

Characterizing files in the modern Gnutella network

Daniel Stutzbach · Shanyu Zhao · Reza Rejaie

© Springer-Verlag 2007

Abstract The Internet has witnessed an explosive increase in the popularity of Peer-to-Peer (P2P) file-sharing applications during the past few years. As these applications become more popular, it becomes increasingly important to characterize their behavior in order to improve their performance and quantify their impact on the network. In this paper, we present a measurement study on characteristics of available files in the modern Gnutella system. We develop two new methodologies to capture accurate snapshots of available files in a large-scale P2P system. These methodologies were implemented in a parallel crawler that captures the entire overlay topology of the system where each peer in the overlay is annotated with its available files. We have captured more than 50 snapshots of the Gnutella system that span over 1 year period. Using these snapshots, we conduct three types of analysis on available files: (1) Static analysis, (2) Topological analysis, and (3) dynamic analysis. Our results reveal several interesting proper-

ties of available files in Gnutella that can be leveraged to improve the design and evaluation of P2P file-sharing applications.

1 Introduction

During the past few years, the Internet has witnessed an explosive increase in the popularity of Peer-to-Peer (P2P) file-sharing applications. Today's popular P2P file-sharing applications such as eDonkey, FastTrack, and Gnutella have more than one million users each at any point of time [2] and make up a significant portion of Internet traffic [3]. These applications are primarily used for exchanging multimedia files among participating peers. Participating peers collectively form an overlay that is used to search for desired files among those available throughout the system. As file sharing applications become more popular, characterizing their behavior becomes increasingly important because it reveals performance limitations of these applications in practice as well as their impact on the network. To fully characterize the behavior of file-sharing applications, three equally important and related aspects of these applications should be examined through measurement: (1) overlay topology [4,5], (2) query workload [6], and (3) available files [7]. In particular, characterizing available files among participating peers is valuable for several reasons. First, it reveals the properties, distribution, and heterogeneity of the contributed resources (i.e., storage space and available files) by individual users in the system. Second, it allows us to identify any potential design anomaly that might be exposed in a practical setting or any opportunity that can be used to improve performance of these systems. Third, collected traces and

This paper extends and supplants the earlier version of this paper presented at MMCN 2006 [1]. This material is based upon work supported by the National Science Foundation (NSF) under Grant No. Nets-NBD-0627202, CAREER Award CNS-0448639, and an unrestricted gift from Cisco Systems. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF or Cisco.

D. Stutzbach · S. Zhao · R. Rejaie (✉)
Department of Computer and Information Science,
University of Oregon, Eugene, OR, USA
e-mail: reza@cs.uoregon.edu

D. Stutzbach
e-mail: agthorr@cs.uoregon.edu

S. Zhao
e-mail: szhao@cs.uoregon.edu

derived characteristics of available files through measurement can be also used to conduct more realistic simulations or analytical modeling of available files in P2P systems.

During the past few years, a handful of previous studies have characterized the distribution of shared files in various P2P file-sharing applications [7–10]. While these studies shed an insightful light on the characteristics of files in file-sharing applications, they have several limitations. First, almost all the previous studies have focused on a small subset of peer population in their target file-sharing applications (i.e., less than 20k peers). However, to our knowledge, none of these studies have verified whether the captured peers are indeed representative of the entire population. Second, many of the previous studies (except [7,8]) are more than 3 years old and thus rather outdated. During the past few years, P2P file-sharing applications have significantly grown in size and have incorporated new features. In particular, most popular file sharing applications have adopted a two-tier architecture along with new search mechanisms (e.g., dynamics querying in Gnutella) to improve their scalability. However, the effect of these new features on the characteristics of files has not been studied. Third, previous studies have not examined the dynamics of file characteristics over time nor any possible correlation between the overlay topology and file distribution.

In this paper, we empirically characterize available files across reachable peers in the modern Gnutella network. We present two new measurement methodologies to capture an accurate snapshot of the system at a particular point of time. A snapshot contains participating peers in the system, available files at each peer and the pair-wise connectivity among peers (i.e., the overlay topology). We have developed a fast P2P crawler called Cruiser [11]. Using Cruiser, we have captured more than 50 snapshots of the files available in Gnutella that span over 1 year and each snapshot contains more than 100 million distinct files in each snapshot. Using these snapshots, we conduct the following analysis on shared files in Gnutella:

1. *Static analysis.* We examine properties of contributed resources (i.e., files and storage space) by participating peers in individual snapshots of the system.
2. *Topological analysis.* We investigate any potential correlation between the pattern of file distribution among peers and the overlay topology.
3. *Dynamic analysis.* We study variations in the popularity of individual files in the system, the variations in available files at individual peers, and the percentage of changes in the top- N files over time.

We also leverage an unbiased sampling technique for P2P system that we developed in our earlier work [12,13] to select a representative subset of peers. This sampling technique is implemented in a light-weight sampling tool, called *ion-sampler*. We use *ion-sampler* to derive static properties of available files in Gnutella and validate the similar characteristics that are derived from complete snapshots captured by Cruiser.

Our main findings can be summarized as follows: (1) Free riding has significantly decreased among Gnutella users during the past few years and is significantly lower than other P2P file-sharing applications such as eDonkey. (2) The number of shared files and contributed storage space by individual peers are heavily skewed. (3) The popularity of individual files is also heavily skewed, with a handful of files being extremely popular while a majority of files are very unpopular. (4) The most popular file type is the MP3 file, which accounts for over two-thirds of all files and one-third of all bytes. Both the popularity and space occupied by video files has tripled over the past few years. Furthermore, the number of video files are less than one-tenth of audio files but they occupy 25% more bytes. Overall, 93% of bytes in the system are occupied by audio and video files. (5) Files are randomly distributed throughout the overlay, and there is no strong correlation between the available files at peers that are one, two, or three hops apart in the overlay topology. (6) The set of files shared by individual peers changes slowly, with the median peer making a change every few days. Over the entire system, more popular files experience larger variations in their popularity. Furthermore, the popularity of most files follows a roughly similar pattern of initial increase for a few months until they reach their peak popularity and then gradual decrease after staying at the peak popularity for a period of time. (7) The percentage of changes in the top- N files over time does not strongly depend on N . The top- N list changes roughly 20% over 1 week. Furthermore, 60% of files that remain in the top- N list for at least 1 week, will remain in the top- N for at least 10 weeks.

Why characterize Gnutella? We conducted our empirical study of Gnutella based on a number of considerations. First, Gnutella is one of the most popular P2P file-sharing networks on the Internet [2], with millions of simultaneous users. Second, Gnutella has a protocol hook that allows a list of shared files to be accurately and efficiently extracted from a peer. Finally, Gnutella is one of the most studied P2P systems in the literature. This enables us to compare and contrast the behavior of modern Gnutella with earlier empirical studies of Gnutella and gain insights on changes in the system.

The rest of this paper is organized as follows. We provide an overview of modern Gnutella in Sect. 2. Section 3 describes the challenges in capturing accurate snapshots and describes our measurement methodology and tools. Sections 4, 5, and 6 present static analysis, topological analysis, and dynamics analysis of files in the Gnutella network, respectively. Section 7 provides an overview of previous studies. Finally, Sect. 8 concludes the paper and sketches our future plans.

2 Overview of modern Gnutella

Gnutella is widely regarded as the first fully decentralized P2P file-sharing system. However, it has evolved considerably since its initial release in early 2000, and grown dramatically in size (especially over the last couple of years [4, 14]). Today, Gnutella is one of the largest P2P networks in operation [2]. Similar to many unstructured P2P networks, each Gnutella peer joins the network by establishing TCP connections to several existing peers in the system. In the original Gnutella protocol, participating peers form an unstructured overlay topology that is used by each peer to perform flood-based searches among other peers. To improve the scalability of the original protocol, modern Gnutella clients adopt a two-tier overlay structure along with a dynamic query distribution mechanism [15].

Two aspects of Gnutella are pertinent to our study. First, because one of our goals is to examine correlations between the distribution of shared files and the overlay topology, a general understanding of Gnutella's overlay topology is required. Second, we use Gnutella's Browse-Host extension [16] to acquire the list of files shared by each peer.

