

Számítógépes Hálózatok

9. Előadás: ++Szállítói réteg

Based on slides from **Zoltán Ács ELTE** and D. Choffnes Northeastern U., Philippa Gill from StonyBrook University , Revised Spring 2016 by S. Laki

További protokollok

Internet Control Message Protocol

3

FELADATA

- Váratlan események jelentése

HASZNÁLAT

- Többféle *ICMP*-üzenetet definiáltak:
 - ▣ Elérhetetlen cél;
 - ▣ Időtúllépés;
 - ▣ Paraméter probléma;
 - ▣ Forráslefojtás;
 - ▣ Visszhang kérés;
 - ▣ Visszhang válasz;
 - ▣ ...

Internet Control Message Protocol

4

- *Elérhetetlen cél* esetén a csomag kézbesítése sikertelen volt.
 - ▣ **Esemény lehetséges oka:** Egy nem darabolható csomag továbbításának útvonalán egy „kis csomagos hálózat” van.
- *Időtúllépés* esetén az IP csomag élettartam mezője elérte a 0-át.
 - ▣ **Esemény lehetséges oka:** Torlódás miatt hurok alakult ki vagy a számláló értéke túl alacsony volt.
- *Paraméter probléma* esetén a fejrészben érvénytelen mezőt észleltünk.
 - ▣ **Esemény lehetséges oka:** Egy az útvonalon szereplő router vagy a hoszt IP szoftverének hibáját jelezheti.

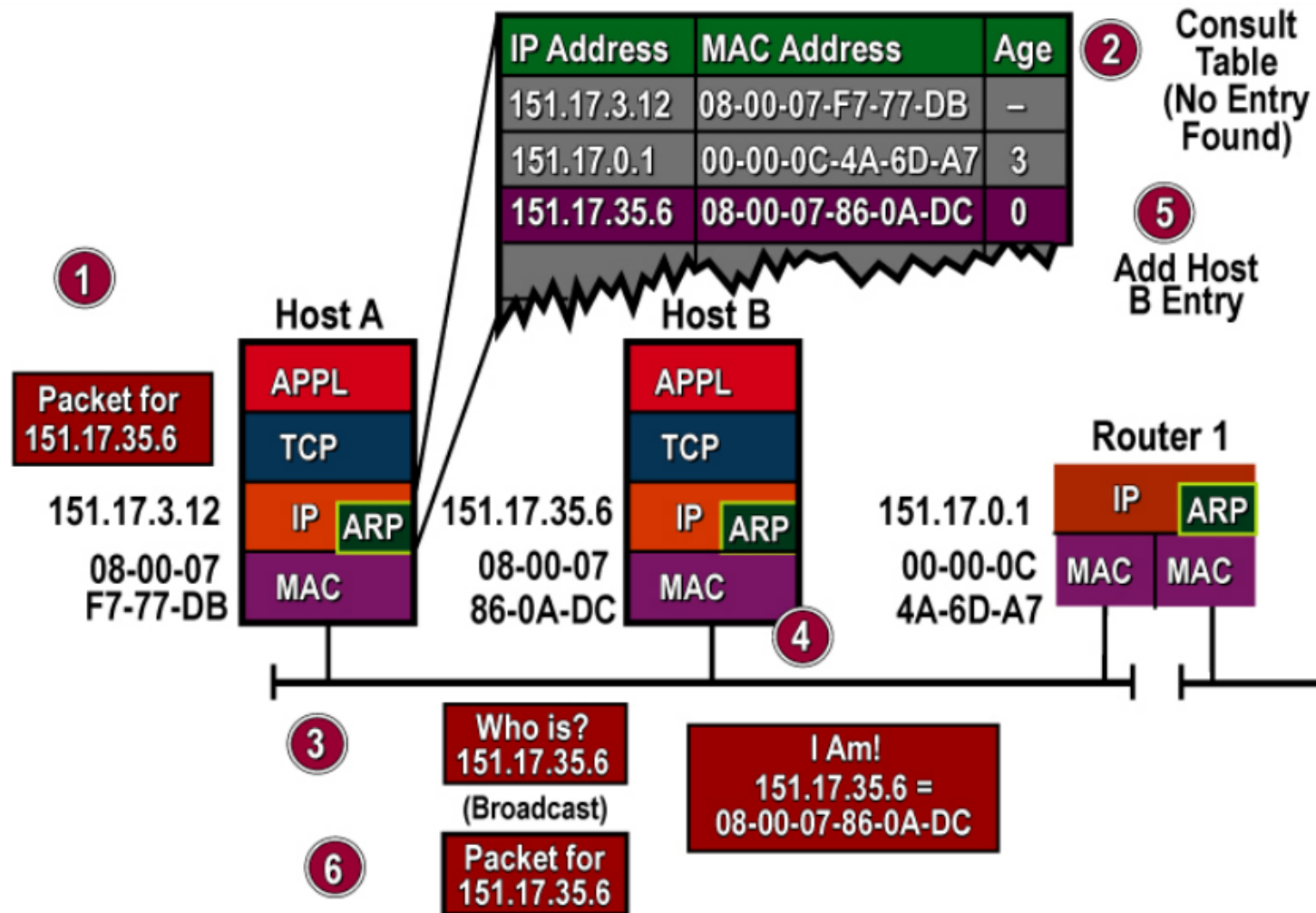
Internet Control Message Protocol

5

- Forráslefojtás esetén lefojtó csomagot küldünk.
 - ▣ **Esemény hatása:** A fogadó állomásnak a forgalmazását lassítania kellett.
- Visszhang kérés esetén egy hálózati állomás jelenlétét lehet ellenőrizni.
 - ▣ **Esemény hatása:** A fogadónak vissza kell küldeni egy visszhang választ.
- Átirányítás esetén a csomag rosszul irányítottságát jelzik.
 - **Esemény kiváltó oka:** Router észleli, hogy a csomag nem az optimális útvonall.

Address Resolution Protocol

6



Address Resolution Protocol

7

FELADATA

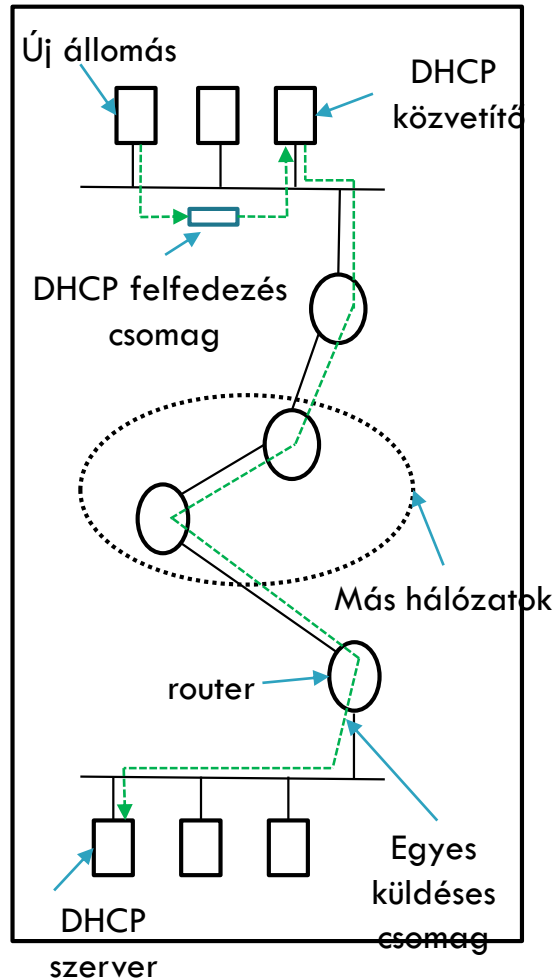
- Az IP cím megfeleltetése egy fizikai címnek.

HOZZÁRENDELÉS

- Adatszóró csomag kiküldése az *Ethernetre* „Ki-é a 192.60.34.12-es IP-cím?” kérdéssel az alhálózaton, és mindenegyes hoszt ellenőrzi, hogy övé-e a kérdéses IP-cím. Ha egyezik az IP a hoszt saját IP-jével, akkor a saját *Ethernet* címével válaszol. Erre szolgál az ARP.
- Opcionális javítási lehetőségek:
 - ▣ a fizikai cím IP hozzárendelések tárolása (*cache használata*);
 - ▣ Leképezések megváltoztathatósága (*időhatály bevezetése*);
- Mi történik távoli hálózaton lévő hoszt esetén?
 - ▣ A router is válaszoljon az ARP-re a hoszt alhálózatán. (*proxy ARP*)
 - ▣ Alapértelmezett Ethernet-cím használata az összes távoli forgalomhoz

Reverse Address Resolution Protocol

8



Reverse Address Resolution Protocol

FELADATA

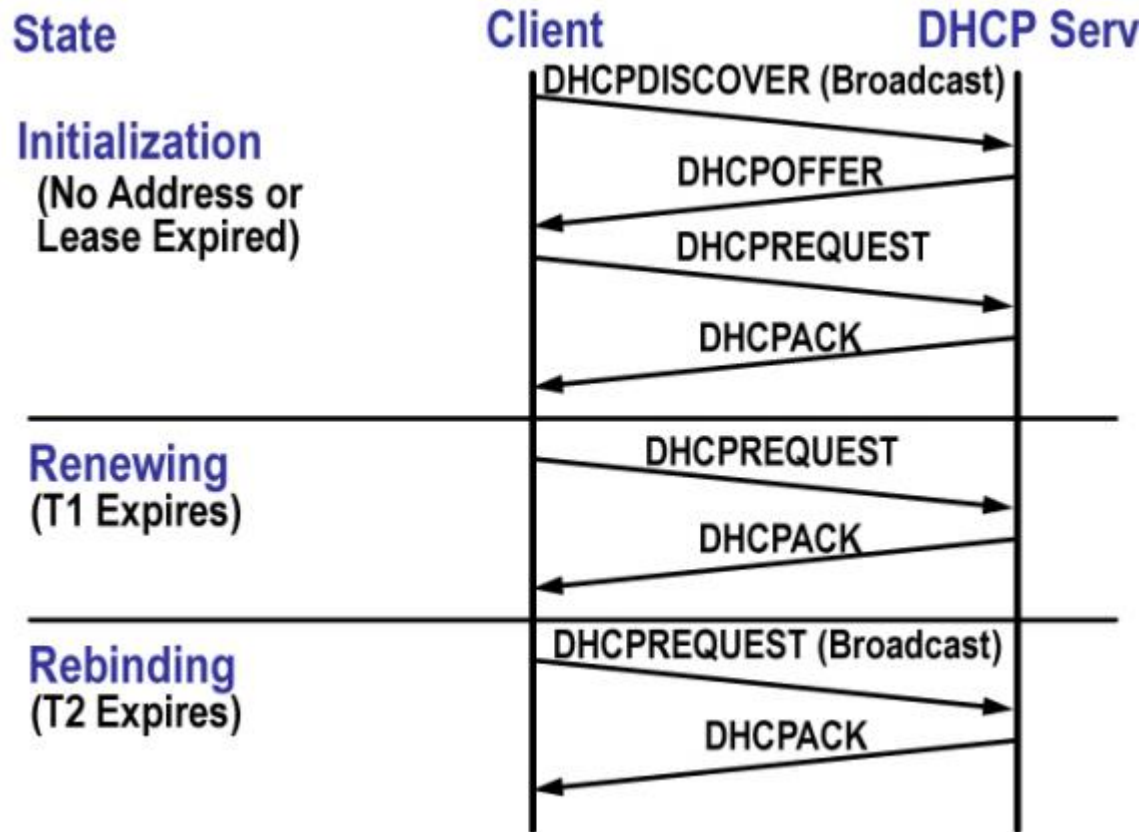
- A fizikai cím megfeleltetése egy IP címnek

