

# On Lifetime-Based Node Failure and Stochastic Resilience of Decentralized Peer-to-Peer Networks

Derek Leonard, *Student Member, IEEE*, Zhongmei Yao, *Student Member, IEEE*,  
Vivek Rai, and Dmitri Loguinov, *Member, IEEE*

**Abstract**—To model P2P networks that are commonly faced with high rates of churn and random departure decisions by end-users, this paper investigates the resilience of random graphs to lifetime-based node failure and derives the expected delay before a user is forcefully isolated from the graph and the probability that this occurs within his/her lifetime. Using these metrics, we show that systems with heavy-tailed lifetime distributions are more resilient than those with light-tailed (e.g., exponential) distributions and that for a given average degree,  $k$ -regular graphs exhibit the highest level of fault tolerance. As a practical illustration of our results, each user in a system with  $n = 100$  billion peers, 30-minute average lifetime, and 1-minute node-replacement delay can stay connected to the graph with probability  $1 - 1/n$  using only 9 neighbors. This is in contrast to 37 neighbors required under previous modeling efforts. We finish the paper by observing that many P2P networks are *almost surely* (i.e., with probability  $1 - o(1)$ ) connected if they have no isolated nodes and derive a simple model for the probability that a P2P system partitions under churn.

**Index Terms**—Peer-to-Peer Networks, Stochastic Resilience, Lifetime Node Failure, User Isolation, Network Disconnection.

## I. INTRODUCTION

RESILIENCE of both random graphs [6] and different types of deterministic networks [7], [20] has been a topic of enduring popularity in research literature. A classical problem in this line of study is to understand failure conditions under which the network disconnects and/or starts to offer noticeably lower performance (such as increased routing distance) to its users. To this end, many existing models assume uniformly random edge/node failure and examine the conditions under which each user [38], certain components [6], or the entire graph [7], [25] stay connected after the failure.

Analysis of current P2P networks often involves the same model of uniform, concurrent node failure and varies from single-node isolation [18], [38] to disconnection of the entire graph [3], [13], [16], [28]. However, it is usually unclear how to accurately estimate failure probability  $p$  in real P2P systems<sup>1</sup> such as KaZaA or Gnutella and whether large quantities

Supported by NSF grants CCR-0306246, ANI-0312461, CNS-0519442, and CNS-0434940.

Derek Leonard, Zhongmei Yao, and Dmitri Loguinov are with the Department of Computer Science, Texas A&M University, College Station, TX 77843, USA (e-mail: dleonard@cs.tamu.edu; mayyao@cs.tamu.edu; dmitri@cs.tamu.edu).

Vivek Rai is with the Department of Mathematics and Computer Science, Vrije Universiteit, Amsterdam, The Netherlands (email: vivekr@few.vu.nl).

<sup>1</sup>In the absence of a better estimate, value  $p = 1/2$  is often used for illustration purposes [38].

of users indeed experience simultaneous failure in practice. Additional P2P resilience studies examine the required rate of neighbor replacement to avoid isolation [24], derive the delay before the system recovers from inconsistencies [29], analyze network connectivity assuming the existence of an adversary [12], [34], examine how stabilization delays affect the consistency of Chord's finger table [19], and give specific recommendations for improving routing resilience under user churn [32].

In contrast to the traditional studies above, recent research [5] suggests that realistic models of P2P node failure should consider the inherent behavior of Internet users who join and leave the system *asynchronously* and base their decisions on a combination of complex (often unmeasurable) factors including attention span, browsing habits, and altruistic inclinations [15]. To examine the behavior of such systems, this paper introduces a simple node-failure model based on user lifetimes and studies the resilience of  $k$ -regular P2P networks in which nodes stay online for random periods of time. In this model, each arriving user is assigned a random lifetime  $L_i$  drawn from some distribution  $F(x)$ , which reflects the behavior of the user and represents the duration of his/her services (e.g., forwarding queries, sharing files) to the P2P community.

We begin our analysis with the *passive* lifetime model in which the failed neighbors are not continuously replaced. We observe that even in this case, a large fraction of nodes are able to stay online for their entire lifespan without suffering an isolation. Through this relatively simple model, we also show that depending on the tail-weight of the lifetime distribution, the probability of individual node isolation can be made arbitrarily small *without* increasing node degree.

While the passive model certainly allows P2P networks to evolve as long as the newly arriving nodes replenish enough broken links in the system, a much more resilient approach is to require that each user utilize a special neighbor-recovery strategy that can repair the failed segments of the graph and maintain constant degree at each node. We thus subsequently study the *active* model, where each failed neighbor is replaced by another node after some random search delay. For this scenario, we derive both the expected time to isolation  $E[T]$  and the probability  $\pi$  that this event occurs within the lifetime of a user. We conclude our analysis of the active model by deriving an upper bound on  $\pi$  with heavy-tailed lifetimes and comparing the resilience of  $k$ -regular networks to that of degree-irregular systems.

We finish the paper by bridging the gap between local

resilience (i.e., isolation of a single node from the graph) and global resilience (i.e., network disconnection) and show that tightly-connected structures (such as DHTs and many  $k$ -regular random graphs) partition *with at least one isolated node* with probability  $1 - o(1)$  as the size of the network  $n \rightarrow \infty$ . This result demonstrates that metric  $\pi$  solely determines the probability that an evolving P2P network partitions under churn and that disconnection of such graphs for sufficiently small  $\pi$  almost surely involves a *single* node.

The rest of the paper is organized as follows. Section II introduces the lifetime model and discusses our assumptions. Section III studies the passive model while Sections IV-V analyze the active model. Section VI extends our node isolation results to network partitioning and Section VII concludes the paper.

## II. LIFETIME-BASED NODE FAILURE

In this section, we introduce our model of node failure and explain the assumptions used later in the paper.

### A. Lifetime Model and Paper Overview

In the discussion that follows, we consider  $k$ -regular P2P graphs and analyze the probability that a randomly selected node  $v$  is forced to disconnect from the system because all of its neighbors have simultaneously departed and left it with no way to route within the graph. For each user  $i$  in the system, let  $L_i$  be the amount of time that the user stays in the network searching for content, browsing for information, or providing services to other peers.

It has been observed that the distribution of user lifetimes in real P2P systems is often heavy-tailed (i.e., Pareto) [8], [35], where most users spend minutes per day browsing the network while a handful of other peers exhibit server-like behavior and keep their computers logged in for weeks at a time. To allow arbitrarily small lifetimes, we use a shifted Pareto distribution  $F(x) = 1 - (1 + x/\beta)^{-\alpha}$ ,  $x > 0, \alpha > 1$  to represent heavy-tailed user lifetimes, where scale parameter  $\beta > 0$  can change the mean of the distribution without affecting its range  $(0, \infty]$ . Note that the mean of this distribution  $E[L_i] = \beta/(\alpha - 1)$  is finite only if  $\alpha > 1$ , which we assume holds in the rest of the paper. While our primary goal is the study of human-based P2P systems, we also aim to keep our results universal and applicable to other systems of non-human devices and software agents where the nodes may exhibit non-Pareto distributions of  $L_i$ . Thus, throughout the paper, we allow a variety of additional user lifetimes ranging from heavy-tailed to exponential.

The most basic question a joining user may ask about the resilience of lifetime-based P2P systems is *what is the probability that I can outlive all of my original neighbors?* We call this model “passive” since it does not involve any neighbor replacement and study it in fair detail in the next section. This model arises when the search time  $S$  to find a neighbor replacement is prohibitively high (i.e., significantly above  $E[L_i]$ ) or when peers intentionally do not attempt to repair broken links. If degree  $k$  is sufficiently large, it is intuitively clear that a given node  $v$  is not likely to out-survive

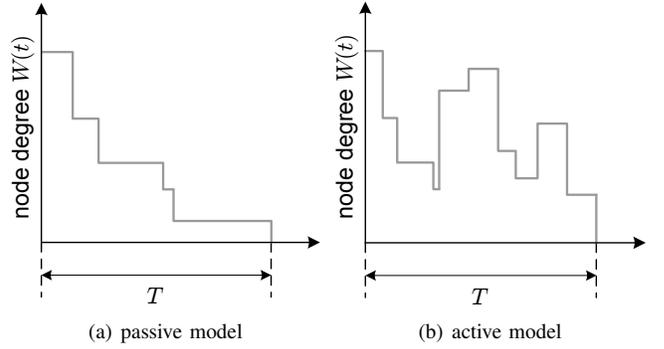


Fig. 1. Degree evolution process leading to isolation under (a) passive and (b) active models.

$k$  other peers; however, it is interesting to observe that Pareto distributions of  $L_i$  make this probability significantly smaller compared to the “baseline” exponential case.

In a later part of the paper, we allow users to randomly (with respect to the lifetime of other peers) search the system for new neighbors once the failure of an existing neighbor is detected. We call this model “active” to contrast the actions of each user  $t$  with those in the passive model. Defining  $W(t)$  to be the degree of  $v$  at time  $t$ , the difference between the passive and active models is demonstrated in Fig. 1, which shows the evolution of  $W(t)$  and the isolation time  $T$  for both models.

### B. Modeling Assumptions

To keep the derivations tractable, we impose the following restrictions on the system. We first assume that  $v$  joins a network that has evolved sufficiently long so as to overcome any transient effects and allow asymptotic results from renewal process theory to hold. This assumption is usually satisfied in practice since P2P systems continuously evolve for hundreds of days or weeks before being restarted (if ever) and the average lifetime  $E[L_i]$  is negligible compared to the age of the whole system when any given node joins it.

Our second modeling assumption requires certain stationarity of lifetime  $L_i$ . This means that users joining the system at different times of the day or month have their lifetimes drawn from the same distribution  $F(x)$ . While it may be argued that users joining late at night browse the network longer (or shorter) than those joining in the morning, our results below can be easily extended to non-stationary environments and used to derive upper/lower bounds on the performance of such systems.

Finally, we should note that these stationarity assumptions do not apply to the number of nodes  $n$ , which we allow to vary with time according to *any* arrival/departure process as long as  $n \gg 1$  stays sufficiently large. We also allow arbitrary routing changes in the graph over time and are not concerned with the routing topology or algorithms used to forward queries. Thus, our analysis is applicable to both structured (i.e., DHTs) and unstructured P2P systems.