Two-tier topology. As shown in Fig. 1, modern Gnutella clients implement a two-tiered overlay structure by dividing peers into two groups: *ultrapeers* (or super-peers) and *leaf peers*. Each ultrapeer neighbors with several other ultrapeers to form the top-level overlay. The majority of the peers are leaf peers and connect

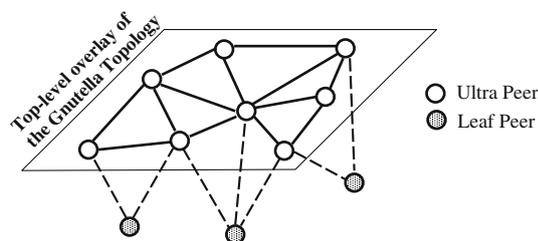


Fig. 1 Two-tier topology of modern Gnutella

to the overlay through a few (two to three) ultrapeers. High-bandwidth, unfirewalled leaf peers become ultrapeers on demand in order to maintain a proper ultrapeer-to-leaf ratio. When a leaf connects to an ultrapeer, it uploads hashes of its filename keywords to that ultrapeer. This allows the ultrapeer to only forward messages to those leaf peers that might have matching files. Leaf peers never forward messages. This approach reduces the number of messages forwarded toward leaf peers which in turn improves the scalability of the system.

The Browse-Host extension. One important reason for selecting modern Gnutella for file characterization is its support for a suite of open and moderately well-documented protocols [17]. The Browse-Host extension [16] is an extension of Gnutella protocol that enables one peer to view the list of files shared (called a *sharing list*) by another peer. This feature allows users with similar interests to learn about new material which may appeal to them. Browse-Host is supported by the two major Gnutella implementations, BearShare, and LimeWire, among others. These two implementations collectively make up roughly 95% of Gnutella ultrapeers [4]. This implies that the Browse-Host extension is able to capture the shared files by most peers in Gnutella.

3 Measurement methodologies

Our goal is to capture a *snapshot* of the Gnutella network at a given point in time that contains (1) addresses of all participating peers and the pair-wise connectivity among them (i.e., an overlay snapshot), and (2) names and sizes of available files at each participating peer in the overlay (i.e., a file snapshot). In essence, we need to capture snapshots of the overlay topology and annotate each peer with its available files. A common approach to capture a snapshot is to use a P2P crawler. Given a set of initial peers, a crawler contacts individual peers to capture their available files and collect information about their neighboring peers in the session. The crawler progressively learns about more peers in the session and contacts them until no other new peer is discovered. However, because the overlay changes as the crawler operates, snapshots captured by a crawler are inherently *distorted* [18]. More specifically, as the crawler explores the overlay, many peers join or leave the system which results in changes in the overlay topology and possibly changes in the set of available files throughout the system. Therefore, the captured snapshot contains some peers that have departed or arrived during the crawl. This problem is further aggravated in large overlays since the arrival of a sufficiently large number of new

peers may significantly increase the duration of a crawl and thus inflate the population of peers in a snapshot.¹

Previous studies implicitly addressed this problem by adopting one of the following ad-hoc sampling schemes to capture a partial snapshot of a P2P system: (1) *Partial Snapshot Through a Short Crawl*: Some studies periodically capture a small subset of participating peers (i.e., a partial snapshot) through a short crawl [10], (2) *Periodic Probe of a Fixed Group*: Some studies identify a subset of participating peers (using a partial crawl or passive monitoring) and periodically probe the same group of peers to collect information about their available files [9]. In the absence of any solid understanding of churn and file characteristics in P2P systems, it is not clear whether these sampling strategies capture a representative population of peers. In our earlier work, we showed that these ad hoc sampling techniques could lead to significant bias toward the properties of the short-lived or high-degree peers [12].

We developed the following two measurement methodologies to capture lists of shared files in the Gnutella network: (1) capturing complete snapshots, (2) capturing unbiased samples. We describe both of these measurement approaches in the following subsections. Capturing complete snapshots is our primary characterization technique, and we use unbiased sampling for validation.

3.1 Capturing complete snapshots

Our goal is to capture the entire population of participating peers in the Gnutella network (i.e., a complete snapshot) within a short period to minimize any potential bias in our characterization. Note that the time required to obtain the list of available files at individual peers is significantly longer than the time to obtain neighbor information. For example, obtaining a list of neighbor peers from a peer may take less than a second whereas the list of available files may take several minutes to download if the peer has many files. In summary, a topology crawl is significantly faster than a content crawl. Therefore, we decouple topology and content crawls to improve the accuracy of captured snapshots and collect our snapshots in three steps as follows. First, we conduct a *topology crawl* to quickly capture all participating peers and their pair-wise connectivity, i.e., capturing the overlay topology. Second, we conduct a *content crawl* and collect the list of files available at each one of the peers identified during the topology crawl. Third, once the content crawl is completed, we conduct another topology crawl in order to identify those long-lived peers

in the initial snapshots that remained in the system during the entire content crawl. At the end of these three steps, we create a snapshot of the overlay topology where each node is annotated with its available file and a label that describes whether it is long or short lived. Since some of the captured peers in the first topology crawl depart the system during the content crawl, the collected content in our measurement is slightly biased toward peers with longer uptime. This motivates our second technique, unbiased sampling, described in the next subsection.

To minimize distortion in captured snapshots, we have developed a parallel P2P crawler, called Cruiser [11], that can crawl an overlay orders of magnitude faster than any previous crawler. Cruiser achieves this goal by significantly increasing the degree of concurrency in the crawling process. Toward this end, Cruiser adopts a master-slave architecture where each slave crawls hundreds of peers simultaneously and the master coordinates among multiple slaves. This architecture allows us to run Cruiser on multiple co-located boxes to further increase the crawling speed. Using six off-the-shelf 1 GHz GNU/Linux boxes in our lab, Cruiser can perform a topology crawl for more than two million Gnutella peers in less than 15 min, and performs a content crawl within 5.5 h, i.e., capturing the annotated snapshot takes 6 h, (15 min + 5.5 h + 15 min). During the content crawl, Cruiser collects the file name and content hash (SHA1) for each shared file on every reachable peer, resulting in a log file more than 10 GB in size.

3.2 Capturing unbiased samples

Capturing complete snapshots is a heavy-weight operation that provides a tremendous amount of information. However, even with a fast crawler, there will inevitably be some distortion as the network changes while the crawler operates. Sampling is a useful, light-weight alternative that does not suffer from distortion. However, commonly used sampling techniques for measuring P2P systems tend to introduce considerable bias for two reasons [12]. First, the dynamic nature of peers (i.e., churn) can bias results toward short-lived peers, much as naively sampling flows in a router can lead to bias toward short-lived flows. Second, the heterogeneous nature of the overlay topology can lead to bias toward high-degree peers, as peers with a high degree are more likely to be discovered. In our prior work [13], we developed a light-weight sampling tool, called `ion-sampler`, for P2P systems that selects peers uniformly at random. The `ion-sampler` tool performs a random walk and selects the peer at the end of the walk as a sample. In static

¹ In the extreme case, the crawler may never terminate since there are always new peers to contact.

