Network Applications: HTTP/1.1/2; Operational Analysis

Qiao Xiang, Congming Gao, Qiang Su

https://sngroup.org.cn/courses/cnnsxmuf25/index.shtml

10/09/2025

Outline

- Admin and recap
- □ HTTP
 - HTTP "acceleration"
 - Operational analysis

Admin

- □ Lab assignment 2 due Oct. 14
- □ Lab assignment 3 to be posted

Outline

- Admin and recap
- □ HTTP/1.0
- > HTTP "acceleration"

Recap: Substantial Efforts to Speedup Basic HTTP/1.0

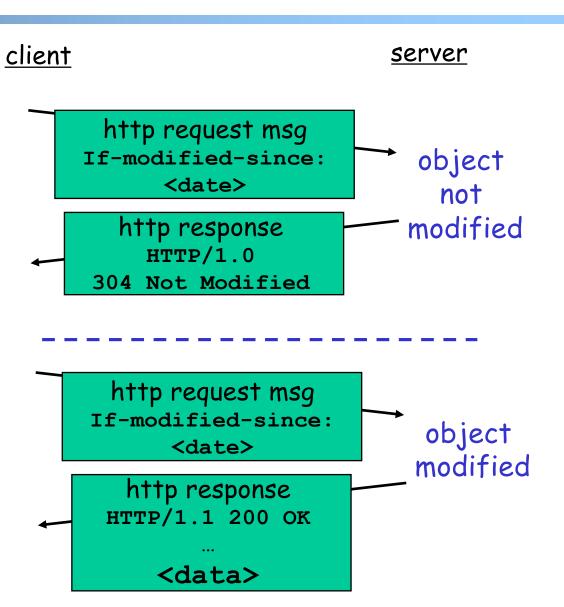
- Reduce the number of objects fetched [Browser cache]
- Reduce data volume [Compression of data]
- Header compression [HTTP/2]
- Reduce the latency to the server to fetch the content [Proxy cache]
- Remove the extra RTTs to fetch an object [Persistent HTTP, aka HTTP/1.1]
- Increase concurrency [Multiple TCP connections]
- Asynchronous fetch (multiple streams) using a single TCP [HTTP/2]
- Server push [HTTP/2]



Browser Cache and Conditional GET

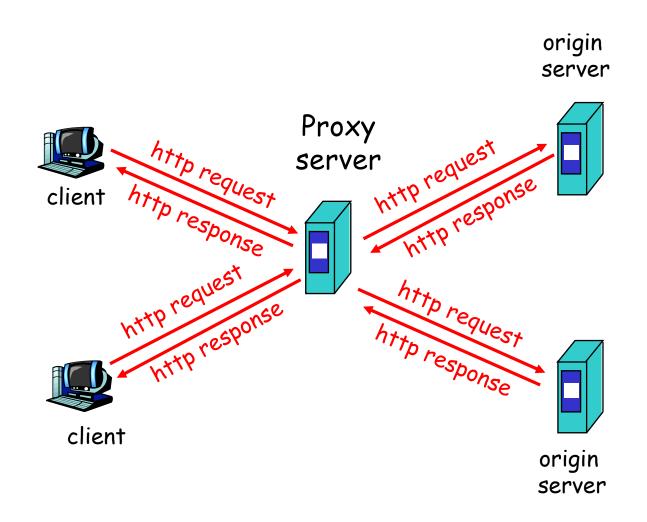
- Goal: don't send object if client has up-to-date stored (cached) version
- server: response contains no object if cached copy upto-date:

HTTP/1.0 304 Not Modified



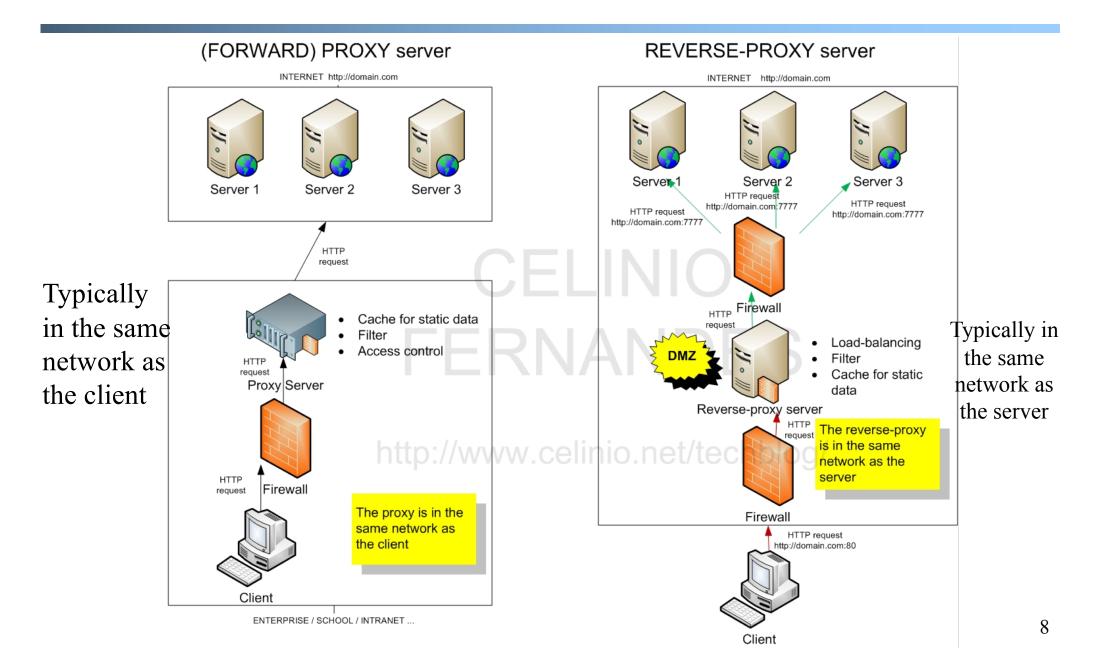
Web Caches (Proxy)

Goal: satisfy client request without involving origin server



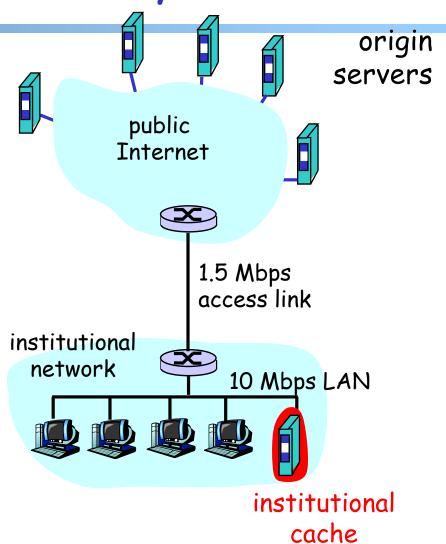
http://www.celinio.net/techblog/?p=1027

Two Types of Proxies



Benefits of Forward Proxy

- Assume: cache is "close" to client (e.g., in same network)
- smaller response time: cache "closer" to client
- decrease traffic to distant servers
 - link out of institutional/local ISP network often is bottleneck



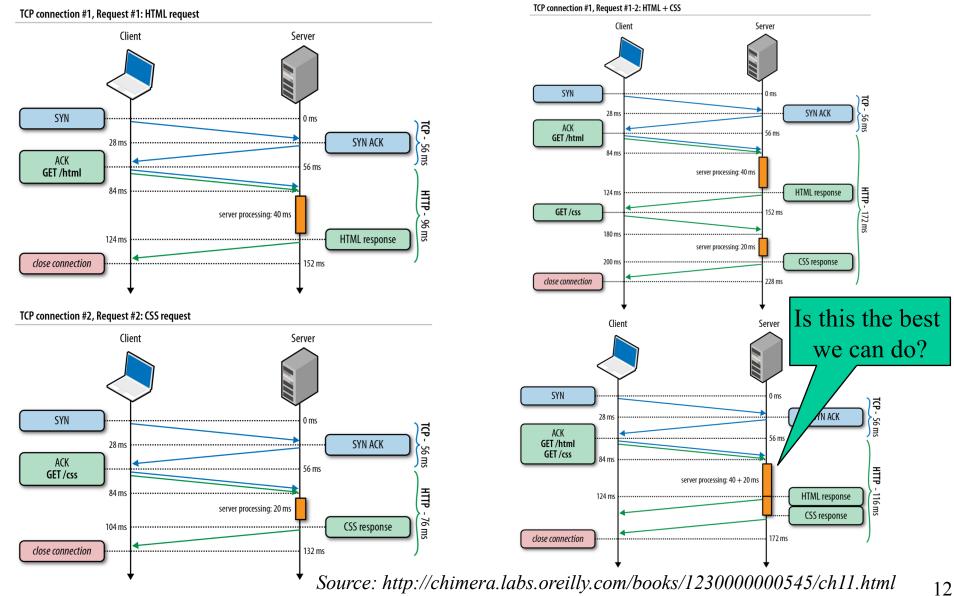
No Free Lunch: Problems of Web Caching

- The major issue of web caching is how to maintain consistency
- □ Two ways
 - o pull
 - Web caches periodically pull the web server to see if a document is modified
 - o push
 - whenever a server gives a copy of a web page to a web cache, they sign a lease with an expiration time; if the web page is modified before the lease, the server notifies the cache

