

# Network Layer

Arnav Gupta

December 5, 2024

## Contents

<b>1</b>	<b>Overview</b>	<b>2</b>
1.1	Best Effort Service . . . . .	3
<b>2</b>	<b>Data Plane: The Internet Protocol</b>	<b>3</b>
2.1	Fragmentation and Assembly . . . . .	4
2.2	IPv4 Addressing . . . . .	5
2.2.1	Subnets . . . . .	5
2.2.2	DHCP . . . . .	5
2.2.3	Hierarchical Addressing . . . . .	6
2.2.4	Network Address Translation (NAT) . . . . .	7
2.2.5	Middleboxes . . . . .	8
2.3	IPv6 Addressing . . . . .	9
<b>3</b>	<b>What's Inside a Router</b>	<b>9</b>
3.1	Switching Fabrics . . . . .	10
3.1.1	Switching via Memory . . . . .	10
3.1.2	Switching via Bus . . . . .	10
3.1.3	Interconnection Network . . . . .	10
3.2	Port Queuing . . . . .	11
3.2.1	Input Port Queuing . . . . .	11
3.2.2	Output Port Queuing . . . . .	11
<b>4</b>	<b>Control Plane: ICMP, Routing</b>	<b>11</b>
4.1	ICMP . . . . .	12
4.2	Routing . . . . .	12
4.2.1	Intra-AS Routing . . . . .	14
4.2.2	Inter-AS Routing: BGP . . . . .	18

# 1 Overview

Network layer protocols transport segments from sending to receiving hosts. The sender encapsulates segments into datagrams and passes them to the link layer. The receiver delivers segments to the transport layer protocol.

Every Internet device uses network layer protocols.

**Routers** examine header fields in all IP datagrams passing through them and move datagrams from input interfaces to output interfaces to transfer datagrams along end-to-end paths.

Key network-layer functions:

- **forwarding**: move packets from a router’s input link to appropriate router output links (done for all datagrams very fast)
- **routing**: network-wide process that determines route taken by packets from source to destination to fill forwarding tables (done in the background, takes longer), uses routing algorithms

The **data plane** is a local, per-router function and determines how a datagram arriving on router input interface is forwarded to router output interface.

The **control plane** uses network-wide logic to determine how a datagram is routed among routers along end-to-end paths from source host to destination host.

Control plane approaches:

- **traditional routing algorithms**: implemented in routers (both planes implemented monolithically within a router)
  - individual routing algorithm components in every router interact in the control plane
- **software-defined networking**: explicitly separate the two planes by implementing the control plane as a service in remote servers
  - remote controller computes and installs forwarding tables in routers

## 1.1 Best Effort Service

The Internet runs on a “best effort” service model, so there are no guarantees on:

- successful datagram delivery to the destination
- timing or order of delivery
- bandwidth available for end-to-end flow

Best effort is simple and has allowed Internet to be widely adopted. Successful provisioning of bandwidth allows performance of real-time applications to be good enough for most of the time.

Replicated, application-layer distributed services connect close to client networks which allows services to be provided from multiple locations.

## 2 Data Plane: The Internet Protocol

Data plane is:

- connectionless (datagram-based)
- best-effort delivery
  - packets can be lost, delivered out of order, or delayed
- a common packet format for IPv4 and new packet format for IPv6
- global addressing for identifying all hosts (ARP)
- sister protocol that performs error reporting and enables signaling between routers: ICMP (v4, v6)

IPv4 datagram has

- IP version number
- header length (bytes)
- type of service
- total datagram length (bytes)
- 16-bit identifier
- flag

- fragmentation/reassembly and offset info
- time to live: remaining max hops
- upper layer protocol (TCP or UDP)
- header checksum
- 32-bit source IP address
- 32-bit destination IP address
- options (timestamp, record route taken, etc)
- payload data

When no options, overhead is 20 bytes. Upper layer protocol can be:

- ICMP: 1
- TCP: 6
- UDP: 17
- IPv4: 4
- IPv6: 41

If header checksum detects an error, datagram is dropped. Must be recomputed at every hop because of TTL and options.

