



# Computer Networks (WS23/24)

## L9: The Transport Layer - Part 3

**Prof. Dr. Lin Wang**

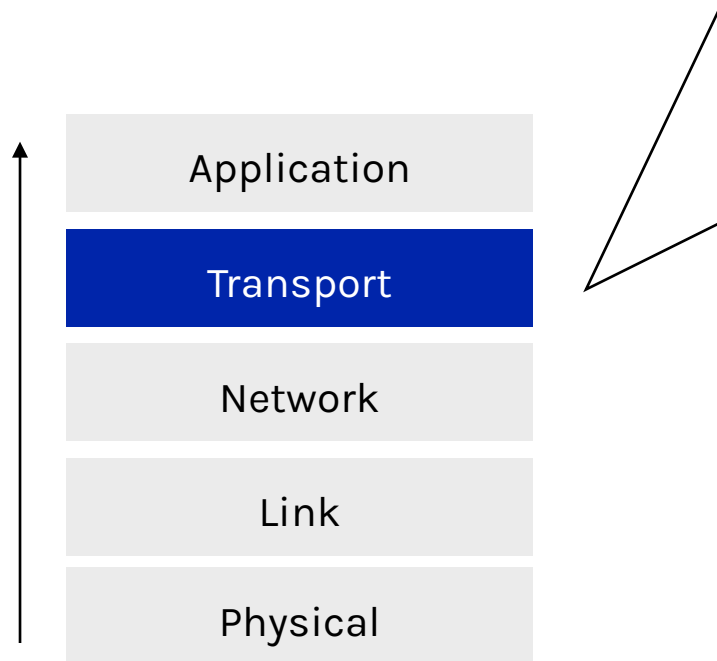
Computer Networks Group (PBNNet)

Department of Computer Science

Paderborn University

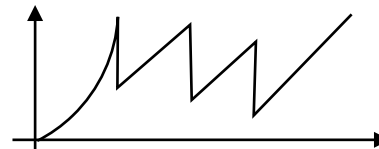


# Learning objectives



## Part 3

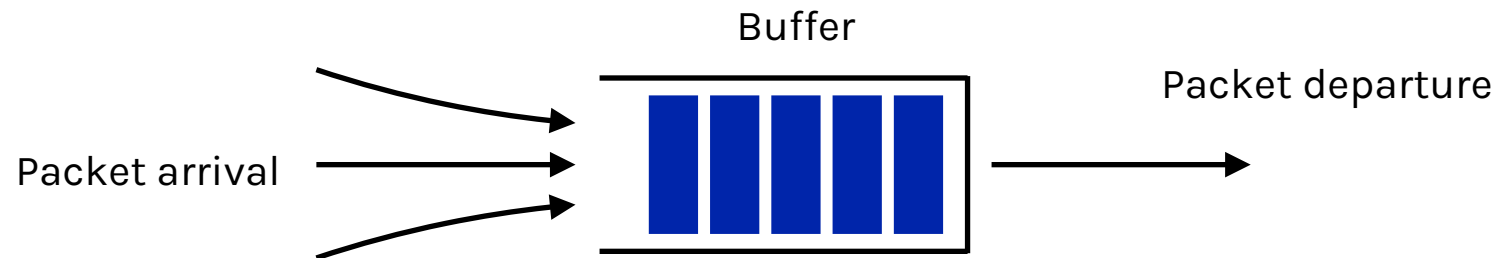
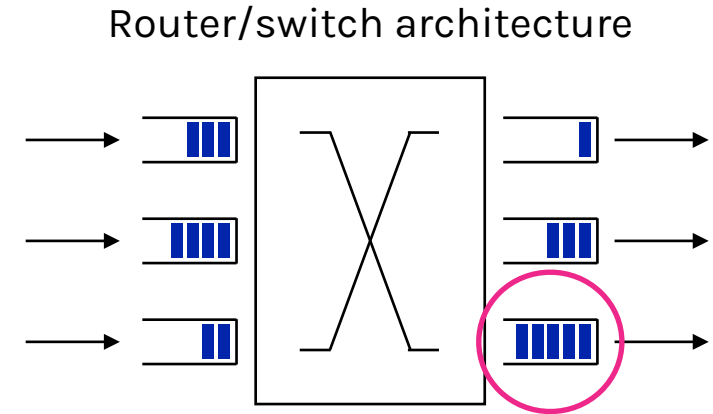
- Network congestion
- Congestion detection and control
- Explicit Congestion Notification (ECN)



# Network Congestion



# Why congestion happens?



Buffer has a **limited** size; buffer **overflow** if too many packets arrive within a short period of time

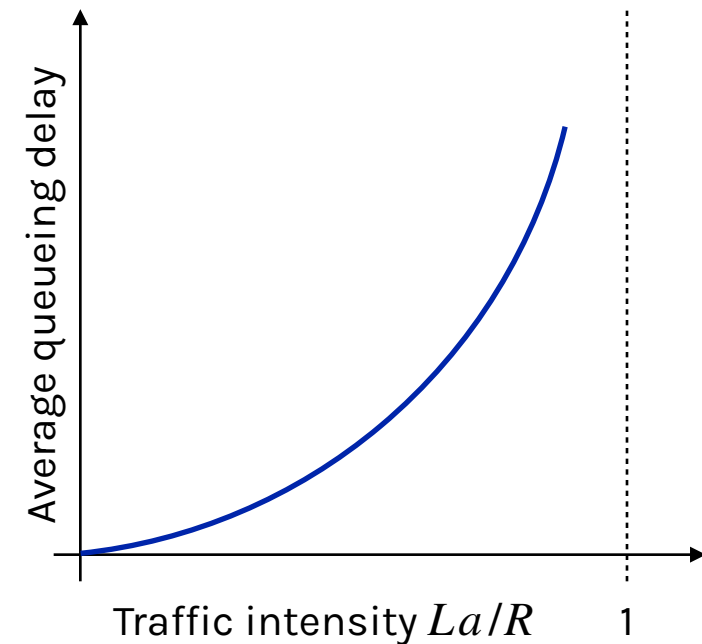
# Quantifying congestion

## Parameters

- Average packet arrival rate  $a$ , transmission rate  $R$ , fixed packet length  $L$
- **Traffic intensity** is given by  $La/R$

When the traffic intensity is  $> 1$ , the queue will increase without bound, so does the queueing delay

When the traffic intensity is  $\leq 1$ , queueing delay depends on the burst size



# Network congestion is not a new problem

## The Internet almost died of congestion in 1986

- Throughput collapsed from 32 Kbps to 40 bps

## Van Jacobson saved us with Congestion Control

- His solution went right into BSD implementation

## Recent resurgence of research interest after brief lag

- New methods (ML-based), new contexts (data centers, 5G), new requirements



# Original network behavior

## Upon connection establishment

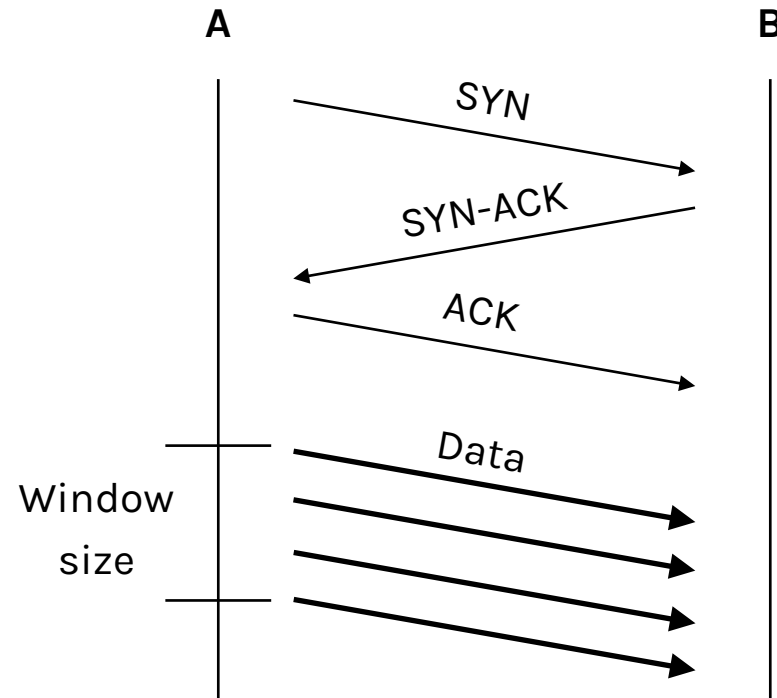
- Send full window of packets

## Upon timer expiration

- Retransmit packet immediately

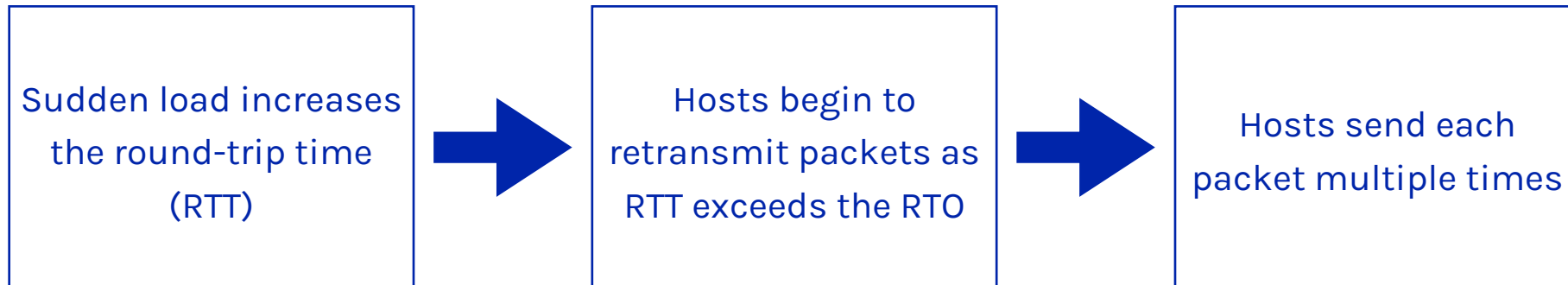
## Outcome

- Sending rate only limited by flow control
- Net effect: window-sized burst of packets



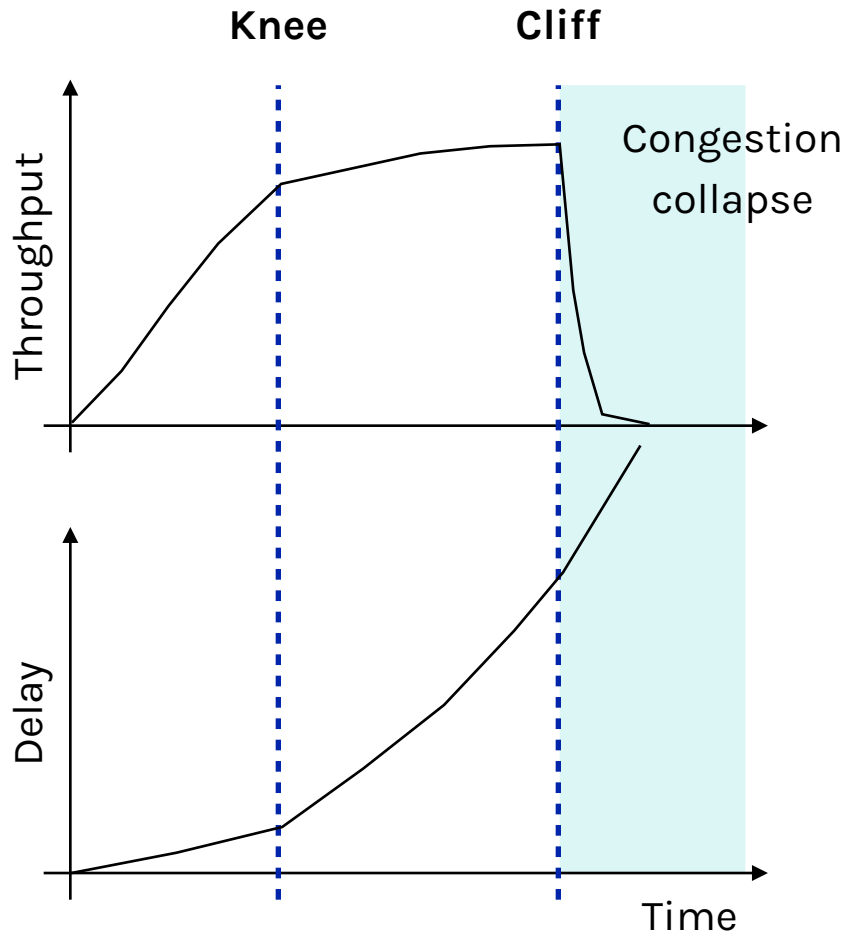
# Congestion collapse

Increase in network load results in a decrease of useful work done





# Congestion collapse analysis



**Knee point:** after which

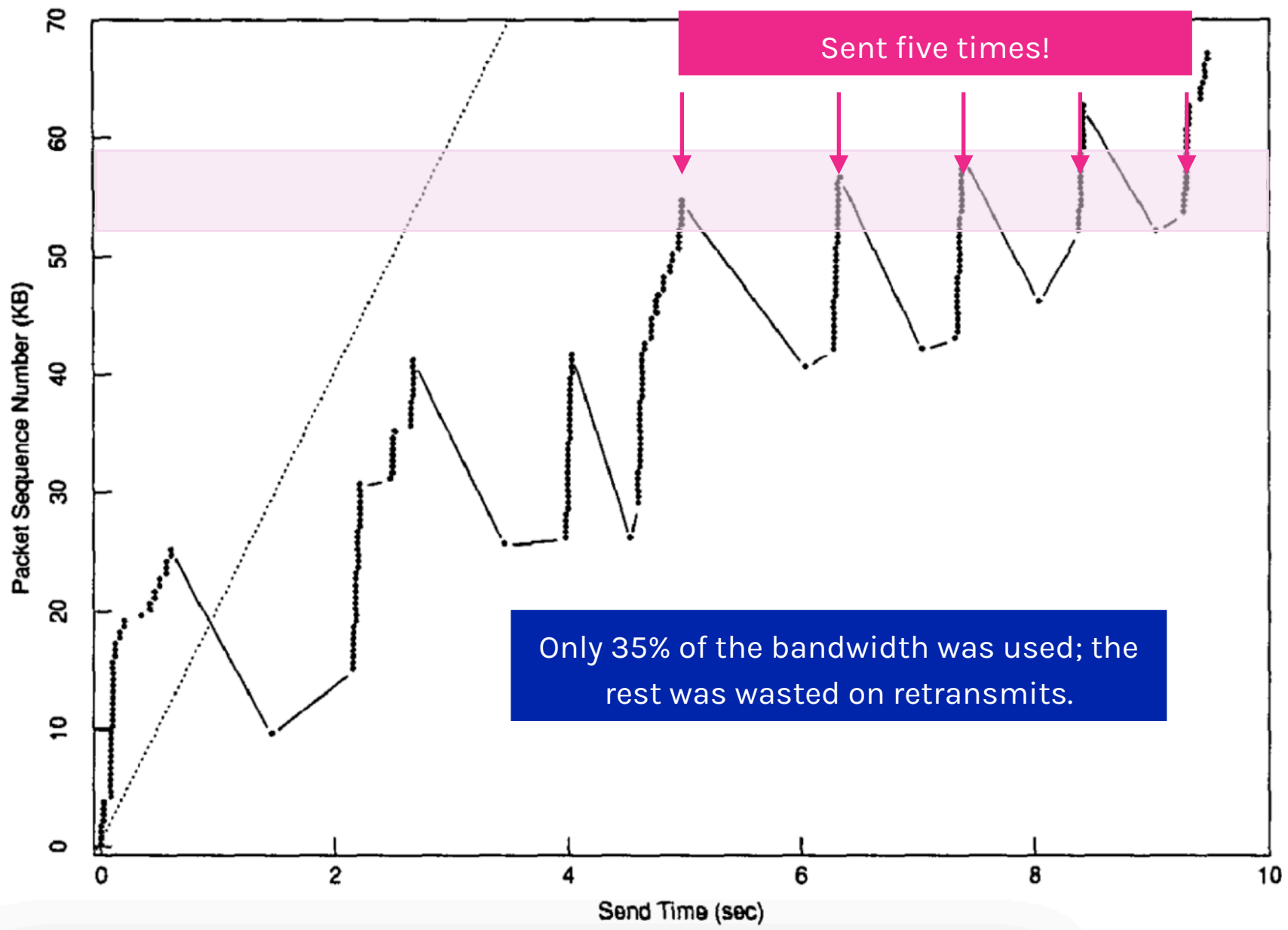
Throughput: increases slowly

Delay: increases quickly

**Cliff point:** after which

Throughput: decreases quickly

Delay: tends to infinity



# Congestion control goals

## Bandwidth estimation

- How to adjust the bandwidth of a single flow to the bottleneck bandwidth?



