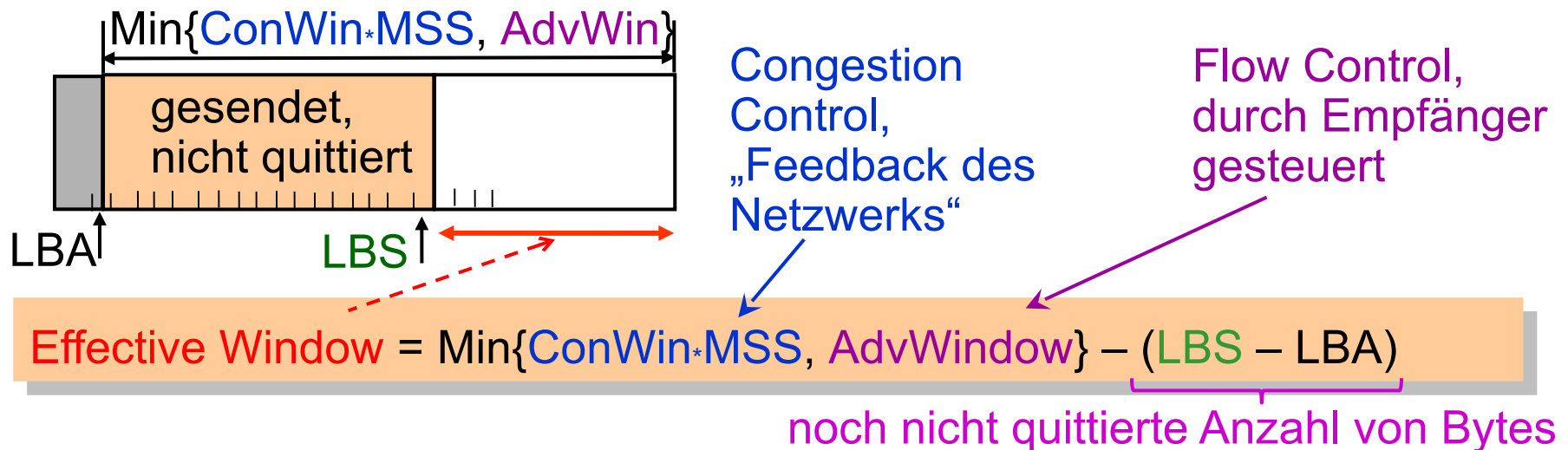


Slow Start (Wiederholung)

- Der Sender verwendet ein **Congestion Window (ConWin)**
- Der Sender darf nie mehr als **ConWin * MSS Bytes** ohne Quittierung senden
- ConWin wird in der Einheit Maximum Segment Size (MSS) angegeben
 - Default nach RFC 1122: MSS = 536 Byte
- Congestion Control und Flow Control sind im Sender gekoppelt
 - Maximal sendbar ist das **Effective Window = Minimum aus ConWin * MSS und AdvWindow**



Prinzip

- Mit einem kleinen Wert des Congestion Windows beginnend, wird dieses exponentiell erhöht, solange keine Überlast auftritt
- Slow Start wird bei Überschreitung eines Grenzwertes beendet
 - **Threshold (Thres):** Grenzwert in Byte (default: 64 kByte)

