

Stable and Scalable Universal Swarms

Ji Zhu
University of Illinois
jizhu1@illinois.edu

Stratis Ioannidis,
Nidhi Hegde
Technicolor
stratis.ioannidis@technicolor.com,
nidhi.hegde@technicolor.com

Laurent Massoulié
MSR-Inria Joint Centre
laurent.massoulie@inria.fr

ABSTRACT

Hajek and Zhu recently showed that the BitTorrent protocol can become unstable when peers depart immediately after downloading all pieces of a file. In light of this result, Zhou *et al.* propose bundling swarms together, allowing peers to exchange pieces across different swarms, and claim that such “universal swarms” can increase BitTorrent’s stability. In this work, we formally characterize the stability region of universal swarms and show that they indeed exhibit excellent properties. In particular, bundling allows a single seed with limited upload capacity to serve an arbitrary number of disjoint swarms if the arrival rate of peers in each swarm is lower than the seed upload capacity. Our result also shows that the stability region is insensitive to peers’ upload capacity, piece selection policies and number of swarms.

Categories and Subject Descriptors: C.2.1 [Computer Communication Networks]: Network Architecture and Design - distributed networks

Keywords: Missing Piece Syndrome; P2P; Stability

1. INTRODUCTION

BitTorrent is one of the most popular peer-to-peer protocols, used by millions of Internet users to share files online. In simple terms, peers interested in downloading a single file from a distinguished user, termed the *seed*, form a so-called *swarm*. Peers in a swarm exchange file *pieces* (or *chunks*) with each other. Each peer thereby acts as both a client and a server, contributing to the aggregate upload capacity of the swarm. A natural question about BitTorrent is what is its *stability region*: assuming the seed’s upload capacity is U pieces per second, what is the largest arrival rate of peers λ that can be supported without the swarm growing to infinity? Intuitively, as every incoming peer increases the swarm’s aggregate upload capacity, one would expect BitTorrent to support high arrival rates.

Determining the stability region of BitTorrent has been an open problem for more than ten years. It was resolved

recently by Hajek and Zhu [3], who showed that the swarm remains stable if $\lambda < U$, and unstable if $\lambda > U$.

In light of this, a series of recent works have focused on extending the stability region of BitTorrent [11, 10, 8, 6, 4]. Zhou *et al.* [10] propose bundling multiple autonomous swarms together into a *universal swarm*: peers immediately depart upon retrieving the file they desire but, while in the system, they store and exchange pieces with peers belonging to *different swarms*. Intuitively, such inter-swarm exchanges utilize bandwidth that would otherwise remain idle. As such, Zhou *et al.* conjecture that a universal swarm with a single seed has an increased stability region compared to autonomous swarms. Sharing pieces with different swarms introduces a trade-off between stability and the average *sojourn time*, *i.e.*, the time peers stay in the system: by consuming bandwidth for pieces they are not interested in, peers may take longer to retrieve the file they desire.

In this paper, we establish that universal swarms have excellent stability properties. In addition, they can be designed so that peers *do not* experience increased delays. In particular, we make the following contributions: First, we formally characterize the stability region of universal swarms, proving necessary and sufficient conditions under which such swarms remain stable. Second, we show that bundling swarms together significantly increases BitTorrent’s stability region, making it scale in the number of bundled swarms. We show that the capacity region under limited seed upload capacity *is insensitive to the number of swarms, as well as the peers’ uploading capacity and piece selection policies*. Third, we show that the increased stability region of universal swarms comes at the cost of increased delays, that scale with the total number of bundled swarms. Finally, to address this, we propose a modified system design that provably extends the stability region *and* does not affect delays.

Our stability result (Theorem 2) is general: in contrast to previous work [3, 11, 7, 10], it covers all work-conserving piece selection policies, determining which pieces peers exchange upon contact with each other. As a consequence, our result has interesting implications for single swarms as well. Moreover, to the best of our knowledge, we are the first to observe and exploit the fact that swarms become *meta-stable* when the swarm seed and peers prioritize rare pieces.

The remainder of the paper is organized as follows. We review related work in the next section. Our model of universal swarm is introduced in Section 3. Earlier results and our new results on stability constitute Sections 4 and 5 respectively. We study sojourn time numerically in Section 6 and conclude in Section 7.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

PODC’13, July 22–24, 2013, Montréal, Québec, Canada.

Copyright 2013 ACM 978-1-4503-2065-8/13/07 ...\$15.00.

2. RELATED WORK

The stability of BitTorrent is determined by a phenomenon called the *missing piece syndrome*, described in detail in Section 4. This was first observed by Hajek and Zhu [3] and, independently, by Norros and Reitu [7] (in the context of a two-piece file sharing system). Hajek and Zhu observed that the phenomenon persists—and the stability region remains the same—when using source coding as well as when prioritizing the transmission of rare pieces.

One way to extend the stability region of BitTorrent is to require that peers remain in the system after they have downloaded all pieces, effectively turning from leechers to seeds. In follow-up work, Zhu and Hajek show that this indeed extends the stability region by a factor that depends on the mean additional time peers spend in the system [11]. By setting the latter to a high enough value, the system remains stable for arbitrary arrival rates. Requiring that peers spend additional time in the system assumes they are altruistic. Universal swarms also presume altruism, though of a different kind: though peers immediately leave upon retrieving a file, they contribute their storage and bandwidth resources to other swarms while in the system, by storing and uploading pieces of files they may not be interested in.

Oguz *et al.* [8] propose an alternative approach to extending BitTorrent’s stability region. The authors consider schemes where (a) newly arriving peers do not accept pieces unless they are sufficiently rare, and (b) peers missing a single piece do not download it before uploading a piece that is sufficiently rare. Crucially, the rarity of a piece is estimated in a distributed fashion, through sampling a small number of peers. The resulting system is stable under arbitrary arrival rates, and also assumes peers are altruistic, as they do not download useful pieces at all opportunities to do so.

Zhou *et al.* [10] propose universal swarms as a means of extending the stability region of BitTorrent. They study stability under a model in which (a) files consist of only one piece and (b) swarms are *seedless* [5]: no seed exists, and all peers arrive already endowed with certain pieces. The authors observe an extension of the stability region under these assumptions; in our work, we consider the more realistic setup of [3, 8, 11], in which files consist of multiple pieces and peers arrive empty, while pieces are injected by a seed.

Bundling across swarms has also been studied in other contexts. Menasche *et al.* [6] show that bundling swarms together improves file availability, while Wu *et al.* [9] show that bundling improves the performance of live streaming. In contrast, our work presents a rigorous analysis of the stability of universal swarms under various sharing policies. Finally, measurement studies [4] establish the prevalence of file bundling in BitTorrent: 50 – 80% of swarms contain bundled files, with content such as movies, TV series and music being bundled more often than games. These real-world measurements attest to the relevance and practical significance of studying universal swarms.

3. SYSTEM DESCRIPTION

We consider a BitTorrent-like file-sharing system, consisting of multiple swarms. Suppose all files are divided into *pieces* of equal size. We denote by $\mathcal{F} = \{1, 2, \dots, K\}$ the set of all pieces of all files. A distinguished peer, the seed, is always present and holds \mathcal{F} . Represent a *file* C to be a non-empty subset of \mathcal{F} : $\emptyset \neq C \subseteq \mathcal{F}$. We refer to the set of peers

interested in downloading file C as *Swarm* C . Assume peers in Swarm C arrive according to a Poisson process with rate λ_C , independent across swarms. We do *not* require different swarms to be disjoint subsets of \mathcal{F} , so our model captures a scenario where arriving peers are interested in multiple files.

Each peer maintains a cache to store pieces it downloads. We assume peers arrive with empty caches, and each peer’s cache is large enough to hold \mathcal{F} . Peers in Swarm C depart *immediately* upon retrieving all pieces in C . We partition peers into *types* according to (a) the swarm they belong to and (b) the set of pieces in their cache. Hence, a peer in Swarm C holding $S \subseteq \mathcal{F}$ is denoted to be of type $\langle C, S \rangle$. Assume the seed is of type $\langle \{\perp\}, \mathcal{F} \rangle$ for some piece $\perp \notin \mathcal{F}$. Denote $n_{\langle C, S \rangle}$ to be the number of type $\langle C, S \rangle$ peers and $\mathbf{n} = (n_{\langle C, S \rangle})$ to be the vector of numbers of peers in all types. The seed uploads pieces at instants of a Poisson process of rate U . At each such instant, the seed contacts a peer selected uniformly at random among all peers across all swarms, and replicates a piece in \mathcal{F} to this peer. Similarly, at instances that follow a Poisson process of rate $\mu > 0$, each peer contacts another peer (also selected uniformly among all peers) and replicates a piece from its cache.