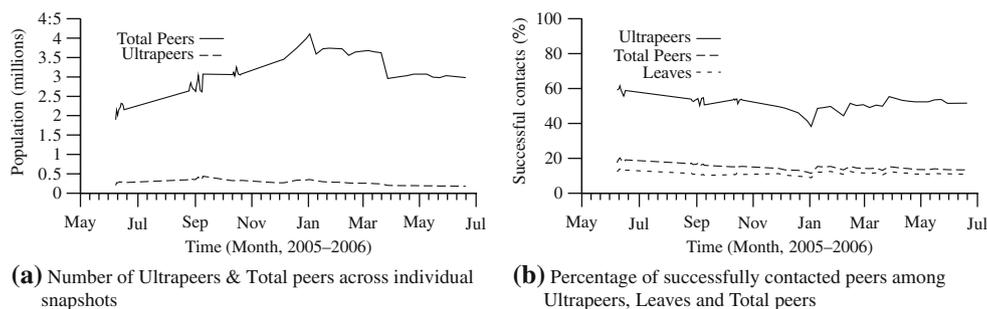


Fig. 2 Population and reachability across captured snapshots

Table 1 The range of percentage of peers that failed during the file crawl for various reasons across all snapshots

Type	TCP refused (%)	Timeout (%)	Conn. lost (%)	App. refused (%)
Ultrapeer	31.30–46.10	1.74–13.25	1.78–6.80	0.97–2.30
Leaf peers	53.53–72.13	1.12–3.85	4.24–23.16	1.27–3.06
All peers	31.63–46.74	1.75–13.03	1.81–7.08	0.97–2.30

graphs, the Metropolis–Hastings method can be used to weight the walk to result in uniform sampling [19–21], eliminating the degree bias. Our prior work has shown [13] that it is also effective in eliminating the degree bias in dynamic networks such as P2P systems. To avoid any bias based on session length, we allow re-sampling; effectively, instead of attempting to sample a peer, we are trying to sample a peer’s state (which may vary with time). The *ion-sampler* tool allows us to gather many samples at once by running many walks in parallel. Our extensive simulations [13] show that the technique yields unbiased samples under a wide variety of peer behavior, degree distributions, and overlay construction techniques. For this study, we use 1,000 parallel walks. We use *ion-sampler* to collect peers, then gather the list of files available from each sampled peer.

While sampling does not suffer from distortion, it provides fewer details than a full crawl. In particular, it does not include a topological snapshot or session length information. Additionally, sampling provides little information about unusual peers, e.g., in the tail of distributions. For these reasons, we only employ sampling to validate the results of our static analysis from complete snapshots.

3.3 Dataset

We have captured more than 50 snapshots of the Gnutella network annotated with the list of files available at each peer. Our snapshots consists of daily snapshots from consecutive days during the following intervals: 6/8/2005–6/18/2005, 8/23/2005–9/9/2005, and 10/11/2005–

10/21/2005.² Furthermore, we have collected 24 snapshots, roughly once per week between 11/29/2005 and 05/30/2006. To avoid any time-of-day effects, each snapshot was captured at approximately the same time-of-day (2 p.m. PDT). These snapshots enable us to examine characteristics of available files in the system over both short and long timescales (i.e., several days and several months). Additionally, we use unbiased samples collected on June 16, 2006 and July 4, 2006 to validate our findings in Sect. 4.

Figure 2a depicts the number of ultrapeers and total peers across captured snapshots during our measurement period. A subset of peers in each snapshot are unreachable by our crawler for one of the following reasons: (1) a firewall or NAT blocking incoming traffic, (2) severe network congestion or an overloaded processor at the peer, (3) the peer departed, or (4) the peer does not support the Browse–Host protocol. Figure 2b shows the percentage of ultrapeers, leaf peers, and the total population in individual snapshots that were successfully contacted by the file crawler. This figure reveals that in average around 50% of ultrapeers, 12% of leaf peers, and 18% of the total population are reachable by the file crawler. Table 1 presents the minimum and maximum value for the percentage of unreachable peers within each group (ultrapeer, leaf, total population) for four different reasons across all snapshots. “TCP refused” column in Table 1 refers to attempts that were rejected by a TCP RST message. “Timeout” column refers to

² We are missing snapshots for 9/6/2005, 9/8/2005, 10/17/2005, and 10/29/2005 due to a mixture of software, network, and power failures.

attempts that did not receive any response. “Conn. lost” column refers to attempts that established a TCP connection, but the connection closed before the file list could be fully retrieved. Finally, “App. refused” refers to attempts where TCP connected, but the request for the list of files was rejected at the application level. In a few cases, our crawler machines were under heavy load by other users during our crawls, leading to a large percentage of timeouts and lost connections. Since ultrapeers are not allowed to be firewalled, any reported connection error for ultrapeers indicates that the contacted peer has departed. However, connection errors for leaf peers might occur due to peer departure or a firewall.³ In our earlier study of Gnutella overlay topology [4], we showed that about half of all leaf peers captured in a topology crawl leave the overlay within a 5 h period. Independent online statistics [22] report that around 70% of leaf peers in the Gnutella network are firewalled. These evidences support the accuracy of the high ratio of connection errors that we experienced for leaf peers (in the column labeled “TCP refused”). In summary, while we aim to capture complete snapshots, we can successfully contact only 20% of all peers in one snapshot (around half a million peers) to obtain their list of available files primarily due to two reasons: (1) the number of departed peers during the long file crawl, and (2) the large number of leaf peers behind firewalls.

3.4 Challenges and problems

We briefly discuss several interesting problems that we experienced during data collection and data processing.

Firewalled peers. Our crawling and sampling tools are unable to contact firewalled peers that cannot accept incoming TCP connections on the Gnutella port. As a consequence, if firewalled peers have significantly different file-sharing characteristics compared to un-firewalled peers, our results may not be representative of the entire Gnutella population. However, we have no reason to believe this is the case.

Low-bandwidth TCP connection. Although Cruiser has a timeout mechanism that closes any idle connections after 20 s, we noticed that some crawls do not complete after the crawling queue becomes empty. Further examinations revealed that around 80 peers in each crawl send their data at an extremely low rate (around 20 bytes per second) which prevents Cruiser from closing their connections. We instructed Cruiser to termi-

nate a crawl a few minutes after its crawling queue becomes empty. Given the negligible number of these misbehaved peers, this should not have any affect on our analysis.

File identity. We use the content hash of a file returned by the target peer to uniquely identify individual files. In our initial measurements, we observed many files with the same name but different content hashes (e.g., setup.exe, login.bmp). This illustrates that the trimmed (or even complete) filenames used by previous studies (e.g., [9]), are not reliable file identifiers. We also discovered around 3,500 files without content hash value in each snapshot and eliminated them from our analysis.

Post-processing. To compute the popularity of individual files in the system, we needed to keep track of more than 100 million distinct files discovered in each snapshot. The large number of files resulted in memory bottlenecks in our analysis. We leveraged the skewed distribution of popularity to address this problem as follows. We divide captured peers in a snapshot into seven segments where each segment contains a random subset of the peers. Then, we calculate the popularity of files within each segment, trimming all files with fewer than ten copies in a segment, and combine all the trimmed segments into a single snapshot. This approximation eliminated several million distinct files and prevented memory bottlenecks during our post-processing. While this prevents us from performing analysis on the least popular files (with fewer than 70 copies in the entire network), it does not significantly affect conducted analysis on more popular files.

4 Static analysis

In this section, we examine the following static characteristics of available files across peers in individual snapshots regardless of their location in the overlay topology: (1) the percentage of free riders, (2) the degree of resource sharing among cooperative peers, (3) the distribution of file popularity, and (4) the distribution of file types. Furthermore, we investigate the variations of these characteristics across our 12-month measurement period and examine to what extent these properties are fundamental rather than ephemeral. We derive these characteristics from both complete snapshots and sampled peers and compare them to verify the accuracy of our results. We also compare our findings with previous studies to illustrate any potential changes in these characteristics of the Gnutella network over the past few years.

³ We are not aware of any reliable technique to distinguish between these two scenarios.

4.1 Percentage of free riders

The success of P2P file-sharing systems depends on the willingness of participating peers to share their files. However, previous studies have frequently reported that participating peers do not have an incentive to contribute their resources (e.g., disk space and network bandwidth) to the system and thus only use resources offered by others, i.e., become “free riders”. In particular, Adar et al. [23] reported that 66% of Gnutella peers were free riders in 2000, while a study by Saroiu et al. [10] found 25% were free riders, with 75% of peers sharing less than 100 files in 2002. A recent study also reported 68% of peers in eDonkey were free riders [7].

