

# Performance Characteristics of Mirror Servers on the Internet

Andy Myers

A thesis submitted in partial fulfillment  
of the requirements for the degree of Master of Science  
in Electrical and Computer Engineering at Carnegie Mellon University

December, 1998

## **Abstract**

As a growing number of web sites introduce mirrors to increase throughput, the challenge for clients becomes determining which mirror will offer the best performance when a document is to be retrieved. In this paper we present findings from measuring 9 clients scattered throughout the United States retrieving over 490,000 documents from 45 production web servers which mirror three different web sites. We have several interesting findings that may aid in the design of protocols for choosing among mirror servers. Though server performance varies widely, we have observed that a server's performance relative to other servers is more stable and is independent of time scale. In addition, a change in an individual server's transfer time is not a strong indicator that its performance relative to other servers has changed. Finally, we have found that clients wishing to achieve near-optimal performance may only need to consider a small number of servers rather than all mirrors of a particular site.

# 1 Introduction

Distributing replicas of servers across the Internet has been employed for many years as a way to increase reliability and performance in the presence of frequent accesses by many clients. The weakness of replication (also known as mirroring) is that clients only have ad hoc mechanisms for selecting among the mirrors. Partridge et al [22] have proposed a scheme called *anycast* that allows a client to automatically reach the replica of a server which is the smallest number of network hops away. Others [8, 5] have observed that static metrics of proximity, such as distance in hops, are less effective at finding a server that will deliver good performance than metrics which take dynamically changing network and server conditions into account.

This paper presents results of a study based on client probing of web servers that we have undertaken to gain insight on approaches to designing server selection mechanisms. Our focus is on characterizing the performance a client receives when transferring documents from mirror servers. We wish to answer three questions:

- Does performance observed by a client vary across mirror servers?
- How dynamic is the set of servers that offer good performance?
- What are effective methods to use in finding a server that will offer good performance for a given client?

To answer the first question, we have looked at the time required to retrieve a document from each mirror server of a site. We have found that the difference in performance between the best and worst servers is typically larger than an order of magnitude, and can grow larger than two orders of magnitude on occasion. This result shows that performance does indeed vary largely from one server to another.

The second question is an attempt to explore how dynamic server performance changes are. By counting the number of servers that a client must visit over time in order to achieve good performance, we can see whether the set of servers that offer good performance at any given time is small or large. We found that the set is usually fairly small, indicating less dynamic behavior.

To answer the third question, we will evaluate two heuristics that a server selection system might potentially employ to select a server for a client. The first heuristic is to make the assumption that a server which offers good performance has a probability of no longer offering good performance which increases with time. We found that this heuristic is not effective. In other words, a server is equally likely to offer good performance whether it offered good performance 2 hours or 2 days ago.

The second heuristic we will consider is to make the assumption that a drop in a server's performance corresponds to a similar drop in its likelihood of offering better performance than the other servers. In this case, we found that the heuristic does hold, but only very mildly. For example, if a server offers better performance than any other server for a given client and subsequently the server's performance drops dramatically, it is likely that the server no longer offers better performance. However, moderate or small performance drops do not indicate whether or not a server still offers better performance than other servers.

Finally, we will consider the effect of document choice on server choice. Though we assume that all mirrors of a server have the same set of documents, it might be the case that some factor

such as document size or popularity would affect the performance of a server. We found that server choice is independent of document choice almost all the time.

To summarize, we have five main results:

- Performance can vary widely from one server to another.
- Clients can achieve near-optimal performance by considering only a few servers out of the whole group of mirrors.
- The probability of any server’s rank change depends very little on the time scale over which the rank change takes place.
- There is a weak but detectable link between a server’s change in transfer time and its change in rank.
- Server choice is independent of document choice in most instances.

We discuss the implications of these results in Section 9.

## 1.1 Related work

Previous work on server selection techniques can be divided into four categories: network-layer server selection systems, application-layer selection systems, metric evaluation, and measurement studies. The first includes work dealing with finding the closest server in terms of number of network hops or in terms of network latency [15, 9, 3, 14, 19, 20, 22]. The second consists of systems that take application performance metrics into account [5, 7, 21, 23, 24, 26, 12]. Most of these systems use a combination of server load and available network throughput to select a server for a client. The third category consists of evaluations of server selection metrics [8, 13, 17, 25]. These studies propose new metrics and test them experimentally.

The fourth category, which includes this work, consists of studies that collect data characterizing the behavior of mirror servers in order to draw conclusions about the design of server selection systems. Bhattacharjee et al [4] measured “server response time,” defined to be the time required to send a query to a server and receive a brief response, using clients at a single site to visit two sets of web sites. While neither set of sites were true mirrors, each set consisted of servers with similar content. Bhattacharjee also measured the throughput between a client and four FTP servers. Carter and Crovella [8] measured ping times and hop counts to 5262 web servers to determine how well one approximated the other. In contrast, our study is on a larger scale, using multiple client sites, a longer measurement period, and a larger number of groups of popular web servers that are true mirrors.

There have been several other web-related measurement studies. Balakrishnan et al [2] analyzed a trace of web accesses to determine how stable network performance is through time and from host to host. Gribble and Brewer [16] looked at users’ web browsing behavior, exploring server response time, burstiness of offered load, and the link between time of day and user activity. Cunha et al [11] also collected user traces via a customized version of Mosaic and looked at a number of factors including document size and popularity. Arlitt and Williamson [1] searched for trends present in a variety of different WWW workloads based on server access logs. Finally, Crovella and Bestavros [10] have found evidence for self-similarity in WWW traffic.

Client Site	Avg. time of one group	Number of fetches	Failure rate
Carnegie Mellon	0:32:49 hours	54695	10.18%
Georgia Tech.	0:23:47	60021	11.55%
ISI	0:36:31	53200	22.13%
U. of California, Berkeley	0:32:33	55062	4.62%
U. of Kentucky	0:31:14	55091	12.76%
U. of Mass., Amherst	1:10:34	36542	10.95%
U. of Texas	0:39:34	51640	4.70%
U. of Virginia	0:19:19	62405	28.88%
Washington U., St. Louis	0:23:16	62187	1.96%

Figure 1: Average time for one round of fetches, number of fetches completed, and failure rate for each client site

The rest of this paper consists of a description of our data collection system (Section 2), a general picture of the data we collected (Sections 3 and 4), a discussion of our findings (Sections 5 through 8), implications of our results (Section 9), and conclusions (Section 10).

## 2 Data collection methodology

At each of nine client sites where we had guest accounts (listed in Figure 1) a perl script periodically fetched documents from each server in three sets of mirrored web sites (the Apache Web Server site, NASA’s Mars site, and News Headlines) listed in Figure 2. The Apache and Mars web sites were true mirrors: each of the servers in one set held the same documents at the same time. However, the News sites were an artificial mirror since they did not contain the same documents. The News servers were picked from Yahoo’s index (<http://www.yahoo.com/>). Current headlines from each of the News sites were fetched and the transfer times were normalized so that all News documents appeared to be 20 KB long. For the Mars and Apache servers, we used five documents ranging in size from 2 KB to 1.3 MB (listed in Figure 3).

Clients visited servers sequentially, fetching all documents from a server before moving on to the next. Similarly, all mirrors of one site were visited before moving on to the next site. For example, a client would start by visiting <http://www.sgi.com/>, the first Mars mirror on the list, and fetching each of the Mars documents from it. Then the client would fetch the Mars documents from the second Mars server, then the third, and so on. When all of the Mars servers had been visited, the client would move on to the Apache mirrors, and finally to the News sites. We refer to the process of visiting all servers and collecting all documents once as a *group* of fetches.