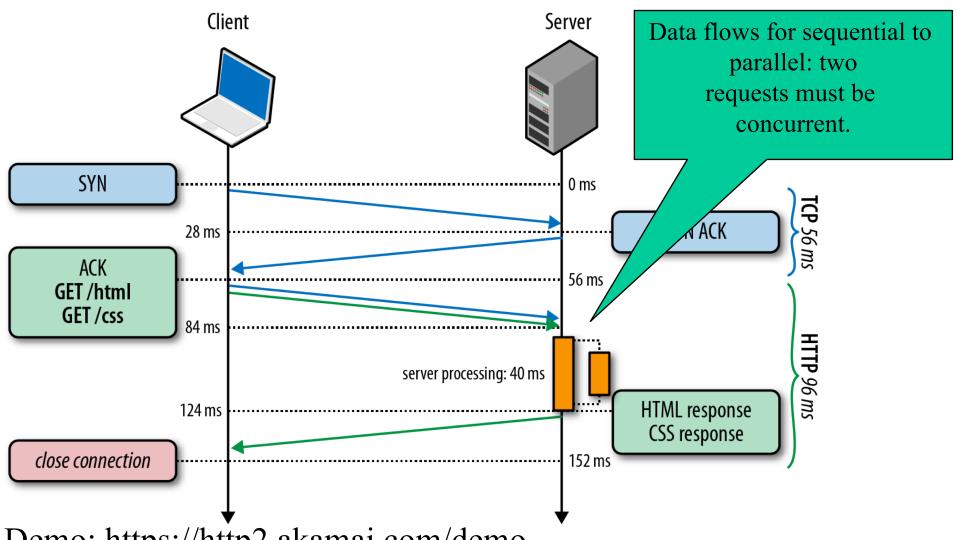
HTTP/1.1: Persistent (keepalive/pipelining) HTTP

- HTTP/1.0 allows a single request outstanding, while HTTP/1.1 allows request pipelining
 - On same TCP connection: server parses request, responds, parses new request, ...
 - Client sends requests for all referenced objects as soon as it receives base HTML
- Benefit
 - Fewer RTTs
 - See Joshua Graessley WWDC 2012 talk: 3x
 within iTunes

HTTP/1.0, Keep-Alive, Pipelining



HTTP/2 Basic Idea: Remove Head-of-Line Blocking in HTTP/1.1

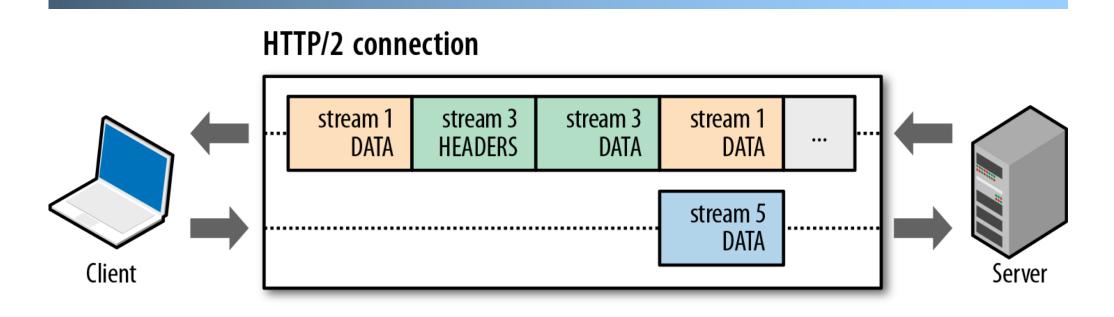


Demo: https://http2.akamai.com/demo

Observing HTTP/2

- export SSLKEYLOGFILE=/tmp/keylog.txt
- □ Start Chrome, e.g.,
 - Mac: /Applications/Google
 Chrome.app/Contents/MacOS/Google Chrome
 - Ubuntu: firefox
- □ Visit HTTP/2 pages, such as https://www.tmall.com
- Wireshark:
 - Mac: Wireshark -> preferences -> protocol -> TSL
 (pre)-master-secret log file name
 - Ubuntu: edit -> perferences -> protocol -> SSL
 (pre)-master-secret log file name

