

Review of Internet Architecture and Protocols

Professor Guevara Noubir
 Northeastern University
 noubir@ccs.neu.edu

Reference Textbooks:
 Computer Networks: A Systems Approach, L. Peterson, B. Davie, Morgan Kaufmann

1

Outline

- Internet Protocol
- Addressing
- IP over LAN
- Routing
- End-to-End protocols
- Naming

Network Security Internet Architecture and Protocols

2

IP Internet

Concatenation of Networks

Protocol Stack

Network Security Internet Architecture and Protocols

3

Internet Control Message Protocol (ICMP) RFC 792

- n Integral part of IP but runs as ProtocolType = 1 using an IP packet
- n Codes:
 - n Echo (ping)
 - n Redirect (from router to inform source host of better route)
 - n Destination unreachable (protocol, port, or host)
 - n TTL exceeded (so datagrams don't cycle forever)
 - n Checksum failed
 - n Reassembly failed

Network Security Internet Architecture and Protocols 7

Global Addresses

- n Properties
 - n globally unique
 - n hierarchical: network + host
- n Dot Notation
 - n 10.3.2.4
 - n 128.96.33.81
 - n 192.12.69.77

A:	0	Network	Host
		14	24

B:	1	0	Network	Host
			14	16

C:	1	1	0	Network	Host
				21	8

D:	1	1	1	0	Group
					28

Multicast

Network Security Internet Architecture and Protocols 8

Datagram Forwarding

- n Strategy
 - n every datagram contains destination's address
 - n if directly connected to destination network, then forward to host
 - n if not directly connected to destination network, then forward to some router
 - n forwarding table maps network number into next hop
 - n each host has a default router
 - n each router maintains a forwarding table
- n Example (R2)

Network Number	Next Hop
1	R3
2	R1
3	interface 1
4	interface 0

Network Security Internet Architecture and Protocols 9

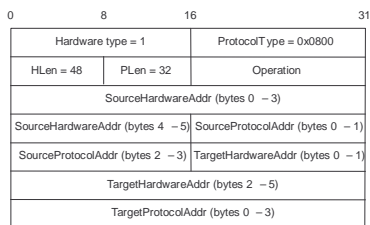
Address Translation

- n Map IP addresses into physical addresses
 - n destination host
 - n next hop router
- n Techniques
 - n encode physical address in host part of IP address
 - n not reasonable
 - n table-based
- n ARP
 - n table of IP to physical address bindings
 - n broadcast request if IP address not in table
 - n target machine responds with its physical address
 - n table entries are discarded if not refreshed

ARP Details

- n Request Format
 - n HardwareType: type of physical network (e.g., Ethernet)
 - n ProtocolType: type of higher layer protocol (e.g., IP)
 - n HLEN & PLEN: length of physical and protocol addresses
 - n Operation: request or response
 - n Source/Target-Physical/Protocol addresses
- n Notes
 - n table entries timeout in about 15 minutes
 - n update table with source when you are the target
 - n update table if already have an entry
 - n do not refresh table entries upon reference

ARP Packet Format

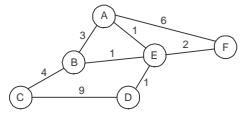


Dynamic Host Configuration Protocol (DHCP)

- IP addresses of interfaces cannot be configured when manufactured (like for Ethernet)
- Configuration is an error-prone process
- Solution: centralize the configuration information in a DHCP server:
 - DHCP server discovery: broadcast a DHCPDISCOVER request
 - Requests are relayed (unicast) to the server by DHCP relays
 - DHCP server broadcast replies with <HWADDR, IPADDR, lease-info>

Routing Overview

- Forwarding vs Routing
 - forwarding: to select an output port based on destination address and routing table
 - routing: process by which routing table is built
- Network as a Graph



- Problem: Find lowest cost path between two nodes
- Factors
 - relatively static: topology
 - dynamic: load

Distance Vector

- Each node maintains a set of triples
 - (Destination, Cost, NextHop)
- Exchange updates directly connected neighbors
 - periodically (on the order of several seconds)
 - whenever table changes (called *triggered* update)
- Each update is a list of pairs:
 - (Destination, Cost)
- Update local table if receive a "better" route
 - smaller cost
 - came from next-hop
- Refresh existing routes; delete if they time out

