

Power-law Revisited: A Large Scale Measurement Study of P2P Content Popularity

György Dán

School of Electrical Engineering
KTH, Royal Institute of Technology
Stockholm, Sweden

Niklas Carlsson

Department of Computer Science
University of Calgary
Calgary, Canada

Abstract—The popularity of contents on the Internet is often said to follow a Zipf-like distribution. Different measurement studies showed, however, significantly different distributions depending on the measurement methodology they followed. We performed a large-scale measurement of the most popular peer-to-peer (P2P) content distribution system, BitTorrent, over eleven months. We collected data on a daily to weekly basis from 500 to 800 trackers, with information about 40 to 60 million peers that participated in the distribution of over 10 million torrents. Based on these measurements we show how fundamental characteristics of the observed distribution of content popularity change depending on the measurement methodology and the length of the observation interval. We show that while short-term or small-scale measurements can conclude that the popularity of contents exhibits a power-law tail, the tail is likely exponentially decreasing, especially over long time intervals.

I. INTRODUCTION

P2P content popularity has received significant research interest. Numerous works measured the instantaneous popularity and the popularity of contents over a time interval [1], [2], [3], [4]. The instantaneous popularity is defined as the number of peers that simultaneously participate in the distribution of the content. It influences the amount of control and data traffic in the overlay, and the efficiency of proximity-aware protocols [5]. Measurement results suggest that the instantaneous popularity of P2P content follows a power-law with an exponential cutoff [4]. The popularity of contents over time is defined as the number of times the contents are downloaded in a time interval. It reflects the amount of data and control traffic (e.g., search) in the overlay, and affects the efficiency of caching for P2P traffic [2], [3]. Several measurements showed that the download popularity follows a Zipf-like distribution, i.e., has a power-law tail [1], [3].

Ideally, measurements of content popularity should be based on probability sampling methods; i.e., the probability at which the units of the population are selected should be known. Probability sampling allows unbiased estimates of the population statistics (e.g., the distribution of the content popularity) to be produced. Probability sampling is, however, difficult to apply to large scale, dynamical systems, like P2P content distribution systems. Instead, measurements are often limited in geographical coverage, in scope and in time. Such opportunity sampling makes it difficult to assess what portion of the population is captured, and in general it is not understood how the sample

statistics relate to the population statistics.

The goal of our work is twofold. First, to show that, contrary to common belief, the popularity of contents in BitTorrent does not obey the power-law over long periods of time, and hence the efficiency of caching and locality-aware content distribution can actually be better than previously thought. Second, to show how the distribution of content popularity depends on the definition of popularity adopted, on the sampling method used and on the measurement interval. We base our findings on a measurement of more than 11 million contents over eleven months in BitTorrent. We observed 40 to 60 million peers on a weekly basis, and a total of more than 8 billion downloads.

II. BACKGROUND AND RELATED WORK

BitTorrent is the dominant peer-to-peer file sharing protocol on the Internet. BitTorrent relies on a set of *trackers*, which maintain state information about all peers currently having pieces of a particular file. The set of these peers is referred to as a *torrent*. Trackers record the number of peers downloading the content (called *leechers*), the number of peers that own the whole content (called *seeds*), and the number of times the contents were downloaded. A client that wants to download a content can learn about leechers and seeds that share a content by contacting a tracker at its *announce URL*. The address of the tracker and the identifier of the content (called *info hash*) is known to the peer from the torrent file, which it typically obtains from a torrent search engine, like *mininova.org*. A tracker can also be contacted at its *scrape URL*, in which case it returns the number of seeds, leechers and completed downloads for a specific torrent, or for all torrents it tracks.

