





Network	Date Range	Files	GUIDs	Records
Gnutella (FOI only)	10/1/2010 – 9/18/2011	139,604	775,941	870,134,671
Gnutella Browse	6/1/2009 – 9/18/2011	87,506,518	570,206	434,849,112
eMule (FOI only)	10/1/2010 – 9/18/2011	29,458	1,895,804	133,925,130
IRC (no file data)	6/2/2011 – 9/18/2011	N/A	N/A	7,272,739
Ares (no file data)	5/31/2011 – 9/18/2011	N/A	N/A	17,706,744

**Table 1: All datasets are observations of CP activity only, but IRC and Ares data do not contain information about files or GUIDs. Except when otherwise stated, a record corresponds to a law enforcement observation and contains date, time, IP address, application-level identifier, geographic location as determined by an IP geolocation database, and a file hash.**

in January 2009, we began deploying a set of forensic tools to investigators in the U.S. and internationally for online investigation of p2p CP trafficking.

Prior to our collaborative efforts, the standard method for online CP investigation was to make isolated cases: leads were not shared among agencies or officers, other than by phone or email. Officers leveraged their own experience to prioritize suspects.

**Tools.** Our suite of tools, called *RoundUp* [14], has enabled seamless sharing of *plain view* observations of online CP and associated activities on various filesharing networks. The shared data, collected in order to make these cases, provide each investigator with a longitudinal view of CP offenders and provide a method of triage for selecting targets for further investigation; and of course, the data enable this study. Because over 2,000 investigators have been trained on our tool to date, and because it is in use by hundreds of investigators daily, the aggregate set of observations we have used for this study is incredibly detailed. The tools are still in use, and currently, law enforcement execute approximately 150 search warrants nationwide per month based on data collected using our tools. We do not, however, present search warrant or arrest data in this study<sup>5</sup>.

**Datasets.** Our datasets, summarized in Table 1, include law enforcement observations from Gnutella and eMule p2p networks. The Gnutella and eMule datasets span a one-year period from October 1, 2010 to September 18, 2011. Each record in these datasets corresponds to a law enforcement observation of a particular peer making available one or more FOI, and minimally contains date, time, IP address, application-level identifier, geographic location as determined by an IP geolocation database, and a file hash.

Most file sharing protocols include an application-level identifier unique to an installation of the application. In both Gnutella and eMule, these identifiers are persistent across users’ sessions, and are referred to as *GUIDs* (globally unique identifiers). Peers on these networks are uniquely identified by their GUID, and we use peer and GUID interchangeably to identify unique running instances of the corresponding p2p software.

All FOI are uniquely identified using hash values; law enforcement manually classify files as FOI from a variety of sources, such as post-arrest forensic analyses. An enormous number of such FOI are shared on Gnutella and eMule. Respectively, there are 139,604 and 29,458 known FOI shared by 775,941 and 1,895,804 GUIDs. Our tool searched only for FOI in a list containing about 384,000 entries; this list was updated several times over the course of this study.

<sup>5</sup>Our study’s procedures were approved by our Institutional Review Boards.

It is a small sample: the National Center for Missing and Exploited Children reports reviewing more than 60 million child pornography images and videos<sup>6</sup>. As such, our work presents only a lower bound on the amount of activity present in these networks.

In a limited fashion, we use two other datasets. Our IRC dataset, based on a more recent tool that we developed, covers a four-month period from June to September 2011. The IRC dataset is a log of IP addresses that were involved in public activity related to the sexual exploitation of children in public chatrooms; no file observations are in this dataset. We also use a dataset of CP-related activity on the Ares p2p network<sup>7</sup> collected using a tool we did not write, but collected by the same law enforcement officers responsible for all data in this paper. The Ares dataset contains only IP addresses and has no information about files shared, but addresses were only logged for peers that shared known FOI.

**Other Details.** Gnutella allows a peer to be *browsed* and thus investigators can enumerate all files shared by peers. Our *Gnutella Browse* dataset consists entirely of peer browses and includes all files a peer is sharing, not just FOI. Some Gnutella peers cannot be browsed; we collected only FOI data from these peers. eMule does not permit browses. Regardless, each of these datasets includes only peers that share one or more FOI; peers without FOI were not logged.

We draw a distinction between a time-limited view of a peer’s shared files and the set of all files with which a given peer was ever observed. We define a GUID’s *library* to be the set of files that were observed being shared by that GUID on a given day. A GUID’s *corpus* is the set of all files shared by that GUID over the entire duration of the study. In both cases, we typically only include FOI, but we make it clear when a corpus or library includes non-FOI observed as the result of a browse.

## 4. AVAILABILITY AND RESILIENCE

In investigating the trafficking of CP on p2p networks, the goal of law enforcement is to prioritize criminals whose arrest will have the greatest impact. But the strategy to achieve this goal depends upon the impact desired: finding contact offenders who go otherwise unreported, finding those who create new CP, and decreasing the availability of FOI on the network are all priorities. In this section, we focus on strategies for reducing the availability of FOI.

Effective CP removal strategies are especially important as a means to prioritize law enforcement’s limited resources and time. After online evidence is collected, days or weeks

<sup>6</sup>See <http://www.missingkids.com/missingkids/servlet/NewsEventServlet?&PageId=4604>.

<sup>7</sup><http://aresgalaxy.sourceforge.net/>

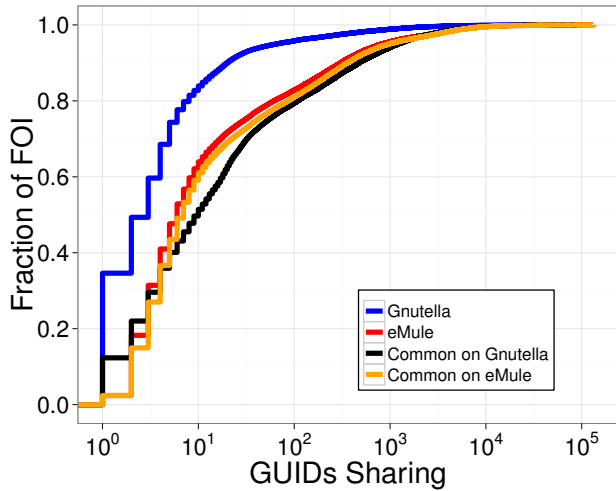


Figure 1: Redundancy of FOI (files of interest) among multiple GUIDs as a CDF. Some files are seen on both networks, but the distribution of these observations is different. The “Common on Gnutella” line shows the CDF of these common files as seen on Gnutella, and similarly for the “Common on eMule” line.

of off-line processes are required in each case until an arrest is made. Additional resources are required to go to trial. It is infeasible for investigators to arrest all users sharing CP and remove all FOI. Investigators need a triage strategy for deciding upon which small fraction of online leads to act.

