

# Long and short paths in uniform random recursive dags

Luc Devroye and Svante Janson

**Abstract.** In a uniform random recursive  $k$ -directed acyclic graph, there is a root, 0, and each node in turn, from 1 to  $n$ , chooses  $k$  uniform random parents from among the nodes of smaller index. If  $S_n$  is the shortest path distance from node  $n$  to the root, then we determine the constant  $\sigma$  such that  $S_n/\log n \rightarrow \sigma$  in probability as  $n \rightarrow \infty$ . We also show that  $\max_{1 \leq i \leq n} S_i/\log n \rightarrow \sigma$  in probability.

## 1. Introduction

A uniform random  $k$ -dag is an infinite directed graph defined as follows. For each of the integers  $1, 2, \dots$ , we pick a random set of  $k$  parents with replacement uniformly from among the smaller non-negative integers. This defines an infinite directed acyclic graph (or, dag) with one root (0), and can be viewed as a (too) simplistic model of the web, a random recursive circuit (Díaz–Sperna–Spirakis–Toran–Tsukiji [12], and Tsukiji–Xhafa [40]), and a generalization of the URRT (uniform random recursive tree), which is obtained for  $k=1$ . (See also the further references in Section 4 to related models.) All the asymptotic results in the paper remain valid when parents are selected without replacement.

The uniform random  $k$ -dag restricted to vertices  $0, 1, \dots, n$ , is denoted by  $U_{k,n}$  or simply  $U_n$ . Indeed, we will take  $k=2$  in the main part of the paper, and point out the obvious modifications needed when  $k>2$  as we proceed. The infinite dag is denoted by  $U_\infty$ .

From a given node  $n$ , let  $\mathcal{P}_n$  be the collection of paths from node  $n$  to the origin. The length of a path  $p \in \mathcal{P}_n$  is  $L(p)$ . One can consider various path lengths:

$$S_n = \min_{p \in \mathcal{P}_n} L(p), \quad R_n^- = L(P_n^-), \quad R_n = L(P_n), \quad R_n^+ = L(P_n^+), \quad L_n = \max_{p \in \mathcal{P}_n} L(p),$$

---

L. Devroye's research was sponsored by NSERC Grant A3456. The research was mostly done at the Institute Mittag-Leffler during the programme *Discrete Probability* held in 2009.

where  $S$ ,  $R$  and  $L$  are mnemonics for shortest, random, and longest, and  $P_n^-, P_n$  and  $P_n^+$  are the paths in  $\mathcal{P}_n$ , where we follow the parent with the smallest index, the first parent and the parent with the largest index, respectively. We can regard  $R_n^-$  and  $R_n^+$  as greedy approximations of  $S_n$  and  $L_n$ , respectively. Note that, at least in a stochastic sense,

$$S_n \leq R_n^- \leq R_n \leq R_n^+ \leq L_n.$$

The length of the longest path is relevant for the time to compute the value of node  $n$  in a random recursive circuit, when nodes know their value only when all parents know their value. However, there are situations in which node values are determined as soon as one parent or a subset of parents know their value—they are called *self-time circuits* by Codenotti–Gemmell–Simon [6]. For the one-parent case, this leads naturally to the study of  $S_n$ . In networks, in general, shortest paths have been of interest almost since they were conceived (Prim [31], Dijkstra [13]). A recent study by physicists (D’Souza–Krapivsky–Moore [16]) predicts, without proof, our Theorem 2 for  $k=2$ .

It is of interest to study the extreme behavior, as measured by

$$\max_{1 \leq l \leq n} S_l, \quad \max_{1 \leq l \leq n} R_l^-, \quad \max_{1 \leq l \leq n} R_l, \quad \max_{1 \leq l \leq n} R_l^+, \quad \text{and} \quad \max_{1 \leq l \leq n} L_l.$$

If we replace max by min in these definitions, we obtain the constant 1, and it is therefore more meaningful to ask for the extreme minimal behavior as defined by

$$\min_{n/2 \leq l \leq n} S_l, \quad \min_{n/2 \leq l \leq n} R_l^-, \quad \min_{n/2 \leq l \leq n} R_l, \quad \min_{n/2 \leq l \leq n} R_l^+, \quad \text{and} \quad \min_{n/2 \leq l \leq n} L_l.$$

So, in all, there are fifteen parameters that could be studied.

We take this opportunity to introduce the label process, which will be referred to throughout the paper. The label of each parent of  $n$  is distributed as  $\lfloor nU \rfloor$ , with  $U$  uniform  $[0, 1]$ . An  $l$ th generation ancestor has a label distributed like

$$\lfloor \dots \lfloor \lfloor nU_1 \rfloor U_2 \rfloor \dots U_l \rfloor \in \lfloor nU_1 U_2 \dots U_l - l, nU_1 U_2 \dots U_l \rfloor,$$

where the  $U_i$ 's are independent and identically distributed (i.i.d.) uniform  $[0, 1]$  random variables.

*The parameter  $R_n$ .* It is clear that  $R_n$  is just the distance from node  $n$  in a URRT to its root. In particular,  $R_n$  and its minimal and maximal versions do not depend upon  $k$ . We dispense immediately with  $R_n$  and its extensions because of well-known results on the URRT obtained via the study of renewal processes (Devroye [11]) and the equivalence between  $R_n$  and the number of records in an i.i.d. sequence of continuous random variables (see, e.g., Rényi [33], Pyke [32], Glick [20]

or Devroye [9]). Only the minimal parameter for  $R_n$  requires a gentle intervention. We know that

$$\frac{R_n}{\log n} \rightarrow 1 \quad \text{in probability,}$$

for example. Furthermore,

$$\frac{R_n - \log n}{\sqrt{\log n}} \xrightarrow{\mathcal{L}} \mathcal{N},$$

where  $\mathcal{N}$  is a standard normal random variable, and  $\xrightarrow{\mathcal{L}}$  denotes convergence in distribution. Furthermore, an explicit tail bound on  $R_n$  will be needed further on in the paper. The maximal value of  $R_l$ ,  $1 \leq l \leq n$ , follows immediately from the theory of extremes of branching random walks (Devroye [8] and Pittel [30]). We summarize the following result.