## III. PASSIVE LIFETIME MODEL

We start by studying the resilience of dynamic P2P systems under the assumption that users do *not* attempt to replace

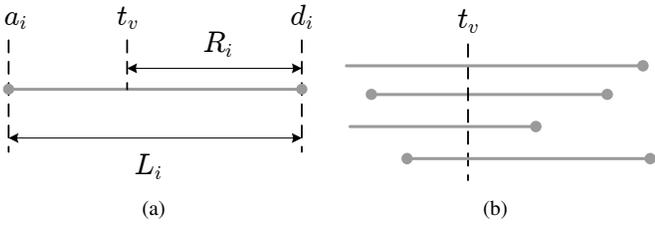


Fig. 2. (a) A neighbor's lifetime  $L_i = d_i - a_i$  and its residual life  $R_i = d_i - t_v$ . (b) Location of  $t_v$  is uniformly random with respect to the lives of all neighbors.

the failed neighbors. As we show below, this analysis can be reduced to basic renewal process theory; however, its application to P2P networks is novel. Due to limited space, we omit the proofs of certain straightforward lemmas and refer the reader to the conference version [23].

### A. Model Basics

We first examine the probability that a node  $v$  can outlive  $k$  randomly selected nodes if all of them joined the system *at the same time*. While the answer to this question is trivial, it provides a lower-bound performance of the system and helps us explain the more advanced results that follow.

*Lemma 1:* The probability that node  $v$  has a larger lifetime than  $k$  randomly selected nodes is  $1/(k+1)$ .

Consider an example of fully-populated Chord [38] with neighbor table size equal to  $\log_2 n$ , where  $n$  is the total number of nodes in the P2P network. Thus, in a system with 1 million nodes, the probability that a randomly selected node outlives  $\log_2 n$  other peers is approximately 4.8%. This implies that with probability 95.2%, a user  $v$  *does not have to replace any of its neighbors to remain online for the desired duration  $L_v$* .

Note, however, that in current P2P networks, it is neither desirable nor possible for a new node  $v$  to pick its neighbors such that their arrival times are exactly the same as  $v$ 's. Thus, when  $v$  joins a P2P system, it typically must randomly select its  $k$  neighbors from the nodes *already present* in the network. These nodes have each been alive for some random amount of time before  $v$ 's arrival, which may or may not affect the remainder of their online presence. In fact, the tail-weight of the distribution of  $L_i$  will determine whether  $v$ 's neighbors are likely to exhibit longer or shorter remaining lives than  $v$  itself.

Throughout the paper, we assume that neighbor selection during join and replacement is *independent* of 1) neighbors' lifetimes  $L_i$  or 2) their current ages  $A_i$ . The first assumption clearly holds in most systems since the nodes themselves do not know how long the user plans to browse the network. Thus, the value of  $L_i$  is generally hard to correlate with any other metric (even under adversarial selection). The second assumption holds in most current DHTs [18], [30], [33], [38] and unstructured graphs [9], [14], [37] since neighbor selection depends on a variety of factors (such as a uniform hashing function of the DHT space [38], random walks [14], interest similarity [37], etc.), none of which are correlated with node age.

The above assumptions allow one to model the time when  $v$  selects each of its  $k$  neighbors to be uniformly random within each neighbor's interval of online presence. This is illustrated in Fig. 2(a), where  $t_v$  is the join time of node  $v$ , and  $a_i$  and  $d_i$  are the arrival and departure times of neighbor  $i$ , respectively. Since the system has evolved for sufficiently long before  $v$  joined, the probability that  $v$  finds neighbor  $i$  at any point within the interval  $[a_i, d_i]$  can be modeled as equally likely. This is schematically shown in Fig. 2(b) for four neighbors of  $v$ , whose intervals  $[a_i, d_i]$  are independent of each other or the value of  $t_v$ .

Next, we formalize the notion of residual lifetimes and examine under what conditions the neighbors are more likely to outlive each joining node  $v$ . Define  $R_i = d_i - t_v$  to be the remaining lifetime of neighbor  $i$  when  $v$  joined the system. As before, let  $F(x)$  be the CDF of lifetime  $L_i$ . Assuming that  $n$  is large and the system has reached stationarity, the CDF of residual lifetimes is given by [31]:

$$F_R(x) = P(R_i < x) = \frac{1}{E[L_i]} \int_0^x (1 - F(z)) dz. \quad (1)$$

For exponential lifetimes, the residuals are trivially exponential using the memoryless property of  $F(x)$ :  $F_R(x) = 1 - e^{-\lambda x}$ ; however, the next result shows that the residuals of Pareto distributions with shape  $\alpha$  are *more* heavy-tailed and exhibit shape parameter  $\alpha - 1$ .

*Lemma 2:* The CDF of residuals for Pareto lifetimes with  $F(x) = 1 - (1 + x/\beta)^{-\alpha}$ ,  $\alpha > 1$  is given by:

$$F_R(x) = 1 - \left(1 + \frac{x}{\beta}\right)^{1-\alpha}. \quad (2)$$

This outcome is not surprising as it is well-known that heavy-tailed distributions exhibit "memory," which means that users who survived in the system for some time  $t > 0$  are likely to remain online for longer periods of time than the arriving users. In fact, the larger the current age of a peer, the longer he/she is expected to remain online. The occurrence of this "heavy-tailed" phenomenon in P2P systems is supported by experimental observations [8] and can also be explained on the intuitive level. If a user  $v$  has already spent 10 hours in the system, it is generally unlikely that he/she will leave the network in the next 5 minutes; however, the same probability for newly arriving peers is substantially higher as some of them depart almost immediately [35].

Since the rest of the derivations in the paper rely on (1), it is important to verify that asymptotic approximations from renewal process theory actually hold in practice, which we accomplish via numerical simulations. We created a hypothetical system with  $n = 1,000$  users and degree  $k = 10$ , in which each node lived for a random duration  $L_i$  and then departed from the system. To prevent the network size from depleting to zero, each failed node was immediately replaced by a fresh node with another random lifetime  $L_j$  (the exact arrival process was not essential and had no effect on the results). For each new arrival  $v$  into the system, we recorded the residual lifetimes of the neighbors that  $v$  randomly selected from the pool of  $n - 1$  online peers.

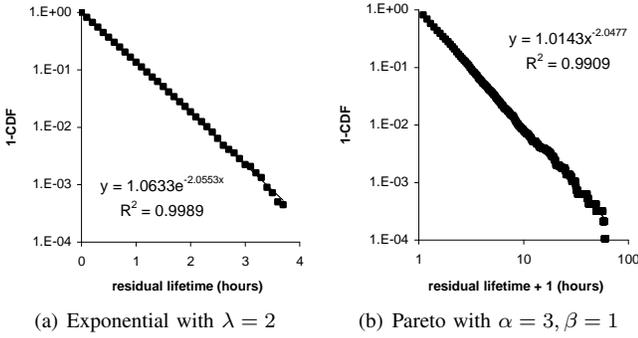


Fig. 3. Residual lifetimes in simulation.

Results of two typical simulations are plotted in Fig. 3 for the exponential and Pareto lifetimes. As we often do throughout the paper, parameters  $\alpha$  and  $\lambda$  are selected so that  $E[L_i]$  is 0.5 hours for both distributions and the scaling parameter  $\beta$  is set to 1 in the Pareto  $F(x)$ . As the figure shows, the residual exponential distribution remains exponential, while the Pareto case becomes more heavy-tailed and indeed exhibits shape parameter  $\alpha - 1 = 2$ . Further notice in the figure that the exponential  $R_i$  are limited by 4 hours, while the Pareto  $R_i$  stretch to as high as 61 hours.

While it is clear that node arrival instants  $t_v$  are uncorrelated with lifespans  $[a_i, d_i]$  of other nodes, the same observation holds for random points  $\tau_i$  at which the  $i$ -th neighbor of  $v$  fails. We extensively experimented with the active model, in which additional node selection occurred at instants  $\tau_i$  and found that all  $R_i$  obtained in this process also followed (1) very well (not shown for brevity).

### B. Resilience Analysis

Throughout the paper, we study resilience of P2P systems using two main metrics – the time before all neighbors of  $v$  are simultaneously in the failed state and the probability of this occurring before  $v$  decides to leave the system. We call the former metric *isolation time*  $T$  and the latter *probability of isolation*  $\pi$ . Recall that the passive model follows a simple pure-death degree evolution process illustrated in Fig. 1(a). In this environment, a node is considered isolated after its last surviving neighbor fails. Thus,  $T$  is equal to the maximum residual lifetime among all neighbors and its expectation can be written as (using the fact that  $T$  is a non-negative random variable) [40]:

$$E[T] = \int_0^{\infty} \left[ 1 - \frac{1}{E[L_i]^k} \left( \int_0^x (1 - F(z)) dz \right)^k \right] dx, \quad (3)$$

which leads to the following two results after straightforward integration.

*Theorem 1:* Assume a passive  $k$ -regular graph. Then, for exponential lifetimes:

$$E[T] = \frac{1}{\lambda} \sum_{i=1}^k \frac{1}{i}, \quad (4)$$

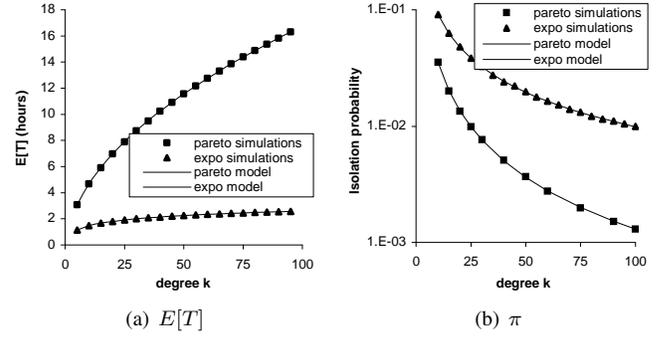


Fig. 4. Comparison of models (4)-(5) and (10)-(11) with simulations.

and for Pareto lifetimes with  $\alpha > 2$ :

$$E[T] = -\beta \left[ 1 + \frac{\Gamma\left(\frac{1}{1-\alpha}\right) k!}{(\alpha-1)\Gamma\left(k+2-\frac{\alpha}{1-\alpha}\right)} \right]. \quad (5)$$

*Proof:* For exponential lifetimes, using  $F(x) = 1 - e^{-\lambda x}$  in (3) and setting  $z = 1 - e^{-\lambda x}$ , we get:

$$\begin{aligned} E[T] &= \frac{1}{\lambda} \int_0^1 \frac{1 - z^k}{1 - z} dz \\ &= \frac{1}{\lambda} \left( z + \frac{z^2}{2} + \frac{z^3}{3} + \dots + \frac{z^k}{k} \right) \Big|_0^1, \end{aligned} \quad (6)$$

which directly leads to (4).

