

# Weak subcritical percolation on finite mean-field graphs.

Olivier Hénard \*

September 17, 2021

## Abstract

We consider the largest connected components in the percolation of a (large) finite vertex-transitive graph. A geometrical condition reminiscent of the one in Nachmias [Nac09] is formulated. Under this condition, the many-to-two formula allows one to compute the size of the largest components in the weak subcritical regime up to fluctuations. The result applies to Hamming graphs up to dimension 3 and expander.

## 1 Introduction

Let  $G_n$  be a vertex-transitive graph on  $n$  vertices, and call  $G_n(p)$  the random subgraph of  $G_n$  obtained by deleting every edge of  $G_n$  independently and with the same probability  $1 - p$  for  $p \in (0, 1)$ . The number of vertices in the connected components of  $G_n(p)$  is one of the simplest functionals of the graph  $G_n$ . Following a long tradition, we study the number of vertices, or sizes, of the largest of these connected components. In the setting of an arbitrary vertex-transitive graph, the problem has been started in the series of papers [BCvdH<sup>+</sup>05a, BCvdH<sup>+</sup>05b, BCvdH<sup>+</sup>06] and continued in [Nac09, vdHN15, HN20]. Specific graphs  $G_n$  had been considered much before: the case of  $K_n(p)$  with  $K_n$  the complete graph on  $n$  vertices is the classical Erdős-Rényi [ER60] random graph model that has been examined in detail by Erdős and Rényi in the '60s. A convenient parametrisation for the probability  $p$  of retaining an edge is in this case

$$p = \frac{\lambda}{n} \quad \text{for } \lambda > 0 \text{ a constant independent of } n \tag{1}$$

For large values of  $n$ , the size of the largest component of  $K_n(p)$  expects a double jump when  $\lambda$  is increased: it is  $\Theta_{\mathbb{P}}(\log(n))$ <sup>1</sup> when  $\lambda < 1$ ,  $\Theta_{\mathbb{P}}(n^{2/3})$  when  $\lambda = 1$ , and  $\Theta_{\mathbb{P}}(n)$  when  $\lambda > 1$ . As for the size of the second largest component of  $K_n(p)$ : it is  $\Theta_{\mathbb{P}}(\log(n))$  when

---

\*LMO, Université Paris-Saclay, 91405 Orsay Cedex, France. Partial support from grant ANR-14-CE25-0014 (ANR GRAAL) and EPSRC grant EP / J004022/2

<sup>1</sup>see Janson [Jan11] for a definition of the probabilistic symbols  $\Theta_{\mathbb{P}}$  and  $o_{\mathbb{P}}(n)$  used in the Introduction only.

$\lambda < 1$ ,  $\Theta_{\mathbb{P}}(n^{2/3})$  when  $\lambda = 1$ , and  $\Theta_{\mathbb{P}}(\log(n))$  again when  $\lambda > 1$ . The three regimes are respectively called the subcritical, the critical and the supercritical regime. A distinctive feature of the critical regime is the non-concentration of the size of the largest component, that weakly converges as  $n \rightarrow \infty$  towards a non-degenerate random variable. Aldous [Ald97] gives a construction of that random variable from the sample path of a Brownian motion with a quadratic drift.

Always in the  $K_n(p)$  case, Bollobás [Bol84] and Łuczak [Łuc90] discovered new regimes of interest around the critical value  $\lambda = 1$ : these regimes are parametrized by sequences  $\lambda = \lambda(n)$  with limit 1. Many qualitative features of the critical case  $\lambda = 1$  are retained when  $\lambda = \lambda(n)$  approaches 1 fast enough:  $|\lambda - 1| = O(n^{-1/3})$ , and this regime is called the critical regime. An essential feature of that regime is the non-concentration of the sizes  $\Theta_{\mathbb{P}}(n^{2/3})$  of the largest components. The two remaining regimes are parametrised by positive sequences  $\varepsilon = \varepsilon(n)$  satisfying <sup>2</sup>

$$\lambda = 1 \pm \varepsilon \quad \text{with} \quad \varepsilon \rightarrow 0 \quad \text{and} \quad \varepsilon^3 n \rightarrow +\infty.$$

Choosing the *minus* sign defines the weak subcritical regime in which the largest component has size  $(2 + o_{\mathbb{P}}(1)) \varepsilon^{-2} \log(\varepsilon^3 n)$ , a quantity that interpolates between  $\log(n)$  and  $n^{2/3}$  for  $\varepsilon$  in the range described above. Choosing the *plus* sign defines the weak supercritical regime, in which the size of the largest component is  $(2 + o_{\mathbb{P}}(1)) \varepsilon n$ , a quantity that interpolates between  $n^{2/3}$  and  $n$ . In the weak supercritical regime, the second largest component is negligible with respect to the largest component, therefore called the giant component.

In this paper, we consider a sequence  $(G_n)_{n \in \mathbb{N}}$  of vertex-transitive graphs on  $n$  vertices in place of the complete graph  $K_n$ . By vertex-transitivity, the degree  $\ell$  of a vertex in  $G_n$  is the same for all the vertices, and we shall (mainly) consider sequences  $\ell = \ell(n)$  diverging with  $n$ . We choose the percolation probability to be

$$p = \frac{\lambda}{\ell - 1} \tag{2}$$

and, again, we consider  $\lambda = \lambda(n)$ . The choice (2) is the natural generalisation of (1): the expected total number of edges in  $G_n(p)$  is for instance in both cases equal to  $\lambda n/2$ . In the same setting, Nachmias [Nac09] investigates the critical and in the weak supercritical regime. Under condition (63) mentioned at the end of this article, he proves that the largest components in the critical regime have  $\Theta_{\mathbb{P}}(n^{2/3})$  many vertices; in the weak supercritical regime, under a slightly stronger condition, he gives a lower bound on the size of the largest component,  $\delta \varepsilon n / \log(\varepsilon^3 n)$  for some fixed constant  $\delta > 0$  independent of  $n$ ; this latter result has been since superseded by the more general work [vdHN15], that catches the correct asymptotic value  $2\varepsilon n(1 + o_{\mathbb{P}}(1))$ . If the choice (2) is natural, it is not always adequate: for many interesting graphs, at the percolation probability  $1/(\ell - 1)$ ,  $|\mathcal{C}_1| = o_{\mathbb{P}}(n^{2/3})$  so this value does not lie in the so-called critical window. A line of research investigates the analogy between the percolation of those graphs with the percolation of  $K_n(p)$  by first defining an appropriate notion of critical probability, see the series [BCvdH<sup>+</sup>05a, BCvdH<sup>+</sup>05b, BCvdH<sup>+</sup>06] and the recent works [vdHN15, HN20] for significant successes.

---

<sup>2</sup> $\varepsilon$  shall always denote a positive quantity in this paper

## 2 Statement of the results and discussion.

Let  $G_n$  be a sequence of vertex transitive graphs on  $n$  vertices. The degree of a vertex is  $\ell = \ell(n) \geq 3$ . We will mainly work under the assumption

$$\ell(n) \rightarrow +\infty \tag{3}$$

although the case where  $\ell(n)$  has a finite limit is also discussed at the end of the article. The edges of  $G_n$  are retained independently with probability  $p$ , and we denote by  $G_n(p)$  the resulting random subgraph of  $G_n$ . We let

$$p_{\pm} = \frac{1 \pm \varepsilon}{\ell - 1}, \text{ for } \varepsilon = \varepsilon(n) \rightarrow 0, \text{ and } \varepsilon^3 n \rightarrow \infty \tag{4}$$

a positive sequence, and we call weak sub-critical the regime in which  $p = p_-$  and weak super-critical regime the regime in which  $p = p_+$

A function that is useful in describing the sizes of the largest components of  $G_n(p)$  is

$$\delta_{\pm}(\varepsilon, \ell - 1) = -\log(1 \pm \varepsilon) - (\ell - 2) \log\left(1 - \frac{\pm \varepsilon}{\ell - 2}\right) \tag{5}$$

It is connected to the tail of the size of a  $\text{Bin}(\ell - 1, p_{\pm})$ -Galton-Watson tree (GW-tree hereafter), see (29) and (30). Taking the  $\ell \rightarrow \infty$  limit in (5) we find that  $\lim_{\ell \rightarrow \infty} \delta_{\pm}(\varepsilon, \ell - 1) = -\log(1 \pm \varepsilon) \pm \varepsilon$ , a function that is ubiquitous in the study of the Erdős-Rényi random graph. It is useful to keep in mind the equivalent

$$\delta_{\pm}(\varepsilon, \ell - 1) \sim \frac{\varepsilon^2}{2} \text{ when } \varepsilon \rightarrow 0 \text{ and } \ell \rightarrow \infty \tag{6}$$

The probability that the non-backtracking random walk on  $G_n$  (a random walk on the vertices of  $G_n$  not allowed to traverse the same edge on two consecutive steps) started at a vertex  $v \in G_n$  returns at  $v$  after  $k \geq 2$  steps is denoted by  $P^k(v, v)$ . By vertex-transitivity, this quantity does not depend on the vertex  $v$ , and we write  $P^k$  for  $P^k(v, v)$ . We also introduce the two parameters

$$t_{\pm} = \delta_{\pm}^{-1} s \text{ and } s = \log(\varepsilon^3 n)$$

Notice that, when  $n$  and  $\ell$  diverge,  $t_{\pm} \sim 2\varepsilon^{-2} \log(\varepsilon^3 n)$ . Our key (asymptotic) condition on the (sequence of) graphs  $G_n$  then writes:

$$\exists \text{ a constant } c < 1/2 \text{ such that: } (t_{\pm})^{1/2} \sum_{k \geq 3} k e^{-ck^2/t_{\pm}} P^k = o\left(\frac{1}{s}\right) \text{ as } n \rightarrow \infty \tag{7}$$

Condition (7) contains two conditions, the one with  $t_+$  and the one with  $t_-$ , we distinguish them by writing  $(7)_+$  and  $(7)_-$  when needed. Condition (7) is a condition on the geometry of the graphs  $G_n$ , where by "geometry" we simply mean the collection of numbers  $(P^k, k \in \mathbb{N})$ . We stress on the implicit dependence of the parameters  $t_{\pm}, P^k$  and  $s$  on  $\varepsilon$  and/or  $n$ . The

quantity  $ke^{-ck^2/t}$  is an upper bound on the expected number of vertices at distance  $k$  from the root in a (critical) GW-tree with finite variance and  $t$  vertices, and the bound is uniform in  $k \in \{1, \dots, t\}$ . See Addario-Berry et al [ABDJ13] for related estimates for GW-trees whose offspring distribution is critical and has a finite variance. The quantity  $t^{1/2}$  that multiplies the sum in (7) is the expected distance of a random vertex to the root in such a GW-tree. We will clarify in Section why  $o(1/s)$  in the RHS of (7) is indeed the precision needed to state the following Theorem :

**Theorem 2.1.** *Assume  $\varepsilon$  satisfies (4),  $\ell$  satisfies (3) and let  $p_- = (1 - \varepsilon)/(\ell - 1)$ . Assume also that the non-backtracking random walk on  $G_n$  satisfies condition (7)<sub>-</sub>. Let  $(\mathcal{P}_j, j \in \mathbb{N})$  be the points, arranged in non-increasing order, of a Poisson point measure with intensity*

$$(4\sqrt{\pi})^{-1}e^{-x}dx \quad (8)$$

on the real line  $\mathbb{R}$ . The sizes  $(|\mathcal{C}_j|, j \in \mathbb{N})$  of the largest components of  $G_n(p_-)$  arranged in non-increasing order converge in distribution for the product topology, in the sense that for each fixed  $k \in \mathbb{N}$ , as  $n \rightarrow \infty$ , we have:

$$\left( \delta_-(\varepsilon, \ell - 1) |\mathcal{C}_j| - \left( \log(\varepsilon^3 n) - \frac{5}{2} \log \log(\varepsilon^3 n) \right), 1 \leq j \leq k \right) \Longrightarrow (\mathcal{P}_j, 1 \leq j \leq k). \quad (9)$$

*Remark 2.2* (On the right-hand side in (9)). The computation:

$$\mathbb{P}(\mathcal{P}_1 < y) = e^{-\int_{[y, \infty)} (4\sqrt{\pi})^{-1} e^{-x} dx} = e^{-(4\sqrt{\pi})^{-1} e^{-y}}, \quad y \in \mathbb{R}$$

ensures the rightmost point of the Poisson measure in (8) exists and is Gumbel distributed. The convergence in distribution of a vector entails the convergence in distribution of its first coordinate. Also, the Gumbel distribution has no atom. Therefore (9) implies, for  $y \in \mathbb{R}$  :

$$\mathbb{P} \left( \delta_-(\varepsilon, \ell - 1) |\mathcal{C}_1| - \left( \log(\varepsilon^3 n) - \frac{5}{2} \log \log(\varepsilon^3 n) \right) > y \right) \rightarrow 1 - e^{-(4\sqrt{\pi})^{-1} e^{-y}}$$

as  $n \rightarrow \infty$ . In the same way, it is possible to write down explicitly the limiting distribution of the size of the  $k$ -th largest component for  $k \in \mathbb{N}$  fixed. Another byproduct is that the (rescaled) spacing between the two largest random variables weakly converges towards an exponential random variable with parameter 1 :  $\mathbb{P}(\delta_-(\varepsilon, \ell - 1) (|\mathcal{C}_1| - |\mathcal{C}_2|) > x) \rightarrow e^{-x}$  as  $n \rightarrow \infty$ .

*Remark 2.3.* (On the left-hand side in (9)) From (6) and (9), we deduce the first order asymptotics:  $|\mathcal{C}_j| = (1 + o_{\mathbb{P}}(1)) \delta_-^{-1}(\varepsilon, \ell - 1) \log(\varepsilon^3 n) = (1 + o_{\mathbb{P}}(1)) 2\varepsilon^{-2} \log(\varepsilon^3 n)$ . Beware one cannot in general replace  $\delta_-(\varepsilon, \ell - 1)$  by  $\varepsilon^2/2$  in (9). Replacing  $\delta_- = \delta_-(\varepsilon, \ell - 1)$  by another quantity  $\tilde{\delta}_-$  on the left-hand side of (9) is possible as long as  $(\delta_- - \tilde{\delta}_-)/\tilde{\delta}_- = o(1/s)$  as  $n \rightarrow \infty$ . From the expansion:

$$\delta_- = \sum_{k \geq 2} \frac{\varepsilon^k}{k} \left( 1 + \frac{(-1)^k}{(\ell - 2)^{k-1}} \right) \quad (10)$$

we see that  $\delta_- = (1 + O(\varepsilon + 1/\ell))\varepsilon^2/2$  as  $\ell \rightarrow \infty$  and  $\varepsilon \rightarrow 0$ . The choice  $\tilde{\delta}_- = \varepsilon^2/2$  is thus possible under the conditions  $(\varepsilon \log(\varepsilon^3 n) = o(1))$  and  $\log(\varepsilon^3 n) = o(\ell)$ .

*Remark 2.4.* (Weakly dependent random variables) For  $v \in G_n$ , call  $\mathcal{C}(v)$  the component that contains the vertex  $v$  in the random subgraph  $G_n(p)$ . The  $n$  coordinates of the vector  $(|\mathcal{C}(v)|, v \in G_n)$  are identically distributed, with the tail of their common distribution given by (31), but not quite independent. Were these random variables also independent, then the statement of Theorem 2.1 would hold with only  $\varepsilon^3 n$  replaced by  $\varepsilon n$  on the LHS of (9). The slight difference reflects the (weak) dependence between the random variables  $(|\mathcal{C}(v)|, v \in G_n)$  that is caused by repetitions : if  $|\mathcal{C}(v)| = t$  for some vertex  $v$  in  $G_n$ , then  $|\mathcal{C}(u)| = t$  for at least  $t - 1$  other vertices  $u$ , the ones that belong to that component.

