

CMSC417 Computer Networks

Michael Li

Contents

| | | |
|----------|---|----------|
| 1 | Foundations | 2 |
| 1.1 | Requirements | 2 |
| 1.2 | Network Architecture | 3 |
| 1.3 | Implementing Network Software | 4 |
| 1.4 | Network Performance | 5 |
| 2 | Routing | 5 |
| 2.1 | Vocabulary | 5 |
| 2.2 | Network as a Graph | 5 |
| 2.2.1 | Distance Vector | 6 |
| 2.2.2 | Distance Vector Notations | 6 |
| 2.2.3 | Distance Vector Protocol | 7 |
| 2.2.4 | Link State | 7 |
| 2.2.5 | Comparison Distance Vector and Link-State Routing | 8 |
| 2.3 | Metrics | 8 |
| 3 | Network Layer | 9 |
| 3.1 | Internet Protocol | 9 |
| 3.1.1 | Vocabulary | 9 |
| 3.1.2 | Service Model | 9 |
| 3.1.3 | Packet Format | 9 |
| 3.1.4 | Fragmentation and Reassembly | 10 |
| 3.2 | Global Addresses | 10 |
| 3.3 | Datagram Forwarding in IP | 10 |
| 3.4 | Subnetting and Classless Addressing | 11 |
| 3.4.1 | Subnetting | 11 |
| 3.4.2 | Classless Addressing (CIDR) | 11 |

1 Foundations

1.1 Requirements

Requirement 1: Scalable Connectivity: Network should be scalable and should be able to connect to a lot of devices. At the lowest level, **nodes** (computers) are connected by **links** (e.g. wire).

- **Point-to-point:** link is only shared between 1 pair of nodes
- **Multiple-access:** link is shared between several nodes

Indirect connectivity: can be achieved in multiple ways:

- **Switched network:** uses nodes that forward data received from one link and out the other link.
 - **Circuit switched:** establishes a stream of bits
 - **Packet switched:** nodes send blocks of data in the form of **packets** and use **store-and-forward** where the node has to receive entire packet, storing it in memory, before forwarding the entire packet. Here the **switches** implement the network and the **hosts** use the network
- **Internetwork (internet):** set of independent networks connected by **routers** that are nodes that act like switches between networks.

In order to provide **host-to-host** connectivity, each node is assigned an **address** that allows other nodes to communicate with the node. Switches and routers use this address to **route** messages to their destination

Requirement 2: Efficiency

Multiplexing: system resources shared among multiple users. Achieved by using

- **Synchronous Time Division Multiplexing (STDM):** divide time into equally sized chunks and round-robin each flow
- **Frequency Division Multiplexing (FDM):** transmit each flow at different frequencies
- STDM and FDM both have issues in that links may be idle during a given time interval or at a particular frequency. Also, both do not support adding new flows easily.
- **Statistical Multiplexing:** links are shared on a time-basis but is on demand, rather than a predetermined time slot. Also limits the size of packets sent over to ensure data is not lost.

With multiplexing, switches may need to buffer packets if the input rate $>$ output rate. This can lead to **congestion** where the buffer is full so it needs to drop packets

Requirement 3: Support for Common Services: networks provide **logical channels** that application processes use to communicate with each other. These channels provide a set of services and are typically implemented on end hosts to keep links simple

Requirement 4: Reliability: need to consider 3 types of errors

- **bit error:** bit is flipped in the physical link
- **packet loss:** packet is dropped or uncorrectable packet is discarded
- **physical crash:** e.g. network shutdown

Requirement 5: Managability: system needs to be scalable (should be able to support more devices)

1.2 Network Architecture

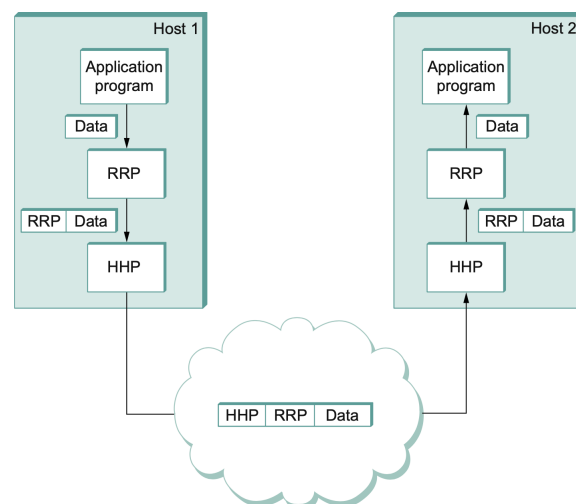
Network Architecture is a blueprint that guides the design and implementation of networks. Uses **layering** and **abstraction** to capture the important aspects of a system. Services at higher **protocols** (layers) are built using the services of the lower protocols. These protocols define 2 different interfaces:

- **service interface**: defines operations local objects can perform on the protocol (e.g. sending and receiving messages)
- **peer interface**: defines the form and meaning of messages exchanged between protocol peers

Peer-to-peer interaction is indirect; each protocol communicates with its peer by passing messages to lower level protocols that deliver the messages to its peer. This ends up creating a **protocol stack**. On the peer side, protocol layers are popped until we reach the original target.

