

1 Uncovering a random tree

2 Benjamin Hackl   

3 Uppsala University, Sweden

4 University of Klagenfurt, Austria

5 Alois Panholzer  

6 TU Wien, Austria

7 Stephan Wagner  

8 Uppsala University, Sweden

9 — Abstract —

10 We consider the process of uncovering the vertices of a random labeled tree according to their
11 labels. First, a labeled tree with n vertices is generated uniformly at random. Thereafter, the
12 vertices are uncovered one by one, in order of their labels. With each new vertex, all edges to
13 previously uncovered vertices are uncovered as well. In this way, one obtains a growing sequence
14 of forests. Three particular aspects of this process are studied in this extended abstract: first the
15 number of edges, which we prove to converge to a stochastic process akin to a Brownian bridge after
16 appropriate rescaling. Second, the connected component of a fixed vertex, for which different phases
17 are identified and limiting distributions determined in each phase. Lastly, the largest connected
18 component, for which we also observe a phase transition.

19 **2012 ACM Subject Classification** Mathematics of computing → Random graphs; Mathematics of
20 computing → Generating functions

21 **Keywords and phrases** Labeled tree, uncover process, functional central limit theorem, limiting
22 distribution, phase transition

23 **Digital Object Identifier** 10.4230/LIPIcs.AofA.2022.3

24 **Funding** *Stephan Wagner*: supported by the Knut and Alice Wallenberg Foundation.

25 **1** Introduction

26 We consider the process of uncovering the vertices of a random tree: starting either from one
27 of the n^{n-2} unrooted or one of the n^{n-1} rooted unordered labeled trees of size n (i.e., with
28 n vertices) chosen uniformly at random, we uncover the vertices one by one in order of their
29 labels. This yields a growing sequence of (rooted) forests induced by the uncovered vertices,
30 and we are interested in the evolution of these forests from the first vertex to the point that
31 all vertices are uncovered.

32 This model is motivated by stochastic models known as coalescent models for particle
33 coalescence, most notably the additive and the multiplicative coalescent [2] and the Kingman
34 coalescent [7]. To make the distinction between these classical coalescent models and
35 our model more explicit, let us briefly revisit the additive coalescent model (see [1]) as a
36 prototypical example. This model describes a Markov process on a state space consisting of
37 tuples (x_1, x_2, \dots) with $x_1 \geq x_2 \geq \dots \geq 0$ and $\sum_{i \geq 0} x_i = 1$ that model the fragmentation of
38 a unit mass into clusters of mass x_i . Pairs of clusters with masses x_i and x_j then merge into
39 a new cluster of mass $x_i + x_j$ at rate $x_i + x_j$. In the corresponding discrete time version of the
40 process, exactly two clusters are merged in every time step. There are various combinatorial
41 settings in which this discrete additive coalescent model appears, for example in the evolution
42 of parking blocks in parking schemes related to “hashing with linear probing” [4], as the
43 dual fragmentation process in the context of the random cutting of trees [6], and in a certain
44 scheme for merging forests by uncovering one edge in every time step [10].



© Benjamin Hackl, Alois Panholzer and Stephan Wagner;
licensed under Creative Commons License CC-BY 4.0

33rd International Conference on Probabilistic, Combinatorial and Asymptotic Methods for the Analysis of
Algorithms (AofA 2022).

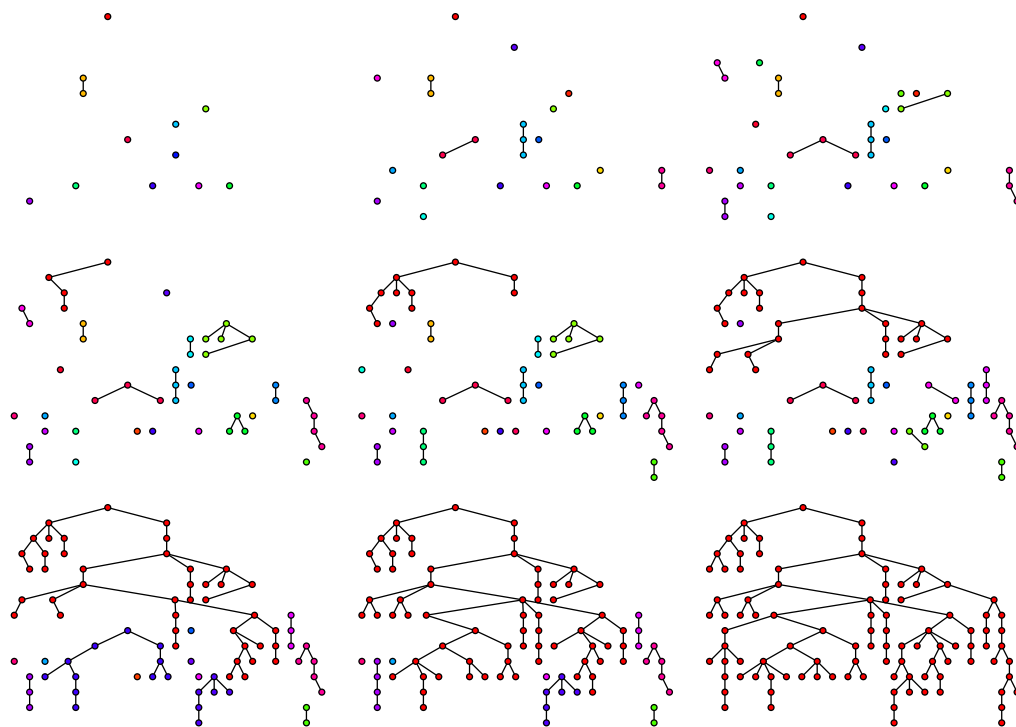
Editor: Mark Daniel Ward; Article No. 3; pp. 3:1–3:17



Leibniz International Proceedings in Informatics
LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

3:2 Uncovering a random tree

45 While the latter incarnation of the additive coalescent in which edges are uncovered
46 successively is very much related in spirit to our vertex uncover model, the underlying
47 processes are fundamentally different: these classical coalescent models rely on the fact that
48 exactly two clusters are merged in every time step, which is not the case in our model. When
49 uncovering a new vertex, a more or less arbitrary number of edges (including none at all)
50 can be uncovered. There are coalescent models like the Λ -coalescent, a generalization of the
51 Kingman coalescent [11], which allow for more than two clusters being merged—however, at
52 present we are not aware of any known coalescent model that is able to capture the behavior
53 of the vertex uncover process.



■ **Figure 1** A few snapshots of the *uncover process* applied to a random labeled tree of size 100. From left to right and top to bottom, there are 12, 23, 34, . . . , 89, and 100 uncovered vertices in the figures, respectively. Vertex labels are omitted for the sake of readability, and vertices are colored per connected component.

54 **Overview.** Different aspects of the uncover process on labeled trees are investigated in this
55 extended abstract. In Section 2, we study the stochastic process given by the number of
56 uncovered edges. The corresponding main result, a full characterization of the process and
57 its limiting behavior, is given in Theorem 5. In this extended abstract, we sketch the proofs
58 of the results in Section 2—full details can be found in Appendix A.

59 Sections 3 and 4 are both concerned with cluster sizes, i.e., with the sizes of the connected
60 components that are created throughout the process. In particular, in Section 3 we shift our
61 attention to rooted labeled trees, to study the behavior of the component containing a fixed
62 vertex. The expected size of the root cluster is analyzed in Theorem 9. Furthermore, we
63 show that the number of rooted trees whose root cluster has a given size is given by a rather
64 simple enumeration formula—which, in turn, manifests in Theorem 11, a characterization

65 of the different limiting distributions for the root cluster size depending on the number of
66 uncovered vertices.

67 Finally, in Section 4 we use the results on the root cluster to draw conclusions regarding
68 the size of the largest cluster in the tree.

