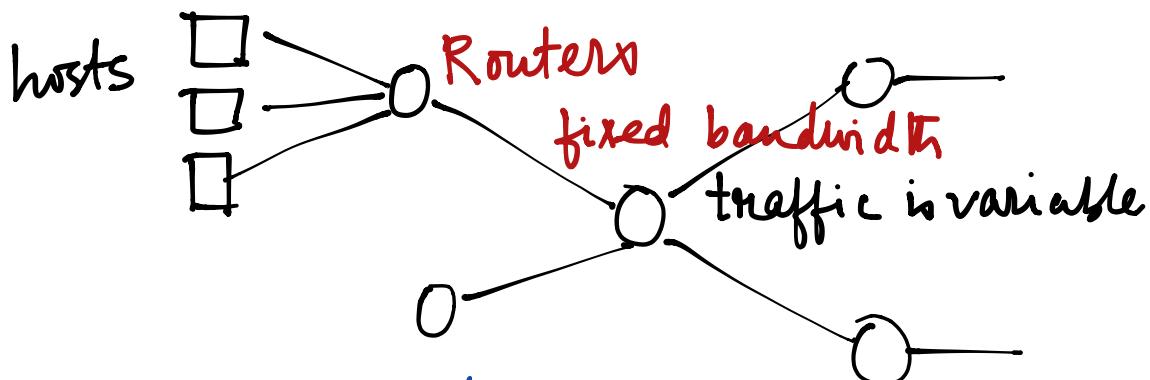


## TCP Congestion Control

Issues : Efficiency and Fairness.



Bandwidth needs to be allocated dynamically to avoid congestion.

### Bandwidth Allocation Model

- recall*
- ① NW layer provides feedback
  - ② TRA layer adjusts the traffic [changes the window size]
  - ③ Efficiency & Fairness

TCP : closed loop - uses network feedback to adjust the rate

host based - network layer devices

just pass the congestion information,  
doesn't adjust rates.

window based - the TCP sliding window  
size is adjusted through the feedback.

## Additive Increase and Multiplicative Decrease

### AIMD

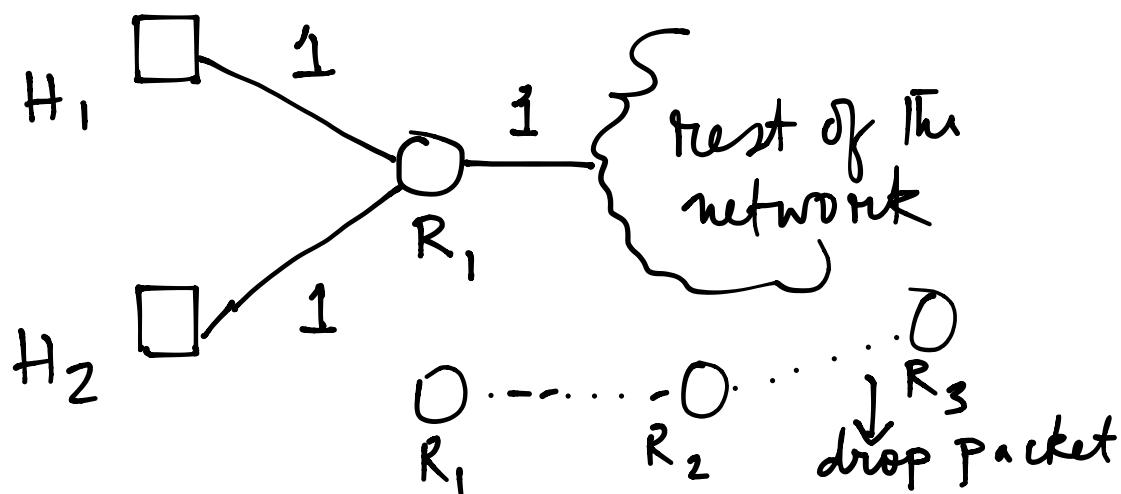
$$w_{t+1} = \begin{cases} w_t + a, & \text{if no congestion} \\ w_t \times b, & \text{otherwise} \end{cases}$$