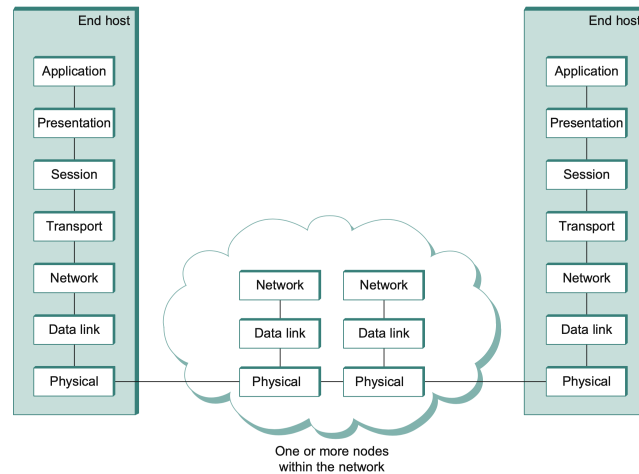
Encapsulation is performed at each protocol level; **headers** are attached to the **payload** (rest of the message) and instruct the peer on how to handle the message (i.e. how the payload is passed to the next higher level in the peer).

A **demux key** is attached to the header, allowing protocols to demultiplex messages to the correct application on the destination host



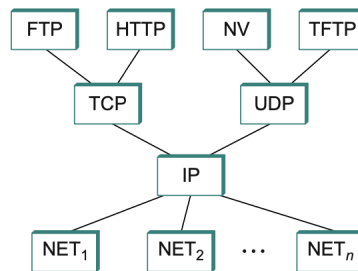
Open Systems Interconnection (OSI): partitions network functionality into 7 layers

- **Physical layer**: handles transmission of raw bits over communication link
- **Data link layer**: aggregates bits into **frames**
- **Network layer**: routes packets to nodes within packet-switched network. Unit of data exchanged is called **packet**
- **Transport layer**: implements process-to-process channels that exchange messages. Unit of data here is called **messages**
- **Session layer**: ties potentially different transportation streams
- **Presentation layer**: concerned with the format of data being exchanged
- **Application layer**: handles high level protocols



Internet Architecture: 4 levels:

- **Network protocols:** lowest level implemented by a combination of hardware and software
- **Internet protocol (IP):** supports multiple networking technologies into 1 internetwork
- **Transmission control protocol (TCP) and User Datagram Protocol (UDP):** TCP provides reliable byte-stream channel and UDP provides unreliable datagram delivery channel
- **Application protocols:** enable interoperation of applications
- One main difference between Internet Architecture and OSI is that layering isn't prevalent for the former; apps can skip layers if desired



1.3 Implementing Network Software

Socket: the point where application attaches to the network. The interface defines how to create sockets, how to attach them to network, how to send and receive messages through sockets, and how to close sockets. To create a socket in TCP:

- invoke `socket()` to create socket
- for the server side, perform a **passive open** (prepared to accept connections but doesn't actually establish them) by invoking `bind()`, `listen()`, and `accept()`
- for the client side, perform an **active open** (says who it wants to communicate with) by invoking `connect()`. Once connected, can invoke `send()` and `recv()` to send and receive messages from specified socket to a buffer

```

// domain: protocol family
// type: semantics of communication
// protocol: protocol used
// returns a handle (identifier for socket)
int socket(int domain, int type, int protocol);

// binds created socket to a specific address
int bind(int socket, struct sockaddr *address, int addr len)
// listen defines how many connections can be pending on a socket
int listen(int socket, int backlog)
// accept carries out passive open
int accept(int socket, struct sockaddr *address, int *addr len)

// connect doesn't return until TCP has established a connection
int connect(int socket, struct sockaddr *address, int addr len)
// once connection is established, applications can send and receive messages
int send(int socket, char *message, int msg len, int flags)
int recv(int socket, char *buffer, int buf len, int flags)

```

1.4 Network Performance

Performance measured by **bandwidth** (number of bits transmitted over time) and **latency** (how long it takes for message to travel)

Bandwidth can be broken down into a sum of

- speed of light propagation (distance/S)
- time to transmit a unit of data (size/bandwidth)
- queuing delays

Delay × bandwidth product: maximum number of bits that can travel in a pipe at once. Tells how many bits must be set before first bit arrives at receiver

2 Routing

2.1 Vocabulary

Forwarding (Data Plane): receiving a packet, looking up its destination on a table, and sending the packet in the appropriate direction.

