

# A note on the scaling limits of random Pólya trees

Bernhard Gittenberger\*

Emma Yu Jin\*

Michael Wallner\*

## Abstract

Panagiotou and Stuffer (arXiv:1502.07180v2) recently proved one important fact on their way to establish the scaling limits of random Pólya trees: a uniform random Pólya tree of size  $n$  consists of a conditioned critical Galton-Watson tree  $C_n$  and many small forests, where with probability tending to one as  $n$  tends to infinity, any forest  $F_n(v)$ , that is attached to a node  $v$  in  $C_n$ , is maximally of size  $|F_n(v)| = O(\log n)$ . Their proof used the framework of a Boltzmann sampler and deviation inequalities.

In this paper, first, we employ a unified framework in analytic combinatorics to prove this fact with additional improvements on the bound of  $|F_n(v)|$ , namely  $|F_n(v)| = \Theta(\log n)$ . Second, we give a combinatorial interpretation of the rational weights of these forests and the defining substitution process in terms of automorphisms associated to a given Pólya tree. Finally, we derive the limit probability that for a random node  $v$  the attached forest  $F_n(v)$  is of a given size.

## 1 Introduction and main results

First, we recall the asymptotic estimation of the number of Pólya trees with  $n$  nodes from the literature [9, 10, 12]. Second, we present Theorem 1.1 that leads to the proof of the scaling limits of random Pólya trees in [11].

**1.1 Pólya trees** A *Pólya tree* is a rooted unlabeled tree considered up to symmetry. The *size* of a tree is given by the number of its nodes. We denote by  $t_n$  the number of Pólya trees of size  $n$  and by  $T(z) = \sum_{n \geq 1} t_n z^n$  the corresponding ordinary generating function. By Pólya's enumeration theory [12], the generating function  $T(z)$  satisfies

$$(1.1) \quad T(z) = z \exp \left( \sum_{i=1}^{\infty} \frac{T(z^i)}{i} \right).$$

The first few terms of  $T(z)$  are then

$$(1.2) \quad T(z) = z + z^2 + 2z^3 + 4z^4 + 9z^5 + 20z^6 + 48z^7 + 115z^8 + 286z^9 + 719z^{10} + \dots,$$

(see OEIS A000081, [13]). By differentiating both sides of (1.1) with respect to  $z$ , one can derive a recurrence relation of  $t_n$  (see [9, Chapter 29] and [10]), which is

$$t_n = \frac{1}{n-1} \sum_{i=1}^{n-1} t_{n-i} \sum_{m|i} m t_m, \quad \text{for } n > 1,$$

and  $t_1 = 1$ . Pólya [12] showed that the radius of convergence  $\rho$  of  $T(z)$  satisfies  $0 < \rho < 1$  and that  $\rho$  is the unique singularity on the circle of convergence  $|z| = \rho$ . Subsequently, Otter [10] proved that  $T(\rho) = 1$  as well as the singular expansion

$$(1.3) \quad T(z) = 1 - b(\rho - z)^{1/2} + c(\rho - z) + \mathcal{O}((\rho - z)^{3/2}),$$

where  $\rho \approx 0.3383219$ ,  $b \approx 2.68112$  and  $c = b^2/3 \approx 2.39614$ .

By transfer theorems [5] he derived

$$t_n = \frac{b\sqrt{\rho}}{2\sqrt{\pi}} \frac{\rho^{-n}}{\sqrt{n^3}} \left( 1 + \mathcal{O}\left(\frac{1}{n}\right) \right).$$

We will see that  $T(z)$  is connected with the *exponential generating function* of Cayley trees. “With a minor abuse of notation” (cf. [7, Ex. 10.2]), Cayley trees belong to the class of *simply generated trees*. Simply generated trees have been introduced by Meir and Moon [8] to describe a weighted version of rooted trees. They are defined by the functional equation

$$y(z) = z\Phi(y(z)), \quad \text{with} \\ \Phi(z) = \sum_{j \geq 0} \phi_j z^j, \quad \phi_j \geq 0.$$

The power series  $y(x) = \sum_{n \geq 1} y_n x^n$  has non-negative coefficients and is the generating function of *weighted simply generated trees*. One usually assumes that  $\phi_0 > 0$  and  $\phi_j > 0$  for some  $j \geq 2$  to exclude the trivial cases. In particular, in the above-mentioned sense, *Cayley trees* can be seen as simply generated trees which are characterized by  $\Phi(z) = \exp(z)$ . It is well known that the number of rooted Cayley trees of size  $n$  is given by  $n^{n-1}$ .

\*Institut für Diskrete Mathematik und Geometrie, Technische Universität Wien, Wiedner Hauptstr. 8–10/104, 1040 Vienna, Austria. Corresponding author: Michael Wallner. Emails: gittenberger@dmg.tuwien.ac.at; yu.jin@tuwien.ac.at; michael.wallner@tuwien.ac.at.

Let

$$C(z) = \sum_{n \geq 0} n^{n-1} \frac{z^n}{n!},$$

be the associated exponential generating function.

Then, by construction it satisfies the functional equation

$$C(z) = ze^{C(z)}.$$

In contrast, Pólya trees are not simply generated (see [4] for a simple proof of this fact). Note that though  $T(z)$  and  $C(z)$  are closely related, Pólya trees are not related to Cayley trees in a strict sense, but to a certain class of weighted unlabeled trees which will be called  $C$ -trees in the sequel and have the ordinary generating function  $C(z)$ . This is precisely the simply generated tree associated with  $\Phi(z) = \exp(z)$ , now in the strict sense of the definition of simply generated trees.