An enormous set of perpetrators are active every day around the world. Even with unlimited resources, U.S. law enforcement can only partially impact file availability. Our results, discussed below, suggest the need for a coordinated international effort.

## 4.1 FOI Redundancy and Availability

Before we further discuss the implications of removing files, we characterize the redundancy and availability of FOI on Gnutella and eMule.

### 4.1.1 File Redundancy Across GUIDs

Many FOI on Gnutella and eMule are not widely redundant among GUIDs within the same network. Figure 1 shows the relative *redundancy* of FOI, which is the number of GUIDs that possess and make available each file. The distribution is presented as a cumulative distribution function (CDF), which shows on the  $y$ -axis the fraction of FOI that are shared by *at most*  $x$  GUIDs. For example, 90% of files on Gnutella were shared by at most 20 GUIDs; 99% of files were shared by at most 1,167 GUIDs; and 99.9% of files were shared by at most 9,129 GUIDs.

Figure 1 also shows the relative redundancy for the subset of FOI appearing on both networks. The set of files common to both networks is significantly more redundantly shared on each network than the set of all files on each network. There is a high degree of FOI overlap among the two networks: 26,136 of the FOI on the eMule network (nearly 89%) were also seen on the Gnutella network, and 97% of Gnutella GUIDs were observed with at least one file that can be found on the eMule network. The overall low redundancy of most files suggests the strategy of prioritizing the investigation of users who possess a large amount of less redundant FOI in order

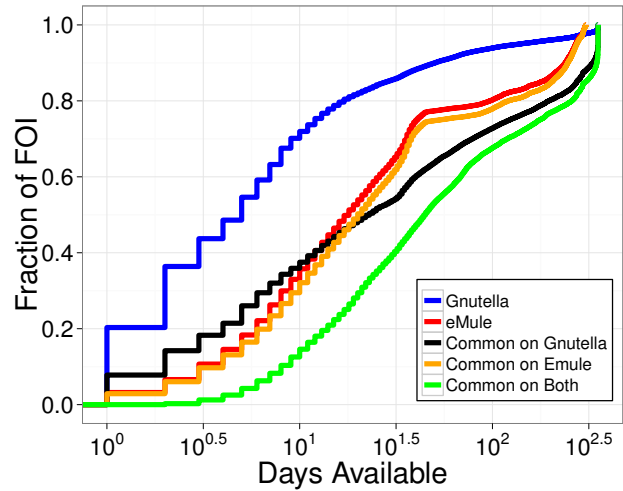


Figure 2: CDF showing the days available per FOI (during 353 days for Gnutella and 329 days for eMule). As in Figure 1, the “Common on Gnutella” line shows the CDF of files common to both networks as seen on Gnutella, and similarly for the “Common on eMule” line. The “Common on Both” line shows these common files available on either network on any given day.

to remove it from the network and prevent its proliferation. An easily intuited proxy for this measure is to target GUIDs who possess large corpora. Since most FOI are relatively less redundant, the GUIDs with the largest libraries likely have the most FOI with low redundancy.

### 4.1.2 File Availability Across Days

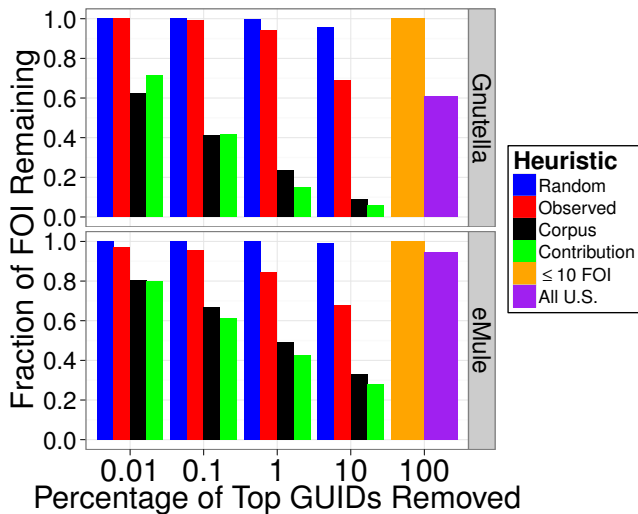
We say a file is *available* on a given day if at least one peer is sharing that file on that day. This approach is simple in that it does not take into account bandwidth and reachability considerations, which are difficult to measure globally. We do not expect this definition to limit the applicability of our results, as the assumption of high bandwidth and reachability is conservative from the perspective of law enforcement.

Figure 2 plots the availability of FOI as a CDF on a semi-log scale. Gnutella files tend to have lower availability than eMule, with 80% of files available for more than one day; about 30% are available for more than 10 days; and about 5% of files are available for more than 100 days. Generally, files that are available for a single day are unique to a specific GUID; files that tend to have longer availability are possessed by many GUIDs, not all of whom are online on a given day. Again we see that the files that are common to both networks are more available than is typical on each individual network: about 30% of these common files are available for more than 100 days. We have also calculated that on a daily basis, an average of 9,712 distinct files are available, with a peak of 32,020 files during our study.

## 4.2 Law Enforcement Strategy

Our *law enforcement model* is as follows. Investigators have a global, historical view of GUIDs and their corpora, including known FOI and other files. Investigators look to reduce FOI *availability*, by arresting the users that correspond to peers and removing their corpora from the network. Investigators aim to remove files from the network completely.

Content can be removed from these networks only by



**Figure 3: The remaining fraction of FOI available at least one day given a percentage of GUIDs removed according to different heuristics: random, number of days observed, corpus size, and contribution to file availability on Gnutella and eMule. We removed the top 0.01, 0.1, 1, and 10 percent of GUIDs according to each heuristic. In Gnutella, the Corpus and Contribution heuristics achieve equal results when 0.113% of GUIDs are removed. Also shown is the impact of removing 100% of peers with 10 or fewer FOI, and 100% of peers in the U.S.**

arresting users and taking their shared libraries offline, as the protocols and implementations inhibit falsifying or polluting content. Our goal is to find out which peers should be removed such that we minimize the number of files that are available at least one day.

In Appendix A of our technical report we show that this problem is NP-Hard [9]. Here, we evaluate four greedy heuristics aimed at reducing the availability of CP by removing peers. Our evaluation consists of removing subsets of peers from the data and examining the effect on availability. Specifically, we examine the following heuristics: (i) removing peers that were *observed* most often, i.e., largest number of days observed; (ii) removing peers with the largest *corpus* size; (iii) removing peers with the largest *contribution* to availability (as defined below); and (iv) removing peers selected randomly, as a baseline. For an arbitrary file on an arbitrary day,  $n$  peers possess that file. We say that each peer provides a *file-contribution* of  $\frac{1}{n}$ th of that file. A peer’s *contribution* to file availability is the sum of the file-contributions of the files in their corpus over the duration of the study.