$w_t$  : window size at  $t$

$a, b$  are typically decided beforehand

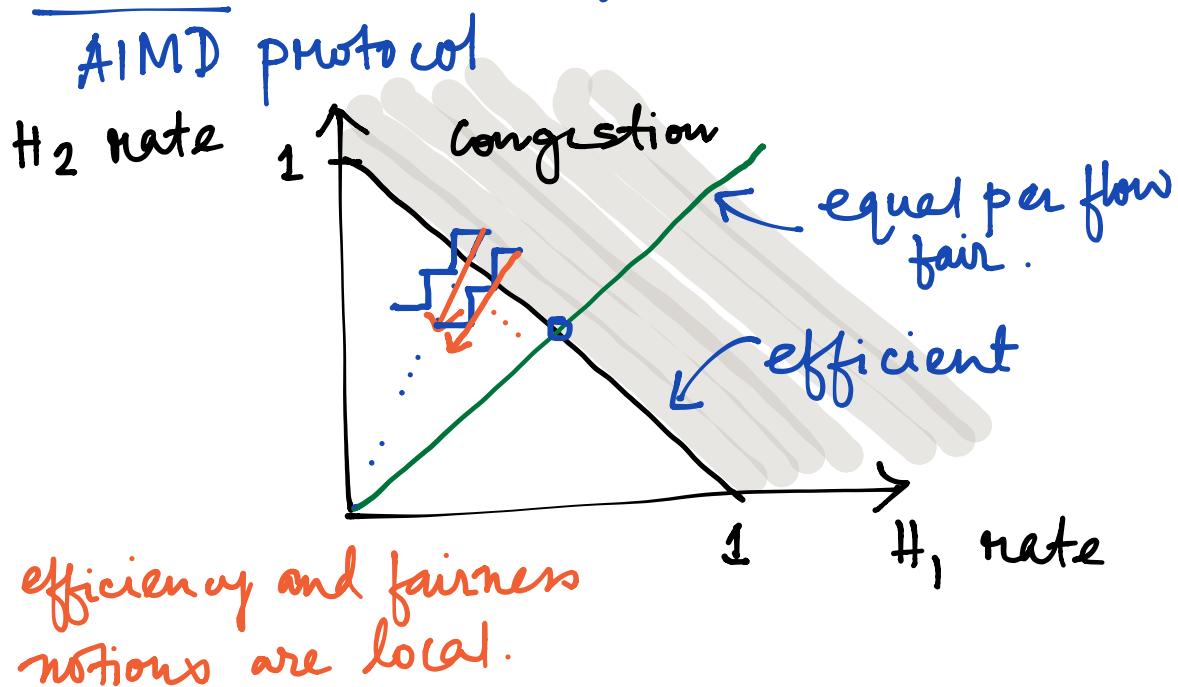
$a = 1$  packet ,  $b = 1/2$  typical

Ex:



Congestion could be at any other router but it will have a back-pressure till the host where the trouble is.

In the example, both hosts start with some Window which can cause congestion.



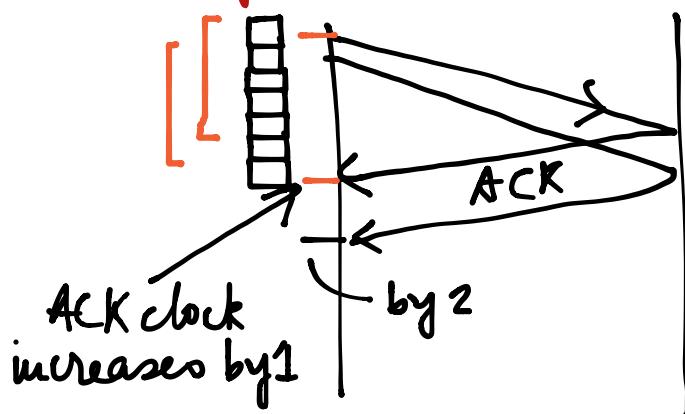
Note : Even if hosts are using dissimilar windows in the beginning, the protocol forces them to be such that the rate becomes equal - even when one uses the whole BW.

