1 2

3

ON THE CONCENTRATION OF THE MAXIMUM DEGREE IN THE DUPLICATION-DIVERGENCE MODELS*

ALAN FRIEZE[†], KRZYSZTOF TUROWSKI[‡], AND WOJCIECH SZPANKOWSKI[§]

Abstract. We present a rigorous and precise analysis of the maximum degree and the average degree in a dynamic duplication-divergence graph model introduced by Solé, Pastor-Satorras et al. in which the graph grows according to a duplication-divergence mechanism, i.e. by iteratively creating a copy of some node and then randomly alternating the neighborhood of a new node with probability p. This model captures the growth of some real-world processes e.g. biological or social networks. In this paper, we prove that for some 0 the maximum degree and the average degree ofa duplication-divergence graph on t vertices are asymptotically concentrated with high probability

11 around t^p and $\max\{t^{2p-1}, 1\}$, respectively, i.e. they are within at most a polylogarithmic factor from 12 these values with probability at least $1 - t^{-A}$ for any constant A > 0.

13 **Key words.** random graphs, dynamic graphs, duplication-divergence model, degree distribu-14 tion, maximum degree, average degree, large deviation

15 AMS subject classifications. 05C07, 05C80, 68R10

16 **1.** Introduction. Studying properties of random graphs is a popular topic of research in computer science and discrete mathematics since the seminal work of Paul 17Erdős and Alfréd Rényi [8]. This model was studied extensively using various prob-18 abilistic and analytic methods. The research mostly concentrated on a few broad 19 topics: distribution of structural properties of graphs (e.g. the number of edges, de-2021 grees of fixed vertex, maximum degree, diameter), the existence of special subgraphs (e.g. motif counting, longest paths, maximum matching, Hamilton cycles), values of 22 well-known combinatorial parameters (e.g. largest independent set, chromatic num-23 ber), or extremal properties (Ramsey- and Turán-type) – see e.g. surveys of results 24in [2, 10, 17, 31]. 25

26The widening array of application domains ranging from biology to finance to 27social science inspired further directions of research: first, there appeared an idea to bring the models to the real-world data and to study important aspects, such as 28 centrality, degree correlation, community detection, or graph compression [19, 20, 24]. 29 Second, more models of random networks were developed e.g. for inhomogeneous ran-30 dom graphs, geometric random graphs, preferential attachment graphs, or duplication 31 32 graphs [4, 10, 31]. Often, these models were inspired by some generation mechanisms 33 (e.g. rich-get-richer), or properties (e.g. scale-free/power-law property) that were

^{*}Submitted to the editors DATE.

The preliminary, partial, and weaker versions of these results appeared in WG 2020 [11] and CO-COON 2021 [12] conference proceedings.

Funding: This work was funded in part by NSF Center for Science of Information Grant CCF-0939370, and NSF Grants DMS1661063 and CCF-2006440.

Krzysztof Turowski's research was funded in part by the Polish National Science Center 2020/39/D/ST6/00419 grant. For the purpose of Open Access, the author has applied a CC-BY public copyright license to any Author Accepted Manuscript (AAM) version arising from this submission. Wojciech Szpankowski's work was funded in part by the Polish National Science Center 2018/31/B/ST6/01294 grant.

[†]Department of Mathematical Sciences, Carnegie Mellon University, Pittsburgh, PA, USA (alan@random.math.cmu.edu).

[‡]Theoretical Computer Science Department, Jagiellonian University, Kraków, Poland (krzysztof.szymon.turowski@gmail.com).

[§]Center for Science of Information, Department of Computer Science, Purdue University, West Lafayette, IN, USA (spa@cs.purdue.edu).

34 claimed at work for the real-world networks [9].

In particular, since the late 1990s attention turned toward dynamics graphs in which the behavior of networks evolves in time, e.g. when sets of vertices and/or edges are functions of time, which is definitely the case for certain biological (e.g. protein-protein networks) and social networks (e.g. graph of citations).

One of the family of such networks is the so-called duplication models [5, 4]. It was observed that the evolutionary dynamics of protein interaction networks can be described by simple duplication and mutation rules [25, 32]. For example, the main mechanisms in such models are duplication and divergence: when vertices arrive one by one, they are created as copies of some already existing node, chosen uniformly at random (duplication), and then the neighborhood is typically altered randomly according to some predefined rules (divergence).

In this paper we study a particular duplication-divergence model, first introduced 46by Solé, Pastor-Satorras et al. [28]. This model is a promising object of inquiry since 47 it has been shown empirically that its degree distribution, small subgraph (graphlets) 48 counts, and the number of symmetries fit very well the structure of some real-world 49biological and social networks, e.g. protein-protein and citation networks. More 50precisely, there exist heuristics to infer the underlying parameters of the model from various biological networks, which enable us to generate similar graphs in terms of degree distributions, k-hop reachability, closeness, betweenness, and graphlet frequency 53 [15, 22] (see also an alternative method of parameter estimation in [29]). It also turns 54out that this model often outperformed alternative ones in terms of systematically 56 replicating the degree distribution, small subgraph (graphlets) counts, and symmetries of the input networks [6, 27, 29]. This suggests a possible real-world significance for 57 the duplication-divergence model, which further motivates the studies of its structural 58 properties.

However, it is also one of the least understood models, much less so than the Erdős-Rényi or preferential attachment models. At the moment there exist only a handful of precise results related to the behavior of the degree distribution of the graphs generated by this model. Our contribution is a step towards closing this gap. In short, we prove an asymptotic tight concentration of two parameters in duplicationdivergence graphs: maximum degree, and average degree (or, equivalently, the number of edges) around their mean values.

The paper is organized as follows: in Section 2 we define formally the duplication-67 68 divergence model, and we present an overview of the previous results related to the properties of the degree distribution. Then, in Section 3 we introduce our result 69 for the maximum degree, with proof split into three parts: in Subsection 3.1 and 70 Subsection 3.2 we prove upper bounds for the degrees of the earliest and later vertices 7172 arriving in the graph, respectively, and in Subsection 3.3 we give a proof of the lower bound for the degree of the first vertices, which is effectively also the lower degree 73 of the maximum degree. Next, we proceed with Section 4, containing the proofs of 74 the upper and the lower bounds for the average degree (or, equivalently, the total 75 number of edges in the graph), respectively. Finally, we offer some further problems 7677 and hypotheses that stem from our current research.

This work is a substantial extension of two conference papers: the one presented at COCOON 2021 [12] which contained the weaker concentration results for maximum degree only for the case $\frac{1}{2} , and the one presented at WG 2020 [11] which$ contained the weaker claims (proved using different methods) for average degree andfor degrees only of the earliest vertices.



Fig. 1: Graph evolution in the duplication-divergence model: new vertices and their parents are marked as white and black squares, respectively; p-edges and r-edges are denoted by dashed and dotted lines.

2. Model definition and earlier work. Throughout the paper we use standard graph notation from [7], e.g. V(G) denotes the vertex set of a graph G, $\deg_G(s)$ is the degree of node s in G, and we write $\Delta(G)$ and D(G) for the maximum degree and the average degree in G. Let also $N_G(s)$ denote the open neighborhood of s in G. All graphs considered in the paper are simple, i.e. without loops or multiple edges.

Additionally, since we are eventually dealing with a probability space over graphs on t vertices, let G_t denote a random variable representing a graph on t vertices. Finally, since we are dealing with graphs growing sequentially, we assume that the vertices are identified with the natural numbers according to their arrival time. For simplicity, we introduce the notation $\deg_t(s)$ for the random variable denoting the degree of vertex s in G_t . Clearly, $\Delta(G_t)$ and $D(G_t)$ are random variables denoting the maximum degree and the average degree in G_t .

Let us now formally define the duplication-divergence model, denoted DD(t, p, r), introduced by Solé et al. [28, 26]. Let G_{t_0} be some graph on $t_0 \leq t$ vertices, with vertices having distinct labels from 1 to t_0 . Now, for every $i = t_0, t_0 + 1, \ldots, t - 1$ we create G_{i+1} from G_i according to the following rules:

- 99 1. we add a new vertex with label i + 1 to the graph,
- 100 2. we choose a vertex u from G_i uniformly at random and we denote u as 101 parent(i + 1),

3. for every vertex v:

102

- (a) if v is adjacent to u ($v \in N_{G_i}(u)$) in G_i , then add an edge between v and i + 1 with probability p,
- 105 (b) if v is not adjacent to u in G_i ($v \notin N_{G_i}(u)$), then add an edge between 106 v and i + 1 with probability $\frac{r}{i}$. Note that this case also occurs when 107 v = u, since $u \notin N_{G_i}(u)$.
- 108 All edge additions are independent Bernoulli random variables.

109 Since both p and $\frac{r}{i}$ for $i = t_0, \ldots, t-1$ are probabilities, we allow the parameter space to be $p \in [0, 1]$ and $r \in [0, t_0]$. 110

There is indeed a duplication-divergence mechanism at work since we can think of 111 the equivalent set of rules in the form "copy a vertex from G_i uniformly at random", 112"remove its neighbors independently at random with probability 1 - p", and "add 113 edges to all other vertices independently at random with probability $\frac{r}{i}$ 114

Throughout the paper we will refer to the standard Big-O Landau notation, as 115popularized e.g. in [13]. Let us recall its basic notion: f(n) = O(g(n)) for some 116functions f and g such that $\exists_{k>0} \exists_{n_0} \forall_{n>n_0} |f(n)| \leq |k \cdot g(n)|$. Additionally, we will use 117 • $f(n) = \Omega(g(n))$ when g(n) = O(f(n)), 118

119

• $f(n) = \Theta(g(n))$ when both f(n) = O(g(n)) and g(n) = O(f(n)), • f(n) = o(g(n)) when f(n) = O(g(n)) but not f(n) = O(g(n))

•
$$f(n) = o(g(n))$$
 when $f(n) = O(g(n))$ but not $f(n) = O(g(n))$.