An alternative measure of availability is *daily redundancy*, the number of peers that share a file on a specific day. The algorithm to optimally reduce the maximum redundancy over all files shared is simple: remove the peers with the largest corpus size first. It is unclear that minimizing redundancy, unless it is to zero (equivalent to unavailability), is useful or effective. To evaluate the effect of reducing redundancy to a small value, we would require reachability, bandwidth, and propagation models of the underlying p2p overlays. Thus, we do not consider daily redundancy further.

#### 4.2.1 Comparison of the Efficiency of Heuristics

Figure 3 compares the effectiveness of each of the above heuristics. Interestingly, removing the peers that were seen

the most often has almost no effect on the availability of FOI. Removing peers by either contribution or corpus size is most effective; these measures are correlated, so their similarity in performance is unsurprising.

The vast majority of files are shared only by a relatively small set of prolific GUIDs. Consider Gnutella (similar trends hold for eMule): If we remove the top 0.01% of 775,941 GUIDs as determined by corpus size, only 59% of the known FOI remain available in the network. In other words, 41% of the unique files on the network are made available by a group of only about 80 GUIDs. The top 0.01% have 3,242 distinct FOI on average, with the top peer possessing about 25,000 FOI. Most of these files, however, are only available for a relatively short amount of time; as Figure 2 shows, only 28% are available for more than 10 days during our study. Some of this is due to the relatively low number of days these high-contributing GUIDs were observed; this also explains the failure of the observed days heuristic. These prolific GUIDs have a worldwide presence. Removing them requires tremendous multi-national cooperation as we discuss below.

#### 4.2.2 Impact of Geography on Availability

Our data are mostly based on the efforts of U.S. law enforcement, and the files they are looking for are arguably tuned to U.S. perpetrators. As law enforcement agents are limited by jurisdiction, the locational diversity of these users provides a resistance to the straightforward approach of prioritizing them. Only a small majority of top Gnutella GUIDs (by corpus size)—57 out of 100—are located in the U.S. The rightmost bar (“All U.S.”) in Figure 3 shows the reduction in availability if we restrict our removal to U.S. GUIDs (that is, GUIDs with an IP located in the U.S.) only. Note that we remove *all* such GUIDs in our analysis, a clearly infeasible approach in practice. Just 30% of files are unavailable (internationally) after removing all GUIDs in the U.S.; removing just the top 0.01% internationally (a group of about 80 GUIDs) has a similar effect, suggesting the utility of a coordinated international approach.

Within the U.S., the problem is similarly large in scope. The top 5% of GUIDs in the U.S. comprises a set of 14,410 GUIDs, each with a corpus of at least 40 known FOI. Due to the weeks of manual effort required for each arrest, the limited resources in the U.S. allow for 3,100 arrests per year for both offline and online offenses [27].

#### 4.2.3 Impact of Low-Sharing GUIDs on Availability

A large portion of GUIDs have comparatively few files. As shown in Figure 4, about 82% have 10 or fewer FOI. There are several reasons peers may appear to have few files. They may have files that are CP, but are not yet known to be FOI. They may be downloading FOI and not subsequently sharing them. They may have downloaded the files incidental to other activities. Finally, they may simply be sharing a smaller library. We expected removal of such low-sharing users to impact file availability significantly, since very many peers possess few files. Contrary to our expectations, removal of these GUIDs sharing few files has essentially no effect on file availability, as shown in the second-rightmost bar in Figure 3 (“≤ 10 FOI”). This result provides further evidence that file availability depends primarily on those GUIDs with the largest corpora, though it does not consider the contribution to redundancy that these low-sharing GUIDs provide.

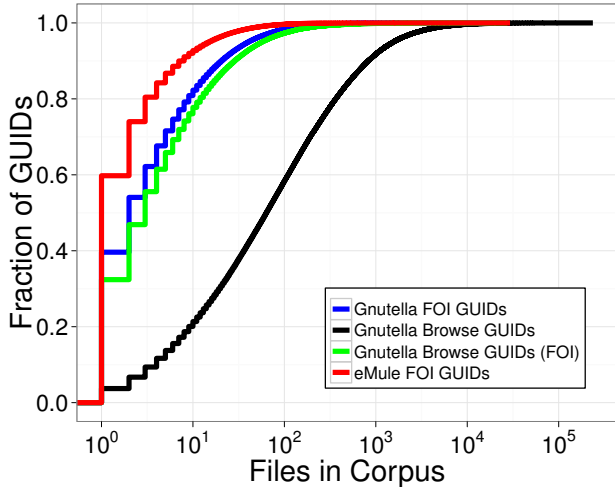


Figure 4: CDF showing the corpus size per GUID, for various measurement types. The black line (“Gnutella Browse GUIDs”) show the corpus size distribution for all files seen at GUIDs whose libraries were browsed, and the corresponding green line shows the distribution of FOIs within those browses. The other two lines show all FOI observed in any manner. (n.b., eMule does not allow browses.) Most GUIDs have very few files in their corpus. We give a week-by-week breakdown of Gnutella library sizes in our technical report [9].

## 5. COMPARING AGGRESSIVE PEERS

In Section 4, we show that strategies for removing content from the entire ecosystem must target offenders from all countries. In the absence of a unified effort—and no such collaboration exists—investigators need a triage strategy. In this section, we characterize triage metrics for local investigators. Ideally, investigators would target the most *dangerous* offenders: those that are personally, physically abusing children. Unfortunately, such information is typically not available until months or years after arrest [26].

In lieu of that ideal, local investigators can target peers that are the most aggressive offenders: peers that exhibit greater evidence of *intent* [11] beyond the average case, which is an important practical legal concern. This includes peers such as those that are online for the longest duration and share the largest number of FOI. Similarly, investigators may target offenders that are conduits between p2p network communities (e.g., by sharing on both eMule and Gnutella), or offenders that seek to escape detection and justice by using Tor or network relays.

We quantify the activity of six subgroups of aggressive peers sharing FOI. We characterize the contribution of each subgroup to the duration of CP availability and the amount of CP content. The subgroups are: (i) the top 10% of GUIDs sharing the largest corpora; (ii) the top 10% of GUIDs seen sharing FOI the most number of days; (iii) the top 10% of GUIDs ranked by the *contribution* metric defined in Section 4.2; (iv) the set of GUIDs sharing FOI on at least two p2p networks (linked by IP address); (v) GUIDs that use a known Tor exit node; (vi) GUIDs sharing FOI that use an IP address that we infer to be a non-Tor relay.