Table 2 presents the range (minimum–maximum values) for the number of ultrapeers (second row), leaf peers (third row) and total peers (last row) across all snapshots. We further divide ultrapeers (rows 4 and 5) and leaf peers (rows 6 and 7) into short and long lived based on their presence in the second topology crawl as we discussed in Sect. 3. This grouping reveals any potential differences in the free riding between ultrapeers and leaves as well as long- versus short-lived peers. For each one of the above groups, the corresponding row in Table 2 presents the range of the following properties across all snapshots (1) the number of contacted peers that provided their sharing list (second column), (2) the percentage of free riders (third column), and (3) the average number of shared files among peers in each group.

Table 2 shows several interesting points as follows. The percentage of free riders in Gnutella has significantly dropped from 25% in 2002 [10] to around 13% among all participating peers (i.e., last row), and it is drastically lower than the 68% recently reported in eDonkey [7]. We speculate that several factors have contributed in the observed drop in the percentage of free riders including the increase in access link bandwidth for average Internet users and active marketing efforts by the Gnutella vendors encouraging their users to share. However, an alternative explanation is that the higher

percentage of free riders observed in previous studies is due to measurement bias.

Table 2 reveals that the percentage of free riding among ultrapeers (9.0–10.5%) is somewhat lower than that in leaf peers (10.8–15.9%). However, on average leaf peers share more files. Since leaf peers constitute a larger portion of the total population (86–90%), their behavior has a greater impact on system performance. Long-lived peers appear to have a slightly higher percentage of free riders compared to short lived peers. The correlation between lifetime and number of files shared is inconclusive. Long lived leaf peers share more files than short lived leaf peers, but long and short lived ultrapeers share about the same number of files.

4.2 Degree of resource sharing among cooperative peers

We now focus our attention on cooperative peers and characterize their willingness to contribute their resources (i.e., both files and storage space). During our analysis, we noticed that the sharing lists of many peers contain duplicate files. This occurs because most Gnutella clients simply put various folders with potentially duplicate content under the sharing folder. We have excluded all the duplicate files from captured sharing lists (which account for around ten million files or roughly 9% of all captured files) in all the reported analysis in this paper.

In the earlier version of this paper [1], we depicted the number of shared files and shared bytes per peer as a histogram on log–log scale and concluded that both distributions resemble a power-law due to their approximately linear shape. Presenting these distributions as histogram results in a “messy”, difficult-to-interpret tail of the distribution, due to the small number of peers in that region. This in turn makes it difficult to properly examine the slope of distribution. Presenting the distribution as a Complementary Cumulative Distribution Function (CCDF) yields a smooth tail without losing any information while retaining the linear slope that a power-law distribution exhibits on a log–log scale (since $\int_x^\infty u^{-\alpha} du \propto x^{-\alpha}$). Toward this end, we compute the CCDF of the number of shared files and shared bytes per peer for individual snapshots. To present the variability of these distributions across all snapshots in a compact and clear fashion, we compute the MIN-CCDF and MAX-CCDF, defined as follows:

$$\text{MIN-CCDF}(x) = \min(\text{CCDF}_i(x) \text{ for all snapshots } i)$$

$$\text{MAX-CCDF}(x) = \max(\text{CCDF}_i(x) \text{ for all snapshots } i)$$

Table 2 Percentage of free riders, its breakdown across different groups, and the mean number of shared files among peer

	Contacted peers	Free riders (%)	Files/peer
Ultra	94,329–188,107	8.96–10.54	324–358
Leaf	253,164–395,682	10.81–15.85	347–406
Long-lived ultra	69,441–147,435	9.27–10.78	320–354
Short-lived ultra	24,888–40,672	8.08–9.77	335–371
Long-lived leaf	166,241–240,969	11.39–16.93	369–437
Short-lived leaf	86,923–154,713	9.93–13.79	316–358
Total	410,252–552,702	10.63–13.41	340–393

Figure 3a and b presents the MIN-CCDF and MAX-CCDF for CCDF of shared files and shared bytes across all snapshots, i.e., the CCDFs from all the snapshots fall in between the MIN-CCDF and MAX-CCDF curves in the corresponding graph. Figure 3a and b also includes the CCDF for shared files and shared bytes per peer based on samples collected with `ion-sampler`, to verify whether the results from complete snapshots are significantly affected by distortion or not (as described in Sect. 3). Note that the MIN-CCDF and MAX-CCDF curves in both figures are very close which indicates that the CCDF for shared files and shared bytes per peer across all snapshots are very stable. Furthermore, the CCDFs from sampling are consistent with the results from complete snapshots, except in the tail of the distribution where the precision of sampling is low.

To verify our earlier conclusion, we recall that power-law distributions are *scale-free*. This means that the “slope” (α) of distribution is consistent over many orders of magnitude on a log–log scale. While the histograms presented in the earlier version of this paper appear approximately linear [1], it is not possible to closely fit a line over the entire range of x values due to the subtle downward curvature. This curvature is more apparent in the CCDFs in Fig. 3a and b. Based on these considerations, we revise our earlier conclusions as follows: the

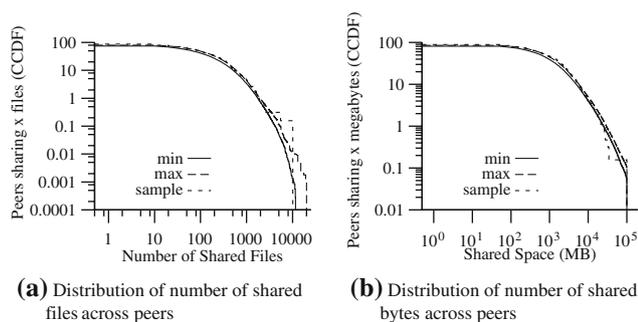


Fig. 3 Distributions of files and bytes

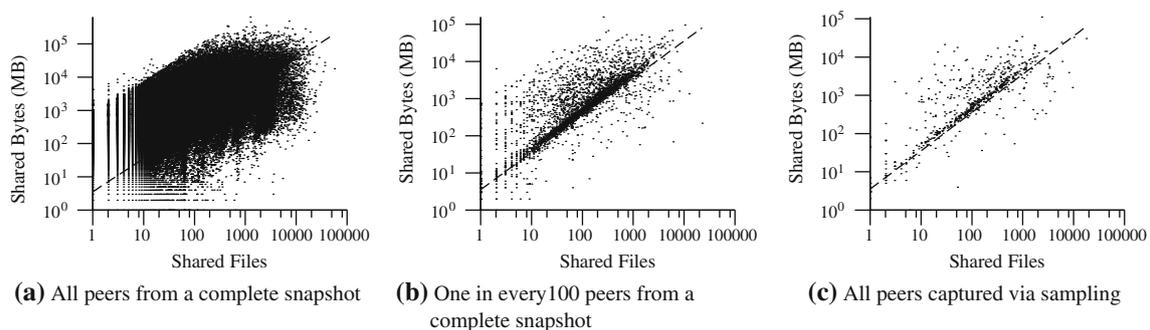


Fig. 4 Scatter-plot of number of shared bytes as a function of number of shared files. Each point represents one peer. The reference line is 3.544 MB per file

distributions of the number of files and bytes shared by individual peers are heavily skewed but do not appear to follow a power-law distribution. Most peers share a moderate amount of files (or bytes) while a few peers contribute an enormous amount of content (or space). The median value for shared files is around 70 files while 0.01% of peers share more than 7,500 files. The median value of shared space is around 650 MB while 0.1% of peers contribute more than 85 GB.

Correlation between contributed files and bytes.

Saroiu et al. [10] reported a strong correlation between the number of shared files and the volume of shared data across Gnutella peers in 2002. Figure 4a shows this correlation as a scatter-plot across all cooperative peers in one snapshot (June 13, 2005). Each point in this figure represents the number of shared files versus the shared disk space for each cooperative peer. We observe a greater variability than what has been reported by Saroiu et al. in 2002. However, this may simply be due to the larger number of data points in our scatter-plot. To verify this issue, we sampled one out of every hundred data points in Fig. 4a and the selected samples are shown in Fig. 4b. Figure 4b reveals a strong correlation between the number of files and the number of bytes contributed by each peer which was obscured in Fig. 4a. We performed a least-squares fit to a line through the origin (0 files = 0 bytes), finding a slope of 3.544 MB per file. This is consistent with the size of a MP3 file that is a few minutes long (i.e., a typical song). Figure 4c shows the correlation between the contributed bytes and files by each peer using the data gathered by `ion-sampler` on June 16, 2006—more than a year after the presented snapshot data in Fig. 4a. Figure 4c fits to the same line and is consistent with the results from complete snapshots.