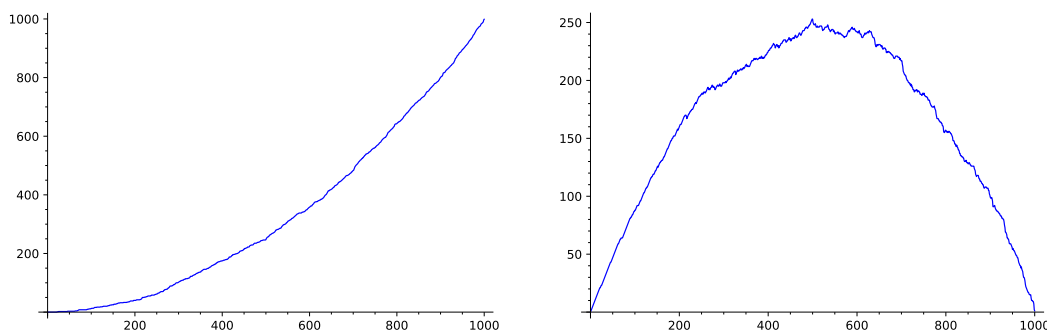
69 **Notation.** Throughout this work we use the notation $[n] = \{1, \dots, n\}$ and $[k, \ell] = \{k, k +$
70 $1, \dots, \ell\}$ for discrete intervals, and $x^{\underline{j}} = x(x-1) \cdots (x-j+1)$ for the falling factorials. The
71 floor and ceiling function are denoted by $\lfloor x \rfloor$ and $\lceil x \rceil$, respectively. Furthermore, we use \mathcal{T}
72 and \mathcal{T}^\bullet for the combinatorial classes of labeled trees and rooted labeled trees, respectively,
73 and \mathcal{T}_n and \mathcal{T}_n^\bullet for the classes of labeled and rooted labeled trees of size n , i.e., with n vertices.
74 Finally, we use $X_n \xrightarrow{d} X$ and $X_n \xrightarrow{p} X$ to denote convergence in distribution resp. probability
75 of a sequence of random variables (r.v.) $(X_n)_{n \geq 0}$ to the r.v. X .

76 2 The number of uncovered edges

77 In this section our main interest is the behavior of the number of uncovered edges in the
78 uncover process. We begin by introducing a formal parameter for this quantity.

79 **► Definition 1.** Let T be a labeled tree with vertex set $V(T) = [n]$. For $1 \leq j \leq n$, we let
80 $k_j(T) := \|T[1, 2, \dots, j]\|$ denote the number of edges in the subgraph of T induced by the
81 vertices in $[j]$. We refer to the sequence $(k_j(T))_{1 \leq j \leq n}$ as the (edge) uncover sequence.

82 We start with a few simple observations. First, for any labeled tree of order n we have
83 $k_1(T) = 0$, as well as $k_n(T) = n - 1$. Second, as any induced subgraph of a tree is a forest,
84 and as forests have the elementary property that the number of edges together with the
85 number of connected components gives the order of the forest, we find that $j - k_j(T)$ is
86 the number of connected components after uncovering the first j vertices of T . Figure 2
87 illustrates the progression of the number of edges and the number of connected components
88 for $1 \leq j \leq 1000$ in a randomly chosen labeled tree on 1000 vertices.



89 **Figure 2** Progression of the number of edges (left) and the number of connected components
90 (right) when sequentially uncovering a random labeled tree on 1000 vertices.

89 Moreover, the fact that the first j vertices of the tree induce a forest also yields the sharp
90 bound $0 \leq k_j(T) \leq j - 1$ for all $1 \leq j \leq n - 1$. The lower bound is attained by the star with
91 central vertex n , and the upper bound is attained by the (linearly ordered) path. We can
92 also observe that as soon as $k_{n-1}(T) > 0$, the set of edges added in the last uncover step is
93 not determined uniquely. Thus, the star with central vertex n is the only tree which is fully
94 determined by its uncover sequence.

3:4 Uncovering a random tree

95 The following theorem provides explicit enumeration formulas for the number of trees
96 with a partially and fully specified uncover sequence, respectively.

97 **► Theorem 2.** *Let r be a fixed positive integer with $1 \leq r < n - 1$, let j_1, j_2, \dots, j_r be
98 positive integers with $1 < j_1 < j_2 < \dots < j_r < n$, and let a_1, a_2, \dots, a_r be a nondecreasing
99 sequence of nonnegative integers satisfying $a_i \leq j_i - 1$ for all $1 \leq i \leq r$. Additionally, let
100 $j_0 = 1$ and $a_0 = 0$. Then, the number of rooted labeled trees T of order n that satisfy
101 $k_{j_i}(T) = a_i$ for all $1 \leq i \leq r$ is given by*

$$102 \quad (n - j_r)^{j_r - a_r - 1} n^{n - j_r - 1} \\ 103 \quad \times \prod_{i=1}^r \left(\sum_{h=0}^{a_i - a_{i-1}} \binom{j_{i-1} - a_{i-1} - 1}{h} \binom{j_i - j_{i-1}}{a_i - a_{i-1} - h} (j_i - j_{i-1})^h j_i^{a_i - a_{i-1} - h} \right). \quad (1)$$

106 We first derive a helpful auxiliary result, namely an explicit formula for the corresponding
107 (multivariate) generating function. The enumeration formula will then follow by extracting
108 the appropriate coefficients.

109 **► Lemma 3.** *In the setting of Theorem 2, the multivariate generating function for the
110 increments in the edge uncover process is given by*

$$111 \quad E_n(z_1, z_2, \dots, z_r) = n^{n - j_r - 1} \prod_{i=1}^r \left(n - j_r + j_i z_i + \sum_{h=i+1}^r (j_h - j_{h-1}) z_h \right)^{j_i - j_{i-1}}. \quad (2)$$

112 *In other words, the coefficient of the monomial $z_1^{a_1} z_2^{a_2 - a_1} \dots z_r^{a_r - a_{r-1}}$ in the expansion of
113 $E_n(z_1, \dots, z_r)$ is the number of labeled trees T of order n with $k_{j_i}(T) = a_i$ for all $1 \leq i \leq r$.*

114 **► Remark 4.** By specifying the integers j_1, j_2, \dots, j_r , the uncover process is effectively
115 partitioned into intervals. This is also reflected by the quantities occurring in the product
116 in (2): the difference $j_i - j_{i-1}$ corresponds to the number of vertices uncovered in the i -th
117 interval, j_i represents the number of vertices uncovered in total up to the i -th interval, and
118 $n - j_r$ corresponds to the number of vertices uncovered in the last interval.

119 **Proof of Lemma 3.** We begin by observing that when the process uncovers the vertex with
120 label j , edges to all adjacent vertices whose label is less than j are uncovered as well. To
121 determine the generating function of the edge increments, we assign the weight $x_i y_j$ to
122 the edge connecting vertex i and vertex j with $i < j$, and then consider the generating
123 function for the tree weight $w(T)$ (which is defined as the product of the edge weights);
124 $E_n(z_1, \dots, z_r) = \sum_{|T|=n} w(T)$.

125 Following Martin and Reiner [9, Theorem 4] or Remmel and Williamson [12, Equation
126 (8)], the generating function of the tree weights $w(T)$ has the explicit formula

$$127 \quad \sum_{|T|=n} w(T) = x_1 y_n \prod_{j=2}^{n-1} \left(\sum_{i=1}^n x_{\min(i,j)} y_{\max(i,j)} \right). \quad (3)$$

128 As initially observed, edges that are counted by $k_{j_i}(T)$ are precisely those that induce a
129 factor y_ℓ for some $\ell \leq j_i$. Thus we make the following substitutions: $x_\ell = 1$ for all ℓ , $y_\ell = z_i$
130 if and only if $j_{i-1} < \ell \leq j_i$ (where¹ $j_0 = 1$), and $y_\ell = 1$ if $\ell > j_r$. To deal with the sum over

¹ Observe that y_1 does not occur, since at least one of the ends of every edge has a label greater than 1.

131 $y_{\max(i,j)}$, observe that we can rewrite it as

$$132 \quad \sum_{i=1}^n y_{\max(i,j)} = n - j_r + \sum_{i=1}^{j_1} y_{\max(i,j)} + \cdots + \sum_{i=j_{r-1}+1}^{j_r} y_{\max(i,j)}.$$

133 In this form, the different values assumed by the sum when j moves through the ranges
 134 $1 < j \leq j_1$, $j_1 < j \leq j_2$, etc. can be determined directly. For some j with $j_{i-1} < j \leq j_i$, the
 135 contribution to the product in (3) is

$$136 \quad n - j_r + j_i z_i + \sum_{h=i+1}^r (j_h - j_{h-1}) z_h,$$

137 and for $j_r < j \leq n - 1$ all y -variables are replaced by 1, so that the contribution to the
 138 product is a factor n . Putting both of these observations together shows that the right-hand
 139 side of (3) can be rewritten as the right-hand side of (2) and thus proves the lemma. ◀

140 With an explicit formula for the generating function of edge increments in the uncover
 141 process at our disposal, an explicit formula for the number of trees with given (partial)
 142 uncover sequence follows as a simple consequence.