**Theorem 1.** *We have*

$$\frac{R_n}{\log n} \rightarrow 1 \quad \text{in probability,}$$

$$\frac{\max_{1 \leq l \leq n} R_l}{\log n} \rightarrow e \quad \text{in probability,}$$

and

$$\lim_{n \rightarrow \infty} \mathbb{P} \left\{ \min_{n/2 \leq l \leq n} R_l \leq 2 \right\} = 1.$$

Finally, for  $t \geq \log n$  integer,

$$\mathbb{P}\{R_n > t\} \leq \exp \left( t - \log n - t \log \frac{t}{\log n} \right).$$

*Proof.* An outline of proof is needed for the third part and the explicit bound in part four. Let us count the number of nodes with index in  $[1, n/2]$  that connect directly to the root. This number is

$$Z = \sum_{l=1}^{n/2} \xi_{1/l},$$

where  $\xi_p$  is Bernoulli( $p$ ). Let  $A$  be the event that no node of index in  $(n/2, n]$  connects to a node counted in  $Z$ . This probability is smaller than

$$\begin{aligned} \mathbb{E} \left\{ \left( 1 - \frac{Z}{n} \right)^{n/2} \right\} &\leq \mathbb{E} \{ e^{-Z/2} \} = \prod_{l=1}^{n/2} \left( 1 - \frac{1}{l} + \frac{1}{\sqrt{el}} \right) \\ &\leq \exp \left( - \sum_{l=1}^{n/2} \frac{1 - 1/\sqrt{e}}{l} \right) \leq \left[ \frac{n}{2} \right]^{-(1-1/\sqrt{e})}. \end{aligned}$$

If the complement of  $A$  holds, then clearly,  $\min_{n/2 \leq l \leq n} R_l \leq 2$ , and thus, we have shown the third part of Theorem 1. Turning to part four, note that

$$R_n \leq \min\{t : nU_1 \dots U_t < 1\},$$

and thus that, for any  $\lambda > 0$ ,

$$\mathbb{P}\{R_n > t\} \leq \mathbb{P}\{nU_1 \dots U_t \geq 1\} \leq \mathbb{E}\{(nU_1 \dots U_t)^\lambda\} = n^\lambda(\lambda+1)^{-t}.$$

Hence,

$$\mathbb{P}\{R_n > t\} \leq \inf_{\lambda > 0} n^\lambda(\lambda+1)^{-t} = \exp\left(t - \log n - t \log \frac{t}{\log n}\right). \quad \square$$

*Conjecture 2.* For all fifteen parameters, generically denoted by  $X_n$ , there exist finite constants  $x = x(k) \geq 0$  such that

$$(1) \quad \frac{X_n}{\log n} \rightarrow x \quad \text{in probability.}$$

*Notation.* The limits in the conjecture are denoted by  $\sigma$ ,  $\rho^-$ ,  $\rho$ ,  $\rho^+$  and  $\lambda$  for  $S_n$ ,  $R_n^-$ ,  $R_n$ ,  $R_n^+$  and  $L_n$ , respectively. For the minimal and maximal versions of these parameters, we will use the subscripts min and max, respectively, as in  $\rho_{\min}^+$  and  $\sigma_{\max}$ , for example.

Let us briefly survey what is known and provide conjectures in the other cases.

*The parameter  $L_n$ .* Tsukiji and Xhafa [40] showed that  $\lambda_{\max} = ke$ .

Since there are (at most)  $k^t$  paths of length  $t$  in  $\mathcal{P}_n$ ,  $\mathbb{P}\{L_n > t\} \leq k^t \mathbb{P}\{R_n > t\}$ , and the Chernoff large deviation bound in Theorem 1 implies that  $\mathbb{P}\{L_n > x \log n\} \rightarrow 0$  if  $x > 1$  and  $k^x e^{x-1-x \log x} < 1$ ; hence  $\lambda$  is at most the largest solution  $x$  of

$$(2) \quad \left(\frac{ke}{x}\right)^x e^{-1} = 1,$$

and thus  $\lambda < \lambda_{\max}$ . We believe that  $\lambda$  is indeed given by (2) based on arguments not unlike the proof of Theorem 3 below. We have no guess at this point about the value of  $\lambda_{\min}$ .

The parameter  $R_n^+$ . In the label process, the parent's index is approximately distributed as  $n \max(U_1, \dots, U_k)$ , where the  $U_i$ 's are i.i.d. uniform  $[0, 1]$  random variables. If  $U$ , as elsewhere in this paper, is uniform  $[0, 1]$ , then the parent's index is thus roughly like  $nU^{1/k}$ . By renewal theory, this implies that

$$\frac{R_n^+}{\log n} \rightarrow k \stackrel{\text{def}}{=} \rho^+ \quad \text{in probability.}$$

(The standard argument is briefly that the index after  $l$  generations is roughly  $nW_1 \dots W_l$ , where  $W_1, W_2, \dots$  are i.i.d. and distributed as  $U^{1/k}$ . Hence,  $R_n^+$  is roughly the smallest  $l$  such that  $\sum_{i=1}^l (-\log W_i) = \log n$ , and thus

$$\frac{R_n^+}{\log n} \rightarrow \frac{1}{\mathbb{E}\{-\log W_1\}} = \frac{k}{\mathbb{E}\{-\log U\}} = k$$

in probability.)

Chernoff's large deviation bound show that  $\rho_{\max}^+$  is at most the unique solution  $x$  of (3) that is above  $k$ :

$$(3) \quad \left(\frac{ke}{x}\right)^x e^{1-k} = 1.$$

We believe that the solution of (3) yields  $\rho_{\max}^+$ . Applying Chernoff to the other tail shows that  $\rho_{\min}^+$  is at least the other solution of (3), as (3) has two solutions, one below  $k$  and one above  $k$ . Furthermore, we believe that this solution of (3) yields  $\rho_{\min}^+$ .

