# A simple branching process approach to the phase transition in $G_{n,p}$

Béla Bollobás[*]

Department of Pure Mathematics and Mathematical Statistics
Wilberforce Road, Cambridge CB3 0WB, UK

b.bollobas@dpmms.cam.ac.uk

Oliver Riordan

Mathematical Institute, University of Oxford
24–29 St Giles', Oxford OX1 3LB, UK

riordan@maths.ox.ac.uk

### Abstract

It is well known that the branching process approach to the study of the random graph $G_{n,p}$ gives a very simple way of understanding the size of the giant component when it is fairly large (of order $\Theta(n)$). Here we show that a variant of this approach works all the way down to the phase transition: we use branching process arguments to give a simple new derivation of the asymptotic size of the largest component whenever $(np-1)^3 n \to \infty$.

## 1 Introduction

Our aim in this note is to show how basic results about the survival probability of branching processes can be used to give an essentially best possible result about the emergence of the giant component in $G_{n,p}$, the random graph with vertex set $[n] = \{1, 2, \ldots, n\}$ in which each edge is present independently with probability $p$. In 1959, Erdős and Rényi [4] showed that if we take $p = p(n) = c/n$ where $c$ is constant, then there is a 'phase transition' at $c = 1$. We write $L_i(G)$ for the number of vertices in the $i$th largest component of a graph $G$. Also, as usual, we say that an event holds *with high probability* or *whp* if its probability tends to 1 as $n \to \infty$. Erdős and Rényi showed that, whp, if $c < 1$ then

---

$L_1(G_{n,c/n})$ is of logarithmic order, if $c = 1$ it is of order $n^{2/3}$, while if $c > 1$ then there is a unique 'giant' component containing $\Theta(n)$ vertices, while the second largest component is much smaller.

In 1984, Bollobás [1] noticed that this is only the starting point, and an interesting question remains: what does the component structure of $G_{n,p}$ look like for $p = (1 + \varepsilon)/n$, where $\varepsilon = \varepsilon(n) \to 0$? He and Łuczak [6] showed that if $\varepsilon = O(n^{-1/3})$ then $G_{n,p}$ behaves in a qualitatively similar way to $G_{n,1/n}$; this range of $p$ is now called the *scaling window* or *critical window* of the phase transition. The range $\varepsilon^3 n \to \infty$ is the *supercritical* regime, characterized by the fact that there is whp a unique 'giant' component that is much larger than the second largest component. The range $\varepsilon^3 n \to -\infty$ is the *subcritical* regime.

In this paper we are interested in the size of the giant component as it emerges. Thus we consider the (weakly) supercritical regime where $p = p(n) = (1 + \varepsilon)/n$, with $\varepsilon = \varepsilon(n)$ satisfying

$$\varepsilon \to 0 \qquad \text{and} \qquad \varepsilon^3 n \to \infty \qquad \text{as } n \to \infty. \tag{1}$$

Our aim here is to use branching processes to give a very simple new proof of the following result, originally due to Bollobás [1] (with a mild extra assumption) and Łuczak [6].

**Theorem 1.** *Under the assumption* (1) *we have*

$$L_1(G_{n,p}) = (2 + o_{\mathrm{p}}(1))\varepsilon n$$

*and* $L_2(G_{n,p}) = o_{\mathrm{p}}(\varepsilon n)$.

Here $o_{\mathrm{p}}(f(n))$ denotes a (random) quantity $X_n$ such that $X_n/f(n)$ tends to 0 in probability: the statement is that for any fixed $\delta > 0$, with probability tending to 1 as $n \to \infty$, $L_1(G_{n,p})$ is in the range $(2 \pm \delta)\varepsilon n$ and $L_2(G_{n,p}) \leqslant \delta \varepsilon n$.

Since the original papers [1, 6] (which in fact gave a more precise bound than that above), many different proofs of many forms of Theorem 1 have been given. For example, Nachmias and Peres [7] used martingale methods to reprove the result as stated here. Pittel and Wormald [8] used counting methods to prove an even more precise result; a simpler martingale proof of (part of) their result is given in [3]. A proof of Theorem 1 combining tree counting and branching process arguments appears in [2]. More recently, aiming for simplicity rather than sharpness, Krivelevich and Sudakov [5] gave a very simple proof of (among other things) a weaker form of Theorem 1, where the size of the giant component is determined only up to a constant factor.