*Remark 2.5.* (Further illustration of the weak dependence) The difference is also visible when computing the expected number of components with size larger than  $t$ . It is in general bounded by (and, for the aforementioned choice of  $t$ , in fact equivalent to)  $(n/t)\mathbb{P}(|\mathcal{C}(v)| > t)$  and *not* to  $n\mathbb{P}(|\mathcal{C}(v)| > t) = \Theta(1)$  as it would be in the independent case: the division by  $t$  precisely accounts for the aforementioned repetitions. The proof will show it is possible to deduce the correct order of magnitude for  $t$  from this first moment argument : setting  $t = \delta^{-1}s$  in the equation  $(n/t)\mathbb{P}(|\mathcal{C}(v)| > t) = 1$ , the LHS is by (31) equivalent to  $(\varepsilon^3 n) \cdot s^{-5/2} \cdot e^{-s}$  up to multiplicative constant, and this is 1 when  $s = \log(\varepsilon^3 n) - 5/2 \log \log(\varepsilon^3 n) + O(1)$ . Notice the  $\log(\varepsilon^3 n)$  term comes from the exponential decay, and the  $\log \log(\varepsilon^3 n)$  from the polynomial correction.

With the help of the estimates on the kernel of the non-backtracking random walk computed in [Nac09], one can check condition (7)<sub>-</sub> holds for some new classes of graphs. A sequence of connected graphs  $G_n$  is called an expander family if the largest eigenvalue in absolute value of the transition matrix of the simple random walk on  $G_n$ , distinct from  $\pm 1$ , is strictly smaller than 1, uniformly in  $n$ . The girth of a graph is the length of a shortest cycle in the graph. The  $d$ -dimensional Hamming graph is the cartesian products of  $d$  complete graphs, it has vertex set  $V = \{1, \dots, n\}^d$  and two vertices are linked by an edge iff the associated  $d$ -tuples differ at a single coordinate. Also  $[x] \in \mathbb{Z}$  denotes the integer part of  $x \in \mathbb{R}$ .

**Proposition 2.6.** *Assume the sequence  $\varepsilon$  satisfies (4). Condition (7) holds for the following graphs  $G_n$  :*

- *the transitive expander graphs with girth  $g = g(n)$  and vertex degree  $\ell = \ell(n)$  that satisfy*

$$\left( \frac{1}{\ell(n) - 1} \right)^{\lfloor g(n)/2 \rfloor} n^{1/3} \log^2(n) = O(1) \tag{11}$$

- *the Hamming graph in dimension  $d = 3$ .*

*Remark 2.7* (Hamming graphs). The Hamming graphs are expanders that in dimension 1 and 2 satisfy condition (11) in the first statement. The Hamming graph in dimension 1 is the complete graph, hence the result on the Erdős-Rényi random graph by Łuczak [Luc90] is recovered, whereas the result in dimension 2 and 3 is new. See Section 2.1 below for recent results on these graphs that are valid in *any* dimension.

Theorem 2.1 (together with the verification of condition (7)) are the highlights of this paper. Slight improvements are conceivable: for instance the assumption of transitivity of the graph could be slightly relaxed to include regular graphs. Similarly, in our theorem,  $P^k$  could possibly be replaced by the (slightly smaller) probability that the random walk draws a self-avoiding loop<sup>3</sup> of length  $k$  : in view of the form of condition (7), one expects little benefit on the applicability of Theorem 2.1 however. Our results follows the line of inquiry set up by Nachmias [Nac09]. Because the precise study of the weak subcritical regime is not relevant for studying the width of the critical window, this work paid little attention to this regime, see the bottom of p.1173 in [Nac09]. That left aside a key feature of the weak subcritical regime, namely that it allows for a sharp estimate of the size of the largest components under general assumptions. That omission was subsequently repaired by the important work [HN20] that tackles the more general setting where the critical probability is *implicitly* defined only. If the largest components in the weak subcritical regime are notably smaller, hence "easier" to deal with, than their counterparts in the critical and supercritical regimes, finding a sharp estimate with the right multiplicative constant may represent a technical challenge, see [HN20] again. Our setting, where the critical probability is explicit, is easier, and allows one to obtain the optimal precision (fluctuations) by carefully controlling the difference between a (conditioned) GW-tree and the component containing a given vertex in  $G_n(p)$ . This crucially requires to work with GW-trees conditioned by the size and not by the height like in [Nac09, HN20]: the former conditioning better approximates the geometry of the largest components in  $G_n(p)$  in our regime. A key technical tool to succeed is the "many-to-two" formula, that we rederive from scratch, see formula (19). If the use of a many-to- $k$  formula seems new in the context of random graphs, the tool, that goes back to back to Ikeda et al [INW69] in the case  $k = 2$ , has already proved very useful in the study of branching Brownian motion. A generic version, the many-to- $k$  (or many-to-few) formula is discussed in Harris and Roberts [HR15]; when  $k = 1$ , it reduces to the standard many-to-one formula, see [ABDJ13] or chapter 12 of the book [LP17].

## 2.1 Previous work

Let us try to summarize the state of affairs concerning the study of the size of the largest components in the weak subcritical regime:

- The complete graph model  $K_n(p)$  : The fluctuations of the random variable  $|\mathcal{C}_j|$ ,  $j \in \mathbb{N}$ , have been identified by Łuczak [Łuc90] in a sharpening of a result by Bollobás [Bol84]. We warn the reader of a small typo in the statement of the result in [Łuc90]:  $\varepsilon^2/2$  should be replaced by  $\delta$  for the result to hold along the critical window, see Remark 1.3 or [BR09] p.50.
- The configuration model: The fluctuations of  $|\mathcal{C}_1|$  are (among other) given in Riordan [Rio12], they involve the same Gumbel distribution as in our result.

---

<sup>3</sup>a path  $v_0, v_1, \dots, v_k = v_0$  where  $0 \leq i < j < k \Rightarrow v_i \neq v_j$

- The Achlioptas bounded-size rules: Riordan and Warnke [RW17] offer a complete study of the transition phase in the context of iteratively constructed random graphs, in which edges come by pair (say) and some amount of choice in the edge to be added is allowed. Again, the authors are able to compute the fluctuations of  $|\mathcal{C}_1|$ , see their Theorem 2.7.
- The  $G_n(p)$  model: Deterministic graphs  $G_n$  distinct from the complete graph require very different methods; upper and lower bounds on  $|\mathcal{C}_j|, j \in \mathbb{N}$ , have been known for some time now under a finite size version of the triangle condition known as the "strong" triangle condition, see Theorem 1.2 in Borgs, Chayes, van der Hofstad and Spencer [BCvdH<sup>+</sup>05a] but these bounds were separated by a multiplicative  $\log(\varepsilon^3 n)$  factor. Only very recently has this gap been reduced to a multiplicative constant factor, see [HN20] :  $|\mathcal{C}_1| = \Theta_{\mathbb{P}}(\varepsilon^{-2} \log(\varepsilon^3 n))$ . A key difficulty in these works is that the critical probability is only defined implicitly; we should also mention in this direction some recent progress concerning the asymptotic expansion of the critical probability in the special case of the Hamming graph in (fixed) dimension  $d \geq 1$ , see [FVDHDDH20].

An important convention concerning the  $\pm$  index: when the index is omitted like in  $\delta, t$ , one should understand that the statement holds for both values  $\delta_{\pm}, t_{\pm}$ . An identity with such a quantity therefore contains two identities: the one with the  $\bullet_+$  index and the one with the  $\bullet_-$  index.

## 2.2 Ideas of proof

The key parameters involved in the study are, in term of  $\delta = \delta(\varepsilon, \ell - 1), \varepsilon$  and  $n$  :

$$t = \delta^{-1}s, s = \log(\varepsilon^3 n) - \frac{5}{2} \log \log(\varepsilon^3 n) + u \quad (12)$$

We do not show the dependence of these parameters on  $\varepsilon$  and  $n$ ;  $u$  will be either a fixed constant, or a very slow function of  $n$ . We denote by  $t : \mathbb{R} \rightarrow \mathbb{R}$  the function defined by  $t = t(u)$  and by  $t^{-1}$  its inverse function. Write  $\mathcal{C}_j$  for the  $j$ -th largest component of  $G_n(p)$  as measured by its total number of vertices  $|\mathcal{C}_j|$  (ties are broken in an arbitrary way). Key to the proof of Theorem 2.1 is the understanding of the following random measure:

$$N(dx) = \sum_{j \geq 0} \delta_{t^{-1}(|\mathcal{C}_j|)}(dx) \quad (13)$$

and of its convergence in distribution in particular. Define  $J = (u_1, u_2)$  with  $-\infty < u_1 \leq u_2 < +\infty$ , and let  $t(J) = (t_1, t_2)$ . We first want to check the convergence of the random variable  $N_-(J)$  (the minus index refers to the subcritical regime) as  $n$  tends to  $\infty$ . Let  $\mathcal{C}(v)$  denote the component of  $G_n(p)$  that contains the vertex  $v$ . Using the estimate  $t_1 \sim t_2$  as  $n \rightarrow \infty$ , we find that the first moment of  $N(J)$  satisfies:

$$\mathbb{E}(N(J)) = \mathbb{E} \left( \sum_{v \in G_n} \frac{\mathbf{1}_{\{|\mathcal{C}(v)| \in t(J)\}}}{|\mathcal{C}(v)|} \right) \sim n \frac{\mathbb{P}(|\mathcal{C}(v)| \in (t_1, t_2))}{t_1} \text{ as } n \rightarrow \infty \quad (14)$$

Another easy fact is that the tail  $\mathbb{P}(|\mathcal{C}(v)| > t)$  is bounded from above by the quantity  $\mathbb{P}(|T^m| > t)$ , where  $T^m$  is a "modified" GW-tree with  $\text{Bin}(\ell - \mathbf{1}_{\{v=\rho\}}, p)$  offspring distribution. The complementary lower bound is the difficult step. Different ideas can be developed to estimate it, let us review three different possibilities.

First, writing  $T^m = T_\ell^m$  to show the dependence of  $T^m$  on  $\ell$ , a step-by-step exploration of  $\mathcal{C}(v)$  reveals that:

$$\mathbb{P}(|T_{\ell-t}^m| \geq t) \leq \mathbb{P}(|\mathcal{C}(v)| \geq t) \leq \mathbb{P}(|T_\ell^m| \geq t)$$

In case  $t$  is small with respect to  $\ell$ , like in the complete graph  $K_n$ , ( $\ell = n - 1$ ), the two bounds are equivalent sequences as  $n \rightarrow \infty$ . The computation of the successive moments of  $|\mathcal{C}(v)|$  does not raise additional difficulty, and one finds sharp asymptotics on  $|\mathcal{C}_1|$ . The same strategy works for "mean-field" model like the configuration model, see [Rio12] around formula (7.5). There are other classes of graphs for which a uniform control of the number of already explored vertices is possible, for instance the Hamming graph in dimension 2 [vdHL10]. Second, Nachmias [Nac09] introduced the idea of pruning off the upper bound tree  $T^m$  from the so-called path-impure vertices to find a sub-tree of  $T^m$  that is stochastically smaller than  $|\mathcal{C}(v)|$ . The path-impure vertices are the vertices that are present in the GW-tree but have no counterpart in the exploration of  $\mathcal{C}(v)$ . The set of path-impure vertices is denoted by  $I^1(T^m)$ . A simple bound is, for  $1 \leq t \leq t'$  :

$$\mathbb{P}(|T^m| \geq t', |I^1(T^m)| \leq t' - t) \leq \mathbb{P}(|\mathcal{C}(v)| \geq t) \leq \mathbb{P}(|T^m| \geq t) \quad (15)$$

Based on this bound, plus the second moment method, Nachmias derives in his Lemma 14 a lower bound on the probability that  $|\mathcal{C}(v)|$  exceeds  $\varepsilon^{-2}$ , that is not sharp in the weak subcritical regime of interest to us (essentially because the method of proof relies on conditioning a GW-tree by the height). The third method is the one developed in this paper, it is based on conditioning the GW-trees by the size. Consider  $J' = (u, \infty)$ . From (14), the estimate (27) on the total progeny of a GW-tree gives that  $\mathbb{E}(N_-(J')) \leq e^{-u}/(4\sqrt{\pi})$ . One asks under what condition  $(1 - o(1))e^{-u}/(4\sqrt{\pi})$  is a lower bound. Examining (27) again, we see that if  $t'$  close to  $t = \delta^{-1}s$  in the sense that  $t' - t = o(\delta^{-1})$  then the ratio  $\mathbb{P}(T^m \geq t')/\mathbb{P}(T^m \geq t)$  has limit 1 as  $n$  diverges. But the difference between  $\mathcal{C}(v)$  and  $T^m$  is controlled by the set  $I^1(T^m)$  of path-impure vertices, so, by the first moment method, we only need to show  $\mathbb{E}(|I^1(T^m)| | |T^m| = t) = o(\delta^{-1})$ . The latter is proved in Proposition 3.8 under condition (7), and the proof of that Proposition is in turn based on a many-to-two formula, Lemma 3.1 with  $k = 2$ , that is indeed our key tool. Under (7) the lower bound matches the upper bound, and:

$$\mathbb{E}(N_-(J)) = (1 + o(1)) \frac{1}{4\sqrt{\pi}} (e^{-u_1} - e^{-u_2}). \quad (16)$$

The rest of the proof is routine: the RHS is also the expectation of a Poisson random variable with parameter the integral over  $J = (u_1, u_2)$  of the intensity measure (8). To claim the convergence in distribution through the method of moments, see *e.g.* Section 6.1 of [JLR11] it remains to check the convergence of the higher factorial moments  $\mathbb{E}(\prod_{0 \leq i < k} (N_-(J) - i))$ , which require a last technical point: (16) has to be proved with  $G \setminus G^0$  the graph induced by

$G$  on  $V(G) \setminus V(G^0)$ , for  $G^0 \subseteq G$  a subgraph with size  $O(t)$ . The next and last section contains the (self-contained) proof; it starts in Section 3.1 with the many-to- $k$  formula. The required estimates on the size of the GW-trees of interest are stated in Section 3.2. In Section 3.3 GW-trees are randomly embedded in the graph  $G_n$  and the number of path-impure vertices in a GW-tree with a given size is estimated. The short subsection 3.4 is a remark on how to translate our estimates in term of modified GW-trees. We start to work with the  $G_n(p)$  model itself only at Section 3.5 where we compute the moment of the number of components in a given interval, and also give the few additional ingredients needed to prove the main theorem. The extension to bounded degree graphs is discussed in Section 3.6. In Section 3.7 we obtain, as a by-product of our analysis in the weak subcritical regime, a lower bound on the expected number of components in certain intervals in the weak *supercritical* regime: yet we fail to get the lower bound on the size of the second largest component that is suggested by this estimate. In the last Section 3.8 we prove the Proposition 2.6 on condition (7).