For  $k=2$ , the parameter  $R_n^+$  is intimately linked to the random binary search tree, which can be grown incrementally by a well-known process described as follows: given an  $n$ -node random binary search tree, sample one of its  $n+1$  external nodes uniformly at random, replace it by node  $n+1$ , and continue. The parent of that node is either its neighbor (in the total ordering) to the left or its neighbor to the right, and in fact, it is the neighbor added last to the tree. But the labels (times of insertion) of the neighbors are uniformly drawn without replacement from  $\{1, \dots, n\}$ , and are thus roughly distributed as  $nU$ , so that the parent of  $n+1$  is roughly distributed as  $n\sqrt{U}$ , because the maximum of two i.i.d. uniform  $[0, 1]$  random variables is distributed as  $\sqrt{U}$ . While this observation is valid for a single  $n$ , it is not true that we can choose parents independently, because in a random binary search tree, a node can be a parent at most twice. However, there is a kinship that allows one to conclude that the behavior should almost be the same. With this in mind,  $\max_{1 \leq l \leq n} R_l^+$  is roughly the height of the random binary search tree,  $R_n^+$  is roughly the depth (distance to the root) of the node of label  $n$  (the  $n$ th node inserted), and  $\min_{n/2 \leq l \leq n} R_l^+$  is very roughly the shortest distance from leaf to root, or fill-up level. These quantities behave in probability as described above, as shown by Devroye [7] and [8], and this explains the values  $\rho_{\max}^+ = 4.31107\dots$ ,  $\rho^+ = 2$  and  $\rho_{\min}^+ = 0.3733\dots$ .

*The parameter  $R_n^-$ .* Arguing as above, the parent's index is approximately distributed as  $n \min(U_1, \dots, U_k)$ . By a property of the uniform (or exponential) distribution, using a sequence of i.i.d. exponential random variables  $E_1, E_2, \dots$ , we have this distributional identity:

$$n \min(U_1, \dots, U_k) \stackrel{\mathcal{L}}{=} n U_1 U_2^{1/2} \dots U_k^{1/k} \stackrel{\mathcal{L}}{=} \exp\left(\log n - \sum_{j=1}^k \frac{E_j}{j}\right).$$

Renewal theory easily gives the law of large numbers and the central limit theorem for  $R_n^-$ . For example,

$$\frac{R_n^-}{\log n} \rightarrow \frac{1}{H_k} \stackrel{\text{def}}{=} \rho^- \quad \text{in probability,}$$

where  $H_k = \sum_{j=1}^k (1/j)$  is the  $k$ th harmonic number. Using large deviation bounds similar to the ones used below in showing part of Theorem 3, one gets that

$$\lim_{n \rightarrow \infty} \mathbb{P}\left\{\max_{1 \leq l \leq n} R_l^- \geq (x + \varepsilon) \log n\right\} = 0$$

for all  $\varepsilon > 0$ , where  $x$  is the solution greater than  $1/H_k$  of

$$1 + f(x) = x \sum_{j=1}^k \log\left(1 + \frac{f(x)}{j}\right),$$

and  $f(x) > 0$  is implicitly defined by

$$\sum_{j=1}^k \frac{1}{j + f(x)} = \frac{1}{x}, \quad x > \frac{1}{H_k}.$$

These equations follow from the obvious Chernoff bound. We conjecture that  $\rho_{\max}^-$  equals this upper bound, but a rigorous proof that  $\rho_{\max}^-$  is indeed as described above is not given in this paper.

*The parameter  $S_n$ .* The most important parameter for computer scientists and combinatorialists is the one in which graph distances are defined by shortest paths, and this leads to the study of  $S_n$ . That was the original motivation of the paper, and we will settle first order asymptotics in this paper. Theorem 1 implies, for example, that with probability tending to one,

$$\min_{n/2 \leq l \leq n} S_l \leq 2.$$

So we turn to  $\sigma$  and  $\sigma_{\max}$ .

Table 1. Proved and conjectured constants (1) for  $k=2$  (conjectured constants are marked with “?”).

	$x_{\min}$	$x$	$x_{\max}$
$\sigma$	0	0.3733...	0.3733...
$\rho^-$	0	0.6666... ( $=\frac{2}{3}$ )	1.6737...?
$\rho$	0	1	2.7182... ( $=e$ )
$\rho^+$	0.3733...?	2	4.3110...?
$\lambda$	?	4.3110...?	5.4365... ( $=2e$ )

**Theorem 3.** *Assume that  $k \geq 2$ . Then  $\sigma = \sigma_{\max}$ , where  $\sigma$  is given by the solution  $x \in (0, 1)$  of*

$$(4) \quad \varphi(x) \stackrel{\text{def}}{=} \left( \frac{ke}{x} \right)^x e^{-1} = 1.$$

[Note that  $\varphi$  is indeed an increasing function on  $(0, 1)$ .]

Observe that Theorem 3 does not extend to  $k=1$ , because in that case,  $S_n \equiv R_n \equiv L_n$ , and similarly for the maximal versions of these parameters, in view of the equivalence with the URRT. Thus,  $S_n / \log n \rightarrow 1$  and  $\max_{1 \leq l \leq n} S_l / \log n \rightarrow e$  in probability.

Table 1 is a table of constants in Conjecture 2 for  $k=2$ . The constants involving  $\sigma$  (top row) are obtained in this paper, while those involving  $\rho$  (third row) are covered by Theorem 1. The constants  $\rho^-$  and  $\rho^+$  follow from ordinary renewal theory. The zeroes in the table follow from Theorem 1. Finally,  $\lambda_{\max}$  is due to Tsukiji and Xhafa [40]. There are thus four conjectured constants, which happen to be one-sided bounds ( $\rho_{\max}^-, \rho_{\min}^+, \rho_{\max}^+, \lambda$ ), and one unknown constant,  $\lambda_{\min}$ .

Table 2 is a table of  $\sigma$ ,  $\rho^-$  and  $\rho_{\max}^-$  for different numbers of parents  $k$ .

## 2. The shortest path length $S_n$

We will establish Theorem 3 in two parts. First we show that for all  $\varepsilon > 0$ ,

$$(5) \quad \lim_{n \rightarrow \infty} \mathbb{P}\{S_n \leq (1 - \varepsilon)\sigma \log n\} = 0,$$

and then that

$$(6) \quad \lim_{n \rightarrow \infty} \mathbb{P}\left\{ \max_{1 \leq l \leq n} S_l \geq (1 + \varepsilon)\sigma \log n \right\} = 0.$$

We only consider the case  $k=2$  since the case  $k > 2$  follows quite easily.

