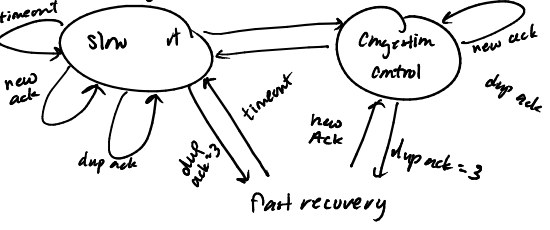


Top Congestion Control

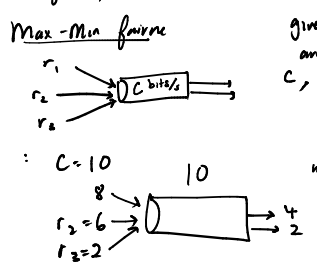


Detecting Congestion
 Packet delays
 Packet loss (fail-safe)
 Duplicate loss: isolated
 Timeout: much more
 (adjust rate)
 Issues w/ TCP

implies: $\sqrt{\frac{3}{2} \frac{1}{RTT_{up}}}$

Implications?
 • TCP unfair if heterogeneous RTTs.
 • High speed TCP numbers unreasonable
 ↳ adapt by router assisted approaches?

Router Assisted
 Ensure flows into flows,
 Ctrl.



① $\frac{C}{3} = 3.33$, but r_3 needs only 2
 • can service r_3 , remove from accounting $C = C - r = 8; N = 2$
 ② $\frac{C}{2} = 4$ can't service all, so give fair s

FR vs FIFO
 (+) Isolation: Cheating flows don't benefit
 (+) bandwidth share doesn't depend on RTT
 (+) flows pick rate adjustment scheme
 (-) more complex than FIFO: per flow queue/state
 additional per-packet bookkeeping

DNS Caching
 • reduces load at all levels, reduces delay for client
 • effective b/c Top level servers rarely change
 • popular sites visited often → local DNS server often has info cached

HTTP (web's app layer protocol, uses TCP as underlying transport protocol)

Non persistent
 each request/response pair over a separate TCP connection
 brand new connection for each object
 ↳ each suffers a delivery delay of 2 RTTs:
 1 to establish connection, 1 to retrieve object

Concurrent Requests & Responses

use multiple connections in parallel network

Persistent Connections

Maintain TCP connection across multiple requests.

Pipelined Requests/Responses

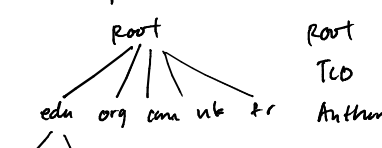
• batch requests & responses to reduce # of packets.
 • multiple requests in 1 TCP segment

Evaluate! Getting n small objects? ← time dominated by latency

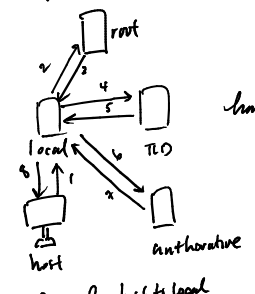
1. one at a time → ~ 2nRTT
 2. M concurrent → ~ 2[n/m]RTT
 3. Persistent → ~ (n+1)RTT
 4. pipelined → ~ 2RTT
 5. pipelined/persistent → ~ 2RTT the first time, RTT after
- Getting N Large Objects
1. 1 at a time: ~ n F/B
 2. M concurrent: ~ [n/m] F/B
 3. Pipelined and/or pers: ~ n x F/B
- ← the only thing that helps is getting more bandwidth

DNS Goals (Application Layer)

• fast lookups, hierarchical, unique



records servers store RRs.
 name, value, type, TTL
 Type (→ Address)
 : hostname
 value: IP address
 type: NS
 name: domain
 val: name of dns server for domain



Ethernet (old) (Link Layer)

Random access protocols

- CSMA
- Carrier sense: check if already sending. If so, wait
- collision detection: listen while transmitting. If collision, abort, send jam signal
- random access: binary exponential backoff. wait random time before trying again

Ethernet frames

how does link layer determine where frames end/begin?

↳ count bytes! → need frames.

Frames: Sentinal bit framing:



bit stuffing!!!

• sender inserts a 0 after 5 1's

• receiver always removes 0 after 5 1's

HTTP Caching

Where?

1. Clients: forward proxies
 ↳ reduce traffic, decrease latency
2. Server: decrease server load network load

HTTP ex: A retrieving files

F & G from site B. RTT

bandwidth b/w A & B is 10Msec = .01sec

Bandwidth b/w A & B is 10Mbps

• Time to retrieve blk files?

① Sequential, non persistent?

