_{n,1} = 0, \Psi_{n,t} = 0$  for  $n < t$ , we obtain

$$(66) \quad \Psi_{n,t} = (n-1)_{t-1} + \frac{1}{h_{n-1}} \sum_{t_1=2}^t \sum_{k=2}^{n-1} \frac{(n-k)_{t_2}}{n-k} \Psi_{k,t_1}, \quad n \geq t \geq 2.$$

This equation is similar to (63). Because of the new factor  $(n-k)_{t_2}$ , we will use

$$(67) \quad (a)_b = \sum_{j=1}^b s(b, j) a^j,$$

where  $s(b, j)$  is the signed Stirling number of the first kind, so that  $|s(b, j)|$  is the total number of permutations of  $[b]$  with  $j$  cycles.

We now repeat the statement of Theorem 1.8.

PROPOSITION 2.15.

$$\mathbb{E}[S_{n,t}^*] = \alpha(t) \log n + O(1), \quad \alpha(t) = \left( h_{t-1} - \sum_{t_1+t_2=t} \frac{(t_1-1)!(t_2-1)!}{(t-1)!} \right)^{-1}.$$

PROOF. The argument is guided by the proof of Proposition 2.14. Given  $\alpha > 0$ , define

$$V_{v,t} = \Psi_{v,t} - \alpha v^{t-1} \log v, \quad v \geq t \geq 2.$$

By (66), we have

$$(68) \quad \begin{aligned} V_{v,t} &= (v-1)_{t-1} + \frac{1}{h_{v-1}} \sum_{t_1=2}^t \sum_{k=2}^{v-1} \frac{(v-k)_{t_2}}{v-k} V_{k,t_1} \\ &\quad + \alpha \left( \frac{1}{h_{v-1}} \sum_{t_1=2}^t \sum_{k=2}^{v-1} \frac{(v-k)_{t_2}}{v-k} k^{t_1-1} \log k - v^{t-1} \log v \right). \end{aligned}$$

Consider the factor by  $\alpha$ . By (67),

$$\begin{aligned} \sum_{k=2}^{v-1} \frac{(v-k)_{t_2}}{v-k} k^{t_1-1} \log k &= \sum_{j=0}^{t_2} s(t_2, j) \Sigma(v, t_1, j), \\ \Sigma(v, t_1, j) &:= \sum_{k=2}^{v-1} (v-k)^{j-1} k^{t_1-1} \log k. \end{aligned}$$

Recalling that  $t_1 > 1$ , we write

$$\begin{aligned} \Sigma(v, t_1, 0) &= \sum_{k=2}^{v-1} \frac{k^{t_1-1} \log k}{v-k} = \sum_{k=2}^{v-1} \frac{k^{t_1-1} (\log v + \log(k/v))}{v-k} \\ &= (\log v) \left( v^{t_1-1} h_{v-1} + \sum_{k=2}^{v-1} \frac{k^{t_1-1} - v^{t_1-1}}{v-k} \right) + \sum_{k=2}^{v-1} \frac{k^{t_1-1} \log(k/v)}{v-k}, \end{aligned}$$

and

$$\begin{aligned} \sum_{k=2}^{v-1} \frac{k^{t_1-1} - v^{t_1-1}}{v-k} &= v^{t_1-1} \left( \int_0^1 \frac{x^{t_1-1} - 1}{1-x} dx + O(v^{-1}) \right) \\ &= v^{t_1-1} \left( - \int_0^1 \sum_{s=0}^{t_1-2} x^s dx + O(v^{-1}) \right) \\ &= -v^{t_1-1} h_{t_1-1} + O(v^{t_1-2}), \end{aligned}$$

while it is easy to see that  $\sum_{k=2}^{v-1} \frac{k^{t_1-1} \log(k/v)}{v-k}$  is of order  $v^{t_1-1} \int_0^1 \frac{x^{t_1-1} \log x}{1-x} dx = O(v^{t_1-1})$ . Therefore,

$$(69) \quad \Sigma(v, t_1, 0) = (h_{v-1} - h_{t_1-1})v^{t_1-1} \log v + O(v^{t_1-1}).$$

