
Transmission Control Protocol (TCP)

Netzwerkgrundlagen (NWG2)

Markus Zeilinger¹

¹FH Oberösterreich
Department Sichere Informationssysteme

Sommersemester 2023



UNIVERSITY
OF APPLIED SCIENCES
UPPER AUSTRIA

*Alle Materialien, die im Rahmen dieser LVA durch den LVA-Leiter zur Verfügung gestellt werden, wie zum Beispiel **Foliensätze, Audio-Aufnahmen, Übungszettel, Musterlösungen, ...** dürfen **ohne explizite Genehmigung** durch den LVA-Leiter **NICHT** weitergegeben werden!*

Grundlagen

TCP Header

TCP Connection Management

TCP Zuverlässigkeit

TCP Flusskontrolle

TCP Segmentation und Optionen

Zusatzmaterial

- ▶ Transmission Control Protocol (TCP) - RFC 793
- ▶ Verbindungsorientierte und zuverlässige Übertragung eines Byte-Stroms (Stream) zwischen Prozessen auf (nicht notwendigerweise) unterschiedlichen Systemen über ein unzuverlässiges Netzwerk.
- ▶ TCP betrachtet von der Anwendung kommende Daten (via write() Call) als prinzipiell unendlichen Byte-Strom (Stream-Orientierung) und kümmert sich selbst um die Einteilung der Daten in Segmente.

Transmission Control Protocol (TCP)

Schlüsselaspekte

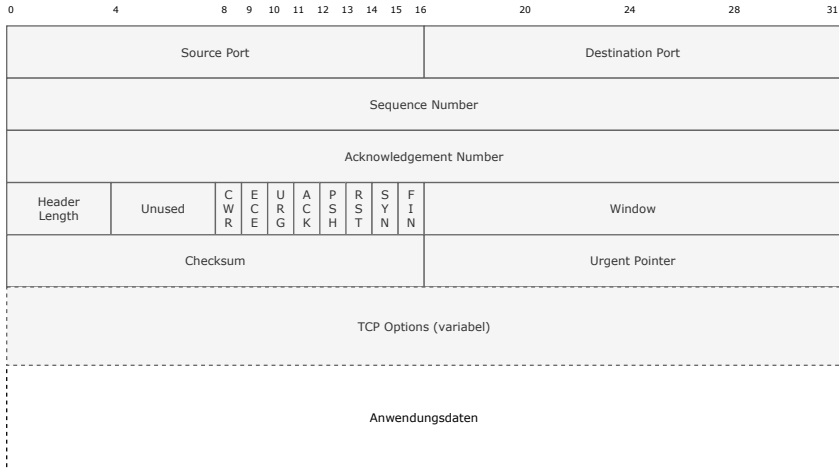
- ▶ **Adressierung über Ports** (lokaler und remote Socket bilden die Kommunikationsendpunkte).
- ▶ **Verbindungsorientierung**: Aufbau, Verwaltung und Abbau von bidirektionalen Verbindungen (**TCP Finite State Machine [FSM]**) (Analogie: Telefonanruf).
- ▶ **Paketierung** (in **Segmente**) des Byte-Stroms (stream-orientiert) und Weiterleitung an IP.
- ▶ **Full-Duplex**, d. h. in einer TCP-Verbindung können **gleichzeitig Daten in beide Richtungen** fließen.
- ▶ **Zuverlässigkeit** (Reliability) durch **Positive Acknowledgment with Retransmission (PAR)**.
- ▶ **Flusskontrolle** (Flow Control) und **Überlastungsüberwachung** (congestion control).

TCP RFCs (Auszug der wichtigsten)

► RFC 7414, A Roadmap for Transmission Control Protocol (TCP) Specification Documents