In order to analyze the dominant singularity of  $T(z)$ , we follow [10, 12], see also [5, Chapter VII.5], and we rewrite (1.2) into

$$(1.4) \quad T(z) = ze^{T(z)}D(z), \quad \text{where}$$

$$D(z) = \sum_{n \geq 0} d_n z^n = \exp\left(\sum_{i=2}^{\infty} \frac{T(z^i)}{i}\right).$$

We observe that  $D(z)$  is analytic for  $|z| < \sqrt{\rho} < 1$  and that  $\sqrt{\rho} > \rho$ . From (1.4) it follows that  $T(z)$  can be expressed in terms of the generating function of Cayley trees: Indeed, assume that  $T(z)$  is a function  $H(zD(z))$  depending on  $zD(z)$ . By (1.4) this is equivalent to  $H(x) = x \exp(H(x))$ . Yet, this is the functional equation for the generating function of Cayley trees. As this functional equation has a unique power series solution we have  $H(x) = C(x)$ , and we just proved

$$(1.5) \quad T(z) = C(zD(z)).$$

Note that  $T(z) = C(zD(z))$  is a case of a super-critical composition schema which is characterized by the fact that the dominant singularity of  $T(z)$  is strictly smaller than that of  $D(z)$ . In other words, the dominant singularity  $\rho$  of  $T(z)$  is determined by the outer function  $C(z)$ . Indeed,  $\rho D(\rho) = e^{-1}$ , because  $e^{-1}$  is the unique dominant singularity of  $C(z)$ .

Let us introduce two new classes of weighted combinatorial structures:  $D$ -forests and  $C$ -trees. We set  $d_n = [z^n]D(z)$  which is the *accumulated weight* of all  $D$ -forests of size  $n$ . These are weighted forests of Pólya trees which are constrained to contain for every Pólya tree at least two identical copies or none. In other

words, if a tree appears in a  $D$ -forest it has to appear at least twice. From (1.2) and (1.4) one gets its first values

$$(1.6) \quad D(z) = \sum_{n=0}^{\infty} d_n z^n$$

$$= 1 + \frac{1}{2}z^2 + \frac{1}{3}z^3 + \frac{7}{8}z^4 + \frac{11}{30}z^5$$

$$+ \frac{281}{144}z^6 + \frac{449}{840}z^7 + \dots$$

The weights are defined in such a way that composition scheme (1.5) is satisfied. In Theorem 1.2 we will make these weights explicit. From (1.4) we can derive a recursion of  $d_n$ . We get

$$d_n = \frac{1}{n} \sum_{i=2}^n d_{n-i} \sum_{\substack{m|i \\ m \neq i}} m t_m, \quad \text{for } n \geq 2,$$

as well as  $d_0 = 1$ , and  $d_1 = 0$ .

The second concept is the one of  $C$ -trees, which are weighted Pólya trees. The weight is defined by the composition (1.5). Let  $c_n = [z^n]C(z) = \frac{n^{n-1}}{n!}$  be the *accumulated weight* of all  $C$ -trees of size  $n$ . In other words, we interpret the exponential generating function of Cayley trees  $C(z)$  as an ordinary generating function of *weighted* objects:

$$C(z) = \sum_{n \geq 0} \frac{n^{n-1}}{n!} z^n.$$

Informally speaking, the composition (1.5) can be interpreted as such that a Pólya tree is constructed from a  $C$ -tree where a  $D$ -forest is attached to each node.

This construction is in general not bijective, because the  $D$ -forests consist of Pólya trees and are not distinguishable from the underlying Pólya tree, see Figure 1. In general there are different decompositions of a given Pólya tree into a  $C$ -tree and  $D$ -forests. Theorem 1.2 will give a probabilistic interpretation derived from the automorphism group of a Pólya tree (see also Example 3.2).

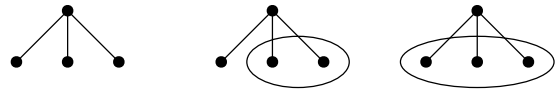


Figure 1: The decomposition of a Pólya tree with 4 nodes into a  $C$ -tree (non-circled nodes) and  $D$ -forests (circled nodes). For this Pólya tree there are 3 different decompositions.

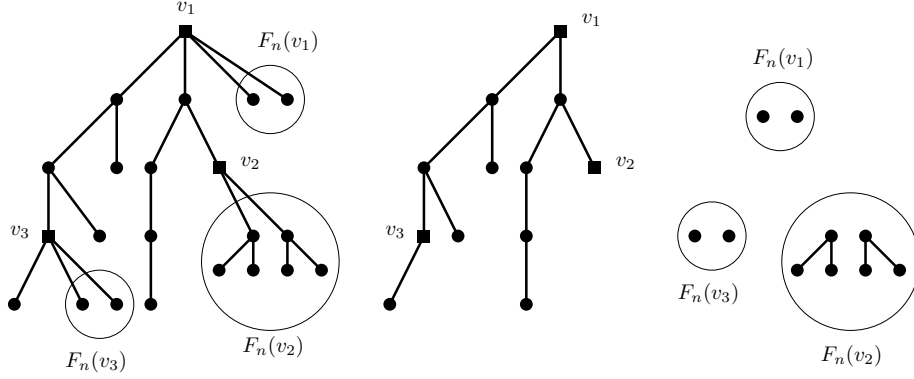


Figure 2: A random Pólya tree  $T_n$  (left), a (possible)  $C$ -tree  $C_n$  (middle) that is contained in  $T_n$  where all  $D$ -forests  $F_n(v)$ , except  $F_n(v_1), F_n(v_2), F_n(v_3)$  (right), are empty.

**1.2 Main results** Consider a random Pólya tree of size  $n$ , denoted by  $T_n$ , which is a tree that is selected uniformly at random from all Pólya trees with  $n$  vertices. We use  $C_n$  to denote the random  $C$ -tree that is contained in a random Pólya tree  $T_n$ . For every vertex  $v$  of  $C_n$ , we use  $F_n(v)$  to denote the  $D$ -forest that is attached to the vertex  $v$  in  $T_n$ , see Figure 2.

