# A Quick Review of Computer Networking
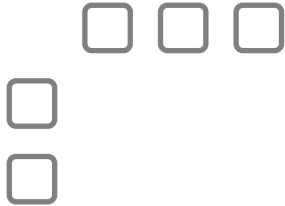
## Architecture, Applications, Transport (TCP), Routing

Antonio Carzaniga
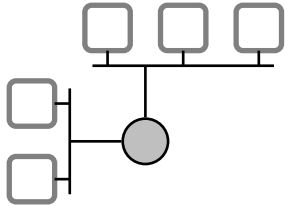
Faculty of Informatics
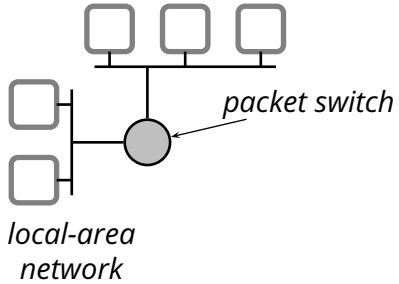Università della Svizzera italiana

February 22, 2021

*packet switch*

*local-area
network*

- The Internet uses *packet switching*

- The Internet uses *packet switching*

- *Packet switch:* a *link-layer switch* or a *router*

- The Internet uses *packet switching*

- *Packet switch:* a *link-layer switch* or a *router*

- *Communication link:* a connection between packet switches and/or end systems

- The Internet uses *packet switching*

- *Packet switch:* a *link-layer switch* or a *router*

- *Communication link:* a connection between packet switches and/or end systems

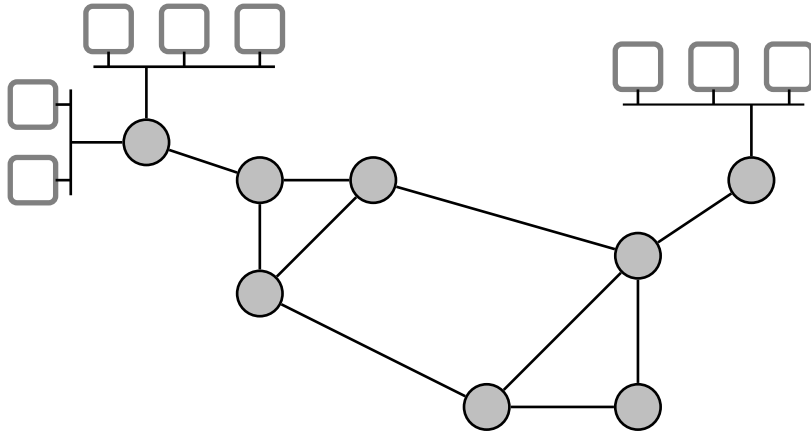- *Route:* sequence of switches that a packet goes through (a.k.a. *path*)

- The Internet uses *packet switching*

- *Packet switch:* a *link-layer switch* or a *router*

- *Communication link:* a connection between packet switches and/or end systems

- *Route:* sequence of switches that a packet goes through (a.k.a. *path*)

- *Protocol:* control the sending and receiving of information to and from end systems and packet switches
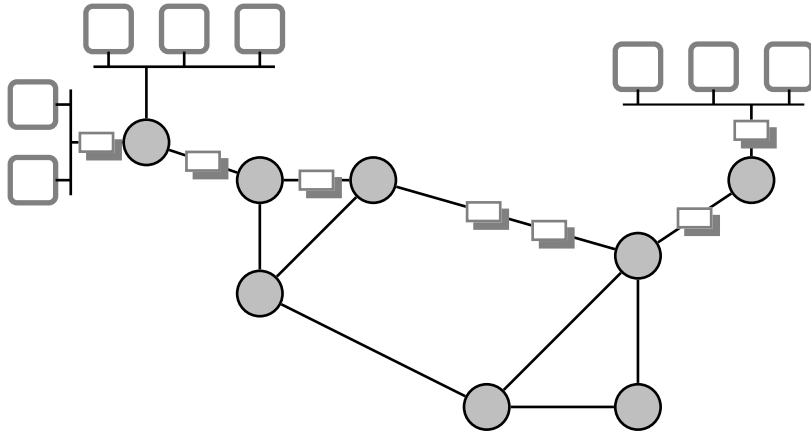
- Various types and forms of medium

- Various types and forms of medium
  - ▶ Fiber-optic cable
  - ▶ Twisted-pair copper wire
  - ▶ Coaxial cable
  - ▶ Wireless local-area links (e.g., 802.11, Bluetooth)
  - ▶ Satellite channel
  - ▶ …

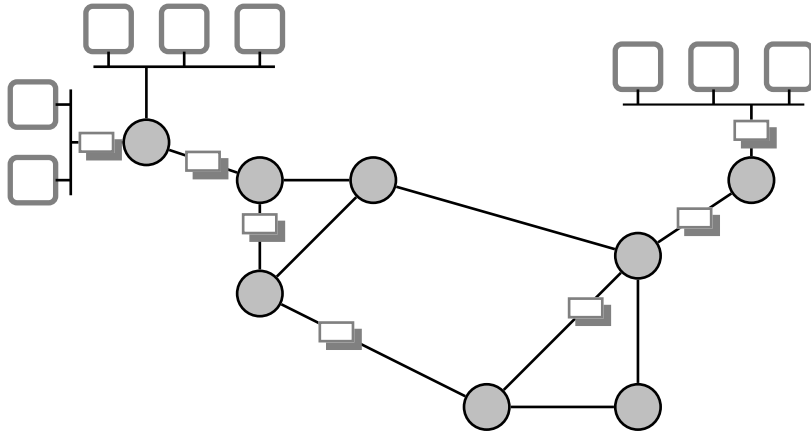- The Internet is a *packet-switched* network

- The Internet is a *packet-switched* network

- Information is transmitted in *packets*

- The Internet is a *packet-switched* network

- Information is transmitted in *packets*

- Switches operate on individual packets

- The Internet is a *packet-switched* network

- Information is transmitted in *packets*

- Switches operate on individual packets

- A switch (router) receives packets and *forwards* them along to other switches or to end systems

- The Internet is a *packet-switched* network
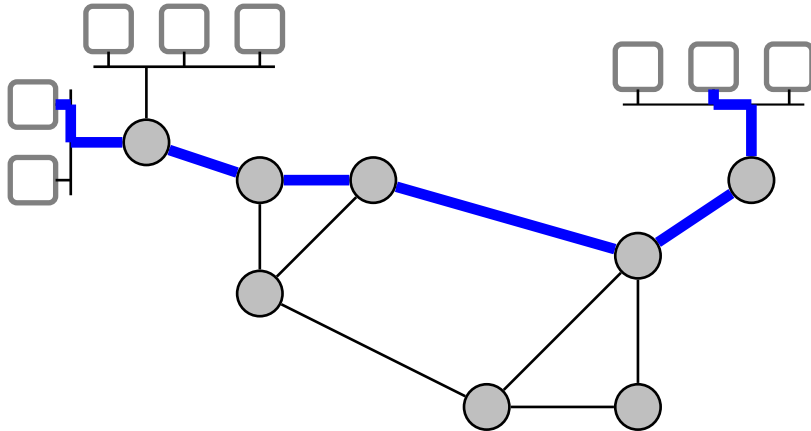
- Information is transmitted in *packets*

- Switches operate on individual packets

- A switch (router) receives packets and *forwards* them along to other switches or to end systems

- Every forwarding decision is taken on the basis of the information contained in the packet

- The telephone network is a typical circuit-switched network
  - ▶ not any more, really, but still…

- The telephone network is a typical circuit-switched network
  - ▶ not any more, really, but still…

- Communication requires a *connection setup* phase in which the network reserves all the necessary resources for that connection (links, buffers, switches, etc.)

- The telephone network is a typical circuit-switched network
  - ▶ not any more, really, but still…

- Communication requires a *connection setup* phase in which the network reserves all the necessary resources for that connection (links, buffers, switches, etc.)

- After a successful setup, the communicating systems are connected by *a set of links dedicated to the connection* for the entire duration of their conversation

- The telephone network is a typical circuit-switched network
  - ▶ not any more, really, but still…

- Communication requires a *connection setup* phase in which the network reserves all the necessary resources for that connection (links, buffers, switches, etc.)

- After a successful setup, the communicating systems are connected by *a set of links dedicated to the connection* for the entire duration of their conversation

- When the conversation ends, the network tears down the connection, freeing the corresponding resources (links, buffers, etc.) for other connections

- Circuit switching requires an expensive setup phase
  - however, once the connection is established, little or no processing is required

- Circuit switching requires an expensive setup phase
  - however, once the connection is established, little or no processing is required

- Packet switching does not incur any setup cost
  - however, it always incurs a significant processing and space overhead, on a per-packet basis
    - *processing cost* for forwarding
    - *space overhead* because every packet must be self-contained

■ Circuit switching admits a straightforward implementation of quality-of-service guarantees

  ▶ network resources are reserved at connection setup time

- Circuit switching admits a straightforward implementation of quality-of-service guarantees
  - ▶ network resources are reserved at connection setup time

- Guaranteeing any quality of service with packet switching is very difficult
  - ▶ no concept of a "connection"
  - ▶ and again, processing, space overhead, etc.

- Circuit switching allows only a limited sharing of communication resources
  - once a connection is established, the resources are blocked even though there might be long silence periods
  - i.e., circuit switching is an inefficient way to use the network

- Circuit switching allows only a limited sharing of communication resources
  - once a connection is established, the resources are blocked even though there might be long silence periods
  - i.e., circuit switching is an inefficient way to use the network

- Packet switching achieves a much better utilization of network resources
  - it is designed specifically to share links
  - it is designed specifically to share links

- Idea: combine the advantages of circuit switching and packet switching

- Idea: combine the advantages of circuit switching and packet switching

- There is a connection setup phase

- Idea: combine the advantages of circuit switching and packet switching

- There is a connection setup phase

- The connection does not create a physical circuit, but rather a "virtual circuit"
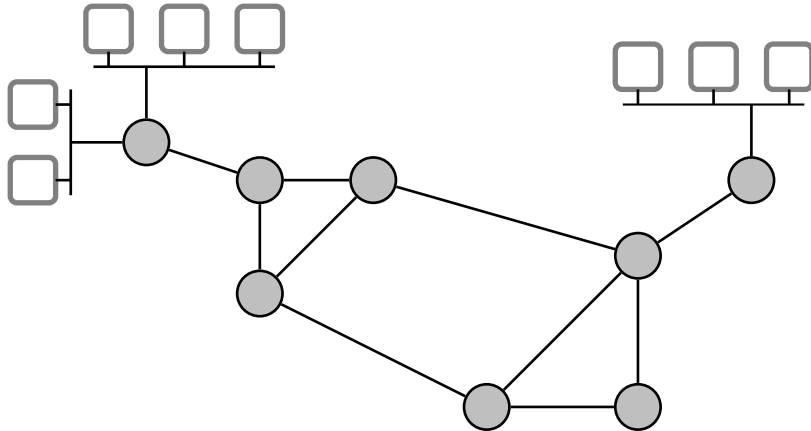
- Idea: combine the advantages of circuit switching and packet switching

- There is a connection setup phase

- The connection does not create a physical circuit, but rather a "virtual circuit"

- Information is sent in packets, so links can be shared more effectively

- Idea: combine the advantages of circuit switching and packet switching

- There is a connection setup phase

- The connection does not create a physical circuit, but rather a "virtual circuit"

- Information is sent in packets, so links can be shared more effectively

- Packets carry a *virtual circuit identifier* instead of the destination address
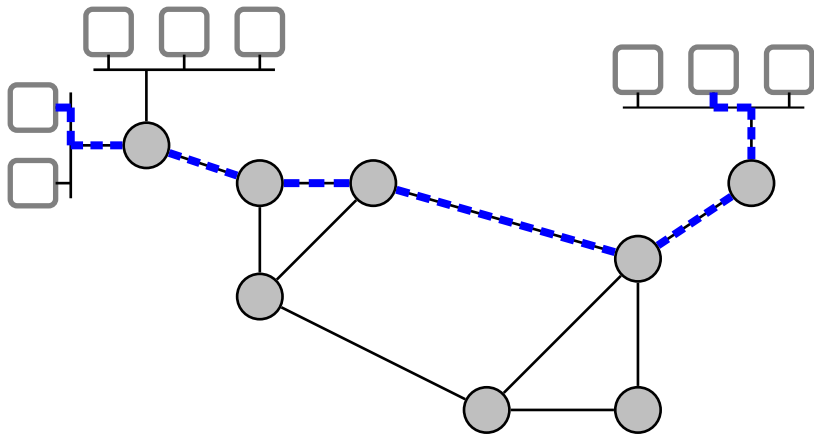
- Idea: combine the advantages of circuit switching and packet switching

- There is a connection setup phase

- The connection does not create a physical circuit, but rather a "virtual circuit"

- Information is sent in packets, so links can be shared more effectively

- Packets carry a *virtual circuit identifier* instead of the destination address
  - ▶ *Important observation:* at any given time there are much fewer *connections* than *destinations*
    - ▶ much faster per-packet processing (forwarding)
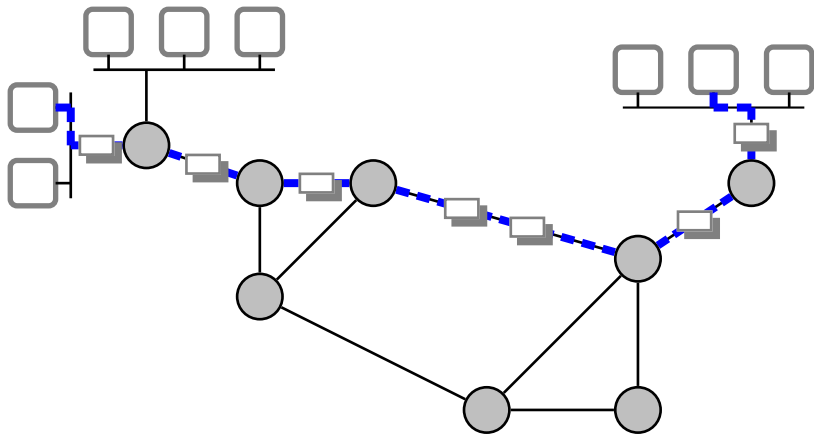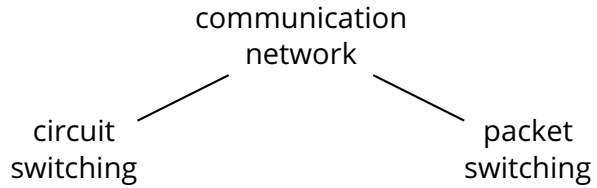    - ▶ lower per-packet space overhead

communication
network

*packet switch*

*local-area network*

■ What kind of *service* does the Internet offer to end systems?

