

# The KaZaA Overlay: A Measurement Study

Jian Liang

Department of Computer and  
Information Science,  
Polytechnic University,  
Brooklyn, NY, USA 11201  
Email: jliang@cis.poly.edu

Rakesh Kumar

Department of Electrical and  
Computer Engineering,  
Polytechnic University,  
Brooklyn, NY, USA 11201  
Email: rkumar04@utopia.poly.edu

Keith W. Ross

Department of Computer and  
Information Science,  
Polytechnic University,  
Brooklyn, NY, USA 11201  
Email: ross@poly.edu

September 15, 2004

## Abstract

Both in terms of number of participating users and in traffic volume, KaZaA is one of the most important applications in the Internet today. Nevertheless, because KaZaA is proprietary and uses encryption, little is understood about KaZaA's overlay structure and dynamics, its messaging protocol, and its index management. We have built two measurement apparatus - the KaZaA Sniffing Platform and the KaZaA Probing Tool - to unravel many of the mysteries behind KaZaA. We deploy the apparatus to study KaZaA's overlay structure and dynamics, its neighbor selection, its use of dynamic port numbers to circumvent firewalls, and its index management. Although this study does not fully solve the KaZaA puzzle, it nevertheless leads to a coherent description of KaZaA and its overlay. Furthermore, we leverage the measurement results to set forth a number of key principles for the design of a successful unstructured P2P overlay. The measurement results and resulting design principles in this paper should be useful for future architects of P2P overlay networks as well as for engineers managing ISPs.

# 1 Introduction

On a typical day, KaZaA has more than 3 million active users sharing over 5,000 terabytes of content. On the University of Washington campus network in June 2002, KaZaA consumed approximately 37% of all TCP traffic, which was more than twice the Web traffic on the same campus at the same time [8]. With over 3 million satisfied users, KaZaA is significantly more popular than Napster or Gnutella ever was. Sandvine estimates that in the US 76% of P2P file sharing traffic is KaZaA/FastTrack traffic and only 8% is Gnutella traffic [23]. Clearly, both in terms of number of participating users and in traffic volume, KaZaA is one of the most important applications ever carried by the Internet. In fact, it can be argued that KaZaA has been so successful that any new proposal for a P2P file sharing system should be compared with the KaZaA benchmark. However, largely because KaZaA is a proprietary protocol which encrypts its signalling messages, little has been known to date about the specifics of KaZaA's overlay, the maintenance of the overlay, and the KaZaA signalling protocol.

In this paper we undertake a comprehensive measurement study of KaZaA's overlay structure and dynamics, its neighbor selection, its use of dynamic port numbers to circumvent firewalls, and its index management. Although this study does not fully solve the KaZaA puzzle, it nevertheless leads to a coherent description of KaZaA and its overlay, while providing many new insights about the details of KaZaA.

To unravel the mysteries of the KaZaA overlay, we developed two measurement apparatus: the KaZaA Sniffing Platform and the KaZaA Probing Tool. The **KaZaA Sniffing Platform** is a set of KaZaA nodes that are forced to interconnect in a controlled manner with one another, while one node is also connected to hundreds of platform-external KaZaA nodes. The KaZaA Sniffing Platform collects KaZaA signalling traffic, from which we can draw conclusions about the structure and dynamics of the KaZaA overlay. The **KaZaA Probing Tool** establishes a TCP connection with any supplied KaZaA node, handshakes with that node, and sends and receives arbitrary encrypted KaZaA messages with the node. It is used for analyzing node availabilities and KaZaA neighbor selection. Both of these apparatus consume limited resources. One of the contributions of this paper is to show how it is possible to obtain extensive overlay information of a large-scale overlay application with a low-cost measurement infrastructure.

We use these tools to obtain insight into the following questions:

- It is well-known that the KaZaA overlay is organized in a two-tier hierarchy consisting of Super Nodes (SNs) in the upper tier and Ordinary Nodes (ONs) in the lower tier. But how many children ONs does a typical SN support? What fraction of the peers in KaZaA are SNs? Are the SNs densely interconnected or sparsely interconnected?

- How long are ON-to-SN connections in the overlay? How long are SN-to-SN connections in the overlay? What is the typical lifetime of a SN?
- How does an ON discover candidate SNs for parenting? Once it has a set of candidate SNs, how does it choose a particular parent among them? In choosing the parent, does it take locality or SN workload into account?
- By allowing peers (ONs and SNs) to select their own server port numbers, KaZaA is more difficult to block with firewalls and NATs. How does KaZaA manage the server port numbers? What fraction of KaZaA nodes are behind NATs?
- What are the characteristics of the protocol that peers use to establish overlay links among themselves?
- How is the file index (relating each file copy to an IP address and port number) organized among the SNs?

In addition to providing novel insights into a remarkably successful P2P system, we leverage our measurement results to set forth a number of key principles for the design of an unstructured P2P overlay. As we'll discuss in Section 5 these principles, including distributed design, exploiting heterogeneity, load balancing, locality, connection shuffling, and firewall/NAT circumvention.

This paper should not only be of interest to P2P designers, but also to engineers at upper- and lower-tier ISPs, who are interested in acquiring a thorough understanding of P2P overlays and traffic. Because P2P file sharing systems can generate vast quantities of traffic, networking engineers, who dimension the network and introduce content distribution devices such as caches, need a basic understanding of how major P2P file sharing systems operate. Although there has been recent work in analyzing the file-sharing workload in KaZaA [8] and [18], to our knowledge we are the first to undertake a comprehensive study of a hierarchical unstructured overlay for a P2P system.

The paper focuses on the KaZaA overlay network and index management. It addresses neither KaZaA's downloading protocol (for example, KaZaA's parallel downloading and request queuing) nor its incentive scheme for encouraging uploading. The paper is complementary to [8] and [18], which focus on KaZaA file-sharing traffic. It is also complementary to a recent measurement study on pollution in P2P file sharing systems [19].

This paper is organized as follows. Section 2 provides an overview of KaZaA. Section 3 describes the measurement apparatus. Section 4 presents our measurement results. Section 5 sets forth basic design principles for unstructured P2P file sharing applications. Section 6 surveys related work. Finally, Section 7 summarizes our findings and concludes.

## 2 Overview of the KaZaA

KaZaA Web site [14] provides a rudimentary description of how KaZaA works. Moreover, various (and often obscure) articles, Web sites, and message boards provide additional scraps of information. In this section we collect and unify this publicly available information. The goal of this section is to *(i)* organize this obscure information in a digestible form for the P2P research community and *(ii)* present a broad-brush picture of KaZaA and its overlay. In the subsequent sections we'll describe our own measurement contributions.

KaZaA peers differ in availability, bandwidth connectivity, CPU power, and NATed access. KaZaA was one of the first P2P systems to exploit this heterogeneity by organizing the peers into two classes, Super Nodes (SNs) and Ordinary Nodes (ONs). SNs are generally more powerful in terms of connectivity, bandwidth, processing, and non-NATed accessibility. As we'll shortly describe, SNs also have greater responsibilities. As shown in Figure 1, each ON has a parent SN. When an ON launches the KaZaA application, the ON chooses a parent SN, maintains a semi-permanent TCP connection with its parent SN, and uploads to this SN the metadata for the files it is sharing.

As with most other P2P file sharing systems, KaZaA maintains a file index that maps file identifiers to the IP addresses. This file index is distributed across the SNs. In particular, each SN maintains a local index for all of its children ONs, so that each SN is similar to a (mini) Napster hub. But in contrast with Napster, a SN is not a dedicated server; instead, it is typically a peer belonging to an individual user.