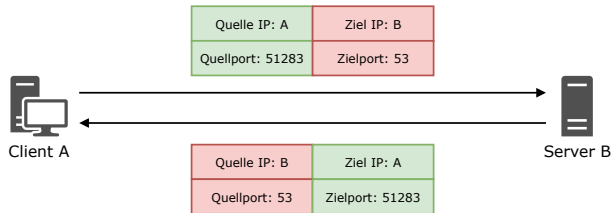
RFC	Jahr	Titel	Beschreibung
793	1981	Transmission Control Protocol	Basisstandard
1122	1989	Requirements for Internet Hosts	Implementierungsdetails zu u. a. TCP
2018	1996	TCP Selective Acknowledgment Options	Selektive Angabe von TCP Segmenten, die neuerlich übertragen werden müssen.
5681	2009	TCP Congestion Control	Beschreibung von vier TCP Congestion Control Algorithmen: Slow Start, Congestion Avoidance, Fast Retransmit und Fast Recovery
6298	2011	Computing TCP's Retransmission Timer	Algorithmus zur Berechnung des Retransmission Timers.
6691	2012	TCP Options and Maximum Segment Size (MSS)	TCP Maximum Segment Size Option
7323	2014	TCP Extensions for High Performance	Mechanismen für TCP auf High-Speed Links (z. B. TCP Window Scale Option).

TCP Header I



TCP Header II

- ▶ **Source (Quell-) Port:** Port des sendenden Prozesses am Quellsystem.
 - ▶ Client → Server: Dynamic/Private Port > 49152
 - ▶ Server → Client: Well-Known oder Registered Port des Services
- ▶ **Destination (Ziel-) Port:** Port des empfangenden Prozesses am Zielsystem.
 - ▶ Client → Server: Well-Known oder Registered Port des Services
 - ▶ Server → Client: Dynamic/Private Port > 49152



TCP Header III

Zusammenhang Bytestrom, Sequence Number und Acknowledgement Number

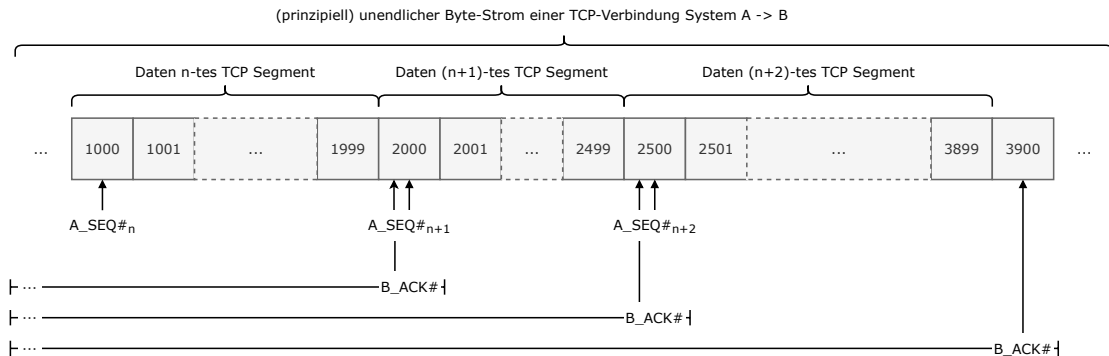
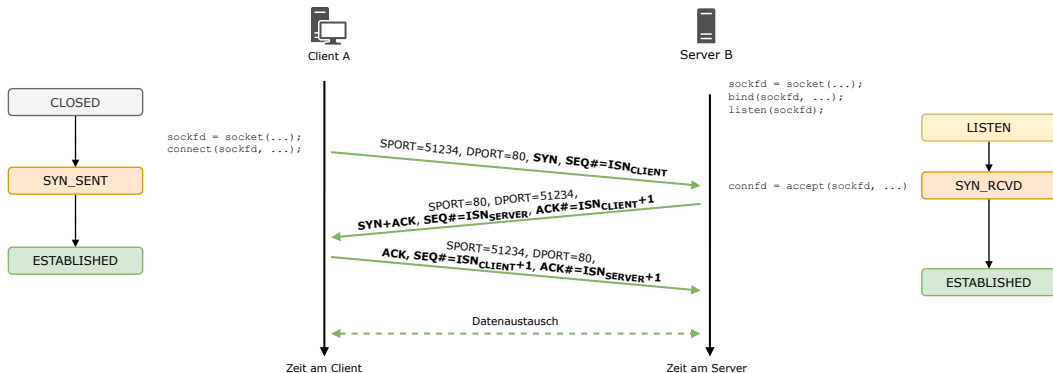


Abbildung 1: Zusammenhang Bytestrom, Sequence Number, Acknowledgement Number (basierend auf [1])

- ▶ **SYN** (**Synchronize**): Signalisiert den Wunsch eines Systems A zum Verbindungsaufbau mit System B.
- ▶ **FIN** (**Finish**): Signalisiert den Wunsch eines Systems A zum Verbindungsabbau mit System B.
- ▶ **RST** (**Reset**): Zurücksetzung (Abbruch) der TCP Verbindung.
- ▶ **ACK** (**Acknowledgement**): Das Segment enthält eine Bestätigung, d. h. der Wert im ACK#-Feld ist gültig.
- ▶ **PSH** (**Push**): Der Empfänger muss die Daten im Segment sofort an die Anwendung weitergeben.

