**Massachusetts Institute of Technology**

# Giant Component in Random Multipartite Graphs with Given Degree Sequences

David Gamarnik [*]       Sidhant Misra [†]

January 23, 2014

## Abstract

We study the problem of the existence of a giant component in a random multipartite graph. We consider a random multipartite graph with $p$ parts generated according to a given degree sequence $n_i^{\mathbf{d}}(n)$ which denotes the number of vertices in part $i$ of the multipartite graph with degree given by the vector $\mathbf{d}$. We assume that the empirical distribution of the degree sequence converges to a limiting probability distribution. Under certain mild regularity assumptions, we characterize the conditions under which, with high probability, there exists a component of linear size. The characterization involves checking whether the Perron-Frobenius norm of the matrix of means of a certain associated edge-biased distribution is greater than unity. We also specify the size of the giant component when it exists. We use the exploration process of Molloy and Reed to analyze the size of components in the random graph. The main challenges arise due to the multidimensionality of the random processes involved which prevents us from directly applying the techniques from the standard unipartite case. In this paper we use techniques from the theory of multidimensional Galton-Watson processes along with Lyapunov function technique to overcome the challenges.

## 1   Introduction

The problem of the existence of a giant component in random graphs was first studied by Erdös and Rényi. In their classical paper [ER60], they considered a random graph model on $n$ and $m$ edges where each such possible graph is equally likely. They showed that if $m/n > \frac{1}{2} + \epsilon$, with high probability as $n \to \infty$ there exists a component of size linear in $n$ in the random graph and that the size of this component as a fraction of $n$ converges to a given constant.

The degree distribution of the classical Erdös-Rényi random graph has Poisson tails. However in many applications the degree distribution associated with an underlying graph does not satisfy this. For example, many so-called "scale-free" networks exhibit power law distribution of degrees. This motivated the study of random graphs generated according

---

[*]Operations Research Center and Sloan School of Management, MIT, Cambridge, MA, 02139, e-mail: `gamarnik@mit.edu`

[†]Department of Electrical Engineering and Computer Science, MIT, Cambridge, MA, 02139, e-mail: `sidhant@mit.edu`

to a given degree sequence. The giant component problem on a random graph generated according to a given degree sequence was considered by Molloy and Reed [MR95]. They provided conditions on the degree distribution under which a giant component exists with high probability. Further in [MR98], they also showed that the size of the giant component as a fraction of the number of vertices converges in probability to a given positive constant. They used an exploration process to analyze the components of vertices of the random graph to prove their results. Similar results were established by Janson and Luczak in [JL08] using different techniques based on the convergence of empirical distributions of independent random variables. There have been several papers that have proved similar results with similar but different assumptions and tighter error bounds [HM12], [BR12], [Rio12]. Results for the critical phase for random graphs with given degree sequences were derived by Kang and Seierstad in [KS08]. All of these results consider a random graph on $n$ vertices with a given degree sequence where the distribution is uniform among all feasible graphs with the given degree sequence. The degree sequence is then assumed to converge to a probability distribution and the results provide conditions on this probability distribution for which a giant component exists with high probability.

In this paper, we consider random *multipartite* graphs with $p$ parts with given degree distributions. Here $p$ is a fixed positive integer. Each vertex is associated with a degree vector $\mathbf{d}$, where each of its component $d_i, i \in [p]$ dictates the number of neighbors of the vertex in the corresponding part $i$ of the graph. As in previous papers, we assume that the empirical distribution associated with the number of vertices of degree $\mathbf{d}$ converges to a probability distribution. We then pose the problem of finding conditions under which there exists a giant component in the random graph with high probability. Our approach is based on the analysis of the Molloy and Reed exploration process. The major bottleneck is that the exploration process is a multidimensional process and the techniques of Molloy and Reed of directly underestimating the exploration process by a one dimensional random walk does not apply to our case. In order to overcome this difficultly, we construct a linear Lyapunov function based on the Perron-Frobenius theorem, a technique often used in the study of multidimensional branching processes. Then we carefully couple the exploration process with some underestimating process to prove our results The coupling construction is also more involved due to the multidimensionality of the process. This is because in contrast to the unipartite case, there are multiple types of clones (or half-edges) involved in the exploration process, corresponding to which pair of parts of the multipartite graph they belong to. At every step of the exploration process, revealing the neighbor of such a clone leads to the addition of clones of *several* types to the component being currently explored. The particular numbers and types of these newly added clones is also dependent on the kind of clone whose neighbor was revealed. So, the underestimating process needs to be constructed in a way such that it simultaneously underestimates the exploration process for each possible type of clone involved. We do this by choosing the parameters of the underestimating process such that for each type of clone, the vector of additional clones which are added by revealing its neighbor is always component wise smaller than the same vector for the exploration process.

All results regarding giant components typically use a configuration model corresponding to the given degree distribution by splitting vertices into clones and performing a uniform matching of the clones. In the standard unipartite case, at every step of the exploration process all available clones can be treated same. However in the multipartite case, this is not the case. For example, the neighbor of a vertex in part 1 of the graph with degree $\mathbf{d}$ can lie in part $j$ only if $d_j > 0$. Further, this neighbor must also have a degree $\hat{\mathbf{d}}$ such

that $\hat{d}_i > 0$. This poses the issue of the graph breaking down into parts with some of the $p$ parts of the graph getting disconnected from the others. To get past this we make a certain irreducibility assumption which we will carefully state later. This assumption not only addresses the above problem, but also enables us to construct linear Lyapunov functions by using the Perron-Frobenius theorem for irreducible non-negative matrices. We also prove that with the irreducibility assumption, the giant component when it exists is unique and has linearly many vertices in each of the $p$ parts of the graph. In [BR12], Bollobas and Riordan show that the existence and the size of the giant component in the *unipartite* case is closely associated with an *edge-biased* branching process. In this paper, we also construct an analogous edge-biased branching process which is now a multi-type branching process, and prove similar results.

Our study of random multipartite graphs is motivated by the fact that several real world networks naturally demonstrate a multipartite nature. The author-paper network, actor-movie network, the network of company ownership, the financial contagion model, heterogenous social networks, etc. are all multipartite [New01], [BEST04], [Jac08]. Examples of biological networks which exhibit multipartite structure include drug target networks, protein-protein interaction networks and human disease networks [GCV+07], [YGC+07], [MBHG06]. In many cases evidence suggests that explicitly modeling the multipartite structure results in more accurate models and predictions.

Random bipartite graphs ($p = 2$) with given degree distributions were considered by Newmann et. al in [NSW01]. They used generating function heuristics to identify the critical point in the bipartite case. However, they did not provide rigorous proofs of the result. Our result establishes a rigorous proof of this result and we show that in the special case $p = 2$, the conditions we derive is equivalent to theirs.

The rest of the paper is structured as follows. In Section 2, we start by introducing the basic definitions and the notion of a degree distribution for multipartite graphs. In Section 3, we formally state our main results. Section 4 is devoted to the description of the configuration model. In Section 5, we describe the exploration process of Molloy and Reed and the associated distributions that govern the evolution of this process. In Section 6 and Section 7, we prove our main results for the supercritical case, namely when a giant component exists with high probability. In Section 8 we prove a sublinear upper bound on the size of the largest component in the subcritical case.

# 2 Definitions and preliminary concepts

We consider a finite simple undirected graph $\mathcal{G} = (\mathcal{V}, \mathcal{E})$ where $\mathcal{V}$ is the set of vertices and $\mathcal{E}$ is the set of edges. We use the words "vertices" and "nodes" interchangeably. A *path* between two vertices $v_1$ and $v_2$ in $\mathcal{V}$ is a collection of vertices $v_1 = u_1, u_2, \ldots, u_l = v_2$ in $\mathcal{V}$ such that for each $i = 1, 2, \ldots, l - 1$ we have $(u_i, u_{i+1}) \in \mathcal{E}$. A component, or more specifically a connected component of a graph $\mathcal{G}$ is a subgraph $\mathcal{C} \subseteq \mathcal{G}$ such that there is a path between any two vertices in $\mathcal{C}$. A family of random graphs $\{\mathcal{G}_n\}$ on $n$ vertices is said to have a giant component if there exists a positive constant $\epsilon > 0$ such that $\mathbf{P}$(There exists a component $\mathcal{C} \subseteq \mathcal{G}_n$ for which $\frac{|\mathcal{C}|}{n} \geq \epsilon) \to 1$. Subsequently, when a property holds with probability converging to one as $n \to \infty$, we will say that the property hold with high probability or w.h.p. for short.

For any integer $p$, we use $[p]$ to denote the set $\{1, 2, \ldots, p\}$. For any matrix $M \in \mathbb{R}^{m \times n}$, we denote by $\|M\| \triangleq \max_{i,j} |M_{ij}|$, the largest element of the matrix $M$ in absolute value.

It is easy to check that $\| \cdot \|$ is a valid matrix norm. We use $\delta_{ij}$ to denote the Kronecker delta function defined by

$$\delta_{ij} = \begin{cases} 1, & \text{if } i = j, \\ 0, & \text{otherwise.} \end{cases}$$

We denote by $\mathbf{1}$ the all ones vector whose dimension will be clear from context.

The notion of an asymptotic degree distribution was introduced by Molloy and Reed [MR95]. In the standard unipartite case, a degree distribution dictates the fraction of vertices of a given degree. In this section we introduce an analogous notion of an asymptotic degree distribution for random multipartite graphs. We consider a random multipartite graph $\mathcal{G}$ on $n$ vertices with $p$ parts denoted by $G_1, \ldots, G_p$. For any $i \in [p]$ a vertex $v \in G_i$ is associated with a "type" $\mathbf{d} \in \mathbb{Z}_+^p$ which we call the "type" of $v$. This means for each $i = 1, 2, \ldots, p$, the node with type $\mathbf{d}$ has $d(i) \triangleq d_i$ neighbors in $G_i$. A degree distribution describes the fraction of vertices of type $\mathbf{d}$ in $G_i$, $i \in [p]$. We now define an *asymptotic degree distribution* as a sequence of degree distributions which prescribe the number of vertices of type $\mathbf{d}$ in a multipartite graph on $n$ vertices. For a fixed $n$, let $\mathcal{D}(n) \triangleq \left( n_i^{\mathbf{d}}(n), \ i \in [p], \mathbf{d} \in \{0, 1, \ldots, n\}^p \right)$, where $n_i^{\mathbf{d}}(n)$ denotes the number of vertices in $G_i$ of type $\mathbf{d}$. Associated with each $\mathcal{D}(n)$ is a probability distribution $\mathbf{p}(n) = \left( \frac{n_i^{\mathbf{d}}(n)}{n}, \ i \in [p], \mathbf{d} \in \{0, 1, \ldots, n\}^p \right)$ which denotes the fraction of vertices of each type in each part. Accordingly, we write $p_i^{\mathbf{d}}(n) = \frac{n_i^{\mathbf{d}}(n)}{n}$. For any vector degree $\mathbf{d}$ the quantity $\mathbf{1}'\mathbf{d}$ is simply the total degree of the vertex. We define the quantity

$$\omega(n) \triangleq \max\{\mathbf{1}'\mathbf{d} : n_i^{\mathbf{d}}(n) > 0 \text{ for some } i \in [p]\}, \tag{1}$$

which is the maximum degree associated with the degree distribution $\mathcal{D}(n)$. To prove our main results, we need additional assumptions on the degree sequence.

**Assumption 1.** The degree sequence $\{\mathcal{D}(n)\}_{n \in \mathbb{N}}$ satisfies the following conditions:

(a) For each $n \in \mathbb{N}$ there exists a simple graph with the degree distribution prescribed by $\mathcal{D}(n)$, i.e., the degree sequence is a *feasible* degree sequence.

(b) There exists a probability distribution $\mathbf{p} = \left( p_i^{\mathbf{d}}, \ i \in [p], \mathbf{d} \in \mathbb{Z}_+^p \right)$ such that the sequence of probability distributions $\mathbf{p}(n)$ associated with $\mathcal{D}(n)$ converges to the distribution $\mathbf{p}$.

(c) For each $i \in [p]$, $\sum_{\mathbf{d}} \mathbf{1}'\mathbf{d} p_i^{\mathbf{d}}(n) \to \sum_{\mathbf{d}} \mathbf{1}'\mathbf{d} p_i^{\mathbf{d}}$.

(d) For each $i, j \in [p]$ such that $\lambda_i^j \triangleq \sum_{\mathbf{d}} d_j p_i^{\mathbf{d}} = 0$, the corresponding quantity $\lambda_i^j(n) \triangleq \sum_{\mathbf{d}} d_j p_i^{\mathbf{d}}(n) = 0$ for all $n$.