- Two end systems can communicate through the Internet, but exactly what kind of communication service is that of the Internet?

- Two end systems can communicate through the Internet, but exactly what kind of communication service is that of the Internet?

- ***Connectionless, "best effort"***
  - ▶ the network accepts "datagrams" for delivery—this is conceptually similar to the postal service
  - ▶ "best effort" really means *unreliable* though not malicious

■ Two end systems can communicate through the Internet, but exactly what kind of communication service is that of the Internet?

■ *Connectionless, "best effort"*
  ▶ the network accepts "datagrams" for delivery—this is conceptually similar to the postal service
  ▶ "best effort" really means *unreliable* though not malicious

■ *Connection-oriented, reliable*
  ▶ virtual duplex communication channel ($A \leftrightarrow B$)—conceptually similar to a telephone service
  ▶ information is transmitted "reliably" and in order

- How reliable is a "reliable" service?

- How reliable is a "reliable" service?

- The term "reliable" means that information will eventually reach its destination if a route is viable within a certain amount of time

- How reliable is a "reliable" service?

- The term "reliable" means that information will eventually reach its destination if a route is viable within a certain amount of time

- The network makes absolutely no guarantees on *latency* (i.e., the time it takes to transmit some information from a source to a destination)

application

| application |
|---|
| transport |

| application |
|:---:|
| transport |
| network |

| application |
| --- |
| transport |
| network |
| link |

| application |
| --- |
| transport |
| network |
| link |
| physical |

- *Application* (e.g., HTTP, SMTP, and DNS)
  - ▶ application functionalities
  - ▶ application messages

- *Application* (e.g., HTTP, SMTP, and DNS)
  - ▶ application functionalities
  - ▶ application messages

- *Transport* (e.g., TCP and UDP)
  - ▶ application multiplexing, reliable transfer (TCP), congestion control (TCP)
  - ▶ datagrams (UDP) or segments (TCP)

- *Application* (e.g., HTTP, SMTP, and DNS)
  - ▶ application functionalities
  - ▶ application messages

- *Transport* (e.g., TCP and UDP)
  - ▶ application multiplexing, reliable transfer (TCP), congestion control (TCP)
  - ▶ datagrams (UDP) or segments (TCP)

- *Network* (IP)
  - ▶ end to end datagram, best-effort service, routing, fragmentation
  - ▶ packets (IP)

- *Application* (e.g., HTTP, SMTP, and DNS)
  - ▶ application functionalities
  - ▶ application messages

- *Transport* (e.g., TCP and UDP)
  - ▶ application multiplexing, reliable transfer (TCP), congestion control (TCP)
  - ▶ datagrams (UDP) or segments (TCP)

- *Network* (IP)
  - ▶ end to end datagram, best-effort service, routing, fragmentation
  - ▶ packets (IP)

- *Link* (e.g., Ethernet and PPP)
  - ▶ point-to-point or local broadcast communication
  - ▶ frames (or packets)

- *Application* (e.g., HTTP, SMTP, and DNS)
  - ▶ application functionalities
  - ▶ application messages

- *Transport* (e.g., TCP and UDP)
  - ▶ application multiplexing, reliable transfer (TCP), congestion control (TCP)
  - ▶ datagrams (UDP) or segments (TCP)

- *Network* (IP)
  - ▶ end to end datagram, best-effort service, routing, fragmentation
  - ▶ packets (IP)

- *Link* (e.g., Ethernet and PPP)
  - ▶ point-to-point or local broadcast communication
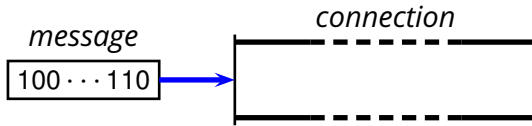  - ▶ frames (or packets)

- *Physical*

# Delay (Latency) and Rate (Throughput)

*connection*

*message*

*connection*

$100 \cdots 110$

*connection*

$100 \cdots 110$

$t_0$
first bit
enters

# Delay (Latency) and Rate (Throughput)

## Delay (Latency) and Rate (Throughput)



*connection*

$100 \cdots 110$

$t_0$
first bit
enters

$t_1$
first bit
exists

$t_2$
last bit
exits

# Delay (Latency) and Rate (Throughput)

# Delay (Latency) and Rate (Throughput)



*Propagation **Delay***         $d_{prop} = t_1 - t_0$        sec

# Delay (Latency) and Rate (Throughput)



*Propagation **Delay***  $d_{prop} = t_1 - t_0$   sec

*Transmission **Rate***  $R = \dfrac{\ell}{t_2 - t_1}$   bits/sec

# Delay (Latency) and Rate (Throughput)



*Propagation **Delay*** $\qquad d_{prop} = t_1 - t_0 \qquad$ sec

*Transmission **Rate*** $\qquad R = \dfrac{\ell}{t_2 - t_1} \qquad$ bits/sec

*Total transfer time* $\qquad d_{end\text{-}end} = d + \dfrac{\ell}{R} \qquad$ sec

$H_1$         $X$         $H_2$

$H_1$ ——— $d_1, R_1$ ——— $X$ ——— $d_2, R_2$ ——— $H_2$

$$d_{end\text{-}end} = d_1 + \frac{\ell}{R_1}$$

$$d_{end\text{-}end} = d_1 + \frac{\ell}{R_1} + d_x$$

$$d_{end\text{-}end} = d_1 + \frac{\ell}{R_1} + d_x + \frac{\ell}{R_2}$$

$$d_{end\text{-}end} = d_1 + \frac{\ell}{R_1} + d_x + \frac{\ell}{R_2} + d_2$$

$$d_{end\text{-}end} = d_1 + \frac{\ell}{R_1} + d_x + \frac{\ell}{R_2} + d_2$$

$$d_{end\text{-}end} = d_1 + \frac{\ell}{R_1} + d_x + \frac{\ell}{R_2} + d_2$$



$$d_{end\text{-}end} = N\left(d_p + \frac{\ell}{R} + d_x\right)$$

$$R_{end\text{-}end} =$$

$$R_{end\text{-}end} = \min\{R_1, R_2\}$$

$$R_{end\text{-}end} = \min\{R_1, R_2\}$$

$$R_{\text{end-end}} = \min\{R_1, R_2\}$$



$$R_{\text{end-end}} = \min\{R_1, R_2, \ldots, R_N\}$$

$$d_x = d_{cpu} + d_{queue}$$

$$d_{queue} = |q|/R_x$$

$d_x$

$H_1$ — $d_1, R_1$ — $X$ — $d_2, R_2$ — $H_2$

queue length

output rate

$d_x = d_{cpu} + d_{queue}$

$d_{queue} = |q|/R_x$

$$d_x = d_{cpu} + d_{queue}$$

queue length

where

$$d_{queue} = |q|/R_x$$

output rate

$$d_x = d_{cpu} + d_{queue}$$

queue length

where

$$d_{queue} = |q|/R_x$$

output rate

… $R_x$ is also the rate at which packets get out of the queue

HTTP

| HTTP | | SMTP |

| HTTP | | SMTP | | DNS |

| HTTP | | SMTP | | DNS |
|------|---|------|---|-----|

| HTTP | SMTP | DNS |
|------|------|-----|

GET / HTTP/1.1...

HTTP/1.1 200...

| HTTP | SMTP | DNS |
|------|------|-----|

```
 ── GET / HTTP/1.1... ──▶

 ◀── HTTP/1.1 200... ──

 ── GET... ──────────▶
```

```
┌─────────────┐   ┌─────────────┐   ┌─────────────┐
│    HTTP     │   │    SMTP     │   │     DNS     │
└─────────────┘   └─────────────┘   └─────────────┘
```

```
│                     │
│── GET / HTTP/1.1... ──→│
│                     │
│←── HTTP/1.1 200... ──│
│                     │
│── GET... ───────────→│
│                     │
│←── HTTP/1.1 200... ──│
│                     │
↓                     ↓
```

| HTTP | | SMTP | DNS |
|---|---|---|---|
| *connection* | | | |

```
GET / HTTP/1.1...  →
←  HTTP/1.1 200...

GET...  →
←  HTTP/1.1 200...
```

```
←  250 Yo...
HELO...  →
   ...
MAIL FROM:...  →
   ...
RCPT TO:...  →
   ...
QUIT  →
←  221 Bye
```

```
┌─────────────────────┐  ┌─────────────────────┐  ┌─────────────────────┐
│        HTTP         │  │        SMTP         │  │        DNS          │
├─────────────────────┤  ├─────────────────────┤  └─────────────────────┘
│     connection      │  │     connection      │
└─────────────────────┘  └─────────────────────┘
```

```
GET / HTTP/1.1...  →                  ←  250 Yo...
←  HTTP/1.1 200...                     HELO...  →
                                           ...
GET...  →                             MAIL FROM:...  →
←  HTTP/1.1 200...                         ...
                                      RCPT TO:...  →
                                           ...
                                      QUIT  →
                                      ←  221 Bye
```

| HTTP | SMTP | DNS |
|------|------|-----|
| *connection* | *connection* | *messages* |

HTTP connection:
- GET / HTTP/1.1... →
- ← HTTP/1.1 200...
- GET... →
- ← HTTP/1.1 200...

SMTP connection:
- ← 250 Yo...
- HELO... →
- ...
- MAIL FROM:... →
- ...
- RCPT TO:... →
- ...
- QUIT →
- ← 221 Bye

DNS:
- root
- .ch
- DNS ←
- unisi.ch
- app.

| HTTP | SMTP | DNS |
|------|------|-----|
| *connection* | *connection* | *messages* |

```
HTTP                          SMTP

┌──────────────────┐  ┌──────────────────┐
│ GET / HTTP/1.1... →│  │←    250 Yo...    │
│                   │  │                   │
│←  HTTP/1.1 200... │  │ HELO... ────────→ │
│                   │  │      ...          │
│                   │  │ MAIL FROM:... ──→ │
│ GET... ─────────→ │  │      ...          │
│                   │  │ RCPT TO:... ────→ │
│←  HTTP/1.1 200... │  │      ...          │
│                   │  │ QUIT ───────────→ │
│                   │  │←     221 Bye      │
└──────────────────┘  └──────────────────┘
```

DNS                messages

root

.ch

DNS ⟶ unisi.ch

app.

## Transport Control Protocol (TCP)

- conntection-oriented (i.e., "connections")

# Transport Layer in the Internet

- ***Transport Control Protocol (TCP)***
  - ▶ conntection-oriented (i.e., "connections")

- ***User Datagram Protocol (UDP)***
  - ▶ connectionless (i.e., "messages")

# Transport Layer in the Internet

- ***Transport Control Protocol (TCP)***
  - ▶ conntection-oriented (i.e., "connections")

- ***User Datagram Protocol (UDP)***
  - ▶ connectionless (i.e., "messages")

- Terminology
  - ▶ transport-layer packets are called ***segments***

# Transport Layer in the Internet

- *Transport Control Protocol (TCP)*
  - ▶ conntection-oriented (i.e., "connections")

- *User Datagram Protocol (UDP)*
  - ▶ connectionless (i.e., "messages")

- Terminology
  - ▶ transport-layer packets are called *segments*

- Basic assumptions on the underlying network layer

# Transport Layer in the Internet

- **_Transport Control Protocol (TCP)_**
    - ▶ conntection-oriented (i.e., "connections")

- **_User Datagram Protocol (UDP)_**
    - ▶ connectionless (i.e., "messages")

- Terminology
    - ▶ transport-layer packets are called **_segments_**

- Basic assumptions on the underlying network layer
    - ▶ every host has one unique _IP address_

# Transport Layer in the Internet

- *Transport Control Protocol (TCP)*
  - conntection-oriented (i.e., "connections")

- *User Datagram Protocol (UDP)*
  - connectionless (i.e., "messages")

- Terminology
  - transport-layer packets are called *segments*

- Basic assumptions on the underlying network layer
  - every host has one unique *IP address*
  - best-effort delivery service
    - no guarantees on the integrity of segments
    - no guarantees on the order in which segments are delivered

# Transport-Layer Value-Added Service

- ***Transport-layer multiplexing/demultiplexing***
  - ▶ i.e., connecting applications as opposed to hosts

# Transport-Layer Value-Added Service

- ***Transport-layer multiplexing/demultiplexing***
  - ▶ i.e., connecting applications as opposed to hosts

- ***Reliable data transfer***
  - ▶ i.e., integrity and possibly ordered delivery

# Transport-Layer Value-Added Service

- ***Transport-layer multiplexing/demultiplexing***
  - ▶ i.e., connecting applications as opposed to hosts

- ***Reliable data transfer***
  - ▶ i.e., integrity and possibly ordered delivery

- ***Connections***
  - ▶ i.e., streams
  - ▶ can be seen as the same as ordered delivery

# Transport-Layer Value-Added Service

- ## *Transport-layer multiplexing/demultiplexing*
  - ▶ i.e., connecting applications as opposed to hosts

- ## *Reliable data transfer*
  - ▶ i.e., integrity and possibly ordered delivery

- ## *Connections*
  - ▶ i.e., streams
  - ▶ can be seen as the same as ordered delivery

- ## *Congestion control*
  - ▶ i.e., end-to-end traffic (admission) control so as to avoid destructive congestions within the network

How do we distinguish all these "connections"?