HOZZÁRENDELÉS

- Az újonnan indított állomás adatszórással csomagot küld ki az *Ethernetre* „A 48-bites Ethernet-címem 14.04.05.18.01.25. Tudja valaki az IP címemet?” kérdéssel az alhálózaton. Az *RARP*-szerver pedig válaszol a megfelelő IP címmel, mikor meglátja a kérést
- Opcionális javítási lehetőségek:
 - *BOOTP* protokoll használata. UDP csomagok használata. Manuálisan kell a hozzárendelési táblázatot karbantartani. (statikus címkiosztás)
 - *DHCP* protokoll használata. Itt is külön kiszolgáló osztja ki a címeket a kérések alapján. A kiszolgáló és a kérő állomások nem kell hogy ugyanazon a *LAN*-on legyenek, ezért *LAN*-onként kell egy *DHCP relay agent*. (statikus és dinamikus címkiosztás)

DHCP: DYNAMIC HOST CONFIGURATION PROTOCOL

10



DHCP

11

- ❑ Lényegében ez már az **Alkalmazási réteg**
 - ▣ de logikailag ide tartozik

- ❑ Segítségével a hosztok automatikusan juthatnak hozzá a kommunikációjukhoz szükséges hálózati azonosítókhoz:
 - ▣ IP cím, hálózati maszk, alapértelmezett átjáró, stb.

- ❑ Eredetileg az RFC 1531 a BOOTP kiterjesztéseként definiálta. Újabb RFC-k: 1541, 2131 (aktuális)

DHCP lehetőségei

12

- IP címek osztása MAC cím alapján DHCP szerverrel
 - ▣ Szükség esetén (a DHCP szerveren előre beállított módon) egyes kliensek számára azok MAC címéhez fix IP cím rendelhető
- IP címek osztása dinamikusan
 - ▣ A DHCP szerveren beállított tartományból „érkezési sorrendben” kapják a kliensek az IP címeket
 - ▣ Elegendő annyi IP cím, ahány gép egyidejűleg működik
- Az IP címeken kívül további szükséges hálózati paraméterek is kioszthatók
 - ▣ Hálózati maszk
 - ▣ Alapértelmezett átjáró
 - ▣ Névkiszolgáló
 - ▣ Domain név
 - ▣ Hálózati rendszerbetöltéshez szerver és fájlnev

DHCP – Címek bérlése

13

- ❑ A DHCP szerver a klienseknek az IP-címeket bizonyos bérleti időtartamra (lease time) adja „bérbe”
 - ▣ Az időtartam hosszánál a szerver figyelembe veszi a kliens esetleges ilyen irányú kérését
 - ▣ Az időtartam hosszát a szerver beállításai korlátozzák
- ❑ A bérleti időtartam lejártá előtt a bérlet meghosszabbítható
- ❑ Az IP-cím explicit módon vissza is adható

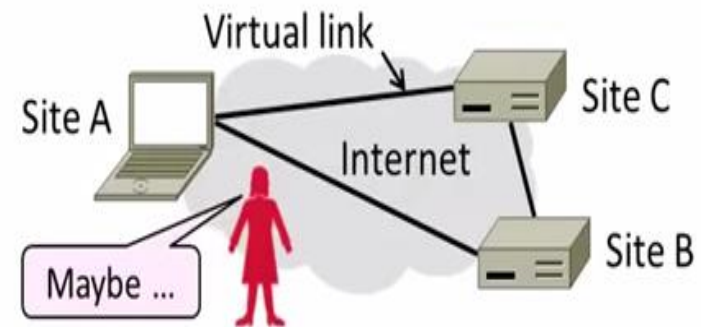
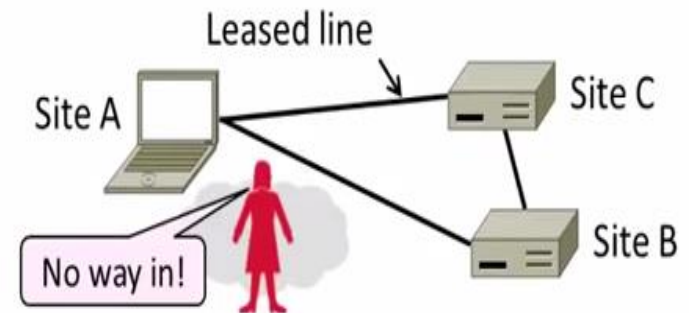
Virtuális magánhálózatok alapok

□ FŐ JELLEMZŐI

- ▣ Mint közeli hálózat fut az interneten keresztül.
- ▣ IPSEC-et használ az üzenetek titkosítására.
- Azaz informálisan megfogalmazva fizikailag távol lévő hosztok egy közös logikai egységet alkotnak.
 - ▣ Például távollévő telephelyek rendszerei.

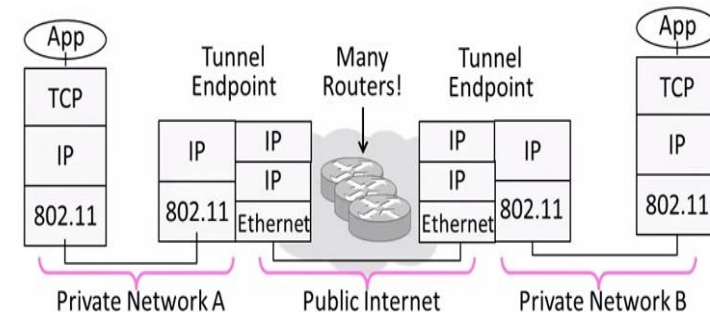
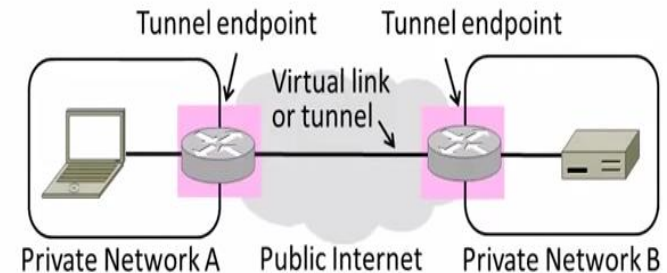
□ ALAPELV

- ▣ Bérelt vonalak helyett használjuk a publikusan hozzáférhető Internet-et.
- ▣ Így az Internettől **logikailag** elkülöníthető hálózatot kapunk. Ezek a virtuális magánhálózatok avagy VPN-ek.
- ▣ A célok közé kell felvenni a külső támadó kizárását.



Virtuális magánhálózatok alapok

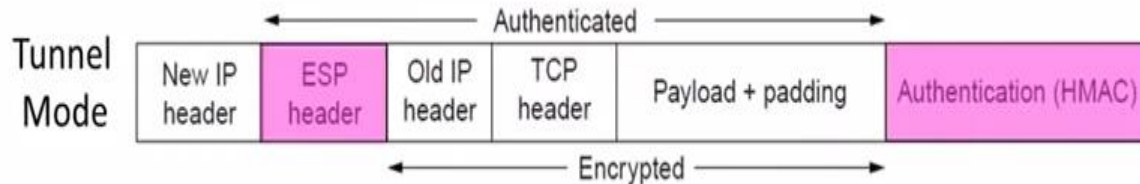
- A virtuális linkeket alagutak képzésével valósítjuk meg.
- **ALAGÚTAK**
 - ▣ Egy magánhálózaton belül a hosztok egymásnak normál módon küldhetnek üzenetet.
 - ▣ Virtuális linken a végpontok beágyazzák a csomagokat.
 - IP az IP-be mechanizmus.
- Az alagutak képzése önmagában kevés a védelemhez. Mik a hiányosságok?
 - ▣ Bizalmasság, autentikáció
 - ▣ Egy támadó olvashat, küldhet üzeneteket.
 - ▣ Válasz: Kriptográfia használata.



Virtuális magánhálózatok alapok

□ IPSEC

- Hosszú távú célja az IP réteg biztonságossá tétele. (bizalmasság, autentikáció)
- Műveletei:
 - Hoszt párok kommunikációjához kulcsokat állít be.
 - A kommunikáció kapcsolatorientáltabbá tétele.
 - Fejlécek és láblécek hozzáadása az IP csomagok védelme érdekében.
- Több módot is támogat, amelyek közül az egyik az **alagút mód**.



Szállítói réteg

17



□ Feladat:

- ▣ Adatfolyamok demultiplexálása

□ További lehetséges feladatok:

- ▣ Hosszú élettartamú kapcsolatok
- ▣ Megbízható, sorrendhelyes csomag leszállítás

- ▣ Hiba detektálás

- ▣ Folyam és torlódás vezérlés

□ Kihívások:

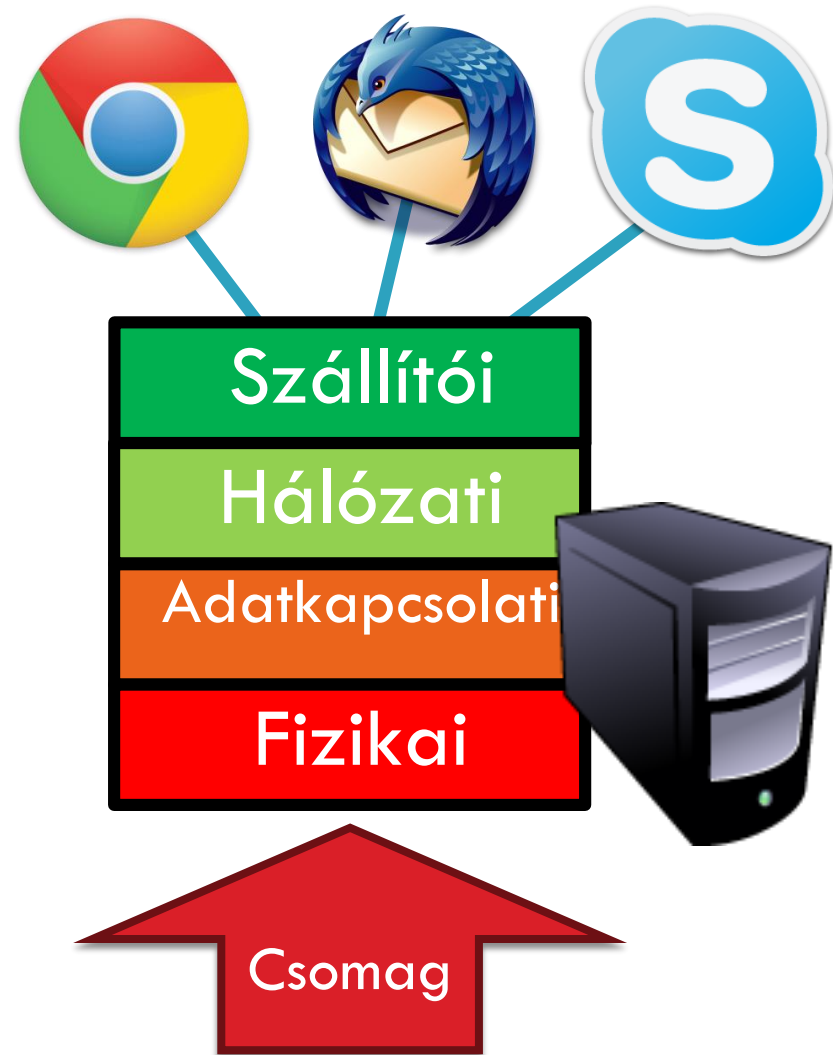
- ▣ Torlódások detektálása és kezelése
- ▣ Fairség és csatorna kihasználás közötti egyensúly