(e) The second moment of the degree distribution given by $\sum_{\mathbf{d}} (\mathbf{1}'\mathbf{d})^2 p_i^{\mathbf{d}}$ exists (is finite) and $\sum_{\mathbf{d}} (\mathbf{1}'\mathbf{d})^2 p_i^{\mathbf{d}}(n) \to \sum_{\mathbf{d}} (\mathbf{1}'\mathbf{d})^2 p_i^{\mathbf{d}}$.

Note that the quantity $\sum_{\mathbf{d}} \mathbf{1}'\mathbf{d} p_i^{\mathbf{d}}(n)$ in condition $(c)$ is simply $\frac{\sum_{v \in \mathcal{G}} deg(v)}{n}$. So this condition implies that the total number of edges is $O(n)$, i.e., the graph is sparse. In condition $(e)$ the quantity $\sum_{\mathbf{d}} (\mathbf{1}'\mathbf{d})^2 p_i^{\mathbf{d}}(n)$ is same as $\frac{\sum_{v \in \mathcal{G}} (deg(v))^2}{n}$. So this condition says that sum of the squares of the degrees is $O(n)$. It follows from condition (c) that $\lambda_i^j < \infty$ and that $\lambda_i^j(n) \to \lambda_i^j$. The quantity $\lambda_i^j$ is asymptotically the fraction of outgoing edges from $G_i$ to $G_j$. For $\mathbf{p}$ to be a valid degree distribution of a multipartite graph, we must have for

4

each $1 \leq i < j \leq p$, $\lambda_i^j = \lambda_j^i$ and for every $n$, we must have $\lambda_i^j(n) = \lambda_j^i(n)$. We have not included this in the above conditions because it follows from condition (a). Condition (d) excludes the case where there are sublinear number of edges between $G_i$ and $G_j$.

There is an alternative way to represent some parts of Assumption 1. For any probability distribution $\mathbf{p}$ on $\mathbb{Z}_+^p$, let $\mathbf{D_p}$ denote the random variable distributed as $\mathbf{p}$. Then (b), (c) and (e) are equivalent to the following.

(b') $\mathbf{D_{p(n)}} \to \mathbf{D_p}$ in distribution.

(c') $\mathbf{E}[\mathbf{1'D_{p(n)}}] \to \mathbf{E}[\mathbf{1'D_p}]$.

(e') $\mathbf{E}[(\mathbf{1'D_{p(n)}})^2] \to \mathbf{E}[(\mathbf{1'D_p})^2]$.

The following preliminary lemmas follow immediately.

**Lemma 1.** *The conditions (b'), (c') and (e') together imply that the random variables $\left\{\mathbf{1'D_{p(n)}}\right\}_{n \in \mathbb{N}}$ and $\left\{\left(\mathbf{1'D_{p(n)}}\right)^2\right\}_{n \in \mathbb{N}}$ are uniformly integrable.*

Then using Lemma 1, we prove the following statement.

**Lemma 2.** *The maximum degree satisfies $\omega(n) = o(\sqrt{n})$.*

*Proof.* For any $\epsilon > 0$, by Lemma 1, there exists $q \in \mathbb{Z}$ such that $\mathbf{E}[(\mathbf{1'D_{p(n)}})^2 \mathbf{1}_{\{\mathbf{1'D}>q\}}] < \epsilon$. Observe that for large enough $n$, we have $\max\{\frac{\omega^2(n)}{n}, \frac{q^2}{n}\} \leq \mathbf{E}[(\mathbf{1'D_{p(n)}})^2 \mathbf{1}_{\{\mathbf{1'D}>q\}}] \leq \epsilon$. Since $\epsilon$ is arbitrary, the proof is complete. $\square$

Let $S \triangleq \{(i,j) \mid \lambda_i^j > 0\}$ and let $N \triangleq |S|$. For each $i \in [p]$, let $S_i \triangleq \{j \in [p] \mid (i,j) \in S\}$.

Note that by condition (a), the set of feasible graphs with the degree distribution is non-empty. The random multipartite graph $\mathcal{G}$ we consider in this paper is drawn uniformly at random among all simple graphs with degree distribution given by $\mathcal{D}(n)$. The asymptotic behavior of $\mathcal{D}(n)$ is captured by the quantities $p_i^{\mathbf{d}}$. The existence of a giant component in $\mathcal{G}$ as $n \to \infty$ is determined by the distribution $\mathbf{p}$.

# 3 Statements of the main results

The neighborhood of a vertex in a random graph with given degree distribution resembles closely a special branching process associated with that degree distribution called the edge-biased branching process. A detailed discussion of this phenomenon and results with strong guarantees for the giant component problem in random unipartite graphs can be found in [BR12] and [Rio12]. The edge biased branching process is defined via the edge biased degree distribution that is associated with the given degree distribution. Intuitively the edge-biased degree distribution can be thought of as the degree distribution of vertices reached at the end point of an edge. Its importance will become clear when we will describe the exploration process in the sections that follow. We say that an edge is of type $(i,j)$ if it connects a vertex in $G_i$ with a vertex in $G_j$. Then, as we will see, the type of the vertex in $G_j$ reached by following a random edge of type $(i,j)$ is $\mathbf{d}$ with probability $\frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}$.

We now introduce the *edge-biased branching process* which we denote by $\mathcal{T}$. Here $\mathcal{T}$ is a multidimensional branching process. The vertices of $\mathcal{T}$ except the root are associated with types $(i,j) \in S$. So other than the root, $\mathcal{T}$ has $N \leq p^2$ types of vertices. The root is assumed to be of a special type which will become clear from the description below. The process starts off with a root vertex $v$. With probability $p_i^{\mathbf{d}}$, the root $v$ gives rise to $d_j$

children of type $(i,j)$ for each $j \in [p]$. To describe the subsequent levels of $\mathcal{T}$ let us consider any vertex with type $(i,j)$. With probability $\frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}$ this vertex gives rise to $(d_m - \delta_{mi})$ children of type $(j,m)$ for each $m \in [p]$. The number of children generated by the vertices of $\mathcal{T}$ is independent for all vertices. For each $n$, we define an edge-biased branching process $\mathcal{T}_n$ which we define in the same way as $\mathcal{T}$ by using the distribution $\mathcal{D}(n)$ instead of $\mathcal{D}$. We will also use the notations $\mathcal{T}(v)$ and $\mathcal{T}_n(v)$ whenever the type of the root node $v$ is specified.

We denote the expected number of children of type $(j,m)$ generated by a vertex of type $(i,j)$ by

$$\mu_{ijjm} \triangleq \sum_{\mathbf{d}} (d_m - \delta_{im}) \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}. \tag{2}$$

It is easy to see that $\mu_{ijjm} \geq 0$. Assumption 1(e) guarantees that $\mu_{ijjm}$ is finite. Note that a vertex of type $(i,j)$ cannot have children of type $(l,m)$ if $j \neq l$. But for convenience we also introduce $\mu_{ijlm} = 0$ when $j \neq l$. By means of a remark we should note that it is also possible to conduct the analysis when we allow the second moments to be infinite (see for example [MR95], [BR12]), but for simplicity, we do not pursue this route in this paper.

Introduce a matrix $M \in \mathbb{R}^N$ defined as follows. Index the rows and columns of the matrix with double indices $(i,j) \in S$. There are $N$ such pairs denoting the $N$ rows and columns of $M$. The entry of $M$ corresponding to row index $(i,j)$ and column index $(l,m)$ is set to be $\mu_{ijlm}$.

**Definition 1.** Let $\mathbf{A} \in \mathbb{R}^{N \times N}$ be a matrix. Define a graph $\mathcal{H}$ on $N$ nodes where for each pair of nodes $i$ and $j$, the directed edge $(i,j)$ exists if and only if $A_{ij} > 0$. Then the matrix $\mathbf{A}$ is said to be *irreducible* if the graph $\mathcal{H}$ is strongly connected, i.e., there exists a directed path in $\mathcal{H}$ between any two nodes in $\mathcal{H}$.

We now state the well known Perron-Frobenius Theorem for non-negative irreducible matrices. This theorem has extensive applications in the study of multidimensional branching processes (see for example [KS66]).

**Theorem 1** (Perron-Frobenius Theorem)**.** *Let* $\mathbf{A}$ *be a non-negative irreducible matrix. Then*

*(a).* $\mathbf{A}$ *has a positive eigenvalue* $\gamma > 0$ *such that any other eigenvalue of* $\mathbf{A}$ *is strictly smaller than* $\gamma$ *in absolute value.*

*(b).* *There exists a left eigenvector* $\mathbf{x}$ *of* $\mathbf{A}$ *that is unique up to scalar multiplication associated with the eigenvalue* $\gamma$ *such that all entries of* $\mathbf{x}$ *are positive.*

We introduce the following additional assumption before we state our main results.

**Assumption 2.** The degree sequence $\{\mathcal{D}(n)\}_{n \in \mathbb{N}}$ satisfies the following conditions.

(a). The matrix $M$ associated with the degree distribution $\mathbf{p}$ is irreducible.

(b). For each $i \in [p]$, $S_i \neq \emptyset$.

Assumption 2 eliminates several degenerate cases. For example consider a degree distribution with $p = 4$, i.e., a 4-partite random graph. Suppose for $i = 1,2$, we have $p_i^{\mathbf{d}}$ is non-zero only when $d_3 = d_4 = 0$, and for $i = 3,4$, $p_i^{\mathbf{d}}$ is non-zero only when $d_1 = d_2 = 0$. In essence this distribution is associated with a random graph which is simply the union of two disjoint bipartite graphs. In particular such a graph may contain more than one giant

component. However this is ruled out under our assumption. Further, our assumption allows us to show that the giant component has linearly many vertices in each of the $p$ parts of the multipartite graph.

Let

$$\eta \triangleq 1 - \sum_{i=1}^{\infty} \mathbf{P}(|\mathcal{T}| = i) = \mathbf{P}(|\mathcal{T}| = \infty). \tag{3}$$

Namely, $\eta$ is the survival probability of the branching process $\mathcal{T}$. We now state our main results.

**Theorem 2.** *Suppose that the Perron Frobenius eigenvalue of $M$ satisfies $\gamma > 1$. Then the following statements hold.*

(a) *The random graph $\mathcal{G}$ has a giant component $C \subseteq \mathcal{G}$ w.h.p. Further, the size of this component $C$ satisfies*

$$\lim_{n \to \infty} \mathbf{P}\left(\eta - \epsilon < \frac{|C|}{n} < \eta + \epsilon\right) = 1, \tag{4}$$

*for any $\epsilon > 0$.*

(b) *All components of $\mathcal{G}$ other than $C$ are of size $O(\log n)$ w.h.p.*

**Theorem 3.** *Suppose that the Perron Frobenius eigenvalue of $M$ satisfies $\gamma < 1$. Then all components of the random graph $\mathcal{G}$ are of size $O(\omega(n)^2 \log n)$ w.h.p.*

The conditions of Theorem 2 where a giant component exists is generally referred to in the literature as the supercritical case and that of Theorem 3 marked by the absence of a giant component is referred to as the subcritical case. The conditions under which giant component exists in random bipartite graphs was derived in [NSW01] using generating function heuristics. We now consider the special case of a bipartite graph and show that the conditions implied by Theorem 2 and Theorem 3 reduce to that in [NSW01]. In this case $p = 2$ and $N = 2$. The type of all vertices $\mathbf{d}$ in $G_1$ are of the form $\mathbf{d} = (0, j)$ and those in $G_2$ are of the form $\mathbf{d} = (k, 0)$. To match the notation in [NSW01], we let $p_1^{\mathbf{d}} = p_j$ when $\mathbf{d} = (0, j)$ and $p_2^{\mathbf{d}} = q_k$ when $\mathbf{d} = (k, 0)$. So $\lambda_1^2 = \lambda_2^1 = \sum_{\mathbf{d}} d_2 p_1^{\mathbf{d}} = \sum_j j p_j = \sum_k k q_k$. Using the definition of $\mu_{1221}$ from equation (2), we get

$$\mu_{1221} = \sum_{\mathbf{d}} (d_1 - \delta_{11}) \frac{d_1 p_2^{\mathbf{d}}}{\lambda_1^2} = \frac{\sum_k k(k-1) q_k}{\lambda_1^2}.$$

Similarly we can compute $\mu_{2112} = \frac{\sum_j j(j-1) p_j}{\lambda_1^2}$. From the definition of $M$,

$$M = \begin{bmatrix} 0 & \mu_{1221} \\ \mu_{2112} & 0 \end{bmatrix}.$$

The Perron-Frobenius norm of $M$ is its spectral radius and is given by $(\mu_{1221})(\mu_{2112})$. So the condition for the existence of a giant component according to Theorem 2 is given by $(\mu_{1221})(\mu_{2112}) - 1 > 0$ which after some algebra reduces to

$$\sum_{j,k} jk(jk - j - k) p_j q_k > 0.$$

This is identical to the condition mentioned in [NSW01]. The rest of the paper is devoted to the proof of Theorem 2 and Theorem 3.