How do we distinguish all these "connections"?
(in this case, connections between the same two hosts)

- Each application running on a host is identified (within that host) by a unique *port number*
    - ▶ port numbers are simply cross-platform process identifiers

- Each application running on a host is identified (within that host) by a unique *port number*
    - ▶ port numbers are simply cross-platform process identifiers

- How do we identify a "connection"?

- Each application running on a host is identified (within that host) by a unique *port number*
  - ▶ port numbers are simply cross-platform process identifiers

- How do we identify a "connection"?
  - ▶ two pairs of *host* and *application* identifiers
  - ▶ i.e., two pairs *(IP-address, port)*

- Each application running on a host is identified (within that host) by a unique *port number*

  - ▶ port numbers are simply cross-platform process identifiers

- How do we identify a "connection"?

  - ▶ two pairs of *host* and *application* identifiers
  - ▶ i.e., two pairs *(IP-address, port)*

- How do we find out which application (host and port number) to connect to?

■ Each application running on a host is identified (within that host) by a unique *port number*

    ▶ port numbers are simply cross-platform process identifiers

■ How do we identify a "connection"?

    ▶ two pairs of *host* and *application* identifiers

    ▶ i.e., two pairs *(IP-address, port)*

■ How do we find out which application (host and port number) to connect to?

    ▶ outside the scope of the definition of the transport layer

    ▶ but of course we can have "well-known" service numbers

- The message format of both UDP and TCP starts with the source and destination port numbers

| 0 | 15 16 | 31 |
|---|---|---|
| *source port* | | *destination port* |

. . .

- The message format of both UDP and TCP starts with the source and destination port numbers

| 0               15 | 16              31 |
|---|---|
| *source port* | *destination port* |

$\cdots$

- E.g.,

| *src port* | *dst port* |
|---|---|
| 1234 | 80 |

A ────────────────────▶ B

- The message format of both UDP and TCP starts with the source and destination port numbers

| 0                   15 | 16                31 |
|---|---|
| *source port* | *destination port* |

· · ·

- E.g.,

*src port*    *dst port*

| 1234 | 80 |
|---|---|



| | |
|---|---|

*src port*    *dst port*

- The message format of both UDP and TCP starts with the source and destination port numbers

| 0               15 | 16              31 |
|---|---|
| *source port* | *destination port* |

· · ·

- E.g.,

*src port*  *dst port*

| 1234 | 80 |
|---|---|

A ⟶ B

A ⟵ B

| 80 | 1234 |
|---|---|

*src port*  *dst port*

- The Internet's primary transport protocol
  - ▶ defined in RFC 793, RFC 1122, RFC 1323, RFC 2018, and RFC 2581

# Transmission Control Protocol

- The Internet's primary transport protocol
  - ▶ defined in RFC 793, RFC 1122, RFC 1323, RFC 2018, and RFC 2581

- Connection-oriented service
  - ▶ endpoints "shake hands" to establish a connection
  - ▶ not a circuit-switched connection, nor a virtual circuit

# Transmission Control Protocol

- The Internet's primary transport protocol
  - ▶ defined in RFC 793, RFC 1122, RFC 1323, RFC 2018, and RFC 2581

- Connection-oriented service
  - ▶ endpoints "shake hands" to establish a connection
  - ▶ not a circuit-switched connection, nor a virtual circuit

- Full-duplex service
  - ▶ both endpoints can both send and receive, at the same time

- *TCP segment:* envelope for TCP data
  - ▶ TCP data are sent within TCP segments
  - ▶ TCP segments are usually sent within an IP packet

- **_TCP segment:_** envelope for TCP data
  - ▶ TCP data are sent within TCP segments
  - ▶ TCP segments are usually sent within an IP packet

- **_Maximum segment size (MSS):_** maximum amount of application data transmitted in a single segment
  - ▶ typically related to the MTU of the connection, to avoid network-level fragmentation (we'll talk about all of this later)

- **_TCP segment:_** envelope for TCP data
  - ▶ TCP data are sent within TCP segments
  - ▶ TCP segments are usually sent within an IP packet

- **_Maximum segment size (MSS):_** maximum amount of application data transmitted in a single segment
  - ▶ typically related to the MTU of the connection, to avoid network-level fragmentation (we'll talk about all of this later)

- **_Maximum transmission unit (MTU):_** largest link-layer frame available to the sender host
  - ▶ *path MTU:* largest link-layer frame that can be sent on all links from the sender host to the receiver host

| 0 | | | 31 |
|---|---|---|---|
| source port | | destination port | |
| sequence number | | | |
| acknowledgment number | | | |
| hdrlen | unused U A P R S F | receive window | |
| Internet checksum | | urgent data pointer | |
| options field | | | |
| data | | | |

- *Source and destination ports:* (16-bit each) application identifiers

- *Source and destination ports:* (16-bit each) application identifiers

- *Sequence number:* (32-bit) used to implement reliable data transfer

- *Acknowledgment number:* (32-bit) used to implement reliable data transfer

- *Source and destination ports:* (16-bit each) application identifiers

- *Sequence number:* (32-bit) used to implement reliable data transfer

- *Acknowledgment number:* (32-bit) used to implement reliable data transfer

- *Receive window:* (16-bit) size of the "window" on the receiver end

# TCP Header Fields

- *Source and destination ports:* (16-bit each) application identifiers

- *Sequence number:* (32-bit) used to implement reliable data transfer

- *Acknowledgment number:* (32-bit) used to implement reliable data transfer

- *Receive window:* (16-bit) size of the "window" on the receiver end

- *Header length:* (4-bit) size of the TCP header in 32-bit words

# TCP Header Fields

- *Source and destination ports:* (16-bit each) application identifiers

- *Sequence number:* (32-bit) used to implement reliable data transfer

- *Acknowledgment number:* (32-bit) used to implement reliable data transfer

- *Receive window:* (16-bit) size of the "window" on the receiver end

- *Header length:* (4-bit) size of the TCP header in 32-bit words

- *Optional and variable-length options field:* may be used to negotiate protocol parameters

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *FIN flag:* (1-bit) used during connection shutdown

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *FIN flag:* (1-bit) used during connection shutdown

- *PSH flag:* (1-bit) "push" flag, used to solicit the receiver to pass the data to the application immediately

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *FIN flag:* (1-bit) used during connection shutdown

- *PSH flag:* (1-bit) "push" flag, used to solicit the receiver to pass the data to the application immediately

- *URG flag:* (1-bit) "urgent" flag, used to inform the receiver that the sender has marked some data as "urgent". The location of this urgent data is marked by the *urgent data pointer* field

- *ACK flag:* (1-bit) signals that the value contained in the *acknowledgment number* represents a valid acknowledgment

- *SYN flag:* (1-bit) used during connection setup and shutdown

- *RST flag:* (1-bit) used during connection setup and shutdown

- *FIN flag:* (1-bit) used during connection shutdown

- *PSH flag:* (1-bit) "push" flag, used to solicit the receiver to pass the data to the application immediately

- *URG flag:* (1-bit) "urgent" flag, used to inform the receiver that the sender has marked some data as "urgent". The location of this urgent data is marked by the *urgent data pointer* field

- *Checksum:* (16-bit) used to detect transmission errors

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*

*application data stream*

|←──────────── 4Kb ────────────→|

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*

*application data stream*

├────────────────────── 4Kb ──────────────────────┤

├─ MSS=1024b ─┤

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*

*application data stream*

├───────────────────────── 4Kb ─────────────────────────┤

| | | | |
|---|---|---|---|

├─ MSS=1024b ─┤

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*

*application data stream*

├─────────────────── 4Kb ───────────────────┤
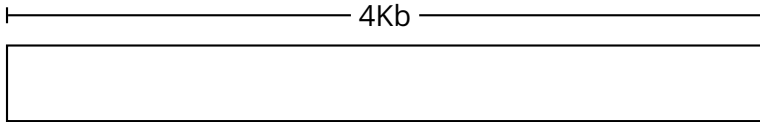
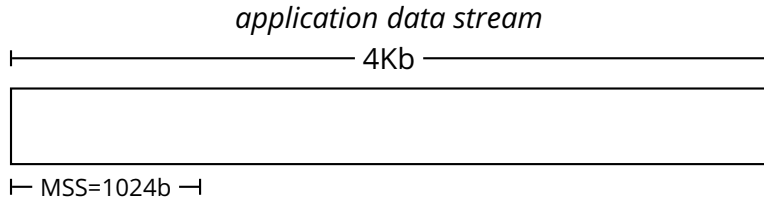├─ MSS=1024b ─┤

1… …1024  1025… 2048  2049… 3072  3073… 4096

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*

*application data stream*



├──────────── 4Kb ────────────┤

├─ MSS=1024b ─┤

1... ...1024  1025... 2048  2049... 3072  3073... 4096

*a TCP segment*

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*
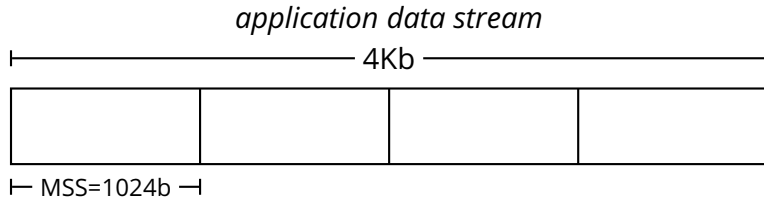
*application data stream*

├─────────────── 4Kb ───────────────┤

├─ MSS=1024b ─┤

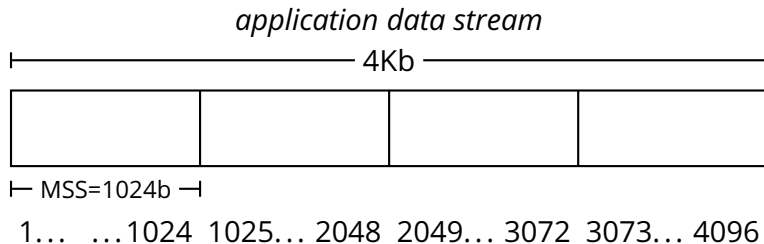1... ...1024  1025... 2048  2049... 3072  3073... 4096

*a TCP segment*

2049

- Sequence numbers are associated with *bytes* in the data stream
  - ▶ not with segments, as we have used them before

- The *sequence number* in a TCP segment indicates *the sequence number of the first byte carried by that segment*
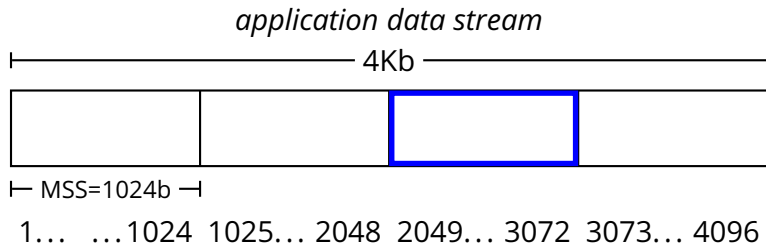
*application data stream*



├────────────────── 4Kb ──────────────────┤

├─ MSS=1024b ─┤

1… …1024  1025… 2048  2049… 3072  3073… 4096

sequence number   *a TCP segment*

2049

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - ▶ TCP acknowledgments are *cumulative*

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
    - ▶ TCP acknowledgments are *cumulative*

A                                                            B

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - ▶ TCP acknowledgments are *cumulative*

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - ▶ TCP acknowledgments are *cumulative*

- An *acknowledgment number* represents the *first sequence number not yet seen by the receiver*
  - ▶ TCP acknowledgments are *cumulative*

A                                                                                    B

——— [$Seq\# = 1200, \ldots$], $size(data) = 1000$ ———→

——— [$Seq\# = 2200, \ldots$], $size(data) = 500$ ———→

←——— [$Seq\# = \ldots, Ack\# = 2700$] ———

# Sequence Numbers and ACK Numbers

- Notice that a TCP connection is a *full-duplex* link
  - therefore, there are ***two streams***
  - two different sequence numbers

■ Notice that a TCP connection is a *full-duplex* link

▶ therefore, there are ***two streams***

▶ two different sequence numbers

E.g., consider a simple "Echo" application:

A                                                     B

■ Notice that a TCP connection is a *full-duplex* link
  ▶ therefore, there are ***two streams***
  ▶ two different sequence numbers

E.g., consider a simple "Echo" application:

A                                                                        B
│                                                                        │
│──── [*Seq#* = 100, *Data* ="C"] ──────────────────────────────▶│
│                                                                        │
│                                                                        │
│                                                                        │
│                                                                        │
│                                                                        │
│                                                                        │
▼                                                                        ▼

## Sequence Numbers and ACK Numbers

- Notice that a TCP connection is a *full-duplex* link
  - ▶ therefore, there are ***two streams***
  - ▶ two different sequence numbers

  E.g., consider a simple "Echo" application:

  A                                                       B

  $\longrightarrow$ [$Seq\# = 100$, $Data =$"C"] $\longrightarrow$

  $\longleftarrow$ [$Ack\# = 101$, $Seq\# = 200$, $Data =$"C"] $\longrightarrow$

# Sequence Numbers and ACK Numbers

■ Notice that a TCP connection is a *full-duplex* link
  ▶ therefore, there are ***two streams***
  ▶ two different sequence numbers
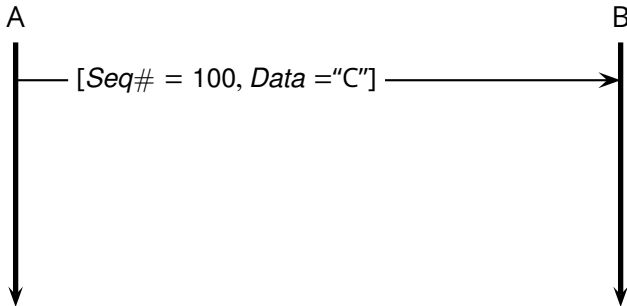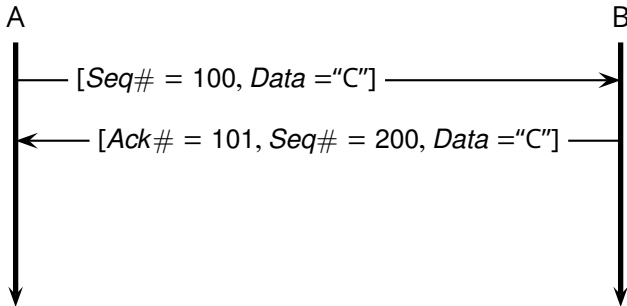
E.g., consider a simple "Echo" application:

A                                                                    B
│                                                                    │
│──── [$Seq\# = 100$, $Data =$"C"] ────────────────────────────────→│
│                                                                    │
│←──── [$Ack\# = 101$, $Seq\# = 200$, $Data =$"C"] ─────────────────│
│                                                                    │
│──── [$Seq\# = 101$, $Ack\# = 201$, $Data =$"i"] ─────────────────→│
│                                                                    │
│                                                                    │
↓                                                                    ↓

# Sequence Numbers and ACK Numbers

- Notice that a TCP connection is a *full-duplex* link
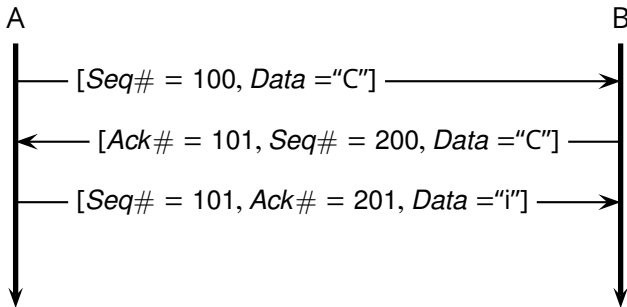  - therefore, there are *two streams*
  - two different sequence numbers
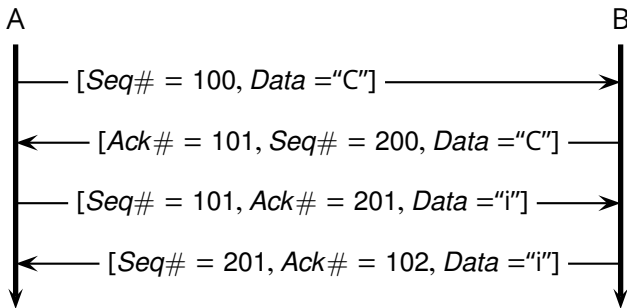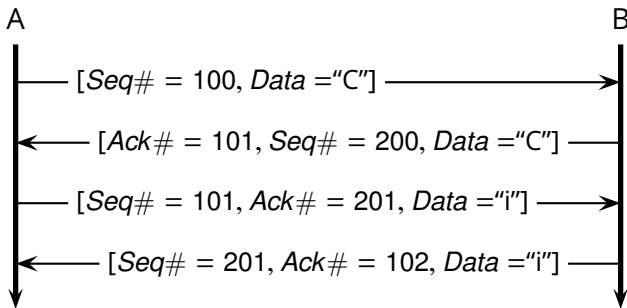
E.g., consider a simple "Echo" application:

```
A                                                              B
│──── [Seq# = 100, Data ="C"] ──────────────────────────────▶│
│                                                              │
│◀──── [Ack# = 101, Seq# = 200, Data ="C"] ─────────────────│
│                                                              │
│──── [Seq# = 101, Ack# = 201, Data ="i"] ─────────────────▶│
│                                                              │
│◀──── [Seq# = 201, Ack# = 102, Data ="i"] ─────────────────│
▼                                                              ▼
```

# Sequence Numbers and ACK Numbers

- Notice that a TCP connection is a *full-duplex* link
  - ▶ therefore, there are *two streams*
  - ▶ two different sequence numbers

E.g., consider a simple "Echo" application:

```
A                                                        B
│                                                        │
│──── [Seq# = 100, Data ="C"] ──────────────────────────▶│
│                                                        │
│◀──── [Ack# = 101, Seq# = 200, Data ="C"] ──────────────│
│                                                        │
│──── [Seq# = 101, Ack# = 201, Data ="i"] ──────────────▶│
│                                                        │
│◀──── [Seq# = 201, Ack# = 102, Data ="i"] ──────────────│
▼                                                        ▼
```

- Acknowledgments are "piggybacked" on data segments

# Reliability and Timeout

- Duplicate acknowledgments to detect lost segments
  - ▶ receiver notices a missing packet → duplicate ACKs → retransmission by sender

- A *timer* to detect lost segments
  - ▶ timeout without an ACK → lost packet → retransmission

- Duplicate acknowledgments to detect lost segments
  - ▶ receiver notices a missing packet → duplicate ACKs → retransmission by sender

- A *timer* to detect lost segments
  - ▶ timeout without an ACK → lost packet → retransmission

- How long to wait for acknowledgments?

# Reliability and Timeout

- Duplicate acknowledgments to detect lost segments
  - ▶ receiver notices a missing packet → duplicate ACKs → retransmission by sender

- A *timer* to detect lost segments
  - ▶ timeout without an ACK → lost packet → retransmission

- How long to wait for acknowledgments?

- Retransmission timeouts should be larger than the round-trip time $RTT = 2L$
  - ▶ as close as possible to the $RTT$

## Reliability and Timeout

- Duplicate acknowledgments to detect lost segments
  - ▶ receiver notices a missing packet → duplicate ACKs → retransmission by sender

- A *timer* to detect lost segments
  - ▶ timeout without an ACK → lost packet → retransmission

- How long to wait for acknowledgments?

- Retransmission timeouts should be larger than the round-trip time $RTT = 2L$
  - ▶ as close as possible to the $RTT$

- TCP controls its timeout by continuously *estimating the current RTT*

- RTT is measured using ACKs
  - ▶ only for packets transmitted once

- Given a single sample $S$ at any given time

- *Exponential weighted moving average (EWMA)*

$$\overline{RTT} = (1 - \alpha)\overline{RTT}' + \alpha S$$

- RTT is measured using ACKs
  - ▶ only for packets transmitted once

- Given a single sample $S$ at any given time

- *Exponential weighted moving average (EWMA)*

$$\overline{RTT} = (1 - \alpha)\overline{RTT}' + \alpha S$$

  - ▶ RFC 2988 recommends $\alpha = 0.125$

# Round-Trip Time Estimation

- RTT is measured using ACKs
  - only for packets transmitted once

- Given a single sample $S$ at any given time

- *Exponential weighted moving average (EWMA)*

$$\overline{RTT} = (1 - \alpha)\overline{RTT}' + \alpha S$$

  - RFC 2988 recommends $\alpha = 0.125$

- TCP also measures the *variability of RTT*

$$\overline{DevRTT} = (1 - \beta)\overline{DevRTT}' + \beta|\overline{RTT}' - S|$$

# Round-Trip Time Estimation

- RTT is measured using ACKs
  - only for packets transmitted once

- Given a single sample $S$ at any given time

- *Exponential weighted moving average (EWMA)*

$$\overline{RTT} = (1 - \alpha)\overline{RTT}' + \alpha S$$

  - RFC 2988 recommends $\alpha = 0.125$

- TCP also measures the *variability of RTT*

$$\overline{DevRTT} = (1 - \beta)\overline{DevRTT}' + \beta|\overline{RTT}' - S|$$

  - RFC 2988 recommends $\beta = 0.25$

- The timeout interval $T$ must be larger than the RTT
  - ▶ so as to avoid unnecessary retransmission

- However, $T$ should not be too far from RTT
  - ▶ so as to detect (and retransmit) lost segments as quickly as possible

- The timeout interval $T$ must be larger than the RTT
  - ▶ so as to avoid unnecessary retransmission

- However, $T$ should not be too far from RTT
  - ▶ so as to detect (and retransmit) lost segments as quickly as possible

- TCP sets its timeouts using the estimated RTT ($\overline{RTT}$) and the variability estimate $\overline{DevRTT}$:

$$T = \overline{RTT} + 4\overline{DevRTT}$$

A simplified TCP sender

- r_send(*data*)

  ---

  **if** (timer not running)
    start_timer()
  u_send([*data*,*next_seq_num*])
  *next_seq_num* ← *next_seq_num* + *length*(*data*)

A simplified TCP sender

- r_send(*data*)
  
  **if** (timer not running)
    start_timer()
  u_send([*data*,*next_seq_num*])
  *next_seq_num* ← *next_seq_num* + *length*(*data*)

- timeout
  
  u_send(pending segment with smallest sequence number)
  start_timer()

A simplified TCP sender

- r_send(*data*)

  **if** (timer not running)
    start_timer()
  u_send([*data,next_seq_num*])
  *next_seq_num* ← *next_seq_num* + *length*(*data*)

- timeout

  u_send(pending segment with smallest sequence number)
  start_timer()

- u_recv([ACK,*y*])

  **if** ($y >$ *base*)
    *base* ← *y*
    **if** (there are pending segments)
      start_timer()
  **else** ...

■ Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
    - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
    - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
    - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

- Arrival of out of order segment with higher-than-expected sequence number (gap detected)

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

- Arrival of out of order segment with higher-than-expected sequence number (gap detected)
  - *Duplicate ACK:* immediately send duplicate ACK

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

- Arrival of out of order segment with higher-than-expected sequence number (gap detected)
  - *Duplicate ACK:* immediately send duplicate ACK

- Arrival of segment that (partially or completely) fills a gap in the received data

# Acknowledgment Generation (Receiver)

- Arrival of in-order segment with expected sequence number; all data up to expected sequence number already acknowledged
  - *Delayed ACK:* wait 500ms for another in-order segment; If that does not arrive, send ACK

- Arrival of in-order segment with expected sequence number. One other in-order segment waiting for ACK (see above)
  - *Cumulative ACK:* immediately send cumulative ACK (for both segments)

- Arrival of out of order segment with higher-than-expected sequence number (gap detected)
  - *Duplicate ACK:* immediately send duplicate ACK

- Arrival of segment that (partially or completely) fills a gap in the received data
  - *Immediate ACK:* immediately send ACK if the packet start at the lower end of the gap

■ u_recv([ACK,*y*])

$\overline{\text{if } (y > base)}$
  *base* ← *y*
  **if** (there are pending segments)
    start_timer()

- u_recv([ACK,*y*])

  **if** (*y* > *base*)

    *base* ← *y*

    **if** (there are pending segments)

      start_timer()

  **else**

    *ack_counter*[*y*] ← *ack_counter*[*y*] + 1

    **if** (*ack_counter*[*y*] = 3)

      u_send(segment with sequence number *y*)

Three-way handshake

Three-way handshake

client                                          server

Three-way handshake



client                                                                                    server

[*SYN*, *Seq# = cli_init_seq*] ⟶

Three-way handshake

client                                                                    server

[$SYN$, $Seq\# = cli\_init\_seq$] ────────────────→

←──── [$SYN$, $ACK$, $Ack\# = cli\_init\_seq + 1$, $Seq\# = srv\_init\_seq$]

Three-way handshake

client                                                                              server

$[SYN, Seq\# = cli\_init\_seq]$ $\longrightarrow$

$\longleftarrow [SYN, ACK, Ack\# = cli\_init\_seq + 1, Seq\# = srv\_init\_seq]$

$[ACK, Seq\# = cli\_init\_seq + 1, Ack\# = srv\_init\_seq + 1] \longrightarrow$

"This is it."
"Okay, Bye now."
"Bye."

# Connection Shutdown

"This is it."
"Okay, Bye now."
"Bye."

client                                          server

"This is it."
"Okay, Bye now."
"Bye."

"This is it."
"Okay, Bye now."
"Bye."

| client | | server |
|--------|--|--------|

[*FIN*] ⟶

⟵ [*ACK*]

"This is it."
"Okay, Bye now."
"Bye."

"This is it."
"Okay, Bye now."
"Bye."

