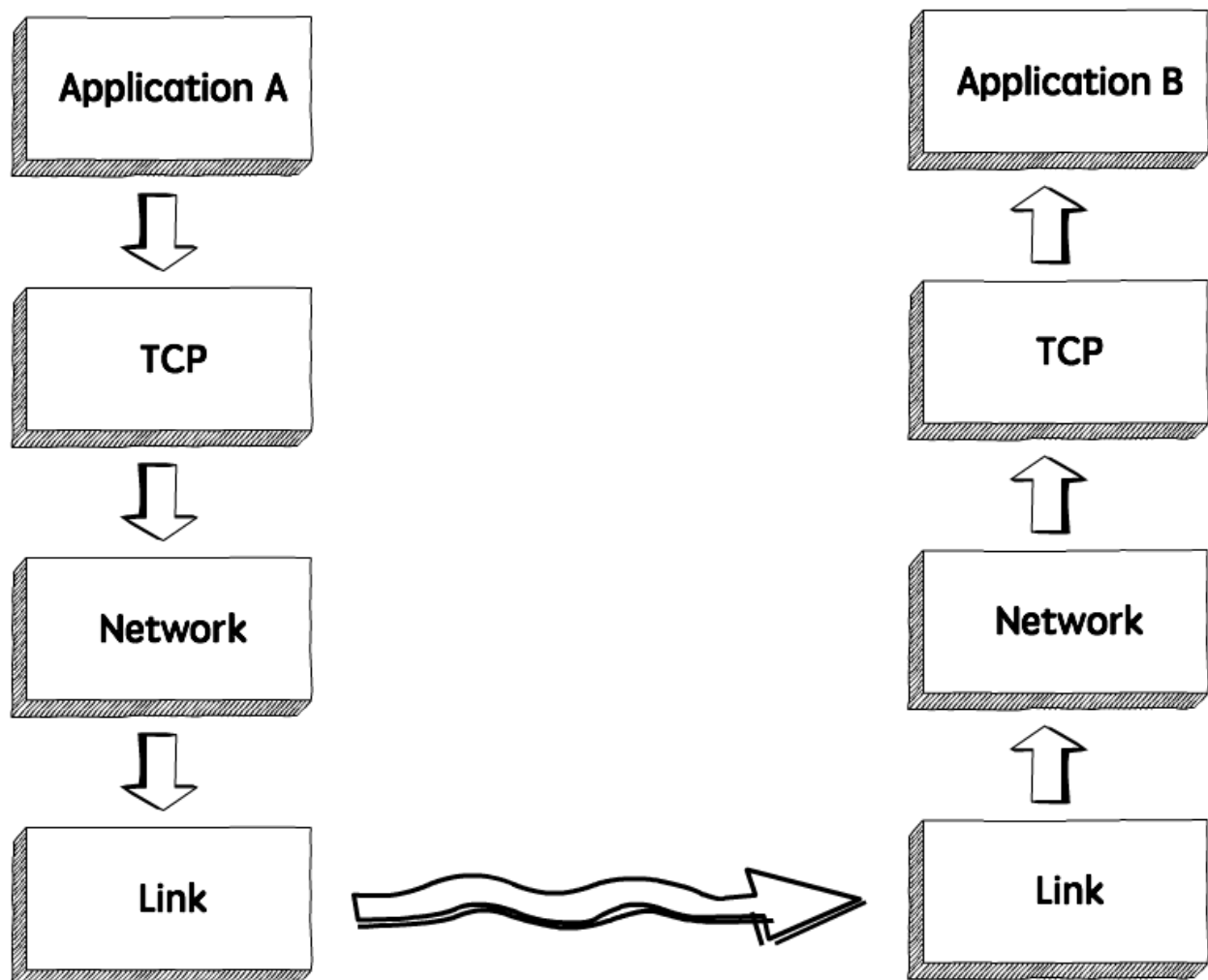


TCP Flow Control

🌐 briantorti.com/tcp-flow-control

June 30, 2017



TCP is the protocol that guarantees we can have a reliable communication channel over an unreliable network. When we send data from a node to another, packets can be lost, they can arrive out of order, the network can be congested or the receiver node can be overloaded. When we are writing an application, though, we usually don't need to deal with this complexity, we just write some data to a socket and **TCP** makes sure the packets are delivered correctly to the receiver node. Another important service that **TCP** provides is what is called *Flow Control*. Let's talk about what that means and how **TCP** does its magic.