## Bandwidth adaptation

- How to adjust the bandwidth of a single flow to the variation of the bottleneck bandwidth?



## Fairness

- How to share bandwidth "fairly" among flows, without overloading the network



# Congestion control vs. flow control

## Flow control

Prevents **one fast sender** from overloading a slow **receiver**

(Solved using a **receiving window**)

## Congestion control

Prevents a **set of senders** from overloading the **network**

(Solved using a **congestion window**)

# Windows at the sender

## Receiving window $RWND$

- How many bytes can be sent without overflowing the receiver buffer?
- Based on the receiver input

## Congestion window $CWND$

- How many bytes can be sent without overflowing the routers?
- Based on network conditions

**Sender window:**  $\min\{RWND, CWND\}$

# Congestion Detection and Control



# Congestion detection approaches

## Approach 1

- Ask the network to tell the source
- But signal itself could be lost

## Approach 2

- Measure packet delay
- But signal is noisy (delay often varies considerably)

## Approach 3

- Measure packet loss (fail-safe signal that TCP has to detect already)

# Two signals for detecting packet loss

Loss detection can be done using ACKs or timeouts: differ in the degree of severity

## Duplicate ACKs

- Mild congestion
- Packets are still making it

## Timeout

- Severe congestion
- Multiple consequent losses



# Reacting to congestion

What increase/decrease function should we use to adjust the *CWND*?



**Bandwidth estimation**  
(slow-start)



**Bandwidth adaptation**  
(additive vs.  
multiplicative)



**Fairness**  
(convergence to fair share)

# Bandwidth estimation

Goal: quickly get a first-order estimate of the available bandwidth

## Intuition

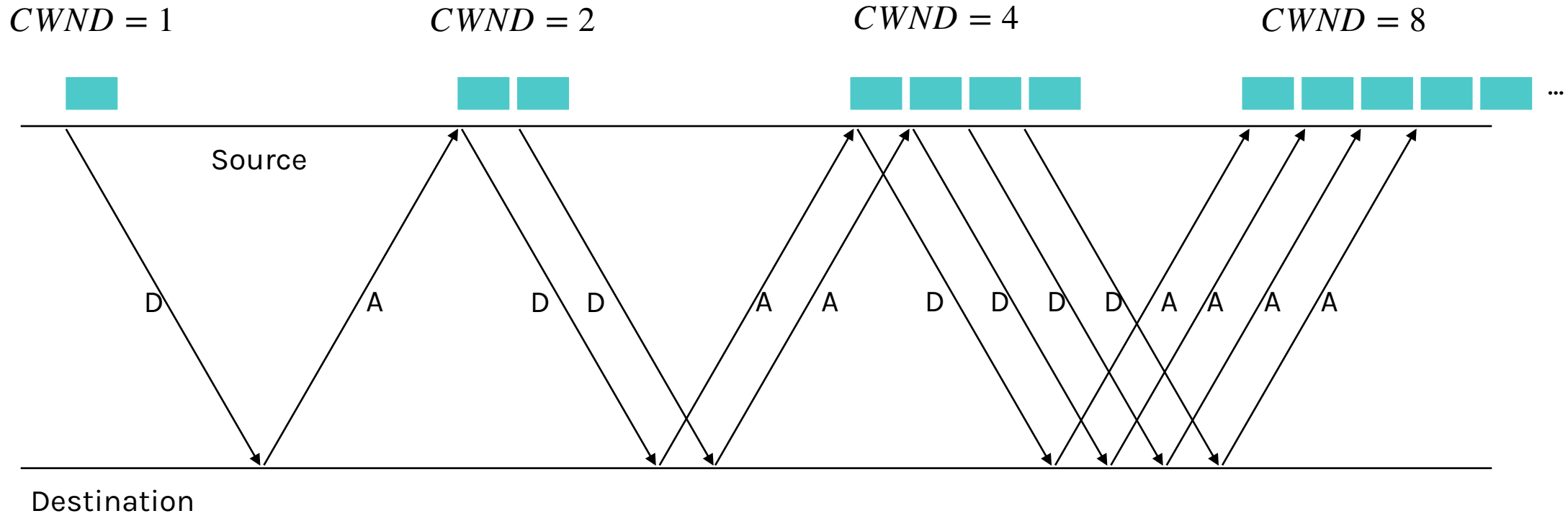
Start slow but rapidly increase until a packet drop occurs

## Increase policy

$CWND = 1$  (initially)

$CWND + = 1$  (upon receipt of an ACK)

# TCP slow-start



*CWND* increases exponentially; slow-start is called like this only because of the small starting point!

# Problem with slow-start

**Problem:** it can result in a full window of packet losses

## **Example**

- Assume that *CWND* is just enough to "fill the pipe"
- After one RTT, *CWND* has doubled; all the excess packets are now dropped

## **Solution**

- We need a more **gentle adjustment algorithm** once we have a rough estimate of the bandwidth

# Bandwidth adaptation

Goal: track the available bandwidth and oscillate round its current value

	Decrease		
Increase		Additive	Multiplicative
Additive		AIAD	AIMD
Multiplicative		MIAD	MIMD

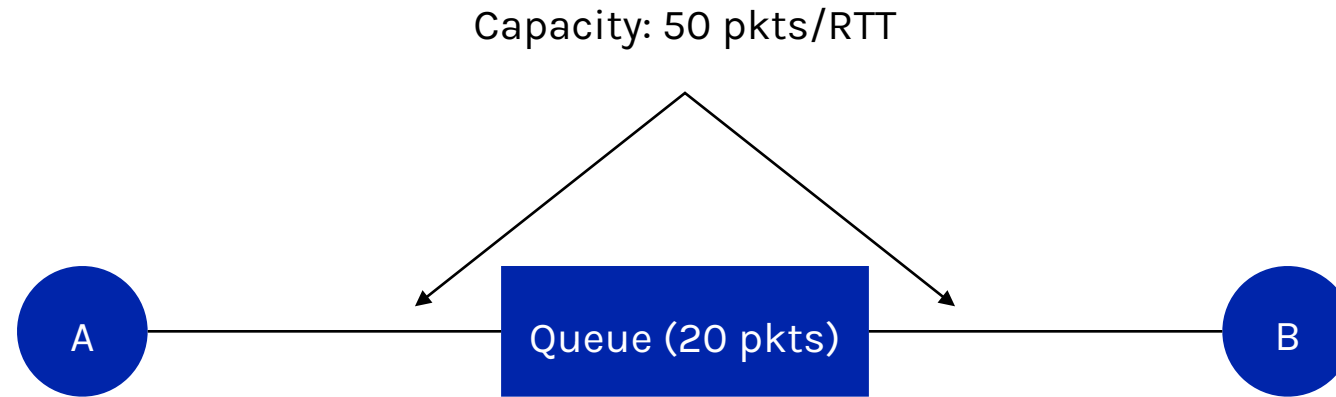
Design space for *CWND* adjustment: 4 alternative designs

# Comparison of design choices

	Increase behavior	Decrease behavior
<b>AIAD</b>	Gentle	Gentle
<b>AIMD</b>	Gentle	Aggressive
<b>MIAD</b>	Aggressive	Gentle
<b>MIMD</b>	Aggressive	Aggressive