CLOSED

# The TCP State Machine (Client)

```
CLOSED
```

application
<u>opens connection</u>
send SYN

```
SYN_SENT
```

CLOSED

application
opens connection
send SYN

SYN_SENT

receive SYN,ACK
send ACK

ESTABLISHED

# The TCP State Machine (Client)

# The TCP State Machine (Client)

application
opens connection
send SYN

CLOSED

wait 30 seconds

TIME_WAIT

SYN_SENT

receive FIN
send ACK

receive SYN,ACK
send ACK

FIN_WAIT_2

ESTABLISHED

receive ACK

application
closes connection
send FIN

FIN_WAIT_1

CLOSED

CLOSED

application
opens server socket

LISTEN

# The TCP State Machine (Server)

# The TCP State Machine (Server)

CLOSED

application
opens server socket

LISTEN

receive SYN
send SYN,ACK

SYN_RCVD

CLOSE_WAIT

receive ACK

receive FIN
send ACK

ESTABLISHED

- A router behaves a lot like a kitchen sink

■ A router behaves a lot like a kitchen sink



max rate = $R$

- A router behaves a lot like a kitchen sink



$\lambda_1 = R/2$

max rate $= R$

throughput $= R/2$

- A router behaves a lot like a kitchen sink



$\lambda_1 = R/2$

$\lambda_2 = R/2$

max rate $= R$

throughput $= R$

- A router behaves a lot like a kitchen sink



$\lambda_3 = R/2$

$\lambda_1 = R/2$             $\lambda_2 = R/2$

max rate $= R$

throughput $= R$

■ A router behaves a lot like a kitchen sink



$\lambda_3 = R/2$

$\lambda_1 = R/2$

$\lambda_2 = R/2$

max rate $= R$

throughput $= R$

- A router behaves a lot like a kitchen sink



$\lambda_3 = R/2$

$\lambda_1 = R/2$     $\lambda_2 = R/2$

max rate = $R$

throughput = $R$

- A router behaves a lot like a kitchen sink



$\lambda_3 = R/2$

$\lambda_1 = R/2$

$\lambda_2 = R/2$

max rate $= R$

throughput $= R$

- A router behaves a lot like a kitchen sink



$\lambda_3 = R/2$

$\lambda_1 = R/2$

$\lambda_2 = R/2$

max rate $= R$

throughput $= R$

- Total latency is the sum of link latency, processing time, and the time that a packet spends in the input queue

$$L = d_{TX} + d_{CPU} + d_q \quad \text{where } d_q = |q|/R$$

■ Total latency is the sum of link latency, processing time, and the time that a packet spends in the input queue

$$L = d_{TX} + d_{CPU} + d_q \quad \text{where } d_q = |q|/R$$

■ *Ideal case:* constant input data rate

$$\lambda_{in} < R$$

In this case the $d_q = 0$, because $|q| = 0$ \hfill (ideal input flow)

- Total latency is the sum of link latency, processing time, and the time that a packet spends in the input queue

$$L = d_{TX} + d_{CPU} + d_q \quad \text{where } d_q = |q|/R$$

- *Ideal case:* constant input data rate

$$\lambda_{in} < R$$

In this case the $d_q = 0$, because $|q| = 0$             (ideal input flow)

- *Extreme case:* constant input data rate

$$\lambda_{in} > R$$

In this case $|q| = (\lambda_{in} - R)t$ and therefore

$$d_q = \frac{\lambda_{in} - R}{R}t$$

■ Steady-state queuing delay

$$d_q = \begin{cases} 0 & \lambda_{in} < R \\ \frac{\lambda_{in} - R}{R} t & \lambda_{in} > R \end{cases}$$

■ Steady-state queuing delay

$$d_q = \begin{cases} 0 & \lambda_{in} < R \\ \frac{\lambda_{in}-R}{R}t & \lambda_{in} > R \end{cases}$$



ideal input flow
$\lambda_{in}$ constant

■ Steady-state queuing delay

$$d_q = \begin{cases} 0 & \lambda_{in} < R \\ \frac{\lambda_{in} - R}{R} t & \lambda_{in} > R \end{cases}$$



ideal input flow
$\lambda_{in}$ constant

realistic input flow
$\lambda_{in}$ variable

■ **Conclusion:** as the input rate $\lambda_{in}$ approaches the maximum throughput $R$, packets will experience very long delays

- **Conclusion:** as the input rate $\lambda_{in}$ approaches the maximum throughput $R$, packets will experience very long delays

- More realistic assumptions and models
  - ▶ finite queue length (buffers) in routers
  - ▶ effects of retransmission overhead
  - ▶ full queues along multi-hops paths

- **Conclusion:** as the input rate $\lambda_{in}$ approaches the maximum throughput $R$, packets will experience very long delays

- More realistic assumptions and models
    - ▶ finite queue length (buffers) in routers
    - ▶ effects of retransmission overhead
    - ▶ full queues along multi-hops paths

- **Conclusion:** as the input rate $\lambda_{in}$ approaches the maximum throughput $R$, packets will experience very long delays

- More realistic assumptions and models
    - ▶ finite queue length (buffers) in routers
    - ▶ effects of retransmission overhead
    - ▶ full queues along multi-hops paths

■ What to do when the network is congested?



$\lambda_3 = R/2$

$\lambda_1 = R/2$                 $\lambda_2 = R/2$

max rate $= R$

throughput $= R$

■ What to do when the network is congested?   ***BACK OFF!***

$\lambda_3 = R/2$

$\lambda_1 = R/2$

$\lambda_2 = R/2$

max rate $= R$

throughput $= R$

■ What to do when the network is congested? **_BACK OFF!_**



$\lambda_3 = R/4$

$\lambda_1 = R/4$

$\lambda_2 = R/4$

max rate $= R$

throughput $= R$

■ What to do when the network is congested?   *BACK OFF!*



$\lambda_3 = R/4$

$\lambda_1 = R/4$                    $\lambda_2 = R/4$

max rate $= R$

throughput $= R$

- What to do when the network is congested?   *BACK OFF!*

■ What to do when the network is congested?   *BACK OFF!*



$\lambda_3 = R/4$

$\lambda_1 = R/4$          $\lambda_2 = R/4$

max rate $= R$

throughput $= R$

**Approach:**

- ***The sender limits its output rate according to the state of the network***
  - ▶ The sender output rate becomes (part of) the input rate for the network ($\lambda_{in}$)

**Approach:**

- *The sender limits its output rate according to the state of the network*
  - ▶ The sender output rate becomes (part of) the input rate for the network ($\lambda_{in}$)

**Ingredients:**

1. How does the sender *measure the state of the network*?
   - ▶ we need *eyes* to see the traffic ahead

**Approach:**

- ■ *The sender limits its output rate according to the state of the network*
  - ▶ The sender output rate becomes (part of) the input rate for the network ($\lambda_{in}$)

**Ingredients:**

1. How does the sender *measure the state of the network*?

   - ▶ we need *eyes* to see the traffic ahead

2. how does the sender *set its output rate*?

   - ▶ we need *accelerator* and *brakes* to speed up or slow down

**Approach:**

- ■ *The sender limits its output rate according to the state of the network*
    - ▶ The sender output rate becomes (part of) the input rate for the network ($\lambda_{in}$)

**Ingredients:**

1. How does the sender *measure the state of the network*?
    - ▶ we need *eyes* to see the traffic ahead

2. how does the sender *set its output rate*?
    - ▶ we need *accelerator* and *brakes* to speed up or slow down

3. how should the sender *control its output rate*?
    - ▶ we need a *brain* and we need to know *how to drive!*

- If all traffic is correctly acknowledged, with fresh acknowledgments, then the sender assumes (quite correctly) that there is no congestion

■ If all traffic is correctly acknowledged, with fresh acknowledgments, then the sender assumes (quite correctly) that there is no congestion

■ Congestion means that some queues overflow in one or more routers between the sender and the receiver

  ▶ the visible effect is that some segments are dropped

■ If all traffic is correctly acknowledged, with fresh acknowledgments, then the sender assumes (quite correctly) that there is no congestion

■ Congestion means that some queues overflow in one or more routers between the sender and the receiver
  ▶ the visible effect is that some segments are dropped

■ Therefore the sender assumes that the network is congested when it (the sender) detects a segment loss
  ▶ duplicate acknowledgements (i.e., NACK)
  ▶ time out (i.e., no ACKs at all)

- The sender maintains a *congestion window* $W$

# Congestion Window (Accelerator/Brakes)

- The sender maintains a *congestion window* $W$

- The congestion window limits the amount of bytes that the sender pushes into the network before blocking waiting for acknowledgments

## Congestion Window (Accelerator/Brakes)

- The sender maintains a *congestion window W*

- The congestion window limits the amount of bytes that the sender pushes into the network before blocking waiting for acknowledgments

$$LastByteSent - LastByteAcked \leq W$$

where

$$W = \min\left(CongestionWindow, ReceiverWindow\right)$$

# Congestion Window (Accelerator/Brakes)

- The sender maintains a ***congestion window*** *W*

- The congestion window limits the amount of bytes that the sender pushes into the network before blocking waiting for acknowledgments

$$LastByteSent - LastByteAcked \leq W$$

where

$$W = \min \left( CongestionWindow, ReceiverWindow \right)$$

- The resulting maximum output rate is roughly

$$\lambda = \frac{W}{2L}$$

- *Additive-increase and multiplicative-decrease*

- *Additive-increase and multiplicative-decrease*

- *Slow start*

# Congestion Control (Brain, Algorithm)

- *Additive-increase and multiplicative-decrease*

- *Slow start*

- *Reaction to timeout events*

# Additive-Increase/Multiplicative-Decrease

- ***How* W *is reduced:*** at every *loss* event, TCP halves the congestion window

# Additive-Increase/Multiplicative-Decrease

- *How $W$ is reduced:* at every *loss* event, TCP halves the congestion window
  - e.g., suppose the window size $W$ is currently 20Kb, and a loss is detected
  - TCP reduces $W$ to 10Kb

# Additive-Increase/Multiplicative-Decrease

- **How $W$ is reduced:** at every *loss* event, TCP halves the congestion window

  - e.g., suppose the window size $W$ is currently 20Kb, and a loss is detected
  - TCP reduces $W$ to 10Kb

- **How $W$ is increased:** at every (good) acknowledgment, TCP increments $W$ by $1MSS/W$, so as to increase $W$ by $MSS$ every round-trip time $2L$. This process is called ***congestion avoidance***

# Additive-Increase/Multiplicative-Decrease

- **How $W$ is reduced:** at every *loss* event, TCP halves the congestion window
  - e.g., suppose the window size $W$ is currently 20Kb, and a loss is detected
  - TCP reduces $W$ to 10Kb

- **How $W$ is increased:** at every (good) acknowledgment, TCP increments $W$ by $1MSS/W$, so as to increase $W$ by $MSS$ every round-trip time $2L$. This process is called ***congestion avoidance***
  - e.g., suppose $W = 14600$ and $MSS = 1460$, then the sender increases $W$ to 16060 after 10 acknowledgments acknowledgments

# Additive-Increase/Multiplicative-Decrease

■ Window size *W* over time



Time

- What is the initial value of $W$?

- What is the initial value of $W$?

- The initial value of $W$ is *MSS*, that is *1 segment*, which is quite low for modern networks

- What is the initial value of *W*?

- The initial value of *W* is *MSS*, that is ***1 segment***, which is quite low for modern networks

- To get quickly to a good throughput level, TCP increases its sending rate exponentially for its first growth phase, up to a ***slow-start threshold*** (***ssthresh***)

- What is the initial value of *W*?

- The initial value of *W* is *MSS*, that is *1 segment*, which is quite low for modern networks

- To get quickly to a good throughput level, TCP increases its sending rate exponentially for its first growth phase, up to a *slow-start threshold* (*ssthresh*)

- After the threshold, TCP proceeds with its linear push

- What is the initial value of $W$?

- The initial value of $W$ is *MSS*, that is ***1 segment***, which is quite low for modern networks

- To get quickly to a good throughput level, TCP increases its sending rate exponentially for its first growth phase, up to a ***slow-start threshold*** (***ssthresh***)

- After the threshold, TCP proceeds with its linear push

- This process is called "slow start" because of the small initial value of $W$

- As we know, three duplicate ACKs are interpreted as a NACK

- As we know, three duplicate ACKs are interpreted as a NACK

- Both timeouts and NACKs signal a loss, but they say different things about the status of the network

- As we know, three duplicate ACKs are interpreted as a NACK

- Both timeouts and NACKs signal a loss, but they say different things about the status of the network

- A *timeout indicates congestion*

- As we know, three duplicate ACKs are interpreted as a NACK

- Both timeouts and NACKs signal a loss, but they say different things about the status of the network

- A *timeout indicates congestion*

- Three (duplicate) ACKs suggest that the network is still able to deliver segments along that path

- As we know, three duplicate ACKs are interpreted as a NACK

- Both timeouts and NACKs signal a loss, but they say different things about the status of the network

- A *timeout indicates congestion*

- Three (duplicate) ACKs suggest that the network is still able to deliver segments along that path

- So, TCP reacts differently to a timeout and to a triple duplicate ACKs

Assuming the current window size is $W = \overline{W}$

Assuming the current window size is $W = \overline{W}$

- *Timeout*
    - ▶ go back to $W = MSS$
    - ▶ set *ssthresh* $= \overline{W}/2$
    - ▶ run *slow start* up to $W = ssthresh$
    - ▶ then proceed with *congestion avoidance*

Assuming the current window size is $W = \overline{W}$

- *Timeout*

    - ▶ go back to $W = MSS$

    - ▶ set *ssthresh* $= \overline{W}/2$

    - ▶ run *slow start* up to $W = ssthresh$

    - ▶ then proceed with *congestion avoidance*

- *NACK* (i.e., triple duplicate-ack)

    - ▶ set *ssthresh* $= \overline{W}/2$

    - ▶ cut $W$ in half: $W = \overline{W}/2$

    - ▶ run *congestion avoidance*, ramping up $W$ linearly

    - ▶ This is called *fast recovery*

Time

Time

SS=*slow start*  CA=*congestion avoidance*

■ Potentially *multiple paths* for the same source/destination

■ Potentially *multiple paths* for the same source/destination
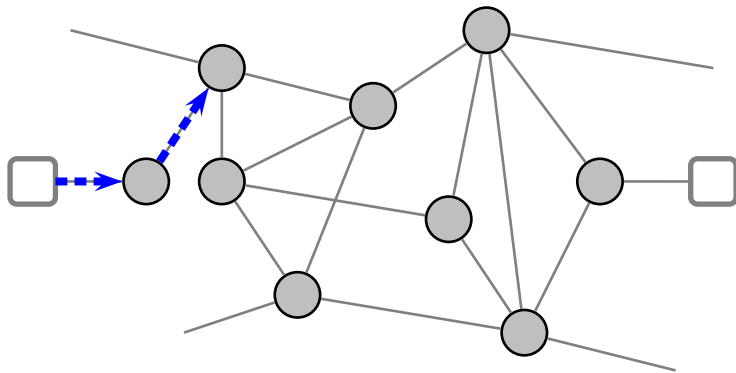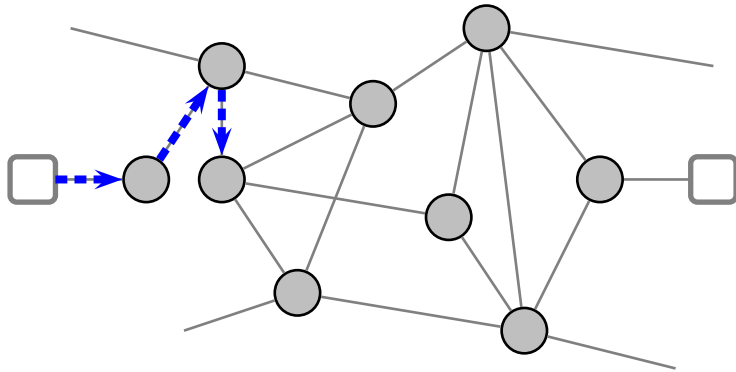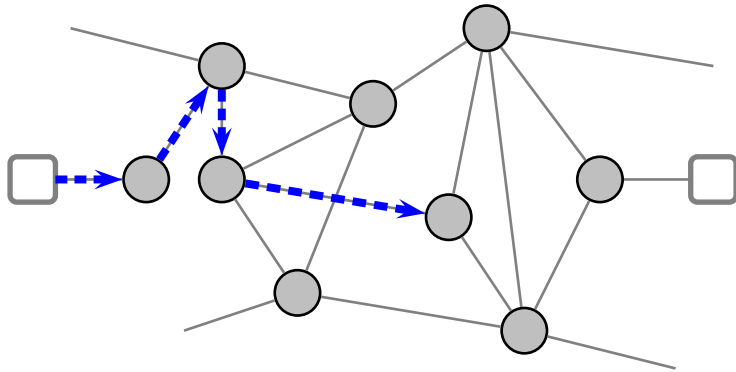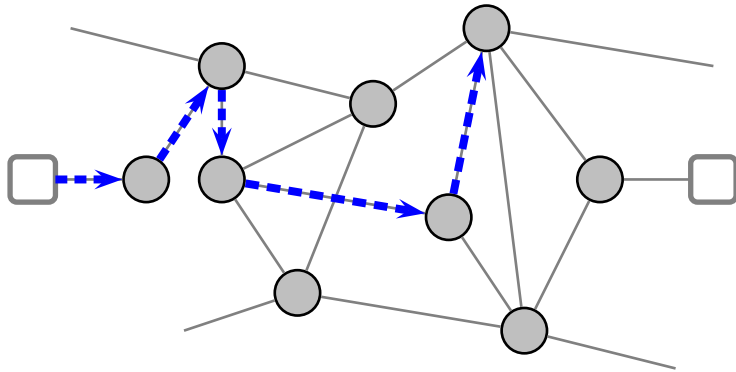
■ Potentially *multiple paths* for the same source/destination

- Potentially *multiple paths* for the same source/destination

■ Potentially *multiple paths* for the same source/destination

■ Potentially *multiple paths* for the same source/destination

- Potentially *multiple paths* for the same source/destination

- Potentially *multiple paths* for the same source/destination

- Potentially *multiple paths* for the same source/destination
- Potentially *asymmetric paths*
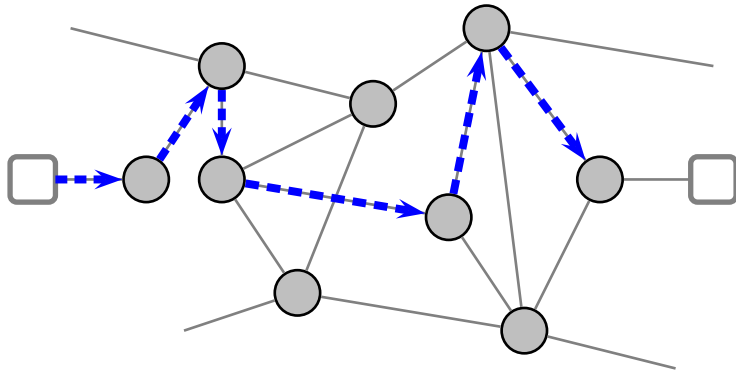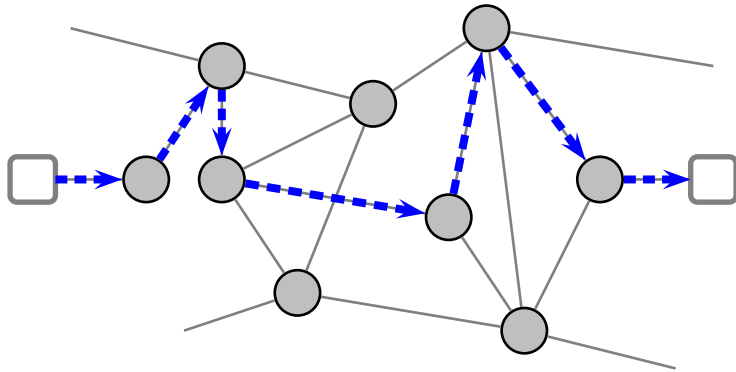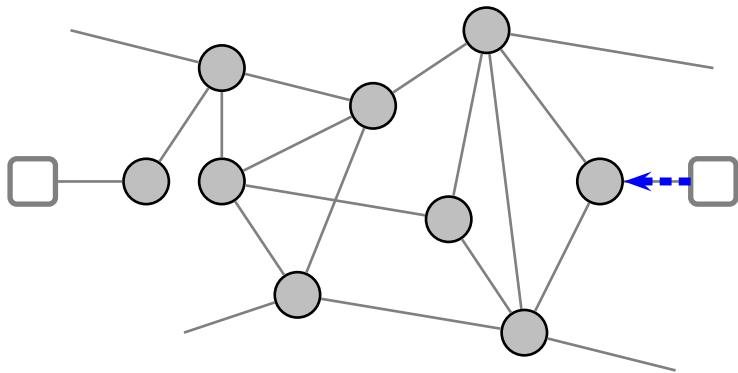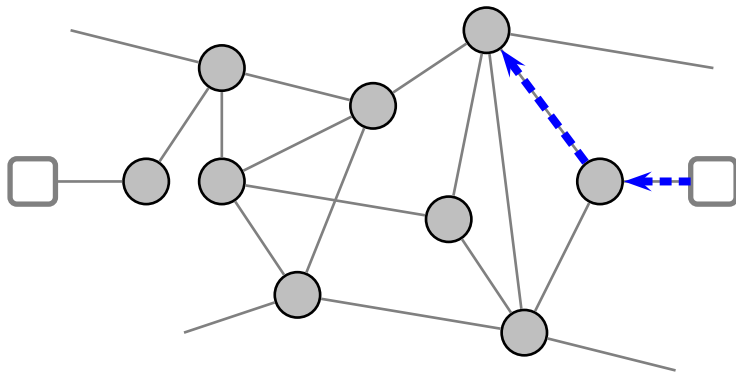
- Potentially *multiple paths* for the same source/destination
- Potentially *asymmetric paths*

- Potentially *multiple paths* for the same source/destination

- Potentially *asymmetric paths*

- Potentially *multiple paths* for the same source/destination
- Potentially *asymmetric paths*
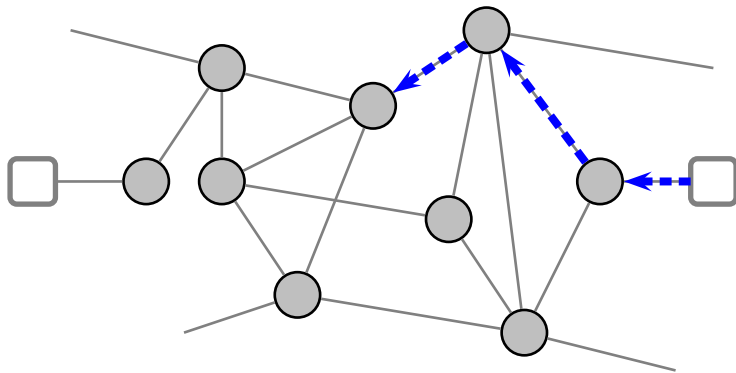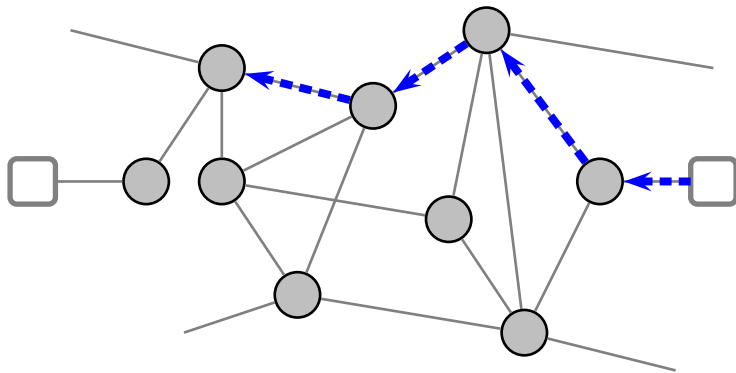
- Potentially *multiple paths* for the same source/destination
- Potentially *asymmetric paths*

- Potentially *multiple paths* for the same source/destination
- Potentially *asymmetric paths*

- *A* sends a datagram to *B*

- *A* sends a datagram to *B*

- The datagram is *forwarded* towards *B*

- *A* sends a datagram to *B*

- The datagram is *forwarded* towards *B*

A — k

B

to: B
…

2
1
3 — e — 4

| forwarding table | |
|---|---|
| dest. | output |
| … | … |
| B | port 4 |
| … | … |

- *Input:* datagram destination

- *Input:* datagram destination

- *Output:* output port

- *Input:* datagram destination

- *Output:* output port

- Simple design: "forwarding table"

- *Input:* datagram destination

- *Output:* output port

- Simple design: "forwarding table"

- Issues

- *Input:* datagram destination

- *Output:* output port

- Simple design: "forwarding table"

- Issues

  - ▶ how big is the forwarding table?

- *Input:* datagram destination

- *Output:* output port

- Simple design: "forwarding table"

- Issues

  - ▶ how big is the forwarding table?
  - ▶ how fast does the router have to forward datagrams?

- *Input:* datagram destination

- *Output:* output port

- Simple design: "forwarding table"

- Issues

  - ▶ how big is the forwarding table?
  - ▶ how fast does the router have to forward datagrams?
  - ▶ how does the router build and maintain the forwarding table?

| router k | |
|---|---|
| A | 2 |
| B | 1 |

- Finding paths through a network

■ Finding paths through a network

- Finding paths through a network



- Example: $a \rightarrow j$?

