

Hybrid TCP-UDP Transport for Web Traffic

Israel Cidon, Amit Gupta, Raphael Rom, and Christoph Schuba

Sun Microsystems Laboratories
901 San Antonio Road
Palo Alto, CA 94303

1 Introduction

As the *World Wide Web* (WWW) and other Internet applications, such as real-time audio and video become more and more popular, both the Internet itself and the most popular sites are suffering from severe congestion. This congestion is perceived by users as service delays at best, and as a lack of service at worst.

The implications for companies depending on the availability of the service are grave: *Connection timed out* is read as *Closed for business*. Given the revenue-generating and often mission-critical service that many Web sites provide, high availability is of great importance and desirable to be maintained, even in the presence of multiple system failures.

Most of the Web traffic today uses the *HyperText Transfer Protocol* (HTTP). Traditionally, HTTP uses the *Transmission Control Protocol* (TCP) [15] as its underlying transport protocol. TCP provides several important services to HTTP, including reliable data transfer and congestion control. Unfortunately, TCP is poorly suited for the short conversations that comprise a significant component of Web traffic. The overhead of setting up and tearing down TCP state amortizes poorly for these small connections.

Many Web sites today use a variety of approaches, including server and content replication and distribution, to improve service latency and reliability. Such features are best served if the HTTP address (URL) is looked at and processed before the final destination server is selected. Current HTTP *redirection* solutions require tearing down and setting up multiple TCP connections. The redirection operation itself is a short-lived HTTP connection. Several commercial products such as the BrightTiger¹ use HTTP redirection for content distribution at the added expense of setting up multiple TCP connections.

Others avoid the overhead associated with HTTP redirection but, consequently, suffer from other limitations such as the need to fully replicate the server content (if URL is not looked at before the destination is selected,) or the need to transfer TCP connections with their associated states between servers (see for example Resonate at <http://www.Resonate.com/> or the Cisco LocalDirector at <http://www.cisco.com/warp/public/751/lodir>).

This paper takes a different approach and tackles the problem at its roots. We propose a dynamic transport protocol selection for reducing service latency, traffic overhead, and server load by giving services a dynamic choice of which transport protocol to use. Today, Web services usually use a single transport protocol, namely TCP, for all traffic; dynamically selecting the underlying transport protocol (e.g., TCP or UDP) can substantially improve service latency and reduce network traffic and server load. Our approach is conservative in the sense that we analyze and test our proposal against current Web traces. The rationale is that in the future load, content, and geographical distribution of servers within a service provider domain (say a video on demand or network computing services etc.) will further increase the need for redirection operations. Similarly, the increasing use of

¹<http://www.brighttiger.com/>

style sheets as well as the increasing use of cache validation also lead to a large number of short HTTP transactions. This will emphasize the overhead reduction even further.

The mechanisms discussed in this paper apply to many client-server protocols; we concentrate on the *HyperText Transport Protocol* (HTTP) [3] as the canonical example.

Available Internet traffic traces [4] show that a large part of Web traffic consists of short HTTP transactions: about 40% of Web transfers can fit in a single 1500 byte datagram, the size of an Ethernet *maximum transmission unit* (MTU), and approximately 63% will fit into two datagrams. In the case of a single MTU, using TCP requires 4-5 times the number of packets compared to the use of UDP. TCP setup also requires more complex state machine operations and setup of data structures at both the server and the client. Consequently, using UDP reduces network traffic as well as client-observed latency.

This hybrid TCP-UDP scheme combines the best of both worlds: for short transactions, it benefits from stateless UDP's low overheads; on the other hand, for large data transfers, it obtains the desired reliability, resequencing, and congestion control from TCP. The hybrid scheme is attractive because it only requires application-level changes, while the operating system kernel code can remain unchanged. Finally, the scheme is incrementally deployable in the current Internet and is fully compatible with the installed base.

Previous research investigated the cost of high TCP overhead for small connections. We address here two of these schemes: T/TCP [1] and HTTP Performance modeling research by Heidemann, Obraczka, and Touch [5].

Another motivation for the design of the hybrid scheme is to guarantee that our scheme will preserve the "TCP-friendly" property. This means that it should not attempt to benefit a particular service (in terms of better response time or throughput) at the expense of the global network or other well behaving users. We maintain the TCP congestion control for all long-term connections. On the other hand, since TCP congestion control is not effective for short-lived transactions (because of the lack of proper round trip adaptation periods), the use of UDP in these instances does not make a difference in that respect.

TCP for Transactions (T/TCP) [1] was developed as a transport protocol for request/reply type message passing protocols. It avoids the overhead that explicit connection setup and tear-down phases impose for small transactions. The designers of T/TCP defined a reliable request/response handshake with exactly one packet in each direction. T/TCP is designed as a backward-compatible extension to the TCP protocol. For Web traffic, T/TCP behavior and performance is similar to that of persistent-HTTP: the first connection between a pair of hosts requires a three-way handshake while successive connections avoid this overhead. T/TCP's requirement for a kernel-level implementation is the primary reason why T/TCP is not widely available[10].

Heidemann et al., [5] evaluates the performance of HTTP over UDP. Unlike our scheme (which uses TCP if needed), [5] relies entirely on using UDP as the transport protocol. UDP was augmented with adaptive retransmission and congestion control mechanisms, thereby mimicking the behavior of TCP similar to the *Asynchronous Reliable Delivery Protocol* (ARDP) [9]. The disadvantage of such a scheme is that we need to develop yet another complex protocol with reliability, resequencing and congestion avoidance duplicating the long tedious process of TCP development. Moreover, if this is not a kernel-level implementation, it might result in a slower and less efficient operation.

The remainder of this paper is organized as follows. Section 2 provides background material and motivation. Section 3 explains our hybrid TCP-UDP scheme and discusses its strengths and weaknesses. In section 4, we evaluate the performance of this hybrid scheme via simulations. We conclude in section 5.

2 Background

In this section, we briefly describe the various transport protocols that are involved in our design, e.g., TCP, UDP, and HTTP.

The *Transmission Control Protocol* (TCP) [15] provides a reliable, connection-oriented data stream delivery service. This reliability comes at a cost, though: TCP requires a three-way handshake for connection setup, normally

a graceful connection tear down and a datagram acknowledgment process. In addition, its reliability, sequencing and congestion control mechanisms introduce extra computation overhead as well. While the setup and tear down costs amortize well over long connections, they are expensive for short connections. Its adaptive congestion control, general resequencing and ARQ mechanisms are excessive for a single or few datagrams exchange.

The *User Datagram Protocol* (UDP) [12] provides a best-effort datagram service and, therefore, can operate in a more efficient manner than TCP. For example, UDP requires no connection setup or tear down, no acknowledgments, and little protocol state machine processing. On the flip side, UDP does not offer any of TCP's reliable delivery or congestion control behavior.

HTTP defines a request/reply protocol, where client applications can request data from servers by providing a *Universal Resource Locator* (URL). HTTP is used to address many different types of resources, including text, image, audio, video, executable files, index search results, and database query results.

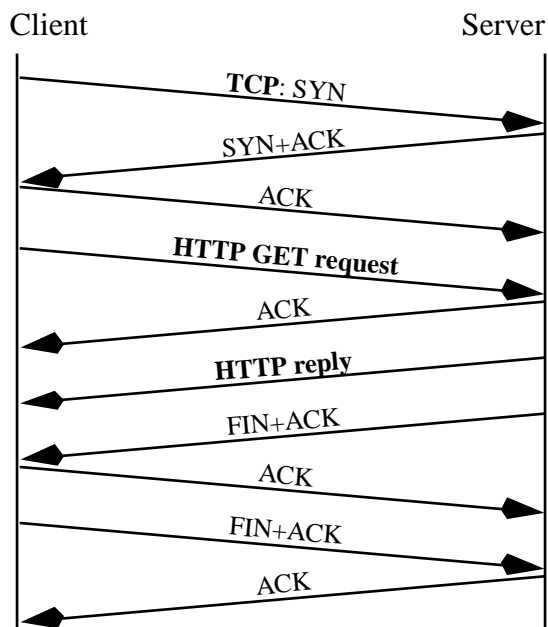


Figure 1. HTTP GET request and reply over TCP

Figure 1 illustrates a typical packet exchange for an HTTP GET request between a client and a server. Time progresses from the top of the figure downwards. First, the client establishes a TCP connection to the server with the three-way handshake (first three packets, labeled SYN, SYN+ACK, and ACK). Once the connection is established, the client transmits the GET request to the server. The GET request contains a URL to be fetched and served (labeled HTTP GET request). The server acknowledges the receipt of the request packet (labeled ACK) before it replies with the requested resource (labeled HTTP reply). The final four packets (FIN+ACK and ACK in both directions) are TCP control packets for graceful connection shutdown.