After all servers were visited, the client would sleep for a random amount of time taken from an exponential distribution with a mean of  $1/2$  hour added to a constant  $1/2$  hour. By scheduling the next group of fetches relative to the previous group’s finish time (rather than its start time), we avoided situations in which multiple fetches from the same client interfered with each other, competing for bandwidth on links near the client.

We introduced the delay between fetches to limit the load our fetches created on client and server sites. A typical group of fetches involved transferring more than 60 MBytes of data to a

Mars sites	
http://mars.sgi.com	http://www.sun.com/mars
http://entertainment.digital.com/mars/JPL	http://mars.novell.com
http://mars.primehost.com	http://mars.hp.com
http://mars.excite.com/mars	http://mars1.demonet.com
http://mars.wisewire.com	http://mars.ihighway.net
http://pathfinder.keyway.net/pathfinder	http://mpfwww.arc.nasa.gov
http://mars.jpl.nasa.gov	http://www.ncsa.uiuc.edu/mars
http://mars.sdsc.edu	http://laguerre.psc.edu/Mars
http://www.ksc.nasa.gov/mars	http://mars.nlanr.net
http://mars.catlin.edu	http://mars.pgd.hawaii.edu
News sites	
http://www.cnn.com/	http://www.nytimes.com/index.gif
http://www.latimes.com/	http://www.washingtonpost.com/
http://www.csmonitor.com/	http://www.usatoday.com/
http://www.abcnews.com/	http://www.msnbc.com/
http://www.s-t.com/	http://nt.excite.com/
http://news.bbc.co.uk/	http://www.newscurrent.com/
http://pathfinder.com/time/daily	http://www.sfgate.com/news/
http://headlines.yahoo.com/Full_Coverage/	http://www.topnews.com/
Apache sites	
http://www.rge.com/pub/infosystems/apache/	http://apache.compuex.com/
http://apache.arctic.org/	http://ftp.epix.net/apache/
http://apache.iquest.net/	http://www.apache.org/
http://apache.utw.com/	http://www.ameth.org/apache/
http://apache.technomancer.com/	http://apache.plinet.com/
http://fanying.eecs.stevens-tech.edu/pub/mirrors/apache/	

Figure 2: Servers visited

client. If the fetches finished in 30 minutes, the average transfer rate would have been 266 Kbps, which is a noticeable share of the traffic on a LAN. The delay between groups of fetches lowered the average resource utilization to roughly half the original average bandwidth.

We used the lynx<sup>1</sup> web browser to perform fetches. Choosing lynx was a compromise between realism and ease of implementation. Lynx is an actual production web browser that people use every day. At the same time, it is easy to control via command line switches, allowing us to run fetches via a perl script. Implementing our own URL fetch code might not have captured the characteristics of actual browsers. Conversely, using a more popular, hence more realistic, browser, e.g. Netscape, would have presented a significant programming challenge.

Our client script would invoke lynx to retrieve a URL and send it to standard output. The number of bytes received by lynx was counted and recorded along with the amount of time the fetch took to complete. If a fetch did not terminate after five minutes, it would be considered unsuc-

<sup>1</sup> Available from <http://lynx.browser.org/>

	URL	Size (bytes)
Mars documents		
0	/nav.html	2967
1	/2001/lander.jpg	70503
2	/mgs/msss/camera/images/12_31_97_release/2303/2303p.jpg	235982
3	/mgs/msss/camera/images/12_31_97_release/2201/2201p.jpg	403973
4	/mgs/msss/camera/images/12_31_97_release/3104/3104p.jpg	1174839
Apache documents		
0	dist/patches/apply_to_1.2.4/no2slash-loop-fix.patch	1268
1	dist/CHANGES_1.2	90631
2	dist/contrib/modules/mod_conv.0.2.tar.gz	74192
3	dist/apache_1.2.6.tar.gz	714976
4	dist/binaries/linux_2.x/apache_1.2.4-i586-whatever-linux2.tar.Z	1299105

Figure 3: URLs of documents fetched from Mars and Apache servers

cessful and the associated lynx process would be killed. We chose five minutes as a compromise between achieving a complete picture of a server’s behavior and forcing groups of fetches to finish in a reasonable amount of time. The observable effects of such a short timeout were a slightly higher failure rate, especially among larger documents. Possible causes for timeouts are network partitions, client errors (lynx might have frozen), server errors (the server might have stopped providing data), or shortages of available bandwidth. In our analysis, we treat these incidents as failures to collect data, rather than as failures of servers.

Fetches could also be unsuccessful if the number of bytes returned was incorrect. We found that the wrong number of bytes usually indicated a temporary failure such as a “server too busy” message although in some cases it signified that the server no longer existed (failed DNS query) or was no longer mirroring data. We assumed that every fetch which returned the proper number of bytes succeeded.

It was more difficult to identify failed fetches from the News sites. Since we were retrieving news headlines, each page’s content was constantly changing so we could not use a hard-coded size to determine success. A simple heuristic that worked well was to assume that all fetches that returned less than 600 bytes were failures. This value was larger than typical error messages (200-300 bytes) and smaller than typical page sizes (as low as 3k on some servers). As with the other servers, fetches lasting five minutes were considered failures.

While our fetch scripts were running, there were multiple occasions on which client machines crashed or were rebooted. To limit the impact of these interruptions, we used the Unix `cron` system to run a “nanny” script every 10 minutes which would restart the fetch script if necessary. This kept all fetch scripts running as often as possible.

## 2.1 Limitations

While our methodology was sufficient to capture the information in which we were most interested, there were some data that we were not able to capture. Because of the relatively large, random

gaps between fetches to the same server, we were unable to capture shorter-term periodic behavior. Further, because each group of fetches finished in a different amount of time because of variations in server load and network congestion, the distribution of fetch interarrivals to a single server from a client was extremely hard to characterize and exploit. Thus, we were unable to map the observed frequency of network conditions to the actual frequency of occurrence of these conditions.

No two fetches from a given client were done simultaneously to prevent the fetches from competing with each other. At the same time, we would like to compare results across servers to rank servers relative to one another. There is a reasonable amount of evidence which suggests that network performance changes over longer time scales [26][2] while our measurements took place over shorter time scales. On average, clients visited all Mars mirrors in just over 17 minutes, all Apache mirrors in under 13 minutes, and all News sites in less than one and a half minutes. Because of these results, we believe that it is valid to treat sequential fetches as occurring simultaneously.

Another artifact of sequential fetches is that periods of network congestion are possibly under-represented in the data. As congestion increases, fetches will take longer. The result is that the number of fetches completed during periods of congestion will be lower than the number completed during periods with less congestion. If periods of congestion are short-lived, only a few fetches will reflect the congestion. If periods of congestion are long-lived, all fetches will take longer but the total number of groups of fetches completed will be smaller.

DNS caching effects could also potentially bias our results. Depending on the DNS workload at a given client site, DNS entries for the servers in our study may or may not remain in the local cache from one group of fetches to another. In fact, cache entries could even be purged within a group of fetches. The DNS lookups added a potentially highly variable amount of time to each fetch we performed. Performing the lookups separately would have been possible, but less realistic.