Our results show that all of these subgroups are more active than a group that consists of all peers that we observed. The exception is the subgroup of GUIDs using non-Tor relays, as

Identifier	Network	
	Gnutella	eMule
All GUIDs	775,941	1,895,804
Multi-Network GUIDs	84,925 (11%)	147,904 (7.8%)
Tor GUIDs	3,666 (0.47%)	16,290 (0.86%)
Tor GUIDs (> 2 days)	2,592 (0.33%)	11,998 (0.63%)
Relayed GUIDs	76,478 (9.9%)	78,223 (4.1%)
Top 10% Observed	84,235 (11%)	190,797 (10%)
Top 10% By Corpus	77,782 (10%)	189,951 (10%)
Top 10% By Contr.	77,595 (10%)	189,581 (10%)

Table 2: Sizes of each GUID subgroup. Definitions of each subgroup appear in this section.

Network	IP Addresses		
	Total	Private	Tor
Gnutella	3,025,530	32,195	7,357
eMule	5,643,350	1,256	21,025
Ares	1,714,894	225	1,799
IRC	88,658	245	746

Table 3: Number of IP addresses per network observed sharing FOI. In the case of IRC, the IP addresses correspond to clients observed in public chat rooms related to child sexual exploitation. The Tor column refers to the number of distinct public IPs where Tor-using GUIDs were seen, including but not limited to known Tor exit nodes.

we explain below. The differences of each subgroup to the set of all GUIDs are significant ( $p < 0.001$ ).

Below we provide characteristics of each subgroup, and details of the behavior of each. For example, we show that GUIDs using Tor to share FOI use Tor irregularly, and therefore their true IP addresses are easily identifiable. Due to space limitations, we omit Gnutella data from some graphs where they largely correspond to the eMule data. The full set of graphs are available in our technical report [9].

### 5.1 Peer Subgroups

The size of each subgroup is shown in Table 2. The size of the top 10% by corpus and observed days subgroups are slightly larger than 10%. This variability is due to ties in the ranked lists of GUIDs. We include all such GUIDs to avoid arbitrary tie-breaking.

#### 5.1.1 Top 10% Groupings

Users can actively participate in p2p networks in two primary ways: by contributing a large number of files or a large amount of time. For example, one peer may share 100 files for a single day, and another may share a single file for 100 days. In the first case, the content is large but other peers have only a limited time to take advantage. In the second case, the content is small but other peers will find it easier to gain access to the content. It is vital for investigators to address both types of activity; the contribution metric balances these two concerns.

For these reasons, we create three subgroups corresponding to the 10% of GUIDs with the largest corpora of files ( $\mathcal{F}$ ), the 10% with the most days observed online ( $\mathcal{D}$ ), and the top 10% of GUIDs when ranked by the contribution metric ( $\mathcal{C}$ ). There is substantial but not overwhelming overlap among these subgroups. The overlap in Gnutella, as defined by Jaccard’s index,  $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ , is  $J(\mathcal{C}, \mathcal{F}) = 0.51$  and  $J(\mathcal{C}, \mathcal{D}) = 0.28$ ; the eMule subgroups overlap similarly.

Networks		IP Addresses Intersection		
A	B	%A	$A \cap B$	%B
Gnutella $\cap$	eMule	6.8%	199,824	3.1%
	IRC	0.1%	3,562	4.1%
	Ares	1.0%	30,596	1.8%
eMule $\cap$	IRC	0.1%	4,654	5.3%
	Ares	0.9%	56,921	3.3%
IRC $\cap$	Ares	2.1%	1,813	0.1%
Intersection of all			308	

**Table 4: Overlap of IP addresses across multiple networks, excluding Tor IPs and private IPs. A small but significant set of IPs were seen across multiple networks, indicating particularly active users.**

### 5.1.2 Multi-Network Peers

Law enforcement are interested in users that are active on multiple p2p networks. Such users are more aggressive in terms of assisting in the distribution and availability of content to two communities, possibly acting as a bridge. We identify the set of GUIDs in Gnutella that are active in another network by finding all IP addresses in our Gnutella dataset that also appear in any of our eMule, Ares, or IRC datasets, and correspondingly in eMule for those that appear in any of the Gnutella, Ares, or IRC datasets. We refer to GUIDs in these sets as *multi-network GUIDs*.

The total number of IPs addresses, private IPs<sup>8</sup>, and IPs used by GUIDs that also used known Tor exit nodes that we observed for each of these networks is shown in Table 3. Generally, private IPs are the result of sub-optimally or misconfigured end-user applications, as opposed to indicating privacy awareness. In contrast, Tor use indicates privacy-aware users. Table 4 shows the size of each pairwise network overlap. For all such intersections, we first remove private IPs and Tor exit nodes (as listed in the Tor consensus files<sup>9</sup>). Of all network pairs, the Gnutella-eMule overlap is the largest.

The union of all three intersections comprises our 84,925 GUID multi-network subgroup for Gnutella. We perform a similar calculation for eMule, resulting in 147,904 GUIDs.

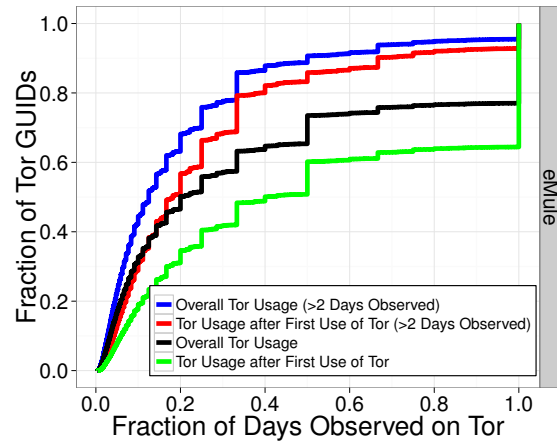
### 5.1.3 Peers that Use Tor

Peers that use Tor are of interest to law enforcement because they are actively masking their identities, thwarting investigations of this crime. Tor does not filter application-level data: GUIDs are passed through to investigators, and thus appear in our dataset as well. We define a GUID as a *Tor GUID* if it was ever observed as having an IP address listed as a Tor exit node in the Tor consensus for the date of the observation. When a Tor GUID’s IP is a known Tor exit node we say that the GUID is *using* Tor. As Table 2 shows, this set is not large on either network: 3,666 GUIDs for Gnutella and 16,290 GUIDs for eMule.

It is striking that the vast majority of Tor GUIDs do not use Tor consistently, which makes it possible to detect their true IP address. In Figure 5, we show the CDFs of *overall Tor usage*. In both networks, only about a quarter of the Tor GUIDs used Tor every time they were observed. More significantly, for these GUIDs, under 40% consistently use

<sup>8</sup>Private IP addresses are those which are non-routable on the public Internet, self-assigned, or otherwise invalid, as defined by RFC 5735.