## 2 Branching process preliminaries

Let us start by recalling some basic concepts and results. The *Galton–Watson branching process* with offspring distribution $Z$ is the random rooted tree constructed as follows: start with a single root vertex in generation 0. Each vertex in generation $t$ has a random number of children in generation $t + 1$, with distribution $Z$. The numbers of children are independent of each other and of the history. It is well known and easy to check that if

$\mathbb{E}[Z] > 1$, then the process *survives* (is infinite) with probability $\varrho$ the unique solution in $(0, 1]$ to $1 - \varrho = f_Z(1 - \varrho)$, where $f_Z$ is the probability generating function of $Z$. When $\mathbb{E}[Z] < 1$, the expectation of the total number of vertices in the branching process is

$$1 + \mathbb{E}[Z] + \mathbb{E}[Z]^2 + \cdots = \frac{1}{1 - \mathbb{E}[Z]}, \tag{2}$$

and in particular the survival probability is 0.

Let us write $\mathcal{T}_{n,p}$ for the *binomial branching process* with parameters $n$ and $p$, i.e., for the branching process as above with offspring distribution $\mathrm{Bi}(n, p)$. Since the generating function of $\mathrm{Bi}(n, p)$ satisfies

$$f(x) = \sum_{k=0}^{n} \binom{n}{k} p^k (1-p)^{n-k} x^k = \left(1 - p(1-x)\right)^n,$$

when $np > 1$ the survival probability $\varrho = \varrho_{n,p}$ satisfies

$$1 - \varrho = (1 - p\varrho)^n.$$

From this it is easy to check that if $\varepsilon = \varepsilon(n) = np - 1 \to 0$ with $\varepsilon > 0$ then

$$\varrho \sim 2\varepsilon \tag{3}$$

as $n \to \infty$.

Conditioning on a suitable branching process dying out (i.e., having finite total size) one obtains another branching process, called the *dual branching process*. In the binomial case, one way to see this is to think of $\mathcal{T}_{n,p}$ as a random subgraph of the infinite $n$-ary rooted tree $\mathcal{T}_{n,1}$ obtained by including each edge independently with probability $p$, and retaining only the component of the root. For a vertex of $\mathcal{T}_{n,1}$ in generation 1 there are three possibilities: it may (i) be *absent*, i.e., not joined to the root, (ii) *survive*, i.e., be joined to the root and have infinitely many descendents, or (iii) *die out*. The probabilities of these events are $1 - p$, $p\varrho$ and $p(1 - \varrho)$, respectively. Let $\mathcal{D}$ denote the event that the process $\mathcal{T}_{n,p}$ dies out, i.e., the total population is finite. Since $\mathcal{D}$ happens if and only if every vertex of $\mathcal{T}_{n,1}$ in generation 1 is absent or dies out, the conditional distribution of $\mathcal{T}_{n,p}$ given $\mathcal{D}$ is the unconditional distribution of $\mathcal{T}_{n,\pi}$, with $\pi = p(1 - \varrho)/(1 - p\varrho)$. Thus the dual of $\mathcal{T}_{n,p}$ is $\mathcal{T}_{n,\pi}$.

Note that when $np - 1 = \varepsilon \to 0$, then

$$1 - n\pi = \frac{1 - p\varrho - np + np\varrho}{1 - p\varrho} \sim np\varrho - (np - 1) - p\varrho \sim \varepsilon$$

as $n \to \infty$. Hence the mean number of offspring in the dual process $\mathcal{T}_{n,\pi}$ is $1 - (1 + o(1))\varepsilon$, and from (2) its expected total size satisfies

$$\mathbb{E}(|\mathcal{T}_{n,\pi}|) \sim \varepsilon^{-1}. \tag{4}$$

Writing $\mathcal{S} = \mathcal{D}^c$ for the event that $\mathcal{T}_{n,p}$ *survives* (is infinite), and $|\mathcal{T}_{n,p}|$ for its total size (number of vertices), it follows that for any integer $L = L(n)$ we have

$$
\begin{aligned}
\mathbb{P}(|\mathcal{T}_{n,p}| \geqslant L) &= \mathbb{P}(\mathcal{S}) + \mathbb{P}(\mathcal{D})\mathbb{P}(|\mathcal{T}_{n,\pi}| \geqslant L) \\
&\leqslant \mathbb{P}(\mathcal{S}) + \mathbb{P}(|\mathcal{T}_{n,\pi}| \geqslant L) \\
&\leqslant (1 + o(1))(2\varepsilon + 1/(\varepsilon L)),
\end{aligned} \tag{5}
$$

with the second inequality following from Markov's inequality.

