

CSE/EE 461

IP/ICMP and the Network Layer

Last Time

- Focus:
 - What to do when one shared LAN isn't big enough?
- Interconnecting LANs
 - Bridges and LAN switches
 - But there are limits ...

Application
Presentation
Session
Transport
Network
Data Link
Physical

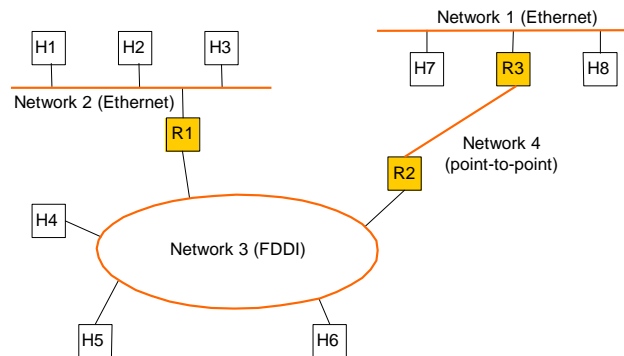
This Lecture

- Focus:
 - How do we build large networks?
- Introduction to the Network layer
 - Internetworks
 - Service models
 - IP, ICMP

Application
Presentation
Session
Transport
Network
Data Link
Physical

Internetworks

- Set of interconnected networks, e.g., the Internet
 - Scale and heterogeneity



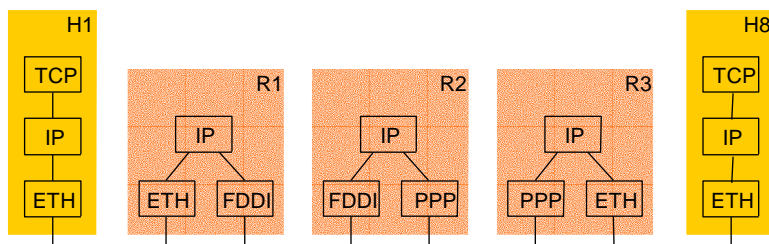
The Network Layer

- Job is to provide end-to-end data delivery between hosts on an internetwork
- Provides a higher layer of addressing

Application
Presentation
Session
Transport
Network
Data Link
Physical

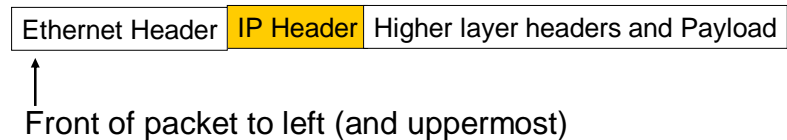
In terms of protocol stacks

- IP is the network layer protocol used in the Internet
- Routers are network level gateways
- Packet is the term for network layer PDUs



In terms of packet formats

- View of a packet on the wire on network 1 or 2
- Routers work with IP header, not higher
 - Higher would be a “layer violation”
- Routers strip and add link layer headers



Network Service Models

- Datagram delivery: postal service
 - Also connectionless, best-effort or unreliable service
 - Network can't guarantee delivery of the packet
 - Each packet from a host is routed independently
 - Example: IP
- Virtual circuit models: telephone
 - Also connection-oriented service
 - Signaling: connection establishment, data transfer, teardown
 - All packets from a host are routed the same way (router state)
 - Example: ATM, Frame Relay, X.25

Datagrams or Virtual Circuits?

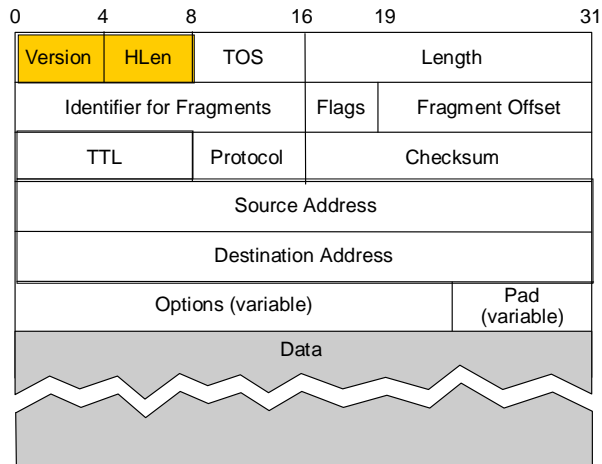
- Pros and Cons?
 - Simplicity/robustness versus stronger resource allocation
- We return to these tradeoffs later
 - Quality of Service (QOS)
 - These issues at the heart of current Internet evolution
 - Intserv (connection oriented) vs Diffserv (“connectionless”)