Intuitively, $f(n) = \Omega(g(n))$ when $\lim_{n \to \infty} \frac{|f(n)|}{|g(n)|} \in [k_1, k_2]$ for some $0 < k_1 < k_2$. Since 121 in the model both p and r (and the order of initial graph G_{t_0}) are constants, the 122 asymptotic results are given exclusively in terms of t. 123

As it was mentioned earlier, there are only a few rigorous results for the DD(t, p, r)124 model and its special cases. For 0 and <math>r = 0, it was proved in [14] that 125asymptotically there exists a phase transition for the limiting distribution of degree 126frequencies: if $p \leq p^*$, then almost all vertices are isolated, i.e. the number of non-127 isolated vertices in G_t is o(t), and if $p > p^*$, then only a constant fraction of vertices 128 (with an explicit constant) are isolated. Moreover, it was proved that for any k the 129 fraction of vertices of degree k in G_t converges to 0, and therefore there is no limiting 130degree distribution for $p > p^*$. From [21] it is known that the number of vertices of 131 degree one in G_t is $\Omega(\log t)$ but again the precise rate of growth of the number of 132vertices with any fixed degree k > 0 is currently unknown. 133

However, also for the same case in [18, 16] it was shown for $p < \exp(-1)$ that 134the (only) connected component in G_t exhibits a power-law property with the scale 135 parameter γ which is the solution of $3 = \gamma + p^{\gamma-2}$. 136

For the general case, the two main parameters under consideration were the degree 137 of fixed vertices $\deg_t(s)$ and the average degree of G_t defined as 138

139
140
$$D(G_t) = \frac{1}{t} \sum_{s=1}^t \deg_t(s)$$

It was shown in [30] that we can solve the recurrence equation for the expected 141 average degree and obtain 142

THEOREM 2.1. For $t \to \infty$ it holds that 143

144
$$\mathbb{E}[D(G_t)] = \begin{cases} \Theta(1) & \text{if } p < \frac{1}{2} \text{ and } r > 0, \\ \Theta(\ln t) & \text{if } p = \frac{1}{2} \text{ and } r > 0, \\ \Theta(t^{2p-1}) & \text{otherwise.} \end{cases}$$

145

In a similar fashion it was shown that the expected degree of a vertex s is given by 146 the following theorem: 147

148 THEOREM 2.2. For $t \to \infty$, it holds that

$$\mathbb{E}[\deg_t(s)] = \begin{cases} \Theta\left(\log\left(\frac{t}{s}\right)\right) & \text{if } p = 0 \text{ and } r > 0, \\ \Theta\left(\left(\frac{t}{s}\right)^p\right) & \text{if } 0 0, \\ \Theta\left(\sqrt{\frac{t}{s}}\log s\right) & \text{if } p = \frac{1}{2} \text{ and } r > 0, \\ \Theta\left(\left(\frac{t}{s}\right)^p s^{2p-1}\right) & \text{otherwise.} \end{cases}$$

150

151 Clearly, the latter result for the earliest vertices implies that the expected maximum 152 degree is $\Omega(t^p)$ for all 0 .

In fact, in [30] the authors obtained more than just Theorem 2.1 and Theorem 2.2, because they derived the exact formulae for both $\mathbb{E}[D(G_t)]$ and $\mathbb{E}[\deg_t(s)]$ with their very convoluted leading coefficients (depending on s, p, r) together with the asymptotics for $\operatorname{Var}[D(G_t)]$ and $\operatorname{Var}[\deg_t(s)]$.

The natural question then is to show that these random variables are concen-157trated, i.e. whether by moving only some small (e.g. polylogarithmic) factor from the 158mean we could observe the polynomial tail decay. Intuitively, for the later vertices 159we should not expect such a phenomenon: since the parent of a new vertex is drawn 160uniformly, and there are two binomial processes on top of it, we expect the degree dis-161 tribution of $\deg_{\ell}(t)$ rather reflect the whole degree distribution, which for some cases 162we know (and for all other we stipulate, based on simulations) is not concentrated. 163However, as we will see in the next sections for the maximum degree and the average 164 degree we can answer this question in the affirmative. 165

166 **3. Maximum degree.** In this section we present our main result concerning the 167 concentration of the maximum degree $\Delta(G_t)$. We formulate it in the next theorem.

 $\Pr[(1-\alpha)t^p < \Delta(G_t) < (1+\alpha)t^p \log^{2-p^2}(t)] = 1 - O(t^{-A})$

168 THEOREM 3.1. Let $0 . Asymptotically for <math>G_t \sim DD(t, p, r)$

171

for any constants
$$\alpha > 0$$
 and $A > 0$.

172We prove separately a lower bound and a matching (within a polylogarithmic factor) upper bound. The main idea of the upper bound proof, presented in the next 173subsection, is as follows: we first in Definition 3.2 introduce auxiliary deterministic 174sequences $(t_i)_{i=0}^k$ and $(X_{t_i})_{i=0}^k$ such that $t_0 < \ldots < t_{k-1} < t \leq t_k$. Although at first glance the dependency between $(t_i)_{i=0}^k$ and $(X_{t_i})_{i=0}^k$ given in this definition could 175176seem very convoluted, the intuition behind it is very simple: by doing this we can 177 prove with little effort that X_{t_i} grows close to t_i^p , provided that we choose the right 178parameters. Indeed, we show that $X_t \leq (1+\alpha)t^p \log^{2-p^2}(t)$ for any constant $\alpha > 0$. 179This way, we want $(X_{t_i})_{i=0}^k$ to be a good (i.e. holding with high probability) upper 180 bound for $\deg_{t_i}(s)$ for all $i = 0, \ldots, k$ and all $s \leq t_0$ (denoted as *early vertices*), which 181in turn should give us a similar lower bound $\deg_t(s)$ in terms of X_t whp. We proceed 182in two major steps: first, by construction, we have $\deg_{t_0}(s) \leq t_0 = X_{t_0}$, and second, 183we prove a bound on $\deg_{t_{i+1}}(s) - \deg_{t_i}(s)$ that ensures it does not exceed $X_{t_{i+1}} - X_{t_i}$ 184with high probability. The latter part is achieved by providing an adequate upper 185186 bounding of $\deg_{t,t_1}(s) - \deg_{t_i}(s)$ by a sum of independent Bernoulli variables, so the Chernoff bound can be employed - and by applying a telescoping sum we establish 187 that $\deg_t(s) \leq X_t$ with high probability for all $s \leq t_0$. Therefore, we find for early 188 vertices s (i.e. $s \leq t_0$) a Chernoff-type bound on the growth of deg_{\u03c0}(s) over an interval 189 190 of certain length h.

The second part of the proof of our upper bound on the maximum degree is 191 inductive: we prove that with high probability for any vertex $s \in (t_i, t_{i+1}]$ it holds 192that $\deg_t(s) \leq \max_{\tau < t_i} \{\deg_t(\tau)\}$, that is, the *later vertices* (that is, for any $s > t_0$) 193 can have maximum degree only with a negligible probability. This proof can also be 194decomposed into three steps: first, we show that a vertex s on its arrival cannot have a 195degree greater than $(1+\varepsilon)(pX_t+r)$ with high probability, and then it cannot increase 196 between time s and t_{i+1} to exceed $X_{t_{i+1}}$. Finally, to proceed from $\deg_{t_i+1}(s) \leq X_{t_i+1}$ 197whp to $\deg_t(s) \leq X_t$ whp we use exactly the same Chernoff bound as for early vertices. 198To prove the lower bound we follow the steps from the upper bound for the early 199

vertices: we show a respective lower Chernoff-type bound on the growth of $\deg_{\tau}(s)$ 200over an interval of certain length h and we combine it with different (but very similar) 201sequences t_i and X_{t_i} , thus proving that in this case $\deg_{\tau}(s) \ge X_{\tau} - \ln^{1+p}(\tau) + 1$ with 202 high probability for all early vertices (that is, $s \leq t_0$), and that $X_t \geq (1 - \alpha)t^p$ for 203any $\alpha > 0$. 204

Note that the asymmetry between the proofs of both bounds stems from the 205fact that for the lower bound we only needed to find an inequality that holds with 206 207 high probability for a single vertex, whereas for the upper bound we had to prove an inequality that holds with high probability for all vertices $s = 1, \ldots, t$. 208

3.1. Upper bound, early vertices $(s \le t_0)$. We begin with the definitions for 209 two auxiliary sequences that we mentioned earlier: 210

DEFINITION 3.2. For any t and the given coefficients $\phi(t)$, $(\beta_i(t))_{i=0}^{k-1}$ and the sequence of positive jumps $(w_i(t))_{i=0}^{k-1}$ we define the sequences $(t_i)_{i=0}^k$ and $(X_{t_i})_{i=0}^k$ 211212and a number $k(t) \in \mathbb{N}$, also implicitly dependent on t as follows: 213

214

 $t_0 = \phi(t), \qquad t_{i+1} = t_i + w_i(t),$ $X_{t_0} = t_0, \qquad X_{t_{i+1}} = X_{t_i} + \beta_i(t) \frac{w_i(t) X_{t_i}}{t_i},$ 215

$$k$$
 is such that $t_{k-1} < t \le t_k$.

Moreover, to prove the desired bounds it would be ultimately necessary that $\phi(t)$ and 218219all $w_i(t)$ tend to infinity with t. For brevity, from now on we assume the dependency on t as implicit and write ϕ , β_i , and w_i instead of $\phi(t)$, $\beta_i(t)$, $w_i(t)$, respectively. 220

Note that inductively from the definition it follows that if $\beta_i \leq 1$, then $X_{t_i} \leq t_i$ 221 for all i = 0, 1, ..., k. 222

Moreover, observe that we do not need to specify the values of X_{τ} for τ other 223 224than $\{t_0, t_1, \ldots, t_k\}$. In the rest of the paper we will be using precisely these values in 225the proofs, so such a definition is sufficient for our purposes. For reader's convenience we shall assume that for any $\tau \in (t_l, t_{l+1})$ for some $l = 0, 1, \ldots, k-1$ the sequence is 226completed in any way such that $X_{t_l} \leq X_{\tau} \leq X_{t_{l+1}}$. 227

Now we analyze the asymptotic properties of these sequences. We start with a 228 simple lower bound: 229

LEMMA 3.3. Assume $\beta_i \ge p - \frac{p(1-p)}{4 \ln t_i}$ and $w_i \le \frac{t_i}{\ln t_i}$. For $t \to \infty$ we have $X_{t_i} \ge t_i^p$ 230 for all i = 0, 1, ..., k. 231

Proof. Let us define $Y_{\tau} = \tau^p$. By definition we know that $X_{t_0} = t_0 \ge Y_{t_0}$. Now, 232

233 let us assume that $X_{t_i} \ge Y_{t_i}$ holds for some $i \ge 0$. Then we have

234
$$Y_{t_{i+1}} - Y_{t_i} = \left((t_i + w_i)^p - t_i^p \right) = t_i^p \left(\left(1 + \frac{w_i}{t_i} \right)^p - 1 \right)$$

235
236
$$\leq t_i^p \left(\frac{pw_i}{t_i} - \frac{p(1-p)w_i^2}{4t_i^2}\right) \leq t_i^p \frac{w_i}{t_i} \left(p - \frac{p(1-p)}{4\ln t_i}\right).$$

since from Taylor expansion it follows that $(1 + x)^p \leq 1 + px - \frac{p(1-p)x^2}{4}$ for any $p \in [0, 1]$ and any $x \in (0, 1)$. Therefore,

239
240
$$Y_{t_{i+1}} - Y_{t_i} \le Y_{t_i} \frac{w_i}{t_i} \left(p - \frac{p(1-p)}{4 \ln t_i} \right) \le X_{t_i} \frac{\beta_i w_i}{t_i} = X_{t_{i+1}} - X_{t_i},$$

so clearly $X_{t_{i+1}} \ge Y_{t_{i+1}}$ holds as well, which completes the inductive step.