HTTP/2 Design: Multi-Streams



Bit		+07	+815	+1623	+2431
0	Length Type			Type	
32		Flags			
40	R	R Stream Identifier			
•••		Frame Payload			

HTTP/2 Binary Framing

https://hpbn.co/http2/

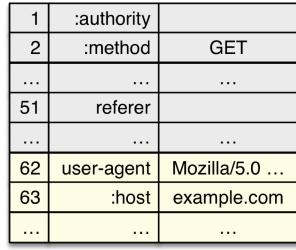
https://tools.ietf.org/html/rfc7540

HTTP/2 Header Compression

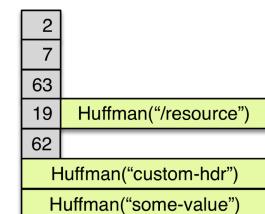
Request headers

:method	GET	
:scheme	https	
:host	example.com	
:path	/resource	
user-agent	Mozilla/5.0	
custom-hdr	some-value	



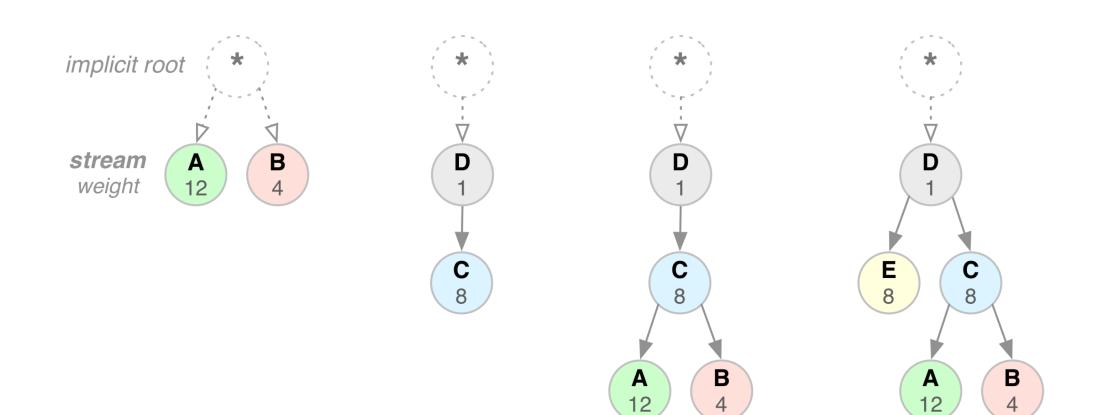


Encoded headers

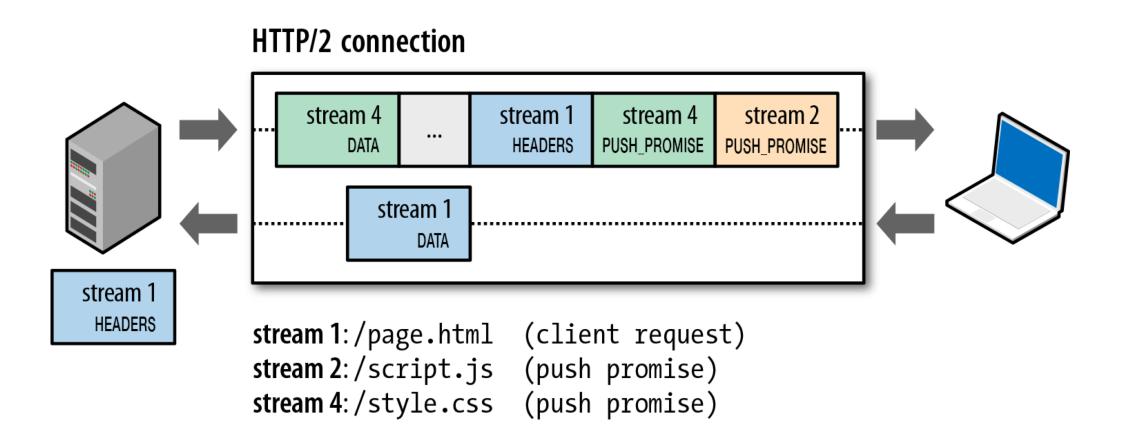


Dynamic table

HTTP/2 Stream Dependency and Weights



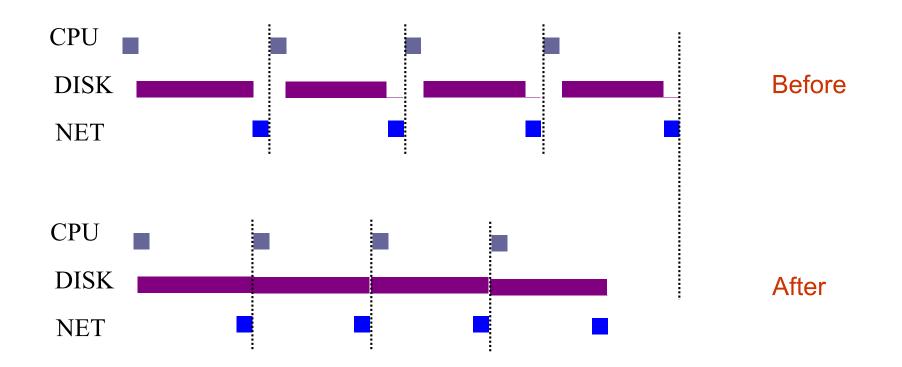
HTTP/2 Server Push

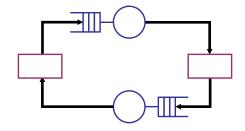


Outline

- Admin and recap
- □ HTTP
 - HTTP "acceleration"
