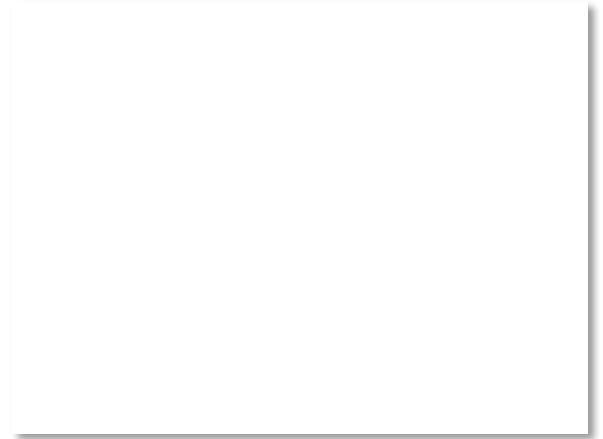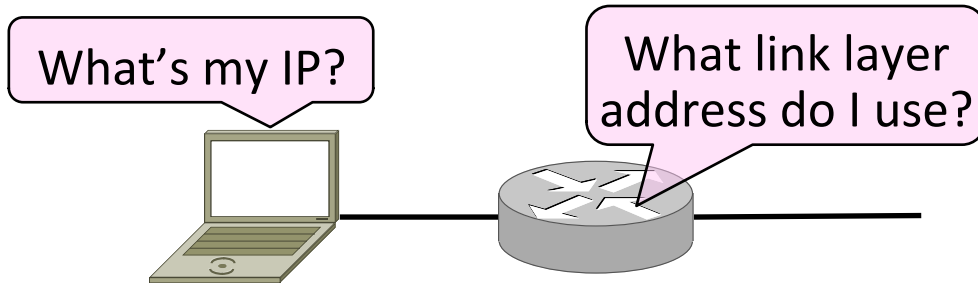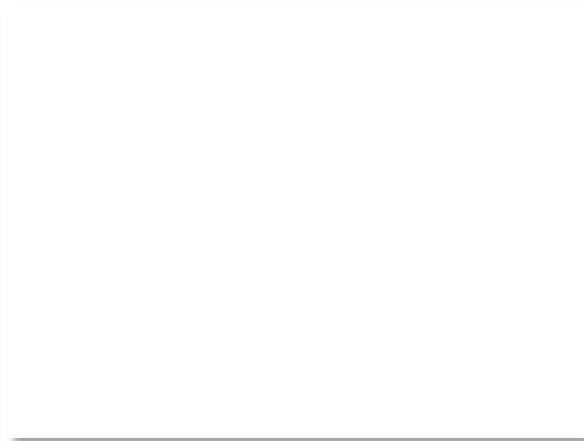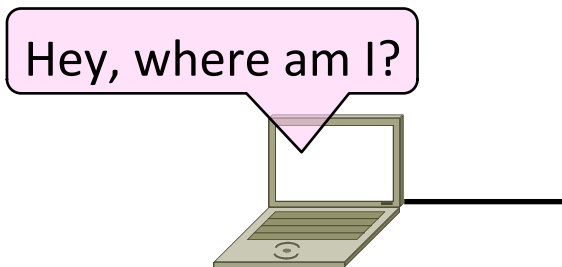# Topic

- Filling in the gaps we need to make for IP forwarding work in practice
  - Getting IP addresses (DHCP) **»**
  - Mapping IP to link addresses (ARP) **»**
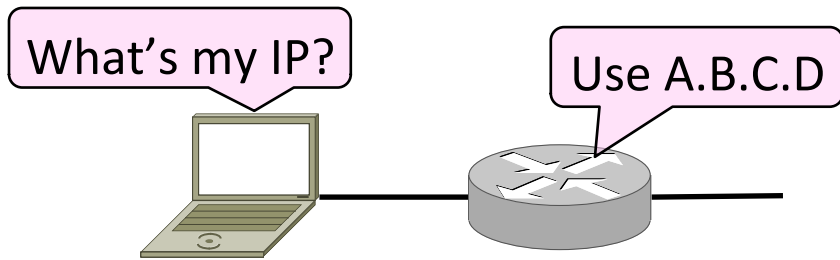
What's my IP?