Suppose that  $j > 0$ . Then

$$\begin{aligned} \Sigma(v, t_1, j) &= v^{t_1+j-1} \left( v^{-1} \sum_{k=1}^{v-1} (1-k/v)^{j-1} (k/v)^{t_1-1} [\log v + \log(k/v)] \right) \\ &= v^{t_1+j-1} \left[ (\log v) \int_0^1 (1-x)^{j-1} x^{t_1-1} dx \right. \\ (70) \quad &\quad \left. + \int_0^1 (1-x)^{j-1} x^{t_1-1} (\log x) dx + O(v^{-1} \log v) \right] \\ &= \frac{(j-1)!(t_1-1)!}{(t_1+j-1)!} \cdot v^{t_1+j-1} \log v + O(v^{t_1+j-2} \log v), \end{aligned}$$

and  $t_1 + j - 1 \leq t_1 + t_2 - 1 = t - 1$ . Combining (69) and (70), and using  $s(b, b) = 1, s(b, 0) = 0$  for  $b > 0$ , we have

$$\begin{aligned} &\sum_{k=2}^{v-1} \frac{(v-k)_{t_2}}{v-k} k^{t_1-1} \log k \\ &= (h_{v-1} - h_{t_1-1})v^{t_1-1} \log v + \frac{(t_2-1)!(t_1-1)!}{(t-1)!} v^{t-1} \log v + O(v^{t-2} \log v). \end{aligned}$$

So, the factor by  $\alpha$  in (68) is

$$\begin{aligned} &\frac{v^{t-1} \log v}{h_{v-1}} \left( h_{v-1} - h_{t-1} + \sum_{t_1=1}^t \frac{(t_2-1)!(t_1-1)!}{(t-1)!} + O(v^{-1}) \right) - v^{t-1} \log v \\ &= \frac{v^{t-1} \log v}{h_{v-1}} \left( -h_{t-1} + \sum_{t_1=1}^t \frac{(t_2-1)!(t_1-1)!}{(t-1)!} + O(v^{-1}) \right). \end{aligned}$$

Consequently the equation (68) becomes

$$\begin{aligned} V_{v,t} &= (v-1)_{t-1} + \alpha \frac{v^{t-1} \log v}{h_{v-1}} \left( -h_{t-1} + \sum_{t_1=1}^t \frac{(t_2-1)!(t_1-1)!}{(t-1)!} + O(v^{-1}) \right) \\ &\quad + \frac{1}{h_{n-1}} \sum_{t_1=2}^t \sum_{k=2}^{v-1} \frac{(v-k)_{t_2}}{v-k} V_{k,t_1} \\ &= O(v^{t-1} \log^{-1} v) + \frac{1}{h_{n-1}} \sum_{t_1=2}^t \sum_{k=2}^{v-1} \frac{(v-k)_{t_2}}{v-k} V_{k,t_1}, \end{aligned}$$

if we select

$$\alpha = \left( h_{t-1} - \sum_{t_1+t_2=t} \frac{(t_1-1)!(t_2-1)!}{(t-1)!} \right)^{-1}.$$

We omit the rest of the proof since it runs just like the final part of the proof of Proposition 2.14.  $\square$

2.12. *Counting the subtrees by the number of their leaves.* Since the tree with  $n$  leaves has  $2n - 1$  vertices, there are exactly  $2n - 1$  subtrees, with the number of leaves ranging, with possible gaps, from 1 to  $n$ . Let  $X_n(t)$  be the number of subtrees with  $t$  leaves; so  $X_n(1) = n$ ,  $X_n(n) = 1$ , and  $X_n(t) = 0$  for  $t > n$ . Now,  $\sum_{t \geq 1} X_n(t) = 2n - 1$ , so  $\{u_n(t)\}_{t \geq 1} := \{\frac{E[X_n(t)]}{2^{n-1}}\}_{t \geq 1}$  is the probability distribution of the number of leaves in the uniformly random subtree, that is, the subtree rooted at the uniformly random vertex of the whole tree. Furthermore

$$(71) \quad E[X_n(t)] = \frac{n}{2h_{n-1}} \sum_{j=1}^{n-1} \frac{E[X_j(t)] + E[X_{n-j}(t)]}{j(n-j)} = \frac{n}{h_{n-1}} \sum_{j=1}^{n-1} \frac{E[X_j(t)]}{j(n-j)}.$$

So, with  $\xi_n(t) := \frac{E[X_n(t)]}{n}$ , and  $h_k := \sum_{j=1}^k \frac{1}{j}$ , we have

$$(72) \quad \xi_n(t) = \frac{1}{h_{n-1}} \sum_{j=t}^{n-1} \frac{\xi_j(t)}{n-j}, \quad n \geq t + 1, \quad \left( \xi_t(t) = \frac{1}{t} \right),$$

and clearly  $u_n(t) = \frac{\xi_n(t)}{2^{-n+1}}$ .

**THEOREM 2.16.** (i)  $\xi_n(t) \in [\frac{1}{t^2}, \frac{1}{th_t}]$ ,  $\frac{1}{t} \leq \sum_{\tau \geq t} \xi_n(\tau) \leq \frac{2}{t}$ , the last bound implying that the sequence of distributions  $\{u_n(t)\}_{t \geq 1}$  is tight.