A. Related work

Zipf's law states that if objects are ranked in order of their frequencies, the frequency for the object with rank r follows a *power-law*, $f_{\text{Zipf}(f_1, \theta)}(r) = f_1 r^{-\theta}$, where θ is the Zipf exponent. A generalization of Zipf's law is the Zipf-Mandelbrot law, for which $f_{\text{MZipf}(f_1, \lambda, \theta)}(r) = f_1 (\lambda + r)^{-\theta}$. The head of the distribution is flattened, i.e., the popularity of the top ranked objects is low compared to Zipf's law, but the tail follows a power-law. The origins of Zipf's law are in linguistics, but linguistics has for some time considered models beyond Zipf's law. A recently proposed model is the

generalized Zipf law [6]:

$$f_{GZipf}(f_1, \lambda, \mu, \theta)(r) = \frac{f_1}{[1 - \lambda/\mu + (\lambda/\mu)e^{(1/\theta)\mu r}]^\theta}, \quad (1)$$

which often captures the head and the tail of the rank frequency statistics better than Zipf's law [6]. Both Zipf's law and the Zipf-Mandelbrot law are limiting cases of this distribution. In particular, when $\mu \ll \lambda$ equation (1) reduces to Zipf-Mandelbrot's law for small r . For large values of r , however, it shows an exponential cutoff. That is, the tail of the distribution decreases faster than a power-law, hence the number of unpopular objects is *significantly lower* than according to Zipf's law or the Zipf-Mandelbrot law.

Zipf-like behavior was observed in the popularity of objects on the Web. Several studies confirmed that Web object popularity can be modeled using Zipf or related distributions [7], [8]. The popularity distribution of user generated contents, such as YouTube files, was found to follow Zipf's law except for the tail [9], to have a flattened head [10], as well as to follow Zipf's law [11]. The different results might be due to the different measurement methodologies: the first two studies were based on crawling, the third study was based on measurements at a university campus.

Measurement studies of the Gnutella and Kazaa P2P file-sharing systems were based on deep-packet inspection [1], [2] or on ultrapeers monitoring search requests [12], [3]. They agree that the rank popularity statistics of the number of downloads or queries of P2P content exhibits a flattened head compared to Zipf's law [1], [2]. Some proposed the Zipf-Mandelbrot law as a suitable model [3], while others used two Zipf curves to fit the body and the tail of the distribution, respectively [12]. However, the instantaneous content popularity shown in [4] based on a small sample of BitTorrent seems to follow Zipf's law with a sharp exponential cutoff. Other BitTorrent measurements focus on a single torrent (e.g., [13]).

Our work is novel in three aspects. First, to the best of our knowledge our measurement is the biggest in the literature in terms of the number of contents and peers observed, and the geographical and temporal coverage. Second, we show on the same measurement data set that content popularity shows different characteristics depending on its definition, on the measurement methodology and the length of the measurement period. Third, we treat our data set as a sample of a population, and test the validity of several hypothesis about the population-wide popularity distribution.

III. MEASUREMENT METHODOLOGY AND DATA

A. Sampling methods of BitTorrent

Measuring P2P content popularity can be done in a content-centric or in a peer-centric way. A content-centric measurement identifies contents first and then the peers that are interested in the individual contents; e.g., in BitTorrent, one obtains the info hashes of contents and the URLs of trackers that track the contents (e.g., [4]). A peer-centric measurement identifies peers first and then the contents they are interested in; e.g., via deep-packet inspection at a router (e.g., [1], [2])

or by monitoring overlay traffic (e.g., [3]). In both cases, the measured content popularity is a sample of the population-wide content popularity. In the following, we describe three practical opportunity sampling methods to measure BitTorrent popularity (*Mininova*, *PirateBay*, *PropPeer*), and two impractical probability sampling methods (*PropTor*, *UnifTor*).

Mininova: This content-centric sample is limited to the torrents that can be found on *mininova.org*, which was the most popular torrent search engine according to *www.alexa.com* on 1 Aug. 2008 (Alexa-rank of 75).

PirateBay: This content-centric sampling is limited to the torrents that are tracked by the tracker with most torrents throughout our measurement, *PirateBay*.

PropPeer: This peer-centric sampling consists of observing n_P samples. The samples belong to a torrent with a probability proportional to the popularity of the torrent. We count the samples observed in every torrent. *PropPeer* sampling can resemble uniform sampling or a local sampling (deep-packet inspection) of the peers, and it captures an (unknown) fraction of all P2P traffic.

PropTor: This content-centric sampling consists of observing n_T torrents at random with probabilities proportional to the popularities of the torrents. If we observe a torrent then we can measure the total popularity of the torrent.