143 **Proof of Theorem 2.** It remains to extract the coefficient of $z_1^{a_1} z_2^{a_2 - a_1} \cdots z_r^{a_r - a_{r-1}}$, which
 144 is done step by step, starting with z_1 . ◀

145 2.1 A closer look at the stochastic process

146 The exceptionally nice formula for the generating function of edge increments can be used to
 147 study the stochastic process that describes the number of uncovered edges in more detail. Let
 148 the sequence of random variables $(K_j^{(n)})_{1 \leq j \leq n}$ be the discrete stochastic process modeling
 149 the number of uncovered edges after uncovering the first j vertices in a random labeled tree
 150 of size n , chosen uniformly at random. The expected number of uncovered edges can be
 151 determined by a simple argument: with j uncovered vertices, $\binom{j}{2}$ of the $\binom{n}{2}$ possible positions
 152 for the edges have been uncovered. As every position is, due to symmetry and the uniform
 153 choice of the labeled tree, equally likely to hold one of the $n - 1$ edges, we find

$$154 \quad \mathbb{E}K_j^{(n)} = (n - 1) \frac{\binom{j}{2}}{\binom{n}{2}} = \frac{j(j - 1)}{n}. \quad (4)$$

155 To motivate our investigations further, consider the illustrations in Figure 3. The rescaled
 156 deviation from the mean is reminiscent of a stochastic process known as Brownian bridge.

157 In order to define this process formally, recall first that the Wiener process $(W(t))_{t \in [0,1]}$
 158 is the unique stochastic process that satisfies

- 159 ■ $W(0) = 0$,
- 160 ■ W has independent, stationary increments,
- 161 ■ $W(t) \sim \mathcal{N}(0, t)$ for all $t > 0$, and
- 162 ■ $t \mapsto W(t)$ is almost surely continuous,

163 see [8, Definition 21.8]. A Brownian bridge can then be defined by setting

$$164 \quad B(t) = (1 - t)W(t/(1 - t)), \quad (5)$$

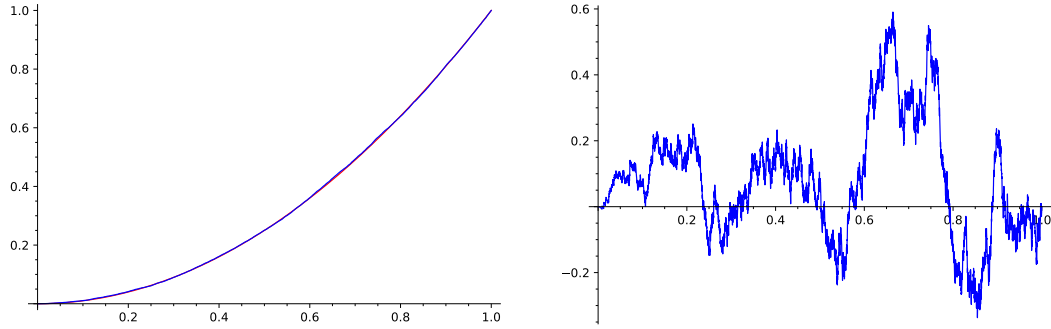
165 see e.g. [13, Exercise 3.10]. The term “bridge” results from the fact that we have $B(0) =$
 166 $B(1) = 0$.

3:6 Uncovering a random tree

167 While the (normalized) deviation from the mean looks like it might converge to a Brownian
 168 bridge, we will prove that this is only *almost* the case. The following theorem characterizes
 169 the stochastic process. For technical purposes, we set $K_0^{(n)} = 0$ and introduce the linearly
 170 interpolated process $(\tilde{K}_t^{(n)})_{t \in [0,1]}$, where

$$171 \quad \tilde{K}_t^{(n)} := (1 + \lfloor tn \rfloor - tn)K_{\lfloor tn \rfloor}^{(n)} + (tn - \lfloor tn \rfloor)K_{\lfloor tn \rfloor + 1}^{(n)}, \quad (6)$$

172 which by construction has continuous paths.



■ **Figure 3** A path of the rescaled stochastic process $(K_{\lfloor tn \rfloor}^{(n)} / (n-1))_{t \in [0,1]}$ (left-hand side) and the corresponding (rescaled) deviation $(\frac{\tilde{K}_t^{(n)} - t^2 n}{\sqrt{n}})_{t \in [0,1]}$ for a random labeled tree of size $n = 10000$.

173 ► **Theorem 5.** Let $(Z^{(n)}(t))_{t \in [0,1]}$ be the continuous stochastic process resulting from centering
 174 and rescaling the linearly interpolated process $(\tilde{K}_t^{(n)})_{t \in [0,1]}$ in the form of

$$175 \quad Z^{(n)}(t) := \frac{\tilde{K}_t^{(n)} - t^2 n}{\sqrt{n}},$$

176 and let $(W(t))_{t \in [0,1]}$ be the standard Wiener process. Then, for $n \rightarrow \infty$, the rescaled process
 177 converges weakly with respect to the sup-norm on $C([0,1])$ to a limiting process $Z^\infty(t)$ that
 178 is given by

$$179 \quad Z^\infty(t) = (1-t)W(t^2/(1-t)). \quad (7)$$

180 Furthermore, for $s, t \in [0,1]$ with $s < t$, the limiting process satisfies

$$181 \quad \mathbb{E}Z^\infty(t) = 0, \quad \mathbb{V}Z^\infty(t) = t^2(1-t), \quad \text{and} \quad \text{Cov}(Z^\infty(s), Z^\infty(t)) = s^2(1-t). \quad (8)$$

182 ► **Remark 6.** While the limiting process $(Z^\infty(t))_{t \in [0,1]}$ is not a Brownian bridge (the cor-
 183 responding variances and covariances as given in (8) do not match), it is closely related.
 184 Comparing the characterization of $Z^\infty(t)$ in (7) to (5), we see that the processes only differ
 185 by the square in the numerator of the argument of the Wiener process.

186 Two main ingredients are required to prove this result: a uniform tightness bound on the
 187 one hand, and information on the finite-dimensional joint distributions of $(\tilde{K}_t^{(n)})_{t \in [0,1]}$ on
 188 the other hand.

189 To prove tightness, let us begin by revisiting (2). Given Cayley's well-known tree
 190 enumeration formula, the corresponding probability generating function for the complete
 191 uncover sequence, i.e., when we choose our integer vector as $\mathbf{j} = (2, 3, \dots, n-1)$, is

$$192 \quad P_n(z_2, \dots, z_{n-1}) = \prod_{i=2}^{n-1} \left(\frac{1}{n} + \frac{i}{n} z_i + \sum_{h=i+1}^{n-1} \frac{1}{n} z_h \right). \quad (9)$$

193 This suggests modeling the process with $n-2$ independent random variables, each representing
 194 an edge increment². The factorization suggests that the j -th increment (which corresponds to
 195 the factor with $i = j + 1$) happens with probability $(j + 1)/n$ when the vertex with label $j + 1$
 196 is uncovered, or with probability $1/n$ every time any of the subsequent vertices are uncovered.
 197 This probabilistic point of view can be used to construct a recursive characterization for the
 198 number of uncovered edges, namely³

$$199 \quad K_{j+1}^{(n)} = K_j^{(n)} + \text{Ber}\left(\frac{j+1}{n}\right) + \text{Bin}\left(j-1 - K_j^{(n)}, \frac{1}{n-j}\right). \quad (10)$$

200 The Bernoulli variable models the probability that the j -th edge increment is added when
 201 uncovering the vertex with label $j + 1$, and the binomial variable models all of the remaining,
 202 not yet uncovered edge increments.

203 Now let us consider a centered and rescaled version of the process $(K_j^{(n)})_{1 \leq j \leq n}$ by defining

$$204 \quad Y_j^{(n)} := \frac{K_j^{(n)} - \frac{j(j-1)}{n}}{n-j}. \quad (11)$$

205 With the help of the recursive description in (10), one can show that $(Y_j^{(n)})_{1 \leq j \leq n-1}$ is a
 206 martingale, see Appendix A.1.

207 We are now ready to state and prove the first required ingredient.

208 ► **Lemma 7.** *For any real $C > 1$ and any positive integer n , the random variable $Z^{(n)}(t)$
 209 satisfies the tightness bound*

$$210 \quad \mathbb{P}\left(\sup_{t \in [0,1]} |Z^{(n)}(t)| \geq C\right) \leq 4(C-1)^{-2}, \quad (12)$$

211 so that for $C \rightarrow \infty$, the probability for the process to exceed C in absolute value converges to
 212 0 uniformly in terms of n .