Table 2. Constants for different  $k$ .

$k$	$\sigma$	$\rho^-$	$\rho_{\max}^-$
2	0.3733...	0.6666...	1.6737...
3	0.3040...	0.5454...	1.3025...
4	0.2708...	0.48	1.1060...
5	0.2503...	0.4379...	0.9818...
6	0.2361...	0.4081...	0.8951...
7	0.2254...	0.3856...	0.8305...
8	0.2170...	0.3679...	0.7800...
9	0.2102...	0.3534...	0.7393...
10	0.2045...	0.3414...	0.7057...
11	0.1996...	0.3311...	0.6773...
12	0.1954...	0.3222...	0.6531...
13	0.1916...	0.3144...	0.6318...
14	0.1883...	0.3075...	0.6132...
15	0.1854...	0.3013...	0.5966...
16	0.1827...	0.2957...	0.5816...
17	0.1802...	0.2907...	0.5683...
18	0.1780...	0.2861...	0.5560...
19	0.1760...	0.2818...	0.5448...
20	0.1740...	0.2779...	0.5346...
21	0.1723...	0.2743...	0.5251...
22	0.1706...	0.2709...	0.5164...
23	0.1691...	0.2677...	0.5083...
24	0.1676...	0.2648...	0.5007...
25	0.1663...	0.2620...	0.4936...
26	0.1650...	0.2594...	0.4868...
27	0.1638...	0.2569...	0.4805...
28	0.1626...	0.2546...	0.4747...
29	0.1615...	0.2524...	0.4690...
30	0.1604...	0.2503...	0.4638...
35	0.1559...	0.2411...	0.4409...
40	0.1521...	0.2337...	0.4225...
45	0.1490...	0.2275...	0.4074...
50	0.1463...	0.2222...	0.3946...

**Lemma 4.** Let  $G_a$  be gamma( $a$ ), with  $a \geq 1$ . Then

$$\frac{\mathbb{P}\{G_a \geq x\}}{\frac{x^{a-1}e^{-x}}{\Gamma(a)}} \leq \frac{1}{1 - \frac{a-1}{x}}, \quad x > a-1,$$

and

$$\frac{\mathbb{P}\{G_a \leq x\}}{\frac{x^{a-1}e^{-x}}{\Gamma(a)}} \leq \frac{1}{\frac{a-1}{x} - 1}, \quad x < a-1.$$

*Proof.* The gamma density is  $f(y) = y^{a-1}e^{-y}/\Gamma(a)$ . It is log-concave for  $a \geq 1$ , and thus, a first-term Taylor series bound yields the inequality

$$f(y) \leq f(x)e^{(y-x)(\log f)'(x)} = f(x)e^{(y-x)((a-1)/x-1)}.$$



Integrating the upper bound over  $[x, \infty)$  or  $(-\infty, x]$  then immediately yields the results.  $\square$

From node  $n$ , we can consider the index of the first of the  $2^l$   $l$ th level ancestors, which is distributed as

$$[\dots[ [nU_1]U_2] \dots U_l] \geq nU_1U_2 \dots U_l - l \stackrel{\mathcal{L}}{=} n \exp(-G_l) - l,$$

where  $\stackrel{\mathcal{L}}{=}$  denotes equality in distribution, and  $G_l$  is  $\text{gamma}(l)$ . If these indices are  $I_1, \dots, I_{2^l}$ , then we have

$$\begin{aligned} \mathbb{P}\{S_n \leq l\} &= \mathbb{P}\left\{ \min_{1 \leq i \leq 2^l} I_i = 0 \right\} \\ &\leq 2^l \mathbb{P}\{I_1 = 0\} \\ &\leq 2^l \mathbb{P}\{n \exp(-G_l) - l \leq 0\} \\ &= 2^l \mathbb{P}\left\{ G_l \geq \log \frac{n}{l} \right\} \\ &\leq \frac{2^l (\log(n/l))^{l-1} e^{-\log(n/l)}}{\Gamma(l) \left(1 - \frac{l-1}{\log(n/l)}\right)} \quad \left( \text{if } \log \frac{n}{l} \geq l-1 \right) \\ &\leq \frac{l^{3/2} (2 \log n)^l e^{-\log n}}{\left(\frac{l}{e}\right)^l \left(1 - \frac{l-1}{\log(n/l)}\right)}. \end{aligned}$$

Set  $l = \lfloor t \log n \rfloor$  for  $t \in (0, 1)$ , and note that the upper bound is

$$\Theta(\log^{3/2} n) \times \varphi(t)^{\log n},$$

where  $\varphi(t) = (2e/t)^t / e$  is as in (4). We have  $\varphi(\sigma) = 1$  for  $\sigma = 0.3733\dots$ , and thus  $\varphi(t) < 1$  for  $0 < t < \sigma$ . Thus, we have shown (5): for all  $\varepsilon > 0$ ,

$$\mathbb{P}\{S_n \leq (\sigma - \varepsilon) \log n\} = o(1).$$

Although we will not need it directly, we will also deal with the upper bound on  $S_n$ . This can be done in a number of ways, but the shortest route is perhaps via the great-grandparent strategy that jumps  $l$  generations at a time, where  $l$  now is a large but fixed integer. We denote this by  $l$ -GGP. We associate with each node  $n$  two independent uniform  $[0, 1]$  integers  $U$  and  $V$  and let the parent labels be  $\lfloor nU \rfloor$  and  $\lfloor nV \rfloor$ . Let  $A_n$  be the event that any of the  $2^l$  ancestors of node  $n$  coincide. It is clear that  $\mathbb{P}\{A_n\} \rightarrow 0$  as  $n \rightarrow \infty$ . As an ancestor label is described by

$$[\dots[ [nU_1]U_2] \dots U_l] \leq nU_1U_2 \dots U_l \stackrel{\mathcal{L}}{=} n \exp(-G_l),$$