Now we prove an upper bound on X_t .

243 LEMMA 3.4. Assume that $\phi \geq \ln t$, $\beta_i \leq p + \frac{1}{2\ln t_i}$ and $w_i \leq \frac{t_i}{\ln t_i}$. It holds 244 asymptotically as $t \to \infty$ that $X_{t_i} \leq \phi^{1-p} t_i^p \ln t_i$ for all $i = 0, 1, \ldots, k$.

245 Proof. We again proceed by induction with $Y_{\tau} = \phi^{1-p}\tau^p \ln \tau$. Clearly, $X_{t_0} = t_0 \leq Y_{t_0} = t_0 \ln t_0$. Directly from the definition we get

247
$$Y_{t_{i+1}} - X_{t_{i+1}} = Y_{t_{i+1}} - X_{t_i} \left(1 + \frac{\beta_i w_i}{t_i} \right)$$

248
$$\geq \phi^{1-p} t_{i+1}^p \ln t_{i+1} - \phi^{1-p} t_i^p \ln t_i \left(1 + \frac{\beta_i w_i}{t_i} \right)$$

249
$$\geq \phi^{1-p} t_i^p \ln t_i \left(\left(\frac{t_{i+1}}{t_i} \right)^p \left(\frac{\ln t_{i+1}}{\ln t_i} \right) - 1 - \frac{\beta_i w_i}{t_i} \right)$$

250
251
$$= \phi^{1-p} t_i^p \ln t_i \left(\left(1 + \frac{w_i}{t_i} \right)^p \left(1 + \frac{\ln(1+w_i/t_i)}{\ln t_i} \right) - 1 - \frac{\beta_i w_i}{t_i} \right).$$

Now we use the inequalities derived from the respective Taylor expansions: $(1 + 253 \quad x)^p \ge 1 + px - \frac{p(1-p)x^2}{2} \ge 1$ and $\ln(1+x) \ge x - \frac{x^2}{2} \ge 0$, true for any $p \in [0,1]$ and any $x \in (0,1)$. In particular, in our case $x = \frac{w_i}{t_i} \le \frac{1}{\ln \ln t} = o(1)$. Therefore

255
$$Y_{t_{i+1}} - X_{t_{i+1}} \ge \phi^{1-p} t_i^p \ln t_i \left(\frac{(p - \beta_i)w_i}{t_i} + \left(1 + \frac{pw_i}{t_i} \right) \left(\frac{w_i}{t_i \ln t_i} - \frac{w_i^2}{2t_i^2 \ln t_i} \right)$$

256
$$-\frac{p(1-p)w_i^2}{2t_i^2} \left(1 + \frac{w_i}{t_i \ln t_i} - \frac{w_i^2}{2t_i^2 \ln t_i}\right) \right)$$

257
$$\geq \phi^{1-p} t_i^{p-1} \ln t_i \cdot w_i \left(-\frac{1}{2 \ln t_i} + \frac{1}{\ln t_i} - \frac{w_i}{8t_i} \left(1 + \frac{w_i}{t_i \ln t_i} - \frac{w_i^2}{2t_i^2 \ln t_i} \right) \right)$$

258
259
$$\geq \phi^{1-p} t_i^{p-1} \cdot w_i \left(\frac{3}{8} - \frac{1}{8 \ln t_i} \left(\frac{w_i}{t_i} - \frac{w_i^2}{2t_i^2}\right)\right)$$

and for sufficiently large t the last expression is clearly non-negative since $\frac{w_i}{t_i} \leq \frac{1}{\ln t_i} \leq \frac{1}{\ln t_0} \leq \frac{1}{\ln \phi} \leq \frac{1}{\ln \ln t} \to 0$, which completes the proof.

Next, we need some bounds on $\deg_{\tau}(s)$ holding with high probability to match with the sequence X_{τ} . Let us begin with the following estimate:

This manuscript is for review purposes only.

LEMMA 3.5. For any $\phi \leq \tau \leq t$ and any $0 \leq d \leq h$ it is true that 264

$$\Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge d \mid \deg_{\tau}(s)\right] \le \exp\left(d\ln\frac{\exp(1) \cdot h(p\deg_{\tau}(s) + pd + r)}{d\tau}\right)$$

Proof. First, it follows from the definition of the model that $\deg_{\tau+i+1}(s) =$ 267 $\deg_{\tau+i}(s) + I_{\tau+i}$ for i = 0, 1, ..., h - 1 where $I_{\tau+i} \sim Be(q_{\tau+i})$ for some $q_{\tau+i} \in [0, 1]$. 268The probability $q_{\tau+i}$ of adding an edge between s and $\tau+i+1$ is just a sum of 269probabilities of two events: 270

- 1. when parent $(\tau + i + 1) \in N_{G_{\tau+i}}(s)$ holds, i.e. with probability $\frac{\deg_{\tau+i}(s)}{\tau+i}$ (since 271we draw the parent uniformly), we add an edge with probability p – so the 272whole event has probability $\frac{p \deg_{\tau+i}(s)}{\tau+i}$ 273
- 2. when $\operatorname{parent}(\tau + i + 1) \notin N_{G_{\tau+i}}(s)$ holds, i.e. with probability $1 \frac{\operatorname{deg}_{\tau+i}(s)}{\tau+i}$, we add an edge with probability $\frac{r}{\tau+i}$ so the whole event has probability 274275 $\frac{r}{\tau+i}\left(1-\frac{\deg_{\tau+i}(s)}{\tau+i}\right).$

Both events are disjoint, so we obtain $q_{\tau+i} = \frac{p \deg_{\tau+i}(s)+r}{\tau+i} - \frac{r \deg_{\tau+i}(s)}{(\tau+i)^2} \leq \frac{p \deg_{\tau+i}(s)+r}{\tau+i}$. Next, we note that the degree grows by at least d if there is a subsequence of d277278 successes i_1, i_2, \ldots, i_d with only failures between them: 279

280
$$\Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge d \mid \deg_{\tau}(s)\right]$$
281
$$= \sum_{\substack{0 \le i_1 < \dots < i_d < h}} \Pr\left[\bigcup_{j \in \{i_1, \dots, i_d\}} I_{\tau+j} \cup \bigcup_{j \in [0, i_d] \setminus \{i_1, \dots, i_d\}} \neg I_{\tau+j}\right]$$
282
$$= \sum_{\substack{0 \le i_1 < \dots < i_d < h}} \prod_{j \in \{i_1, \dots, i_d\}} \Pr[I_{\tau+j} \mid \deg_{\tau+j}(s)] \prod_{j \in [0, i_d] \setminus \{i_1, \dots, i_d\}} \Pr[\neg I_{\tau+j} \mid \deg_{\tau+j}(s)].$$

Now observe that $\Pr[\neg I_{\tau+j} | \deg_{\tau+j}(s)] \leq 1$ for any j and $\Pr[I_{\tau+i_j} | \deg_{\tau+i_j}(s)] \leq 1$ 284 $\frac{p(\deg_{\tau}(s)+j-1)+r}{\tau+i_{j}}$ for $j=1,2,\ldots,d$ since j-th success occurs after exactly j-1 suc-285cesses, i.e. when the degree of the vertex s is exactly equal to $\deg_{\tau}(s) + j - 1$. Thus 286

287
$$\Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge d \mid \deg_{\tau}(s)\right] \le \sum_{\substack{0 \le i_1 < \dots < i_d < h \ j=1}} \prod_{j=1}^d \frac{p(\deg_{\tau}(s) + j - 1) + r}{\tau + i_j}$$
288
$$\le \binom{h}{d} \max_{0 \le i_1 < \dots < i_d < h} \left\{ \prod_{j=1}^d \frac{p(\deg_{\tau}(s) + j - 1) + r}{\tau + i_j} \right\}.$$

One can easily spot that the maximum occurs in the case when $i_j = j - 1$ for all 290 $j = 1, 2, \ldots, d$. This, coupled with a simple upper bound on the value of the binomial 291 coefficient, leads us to the final result 292

293
$$\Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge d \mid \deg_{\tau}(s)\right] \le \frac{h^d \exp(d)}{d^d} \prod_{j=0}^{d-1} \frac{p(\deg_{\tau}(s)+j)+r}{\tau+j}$$

294
$$\leq \exp\left(d\ln h - d\ln d + d + d\ln \frac{p(\deg_{\tau}(s) + d) + r}{\tau}\right)$$

$$\leq \exp\left(d\ln\frac{\exp(1)\cdot h(p\deg_{\tau}(s)+pd+r)}{d\tau}\right).$$

This manuscript is for review purposes only.