### 3 Proof

The set of integers is denoted by  $\mathbb{Z}$ , the subset of non-negative integers  $\{0, 1, 2, \dots\}$  by  $\mathbb{Z}^+$ , the subset of positive integers  $\{1, 2, \dots\}$  by  $\mathbb{N}$  and the set of real numbers by  $\mathbb{R}$ . Unless explicitly advertised, all the limits and asymptotic are as  $n$  the size of the graph  $G_n$  goes to  $\infty$ . Also, we use the Landau notation  $o$  and  $O$  (but no more probabilistic counterpart from now on). Sums over an empty set are 0, and products over an empty set are 1. If  $G = (V, E)$  is a graph,  $|G|$  will denote its number of vertices, or size, of  $G$ . For  $V^0 \subseteq V$ , the graph induced by  $G$  on  $V^0$  is the graph with vertex set  $V^0$  and edge set the restriction of the original edge set  $E$  to  $V^0 \times V^0$ . Assume another graph  $G' = (V', E')$  is given. A map  $f : V \rightarrow V'$  is called a graph homomorphism when adjacent vertices in  $G$  are mapped to adjacent vertices in  $G'$ ,  $(u, v) \in E \Rightarrow (f(u), f(v)) \in E'$  for each  $u, v \in V$ . It is called a graph isomorphism when  $f$  is a bijection from  $V$  to  $V'$  and two vertices are adjacent in  $G$  iff their images are adjacent in  $G'$ , that is  $(u, v) \in E$  iff  $(f(u), f(v)) \in E'$ .

We may need to attach several distinguished vertices to a graph, that we shall call the pointed vertices: for  $k \in \mathbb{N}$ , and  $u_1, u_2, \dots, u_k \in V$ ,  $G = (V, E, u_1, u_2, \dots, u_k)$  is a pointed graph, and if  $G' = (V', E', u'_1, u'_2, \dots, u'_k)$  for  $u'_1, u'_2, \dots, u'_k \in V'$  is another pointed graph, we say a graph-homomorphism  $f$  from  $(V, E)$  to  $(V', E')$  preserves the pointed vertices when  $f(u_i) = u'_i$  for  $1 \leq i \leq k$ . The trees we will encounter will be planar and rooted. Such trees are embedded in the so-called Ulam tree: this is the graph with vertex set the finite sequences of integers

$$\mathcal{U} = \bigcup_{n \geq 0} \mathbb{N}^n.$$

The root of the Ulam tree is the vertex  $\mathbb{N}^0$ , that we shall denote by  $\rho$ . A vertex distinct from the root,  $u = (u(1), u(2), \dots, u(k)) \in \mathcal{U}$ ,  $k \geq 1$ , has a unique father  $a(u) := (u(1), u(2), \dots, u(k-1))$ , and there is one edge between every vertex distinct from the root and its father. Notice that the Ulam tree is a locally-infinite. The integer  $k$  is the generation of  $u$ , denoted by  $|u|$ . By convention,  $|\rho| = 0$ . For  $i \leq k$ , define  $u^{(i)} = (u(1), u(2), \dots, u(i))$  the

ancestor of  $u$  at generation  $i$ . When  $u$  is an ancestor of  $v$  we write  $u \preceq v$ , and  $u \prec v$  when also  $u \neq v$ ; in the latter case,  $u$  is called a strict ancestor of  $v$ .  $\preceq$  defines a partial order on  $\mathcal{U}$  called the ancestral order. A subset  $\mathbf{t}$  of  $\mathcal{U}$  is a planar rooted tree when: (i)  $\mathbf{t}$  contains  $\rho$ . (ii)  $v = (v(1), \dots, v(k-1), v(k)) \in \mathbf{t}$  implies  $(v(1), \dots, v(k-1), i) \in \mathbf{t}$  for any  $i \in \{1, \dots, v(k)\}$  (iii)  $v = (v(1), \dots, v(k)) \in \mathbf{t}$  implies  $v^{(j)} = (v(1), \dots, v(j)) \in \mathbf{t}$ , for any  $j \in \{1, \dots, k\}$ . A vertex  $v \in \mathbf{t}$  is a leaf when its set of children  $c_{\mathbf{t}}(v) := \{w \in \mathbf{t}, a(w) = v\}$  is empty. We write  $c(v)$  for  $c_{\mathbf{t}}(v)$  when  $\mathbf{t}$  is clear from the context. The number  $|c(v)|$  of children of  $v \in \mathbf{t}$  is called the outdegree of  $v$  in  $\mathbf{t}$ . We call  $\mathcal{T}$  the set of planar rooted trees.

There is a second natural order defined on  $\mathcal{U}$ : the breadth-first order. We write  $u \prec_{\text{bfs}} v$  when  $|u| < |v|$ , or  $|u| = |v|$  and there exists  $i < |u|$  such that  $(j \leq i \Rightarrow u(j) = v(j))$  and  $u(i+1) < v(i+1)$ . We write  $u \preceq_{\text{bfs}} v$  when  $u \prec_{\text{bfs}} v$  or  $u = v$ . Unlike the ancestral order, the breadth-first order is a total order. Again, we shall not distinguish between the order  $\preceq_{\text{bfs}}$  and its restriction to  $\mathbf{t} \in \mathcal{T}$ .

If  $k$  is an integer and  $u_1, \dots, u_k$  are  $k$  distinct vertices of  $\mathbf{t} \in \mathcal{T}$  distinct from the root  $\rho$ , we call  $\mathbf{t}_k = (\mathbf{t}, u_1, \dots, u_k)$  a pointed planar rooted tree, and  $\mathcal{T}_k$  the set of pointed planar rooted trees. Let us stress that the pointed vertices come in a specific order:  $(\mathbf{t}, u_1, u_2)$  and  $(\mathbf{t}, u_2, u_1)$  are for instance two distinct elements of  $\mathcal{T}_2$ . In case every vertex in  $\mathbf{t}_k \in \mathcal{T}_k$  is an ancestor of a pointed vertex (possibly itself), we say that  $\mathbf{t}_k$  is spanned by its pointed vertices, or simply that  $\mathbf{t}_k$  is spanned. Notice that the set of pointed vertices of a spanned tree  $\mathbf{t}_k$  contains its set of leaves. Two pointed planar rooted trees  $\mathbf{t}_k, \mathbf{t}'_k \in \mathcal{T}_k$  are equivalent when there exists a graph-isomorphism between  $\mathbf{t}_k$  and  $\mathbf{t}'_k$  that preserves the root and the pointed vertices. We call  $\overline{\mathcal{T}}_k^s$  the set of equivalence classes of spanned pointed rooted trees (we emphasise that those trees are no more ordered.)

For  $\mathbf{t}_k \in \mathcal{T}_k$  and  $0 \leq i \leq k$ , let  $V_i \subseteq V$  be the set of vertices with precisely  $i$  children that are ancestors of pointed vertices:  $V_i = \{v \in \mathbf{t}, |\{w \in c(v), \exists j, w \preceq u_j\}| = i\}$ . In the special case  $\mathbf{t}_k \in \overline{\mathcal{T}}_k^s$ , every vertex is the ancestor of a pointed vertex, and  $V_i$  is the subset of the vertices with  $i$  children. The following is a partition of the set of vertices,  $V = \bigcup_{0 \leq i \leq k} V_i$ , and we set

$$\ell_i = |V_i|. \quad (17)$$

The definitions of  $\ell_i$  and  $V_i$  extend to a tree  $\mathbf{t}_k \in \overline{\mathcal{T}}_k^s$ .

Let  $\mathbf{p} = (p_k, k \geq 0)$  be a distribution on the non negative-integers, and  $(X_u, u \in \mathcal{U})$  be a collection of i.i.d. random variables indexed by  $\mathcal{U}$  with distribution  $\mathbf{p}$ . Call  $T$  the random tree in which the number of children of a vertex  $u$  is given by  $X_u$  in  $T$ , provided  $X_v \neq 0$  for every strict ancestor  $v$  of  $u$ . The tree  $T$  is distributed as the GW-tree with offspring distribution  $\mathbf{p}$ . Formally,  $\rho \in T$  and for every  $u \in \mathcal{U} \setminus \{\rho\}$ ,  $u \in T$  iff for every  $1 \leq i \leq |u|$

$$u(i) \leq X_{a(u^{(i)})}.$$

### 3.1 Many-to- $k$ formula

Certain functions on  $\mathcal{T}$  can be decomposed as a sum over the different  $k$ -tuples of the vertices of  $\mathbf{t}$ . The expected value of such functions evaluated on GW-trees is then computed using a many-to- $k$  formula (19).

Consider  $\mathbf{p} = (\mathbf{p}_k, k \geq 0)$  a probability distribution on the set  $\mathbb{Z}^+$ . Let  $k_{\max}(\mathbf{p}) = \max \{k \geq 0, \sum_k k^\ell \mathbf{p}_k < \infty\}$ , and  $k_0(\mathbf{p}) = \max \{k \geq 0, \mathbf{p}_k \neq 0\} \in \mathbb{Z}^+ \cup \{\infty\}$ . If  $i \in \mathbb{Z}^+$  satisfies  $i \leq \min \{k_{\max}, k_0\}$ , the  $i$ -th factorial moment and the  $i$ -th size-biased probability distribution  $\mathbf{p}^{(i)} = (\mathbf{p}_k^{(i)}, k \geq 0)$  are defined by:

$$m_i = \sum_{k \geq i} \left[ \prod_{0 \leq j < i} (k - j) \right] \mathbf{p}_k, \text{ and } \mathbf{p}_k^{(i)} = \frac{\prod_{0 \leq j < i} (k - j)}{m_i} \mathbf{p}_k \quad (18)$$

According to our convention that products over an empty set are 1,  $\mathbf{p}^{(0)} = \mathbf{p}$ . Also, for  $k < i$ ,  $\mathbf{p}_k^{(i)} = 0$  since the product in the RHS of (19) contains a null factor. If  $i > k_0(\mathbf{p})$   $m_i = 0$ , and  $\mathbf{p}^{(i)}$  is not defined. Fix  $\mathbf{t}_k = (\mathbf{t}, v_1, \dots, v_k) \in \overline{\mathcal{T}}_k$  with  $k \leq k_{\max}(\mathbf{p}), k_0(\mathbf{p})$ . From  $\mathbf{p}$  and  $\mathbf{t}_k$ , we build  $T_k = T_k(\mathbf{t}_k) = (T, u_1, \dots, u_k) \in \overline{\mathcal{T}}_k$  a random tree that contains  $\mathbf{t}_k$  as a subtree. Formally, there is a copy <sup>4</sup>  $T' \subseteq T$  of  $\mathbf{t}$  embedded in  $T$ ; for a vertex in  $T'$  the offspring distribution is  $\mathbf{p}^{(i)}$  the  $i$ -th size-biased distribution, where  $i$  is the number of children of the corresponding vertex in  $\mathbf{t}$ . For a vertex in  $T \setminus T'$ , the offspring distribution is  $\mathbf{p}^{(0)} = \mathbf{p}$ . A precise definition uses an induction: The trees  $T'$  and  $T, T' \subseteq T$ , and the graph isomorphism  $\omega : \mathbf{t} \rightarrow T'$  are defined by the following steps:

- Initialization: The root  $\rho$  belongs to  $T$  and  $T'$ , it is the image of  $\rho_{\mathbf{t}}$  the root in  $\mathbf{t}$  :  $\rho = \omega(\rho_{\mathbf{t}})$ .
- Induction 1 : If  $u \in T'$ , its number of children in  $T$  is distributed as  $\mathbf{p}^{(i)}$  where  $i = |c(v)|$  is the number of children of  $v = \omega^{-1}(u)$  in  $\mathbf{t}$ . Furthermore, if  $(v_j, 1 \leq j \leq i)$  are the ordered children of  $v$  in  $\mathbf{t}$ , then  $(\omega(v_j), 1 \leq j \leq i)$  is a random sequence of  $i$  distinct children <sup>5</sup> of  $u$  in  $T$ , with uniform distribution. These  $i$  children of  $u$  belong to  $T'$ , and the remaining children of  $u$  belong to  $T \setminus T'$ .
- Induction 2 : If  $u \in T \setminus T'$ , its number of children is distributed as  $\mathbf{p} = \mathbf{p}^{(0)}$ , all of them belong to  $T \setminus T'$ .
- Induction 3 : distinct vertices in  $T$  have an independent number of children.
- The pointed vertices in  $T_k$  are  $(u_1, \dots, u_k) = (\omega(v_1), \dots, \omega(v_j))$

In the next formula,  $m_i = m_i(\mathbf{p})$  is the  $i$ -th factorial moment of the distribution  $\mathbf{p}$ , see (18), and  $\ell_i = \ell_i(\mathbf{t}_k)$  simply counts the number of vertices with precisely  $i$  children since  $\mathbf{t}_k \in \mathcal{T}_k^s$  or  $\mathbf{t}_k \in \overline{\mathcal{T}}_k^s$ , see (17).

**Lemma 3.1** (Many-to- $k$  formula). *Let  $\mathbf{p}$  be a distribution on  $\mathbb{Z}^+$ , and let  $k \in \mathbb{N}$  satisfy  $k \leq k_{\max}(\mathbf{p}), k_0(\mathbf{p})$ . For  $F$  a non-negative measurable function on the set of pointed trees  $\overline{\mathcal{T}}_k$ ,*

<sup>4</sup> $T'$  does not belong to  $\mathcal{T}$  since item (ii) in the definition of  $\mathcal{T}$  is not satisfied, so the word "tree" that we use here is an abuse

<sup>5</sup>this is a.s. possible since the number of children of  $v$  in  $T$ , with distribution  $\mathbf{p}^{(i)}$ , is a.s.  $\geq i$

and  $T \in \mathcal{T}$  a GW-tree with offspring distribution  $\mathfrak{p}$ , it holds:

$$\mathbb{E} \left( \sum_{u_1, \dots, u_k} F(T, u_1, \dots, u_k) \right) = \sum_{\mathfrak{t}_k} \left( \prod_{i \geq 1} m_i^{\ell_i} \right) \mathbb{E}(F(T_k(\mathfrak{t}_k))) \quad (19)$$

where:

- the first sum is over the vertices  $u_1, \dots, u_k$  of  $T$  such that no one is an ancestor of the other, and the second sum is over the trees  $\mathfrak{t}_k$  in  $\overline{\mathcal{T}}_k^s$  that are spanned by  $k$  pointed leaves,
- or the first sum is over the pairwise distinct vertices  $u_1, \dots, u_k$  of  $T$ , and the second sum is over the trees  $\mathfrak{t}_k$  in  $\overline{\mathcal{T}}_k^s$

The so-called many-to-one formula is the  $k = 1$  case. The generation of the leaf uniquely specifies a tree in  $\overline{\mathcal{T}}_1^s$ , which identifies the latter set with the set of non-negative integers  $\mathbb{Z}^+$ , while the product on the RHS reduces to  $m_1^{\ell_1} = m_1^{h_1}$ .

The restriction to  $k \leq k_0(\mathfrak{p})$  is for the sake of simplicity: in case  $k > k_0(\mathfrak{p})$ , the correct formula is obtained by discarding those trees  $\mathfrak{t}_k$  that have vertices with outdegree larger than  $k_0$  in the sums (i) and (ii).

The sequence  $(\ell_i, i \geq 1)$  has finitely many non-null terms, hence the product  $(\prod_{i \geq 1} m_i^{\ell_i})$  is well defined.