Internet Protocol (IP)

- IP (RFC791) defines a “best effort” service
 - May be loss, reordering, duplication, and errors!
 - Currently IPv4 (IP version 4), IPv6 on the way
- Routers forward packets using predetermined routes
 - Routing protocols (RIP, OSPF, BGP) run between routers to maintain routes (routing table, forwarding information base)
- Global, hierarchical addresses, not flat addresses
 - 32 bits in IPv4 address; 128 bits in IPv6 address
 - ARP (Address Resolution Protocol) maps IP to MAC addresses

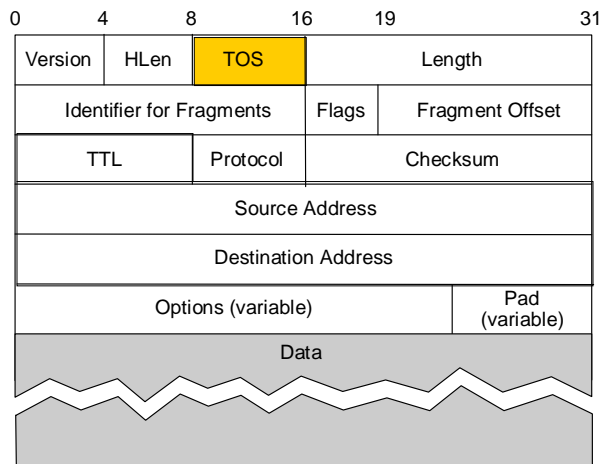
IPv4 Packet Format

- Version is 4
- Header length is number of 32 bit words
- Limits size of options



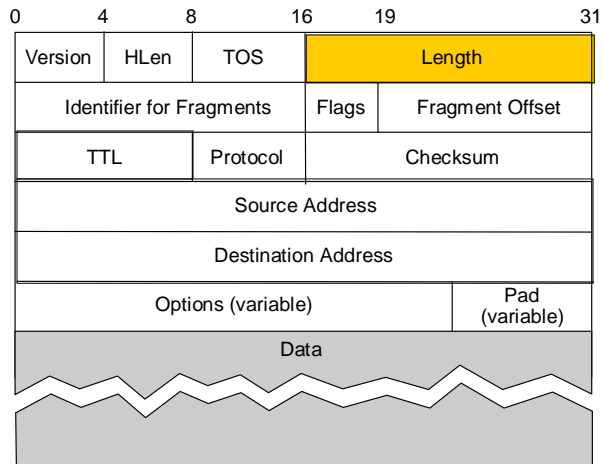
IPv4 Header Fields ...

- Type of Service
- Abstract notion, never really worked out
 - Routers ignored
- But now being redefined for Diffserv



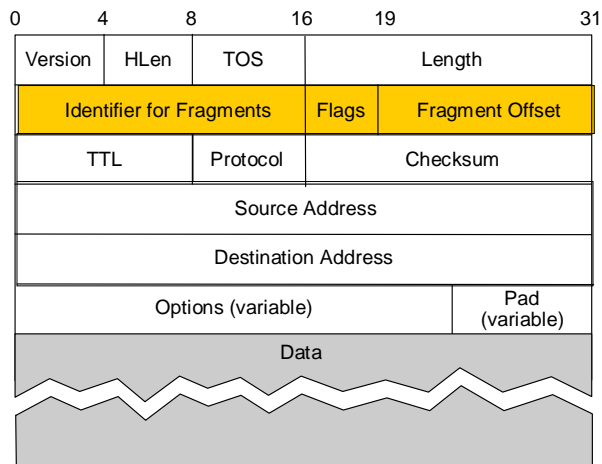
IPv4 Header Fields ...

- Length of packet
- Min 20 bytes, max 65K bytes (limit to packet size)



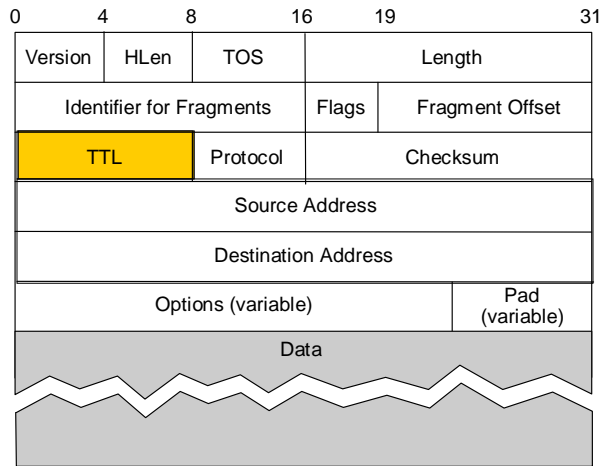
IPv4 Header Fields ...