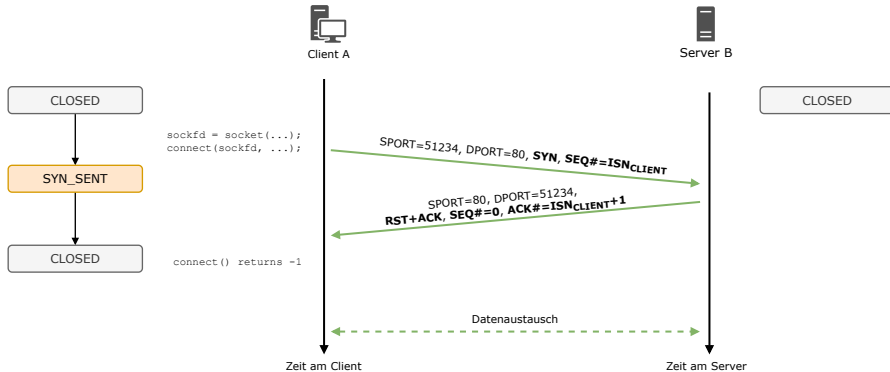
TCP Connection Management I

TCP Verbindungsaufbau I



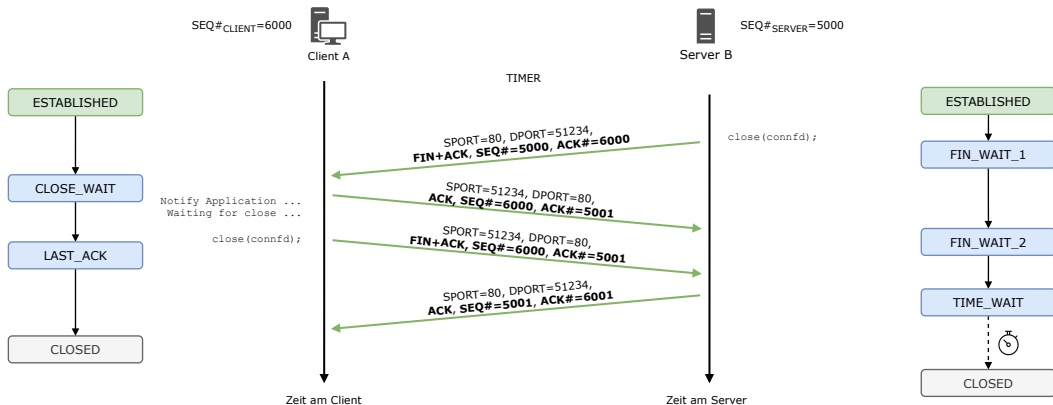
TCP Connection Management III

TCP Verbindungsaufbau II



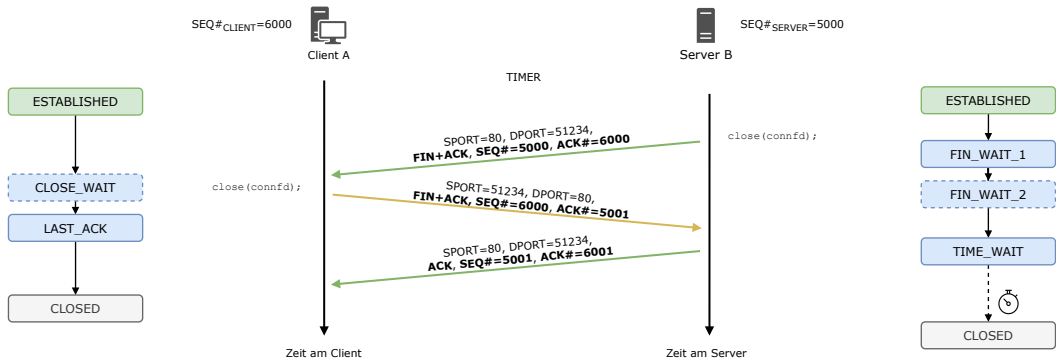
TCP Connection Management V

TCP Verbindungsabbau I

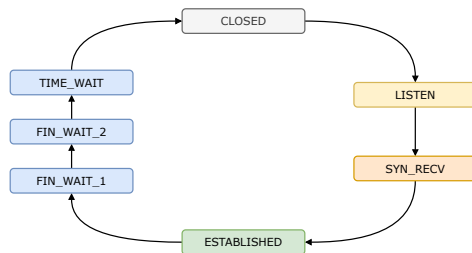
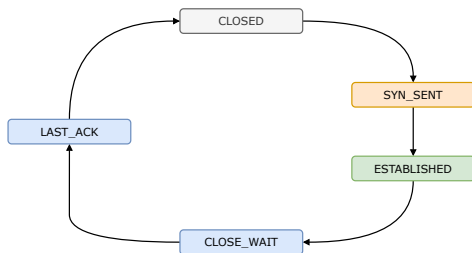


TCP Connection Management V

TCP Verbindungsabbau II



- Typische TCP Zustandsabfolge für Clients (links) und Server (rechts) (basierend auf [1]):



- ▶ Zuverlässigkeit (Reliability) wird bei TCP über **Positive Acknowledgement with Retransmission (PAR)** realisiert.
- ▶ **Positive Acknowledgement**
 - ▶ Der Empfänger gibt dem Sender positive Rückmeldung über erfolgreich empfangene Bytes (Acknowledgement Number, ACK#).
 - ▶ **Cumulative Acknowledgement** = TCP bestätigt immer alles bis zum letzten korrekt empfangenen Byte, d. h. auch schon alle zuvor bestätigten.
- ▶ **Retransmission**
 - ▶ Erhält ein Sender innerhalb einer definierten Frist (**Retransmission Timeout, RTO**) keine Bestätigung für in einem Segment enthaltene Daten, schickt er das Segment erneut.

TCP Zuverlässigkeit (Reliability) II

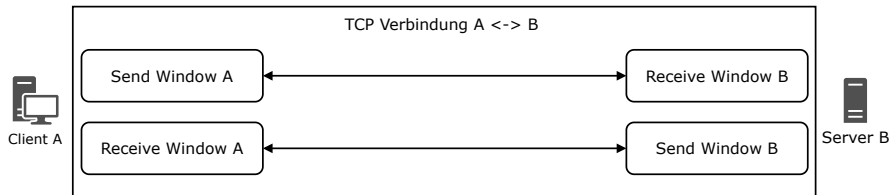
Retransmission Timeout (RTO)

- ▶ Die Berechnung des **TCP Retransmission Timeouts (RTO)** ist in RFC 6298 [3] beschrieben.