How to select among these? Need to weight in another problem: fairness

# Example: one flow with AIMD



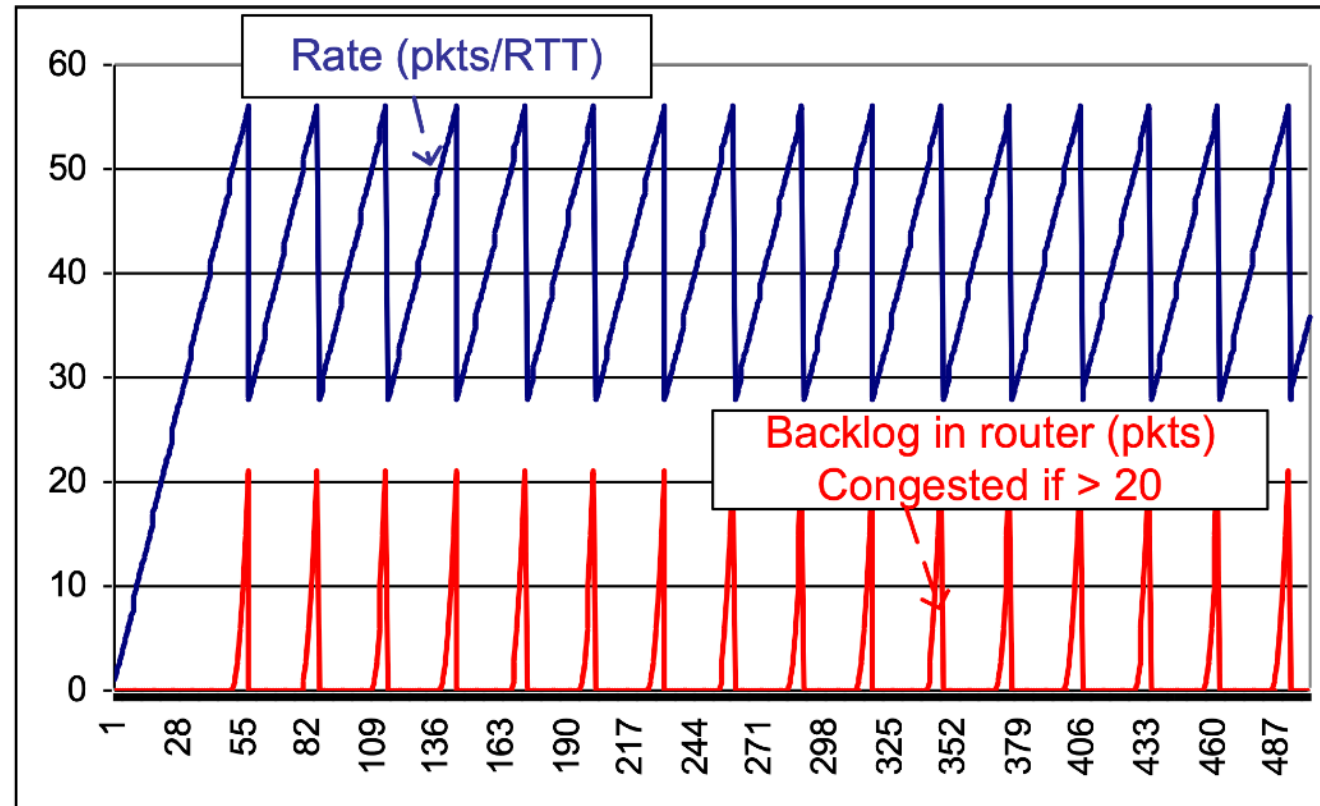
**Without congestion**

*CWND* increases by one packet every RTT

**Upon congestion**

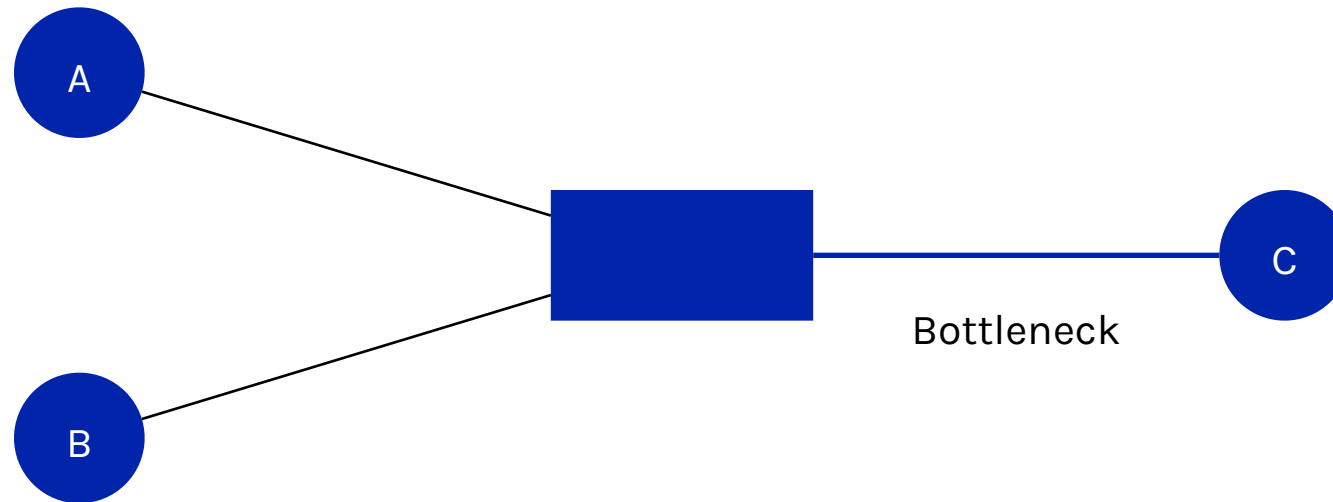
*CWND* decreases by a factor of 2

# Example: one flow with AIMD

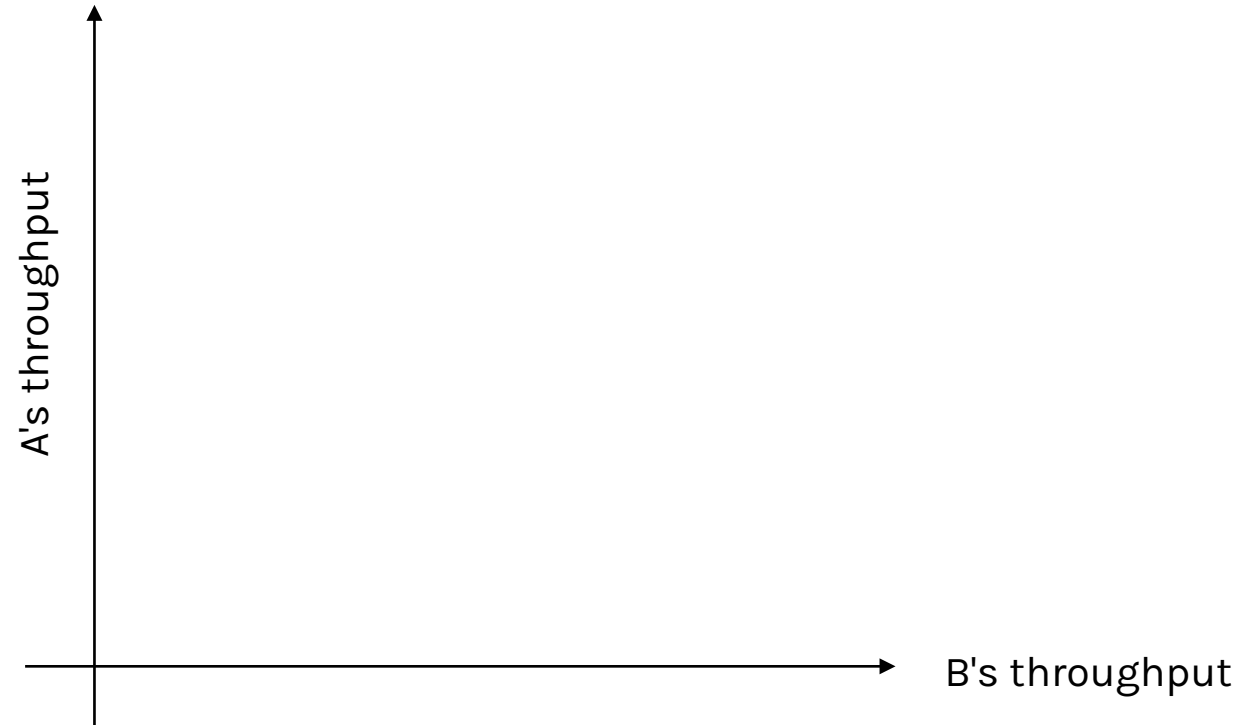




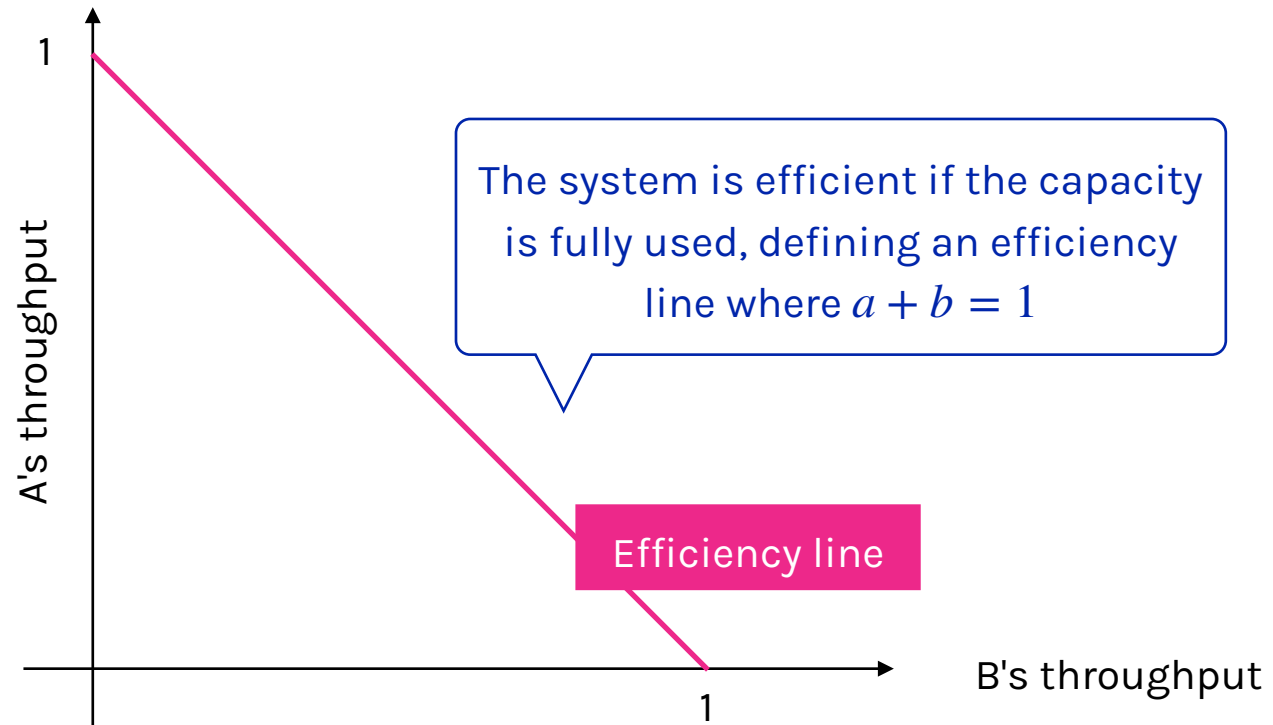
## Example: two flows sharing a bottleneck



