

# Is Deviant Behaviour the Norm on P2P File-Sharing Networks?

Daniel Hughes<sup>1</sup>, Stephen Gibson<sup>2</sup>, James Walkerdine<sup>1</sup>, Geoff Coulson<sup>1</sup>

<sup>1</sup>Computing Department, <sup>2</sup>Psychology Department, Lancaster University, Lancaster, UK  
{ hughesdr, walkerdi, geoff } @comp.lancs.ac.uk | s.j.gibson1@lancaster.ac.uk

## Abstract

*Major international law-enforcement initiatives are underway to fight the distribution of illegal pornography via the Internet. In this paper we examine the role of peer-to-peer (P2P) file-sharing networks in illegal pornography distribution. First, we investigate the contention that these networks are especially implicated in illegal pornography distribution. Our finding is that this conventional wisdom is in fact flawed: while we confirm that P2P networks are indeed used for the distribution of this material, we also find that the vast majority of it is produced and consumed by a tiny minority of P2P users who, furthermore, have little or no interaction with the wider law-abiding P2P community. On the basis of this finding, we outline a socio-technical approach through which P2P communities (which are in general as opposed to illegal pornography as the rest of the population) might themselves collectively subvert the activities of the disconnected minority that deal in illegal pornography. We believe that such a ‘self-policing’ approach is potentially far more realistic and effective than the commonly-proposed “blunt instrument” measures of attempting to close down or centrally police P2P communities.*

## 1. Introduction

Since the release of Napster [1] in 1999, peer-to-peer (P2P) file-sharing communities have been growing extremely rapidly. Today, several P2P networks (e.g. [2], [3], [4]) boast users numbering in the millions. Due to both scalability concerns and legal issues, today’s P2P networks have moved away from the semi-centralized approach typified by Napster, towards more scalable and anonymous P2P architectures [5]. These decentralized networks, because they exist in the absence of any central authority, provide a new and interesting context for the expression of human social behaviour.

However, the activities of the individuals who compose P2P communities are sometimes at odds with what is considered acceptable by authorities in the ‘real world’. The most obvious example of this is the phenomenon of copyright infringement, and the resulting copyright enforcement activities by organizations like the Recording Industry Association of America (RIAA) [6]. This issue has generated lively debate in the P2P research community, and

has given rise to a significant body of work devoted, on one side of the fence, to policing user behavior [7] and, on the other side, to creating radically anonymous networks in which the enforcement of any form of control is impossible [8]. Copyright infringement has also exercised the legal world in the shape of recent challenges to the legality of P2P file-sharing systems [9]. More broadly, fundamental questions have been raised about the nature of copyright law and its enforcement, which, despite vigorous debate, have yet to be resolved.

A more clear-cut example of online activity which society finds unacceptable is the use of Internet applications (including P2P networks) to distribute illegal pornography. Major law-enforcement efforts are currently underway in both Europe and North America [10] that target the distributors of such material; and these have resulted in a number of high-profile prosecutions. In addition, there have been large-scale public awareness campaigns [11] regarding the dangers that the Internet poses to children. Furthermore, a recent attempt by the California legislature to outlaw P2P file-sharing listed amongst its justifications the sharing of illegal pornographic material [9].

To gauge the nature and extent of P2P-based distribution of illegal pornography, we report in this paper on an extensive analysis of pornography-related resource-discovery traffic in the Gnutella P2P network [12]. We chose Gnutella because it is a good example of a large-scale, decentralized, anonymous, P2P file-sharing system; and it also has a well-studied user-base and an open protocol specification. Specifically, Gnutella was chosen over Fastrack [2] and eDonkey [3] as these latter networks feature username and password authentication and therefore cannot be considered anonymous. However, while our experiments only address Gnutella, there is no reason to suppose that Gnutella users download any more or less pornography than users of other anonymous P2P networks; therefore our results may be considered indicative of what one might expect elsewhere.

From the results of our analysis, we do indeed find that a small yet significant proportion of Gnutella activity relates to illegal pornography. But does this imply that such activity is widespread in the file-sharing population? On the contrary, we also find that this activity is carried out by a small, yet particularly active sub-community of users, and that searching for and distributing illegal pornography is *not* a behavioral norm.