Let  $L_n$  be the maximal size of a  $D$ -forest contained in  $T_n$ , that is,  $|F_n(v)| \leq L_n$  holds for all  $v$  of  $C_n$  and the inequality is sharp. For the upper bound see also [11, Eq. (5.5)].

**THEOREM 1.1.** For  $0 < s < 1$ ,

$$(1.7) \quad \begin{aligned} (1 - (\log n)^{-s}) \left( \frac{-2 \log n}{\log \rho} \right) &\leq L_n \leq \\ (1 + (\log n)^{-s}) \left( \frac{-2 \log n}{\log \rho} \right) & \end{aligned}$$

holds with probability  $1 - o(1)$ .

Our first main result is a new proof of Theorem 1.1 by applying the unified framework of Gourdon [6]. Our second main result is a combinatorial interpretation of all weights on the  $D$ -forests and  $C$ -trees in terms of automorphisms associated to a given Pólya tree.

Let  $c_{n,k}$  denote the cumulative weight of all  $C$ -trees of size  $k$  that are contained in Pólya trees of size  $n$ . By  $t_{c,n}(u)$  and  $T_c(z, u)$  we denote the corresponding generating function and the bivariate generating function of  $(c_{n,k})_{n,k \geq 0}$ , respectively, that is,

$$\begin{aligned} t_{c,n}(u) &= \sum_{k=1}^n c_{n,k} u^k \quad \text{and} \\ T_c(z, u) &= \sum_{n \geq 0} t_{c,n}(u) z^n. \end{aligned}$$

Note that  $c_{n,k}$  is in general not an integer. By marking the nodes of all  $C$ -trees in Pólya trees we find a functional equation for the bivariate generating function  $T_c(z, u)$ , which is

$$(1.8) \quad \begin{aligned} T_c(z, u) &= zu \exp(T_c(z, u)) \exp\left(\sum_{i=2}^{\infty} \frac{T(z^i)}{i}\right) \\ &= zu \exp(T_c(z, u)) D(z). \end{aligned}$$

For a given permutation  $\sigma$  let  $\sigma_1$  be the number of fixed points of  $\sigma$ . Our second main result is the following:

**THEOREM 1.2.** Let  $\mathcal{T}$  be the set of all Pólya trees, and  $\text{MSET}^{(\geq 2)}(\mathcal{T})$  be the multiset (or forest) of Pólya trees where each tree appears at least twice if it appears at all. Then the cumulative weight  $d_n$  (defined in (1.6)) of all such forests of size  $n$  satisfies

$$d_n = \sum_{\substack{F \in \text{MSET}^{(\geq 2)}(\mathcal{T}) \\ |F|=n}} \frac{|\{\sigma \in \text{Aut}(F) \mid \sigma_1 = 0\}|}{|\text{Aut}(F)|}$$

where  $\text{Aut}(F)$  is the automorphism group of  $F$  (see Definition 3.1 section 3). Furthermore, the polynomial associated to  $C$ -trees in Pólya trees of size  $n$  is given by

$$\begin{aligned} t_{c,n}(u) &= \sum_{T \in \mathcal{T}, |T|=n} t_T(u), \quad \text{where} \\ t_T(u) &= \frac{1}{|\text{Aut}(T)|} \sum_{\sigma \in \text{Aut}(T)} u^{\sigma_1}. \end{aligned}$$

In particular, for all  $T \in \mathcal{T}$ , it holds that  $t'_T(1) = |\mathcal{P}(T)|$  where  $\mathcal{P}(T)$  is the set of all trees which are obtained by pointing (or coloring) one single node in  $T$ .

For a given Pólya tree  $T$  the polynomial  $t_T(u)$  gives rise to a probabilistic interpretation of the composition

scheme (1.5). For a given tree  $T$ , the weight of  $u^k$  in the polynomial  $t_T(u)$  is the probability that the underlying  $C$ -tree is of size  $k$ . In other words,  $t_T(u)$  is the probability generating function of the random variable  $C_T$  of the number of  $C$ -tree nodes in the tree  $T$  defined by

$$(1.9) \quad \mathbb{P}(C_T = k) := [u^k]t_T(u).$$

This random variable  $C_T$  is a refinement of  $T_n$  in the sense that

$$\mathbb{P}(C_T = k) = \mathbb{P}(|C_n| = k \mid T_n = T).$$

Finally, we derive the limiting probability that for a random node  $v$  the attached forest  $F_n(v)$  is of a given size. This result is consistent with the Boltzmann sampler from [11]. The precise statement of our third main result is the following:

**THEOREM 1.3.** *The generating function  $T^{[m]}(z, u)$  of Pólya trees, where each vertex is marked by  $z$ , and each weighted  $D$ -forest of size  $m$  is marked by  $u$ , is given by*

$$(1.10) \quad T^{[m]}(z, u) = C(uzd_m z^m + z(D(z) - d_m z^m)),$$

where  $d_m = [z^m]D(z)$ . The probability that the  $D$ -forest  $F_n(v)$  attached to a random  $C$ -tree node  $v$  is of size  $m$  is given by

$$\mathbb{P}(|F_n(v)| = m) = \frac{d_m \rho^m}{D(\rho)} (1 + \mathcal{O}(n^{-1})).$$

**1.3 Paper outline** The paper is organized as follows. In Section 2 we prove Theorem 1.1 and discuss the size of the  $C$ -tree  $C_n$  in a random Pólya tree  $T_n$ . In Section 3 we prove Theorems 1.2 and 1.3. In Section 4 we conclude with final remarks.

## 2 The maximal size of a $D$ -forest

We will use the generating function approach from [6] to analyze the maximal size  $L_n$  of  $D$ -forests in a random Pólya tree  $T_n$ , which provides a new proof of Theorem 1.1. Following the same approach, we can establish a central limit theorem for the random variable  $|C_n|$ , which has been done in [14] for the more general random  $\mathcal{R}$ -enriched trees.