RTT (2 SYN/ACK + 2 DATA/ACK)

+ 2 (Transmission time)

= 4 x .01 + 2 (1mb / 10Mbps) = .04 + .2 = .24

② Concurrent, nonpersistent

2RTT (1 SYN/ACK + 1 DATA/ACK)

+ 1 Transmission time

Switched Ethernet: concurrent communication

Mac Addresses (link layer)

- associated w/ network adapter
- flat name space of 48 bits in HEX.
- social security #!
- portable, stays the same....
- * used to get packet b/w interfaces on the same network

vs... IP Addresses

• configured, learned dynamically, partial address used to get a packet to dest. subnet.

"Flooding" of broadcast Ethernet

- doesn't use LS/DV b/c not scalable (MAC addresses cannot aggregate like IP); plug & play!

How it works: sender transmits frame to broadcast link, frame contains MAC address.

↳ if dest. matches receiver's MAC addy, or broadcast (FF:FF:FF:FF:FF:FF)

PASS FRAME TO NETWORK LAYER!

SPANNING TREE APPROACH

* Messages: (y, d, x)

propose y as root, distance d, and from node X

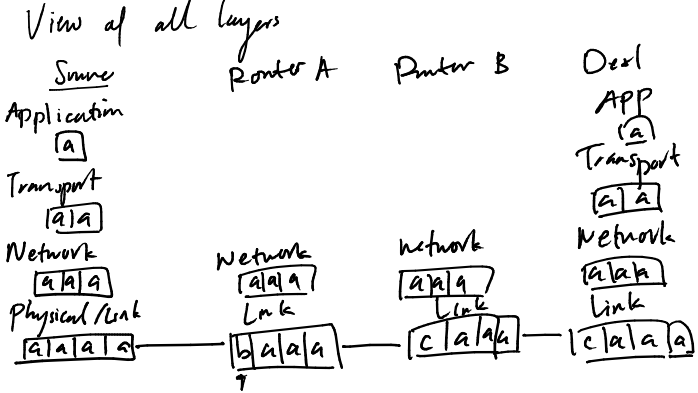
• switches elect node w/ smallest identifier as root

How? Switched Ether

① build spanning tree for loop free flooding

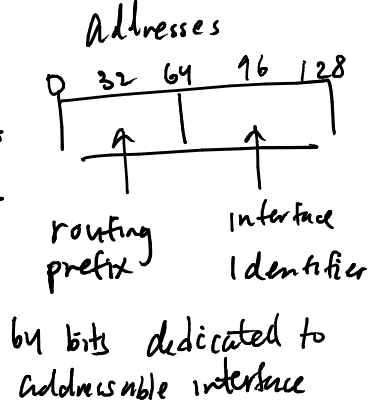
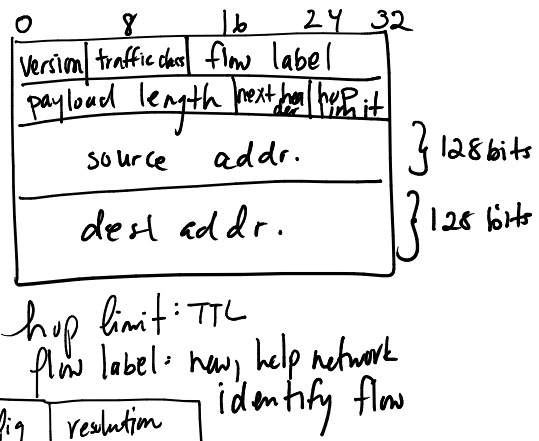
② "self learning" switches....

optimizes... switches can each dest w/o flooding



IPv6
 why? exhausted space.
 Differences?
 • addresses 128 bits, not 32 bits
 • headers different
 • address management

Same?
 • routing protocols
 • longest prefix / shortest path



Naming

- Application layer: URLs, domain names
- Network: IP addresses: host's network location
- Link layer: MAC addresses: host identifier

all 3 for E2E

| Layer | ex | Structure | Config | resolution |
|---------|-------------|----------------|------------|------------|
| APP | bbc.com | organizational | manual | ↕ DNS |
| Network | 123.45.6.78 | topological | DHCP | ↕ ARP |
| Link | 45-CC-42... | Flat | hard coded | |

ARP & DHCP Link Layer Discovery Protocols used by Network Layer

2 functions:

- ① Discovery of local end hosts, for communication
 bit hosts on same LAN
- ② Bootstrap communication w/ remote hosts
 • what's my IP address?
 • who/where is my DNS server?
 • who/where is my first hop router?

ideas: broadcast, soft state, caching

DHCP: used by host to discover:
 - netmask, IP address for DNS server(s), router
 P addresses, delays for f.h

Discovery Mechanisms
ARP/DHCP: broadcast
 • flooding doesn't scale!
 • zero configuration
 • no centralized point of failure

DNS:
 • scalable (no floods)
 • manual config: (local, root servers, etc)

ARP table: IP address → MAC addr.
 • maintained by every host
 • consult table when sending packet, broadcast IP address, receiver responds w/ MAC address to IP addr.