- ❑ UDP
- ❑ TCP
- ❑ Torlódás vezérlés
- ❑ TCP evolúciója
- ❑ A TCP problémái

Multiplexálás

19

- ❑ Datagram hálózat
 - ❑ Nincs áramkör kapcsolás
 - ❑ Nincs kapcsolat
- ❑ A kliensek számos alkalmazást futtathatnak egyidőben
 - ❑ Kinek szállítsuk le a csomagot?
- ❑ IP fejléc “protokoll” mezője
 - ❑ 8 bit = 256 konkurens folyam
 - ❑ Ez nem elég...
- ❑ Demultiplexálás megoldása a szállítói réteg feladata



Forralom demultiplexálása

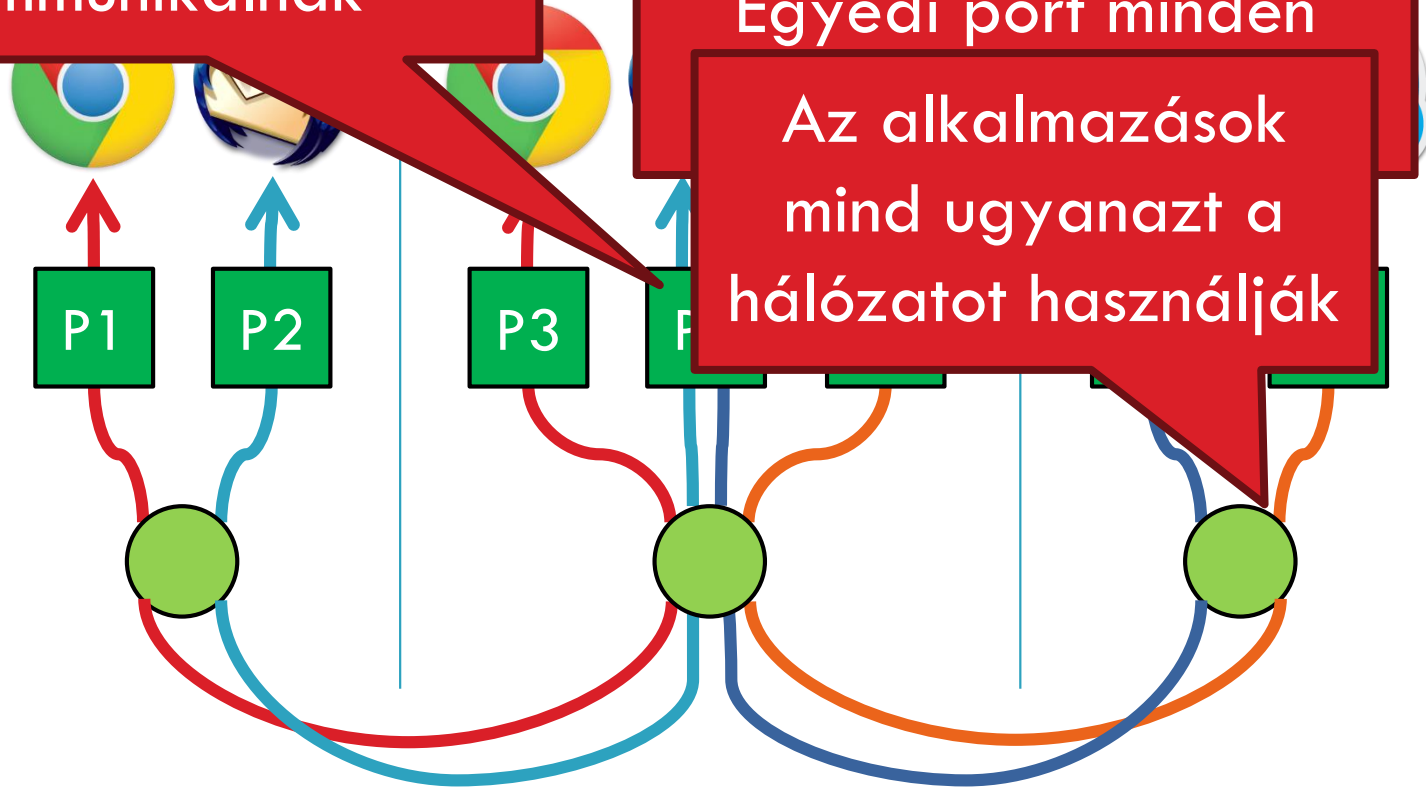
20

A szerver alkalmazások
számos klienssel
kommunikálnak

Alkalmazási

Szállítói

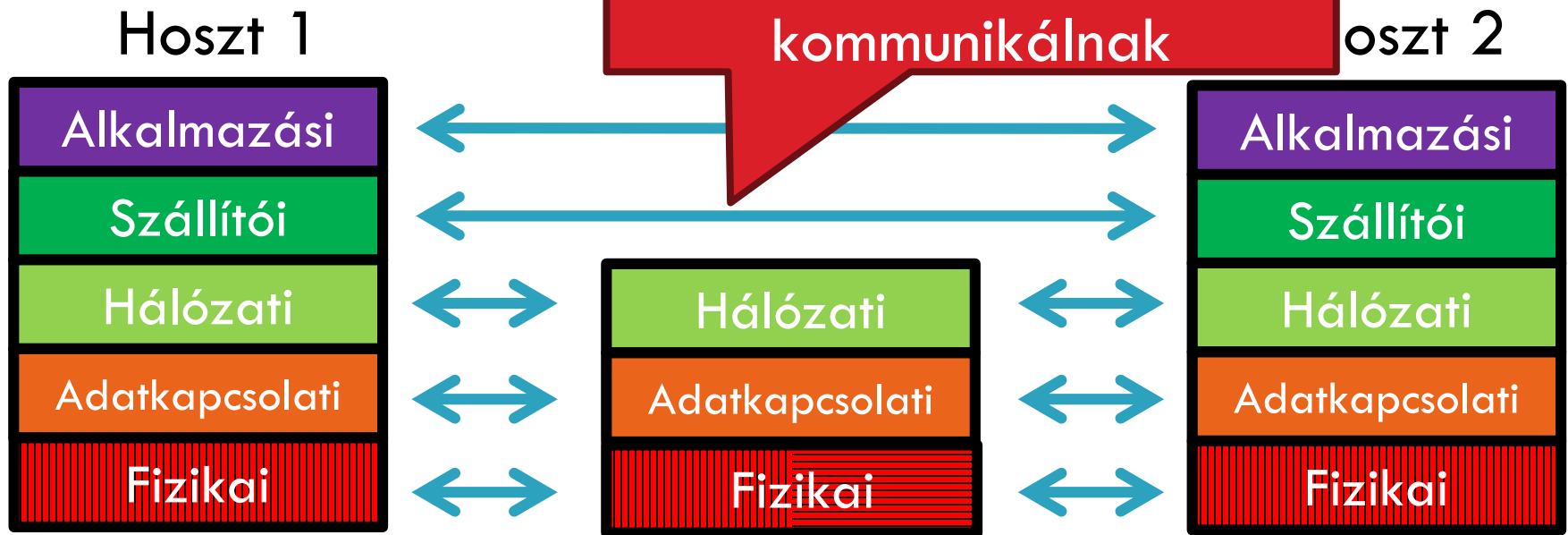
Hálózati



Végpontok azonosítása: $\langle \text{src_ip}, \text{src_port}, \text{dest_ip}, \text{dest_port}, \text{proto} \rangle$
ahol src_ip , dst_ip a forrás és cél IP cím,
 src_port , dest_port forrás és cél port, proto pedig UDP vagy TCP.

Réteg modellek

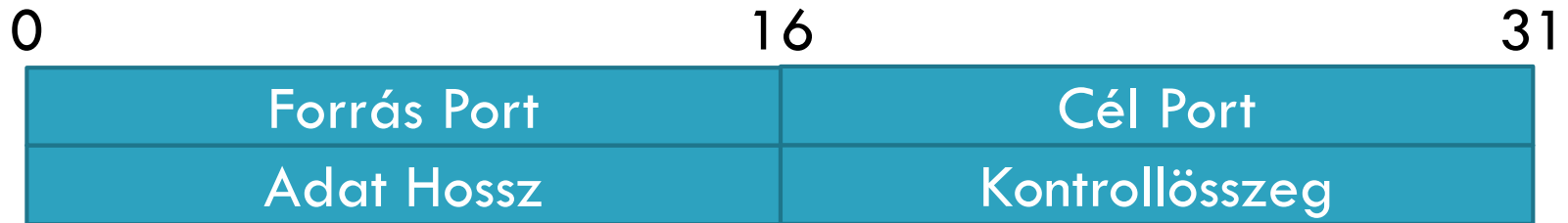
21



- A legalacsonyabb szintű végpont-végpont protokoll
 - ▣ A szállítói réteg fejlécei csak a forrás és cél végpontok olvassák
 - ▣ A routerek számára a szállítói réteg fejléce csak szállítandó adat (payload)

User Datagram Protocol (UDP)

22



- ❑ 8 bájtos UDP fejléc
- ❑ Egyszerű, kapcsolat nélküli átvitel
 - ▣ C socketek: SOCK_DGRAM
- ❑ Port számok teszik lehetővé a demultiplexálást
 - ▣ 16 bit = 65535 lehetséges port
 - ▣ 0 port nem engedélyezett
- ❑ Kontrollösszeg hiba detektáláshoz
 - ▣ Hibás csomagok felismerése
 - ▣ Nem detektálja az elveszett, duplikátum és helytelen sorrendben beérkező csomagokat (UDP esetén nincs ezekre garancia)