(ii) For  $q \in (0, 1)$ ,  $F_n(q) := \sum_{t \geq 1} q^t \xi_n(t)$  decreases with  $n$ . Consequently, the sequence of distributions  $\{u_n(t)\}_{t \geq 1}$  converges to a proper distribution  $\{u(t)\}_{t \geq 1}$ .

(iii) However, the expected size of the uniformly random subtree is asymptotic to  $\frac{3}{2\pi^2} \log^2 n$ .

We conjecture that (ii) can be improved to the stronger assertion that  $\xi_n(t)$  is decreasing with  $n$ , for each  $t$ . We are grateful to Huseyin Acan [1] for numerically verifying this for  $n$  and  $t$  below 1000.

**PROOF.** (i) Let us show that  $\xi_n(t) \geq \frac{1}{t^2}$  for  $n \geq t > 1$ . By (71), we have  $\xi_t(t) = \frac{1}{t}$  and  $\xi_{t+1}(t) = \frac{1}{th_t}$ , both above  $\frac{1}{t^2}$ . Suppose that  $n \geq t + 1$  is such that  $\xi_j(t) \geq \frac{1}{t^2}$  for all  $j \in [t, n]$ .

This is true for  $n = t + 1$ . For  $n > t + 1$ ,

$$\begin{aligned} \xi_n(t) &\geq \frac{\xi_t(t)}{h_{n-1}(n-t)} + \frac{1}{t^2 h_{n-1}} \sum_{j=t+1}^{n-1} \frac{a}{n-j} = \frac{1}{h_{n-1}(n-t)t} + \frac{h_{n-1-t}}{t^2 h_{n-1}} \\ &= \frac{1}{t^2} + \frac{1}{h_{n-1}(n-t)t} + \frac{h_{n-1-t} - h_{n-1}}{t^2 h_{n-1}} \\ &\geq \frac{1}{t^2} + \frac{1}{h_{n-1}(n-t)t} - \frac{1}{t^2 h_{n-1}} \cdot \frac{t}{n-t} = \frac{1}{t^2}, \end{aligned}$$

which completes the induction step. The proof of  $\xi_n(t) \leq \frac{1}{th_t}$  is similarly reduced to showing that  $\frac{(n-1)h_t}{(n-t)th_{n-1}} \leq 1$  for  $n > t + 1$ . This is so, as the fraction is at most  $\frac{h_t}{h_{t+1}} \cdot \frac{t+1}{2t}$ .

Let us prove that  $\frac{1}{t} \leq \sum_{\tau \geq t} \xi_n(\tau) \leq \frac{2}{t}$ . Introduce  $Y_n(t) = \sum_{\tau \geq t} X_n(\tau)$ , the total number of subtrees with at least  $t$  leaves, and  $\eta_n(t) := \frac{\mathbb{E}[Y_n(t)]}{n} = \sum_{\tau \geq t} \xi_n(\tau)$ ; so  $\eta_n(1) = \frac{2n-1}{n}$ , and  $\eta_n(n) = \frac{1}{n}$ . Analogously to (71), we have

$$\eta_n(t) = \frac{1}{h_{n-1}} \sum_{j=t}^{n-1} \frac{\eta_j(t)}{n-j}, \quad n \geq t + 1.$$

We need to show that  $\eta_n(t) \leq \frac{2}{t}$  for all  $n \geq t$ . It suffices to consider  $n > t > 1$ . Suppose that for some  $n \geq t$  and all  $j \in [t, n]$  we have  $\eta_j(t) \leq \frac{2}{t}$ . This is definitely true for  $n = t$ . Then

$$\eta_{n+1}(t) = \frac{1}{h_n} \sum_{j=t}^n \frac{\eta_j(t)}{n+1-j} \leq \frac{2}{th_n} \sum_{j=t}^n \frac{1}{n+1-j} = \frac{2h_{n+1-t}}{th_n} \leq \frac{2}{t},$$

which completes the inductive proof of  $\eta_n(t) \leq \frac{2}{t}$ .