■ The network is modeled as a graph

$$G = (V, E)$$

■ The network is modeled as a graph

$$G = (V, E)$$

▶ *V* is a set of ***vertices*** representing the routers

# Graph Model

■ The network is modeled as a graph

$$G = (V, E)$$

- ▶ *V* is a set of ***vertices*** representing the routers

- ▶ $E \subseteq V \times V$ is a set of ***edges*** representing communication links
  - ▶ e.g., $(u, v) \in E$ iff router $u$ is on the same subnet as $v$

■ The network is modeled as a graph

$$G = (V, E)$$

- ▶ $V$ is a set of ***vertices*** representing the routers

- ▶ $E \subseteq V \times V$ is a set of ***edges*** representing communication links
  - ▶ e.g., $(u, v) \in E$ iff router $u$ is on the same subnet as $v$

- ▶ $G$ is assumed to be an ***undirected graph***, meaning that ***links are bidirectional***
  - ▶ i.e., $(u, v) \in E \Leftrightarrow (v, u) \in E$ for all $u, v \in N$

■ The network is modeled as a graph

$$G = (V, E)$$

▶ $V$ is a set of **vertices** representing the routers

▶ $E \subseteq V \times V$ is a set of **edges** representing communication links
   ▶ e.g., $(u, v) \in E$ iff router $u$ is on the same subnet as $v$

▶ $G$ is assumed to be an **undirected graph**, meaning that **links are bidirectional**
   ▶ i.e., $(u, v) \in E \Leftrightarrow (v, u) \in E$ for all $u, v \in N$

▶ A **cost** function $c : E \to \mathbb{R}$
   ▶ costs are always positive: $c(e) > 0$ for all $e \in E$
   ▶ links are symmetric: $c(u, v) = c(v, u)$ for all $u, v \in N$

- For every router $u \in V$, for every other router $v \in V$, compute the path
  $P_{u \to v} = u, x_1, x_2, \ldots, x_n, v$ such that

- For every router $u \in V$, for every other router $v \in V$, compute the path $P_{u \rightarrow v} = u, x_1, x_2, \ldots, x_n, v$ such that

  - $P_{u \rightarrow v}$ is completely contained in the network graph $G$. I.e., $(u, x_1) \in V, (x_1, x_2) \in V, \ldots, (x_n, v) \in V$

## Routing in the Graph Model

- For every router $u \in V$, for every other router $v \in V$, compute the path $P_{u \to v} = u, x_1, x_2, \ldots, x_n, v$ such that

  - $P_{u \to v}$ is completely contained in the network graph $G$. I.e., $(u, x_1) \in V, (x_1, x_2) \in V, \ldots, (x_n, v) \in V$

  - $P_{u \to v}$ is a *least-cost path*, where the cost of the path is $c(P_{u \to v}) = c(u, x_1) + c(x_1, x_2) + \ldots + c(x_n, v)$

# Routing in the Graph Model

- For every router $u \in V$, for every other router $v \in V$, compute the path $P_{u \to v} = u, x_1, x_2, \ldots, x_n, v$ such that

  - $P_{u \to v}$ is completely contained in the network graph $G$. I.e., $(u, x_1) \in V, (x_1, x_2) \in V, \ldots, (x_n, v) \in V$

  - $P_{u \to v}$ is a *least-cost path*, where the cost of the path is $c(P_{u \to v}) = c(u, x_1) + c(x_1, x_2) + \ldots + c(x_n, v)$

- Compile $u$'s forwarding table by adding the following entry:

$$A(v) \to I_u(x_1)$$

  - $A(v)$ is the address (or set of addresses) of router $v$
  - $I_u(x_1)$ is the interface that connects $u$ to the first next-hop router $x_1$ in $P_{u \to v} = u, x_1, x_2, \ldots, x_n, v$

- Example: $a \rightarrow j$

- Example: $a \rightarrow j$

  ▶ least-cost path is $P_{a \rightarrow j} = a, e, b, f, j$

- Example: $a \rightarrow j$

    ▶ least-cost path is $P_{a \rightarrow j} = a, e, b, f, j$

    ▶ $a$'s forwarding table will contain an entry $\boxed{j \rightarrow 2}$ since $I_a(e) = 2$

- There are two main strategies to implement a routing algorithm

- There are two main strategies to implement a routing algorithm

- ***Link-state routing***

- There are two main strategies to implement a routing algorithm

- ***Link-state routing***
  - ▶ global view of the network
  - ▶ local computation of least-cost paths

■ There are two main strategies to implement a routing algorithm

■ *Link-state routing*

    ► global view of the network

    ► local computation of least-cost paths

■ *Distance-vector routing*

■ There are two main strategies to implement a routing algorithm

■ *Link-state routing*

- ▶ global view of the network
- ▶ local computation of least-cost paths

■ *Distance-vector routing*

- ▶ local view of the network
- ▶ global computation of least-cost paths

- Router $u$ maintains a complete view of the network graph $G$ (including all links and their costs)

- Router *u* maintains a complete view of the network graph *G* (including all links and their costs)

    - ▶ every router *v* advertises its adjacent links (their costs) to every other router in the network; this information is called *link state*

    - ▶ ***link-state advertisements (LSAs)*** are broadcast through the entire network

- Router *u* maintains a complete view of the network graph *G* (including all links and their costs)

  ▶ every router *v* advertises its adjacent links (their costs) to every other router in the network; this information is called *link state*

  ▶ ***link-state advertisements (LSAs)*** are broadcast through the entire network

  ▶ routers collect link-state advertisements from other routers, and they use them to compile and maintain a complete view of *G*

- Router $u$ maintains a complete view of the network graph $G$ (including all links and their costs)

  - ▶ every router $v$ advertises its adjacent links (their costs) to every other router in the network; this information is called *link state*

  - ▶ ***link-state advertisements (LSAs)*** are broadcast through the entire network

  - ▶ routers collect link-state advertisements from other routers, and they use them to compile and maintain a complete view of $G$

- Using its local representation of $G$, router $u$ computes the least-cost paths from $u$ to every other router in the network

# Link-State Routing

- Router *u* maintains a complete view of the network graph *G* (including all links and their costs)

  - ▶ every router *v* advertises its adjacent links (their costs) to every other router in the network; this information is called *link state*

  - ▶ ***link-state advertisements (LSAs)*** are broadcast through the entire network

  - ▶ routers collect link-state advertisements from other routers, and they use them to compile and maintain a complete view of *G*

- Using its local representation of *G*, router *u* computes the least-cost paths from *u* to every other router in the network

  - ▶ the computation is local

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$
$LSA_h = \{(h, e, 1), (h, f, 4), (h, j, 14)\}$

**Link-State Advertisements**

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$
$LSA_h = \{(h, e, 1), (h, f, 4), (h, j, 14)\}$

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$
$LSA_h = \{(h, e, 1), (h, f, 4), (h, j, 14)\}$
$LSA_d = \{(d, a, 1), (d, g, 1), (d, e, 3)\}$

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$
$LSA_h = \{(h, e, 1), (h, f, 4), (h, j, 14)\}$
$LSA_d = \{(d, a, 1), (d, g, 1), (d, e, 3)\}$

## Link-State Advertisements

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$
$LSA_h = \{(h, e, 1), (h, f, 4), (h, j, 14)\}$
$LSA_d = \{(d, a, 1), (d, g, 1), (d, e, 3)\}$
$LSA_f = \{(f, c, 1), (f, b, 1), (f, e, 3), (f, h, 4), (f, j, 2)\}$

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$
$LSA_h = \{(h, e, 1), (h, f, 4), (h, j, 14)\}$
$LSA_d = \{(d, a, 1), (d, g, 1), (d, e, 3)\}$
$LSA_f = \{(f, c, 1), (f, b, 1), (f, e, 3), (f, h, 4), (f, j, 2)\}$

$LSA_a = \{(a, b, 3), (a, e, 1), (a, d, 1)\}$
$LSA_h = \{(h, e, 1), (h, f, 4), (h, j, 14)\}$
$LSA_d = \{(d, a, 1), (d, g, 1), (d, e, 3)\}$
$LSA_f = \{(f, c, 1), (f, b, 1), (f, e, 3), (f, h, 4), (f, j, 2)\}$
...

# Link-State Routing Ingredients

What do we need to implement link-state routing?

What do we need to implement link-state routing?

- Every router sends its LSA to every other router in the network, so we need a ***broadcast routing scheme***

# Link-State Routing Ingredients

What do we need to implement link-state routing?

- Every router sends its LSA to every other router in the network, so we need a ***broadcast routing scheme***

- Once we have all the LSAs from every router, and therefore we complete knowledge of *G*, we need an ***algorithm to compute least-cost paths in a graph***

- **■ *Flooding***
  - ▶ every router forwards a broadcast packet to every adjacent router, except the one that sent the packet

- ***Flooding***
  - ► every router forwards a broadcast packet to every adjacent router, except the one that sent the packet

- Simple and elegant

- *Flooding*
  - ► every router forwards a broadcast packet to every adjacent router, except the one that sent the packet

- Simple and elegant

- Correct w.r.t. the broadcast requirement: a broadcast packet will eventually reach every router

- *Flooding*
  - ▶ every router forwards a broadcast packet to every adjacent router, except the one that sent the packet

- Simple and elegant

- Correct w.r.t. the broadcast requirement: a broadcast packet will eventually reach every router

- Any problem with this solution?

- ***Flooding***
  - ▶ every router forwards a broadcast packet to every adjacent router, except the one that sent the packet

- Simple and elegant

- Correct w.r.t. the broadcast requirement: a broadcast packet will eventually reach every router

- Any problem with this solution?

  - ▶ cycles in the network create *packet storms*

■ *Reverse-path broadcast*

- ▶ every router forwards a broadcast packet to every adjacent router, except the one where it received the packet router

- ▶ a router *u* accepts a broadcast packet *p* originating at router *s* only if *p* arrives on the link that is on the direct (unicast) path from *u* to *s*

- ### *Reverse-path broadcast*
  - ▶ every router forwards a broadcast packet to every adjacent router, except the one where it received the packet router
  - ▶ a router *u* accepts a broadcast packet *p* originating at router *s* only if *p* arrives on the link that is on the direct (unicast) path from *u* to *s*

- Correct w.r.t. the broadcast requirement: a broadcast packet will eventually reach every router

- ■ *Reverse-path broadcast*
  - ▶ every router forwards a broadcast packet to every adjacent router, except the one where it received the packet router
  - ▶ a router $u$ accepts a broadcast packet $p$ originating at router $s$ only if $p$ arrives on the link that is on the direct (unicast) path from $u$ to $s$

- ■ Correct w.r.t. the broadcast requirement: a broadcast packet will eventually reach every router

- ■ No packet storms even in the presence of cycles in $G$

- ***Reverse-path broadcast***
  - ▶ every router forwards a broadcast packet to every adjacent router, except the one where it received the packet router
  - ▶ a router $u$ accepts a broadcast packet $p$ originating at router $s$ only if $p$ arrives on the link that is on the direct (unicast) path from $u$ to $s$

- Correct w.r.t. the broadcast requirement: a broadcast packet will eventually reach every router

- No packet storms even in the presence of cycles in $G$

- Any problem with this solution?

- **■ *Reverse-path broadcast***
  - ▶ every router forwards a broadcast packet to every adjacent router, except the one where it received the packet router
  - ▶ a router $u$ accepts a broadcast packet $p$ originating at router $s$ only if $p$ arrives on the link that is on the direct (unicast) path from $u$ to $s$

- ■ Correct w.r.t. the broadcast requirement: a broadcast packet will eventually reach every router

- ■ No packet storms even in the presence of cycles in $G$

- ■ Any problem with this solution?
  - ▶ it requires (unicast) routing information
  - ▶ so it is obviously useless to implement a routing algorithm

- *Sequence-number controlled flooding*

## ■ *Sequence-number controlled flooding*

- ▶ the originator $s$ of a broadcast packet marks the packet with a sequence number $n_s$

### ■ *Sequence-number controlled flooding*

- ▶ the originator $s$ of a broadcast packet marks the packet with a sequence number $n_s$

- ▶ every router $u$ stores the most recent sequence number seen from each source router. Let's assume that $u$ has seen sequence numbers from $s$ up to $n_s$

■ *Sequence-number controlled flooding*

▶ the originator $s$ of a broadcast packet marks the packet with a sequence number $n_s$

▶ every router $u$ stores the most recent sequence number seen from each source router. Let's assume that $u$ has seen sequence numbers from $s$ up to $n_s$

▶ a router accepts a broadcast packet $p$ originating at $s$ only if $p$ carries a sequence number $seq(p)$ that is higher than the most recent one seen from $s$: $seq(p) > n_s$

- **Sequence-number controlled flooding**
  - the originator $s$ of a broadcast packet marks the packet with a sequence number $n_s$
  - every router $u$ stores the most recent sequence number seen from each source router. Let's assume that $u$ has seen sequence numbers from $s$ up to $n_s$
  - a router accepts a broadcast packet $p$ originating at $s$ only if $p$ carries a sequence number $seq(p)$ that is higher than the most recent one seen from $s$: $seq(p) > n_s$
  - accepted packets are forwarded to every adjacent router, except the previous-hop router

■ *Sequence-number controlled flooding*

- ▶ the originator $s$ of a broadcast packet marks the packet with a sequence number $n_s$
- ▶ every router $u$ stores the most recent sequence number seen from each source router. Let's assume that $u$ has seen sequence numbers from $s$ up to $n_s$
- ▶ a router accepts a broadcast packet $p$ originating at $s$ only if $p$ carries a sequence number $seq(p)$ that is higher than the most recent one seen from $s$: $seq(p) > n_s$
- ▶ accepted packets are forwarded to every adjacent router, except the previous-hop router
- ▶ $u$ updates its table of sequence numbers $n_s \leftarrow seq(p)$

- *Scalability*
  - ▶ hundreds of millions of hosts in today's Internet

- *Scalability*
    - ▶ hundreds of millions of hosts in today's Internet
    - ▶ transmitting routing information (e.g., LSAs) would be too expensive

- *Scalability*
  - ▶ hundreds of millions of hosts in today's Internet
  - ▶ transmitting routing information (e.g., LSAs) would be too expensive
  - ▶ forwarding would also be too expensive

- *Scalability*
  - ▶ hundreds of millions of hosts in today's Internet
  - ▶ transmitting routing information (e.g., LSAs) would be too expensive
  - ▶ forwarding would also be too expensive

- *Administrative autonomy*

- *Scalability*
  - ▶ hundreds of millions of hosts in today's Internet
  - ▶ transmitting routing information (e.g., LSAs) would be too expensive
  - ▶ forwarding would also be too expensive