To reach destinations
 1) need own IP: DHCP
 need local DNS: DHCP

Send Packet
 • Same subnet: ARP
 Use MAC addr of dest
 • Other subnet: ARP+DHCP
 (use MAC address of first hop router)

Datacenters

- scalability baseline req.
- more emphasis on performance
- less on heterogeneity & interoperability

Solutions
 • extend DV/LS?
 (-) Scales poorly... N destinations, O(N) routing entries/m.syp

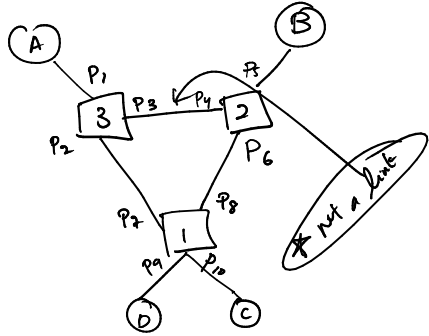
Recap

MAC?
 own: hard coded
 others: ARP (given IP addr)

IP.
 own: DHCP
 others: DNS (given domain name)

use netmask to tell if host on same or diff subnet, or DHCP

Spanning Tree, Self Learning



Switch A wants to talk to B (cache empty) ... state of forwarding table after ARP request complete?

| S1 | | S2 | | S3 | |
|------|------|----|----|----|----|
| dest | port | d | P | d | P |
| A | p7 | A | P6 | A | P1 |

After B responds?

| S1 | | S2 | | S3 | |
|------|------|----|----|----|----|
| dest | port | d | P | d | P |
| A | p7 | A | P6 | A | P1 |
| B | p8 | B | P5 | B | P2 |

C sends to A

| S1 | | S2 | | S3 | |
|------|------|----|----|----|----|
| dest | port | d | P | d | P |
| B | p8 | B | P5 | B | P2 |
| C | p9 | C | P6 | C | P2 |

(all newest entries move up, replaces oldest b/c max 2 entries in this problem)

HTTP continued

(a) sequential, 1 persistent connection

3 RTT (1 syn/ack + 2 Data/Ack)

+ 2 transmission time = $.03 + .02 = .23 \text{ sec}$

(d) pipelined, persistent

2 RTT (1 syn/ack + 1 data/ack) + 2 transmission time

= $.02 + .2 = .22$

929

(1) self discovery (DHCP)

IP | UDP | DHCP Disc.

IP: source: 0.0.0.0

Dest: 255.255.255.255

add a link layer frame
broadcast!

LL | IP | UDP | DHCP disc.

mac addr: FF:FF:...

(2) Machine running DHCP server

prepares offer w/

IP, DNS IP, Default Gateway IP, subnet mask

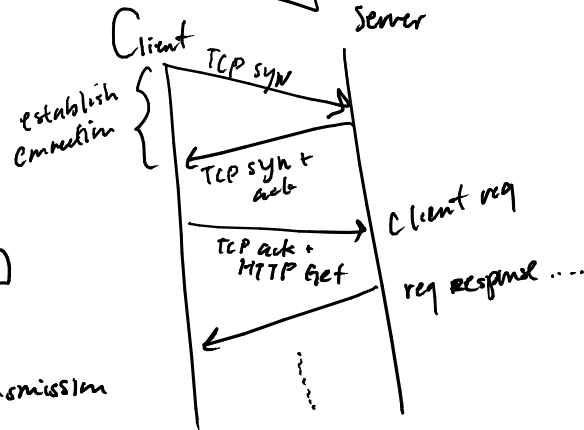
..... Client accepts offer by broadcasting 'request' msg, server ACKs

(3) connect to google?

- need MAC address of router

↳ broadcast ARP req, Gateway Router responds

LL | ARP



(4) get dest ip address! (DNS)

DNS response w/ Google's IP address. (unless cached)

(5) use HTTP to communicate

TCP is transport layer protocol used

LL | IP | TCP | HTTP

ethernet | IP | TCP | HTTP

Router's mac

source: me
dst: google IP

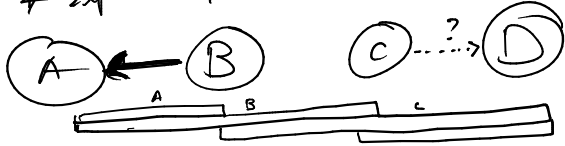
Wireless

* Hidden terminals *



- cannot use carrier sense!
Needs to sense at receiver

* Exposed Terminals



- if B talks to A, C doesn't transmit to D b/c of carrier sense... it would have worked!

MACA: multiple access w/ collision avoidance

- if no CTS, assume collision
- if you hear a CTS, wait for ACK.
- if hear RTS, no CTS, send

Gen Protocol Rules for BGP:

