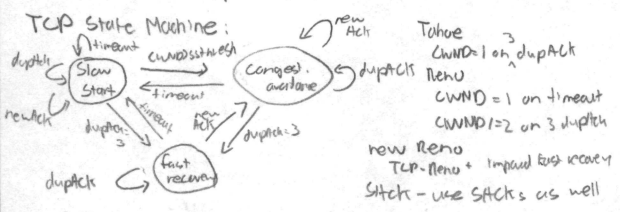


Congestion Control

Internet traffic is bursty, may cause delays/drops
 TCP has endpoints adjust rate, feedback by ACK/drops
 Goals: Discover bottleneck bandwidth, adjust to variations in bandwidth and share bandwidth between flows
 TCP uses dynamic adjustment for rates but assumes that endpoints will be good citizens

Congestion Window (CWND): Bytes w/o overflowing routers
 Advertised window (RWND): Receiver → sender, buffers
 Sender-side window = $\min(\text{CWND}, \text{RWND})$

- ① Slow start - Grow window exponentially by +1 for each ACK received, until a loss is experienced or ssthresh is reached
 - ② AIMD - Increase window by $\frac{1}{\text{dupack}}$ per ACK and $\text{CWND} = \text{CWND}/2$ on packet loss
- On timeout, set $\text{ssthresh} = \text{CWND}/2$, $\text{CWND} = 1$
 AIMD allows us to achieve efficiency and fairness with dupACKs, if 3 then $\text{ssthresh} = \text{CWND}/2$, $\text{CWND} = \text{CWND}/2$
 Fast Recovery - If dupack count = 3, $\text{ssthresh} = \text{CWND}/2$, $\text{CWND} = \text{ssthresh} + 3$
 $\text{CWND} + 1$ for each dupack and $\text{CWND} = \text{ssthresh}$ w/ new ACK to exit



TCP Throughput Equation
 Average $\frac{1}{\text{RTT}}$ packet, $\frac{1}{\text{RTT}}$ between dupacks
 Throughput = $\frac{1}{\text{RTT}}$

- ① TCP unfair if we have different RTT's
- ② TCP requires extremely low dup rate to have high speed, pipelined. we can increase AIMD faster after thresh, use multiple, or vector-assisted
- ③ Can follow TCP throughput to send at constant rate

TCP Problems
 ① Will cut rate even for corruption rather than congestion
 ② Short flows never leave slow start, hard to transfer volume
 ③ TCP fills up queues by default, which delays everyone
 ④ Can cheat by increasing window risks, open mult. connections, start w/ large CND
 ⑤ CC = reliability (maybe try ssthresh?), complicates evolution

Routers can fix by telling endpoints about congestion, give rate to send it and by enforcing fair sharing

Routers and Fairness
 Packets are classified into flows and served fairly on definition
 Max-min fairness: $a_i = \min(f_i, r_i)$, f such that $\sum(a_i) = C$
 So if you don't get full demand, no one gets more

Fair Queuing: Compute time for last bit to leave router if served bit by bit and serve by order of deadlines
 Allows cheat prevention, bandwidth not dep on RTT, flows can pick any rate adjustment scheme they want, but more complex

ECN: Packets have "red", f fair share, endpoints get info to f
 ECN: single bit in header for congestion, can interpret as drop
 Helps distinguish corruption (congestion and early indicate congest.)
 Perhaps one can also charge money per ECN
 Highlights importance of performance, shows we can do better
 RC3: Transmit high but use priorities (CWND #1)

Application Layer

Host Names and Addresses (DNS)
 Addresses used by protocols ("where"), host names by people ("who")
 DNS maps one to other, directory service that decouples them
 we want scalability, availability, correctness and performance

Partition namespace, distribute partition administration and resolution
 Hierarchical namespace, administration and servers
 Namespace: "top level domains", domain = subtree, name = leaf to root path
 Administration: Each zone corresponds to admin for that part
 Servers: Root server → top-level domain servers → authoritative DNS serv
 All know DNS addr, DNS knows TTL, authoritative DNS stores name-to-addr ("resource records") for own, anyone can discover server in protocol.
 Root server replicated by anycasting, deliver packet to closest machine
 DNS resource records → (name, value, type, TTL)
 Type A (Address) name = hostname, value = IP addr, type = MX (mail exchange)
 Type NS (Name server) name = domain, value = name of domain DNS server, value = name of mail server
 Inset RR by register domain name w/ IP addr name of auth name server
 then RR added to TLD server, store RR in your own server