- ▶ **Problemfeld** der RTO Berechnung: RTO zu kurz → unnötige Retransmissions, RTO zu lang → unnötig hohe Latenzen
- ▶ Die Basis bilden **laufende RTT Messungen** ($SampleRTT$, R').
 - ▶ TCP RTT = Zeit vom Start der Sendeoperation für ein Segment bis zum Empfang der Bestätigung dafür.
- ▶ **Startwert** für RTO = 1 Sekunde, Erreichen des RTOs → **Verdopplung** des RTOs .
- ▶ Retransmission Timer läuft **nicht für jedes Segment** einzeln (zu großer Overhead), sondern immer für eines, für das noch eine Bestätigung ausständig ist.

- ▶ Eine Bestätigung kann **Daten eines oder auch mehrere Segmente** umfassen.
- ▶ Mögliche **Fälle** und **Strategien** (u. a. RFC 5681 [4]):
 1. Segment mit erwarteter SEQ# trifft ein + alle Daten davor sind schon bestätigt → **Delayed Acknowledgement** = Bestätigung darf um 500 ms verzögert werden.
 2. Segment mit erwarteter SEQ# trifft ein + Daten eines weiteren Segments sind noch nicht bestätigt → sofortiges Senden einer Bestätigung für damit zwei Segmente.
 3. Segment mit unerwarteter SEQ# trifft ein (Lücke) → **Duplicate Acknowledgement**, d. h. dreimaliges Bestätigen bis zum Byte unmittelbar vor der Lücke (bewirkt **Fast Retransmit** ohne Abwarten des RTOs).
 4. Segment mit Daten in einer Lücke trifft ein → sofortiges Senden einer Bestätigung soweit als möglich.

