

Unit1 Introduction 协议栈

Network classification by distance

PAN (1m-10m, e.g., room), LAN (10m-1km, e.g., building/campus), MAN (1km-100km, e.g., town/country), WAN (100km-1000km, e.g., continent), Internet (10000km, e.g., planet)

What is Internet? 组成, 架构, 服务

Components: host/end systems + communication links + routers. Architecture: network of networks, loosely hierarchical, public Internet versus private intranet. Service: communication infrastructure, reliable/best effort data delivery.

Network architecture 架构

Client/server model: client host requests, receives service from always-on server. Peer-peer model: minimal or no use of dedicated servers, end systems interact and run programs that perform both client and server functions.

Access types 接入方式

Dial-up modem: telephony infrastructure, share physical line (surf or phone). Digital subscriber line (DSL): telephony infrastructure, dedicated physical line to center office. Cable modem: cable TV infrastructure, homes share access to router. Ethernet: end device + switch + router. Wireless: shared wireless access network, via base station/access point, e.g. Wi-Fi (802.11b/g), WiMAX (wireless interoperability for microwave access, IEEE 802.16), LTE (long term evolution)

Link types 链路分类

Guided media: signals propagate in solid media, e.g., copper (铜线), fiber (光纤/纤维), coax (同轴). Specific type (guided): twisted pair (双绞线), coaxial cable (同轴电缆), fiber optic cable (光纤). Unguided media: signals propagate freely, e.g., radio. Radio link types: terrestrial microwave, LAN (e.g., Wi-Fi), wide-area (e.g., cellular), satellite.

Network core 核心的架构, 实现方式

Architecture: mesh of interconnected routers. Approaches: circuit switching: dedicated circuit per call. packet switching: data sent through net in discrete "chunks" (Reality: pure①/pure②/mixture③+④)

Circuit switching

Feature: ①End-end resources (like bandwidth, switch capability) reserved for "calls" ②no sharing ③provide guaranteed service (a constant speed). Bandwidth divide: FDM (Frequency-Division Multiplexing), TDM (Time-Division Multiplexing)

Packet switching (datagram networks + virtual circuit networks)

Feature: ①Each end-end data stream divided into packets. ②Each packet uses full link bandwidth ③resource contention (no admission control, congestion, store and forward). Store and forward: A packet (size L) transmit through a link (bandwidth R) with 2 routers in the link, need 3L/R secs (store all the packet then push out at a router). Comments: ①great for burst data (resource sharing, no call setup) ②excessive congestion (packet delay/loss, need protocols for reliable transfer and congestion control). Virtual-Circuit Packet Switching (+Unit2): ①Data is transmitted as packets. ②Packets from one packet stream are sent along a pre-established path according to VC identifier, call setup for each call before data can flow and teardown afterwards ③Packets from different virtual circuits may be interleaved ④every router on source-destination path maintains "state" for each passing connection ⑤Guarantees in-sequence delivery of packets.

Delay

Nodal delay = processing + queuing + transmission (push out packet, Size/Bandwidth) + propagation Queuing delay: traffic intensity = (packet length * average packet arrival rate)

)/(bandwidth), this value -> 1, queuing delay becomes large. Other delay: ①purposefully delay (determined by protocol) ②packetization delay (in Voice over-IP (VoIP) applications)

Why layering?

①Explicit structure allows identification, relationship of complex system's pieces. ②modularization eases maintenance, updating of system. Lead to some problems: ③Functionality may be duplicated. ④One layer may need information present only in another layer.

Protocol stack

Internet protocol stack: ①Application (message): supporting network applications (e.g., FTP, SMTP, HTTP) ②transport (header[port]+QoS->segment): process-process data transfer (e.g., TCP, UDP) ③network (header[ip address]+QoS->datagram): routing of datagrams from source to destination (e.g., IP, routing protocols) ④link (header[physical address]+QoS+tail->frame): data transfer between neighboring network elements (e.g., PPP, Ethernet) ⑤Physical: bits "on the wire" OSI model: ①... ②presentation: allow applications to interpret meaning of data, e.g., encryption, compression, machine-specific conventions ③session: synchronization, checkpointing, recovery of data exchange ④... ⑤... ⑥... ⑦... ⑧... ⑨... ⑩...