What link layer address do I use?

# Getting IP Addresses

- Problem:
  - A node wakes up for the first time …
  - What is its IP address? What's the IP address of its router? Etc.
  - At least Ethernet address is on NIC

  Hey, where am I?

# Getting IP Addresses (2)

1. Manual configuration (old days)
   – Can't be factory set, depends on use
2. A protocol for automatically configuring addresses (DHCP) **»**
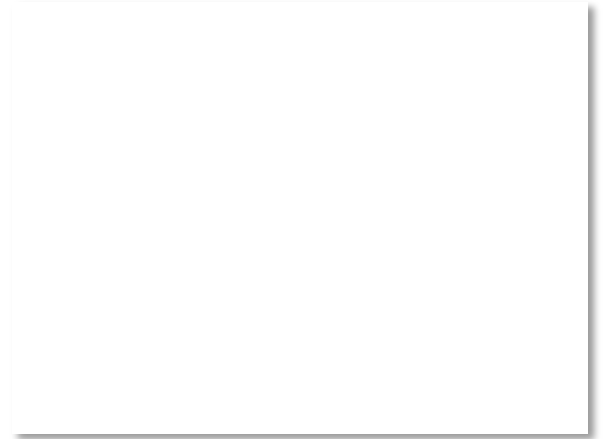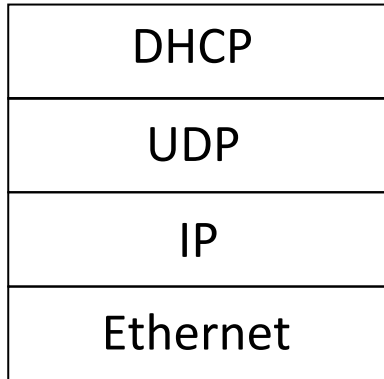   – Shifts burden from users to IT folk

What's my IP?

Use A.B.C.D

# DHCP

- DHCP (Dynamic Host Configuration Protocol), from 1993, widely used

- It leases IP address to nodes
- Provides other parameters too
  - Network prefix
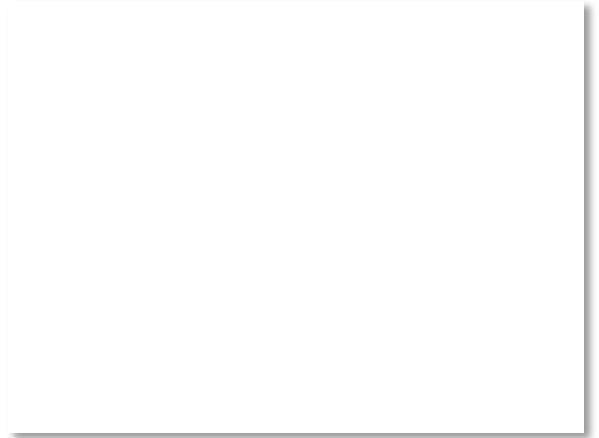  - Address of local router
  - DNS server, time server, etc.