In Figure 1, only two packets seem to be *useful* (i.e., carrying data): one is the HTTP GET request, and the other one the HTTP reply. All other packets represent TCP overhead. Under the HTTP protocol Version 1.0, each transfer requires a separate TCP connection. HTTP Version 1.1 introduced *persistent HTTP connections* to address this problem. Figure 2 illustrates persistent HTTP: multiple HTTP GET requests (and their responses) use the same TCP connection.

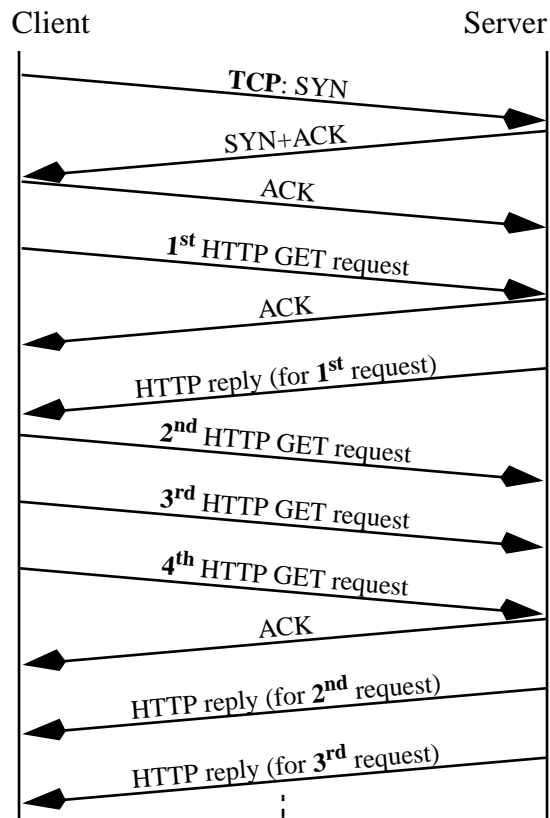


Figure 2. Example of multiple HTTP GET requests and replies over a persistent HTTP connection

As of end-1998, the use of persistent HTTP is limited by several factors. About 50% of installed HTTP browsers and a large portion of HTTP servers are not capable of using persistent HTTP.

One may expect that all pages for a Web site (i.e., with a common server identifier) reside on the same server. For example, the URLs `http://www.sunlabs.com/people/index.html` and the URL `http://www.sunlabs.com/research/index.html` share the server identifier `www.sunlabs.com` and we may expect that both URLs will be fetched from the same Web server. In practice, many Web sites are set up to use different servers for different URLs. In the example described above, the first URL may be served by the server `www.people.sunlabs.com` while the second is served by `www.research.sunlabs.com`. We call this approach *heterogeneous content provisioning*.

A site may be organized in this manner because providing the client with a common site access address is appealing for administrative and business reasons. Furthermore, the nature of the data and its amount might call for its distribution (static or dynamic) among multiple servers.

There are many reasons why a site may be organized in this manner. Firstly, providing the client with a common site access address is appealing for business reasons. It is easier for clients to remember a single commonly used site name and it gives the service provider the flexibility to index his Web services independently of their physical location. Secondly, the nature of the data and its amount might call for its distribution (static or dynamic) among multiple servers (for example, video clips might be served best by a customized video server.) Such heterogeneous provisioning may be dictated by administrative concerns (for example, multiple administrative domains). Also, a

Web site may include content from multiple sources and of variable popularity; it may not be attractive to move all content to a single server, or to several fully-replicated servers.

The HTTP REDIRECT mechanism (illustrated in Figure 3) is used to support heterogeneous content. HTTP server A hands off HTTP requests from clients to server B by sending a REDIRECT as a response to an HTTP GET request. Note that the HTTP REDIRECT is part of a short HTTP transfer. Using TCP for an HTTP request–HTTP REDIRECT pair requires at least 7 packets (usually 9-10 packets), while two packets suffice if UDP is used as the underlying transport protocol.

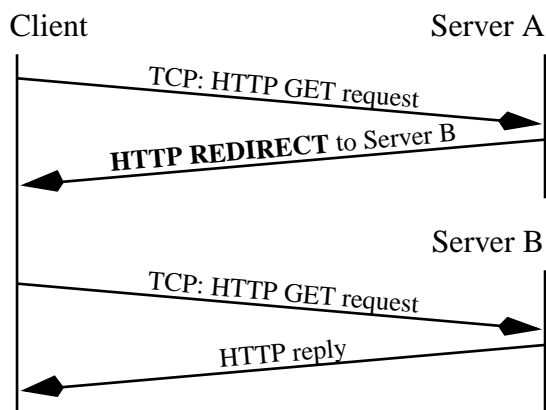


Figure 3. HTTP REDIRECT by Server A to Server B

We already mentioned that TCP transport is expensive for small payloads. Here is another example of a scenario in which UDP is a sufficient and more efficient choice for HTTP data transport than TCP. Section 4 then presents our analysis and simulation results supporting this claim, even in the presence of packet loss and corruption.

Although resources have uniform identifiers (URLs) and share a common prefix, they are usually not accumulated at a single location. Therefore, most HTTP servers operate as an intermediary, single point of access that retrieves the requested resources from a variety of other services. For example, if a requested file is not on the local file system of the HTTP server, before it can be served to the HTTP client, it first needs to be retrieved through a network file system. Database queries (results) are relayed by the HTTP server to (from) the appropriate database system. There are often multiple dedicated servers used behind the uniform front end of the HTTP server. This discussion illustrates how resources are often not arranged according to the views presented to clients. Rather, resources are distributed according to administrative structures and control (e.g., file systems that represent project file space of different engineering groups), functional criteria (e.g., customer database queries are processed by the database server, whereas files are served by a network file system), or because their collective size exceeds storage available locally.

Instead of replicating and moving data to a single HTTP server, one can imagine using multiple servers, each located closer to a subset of the data. Servers would serve only the data that is locally available to them. Files would be organized preserving administrative and organizational structures. We call this approach *heterogeneous content provisioning*. Its implementation requires the availability of multiple HTTP servers that provide to clients the illusion of a single service. The HTTP protocol includes a protocol message for the redirection of HTTP requests to other servers. It is called an HTTP REDIRECT and is illustrated in Figure 3. HTTP server A hands off HTTP requests from clients to server B by sending a REDIRECT as a response to an HTTP GET request. Using TCP for an HTTP request–HTTP REDIRECT pair requires a total of at least 7, and usually 10 TCP packets. Using UDP, two datagrams are sufficient and the overall number of bytes sent over the network is smaller with UDP.

3 Design

A good solution for reducing the TCP overhead for small connections should have the following characteristics:

- transparency to users,
- backward compatibility,
- opportunity for incremental deployment, and
- low runtime overhead in space and time (if any)

3.1 Proposal: Hybrid TCP-UDP Transport

One approach is to use UDP instead of TCP as the transport protocol for HTTP traffic. However, in the presence of unreliable network layer communications, reliability services and congestion management need to be added to UDP to make this a viable proposal. We decided against putting such functionality into UDP.

