

Dissecting BitTorrent: Five Months in a Torrent’s Lifetime

M. Izal, G. Urvoy-Keller, E.W. Biersack, P.A. Felber, A. Al Hamra, and L. Garcés-Erice

Institut Eurecom, 2229, route des Crêtes, 06904 Sophia-Antipolis, France
{izal,urvoy,erbi,felber,alhamra,garces}@eurecom.fr

Abstract. Popular content such as software updates is requested by a large number of users. Traditionally, to satisfy a large number of requests, larger server farms or mirroring are used, both of which are expensive. An inexpensive alternative are peer-to-peer based replication systems, where users who retrieve the file, act simultaneously as clients and servers. In this paper, we study BitTorrent, a new and already very popular peer-to-peer application that allows distribution of very large contents to a large set of hosts. Our analysis of BitTorrent is based on measurements collected on a five months long period that involved thousands of peers. We assess the performance of the algorithms used in BitTorrent through several metrics. Our conclusions indicate that BitTorrent is a realistic and inexpensive alternative to the classical server-based content distribution.

1 Introduction

BitTorrent [4] is a file distribution system based on the peer-to-peer (P2P) paradigm. BitTorrent has quickly emerged as a viable and popular alternative to file mirroring for the distribution of large content, as testified by the numerous Web sites that host active “torrents” (e.g., <http://f.scarywater.net/>).

We have conducted a comprehensive analysis of BitTorrent to assess its performance. To that end, we have used two sources of information. First, we have obtained the “tracker” log of arguably the most popular torrent (BitTorrent session) so far—the 1.77GB Linux Redhat 9 distribution—for its first 5 months of activity. The log contains statistics for more than 180,000 clients, and most interestingly, it clearly exhibits an initial flash-crowd period with more than 50,000 clients initiating a download in the first five days. This first source of information allows us to estimate the global efficiency of BitTorrent, the macroscopic behavior of clients, and the scalability of a P2P application under flash-crowd conditions. Our second source of information consists of data collected with a modified client that participated to the same torrent downloading Redhat 9. This second log allows us to study the direct interactions between the clients. The remaining of the paper is organized as follows. In Section 2, we present the main features of BitTorrent. In Section 3, we review the related work. In Section 4, we present the results obtained from the tacker log and in Section 5 the

conclusions obtained from our client log. We conclude with future directions in Section 6.

2 BitTorrent

BitTorrent is a P2P application that capitalizes the *resources* (access bandwidth and disk storage) of peer nodes to efficiently distribute large contents. There is a separate torrent for each file that is distributed. Unlike other well-known P2P applications, such as Gnutella or Kazaa, which strive to quickly *locate* hosts that hold a given file, the sole objective of BitTorrent is to quickly *replicate* a single large file to a set of clients. The challenge is thus to maximize the speed of replication.

A torrent consists of a central component, called *tracker* and all the currently active peers. BitTorrent distinguishes between two kinds of peers depending on their download status: clients that have already a complete copy of the file and continue to serve other peers are called *seeds*; clients that are still downloading the file are called *leechers*. The tracker is the only centralized component of the system. The tracker is not involved in the actual distribution of the file; instead, it keeps meta-information about the peers that are currently active and acts as a rendez-vous point for all the clients of the torrent.

A user joins an existing torrent by downloading a *torrent* file (usually from a Web server), which contains the IP address of the tracker. To initiate a new torrent, one thus needs at least a Web server that allows to discover the tracker and an initial seed with a complete copy of the file. To update the tracker's global view of the system, active clients periodically (every 30 minutes) report their state to the tracker or when joining or leaving the torrent. Upon joining the torrent, a new client receives from the tracker a list of active peers to connect to. Typically, the tracker provides 50 peers chosen at random among active peers while the client seeks to maintain connections to 20–40 peers. If ever a client fails to maintain at least 20 connections, it recontacts the tracker to obtain additional peers. The set of peers to which a client is connected is called its *peer set*.