Example

Table for node B

Destination	Cost	NextHop
A	1	A
C	1	C
D	2	C
E	2	A
F	2	A
G	3	A

Network Security Internet Architecture and Protocols 16

Routing Information Protocol (RIP)

- Uses Bellman-Ford's algorithm
- Protocol over UDP, port 520
- Distance-vector protocol
- Protocol overview:
 - Init: send a request packet over all interfaces
 - On response reception: update the routing table
 - On request reception:
 - if request for complete table (*address family=0*) send the complete table
 - else send reply for the specified address (*infinity=16*)
 - Regular routing updates:
 - every 30 seconds part/entire routing table is sent (broadcast) to neighboring routers
 - Triggered updates: on metric change for a route
 - Simple authentication scheme

Network Security Internet Architecture and Protocols 17

Link State

- Strategy
 - send to all nodes (not just neighbors) information about directly connected links (not entire routing table)
- Link State Packet (LSP)
 - id of the node that created the LSP
 - cost of link to each directly connected neighbor
 - sequence number (SEQNO)
 - time-to-live (TTL) for this packet

Network Security Internet Architecture and Protocols 18



Link State (cont)

- n Reliable flooding
 - n store most recent LSP from each node
 - n forward LSP to all nodes but one that sent it
 - n do not forward already received LSPs
 - n generate new LSP periodically
 - n increment SEQNO
 - n start SEQNO at 0 when reboot
 - n decrement TTL of each stored LSP
 - n discard when TTL=0



Route Calculation

- n Dijkstra's shortest path algorithm
 - n Let
 - n N denotes set of nodes in the graph
 - n $I(i, j)$ denotes non-negative cost (weight) for edge (i, j)
 - n s denotes this node
 - n M denotes the set of nodes incorporated so far
 - n $C(n)$ denotes cost of the path from s to node n
- ```

M = {s}
for each n in N - {s}
 C(n) = I(s, n)
while (N != M)
 M = M union {w} such that C(w) is the minimum for
 all w in (N - M)
 for each n in (N - M)
 C(n) = MIN(C(n), C(w) + I(w, n))

```

---

---

---

---

---

---

---

---



## Open Shortest Path First

- n IP protocol (not over UDP), reliable (sequence numbers, acks)
- n Protocol overview: link state protocol
  - n The link status (cost) is sent/forwarded to all routers (LSP)
  - n Each router knows the exact topology of the network
  - n Each router can compute a route to any address
  - n simple authentication scheme
- n Advantages over RIP
  - n Faster to converge
  - n The router can compute multiple routes (e.g., depending on the type of services, load balancing)
  - n Use of multicasting instead of broadcasting (concentrate on OSPF routers)

---

---

---

---

---

---

---

---

## Popular Interior Gateway Protocols

- n RIP: Route Information Protocol
  - n distributed with Unix
  - n distance-vector algorithm
  - n based on hop-count
- n OSPF: Open Shortest Path First
  - n more recent Internet standard
  - n uses link-state algorithm
  - n supports load balancing
  - n supports basic integrity check

---

---

---

---

---

---

---

---

## How to Make Routing Scale

- n Flat versus Hierarchical Addresses
- n Inefficient use of Hierarchical Address Space
  - n class C with 2 hosts ( $2/256 = 0.78\%$  efficient)
  - n class B with 255 hosts ( $255/65536 = 0.39\%$  efficient)
- n Still Too Many Networks
  - n routing tables do not scale
  - n route propagation protocols do not scale

---

---

---

---

---

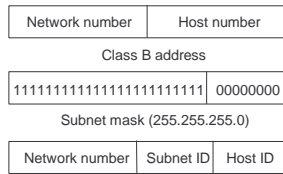
---

---

---

## Subnetting

- n Add another level to address/routing hierarchy: *subnet*
- n *Subnet masks* define variable partition of host part
- n Subnets visible only within site




