



15-441 15-641 Computer Networking

Lecture 6 – The Internet Protocol Peter Steenkiste

Fall 2016

www.cs.cmu.edu/~prs/15-441-F16

Outline



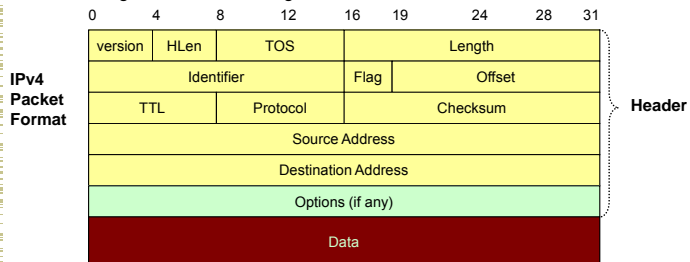
- The IP protocol
 - IPv4
 - IPv6
- IP in practice
 - Network address translation
 - Address resolution protocol
 - Tunnels

2

IP Service Model

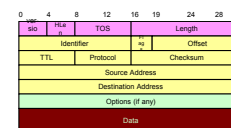


- Low-level communication model provided by Internet
- Datagram
 - Each packet self-contained
 - All information needed to get to destination
 - No advance setup or connection maintenance
 - Analogous to letter or telegram



3

IPv4 Header Fields

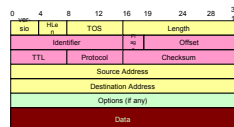


- Version: IP Version
 - 4 for IPv4
- HLen: Header Length
 - 32-bit words (typically 5)
- TOS: Type of Service
 - Priority information
- Length: Packet Length
 - Bytes (including header)
- Header format can change with versions
 - First byte identifies version
- Length field limits packets to 65,535 bytes
 - In practice, break into much smaller packets for network performance considerations

4

IPv4 Header Fields

- Identifier, flags, fragment offset → used for fragmentation
- Time to live
 - Must be decremented at each router
 - Packets with TTL=0 are thrown away
 - Ensure packets exit the network
- Protocol
 - Demultiplexing to higher layer protocols
 - TCP = 6, ICMP = 1, UDP = 17...
- Header checksum
 - Ensures some degree of header integrity
 - Relatively weak – 16 bit
- Source and destination IP addresses
- Options
 - E.g. Source routing, record route, etc.
 - Performance issues
 - Poorly supported



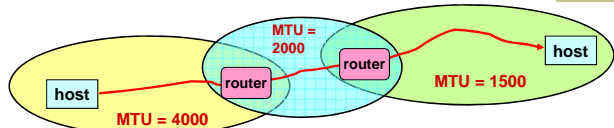
5

IP Delivery Model

- **Best effort service**
 - Network will do its best to get packet to destination
- Does NOT guarantee:
 - Any maximum latency or even ultimate success
 - Informing the sender if packet does not make it
 - Delivery of packets in same order as they were sent
 - Just one copy of packet will arrive
- Implications
 - Scales very well (really, it does)
 - Higher level protocols must make up for shortcomings
 - Reliably delivering ordered sequence of bytes → TCP
 - Some services not feasible (or hard)
 - Latency or bandwidth guarantees

6

IP Fragmentation



- Every network has own Maximum Transmission Unit (MTU)
 - Largest IP datagram it can carry within its own packet frame
 - E.g., Ethernet is 1500 bytes
 - Don't know MTUs of all intermediate networks in advance
- IP Solution
 - When hit network with small MTU, router fragments packet
 - Destination host reassembles the packet – why?