Instead, we use a hybrid approach where short connections are served using UDP and large connections use TCP as its transport protocol. In this manner, the TCP overhead for short connections is avoided, but the benefits from TCP's well-tuned timers, retransmission, congestion control, and error recovery mechanisms are preserved. In this scheme, clients first attempt to use UDP as their transport protocol, and fall-back to using TCP, if UDP turns out to be the wrong choice for the requested URL. The fall back mechanism provides the following guarantees:

- If any of the initial UDP packets are lost, the loss is gracefully handled by switching to TCP.
- If the contacted HTTP server does not implement the hybrid scheme, the client will re-try with TCP.

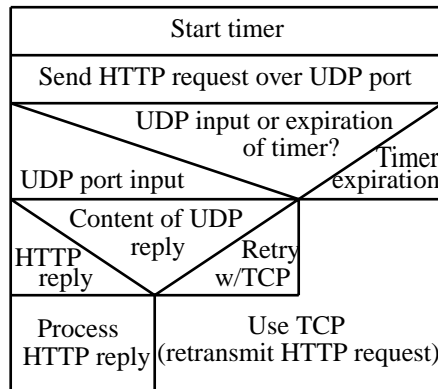


Figure 4. Algorithm executed at HTTP client application to receive HTTP service over either TCP or UDP

Figure 4 presents this algorithm in more detail. A *hybrid* capable HTTP client sends the HTTP GET request using UDP as the underlying transport protocol. The client starts a timer at the time the request is sent.

When the server processes the client's request, it can choose from one of the following alternatives:

1. If the response is small enough (say, it fits into one datagram), the server returns it using UDP. Small Web pages and most cache validation and HTTP REDIRECT responses fall into this category.

- The server can ask the client to re-try using TCP if the reply is too large to fit into a single UDP packet. At this time, the server can also ask the client to try a different URL. This approach is useful if the server generated the reply dynamically and attempts to avoid generating the reply a second time for the subsequent re-try with TCP. This extension can be implemented with a new HTTP return code.

At the client side, one of the following three events can happen:

- The client gets a response from the server. If the reply contains the desired HTTP reply, the client processes the data. If the server asks the client to re-try (a different URL and/or using TCP), the client does so.
- If the server does not handle HTTP packets sent over UDP, the client may get an ICMP (*Internet Control Message Protocol* [14]) error message (destination unreachable/protocol unreachable). In this case, the client should re-try using TCP.
- If the timer expires, the client should re-try with TCP.

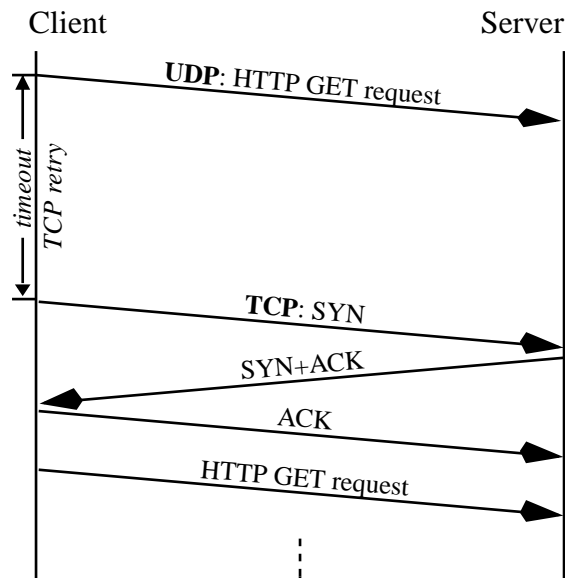


Figure 5. Implicit fall-back from UDP to TCP: One of the UDP request/response packets got lost or the server is not UDP capable

Figure 5 illustrates the packet exchange in case the timer expires. This timeout feature provides reliability and backwards compatibility with servers that do not use the hybrid TCP-UDP scheme. We recommend that this timeout interval be set the same as the corresponding TCP initial timeout interval. As the upgraded clients can work with servers that do not implement the hybrid scheme, this fall-back mechanism supports gradual, incremental deployment in the Internet.

Figure 6 demonstrates the packet exchange where the server requests the client to resend the HTTP request over TCP.

For heterogeneous content, this hybrid scheme presents another interesting option (bypassing the client completely): as the HTTP server receives HTTP requests (with the URL) without any initial state establishment, the

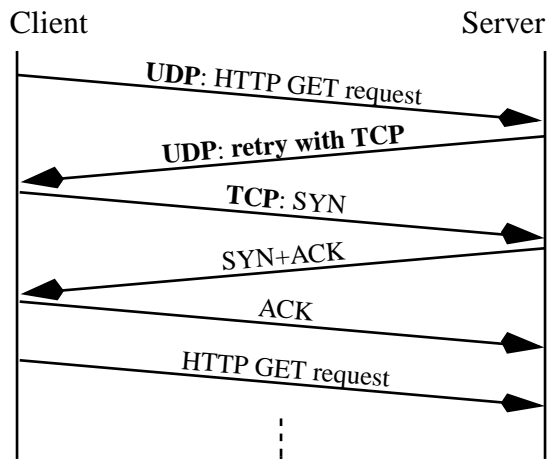


Figure 6. Explicit fall-back from UDP to TCP: Server is UDP capable, but the requested resource warrants use of TCP

HTTP request can be forwarded and served by another server (which masquerades as the first server in its response) without further client involvement.

For heterogeneous content, our hybrid scheme offers a new mechanism to redirect requests; this new mechanism is transparent to clients. When the HTTP server receives an HTTP request over UDP, the request can be forwarded and served by another server without further client involvement. The server that sends the HTTP response must masquerade as the first server, or the client would not accept the HTTP reply. It is interesting to see that we can forward a single HTTP request through a chain of servers before we find one that can respond to the client, and it would still be completely transparent to the clients.

The use of UDP makes this optimization possible without any kernel-level changes in the server operating systems. This scheme can be used to improve Web cache validation, robust anycasting, and transparent proxies.

Furthermore, we can use HTTP proxies to incrementally deploy the hybrid TCP-UDP scheme without modifying the installed client browsers.

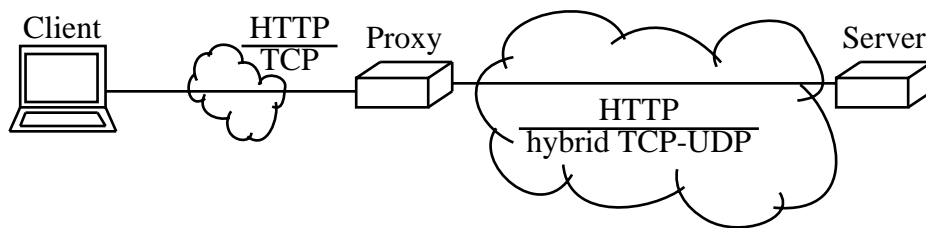


Figure 7. Incremental deployment of hybrid TCP-UDP scheme

The hybrid TCP-UDP scheme can be deployed incrementally by using proxy servers. Figure 7 illustrates a scenario where clients have not yet upgraded, but where the HTTP retrieval in the wide area takes advantage of our proposal.

Using UDP as the underlying transport protocol does not raise new security concerns. While it is relatively (compared to TCP) easier for intermediate routers to alter HTTP replies if they are sent over UDP, the same basic problem exists in TCP and is referred to as *TCP hijacking* [2, 6]. It can be solved using end-to-end data

integrity security services, such as cryptographically signed digital signatures protecting the payload and control information of communications.

The same argument can be brought forward against the possibility of *smurf*-style denial of service attacks, where a third party causes many servers to reply to a single client by spoofing its source address on the HTTP/UDP request message, causing the client to be overloaded. [16, §4] classifies a variety of solutions against this type of attack.

The hybrid TCP-UDP scheme improves the utility of HTTP for all three entities involved in HTTP transfers: clients, servers, and the network. The scheme achieves:

- reduced browsing latency for clients because latency introduced by the 3-way handshake for TCP connection establishment is avoided,
- reduced load for servers because fewer TCP connections are set-up, maintained, torn down, and because fewer packets are processed, and
- reduced network traffic because fewer TCP control packets are needed.