- ▶ Flusskontrolle ermöglicht es einem Empfänger einen Sender in dessen Übertragungsgeschwindigkeit zu steuern.
 - ▶ Ziel: Ein schneller Sender soll einen langsamen Empfänger nicht überlasten.
- ▶ TCP verwendet dafür einen Sliding-Window-Ansatz.
- ▶ Der Empfänger signalisiert dem Sender über das Window-Feld im TCP Header wie viele Bytes der Sender senden darf (in einer beliebigen Anzahl von Segmenten), bevor er längstens auf eine Bestätigung durch den Empfänger warten muss.
- ▶ Der Empfänger verwaltet dafür eine Receive Window, der Sender ein Send Window.
- ▶ TCP ist Full-Duplex → beide Teilnehmer sind Sender und Empfänger gleichzeitig → beide haben je ein Receive und ein Send Window.

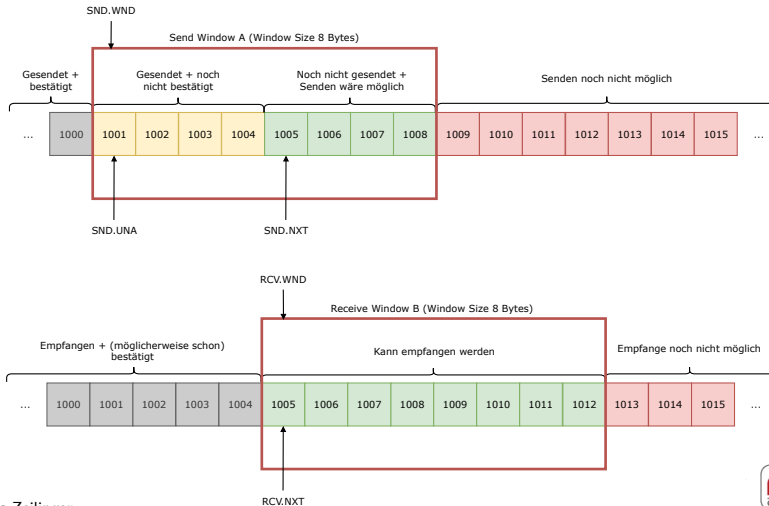
- ▶ Server B bestimmt die Größe seines Receive Windows und damit die Größe des Send Windows von Client A.
- ▶ Client A bestimmt die Größe seines Receive Windows und damit die Größe des Send Windows von Server B.



- ▶ Der TCP Standard (RFC 793, Section 3.2) definierte eine Reihe von Variablen zur Verwaltung einer TCP Verbindung in einem Transmission Control Block (TCB).
- ▶ Für die Flusskontrolle sind dabei u. a. relevant:
 - ▶ `SND.WND` = Größe des Send Windows in Bytes.
 - ▶ `SND.UNA` = Nummer des ersten schon gesendeten aber noch nicht bestätigten Bytes (Send Unacknowledged).
 - ▶ `SND.NXT` = Nummer des nächsten zu sendenden Bytes.
 - ▶ `RCV.WND` = Größe des Receive Windows in Bytes.
 - ▶ `RCV.NXT` = Nummer des nächsten erwarteten Bytes.

TCP Flusskontrolle IV

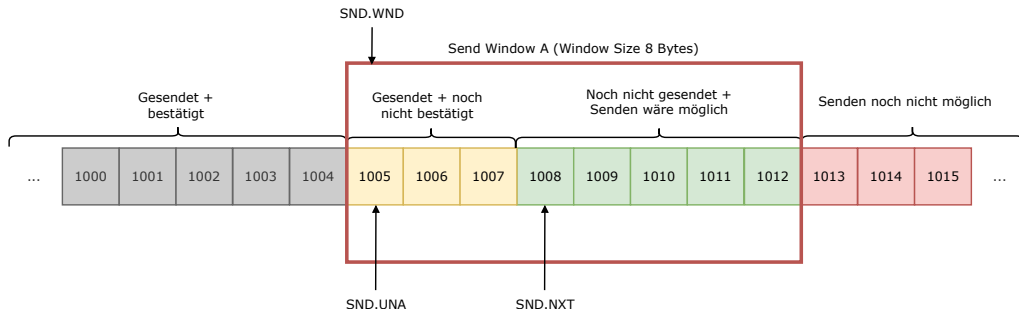
Beispiel I - Ausgangslage TCP Verbindung zw. Client A und Server B



