

# Transport Layer II

Prof. Ling Li | Dr. Nadith Pathirage | Lecture 07

Semester 1, 2021





# Flow and **Congestion** Control

- Flow Control
  - Transport Layer
  - Data Link Layer
- Congestion Control
  - Warning-bit
  - Choke-packets
  - Load-shedding
  - RED

# Flow Control – Data Link Layer

- Common flow and error control **techniques at the Data Link Layer:**

1. **Stop-and-Wait** ARQ

2. **Go-back-N** ARQ

3. **Select-Reject** ARQ

} Sliding Window  
Methods

# Flow Control – Transport Layer

## ■ Why limit flow?

- Source may send frames/packets at a rate that is faster than the processing speed of the destination host
- The buffer at the destination could be full because
  - ✓ Higher level protocol is not ready
  - ✓ Outgoing I/O port not ready
  - ✓ Destination protocols cannot process PDU as fast



**Prevents  
Buffer  
Overflow**

“limits the amount of data being transmitted so that destination host can ‘cope’ ”

# Flow Control – Transport Layer

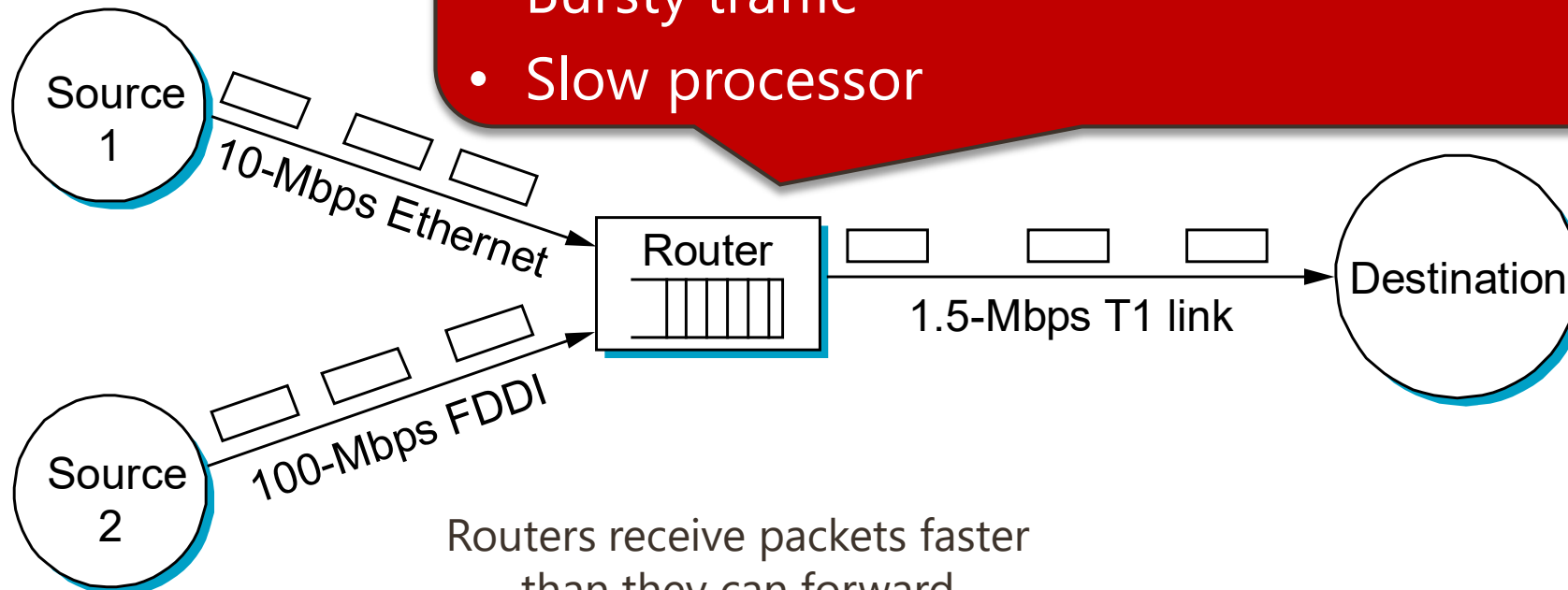
---

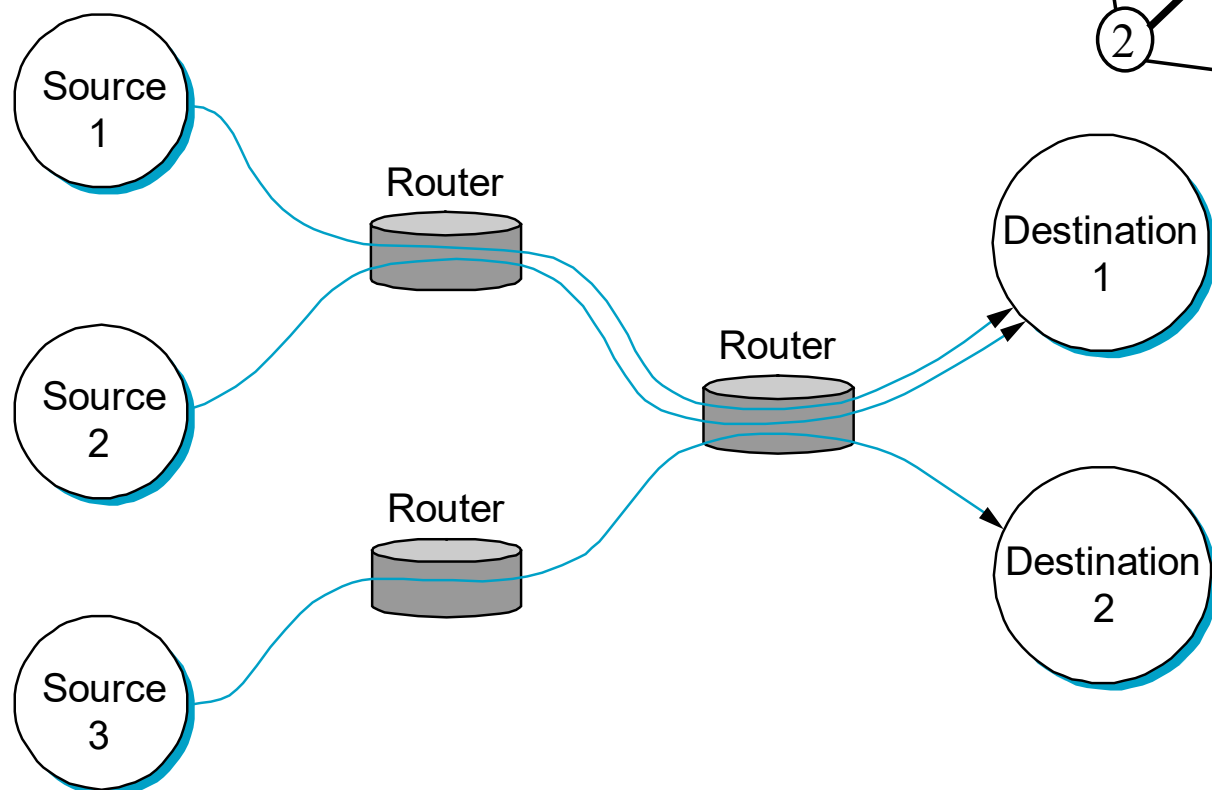
- Also used Sliding Window methods
- **Decouple Acknowledgement with available buffer**
- E.g., TCP implementation

# Congestion

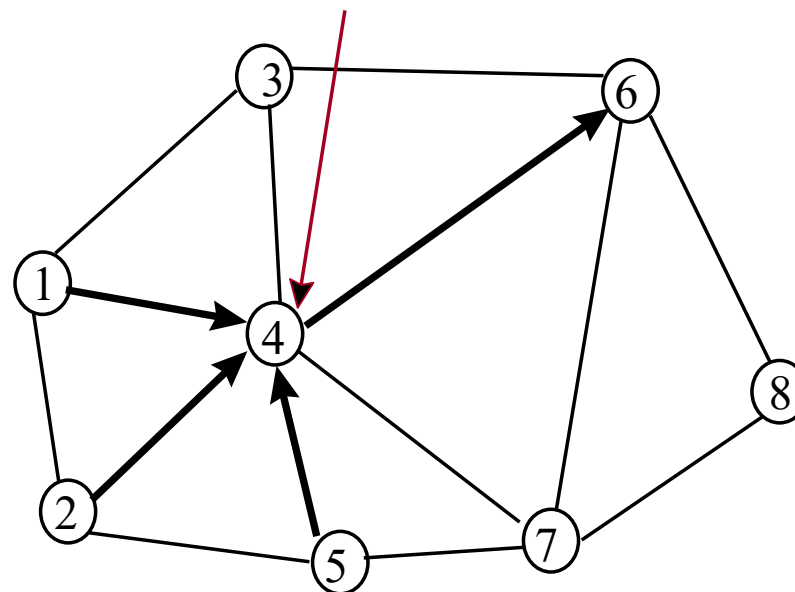
- One part of the subnet (e.g. one or more routers in an area) becomes overloaded, **congestion** occurs

- Packet arrival rate exceeds the outgoing link capacity
- Insufficient memory to store arriving packets
- Bursty traffic
- Slow processor





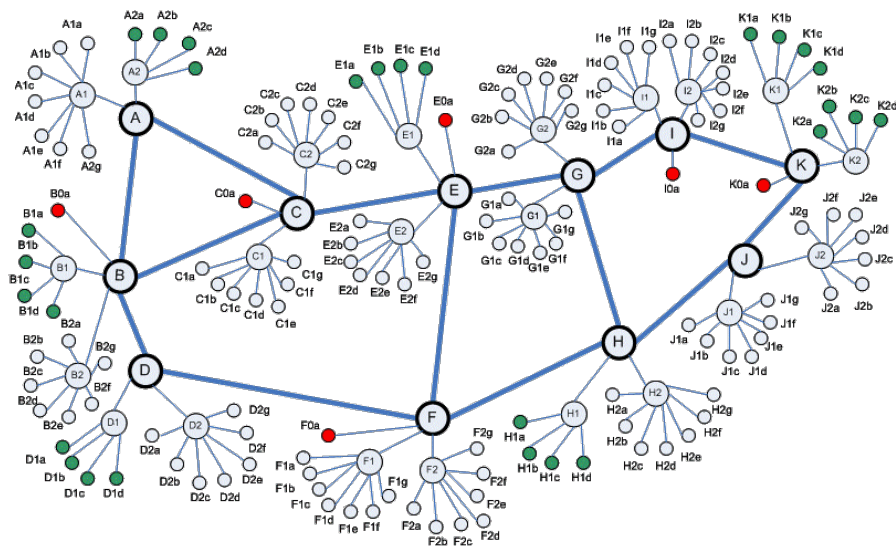
## Congestion



Congestion control can be distinguished from **routing** in that sometimes there is no way to *'route around'* a congested router.

# Flow vs Congestion Control

## ■ Congestion Control - Global Issue



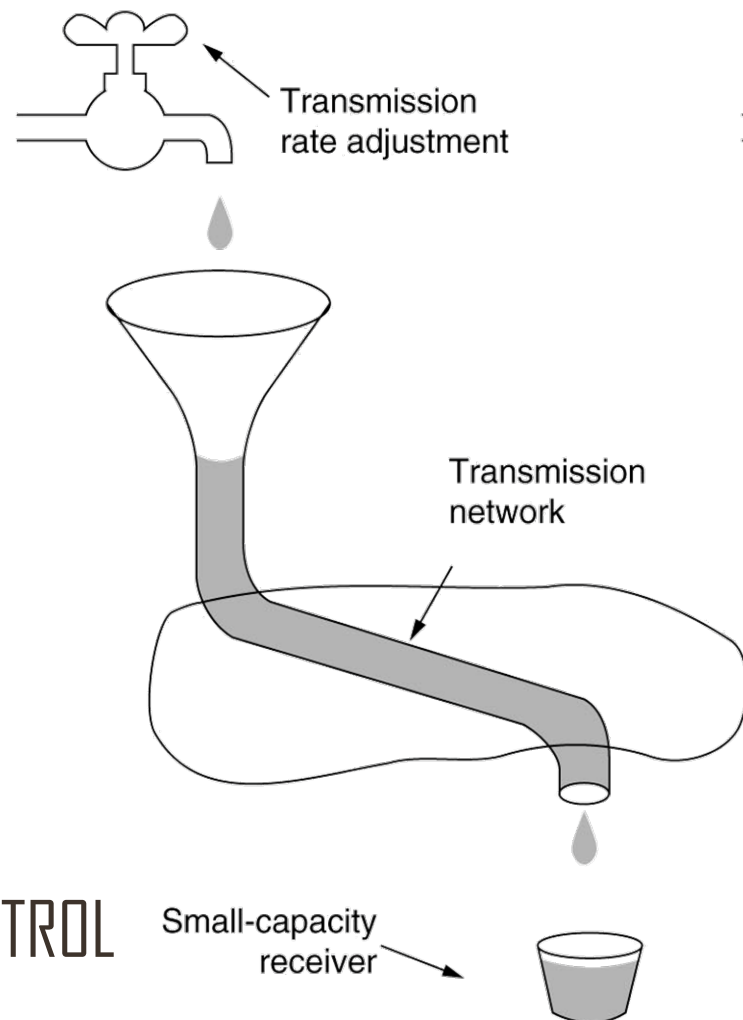
Make sure the subnet is able to carry the offered traffic; **involves** all the **host**, **routers**, **store-and-forwarding** processing, etc. within the subnet