The clients involved in a torrent cooperate to replicate the file among each other using *swarming* techniques: the file is broken into equal size chunks (typically 256kB each) and the clients in a peer set exchange chunks with one another. The swarming technique allows the implementation of parallel download [7] where different chunks are simultaneously downloaded from different clients. Each time a client obtains a new chunk, it informs all the peers it is connected with. Interactions between clients are primarily guided by two principles. First, a peer preferentially sends data to peers that reciprocally sent data to him. This “tit-for-tat” strategy is used to encourage cooperation and ban “free-riding” [1]. Second, a peer limits the number of peers being served simultaneously to 4 peers and continuously looks for the 4 best downloaders (in terms of the rate achieved) if it is a seed or the 4 best uploaders if it is a leecher.

BitTorrent implements these two principles, using a “choke/unchoke” policy. “Choking” is a temporary refusal to upload to a peer. However, the connection

is not closed and the other party might still upload data. A leecher services the 4 best uploaders and chokes the other peers. Every 10 seconds, a peer re-evaluates the upload rates for all the peers that transfer data to him. There might be more than 4 peers uploading to him since first, choking is not necessarily reciprocal and second, peers are not synchronized.¹ He then chokes the peer, among the current top 4, with the smallest upload rate if another peer offered a better upload rate. Also, every 3 rounds, that is every 30 seconds, a peer performs an *optimistic unchoke*, and unchokes a peer regardless of the upload rate offered. This allows to discover peers that might offer a better service (upload rate). Seeds essentially apply the same strategy, but based solely on download rates. Thus, seeds always serve the peers to which the download rate is highest.

Another important feature of BitTorrent is the chunk selection algorithm. The main objective is to consistently maximize the entropy of each chunk in the torrent. The heuristic used to achieve this goal is that a peer always seeks to upload the chunk the least duplicated among the chunks it needs in its peer set (keep in mind that peers only have a local view of the torrent). This policy is called the *rarest first policy*. There exists an exception to the rarest first policy when a peer joins a torrent and has no chunks. Since this peer needs to quickly obtain a first chunk (through optimistic unchoke), it should not ask for the rarest chunk because few peers hold this chunk. Instead, a newcomer uses a random first policy for the first chunk and then turns to the rarest first policy for the next ones.

3 Previous Work

Approaches to replicate contents to a large set of clients can be classified as client side and server side approaches. The first client-side approach was to cache contents already downloaded by clients of a given network. A symmetric approach on the server side is to transparently redirect clients to a set of mirror sites run by a content provider (e.g. Akamai).

The peer-to-peer paradigm has been applied to obtain client side and server side solutions. On the server side, one finds proposals like [5] where overlay nodes dynamically organize themselves so as to form a tree with maximum throughput. FastReplica [2] is designed to efficiently replicate large contents to a set of well-known and stable overlay nodes.

On the client side, one finds many proposals to build application-layer multicast services [8, 6, 3]. A solution similar to BitTorrent is Slurpie [9]. Slurpie aims at enforcing cooperation among clients to alleviate the load of a Web server that is the primary source of the content. The algorithms of Slurpie are much more complex than the ones of BitTorrent and require to estimate the number of peers in the Slurpie network. The expected improvements over BitTorrent is less load on the topology server (equivalent to the BitTorrent tracker) and on the primary

¹ For instance, a peer that was offering a very good upload rate has decided to choke this connection just 1 second before the reevaluation on the other side and thus he might still be in the top 4 list for the next 10 seconds and received service.

source (original seed) through a back-off algorithm. Results are promising since Slurpie is able to outperform BitTorrent in a controlled environment. Still, the actual performance of Slurpie in case of flash crowds and for a large number of clients is unknown.

4 Tracker Log Analysis

The tracker log covers a period of 5 months from April to August 2003. The corresponding torrent has as content the 1.77 GB Linux Redhat 9 distribution. 180,000 clients participated to this torrent with a peak of 51,000 clients during the first five days (see Figures 4 and 4). These first five days clearly exhibits a flash-crowd. As clients periodically report to the tracker their current state, along with the amount of bytes they have uploaded and downloaded, the tracker log allows us to observe the global evolution of the file replication process among peers.