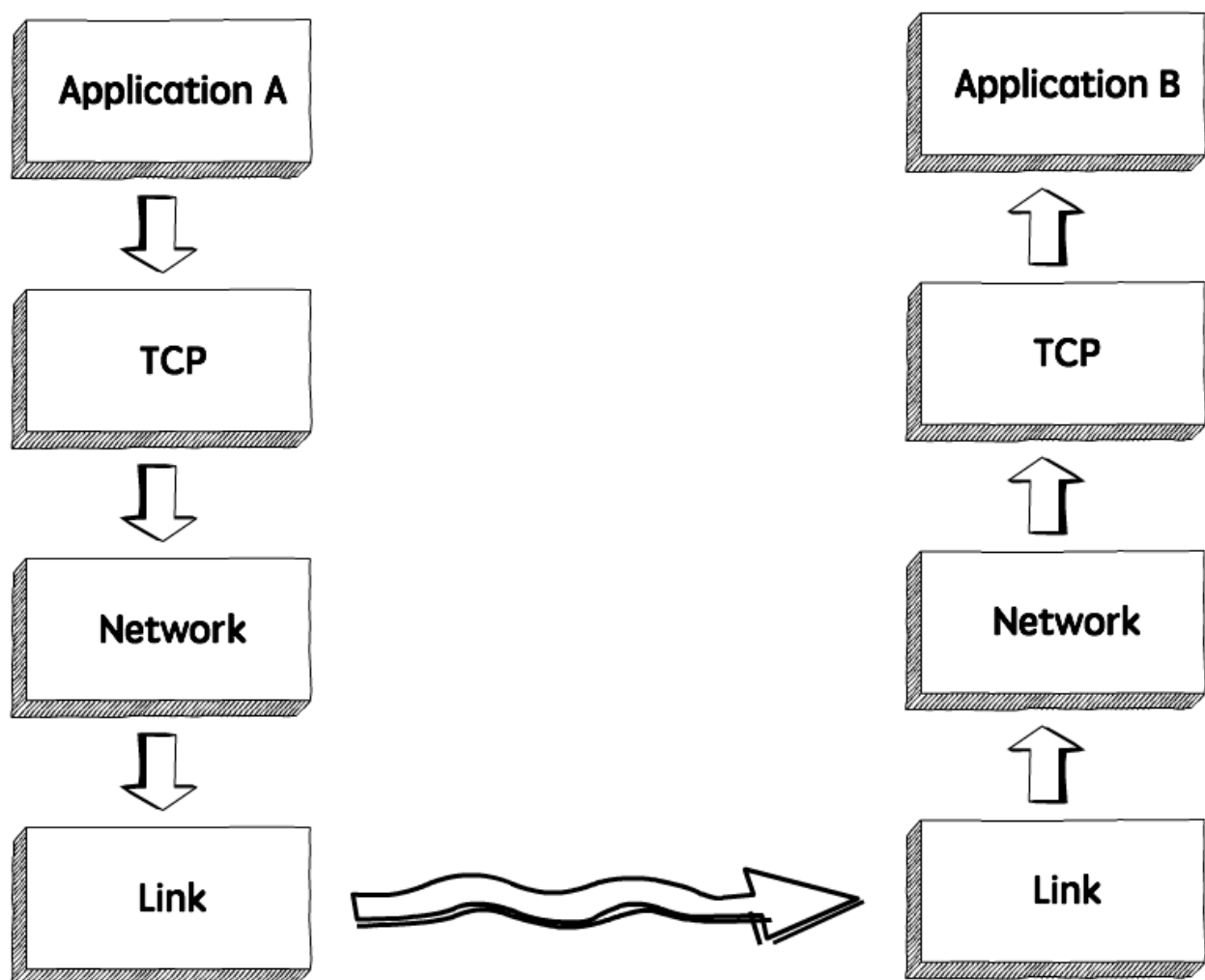
What is Flow Control (and what it's not)

Flow Control basically means that **TCP** will ensure that a sender is not overwhelming a receiver by sending packets faster than it can consume. It's pretty similar to what's normally called *Back pressure* in the Distributed Systems literature. The idea is that a node receiving data will send some kind of feedback to the node sending the data to let it know about its current condition.

It's important to understand that this is **not** the same as *Congestion Control*. Although there's some overlap between the mechanisms **TCP** uses to provide both services, they are distinct features. Congestion control is about preventing a node from overwhelming the network (i.e. the links between two nodes), while Flow Control is about the end-node.