We know from [16] that for each file an ON is sharing, the metadata that the ON uploads to its parent SN includes: the **file name**, the **file size**, the **ContentHash**, and the **file descriptors** (for example, artist name, album name, and text entered by users). The file descriptors are used for keyword matches during querying. The ContentHash plays an important role in the KaZaA architecture. KaZaA hashes every file to a hash signature, which becomes the ContentHash of the file. The ContentHash is the only identifier used to identify a file in an HTTP download request. If a download from a specific peer fails, the ContentHash enables the KaZaA client to automatically search for the specific file, without issuing a new keyword query.

When a user wants to locate files, the user's ON sends a query with keywords over the TCP connection to its parent SN. For each match in its database, the SN returns the IP address, server port number, and metadata corresponding to the match. As shown in Figure 1, each SN also maintains long-lived TCP connections with other SNs, creating an overlay network among the SNs. When a SN receives a query, it may forward the query to one or more of the SNs to which it is connected. A given query will in general visit a small subset of the SNs, and hence will obtain the metadata information of a small subset of all the ONs.

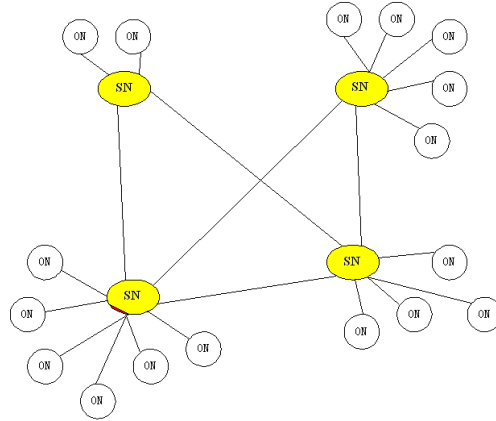


Figure 1: *KaZaA's two-tier overlay network.*

The FastTrack File Format project [16] has determined the syntax and semantics of KaZaA system files. From this project, we know that a KaZaA peer has the following software components:

1. The KaZaA Media Desktop (KMD).
2. Software environment information stored in the Windows Registry. Included in this environment information is a list of up to 200 SNs, which we refer to as the **SN list cache**. For each the 200 SNs in cache, the list includes a number of attributes including the SN IP address and port number.
3. DBB files, with each DBB file containing metadata for the files that the peer is willing to share. An active KaZaA process permanently monitors the local folders that are shared; file add, delete, is reflected in the DBB file.
4. DAT files, with each file containing a partially downloaded file. A DAT file grows in size as more data is retrieved. Once all the file data is retrieved, the DAT file is renamed to the original file which was intended to be downloaded.

Each KaZaA peer exchanges four different types of TCP traffic with other peers in the network:

1. Signaling traffic, which includes handshaking traffic for connection establishment between peers; metadata extracted from the DBB files, uploaded from ONs to SNs; supernode lists; and queries and replies. All signaling traffic is encrypted.
2. File transfer traffic (e.g., MP3s, videos, etc.) transferred directly among the peers without passing through intermediate SNs. File transfers are not encrypted and are sent within HTTP messages.
3. Commercial advertisements, sent over HTTP.

#### 4. Instant messaging traffic, encoded as Base64.

The KaZaA ON-SN and SN-SN signalling messages are encrypted. The impressive giFT project [1] has reverse-engineered KaZaA’s encryption algorithms, so that users of giFT-FastTrack can search and download files from the KaZaA network. Some of our own measurement tools incorporate the encryption/decryption code provided through the giFT project. The Sig2dat tool project [29] makes available a tool for obtaining the KaZaA ContentHash of any file. This tool is increasingly being used by KaZaA users, who post file names and corresponding ContentHash values on Web sites and message boards. This helps in countering pollution attacks, wherein bogus files are intentionally placed in the network by competing interests [22] [19].

Many users today use KaZaA-Lite [15], an unofficial copy of KMD, rather than the KaZaA client(KMD) distributed by Sharman. Each KaZaA-Lite client emulates Sharman’s KMD and participates in the KaZaA network. During the search process, a KaZaA-Lite ON first sends its query to the SN to which it is connected. We have learned from our own measurement work that after receiving all the replies from its parent SN, the ON disconnects and connects with a new SN, and re-sends the query to the new SN. During a specific search, the ON may connect to many SNs. The ordinary node typically maintains the TCP connection with the last SN in the sequence of connections, until another search is performed. Our measurement work has determined that during each hop, the ON re-sends its meta-data to the new SN, and the previous SN removes the ONs meta-data.

The use of the terms “KaZaA protocol” and “KaZaA overlay” is convenient but somewhat incorrect. More precisely, KaZaA is just one of several clients that use the FastTrack protocol and participate in the FastTrack overlay. In addition to KaZaA, Grokster and iMesh are two other clients that currently participate in the FastTrack overlay network. All three clients use the same protocol as KaZaA. Each KaZaA-Lite client also emulates FastTrack protocol and participates in the FastTrack network. When we say “KaZaA” in this paper, we are actually referring to the FastTrack network and all of its clients.

### 3 Measurement Apparatus

This section describes the measuring apparatus we use to conduct our measurements on the KaZaA’s overlay.

#### 3.1 The KaZaA Sniffing Platform

As shown in Figure 2, we built a sniffing platform consisting of three workstations, each with a KMD version 2.0 client installed. We patiently waited until KaZaA promoted

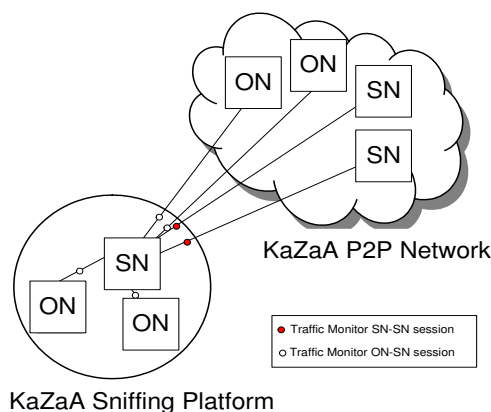


Figure 2: *The KaZaA Sniffing Platform.* The SN typically has hundreds of connections to external peers.

one of the three nodes to a SN. At startup all three workstations functioned as ONs in the KaZaA network. These workstations enjoyed high bandwidth connectivity with sufficient hardware resources. When one of the workstations was promoted to a SN, we manipulated the Windows Registries in the other two Platform ONs in a manner that forced them to adopt the Platform SN as their parent. (How we did this will soon become clear.)

As shown in Figure 2, the Platform SN also connects to platform-external KaZaA ONs and SNs. We deploy software traffic monitors around the Platform SN to capture all of the inbound and outbound signalling traffic. We then do an offline traffic analysis based on our understanding of the KaZaA signalling protocol.

The KaZaA Sniffing Platform was installed in two different subnets: one connected to Polytechnic University campus network; the other connected to a residential cable access network. In this manner we can take snapshots of the KaZaA network from two entirely different types of network access. These locations were chosen because they are representative of type of subnets that the majority of KaZaA peers connect through, i.e., campus networks and residential access.