4 Performance Evaluation

In this section, we quantify the gains obtained by using the hybrid TCP-UDP scheme; our analysis concentrates on answering the following questions:

- How do the use of persistent HTTP and the choice of connection parameters (e.g., max keep-alive time, number of concurrent connections) affect the savings resulting from the use of the hybrid scheme?
- How do other network-related parameters, including network loss rates and MTU size, affect the savings?
- How are these savings then affected by various parameter choices for the hybrid TCP-UDP scheme, as described in Section 4.2? We are especially interested in the *max-udp* parameter. UDP is used only when the HTTP response is small enough to fit within *max-udp* packets.

Answering these questions, the following sections present the results from several simulations on our C++/Tcl based simulator. We performed several simulation experiments to answer the above-mentioned questions and we present here the results from these experiments. We ran these simulations on our own C++/Tcl-based simulator. The goal was to make the experiments realistic so that the results obtained can be transposed to implementations. For this reason, we decided against using synthetic workloads. Instead, we relied on HTTP trace data to drive our simulations. This trace data consists of 18 days' worth of HTTP traces gathered from U.C. Berkeley's Home IP service [4] and of traffic traces obtained within Sun Microsystems' internal network. The traffic for port 80 (HTTP) was recorded. All other protocols (or ports) were excluded from these traces. The traces [4] amounted to 9 million connections over 18 days, thereby providing us with about 500K HTTP transfers per day.

To get an initial feeling for the relevance of the approach we first derived some basic statistics from the traces. First, we derived the Cumulative Distribution Function (CDF) of the size of the replied data in the servers' response. A curve-fitting experiment indicated that a shifted exponential distribution of the type $1 - \exp[-.00035(x - 80)]$ (where x is the reply size in bytes, and 80 is the size of the shortest reply) provides a very close approximation, in particular for response sizes smaller than 4 Kilobytes. These data are depicted in Figure 8. A noteworthy result of this distribution is that approximately 40% of the replies would fit into a single 1500 Byte packet. These results are in line with those reported in [7].

We also examined the statistics of HTTP transactions with respect to persistence. Figures 9 and 10 depict the histogram and CDF of the HTTP transactions that took place over the entire measurement period. Noteworthy here

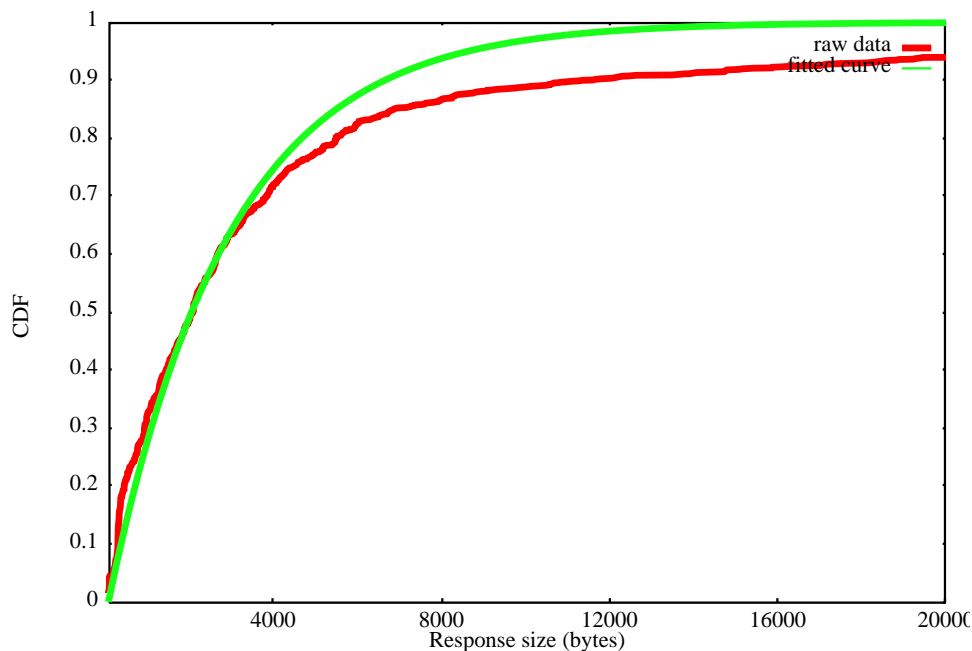


Figure 8. CDF of reply size

is the dominance of the occasional transaction. For example, 88% of the source destination pairs conducted 40 or fewer transactions over the entire period, and 98% conducted fewer than 100 transactions. Examining the inter-arrival times of these pairs reveals that these transactions were randomly distributed over the entire measurement period.

These results are very encouraging. They indicate that for a vast majority of the transactions, persistence would not help because the underlying TCP connection would not stay open long enough to allow its reuse. Therefore, most connections would have to be reestablished, and, given the distribution of the reply size, most connections would complete with a very small number of reply packets.

With these encouraging initial results, we conducted more elaborate experiments testing the approach under a variety of parameter changes as indicated above.

4.1 Evaluation Metrics

The following metrics are used for the performance evaluation:

Number of bytes/packets transferred This metric captures the network load, as well as one aspect of the load on the servers and the clients.

Browsing latency This metric captures the performance as observed by the clients.

Number of connections set up This metric captures one aspect of the load on the servers (extra work required to set up and maintain TCP connection state, as well as extra cost in searching *Process Control Block* (PCB) lists).

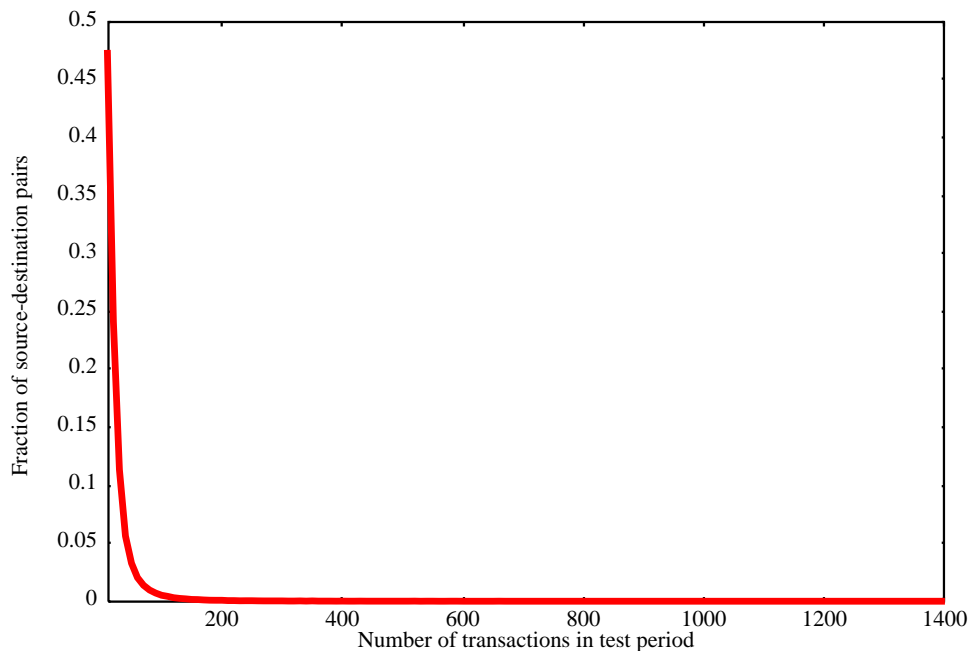


Figure 9. Histogram of transaction densities

The choice of our third metric, number of connections set up, was motivated by the following reasons. Our packet traces were HTTP logs. It was impossible for us to obtain the inter-packet traffic patterns that are needed for the first two metrics. Furthermore, the number of connections is the appropriate metric for evaluating HTTP redirection (as well as cache validation). These requests, and the corresponding responses, tend to be small enough to fit into the minimum-MTU IP packets.² Also, this metric is linked to the other two metrics. If each HTTP transfer requires a new TCP connection, more bytes/packets will be seen on the network (connection setup and teardown packets) and clients will observe extra latency (due to connection setup delays as well as slow-start latency).