213 **Sketch of proof.** One first shows that $\sup_{t \in [0,1]} |Z^{(n)}(t)|$ can be bounded in terms of the
 214 deviation of the discrete process $(K_j^{(n)})_{0 \leq j \leq n}$ from its mean. The bound then follows after
 215 expressing the discrete process in terms of the martingale $Y_j^{(n)}$ constructed above, partitioning
 216 the discrete interval $[0, n]$ appropriately and applying both Doob's L^p -inequality and the
 217 union bound. See Appendix A.2 for details. ◀

218 ► **Lemma 8.** *Let r be a fixed positive integer, and let $\mathbf{t} = (t_1, \dots, t_r) \in (0, 1)^r$. Then for
 219 $n \rightarrow \infty$, the random vector*

$$220 \quad \mathbf{K}_{[\mathbf{t}n]}^{(n)} := (K_{[t_1 n]}^{(n)}, K_{[t_2 n]}^{(n)}, \dots, K_{[t_r n]}^{(n)})$$

221 converges, after centering and rescaling, for $n \rightarrow \infty$ in distribution to a multivariate normal
 222 distribution,

$$223 \quad \frac{\mathbf{K}_{[\mathbf{t}n]}^{(n)} - \mathbb{E}\mathbf{K}_{[\mathbf{t}n]}^{(n)}}{\sqrt{n}} \xrightarrow[n \rightarrow \infty]{d} \mathcal{N}(\mathbf{0}, \Sigma),$$

² We explicitly model edge increments here instead of edges, because with this approach we do not need to care about *which* edge is being uncovered. Our model explicitly only captures the behavior of the number of uncovered edges.

³ We slightly abuse notation: formally, we would need to introduce auxiliary variables that are distributed according to the specified binomial and Bernoulli distributions.

3:8 Uncovering a random tree

224 where the expectation vector $\mathbb{E}\mathbf{K}_{\lfloor tn \rfloor}^{(n)}$ satisfies

$$225 \quad \mathbb{E}\mathbf{K}_{\lfloor tn \rfloor}^{(n)} = n(t_1^2, t_2^2, \dots, t_r^2) + O(1), \quad (13)$$

226 and the entries of the variance-covariance matrix $\Sigma = (\sigma_{i,j})_{1 \leq i, j \leq r}$ are

$$227 \quad \sigma_{i,j} = \begin{cases} t_i^2(1-t_j) & \text{if } i \leq j, \\ t_j^2(1-t_i) & \text{if } i > j. \end{cases} \quad (14)$$

228 **Sketch of proof.** The factorization of the probability generating function (9) implies that the
 229 distribution of the corresponding random vector $\Delta_{\lfloor tn \rfloor}^{(n)} = (K_{\lfloor t_1 n \rfloor}^{(n)}, K_{\lfloor t_2 n \rfloor}^{(n)} - K_{\lfloor t_1 n \rfloor}^{(n)}, \dots, K_{\lfloor t_r n \rfloor}^{(n)} -$
 230 $K_{\lfloor t_{r-1} n \rfloor}^{(n)})$ is a marginal distribution of the sum of r independent, multinomially distributed
 231 random vectors. By the multivariate central limit theorem, $\Delta_{\lfloor tn \rfloor}^{(n)}$ converges to a multivariate
 232 normal distribution—and as a consequence, so does $\mathbf{K}_{\lfloor tn \rfloor}^{(n)}$.

233 The variance-covariance matrix of the centered and rescaled vector can be obtained, for
 234 example, by using the martingale constructed above. See again Appendix A.2 for details. ◀

235 **Proof of Theorem 5.** The proof relies on the well-known result asserting that given tightness
 236 of the sequence of corresponding probability measures as well as convergence of the finite-
 237 dimensional probability distributions, a sequence of stochastic processes converges to a
 238 limiting process (see [3, Theorem 7.1, Theorem 7.5]).

239 Tightness is implied by the uniform bound (12) derived in Lemma 7. The (limiting)
 240 behavior of the finite-dimensional distributions for the original process $(K_{\lfloor tn \rfloor}^{(n)})_{t \in [0,1]}$ is
 241 characterized by Lemma 8. This characterization carries over to the linearly interpolated
 242 process by an application of Slutsky's theorem [8, Theorem 13.18] after observing that

$$243 \quad \mathbb{P}\left(\left|Z^{(n)}(t) - \frac{K_{\lfloor tn \rfloor}^{(n)} - t^2 n}{\sqrt{n}}\right| > \varepsilon\right) = \mathbb{P}\left(\frac{|\tilde{K}_t^{(n)} - K_{\lfloor tn \rfloor}^{(n)}|}{\sqrt{n}} > \varepsilon\right) \leq \frac{\mathbb{E}((\tilde{K}_t^{(n)} - K_{\lfloor tn \rfloor}^{(n)})^2)}{n\varepsilon^2}$$

$$244 \quad \leq \frac{\mathbb{E}((K_{\lfloor tn \rfloor + 1}^{(n)} - K_{\lfloor tn \rfloor}^{(n)})^2)}{n\varepsilon^2} \xrightarrow{n \rightarrow \infty} 0,$$

246 as a mechanical computation shows that $\mathbb{E}((K_{\lfloor tn \rfloor + 1}^{(n)} - K_{\lfloor tn \rfloor}^{(n)})^2) = O(1)$.

247 Note that as the finite-dimensional distributions converge to Gaussian distributions, the
 248 limiting process $(Z^\infty(t))_{t \in [0,1]}$ is Gaussian itself—which means that it is fully characterized by
 249 its first and second order moments. As a consequence of Lemma 8, we find for all $s, t \in [0, 1]$
 250 with $s < t$ that

$$251 \quad \mathbb{E}Z^\infty(t) = 0, \quad \mathbb{V}Z^\infty(t) = t^2(1-t), \quad \text{Cov}(Z^\infty(s), Z^\infty(t)) = s^2(1-t).$$

252 It can be checked that if $(W(t))_{t \in [0,1]}$ is a standard Wiener process, the Gaussian process
 253 $((1-t)W(t^2/(1-t)))_{t \in [0,1]}$ has the same first and second order moments and therefore also
 254 the same distribution as Z^∞ . ◀

3 Size of the root cluster

256 We now shift our attention from the number of uncovered edges to the sizes of the connected
 257 components (or *clusters*) appearing in the graph throughout the uncover process. It will
 258 prove convenient to change our tree model to *rooted* labeled trees, as the nature of rooted
 259 trees allows us to focus our investigation on one particular cluster – the one containing the

260 root vertex. In case the root vertex has not yet been uncovered, we will consider the size of
 261 the root cluster to be 0. Formally, we let the random variable $R_n^{(k)}$ be the size of the root
 262 cluster of a (uniformly) random rooted labeled tree of size n with k uncovered vertices.

263 Using the symbolic method for labeled structures (cf. [5, Chapter II]), we can set up a
 264 formal specification for the corresponding combinatorial classes and subsequently extract
 265 functional equations for the associated generating functions. Let \mathcal{T}^\bullet be the class of rooted
 266 labeled trees, and let \mathcal{G} be a refinement of \mathcal{T}^\bullet where the vertices can either be covered or
 267 uncovered, and where uncovered vertices are marked with a marker U . Finally, let \mathcal{F} be
 268 a further refinement of \mathcal{G} where all uncovered vertices in the root cluster are additionally
 269 marked with marker V . A straightforward “top-down approach”, i.e., a decomposition of the
 270 members of the tree family w.r.t. the root vertex, yields the formal specification

$$271 \quad \mathcal{F} = \mathcal{Z} * \text{SET}(\mathcal{G}) + \mathcal{Z} \times \{U, V\} * \text{SET}(\mathcal{F}), \quad \mathcal{G} = \mathcal{Z} * \text{SET}(\mathcal{G}) + \mathcal{Z} \times \{U\} * \text{SET}(\mathcal{G}).$$

272 Introducing the corresponding exponential generating functions $F := F(z, u, v)$ and
 273 $G := G(z, u)$, we obtain the characterizing equations

$$274 \quad F = ze^G + zuve^F, \quad G = z(1 + u)e^G. \tag{15}$$

275 Of course, $G(z, u) = T^\bullet(z(1 + u))$, where T^\bullet is the exponential generating function associated
 276 with \mathcal{T}^\bullet , the Cayley tree function. Starting with (15), the following results on $R_n^{(k)}$ can be
 277 deduced.