For Pareto lifetimes, substituting  $F(x) = 1 - (1 + x/\beta)^{-\alpha}$  into (3) and setting  $y = 1 + x/\beta$ , we obtain:

$$\begin{aligned} E[T] &= \beta \int_1^{\infty} \left[ 1 - (1 - y^{1-\alpha})^k \right] dy \\ &= \beta y \left[ 1 - {}_2F_1\left(\frac{1}{1-\alpha}, -k; \frac{2-\alpha}{1-\alpha}; y^{1-\alpha}\right) \right] \Big|_1^{\infty}, \end{aligned} \quad (7)$$

where  ${}_2F_1(a, b; c; z)$  is the Gauss hypergeometric function, which for  $z = 0$  is always 1 and for  $z = 1$  is [11]:

$${}_2F_1(a, b; c; 1) = \frac{\Gamma(c)\Gamma(c-b-a)}{\Gamma(c-a)\Gamma(c-b)}. \quad (8)$$

Expanding (8) and keeping in mind that  $\Gamma(k+1) = k!$ , we get (5). ■

Note that the gamma function in the numerator of (5) is negative due to  $\alpha > 1$ , which explains the  $-\beta$  term outside the brackets. Simulation results of (4)-(5) are shown in Fig. 4(a) for the average lifetime  $E[L_i]$  equal to 0.5 hours. Note that in the figure, simulations are plotted as isolated points and the two models as continuous lines. As the figure shows, simulation results for both exponential and Pareto distributions match the corresponding model very well. We also observe that for the same degree and average lifetime, Pareto nodes exhibit longer average times to isolation. For  $k = 10$ ,  $E[T]$  is 1.46 hours given exponential lifetimes and 4.68 hours given Pareto lifetimes. This difference was expected since  $T$  is determined by the residual lives of the neighbors, who in the Pareto case have large  $R_i$  and stay online longer than newly arriving peers.

We next focus on the probability that isolation occurs within the lifespan of a given user. Consider a node  $v$  with lifetime  $L_v$ . This node is forced to disconnect from the system only

if  $L_v$  is greater than  $T$ , which happens with probability  $\pi = P(T < L_v) = \int_0^\infty F_T(x)f(x)dx$ , where  $F_T(x) = F_R(x)^k$  is the CDF of time  $T$  and  $f(x)$  is the PDF of user lifetimes. This leads to:

$$\pi = \frac{1}{E[L_i]^k} \int_0^\infty \left( \int_0^x (1 - F(z))dz \right)^k f(x)dx. \quad (9)$$

Next, we study two distributions  $F(x)$  and demonstrate the effect of tail-weight on the local resilience of the system.

*Theorem 2:* Assume a passive  $k$ -regular graph. Then, for exponential lifetimes:

$$\pi = \frac{1}{k+1} \quad (10)$$

and for Pareto lifetimes with  $\alpha > 1$ :

$$\pi = \frac{\Gamma(1 + \frac{\alpha}{\alpha-1})k!}{\Gamma(k+1 + \frac{\alpha}{\alpha-1})}. \quad (11)$$

*Proof:* For exponential lifetimes, straightforward integration of (9) readily leads to the desired result:

$$\pi = \int_0^\infty \lambda e^{-\lambda x} (1 - e^{-\lambda x})^k dx = \frac{1}{k+1}. \quad (12)$$

For Pareto lifetimes, we have:

$$\pi = \frac{\alpha}{\beta} \int_0^\infty \left(1 + \frac{x}{\beta}\right)^{-\alpha-1} F_R^k(x) dx. \quad (13)$$

Reorganizing the terms in (13), using  $F_R(x)$  from (2), and setting  $w = (1 + x/\beta)^{1-\alpha}$ :

$$\begin{aligned} \pi &= \frac{\alpha}{\alpha-1} \int_0^1 w^{1/(\alpha-1)} (1-w)^k dw \\ &= \left( {}_2F_1\left(\frac{\alpha}{\alpha-1}, -k; \frac{2\alpha-1}{\alpha-1}; w\right) w^{\alpha/(\alpha-1)} \right) \Big|_0^1. \end{aligned} \quad (14)$$

Using (8) and the properties of Gauss hypergeometric functions discussed in the previous proof, we immediately obtain (11). ■

The exponential part of this theorem was expected from the memoryless property of exponential distributions [31], [40]. Hence, when a new node  $v$  joins a P2P system with exponentially distributed lifetimes  $L_i$ , it will be forced to disconnect if and only if it can outlive  $k$  other random nodes that started at the same time  $t_v$ . From Lemma 1, we already know that this happens with probability  $1/(k+1)$ .

The accuracy of (10)-(11) is shown in Fig. 4(b), which plots  $\pi$  obtained in simulations together with that predicted by the models. The simulation again uses a hypothetical P2P system with  $n = 1,000$  nodes and  $E[L_i] = 0.5$ . As the figure shows, simulations agree with predicted results well.

### C. Discussion

Notice from Fig. 4(b) that the Pareto  $\pi$  decays more quickly and always stays lower than the exponential  $\pi$ . To better understand the effect of  $\alpha$  on the isolation probability in the rather cryptic expression (11), we first show that for all choices of  $\alpha$ , Pareto systems are more resilient than exponential. We

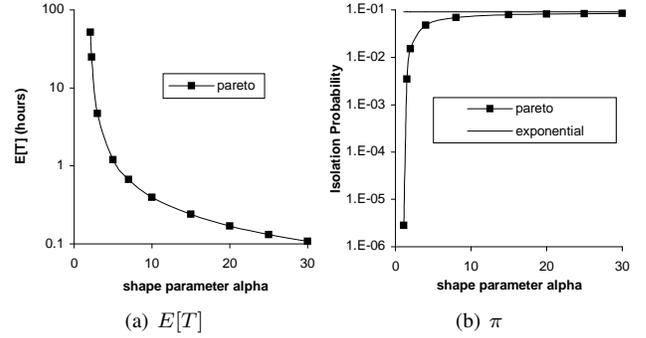


Fig. 5. Impact of tail weight on the average time to isolation and probability  $\pi$  for  $k = 10$  and Pareto lifetimes.

then show that as  $\alpha \rightarrow \infty$ , (11) approaches from below its upper bound (10).

Setting  $c = \Gamma(1 + \frac{\alpha}{\alpha-1})$ , re-write (11) expanding the gamma function in the denominator:

$$\pi = \frac{ck!}{(k + \frac{\alpha}{\alpha-1})!} \approx c \left( k + 1 + \frac{1}{2(\alpha-1)} \right)^{-\alpha/(\alpha-1)} \quad (15)$$

and notice that (15) always provides a faster decay to zero as a function of  $k$  than (10). For the Pareto example of  $\alpha = 3$  shown in Fig. 4(b),  $\pi$  follows the curve  $(k + 1.25)^{-1.5}$ , which decays faster than the exponential model by a factor of  $\sqrt{k}$ . This difference is even more pronounced for distributions with heavier tails. For example, (11) tends to zero as  $(k + 1.5)^{-2}$  for  $\alpha = 2$  and as  $(k + 6)^{-11}$  for  $\alpha = 1.1$ . The effect of tail-weight on isolation dynamics is shown in Fig. 5, where small values of  $\alpha$  indeed provide large  $E[T]$  and small  $\pi$ . Fig. 5(b) also demonstrates that as shape  $\alpha$  becomes large, the Pareto distribution no longer exhibits its “heavy-tailed” advantages and is essentially reduced to the exponential model. This can also be seen in (15), which tends to  $1/(k+1)$  for  $\alpha \rightarrow \infty$ .

Given the above discussion, it becomes apparent that it is possible to make  $\pi$  arbitrarily small with very heavy-tailed distributions (e.g.,  $\alpha = 1.05$  and  $k = 20$  produce  $\pi = 3.7 \times 10^{-12}$ ). While these results may be generally encouraging for networks of non-human devices with controllable characteristics, most current peer-to-peer systems are not likely to be satisfied with the performance of the passive model since selection of  $\alpha$  is not possible in the design of a typical P2P network and isolation probabilities in (11) are unacceptably high for  $\alpha > 1.5$ . The second problem with the passive framework is that its application to real systems requires *accurate* knowledge of the shape parameter  $\alpha$ , which may not be available in practice.

We overcome both problems in the next two section, where we show that active node replacement significantly increases resilience and that all Pareto distributions have a reasonably tight upper bound on  $\pi$  that does not depend on  $\alpha$ .

## IV. ACTIVE LIFETIME MODEL: STATIONARY ANALYSIS

To reduce the rate of isolation and repair broken routes in P2P networks, previous studies have suggested distributed recovery algorithms in which failed neighbors are dynamically

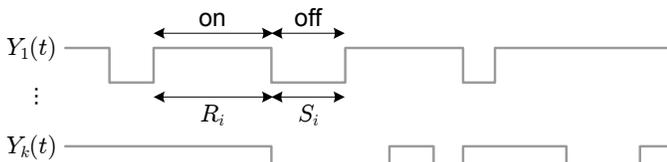


Fig. 6. On/off processes  $\{Y_i(t)\}_{i=1}^k$  depicting neighbor failure and replacement.

replaced with nodes that are still alive [32]. In this section, we offer a model for this strategy, derive the expected value of  $T$  using stationary techniques in renewal process theory, and analyze performance gains of this framework compared to the passive case. In the next section, we apply the theory of rare events for mixing processes to  $W(t)$  and derive a reasonably good upper bound on  $\pi$ .