*Proof of Theorem 1.1.* In (5.5) of [11], only an upper bound of  $L_n$  is given. By directly applying Gourdon's results (Theorem 4 and Corollary 3 of [6]) for the super-critical composition schema, we find that for any positive  $m$ ,

$$\mathbb{P}[L_n \leq m] = \exp\left(-\frac{c_1 n}{m^{3/2}} \rho^{m/2}\right) (1 + \mathcal{O}(\exp(-m\varepsilon))),$$

$$c_1 \sim \frac{b}{2\sqrt{\pi}(1 - \sqrt{\rho})(D(\rho) + \rho D'(\rho))},$$

as  $n \rightarrow \infty$ . Moreover, the maximal size  $L_n$  satisfies asymptotically, as  $n \rightarrow \infty$ ,

$$\mathbb{E}L_n = -\frac{2 \log n}{\log \rho} - \frac{3}{2} \frac{2}{\log \rho} \log \log n + \mathcal{O}(1) \quad \text{and}$$

$$\text{Var } L_n = \mathcal{O}(1).$$

By using Chebyshev's inequality, one can prove that  $L_n$  is highly concentrated around the mean  $\mathbb{E}L_n$ . We set  $\varepsilon_n = (\log n)^{-s}$  where  $0 < s < 1$  and we get

$$\mathbb{P}(|L_n - \mathbb{E}L_n| \geq \varepsilon_n \cdot \mathbb{E}L_n) \leq \frac{\text{Var } L_n}{\varepsilon_n^2 \cdot (\mathbb{E}L_n)^2} = o(1),$$

which means that Relation (1.7) holds with probability  $1 - o(1)$ .  $\square$

It was shown in [14] that the size  $|C_n|$  of the  $C$ -tree  $C_n$  in  $T_n$  satisfies a central limit theorem and  $|C_n| = \Theta(n)$  holds with probability  $1 - o(1)$ . In particular see [14, Eq. (3.9) and (3.10)], and [11, Eq. (5.6)]. The precise statement is the following.

**THEOREM 2.1.** *The size of the  $C$ -tree  $|C_n|$  in a random Pólya tree  $T_n$  of size  $n$  satisfies a central limit theorem where the expected value  $\mathbb{E}|C_n|$  and the variance  $\text{Var } |C_n|$  are asymptotically*

$$\mathbb{E}|C_n| = \frac{2n}{b^2 \rho} (1 + \mathcal{O}(n^{-1})), \quad \text{and}$$

$$\text{Var } |C_n| = \frac{11n}{12b^2 \rho} (1 + \mathcal{O}(n^{-1})).$$

Furthermore, for any  $s$  such that  $0 < s < 1/2$ , with probability  $1 - o(1)$  it holds that

$$(2.11) \quad (1 - n^{-s}) \frac{2n}{b^2 \rho} \leq |C_n| \leq (1 + n^{-s}) \frac{2n}{b^2 \rho}.$$

Random Pólya trees belong to the class of random  $\mathcal{R}$ -enriched trees and we refer the readers to [14] for the proof of Theorem 2.1 in the general setting. Here we provide a proof of Theorem 2.1 to show the connection between a bivariate generating function and the normal distribution and to emphasize the simplifications for the concrete values of the expected value and variance in this case.

*Proof of Theorem 2.1 (see also [14]).* It follows from [3, Th. 2.23] that the random variable  $|C_n|$  satisfies a central limit theorem. In the present case, we set  $F(z, y, u) = zu \exp(y)D(z)$ . It is easy to verify that  $F(z, y, u)$  is an analytic function when  $z$  and  $y$  are near 0 and that  $F(0, y, u) \equiv 0$ ,  $F(x, 0, u) \neq 0$  and all coefficients  $[z^n y^m]F(z, y, 1)$  are real and non-negative. From [3, Th. 2.23] we know that  $T_c(z, u)$  is the unique solution of the functional identity  $y = F(z, y, u)$ . Since all coefficients of  $F_y(z, y, 1)$  are non-negative and the

coefficients of  $T(z)$  are positive as well as monotonically increasing, this implies that  $(\rho, T(\rho), 1)$  is the unique solution of  $F_y(z, y, 1) = 1$ , which leads to the fact that  $T(\rho) = 1$ . Moreover, the expected value is

$$\begin{aligned} \mathbb{E}|C_n| &= \frac{nF_u(z, y, u)}{\rho F_z(z, y, u)} \\ &= \frac{[z^n] \partial_u T_c(z, u)|_{u=1}}{[z^n] T(z)} \\ &= \left( [z^n] \frac{T(z)}{1 - T(z)} \right) ([z^n] T(z))^{-1} \\ &= \frac{2n}{b^2 \rho} \left( 1 + \mathcal{O}\left(\frac{1}{n}\right) \right). \end{aligned}$$

The asymptotics are directly derived from (1.3). Likewise, we can compute the variance

$$\begin{aligned} \text{Var}|C_n| &= \frac{[z^n] T(z) (1 - T(z))^{-3}}{[z^n] T(z)} - (\mathbb{E}|C_n|)^2 \\ &= \frac{11n}{12b^2 \rho} (1 + \mathcal{O}(n^{-1})). \end{aligned}$$

Furthermore,  $|C_n|$  is highly concentrated around  $\mathbb{E}|C_n|$ , which can be proved again by using Chebyshev's inequality. We set  $\varepsilon_n = n^{-s}$  where  $0 < s < 1/2$  and get

$$\begin{aligned} \mathbb{P}(|C_n| - \mathbb{E}|C_n| \geq \varepsilon_n \cdot \mathbb{E}|C_n|) &\leq \frac{\text{Var}|C_n|}{\varepsilon_n^2 \cdot (\mathbb{E}|C_n|)^2} \\ &= \mathcal{O}(n^{2s-1}) = o(1), \end{aligned}$$

which yields (2.11).  $\square$

As a simple corollary, we also get the total size of all weighted  $D$ -forests in  $T_n$ . Let  $\mathcal{D}_n$  denote the union of all  $D$ -forests in a random Pólya tree  $T_n$  of size  $n$ .