## MIMD endnotes

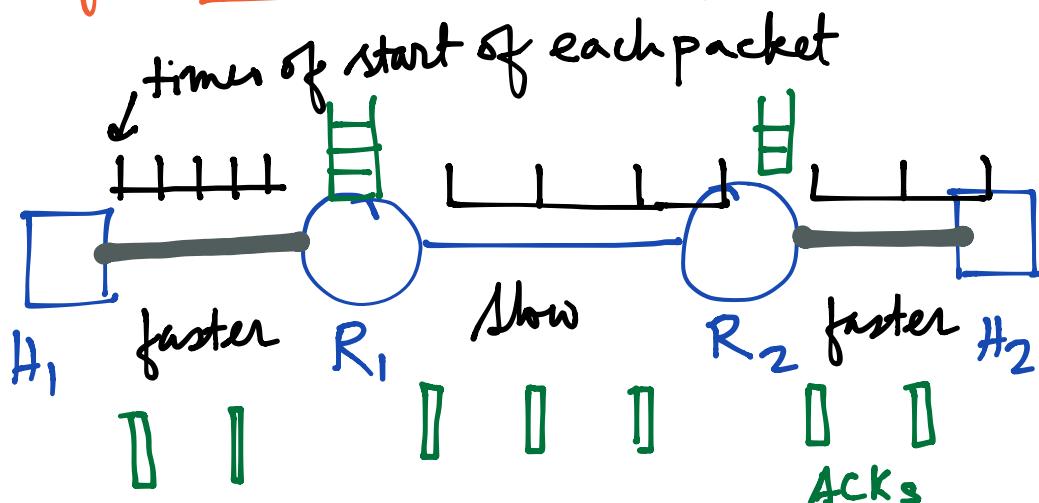
- ① decentralized
- ② convergence is guaranteed. not the case for MIMD, SIMD, SIMD
- ③ Needs simple binary feedback from the network

Feedback signal	Pros/cons	Example
packet loss	Easy to detect Slow convergence	TCP NewReno Cubic TCP
packet delay	Early detection computation-heavy and may be inaccurate	Compound TCP
winter feedback	Early detection require winter to send info.	TCPs with explicit congestion notification

## Sliding Window ACK clocking



Why ACK clocking is useful?



the window  
will slide at this rate - automatically adjusting  
the rate of transmission.

- network smooths out the traffic
- ACK clocking ensures no routers  
are having buffer overflow.

TCP uses a sliding window because of ACK clocking

Sliding window controls how many segments in a window [the size is variable]

- called **congestion window** or **cwnd**
  - Rate =  $\frac{\text{cwnd}}{\text{RTT}}$
- 

### More details of TCP

How does TCP implement AIMD?

Stage 1: SLOW START  $\rightarrow$  AIMD

- Recor**
- Sender uses cwnd
  - Sender uses packet loss as an indication of congestion
  - TCP handles wide variety of rates and RTTs (need to dynamically adjust)

Problem : there is a  $cwnd_{OPT}$ . but how should TCP know it ?

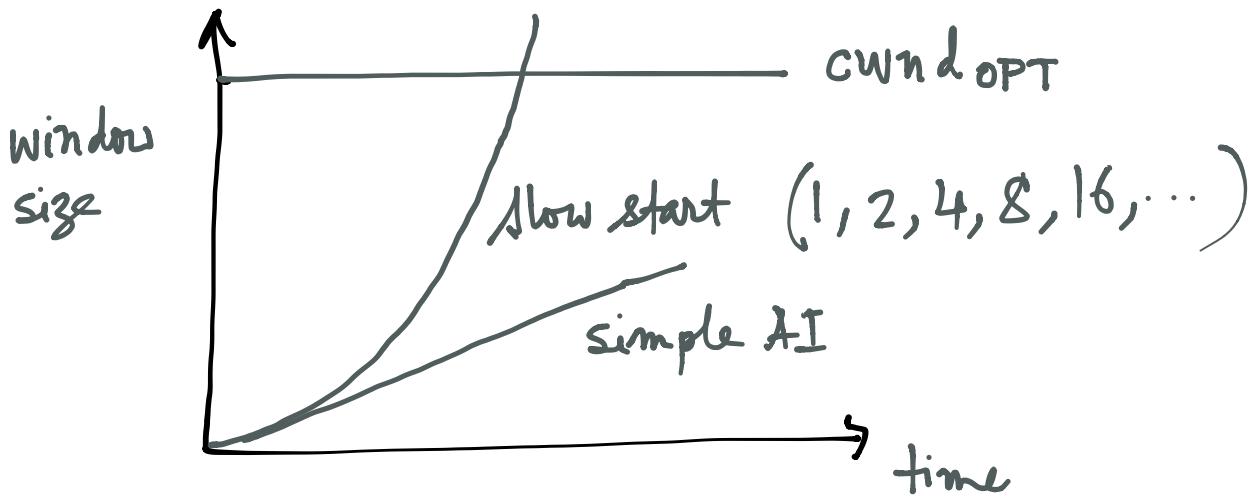
- ① Fixed window size - could be large or small
- ② AI from 1 could be too slow to reach the optimal