# 4  Configuration Model

The configuration model [Wor78], [Bol85], [BC78] is a convenient tool to study random graphs with given degree distributions. It provides a method to generate a multigraph from the given degree distribution. When conditioned on the event that the graph is simple, the resulting distribution is uniform among all simple graphs with the given degree distribution. We describe below the way to generate a configuration model from a given multipartite degree distribution.

1. For each of the $n_i^{\mathbf{d}}(n)$ vertices in $G_i$ of type $\mathbf{d}$ introduce $d_j$ clones of type $(i, j)$. An ordered pair $(i, j)$ associated with a clone designates that the clones belongs to $G_i$ and has a neighbor in $G_j$. From the discussion following Assumption 1, the number of clones of type $(i, j)$ is same as the number of clones of type $(j, i)$.

2. For each pair $(i, j)$, perform a uniform random matching of the clones of type $(i, j)$ with the clones of type $(j, i)$.

3. Collapse all the clones associated with a certain vertex back into a single vertex. This means all the edges attached with the clones of a vertex are now considered to be attached with the vertex itself.

The following useful lemma allows us to transfer results related to the configuration model to uniformly drawn simple random graphs.

**Lemma 3.** *If the degree sequence $\{\mathcal{D}(n)\}_{n \in \mathbb{N}}$ satisfies Assumption 1, then the probability that the configuration model results in a simple graph is bounded away from zero as $n \to \infty$.*

As a consequence of the above lemma, any statement that holds with high probability for the random configuration model is also true with high probability for the simple random graph model. So we only need to prove Theorem 2 and Theorem 3 for the configuration model.

The proof of Lemma 3 can be obtained easily by using a similar result on directed random graphs proved in [COC13]. The specifics of the proof follow.

*Proof of Lemma 3.* In the configuration model for multipartite graphs that we described, we can classify all clones into two categories. First, the clones of the kind, $(i, i) \in S$ and the clones of the kind $(i, j) \in S$, $i \neq j$. Since the outcome of the matching associated with each of the cases is independent, we can treat them separately for this proof. For the first category, the problem is equivalent to the case of configuration model for standard unipartite graphs. More precisely, for a fixed $i$, we can construct a standard degree distribution $\tilde{\mathcal{D}}(n)$ from $\mathcal{D}(n)$ by taking the $i^{th}$ component of the corresponding vector degrees of the latter. By using Assumptions 1, our proof then follows from previous results for unipartite case.

For the second category, first let us fix $(i, j)$ with $i \neq j$. Construct a degree distribution $\mathcal{D}_1(n) = (n^k(n), \ k \in [n])$ where $n^k(n)$ denotes the number of vertices of degree $k$ by letting $n^k(n) = \sum_{\mathbf{d}} \mathbf{1}\{d(j) = k\} n_i^{\mathbf{d}}$. Construct $\mathcal{D}_2(n)$ similar to $\mathcal{D}_1(n)$ by interchanging $i$ and $j$. We consider a bipartite graph where degree distribution of the vertices in part $i$ is given by $\mathcal{D}_i(n)$ for $i = 1, 2$. We form the corresponding configuration model and perform the usual uniform matching between the clones generated from $\mathcal{D}_1(n)$ with the clones generated from $\mathcal{D}_2(n)$. This exactly mimics the outcome of matching that occurs in our original multipartite configuration model between clones of type $(i, j)$ and $(j, i)$. With this formulation, the problem of controlling number of double edges is very closely related to a similar problem concerning the configuration model for directed random graphs which

was studied in [COC13]. To precisely match their setting, add "dummy" vertices with zero degree to both $\mathcal{D}_1(n)$ and $\mathcal{D}_2(n)$ so that they have exactly $n$ vertices each and then arbitrarily enumerate the vertices in each with indices from $[n]$. From Assumption 1 it can be easily verified that the degree distributions $\mathcal{D}_1(n)$ and $\mathcal{D}_2(n)$ satisfy Condition 4.2 in [COC13]. To switch between our notation and theirs, use $\mathcal{D}_1(n) \to M^{[n]}$ and $\mathcal{D}_2(n) \to D^{[n]}$. Then Theorem 4.3 in [COC13] says that the probability of having no self loops and double edges is bounded away from zero. In particular, observing that self loops are irrelevant in our case, we conclude that $\lim_{n \to \infty} \mathbf{P}(\text{No double edges}) > 0$. Since the number of pairs $(i, j)$ is less than or equal to $p(p-1)$ which is a constant with respect to $n$, the proof is now complete. $\square$

# 5 Exploration Process

In this section we describe the exploration process which was introduced by Molloy and Reed in [MR95] to reveal the component associated with a given vertex in the random graph. We say a clone is of type $(i, j)$ if it belongs to a vertex in $G_i$ and has its neighbor in $G_j$. We say a vertex is of type $(i, \mathbf{d})$ if it belongs to $G_i$ and has degree type $\mathbf{d}$. We start at time $k = 0$. At any point in time $k$ in the exploration process, there are three kinds of clones - 'sleeping' clones, 'active' clones and 'dead' clones. For each $(i, j) \in S$, the number of active clones of type $(i, j)$ at time $k$ are denoted by $A_i^j(k)$ and the total number of active clones at time $k$ is given by $A(k) = \sum_{(i,j) \in S} A_i^j(k)$. Two clones are said to be "siblings" if they belong to the same vertex. The set of sleeping and awake clones are collectively called 'living' clones. We denote by $L_i(k)$ the number of living clones in $G_i$ and $L_i^j(k)$ to be the number of living clones of type $(i, j)$ at time $k$. It follows that $\sum_{j \in [p]} L_i^j(k) = L_i(k)$. If all clones of a vertex are sleeping then the vertex is said to be a sleeping vertex, if all its clones are dead, then the vertex is considered dead, otherwise it is considered to be active. At the beginning of the exploration process all clones (vertices) are sleeping. We denote the number of sleeping vertices in $G_i$ of type $\mathbf{d}$ at time $k$ by $N_i^{\mathbf{d}}(k)$ and let $N_S(k) = \sum_{i,\mathbf{d}} N_i^{\mathbf{d}}(k)$. Thus $N_i^{\mathbf{d}}(0) = n_i^{\mathbf{d}}(n)$ and $N_S(0) = n$. We now describe the exploration process used to reveal the components of the configuration model.

**Exploration Process.**

1. *Initialization*: Pick a vertex uniformly at random from the set of all sleeping vertices and and set the status of all its clones to active.

2. Repeat the following two steps as long as there are active clones:

    (a). Pick a clone uniformly at random from the set of active clones and kill it.

    (b). Reveal the neighbor of the clone by picking uniformly at random one of its candidate neighbors. Kill the neighboring clone and make its siblings active.

3. If there are alive clones left, restart the process by picking an alive clone uniformly at random and setting all its siblings to active, and go back to step 2. If there are no alive clones, the exploration process is complete.

Note that in step 2(b), the candidate neighbors of a clones of type $(i, j)$ are the set of alive clones of type $(j, i)$.

The exploration process enables us to conveniently track the evolution in time of the number of active clones of various types. We denote the change in $A_i^j(k)$ by writing

$$A_i^j(k+1) = A_i^j(k) + Z_i^j(k+1), \quad (i,j) \in S.$$

Define $\mathbf{Z}(k) \triangleq \left( Z_i^j(k), \ (i,j) \in S \right)$ to be the vector of changes in the number of active clones of all types. To describe the probability distribution of the changes $Z_i^j(k+1)$, we consider the following two cases.

*Case 1:* $A(k) > 0$.

Let $E_i^j$ denote the event that in step 2-(a) of the exploration process, the active clone picked was of type $(i,j)$. The probability of this event is $\frac{A_i^j(k)}{A(k)}$. In that case we kill the clone that we chose and the number of active clones of type $(i,j)$ reduces by one. Then we proceed to reveal its neighbor which of type $(j,i)$. One of the following events happen:

(i). $E_a$: the neighbor revealed is an active clone. The probability of the joint event is given by

$$\mathbf{P}(E_i^j \cap E_a) = \begin{cases} \dfrac{A_i^j(k)}{A(k)} \dfrac{A_j^i(k)}{L_j^i(k)} & \text{if } i \neq j, \\[2ex] \dfrac{A_i^i(k)}{A(k)} \dfrac{A_i^i(k)-1}{L_i^i(k)-1} & \text{if } i = j. \end{cases}$$

Such an edge is referred to as a back-edge in [MR95]. The change in active clones of different types in this joint event is as follows.

- If $i \neq j$,

$$Z_i^j(k+1) = Z_j^i(k+1) = -1,$$
$$Z_l^m(k+1) = 0, \quad \text{otherwise .}$$

- If $i = j$,

$$Z_i^i(k+1) = -2,$$
$$Z_l^m(k+1) = 0, \quad \text{otherwise .}$$

(ii). $E_s^{\mathbf{d}}$: The neighbor revealed is a sleeping clone of type $\mathbf{d}$. The probability of this joint event is given by

$$\mathbf{P}(E_i^j \cap E_s^{\mathbf{d}}) = \frac{A_i^j(k)}{A(k)} \frac{d_i N_j^{\mathbf{d}}(k)}{L_j^i(k) - \delta_{ij}}.$$

The sleeping vertex to which the neighbor clone belongs is now active. The change in the number of active clones of different types is governed by the type $\mathbf{d}$ of this new active vertex. The change in active clones of different types in this event are as follows.

- If $i \neq j$,

$$Z_i^j(k+1) = -1,$$
$$Z_j^m(k+1) = d_m - \delta_{im},$$
$$Z_l^m(k+1) = 0, \quad \text{otherwise.}$$

10

- If $i = j$,

$$Z_i^i(k+1) = -2 + d_i,$$
$$Z_i^m(k+1) = d_m, \text{ for } m \neq i,$$
$$Z_l^m(k+1) = 0, \text{ otherwise .}$$

Note that the above events are exhaustive, i.e.,

$$\sum_{i,j \in S} \sum_{\mathbf{d}} \mathbf{P}(E_i^j \cap E_s^{\mathbf{d}}) + \sum_{i,j \in S} \mathbf{P}(E_i^j \cap E_a) = 1.$$

*Case 2:* $A(k) = 0$.

In this case, we choose a sleeping clone at random and make it and all its siblings active. Let $E_i^j$ be the event that the sleeping clone chosen was of type $(i, j)$. Further let $E^{\mathbf{d}}$ be the event that this clone belongs to a vertex of type $(i, \mathbf{d})$. Then we have

$$\mathbf{P}(E_i^j \cap E^{\mathbf{d}}) = \frac{L_i^j(k)}{L(k)} \frac{d_j N_i^{\mathbf{d}}(k)}{L_i^j(k)} = \frac{d_j N_i^{\mathbf{d}}(k)}{L(k)}.$$

In this case the change in the number of active clones of different types is given by

$$Z_i^m(k+1) = d_m, \text{ for } m \in S_i,$$
$$Z_{i'}^{m'}(k+1) = 0, \text{ otherwise.}$$

We emphasize here that there are two ways in which the evolution of the exploration process deviates from that of the edge-biased branching process. First, a back-edge can occur in the exploration process when neighbor of an active clone is revealed to be another active clone. Second, the degree distribution of the exploration process is time dependent. However, close to the beginning of the process, these two events do not have a significant impact. We exploit this fact in the following sections to prove Theorem 2 and 3.

# 6 Supercritical Case

In this section we prove the first part of Theorem 2. To do this we show that the number of active clones in the exploration process grows to a linear size with high probability. Using this fact, we then prove the existence of a giant component. The idea behind the proof is as follows. We start the exploration process described in the previous section at an arbitrary vertex $v \in \mathcal{G}$. At the beginning of the exploration process, i.e. at $k = 0$, we have $N_j^{\mathbf{d}}(0) = np_j^{\mathbf{d}}(n)$ and $L_i^j(0) = n\lambda_i^j(n)$. So, close to the beginning of the exploration, a clone of type $(i, j)$ gives rise to $d_m - \delta_{im}$ clones of type $(j, m)$ with probability close to $\frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_j^i(n)}$ which in turn is close to $\frac{d_i p_j^{\mathbf{d}}}{\lambda_j^i}$ for large enough $n$. If we consider the exploration process in a very small linear time scale, i.e. for $k < \epsilon n$ for small enough $\epsilon$, then the quantities $\frac{d_i N_j^{\mathbf{d}}(k)}{L_j^i(k) - \delta_{ij}}$ remain close to $\frac{d_i p_j^{\mathbf{d}}}{\lambda_j}$ and the quantities $\frac{A_j^i(k)}{L_j^i(k) - \delta_{ij}}$ are negligible. We use this observation to construct a process which underestimates the exploration process in some appropriate sense but whose parameters are time invariant and "close" to the initial degree distribution. We then use this somewhat easier to analyze process to prove our result.

