# Hello!

---

Connection Established!

# TCP - Architecture, Enhancements & Tuning

**Dinesh Dhakal**

Site Reliability Engineer

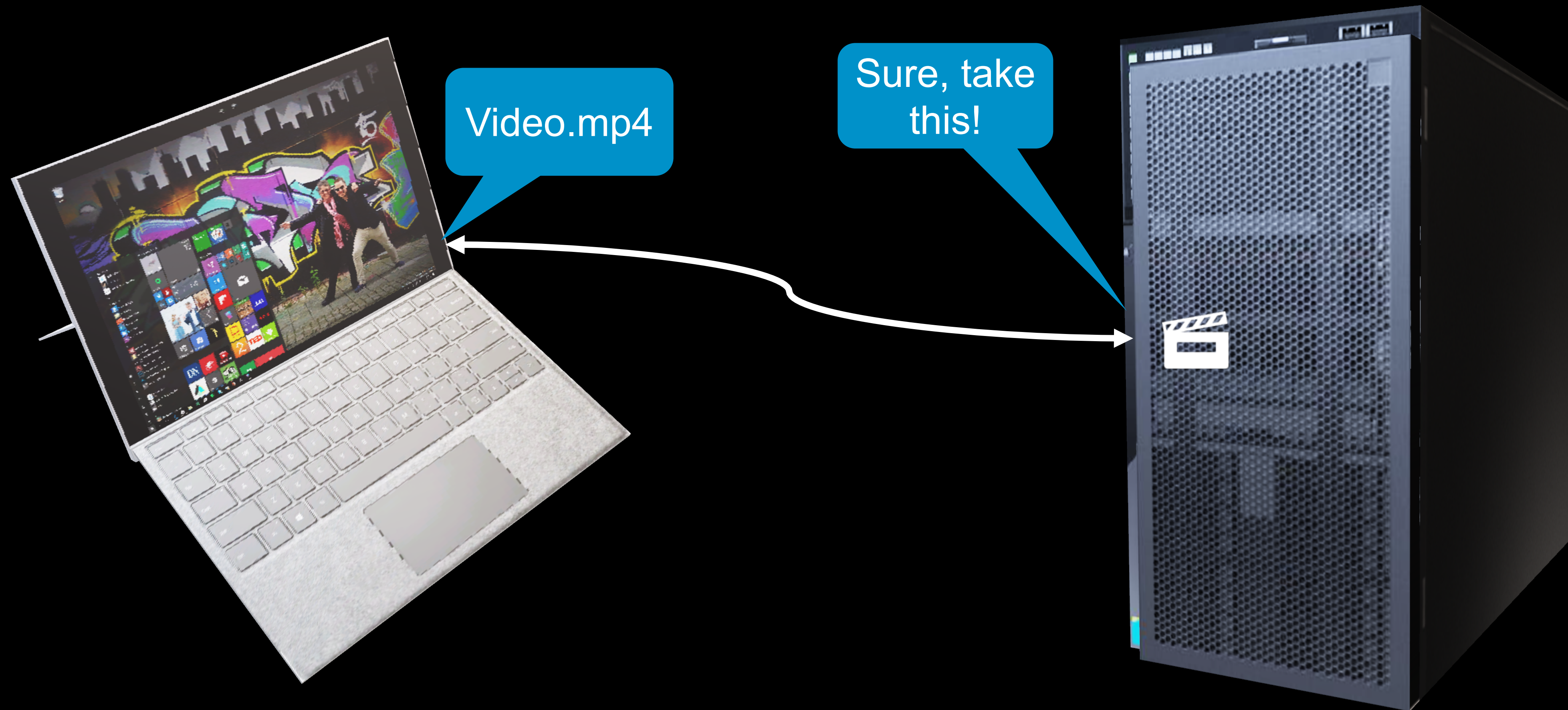# Today's agenda

| | |
|---|---|
| 11:00 | Introductions |
| 11:05 | Core Functionality |
| 11:15 | Enhancements and Extensions |
| 11:30 | Tuning of TCP Parameters on Linux |
| 11:40 | The March Ahead |
| 11:45 | Q&A |