We convert this metric to a unitless ratio (fraction of connections saved) by performing each experiment twice. If an instance requires 100 distinct TCP connections under the current system, and it requires only 80 distinct TCP connections under the proposed hybrid TCP-UDP scheme, we obtain the resulting fraction of $(100 - 80)/100 = 20\%$.

4.2 Experiments

We performed many sets of experiments, each time varying one of four workload parameters:

Persist This parameter describes the fraction of requests from the clients that support persistent HTTP. The parameter value ranges from 0.0 to 100.0; a value of 70.0 implies that 70% of requests originate from clients that support persistent HTTP.

Loss-rate This parameter describes the fraction of packets that are lost in transit. The parameter value ranges from 0.0 to 1.0; a value of 0.10 implies that with 10% probability a packet will be dropped in the network.

²The Internet Protocol requires the underlying network to support IP packets at least as large as 576 bytes [13, §3.1].

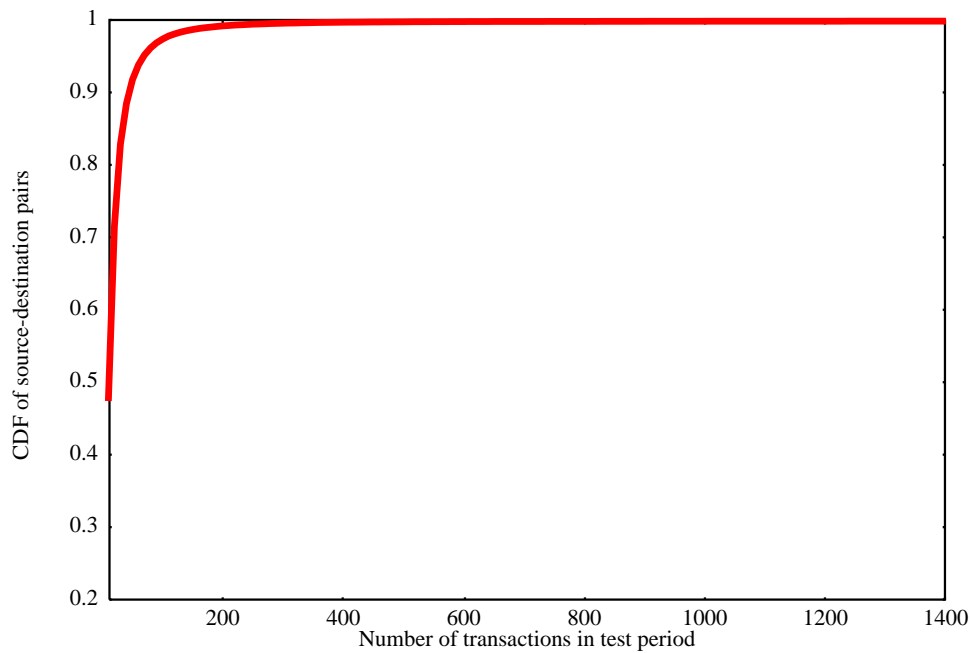


Figure 10. CDF of transaction densities

For these simulations, we assumed that packet losses are independent (not bursty, zero correlation). This is a pessimistic assumption because it overstates the number of connections that will lose at least one packet, thereby diminishing the overall benefits of the hybrid TCP-UDP scheme.

Max-udp This parameter controls the server policy for choosing between use of UDP and TCP in the hybrid TCP-UDP scheme. A value of $\text{max-udp}=4$ implies the server will try to use UDP for all conversations where it can send the entire data in up to 4 packets. Otherwise, it directs the client to use TCP.

MTU This parameter describes the *Maximum Transmission Unit* (MTU) for packets in the network. An MTU size of 1460 implies that the network MTU allows for up to 1460 bytes payload (packet size minus TCP header, IP header, and link-level headers). We assume that the end hosts send MTU-sized packets whenever possible. The larger the MTU, the more likely it is that the hybrid TCP-UDP scheme would avoid setting up a TCP connection (we avoid a TCP connection when HTTP data size is less than $\text{MTU} \cdot \text{max-udp}$, and the network does not lose any of these packets).

In our simulations, we included further parameters for controlling the behavior of persistent HTTP. A keep-alive parameter was set to 60.0, meaning that a persistent HTTP connection would be closed after 60 seconds of inactivity. Also, a max-connect parameter was set to 1024. Thus, the server will only support 1024 concurrently active connections. Inactive connections were closed according to a *Least Recently Used* (LRU) policy.

4.3 Results

The simulation results were reasonably similar for the various traffic traces. For all graphs in this Section, the x-axis quantifies the varied parameter (e.g., loss rate) and the y-axis quantifies the evaluation metric (e.g., fractions of connections saved, expressed as a percentage).

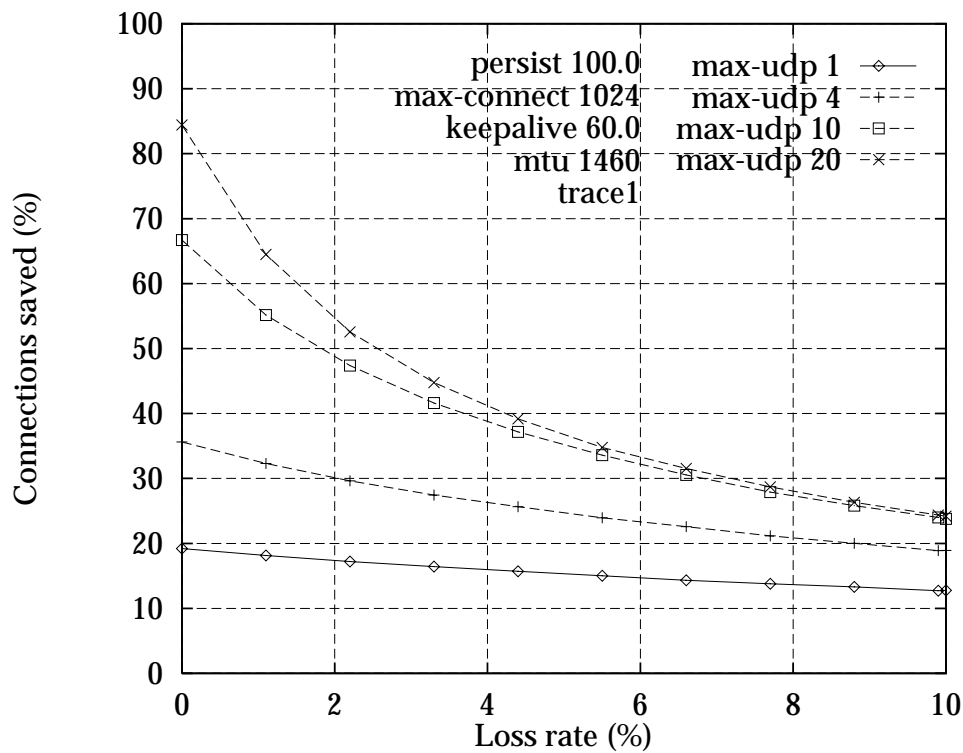


Figure 11. Effect of loss rate and max-udp parameter

Loss Rate and Max-Udp

The graph in Figure 11 shows the effect of the network packet loss rate and the setting of max-udp parameter in the server on the overall performance of the hybrid TCP-UDP scheme. For this experiment, we assumed an MTU of 1460 bytes, as well as full persistence (a very conservative assumption), i.e., all clients support persistent HTTP, and all requests are multiplexed on existing TCP connections wherever possible. As the graph shows, with the very conservative policy of setting max-udp value to 1 (use UDP only if all data from server fits into one packet), the hybrid scheme gains as much as 18%-19% with the typical (around 1%) loss rate observed in many intranets, as well as on the Internet. Even with very high packet losses (around 10% loss rate), the hybrid scheme outperforms the traditional mode by 12%-13%. The graph also depicts the additional benefits that using a higher max-udp value would provide: up to 85% improvement over the current schemes. Even with a max-udp setting of 4 to 10, we can expect to see benefits of 35% to 66%, even under the assumption that persistent HTTP is used exclusively.