- *Administrative autonomy*
  - ▶ one organization might want to run a distance-vector routing protocol, while another might want to run a link-state protocol

- *Scalability*

  ▶ hundreds of millions of hosts in today's Internet

  ▶ transmitting routing information (e.g., LSAs) would be too expensive

  ▶ forwarding would also be too expensive

- *Administrative autonomy*

  ▶ one organization might want to run a distance-vector routing protocol, while another might want to run a link-state protocol

  ▶ an organization might not want to expose its internal network structure

- Today's Internet is organized in *autonomous systems (ASs)*
  - ▶ independent administrative domains

■ Today's Internet is organized in ***autonomous systems (ASs)***

  ▶ independent administrative domains

■ ***Gateway routers*** connect an autonomous system with other autonomous systems

- Today's Internet is organized in ***autonomous systems (ASs)***

  - ▶ independent administrative domains

- ***Gateway routers*** connect an autonomous system with other autonomous systems

- An *intra-autonomous system routing protocol* runs within an autonomous system (e.g., OSPF)

  - ▶ this protocol determines internal routes

    - ▶ internal router ↔ internal router
    - ▶ internal router ↔ gateway router
    - ▶ gateway router ↔ gateway router

- An *inter-autonomous system routing protocol* determines routing at the autonomous-system level

- An *inter-autonomous system routing protocol* determines routing at the autonomous-system level



At AS3:
AS1 →

■ An *inter-autonomous system routing protocol* determines routing at the autonomous-system level



At AS3:
AS1 → AS1;

- An *inter-autonomous system routing protocol* determines routing at the autonomous-system level



At AS3:
AS1 → AS1; AS2 →

■ An *inter-autonomous system routing protocol* determines routing at the autonomous-system level



At AS3:
AS1 → AS1; AS2 → AS2;

- An *inter-autonomous system routing protocol* determines routing at the autonomous-system level



At AS3:
AS1 → AS1; AS2 → AS2; AS4 →

■ An *inter-autonomous system routing protocol* determines routing at the autonomous-system level



At AS3:
AS1 → AS1; AS2 → AS2; AS4 → AS1.

# Hierarchical Routing

- All routers within an AS compute their *intra-AS* routing information
  - using an *intra-doman* routing protocol

# Hierarchical Routing

- All routers within an AS compute their *intra-AS* routing information
  - using an *intra-doman* routing protocol

- Gateway routers figure out *inter-AS* routing information
  - using an *inter-domain* routing protocol

- All routers within an AS compute their *intra-AS* routing information
  - using an *intra-doman* routing protocol

- Gateway routers figure out *inter-AS* routing information
  - using an *inter-domain* routing protocol

- *inter-AS* routing information is propagated within an AS
  - using an appropriate protocol

# Hierarchical Routing

- All routers within an AS compute their *intra-AS* routing information
  - using an *intra-doman* routing protocol

- Gateway routers figure out *inter-AS* routing information
  - using an *inter-domain* routing protocol

- *inter-AS* routing information is propagated within an AS
  - using an appropriate protocol

- Both *inter-AS* and *intra-AS* routing information is used to compile the forwarding tables

- Destinations within the same autonomous system are reached as usual

- Destinations within the same autonomous system are reached as usual

- What about a destination $x$ outside the autonomous system?

- Destinations within the same autonomous system are reached as usual

- What about a destination $x$ outside the autonomous system?
    - *inter-AS* information is used to figure out that $x$ is reachable through gateway $G_x$

■ Destinations within the same autonomous system are reached as usual

■ What about a destination $x$ outside the autonomous system?

  ▶ *inter-AS* information is used to figure out that $x$ is reachable through gateway $G_x$

  ▶ *intra-AS* information is used to figure out how to reach $G_x$ within the AS

■ Destinations within the same autonomous system are reached as usual

■ What about a destination $x$ outside the autonomous system?

  ▶ *inter-AS* information is used to figure out that $x$ is reachable through gateway $G_x$

  ▶ *intra-AS* information is used to figure out how to reach $G_x$ within the AS

  ▶ what if $x$ is reachable through multiple gateway routers $G_x, G'_x, \ldots$?

■ Destinations within the same autonomous system are reached as usual

■ What about a destination $x$ outside the autonomous system?

▶ *inter-AS* information is used to figure out that $x$ is reachable through gateway $G_x$

▶ *intra-AS* information is used to figure out how to reach $G_x$ within the AS

▶ what if $x$ is reachable through multiple gateway routers $G_x$, $G'_x$, . . .?

▶ use *intra-AS* routing information to determine the costs of the (least-cost) paths to $G_x$, $G'_x$, . . .

▶ "hot-potato" routing: send it through the closest gateway

- *Administrative autonomy*

- *Administrative autonomy*
  - ► each autonomous system decides what intra-AS routing to use

- *Administrative autonomy*
    - ▶ each autonomous system decides what intra-AS routing to use
    - ▶ an autonomous system needs to expose only minimal information about the internal structure of its network
        - ▶ essentially only (sub)net addresses

- *Administrative autonomy*

  - ▶ each autonomous system decides what intra-AS routing to use

  - ▶ an autonomous system needs to expose only minimal information about the internal structure of its network

    - ▶ essentially only (sub)net addresses

- *Scalability*

# Benefits of Hierarchical Routing

- *Administrative autonomy*
  - ▶ each autonomous system decides what intra-AS routing to use
  - ▶ an autonomous system needs to expose only minimal information about the internal structure of its network
    - ▶ essentially only (sub)net addresses

- *Scalability*
  - ▶ routers within an autonomous system need to know very little about the internal structure of other autonomous systems

- *Administrative autonomy*

  - ▶ each autonomous system decides what intra-AS routing to use

  - ▶ an autonomous system needs to expose only minimal information about the internal structure of its network

    - ▶ essentially only (sub)net addresses

- *Scalability*

  - ▶ routers within an autonomous system need to know very little about the internal structure of other autonomous systems

    - ▶ essentially only (sub)net addresses

- *Administrative autonomy*
    - ▶ each autonomous system decides what intra-AS routing to use
    - ▶ an autonomous system needs to expose only minimal information about the internal structure of its network
        - ▶ essentially only (sub)net addresses

- *Scalability*
    - ▶ routers within an autonomous system need to know very little about the internal structure of other autonomous systems
        - ▶ essentially only (sub)net addresses

- External subnet addresses are likely to "aggregate" in groups that admit compact representations
    - ▶ this process is called *supernetting*

- The ***Border Gateway Protocol (BGP)*** is the inter-AS routing protocol in today's Internet

- The *Border Gateway Protocol (BGP)* is the inter-AS routing protocol in today's Internet

    - provides reachability information from neighbor ASs

# Inter-AS Routing in the Internet

- The ***Border Gateway Protocol (BGP)*** is the inter-AS routing protocol in today's Internet
  - ▶ provides reachability information from neighbor ASs
  - ▶ transmits reachability information to all internal routers within an AS

- The *Border Gateway Protocol (BGP)* is the inter-AS routing protocol in today's Internet

  - provides reachability information from neighbor ASs

  - transmits reachability information to all internal routers within an AS

  - determines good routes to all outside subnets

- The **_Border Gateway Protocol (BGP)_** is the inter-AS routing protocol in today's Internet

  - ▶ provides reachability information from neighbor ASs

  - ▶ transmits reachability information to all internal routers within an AS

  - ▶ determines good routes to all outside subnets

    - ▶ based on reachability information

# Inter-AS Routing in the Internet

- The ***Border Gateway Protocol (BGP)*** is the inter-AS routing protocol in today's Internet
    - provides reachability information from neighbor ASs
    - transmits reachability information to all internal routers within an AS
    - determines good routes to all outside subnets
        - based on reachability information
        - based on *policies*

# Inter-AS Routing in the Internet

- The ***Border Gateway Protocol (BGP)*** is the inter-AS routing protocol in today's Internet

  - ▶ provides reachability information from neighbor ASs

  - ▶ transmits reachability information to all internal routers within an AS

  - ▶ determines good routes to all outside subnets
    - ▶ based on reachability information
    - ▶ based on *policies*

  - ▶ BGP is a *path-vector* protocol

- **BGP session:** a semi-permanent connection between two routers

# BGP Architecture and Terminology

- *BGP session:* a semi-permanent connection between two routers

- *BGP peers:* two routers engaged in a BGP session
    - ▶ BGP sessions are established over TCP

# BGP Architecture and Terminology

- *BGP session:* a semi-permanent connection between two routers

- *BGP peers:* two routers engaged in a BGP session
  - ▶ BGP sessions are established over TCP

- *BGP external session (eBGP):* a session across two autonomous systems

# BGP Architecture and Terminology

- **BGP session:** a semi-permanent connection between two routers

- **BGP peers:** two routers engaged in a BGP session
  - ▶ BGP sessions are established over TCP

- **BGP external session (eBGP):** a session across two autonomous systems

- **BGP internal session (iBGP):** a session within an autonomous system
  - ▶ note that internal sessions carry *inter-AS* information
  - ▶ *intra-AS* routing uses a separate protocol (e.g., OSPF)

- **■ *BGP advertisement:*** a router advertises routes to networks, much like an entry in a distance-vector

  - ▶ destinations are denoted by address *prefixes*

- *BGP advertisement:* a router advertises routes to networks, much like an entry in a distance-vector

  ▶ destinations are denoted by address *prefixes*

  ▶ an AS may or may not forward an advertisement for a foreign network; doing so means being willing to carry traffic for that network

- *BGP advertisement:* a router advertises routes to networks, much like an entry in a distance-vector

  ▶ destinations are denoted by address *prefixes*

  ▶ an AS may or may not forward an advertisement for a foreign network; doing so means being willing to carry traffic for that network

  ▶ this is where a router may aggregate prefixes (a.k.a., "supernetting")
    E.g.,

$$\left. \begin{array}{l} 128.138.242.0/24 \\ 128.138.243.0/24 \end{array} \right\} \rightarrow 128.138.242.0/23$$

■ *BGP advertisement:* a router advertises routes to networks, much like an entry in a distance-vector

  ▶ destinations are denoted by address *prefixes*

  ▶ an AS may or may not forward an advertisement for a foreign network; doing so means being willing to carry traffic for that network

  ▶ this is where a router may aggregate prefixes (a.k.a., "supernetting")
    E.g.,

$$\left.\begin{array}{l} 128.138.242.0/24 \\ 128.138.243.0/24 \end{array}\right\} \rightarrow 128.138.242.0/23$$

$$\left.\begin{array}{l} 191.224.128.0/22 \\ 191.224.136.0/21 \\ 191.224.132.0/22 \end{array}\right\} \rightarrow$$

- **BGP advertisement:** a router advertises routes to networks, much like an entry in a distance-vector

  - ▶ destinations are denoted by address *prefixes*

  - ▶ an AS may or may not forward an advertisement for a foreign network; doing so means being willing to carry traffic for that network

  - ▶ this is where a router may aggregate prefixes (a.k.a., "supernetting")
    E.g.,

$$
\left.
\begin{array}{l}
128.138.242.0/24 \\
128.138.243.0/24
\end{array}
\right\} \rightarrow 128.138.242.0/23
$$

$$
\left.
\begin{array}{l}
191.224.128.0/22 \\
191.224.136.0/21 \\
191.224.132.0/22
\end{array}
\right\} \rightarrow 191.224.128.0/20
$$

- **_Autonomous system number (ASN):_** a unique identifier for each AS (with more than one gateway)

- *Autonomous system number (ASN):* a unique identifier for each AS (with more than one gateway)

- *BGP attributes:* a route advertisement includes a number of attributes
  - *AS-PATH:* sequence of ASNs through which the advertisement has been sent

- ***Autonomous system number (ASN):*** a unique identifier for each AS (with more than one gateway)

- ***BGP attributes:*** a route advertisement includes a number of attributes

  - *AS-PATH:* sequence of ASNs through which the advertisement has been sent
  - *NEXT-HOP:* specifies the interface (IP address) to use to forward packets towards the advertised destination
    - used to resolve ambiguous cases where an AS can be reached through multiple gateways (interfaces)

- *Autonomous system number (ASN):* a unique identifier for each AS (with more than one gateway)

- *BGP attributes:* a route advertisement includes a number of attributes
  - *AS-PATH:* sequence of ASNs through which the advertisement has been sent
  - *NEXT-HOP:* specifies the interface (IP address) to use to forward packets towards the advertised destination
    - used to resolve ambiguous cases where an AS can be reached through multiple gateways (interfaces)

- *BGP import policy:* used to decide whether to accept or reject the route advertisement
  - e.g., a router may not want to send its traffic through one of the AS listed in *AS-PATH*

1. Router preference: routes are ranked according to a *preference* value
   - ▶ configured at the router
   - ▶ or learned from another router within the same AS
   - ▶ essentially a configuration parameter for the AS

1. Router preference: routes are ranked according to a *preference* value
   - ▶ configured at the router
   - ▶ or learned from another router within the same AS
   - ▶ essentially a configuration parameter for the AS

2. Shortest AS-PATH

1. Router preference: routes are ranked according to a *preference* value
   - ▶ configured at the router
   - ▶ or learned from another router within the same AS
   - ▶ essentially a configuration parameter for the AS

2. Shortest AS-PATH

3. Closest NEXT-HOP router

1. Router preference: routes are ranked according to a *preference* value
   - ▶ configured at the router
   - ▶ or learned from another router within the same AS
   - ▶ essentially a configuration parameter for the AS

2. Shortest AS-PATH

3. Closest NEXT-HOP router

4. …