11

We now get into the specific details of the proof. We define a stochastic process $B_i^j(k)$ which we will couple with $A_i^j(k)$ such that $B_i^j(k)$ underestimates $A_i^j(k)$ with probability one. We denote the evolution in time of $B_i^j(k)$ by

$$B_i^j(k+1) = B_i^j(k) + \hat{Z}_i^j(k+1), \quad (i,j) \in S.$$

To define $\hat{Z}_i^j(k+1)$, we choose quantities $\pi_{ji}^{\mathbf{d}}$ satisfying

$$0 \le \pi_{ji}^{\mathbf{d}} < \frac{d_i p_j^{\mathbf{d}}}{\lambda_j^i}, \quad p_j^{\mathbf{d}} > 0, \tag{5}$$

$$\sum_{\mathbf{d}} \pi_{ji}^{\mathbf{d}} = 1 - \gamma, \tag{6}$$

for some $0 < \gamma < 1$ to be chosen later.

We now show that in a small time frame, the parameters associated with the exploration process do not change significantly from their initial values. This is made precise in Lemma 4 and Lemma 5 below. Before that we first introduce some useful notation to describe these parameters for a given $n$ and at a given step $k$ in the exploration process. Let $M(n)$ denote the matrix of means defined analogous to $M$ by replacing $\frac{d_i p_j^{\mathbf{d}}}{\lambda_j^i}$ by $\frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_j^i(n)}$. Also for a fixed $n$, define $M_k(n)$ similarly by replacing $\frac{d_i p_j^{\mathbf{d}}}{\lambda_j^i}$ by $\frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k)-\delta_{ij}}$. Note that $M_0(n) = M(n)$. Also from Assumption 1 it follows that $\frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} \to \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}$ and that $M(n) \to M$.

**Lemma 4.** *Given $\delta > 0$, there exists $\epsilon > 0$ and some integer $\hat{n}$ such that for all $n \ge \hat{n}$ and for all time steps $k \le \epsilon n$ in the exploration process we have $\sum_{\mathbf{d}} \left| \frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k)-\delta_{ij}} - \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j} \right| < \delta$.*

*Proof.* Fix $\epsilon_1 > 0$. From Lemma 1 we have that that random variables $\mathbf{1}'\mathbf{D}_{\mathbf{p}(n)}$ are uniformly integrable. Then there exists $q \in \mathbb{Z}$ such that for all $n$ we have $\sum_{\mathbf{d}} d_i p_j^{\mathbf{d}}(n) \mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}} < \epsilon_1$. Since $0 \le \frac{N_j^{\mathbf{d}}(k)}{n} \le \frac{N_j^{\mathbf{d}}(0)}{n} = p_j^{\mathbf{d}}(n)$, we have $\sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}} \left| d_i p_j^{\mathbf{d}}(n) - d_i \frac{N_j^{\mathbf{d}}(k)}{n} \right| < \epsilon_1$. For each time step $k \le \epsilon n$ in the exploration process we have $\frac{N_j^{\mathbf{d}}(k)}{n} \ge \frac{N_j^{\mathbf{d}}(0)}{n} - \epsilon$. So for small enough $\epsilon$, we can make $\sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}\le q\}} \left| d_i \frac{N_j^{\mathbf{d}}(k)}{n} - d_i p_j^{\mathbf{d}}(n) \right| < \epsilon_1$. Additionally, $L_i^j(k)$ can change by at most two at each step. So $\left| \frac{L_i^j(k)-\delta_{ij}}{n} - \lambda_i^j(n) \right| \le 2\epsilon$. So for small enough $\epsilon$, for every $(i,j) \in S$ we have $\frac{n}{L_i^j(k)-\delta_{ij}} - \frac{1}{\lambda_i^j(n)} < \epsilon_1$. Now we can bound

$$\sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}} \left| \frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k)-\delta_{ij}} - \frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} \right| \tag{7}$$

$$\le \sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}} \left( \left| \frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k)-\delta_{ij}} - \frac{d_i N_j^{\mathbf{d}}(k)}{n\lambda_i^j(n)} \right| + \left| \frac{d_i N_j^{\mathbf{d}}(k)}{n\lambda_i^j(n)} - \frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} \right| \right)$$

$$\le \sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}} \frac{d_i N_j^{\mathbf{d}}(k)}{n} \epsilon_1 + \frac{\epsilon_1}{\lambda_i^j(n)}$$

$$\le \delta/4,$$

12

where the last inequality can be obtained by choosing small enough $\epsilon_1$. Since $q$ is a constant, by choosing small enough $\epsilon$ we can ensure that $\sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}\leq q\}} \left| \frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k)-\delta_{ij}} - \frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} \right| \leq \delta/4$. Additionally from Assumption 1, for large enough $n$ we have $\sum_{\mathbf{d}} \left| \frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} - \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j} \right| < \delta/2$. The lemma follows by combining the above inequalities. $\square$

**Lemma 5.** *Given $\delta > 0$, there exists $\epsilon > 0$ and some integer $\hat{n}$ such that for all $n \geq \hat{n}$ and for all time steps $k \leq \epsilon n$ in the exploration process we have $||M_k(n) - M|| \leq \delta$.*

*Proof.* The argument is very similar to the proof of Lemma 4. Fix $\epsilon_1 > 0$. From Lemma 1 we know that the random variables $(\mathbf{1}'\mathbf{D}_{\mathbf{p}(n)})^2$ are uniformly integrable. It follows that there exists $q \in \mathbb{Z}$ such that for all $n$, we have $\mathbf{E}[(\mathbf{1}'\mathcal{D}(n))^2 \mathbf{1}_{\{(\mathbf{1}'\mathcal{D}(n))>q\}}] \leq \epsilon_1$. From this we can conclude that for all $i, j, m$ we have $\sum_{\mathbf{d}}(d_m - \delta_{im})d_i p_j^{\mathbf{d}}(n)\mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}} \leq \epsilon_1$. Since $\frac{N_j^{\mathbf{d}}(0)}{n} - \epsilon \leq \frac{N_j^{\mathbf{d}}(k)}{n} \leq \frac{N_j^{\mathbf{d}}(0)}{n} = p_j^{\mathbf{d}}(n)$, we have

$$|\sum_{\mathbf{d}}(d_m - \delta_{im})d_i p_j^{\mathbf{d}}(n)\mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}} - \sum_{\mathbf{d}}(d_m - \delta_{im})\frac{d_i N_j^{\mathbf{d}}(n)}{n}\mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}}| \leq \epsilon_1. \qquad (8)$$

Also $L_i^j(k)$ can change by at most $2\epsilon n$. So, for small enough $\epsilon$, by an argument similar to the proof of Lemma 4, we can prove analogous to (7) that

$$\left| \sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}}(d_m - \delta_{im})\frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k) - \delta_{ij}} - \sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}>q\}}(d_m - \delta_{im})\frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} \right| \leq \frac{\delta}{4}. \qquad (9)$$

By choosing $\epsilon$ small enough, we can also ensure

$$\left| \sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}\leq q\}}(d_m - \delta_{im})\frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k) - \delta_{ij}} - \sum_{\mathbf{d}} \mathbf{1}_{\{\mathbf{1}'\mathbf{d}\leq q\}}(d_m - \delta_{im})\frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} \right| \leq \frac{\delta}{4}. \qquad (10)$$

Since $M(n)$ converges to $M$ we can choose $\hat{n}$ such that $||M(n) - M|| \leq \frac{\delta}{2}$. By combining the last two inequalities, the proof is complete. $\square$

**Lemma 6.** *Given any $0 < \gamma < 1$, there exists $\epsilon > 0$, an integer $\hat{n} \in \mathbb{Z}$ and quantities $\pi_{ij}^{\mathbf{d}}$ satisfying (5) and (6) and the following conditions for all $n \geq \hat{n}$:*

*(a) For each time step $k \leq \epsilon n$,*

$$\pi_{ji}^{\mathbf{d}} < \frac{d_i N_j^{\mathbf{d}}(k)}{L_j^i(k) - \delta_{ij}}, \qquad (11)$$

*for each $(i, j) \in S$.*

*(b) The matrix $\hat{M}$ defined analogous to $M$ by replacing $\frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}$ by $\pi_{ji}^{\mathbf{d}}$ in (2) satisfies*

$$||\hat{M} - M|| \leq err(\gamma), \qquad (12)$$

*where $err(\gamma)$ is a term that satisfies $\lim_{\gamma \to 0} err(\gamma) = 0$.*

*Proof.* Choose $q = q(\gamma) \in \mathbb{Z}$ such that $\sum_d \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j} \mathbf{1}_{\{\mathbf{1}'\mathbf{d} > q\}} \leq \gamma/2$. Now choose $\pi_{ji}^{\mathbf{d}}$ satisfying (5) and (6) such that $\pi_{ji}^{\mathbf{d}} = 0$ whenever $\mathbf{1}'\mathbf{d} > q$. Using Lemma 4, we can now choose $\hat{n}$ and $\epsilon$ such that for every $(i,j) \in S$ and $\mathbf{d}$ such that $\mathbf{1}'\mathbf{d} \leq q$, (11) is satisfied for all $n \geq \hat{n}$ and all $k \leq \epsilon n$. The condition in part (a) is thus satisfied by this choice of $\pi_{ji}^{\mathbf{d}}$.

For any $\gamma$, let us denote the choice of $\pi_{ji}^{\mathbf{d}}$ made above by $\pi_{ji}^{\mathbf{d}}(\gamma)$. By construction, whenever $M_{ijlm} = 0$, we also have $\hat{M}_{ijlm} = 0$. Suppose $M_{ijjm} = \sum_{\mathbf{d}}(d_m - \delta_{im})\frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j} > 0$. Also, by construction we have $0 \leq \pi_{ji}^{\mathbf{d}}(\gamma) < \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}$ and that $\pi_{ji}^{\mathbf{d}}(\gamma) \to \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}$ as $\gamma \to 0$. Let $X_\gamma$ be the random variable that takes the value $(d_m - \delta_{im})$ with probability $\pi_{ji}^{\mathbf{d}}(\gamma)$ and 0 with probability $\gamma$. Similarly, let $X$ be the random variable that takes the value $(d_m - \delta_{im})$ with probability $\frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}$. Then, from the above argument have $X_\gamma \to X$ as $\gamma \to 0$ and that the random variable $X$ dominates the random variable $X_\gamma$ for all $\gamma \geq 0$. Note that $X$ is integrable. The proof of part (b) is now complete by using the Dominated Convergence Theorem. $\qquad\square$

Assume that the quantities $\epsilon$ and $\pi_{ij}^{\mathbf{d}}$ have been chosen to satisfy the inequalities (11) and (12). We now consider each of the events that can occur at each step of the exploration process until time $\epsilon n$ and describe the coupling between $Z_i^j(k+1)$ and $\hat{Z}_i^j(k+1)$ in each case.

*Case 1:* $A(k) > 0$.

Suppose the event $E_i^j$ happens. We describe the coupling in case of each of the following two events.

(i). $E_a$: the neighbor revealed is an active clone. In this case we simply mimic the evolution of the number of active clones in the original exploration process. Namely, $\hat{Z}_l^m(k+1) = Z_l^m(k+1)$ for all $l, m$.

(ii). $E_s^{\mathbf{d}}$: The neighbor revealed is a sleeping clone of type $\mathbf{d}$. In this case, we split the event further into two events $E_{s,0}^{\mathbf{d}}$ and $E_{s,1}^{\mathbf{d}}$, that is $E_{s,0}^{\mathbf{d}} \cup E_{s,1}^{\mathbf{d}} = E_s^{\mathbf{d}}$ and $E_{s,0}^{\mathbf{d}} \cap E_{s,1}^{\mathbf{d}} = \emptyset$. In particular,

$$\mathbf{P}(E_{s,0}^{\mathbf{d}}|E_i^j \cap E_s^{\mathbf{d}}) = \frac{\pi_{ji}^{\mathbf{d}}(L_j^i(k) - \delta_{ij})}{d_i N_j^{\mathbf{d}}(k)}$$

$$\mathbf{P}(E_{s,1}^{\mathbf{d}}|E_i^j \cap E_s^{\mathbf{d}}) = 1 - \mathbf{P}(E_{s,0}^{\mathbf{d}}|E_i^j \cap E_s^{\mathbf{d}}).$$

For the above to make sense we must have $\pi_{ji} \leq \frac{d_i N_j^{\mathbf{d}}(k)}{L_j^i(k) - \delta_{ij}}$ which is guaranteed by our choice of $\pi_{ij}^{\mathbf{d}}$. We describe the evolution of $B_i^j(k)$ in each of the two cases.