This result has an important technical implication: it opens up avenues for the suppression of illegal pornographic activity in P2P networks that are more subtle and yet potentially more effective than the commonly-proposed “blunt instrument” approaches of attempting to close down or centrally police P2P communities. In particular, it opens the way for technical mechanisms that can empower and facilitate the law-abiding majority in P2P communities to self-police their own communities and thereby subvert illegal activities that they find just as abhorrent as the rest of society.

The remainder of this paper is structured as follows: Section 2 provides necessary background on the operation of Gnutella, and section 3 provides background on the known socio-psychological effects of anonymity on online behaviour. Section 4 then describes our experimental set-up and the results of our experimental analysis. Following this, section 5 discusses the implications of our findings and section 6 speculates on technical mechanisms to discourage illegal pornography in a self-policing manner along the lines suggested above. Finally, section 7 suggests directions for future work.

## 2. The Gnutella Network

Gnutella is an open distributed protocol designed to support the discovery and transfer of files among its users. Gnutella and similar decentralized file-sharing systems are considered to be more anonymous than earlier semi-centralized systems such as Napster [1] as these earlier systems made use of 3<sup>rd</sup> party indexing servers to store information about each peer and the files it was making available to the network. In entirely decentralized networks like Gnutella, no such entity has knowledge of the peers or files available on the network.

In technical terms, Gnutella builds an unstructured, decentralized, ‘Cayley-Tree’ network [5]. As with any decentralized P2P network, participating peers are required to forward network maintenance and file discovery messages, and to share files on the network. The protocol itself is very simple and uses just five message types as shown in Table 1.

**Table 1 – The Gnutella protocol messages**

Message	Description
PING	This message is flooded onto the network to discover peers. Peers that are willing to accept a connection to the sender respond with a PONG.
PONG	The response to a PING. Contains connection information and data regarding the number and size of files the sending peer is sharing.
QUERY	A search message with a plain-text payload. If a peer receiving a QUERY has matching data, it generates a QUERYHIT.

QUERYHIT	A response to a QUERY. Contains the information needed to acquire the requested data.
PUSH	A mechanism to support downloads from firewalled peers. If both peers are behind a firewall, no file transfer is possible.

Having connected to the Gnutella network using PING and PONG, a peer’s subsequent activities fall into two distinct phases: *i*) discovering resources, and *ii*) transferring resources (files).

To discover resources, a requesting peer forwards a QUERY message to its neighbors; each neighbor then forwards this message to some of *its* neighbors, and so on, thus ‘flooding’ the query onto the network. If a peer is able to satisfy an incoming query (i.e. it is sharing a file which matches a search-term contained in the QUERY message) it responds by sending a QUERYHIT message back along the same path. QUERYHITs contain the information required to subsequently acquire the requested file: i.e., the network address and port of the responding peer.

Having received one or more appropriate QUERYHITs, the requester selects a suitable peer, opens an HTTP connection to it, and downloads the target file directly using HTTP GET. Thus file transfer itself takes place outside of Gnutella proper.

## 3. The Effects of Anonymity on Online Behaviour

Recently, a significant amount of research has been devoted to the effects of anonymity and perceived unidentifiability in ‘computer mediated communication’ (CMC) [13]. Some researchers have argued that anonymity results in a greater likelihood of engaging in ‘deviant’ or ‘disinhibited’ online behaviour [14]. For example, [15] found that well over half of individuals visiting a website for non-pornographic material nevertheless attempted to access pornography when presented with an opportunity to do so. Such findings suggest that when online people may find it harder to resist the temptation to engage in behaviour that may ordinarily incur strong social disapproval or sanction.

Other researchers, however, have suggested that the consequences of anonymity in CMC may be better understood in terms of *group-specific social norms* [16]. According to this view, anonymity will only lead to deviant or illegal behaviour (as defined by general societal norms) if the norms appropriate to the particular situation allow for it. Thus, an individual will *not* necessarily be more likely to engage in behaviour that runs counter to general social norms when anonymous online, but will in fact be more likely to engage in behaviour which *conforms* to group-specific social norms. As just about any group identity may be relevant to self-definition when online (e.g. Internet user,

Canadian, Star Trek fan, person with a sexual preference for children), the type of behaviour possible is almost infinitely variable depending on the specific behavioural norms associated with the relevant group.

