# Tracking the Evolution of Web Traffic: 1995-2003[*]

Félix Hernández-Campos    Kevin Jeffay    F. Donelson Smith
*University of North Carolina at Chapel Hill*
*Department of Computer Science*
*Chapel Hill, NC 27599-3175 USA*
*http://www.cs.unc.edu/Research/dirt*

## Abstract

*Understanding the nature and structure of web traffic is essential for valid simulations of networking technologies that affect the end-to-end performance of HTTP connections. We provide data suitable for the construction of synthetic web traffic generators and in doing so retrospectively examine the evolution of web traffic. We use a simple and efficient analysis methodology based on the examination of only the TCP/IP headers of one-half (server-to-client) of the HTTP connection. We show the impact of HTTP protocol improvements such as persistent connections as well as modern content structure that reflect the influences of "banner ads," server load balancing, and content distribution networks. Lastly, we comment on methodological issues related to the acquisition of HTTP data suitable for performing these analyses, including the effects of trace duration and trace boundaries.*

## 1. Introduction

Since the mid-1990s, web traffic has been the dominant traffic type on the Internet, representing more bytes, packets, and flows on the Internet than any other single traffic class. Given this dominance, when performing network experiments and simulations involving end-to-end performance issues, it is essential that one consider both the effects of web traffic on the mechanism/protocol under study as well as the effects of the technology under study on the performance of web traffic.

However, since the mid-1990s the web has evolved from a simple hypertext document delivery system to a sophisticated client-server system for delivering a vast array of static and dynamic media. The HTTP protocol is now routinely used to deliver content once carried on more specialized application-level protocols. For example, the web is now the *de facto* user-interface for many remote data processing systems, commercial transactions, and email, news, and instant messaging systems. Our goal is to discover and document the evolving nature and structure of web traffic and, by doing so, inform researchers who need accurate characterizations of web traffic in order to simulate "Internet traffic."

We report on the analysis of nearly 1 terabyte of TCP/IP header traces collected in 1999, 2001, and 2003 from the gigabit link connecting the University of North Carolina at Chapel Hill (UNC) to its Internet service provider. In addition, we compare our results to smaller but similar measurements taken by other researchers in the 1995 to 1998 time frame. Beyond documenting the evolution of the web, we contribute to the simulation community:

- Empirical data suitable for constructing traffic generating models of contemporary web traffic,
- New characterizations of TCP connection usage showing the effects of HTTP protocol improvement, notably persistent connections, and
- New characterizations of web usage and content structure that reflect the influences of "banner ads," server load balancing, and content distribution.

A novel aspect of this study is a demonstration that a relatively light-weight methodology based on passive tracing of only TCP/IP headers from one direction of the TCP connection (*e.g.*, TCP/IP headers from packets flowing from web servers to web clients) was sufficient for providing detailed data about web traffic. These data were obtained without ever examining HTTP headers thus allowing us to address users' privacy concerns by simply anonymizing source and destination addresses in our traces. Moreover, unidirectional tracing of only TCP/IP headers greatly reduced the processing complexity and storage overhead of the tracing effort (compared to capturing HTTP headers).

In this paper we present retrospective and new analyses of web traffic and investigate some methodological issues concerning the acquisition of the data. Specifically, to assess the evolution of the web we present a comparison of the UNC data to those obtained in the seminal Mah [10], and Barford, Crovella, *et al.* [2-4, 7], measurement studies. We also report new web-page-level and user-level statistics including the use of primary and non-primary web servers to deliver content, the number of request/response exchanges (*i.e.*, the number of "objects") per page, and the number of page requests per user.

Finally, we present an analysis of the impact of the tracing duration on our ability to fully capture the complete distribution of a variety of statistics. The issue here is to determine (qualitatively) how sensitive various types of characterizations are to the length of the trace. This issue is important both for reasons of ensuring correctness/accuracy

of the computed distributions and for assessing the computation and storage requirements to acquire and process traces.

Our primary observations and conclusions are:

- The size of HTTP requests has been steadily increasing,
- The majority of HTTP responses are decreasing in size, while the very largest responses are increasing in size,
- Web page complexity, as measured by the number of objects per page and the numbers of distinct servers providing content per page, is increasing, and
- When measuring HTTP traffic on high-speed links, 90-second observation intervals are sufficient to capture reasonable estimates of most structural properties (*i.e.*, non-user behavioral properties) of HTTP connections.

In total these results demonstrate that usage of the web by both consumers and content providers has evolved significantly and make a compelling case for continual monitoring of web traffic and updating of models of web traffic. Moreover, our work demonstrates that such continual monitoring can be performed with very modest effort.

The remainder of this paper presents these data. Section 2 briefly reviews related work in web traffic modeling. Section 3 briefly reviews the tracing procedure and characteristics of the UNC network. Section 4 presents new distributions of web page and user-level statistics for data collected at UNC in 1999, 2001, and 2003. For a more historical (and methodological) perspective, Section 5 compares these UNC data with those obtained by Mah, Barford, and Crovella, *et al*. in 1995 and 1998. Section 6 presents the results of an analysis into the impact of tracing duration on the reported distributions.

## 2. Related Work

Web traffic generators in use today are usually based on data from the two pioneering measurement projects[1] that focused on capturing web-browsing behaviors: the Mah [10], Barford, and Crovella, *et al*., [2, 4, 7] studies. Traffic generators based on both of these sources have been built into the widely used *ns* network simulator [5] that has been used in a number of studies related to web-like traffic, *e.g.*, [8, 11]. These models have also been used to generate web-like traffic in laboratory networks [3, 6]. For both sets of measurements, the populations of users were highly distinctive and the sizes of the traces gathered were relatively small. Mah captured data from a user population of graduate students in the Computer Science Department at UC Berkeley. His results were based on analysis of approximately 1.7 million TCP segments carrying HTTP protocols. The measurement programs by Barford and Crovella reflected a user population consisting primarily of undergraduate students in the Computer Science Department at Boston University and in aggregate represented around 1 million references to web objects. In addition, both sets of data are now quite old. The Mah data were collected in 1995 and the Barford and Crovella, *et al*., data in 1995 and 1998. It is espe-

cially important to note that these studies were conducted before significant deployment of HTTP version 1.1 protocol implementations. For comparison, our study involved traces consisting of over 1.6 billion TCP segments generated by a user population of approximately 35,000 users and representing the transfer of almost 200 million web objects.

## 3. Data Sets Considered

The data used in our study were obtained using the methods described in [12]. Briefly, our analysis method consists of analyzing unidirectional traces of TCP/IP headers sent from web servers to clients (browsers) in order to infer application-level characteristics of HTTP. In particular, we exploit properties of TCP's sequence number and acknowledgement number increases to determine request and responses sizes [12].

