

# A Content-based View of Real-Life BitTorrent Swarms

**Gábor Simon, Róbert Dávid**

Department of Automation and Applied Informatics, Budapest University of  
Technology and Economics  
Magyar tudósok körútja 2, H-1117 Budapest, Hungary  
{simon.gabor, davidrobi}@aut.bme.hu

*Abstract: BitTorrent network offers, as well as other peer-to-peer (P2P) file-sharing solutions, a way to let simple Internet edge nodes (peers) to share data with other peers. Distribution groups (swarms) of this network are enhanced to make data content of a file group available to all joined peers in the swarm. However, the primary needs of users are about acquiring a specific content. This content can be represented in very different forms within the digital world as well as within the file-sharing networks. In this paper we present the impact of the semantic gap between the content-based demands and the file-based distribution solutions through measurement results from a real-life BitTorrent community. Besides this we point out the typical structures built by different manifests of the same content in BitTorrent network. As a result we will be able to propose some guidelines to enhance the content-awareness of tomorrow's P2P solutions.*

*Keywords: BitTorrent; Swarms; Content distribution; Metadata; File-sharing*

## 1 Introduction

### 1.1 Bittorrent

Although there are differences among Internet traffic measurement methodologies, all major research organizations agree that BitTorrent gained a significant share in Internet traffic during the last few years [1][2][3][4]. The cause can be found in the design of the BitTorrent protocol, which offers efficient distribution using dynamic load balancing and tit-for-tat policies. The overall performance of transfer is increasing, when a downloader (*leecher*) joins the downloading group (*swarm*). Members of a swarm require the same piece of data specified by the so-called torrent file. Torrent files hold information about one or more files, aimed to be distributed together. When a leecher finishes downloading all pieces (*chunks*) of each and every file specified in torrent, it becomes a seeder:

a peer storing content as well as uploading it to other leechers. Performance improvements are achieved by balancing algorithms controlling piece-transfer. The algorithms choose the next downloadable piece by considering the overall availability of that piece in the swarm. The selection of the download source is controlled by the choke/unchoke policy considering the upload/download rate and the number of active (download-upload) connections to avoid freeriding. The communication in the swarm is controlled by a special node called tracker. Trackers are also supervising the join process, performing authentication, relaying IP addresses and keeping track of joined peers.

## 1.2 Basic Categories

We are going to observe these swarms from a content-based view. Before that we should clarify the meaning of the category content and its relation to files. *Content* is a piece of information, result of a creative effort. Content can be stored in different ways - in the field of information technology, these are called *releases*. One release can consist of one or more files, or can be a part of a large archive file. Contents can be stored in more than one release.

Swarms do not maintain any additional metadata about the distributed contents, although the files are definitely parts of one or more contents and the users have demands about contents. However, swarms are only the core elements of the BitTorrent infrastructure: there are BitTorrent search sites and index sites where users can choose a suitable swarm to join to. These sites offer similar metadata to support the decision: torrent-file name, names and sizes of files distributes in the swarm and some statistics about swarms, but they do not maintain any content information and content-release connections.

In the following sections, we present measurement results indicating the consequence of this shortcoming and reveal some typical phenomena about content distribution in BitTorrent network.

## 2 Backgrounds and Related Work

In one of our previous studies [5] we have constructed a simple category system to make simpler the content-based comparison between different networks. Although we discuss here only the BitTorrent network, we are going to use these categories in the followings. In that paper we have also extended the scope of some file-based metrics and we have introduced content-based metrics.

Former measurement studies on BitTorrent network [6][7][8][9][10][11] focused on torrent-related issues, and revealed some significant phenomena about file distribution. In [8] tracker logs were analyzed and the so-called flashcrowd effect

for a single-file was shown. In [7] a real BitTorrent community was monitored and several issues were addressed, such as integrity, availability, download performance and flashcrowds. Some other studies have examined the distribution performance of BitTorrent and the traffic share of P2P applications. In [12] six different torrent communities were monitored and several behavioral characteristics were revealed regarding freeriding and share ratio. In [13] a fluid model was introduced to estimate the download time of a single file, however the results is contradicted with the results in [7]. This indicates the importance of accurate measurement studies to validate P2P models.