This lemma gives a far better bound than the simple estimation $\deg_{\tau+h}(s) \leq \deg_{\tau}(s) +$ 297h (e.g. used in [12]). However, it is still too coarse to obtain a desired upper bound 298that could be coupled with the sequence X_{τ} . But we can still use it to kickstart the 299Chernoff bound by bounding the probabilities of all Bernoulli variables: 300

LEMMA 3.6. For $\ln^{1+p} t \leq \tau \leq t$, $\varepsilon = \frac{1}{5 \ln \tau}$ with $h \leq \frac{\varepsilon \tau}{p(1+2\varepsilon) \exp(2)}$ it holds for any 301 constant A > 0 that 302

$$\underset{304}{\overset{303}{=}} \Pr\left[\max_{j=0,\dots,h-1}\left\{\frac{p\deg_{\tau+j}(s)+r}{\tau+j}\right\} \ge (1+\varepsilon)\frac{pX_{\tau}+r}{\tau} \ \left| \deg_{\tau}(s) \le X_{\tau} \right] = O(t^{-A}).$$

Proof. Substituting $d = \varepsilon X_{\tau}$ in Lemma 3.5 we get asymptotically as $t \to \infty$ that 305

306
$$\Pr\left[\frac{p \deg_{\tau+h}(s) + r}{\tau + h} \ge (1 + \varepsilon) \frac{p X_{\tau} + r}{\tau} \middle| \deg_{\tau}(s) \le X_{\tau}\right]$$

307
$$\leq \Pr[\deg_{\tau+h}(s) \geq (1+\varepsilon)X_{\tau} \mid \deg_{\tau}(s) \leq X_{\tau}]$$

308
$$\leq \Pr[\deg_{\tau+h}(s) - \deg_{\tau}(s) \geq \varepsilon X_{\tau} \mid \deg_{\tau}(s) \leq X_{\tau}$$

$$\begin{aligned}
& \leq \Pr[\deg_{\tau+h}(s) \ge (1+\varepsilon)X_{\tau} \mid \deg_{\tau}(s) \le X_{\tau} \\
& \leq \Pr[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge \varepsilon X_{\tau} \mid \deg_{\tau}(s) \le \varepsilon \\
& \leq \exp\left(\varepsilon X_{\tau} \ln \frac{\exp(1) \cdot h(pX_{\tau} + p\varepsilon X_{\tau} + r))}{\varepsilon X_{\tau} \cdot \tau}\right)
\end{aligned}$$

310
$$\leq \exp\left(\varepsilon X_{\tau} \ln \frac{\exp(1) \cdot hp(1+2\varepsilon)X_{\tau})}{\varepsilon X_{\tau} \cdot \tau}\right) \leq \exp(-\varepsilon X_{\tau})$$

$$\begin{aligned}
& 311 \\
& 312 \\
& 312 \\
& \leq \exp\left(-\frac{\max\{\ln^{1+p}t,\tau^p\}}{5\ln\tau}\right) \leq \exp\left(-\frac{\ln t\cdot\tau^{p^2/(1+p)}\}}{5\ln\tau}\right) \leq t^{-A-1},
\end{aligned}$$

for any constant A > 0. In the fourth line we applied inequality $r \leq p \varepsilon X_{\tau}$. Moreover, 313 in the last line we used the facts that $X_{\tau} \geq \max\{\phi, \tau^p\}$ and $\max\{a, b\} \geq a^{\gamma} b^{1-\gamma}$ for 314 any a, b > 0 and $\gamma \in [0, 1]$. 315

To complete the proof it is sufficient to use a union bound over all values up to 316 h = O(t).Π 317

318 Let us now proceed with providing a Chernoff-type bound on the growth of the degree of a given early vertex: 319

LEMMA 3.7. Let $1 \le s \le \tau \le t$ such that $\tau \ge \phi = \ln^{1+p} t$. Then for any A > 0 it 320 is true that 321

322
323
$$\Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge \frac{3A(1+\delta)}{\delta^2} \ln t \ \left| \deg_{\tau}(s) \le X_{\tau} \right] = O(t^{-A}),$$

with $\varepsilon = \delta = \frac{1}{5 \ln \tau}$, and $h = \frac{3A\tau \ln t}{\delta^2 (1+\varepsilon)(pX_\tau+r)}$. 324

Proof. Let us first define an event 325

$$\mathcal{D}_{\varepsilon}(\tau,h) = \left[\max_{\substack{j=0,\dots,h-1\\ \tau+j}} \left\{ \frac{p \deg_{\tau+j}(s) + r}{\tau+j} \right\} \ge (1+\varepsilon) \frac{p X_{\tau} + r}{\tau} \left| \deg_{\tau}(s) \le X_{\tau} \right].$$

329
$$\Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge d \mid \deg_{\tau}(s) \le X_{\tau}\right]$$

$$\stackrel{330}{\le} \Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge d \mid \deg_{\tau}(s) \le X_{\tau}, \neg \mathcal{D}_{\varepsilon}(\tau, h)\right] + \Pr[\mathcal{D}_{\varepsilon}(\tau, h)],$$

Let us estimate the probability of the second event. If $h = \frac{3A\tau \ln t}{\delta^2(1+\varepsilon)(pX_\tau+r)}$ and 332 $\varepsilon = \frac{1}{5 \ln \tau}$, then the condition $h \le \frac{\varepsilon \tau}{p(1+\varepsilon) \exp(2) \ln \tau}$ is met since for some constant C > 0333 we have 334

335
$$h \leq \frac{C\tau \ln t}{\delta^2 X_{\tau}} = \frac{C\tau \ln t \cdot 25 \ln^2 \tau}{\max\{\ln^{1+p} t, \tau^p\}} = \frac{C\tau \ln t \cdot 25 \ln^2 \tau}{\ln t \cdot \tau^{p^2/(1+p)}} \leq \frac{\tau}{\ln^2 \tau}$$

$$\leq \frac{\tau}{p \cdot 2 \exp(2) \cdot 5 \ln \tau} \leq \frac{\varepsilon\tau}{p(1+\varepsilon) \exp(2) \ln \tau}$$

$$\leq \frac{\tau}{p \cdot 2 \exp(2) \cdot 5 \ln \tau} \leq \frac{\varepsilon}{p(1+\varepsilon)}$$

and from Lemma 3.6 we obtain that $\Pr[\mathcal{D}_{\varepsilon}(\tau,h)] = O(t^{-A})$. Here we again used the 338 facts that $X_{\tau} \ge \max\{\phi, \tau^p\}$ and $\max\{a, b\} \ge a^{\gamma} b^{1-\gamma}$ for any a, b > 0 and $\gamma \in [0, 1]$. 339

Thus, it is sufficient to bound $\deg_{\tau+h}(s) - \deg_{\tau}(s)$ with high probability when 340 $\mathcal{D}_{\varepsilon}(\tau, h)$ does not hold, that is, when for all $i = 1, \ldots, h$ it is true that 341

$$\frac{\deg_{\tau+i}(s)}{\tau+i} < (1+\varepsilon)\frac{X_{\tau}}{\tau}.$$

It follows that $I_{\tau+i} = \deg_{\tau+i+1}(s) - \deg_{\tau+i}(s)$ is stochastically dominated by inde-344 pendent random variables $I_{\tau+i}^* \sim Be\left((1+\varepsilon)\frac{pX_{\tau}+r}{\tau}\right)$ for any $i = 0, 1, \ldots, h-1$ – since in the case of Bernoulli variables $Be(p_1)$ is stochastically dominated by $Be(p_2)$ 345346 whenever $p_1 \leq p_2$. This way we can eliminate dependencies – the outcome of each 347 I_{τ} influences the distributions for $I_{\tau'}$, $\tau' > \tau$ – and work with independent variables 348 349 $I_{\tau+i}^*$.

350 Now, since the new variables are both Bernoulli and independent, we can use the well-known left tail Chernoff bound for binomial setting from [10] (see Corollary 21.7) 351which states that for any $\delta \in (0, 1)$ 352

$$\Pr\left[\sum_{i=0}^{h-1} I_{\tau+i}^* \ge (1+\delta) \mathbb{E}\left[\sum_{i=0}^{h-1} I_{\tau+i}^*\right]\right] \le \exp\left(-\frac{\delta^2}{3} \mathbb{E}\left[\sum_{i=0}^{h-1} I_{\tau+i}^*\right]\right)$$

355 and therefore

356
$$\Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \ge (1+\delta)(1+\varepsilon)\frac{h(pX_{\tau}+r)}{\tau} \middle| \deg_{\tau}(s) \le X_{\tau}, \neg \mathcal{D}_{\varepsilon}(\tau,h) \right]$$

357
$$\le \exp\left(-\frac{h\delta^{2}(1+\varepsilon)(pX_{\tau}+r)}{3\tau}\right).$$

To finish the proof it is sufficient to see that
$$h = \frac{3A\tau \ln t}{\delta^2(1+\varepsilon)(pX_\tau+r)}$$
 gives the required

 $O(t^{-A})$ bound in the last equation. 360 Finally, we proceed with the proof of the main result of this section. 361

THEOREM 3.8. For $G_t \sim DD(t, p, r)$ with $0 and <math>s \in [1, \ln^{1+p} t]$ it holds 362 363 asymptotically that

364
365
$$\Pr\left[\deg_t(s) \ge (1+\alpha) t^p \ln^{2-p^2} t\right] = O(t^{-A})$$

for any constants $\alpha > 0$ and A > 0. 366

367 Proof. Throughout the proof we will use sequences
$$(t_i)_{i=0}^k$$
 and $(X_{t_i})_{i=0}^k$ with
368 $\phi = \ln^{1+p} t, \ \beta_i = p + \frac{1}{2\ln t_i}, \ w_i = \frac{3(A+1)t_i \ln t}{\delta^2(1+\varepsilon)(pX_{t_i}+r)}, \ \text{and} \ \varepsilon = \delta = \frac{1}{5\ln t_i}.$

Now let us define events $\mathcal{A}_i(s) = [\deg_{t_i}(s) < X_{t_i}]$ for $i = 0, \ldots, k$. Clearly, $\mathcal{A}_0(s)$ 372 373 holds since by definition of X_{t_0} we have $\deg_{t_0}(s) < t_0 = X_{t_0}$.

Suppose that $\mathcal{A}_i(s)$ holds. Then we can apply Lemma 3.7 with $\tau = t_i$ and $h = w_i$: 374

375
$$\Pr[\neg \mathcal{A}_{i+1}(s) | \mathcal{A}_i(s)] = \Pr[\deg_{t_{i+1}}(s) \ge X_{t_{i+1}} | \deg_{t_i}(s) < X_{t_i}]$$

376
$$\leq \Pr[\deg_{t_{i+1}}(s) - \deg_{t_i}(s) \ge X_{t_{i+1}} - X_{t_i} | \deg_{t_i}(s) < X_{t_i}]$$

377
$$= \Pr\left[\deg_{t_{i+1}}(s) - \deg_{t_i}(s) \ge \beta_i \frac{w_i X_{t_i}}{t_i} \left| \deg_{t_i}(s) < X_{t_i} \right| \right]$$