To gain a deeper insight about the correlation between contributed files and bytes, we derive the CCDF distribution of file size among all files in individual snapshots. Figure 5a depicts the CCDF distribution of file size

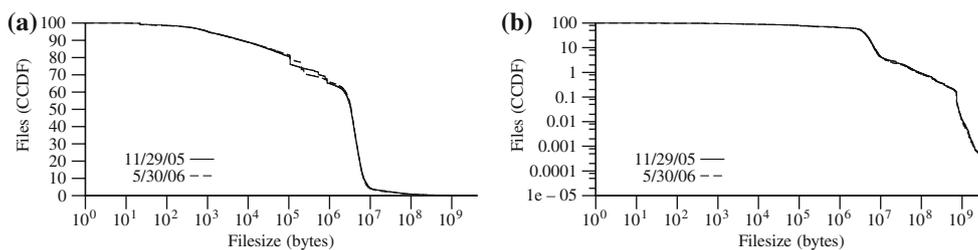


Fig. 5 Distribution of filesize across all files

in two snapshots that are 6 months apart in a linear-log scale. Figure 5a shows that the distribution of file size has very little change between the two snapshots that are 6 months apart. This figure also indicates that 60% of all files are a few megabytes in size (i.e., the graph has a knee between 1 and 10 MB). Furthermore, the majority of other files (around 40%) are smaller than 1 MB in size. Figure 5b shows the same CCDFs for file size distribution of two snapshots on a log-log scale. This figure illustrates that a very small portion of files (around 0.1%) have a size larger than 1 GB.

4.3 File popularity distribution

The distribution of popularity for individual files throughout the system is an important property that shows the degree of similarity and thus availability of individual files among participating peers. Chu et al. [9] showed that the file popularity follows a log-quadratic distribution, which can be viewed as a second-order power-law distribution, among Gnutella peers in 2001. Furthermore, Fessant et al. [7] recently reported a Zipf distribution for the file popularity in eDonkey. However, none of these studies have captured a large number of peers.

Figure 6 shows the range of CCDF for file popularity as a function of its rank (in log-log scale) among cooper-

ative peers across all snapshots. Each snapshot typically contains more than 800 terabytes worth of content in more than 100 million unique files based on information from 0.4 million peers, constituting 18.5% of identified peers. If we assume that unreachable peers have similar profiles, the volume of available content in the Gnutella network is around 4,400 terabytes. Figure 6 illustrates two points: (1) file popularity is heavily skewed, and (2) the distribution of file popularity has remained very stable across the 12-month measurement period. A small number of files are very popular while for most files only few copies are available in the network. Again, we also present data collected using unbiased sampling, which closely matches the snapshot data except for very unpopular files where sampling has poor precision.

4.4 File type analysis

We have also examined the distribution of available files among Gnutella peers across different types of video and audio formats. This basically illustrates what types of content are available in the system and thus exchanged among peers. Chu et al. [9] conducted similar analysis for Gnutella peers in 2001 and reported that audio files constitute 67.2% of files and 79.2% of bytes. However, video files were significantly less popular and only contributed 2.1% of files and 19.1% of bytes. Using our snapshots, we analyze the various types of audio and video files based on file extensions.

Figure 7a and b depicts the contribution of the most common file types in the total number of files and the total number of bytes across cooperative peers, respectively. These box-and-whisker plots present the variations of the contribution by different file types across all snapshots by showing the 75th percentile (the top of the box), 25th percentile (the bottom of the box), the median value (the line through the box), the maximum value (the top of the line), and the minimum value (the bottom line). The distribution of file types across all snapshots seems fairly consistent with respect to both

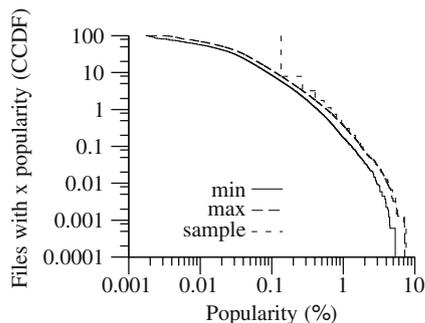


Fig. 6 Distribution of file popularity across unique files

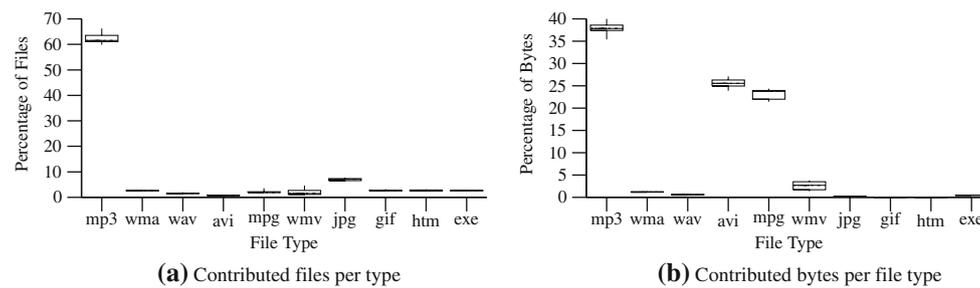


Fig. 7 Breakdown of file types based on the contributed number of files or bytes by each file type

files and bytes (i.e., all boxes are small). Note that except for MP3 files, each file type individually makes up at most a few percent of all files. MP3 audio files are significantly more popular than any other file type, making up two-thirds of all files and occupying more than one-third of all bytes in the system.

Although non-media files (jpg, gif, htm, exe, txt) are among the top-ten most popular types, audio and video files collectively occupy more than 93% of bytes and make up more than 73% of all files in the system. These figures also reveal that video files (avi, mpg, wmv) are significantly larger than other file types in the system. The audio files (MP3, wma, wav) account for 67% of files and 40% of bytes whereas video files constitute around 6% of files but 52.5% of bytes among Gnutella peers (i.e., most files are audio files, but most bytes are in video files). Comparing to the reported results by Chu et al. in 2001 [9], video files have become three times more popular and occupy almost three times more space in the system.

We are unable to apply our sampling technique for file type analysis because it samples random *peers* not random *files*. Consequently, analyzing the file data gathered with *ion-sampler* shows significant correlations based on the preferences of the peers collected. For example, in one instance we saw that 16% of files had the jpg extension, due to sampling two peers with large jpg archives.

5 Topological analysis

In this section, we investigate whether copies of a given file are located at close-by peers (i.e., *topologically clustered*) or at randomly scattered peers throughout the overlay. Understanding this issue would be useful in the design and evaluations of search techniques in file-sharing applications. There are two factors that affect the location of available copies of individual files throughout the overlay: (1) scoped search and (2) the dynamics of peer participation (or churn). To accommodate scalability, each searching peer often examines available files among nearby peers (i.e., conducts scoped search). This suggests that a single copy of a file gradually diffuses outward from the original location to other nearby peers and some type of topological clustering might exist. However, due to the dynamics of peer participation, the location of a peer in the overlay changes over time and could prevent such clustering. The key question is which one of these two factors has the dominating effect. To answer this question, we generate an annotated snapshot of the overlay topology where each peer is annotated with its available files. Then, we explore topological clustering from the following two perspectives.