About security

Denial of service (DoS): ①an attack against any system component that attempts to force that system component to limit, or even halt, normal services. ②only from one host or network node. Distributed denial of service (DDoS): ③more than one attack source. ④consume the resource of target host so that normal service cannot be provided.

Approach: (attacker) -> n (masters) -> n*m (slaves) -> 1 (target) [HW3]: Why hard to defend, IP spoofing: send packet with false source address

Unit2 IP Technology IPv4, IPv6

Key functions of network layer 网络层功能

①forwarding: move packets from router's input to appropriate router output (IP protocol) ②routing: determine route taken by packets from source to destination (routing algorithms)

Virtual circuit (VC network (+Unit1) Connection setup, forward, route Function: Provides network-layer connection service (analogous to TCP) Different to TCP: ①service: host-to-host according to IP address (TCP: port to port) ②no choice: network provides one or the other ③Implementation: in network core and end systems. VC: path from source to destination + VC number + entries in forwarding tables. VC number: can be changed on each link [HWS: not same VC number]. Signaling protocols: used to setup, maintain, teardown VC in ATM

Datagram network

Function: Provides network-layer connectionless service (analogous to UDP). Feature: ①no call setup at network layer ②no state about end-to-end connections at router ③no network-level concept of "connection" ④packets forwarded using destination host address ⑤packets between same source-destination pair may take different paths

Network layer protocol

IP (Internet Protocol), ARP (Address Resolution Protocol): IP address -> physical address ①ARP table (IP, MAC, TTL) ②TTL: times out, delete the mapping ③Broadcasts ARP query contains destination IP (MAC, FF-FF-FF-FF-FF-FF) ④destination replies MAC unicast. RARP (Reverse Address Resolution Protocol): physical address -> IP address.

ICMP (Internet Control Message Protocol): used to communicate network-level information (error reporting e.g., unreachable; echo request/reply e.g., ping), ICMP messages carried in the data portion of IP datagrams. IGMP (Internet Group Management Protocol): Host uses IGMP to announce participation in multicast (more see Unit9)

IPv4

Header length: 20B (32b*5) MTU (maximum transmission unit): largest possible link-level frame. MTU=header+data. Fragmentation: "reassembled" only at final destination IP header bits used to identify, order related fragments. e.g., 4000 byte datagram, MTU=1500 B ①1480B data area per fragment ②offset = 1480/8 for second fragment (0, 185, 370) ③the last fragment fragflag = 0, others = 1. Special IP: ①all0 (startup, source) ②all0+hid (a host of this network, source) ③all1 (local/limited broadcast, destination) ④nid+all1 (directed broadcast, destination) ⑤nid+all0 (network itself/directed broadcast, destination) ⑥127+notall0/1 (for loopback, source/destination) Classful address schema: ①Adventurer: A router can keep one routing entry per network instead of per destination host. Disadvantage: Requiring a unique prefix for each physical network would exhaust the address space quickly as the Internet proliferates. Solution: unnumbered point-to-point links, proxy ARP, and subnet advertisement. Subnet: ①Qdevice interfaces with same subnet part of IP address. ②can physically reach each other without intervening router Classless Inter-Domain Routing (CIDR) was invented to use address space more efficiently. Notation: a.b.c.d/x. Longest Match: subnet mask AND des IP -> network id, choose the longest match. ICANN: Internet Corporation for Assigned Names and Numbers, only largest ISPs need to contact. Address classification: A: 1.0.0.0-127.0.0.0, B: 128.0.0.0-191.255.0.0, C: 192.0.0.0-223.255.0.0, D: 224.0.0.0-239.255.255.0 (for multicast), reserved: 240.0.0.0-239.255.255.255. Private Address: 10.0.0.0/8, 172.16.0.0/12, 192.168.0.0/16, 169.254.0.0/16.