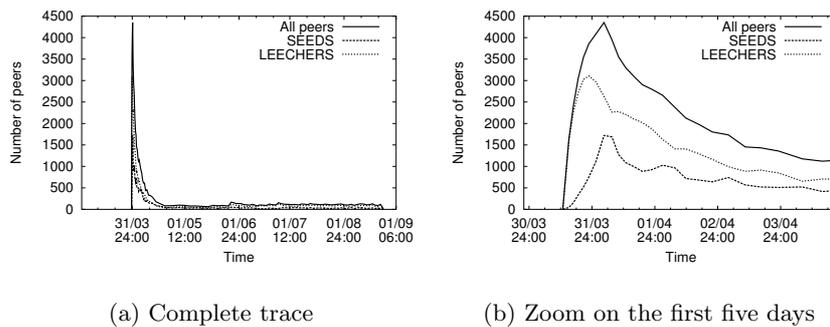


Fig. 1. Number of active peers over time

4.1 Global Performance

Analyzing the tracker log, our first finding is that BitTorrent clients are altruistic in the sense that they actively send data to other clients, both as leechers and as seeds. Altruism is enforced during the download phase by the tit-for-tat policy, as a selfish client will be served with a very low priority. Once they become seed, the peers remain connected for another six and a half hours on average. This “social” behavior can be explained by two factors: first, the client must be explicitly terminated after completion of the download, which might well happen while the user is not at his computer, e.g., overnight; second, as the

content being replicated is perfectly legal, the user has no particular incentive to quickly disconnect from the torrent. In fact, the presence of seeds is a key feature, since it greatly enhances the upload capacity of this torrent and the ability to scale to large client populations. Over the 5 months period covered by the log file, we observed that the seeds have contributed more than twice the amount of data sent by leechers (see Figure 2). We also observed that the proportion of seeds is consistently higher than 20%, with a peak at 40% during the first 5 days (see Figure 3). This last figure clearly illustrates that BitTorrent can sustain a high flash-crowd since it quickly creates new seeds. To put it differently, in situations where new peers arrive at a high rate, the resources of the system are not divided evenly between clients, which would result, like in a processor sharing queue under overload, in no peers completing the download. On the contrary, older peers have a higher priority since they hold more chunks than younger peers, which gives them more chance to complete the download and become seeds for newcomers. Obviously, this strategy benefits from the cooperation of users that let their clients stay as seeds for long periods of time.

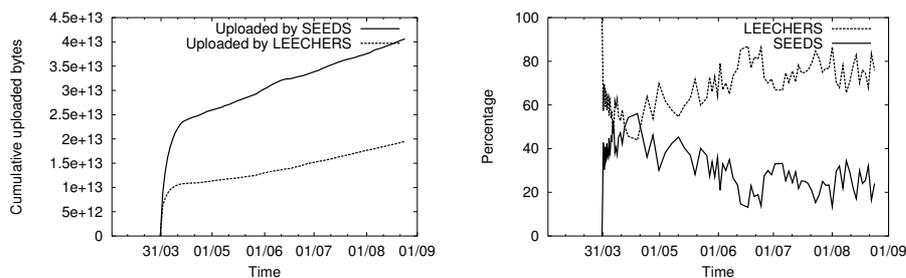


Fig. 2. Volumes uploaded by seeds and **Fig. 3.** Proportions of seeds and leechers

4.2 Individual Session Performance

A fundamental performance metric for BitTorrent is the average download rate of leechers. Over the 5 months period covered by the tracker log, we observed an average download rate consistently above 500kb/s. This impressive figure indicates that most of the BitTorrent clients have good connectivity (at least ADSL), which is not surprising given the total size of the file. Moreover, BitTorrent exhibits good scalability: during the initial flash-crowd, the average download rate was close to 600kb/s and the aggregate download rate of all simultaneously active clients was more than 800Mb/s.