# The Network Stack

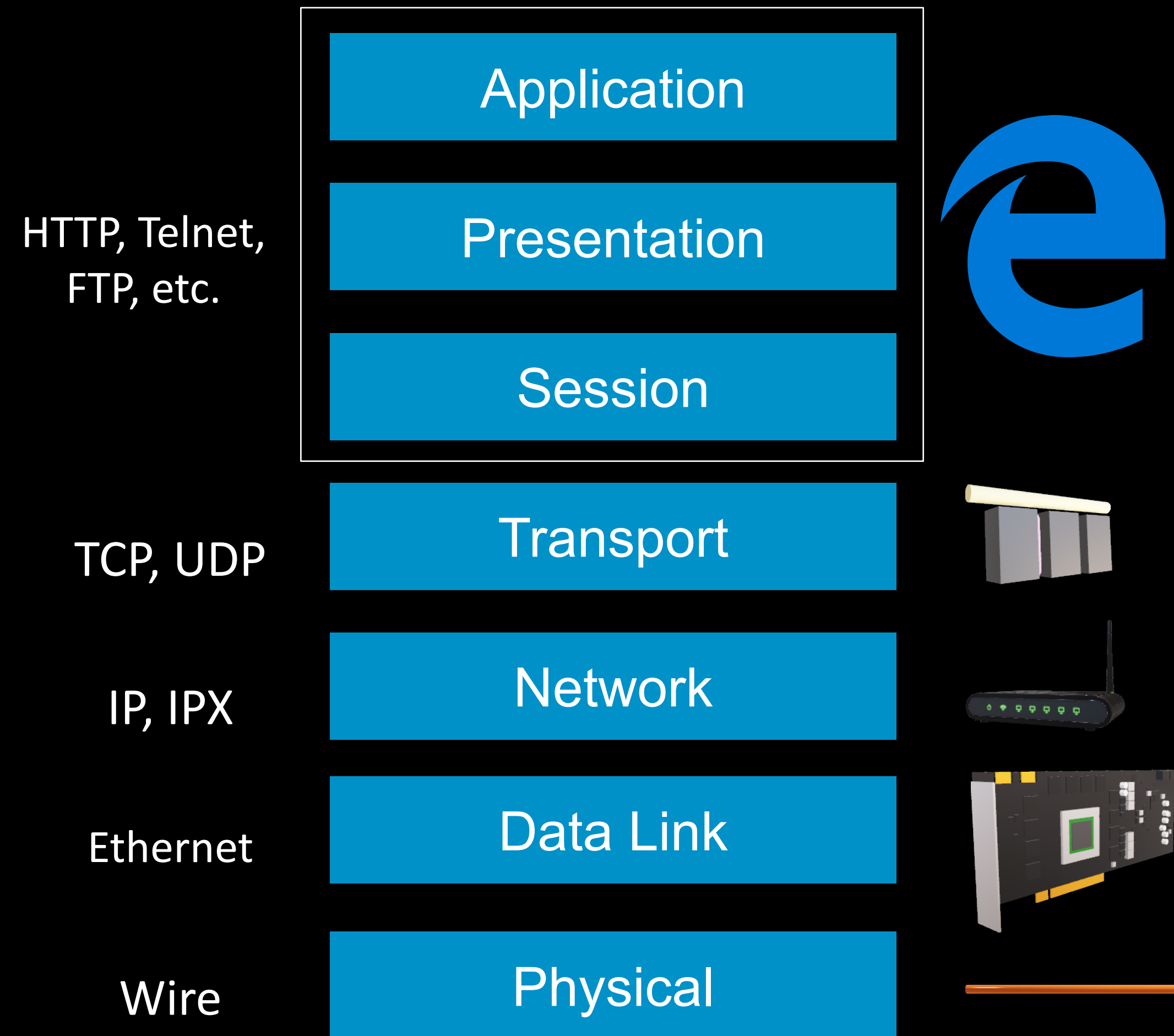| | | |
|---|---|---|
| | **Application** | |
| HTTP, Telnet, FTP, etc. | **Presentation** | |
| | **Session** | |
| TCP, UDP | **Transport** | |
| IP, IPX | **Network** | |
| Ethernet | **Data Link** | |
| Wire | **Physical** | |

# What should the Transport layer do?

## Problems

- Applications send byte streams

- Underlying IP network is stateless

- Devices are of varied capabilities

- Multiple processes need reliable communication

- Cannot control all the variables

## Requirements

- Ordered Segmentation

- Stateful Communication

- Flow Control

- Multiplexing

- Reliability and Congestion Control
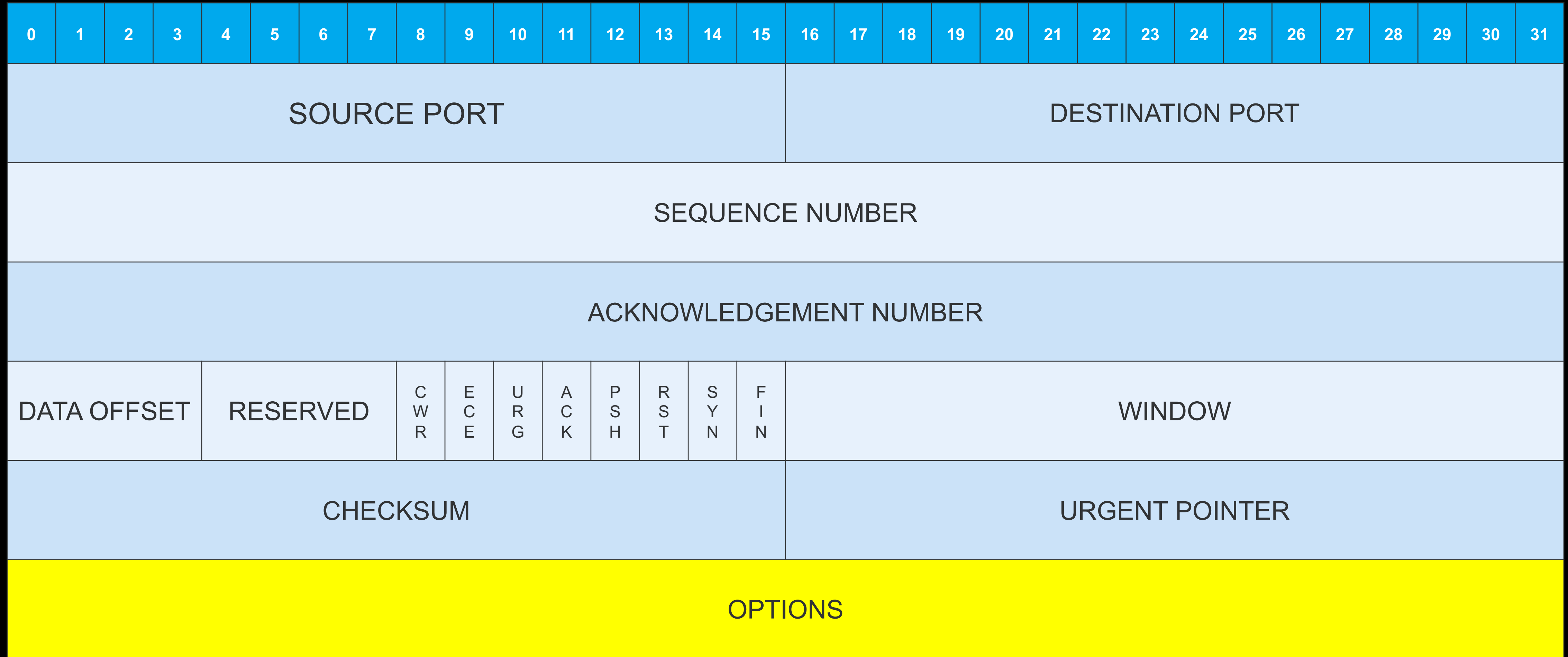
# TCP – Architecture