# DHCP Protocol Stack

- DHCP is a client-server application
  - Uses UDP ports 67, 68

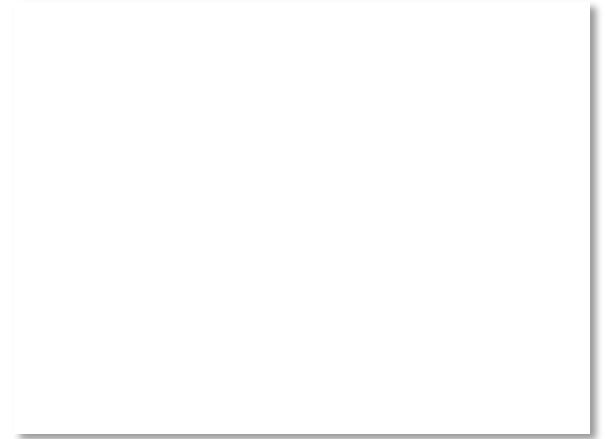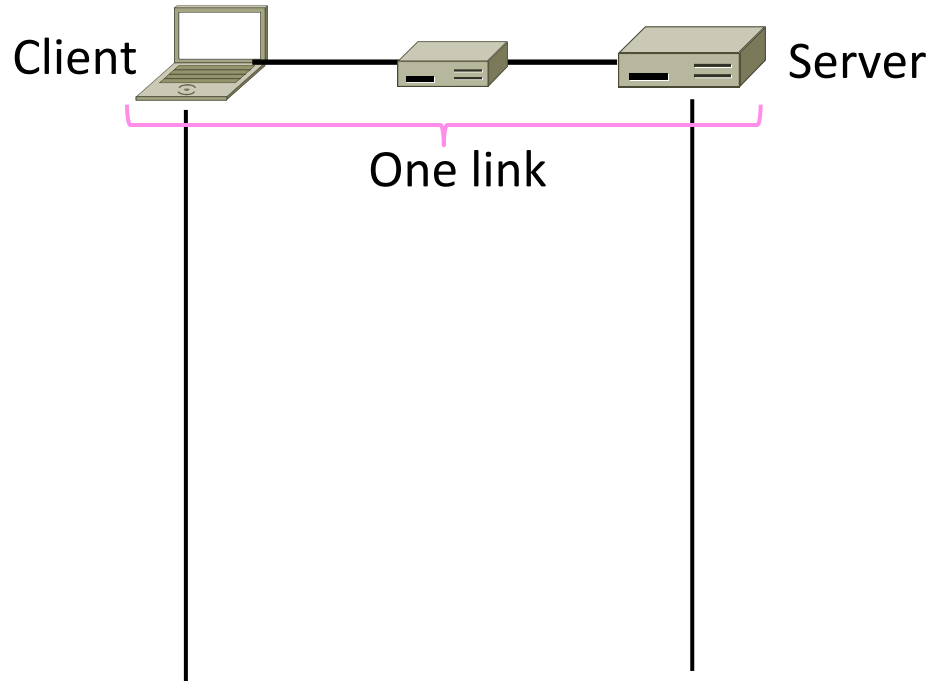| DHCP |
|:---:|
| UDP |
| IP |
| Ethernet |

# DHCP Addressing

- Bootstrap issue:
  - How does node send a message to DHCP server before it is configured?

- Answer:
  - Node sends <u>broadcast</u> messages that delivered to all nodes on the network
  - <u>Broadcast address</u> is all 1s
  - IP (32 bit): 255.255.255.255
  - Ethernet (48 bit): ff:ff:ff:ff:ff:ff

# DHCP Messages

Client    Server

One link

# DHCP Messages (2)

# DHCP Messages (3)

- To renew an existing lease, an abbreviated sequence is used:
  - REQUEST, followed by ACK

- Protocol also supports replicated servers for reliability

# Sending an IP Packet

- Problem:
  - A node needs Link layer addresses to send a frame over the local link
  - How does it get the destination link address from a destination IP address?

# ARP (Address Resolution Protocol)

- Node uses to map a local IP address to its Link layer addresses

Link layer

| Source Ethernet | Dest. Ethernet | Source IP | Dest. IP | Payload ... |
|---|---|---|---|---|

From NIC

From ARP

From DHCP

# ARP Protocol Stack

- ARP sits right on top of link layer
  - No servers, just asks node with target IP to identify itself
  - Uses broadcast to reach all nodes

| ARP |
|---|
| Ethernet |

# ARP Messages

Node ▭━━▭━━▭ Target

One link

# ARP Messages (2)



Node

Target

REQUEST

Broadcast

Who has IP 1.2.3.4?