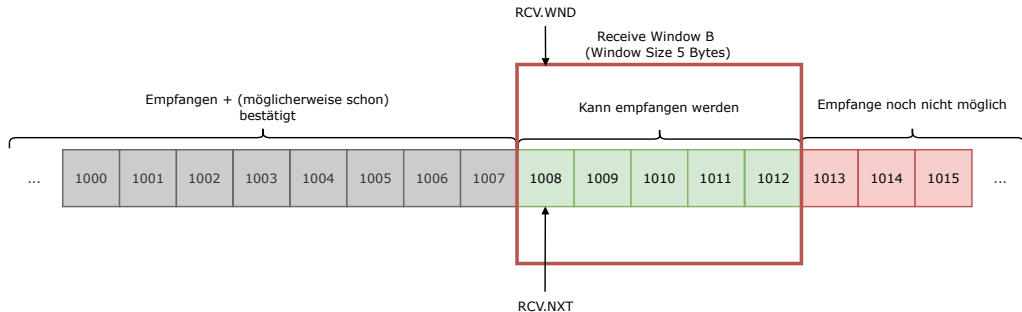
TCP Flusskontrolle V

Beispiel II

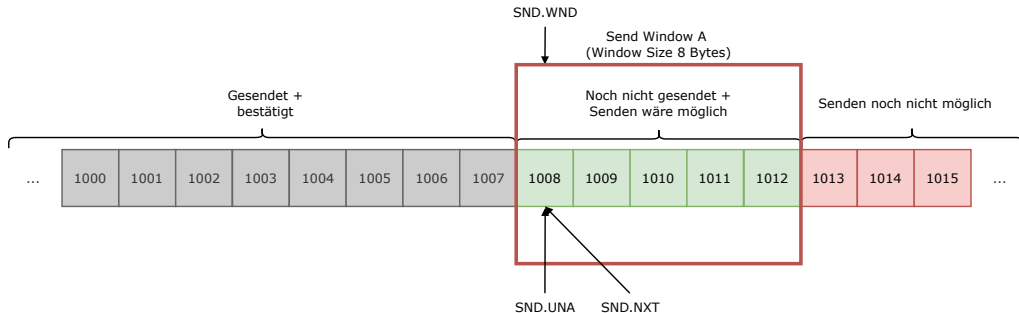
- ▶ Server B schickt eine Bestätigung bis inkl. Byte 1004 + Client A empfängt diese + Client A sendet 3 Byte Daten.



- ▶ Server B empfängt diese Daten + schickt eine Bestätigung bis inkl. Byte 1007 + schickt eine Aktualisierung der Fenstergröße auf 5 Bytes.



- ▶ Client A empfängt die Bestätigung und die Aktualisierung der Fenstergröße.



- ▶ Wann ist eine TCP Segment "voll"?
 - ▶ TCP orientiert sich dabei daran, was sein Gegenüber maximal empfangen kann.
- ▶ TCP Maximum Segment Size (MSS) Option (RFC 6691)
 - ▶ Die Teilnehmer einer TCP Verbindung signalisieren sich im TCP 3-Wege-Handshake die maximale Größe eines für sie empfangbaren TCP Segments.
 - ▶ Diese Größe hängt von der verwendeten Netzwerktechnologie (z. B. Ethernet, WLAN, ...) ab.
 - ▶ Beispiel: Ethernet + IPv4 → MSS 1460 Bytes (Ethernet MTU¹ 1500 Bytes - 20 TCP Header - 20 Byte IPv4 Header = 1460 Byte)

¹ MTU = Maximum Transmission Unit = größt mögliche Datenmenge, die über eine Netzwerktechnologie übertragen werden kann (z. B. Ethernet 1500 Bytes) → VO Teil zu Ethernet

TCP Window Scale Option (RFC 7323)

- ▶ Das Window-Feld im TCP Header hat eine Länge von 16 Bit → die maximale Fenstergröße ist damit $2^{16} = 65535$ Bytes.
- ▶ Für moderne Netzwerke (z. B. mit Datenraten > 1 Gbps) ist das klein und ineffizient.
- ▶ Die **TCP Window Scale Option** (RFC 7323) ermöglicht es, den Wert im Window-Feld um einen **Faktor 2^x mit $x \leq 14$** zu skalieren.
- ▶ Die maximal mögliche Fenstergröße ist somit $(2^{16} - 1) \cdot 2^{14} \approx 1$ GB.
- ▶ Die Verfügbarkeit der Option und der Skalierungsfaktor werden im SYN-Segment des TCP 3-Wege-Handshakes signalisiert.