*Proof.* Let  $\mathfrak{t}_k^0 = (\mathfrak{t}^0, v_1, \dots, v_k) \in \mathcal{T}_k$  be a pointed planar rooted tree. We first check (19) for  $F = \mathbf{1}_{\mathfrak{t}_k^0}$ . The subtree of  $\mathfrak{t}_k^0$  spanned by  $v_1, \dots, v_k$  is denoted  $\mathfrak{t}_k \in \overline{\mathcal{T}}_k^s$ . For  $i \geq 0$ , set  $V_i^0 = V_i(\mathfrak{t}_k^0)$ ,  $\ell_i^0 = \ell_i(\mathfrak{t}_k^0)$  and  $\ell_i = \ell_i(\mathfrak{t}_k)$ . For  $i \geq 1$ , notice that  $\ell_i^0 = \ell_i$ . Using the definition (18) of  $\mathfrak{p}^{(i)}$ , we find:

$$\begin{aligned} \mathbb{P}(T = \mathfrak{t}^0) &= \prod_{u \in V(\mathfrak{t}^0)} \mathfrak{p}_{|c(u)|} \\ &= \left( \prod_{1 \leq i \leq k_0} m_i^{\ell_i} \right) \prod_{0 \leq i \leq k_0} \prod_{u \in V_i^0} \frac{1}{\prod_{0 \leq j < i} (|c(u)| - j)} \mathfrak{p}_{|c(u)|}^{(i)} \\ &= \left( \prod_{1 \leq i \leq k_0} m_i^{\ell_i} \right) \mathbb{P}(T_k(\mathfrak{t}_k) = \mathfrak{t}_k^0) \end{aligned} \quad (20)$$

Formula (20) is formula (19) with the choice  $F = \mathbf{1}_{\mathfrak{t}_k^0}$ , since one single tree, the tree  $\mathfrak{t}_k$  defined above, contributes to the sum on the RHS of (19). A function  $F$  on  $\mathcal{T}_k$  may be decomposed as a sum of indicator functions, and relation (19) is linear in  $F$ : Formula (20) therefore implies (19), with the sum as in (ii). Adding the restriction that, in the collection  $v_1, \dots, v_k$ , no one is an ancestor of another is equivalent to summing over the trees  $\mathfrak{t}_k$  spanned by pointed leaves, giving (i).  $\square$

### 3.2 Asymptotics for the size of GW-trees

We restrict in this section our views to a particular class of GW-trees. These are, for  $\varepsilon > 0$  and  $\ell \geq 2$  an integer

$$T = T_{\pm} \in \mathcal{T} \text{ the GW-tree with } \text{Bin}(\ell - 1, p_{\pm}) \text{ offspring distribution.} \quad (21)$$

Although the notation does not show it,  $T$  depends on  $\varepsilon$  and  $\ell$ . The next Lemma bounds the probability that a forest of GW-trees distributed as (21) has a given size in term of the probability that a single GW-tree has that same size. We stress the statement is valid for  $\varepsilon$  fixed (i.e., condition (4) is not assumed).

**Lemma 3.2.** *Let  $(T_i, i \geq 1)$  be independent GW-trees distributed as  $T = T_{\pm} \in \mathcal{T}$  the GW-tree in (21). Let  $c < 1/2$ . There exists  $\ell_0 = \ell_0(c)$  such that, for  $\ell \geq \ell_0$ , and for  $1 \leq h \leq j$ :*

$$\mathbb{P} \left( \sum_{1 \leq i \leq h} |T_i| = j \right) \leq h \left( \frac{1-p}{1 \pm \varepsilon} \right)^{h-1} e^{-c \frac{h(h-1)}{j}} \mathbb{P}(|T| = j) \quad (22)$$

*Remark 3.3.* For *critical* GW-trees with a finite variance, similar bounds may be deduced from the estimates in Addario-Berry, Devroye and Janson [ABDJ13], see formulae (17) and (19) in this article.

*Proof.* Let  $(Z_t, t \geq 0)$  be a random walk started at  $Z_0 = 0$  with independent increments distributed as  $X - 1, X$  a  $\text{Bin}(\ell - 1, p)$  random variable. For  $h \in \mathbb{N}$ , we let  $H_{-h}(Z) = \inf \{j \geq 1, Z_j = -h\}$  be the hitting time of  $-h$ . Also, let  $T \in \mathcal{T}$  be a GW-tree with offspring distribution the distribution of  $X$ . If  $\rho = v_1, \dots, v_{|T|}$  is the sequence of the vertices of  $T$  arranged in breadth-first order, then there is the identity in distribution

$$(Z_j, 1 \leq j \leq H_{-1}(Z)) = \left( \left| \bigcup_{i \leq j} \{c(v_i)\} \setminus \bigcup_{i \leq j} \{v_i\} \right|, 1 \leq j \leq |T| \right) \quad (23)$$

where the RHS counts the number of vertices among the children of the vertices  $v_1 \dots v_j$  that do not belong to  $v_1 \dots v_j$ . The identity (23) implies in particular:

$$|T| = H_{-1}(Z) \quad (24)$$

which extends in a straightforward way to a collection of  $h \geq 1$  independent GW-trees  $(T_i, 1 \leq i \leq h)$  distributed as  $T$ :

$$\sum_{1 \leq i \leq h} |T_i| = H_{-h}(Z)$$

We also use the following combinatorial identity, known as Spitzer lemma [Spi56], that connects the distribution of the hitting time of a random walk with the marginal distribution of that random walk <sup>6</sup>

$$j \mathbb{P}(|H_{-h}(Z)| = j) = h \mathbb{P}(Z_j = -h)$$

---

<sup>6</sup>the proof only requires invariance by cyclic shift of the distribution of the increments of the random walk, see *e.g.* Pitman [Pit02]

We deduce from this and the convolution property of the Binomial distribution that:

$$\begin{aligned}
\mathbb{P}\left(\sum_{1 \leq i \leq h} |T_i| = j\right) &= \frac{h}{j} \mathbb{P}(Z_j = -h) \\
&= \frac{h}{j} \mathbb{P}(\text{Bin}((\ell-1)j, p) = j-h) \\
&= \frac{h}{j} \binom{(\ell-1)j}{j-h} p^{j-h} (1-p)^{(\ell-1)j-(j-h)}
\end{aligned} \tag{25}$$

Expanding the binomial coefficient in the last expression allows to relate the above probability to the same probability when  $h$  is set to 1,  $\mathbb{P}(|T_1| = j)$  :

$$\begin{aligned}
h \left(\frac{1-p}{p}\right)^{h-1} \frac{\prod_{1 \leq i < h} (j-i)}{\prod_{j-h \leq i < j-1} (\ell-1)j-i} \cdot \left(\frac{1}{j} \frac{\prod_{0 \leq i < j-1} (\ell-1)j-i}{(j-1)!} p^{j-1} (1-p)^{(\ell-1)j-(j-1)}\right) \\
= h \left(\frac{1-p}{p}\right)^{h-1} \frac{\prod_{1 \leq i < h} (j-i)}{\prod_{j-h \leq i < j-1} (\ell-1)j-i} \mathbb{P}(|T_1| = j)
\end{aligned}$$

Using the definition of  $p = p_{\pm}$  in (4), one can estimate the prefactor:

$$\begin{aligned}
&h \left(\frac{1-p}{p}\right)^{h-1} \frac{\prod_{1 \leq i < h} (j-i)}{\prod_{j-h \leq i < j-1} (\ell-1)j-i} \\
&= h \left(\frac{1-p}{1 \pm \varepsilon}\right)^{h-1} \frac{\prod_{1 \leq i < h} \left(1 - \frac{i}{j}\right)}{\prod_{j-h \leq i < j-1} \left(1 - \frac{i}{j(\ell-1)}\right)} \\
&= h \left(\frac{1-p}{1 \pm \varepsilon}\right)^{h-1} e^{-\frac{h(h-1)}{2j} (1-O(\frac{1}{\ell}))}
\end{aligned} \tag{26}$$

For  $\ell$  large enough,  $O(\frac{1}{\ell}) < 1 - 2c$ , and the inequality (22) is proved.  $\square$

Another key ingredient is a precise estimate on the tail of the size of a single tree  $T$ . Such estimates generally follow from two ingredients: a local central limit theorem on the random walk, plus a centering (or tilting) operation, see Riordan [Rio12] for an implementation of that combination. Our example allows for direct computations. Unlike the previous Lemma, we now assume that  $\varepsilon$  and  $\ell$  depend on  $n$  in a way specified by conditions (3) and (4)

**Proposition 3.4.** *Let  $T = T_{\pm} \in \mathcal{T}$  be the GW-tree given by (21), with sequences  $\varepsilon$  and  $\ell$  satisfying (4) and (3) respectively. Let  $v$  be a sequence that satisfies  $v_n = o(\log \log(\varepsilon^3 n))$ . Then, for  $t = t(u)$  given by (12), it holds that:*

$$\frac{n}{t} \mathbb{P}(t \leq |T| < \infty) = (1 + o(1)) \frac{1}{4\sqrt{\pi}} e^{-u} \text{ as } n \rightarrow \infty \tag{27}$$

with  $o(1)$  uniform over the sequences  $u = (u_n)$  such that  $|u_n| \leq v_n, n \in \mathbb{N}$ .

Under (4), the sequence  $(\varepsilon^3 n)_{n \geq 1}$  diverges and the condition  $|u_n| \leq v_n$  allows for constant sequences  $(u_n)$ . Also, the restriction to finite trees on the LHS of (27) is necessary in the case of supercritical GW-trees only, since subcritical GW-trees are a.s. finite. Last, (27) entails  $|T_+|$  and  $|T_-|$  have equivalent tail (not equal, though).

*Proof.* We first consider a given GW-tree, associated with a fixed  $n$  and  $\ell$ . The symbol  $o_j(1)$  stand for a sequence with null limit as  $j \rightarrow \infty$ . Formula (25) with  $h = 1$  reads:

$$\mathbb{P}(|T| = j) = \frac{1}{j} \binom{(\ell-1)j}{j-1} p^{j-1} (1-p)^{(\ell-1)j-(j-1)} \quad (28)$$

Stirling formula under the form

$$j! = (1 + o_j(1)) \sqrt{2\pi} j^{j+1/2} e^{-j}$$

allows to approximate the Binomial coefficient and to estimate:

$$\begin{aligned} \mathbb{P}(|T| = j) &= \frac{1}{j} \frac{((\ell-1)j)!}{(j-1)!((\ell-2)j+1)!} p^{j-1} (1-p)^{(\ell-1)j-(j-1)} \\ &= (1 + o_j(1)) \frac{1}{j} \frac{1}{\sqrt{2\pi}} \frac{((\ell-1)j)^{(\ell-1)j+1/2}}{(j-1)^{j-1/2}((\ell-2)j+1)^{(\ell-2)j+3/2}} p^{j-1} (1-p)^{(\ell-1)j-(j-1)} \\ &= (1 + o_j(1)) \frac{1}{j} \frac{1}{\sqrt{2\pi}} \left[ \frac{j^{j-1}}{(j-1)^{j-1/2}} \left( \frac{(\ell-1)j}{(\ell-2)j+1} \right)^{(\ell-2)j+3/2} \right] (1 \pm \varepsilon)^{j-1} (1-p)^{(\ell-2)j+1} \end{aligned}$$

The term under bracket requires some care, and may be evaluated using:

$$\begin{aligned} &\frac{j^{j-1}}{(j-1)^{j-1/2}} \left( \frac{(\ell-1)j}{(\ell-2)j+1} \right)^{(\ell-2)j+3/2} \\ &= j^{-1/2} \left( \frac{j}{j-1} \right)^{j-1/2} \left( \frac{(\ell-2)j}{(\ell-2)j+1} \right)^{(\ell-2)j+3/2} \left( \frac{(\ell-1)j}{(\ell-2)j} \right)^{(\ell-2)j+3/2} \\ &= j^{-1/2} \left( 1 + \frac{1}{j-1} \right)^{j-1/2} \left( 1 - \frac{1}{(\ell-2)j+1} \right)^{(\ell-2)j+3/2} \left( 1 + \frac{1}{\ell-2} \right)^{(\ell-2)j+3/2} \\ &= j^{-1/2} (e^1 + o_j(1)) (e^{-1} + o_j(1)) \left( 1 + \frac{1}{\ell-2} \right)^{(\ell-2)j+3/2} \end{aligned}$$

This entails

$$\begin{aligned} \mathbb{P}(|T| = j) &= (1 + o_j(1)) \frac{1}{\sqrt{2\pi}} \frac{1}{j^{3/2}} \left( 1 + \frac{1}{\ell-2} \right)^{(\ell-2)j+3/2} (1 \pm \varepsilon)^{j-1} (1-p)^{(\ell-2)j+1} \\ &= (1 + o_j(1)) q_1 \frac{1}{\sqrt{2\pi}} \frac{1}{j^{3/2}} q_2^j \quad (29) \end{aligned}$$

with the notation

$$q_1 = \left( 1 + \frac{1}{\ell-2} \right)^{3/2} (1 \pm \varepsilon)^{-1} (1-p) \text{ and } q_2 = \left( 1 + \frac{1}{\ell-2} \right)^{\ell-2} (1 \pm \varepsilon) (1-p)^{\ell-2}$$

We now recognise that

$$q_2 = e^{-\delta}, \text{ with } \delta = \delta(\ell - 1, \varepsilon) \text{ defined in (5)} \quad (30)$$

by the following computation:

$$\begin{aligned} \left(1 + \frac{1}{\ell - 2}\right) (1 - p) &= \left(1 + \frac{1}{\ell - 2}\right) \left(1 - \frac{1 \pm \varepsilon}{\ell - 1}\right) \\ &= 1 + \frac{1}{\ell - 2} - \frac{1 \pm \varepsilon}{\ell - 1} - \frac{1 \pm \varepsilon}{(\ell - 1)(\ell - 2)} \\ &= 1 + \frac{(\ell - 1) - (1 \pm \varepsilon)(\ell - 2)}{(\ell - 1)(\ell - 2)} - \frac{1 \pm \varepsilon}{(\ell - 1)(\ell - 2)} \\ &= 1 + \frac{(\ell - 1)(1 - (1 \pm \varepsilon)) + (1 \pm \varepsilon)}{(\ell - 1)(\ell - 2)} - \frac{1 \pm \varepsilon}{(\ell - 1)(\ell - 2)} \\ &= 1 - \frac{\pm \varepsilon}{\ell - 2} \end{aligned}$$

Summing the equivalents in (29) we find that:

$$\mathbb{P}(|T| \geq j) = (1 + o(1)) \frac{q_1}{\sqrt{2\pi}} \sum_{i \geq j} i^{-3/2} e^{-\delta i}$$

We now set  $j = t$  and assume the parameters  $\varepsilon$  and  $\ell$  depend on  $n$ , with  $\varepsilon \rightarrow 0$  and  $\ell \rightarrow \infty$ . The next  $o(1)$  are  $o_n(1)$  as usual. Plainly,  $q_1 = 1 + o(1)$ . The map  $x \mapsto x^{-3/2} e^{-\delta x}$  is non increasing on  $[t, \infty)$ , this implies:

$$\frac{1}{\sqrt{2\pi}} \sum_{i \geq t} i^{-3/2} e^{-\delta i} = (1 + o(1)) \frac{1}{\sqrt{2\pi}} \int_{x \geq t} x^{-3/2} e^{-\delta x} dx$$

Integrating the latter expression we find that

$$\mathbb{P}(t \leq |T| < \infty) = (1 + o(1)) \frac{1}{\sqrt{2\pi}} \delta^{-1} t^{-3/2} e^{-\delta t} \quad (31)$$

The upper bound follows from the bound  $x^{-3/2} \leq t^{-3/2}$  in the integrand, and a lower bound is given by the same integral over  $[t, t']$  instead of  $[t, \infty)$ . Then we choose  $t'$  in such a way that  $\delta(t' - t) \rightarrow \infty$  with  $t' = (1 + o(1))t$ . To complete the proof, it remains to expand  $t = \delta^{-1}s$ :

$$\frac{n}{t} \delta^{-1} t^{-3/2} e^{-\delta t} = \frac{n}{\delta^{-1}s} \delta^{-1} (\delta^{-1}s)^{-3/2} e^{-\delta(\delta^{-1}s)} = \frac{n \delta^{3/2}}{s^{5/2}} e^{-s}$$

and then  $s = s(u)$ , using also  $\delta = (1 + o(1))\varepsilon^2/2$  from (5) (recall  $\varepsilon \rightarrow 0$  and  $\ell \rightarrow \infty$ ) as well as  $\varepsilon^3 n \rightarrow \infty$

$$\begin{aligned} &2^{-3/2} \frac{\varepsilon^3 n}{(\log(\varepsilon^3 n) - \frac{5}{2} \log(\log(\varepsilon^3 n)) + u)^{5/2}} e^{-(\log(\varepsilon^3 n) - \frac{5}{2} \log(\log(\varepsilon^3 n)) + u)} \\ &= (2^{-3/2} + o(1)) e^{-u} \end{aligned}$$

with  $o(1)$  in the last expression uniform in  $|u| \leq v$  with  $v = o(\log \log(\varepsilon^3 n))$ . Bearing in mind the  $1/\sqrt{2\pi}$  factor in (31), the Proposition is proved.  $\square$

*Remark 3.5* (Duality of GW-trees.). Let  $p_+ = (1 + \varepsilon)/(\ell - 1)$  and  $\bar{p}_- = (1 - \bar{\varepsilon})/(\ell - 1)$ . We have, from (5), that  $p_+ (1 - p_+)^{\ell-2} = \bar{p}_- (1 - \bar{p}_-)^{\ell-2}$  if and only if  $\delta_+(\varepsilon, \ell - 1) = \delta_-(\bar{\varepsilon}, \ell - 1)$ . For a fixed  $\varepsilon > 0$ , the latter equation has a unique solution  $\bar{\varepsilon} > 0$ , from the definition of the function  $\delta_{\pm}$ . (28) then implies that, for the GW-trees  $T_+$  and  $\bar{T}_-$  respectively associated with  $p_+$  and  $\bar{p}_-$  as above:

$$(T_+ \mid |T_+| < \infty) \text{ is distributed as } \bar{T}_-$$

It is a general fact that supercritical GW-trees conditioned on being finite have the same distribution as certain subcritical GW-trees. The specificity of the Binomial GW-trees is that the corresponding subcritical GW-tree again are Binomial GW-trees.