REPLY

I do at 1:2:3:4:5:6

# Discovery Protocols

- Help nodes find each other
  - There are more of them!
    - E.g., zeroconf, Bonjour

- Often involve broadcast
  - Since nodes aren't introduced
  - Very handy glue

# Other Aspects of Forwarding

- It's not all about addresses …

# Other Aspects (2)

- Decrement TTL value
  - Protects against loops
- Checks header checksum
  - To add reliability
- Fragment large packets
  - Split to fit it on next link
- Send congestion signals
  - Warns hosts of congestion
- Generates error messages
  - To help mange network
- Handle various options

Coming later

# Topic

- ## How do we connect networks with different maximum packet sizes?
  - Need to split up packets, or discover the largest size to use

# Packet Size Problem

- Different networks have different maximum packet sizes
  - Or MTU (<u>Maximum Transmission Unit</u>)
  - E.g., Ethernet 1.5K, WiFi 2.3K

- Prefer large packets for efficiency
  - But what size is too large?
  - Difficult because node does not know complete network path

# Packet Size Solutions

- Fragmentation (now)
  - Split up large packets in the network if they are too big to send
  - Classic method, dated

- Discovery (next)
  - Find the largest packet that fits on the network path and use it
  - IP uses today instead of fragmentation

# IPv4 Fragmentation

- Routers fragment packets that are too large to forward
- Receiving host reassembles to reduce load on routers

# IPv4 Fragmentation Fields

- Header fields used to handle packet size differences
  - Identification, Fragment offset, MF/DF control bits

# IPv4 Fragmentation Procedure

- Routers split a packet that is too large:
  - Typically break into large pieces
  - Copy IP header to pieces
  - Adjust length on pieces
  - Set offset to indicate position
  - Set MF (More Fragments) on all pieces except last

- Receiving hosts reassembles the pieces:
  - Identification field links pieces together, MF tells receiver when it has all pieces

# IPv4 Fragmentation (2)

Before
MTU = 2300

ID = 0x12ef
Data Len = 2300
Offset = 0
MF = 0

(Ignore length of headers)

After
MTU = 1500

ID =
Data Len =
Offset =
MF =

ID =
Data Len =
Offset =
MF =

# IPv4 Fragmentation (3)

Before
MTU = 2300

ID = 0x12ef
Data Len = 2300
Offset = 0
MF = 0

After
MTU = 1500

ID = 0x12ef
Data Len = 1500
Offset = 0
MF = 1

ID = 0x12ef
Data Len = 800
Offset = 1500
MF = 0

# IPv4 Fragmentation (4)

- It works!
  - Allows repeated fragmentation

- But fragmentation is undesirable
  - More work for routers, hosts
  - Tends to magnify loss rate
  - Security vulnerabilities too

# Path MTU Discovery

- Discover the MTU that will fit
  - So we can avoid fragmentation
  - The method in use today

- Host tests path with large packet
  - Routers provide feedback if too large; they tell host what size would have fit

# Path MTU Discovery (2)

Packet (with length)

1400

MTU=1400  MTU=1200 bytes  MTU=900

Source  Destination

# Path MTU Discovery (3)

# Path MTU Discovery (4)

- Process may seem involved
  - But usually quick to find right size

- Path MTU depends on the path and so can change over time
  - Search is ongoing

- Implemented with ICMP (next)
  - Set DF (Don't Fragment) bit in IP header to get feedback messages

# Topic

- What happens when something goes wrong during forwarding?
  - Need to be able to find the problem

# Internet Control Message Protocol

- ICMP is a companion protocol to IP
  - They are implemented together
  - Sits on top of IP (IP Protocol=1)

- Provides error report and testing
  - Error is at router while forwarding
  - Also testing that hosts can use

# ICMP Errors

- When router encounters an error while forwarding:
  - It sends an ICMP error report back to the IP source address
  - It discards the problematic packet; host needs to rectify

# ICMP Message Format

- Each ICMP message has a Type, Code, and Checksum

- Often carry the start of the offending packet as payload

- Each message is carried in an IP packet

# ICMP Message Format (2)

- Each ICMP message has a Type, Code, and Checksum
- Often carry the start of the offending packet as payload
- Each message is carried in an IP packet

Portion of offending packet,
starting with its IP header

| Src=router, Dst=A<br>Protocol = 1 | Type=X, Code=Y | Src=A, Dst=B<br>XXXXXXXXXXXXXX |
|---|---|---|
| IP header | ICMP header | ICMP data |