NAT (Network Address Translation) 16b5b5

Motivation: local network uses just one IP address as far as outside world is concerned. Benefit: ①simple gateway between Internet and private network ②simple security due to stateful filter implementation ③privacy and topology hiding. Argument: Routers should only process up to layer 3 (but NAT provide services of transport layer) ④lead to NAT traversal problem -> Solution: statically configure, universal plug and play (UPnP) Internet gateway device (IGD) protocol and relaying (through another site)

IPv6

Motivation: ③2-bit address space soon to be completely allocated (Approach you slow the consumption rate: Dial-access/PPP/DHCP, Strict allocation policies, CIDR, NAT) ②helps speed processing/forwarding ③to facilitate QoS/resource allocation. Difference: [HW6: benefits] ①Options indicated by "Next Header" field. [tip: each header 4 Bytes, order: hop-by-hop, Routing, Fragment, Authentication, Encryption, Destination, only hop-by-hold is processed at a hop, for routing header: A send to D through B, C -> src A, des B, routing header: C, D]; ②Header Checksum eliminated to reduce processing time at each hop; ③Fragmentation: (a)move to extension header (b)fragmentation is end-to-end function, no fragmentation occurs in intermediate routers (c)use guaranteed minimum MTU of 1280 octets (8bits) [tip: MTU is 1280 for ipv6, 68 for ipv4, if MTU < 1280, use link-specific fragmentation at end device] or perform Path MTU Discovery [tip: send a specific packet, ICMP "packet too big" message would occur if packet is too big] to identify the minimum MTU along the path to the destination. [Problem: routers cannot be changed as easily as those in IPv4 because a change in a route can also change the path MTU] ④TTL -> Hop Limit ⑤Protocol -> Next Header. ⑥Service Type -> Traffic Class ⑦Addresses increased 32 bits -> 128 bits ⑧Flow Label added, identify datagrams in same "flow" (e.g., two applications that need to send video can establish a flow on which QoS is guaranteed). [HW: Private addresses]. Transition [HW2]: IPv6 deployments will occur piecemeal from the edge. Co-Existence (3种) Techniques: ①Dual-stack techniques: to allow IPv4 and IPv6 to co-exist in the same devices and networks ②Tunneling techniques: IPv6 carried as payload in IPv4 datagram among IPv4 routers ③translation techniques: to allow IPv6-only devices to communicate with IPv4-only devices

Unit3 Network Switching

Switching

Switch type: ①Memory ②bus ③cross bus. Output port queueing: buffering when arrival rate via switch exceeds output line speed. Input Port Queueing: fabric slower than input ports combined [HW8: queue]

Devices [HW2: comparison]

Hubs: no buffer, no filtering, no redirection, no CSMA/CD Switch: store forward, no collisions, full duplex, network is restricted to a spanning tree in order to prevent the cycling of broadcast storm, maintain switch tables, implement filtering, learning algorithms Bridge: only has one incoming and one outgoing port, perform in software (switch hardware) Router: provide firewall protection and allow the network to be built with a rich topology, maintain routing tables, implement routing algorithms.

[HW11: IP over ATM/IP over SDH/IP over WDM]

[HW9&10: MPLS]

Unit4 Transport Layer TCP, UDP, GBN, SR

Transport-layer protocols

TCP (transmission control protocol): Header length: 20B, reliable, in-order delivery, congestion control, flow control, connection-oriented, integrity checking. UDP (user datagram protocol): Unreliable, unordered delivery, error checking. ①Header length: 8B ②used for streaming multimedia apps (loss tolerant, rate sensitive) ③DNS, SNMP

TCP & UDP Segment structure

① Pipelining protocols [HW16: Compare GBN/SR/TCP] ② Stop and wait ARQ

GBN (go-back-N/sliding window protocol): receiver only send cumulative ACKs, drop unexpected packets; sender sets timer for oldest unACKed packet, and will retransmit a series packets if a former packet is lost SR (selective repeat): receiver buffers and ACK each packets, sender sets timer for each unACKed packet, only retransmit packets which are in error. (Requirement: window size <= half of seq # size, if not, can't distinguish new packet and retransmission)

TCP [HW16: Compare GBN/SR/TCP]

Seq & ACK (Telnets): initial number (given or random) for client and server (seq1, seq2) Seq for expect segment seq (rcvseq+rcvdata size). [e.g., client seq=42, ACK=79, data=C, server seq=79, ACK=43, data=C, client seq=43, ACK=80], two C for echo reply. Fast retransmit: if sender receives 3 client seq=43s, resend segment before timer expires. Three-way handshake: ①client host sends TCP SYN segment to server (SYN=seq=C, isn) ②server host receives SYN, replies with SYNACK segment (SYN=1, seq=C+1, isn=ACK=C+1) ③client receives SYNACK, replies with ACK segment (SYN=0, seq=C+1, ACK=C+1) ④Close connection: ①client send FIN ②server ACK ③server send FIN ④client ACK. Flow control: receiver send its rcv window size in the segment back to sender.