It is natural to assume that node failure in P2P networks can be detected through some keep-alive mechanism, which includes periodic probing of each neighbor, retransmission of lost messages, and timeout-based decisions to search for a replacement. We do not dwell on the details of this framework and assume that each peer  $v$  is capable of detecting neighbor failure through some transport-layer protocol. The second step after node failure is detected is to repair the “failed” zone of the DHT and restructure certain links to maintain consistency and efficiency of routing (non-DHT systems may utilize a variety of random neighbor-replacement strategies [9], [14], [37]). We are not concerned with the details of this step either and generically combine both failure detection and repair into a random variable called  $S_i$ , which is the total “search” time for the  $i$ -th replacement in the system.<sup>2</sup>

#### A. Preliminaries

In the active model, each neighbor  $i$  ( $1 \leq i \leq k$ ) of node  $v$  is either alive at any time  $t$  or its replacement is being sought from among the remaining nodes in the graph. Thus, neighbor  $i$  can be considered in the *on* state at time  $t$  if it is alive or in the *off* state otherwise. This neighbor failure/replacement procedure can be modeled as an on/off process  $Y_i(t)$ :

$$Y_i(t) = \begin{cases} 1 & \text{neighbor } i \text{ alive at } t \\ 0 & \text{otherwise} \end{cases}. \quad (16)$$

This framework is illustrated in Fig. 6, which shows the evolution of  $k$  neighbor processes  $Y_1(t), \dots, Y_k(t)$ . Using this notation, the degree of node  $v$  at time  $t$  is equal to  $W(t) = \sum_{i=1}^k Y_i(t)$ . Similar to our definition in Section II-A, a node is isolated at such time  $T$  when all of its neighbors are simultaneously in the *off* state (see Fig. 1(b)). Thus, the maximum time a node can spend in the system before it is forced to disconnect can be formalized as the *first hitting time* of process  $W(t)$  on level 0:

$$T = \inf(t > 0 : W(t) = 0 | W(0) = k). \quad (17)$$

<sup>2</sup>Note that we assume that replacement only occurs *after* a failure of an out-degree neighbor. This behavior occurs in unstructured P2P systems and certain DHTs (such as Randomized Chord) where new arrivals do not replace the out-degree neighbors of existing nodes.

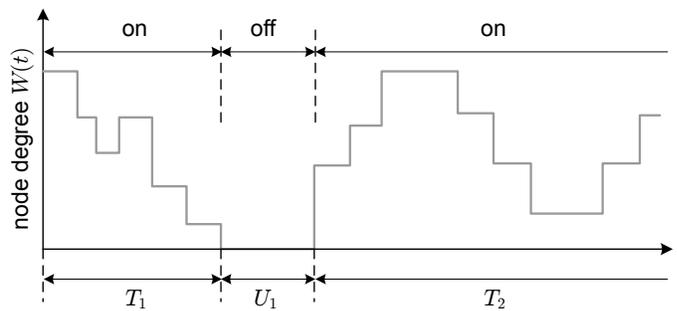


Fig. 7. On/off regenerative model of process  $W(t)$ .

Notice that under proper selection of the tail-weight of the lifetime distribution (i.e., the length of *on* periods),  $W(t)$  becomes a super-position of heavy-tailed on/off processes and may exhibit self-similarity for sufficiently large  $k$  [17], [21], [39]. Due to limited space, we omit log-log variance plots that confirm this effect, but note that to our knowledge, the fact that node degree in P2P networks may become self-similar has not been documented before.

#### B. Expected Time to Isolation

In what follows in the rest of this section, we apply stationary renewal process techniques to  $W(t)$  and derive a closed-form expression for  $E[T]$ .

*Theorem 3:* Assuming asymptotically small search delays and the active lifetime model, the expected time a node can remain in the system before becoming isolated is:

$$E[T] \approx \frac{E[S_i]}{k} \left[ \left( 1 + \frac{E[R_i]}{E[S_i]} \right)^k - 1 \right], \quad (18)$$

where  $E[S_i] < \infty$  is the mean search time and  $E[R_i] < \infty$  is the expected residual lifetime.

*Proof:* To simplify the analysis, we also view  $W(t)$  as an alternating on/off process, where each *on* period corresponds to  $W(t) > 0$  and each *off* period corresponds to  $W(t) = 0$ . Let  $U_j$  be the time that  $W(t)$  spends in its  $j$ -th *off* state and  $T_j$  represent the time in the  $j$ -th *on* state. As shown in Fig. 7,  $W(t)$  alternates between its on/off cycles and our goal is to determine the expected length of the first *on* period  $T_1$ . The proof consists of two parts: we first argue that the length of cycle  $T_1$  is similar to that of the remaining cycles  $T_j, j \geq 2$ , and then apply Smith’s theorem to  $W(t)$  to derive  $E[T_j], j \geq 2$ .

First, notice that cycle  $T_1$  is different from the other *on* periods since it always starts from  $W(t) = k$ , while the other *on* cycles start from  $W(t) = 1$ . However, since we already assumed that the search times are sufficiently small,  $W(t)$  at the beginning of each *on* period almost immediately “shoots back” to  $W(t) = k$ . This can be shown using arguments from large-deviations theory [36], which derives bounds on the return time of the system from very rare states back to its “most likely” state. This generally makes cycles  $T_1$  and  $T_j$  ( $j \geq 2$ ) different by a value that is negligible compared to  $E[T]$  in all real-life situations (see examples after the proof).

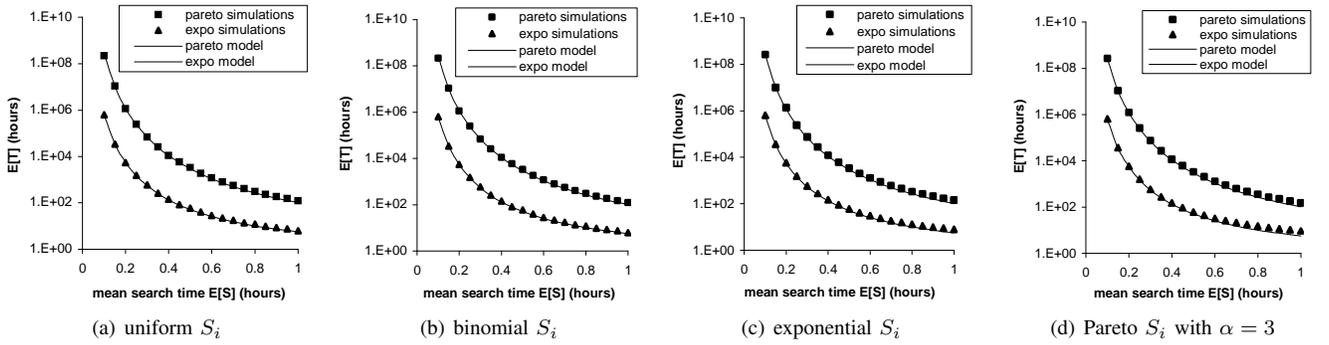


Fig. 8. Comparison of model (18) to simulation results with  $E[L_i] = 0.5$  and  $k = 10$ .

Second, approximate a superposition of renewal processes  $W(t) = \sum_i Y_i(t)$  as regenerative and assume that points  $\tau_j$  when  $W(t)$  goes into the  $j$ -th *off* state (i.e., makes a transition from 1 to 0) are its regenerative instances. Then, applying Smith's theorem [31] to  $W(t)$ , the probability of finding it in an isolated state at any random time  $t \gg 0$  is given by:

$$P(W(t) = 0) = \frac{E[U_j]}{E[T_j] + E[U_j]}. \quad (19)$$

Notice that (19) can also be expressed as the probability of finding all  $k$  neighbors in their *off* state:

$$P(W(t) = 0) = \left( \frac{E[S_i]}{E[S_i] + E[R_i]} \right)^k. \quad (20)$$

Combining the last two equations and solving for  $E[T_j]$ , we get:

$$E[T_j] = E[U_j] \left[ \left( \frac{E[R_i] + E[S_i]}{E[S_i]} \right)^k - 1 \right]. \quad (21)$$

Next, we compute  $E[U_j]$ . As before, suppose that the first instant of the  $j$ -th *off* cycle of  $W(t)$  starts at time  $\tau_j$ . At this time, there are  $k - 1$  already-failed neighbors still “searching” for their replacement and one neighbor that just failed at time  $\tau_j$ . Thus,  $U_j$  is the minimum time needed to find a replacement for the last neighbor or for one of the on-going searches to complete.

More formally, suppose that  $V_1, \dots, V_{k-1}$  are (from the Palm-Khintchine approximation) *i.i.d.* random variables that represent the remaining replacement delays of the  $k - 1$  already-failed neighbors and  $S_k$  is the replacement time of the last neighbor. Then duration  $U_j$  of the current *off* cycle is  $U_j = \min\{V_1, \dots, V_{k-1}, S_k\}$ . Assuming that  $F_S(x)$  is the CDF of search times  $S_i$ , the distribution of  $V = \min\{V_1, \dots, V_{k-1}\}$  is given by [40]:

$$F_V(x) = 1 - \left( 1 - \frac{1}{E[S_i]} \int_0^x (1 - F_S(z)) dz \right)^{k-1}. \quad (22)$$

Notice that  $U_j$  can also be written as  $U_j = \min\{V, S_k\}$  and its expectation is:

$$E[U_j] = \int_0^\infty (1 - F_S(x)) (1 - F_V(x)) dx. \quad (23)$$

Substituting (22) into (23) and setting  $y = \int_0^x (1 - F_S(z)) dz$ , we get:

$$E[U_j] = \int_0^\infty (1 - F_S(x)) \left( 1 - \frac{y}{E[S_i]} \right)^{k-1} dx. \quad (24)$$

Integrating (24) leads to:

$$E[U_j] = \int_0^{E[S_i]} \left( 1 - \frac{y}{E[S_i]} \right)^{k-1} dy = \frac{E[S_i]}{k}, \quad (25)$$

which gives us (18).  $\blacksquare$

Note that for small search delays this result does not generally depend on the distribution of residual lifetimes  $R_i$  or search times  $S_i$ , but only on their expected values. Fig. 8 shows  $E[T]$  obtained in simulations in a system with 1,000 nodes,  $k = 10$ , and four different distributions of search delay. In each part (a)-(d) of the figure, the two curves correspond to exponential and Pareto lifetimes with mean 30 minutes (as before, the models are plotted as solid lines and simulations are drawn using isolated points). Notice in all four subfigures that the model tracks simulation results for over 7 orders of magnitude and that the expected isolation time is in fact not sensitive to the distribution of  $S_i$ .