# Example ICMP Messages

| Name | Type / Code | Usage |
|------|-------------|-------|
| Dest. Unreachable (Net or Host) | 3 / 0 or 1 | Lack of connectivity |
| Dest. Unreachable (Fragment) | 3 / 4 | Path MTU Discovery |
| Time Exceeded (Transit) | 11 / 0 | Traceroute |
| Echo Request or Reply | 8 or 0 / 0 | Ping |

Testing, not a forwarding error: Host sends Echo
Request, and destination responds with an Echo Reply

# Traceroute

- IP header contains TTL (Time to live) field
  - Decremented every router hop, with ICMP error if it hits zero
  - Protects against forwarding loops

| Version | IHL | Differentiated Services | | | Total length | |
|---------|-----|------------------------|---|---|--------------|---|
| Identification | | | | DF MF | Fragment offset | |
| Time to live | | Protocol | | | Header checksum | |
| Source address | | | | | | |
| Destination address | | | | | | |
| Options (0 or more words) | | | | | | |

# Traceroute (2)

- Traceroute repurposes TTL and ICMP functionality
  - Sends probe packets increasing TTL starting from 1
  - ICMP errors identify routers on the path



Local Host

Remote Host

1 hop
2 hops
3 hops
N-1 hops
N hops

. . .

# Topic

- IP version 6, the future of IPv4 that is now (still) being deployed

Why do I want IPv6 again?

# Internet Growth

- At least a billion Internet hosts and growing …

- And we're using 32-bit addresses!

**Internet Domain Survey Host Count**



Source: Internet Systems Consortium (www.isc.org)

# The End of New IPv4 Addresses

- Now running on leftover blocks held by the regional registries; much tighter allocation policies

Exhausted on 4/11 and 9/12!

ARIN (US, Canada)

APNIC (Asia Pacific)

RIPE (Europe)

IANA (All IPs)

ISPs

Companies

LACNIC (Latin Amer.)

AfriNIC (Africa)

Exhausted on 2/11!

End of the world ? 12/21/12?

# IP Version 6 to the Rescue

- Effort started by the IETF in 1994
  - Much larger addresses (128 bits)
  - Many sundry improvements

- Became an IETF standard in 1998
  - Nothing much happened for a decade
  - Hampered by deployment issues, and a lack of adoption incentives
  - Big push ~2011 as exhaustion looms

# IPv6 Deployment

Percentage of users accessing Google via IPv6

Time for growth!



Source: Google IPv6 Statistics, 30/1/13

# IPv6

- Features large addresses
  - 128 bits, most of header
- New notation
  - 8 groups of 4 hex digits (16 bits)
  - Omit leading zeros, groups of zeros

Ex:  2001:0db8:0000:0000:0000:ff00:0042:8329
  →

32 bits

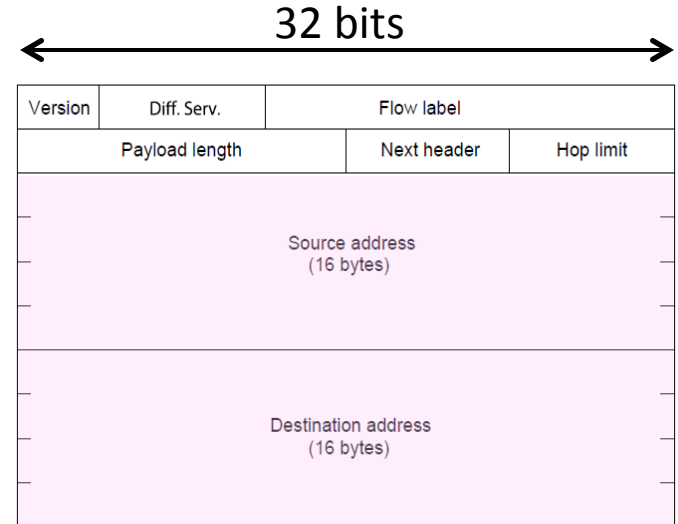| Version | Diff. Serv. | Flow label | | |
| Payload length | | | Next header | Hop limit |