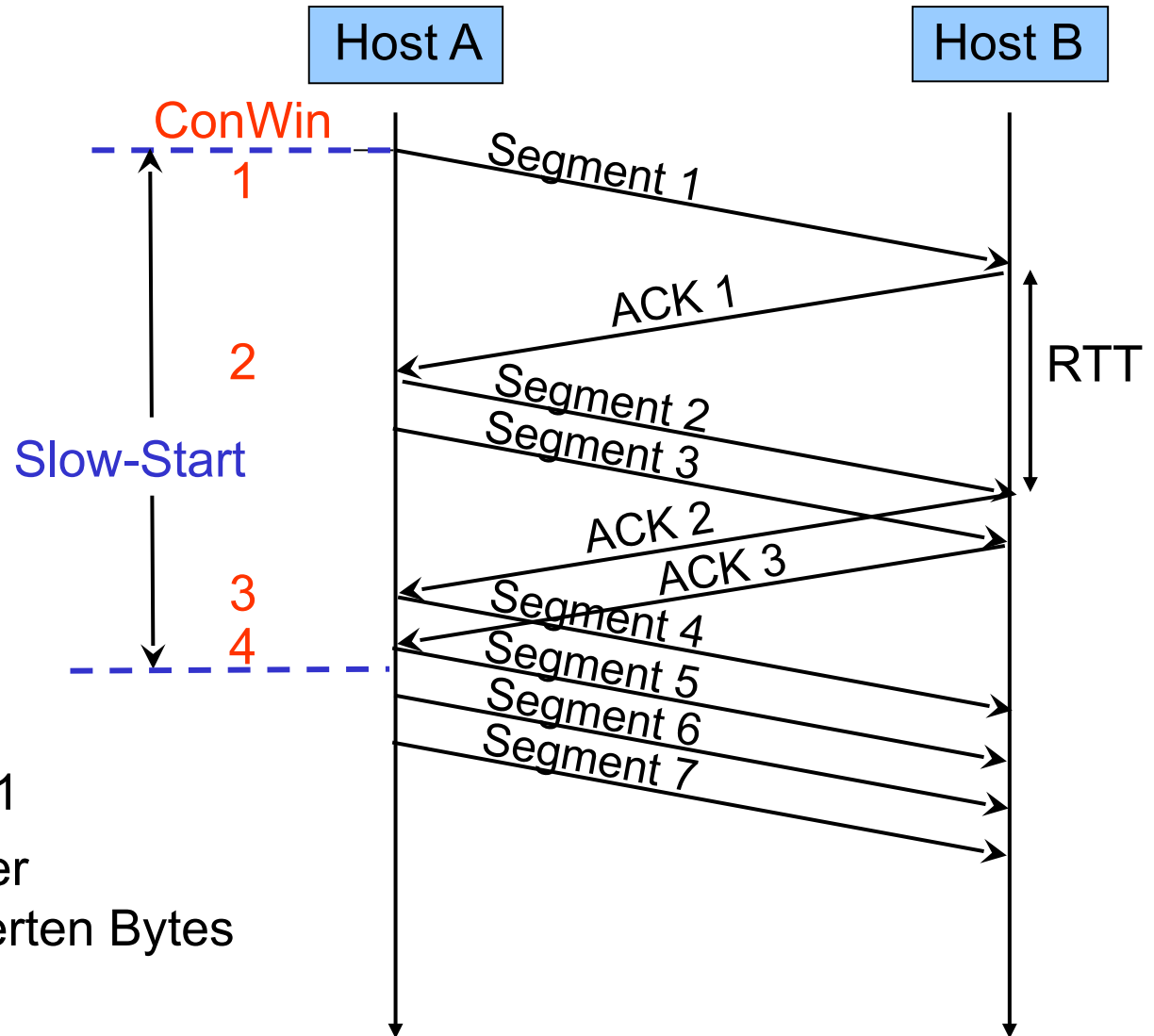
Verfahren

- **Start:** Beim Verbindungsstart oder nach einem Segmentverlust
 $\text{ConWin} = 1 \text{ MSS}$
- **Eingehendes ACK:**
 $\text{ConWin} = \text{ConWin} + 1 \text{ MSS}$
- **Stop:** Slow Start endet, falls ConWindow den Grenzwert überschreitet
 $\text{ConWin} > \text{Thres} \rightarrow \text{stop Slow Start}$

Initialisierung:

- $\text{AdvWin} > 50 \text{ MSS}$
- $\text{Tres} = 3 \text{ MSS}$
- kein delayed Ack

- Jedes ACK erhöht das ConWin um 1
- unabhängig von der Anzahl der quittierten Bytes



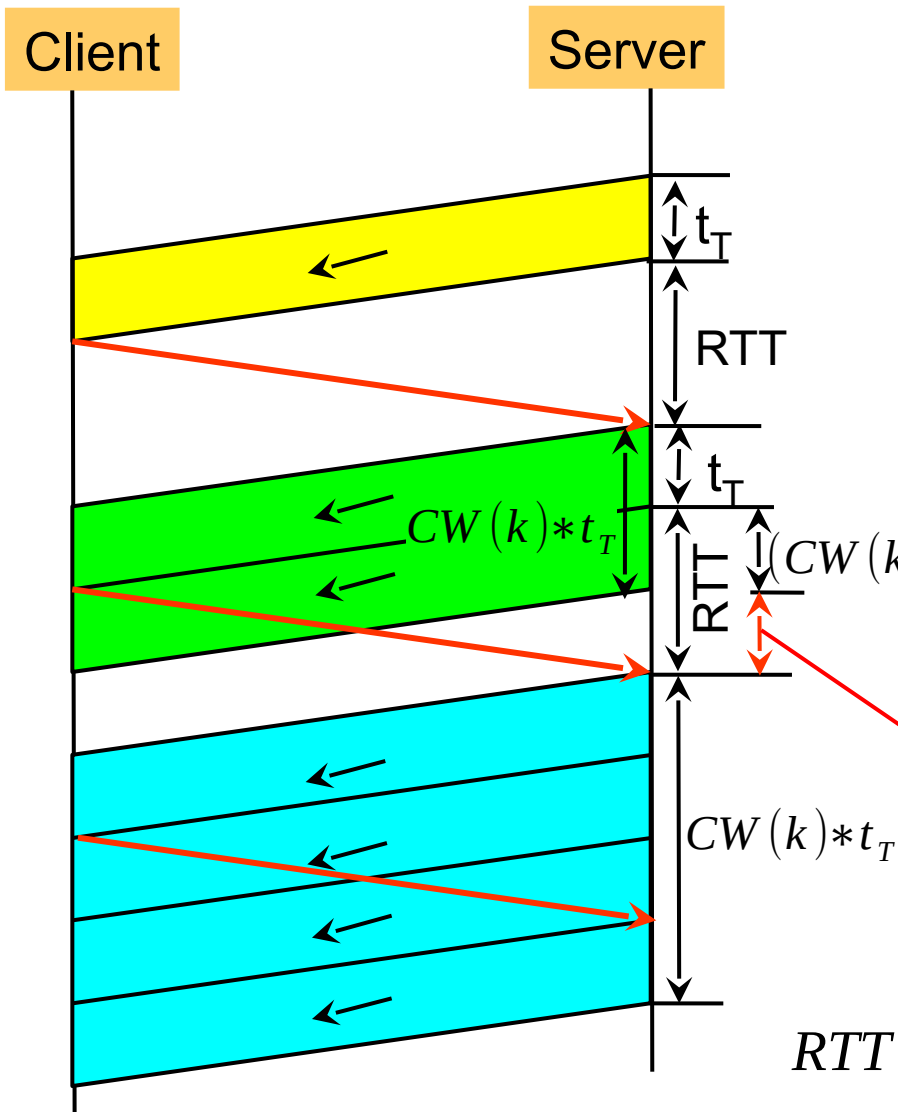
Slow Start – Analytischer Ansatz

Variablen

- $CW(t)$ = Congestion Window zum Zeitpunkt t
- $AW(t)$ = Advertised Window zum Zeitpunkt t
- D = Gesamte zu übertragende Datenmenge in bit
- $EFF(t)$ = Größe des Effective Window
- $k \in \{1, 2, 3, \dots\}$ = Anzahl von Fensterzyklen für Übertragung
- S = Anzahl der zur Übertragung benötigten Segmente D/MSS
- t_T = Zeiteinheit für Slots auf Graphen ($= MSS/v_B$)
- Q = Anzahl an Fensterzyklen während derer auf ACK gewartet wird
- $thres$ = Grenze für Übergang von Slow Start zu Congestion Avoidance
- RTT = Round Trip Time, v_B = Leitungsgeschwindigkeit
- K = Gesamtzahl der Zyklen, MSS = Maximum Segment Size