we define

$$Z_l = \min_{p \in \mathcal{P}} \prod_{e \in p} U_e,$$

where  $\mathcal{P}$  is the collection of all paths of length  $l$  above node  $n$ , and each  $p \in \mathcal{P}$  consists of edges  $e$  that each have an independent uniform random variable associated with it. If  $\varepsilon > 0$  and  $n$  is greater than some  $n_\varepsilon$ , then the  $l$ -GGP gives a node with label less than  $Z_l n$  with probability greater than  $1 - \varepsilon$  [failure would imply that  $A_n$  holds]. Define

$$Z_l^{(\varepsilon)} = \begin{cases} Z_l, & Z_l > b, \\ 1, & Z_l \leq b, \end{cases}$$

where  $b$  is chosen such that  $\mathbb{P}\{Z_l \leq b\} = \varepsilon$ . As long as the label stays above  $n_\varepsilon$ , one can dominate the labels in the  $l$ -GGP by multiplying  $n$  with successive independent copies of  $Z_l^{(\varepsilon)}$ . Let  $T_n$  be the number of steps until the label in  $l$ -GGP reaches  $n_\varepsilon$  or less. Renewal theory shows that with probability tending to one,

$$T_n \leq \frac{(1+\varepsilon) \log n}{\mathbb{E}\{-\log Z_l^{(\varepsilon)}\}}.$$

Because the  $l$ -GGP takes  $l$  steps at a time, and because a node with label  $n_\varepsilon$  is not further than  $n_\varepsilon$  away from the origin, we see that with probability tending to one,

$$S_n \leq n_\varepsilon + \frac{l(1+\varepsilon) \log n}{\mathbb{E}\{-\log Z_l^{(\varepsilon)}\}} \leq \frac{l(1+2\varepsilon) \log n}{\mathbb{E}\{-\log Z_l^{(\varepsilon)}\}}.$$

Uniform integrability implies that

$$\lim_{\varepsilon \downarrow 0} \mathbb{E}\{-\log Z_l^{(\varepsilon)}\} = \mathbb{E}\{-\log Z_l\}.$$

Therefore, for any (new, fresh)  $\varepsilon > 0$  and  $l \geq 1$ , with probability going to one,

$$S_n \leq \frac{l(1+\varepsilon) \log n}{\mathbb{E}\{-\log Z_l\}}.$$

Observe that

$$\frac{-\log Z_l}{l} \leq \frac{1}{l} \max_{p \in \mathcal{P}} \sum_{e \in p} E_u,$$

where the  $E_u$  are i.i.d. exponential random variables. From the theory of branching random walks, it is easy to verify (see, e.g., Biggins [3], or Devroye [7] and [8]) that, as  $l \rightarrow \infty$ ,

$$\frac{1}{l} \max_{p \in \mathcal{P}} \sum_{e \in p} E_u \rightarrow \frac{1}{\sigma}$$

in probability. Hence,

$$\liminf_{l \rightarrow \infty} \frac{-\mathbb{E}\{\log Z_l\}}{l} \geq \frac{1}{\sigma},$$

and thus, by choosing  $l$  large enough, we see that with probability tending to one,

$$S_n \leq (1+2\varepsilon)\sigma \log n.$$

This concludes the proof of the first part of Theorem 3.

The next section requires an explicit rate of convergence. To this end, still restricting ourselves to  $k=2$  only, let  $Z_{l,1}^{(\varepsilon)}, Z_{l,2}^{(\varepsilon)}, \dots$  be i.i.d. copies of  $Z_l^{(\varepsilon)}$ , and note that,

$$T_n \leq \min\left\{t : nZ_{l,1}^{(\varepsilon)} \dots Z_{l,t}^{(\varepsilon)} < 1\right\} = \min\left\{t : \log \frac{1}{Z_{l,1}^{(\varepsilon)}} + \dots + \log \frac{1}{Z_{l,t}^{(\varepsilon)}} > \log n\right\}.$$

Set  $\mu = \mathbb{E}\{\log(1/Z_l^{(\varepsilon)})\}$ . Then, assuming that  $\delta^* \in (0, \frac{1}{2})$  and that  $\delta \in (\delta^*, 2\delta^*)$  is such that  $m = (1/\mu + \delta) \log n$  is integer-valued,

$$\begin{aligned} \mathbb{P}\{T_n > m\} &\leq \mathbb{P}\left\{\log \frac{1}{Z_{l,1}^{(\varepsilon)}} + \dots + \log \frac{1}{Z_{l,m}^{(\varepsilon)}} < \log n\right\} \\ &= \mathbb{P}\left\{\log \frac{1}{Z_{l,1}^{(\varepsilon)}} + \dots + \log \frac{1}{Z_{l,m}^{(\varepsilon)}} - m\mu < -\delta\mu \log n\right\}. \end{aligned}$$

Let  $p > 2$  be a fixed number. Rosenthal's inequality (Rosenthal [34], Fuk-Nagaev [18], see also Petrov [29]) states that there is a constant  $C_p$  with the following property. If  $\{X_n\}_{n \geq 1}$  is a sequence of centered and independent random variables, and if  $Y_n = X_1 + \dots + X_n$ , and if  $\mathbb{E}\{|X_n|^p\} < \infty$  for all  $n$ , then

$$\mathbb{E}\{|Y_n|^p\} \leq C_p \left( \sum_{j=1}^n \mathbb{E}\{|X_j|^p\} + (\mathbb{V}\{Y_n\})^{p/2} \right).$$

For i.i.d. random variables with  $X_1 = X$ , we have

$$\mathbb{E}\{|Y_n|^p\} \leq C_p (n\mathbb{E}\{|X|^p\} + n^{p/2}(\mathbb{E}\{X^2\})^{p/2}) \leq 2C_p \max(n, n^{p/2})\mathbb{E}\{|X|^p\}.$$

Applied to our situation with  $p=4$ , using Markov's inequality, we have

$$\begin{aligned} \mathbb{P}\{T_n > m\} &\leq (\delta\mu \log n)^{-4} \mathbb{E}\left\{\left(\log \frac{1}{Z_{l,1}^{(\varepsilon)}} + \dots + \log \frac{1}{Z_{l,m}^{(\varepsilon)}} - m\right)^4\right\} \\ &\leq 2C_4 (\delta\mu \log n)^{-4} m^2 \mathbb{E}\left\{\left|\log \frac{1}{Z_l^{(\varepsilon)}} - \mu\right|^4\right\} \leq C (\log n)^{-2} \delta^{*-4}, \end{aligned}$$

where  $C$  depends upon  $\varepsilon$  and  $l$  only. The remainder of the argument involving an appropriate choice of  $l$  remains valid, and we can conclude that for any  $\varepsilon > 0$ ,

$$(7) \quad \mathbb{P}\{S_n > (\sigma + \varepsilon) \log n\} = O\left(\frac{1}{\log^2 n}\right),$$

with room to spare.