<sup>9</sup>Consensus files contain the list of IPs addresses acting as exit nodes on a daily basis; see <https://metrics.torproject.org/data.html>



**Figure 5: CDF of Tor usage per GUID for eMule. GUIDs do not use Tor consistently after first being observed at a Tor IP. Under 40% of Tor GUIDs consistently used Tor after first being observed using it. When considering only Tor GUIDs seen on >2 days (which comprise about 70% of all Tor GUIDs), the rate falls to below 10%. The Gnutella data show similar characteristics.**

Tor after their first use of Tor.

When we examine these 40% of nodes that used Tor consistently, we found that most were observed on the Gnutella and eMule networks for only one or two days. Therefore, we recomputed the distribution of Tor usage for the subset of Tor GUIDs observed three or more days, which is over 70% of all Tor GUIDs. We again also computed the CDFs of Tor usage after first using Tor. The resulting CDFs are the upper lines in Figure 5. In sum, over 90% of GUIDs using Tor for more than two days on eMule and Gnutella are easily linked back to a non-Tor IP address, one that is most likely their real location.

This irregular use could be due to ignorance of how Tor works, careless configuration, or frustration with the lower throughput of Tor. It is well known that Tor’s *design* does not offer technical protection to p2p users because it does not hide identifying application-level data [16]. In contrast, we provide the first empirical evidence that Tor *users* do not use the software consistently, even those with a strong reason to so. Regardless of the quality of Tor’s security, this evidence strongly suggests that its usability (its interface, its effects upon perceived speed, or some other factor) is lacking. We conclude that the use of Tor, as observed in practice, poses only a small hurdle to investigators. Reports by the Tor developers that “Journalists use Tor to communicate more safely with whistleblowers and dissidents”<sup>10</sup> should give one pause, as there is no evidence that those groups are significantly more or less tech-savvy than the users we study.

### 5.1.4 Peers that Use Suspected Relays

The final subgroup we identify is a set of peers that are using IPs that we suspect are relays (other than Tor exit nodes). To create this subgroup, we first collected the set of non-Tor IP addresses used by GUIDs that also used a Tor exit node. We discard the IPs that hosted fewer than four GUIDs (267,035 in the case of Gnutella, and 1,671,419

<sup>10</sup>Quoted from <https://www.torproject.org/about/overview.html.en>

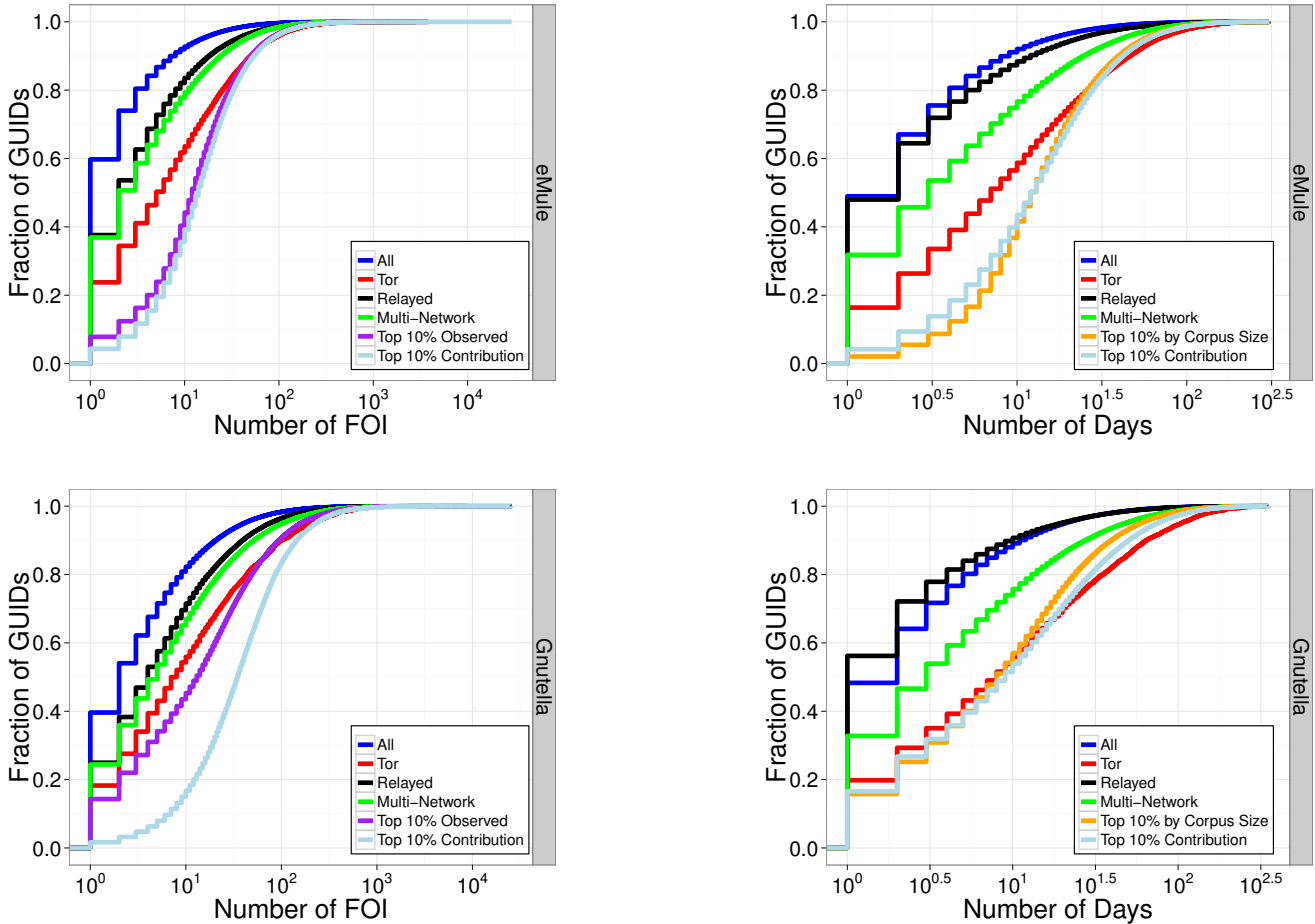


Figure 6: Characterizations, as CDFs, of per-GUID corpora and days observed for eMule and Gnutella. The aggressive subgroups, sans relayed, appear to be more active in their trafficking, having more FOI and uptime than the average peer sharing FOI.

for eMule), and we nominate the remaining IPs as potential relays. Finally, we create the subgroup of *relayed GUIDs* as the set of GUIDs seen using the potential relays. We cannot validate these GUIDs as having definitely used relays; for example, it may be the potential relays are IP addresses that get reassigned frequently. However, we consider their use of these shared IPs sufficient to define them as a distinct set.