# System behavior

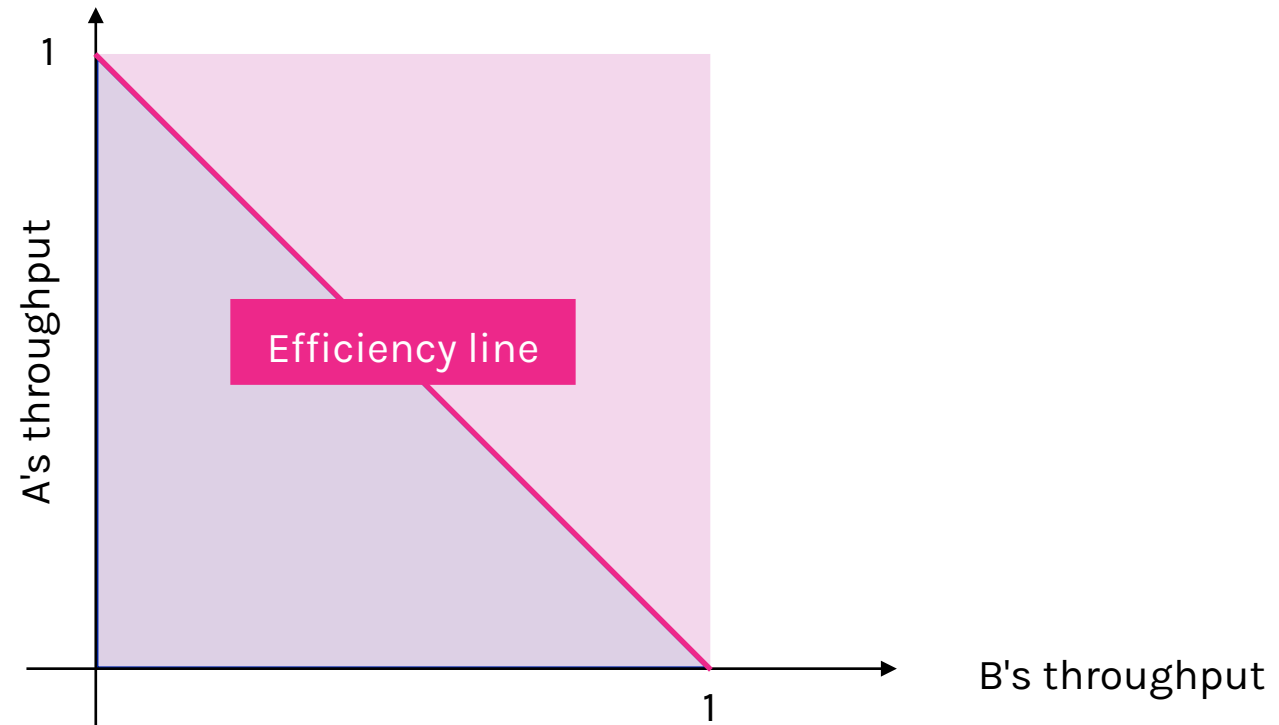


# System behavior: efficiency line

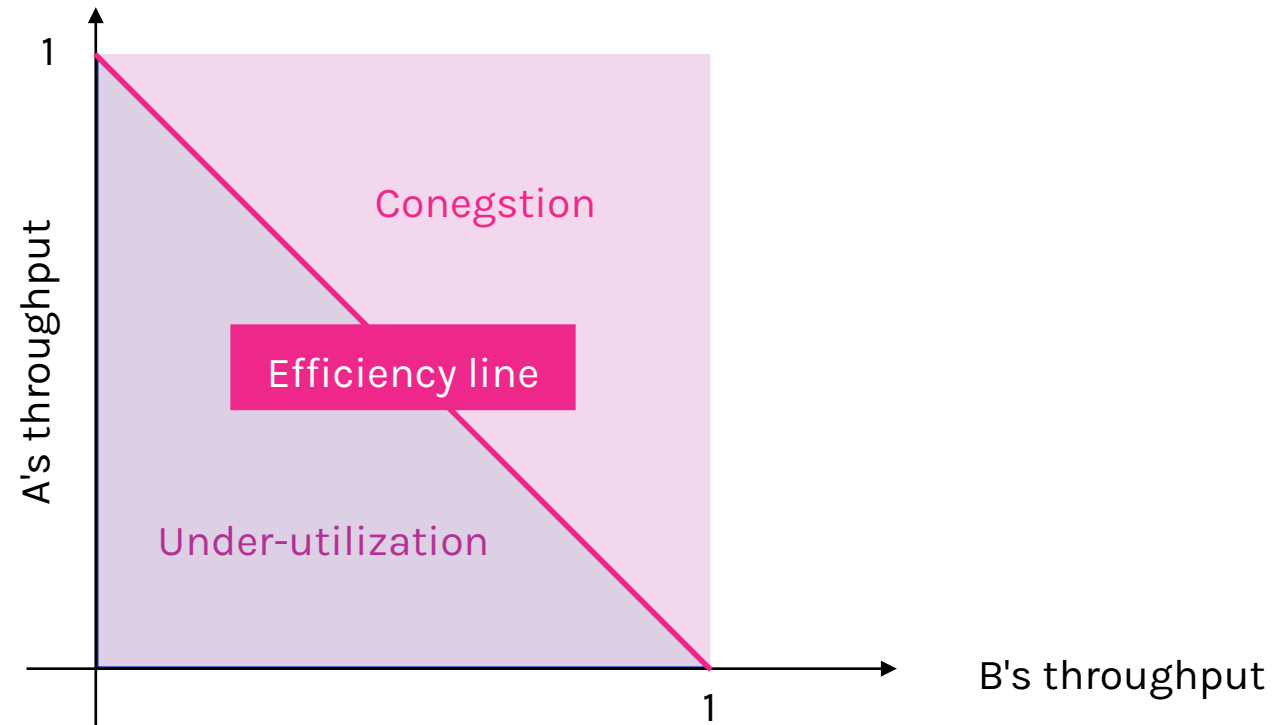


**Goal of congestion control:** bring the system as close as possible to this line and stay there

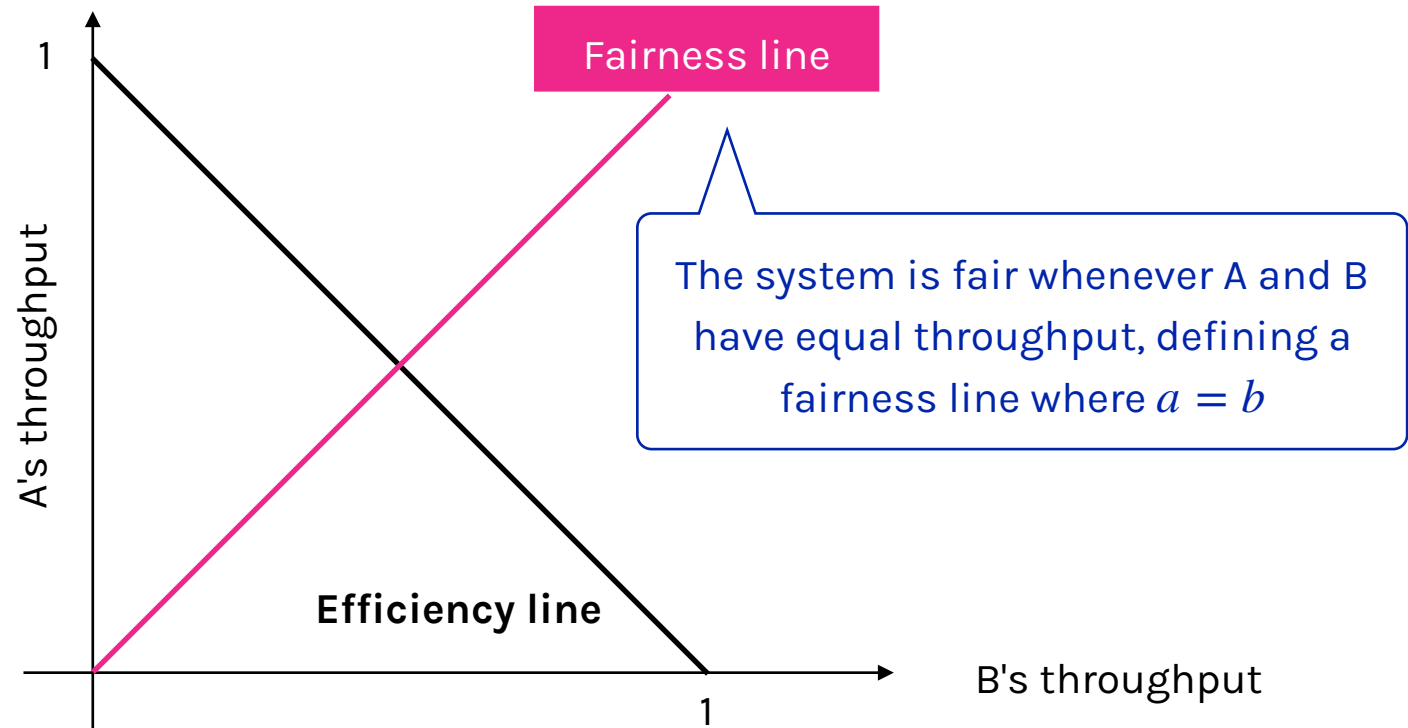
# System behavior: efficiency line



# System behavior: efficiency line

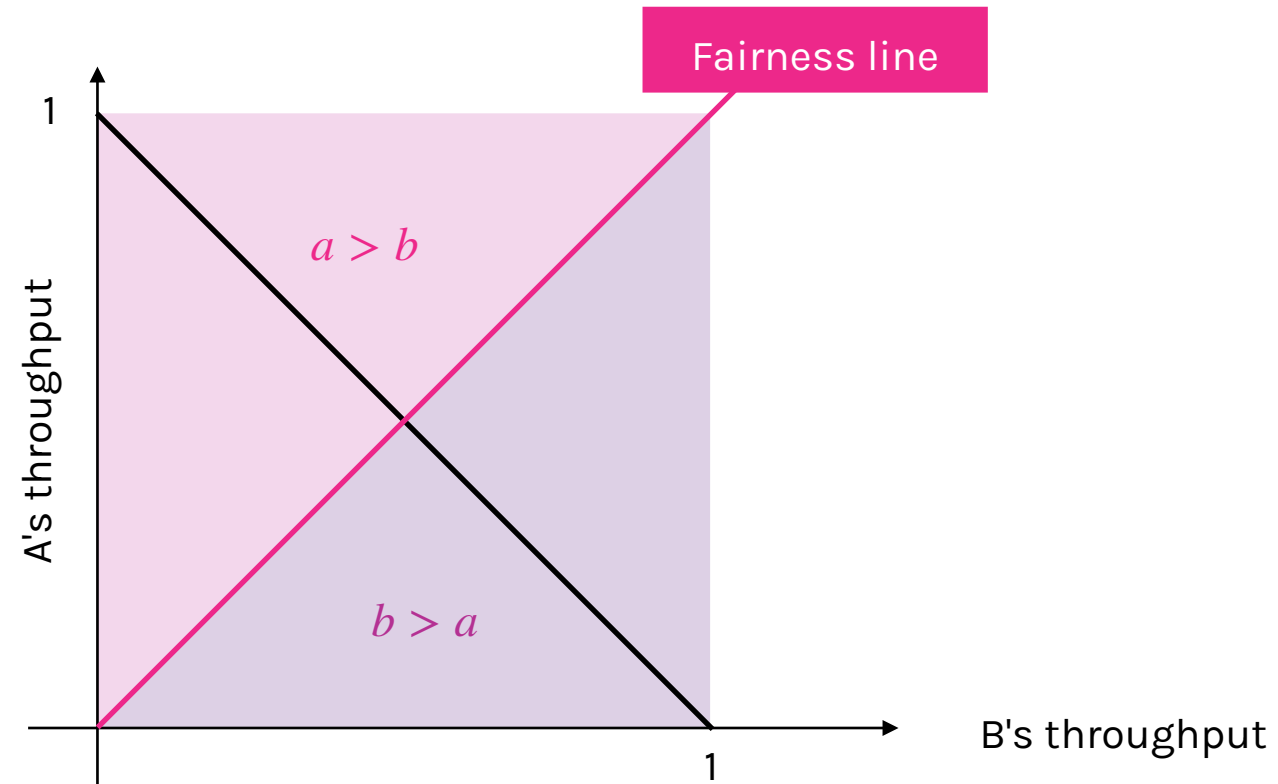


# System behavior: fairness line

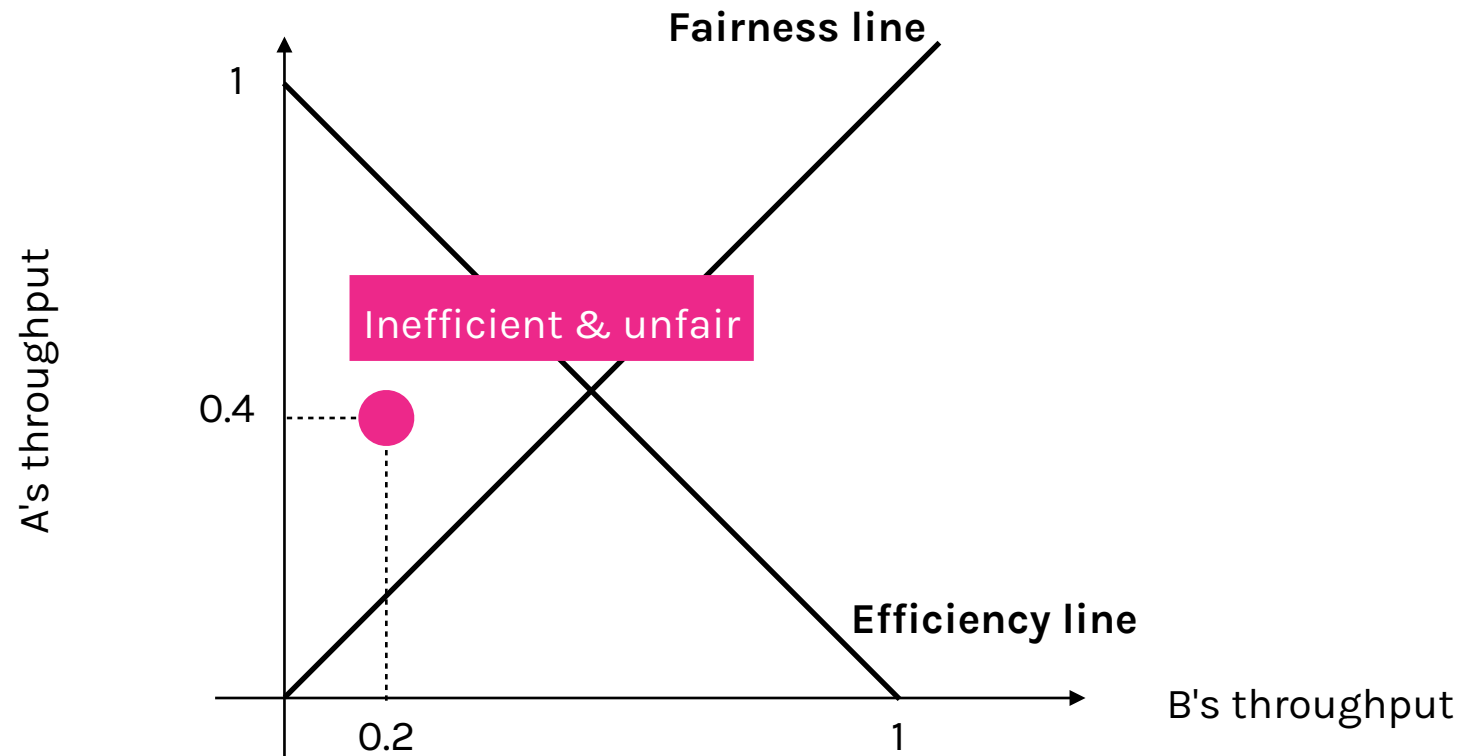


**Goal of congestion control:** bring the system as close as possible to this line and stay there

# System behavior: fairness line

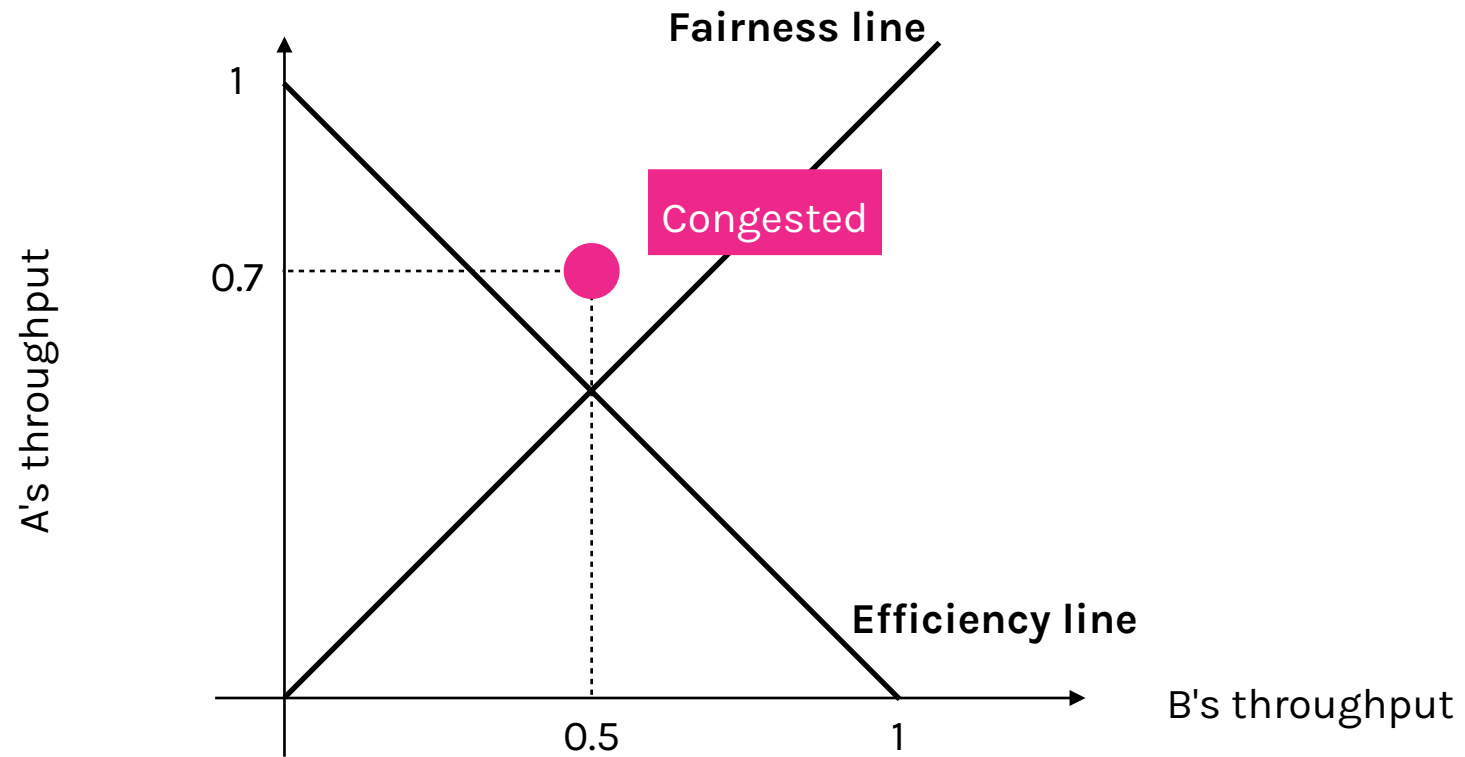


# System behavior: possible states

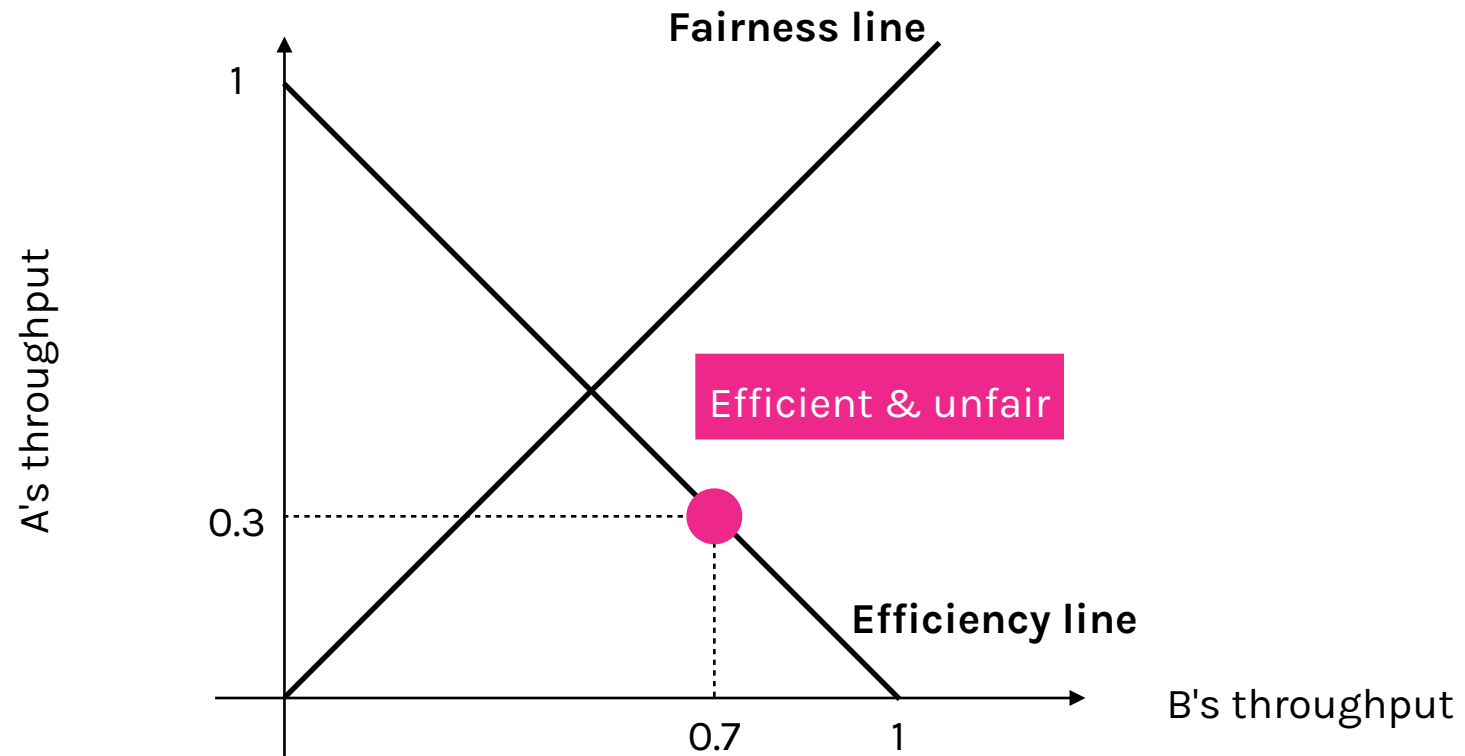




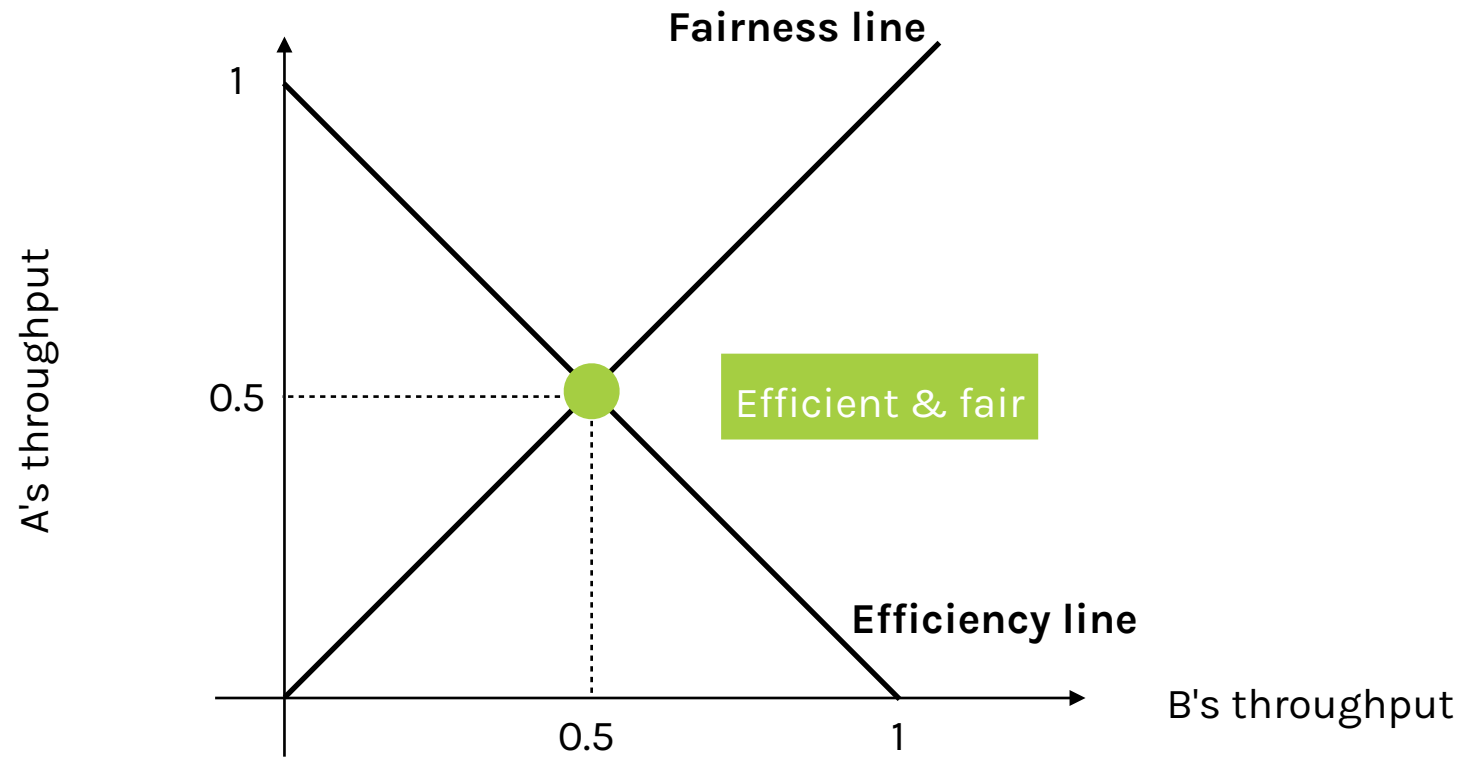
# System behavior: possible states



# System behavior: possible states



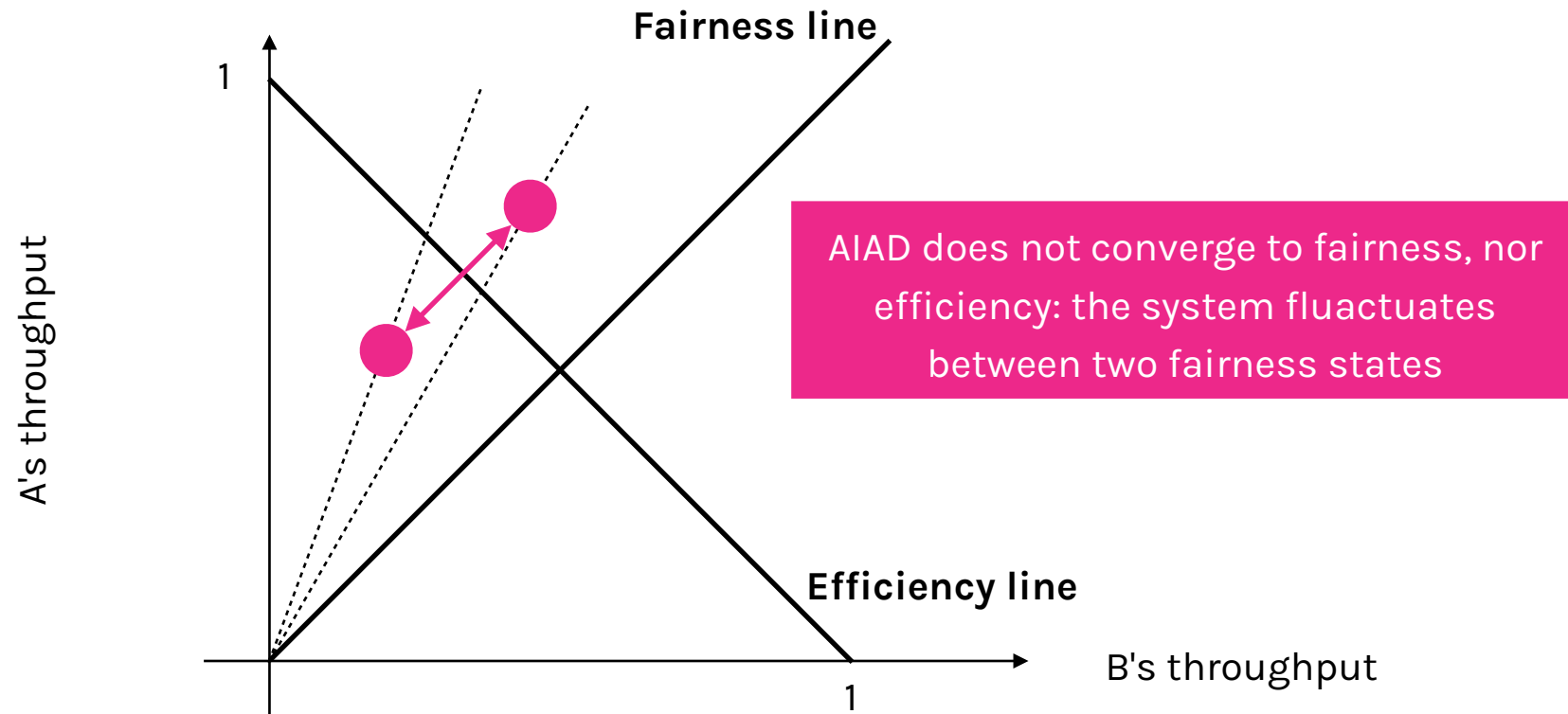