# TCP Core concepts

## Requirements

- Ordered Segmentation

- Stateful Communication

- Flow Control

- Multiplexing

- Reliability and Congestion Control

## How TCP addresses it

- Sequence Numbers

- Connections

- TCP Window Size

- Port Numbers

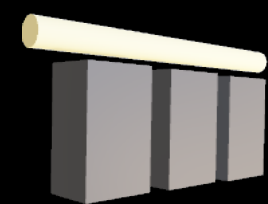- Acknowledgements and Retransmissions

# TCP Header

| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 |
|---|---|---|---|---|---|---|---|---|---|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|----|
| SOURCE PORT | | | | | | | | | | | | | | | | DESTINATION PORT | | | | | | | | | | | | | | | |
| SEQUENCE NUMBER | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| ACKNOWLEDGEMENT NUMBER | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |
| DATA OFFSET | | | | RESERVED | | | | CWR | ECE | URG | ACK | PSH | RST | SYN | FIN | WINDOW | | | | | | | | | | | | | | | |
| CHECKSUM | | | | | | | | | | | | | | | | URGENT POINTER | | | | | | | | | | | | | | | |
| OPTIONS | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | | |

# TCP Segments

101110101011101000010101010100010100100010001000100101010010101010101…

| 8080 | 22003 |
|------|-------|
| SEQ : 1101 | |
| ACK : 2201 | |
| Data : 101110101… | |

| 8080 | 22003 |
|------|-------|
| SEQ : 1102 | |
| ACK : 2202 | |
| Data : 101110101… | |

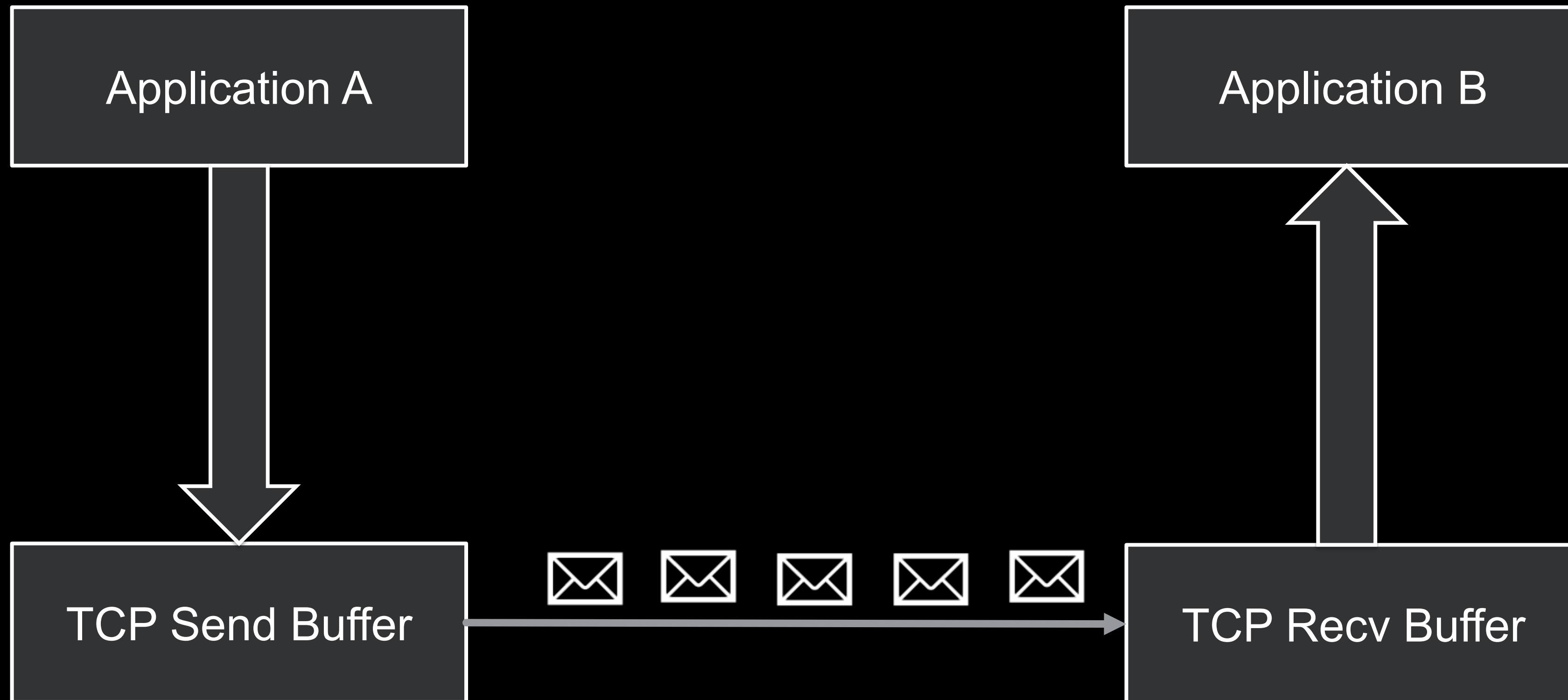| 8080 | 22003 |
|------|-------|
| SEQ : 1103 | |
| ACK : 2203 | |
| Data : 101110101… | |