### 3. The maximal shortest path length

The purpose of this section is to show (6). We let  $\sigma$  be as in the first part of the proof, and let  $\varepsilon > 0$  be arbitrary. Fix  $n$  large enough. From (7),

$$\mathbb{E}\left\{\left|\left\{j: \frac{n}{2} \leq j \leq n \text{ and } S_j > (\sigma + \varepsilon) \log n\right\}\right|\right\} = O\left(\frac{n}{\log^2 n}\right),$$

and thus  $\mathbb{P}\{A(n)\} = O(1/\log^2 n)$ , where

$$A(n) \stackrel{\text{def}}{=} \left[\left|\left\{j: \frac{n}{2} \leq j \leq n \text{ and } S_j > (\sigma + \varepsilon) \log n\right\}\right| > \frac{n}{4}\right].$$

If we take an incremental view of the process of adding edges, then a node with index in  $[n, 2n]$  selects a parent of depth  $\leq (\sigma + \varepsilon) \log n$  and index  $\geq n/2$  with probability  $\geq \frac{1}{8}$  if  $A(n)$  fails to hold. It is this observation that will allow us to uniformly bound all depths by something close to  $(\sigma + \varepsilon) \log n$ .

Consider the indices in dyadic groups,  $\{2^{r-1} + 1, \dots, 2^r\}$ ,  $r \geq 1$ . We recall from a comparison with the URRT, that  $S_n \leq R_n$  and thus that  $\max_{1 \leq j \leq n} S_j \leq \max_{1 \leq j \leq n} R_j$ , and that (see Theorem 1)

$$\mathbb{P}\left\{\max_{1 \leq j \leq n} R_j > 2e \log n\right\} \leq n^{-2e \log 2} < n^{-3}.$$

Thus, for  $\gamma > 0$  small enough,

$$\mathbb{P}\left\{\max_{1 \leq j \leq \lfloor n^\gamma \rfloor} S_j > (\sigma + \varepsilon) \log n\right\} = O(n^{-3\gamma}) = o(1).$$

It remains to show that

$$\mathbb{P}\left\{\max_{n^\gamma \leq j \leq n} S_j > (\sigma + \varepsilon) \log n\right\} = o(1).$$

Consider the event

$$B(r) = \bigcup_{r' \leq s \leq r} A(2^s),$$

where  $r'$  is the largest integer such that  $2^{r'} < n^\gamma$ . Clearly,

$$\mathbb{P}\{B(r)\} = O\left(\frac{1}{r'}\right) = O\left(\frac{1}{\log n}\right).$$

On the complement,  $B(r)^c$ , intersected with  $[\max_{1 \leq j \leq \lfloor n^\gamma \rfloor} S_j \leq (\sigma + \varepsilon) \log n]$ , we look at the process started at a node  $m \leq n$  and assume that its index  $m$  is in  $\{2^r + 1, \dots, 2^{r+1}\}$ . That process is looked at as a binary tree of consecutive parents, and will be cut off at height  $h = \lfloor 10 \log \log n \rfloor$ . There may be duplicate parents (in which case the tree degenerates to a dag), so we need to be a bit careful. If any parent in the tree is selected with index  $\leq 2^{r'} < n^\gamma$ , then  $S_m \leq (\sigma + \varepsilon) \log n + h$ , and thus, we can assume that in this “tree” any node  $j$  selects its parent uniformly in the range  $(2^{r'}, j)$ . At any stage, by our assumption, the probability of picking a parent  $i$  having  $S_i \leq (\sigma + \varepsilon) \log n$  is at least  $\frac{1}{8}$  (and this is why we needed the dyadic trick, so that we can make this statement regardless of the choice of  $i$  within the range  $(2^{r'}, n]$ ). We claim that this “tree” has at least  $2^{h-1}$  leaves or reaches  $[1, 2^{r'}]$  with overwhelming probability. To see this, note that a node  $j$  in it picks a node already selected with probability not exceeding  $2^h/j$ . But the index  $j$  is stochastically larger than

$$X_h \stackrel{\text{def}}{=} [\dots [mU_1]U_2] \dots U_h]$$

by our remarks about the labeling process. The probability that there are in fact at least two such unwanted parent selections (but none of them less than  $n^\gamma$ ) in that “tree” is not more than

$$(8) \quad 2^{2h+2} \times \mathbb{E}^2 \left\{ \frac{2^h}{X_h} \mathbb{1}_{[X_h \geq n^\gamma]} \right\} \leq 2^{4h+2} \times \mathbb{E}^2 \left\{ \frac{1}{X_h} \mathbb{1}_{[X_h \geq n^\gamma]} \right\}.$$