---

---

---

---

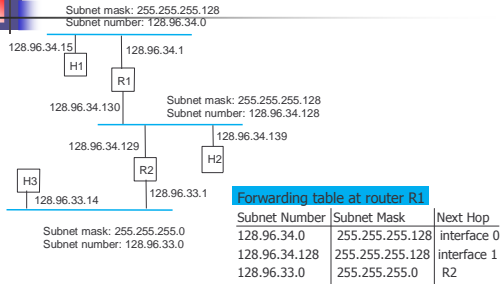
---

---

---

---

## Subnet Example



---

---

---

---

---

---

---

---

## Forwarding Algorithm

```
D = destination IP address
for each entry (SubnetNum, SubnetMask, NextHop)
 D1 = SubnetMask & D
 if D1 = SubnetNum
 if NextHop is an interface
 deliver datagram directly to D
 else
 deliver datagram to NextHop
```

- n Use a default router if nothing matches
- n Not necessary for all 1s in subnet mask to be contiguous
- n Can put multiple subnets on one physical network
- n Subnets not visible from the rest of the Internet

---

---

---

---

---

---

---

---

## Supernetting: Restructuring IP Addresses

- n Assign block of contiguous network numbers to nearby networks
- n Called CIDR: Classless Inter-Domain Routing
- n Represent blocks with a single pair (first\_network\_address, count)
- n Restrict block sizes to powers of 2
  - n E.g., 192.4.16 – 192.4.31: /20
- n Use a bit mask (CIDR mask) to identify block size
- n All routers must understand CIDR addressing

---

---

---

---

---

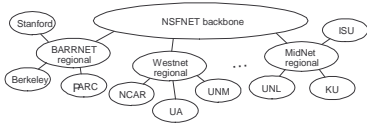
---

---

---

## Internet Structure

### Past



---

---

---

---

---

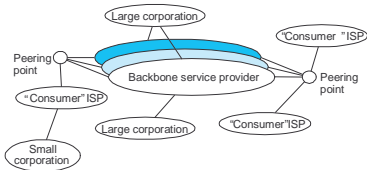
---

---

---

## Internet Structure

### Yesterday



---

---

---

---

---

---

---

---

## Route Propagation

- n Know a smarter router
  - n hosts know local router
  - n local routers know site routers
  - n site routers know core router
  - n core routers know almost everything
- n Autonomous System (AS)
  - n corresponds to an administrative domain
  - n examples: University, company, backbone network
  - n assign each AS a 16-bit number
- n Two-level route propagation hierarchy
  - n interior gateway protocol (each AS selects its own)
  - n exterior gateway protocol (Internet-wide standard)

---

---

---

---

---

---

---

---

## EGP: Exterior Gateway Protocol

- Overview
  - designed for tree-structured Internet
  - concerned with *reachability*, not optimal routes
- Protocol messages
  - neighbor acquisition: one router requests that another be its peer; peers exchange reachability information
  - neighbor reachability: one router periodically tests if the another is still reachable; exchange HELLO/ACK messages; uses a k-out-of-n rule
  - routing updates: peers periodically exchange their routing tables (distance-vector)

---

---

---

---

---

---

---

---

## BGP-4: Border Gateway Protocol

- AS Types
  - stub AS: has a single connection to one other AS
    - carries local traffic only
  - multihomed AS: has connections to more than one AS
    - refuses to carry transit traffic
  - transit AS: has connections to more than one AS
    - carries both transit and local traffic
- Each AS has:
  - one or more border routers
  - one BGP *speaker* that advertises:
    - local networks
    - other reachable networks (transit AS only)
    - gives *path* information

---

---

---

---

---

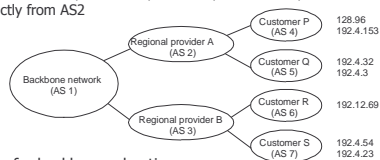
---

---

---

## BGP Example

- Speaker for AS2 advertises reachability to P and Q
  - network 128.96, 192.4.153, 192.4.32, and 192.4.3, can be reached directly from AS2



- Speaker for backbone advertises
  - networks 128.96, 192.4.153, 192.4.32, and 192.4.3 can be reached along the path (AS1, AS2).
- Speaker can cancel previously advertised paths