The piece replicated when a source (either a peer or the seed) contacts a receiver is determined by the source’s piece selection policy. We consider a broad class of work-conserving piece selection policies satisfying the following:

ASSUMPTION 1. *If a source in type $\langle C, S \rangle$ contacts a receiver in type $\langle C', S' \rangle$, no piece is replicated if $S \subseteq S'$. Otherwise exactly one piece in $S \setminus S'$ is replicated, with piece $i \in S \setminus S'$ replicated with probability $h_{\langle C, S \rangle}(i, \langle C', S' \rangle, \mathbf{n}) \in [0, 1]$, determined by the types of the source and the receiver, the piece id, and the current state \mathbf{n} . Function $h_{\langle C, S \rangle}$, also referred to as the policy, satisfies:*

$$\begin{aligned} \sum_{i \notin S \setminus S'} h_{\langle C, S \rangle}(i, \langle C', S' \rangle, \mathbf{n}) &= 0, \\ \sum_{i \in S \setminus S'} h_{\langle C, S \rangle}(i, \langle C', S' \rangle, \mathbf{n}) &= 1 \text{ if } S \not\subseteq S'. \end{aligned}$$

Suppose a type $\langle C, S \rangle$ source contacts a type $\langle C', S' \rangle$ receiver. Examples of work-conserving policies considered are: **Random Novel [RN]**: If $S \setminus S' \neq \emptyset$, the source replicates a piece chosen uniformly from $S \setminus S'$.

Rarest First [RF]: Define the *availability* of a piece $i \in \mathcal{F}$ to be the number of peers holding it. The source replicates the piece in $S \setminus S'$ that has the least availability, with ties broken randomly.

Priority Rarest First [PRF]: The source prioritizes pieces within the swarm of the receiver: if $(S \setminus S') \cap C' \neq \emptyset$, it replicates the piece in $(S \setminus S') \cap C'$ that has the least availability; if $(S \setminus S') \cap C'$ is empty but $S \setminus S'$ is not, the source reverts to RF. **Priority Random Novel [PRN]** is defined similarly.

Notice that sources of the same type apply the same policy. The piece selection policy of the system is denoted by a tuple of $h_{\langle C, S \rangle}$ indexed by each $\langle C, S \rangle$, where all sources in type $\langle C, S \rangle$ apply the policy $h_{\langle C, S \rangle}$ in the tuple. Different policies h can co-exist across types: *e.g.*, the seed may implement a random novel policy, while peers implement priority rarest first. Contrary to random novel, the RF and PRF policies depend on the system state \mathbf{n} , and require knowledge of a global property; as such, they are harder to implement in a distributed fashion. In a centralized setting, which includes most present BitTorrent implementations, the availability is monitored by a distinguished peer called the swarm tracker. Alternatively, distributed

techniques such as gossiping or sampling can be used to estimate availability. Other definitions are as follows:

DEFINITION 1. Define $\lambda_{total} := \sum_{C: C \in \mathcal{C}} \lambda_C$ to be the total arrival rate. Define $n_S := \sum_{C: C \in \mathcal{C}} n_{\langle C, S \rangle}$ to be the number of peers holding the set of pieces S .

Let $\mathcal{C} := \{C : C \in 2^{\mathcal{F}} \setminus \{\emptyset\}, \lambda_C > 0\}$ be the set of swarms. Define $\mathcal{T} := \{\langle C, S \rangle : C \in \mathcal{C}, S \in 2^{\mathcal{F}} \setminus \{\mathcal{F}\}, C \not\subseteq S\}$ to be the set containing all peer types for peers in the system, and $\tilde{\mathcal{T}} = \mathcal{T} \cup \{\{\perp\}, \mathcal{F}\}$ to be the extended set of peer types. Define $\mathcal{D} = \mathbb{N}^{|\tilde{\mathcal{T}}|}$ to be the set of all possible vectors \mathbf{n} .

Our main stability result (Theorem 2) assumes that the seed applies RN, while peers apply piece selection policies satisfying Assumption 1.

3.1 Markov Process Description and Stability

The system evolution is described by a Markov process $\{\mathbf{n}(t)\}_{t \in \mathbb{R}_+}$ with state space \mathcal{D} . The transition rates of the process depend on how pieces are uploaded.

Assume that the seed implements RN, and peers apply policies as in Assumption 1. Given a state \mathbf{n} , let $T_C(\mathbf{n})$ be the new state resulting from the arrival of a new peer in swarm C . Given $\langle C, S \rangle \in \mathcal{T}$ such that $i \notin S$, and a state \mathbf{n} such that $n_{\langle C, S \rangle} \geq 1$, let $T_{\langle C, S \rangle, i}(\mathbf{n})$ denote the new state resulting from a type $\langle C, S \rangle$ peer downloading piece i . The positive entries of the generator matrix $Q = (q(\mathbf{n}, \mathbf{n}') : \mathbf{n}, \mathbf{n}' \in \mathcal{D})$ are given by:

$$q(\mathbf{n}, T_C(\mathbf{n})) = \lambda_C$$

$$q(\mathbf{n}, T_{\langle C, S \rangle, i}(\mathbf{n})) = \frac{n_{\langle C, S \rangle}}{n} \left[\frac{U}{K - |S|} + \mu \sum_{v \in \mathcal{T}} n_v h_v(i, \langle C, S \rangle, \mathbf{n}) \right]$$

We follow the usual definitions of stability and instability for Markov processes [2]:

DEFINITION 2. The system is unstable if it is transient and the number of peers converges to infinity with probability one; and the system is stable if it is positive recurrent and it has a finite mean number of peers in equilibrium.

4. SINGLE SWARM AND MISSING PIECE SYNDROME

Hajek and Zhu study the above system in the case of a single, autonomous swarm, *i.e.*, a system in which all peers are interested in downloading the file $C = \mathcal{F}$. They determine the stability region under the RN policy:

THEOREM 1 (HAJEK AND ZHU [3]). Consider a single swarm of peers requesting all pieces in \mathcal{F} , in which both the seed and peers follow the random novel piece selection policy. The system is stable if $\lambda_{\mathcal{F}} < U$, and unstable if $\lambda_{\mathcal{F}} > U$.

Hajek and Zhu further establish that the so-called missing piece syndrome [3, 7] is the reason of instability when $\lambda_{\mathcal{F}} > U$. This phenomenon arises when there is a large number of peers in the system that store all pieces in \mathcal{F} except for one missing piece. When this set of peers, termed the *one-club*, is large enough, most of the contacts of new peers arriving in the system will be with such peers. The new peers thus quickly retrieve all pieces except the missing piece, thus joining the one-club set. Since peers holding the missing piece are few, departures from the one-club are mostly due to uploads by the seed; as a result, the departure rate of the one-club is close to the seed upload rate U .

Since $\lambda_{\mathcal{F}} > U$, the rate of growth of peers in the one-club is positive, causing the size of this set to increase to infinity and resulting to instability.

Theorem 1 has an immediate corollary in the case of multi-swarm systems. In particular, suppose that each swarm operates in an *autonomous mode*, independently and in isolation of other swarms. More specifically, peers in swarm $C \in \mathcal{C}$ contact and exchange pieces only with other peers in the same swarm. In addition, the seed divides its upload capacity across different swarms (possibly unevenly), serving each with an appropriate fraction of its total capacity. Finally, peers in swarm C store and exchange only pieces in set C . Theorem 1 directly applies to each such system and, thus, it is easy to verify the following:

COROLLARY 1. Consider a multi-swarm system operating in autonomous mode. The seed can allocate its upload capacity so that the system is stable if $\sum_{C \in \mathcal{C}} \lambda_C < U$; the system is unstable for all allocations of the seed's upload capacity if $\sum_{C \in \mathcal{C}} \lambda_C > U$.

Note that the corollary assumes a static allocation of a seed's rate to each swarm. Though formally studying the capacity of a dynamic allocation is beyond the scope of this work, we observe through simulations that an allocation that is proportional to the size of each swarm does not improve stability or prevent the missing piece syndrome (see Figure 1(a)).

In what follows, we study multi-swarm systems that operate as described formally in Section 3. To distinguish this system, we refer to it as a multi-swarm system in *universal mode* or, simply, a *universal swarm*.

5. STABILITY REGION OF A UNIVERSAL SWARM

The following is our main result regarding the stability region of the Markov process defined in Section 3.1:

THEOREM 2. If the seed implements RN and peers implement work-conserving piece selection policies in Assumption 1, the system is

- i) unstable if $\max_{i: i \in \mathcal{F}} \sum_{C: i \in C} \lambda_C > U$,
- ii) stable if $\max_{i: i \in \mathcal{F}} \sum_{C: i \in C} \lambda_C < U$.

Beyond considering universal swarms, Theorem 2 extends Theorem 1 to the case where peers implement arbitrary piece selection policies. In fact, an immediate corollary of Theorem 2, applied to the single swarm setup, is that the stability region of Theorem 1 extends to such policies as well. The theorem assumes that the seed uses the RN policy; our numerical evaluations in Section 7 suggest that other seed policies (*e.g.*, RF) also yield the same stability region.