**COROLLARY 2.1.** *The size of weighted  $D$ -forests in a random Pólya tree of size  $n$  satisfies a central limit theorem where the expected value  $\mathbb{E}|\mathcal{D}_n|$  and the variance  $\text{Var}|\mathcal{D}_n|$  are asymptotically*

$$\begin{aligned} \mathbb{E}|\mathcal{D}_n| &= n \left( 1 - \frac{2}{b^2 \rho} \right) (1 + \mathcal{O}(n^{-1})), \quad \text{and} \\ \text{Var}|\mathcal{D}_n| &= \frac{11n}{12b^2 \rho} (1 + \mathcal{O}(n^{-1})). \end{aligned}$$

Theorem 2.1 and Corollary 2.1 tell us that a random Pólya tree  $T_n$  consists mostly of a  $C$ -tree (proportion  $\frac{2}{b^2 \rho}$  comprising  $\approx 82.2\%$  of the nodes) and to a small part of  $D$ -forests (proportion  $1 - \frac{2}{b^2 \rho}$  comprising  $\approx 17.8\%$  of the nodes). Furthermore, the average size of a  $D$ -forest  $F_n(v)$  attached to a random  $C$ -tree vertex in  $T_n$  is  $\frac{b^2 \rho}{2} - 1 \approx 0.216$ , which indicates that on average the  $D$ -forest  $F_n(v)$  is very small, although the maximal size of all  $D$ -forests in a random Pólya tree  $T_n$  reaches  $\Theta(\log n)$ .

**REMARK 2.1.** *Let us describe the connection of (1.5) to the Boltzmann sampler from [11]. We know that  $F(z, y, 1) = z\Phi(y)D(z)$  where  $\Phi(x) = \exp(x)$  and  $y = T(z)$ . By dividing both sides of this equation by  $y = T(z)$ , one obtains from (1.4) that*

$$1 = \frac{zD(z)}{T(z)} \exp(T(z)) = \exp(-T(z)) \sum_{k \geq 0} \frac{T^k(z)}{k!},$$

which implies that in the Boltzmann sampler  $\Gamma T(x)$ , the number of offspring contained in the  $C$ -tree  $C_n$  is Poisson distributed with parameter  $T(x)$ . As an immediate result, this random  $C$ -tree  $C_n$  contained in the Boltzmann sampler  $\Gamma T(\rho)$  is a critical Galton-Watson tree since the expected number of offspring is  $F_y(z, y, 1) = 1$  which holds only when  $(z, y) = (\rho, 1)$ .

### 3 $D$ -forests and $C$ -trees

In order to get a better understanding of  $D$ -forests and  $C$ -trees, we need to return to the original proof of Pólya on the number of Pólya trees [12]. The important step is the treatment of tree automorphisms by the cycle index. Let us recall what it means that two graphs are isomorphic.

**DEFINITION 3.1.** *Two graphs  $G_1$  and  $G_2$  are isomorphic if there exists a bijection between the vertex sets of  $G_1$  and  $G_2$ ,  $f : V(G_1) \rightarrow V(G_2)$  such that two vertices  $v$  and  $w$  of  $G_1$  are adjacent if and only if  $f(v)$  and  $f(w)$  are adjacent in  $G_2$ . If  $G_1 = G_2$  we call the bijection  $f$  an automorphism. The automorphism group of the graph  $G_1$  is denoted by  $\text{Aut}(G_1)$ .*

For any permutation  $\sigma$ , let  $\sigma_i$  be the number of cycles of length  $i$  of  $\sigma$ . We define the *type* of  $\sigma$ , to be the sequence  $(\sigma_1, \sigma_2, \dots, \sigma_k)$  if  $\sigma \in S_k$ . Note that  $k = \sum_{i=1}^k i\sigma_i$ .

**DEFINITION 3.2. (CYCLE INDEX)** *Let  $G$  be a subgroup of the symmetric group  $S_k$ . Then, the cycle index is*

$$Z(G; s_1, s_2, \dots, s_k) = \frac{1}{|G|} \sum_{\sigma \in G} s_1^{\sigma_1} s_2^{\sigma_2} \dots s_k^{\sigma_k}.$$

Now we are ready to prove Theorem 1.2.

**3.1 Proof of Theorem 1.2** By Pólya's enumeration theory [12], the generating function  $T(z)$  satisfies the functional equation

$$\begin{aligned} T(z) &= z \sum_{k \geq 0} Z(S_k; T(z), T(z^2), \dots, T(z^k)) \\ &= z \sum_{k \geq 0} \frac{1}{k!} \sum_{\sigma \in S_k} (T(z))^{\sigma_1} (T(z^2))^{\sigma_2} \dots (T(z^k))^{\sigma_k}, \end{aligned}$$

which can be simplified to (1.1), the starting point of our research, by a simple calculation. However, this shows that the generating function of  $D$ -forests from (1.4) is given by

$$\begin{aligned} D(z) &= \exp\left(\sum_{i=2}^{\infty} \frac{T(z^i)}{i}\right) \\ &= \sum_{k \geq 0} Z(S_k; 0, T(z^2), \dots, T(z^k)) \\ &= \sum_{k \geq 0} \frac{1}{k!} \sum_{\sigma \in S_k, \sigma_1=0} (T(z^2))^{\sigma_2} \dots (T(z^k))^{\sigma_k}. \end{aligned}$$

This representation enables us to interpret the weights  $d_n$  of  $D$ -forests of size  $n$ : A  $D$ -forest of size  $n$  is a multiset of  $k$  Pólya trees, where every tree occurs at least twice. Its weight is given by the ratio of fixed point free automorphisms over the total number of automorphisms. Equivalently, it is given by the number of fixed point free permutations  $\sigma \in S_k$  of these trees rescaled by the total number of orderings  $k!$ .