# Connection Establishment – 3 way handshake

# TCP Sockets

# Flow Control – Sliding Window

# Retransmission



SEQ: 1101

ACK: 1101

SEQ: 1102

Retransmission timer

SEQ: 1102

ACK: 1102

# Enhancements

# Slow Start Phase



Receiver

Sender

cwnd = 1

SEQ: 1101

ACK: 1102

SEQ: 1102

cwnd = 2

SEQ: 1103

ACK: 1103

ACK: 1104

SEQ: 1104 , 1105,1006,1007

Cwnd=4

# Congestion Avoidance

Receiver

Sender

Cwnd=4

SEQ: 1104 , 1105,1006,1007

ACK: 1105,  1106, 1107, 1008
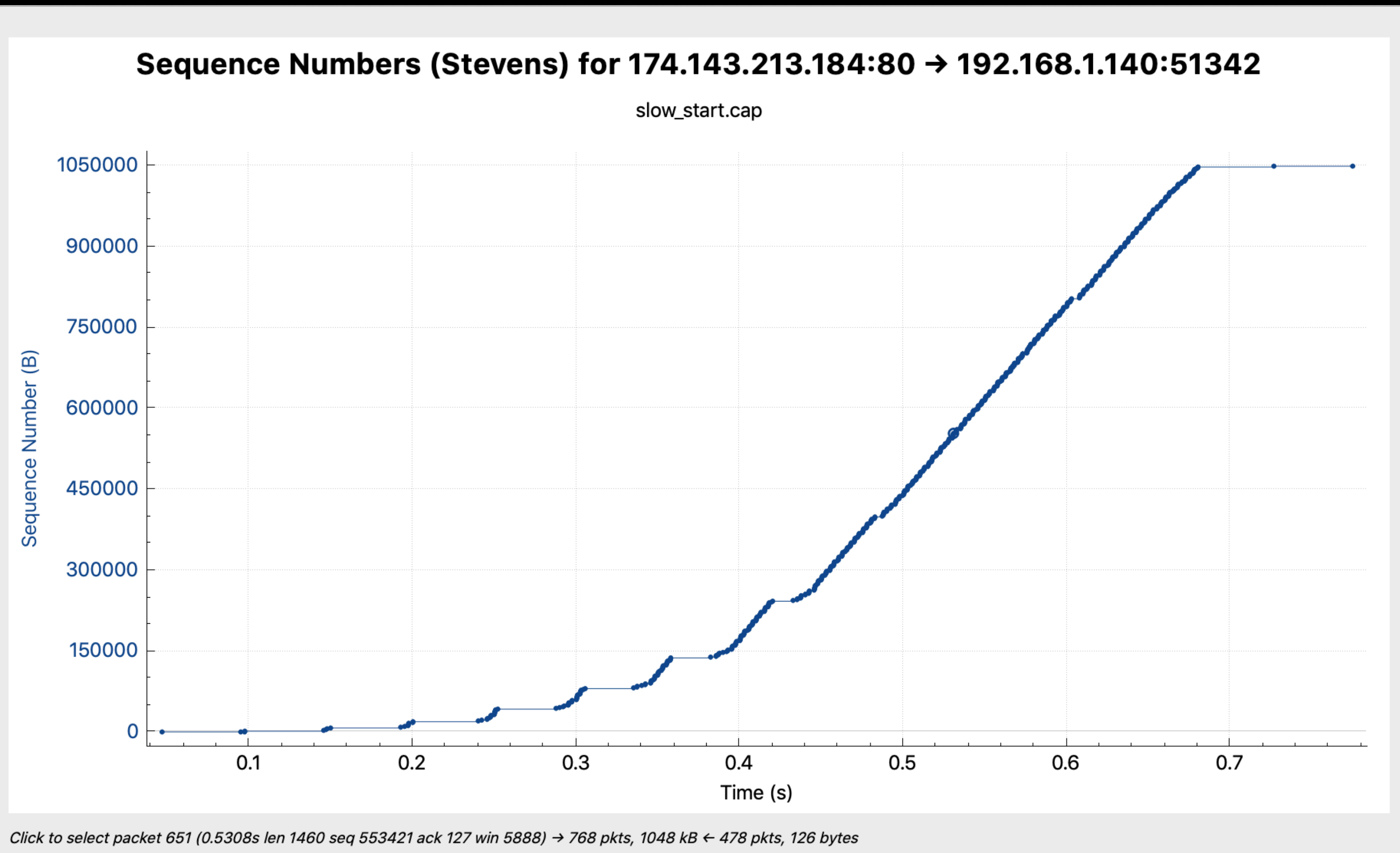
SEQ: 1108 , 1109, 1010, 1011, 1012
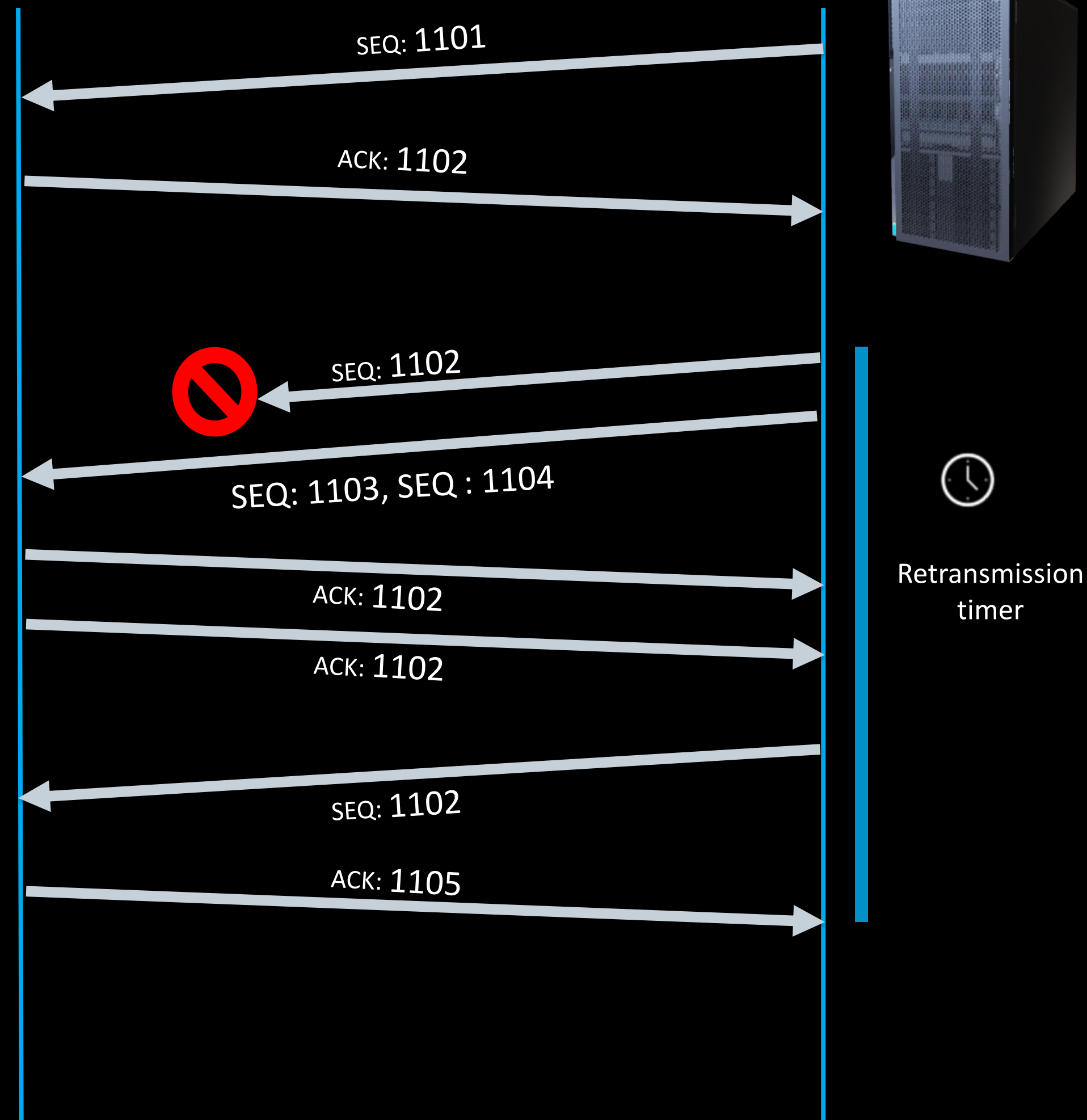
Cwnd=5