Note that the theorem implies that bundling swarms together yields a significant increase in the stability region. Observe that, when the files $C \in \mathcal{C}$ are disjoint, Theorem 2(ii) becomes $\max_{C \in \mathcal{C}} \lambda_C < U$. This defines a larger stability region than the one of autonomous mode, given by Corollary 1. In particular, by bundling swarms together, the stability region scales extremely well as the number of swarms increases: a single seed can support an *unbounded number* of swarms with constant arrival rate, with no effect on the stability region! However, as discussed in Section 7, bundling swarms together comes at the cost of increased delays. Hence, the number of swarms cannot be arbitrarily

large in practice. In Section 7.4, we address this by proposing a *hybrid system* that, by alternating between the universal and autonomous mode, maintains the same stability region as a universal swarm while also ensuring small delays for large numbers of swarms.

The stability region is insensitive to peer piece selection policies. However, policies can be quite different w.r.t. other performance metrics. For example, as we will see in Section 7, policies can differ drastically in how quickly they stabilize the system when operating within the stability region, as well as in how quickly the missing piece syndrome manifests under the condition of Theorem 2(ii).

6. LYAPUNOV FUNCTION FOR PROOF OF Theorem 2(ii)

We provide the detailed proof of Theorem 2(i) in Appendix A. In this section we describe the Lyapunov function applied for the proof of Theorem 2(ii). We establish the stability by Foster’s criterion [1]. Specifically, we use the result below, itself a standard version of the criterion:

LEMMA 1. *The Markov process \mathbf{n} is stable, if there exist a non-negative function $W(\mathbf{n})$ on space \mathcal{D} , with $\{\mathbf{n} : W(\mathbf{n}) \leq c\}$ a finite set for any $c \geq 0$, and constants $n_o \geq 0, \xi > 0$, such that $\forall n \geq n_o, QW \leq -\xi n < 0$, where Q denotes the infinitesimal generator of the process. We then say that W is a valid Lyapunov function.*

We devote the remainder of this section to the description of our Lyapunov function, and show in Appendix B that it satisfies the above Foster criterion.

DEFINITION 3. $\Delta := U - \max_{i:i \in \mathcal{F}} \sum_{C:i \in C} \lambda_C > 0$.

DEFINITION 4. $\mathcal{E}_C := \{\langle C', S' \rangle : C' \not\subseteq C, S' \subseteq C\}$ is set of types of peers which may hold the set of pieces C in the future, $E_C := \sum_{\langle C', S' \rangle \in \mathcal{E}_C} n_{\langle C', S' \rangle}$ is the number of peers with types in \mathcal{E}_C :

DEFINITION 5. $\mathcal{H}_C := \bigcup_{i \in \mathcal{F} \setminus C} \{\langle C', C \cup \{i\} \rangle : C' \not\subseteq C \cup \{i\}\}$ is set of types of peers which hold one more piece than C , $H_C := \sum_{\langle C', S' \rangle \in \mathcal{H}_C} n_{\langle C', S' \rangle}$ is the number of peers with types in \mathcal{H}_C .

DEFINITION 6. $r \in (0, \frac{1}{2})$, $d \in (1, \infty)$, $\beta \in (0, \frac{1}{2})$, $\alpha := 8K(\lambda_{total} + \Delta/2)/U$ are four constants, and

$$\phi(x) := \begin{cases} (2d + \frac{1}{2\beta} - x) & 0 \leq x \leq 2d \\ \frac{\beta}{2}(x - 2d - \frac{1}{\beta})^2 & 2d < x \leq 2d + \frac{1}{\beta} \\ 0 & x > 2d + \frac{1}{\beta} \end{cases}$$

DEFINITION 7. *The Lyapunov function W is defined as $W := \sum_C r^{|C|} T_C$, where*

$$T_C := \begin{cases} \frac{1}{2} E_C^2 + \alpha E_C \phi(H_C), & |C| \leq K - 2, \\ \frac{1}{2} E_C^2, & |C| = K - 1 \end{cases}$$

Function ϕ is affine in the range $[0, 2d]$, and thus the Lyapunov function W is quadratic in the state variable \mathbf{n} in a corresponding range. However, we need to go beyond quadratic functions in order to apply the Foster criterion. We show in Appendix B that for suitable choices of constants r, d, β the function W satisfies Foster’s criterion.

7. EVALUATION

In this section, we evaluate the performance of universal swarms through simulations, studying swarm behavior for different piece selection policies, as well as the dependence of the sojourn time in swarm parameters. In what follows, we denote by RF, RN, PRN and PRF the piece selection policies rarest first, random novel, priority random novel and priority rarest first. Note that PRF and PRN reduce to RF and RN when the system operates in autonomous mode.

7.1 Autonomous vs. Universal Swarms

We begin by validating Theorems 1 and 2. We do so by studying the evolution of the system size n in autonomous and universal mode for a system comprising 3 swarms, each requesting a different 3-piece set. The seed rate is $U = 3.1$ and the arrival rate in each swarm is $\lambda = 3.0$; note that Theorems 1 and 2 imply that the autonomous mode is unstable while the universal system is stable, in this regime.

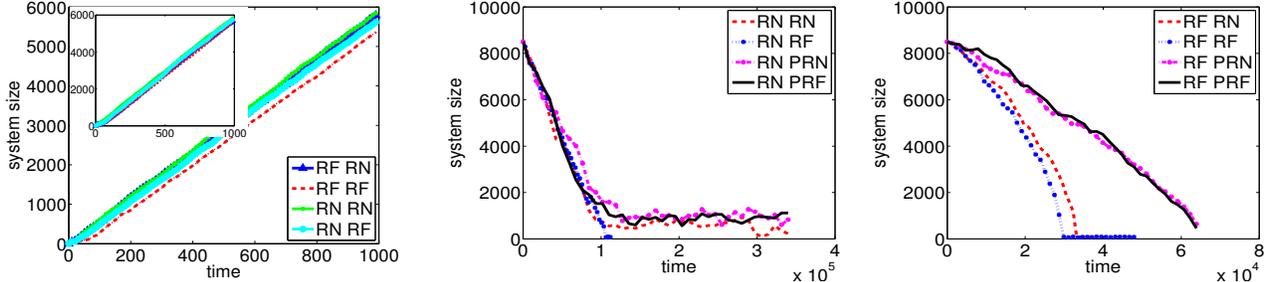
Figure 1(a) shows the evolution of the system size in autonomous mode, when the seed statically allocates 1/3 of its upload rate to each swarm, for different combinations of policies at the seed and the peers. All simulations start from an empty system. Even though applying RF at both the seed and the peers leads to a slightly smaller system size, the missing piece syndrome manifests in all four cases. We repeat these experiments with the seed allocating its rate dynamically, so that each swarm receives pieces at a rate proportional to its size. The results (inset of Figure 1(a)) show that instability persists in this setup too.

We repeat these experiments in universal mode, starting the system from an initial state comprising 8500 peers forming a *one-club*: all peers belong to the same swarm, and store in their cache all nine pieces *except for one common piece they request*. Figure 1(b) shows the system evolution when the seed applies RN; indeed the system stabilizes after 10^5 time units, confirming Theorem 2. The system stabilizes faster (in the order of 10^4 time units) when the seed applies RF, as seen in Figure 1(c), with RF at both seed and peers stabilizing the system the fastest (in roughly $3 \cdot 10^4$ time units). Interestingly, prioritizing pieces at peers (through either PRN or PRF) leads to *slower* stabilization: this is precisely because these policies reduce piece diversity.

7.2 Instability and Meta-Stability in Universal Mode

We next consider the same experiments as above with $U = 2.9$. As the arrival rate at each swarm is $\lambda = 3$, Theorem 2 stipulates that when the seed applies RN the system is unstable. A question we wish to address here is how quickly the missing piece syndrome manifests in this case, depending on the piece selection policies.

To evaluate this, we conduct the following experiments. We start our simulations with initial system size n_0 , where the initial state comprises all peers forming a one-club (*i.e.*, storing all pieces but one). We then terminate the simulation when either the system size increases to the threshold $\max(2000 + n_0, 2n_0)$ or the simulation time reaches 10^7 , whichever occurs first. We first conduct this experiment with an empty initial state $n_0 = 0$; if the experiment does not reach the threshold in 10^7 units, we increase n_0 by 100 and repeat the experiment. This way, we identify the *critical one-club size*: if the system reaches a state with a one-club of that size, it becomes unstable.



(a) Autonomous mode, static and dynamic (inset) allocation.

(b) Universal mode, RN at seed.

(c) Universal mode, RF at seed.

Figure 1: System size VS time. (“RN RF” means RN at seed and RF at peers, other legends follow similarly.)