UDP felhasználások

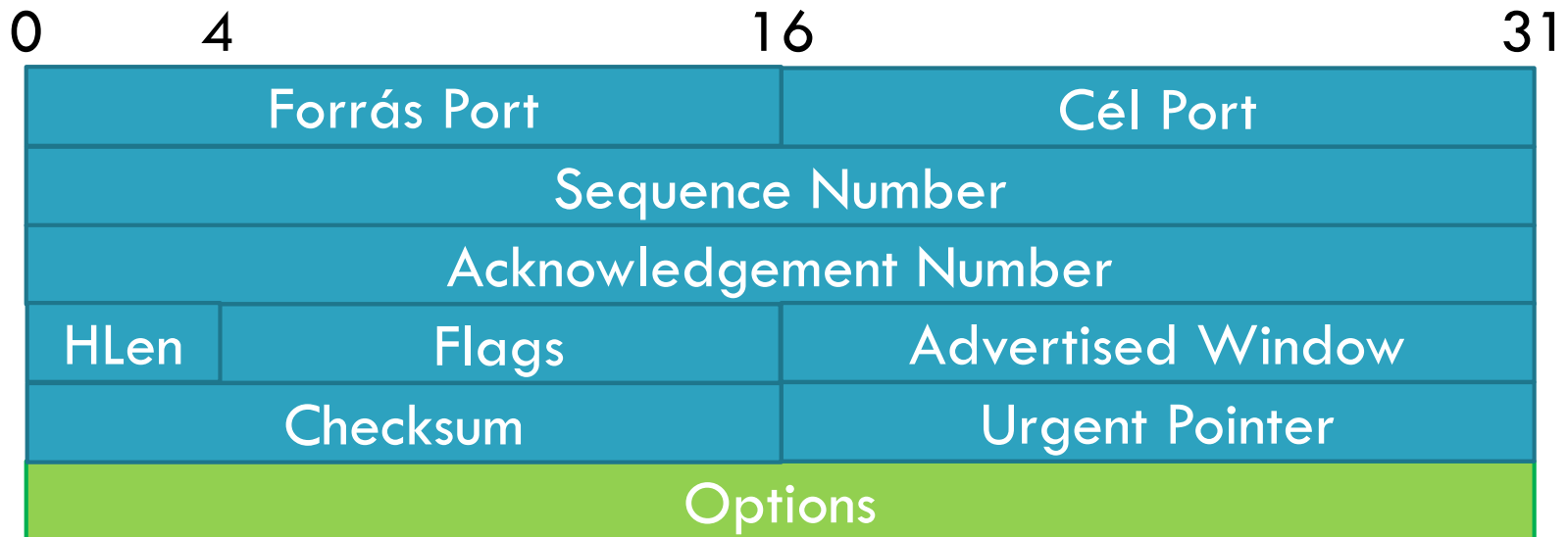
23

- ❑ A TCP után vezették be
 - ▣ Miért?
- ❑ Nem minden alkalmazásnak megfelelő a TCP
- ❑ UDP felett egyedi protokollok valósíthatók meg
 - ▣ Megbízhatóság? Helyes sorrend?
 - ▣ Folyam vezérlés? Torlódás vezérlés?
- ❑ Példák
 - ▣ RTMP, real-time média streamelés (pl. hang, video)
 - ▣ Facebook datacenter protocol

Transmission Control Protocol

24

- ❑ Megbízható, sorrend helyes, két irányú bájtfolyamok
 - ▣ Port számok a demultiplexáláshoz
 - ▣ Kapcsolat alapú
 - ▣ Folyam vezérlés
 - ▣ Torlódás vezérlés, fair viselkedés
- ❑ 20 bájtos fejléc + options fejlécek



Kapcsolat felépítés

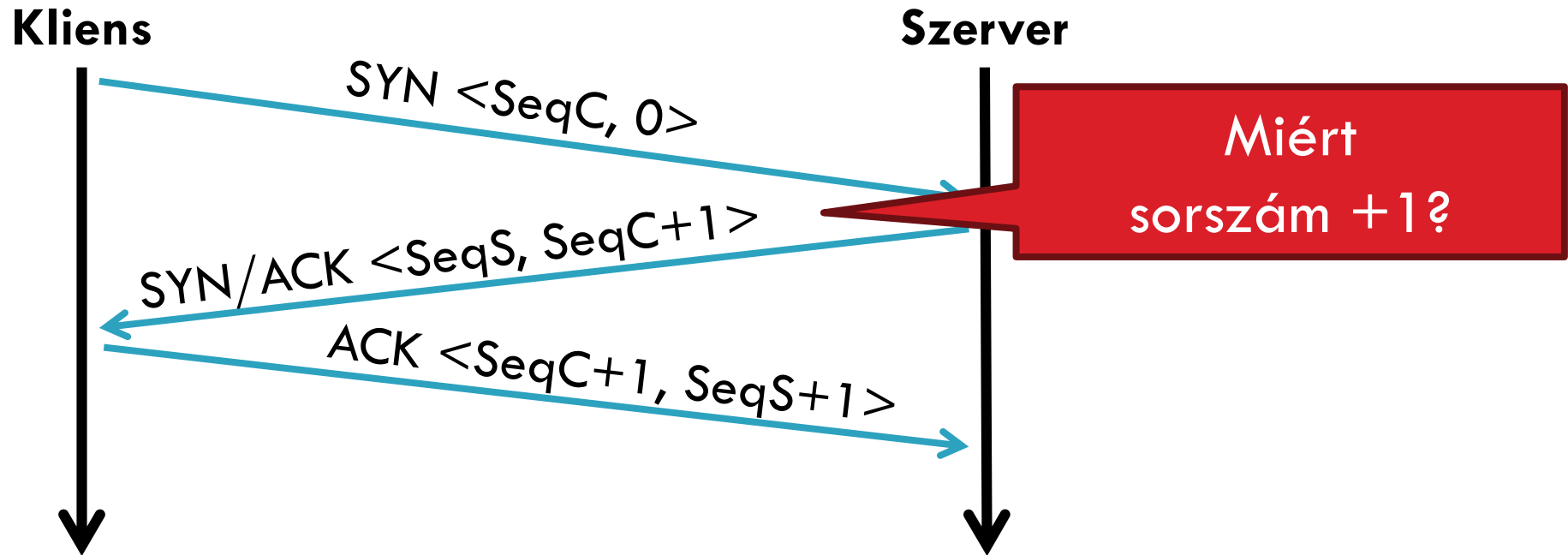
25

- ❑ Miért van szükség kapcsolat felépítésre?
 - ❑ Állapot kialakítása mindkét végponton
 - ❑ Legfontosabb állapot: sorszámok/sequence numbers
 - Az elküldött bájtok számának nyilvántartása
 - Véletlenszerű kezdeti érték
- ❑ Fontos TCP flag-ek/jelölő bitek (1 bites)
 - ❑ SYN – szinkronizációs, kapcsolat felépítéshez
 - ❑ ACK – fogadott adat nyugtázása
 - ❑ FIN – vége, kapcsolat lezárásához

Three Way Handshake

Három-utas kézfogás

26



□ Mindkét oldalon:

- ▣ Másik fél értesítése a kezdő sorszámról
- ▣ A másik fél kezdő sorszámának nyugtázása

Kapcsolat felépítés problémája

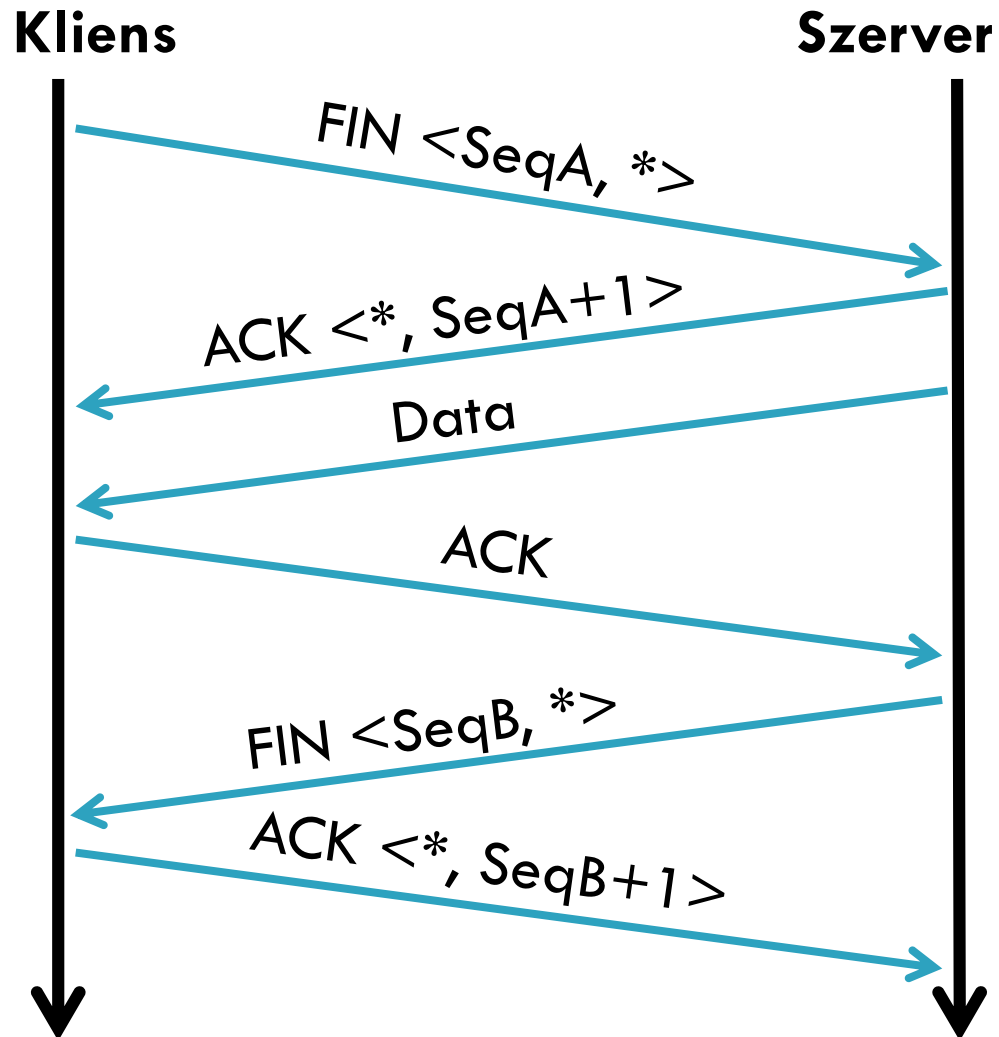
27

- ❑ Kapcsolódási zűrzavar
 - ▣ Azonos hoszt kapcsolatainak egyértelműsítése
 - ▣ Véletlenszerű sorszámmal - biztonság
- ❑ Forrás hamisítás
 - ▣ Kevin Mitnick
 - ▣ Jó random szám generátor kell hozzá!
- ❑ Kapcsolat állapotának kezelése
 - ▣ Minden SYN állapotot foglal a szerveren
 - ▣ SYN flood = denial of service (DoS) támadás
 - ▣ Megoldás: SYN cookies

Kapcsolat lezárása

28

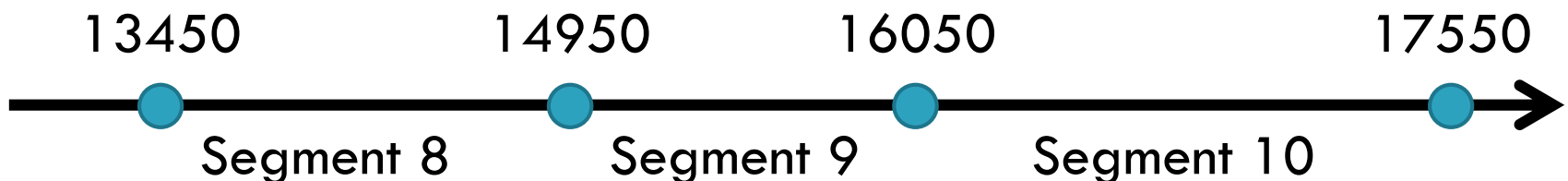
- Mindkét oldal kezdeményezheti a kapcsolat bontását
- A másik oldal még folytathatja a küldést
 - ▣ Félig nyitott kapcsolat
 - ▣ `shutdown()`
- Az utolsó FIN nyugtázása
 - ▣ Sorszám + 1
- Mi történik, ha a 2. FIN elveszik?