BitTorrent is thus able to capitalize the bandwidth of end system hosts to achieve a very high aggregate throughput. Still, the final objective of BitTorrent is to replicate a given content over all the peers in a torrent. We thus need to know

how many peers eventually completed the download and the fraction of bytes that the peers that did not complete the transfer downloaded. Since BitTorrent allows users to suspend and resume the download at any time, a peer might complete the download after one or several sessions. In the tracker log, a session id identifies the session (hash of the IP address of the host and its current time), along with the IP address of the host as seen by the tracker. Thus, identifying single session downloads is easy based on the session id. To reconstruct multi-sessions downloads, we assumed that the IP address of the host does not change from one session to another. NATs often used port numbers to disambiguate connexions and not IP addresses, so our assumption can hold even in this case. Of course, if two or more hosts behind a firewall simultaneously participate to the same torrent, it might be difficult to disambiguate them. In addition, since a peer provides the amount of bytes it has downloaded in each of its reports (sent every 30 minutes), we mated two sessions with the same IP if the amount of bytes of the last report for the first session is close to the amount of bytes of the first report of the second session. These two values are not necessarily the same if for instance, the user disconnected the BitTorrent client improperly which results in the latter not sending its disconnection report with its current amount of bytes downloaded. A set of sessions form a multi-sessions download if the download of the file is completed during the last session (on average, we found that a multi-session download consists of 3.9 sessions). We also observed some sessions that started with already 100% of the file downloaded. This is the case if the user is kind enough to rejoin the torrent after the completion of the download to act as a seed. We thus categorize the sessions into four categories : single session download, multi-session download, seed sessions and session that never complete (incomplete downloads). We provide statistics for these sessions in Table 1.

Type	Number of sessions	Total down. (TB)	Total up. (TB)	Down. /session (MB)	Up. /session (MB)	Down. rate (kb/s)	Up. rate (kb/s)	Down. Duration (hours)
Single session downloads	20584	34.7	28.5	1765.7	1450.2	1331.1	325.2	8.1
Multi-session downloads	14285	6.2	4.3	455.2	319.2	390.6	145.0	19.2
Incomplete downloads	138423	11.2	9.1	84.5	69.1	114.7	50.5	-
Seed sessions	8555	-	12.7	-	1556.6	-	223.6	-
Total	181847	52.1	54.6	-	-	-	-	-

Table 1. Sessions types and their characteristics

Overall, Table 1 indicates that only 19% of the sessions are part of a transfer that eventually completed. However, the 81% of sessions that did not complete the transfer represent only 21.5% of the total amount of downloaded bytes. Moreover, the incomplete sessions uploaded almost as much bytes as they downloaded. It is difficult to assess the reason behind the abortion of a transfer. It

might be because the user eventually decides he is not interested in the file or because of a poor download performance. To investigate this issue, we look at the statistics (durations and volumes) for sessions that start with 0% of the file during the first five days. We classify these sessions into two categories denoted as “completed” and “non-completed”. The completed category corresponds to the single session downloads defined previously while the non-completed category corresponds to transfers that never complete or the first session of transfers that will eventually complete. We present on Figure 4 the complementary cumulative distribution functions for the download volumes and session durations for the completed and non-completed sessions. We can observe from this figure that 90% of the non-completed sessions remain less than 10,000 seconds in the torrent (and 60% less than 1000 seconds) and retrieve less than 10% of the file. Most of the non-completed sessions are thus short and little data is downloaded.