We shall use one further property of $\mathcal{T}_{n,p}$, which can be proved in a number of simple ways. Suppose, as above, that $\varepsilon = np - 1 \to 0$, and let $M = M(n)$ satisfy $\varepsilon M \to \infty$. Let $w(\mathcal{T})$ denote the *width* of a rooted tree $\mathcal{T}$, i.e., the maximum (supremum) of the sizes of the generations. Then

$$
\mathbb{P}\big(\{w(\mathcal{T}_{n,p}) \geqslant M\} \cap \mathcal{D}\big) = o(\varepsilon). \tag{6}
$$

To see this, consider testing whether the event $\mathcal{W}_M = \{w(\mathcal{T}_{n,p}) \geqslant M\}$ holds by constructing $\mathcal{T}_{n,p}$ generation by generation, stopping at the first (if any) of size at least $M$. If such a generation exists then (since the descendents of each vertex in this generation form independent copies of $\mathcal{T}_{n,p}$), the conditional probability that the process dies out is at most $(1 - \varrho)^M \leqslant e^{-\varrho M} \to 0$. Hence

$$
\mathbb{P}(\mathcal{D} \mid \mathcal{W}_M) = o(1). \tag{7}
$$

Thus

$$
\mathbb{P}(\mathcal{W}_M) \sim \mathbb{P}(\mathcal{S} \cap \mathcal{W}_M) \leqslant \mathbb{P}(\mathcal{S}) \sim 2\varepsilon,
$$

which with (7) gives (6).

# 3   Application to $G_{n,p}$

The binomial branching process is intimately connected to the component exploration process in $G_{n,p}$. Given a vertex $v$ of $G_{n,p}$, let $C_v$ denote the component of $G_{n,p}$ containing $v$, and let $\mathcal{T}_v$ be the random tree obtained by exploring this component by breadth-first search. In other words, starting with $v$, find all its neighbours, $v_1, \ldots, v_\ell$, say, next find all the neighbours of $v_1$ different from the vertices found so far, then the new neighbours of $v_2$, and so on, ending the second stage with the new neighbours of $v_\ell$. The third stage consists of finding all the new neighbours of the vertices found in the second stage, and so on. Eventually we build a tree $\mathcal{T}_v$, which is a spanning tree of $C_v$.

Note that our notation suppresses the fact that the distributions of $\mathcal{T}_v$ and of $C_v$ depend on $n$ and $p$. In the next lemma, as usual, $|H|$ denotes the total number of vertices in a graph $H$.

**Lemma 2.** (i) *For any $n$ and $p$, the random rooted trees $\mathcal{T}_v$ and $\mathcal{T}_{n,p}$ may be coupled so that $\mathcal{T}_v \subseteq \mathcal{T}_{n,p}$.*

(ii) *For any $n$, $k$ and $p$ there is a coupling of the integer-valued random variables $|C_v|$ and $|\mathcal{T}_{n-k,p}|$ so that either $|C_v| \geqslant |\mathcal{T}_{n-k,p}|$ or both are at least $k$.*

*Proof.* For the first statement we simply generate $\mathcal{T}_v$ and $\mathcal{T}_{n,p}$ together, always adding fictitious vertices to the vertex set of $G_{n,p}$ for the branching process to take from, so that in each step a vertex has $n$ potential new neighbours (some fictitious) each of which it is joined to with probability $p$. All the descendants of the fictitious vertices are themselves fictitious.

To prove (ii) we slightly modify the exploration, to couple a tree $\mathcal{T}_v'$ contained within $C_v$ with $\mathcal{T}_{n-k,p}$ such that one of two alternatives holds: either $\mathcal{T}_v' \supseteq \mathcal{T}_{n-k,p}$, or else both $\mathcal{T}_v'$ and $\mathcal{T}_{n-k,p}$ have at least $k$ vertices. Indeed, construct $\mathcal{T}_v'$ exactly as $\mathcal{T}_v$, except that at each step at the start of which we have not yet reached more than $k$ vertices, we test for edges from the current vertex to exactly $n - k$ potential new neighbours. Since $|C_v| \geqslant |\mathcal{T}_v'|$, this coupling gives the result. $\qquad\square$