Finally, we must consider inter-client effects. Because each client's fetches are independently scheduled, two clients could wind up visiting the same server at the same time. We will refer to such an incident as a *collision*. We believe that collisions have a negligible effect on fetch times. Further, less than 10% of all fetches were involved in collisions.

### 3 Data characteristics

All clients began fetching documents on the afternoon of Thursday, April 23, 1998 and continued until the morning of Thursday, May 14, 1998. During this 3 week period, there were a total of 490843 fetches made. By data set, there were 287209 fetches to Mars servers, 157762 to Apache servers, and 45872 to News servers. The much lower number for the News data is mostly due to the fact that we only fetched one document from each News site compared to five from each Mars and Apache site. We can estimate the number of times each set of servers was visited by dividing the number of fetches by the number of combinations of servers and documents. For Mars, we divide 287209 by 100 (20 servers x 5 documents) to find that the Mars servers were visited 2872 times. Similarly, we see that Apache servers were visited 2868 times and News servers were visited 2867 times.

The slightly lower number of visits to Apache and News sites is a product of the way the client fetch script reacted to crashes. When a client was restarted, it began fetching documents from the first server on its list rather than starting at the place where the last series of fetches left off. Since

clients visited Mars sites first, then Apache sites, and finally News sites, it is not surprising that there are more fetches to Mars sites than to Apache sites and more fetches to Apache sites than to News sites.

The number of fetches performed and the average length of time that one group of fetches took to complete at each client site can be found in Figure 1. As expected, sites with longer group fetch times completed fewer fetches. We believe the differences across clients reflect variation in the amount of available bandwidth and machine speed at each client site.

Figure 1 also shows the percentage of fetches that were classified as failures (because timeouts and improper amounts of data returned). By client, the proportion of failures ranged from 1.96% to 22.13% of fetches. Considering the loss rate by server set, we see that Mars servers failed 5.85% of the time, News servers failed 9.49% of the time, and Apache servers failed 24.23% of the time. As far as we can tell, the differences in failure rates across types of mirrors are not the result of using one brand of web server or another. However, we did notice that three Apache servers consistently timed out for some clients while they succeeded a reasonable amount of time for other clients. These three servers account for most of the Apache servers' comparatively high failure rate.

### 3.1 Ranks

Throughout this paper, we use *rank* to compare servers' performance. In this section we explain how ranks are computed and give some insight into what differences in rank mean. A ranking of servers is computed for each data set (Mars, News, or Apache) for each group of fetches at each client. Recall that after each group of fetches, a client has performance data for each web server. For each document, we can order the servers by their fetch times from lowest to highest, discarding those servers whose fetches failed. A server's rank is merely its place in this order. The server which comes first in the order has the highest rank (0), the server which comes next has a rank of 1, and so on. In our terminology, lower ranks correspond to better server performance. In summary, each successful group of fetches generates one set of ranks for each of the 11 documents: 5 sets for Mars documents, 5 for Apache documents, and one for the News document.

There is some inaccuracy in our method of ranking servers: The tacit assumption in computing ranks is that the fetch times being compared were generated under identical conditions. As we have discussed in Section 2.1, this is not possible, but we believe that network conditions do not change a significant amount between the first and last fetch of a document from a group of servers.

Ranks are not significant performance indicators by themselves. Ranks will not say whether or not the difference in performance between servers is negligible. But in the data that we collected, we have found a very strong link between noticeable differences in performance and differences in rank.

Figure 4 plots the normalized, average increase in transfer time vs. server rank for document 4 of the Mars data set. It was produced by averaging the differences of all pairs of servers with ranks  $i$  and  $i - 1$  in each group. The graph shows a definite rise in transfer time as rank increases. For example, we see that on average, a server with a rank of 4 has twice the transfer time of a server with a rank of 0. Further, the server with the largest rank (17) takes more than 30 times as long to transfer a document as the best server, and it takes more than 3 times as long to deliver a document as a server with a rank of 14.

The primary point of Figure 4 is that rank changes usually correspond to noticeable perfor-

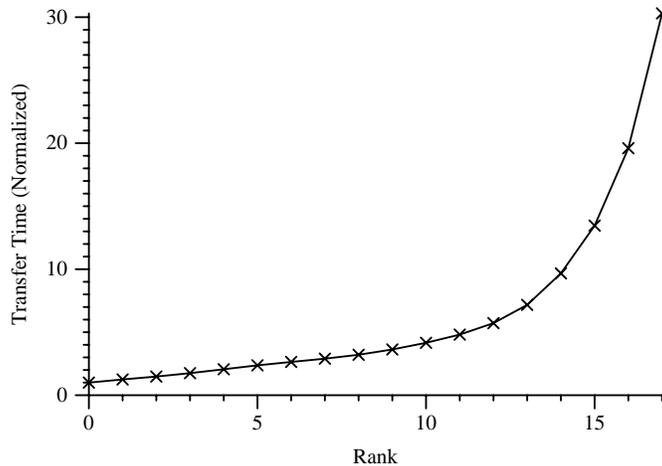


Figure 4: Average amount of separation between rank values for Mars servers, document 4, all clients aggregated.

mance changes for document 4 of the Mars set. All other documents from Mars, Apache, and News produced similar graphs, though the Apache and News data tended to have much larger differences in performance. This gives us confidence that ranks are a reasonable way to talk about the relative performance of servers.

## 4 Summary statistics and distributions

Our data consists of random samples (as we note in the next section, there is almost no significant sequential correlation in our samples) where each sample consists of a transfer time from a client to a server and its ranking relative to the other transfers in its group of fetches. This section summarizes these samples in terms of general statistics and analytic distributions. Conceptually, the analysis gives some insight into what a random client can expect from a random mirror site for different sizes and kinds of documents. There are two main results here. First, transfer times and server rankings exhibit considerable variability. Second, transfer times, taken to a fractional power, are well fit by an exponential distribution.

The analysis is from the point of view of a random client site (from Figure 1) attempting to fetch a particular document from a set of mirror sites (Figure 2.) There are 11 different combinations here (Apache and Mars each serve five different documents while News serves one virtual document.) For each of these combinations, we examine the transfer times and corresponding ranks for all the client fetches of the document to the set of mirror sites. In effect, we factor out the set of mirrors and the document size here by doing this.

Figure 5 presents the summary statistics of transfer times and ranks for each of the combinations. Notice that mean transfer times as well as standard deviations increase with increasing document size. Further, transfer times are highly variable — standard deviations are about as large as means, and we see maxima and minima near the limits we placed on observed transfer times (300 seconds.) It is important to note that the maximum transfer time of 638.98 seconds for the

Dataset/Doc	Transfer time (seconds)					Ranks				
	Mean	StdDev	Median	Min	Max	Mean	StdDev	Median	Min	Max
Apache/0	1.9632	5.8366	.7	0.1000	230.5100	4.2790	2.9834	4	0	10
Apache/1	3.9112	7.9753	2	0.3800	297.7000	4.2737	2.9610	4	0	10
Apache/2	3.2929	6.3993	1.7	0.3000	293.9000	4.1372	2.8277	4	0	10
Apache/3	15.4776	18.2385	10.7	1.3000	299.9000	3.9916	2.7643	4	0	10
Apache/4	23.1960	22.9257	17.9	2.2000	298.2000	3.7789	2.6920	4	0	10
Mars/0	1.5416	4.6808	0.7	0.1000	296.6000	8.2060	5.0596	8	0	17
Mars/1	2.6929	6.5319	1.3	0.1000	292.6000	8.1667	5.0496	8	0	17
Mars/2	5.8062	9.4102	3.3	0.3000	290.5000	8.1287	5.0350	8	0	17
Mars/3	8.7380	12.3967	5.3	0.6000	297.3000	8.0995	5.0268	8	0	17
Mars/4	19.9019	23.5427	13.9	1.6000	298.2000	7.9213	4.9654	8	0	17
News/0	3.8185	11.8028	1.06	0.1200	638.9800	6.4049	4.0692	6	0	14

Figure 5: Summary statistics of transfer time and corresponding ranks.