- Fragment fields
- Different LANs have different frame size limits
- May need to break large packet into smaller fragments



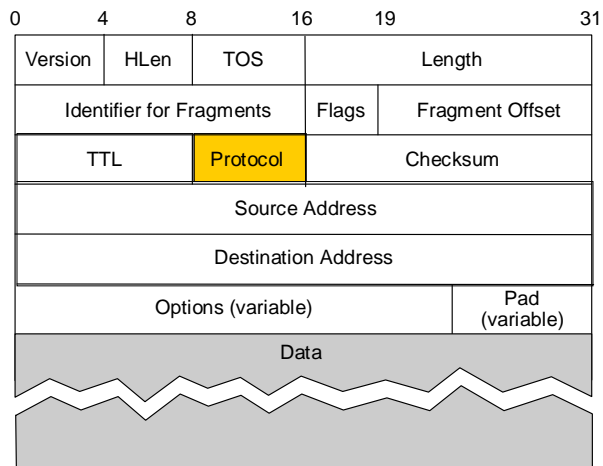
IPv4 Header Fields ...

- Time To Live
- Decremented by router and packet discarded if = 0
- Prevents immortal packets



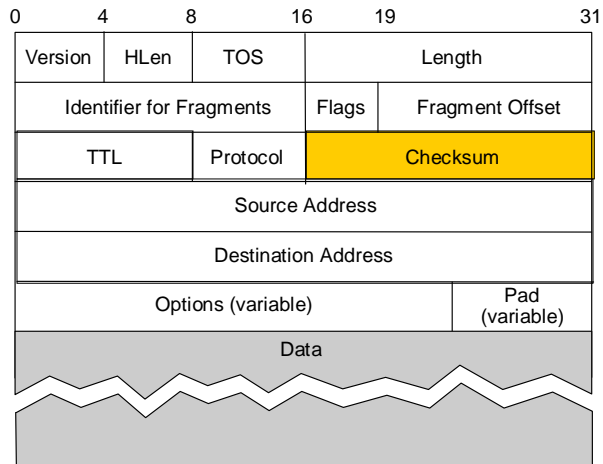
IPv4 Header Fields ...

- Identifies higher layer protocol
 - E.g., TCP, UDP



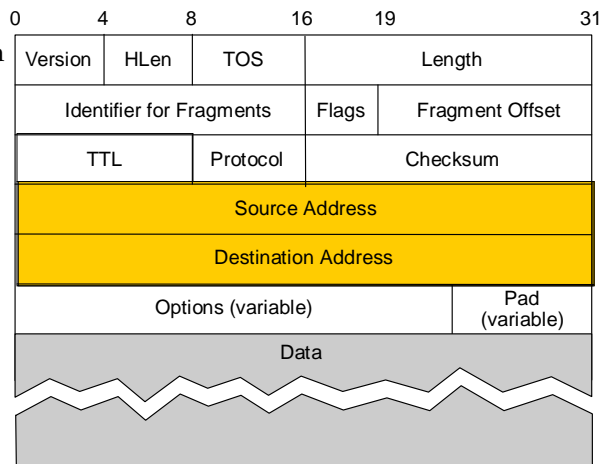
IPv4 Header Fields ...

- Header checksum
- Recalculated by routers (TTL drops)
- Doesn't cover data
- Disappears for IPv6



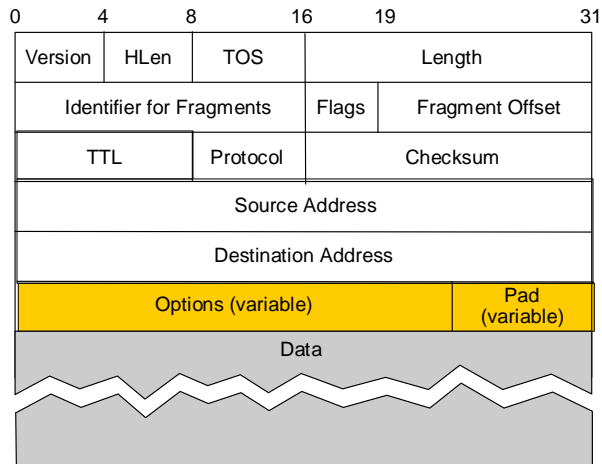
IPv4 Header Fields ...