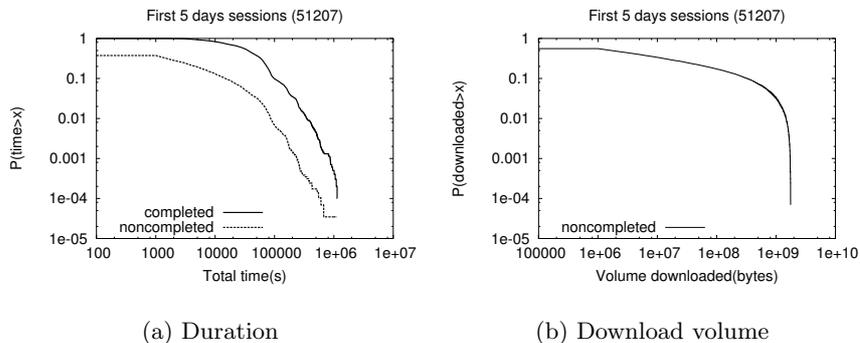


Fig. 4. Single sessions statistics for sessions seen during the first 5 days

Let us now further focus on the performance achieved by single session downloads. The average download rate of these 20,584 sessions is close to 1.3Mb/s over the whole trace, which is larger than the download rate averaged over all sessions (500kb/s). The average download time for single session download is about 29,000 seconds overall. Such values reveal a high variability in the download rates achieved by these sessions since if they had all achieved a download rate of 1.3Mb/s, the average download time would be $\frac{1.77\text{GB}}{1.3\text{Mb/s}} \sim 10,000$ seconds. This is confirmed by the distribution of these download throughputs (see Figure 4.2) that clearly exhibits a peak around 400kb/s (a typical value for ADSL lines), a value significantly smaller than the mean.

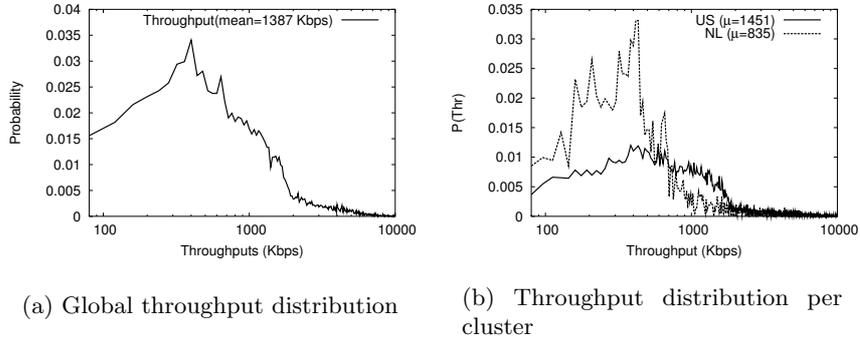


Fig. 5. Throughputs for single session downloads

4.3 Geographical Analysis

The tracker log provides the IP addresses of the clients. Based on these addresses, we have used NetGeo (<http://www.caida.org/tools/utilities/netgeo/>) to obtain an estimation of the origin country of the peers that participated to the torrent. The estimation might be imprecise as NetGeo is not maintained any more. Overall, we were not able to obtain information for around 10% of the hosts. In Table 2 we indicate the top five countries for the first 5 days, the first 4 weeks and the complete trace. We can observe that this set of countries, as well as the ranking is consistently the same for all three time scales. Most of the clients are US clients while Europe is represented only through the Netherlands. Still, for Netherlands, 50% of the hosts originate from ripe.net (an information also provided by NetGeo), which acts as the Regional Internet Registry for the whole of Europe, Middle East, Africa and part of Asia also. Thus we can guess that most of these peers are spread all over Europe.

Countries	First week	First four weeks	Complete trace
United States	44.6%	44.3%	32%
Netherlands (Europe)	14.9%	15.3%	23.9%
Australia	12.7%	12.6%	17.8%
Canada	5.7%	5.8%	4.9%
% of total hosts	77.9%	78%	78.6%

Table 2. Geographical origins of peers

We next investigate the relative performance of the four clusters identified above (US, NL, AU and CN). To do so, we consider again the 20,584 peers

that complete the download in a single session. 17,000 peers out of the 20,584 initial peers, are in the four clusters (11498 in US cluster, 2114 in NL, 1995 in AU and 1491 in CM). We plot in Figure 4.2 the distribution of download throughputs achieved by the peers in the NL and US clusters (the AU and CN clusters are highly similar to the US cluster and not depicted for sake of clarity). Figure 4.2 reveals that the download throughput of the hosts in the NL cluster is significantly smaller than the throughput of the hosts in the US cluster (where more mass is on the right side of the curve). This can indicate that clients in the US have, in general, better access links than in Europe.