As in the passive model, the exponential distribution of lifetimes provides a lower bound on the performance of any Pareto system since exponential  $E[R_i]$  are always smaller than the corresponding Pareto  $E[R_i]$ . Further observe that the main factor that determines  $E[T]$  is the *ratio* of  $E[R_i]$  to  $E[S_i]$  and not their individual values. Using this insight and Fig. 8, we can conclude that in systems with 10 neighbors and expected search delay at least 5 times smaller than the mean lifetime,  $E[T]$  is at least *one million* times larger than the mean session length of an average user. Furthermore, this result holds for exponential as well as Pareto distributions with arbitrary  $\alpha$ . This is a significant improvement over the results of the passive model in Section III.

### C. Example

We next rephrase the results obtained in the previous subsection in more practical terms and study several examples that demonstrate the usage of model (18). For the sake of this example, suppose that each node relies on a keep-alive

TABLE I  
EXPECTED TIME  $E[T]$  FOR  $E[R_i] = 1$  HOUR

Timeout $\delta$	$k = 20$	$k = 10$	$k = 5$
20 sec	$10^{41}$ years	$10^{17}$ years	188,034 years
2 min	$10^{28}$ years	$10^{11}$ years	282 years
45 min	404,779 years	680 days	49 hours

protocol with timeout  $\delta$ . Then the distribution of failure-detection delays is uniform in  $[0, \delta]$  depending on when the neighbor died with respect to the nearest ping message. The actual search time  $d$  to find a replacement may be determined by the average number of application-layer hops between each pair of users, may depend on the bootstrap delay if a centralized server is used for neighbor replacement, or may be simply zero if users select backup neighbors from a local cache. Using the notation above, we have the following restatement of the previous theorem.

*Corollary 1:* The expected isolation time in P2P systems with active neighbor replacement is given by:

$$E[T] = \frac{\delta + 2d}{2k} \left( 1 + \frac{2E[R_i]}{\delta + 2d} \right)^k. \quad (26)$$

Consider a system with the average replacement delay  $d = 2$  seconds and  $E[R_i] = 1$  hour. Table I shows the expected time to isolation for several values of timeout  $\delta$  and degree  $k$ . For small keep-alive delays (2 minutes or less), even  $k = 5$  provides longer expected times to isolation than the lifetime of any human being. Also notice that for  $\delta = 2$  minutes, degree  $k = 20$  provides more years before expected isolation than there are molecules in a glass of water [4].

Since routing delay  $d$  in the overlay network is generally much smaller than keep-alive timeout  $\delta$ , the diameter of the graph does not usually contribute to the resilience of the system. In other cases when  $d \log n$  is comparable to  $\delta$ , P2P graphs with smaller diameter may exhibit higher resilience as can be observed in (26).

#### D. Real P2P Networks

Finally, we address the practicality of the examples shown in the paper so far. In particular we consider the results shown in [8], which suggest that the empirical distribution of user lifetimes in real P2P networks follows a Pareto distribution with shape parameter  $\alpha = 1.06$ . Such heavy-tailed distributions result in  $E[R_i] = E[T] = \infty$  and do not lead to much interesting discussion. At the same time, notice that while it is hypothetically possible to construct a P2P system with  $\alpha = 1.06$ , it can also be argued that the measurement study in [8] sampled the *residuals* rather than the actual lifetimes of the users. This is a consequence of the “snapshots” taken every 20 minutes, which missed all peers with  $L_i < 20$  minutes and shortened the lifespan of all remaining users by random amounts of time. As such, these results point toward  $\alpha = 2.06$ , which is a much more realistic shape parameter even though it still produces enormous  $E[T]$  for all feasible values of  $E[S_i]$ . This is demonstrated for model (26) in Table II where the expected lifetime of each user is only double that in Table I,

TABLE II  
EXPECTED TIME  $E[T]$  FOR PARETO LIFETIMES WITH  $\alpha = 2.06$   
( $E[L_i] = 0.93$  HOURS,  $E[R_i] = 16.6$  HOURS)

Timeout $\delta$	$k = 10$	$k = 5$	$k = 2$
20 sec	$10^{29}$ years	$10^{11}$ years	4.7 years
2 min	$10^{23}$ years	$10^8$ years	336 days
45 min	$10^{11}$ years	1,619 years	16 days

but  $E[T]$  is 5-12 orders of magnitude larger. This is a result of  $E[R_i]$  rising from 1 hour in the former case to 16.6 hours in the latter scenario.

#### V. ACTIVE LIFETIME MODEL: TRANSIENT ANALYSIS

Given the examples in the previous section, it may at first appear that  $\pi$  must automatically be very small since  $E[T]$  is so “huge” under all practical conditions. However, in principle, there is a possibility that a large mass of  $T$  is concentrated on very small values and that a handful of extremely large values skew the mean of  $T$  to its present location. We additionally are interested in more than just knowing that  $\pi$  is “small” – we specifically aim to understand the order of this value for different  $E[S_i]$ .

As in previous sections, let  $L_v$  denote the lifetime of  $v$  and  $T$  the random time before  $v$ ’s neighbors force an isolation. Notice that  $\pi = P(T < L_v) = \int_0^\infty F_T(t) f(t) dt$  is an integral of the CDF function  $F_T(t) = P(T < t)$  of the first hitting time of process  $W(t)$  on level 0. The exact distribution of  $T$  is difficult to develop in closed-form since it depends on *transient* properties of a complex process  $W(t)$ . To tackle this problem, we first study the asymptotic case of  $E[S_i] \ll E[R_i]$  and apply results from the theory of rare events for Markov jump processes [2], [36] to derive a very accurate formula for  $\pi$  assuming exponential lifetimes. We then use this result to upper-bound the Pareto version of this metric.

##### A. Exponential Lifetimes

We start with exponential lifetimes and assume reasonably small search times. For  $E[S_i]$  larger than  $E[R_i]$ , accurate isolation probabilities are available from the passive model in Section III.

*Theorem 4:* For exponential lifetimes  $L_i$  and exponential search delays  $S_i$ , the probability of isolation converges to the following as  $E[S_i] \rightarrow 0$ :

$$\pi = \frac{E[L_i]}{E[T]}. \quad (27)$$

*Proof:* Given exponential lifetimes and search delays, notice that  $W(t)$  can be viewed as a continuous-time Markov chain, where the time spent in each state  $j$  before making a transition to state  $j-1$  is the minimum of exactly  $j$  exponential variables (i.e., the time to the next failure). Assume that the CDF of  $R_i$  is  $F_R(x) = 1 - e^{-\lambda x}$ , where  $\lambda = 1/E[L_i]$ . Then the CDF of  $\min\{R_1, \dots, R_j\}$  is  $1 - (1 - F_R(x))^j = 1 - e^{-\lambda j x}$ , which is another exponential variable with rate  $\lambda j$ . Next notice that the delays before  $W(t)$  makes a transition from state  $j$

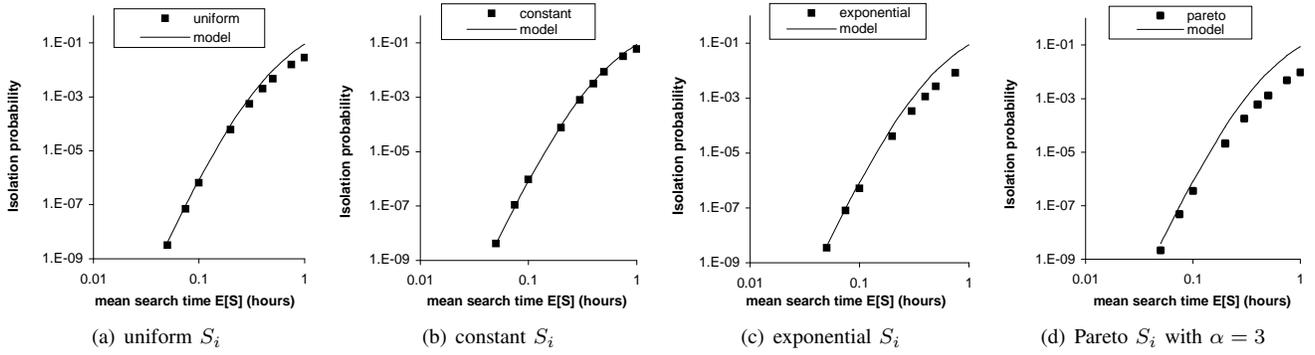


Fig. 9. Comparison of model (27) to simulation results for exponential lifetimes with  $E[L_i] = 0.5$  and  $k = 10$ .

to  $j + 1$  (i.e., upon recovering a neighbor) are given by the minimum of  $k - j$  residual search times, which is yet another exponential random variable with rate  $(k - j)\mu$ , where  $\mu = 1/E[S_i]$ .