## 2.1 Fragmentation and Assembly

Link layer protocols have maximum transfer unit size, which is the largest possible data size in a frame.

Large IP datagrams can be divided (fragmented) within a network, and then reassembled at the final destination. IP header bits are used to identify and order related fragments.

A receiver cannot hold fragments forever and fragments can arrive out of order, so loss of fragments can mean loss of entire datagram. The receiver starts a timer when the first fragment of a datagram arrives. If the timer expires before all the fragments are received, those already received are discarded.

Fragmentation complicated routers and end-systems, which is used by attackers.

## 2.2 IPv4 Addressing

**IP address:** 32-bit identifier associated with each host or router interface (about 4 billion total)

**Interface:** connection between host/router and physical link

- routers typically have multiple interfaces
- hosts typically have 1 or 2 interfaces

### 2.2.1 Subnets

The addressing scheme is **Classless InterDomain Routing (CIDR)**:

- IP addresses have a subnet part and a host part
- the subnet portion of the address has arbitrary length
- the address has format a.b.c.d/x where x is the prefix, the number of bits in the subnet portion

**Subnet:** set of interfaces that have IP addresses with the same prefix and same subnet portion, and can physically reach each other without passing through an intervening router

A subnet mask is x 1 bits followed by 0 bits, which is bitwise ANDed to the IP address. With masking, for the same subnet, can send a datagram directly to the destination and for a different subnet, send the datagram to the router.

### 2.2.2 DHCP

An IP address can be hard-coded by sysadmin in the config file or using **DHCP (Dynamic Host Configuration Protocol)** which dynamically gets address from a server when the host joins the network.

DHCP allows:

- addresses to be renewed on use
- reused since only holding addresses while host connected/on

For a client-server protocol using UDP, DHCP is a network function implemented as an application protocol:

- host broadcasts DHCP discover message

- DHCP server responds with DHCP offer message
- host requests IP address with a DHCP request message
- DHCP server sends the address with a DHCP ack message

Only last 2 steps are needed if the client remembers and wishes to reuse a previously allocated network address.

A DHCP server is co-located in the router, serving all subnets to which the router is attached.

DHCP also returns:

- address of first hop router for client
- name and IP address of DNS server
- address prefix (indicating network vs host portion of address)

For a network to get the subnet part of the IP address, it gets allocated a portion of its provider ISP's address space.

### 2.2.3 Hierarchical Addressing

Allows efficient advertisement of routing info.

When the router receives a datagram with some destination address:

1. mask the destination address with the mask for that row of the forwarding table
2. check if the results correspond to the value in the table, if so then remember as candidate for forwarding, and then regardless continue to the next row
  - (a) if there are no candidates, datagram sent to otherwise output interface
  - (b) if there is one candidate, datagram sent to the corresponding output interface
  - (c) if there are multiple candidates (due to multiple routes to some host), the most specific one is taken: **longest prefix rule**

To get a block of addresses, ICANN (Internet Corporation for Assigned Names and Numbers) allocates IP addresses through 5 regional registries (RRs) and manages DNS.

DHCP and NAT help with IPv4 address space exhaustion. IPv6 has 128-bit address space.

#### 2.2.4 Network Address Translation (NAT)

Limiting the number of addresses and re-using IP addresses in a smart way.

Addresses in a private address space are not routable outside it and can be reused as much as desired.

Advantages:

- range of addresses not needed from ISP: just one IP address for all devices
- can change addresses of devices in local network without notifying outside world
- can change ISP without changing addresses of devices in local network
- devices inside local net not explicitly addressable or visible by outside world

Host creates IP datagrams with source and destination IP address that carries 1 transport layer segment. Each segment has source and destination port numbers (Layer 4 entity).

Client knows the port number at the server for the service it needs and OS selects a unique source port number.

Receiving host uses IP addresses and port numbers to direct segment to appropriate process (via sockets).