By deploying the KaZaA Sniffing Platform in conjunction with the encryption/decryption tools from the giFT project [1], we have been able to determine that KaZaA nodes frequently exchange with each other lists of SNs. In particular, when an ON connects with a parent SN, the SN immediately pushes to the ON a **SN refresh list**, which consists of the IP addresses, port numbers and workload values of up to 200 SNs. The first entry in the SN refresh list is the parent SN that is sending the list. When an ON receives a SN refresh list from its parent SN, the ON will typically purge some of the

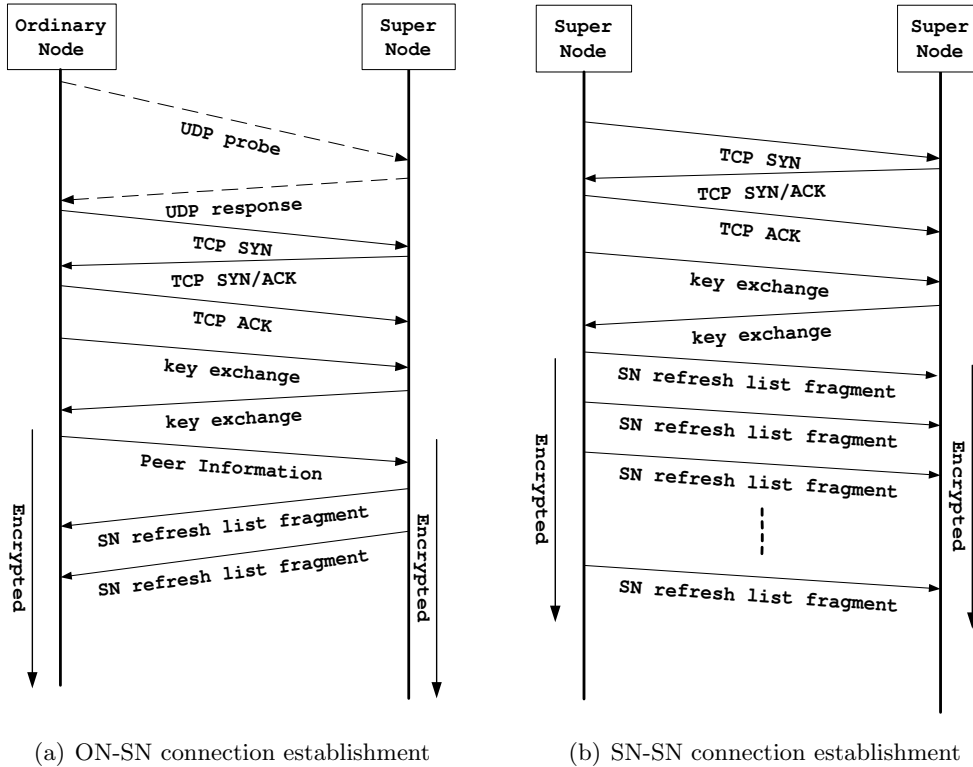


Figure 3: *Connection establishment protocol between peers in the KaZaA overlay.*

entries from its SN list cache and add entries sent by the parent SN. Neighboring SNs in the overlay also exchange SN refresh lists. By frequently exchanging SN refresh lists, nodes maintain up-to-date lists of active SNs.

### 3.2 Overlay Probing Tool

By deploying the KaZaA Sniffing Platform in conjunction with the encryption/decryption tools from the giFT project [1], we have been able to determine the sequence and semantics of many of the KaZaA messages, as well as the sequence of events that ensue during overlay link establishment between ON and SN and between SN and SN.

When a peer launches the KaZaA client, the first task of the client is to choose a parent SN and establish an overlay link (i.e., TCP connection) with it. To this end, we have discovered that the following steps are taken:

- As described in Section 2, the ON has a SN list cache. The ON chooses several (typically 5) candidate SNs from the list and probes the candidates by sending one UDP packet to each candidate. The ON then receives UDP responses from a subset of these candidates.
- To each SN from which it receives a UDP response, the ON attempts to establish



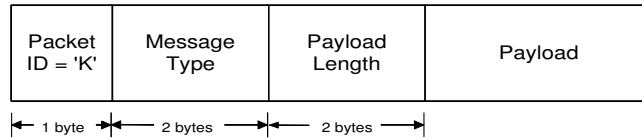


Figure 4: *KaZaA's signalling message format*. 5 bytes of header with variable length payload.

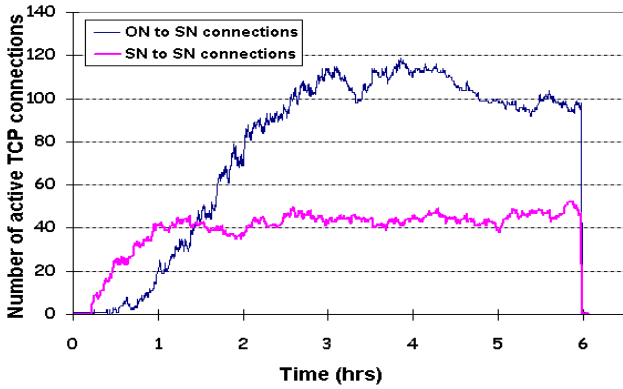
a TCP connection. For each such connection, the SN and ON will exchange encryption key material, the ON will send peer information, and the SN will send a SN refresh list. Included in the peer information is the local IP address, service port number and username. The ON may be behind a NAT in which case the local IP address is a private address and different from the NAT's address.

- The ON will then select one of the SNs and disconnect from the other SNs. The one remaining SN becomes the ON's parent SN.
- The ON can then send query messages to its chosen SN.

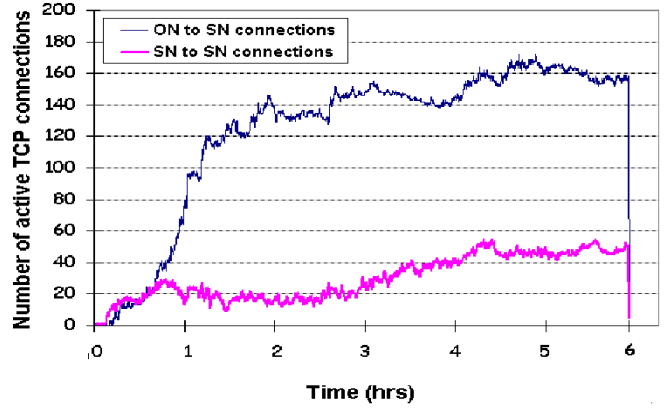
The timing diagram for both the ON-SN and SN-SN overlay connection establishment is shown in Figure 3. Note that the procedure for establishing a SN-SN overlay link differs from that of establishing ON-SN overlay link.

We also determined the structure of signalling messages exchanged between KaZaA peers. This structure is presented in Figure 4. Each message begins with the identifier "K", which is then followed by a message-type field (two bytes), a payload-length field (two bytes), and the payload itself.

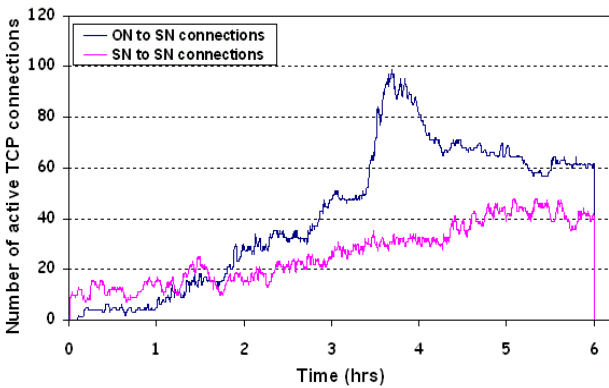
Based on our understanding of the KaZaA signaling protocol and message structure, we have developed the KaZaA Overlay Probing Tool. This tool can fully emulate the behavior of the KMD client for initiating new connections and exchanging signalling messages with other KaZaA nodes. We use this tool to (1) probe whether any arbitrarily specified SN in the overlay is alive and (2) to retrieve the SN refresh list sent from probed SN, and (3) obtain the workload of the probed SN. As discussed in Section 4, we use this list to test our hypothesis of locality awareness as a selection criterion in forming new SN-ON overlay links in KaZaA.