How it works

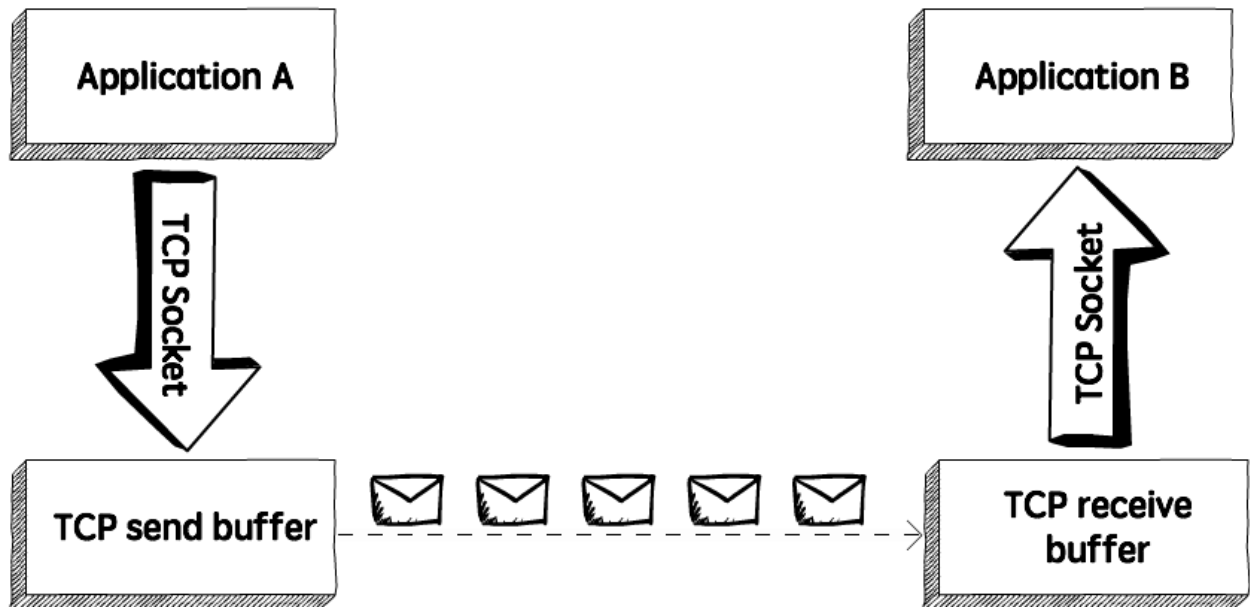
When we need to send data over a network, this is normally what happens.



The sender application writes data to a socket, the transport layer (in our case, **TCP**) will wrap this data in a segment and hand it to the network layer (e.g. **IP**), that will somehow route this packet to the receiving node.

On the other side of this communication, the network layer will deliver this piece of data to **TCP**, that will make it available to the receiver application as an exact copy of the data sent, meaning it will not deliver packets out of order, and will wait for a retransmission in case it notices a gap in the byte stream.

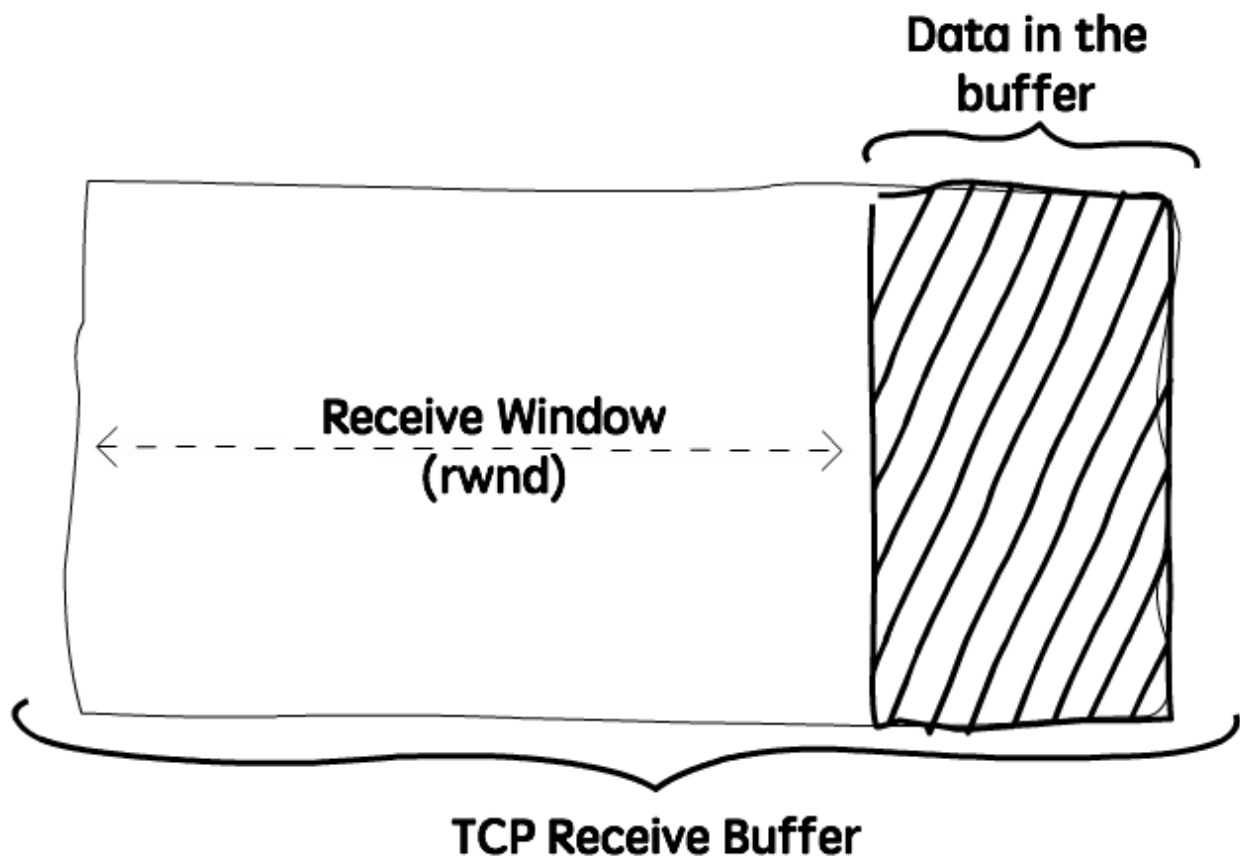
If we zoom in, we will see something like this.