We have

$$\begin{aligned} \mathbb{E} \left\{ \frac{1}{X_h} \mathbb{1}_{[X_h \geq n^\gamma]} \right\} &= \int_0^\infty \mathbb{P} \left\{ \frac{1}{X_h} \mathbb{1}_{[X_h \geq n^\gamma]} > t \right\} dt \\ &= \int_0^{1/n^\gamma} \mathbb{P} \left\{ X_h < \frac{1}{t} \right\} dt \\ &\leq \int_0^{1/n^\gamma} \mathbb{P} \left\{ mU_1 \dots U_h < h + \frac{1}{t} \right\} dt \\ &= \int_0^{1/n^\gamma} \mathbb{P} \left\{ \log \frac{m}{h+1/t} < G_h \right\} dt \\ &= \int_0^{1/n^\gamma} \int_{\log_+(m/(h+1/t))}^\infty \frac{y^{h-1} e^{-y}}{\Gamma(h)} dy dt \end{aligned}$$

$$\begin{aligned}
&= \int_0^\infty \frac{y^{h-1} e^{-y}}{\Gamma(h)} \min\left(n^{-\gamma}, \frac{1}{(me^{-y} - h)_+}\right) dy \\
&\leq \int_0^{\log(m/2h)} \frac{2y^{h-1}}{\Gamma(h)m} dy + n^{-\gamma} \int_{\log(m/2h)}^\infty \frac{y^{h-1} e^{-y}}{\Gamma(h)} dy \\
&\leq \frac{2(\log n)^h}{m h!} + \frac{n^{-\gamma} (\log n)^{h-1} 4h}{\Gamma(h)m} \\
&\quad (\text{for } n \text{ large enough, by Lemma 4}) \\
&= O\left(\frac{n^{o(1)}}{m}\right) \\
&= O(m^{-1+o(1)}).
\end{aligned}$$

Thus, our probability (8) is not larger than  $O(m^{-2+o(1)})$ . If there is only one unwanted parent selection and we avoid indices below  $n^\gamma$ , and considering that the first parent selection at the root node is always good, we see that at least half of the  $2^h$  potential leaves are in fact realized. Each of these leaves makes two independent parent selections. The probability that all these leaves avoid parents  $j$  with  $S_j < (\sigma + \varepsilon) \log n$  is at most  $(\frac{7}{8})^{2^{h-1}} = o(n^{-2})$ . If there is a connection, however, to such a parent of low depth, then the root has shortest path length at most  $h+1$  more than  $(\sigma + \varepsilon) \log n$ . Hence, if  $\mathcal{E}_m$  is the event  $[S_m > (\sigma + \varepsilon) \log n + h + 1]$ , then

$$\mathbb{P}\left\{\mathcal{E}_m \cap B(r)^c \cap \left[\max_{1 \leq j \leq \lfloor n^\gamma \rfloor} S_j \leq (\sigma + \varepsilon) \log n\right]\right\} = O(m^{-2+o(1)}).$$

Thus

$$\begin{aligned}
\mathbb{P}\left\{\max_{n^\gamma \leq j \leq n} S_j > (\sigma + \varepsilon) \log n + h + 1\right\} &= \mathbb{P}\left\{\bigcup_{m \geq n^\gamma}^n \mathcal{E}_m\right\} \\
&\leq \mathbb{P}\left\{\max_{1 \leq j \leq \lfloor n^\gamma \rfloor} S_j > (\sigma + \varepsilon) \log n\right\} \\
&\quad + \mathbb{P}\{B(r)\} + \sum_{m \geq n^\gamma}^n m^{-2+o(1)} \\
&= O(n^{-3\gamma}) + O\left(\frac{1}{r'}\right) + n^{-\gamma+o(1)} \\
&= O\left(\frac{1}{\log n}\right).
\end{aligned}$$

This concludes the proof of Theorem 3.

#### 4. Bibliographic remarks and possible extensions

The study of the URRT goes back as far as Na–Rapoport [28] and Meir–Moon [27]. Single nonuniform parent selections have been considered as early as 1987 by Szymański. Szymański [37] showed that if a parent is selected with probability proportional to its degree, then with high probability there is a node of degree  $\Omega(\sqrt{n})$ . This is nothing but the preferential attachment model of Barabasi and Albert (see Albert–Barabasi–Jeong [2], or Albert–Barabasi [1]), which for a single parent is a special case of the linear recursive trees or plane-oriented recursive tree. For this model, the parameter  $R_n$  was studied by Mahmoud [24], and the height by Pittel [30] and Biggins–Grey [4], and in a rather general setting by Broutin–Devroye [5]: the height is in probability  $(1.7956\dots + o(1)) \log n$ . The profile (number of nodes at each depth level) was studied by Hwang [21], [22] and Sulzbach [36].

One can ask the questions studied in the present paper for these more general models.

Various aspects of URRT's besides the depth and height have been studied by many researchers. These include the degrees of the nodes, the profile, sizes of certain subtrees of certain nodes, the number of leaves, and so forth. Surveys and references can be found in the book by Mahmoud [25] or the paper by Devroye [10]. Specific early papers include Timofeev [39], Gastwirth [19], Dondajewski–Szymański [14], Mahmoud [23], Mahmoud–Smythe [26], Smythe–Mahmoud [35], Szymański [38], and the most recent contributions include Fuchs–Hwang–Neininger [17], and Drmota–Janson–Neininger [15]. One may wonder how the profiles behave for uniform random  $k$ -dags.

#### References

1. ALBERT, R. and BARABASI, A., Emergence of scaling in random networks, *Science* **286** (1999), 509–512.
2. ALBERT, R., BARABASI, A. and JEONG, H., Diameter of the World-Wide Web, *Nature* **401** (1999), 130–131.
3. BIGGINS, J. D., Chernoff's theorem in the branching random walk, *J. Appl. Probab.* **14** (1977), 630–636.
4. BIGGINS, J. D. and GREY, D. R., A note on the growth of random trees, *Statist. Probab. Lett.* **32** (1997), 339–342.
5. BROUTIN, N. and DEVROYE, L., Large deviations for the weighted height of an extended class of trees, *Algorithmica* **46** (2006), 271–297.
6. CODENOTTI, B., GEMMELL, P. and SIMON, J., Average circuit depth and average communication complexity, in *Third European Symposium on Algorithms*, pp. 102–112, Springer, Berlin, 1995.
7. DEVROYE, L., A note on the height of binary search trees, *J. Assoc. Comput. Mach.* **33** (1986), 489–498.