Persistent HTTP and Loss Rate

The graphs in Figures 12 and 13 show the effect of persistent-HTTP and network packet loss rates on the overall performance of the hybrid TCP-UDP scheme. The max-udp parameter is set to the value 1 in Figure 12 and the value 4 in Figure 13; the MTU remains 1460 bytes. We performed this experiment to evaluate the performance in the presence of heterogeneous clients (some support persistent HTTP; others do not), as well as to see the effect of the max-udp parameter in this system. As both figures show, the performance is dominated by the presence of

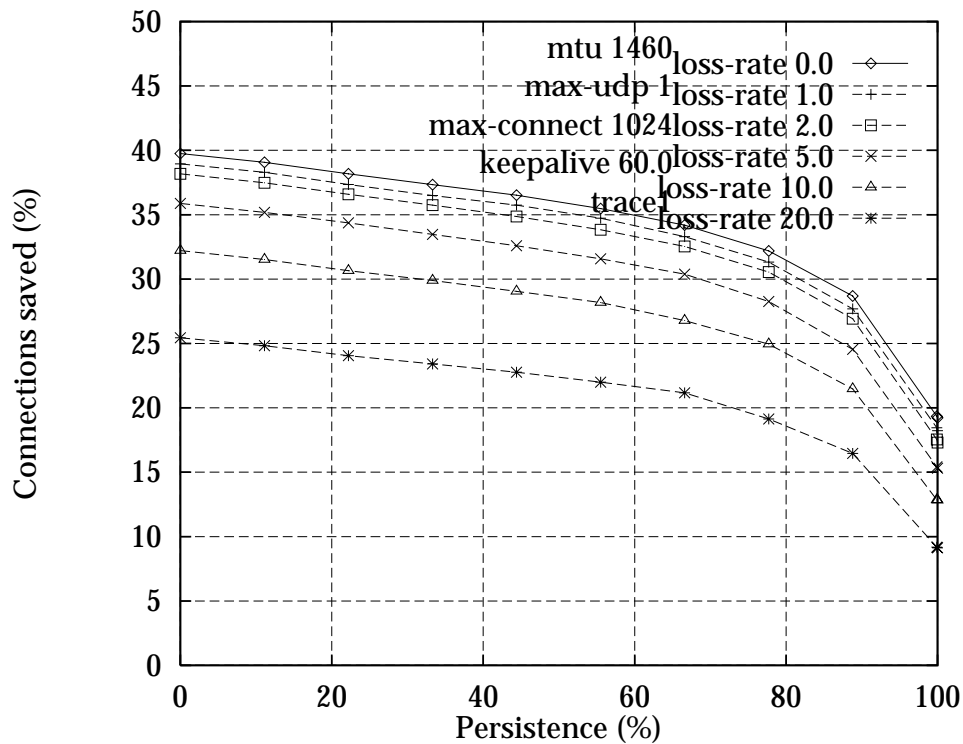


Figure 12. Persistence and loss rate (a)

non-persistent clients. Even with 80% persistence (i.e., only 20% do not support persistence), the results are close to the performance with zero persistence. As expected, the results are much better for max-udp of 4 (as compared to max-udp value of 1). As Figure 12 shows, with max-udp set to 1 and with zero persistence, the hybrid scheme provides a gain of 25%-40% over the traditional schemes, depending on the overall packet loss rate. The gains decrease very slightly as the persistence parameter increases to well beyond 80%. Even in the presence of 90% to full persistence, the hybrid scheme provides performance gains of 15% to 30%, except when the packet loss rate goes up to 20% (i.e., pathological conditions). With the max-udp set to 4, the performance gains tend to be around 70% for moderate loss rates, decreasing to around 60% for 70%-80% persistence. Even with 10% packet loss rate, the hybrid scheme gains around 40% at 80% persistence.

Max-udp and Loss Rate

Figure 14 illustrates the effect of the max-udp parameter and the network packet loss rate values on the overall performance of the hybrid TCP-UDP scheme. We chose worst-case values for the MTU and the persistence parameters: the MTU parameter is set to the value 460 which reflects the lower MTU settings that some people use over 28.8K modems and other low-speed links. Most of the Internet links, as well as higher-speed links to end-users, however, use higher MTU values (typically around 1.5 Kilobytes)[8]. We also assume full persistence (i.e., all HTTP transfers are multiplexed on existing TCP connections whenever possible). We performed this experiment to evaluate the effect of max-udp parameter on the system performance, with varying loss rates. As the graph shows, the performance gains due to the hybrid TCP-UDP scheme increase steadily with increasing value(s) of the max-udp parameters, and are slightly reduced under reasonable loss (around 1% packet loss rate). Even in

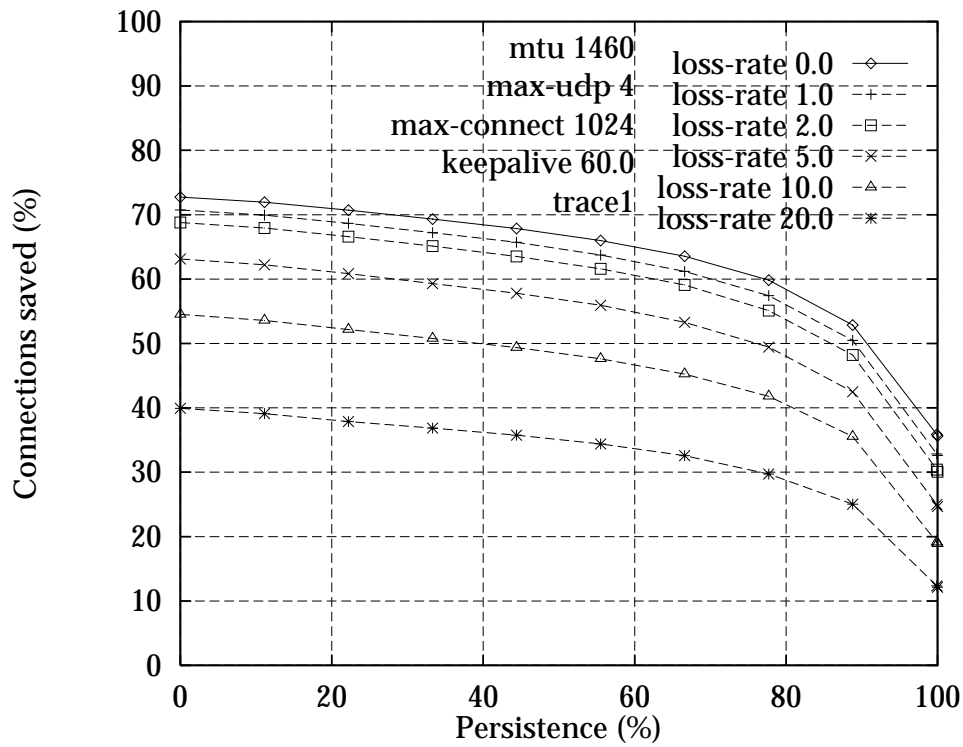


Figure 13. Persistence and loss rate (b)

the presence of relatively high loss rates (about 5%), we observe 15% gains with the max-udp equal to 3. This result is especially interesting in the context of recent research in increasing the TCP initial window (to 4 packets or 4Kb, whichever is lower).

MTU and Loss Rate

The graphs in Figure 15 and Figure 16 show the effect of the network MTU values and packet loss rates on the overall performance of the hybrid TCP-UDP scheme. The max-udp parameter is set to the value 4 in Figure 15 and the value 1 in Figure 16; the persistence parameter is set to the value 70% in Figure 15 and the value 100% (full persistence) in Figure 16. We performed this experiment to evaluate the performance, for varying MTU values and loss rates, in the presence of heterogeneous clients (some support persistent HTTP; others do not), as well as to see the effect of the max-udp parameter in this system. As both figures show, higher MTU values lead to increased performance gains due to the hybrid TCP-UDP scheme; these gains are reduced somewhat due to the packet losses in the network, but the gains remain high even with significant (around 2%) packet loss rates. As Figure 15 shows, the gains are much higher with heterogeneous clients (70% persistence); setting max-udp to 4 results improves the performance gains, even in the presence of pathological network behavior (10% to 20% packet loss rate).