## 5.2 A Comparison of Peer Behavior

There are substantive and statistically significant differences among the subgroups in terms of per-GUID corpora and number of days observed. These differences can be seen in Figure 6 and are summarized in Table 5. In particular, the subgroups generally have a larger corpus size and more days observed online than the set of all GUIDs. The three top-10% subgroups show this effect most strongly, but the Tor subgroup and multi-network subgroups show similar effects. Notably, these latter two subgroups are selected independently of corpus size and days online. This result confirms a hypothesis that tech-savvy groups, whether through Tor or multi-network use, are more active.

The set of GUIDs in the top 10% contribution subgroup represent a combination of the other aggressiveness metrics. This result can be viewed by comparing CDFs in the figure, or

by comparing means in the table. For example, the top 10% contribution subgroup’s mean corpus size is higher than the top 10% observed, and its mean number of days observed is higher than the top 10% corpus subgroup. The contribution metric could easily be parameterized to weight observations more heavily, though we do not show such results here.

The relayed subgroup in general has larger number of FOI than the all group, and appears online more days on average than the all group in eMule. However, the relayed subgroup shows fewer days observed online than the all group in the Gnutella network. This result suggests that either this subgroup, as we’ve defined it, is not particularly active, or that our process for identifying non-Tor relays is faulty. It may also be that the peers in the relayed subgroup are more successful at aliasing themselves as different GUIDs that appear on the network fewer number of days each. In the following section, we examine the general problem of user aliasing in this data set.

## 6. ANALYSIS OF USER ALIASING

The relationship between p2p network GUIDs and real users is not one-to-one in our dataset. In fact, it is possible for a single user to correspond to multiple, distinct GUIDs. We refer to this phenomenon as *user aliasing*, and for some



GUID Group		Mean Value (99% CI)	
		Corpus Size	Days Obs.
Gnutella	All	10.9 (10.7, 11.1)	5.2 (5.2, 5.2)
	For	43.9 (39.0, 49.6)	23.4 (21.8, 25.1)
	Relayed	18.9 (18.3, 19.5)	4.8 (4.7, 4.9)
	Multi-Network	25.9 (24.9, 27.0)	10.8 (10.6, 11.0)
	Top 10% Obs.	41.8 (40.7, 43.0)	28.7 (28.5, 29.0)
	Top 10% Corp.	75.9 (74.3, 77.7)	16.2 (16.0, 16.5)
	Top 10% Contr.	69.1 (67.6, 70.9)	19.5 (19.3, 19.8)
eMule	All	4.3 (4.3, 4.4)	4.1 (4.1, 4.1)
	For	21.2 (19.9, 22.5)	17.4 (16.9, 18.0)
	Relayed	9.2 (8.9, 9.6)	5.5 (5.4, 5.6)
	Multi-Network	10.8 (10.6, 11.0)	9.5 (9.4, 9.7)
	Top 10% Obs.	23.5 (23.2, 23.8)	22.3 (22.2, 22.4)
	Top 10% Corp.	27.8 (27.4, 28.5)	18.7 (18.6, 18.8)
	Top 10% Contr.	25.8 (25.4, 26.5)	19.0 (18.9, 19.1)

**Table 5: The expected value and 99% confidence interval of each characteristic for each subgroup of GUIDs. Each subgroup’s mean differs from the mean of the “All” group. Each such difference is statistically significant ( $p < 0.001$ ), as determined by a computational permutation test ( $R = 10,000$ ). Confidence intervals are computed by bootstrap ( $R = 10,000$ ).**

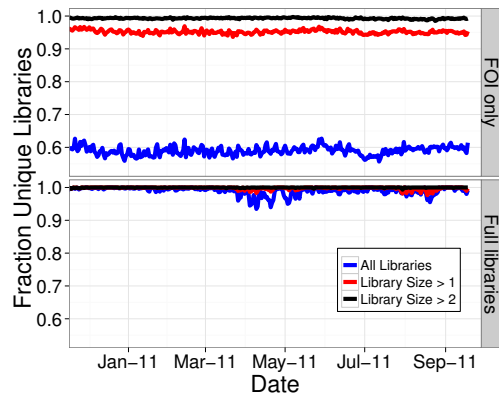
users it is intentional. In this section, we examine observable user aliasing, and we also attempt to quantify its effects upon the analyses in the previous sections. In sum, we find that GUIDs that share at least three FOI any given day generally have distinct libraries. In Gnutella, we can compare all files shared by a GUID, and in that case users sharing a library of at least two files are generally distinct on a given day. We also find little evidence to suggest users are changing their GUIDs and then continuing to share the same library or a portion of it later that day. Parallel results generally held for eMule, though without the ability to browse eMule user libraries, we are less certain of that result.

The true user aliasing rate in our data is unknowable to us. However, the reasons for deliberate aliasing can be enumerated: (i) if a user has two computers (or multiple accounts on a single computer), each with an installation of Gnutella, he will control two unique GUIDs; and (ii) a user may reinstall or upgrade their p2p client on a single computer or otherwise modify their GUID over time. We have no way of detecting the first case from only network data; however, the second case can be detected if the user does not alter what files they are sharing, as the file library acts as a kind of signature for the user. It is this latter case that we evaluate in the remainder of this section.

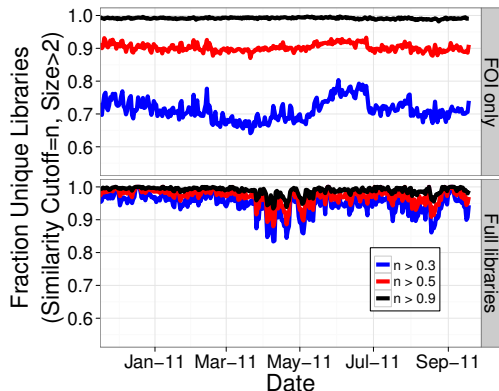
Most users, as identified by GUIDs, are seen with very small libraries of a single file or two. This fact is illustrated in Figure 4 in Section 4 (and in a week-by-week breakdown in our technical report [9]). We posit that such small libraries are not particularly differentiable. By excluding them, we can determine a lower bound on the user aliasing of type (ii) that may be occurring.

**Analysis and Results.** We computed day-to-day similarities between Gnutella libraries to determine a lower bound on user aliasing, or alternatively, an upper bound on the number of unique libraries present in the dataset. Generally, we found most libraries to be distinct.