From now on we take $p = p(n) = (1 + \varepsilon)/n$, where $\varepsilon = \varepsilon(n)$ satisfies (1). We start by using the two couplings described above to give bounds on the expected number of vertices in large components. In both lemmas, $N_{[L,n]}$ denotes the number of vertices of $G_{n,p}$ in components with between $L$ and $n$ vertices (inclusive); $\mathbb{P}_{n,p}$ and $\mathbb{E}_{n,p}$ denote the probability measure and expectation associated to $G_{n,p}$. Note that since all vertices are equivalent, $\mathbb{E}_{n,p}(N_{[L,n]}) = n\mathbb{P}_{n,p}(|C_v| \geqslant L)$ for any fixed vertex $v$ of $G_{n,p}$.

**Lemma 3.** *Suppose that $L = L(n) = o(\varepsilon n)$. Then $\mathbb{P}_{n,p}(|C_v| \geqslant L) \geqslant (2 + o(1))\varepsilon$. Equivalently, $\mathbb{E}_{n,p}(N_{[L,n]}) \geqslant (2 + o(1))\varepsilon n$.*

*Proof.* Taking $k = L$ in Lemma 2(ii),

$$\begin{aligned} \mathbb{P}_{n,p}(|C_v| \geqslant L) &\geqslant \mathbb{P}(|\mathcal{T}_{n-L,p}| \geqslant L) \\ &\geqslant \mathbb{P}(\mathcal{T}_{n-L,p} \text{ survives}) \sim 2\big((n-L)p - 1\big) \sim 2\varepsilon, \end{aligned}$$

where the approximation steps follow from (3) and the assumption on $L$. $\qquad\square$

**Lemma 4.** *Suppose that $L = L(n)$ satisfies $\varepsilon^2 L \to \infty$. Then $\mathbb{E}_{n,p}(N_{[L,n]}) \leqslant (2 + o(1))\varepsilon n$.*

*Proof.* By Lemma 2(i) and (5),

$$\mathbb{P}_{n,p}(|C_v| \geqslant L) \leqslant \mathbb{P}(|\mathcal{T}_{n,p}| \geqslant L) \leqslant (1 + o(1))(2\varepsilon + 1/(\varepsilon L)) \sim 2\varepsilon.$$

$\qquad\square$

Together these lemmas show that the expected number of vertices in components of size at least $n^{2/3}$, say, is asymptotically $2\varepsilon n$. Two tasks remain: to establish concentration, and to show that most vertices in large components are in a single giant component. For the first task, one can simply count tree components. (This is a little messy, but theoretically trivial. The difficulties in the original papers [1, 6] stemmed from the fact that non-tree components had to be counted as well. What is surprising is that here it suffices to count tree components.) Indeed, applying the first and second moment methods to the number $N$ of vertices in tree components of size at most $n^{2/3}/\omega$, where $\omega = \omega(n) \to \infty$ sufficiently slowly, shows that this number is within $o_p(\varepsilon n)$ of $(1 - \varrho)n$, reproving Lemma 4 and (together with Lemma 3) giving the required concentration. See [2] for a version of this argument with a (best possible) $O_p(\sqrt{n/\varepsilon})$ error term. Since the calculations, though requiring no combinatorial ideas, are somewhat lengthy, we take a different approach here.

**Lemma 5.** *Suppose that $L = L(n)$ satisfies $\varepsilon^2 L \to \infty$ and $L = o(\varepsilon n)$. Then*

$$N_{[L,n]}(G_{n,p}) = (2 + o_{\mathrm{p}}(1))\varepsilon n.$$

*Proof.* Let $N = N_{[L,n]}(G_{n,p})$ be the number of vertices of $G_{n,p}$ in components of size at least $L$. From Lemmas 3 and 4 the expectation $\mathbb{E}[N]$ of $N$ satisfies $\mathbb{E}[N] \sim 2\varepsilon n$, so it suffices to show that

$$\mathbb{E}[N^2] \leqslant (4 + o(1))\varepsilon^2 n^2. \tag{8}$$

Fix a vertex $v$ of $G_{n,p}$. Let us reveal a tree $\mathcal{T}'_v$ spanning a subset $C'_v$ of $C_v$ by exploring using breadth-first search as before, except that we stop the exploration if at any point
(i) we have reached $L$ vertices in total, or
(ii) there are $\lceil \varepsilon L \rceil$ vertices that have been reached (found as a new neighbour of an earlier vertex) but not yet explored (tested for new neighbours).