Recently, Go et al. [11] published an extensive measurement study about performance of BitTorrent-like systems. They discovered several limitations in current design of BitTorrent network. First of all, the peer arrival rates are exponentially decreasing, thus, torrents are facing relatively fast extinction after flashcrowd phase. Secondly, download performance in the network is fluctuating significantly, along with the number of online peers. Finally, due to the instantly rewarding tit-for-tat policies the network provide unfair services to the peers and they do not encourage high-speed seeders to stay in the swarm. The authors propose an inter-torrent collaboration strategy that uses the trivial relation between torrents: two torrents are related when they have had common downloaders. Based on this relation, a new kind of tit-for-tat policy can be implemented, which is aware of the group of torrents downloading or seeding by the same peer and encourages the long-time seeding. In [14], our working group also emphasized the importance of inter-torrent (more likely *inter-swarm*) relations, but we considered content-based relations, because of the content-related needs of users. However, our approach requires content-based metadata support, which makes it a future solution, considering the low number of content-aware real-life file-sharing communities.

### 3 Content-Aware BitTorrent Communities

As it mentioned before, most BitTorrent sites does not maintain content-related metadata, only some extracted property of the torrent along with its tracker statistics. The underlying cause is that the whole BitTorrent infrastructure does not maintain this missing piece of information either. As the communities have realized this shortcoming, users have begun to utilize the community tools provided by BitTorrent sites to build up connections between different torrents distributing the same or related content. The most common form of such cross-referencing was to place links in the forum of torrents to refer other torrents with the same content but with better download performance or media quality as well as to torrents with related content. Most communities are improving their sites to efficiently utilize this community-provided information by swarm selection.

### 3.1 Isohunt

Isohunt [15] is one of the major torrent sites indexing more than one million torrents. The site recently introduced a novel feature: a catalogue of releases. The catalogue currently contains more than 45,000 items. An item can be categorized and several metadata entries can be filled, such as general description, links for identifying the content and additional links to identify the release in different file-sharing networks. Although this approach is quite innovative, it does not improve the quality of the search performance due to the lack of conventions or schemas regarding metadata entries and the low number of filled metadata fields.

### 3.2 Mininova

Mininova [16] is another significant torrent site with over 650,000 indexed torrents. This site offers a simple way to connect torrents distributing the same movie content. A special link can be assigned to each torrent referencing to an entry in the largest movie database (IMDB [17]). As a result, torrents with the same reference are distributing the releases of the same content.

In our measurement process we extracted these semantic link information from the Mininova dataset in order to mark the content-wise related swarm groups.

## 4 Measurement Results

### 4.1 Basic Measurement Statistics

Our sample contained over 17000 torrents (swarms), connecting about 860000 peers. More basic statistics about our sample can be found in Table 1.

Table 1  
Basic measurement statistics

No. of torrents involved	17583
No. of contents involved	5318
Involved torrents – all movie torrents percentage	9.87 %
Involved torrents – all indexed torrents percentage	2.52 %
Average no. of seeders	23.3
Average no. of leechers	25.6
Average no. of swarm members	48.9
Average seeder-leecher ratio	1.05
Average no. of swarms in meta-swarm	3

## 4.2 Seeder-Leecher Ratio

The one of the most commonly used indicators of swarm performance is the seeder-leecher ratio. Our first step was to determine if content-based view of the swarms improved this metric. Figure 1 shows the number of seeders and leechers in each torrent swarm in a log-log coordinate system. The linear regression for seeders is

$$y = 0,7218x + 5,8031 \quad (1)$$

where  $x$  is the number of leechers, so on average, there are 72 seeders for 100 leechers.

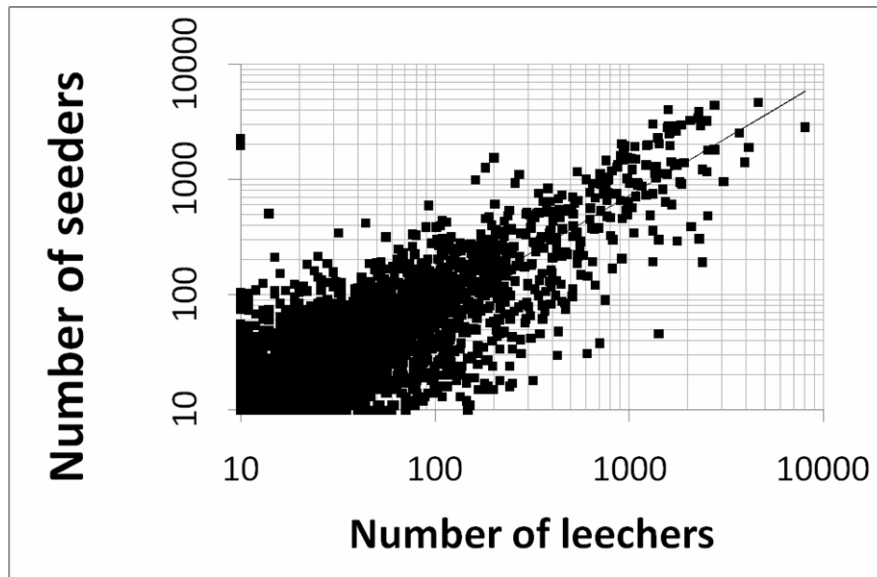


Figure 1  
Number of seeders and leechers in torrents (log-log coordinates)

Figure 2 shows the same for metaswarms. In this case, the regression is

$$y = 0,9539x - 3,4426 \quad (2)$$

This suggests that there are 95 seeders for 100 leechers of the same content.

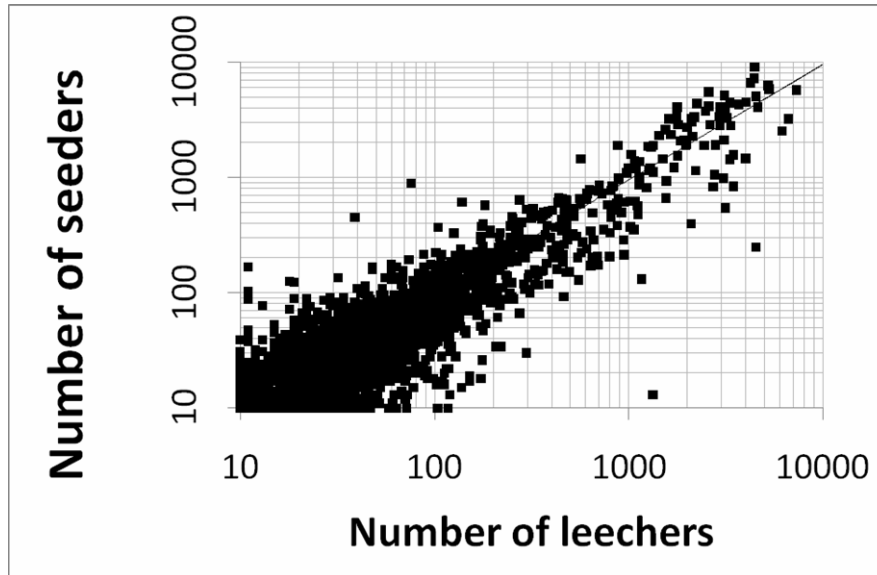


Figure 2  
Number of seeders and leechers in metaswarms (log-log coordinates)

Note that both number is worse then the average seeder:leecher ratio of the torrents. This is because torrents with few peers and small seeder-leecher ratio are rare, people tend to avoid such torrents, and torrents with only a few number of peers may have very high seeder-leecher ratios, and that distorts the average.

### 4.3 Metaswarm Structure

Critical issue of the bittorrent network is that usually for each metaswarm there is only one swarm that gathers most of the peers for that content, other swarms only have few peers. To show this, we have sorted torrents by number of peers and calculated their average contribution to their metaswarm. This can be seen on Figures 3 and 4. On average, more then 80% of all peers is connected to the largest swarm, and the 5th largest swarm is less then 5% as big as the metaswarm, however, torrents smaller then the first five largest can still hold up to 20% of all peers, and these peers will see an overall worse download performance. It is desired to guide users into the larger swarms of the same content.

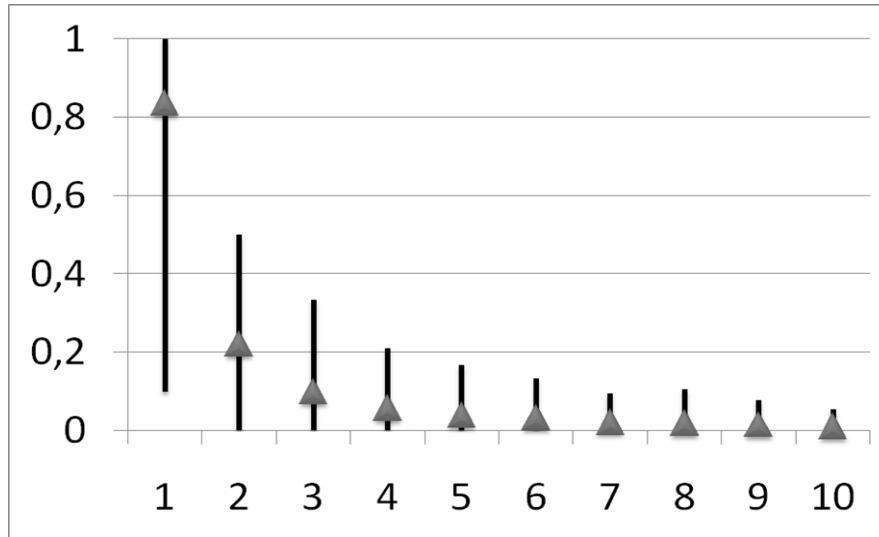


Figure 3

Torrents sorted by number of peers descending, showing their minimum, maximum and average percentage of peers in the metaworm

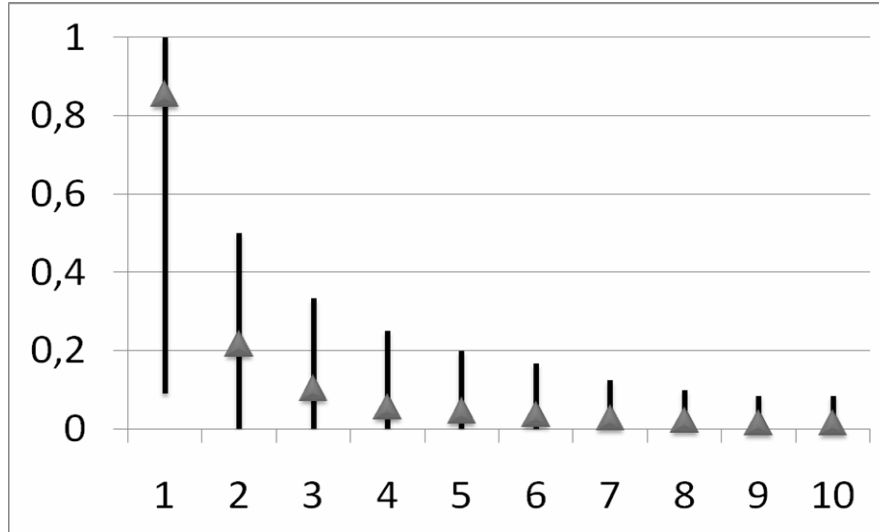


Figure 4

Torrents sorted by number of seeders descending, showing their minimum, maximum and average percentage of seeders in the metaworm

Bittorrent's current mechanism for that is the following: New downloaders will join the larger swarms, so torrents with few peers will die out when their seeders

are quitting. We examined this process: We compared the size of the swarms against the age of the torrents, as illustrated in Figure 5. The five lines on the diagram means the average ratio of torrents smaller then 20%-40%-60%-80%-99% of the largest torrent of the metaswarm. The result is quite surprising: Over 30 months, the ratio of small torrents decreased only around 10%, so the mechanism is clearly not effective, bittorrent needs a better solution for that problem.

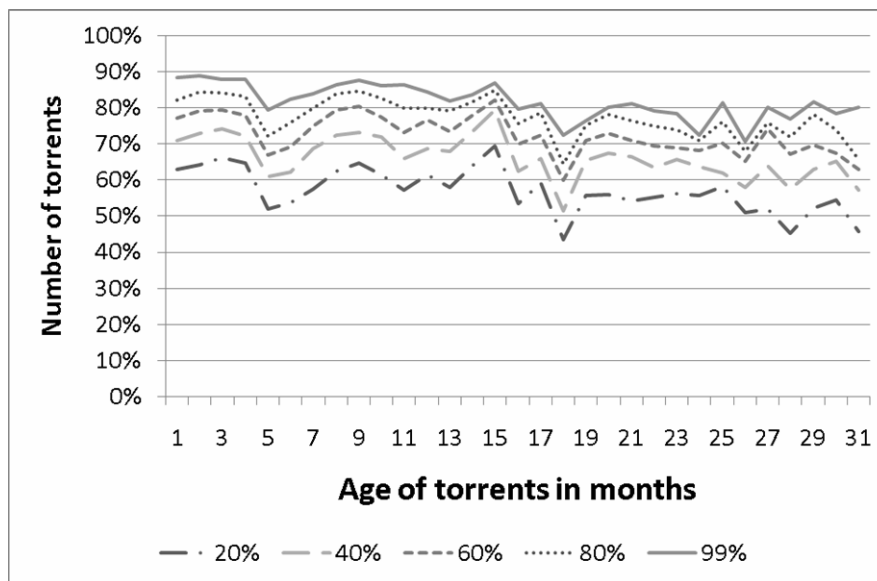


Figure 5  
Ratio of torrents of a metaswarm compared against the age of the torrents

## 5 Limitations and Future Works

There are some issues regarding the design of Mininova's IMDB reference system. First of all, it works only for movie torrents. However, similar feature could be implemented easily involving different metadata stores for different media types (e.g. MusicBrainz [18] for music). Secondly, there is no validation for reference entries, thus, invalid values (e.g. IMDB rating instead of IMDB id) are also accepted. Furthermore, the reference entry is optional. Finally, the relation between a torrent and a content can be more complex than it could be expressed with one-to-one mapping (a single hyperlink in MiniNova). A torrent can distribute more than one release as well as only a part of a release. Issues



regarding the relativistic nature of contents and releases are discussed in details in [14].