1

- **Flow Control** – Scope is **point-to-point**;
- involves just **sender** and **receiver**

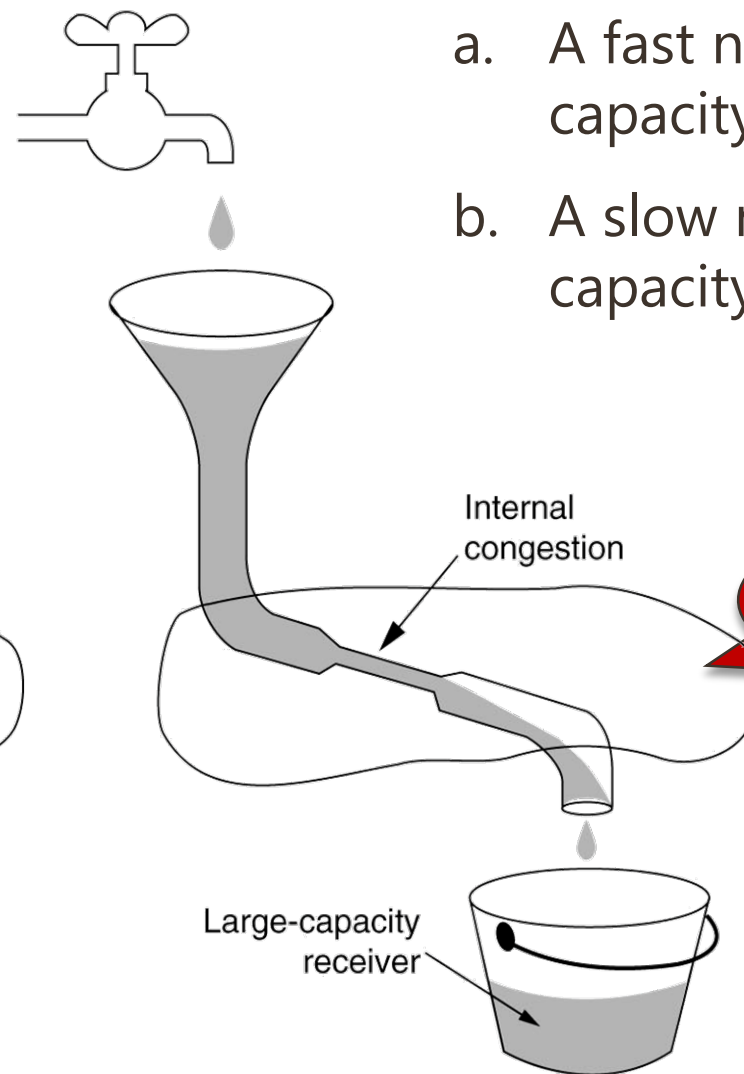


# Flow vs Congestion Control



FLOW CONTROL

(a)



Congestion CONTROL

- a. A fast network feeding a low-capacity receiver
- b. A slow network feeding a high-capacity receiver

reduced throughput  
long response time

# Congestion Control - Solutions

- Congestion Control is concerned with efficiently using a network at high load
- Several Solutions:



**1. Warning bit**

**2. Choke Packets**

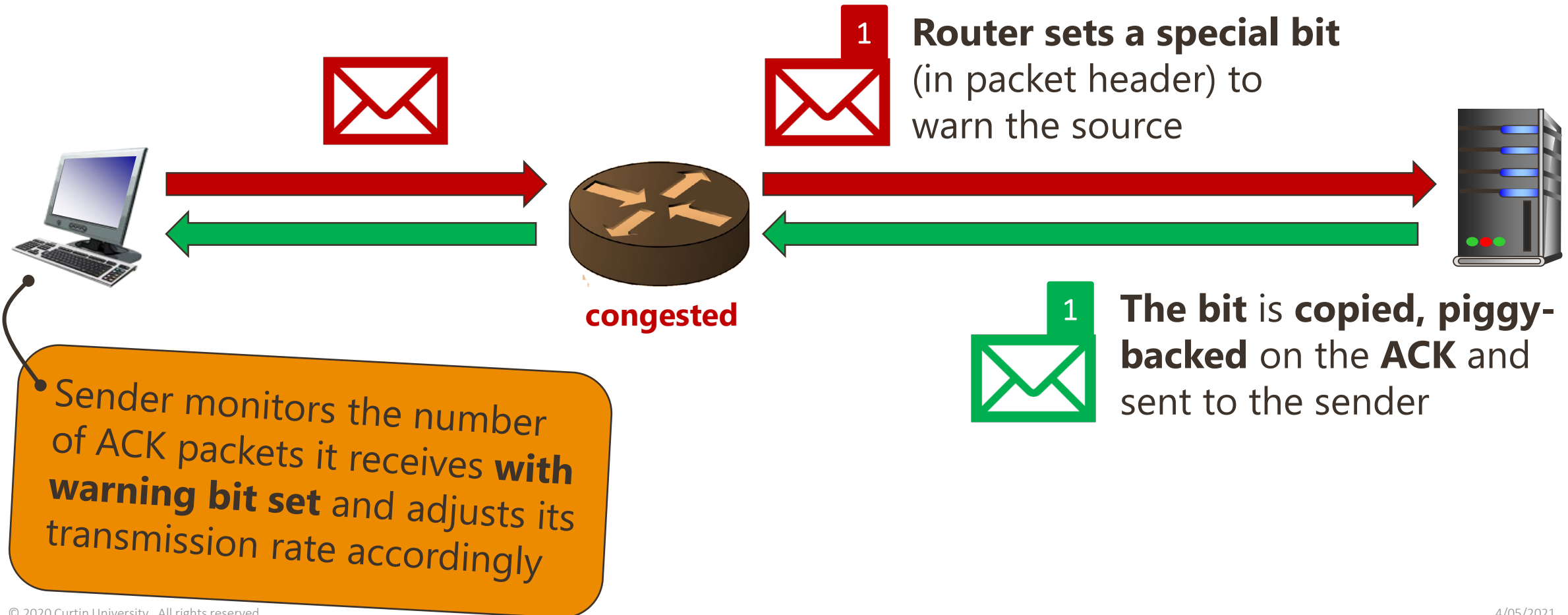
**3. Load Shedding**

**Congestion  
Detection &  
Recovery**

**4. Random Early Discard (RED)**

**Congestion  
Avoidance**



# 1. Warning Bit




## 2. Choke Packets

- A **more direct way** of telling the **source to slow down**

**ICMP Source  
Quench Packet**

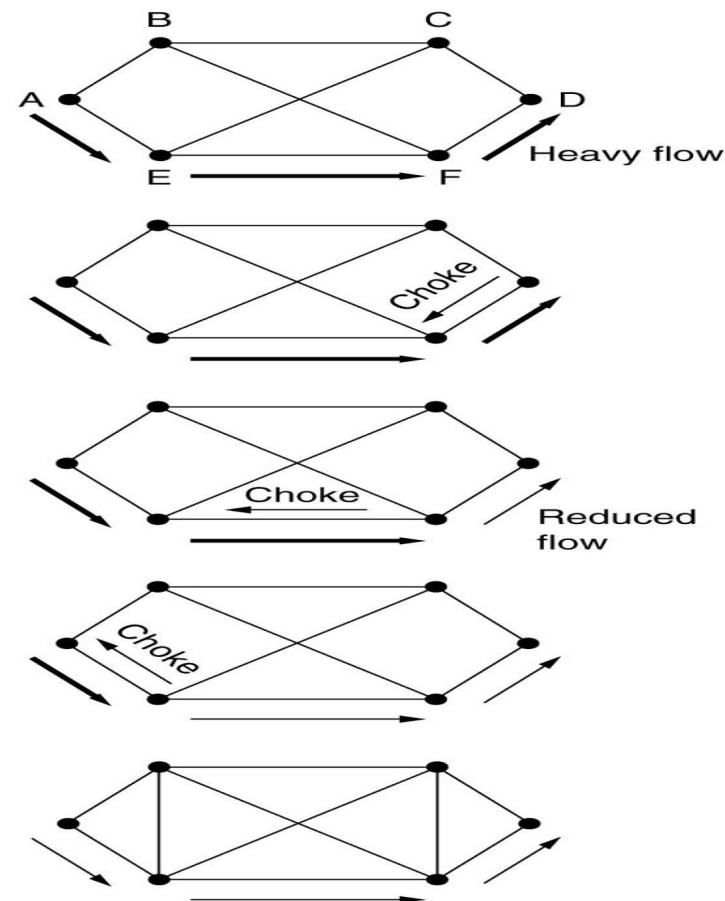
- A choke packet is a **control packet** generated at a congested node, transmitted to restrict traffic flow
- On receiving the choke packet, 
  - Source will **reduce transmission rate** by a 

## 2. Choke Packets - Hop-by-hop

- **Over long distances or at high speeds choke packets are not very effective.**
  - A more efficient method: **Send the choke packets hop-by-hop**
- 
- This requires **each hop to reduce** its transmission even before the choke packet arrive at the source

Hop-by-Hop



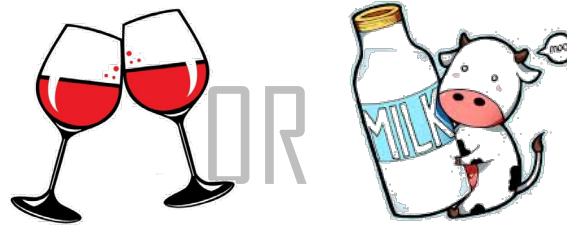


- a. A choke packet that **affects only the source**
- b. A choke packet that **affects each hop** it passes through (Tanenbaum)

# 3. Load Shedding

- When **buffers become full**, routers simply **discard packets**
- Which packet is chosen to be the victim depends on:

## 1. Wine or Milk policy

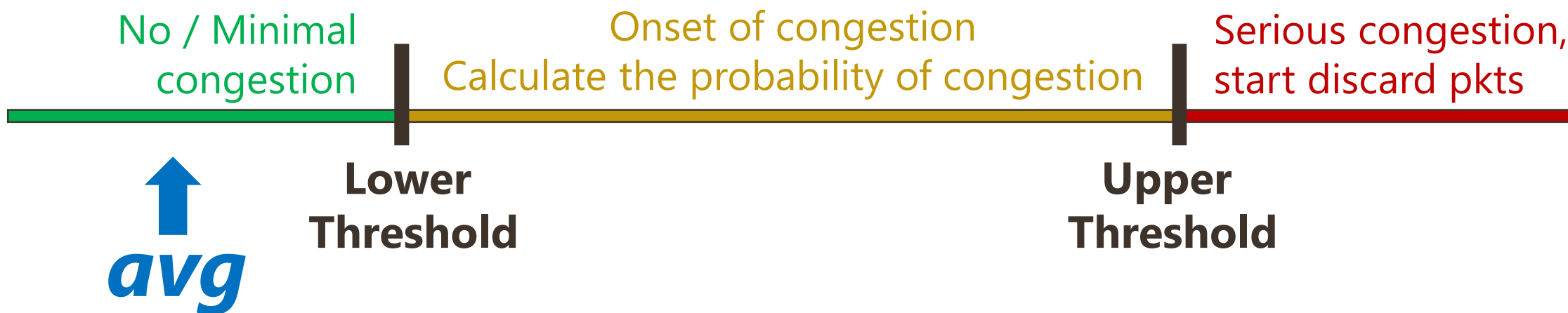


- ✗ File transfer cannot discard older packets: *will cause a gap in the received data.*
- ✓ Real-time voice or video may throw away old data and keep new packets.