As a system that offers almost total anonymity, Gnutella is an ideal environment in which to evaluate these two competing theories of online behaviour. According to the first class of theories (which emphasise the generally negative effects of anonymity), no clear pattern should be detectable in the way users search for and serve illegal pornographic material. Users who provide and access such material will simply be acting in an individually disinhibited way, and any user of Gnutella is therefore a potential user of such material. In contrast, according to the second class of theories (those which draw attention to the importance of group-specific social norms), there should be a clear pattern detectable in the behaviour of those who search for and serve illegal material. Specifically, such users would be expected to form a distinct ‘sub-class’ within the wider class of Gnutella users.

For example, the effects of online anonymity on the behaviour of someone with a sexual interest in children would be to facilitate their downloading of images of childhood sexual abuse to the extent that such behaviour is normative for someone with such a sexual preference. Conversely, anonymity would *not* be expected to produce such behaviour in someone who does not identify themselves as having such a sexual preference. The crucial distinction is that in the former case the simple act of using Gnutella (or the Internet in general) is assumed to produce deviant behaviour, whereas in the latter case the anonymity associated with online activity merely facilitates the influence of group norms which are more-or-less already inscribed.

If it is true that the anonymity in P2P networks has generally negative effects upon user behaviour, then this would lend support to those who argue for wide-ranging legal restrictions on P2P technology, as was argued in recent legal action in California [9]. If, however, this behaviour is due to the influence of pre-inscribed group norms, it may be that the sharing of illegal sexual material on P2P file sharing networks merely reflects deeper issues in society, and that more subtle approaches to discouraging it are required.

#### 4. The Experiments

Our experimental work was based on intercepting and analysing QUERY and QUERYHIT messages on the Gnutella network. Essentially, analysing QUERY messages tells us what people are searching for, and analysing QUERYHIT messages tells us what people are offering to provide.

As each peer in Gnutella participates in routing all network messages, we can intercept these messages simply by deploying a modified peer onto the network which logs all the QUERY and QUERYHIT messages it routes. Using such a peer (based on the Jtella classes [17]), we monitored Gnutella traffic over a one month period between February 27<sup>th</sup> and March 27<sup>th</sup> 2005. We maximized the typicality of our sample base by connecting to the network as an ultra-peer [19], by maintaining a large number of incoming and outgoing connections, and by periodically re-connecting to different areas of the network.

The legality of various types of pornographic material varies from country to country, and even within countries there are disagreements over the precise letter of the law [20]. Therefore, for the purposes of the present study we limited our definition of illegality to those materials depicting *practices* which are clearly illegal under UK and international law. These practices are rape, incest, bestiality and the sexual abuse of children.

Our experiments addressed two questions, which are addressed in sections 4.1 and 4.2 respectively:

1. What proportion of Gnutella traffic relates to illegal pornography?
2. Is this the activity of a deviant sub-community?

#### 4.1 What proportion of Gnutella traffic relates to illegal pornography?

To answer this question we examined samples from three Saturdays that fell within our monitoring period (the Saturdays of March 5<sup>th</sup>, 12<sup>th</sup> and 19<sup>th</sup>). We focused on Saturdays due to the relatively higher level of traffic observed during weekends. From these samples, we randomly extracted 10,000 QUERYs and 10,000 QUERYHITs, and then manually classified each of these as relating to either illegal pornographic material (as defined above) or to ‘other’ material. The samples used in this classification, along with all the rest of the raw data, are available online at <http://polo.lancs.ac.uk/p2p/deviant>.

To ensure accurate classification, the messages were independently classified by two independent reviewers. Our approach was to classify messages as relating to illegal pornographic material if they could *only* be interpreted as referring to such material. Despite this conservative approach, some level of misclassification was inevitable due to the nature of plain-text searches. For example, while it is possible that a query for ‘*young girl*’ may be intended to retrieve illegal material, it could also refer to legal material (for example a song) and therefore such messages were not placed in the ‘illegal’ category.