Table 1: Critical One Club Size

Policy		Critical	Final	Final One	Sim.
Seed	Peer	n_0	Size	Club Ratio	Duration
RN	PRN	0	2000	95.6%	13181
	PRF	0	2000	95.4%	17211
	RN	0	2000	93.3%	13603
	RF	500	2500	94.7%	22655
RF	PRN	2100	4200	98.1%	74655
	PRF	2000	4000	98.0%	51415
	RN	8000	16000	99.4%	283197
	RF	8100	16200	99.4%	323738

Our simulation results for the case where the seed applies RN are summarized in the top half of Table 1. We see that the missing piece syndrome indeed manifests at the critical initial conditions, with the one-club comprising more than 90% of the peer population at termination time. When peers use any policy other than RF, the critical one-club size is 0. In contrast, when peers use RF, the syndrome manifests only when $n_0 = 500$; indeed, using RF improves the diversity of pieces in the system, which in turns makes reaching a critical one-club size more difficult. This behavior becomes even more striking when the seed uses RF: as shown in the bottom half of Table 1, piece diversity is so high that critical one-club sizes lie between 2 and 8 thousand peers.

Crucially, in all simulations starting from an initial size below the critical value, we observe interestingly that the system size actually *decreases* to a size below 200 and lingers around this value for the entire 10^7 time units! This implies that, though the system is clearly not stable in any of the cases in Table 1, applying RF at the seed or peers yields *meta-stability*: although there exists a critical one-club size, its value is so high that it is quite hard to reach from the “typical” size at which the system operates most of the time (~ 200 peers in our simulations). A natural question to ask is what is the critical value n_0 , as well as what is the “typical” size at which a meta-stable system operates most of the time; we revisit these questions below, in Section 7.4.

7.3 Average Sojourn Time

We next turn our attention to the sojourn time. First we study a universal system comprising 3 swarms with 3 pieces each. The seed rate is fixed at $U = 3.0$ and the swarm arrival

rate varies as $\lambda = U(1 - \frac{1}{2^i})$, for $i = 1, \dots, 10$, remaining within the stability region but approaching U from below. In Figure 2(a), we plot the average sojourn time for different piece selection policies as a function of $1/(U - \lambda)$. We observe that, as λ approaches U , the sojourn time under the RN policy at the seed increases considerably, with the exception of the RN-RF case, *i.e.*, when peers use RF. In all four cases for which the seed uses RF, the sojourn time remains practically constant as λ approaches U . This is consistent with the fact that, by meta-stability, when the seed uses RF the system size remains small most of the time even if $\lambda > U$; as such, there is no sharp increase in the sojourn time as λ approaches U from below.

We next study how the sojourn time scales with the number of swarms L . In Figure 2(b), we plot the average sojourn time vs L for the case where each swarm comprises peers requesting a k -piece file, for $k \in \{10, 30, 60\}$. The total number of pieces is $K = kL$. Across all values of k , the average sojourn time increases linearly as L increases. Similarly, the sojourn time also increases proportionally to k . Thus, the increased stability offered by bundling swarms together comes at the cost of increased delays; we address this in the next section by showing that delays can be suppressed for a wide range values of L by using a *hybrid* approach, alternating between universal and autonomous mode.

7.4 Stable, Low Sojourn Universal Swarms

Our simulations suggest that, in a meta-stable swarm, there are two important system sizes: the *operating* size n_{op} , which is the size around which the system stays most of the time, and the *critical* size n_0 , which is the size of a one-club that, once attained, leads the system to instability. If the two sizes are sufficiently far apart from each other, the system will exhibit meta-stability: when $N \approx n_{op}$, it will take a long time for N to reach n_0 , from which the missing piece syndrome manifests.

Calculating exactly n_{op} and n_0 is quite challenging, and is beyond the scope of this paper. Nevertheless, in Appendix C, we derive some simple estimates of n_{op} and n_0 when (a) the system comprises of a single swarm, (b) $\lambda > U$, and (c) both the seed and peers use RF. In particular, we show that in such a single swarm system:

$$n_{op} \approx \frac{\lambda K}{\mu}, \text{ and } n_0 \approx \frac{\lambda(K-1)}{\mu} \left[\frac{1}{2} \frac{U(K-1)}{\lambda-U} - 1 \right]. \quad (1)$$

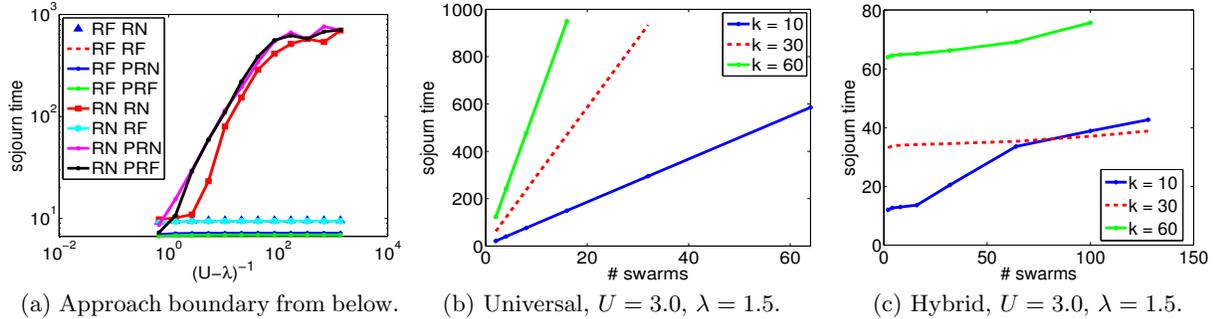


Figure 2: Average sojourn time for universal swarms.

Using these two estimates, we propose a *hybrid system* that attains the increased stability region of the universal swarm, while also ensuring that the sojourn times remain small for a wide range of values of L . The hybrid system alternates between the autonomous mode, whereby swarms operate in isolation while sharing a U/L portion of the seed’s capacity, and the universal mode, where swarms are bundled together. In particular, consider a system with L swarms, each requesting a file of $k = K/L$ pieces. The system alternates between the two modes according to the following rules: (a) If in autonomous mode, the system switches to universal mode if any single swarm has size $\geq n_{op} + \max(n_0, 2n_{op})$; (b) if in universal mode, the system switches to autonomous mode if each piece requested by a swarm is held by at least $\max(n_{op}/10, 1)$ peers within the swarm. Values n_{op}, n_0 are computed by (1), assuming an upload rate U/L and a number of pieces k . Intuitively, the universal mode is applied when there is strong evidence that the missing piece syndrome is manifesting, as the swarm size becomes greater than $n_{op} + n_0$. The system reverts to an autonomous mode when there is enough diversity in each swarm—each piece is held by at least the one tenth of the peer population.

The hybrid system switches to the universal mode when the system size becomes large, so it exhibits the increased stability region of universal swarms described in Theorem 2. Figure 2(c) shows the sojourn time of a hybrid system as L increases. In contrast to Figure 2(b), for $k = 30$ and $k = 60$, the sojourn time stays close to the value attained when $L = 1$ (~ 33 and ~ 64 time units, respectively). For $k = 10$, the sojourn time starts increasing linearly after $L = 12$.

These improved sojourn times appear precisely because of the meta-stability. Indeed, swarms operate fine most of the time without the intervention of other swarms, and this is why they experience the same delay as if $L = 1$. As $U/L < \lambda$, the autonomous mode is unstable; however, at the few (and rare) occasions when the missing piece syndrome manifests, bundling swarms together ensures the system quickly stabilizes and reverts to its operating size.

The knee observed for $k = 10$ suggests that this behavior cannot be sustained for arbitrarily large L . Equation (1) can help us give an approximate answer to how large L can be. Indeed, for the system to be meta-stable, the critical one-club size must be significantly larger than the operating size. Requiring that $n_0 > 2n_{op}$, so that the missing piece syndrome rarely manifests, and taking $K/(K-1) \approx 1$, gives the following heuristic for metastability when $L=1$: $K \frac{U}{\lambda-U} > 6$.

Considering now $L > 1$ swarms in autonomous mode, each requesting $k = K/L$ pieces. Each swarm gets a U/L upload rate in autonomous mode. Then, the above condition becomes $L < \frac{U}{6(\lambda-\frac{U}{L})}k \approx \frac{U}{6\lambda}k$. In other words, the hybrid system can support a number of swarms L with small delay so long as L is of the order of k , the number of pieces in each swarm. As the number of pieces in a file typically numbers in the thousands, this implies that the above system can sustain low sojourn times for a large number of swarms.

8. CONCLUSION

In this paper, we formally characterized the stability region of universal swarms. Our simulations reveal an interesting relationship between stability and piece selection policies. Though piece selection policies may share the same stability region they can differ in their ability to resist the missing piece syndrome. This intuition helps us design a hybrid system that achieves simultaneously a large stability region and low sojourn times; establishing its properties analytically, and investigating other hybrid designs, remains an open question.

9. ACKNOWLEDGMENTS

This work was supported by the National Science Foundation under Grant NSF CCF 10-16959, and was partially funded by the European Commission under the FIRE SCAMPI (FP7- IST-258414) project.