UnifTor: This content-centric sampling consists of observing n_T torrents at random with uniform probabilities. If we observe a torrent we can measure the total popularity of the torrent. We use the last two sampling methods to understand what sampling method our measurement corresponds to.

B. Measurement data set

On 31 Aug. 2008, 15 Oct. 2008 and 31 Aug. 2009 we performed screen-scrapes of *mininova.org*, which was the most popular torrent search engine at the beginning of our measurement period. From the screen-scrapes we obtained the announce URLs of 1690 trackers and the info hash of 1.24 million contents. We constructed the scrape URLs of the 1690 trackers and scraped the trackers. We did not specify any info hash, so the trackers returned the scrape information for all torrents that they were tracking. This allowed us to efficiently obtain the number of leechers, seeds, and completed downloads as seen by the trackers that we determined via the screen-scrape of *mininova*. We performed the tracker-scrapes weekly between 15 Sept. 2008 and 17 Aug. 2009, and daily from 18 Sept. 2008 to 18 Oct. 2008.

We removed redundant tracker information for trackers that share information about the same swarms of peers, and identified 721 unique, responsive trackers. All scrapes were performed at 8pm GMT. The scrapes of all trackers were done simultaneously; obtaining the biggest scrape took less than half an hour. Due to the short scrape duration our data is a sequence of 31 daily and 49 weekly simultaneous snapshots of the information stored on 721 trackers world-wide.

IV. INSTANTANEOUS AND DOWNLOAD POPULARITY

In this section we present the rank popularity statistics of the instantaneous and the download popularity. We start with the

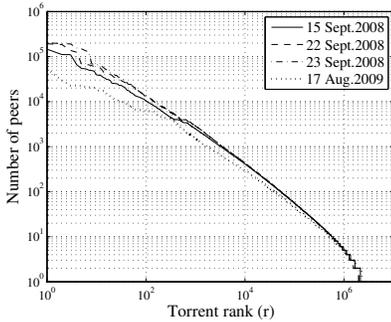


Fig. 1. Rank popularity statistics of the number of peers, observed on four different dates.

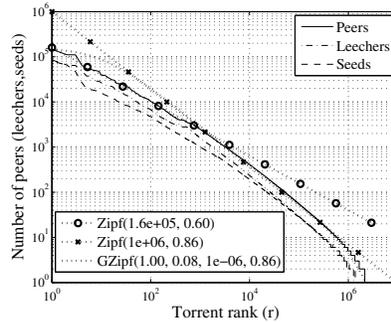


Fig. 2. Number of leechers, seeds and peers on 15 Sept. 2008. Number of peers is fitted with two Zipf curves, and a generalized Zipf curve.

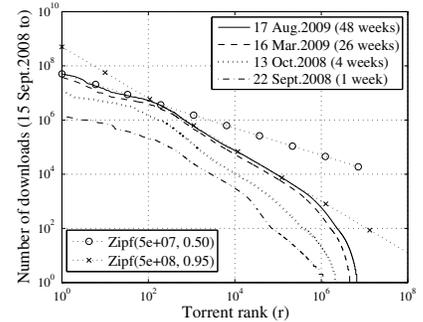


Fig. 3. Rank popularity statistics of the number of downloads over four intervals starting on 15 Sept. 2008.

instantaneous popularity; i.e., the concurrent number of peers participating in the distribution of the contents. This definition of popularity was used in [4], [5]. Figure 1 shows the rank popularity plots of the number of peers on four dates during our measurement. The curves show similar characteristics, and at a first look they follow Zipf's law. A closer inspection reveals, however, that each curve consists of a head, a trunk and a tail with different properties. The difference between the three regions is most obvious for the curve of 17 Aug. 2009. While the head and the trunk seem to follow Zipf's law with different exponents, the tail seems to decrease exponentially instead of according to a power-law.