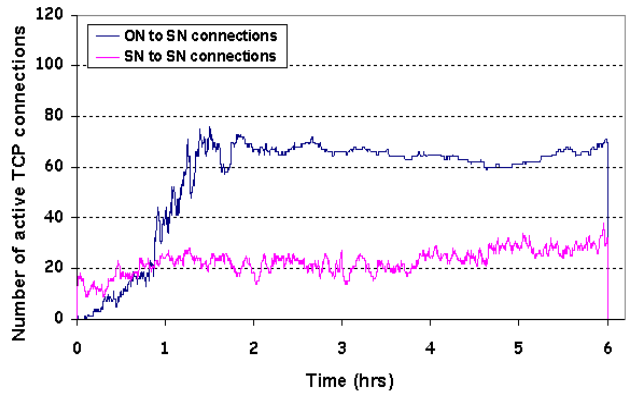
(a) Polytechnic campus - session evolution, Aug. 22, 2003



(b) Polytechnic campus - session evolution, Oct 24, 2003



(c) Residential Cable Modem - session evolution, Aug. 23, 2003



(d) Residential Cable Modem - session evolution, Oct 25, 2003

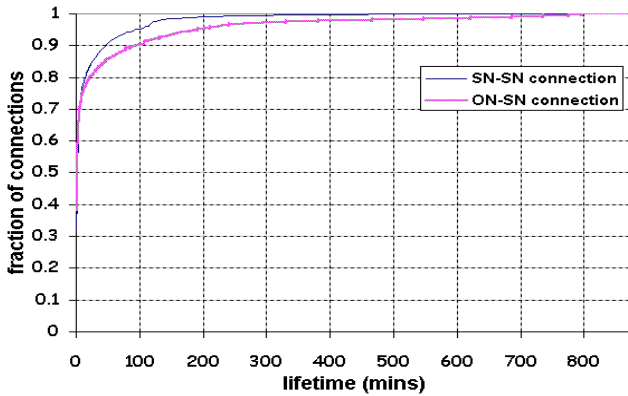
Figure 5: *Evolution of SN-SN and ON-SN connections with time.*

## 4 Measurement Results

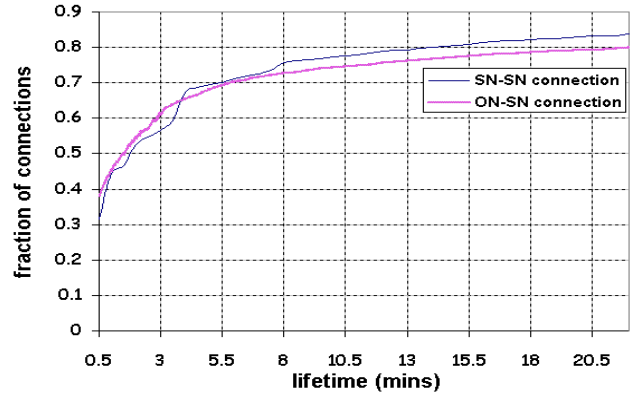
### 4.1 Overlay Structure and Dynamics

#### 4.1.1 Structural Properties of the Overlay

We first explored the degree of connectivity of a typical SN. Specifically, from the SN in the KaZaA Sniffing Platform, we study the number of simultaneous connections to platform-external ONs and SNs. The Platform SN does not offer any files for sharing. From the moment when a Platform node is promoted to a SN, we monitor the number TCP connections emanating from the SN. We did this experiment in the two environments – Polytechnic campus and broadband residential access network – over different time periods. Figure 5 shows that in every case, the number of connections begins at one and climbs to a threshold, around which it subsequently vacillates. For the number of simultaneous SN-SN connections, this threshold is almost always in the 40-50 range.



(a) full duration plot



(b) close-up of the plot

Figure 6: *Connection lifetime distributions.* The plots on the left are for the full duration of trace. The plots on the right are the corresponding close-ups for shorter duration which show the distribution more clearly for connections of lower lifetimes.

For the number of simultaneous SN-ON connections, depending on the day, this threshold is in the 100-160 connection range for the Platform SN on the Polytechnic campus and in the 55-70 range for the Platform SN in the residential access network. Since on a typical day there are roughly 3 million peers, we therefore speculate that there are on the order of 25,000-40,000 SNs in the KaZaA overlay, with the number varying with the time of day. This claim has also been corroborated by a complimentary measurement study reported in [19]. Combining these estimates, we conclude that the SN-SN overlay network is very sparsely connected, with each SN connected to about 0.1% of other SNs in the overlay.

#### 4.1.2 Overlay Dynamics

Our measurement study has determined the KaZaA overlay is highly dynamic. Although the number of simultaneous connections vacillates around a threshold, as observed in Figure 5, the individual connections change frequently.

Using the KaZaA Sniffing Platform we performed measurements on the duration of ON-SN TCP connections and SN-SN TCP connections spread over seven days. We show one such representative measurement done on Oct. 24, 2003 in Figure 6 at Polytechnic University. Here we monitored over a period of 12 hours a total of 5,206 ON connections and 3,850 SN with our Platform SN. We plot the distribution of connection lifetime for these two types of TCP connections in Figure 6. The average durations of ON-SN connections and SN-SN connections are 34 mins and 11 mins, respectively. We also observe that a remarkable 32% of the SN-SN connections and 38% of the ON-SN connections lasted for less than 30 seconds. Among connections that last for at least

30 seconds, the average durations of ON-SN connections and SN-SN connections are 57 mins and 23 mins, respectively.

We attribute the large number of short lifetime ON-SN connections to two factors. First, as discussed in Section 3, at startup an ON probes candidate SNs listed in its SN refresh list with UDP packets for possible connections. The ON then initiates simultaneous TCP connections with the available SNs in its SN refresh list. Out of these successful connections, the ON selects one SN as the final choice and it disconnects from other SNs. Hence the ON-SN connection establishment process generates many short-lived ON-SN connections. A second reason for short-lived ON-SN connections is that many ONs are KaZaA-Lite clients. As described in the Introduction, KaZaA-Lite clients hop supernodes during the query process. Each such hop generates a short-lived connection.

We conjecture the short lifetime of SN-SN connections is due to (1) SNs searching for other SNs with currently small workloads, (2) long-term connection shuffling, to allow users to query a large set of SNs over long time scales and (3) at times, SNs in the overlay connect to each other just for the purpose of exchanging SN lists. The shuffling of neighbor SN-SN connections allows a larger range of the network to be explored, for example, when searching takes place over hours or days for download lists and fragments of large files (such as movies).

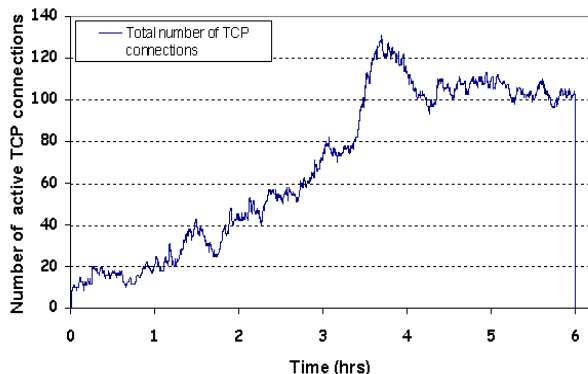
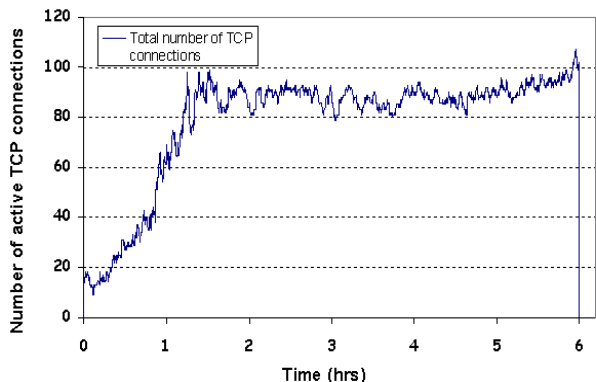
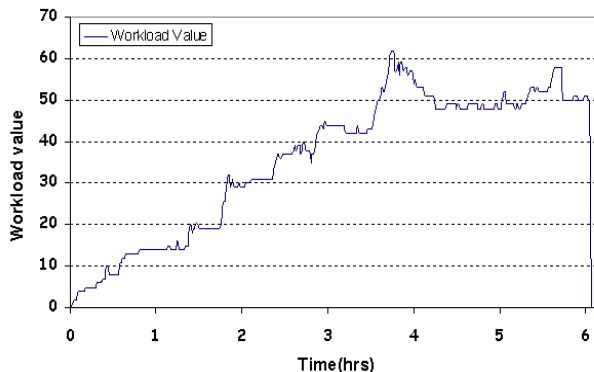
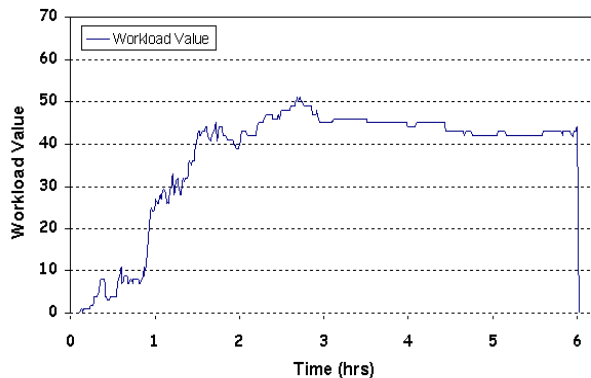
## 4.2 Parent Selection

