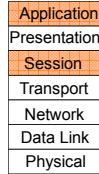


The Web as a Shim on TCP

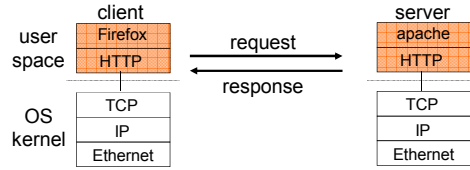
- HTTP and the Web (but not HTML)
- Focus
 - How do Web transfers work?
- Topics
 - HTTP, HTTP1.1
 - Performance Improvements
 - Protocol Latency
 - Caching



djw // CSE 461, Fall 2009

Lhttp.1

Web Protocol Stacks



- To view the URL <http://server/page.html> the client makes a TCP connection to port 80 of the server, by it's IP address, sends the HTTP request, receives the HTML for page.html as the response, repeats the process for inline images, and displays it.

djw // CSE 461, Fall 2009

Lhttp.2

HTTP Request/Response

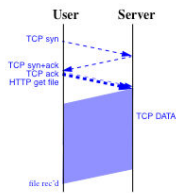


FIGURE 3. HTTP File Transfer

1 RTT channel OPEN
 0.5 RTT send request
 0.5 RTT file starts to arrive
 Ftrans time to transmit the file
 2 RTT + Ftrans = time to get a file in HTTP

djw // CSE 461, Fall 2009

Lhttp.3

Simple HTTP 1.0



- HTTP is a tiny, text-based language
- The GET method requests an object
- There are HTTP headers, like "Content-Length:", etc.
- Try "telnet server 80" then "GET index.html HTTP/1.0"
 - Other methods: POST, HEAD, ... google for details

djw // CSE 461, Fall 2009

Lhttp.4

HTTP Request/Response in Action

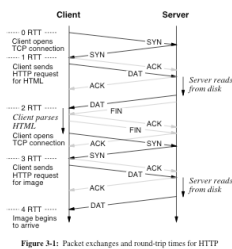


Figure 3-1: Packet exchanges and round-trip times for HTTP

- Problem is that:
 - Web pages are made up of many files. Most are very small (< 10k)
 - files are mapped to connections
- For each file
 - Setup/Tear-down
 - Time-Wait table bloat
 - 2RTT "first byte" latency
 - Slow Start+ AIMD Congestion Avoidance

The goals of HTTP and TCP protocols are not aligned!

TCP Behavior for Short Connections

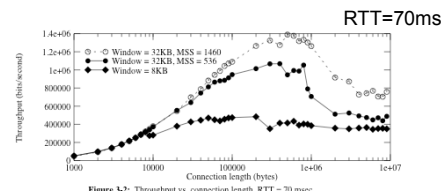


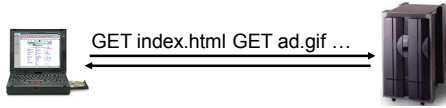
Figure 3-2: Throughput vs. connection length, RTT = 70 msec

Figure 3-2 shows that, in the remote case, using a TCP connection to transfer only 2 Kbytes results in a throughput less than 10% of best-case value. Even a 20 Kbyte transfer achieves only about 50% of the throughput available with a reasonable window size. This reduced throughput translates into increased latency for document retrieval. The figure also shows that, for this 70 msec RTT, use of too small a window size limits the throughput no matter how many bytes are transferred.

djw // CSE 461, Fall 2009

Lhttp.6

HTTP1.1: Persistent Connections

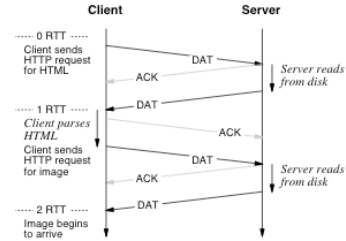


- Idea: Use one TCP connection for multiple page downloads (or just HTTP methods)
- Q: What are the advantages?
- Q: What are the disadvantages?
 - Application layer multiplexing

djw // CSE 461, Fall 2009

Lhttp7

HTTP/1.1



djw // CSE 461, Fall 2009

Lhttp8

Effect of Persistent HTTP

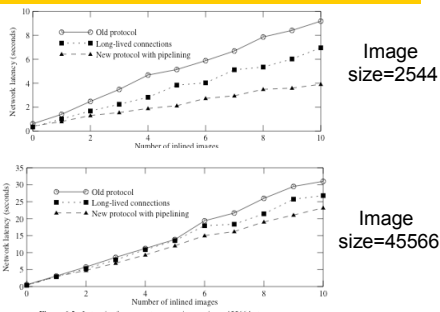


Figure 6-2: Latencies for a remote server, image size = 4556 bytes

djw // CSE 461, Fall 2009

Lhttp9

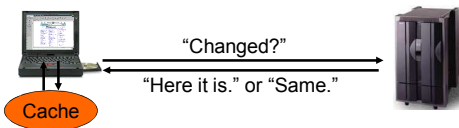
Caching

- It is faster and cheaper to get data that is closer to here than closer to there.
- “There” is the origin server. 2-5 RTT
- “Here” can be:
 - Local browser cache (file system) (1-10ms)
 - Client-side proxy (institutional proxy) (10-50)
 - Content-distribution network (CDN -- “cloud” proxies) (50-100)
 - Server-side proxy (reverse proxy @ origin server) (2-5RTT)

djw // CSE 461, Fall 2009

Lhttp10

Browser Caches



- Bigger win: avoid repeated transfers of the same page
- Check local browser cache to see if we have the page
- GET with If-Modified-Since makes sure it's up-to-date

djw // CSE 461, Fall 2009

Lhttp11

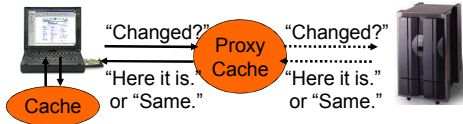
Consistency and Caching Directives

- Browsers typically use heuristics
 - To reduce server connections and hence realize benefits
 - Check freshness once a “session” with GET If-Modified-Since and then assume it's fresh the rest of the time
 - Possible to have inconsistent data.
- Key issue is knowing when cached data is fresh/stale
 - Otherwise many connections or the risk of staleness
- Caching directives provide hints
 - Expires: header is basically a time-to-live
 - Also indicate whether page is cacheable or not

djw // CSE 461, Fall 2009

Lhttp12

Proxy Caches

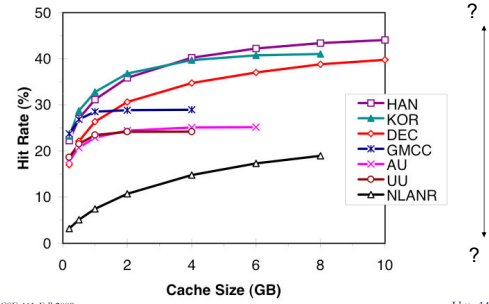


- Insert further levels of caching for greater gain
- Share proxy caches between many users (not shown)
 - If I haven't downloaded it recently, maybe you have
- Your browser has built-in support for this

djw // CSE 461, Fall 2009

Lhttp13

Proxy Cache Effectiveness



djw // CSE 461, Fall 2009

Lhttp14