5 Client Log Analysis

To better observe the individual behavior of a BitTorrent peer, we have run an instrumented client behind a 10Mb/s campus network access link. Our client joined the torrent approximatively in the middle of the 5 months period (i.e., far after the initial flash crowd). We experienced a transfer time of approximately 4,500 seconds (i.e., much lower than the average download times previously mentioned) and our client remained connected as a seed for 13 hours. Our client logged detailed information about the upload or download of each individual chunk. In Figure 6, we represent the number of peers with whom our client is trading. At time 0, our client knows around 40 peers whose addresses have been provided by the tracker. We next continuously discover new peers (peers that contacted us) up to the end of the download (at time 4,500 seconds) where we observe a sudden decrease of the number of peers, most probably because the seeds we were connected to have closed their connection to us as soon as we have completed the download. After the download, we stay connected as seed to between 80 and 95 leechers (4 of them being served while the others are choked).

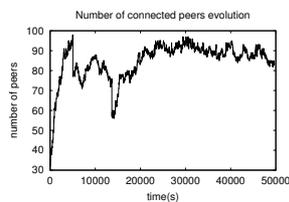


Fig. 6. Number of peers during and after download

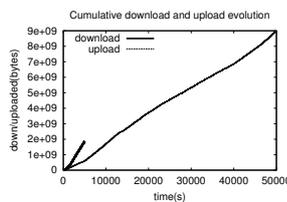


Fig. 7. Complete torrent

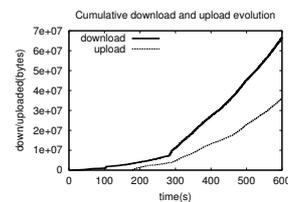


Fig. 8. First 10 minutes of download

Figures 7 and 8 show the amount of bytes downloaded and uploaded with respect to time for the complete trace and the first 10 minutes of the trace respectively. From these figures, we can draw several conclusions:

- There is a warm-up period (around 100 seconds) to obtain some first chunks. But as soon as the client has obtained a few chunks, it is able to start

uploading to others peers. This means that the rarest first policy works well, otherwise our client would maybe not find anyone interested in its first chunks.

- The download and upload rates are (positively) correlated, which indicates that the tit-for-tat policy works. It also indicates that we always find peers interested in our chunks and peers from which we can download chunks. Otherwise, we would observe some periods where the curves are flat, indicating that our client is stalled. Also, the way peer sets are built with “old” and “new” peers mixed together must have a positive impact on the efficiency of BitTorrent. Indeed, a torrent can be viewed as a collection of interconnected sets of peers. A peer that joins a torrent will obviously be the youngest peer in its initial set, but it may later be contacted by younger peers. It may also be contacted by older peers since peers try to keep contact with a minimum number of peers (generally 20) throughout their lifetimes in the torrent. This temporal diversity is a key component of the efficiency of BitTorrent, since it guarantees with high probability that a given peer will find other peers (younger or older) that hold some missing chunks of the file.
- It takes twice as much the download period to upload the same quantity of bytes (1.77GB) to the system. This illustrates the importance of peers staying as seeds once they have completed the download. This means also that we downloaded at a faster rate than we uploaded. This is due to the fact that we have a high speed link and since seeds always seek for the fastest downloaders, we should be consistently favored by the seeds that serve us.

We further investigated the type of clients we have been trading with. Overall, we found that approximately 40% of the file was provided by seeds and 60% by leechers. We also observed that more than 85% of the total file was sent by only 20 peers, including the 8 seeds that provided 40% of the file. An interesting remark is that these 20 top uploaders were not in our initial peer set provided by the tracker, but contacted us later. We can thus conjecture that to obtain the best possible performance with BitTorrent, clients should not be behind firewalls or NATs that prevent inbound connections.