Here we provide results for three sets of traces: one set taken in Fall 1999, one in Spring 2001, and one in Spring 2003. The Fall 1999 traces were collected during six one-hour sampling periods (8:30-9:30 AM, 11:00-12:00 noon, and 1:30-2:30, 4:00-5:00, 7:30-8:30, and 10:00-11:00 PM) over seven consecutive days. This set of 42 traces will be referred to "UNC 99." In the following, a "trace set" consists of all the TCP/IP headers collected during these sampling intervals.

The Spring 2001 traces were collected during three 4-hour sampling periods each day for seven consecutive days. The sampling periods were 8:00-12:00 noon, 1:00-4:00PM, and 7:30-11:30PM giving a total of 21 4-hour traces. This set of traces will be referred to as "UNC 01."

The third trace set ("UNC 03") consists of 56 one-hour traces collected at 5, 6, 10, and 11 AM, and 3, 4, 9:30 and 10:30 PM, also collected over seven consecutive days.

When the UNC 99 traces were gathered, our campus was connected to the ISP by an OC-3 (155 Mbps) full-duplex, ATM link. This link carried all network traffic between the campus and the "public" Internet. Traffic between the campus and Internet 2 sites was routed over a separate OC-3 link. We placed the monitor on the OC-3 link to the "public" Internet.

When the 2001 traces were collected 18 months later, the ISP link had been upgraded to OC-48 (2.4 Gbps) and used Cisco-proprietary DPT technology. Fortunately, all the traffic between the campus and the Internet traversed a single dedicated full-duplex Gigabit Ethernet link from the campus aggregation switch to the edge router with the DPT interface. In this configuration, both "public" Internet and Internet 2 traffic are merged in one Gigabit Ethernet link. For the 2001 and 2003 traces we placed a monitor on this link.

## 4. Analysis of TCP Connections Used for HTTP

### 4.1 Request and Response Data Sizes

Given the analysis methods described above, it is straightforward to compute empirical distributions of HTTP request and response sizes. Figures 1 and 2 give the cumulative distribution function (CDF) and complementary CDF
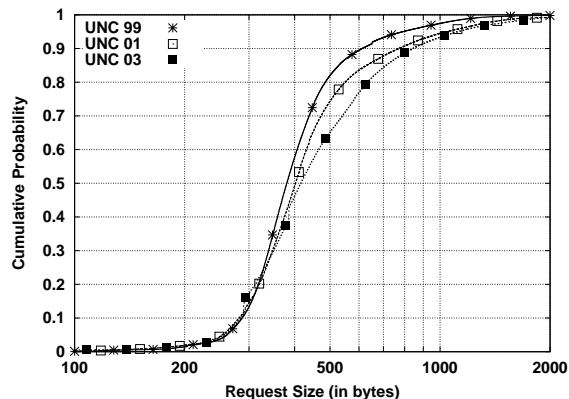
---

[1] Arlitt and Williamsom presented in [1] an earlier study of web traffic conducted using modest-size data sets.

**Figure 1:** Cumulative distribution of request data sizes (100 − 2,000 bytes).
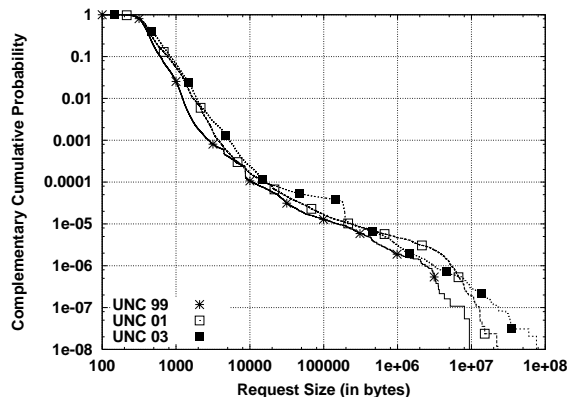


**Figure 2:** Complementary cumulative distribution of request data sizes greater than 100 bytes.
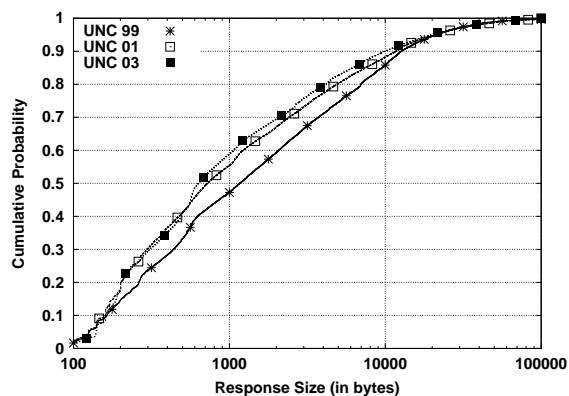


**Figure 3:** Cumulative distribution of response data sizes (100 − 100,000 bytes).
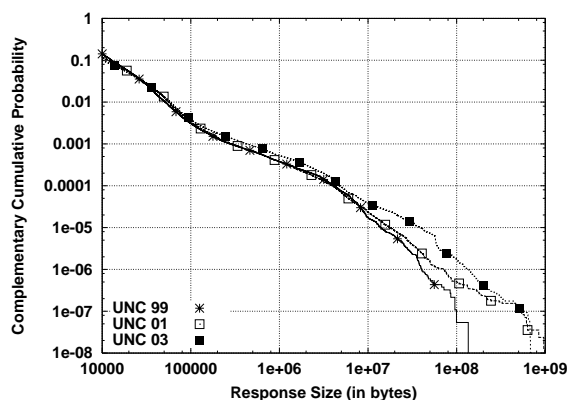


**Figure 4:** Complementary cumulative distribution of response data sizes greater than 10,000 bytes.

(CCDF), respectively, for HTTP request sizes while Figures 3 and 4 give the CDF and CCDF, respectively, for HTTP response sizes. These results confirm the observation that HTTP requests are becoming larger over time in both the body and tail of the distribution. In contrast, HTTP responses show a trend toward a larger proportion of smaller responses in the range of 200 bytes to 10,000 bytes. For example, in the 1999 traces about 47% of responses were 1,000 bytes or smaller while in the 2003 traces, about 59% of the responses were 1,000 bytes or less. The tail of the distribution shows a slight trend toward increased frequency of very large responses.

We also observe that the CCDFs for the 2001 and 2003 response sizes are approximately linear over nearly five orders of magnitude, which is consistent with a very heavy-tailed distribution. Note also that the tails of the response size distributions exhibit a systematic "wobbling," with knees around 100 KB and 3 MB. This phenomenon was studied in detailed in [9] using data from 2000 and 2001.