---

Slow Start (counterintuitive name)

does it with exponential growth till a threshold called slow start threshold

$ss\_thres$



The window size is doubled every time

but if the exponentiation is continued it will congest the network hugely since the increase of window size is large

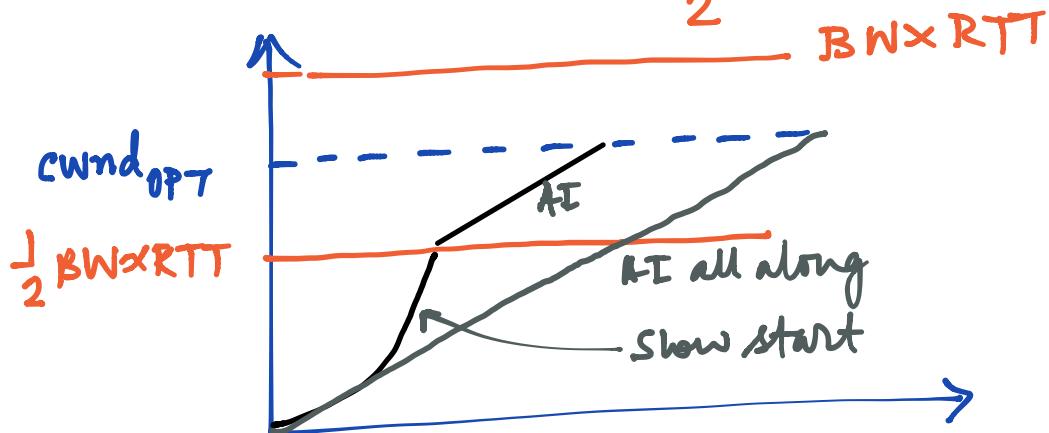
switches to AI in advance

RTT is usually calculated between the packet and ACK

$BW \times RTT \approx$  reasonable proxy for the max window size

Threshold heuristic

$$ss\text{thres} = \frac{BW \times RTT}{2}$$



## Example: TCP Tahoe

- Start with cwnd = 1

when  $cwnd < ssthresh$

- $cwnd += 1$  for every ACK

[when all ACKs for a window comes, the cwnd doubles]

when  $cwnd \geq ssthresh$

- $cwnd += \frac{1}{cwnd}$  packets per ACK

[increase one packet for the whole window]

### AI step

ssthresh is also dynamic

if packet loss occurs

$$ssthresh = \frac{cwnd}{2}$$

begin slow start after timeout  
with  $cwnd = 1$

} not  
an MD  
step

### Why not MD instead?

- loss of packet is detected via timeouts
  - which is a pessimistic estimate
- ACK clock synchronization is lost.

Therefore, no way to get back to the state where the congestion occurred.

If we want to retain the ACK clock sync and do an MD, then we need to detect loss before a timeout

(avoidance rather a recovery - TCP Reno)

- Slow start is a reasonable compromise.

---

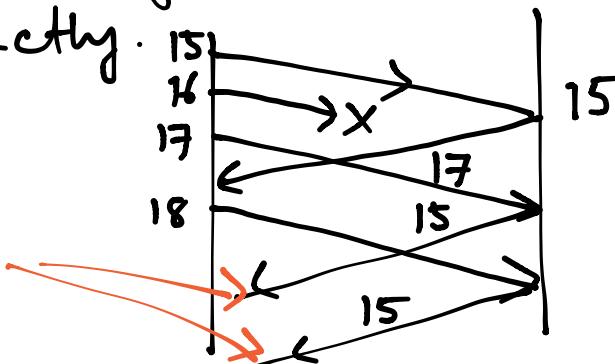
Stage 2: FAST RETRANSMIT & FAST RECOVERY  
(MD part)

TCP infers the loss from ACKs

uses cumulative ACKing

- sends the highest in-order seq. #  
e.g., if packets till 15 was received correctly, it will send that ACK until 16 is received, even if 17, 18, ... are received correctly.