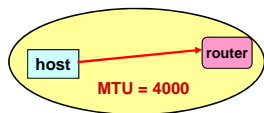
7

Fragmentation Related Fields

- Length
 - Length of IP fragment
- Identification
 - To match up with other fragments
- Flags
 - Don't fragment flag
 - More fragments flag
- Fragment offset
 - Where this fragment lies in entire IP datagram
 - Measured in 8 octet units (13 bit field)

8

IP Fragmentation Example #1

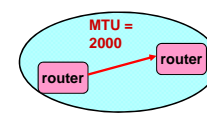


Length = 3820, M=0

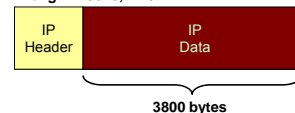


9

IP Fragmentation Example #2

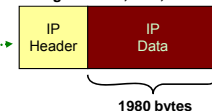


Length = 3820, M=0



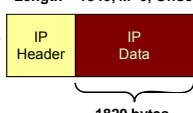
3800 bytes

Length = 2000, M=1, Offset = 0



1980 bytes

Length = 1840, M=0, Offset = 1980



1820 bytes

10

Fragmentation is Harmful

- Uses resources poorly
 - Forwarding costs per packet
 - Best if we can send large chunks of data
 - Worst case: packet just bigger than MTU
- Poor end-to-end performance
 - Loss of a fragment
- Path MTU discovery protocol → determines minimum MTU along route
 - Uses ICMP error messages
- Common theme in system design
 - Assure correctness by implementing complete protocol
 - Optimize common cases to avoid full complexity

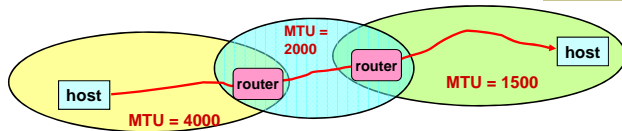
11

Internet Control Message Protocol (ICMP)

- Short messages used to send error & other control information
- Some functions supported by ICMP:
 - Ping request /response: check whether remote host reachable
 - Destination unreachable: Indicates how packet got & why couldn't go further
 - Flow control: Slow down packet transmit rate
 - Redirect: Suggest alternate routing path for future messages
 - Router solicitation / advertisement: Helps newly connected host discover local router
 - Timeout: Packet exceeded maximum hop limit
- How useful are they functions today?

12

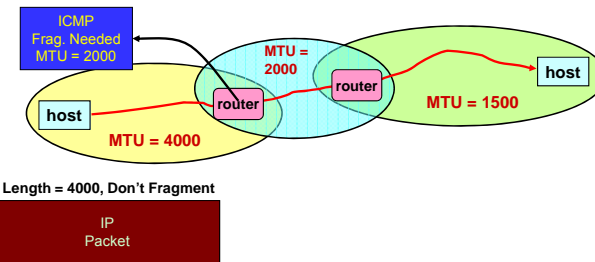
IP MTU Discovery with ICMP



- Typically send series of packets from one host to another
- Typically, all will follow same route
 - Routes remain stable for minutes at a time
- Makes sense to determine path MTU before sending real packets
- Operation: Send max-sized packet with "do not fragment" flag set
 - If encounters problem, ICMP message will be returned
 - "Destination unreachable: Fragmentation needed"
 - Usually indicates MTU problem encountered
- ICMP abuse? Other solutions?

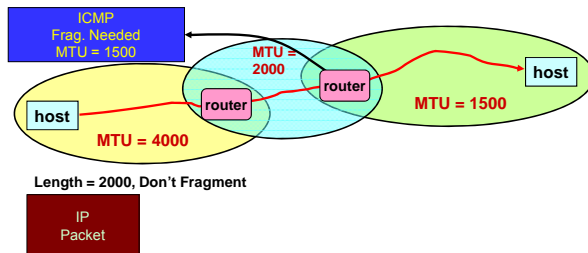
13

IP MTU Discovery with ICMP



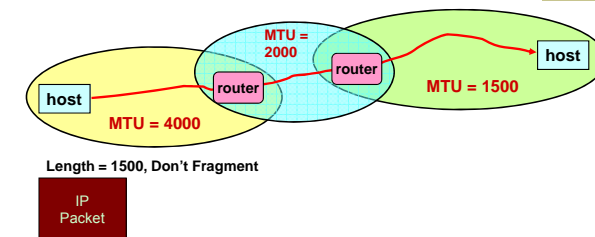
14

IP MTU Discovery with ICMP



15

IP MTU Discovery with ICMP



- When successful, no reply at IP level
 - "No news is good news"
- Higher level protocol might have some form of acknowledgement