- Source/destination addresses
 - Not Ethernet
- Unchanged by routers
- Not authenticated by default



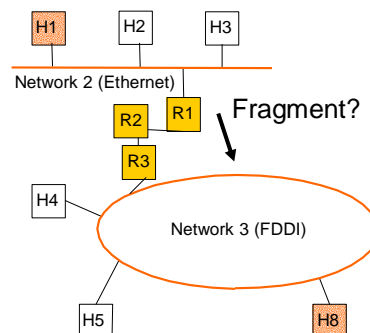
IPv4 Header Fields ...

- IP options indicate special handling
 - Timestamps
 - "Source" routes
- Rarely used ...



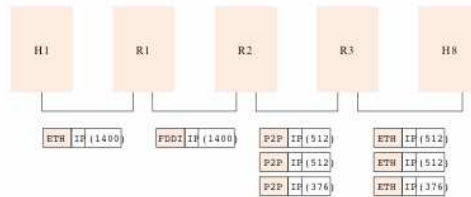
Fragmentation Issue

- Different networks may have different frame limits (MTUs)
 - Ethernet 1.5K, FDDI 4.5K
- Don't know if packet will be too big for path beforehand
 - IPv4: fragment on demand and reassemble at destination
 - IPv6: network returns error message so host can learn limit



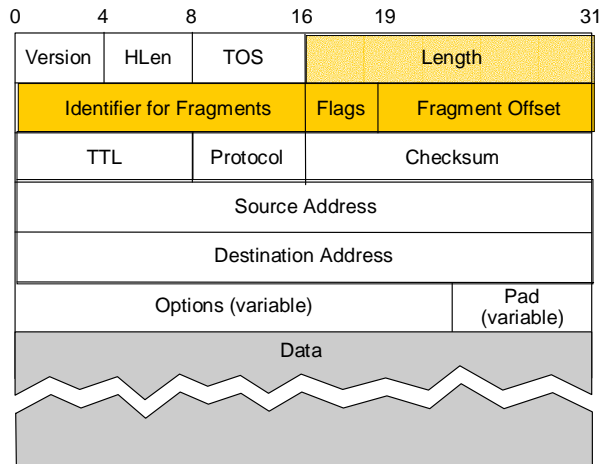
Fragmentation and Reassembly

- Each network has some MTU
- Strategy
 - fragment when necessary (MTU < Datagram)
 - try to avoid fragmentation at source host
 - refragmentation is possible
 - fragments are self-contained datagrams
 - use CS-PDU (not cells) for ATM
 - delay reassembly until destination host
 - do not recover from lost fragments
- Example

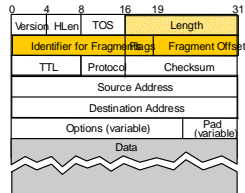


Fragment Fields

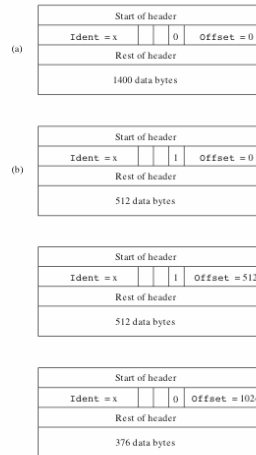
- Fragments of one packet identified by (source, dest, frag id) triple
 - Make unique
- Offset gives start, length changed
- Flags are More Fragments (MF) Don't Fragment (DF)



Fragmenting a Packet



Packet Format



Fragment Considerations

- Relating fragments to original datagram provides:
 - Tolerance of loss, reordering and duplication
 - Ability to fragment fragments
- Reassembly done at the endpoint
 - Puts pressure on the receiver
- Consequences of fragmentation:
 - Loss of any fragments causes loss of entire packet
 - The packet train and buffer overflow
 - Need to time-out reassembly when any fragments lost