4.4 Summary of Simulation Results

The hybrid TCP-UDP scheme provides significant performance gains over the traditional use of TCP as the only transport protocol for HTTP traffic. We observed performance gains exceeding 20-25% with persistent HTTP

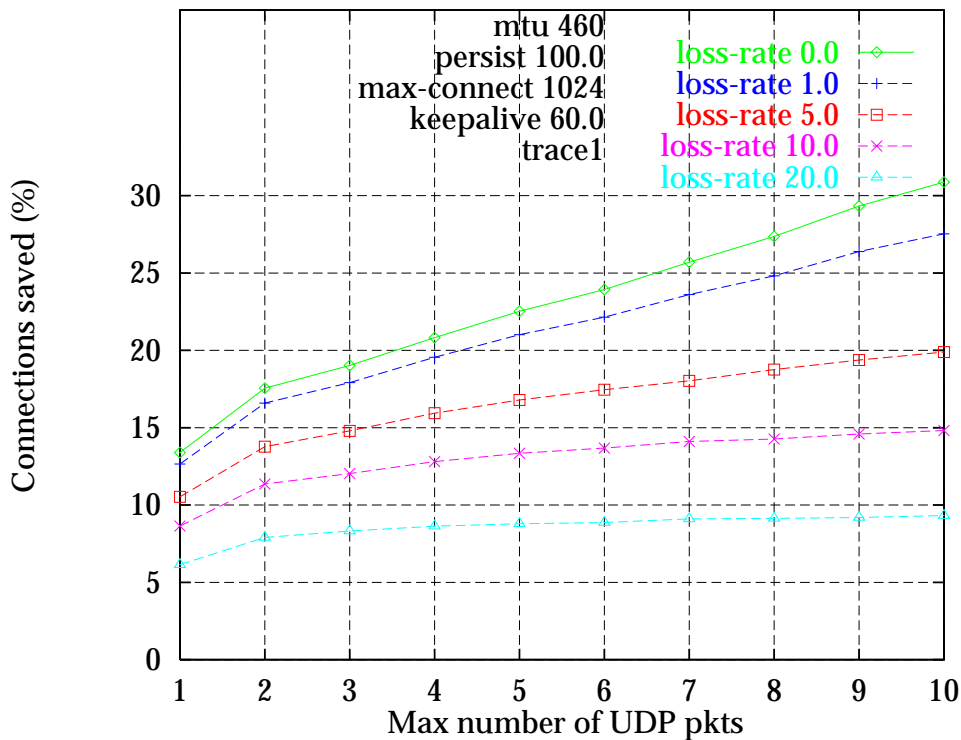


Figure 14. Max-udp and loss rate

clients, and over 40-50% with clients without persistent HTTP.

In heterogeneous environments, the system performance was dominated by the presence of non-persistent clients. Even with 70% to 80% persistence, the performance gains of the hybrid scheme were close to that of a system with zero persistence.

The hybrid TCP-UDP scheme’s performance gains were reduced somewhat by network packet losses, but the gains remained significant even under pathological network behavior (10% to 20% packet loss rate).

It is beneficial to try sending more than 1 packet via UDP, especially if the network MTU is small (around 500 bytes). However, we must appropriately consider and trade-off the increased likelihood of network congestion due to the increase in the max-udp parameter.

5 Conclusions and Outlook

This paper introduced and analyzed a hybrid TCP-UDP transport layer scheme for HTTP. The hybrid TCP-UDP scheme combines the low cost of using UDP with the high reliability and congestion control features of TCP. It benefits from the low overhead of using UDP for short transfers and, for large transfers, it is able to use TCP for reliable delivery and good congestion behavior. Our simulation experiments verified these gains. We observed performance gains exceeding 25% with persistent HTTP clients, and over 50% for clients without persistent HTTP, even in the presence of high network packet losses. The performance gains for mixed environments (i.e., clients using persistent HTTP, as well as traditional HTTP) were similar to those of a system where no clients used persistent HTTP. The hybrid scheme is attractive because it only requires application-level changes, while the operating system kernel code can remain unchanged. Finally, the scheme is incrementally deployable in the

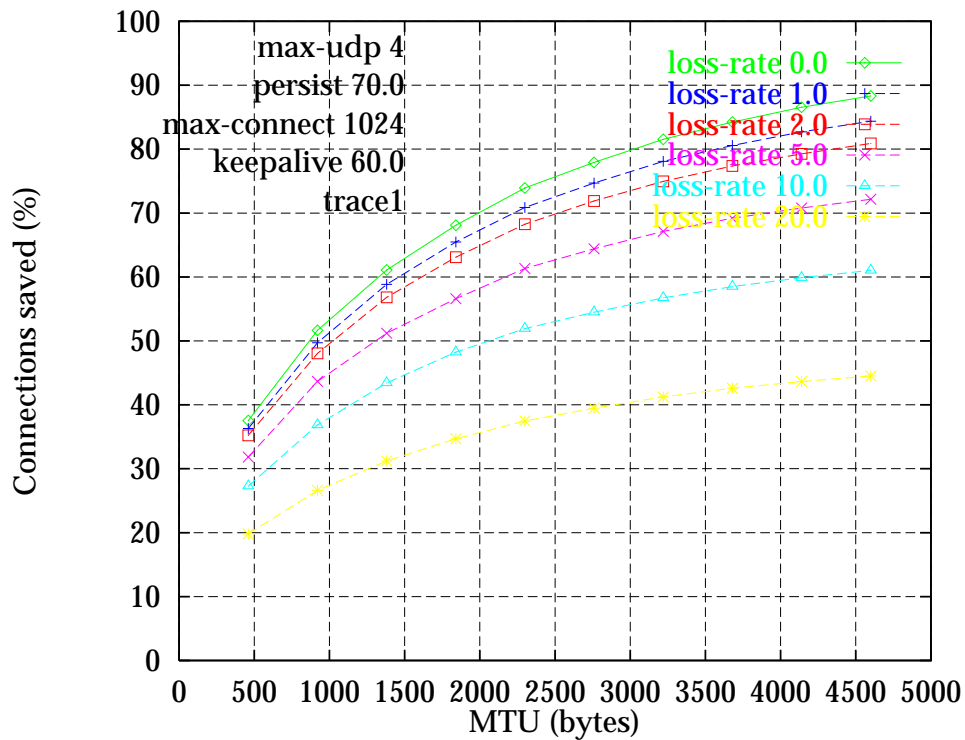


Figure 15. Maximum Transmission Unit (MTU) and loss rate (a)

current Internet and it is fully compatible with the deployed base.

We are currently exploring some promising optimizations to further improve system performance. First, the clients can use some heuristics to guess (on a per-transfer or per-session basis) if they can use UDP for successful transmission. For example, the clients can keep track of whether the server supports the hybrid scheme at all, or predict the file size based on the size of the (expired) cached copy of a previously retrieved version, or assume that the Postscript and PDF files tend to be large, while README files tend to be small. Second, the servers can adapt the max-udp parameter based on the observed network loss rates. Third, the clients can avoid redundant data retransmissions by using sub-ranges appropriately.

We also expect that further performance gains can be observed, if recent trends in Web traffic continue. The increasing use of style sheets will lead to smaller HTTP transfers, and so will the increasing use of cache validation and HTTP redirection. The performance gains can be further increased by setting the server policy to use UDP more aggressively (for example, when the data can fit into 4 packets, analogous to some recent proposals to increase the initial TCP window size[11]).

Acknowledgments

We are grateful to Geoffrey Baehr, Charles Perkins, Erik Guttman, and Abhishek Chauhan for their insightful comments on the drafts of this paper.