Initialwerte

- $CW(0) = 1$



- Im k -ten Fenster können $CW(k) = 2^{k-1}$ Segmente gesendet werden
- Das ACK für das erste Segment des Fensters erreicht den Server nach der Dauer $RTT + t_T$
- Bevor das nächste Fenster gesendet wird, muss der Sender auf dieses ACK warten, falls $CW(k) * t_T < RTT + t_T$
- Die Wartedauer bis zum Empfang des ACKS beträgt:

$$RTT - (CW(k) - 1) * t_T = RTT - (2^{k-1} - 1) * t_T$$

- Sind genügend Daten zu senden, so muss der Server nur während der ersten Q Fensterzyklen auf ein ACK warten, wobei

$$Q = \max \{k : CW(k) * t_T \leq RTT + t_T\} = \max \{k : 2^{k-1} * t_T \leq RTT + t_T\}$$

$$= \max \left\{ k : 2^{k-1} \leq 1 + \frac{RTT}{t_T} \right\} = \max \left\{ k : k \leq 1 + \log_2 \left(1 + \frac{RTT}{t_T} \right) \right\}$$

Lies: maximales k
das die Bedingung
 $k \leq 1 + [\dots]$ erfüllt

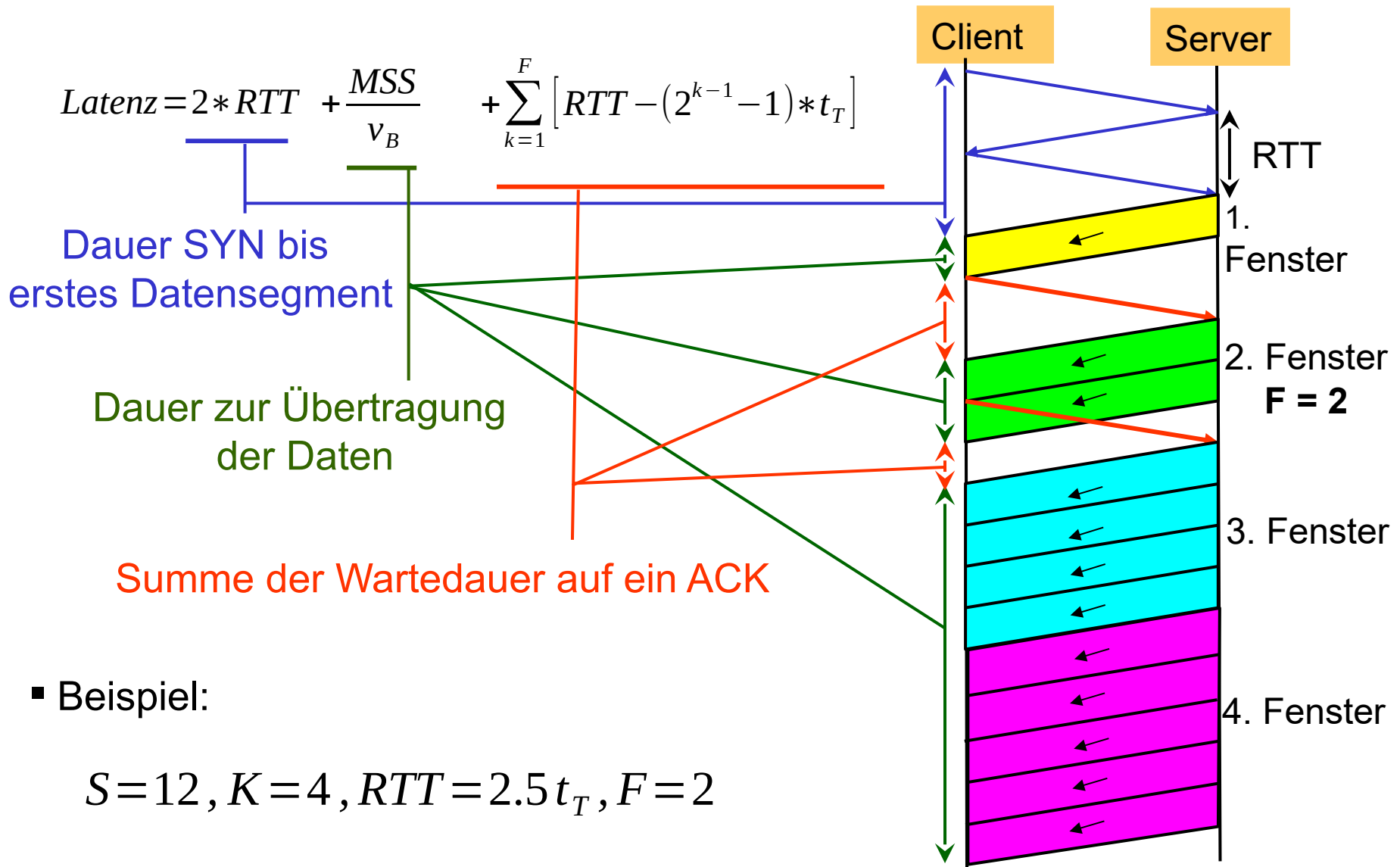
- Es werden insgesamt K Fenster-Zyklen an Daten übertragen.
 - Der Server muss nur während der ersten

$$F = \min \{Q, K - 1\}$$

Fenster-Zyklen warten

- Beispiel

$$S = 12, K = 4, RTT = 2.5 t_T \rightarrow Q = \max \left\{ k : k \leq 1 + \overbrace{\log_2(3.5)}^{1.8} \right\}$$
$$F = 2$$



▪ Beispiel:

$$S = 12, K = 4, RTT = 2.5 t_T, F = 2$$

- Es gilt die geschlossene Form die wie folgt vereinfacht werden kann:

$$\begin{aligned} \text{Latenz} &= 2 * RTT + \frac{D}{v_B} + \sum_{k=1}^F \left[RTT - (2^{k-1} - 1) * t_T \right] \\ &= 2 * RTT + \frac{D}{v_B} + F * (t_T + RTT) - t_T \sum_{k=1}^F 2^{k-1} \\ &= 2 * RTT + \frac{D}{v_B} + F * (t_T + RTT) - t_T (2^F - 1) \end{aligned}$$