More precisely, we stop as soon as condition (i) or (ii) holds, even if this is partway through revealing a generation of $\mathcal{T}'_v$, or indeed partway through revealing the new neighbours of a vertex. We call a vertex reached but not (fully) explored a *boundary vertex*, and note that there are at most $\lceil \varepsilon L \rceil \leqslant 2\varepsilon L$ boundary vertices. Let $\mathcal{A}$ be the event that we stop for reason (i) or (ii), rather than because we have revealed the whole component:

$$\mathcal{A} = \{ \text{ the exploration stops due to (i) or (ii) holding } \}.$$

Note that if $|C_v| \geqslant L$, then $\mathcal{A}$ holds.

As before, we may couple $\mathcal{T}'_v$ with $\mathcal{T}_{n,p}$ so that $\mathcal{T}'_v \subseteq \mathcal{T}_{n,p}$. Since the boundary vertices correspond to a set of vertices of $\mathcal{T}_{n,p}$ contained in two consecutive generations, if $\mathcal{A}$ holds, then either $|\mathcal{T}_{n,p}| \geqslant L$ or $w(\mathcal{T}_{n,p}) \geqslant \varepsilon L/2$. From (5) and (6) it follows that $\mathbb{P}(\mathcal{A}) \leqslant (2 + o(1))\varepsilon$.

Since all vertices are equivalent and $|C_v| \geqslant L$ implies that $\mathcal{A}$ holds, we have

$$\mathbb{E}[N^2] = n\mathbb{E}[1_{|C_v| \geqslant L} N] \leqslant n\mathbb{E}[1_\mathcal{A} N] = n\mathbb{P}(\mathcal{A})\mathbb{E}[N \mid \mathcal{A}] \leqslant (2 + o(1))\varepsilon n\mathbb{E}[N \mid \mathcal{A}]. \tag{9}$$

Suppose that $\mathcal{A}$ does hold. Given any vertex $w \notin C'_v$, we explore from $w$ as usual, but within $G' = G_{n,p} \setminus V(C'_v)$, coupling the resulting tree $\mathcal{T}'_w$ with $\mathcal{T}_{n,p}$ so that $\mathcal{T}'_w \subseteq \mathcal{T}_{n,p}$. Let $C'_w$ be the component of $w$ in $G'$, so $C'_w$ is spanned by $\mathcal{T}'_w$. Let $\mathcal{S}$ be the event that (this final copy of) $\mathcal{T}_{n,p}$ is infinite, and let $\mathcal{D} = \mathcal{S}^c$. Note that $C'_w \subseteq C_w$, and that the two are equal unless there is an edge from $C'_w$ to some boundary vertex. Since there are at most $2\varepsilon L$ boundary vertices, this last event has conditional probability at most $2\varepsilon L|C'_w|p \leqslant 3\varepsilon L|C'_w|/n$, say. Since $|C'_w| \leqslant |\mathcal{T}_{n,p}|$, it follows that

$$\begin{aligned}
\mathbb{P}(|C_w| \geqslant L \mid \mathcal{A}) &\leqslant \mathbb{P}(\mathcal{S}) + \mathbb{P}(\mathcal{D})\mathbb{P}(|C'_w| \geqslant L \mid \mathcal{D}) + 3\mathbb{P}(\mathcal{D})\varepsilon Ln^{-1}\mathbb{E}[|C'_w| \mid \mathcal{D}] \\
&\leqslant \mathbb{P}(\mathcal{S}) + \mathbb{P}(|\mathcal{T}_{n,p}| \geqslant L \mid \mathcal{D}) + 3\varepsilon Ln^{-1}\mathbb{E}[|\mathcal{T}_{n,p}| \mid \mathcal{D}] \\
&\leqslant \mathbb{P}(\mathcal{S}) + (L^{-1} + 3\varepsilon Ln^{-1})\mathbb{E}[|\mathcal{T}_{n,p}| \mid \mathcal{D}],
\end{aligned}$$

by Markov's inequality. Since, by (4), the final expectation above is $\sim \varepsilon^{-1}$, and our assumptions give that both $L^{-1}$ and $3\varepsilon Ln^{-1}$ are $o(\varepsilon^2)$, we see that $\mathbb{P}(|C_w| \geqslant L \mid \mathcal{A}) \leqslant (2 + o(1))\varepsilon$. Hence, recalling that there are at most $L$ vertices in $C'_v$,

$$\mathbb{E}[N \mid \mathcal{A}] \leqslant L + (n - L)\mathbb{P}(|C_w| \geqslant L \mid \mathcal{A}) \leqslant L + (2 + o(1))\varepsilon n \sim 2\varepsilon n.$$

Combined with (9) this gives (8). □

To complete the proof of our main result, it remains only to show that almost all vertices in large components are in a single giant component. For this we use a simple form of the classical sprinkling argument of Erdős and Rényi [4].