 - Operational analysis

Goal: Best Server Design Limited Only by Resource Bottleneck





Some Questions

- When is CPU the bottleneck for scalability?
 - So that we need to add helper threads
- □ How do we know that we are reaching the limit of scalability of a single machine?
- These questions drive network server architecture design
- Some basic performance analysis techniques are good to have

Background: Little's Law (1961)

- □ For any system with no or (low) loss.
- Assume
 - o mean arrival rate λ , mean time R at system, and mean number Q of requests at system

R, Q

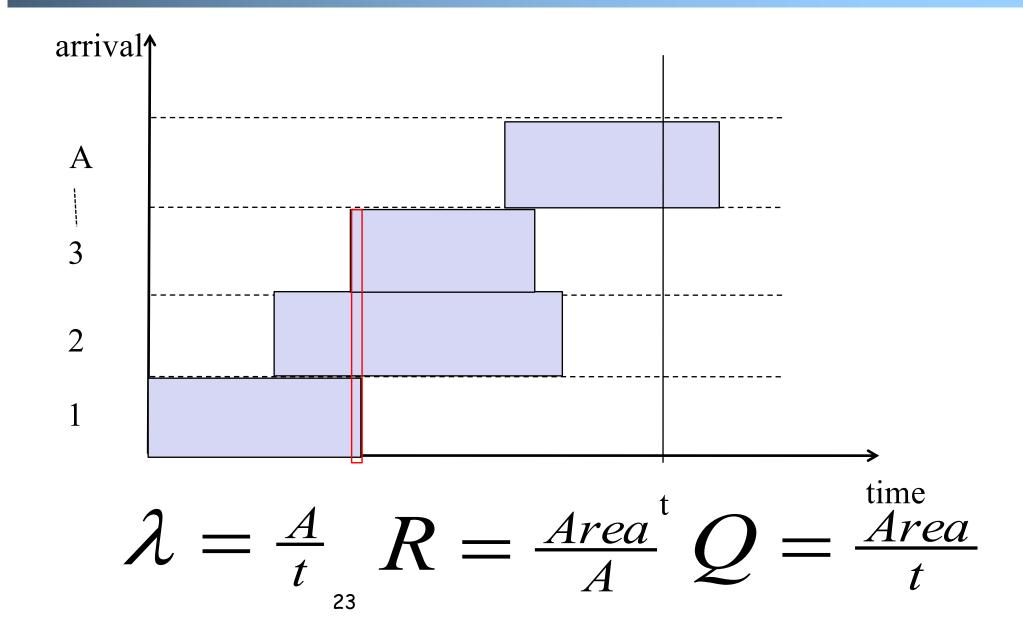
 \square Then relationship between Q, λ , and R:

$$Q = \lambda R$$

Example: XMU admits 3000 students each year, and mean time a student stays is 4 years, how many students are enrolled?

