



TCP Reno

window-based Four Stage:

Loss-based ① slow start

② congestion avoidance

③ fast retransmission

④ fast recovery

ideal TCP self-clocking

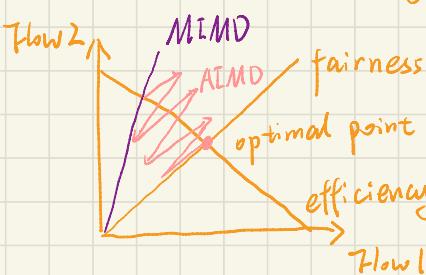
For each ack, cwnd + 1

↓

double cwnd in each RTT

cwnd = packets in flight

= bandwidth-delay product



Why MIMD cannot reach optimal point?

Cons

① Bursty traffic when ACKs come quickly

② Under-utilization in lossy networks
(why for Reno but not for others)

③ high delay from deep buffers ← bufferbloat window-based?

④ synchronize flows ← caused by AIMD?

Once congestion

cwnd = cwnd / 2

then in each RTT, cwnd + 1

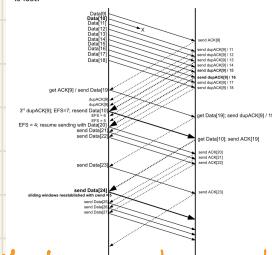
Resend Data[n]

when receiving 3 dup ACK[n-1]

lost Data[n] but received n+1, n+2

Received	Response
Data[1]	ACK[1]
Data[3]	ACK[1]
Data[4]	ACK[1]
Data[2]	ACK[4]
Data[6]	ACK[4]
Data[5]	ACK[6]

Here is a diagram illustrating Fast Recovery for cwnd=10. Data[10] is lost.



why
Reno?

Active Queue Management

goal — reduce excessive queuing delays \Rightarrow buffer bloat

three drop modes:

① Droptail (FIFO)



TCP Reno connections are happiest when the queue capacity at the bottleneck router exceeds the bandwidth \times delay transit capacity.

② Drophead

\rightarrow to inform sender congestion
as early as possible

why happiest? If not, it cannot fully utilize the bw?

③ Random

RED

(Random Early Detection)

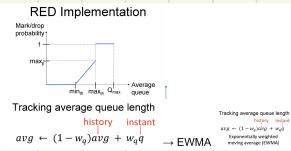
The actual RED algorithm does two things. First, the base drop probability p_{base} rises steadily from a minimum queue threshold q_{min} to a maximum queue threshold q_{max} (these might be 40% and 80% respectively of the absolute queue capacity); at the maximum threshold, the drop probability is still quite small. The base probability p_{base} increases linearly in this range according to the following formula, where p_{max} is the maximum RED-drop probability; the value for p_{max} proposed in [F93] was 0.02.

$$p_{base} = p_{max} \times (\text{avg_queuesize} - q_{min}) / (q_{max} - q_{min})$$

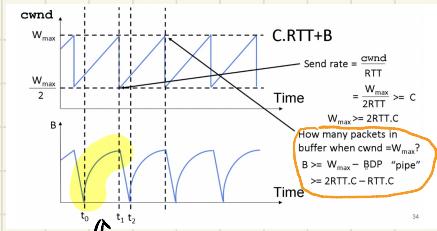
Second, as time passes after a RED drop, the actual drop probability p_{actual} begins to rise, according to the next formula:

$$p_{actual} = p_{base} / (1 - count \times p_{base})$$

Here, count is the number of packets sent since the last RED drop. With count=0 we have $p_{actual} = p_{base}$, but p_{actual} rises from then on with a RED drop guaranteed within the next $1/p_{base}$ packets. This provides a



$$\text{Send rate} = \frac{\text{cwnd}}{\text{RTT}}$$



- Queue increases linearly.
 - but queuing delay $\uparrow \rightarrow$ RTT \uparrow
 - Increase the same amount, it takes more time

CoDel

"Sojourn Time"

measure queue utilization with the time the packets spend in the queue

Implementing CoDel (2012)

- Buffer overflow \Rightarrow drop packet (as per normal)
- Track local minimum queue delay (time)
- Start dropping when delay > target
- Next drop time is decreased in inverse proportion to the square root of the number of drops since we started dropping
- When delay < target, stop dropping



Queue size is not a good measure?

measure the minimum value of queue utilization in a interval

$$\text{drop time} = \langle \frac{1}{\sqrt{2}}, \frac{1}{\sqrt{3}}, \frac{1}{\sqrt{4}} \rangle \times \text{Interval}$$

a little larger than $\text{RTT}_{\text{no load}}$

How can we measure fairness?

$$\text{Jain's index: } f = \frac{\left(\sum_{i=1}^N x_i\right)^2}{\left(N \sum_{i=1}^N x_i^2\right)}$$

Max-min fairness: Zero-sum game

Allocation that maximizes the minimum throughput among all flows

fairness { per-flow state
long flow vs short flow
Bursty flow
flow last time

TCP Vegas

delay-based

$$\text{expected throughput} = \frac{\text{cwnd}}{\text{RTT}}$$



= sending rate?

[Westwood+](#) below. TCP Vegas estimates $\text{RTT}_{\text{noLoad}}$ by the minimum RTT (RTT_{min}) encountered during the connection. The "ideal" cwnd that just saturates the bottleneck link is $\text{BWE} \times \text{RTT}_{\text{noLoad}}$. Note that BWE will be much more volatile than RTT_{min} ; the latter will typically reach its final value early in the connection, while BWE will fluctuate up and down with congestion (which will also act on RTT, but by increasing it).

As in 8.3.2 - RTT Calculations, any TCP sender can estimate queue utilization as

$$\text{queue_use} = \text{cwnd} - \text{BWE} \times \text{RTT}_{\text{noLoad}} = \text{cwnd} \times (1 - \text{RTT}_{\text{noLoad}}/\text{RTT}_{\text{actual}})$$

TCP Vegas then adjusts cwnd regularly to maintain the following:

$$\alpha \leq \text{queue_use} \leq \beta$$

which is the same as

$$\text{BWE} \times \text{RTT}_{\text{noLoad}} + \alpha \leq \text{cwnd} \leq \text{BWE} \times \text{RTT}_{\text{noLoad}} + \beta$$

Typically $\alpha = 2-3$ packets and $\beta = 4-6$ packets. We increment cwnd by 1 if cwnd falls below the lower limit (eg if BWE has increased). Similarly, we decrement cwnd by 1 if BWE drops and cwnd exceeds $\text{BWE} \times \text{RTT}_{\text{noLoad}} + \beta$.

limit cwnd close to the number of packets currently in flight
(measure and control the amount of extra data in transit)

queue size

BBR

main idea — manipulate sending rate
rate-based

BWE (Bandwidth Estimate)



maximum rate recorded over the past ten RTTs



How to measure?

PROBE-BW

STARTUP

$$\hookrightarrow \text{pacing-gain} = 2.89 (3/\log_2)$$

\hookrightarrow Until an additional RTT yields no improvement in BWE

DRAIN

$$\hookrightarrow \text{pacing-rate} = 1/2.89$$

Behaviour ① Set base sending rate to BWE (Update every 8 RTTs) Till in-flight pkts



purpose ② Set cloud target to $2 \times \text{BWE} \times \text{RTT}_{\min}$

why?



To let bottleneck queue = transit capacity

WHY

$\text{BWE} \uparrow \leftarrow \text{if } \text{BWB} \downarrow, 0.75 \rightarrow 1.25 \rightarrow 1.0 ?$



$$\text{new rate} = 1.25 \times \text{BWE}$$

(last 1 RTT)



After STARTUP,

PROBE-BW and
PROBE-RTT, which
comes first? (negligible "wasted" throughput)



IF \rightarrow RTT is also 1.25 times larger



$$\text{pacing-gain} = 0.75$$

$$\text{new rate} = 0.75 \times \text{BWE}$$

(last 1 RTT)



$$\text{pacing-gain} = 1.0$$

$$(\text{last 6 RTT})$$



BBR is TCP friendly?

200ms out of 10 sec)

\uparrow
check period

In any RTT, we can either measure

bottleneck bandwidth or RTT_{\min}



To test BWE, packets in flight \rightarrow use ACK rate to test?
 \rightarrow "transit capacity" here is only this flow?

To test RTT_{\min} , packets in flight $<$ transit capacity

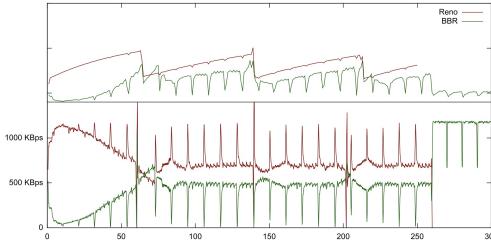
Here is a concrete example of BWE increase. To simplify the analysis, we will assume TCP BBR's FlightSize is $BWE \times RTT_{min}$, dropping the factor of 2. Suppose a TCP BBR connection and a TCP Reno connection share a bottleneck link with a bandwidth of 2 packets/ms. The RTT_{min} ($= RTT_{noLoad}$) of each connection is 80 ms, making the transit capacity 160 packets. Finally, suppose that each connection has 80 packets in flight, exactly filling the transit capacity but with no queue utilization (so $RTT_{min} = RTT_{actual}$). Over the course of the eight-RTT `pacing_gain` cycle, the Reno connection's `cwnd` rises by 8, to 88 packets. This means the total queue utilization is now 8 packets, divided on average between BBR and Reno in the proportion 80 to 88.

Now the BBR cycle with `pacing_gain`=1.25 arrives; for the next RTT, the BBR connection has $80 \times 1.25 = 100$ packets in flight. The total number of packets in flight is now 188. The RTT climbs to $188/2 = 94$ ms, and the next BBR BWE measurement is 100 packets in 94 ms, or 1.064 packets/ms (the precise value may depend on exactly when the measurement is recorded). For the following RTT, `pacing_gain` drops to 0.75, but the higher BWE persists. For the rest of the `pacing_gain` cycle, TCP BBR calculates a base rate corresponding to 1.064 × 80 = 85 packets in flight per RTT, which is close to the TCP Reno `cwnd`. See also exercise 14.0.

competitive rate. Suppose, for example, that in the BBR-vs-Reno scenario above, Reno has gobbled up a total of 240 spots in the bottleneck queue, thus increasing the RTT for both connections to $(240+80)/2 = 160$. During a PROBE_RTT cycle, TCP BBR will drop its link utilization essentially to zero, but TCP Reno will still have 240 packets in transit, so TCP BBR will measure RTT_{min} as $240/2 = 120$ ms. After the PROBE_RTT phase is over, TCP BBR will increase its sending rate by 50% over what it had been when RTT_{min} was 80.

The lower part of the diagram shows each connection's share of the 10 Mbps (1.25 MBps) bottleneck bandwidth. The upper part shows the number of packets "in flight" (for TCP Reno, outside of Fast Recovery, that is of course `cwnd`). The Reno sawtooth pattern is clearly visible.

A dominant feature of the graph is the spikes every 10 seconds (down for BBR, correspondingly up for Reno) caused by TCP BBR's periodic PROBE_RTT mode.



Known Issues: BBR

- Not friendly with Reno/CUBIC
- Impossible to estimate RTT_{min} ⇒ $cwnd = 2 \times BDP \Rightarrow$ High delay
- Loss rate can be quite high for small buffers
- PROBE_RTT causes significant throughput reduction