TCP stores the data it needs to send in the *send buffer*, and the data it receives in the *receive buffer*. When the application is ready, it will then read data from the receive buffer.

Flow Control is all about making sure we don't send more packets when the receive buffer is already full, as the receiver wouldn't be able to handle them and would need to drop these packets.

To control the amount of data that **TCP** can send, the receiver will advertise its *Receive Window (rwnd)*, that is, the spare room in the receive buffer.



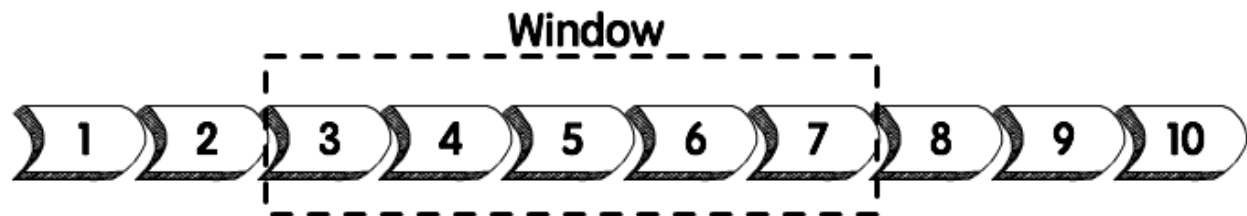
Every time **TCP** receives a packet, it needs to send an **ack** message to the sender, acknowledging it received that packet correctly, and with this **ack** message it sends the value of the current receive window, so the sender knows if it can keep sending data.

The sliding window

TCP uses a sliding window protocol to control the number of bytes in flight it can have. In other words, the number of bytes that were sent but not yet **acked**.

Let's say we want to send a 150000 bytes file from node A to node B. **TCP** could break this file down into 100 packets, 1500 bytes each. Now let's say that when the connection between node A and B is established, node B advertises a receive window of 45000 bytes, because it really wants to help us with our math here.

Seeing that, **TCP** knows it can send the first 30 packets ($1500 * 30 = 45000$) before it receives an acknowledgment. If it gets an **ack** message for the first 10 packets (meaning we now have only 20 packets in flight), and the receive window present in these **ack** messages is still 45000, it can send the next 10 packets, bringing the number of packets in flight back to 30, that is the limit defined by the receive window. In other words, at any given point in time it can have 30 packets in flight, that were sent but not yet **acked**.



Example of a sliding window. As soon as packet 3 is acked, we can slide the window to the right and send the packet 8.

Now, if for some reason the application reading these packets in node B slows down, **TCP** will still **ack** the packets that were correctly received, but as these packets need to be stored in the receive buffer until the application decides to read them, the receive window will be smaller, so even if **TCP** receives the acknowledgment for the next 10 packets (meaning there are currently 20 packets, or 30000 bytes, in flight), but the receive window value received in this **ack** is now 30000 (instead of 45000), it will not send more packets, as the number of bytes in flight is already equal to the latest receive window advertised.

The sender will always keep this invariant:

$$\text{LastByteSent} - \text{LastByteAcked} \leq \text{ReceiveWindowAdvertised}$$

Visualizing the Receive Window

Just to see this behavior in action, let's write a very simple application that reads data from a socket and watch how the receive window behaves when we make this application slower. We will use **Wireshark** to see these packets, **netcat** to send data to this application, and a **go** program to read data from the socket.