10. REFERENCES

- [1] F. Foster. On the stochastic matrices associated with certain queuing processes. *The Annals of Mathematical Statistics*, 24(3):355–360, 1953.
- [2] B. Hajek. An exploration of random processes for engineers. December 20, 2011.
- [3] B. Hajek and J. Zhu. The missing piece syndrome in peer-to-peer communication. *Stochastic Systems*, 1(2):246–273, 2011.
- [4] J. Han, T. Chung, S. Kim, T. T. Kwon, H.-c. Kim, and Y. Choi. How prevalent is content bundling in bittorrent. In *ACM SIGMETRICS*, 2011.
- [5] L. Massoulié and M. Vojnović. Coupon replication systems. In *ACM SIGMETRICS*, 2005.
- [6] D. S. Menasche, A. A. Rocha, B. Li, D. Towsley, and A. Venkataramani. Content availability and bundling in swarming systems. In *ACM CoNEXT*, 2009.
- [7] I. Norros and H. Reitu. On the stability of two-chunk file sharing systems. *Queueing Systems*, 3, 2011.
- [8] B. Oğuz, V. Anantharam, and I. Norros. Stable, distributed p2p protocols based on random peer sampling. In *IEEE Allerton*, pages 915–919, 2012.

- [9] D. Wu, Y. Liu, and K. W. Ross. Queuing network models for multi-channel P2P live streaming systems. In *IEEE INFOCOM*, 2009.
- [10] X. Zhou, S. Ioannidis, and L. Massoulié. On the stability and optimality of universal swarms. In *ACM SIGMETRICS*, 2011.
- [11] J. Zhu and B. Hajek. Stability of a peer-to-peer communication system. *IEEE Transactions on Information Theory*, 2012.

APPENDIX

A. PROOF OF Theorem 2(i)

If $K = |\mathcal{F}| = 1$, Theorem 2 follows because \mathbf{n} is an $M/M/1$ queue, with arrival rate $\lambda_{\mathcal{F}}$ and service rate U .

Assume that $K \geq 2$ and that, *w.l.o.g.*, $\lambda_1 := \sum_{C:1 \in C} \lambda_C > U$. Notice that \mathbf{n} is irreducible; as such, to show it is transient, it suffices to show that there exists a transient state [2, Proposition 6.3.5]. We show that the initial state where many peers are missing piece one is a transient state: starting from this state, the number of peers converges to infinity with a positive probability. Transience directly implies that the number of peers converges to infinity with probability one, as (a) there is a finite number of states where the number of peers is bounded by a constant, and (b) the probability for \mathbf{n} to stay in any finite set is zero. To construct a transient initial state, we assume:

ASSUMPTION 2. *Select positive values $\epsilon, \xi, \rho, \epsilon_o, B, N_o$ so that*

$$3\epsilon < \lambda_1 - U, \xi < 0.5, \epsilon > 4K\xi U, \quad (2)$$

$$\rho := 2\xi(K-1) < 0.5, \epsilon_o < \xi(\lambda_1 - U - 3\epsilon), \quad (3)$$

$$e^{\lambda_{total}[2^{(K-1)/\mu+1}]2^{-B}} \leq 0.1(1 - 2^{-\epsilon_o}), \quad (4)$$

$$64K^2\xi U \leq 0.2B(\epsilon - 4K\xi U), \quad (5)$$

$$\lambda_{total} \leq 0.2B\epsilon, U \leq 0.2B\epsilon, \quad (6)$$

$$N_o > 6B, B+1 < \xi(N_o - 3B - 1). \quad (7)$$

We partition the set of peers in two classes: *one-club peers*, which are the peers having all pieces in $\mathcal{F} \setminus \{1\}$, and *young peers*, which are peers missing at least one piece from $\mathcal{F} \setminus \{1\}$. We also refer to peers holding piece 1 as *infected peers*; note that infected peers are necessarily young.

DEFINITION 8. *Define the following random processes:*

- A_t : cumulative number of arrivals of peers wanting to download piece one, up to time t
- N_t : number of peers at time t
- Y_t : number of young peers at time t
- D_t : cumulative number of uploads of piece one by infected peers, up to time t
- Z_t : cumulative number of uploads of piece one by the seed, up to time t

We construct the following initial state:

ASSUMPTION 3. *At $t = 0$, $N = N_o$ and all N peers are one-club peers.*

Let τ be the extended stopping time defined by $\tau = \min\{t \geq 0 : Y_t \geq \xi N_t\}$, with the usual convention that $\tau = \infty$ if $Y_t < \xi N_t$ for all t .

LEMMA 2. *Under Assumption 2 and Assumption 3,*

$$P\{A_t > -B + (\lambda_1 - \epsilon)t \text{ for all } t \in [0, \tau]\} \geq 0.9, \quad (8)$$

$$P\{Z_t < B + (U + \epsilon)t \text{ for all } t \in [0, \tau]\} \geq 0.9, \quad (9)$$

$$P\{Y_t < B + \epsilon_o t \text{ for all } t \in [0, \tau]\} \geq 0.9, \quad (10)$$

$$P\{D_t < B + \epsilon t \text{ for all } t \in [0, \tau]\} \geq 0.9. \quad (11)$$

The proof of Lemma 2 is given after this section. Lemma 2 implies Lemma 3, thereby Theorem 2(i) follows because the state in Assumption 3 is a transient state.

LEMMA 3. $P\{\tau = \infty \text{ and } \lim_{t \rightarrow \infty} N_t = +\infty\} \geq 0.6$

PROOF. let \mathcal{Z} be the intersection of the four events on the left sides of (8)-(11). We have $P(\mathcal{Z}) \geq 0.6$. Note that $N_o + A_t - D_t - Z_t$ is no larger than the number of peers wanting to download piece one at t . So, on \mathcal{Z} ,

$$N_t \geq N_o + A_t - D_t - Z_t > N_o - 3B + (\lambda_1 - U - 3\epsilon)t$$

for all $t \in [0, \tau]$. We claim that τ cannot be finite on \mathcal{Z} . Otherwise, at time τ ,

$$\begin{aligned} \xi &\leq \frac{Y_\tau}{N_\tau} < \frac{B+1 + \epsilon_o\tau}{N_o - 3B - 1 + (\lambda_1 - U - 3\epsilon)\tau} \\ &\leq \max\left\{\frac{B+1}{N_o - 3B - 1}, \frac{\epsilon_o}{\lambda_1 - U - 3\epsilon}\right\} < \xi, \end{aligned}$$

a contradiction. Thus, on \mathcal{Z} we have $\tau = \infty$. So $N_t \rightarrow \infty$, and $\forall t, Y_t/N_t < \xi$. Thus, \mathcal{Z} is a subset of the event on the left side of Lemma 3, and so Lemma 3 follows. \square

A.1 Proof of Lemma 2

Define an alternative process to be a process the same as the original process \mathbf{n} , but it terminates at time $\tau = \min\{t \geq 0 : Y_t \geq \xi N_t\}$. It is sufficient to prove for the alternative process. In the alternative process, A is a Poisson process with rate λ_1 , and Z is stochastically dominated by a Poisson process with rate U . Thus, both (8) and (9) follow from (6) and the consequence of Kingman's moment bound below.

PROPOSITION 1 (COROLLARY 6.1 IN [3]). *Let C be a compound Poisson process with $C_0 = 0$, with jump times given by a Poisson process of rate α , and jump sizes having mean m_1 and mean square value m_2 . Then for all $B > 0$ and $\epsilon > \alpha m_1$, $P\{C_t < B + \epsilon t \text{ for all } t\} \geq 1 - \frac{\alpha m_2}{2B(\epsilon - \alpha m_1)}$.*

The inequality (10) follows from (4), Lemmas 4 and 5:

LEMMA 4 (LEMMA 6.2 IN [3]). *Let M denote the number of customers in an $M/GI/\infty$ queue, with arrival rate λ and mean service time m . Suppose that $M_0 = 0, B, \epsilon > 0$, then $P(M_t \leq B + \epsilon t \text{ for all } t) \geq 1 - \frac{\epsilon \lambda^{(m+1)} 2^{-B}}{1 - 2^{-\epsilon}}$.*

LEMMA 5. *The process Y is stochastically dominated by the number of customers in one $M/GI/\infty$ queue initially empty, arrival rate λ_{total} , and service times $\sim \text{Gamma}(K-1, 2/\mu)$.*

PROOF. We construct one $M/GI/\infty$ queue on the same probability space as the alternative system, so that $Y_t \leq$ the number of peers in the $M/GI/\infty$ queue. Let the $M/GI/\infty$ queue have the same arrival process as the alternative system—*i.e.*, a Poisson process of rate λ_{total} . As $\xi < 0.5$, for any young peer, the intensity of downloads from the one-club is always greater than or equal to $\mu/2$ for the alternative system. Suppose thus that each young peer has an internal Poisson clock, which ticks at rate $\mu/2$, and is such that whenever the internal clock of a young peer ticks, that young peer downloads a piece from the one-club. We declare that a peer remains in the $M/GI/\infty$ system until its internal clock ticks $K-1$ times. This gives the desired service time distribution, and the service times of different peers in the $M/GI/\infty$ are independent. A young peer may leave the group of young