Per-file perspective. We conduct trace-driven simulations of flood-based querying over our annotated overlay topology. Figure 8a depicts the CDF of the minimum

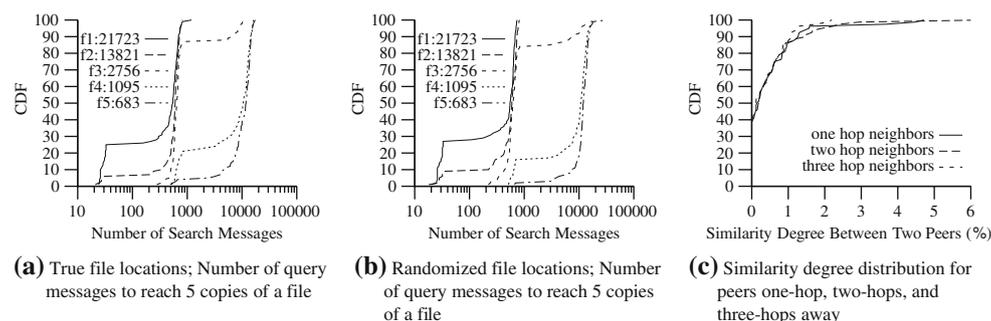


Fig. 8 Figures demonstrated the lack of a strong correlation between topology and files

number of query messages to find five copies of a target file from 100 randomly selected peers in the overlay. The lines in this figure correspond to five target files with different popularity from our static analysis. The abrupt steps in the CDF indicates one-hop increases in the scope of the search. Note that each CDF has at most two steps. This implies that most of the 100 randomly selected peers find all 5 copies of each target file among other peers that are either n or $n + 1$ hops away where n inversely depends on the log of the target file's popularity.

Clearly, more search messages are required for less popular files. If significant topological clustering exists, a few searches will complete using very few messages (i.e., fewer hops) while most searches will require a much larger number of messages (i.e., many more hops). However, the pattern of increase in the number of search messages for less popular files in Fig. 8a indicates that files are randomly distributed.

To verify this conclusion, we randomized the placement of available files without changing their popularity, guaranteeing that no topological clustering exists in the overlay. Figure 8b shows the CDF of the number of required query messages to reach five copies of the same target files over the annotated topology with randomized file placement. The high similarity between Fig. 8b and a supports our conclusion that no significant topological clustering exists in the locations of a given file throughout the overlay. It is worth noting that our annotated overlay topology is incomplete due to the unreachable peers behind firewalls and departed peers. This implies that the number of messages are likely to be smaller in practice. However, the unreachable peers should not significantly affect topological clustering. More specifically, if a file is available only in a certain region of the overlay, the clustering property is preserved even if information is absent for some fraction of the peers in that region.

Per-peer perspective. We also investigate the possibility of any topological clustering from per-peer's perspective by examining the degree of similarity between available content at a random peer and all its one-, two-, and three-hop neighbors. The similarity degree between two peers is defined as the number of common files divided by the smaller size of the sharing lists for the two peers. Figure 8c shows a separate CDF for the average similarity degrees of 100 randomly selected peers with their one-, two- and three-hop neighbors. If topological clustering were present, the similarity degree would decrease as the distance between comparing peers increases. However, since the distributions are nearly identical, Fig. 8c clearly illustrates that there is no

correlation between the similarity degree and distance between two peers.

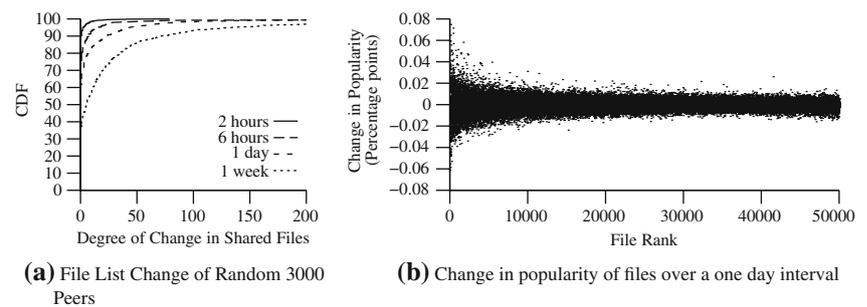
These results suggest that churn is the dominant factor in determining the distribution of files throughout the overlay. In our prior work, we observed that more than half of the Gnutella peers in a given snapshot will depart within 5 h [24]. Any time a peer joins the overlay, it attaches itself as a leaf peer to several random ultrapeers. Furthermore, qualified leaf peers may become ultrapeers in order to maintain a proper ultrapeer-to-leaf ratio. Hence, the rapid changes in the overlay topology prevents formation of such topological clustering. This finding is important for two reasons: (1) measurement studies may sample the list of files from random peers in any part of the overlay topology, and do not need to capture the entire network; (2) simulation studies may randomly distribute available files among participating peers regardless of location. However, the number of files per peer and the total number of copies of each file should follow appropriate distributions, examined in Sect. 4. While previous studies have frequently assumed these properties, to our knowledge, they had not previously been empirically verified.

6 Dynamic analysis

In this section, we turn our attention to the dynamic properties of available files in Gnutella. More specifically, we investigate how various properties of available files change over time. Prior studies examined changes in the popularity of queries and exchanged files in P2P file-sharing systems (e.g., [6,25]). However, to our knowledge, no study has previously explored the dynamic characteristics of stored files in P2P systems. For these analysis, we leverage our three sets of daily snapshots that were captured once every 2 months, as well as several months of weekly snapshots. This dataset allows us to explore dynamic properties over short timescales (i.e., hours and days) as well as long timescales (i.e., months). We explore the following three dynamic properties of shared files by peers in Gnutella: (1) variations in shared files by individual peers, (2) variations in popularity of individual files, (3) long-term evolution of file popularity, and (4) variations in top- N files.

6.1 Variations in shared files by individual peers

Our goal is to determine how rapidly the available files at individual peers change with time. These dynamics show whether past information about the available files at individual peers can be reliably used in future searches.

Fig. 9 Changes in popularity

There are two types of change that can occur to the list of shared files at each peer. First, the user may add new files, either by downloading them from other peers or by manually adding them to the shared folder. Second, the user may remove files, by moving (or deleting) those files from the sharing folder. We note that dynamic IP addresses may introduce error into our results when a peer departs from the system and another peer joins the system with the same IP address. However, our prior study on churn revealed that such events are not common [24]. We define the total number of added and deleted files at a single peer as the *degree of change* to quantify both types of changes in shared files at each peer.

Figure 9a depicts the CDF of the degree of change for 3,000 randomly selected peers over 4 different timescales: 2 h, 6 h, 1 day and 1 week. This figure reveals that 36% of monitored peers do not change their sharing lists over a 1-week interval. That number increases to 69, 80 and 90% over 1 day, 6 h and 2 h intervals, respectively. On the other hand, 90% of peers change less than 10 files over a 6 h period, less than 25 files over a 1 day period, and less than 80 files over a week. Given the average number of shared files by each peer (which is around 350 from Table 2), this result indicates that the variations of shared files by individual peers is rather small especially over several days. While this finding is rather intuitive and expected, Figure 9a allows us to quantify the distribution of the degree of change across different timescales. Finally, this result implies that caching information about the available files at other peers (especially over the timescale of few days) can be a highly effective bandwidth-saving strategy in P2P systems.

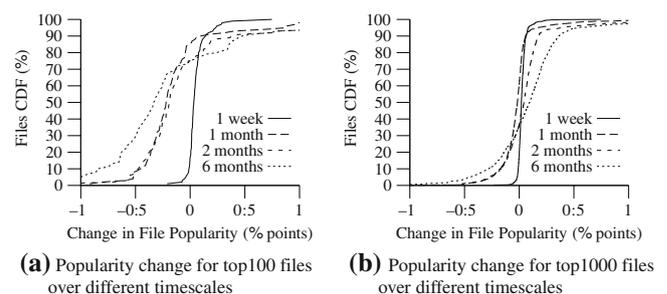
6.2 Variations in popularity of individual files

In this subsection, we focus on the effect of changes in shared files by each peer on the popularity of individual files across the system. Understanding the dynamics of file popularity can be used to determine how often the popularity of available files should be sampled. To eliminate the effects of a varying peer population across

different snapshots, we define the *popularity* of a file as the percentage of successfully contacted peers with the file. Given the random distribution of files among peers, the popularity can be interpreted as the probability of having that file at a random peer. We define the *change in popularity* of a given file over interval τ as the difference between its popularity at the beginning and the end of such an interval (in percentage points).