(ii) For  $n \geq 2$ , we have

$$\begin{aligned} F_n(q) &= \sum_{t \geq 1} q^t \xi_n(t) = \sum_{t \geq 1} \frac{q^t}{h_{n-1}} \sum_{j=t}^{n-1} \frac{\xi_j(t)}{n-j} \\ &= \frac{1}{h_{n-1}} \sum_{j=1}^{n-1} \frac{1}{n-j} \sum_{t=1}^j q^t \xi_j(t) = \frac{1}{h_{n-1}} \sum_{j=1}^{n-1} \frac{F_j(q)}{n-j}. \end{aligned}$$

Therefore,

$$\begin{aligned} F_{n+1}(q) &= \frac{1}{h_n} \sum_{j=1}^n \frac{F_j(q)}{n+1-j} = \frac{1}{h_n} \left( \frac{F_1(q)}{n} + \sum_{j=2}^n \frac{F_j(q)}{n+1-j} \right) \\ &\leq \frac{1}{h_n} \left( \frac{F_1(q)}{n} + \sum_{j=1}^{n-1} \frac{F_j(q)}{n-j} \right) \leq \frac{1}{h_n} \left( \frac{q}{n} + h_{n-1} F_n(q) \right) \leq F_n(q), \end{aligned}$$

since  $F_n(q) \geq q \xi_n(t) = q$ , and  $h_n - h_{n-1} = \frac{1}{n}$ . Therefore, for each  $q \in (0, 1)$  there exists  $F(q) = \lim_{n \rightarrow \infty} \sum_{t \geq 1} q^t \xi_n(t)$ , implying existence of  $\lim_{n \rightarrow \infty} \sum_{t \geq 1} q^t u_n(t) = 2F(q)$ . For any weakly convergent subsequence of the distributions  $\{u_{n_i}(t)\}_{t \geq 1}$ , for each  $x$  in the unit disc there exists  $\lim_{n_i \rightarrow \infty} \sum_{t \geq 1} x^t u_{n_i}(t)$ , dependent on the subsequence, which is analytic within the disc. All these limits coincide for  $x \in (0, 1)$ , whence they coincide for all  $x$  within the disc, whence on the whole disc. Since the characteristic function determines the distribution uniquely, we see that the whole sequence of the distributions converges to a proper distribution.

(iii)  $Z_n := \sum_{t \geq 1} t X_n(t)$  is the total number of the leaves, each leaf counted as many times as the number of the subtrees rooted at the vertices along the path from the root to the leaf, which is distributed as 1 plus  $L_n$ , the edge-length of the path to the random leaf. Therefore,  $\frac{\mathbb{E}[Z_n]}{2n-1} = \frac{n}{2n-1}(1 + \mathbb{E}[L_n])$ , and it remains to use Proposition 2.9.  $\square$

**3. Other methods.** Our results here demonstrate that the analysis of recursions method is very effective at deriving sharp asymptotics for the questions addressed here. However, there are many other aspects of the model that could be studied, and a wide variety of familiar general modern probabilistic techniques that could be applied. We indicate such possibilities briefly below—see the preprint [4] for a more comprehensive account.

It is intuitively clear that for our tree model (call it  $\mathcal{T}_n$ , say) there should be two  $n \rightarrow \infty$  limit structures:

(a) A scaling limit process, which is a fragmentation of the unit interval via some sigma-finite splitting measure.

(b) A fringe process, which is the local weak limit relative to a random leaf, describable as some marked branching process. This starts with an explicit description of the limit distribution  $\{u(t)\}_{t \geq 1}$  in Theorem 2.16.

Less intuitive is a piece of structure theory. In our discrete-time model there is no simple connection between  $\mathcal{T}_n$  and  $\mathcal{T}_{n+1}$ . In the continuous-time model with our chosen rates  $h_{n-1}$  (and this is the reason for that particular choice), [4] shows there is a nonobvious consistency property under a “delete a random leaf and prune” operation. This enables an inductive construction of a process  $(\mathcal{T}_n, n = 2, 3, 4, \dots)$ , which in turn suggests the possibility of a.s. limit theorems.

Readers may notice that the issues above are somewhat analogous to those arising in the theory surrounding the Brownian continuum random tree (CRT) as a limit of certain other random tree models [2, 9]. However, in contrast to the CRT setting, the two limit processes above would not capture the asymptotics of the quantities studied in this article.

**Acknowledgments** We thank Svante Janson for catching several errors in earlier versions. We are grateful to Huseyin Acan for the numerics below Theorem 2.16. We thank the referees for the time and effort to repeatedly read the paper and to provide expert critical advice. We thank the Editors for the efficient reviewing process.