Sorszámok tere

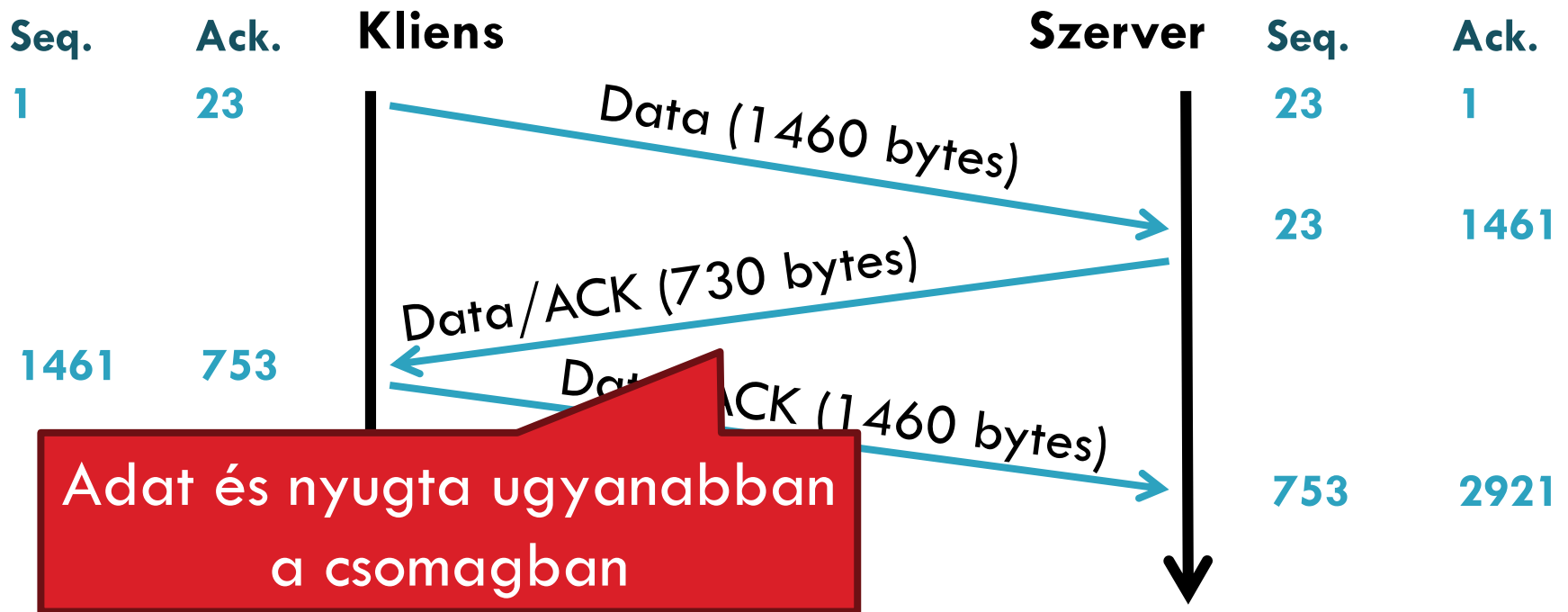
29

- ❑ A TCP egy absztrakt bájt folyamatot valósít meg
 - ❑ A folyam minden bájtja számozott
 - ❑ 32-bites érték, körbefordul egy idő után
 - ❑ Kezdetben, véletlen érték a kapcsolat felépítésénél.
- ❑ A bájt folyamat szegmensekre bontjuk (TCP csomag)
 - ❑ A méretét behatárolja a Maximum Segment Size (MSS)
 - ❑ Úgy kell beállítani, hogy elkerüljük a fregmentációt
- ❑ Minden szegmens egyedi sorszámmal rendelkezik



Kétirányú kapcsolat

30



- Mindkét fél küldhet és fogadhat adatot
 - ▣ Különböző sorszámok a két irányba

Folyam vezérlés

31

- ❑ Probléma: Hány csomagot tud a küldő átvinni?
 - ▣ Túl sok csomag túlterhelheti a fogadót
 - ▣ A fogadó oldali puffer-méret változhat a kapcsolat során
- ❑ Megoldás: csúszóablak
 - ▣ A fogadó elküldi a küldőnek a pufferének méretét
 - ▣ Ezt nevezzük meghirdetett ablaknak: **advertised window**
 - ▣ Egy n ablakmérethez, a küldő n bájtot küldhet el ACK fogadása nélkül
 - ▣ Minden egyes ACK után, léptetjük a csúszóablakot
- ❑ Az ablak akár nulla is lehet!

Folyam vezérlés - csúszóablak

32

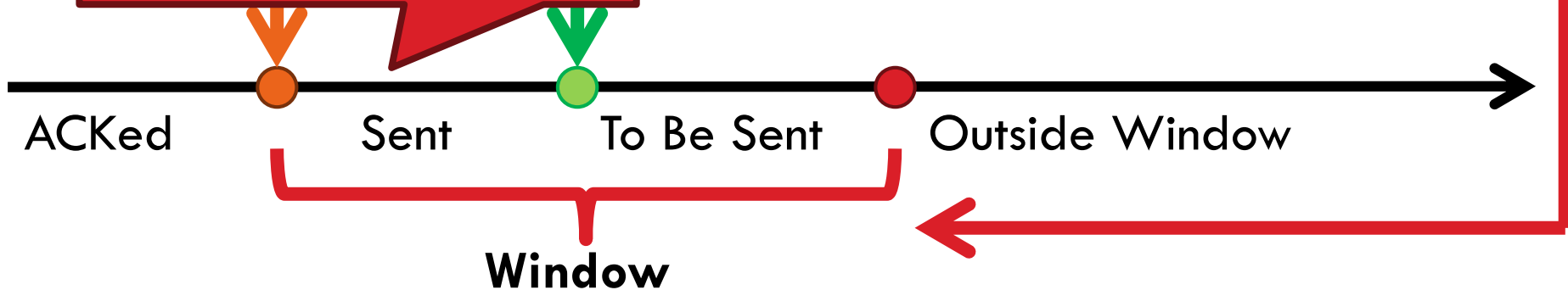
Packet Sent

Src. Port		Dest. Port	
Sequence Number			
Acknowledgement Number			
HL	Flags	Window	
Checksum		Urgent Pointer	

Packet Received

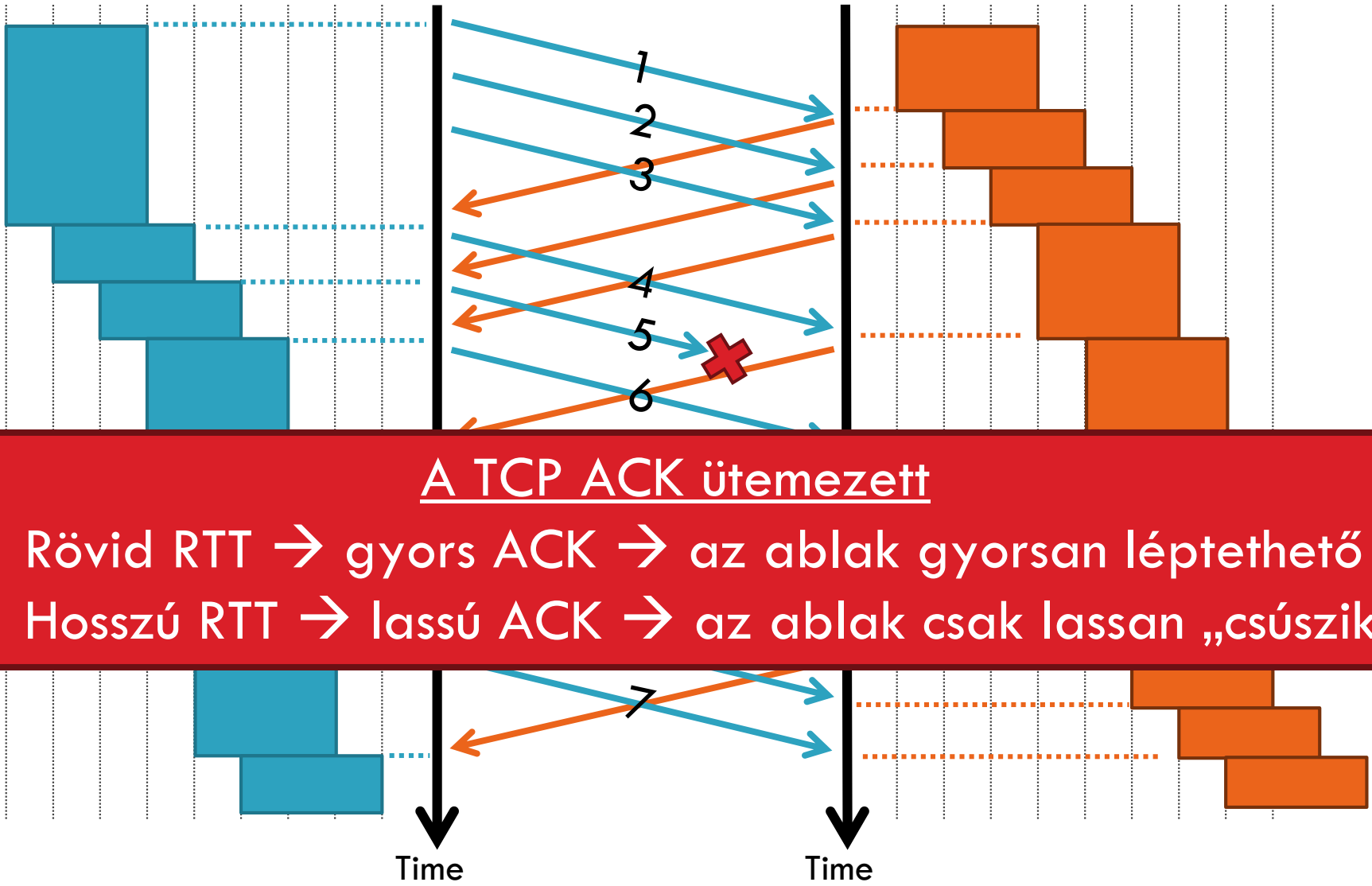
Src. Port		Dest. Port	
Sequence Number			
Acknowledgement Number			
HL	Flags	Window	
Checksum		Urgent Pointer	

Pufferelni kell a nyugtáig



Csúszóablak példa

33



Megfigyelések

34

- Átvitel arányos $\sim w/\text{RTT}$
 - ▣ w : küldési ablakméret
 - ▣ RTT: körülfordulási idő
- A küldőnek pufferelni kell a nem nyugtázott csomagokat a lehetséges újraküldések miatt
- A fogadó elfogadhat nem sorrendben érkező csomagokat, de csak amíg az elfér a pufferben

Mit nyugtázhat a fogadó?

35

1. Minden egyes csomagot
2. Használhat *kumulált nyugtát*, ahol egy n sorszámú nyugta minden $k \leq n$ sorszámú csomagot nyugtáz
3. Használhat *negatív nyugtát* (NACK), megjelölve, hogy mely csomag nem érkezett meg
4. Használhat *szelektív nyugtát* (SACK), jelezve, hogy mely csomagok érkeztek meg, akár nem megfelelő sorrendben
 - SACK egy TCP kiterjesztés
 - SACK TCP

Sorszámok

36

- 32 bites, unsigned
 - ▣ Miért ilyen nagy?
- A csúszó-ablakhoz szükséges...
 - ▣ $|\text{sorszámok tere}| > 2 * |\text{Küldő ablak mérete}|$
 - ▣ $2^{32} > 2 * 2^{16}$
- Elkóborolt csomagok kivédése
 - ▣ IP csomagok esetén a maximális élettartam (MSL) of 120 mp
 - Azaz egy csomag 2 percig bolyonghat egy hálózatban

Buta ablak szindróma

37

- Mi van, ha az ablak mérete nagyon kicsi?
 - ▣ Sok, apró csomag. A fejlécek dominálják az átvitelt.