# Congestion Control Enhancements

- Slow Start

  - Slow start when Congestion window (cwnd) < slow start threshold (ssthresh)

    - Typically, ssthresh starts at 65535 bytes.

    - cwnd += min (N, SMSS)                              SMSS – Sender Max Segment Size

- Congestion Avoidance

  - Congestion avoidance when cwnd > ssthresh

    - On ACK: cwnd += SMSS*SMSS/cwnd

  - ssthresh = min(cwnd,rwnd) / 2 when congestion

# TCP Slow Start



Sequence Numbers (Stevens) for 174.143.213.184:80 → 192.168.1.140:51342

slow_start.cap

*Click to select packet 651 (0.5308s len 1460 seq 553421 ack 127 win 5888) → 768 pkts, 1048 kB ← 478 pkts, 126 bytes*

# Fast Recovery

- Receiver sends duplicate ack → Segments have left the network

- Artificially inflates the cwnd as segments sent are *assumed* to have left network

  - cwnd = ssthresh + 3 * SMSS

  - Every Additional ack : cwnd = cwnd + SMSS

- When a new segment is acknowledged –

  - Cwnd = ssthresh

# Loss Recovery Enhancements

- TCP Selective Acknowledgment Options

  - Informs the sender about OOR segments received

  - Uses the TCP options fields to acknowledge the received segments

- Partial acks

  - Aims to reduce the number of duplicate acks needed for retransmit

  - Specifically useful for cases of continuous packet loss

  - Every partial ack in the gap triggers retransmit of next unacked segment

# TCP Tuning

# Bandwidth Delay Product

- The amount of data that can be in transit in the network

- Product of Bandwidth and Delay (RTT)