**Table 1: Traffic relating to illegal material: reviewer 1**

	5 <sup>th</sup> March		12 <sup>th</sup> March		19 <sup>th</sup> March	
<b>QUERY</b>	1.2%	122	1.6%	156	1.6%	158
<b>QUERYHIT</b>	2.1%	206	3.0%	295	2.0%	195

The results of the classification are given in Tables 1 and 2. It can be seen that there is a high degree of correlation between the independent reviewers' classifications ( $p=0.3$  for QUERYS; and  $p=0.8$  for QUERYHITS).

**Table 2: Traffic relating to illegal material: reviewer 2**

	5 <sup>th</sup> March		12 <sup>th</sup> March		19 <sup>th</sup> March	
<b>QUERY</b>	1.4%	142	1.8%	184	1.7%	174
<b>QUERYHIT</b>	2.3%	234	3.0%	297	2.1%	208

The results were as follows:

- An average of 1.6% of QUERY messages were classified as relating to illegal pornography. The minimum value we observed was 1.2% on March 5<sup>th</sup>, rising to 1.8% on March 12<sup>th</sup>. The standard deviation between samples was 0.2%.
- An average of 2.4% of QUERYHIT messages relate to illegal pornography. The minimum value observed was 2% on March 19<sup>th</sup>, rising to 3% on March 12<sup>th</sup>. The standard deviation was 0.7%.

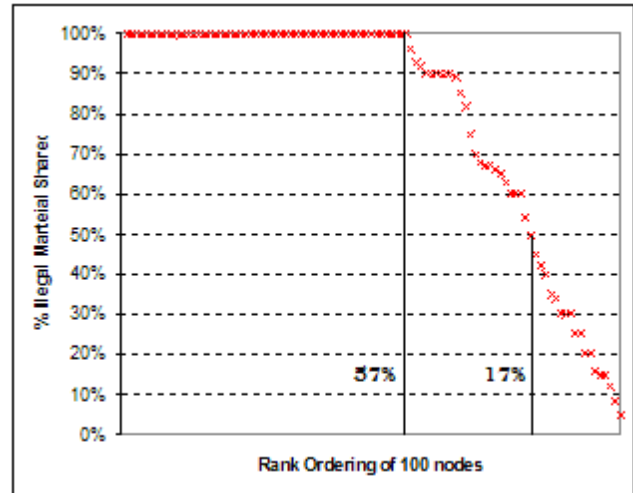
The disparity between the numbers of QUERY and QUERYHIT messages is primarily due to the fact that QUERYHIT messages refer to multiple files which may have matched a single received QUERY.

#### 4.2 Is this activity the result of a deviant sub-community?

To assess whether or not individuals who share illegal pornography form a deviant sub-community within the wider Gnutella community, we first produced a ranked list of the top 20 pornography-related search terms. From this list, we identified peers who responded with QUERYHITS on our selected dates (the Saturdays within our sampling period). This yielded a list of peers which could reasonably be assumed to be distributing illegal material.

We then selected 100 hosts at random from the above set and determined whether or not they participated in sharing other, legal, material. Figure 2 shows the proportion of illegal material that was being served by these 100 peers over the one month period. As can be seen, the majority of peers that share illegal pornography (57%) share no other

material whatsoever, while only 17% share less than 50% of illegal materials. Table 4 shows the in-between points.



**Figure 2: Percentage of illegal material shared**

**Table 4: Percentage of illegal material shared**

Illegal QUERYHITS	Total % of peers sharing
100%	57%
> 90%	66%
> 50%	83%
> 25%	91%

Unfortunately, it is not possible to associate QUERY traffic with their originating peer in the same way that it is with QUERYHIT traffic [12]; and therefore it is not possible to ascertain whether the peers serving illegal material are the same as those generating QUERY messages searching for this material.

#### 5. Discussion

We found that 1.6% of search traffic and 2.4% of response traffic was related to illegal pornography. While this is a small proportion, it remains significant, particularly given the large size of the Gnutella network. We also found strong evidence that those peers who share illegal pornography form a deviant sub community: 57% of peers that share such material share no other material, while only 17% share less than 50% illegal material.

This second finding lends clear support to those socio-psychological theories that emphasise the importance of group-specific social norms. Individuals do *not* automatically engage in 'disinhibited' behaviour simply because the opportunity arises; rather, such behaviour is

only facilitated in members of a group for whom it is already normative.