Source address
(16 bytes)

Destination address
(16 bytes)

# IPv6 (2)

- Lots of other, smaller changes
  - Streamlined header processing
  - Flow label to group of packets
  - Better fit with "advanced" features (mobility, multicasting, security)

32 bits

| Version | Diff. Serv. | Flow label | | |
|---------|-------------|------------|---|---|
| Payload length | | | Next header | Hop limit |

Source address
(16 bytes)

Destination address
(16 bytes)

# IPv6 Transition

- The Big Problem:
  - How to deploy IPv6?
  - Fundamentally incompatible with IPv4

- Dozens of approaches proposed
  - Dual stack (speak IPv4 and IPv6)
  - Translators (convert packets)
  - Tunnels (carry IPv6 over IPv4) »

# Tunneling

- Native IPv6 islands connected via IPv4
  - Tunnel carries IPv6 packets across IPv4 network

# Tunneling (2)

- Tunnel acts as a single link across IPv4 network

# Tunneling (3)

- Tunnel acts as a single link across IPv4 network
  - Difficulty is to set up tunnel endpoints and routing

# Topic

- What is NAT (Network Address Translation)? How does it work?
  - NAT is widely used at the edges of the network, e.g., homes



I'm a NAT box too!

Internet

# Layering Review

- Remember how layering is meant to work?
  - "Routers don't look beyond the IP header." Well …

# Middleboxes

- Sit "inside the network" but perform "more than IP" processing on packets to add new functionality
  - NAT box, Firewall / Intrusion Detection System

# Middleboxes (2)

- Advantages
  - A possible rapid deployment path when there is no other option
  - Control over many hosts (IT)

- Disadvantages
  - Breaking layering interferes with connectivity; strange side effects
  - Poor vantage point for many tasks

# NAT (Network Address Translation) Box

- NAT box connects an internal network to an external network
  - Many internal hosts are connected using few external addresses
  - Middlebox that "translates addresses"

- Motivated by IP address scarcity
  - Controversial at first, now accepted

# NAT (2)

- Common scenario:
  - Home computers use "private" IP addresses
  - NAT (in AP/firewall) connects home to ISP using a single external IP address

Unmodified computers at home

Looks like one computer outside



NAT box

ISP

# How NAT Works

- Keeps an internal/external table
  - Typically uses IP address + TCP port
  - This is address and port translation

<span style="color:magenta">What host thinks</span>    <span style="color:magenta">What ISP thinks</span>

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |
| 192.168.1.13 : 1234 | 44.25.80.3 : 1501 |
| 192.168.2.20 : 1234 | 44.25.80.3 : 1502 |

- Need ports to make mapping 1-1 since there are fewer external IPs

# How NAT Works (2)

- Internal → External:
  - Look up and rewrite Source IP/port

Internal source

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |

External destination
IP=X, port=Y

Src =

Dst =

NAT box

Src =

Dst =

# How NAT Works (3)

- External → Internal
  - Look up and rewrite Destination IP/port

Internal
destination

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | 44.25.80.3 : 1500 |

External
source
IP=X, port=Y

NAT box

Src =

Dst =

Src =

Dst =

# How NAT Works (4)

- Need to enter translations in the table for it to work
  - Create external name when host makes a TCP connection

Internal source

| Internal  IP:port | External  IP : port |
|---|---|
| 192.168.1.12 : 5523 | |

External destination
IP=X, port=Y

NAT box

Src =

Dst =

Src =

Dst =

# NAT Downsides

- Connectivity has been broken!
    - Can only send incoming packets after an outgoing connection is set up
    - Difficult to run servers or peer-to-peer apps (Skype) at home

- Doesn't work so well when there are no connections (UDP apps)

- Breaks apps that unwisely expose their IP addresses (FTP)

# NAT Upsides

- Relieves much IP address pressure
  - Many home hosts behind NATs
- Easy to deploy
  - Rapidly, and by you alone
- Useful functionality
  - Firewall, helps with privacy

- Kinks will get worked out eventually
  - "NAT Traversal" for incoming traffic