Here's the simple **go** program that reads and prints the data received:

```
package main

import (
    "bufio"
    "fmt"
    "net"
)

func main() {
    listener, _ := net.Listen("tcp", "localhost:3040")
    conn, _ := listener.Accept()

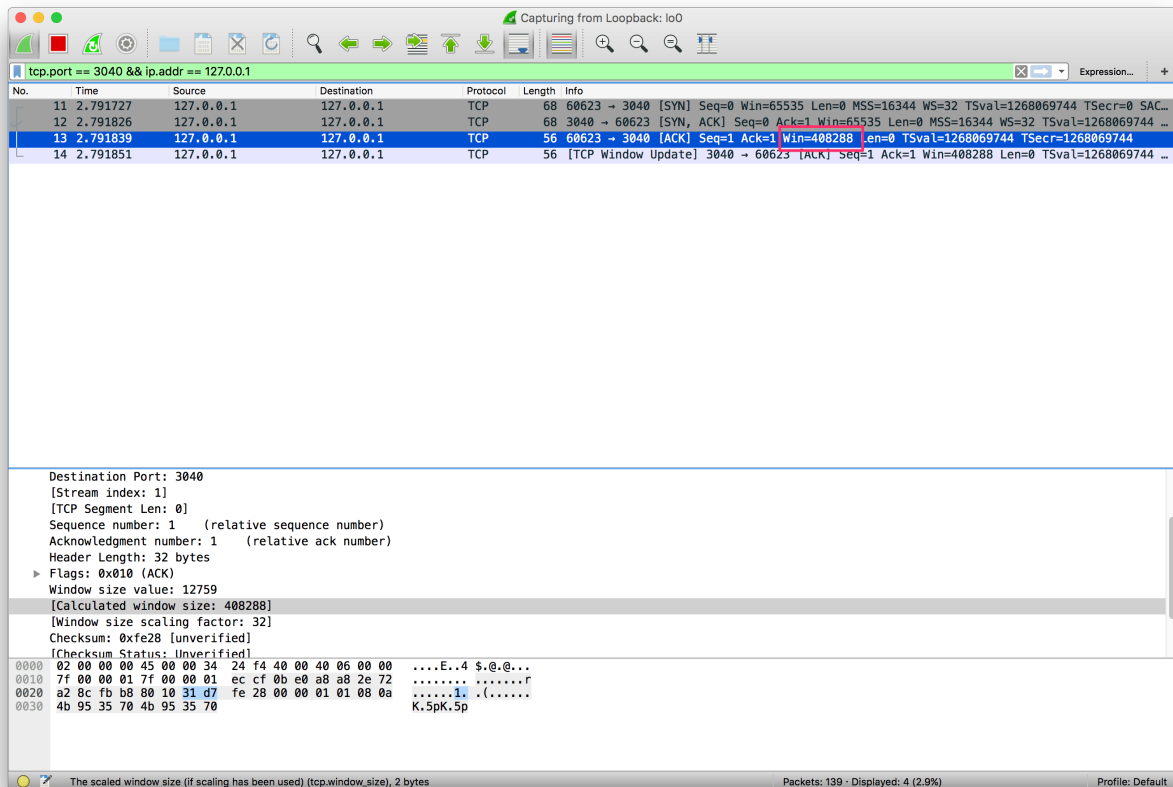
    for {
        message, _ := bufio.NewReader(conn).ReadBytes('\n')
        fmt.Println(string(message))
    }
}
```

This program will simply listen to connections on port **3040** and print the string received.

We can then use **netcat** to send data to this application:

```
$ nc localhost 3040
```

And we can see, using **Wireshark**, that the connection was established and a window size advertised:



Click on the image to enlarge it.

Now let's run this command to create a stream of data. It will simply add the string "foo" to a file, that we will use to send to this application:

```
$ while true; do echo "foo" > stream.txt; done
```

And now let's send this data to the application:

```
tail -f stream.txt | nc localhost 3040
```

Now if we check **Wireshark** we will see a lot of packets being sent, and the receive window being updated:

Capturing from Loopback: lo0

tcp.port == 3040 && ip.addr == 127.0.0.1

No.	Time	Source	Destination	Protocol	Length	Info
291	164.441892	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58053 Ack=1 Win=408288 Len=4 TSval=1268613879 TSecr=126...
291	164.441911	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58057 Win=350240 Len=0 TSval=1268613879 TSecr=12686138...
291	164.442135	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58057 Ack=1 Win=408288 Len=4 TSval=1268613880 TSecr=126...
291	164.442152	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58061 Win=350240 Len=0 TSval=1268613880 TSecr=12686138...
291	164.442313	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58061 Ack=1 Win=408288 Len=4 TSval=1268613880 TSecr=126...
291	164.442330	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58065 Win=350208 Len=0 TSval=1268613880 TSecr=12686138...
291	164.442412	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58065 Ack=1 Win=408288 Len=4 TSval=1268613880 TSecr=126...
291	164.442432	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58069 Win=350208 Len=0 TSval=1268613880 TSecr=12686138...
291	164.442620	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58069 Ack=1 Win=408288 Len=4 TSval=1268613880 TSecr=126...
291	164.442637	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58073 Win=350208 Len=0 TSval=1268613880 TSecr=12686138...
291	164.442858	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58073 Ack=1 Win=408288 Len=4 TSval=1268613880 TSecr=126...
291	164.442875	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58077 Win=350208 Len=0 TSval=1268613880 TSecr=12686138...
291	164.443068	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58077 Ack=1 Win=408288 Len=4 TSval=1268613880 TSecr=126...
291	164.443085	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58081 Win=350208 Len=0 TSval=1268613880 TSecr=12686138...
291	164.443259	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58081 Ack=1 Win=408288 Len=4 TSval=1268613881 TSecr=126...
291	164.443277	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58085 Win=350208 Len=0 TSval=1268613881 TSecr=12686138...
291	164.443403	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58085 Ack=1 Win=408288 Len=4 TSval=1268613881 TSecr=126...
291	164.443420	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58089 Win=350208 Len=0 TSval=1268613881 TSecr=12686138...
291	164.443443	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58089 Ack=1 Win=408288 Len=4 TSval=1268613881 TSecr=126...
291	164.443456	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58093 Win=350208 Len=0 TSval=1268613881 TSecr=12686138...
291	164.443674	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58093 Ack=1 Win=408288 Len=4 TSval=1268613881 TSecr=126...
291	164.443687	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=58097 Win=350176 Len=0 TSval=1268613881 TSecr=12686138...
291	164.443791	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=58097 Ack=1 Win=408288 Len=4 TSval=1268613881 TSecr=126...

Source Port: 3040
Destination Port: 61435
[Stream index: 6]
[TCP Segment Len: 0]
Sequence number: 1 (relative sequence number)
Acknowledgment number: 126193 (relative ack number)
Header Length: 32 bytes
Flags: 0x010 (ACK)
Window size value: 8815
[Calculated window size: 282080]
[Window size scaling factor: 32]
Checksum: 0xefe28 [unverified]
0000 02 00 00 00 45 00 00 34 42 51 40 00 00 06 00 00E..4 BQ@.@..
0010 7f 00 00 01 7f 00 00 01 0b e0 ef fb 08 61 4a ffAJ..
0020 a5 9c 53 7c 80 10 22 6f fe 28 00 00 01 01 08 0a ..S|..o .(.....
0030 4b 9d f4 ba 4b 9d f4 ba K...K...

The scaled window size (if scaling has been used) (tcp.window_size), 2 bytes

Packets: 63288 - Displayed: 62778 (99.2%)

Profile: Default

Capturing from Loopback: lo0

tcp.port == 3040 && ip.addr == 127.0.0.1

No.	Time	Source	Destination	Protocol	Length	Info
630	175.761551	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126149 Ack=1 Win=408288 Len=4 TSval=1268623493 TSecr=12...
630	175.761622	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126153 Win=282144 Len=0 TSval=1268623493 TSecr=1268623...
630	175.762723	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126153 Ack=1 Win=408288 Len=4 TSval=1268623494 TSecr=12...
630	175.762781	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126157 Win=282144 Len=0 TSval=1268623494 TSecr=1268623...
630	175.776996	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126157 Ack=1 Win=408288 Len=4 TSval=1268623507 TSecr=12...
630	175.777028	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126161 Win=282112 Len=0 TSval=1268623507 TSecr=1268623...
630	175.777126	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126161 Ack=1 Win=408288 Len=4 TSval=1268623507 TSecr=12...
630	175.777148	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126165 Win=282112 Len=0 TSval=1268623507 TSecr=1268623...
630	175.780480	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126165 Ack=1 Win=408288 Len=4 TSval=1268623510 TSecr=12...
630	175.780514	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126169 Win=282112 Len=0 TSval=1268623510 TSecr=1268623...
631	175.782776	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126169 Ack=1 Win=408288 Len=4 TSval=1268623512 TSecr=12...
631	175.782803	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126173 Win=282112 Len=0 TSval=1268623512 TSecr=1268623...
631	175.788081	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126173 Ack=1 Win=408288 Len=4 TSval=1268623517 TSecr=12...
631	175.788126	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126177 Win=282112 Len=0 TSval=1268623517 TSecr=1268623...
631	175.794693	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126177 Ack=1 Win=408288 Len=4 TSval=1268623523 TSecr=12...
631	175.794726	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126181 Win=282112 Len=0 TSval=1268623523 TSecr=1268623...
631	175.801599	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126181 Ack=1 Win=408288 Len=4 TSval=1268623528 TSecr=12...
631	175.801629	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126185 Win=282112 Len=0 TSval=1268623528 TSecr=1268623...
631	175.806760	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126185 Ack=1 Win=408288 Len=4 TSval=1268623533 TSecr=12...
631	175.806814	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126189 Win=282112 Len=0 TSval=1268623533 TSecr=1268623...
631	195.467968	127.0.0.1	127.0.0.1	TCP	60	61435 → 3040 [PSH, ACK] Seq=126189 Ack=1 Win=408288 Len=4 TSval=1268643002 TSecr=12...
631	195.468012	127.0.0.1	127.0.0.1	TCP	56	3040 → 61435 [ACK] Seq=1 Ack=126193 Win=282080 Len=0 TSval=1268643002 TSecr=1268643...