One crucial characteristic of a two-tier overlay is the criteria that an ON employs to select a parent SN. In this section we describe results from our experiments on determining the prominent factors influencing the choice of a parent SN.

As discussed in Section 3, when an ON establishes an overlay link with a parent SN, the ON receives a list of 200 SNs from the parent SN. As already discussed, this list is a subset of all the SNs that the parent maintains in a cache. The specific SNs included in this list restricts the ON's future choices about which SNs to connect to; this in turn affects the overlay topology. Based on our measurements, we hypothesize that KaZaA peers mainly use two criteria for ON-to-SN and SN-to-SN neighbor selection, namely, workload and Locality.

### 4.2.1 Workload

Recall that each ON maintains a SN list Cache in the Windows Registry. For each SN in the list, this list includes four attributes: SN IP address, SN port number, SN workload, and timestamp. We now investigate the how the SN workload values in the list influence parent selection. The exact definition of SN workload is unknown, but as shown in Figure 7 there is a clear correlation between the workload and number



(a) ON-SN connection establishment

(b) SN-SN connection establishment

Figure 7: Correlation between the workload value for SN and the number of TCP connections to the SN in consideration.

of connections that the SN is supporting. For our Platform SN, Figure 7 plots the evolution of the number of on-going TCP connections as well as the evolution of the Platform SN workload. The similarity of these evolution plots is remarkable.

Recall that during startup, the ON chooses a subset of SNs (usually five SNs) from the SN cache list as candidates for a parent SN. We hypothesize that an ON takes SN workload into account when choosing the candidates. To test this hypothesis, we force an ON in the Platform to connect to and then disconnect from the overlay. Every time the ON attempts to connect, it chooses a subset of SNs from its SN list cache in the Windows Registry and attempts connections as discussed in Section 3. We sniff this signalling traffic and determine this SN subset; we also extract the 200 SNs in the SN cache list in the ON from the Windows Registry. We then calculate the average workload of the chosen subset of SNs and the average workload of the SNs in the SN cache list. We repeat this measurement every half hour. Figure 10 presents the results. Clearly, the KaZaA client displays a marked preference for SNs with low values for the

workload.

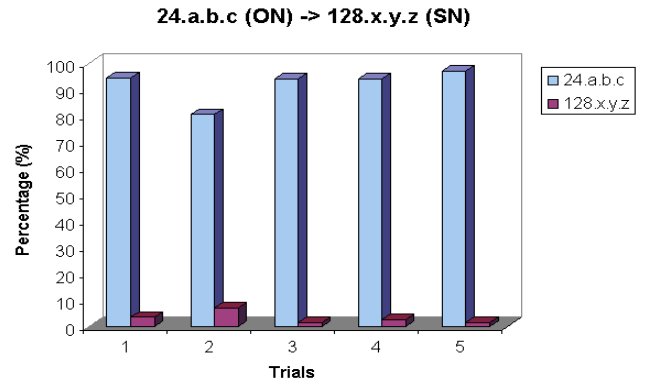
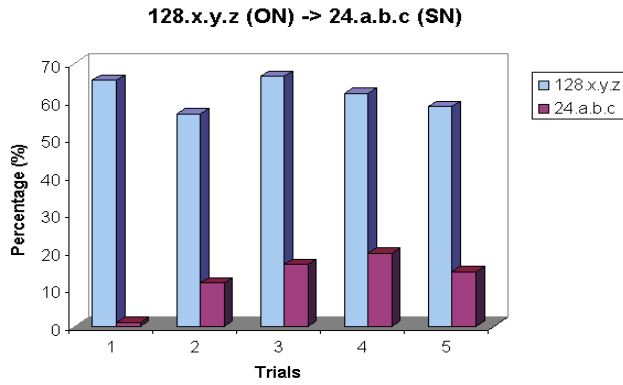
### 4.2.2 Locality

We hypothesize that an ON takes locality into account when selecting a parent SN, and that SNs take locality into account when selecting neighboring SNs in the overlay. We have performed two experiments to investigate locality. The first experiment uses Ping to measure the round-trip time (RTT) from the Platform SN to the platform-external ONs and SNs to which it connects. Figure 9 shows the distribution of these RTTs. We observe that about 60% of the SN-SN connections have RTTs less than 50 msec. It is instructive to compare these values with some typical RTT values for IP datagrams in the Internet. Transatlantic traffic between U.S East Coast and Europe experiences a RTT of about 100 ms, while the RTT for traffic between North America and Asia is approximately 180 ms [7]. Also it can be observed that almost 40% of the ON-SN connections have RTTs less than 5 msec, with the the other 60% having RTTs more or less uniformly distributed over hundreds of milliseconds.

The second locality experiment is based on IP prefixes. This measurement is made possible with the use of the KaZaA Probing Tool, discussed in Section 3, which can connect to a pre-specified SN and retrieve its SN refresh list. In this experiment we install the KaZaA probing tool in two nodes in the US, one with a 128/8 prefix and the other with a 24/8 prefix. From each of these nodes we connect to five SNs with prefix 128/8, to another five SNs with prefix 24/8 and again to another five SNs with prefix 213/8. The first two groups of these SNs are in the US and the third group is in Sweden.

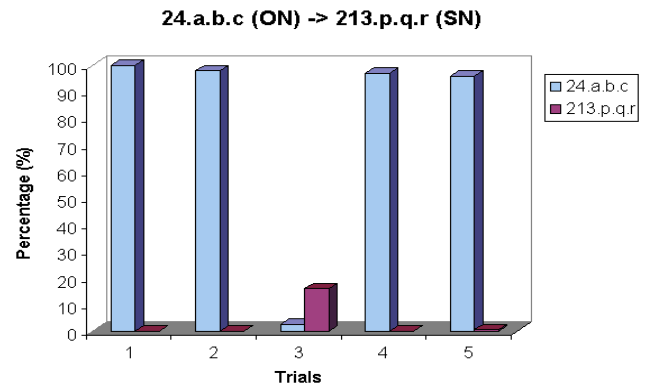
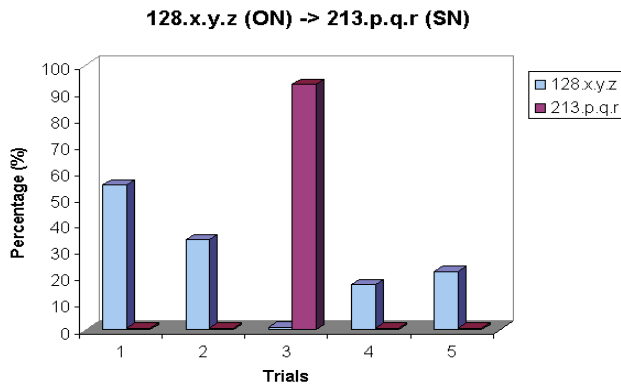
We can see from Figure 4.2.2 that a high percentage of SNs in the SN lists have similar IP prefixes as the child ON. Thus, when a SN prepares a SN list for an ON, it appears that the SN includes in the list SNs that topologically close to the ON. The percentage is slightly less in Figure 8(c). This is likely because the SNs based in European countries tend to have less knowledge of SNs based in the U.S. and thus are not able to include as many SNs matching the IP prefix of the connecting child ON. We have also discovered from other experiments that 24/8 prefix has a very high density of SNs. This is the reason we see even the Swedish SN supplying a high percentage of SNs matching the IP prefix connecting ON.