All devices in local network share just one IPv4 address to the outside world:

- all datagrams leaving local network share same source NAT IP address but different source port numbers
- datagrams with source or destination in local network have subnetted address for source and destination, as usual

The NAT router must transparently:

- for outgoing datagrams, replace source IP address and port number of every outgoing datagram to NAT IP address and new port number

- remote clients and servers will respond using the NAT IP address and new port number as destination address
- remember in the NAT translation table every pair of source IP address and port number to NAT IP address and new port number translation mapping
- for incoming datagrams, replace the NAT IP address and new port number in the destination fields of every incoming datagram with the corresponding source IP address and port number stored in the NAT table

The 16-bit port number field allows 60000 simultaneous connections with a single public IP address.

The router keeps NAT entries in the translation table for a configurable length of time. For TCP connections, default timeout is 24 hours. Since UDP is not connection based, default timeout is 5 minutes.

NAT is controversial since:

- routers should only process up to layer 3
- address shortage should be solved by IPv6
- violates end-to-end argument since port number is manipulated by network-layer device
- can be tricky if client wants to connect to the server behind NAT

### 2.2.5 Middleboxes

Any intermediate box performing functions apart from normal, standard functions of an IP router on the data path between a source host and destination host.

Includes NAT, application-specific, firewalls, intrusion detection systems, load balancers, and caches.

The internet has a thin waist, since there is a single network layer protocol: IP that must be implemented by every Internet-connected device (compared to many protocols in other layers).

Middleboxes give love handles that operate inside the network.



## 2.3 IPv6 Addressing

Possibly not enough 32-bit IPv4 addresses. Also, IPv4 is slow (variable length header). IPv6 allows different network-layer treatment of flows and better mobility management.

IPv6 datagram has:

- IP version
- priority among datagrams in flow
- flow label: identify datagrams in same flow
- payload length
- next header
- hop limit
- 128-bit source address
- 128-bit destination address
- data payload

Compared to IPv4, has no checksum, fragmentation/reassembly, and options.

Not all routers can be upgraded simultaneously so must operate with mixed IPv4 and IPv6.

**Tunneling:** IPv6 datagram is carried as payload in IPv4 datagram among IPv4 routers (packet within a packet)

44.5% of clients access services via IPv6, so takes time to deploy.

## 3 What's Inside a Router

High-level view of generic router architecture has **routing processor** (control plane) and **high-speed switching fabric** (data plane).

Input ports have a physical layer, link layer, and decentralized switching (using header field values, lookup output port using forwarding table in input port memory).

**Destination-based forwarding:** forward based only on destination IP address

**Generalized forwarding:** forward based on any set of header field values

### 3.1 Switching Fabrics

Transfer packets from input link to appropriate output link.

**Switching rate:** rate at which packets can be transferred from inputs to outputs, measured as multiple of input/output line rate

For  $N$  inputs, switching rate of  $N$  times the line rate is desirable.

Major types of switching fabrics are:

- memory
- bus
- interconnection network

#### 3.1.1 Switching via Memory

Used traditionally, with switching under direct control of CPU.

The packet is copied into system memory and speed is limited by memory bandwidth (2 bus crossings per datagram).

#### 3.1.2 Switching via Bus

Datagram from input port memory to output port memory via shared bus.

**Bus contention:** switching speed is limited by bus bandwidth

#### 3.1.3 Interconnection Network

Initially developed to connect processors in multiprocessor.

**Multistage switch:**  $n \times n$  switch from multiple stages of smaller switches

With parallelism:

- fragment datagram into fixed length cells on entry
- switch cells through the fabric and reassemble datagram at exit

Can scale by using multiple switching planes in parallel.

Cisco CRS router:

- basic unit has 8 switching planes
- each plane has a 3 stage interconnection network
- up to 100s of Tbps switching capacity

## 3.2 Port Queuing

### 3.2.1 Input Port Queuing

If switch fabric slower than input ports combined, queuing may occur at input queues. This can lead to queuing delay and loss due to input buffer overflow.