We investigate these three regions of the rank popularity plot in Figure 2, which shows the rank popularity statistics of the number of peers, leechers and seeds observed on 15 Sept. 2008. Out of 5.23×10^6 torrents 2.93×10^6 were active, i.e., had at least one peer. The total number of peers was 4.2×10^7 . The rank popularity statistics of the number of leechers and the number of seeds is similar in shape to that of the number of peers. We use the number of peers to analyze the difference between the behavior of the head, the trunk and the tail of the distributions. We fitted a Zipf distribution to the head and the trunk of the measured distribution. The difference between the two fitted curves is significant, both in terms of the Zipf exponents, the domains they span and the maximum number of peers they predict. We also fitted a generalized Zipf distribution to the measured distribution. The generalized Zipf distribution does not capture the power-law of the head of the measured distribution, but it captures the trunk and the tail behavior. Consequently, the tail of the rank popularity statistics decreases exponentially. The tail of the distribution (above rank 5×10^5) represents 2.43×10^6 torrents and 8.2×10^6 peers, which is about 20 % of all peers.

We continue the analysis with the *download popularity*, i.e., the number of times the individual contents were downloaded in a time interval. This was the definition of popularity considered in [1], [2], [12]. Figure 3 shows the rank popularity plots of the number of downloads for four time intervals starting on 15 Sept. 2008. The total number of downloads were 2.23×10^8 , 1.31×10^9 , 5.86×10^9 and 8.34×10^9 over 1, 4, 26 and 48 weeks, respectively.

The curves for different intervals show similar characteristics. Comparing the curves we see that the trunks' slopes are

almost the same. Surprisingly, the heads' slopes increase as the interval gets longer. One would expect that the number of contents with approximately equally many downloads would increase over time, and hence the head of the distribution would become flatter as the time interval increases. Even more interesting is that the exponential cutoff of the tail is barely visible for the 1 week interval, but is very pronounced for longer intervals, e.g., for 48 weeks. We fitted two Zipf curves to the head and the trunk of the number of downloads over 48 weeks. The difference between the Zipf exponents for the two regions is almost a factor of two. The tail of the distribution shows an exponential cutoff, but still accounts for a significant part of the distribution: for 48 weeks there are 5.95×10^6 torrents above rank 5×10^5 with a total of 9.62×10^8 downloads, i.e., 11 % of all downloads.

The rank popularity data presented above raise two important questions. *First, are the characteristics of the rank popularity statistics an artifact of our measurement methodology? Second, would we observe similar characteristics if we had followed another measurement methodology?*

V. POWER-LAW OR EXPONENTIAL-CUTOFF?

We address the first question in the following.

Instantaneous popularity: We use the data from 15 Sept. 2008 shown in Figure 2 to test two hypotheses.

Power-Law Trunk Hypothesis (PLTH): The population-wide rank popularity statistics follow the Zipf curve fitted to the trunk of the distribution, but our measurement failed to capture the distribution's head and tail. If PLTH is true then the most popular torrent would have about 10^6 peers, the number of active torrents would be 9.5×10^6 and the total number of peers would be 6.1×10^7 . Based on the hypothetical number of peers and torrents we cannot reject PLTH.

Exponential Cutoff Hypothesis (EXCH): The population-wide rank popularity statistics exhibits an exponential cutoff (i.e., no power-law tail). To see if the exponential cutoff observed on our data is due to our sampling method, we sampled a hypothetical double-Zipf distribution ($\min(f_{\text{Zipf}(1.6e+5,0.6)}(r), f_{\text{Zipf}(1e6,0.86)}(r))$) fitted to the measured data. We took $n_T = 2.93 \times 10^6$ samples (i.e., the measured number of active torrents) from the double-Zipf distribution according to the *PropTor* and the *UnifTor* sampling methods, and recorded the discovered torrents. Figure 4 shows the double-Zipf distribution, the measured statistics, and the

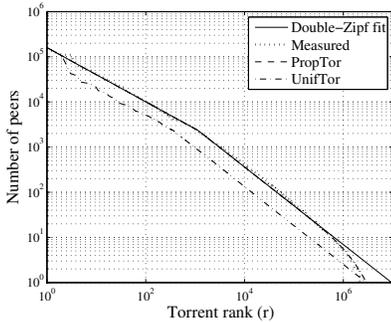


Fig. 4. Sampling from the hypothetical Double-Zipf distributed rank popularity statistics for 15 Sept. 2008

results of the sampling. The curve for *UnifTor* sampling lies well below the hypothetical popularity and does not exhibit the exponential cutoff. The curve for *PropTor* sampling shows, however, a very good match with our measurement data, with a Pearson product-moment correlation coefficient of 0.99. The observed number of peers matches as well; 4.23×10^7 for our measurement and 4.02×10^7 using *PropTor* sampling. Hence, we conclude that either (i) EXCH holds, or (ii) our sample of the trackers closely resembles *PropTor* sampling. Note that if EXCH does not hold, the number of active torrents is $9.5/2.93 \approx 3$ times higher, and the number of active peers is $5.5/4.23 \approx 1.3$ times higher than what we measured.