We also want to assess the efficiency of the tit-for-tat policy. A good tit-for-tat policy should enforce clients to exchange chunks with one another, but with enough flexibility so as not to artificially block transfers when data does not flow at the same rate in both directions. BitTorrent avoids this type of problem by chocking/unchoking connections every 10 seconds, a long period at the time scale of a single TCP connection. We have studied the correlations between upload and download throughput, as well as between upload and download traffic volumes. Results show that, while traffic *volumes* are correlated (positive correlation close to 0.5), the upload and download *throughputs* are not correlated (close to 0), which means that BitTorrent is flexible. We also computed the correlations between download volumes and download throughputs on one side and upload volumes and upload throughputs on the other side. Both are correlated with value 0.9 and 0.5, respectively. The high correlation observed between downloaded throughputs and volumes is probably due to the fact that the top

three downloaders are seeds that provide 29% of the file. The upload throughput and volume are also correlated because, once our client becomes a seed, it continuously seeks for the best downloaders, which is not the case during the download phase where a peer primarily seeks for peers that have some chunks he is interested in.

6 Conclusion

Large content replication has become a key issue in the Internet; e.g. big companies are interested in a replication service to update simultaneously a large set of hosts (e.g., virus patches or software updates). BitTorrent is a very popular peer-to-peer application targeted at large file replication. We have extensively analyzed a large torrent (with up to a few thousands simultaneously active clients) over a large period of time. The performance achieved, in terms of the throughput per client during download and the ability to sustain high flash-crowd, demonstrate that BitTorrent is highly effective.

There still remain open questions in the area of large content replication such as: (i) what is the optimal (at least theoretically) replication policy, in terms of response times, for a given user profile (arrival rates, willingness to stay as seeds), (ii) how to build a robust replication service, i.e. to ensure that all machines eventually complete the downloads (a must for corporate usage), (iii) how to efficiently protect these applications against denial of service and malicious peers.

Acknowledgment

We are extremely grateful to Eike Frost for kindly providing us with the tracker log of the torrent analyzed in this paper.

References

1. E. Adar and B. A. Huberman, "Free Riding on Gnutella", *First Monday*, 5(10), October 2000.
2. L. Cherkasova and J. Lee, "FastReplica: Efficient Large File Distribution within Content Delivery Networks (USITS 2003)", In *Proceedings of the 4th USENIX Symposium on Internet Technologies and Systems*, March 2003.
3. Y.-H. Chu, S. G. Rao, and H. Zhang, "A case for end system multicast", In *ACM SIGMETRICS 2000*, pp. 1–12, Santa Clara, CA, USA, June 2000.
4. B. Cohen, "Incentives to Build Robustness in BitTorrent", <http://bitconjurer.org/BitTorrent/bittorrentecon.pdf>, May 2003.
5. J. Jannotti, D. K. Gifford, and K. L. Johnson, "Overcast: Reliable Multicasting with an Overlay Network", In *Proc. 4-th Symp. on Operating Systems Design and Implementation*, Usenix, October 2000.
6. S. Ratnasamy, M. Handley, R. M. Karp, and S. Shenker, "Application-Level Multicast Using Content-Addressable Networks", In *NGC 2001*, pp. 14–29, November 2001.

7. P. Rodriguez and E. W. Biersack, "Dynamic Parallel-Access to Replicated Content in the Internet", *IEEE/ACM Transactions on Networking*, 10(4):455–464, August 2002.
8. A. Rowstron et al., "SCRIBE: The design of a large scale event notification infrastructure", In *Proc. NGC 2001*, November 2001.
9. R. Sherwood, R. Braud, and B. Bhattacharjee, "Slurpie: A Cooperative Bulk Data Transfer Protocol", In *Proceedings of IEEE INFOCOM*, March 2004.