### 4.2  User and Web Content Characterizations

Because we do not have access to any of the HTTP protocol headers, we must use heuristics to *infer* characteristics of user and browser behavior from the analysis of TCP connections. The first step in this process is to aggregate TCP connection traces by unique client IP addresses. We then create a time-sorted summary of the TCP connection activ-

ity between each individual client and the server(s) that client used. We assume that in the vast majority of cases a client IP address identifies a single human user running one or more browser instances on a personal computer or workstation. Although we know that there are times when multiple users concurrently run browsers on a shared compute server (single client IP address), we believe this to be rare on our campus where there are basically one or more computers for each Internet user. Even though the vast majority of computers on our campus have IP addresses assigned by DHCP, we have confirmed that the reuse of a given IP address on different machines during a single trace is rare because leases last eight hours or more. Further, many of the larger DHCP servers maintain a fixed mapping of IP address assignments to Ethernet MAC addresses.

The time-sorted summary of TCP connections used by a client contains the connection start times, the server IP address for each connection, the beginning time and size in bytes of each request in all connections, the beginning and ending times and size of each response, and the ending time of the connection. We then use this time-ordered information to infer certain characteristics of the activity of each user or the browser software.

Using a heuristic approach similar to those developed originally by Mah [10] and Barford and Crovella [2-4], we attempted to identify points in each client's activity that are likely to mark a request for a new (or refreshed) page. We

3

use the term "page" as a convenient label for a web object referenced in a "top-level" sense, *i.e.*, not referenced through interpreting references found internal to some other object (*e.g.*, embedded references in HTML). We also use the term "object" synonymously with a response from a web server. Server responses that are error reports (*e.g.*, "404 – Not found") or responses to conditional-GET requests (*e.g.*, "304 – Not Modified") are counted as objects (or pages) in this discussion. We assume that page references normally occur after some period of idle or "think" time at the client, *e.g.*, the time a user spends digesting the contents of one browser display and selecting (or entering) a link to a new page. This same model of a page request following an idle period also captures the behavior of periodically refreshed pages.

We define an idle period heuristically by examining the time-ordered set of TCP connections used by a client. We identify periods in which the client either has no established TCP connections or where no established connection has an active request/response exchange in progress. We consider a request/response exchange to be active from time the request begins until the corresponding response ends. If any such period persists for longer than a time threshold, it is classified as an idle period. We found empirically that a threshold of 1 second works well for distinguishing idle periods (as did Mah and Barford and Crovella). It is important to note that this approach works only on traces for which we can be reasonably certain that *all* the TCP connections for a given browser appear in the traces.

We consider the initial request/response exchange following an idle period to be for the "top-level" page object (typically HTML) and all the subsequent request/response exchanges before the next idle period to be for the "embedded" object references (typically images) within the initial page object. The server IP address involved in the request/response exchange for the top-level page object is considered to be the *primary* server for the page. All server IP addresses *not equal to the primary IP address* involved in subsequent request/response exchanges for objects related to that page are considered to be *non-primary* servers. Embedded objects may come from either the primary or non-primary server(s).

Using this approach we obtained a number of distributions that characterize user browsing activities and web content. Distributions for user "think" time and the number of unique server IP addresses per page are reported in [12] for the 1999 traces and are not repeated here for space considerations. The distribution of server IP addresses reflects the ways page content is obtained dynamically from a number of sources including advertisements from agency sites and explicit content distribution services (*e.g.*, Akamai).

We now report briefly on additional interesting results that provide some insight into how users access web objects and how web content is organized and distributed among servers. In interpreting these data it is important to keep in mind that all the web traffic we observed represents only objects that were not be obtained from local browser caches. The six-year time span of our three data sets is used to draw contrasts in the evolution of web traffic. In addition, in Section 6 we illustrate the effect of longer tracing intervals by comparing the 1999 and 2003 data sets (one-hour-long traces) with the 2001 data set (4-hour-long traces).

Figures 5 and 6 give the CDF and CCDF, respectively, of the number of top-level page requests per unique client (browser) IP address observed in the traces. The same IP address appearing in different traces is counted as a different address and hence as a different "user instance." We counted 140,522 client addresses for the 1999 traces, 238,287 for 2001, and 532,555 for 2003. If we assume that a unique client IP address represents a single "user instance," by the reasoning given earlier, we can infer characteristics of user browsing activity. A slight majority of users were observed to make more than 10 top-level page requests during the one-hour tracing intervals used in 1999. There were some users, however, (about 5%) that made more than 100 page requests in one hour. Notice that these distributions change noticeably for 2001 because we traced for 4-hour intervals. For example, in 2001 we found that about 65% of the identified users request more than 10 pages in a 4-hour interval and 15% request more than 100 pages. Figure 6 likely indicates the existence of pages that are automatically refreshed at very short intervals (but greater than one second — the threshold for identifying page boundaries).

Figures 7 and 8 give the CDF and CCDF, respectively, of consecutive top-level page requests by a given user (unique client IP address) to the *same* primary server IP address. The plot of the body of the distributions shows that an increasingly large percentage of consecutive top-level page references made by the same user go to a different server IP address than the immediately prior reference (from 68% in 1999 to 77% in 2003). We believe these results reflect the basic organization of web servers for a web site into "server farms" for load balancing and content distribution. The presence of a heavy tail in the distribution of consecutive page requests to the same server is consistent with the earlier observation of automatically refreshed pages.

Web page structure is reflected to some extent in Figure 9 that shows the distribution of objects (top-level plus embedded) per page as inferred from our analysis. While the proportion of quite simple pages appears to be large (over 40% of pages have no embedded objects and 67-75% of pages are composed of three or fewer objects), there are significant numbers of pages with complex structure. It is again important to keep in mind that many objects in pages that a user views often may be cacheable (in browser caches) and not be included in our traces.

Because embedded page content is often from a number of sources including advertisements from agency sites and explicit content distribution sites, it is interesting to see if we find differences in the characteristics between primary and non-primary servers. Recall that a primary server is defined as the server from which the top-level page is requested and a non-primary server is any server other than the primary server from which embedded content is obtained. Figure 10 gives the distribution of the number of objects per page (top-level plus embedded) from primary and non-primary servers (the top-level object comes from the primary server by definition). We find that about 76% of pages in the 1999
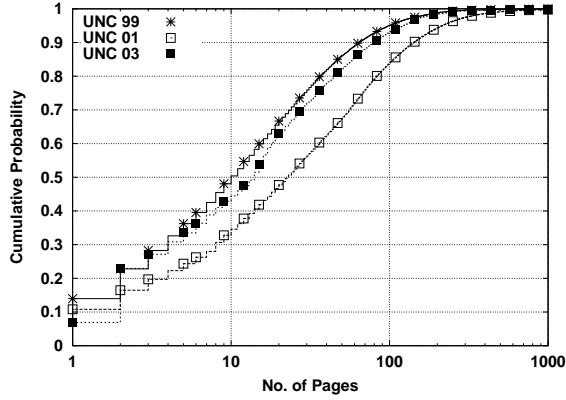
**Figure 5**: Cumulative distribution of the number of page requests made from unique client IP addresses during a tracing interval (1 hour in 1999 and 2003, 4 hours in 2001).
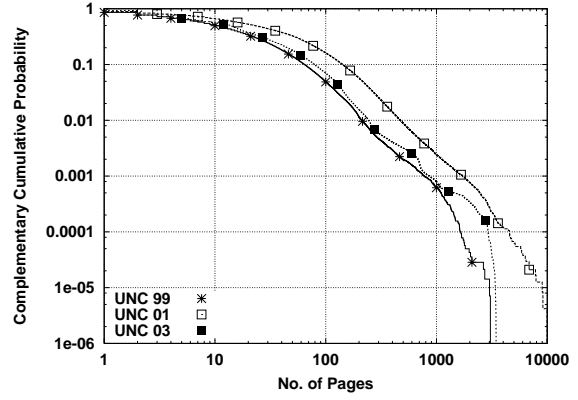