**Head of the Line (HOL) blocking:** queued datagram at front of queue prevents others in queue from moving forward

### 3.2.2 Output Port Queuing

Buffering occurs when arrival rate via switch exceeds output line speed.

Queuing delay (and loss) due to output port buffer congestion (overflow).

**Buffering** is required when datagrams arrive from fabric faster than the link transmission rate. Must have **drop policy** to decide which datagrams to drop if no free buffers.

**Scheduling discipline** chooses among queued datagrams for transmission. With priority scheduling, this decides who gets best performance.

## 4 Control Plane: ICMP, Routing

Even though best effort, IP attempts to avoid errors and report problems when they occur.

IP does introduce errors or ignore all errors.

Errors detected can be:

- corrupted header bits: header checksum
- illegal addresses: routing tables
- routing loop: TTL field
- fragment loss: timeout

## 4.1 ICMP

Internet Control Message Protocol is a separate protocol for errors reporting and information. Required part of IP (just above IP, layer 3.5) and sends error message to original source.

Used by hosts and routers to communicate network-level info like error reporting (unreachable host, network, port, protocol, etc.) and uses echo request/reply (used by ping).

Network-layer above IP so ICMP messages carried in IP datagrams.

ICMP message has a type, code and first 8 bytes of IP datagram causing error.

IP datagram header contains a bit to specify no fragmentation allowed, which can be bit 0 → must be zero, bit 1 → don't fragment, bit 2 → more fragments.

ICMP sends an error message when fragmentation required but not permitted. This is done by probing to find the largest MTU that does not generate an error message. This MTU is not guaranteed if routes change.

For traceroute (provides delay measurement from source to router):

- source sends sets of UDP segments to destination with an unknown port number where each set has an increasing TTL starting from 1
- datagram in set  $n$  arrives to router  $n$  where router discards datagram and sends source ICMP message which possibly includes name of router and IP address
- when ICMP message arrives at source, it records RTTs

The stopping criteria for traceroute is that the UDP segment eventually arrives at destination host. The destination returns ICMP “port unreachable” message and so the source stops.

## 4.2 Routing

Goal: determine good paths/routes from sending host to receiving hosts, through network of routers

**Path**: sequence of routers packets traverse from given initial source host to final destination host

A good path can have least cost, fastest, least congested, or other criteria.

**Broadcast routing:** route packet from a source to all nodes in the network

**Flooding:** each node sends packet on all outgoing links and discard packets received a second time

**Spanning Tree routing:** send packet along a tree that includes all nodes in the network

Internet too big for flat routing:

- can't store all destinations in routing tables
- each network within the Internet may want to control routing in its own network

Approach: aggregate routers into regions known as autonomous systems (AS) aka domains

**Intra-AS:** routing within same AS

- all routers in AS must run same intra-domain protocol
- routers in different AS can run different intra-domain routing protocols
- **gateway router:** at edge of AS, has link to routers in other ASs

**Inter-AS:** routing among ASs

- gateways perform inter-domain routing (and intra-domain routing)

Forwarding table configured by intra-AS and inter-AS routing algorithms:

- intra-AS routing determine entries for destinations within AS
- inter-AS and intra-AS determine entries for external destinations

Intra-AS routing protocols are **Interior Gateway Protocols (IGP):**

- RIP: Routing Information Protocol
- OSPF: Open Shortest Path First
- IGRP: Interior Gateway Routing Protocol (Cisco proprietary)

Only Inter-AS routing protocol is BGP (Border Gateway Protocol).

**Routing algorithm:** algorithm that finds least-cost path, used to fill routing tables

Routing algorithms can be classified by information and static/dynamic.

Some global info vs more local info:

- global
  - all routers flood the network with info about link state so a complete map can be created
  - routes then computed
  - link state algorithms
    - \* nodes send info about their distances from their neighbours to all other nodes
- local
  - each router initially knows physically connected neighbours, exchanges part of its routing table with its neighbours
  - iterative process of computation
  - distance vector algorithms
    - \* nodes send info about their distances from all other nodes to their neighbours

Static vs dynamic:

- static: routes change slowly over time
- dynamic: routes change more quickly, periodic update in response to link cost changes

#### 4.2.1 Intra-AS Routing

All routers identical, network considered flat.