- Lényegében olyan, mintha bájtónként küldenénk az üzenetet...
 1. `for (int x = 0; x < strlen(data); ++x)`
 2. `write(socket, data + x, 1);`

Nagle algoritmus

38

1. Ha az ablak \geq MSS és az elérhető adat \geq MSS:

Küldjük el az adatot

Egy teljes csomag küldése

2. Különben ha van nem nyugtázott adat::

Várakoztassuk az adatot egy pufferben, amíg nyugtát nem kapunk

3. Különben: küldjük az adatot

Küldünk egy nem teljes csomagot, ha nincs más

□ Probléma: Nagle algoritmus késlelteti az átvitelt

▣ Mi van, ha azonnal el kell küldeni egy csomagot?

1. `int flag = 1;`
2. `setsockopt(sock, IPPROTO_TCP, TCP_NODELAY, (char *) &flag, sizeof(int));`

Hiba detektálás

39

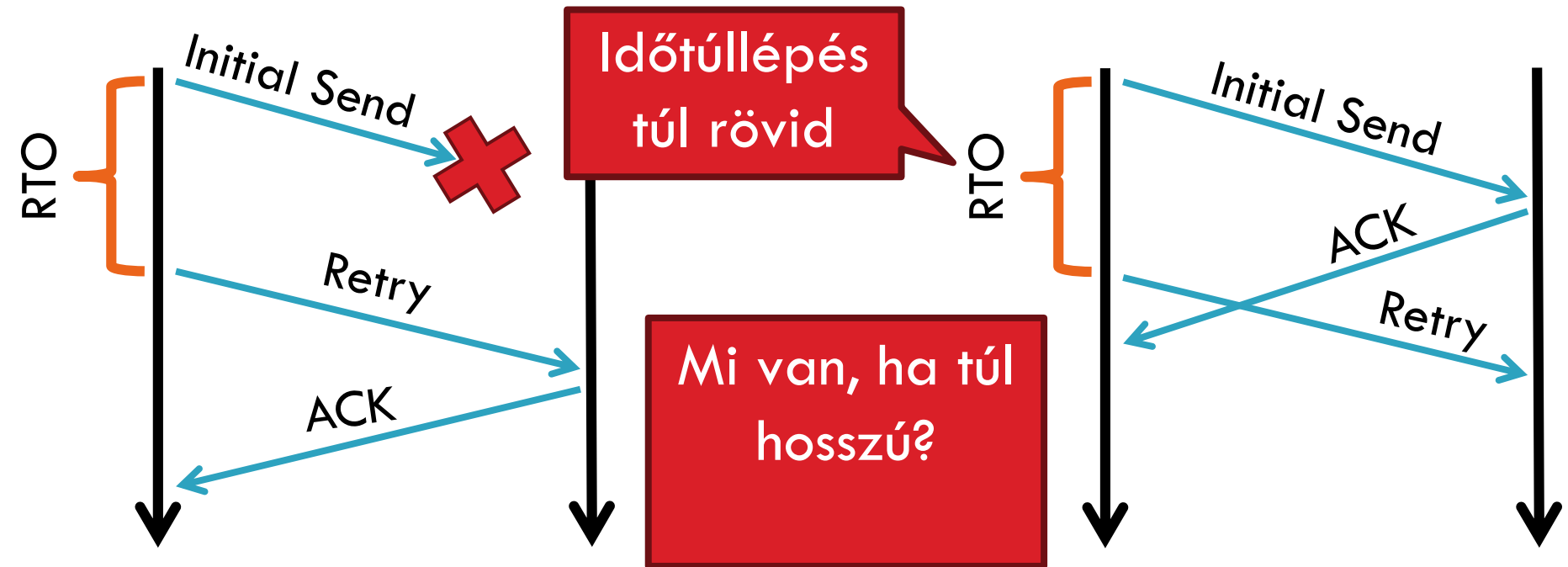
- ❑ A kontrollösszeg detektálja a hibás csomagokat
 - ▣ Az IP, TCP fejlécből és az adatból számoljuk
- ❑ A sorszámok segítenek a sorrendhelyes átvitelben
 - ▣ Duplikátumok eldobása
 - ▣ Helytelen sorrendben érkező csomagok sorba rendezése vagy eldobása
 - ▣ Hiányzó sorszámok elveszett csomagot jeleznek
- ❑ A küldő oldalon: elveszett csomagok detektálása
 - ▣ Időtúllépés (timeout) használata hiányzó nyugtákhoz
 - ▣ Szükséges az RTT becslése a időtúllépés beállításához
 - ▣ Minden nem nyugtázott csomagot pufferelni kell a nyugtáig

Retransmission Time Outs (RTO)

Időtűllépés az újraküldéshez

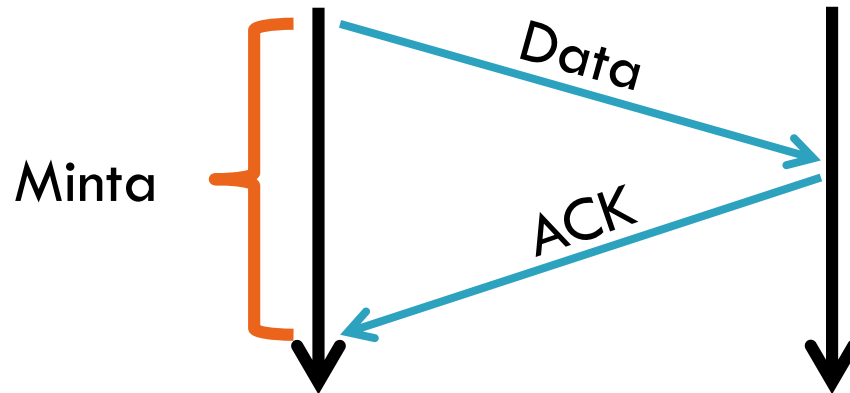
40

- Probléma: Időtűllépés RTT-hez kapcsolása



Round Trip Time becslés

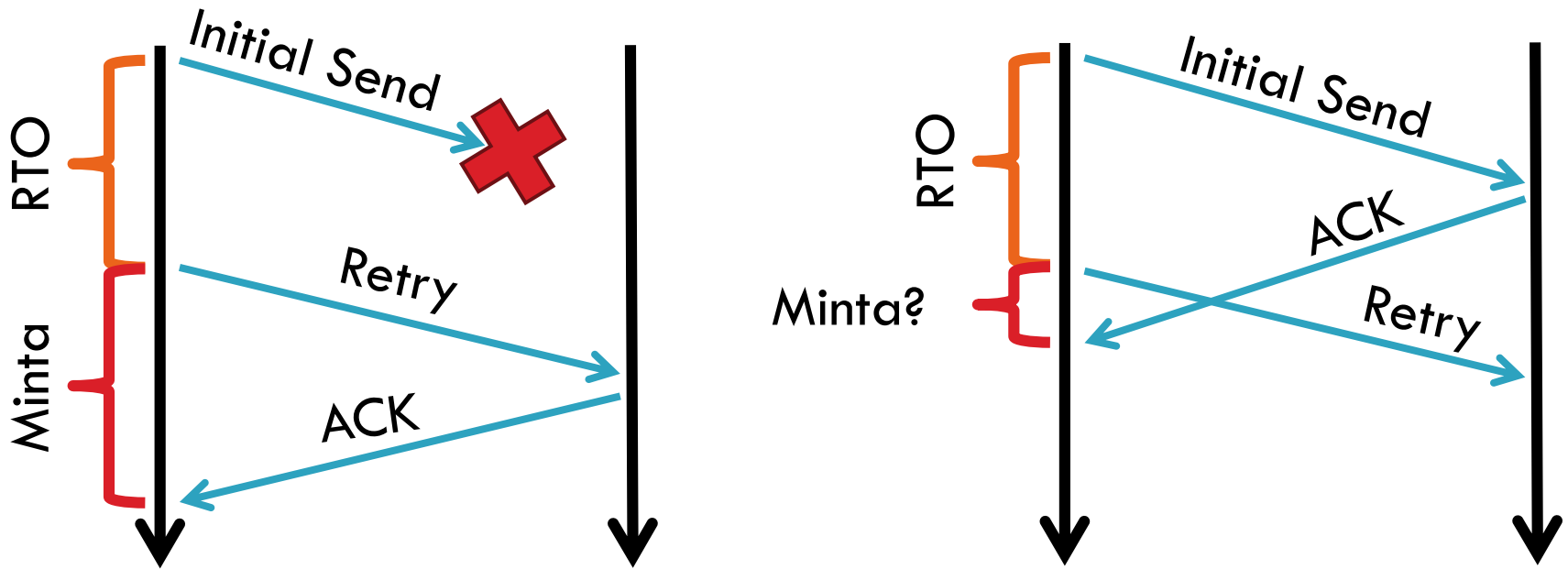
41



- Az eredeti TCP RTT becselője:
 - ▣ RTT becslése mozgó átlaggal
 - ▣ $\text{new_rtt} = \alpha (\text{old_rtt}) + (1 - \alpha)(\text{new_sample})$
 - ▣ Javasolt α : 0.8-0.9 (0.875 a legtöbb TCP esetén)
- $\text{RTO} = 2 * \text{new_rtt}$ (a TCP konzervatív becslése)

Az RTT minta félre is értelmezhető

42



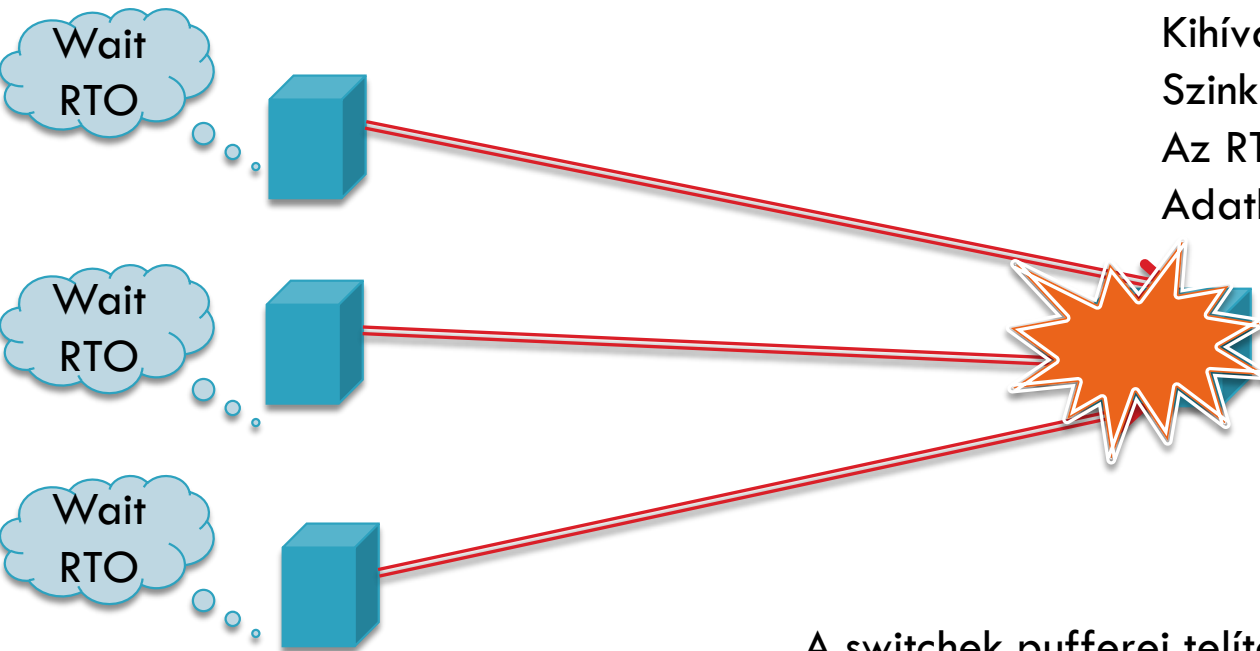
- **Karn algoritmus**: dobjuk el azokat a mintákat, melyek egy csomag újraküldéséből származnak

RTO adatközpontokban???