(a). $E_{s,0}^{\mathbf{d}}$: in this case set $\hat{Z}_l^m(k+1) = Z_l^m(k+1)$ for all $l, m$.

(b). $E_{s,1}^{\mathbf{d}}$: In this case, we mimic the evolution of the active clones of event $E_a$ instead of $E_s^{\mathbf{d}}$. More specifically,

 - If $i \neq j$,

$$\hat{Z}_i^j(k+1) = \hat{Z}_j^i(k+1) = -1,$$
$$\hat{Z}_l^m(k+1) = 0, \quad \text{otherwise .}$$

14

- If $i = j$,

$$\hat{Z}_i^i(k+1) = -2,$$
$$\hat{Z}_l^m(k+1) = 0, \quad \text{otherwise} .$$

*Case 2:* $A(k) = 0$.

Suppose that event $E_i^j \cap E^{\mathbf{d}}$ happens. In this case we split $E^{\mathbf{d}}$ into two disjoint events $E_0^{\mathbf{d}}$ and $E_1^{\mathbf{d}}$ such that

$$\mathbf{P}(E_0^{\mathbf{d}}|E_i^j \cap E^{\mathbf{d}}) = \frac{\pi_{ij}^{\mathbf{d}}(L_j^i(k) - \delta_{ij})}{d_j N_i^{\mathbf{d}}(k)}$$
$$\mathbf{P}(E_1^{\mathbf{d}}|E_i^j \cap E^{\mathbf{d}}) = 1 - \mathbf{P}(E_0^{\mathbf{d}}|E_i^j \cap E^{\mathbf{d}}).$$

Again, the probabilities above are guaranteed to be less than one for time $k \le \epsilon n$ because of the choice of $\pi_{ij}^{\mathbf{d}}$. The change in $B_i^j(k+1)$ in case of each of the above events is defined as follows.

(a) $E_0^{\mathbf{d}}$.
- If $i \ne j$,

$$\hat{Z}_j^i(k+1) = -1,$$
$$\hat{Z}_i^m(k+1) = d_m - \delta_{im},$$
$$\hat{Z}_l^m(k+1) = 0, \quad \text{for } l \ne j.$$

- If $i = j$,

$$\hat{Z}_i^i(k+1) = -2 + d_i,$$
$$\hat{Z}_i^m(k+1) = d_m, \text{ for } m \ne i,$$
$$\hat{Z}_l^m(k+1) = 0, \text{ for } l \ne i.$$

(b) $E_1^{\mathbf{d}}$.
- If $i \ne j$,

$$\hat{Z}_i^j(k+1) = \hat{Z}_j^i(k+1) = -1,$$
$$\hat{Z}_l^m(k+1) = 0, \quad \text{otherwise} .$$

- If $i = j$,

$$\hat{Z}_i^i(k+1) = -2,$$
$$\hat{Z}_l^m(k+1) = 0, \quad \text{otherwise} .$$

This completes the description of the probability distribution of the joint evolution of the processes $A_i^j(k)$ and $B_i^j(k)$.

Intuitively, we are trying to decrease the probability of the cases that actually help in the growth of the component and compensate by increasing the probability of the event which hampers the growth of the component (back-edges). From the description of the the coupling between $Z_i^j(k+1)$ and $\hat{Z}_i^j(k+1)$ it can be seen that for time $k < \epsilon n$, with probability one we have $B_i^j(k) \le A_i^j(k)$.

Our next goal is to show that for some $(i,j) \in S$ the quantity $B_i^j(k)$ grows to a linear size by time $\epsilon n$. Let $H(k) = \sigma(\{A_i^j(r), B_i^j(r), \quad (i,j) \in S, \ 1 \le r \le k\})$ denote the filtration of the joint exploration process till time $k$. Then the expected conditional change in $B_i^j(k)$ can be computed by considering the two cases above. First suppose that at time step $k$ we have $A(k) > 0$, i.e., we are in Case 1. We first assume that $i \ne j$. Note that the only events that affect $\hat{Z}_i^j(k+1)$ are $E_i^j$ and $E_m^i$ for $m \in [p]$. Then,

$$\mathbf{E}[\hat{Z}_i^j(k+1)|H(k)] = \mathbf{P}(E_i^j|H(k)) \ \mathbf{E}[\hat{Z}_i^j(k+1)|H(k), E_i^j] \qquad (13)$$
$$+ \sum_m \mathbf{P}(E_m^i \cap E_a|H(k)) \ \mathbf{E}[\hat{Z}_i^j(k+1)|H(k), E_m^i \cap E_a]$$
$$+ \sum_{m,\mathbf{d}} \mathbf{P}(E_m^i \cap E_{s0}^{\mathbf{d}}|H(k)) \ \mathbf{E}[\hat{Z}_i^j(k+1)|H(k), E_m^i \cap E_{s0}^{\mathbf{d}}]$$
$$+ \sum_{m,\mathbf{d}} \mathbf{P}(E_m^i \cap E_{s1}^{\mathbf{d}}|H(k)) \ \mathbf{E}[\hat{Z}_i^j(k+1)|H(k), E_m^i \cap E_{s1}^{\mathbf{d}}].$$

The event $E_m^i \cap E_a$ affects $\hat{Z}_i^j(k+1)$ only when $m = j$, and in this case, $\hat{Z}_i^j(k+1) = -1$. The same is true for the event $E_m^i \cap E_{s1}^{\mathbf{d}}$. In the event $E_m^i \cap E_{s0}^{\mathbf{d}}$, we have $\hat{Z}_i^j(k+1) = d_j - \delta_{jm}$. Using this, the above expression is

$$= \frac{A_i^j(k)}{A(k)}(-1) \ + \ \frac{A_j^i(k)}{A(k)}\frac{A_i^j(k)}{L_i^j(k)}(-1) + \sum_{m,\mathbf{d}} \frac{A_m^i(k)}{A(k)}\pi_{im}^{\mathbf{d}}(d_j - \delta_{jm})$$
$$+ \sum_{\mathbf{d}} \frac{A_j^i(k)}{A(k)}\left(\frac{d_j N_i^{\mathbf{d}}(k)}{L_i^j(k)} - \pi_{ij}^{\mathbf{d}}\right)(-1)$$
$$= \frac{A_i^j(k)}{A(k)}(-1) + \frac{A_j^i(k)}{A(k)}\frac{A_i^j(k)}{L_i^j(k)}(-1) + \sum_m \frac{A_m^i(k)}{A(k)}\left(\sum_{\mathbf{d}} \pi_{im}^{\mathbf{d}}(d_j - \delta_{jm})\right)$$
$$+ \sum_{\mathbf{d}} \frac{A_j^i(k)}{A(k)}\left(\frac{d_j N_i^{\mathbf{d}}(k)}{L_i^j(k)} - \pi_{ij}^{\mathbf{d}}\right)(-1).$$

$$= \frac{A_i^j(k)}{A(k)}(-1) + \frac{A_j^i(k)}{A(k)}\left(\frac{A_i^j(k)}{L_i^j(k)} + \sum_{\mathbf{d}}\left(\frac{d_j N_i^{\mathbf{d}}(k)}{L_i^j(k)}\right) - \sum_{\mathbf{d}} \pi_{ij}^{\mathbf{d}}\right)(-1)$$
$$+ \sum_m \frac{A_m^i(k)}{A(k)}\left(\sum_{\mathbf{d}} \pi_{im}^{\mathbf{d}}(d_j - \delta_{jm})\right)$$

$$= \frac{A_i^j(k)}{A(k)}(-1) + \frac{A_j^i(k)}{A(k)}(-\gamma) + \sum_m \frac{A_m^i(k)}{A(k)}\left(\sum_{\mathbf{d}} \pi_{im}^{\mathbf{d}}(d_j - \delta_{jm})\right),$$

where the last equality follows from (6). Now suppose that at time $k$ we have $A(k) = 0$,

16

i.e., we are in Case 2. In this case, we can similarly compute

$$\mathbf{E}[\hat{Z}_i^j(k+1)|H(k)] = \mathbf{P}(E_i^j|H(k))\,\mathbf{E}[\hat{Z}_i^j(k+1)|H(k),E_i^j]$$
$$+ \sum_{m,\mathbf{d}} \mathbf{P}(E_i^m \cap E^\mathbf{d} \cap E_0^\mathbf{d}|H(k))\,\mathbf{E}[\hat{Z}_i^j(k+1)|H(k),E_i^m \cap E^\mathbf{d} \cap E_0^\mathbf{d}]$$
$$+ \sum_{m,\mathbf{d}} \mathbf{P}(E_i^m \cap E^\mathbf{d} \cap E_1^\mathbf{d}|H(k))\,\mathbf{E}[\hat{Z}_i^j(k+1)|H(k),E_i^m \cap E^\mathbf{d} \cap E_1^\mathbf{d}].$$

Using the description of the coupling in Case 2, the above expression is

$$= \frac{L_j^i(k)}{L(k)}(-1) + \sum_m \frac{L_i^m(k)}{L(k)} \sum_d \pi_{mi}^\mathbf{d}(d_j - \delta_{jm}) + \sum_\mathbf{d} \frac{L_i^j(k)}{L(k)}\frac{d_j N_i^\mathbf{d}(k)}{L_i^j(k)}\left(1 - \frac{\pi_{ji}^\mathbf{d} L_i^j(k)}{d_j N_i^\mathbf{d}(k)}\right)$$
$$= \frac{L_j^i(k)}{L(k)}(-1) + \frac{L_j^i(k)}{L(k)}(-\gamma) + \sum_m \frac{L_i^m(k)}{L(k)} \sum_d \pi_{mi}^\mathbf{d}(d_j - \delta_{jm}).$$

For the case $i = j$, a similar computation will reveal that we obtain very similar expressions to the case $i \neq j$. We give the expressions below and omit the computation. For Case 1, $A(k) > 0$,

$$\mathbf{E}[\hat{Z}_i^i(k+1)|H(k)] = \frac{A_i^i(k)}{A(k)}(-1) + \frac{A_i^i(k)}{A(k)}(-\gamma) + \sum_m \frac{A_m^i(k)}{A(k)}\left(\sum_\mathbf{d} \pi_{im}^\mathbf{d}(d_i - \delta_{im})\right).$$

and for Case 2, $A(k) = 0$,

$$\mathbf{E}[\hat{Z}_i^i(k+1)|H(k)] = \frac{L_i^i(k)}{L(k)}(-1) + \frac{L_i^i(k)}{L(k)}(-\gamma) + \sum_m \frac{L_i^m(k)}{L(k)} \sum_d \pi_{mi}^\mathbf{d}(d_i - \delta_{im}).$$

Define the vector of expected change $\mathbf{E}[\hat{\mathbf{Z}}(k+1)|H(k)] \triangleq \left(\mathbf{E}[Z_i^j(k+1)|H(k)],\ (i,j) \in S\right)$. Also define $\mathbf{A}(k) = \left(\frac{A_i^j(k)}{A(k)},\ (i,j) \in S\right)$ if $A(k) > 0$ and $\mathbf{A}(k) = \left(\frac{L_i^j(k)}{L(k)},\ (i,j) \in S\right)$ if $A(k) = 0$. Let $Q \in \mathbb{R}^{N \times N}$ be given by

$$Q_{ijji} = 1,\ \text{for } (i,j) \in S,$$
$$Q_{ijlm} = 0,\ \text{otherwise }.$$

Then we can write the expected change of $B_i^j(k)$ compactly as

$$\mathbf{E}[\hat{\mathbf{Z}}(k+1)|H(k)] = \left(\hat{M} - \gamma Q - I\right)\mathbf{A}(k). \tag{14}$$