  - 1 Mbps X 70ms = 0.88 MByte

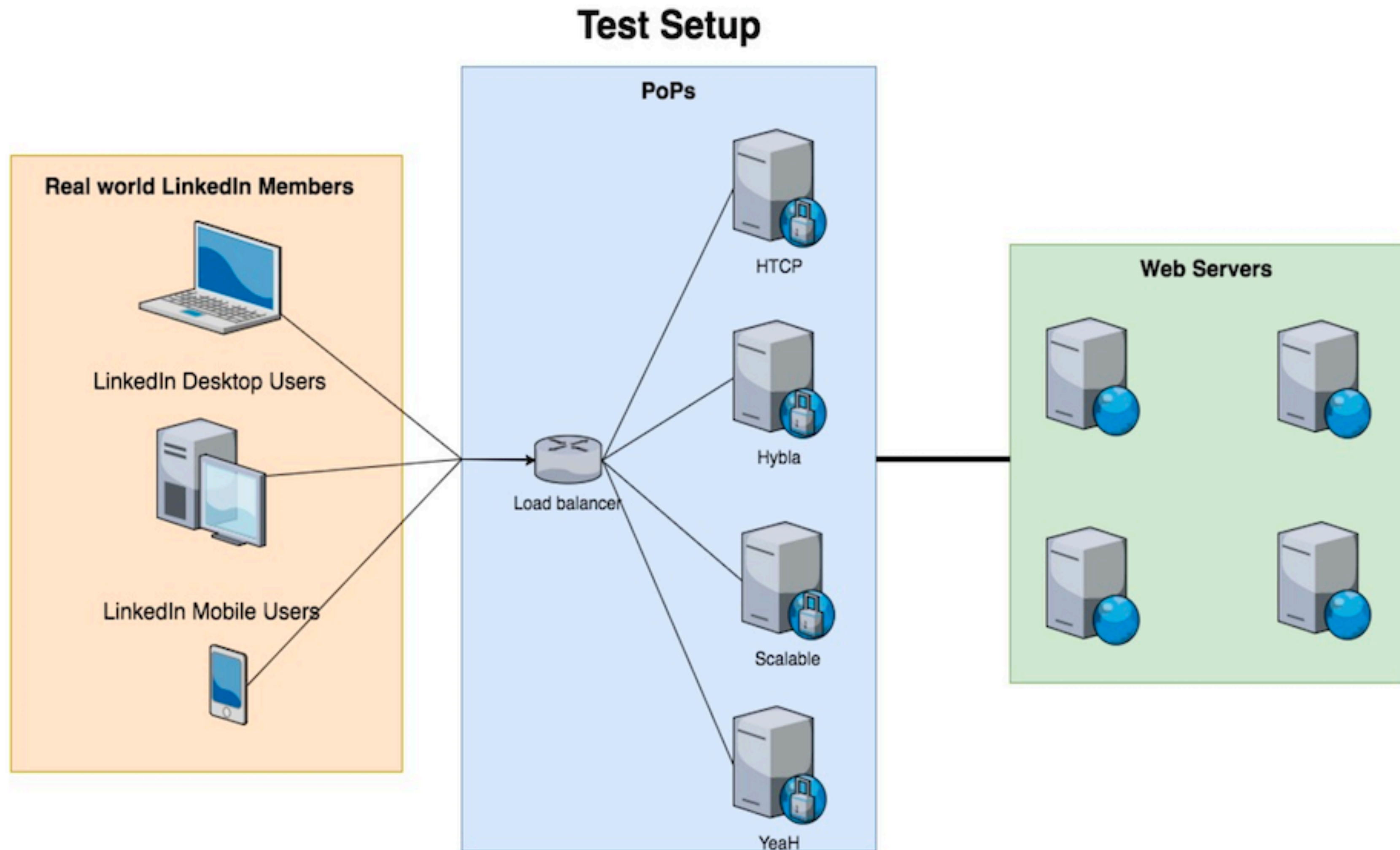- Buffer sizes can be appropriately tuned to gain max utilization of bandwidth