Thus, from our two locality experiments, involving RTT and prefix matching, it appears that KaZaA takes locality into account as it dynamically constructs the overlay network. However, it is not clear whether KaZaA uses RTTs, prefixes, or some combination of the two. Although locality helps to confine the KaZaA traffic within nearby ASes, it also means that the search results tend to be localized.



(a) The child ON IP address is from 128/8 and of parent SN is from 24/8

(b) The child ON IP address is from 24/8 and of parent SN is from 128/8



(c) The child ON IP address is from 128/8 and parent SNs are from 213/8

(d) The child ON IP address is from 24/8 and parent SNs are from 213/8

Figure 8: *IP prefix locality*. Shows the percentage of SNs in the SN list having common IP prefix with the child ON and the parent SN.

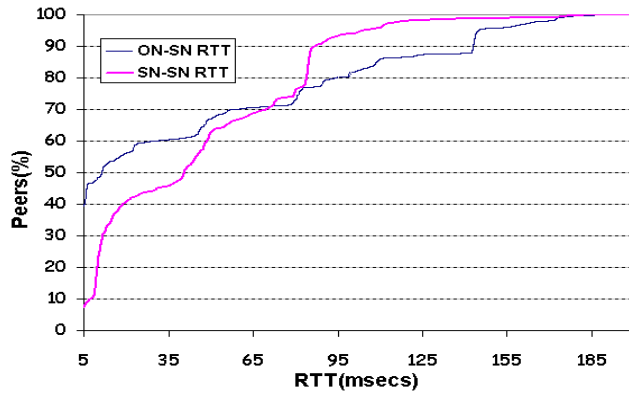


Figure 9: *Round-Trip Time measurement.* CDF of RTTs between supernode neighbors. The X-axis is the RTT values and the Y-axis is the corresponding percentage of neighbors that have a RTT less than that value.

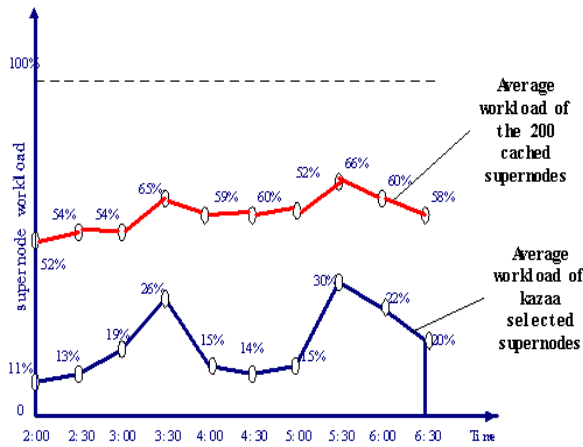


Figure 10: *Preference for Supernodes with less Workload.* The top curve shows the average of the workloads listed for each entry in the SN list cache. The bottom curve shows the average of the workloads of the SN chosen by the KaZaA client to be probed.



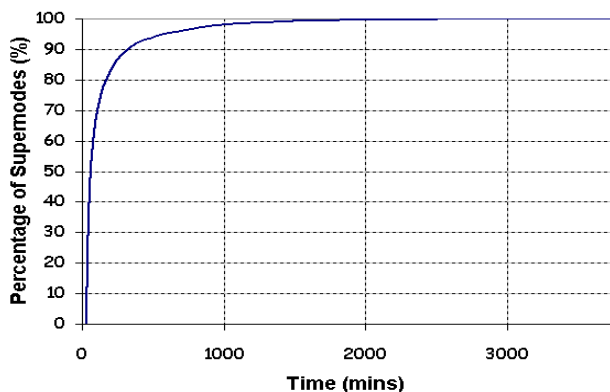


Figure 11: *CDF as percentage of lifetime of supernodes.* 965 supernodes were monitored over 65hrs for their lifetime. The average lifetime was found to be around 149 mins (2.5 hrs)

### 4.3 Supernode Lifetime

Figure 11 shows distribution of lifetime for 965 unique SNs, monitored over a period of 65 hours. The data for this figure was obtained with the KaZaA Probing Tool (at Polytechnic University), which repeatedly sent probing packets to each of the 965 SNs every five minutes. From our experiment we determined the average lifetime of a supernode in the KaZaA overlay to be 149 mins ( $\approx 2.5$ hours).

### 4.4 Firewall Evasion and NAT Circumvention

Earlier KMD clients used the default port number 1214. Thus, with the earlier versions, whenever Peer A wanted to connect with Peer B, it connected to port 1214. With this fixed port number, administrators of campus and corporate networks could easily configure their firewalls to prevent internal KaZaA peers from connecting to external KaZaA peers. Later versions (KMD v2.0+ and KaZaA-lite) employ dynamic port numbers to evade firewalls. Here, each peer chooses its own random port number and advertises it to other peers in the overlay. Specifically, when an ON establishes a link with a parent SN, it informs the parent SN of its port number. Furthermore, the SN refresh lists sent among the peers also advertise the port numbers of the SNs. Based on the 19,637 SN addresses we collected, we found only 707 SNs (3.6%) use the default 1214 port. 18,887 SNs (96.3%) use the non-default port numbers from 1024 to 65535. 10 SNs even use the 80 port number. Since the KaZaA port numbers are dynamic, it is very difficult to block KaZaA connections, unless a very rigid filtering policy is employed at the firewall.

The reality of today's Internet is that a large fraction of peers reside behind NATs. In fact our measurements indicate that roughly 30% of the KaZaA peers are behind

NATs. This is problematic for P2P networks, where in principle, every peer should be capable of serving (i.e., uploading) files. In particular, if peer A wants to download a file from a NATed peer B, peer A cannot initiate a direct TCP connection to peer B and send a request for a file. KaZaA’s two-tier hierarchy provides a mechanism to partially solve this problem. In KaZaA, when peer A sees that peer B has a private NAT address, instead of sending a request directly to peer B, it sends the request to peer B’s parent SN. The parent SN then sends a message to peer B, indicating that it should initiate a connection directly back to peer A. With this connection in place, the file can be sent over the connection from peer B to peer A. This technique of using an intermediate peer which already has a TCP connection in place to the NATed peer is called **connection reversal** [17]. Our measurement studies have shown that KaZaA implements connection reversal.

## 4.5 Index Management

As described in Section 2, on joining the KaZaA network, an ON uploads metadata information contained in the DBB file to its parent SN. We did experiments to measure the distribution of the amount of this metadata uploaded to the Platform SN (at Polytechnic University) from the connected ON sessions. The trace combines experimental data from a total of 894 ON-SN sessions. Figure 12 shows the cumulative distribution function of the meta-data uploaded to the Platform SN with respect to the percentage of ON-SN connections responsible for it. It can be observed from the plot that 13% of the ON peers are responsible for over 80% of the meta-data uploaded. It is interesting to compare this data with results reported in [25], wherein on University of Washington campus, 8.6% of KaZaA peers were serving 80% of the requests.

Another important index management strategy KaZaA SNs employ is to purge meta-data of any child ON as soon as it disconnects from the parent. We verified this by having a child connect to a SN, upload meta data to the SN, and then later disconnect from the SN. When querying the SN, from a second ON, for the meta data, the second ON would get a response as long as the first ON remained connected to the SN. As soon as the first ON disconnected, the SN ceased to receive responses for its query.