16

Important Concepts



- Base-level protocol (IP) provides minimal service level
 - Allows highly decentralized implementation
 - Each step involves determining next hop
 - Most of the work at the endpoints
- ICMP provides low-level error reporting
- IP forwarding → global addressing, alternatives, lookup tables
- IP addressing → hierarchical, CIDR
- IP service → best effort, simplicity of routers
- IP packets → header fields, fragmentation, ICMP

17

Outline



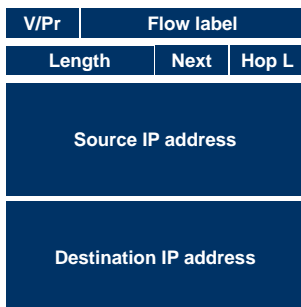
- The IP protocol
 - IPv4
 - IPv6
- IP in practice
 - Network address translation
 - Address resolution protocol
 - Tunnels

18

IPv6



- “Next generation” IP.
- Most urgent issue: increasing address space.
 - 128 bit addresses
- Simplified header for faster processing:
 - No checksum (why not?)
 - No fragmentation (really?)
- Support for guaranteed services: priority and flow id
- Options handled as “next header”
 - reduces overhead of handling options



19

IPv6 Address Size Discussion



- Do we need more addresses? Probably, long term
 - Big panic in 90s: “We’re running out of addresses!”
 - Big worry: Devices. Small devices. Cell phones, toasters, everything.
- 128 bit addresses provide space for structure (good!)
 - Hierarchical addressing is much easier
 - Assign an entire 48-bit sized chunk per LAN – use Ethernet addresses
 - Different chunks for geographical addressing, the IPv4 address space,
 - Perhaps help clean up the routing tables - just use one huge chunk per ISP and one huge chunk per customer.



20

IP Router Implementation: Fast Path versus Slow Path



- Common case: Switched in silicon (“fast path”)
 - Almost everything
- Weird cases: Handed to CPU (“slow path”, or “process switched”)
 - Fragmentation
 - TTL expiration (traceroute)
 - IP option handling
- Slow path is evil in today’s environment
 - “Christmas Tree” attack sets weird IP options, bits, and overloads router
 - Developers cannot (really) use things on the slow path
 - Slows down their traffic – not good for business
 - If it became popular, they are in trouble!

21

IPv6 Header Cleanup: Options



- 32 IPv4 options → variable length header
 - Rarely used
 - No development / many hosts/routers do not support
 - Worse than useless: Packets w/options often even get dropped!
 - Processed in “slow path”.
- IPv6 options: “Next header” pointer
 - Combines “protocol” and “options” handling
 - Next header: “TCP”, “UDP”, etc.
 - Extensions header: Chained together
 - Makes it easy to implement host-based options
 - One value “hop-by-hop” examined by intermediate routers
 - E.g., “source route” implemented only at intermediate hops

22

IPv6 Header Cleanup: “no”



- No checksum
 - Motivation was efficiency: If packet corrupted at hop 1, don’t waste b/w transmitting on hops 2..N.
 - Useful when corruption frequent, b/w expensive
 - Today: corruption is rare, bandwidth is cheap
- No fragmentation
 - Router discard packets, send ICMP “Packet Too Big”
 - host does MTU discovery and fragments
 - Reduced packet processing and network complexity.
 - Increased MTU a boon to application writers
 - Hosts can still fragment - using fragmentation header. Routers don’t deal with it any more.

23

Migration from IPv4 to IPv6



- Interoperability with IP v4 is necessary for incremental deployment.
 - No “flag day”
- Fundamentally hard because a (single) IP protocol is critical to achieving global connectivity across the internet
- Process uses a combination of mechanisms:
 - Dual stack operation: IP v6 nodes support both address types
 - Tunnel IP v6 packets through IP v4 clouds
 - IPv4-IPv6 translation at edge of network
 - NAT must not only translate addresses but also translate between IPv4 and IPv6 protocols
 - IPv6 addresses based on IPv4 – no benefit!
- 20 years later, this is still a major challenge!