- Beispiel:

$$S = 12, D = S * MSS = 12 \text{ MSS}, K = 4, F = 2$$

$$t_T = \frac{MSS}{v_B}, RTT = 2.5 t_T = 2.5 \frac{MSS}{v_B}$$

$$\text{Latenz} = \frac{MSS}{v_B} (5 + 12 + 2(1 + 2.5) - 3) = 21 \frac{MSS}{v_B} = 1.75 \frac{D}{v_B}$$

- Minimale Latenz: Congestion Window schränkt Senderate nicht ein

$$\text{minimale Latenz} = 2 * RTT + \frac{D}{V_B}$$

- Beispiele: mit MSS = 536 Bytes, SIS=Slow Start

	RTT= 100 ms, D = 100 Kbyte, K = 8			RTT= 100 ms, D = 5 Kbyte, K = 4			RTT= 1 s, D = 100 Kbyte, K = 8		
V_B [Mbit/s]	F	minimale Latenz	Latenz mit SIS	F	Minimale Latenz	Latenz mit SIS	F	Minimale Latenz	Latenz mit SIS
0,1	2	8,4 s	8,5 s	2	0,6 s	0,76 s	5	10,2 s	14,1 s
1	5	1 s	1,5 s	3	0,24 s	0,52 s	7	2,8 s	9,3 s
10	7	0,28 s	0,93 s	3	0,2 s	0,5 s	7	2,1 s	9,0 s

- Negativer Einfluss von Slow Start bei sehr großer RTT

Congestion Avoidance

Beobachtung

- Nicht jeder Verlust eines Segmentes beruht auf einer Überlastsituation
- Der Empfang von „duplicate ACKs“ deutet auf einzelne Segmentverluste hin

Forderung

- Vermeidung langer Timeout-Intervalle bei einzelnen Segmentfehlern
 - Lösung: Fast-Retransmit; Neuübertragung einzelner Segmente
- Liegt keine Überlast vor kann die Slow Start Phase entfallen
 - Lösung: Fast-Recovery; zurücksetzen des ConWin auf Thres

Verfahren: **Additive Increase - Multiplicative Decrease (AIMD)**

ConWin > Thres (=Zustand CA):

- Neues ACK vor Ablauf des Retransmission Timers

$$\mathbf{ConWin = Min\{ ConWin + 1 / ConWin, TCP_MaxWin\}}$$

- Timeout: Retransmission Timer abgelaufen → Überlast angenommen

$$\mathbf{Thres = max\{ ConWin/2, 2 MSS\}} \text{ und } \mathbf{Retransmission}$$

$$\mathbf{ConWin = 1} \text{ (} \rightarrow \text{ zurück zu Phase Slow Start)}$$