Independently, we also attempted to observe whether SN-SN traffic between neighbors in KaZaA involves exchange of meta-data information collected from their respective child ONs. We then sniffed traffic between our Platform SN and other neighboring SNs of the KaZaA overlay, and no index information exchange was observed.

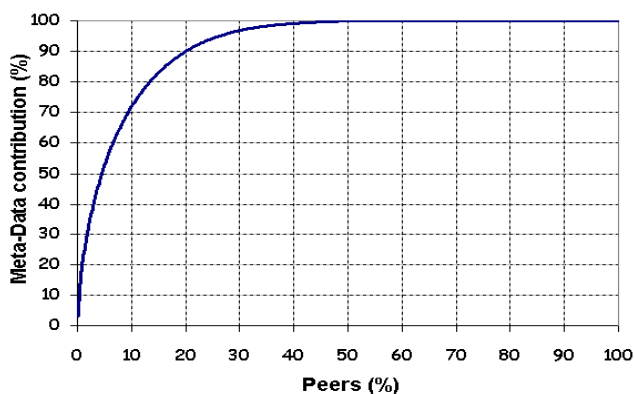


Figure 12: *CDF of amount of meta-data uploaded to a SN.*

## 5 Basic Design Principles for Unstructured P2P File Sharing Systems

KaZaA is one of the most successful large-scale P2P applications to date. Its success is partially due to many of the design decisions that were made for its overlay. We now leverage our measurement results to set forth a number of key principles for the design of an unstructured P2P overlay. These principles can serve as broad guidelines to advance the state of the art in design of P2P systems.

1. **Distributed Design:** Unlike Napster, KaZaA does not rely on infrastructure servers; essentially all of its nodes run on user peers. A distributed design has a number of advantages: no need for infrastructure equipment and maintenance; resilience to faults; and resilience to legal attacks (that is, KaZaA can not be brought down as Napster was by simply unplugging a central server).
2. **Exploiting Heterogeneity:** Peers differ in availability, bandwidth connectivity, CPU power, and NATed access. KaZaA was one of the first P2P systems to exploit this heterogeneity by organizing the peers into two classes, Super Nodes (SNs) and Ordinary Nodes (ONs). SNs are generally more powerful in terms of availability, bandwidth, processing, and non-NATed accessibility. The powerful peers should naturally bear more of the workload due to signaling and querying traffic. KaZaA assigns more responsibilities to the SNs. In particular, the SNs process, distribute, and respond to query traffic; and the SNs process index-maintenance traffic and overlay maintenance traffic.
3. **Load Balancing:** In the a two-layer hierarchical design such as KaZaA, the upper-tier nodes (i.e., the SNs) process the large majority of the signaling traffic (overlay-maintenance traffic and index-maintenance traffic) as well as the query

traffic. In order to not overwhelm any subset of SNs, care should be taken to balance the load of this traffic across all the upper-tier nodes. However, it is difficult to equally balance this traffic as peer behavior is unpredictable. To achieve approximate balance, the overlay can be designed so that each SN has roughly the same degree in the overlay (that is, has roughly the same number of TCP connections to ON and SN neighbors). From our measurement results in Section 4 we found that an ON selects a parent SN with a relatively small value of “workload”. Since we also showed there is a strong correlation between workload and the degree of the SN, KaZaA attempts to distribute the signaling and querying load across the SNs.

4. **Locality in Neighbor Selection:** Ideally, the neighbors in an overlay network should be close in terms of latency and network topology [4]. Constructing overlay links with short RTTs helps to reduce query/response delays. Constructing overlays so that neighboring peers (e.g., a SN and its children) are topological close helps to confine query and download traffic within an AS (or within nearby ASes) We saw in section 4.2.2 that locality in the form of a common IP prefix and short RTTs play in determining an ON’s parent SN as well as in the selection of SN-SN links. Although locality considerations help to confine query and download traffic to regions, we also note it may have an adverse impact on the content availability in the file-sharing system; for example, if a query from the US only reaches US SNs, which in turn only have US peers for children, then the query may not be able to locate obscure African content. Thus the advantages of locality has to be weighed against to the need for high content availability.
5. **Connection Shuffling:** In a P2P file sharing system, a user may repeatedly query for an obscure file over a period of days. By shuffling the links in the overlay, a larger set of SNs can be visited for an extended search period. Similarly, during a download, a peer server providing the download may disconnect in the middle of the download. Many P2P file sharing systems, including KaZaA, automatically and repeatedly search for a new copy of the file. Shuffling the links in the overlay helps the P2P system to find a replacement copy to complete the download. We discussed in section 4.1.2 that SNs engage in shuffling of neighbor connections. Thus, a design principle in building unstructured P2P file sharing systems with limited scope queries (such as KaZaA) is that of overlay shuffling, which broadens the scope of a query when repeatedly querying over long time periods.
6. **Efficient gossiping algorithms:** In a two-tier distributed P2P system, it is critical that the SNs learn about the other SNs in the network, so that they can shuffle connections as well as find new SNs when existing connections leave. Thus the SNs need to *gossip* SN lists to each other. On one hand, it is desirable

that the SNs frequently exchange gossiping information, so that the SNs have reasonably accurate knowledge about the SNs that are currently operating. On the other hand, it is desirable to minimize the amount of overlay maintenance traffic, as this traffic can be a burden on the SNs. Our experiments have enabled us to determine that KaZaA’s signalling protocol makes explicit provisions to this end. Specifically, one of the fields in the SN refresh list is the “freshness” field. This value enables the peers (ONs and SNs) to estimate the freshness of the SN availability information. We have observed that a newly connecting ON completely ignores SNs that have a low value of “freshness” regardless of the proximity or workload values of the SN.

7. **Firewall avoidance and NAT circumvention:** Although not typically addressed in the P2P research literature, a P2P file sharing will not thrive unless it takes explicit measures to deal with firewalls and NATs, which are prevalent throughout the Internet. As discussed in Section 4, KaZaA uses dynamic port numbers along with its hierarchical design to avoid firewall blocking. Furthermore, it uses connection reversal to allow NATed peers to share files.

## 6 Related Work

Recently there have been a number of P2P measurement studies, although to our knowledge none has carefully examined KaZaA. The identification of P2P specific traffic is considered in [27] and [12]. The accurate signature based techniques discussed in [27] could be deployed by an ISP to identify and filter illicit P2P traffic. An analysis of P2P traffic by measuring flow-level information collected at multiple border routers across a large ISP-network is done in [28]. By measuring KaZaA traffic in the University of Washington campus, the paper [8] studies file-sharing workloads and develops models for multimedia workload. A recent paper that characterizes P2P traffic is [11]. This last measurement of P2P traffic was done at the link level by reverse engineering the protocols of P2P applications and identifying characteristic strings in the payload.

A crawling system was previously developed for the Gnutella P2P network [24]. See also [10] for some additional work on crawling Gnutella and Napster. In [19], a KaZaA overlay crawler is developed, which is used to study the extent of polluted content in KaZaA.

There has also been some recent measurement work on the spread of spyware in P2P systems. In [26] the authors develop signatures for popular spyware and obtain traces of network activity within the University of Washington campus to quantify the spreading of spyware.

Scalability issues in the Gnutella network are studied in [21] and [2]. The latter presents a detailed analysis of how scalability is improved with flow control, dynamic

topology adaptation, one-hop replication, and node heterogeneity. The paper [20] studies the general problem of search and replication strategies in unstructured P2P networks. It proposes a query algorithm based on multiple random walks which is shown to be as fast as the Gnutella's query flooding method but reduces the network traffic by up-to two orders of magnitude. The paper [3] evaluates and compares different replication strategies in unstructured P2P networks. It identifies that uniform and proportional replication strategies yield sub-optimal performance and then propose an optimal replication strategy based on square root replication.