Download popularity: For the download popularity we consider two time intervals starting on 15 Sept. 2008: the 48 weeks interval shown in Figure 3, and the 4 weeks interval shown in Figure 5. The total number of downloads over 4 weeks was 1.31×10^9 , the number of active torrents was 2.29×10^6 . We fitted Zipf curves to the head and the trunk of the distribution, and fitted a generalized Zipf curve to the entire distribution. The generalized Zipf curve shows an excellent fit with the entire distribution. We test two hypotheses.

Double Power-Law Hypothesis (DPLH): The population-wide rank popularity statistics follow the Zipf curves fitted to the head and the trunk of the measured distribution, but our measurement failed to capture the distribution’s tail. If DPLH holds then there should be 1.77×10^7 active torrents based on the 4 weeks interval, and 1.43×10^9 active torrents based on the 48 weeks interval. If we compare these numbers to the hypothetical number of active torrents predicted by the PLTH for 15 Sept. 2008 (9.5×10^6), we see that DPLH for the number of downloads would require too many torrents to have nonzero downloads over 48 weeks. Hence, for the distribution of the number of downloads we reject the DPLH for 48 weeks but cannot reject it for 4 weeks.

Exponential Cutoff Hypothesis (EXCH): The good match between our data and the tail of the generalized Zipf distribution in Figure 5 suggests an exponential cutoff for the 4 weeks interval. To see whether the exponential cutoff could eventually be an artifact of our measurement methodology, we performed the same experiment as for the instantaneous popularity. We took $n_T = 2.29 \times 10^6$ samples from the DPLH hypothetical distribution according to *PropTor* and *UnifTor* sampling.

Figure 6 shows the results. The curve for *UnifTor* sampling

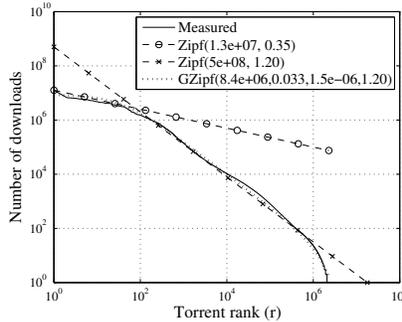


Fig. 5. Number of downloads over 4 weeks starting on 15 Sept. 2008. Fitted with two Zipf curves, and a generalized Zipf curve.

lies well below the hypothetical popularity and does not exhibit the exponential cutoff. Surprisingly, the curve for *PropTor* sampling shows a good match with our measurement data. The Pearson product-moment correlation coefficient is 0.99, and the observed number of downloads is 1.31×10^9 for our measurement and 1.21×10^9 using *PropTor* sampling. Hence, for the 4 weeks interval either EXCH holds or DPLH holds and our sampling resembles *PropTor* sampling. Note that if EXCH does not hold, there are $17.7/2.29 \approx 7.8$ times more active torrents world-wide than what we measured. For the 48 weeks interval we rejected DPLH, hence either EXCH holds or the tail of the population-wide distribution follows Zipf’s law but with a higher exponent than that of the trunk. In the latter case our sampling resembles *PropTor* sampling. In both cases, the number of non-cacheable (1 download only) contents is *orders of magnitude less* than under the DPLH hypothesis.

VI. THE IMPACT OF SAMPLING

In the following we address the second question, i.e., we investigate how the measured distribution of content popularity depends on the measurement methodology.