43

- ❑ TCP Incast probléma – pl. Hadoop, Map Reduce, HDFS, GFS

Sok szimultán küldő egy fogadóhoz



Kihívás:

Szinkronizáció megtörése

Az RTO becslést WAN-ra tervezték

Adatközpontban sokkal kisebb RTT van

1-2ms vagy kevesebb

A switchek pufferei telítődnek és csomagok vesznek el!

Nyugta nem megy vissza ☹

Mi az a torlódás?

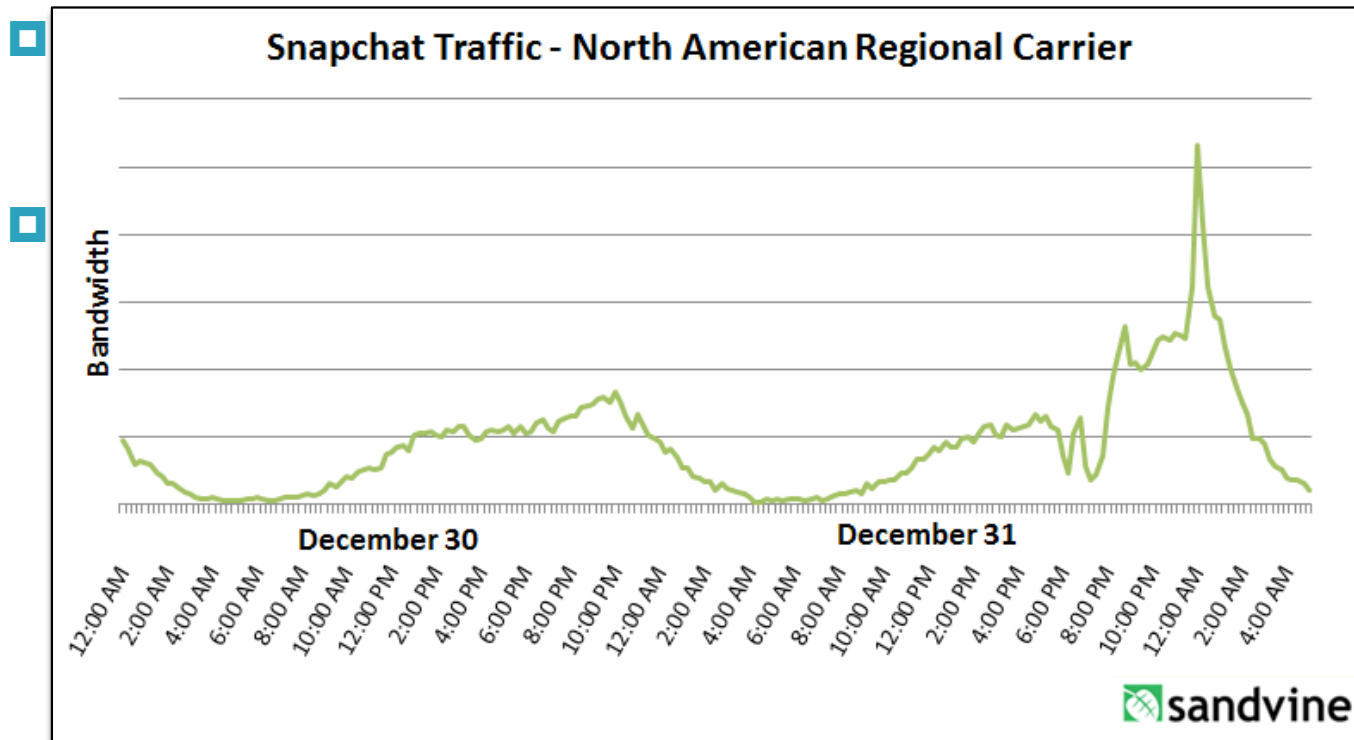
44

- A hálózat terhelése nagyobb, mint a kapacitása
 - ▣ A kapacitás nem egyenletes a hálózatban
 - Modem vs. Cellular vs. Cable vs. Fiber Optics
 - ▣ Számos folyam verseng a sávszélességért
 - otthoni kábel modem vs. corporate datacenter
 - ▣ A terhelés időben nem egyenletes
 - Vasárnap este 10:00 = Bittorrent Game of Thrones

Mi az a torlódás?

45

- A hálózat terhelése nagyobb, mint a kapacitása
 - ▣ A kapacitás nem egyenletes a hálózatban
 - Modem vs. Cellular vs. Cable vs. Fiber Optics



Miért rossz a torlódás?

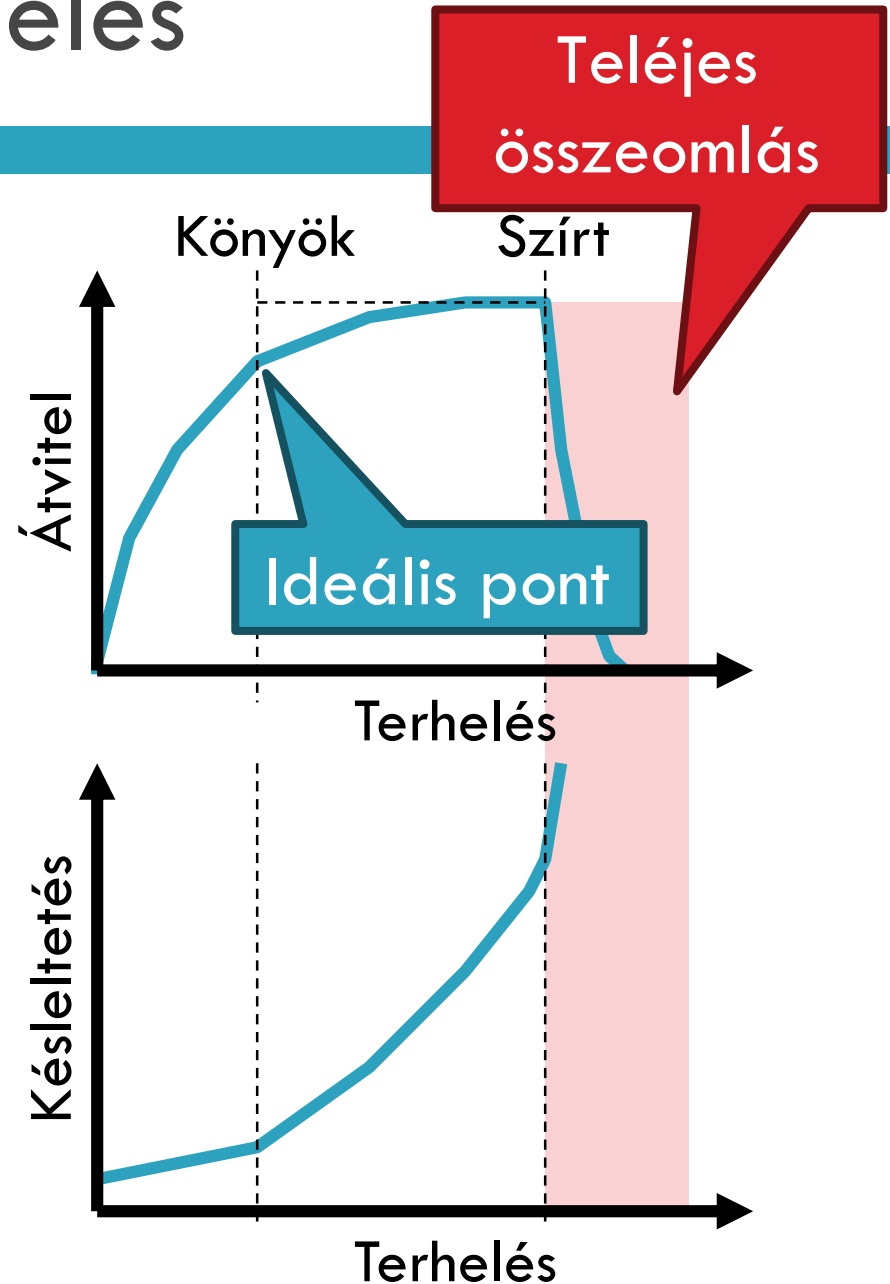
46

- ❑ Csomagvesztést eredményez
 - A routerek véges memóriával (puffer) rendelkeznek
 - Önhasonló Internet forgalom, nincs puffer, amiben ne okozna csomagvesztést
 - Ahogy a routerek puffere elkezd telítődni, csomagokat kezd eldobni... (RED)
- ❑ Gyakorlati következmények
 - A routerek sorai telítődnek, **megnövekedett késleltetés**
 - Sáv szélesség pazarlása az **újraküldések miatt**
 - Alacsony hálózati átvitel (goodput)