Table 2: Specification of comparison system

Alternative system	Comparison system
The seed creates infected peers at a rate $\leq \xi U_s$.	The seed creates infected peers at rate ξU_s .
An infected peer creates infected peers at rate $\leq \xi \mu$.	An infected peer creates infected peers at rate $\xi \mu$.
An infected peer uploads piece one to one-club at a rate $\leq \mu$.	An infected peer uploads piece one to one-club at rate μ .
Any infected peer needs at most $K - 1$ additional pieces, at rate $\geq \mu/2$.	A new infected peer must get $K - 1$ additional pieces, at rate $\mu/2$.

peers (depart or join the one-club) sooner than it leaves the $M/GI/\infty$ system, because a young peer in the alternative system can possibly download pieces at times when its internal clock doesn't tick. But if a peer is still a young peer in the alternative system, it is in the $M/GI/\infty$ system. \square

To prove (11), consider the stochastic system described in Table 2, which we call the *comparison system*. It should be clear to the reader that both the alternative system and the comparison system can be constructed on the same underlying probability space; in particular, such a construction can be done so that any infected peer in the alternative system at a given time is also in the comparison system. To enforce this, when a peer becomes infected in the alternative system, we require that (a) it also arrives to the comparison system, (b) it discards all pieces it may have downloaded before becoming infected, and (c) it subsequently ignores all opportunities to download except those occurring at the times its internal Poisson clock with ticking rate $\mu/2$ ticks. Because infected young peers may stay longer in the comparison system than in the alternative system, some of the peers in the comparison system correspond to peers that already departed from the alternative system. There can also be some infected peers in the comparison system that never existed in the alternative system because the higher arrival rate of infected peers in the comparison system.

In all cases, any infected peer in the alternative system is also in the comparison system. That is, any of the following events occurring in the alternative system also occurs in the comparison system: (a) the seed creates an infected peer, (b) an infected peer creates an infected peer, and (c) an infected peer replicates piece one to a one-club peer. Events of types (b) and (c) correspond to the two possible ways that infected peers can upload piece one. Therefore, this property implies the Lemma 6, where \widehat{D} is the cumulative number of uploads of piece one by infected peers, up to time t , in the comparison system.

LEMMA 6. *The process $(D_t : t \geq 0)$ is stochastically dominated by $(\widehat{D}_t : t \geq 0)$.*

LEMMA 7. *$(\widehat{D}_t : t \geq 0)$ can be stochastically dominated by a compound Poisson process $(\widetilde{D}_t : t \geq 0)$, with arrival rate of batches $= \xi U$, and first and second moments of batch sizes bounded by $4K$ and $64K^2$, respectively*

PROOF. Identify two kinds of infected peers in the comparison system—the *root peers*, which are those created by the seed, and the infected peers created by other infected

peers. Assume each root peer signs uniquely on its piece one received from the seed, and the signature is inherited by all copies of piece one replicated from the root peer. Partition the uploads of piece one according to their signatures. Let $(\widetilde{D}_t : t \geq 0)$ denote a new process which results when all of the uploads of piece one signed by a root peer (in the comparison system) are counted at the arrival time of the root peer. Since \widetilde{D} counts the same events as \widehat{D} , but does so earlier, $\widehat{D}_t \leq \widetilde{D}_t$ for all $t \geq 0$. It is sufficient to prove (11) with D replaced by \widetilde{D} .

The random process \widetilde{D} is a compound Poisson process. Jumps occur at the arrival times of root peers, which form a Poisson process of rate ξU . Let J be the size of the jump of \widetilde{D} associated with a typical root peer. Then $J = J_1 + J_2$, where (a) J_1 is the number of infected peers holding piece one signed by the root peer, not counting the root peer itself, and (b) J_2 is the number of uploads of piece one to the one club by the root peer and these J_1 peers. Consider one $M/GI/1$ queueing system with arrival rate $\xi \mu$ and service times following the distribution of a random variable $\widetilde{X} \sim \text{Gamma}(K - 1, 2/\mu)$. Then the sum of all the times that the root peer and these J_1 peers are in the comparison system is the same as the duration, L , of a busy period of the $M/GI/1$ queue. And J_1 has the same distribution as the number of customers in a busy period of the $M/GI/1$ queue, not counting the customer who started the busy period. Note that ρ in (3) is the load factor for the $M/GI/1$ queue: $\rho = \xi \mu E[\widetilde{X}]$. Apply results on busy periods of $M/GI/1$ queueing system (see [3, Lemma 6.3]) and Assumption 2, follow a similar argument as in [3, Page 259], we have $E[J] = \frac{1 + \mu E[\widetilde{X}]}{1 - \rho} - 1 \leq 4K$, and $E[J^2] \leq 2\{E[J_1^2] + E[J_2^2]\} \leq 64K^2$. Lemma 7 follows. \square

Hence, (11) with D replaced by \widetilde{D} follows from Proposition 1 and (5). Lemma 2 therefore follows.

B. PROOF OF LYAPUNOV FUNCTION VALIDITY

We show here that W in Definition 7 is a valid Lyapunov function. Define:

DEFINITION 9. $M_\phi := 3d + \frac{1}{\beta}$. Note $M_\phi > \max_x \phi(x)$ and $M_\phi > \min\{x : \phi(x) = 0\} + d > 1$.

DEFINITION 10. Define $D_{\langle C, S \rangle}$ to be the transition rate for type $\langle C, S \rangle$ peers to download pieces.

Define $D_S := \sum_{C: C \in \mathcal{C}, C \not\subseteq S} D_{\langle C, S \rangle}$.

Define $D_{total} := \sum_{S: S \subseteq \mathcal{F}} D_S$.

Note that $D_{\langle C, S \rangle}, D_S, \widetilde{D}_{total}$ are functions of the state \mathbf{n} .

The following two lemmas obviously hold:

LEMMA 8. *Function ϕ verifies $\phi'(x) = -1$ for $0 \leq x \leq 2d$ and $\phi'(x) = 0$ for $x \geq 2d + 1/\beta$. And ϕ' increases linearly from -1 to 0 in $[2d, 2d + 1/\beta]$. $\forall x \geq 0, \phi'(x) \in [-1, 0]$.*

LEMMA 9. $D_S \leq U + \mu \min\{n_S, n - n_S\}$, $D_{total} \leq U + n\mu$, $D_S \geq (U + H_S \mu) n_S / n$.

In the proof, we consider the following two classes of states, where σ is to be selected within $\sigma \in (0, 1/2)$. The classes overlap and their union includes every non-zero state:

DEFINITION 11. *Class I is the set of states \mathbf{n} such that there exists $S \subsetneq \mathcal{F}$, and $n_S/n > 1 - \sigma$; Class II is the set of states \mathbf{n} such that there exist $C_1, C_2 \subsetneq \mathcal{F}$, so that, $nc_1/n > \sigma/2^K$ and $nc_2/n > \sigma/2^K$.*

For a specific σ , a state \mathbf{n} can either be Class I or Class II, or both. The main idea of the proof is to show that W is a valid Lyapunov function for an appropriate choice of (r, d, β, σ) . The given parameters of the network, $\lambda_C (C \in \mathcal{C}), U$ and μ , are treated as constants. Functions on the state space may or may not depend on the variables r, d, β and σ . It is convenient to adopt the big theta notation $\Theta(\cdot)$, with the understanding that it is uniform in these variables; this is summarized in the following definitions.

DEFINITION 12. *Given functions f and g on the state space \mathcal{D} , we say $f = \Theta(g)$ if there exist $k_1, k_2, n_0 > 0$, not dependent on (r, d, β, σ) , such that $k_1|g(\mathbf{n})| \leq |f(\mathbf{n})| \leq k_2|g(\mathbf{n})|$ for all \mathbf{n} such that $n > n_0$.*

For example, $2 = \Theta(1)$, $\lambda_{total}n = \Theta(n)$, $d = \Theta(d)$, $1 \leq \Theta(n)$ and $\Theta(n) - \Theta(n)/2 = \Theta(n)$. Notice that d and $\Theta(1)$ cannot be compared. Similarly, we adopt notions of “small enough” and “large enough” that are uniform in (r, d, β, σ) :

DEFINITION 13. *We say that “condition A is true if $x > 0$ is small enough” if there exists a constant $k > 0$, not depending on (r, d, β, σ) , such that A is true for any $x \in (0, k)$. Similarly, we say that “condition A is true if $x > 0$ is large enough” if there exists a constant $k > 0$, not depending on (r, d, β, σ) , such that A is true for any $x \in (k, \infty)$.*

We identify an approximation to the drift of W . Notice that the infinitesimal generator Q is linear, so that $Q(W) = \sum_C r^{|C|} Q(T_C)$, with $Q(T_C) = \frac{1}{2}Q(E_C^2) + \alpha Q(E_C\phi(H_C))$.