Figure 7 shows a comparison of Gnutella GUID libraries, plotting the fraction of GUIDs with libraries that are a unique collection of files. In the upper portion of the figure, a comparison is made of just the files of interest at each



**Figure 7: Fraction of Gnutella GUIDs with unique libraries on specific days, where uniqueness is defined as libraries that completely match. When considering libraries of at least two FOI, approximately 95% are unique. Similar results hold for eMule. When considering full (browsed) libraries, over 93% are unique.**



**Figure 8: Fraction of Gnutella GUIDs with a unique library, where uniqueness is defined as there being no other library with a similarity greater than  $n$ . The similarity of two libraries is defined their Jaccard index,  $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ . On most days, 90% of libraries have no more than half their files in common.**

GUID; the lower portion compares all files in the library of each GUID (from a Gnutella browse request). GUIDs that have tens of files or more are easy to distinguish from others.

Figure 7 shows that in general, GUIDs with a single file are easily aliased with other GUIDs with the same single file: only about 58% of GUIDs have unique libraries on a given day of our dataset. Among the 40% of Gnutella GUIDs that have two or more FOI, over 95% have unique libraries. Among the 25% of GUIDs with three or more FOI, over 99% have distinct libraries.

Fewer aliases are present when comparisons can be made of the complete libraries, as is possible with Gnutella browse information, by including all files, not just FOI. This is illustrated in the lower portion of Figure 7. Note that GUIDs with a single FOI typically possess more than one file, and thus they are more likely to be unique. Typically, GUIDs seen with two or more files in their library had a unique library about 95% of the time; GUIDs with three or more files were unique over 99% of the time.

The above data suggest that we can treat GUIDs as

uniquely distinguishable when their libraries contain at least two FOI or when we consider all files that they share. The analysis also suggests that users are rarely if ever changing their GUID and appearing on the same day with the same library. They would appear as aliases if so, and if this was common, the fraction of unique libraries would be lower.

Based on a similar analysis, we also make the claim that there is no compelling evidence that many users are changing GUIDs appearing on the network that day and preserving only *most* of their shared libraries. Figure 8 quantifies the uniqueness of partial and complete libraries using the Jaccard index:  $J(A, B) = \frac{|A \cap B|}{|A \cup B|}$ . In the upper portion of Figure 8, we see that for GUIDs with at least three FOI, approximately 90% of libraries have no more than half their files in common on most days of our study. In the lower portion of Figure 8, we compare all files in each GUID’s library, not just FOI. Here we see on most days, 85% of libraries have no more than 30% of their files in common.

A limitation of our calculations above is that we compare GUIDs only within a day’s time. We haven’t computed all-pairs, all-times equivalence or similarity among GUID libraries across multiple days because the computation is too lengthy to handle in a reasonable timescale for our dataset.

## 7. MEASUREMENT LIMITATIONS

The limitations of our study prevent us from providing more than conservative lower bounds on the observable activity of CP perpetrators. First, our set of known FOI is likely biased towards files and filenames shared by traffickers in the U.S. Traffickers in other countries are likely underestimated by our study. Second, all of our records would ideally be associated with a browse, in other words, a complete listing of the peer’s current files. eMule does not support browse functionality at all, and investigators do not browse all Gnutella peers on all days. For example, a peer may be identified as having file *A* on day 1 and day 3, but that file is not seen on day 2 because the appropriate keyword or hash search was not run. As a result, we may be underestimating the amount of CP content possessed by each peer as well as the number of days they are online. Third, peers that are online more often are also more likely to be found using a search. We might be underestimating the number of peers that are rarely online and have few files.

On the other hand, one user might have one or more installations of the p2p client software, with each installation showing up as a different GUID. Hence, the number of GUIDs in these networks serves as a rough upper bound on the number of users (for the FOI we knew about).

We also note that before, during, and after the collection of the datasets we analyze, law enforcement were and are active in investigating and arresting CP traffickers. We do not know which peers were removed from the network, and we do not take these removals into account in our analyses. The specific metrics we report on do not rely on linking arrested users to a search warrant and the outcome of a subsequent trial.

## 8. RELATED WORK

**Ecosystems & Underground Economies.** Our work is similar in theme to a body of work exploring economic characteristics of network-based ecosystems [2, 5, 10, 13, 18]. For example, the irregular use of Tor by the peers in our

dataset might be explained by recent work showing that users abandon privacy for short-term benefits [1].

**Content Availability in P2P Systems.** A large body of related work on p2p systems investigates availability, performance, and issues related to the use of incentives [3, 4, 6, 15, 17, 19, 20, 30]. Unlike our work, these studies mostly focus on understanding and analyzing the unique properties of p2p networks and their users’ behavior, and do not specifically target CP or separate aggressive subgroups.

**CP Trafficking in P2P Systems.** Prior studies of CP-related trafficking on the Internet have a limited scope. They are mostly indicative of the alarming presence of contraband rather than comprehensively quantifying how the files are being shared [8, 12, 21, 22]. All previous work focused on CP (rather than copyright violations) is based on only CP-related search terms rather than verified content [7, 8, 12, 23, 24].

The exception is our own prior work [14], where we analyze CP-related activity on Gnutella during a five-month period with no overlap with the study in this paper. In that work, we show that the correspondence between IP addresses and application-level identifiers is not one-to-one, and then propose proactive methods of differentiating the end hosts. In contrast, our focus in this work is on reducing availability and characterizing peer behavior.

## 9. CONCLUSIONS AND FUTURE WORK

The criminal trafficking of CP on p2p networks is widespread, with no easy answers for law enforcement looking to curtail it. The diversity in peers’ location, the redundancy of their libraries, and the many p2p networks, coupled with limited law enforcement resources, dictate triage as a strategy. Specifically, investigators should carefully choose peers to investigate and remove from p2p networks.

We have shown that although naive approaches to triage are ineffective and optimal approaches are NP-Hard, tractable heuristics yield reasonable and useful results. Further, the use of these heuristics are complemented by our discovery of aggressive subgroups of CP traffickers, where such groups correspond to aspects of the heuristics we identified. Prioritizing enforcement in these groups is both effective and easily understandable by LE and policymakers alike.

Further, we have found no significant evidence of users attempting to hide by altering their visible file libraries: peers’ libraries are largely unique, strongly implying a unique user behind each such library. Some users do use Tor, but surprisingly, most do so inconsistently, making the investigation of such users straightforward.

It is an open question as to whether network-observable behaviors, such as interest in particular types of imagery, correlate with off-line behaviors of interest to LE, such as child molestation. In ongoing and future interdisciplinary work, we will explore this interesting question.

**Acknowledgements.** This work is based in part upon the following awards. UMass personnel were supported in part by CNS-1018615 (re: CP p2p trafficking), CNS-0905349 (re: geographic/mobile analysis in techreport), or 2008-CE-CX-K005 (re: tool development). Wolak was supported in part by CNS-1016788 (re: CP p2p trafficking). Albrecht was supported in part by CAREER award CNS-0845349. We are grateful for the comments and assistance of Hanna Wallach.