Unit5 Congestion Control [HW17: Compare flow/congestion]

Approach: ①Qend-end: end-system observed loss, delay; approach taken by TCP ②Network-assisted: routers provide feedback to end systems (e.g. IBM SNA, DECbit, TCP/IP ECN and ATM)

TCP congestion control

Feature: AIMD (additive increase, multiplicative decrease): increase cwnd by 1 MSS every RTT until loss detected (CA mode), cut cwnd in half after loss [tip: sending rts=cwnd/RTT]. TCP congestion control: ①When cwnd is below Threshold, sender in slow-start (double cwnd every RTT phase, window grows exponentially. ②When cwnd is above Threshold, sender is in CA (congestion-avoidance) (cwnd+= every RTT phase, rwnd set to Threshold. ③When a triple duplicate ACK occurs, Threshold set to cwnd/2 and cwnd set to 1 MSS. Tail-drop policy cause global synchronization [HW20], reason: under a tail-drop policy, the router will discard one segment from N connections rather than N segments from one connection, the simultaneous loss causes all N instances of TCP to enter slow-start at the same time and throughput decreases suddenly, after the network recovers, throughput will suddenly increase a lot. Solution -> RED (Random Early detection): instead of waiting until the queue overflows, a router slowly and randomly drops datagrams as congestion increases. Throughput = W/RTT * (1+1/2)/2 = 0.75 W/RTT (W is window size when loss occurs) Fairness: for idealized conditions (same MSS and RTT), TCP is fair. In practice, those sessions with smaller RTT are able to grab the available bandwidth at that link more quickly as the link becomes free. Moreover, consider UDP and parallel TCP connections. Explicit feedback mechanisms: selective acknowledgment (SACK), explicit congestion notification (ECN)

Unit6 Multimedia QoS 区分服务模式

Multimedia applications (delay sensitive, loss tolerant)

Stored streaming: e.g., YouTube, media stored at source, transmitted to client, client playback begins before data has arrived. Constraint: in time for playback ②live streaming: e.g., IPTV, can't fast forward ③real time interactive: e.g., IP telephony. Approach: ①Integrated services philosophy: fundamental changes in Internet ②Differentiated services philosophy: fewer changes to Internet infrastructure ③Loose-guarantee: no major changes, provide more bandwidth when needed (e.g., CDN) [HW23], application-layer multicast [tip: Hard guarantee: receive QoS with certainty. Soft guarantee: with high probability]

Supporting Multimedia applications

①Approach, Unit of allocation, Guarantee, Mechanisms ②making the best of best-of service, none, none or soft, application layer support, CDN, over-provisioning ③Differentiated QoS, classes of flows, none or soft, policing, scheduling ④guaranteed QoS, individual flows, soft or hard, once a flow is admitted, policing, scheduling, call admission and signaling.

Streaming Multimedia: UDP or TCP?

UDP: sends at rate appropriate, send rate=encoding rate=constant rate, fill rate=constant rate-packet loss. TCP: send at maximum possible rate, fill rate fluctuates due to TCP congestion control, HTTP/TCP passes more easily through firewalls. [Tip: Handle different client receive rate capabilities: server stores, transmits multiple copies of video, encoded at different rates]

Principles for QoS Guarantees

①packet classification ②isolation: Scheduling and policing ③high resource utilization ④call admission

Scheduling Mechanisms

①FIFO ②Priority ③Round robin scheduling / fair queuing: cyclically scan class queues ④Weighted Fair Queuing (WFQ): each class gets weighted amount of service in each cycle

Policing Mechanisms [HW25]

Token Bucket: bucket can hold b tokens, tokens generated at rate r token/sec unless bucket full, the max number of packets to send in a given time T is (r*T+b), if scheduling is WFQ, the max delay is (b/r+WQ/sum(wi)).