- 3 dACKs empfangen: → Einzelne Segmentfehler angenommen

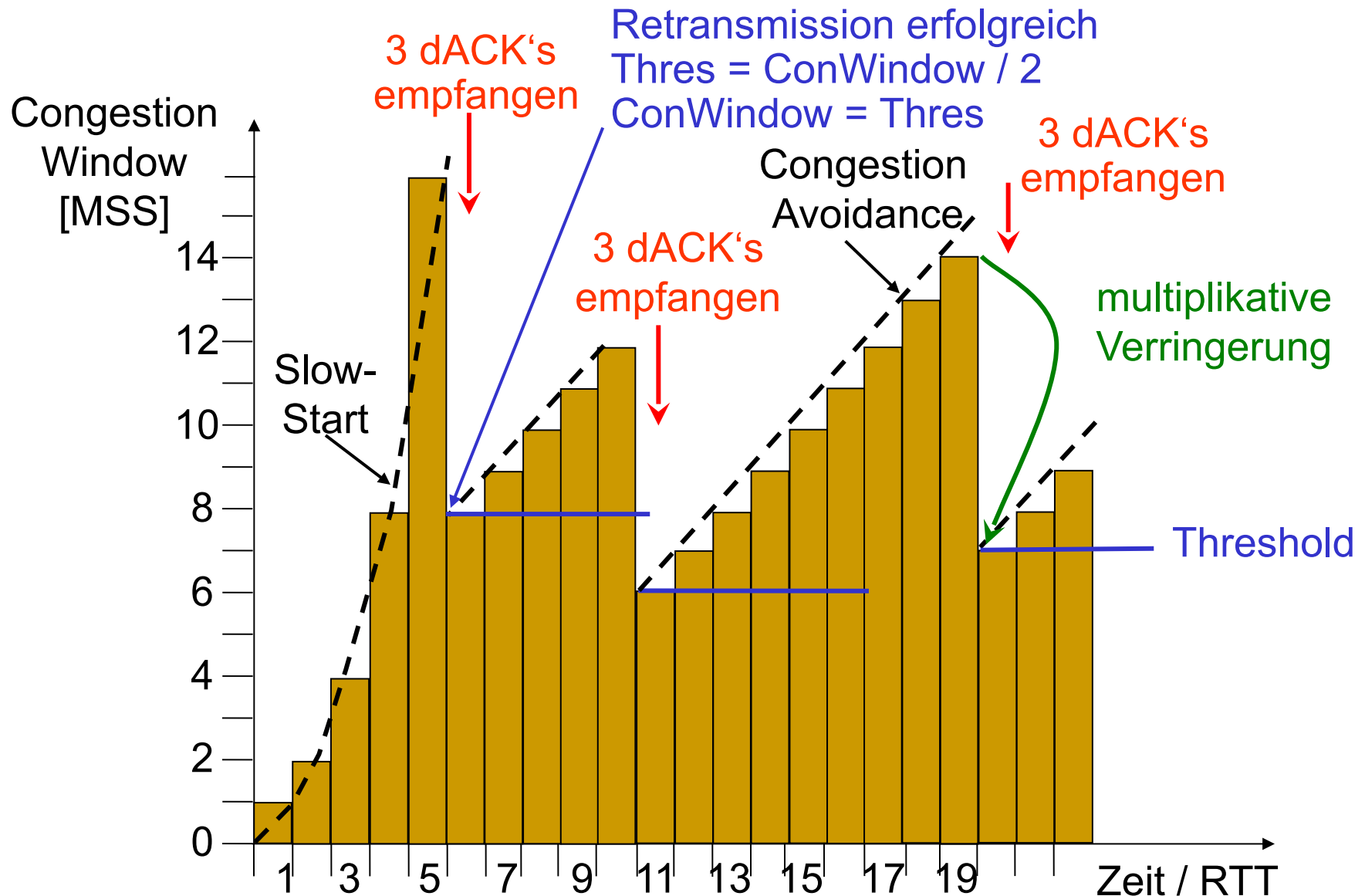
Fast Retransmit

$$\mathbf{Thres = max\{ ConWin/2, 2 MSS\}}$$

- Ack für Retransmit (oder neues ACK): → Fast Recovery

$$\mathbf{ConWin = Thres}$$

Zustand: Congestion Avoidance



Zustandstabelle (TCP RENO)

GIT (WS 2023/24)

Zustand	Ereignis	TCP Sender Aktion	Kommentar
Slow Start (SS)	neues ACK für nicht quittierte Daten	$\text{ConW} = \text{ConW} + \text{MSS}$, If ($\text{ConW} > \text{Threshold}$) → Zustand CA	Das ConW wird je RTT verdoppelt
Con. Avoid. (CA)	neues ACK	$\text{ConW} = \text{ConW} + 1 / \text{ConW}$	Additive Erhöhung, ConW $\approx 1 \text{ MSS}$ je RTT erhöhen
SS oder CA	Duplicate ACK (dACK)	erhöhe dACK Zähler für jedes dupliziert quittierte Segment	ConW und Threshold bleiben erhalten
	3 dACKs	Überlast → Zustand FR	Einzelsegmentverlust
SS / CA / FR	Timeout	$\text{Thr}^* = \max\{\text{ConW}/2, 2 \text{ MSS}\}$ $\text{ConW} = 1$ → Zustand SS	Überlast des Netzes
Fast Recovery (FR)	drittes dACK	$\text{Thresh} = \text{Thr}^* = \max\{\text{ConW}/2, 2\}$ Fast-Retransmission $\text{ConW} = \text{Thresh} + 3 \text{ MSS}$	Fast Retransmit und verringertem neuem Threshold.
FR	dACK	$\text{ConW} = \text{ConW} + \text{MSS}$ nächstes Segment senden	DACK → ein Segment hat den Empfänger erreicht
FR	ACK für Re-transmission	$\text{ConW} = \text{Thr}^*$ → Zustand CA	wiederholtes Segment quittiert, FR beendet

Congestion Avoidance – Analytischer Ansatz

Mathematische Betrachtung

$$CW_0 = Thres$$

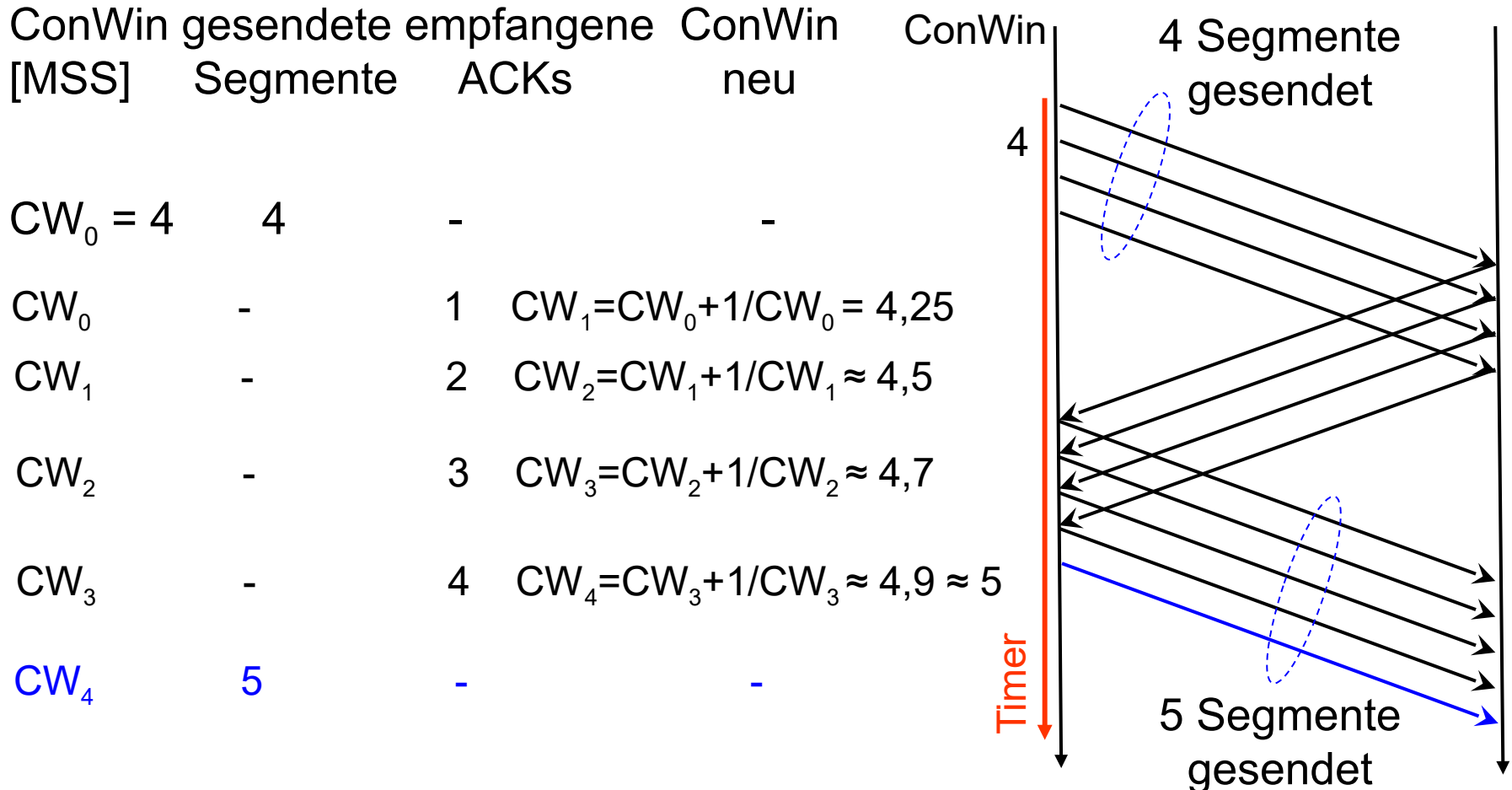
- If ACK received:

$$CW(t) = CW(t-1) + \frac{1}{CW(t-1)}$$

- CW wächst um $1/CW$ mit jedem eintreffenden ACK
- Hinweis: Sie dürfen diese Werte in der Klausur abrunden und das CW inkrementieren sobald „alle“ ACKs angekommen sind (=Nachkommastellen ignorieren)
- Beispiel: $CW = 4 \rightarrow 4$ Segmente gesendet, nach Empfang von Acknowledgement 4 $\rightarrow CW = 5$

Beispiel:

- Annahme: TCP sendet nur Segmente von MSS Bytes



Netzwerkunterstützte Congestion Control

Problem:

- Ende-zu-Ende-Überlastkontrolle schließt durch Paketverluste auf Überlast
- Paketwiederholungen sind für verzögerungszeitsensitive Anwendungen ungünstig

Lösung:

Aktives Warteschlangenmanagement:

- Verhindern, dass eine Warteschlange in den Überlastbereich gelangt
- **Random Early Detection**
 - Netzwerkgestützte Überlastcontroller:
 - Netzwerk teilt Endsystemen eine bestehende Überlast mit
 - Explicit Congestion Notification

Prinzip

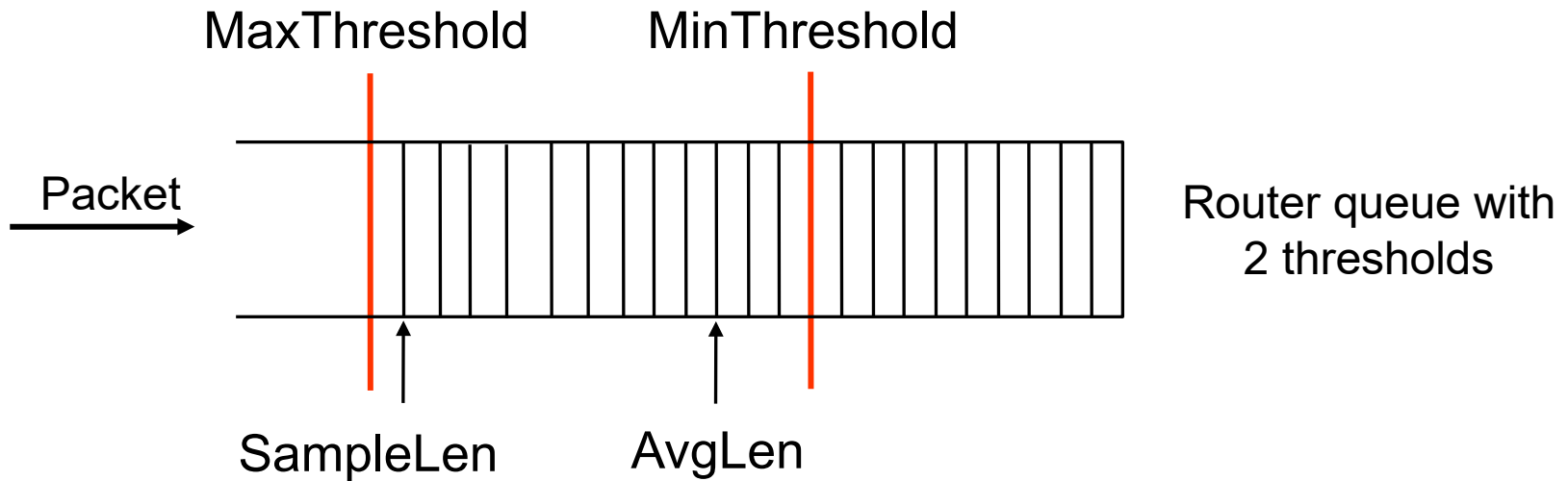
- Implizite Benachrichtigung der Übertragungspartner
 - Verwerfen von Paketen durch den Router (führt zu Timeouts)
→ Halbierung des Congestion Windows, Reduktion der Senderate
- Explizite Benachrichtigung durch Markieren von Paketen möglich

Frühzeitiges Verwerfen

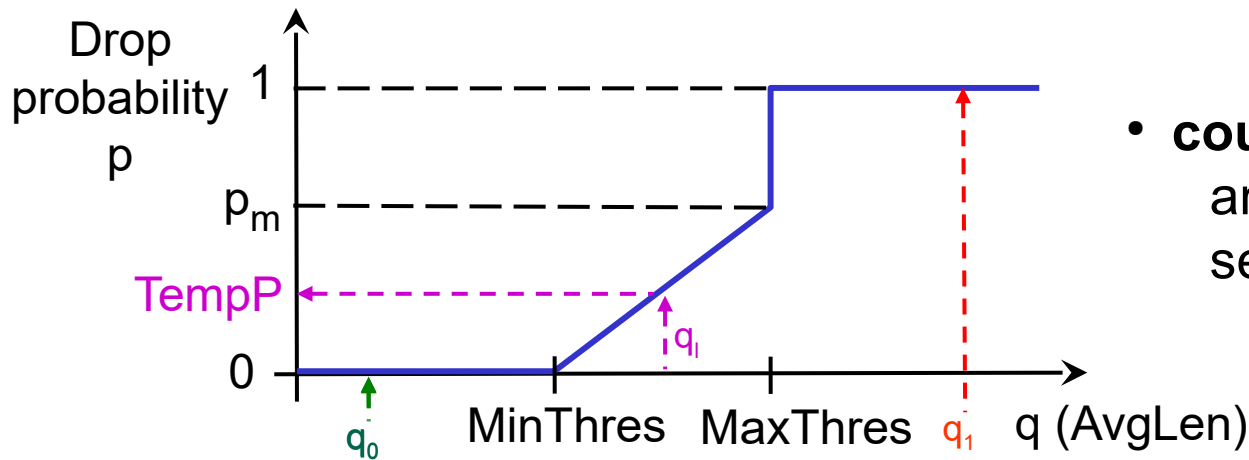
- Anstatt darauf zu warten, dass die Queue vollständig gefüllt ist wird jedes Paket mit einer **drop probability** verworfen sobald die Queue eine definierte Länge überschreitet