Cost between links defined by network operator, and costs should be additive.

Most routing algorithms use a shortest path algorithm, such as Bellman-Form and Dijkstra's.

Can define cost as:

- fixed quantity based on link rate, propagation delay, or some combination
- time-varying quantity based on average traffic on a link, buffer occupancy, delay, and error conditions

- dangerous since oscillations possible

Dijkstra's link-state routing algorithm is

- centralized: network topology, reachability, link costs known to all nodes
  - accomplished via link state broadcast
  - all nodes have same info
- each router uses Dijkstra's algorithm to compute least cost paths from source to all other nodes
  - gives forward table for that node
- $D(v)$  is current estimate of cost of least-cost path from source to destination  $v$
- $p(v)$  is predecessor node along path from source to  $v$
- $N'$  is set of nodes whose least-cost path definitively known

Open Shortest Path First (OSPF) routing:

- open means publicly available
- state of a link is a description of interface on router and its relationship to its neighbouring routers
- description of interface would include IP address of interface, mask, subnet, type of network it is connected to, routers connected to that network, and more
- collection of link-states forms link-state DB
- each router floods OSPF link-state advertisements (over IP) to all other routers in entire AS
- multiple link costs metrics possible link bandwidth and delay
- each router has full topology (subnets) and uses Dijkstra's algorithm to compute forwarding table
- security: all OSPF messages authenticated (to prevent malicious intrusion)
- link state packet (LSP) contains
  - id of node that created LSP

- cost of link to each directly connected neighbour (typically link rate)
- sequence number
- TTL for this packet
- advertisements disseminated to entire AS (via flooding, at least once every 30 min or when changes)
- send new LSP with infinite cost to signal a link down
- reliable flooding
  - generate new LSP periodically, increment sequence number
  - store most recent LSP from each node, forward LSP to all nodes but one that sent it, throw all new copies of same LSP (do not forward)
  - decrement TTL of each stored LSP regularly, discard when TTL=0

**Two-level hierarchy:** local area and backbone

- link-state advertisements flooded only in area, or backbone
- each node has detailed area topology, only knows direction to reach other destinations

**Area border routers:** summarize distances to destinations in own area, advertise in backbone

**Boundary router:** connects to other ASs

**Local router:** flood LS in area only, compute routing within area, forward packets to outside via area border router

**Backbone router:** runs OSPF limited to backbone

With distance-vector algorithm, a router never builds a complete map of the network, just neighbours.

- a node computes distance vector whenever it receives something from one of its neighbours
- a node sends distance to neighbours whenever it changes

Distributed Bellman Ford used in Distance Vector algorithm:



- from time-to-time, each node sends distance vector estimate to neighbours
- asynchronous and iterative
  - each local iteration caused by local link cost change or distance vector update message from neighbours
- whenever a node receives a new distance vector estimate from a neighbour, it updates its own distance vector using the Bellman Ford equation
- under minor, natural conditions, estimate converges to actual least cost
- each node notifies neighbours only when distance vector changes
  - neighbours notify neighbours only if necessary

With distance vector algorithm, iterative communication, so computation steps diffuse information through network. Good news travels fast, bad news travels slow.