*Proof of Theorem 1.* It will be convenient to write $\varepsilon = \omega n^{-1/3}$, with $\omega = \omega(n) \to \infty$ and $\omega = o(n^{1/3})$. Also, let $\omega' \to \infty$ *slowly*, say with $\omega' = o(\log \log \omega)$.

Set $L = \varepsilon n/\omega'$. By Lemma 5 there are in total at most $(2 + o_{\mathrm{p}}(1))\varepsilon n$ vertices in components of size larger than $L$, so $L_1(G_{n,p}) + L_2(G_{n,p}) \leqslant (2 + o_{\mathrm{p}}(1))\varepsilon n$. It remains only to show that

$$L_1(G_{n,p}) \geqslant (2 - o_{\mathrm{p}}(1))\varepsilon n. \tag{10}$$

To see this, set $p_1 = n^{-4/3}$, and define $p_0$ by $p_0 + p_1 - p_0 p_1 = p$, so that if first we choose the edges with probability $p_0$ and then (we sprinkle some more) with probability $p_1$, then the random graph we get is exactly $G_{n,p}$. Since $np_0 - 1 = (1 + o(1))\varepsilon$, for any $\delta > 0$ Lemma 5 shows that with probability $1 - o(1)$ the graph $G_{n,p_0}$ has at least $(2 - \delta)\varepsilon n$ vertices in components of size at least $L$.

Let $U_1, \ldots, U_\ell$ be the vertex sets of the components of $G_{n,p_0}$ of size at least $L$, and let $U$ be their union. The probability that no edge sprinkled with probability $p_1$ joins $U_1$ to $U_j$ is

$$(1 - p_1)^{|U_1||U_j|} \leqslant e^{-p_1 L^2} = \exp\left(-n^{-4/3}\omega^2 n^{4/3}/(\omega')^2\right),$$

so the expected number of vertices of $U$ not contained in the component of $G_{n,p}$ containing $U_1$ is at most

$$\sum_{j=2}^{\ell} \exp\left(-(\omega/\omega')^2\right)|U_j| = o(|U|).$$

Consequently, with probability $1 - o(1)$ all but at most $\delta|U|$ vertices of $U$ are contained within a single component of $G_{n,p}$, in which case $L_1(G_{n,p}) \geqslant (1 - \delta)(2 - \delta)\varepsilon n$. Since $\delta > 0$ was arbitrary, (10) follows, completing the proof. □

To conclude, let us remark that although Theorem 1 is a key result about the phase transition, as discussed in the introduction it is far from the final word on the topic.

# References

[1] B. Bollobás, The evolution of random graphs, *Trans. Amer. Math. Soc.* **286** (1984), 257–274.

[2] B. Bollobás and O. Riordan, Random graphs and branching processes, in *Handbook of large-scale random networks*, Bolyai Soc. Math. Stud **18**, B. Bollobás, R. Kozma and D. Miklós eds (2009), pp. 15–115.

[3] B. Bollobás and O. Riordan, Asymptotic normality of the size of the giant component via a random walk *J. Combinatorial Theory B* **102** (2012), 53–61.

[4] P. Erdős and A. Rényi, On the evolution of random graphs, *Magyar Tud. Akad. Mat. Kutató Int. Közl.* **5** (1960), 17–61.

[5] M. Krivelevich and B. Sudakov, The phase transition in random graphs — a simple proof, preprint (2012) `arXiv:1201.6529v4`

[6] T. Łuczak, Component behavior near the critical point of the random graph process, *Random Structures Algorithms* **1** (1990), 287–310.

[7] A. Nachmias and Y. Peres, Component sizes of the random graph outside the scaling window, *ALEA Lat. Am. J. Probab. Math. Stat.* **3** (2007), 133–142.

[8] B. Pittel and C. Wormald, Counting connected graphs inside-out, *J. Combinatorial Theory B* **93** (2005), 127–172.