## 2. Alternatively, implement an **Intelligent Discard Policy** or get the application to mark packets with discard priority

## 4. Random Early Discard (RED)

- This is a proactive approach in which the router **discards** one or more packets *before the buffer becomes completely full*
- Each time a packet arrives, the **RED algorithm** computes the average **queue length, *avg***





# Transport Control Protocol (TCP)

- Fundamentals
- TCP Header
  - Flags (SYN, FIN, ACK, RST)
  - Flag (URG, PSH) – in depth
  - TCP Options
  - Window Size (Dynamic Buffer Management)
- TCP Flow Control
  - Dynamic Buffer Management

# TCP

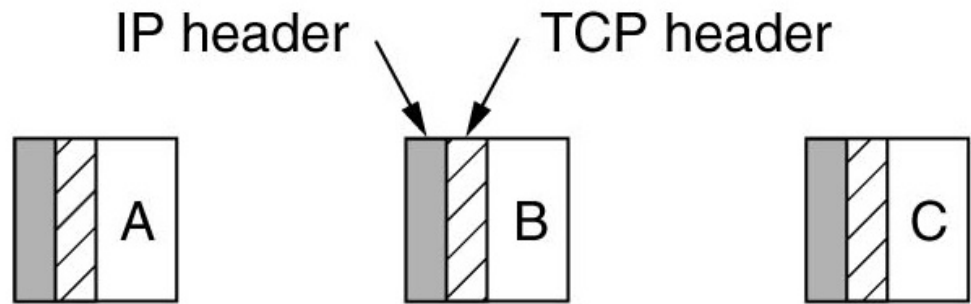
- The (other) **main** transport **protocol** used in the **Internet**
  - ✓ **Connection-oriented** protocol
  - ✓ RFC 793 (formal), RFC 1122 & 1323 (bug fixes)
  - ✓ Provide a reliable end-to-end communication over an unreliable internetwork
- **Connections** are:
  - ✓ **Full duplex** and point-to-point
  - ✓ A **byte stream** not a message stream



**No support for Multicasting or Broadcasting !**

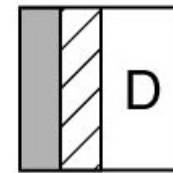


# TCP: Header



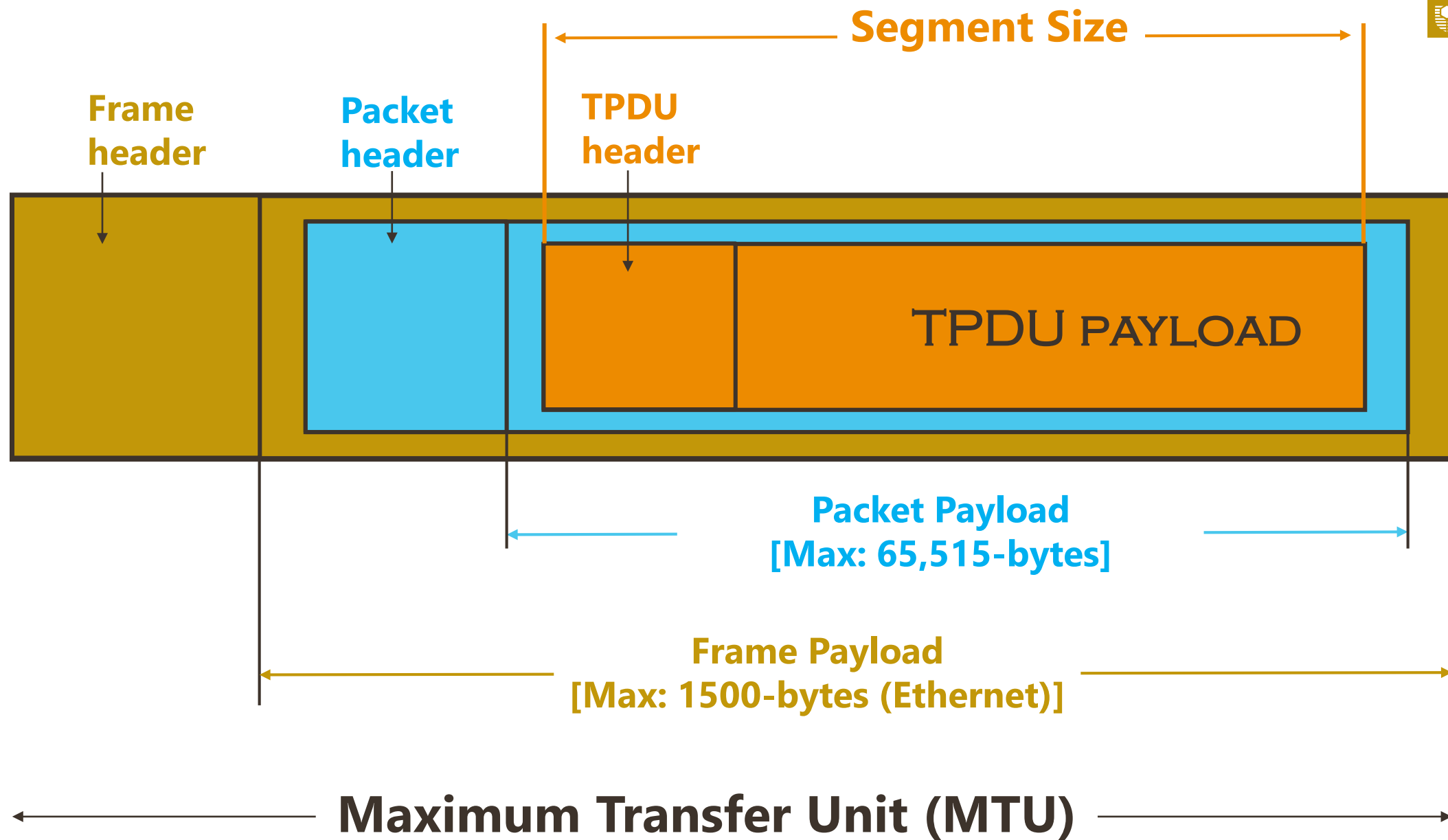
(a)

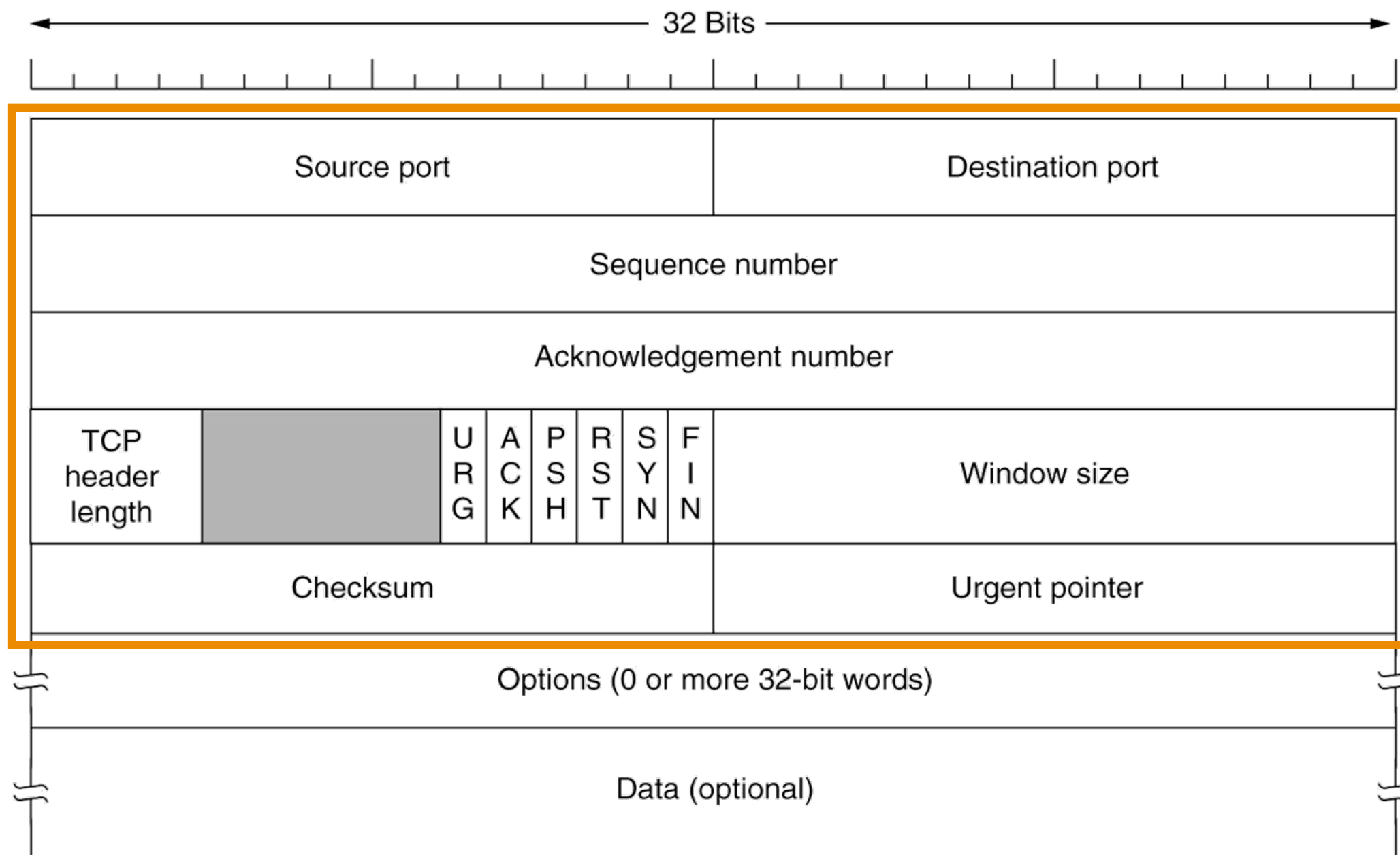
**512-byte segments**  
sent as a separated IP  
datagrams



(b)

**2048-bytes of data**  
received to the application  
in a single **READ** call





# TCP HEADER

**Fixed  
20-bytes  
header**

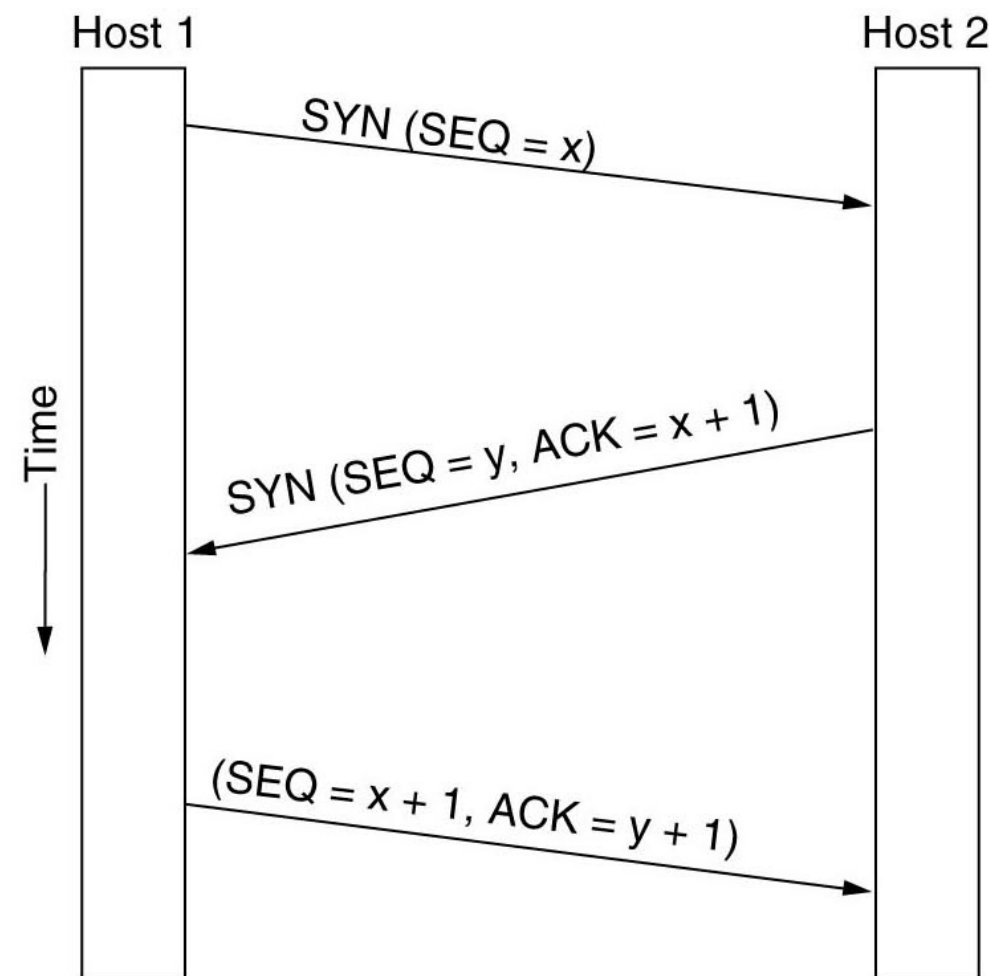
# TCP: Header – cont.

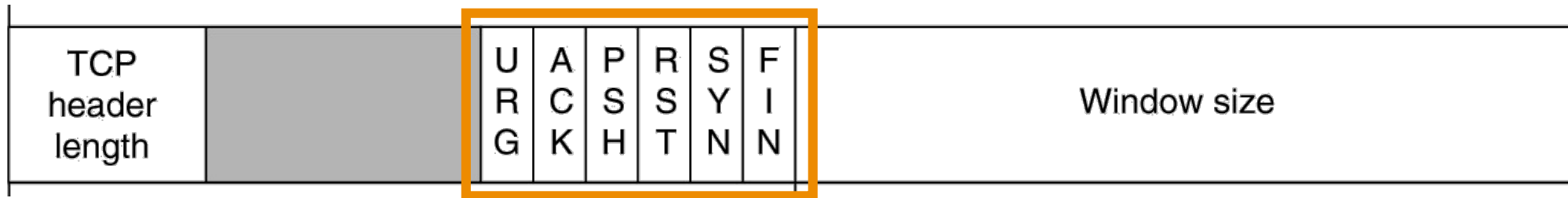
## ▪ Sequence Number (32-bit)

- sequence number of the first data octet (byte) in this segment
- if SYN is set this field is the initial sequence number (ISN) and the first data octet is ISN+1

## ▪ Acknowledgement Number (32-bit)

- sequence number of the next data octet the TCP entity expects to receive. May be piggybacked!!
- NOTE that TCP is byte or stream oriented.

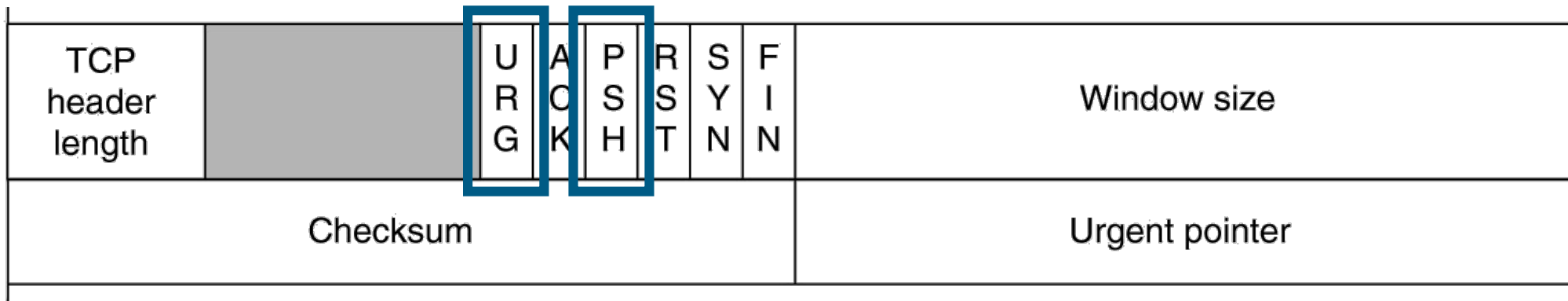




- ✓ **URG:** Urgent pointer field significant. Inform the destination TCP user that 'urgent' data is arriving.
- ✓ **ACK:** Acknowledgement field significant.
- ✓ **PSH:** Push function. A TCP user can require TCP to send (receive) all outstanding data up to and including that labelled with a **PUSH** flag.
- ✓ **RST:** Reset the connection.
- ✓ **SYN:** Synchronize the sequence numbers. Used to establish connections.
- ✓ **FIN:** No more data from sender. Used to release connections.

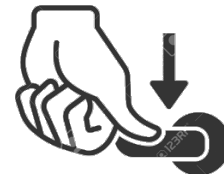
## FLAGS



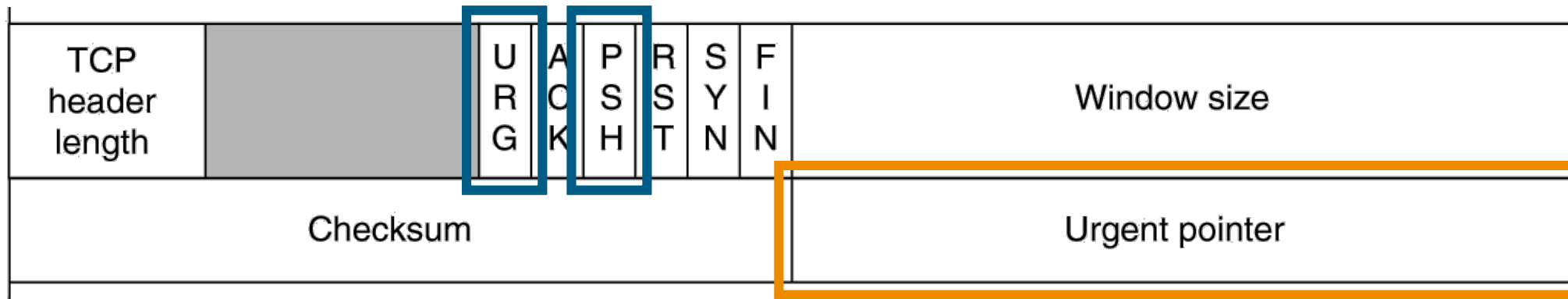


When application passes data to TCP, TCP may **send** it **immediately** or **buffer it (at both sides)**  
(in order to collect a larger amount to send at once)

## FLAGS

- **PUSH Flag** – (used by application) -> Force/Flush data out  **at both sides!**
- **URGENT Flag** – (used by application) -> Cause TCP to stop accumulating data and transmit everything it has for the connection immediately

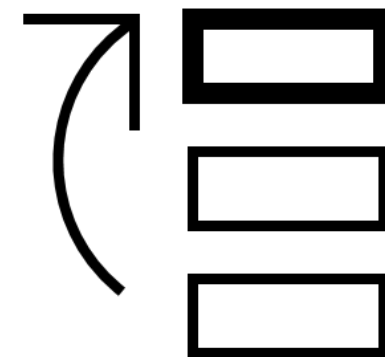
when urgent data is received, receiving application is **interrupted** !  
(stops whatever it was doing) & read the data stream to find the urgent data

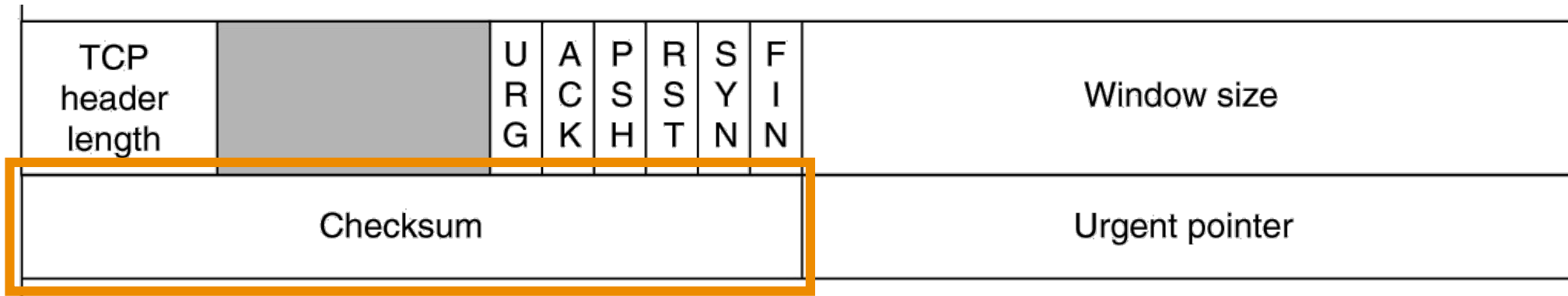


## ■ Urgent Pointer

- ✓ **Points to the last octet** in a sequence of urgent data. Allows the receiver to know how much urgent data is coming

# URG POINTER



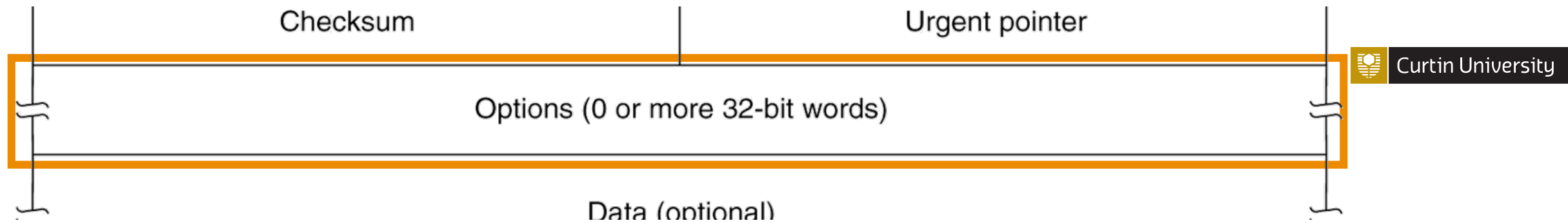


## ▪ Checksum:

- ✓ header
- ✓ data
- ✓ conceptual **pseudo-header**

**CHECKSUM**

**source & destination addresses, segment length.**  
**This provides protection from mis-delivery**



## ✓ Maximum Segment Size (MSS)

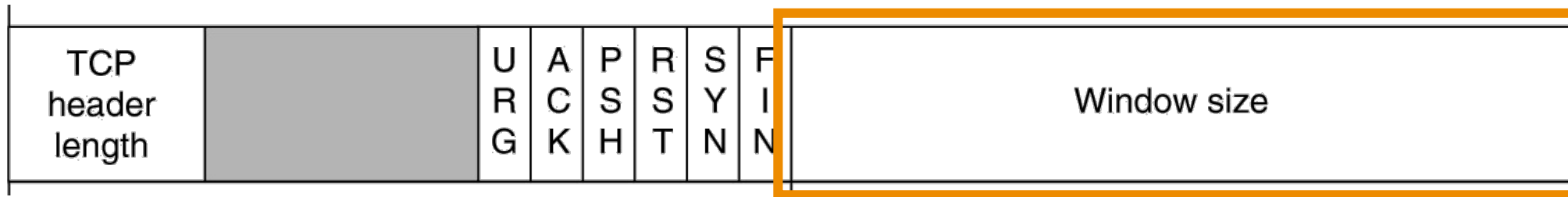
## ✓ Window Scale Factor:

Window is multiplied by  $2^F$  where  $F$  is the window scale factor.

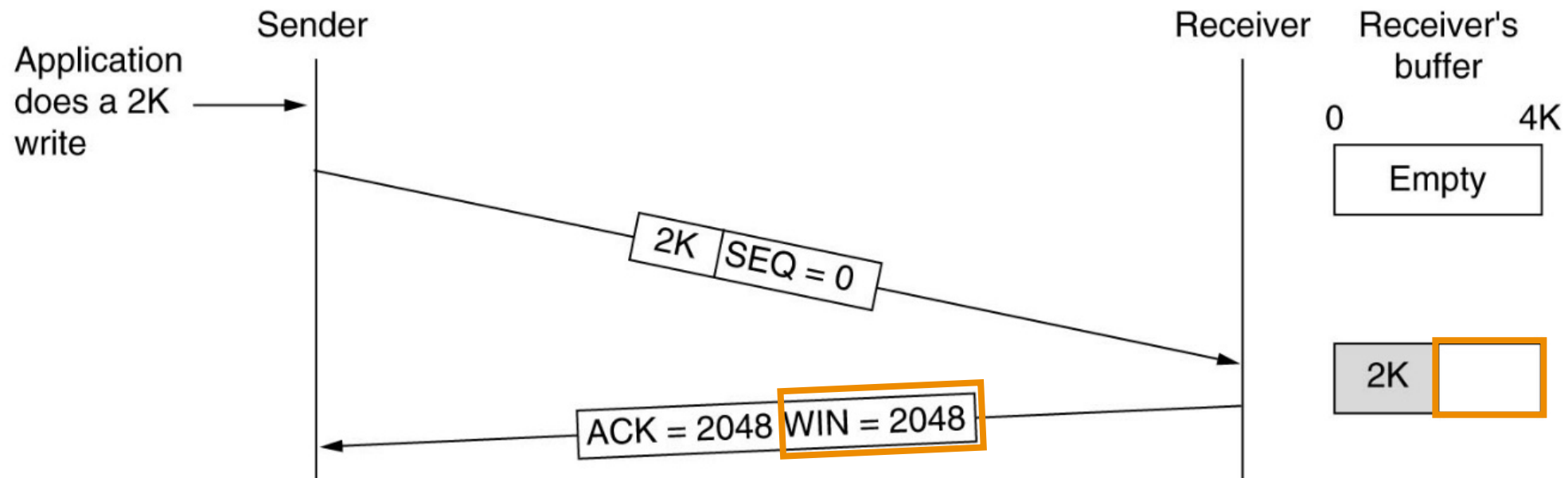
## ✓ Timestamp

Any outgoing packet with a timestamp will cause the ACK to carry a timestamp echo with the same value. Can be used to calculate round trip time

**OPTIONS**



- ✓ Used in flow control
- ✓ **Variable-sized** Sliding Window

**WND SIZE**



# Window Size

---

- **Window Management:** not tied to ACKs  
*(differs from data link protocols)*
- Receive entity will **advertise** window segment that it can receive & buffer
- Each entity can **alter size** of the other's sending window **dynamically** using the segment's Window field.

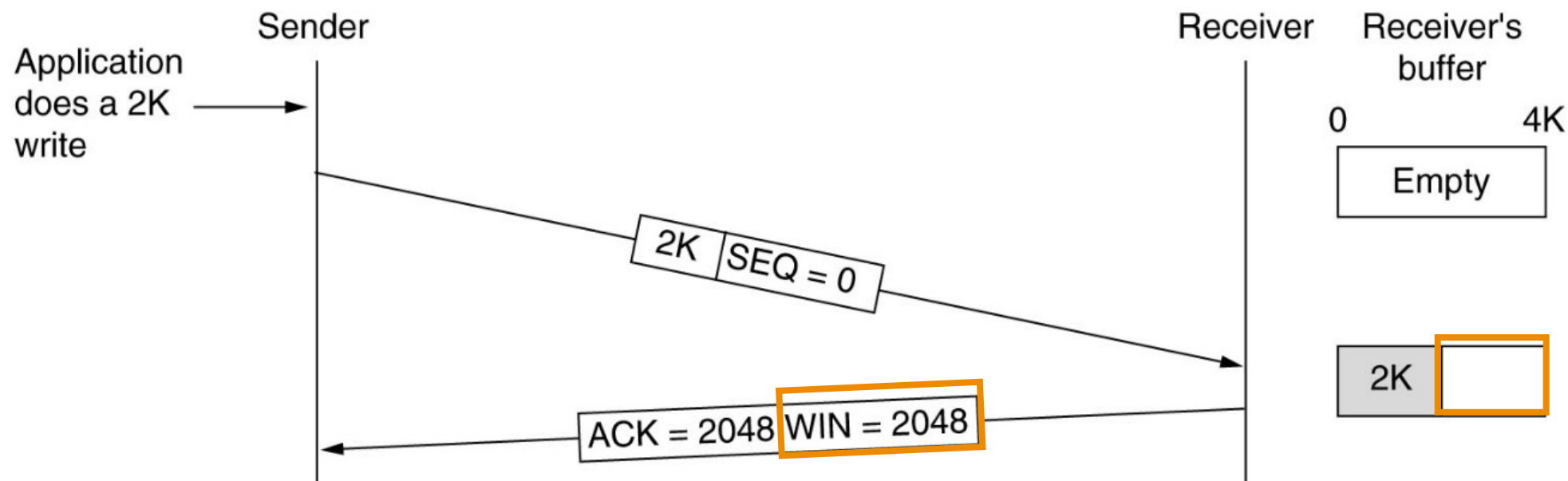
**DYNAMIC  
BUFFER  
MANAGEMENT**

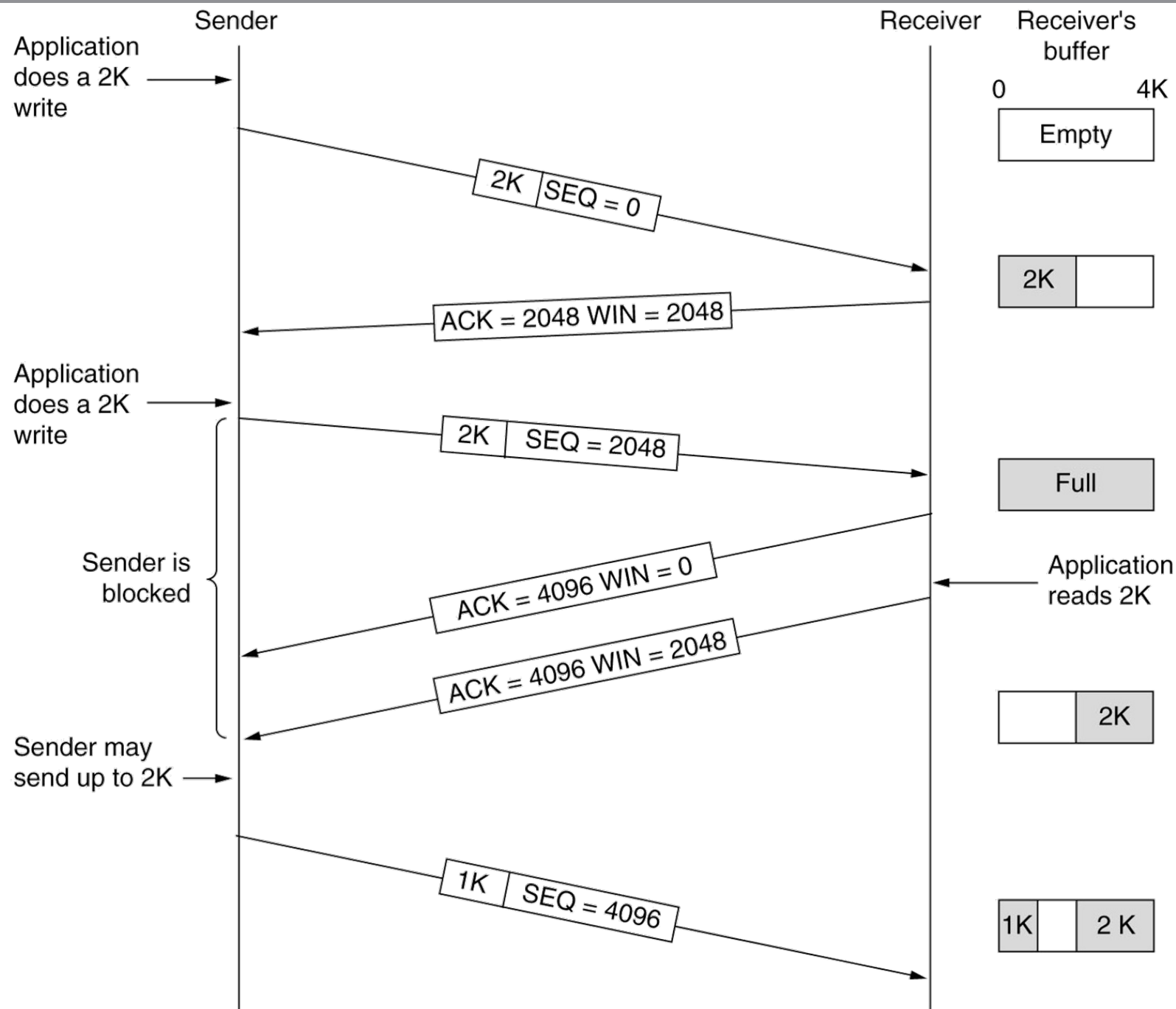
# TCP: Flow Control

Each entity implement flow control using a **credit mechanism**, also called a **window advertisement**.

A **credit specifies** the maximum number of bytes the entity sending this segment can receive and buffer from the other entity

**DYNAMIC  
BUFFER  
MANAGEMENT**





## DYNAMIC BUFFER MANAGEMENT

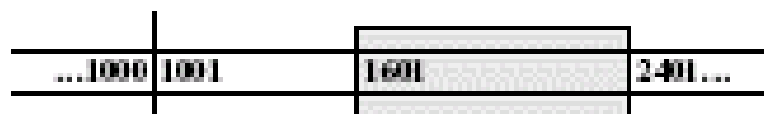


TCP flow control the **sequence number** refers to **byte sequence** instead of packet (or segment) sequences.

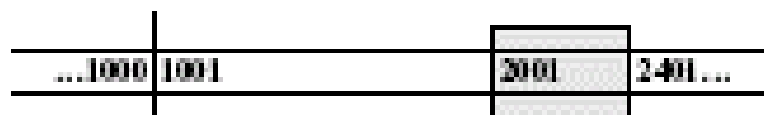
## Transport Entity A



A may send 1400 octets



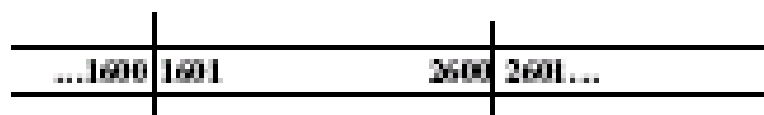
A shrinks its transmit window with each transmission



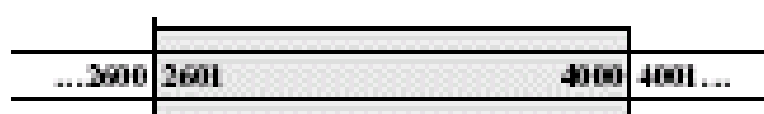
A adjusts its window with each credit



A exhausts its credit



A receives new credit



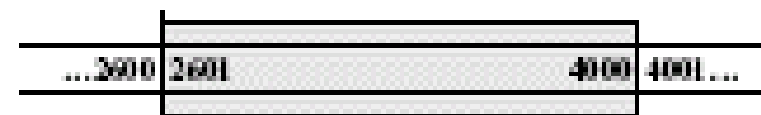
## Transport Entity B



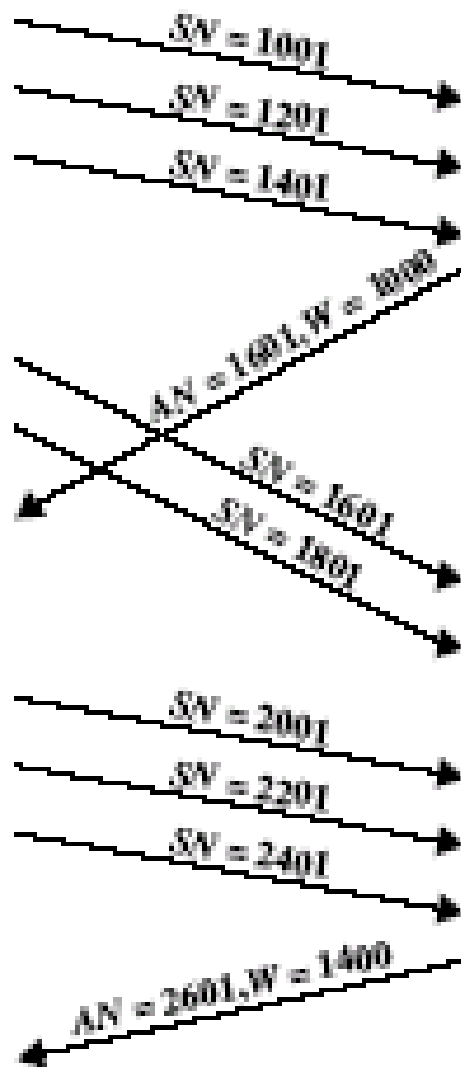
B is prepared to receive 1400 octets, beginning with 1001



B acknowledges 3 segments (600 octets), but is only prepared to receive 200 additional octets beyond the original budget (i.e., B will accept octets 1601 through 2600)



B acknowledges 5 segments (1000 octets) and restores the original amount of credit



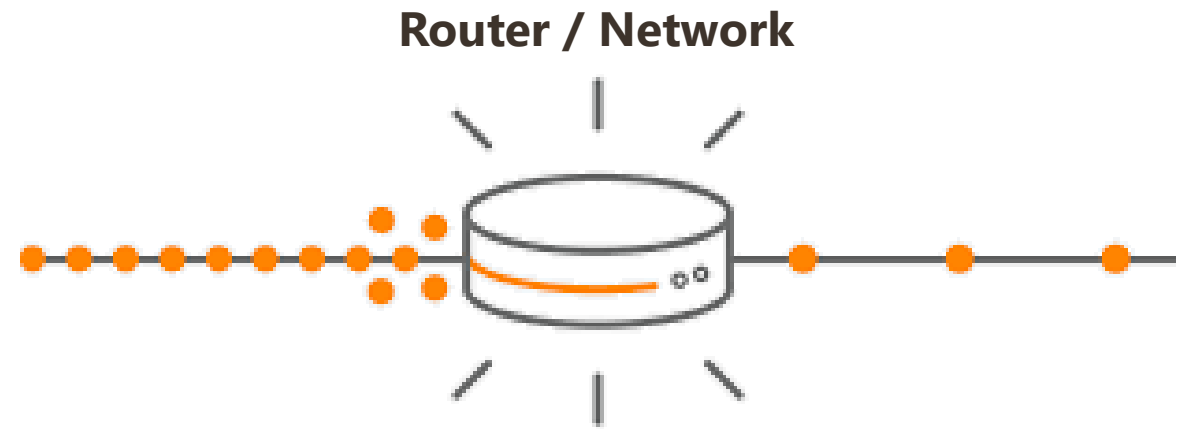
# TCP Congestion Control

- **Network layer** *also tries* to manage congestion, **but difficult**

- TCP does the **heavy lifting**



**“The law of conservation of packets:  
cannot inject new packets into the  
network until the old one leaves”**







# TCP Congestion Control

- **TCP dynamically manipulates the window size**

- ✓ Detecting congestion
- ✓ Prevent congestion (try)
- ✓ React to congestion

- **Detecting Congestion, Difficult?**

- Old days : **transmission error** or **packet discard** 
- Nowadays : **most transmission timeouts** on the Internet are due to congestion 

# Prevent Congestion

“ultimately, congestion can only be controlled by limiting the total amount of data entering the internet to the amount that the internet can carry”