**Poisoned reverse:** if  $z$  routes through  $y$  to get to  $x$ ,  $z$  tells  $y$  its distance to  $x$  is infinite (so  $y$  won't route to  $x$  via  $z$ )

#### **Routing Information Protocol:**

- included in BSD-UNIX distribution in 1982
- distance vector algorithm
  - distance metric is number of hops (max 15), each link has cost 1
  - distance vectors exchanged with neighbours every 30 secs in response message (aka advertisement)
  - each advertisement is a list of up to 25 destination subnets as well as sender's distance to each subnet
- if no advertisement heard after 180 secs, neighbour/link declared dead
  - routes via neighbour invalidated
  - new advertisements sent to neighbours
  - neighbours in turn send out new advertisements (if tables changed)
  - link failure info quickly propagates to entire net
  - poison reverse used to prevent ping-pong loops

Problems with distance vector routing:

- convergence is slow
- loops can be formed, due to routing table inconsistency: packets being forwarded from router to router and never reach destination
- loops might last until convergence or count to infinity problem

OSPF is better than RIP in stability (robustness to network changes) since distance vector converges slowly compared to link state. But, OSPF creates more routing traffic (flooding).

#### 4.2.2 Inter-AS Routing: BGP

De facto inter-domain routing protocol. Allows subnet to advertise existence and the destinations it can reach, to rest of Internet.

BGP provides each AS a means to:

- eBGP: obtain subnet reachability info from neighbouring ASs
- iBGP: propagate reachability info to all AS-internal routers
- determine good routes to other networks based on reachability info and policy

**BGP session:** two BGP routers exchange BGP messages over semi-permanent TCP connection

- advertising paths to different destination network prefixes (BGP is a path vector protocol)
  - advertised route has prefix (destination being advertised) and AS-PATH (list of ASs through which prefix advertisement has passed) and NEXT-HOP (indicates specific internal-AS router to next-hop AS)

**Policy-based routing:**

- gateway receiving route advertisement uses import policy to accept/decline path
- AS policy also determines whether to advertise path to other neighbouring ASs

Gateway router may learn about multiple paths to destination, selects route based on:

- local preference value attribute: policy decision
- shortest AS-PATH
- additional criteria

**BGP messages** are exchanged between peers over TCP connection:

- OPEN: opens TCP connection to remote BGP peer and authenticates sending BGP peer
- UPDATE: advertises new path (or withdraws old)
- KEEPALIVE: keeps connection alive in absence of UPDATEs, also ACKs OPEN request
- NOTIFICATION: reports errors in previous message, also used to close connection

Policy:

- inter-AS: admin wants control over how its traffic is routed, who routes through its network
- intra-AS: single admin, so policy less of an issue

Scale:

- hierarchical routing saves table size, reduced update traffic

Performance:

- inter-AS: policy dominates over performances
- intra-AS: can focus on performance

**Hot Potato Routing**: choose local gateway that has least intra-domain cost

ISP only wants to route traffic to/from its customer networks (not carry traffic between other ISPs).

### 4.3 Software Defined Networking

Internet network layer: historically implemented via distributed, per-router control approach:

- monolithic router contain switching hardware, runs proprietary implementation of Internet standard protocols (IP, RIP, IS-IS, OSPF, BGP) in proprietary router OS
- different middleboxes for different network layer functions: firewalls, load balancers, NAT boxes

In per-router control plane, individual routing algorithm components in each and every router interact in the control plane to compute forwarding tables.

In SDN control plane, remote controller computes and installs forwarding tables in routers.

SDN is logically centralized since easier network management (avoid router misconfigurations and greater flexibility of traffic flows), table-based forwarding allows programming routers (centralized programming easier than distributed programming), and open (non-proprietary) implementation of control plane (foster innovation).

With traditional routing:

- link weights are control knobs, so not much control.
- load balancing not possible
- can't route along different paths from same node

Each router contains a forwarding table: *match plus action* abstraction, so match bits in arriving packet then take action

- destination-based forwarding: forward based on destination IP address
- generalized forwarding: many header fields can determine action and many actions possible

**Flow**: defined by header field values (in link, network, and transport-layer fields)

**Generalized forwarding**: simple packet-handling rules

- match: pattern values in packet header fields
- actions: for matched packet  $\rightarrow$  drop, forward, modify matched packet or send matched packet to controller

Orchestrated tables can create network-wide behaviour.

Intelligence for 20th century phone net was at network switches, for early Internet at edge computing, and for modern Internet at programmable network devices and application-level infrastructure at edge.