Differentiated services [HW21]

Edge router: per-flow traffic management, packet classification and traffic conditioning, marks packets as in-profile and out-of-profile. Core router: per class traffic management, buffering and scheduling based on marking at edge, forwarding, preference given to in-profile packets. [tip: 1 packet is marked in the 8-bit (6 bits used for Differentiated Service Code Point (DSCP) determine Per-hop Behavior (PHB), 2 bits not used) Type of service (TOS) in IPv4, and 8-bit Traffic Class in IPv6] [PHB result in a different observable (measurable) forwarding performance behavior, PHB does not specify what mechanisms to use to ensure required PHB performance behavior. Two type: Expedited Forwarding (Premium): pkt departure rate of a class equals or exceeds specified rate. Assured Forwarding: define 4 classes of traffic]

Integrated Services [HW21]

RSVP [HW24]

Unit8 Internal Routing 路由算法 LS, DV, RIP, OSPF

Routing algorithms classification

①Global: all routers have complete topology, link cost information, e.g., LS. Decentralized: router knows physically-connected neighbors, link costs to neighbors, iterative process of computation, exchange of information with neighbors, e.g., DV. ②Static: change slowly, only for simplest cases. Dynamic: periodic update in response to topology or link cost changes, necessary in large internets. ③Load-sensitive or load-insensitive (today's algorithms)

LS (Link state/Dijkstra) [HW26]

Complexity: with n nodes, E links, O(nE) msgs sent, n(n-1)/2 comparisons: O(n^2), possible. O(nlogn). Oscillations possible [HW28]. Robustness (鲁棒性): node can advertise incorrect link cost, each node computes only its own table, somewhat separated route calculations providing a degree of robustness

DV (Distance vector) [HW26]

May be routing loops, "count to infinity" problem [HW29]. Solution: Poisoned reverse: [HW30]: to avoid routing loops, when a router find a subnet is not alive, it will set the cost infinite (e.g., 16) when broadcast to other routers instead of deleting it immediately.

RIP (Routing Information Protocol) [HW27: Compare] -> DV

Basic idea: ①Distance vectors: exchanged among neighbors every 30 sec via Response Message ②QoS advertisement: list of up to 25 destination subnets within AS. Timer: 30s for routing-update, 180s for time out, 120s for garbage collection (delete route) Disadvantages: [HW31]

OSPF (Open Shortest Path First) [HW27: Compare] -> LS

Scale: 150-500 routers/Area. Basic idea: ①Distributed replicated database model

(Each router builds a topology database describes complete routing topology) ②Link state database (identical for all the routers) ③LSA (Information about adjacencies sent to all routers) ④A "shortest path" algorithm is used to find best route (dijkstra) (Converge as quickly as databases can be updated, every router calculate itself routing table independently) Two-level hierarchy: local area, backbone. Link state advertisement (LSA) is bounded by area. Advantage: security, load balancing, type of service (TOS) routing, integrated unicast and multicast support, hierarchical.

Unit7 External Routing 区域间选路 BGP, 熟土豆, 冷土豆

EGP (interior gateway protocol): OSPF, IS-IS, RIP, EIGRP (cisco). EGP (exterior gateway protocol): BGP. IXP (Internet exchange point): a physical infrastructure through which ISPs and CDN's exchange Internet traffic between their networks. ASN (autonomous system number): 自治系统号

Why do we need EGP? [HW33: IntraAS/InterAS routing]

①Scalability (hierarchy, limit scope of failure) ②Flexibility in choosing routes ③Define administrative boundary ④Policy (control reachability to prefixes) ⑤Interconnections type

Transit ①Peering ②If a peer with B, B peer with C, A's customer could not send data to C directly. [HW34]

BGP (Border Gateway Protocol) v4 -> BGP?