News/0 dataset is due to our normalizing the transfer times for News documents according to their size to approximate always fetching a 20 KB document. In some cases, particularly slow fetches can result in normalized transfer times exceeding 300 seconds. This is rare.

Figure 5 also shows statistics of ranks. An interesting observation here is that the standard deviation of ranks, although quite large, does not bode disaster for server selection algorithms. A random selection is likely to result in an average server. Further, it may well be the case that some servers vary less in their ranking than others – for example, the rankings of a few good servers may very well remain stable while the remaining servers have more drastically varying rankings. The reader may notice that the median and maximum ranks noted in the table are low - for example, there are 16 News servers, yet the maximum rank observed is 14. This effect is due to the fact that in any group there is likely to be at least one failed fetch which results in a truncation of ranks.

While summary statistics provide some insight on the performance, both absolute (transfer time) and relative (rank) a client can expect to receive, they provide a very limited view of the distribution of these quantities. To better understand the distribution of transfer times, we attempted to fit a variety of analytic distributions to the data. The quality of such a fit can be determined by a quantile-quantile plot, which plots the quantiles of the data versus the quantiles of the distribution [18, pp. 196-200]. A good fit results in a straight line, regardless of the parameters chosen for the analytic distribution.

We tried normal, exponential, and Poisson distributions. None of these fit the transfer time data very well, especially at the tails. The distribution of transfer times is heavy-tailed compared to these distributions. Next, we tried the log-normal distribution by testing if the logarithms of our data points were normally distributed. Figure 6(a) shows the quantile-quantile plot of the logarithm of the Mars/1 dataset versus normal and is representative of the other datasets. Generally, log-normal was much better than the earlier distributions. This result agrees with Balakrishnan et al [2], who also found that a single client’s observed throughput can be modeled reasonably well by a log-normal distribution.

We next tried a power transformation — raising the data to a fractional power — and seeing if the transformed data could be fitted with a common analytic distribution. This turned out to provide the best results. For example, in Figure 6(b) we have raised the data to the 0.15 power and plotted

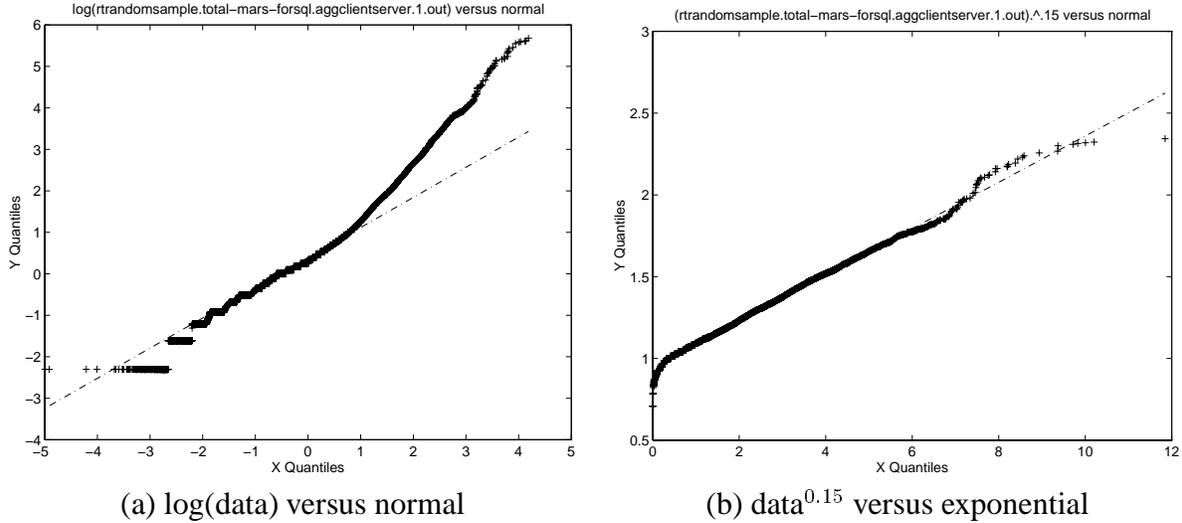


Figure 6: Quantile-Quantile plots comparing distribution of transfer times in Mars/1 to (a) log-normal and (b) power-exponential.

the quantiles of the transformed data versus the quantiles of an exponential distribution. The near-perfect linearity makes it clear that this “power-exponential” distribution is a particularly good fit for Mars/1. Power-exponential also fits the other datasets exceedingly well. Some caution must be used here, however. Because transfer times were artificially truncated at 5 minutes, we do not have an accurate picture of the tail of the distribution. It may be the case that the actual distribution of server transfer times is much more heavy-tailed, meaning that the power-exponential distribution may not fit this data as well as it seems to.

## 5 The time scale of rank changes

Once a client has found a highly ranked server, the client is interested in how long that server is likely to maintain a high rank among the servers the client could fetch from. Fundamentally, it is the time scale over which significant rank changes occur that is important. In this section, we show that most rank changes are small, even over relatively long time scales. Good servers remain good for long periods of time. Indeed, the probability of rank changes depends very little on the time scale over which they occur.

Given periodically sampled ranks, the natural way to study change on different time scales would be via frequency domain or time series analysis [6]. However, as we discussed in Section 2, our data was collected at exponentially distributed intervals, making this difficult. The transfer time data could be resampled periodically and new rankings computed, but such resampling is complex and since signal reconstruction from non-periodic samples is an area of current research, such an approach would be questionable as well as difficult to understand. We did informally try this method and found results similar to those presented here.

Our approach was to estimate the cumulative probability of rank changes over increasing time scales. Consider a single mirror server. From the point of view of a single client using the set of mirrors, we have a sequence of time-stamped samples of that server’s rank (as well as transfer

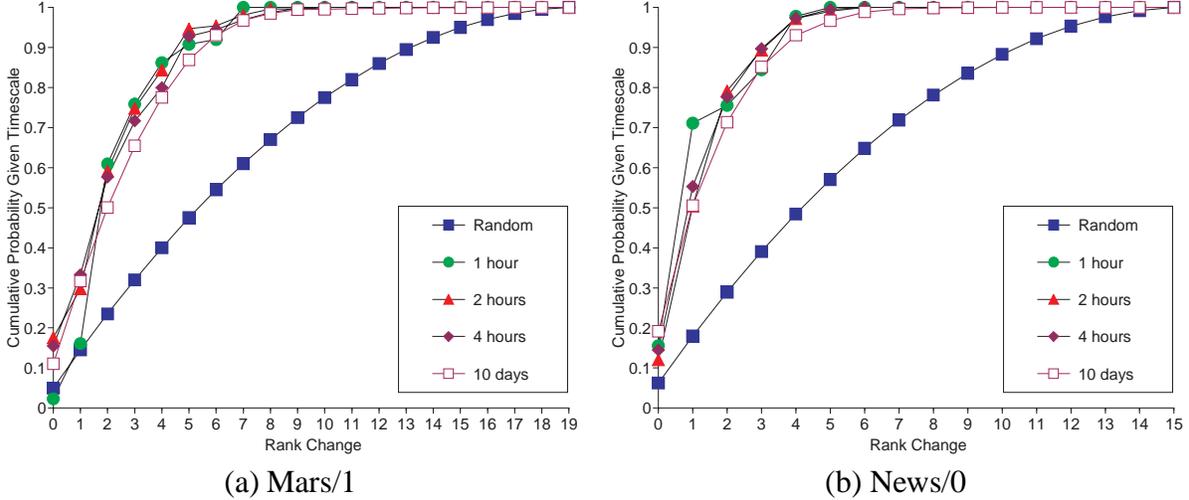


Figure 7:  $P[|r_{t+w} - r_t| \leq R \mid \text{sample period} \leq w \leq W]$  for (a) Mars/1 and News/0 plotted for several different values of  $W$  where  $R$  is the rank change and  $W$  is the maximum time period. Other Mars plots are similar.