The recent paper [31] studies the advantages of design of unstructured P2P systems based on super-peers (SNs in KaZaA). The paper explores a number of important questions such as the potential drawbacks of super-peers, ways of improving reliability of super-peers, and the maximum number of children a super-peer should entertain to optimize efficiency. The paper [30] argues for a hybrid architecture for P2P systems, whereby structured search techniques are used to index and locate rare items, and flooding based techniques are used for locating highly replicated content.

## 7 Summary

To conclude this paper with a short summary of our findings. The supernodes form the backbone of the KaZaA network. There are roughly 30,000 supernodes; the average supernode lifetime is about 2.5 hours, although these lifetimes greatly vary across supernodes. Each supernode maintains a list of SNs it believes to be up. The SNs frequently exchange subsets of these lists with each other. Thus, the KaZaA backbone is self-organizing and is managed with a distributed, but proprietary, gossip algorithm. SNs establish both short-lived and long-lived TCP connections with each other. The SNs shuffle the long-lived connections (with average duration of 23 minutes), which improves search performance when search is carried out over long periods (hours) as it is often done in P2P file sharing. Each SN has about 40-60 connections to other SNs at any given time. Each SN has about 60-150 children ONs at any given time. Each SN maintains an index, storing the metadata of the files its children are sharing. SNs do not exchange metadata with each other.

When a user first acquires a KaZaA client, the client comes pre-installed with a list of candidate SNs. When the client is executed, the client connects with one or more the SNs in this list and obtains refresh lists. It appears that the entries in these refresh lists are biased with locality; the provided SNs are close to the ON with respect to various locality metrics. When an ON obtains a new list of SNs, it modifies its own cached list. Thus the ON-to-SN connections are formed in a decentralized, distributed manner and appear to take locality into account.

KaZaA has a life of its own, without requiring any intervention from a centralized

authority. Unlike Napster, KaZaA cannot be shut down by simply pulling the plug on a centralized server farm. Thus, KaZaA will likely persist for the foreseeable future. Many design decisions taken by the creators of KaZaA (and KaZaA-lite) seem to have been done without careful consideration. We conjecture that there is significant room for improving the search performance in two-tier unstructured P2P file sharing systems.

## Acknowledgment

We gratefully acknowledge help of ZhongQiang Chen for valuable feedback and discussions.

## References

- [1] A GiFT plugin for FastTrack, <http://developer.berlios.de/projects/gift-fasttrack/>
- [2] Y. Chawathe, S. Ratnasamy, L. Breslau, N. Lanham, S. Shenker, “*Making Gnutella-like P2P Systems Scalable*,” ACM SIGCOMM, Karlsruhe, Germany, August 2003.
- [3] E. Cohen, S. Shenker, “*Replication strategies in unstructured peer-to-peer networks*,” Proceedings of the ACM SIGCOMM Conference, Pittsburgh, PA, USA, August 2002.
- [4] L. Garces-Erce, K.W. Ross, E. Biersack, P. Felber, G. Urvoy-Keller, “*TOPLUS: Topology Centric Lookup Service*”, Fifth International Workshop on Networked Group Communications (NGC’03), Munich, Germany, September 2003.
- [5] The Gnutella protocol specification v4.0. <http://dss.clip2.com/GnutellaProtocol04.pdf/>
- [6] giFT-FastTrack plugin for giFT, <http://giftproject.org>
- [7] Global Crossing’s IP Network Performance, <http://ipstats.globalcrossing.net/>
- [8] K. P. Gummadi, R. J. Dunn, S. Saroiu, S. D. Gribble, H. M. Levy, and J. Zahorjan, “*Measurement, Modeling, and Analysis of a Peer-to-Peer File-Sharing Workload*,” Proceedings of the 19th ACM Symposium on Operating Systems Principles (SOSP-19), New York, USA, October 2003.
- [9] K. Gummadi, S. Saroiu, S. Gribble, “*King estimating Latency between Arbitrary Internet End Hosts*,” Proceedings of the SIGCOMM Internet Measurement Workshop, Marseille, France, November 2002.

- [10] K.P. Gummadi, S. Saroiu and S.D. Gribble, “*A Measurement Study of Peer-to-Peer File Sharing Systems,*” Proceedings of Multimedia Computing and Networking, January 2002 (MMCN’02), San Jose, CA, USA.
- [11] T. Karagiannis, A. Broido, N. Brownlee, K. Claffy, M. Faloutsos, “*Is P2P dying or just hiding?*,” IEEE Globecom, Dallas, Texas, November 2004.
- [12] T. Karagiannis, A. Broido, N. Brownlee, K. Claffy, M. Faloutsos, “*File-sharing in the Internet: A characterization of P2P traffic in the backbone,*” Technical report, November, 2003.
- [13] KaZaA Hack 2.5, <http://www.kazaahack.net/home.html>
- [14] KaZaA Homepage, <http://www.kazaa.com>
- [15] KaZaA Lite 2.10, <http://www.k-lite.tk/>
- [16] KaZaA P2P FastTrack File Formats <http://home.hetnet.nl/~frejon55/>
- [17] J. Kurose, K.W. Ross, “*Computer Networking: A Top-Down Approach Featuring the Internet,*” Addison-Wesley, 2005.
- [18] N. Leibowitz, M. Ripeanu, and A. Wierzbicki, “*Deconstructing the Kazaa Network,*” 3rd IEEE Workshop on Internet Applications (WIAPP’03), Santa Clara, CA, USA, June 2003.
- [19] J. Liang, R. Kumar, Y. Xi, K.W. Ross, “Pollution in P2P File Sharing Systems,” submitted, 2004.
- [20] Q. Lv, P. Cao, E. Cohen, K. Li, S. Shenker, “*Search and Replication in Unstructured Peer-to-peer Networks,*” 16th Annual ACM International Conference on Supercomputing, New York City, NY, USA, June, 2002.
- [21] Q. Lv, S. Ratnasamy, S. Shenker, “*Can Heterogeneity Make Gnutella Scalable?*”, The 1st International Workshop on Peer-to-Peer Systems (IPTPS), Cambridge, MA, USA, March 2002.
- [22] Overpeer Inc, <http://www.overpeer.com>
- [23] Regional characteristics of P2P, <http://www.sandvine.com>
- [24] M. Ripeanu, I. Foster, and A. Iamnitchi, “*Mapping the Gnutella network: Properties of large-scale peer-to-peer systems and implications for system design,*” IEEE Internet Computing Journal, vol. 6, no. 1, 2002.



- [25] S. Saroiu, K. P. Gummadi, R. J. Dunn, S. D. Gribble, and H. M. Levy, “*An Analysis of Internet Content Delivery Systems*,” Proceedings of the 5th Symposium on Operating Systems Design and Implementation (OSDI 2002), Boston, MA, USA, December 2002.
- [26] S. Saroiu, S.D. Gribble and Henry M. Levy, “*Measurement and Analysis of Spyware in a University Environment*,” Proceedings of the First Symposium on Networked Systems Design and Implementation (NSDI '04), San Francisco, CA, March 2004.
- [27] S. Sen, O. Spatcheck and D. Wang, “*Accurate, Scalable In-Network Identification of P2P Traffic Using Application Signatures*,” Proceedings International WWW Conference, New York, USA, May 2004.
- [28] S. Sen and J. Wang, “*Analyzing Peer-to-Peer Traffic Across Large Networks*,” ACM/IEEE Transactions on Networking, Vol. 12, No. 2, April 2004.
- [29] Sig2dat tool for FastTrack network, <http://www.geocities.com/vlaibb/tools.html>
- [30] B. Thau Loo, R. Huebsch, I. Stoica, J. Hellerstein, “*The Case for a Hybrid P2P Search Infrastructure*,” The 3rd International Workshop on Peer-to-Peer Systems (IPTPS), San Diego, CA, USA, February 2004.
- [31] B. Yang, H. Garcia-Molina, “*Designing a Super-Peer Network*,” 19th International Conference on Data Engineering, Bangalore, India, March 2003.