Our findings also have significant implications for the debate over the legality and future survival of P2P networks. As mentioned in section 1, there is a growing trend of legal action targeting the online distribution of illegal pornography. This, together with the significant amount of this material available on P2P networks, means that P2P file-sharing is likely to be increasingly targeted by such activities.

However, our research has shown that those responsible for the distribution of the illegal material are a small and separate sub-community. Our findings suggest that no action need be taken with regard to P2P networks as a whole *if* it possible to effectively target this sub-community without encroaching on the wider P2P community [21]. Furthermore, recent research suggests that a significant number of users are migrating to P2P networks that are harder to police [26]. Legal attacks such as [9] may simply accelerate this process, forcing the deviant sub-community onto networks where enforcement may be more difficult or even impossible.

## 6. A Suggested Response

We believe that there is considerable potential in developing mechanisms through which P2P file-sharing communities can *police themselves*. This will leave untouched the (legal) activities of the majority while subverting the illegal activities of the deviant minority.

In particular, we see potential in setting up a clear demarcation between the general Gnutella community and the deviant sub-community, thus encouraging in non-deviant individuals a feeling of identification with the identity ‘Gnutella user’ in explicit opposition to ‘deviant Gnutella user’ [22]. On the basis of the social psychology work reviewed above, this would be predicted to increase the likelihood of general users taking collective responsibility for policing the network and encouraging pro-social behaviour.

More concretely, such an approach could be implemented on Gnutella (or similar P2P networks) by allowing each peer to locally specify a list of search terms it does not approve of. Then, whenever a peer receives a corresponding QUERY message, the message could be routed differentially. If the topic is judged particularly offensive, the QUERY could be discarded entirely; or in the case of material that is considered only somewhat objectionable, the ‘time-to-live’ of the message could be reduced, thus limiting its propagation through the network. As search is performed by flooding and multiple potential paths usually exist between a searching peer and peers with matching resources, whenever a node discards, or limits the

propagation of a QUERY message, the originating peer will likely still receive a response via another route. However, the emergent result of a significant number of peers implementing such measures would be to limit the usefulness of the network for deviant users, while files which most users find acceptable would remain available (in proportion to the extent to which the P2P community approves of them).

Furthermore, as the deviant sub-community does not typically contribute to the network in other ways, constraining this community would not have a detrimental effect on the greater Gnutella community. On the negative side, such a system would be vulnerable to a Sybil Attack [28] (i.e. an attack in which a large number of malicious entities are deployed throughout a network, behaving in such a way as to cause undesirable emergent effects). This, however, is already true of the Gnutella network itself. Furthermore, for such an attack to be effective on a network the size of Gnutella a phenomenal amount of resources would be required.

## 7. Future Work

We believe that there are significant benefits to be gained from extending our study, both over a longer time period and by increasing the depth of analysis. First, the present study provides only a ‘snapshot’ of the current situation. Extending it over a longer period, e.g. 12 months, would expose any underlying trends that may have been missed. More interestingly, though, there are other phenomena that a more longitudinal study might illuminate—for example, do high-profile prosecutions [10] or public awareness campaigns [11] actually reduce the level of illegal pornography being shared? Such questions may prove important in determining effective law enforcement approaches. It would also be possible to extend our study to obtain more information about the composition of the deviant sub-community. For example, does the volume of material being shared by a peer relate to its geographical location? Do cultural attitudes and the laws of a peer’s host-country play a significant role in shaping online behaviour? Anonymous file-sharing networks are an ideal environment in which to evaluate CMC-related questions of this type.