## REFERENCES

- [1] ACAN, H. Personal communication, 05/10/2023.
- [2] ALDOUS, D. (1991). The continuum random tree. II. An overview. In *Stochastic Analysis (Durham, 1990)*. London Mathematical Society Lecture Note Series **167** 23–70. Cambridge Univ. Press, Cambridge. MR1166406 <https://doi.org/10.1017/CBO9780511662980.003>
- [3] ALDOUS, D. (1996). Probability distributions on cladograms. In *Random Discrete Structures (Minneapolis, MN, 1993)*. IMA Vol. Math. Appl. **76** 1–18. Springer, New York. MR1395604 [https://doi.org/10.1007/978-1-4612-0719-1\\_1](https://doi.org/10.1007/978-1-4612-0719-1_1)
- [4] ALDOUS, D. and JANSON, S. (2023). The critical beta-splitting random tree model II: Overview and open problems. Available at [arXiv:2303.02529](https://arxiv.org/abs/2303.02529).
- [5] ALDOUS, D. J. (2001). Stochastic models and descriptive statistics for phylogenetic trees, from Yule to today. *Statist. Sci.* **16** 23–34. MR1838600 <https://doi.org/10.1214/ss/998929474>
- [6] CURTISS, J. H. (1942). A note on the theory of moment generating functions. *Ann. Math. Stat.* **13** 430–433. MR0007577 <https://doi.org/10.1214/aoms/1177731541>
- [7] DRMOTA, M. (2009). *Random Trees: An Interplay Between Combinatorics and Probability*. Springer, Vienna. MR2484382 <https://doi.org/10.1007/978-3-211-75357-6>
- [8] DYER, M., FRIEZE, A. and PITTEL, B. (1993). The average performance of the greedy matching algorithm. *Ann. Appl. Probab.* **3** 526–552. MR1221164

- [9] EVANS, S. N. (2008). *Probability and Real Trees. Lecture Notes in Math.* **1920**. Springer, Berlin. Lectures from the 35th Summer School on Probability Theory held in Saint-Flour, July 6–23, 2005. [MR2351587](#) <https://doi.org/10.1007/978-3-540-74798-7>
- [10] GRAHAM, R. L., KNUTH, D. E. and PATASHNIK, O. (1994). *Concrete Mathematics: A Foundation for Computer Science*, 2nd ed. Addison-Wesley Company, Reading, MA. [MR1397498](#)
- [11] IKSANOV, A. (2024). Another proof of CLT for critical beta-splitting tree. Unpublished.
- [12] JANSON, S. (2012). Simply generated trees, conditioned Galton-Watson trees, random allocations and condensation: Extended abstract. In *23rd Intern. Meeting on Probabilistic, Combinatorial, and Asymptotic Methods for the Analysis of Algorithms (AofA'12)*. *Discrete Math. Theor. Comput. Sci. Proc.*, AQ 479–490. Assoc. Discrete Math. Theor. Comput. Sci., Nancy. [MR2957350](#)
- [13] KOLESNIK, B. (2024). Critical beta-splitting, via contraction. Available at [arXiv:2404.16021](#).
- [14] LAMBERT, A. (2017). Probabilistic models for the (sub)tree(s) of life. *Braz. J. Probab. Stat.* **31** 415–475. [MR3693976](#) <https://doi.org/10.1214/16-BJPS320>
- [15] MAHMOUD, H. M. and PITTEL, B. (1989). Analysis of the space of search trees under the random insertion algorithm. *J. Algorithms* **10** 52–75. [MR0987097](#) [https://doi.org/10.1016/0196-6774\(89\)90023-0](https://doi.org/10.1016/0196-6774(89)90023-0)
- [16] MUKHERJEA, A., RAO, M. and SUEN, S. (2006). A note on moment generating functions. *Statist. Probab. Lett.* **76** 1185–1189. [MR2270543](#) <https://doi.org/10.1016/j.spl.2005.12.026>
- [17] PITTEL, B. (1987). An urn model for cannibal behavior. *J. Appl. Probab.* **24** 522–526. [MR0889816](#) <https://doi.org/10.1017/s0021900200031156>
- [18] PITTEL, B. (1990). On tree census and the giant component in sparse random graphs. *Random Structures Algorithms* **1** 311–342. [MR1099795](#) <https://doi.org/10.1002/rsa.3240010306>
- [19] PITTEL, B. (1999). Normal convergence problem? Two moments and a recurrence may be the clues. *Ann. Appl. Probab.* **9** 1260–1302. [MR1728562](#) <https://doi.org/10.1214/aoap/1029962872>
- [20] PITTEL, B. and POOLE, D. (2016). Asymptotic distribution of the numbers of vertices and arcs of the giant strong component in sparse random digraphs. *Random Structures Algorithms* **49** 3–64. [MR3521273](#) <https://doi.org/10.1002/rsa.20622>
- [21] YAKYMIV, A. L. (2011). A generalization of the Curtiss theorem for moment generating functions. *Mat. Zametki* **90** 947–952. [MR2962966](#) <https://doi.org/10.1134/S0001434611110290>