times.) Now extract all the pairs of rank samples that are four or fewer hours apart. For each pair, subtract the earlier rank from the later rank and take the absolute value. Count the number of occurrences of each of the possible rank changes. Accumulating these in the appropriate order gives an estimate of the cumulative probability of rank changes given measurements four or fewer hours apart.

We may find that rank changes of three or fewer are 80% probable given time scales of four or fewer hours. Notationally, we express this as  $P[|r_{t+w} - r_t| \leq R \mid \text{sample period} \leq w \leq W] = 0.8$ , where  $R = 3$  is the rank change,  $W = 4$  hours is the maximum time scale and the  $rs$  are our rank samples. For each combination of  $W$  and  $R$  examined, we use a randomly selected 10,000 samples to assure a tight estimate of the probability. Further, we aggregate the probabilities across all clients for each dataset and document to obtain the point of view of a random client interacting with a random server within the group. Finally, it is important to note that we are limited by our average sampling interval of one hour — we cannot discern behavior for  $W < 1$  hour.

Figure 7 shows representative plots of the cumulative probability for the (a) Mars/1 and (b) News/0 datasets. The way to read these plots is to pick a time scale, follow the corresponding curve horizontally to the maximum rank change that is of interest, and then read the cumulative probability from the vertical axis. For example, considering the Mars/1 data set, we see that for time scales of two (or fewer) hours, rank changes of four (or fewer) occur with probability 0.9. The graphs also include the curve that would result if rankings were simply uniformly distributed random permutations.

It is clear from the graphs that most rank changes are small. The 10 day curves cover the vast majority of the data, and we can see that the smallest 25% of possible rank changes account for about 90% of rank changes.

The graphs also show that rank changes depend very little on the time scales over which they occur. If there was a strong dependency, the curves for the different time scales would be more widely separated. We can see that the curves for increasingly longer time scales slowly move to

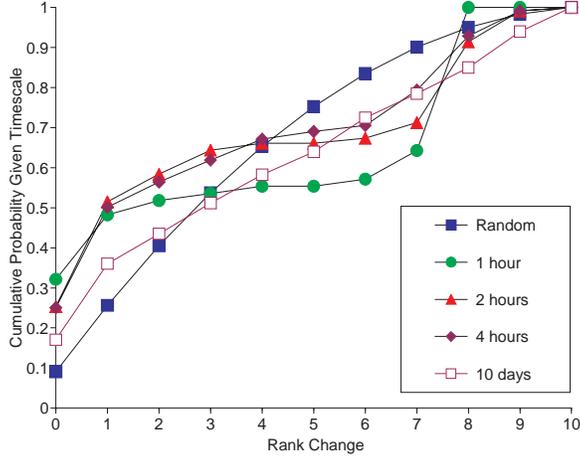


Figure 8:  $P[|r_{t+w} - r_t| \leq R \mid \text{sampleperiod} \leq w \leq W]$  for Apache/1, plotted for several different values of  $W$ . Other Apache plots are roughly similar, and all Apache plots differ significantly from Mars or News plots.

the right (toward larger rank changes), but the effect is very marginal. This is a very promising result. If a client can find a good server, it is highly likely that it will remain good for quite some time.

The graphs of Figure 7 are representative for the Mars and News datasets. Unfortunately, the Apache data shows very different behavior, as can be seen in Figure 8, which shows a cumulative probability plot for a representative, Apache/1. Here, we don't see the quick rise of the curves, so large rank changes are relatively much more probable than with the Mars and News data. Further, since the curves do not hug each other very closely, there is more dependence on the time scale. At this point, we do not understand why the Apache data is so different. The clearest distinguishing characteristic of the Apache sites is that they tend to run non-commercial web servers (the Apache web server) while the Mars and News sites tend to run commercial web servers (Netscape and Microsoft servers.) We have no evidence that this difference causes the discrepancy, however.

## 6 Changes in transfer time and rank

A client using a highly ranked server is interested in warning signs that may indicate that the server's ranking has changed dramatically. The client cannot measure rankings without measuring all of the mirror servers; it can only observe the transfer times it is experiencing on the currently chosen server. The natural question then is what, if any, relationship exists between the changes in transfer time a client observes and the changes in rank the server experiences. Our study shows that while a relationship does exist, it is very marginal.

Our approach was to estimate the cumulative probability of rank changes over increasing changes in observed transfer times. Consider a single mirror server. From the point of view of a single client using the set of mirrors, we have a sequence of samples of that server's transfer times and their corresponding ranks. We form the cross product of these samples and select a random subset of 100,000 of these sample pairs. For each pair of samples in the subset, we subtract

Dataset/Doc	Changes in transfer time (seconds)					Changes in rank				
	Mean	StdDev	Median	Min	Max	Mean	StdDev	Median	Min	Max
Apache/5	0.0039	8.4087	0	-123.8000	226.4100	0.0091	4.2022	0	-10	10
Apache/6	-0.0810	10.3995	0	-295.7300	267.8000	0.0010	4.1948	0	-10	10
Apache/7	-0.0503	9.0000	0	-292.5000	205.9000	-0.0621	4.0177	0	-10	10
Apache/8	-0.5457	25.4940	0	-285.3000	276.5000	-0.0196	3.8818	0	-10	10
Apache/9	-0.1912	31.8086	0.1	-278.0100	287.7000	-0.0367	3.8072	0	-10	10
Mars/0	0.1068	6.0450	0	-227.9100	221.4600	0.0244	7.1711	0	-17	17
Mars/1	0.1218	8.1173	0	-184.0000	232.5900	0.1330	7.0711	0	-17	17
Mars/2	0.1189	14.3483	0	-285.2000	287.4000	-0.0685	7.1195	0	-17	17
Mars/3	-0.0226	17.5260	0	-253.6000	282.1000	-0.0038	7.0849	0	-17	17
Mars/4	0.3308	34.5870	0	-286.9000	288.6000	0.0194	7.0992	0	-17	17
News/0	0.0282	17.1793	0	-298.8300	293.8300	-0.0316	5.8363	0	-14	14

Figure 9: Summary statistics of changes in transfer time and changes in corresponding ranks.

the transfer times and ranks.

Figure 9 shows the summary statistics of these changes in transfer time and corresponding rank. We see that the mean and median changes in both quantities are almost exactly zero. The distributions of these changes are also quite symmetric about zero. For this reason, we concentrate on absolute changes in transfer time and rank.