HTTP and the Web

Uniform Resource Locator (URL) protocol://host-name[:port]/direct-path/resource
 protocol: http, ftp, nat: DNS, IP, port: default to standard, direct path: browser, file sys resource
 HTTP (TCP, Port 80)

Client-server architecture (Server "always on" and "well known")
 Synchronous request-reply protocol, stateless, HTML format
 ① Est connection (TCP SYN, SYNACK) ② Client request (TCP ACK + HTTP GET)
 Request line: method, resource, ... and protocol version
 Request headers: provide information or modify request
 Body: optional data (e.g. to "POST" data to the server)
 ③ Request response ④ TCP FIN/FINACK
 Status line: protocol version, status code, status phrase
 Response headers: provide information, Body: optional data

Statelessness and Cookies
 Treat each request-response independently
 Improves scalability on server-side (failure easier, higher utilization)
 but some applications need persistent state → cookies
 Cookies: client stores small state for server and sends in future
 Performance goals (User blind availability) content provider (not user, network oriented)
 Improve HTTP to compensate for TCP, caching, replication, economics of scale

Concurrent Requests and Responses
 Use multiple connections in parallel w/o regard to order
 Good for client and content provider, bad for network

Persistent Connections
 Maintain TCP connection across multiple requests
 Avoid setup/teardown overhead, get accurate RTT, allow TCP CND ↑
 This is the default in HTTP/1.1

Pipelined Requests and Responses
 Batch together requests and responses to reduce number of packets,
 allowing for multiple requests in one TCP segment
 Comparison Between Small and Large Objects

Small	vs.	Large
Time dominated by latency		Time dominated by bandwidth
Measurement $\sim 2nRTT$		$\sim nF/B$
Persistent $\sim (n+1)RTT$		$\sim [n/m] F/B$
Pipelined $\sim 2RTT$		$\sim nF/B$
		$\sim 2RTT$, then RTT

Caching
 Modify GET to include "If-modified-since", reply not or OK if latest
 Expires (how long cache TTL), NO-CACHE (always ask server)
 Reverse Proxies: close to server to decrease server load, content provider
 Forward Proxies: close to clients to reduce network delay, latency, ISP
 Replication: Copy across many machines, load balance by DNS response
 Content Dist Networks: caching and replication as a service
 pull caching at client proxy and push replication "idle material"
 In cost-effective content delivery, usually see multiple sites based on same shared physical infrastructure allowing use of statistical mult., economies of scale, amortization of human operator costs

Data Link Layer

Broadcast Medium Access
Random Access MAC Protocols
 Transmit w/o coordination but know how to detect/recover from collisions
 Aloha signaling: site sends packet to hub via wire, hub broadcasts packet
Carrier Sense Multiple Access (CSMA)
 Transmit entire frame if channel idle else defer
 CSMA/CD: collision detected in channel time and abort
 Limits on max wire length
 Carrier sense: wait for idle link
 Collision detection: listen while transmit, jamming signal if collision
 Random access: binary exponential back-off - (0, 1, 2, 3) s delay with collision
 $E \sim \frac{P/B}{P/B + K}$ E as bandwidth increases, high-speed LAN → switched

Switched Ethernet

Point-to-point links from switches to hosts
Frame: Preamble (7 bytes clock sync, 1 byte start of frame), Data (CRC), Trailer (CRC)
 Preamble: 7 bytes clock sync, 1 byte start of frame
 Address: 6 bytes Type: 2 bytes protocol Data: 46-1500 bytes
 CRC: 4 bytes error
 Two hosts must exchange frames, want to know where frames start and end by link layer
 Sentinel bits (01111111 start, 11111111 end)
 bit stuffing, sender inserts 0 after 5 1's
 Receiver removes 0 after 5 1's

Medium Access Control (MAC) Address
 Flat namespace 48 bits (6 in hex)
 Unique, hard-coded in network adapter, 16 bits
 First 24 bits by vendor, last 24 assigned by vendor

Routing in Ethernet
 No DVLS due to scalability issues and broadcast Ethernet compatibility, play-n-play
 In broadcast, sender transmits frame on broadcast link and receiver link layer passes frame to network layer if destination MAC or port: FF:FF:FF:FF:FF:FF

Spanning Tree Protocol
 Need to pick a root and compute shortest paths to root to form spanning tree
 Send messages (Y, d, X) from X, Y bits with distance d

① Each switch purports self as root
 ② Switches update view of root
 If Y's id < current root's id root = Y and add 1 to shortest distance from neighbor
 Root switch sends periodic announcement
 Timeout (soft state) if no word from root