Instantaneous popularity: We applied the five sampling methods described in Section III-A to the measured popularity distribution of 15 Sept. 2008. The results are shown in Figure 7. The number of torrents in the *Mininova* sample is 9.7×10^5 , out of which 4.95×10^5 were active. The number of torrents in the *PirateBay* sample is 6.64×10^5 , out of which 6.55×10^5 were active. For *PropPeer* we used $n_P = 4.23 \times 10^5$; i.e., 1 % of the total number of peers in our measured data set. For *PropTor* and *UnifTor* we used $n_T = 6.55 \times 10^5$; i.e., the number of active torrents in the *PirateBay* sample.

The *Mininova*, *UnifTor* and *PropTor* samples overrepresent the most popular torrents. The *PirateBay* sample resembles the shape of the original distribution much closer. One would expect the *PropPeer* sample to have the same shape as the complete distribution, but it does not exhibit the exponential cutoff. The exponential cutoff disappears as the sample size decreases, which is in accordance with observations on the word frequency in the corpora of natural languages [14]. Hence, *PropPeer* gives the impression that the power-law holds for the tail of the distribution, even if the population-wide distribution has an exponential cutoff (i.e., EXCH holds).

Download popularity: We applied the same five sampling methods to the number of downloads over 4 weeks starting

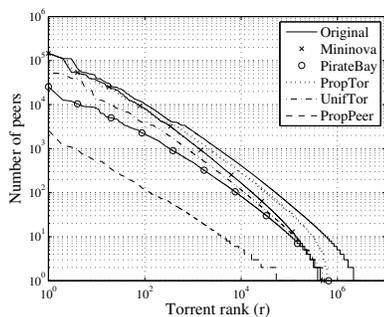


Fig. 7. Rank popularity distribution of the number of peers for five methods of sampling the measured statistics on 15 Sept. 2008

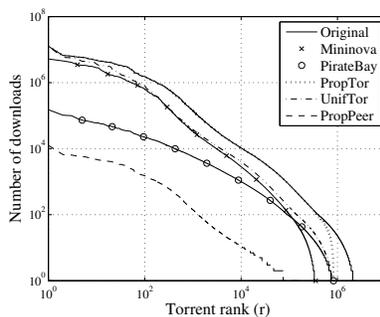


Fig. 8. Rank popularity statistics for five methods of sampling the measured number of downloads in 4 weeks starting on 15 Sept. 2008.

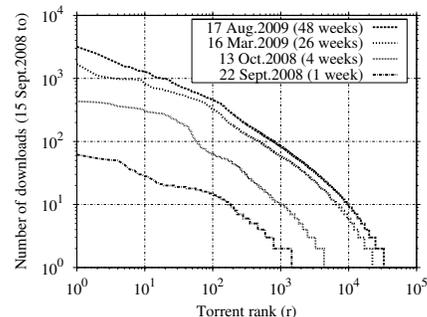


Fig. 9. Rank popularity statistics of the measured number of downloads over four intervals starting on 15 Sept. 2008 on campus.

on 15 Sept. 2008. Figure 8 shows the results. The number of torrents in the *Mininova* sample is 6.30×10^5 , out of which 3.53×10^5 were active (at least one download). The number of torrents in the *PirateBay* sample is 1.69×10^6 , out of which 8.29×10^5 were active. For *PropPeer* we used $n_P = 1.31 \times 10^6$, i.e., 0.1 % of the total number of downloads. For *PropTor* and *UnifTor* we used $n_T = 8.29 \times 10^5$, i.e., the number of active torrents in the *PirateBay* sample.

The *Mininova*, *UnifTor* and *PropTor* samples overrepresent the torrents with most downloads, similar to when used to sample the number of peers. Surprisingly, unlike for the instantaneous popularity, the *PirateBay* sample differs significantly from the original popularity distribution; one can hardly identify the trunk of the distribution. The *PropPeer* sample captures the head and the trunk of the distribution, but it again fails to capture the exponential cutoff, hence giving the impression that the tail follows a power-law. As *PropPeer* is closely related to deep-packet inspection, our results show that local measurements should be used with care to infer the characteristics of global popularity. Similarly, small scale uniform sampling of the peer population does not reveal all characteristics of content popularity.

Local file popularity: Finally, we consider the sample that would be observed from a local organization. For this experiment we captured every HTTP-tracker request at the University of Calgary campus with roughly 33,000 students and staff (during the same measurement period as our primary data set). Figure 9 shows the number of download completions as reported by local university clients. Curves are shown for the same time-durations as in Figure 3. Comparing these two figures, we note that the popularity shows similar characteristics (but at a smaller scale) as observed in our data set. Clearly, the sampling methodology is important when trying to capture the wide-area popularity characteristics.