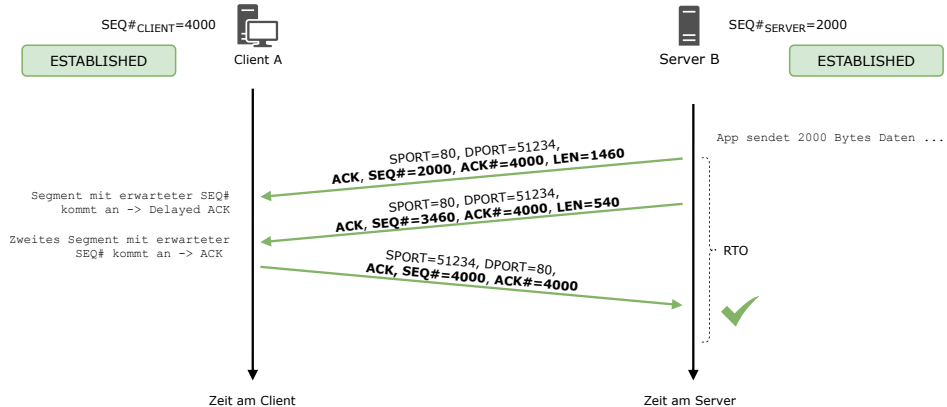
Robustness Principle (John Postel [2])

[...] TCP implementations will follow a general principle of robustness: be conservative in what you do, be liberal in what you accept from others. [...]

- ▶ Was bedeutet das Robustness Principle?
- ▶ Was ist aus heutiger Sicht davon zu halten (s. [6])?