**Figure 6:** Complementary cumulative distribution of the number of page requests made from unique client IP addresses during a tracing interval.
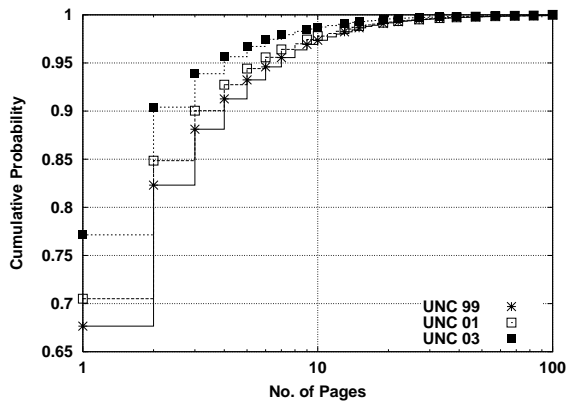


**Figure 7:** Cumulative distribution of the number of consecutive page requests by one user to the same primary server IP address
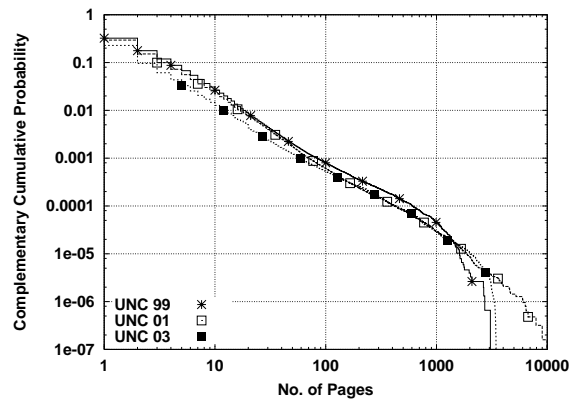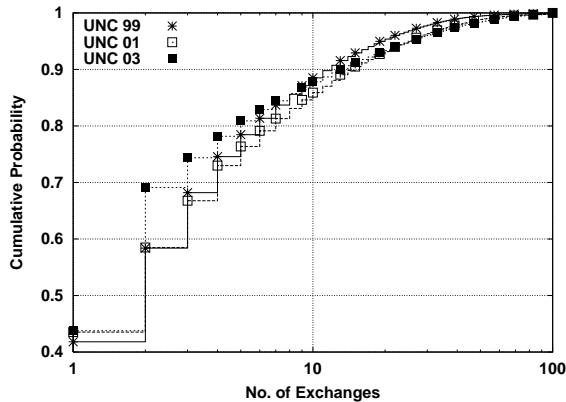


**Figure 8**: Complementary cumulative distribution of the number of consecutive page requests by one user to the same primary server IP address.



**Figure 9**: Cumulative distribution of the number of objects (top-level plus embedded) per page.



**Figure 10:** Cumulative distribution of the number of objects requested from primary and non-primary servers.

and 2001 data sets require only one object from the primary server, and this percentage increases to 84% in 2003. We also observe that around 40% of pages required more than one object from a non-primary server, further reflecting the popularity of multi-server website and content-distribution networks. Furthermore, we see an increase in the number of pages that have a large number of embedded objects from non-primary servers. For example, in 1999, 34% of pages

required more than one object from a non-primary server, compared with 40% in 2001 and 2003.

Figures 11 and 12 give the CDF and CCDF, respectively, for request sizes from primary and non-primary servers. Figures 13 and 14 show the corresponding response sizes. We find only minor differences in request and response sizes between primary and non-primary servers and all of them are in the tail of the distributions. More very large requests are sent to primary servers compared to non-primary
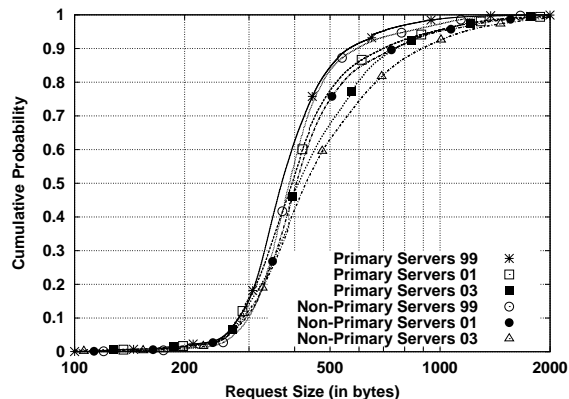
**Figure 11:** Cumulative distribution of request data sizes greater than 100 bytes for primary and non-primary servers.
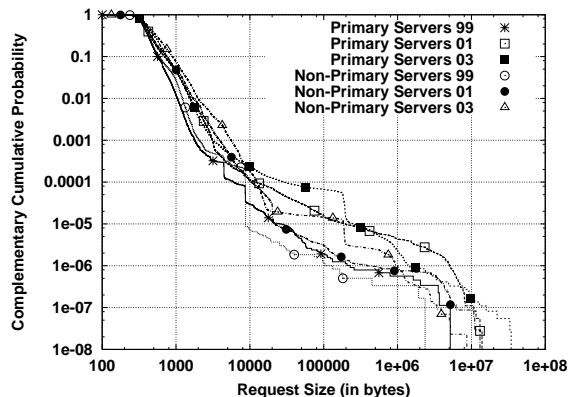


**Figure 12**: Complementary cumulative distribution of request data sizes greater than 100 bytes for primary and non-primary servers.
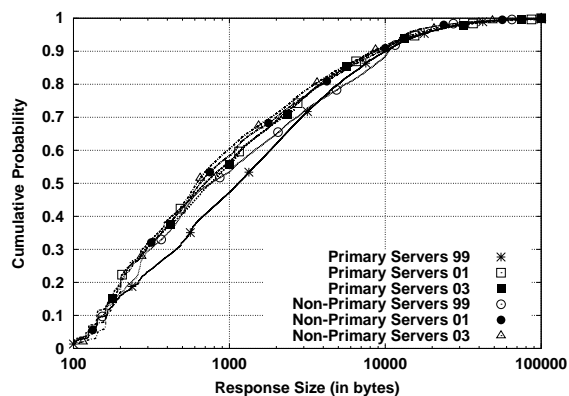