Flooding/Fowarding on Spanning Tree
 Flood packets to all parts on spanning tree
 Learn paths based on flooded packets for efficiency
 Map src MAC to incoming port in switch table and store with TTL

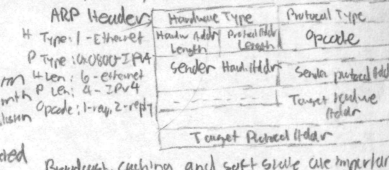
Ethernet → zero config, simple but inefficient, unpredictable

ARP and DHCP

ARP and DHCP
 Address Resolution Protocol and Dynamic Host Config Protocol are link layer discovery protocols that allow discovery of local endpoints, common to single or multiple hosts

DHCP
 Client to discover own IP addr, network, DNS server IP
 ① Client broadcast discovery FF:FF:FF:FF:FF:FF
 one or more local DHCP servers have info w/ IP addr pool, network, DNS server (UDP Port 67)
 ② One or more DHCP servers responds with "offer" (proposed IP, lease time)
 ③ Client broadcasts DHCP request message to specify accepted offer and echo parameters
 ④ Selected DHCP responds w/ ACK
 "soft state": addr allocations have lease period, server sets timer and client must request refresh, reset timer if refresh else reclaim address

ARP
 Table with (IP addr → MAC addr) pairs
 consult when sending packet else broadcast
 If dest is remote, look up first-hop router MAC (IPy, DHCP) network also by DHCP



Broadcast, caching and soft state are important

Steps in end-to-end communication

- ① Use DHCP to discover own IP addr, DNS IP, Router IP
- ② Resolve domain name by having local DNS server resolve DNS Eth, IP, UDP, DNS
 Use ARP to get MAC addr of router Eth, ARP
 Then send the DNS query to your DNS server
- ③ DNS server returns appropriate IP addresses
- ④ HTTP to get content from the server Eth, IP, TCP, HTTP

Wireless

Frequency and wavelength affect distance (free-space loss), penetration/retention, antenna size, energy propagation characteristics: broadcast medium, cannot receive while transmitting, signals don't always arrive intact

Free space Path Loss: $\left(\frac{4\pi d}{c}\right)^2 = \left(\frac{4\pi d f}{c}\right)^2$

Also consider multipath effects

The lower SNR, higher the bit error (power vs. info)

If loss, lower bitrate if free-space/multi-path fading, raise if external interference (chirp burst)

802.11 Architecture

Access points (AP) for specific channels
Broadcast beacon messages w/ SSID (same set ID) and MAC address periodically
Hosts scan all channels to discover AP's, then associate

CSMA/CA

Hidden Terminals

A, B, C
A and C cannot hear each other, so carrier sense ineffective

Exposed Terminal

A, B, C, D
B and C can hear other and may refrain from transmitting even though zero chance of collision from C's

No concept of global collision, collision only at receiver

Goal: Check if receiver can hear, tell initiating sender to stop

① Sender sends Request to Send (RTS) frame giving length of transmission and destination

② Receiver responds with Clear to Send (CTS) frame

Sender who hear shut up until ACK → over

③ Sender sends data ④ Receiver finishes w/ ACK

Can also prevent collisions by partitioning frequency spectrum into multiple channels, or partition space into non-overlapping cells

Multi-hop Wireless Ad-hoc Networks

Ideally w/ self-healing and multi-path routing

① Node cannot use all links at once due to sharing

② Overlapping ranges further decrease utilization, contention from RTS and CTS

③ TCP ACKs produce bidirectional traffic

Static and forward technique halves hop bandwidth, doubles latency and increases interference

Datacenters

Organization: servers in racks, headed by 'top of rack' switch

Aggregation switches connect TOR switches, to outside by core switches → 2x redundancy for fault tolerance

Why?: Scale changes, new applications: "big data", customer-facing revenue generating services

Same model: Cloud, and multi-tenancy

Scale: Need scalable design (no flooding), low cost, high utilization, tolerate frequent failure, automation

Same model: performance guarantees, isolation guarantees, portability (achieve through network virtualization)

Applications: high bandwidth "broadband", low latency is critical, must-also ("for") latency critical

Partition Aggregate "North-South"

Ext clients and datacenter query response

Front-end web servers, middle app servers, back-end database

Map Reduce "East-West"

Traffic between servers in datacenter ("big data" computation) may shift on small timescales (< minutes)

"Elephants" vs. "Mice" flows (55% flows, 1M, 5w the 35% bytes)

Bandwidth

Ideally, have all servers talk to others at full access (link rate solution: Build network of switches "Fabric") High bisection bandwidth

(split return to two, minimum bandwidth between any two)

Full bisection bandwidth b/w band in N node is $\frac{1}{2}$ times bandwidth of a single link

"scale up" is inefficient for traditional tree topology

"scale out" 'Clos' topology - multi-stage network where all switches have $\frac{1}{2}$ ports up, $\frac{1}{2}$ down all are spread

To be efficient, routing must be able to explore all paths

Now have single administrative domain, central core returns and endpoints, placement of traffic source/sink

Datacenters (cont.)