We notice the following stability property of the  $\text{Bin}(\ell - 1, p)$  distribution under size-biasing.

**Lemma 3.6.** *Let  $\ell \in \mathbb{N}$ ,  $i \in \{0, 1, \dots, \ell - 1\}$ , and  $p \in (0, 1]$ . If  $p$  is the  $\text{Bin}(\ell, p)$  distribution, then  $\left(p_{i+k}^{(i)}, k \in \mathbb{Z}^+\right)$  is the  $\text{Bin}(\ell - i, p)$  distribution.*

*Proof.* The generating function of  $p$  the  $\text{Bin}(\ell, p)$  distribution is given by:

$$\sum_{k \geq 0} s^k p_k = (1 - p(1 - s))^\ell$$

The  $i$ -th factorial moment of  $p$  is:

$$m_i = p^i \prod_{0 \leq j < i} (\ell - j)$$

and the size-biased distribution  $p^{(i)}$  satisfies:

$$\begin{aligned} p_k^{(i)} &= \frac{\prod_{0 \leq j < i} (k - j)}{p^i \prod_{0 \leq j < i} (\ell - j)} \binom{\ell}{k} p^k (1 - p)^{\ell - k} \\ &= \binom{\ell - i}{k - i} p^{k-i} (1 - p)^{(\ell - i) - (k - i)} \end{aligned}$$

Therefore  $\left(p_{i+k}^{(i)}, k \in \mathbb{Z}^+\right)$  is the  $\text{Bin}(\ell - i, p)$  distribution. Consider a tree  $\mathbf{t}_k \in \overline{\mathcal{T}}_k^s$ , and the corresponding random tree  $T_k(\mathbf{t}_k)$ . The next Lemma estimates the number of vertices in the latter tree, that is the random variable  $|T_k(\mathbf{t}_k)|$ . We set  $h + 1 := |\mathbf{t}_k|$ . Recall the definition of  $\ell_i$  in (17). We have

$$h = \sum_{i \geq 1} i \ell_i \tag{32}$$

since both sides count the number of non-root vertices in  $\mathbf{t}_k$ . □

**Lemma 3.7.** *Let  $T = T_{\pm} \in \mathcal{T}$  be the GW-tree in (21). Let  $j \geq 3$ , and  $c < 1/2$ . There exists  $\ell_0 = \ell_0(c)$  such that, for  $\ell \geq \ell_0$*

$$\mathbb{P}(|T_k(\mathbf{t}_k)| = j) \leq (h + 1)(1 \pm \varepsilon)^{-h} e^{-c \frac{h(h+1)}{j}} \mathbb{P}(|T| = j) \tag{33}$$

*Proof.* We first compare  $p^{(i)}$  and  $p$ , for  $p$  the  $\text{Bin}(\ell - 1, p)$  distribution. Using Lemma 3.6, and expanding the binomial coefficient, we find, for  $1 \leq k \leq \ell - i$ , that:

$$\begin{aligned} p_{k+i}^{(i)} &= \binom{\ell - 1 - i}{k} p^k (1 - p)^{\ell - 1 - i - k} \\ &= \frac{\prod_{0 \leq j < i} (\ell - 1 - j - k)}{\prod_{0 \leq j < i} (\ell - 1 - j)} \cdot \frac{1}{(1 - p)^i} \cdot \binom{\ell - 1}{k} p^k (1 - p)^{\ell - 1 - k} \\ &\leq \frac{1}{(1 - p)^i} \cdot p_k \end{aligned} \tag{34}$$

where the inequality follows by bounding the first factor of the product by 1.

Let  $\mathbf{t}_k = (\mathbf{t}, u_1, \dots, u_k) \in \mathcal{T}_k$ . Denote by  $(v_j)_{1 \leq j \leq h+1}$  the ancestors (large or strict) of the  $k$  pointed vertices (including the pointed vertices themselves) in  $\mathbf{t}_k$  ranked by breadth-first order. (This definition of  $h$  is consistent with the one in (32)). Consider, for  $i \in \{1, \dots, h+1\}$ ,  $U_i$  the set of vertices of  $\mathbf{t}_k$  whose most recent common ancestor in  $\{(v_j)_{1 \leq j \leq h+1}\}$  is  $v_i$ , and define  $\mathbf{t}^{[i]} \in \mathcal{T}$  the tree induced by  $\mathbf{t}$  on  $U_i$  and rooted at  $v_i$ . We apply this construction to  $T_k(\mathbf{t}_k) \in \mathcal{T}_k$ , shortened in  $T_k$  in the following lines. Recall  $T'$  is the random tree embedded in  $T_k$ . Conditionally on  $T'$ , the random tree  $T_k^{[i]}$  is a GW-tree with offspring distribution  $p$ , except for the root that has offspring distribution  $p_{j+}^{(j)}$ , for  $j = |c_{T'}(v_i)|$  the number of children of  $v_i$  in  $T'$ . With the help of (34), we see that:

$$\mathbb{P}\left(T_k^{[i]} = \mathbf{t}\right) \leq (1 - p)^{-|c_{T'}(v_i)|} \mathbb{P}(T = \mathbf{t})$$

Now, the trees  $T_k^{[i]}$  are independent. Let  $(T_i)_{1 \leq i \leq h+1}$  be a collection of independent trees distributed as  $T$ . The latter identity and  $\sum_{1 \leq i \leq h+1} |c_{T'}(v_i)| = h$  together imply that

$$\mathbb{P}(|T_k| = j) = \mathbb{P}\left(\sum_{1 \leq i \leq h+1} |T_k^{[i]}| = j\right) \leq (1 - p)^{-h} \mathbb{P}\left(\sum_{1 \leq i \leq h+1} |T_i| = j\right)$$

Combined with (22) and the inequality  $1 - p \leq 1$ , the latter gives (33) □

### 3.3 The number of path-impure vertices in one large GW-tree.

Let  $G_n$  be a vertex transitive graph on  $n$  vertices with a pointed vertex called the root, and let  $\mathbf{t} \in \mathcal{T}$ . Call  $\ell$  the common degree of the vertices in  $G_n$ , and assume that the number of children of every vertex in  $\mathbf{t}$  is  $\leq \ell - 1$ , except the root of  $\mathbf{t}$  that may have up to  $\ell$  children. Conditionally on  $\mathbf{t}$ , we first define  $\iota : \mathbf{t} \rightarrow G_n$  a random graph homomorphism by induction:

- $\iota(\rho)$  is a random vertex  $x$  in  $G_n$
- If  $v \in \mathbf{t}$  has  $i$  children, denoted by  $v_1, \dots, v_i$ , then  $(\iota(v_j), 1 \leq j \leq i)$  is a random, uniformly distributed, sequence of  $i$  distinct elements of the set of neighbours of  $\iota(v)$  in  $G_n$ , also distinct from  $\iota(a(v))$  in case  $v \neq \rho$ .

- The vertices of  $\mathbf{t}$  are seen in the breadth-first order.

Conditionally on  $\mathbf{t}$  and  $\iota$ , a vertex  $w \in \mathbf{t}$  is called *impure* if there exists a vertex  $v$  smaller than  $w$  in the *breadth-first* order,  $v \prec_{\text{bfs}} w$ , such that  $\iota(v) = \iota(w)$ . In that case, say that  $v$  makes  $w$  impure; also, a vertex  $w \in \mathbf{t}$  is called *path-impure* if it has an ancestor in  $\mathbf{t}$  that is impure. This means that there exists  $v \preceq w$ , and a vertex  $u \prec_{\text{bfs}} v$  such that  $\iota(u) = \iota(v)$ ; say that  $u$  makes  $w$  path-impure in this case. We denote by  $I^1(\mathbf{t})$  the subset of path-impure vertices of  $\mathbf{t}$ .

Fix a graph  $G^0 \subseteq G = G_n$ . Define  $G \setminus G^0$  the graph induced by  $G$  on the vertex set  $V(G) \setminus V(G^0)$ . Conditionally on  $\mathbf{t}$  and  $\iota$ , a vertex of  $\mathbf{t}$  is called  $G^0$ -impure if it is mapped by  $\iota$  to a vertex in  $G^0$ , and is called  $G^0$ -path-impure if it has an ancestor (strict or large) in  $\mathbf{t}$  that is  $G^0$ -impure. We denote by  $I^0(\mathbf{t})$  the subset of  $G^0$ -path-impure vertices of  $\mathbf{t}$ .

For  $v$  a vertex of  $G$ , we let  $\mathcal{C}^0(v)$  be the component that contains  $v$  in the percolation of  $G \setminus G^0$ . We set, for  $\mathbf{t} \in \mathcal{T}$ ,

$$I(\mathbf{t}) = I^0(\mathbf{t}) \cup I^1(\mathbf{t}) \quad (35)$$

We shall consider successively the expected number of path-impure vertices and of  $G^0$ -path-impure vertices in a large GW-tree. To bound the expected number of path-impure vertices, we need the special case  $k = 2$  in formula (19) (many-to-two formula).

**Proposition 3.8.** *Let  $T \in \mathcal{T}$  be the GW-tree given by (21). Assume  $\varepsilon$  satisfies (4) and  $\ell$  satisfies (3). Let also  $I^1 = I^1(T)$  be the subset of the path-impure vertices of  $T$ , and let  $P^k$  be the kernel of the non-backtracking random walk on  $G_n$ . Let  $c < 1/2$ . There exists  $\ell_0 = \ell_0(c)$  such that, for  $\ell \geq \ell_0$  and for any  $j \geq 1$*

$$\mathbb{E} (|I^1(T)| \mid |T| = j) \leq \frac{\pi^{1/2}}{2^{5/2} c^{3/2}} \left( \sum_{k \geq 3} k e^{-c \frac{k^2}{j}} P^k \right) j^{3/2} \quad (36)$$

Notice that, for the RHS of (36) to be  $o(j)$ , we need the term in parenthesis to be  $o(j^{-1/2})$ . We first fix some notation. Let  $\mathbf{t}_2 \in \overline{\mathcal{T}}_2^s$  be a tree spanned by 2 pointed vertices. We denote by  $h_0$  the generation of the most common ancestor of the two pointed vertices, and by  $h_0 + h_1$  and  $h_0 + h_2$  the generations of the two pointed vertices, with  $h_0 + h_1 \leq h_0 + h_2$ . In this way, the triplet  $\mathbf{h} = (h_0, h_1, h_2)$  uniquely defines a tree  $\mathbf{t}_2 \in \overline{\mathcal{T}}_2^s$ , and we abuse notation by writing  $T_2(\mathbf{h})$  for  $T_2(\mathbf{t}_2)$  the associated GW-tree in this case. Let us point that the number  $h$  of non-root vertices in  $\mathbf{t}_2$  then satisfies:

$$h = h_0 + h_1 + h_2 \quad (37)$$

Proof. Let  $j \in \mathbb{N}$ . We use the many-to-two formula in Lemma 3.1 with the index of summation (ii). For  $v, w$  distinct vertices of  $T$ , the choice

$$F(T, v, w) = \mathbf{1}\{w \text{ makes } v \text{ path-impure in } T, |T| = j\}$$

allows to estimate the size of  $I^1(T)$  the subset of the path-impure vertices in the GW-tree  $T$ . (The function  $F$  is, through  $\iota$ , a random function, but the many-to-two formula still holds

for such a function.) Recall the notation  $T_2(\mathbf{h}) = (T(\mathbf{h}), u_1, u_2)$  and the notation  $u_0$  for the most recent common ancestor of  $u_1$  and  $u_2$ . Recall the definition of  $h$  in (37) and the equality (32) on  $\ell_1 + 2\ell_2$ . We point out that  $m_1 = p(\ell - 1) = 1 \pm \varepsilon$  and  $m_2 = p^2(\ell - 1)(\ell - 2) \leq (p(\ell - 1))^2 \leq (1 \pm \varepsilon)^2$ , so the product that appears in (19) satisfies

$$\prod_{1 \leq i \leq k_0} m_i^{\ell_i} = m_1^{\ell_1} m_2^{\ell_2} \leq (1 \pm \varepsilon)^{\ell_1 + 2\ell_2} = (1 \pm \varepsilon)^h. \quad (38)$$

Recall  $u_1^{(i)}$  denotes the ancestor of  $u_1$  at generation  $i$ . Fix  $j$  and apply (19) to find:

$$\begin{aligned} \mathbb{E}(|I^1|, |T| = j) &\leq \mathbb{E} \left( \sum_{v \neq w \in T \setminus \{\rho\}} F(T, v, w) \right) \quad (39) \\ &\leq \sum_{h_0, h_1, h_2} (1 \pm \varepsilon)^h \mathbb{P}(u_2 \text{ makes } u_1 \text{ path-impure in } T_2(\mathbf{h}), |T_2(\mathbf{h})| = j) \\ &\leq \sum_{h_0, h_1, h_2, i} (1 \pm \varepsilon)^h \mathbb{P}(u_2 \prec_{\text{bfs}} u_1^{(i)}, \iota(u_2) = \iota(u_1^{(i)}), |T_2(\mathbf{h})| = j) \\ &\leq \sum_{h_0, h_1, h_2, i} (1 \pm \varepsilon)^h \mathbf{1}_{\{h_0 + h_2 \leq i \leq h_0 + h_1\}} \mathbb{P}(\iota(u_2) = \iota(u_1^{(i)}), |T_2(\mathbf{h})| = j) \\ &= \sum_{h_0, h_1, h_2, k} (1 \pm \varepsilon)^h \mathbf{1}_{\{h_2 + h_2 \leq k \leq h_2 + h_1\}} P^k \mathbb{P}(|T_2(\mathbf{h})| = j). \quad (40) \end{aligned}$$