To bound the CDF of  $T$ , one approach is to utilize the classical analysis from Markov chains that relies on numerical exponentiation of transition (or rate) matrices; however, it does not lead to a closed-form solution for  $P(T < L_v)$ . Instead, we apply a result for rare events in Markov chains due to Aldous *et al.* [2], which shows that  $T$  asymptotically behaves as an exponential random variable with mean  $E[T]$ :

$$|P(T > t) - e^{-t/E[T]}| \leq \frac{\tau}{E[T]}, \quad (28)$$

where  $E[T]$  is the expected time between the visits to the rare state 0 and  $\tau$  is the relaxation time of the chain. Re-writing (28) in terms of  $F_T(t) = P(T < t)$  and applying Taylor expansion to  $e^{-t/E[T]}$ :

$$\frac{t - \tau}{E[T]} \leq F_T(t) \leq \frac{t + \tau}{E[T]}. \quad (29)$$

Next, recall that relaxation time  $\tau$  is the inverse of the second largest eigenvalue of  $-Q$ , where  $Q$  is the rate matrix of the chain. For birth-death chains, matrix  $Q$  is tri-diagonal with  $Q(i, i) = -\sum_{j \neq i} Q(i, j)$ :

$$Q = \begin{bmatrix} -k\mu & k\mu & \cdots & 0 \\ \lambda & -\lambda - (k-1)\mu & (k-1)\mu & \\ 0 & \cdots & \cdots & \mu \\ 0 & \cdots & k\lambda & -k\lambda \end{bmatrix}. \quad (30)$$

We treat state  $W(t) = 0$  as non-absorbing and allow the chain to return back to state 1 at the rate  $k\mu$ . Then, the second largest eigenvalue of this matrix is available in closed-form (e.g., [22]) and equals the sum of individual rates:  $\lambda_2 = 1/\tau = \lambda + \mu$ . Noticing that:

$$\tau = \frac{1}{\lambda + \mu} = \frac{1}{1/E[L_i] + 1/E[S_i]} \approx E[S_i], \quad (31)$$

we conclude that  $\tau$  is on the order of  $E[S_i]$  and is generally very small. Writing  $\tau \approx E[S_i]$  and integrating the upper bound of (29) over all possible values of lifetime  $t$ , we get:

$$\pi \leq \int_0^{\infty} \frac{(t + \tau)f(t)dt}{E[T]} = \frac{E[L_i] + E[S_i]}{E[T]}. \quad (32)$$

We similarly obtain a lower bound on  $\pi$ , which is equal to  $(E[L_i] - E[S_i])/E[T]$ . Neglecting small  $E[S_i]$ , observe that both bounds reduce to (27). ■

Interestingly, for non-exponential, but asymptotically small search delays,  $W(t)$  can usually be approximated by an equivalent, but quickly-mixing process. Further note that bounds similar to (29) are reasonably accurate regardless of the distribution of  $S_i$  [1]. This is demonstrated in Fig. 9 using four distributions of search time – exponential with rate  $\lambda = 1/E[S_i]$ , constant equal to  $E[S_i]$ , uniform in  $[0, 2E[S_i]]$ , and Pareto with  $\alpha = 3$ . As shown in the figure, all four cases converge with acceptable accuracy to the asymptotic formula (27) and achieve isolation probability  $\pi \approx 3.8 \times 10^{-9}$  when the expected search time reduces to 3 minutes. Also notice in the figure that for all values of  $E[S_i]$  and all four search delay distributions, model (27) provides an *upper* bound on the actual  $\pi$ , which means that real systems are *more* resilient than suggested by our model.

### B. Heavy-Tailed Lifetimes

Although it would be nice to obtain a similar result  $\pi \approx E[L_i]/E[T]$  for the Pareto case, unfortunately the situation with a superposition of heavy-tailed on/off processes is different since  $W(t)$  is slowly mixing and the same bounds no longer apply. Intuitively, it is clear that large values of  $E[T]$  in the Pareto case are caused by a handful of users with enormous isolation delays, while the majority of remaining peers acquire neighbors with short lifetimes and suffer isolation almost as quickly as in the exponential case.

Consider an example that illustrates this effect and shows that huge values of  $E[T]$  in Pareto systems have little impact on  $\pi$ . For 10 neighbors,  $\alpha = 3$ ,  $\lambda = 2$  ( $E[L_i] = 30$  minutes), and constant search time  $s = 6$  minutes, the Pareto  $E[T]$  is larger than the exponential  $E[T]$  by a factor of 865. However, the ratio of their isolation probabilities is only 5.7. For  $\alpha = 2.5$  and  $\lambda = 1.5$  ( $E[L_i] = 40$  minutes), the expected times to isolation differ by a factor of  $8.1 \times 10^5$ , but the ratio of their  $\pi$  is only 7.5.

It may be possible to derive an accurate approximation for Pareto  $\pi$ ; however, one may also argue that the usefulness of such a result is limited given that shape parameter  $\alpha$  and the distribution of user lifetimes (lognormal, Pareto, etc.) are often not known accurately. We leave the exploration of this

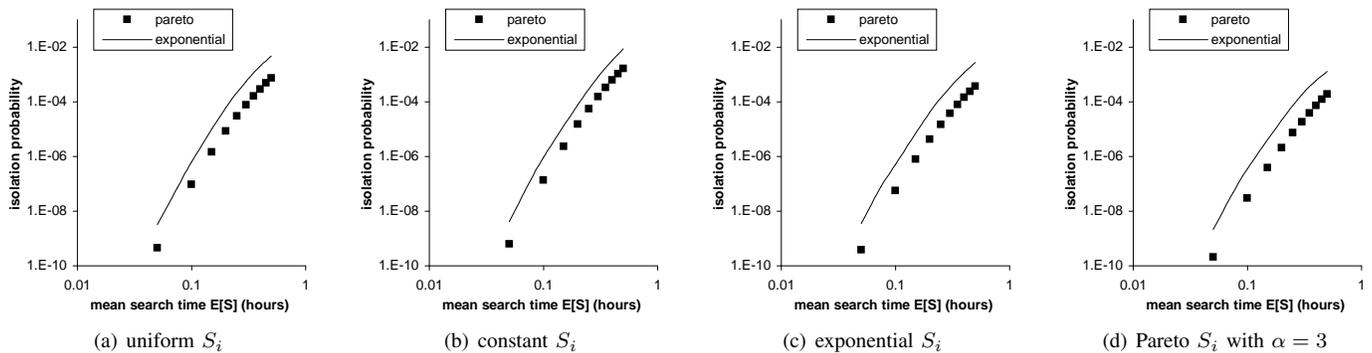


Fig. 10. Exponential  $\pi$  as an upper bound for Pareto  $\pi$  in simulations with  $E[L_i] = 0.5$  hours and  $k = 10$ .

problem for future work and instead utilize the exponential metric (27) as an upper bound on  $\pi$  in systems with sufficiently heavy-tailed lifetime distributions. One example for  $E[L_i] = 30$  minutes and  $k = 10$  is illustrated in Fig. 10, where the exponential  $\pi$  indeed tightly upper-bounds the Pareto  $\pi$  over the entire range of search delays and their distributions. In fact, it can be seen in each individual figure that the ratio between the two metrics does not depend on  $E[S_i]$  and that both curves decay at the same rate as  $E[S_i] \rightarrow 0$ .

The above observations are summarized in the next result, which formally follows from the fact that heavy-tailed  $L_i$  imply stochastically larger residual lifetimes  $R_i$  and a straightforward expansion of  $E[T]$  in (27).

*Corollary 2:* For an arbitrary distribution of search delays and any lifetime distribution  $F(x)$  with an exponential or heavier tail, which includes Pareto, lognormal, Weibull, and Cauchy distributions, the following upper bound holds:

$$\pi \leq \frac{\rho k}{(1 + \rho)^k - 1}, \quad (33)$$

where  $\rho = E[L_i]/E[S_i]$  is the ratio of the mean user lifetime to the mean search delay.

For example, using 30-minute average lifetimes, 9 neighbors per node, and 1-minute average node replacement delay, the upper bound in (33) equals  $1.02 \times 10^{-11}$ , which allows the joining users in a 100-billion node network to stay connected to the graph for their entire lifespans with probability  $1 - 1/n$ . Using the uniform failure model of prior work and  $p = 1/2$ , each user required 37 neighbors to achieve the same  $\pi$  regardless of the actual dynamics of the system.

Even though exponential  $\pi$  is often several times larger than the Pareto  $\pi$  (the exact ratio depends on shape  $\alpha$ ), it turns out that the difference in node degree needed to achieve a certain level of resilience is usually negligible. To illustrate this result, Table III shows the minimum degree  $k$  that ensures a given  $\pi$  for different values of search time  $E[S_i]$  and Pareto lifetimes with  $\alpha = 2.06$  (to maintain the mean lifetime 30 minutes, the distribution is scaled using  $\beta = 0.53$ ). The column “uniform  $p = 1/2$ ” contains degree  $k$  that can be deduced from the  $p$ -percent failure model (for  $p = 1/2$ ) discussed in previous studies [38]. Observe in the table that the exponential case in fact provides a tight upper bound on the actual minimum degree and that the difference between the two cases is at most 1 neighbor.

TABLE III

MINIMUM DEGREE NEEDED TO ACHIEVE A CERTAIN  $\pi$  FOR PARETO LIFETIMES WITH  $\alpha = 2.06$  AND  $E[L_i] = 0.5$  HOURS

$\pi$	Uniform $p = 1/2$	Lifetime P2P	Mean search time $E[S_i]$		
			6 min	2 min	20 sec
$10^{-6}$	20	Bound (33)	10	7	5
		Simulations	9	6	4
$10^{-9}$	30	Bound (33)	14	9	6
		Simulations	13	8	6
$10^{-12}$	40	Bound (33)	18	12	8
		Simulations	17	11	7

### C. Irregular Graphs

The final issue addressed in this section is whether P2P networks can become more resilient if node degree is allowed to vary from node to node. It is sometimes argued [10], [35] that graphs with a heavy-tailed degree distribution exhibit highly resilient characteristics and are robust to node failure. Another question raised in the literature is whether DHTs are more resilient than their unstructured counterparts such as Gnutella. In this section we prove that, given the assumptions used so far in the paper,  $k$ -regular graphs offer the highest local resilience among all systems with a given average degree. This translates into “optimality” of DHTs as long as they can balance their zone-sizes and distribute degree evenly among the peers.

Consider a P2P system in which node degrees  $k_1, \dots, k_n$  are drawn from an arbitrary distribution with mean  $E[k_i]$ . Using Jensen’s inequality for convex functions and the upper bound in (33), the following result follows effortlessly.

*Theorem 5:* Assuming that lifetimes are independent of node degree and are not used in the neighbor-selection process, regular graphs are the most resilient for a given average degree  $E[k_i]$ .