Forwarding

Per-processor load balancing spreads round-robin speeds traffic equally but incurs badly w/ TCP due to dupl/reordering, Multipath TCP

Per-flow load balancing (ECMP, equal cost multi-path)

Hash per-tuple allows following of single path but can be non-optimal due to elephants

Routing

Extending DUAL w/ ECMP allows simple reuse of skulls but is not scalable N → millions

① Topology-aware Addressing

Addresses embed location in regular topology "load coded" topology allows no computation instant entries

→ localized link failure detection

But VM migration or inconsistency may be problematic

② Centralize and Share Routes

No entries in switch, routing = broadcast from controller → server switches simple and scalable, end-points control route selection

scalability/delay/route issue

Transport Protocol Design

want low queue occupancy and high throughput

DC-TCP

Use ECN (treat like dup), react to event not presence of cong.

Switch: set ECN bit if queue len > K (independent)

Sender: Maintain avg of flight packets marked (d)

and adapt window to: $W \leftarrow (1 - \frac{d}{2}) W$

At each RTT, $F = \frac{\# \text{ marked pkts}}{\text{Total pkts}} \Rightarrow d \leftarrow (1 - \frac{d}{2}) + \frac{d}{2} F$

with $d = 1$, proves TCP-like behavior

Overall, ↓ latency by avoiding large buildup in queues,

maintain high throughput by proactive-retrying and redies

Variable in sending rates to ↓ queue buildup

Flow-completion Time (FCT) [pFabric]

Use priorities (remaining flow size)

switches have small queues + send high/drop low priority

server's transmit/retransmit start at high, drop on transmit

Similar to job scheduling where "shortest job first" is a heuristic to approximate as remaining flow

IPV6

IPv4 address space exhaustion is primary motivation

w/ 128-bit rather than 32-bit (8 16-bit blocks), recursion

Routing (addr diff), Forwarding (addr), but least prefix matching (same)

Version	Traffic class	Flow Label	next header	hop limit
8	6	20	8	8
16	6	20	8	8
32	6	20	8	8
48	6	20	8	8
64	6	20	8	8
80	6	20	8	8
96	6	20	8	8
112	6	20	8	8
128	6	20	8	8

No fragmentation, checksums, optional state in extension headers (no options)

Addressing

64 bits addressable interface and 64 bits interface identifier

1st gives one subset, longest-prefix matching still applies

Special: host-local (localnet) ::1 ≡ 127.0.0.1

link-local (not routed): fe80::/10 ≡ 168.254.0.0/16

site-local (not routed globally)

global unicast

multi-cast

EUI-64 Identifier: **ab** → **abff:feab**, flip 7th bit of a

Neighbor Discovery Protocol

ICMPv6 messages for neighbor solicitation/advertisement, redirection to get Layer-2 address

Stateless Address Auto Configuration (SLAAC)

Network (top 64) obtained through Router Advertisements,

Interface Identifier (bottom 64) by EUI-64 or On-Link (Stable/Full)

Software-Defined Networking

Internet built as an artifact w/ complexity "network management" requires simplifying ideas to build an academic discipline

Define abstractions for the control plane

want isolation (VLANs), Access Control (ACL: header, out)

and traffic engineering (centralized computation)

Abstraction → interfaces → modularity

Control Plane Abstractions

① Be compatible with low-level hardware/software

want independence of proprietary HW/SW (forwarding)

OpenFlow (header, action) as standardized interface

switches accept ext control msg, standard flow entry format

② Make decisions based on entire network

Abstract (network state) as global network view

by annotated network graph through API (NETOS)

Runs on servers ("controllers") w/ info to form view and

entry to forwarding

Control by functional view, just a program → graph algo.

③ Compute configuration of each physical device

Abstract view of network to simplify configuration

Control prog → virtualization layer → Network OS

SDN separates data and control planes to reduce

complexity (simplifies interface, allows reusability)

Routing (graph, compute, give to SDN platform, reconfig)

Access Control (Control)

Control program → SDN platform

ACL