Source Port: 3040
Destination Port: 61435
[Stream index: 6]
[TCP Segment Len: 0]
Sequence number: 1 (relative sequence number)
Acknowledgment number: 126193 (relative ack number)
Header Length: 32 bytes
Flags: 0x010 (ACK)
Window size value: 8815
[Calculated window size: 282080]
[Window size scaling factor: 32]
Checksum: 0xefe28 [unverified]
0000 02 00 00 00 45 00 00 34 42 51 40 00 00 06 00 00E..4 BQ@.@..
0010 7f 00 00 01 7f 00 00 01 0b e0 ef fb 08 61 4a ffAJ..
0020 a5 9c 53 7c 80 10 22 6f fe 28 00 00 01 01 08 0a ..S|..o .(.....
0030 4b 9d f4 ba 4b 9d f4 ba K...K...

The scaled window size (if scaling has been used) (tcp.window_size), 2 bytes

Packets: 63178 - Displayed: 62778 (99.4%)

Profile: Default

The application is still fast enough to keep up with the work, though. So let's make it a bit slower to see what happens:


```

package main

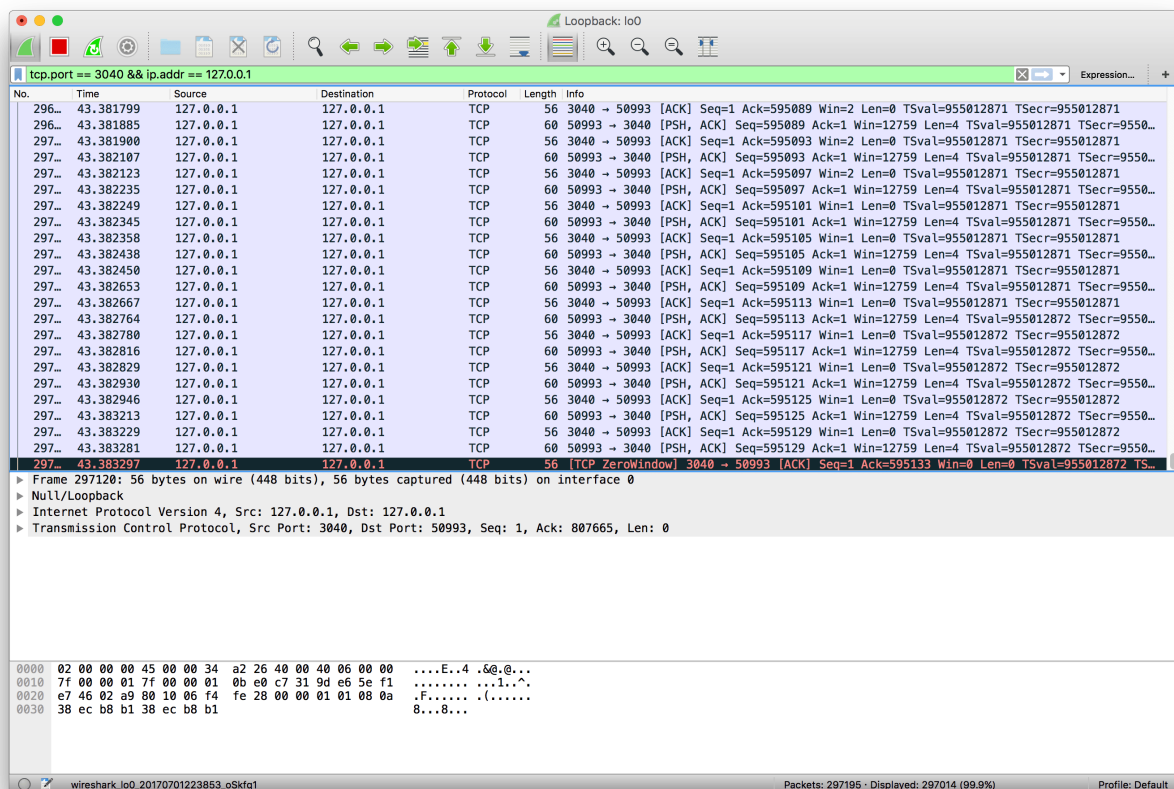
import (
    "bufio"
    "fmt"
    "net"
    "time"
)

func main() {
    listener, _ := net.Listen("tcp", "localhost:3040")
    conn, _ := listener.Accept()

    for {
        message, _ := bufio.NewReader(conn).ReadBytes('\n')
        fmt.Println(string(message))
        +
        time.Sleep(1 * time.Second)
    }
}

```

Now we are sleeping for 1 second before we read data from the receive buffer. If we run **netcat** again and observe **Wireshark**, it doesn't take long until the receive buffer is full and **TCP** starts advertising a 0 window size:

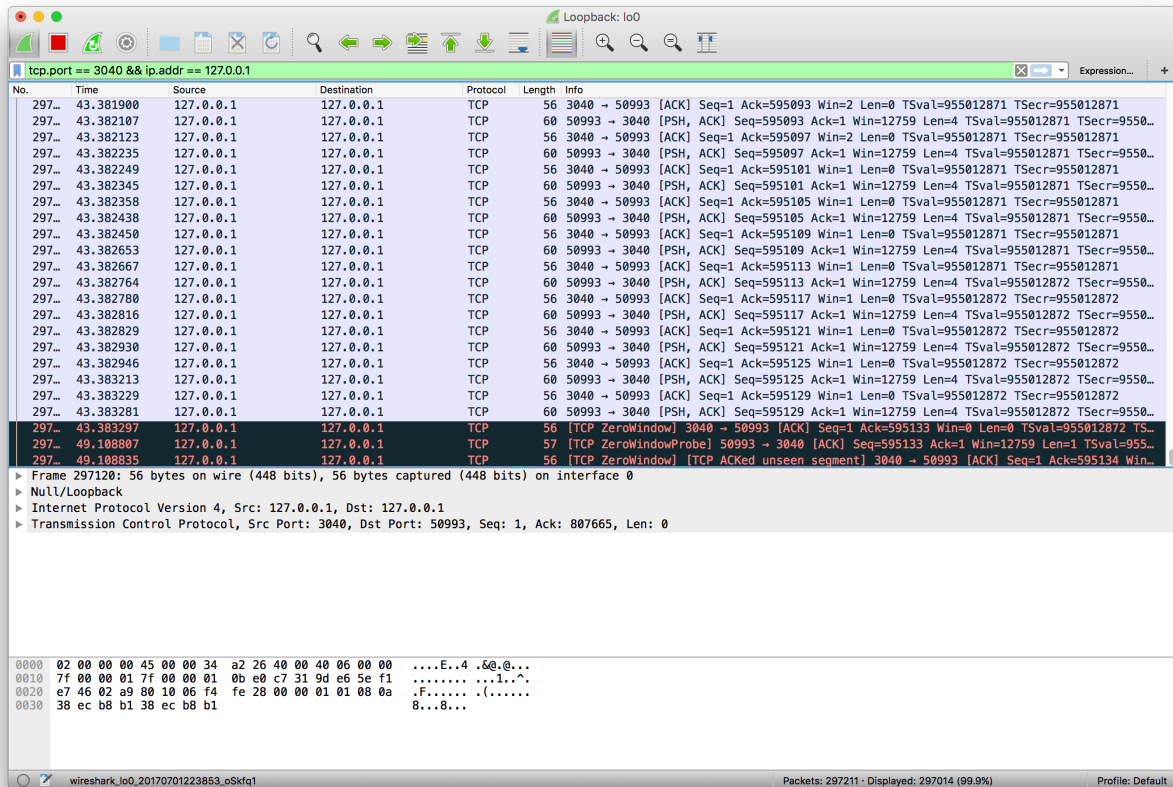


At this moment **TCP** will stop transmitting data, as the receiver's buffer is full.

The persist timer

There's still one problem, though. After the receiver advertises a zero window, if it doesn't send any other **ack** message to the sender (or if the **ack** is lost), it will never know when it can start sending data again. We will have a deadlock situation, where the receiver is waiting for more data, and the sender is waiting for a message saying it can start sending data again.

To solve this problem, when **TCP** receives a zero-window message it starts the *persist timer*, that will periodically send a small packet to the receiver (usually called **WindowProbe**), so it has a chance to advertise a nonzero window size.



When there's some spare space in the receiver's buffer again it can advertise a non-zero window size and the transmission can continue.

Recap

- **TCP**'s flow control is a mechanism to ensure the sender is not overwhelming the receiver with more data than it can handle;
- With every **ack** message the receiver advertises its current receive window;
- The receive window is the spare space in the receive buffer, that is, $rwnd = ReceiveBuffer - (LastByteReceived - LastByteReadByApplication)$;
- **TCP** will use a sliding window protocol to make sure it never has more bytes in flight than the window advertised by the receiver;
- When the window size is 0, **TCP** will stop transmitting data and will start the persist timer;

- It will then periodically send a small **WindowProbe** message to the receiver to check if it can start receiving data again;
- When it receives a non-zero window size, it resumes the transmission.

If you want to learn more about TCP (and *a lot* more), the book [Computer Networking: A Top-Down Approach](#) is a great resource.

[Get PDF](#)

Interested in learning Kubernetes?

I just published a new book called [Kubernetes in Practice](#), you can use the discount code **blog** to get 10% off.

Get fresh articles in your inbox

If you liked this article, you might want to subscribe. If you don't like what you get, unsubscribe with one click.