Figure 9b depicts the change in popularity of 50,000 randomly selected files over a 1 day interval as a function of their popularity rank in a scattered-dot plot where each dot corresponds to a particular file. The number of copies for the highest and lowest ranked files in this figure are roughly 28,000 and 130, respectively. This figure clearly demonstrates the effect of file popularity on the variations of its popularity over a 1-day period. This figure shows that the most popular files (top 15,000) exhibit significantly larger variations (increase or decrease) in their popularity compared to the rest of the files. Note that the changes in popularity for most of the files are within 0.02 points. While the variations in file popularity rapidly drop with file rank among the top 15,000 files, it becomes relatively stable across the least popular files and remains below 0.01 points. In summary, a group of the most popular files experience wider variations in their popularity than unpopular files.

To study popularity dynamics in further detail, we focus on the top 100 and top 1,000 most popular files in a particular snapshot and examine their popularity variations over different timescales. Figure 10a and b

**Fig. 10** Changes in file popularity over different timescales

plots the CDF graph of the change in popularity for the top 100 and top 1,000 files, respectively, for intervals of 1 week, 1 month, 2 months, and 6 months. Figure 10a and b individually show that the range of variations in popularity of top- N files expanded with time. However, comparing these figures clearly illustrates that (1) for any given timescale, more popular files exhibit significantly larger variations in their popularity, and (2) the popularity of more popular files changes more rapidly with time.

6.3 Long-term evolution in popularity

Two other interesting questions are “whether the long-term evolution of file popularity follows a particular pattern?”, and “what does that pattern look like?” Intuitively, the popularity of a new file should gradually increase with some average rate until it reaches its peak popularity. Then, it remains close to its peak popularity for a while before its popularity gradually drops with time. However, the rate of change in popularity, the peak popularity and the time that a file remains around its peak popularity could significantly vary across different files.

To examine the pattern of long-term evolution in file popularity, we track the pattern of changes in the popularity of top-ten files in the first, middle, and last snapshots across all other snapshots over our 1-year measurement period. Figure 11a shows the evolution of the popularity in the top-ten file in our first snapshot (captured in June 8, 2005) during the 1-year interval *after* they appear in the top-ten list. Despite the fluctuations in the popularity of individual files, popularity of all files exhibit an overall downward trend. Figure 11c depicts the evolution of popularity for the top-ten files in our last snapshot during the 1-year period *before* they appear in the top-ten list. This figure shows that most files transition from very low popularity to top-ten list in the span of a few months. Some of these files experienced a rather

long periods of relative obscurity before making rapid gains in popularity. However, two files transitioned gradually from obscurity to popularity over the course of the year. Also note that two of the top-ten files in the last snapshot appears to have passed their peak popularity. Figure 11b provides a third point of view, examining the pattern of evolution in the popularity of top-ten files in the middle of our measurement window. This figure illustrates how the popularity of top-ten files have changed during a 6-month period before and after they become top-ten. This figure shows that one of the top-ten files in this snapshot was more popular at the beginning of our measurement and has been gradually losing popularity, while the other top-ten files are newcomers that have been gaining popularity during the same period. Afterward, they all gradually lose their popularity over the following 6-month period.

In summary, these figures show that files rise in popularity at varying (and often rapid) rates until they reach their peak popularity. Then, they enter a period of slow decay where their popularity gradually drops. The top-ten most popular files at each snapshot is composed of three groups of files: (1) files whose popularity is still rising toward their peak, (2) files whose popularity is around their peak, and (3) files whose popularity is dropping.

6.4 Variations in top- N files

We examine how the changes in file popularity affects the identity of the top- N files in the system over both short and long timescales. In essence, this explores the rather subtle relationship between the popularity and the ranking of top- N files. Figure 12a depicts the percentage of the top- N files on day $x - 1$ (starting at 06/08/2005) that remain in the top- N files on day x (i.e., daily drift in the top- N list) for four values of N . Note that the y-axis begins at 89%, indicating that the top- N list is highly stable from one day to the next. The changes in the top- N

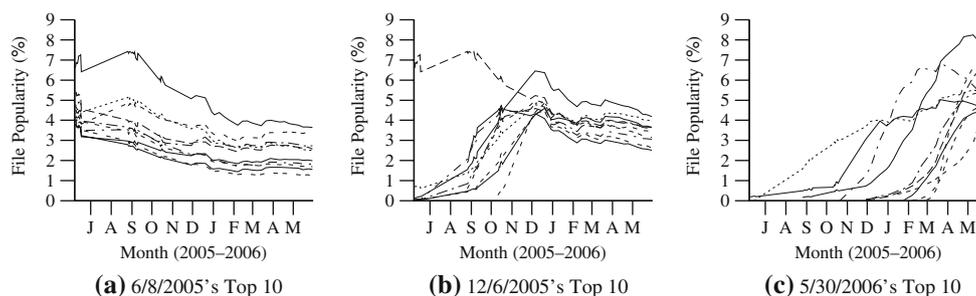
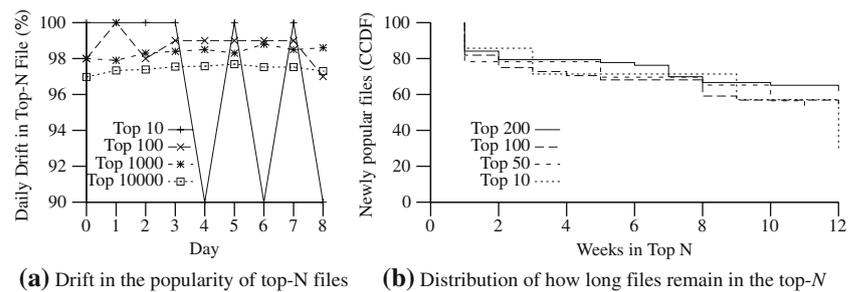


Fig. 11 Changes in the identity of Top- N files

Fig. 12 Distribution of how long files remain in the top- N



list are more visible for smaller values of N mainly due to two reasons: (1) the sensitivity to the small number of files (one change in the list results in $\frac{1}{N}$ variations) coupled with (2) the noisy variations in the popularity of individual files that can easily change a couple of files at the bottom of the top- N list. The main conclusion from this figure is that the identity of the top- N list remains highly stable across consecutive days for different values of N . Its minor variations is inversely proportional to N .

We also investigate the stability the top- N list over longer time scales, by examining how long files remain in the top- N . We use the “create-based method” [10] to avoid bias toward files that remain in the top- N files for a short period of time. Toward this end, we divide our measurement window into two halves. For any file that enters the top- N during the first half of the window, we measure the number of weeks that a file remains in the top- N . This technique provides us with unbiased measurements for files that remain in the top- N for up to one half of the window (or 12 weeks). We can count how many files remain in the top- N for more than 12 weeks, but cannot accurately compute their distribution.

Figure 12b depicts the CCDF for the percentage of changes in the top- N files over time for different values of N . This figure shows that approximately 20% of files remain in the top- N for only 1 week. However, approximately 60% of files that remain in the top- N list for at least one week, will continue to remain in the top- N for at least 10 weeks. More interestingly, the percentage of changes in the top- N files over time does not appear to have a strong dependency to N . Note that over a given period, the number of changed files in the top- N list is proportional to N . Therefore, the *percentage* of changes in the top- N files is relatively similar for different values of N .

7 Related work

Several measurement studies have examined different properties of P2P file-sharing networks including (1) dynamics of peer participations (i.e., churn) [24,26],

(2) overlay topology structure [4,27,28,5], (3) query traffic [6], (4) data traffic [29,25,30], and (5) shared files [7,8]. We are aware of only two other studies that focus on the characteristics of files shared by users. First, Fessant et al. [7] examined characteristics of available files, using data collected from 12,000 eDonkey clients over a 3-day period in 2003. They showed that the popularity of files stored in file-sharing systems is heavily skewed, following a Zipf distribution. When two peers have ten files in common, there is an 80% chance they have at least one more file in common. The probability is close to 100% if they have at least 50 files in common. Second, Liang et al. [8] recently analyzed the nature and magnitude of deliberately corrupted files (“pollution”) in Kazaa. To combat P2P sharing of copyrighted content, some companies intentionally inject decoy files, which have the same file name as a popular song. Liang et al. developed a multi-threaded crawler that queries all 30,000 Kazaa super-nodes for 7 popular songs over the course of 1 hour. They showed that the popularity of different versions of a song also follows a Zipf distribution. For most of the seven popular songs, over 50% of the copies are polluted.