Fix $\delta > 0$. Let $\gamma$ be small enough such that the function $err(\gamma)$ in (12) satisfies $err(\gamma) \leq \delta$. Using Lemma 6 we can choose $\epsilon$ and $\pi_{ij}^\mathbf{d}$ satisfying (11) and (12). In particular, we have $||\hat{M} - M|| \leq \delta$. For small enough $\delta$, both $M$ and $\hat{M}$ have strictly positive entries in the exact same locations. Since $M$ is irreducible, it follows that $\hat{M}$ is irreducible. The Perron-Frobenius eigenvalue of a matrix which is the spectral norm of the matrix is a continuous function of its entries. For small enough $\delta$, the Perron-Frobenius eigenvalue of $\hat{M}$ is bigger than 1, say $1 + 2\zeta$ for some $\zeta > 0$. Let $\mathbf{z}$ be the corresponding left eigenvector with all positive entries and let $z_m \triangleq \min_{(i,j) \in S} z_i^j$ and $z_M \triangleq \max_{(i,j) \in S} z_i^j$. Define the random

17

process $W(k) \triangleq \sum_{(i,j) \in S} z_i^j B_i^j(k)$. Then setting $\Delta W(k+1) = W(k+1) - W(k)$, from (14) we have

$$\mathbf{E}[\Delta W(k+1)|H(k)] = \mathbf{z}'\mathbf{E}\hat{\mathbf{Z}}(k+1)$$
$$= \mathbf{z}'\left(\hat{M} - I\gamma Q\right)\mathbf{A}(k)$$
$$= 2\zeta\mathbf{z}'\mathbf{A}(k) - \gamma\mathbf{z}'Q\mathbf{A}(k).$$

The first term satisfies $2\zeta z_m \leq 2\zeta\mathbf{z}'\mathbf{A}(k) \leq 2\zeta z_M$. This is because $\mathbf{1}'\mathbf{A}(k) = 1$ and hence $\mathbf{z}'\mathbf{A}(k)$ is a convex combination of the entries of $\mathbf{z}$. By choosing $\gamma$ small enough, we can ensure $\gamma\mathbf{z}'Q\mathbf{A}(k) \leq \zeta z_m$. Let $\kappa = \zeta z_m > 0$. Then, we have

$$\mathbf{E}[\Delta W(k+1)|H(k)] \geq \kappa. \tag{15}$$

We now use a one-sided Hoeffding bound argument to show that with high probability the quantity $W(k)$ grows to a linear size by time $\epsilon n$. Let $X(k+1) = \kappa - \Delta W(k+1)$. Then

$$\mathbf{E}[X(k+1)|H(k)] \leq 0. \tag{16}$$

Also note that $|X(k+1)| \leq c\omega(n)$ almost surely, for some constant $c > 0$.

For any $B > 0$ and for any $-B \leq x \leq B$, it can be verified that

$$e^x \leq \frac{1}{2}\frac{e^B + e^{-B}}{2} + \frac{1}{2}\frac{e^B - e^{-B}}{2}x \leq e^{\frac{B^2}{2}} + \frac{1}{2}\frac{e^B - e^{-B}}{2}x.$$

Using the above, we get for any $t > 0$,

$$\mathbf{E}[e^{tX(k+1)}|H(k)] \leq e^{\frac{t^2c^2\omega^2(n)}{2}} + \frac{1}{2}\frac{e^{tc\omega(n)} - e^{-tc\omega(n)}}{2}\mathbf{E}[X(k+1)|H(k)] \leq e^{\frac{t^2c^2\omega^2(n)}{2}},$$

where the last statement follows from (16). We can now compute

$$\mathbf{E}[e^{t\sum_{k=0}^{\epsilon n-1} X(k+1)}] = \prod_{k=0}^{\epsilon n-1}\mathbf{E}[e^{tX(k+1)}|H(k)] \leq e^{\frac{t^2c^2\omega^2(n)\epsilon n}{2}}.$$

So,

$$\mathbf{P}\left(\sum_{k=0}^{\epsilon n-1} X(k+1) > \epsilon\kappa n/2\right) = \mathbf{P}(e^{t\sum_{k=0}^{\epsilon n-1} X(k+1)-t\epsilon\kappa n/2} > 1) \leq e^{-\frac{t\epsilon\kappa n}{2}+\frac{t^2c^2\omega^2(n)\epsilon n}{2}}.$$

Optimizing over $t$, we get

$$\mathbf{P}\left(\sum_{k=0}^{\epsilon n-1} X(k+1) > \epsilon\kappa n/2\right) \leq e^{-\frac{\kappa^2\epsilon n}{8c^2\omega^2(n)}} = o(1),$$

which follows by using Lemma 2. Substituting the definition of $X(k+1)$,

$$\mathbf{P}\left(W(\epsilon n) < \frac{\kappa\epsilon n}{2}\right) = o(1). \tag{17}$$

Recall that $W(k) = \sum_{(i,j)\in S} z_i^j B_i^j(k) \leq Nz_M \max_{(i,j)\in S} B_i^j(k) \leq Nz_M \max_{(i,j)\in S} A_i^j(k)$. Define $\mu \triangleq \frac{\kappa\epsilon}{2Nz_M}$. Then it follows from (17) that there exists a pair $(i', j')$ such that

$$A_{i'}^{j'}(\epsilon n) > \mu n, \quad \text{w.p} \quad 1 - o(1).$$

18

Using the fact that the number of active clones grows to a linear size we now show that the corresponding component is of linear size. To do this, we continue the exploration process in a modified fashion from time $\epsilon n$ onwards. By this we mean, instead of choosing active clones uniformly at random in step $2(a)$ of the exploration process, we now follow a more specific order in which we choose the active clones and then reveal their neighbors. This is still a valid way of continuing the exploration process. The main technical result required for this purpose is Lemma 7 below.

**Lemma 7.** *Suppose that after $\epsilon n$ steps of the exploration process, we have $A_{i'}^{j'}(\epsilon n) > \mu n$ for some pair $(i', j')$. Then, there exists $\epsilon_1 > \epsilon$ and $\delta_1 > 0$ for which we can continue the exploration process in a modified way by altering the order in which active clones are chosen in step $2(a)$ of the exploration proces such that at time $\epsilon_1 n$, w.h.p. for all $(i, j) \in S$, we have $A_i^j(\epsilon_1 n) > \delta_1 n$.*

The above lemma says that we can get to a point in the exploration process where there are linearly many active clones of *every* type. An immediate consequence of this is the Corollary 1 below. We remark here that Corollary 1 is merely one of the consequences of Lemma 7 an can be proved in a much simpler way. But as we will see later, we need the full power of Lemma 7 to prove Theorem 2-(b).

**Corollary 1.** *Suppose that after $\epsilon n$ steps of the exploration process, we have $A_{i'}^{j'}(\epsilon n) > \mu n$ for some pair $(i', j')$. Then there exists $\delta_2 > 0$ such that w.h.p., the neighbors of the $A_{i'}^{j'}$ clones include at least $\delta_2 n$ vertices in $G'_j$.*

Before proving Lemma 7, we state a well known result. The proof can be obtained by standard large deviation techniques. We omit the proof.

**Lemma 8.** *Fix $m$. Suppose there are there are $n$ objects consisting of $\alpha_i n$ objects of type $i$ for $1 \le i \le m$. Let $\beta > 0$ be a constant that satisfies $\beta < \max_i \alpha_i$. Suppose we pick $\beta n$ objects at random from these $n$ objects without replacement. Then for given $\epsilon' > 0$ there exists $z = z(\epsilon', m)$ such that,*

$$\mathbf{P}\left(\left|\frac{\#\text{objects chosen of type } i}{n} - \alpha_i \beta\right| > \epsilon'\right) < z^n.$$

*Proof of Lemma 7.* The proof relies on the fact that the matrix $M$ is irreducible. If we denote the underlying graph associated with $M$ by $\mathcal{H}$, then $\mathcal{H}$ is strongly connected. We consider the subgraph $\mathcal{T}_{i'}^{j'}$ of $\mathcal{H}$ which is the shortest path tree in $\mathcal{H}$ rooted at the node $(i', j')$. We traverse $\mathcal{T}_{i'}^{j'}$ breadth first. Let $d$ be the depth of $\mathcal{T}_{i'}^{j'}$. We continue the exploration process from this point in $d$ stages $1, 2, \ldots, d$. Stage 1 begins right after time $\epsilon n$. Denote the time at which stage $l$ ends by $\epsilon_l n$. For convenience, we will assume a base stage 0, which includes all events until time $\epsilon n$. For $1 \le l \le d$, let $\mathcal{I}_l$ be the set of nodes $(i, j)$ at depth $l$ in $\mathcal{T}_{i'}^{j'}$. We let $\mathcal{I}_0 = \{(i', j')\}$.

We will prove by induction that for $l = 0, 1, \ldots, d$, there exists $\delta^{(l)} > 0$ such that at the end of stage $l$, we have w.h.p., $A_i^j > \delta^{(l)} n$ for each $(i, j) \in \bigcup_{x=0}^{l} \mathcal{I}_x$. Note that at the end of stage 0 we have w.h.p. $A_{i'}^{j'} > \mu n$. So we can choose $\delta^{(0)} = \mu$ to satisfy the base case of the induction. Suppose $|\mathcal{I}_l| = r$. Stage $l + 1$ consists of $r$ substages, namely $(l + 1, 1), (l + 1, 2), \ldots, (l + 1, r)$ where each substage addresses exactly one $(i, j) \in \mathcal{I}_l$. We start stage $(l + 1, 1)$ by considering any $(i, j) \in \mathcal{I}_l$. We reveal the neighbors of $\alpha \delta^{(l)} n$ clones among the $A_i^j > \delta^{(l)} n$ clones one by one. Here $0 < \alpha < 1$ is a constant that will describe

19

shortly. The evolution of active clones in each of these $\alpha\delta^{(l)}n$ steps is identical to that in the event $E_i^j$ in Case 1 of the original exploration process. Fix any $(j,m) \in \mathcal{I}_{l+1}$. Note that $M_{ijjm} > 0$ by construction of $\mathcal{T}_{i'}^{j'}$. So by making $\epsilon$ and $\epsilon_1, \ldots, \epsilon_l$ smaller if necessary and choosing $\alpha$ small enough, we can conclude using Lemma 5 that for all time steps $k < \epsilon_l n + \alpha\delta^{(l)}n$ we have $||M_k(n) - M|| < \delta$ for any $\delta > 0$. Similarly, by using Lemma 4, we get

$$\sum_{\mathbf{d}} \left( -\frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k) - \delta_{ij}} + \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j} \right) = \frac{A_i^j(k) - \delta_{ij}}{L_i^j(k) - \delta_{ij}} \leq \sum_{\mathbf{d}} \left| \frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k) - \delta_{ij}} - \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j} \right| < \delta. \qquad (18)$$

By referring to the description of the exploration process for the event $E_i^j$ in *Case* 1, the expected change in $Z_j^m(k+1)$ during stage $(l+1, 1)$ can be computed similar to (13) as

$$\mathbf{E}[Z_j^m(k+1)|H(k)] = \frac{A_j^i(k) - \delta_{ij}}{L_i^j(k) - \delta_{ij}}(-\delta_{im}) + \sum_{\mathbf{d}} \frac{d_i N_j^{\mathbf{d}}(k)}{L_i^j(k) - \delta_{ij}}(d_m - \delta_{im})$$

$$= (M_k(n))_{ijjm} - \frac{A_j^i(k) - \delta_{ij}}{L_i^j(k) - \delta_{ij}}(-\delta_{im})$$

$$\overset{(a)}{\geq} M_{ijjm} - 2\delta \overset{(b)}{\geq} \delta,$$

where $(a)$ follows from (18) and $(b)$ can be guaranteed by choosing small enough $\delta$. The above argument can be repeated for each $(j,m) \in \mathcal{I}_{l+1}$. We now have all the ingredients we need to repeat the one-sided Hoeffding inequality argument earlier in this section. We can then conclude that there exists $\delta_j^m > 0$ such that w.h.p. we have at least $\delta_j^m n$ active clones of type $(j,m)$ by the end of stage $(l+1, 1)$. By the same argument, this is also true for all children of $(i,j)$ in $\mathcal{T}_{i'}^{j'}$. Before starting stage $S_{l+1}^2$, we set $\delta^{(l)} = \min\{(1-\alpha)\delta^{(l)}, \delta_{j_1}^m\}$. This makes sure that at every substage of stage $l$ we have at least $\delta^{(l)}n$ clones of each kind that has been considered before. This enables us to use the same argument for all substages of stage $l$. By continuing in this fashion, we can conclude that at the end of stage $l+1$ we have $\delta^{(l+1)}n$ clones of each type $(i,j)$ for each $(i,j) \in \bigcup_{x=1}^{l+1} \mathcal{I}_x$ for appropriately defined $\delta^{(l+1)}$. The proof is now complete by induction. $\square$

*Proof of Corollary 1.* Consider any $j \in [p]$. We will prove that the giant component has linearly many vertices in $G_j$ with high probability.