8. DEVROYE, L., Branching processes in the analysis of the heights of trees, *Acta Inform.* **24** (1987), 277–298.
9. DEVROYE, L., Applications of the theory of records in the study of random trees, *Acta Inform.* **26** (1988), 123–130.
10. DEVROYE, L., Branching processes and their applications in the analysis of tree structures and tree algorithms, in *Probabilistic Methods for Algorithmic Discrete Mathematics*, Algorithms Combin. **16**, pp. 49–314, Springer, Berlin, 1998.
11. DEVROYE, L., Universal limit laws for depths in random trees, *SIAM J. Comput.* **28** (1999), 409–432.
12. DÍAZ, J., SERNA, M. J., SPIRAKIS, P., TORAN, J. and TSUKIJI, T., On the expected depth of Boolean circuits, *Technical Report LSI-94-7-R*, Universitat Politècnica de Catalunya, Departament de Llenguatges i Sistemes Informàtics, Barcelona, 1994.
13. DIJKSTRA, E. W., A note on two problems in connexion with graphs, *Numer. Math.* **1** (1959), 269–271.
14. DONDAJEWSKI, M. and SZYMAŃSKI, J., On the distribution of vertex-degrees in a strata of a random recursive tree, *Bull. Acad. Polon. Sci. Sér. Sci. Math.* **30** (1982), 205–209.
15. DRMOTA, M., JANSON, S. and NEININGER, R., A functional limit theorem for the profile of search trees, *Ann. Appl. Probab.* **18** (2008), 288–333.
16. D’SOUZA, R. M., KRAPIVSKY, P. L. and MOORE, C., The power of choice in growing trees, *Eur. Phys. J. B* **59** (2007), 535–543.
17. FUCHS, M., HWANG, H.-K. and NEININGER, R., Profiles of random trees: Limit theorems for random recursive trees and binary search trees, *Algorithmica* **46** (2006), 367–407.
18. FUK, D. K. and NAGAEV, S. V., Probability inequalities for sums of independent random variables, *Teor. Veroyatnost. i Primenen.* **16** (1971), 660–675 (Russian). English transl.: *Theory Probab. Appl.* **16** (1971), 643–660.
19. GASTWIRTH, J. L., A probability model of a pyramid scheme, *Amer. Statist.* **31** (1977), 79–82.
20. GLICK, N., Breaking records and breaking boards, *Amer. Math. Monthly* **85** (1978), 2–26.
21. HWANG, H.-K., Profiles of random trees: plane-oriented recursive trees (Extended Abstract), in *2005 International Conference on Analysis of Algorithms*, pp. 193–200, Assoc. Discrete Math. Theor. Comput. Sci., Nancy, 2005.
22. HWANG, H.-K., Profiles of random trees: plane-oriented recursive trees, *Random Structures Algorithms* **30** (2007), 380–413.
23. MAHMOUD, H. M., Limiting distributions for path lengths in recursive trees, *Probab. Engrg. Inform. Sci.* **5** (1991), 53–59.
24. MAHMOUD, H. M., Distances in random plane-oriented recursive trees, *J. Comput. Appl. Math.* **41** (1992), 237–245.
25. MAHMOUD, H. M., *Evolution of Random Search Trees*, Wiley, New York, 1992.
26. MAHMOUD, H. M. and SMYTHE, R. T., On the distribution of leaves in rooted subtrees of recursive trees, *Ann. Appl. Probab.* **1** (1991), 406–418.
27. MEIR, A. and MOON, J. W., On the altitude of nodes in random trees, *Canad. J. Math.* **30** (1978), 997–1015.



28. NA, H. S. and RAPOPORT, A., Distribution of nodes of a tree by degree, *Math. Biosci.* **6** (1970), 313–329.
29. PETROV, V. V., *Limit Theorems of Probability Theory: Sequences of Independent Random Variables*, Oxford University Press, Oxford, 1995.
30. PITTEL, B., Note on the heights of random recursive trees and random  $m$ -ary search trees, *Random Structures Algorithms* **5** (1994), 337–347.
31. PRIM, R. C., Shortest connection networks and some generalizations, *Bell System Tech. J.* **36** (1957), 1389–1401.
32. PYKE, R., Spacings, *J. Roy. Statist. Soc. Ser. B* **27** (1965), 395–445.
33. RÉNYI, A., Théorie des éléments saillants d’une suite d’observations, in *Colloquium on Combinatorial Methods in Probability Theory*, pp. 104–115, Matematisk Institut, Aarhus Universitet, Aarhus, 1962.
34. ROSENTHAL, H. P., On the subspaces of  $L^p$  ( $p > 2$ ) spanned by sequences of independent random variables, *Israel J. Math.* **8** (1970), 273–303.
35. SMYTHE, R. T. and MAHMOUD, H. M., A survey of recursive trees, *Teor. ĭmovĭr. Mat. Stat.* **51** (1994), 1–29 (Ukrainian). English transl.: *Theory Probab. Math. Statist.* **51** (1995), 1–27 (1996).
36. SULZBACH, H., A functional limit law for the profile of plane-oriented recursive trees, in *Fifth Colloquium on Mathematics and Computer Science*, pp. 339–350, Assoc. Discrete Math. Theor. Comput. Sci., Nancy, 2008.
37. SZYMAŃSKI, J., On a nonuniform random recursive tree, in *Random Graphs ’85 (Poznań, 1985)*, North-Holland Math. Stud. **144**, pp. 297–306, North-Holland, Amsterdam, 1987.
38. SZYMAŃSKI, J., On the maximum degree and height of a random recursive tree, in *Random Graphs ’87 (Poznań, 1987)*, pp. 313–324, Wiley, Chichester, 1990.
39. TIMOFEEV, E. A., Random minimal trees, *Teor. Veroyatnost. i Primenen.* **29** (1984), 134–141 (Russian). English transl.: *Theory Probab. Appl.* **29** (1985), 134–141.
40. TSUKIJI, T. and XHAFI, F., On the depth of randomly generated circuits, in *Algorithm—ESA ’96*, Lecture Notes in Comput. Sci. **1136**, pp. 208–220, Springer, Berlin, 1996.

Luc Devroye  
 School of Computer Science  
 McGill University  
 3450 University Street  
 Montreal, QC H3A 2K6  
 Canada  
[luc@cs.mcgill.ca](mailto:luc@cs.mcgill.ca)

Svante Janson  
 Department of Mathematics  
 Uppsala University  
 P.O. Box 480  
 SE-751 06 Uppsala  
 Sweden  
[Svante.Janson@math.uu.se](mailto:Svante.Janson@math.uu.se)

*Received June 1, 2009*  
*published online February 25, 2010*