DEFINITION 14. *Define $\mathcal{Q}W$, an approximation of $Q(W)$, as $\mathcal{Q}W := \sum_C r^{|C|} \mathcal{Q}T_C$, with $\mathcal{Q}T_C := E_C Q(E_C) + \alpha E_C Q(\phi(H_C))$ if $|C| \leq K - 2$, $\mathcal{Q}T_C := E_C Q(E_C)$ if $|C| = K - 1$.*

Our proof relies on Lemmas 10 to 13, which are conditioned on (r, d, β, σ) , to bound terms in $Q(W)$ and $\mathcal{Q}W$. The proofs of these lemmas are provided after this section.

LEMMA 10. *Bound for the approximation error:*

$$|Q(W) - \mathcal{Q}W| \leq M_\phi(D_{total} + 1)\Theta(1). \quad (12)$$

LEMMA 11. *If d is large enough, $\forall C \subsetneq \mathcal{F}$, $Q(E_C) \leq \Theta(1)$, $Q(\phi(H_C)) \leq M_\phi\Theta(1)$, and $\mathcal{Q}T_C \leq M_\phi\Theta(E_C) \leq M_\phi\Theta(n)$.*

LEMMA 12. *If d is large enough, $\sigma M_\phi, \beta$ are small enough, for any $S \subsetneq \mathcal{F}$ and any nonzero state \mathbf{n} such that $n_S/n > 1 - \sigma$, $\mathcal{Q}T_S \leq -\frac{1}{2}\Delta E_S$.*

Lemmas 11 and 12 imply Lemma 13:

LEMMA 13. *If d is large enough, $\beta, rM_\phi, \sigma M_\phi r^{-K}$ are small enough, (a) on Class I, $\mathcal{Q}W \leq -r^K\Theta(n)$; (b) on Class II, $\mathcal{Q}W \leq -r^K\sigma^3\Theta(n^2) + M_\phi\Theta(n)$.*

Notice that the conditions of Lemmas 10 to 13 are consistent with each other, so we claim

LEMMA 14. *There exists (r, d, β, σ) satisfying all conditions of Lemmas 10 to 13, such that W is a valid Lyapunov function.*

PROOF. On Class I, $D_{total} = D_S + \sum_{C:C \neq S} D_C \leq 2U + 2(n - n_S)\mu \leq 2U + 2\sigma n\mu = \Theta(1) + 2\sigma\Theta(n)$ by Lemma 9. So Lemma 10 implies that on Class I, $|Q(W) - \mathcal{Q}W| \leq \sigma M_\phi\Theta(n) + M_\phi\Theta(1)$. Combined with Lemma 13(a), on Class I, $Q(W) \leq -r^K\Theta(n) + M_\phi\Theta(1)$ if $\sigma M_\phi r^{-K}$ is small enough. On Class II, $D_{total} \leq U + n\mu = \Theta(n)$, so Lemma 10 implies that $|Q(W) - \mathcal{Q}W| \leq M_\phi\Theta(n)$. On Class II, Combined with Lemma 13(b), $Q(W) \leq -r^K\sigma^3\Theta(n^2) + M_\phi\Theta(n)$. Thus, if (r, d, β, σ) satisfies conditions of Lemmas 10 to 13, and $\sigma M_\phi r^{-K}$ is small enough, $\exists \zeta > 0$ such that $Q(W) \leq -\zeta n$ whenever n is larger than some constant. Thus W is a Lyapunov function. \square

We provide proofs for Lemmas 10 to 13 below. Additional definitions are applied:

DEFINITION 15. *Define Γ_{J_1, J_2} for $J_1, J_2 \in \mathcal{T}$, $J_1 \neq J_2$ to be the transition rate for type J_1 peers to become type J_2 peers.*

Define $\Gamma_{\mathcal{X}_1, \mathcal{X}_2} := \sum_{J_1: J_1 \in \mathcal{X}_1} \sum_{J_2: J_2 \in \mathcal{X}_2} \Gamma_{J_1, J_2}$, where $\mathcal{X}_1 \cap \mathcal{X}_2 = \emptyset$.

Define $D_{\mathcal{H}_S} := \sum_{J: J \in \mathcal{H}_S} D_J$, for \mathcal{H}_S in Definition 5.

Define $A_{\mathcal{H}_S}$ to be the total transition rate for peers to join the group of peers with types in \mathcal{H}_S .

Let $\mathcal{P}_S := \{ \langle C, S \rangle : C \not\subseteq S, C \in \mathcal{C} \}$ be the set of types of peers holding piece set S .

Notice that Γ, D, A are all functions of \mathbf{n} .

B.1 Proof of Lemma 10

We compare $Q(W)$ and $\mathcal{Q}W$ term by term. Consider terms $Q(T_C)$ and $\mathcal{Q}T_C$. First, assume $|C| \leq K - 2$. Because α is fixed, we have $|Q(T_C) - \mathcal{Q}T_C| \leq a_1 + \alpha(a_2 + a_3)$,

$$a_1 = \left| \frac{1}{2}Q(E_C^2) - E_C Q(E_C) \right| \leq \lambda_{total} + D_{total},$$

$$a_2 = |Q(E_C\phi(H_C)) - Q(E_C)\phi(H_C) - E_C Q(\phi(H_C))|,$$

$$a_3 = |Q(E_C)\phi(H_C)| \leq M_\phi(\lambda_{total} + D_{total}).$$

The only way E_C and $\phi(H_C)$ can simultaneously change is that some peer with type in \mathcal{E}_C becomes a peer with type in \mathcal{H}_C , causing E_C to decrease by one, and $\phi(H_C)$ to decrease by at most one, so $a_2 \leq \Gamma_{\mathcal{E}_C, \mathcal{H}_C}$. Notice the fact that $\Gamma_{\mathcal{E}_C, \mathcal{H}_C} \leq D_{total}$, we have $\forall |C| \leq K - 2$, $|Q(T_C) - \mathcal{Q}T_C| \leq M_\phi(D_{total} + 1)\Theta(1)$. Secondly, assume $|C| = K - 1$. Then, $|Q(T_C) - \mathcal{Q}T_C| = a_1 \leq \lambda_{total} + D_{total} \leq M_\phi(D_{total} + 1)\Theta(1)$. There are only finitely many terms of T_C in W (2^K in total), and notice that $r < 1$, Lemma 10 follows. \square

B.2 Proof of Lemma 11

The upper bound for the drift of E_C is obvious: $Q(E_C) \leq \lambda_{total} = \Theta(1)$. Next consider $Q(\phi(H_C))$. Because ϕ is a decreasing function, only the rate for H_C to decrease contributes to the positive part in the drift of $\phi(H_C)$, so to consider an upper bound of $Q(\phi(H_C))$ it suffices to consider the rates of transitions which reduce H_C . There is only one way for H_C to decrease: peers with types in \mathcal{H}_C to download a novel piece – with aggregate rate $D_{\mathcal{H}_C}$. Each peer with type in \mathcal{H}_C downloading a novel piece can cause $\phi(H_C)$ to increase at most one. Thus, an upper bound for the drift of $\phi(H_C)$ is $Q(\phi(H_C)) \leq D_{\mathcal{H}_C} \leq U + H_C\mu = \Theta(1) + H_C\Theta(1)$.

We can choose d large enough, i.e. $d > 1$, so $M_\phi > 2d + 1/\beta + 1$. Thus $Q(\phi(H_C))$ vanishes when $H_C > M_\phi$, because $\phi(H_C)$ vanishes when $H_C > 2d + 1/\beta$ and the decreasing of H_C when state changes is bounded by $1 < d$. Hence $Q(\phi(H_C)) \leq M_\phi\Theta(1)$, because $M_\phi > 1$.

Finally, the bound on $\mathcal{Q}T_C$ follows from the other two bounds already proved. \square

B.3 Proof of Lemma 12

Assume that $n_S/n > 1 - \sigma$, where $0 < \sigma < \frac{1}{4}$ is to be specified. We consider two cases. Suppose first $|S| = K - 1$. Peers with cache S only miss one piece $i \in \mathcal{F} \setminus S$, and $\mathcal{Q}T_S = E_S Q(E_S) \leq E_S [\sum_{C:i \in C} \lambda_C - U(1 - \sigma)] \leq -\frac{1}{2} \Delta E_S$, if σ is set to be small enough: $\sigma < \Delta/(2U)$. Lemma 12 follows.