# Buffers

- TCP Buffer sizes can be tuned for optimal use of Bandwidth

- net.core.rmem_max = 268435456

- net.core.wmem_max = 268435456

- net.ipv4.tcp_rmem = 4096 87380 134217728
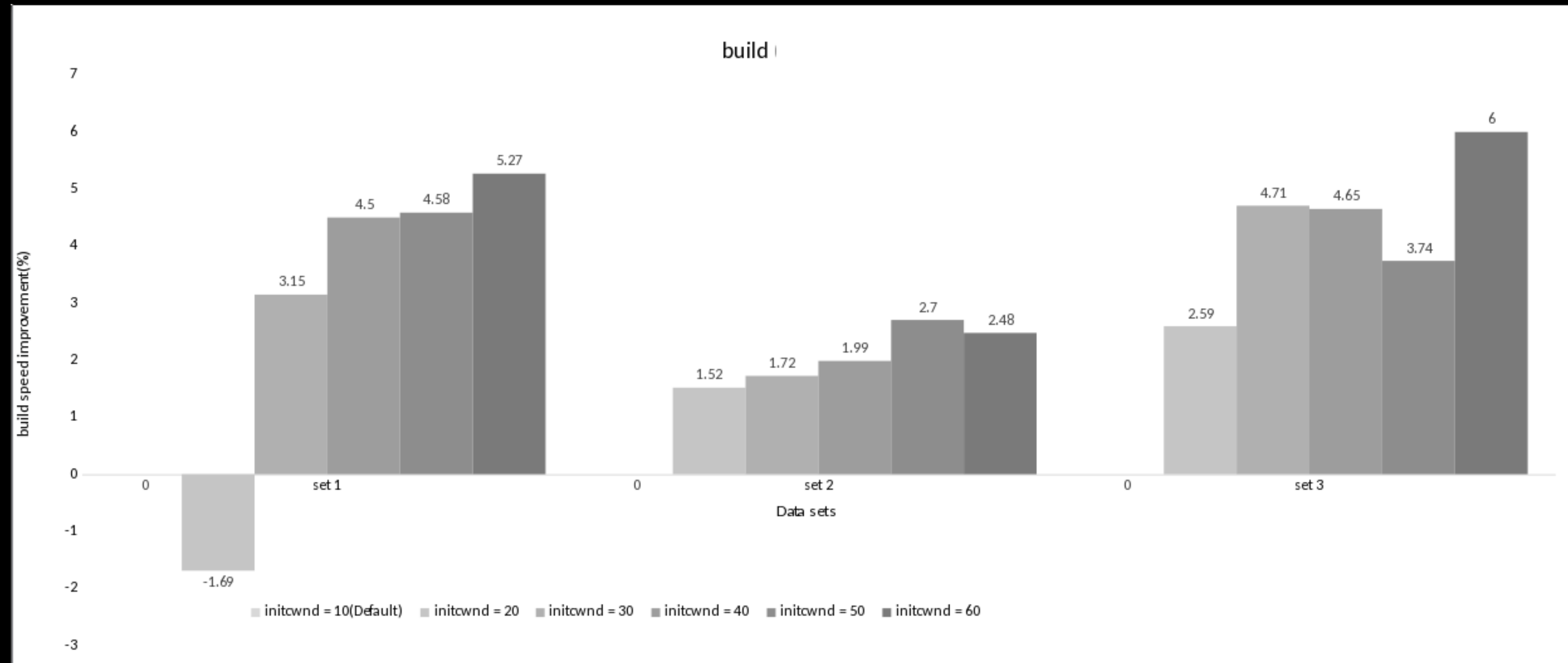  net.ipv4.tcp_wmem = 4096 65536 134217728

# Some more Parameters…

- Enable Selective Ack
  net.ipv4.tcp_sack = 1

- Enable Window Scaling -
  net.ipv4.tcp_window_scaling = 1

- MTU probing
  net.ipv4.tcp_mtu_probing = 0

# Test Setup

# Initial Congestion Window

- Increasing initcwnd can reduce the number of Round Trips thus increasing performance
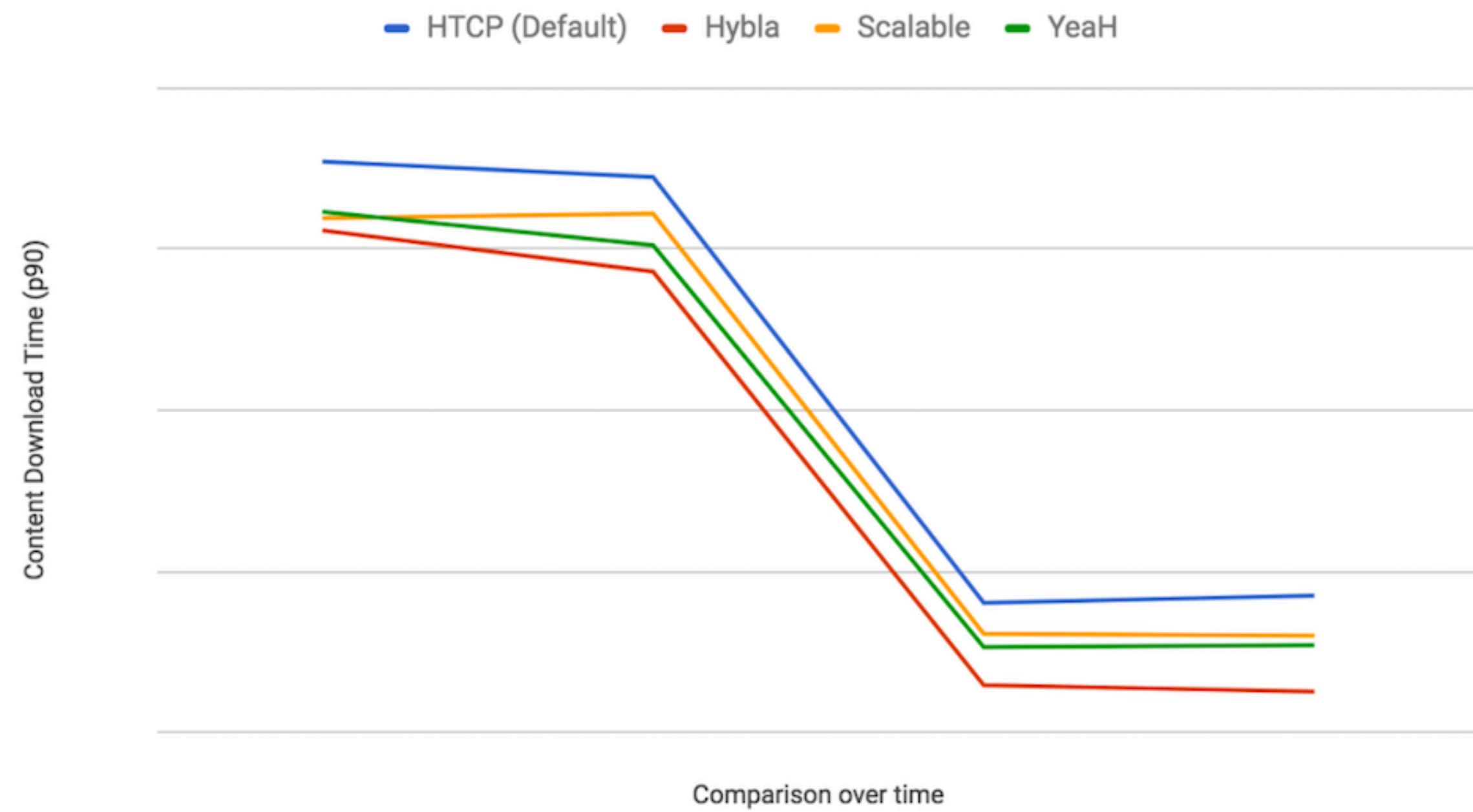
# Congestion Control Algorithms