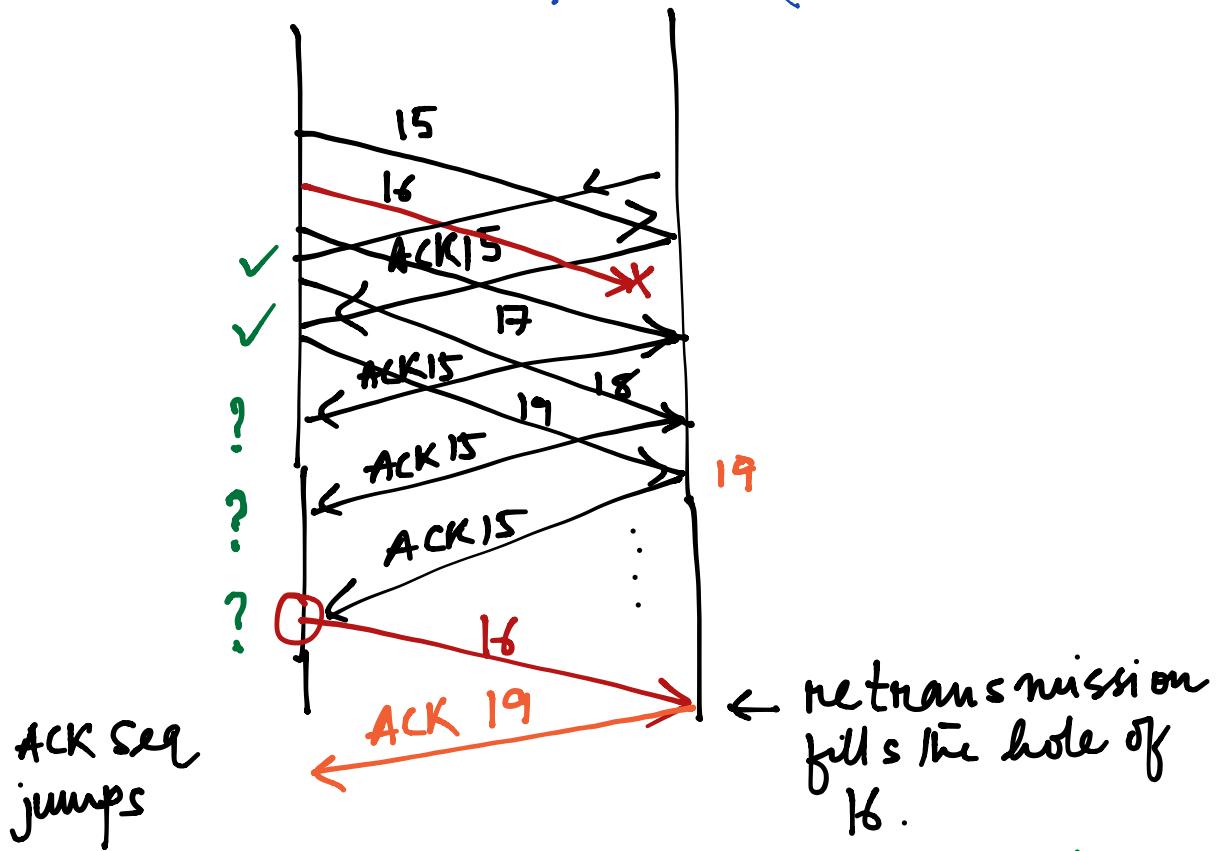
duplicate  
ACKs



Duplicate ACKs hints that some data didn't reach. - i.e., The next packet

## Fast Retransmit

- Treat 3 duplicate ACKs as the loss of the next packet.
- Retransmit that packet. (16 in our example)



Fast retransmit repairs single segment loss (can improvise this for multi segment los)

