# TCP

The *Transmission Control Protocol* (TCP) is the transport layer
for a reliable, connection-oriented protocol

# TCP

The *Transmission Control Protocol* (TCP) is the transport layer for a reliable, connection-oriented protocol

Often called "TCP/IP"

# TCP

The *Transmission Control Protocol* (TCP) is the transport layer for a reliable, connection-oriented protocol

Often called "TCP/IP"

It is *hugely* more complicated than UDP as it must create a reliable transport from the unreliable IP it runs on

# TCP

The *Transmission Control Protocol* (TCP) is the transport layer for a reliable, connection-oriented protocol

Often called "TCP/IP"

It is *hugely* more complicated than UDP as it must create a reliable transport from the unreliable IP it runs on

There is a lot of complication to deal with the error cases, such as packet loss and packet duplication

# TCP

The *Transmission Control Protocol* (TCP) is the transport layer for a reliable, connection-oriented protocol

Often called "TCP/IP"

It is *hugely* more complicated than UDP as it must create a reliable transport from the unreliable IP it runs on

There is a lot of complication to deal with the error cases, such as packet loss and packet duplication

There is overhead in setting up (and taking down) the connection to manage these mechanisms

# TCP

The *Transmission Control Protocol* (TCP) is the transport layer for a reliable, connection-oriented protocol

Often called "TCP/IP"

It is *hugely* more complicated than UDP as it must create a reliable transport from the unreliable IP it runs on

There is a lot of complication to deal with the error cases, such as packet loss and packet duplication

There is overhead in setting up (and taking down) the connection to manage these mechanisms

And more to complexity improve performance and flow control

# TCP

The *Transmission Control Protocol* (TCP) is the transport layer for a reliable, connection-oriented protocol

Often called "TCP/IP"

It is *hugely* more complicated than UDP as it must create a reliable transport from the unreliable IP it runs on

There is a lot of complication to deal with the error cases, such as packet loss and packet duplication

There is overhead in setting up (and taking down) the connection to manage these mechanisms

And more to complexity improve performance and flow control

A lot of state about each connection needs to be stored by the OS

# TCP

The basis of the reliability is the use of acknowledgement (ACK) packets for every packet sent

# TCP

The basis of the reliability is the use of acknowledgement (ACK) packets for every packet sent

If host A sends host B a packet, B must send an ACK packet back to A to inform it of the safe arrival of the packet

# TCP

The basis of the reliability is the use of acknowledgement (ACK) packets for every packet sent

If host A sends host B a packet, B must send an ACK packet back to A to inform it of the safe arrival of the packet

If A does not get an ACK, it resends the packet

# TCP

The basis of the reliability is the use of acknowledgement (ACK) packets for every packet sent

If host A sends host B a packet, B must send an ACK packet back to A to inform it of the safe arrival of the packet

If A does not get an ACK, it resends the packet

But ACKs on their own do not solve all the problem

This is due to the *Two Armies Problem*: suppose two armies A and B wish to coordinate an attack on C

This is due to the *Two Armies Problem*: suppose two armies A and B wish to coordinate an attack on C

A sends a message to B: "attack at dawn"

# TCP

This is due to the *Two Armies Problem*: suppose two armies A and B wish to coordinate an attack on C

A sends a message to B: "attack at dawn"

How does A know that B got the message? A cannot safely attack until it knows B is ready

# TCP

This is due to the *Two Armies Problem*: suppose two armies A and B wish to coordinate an attack on C

A sends a message to B: "attack at dawn"

How does A know that B got the message? A cannot safely attack until it knows B is ready

So B sends an acknowledgement to A: "OK"

# TCP

This is due to the *Two Armies Problem*: suppose two armies A and B wish to coordinate an attack on C

A sends a message to B: "attack at dawn"

How does A know that B got the message? A cannot safely attack until it knows B is ready

So B sends an acknowledgement to A: "OK"

But the ACK might be intercepted and A might not get the ACK

B can't attack until it knows A got the ACK

# TCP

B can't attack until it knows A got the ACK

So A should send an ACK for the ACK back to B

# TCP

B can't attack until it knows A got the ACK

So A should send an ACK for the ACK back to B

But this might not get through…

# TCP

B can't attack until it knows A got the ACK

So A should send an ACK for the ACK back to B

But this might not get through. . .

For full reliability it looks like we need an infinite regress!