Little's Law: Proof

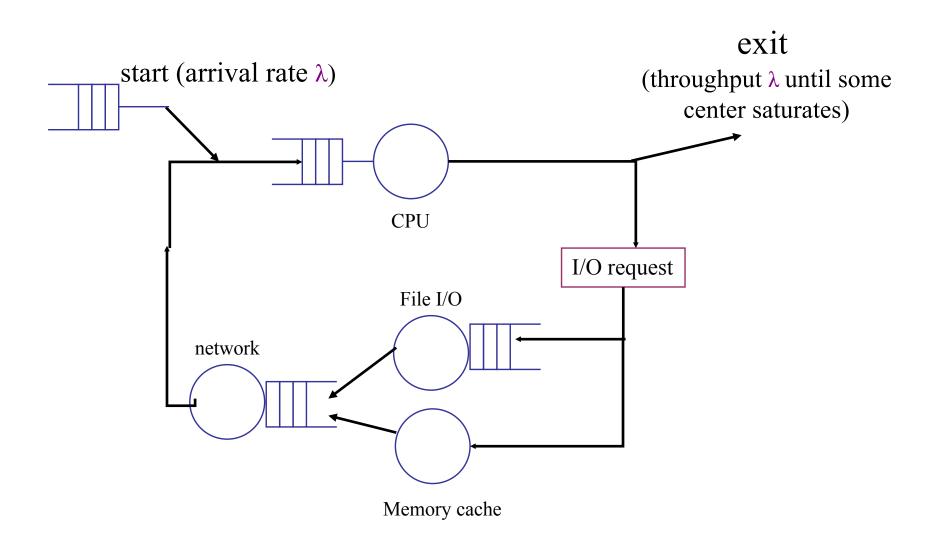
$Q = \lambda R$



Operational Analysis

- Relationships that do not require any assumptions about the distribution of service times or inter-arrival times
 - Hence focus on measurements
- □ Identified originally by Buzen (1976) and later extended by Denning and Buzen (1978).
- □ We touch only some techniques/results
 - o In particular, bottleneck analysis
- More details see linked reading

Under the Hood (An example FSM)



Operational Analysis: Resource Demand of a Request

CPU V_{CPU} visits for S_{CPU} units of resource time per visit Network V_{Net} visits for S_{Net} units of resource time per visit Disk V_{Disk} visits for S_{Disk} units of resource time per visit Memory V_{Mem} visits for S_{Mem} units of resource time per visit

Operational Quantities

- T: observation interval
- Bi: busy time of device i
- □ i = 0 denotes system

Ai: # arrivals to device i

Ci: # completions at device i

arrival rate
$$\lambda_i = \frac{A_i}{T}$$

Throughput
$$X_i = \frac{C_i}{T}$$

Utilization
$$U_i = \frac{B_i}{T}$$

Mean service time
$$S_i = \frac{B_i}{C_i}$$

Utilization Law

Utilization
$$U_i = \frac{B_i}{T}$$

$$= \frac{C_i}{T} \frac{B_i}{C_i}$$

$$=X_{i}S_{i}$$

- The law is independent of any assumption on arrival/service process
- Example: Suppose NIC processes 125 pkts/sec, and each pkt takes 2 ms. What is utilization of the network NIC?

<u>Deriving Relationship Between</u> R, U, and S for one Device

Assume flow balanced (arrival=throughput), Little's Law:

$$Q = \lambda R = XR$$

□ Assume PASTA (Poisson arrival--memory-less arrival--sees time average), a new request sees Q ahead of it, and FIFO

$$R = S + QS = S + XRS$$

According to utilization law, U = XS

$$R = S + UR \longrightarrow R = \frac{S}{1-U}$$

Forced Flow Law

Assume each request visits device i Vi times

Throughput
$$X_i = \frac{C_i}{T}$$

$$= \frac{C_i}{C_0} \frac{C_0}{T}$$

$$=V_iX$$

Bottleneck Device

Utilization
$$U_i = X_i S_i$$

$$= V_i X S_i$$

$$= XV_i S_i$$

- Define Di = Vi Si as the total demand of a request on device i
- □ The device with the highest Di has the highest utilization, and thus is called the bottleneck

Bottleneck vs System Throughput

Utilization
$$U_i = XV_iS_i \le 1$$

$$\rightarrow X \leq \frac{1}{D_{\text{max}}}$$