278 ► **Theorem 9.** *The expectation $\mathbb{E}(R_n^{(k)})$ is, for $0 \leq k \leq n$ and $n \geq 1$, given by*

$$279 \quad \mathbb{E}(R_n^{(k)}) = \sum_{j=1}^k \frac{j k^j}{n^j}. \tag{16}$$

280 *Depending on the growth of $k = k(n)$, $\mathbb{E}(R_n^{(k)})$ has the following asymptotic behavior:*

$$281 \quad \mathbb{E}(R_n^{(k)}) \sim \begin{cases} \frac{k}{n}, & \text{for } k = o(n), \quad (k \text{ small}), \\ \frac{\alpha}{(1-\alpha)^2}, & \text{for } k \sim \alpha n, \text{ with } 0 < \alpha < 1, \quad (k \text{ in central region}), \\ \frac{n^2}{d^2}, & \text{for } k = n - d, \text{ with } d = \omega(\sqrt{n}) \text{ and } d = o(n), \\ & (k \text{ subcritically large}), \\ \kappa n, & \text{with } \kappa = 1 - ce^{\frac{c^2}{2}} \int_c^\infty e^{-\frac{t^2}{2}} dt, \\ & \text{for } k = n - d, \text{ with } d \sim c\sqrt{n} \text{ and } c > 0, \quad (k \text{ critically large}), \\ n - \sqrt{\frac{\pi}{2}}d\sqrt{n}, & \text{for } k = n - d, \text{ with } d = o(\sqrt{n}), \quad (k \text{ supercritically large}). \end{cases}$$

282 **Proof.** After introducing $E := E(z, u) = \frac{\partial}{\partial v} F(z, u, v)|_{v=1} = \sum_{n,k} \frac{n^{n-1}}{n!} \binom{n}{k} \mathbb{E}(R_n^{(k)})$, consid-
 283 ering the partial derivative of (15) with respect to v yields

$$284 \quad E = \frac{1}{1 - \frac{u}{1+u}G} - 1.$$

285 Extracting coefficients and using $\mathbb{E}(R_n^{(k)}) = \frac{[z^n u^k]E}{[z^n u^k]G}$, we obtain (16). In order to analyze
 286 the asymptotic behavior of $\mathbb{E}(R_n^{(k)})$, the integral representation

$$287 \quad \mathbb{E}(R_n^{(k)}) = \int_0^\infty (x - 1) e^{-x} \left(1 + \frac{x}{n}\right)^k dx,$$

288 which can be verified in a straightforward way, turns out to be advantageous. Expanding
 289 the integrand and distinguishing several cases yields the asymptotic results given in the
 290 theorem. ◀

3:10 Uncovering a random tree

291 We can even obtain the exact distribution of $R_n^{(k)}$. There are two different approaches
 292 we want to briefly sketch: for one, an explicit formula for the generating function $F =$
 293 $F(z, u, v)$ can be found either from manipulating the recursive description (15), or directly
 294 by decomposing \mathcal{F} as a tree forming the uncovered root cluster with a forest with covered
 295 roots attached. Either way, this yields

$$296 \quad F = T^\bullet(vXe^{-X}) + \frac{G}{1+u}, \quad \text{with} \quad X = \frac{uG}{1+u}.$$

297 Extracting coefficients via an application of the Lagrange inversion formula (see, e.g., [5,
 298 Theorem A.2]) then yields an explicit formula for $F_{n,k,m} = n![z^n u^k v^m]F(z, u, v)$, i.e., the
 299 number of labeled rooted trees with n vertices of which k are uncovered and m belong to the
 300 root cluster (for $0 \leq m \leq k \leq n$ and $n \geq 1$):

$$301 \quad F_{n,k,m} = \begin{cases} \binom{n-1}{k} n^{n-1}, & m = 0, \\ \binom{n}{m} \binom{n-m-1}{k-m} n^{n-k-1} m^m (n-m)^{k-m}, & m \geq 1. \end{cases}$$

302 From this formula, the probabilities $\mathbb{P}(R_n^{(k)} = r)$ given below can be obtained directly.

303 Alternatively, there is also a more combinatorial approach to determine these probabilities:
 304 there is an elementary formula enumerating trees where a specified set of vertices forms a
 305 cluster.

306 \triangleright **Claim 10.** The number of trees on $[n]$ that do not have any edges between the vertex sets
 307 $[r]$ and $[r+1, k]$, where additionally the induced subgraph on $[r]$ is a tree itself, is given by

$$308 \quad r^{r-1} n^{n-k-1} (n-k)(n-r)^{k-r-1}. \quad (17)$$

309 Proof of Claim 10. By Cayley's tree enumeration formula, the number of labeled trees on $[r]$
 310 is r^{r-2} . Hence, the statement of the claim is equivalent to

$$311 \quad r n^{n-k-1} (n-k)(n-r)^{k-r-1} \quad (18)$$

312 enumerating all rooted forests on $[n]$ whose roots are the vertices in $[r]$ that do not have any
 313 edges between $[r]$ and $[r+1, k]$. This can be proved bijectively with a construction similar
 314 to the Prüfer code, or by an application of Kirchhoff's Matrix-Tree Theorem [15, Theorem
 315 5.6.8]. Details can be found in the full version of this extended abstract. \triangleleft

316 As a consequence, the probability $\mathbb{P}(R_n^{(k)} = r)$ can be obtained by multiplying the number of
 317 these trees with $r \binom{k}{r}$ (to account for which vertices in $[k]$ form the root cluster and for the
 318 choice of the root), and then normalizing by n^{n-1} , the number of labeled rooted trees on n
 319 vertices.

320 \blacktriangleright **Theorem 11.** *The exact distribution of $R_n^{(k)}$ is characterized by the following probability*
 321 *mass function (p.m.f.), which is given as follows for $0 \leq m \leq k \leq n$ and $n \geq 1$ (and 0*
 322 *otherwise):*

$$323 \quad \mathbb{P}(R_n^{(k)} = m) = \begin{cases} 1 - \frac{k}{n}, & \text{for } m = 0, \\ \frac{m^m (n-k)(n-m)^{k-m-1}}{n^k} \binom{k}{m}, & \text{for } 1 \leq m \leq k < n, \\ 1, & \text{for } m = k = n. \end{cases}$$

324 *Depending on the growth of $k = k(n)$, we obtain the following limiting behavior:*

325 ■ *k* small, i.e., $k = o(n)$:

326
$$R_n^{(k)} \xrightarrow{d} 0.$$

327 ■ *k* in central region, i.e., $k \sim \alpha n$ with $0 < \alpha < 1$:

328 $R_n^{(k)} \xrightarrow{d} R_\alpha$, where the discrete r.v. R_α is characterized by its p.m.f.

329
$$\mathbb{P}(R_\alpha = m) =: p_m = \begin{cases} 1 - \alpha, & m = 0, \\ \frac{m^m}{m!} (1 - \alpha) \alpha^m e^{-\alpha m}, & m \geq 1, \end{cases}$$

330

331 or alternatively by the probability generating function $p(v) = \sum_{m \geq 0} p_m v^m = \frac{1 - \alpha}{1 - T^*(v \alpha e^{-\alpha})}$.

332 ■ *k* subcritically large, i.e., $k = n - d$ with $d = \omega(\sqrt{n})$ and $d = o(n)$:

333
$$\left(\frac{d}{n}\right)^2 \cdot R_n^{(k)} \xrightarrow{d} \text{GAMMA}\left(\frac{1}{2}, \frac{1}{2}\right),$$

334

335 where $\text{GAMMA}\left(\frac{1}{2}, \frac{1}{2}\right)$ is a Gamma-distribution characterized by its density $f(x) = \frac{1}{\sqrt{2\pi x}} e^{-\frac{x}{2}}$,
336 for $x > 0$.

337 ■ *k* critically large, i.e., $k = n - d$ with $d \sim c\sqrt{n}$ and $c > 0$:

338
$$\frac{1}{n} \cdot R_n^{(k)} \xrightarrow{d} R(c),$$

339

340 where the continuous r.v. $R(c)$ is characterized by its density $f_c(x) = \frac{1}{\sqrt{2\pi}} \frac{c}{\sqrt{x(1-x)}^{\frac{3}{2}}} e^{-\frac{c^2 x}{2(1-x)}}$,
341 for $0 < x < 1$.

342 ■ *k* supercritically large, i.e., $k = n - d$ with $d = \omega(1)$ and $d = o(\sqrt{n})$:

343
$$\frac{1}{d^2} \cdot (n - R_n^{(k)}) \xrightarrow{d} D,$$

344

345 where the continuous r.v. D is characterized by its density $f(x) = \frac{1}{\sqrt{2\pi x^{\frac{3}{2}}}} e^{-\frac{1}{2x}}$, $x > 0$.

346 ■ *k* supercritically large with fixed difference, i.e., $k = n - d$ with d fixed:

347
$$n - d - R_n^{(k)} \xrightarrow{d} D(d),$$

348

349 where the discrete r.v. $D(d)$ is characterized by the p.m.f.