# System behavior: possible states



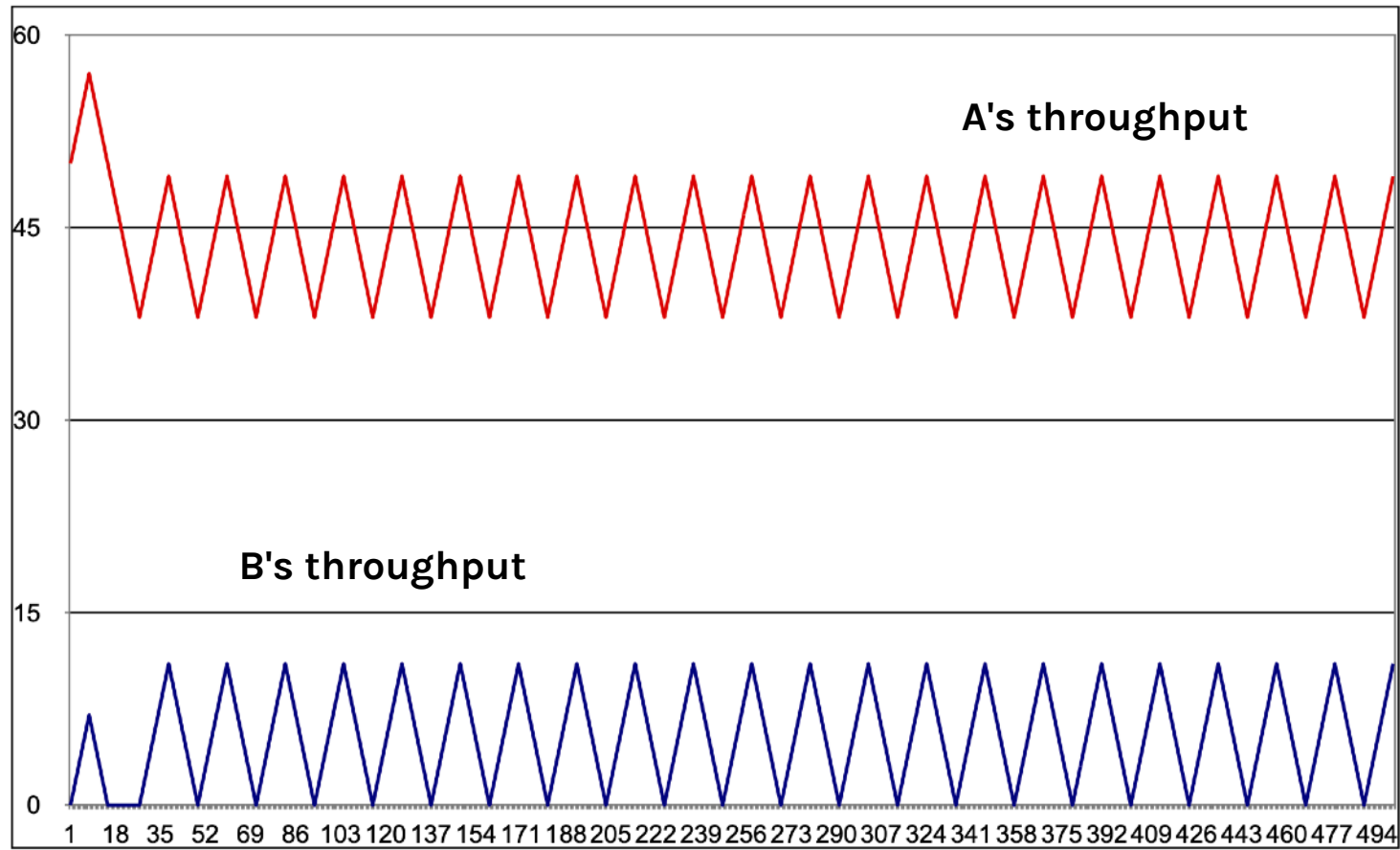
# Comparison of design choices

	Increase behavior	Decrease behavior
<b>AIAD</b>	Gentle	Gentle
<b>AIMD</b>	Gentle	Aggressive
<b>MIAD</b>	Aggressive	Gentle
<b>MIMD</b>	Aggressive	Aggressive

# AIAD fairness



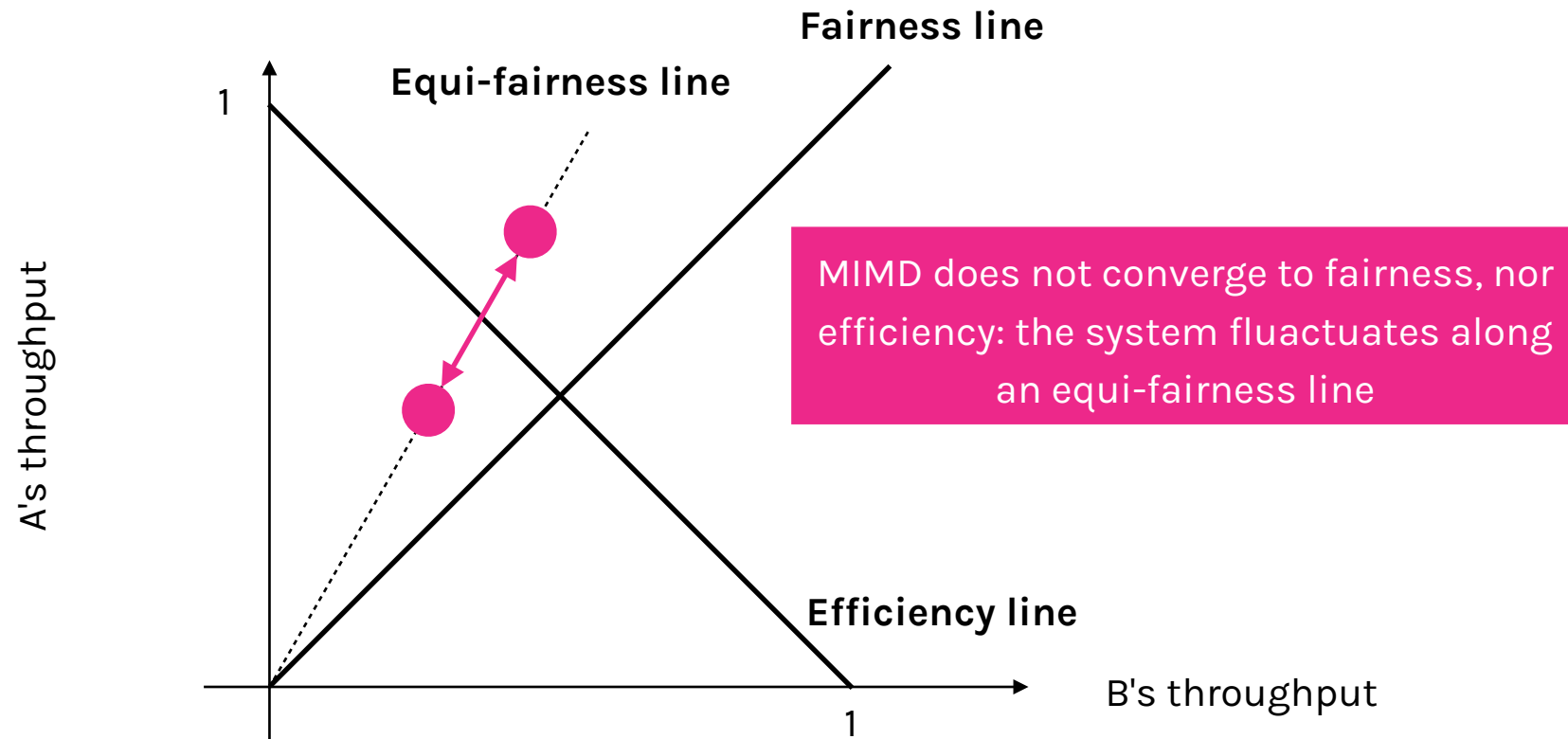
# AIAD fairness



# Comparison of design choices

	Increase behavior	Decrease behavior
AIAD	Gentle	Gentle
AIMD	Gentle	Aggressive
MIAD	Aggressive	Gentle
MIMD	Aggressive	Aggressive

# MIMD fairness

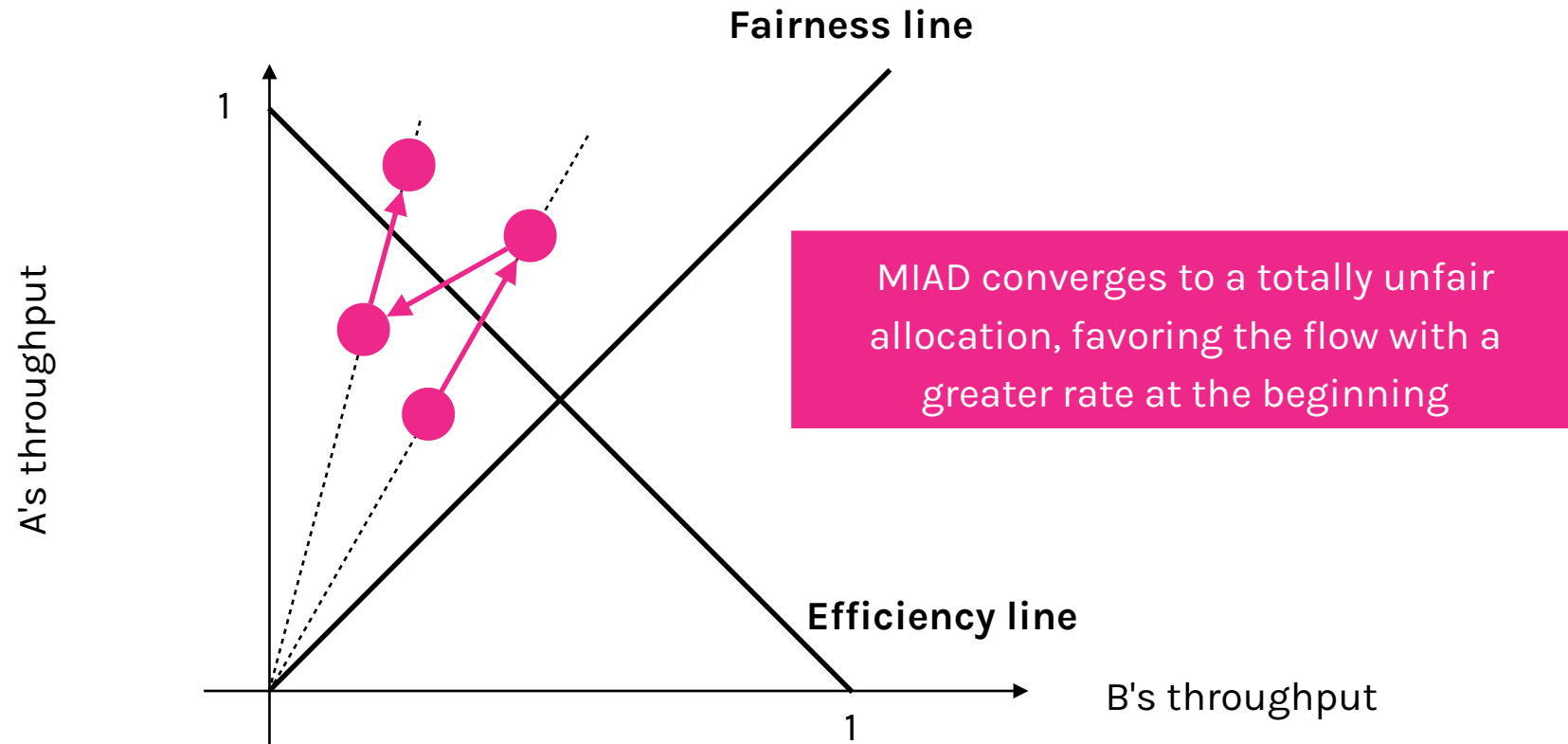




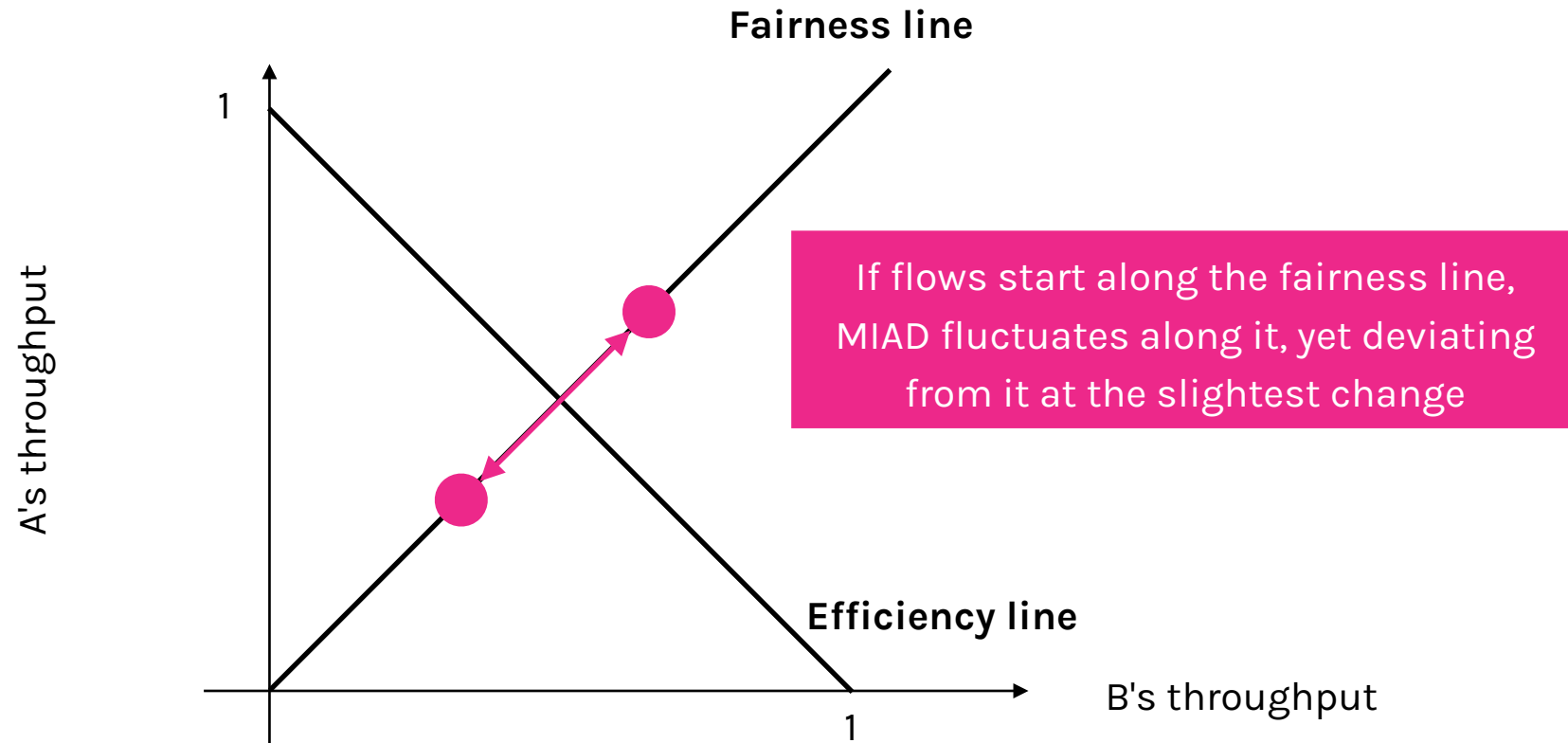
# Comparison of design choices

	Increase behavior	Decrease behavior
AIAD	Gentle	Gentle
AIMD	Gentle	Aggressive
MIAD	Aggressive	Gentle
MIMD	Aggressive	Aggressive