Layer: use reliable transport i.e., TCP. Function: ①Obtain subnet reachability information from neighboring ASs. ②Propagate reachability information to all AS-internal routers. ③Determine "good" routes to subnets based on reachability information and policy. Neighbor relationships: ①EBGP session spans two ASs, to share connectivity information across AS Network Layer Reachability Information (NLRI). ②IBGP session between routers in the same AS, carrying information within an AS. Handle prefix: because of longest match principle, if AS has 3 subnets, 138.16.64/24, 138.16.65/24 and 138.16.66/24, it will aggregate prefixes to 138.16.64/22 in its advertisement, because another AS will advertise 138.16.67/24. AS-PATH contains ASs through which prefix advertisement has passed, e.g., AS2 receive prefix from AS1, when AS2 advertise to AS3, AS-PATH=AS2 AS1. To prevent loop, AS will never accept a route containing AS itself. Next-hop: the router interface that begins the AS-PATH and indicates specific internal AS router to next-hop AS. Every time a IP address of the border router (when the announcement is in an AS, the Next-Hop not change). Route selection: ①local preference value attribute: policy decision up to AS's network administrator (Highest values are selected) ②Shortest AS-PATH (Distance vector algorithm; Distance metric: # AS hops, NOT # router hops) ③Closest Next-Hop router/hot potato routing (Least-cost path determined by intra-AS algorithm) ④Additional criteria (can be more complicated) Hot potato routing & Cold potato routing [HW35&36]

Unit8 Virtual Private Networks (VPN)

Goal: to keep internal datagrams private while still allowing external communication Main benefit: reduce cost. Other benefit: Scalability NAS (Network Access Servers): a device that interfaces between an access network and a packet-switched network, serve as a tunnel endpoint in a remote access VPN.

Types

①Site-to-site: allow connectivity between an organization's (or organizations') geographically dispersed sites (such as a head office and branch offices). Two types of site-to-site: Intranet: Allow connectivity between sites of a single organization. Extranet: Allow connectivity between organizations such as business partners or a business and its customers. ②Remote access: allow mobile or home-based users to access an organization's resources remotely.

Protocols for site-to-site [HW37-41]

①IPsec (IP security) ②GRE (Generic Routing Encapsulation) ③L2TPv3 (Layer Two Tunneling Protocol version 3) ④Q-in-Q (IEEE 802.1Q tunneling) ⑤MPLS (multiprotocol label switching) LSPs

Protocols for remote access [HW37-41]

①L2F (The Layer Two Forwarding) Protocol ②PPTP (The Point-to-Point Tunneling Protocol) ③L2TPv2/L2TPv3 ④IPsec ⑤SSL (Secure Sockets Layer)

Most popular: PPTP, L2TP and IPsec

Protocol Classified by layer

①Application to Application: SSL ②End to End: IPsec Transport Mode ③Gateway to Gateway: PPTP, L2TP/IPsec, IPsec Tunnel Mode ④Client to Gateway: L2TP/IPsec

VPN Critical Functions

①Authentication: validates that the data was sent from the sender. ②Access control: limiting unauthorized users from accessing the network. ③Confidentiality: preventing the data to be read or copied as the data is being transported. ④Data integrity: ensuring that the data has not been altered AS: Before sending IPsec datagrams from source entity to destination entity, the source and destination entities create a network-layer logical connection - called a security association

Unit9 Multicasting 多播, 生成树

In-network duplication

①Uncontrolled flooding: when node receives broadcast packet, sends copy to all neighbors (except the source neighbor) Problem: cycles, broadcast storms. ②Controlled flooding: node only broadcast packet if it hasn't broadcast same packet before. Approach: sequence-number-controlled / reverse path forwarding (RPF) / reverse path broadcast (RPB) ③Spanning tree: no redundant packets received by any node [tip: a node need not be aware of the entire tree, simply needs to know its spanning-tree neighbors.]

Multicast: one sender to many receivers

Control scope: ①IP's TTL (Time-To-Live) field ②Administrative scoping. Local -> IGMP: ③When it joins a group, host sends message (REPORT) declaring membership. ④Multicast router periodically polls (QUERY) a host to determine if any host on the network is still a member of a group Wide area (among routers) -> multicast trees: local router interacts with other routers to receive multicast datagram flow

Approaches for building multicast trees

①Source-based tree: one tree per source, e.g., shortest path trees [MOSPF (Multicast extensions to OSPF), reverse path forwarding [DVMRP (Distance-Vector Multicast Routing Protocol), PIM-DM (Protocol Independent Multicast-Dense Mode)] ②Group-shared tree: group uses one tree, e.g., minimal spanning [Steiner], center-based trees [CBT (Core-Based Trees), PIM-SM (Protocol Independent Multicast-Sparse Mode)]

Homework 1. 35users, active 10% of time, probability > 10 active at same time P = 1- sum n. from 0 to 10 (C(n,35) * 0.9^n * (35-n) * 0.1^n)

Complete hubs, switches, bridges and routers