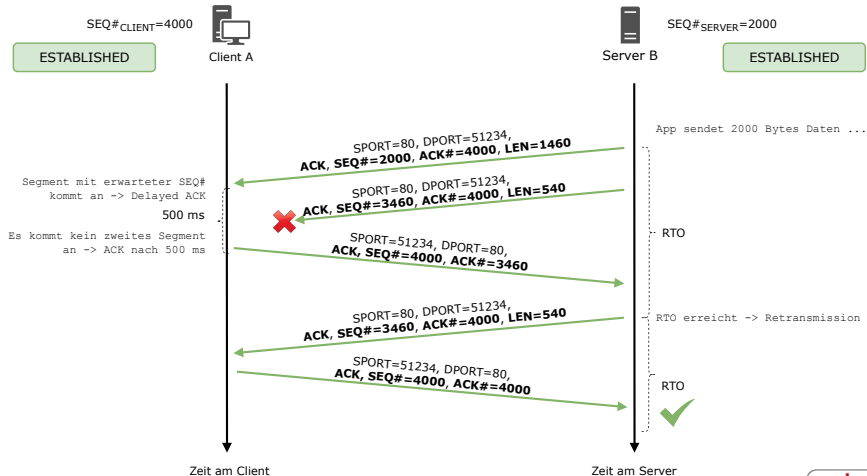
TCP Zuverlässigkeit (Reliability) IV

Beispiel 1 - Alles normal



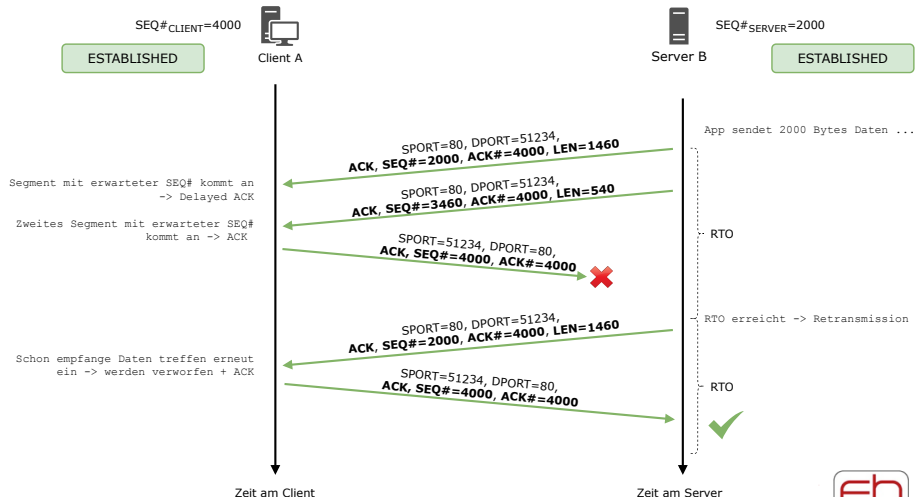
TCP Zuverlässigkeit (Reliability) V

Beispiel 2 - Data lost



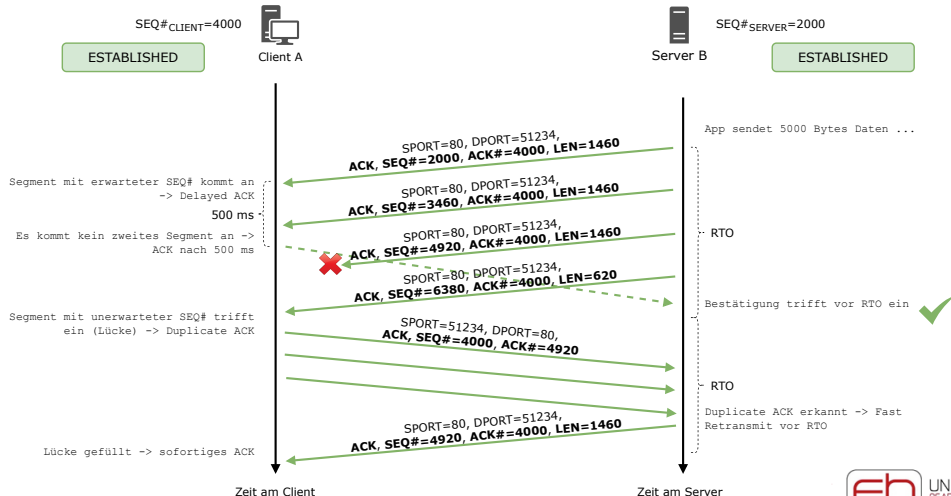
TCP Zuverlässigkeit (Reliability) VI

Beispiel 3 - ACK lost



TCP Zuverlässigkeit (Reliability) VII

Beispiel 4 - Duplicate ACK



- [1] J. Kurose und K. Ross, Computer Networking - A Top-Down Approach, 7. Aufl. Pearson Education, 2017, S. 852, ISBN: 1-292-15359-8.
- [2] J. Postel, „Internet Protocol,“ Sep. 1981, S. 45. DOI: [10.17487/rfc0791](https://doi.org/10.17487/rfc0791). Adresse: <https://www.rfc-editor.org/info/rfc0791>.
- [3] M. Sargent, J. Chu, D. V. Paxson und M. Allman, Computing TCP's Retransmission Timer, RFC 6298, 2011. DOI: [10.17487/RFC6298](https://doi.org/10.17487/RFC6298). Adresse: <https://rfc-editor.org/rfc/rfc6298.txt>.
- [4] E. Blanton, D. V. Paxson und M. Allman, TCP Congestion Control, RFC 5681, 2009. DOI: [10.17487/RFC5681](https://doi.org/10.17487/RFC5681). Adresse: <https://rfc-editor.org/rfc/rfc5681.txt>.

- [5] R. T. Braden, Requirements for Internet Hosts - Communication Layers, RFC 1122, 1989. DOI: [10.17487/RFC1122](https://doi.org/10.17487/RFC1122). Adresse: <https://rfc-editor.org/rfc/rfc1122.txt>.
- [6] E. Allman, „The Robustness Principle Reconsidered,“ Commun. ACM, Jg. 54, Nr. 8, S. 40–45, Aug. 2011, ISSN: 0001-0782. DOI: [10.1145/1978542.1978557](https://doi.org/10.1145/1978542.1978557). Adresse: <https://doi.org/10.1145/1978542.1978557>.