# TCP

TCP avoids the Two Armies Problem by using timeouts and packet retransmissions

TCP avoids the Two Armies Problem by using timeouts and packet retransmissions

For every packet:

# TCP

TCP avoids the Two Armies Problem by using timeouts and packet retransmissions

For every packet:

A starts a *retransmission timer* when it sends to B

# TCP

TCP avoids the Two Armies Problem by using timeouts and packet retransmissions

For every packet:

A starts a *retransmission timer* when it sends to B

If the timer runs out before it gets an ACK, it resends the packet and restarts the timer

# TCP

TCP avoids the Two Armies Problem by using timeouts and packet retransmissions

For every packet:

A starts a *retransmission timer* when it sends to B

If the timer runs out before it gets an ACK, it resends the packet and restarts the timer

Repeat until A gets an ACK (or A gives up)

# TCP

Problems to solve include:

# TCP

Problems to solve include:

- how long to wait before a resend? This might be a slow but otherwise reliable link and resending will just clog the system with extra duplicate packets

# TCP

Problems to solve include:

- how long to wait before a resend? This might be a slow but otherwise reliable link and resending will just clog the system with extra duplicate packets

- how many times to resend before giving up? It might be the destination has gone away entirely (perhaps disconnected or crashed)

# TCP

Problems to solve include:

- how long to wait before a resend? This might be a slow but otherwise reliable link and resending will just clog the system with extra duplicate packets
- how many times to resend before giving up? It might be the destination has gone away entirely (perhaps disconnected or crashed)
- how long B should wait before sending the ACK? You can *piggyback* an ACK on an ordinary data packet, so it may be better for B to wait until some data is ready to be returned rather than sending an otherwise empty ACK. This saves on packets sent

# TCP

Problems to solve include:

- how long to wait before a resend? This might be a slow but otherwise reliable link and resending will just clog the system with extra duplicate packets

- how many times to resend before giving up? It might be the destination has gone away entirely (perhaps disconnected or crashed)

- how long B should wait before sending the ACK? You can *piggyback* an ACK on an ordinary data packet, so it may be better for B to wait until some data is ready to be returned rather than sending an otherwise empty ACK. This saves on packets sent

- IP datagrams can arrive out of order, so we need some way to recognise which ACK goes with which packet

Other problems TCP also needs to address include:

Other problems TCP also needs to address include:

- how to maintain order in the data? IP datagrams can arrive out of order, so we need some way of reassembling the original data stream in the correct order

# TCP

Other problems TCP also needs to address include:

- how to maintain order in the data? IP datagrams can arrive out of order, so we need some way of reassembling the original data stream in the correct order
- how to manage duplicates? Resends can produce duplicate packets (if the original was not actually lost) so we need some way to recognise and discard extra copies

# TCP

Other problems TCP also needs to address include:

- how to maintain order in the data? IP datagrams can arrive out of order, so we need some way of reassembling the original data stream in the correct order

- how to manage duplicates? Resends can produce duplicate packets (if the original was not actually lost) so we need some way to recognise and discard extra copies

- Flow control: how to increase the rate of sending packets when things are going well, and decrease the rate when they are not

# TCP

TCP packets are often called *segments*

TCP packets are often called *segments*

(Reminder: "segment", "packet", "datagram", "frame" all mean pretty much the same thing, just in different layers)

# TCP

TCP packets are often called *segments*

(Reminder: "segment", "packet", "datagram", "frame" all mean pretty much the same thing, just in different layers)

A TCP header is complicated as it must address many complex issues

# TCP



TCP header

# TCP

- Ports: identical to UDP (on purpose: actually UDP copied TCP)

# TCP

- Ports: identical to UDP (on purpose: actually UDP copied TCP)
- Two 32 bit values: *sequence* and *acknowledgement*

These numbers are the heart of TCP's reliability

Sequence numbers

These numbers are the heart of TCP's reliability

Every **byte** in a TCP connection is numbered

These numbers are the heart of TCP's reliability

Every **byte** in a TCP connection is numbered

The 32 bit sequence number starts at some random value and increases by 1 for each byte sent

These numbers are the heart of TCP's reliability

Every **byte** in a TCP connection is numbered

The 32 bit sequence number starts at some random value and increases by 1 for each byte sent

So if a segment contains 10 bytes of data, the sequence number on the next segment sent will be 10 greater

The sequence number in the header is the number of the first byte of data in the segment

The sequence number in the header is the number of the first byte of data in the segment

The destination acknowledges those bytes it has received by filling in the ACK field with the appropriate byte number and setting the ACK flag

The reverse connection from destination to source has its own sequence number as TCP is fully duplex

# TCP
## Sequence numbers

The reverse connection from destination to source has its own sequence number as TCP is fully duplex

Everything we say here is true for data travelling in the reverse direction: the reverse traffic has its own independent sequence numbers and flow control

Note that a destination might not immediately get the whole segment that was sent due to fragmentation in the IP layer

Note that a destination might not immediately get the whole segment that was sent due to fragmentation in the IP layer

IP must wait for all the fragments and reconstruct the segment before it can pass it on to TCP and then TCP can send the ACK

Note that a destination might not immediately get the whole segment that was sent due to fragmentation in the IP layer