At the third line, we use the definition of path-impurity of  $u_1$  in term of its ancestors. For any two vertices  $u$  and  $v$ ,  $u \prec_{\text{bfs}} v$  implies  $|u| \leq |v|$ , whence the inequality at the fourth line. We set  $k = h_2 + i - h_0$  at the fifth line,  $k$  is the graph distance between the vertices  $u_2$  and  $u_1^{(i)}$ . Also we use that, for two vertices of  $T_2(\mathbf{h})$  that are mapped to the same vertex of  $G_n$  by  $\iota$ , the image in  $G_n$  of the unique path in  $T_2(\mathbf{h})$  between these vertices is distributed as a loop of the non-backtracking random walk: this follows by construction of  $\iota$ . We now fix  $k$ , and, motivated by (33), compute a sum over  $h_0, h_1, h_2$  in (40): Fix  $c > 0$  an arbitrary positive number, and set

$$A = \sum_{h_0, h_1, h_2} \mathbf{1}_{\{h_2 + h_2 \leq k \leq h_2 + h_1\}} (h_0 + h_1 + h_2) e^{-c \frac{(h_0 + h_1 + h_2)^2}{j}}.$$

First we can sum over  $h_0 \geq 1$ , using a simple comparison with an integral. This leaves

$$A \leq \frac{j}{2c} \sum_{h_1, h_2} \mathbf{1}_{\{h_2 + h_2 \leq k \leq h_2 + h_1\}} e^{-c \frac{(h_1 + h_2)^2}{j}}.$$

There remains two sums to perform. Set  $i = h_1 + h_2$ ; due to the restriction,  $h_2 + h_2 \leq k$  a given value of  $i$  appears at most  $k/2$  times:

$$A \leq \frac{j}{2c} \sum_{h_1, h_2} \mathbf{1}_{\{h_2 + h_2 \leq k \leq h_2 + h_1\}} e^{-c \frac{(h_1 + h_2)^2}{j}} \leq \frac{j}{2c} \sum_{i \geq k} \frac{k}{2} e^{-c \frac{i^2}{j}}$$

The last sum is estimated writing  $i = k + i'$

$$\sum_{j' \geq 0} e^{-c \frac{(k+i')^2}{j}} \leq e^{-c \frac{k^2}{j}} \sum_{i' \geq 0} e^{-c \frac{(i')^2}{j}} \leq \left( \frac{j\pi}{2c} \right)^{1/2} e^{-c \frac{k^2}{j}}$$

so that

$$A \leq \frac{\pi^{1/2}}{2^{5/2} c^{3/2}} k e^{-c \frac{k^2}{j}} j^{3/2}.$$

Inserting (33) into (40) and using the bound on  $A$ , we finally deduce for any  $c < 1/2$ :

$$\mathbb{E}(|I^1|, |T| = j) \leq \frac{\pi^{1/2}}{2^{5/2} c^{3/2}} \left( \sum_{k \geq 3} k e^{-c \frac{k^2}{j}} P^k \right) j^{3/2} \mathbb{P}(|T| = j)$$

and this is estimate (36). In the latter formula, the sum over  $k$  starts from  $k = 3$  because the non-backtracking walk can not do shorter loops. The Proposition now follows by definition of the conditional expectation.

To bound the expected number of  $G^0$ -path-impure vertices, we only need the special case  $k = 1$  in formula (19)

**Proposition 3.9.** *Let  $T \in \mathcal{T}$  be the GW-tree given by (21). Let  $I^0 = I^0(T)$  be the subset of the  $G^0$ -path-impure vertices of  $T$  constructed from the random homomorphism  $\iota$  from  $T$  to  $G_n$ . Let  $c < 1/2$ . There exists  $\ell_0 = \ell_0(c)$  such that, for  $\ell \geq \ell_0$ , and for any  $j \geq 1$*

$$\mathbb{E}(|I^0(T)| \mid |T| = j) \leq \frac{\pi^{1/2} j^{3/2} |G^0|}{2c^{3/2} n} \quad (41)$$

Notice that, whenever  $\iota(\rho) \in G^0$ , we have  $I^0(T) = T$ . The fact that  $\iota(\rho)$  is random is therefore important to avoid starting from  $G^0$  too often.

For a tree  $\mathbf{t}_1 \in \overline{\mathcal{T}}_1^s$  spanned by 1 pointed vertex, we denote by  $h$  the generation of the pointed vertex, which uniquely defines the tree  $\mathbf{t}_1 \in \overline{\mathcal{T}}_1^s$ . We abuse notation by writing  $T_1(h)$  for  $T_1(\mathbf{t}_1)$  in this case.

*Proof.* We do the choice

$$F(T, v) = \mathbf{1} \{v \text{ } G^0\text{-path-impure in } T, |T| = j\}$$

and we use the many-to-one formula: this formula involves the tree  $T_1(h)$  with one single pointed vertex at generation  $h$ . We find, after (38), that the product  $\prod_{1 \leq i \leq k_0} m_i^{\ell_i}$  in (19) simplifies to:

$$m_1^{\ell_1} = (1 \pm \varepsilon)^{h_1}$$

Recall  $u_1^{(i)}$  is the ancestor of  $u_1$  at generation  $i$ . We apply (19) to the GW-tree  $T$  to find:

$$\begin{aligned}
\mathbb{E}(|I^0(T)|, |T| = j) &\leq \mathbb{E}\left(\sum_{v \in T \setminus \{\rho\}} F(T, v)\right) \\
&= \sum_h (1 \pm \varepsilon)^h \mathbb{P}(u_1 G^0 \text{-path-impure in } T_1(h), |T_1(h)| = j) \\
&\leq \sum_{i \leq h} (1 \pm \varepsilon)^h \mathbb{P}(\iota(u_1^{(i)}) \in G^0, |T_1(h)| = j) \\
&= \sum_{i \leq h} (1 \pm \varepsilon)^h \mathbb{P}(\iota(u_1^{(i)}) \in G^0) \mathbb{P}(|T_1(h)| = j) \\
&\leq \frac{|G^0|}{n} \sum_h (1 \pm \varepsilon)^h h \mathbb{P}(|T_1(h)| = j)
\end{aligned} \tag{42}$$

For the last estimate, we used that  $\mathbb{P}(\iota(u_1^{(i)}) \in G^0) = |G^0|/n$  for any  $1 \leq i \leq h$ , which holds because  $\iota(\rho)$  is a random vertex in  $G$ . We now use Lemma 3.7 :

$$\begin{aligned}
\sum_h (1 \pm \varepsilon)^h h \mathbb{P}(|T_1(h)| = j) &\leq \sum_h h^2 e^{-ch^2/j} \mathbb{P}(|T| = j) \\
&\leq \frac{j}{2c} \sum_{h \geq 1} e^{-ch^2/j} \mathbb{P}(|T| = j) \\
&\leq \frac{\pi^{1/2}}{2c^{3/2}} j^{3/2} \mathbb{P}(|T| = j)
\end{aligned} \tag{43}$$

The estimate (41) follows from (42) and (43).  $\square$

We now bound  $I^1(T)$  and  $I^0(T)$  using the basic

**Lemma 3.10.** *Let  $X_n \geq 0$  be a sequence of non-negative random variables with a finite first moment, and  $d_n$  be a sequence such that  $\mathbb{E}(X_n) = o(d_n)$ . There exists a sequence  $b_n$  satisfying*

$$b_n = o(d_n) \text{ and } \mathbb{P}(X_n \geq b_n) = o(1)$$

*Proof.* Set  $a_n = \mathbb{E}(X_n)$  and choose  $b_n = \sqrt{a_n d_n} = o(d_n)$ . By Markov inequality,  $\mathbb{P}(X_n \geq b_n) \leq \mathbb{E}(X_n)/b_n = (a_n/d_n)^{1/2} = o(1)$ .  $\square$

Lemma 3.10 entails the following Corollary to Proposition 3.8.

**Corollary 3.11.** *Let  $t$  be given by (12) and  $\beta'$  be a non-negative sequence such that  $t + \beta' \sim t$ . In the setting of Proposition 3.8, and under condition (7), there exists a sequence  $\beta_1$  such that*

$$\beta_1 = o(\delta^{-1}) \text{ and } \sup_{t \leq t' \leq t + \beta'} \mathbb{P}(|I^1| \geq \beta_1 \mid |T| = t') = o(1). \tag{44}$$

*Proof.* Using condition (7) and the estimate (36) we find the bound:

$$\sup_{t \leq t' \leq t + \beta'} \mathbb{E}(|I^1| \mid |T| = t') \leq \frac{\pi^{1/2}}{2^{5/2} c^{3/2}} \left( \sum_{k \geq 3} k e^{-c \frac{k^2}{t'}} P^k \right) (t + \beta')^{3/2}$$

The assumption  $t + \beta' \sim t$  and the condition (7) ensure the RHS is  $o(\delta^{-1})$ . The existence of a sequence  $\beta_1$  satisfying (44) now follows from Lemma 3.10.  $\square$

**Corollary 3.12.** *Let  $t$  be given by (12), and  $\beta'$  be a non-negative sequence such that  $t + \beta' \sim t$ . In the setting of Proposition 3.9 if  $\varepsilon$  satisfies (4), there exists a sequence  $\beta_0$  such that*

$$\beta_0 = o(\delta^{-1}) \text{ and } \sup_{t', G^0} \mathbb{P}(|I^0| \geq \beta_0 \mid |T| = t') = o(1) \quad (45)$$

with the sup over  $(t', G^0)$  such that  $t \leq t' \leq t + \beta'$  and  $G^0 \leq ct$  for  $c$  a finite constant.

*Proof.* We set  $j = t$  in (41) and observe, using (12) that

$$\frac{t^{5/2}}{n} = \frac{\delta^{-5/2} s^{5/2}}{n} = o(\delta^{-1}).$$

Together with (41), we obtain:  $\sup_{t', G^0} \mathbb{E}(|I^0(T)| \mid |T| = t') = o(\delta^{-1})$  with the sup as indicated above. Lemma 3.10 now applies to give (45).  $\square$

The following Proposition is key to the proof of Theorem 2.1. Notice Proposition 3.4 makes a similar statement without taking into account the path-impure vertices.

**Proposition 3.13.** *Let  $T = T_{\pm} \in \mathcal{T}$  be the GW-tree in (22). Assume  $\varepsilon$  satisfies (4),  $\ell$  satisfies (3), and  $t = t(u)$  is given by (12). There exists a sequence  $\beta$  such that  $\beta = o(\delta^{-1})$  and*

$$\frac{n}{t} \mathbb{P}(t + \beta \leq |T| < \infty, |I(T)| \leq \beta) = (1 + o(1)) \frac{1}{4\sqrt{\pi}} e^{-u} \quad (46)$$

with  $o(1)$  uniform over the sequences  $|u| \leq v$  such that  $v(n) = o(\log \log(\varepsilon^3 n))$ .

Crucial in this estimate is the choice of  $\beta$ : it should be larger than the typical values of  $I(T)$ , but small enough so the replacement of  $t$  by  $t + \beta$  on the LHS in (27) and (46) does not change the limit in the RHS.

*Proof.* We set  $\beta = \beta_0 + \beta_1$  given by the two Corollaries, and we notice that  $t + \beta = \delta^{-1}(s + o(1))$ . Also, we set  $u' = (\log \log(\varepsilon^3 n))^{1/2}$  and<sup>7</sup>  $\beta' = \delta^{-1} u'$ . With this definition, we have that  $t + \beta' \sim t + \beta \sim t$ . Proposition 3.4 applies and we find that the three quantities:

$$\mathbb{P}(t + \beta \leq |T| \leq t + \beta'), \mathbb{P}(t + \beta \leq |T| < \infty), \text{ and } \mathbb{P}(t \leq |T| < \infty) \quad (47)$$

are equivalent as  $n \rightarrow \infty$ . We have the lower and upper bounds:

$$\begin{aligned} (1 - \sup_{t + \beta \leq t' \leq t + \beta'} \mathbb{P}(|I| > \beta \mid |T| = t')) \mathbb{P}(t + \beta \leq |T| \leq t + \beta') &\leq \mathbb{P}(t + \beta \leq |T| < \infty, |I| \leq \beta) \\ &\leq \mathbb{P}(t + \beta \leq |T| < \infty) \end{aligned}$$

By (44), (45), the definition of  $\beta$  and (47), the two bounds are equivalent as  $n \rightarrow \infty$ , moreover  $(n/t) \mathbb{P}(t + \beta \leq |T| < \infty, |I| \leq \beta) = (1 + o(1)) e^{-u} / (4\sqrt{\pi})$ .  $\square$

<sup>7</sup>any sequence  $u'$  satisfying  $1 \ll u' \ll \log \log(\varepsilon^3 n)$  works as well

### 3.4 Modified GW-trees.

In the modified GW-tree  $T_{\pm}^m \in \mathcal{T}$ , every vertex  $v$  has a random independent number of offspring distributed as  $\text{Bin}(\ell - \mathbf{1}_{v \neq \rho}, p)$ . It is called "modified" because the distinct offspring distribution is different at the root. In the case (3) that  $\ell$  diverges, this change affects the asymptotic of the tree only slightly, and the following dominations hold between the tails of  $|T^m|$  and  $|T|$ .

**Lemma 3.14.** *Let  $j \in \mathbb{N}$ . The tail of the random variables  $|T|$  and  $|T^m|$  satisfy:*

$$\mathbb{P}(|T| \geq j) \leq \mathbb{P}(|T^m| \geq j) \leq (1 + O(1/\ell))\mathbb{P}(|T| \geq j) \quad (48)$$

*Proof.* We call natural coupling of  $T$  and  $T^m$  the coupling that uses the same Bernoulli random variables to define the Binomial number of offspring at each vertex. In this coupling, only the number of children of the root may differ, by at most 1. We have  $T \subseteq T^m$ , whence the lower bound in (50). If we consider  $T_1$  and  $T_2$  two independent copies of  $T$  and  $B$  an independent random variable such that  $\mathbb{P}(B = 1) = 1 - \mathbb{P}(B = 0) = 1/(\ell - 1)$  then we have

$$|T^m| = |T_1| + B|T_2|. \quad (49)$$

in this coupling. We also point out that  $\mathbb{P}(|T_1| + |T_2| \geq j) \leq 2(1 + O(\varepsilon))\mathbb{P}(|T| \geq j)$  follows from (23) with  $h = 2$ . These two equations entail the upper bound in (48).

$$\begin{aligned} \mathbb{P}(|T^m| \geq j) &\leq \mathbb{P}(|T_1| \geq j) + \mathbb{P}(B = 1)\mathbb{P}(|T_1| + |T_2| \geq j) \\ &\leq (1 + O(1/\ell))\mathbb{P}(|T| \geq j) \end{aligned}$$

□

Recall the definition of  $\mathcal{C}^0(v)$  a few lines before (35).

**Lemma 3.15.** *There exists a coupling in which, if  $v$  is a uniformly chosen random vertex independent of the percolation of  $G$ :*

$$|T^m \setminus I(T^m)| \leq |\mathcal{C}^0(v)| \leq |T^m| \quad (50)$$

This is essentially the statement of Proposition 11 in [Nac09], with the only difference that we take  $v$  a uniformly chosen random vertex in  $G$  (consider the case of a fixed, deterministic  $v \in G^0$  to see why this is needed). We do not repeat the proof. The lower bound in (50) may be strict: this is because a vertex  $w \in T$  can be path-impure because of a vertex  $v \in T$  that is itself path-impure.