# MIAD fairness



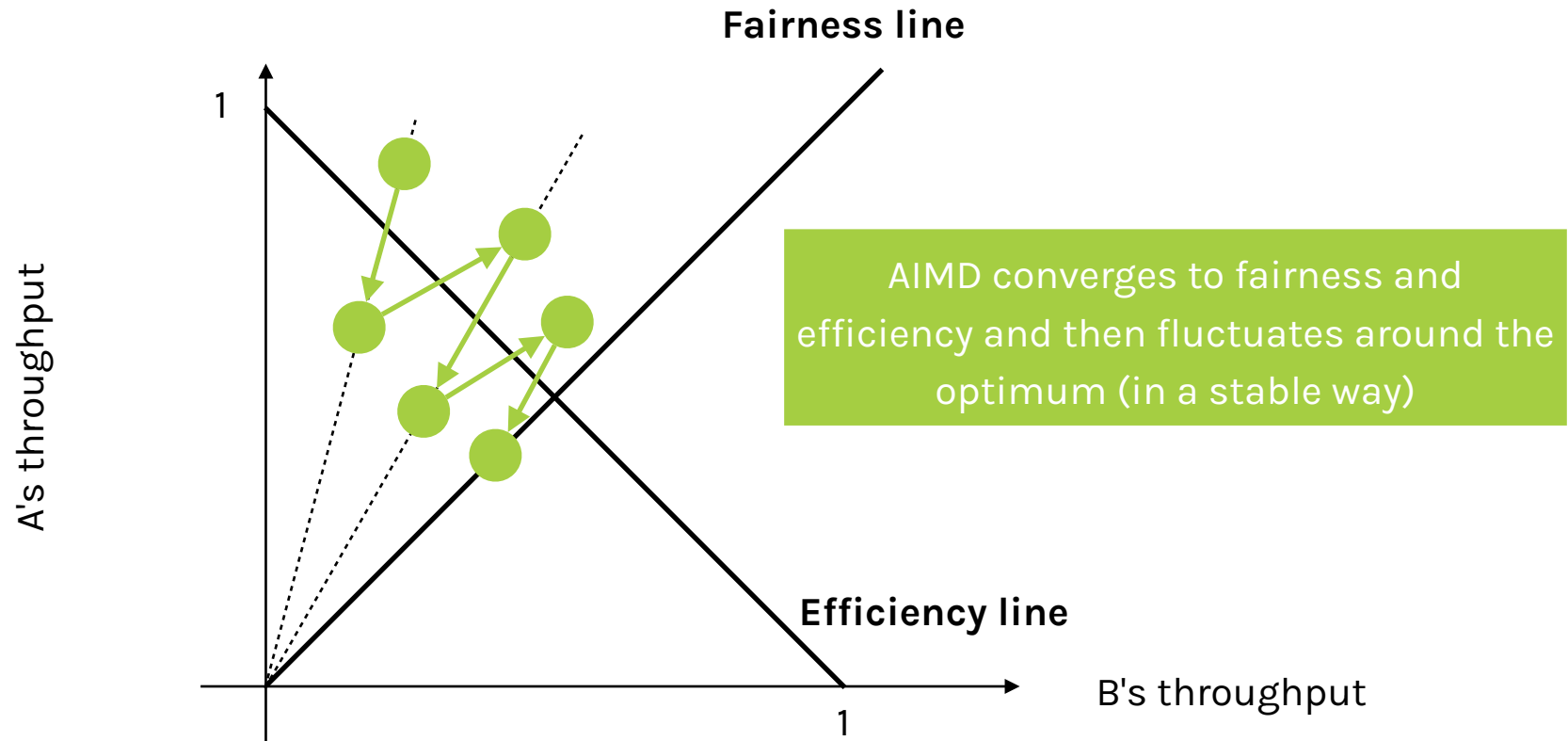
# MIAD fairness



# Comparison of design choices

	Increase behavior	Decrease behavior
AIAD	Gentle	Gentle
AIMD	Gentle	Aggressive
MIAD	Aggressive	Gentle
MIMD	Aggressive	Aggressive

# AIMD fairness



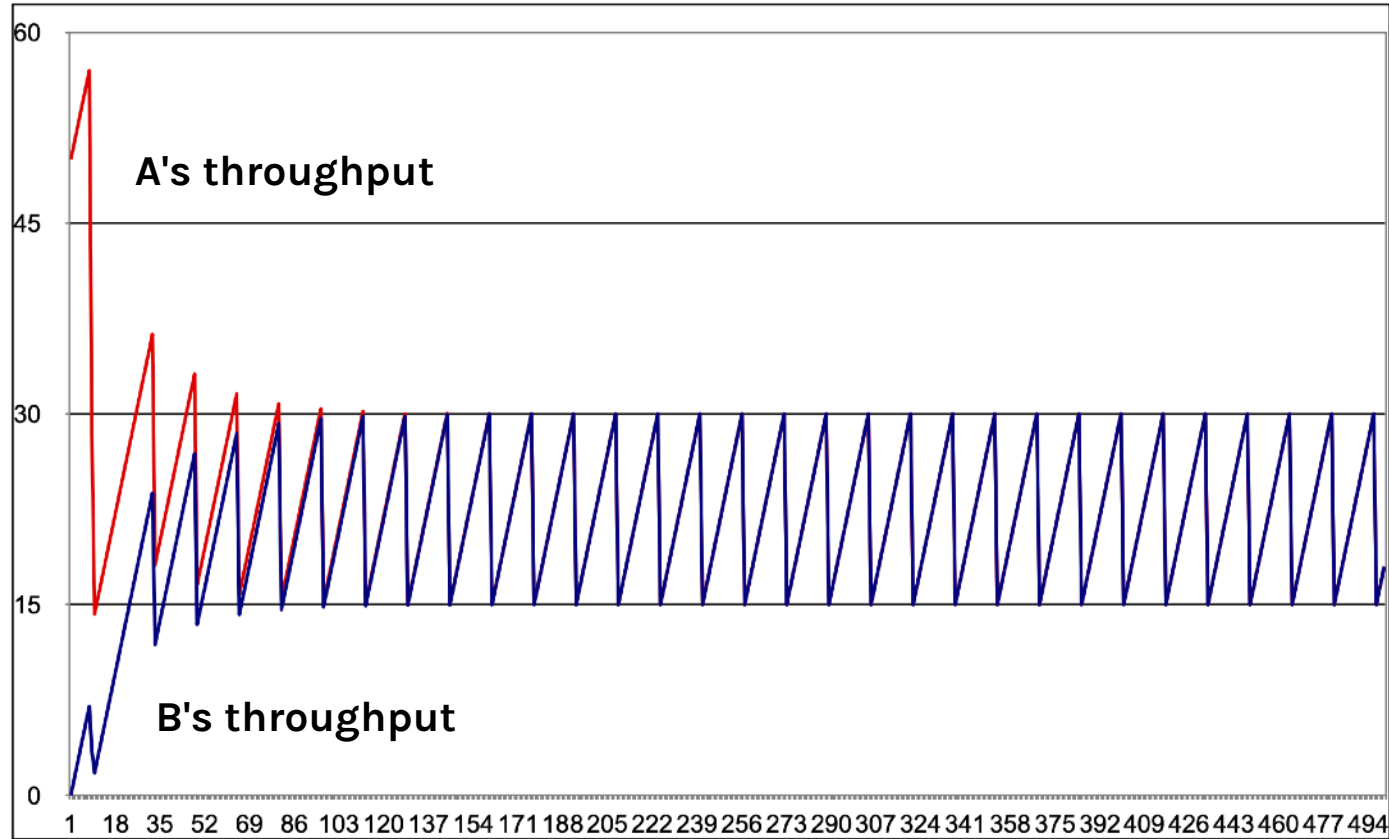
# AIMD fairness

## **Intuition**

During increase, both flows gain bandwidth at the same rate

During decrease, the flow with a higher rate releases more

# AIMD fairness



# TCP implements AIMD

## Implementation

After each ACK: increment  $CWND$  by  $1/CWND$   
(linear increase of max. 1 per RTT)

## Question

When does a sender leave slow-start and start AIMD?

Introduce a new slow-start threshold ( $ssthresh$ ), adapt it in function of congestion, e.g.,

$ssthresh = CWND/2$  (on timeout)



# TCP congestion control

## Initially:

```
    cwnd = 1  
    ssthresh = infinite
```

## New ACK received:

```
    if (cwnd < ssthresh):
```

```
        /* Slow Start */
```

```
        cwnd = cwnd + 1
```

```
    else:
```

```
        /* Congestion Avoidance */
```

```
        cwnd = cwnd + 1/cwnd
```

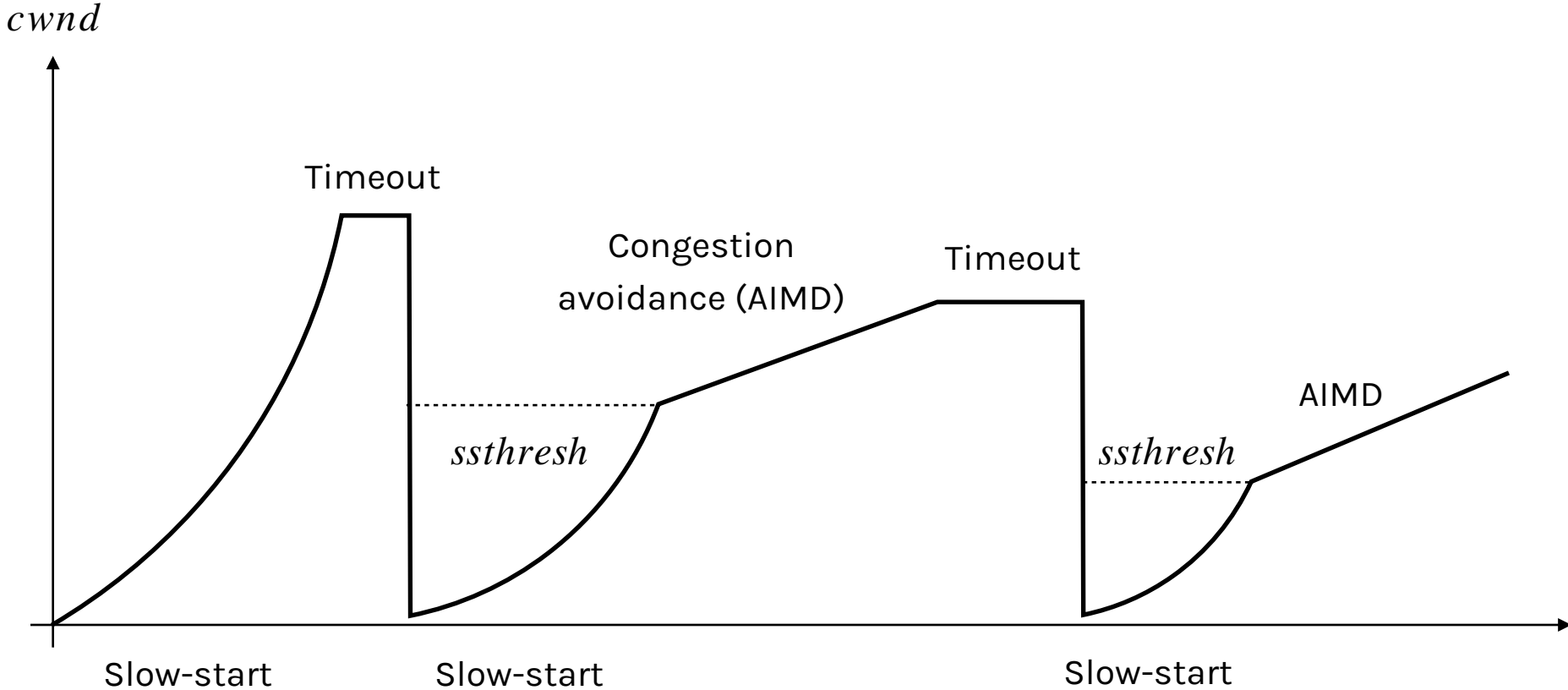
## Timeout:

```
    /* Multiplicative decrease */
```

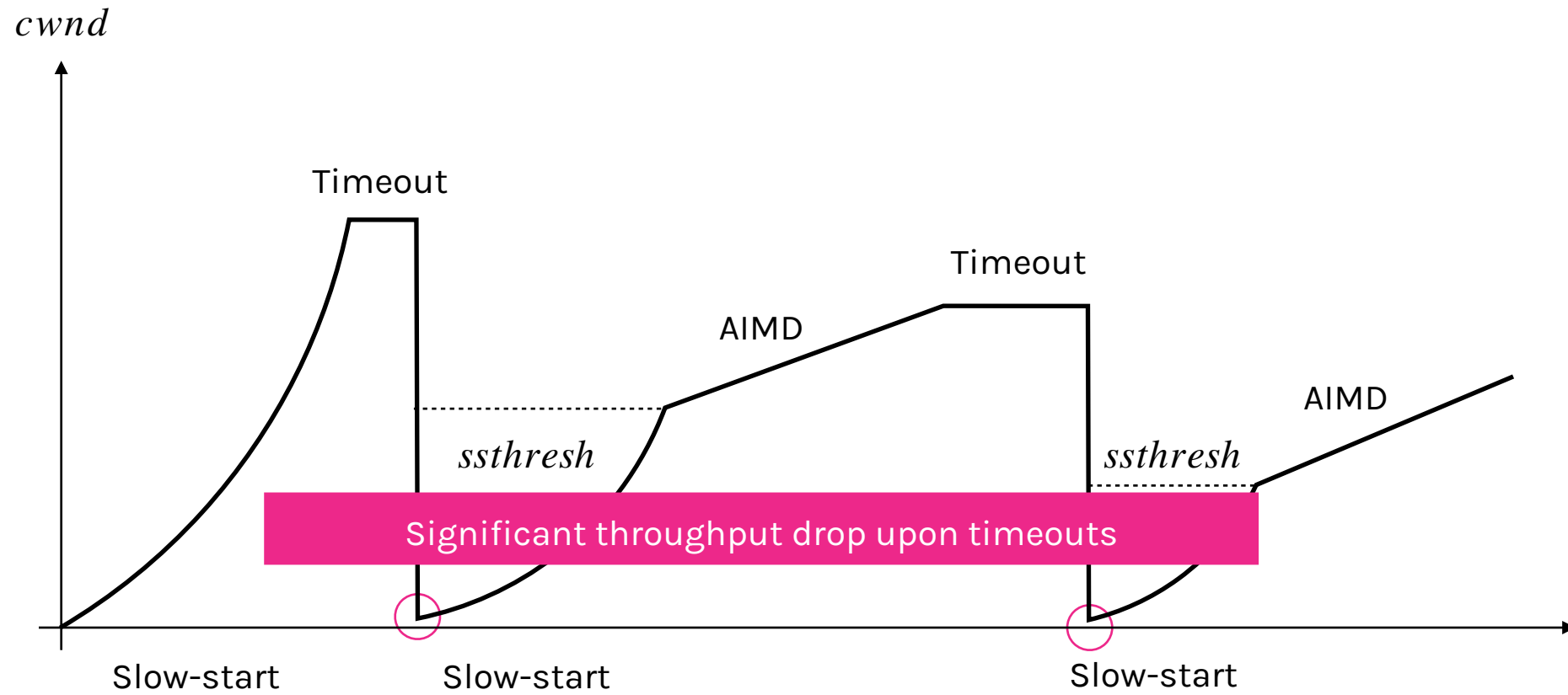
```
    ssthresh = cwnd/2
```

```
    cwnd = 1
```