IP must wait for all the fragments and reconstruct the segment before it can pass it on to TCP and then TCP can send the ACK

And this can play havoc with TCP's timers

Note that a destination might not immediately get the whole segment that was sent due to fragmentation in the IP layer

IP must wait for all the fragments and reconstruct the segment before it can pass it on to TCP and then TCP can send the ACK

And this can play havoc with TCP's timers

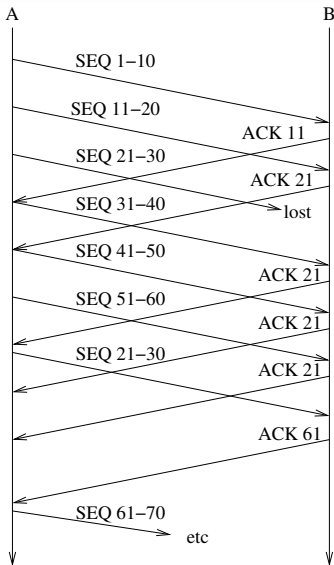Another reason to avoid fragmentation

# TCP
## Sequence numbers

The returning ACK field contains the sequence number of the next byte the destination expects to receive, e.g., if the sequence number is 20001 and 14 bytes are received it returns 20015 in the ACK field

The returning ACK field contains the sequence number of the next byte the destination expects to receive, e.g., if the sequence number is 20001 and 14 bytes are received it returns 20015 in the ACK field

ACKs can be *piggybacked* on normal returning data packets, they don't need to be separate packets

# TCP
## Sequence numbers

The returning ACK field contains the sequence number of the next byte the destination expects to receive, e.g., if the sequence number is 20001 and 14 bytes are received it returns 20015 in the ACK field

ACKs can be *piggybacked* on normal returning data packets, they don't need to be separate packets

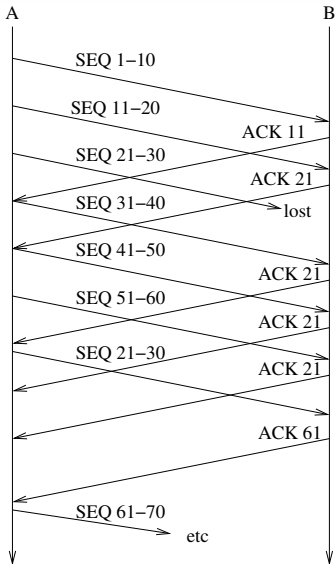This helps reduce the amount of network traffic

# TCP



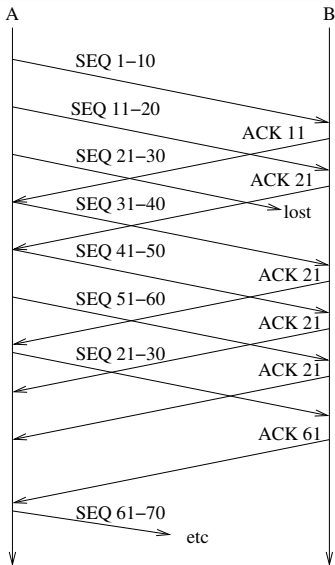A is sending 10 byte segments to B, and B is ACKing them;

ACKing lost segments

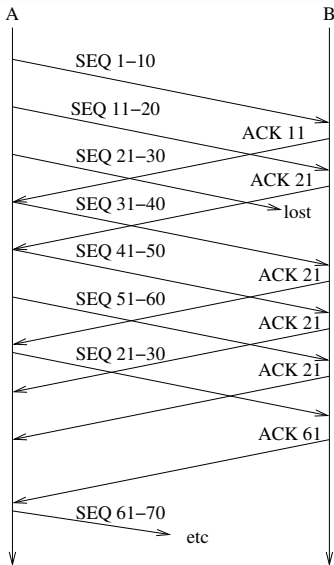The segment containing bytes 21-30 is lost;

ACKing lost segments

# TCP



When B next gets a segment
it still ACKS with 21: that's
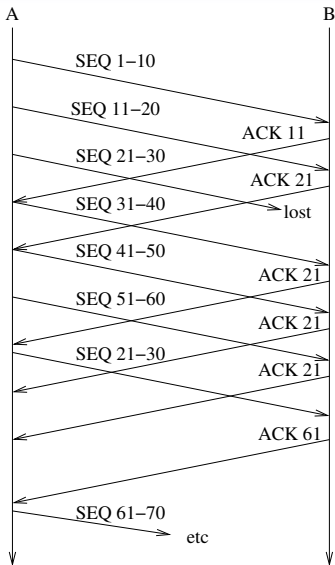the byte it wants next;