---

---

---

---


---

---

---

---

# END TO END PROTOCOLS



34

---

---

---

---

---

---

---

---

# End-to-End Protocols

- n Goal: turn host-to-host packet delivery into **process-to-process** communication channel
- n Underlying best-effort network
  - n drop messages
  - n re-orders messages
  - n limits messages to some finite size
  - n delivers messages after an arbitrarily long delay
  - n delivers duplicate copies of a given message
- n Common end-to-end services
  - n guarantee message delivery
  - n deliver messages in the same order they are sent
  - n deliver at most one copy of each message
  - n support arbitrarily large messages
  - n support synchronization
  - n allow the receiver to flow control the sender
  - n support multiple application processes on each host

Network Security      Internet Architecture and Protocols      35

---

---

---

---

---

---

---

---

# Types of End-to-End Protocols

- n Simple asynchronous demultiplexing service (e.g., UDP)
- n Reliable byte-stream service (e.g., TCP)
- n Request reply service (e.g., RPC)

Network Security      Internet Architecture and Protocols      36

---

---

---

---

---

---

---

---

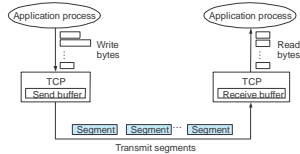
## Simple Demultiplexor (UDP)

- Unreliable and unordered datagram service
  - Adds multiplexing
  - No flow control
  - Endpoints identified by ports
    - servers have *well-known* ports (e.g., DNS: port 53, talk: 517)
    - see `/etc/services` on Unix
  - Header format
- ```

0          16          31
+-----+-----+
| SrcPort | DestPort |
+-----+-----+
| Checksum | Length  |
+-----+-----+
| Data     |
+-----+
            
```
- Optional checksum
 - pseudo header + UDP header + data
 - Pseudo header = protocol number, source IP addr, dest IP addr, UDP length

TCP Overview

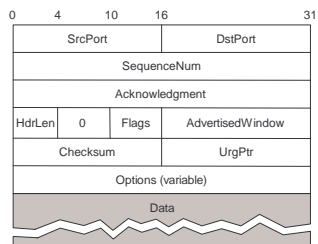
- Reliable
- Connection-oriented
- Byte-stream
 - app writes bytes
 - TCP sends *segments*
 - app reads bytes
- Full duplex
- Flow control: keep sender from overrunning receiver
- Congestion control: keep sender from overrunning network



Data Link Versus Transport

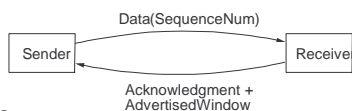
- Potentially connects many different hosts
 - need explicit connection establishment and termination
- Potentially varying RTT
 - need adaptive timeout mechanism
- Potentially long delay in network
 - need to be prepared for arrival of very old packets
- Potentially varying capacity at destination
 - need to accommodate different node capacity
- Potentially varying network capacity
 - need to be prepared for network congestion

Segment Format



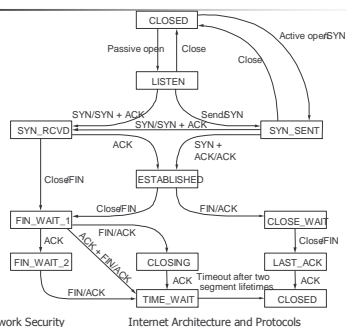
Segment Format (cont)

- n Each connection identified with 4-tuple:
 - n (SrcPort, SrcIPAddr, DsrPort, DstIPAddr)
- n Sliding window + flow control
 - n acknowledgment, SequenceNum, AdvertisedWindow



- n Flags
 - n SYN, FIN, RESET, PUSH, URG, ACK
- n Checksum
 - n pseudo header + TCP header + data

State Transition Diagram





Sliding Window in TCP

- n Purpose:
 - n Guarantees a reliable delivery of data (ARQ)
 - n Ensures that data is delivered in order (SeqNum)
 - n Enforces flow-control between sender and receiver (AdvertisedWindow field)