After taking absolute values, we count occurrences of value pairs to estimate the joint cumulative probability of absolute changes in rank and absolute changes in transfer time,  $P[|r_{t_i} - r_{t_j}| \leq R \wedge |d_{t_i} - d_{t_j}| \leq D]$  where  $R$  is the rank change and  $D$  is the change in transfer time. Since changes in rank are categorical, we can then trivially compute the cumulative probability of an absolute change in rank *given* an absolute change in transfer time of  $D$  or smaller. Notationally, this is  $P[|r_{t_i} - r_{t_j}| \leq R \mid |d_{t_i} - d_{t_j}| \leq D]$ . We aggregate the probabilities from all clients for each dataset and document to obtain the point of view of a random client interacting with a random server within the set of mirrors. The reader may object that this scheme also aggregates changes happening at all time scales. This is true. However, recall from Section 5 that changes in rank are virtually independent of time scale.

Figure 10 shows representative plots of the cumulative probability for the (a) Apache/4 and (b) News/0 datasets. The plots for all of the datasets are similar. The way to read these plots is to pick a change in duration, follow the corresponding curve horizontally to the maximum rank change that is of interest, and then read the cumulative probability from the vertical axis. For example, considering the News/0 data set, we see that for a transfer time change of 128 seconds or less, 90% of rank changes are of seven or less.

We can see that large changes in transfer time are more likely than small changes to indicate large rank changes. The curves for increasingly larger changes in transfer time shift toward the right (toward larger rank changes.) However, the difference is slight. For example, consider Apache/4: A rank change of three or smaller is 90% probable with a change in transfer time of one second or smaller, while a change of transfer time of up to 128 seconds reduces the probability only to 80%. This is typical of the Apache data, and the relationship is even less pronounced for the other data.

Another way to see the limited relationship of changes of rank to changes in transfer time

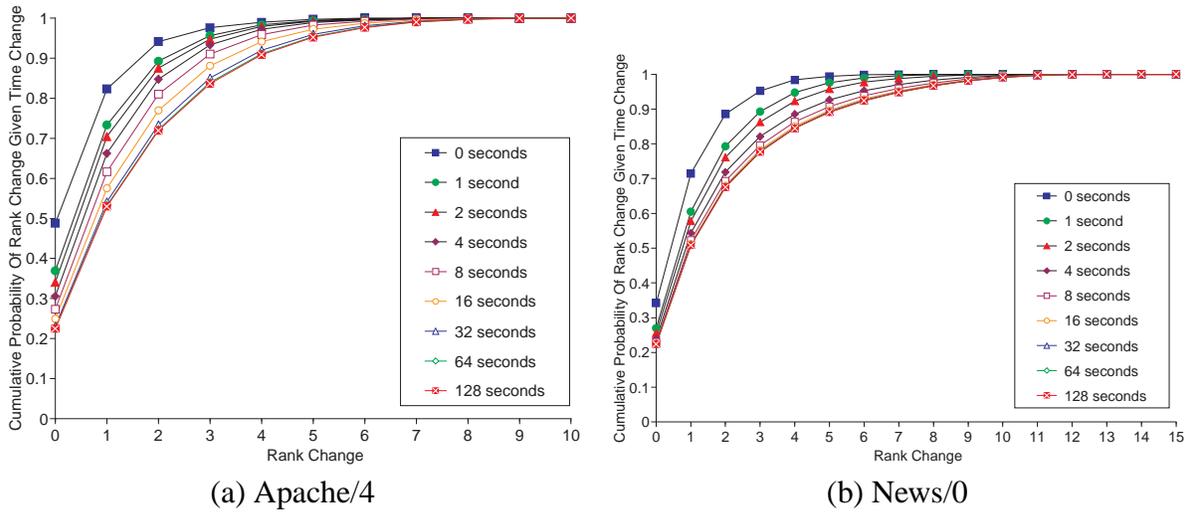


Figure 10: Cumulative probability of rank change given changes in transfer time less than  $D$  ( $P[|r_{t_i} - r_{t_j}| \leq R \mid |d_{t_i} - d_{t_j}| \leq D]$ ) for (a) Apache/4 and (b) News/0, plotted for several different values of  $D$ . All other plots are similar.

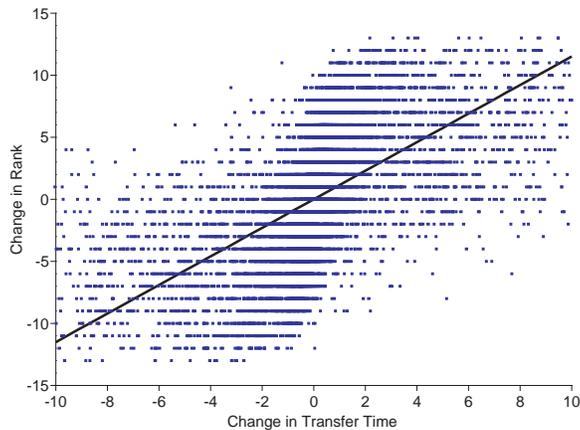


Figure 11: Changes in rank versus changes in transfer time (limited to +/- 10 seconds) for News/0 dataset. Note the inferiority of linear fit ( $R^2 = 0.36$ .) There is little relationship between changes in transfer time and changes in ranking.

is to plot rank changes against their corresponding transfer time changes. Figure 11 shows a representative plot for the News/0 data, where we have focused on transfer time changes in the  $[-10, 10]$  range. We have fit a least squares line to the data and have found that the relationship is marginal at best. The  $R^2$  value for the line is only 0.36. For a wider range of transfer times, the fit is even worse. Clearly, there is only a limited relationship between changes in transfer time and changes in rank.

## 7 Small server sets

The observation in Section 5 that most rank changes are small leads us to ask how many servers must a client consider to achieve optimal performance. If server ranks never changed, a client would only need to use one server, the one with the best rank. But because server ranks do change, a client will need to evaluate multiple servers to find the current best. We have found that a client needs to evaluate a very small number of servers, usually less than half the total number of servers, to achieve near-optimal performance. In this section, we define a server’s performance to be near-optimal, or “good,” if it can deliver a document in no longer than 10% more than the time it takes the current best server to deliver the same document.

We define a *server set* to be the minimum subset of servers from a group of mirrors that can provide near-optimal performance to a client. If a server set contains all the mirrors of a site, it means that a client must consider all mirrors when choosing a server. From the data we have, we can build a server set for each client-document combination using a straightforward greedy algorithm: In each group of fetches, all servers that deliver good performance are marked. The number of marks that each server accrues over all groups is computed, and the server,  $s$ , with the highest total, is added to the server set. The groups where  $s$  exhibited good performance are discarded, and the procedure is repeated on the remaining groups. The algorithm terminates when there are no groups left.

Figure 12 shows the composition of the server sets for 10 data sets composed of the five documents from U. Mass’s Apache data and the five documents from Washington U.’s Mars set. Each stripe from each column represents the proportion of time that one server offers good performance. For example, the first column of the graph shows the server set for the Wash. U. client’s fetches of document 0 from the Mars sites. Each stripe in that column corresponds to one server. For purposes of this discussion, it does not matter which server maps to which stripe. What is significant is the size of each stripe, which shows how often the corresponding server is able to deliver good performance. The first stripe in the column shows that one server is good in almost 70% of its fetches. The next stripe represents a server that is good in a little more than 10% of the remaining fetches.