Suppose now that $|S| \leq K - 2$, then by [11, Lemma 19],

$$\begin{aligned} Q(\phi(H_S)) &\leq \phi'(H_S)Q(H_S) + \beta/2(A_{\mathcal{H}_S} + D_{\mathcal{H}_S}) \\ &= \phi'(H_S)(A_{\mathcal{H}_S} - D_{\mathcal{H}_S}) + \beta/2(A_{\mathcal{H}_S} + D_{\mathcal{H}_S}). \end{aligned}$$

Substitute the above inequality into $Q(E_S) + \alpha Q(\phi(H_S))$, which is one component in Definition 14, apply Lemma 9,

$$\begin{aligned} Q(E_S) + \alpha Q(\phi(H_S)) &\leq \bar{\varsigma} - \varpi, \text{ where} \quad (13) \\ \begin{cases} \bar{\varsigma} := \lambda_{total} - \frac{1}{2} D_S + \alpha \phi'(H_S) A_{\mathcal{H}_S} \\ \varpi := \frac{1}{2} D_S + \alpha \phi'(H_S) D_{\mathcal{H}_S} - \alpha \beta/2(A_{\mathcal{H}_S} + D_{\mathcal{H}_S}) \end{cases} \end{aligned}$$

We claim first that $\bar{\varsigma} \leq -\Delta/2$. To prove, suppose first that $H_S < d$: Then, $\phi'(H_S) = -1$ and because the seed applies RN, $A_{\mathcal{H}_S} \geq \Gamma_{\mathcal{P}_S, \mathcal{H}_S} \geq U(1 - \sigma)/K \geq U/(4K)$. So $\bar{\varsigma} \leq \lambda_{total} - \alpha U/(4K) \leq -\Delta/2$. Suppose secondly that $H_S \geq d$: Then, $D_S \geq d\mu(1 - \sigma) \geq \frac{1}{2} d\mu$. So $\bar{\varsigma} \leq \lambda_{total} - \frac{1}{2} d\mu \leq -\frac{1}{2} \Delta$, for d large enough: $d > 2(\lambda_{total} + \Delta/2)/\mu$.

We further claim that $\varpi \geq 0$. To prove, let ω_S be the number of peers holding pieces not in S . For $\sigma < \frac{1}{2}$, we have

$$D_S \geq (U + \omega_S \mu)(1 - \sigma) \geq \frac{1}{2}(U + \omega_S \mu). \quad (14)$$

Notice that pieces novel to peers with types in \mathcal{H}_S are not contained in S . The number of peers that can upload pieces to peers with types in \mathcal{H}_S is no larger than ω_S , so

$$D_{\mathcal{H}_S} \leq (U + \omega_S \mu)\sigma \leq 2\sigma D_S \quad (15)$$

In addition, $A_{\mathcal{H}_S} = \Gamma_{\mathcal{P}_S, \mathcal{H}_S} + \Gamma_{I_S, \mathcal{H}_S}$, where $I_S := \{(C', S') : |S'| = |S|, S' \neq S\}$. The number of peers with types in I_S is no larger than ω_S . Therefore $\Gamma_{I_S, \mathcal{H}_S} \leq U + \omega_S \mu \leq 2D_S$ by (14). And notice $\Gamma_{\mathcal{P}_S, \mathcal{H}_S} \leq D_S$, we have $A_{\mathcal{H}_S} \leq 3D_S$. Combine it with Lemma 9, and (15), we have $\varpi \geq \frac{1}{2} D_S - 2\alpha D_{\mathcal{H}_S} - \frac{\alpha\beta}{2} A_{\mathcal{H}_S} \geq (\frac{1}{2} - 4\alpha\sigma - \frac{3}{2}\alpha\beta) D_S$. We can set $4\alpha\sigma < \frac{1}{4}$, $\frac{3}{2}\alpha\beta < \frac{1}{4}$ so that $\varpi \geq 0$ follows.

Therefore, $\bar{\varsigma} \leq -\frac{1}{2} \Delta$ and $\varpi \geq 0$ imply that, when $n_S/n > 1 - \sigma$ and $|S| \leq K - 2$, $\mathcal{Q}T_S = [Q(E_S) + \alpha Q(\phi(H_S))]E_S \leq -\frac{1}{2} \Delta E_S$ and Lemma 12 follows. \square

B.4 Proof of Lemma 13

First, consider Lemma 13(a). Since there are only finitely many types, we can fix a set $S \subsetneq \mathcal{F}$ and consider the set of Class I states \mathbf{n} for which $n_S/n > 1 - \sigma$. Since $\sigma \in (0, \frac{1}{2})$, $E_S > \frac{1}{2}n$. By Definition 3, $\Delta > 0$. By Lemma 12,

$$\mathcal{Q}T_S \leq -\frac{1}{4} \Delta n = -\Theta(n). \quad (16)$$

Consider two conditions: (a) for type C with $|C| > |S|$, E_C may be larger than E_S . Lemma 11 and (16) imply $r^{|C|} \mathcal{Q}T_C \leq r M_\phi r^{|S|} \Theta(n) < 2^{-K-1} r^{|S|} |\mathcal{Q}T_S|$, if $r M_\phi$ is selected to be small enough; (b) for type C with $|C| \leq |S|$ but $C \neq S$, $E_C \leq \sigma n$. Note that $r^{|C|} r^K \leq r^{|S|}$, Lemma 11 and (16) imply $r^{|C|} \mathcal{Q}T_C \leq r^{|C|} M_\phi \Theta(E_C) \leq \sigma M_\phi r^{-K} r^{|S|} \Theta(n) < 2^{-K-1} r^{|S|} |\mathcal{Q}T_S|$, if $\sigma M_\phi r^{-K}$ is selected to be small enough.

Notice that $E_S \geq n/2$, apply Lemma 12, we have $\mathcal{Q}W = r^{|S|} \mathcal{Q}T_S + \sum_{C:C \neq S} r^{|C|} \mathcal{Q}T_C \leq r^{|S|} \mathcal{Q}T_S + \frac{1}{2} r^{|S|} |\mathcal{Q}T_S| \leq -\frac{1}{8} r^{|S|} \Delta n \leq -r^K \Theta(n)$, which proves Lemma 13(a).

Second, consider Lemma 13(b). Suppose $C_1 \not\subseteq C_2$ and consider the set of Class II states \mathbf{n} such that $n_{C_1}/n > \eta$, $n_{C_2}/n > \eta$, where $\eta = \sigma/2^K$. In such states: $D_{C_2} \geq n_{C_2} n_{C_1} \mu/n \geq \mu \eta^2 n \in \sigma^2 \Theta(n)$. Notice that $E_{C_2} \geq n_{C_2} \geq \eta n$, we have $E_{C_2} Q(E_{C_2}) \leq E_{C_2} (\lambda_{E_{C_2}} - D_{C_2}) \leq -\sigma^3 \Theta(n^2) + \Theta(n)$. Lemma 11 indicates $E_{C_2} Q(\phi(H_{C_2})) \leq M_\phi \Theta(n)$, so

$$\mathcal{Q}T_{C_2} = E_{C_2} Q(E_{C_2}) + \alpha E_{C_2} Q(\phi(H_{C_2})) \leq -\sigma^3 \Theta(n^2) + M_\phi \Theta(n).$$

Obviously the above inequality works for the case $|C_2| = K - 1$ where the ϕ term does not exist too. Apply the above inequalities and the result $\forall C, \mathcal{Q}T_C \leq M_\phi \Theta(n)$ from Lemma 11, we claim that, over the set of all Class II states, $\mathcal{Q}W \leq r^{|C_2|} \mathcal{Q}T_{C_2} + \sum_{C:C \neq C_2} \mathcal{Q}T_C \leq -r^K \sigma^3 \Theta(n^2) + M_\phi \Theta(n)$, which proves Lemma 13(b).

C. OPERATING AND CRITICAL SIZE

Consider a single swarm, where peers arrive with rate λ and wish to download K pieces. Assume that the seed has upload rate $U < \lambda$, so that the system is unstable, and peers have upload rate μ . When the system is in the operating state, we expect that the diversity of pieces is high enough, so that every contact a peer makes leads to a piece download. Under this assumption, as a peer wishes to download K pieces, the expected sojourn time is K/μ . By Little's law, as the arrival rate is λ , an approximation of the operating size is therefore $n_{op} \approx \frac{\lambda K}{\mu}$.

Estimating the critical size of a one-club requires a more involved argument. Suppose that a one-club with size B has formed. For B large, peers outside the one club download pieces from the one club at a rate close to μ . As such, the expected time it takes a new peer to be converted to a one club peer is approximately $(K - 1)/\mu$; hence, the number of *young peers*, *i.e.*, peers outside the one-club, is approximately $\lambda(K - 1)/\mu$. Young peers can become *infected*, *i.e.*, obtain the piece the one-club peers are missing. As the seed uses the RF policy, the rate with which the seed infects young peers is $U \frac{\lambda(K-1)/\mu}{B + \lambda(K-1)/\mu}$. Ignoring the fact that young peers may also infect other young peers, and assuming that an infection occurs at an instant sampled uniformly at random within a young peer's lifetime, each infected peer stays for $(K - 1)/2\mu$ time units before it departs, in expectation. Then, the drift of the one club is roughly $\Delta B = \lambda - U - U \frac{\lambda(K-1)/\mu}{B + \lambda(K-1)/\mu} \frac{(K-1)}{2\mu} \mu$. Requiring the drift to be zero and solving for B , we obtain (1).