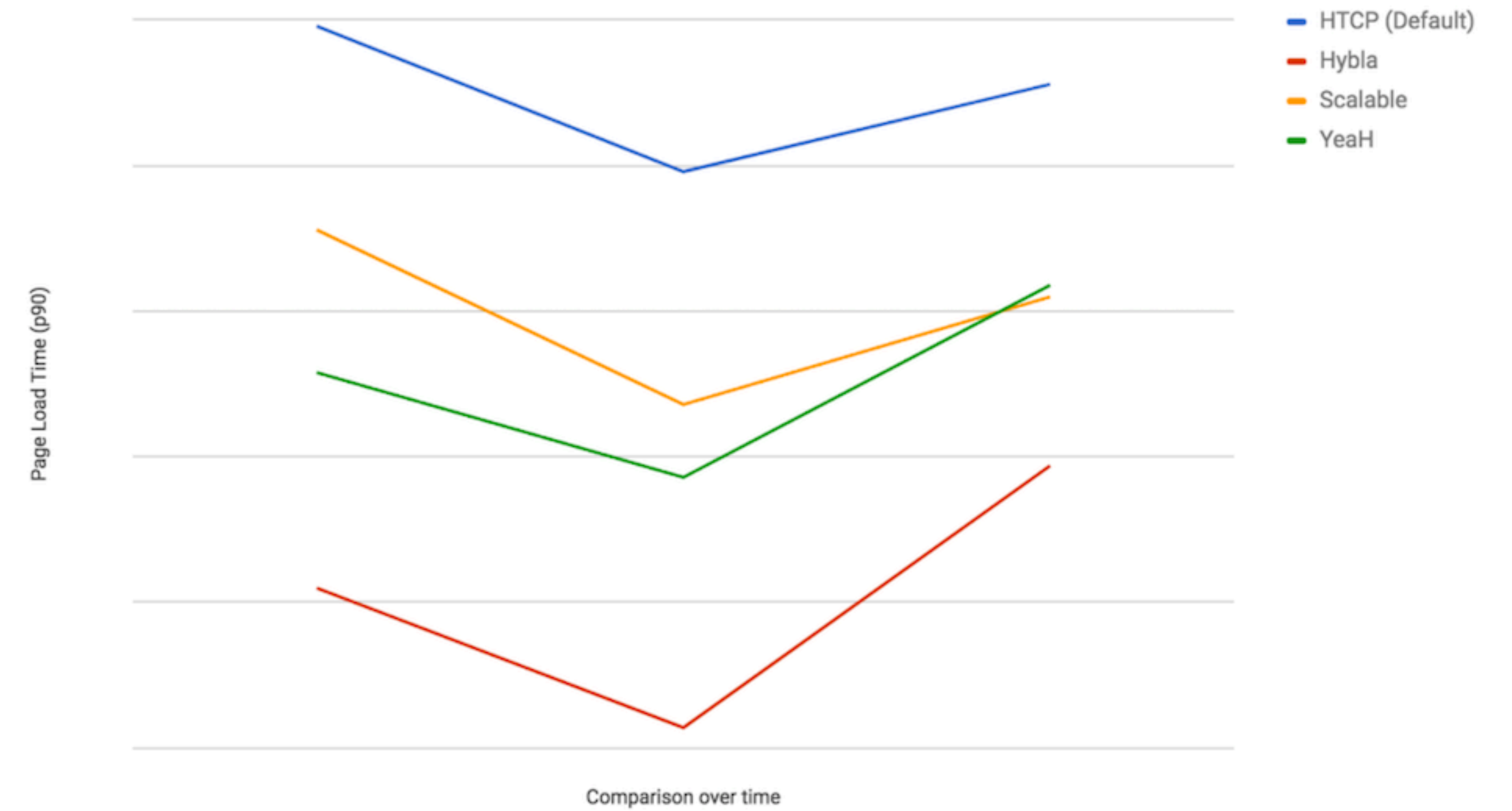| Algorithm | What it does best |
| --- | --- |
| TCP-Hybla | Built for networks with long round trip delays. Window update is based on a ratio of current RTT and a reference RTT0. |
| TCP-Scalable | Built for performance on high-speed, wide area networks. Window updates use fixed increase and decrease parameters. |
| TCP-YeaH | Built to be fair, efficient, and prevent Lossy-Link penalties. Switches between fast and slow modes, based on an estimate of queued packets. |
| HTCP | Built for long distance, high-speed transmission. Window updates are based on time since last loss event. This is the default algorithm on our Linux machines. |

# Site Speed Improvements

# The March Ahead

# QUIC

- Intended to eventually replace TCP and TLS on the web

- Provides security features like authentication and encryption, that are typically handled by a higher layer protocol

- Establishes multiple connections over UDP

- Avoids head of line blocking by using multiple HTTP streams mapped to multiple QUIC connections

# SCTP

- TCP provides both reliable data transfer and strict transmission ordered delivery of data

- Head-of-line blocking in TCP causes delays

- SCTP is a message based reliable protocol

- Reliable transmission of both ordered and unordered data streams.

- Multihoming support and transparent fail over

# Thank You