The distribution and number of stripes show that client sites do not have to consider every server in the mirror set to achieve good performance. Rather, a small number of servers can provide good performance for a significant fraction of all fetches. Looking at the Washington U. data, we see that for documents 1 through 4, the client can receive good performance over 90% of the time by considering only 2 servers out of the group of 20. For document 0, the client would need to consider 5 servers to achieve good performance more than 90% of the time. On the other hand, the client at U. Mass. requires more servers to achieve good performance when fetching from

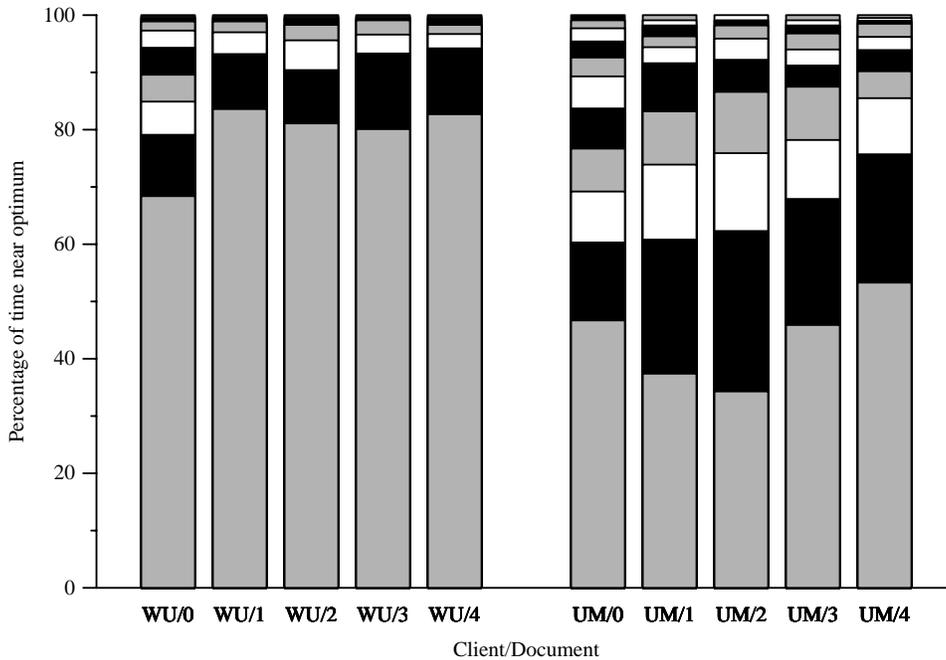


Figure 12: Server sets for two client-data combinations: Wash. U.'s Mars set and U. Mass.'s Apache

Apache servers. Seven servers are required for the first document while 5 are required for the other documents. This is a much higher proportion of servers than for the Washington U. client (7 out of 11 vs. 5 out of 20).

Figure 13 summarizes our findings over all documents. On average, less than half of all servers need to be considered to achieve good performance most (or even all) of the time. This result implies that when designing a server selection mechanism, it is unnecessary to assume that all clients will need to contact or evaluate all servers.

## 8 Server choice and document choice

The reader may have noticed that in Figure 12, the composition of server sets obviously varies from document to document. This seems to suggest that in some cases, a server that provides good performance for one document does not provide good performance for another document. However, further examination reveals that document choice has at best a weak effect on server choice.

Recall that a server set is the *smallest* set of servers that provide good performance for a given client. Other servers not in the server set could provide good performance at any given moment. For example, there are cases in which more than one collection of servers can be a server set. If two servers, A and B, provide good performance at exactly the same moment, then two server sets are possible: one using A and the other using B. Thus, it is unwise to rely on apparent differences in server sets as an indicator of differences in server performance.

Doc.	Avg. for 90% Good	Avg. for 100% Good
Mars (20 Servers)		
0	3.44	8.57
1	2.67	5.83
2	2.56	5.83
3	2.67	5.67
4	2.22	5.60
Apache (11 servers)		
0	3.89	6.25
1	3.00	5.20
2	3.11	5.25
3	3.00	5.80
4	3.00	6.00
News (16 servers)		
0	2.44	5.88

Figure 13: Average (taken over all clients) number of servers required to achieve good performance in 90% and 100% of fetches

$j$	$i$				
	0	1	2	3	4
0	0.55%	4.11%	4.11%	6.85%	6.85%
1	9.32%	0.00%	0.00%	0.82%	0.82%
2	9.86%	0.00%	0.00%	0.27%	0.27%
3	10.68%	0.55%	0.55%	0.00%	0.00%
4	8.77%	0.27%	0.27%	0.00%	0.00%

Figure 14: Percentage of time that good performance is not achieved using the top 5 servers from the server set of document  $i$  to fetch document  $j$

Figure 14 shows how using one document’s server set to fetch another document affects performance. The table was built by counting how often the top 5 servers from document  $i$ ’s server set are able to offer good performance for document  $j$  for every  $i, j \in [0, 4]$ . Though this data is generated from the Mars data at client site U. Va, all other combinations of clients and web sites produced similar results. The entry at column  $c$  and row  $r$  in the table is the percentage of fetches for which the server set for document  $c$  was **not** able to provide good performance for document  $r$ . For example, we can see that using the server set for document 4 would lead to good performance in over 99% of fetches.

We used only the top 5 servers from each server set so that all sets of servers considered would be the same size. Server sets for documents 2 through 4 only contained 5 servers, so they were unaffected. Document 0’s server set, however, contained 7 servers. The most immediate effect is that in the table above, the (truncated) server set for document 0 failed to provide good performance 0.55% of the time.

Measuring how well one document’s server set would do to fetch another is a much more reasonable way to judge the differences in server performance among documents. It can directly show how often a server identified as good for one document is not actually good for another document. In Figure 14, we can see that most often, performance remains good across server sets. Ignoring data from the first row and first column, we see that instances when one document’s server set does not offer good performance for another document are very rare.

Looking at the table’s first row and the first column, which correspond to document 0 and server set 0 respectively, we see that good performance is achieved less frequently. The servers which offer good performance for document 0 are at least partially different from the servers that offer good performance for other documents. This indicates that there might be some link between document choice and server choice. In all client-site combinations, we observed that the first document had a noticeably different set of good servers than the other documents.

In both the Apache and Mars data, the first document is also the smallest (about 2 KB). We believe the dependence is more a function of document size than the specific documents being fetched, but further study using a larger variety of documents is required to verify this. We can explain the effect of document size on server choice if we assume that the network (and not the server) is the bottleneck. For smaller documents, the transfer time depends more on the round trip time between the client and server. The smallest documents fit in one or two packets so the client-server conversation lasts only a few round trip times. For larger documents, the amount of bandwidth available on the path between the client and server becomes the important factor as the network “pipe” is packed with as much data as possible. In this scenario, one property of a server (the round trip time between it and the client) would dominate for small documents and another property (the throughput between the client and server) would dominate for larger documents.

Regardless of the cause, the effect is not extremely significant. First of all, at most 11% of fetches were adversely affected by the difference in server sets. In these fetches, the increase in transfer time was less than 25% above optimal on average. Also note that these performance penalties are on top of a rather small transfer time (about 1 second), so the actual penalties are on the order of hundreds of milliseconds. Thus there is little cause for concern over using only one server set for all document sizes will lead to bad performance.