Let  $\mathcal{T}$  be the set of all Pólya trees and  $\text{MSET}^{(\geq 2)}(\mathcal{T})$  be the multiset of Pólya trees where each tree appears at least twice if it appears at all. Combinatorially, this is a forest without unique trees. Then, their weights are given by

$$d_n = \sum_{\substack{F \in \text{MSET}^{(\geq 2)}(\mathcal{T}) \\ |F|=n}} \frac{|\{\sigma \in \text{Aut}(F) \mid \sigma_1 = 0\}|}{|\text{Aut}(F)|}.$$

**EXAMPLE 3.1.** *The smallest  $D$ -forest is of size 2, and it consists of a pair of single nodes. There is just one fixed point free automorphism on this forest, thus  $d_2 = 1/2$ . For  $n = 3$  the forest consists of 3 single nodes. The fixed point free permutations are the 3-cycles, thus  $d_3 = 2/6 = 1/3$ . The case  $n = 4$  is more interesting. A forest consists either of 4 single nodes, or of 2 identical trees, each consisting of 2 nodes and one edge. In the first case we have 6 4-cycles and 3 pairs of transpositions. In the second case we have 1 transposition swapping the two trees. Thus,  $d_4 = \frac{6+3}{24} + \frac{1}{2} = \frac{7}{8}$ .*

These results also yield a natural interpretation of  $C$ -trees. We recall that by definition

$$T_c(z, u) = \sum_{n \geq 0} t_{c,n}(u) z^n,$$

where  $t_{c,n}(u) = \sum_k c_{n,k} u^k$  is the polynomial marking the  $C$ -trees in Pólya trees of size  $n$ . From the decom-

positions (1.5) and (1.8) we get the first few terms:

$$\begin{aligned} t_{c,1}(u) &= u, \\ t_{c,2}(u) &= u^2, \\ t_{c,3}(u) &= \frac{3}{2}u^3 + \frac{1}{2}u, \\ t_{c,4}(u) &= \frac{8}{3}u^4 + u^2 + \frac{1}{3}u. \end{aligned}$$

Evaluating these polynomials at  $u = 1$  obviously returns  $t_{c,n}(1) = t_n$ , which is the number of Pólya trees of size  $n$ . Their coefficients, however, are weighted sums depending on the number of  $C$ -tree nodes. For a given Pólya tree there are in general several ways to decide what is a  $C$ -tree node and what is a  $D$ -forest node. The possible choices are encoded in the automorphisms of the tree, and these are responsible for the above weights as well.

Let  $T$  be a Pólya tree, and  $\text{Aut}(T)$  be its automorphism group. For an automorphism  $\sigma \in \text{Aut}(T)$  the nodes which are fixed points of  $\sigma$  are  $C$ -tree nodes. All other nodes are part of  $D$ -forests. Summing over all automorphisms and normalizing by the total number gives the  $C$ -tree generating polynomial for  $T$ :

$$(3.12) \quad \begin{aligned} t_T(u) &= Z(\text{Aut}(T); u, 1, \dots, 1) \\ &= \frac{1}{|\text{Aut}(T)|} \sum_{\sigma \in \text{Aut}(T)} u^{\sigma_1}. \end{aligned}$$

The polynomial of  $C$ -trees in Pólya trees of size  $n$  is then given by

$$t_{c,n}(u) = \sum_{T \in \mathcal{T}, |T|=n} t_T(u).$$

**EXAMPLE 3.2.** *For  $n = 3$  we have 2 Pólya trees, namely the chain  $T_1$  and the cherry  $T_2$ . Thus,  $\text{Aut}(T_1) = \{id\}$ , and  $\text{Aut}(T_2) = \{id, \sigma\}$ , where  $\sigma$  swaps the two leaves but the root is unchanged. Thus,*

$$\begin{aligned} t_{T_1}(u) &= u^3, \\ t_{T_2}(u) &= \frac{1}{2}(u^3 + u). \end{aligned}$$

*For  $n = 4$  we have 4 Pólya trees shown in Figure 3. Their automorphism groups are given by*

$$\begin{aligned} \text{Aut}(T_1) &= \{id\}, \\ \text{Aut}(T_2) &= \{id\}, \\ \text{Aut}(T_3) &= \{id, (v_3 v_4)\} \cong S_2, \\ \text{Aut}(T_4) &= \{id, (v_2 v_3), (v_3 v_4), (v_2 v_4), \\ &\quad (v_2 v_3 v_4), (v_2 v_4 v_3)\} \cong S_3. \end{aligned}$$

This gives

$$\begin{aligned} t_{T_1}(u) &= u^4, \\ t_{T_2}(u) &= u^4, \\ t_{T_3}(u) &= \frac{1}{2}(u^4 + u^2), \\ t_{T_4}(u) &= \frac{1}{6}(u^4 + 3u^2 + 2u). \end{aligned}$$

This enables us to give a probabilistic interpretation of the composition scheme (1.5). For a given tree the weight of  $u^k$  is the probability that the underlying  $C$ -tree is of size  $k$ . In particular,  $T_1$  and  $T_2$  do not have  $D$ -forests. The tree  $T_3$  consists of a  $C$ -tree with 4 or with 2 nodes, each case with probability  $1/2$ . In the second case, as there is only one possibility for the  $D$ -forest, it consists of the pair of single nodes which are the leaves. Finally, the tree  $T_4$  has either 4  $C$ -tree nodes with probability  $1/6$ , 2 with probability  $1/2$ , or only one with probability  $1/3$ . These decompositions are shown in Figure 1.