Pruning off the path-impure vertices does not necessarily preserve the inclusion of trees:  $T \subseteq T^m$  does not imply in general  $T \setminus I(T) \subseteq T^m \setminus I(T^m)$ . However, in the coupling (49),  $T$  and  $T^m$  agree in distribution on the event  $B = 0$ , and we always have the lower bound:

$$\begin{aligned} \mathbb{P}(|T^m \setminus I(T^m)| \geq j) &\geq \mathbb{P}(|T \setminus I(T)| \geq j, B = 0) \\ &\geq (1 - O(1/\ell))\mathbb{P}(|T| \geq j + \beta, |I(T)| \leq \beta) \end{aligned}$$

### 3.5 The number of components of $G_n(p)$ with size in a given interval.

The results collected so far are applied in this section to the random graph of interest. For  $G^0 \subseteq G = G_n$  recall  $G \setminus G^0$  is the graph induced by  $G$  on the vertex set  $V(G) \setminus V(G^0)$ . Call  $(\mathcal{C}_j^0, j \geq 1)$  the largest components of  $(G \setminus G^0)(p)$  arranged in non-increasing order of size. Recall the definition of the map  $\mathbb{R} \rightarrow \mathbb{R}, u \mapsto t(u)$  in (12), and set  $t^{-1}$  for the inverse function of  $t$ . A point measure recording the sizes of the components in  $(G \setminus G^0)(p)$  is defined by:

$$N^0 = \sum_{j \geq 1} \delta_{t^{-1}(|\mathcal{C}_j^0|)}$$

and we set  $N$  for  $N^0$  when  $G^0$  is the empty graph. The goal of this section is to show the convergence, in a sense to be precised, of the random measure  $N$  towards a Poisson point measure whose intensity measure  $\mu$  is given by (8). A first step is to compute the first moment of the random measure  $N^0$  evaluated on intervals. Call an interval  $J = (u_1, u_2)$  bounded when  $-\infty < u_1 \leq u_2 < \infty$ , and bounded from the left when  $-\infty < u_1 \leq u_2 \leq \infty$ . In the following, we set  $t_1 = t(u_1)$  and  $t_2 = t(u_2)$ . We stress that the next two propositions are concerned with the weak subcritical regime.

**Proposition 3.16.** *Let  $J$  be a bounded interval. Assume  $\varepsilon$  satisfies (4),  $\ell$  satisfies (3),  $t$  is given by (13) and  $P^k$  satisfies condition (7). The first moment of the random variable  $N_-^0(J)$  satisfies:*

$$\mathbb{E}(N_-^0(J)) \rightarrow \int_J \mu(dx) \text{ as } n \rightarrow \infty$$

and the convergence is uniform over the graphs  $G^0 \subseteq G$  such that  $|G^0| \leq Ct$  for any constant  $C$  independent of  $n$ .

*Proof.* Let  $j, j' \in \mathbb{N}$  with  $j \leq j'$ . The inequality in Lemma 3.15 :

$$\mathbb{P}(|T^m \setminus I(T^m)| \geq j) \leq \mathbb{P}(|\mathcal{C}^0(v)| \geq j) \leq \mathbb{P}(|T^m| \geq j)$$

can also be written<sup>8</sup> in term of  $T$ :

$$(1 - O(1/\ell))\mathbb{P}(|T \setminus I(T)| \geq j) \leq \mathbb{P}(|\mathcal{C}^0(v)| \geq j) \leq (1 + O(1/\ell))\mathbb{P}(|T| \geq j)$$

using the natural coupling of  $T$  and  $T^m$  for the lower bound, see the proof of Lemma 3.14, and again Lemma 3.14 for the upper bound. We can now put the pieces together. We multiply the last inequality by  $n/j$  and then observe that  $\mathbb{P}(|T \setminus I(T)| \geq j) \geq \mathbb{P}(j' \leq |T| < \infty, |I(T)| \leq j' - j)$ . Also we choose  $j = t$  and  $j' = t + \beta$ . The estimates in (27),(46) entail that

$$\lim_{n \rightarrow \infty} \frac{n}{t} \mathbb{P}(|\mathcal{C}_-^0(v)| \geq t) = \frac{1}{4\sqrt{\pi}} e^{-u} \quad (51)$$

The first moment of  $N_-^0(J)$  then satisfies:

---

<sup>8</sup>this step could have been avoided by stating a many-to-two formula for modified GW-trees

$$\begin{aligned}\mathbb{E}(N_-^0(J)) &= \mathbb{E}\left(\sum_v \frac{\mathbf{1}_{\{|\mathcal{C}^0(v)| \in (t_1, t_2)\}}}{|\mathcal{C}^0(v)|}\right) = (1 + o(1)) \frac{n}{t_1} \mathbb{P}(|\mathcal{C}^0(v)| \in (t_1, t_2)) \\ &= (1 + o(1)) \frac{1}{4\sqrt{\pi}} (e^{-u_1} - e^{-u_2})\end{aligned}\quad (52)$$

when  $n \rightarrow \infty$ . Notice we used  $t_1 \sim t_2$  for the second equality. The RHS of (52) is the integral of  $\mu$  given by (8) on  $J$ , which concludes the proof.  $\square$

Recall that  $N$  stands for  $N^0$  when  $G^0$  is the empty graph.

**Proposition 3.17.** *Let  $J$  be a bounded interval. Assume  $\varepsilon$  satisfies (4),  $\ell$  satisfies (3),  $t$  is given by (12) and  $P^k$  satisfies condition (7). The second factorial moment of  $N_-(J)$  satisfies:*

$$\mathbb{E}(N_-(J)(N_-(J) - 1)) \rightarrow \left(\int_J \mu(dx)\right)^2 \quad (53)$$

*Proof.* We start by writing:

$$\begin{aligned}\mathbb{E}(N(J)(N(J) - 1)) &= \mathbb{E}\left(\sum_w \frac{\mathbf{1}_{\{|\mathcal{C}(w)| \in t(J)\}}}{|\mathcal{C}(w)|} (N(J) - 1)\right) \\ &= n \sum_{j \in t(J)} \frac{\mathbb{E}(N(J) - 1, |\mathcal{C}(w)| = j)}{j}\end{aligned}$$

where the sum at the first equality is over the vertices  $w$  of  $G_n$ , and  $w$  may be any vertex (by vertex transitivity) at the second equality. The latter expression may be written

$$(1 + o(1)) \frac{n}{t_1} \sum_{G^0 \subseteq G_n, |G^0| \in (t_1, t_2)} \mathbb{E}(N(J) - 1 \mid \mathcal{C}(w) = G^0) \mathbb{P}(\mathcal{C}(w) = G^0) \quad (54)$$

In term of  $N_-^0$ , this also writes:

$$\mathbb{E}(N_-(J) - 1 \mid \mathcal{C}(w) = G^0) = \mathbb{E}(N_-^0(J)) = (1 + o(1)) \int_J \mu(dx) \quad (55)$$

using Proposition 3.16 for the latter identity, where  $o(1)$  is uniform over the graphs  $G^0 \subseteq G_n$  such that  $|G^0| \in (t_1, t_2)$ . Putting this into (54) and using (55) with  $G^0$  reduced to the empty graph, we conclude that (53) holds.  $\square$

Let  $\mathcal{J}$  denote the set of finite unions of intervals bounded from the left.

**Lemma 3.18.** *Propositions 3.16 and 3.17 are valid for  $J \in \mathcal{J}$ .*

We stress that Propositions 3.16 and 3.17 are stated in the weak subcritical regime. Proposition 3.16 does not hold for every  $J \in \mathcal{J}$  in the weak supercritical regime.

*Proof.* We start with Proposition 3.16. First consider the case when  $J$  is an interval,  $J = (u_1, u_2)$  with  $u_2 = \infty$ . The equivalent  $t_1 \sim t_2$  does no more hold in the equality (52). An inequality replaces that equality:

$$\mathbb{E}(N_-^0(J)) \leq \frac{n}{t_1} \mathbb{P}(|\mathcal{C}^0(v)| \in J) = (1 + o(1)) \frac{1}{4\sqrt{\pi}} e^{-u_1}. \quad (56)$$

The converse inequality is proved by approximating  $J$  by an increasing sequence of finite intervals, and we obtain  $\liminf_{n \rightarrow \infty} \mathbb{E}(N_-^0(J)) \geq e^{-u_1}/4\sqrt{\pi}$ . Linearity of the expectation entails that (52) extends to an arbitrary  $J \in \mathcal{J}$ . We turn to Proposition 3.17. First consider  $J = (u_1, u_2)$  with  $u_2 = \infty$ . As before, a lower bound is achieved by approximation through an increasing sequence of finite intervals:  $\liminf_{n \rightarrow \infty} \mathbb{E}(N_-(J)(N_-(J) - 1)) \geq (e^{-u_1}/4\sqrt{\pi})^2$ . For the upper bound, we notice that  $|\mathcal{C}^0(v)| \leq |\mathcal{C}(v)|$  and this entails, since  $u_2 = \infty$ , that:

$$\mathbb{E}(N^0(J)) \leq \frac{n}{t_1} \mathbb{P}(|\mathcal{C}^0(v)| \in J) \leq \frac{n}{t_1} \mathbb{P}(|\mathcal{C}(v)| \in J).$$

We compute as in (54)

$$\begin{aligned} \mathbb{E}(N(J)(N(J) - 1)) &\leq \frac{n}{t_1} \sum_{|G^0| \in \mathfrak{t}(J)} \mathbb{E}(N^0(J)) \mathbb{P}(\mathcal{C}(w) = G^0) \\ &\leq \left(\frac{n}{t_1}\right)^2 \mathbb{P}(|\mathcal{C}(v)| \in \mathfrak{t}(J)) \sum_{|G^0| \in \mathfrak{t}(J)} \mathbb{P}(\mathcal{C}(w) = G^0) \\ &\leq \left(\frac{n}{t_1} \mathbb{P}(|\mathcal{C}(v)| \in \mathfrak{t}(J))\right)^2 \end{aligned}$$

But we know from (56) with  $G^0$  the empty graph that  $(n/t_1) \mathbb{P}(|\mathcal{C}(v)| \in \mathfrak{t}(J)) = (1 + o(1))e^{-u_1}/4\sqrt{\pi}$  in the weak subcritical regime, and the upper bound  $\limsup_{n \rightarrow \infty} \mathbb{E}(N_-(J)(N_-(J) - 1)) \leq (e^{-u_1}/4\sqrt{\pi})^2$  follows. The case of an arbitrary  $J \in \mathcal{J}$  is similar.  $\square$

Let  $\mathcal{M}(\mathbb{R})$  be the set of locally finite measures on the Borel sigma-algebra of  $\mathbb{R}$ . A measure  $M \in \mathcal{M}(\mathbb{R})$  is called a point measure when  $M(J)$  takes values in  $\mathbb{N}$  for any  $J$  bounded Borel set; a point measure is further called simple when  $M(\{x\}) \in \{0, 1\}$  for any  $x \in \mathbb{R}$ . A random element of  $\mathcal{M}(\mathbb{R})$  is called a random measure. A sequence  $(M_n)_{n \in \mathbb{N}}$  of random measures weakly converges (resp. vaguely converges) towards a random measure  $M$  when the sequence of random variables  $\int M_n(dx) f(x)$  weakly converge to  $\int M(dx) f(x)$  for each  $f$  continuous and bounded (resp. continuous bounded and compactly supported). Proposition 16.17 in [Kal02], reproduced below, gives a criterion for the vague convergence of probability measures in term of the void probabilities.

**Proposition 3.19.** *Let  $(M_n)_{n \geq 1}$  be a sequence of random point measures on  $\mathbb{R}$ , and let  $M$  be a random simple point measure. Then  $(M_n)_{n \geq 1}$  vaguely converges to  $M$  if the following two conditions hold:*

- $\lim_{n \rightarrow \infty} \mathbb{P}(M_n(J) = 0) = \mathbb{P}(M(J) = 0)$ , for any  $J$  finite union of bounded intervals of  $\mathbb{R}$ .

- $\limsup_{n \rightarrow \infty} \mathbb{E}(M_n(K)) \leq \mathbb{E}(M(K))$ , for any compact set  $K$ .

Theorem 2.1 requires the weak convergence of the measures restricted to intervals bounded from the left. The latter may be turned into vague convergence by compactification of the space.

*Proof of Theorem 2.1.* To stress on the dependence in  $n$ , we write  $N_n$  for  $N$  in this proof. Let  $k \geq 1$ , and  $J \in \mathcal{J}$ . By an induction argument, we arrive at the following generalisation of Proposition 3.17 : the  $k$ -th factorial moment satisfies

$$\mathbb{E} \left( \prod_{0 \leq i \leq k-1} (N_n(J) - i) \right) = (1 + o(1)) \left( \int_J \mu(dx) \right)^k \quad (57)$$

For  $k \in \mathbb{N}$ , the  $k$ -th moment of a random variable is a linear combination of the  $i$ -th factorial moments for  $1 \leq i \leq k$ , hence the convergence (57) of the factorial moments entails that of the usual moments. Moreover, the Poisson distribution is uniquely determined by its moments. By the method of moments, see *e.g.* Section 6.1 of [JLR11] and the Theorem 6.1 in particular,  $N_n(J)$  weakly converges towards the Poisson distribution with parameter  $\int_J \mu(dx)$ . This implies in particular the convergence of the void probabilities:

$$\lim_{n \rightarrow \infty} \mathbb{P}(N_n(J) = 0) = e^{-\int_J \mu(dx)} \quad (58)$$

The proof is concluded applying Proposition 3.19 to  $N_n^\phi$  the push-forward of  $N_n$  by an increasing diffeomorphism  $\phi : \mathbb{R} \rightarrow (-\infty, 0)$ , *e.g.*  $x \mapsto -e^{-x}$ . Call  $\mu^\phi$  the push-forward measure  $\mu$  by the map  $\phi$ . Let  $J$  be a finite union of bounded intervals. Equation (58) is equivalent to:

$$\lim_{n \rightarrow \infty} \mathbb{P}(N_n^\phi(J) = 0) = e^{-\int_J \mu^\phi(dx)}$$

Let  $\varepsilon > 0$ , and  $K$  be a compact set of the real line. For  $x \in \mathbb{R}$ , set  $d(x, K) = \inf\{d(x, y), y \in K\}$  for the distance of  $x$  to  $K$ . For  $\eta > 0$ , the set  $O_\eta = \{x; d(x, K) < \eta\}$  is an open set of the real line, hence it can be written as a union of open intervals, that is furthermore finite. By monotone convergence there exists  $\eta > 0$  so that  $\mu^\phi(O_\eta \setminus K) < \varepsilon$ . Proposition 3.16 applies:

$$\mathbb{E}(N_n^\phi(K)) \leq \mathbb{E}(N_n^\phi(O_\eta)) \rightarrow \int_{O_\eta} \mu^\phi(dx) \leq \int_K \mu^\phi(dx) + \varepsilon$$

and this proves, since  $\varepsilon$  is arbitrarily small, that

$$\limsup_{n \rightarrow \infty} \mathbb{E}(N_n^\phi(K)) \leq \int_K \mu^\phi(dx), \quad (59)$$

for  $K$  a compact set. (58) and (59) are the assumptions to apply Proposition 3.19, which entails the vague convergence of  $N_n^\phi$  towards the Poisson measure with intensity  $\mu^\phi$ . The latter in turn is equivalent to the weak convergence of  $N_n$  to the Poisson measure with intensity  $\mu$ , when both measures are restricted to intervals bounded from the left. This is the statement of Theorem 1.1.  $\square$