Let $\mathbf{d}$ be such that $p_j^{\mathbf{d}} > 0$ and let $d_i > 0$ for some $i \in [p]$. This means in the configuration model, each of these type $\mathbf{d}$ vertices have at least one clones of type $(j,i)$. Continue the exploration process as in Lemma 7. For small enough $\epsilon_1$ there are at least $n(p_j^{\mathbf{d}} - \epsilon_1)$ of type $(j,i)$ clones still unused at time $\epsilon_1 n$. From Lemma 7, with high probability we have at least $\delta_1 n$ clones of type $(i,j)$ at this point. Proceed by simply revealing the neighbors of each of these. Form Lemma 8, it follows that with high probability, we will cover at least a constant fraction of these clones which correspond to a linear number of vertices covered. Each of these vertices are in the giant component and the proof is now complete. $\square$

We now prove part(b) of Theorem 2. Part (a) will be proved in the next section. We use the argument by Molloy and Reed, except for the multipartite case, we will need the help of Lemma 7 to complete the argument.

*Proof of Theorem 2 (b).* Consider two vertices $u, v \in \mathcal{G}$. We will upper bound the probability that $u$ lies in the component $C$, which is the component being explored at time $\epsilon n$ and $v$

lies in a component of size bigger than $\beta \log n$ other than $C$. To do so start the exploration process at $u$ and proceed till the time step $\epsilon_1 n$ in the statement of Lemma 7. At this time we are in the midst of revealing the component $C$. But this may not be the component of $u$ because we may have restarted the exploration process using the "Initialization step" at some time between 0 and $\epsilon_1 n$. If it is not the component of $u$, then $u$ does not lie in $C$. So, let us assume that indeed we are exploring the component of $u$. At this point continue the exploration process in a different way by switching to revealing the component of $v$. For $v$ to lie in a component of size greater than $\beta \log n$, the number of active clones in the exploration process associated with the component of $v$ must remain positive for each of the first $\beta \log n$ steps. At each step choices of neighbors are made uniformly at random. Also, from Lemma 7, $C$ has at least $\delta_1 n$ active clones of each type. For the component of $v$ to be distinct from the component of $u$ this choice must be different from any of these active clones of the component of $u$. So it follows that the probability of this event is bounded above by $(1 - \delta_1)^{\beta \log n}$. For large enough $\beta$, this gives

$$\mathbf{P}(C(u) = C, \ C(v) \neq C, \ |C(v)| > \beta \log n) = o(n^{-2}).$$

Using a union bound over all pairs of vertices $u$ and $v$ completes the proof. $\square$

# 7  Size of the Giant Component

In this section we complete the proof of Theorem 2-$(a)$ regarding the size of the giant component. For the unipartite case, the first result regarding the size of the giant component was obtained by Molloy and Reed [MR98] by using Wormald's results [Wor95] on using differential equations for random processes. As with previous results for the unipartite case, we show that the size of the giant component as a fraction of $n$ is concentrated around the survival probability of the edge-biased branching process. We do this in two steps. First we show that the probability that a certain vertex $v$ lies in the giant component is approximately equal to the probability that the edge-biased branching process with $v$ as its root grows to infinity. Linearity of expectation then shows that the expected fraction of vertices in the giant component is equal to this probability. We then prove a concentration result around this expected value to complete the proof of Theorem 2. These statements are proved formally in Lemma 10.

Before we go into the details of the proof, we first prove a lemma which is a very widely used application of Azuma's inequality.

**Lemma 9.** *Let $\mathbf{X} = (X_1, X_2, \ldots, X_t)$ be a vector valued random variable and let $f(\mathbf{X})$ be a function defined on $\mathbf{X}$. Let $\mathcal{F}_k \triangleq \sigma(X_1, \ldots, X_k)$. Assume that*

$$|\mathbf{E}(f(\mathbf{X})|\mathcal{F}_k) - \mathbf{E}(f(\mathbf{X})|\mathcal{F}_{k+1})| \leq c.$$

*almost surely. Then*

$$\mathbf{P}(|f(\mathbf{X}) - \mathbf{E}[f(\mathbf{X})]| > s) \leq 2e^{-\frac{s^2}{2tc^2}}.$$

*Proof.* The proof of this lemma is a standard martingale argument. We include it here for completeness. Define the random variables $Y_0, \ldots, Y_t$ as

$$Y_k = \mathbf{E}(f(\mathbf{X})|\mathcal{F}_k).$$

The sequence $\{Y_k\}$ is a martingale and $|Y_k - Y_{k+1}| \leq c$ almost surely. Also $Y_0 = f(\mathbf{X})$ and $Y_t = \mathbf{E}[f(\mathbf{X})]$. The lemma then follows by applying Azuma's inequality to the martingale sequence $\{Y_k\}$. $\qquad\square$

**Lemma 10.** *Let $\epsilon > 0$ be given. Let $v \in \mathcal{G}$ be chosen uniformly at random. Then for large enough $n$, we have*

$$|\mathbf{P}(v \in C) - \mathbf{P}(|\mathcal{T}| = \infty)| \leq \epsilon.$$

*Proof.* We use a coupling argument similar to that used by Bollobas and Riordan [BR12] where it was used to prove a similar result for "local" properties of random graphs. We couple the exploration process starting at $v$ with the branching process $\mathcal{T}_n(v)$ by trying to replicate the event in the branching process as closely as often as possible. We describe the details below.

The parameters of the distribution associated with $\mathcal{T}_n$ is given by $\frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)}$. In the exploration process, at time step $k$ the corresponding parameters are given by $\frac{d_i N_j^{\mathbf{d}}(k)}{L_j^i(k) - \delta_{ij}}$ (see Section 5). We first show that for each of the first $\beta \log n$ steps of the exploration process, these two quantities are close to each other. The quantity $d_i N_j^{\mathbf{d}}(k)$ is the total number of sleeping clones at time $k$ of type $(j, i)$ in $G_j$ that belong to a vertex of type $\mathbf{d}$. At each step of the exploration process the total number of sleeping clones can change by at most $\omega(n)$. Also $L_i^j(k)$ is the total number of living clones of type $(j, i)$ in $G_j$ and can change by at most two in each step.

Then initially for all $(i, j)$ we have $L_i^j(0) = \Theta(n)$ and until time $\beta \log n$ it remains $\Theta(n)$. Therefore,

$$\sum_{i,j,\mathbf{d}} \left| \frac{d_i N_j^{\mathbf{d}}(k+1)}{L_j^i(k+1) - \delta_{ij}} - \frac{d_i N_j^{\mathbf{d}}(k)}{L_j^i(k) - \delta_{ij}} \right| \leq \sum_{i,j,\mathbf{d}} \left| \frac{d_i N_j^{\mathbf{d}}(k+1) - d_i N_j^{\mathbf{d}}(k)}{L_j^i(k) - \delta_{ij}} \right|$$
$$+ \left| \frac{d_i N_j^{\mathbf{d}}(k+1)}{L_j^i(k) - \delta_{ij}} - \frac{d_i N_j^{\mathbf{d}}(k+1)}{L_j^i(k+1) - \delta_{ij}} \right|.$$

From the explanation above, the first term is $O(\omega(n)/n)$ and the second term is $O(1/n)$. Recall that $\frac{d_i N_j^{\mathbf{d}}(0)}{L_j^i(0)} = \frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)}$. From this we can conclude by using a telescopic sum and triangle inequality that for time index $k \leq \beta \log n$,

$$\sum_{i,j,\mathbf{d}} \left| \frac{d_i N_j^{\mathbf{d}}(k)}{L_j^i(k) - \delta_{ij}} - \frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} \right| = O(k\omega(n)/n) = O(\omega(n)\log n/n).$$

So the total variational distance between the distribution of the exploration process and the branching process at each of the first $\beta \log n$ steps is $O(\omega(n)\log n/n)$. We now describe the coupling between the branching process and the exploration process. For the first time step, note that the root of $\mathcal{T}_n$ has type $(i, \mathbf{d})$ with probability $p_i^{\mathbf{d}}$. We can couple this with the exploration process by letting the vertex awakened in the "Initialization step" of the exploration process to be of type $(i, \mathbf{d})$. Since the two probabilities are the same, this step of the coupling succeeds with probability one. Suppose that we have defined the coupling until time $k < \beta \log n$. To describe the coupling at time step $k + 1$ we need to consider the case of two events. The first is the event when the coupling has succeeded until time $k$, i.e., the two processes are identical. In this case, since the total variational distance between

the parameters of the two processes is $O(\omega(n)\log n/n)$ we perform a maximal coupling, i.e., a coupling which fails with probability equal to the total variational distance. For our purposes, we do not need to describe the coupling at time $k+1$ in the event that the coupling has failed at some previous time step. The probability that the coupling succeeds at each of the first $\beta \log n$ steps is at least $(1 - O(\omega(n)\log n/n))^{\beta \log n} = 1 - O(\omega(n)(\log n)^2/n) = 1 - o(1)$. We have shown that the coupling succeeds till time $\beta \log n$ with high probability. Assume that it indeed succeeds. In that case the component explored thus far is a tree. Therefore, at every step of the exploration process a sleeping vertex is awakened because otherwise landing on an active clone will result in a cycle. This means if the branching process has survived up until this point, the corresponding exploration process has also survived until this time and the component revealed has at least $\beta \log n$ vertices. Hence,

$$\mathbf{P}(|C(v)| > \beta \log n) = \mathbf{P}(|\mathcal{T}_n| > \beta \log n) + o(1).$$

But Theorem 2 ($b$) states that with high probability, there is only one component of size greater than $\beta \log n$, which is the giant component, i.e.,

$$\mathbf{P}(v \in C) = \mathbf{P}(|C(v)| > \beta \log n) + o(1) = \mathbf{P}(|\mathcal{T}_n| > \beta \log n) + o(1).$$

So, for large enough $n$, we have $|\mathbf{P}(v \in C) - \mathbf{P}(|\mathcal{T}_n| > \beta \log n)| \leq \epsilon/2$. The survival probability of the branching process $\mathcal{T}$ is given by

$$\mathbf{P}(|\mathcal{T}| = \infty) = 1 - \sum_{i=1}^{\infty} \mathbf{P}(|\mathcal{T}| = i).$$

Choose $K$ large enough such that $|\mathbf{P}(|\mathcal{T}| \geq K) - \mathbf{P}(|\mathcal{T}| = \infty)| \leq \epsilon/4$. Also, since $\frac{d_i p_j^{\mathbf{d}}(n)}{\lambda_i^j(n)} \to \frac{d_i p_j^{\mathbf{d}}}{\lambda_i^j}$ for all $i, j, \mathbf{d}$, from the theory of branching processes, for large enough $n$,

$$|\mathbf{P}(|\mathcal{T}_n| \geq K) - \mathbf{P}(|\mathcal{T}| \geq K)| \leq \epsilon/4,$$
$$|\mathbf{P}(|\mathcal{T}_n| = \infty) - \mathbf{P}(|\mathcal{T}| = \infty)| \leq \epsilon/2.$$

Since for large enougn $n$, we have $\mathbf{P}(|\mathcal{T}_n| = \infty) \leq \mathbf{P}(|\mathcal{T}_n| > \beta \log n) \leq \mathbf{P}(|\mathcal{T}_n| \geq K)$, the proof follows by combining the above statements. $\qquad\square$

Now what is left is to show that the size of the giant component concentrates around its expected value.

*Proof of Theorem 2 (a) - (size of the giant component).* From the first two parts of Theorem 2, with high probability we can categorize all the vertices of $\mathcal{G}$ into two parts, those which lie in the giant component, and those which lie in a component of size smaller than $\beta \log n$, i.e., in small components. The expected value of the fraction of vertices in small components is $1 - \eta + o(1)$. We will now show that the fraction of vertices in small components concentrates around this mean.

Recall that $cn \triangleq n \sum_{i \in [p], \mathbf{d} \in D} \mathbf{1}'\mathbf{d}\, p_i^{\mathbf{d}}$ is the number of edges in the configuration model. Let us consider the random process where the edges of the configuration model are revealed one by one. Each edge corresponds to a matching between clones. Let $E_i\ 1 \leq i \leq cn$ denote the (random) edges. Let $N_S$ denote the number of vertices in small components, i.e., in components of size smaller than $\beta \log n$. We wish to apply Lemma 9 to obtain the desired concentration result for which we need to bound $|\mathbf{E}[N_S|E_1, \ldots, E_k] - \mathbf{E}[N_S|E_1, \ldots, E_{k+1}]|$.