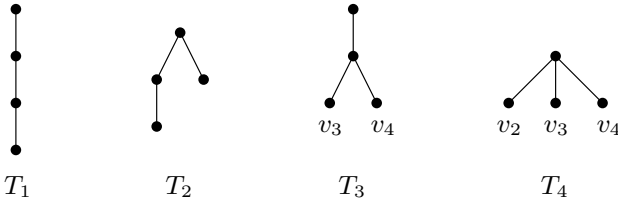


Figure 3: All Pólya trees of size 4.

In the same way as we got the composition scheme in (1.5), we can rewrite  $T_c(z, u)$  from (1.8) into  $T_c(z, u) = C(uzD(z))$ . The expected total weight of all  $C$ -trees contained in all Pólya trees of size  $n$  is the  $n$ -th coefficient of  $T_c(z)$ , which is

$$\begin{aligned} (3.13) \quad T_c(z) &:= \left. \frac{\partial}{\partial u} T_c(z, u) \right|_{u=1} \\ &= \frac{T(z)}{1 - T(z)} \\ &= z + 2z^2 + 5z^3 + 13z^4 + 35z^5 + \dots \end{aligned}$$

Let us explain why these numbers are integers, although the coefficients of  $t_{c,n}(u)$  are in general not. We will show an even stronger result. Let  $T$  be a tree and  $\mathcal{P}(T)$  be the set of all trees with one single pointed (or colored) node which can be generated from  $T$ .

LEMMA 3.1. For all  $T \in \mathcal{T}$  it holds that  $t'_T(1) = |\mathcal{P}(T)|$ .

*Proof.* From (3.12) we get that

$$t'_T(1) = \sum_{\sigma \in \text{Aut}(T)} \frac{\sigma_1}{|\text{Aut}(T)|}$$

is the expected number of fixed points in a uniformly at random chosen automorphism of  $T$ . The associated random variable  $C_T$  is defined in (1.9). We will prove  $\mathbb{E}(C_T) = |\mathcal{P}(T)|$  by induction on the size of  $T$ .

The most important observation is that only if the root of a subtree is a fixed point, its children can also be fixed points. Obviously, the root of the tree is always a fixed point.

For  $|T| = 1$ , the claim holds as  $\mathbb{E}(C_T) = 1$  and there is just one tree with a single node and a marker on it. For larger  $T$  consider the construction of Pólya trees. A Pólya tree consists of a root  $T_0$  and its children, which are a multiset of smaller trees. Thus, the set of children is of the form

$$\{T_{1,1}, \dots, T_{1,k_1}, T_{2,1}, \dots, T_{2,k_2}, \dots, T_{r,1}, \dots, T_{r,k_r}\},$$

with  $T_{i,j} \in \mathcal{T}$ , and where trees with the same first index are isomorphic. On the level of children, the possible behaviors of automorphisms are permutations within the same class of trees. In other words, an automorphism may interchange the trees  $T_{1,1}, \dots, T_{1,k_1}$  in  $k_1!$  many ways, etc. Here the main observation comes into play: only subtrees of which the root is a fixed point might also have other fixed points. Thus, the expected number of fixed points is given by the expected number of fixed points in a random permutation of  $S_{k_i}$  times the expected number of fixed points in  $T_{k_i}$ . By linearity of expectation we get

$$\mathbb{E}(C_T) = \mathbb{E}(C_{T_0}) + \sum_{i=0}^r \underbrace{\mathbb{E}(\# \text{ fixed points in } S_{k_i})}_{=1} \mathbb{E}(C_{T_i}),$$

where  $\mathbb{E}(C_{T_i}) = \mathbb{E}(C_{T_{i,j}})$  for all  $1 \leq j \leq k_i$  and  $\mathbb{E}(C_{T_0}) = 1$  because the root is a fixed point of any automorphism. Since the expected number of fixed points for each permutation is 1, we get on average 1 representative for each class of trees. This is exactly the operation of labeling one tree among each equivalence class. Finally, by induction the claim holds.

This completes the proof of Theorem 1.2.  $\square$

As an immediate consequence of Lemma 3.1,  $t'_{c,n}(1)$  counts the number of Pólya trees with  $n$  nodes and a single labeled node (see OEIS A000107, [13]). This also explains the construction of non-empty sequences of trees in (3.13): Following the connection [1, pp. 61–62]

one can draw a path from the root to each labeled node. The nodes on that path are the roots of a sequence of Pólya trees.

REMARK 3.1. *Note that Lemma 3.1 also implies that the total number of fixed points in all automorphisms of a tree is a multiple of the number of automorphisms.*

REMARK 3.2. *Lemma 3.1 can also be proved by considering cycle-pointed Pólya trees; see [2, Section 3.2] for a full description. Let  $(T, c)$  be a cycle-pointed structure considered up to symmetry where  $T$  is a Pólya tree and  $c$  is a cycle of an automorphism  $\sigma \in \text{Aut}(T)$ . Then, the number of such cycle-pointed structures  $(T, c)$  where  $c$  has length 1 is exactly the number  $t'_T(1)$ .*

Let us analyze the  $D$ -forests in  $T_n$  more carefully. We want to count the number of  $D$ -forests that have size  $m$  in a random Pólya tree  $T_n$ . Therefore, we label such  $D$ -forests with an additional parameter  $u$  in (1.5). From the bivariate generating function (1.10) we can recover the probability  $\mathbb{P}[|F_n(v)| = m]$  to generate a  $D$ -forest of size  $m$  in the Boltzmann sampler from [11].