## 10. REFERENCES

- [1] A. Acquisti and J. Grossklags. Privacy and rationality in individual decision making. *IEEE Security & Privacy*, 3:26–33, January 2005.
- [2] J. Caballero, C. Grier, C. Kreibich, and V. Paxson. Measuring pay-per-install: the commoditization of malware distribution. In *Proc. USENIX Security*, pages 1–13, Aug. 2011.
- [3] M. Cha, H. Kwak, P. Rodriguez, Y.-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*, pages 1–14, 2007.
- [4] N. Christin, A. Weigend, and J. Chuang. Content availability, pollution and poisoning in file sharing peer-to-peer networks. In *Proc. ACM Electronic Commerce (EC)*, pages 68–77, 2005.
- [5] J. Franklin, V. Paxson, A. Perrig, and S. Savage. An inquiry into the nature and causes of the wealth of internet miscreants. In *Proc. ACM CCS*, pages 375–388, Oct. 2007.
- [6] 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. *SIGOPS Oper. Syst. Rev.*, 37:314–329, Dec. 2003.
- [7] D. Hughes, P. Rayson, J. Walkerdine, K. Lee, P. Greenwood, A. Rashid, C. May-Chahal, and M. Brennan. Supporting law enforcement in digital communities through natural language analysis. In *Proc. Intl Workshop Computational Forensics*, pages 122–134, Berlin, Heidelberg, 2008.
- [8] D. Hughes, J. Walkerdine, G. Coulson, and S. Gibson. Peer-to-peer: is deviant behavior the norm on P2P file-sharing networks? *IEEE Distributed Systems Online*, 7(2), Feb. 2006.
- [9] R. Hurley et al. Measurement and analysis of child pornography trafficking on p2p networks. Technical report, University of Massachusetts Amherst, May 2013. UM-CS-2013-007.
- [10] C. Kanich, N. Weavery, D. McCoy, T. Halvorson, C. Kreibich, K. Levchenko, V. Paxson, G. Voelker, and S. Savage. Show me the money: characterizing spam-advertised revenue. In *Proc. USENIX Security*, Aug 2011.
- [11] O. S. Kerr. *Computer Crime Law*. West / Thompson Reuters, St. Paul, Minnesota, 2nd edition, 2009.
- [12] M. Latapy, C. Magnien, and R. Fournier. Quantifying paedophile queries in a large p2p system. In *Prof. IEEE INFOCOM*, pages 401–405, April 2011.
- [13] K. Levchenko, A. Pitsillidis, N. Chachra, B. Enright, M. Félegyházi, C. Grier, T. Halvorson, C. Kanich, C. Kreibich, H. Liu, D. McCoy, N. Weaver, V. Paxson, G. Voelker, and S. Savage. Click trajectories: End-to-end analysis of the spam value chain. In *Proc. IEEE Symp. Security & Privacy*, pages 431–446, Nov 2011.
- [14] M. Liberatore, B. Levine, and C. Shields. Strengthening Forensic Investigations of Child Pornography on P2P Networks. In *Proc. ACM CoNext*, Nov 2010.
- [15] X. Lou and K. Hwang. Collusive Piracy Prevention in P2P Content Delivery Networks. *IEEE Transactions on Computers*, 58(7):970–983, July 2009.
- [16] P. Manils, A. Chaabane, S. Le Blond, M. Kaafar, C. Castelluccia, A. Legout, and W. Dabbous. Compromising Tor Anonymity Exploiting P2P Information Leakage. In *Proc. HotPets*. (See also <http://blog.torproject.org/blog/bittorrent-over-tor-isnt-good-idea>), July 2010.
- [17] D. Menasche, A. Rocha, B. Li, D. Towsley, and A. Venkataramani. Content availability and bundling in swarming systems. In *Proc. ACM CoNext*, pages 121–132, 2009.
- [18] M. Motoyama, D. McCoy, K. Levchenko, S. Savage, and G. Voelker. Dirty jobs: the role of freelance labor in web service abuse. In *Proc. USENIX Security*, Aug 2011.
- [19] J. Otto, M. Sánchez, D. Choffnes, F. Bustamante, and G. Siganos. On blind mice and the elephant: understanding the network impact of a large distributed system. In *Proc. ACM Sigcomm*, pages 110–121, Aug 2011.
- [20] M. Piatek, T. Isdal, T. Anderson, A. Krishnamurthy, and A. Venkataramani. Do incentives build robustness in bit torrent. In *Proc. USENIX NSDI*, Apr 2007.
- [21] J. Prichard, P. Watters, and C. Spiranic. Internet subcultures and pathways to the use of child pornography. *Computer Law and Security Review*, 27(6):585–600, 2011.
- [22] J. Ropelato. Internet pornography statistics. <http://internet-filter-review.toptenreviews.com/internet-pornography-statistics.html>, 2007.
- [23] M. Rutgaizer, Y. Shavitt, O. Vertman, and N. Zilberman. Detecting pedophile activity in bittorrent networks. In *Proc. PAM Conf.*, Vienna, Austria, March 2012.
- [24] C. Steel. Child pornography in peer-to-peer networks. *Child Abuse & Neglect*, 33(8):560–568, 2009.
- [25] D. Stutzbach and R. Rejaie. Understanding churn in peer-to-peer networks. In *Proc. ACM IMC*, pages 189–202, Aug 2006.
- [26] United States Sentencing Commission. Public hearing on federal child pornography crimes. [http://www.usssc.gov/Legislative\\_and\\_Public\\_Affairs/Public\\_Hearings\\_and\\_Meetings/20120215-16/Agenda\\_15.htm](http://www.usssc.gov/Legislative_and_Public_Affairs/Public_Hearings_and_Meetings/20120215-16/Agenda_15.htm), February 15, 2012.
- [27] U.S. Dept. of Justice. The National Strategy for Child Exploitation Prevention and Interdiction: A Report to Congress. <http://www.projectsafefchildhood.gov/docs/natstrategyreport.pdf>, August 2010.
- [28] R. J. Walls, B. N. Levine, M. Liberatore, and C. Shields. Effective Digital Forensics Research is Investigator-Centric. In *Proc. USENIX Workshop on Hot Topics in Security (HotSec)*, August 2011.
- [29] J. Wolak, D. Finkelhor, and K. Mitchell. Child-Pornography Possessors Arrested in Internet-Related Crimes: Findings From the NJOV Study. Technical report, National Center for Missing & Exploited Children, 2005.
- [30] C. Zhang, P. Dhungel, D. Wu, and K. Ross. Unraveling the BitTorrent Ecosystem. *IEEE Transactions on Parallel and Distributed Systems*, 22(7):1164–1177, July 2011.