Routing (Control Plane): process that builds **forwarding tables**.

Forwarding table: contains information regarding mappings from network prefix to an outgoing interface to aid in forwarding packet.

Routing table: contains mappings of network prefixes to next hops.

Routing domain: internetwork where all routers are under the same administrative control.

2.2 Network as a Graph

Entire network can be represented as a graph (nodes, network links as edges with cost). The goal of **routing** is to find the lowest-cost path between any 2 nodes. Algorithms that calculate the shortest path usually do so in a dynamic, distributed manner, but this can lead to some confusion when different paths thinking that they have the shortest path available. There are 2 main routing protocols: **Distance vector** and **Link state**.

2.2.1 Distance Vector

Each node constructs an array (vector) that contains the cost to all other nodes and distributes this vector to its immediate neighbor and the vector is updated based on which nodes can be reached by the neighbor node. Links that are down are assigned an infinite cost.

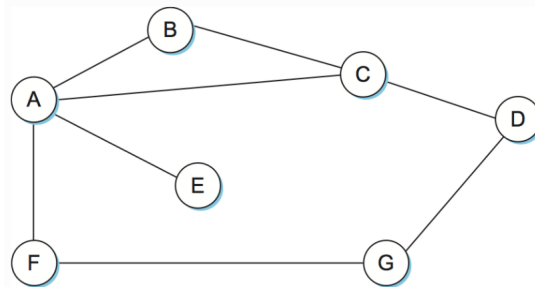


Figure 85. Distance-vector routing: an example network.

Each routing table contains entries (destination, cost, next hop). Routers exchange information (destination, cost) with immediate neighbors and the table is built from a sequence of exchanges. **Note:** nodes propagate a global view to its neighbors.

End result is that each node has the same routing table so routing information is consistent (**convergent**). Note that no one node contains all information in this table, each node only knows about the content of its own routing table. There are 2 ways in which updates occur:

- **Periodic update:** each node automatically sends an update message, even if nothing has changed, to let other nodes know that it is still running and to allow other nodes check if the current route becomes unviable.
- **Triggered update:** occurs when a node notices a link failure or receives an update from one of its neighbors that causes 1 of the routes in the routing table to change. This then triggers an update to the node's other neighbors.

Count to infinity problem: example - from the graph above, if the link A to E goes down, A notifies its neighbors that its distance to E is infinity, but B and C claim that they have a distance of 2 away from E. This then leads to the conclusion that A is 3 away from E, and then B and C is 4 away, leading to a distance of infinity. Potential solutions include:

- **Set a maximum** (e.g. 16) to represent infinity. Once we reach that many hops, we claim that the distance is infinity. This solution is restrictive in that it isn't scalable.
- **Split horizon:** when a node sends updates, it does not send routes it learned from a specific neighbor to that same neighbor. However, this solution isn't scalable for loops involving more than 2 nodes.
- **Poisoned reverse:** report failed links as ∞ which helps with routing loops involving 2 nodes. Doesn't work well loops involving ≥ 2 nodes.

Routing Information Protocol (RIP): canonical example of routing protocol built on distance vector algorithm

2.2.2 Distance Vector Notations

$x.C(y)$ link cost from node x to node y . If they are not neighbors, value is ∞

$x.D(z \rightarrow y)$ distance known to x to travel from z , a neighbor of x , to y

$x.N$ set of neighbors to x

\mathbb{N} set of all nodes

$x \rightarrow y : M$ node x sends a message M to y

At the start, routing table consists of values:

- if $y \in x.N \implies x.D(y \rightarrow y) = x.C(y)$
- if $y \notin x.N \implies \forall v \in x.N, x.D(v \rightarrow y) = \infty$

Entries consist of:

- destination
- cost
- next hop

At the end, routing table consists of values:

$$\bullet \underbrace{x.D(z \rightarrow y)}_{z \in x.N} = x.C(z) + \underbrace{\min_{v \in z.N} (z.D(v \rightarrow y))}$$

2.2.3 Distance Vector Protocol

x sends an update message, $M = \forall d \in \mathbb{N}, (d, \underbrace{\min_{v \in x.N} (x.D(v \rightarrow d))})$, to all $v \in x.N$