VII. CONCLUSION

Based on a large scale measurement of BitTorrent content popularity performed over eleven months we showed that previous beliefs about P2P content popularity do not hold on a global scale. We found that Zipf's law might describe the instantaneous popularity. The download popularity over long time intervals does not follow a power-law, but neither small scale nor short term measurements would be able to capture

the exponential cutoff. We showed how methods of sampling the same global popularity distribution affect the observed content popularity, and provided insights about the limitations and biases of the different methodologies.

VIII. ACKNOWLEDGEMENTS

The authors would like to thank Martin Arlitt for collecting the University data set. Financial support for this research was provided by the Information Circle of Research Excellence (iCore) in the Province of Alberta, Canada.

REFERENCES

- [1] K. Gummadi, R. Dunn, S. Saroiu, S. Gribble, H. Levy, and J. Zahorjan, "Measurement, modeling, and analysis of a peer-to-peer file-sharing workload," in *Proc. SOSP*, October 2003, pp. 314–329.
- [2] A. Wierzbicki, N. Leibowitz, M. Ripeanu, and R. Woźniak, "Cache replacement policies for P2P file sharing protocols," *Euro. Trans. on Telecomm.*, vol. 15, pp. 559–569, 2004.
- [3] M. Hefeeda and O. Saleh, "Traffic modeling and proportional partial caching for peer-to-peer systems," *IEEE/ACM Trans. on Networking*, vol. 16, no. 6, pp. 1447–1460, 2008.
- [4] L. Guo, S. Chen, Z. Xiao, E. Tan, X. Ding, and X. Zhang, "Measurement, Analysis, and Modeling of BitTorrent-like Systems," in *Proc. ACM IMC*, Oct. 2005.
- [5] D. Choffnes and F. Bustamante, "Taming the torrent: A practical approach to reducing cross-ISP traffic in P2P systems," in *Proc. of ACM SIGCOMM*, Aug. 2008.
- [6] M. Montemurro, "Beyond the Zipf-Mandelbrot law in quantitative linguistics," *Physica A.: Statistical Mechanics and its Applications*, vol. 300, no. 3–4, pp. 567–578, 2001.
- [7] L. Breslau, P. Cao, L. Fan, G. Phillips, and S. Shenker, "Web caching and Zipf-like distributions: Evidence and implications," in *Proc. of IEEE INFOCOM*, 1999, pp. 126–134.
- [8] A. Mahanti, C. Williamson, and D. Eager, "Traffic analysis of a Web proxy caching hierarchy," *IEEE Network*, vol. 14, no. 3, pp. 16–23, May/June 2000.
- [9] X. Cheng, C. Dale, and J. Lui, "Understanding the characteristics of Internet short video sharing: Youtube as a case study," in *Proc. IWQoS*, June 2008.
- [10] M. Cha, H. Kwak, P. Rodriguez, Y. Ahn, and S. Moon, "I Tube, You Tube, Everybody Tubes: Analyzing the World's Largest User Generated Content Video System," in *Proc. ACM IMC*, October 2007.
- [11] P. Gill, M. Arlitt, Z. Li, and A. Mahanti, "YouTube Traffic Characterization: A View from the Edge," in *Proc. ACM IMC*, October 2007.
- [12] A. Klemm, C. Lindemann, M. K. Vernon, and O. P. Waldhorst, "Characterizing the query behavior in peer-to-peer file sharing systems," in *Proc. of ACM IMC*, Oct. 2004, pp. 55–67.
- [13] M. Izal, G. Urvoy-Keller, E. Biersack, P. Felber, A. Al Hamra, and L. Garcés-Ericc, "Dissecting bittorrent: Five months in a torrent's lifetime," in *Proc. of Passive and Active Measurement (PAM)*, 2004.
- [14] O. Tripp and D. G. Feitelson, "Zipf's law revisited," School of Computer Science and Engineering, The Hebrew University of Jerusalem, Tech. Rep. 2007-115, Aug 2007.