A few other studies have examined the files shared by users as part of broader measurement studies on P2P systems. In 2001, Chu et al. [9] studied peer churn and the distribution of file popularity. They found that file popularity follows a log-quadratic distribution (which can be thought of as a second-order Zipf distribution). Saroiu et al. [10] examined many characteristics of peers in Napster and Gnutella, such as their bottleneck bandwidth, latency, uptime, and number of shared files in 2001. They found that the number of shared files was heavily skewed.

Our study differs from the few previous studies on the shared files in P2P systems in at least two ways. First, we used recent and accurate snapshots of the Gnutella network with a significantly larger population of peers (i.e., more than a million concurrent peers). Second, we presented two properties of the shared files that have not previously been studied: (1) the correlation between shared files among peers and the overlay

topology structure, and (2) the variations in popularity of shared files across participating peers over time.

Another group of studies passively captured P2P traffic at a router to characterize exchanged files among peers. Gummadi et al. [29] analyzed a 200-day trace of Kazaa traffic collected at the University of Washington, demonstrating that file transfers in Kazaa do not follow a Zipf distribution and argued that this difference is due to the “fetch-at-most-once” nature of downloads in file-sharing applications. Another analysis of Kazaa traffic was conducted by Leibowitz et al. [25] at a large Israeli ISP. They examined the changing popularity of exchanged files among peers and showed that the data-sharing graph exhibits small-world properties. Note that the pattern of exchanged files among peers affects the characteristics of shared files throughout the system, but is subject to shorter-term trends. In contrast, the files shared by a peer may be the result of transfers over the course of months or years, followed by a gradual pruning of unwanted files. In summary, these studies on exchanged files are closely related and complement our work.

8 Conclusions and Future Work

This paper presented a measurement-based characterization of available files in the Gnutella file-sharing application. We discussed the challenges in capturing an accurate snapshot of available files in P2P file-sharing applications, and then presented two new measurement methodologies to achieve this goal. We used our parallel crawl and our peer sampler to obtain fairly accurate snapshots of available files among peers in the Gnutella network along with the connectivity among peers. Using these snapshots, we conducted three types of analysis and provided a better understanding of the distribution, correlation and dynamics of available files throughout the system.

We plan to continue this work in the following directions: We continue to collect many more snapshots to characterize properties of Gnutella files over longer timescales. Furthermore, we plan to further examine the models that properly capture each characteristics of available files in Gnutella. These models can be used in simulation based evaluations of file sharing applications. Finally, we are in the process of making our dataset publicly available to other researchers.

References

- Zhao, S., Stutzbach, D., Rejaie, R.: Characterizing files in the modern Gnutella network: a measurement study. In: *Multimedia Computing and Networking*, San Jose, CA (2006)
- slyck.com. <http://www.slyck.com>
- Saroiu, S., Gummadi, K.P., Dunn, R.J., Gribble, S.D., Levy, H.M.: An analysis of Internet content delivery systems. In: *OSDI*, Boston, MA, pp. 315–327 (2002)
- Stutzbach, D., Rejaie, R., Sen, S.: Characterizing unstructured overlay topologies in modern P2P file-sharing systems. In: *Internet Measurement Conference*, Berkeley, CA, pp. 49–62 (2005)
- Ripeanu, M., Foster, I., Iamnitchi, A.: Mapping the Gnutella network: properties of large-scale peer-to-peer systems and implications for system design. *IEEE Internet Comput. J.* **6**(1), 2002
- Klemm, A., Lindemann, C., Vernon, M., Waldhorst, O.P.: Characterizing the query behavior in peer-to-peer file sharing systems. In: *Internet Measurement Conference*, Taormina, Italy, pp. 55–67 (2004)
- Fessant, F.L., Handurukande, S., Kermarrec, A.-M., Massoulié, L.: Clustering in peer-to-peer file sharing workloads. In: *International Workshop on Peer-to-Peer Systems*, 2004
- Liang, J., Kumar, R., Xi, Y., Ross, K.W.: Pollution in P2P file sharing systems. In: *INFOCOM*, Miami, FL, pp. 1174–1185 (2005)
- Chu, J., Labonte, K., Levine, B.N.: Availability and locality measurements of peer-to-peer file systems. In: *ITCom: Scalability and Traffic Control in IP Networks II Conferences* (2002)
- Saroiu, S., Gummadi, P.K., Gribble, S.D.: Measuring and analyzing the characteristics of Napster and Gnutella hosts. *Multimedia Syst. J.* **9**, 170–184 (2003)
- Stutzbach, D., Rejaie, R.: Capturing accurate snapshots of the Gnutella network. In: *Global Internet Symposium*, Miami, FL, pp. 127–132 (2005)
- Stutzbach, D., Rejaie, R., Duffield, N., Sen, S., Willinger, W.: Sampling techniques for large, dynamic graphs. In: *Global Internet Symposium*, Barcelona, Spain (2006)
- Stutzbach, D., Rejaie, R., Duffield, N., Sen, S., Willinger, W.: On unbiased sampling for unstructured peer-to-peer networks. In: *Internet Measurement Conference*, Rio de Janeiro, Brazil, pp. 27–40 (2006)
- Rasti, A., Stutzbach, D., Rejaie, R.: On the long-term evolution of the two-tier Gnutella overlay. In: *Global Internet*, Barcelona, Spain (2006)
- Fisk, A.: Gnutella dynamic query protocol v0.1.” *Gnutella Developer’s Forum*, 2003
- Gnutella Developer Forum: Browse Host extension. http://www.the-gdf.org/wiki/index.php?title=Browse_Host_Extension
- Gnutella Developer Forum. <http://www.the-gdf.org>, 2005
- Stutzbach, D., Rejaie, R.: Evaluating the accuracy of captured snapshots by peer-to-peer crawlers. In: *Passive and Active Measurement Workshop*, Extended Abstract, Boston, MA, pp. 353–357 (2005)
- Chib, S., Greenberg, E.: Understanding the Metropolis–Hastings algorithm. *Am. Stat.* **49**:327–335 (1995)
- Hastings, W.: Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **57**, 97–109 (1970)
- Metropolis, N., Rosenbluth, A., Rosenbluth, M., Teller, A., Teller, E.: Equations of state calculations by fast computing machines. *J. Chem. Phys.* **21**, 1087–1092 (1953)
- Free Peers, Inc.: BearShare network statistics. <http://www.bearshare.com/stats>, October 2005
- Adar, E., Huberman, B.A.: Free riding on Gnutella First. *Monday* **5**, October 2000
- Stutzbach, D., Rejaie, R.: Understanding churn in peer-to-peer networks. In: *Internet Measurement Conference*, Rio de Janeiro, Brazil, pp. 189–202, 2006

25. Leibowitz, N., Ripeanu, M., Wierzbicki, A.: Deconstructing the Kazaa Network. In: *WIAPP 2003*
26. Bhagwan, R., Savage, S., Voelker, G.: Understanding availability. In: *International Workshop on Peer-to-Peer Systems*, (2003)
27. Stutzbach, D., Rejaie, R.: Characterizing the two-tier Gnutella topology. In: *SIGMETRICS, Extended Abstract*, Banff, AB, Canada (2005)
28. Adamic, L.A., Lukose, R.M., Huberman, B., Puniyani, A.R.: Search in power-law networks. *Phys. Rev. E* **64**(46135), (2001)
29. Gummadi, K.P., Dunn, R.J., Saroiu, S., Gribble, S.D., Levy, H.M., Zahorjan, J.: Measurement, modeling, and analysis of a peer-to-peer file-sharing workload. In: *SOSP*, Bolton Landing, NY, pp. 314–329 (2003)
30. Sen, S., Wang, J.: Analyzing peer-to-peer traffic across large networks. *IEEE/ACM Trans. Netw.* **12**, 219–232 (2004)