350
$$\mathbb{P}(D(d) = j) =: p_j = e^{-d} \cdot \frac{d(d+j)^{j-1}}{j!} \cdot e^{-j}, \quad j \geq 0,$$

351

352 or alternatively via the probability generating function $p(v) = \sum_{j \geq 0} p_j v^j = e^{d(T^*(\frac{v}{e}) - 1)}$.

353 **Proof.** The probability mass function of $R_n^{(k)}$ follows from the considerations made before
354 the statement of the theorem. Due to its explicit nature, the limiting distribution results
355 stated in Theorem 11 can be obtained in a rather straightforward way by applying Stirling's
356 formula for the factorials after distinguishing several cases. ◀

357 ▶ **Remark 12.** Of course, for labeled trees, the distribution of $R_n^{(k)}$ matches with the dis-
358 tribution of the cluster size of a random vertex. Furthermore, by conditioning, one can
359 easily transfer the results of Theorem 11 to results for the size $S_n^{(k)}$ of the cluster of the k -th
360 uncovered vertex: $\mathbb{P}(S_n^{(k)} = m) = \mathbb{P}(R_n^{(k)} = m | R_n^{(k)} > 0) = \frac{n}{k} \cdot \mathbb{P}(R_n^{(k)} = m)$, for $m \geq 1$.

361 **4 Size of the largest uncovered component**

362 With knowledge about the behavior of the root cluster at our disposal, we return to non-
 363 rooted labeled trees and study the size of the largest cluster. To this aim, we introduce the
 364 random variable $X_{n,r}^{(k)}$ which models the number of components of size r after uncovering the
 365 vertices 1 to k in a uniformly random labeled tree of size n .

366 Formally, $X_{n,r}^{(k)}: \mathcal{T}_n \rightarrow \mathbb{Z}_{\geq 0}$. Note that we have, for all labeled trees $T \in \mathcal{T}_n$,

$$367 \sum_{r=1}^n r \cdot X_{n,r}^{(k)}(T) = k. \tag{19}$$

368 **► Theorem 13.** *Let $n, k, r \in \mathbb{Z}_{\geq 0}$ with $0 \leq r \leq k \leq n$. The expected number of connected*
 369 *components of size r after uncovering k vertices of a labeled tree of size n chosen uniformly*
 370 *at random is*

$$371 \mathbb{E}X_{n,r}^{(k)} = \binom{k}{r} \left(\frac{r}{n}\right)^{r-1} \left(1 - \frac{k}{n}\right) \left(1 - \frac{r}{n}\right)^{k-r-1}. \tag{20}$$

372 **Sketch of proof.** Observe that $X_{n,r}^{(k)}$ can be written as a sum of Bernoulli random variables

$$373 X_{n,r}^{(k)} = \sum_{\substack{S \subseteq [k] \\ |S|=r}} X_{n,S}^{(k)},$$

374 with $X_{n,S}^{(k)}$ being 0 or 1 depending on whether or not the vertices in S form a cluster after k
 375 uncover steps. By symmetry and linearity of the expected value, we have

$$376 \mathbb{E}X_{n,r}^{(k)} = \sum_{\substack{S \subseteq [k] \\ |S|=r}} \mathbb{E}X_{n,S}^{(k)} = \binom{k}{r} \mathbb{E}X_{n,[r]}^{(k)}.$$

377 A formula for the expected value on the right-hand side follows from Claim 10, and thus
 378 proves the theorem. ◀

379 In the spirit of the observation in (19), the formula in Claim 10 provides a combinatorial
 380 proof for the following summation identity.

381 **► Corollary 14.** *Let $n, k \in \mathbb{Z}_{\geq 0}$ with $0 \leq k \leq n$. Then, the identity*

$$382 \sum_{r=1}^k \binom{k}{r} r^r n^{n-k-1} (n-r)^{k-r-1} (n-k) = kn^{n-2} \tag{21}$$

383 *holds.*

384 **Proof.** The right-hand side enumerates the vertices in $[k]$ in all labeled trees on n vertices.
 385 The left-hand side does the same, with the summands enumerating the vertices in connected
 386 components of size r . ◀

387 **► Remark 15.** Observe that the identity in (21) can be rewritten as

$$388 \sum_{r=1}^k \binom{k}{r} r^r (n-r)^{k-r-1} = \frac{k}{n-k} n^{k-1},$$

389 which is a specialized form of Abel's Binomial Theorem—a classical, and well-known result;
 390 see, e.g., [14].

391 For a tree $T \in \mathcal{T}$, let $c_{\max}^{(k)}(T)$ denote the largest connected component of T after
 392 uncovering the first k vertices.

393 ► **Theorem 16.** *Let $n \in \mathbb{Z}_{\geq 0}$, and let $T_n \in \mathcal{T}_n$ be a tree chosen uniformly at random. Then
 394 the behavior of the random variable $c_{\max}^{(k)}(T_n)$ as $n \rightarrow \infty$ can be described as follows:*

- 395 ■ for $k = n - d$ with $d = \omega(\sqrt{n})$ (subcritical case), we have $c_{\max}^{(k)}(T_n)/n \xrightarrow{P} 0$.
- 396 ■ for $k = n - d$ with $d \sim c\sqrt{n}$ for a constant c (critical case), the rescaled random variable
 397 $c_{\max}^{(k)}(T_n)/n$ converges weakly to a (non-degenerate) continuous limiting distribution,
- 398 ■ for $k = n - d$ with $d = o(\sqrt{n})$ (supercritical case), we have $c_{\max}^{(k)}(T_n)/n \xrightarrow{P} 1$. With high
 399 probability, there is one “giant” component whose size is asymptotically equal to n .

400 **Sketch of proof.** For the subcritical case, we use the expected root cluster size from Theo-
 401 rem 9. Since a cluster of size r contains the root with probability $\frac{r}{n}$, we have

$$402 \frac{n^2}{d^2} \sim \mathbb{E}R_n^{(n-d)} = \sum_{r=0}^{n-d} \mathbb{E}X_{n,r}^{(n-d)} \cdot r \cdot \frac{r}{n} \geq \sum_{r=m}^{n-d} \mathbb{E}X_{n,r}^{(n-d)} \frac{r^2}{n} \geq \frac{m^2}{n} \sum_{r=m}^{n-d} \mathbb{E}X_{n,r}^{(n-d)}$$

$$403 \geq \frac{m^2}{n} \mathbb{P}(c_{\max}^{(n-d)}(T_n) \geq m).$$

405 This implies that

$$406 \mathbb{P}(c_{\max}^{(n-d)}(T_n) \geq m) = O\left(\frac{n^3}{d^2 m^2}\right),$$

407 so if $m = \epsilon n$ for any fixed $\epsilon > 0$, we have $\mathbb{P}(c_{\max}^{(n-d)}(T_n) \geq m) \rightarrow 0$.

408 In the critical case, we first consider the situation that there is a cluster that contains
 409 more than half of the vertices. Clearly, such a cluster must be the largest cluster and the
 410 only cluster of its size. Therefore, if $r = \rho n$ with $\rho > \frac{1}{2}$, we have

$$411 \mathbb{P}(c_{\max}^{(k)}(T_n) = r) = \mathbb{E}X_{n,r}^{(k)} \sim \frac{1}{n} \frac{e^{-\frac{c^2}{2} \frac{\rho}{1-\rho}}}{\sqrt{2\pi}} \frac{c}{(\rho(1-\rho))^{3/2}},$$

412 using Theorem 13 and Stirling’s approximation. Thus, for $\rho > \frac{1}{2}$,

$$413 \mathbb{P}(c_{\max}^{(k)}(T_n) \geq \rho n) \rightarrow \int_{\rho}^1 \frac{e^{-\frac{c^2}{2} \frac{t}{1-t}}}{\sqrt{2\pi}} \frac{c}{(t(1-t))^{3/2}} dt.$$

414 For $\rho \leq \frac{1}{2}$, we can modify this argument with a generalized version of Theorem 13 for
 415 several clusters and the inclusion-exclusion principle to prove convergence of $c_{\max}^{(k)}(T_n)/n$ to
 416 a continuous random variable with support $[0, 1]$. Details are left to the full version.