**Figure 13**: Cumulative distribution of response data sizes (100 – 100,000 bytes) for primary and non-primary servers.
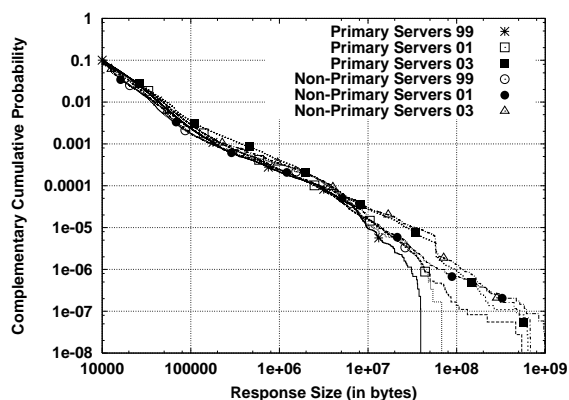


**Figure 14:** Complementary cumulative distribution of response data sizes greater than 10,000 bytes for primary and non-primary servers.

servers, while non-primary servers appear to have a slightly greater proportion of very large responses.

### 4.3 Limitations of the Methodology

Our analysis methods are based on making inferences from the limited information available in the TCP/IP proto- col header for one direction of a TCP connection. There are a number of inherent limitations and uncertainties that arise when making these inferences. However, the degree of un- certainty in the results is not uniform. For characterizations of TCP connection-level properties such as the sizes and numbers of request/response exchanges, the methodology should produce very good results. For other characterizations of the Web, especially those that depend on identifying pages or classifying requests and responses as belonging to primary or non-primary servers, there is greater uncertainty. We have identified four classes of issues that contribute to uncertainty in the results: pipelined exchanges, user/browser interactions (such as using the "Stop" and "Reload" browser buttons), browser and proxy caches, and TCP segment proc- essing to deal with packet loss, duplication and reordering in the network. Each of these is discussed fully in [12].

## 5. Comparison with the Mah, Barford, and Crovella, *et al.*, Studies

In this section we compare our empirical distributions with the published results from the Mah, Barford, and Crov-

ella, *et al.*, studies. Barford, *et al.*, presented in [4] two data sets of HTTP traffic, collected in 1995 ("W95") and 1998 ("W98"). (The well-known SURGE traffic generator for web workloads [2] is based on the W98 data.) The data from 1995 was acquired by instrumenting the web browsers in a computer lab used by computer science students at Boston University. The 1998 data set was acquired by instrument- ing a transparent proxy server used by a similar group of students at Boston University. Mah's data was acquired in September 1995, using packet tracing on an Ethernet seg- ment in the Computer Science department at the University of California at Berkeley. Both studies reflect relatively small and homogeneous populations of users.

A common element in all three studies is the distribu- tion of response sizes. Summary statistics for these distri- butions are given in Table 1. The number of samples con- sidered in our distributions is two orders of magnitude larger than that of Barford and Crovella, and four orders of magni- tude larger than that of Mah. The maximum response size observed in our study is significantly larger.

Barford, *et al*. present in [4] hybrid lognormal-Pareto models of response sizes for both data sets that accurately match their empirical data. Figure 15 compares the body of these models (the lognormal part) with the distributions obtained from our data. We believe the striking differences seen in this figure are a clear reflection of how web objects are evolving over time and make a clear statement about the
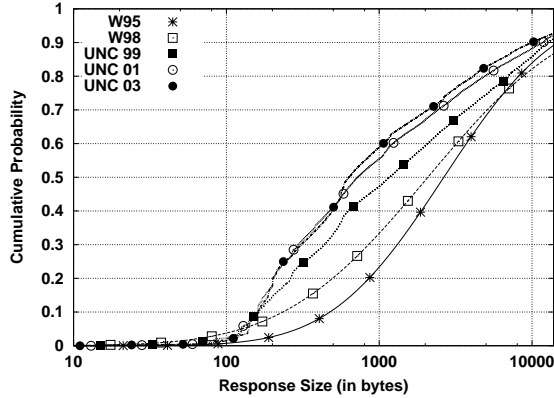
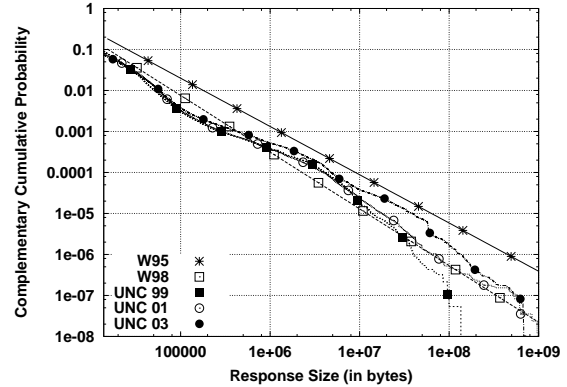**Figure 15:** Cumulative distribution of response data sizes, SURGE *v.* UNC traces.



**Figure 16:** Complementary cumulative distribution of response data sizes greater than 10,000 bytes, SURGE *v.* UNC traces.
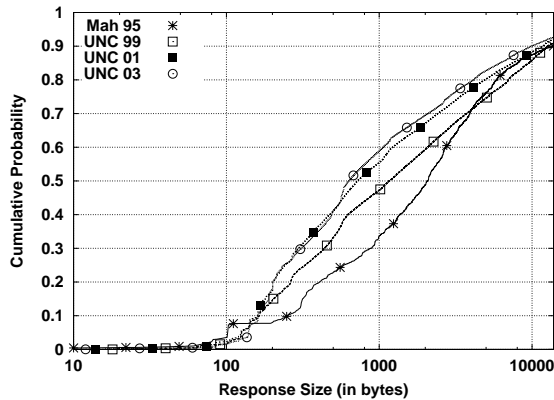


**Figure 17:** Cumulative distribution of response data sizes, Mah *v.* UNC traces.
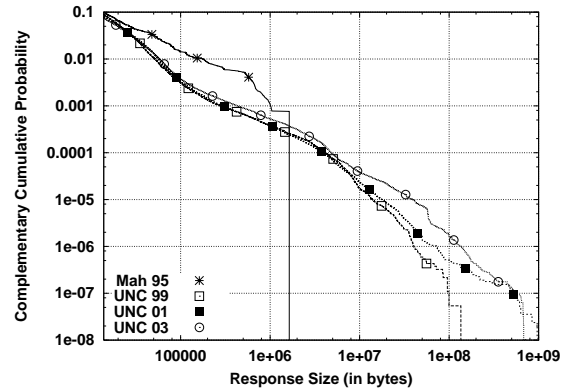


**Figure 18:** Complementary cumulative distribution of response data sizes greater than 10,000 bytes, Mah *v.* UNC traces.

need for continual updating of web models. For instance, the percentage of responses that had a length of 1,000 bytes or less increased from 23% as reported in W95 to 60% as reported in our 2003 data with the W98, UNC99, and UNC01 results falling between them. A similar result is presented in Figure 17, which compares the empirical distribution of response sizes in Mah's study with UNC data.