When x receives update from neighbor v :

- Did $x.C(v)$ change? If yes, $\forall y \in \mathbb{N}, x.D(v \rightarrow y) = \underbrace{x.D(v \rightarrow y) + \delta}_{\text{old val} + \text{change}}$
 - Did the lowest cost to y change? If yes, update the routing table.
- Did $\underbrace{\min_{z \in v.N} (v.D(z \rightarrow y))}_{\text{best distance from } v \rightarrow y}$ change?
 - If yes, update $x.D(v \rightarrow y) = \underbrace{x.C(v) + \delta}_{\text{link cost} + \text{new change}}$
 - Did the lowest cost to y change? If yes, update the routing table.

2.2.4 Link State

Since each node knows how to reach its neighbors, we can make it so that this information is known to each node to create a complete map of the network at each node. **Note:** nodes propagate local views to all nodes so each has enough information to compute routes locally and can detect routing loops from failures. Relies on

- **Reliable flooding:** making sure that all nodes get the same copy of the link-state information from all the other nodes. Nodes send their link-state information (**link-state packet (LSP)**) to all of their neighbors who then forward to their neighbors until all nodes have been reached. LSPs contain
 - ID of node that created the LSP
 - list of connected neighbors of that node with the link cost to each one
 - sequence number
 - time to live for this packet

First 2 items enable route calculation. Last 2 items make flooding more reliable (ensuring that each node has the most recent copy of the information since there might be multiple LSPs from different nodes traversing the network)

Nodes receive **ACK** (acknowledgement) packets when they send an LSP. The LSP is retransmitted if no ACK is received.

How flooding works: consider a node X that received a copy of an LSP that originated from node Y. If X already has a stored copy of an LSP from Y. If not, X stores the LSP. Otherwise X compares the new and old LSP sequence numbers. If the new LSP has a larger sequence (more recent), then it replaces the old LSP and X sends a copy of this LSP to all of its neighbors, except the neighbor that sent the new LSP, who do the same check.

LSPs are generated either on a **periodic** basis or **topology change** (a failed link is detected), similar to that of Distance Vectors.

Optimizations for LSPs:

- set long periodic timers for updates (reduces overhead for creating LSPs unless absolutely necessary)
 - Each time a node generates a new LSP, it increments the sequence by 1. If the node crashes, it starts at 0 when it comes back up. This ensures that newer information replaces older information
 - LSPs carry a time to live (ensures that old information is eventually removed). A node always decrements the time to live of a newly received LSP before flooding its neighbors. TTL is also decremented while a node holds the LSP. This ensures that when a node comes back online, there will be a point where we can use its LSP (all other LSPs with the same ID will have died at some point).
- **Route Calculation:** once a node has a copy of the LSP from every other node, it can compute a graph of the entire network and can determine the best route, using Dijkstra's.

Link-state routing stabilizes quickly and responds rapidly to topology changes. However, each node ends up being very large with all of the LSPs it must hold.

2.2.5 Comparison Distance Vector and Link-State Routing

Both are distributed routing algorithms but in distance vector, each node only communicates with its direct neighbor and tells the neighbor everything it knows. In link-state, each node communicates with all other nodes but only tells them what it knows for sure.

2.3 Metrics

Measuring cost of link was done in a variety of ways

- Number of packets queued for waiting at each link (larger queue is assigned a larger cost weight). Not a good measurement since this focuses on queue size rather than distance and time.
- Consider both bandwidth and latency. Each packet is given a timestamp ArrivalTime and DepartTime and Delay is calculated as $(\text{DepartTime} - \text{ArrivalTime}) + \text{TransmissionTime} + \text{Latency}$, at the other node. Not a good measurement since links will become idle when they have a high load.
- Compress dynamic range of metric to smooth the variation of the metric over time. Smoothing was achieved by limiting how much the metric could change from 1 measurement cycle to the next and by calculating the average amount of utilization for a link.

3 Network Layer

3.1 Internet Protocol

3.1.1 Vocabulary

Internetwork: arbitrary collection of networks interconnected to provide a host-to-host packet delivery service

Router: node that interconnects networks

Service Model: host-to-host service one wants to provide

Datagram: packet sent in a connectionless manner over a network **Maximum Transmission Unit (MTU):** largest IP datagram it can carry in a frame (usually has to be smaller than max packet size so that it fits in the payload of the link-layer frame)

3.1.2 Service Model

In order to provide host-to-host service, the underlying physical networks must be able to support it. Service model can be split into 2 parts

- **Addressing scheme:** provides a way to identify all hosts in the internetwork
- **Datagram (connectionless) model:** makes **best effort** to deliver datagrams (ie does not guarantee datagrams are sent successfully). This is sometimes called **unreliable** service.

