

Typical Workload (Web Pages)



- Multiple (typically small) objects per page
- File sizes
- Heavy-tailed
 - · Pareto distribution for tail
 - Lognormal for body of distribution
- Embedded references
- •Lots of small objects versus TCP
- 3-way handshake
- Lots of slow starts
- Extra connection state
- Number of embedded objects also Pareto $Pr(X>x) = (x/x_m)^{-k}$
- · This plays havoc with performance. Why?
- · Solutions?

6

HTTP 0.9/1.0



- One request/response per TCP connection
 - Simple to implement
- Short transfers are very hard on TCP
 - Multiple connection setups → three-way handshake each time
 - · Several extra round trips added to transfer
 - Many slow starts low throughput because of small window
 - · Never leave slow start for short transfers
 - Loss recovery is poor when windows are small
 - Lots of extra connections increase server state and processing overhead

7

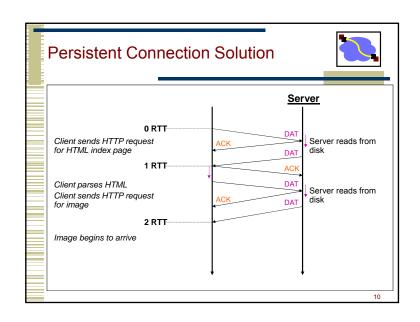
Single Transfer Example Client Server 0 RTT SYN Client opens TCP SYN 1 RTT connection Client sends HTTP request Server reads from for HTML 2 RTT Client parses HTML Client opens TCP connection SYN 3 RTT Client sends HTTP request Server reads from 4 RTT DAT Image begins to arrive

HTTP 1.1



- Multiplex multiple transfers onto one TCP connection
 - Avoid handshake and slow start when getting multiple objects from same server
- Transfers can be pipelined, i.e., multiple outstanding requests
 - · Requests are handled in FIFO order by server
- How to identify requests/responses?
 - Delimiter → Server must examine response for delimiter string
 - Content-length and delimiter → Must know size of transfer in advance
 - Block-based transmission → send in multiple length delimited blocks
 - Store-and-forward → wait for entire response and then use content-length

9

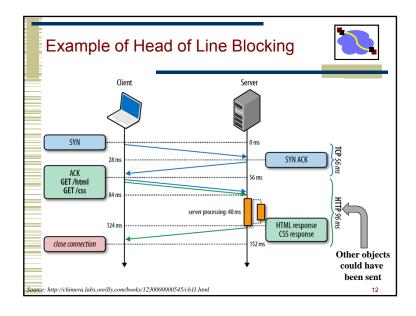


Some Challenges with HTTP 1.1



- Head of line blocking: "slow" objects can delay all requests that follow
 - E.g., objects from disk versus objects in cache
 - Single "slow" object can delay many "fast" objects
- Browsers open multiple TCP connections to achieve parallel transfers
 - Increases load on servers and network
- · HTTP headers are big
 - · Cost higher for small objects
- · Embedded objects add RTT
 - · With small objects, RTT dominates

11

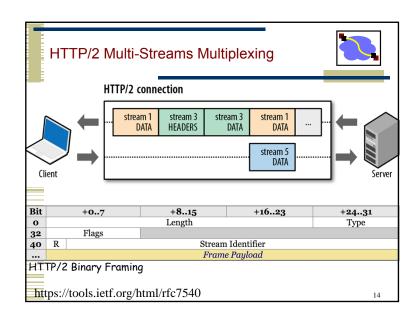


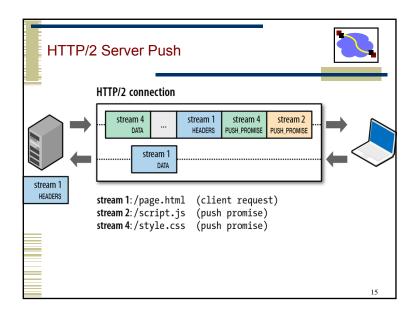
HTTP 2.0 to the Rescue

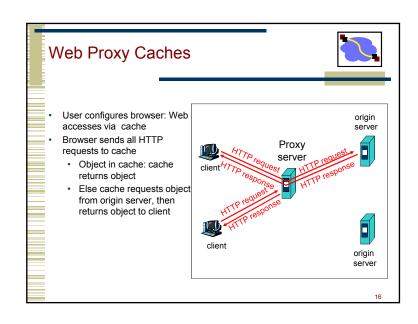


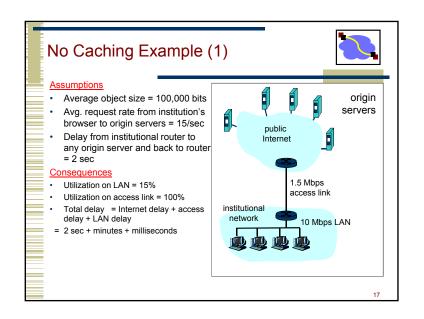
- Can multiplex many requests over a TCP connection AND
 - Responses are carried over flow controlled streams avoids HOL blocking
 - Streams can be prioritized by client based on how critical they are to rendering
 - ≈ multiple prioritized parallel TCP streams
 - Also: fewer handshakes and more traffic (help congestion control)
- HTTP headers are compressed
- A PUSH features allows server to push embedded objects to the client without waiting for a client request
 - · Avoids an RTT
 - · What is the challenge?
- Default is to use TLS fall back on 1.1 otherwise

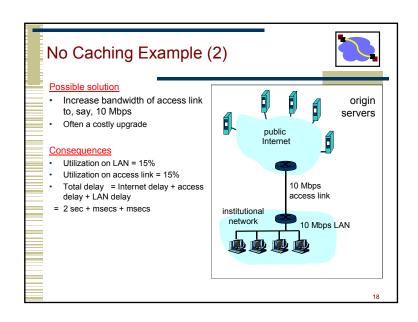
13

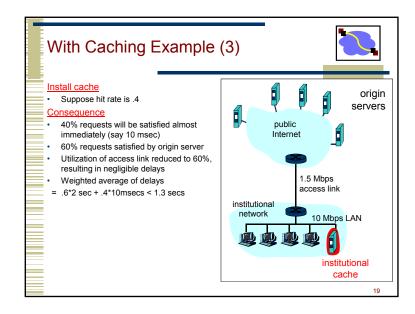


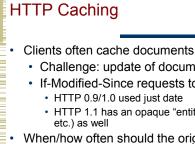


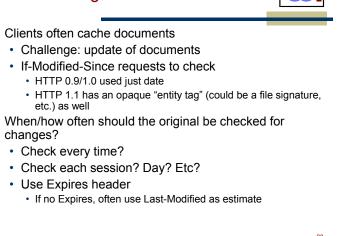












Problems · Fraction of HTTP objects that are cacheable is dropping Whv? · Major exception? This problem will not go away Dynamic data → stock prices, scores, web cams CGI scripts → results based on passed parameters Other less obvious examples SSL → encrypted data is not cacheable Most web clients don't handle mixed pages well → many generic objects transferred with SSL Cookies → results may be based on past data Hit metering → owner wants to measure # of hits for revenue, etc. What will be the end result?