# TCP Tahoe



# TCP Tahoe



# Recall: two signals for detecting packet loss

Loss detection can be done using ACKs or timeouts: differ in the degree of severity

## Duplicate ACKs

- Mild congestion
- Packets are still making it

## Timeout

- Severe congestion
- Multiple consequent losses

# TCP fast retransmit and fast recovery

## Fast retransmit

TCP automatically resends a segment after receiving 3 duplicate ACKs for it

## Fast recovery

After a fast retransmit, TCP switches back to AIMD, without going all way back to 1

# TCP congestion control (with fast recovery)

## Initially:

```
cwnd = 1
ssthresh = infinite
```

## New ACK received:

```
if (cwnd < ssthresh):
    /* Slow Start */
    cwnd = cwnd + 1
else:
    /* Congestion Avoidance */
    cwnd = cwnd + 1/cwnd
```

```
dup_ack = 0
```

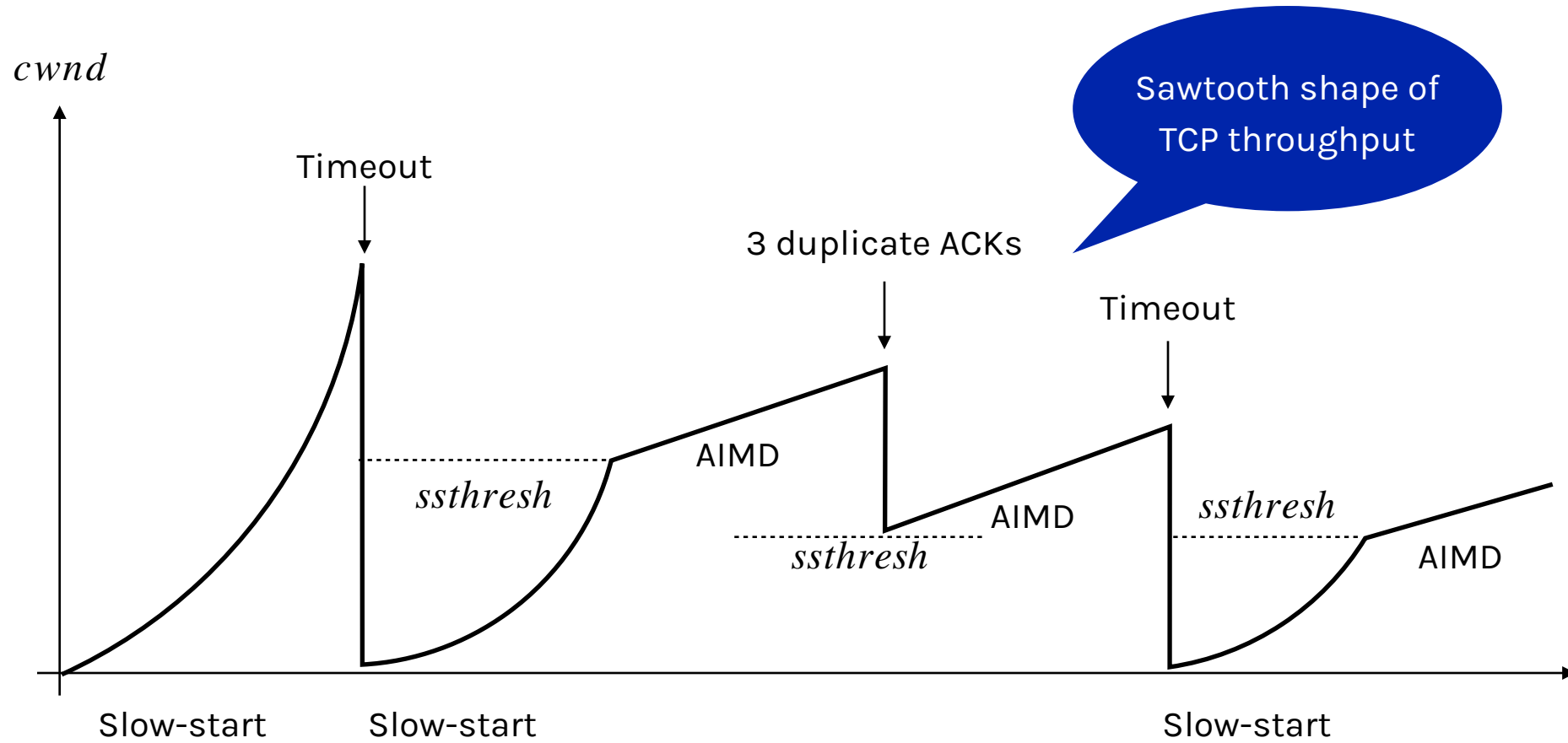
## Timeout:

```
/* Multiplicative decrease */
ssthresh = cwnd/2
cwnd = 1
```

## Duplicate ACKs received:

```
dup_ack ++
if (dup_ack >= 3):
    /* Fast recovery */
    ssthresh = cwnd/2
    cwnd = ssthresh
```

# TCP Reno



# TCP CUBIC widely used

```
linwang@Lins-MacBook-Pro:~  
net.inet.tcp.log.rate_current: 0  
net.inet.tcp.log.enable: 0  
net.inet.tcp.log.enable_usage: connection:0x00000001 rtt:0x00000002 ka:0x00000004 log:0x00000008 loop:0x00000010 local:0x00000020 gw:0x00000040 syn:0x00000100 fin:0x00000200 rst:0x00000400 dropncp:0x00001000 droppcb:0x00002000 droppkt:0x00004000 fswflow:0x00008000 state:0x00010000 synrxmt:0x00020000 output:0x00040000  
net.inet.tcp.cubic_tcp_friendliness: 0  
net.inet.tcp.cubic_fast_convergence: 0  
net.inet.tcp.cubic_use_minrtt: 0  
net.inet.tcp.cubic_minor_fixes: 1  
net.inet.tcp.cubic_rfc_compliant: 1  
net.inet.tcp.bg_target_delay: 40  
net.inet.tcp.bg_allowed_increase: 8  
net.inet.tcp.bg_tether_shift: 1  
net.inet.tcp.bg_ss_fltsz: 2  
net.inet.tcp.ledbat_plus_plus: 1  
net.inet.tcp.rledbat: 1  
net.inet.tcp.cc_debug: 0  
net.inet.tcp.newreno_sockets: 0  
net.inet.tcp.background_sockets: 15  
net.inet.tcp.use_ledbat: 0  
net.inet.tcp.cubic_sockets: 50  
net.inet.tcp.use_newreno: 0  
net.inet.tcp.mptcp_preferred_version: 1  
net.inet.tcp.backoff_maximum: 65536  
net.inet.tcp.ecn_timeout: 60  
net.inet.tcp.disable_tcp_heuristics: 0  
net.inet.tcp.mptcp_version_timeout: 1440  
net.inet.tcp.clear_tfo_cache: 0  
net.inet.tcp.flow_control_response: 1  
net.inet.tcp.log_in_vain: 0  
net.inet.tcp.ack_strategy: 1
```

```
lin.wang@erdos: ~  
tcp_allowed_congestion_control  
tcp_app_win  
tcp_autocorking  
tcp_available_congestion_control  
tcp_available_ulp  
tcp_base_mss  
tcp_challenge_ack_limit  
tcp_comp_sack_delay_ns  
tcp_comp_sack_nr  
tcp_comp_sack_slack_ns  
tcp_congestion_control  
tcp_dsack  
tcp_early_demux  
tcp_early_retrans  
tcp_ecn  
tcp_ecn_fallback  
tcp_fack  
tcp_fastopen  
tcp_fastopen_blackhole_timeout_sec  
tcp_fastopen_key  
tcp_fin_timeout  
tcp_frto  
tcp_fwmark_accept  
tcp_invalid_ratelimit  
tcp_keepalive_intvl  
tcp_keepalive_probes  
tcp_keepalive_time  
tcp_l3mdev_accept  
tcp_limit_output_bytes  
tcp_min_tso_segs  
tcp_moderate_rcvbuf  
tcp_mtu_probe_floor  
tcp_mtu_probing  
tcp_no_metrics_save  
tcp_no_ssthresh_metrics_save  
tcp_notsent_lowat  
tcp_orphan_retries  
tcp_pacing_ca_ratio  
tcp_pacing_ss_ratio  
tcp_probe_interval  
tcp_probe_threshold  
tcp_recovery  
tcp_reordering  
tcp_retrans_collapse  
tcp_retries1  
tcp_retries2  
tcp_rfc1337  
tcp_rmem  
tcp_rx_skb_cache  
tcp_sack  
tcp_slow_start_after_idle  
tcp_stdurg  
tcp_syn_retries  
tcp_synack_retries  
tcp_syncookies  
tcp_thin_linear_timeouts  
tcp_timestamps  
tcp_tso_win_divisor  
lin.wang@erdos:~$ cat /proc/sys/net/ipv4/tcp_congestion_control  
cubic  
lin.wang@erdos:~$
```

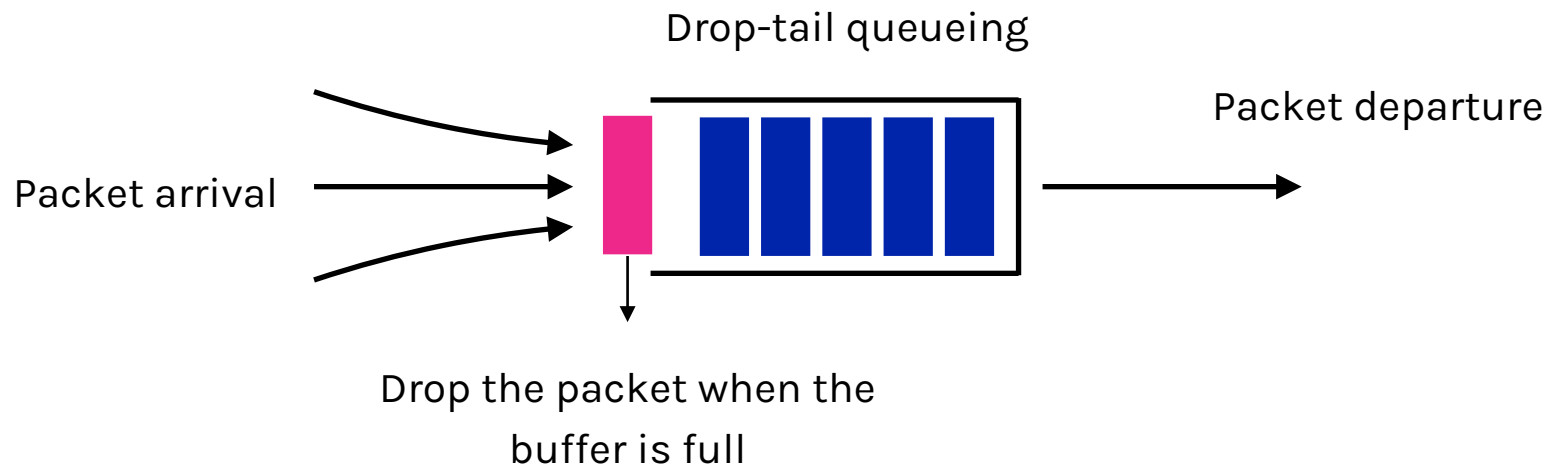
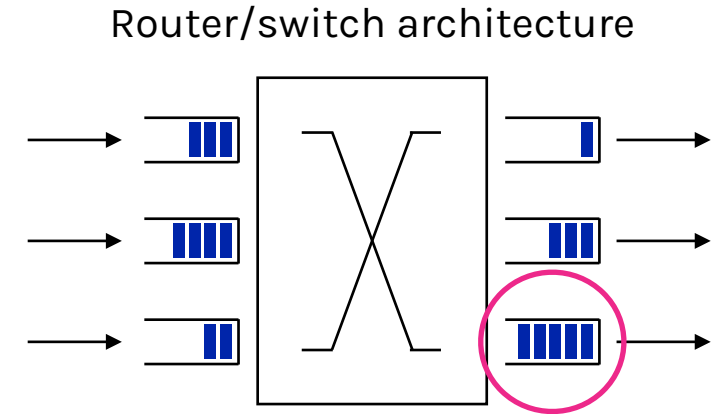


# Explicit Congestion Notification (ECN)

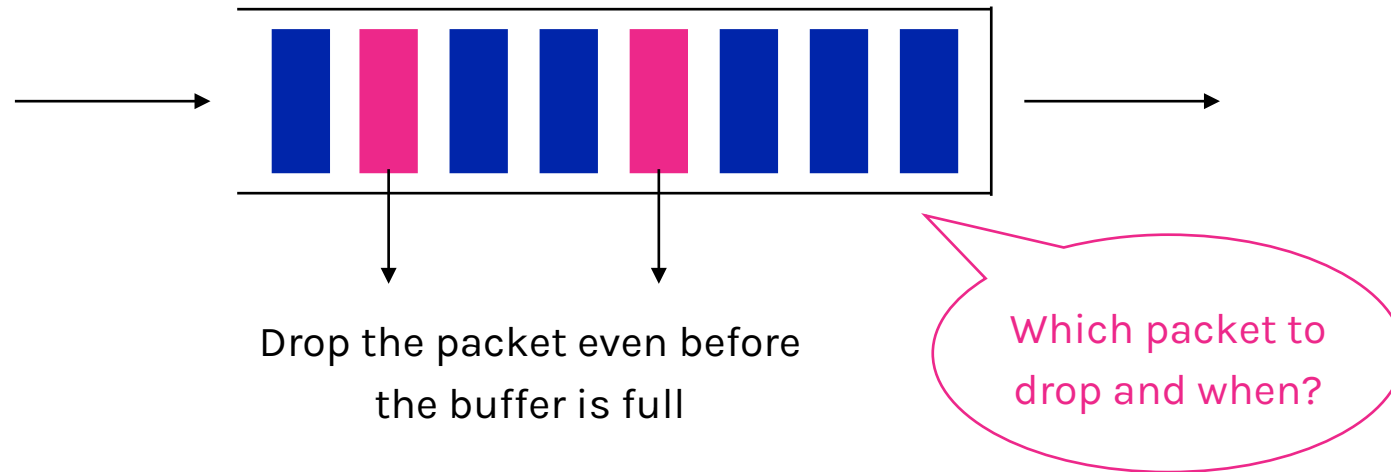


# Why congestion happens?

Can we avoid the buffer to be full completely?