417 Finally, in the supercritical case, we recall the corresponding case for the size of the root
 418 cluster from Theorem 9. Using Markov’s inequality yields, for any $\epsilon > 0$,

$$419 \mathbb{P}(n - R_n^{(k)} \geq \epsilon n) \leq \frac{n - \mathbb{E}(R_n^{(k)})}{\epsilon n} \sim \frac{d\sqrt{n}}{\epsilon n} \xrightarrow{n \rightarrow \infty} 0.$$

420 Thus, the root cluster is the largest cluster of size $\sim n$ with high probability. Translating
 421 this from rooted to unrooted trees proves the theorem. ◀

422 — References —

- 423 1 David Aldous and Jim Pitman. The standard additive coalescent. *Ann. Probab.*, 26(4):1703–
424 1726, 1998. doi:10.1214/aop/1022855879.
- 425 2 Jean Bertoin. *Random fragmentation and coagulation processes*, volume 102 of *Cambridge*
426 *Studies in Advanced Mathematics*. Cambridge University Press, Cambridge, 2006. doi:
427 10.1017/CB09780511617768.
- 428 3 Patrick Billingsley. *Convergence of probability measures*. John Wiley & Sons, Inc., New York,
429 second edition, 1999. doi:10.1002/9780470316962.
- 430 4 P. Chassaing and G. Louchard. Phase transition for parking blocks, Brownian excursion and
431 coalescence. *Random Structures Algorithms*, 21(1):76–119, 2002. doi:10.1002/rsa.10039.
- 432 5 Philippe Flajolet and Robert Sedgewick. *Analytic combinatorics*. Cambridge University Press,
433 Cambridge, 2009. doi:10.1017/CB09780511801655.
- 434 6 Svante Janson. Random cutting and records in deterministic and random trees. *Random*
435 *Structures Algorithms*, 29(2):139–179, 2006. doi:10.1002/rsa.20086.
- 436 7 J. F. C. Kingman. The coalescent. *Stochastic Process. Appl.*, 13(3):235–248, 1982. doi:
437 10.1016/0304-4149(82)90011-4.
- 438 8 Achim Klenke. *Probability theory—a comprehensive course*. Universitext. Springer, Cham,
439 2020. Third edition. doi:10.1007/978-3-030-56402-5.
- 440 9 Jeremy L. Martin and Victor Reiner. Factorization of some weighted spanning tree enumerators.
441 *J. Combin. Theory Ser. A*, 104(2):287–300, 2003. doi:10.1016/j.jcta.2003.08.003.
- 442 10 Jim Pitman. Coalescent random forests. *J. Combin. Theory Ser. A*, 85(2):165–193, 1999.
443 doi:10.1006/jcta.1998.2919.
- 444 11 Jim Pitman. Coalescents with multiple collisions. *Ann. Probab.*, 27(4):1870–1902, 1999.
445 doi:10.1214/aop/1022677552.
- 446 12 Jeffery B. Remmel and S. Gill Williamson. Spanning trees and function classes. *Electron. J.*
447 *Combin.*, 9(1):Research Paper 34, 24, 2002. URL: [http://www.combinatorics.org/Volume_](http://www.combinatorics.org/Volume_9/Abstracts/v9i1r34.html)
448 [9/Abstracts/v9i1r34.html](http://www.combinatorics.org/Volume_9/Abstracts/v9i1r34.html).
- 449 13 Daniel Revuz and Marc Yor. *Continuous martingales and Brownian motion*, volume 293 of
450 *Grundlehren der mathematischen Wissenschaften*. Springer-Verlag, Berlin, third edition, 1999.
451 doi:10.1007/978-3-662-06400-9.
- 452 14 John Riordan. *Combinatorial identities*. John Wiley & Sons, Inc., New York, 1968.
- 453 15 Richard P. Stanley. *Enumerative combinatorics. Vol. 2*, volume 62 of *Cambridge Studies*
454 *in Advanced Mathematics*. Cambridge University Press, Cambridge, 1999. doi:10.1017/
455 CB09780511609589.

A Additional details on the proof of Theorem 5

A.1 Computations related to the martingale

With the help of the recursive description in (10), we can show that $(Y_j^{(n)})_{1 \leq j \leq n-1}$ is a martingale by computing

$$\begin{aligned} \mathbb{E}(Y_{j+1}^{(n)} | Y_j^{(n)}) &= \frac{\mathbb{E}(K_{j+1}^{(n)} | K_j^{(n)})}{n-j-1} - \frac{j(j+1)}{n(n-j-1)} = \frac{K_j^{(n)} + \frac{j+1}{n} + \frac{j-1-K_j^{(n)}}{n-j}}{n-j-1} - \frac{j(j+1)}{n(n-j-1)} \\ &= \frac{K_j^{(n)}}{n-j} - \frac{(j-1)j}{n(n-j)} = Y_j^{(n)}. \end{aligned}$$

We can also give an explicit expression for the variance of $Y_k^{(n)}$: recall that by (11), we have $\mathbb{V}Y_k^{(n)} = (n-k)^{-2} \mathbb{V}K_k^{(n)}$. Then, with (10) and the laws of total variance and total expectation we find the recurrence

$$\mathbb{V}K_{k+1}^{(n)} = \left(1 - \frac{1}{n-k}\right)^2 \mathbb{V}K_k^{(n)} + \frac{(n-k-1)(2n-k-1)k}{(n-k)n^2},$$

for $1 \leq k < n-1$ and $\mathbb{V}K_1^{(n)} = 0$. This allows us to conclude that

$$\mathbb{V}K_k^{(n)} = \sum_{j=1}^{k-1} \left(\frac{n-k}{n-j-1}\right)^2 \frac{(n-j-1)(2n-j-1)j}{(n-j)n^2} = \frac{k(k-1)(n-k)}{n^2}, \quad (22)$$

where the sum can be evaluated with the help of partial fractions and telescoping.

A.2 Proofs of auxiliary results

Proof of Lemma 7. In order to obtain the tightness condition, we show first that it can be reduced to an inequality for the martingale from the previous section. To this end, let us write $tn = j + \eta$, with $j \in \mathbb{Z}$ and $\eta \in [0, 1)$. A simple calculation shows that

$$\begin{aligned} Z^{(n)}(t) &= \frac{\tilde{K}_t^{(n)} - t^2 n}{\sqrt{n}} \\ &= \frac{(1-\eta)K_j^{(n)} + \eta K_{j+1}^{(n)} - (j+\eta)^2/n}{\sqrt{n}} \\ &= \frac{(1-\eta)(K_j^{(n)} - j(j-1)/n) + \eta(K_{j+1}^{(n)} - j(j+1)/n) - (j+\eta)^2/n}{\sqrt{n}} \\ &= (1-\eta) \frac{K_j^{(n)} - j(j-1)/n}{\sqrt{n}} + \eta \frac{K_{j+1}^{(n)} - j(j+1)/n}{\sqrt{n}} - \frac{j+\eta^2}{n^{3/2}}. \end{aligned}$$

The final fraction is bounded by 1, since $j + \eta^2 \leq j + \eta = tn \leq n$. It follows that

$$\sup_{t \in [0,1]} |Z^{(n)}(t)| \leq \sup_{0 \leq j \leq n} \left| \frac{K_j^{(n)} - j(j-1)/n}{\sqrt{n}} \right| + 1,$$

so

$$\begin{aligned} \mathbb{P}\left(\sup_{t \in [0,1]} |Z^{(n)}(t)| \geq C\right) &\leq \mathbb{P}\left(\sup_{0 \leq j \leq n} \left| \frac{K_j^{(n)} - j(j-1)/n}{\sqrt{n}} \right| \geq C-1\right) \\ &= \mathbb{P}\left(\sup_{1 \leq j \leq n-1} \left| \frac{Y_j^{(n)}(n-j)}{\sqrt{n}} \right| \geq C-1\right). \end{aligned} \quad (23)$$

3:16 Uncovering a random tree

486 Note here that we need not consider $j = 0$ and $j = n$ in the supremum, since $K_j^{(n)} -$
487 $j(j-1)/n = 0$ in either case. Since $(Y_j^{(n)})_{1 \leq j \leq n-1}$ is a martingale, we can use Doob's L^p -
488 inequality [8, Theorem 11.2]. For any real $C > 0$ and any fixed integer k with $1 \leq k \leq n-1$,
489 we have

$$490 \quad \mathbb{P}\left(\sup_{1 \leq j \leq k} |Y_j^{(n)}| \geq C\right) \leq \frac{\mathbb{V}Y_k^{(n)}}{C^2} = \frac{k(k-1)}{C^2(n-k)n^2}.$$