378
$$= \Pr\left[\deg_{t_i}(s) - \deg_{t_i}(s) \ge \frac{\beta_i X_{t_i}}{2} \frac{3(A+1)}{2} \ln t \left| \deg_{t_i}(s) < x_{t_i} \right| \right]$$

$$378 \qquad = \Pr\left[\deg_{t_{i+1}}(s) - \deg_{t_i}(s) \ge \frac{p + t - t_i}{(1+\varepsilon)(pX_{t_i}+r)} \frac{\delta(t-1-2)}{\delta^2} \ln t \left| \deg_{t_i}(s) < X_{t_i} \right| \right]$$

$$379 \qquad \leq \Pr\left[\deg_{t_{i+1}}(s) - \deg_{t_i}(s) \ge \frac{3(A+1)(1+\delta)}{\delta^2} \ln t \ \left| \deg_{t_i}(s) < X_{t_i} \right| = O(t^{-A-1})\right]$$

where we used the fact that asymptotically as $t \to \infty$ 381

$$\frac{\beta_i X_{t_i}}{(1+\delta)(1+\varepsilon)(pX_{t_i}+r)} = \frac{\beta_i X_{t_i}}{X_{t_i} \left(p+\delta\left(p+p\varepsilon+\frac{r(1+\varepsilon)}{X_{t_i}}\right)+\varepsilon\left(p+\frac{r}{X_{t_i}}\right)\right)+r}$$

$$\geq \frac{p_{t}}{p + \delta\left(p + p\varepsilon + \frac{r(1+\varepsilon)}{X_{t_{i}}}\right) + \varepsilon\left(p + p\varepsilon + \frac{r(1+\varepsilon)}{X_{t_{i}}}\right) + \frac{r}{X_{t_{i}}}}$$

$$\geq \frac{p + \frac{1}{2\ln t_i}}{p + \delta + \varepsilon + \frac{r}{X_{t_i}}} \ge 1,$$

where in the denominator of the first inequality we used the facts that $p + \frac{r}{X_{t_i}} \leq$ 386 $p + p\varepsilon + \frac{r(1+\varepsilon)}{X_{t_i}} = p + o(1) \le 1$ for any constants $0 when <math>t \to \infty$. 387 Next, we get 388

389
$$\Pr[\deg_t(s) \ge X_{t_k}] \le \Pr[\deg_{t_k}(s) \ge X_{t_k}] = \Pr[\neg \mathcal{A}_k(s)]$$
390
$$\le \sum_{i=0}^{k-1} \Pr[\neg \mathcal{A}_{i+1}(s) | \mathcal{A}_i(s)] + \Pr[\neg \mathcal{A}_0(s)] = \sum_{i=0}^{k-1} O(t^{-A-1}) = O(t^{-A})$$
391

since asymptotically it is true that $w_i \ge 1$ for all $i = 0, \ldots, k$, and therefore $k \le t$. 392

To complete the proof it is sufficient to note that $t_k = t_{k-1}(1+\alpha) \leq (1+\alpha)t$ and thus $X_{t_k} \leq (1+\alpha)t^p \ln^{2-p^2} t$ for any constant $\alpha > 0$. 393 394

3.2. Upper bound, late vertices $(s > t_0)$. In the second part of the proof we 395 also use the sequences $(t_i)_{i=0}^k$ and $(X_{t_i})_{i=0}^k$ as defined in Definition 3.2. Moreover, 396 throughout this section we use the same constants as in the proof of Theorem 3.8: $\phi = \ln^{1+p} t, \ \beta_i = p + \frac{1}{2\ln t_i} \ \text{and} \ w_i = \frac{3(A+1)t_i \ln t}{\delta^2(1+\varepsilon)(pX_{t_i}+r)}.$ 397 398

The proof consists of showing that for $s \in [t_i, t_{i+1})$ for some $i = 0, 1, \ldots, k-1$ 399 the degree graph (i.e. $\deg_s(s)$) is with high probability significantly smaller than its 400 corresponding $X_{t_{i+1}}$. Furthermore, we show that the increase in the degree between 401 $\deg_s(s)$ and $\deg_{t_{i+1}}(s)$ with high probability cannot compensate for this difference. 402

11

Thus, X_t (or, to be more precise, X_{t_k}) gives us a good upper bound on deg_t(s) for all 403 404 s – and therefore also we obtain an upper bound for $\Delta(G_t)$.

Let us introduce auxiliary events $\mathcal{B}_l(s) = \bigcup_{\tau=1}^s \mathcal{A}_l(\tau) = [\deg_{t_l}(\tau) \leq X_{t_l} \text{ for all } \tau \leq s \leq t_l]$ where $\mathcal{A}_i(s)$ is, as before, the event that $\deg_{t_i}(s) \leq X_{t_i}$ for a fixed $s \leq t_i$. 405 406

LEMMA 3.9. Let $s \in (t_l, t_{l+1}]$ for some $l = 0, 1, \ldots, k-1$. Then, for any constants 407 $\varepsilon > 0$ and A > 0408

$$\Pr\left[\deg_s(s) \ge (1+\varepsilon)(pX_{t_{l+1}}+r) \mid \mathcal{B}_l(t_l) \land \mathcal{B}_{l+1}(s-1)\right] = O(t^{-A}).$$

Proof. First, we notice the fact that $\max\{\deg_{t_{l+1}}(\tau): 1 \leq \tau \leq s-1\} \leq X_{t_{l+1}}$ guarantees that $\max\{\deg_s(\tau): 1 \leq \tau \leq s-1\} \leq X_{t_{l+1}}$. Therefore, $\deg_s(s)$ is stochastically dominated by $A_s \sim Bin(X_{t_{l+1}}, p) + Bin(s-1, \frac{r}{s-1})$ and we directly obtain the 411 412413 result using the Chernoff bound with $\mathbb{E}[A_s] = pX_{t_{l+1}} + r$: 414

415
$$\Pr\left[\deg_s(s) \ge (1+\varepsilon)(pX_{t_{l+1}}+r) \left| \mathcal{B}_l(t_l) \wedge \mathcal{B}_{l+1}(s-1) \right] \right]$$

418

 $\leq \exp\left(-\frac{\varepsilon^2}{\varepsilon+2}(pX_{t_{l+1}}+r)\right) \leq t^{-A},$ asymptotically for any constants $\varepsilon, A > 0$ since $X_{t_{l+1}} \ge \ln^{1+p} t$.

Note that the result implies that with high probability at most slightly more than 419a p fraction of the maximum degree is already present at time s. Therefore, we are 420 interested in bounding the remaining part of the degree, i.e. $\deg_{t_{l+1}}(s) - \deg_s(s)$, by 421 something smaller than the remaining fraction of the maximum degree. 422

LEMMA 3.10. Let $s \in (t_l, t_{l+1}]$ for some $l = 0, 1, \ldots, k-1$. Then, for any constant 423 $\alpha > 0$ and A > 0424

425
426
$$\Pr\left[\deg_{t_{l+1}}(s) - \deg_s(s) \ge \alpha X_{t_{l+1}} \ |\mathcal{B}_l(t_l) \land \mathcal{B}_{l+1}(s-1)\right] = O(t^{-A}).$$

Proof. We use Lemma 3.5 with $d = \alpha X_{t_{l+1}}$ to obtain asymptotically as $t \to \infty$ 427 428that for any A > 0 it holds that

429
$$\Pr\left[\deg_{t_{l+1}}(s) - \deg_s(s) \ge \alpha X_{t_{l+1}} | \mathcal{B}_l(t_l) \land \mathcal{B}_{l+1}(s-1) \right]$$

430
$$= \Pr\left[\deg_{t_{l+1}}(s) - \deg_s(s) \ge \alpha X_{t_{l+1}} \right] \le \Pr\left[\deg_{s+w_l}(s) - \deg_s(s) \ge \alpha X_{t_{l+1}} \right]$$

431
$$\le \exp\left(\alpha X_{t_{l+1}} \ln \frac{\exp(1) \cdot w_l p(1+2\alpha) X_{t_{l+1}}}{\alpha X_{t_{l+1}} \cdot s}\right)$$

431
$$\leq \exp\left(\alpha X_{t_{l+1}} \ln \frac{\exp(1) \cdot w_l p(1)}{\alpha X_{t_{l+1}}}\right)$$

432
$$\leq \exp\left(\alpha X_{t_{l+1}}\left(\frac{\exp(1)\cdot(1+2\alpha)\cdot 3(A+1)}{\alpha(1+\alpha)} + \ln\frac{\ln t}{\delta^2(X_{t_l}+r/p)}\right)\right)$$

433
$$\leq \exp\left(\alpha X_{t_{l+1}}\left(\Theta(1) + \ln\frac{25\ln t \cdot \ln^2 t_l}{\max\{\ln^{1+p} t, t_l^p\}}\right)\right)$$

434
435
$$\leq \exp\left(\alpha \ln^{1+p} t\left(\Theta(1) + \ln \frac{25 \ln^2 t_l}{t_l^{p^2/(1+p)}}\right)\right) \leq \exp(-A \ln t) \leq t^{-A}$$

as needed. 436

To proceed we need the following two lemmas. 437

LEMMA 3.11. Let $s \in (t_l, t_{l+1}]$ for some $l = 0, 1, \ldots, k-1$. Then asymptotically 438 as $t \to \infty$, for any constant A > 0 it holds that 439

440
441
$$\Pr\left[\deg_{t_{l+1}}(s) \ge X_{t_{l+1}} | \mathcal{B}_l(t_l) \land \mathcal{B}_{l+1}(s-1)\right] = O(t^{-A}).$$

Proof. We combine Lemma 3.9 with $\varepsilon = \frac{1-p}{4p}$ and Lemma 3.10 with $\alpha = \frac{1-p}{2}$ to 442 obtain 443

444
$$\Pr\left[\deg_{t_{l+1}}(s) \ge X_{t_{l+1}} \mid \mathcal{B}_l(t_l) \land \mathcal{B}_{l+1}(s-1)\right]$$

445
$$\leq \Pr\left[\deg_{s}(s) \geq \left(1 + \frac{1-p}{4p}\right)(pX_{t_{l+1}} + r) \left|\mathcal{B}_{l}(t_{l}) \wedge \mathcal{B}_{l+1}(s-1)\right] + \Pr\left[\deg_{t_{l+1}}(s) - \deg_{s}(s) \geq \frac{1-p}{2}X_{t_{l+1}} \left|\mathcal{B}_{l}(t_{l}) \wedge \mathcal{B}_{l+1}(s-1)\right] = O(t^{-A}).$$

446
447
$$+ \Pr\left[\deg_{t_{l+1}}(s) - \deg_s(s) \ge \frac{1-P}{2} X_{t_{l+1}} \left| \mathcal{B}_l(t_l) \wedge \mathcal{B}_{l+1}(s-1) \right| = O(t^{-A}). \quad \Box$$

448 LEMMA 3.12. Let $s \in (t_l, t_{l+1}]$ for some $l = 0, 1, \ldots, k-1$. Then asymptotically as $t \to \infty$, for any constant A > 0 it holds that 449

$$\Pr\left[\neg \mathcal{B}_{l+1}(t_{l+1}) | \mathcal{B}_l(t_l)\right] = O(t^{-A})$$

Proof. Let l be the first value for which the lemma does not hold. Then, from 452 Lemma 3.11 we get that for any constant A > 0 it holds that 453