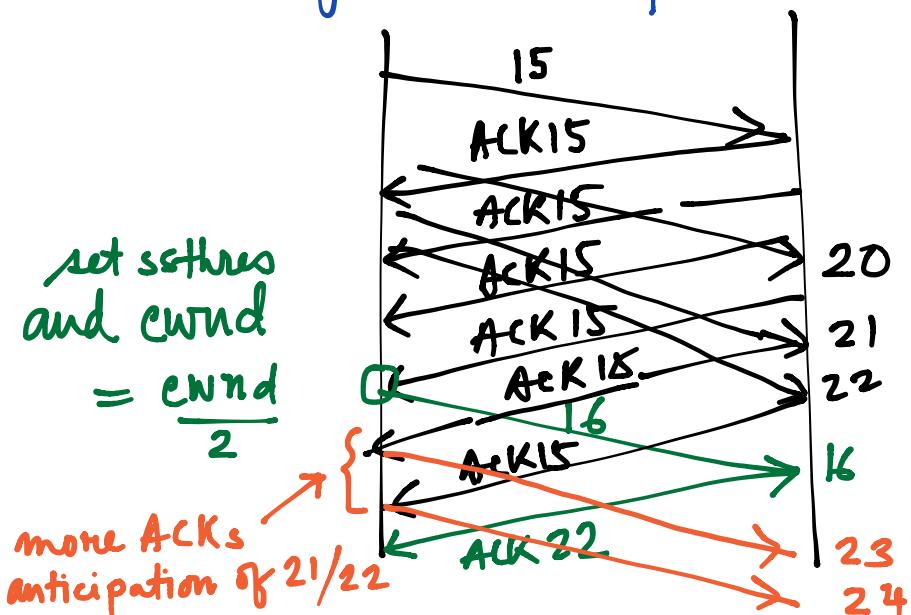
- Still quiet time at sender/receiver when the retransmission is taking place.
- MD pending.

Observation: the duplicate ACKs say

- ① the next packet didn't arrive
- ② some later packet arrived. - i.e. no queuing overflow

## Fast Recovery

- First fast retransmit, then MD.
- pretend duplicate ACKs are expected ACKs for the next packets.



fast recovery

With fast retransmit, repairs a single segment loss quickly and keeps the ACK clock running.

While the 3<sup>rd</sup> duplicate ACK comes, it reduces the window by  $\frac{1}{2}$ , not reset to 1.

---

Example: TCP Reno uses slow start, fast retransmit, fast recovery.

---

Endnotes:

- Tahoe, Reno repairs one loss per RTT
- For multiple losses other ACK heuristics are used
  - NewReno, e.g.
  - Selective ACK (SACK) another approach [range of bytes received]

## Explicit Congestion Notification

Routers explicitly notifying hosts about congestion.

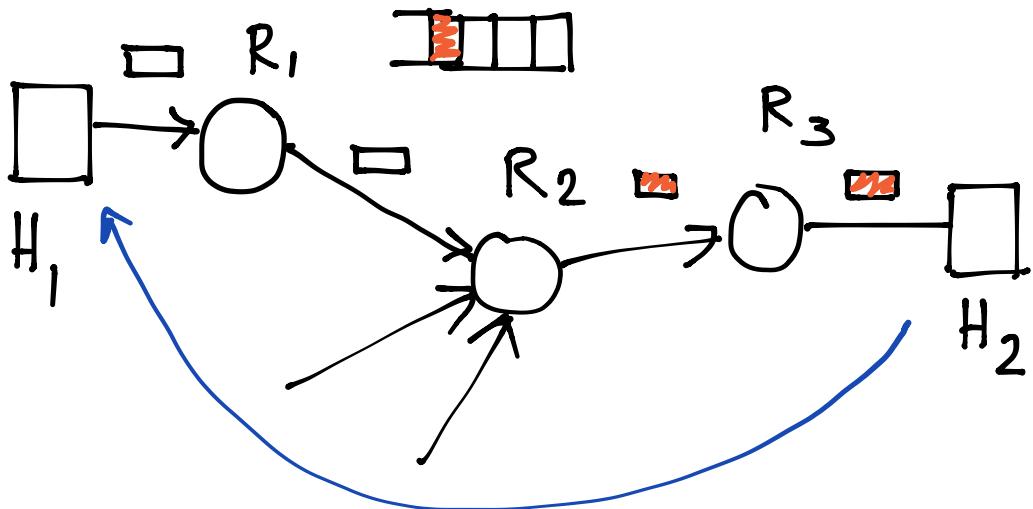
- congestion control - driving the network towards congestion and then recovering from it.
- Congestion avoidance - shape the traffic in such a way that congestion doesn't occur.

---

### Strategy of ECN

- Router detects congestion through its queue length
- when congested it marks every packet through a flag in the IP header (ECN)

EX .



receiving host sends the congestion notification when it sends ACK

- $H_1$  gets this and slows down.

### Advantages

- congestion detected early, no retransmission needed
- no extra packet needs to be sent

### disadvantage

- routers need to be upgraded