- **In Data Link Control Protocol**

- ✓ Sliding window mechanism helps **to pace the sender**

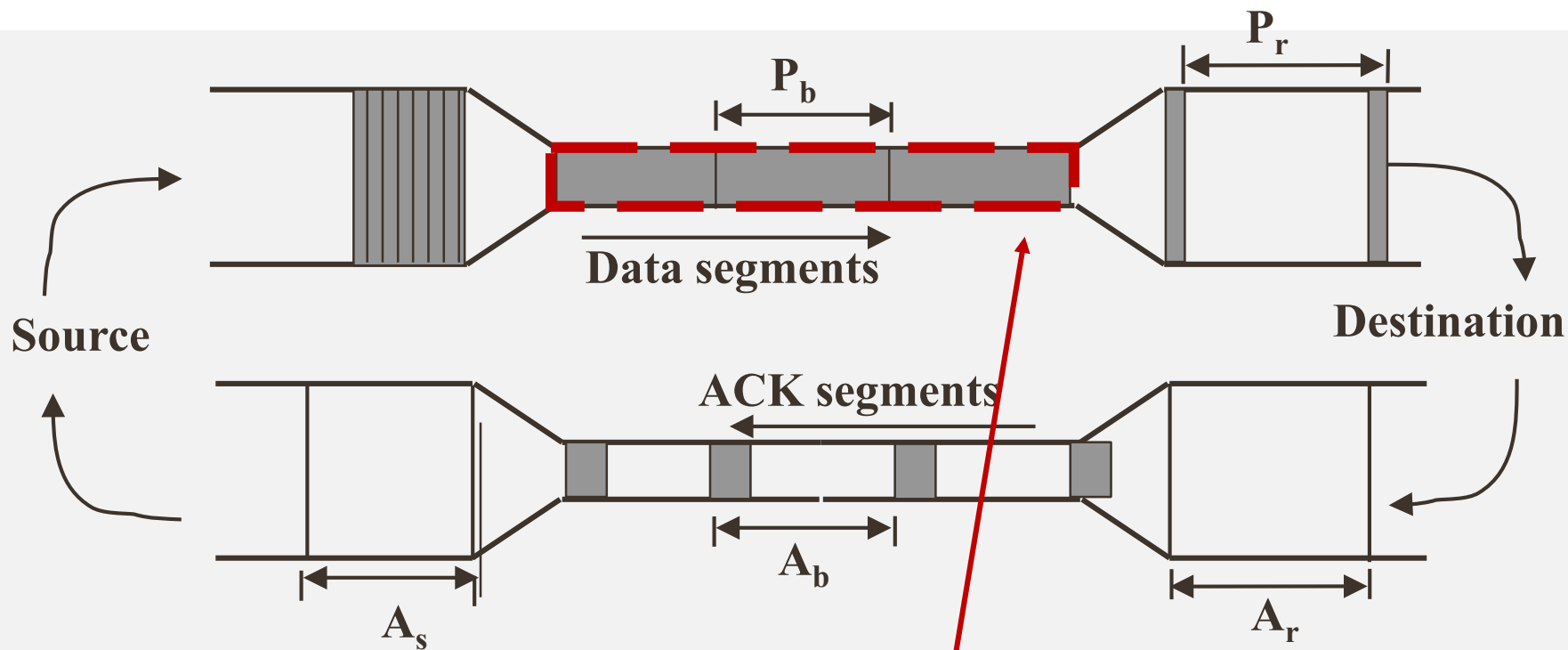
- **In TCP**

- same **pacing effect** as data link protocol
- **Data Sending Rate = Ack Arrival Rate**, once any initial credit is used up



# Flow Determined by **Network**

- $P_b$  = time of minimum segment spacing on the slowest link



Congestion  
CONTROL

Rate of ACK arrival is dependent on round trip time

1. bottlenecks in the **network**
2. bottlenecks at the **destination**



# TCP Segment Pacing

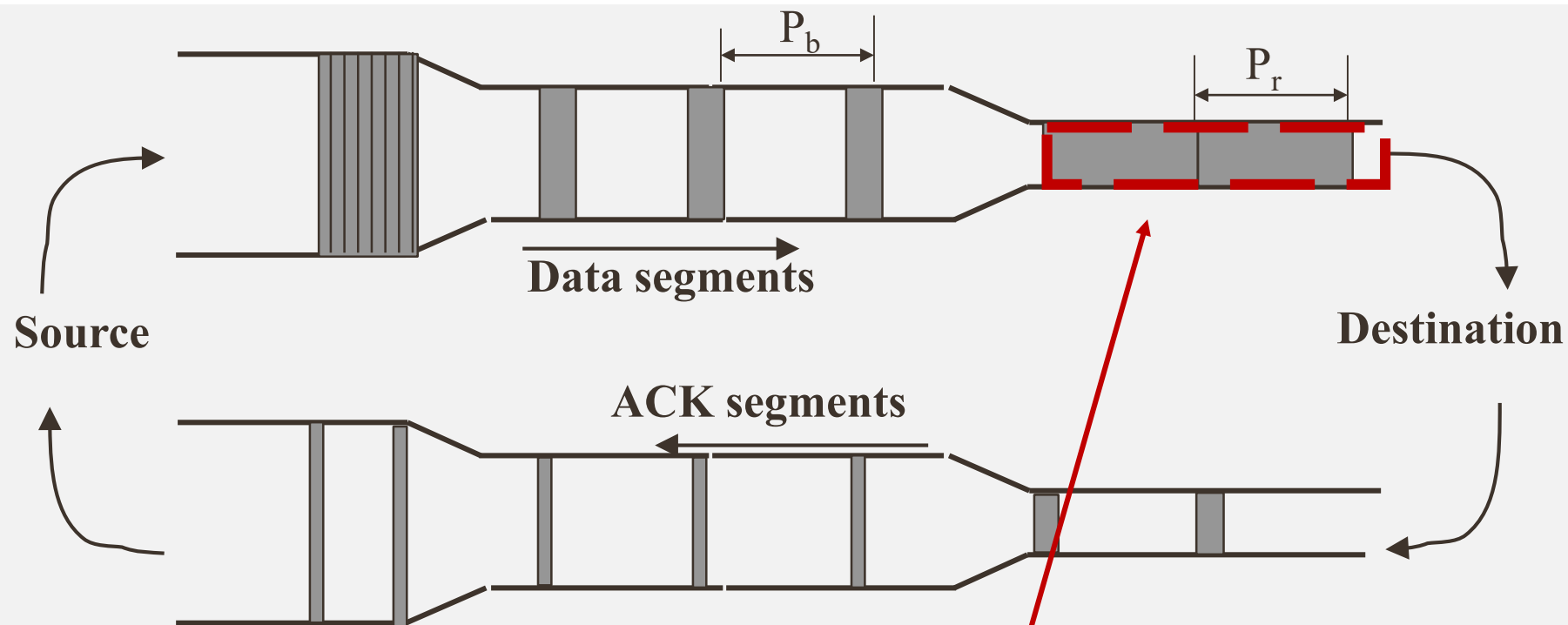
- The **thickness** of the pipe is **proportional to the data rate**
- **Source** and **destination** are on **high capacity networks**
- **Lower speed links create bottlenecks**
- Segment represented by rectangle areas - spreads out if data rate is low
- The wider spacing is preserved at the destination even though the data rate on the final link is high, therefore the segment spacing at the receiver is:  **$P_r = P_b$**

# TCP Segment Pacing – cont.

- In **steady state**,
  - ✓ sender's segment rate = arrival rate of ACKs
  - ✓ sender's rate is determined by the slowest link on the path
- Sending TCP entity automatically **senses and regulates flow**
  - ✓ self-clocking
- Self-clocking works equally well with the bottleneck at receiver
  - ✓ ACKs are sent out at a rate equal to the absorption capacity of the destination

# Flow Determined by **Destination system**

- $P_b$  = time of minimum segment spacing on the slowest link



Rate of ACK arrival is dependent on round trip time

1. bottlenecks in the **network**
2. bottlenecks at the **destination**

# TCP Congestion - Solutions

---

1. **Segment Pacing Effect** – *discussed before*
2. **Slow Start**
3. **Dynamic Window Sizing on Congestion**

## 2. Slow Start

$$\text{awnd} = \min\{\text{credit}, \text{cwnd}\}$$

- **Each sender maintains two windows**
  - **awnd** - Window the receiver has granted
  - **cwnd** - Congestion window  
each window reflects the number of bytes the sender may transmit
- **credit** - amount of unused credit granted in the most recent ACK in segments

## 2. Slow Start – cont.

- **After connection is established**

- sender initializes **cwnd** to the size of the maximum segment (normally **cwnd=1**) and sends one maximum segment
- if segment is ACKed before timeout, **cwnd** is increase to two maximum segment size and two segments are sent

- **As each segments is ACKed,**

- the **cwnd** is increased by one maximum segments size

- **The idea:**

- if bursts of size **n** bytes work fine but a burst of **n+n** bytes gives timeout, the **cwnd** is set to n bytes to avoid congestion

**cwnd** grows exponentially until either a receiver's window is reached, or timeout occur

### 3. Dynamic Window Sizing on Congestion

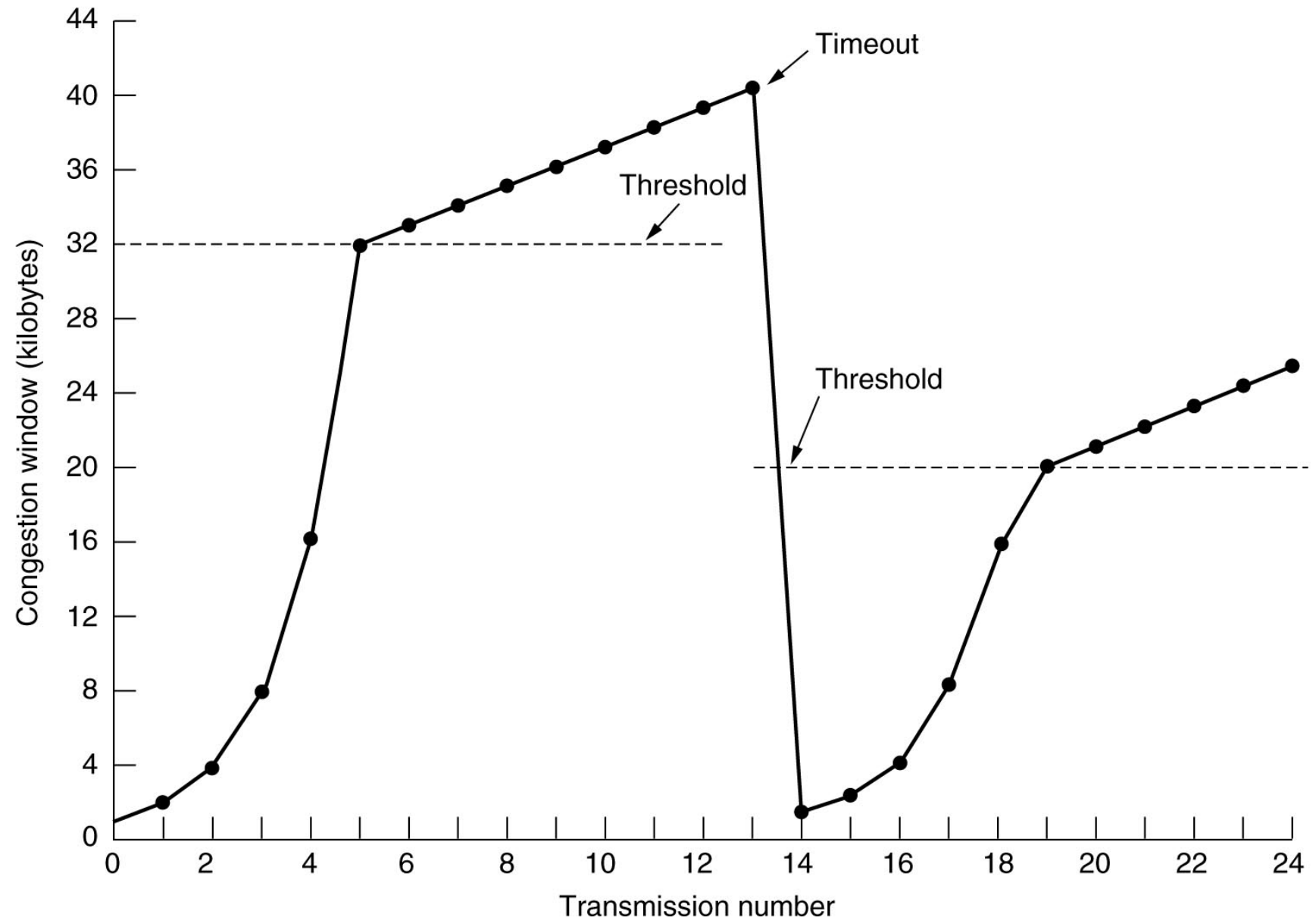
- An internet congestion control algorithm
- if **congestion causes** a timeout, **cut back flow, ramp up slowly**

- **when timeout occurs** 

- **ssthresh** = slow start threshold
- set **ssthresh** = **cwnd**/2
- set **cwnd**=1 and **performs slow-start until cwnd = ssthresh**
- from then on **linear grow**, increase **cwnd by 1** (segment) **after every ACK** until the receiver's window is reached or timeout occur



# Dynamic Window Sizing



# TCP Timer Management

---

- TCP uses multiple timer (at least conceptually) to do its work.
- Retransmission timer
  - How long should the timeout interval be?
- Determine the Round-Trip Time (***RTT***), time between sending a segment and receiving its ACK, is tricky.
- Even if ***RTT*** is known, deciding the timeout interval is also difficult
  - if timeout is set too short (say T1) - unnecessary retransmission
  - if timeout is set too long (say T2) – performance suffer due to long retransmission delay

# TCP Timer Management – cont.

- Highly dynamic algorithm that constantly adjust the timeout interval, based on continuous measurements of the network performance.
- Jacobson's algorithm (1988) – **RTT** variance estimation:
  - ✓ generally **used by TCP**
  - ✓ for each connection, determine best retransmission timer (**RTT**) by using an estimate of **RTT** which includes an estimate of the variance
  - ✓ use the *mean deviation* (not standard deviation) and *exponential averaging*

# TCP Timer Management – cont.

- for each connection, TCP entity maintains a variable  $SRTT$ , that is the best current estimate of  $RTT$  to the destination.
- if an ACK gets back before timeout, TCP measures  $RTT$  for the ACK, and updates  $SRTT$ :

$$SRTT(k+1) = (1-\alpha) \times SRTT(k) + \alpha \times RTT(k+1)$$

where  $\alpha$  is a smoothing factor that determines how much weight is given to the old value. Typically  $\alpha = 1/8$

- even with good value of  $SRTT$ , selecting the retransmission timeout is still difficult

- Normally, TCP uses  $\beta \times SRTT$ 
  - Trick is selecting  $\beta$
  - it was initially constant value  $\beta = 2$  (inflexible)
- Jacobson proposed making  $\beta$  roughly proportional to the *mean deviation* (not standard deviation) of the acknowledgment arrival time probability density function

- keeping track of another smooth variable,  $D$

$$D = \alpha \times D + (1-\alpha) |SRTT - RTT|$$

where  $\alpha$  may or may not be the same value used to smooth  $SRTT$

- most TCP implementations now uses Jacobson's algorithm to set the timeout interval.

$$RTO = SRTT + 4 \times D$$

# Karn's Algorithm

What happen during timeout retransmission

Which RTT samples should be used as input to Jacobson's algorithm



- **Karn's proposed a simple fix to Jacobson algorithm**

- Do not update **SRTT** on any segments that have been retransmitted
- Timeout is doubled on each failure until the segments get through the first time



# Exponential Backoff

- What retransmission timer (**RTO**) value should be used on a retransmitted segment?