Eigenschaften

- Verhindert Synchronisation der Sendepartner
- Reduziert die mittlere Queuelänge im Router



- SampleLen ist die Länge der Queue zum Ankunftszeitpunkt eines Pakets
- Average queue length (AvgLen)
 - $AvgLen = (1 - Weight) * AvgLen + Weight * SampleLen$
 - $0 < Weight < 1$ (usually 0.002)
- AvgLen wird meist für jedes ankommende Paket errechnet



- **count** = Anzahl der angekommenen Pakete seit letztem Drop

Average queue length q	Drop probability p
$0 \leq q \leq \text{MinThres}$	$p = 0$
$\text{MinThres} < q < \text{MaxThres}$	p steigt linear bis p_m an
$q \geq \text{MaxThres}$	$p = 1$ (all packets are dropped)

• Berechnung der Wahrscheinlichkeit P

$$\text{TempP} = p_m * (q - \text{MinThres}) / (\text{MaxThres} - \text{MinThres})$$

$$P = \text{TempP} / (1 - \text{count} * \text{TempP})$$

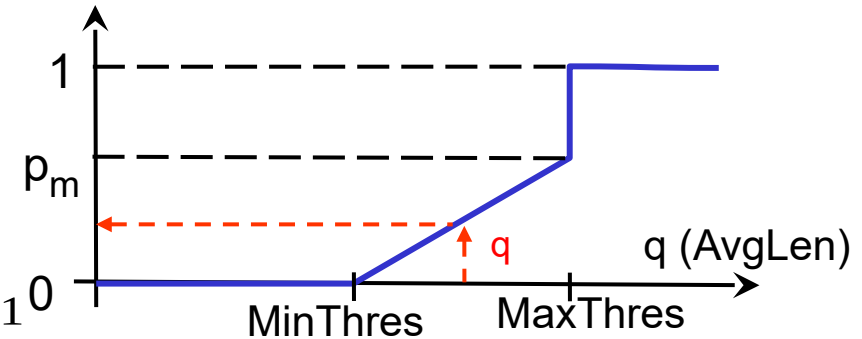
$$p_m = 0,02$$

$$q = \frac{MaxThres + MinThres}{2} \approx const, da \text{ Mittelwert}$$

$$TempP = p_m \frac{q - MinThres}{MaxThres - MinThres} = 0,5 \quad p_m = 0,01$$

$$P = \frac{TempP}{1 - count \cdot TempP} = \frac{0,01}{1 - 0,01 \cdot count}$$

$$count = 0$$



wäre $P = TempP = 0,01$,
würde im Mittel jedes
100'tes Paket gelöscht

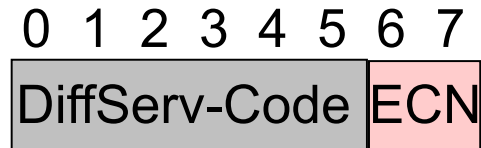
Anzahl Pakete	count	P	Aktion
1	0	0,01	nicht löschen
50	49	1/49	nicht löschen
51	50	0,02	nicht löschen
99	98	0,5	nicht löschen
100	99	1,0	löschen

- count: Number of arriving packets since the last packet has been dropped
- Aktion: drop packet or accept packet

- Drop-Wahrscheinlichkeit für einen Flow ist grob proportional zum Anteil der Bandbreite die der Flow bekommt
- **MaxP** (p_m) ist typischerweise auf 0.02 gesetzt, wenn also die durchschnittliche Queuelänge zwischen beiden Thresholds liegt, wird ca. jedes 50. Paket verworfen
- Wenn der Verkehr „bursty“ ist, sollte **MinThreshold** ausreichend groß sein um ausreichend hohe Verbindungsauslastung zu gewährleisten
- Die Differenz zwischen zwei Thresholds sollte größer sein als der typische Zuwachs der durchschnittlichen Queuelänge während einer RTT (oft: $\text{MaxThreshold} = 2 * \text{MinThreshold}$)
- Weitere Anmerkungen:
 - RED erfordert komplexes Queue-Management
 - Ist heutzutage in quasi allen Routern implementiert

- Router im Netz zeigt Überlast im TOS-Feld des IP-Headers an
- Setzt ein Warteschlangen-Management im Router voraus (z.B. RED)
 - Signalisierung der Überlast muss vor einem Paketverlust erfolgen
 - Im Überlastfall werden die IP-Datagramme markiert
- Die ECN-Verwendung muss im IP-TOS-Feld angezeigt werden
 - ECN-Capable Transport (ECT) Bits

TOS-Feld im IP-Header:



0 0	kein ECT	Sender verwendet kein ECN
0 1	ECT(1)	Sender zeigt ECN - Verwendung an
1 0	ECT(0)	alternative ECN - Anzeige
1 1	CE	Überlastanzeige (Congestion Experiences)

Erweiterung des TCP-Headers (2 Bit des reservierten Bereiches)

- ECN-Echo-Flag (ECE)
- Congestion-Window-Reduced Flag (CWR)