454
$$\Pr\left[\neg \mathcal{B}_{l+1}(t_{l+1}) | \mathcal{B}_{l}(t_{l}) \land \mathcal{B}_{l+1}(t_{l})\right] = \sum_{s=t_{l}}^{t_{l+1}-1} \Pr\left[\neg \mathcal{B}_{l+1}(s+1) | \mathcal{B}_{l}(t_{l}) \land \mathcal{B}_{l+1}(s)\right]$$

455
$$= \sum_{s=t_l}^{t_{l+1}} \Pr[\neg \mathcal{A}_{l+1}(s+1) | \mathcal{B}_l(t_l) \wedge \mathcal{B}_{l+1}(s)] = O(t^{-A}).$$

From Theorem 3.8 we know that $\Pr[\mathcal{B}_0(t_0)] = 1 - O(t^{-A})$. Recall that by our 457458

assumption $\Pr[\neg \mathcal{B}_{i+1}(t_{i+1})|\mathcal{B}_i(t_i)] = 1 - O(t^{-A})$ for all $i = 0, 1, \ldots, l-1$, so it follows that $\Pr[\mathcal{B}_i(t_i)] = 1 - O(t^{-A})$ for all $i = 0, 1, \ldots, l$. We use this fact, combined with 459 the observation that $\mathcal{B}_l(t_l) \subseteq \mathcal{A}_l(s)$ and Theorem 3.8 to get 460

461
$$\Pr\left[\neg \mathcal{B}_{l+1}(t_l) | \mathcal{B}_l(t_l)\right] \le \sum_{s=1}^{t_l} \Pr\left[\neg \mathcal{A}_{l+1}(s) | \mathcal{B}_l(t_l)\right]$$

462
$$\leq \sum_{s=1}^{t_l} \frac{\Pr[\neg \mathcal{A}_{l+1}(s) \land \mathcal{B}_l(t_l)]}{\Pr[\mathcal{B}_l(t_l)]} \leq \sum_{s=1}^{t_l} \frac{\Pr[\neg \mathcal{A}_{l+1}(s) \land \mathcal{A}_l(s)]}{\Pr[\mathcal{B}_l(t_l)]}$$

463
464
$$\leq \sum_{s=1}^{t_l} \frac{\Pr[\neg \mathcal{A}_{l+1}(s) | \mathcal{A}_l(s)]}{\Pr[\mathcal{B}_l(t_l)]} = \sum_{s=1}^{t_l} \frac{O(t^{-A})}{1 - O(t^{-A})} = O(t^{-A}).$$

Finally, for any events E_1 , E_2 , E_3 we have 465

466
$$\Pr[\neg E_1 | E_2] = \Pr[\neg E_1 \land E_3 | E_2] + \Pr[\neg E_1 \land \neg E_3 | E_2]$$

467

$$\leq \Pr[\neg E_1 | E_3 \land E_2] + \Pr[\neg E_3 | E_2].$$

We substitute $E_1 = \mathcal{B}_{l+1}(t_{l+1}), E_2 = \mathcal{B}_l(t_l)$ and $E_3 = \mathcal{B}_{l+1}(t_l)$ to obtain the final 469result. 470

This manuscript is for review purposes only.

Finally, we present the main result of this section. 471

THEOREM 3.13. For $G_t \sim DD(t, p, r)$ with $0 and any constants <math>\alpha, A > 0$ 472 it holds asymptotically that 473

474
475
$$\Pr\left[\Delta(G_t) \ge (1+\alpha)t^p \ln^{2-p^2} t\right] = O(t^{-A}).$$

Proof. From Lemma 3.4 we know that $X_{t_k} \leq (1+\alpha)t^p \ln^{2-p^2} t$ holds asymptoti-476cally. It follows that in this case 477

478
$$\Pr\left[\Delta(G_t) \ge (1+\alpha)t^p \ln^{2-p^2} t\right] \le \Pr[\Delta(G_t) \ge X_{t_k}] \le \Pr[\neg \mathcal{B}_k(t_k)]$$
479
$$\le \sum_{k=1}^{k-1} \Pr[\neg \mathcal{B}_{l+1}(t_{l+1})|\mathcal{B}_l(t_l)] + \Pr[\neg \mathcal{B}_0(t_0)].$$

479
480
$$\leq \sum_{l=0}^{n-1} \Pr[\neg \mathcal{B}_{l+1}(t_{l+1}) | \mathcal{B}_l(t_l)] + \Pr[\neg \mathcal{B}_0(t_{l+1}) | \mathcal{B}_l(t_l)] + \Pr[\neg \mathcal{B}_0(t_{l+1}) | \mathcal{B}_l(t_l)] + \Pr[\neg \mathcal{B}_l(t_l) | \mathcal{B}_l(t_$$

Now, from Theorem 3.8 and Lemma 3.12 we know that both $\Pr[\neg \mathcal{B}_0(t_0)] =$ 481 $O(t^{-A})$ and $\Pr[\neg \mathcal{B}_{l+1}(t_l) | \mathcal{B}_l(t_l)] = O(t^{-A})$ for any A > 0, respectively. Putting 482this all together with the fact that asymptotically as $t \to \infty$ it holds that $k \leq t$ we 483 obtain the final result. Π 484

3.3. Lower bound. Here we proceed analogously to the case of the upper bound 485for early vertices. We provide an appropriate Chernoff-type bound for the degree of 486 487 a given vertex with respect to some deterministic sequence. Then we again use a special sequence, which has the desired rate of growth and serves as a lower bound 488 on $\deg_t(s)$. Note that we don't need to extend our analysis for the late vertices since 489 a lower bound for the degree of any vertex s at time t is also a lower bound for the 490minimum degree of G_t . 491

Now, we note that if we start the whole process from a non-empty graph, then 492493 there exists $s \in [1, t_0]$ such that $\deg_{t_0}(s) \ge 1$. Moreover, even if the starting graph is empty, but r > 0, then with high probability there exists a vertex with positive 494degree, as the probability of adding another isolated vertex to an empty graph on t495 vertices is at most $(1-\frac{r}{t})^t \leq \exp(-r)$, so within first $\frac{A}{r} \ln t$ vertices for any A > 0 we 496have a non-isolated vertex with probability at least $1 - O(t^{-A})$. Of course, if we start 497from an empty graph and r = 0, then for any p there is no edge in the duplication 498499process. However, in this case it trivially follows that $\Delta(G_t) = 0$, so we omit this case in further analysis. 500

That said, let us now proceed with the aforementioned Chernoff-type lower bound 501for the degree of a given early vertex: 502

LEMMA 3.14. Let $1 \le s \le \tau \le t$ such that $\tau \ge \phi = \ln^{1+p} t$. Then for any A > 0503 it is true that 504

505
506
$$\Pr\left[\deg_{\tau+h}(s) - \deg_{\tau}(s) \le \frac{2A(1-\delta)}{\delta^2} \ln t \ \left| \deg_{\tau}(s) \ge X_{\tau} \right| = O(t^{-A}),$$

with $\varepsilon = \delta = \frac{p(1-p)}{8\ln \tau}$ and $h = \frac{2A\tau \ln t}{\delta^2(1-\varepsilon)(pX_{\tau}+r)}$ 507

Proof. Let us recall (as in the proof of Lemma 3.5) that for i = 0, 1, ..., h - 1 we 508 have $\deg_{\tau+i+1}(s) = \deg_{\tau+i}(s) + I_{\tau+i}$ where $I_{\tau+i} \sim Be(q_{\tau+i})$ for $q_{\tau+i} = \frac{p \deg_{\tau+i}(s) + r}{\tau+i}$ 509

510 $\frac{r \deg_{\tau+i}(s)}{(\tau+i)^2}$. Also clearly $\deg_{\tau+i}(s) \ge \deg_{\tau}(s)$ for any $i = 0, 1, \dots, h$, so we have

511
$$q_{\tau+i} = \frac{p \deg_{\tau+i}(s) \left(1 - \frac{r}{p(\tau+i)}\right) + r}{\tau+i} \ge \frac{p \deg_{\tau}(s) \left(1 - \frac{r}{p\tau}\right) + r}{\tau+h}$$

$$\sum_{512} \sum \frac{pX_{\tau} \left(1-\varepsilon^2\right)+r}{\tau(1+\varepsilon)} \ge (1-\varepsilon)\frac{pX_{\tau}+r}{\tau}$$

514 since for $\varepsilon = \frac{p(1-p)}{8 \ln \tau}$ it holds that $h \leq \varepsilon t$ and $\varepsilon^2 \geq \frac{r}{p\tau}$. Therefore for any i =515 $0, 1, \ldots, h-1$ we know that $I_{\tau+i}$ stochastically dominates $I_{\tau+i}^* \sim Be\left((1-\varepsilon)\frac{pX_{\tau}+r}{\tau}\right)$. 516 As in the proof of the upper bound, the new variables are both Bernoulli and 517 independent. So this time we can use the right tail Chernoff bound for binomial 518 setting from [10] (see Corollary 21.7) which states that for any $\delta \in (0, 1)$

519
$$\Pr\left[\sum_{i=0}^{h-1} I_{\tau+i}^* \le (1-\delta) \mathbb{E}\left[\sum_{i=0}^{h-1} I_{\tau+i}^*\right]\right] \le \exp\left(-\frac{\delta^2}{2} \mathbb{E}\left[\sum_{i=0}^{h-1} I_{\tau+i}^*\right]\right)$$
520

521 and therefore

522
$$\Pr\left[\deg_{\tau+h}(s) \le \deg_{\tau}(s) + (1-\delta)(1-\varepsilon)\frac{h(pX_{\tau}+r)}{\tau}\right] \le \exp\left(-\frac{h\delta^2(1-\varepsilon)(pX_{\tau}+r)}{2\tau}\right)$$

524 as clearly $\Pr[\deg_{\tau+h}(s) - \deg_{\tau}(s) \le k] = \Pr\left[\sum_{i=0}^{h-1} I_{\tau+i} \le k\right] \le \Pr\left[\sum_{i=0}^{h-1} I_{\tau+i}^* \le k\right]$ 525 for any k, due to the stochastic dominance.

To finish the proof it is sufficient to see that $h = \frac{2A\tau \ln t}{\delta^2(1-\varepsilon)(pX_\tau+r)}$ gives the required 527 $O(t^{-A})$ bound in the last equation.

In the following, we again use sequences $(t_i)_{i=1}^k$ and $(X_{t_i})_{i=1}^k$ from Definition 3.2. Let us also define $C_i(s) = [\deg_{t_i}(s) > X_{t_i} - \phi + 1]$ for a fixed $s \le t_i$. Now we are in the position to proceed with the main theorem of this section:

531 THEOREM 3.15. For $G_t \sim DD(t, p, r)$ with 0 there exists s such that it532 holds asymptotically that

$$\Pr\left[\deg_t(s) < (1-\alpha)t^p\right] = O(t^{-A})$$

535 for any constants $\alpha, A > 0$.