3.1.3 Packet Format

Similar to regular packet where it consists of **headers** and a data **payload**. The headers is split into 32-bit boundaries, for higher network layers to use.

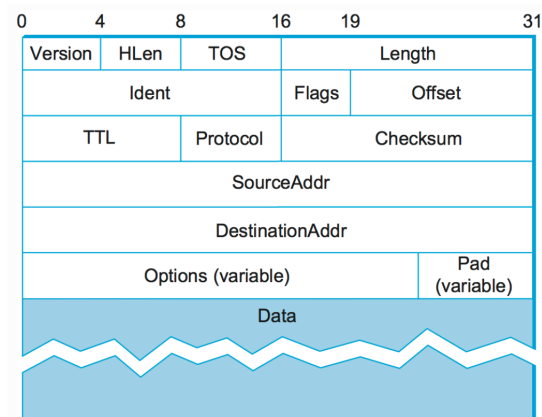


Figure 72. IPv4 packet header.

Looking at the various IP header fields:

- **Version:** specifies which version of IP is being used.
- **Hlen:** specifies the length of the header in 32-bit words.
- **TOS:** type of service. Allows packets to be treated differently based on application needs.
- **Length:** size of entire packet (including header) as bytes.
- Note that the second word is used for fragmentation and reassembly (explained below)
- **TTL:** time to live (used to discard packets that have been stuck in the routing loop)
- **Protocol:** demultiplexing key used to identify which higher protocol to send to (UDP, TCP)

- **Checksum:** used to check for bit corruption in the header
- **SourceAddr:** which address packet came from (allows recipients decide if they want to accept)
- **DestinationAddr:** which address packet is sent to
- **Options:** additional options may be added (can be determined using **HLen**)

3.1.4 Fragmentation and Reassembly

Biggest issue is that each network technology has its own idea of how large a packet is. Thus, packets must be fragmented and reassembled when they don't fit in the network technology.

Fragmentation occurs when a router receives a datagram that it wants to forward over the network with a lower MTU. The receiving host will recognize which fragments belong together through its **Ident** header field. If all fragments don't arrive, the host discards the packet entirely and does not attempt recover missing fragments.

The **Flag** header field is used to indicate that there are more packets coming. The **Offset** header field is used to identify which partition of the packet is being sent over.

3.2 Global Addresses

IP addresses are **hierarchical**, consisting of a **network** part (identifies the network thae host is attached to) and a **host** part (unique host on the network). Addresses are written as 4 octets, or **dotted quad** notation and can be broken down into 5 classes subnets, in the classical model:

- A: $0 - 127.y.x.w$ 0, $\underbrace{7b}_{\text{network}}, \underbrace{24b}_{\text{host}}$
- B: $128 - 191.y.x.w$ 10, $\underbrace{14b}_{\text{network}}, \underbrace{16b}_{\text{host}}$
- C: $192 - 223.y.x.w$ 110, $\underbrace{21b}_{\text{network}}, \underbrace{8b}_{\text{host}}$
- D: $224 - 239.y.x.w$
- D: $240 - 255.y.x.w$

Classical representation runs out of space rapidly so we can instead break larger subnets into smaller subnets using a **network mask** (e.g. 255.255.0.0). Then:

$$\text{mask \& address} = 128.8.0.0 \text{ base address}$$

3.3 Datagram Forwarding in IP

Important things to note:

- Every IP datagrams contains the IP address of the destination host
- Network part of the IP address uniquely identifies a physical network that is connected to other networks

Datagrams are sent from source host to destination host, passing through a series of routers. Nodes first try to establish if they are in the same physical network as the destination (comparing the network part of the address). If they are not in the same network, the datagram is sent to the **next hop router** that consults a forwarding table.

3.4 Subnetting and Classless Addressing

3.4.1 Subnetting

Reduces the total number of network numbers assigned by taking an IP network number and allocating it to several **subnets**. This is done using a **network mask**.

3.4.2 Classless Addressing (CIDR)

Coalesces several classes of addresses into a **supernet**. This is more efficient than subnetting by giving networks address space as needed (don't over allocate). CIDR blocks define a **length**

128.8.13.0/24

Additional notes:

- IP address are /32 CIDR blocks
- Number of address in a /N block is 2^{32-N}
- Base address (e.g. 128.3.0.0 are usually **gateways**
- Highest address (e.g. 128.3.255.255) are called **broadcast addresses**