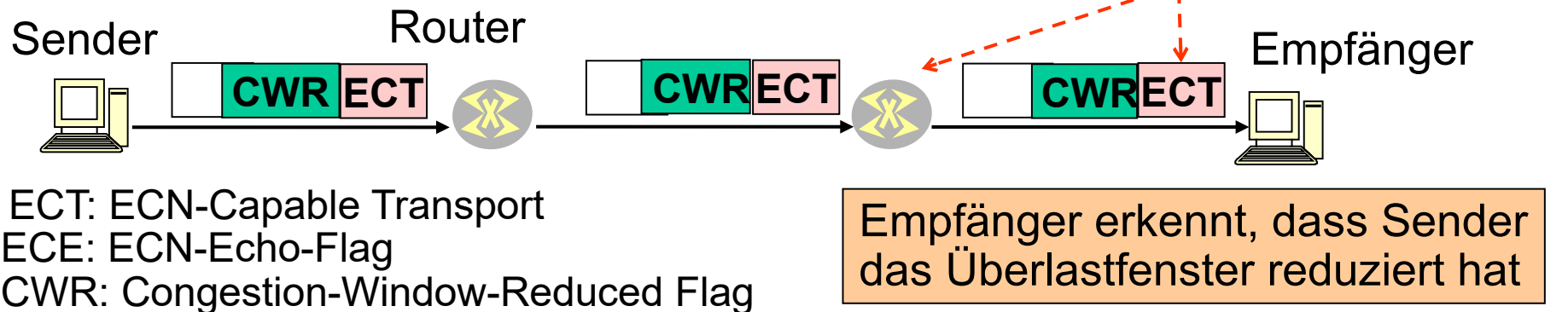
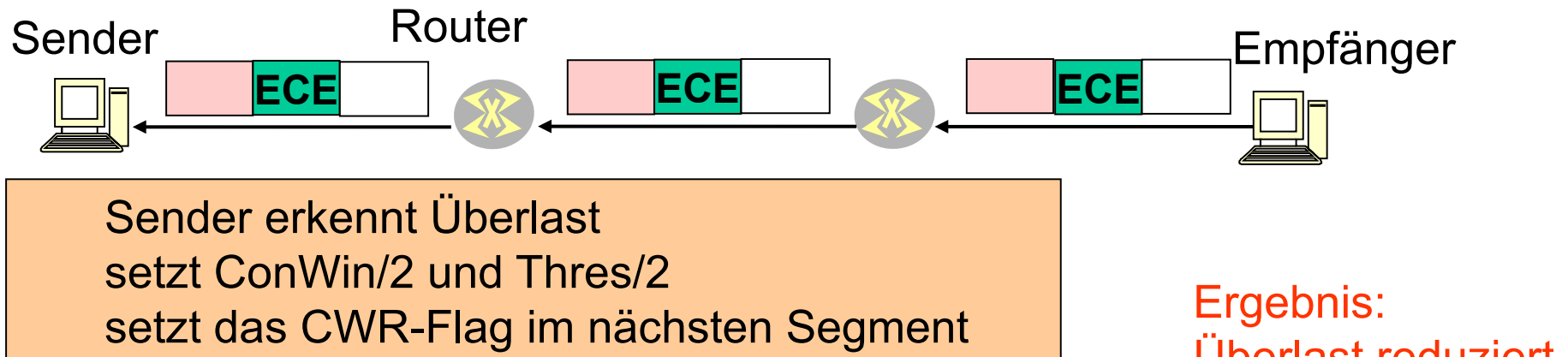
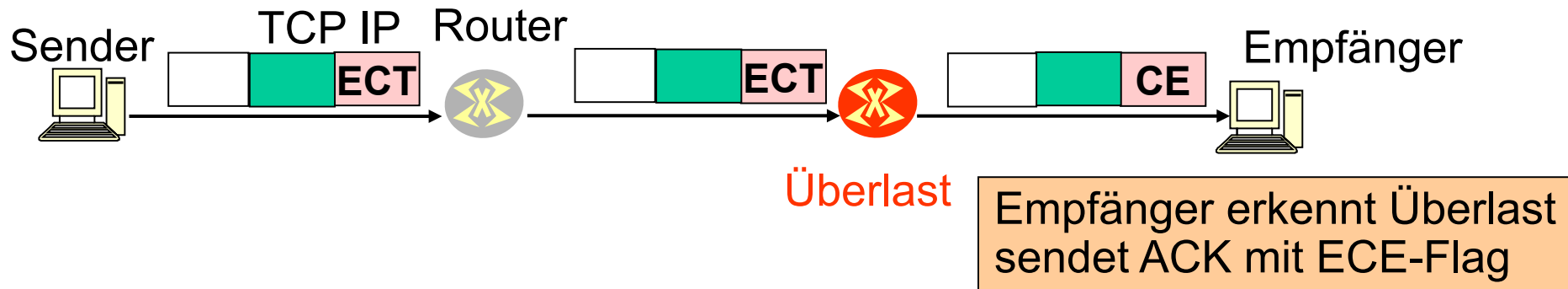
Source-Port			Destination-Port		
Sequence Number					
Acknowledgement Number					
Header-length		ECW C E R	Flags	Window Size	
TCP Checksum			Urgent Pointer		

Verbindungsaufbau

- ECN-Verwendung wird durch SYN+ECE+CWR signalisiert

Congestion Control:

- TCP-Sender:
 - setzt ECT(0) oder ECT(1) im ECN-Feld des IP-Headers
- Router
 - setzt im Überlastfall ECN-Feld des IP-Headers auf CE
- TCP Empfänger
 - Erkennt Überlast und sendet die folgenden ACKs mit gesetztem ECE-Flag, bis er ein Sender-Segment mit gesetztem CWR-Flag erhält
- TCP-Sender empfängt ACK mit ECE-Flag
 - halbiert das Überlastfenster und den Grenzwert (Thres)
 - setzt das CWR-Flag im nächsten Segment



ECT: ECN-Capable Transport

ECE: ECN-Echo-Flag

CWR: Congestion-Window-Reduced Flag