## **9 Implications for server selection systems**

The observations about the properties of mirror servers that we have presented may be useful when designing server selection systems. However, our measurements were made in the absence of a selection mechanism. The introduction of a systematic way to evaluate servers may alter the performance of the servers significantly. For example, if all clients employ a load balancing algorithm, the correlation of performance among the servers may increase. Still, our observations do supply a picture of the network that can be used as a starting point in designing a system.

We have pointed out that the difference in performance from one mirror server to another is quite significant. This implies that choosing the right server has the potential to significantly improve client performance. We have also seen that most server sets usually contain more than one server, indicating that the best server for a given client changes over time. Server selection needs to take place periodically to achieve good performance. But because server sets are also

usually small, the server selection task is potentially a very lightweight operation.

We have observed that server rank changes do not depend significantly on time scale, implying that a ranking of servers from two hours ago is as useful as a ranking from two days ago. In other words, all performance results older than an hour are equally useful. Because of our experimental design, we cannot say anything about performance results younger than an hour.

For the News and Mars data sets, we have found that most rank changes are small, implying that a client may assume with a reasonable amount of confidence that a server which delivered good performance during the last fetch will have acceptable performance during the next fetch even if the two fetches are far apart in time. For these data sets, the benefits of server selection may be outweighed by the cost of evaluating servers. However, this does not hold for the Apache set, where ranks are less stable.

We have found a weak link between a change in a server's performance and a change in the server's rank. If the performance that a server can offer a client degrades massively, then it can be inferred that the server's rank has dropped and a new server should be selected for the client. However, for smaller performance drops, we cannot assume that a corresponding drop in rank has taken place.

Finally, protocols probably do not have to make allowances for picking a server based on the document that is being fetched. While we have noticed that there is a difference between the good servers for the smallest Mars and Apache documents and other documents' good servers, the difference in performance, both in relative and absolute terms, is not very large.

## 10 Conclusion

We have presented measurements of the performance of replicated web servers which have ramifications for server selection system designs. We have found that though transfer times are highly variable, server ranks are fairly stable over time. Further, clients can receive good performance if they use a carefully chosen subset of a group of mirrors. To further substantiate our conclusions, future work includes collecting longer traces, trying other mirror sets, and exploring shorter time scales.

The data collected for this study is available on the World Wide Web at <http://www.cs.cmu.edu/~acm/research/anycast.html>.

## 11 Acknowledgements

Peter Dinda has played a critical role in a large part of this work. He has contributed both a deep knowledge of statistics as well as a fine sense of humor to the enterprise. My advisor, Hui Zhang, has also been pivotal, offering detailed comments and both monetary and moral support.

We would like to thank the client sites (Berkeley, Georgia Tech., ISI, the University of Kentucky, the University of Massachusetts, the University of Texas, the University of Virginia, and Washington University in St. Louis), who gave us guest accounts and tolerated our experiments. We would also like to thank the server administrators for their support and cooperation.

## References

- [1] M. Arlitt and C. L. Williamson. Web server workload characterization: The search for invariants. In *Proceedings of ACM SIGMETRICS '96*, 1996.
- [2] H. Balakrishnan, S. Seshan, M. Stemm, and R. H. Katz. Analyzing stability in wide-area network performance. In *Proceedings of ACM SIGMETRICS Conference on Measurement and Modeling of Computer Systems*, June 1997.
- [3] E. Basturk, R. Engel, R. Haas, D. Kandlur, V. Peris, and D. Saha. Using network layer anycast for load distribution in the Internet. Technical report, IBM T.J. Watson Research Center, 1997.
- [4] S. Bhattacharjee, M. H. Ammar, and E. W. Zegura. Application-layer anycasting. Technical Report GIT-CC-96/25, Georgia Institute of Technology, 1996.
- [5] S. Bhattacharjee, M. H. Ammar, E. W. Zegura, V. Shah, and Z. Fei. Application-layer anycasting. In *Proceedings of INFOCOM '97*, 1997.
- [6] G. E. P. Box, G. M. Jenkins, and G. Reinsel. *Time Series Analysis: Forecasting and Control*. Prentice Hall, 3rd edition, 1994.
- [7] BrightTiger. ClusterCats. Available at <http://www.brighttiger.com/products/default.htm>.
- [8] R. L. Carter and M. E. Crovella. Dynamic server selection using bandwidth probing in wide-area networks. Technical Report BU-CS-96-007, Boston University, March 1996.
- [9] Cisco. DistributedDirector. Available at [http://www.cisco.com/warp/public/751/distdir/dd\\_wp.htm](http://www.cisco.com/warp/public/751/distdir/dd_wp.htm).
- [10] M. Crovella and A. Bestavros. Self-similarity in World Wide Web traffic: Evidence and possible causes. *IEEE/ACM Transactions on Networking*, 5(6):835–846, December 1997.
- [11] C. R. Cunha, A. Bestavros, and M. E. Crovella. Characteristics of WWW client-based traces. Technical Report BU-CS-95-010, Boston University, 1995.
- [12] C. Yoshikawa, B. Chun, P. Eastham, A. Vahdat, T. Anderson, and D. Culler. Using smart clients to build scalable services. In *Proceedings of USENIX '97*, January 1997.
- [13] Z. Fei, S. Bhattacharjee, E. W. Zegura, and M. H. Ammar. A novel server selection technique for improving the response time of a replicated service. In *Proceedings of INFOCOM '98*, March 1998.
- [14] P. Francis. A Call for an Internet-wide Host Proximity Service (HOPS). <http://www.ingrid.org/hops/wp.html>.
- [15] P. Francis, S. Jamin, V. Paxson, and L. Zhang. Internet distance maps (IDMaps). Available at <http://idmaps.eecs.umich.edu/>.

- [16] S. D. Gribble and E. A. Brewer. System design issues for Internet middleware services: Deductions from a large client trace. In *Proceedings of the USENIX Symposium on Internet Technologies and Systems*, 1997.
- [17] J. D. Guyton and M. F. Schwartz. Locating nearby copies of replicated Internet servers. Technical Report CU-CS-762-95, University of Colorado at Boulder, 1993.
- [18] R. Jain. *The Art of Computer Systems Performance Analysis*. John Wiley and Sons, Inc., 1991.
- [19] D. B. Johnson and S. E. Deering. Reserved IPv6 subnet anycast addresses. Internet Draft available from <http://www.ietf.org/internet-drafts/draft-ietf-ipngwg-resv-anycast-01.txt>, 1998.
- [20] B. N. Levine and J. J. Garcia-Luna-Aceves. Improving Internet multicast with routing labels. In *Proceedings of ICNP '97*, 1997.
- [21] WindDance Networks. WebChallenger. Available at <http://www.winddancenet.com/products/challenger.html>.
- [22] C. Partridge, T. Mendez, and W. Milliken. Request for comments 1546: Host anycasting service, November 1993.
- [23] J. Rosenberg and H. Schulzrinne. Internet telephony gateway location. In *Proceedings of INFOCOM '98*, March 1998.
- [24] J. Rosenberg, H. Schulzrinne, and B. Suter. Wide area network service location. Internet Draft available from <http://www.ietf.org/internet-drafts/draft-ietf-svrloc-wasrv-01.txt>, 1997.
- [25] M. Sayal, Y. Breitbart, P. Scheuermann, and R. Vingralek. Selection algorithms for replicated web servers. In *Proceedings of the Workshop on Internet Server Performance '98*, June 1998.
- [26] S. Seshan, M. Stemm, and R. H. Katz. SPAND: Shared passive network performance discovery. In *Proceedings of USITS '97*, 97.