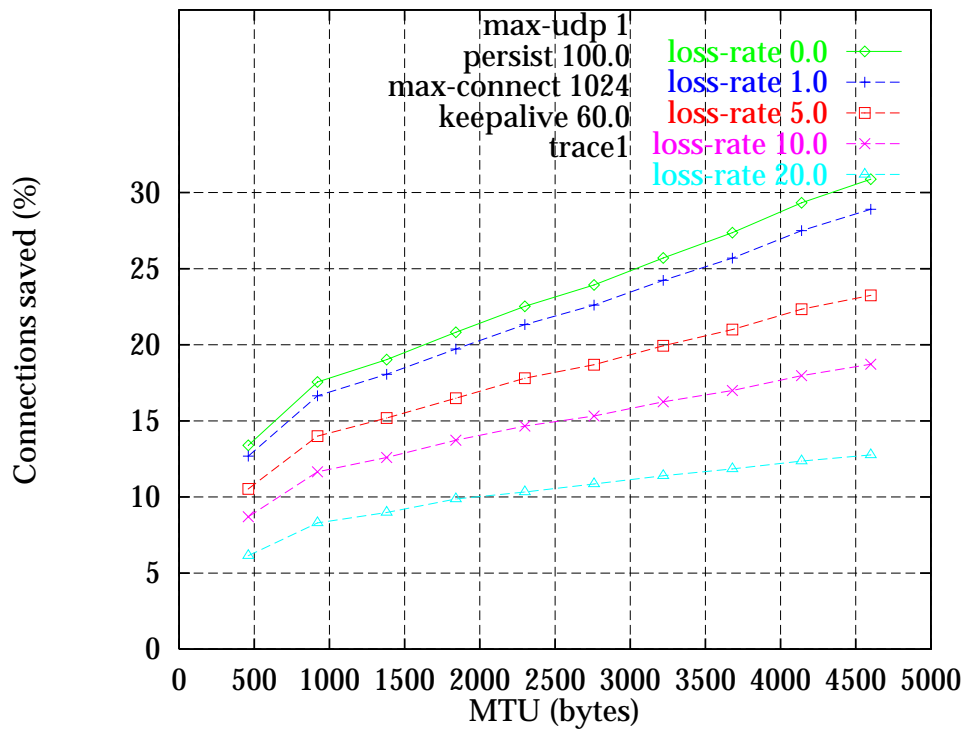


Figure 16. Maximum Transmission Unit (MTU) and loss rate (b)

References

- [1] Bob Braden. *RFC-1644 T/TCP - TCP Extensions for Transactions*. Network Working Group, July 1994.
- [2] CERT. *IP Spoofing Attacks and Hijacked Terminal Connections, CA-95:01*. Computer Emergency Response Team, Carnegie Mellon University, Pittsburgh, Pennsylvania, January 1995.
- [3] Roy T. Fielding, Jim Gettys, Jeffrey C. Mogul, Henrik Frystyk Nielsen, and Tim Berners-Lee. *RFC-2068 Hypertext Transfer Protocol – HTTP/1.1*. Network Working Group, January 1997.
- [4] Steven D. Gribble. UC Berkeley Home IP HTTP traces. July 1997. Available at <http://www.acm.org/sigcomm/ITA/>.
- [5] John Heidemann, Katia Obraczka, and Joe Touch. Modeling the Performance of HTTP Over Several Transport Protocols. *IEEE/ACM Transactions on Networking*, 5(5):616–630, October 1997.
- [6] Laurent Joncheray. A Simple Active Attack Against TCP. In *Proceedings of the 5th UNIX Security Symposium*, pages 7–19, Salt Lake City, Utah, June 1995. USENIX.
- [7] Bruce Mah. An Empirical Model of HTTP Network Traffic. In *Proceedings of Infocom'97*, pages 593–600, Kobe, Japan, April 1997.
- [8] Jeffrey Mogul and Steve Deering. *RFC-1191 Path MTU Discovery*. SRI International, Menlo Park, CA, November 1990.

- [9] B. Clifford Neuman. *The Virtual System Model: A Scalable Approach to Organizing Large Systems*. PhD thesis, University of Seattle, Washington, 1992.
- [10] Venkata N. Padmanabhan. Private communication, 1997.
- [11] Kedarnath Poduri and Kathleen Nichols. *RFC-2415 Simulation studies of increased initial TCP window size*. Network Working Group, September 1998.
- [12] Jon Postel, editor. *RFC-768 User Datagram Protocol*. Network Information Center, August 1980.
- [13] Jon Postel, editor. *RFC-791 Internet Protocol*. Information Science Institute, University of Southern California, September 1981.
- [14] Jon Postel, editor. *RFC-792 Internet Control Message Protocol*. Information Sciences Institute, University of Southern California, September 1981.
- [15] Jon Postel, editor. *RFC-793 Transmission Control Protocol*. Information Sciences Institute, University of Southern California, September 1981.
- [16] Christoph L. Schuba, Ivan V. Krsul, Markus G. Kuhn, Eugene H. Spafford, Aurobindo Sundaram, and Diego Zamboni. Analysis of a Denial of Service Attack on TCP. In *Proceedings of the Symposium on Security and Privacy*, pages 208–223, Oakland, California, May 1997. IEEE.

About the Authors

Israel Cidon received the B.Sc. and D.Sc. degrees in electrical engineering from the Technion (Israel Institute of Technology), Haifa, Israel. Dr. Cidon was one of the (four) founding managers of Micronet. He was responsible for the product hardware development at both Micronet and Mechkam Technologies. From 1985 until 1990 Dr. Cidon managed the Network Architectures and Algorithms Group at IBM T.J. Watson Research Center and served as a Research Staff Member. During his time at the Technion (1990-1993), Israel was one of the founders (along with Prof. Rom and Prof. Sidi) of the Computer Networking Laboratory at the Electrical Engineering Dept. and served as its director in the first year. In 1994-95, he was the manager of the High Speed Networking Research group at Sun Microsystems Laboratories. Currently, Israel holds the position of associate professor at the department of electrical engineering at the Technion.

Amit Gupta received his Bachelors in Computer Science from the Indian Institute of Technology (IIT) in Delhi, India, and his Ph.D. in Computer Science from the University of California at Berkeley. During his doctoral research, Amit designed and built network architectures and protocols for delivering performance guarantees for network services. Dr. Gupta is currently Chief Technology Officer at Andale. Prior to co-founding Andale, Amit was at Sun Microsystems' visionary Sun Labs group working to define innovative networking and security technologies. He designed and built new network architectures and protocols for highly reliable and scalable web services. He also developed new technologies for securing large network-based services from unauthorized and malicious users. Amit has more than 25 patents pending for distributed systems technologies.

Raphael Rom received the B.Sc. and M.Sc. degrees in electrical engineering from the Technion (Israel Institute of Technology), Haifa, Israel, and the Ph.D. degree in computer science from the University of Utah, Salt Lake City, UT. He was a Senior researcher on the research staff of SRI International in California, and subsequently joined the Faculty of Electrical Engineering in the Technion, Haifa, Israel. Since 1989, he has been with Sun Microsystems where he lead and managed the high speed networking group of Sun Labs and is engaged in modeling and analysis of communication networks. Dr. Rom held visiting positions in IBM T.J. Watson Research Center and Stanford University. Dr. Rom is the coauthor of the book "Multiple Access Protocols: Performance and Analysis" and has coautherd a large number of technical and scientific publications. His areas of interest are algorithms for, and performance analysis of data communication and wireless networks and the design of general data communication systems.

Christoph L. Schuba is Staff Engineer in the Security Research Group at Sun Microsystems Inc.'s Sun Labs, Mountain View, California. He studied mathematics and management information systems at the University of Heidelberg and the University of Mannheim in Germany. As a Fulbright scholar, he earned an M.S. in Computer Science in 1993 from Purdue University, where he was a member of the COAST/CERIAS research group. He received a Ph.D. from Purdue in 1997, performing part of his dissertation research in the Computer Science Laboratory at the Xerox Palo Alto Research Center (PARC). Recently, Christoph taught the graduate course in computer and network security at the Mathematics and Computer Science Department at San Jose State University. Christoph's research interests include network and computer system security and distributed systems.