24

Outline

- The IP protocol
 - IPv4
 - IPv6
- IP in practice
 - Network address translation
 - Address resolution protocol
 - Tunnels

25

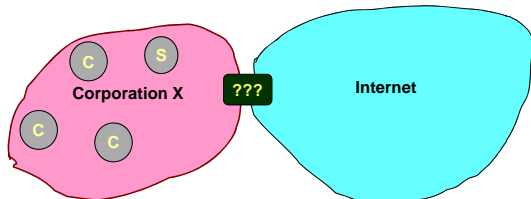
Altering the Addressing Model

- Original IP Model: Every host has unique IP address
- This has very attractive properties ...
 - Any host can communicate with any other host
 - Any host can act as a server
 - Just need to know host ID and port number
- ... but the system is open – complicates security
 - Any host can attack any other host
 - It is easy to forge packets
 - Use invalid source address
- ... and it places pressure on the address space
 - Every host requires “public” IP address

26

Challenges When Connecting to Public Internet

C: Client
S: Server

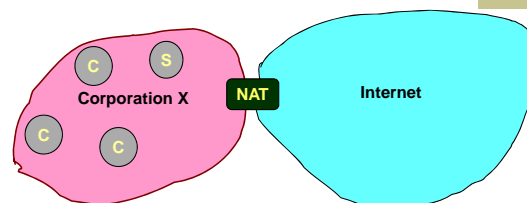


- Not enough IP addresses for every host in organization
 - Increasingly hard to get large address blocks
- Security
 - Don't want every machine in organization known to outside world
 - Want to control or monitor traffic in / out of organization

27

But not All Hosts are Equal!

C: Client
S: Server



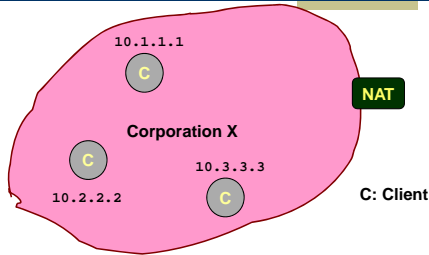
- Most machines within organization are used by individuals
 - For most applications, they act as clients
- Only a small number of machines act as servers for the entire organization
 - E.g., mail server, web, ..
 - All traffic to outside passes through firewall

(Most) machines within organization do not need public IP addresses!

28

Reducing Address Use: Network Address Translation

- Within organization:
 - assign each host a private IP address
 - IP addresses blocks 10/8 & 192.168/16 are set aside for this
 - Route within organization by IP protocol
 - Can do subnetting, ..
- The NAT translates between public and private IP addresses as packets travel to/from the public Internet
 - It does not let any packets from internal nodes "escape"
 - Outside world does not need to know about internal addresses

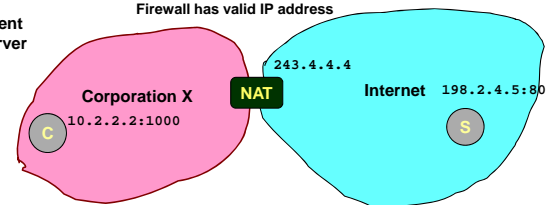


29

NAT: Opening Client Connection

C: Client
S: Server

Firewall has valid IP address



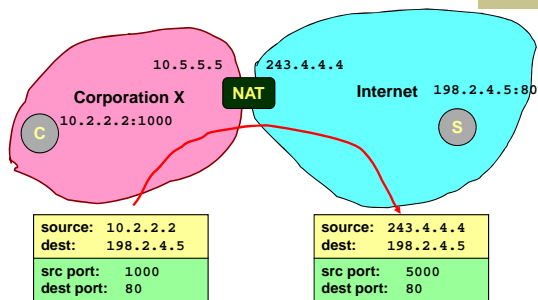
- Client 10.2.2.2 wants to connect to server 198.2.4.5:80
 - OS assigns ephemeral port (1000)
- Connection request intercepted by firewall
 - Maps client to port of firewall (5000)
 - Creates NAT table entry