These implementation limitations cause some further limitations in the measurement strategy as well. As a start, we only considered torrents with IMDB reference entry filled. This is a special set of torrents, since they generally distribute only one release of a single movie content. As a result, either torrents without IMDB reference or swarms distributing multiple contents are not involved. Furthermore, during the aggregation step we assumed that the peer sets of the same content are distinct and can be aggregated using set union operation, although it is not likely for a peer to distribute more than one release of a content, thus, we do not bring significant distortion. However, there are other issues about this kind of aggregation, because we assumed that the merge of the peer sets meets the user preferences. Migration of users, who have chosen their release for its specific, optimized property (low-bandwidth, low-size), will likely fail. An ideal swarm structure that would not only suit the user needs, but also maximizes the utilization of peer resources is described in [14].

Currently we follow two different research directions in our working group. As we are aware of that the most of the formerly mentioned limitations are derived from the moderate way of implementation of the relation management system, we are going to design a more extensive content management layer for P2P application in order to gain more accurate measurement results. On the other hand, we are currently gathering data snapshots from Mininova site periodically in order to observe the dynamic behavior of contents and meta-swarms.

### **Conclusion**

Although the content-release management system of Mininova is only a simple example of building inter-swarm relations, it offers a unique way to get a content-based view of Bittorrent swarms. To take this chance we extracted the swarm statistics along the IMDB references from the Mininova site in order to determine swarms distributing the same content. The results indicated that the most of the peers of an ordinary meta-swarm are joined to one of the few major swarms. Remaining swarms have only a few peers each. The current meta-swarm structure carries a significant performance leakage as the resources of the peers in minor swarms are unavailable for the most of the peers downloading the same content. Unfortunately, since the Bittorrent has no strategy to shape meta-swarm structure, the share of minor swarms does not decrease over time.

We hope that these results will contribute to the better understanding of content distribution in Bittorrent network and encourage the use of content-based aspect in further research activities.

### **References**

- [1] Ipoque, Internet Study 2007: Data about P2P, VoIP, Skype, File Hosters like RapidShare and Streaming Services like YouTube.

- [2] CacheLogic; <http://www.cachelogic.com/>.
- [3] Sprint, Sprint Academic Research Group: Packet Trace Analysis.
- [4] Ellacoya Networks; <http://www.ellacoya.com>.
- [5] R. David and G. Simon, A Content-Based View of Peer-to-Peer File-Sharing Networks, MicroCAD, 2008.
- [6] B. Cohen, Incentives Build Robustness in BitTorrent, Workshop on Economics of Peer-to-Peer Systems, vol. 6, 2003.
- [7] J.A. Pouwelse et al., The bittorrent p2p file-sharing system: Measurements and analysis, International Workshop on Peer-to-Peer Systems (IPTPS), 2005.
- [8] M. Izal et al., Dissecting BitTorrent: Five Months in a Torrent's Lifetime, Passive and Active Measurements, vol. 2004, 2004.
- [9] T. Karagiannis et al., Is P2P dying or just hiding, IEEE Globecom, 2004.
- [10] A. Bellissimo, B.N. Levine, and P. Shenoy, Exploring the use of BitTorrent as the basis for a large trace repository, Technical Report 04-41, University of Massachusetts Amherst, 2004.
- [11] L. Guo et al., A performance study of BitTorrent-like peer-to-peer systems, Selected Areas in Communications, IEEE Journal on, vol. 25, 2007, o. 155—169.
- [12] M. Ripeanu et al., Gifting technologies: A BitTorrent case study, First Monday, vol. 11, 2006.
- [13] D. Qiu and R. Srikant, Modeling and performance analysis of BitTorrent-like peer-to-peer networks, SIGCOMM Comput. Commun. Rev., vol. 34, 2004, o. 367—378.
- [14] G. Simon and R. David, Design Challenges of Information Driven P2P Networks, Automation and Applied Computer Science Workshop (AACS), 2007, o. 147-158.
- [15] Isohunt, Releases; <http://isohunt.com/release>.
- [16] Mininova, Mininova : The ultimate BitTorrent source!; <http://www.mininova.org/>.
- [17] IMDB, The Internet Movie Database (IMDb); <http://www.imdb.com/>.
- [18] MetaBrainz-Foundation, MusicBrainz; <http://musicbrainz.org/>.