Barford, *et al.* model the tail of the distribution of response sizes using Pareto distributions with different parameters for each of the sets. Figure 16 compares these tails with UNC data. Our data matches the analytic model they fit to the W98 data remarkably well. It is interesting to note that this model matched their empirical data up to sizes of $10^5$ (based on their observations of response sizes slightly larger than about $10^6$ bytes). It did, however, predict very

accurately the heavy tail we find in our data up to sizes of $10^9$ bytes. Considering the differences in the sizes of the data sets, the user populations, and the years the data were obtained, this is a strong confirmation of their results. However, this analytical fitting does not capture the wobbling of the tail, that is more properly modeled using a mixture of distributions, as described in [9]. The comparison with the tail of Mah's distribution as given in Figure 18 shows that the heavy tail properties of response sizes were not adequately captured in his empirical distributions.

Another set of distributions common to the three studies is the number of top-level and embedded objects per page. The analytical distributions derived by Barford and Crovella and used in two different versions of the SURGE workload generator [2, 3] are compared with our results in Figure 19. (Their analytical distributions are Pareto so we evaluated them for integer values.) Their two distributions are quite different from each other, with the first one (SURGE98) being much lighter than the second (SURGE99). The model in SURGE99 [3] is remarkably close to our data for values beyond ten objects per page but differs substantially for 5 or fewer objects per page. This may reflect some differences in the methods and heuristics used to identify pages. The corresponding distribution from Mah's study, shown in Figure 20, differs substantially from the empirical distribu-

**Table 1:** Summary data for response size distributions. (All sizes are in bytes.)

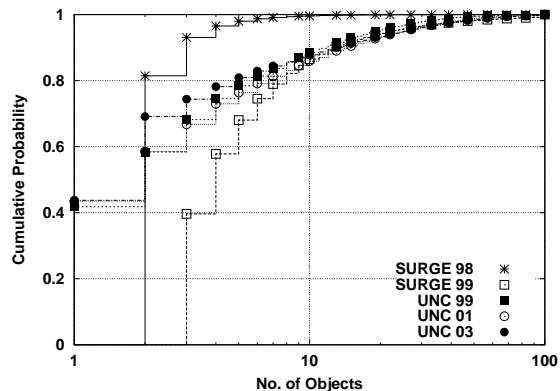| Data Set | Sample Size (Number of responses) | Min Response Size | Max Response Size | Mean Response Size | Median Response Size |
|---|---|---|---|---|---|
| W95 | 269,811 | 3 | 20,135,435 | 14,826 | 2,245 |
| W98 | 66,988 | 1 | 4,092,928 | 7,247 | 2,416 |
| Mah 95 | 5,300 | 62 | 8,146,796 | 10,664 | 2,035 |
| UNC99 | 18,526,201 | 1 | 135,294,044 | 6,734 | 1,164 |
| UNC01 | 84,343,238 | 1 | 984,871,070 | 6,397 | 722 |
| UNC03 | 96,836,703 | 1 | 718,067,386 | 7,296 | 632 |

7

**Figure 19:** Cumulative distribution of number of objects per page, SURGE *v.* UNC traces.
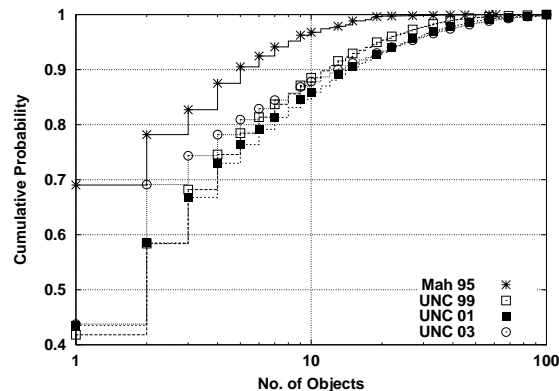


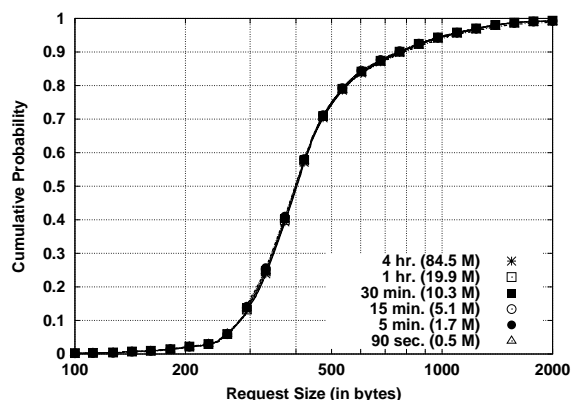**Figure 20:** Cumulative distribution of number of objects per page, Mah *v.* UNC traces.



**Figure 21:** Cumulative distribution of request data sizes greater than 100 bytes for sub-sampled traces. The legend shows the number of sample values in each empirical distribution.
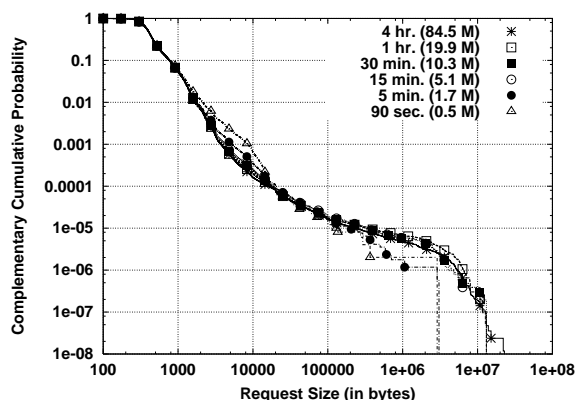


**Figure 22**: Complementary cumulative distribution of request data sizes greater than 100 bytes for sub-sampled traces. The legend shows the number of sample values in the empirical distribution.

tions computed from our traces for all values. This difference likely reflects the increase in complexity in layout of web pages that took place in the past 8 years.

## 6. Sampling Issues

In collecting traces for large-scale studies of Internet traffic there are important issues related to the number and duration of tracing intervals. All traffic between the Internet and the UNC-CH campus passes through a single Gigabit Ethernet link that is shared by a large user population. When we are tracing on this link, storage for traces is a potential concern. For example, each 2001 trace spanned a 4-hour interval. During heavy traffic periods, a single trace of just TCP/IP headers in one direction (68 bytes per packet) consumes over 30 Gigabytes of storage. In contrast to our 4-hour traces, most of the NLANR/MOAT trace collection [13] is for 90-second intervals. For tracing the UNC-CH Internet link, a 90-second interval would require storing only about 200 MB for each of the inbound and outbound traces, a very substantial reduction.