491 With this, we have all required prerequisites to prove tightness of $Z^{(n)}(t)$. We partition the
492 interval over which the supremum is taken in (23), apply the martingale inequality, and then
493 obtain the desired result after summing over all these upper bounds. For every integer $i > 0$,
494 let $I_i^{(n)} := [2^{-i}n, 2^{-i+1}n] \cap \mathbb{Z}$. We find

$$495 \quad \mathbb{P}\left(\sup_{n-j \in I_i^{(n)}} \left| \frac{Y_j^{(n)}(n-j)}{\sqrt{n}} \right| \geq C-1\right) \leq \mathbb{P}\left(\sup_{n-j \in I_i^{(n)}} |Y_j^{(n)}| 2^{-i+1} \sqrt{n} \geq C-1\right)$$

$$496 \quad = \mathbb{P}\left(\sup_{n-j \in I_i^{(n)}} |Y_j^{(n)}| \geq \frac{2^{i-1}(C-1)}{\sqrt{n}}\right)$$

$$497 \quad \leq \frac{n}{2^{2i-2}(C-1)^2} \mathbb{V}(Y_{n-[2^{-i}n]}^{(n)})$$

$$498 \quad \leq \frac{n}{2^{2i-2}(C-1)^2} \cdot \frac{2^i}{n} = \frac{4}{2^i(C-1)^2},$$

499 where in the last inequality we bounded the variance as follows:

$$500 \quad \mathbb{V}(Y_{n-[2^{-i}n]}^{(n)}) = \frac{\mathbb{V}(K_{n-[2^{-i}n]}^{(n)})}{[2^{-i}n]^2} = \frac{(n-[2^{-i}n])(n-[2^{-i}n]-1)[2^{-i}n]}{n^2[2^{-i}n]^2} \leq \frac{2^i}{n}.$$

501 Finally, the union bound together with the observation that $\sum_{i \geq 1} \frac{4}{2^i(C-1)^2} = 4(C-1)^{-2}$
502 yields the upper bound in (12) and therefore completes the proof. \blacktriangleleft

503 **Proof of Lemma 8.** As a consequence of Lemma 3 and Cayley's well-known enumeration
504 formula for labeled trees of size n , we find that the probability generating function of the
505 number of edge increments after $1 < j_1 < j_2 < \dots < j_r < n$ steps, respectively, is given by

$$506 \quad P_n(z_1, z_2, \dots, z_r) = \frac{E_n(z_1, z_2, \dots, z_r)}{n^{n-2}}$$

$$507 \quad = \prod_{i=1}^r \left(1 - \frac{j_r}{n} + \frac{j_i}{n} z_i + \sum_{h=i+1}^r \frac{j_h - j_{h-1}}{n} z_h\right)^{j_i - j_{i-1}}, \quad (24)$$

508 where $j_0 = 1$ for the sake of convenience. Now observe that $\Delta_j^{(n)}$ can be seen as a marginal
509 distribution of the sum of r independent, multinomially distributed random vectors: write
510 $t_i = j_i/n$ and consider $M_j \sim \text{Multi}(j_i - j_{i-1}, \mathbf{p}_i)$ where

$$511 \quad \mathbf{p}_i = (p_{i,0}, p_{i,1}, \dots, p_{i,r}) \in [0, 1]^r \quad \text{such that} \quad p_{i,h} = \begin{cases} 1 - t_r & \text{if } h = 0, \\ 0 & \text{if } 0 < h < i, \\ t_i & \text{if } h = i, \\ t_h - t_{h-1} & \text{otherwise.} \end{cases} \quad (25)$$

512 By construction, the probability generating function of M_i is then given by

$$513 \quad \left((1 - t_r)z_0 + t_i z_i + \sum_{h=i+1}^r (t_h - t_{h-1}) z_h\right)^{j_i - j_{i-1}},$$

516 so that the probability generating function of the sum $M_1 + \dots + M_r$ is a product that is
 517 very similar (and actually equal if we set $z_0 = 1$, which corresponds to marginalizing out the
 518 first component) to (24). In order to make the following arguments formally easier to read,
 519 and as the first component is not relevant for us at all, we slightly abuse notation and let M_i
 520 for $1 \leq i \leq r$ denote the corresponding marginalized multinomial distributions instead.

521 For the sake of convenience, we make a slight adjustment: instead of fixing the integer
 522 vector $\mathbf{j} = (j_1, \dots, j_r)$, we fix $\mathbf{t} = (t_1, \dots, t_r)$ with $0 < t_1 < \dots < t_r < 1$ and define $\mathbf{j} = \lfloor \mathbf{t}n \rfloor$.
 523 Here, n is considered to be sufficiently large so that the conditions for the corresponding
 524 integer vector, $1 < \lfloor t_1 n \rfloor < \dots < \lfloor t_r n \rfloor < n$, are still satisfied.

525 By the multivariate central limit theorem, it is well-known that a multinomially distributed
 526 random vector $M \sim \text{Multi}(n, \mathbf{p})$ converges, for $n \rightarrow \infty$ and after appropriate scaling, in
 527 distribution to a multivariate normal distribution,

$$528 \quad \frac{M - n\mathbf{p}}{\sqrt{n}} \xrightarrow{d} \mathcal{N}(\mathbf{0}, \text{diag}(\mathbf{p}) - \mathbf{p}^\top \mathbf{p}). \tag{26}$$

529 As a consequence, we find that

$$\begin{aligned} 530 \quad \frac{\Delta_{\lfloor \mathbf{t}n \rfloor}^{(n)} - \mathbb{E}\Delta_{\lfloor \mathbf{t}n \rfloor}^{(n)}}{\sqrt{n}} &= \frac{(M_1 + \dots + M_r) - \mathbb{E}(M_1 + \dots + M_r)}{\sqrt{n}} \\ 531 \quad &= (\sqrt{t_1} + O(n^{-1})) \frac{M_1 - \mathbb{E}M_1}{\sqrt{\lfloor t_1 n \rfloor}} + \dots \\ 532 \quad &\quad + (\sqrt{t_r - t_{r-1}} + O(n^{-1})) \frac{M_r - \mathbb{E}M_r}{\sqrt{\lfloor t_r n \rfloor - \lfloor t_{r-1} n \rfloor}} \\ 533 \quad &\xrightarrow[n \rightarrow \infty]{d} \sqrt{t_1} \mathcal{N}(\mathbf{0}, \Sigma_1) + \dots + \sqrt{t_r - t_{r-1}} \mathcal{N}(\mathbf{0}, \Sigma_r) \\ 534 \quad &= \mathcal{N}(\mathbf{0}, t_1 \Sigma_1 + \dots + (t_r - t_{r-1}) \Sigma_r), \\ 535 \end{aligned}$$

536 where the variance-covariance matrices are given by

$$537 \quad \Sigma_j = \text{diag}(\mathbf{p}_j) - \mathbf{p}_j^\top \mathbf{p}_j.$$

538 By a straightforward (linear) transformation consisting of taking partial sums, the random
 539 vector of increments $\Delta_{\lfloor \mathbf{t}n \rfloor}^{(n)}$ can be transformed into $\mathbf{K}_{\lfloor \mathbf{t}n \rfloor}^{(n)}$. This proves that $\mathbf{K}_{\lfloor \mathbf{t}n \rfloor}^{(n)}$ converges,
 540 after centering and rescaling, to a multivariate normal distribution.

541 The entries of the corresponding variance-covariance matrix can either be determined
 542 mechanically from the entries of $t_1 \Sigma_1 + \dots + (t_r - t_{r-1}) \Sigma_r$ by taking the partial summation
 543 into account, or alternatively, our observations concerning the martingale from Section A.1
 544 can be used. In particular, using (11), we find, for fixed $s, t \in [0, 1]$ with $s < t$, that

$$\begin{aligned} 545 \quad \text{Cov}\left(\frac{K_{\lfloor sn \rfloor}^{(n)} - \mathbb{E}K_{\lfloor sn \rfloor}^{(n)}}{\sqrt{n}}, \frac{K_{\lfloor tn \rfloor}^{(n)} - \mathbb{E}K_{\lfloor tn \rfloor}^{(n)}}{\sqrt{n}}\right) &= \frac{(n - \lfloor tn \rfloor)(n - \lfloor sn \rfloor)}{n} \mathbb{E}(Y_{\lfloor sn \rfloor}^{(n)} Y_{\lfloor tn \rfloor}^{(n)}) \\ 546 \quad &= (n(1-t)(1-s) + O(1)) \mathbb{E}(Y_{\lfloor sn \rfloor}^{(n)})^2 \\ 547 \quad &= s^2(1-t) + O(n^{-1}), \\ 548 \end{aligned}$$

549 where we made use of the martingale property, and that the second moment of $Y_j^{(n)}$ is equal
 550 to the variance $n^{-2}j(j-1)/(n-j)$. Ultimately, this verifies (14) and thus completes the
 551 proof. \blacktriangleleft