## 7. References

[01] Napster. <http://www.napster.com>

[02] Kazaa. <http://www.kazaa.com>

[03] eDonkey. <http://www.edonkey.com>

[04] Limewire. <http://www.limewire.com>

- [05] "Dependability Properties of P2P Architectures", Walkerdine, J., Melville, L., Sommerville, I. Proc. 2nd IEEE International Conference on Peer-to-Peer computing (P2P'02), Linköping, Sweden, September, 2003.
- [06] The Recording Industry Association of America (RIAA) <http://www.riaa.com>.
- [07] "An Architecture for Peer-to-Peer Economies", Strulo, B., Smith, A., Farr, J., Proc. 3rd IEEE International Conference on Peer-to-Peer computing (P2P'03). Linköping, Sweden, September, 2003.
- [08] "Freenet: a Distributed Anonymous Information Storage and Retrieval System", Clark, I., Sandberg, O., Wiley, B., Hong, T., International Workshop on Design Issues in Anonymity and Unobservability, LNCS 2009, Springer, New York, 2001.
- [09] "A Bill to Outlaw the Selling, Advertising, and Distribution of Peer-to-Peer (P2P) File-sharing Software", A proposed change to the California legal code, US Senate, January 2005, [http://info.sen.ca.gov/pub/bill/sen/sb\\_0051-0100/sb\\_96\\_bill\\_20050114\\_introduced.html](http://info.sen.ca.gov/pub/bill/sen/sb_0051-0100/sb_96_bill_20050114_introduced.html).
- [10] "Hundreds targeted over child porn", BBC News, March 05 <http://news.bbc.co.uk/1/hi/world/europe/4354573.stm>
- [11] The Internet Watch Foundation <http://www.iwf.org.uk/>
- [12] "The Gnutella Protocol Specification v0.4". [http://www9.limewire.com/developer/gnutella\\_protocol\\_0.4.pdf](http://www9.limewire.com/developer/gnutella_protocol_0.4.pdf), 2000.
- [13] "Identifiability and Self-Presentation: Computer-Mediated Communication and Intergroup Interaction", Douglas, K. M., McGarty, C., British Journal of Social Psychology, 40, 399-416, 2001.
- [14] "Causes and Implications of Disinhibited Behavior on the Internet", Joinson, A., in J. Gackenbach (Ed.), "Psychology and the Internet: Intrapersonal, Interpersonal, and Transpersonal Implications", London, Academic Press, 1998.
- [15] "A Criminological Internet "Sting": Experimental Evidence of Illegal and Deviant Visits to a Website Trap", Demetious, C., Silke, A., British Journal of Criminology, 43, 213-222, 2003.
- [16] "Social Identity, Normative Content and 'Deindividuation' in Computer-Mediated Groups", Postmes, T., Spears, R., Lea, M., in N. Ellemers, R. Spears, & B. Doosje (Eds.), "Social identity: Context, commitment, content", Oxford: Blackwell, 1999.
- [17] "JTella", <http://jtella.sourceforge.net/>.
- [18] "Lancaster's Peer-to-Peer Website", <http://polo.lancs.ac.uk/p2p>.
- [19] "The Gnutella Protocol Specification v0.6", <http://rfc-gnutella.sourceforge.net/>.
- [20] "Computer Pornography: A Comparative Study of the US and UK Obscenity Laws and Child Pornography Laws in Relation to the Internet", Akdeniz, Y., International Review of Law, Computers & Technology, 10, 235-263, 1996.
- [21] "Internet Content Regulation: UK Government and the Control of Internet Content", Akdeniz, Y., Computer Law & Security Report, 17, 303-317, 2001.
- [22] "Social Categorization and Intergroup Behaviour", Tajfel, H., Billig, M.G., Bundy, R.P., Flament, C., European Journal of Social Psychology, 1, 149-178, 1971.
- [23] "Improving Quality of Service on Gnutella", Hughes, D., Warren, I., Coulson, G., Technical Report COMP-005-2004, Computing Department, Lancaster University, 2004.
- [24] "Free Riding on Gnutella Revisited: the Bell Tolls?", Hughes, D., Coulson, G., Walkerdine J., IEEE Distributed Systems Online, June 2005.
- [25] "Free Riding on Gnutella", Adar, E., Huberman, B., First Monday, October 2000. [http://www.firstmonday.dk/issues/issue5\\_10](http://www.firstmonday.dk/issues/issue5_10).
- [26] "Is P2P Dying or Just Hiding?", Karagiannis, T., Broido, A., Brownlee, N., Faloutsos, M., Proc. Globecom 2004, Dallas, U.S., December 2004.
- [27] "Looking at the Server-Side of Peer-to-Peer Systems" Qiao Y., Lu D., Bustamante F. E., Dinda P., Proc. 7th Workshop on Languages, Compilers and Run-time Support for Scalable Systems, October 2004.
- [28] "The Sybil Attack" Douceur, J., Proc. 1<sup>st</sup> Internet Workshop on Peer-to-Peer Systems (IPTPS02), Cambridge, MA (USA), March 2002.