### 3.6 Adaptation to the case of graphs with bounded degrees

Interesting examples of expander graphs with bounded degrees are known (the Ramanujan graphs constructed in [LPS88] for instance), and we would like them to be included in our analysis. This requires modifications in the both the statement and the proof of Theorem 2.1. Precisely, assume

$$\ell(n) \geq 3, \quad \lim_{n \rightarrow \infty} \ell(n) = L < \infty$$

which means that  $\ell(n)$  is constant equal to  $L$  for  $n$  large enough. Theorem (2.1) then holds with the intensity of the Poisson point measure in that Theorem replaced by its multiple

$$\frac{1}{4\sqrt{\pi}} \frac{L}{L-1} e^{-x} dx$$

Observe that, when  $L = \infty$ , the convention  $L/(L-1) = 1$  allows to recover the original intensity measure. If  $\ell = \ell(n)$  and  $\varepsilon = \varepsilon(n)$  satisfy (4) and  $n$  goes to infinity, one has

$$\delta_{\pm}(\varepsilon, \ell - 1) \sim \frac{(L-1)(L-3)}{(L-2)^2} \frac{\varepsilon^2}{2} \text{ as } n \rightarrow \infty :$$

and the first order of the size of the largest components is now  $2\varepsilon^{-2} \log(\varepsilon^3 n) (L-2)^2 / ((L-1)(L-3))$ . To prove Theorem 2.1 in this new setting, we first observe using the decomposition (51) and the exact computation (31) for the size of the union of  $h = 2\text{GW}$  trees that the estimate (50) on the tail of the modified GW tree has to be replaced by

$$\frac{n}{t} \mathbb{P}(|T_-^m| \geq t) = (1 + o(1)) \frac{1}{4\sqrt{\pi}} \frac{L}{L-1} e^{-u} \text{ as } n \rightarrow \infty \quad (60)$$

in other words,  $T_-$  and  $T_-^m$  have no more equivalent tails (in the scale  $t = t(u)$ ). We then need an estimate similar to (60) with  $T_-^m \setminus I(T_-^m)$  in place of  $T_-^m$  in order to prove the analogue of (51). There is no way round but to find a many-to- $k$  formula in the context of modified GW-trees. This modification is achieved as follows: in Lemma 3.1, if  $\tilde{m}_i$  denotes the  $i$ -th factorial moment of the offspring distribution at the root, the RHS in (22) is multiplied by  $\tilde{m}_{c_{\text{tk}}(\rho)} / m_{c_{\text{tk}}(\rho)}$  (remember that the root only has a distinct offspring distribution in modified GW-trees). In our case, the ratio  $\tilde{m}_i / m_i$  is  $\ell / (\ell - i)$ , and the upper bounds in the estimates (36) and (41) are multiplied by a positive constant independent of  $n$ . Since condition (7) is not sensitive to constants, we conclude that, under (7)

$$\frac{n}{t} \mathbb{P}(|T_-^m \setminus I(T_-^m)| \geq t) = (1 + o(1)) \frac{1}{4\sqrt{\pi}} \frac{L}{L-1} e^{-u} \text{ as } n \rightarrow \infty \quad (61)$$

and the rest of the proof follows unchanged from this point on.

### 3.7 A remark on the second largest component in the weak super-critical regime.

Recall the definition of the quantity  $\bar{\varepsilon}$  in Remark 3.5. There is a conjectured parallel, known as the discrete duality principle, between the largest components in  $G_n(\bar{p}_-)$ ,  $\bar{p}_- = (1 - \bar{\varepsilon}) / (\ell - 1)$

and the largest components from the second one in  $G_n(p_+)$   $p_+ = (1 + \varepsilon)/(\ell - 1)$ . This principle has been proved for a few graphs, among which the complete graph [NP07] and the configuration model [Rio12]. The proof usually relies on the possibility to characterise the (random) graph induced by  $G_n(p_+)$  on the complement of  $\mathcal{C}_1$  in a simple way. We believe that an analogous result should hold in our case, yet the assumption we make or the methods we use in this paper only give a lower bound for the expected number of components in certain intervals, as shown below. The random measure of interest is again

$$N = \sum_{j \geq 1} \delta_{t^{-1}(|\mathcal{C}_j|)}$$

and the lower bound is on the expected number of components with size in the interval  $t(u_1, u_2) = (t_1, t_2)$ , as follows.

**Proposition 3.20.** *Assume  $\varepsilon$  satisfies (4),  $\ell$  satisfies (3), and let  $p_+ = (1 + \varepsilon)/(\ell - 1)$ . Assume the non-backtracking random walk on  $G_n$  satisfies (7). It holds, for  $J = (u_1, u_2)$ ,  $-\infty < u_1 \leq u_2 < \infty$ , that:*

$$\mathbb{E}(N(J)) \geq (1 + o(1)) \frac{1}{4\sqrt{\pi}} (e^{-u_1} - e^{-u_2})$$

as  $n \rightarrow \infty$ .

*Proof.* The set  $I^1(T)$  is the subset of path-impure vertices in  $T$  defined from the random homomorphism  $\iota$  from  $T$  to  $G_n$ . For  $\beta = o(\delta^{-1})$  as in Corollary 3.11, we have:

$$\begin{aligned} \mathbb{E}(N(J)) &= (1 + o(1)) \frac{n}{t_1} \mathbb{P}(|\mathcal{C}(v)| \in (t_1, t_2)) \\ &\geq (1 + o(1)) \frac{n}{t_1} \mathbb{P}(|T| \in (t_1 + \beta, t_2), |I^1| \leq \beta) \\ &= (1 + o(1)) \frac{1}{4\sqrt{\pi}} (e^{-u_1} - e^{-u_2}) \end{aligned}$$

using  $t_1 \sim t_2$  at the first and second line and Proposition 3.4 (with condition (7)) at the second line.  $\square$

With respect to the computation (52), we obtain an inequality in place of an equality at the second line. Let  $\eta > 0$  and  $j \in \mathbb{N}$  be fixed, independent of  $n$ . The lower bound above suggests that, with high probability as  $n \rightarrow \infty$ , the  $j$ -th largest component has size at least  $\delta_+(\varepsilon, \ell - 1)^{-1} (\log(\varepsilon^3 n) - (5/2 + \eta) \log \log(\varepsilon^3 n))$

$$\mathbb{P} \left( \frac{\delta_+(\varepsilon, \ell - 1) |\mathcal{C}_j| - \log(\varepsilon^3 n)}{\log \log(\varepsilon^3 n)} > -(5/2 + \eta) \right) = 1 - o(1) \quad (62)$$

To prove this statement, it would be enough to have an upper bound on the second (factorial) moment. For  $v, w \in V(G_n)$ , we write  $v \not\sim w$  if  $v$  and  $w$  are not connected by a path of open edges. Now,  $N(J)(N(J) - 1)$  counts the number of ordered pair of distinct components, and

It follows that if  $V$  and  $W$  stand for two independent uniform vertices in  $V(G_n)$  under  $\mathbb{P}$ ,

$$\mathbb{E}(N(J)(N(J) - 1)) = (1 + o(1)) \left( \frac{n}{t_1} \right)^2 \mathbb{P}(|\mathcal{C}(V)| \in t(J), |\mathcal{C}(W)| \in t(J), V \not\sim W)$$

We did not find <sup>9</sup> an obvious way to bound  $\mathbb{P}(|\mathcal{C}(V)| \in t(J), |\mathcal{C}(W)| \in t(J), V \not\sim W)$  by  $\mathbb{P}(|\mathcal{C}(V)| \in t(J))^2$ , which is the first step to implement the second moment method and conclude to (67). It may be that further conditions are necessary to prove (62).

### 3.8 Verification of condition (7)

We rely on [Nac09] to check condition (7). To that aim, it is useful first to relate (7) with the assumption

$$n^{1/3} \sum_{k=1}^{n^{1/3}} k P^k = O(1) \tag{63}$$

made by Nachmias in the study of the critical regime. Assumption (63) alone is not enough to check condition (7). One also needs the following condition <sup>10</sup>: there exists a finite constant  $c$  independent of  $n$  such that for  $n$  large enough,

$$P^k \leq \frac{c}{n} \text{ for } k \geq n^{1/3}. \tag{64}$$

**Lemma 3.21.** *Assume the sequence  $\varepsilon$  satisfies (4). Conditions (63) and (64) imply (7).*

*Proof.* The first  $n^{1/3}$  terms in the sum in (7) are bounded using Nachmias condition (63)

$$t^{1/2} \sum_{k=1}^{n^{1/3}} k e^{-c \frac{k^2}{t}} P^k \leq t^{1/2} n^{-1/3} \left( n^{1/3} \sum_{k=1}^{n^{1/3}} k P^k \right) = O \left( \frac{s^{1/2}}{(\delta^{3/2} n)^{1/3}} \right)$$

and it is simple to check that  $s^{1/2} / (\delta^{3/2} n)^{1/3} = o(1/s)$ . For the subsequent terms in the sum, we have from assumption (64) that:

$$t^{1/2} \sum_{k=n^{1/3}}^{\infty} k e^{-c \frac{k^2}{t}} P^k = O \left( \frac{t^{3/2}}{n} \right)$$

with room to spare. Then

$$\frac{t^{3/2}}{n} = \frac{(\delta^{-1} s)^{3/2}}{n} = \frac{s^{3/2}}{\delta^{3/2} n} = o \left( \frac{1}{s} \right) \tag{65}$$

---

<sup>9</sup>The possibility of closed edges with endvertices in both  $|\mathcal{C}(V)|$  and  $|\mathcal{C}(W)|$  prevents us from using the van den Berg-Kesten-Reimer (BKR) inequality

<sup>10</sup>In practice, checking (64) usually does not raise additional difficulties with respect to (63). In [Nac09] for instance, (63) is checked by proving in the first place that the bound on  $P^k$  in (64) holds for  $k \geq \log(n)$  which entails (64).

using for the last estimate that  $\delta^{3/2}n = (2^{-3/2} + o(1))\varepsilon^3n$ , as follows from (10), and then  $\log^{5/2}(\varepsilon^3n) = o(\varepsilon^3n)$ , as follows from (4). Both terms in the sum in the LHS of (7) are  $o(1/s)$ , hence the condition is satisfied.  $\square$

Estimates on the kernel of the non-backtracking walk computed in [Nac09] then yield Proposition 2.6.

*Proof of Proposition 2.6.* For the two graphs, condition (63) is checked in Theorem 2 and, under condition 11, in Theorem 6 of [Nac09]. Also, condition (64) is checked along the proofs of these two theorems: see p.1177 for the expander graphs and p.1178 for the Hamming graph in that same reference. Lemma 3.21 concludes the proof.  $\square$

In the case of the Hamming graph in dimension 1, 2 and 3, one may check (7) by hand, without the intermediate step of checking condition (63). Linking conditions (7) and (63) allows us not to display tedious but straightforward computations.

**Acknowledgments.** We are grateful to M.J. Luczak for asking us a question about the percolation of Hamming graphs that started this work and for sharing her expertise in the field.

## References

- [ABDJ13] Louigi Addario-Berry, Luc Devroye, and Svante Janson. Sub-Gaussian tail bounds for the width and height of conditioned Galton–Watson trees. *The Annals of Probability*, 41(2):1072–1087, 2013.
- [Ald97] David Aldous. Brownian excursions, critical random graphs and the multiplicative coalescent. *Ann. Probab.*, 25(2):812–854, 1997.
- [BCvdH<sup>+</sup>05a] Christian Borgs, Jennifer T. Chayes, Remco van der Hofstad, Gordon Slade, and Joel Spencer. Random subgraphs of finite graphs. I. The scaling window under the triangle condition. *Random Structures Algorithms*, 27(2):137–184, 2005.
- [BCvdH<sup>+</sup>05b] Christian Borgs, Jennifer T. Chayes, Remco van der Hofstad, Gordon Slade, and Joel Spencer. Random subgraphs of finite graphs. II. The lace expansion and the triangle condition. *Ann. Probab.*, 33(5):1886–1944, 2005.
- [BCvdH<sup>+</sup>06] Christian Borgs, Jennifer T. Chayes, Remco van der Hofstad, Gordon Slade, and Joel Spencer. Random subgraphs of finite graphs. III. The phase transition for the  $n$ -cube. *Combinatorica*, 26(4):395–410, 2006.
- [Bol84] Béla Bollobás. The evolution of random graphs. *Trans. Amer. Math. Soc.*, 286(1):257–274, 1984.

- [BR09] Béla Bollobás and Oliver Riordan. Random graphs and branching processes. In *Handbook of large-scale random networks*, volume 18 of *Bolyai Soc. Math. Stud.*, pages 15–115. Springer, Berlin, 2009.
- [ER60] Paul Erdős and Alfréd Rényi. On the evolution of random graphs. *Magyar Tud. Akad. Mat. Kutató Int. Közl.*, 5:17–61, 1960.
- [FVDHDDHH20] Lorenzo Federico, Remco Van Der Hofstad, Frank Den Hollander, and Tim Hulshof. Expansion of percolation critical points for hamming graphs. *Combinatorics, Probability and Computing*, 29(1):68–100, 2020.
- [HN20] Tim Hulshof and Asaf Nachmias. Slightly subcritical hypercube percolation. *Random Structures & Algorithms*, 56(2):557–593, 2020.
- [HR15] Simon Harris and Matthew Roberts. The many-to-few lemma and multiple spines. *Annales de l’Institut Henri Poincaré, Probabilités et Statistiques*, 2015.
- [INW69] Nobuyuki Ikeda, Masao Nagasawa, and Shinzo Watanabe. Branching Markov processes iii. *Journal of Mathematics of Kyoto University*, 9(1):95–160, 1969.
- [Jan11] Svante Janson. Poset limits and exchangeable random posets. *Combinatorica*, 31(5):529–563, 2011.
- [JLR11] Svante Janson, Tomasz Łuczak, and Andrzej Ruciński. *Random graphs*, volume 45. John Wiley & Sons, 2011.
- [Kal02] Olav Kallenberg. *Foundations of modern probability*. Probability and its Applications (New York). Springer-Verlag, New York, second edition, 2002.
- [LP17] Russell Lyons and Yuval Peres. *Probability on trees and networks*, volume 42. Cambridge University Press, 2017.
- [LPS88] A. Lubotzky, R. Phillips, and P. Sarnak. Ramanujan graphs. *Combinatorica*, 8(3):261–277, 1988.
- [Łuc90] Tomasz Łuczak. Component behavior near the critical point of the random graph process. *Random Structures & Algorithms*, 1(3):287–310, 1990.
- [Nac09] Asaf Nachmias. Mean-field conditions for percolation on finite graphs. *Geom. Funct. Anal.*, 19(4):1171–1194, 2009.
- [NP07] Asaf Nachmias and Yuval Peres. Component sizes of the random graph outside the scaling window. *ALEA Lat. Am. J. Probab. Math. Stat.*, 3:133–142, 2007.

- [Pit02] J. Pitman. Combinatorial stochastic processes. Lectures from the 32nd Summer School on Probability Theory held in saint-flour, july 7–24, 2002. *Lecture Notes in Mathematics*, 2002.
- [Rio12] Oliver Riordan. The phase transition in the configuration model. *Combinatorics, Probability and Computing*, 21(1-2):265–299, 2012.
- [RW17] O. Riordan and L. Warnke. The phase transition in bounded-size Achlioptas processes. *ArXiv e-prints*, April 2017.
- [Spi56] Frank Spitzer. A combinatorial lemma and its application to probability theory. *Transactions of the American Mathematical Society*, 82(2):323–339, 1956.
- [vdHL10] Remco van der Hofstad and Malwina J. Luczak. Random subgraphs of the 2D Hamming graph: the supercritical phase. *Probab. Theory Related Fields*, 147(1-2):1–41, 2010.
- [vdHN15] Remco van der Hofstad and Asaf Nachmias. Hypercube percolation. *Journal of the European Mathematical Society*, 2015.