534

536 Proof. Let us use $\phi = \ln^{1+p} t$, $\beta_i = p - \frac{p(1-p)}{4 \ln t_i}$ and $w_i = \frac{2(A+1)t_i \ln t}{\delta^2(1-\varepsilon)(pX_{t_i}+r)}$ with 537 $\delta = \varepsilon = \frac{p(1-p)}{8 \ln t_i}$.

Suppose that $C_i(s)$ holds. Then we can apply Lemma 3.14 with $\tau = t_i$ and $h = w_i$:

539
$$\Pr[\neg \mathcal{C}_{i+1}(s) | \mathcal{C}_i(s)] = \Pr[\deg_{t_{i+1}}(s) \le X_{t_{i+1}} | \deg_{t_i}(s) > X_{t_i} - \phi + 1]$$

540
$$\leq \Pr[\deg_{t_{i+1}}(s) - \deg_{t_i}(s) \leq X_{t_{i+1}} - X_{t_i} \mid \deg_{t_i}(s) > X_{t_i} - \phi + 1]$$

541
$$= \Pr\left[\deg_{t_{i+1}}(s) - \deg_{t_i}(s) \le \beta_i \frac{\omega_i x_{t_i}}{t_i} \mid \deg_{t_i}(s) > X_{t_i} - \phi + 1\right]$$

542
$$\leq \Pr\left[\deg_{t_{i+1}}(s) - \deg_{t_i}(s) \le \frac{2(A+1)(1-\delta)}{\delta^2} \ln t \ \middle| \ \deg_{t_i}(s) > X_{t_i} - \phi + 1 \right]$$

$$543_{544} = O(t^{-A-1}),$$

where we used the fact that asymptotically as $t \to \infty$ it holds that 545

$$\frac{\beta_i X_{t_i}}{(1-\delta)(1-\varepsilon)(pX_{t_i}+r)} \le \frac{p - \frac{p(1-p)}{4\ln t_i}}{p(1-\delta-\varepsilon)} = 1.$$

548 Next, we get

549
$$\Pr[\deg_t(s) \le X_{t_k} - \phi + 1] \le \Pr[\deg_{t_k}(s) \le X_{t_k} - \phi + 1] = \Pr[\neg \mathcal{C}_k(s)]$$

550
$$\le \sum_{k=1}^{k-1} \Pr[\neg \mathcal{C}_{i+1}(s) | \mathcal{C}_i(s)] + \Pr[\neg \mathcal{C}_0(s)] = \sum_{k=1}^{k-1} O(t^{-A-1}) = O(t^{-A}),$$

550
$$\leq \sum_{i=0} \Pr[\neg \mathcal{C}_{i+1}(s) | \mathcal{C}_i(s)] + \Pr[\neg \mathcal{C}_0(s)] = \sum_{i=0} O(t^{-A-1}) =$$

552 since asymptotically it is true that $w_i \ge 1$ for all $i = 0, \ldots, k$, and therefore $k \le t$. To complete the proof it is sufficient to note that $t \le t_k \le (1+\alpha)t_{k-1} \le (1+\alpha)t$ 553for any constant $\alpha > 0$ and thus $X_{t_k} \leq (1 + \alpha)t^p$. 554Π

4. Average degree. Now let us proceed to the results on the average degree of 556 G_t defined as

557
558
$$D(G_t) = \frac{1}{t} \sum_{s=1}^t \deg_t(s).$$

559 First, we recall from [30, Theorem 9(iii)] that for any $\tau = t_0, \ldots, t-1$ it holds asymptotically (i.e. when $t_0 \to \infty$) that 560

$$561 \quad \mathbb{E}[\deg_{\tau}(\tau)] = \begin{cases} D(G_{t_0}) \frac{p\Gamma(t_0)\Gamma(t_0+1)}{\Gamma(t_0+c_3)\Gamma(t_0+c_4)} \tau^{2p-1}(1+o(1)) & \text{if } p \leq \frac{1}{2}, r = 0, \\ r(1+o(1)) & \text{if } p = 0, r > 0, \\ \left(\frac{r(1-p)}{p(1-2p)} - \frac{r}{p}\right) (1+o(1)) & \text{if } 0 0, \\ r\log\tau(1+o(1)) & \text{if } p = \frac{1}{2}, r > 0, \\ \left(D(G_{t_0}) + \frac{2rt_0}{t_0^2+2pt_0-2r} \,_3F_2\left[\frac{t_0+1,t_0+1,1}{t_0+c_4+1};1\right]\right) \\ \frac{p\Gamma(t_0)\Gamma(t_0+1)}{\Gamma(t_0+c_3)\Gamma(t_0+c_4)} \tau^{2p-1}(1+o(1)) & \text{if } p > \frac{1}{2}, \end{cases}$$

where $D(G_{t_0})$ is the average degree of the initial graph G_{t_0} and 563

564
$${}_{3}F_{2}\left[{}^{a_{1},a_{2},a_{3}}_{b_{1},b_{2}};z\right] = \sum_{l=0}^{\infty} \frac{(a_{1})_{l}(a_{2})_{l}(a_{3})_{l}}{(b_{1})_{l}(b_{2})_{l}} \frac{z^{l}}{l!}$$

is the generalized hypergeometric function with $(a)_l = a(a+1) \dots (a+l-1), (a)_0 = 1$ 565the rising factorial (see [1] for details). 566

In short, if we omit constant factors, there are three regimes of growth: constant, 567 $\ln t$, and t^{2p-1} . We need to find the proper high probability bound for each case 568 569 separately, however it turns out that the proofs are very similar.

4.1. Upper bound. Now we may proceed to the main result of this section: 570the upper bound for the average degree of G_t . It turns out that there are exactly two 571regimes with somewhat different behavior: 572

573 THEOREM 4.1. Asymptotically for
$$G_t \sim DD(t, p, r)$$
 it holds that

574
$$\Pr[D(G_t) \ge A C \ln t] = O(t^{-A}) \qquad \text{for } p \le \frac{1}{2}$$

575
576
$$\Pr[D(G_t) \ge C t^{2p-1}] = O(t^{-A})$$
 for $p > \frac{1}{2}$.

for some fixed constant C > 0 and any A > 0. 577

Proof. For simplicity, we will work with the total number of edges $\tau D(G_{\tau})$ instead 578 of $D(G_{\tau})$. Clearly, for any $\tau = t_0, \ldots, t-1$ it holds that 579

580
$$(\tau+1)D(G_{\tau+1}) - \tau D(G_{\tau}) = 2\deg_{\tau+1}(\tau+1),$$

 $\deg_{\tau+1}(\tau+1) \sim Bin(\deg_{\tau}(\operatorname{parent}(\tau+1)), p) + Bin(\tau - \deg_{\tau}(\operatorname{parent}(\tau+1)), r/\tau).$ 381

Therefore, we can use Chernoff bound to obtain for any $\delta \geq 0$ 583

584585

586

$$\Pr\left[(\tau+1)D(G_{\tau+1}) - \tau D(G_{\tau}) \ge 2(1+\delta) \mathbb{E}[\deg_{\tau+1}(\tau+1)]\right]$$
$$\le \exp\left(-\frac{2\delta^2}{2+\delta} \mathbb{E}[\deg_{\tau+1}(\tau+1)]\right).$$

Now, for $p > \frac{1}{2}$ we know that $\mathbb{E}[\deg_{\tau}(\tau)] \leq C^* \tau^{2p-1}$ for some constant $C^* > 0$. 587 Thus, it is sufficient to set $t_0 = t^{p/3}$ and $\delta = \sqrt{\frac{3(A+1)\ln t}{2C^*\tau^{2p-1}}} = o(1)$ for all $\tau = t_0, \ldots, t-1$ 588 to get 589

590
$$\Pr\left[(\tau+1)D(G_{\tau+1}) - \tau D(G_{\tau}) \ge 2(1+\delta)C^* \tau^{2p-1}\right] = O(t^{-A-1}),$$

and by summing over all τ that no event from polynomial tails happens we obtain

593 Pr
$$[tD(G_t) \ge C t^{2p}] \le \Pr\left[tD(G_t) - t_0D(G_{t_0}) \ge \sum_{i=t_0}^{t-1} 2(1+\delta)C^* \tau^{2p-1}\right] = O(t^{-A}),$$

594

for any constant $C \ge t^{-2p} \sum_{i=t_0}^{t-1} 2(1+\delta)C^* \tau^{2p-1} + t^{-2p}t_0 D(G_{t_0})$ – and such constant indeed exists since it is not hard to verify that the latter sum is finite. 595 596

In all cases $0 it turns out that <math>\sqrt{\frac{3(A+1)\ln t}{2C^*\tau^{2p-1}}} \to \infty$. However, for 0 ,597 r > 0 we have $\mathbb{E}[\deg_{\tau}(\tau)] \leq C^* \ln \tau$ for some constant $C^* > 0$, and we can assume 598 $\delta \rightarrow \infty$ such that 599

$$\frac{1+\delta}{2} \le \frac{\delta^2}{2+\delta} = \frac{(A+1)\ln t}{2C^*\ln \tau},$$

so therefore 602

603
$$\Pr[(\tau+1)D(G_{\tau+1}) - \tau D(G_{\tau}) \ge 2(A+1)\ln t] = O(t^{-A-1}),$$

604
$$\Pr[tD(G_t) \ge ACt\ln t] \le \Pr\left[tD(G_t) - t_0D(G_{t_0}) \ge \sum_{i=t_0}^{t-1} 2(A+1)\ln i\right] = O(t^{-A}),$$

605

606

for some constant $C \ge 2 + \frac{t_0}{At \ln t} D(G_{t_0})$ when $t_0 = t^{1/3}$. Finally, let us study the case 0 , <math>r = 0. Again we know that $\mathbb{E}[\deg_{\tau}(\tau)] \le \frac{1}{2}$. 607 $C^* \tau^{2p-1}$ for some constant $C^* > 0$. Again, we can assume 608

$$\frac{1+\delta}{2} \le \frac{\delta^2}{2+\delta} = \frac{(A+1)\ln t}{2C^*\tau^{2p-1}},$$

so by a similar reasoning as before we get 611

612
$$\Pr[(\tau+1)D(G_{\tau+1}) - \tau D(G_{\tau}) \ge 2(A+1)\ln t] = O(t^{-A-1}),$$

613
$$\Pr[tD(G_t) \ge ACt \ln t] \le \Pr\left[tD(G_t) - t_0D(G_{t_0}) \ge \sum_{i=t_0}^{t-1} 2(A+1)\ln t\right] = O(t^{-A}),$$
614

for sufficiently large constant C when $t_0 = t^{1/3}$.