ACKing lost segments

# TCP



While the ACK travels back to A, A is still sending new data;
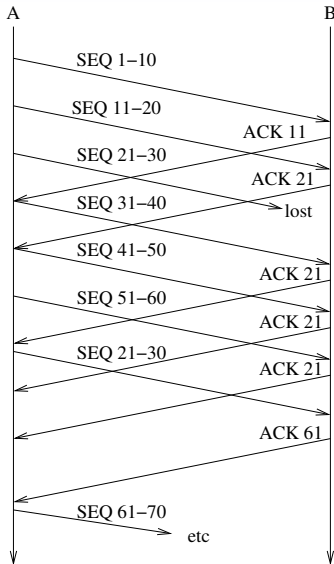
ACKing lost segments

# TCP



A          B

SEQ 1–10

SEQ 11–20

ACK 11

SEQ 21–30

ACK 21

SEQ 31–40          lost

SEQ 41–50

ACK 21

SEQ 51–60

ACK 21

SEQ 21–30

ACK 21

ACK 61

SEQ 61–70          etc

Eventually A gets duplicate ACKs from B: this is a sign of a problem;

ACKing lost segments

# TCP



A resends bytes 21-30;

ACKing lost segments

# TCP



A | B

SEQ 1-10
SEQ 11-20
ACK 11
SEQ 21-30
ACK 21
SEQ 31-40 → lost
SEQ 41-50
ACK 21
SEQ 51-60
ACK 21
SEQ 21-30
ACK 21
ACK 61
SEQ 61-70 → etc

When B gets these bytes it can ACK all the way up to 60

ACKing lost segments

In fact this diagram is not realistic: it is over-simplified to fit on the slide

In fact this diagram is not realistic: it is over-simplified to fit on the slide

TCP specifies that A should continue until it get gets *three duplicate ACKs* (i.e., four ACKs with the same sequence number, not piggybacked on data and not changing the advertised window) before resending

In fact this diagram is not realistic: it is over-simplified to fit on the slide

TCP specifies that A should continue until it get gets *three duplicate ACKs* (i.e., four ACKs with the same sequence number, not piggybacked on data and not changing the advertised window) before resending

This is to avoid triggering resends too easily, e.g., it might be just a case of A's packets being slightly reordered in transit, where a resend is not actually required (remember TCP runs on top of the unreliable IP)

In fact this diagram is not realistic: it is over-simplified to fit on the slide

TCP specifies that A should continue until it get gets *three duplicate ACKs* (i.e., four ACKs with the same sequence number, not piggybacked on data and not changing the advertised window) before resending

This is to avoid triggering resends too easily, e.g., it might be just a case of A's packets being slightly reordered in transit, where a resend is not actually required (remember TCP runs on top of the unreliable IP)

**Exercise** When might we receive many ACKs with the same sequence number, but nothing is in error?

The sequence number wraps around after
$2^{32} - 1 = 4294967295$ bytes

The sequence number wraps around after
$2^{32} - 1 = 4294967295$ bytes

This is under 10 seconds for a 10Gb/s Ethernet

The sequence number wraps around after
$2^{32} - 1 = 4294967295$ bytes

This is under 10 seconds for a 10Gb/s Ethernet

Additional mechanisms to extend the count have had to be
devised in the light of modern fast networks

# TCP
## Sequence numbers

The sequence number wraps around after
$2^{32} - 1 = 4294967295$ bytes

This is under 10 seconds for a 10Gb/s Ethernet

Additional mechanisms to extend the count have had to be devised in the light of modern fast networks

**Exercise** E.g., using the TCP header timestamp option. Read about PAWS

# TCP
## Sequence numbers

The sequence number wraps around after
$2^{32} - 1 = 4294967295$ bytes

This is under 10 seconds for a 10Gb/s Ethernet

Additional mechanisms to extend the count have had to be devised in the light of modern fast networks

**Exercise** E.g., using the TCP header timestamp option. Read about PAWS

Much more on SEQ and ACKing later, but note that sequence numbers solve the segment ordering problem, too

Back to the TCP header

Back to the TCP header

- 4 bits header length: measured in 32 bit words: the header can have options, so is of variable length

# TCP

Back to the TCP header

- 4 bits header length: measured in 32 bit words: the header can have options, so is of variable length

So maximum is 60 bytes. Minimum is the fixed part: 20 bytes

# TCP

Back to the TCP header

- 4 bits header length: measured in 32 bit words: the header can have options, so is of variable length

So maximum is 60 bytes. Minimum is the fixed part: 20 bytes

- Many flags performing various functions

# TCP

Most of these will be described in more detail as we go along:

- URG: urgent data
- ACK: the acknowledgement field is active
- PSH: push this data to the application as fast as possible
- RST: reset (break) the connection
- SYN: synchronise a new connection
- FIN: finish a connection
- ECE: congestion notification
- CWR: congestion window reduced
- 4 reserved bits, set to 0