NAMING

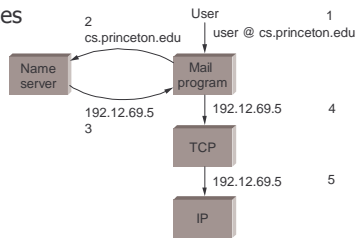


Naming in the Internet

- n Hosts
 - `cheltenham.cs.princeton.edu` → `192.12.69.17`
 - `192.12.69.17` → `80:23:A8:33:5B:9F`
- n Files
 - `/usr/llp/tmp/foo` → `(server, fileid)`
- n Users
 - `Larry Peterson` → `llp@cs.princeton.edu`

Examples (cont)

Mailboxes

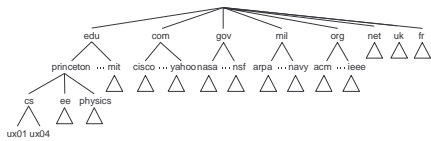


Services

nearby ps printer with short queue and 2MB

Domain Naming System

Hierarchy

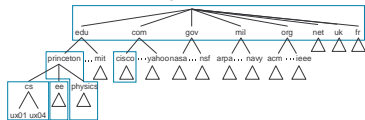


Name

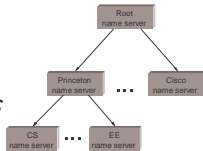
chinstrap.cs.princeton.edu

Name Servers

Partition hierarchy into zones



Each zone implemented by two or more *name servers*



Resource Records

- Each name server maintains a collection of *resource records* (Name, Value, Type, Class, TTL)
- Name/Value: not necessarily host names to IP addresses
- Type
 - A: Value is an IP address
 - NS: Value gives domain name for host running name server that knows how to resolve names within specified domain.
 - CNAME: Value gives canonical name for particle host; used to define aliases.
 - MX: Value gives domain name for host running mail server that accepts messages for specified domain.
- Class: allow other entities to define types
 - IN: Means Internet
- TTL: how long the resource record is valid

Root Server

```
(princeton.edu, cit.princeton.edu, NS, IN)
(cit.princeton.edu, 128.196.128.233, A, IN)

(cisco.com, thumper.cisco.com, NS, IN)
(thumper.cisco.com, 128.96.32.20, A, IN)

...
```

Princeton Server

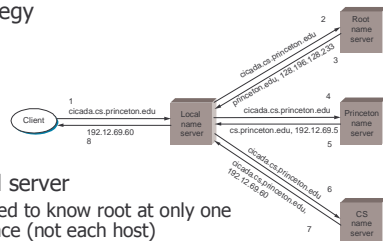
```
(cs.princeton.edu, optima.cs.princeton.edu, NS, IN)
(optima.cs.princeton.edu, 192.12.69.5, A, IN)
(ee.princeton.edu, helios.ee.princeton.edu, NS, IN)
(helios.ee.princeton.edu, 128.196.28.166, A, IN)
(jupiter.physics.princeton.edu, 128.196.4.1, A, IN)
(saturn.physics.princeton.edu, 128.196.4.2, A, IN)
(mars.physics.princeton.edu, 128.196.4.3, A, IN)
(venus.physics.princeton.edu, 128.196.4.4, A, IN)
```

CS Server

```
(cs.princeton.edu, optima.cs.princeton.edu, MX, IN)
(cheltenham.cs.princeton.edu, 192.12.69.60, A, IN)
(che.cs.princeton.edu, cheltenham.cs.princeton.edu,
 CNAME, IN)
(optima.cs.princeton.edu, 192.12.69.5, A, IN)
(opt.cs.princeton.edu, optima.cs.princeton.edu,
 CNAME, IN)
(baskerville.cs.princeton.edu, 192.12.69.35, A, IN)
(bas.cs.princeton.edu, baskerville.cs.princeton.edu,
 CNAME, IN)
```

Name Resolution

Strategy



Local server

- need to know root at only one place (not each host)
- site-wide cache

Summary

Multi-layer stack of protocols:

- Link Layer: ethernet (IEEE802.3), FDDI, ATM, wlan (IEEE802.11)
- Network Layer:
 - Internet Protocol (IP) is a focal point
 - Routing protocols: RIP, OSPF, BGP-4
- Transport Layer: UDP, TCP
- Naming: DNS