**3.2 Proof of Theorem 1.3** The first result is a direct consequence of (1.5), where only vertices with weighted  $D$ -forests of size  $m$  are marked. For the second result we differentiate both sides of (1.10) and get

$$T_u^{[m]}(z, 1) = \frac{T(z)}{1 - T(z)} \frac{d_m z^m}{D(z)} = T_c(z) \frac{d_m z^m}{D(z)}.$$

Then, the sought probability is given by

$$\begin{aligned} \mathbb{P}[|F_n(v)| = m] &= \frac{[z^n] T_u^{[m]}(z, 1)}{[z^n] T_c(z)} \\ &= \frac{d_m \rho^m}{D(\rho)} (1 + \mathcal{O}(n^{-1})). \end{aligned}$$

For the last equality we used the fact that  $D(z)$  is analytic in a neighborhood of  $z = \rho$ .

Let  $P_n(u)$  be the probability generating function for the size of a weighted  $D$ -forest  $F_n(v)$  attached to a vertex  $v$  of  $C_n$  in a random Pólya tree  $T_n$ . From the previous theorem it follows that

$$\begin{aligned} P_n(u) &= \sum_{m \geq 0} \frac{[z^n] T_u^{[m]}(z, 1)}{[z^n] T_c(z)} u^m \\ &= \frac{[z^n] T_c(z) \frac{D(zu)}{D(z)}}{[z^n] T_c(z)} \\ &= \frac{D(\rho u)}{D(\rho)} (1 + \mathcal{O}(n^{-1})). \end{aligned}$$

This is exactly [11, Eq. (5.2)].  $\square$

Summarizing, we state the asymptotic probabilities that a weighted  $D$ -forest  $F_n(v)$  in  $T_n$  has size equal to or greater than  $m$ .

$m$	$\mathbb{P}[ F_n(v)  = m] \approx$	$\mathbb{P}[ F_n(v)  \geq m] \approx$
0	0.9197	1.0000
1	0.0000	0.0803
2	0.0526	0.0803
3	0.0119	0.0277
4	0.0105	0.0161
5	0.0015	0.0060
6	0.0027	0.0041
7	0.0003	0.0014

Table 1: The probability that a weighted  $D$ -forest  $F_n(v)$  has size equal to or greater than  $m$  when  $0 \leq m \leq 7$ .

## 4 Conclusion and perspectives

In this paper we provide an alternative proof of the maximal size of  $D$ -forests in a random Pólya tree. We interpret all weights on  $D$ -forests and  $C$ -trees in terms of automorphisms associated to a Pólya tree, and we derive the limiting probability that for a random node  $v$  the attached  $D$ -forest  $F_n(v)$  is of a given size.

Our work can be extended to  $\Omega$ -Pólya trees: For any  $\Omega \subseteq \mathbb{N}_0 = \{0, 1, \dots\}$  such that  $0 \in \Omega$  and  $\{0, 1\} \neq \Omega$ , an  $\Omega$ -Pólya tree is a rooted unlabeled tree considered up to symmetry and with outdegree set  $\Omega$ . When  $\Omega = \mathbb{N}_0$ , a  $\mathbb{N}_0$ -Pólya tree is a Pólya tree. In view of the connection between Boltzmann samplers and generating functions, it comes as no surprise that the “colored” Boltzmann sampler from [11] is closely related to a bivariate generating function. But the unified framework in analyzing the (bivariate) generating functions offers stronger results on the limiting distributions of the size of the  $C$ -trees and the maximal size of  $D$ -forests.

The next step is the study of shape characteristics of  $D$ -forests like the expected number of (distinct) trees. The  $C$ -tree is the simply generated tree within a Pólya tree and therefore its shape characteristics is well-known – when conditioned on its size. Moreover,  $D$ -forests certainly show a different behavior and, though they are fairly small, they still have significant influence on the tree. We will address these and other questions in the full version of this work.

**Acknowledgements:** This work was supported by the SFB project F50-03 “Combinatorics of Tree-Like Structures and Enriched Trees”. We also thank the three referees for their feedback.



## References

- [1] F. Bergeron, G. Labelle and P. Leroux. *Combinatorial Species and Tree-Like Structures*. Cambridge, 1998.
- [2] M. Bodirsky, É. Fusy, M. Kang and S. Vigerske. Boltzmann samplers, Pólya theory, and cycle pointing. *SIAM J. Comput.*, 40(3):721–769, 2011.
- [3] M. Drmota. *Random Trees. An Interplay Between Combinatorics and Probability*. Springer Verlag, 2008.
- [4] M. Drmota and B. Gittenberger. The shape of unlabeled rooted random trees. *European Journal of Combinatorics*, 31(8):2028–2063, 2010.
- [5] P. Flajolet and R. Sedgewick. *Analytic Combinatorics*. Cambridge University Press, 2009.
- [6] X. Gourdon. Largest components in random combinatorial structures. *Discrete Mathematics*, 180:185–209, 1998.
- [7] Svante Janson. Simply generated trees, conditioned Galton-Watson trees, random allocations and condensation. *Probab. Surv.*, 9:103–252, 2012.
- [8] A. Meir and J.W. Moon. On the altitude of nodes in random trees. *Canad. J. Math.*, 30(5):997–1015, 1978.
- [9] A. Nijenhuis and H.S. Wilf. *Combinatorial Algorithms*. Academic Press, Inc. [Harcourt Brace Jovanovich, Publishers], New York-London, second edition, 1978. For computers and calculators, Computer Science and Applied Mathematics.
- [10] R. Otter. The number of trees. *Ann. of Math.*, 49(2):583–599, 1948.
- [11] K. Panagiotou and B. Stuffer. Scaling limits of random Pólya trees. *Preprint: arXiv:1502.07180v2*, 2015.
- [12] G. Pólya. Kombinatorische Anzahlbestimmungen für Gruppen, Graphen und chemische Verbindungen. *Acta Mathematica* 68(1):145–254, 1937.
- [13] N.J.A. Sloane, *The On-line Encyclopedia of Integer Sequences (OEIS)*.
- [14] B. Stuffer. Random enriched trees with applications to random graphs. *Preprint: arXiv:1504.02006v6*, 2015.