For this reason, we decided to analyze the effects on the quality of our results if we were to use shorter tracing intervals. We examined three inter-related issues:

- Can we obtain a sufficiently large sample with a small number of short traces?

- How does the length of tracing intervals affect the overall empirical distribution shapes?
- Should we include in the empirical distributions the data from incomplete TCP connections at the beginning and end of traces and does the length of the tracing interval matter for deciding?

Our 2001 trace collection is comprised of 21 4-hour traces taken at three intervals on each of seven days. Each of the 4-hour traces was then sub-sampled by taking its initial one-hour slice (using *tcpslice* to truncate the *tcpdump* trace) to produce another set of 21 traces. The one-hour traces were then sub-sampled by taking the initial 30-minute slice to produce a third set of 21 traces. This procedure was repeated to obtain the initial 15-minute slices of the 30-minute traces, the initial 5-minute slices of the 15-minute traces, and the initial 90-second slices of the 5-minute traces. Each set of 21 traces was then processed separately to compute some of the distributions described above. Figures 21 and 22 give the CDF and Complementary CDF (CCDF), respectively, for request sizes while Figures 23 and 24 give the CDF and CCDF, respectively, for response sizes.

Contrary to our expectations, we found that the 90-second traces produce results for these distributions that are virtually indistinguishable from the 4-hour traces up to the 99.5 percentile for request sizes and the 99.95 percentile for response sizes. Clearly, if one is interested in the details of
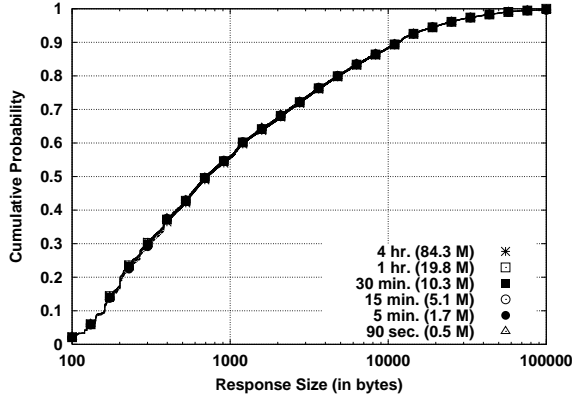
**Figure 23**: Cumulative distribution of response data sizes (100 – 100,000 bytes) for sub-sampled traces. The legend shows the number of sample values in each empirical distribution.
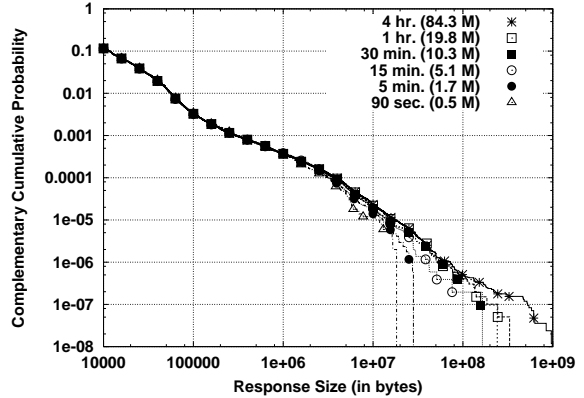


**Figure 24**: Complementary cumulative distribution of response data sizes greater than 10,000 bytes for sub-sampled traces. The legend shows the number of sample values in each empirical distribution.
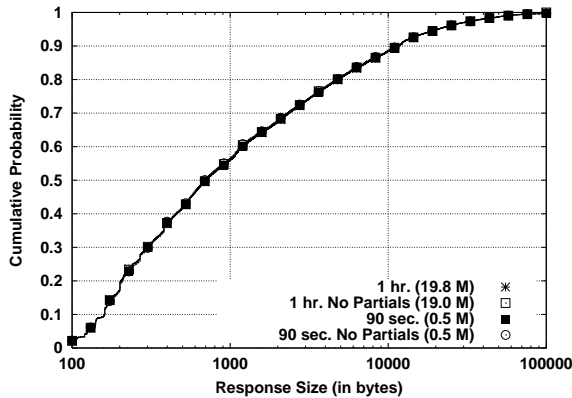


**Figure 25:** Cumulative distribution of response data sizes (100 – 100,000 bytes) for complete and partial responses. The legend shows the number of sample values in each empirical distribution.
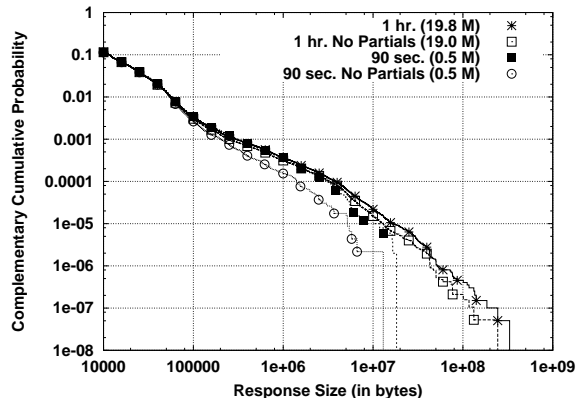


**Figure 26**: Complementary cumulative distribution of response data sizes greater than 10,000 bytes for complete and partial responses. The legend shows the number of sample values in each distribution.

the distribution tail (the extreme values), occasionally tracing for longer intervals may add value. Note that because we are tracing a high-speed link shared by a large user population, we were able to obtain about 500,000 sample values in 21 traces of 90-seconds duration. It appears reasonable to conclude that much shorter traces of 5 minutes or even 90 seconds are appropriate so long as they contain an adequate number of sample values. Accurate modeling of the "moderate" tail (see [9]) of the distributions does require larger tracing intervals, but 15 minutes seems to suffice provided a large sample is obtained (*e.g.*, by aggregating a number of intervals as was done in Figures 21-24).

We next compared the results for the response size distributions for the one-hour trace collection and the 90-second trace collection by altering the way truncated TCP connections at the beginning and end of the trace are treated. In one case we retained the data for incomplete responses and in the other we used only those responses known to be complete (both the initial SYN and the ending FIN for the connection were in the trace). Figures 25 and 26 show these results. Again we found that there was some effect only beyond the 99.5 percentile and, furthermore, counting the partial responses for the 90-second traces actually produces a distribution slightly closer to the one-hour traces.

Finally we should note that these observations hold only for distributions that are not used to characterize user activities over time. For example, Figures 27 and 28 show the distribution of page requests per unique user during a tracing interval (identified as unique client IP addresses as described above). Clearly, the length of the tracing interval dramatically alters the characterizations one obtains about any user level activities that span significant amounts of time.