TCP source increases its **RTO** value each time the same segment is retransmitted - backoff process

- After the first retransmission of a segment
  - wait for a longer time before performing a second retransmission

$$\text{RTO} = q \times \text{RTO}$$

- **RTO** grows exponentially after each backoff (common value of  $q$  is 2)



# TCP

## Implementation Policies

- Send Policy
  - Silly Window Syndrome
- Deliver Policy
- Accept Policy
- Retransmit Policy
- Acknowledge Policy



# TCP: Implementation Policies

---

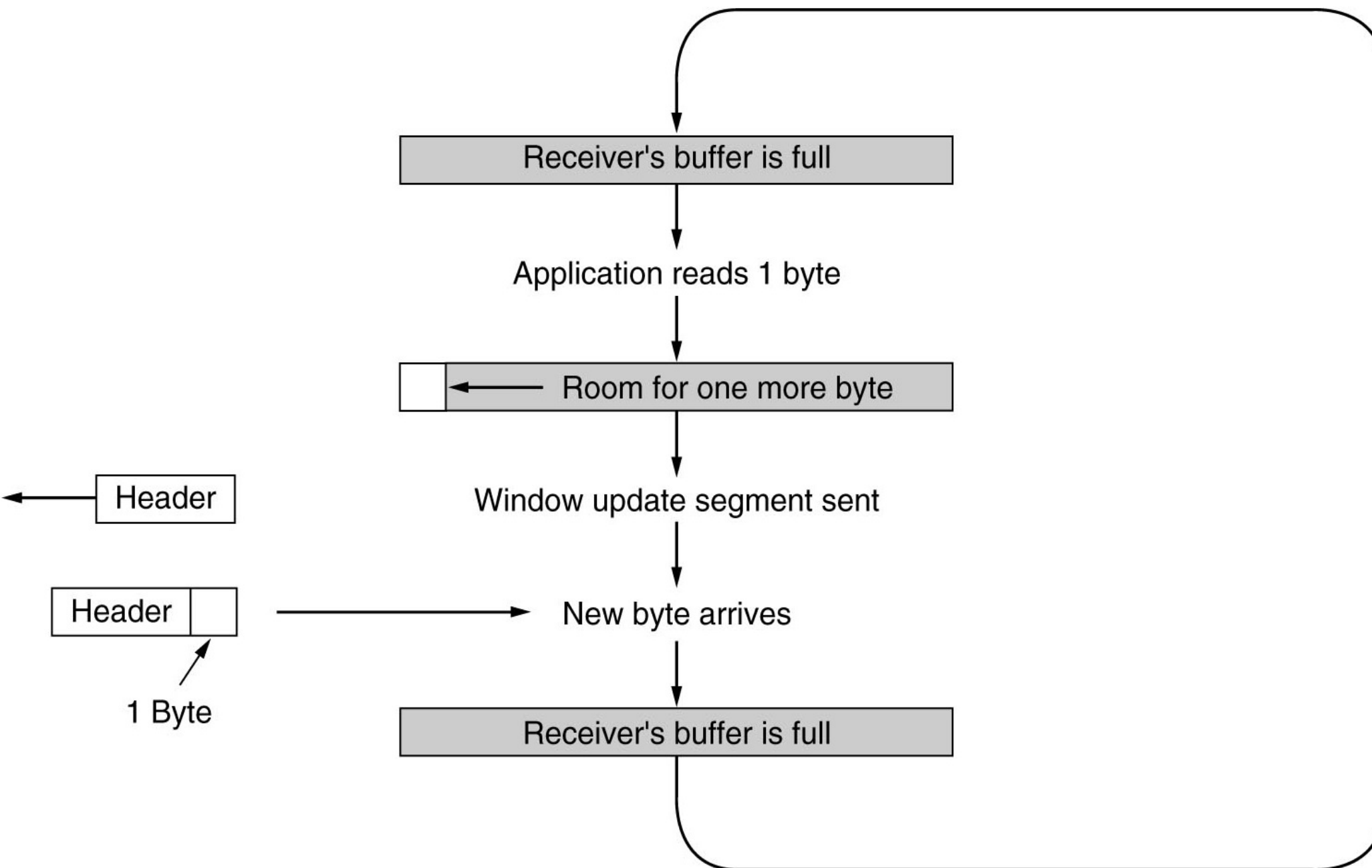
- The standard provide a precise specification of the protocol to be used between TCP entities
  
- Policies defined for the protocol
  1. **Send policy**
  2. **Deliver policy**
  3. **Accept policy**
  4. **Retransmit policy**
  5. **Acknowledge policy**

# 1. Send Policy

In the **absence of PUSH** data and a closed transmission window, a sending TCP entity is free to **transmit data at its own convenience**

- depend on performance considerations
  - ✓ if **infrequent and large** (*buffer @ sender*)
    - low overheads
    - slow response
  - ✓ if **frequent and small** (*buffer @ receiver*)
    - quick response
    - high overheads
    - silly window syndrome

# Silly Window Syndrome



## 2. Deliver Policy

**Too frequent delivery means too many OS interrupts**

- Arriving data are stored in deliver buffer

- ✓ **if PUSH flag is set**

Data along with any other data in the deliver buffer are submitted to the destination application in a *RECEIVE* command.

- ✓ **if PUSH flag is not set**

TCP may wait, *e.g. to avoid excess interrupts*

- ✓ **if URG flag is set**

The receiving application is signaled that urgent data is present

# 3. Accept Policy


---

## ▪ In Order

- ✓ **Discard** out of order segments

## ▪ In Window

- ✓ **Accept** all segments within the receive window
- ✓ Complex acceptance test and sophisticated storage



out of order  
segments are  
accepted

# 4. Retransmit Policy

TCP maintains a **queue of segments** that have been sent but **not ACKed**

- **First Only** suited for **in-window accept policy**
  - ✓ one retransmission timer for the entire queue
  - ✓ if an ACK is received, the segment/s removed and timer reset
  - ✓ if timer expires, first segment in the queue is retransmitted
- **Batch** suited for **in-order accept policy**
  - ✓ same as above, except when timer expires, retransmit all segments in the queue
- **Individual** suited for **in-window accept policy**
  - ✓ one timer for each segment

Selective-  
Repeat ARQ

Go-Back-N  
ARQ

Selective-  
Repeat ARQ

# 5. Acknowledge Policy

---

- **Immediate**

- **Cumulative**

- ✓ wait for an outbound segment, piggyback the ACK
- ✓ Timer to avoid long delay

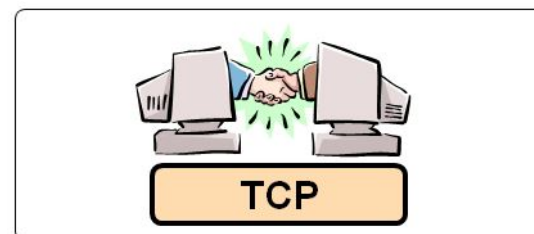
# User Datagram Protocol (UDP)

- Fundamentals
- Applications
  - DNS
  - RPC

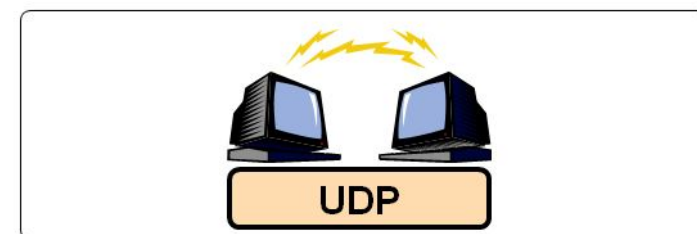


# Protocols: UDP [RFC 768]

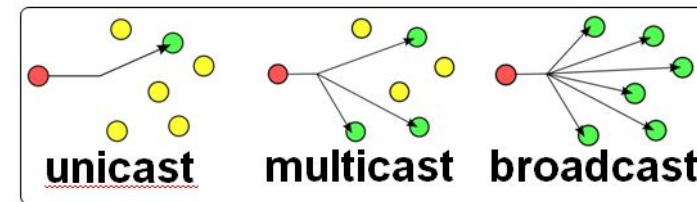
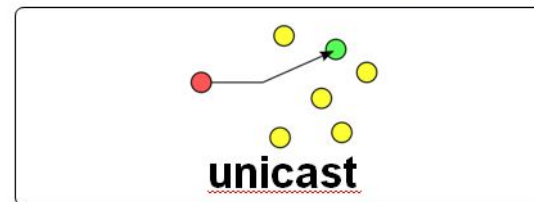
- **TCP: Reliable ordered delivery** of packets
  - Error detection, retransmissions and acknowledgements.
  - TCP is strictly used for **point to point**
  - TCP segments the data before sending to Network layer – Stream oriented
- **UDP: Unreliable delivery** of packets
  - Connectionless protocol
  - No retransmissions, acknowledgements
  - UDP does not segment the data – message oriented



- **Slower but reliable transfers**
- **Typical applications:**
  - Email
  - Web browsing



- **Fast but non-guaranteed transfers ("best effort")**
- **Typical applications:**
  - VoIP
  - Music streaming