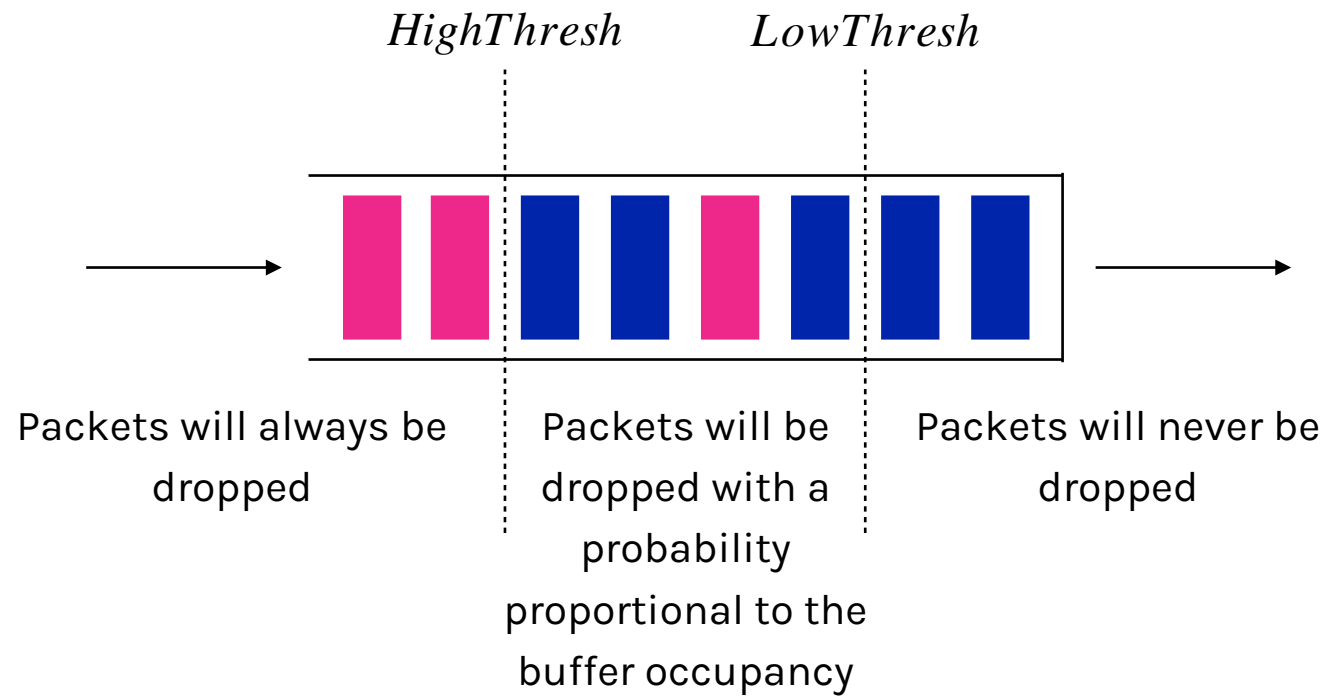


# Active queue management (AQM)



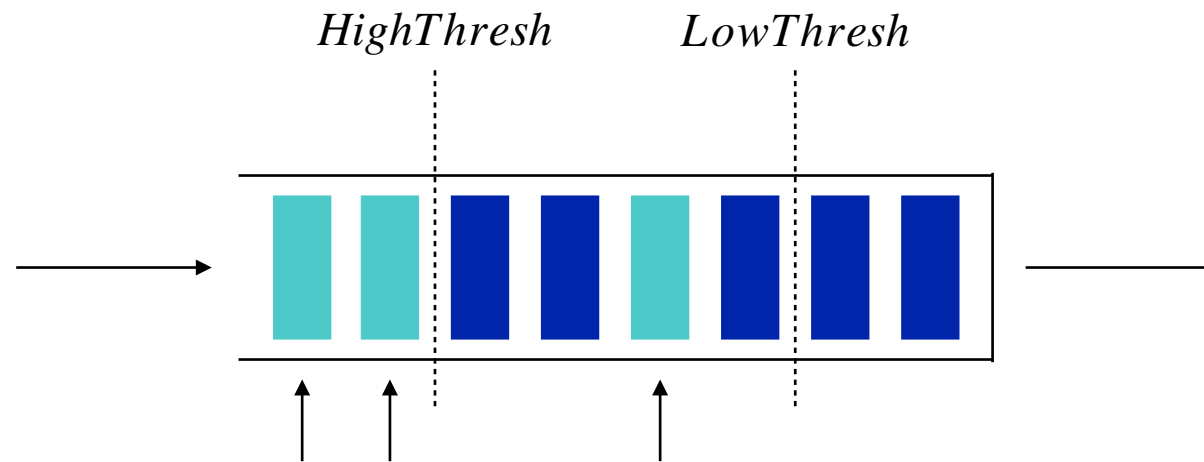
The dropped packet will serve as a signal for the sender to adjust its sending rate

# AQM policy: random early detection (RED)



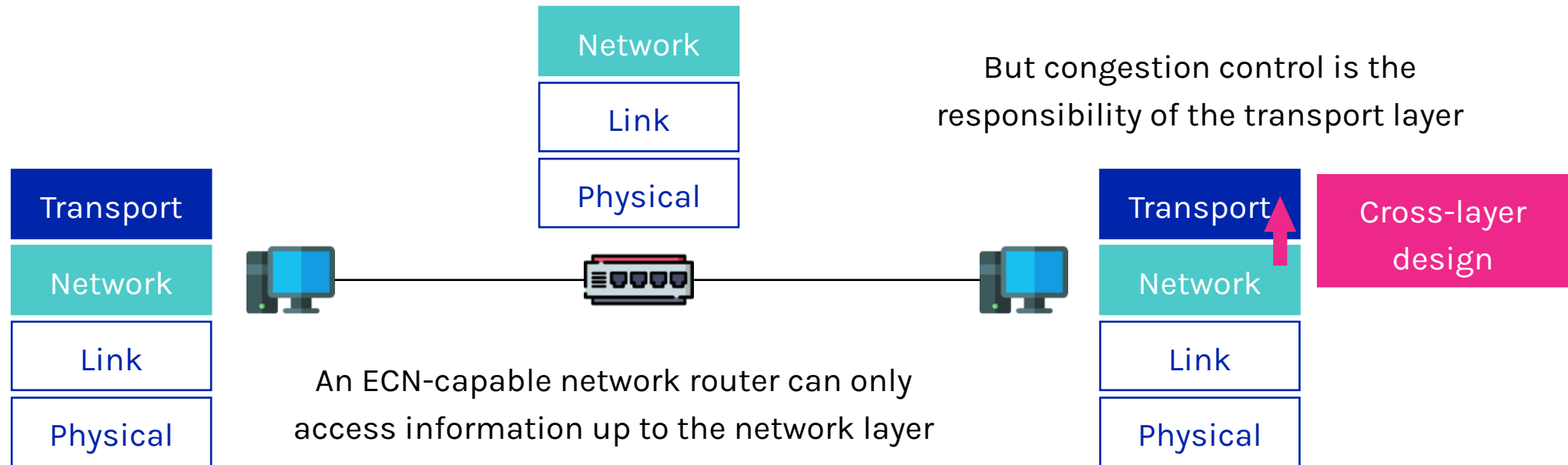
**"Why should I drop perfectly good packets  
when I still have free buffer space?"**

# Explicit congestion notification (ECN)



Instead of dropping, **mark** these packets and **notify the sender** to reduce its sending rate

# ECN mechanism



# ECN mechanism: IP header

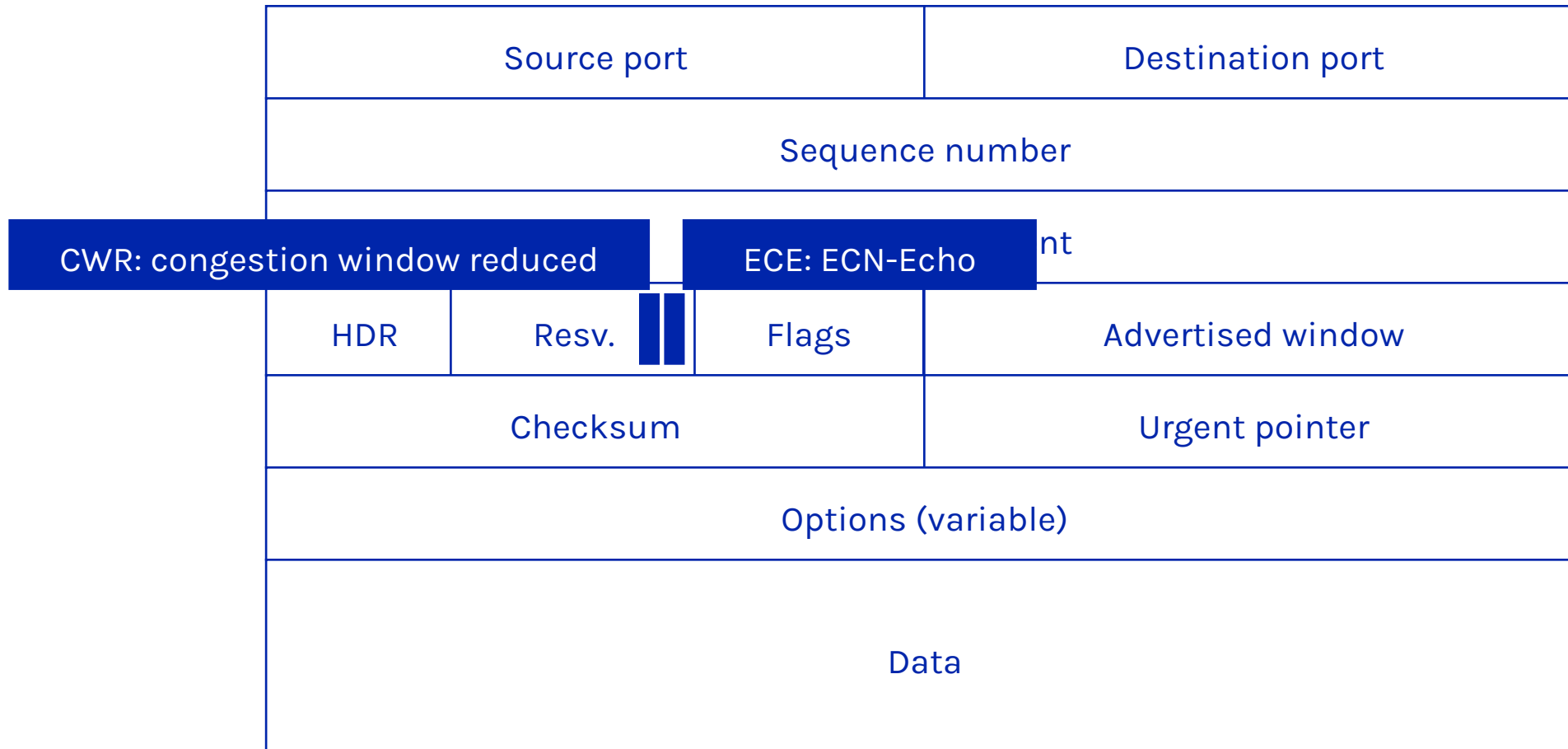
Last two bits of ToS for ECN: 00 (Not-ECT), 01/10 (ECT), 11 (CE)

Ver	HLEN	ToS	Total length	
Identification		Flags	Fragment offset	
Time to live	Protocol	Header checksum		
Source IP address				
Destination IP address				
Options (if any)				
Payload (data)				

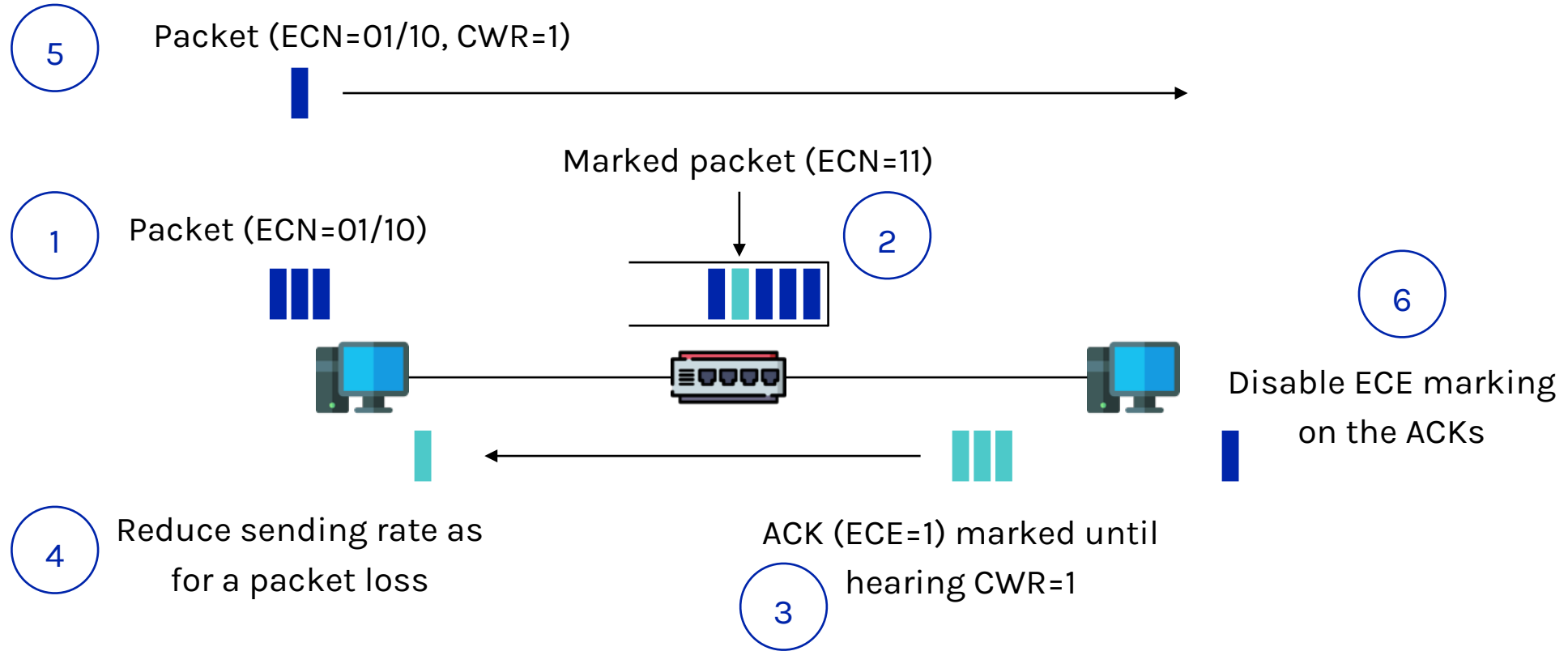
ECT: ECN-capable transport      CE: Congestion experienced



# ECN mechanism: TCP header



# ECN example



# Summary

## Network congestion

- Why congestion?
- Congestion quantification
- Congestion collapse of the Internet
- Congestion control goals
- Congestion window

## Congestion detection and control

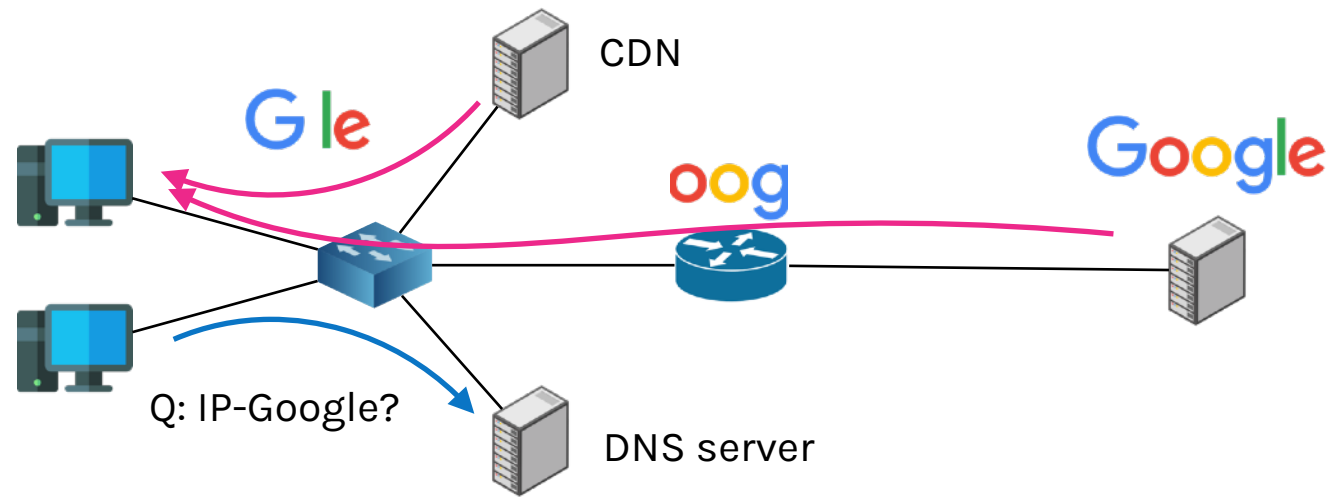
- Congestion signals
- Bandwidth estimation (slow-start)

- Bandwidth adaptation
- Fairness
- AIMD
- Fast recovery

## Explicit Congestion Notification

- AQM + RED
- ECN

# Next time: application layer



How does the Web work?

## Further reading material

**Andrew S. Tanenbaum, David J. Wetherall. Computer Networks (5th edition).**

- Section 6.5.10: TCP Congestion Control

**Larry Peterson, Bruce Davie. Computer Networks: A Systems Approach.**

- Chapter 6: Congestion Control