4.2. Lower bound. We now turn our attention to establishing the corresponding lower bound. Note that since $\mathbb{E}[D(G_t)] = O(\log t)$ for $p \leq \frac{1}{2}$, the lower polynomial tail is trivial in this range since all smaller values are within the polylogarithmic distance from the mean. However, we can investigate the case $p > \frac{1}{2}$.

620 THEOREM 4.2. For $G_t \sim DD(t, p, r)$ with $p > \frac{1}{2}$ asymptotically it holds that

$$\Pr[D(G_t) \le C t^{2p-1}] = O(t^{-A}).$$

623 for some fixed constant C > 0 and any A > 0.

624 *Proof.* Similarly as before, we invoke the appropriate Chernoff bound for $\delta \in (0, 1)$

625
$$\Pr\left[(\tau+1)D(G_{\tau+1}) - \tau D(G_{\tau}) \le 2(1-\delta)\mathbb{E}[\deg_{\tau+1}(\tau+1)]\right]$$

626

18

$$\leq \exp\left(-\delta^2 \mathbb{E}[\deg_{\tau+1}(\tau+1)]\right).$$

For $p > \frac{1}{2}$ it is true that $\mathbb{E}[\deg_{\tau}(\tau)] \ge C^* \tau^{2p-1}$ for some constant $C^* > 0$. Thus, it is sufficient to set $t_0 = t^{p/3}$ and $\delta = \sqrt{\frac{(A+1)\ln t}{C^* \tau^{2p-1}}} \le \frac{1}{2}$ for all $\tau = t_0, \ldots, t-1$ to get

$$\Pr\left[(\tau+1)D(G_{\tau+1}) - \tau D(G_{\tau}) \le 2(1-\delta) C^* \tau^{2p-1}\right] = O(t^{-A-1}),$$

632 which leads us to

633 634

635

$$\Pr\left[tD(G_t) - t_0D(G_{t_0}) \le C t^{2p}\right]$$

$$\le \Pr\left[tD(G_t) - t_0D(G_{t_0}) \le \sum_{i=t_0}^{t-1} 2(1-\delta)C^* \tau^{2p-1}\right] = O(t^{-A})$$

636 for any constant $0 < C \leq t^{-2p} \sum_{i=t_0}^{t-1} 2(1-\delta)C^* \tau^{2p-1} + t^{-2p}t_0 D(G_{t_0})$ – and such 637 constant indeed exists since it is not hard to verify that the latter sum is non-zero 638 and finite when $t_0 = t^{1/3}$.

639 **5. Further challenges.** In this paper we focus on deriving large deviations 640 for the average and the maximum degree in the duplication-divergence networks. 641 By a simple martingale argument one can show that $\Delta(G_t)/t^p$ converges to some 642 random variable Δ . However, it is still worth asking whether Δ has finite support 643 (e.g. dependent only on p and r, but not on t).

A natural next challenge would be to obtain the exact asymptotic formula for the 644 whole degree distribution. For example, there is an open question whether DD(t, p, r)645 graphs are scale-free, i.e. they have $\Theta(k^{-\gamma})$ fraction of vertices with degree k. A first 646 647 step towards this goal was already done for r = 0 in [18, 16], where it was proved that this property indeed holds for the (only) giant component $p < e^{-1}$. However, it was 648 noticed in [14] that for r = 0 and all 0 such phenomenon does not appear in649 the whole graph, since almost all vertices are isolated, thus for any k > 0 the fraction 650 of vertices of degree k tends to 0 as $t \to \infty$. 651

Finally, finding good bounds on the concentration of both $D(G_t)$ and $\Delta(G_t)$ is only the step towards the full understanding of this model, as we still do not know for example how symmetric such networks are. This, in turn, we believe could help find good compression algorithms for these types of networks, as was the case with other graph models [3, 23].

REFERENCES

- 658 [1] M. ABRAMOWITZ AND I. STEGUN, Handbook of mathematical functions: with formulas, graphs, 659 and mathematical tables, vol. 55, Dover Publications, 1972.
- 660 B. BOLLOBÁS, Random graphs, Cambridge University Press, 2001. [2]
- 661 [3] F. CHIERICHETTI, R. KUMAR, S. LATTANZI, A. PANCONESI, AND P. RAGHAVAN, Models for the 662 compressible web, SIAM Journal on Computing, 42 (2013), pp. 1777–1802.
- 663 [4] F. CHUNG AND L. LU, Complex graphs and networks, no. 107 in CBMS Regional Conference 664Series in Mathematics, American Mathematical Society, 2006.
- 665 [5] F. CHUNG, L. LU, T. G. DEWEY, AND D. GALAS, Duplication models for biological networks, 666 Journal of Computational Biology, 10 (2003), pp. 677-687.
- [6] R. COLAK, F. HORMOZDIARI, F. MOSER, A. SCHÖNHUTH, J. HOLMAN, M. ESTER, AND S. C. 667 668 SAHINALP, Dense graphlet statistics of protein interaction and random networks, in Bio-669 computing 2009, World Scientific Publishing, Singapore, 2009, pp. 178-189.
- 670 R. DIESTEL, Graph Theory, Springer, 2005. [7]
- 671 P. ERDŐS AND A. RÉNYI, On random graphs I, Publicationes Mathematicae, 6 (1959), pp. 290-[8] 672 297.
- [9] M. FALOUTSOS, P. FALOUTSOS, AND C. FALOUTSOS, On power-law relationships of the internet 673 674topology, ACM SIGCOMM Computer Communication Review, 29 (1999), pp. 251–262.
- 675 [10] A. FRIEZE AND M. KAROŃSKI, Introduction to Random Graphs, Cambridge University Press, 676 2016.
- 677 [11] A. FRIEZE, K. TUROWSKI, AND W. SZPANKOWSKI, Degree distribution for duplication-divergence graphs: large deviations, in Graph-Theoretic Concepts in Computer Science: 46th Inter-678679 national Workshop, WG 2020, Leeds, UK, June 24–26, 2020, Revised Selected Papers, 680 Springer, 2020, pp. 226–237.
- 681 [12] A. FRIEZE, K. TUROWSKI, AND W. SZPANKOWSKI, The concentration of the maximum degree in the duplication-divergence models, in 27th International Conference on Computing 682 683 and Combinatorics, COCOON 2021, Tainan, Taiwan, October 24-26, 2021, Proceedings, 684 Springer, 2021, pp. 413-424.
- 685[13] R. GRAHAM, D. KNUTH, AND O. PATASHNIK, Concrete Mathematics: A Foundation for Com-686puter Science, Addison-Wesley Professional, 1994.
- 687 [14] F. HERMANN AND P. PFAFFELHUBER, Large-scale behavior of the partial duplication random 688 graph, ALEA, 13 (2016), pp. 687-710.
- [15] F. HORMOZDIARI, P. BERENBRINK, N. PRŽULJ, AND S. C. SAHINALP, Not all scale-free networks 689 690 are born equal: the role of the seed graph in PPI network evolution, PLoS Computational 691 Biology, 3 (2007), p. e118.
- [16] P. JACQUET, K. TUROWSKI, AND W. SZPANKOWSKI, Power-law degree distribution in the con-692 693 nected component of a duplication graph, in 31st International Conference on Probabilistic, 694 Combinatorial and Asymptotic Methods for the Analysis of Algorithms (AofA 2020), 2020. 695
 - S. JANSON, T. ŁUCZAK, AND A. RUCIŃSKI, Random graphs, John Wiley & Sons, 2011. [17]
- 696 [18] J. JORDAN, The connected component of the partial duplication graph, ALEA – Latin American 697 Journal of Probability and Mathematical Statistics, 15 (2018), pp. 1431-1445.
- 698 [19] B. KAMINSKI, P. PRAŁAT, AND F. THÉBERGE, Mining complex networks, CRC Press, 2021.
- [20] V. LATORA, V. NICOSIA, AND G. RUSSO, Complex networks: principles, methods and applica-699 700 tions, Cambridge University Press, 2017.
- 701 [21] S. LI, K. P. CHOI, AND T. WU, Degree distribution of large networks generated by the partial 702 duplication model, Theoretical Computer Science, 476 (2013), pp. 94–108.
- 703[22] S. LI, K. P. CHOI, T. WU, AND L. ZHANG, Maximum likelihood inference of the evolutionary 704 history of a ppi network from the duplication history of its proteins, IEEE/ACM Transac-705 tions on Computational Biology and Bioinformatics, 10 (2013), pp. 1412-1421.
- 706 [23] T. ŁUCZAK, A. MAGNER, AND W. SZPANKOWSKI, Compression of preferential attachment 707 graphs, in 2019 IEEE International Symposium on Information Theory (ISIT), IEEE, 2019, 708 pp. 1697-1701.
- 709 [24]M. NEWMAN, Networks: An Introduction, Oxford University Press, 2010.
- 710 [25]S. OHNO, Evolution by gene duplication, Springer-Verlag, Berlin-Heidelberg, 1970.
- 711 [26]R. PASTOR-SATORRAS, E. SMITH, AND R. SOLÉ, Evolving protein interaction networks through 712 gene duplication, Journal of Theoretical Biology, 222 (2003), pp. 199-210.
- 713 [27] M. SHAO, Y. YANG, J. GUAN, AND S. ZHOU, Choosing appropriate models for protein-protein 714interaction networks: a comparison study, Briefings in Bioinformatics, 15 (2013), pp. 823-715838.
- [28] R. SOLÉ, R. PASTOR-SATORRAS, E. SMITH, AND T. KEPLER, A model of large-scale proteome 716 717evolution, Advances in Complex Systems, 5 (2002), pp. 43-54.
- [29] J. K. SREEDHARAN, K. TUROWSKI, AND W. SZPANKOWSKI, Revisiting parameter estimation in 718 719biological networks: Influence of symmetries, IEEE/ACM Transactions on Computational

A. FRIEZE, K. TUROWSKI, AND W. SZPANKOWSKI

- 720 Biology and Bioinformatics, 18 (2020), pp. 836–849.
- [30] K. TUROWSKI AND W. SZPANKOWSKI, Towards degree distribution of a duplication-divergence 721 [60] R. FORGWAR MED W. DEFINITIONSIR, FORGER and Regree and Forger of a depresent and spectrum and s 722
- 723 724 2016.
- 725[32] J. ZHANG, Evolution by gene duplication: an update, Trends in Ecology & Evolution, 18 (2003), 726pp. 292–298.