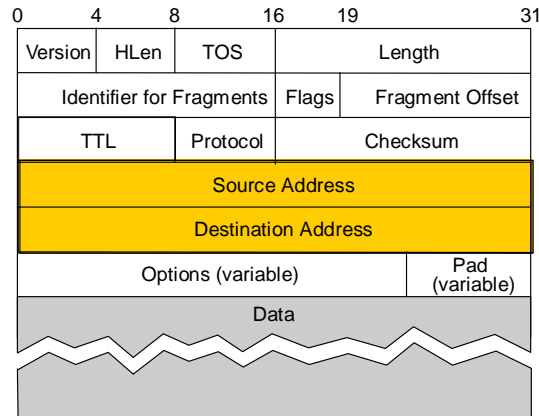
Fragmentation Issues Summary

- Causes inefficient use of resources within the network
 - BW, CPU
 - Eg, App sends 1024 bytes across ARPANET (1007 MTU)
 - 1024 + 40 for TCP/IP header
 - Frag 1 == 1000, Frag 2 == 84
 - Should have sent 1006 bytes!
- Higher level protocols must retransmit entire datagram
 - Really hard with “guaranteed packet loss”
- Efficient reassembly is hard
 - Lots of special cases
 - (think linked lists)

Avoiding Fragmentation

- Always send small datagrams
 - Might be too small
- “Guess” MTU of path
 - Use DF flag. May have large startup time
- Discover actual MTU of path
 - One RT delay w/help, much more w/o.
 - “Help” requires router support
- Guess or discover, but be willing to accept your mistakes

What is an Internet Address?

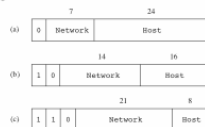


Global Addresses

- Properties
 - globally unique
 - hierarchical: network + host

1. Small number of large networks
2. Modest # of medium sized networks
3. Many small networks

- Format



CLASS	SIZE	NUMBER
A	2G	126
B		
C	254	2M

- Dot notation

- 10.3.2.4
- 128.96.33.81
- 192.12.69.77

Original Rationale: Beware the Routing Tables

1. You don't care about most networks.
2. The few networks you do care about, you care about them a lot.
3. Not many routing table entries get you "closer" to a lot of the hosts

Datagram Forwarding

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

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

Address Translation

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

ARP Packets

0	8	16	31
Hardware type = 1		ProtocolType = 0x0800	
HLEN = 48	PLEN = 32		Operation
SourceHardwareAddr (bytes 0-3)			
SourceHardwareAddr (bytes 4-5)		SourceProtocolAddr (bytes 6-7)	
SourceProtocolAddr (bytes 8-9)		TargetHardwareAddr (bytes 10-15)	
TargetHardwareAddr (bytes 16-21)			
TargetProtocolAddr (bytes 22-31)			

- HardwareType: type of physical network (e.g., Ethernet)
- ProtocolType: type of higher layer protocol (e.g., IP)
- HLEN & PLEN: length of physical and protocol addresses
- Operation: request or response
- Source/Target Physical/Protocol addresses

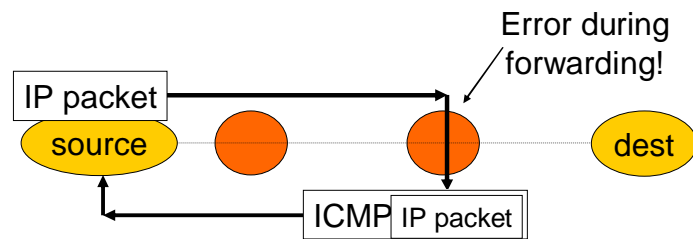
- Notes

- table entries timeout in about 10 minutes
- update table with source when you are the target
- update table if already have an entry
- do not refresh table entries upon reference

ICMP

- What happens when things go wrong?
 - Need a way to test/debug a large, widely distributed system
- ICMP = Internet Control Message Protocol (RFC792)
 - Companion to IP – required functionality
- Used for error and information reporting:
 - Errors that occur during IP forwarding
 - Queries about the status of the network

ICMP Generation



Common ICMP Messages

- Destination unreachable
 - "Destination" can be host, network, port or protocol
- Packet needs fragmenting but DF is set
- Redirect
 - To shortcut circuitous routing
- TTL Expired
 - Used by the "traceroute" program
- Echo request/reply
 - Used by the "ping" program
- Cannot Fragment
- Busted Checksum

- ICMP messages include portion of IP packet that triggered the error (if applicable) in their payload

ICMP Restrictions

- The generation of error messages is limited to avoid cascades ... error causes error that causes error!
- Don't generate ICMP error in response to:
 - An ICMP error
 - Broadcast/multicast messages (link or IP level)
 - IP header that is corrupt or has bogus source address
 - Fragments, except the first
- ICMP messages are often rate-limited too.

Key Concepts

- Network layer provides end-to-end data delivery across an internetwork, not just a LAN
 - Datagram and virtual circuit service models
 - IP/ICMP is the network layer protocol of the Internet
- Up next: More detailed look at routing and addressing