In the term $\mathbf{E}[N_S|E_1, \ldots, E_{k+1}]$, let $E_{k+1}$ be the edge $(x, y)$. The expectation is taken over all possible outcomes of the rest of the edges with $E_{k+1}$ fixed to be the edge $(x, y)$. In the first term $\mathbf{E}[N_S|E_1, \ldots, E_k]$, after $E_1, \ldots, E_k$ are revealed, the expectation is taken over the rest of of the edges, which are chosen uniformly at random among all possible edges. All outcomes are equally likely. We construct a mapping from each possible outcome to an outcome that has $E_{k+1} = (x, y)$. In particular, if the outcome contains the edge $(x, y)$ we can map it to the corresponding outcome with $E_{k+1} = (x, y)$ by simply cross-switching the positions of $(x, y)$ with the edge that occured at $k + 1$. This does not change the value of $N_S$ because it does not depend on the order in which the matching is revealed. On the other hand, if the outcome does not contain $(x, y)$, then we map it to one of the outcomes with $E_{k+1} = (x, y)$ by switching the two edges connected to the vertices $x$ and $y$. We claim that switching two edges in the configuration model can change $N_S$ by at most $4\beta \log n$. To see why observe that we can split the process of cross-switching two edges into four steps. In the first two steps we delete each of the two edges one by one and in the next two steps we put them back one by one in the switched position. Deleting an edge can increase $N_S$ by at most $2\beta \log n$ and can never reduce $N_S$. Adding an edge can decrease $N_S$ by at most $2\beta \log n$ and can never increase $N_S$. So cross-switching can either increase or decrease $N_S$ by at most $4\beta \log n$. Using this we conclude

$$|\mathbf{E}[N_S|E_1, \ldots, E_k] - \mathbf{E}[N_S|E_1, \ldots, E_{k+1}]| \leq 4\beta \log n.$$

We now apply Lemma 9 to obtain.

$$\mathbf{P}\left(\frac{1}{n}(N_S - (1 - \eta)) > \delta\right) < e^{-\frac{n^2 \delta^2}{8n\beta \log n}} = o(1).$$

Since with high probability, the number of vertices in the giant component is $n - N_S$, the above concentration result completes the proof. $\square$

# 8   Subcritical Case

In this section we prove Theorem 3. The idea of the proof is quite similar to that of the supercritical case. The strategy of the proof is similar to that used in [MR95]. More specifically, we consider the event $E_v$ that a fixed vertex $v$ lies in a component of size greater than $\zeta\omega(n)^2 \log n$ for some $\zeta > 0$. We will show that $\mathbf{P}(E_v) = o(n^{-1})$. Theorem 3 then follows by taking a union bound over $v \in \mathcal{G}$.

Assume that we start the exploration process at the vertex $v$. For $v$ to lie in a component of size greater than $\zeta\omega(n)^2 \log n$ the exploration process must remain positive for at least $\zeta\omega(n)^2 \log n$ time steps, at each step of the exploration process, at most one vertex is new vertex is added to the component being revealed. This means at time $\zeta\omega(n)^2 \log n$ we must have $A\left(\zeta\omega(n)^2 \log n\right) > 0$, where recall that $A(k)$ denotes the total number of active clones at time $k$ of the exploration process.

Let $H(k) = \sigma(\{A_i^j(r), \quad (i, j) \in S, \ 1 \leq r \leq k\})$ denote the filtration of the exploration process till time $k$. We will assume that $A(k) > 0$ for $0 < k \leq \zeta\omega(n)^2 \log n$ and upper bound $\mathbf{P}(A(\zeta\omega(n)^2 \log n) > 0)$. We first compute the expected conditional change in the number of active clones at time $k$ for $0 \leq k \leq \zeta\omega(n)^2 \log n$ by splitting the outcomes into

24

the several possible cases that affects $\hat{Z}_i^j(k+1)$ as in (13).

$$
\begin{aligned}
\mathbf{E}[Z_i^j(k+1)|H(k)] &= \mathbf{P}(E_i^j|H(k))\ \mathbf{E}[Z(k+1)|H(k), E_i^j] \\
&\quad + \sum_{m,\mathbf{d}} \mathbf{P}(E_m^i \cap E_a|H(k))\ \mathbf{E}[Z(k+1)|H(k), E_m^i \cap E_a] \\
&\quad + \mathbf{P}(E_m^i \cap E_s^{\mathbf{d}}|H(k))\ \mathbf{E}[Z(k+1)|H(k), E_m^i \cap E_s^{\mathbf{d}}]
\end{aligned}
$$

$$
\begin{aligned}
&= \frac{A_i^j(k)}{A(k)}(-1) \ + \ \sum_m \frac{A_m^i(k)}{A(k)} \frac{A_i^m(k)}{L_i^m(k)}(-\delta_{mj}) \\
&\quad + \sum_{m,\mathbf{d}} \frac{A_m^i(k)}{A(k)} \frac{d_m N_i^{\mathbf{d}}(k)}{L_i^m(k)}(d_j - \delta_{jm}) \\
&= -\frac{A_i^j(k)}{A(k)} - \frac{A_j^i(k)}{A(k)} \frac{A_i^j(k)}{L_i^j(k)} + \sum_m \frac{A_m^i(k)}{A(k)} \sum_{\mathbf{d}} \frac{d_m N_i^{\mathbf{d}}(k)}{L_i^m(k)}(d_j - \delta_{jm}).
\end{aligned}
$$

We proceed with the proof in a similar fashion to the proof of the supercritical case. Let $\mathbf{E}[\hat{\mathbf{Z}}(k+1)|H(k)] = (\mathbf{E}[Z_i^j(k+1)|H(k)],\ (i,j) \in S)$ and define the vector quantity $\mathbf{A}(k) = \left( \frac{A_i^j(k)}{A(k)},\ (i,j) \in S \right)$. Also define the matrix $Q(k) \in \mathbb{R}^{N \times N}$ where rows and columns are indexed by double indices and for each $(i,j) \in S$, and

$$
Q_{ijji}(k) = -\frac{A_i^j(k)}{L_i^j(k) - \delta_{ij}},
$$
$$
Q_{ijlm}(k) = 0 \ \text{ for } (l,m) \neq (j,i).
$$

Then the expected change in the number of active clones of various types can be compactly written as

$$
\mathbf{E}[\hat{\mathbf{Z}}(k+1)|H(k)] = (M(k) - I + Q(k))\, \mathbf{A}(k).
$$

As the exploration process proceeds, the matrix $M(k)$ changes over time. However for large enough $n$, it follows from Lemma 5 that the difference between $M(k)$ and $M$ is small for $0 \leq k \leq \frac{1}{2}\zeta\omega(n)^2 \log n$. In particular given any $\epsilon > 0$, for large enough $n$, we have $||M(k) - M|| < \epsilon$. Also from Lemma 4 we also have $||Q(k)|| < \epsilon$. Let $\mathbf{z}$ be the Perron-Frobenius eigenvector of $M$. By the assumption in Theorem 3, we have

$$
\mathbf{z}'M = (1-\delta)\mathbf{z}',
$$

for some $0 < \delta < 1$, where $(1-\delta) = \gamma$ is the Perron-Frobenius eigenvalue of $M$. Also let $z_m \triangleq \min_i z_i$ and $z_M \triangleq \max_i z_i$. Define the random process

$$
W(k) \triangleq \sum_i z_i A_i(k)
$$

Then the expected conditional change in $W(k)$ is given by

$$
\begin{aligned}
\mathbf{E}(\Delta W(k+1)|H(k)) &= \mathbf{z}'\mathbf{E}\hat{\mathbf{Z}}(k+1) \\
&= \mathbf{z}' (M(k) - I + Q(k))\, \mathbf{A}(k) \\
&= \mathbf{z}'(M - I)\mathbf{A}(k) + \mathbf{z}'(M(k) - M + Q(k))\mathbf{A}(k) \\
&= (-\delta)\mathbf{z}'\mathbf{A}(k) + \mathbf{z}'(M(k) - M + Q(k))\mathbf{A}(k).
\end{aligned}
$$

We can choose $\epsilon$ small enough such that $\mathbf{z}'(M(k) - M + Q(k)) < \frac{1}{2}\delta\mathbf{z}'$, where the inequality refers to element wise inequality. Thus

$$\mathbf{E}(\Delta W(k)|H(k)) < -\frac{1}{2}\delta z' \mathbf{A}(k) < -\frac{1}{2}\delta z_m \triangleq \kappa.$$

We can now repeat the one-sided Hoeffding bound argument following equation (15) in the supercritical case and obtain the following inequality:

$$\mathbf{P}(|W(\alpha) + \kappa\alpha)| > \delta) \leq 2e^{-\frac{\delta^2}{2\alpha\omega^2(n)}}.$$

Setting $\alpha = \zeta\omega^2(n)\log n$ and $\delta = \frac{1}{2}\kappa\alpha$, we get

$$\mathbf{P}(W(\zeta\omega^2(n)\log n) > 0) \leq 2e^{-\frac{\kappa^2\zeta\log n}{8}} = o(n^{-1}),$$

for large enough $\zeta$. We conclude

$$\mathbf{P}(\mathcal{G} \text{ has a component bigger than } \zeta\omega^2(n)\log n \;) < \sum_{v\in\mathcal{G}}\mathbf{P}(C(v) > \zeta\log n) = o(1).$$

This completes the proof of the theorem.

# References

[BC78]    E. A. Bender and E. R. Canfield, *The asymptotic number of labelled graphs with given degree sequences*, Journal of Combinatorial Theory **24** (1978), 296–307.

[BEST04]  M. Boss, H. Elsinger, M. Summer, and S. Thurner, *Network topology of the interback market*, Quantitative Finance **4** (2004), no. 6.

[Bol85]   B. Bollobás, *Random Graphs*, Academic Press (1985).

[BR12]    B. Bollobás and O. Riordan, *An old approach to the giant component problem*.

[COC13]   N. Chen and M. Olvera-Cravioto, *Directed Random Graphs with Given Degree Distributions*, Arxiv.org **1207.2475** (2013).

[ER60]    P. Erdős and A. Rényi, *On the Evolution of Random Graphs*, Magayr Tud. Akad. Mat. Kutato Int. Kozl **5** (1960), 17–61.

[GCV+07]  K. Goh, M. E. Cusick, D. Valle, B. Childs, M. Vidal, and A. Barabasi, *The human disease network*, PNAS **104** (2007), no. 21.

[HM12]    H. Hatami and M. Molloy, *The scaling window for a random graph with a given degree sequence*, Random Structures and Algorithms **41** (2012), 99 – 123.

[Jac08]   M. O. Jackson, *Social and economic networks*, Princeton University Press, 2008.

[JL08]    S. Janson and M. Luczak, *A new approach to the Giant Component Problem*, Random Structures and Algorithms **37** (2008), no. 2, 197–216.

[KS66]    H. Kesten and B. P. Stigum, *A Limit Theorem for Multidimensional Galton-Watson Processes*, The Annals of Mathematical Statistics **37** (1966), no. 5, 1211 – 1223.

[KS08]    M. Kang and T.G. Seierstad, *The critical phase for random graphs with a given degree sequence*, Combinatorics, Probability and Computing **17** (2008), 67–86.

[MBHG06] J. L. Morrision, R. Breitling, D. J. Higham, and D. R. Gilbert, *A lock-and-key model for protein-protein interactions*, Bioinformatics **22** (2006), no. 16.

[MR95] M. Molloy and B. Reed, *A critical point for Random Graphs with a given degree sequence*, Random Structures and Algorithms **6** (1995), 161–180.

[MR98] _____, *The Size of the Largest Component of a Random Graph on a fixed Degree Sequence*, Combinatorics, Probability and Computing **7** (1998), 295–306.

[New01] M.E.J. Newmann, *The structure of scentific collaboration networks*, Proc. Natl. Acad. Sci. USA **98** (2001).

[NSW01] M.E.J. Newmann, S.H. Strogatz, and D.J. Watts, *Random graphs with arbitrary degree distributions and their applications*, Phys. Rev. E **64** (2001), no. 026118.

[Rio12] O. Riordan, *The phase transition in the configuration model*, Combinatorics, Probability and Computing **21** (2012), no. 265–299.

[Wor78] N. C. Wormald, *Some Problems in the Enumeration of Labelled Graphs*, Ph.D. thesis, Newcastle University, 1978.

[Wor95] _____, *Differential Equations for Random Processes and Random Graphs*, Annals of Applied Probability **5** (1995), 1217–1235.

[YGC$^+$07] M. Yildrim, K. Goh, M. E. Cusick, A. Barabasi, and M. Vidal, *Drug-target network*, Nat Biotechnol **25** (2007).