Int Addr	Int Port	NAT Port
10.2.2.2	1000	5000

30

NAT: Client Request

C: Client
S: Server



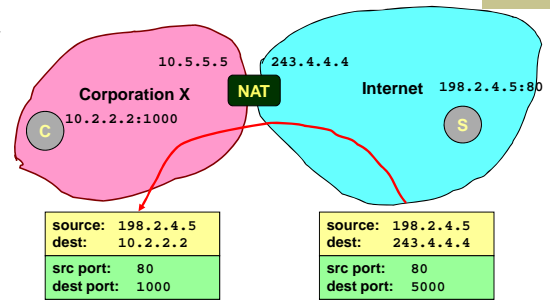
- Firewall acts as proxy for client
 - Intercepts message from client and marks itself as sender

Int Addr	Int Port	NAT Port
10.2.2.2	1000	5000

31

NAT: Server Response

C: Client
S: Server



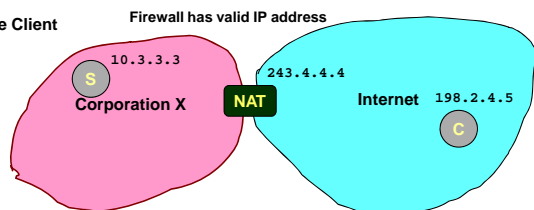
- Firewall acts as proxy for client
 - Acts as destination for server messages
 - Relabels destination to local addresses

Int Addr	Int Port	NAT Port
10.2.2.2	1000	5000

32

NAT: Enabling Servers

C: Remote Client
S: Server



- Use *port mapping* to make servers available

Int Addr	Int Port	NAT Port
10.3.3.3	80	80

- Manually configure NAT table to include entry for well-known port
- External users give address 243.4.4.4:80
- Requests forwarded to server

33

Additional NAT Benefits

- They significantly reduce the need for public IP addresses
- NATs directly help with security
 - Hides IP addresses used in internal network
 - Easy to change ISP: only NAT box needs to have IP address
 - Fewer registered IP addresses required
 - Basic protection against remote attack
 - Does not expose internal structure to outside world
 - Can control what packets come in and out of system
 - Can reliably determine whether packet from inside or outside
- And NATs have many additional benefits
 - NAT boxes make home networking simple
 - Can be used to map between addresses from different address families, e.g. IPv4 and IPv6

34

NAT Challenges

- NAT has to be consistent during a session.
 - Mapping (hard state) must be maintained during the session
 - Recall Goal 1 of Internet: Continue despite loss of networks or gateways
 - Recycle the mapping after the end of the session
 - May be hard to detect
- NAT only works for certain applications.
 - Some applications (e.g. ftp) pass IP information in payload - oops
 - Need application level gateways to do a matching translation
- NATs are a problem for peer-peer applications
 - File sharing, multi-player games, ...
 - Who is server?
 - Need to "punch" hole through NAT

35

Principle: Fate Sharing



- "You can lose state information relevant to an entity's connections if and only if the entity itself is lost"
 - Example: OK to lose TCP state if either endpoint crashes
 - The TCP connection is no longer useful anyway!
- It is NOT okay to lose it if an unrelated entity goes down
 - Example: if an intermediate router reboots
- NATs violate this principle: if a NAT goes down, all communication session it supports are lost!
 - Unless you add redundancy and put state in persistent storage
- Bad news: many stateful "middleboxes" violate this rule
 - Firewalls, mobility services, ... - more on this later
- Good news: today's hardware is very reliable

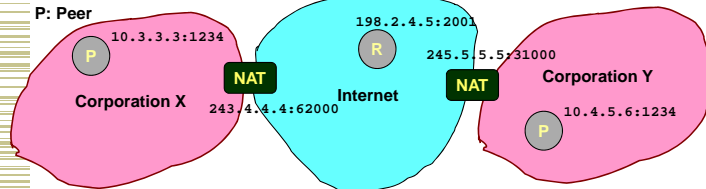
36

Many Options Exist for Peer-Peer



R: Rendezvous server

P: Peer



- NAT recognizes certain protocols and behaves as a application gateway
 - Used for standard protocols such as ftp
- Applications negotiate directly with NAT or firewall – need to be authorized
 - Multiple protocols dealing with different scenarios
- Punching holes in NAT: peers contact each other simultaneously using a known public (IP, port), e.g. used with rendezvous service
 - Use publicly accessible rendezvous service to exchange accessibility information
 - Assumes NATs do end-point independent mapping
- But remains painful!

37