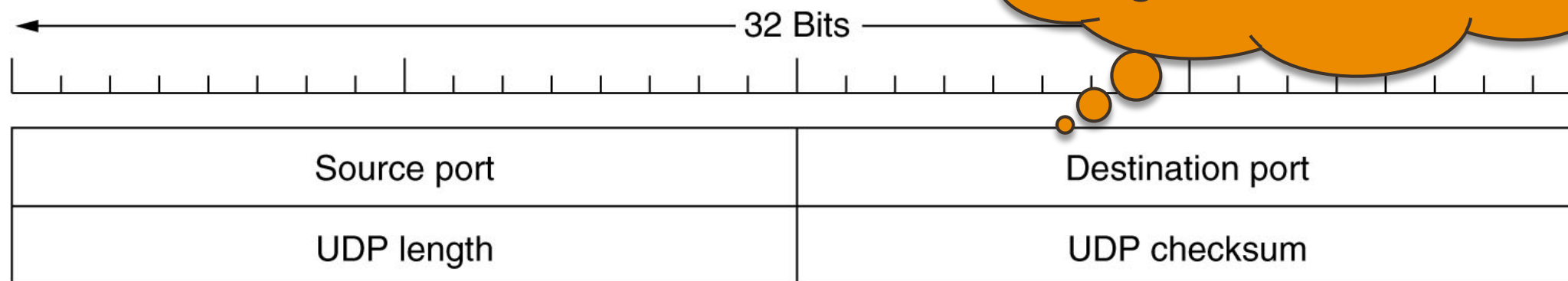


No flow control, error control, or retransmission upon receipt of bad segments

# Why UDP?

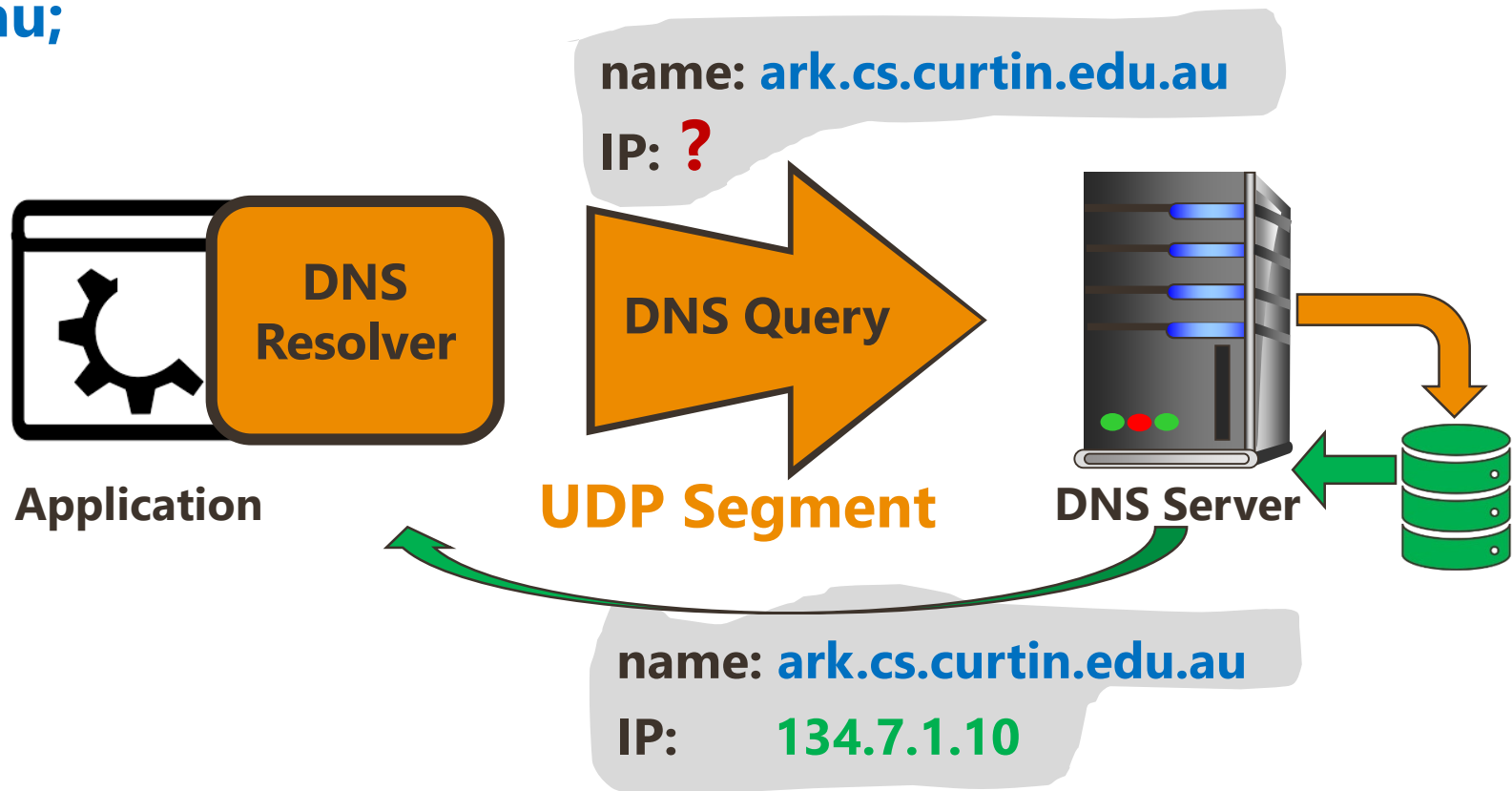
- Lower overhead enables faster transmissions
- **Unicast, Multicast, Broadcast**
- UDP is especially useful in client-server situation
  - ✓ DNS (Domain Name System) uses UDP

TCP protocol requires more information and overhead to guarantee data delivery



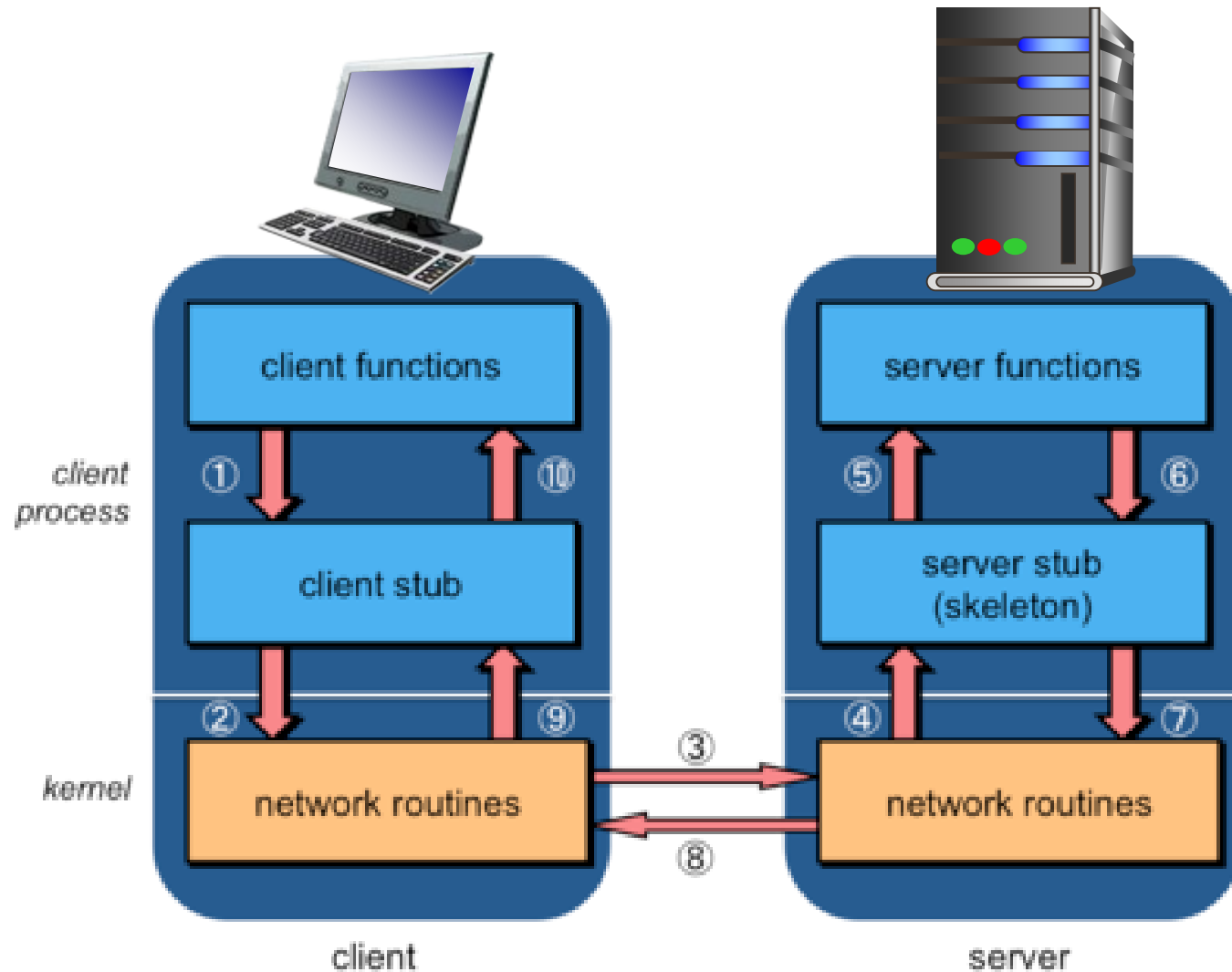
# UDP: DNS (Domain Name System)

- How can the **IP address of the remote host** be found?
- Hierarchical, symbolic addresses are used
  - e.g. **ark.cs.curtin.edu.au**;
  - **curtin.edu.au**



# Remote Procedure Call (RPC) - 1984

- **Allows program to call procedures located on remote hosts**
- Information is transported from the caller to the callee (remote host) in the parameters and can come back in the procedure result
- **Message passing is invisible** to the programmer



# RPC – cont.

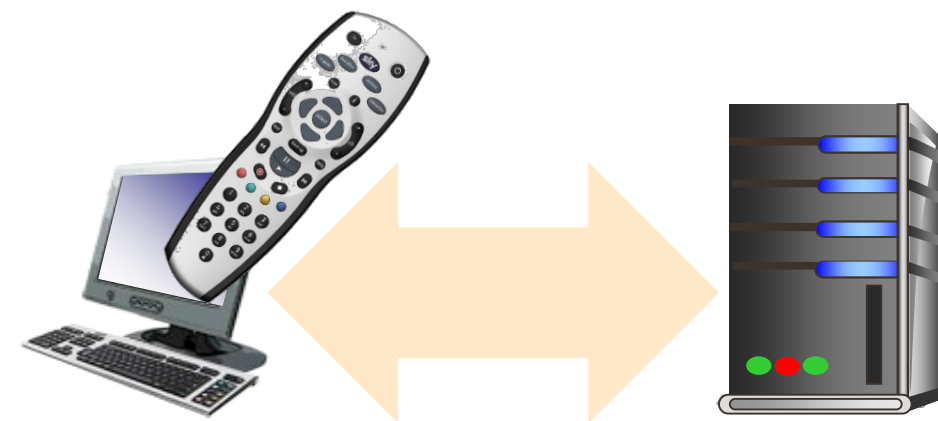
## ▪ Client Stub

- **a small library procedure** that represents the server procedure in the client's address space that the Client program must be bound with.

## ▪ Server Stub

- **a procedure call that represents the client procedure call** in the server's address space that the Server program must be bound with.

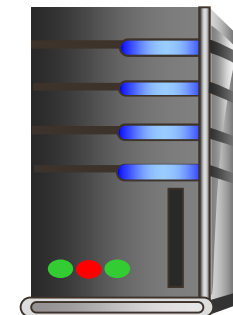
## ▪ UDP is commonly used for RPC



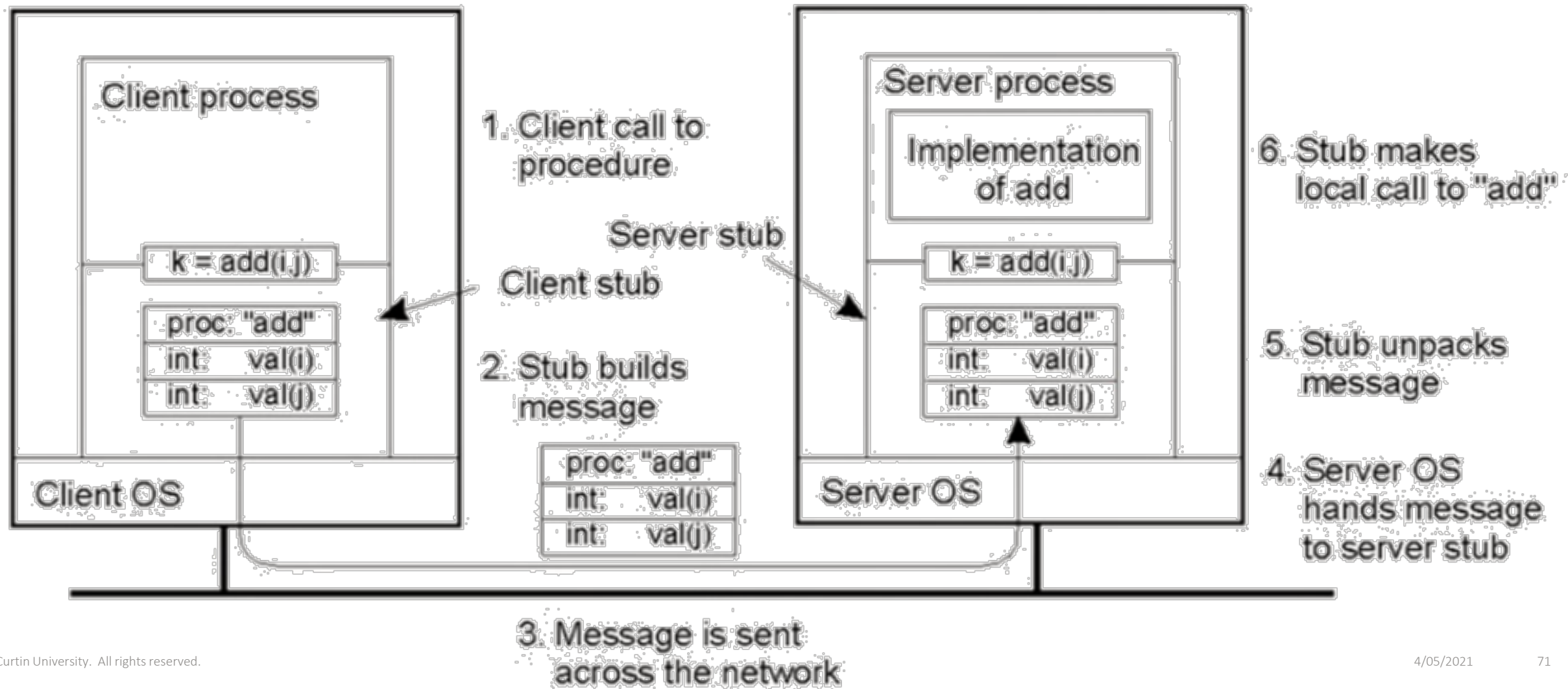
# RPC



Client machine



Server machine



# SUMMARY

## ■ Flow Control and Congestion Control

- Flow Control
  - Transport Layer
  - Data Link Layer
- Congestion Control
  - Warning-bit
  - Choke-packets
  - Load-shedding
  - RED

## ■ TCP

- TCP Header
  - Flags (SYN, FIN, ACK, RST)
  - Flag (URG, PSH) – *in depth*
  - TCP Options
  - Window Size (Dynamic Buffer Management)
- TCP Flow Control
  - Dynamic Buffer Management
- TCP Congestion Control
  - TCP Segment Pacing Effect
  - Slow Start
  - Dynamic Window Sizing
- TCP Timer Management

## • TCP - Implementation Policies

- Send Policy
  - Silly Window Syndrome
- Deliver Policy
- Accept Policy
- Retransmit Policy
- Acknowledge Policy

## • UDP

- Fundamentals
- Applications
  - DNS
  - RPC





Curtin University

# THANK YOU

Make tomorrow better.

314