*Proof:* We proceed by computing the disconnection probability of an arriving user  $i$  who is assigned random degree  $k_i$  from some arbitrary distribution. Notice that  $\pi$  in (27), (33) can be written as some *strictly convex* function:

$$g(k) = \frac{\rho k}{(1 + \rho)^k - 1}. \quad (34)$$

Then, the probability of disconnection averaged over all join-

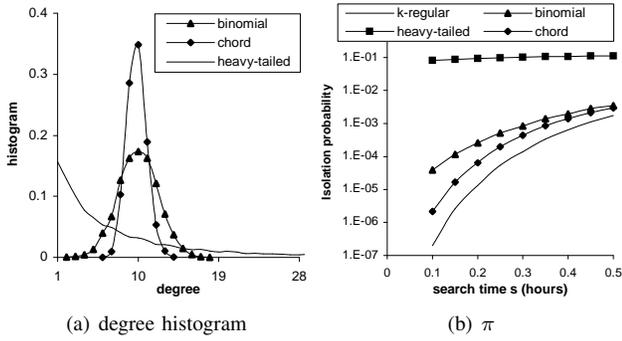


Fig. 11. (a) Degree distribution in irregular graphs. (b) User isolation probability  $\pi$  for irregular graphs (average degree  $E[k_i] = 10$ , lifetimes are Pareto with  $E[L_i] = 0.5$  hours).

ing users is:

$$\pi = \sum_{k=1}^{\infty} g(k)P(k_i = k) = E[g(k_i)]. \quad (35)$$

Using Jensen's inequality [40] for convex functions, we have  $E[g(k_i)] \geq g(E[k_i])$ , which means that (35) is always no less than the same metric in graphs with a constant degree equal to  $E[k_i]$ . ■

To demonstrate the effect of node degree on isolation probability in irregular graphs, we examine three systems with 1,000 nodes: 1) Chord with a random distribution of out-degree, which is a consequence of imbalance in zone sizes; 2) a  $G(n, p)$  graph with binomial degree for  $p = 0.5$ ; and 3) a heavy-tailed graph with Pareto degree for  $\alpha = 2.5$  and  $\beta = 15$ . We selected these parameters so that each of the graphs had mean degree  $E[k_i]$  equal to 10. The distribution of degree in these graphs is shown in Fig. 11(a). Notice that Chord has the lowest variance and its probability mass concentration around the mean is the best of the three systems. The binomial case is slightly worse, while the heavy-tailed graph is the worst. According to Theorem 5, all of these systems should have larger isolation probabilities than those of 10-regular graphs and should exhibit performance inverse proportional to the variance of their degree.

Simulation results of  $\pi$  are shown in Fig. 11(b) for Pareto lifetimes with  $\alpha = 3$  and  $E[L_i] = 0.5$  hours (search times are constant). Observe in the figure that the  $k$ -regular system is in fact better than the irregular graphs and that the performance of the latter deteriorates as  $Var[k_i]$  increases. For  $s = 0.1$  (6 minutes), the  $k$ -regular graph offers  $\pi$  lower than Chord's by a factor of 10 and lower than that in  $G(n, p)$  by a factor of 190. Furthermore, the P2P system with a heavy-tailed degree distribution in the figure exhibits the same poor performance regardless of the search time  $s$  and allows users to become isolated  $10^2 - 10^6$  times more frequently than in the optimal case, all of which is caused by 37% of the users having degree 3 or less.

Thus, in cases when degree is independent of user lifetimes, we find no evidence to suggest that unstructured P2P systems with a heavy-tailed (or otherwise irregular) degree can provide better resilience than  $k$ -regular DHTs.

## VI. GLOBAL RESILIENCE

We finish the paper by analyzing the probability of network partitioning under uniform node failure and showing that this metric has a simple expression for a certain family of graphs, which includes many proposed P2P networks. We then apply this insight to lifetime-based systems and utilize the earlier derived metric  $\pi$  to characterize the evolution of P2P networks under churn.

### A. Classical Result

One may wonder how local resilience (i.e., absence of isolated vertices) of P2P graphs translates into their global resilience (i.e., connectivity of the entire network). While this topic has not received much attention in the P2P community, it has been extensively researched in random graph theory and interconnection networks. Existing results for classical random graphs have roots in the work of Erdős and Rényi in the 1960s and demonstrate that *almost every* (i.e., with probability  $1 - o(1)$  as  $n \rightarrow \infty$ ) random graph including  $G(n, p)$ ,  $G(n, M)$ , and  $G(n, k_{out})$  is connected if and only if it has no isolated vertices [6], i.e.,

$$P(G \text{ is connected}) = P(X = 0) \quad \text{as } n \rightarrow \infty, \quad (36)$$

where  $X$  is the number of isolated nodes after the failure. After some manipulation, this result can be translated to apply to unstructured P2P networks, where each joining user draws some number of random out-degree neighbors from among the existing nodes (see below for simulations that confirm this).

For deterministic networks, connectivity of a graph  $G$  after node/edge failure has also received a fair amount of attention (e.g., [7], [25]). In interconnection networks, exact formulas for the connectivity of deterministic graphs exist [25]; however, they require computation of NP-complete metrics and no closed-form solution is available even for the basic hypercube. However, from the perspective of random graph theory, it has been shown [6], [7] that hypercubes with faulty elements asymptotically behave as random graphs and thus almost surely disconnect with isolated nodes as  $n$  becomes large.

Even though the necessary condition for a deterministic graph  $G$  to satisfy (36) is unknown at this time, sufficient conditions can be extrapolated from the proofs of this relationship for the hypercube [6]. The general requirement on  $G$  is that its expansion (strength of the various cuts) must be no worse than that of the hypercube.<sup>3</sup> Due to limited space and wide variety of deterministic P2P constructs, we do not provide a rigorous re-derivation of this fact, but instead note that Chord [38], Pastry [33], and CAN with the number of dimensions  $d = \Theta(\log n)$  [30] can be directly reduced to hypercubes; de Bruijn graphs [18] exhibit better connectivity than hypercubes [26]; and hybrid networks (such as Symphony [27], Randomized Chord [27], and [3]) generally have connectivity no worse than  $G(n, k_{out})$ .

<sup>3</sup>For each set  $S$  in the original graph  $G$ , its node boundary must satisfy a certain inequality that is an increasing function of  $|S|$  [7]. Graphs that do not fulfill this requirement include trees, cycles, and other weakly connected structures.

TABLE IV  
GLOBAL RESILIENCE FOR  $n = 16,384$  AND OUT-DEGREE  $k = 14$

$p$	$P(G$ is weakly connected)			$P(G$ has no isolated nodes)			$P(G$ disconnects with isolated nodes)		
	Chord	Symphony	Gnutella	Chord	Symphony	Gnutella	Chord	Symphony	Gnutella
0.5	0.99996	0.99768	0.86257	0.99996	0.99768	0.86260	1	1	0.999782
0.55	0.99918	0.98750	0.58042	0.99918	0.98750	0.58064	1	1	0.999476
0.6	0.99354	0.93914	0.17081	0.99354	0.93917	0.17148	1	0.99958	0.999192
0.65	0.95001	0.75520	0.00547	0.95004	0.75527	0.00560	0.99946	0.99971	0.999869
0.7	0.72619	0.31153	0	0.72650	0.31205	0	0.99980	0.99922	1
0.8	0.00040	0	0	0.00043	0	0	0.99997	1	1

We next show simulation results that confirm the application of classical result (36) to two types of DHTs and one unstructured Gnutella-like network (all three are directed graphs). Table IV shows simulations of fully-populated Chord, Symphony, and a  $k$ -regular Gnutella network under uniform  $p$ -percent node failure using 100,000 failure patterns. Note that for directed graphs, (36) applies only to the definition of *weak* connectivity, which means that  $G$  is disconnected if and only if its *undirected* version  $G'$  is and a node  $v \in G$  is isolated if and only if it is isolated in  $G'$ . Since the total degree (which is the sum of out-degree and in-degree) at each node is different between the graphs in the table, their respective disconnection probabilities are different.

As the table shows, (36) holds with high accuracy and a vast majority of disconnections contain at least one isolated node (see the last three columns of the table). Additional simulations show that (36) applies to Pastry, CAN, de Bruijn graphs, Randomized Chord, and degree-irregular Gnutella. We omit these results for brevity.

### B. Lifetime-Based Extension

It is not difficult to notice that (36) holds for *lifetime-based* P2P graphs and that dynamic P2P networks are also much more likely to develop disconnections around single nodes rather than along boundaries of larger sets  $S$ . However, instead of having a single node-failure metric  $p$ , we have a probability of isolation  $\pi$  associated with each joining user  $i$ . Thus, one may ask the question *what is the probability that the system survives  $N$  user joins and stays connected the entire time?* The answer is very simple: assuming  $Z$  is a geometric random variable measuring the number of user joins before the first disconnection of the network, we have for almost every sufficiently large graph:

$$P(Z > N) \approx \left(1 - \frac{\rho k}{(1 + \rho)^k - 1}\right)^N. \quad (37)$$

Simulation results of (37) for a degree-regular, fully-populated CAN (i.e., each failed node is immediately replaced by a new arrival with the same hash index) are shown in Table V using  $N = 1$  million joins and 10,000 iterations per search time, where metric  $q(G) = P(X > 0 | G \text{ is disconnected})$  is the probability that the graph partitions with at least one isolated node and  $r(G)$  is the probability that the largest connected component after the disconnection contains exactly  $n - 1$  nodes. As the table shows, simulations match the model very well and also confirm that the most likely disconnection

TABLE V  
COMPARISON OF  $P(Z > 10^6)$  IN CAN TO MODEL (37) (EXPONENTIAL LIFETIMES WITH  $E[L_i] = 30$  MINUTES, 6 DIMENSIONS, DEGREE  $k = 12$ , AND  $n = 4096$  NODES)

Fixed search time (min)	Actual $P(Z > N)$	Model (37)	$q(G)$	$r(G)$
6	0.9732	0.9728	1	1
7.5	0.8118	0.8124	1	1
8.5	0.5669	0.5659	1	1
9	0.4065	0.4028	1	1
9.5	0.2613	0.2645	1	1
10.5	0.0482	0.0471	1	1

pattern of lifetime-based systems includes at least one isolated node (i.e.,  $q(G) = 1$ ). In fact, the table shows an even stronger result – for reasonably small search delays, network partitioning almost surely affects only *one* node in the system (i.e.,  $r(G) = 1$ ). The same conclusion holds for other P2P graphs, Pareto lifetimes, and random search delays. We omit these results for brevity.

Model (37) suggests that when search delays become very small, the system may evolve for many months or years before disconnection. Consider a 12-regular CAN system with 1-minute search delays and 30-minute average lifetimes. Assuming that  $n = 10^6$  and each user diligently joins the system once per day, the probability that the network can evolve for 2,700 years ( $N = 10^{13}$  joins) before disconnecting for the first time is 0.9956. The mean delay between disconnections is  $E[Z] = 1/\pi$  user joins, or 5.9 million years.