## 7. Summary and Conclusions

Accurate, contemporary characterizations of web traffic are essential for simulations involving realistic Internet traffic. We have reported data on the structure and makeup of web traffic based on a comprehensive study of the usage of the web by the 35,000 person UNC user community. Our method has been to capture unidirectional traces of only the TCP/IP headers flowing from the Internet to the UNC campus on a Gigabit Ethernet link. To date we have acquired over a terabyte of headers and are using these data to construct an (evolving) empirical model of web traffic.

In addition we have performed a retrospective analysis of web traffic comparing UNC data from 1999, 2001, and 2003 with similar data obtained by Mah in 1995 and Barford and Crovella, *et al*. in 1995 and 1998. The results hold nu-
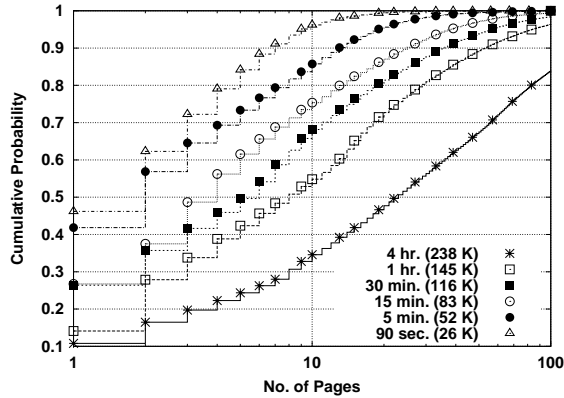
**Figure 27**: Cumulative distribution of the number of page requests made from unique client IP addresses during a tracing interval for sub-sampled traces.
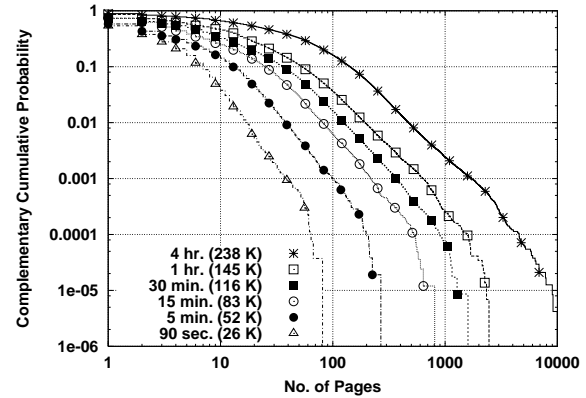


**Figure 28**: Complementary cumulative distribution of the number of page requests made from unique client IP addresses during a tracing interval for sub-sampled traces.

merous insights into the evolution of the web and web content; chief among them being:

- The size of HTTP requests has been steadily increasing. This likely reflects the evolution of the use of the web from a simple vehicle for requests for content to a means for uploading forms, data, files, email (with attachments), and other objects of arbitrary size.

- The majority of HTTP responses (those in the main body of the response-size distribution) are decreasing in size. This likely reflects the increase in the frequency of both small objects (adds, page decorations) as well as non-content server responses (error messages, conditional-GET replies *etc.*).

- The largest HTTP responses (those in the tail of the response-size distribution) are increasing in size. This likely reflects the natural evolution of object sizes delivered via HTTP (*e.g.*, MP3 files, software distribution CD-images, DVD segments, *etc.*).

- Web page complexity, as measured by the number of objects per page and the numbers of distinct servers providing content per page, is increasing. Pages have more objects and more servers (be it through content distribution networks or server farms) are involved in delivering content to the average web page.

In addition, from a methodological standpoint, we have determined that when tracing on a high-speed link shared by a large population (tens of thousands) of users, surprisingly short traces, as short as 90 seconds, produce distributions for many (non-user-behavioral-related) measures of web traffic that are indistinguishable from those obtained from 4-hour traces.

In total these results demonstrate that usage of the web by both consumers and content providers has evolved significantly and make a compelling case for continual monitoring of web traffic and updating of models of web traffic. Moreover, our work demonstrates that such continual monitoring can be performed with very modest effort.

## 8. References

[1]  M. Arlitt and C. Williamson, *A Synthetic Workload Model for Internet Mosaic Traffic*, Proc., 1995 Summer Computer Simulation Conference (SCSC`95), Ottawa, Canada, July 1995, pp. 852-857.

[2]  P. Barford and M. E. Crovella, *Generating Representative Web Workloads for Network and Server Performance Evaluation*, Proceedings of ACM SIGMETRICS '98, 1998, pp. 151-160.

[3]  P. Barford and M. E. Crovella, *A Performance Evaluation of HyperText Transfer Protocols*, Proceedings of ACM SIGMETRICS '99, May 1999, pp. 188-197. (Extended version available as Boston University Technical Report BU-TR-98-016.

[4]  P. Barford, A. Bestavros, A. Bradley, and M. E. Crovella, *Changes in Web Client Access Patterns: Characteristics and Caching Implications*, World Wide Web, Special Issue on Characterization and Performance Evaluation, Vol. 2, 1999, pp. 15-28.

[5]  L. Breslau, D. Estrin, K. Fall, S. Floyd, J. Heidemann, A. Helmy, P. Huang, S. McCanne, K. Varadhan, Y. Xu, and H. Yu, *Advances in Network Simulation*, IEEE Computer, vol. 33 no. 5, May 2000, pp. 59-67.

[6]  M. Christiansen, K. Jeffay, D. Ott, and F. D. Smith, *Tuning RED for Web Traffic*, Proceedings of ACM SIGCOMM 2000, September 2000, pp. 139-150.

[7]  M. Crovella, and A. Bestavros, *Self-Similarity in World Wide Web Traffic: Evidence and Possible Causes*, IEEE/ACM Transactions on Networking, vol. 5, no. 6, December 1997, pp. 835-846.

[8]  W. Feng, D. Kandlur, D. Saha, K. Shin, *Blue: A New Class of Active Queue Management Algorithms*, University of Michigan Technical Report CSE-TR-387-99, April 1999.

[9]  F. Hernández-Campos, J. S. Marron, G. Samorodnitsky, and F. D. Smith, *Variable Heavy Tailed Durations in Internet Traffic, Part I: Understanding Heavy Tails*, Proc., 10[th] IEEE/ACM Intl. Symp. on Modeling, Analysis and Simulation of Computer and Telecommunication Systems (MASCOTS), October 2002, pp. 43-52.

[10] B. Mah. *An Empirical Model of HTTP Network Traffic*, Proc. IEEE INFOCOM '97, April 1997, pp. 592-600.

[11] T. Ott, T. Lakshman, and L. Wong, *SRED: Stabilized RED*, Proceedings IEEE INFOCOM '99, March 1999, pp. 1346-1355.

[12] F.D. Smith, F. Hernandez Campos, K. Jeffay, D. Ott, *What TCP/IP Protocol Headers Can Tell Us About the Web*, Proc. ACM SIGMETRICS '01, June 2001, pp. 245-256.

[13] *http://moat.nlanr.net/Traces*