## VII. CONCLUSION

This paper examined two aspects of resilience in dynamic P2P systems – ability of each user to stay connected to the system in the presence of frequent node departure and partitioning behavior of the network as  $n \rightarrow \infty$ . We found that under all practical search times,  $k$ -regular graphs were much more resilient than traditionally implied [18], [24], [38] and further showed that dynamic P2P networks could almost surely remain connected as long as no user suffered simultaneous neighbor failure. We also demonstrated that varying node degree from peer to peer can have a positive impact on resilience *only* when such decisions are correlated with the users' lifetimes. Future work involves additional analysis of user churn in networks where lifetimes are allowed to be used in neighbor selection.

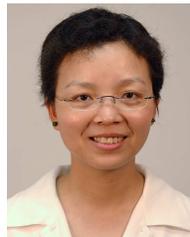
## REFERENCES

- [1] M. Abadi and A. Galves, "Inequalities for the occurrence times of rare events in mixing processes. the state of the art," *Markov Proc. Relat. Fields*, vol. 7, no. 1, pp. 97–112, 2001.
- [2] D. Aldous and M. Brown, "Inequalities for rare events in time-reversible markov chains ii," *Stochastic Processes and their Applications*, vol. 44, no. 1, pp. 15–25, Jan. 1993.
- [3] J. Aspnes, Z. Diamadi, and G. Shah, "Fault tolerant routing in peer to peer systems," in *ACM PODC*, July 2002, pp. 223–232.
- [4] P. Atkins, *Physical Chemistry*. Freeman, NY, 1986.
- [5] R. Bhagwan, S. Savage, and G. Voelker, "Understanding availability," in *IPTPS*, Feb. 2003, pp. 256–267.
- [6] B. Bollobás, *Random Graphs*. Cambridge University Press, 2001.
- [7] Y. Burtin, "Connection probability of a random subgraph of an  $n$ -dimensional cube," *Probl. Pered. Inf.*, vol. 13, no. 2, pp. 147–152, Apr.-June 1977.
- [8] F. Bustamante and Y. Qiao, "Friendships that last: Peer lifespan and its role in p2p protocols," in *IEEE WCW*, Sept. 2003, pp. 233–246.
- [9] Y. Chawathe, S. Ratnasamy, L. Breslau, N. Lanham, and S. Shenker, "Making gnutella like p2p systems scalable," in *ACM SIGCOMM*, Aug. 2003, pp. 407–418.
- [10] R. Cohen, K. Erez, D. ben Avraham, and S. Havlin, "Resilience of the internet to random breakdowns," *Physical Review Letters*, vol. 85, no. 21, pp. 4626–4628, Nov. 2000.
- [11] H. Exton, *Handbook of Hypergeometric Integrals: Theory, Applications, Tables, Computer Programs*. Ellis Horwood, Chichester, England, 1978.
- [12] A. Fiat and J. Saia, "Censorship resistant peer-to-peer content addressable networks," in *ACM-SIAM SODA*, Jan. 2002, pp. 94–103.
- [13] A. Ganesh and L. Massoulié, "Failure resilience in balanced overlay networks," in *Allerton Conference on Communication, Control and Computing*, Oct. 2003.
- [14] C. Gkantsidis, M. Mihail, and A. Saberi, "Random walks in peer-to-peer networks," in *IEEE INFOCOM*, Mar. 2004, pp. 119–130.
- [15] P. Golle, K. Leyton-Brown, I. Mironov, and M. Lillibridge, "Incentives for sharing in peer-to-peer networks," in *ACM Conference on Electronic Commerce*, Oct. 2001, pp. 264–267.
- [16] K. Gummadi, R. Gummadi, S. Gribble, S. Ratnasamy, S. Shenker, and I. Stoica, "The impact of dht routing geometry on resilience and proximity," in *ACM SIGCOMM*, Aug. 2003, pp. 381–394.
- [17] D. Heath, S. Resnick, and G. Samorodnitsky, "Patterns of buffer overflow in a class of queues with long memory in the input stream," *Annals of Applied Probability*, vol. 7, no. 4, pp. 1021–1057, Nov. 1997.
- [18] M. Kaashoek and D. Karger, "Koorde: A simple degree-optimal distributed hash table," in *IPTPS*, Feb. 2003, pp. 98–107.
- [19] S. Krishnamurthy, S. El-Ansarh, E. Aurell, and S. Haridi, "A statistical theory of chord under churn," in *IPTPS*, Feb. 2005, pp. 93–103.
- [20] F. Leighton, B. Maggs, and R. Sitamaran, "On the fault tolerance of some popular bounded-degree networks," in *IEEE FOCS*, Oct. 1995, pp. 542–552.
- [21] W. Leland, M. Taqqu, W. Willinger, and D. Wilson, "On the self-similar nature of ethernet traffic," in *ACM SIGCOMM*, Sept. 1993, pp. 183–193.
- [22] R. Lenin and P. Parthasarathy, "Transient analysis in discrete time of markovian queues with quadratic rates," *Southwest Journal of Pure and Applied Mathematics*, pp. 1–15, July 2000.
- [23] D. Leonard, V. Rai, and D. Loguinov, "On lifetime-based node failure and stochastic resilience of decentralized peer-to-peer networks," in *ACM SIGMETRICS*, June 2005, pp. 26–37.
- [24] D. Liben-Nowell, H. Balakrishnan, and D. Karger, "Analysis of the evolution of peer-to-peer systems," in *ACM PODC*, July 2002, pp. 233–242.
- [25] S. Liu, K.-H. Cheng, and X. Liu, "Network reliability with node failures," *Networks*, vol. 35, no. 2, pp. 109–117, Mar. 2000.
- [26] D. Loguinov, A. Kumar, V. Rai, and S. Ganesh, "Graph-theoretic analysis of structured peer-to-peer systems: Routing distances and fault resilience," in *ACM SIGCOMM*, Aug. 2003, pp. 395–406.
- [27] G. Manku, M. Naor, and U. Weider, "Know thy neighbor's neighbor: the power of lookahead in randomized p2p networks," in *ACM STOC*, June 2004, pp. 54–63.
- [28] L. Massoulié, A.-M. Kermarec, and A. Ganesh, "Network awareness and failure resilience in self-organising overlay networks," in *IEEE Symposium on Reliable Distributed Systems*, Oct. 2003, pp. 47–55.
- [29] G. Pandurangan, P. Raghavan, and E. Upfal, "Building low-diameter peer-to-peer networks," *IEEE Journal on Selected Areas in Communications*, vol. 21, no. 6, pp. 995–1002, Aug. 2003.
- [30] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker, "A scalable content-addressable network," in *ACM SIGCOMM*, Aug. 2001, pp. 161–172.
- [31] S. Resnick, *Adventures in Stochastic Processes*. Birkhäuser, Boston, 2002.
- [32] S. Rhea, D. Geels, T. Roscoe, and J. Kubiawicz, "Handling churn in a dht," in *USENIX Annual Technical Conference*, June/July 2004, pp. 127–140.
- [33] A. Rowstron and P. Druschel, "Pastry: Scalable, decentralized object location and routing for large-scale peer-to-peer systems," in *IFIP/ACM Middleware*, Nov. 2001, pp. 329–350.
- [34] J. Saia, A. Fiat, S. Gribble, A. Karlin, and S. Saroiu, "Dynamically fault-tolerant content addressable networks," in *IPTPS*, Mar. 2002, pp. 270–279.
- [35] S. Saroiu, P. Gummadi, and S. Gribble, "A measurement study of peer-to-peer file sharing systems," in *MMCN*, Jan. 2002, pp. 156–170.
- [36] A. Shwartz and A. Weiss, *Large Deviations for Performance Analysis*. Chapman and Hall, 1995.
- [37] K. Sripanidkulchai, B. Maggs, and H. Zhang, "Efficient content location using interest-based locality in peer-to-peer systems," in *IEEE INFOCOM*, Mar. 2003, pp. 2166–2176.
- [38] I. Stoica, R. Morris, D. Karger, M. Kaashoek, and H. Balakrishnan, "Chord: A scalable peer-to-peer lookup service for internet applications," in *ACM SIGCOMM*, Aug. 2001, pp. 149–160.
- [39] W. Willinger, M. Taqqu, R. Sherman, and D. Wilson, "Self-similarity through high-variability: Statistical analysis of ethernet lan traffic at the source level," in *ACM SIGCOMM*, Aug. 1995, pp. 100–113.
- [40] R. Wolff, *Stochastic Modeling and the Theory of Queues*. Prentice Hall, 1989.



**Derek Leonard** (S'05/ACM S'06) received the B.A. degree (with distinction) in computer science and mathematics from Hendrix College, Conway, Arkansas, in 2002.

Since 2002, he has been a PhD student in the Department of Computer Science at Texas A&M University, College Station, TX. His research interests include peer-to-peer networks, optimization-based graph construction, and large-scale measurement studies of the Internet.



**Zhongmei Yao** (S'06/ACM S'06) received the B.S. degree in engineering from Donghua University (formerly China Textile University), Shanghai, China, in 1997 and the M.S. degree in computer science from Louisiana Tech University, Ruston, Louisiana, in 2004.

She is currently a PhD student in the Department of Computer Science at Texas A&M University, College Station, TX. Her research interests include peer-to-peer systems, Markov chains, and stochastic network modeling.



**Vivek Rai** received the B.Tech. degree in electronics and communication engineering from Indian Institute of Technology, Guwahati, India, in 2001, and the M.E. degree in computer engineering from Texas A&M University, College Station, Texas, in 2004.

Since 2005, he has been a PhD student at the Vrije Universiteit, Amsterdam, The Netherlands. His research interests include peer-to-peer networks, network creation games, and performance modeling of communication networks.



**Dmitri Loguinov** (S'99–M'03) received the B.S. degree (with honors) in computer science from Moscow State University, Moscow, Russia, in 1995 and the Ph.D. degree in computer science from the City University of New York, New York, in 2002.

Since 2002, he has been an Assistant Professor of computer science with Texas A&M University, College Station, TX. His research interests include peer-to-peer networks, video streaming, congestion control, Internet measurement and modeling.