Megnövekedett terhelés

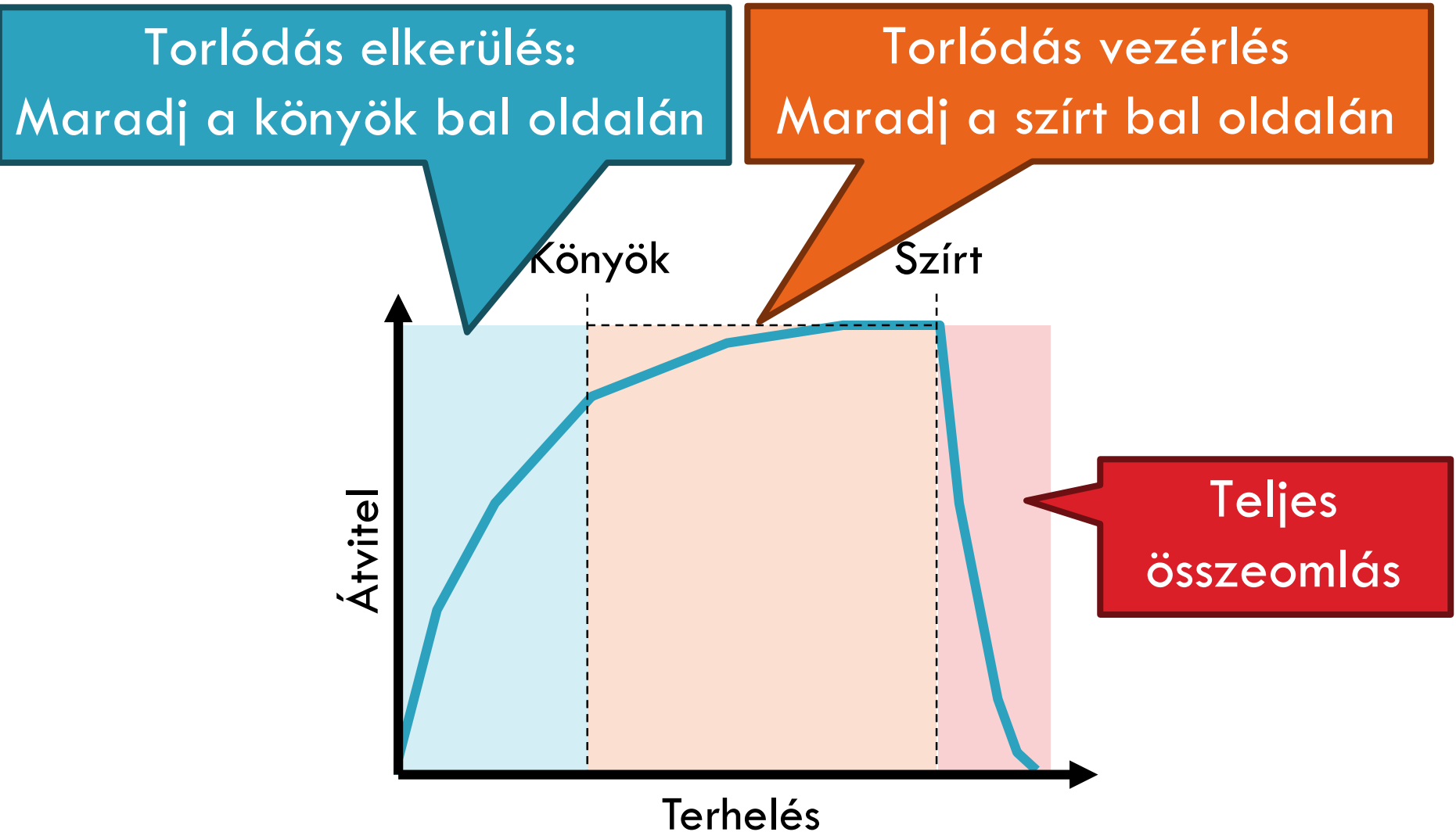
47

- ❑ Könyök („knee”)– a pont, ami után
 - ▣ Az átvitel szinte alig nő
 - ▣ Késleltetés viszont gyorsan emelkedik
- ❑ Egy egyszerű sorban (M/M/1)
 - ▣ Késleltetés = $1 / (1 - \text{utilization})$
- ❑ Szírt („cliff”) – a pont, ami után
 - ▣ Átvitel lényegében leesik 0-ra
 - ▣ A késleltetés pedig $\rightarrow \infty$



Torlódás vezérlés vs torlódás elkerülés

48



Advertised Window

Meghirdetett ablak, újragondolva

49

- ❑ Megoldja-e a torlódás problémáját a TCP esetén a meghirdetett ablak használata?

NEM

- ❑ Ez az ablak csak a fogadót védi a túlterheléstől
- ❑ Egy kellően gyors fogadó kimaxolhatja ezt az ablakot
 - ▣ Mi van, ha a hálózat lassabb, mint a fogadó?
 - ▣ Mi van, ha vannak konkurens folyamatok is?
- ❑ Következmények
 - ▣ Az ablak méret határozza meg a küldési rátát
 - ▣ Az ablaknak állíthatónak kell lennie, hogy elkerüljük a torlódás miatti teljes összeomlást...

Általános megoldások

50

- ❑ Ne csináljunk semmit, küldjük a csomagokat megkülönböztetés nélkül
 - ▣ Nagy csomagvesztés, jósolhatatlan teljesítmény
 - ▣ Teljes összeomláshoz vezethet
- ❑ Erőforrás foglalás
 - ▣ Folyamokhoz előre sáv szélességet allokálunk
 - ▣ Csomagküldés előtt egy tárgyalási szakaszra is szükség van
 - ▣ Hálózati támogatás kell hozzá
- ❑ Dinamikus beállítás
 - ▣ Próbák használata a torlódási szint megbecsléséhez
 - ▣ Gyorsítás, ha torlódási szint alacsony
 - ▣ Lassítás, amint nő a torlódás
 - ▣ Nem rendezett dinamika, elosztott koordináció

TCP Torlódásvezérlés

51

- Minden TCP kapcsolat rendelkezik egy ablakkal
 - ▣ A nem-nyugtázott csomagok számát vezérli
- Küldési ráta $\sim \text{window} / \text{RTT}$
- Ötlet: ablak méretének változtatása a küldési ráta vezérléséhez
- Vezessünk be egy **torlódási ablakot (congestion window)** a küldő oldalon
 - ▣ Torlódás vezérlés egy küldő oldali probléma
 - ▣ Jelölése: cwnd

Két fő komponens

52

1. Torlódás detektálás

- Eldobott csomag egy megbízható jel
 - Késleltetés alapú megoldások – nehéz és kockázatos
- Hogyan detektáljuk a csomag eldobását? Nyugtával
 - Időkorlát lejár ACK fogadása nélkül
 - Számos duplikált ACK jön be sorban (később lesz róla szó)

2. Ráta beállító algoritmus

- *cwnd* módosítása
- Sáv szélesség próba
- Válasz lépés a torlódásra

Ráta vezérlés

53

- Tudjuk, hogy a TCP ACK ütemezett
 - ▣ Torlódás = késleltetés = hosszú várakozás a nyugták között
 - ▣ Nincs torlódás = alacsony késleltetés = gyors ACK
- Alapvető algoritmus
 - ▣ ACK fogadása esetén: növeljük a *cwnd* ablakot
 - Adat leszállítva, valószínűleg gyorsabban is küldhetünk
 - *cwnd* növekedése arányos az RTT-vel
 - ▣ Csomagvesztés esetén: csökkentjük a *cwnd* ablakot
 - Adat elveszett, torlódásnak kell lennie a hálózatban
- Kérdés: milyen függvényt használjuk a növeléshez és csökkentéshez? !!!!

Torlódás vezérlés megvalósítása

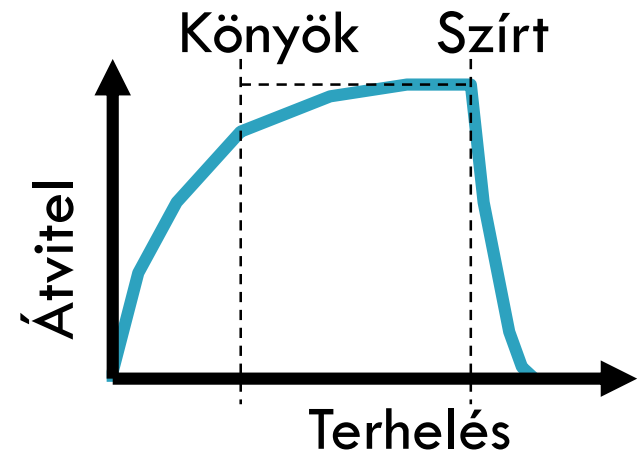
54

- Három változót kell nyilvántartani:
 - ▣ *cwnd*: torlódási ablak
 - ▣ *adv_wnd*: a fogadó meghirdetett ablaka
 - ▣ *ssthresh*: vágási érték (a *cwnd* frissítésére használjuk)
- Küldésnél használjuk: $wnd = \min(cwnd, adv_wnd)$
- A torlódás vezérlés két fázisa:
 1. Lassú indulás („Slow start”) ($cwnd < ssthresh$)
 - Az ún. bottleneck (legsűkebb) sávszélesség meghatározása a cél.
 2. Torlódás elkerülés ($cwnd \geq ssthresh$)
 - AIMD – Additive Increase Multiplicative Decrease

Lassú indulás - Slow Start

55

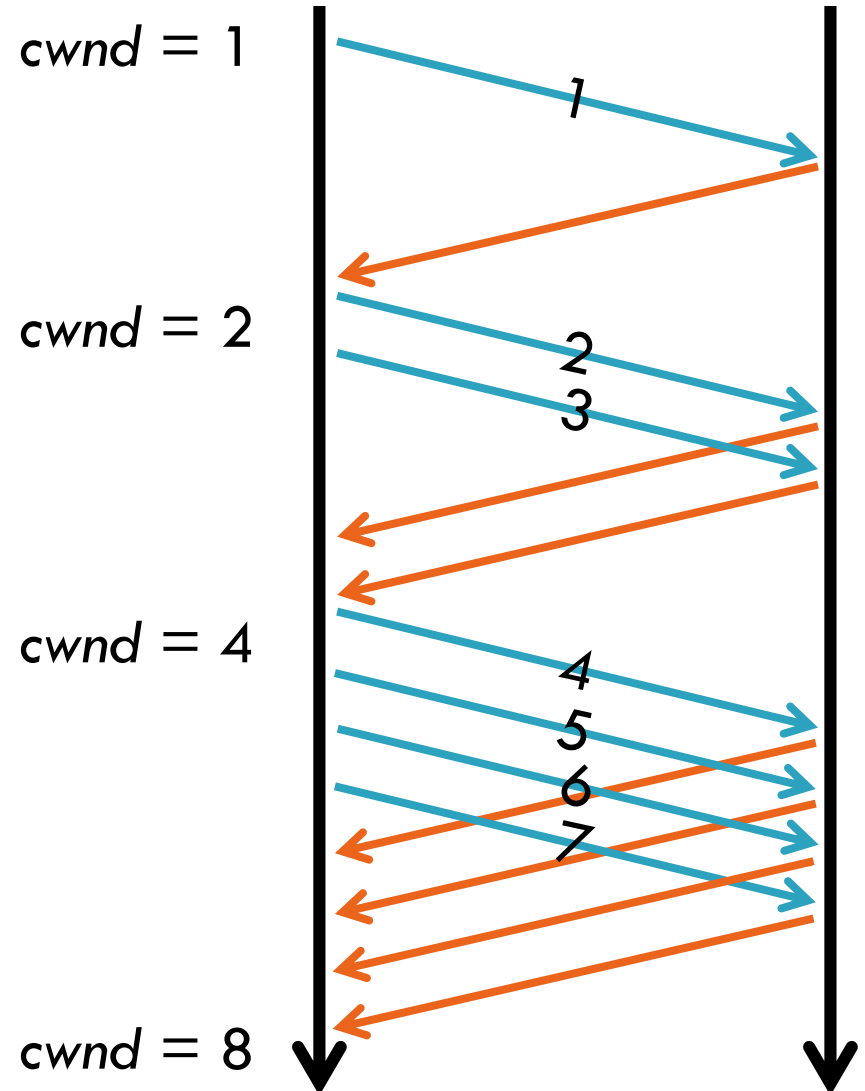
- ❑ Cél, hogy gyorsan elérjük a könyök pontot
- ❑ Egy kapcsolat kezdetén (vagy újraindításakor)
 - ▣ $cwnd = 1$
 - ▣ $ssthresh = adv_wnd$
 - ▣ Minden nyugtázott szegmensre: $cwnd++$
- ❑ Egészen addig amíg
 - ▣ El nem érjük az $ssthresh$ értéket
 - ▣ Vagy csomagvesztés nem történik
- ❑ A Slow Start valójában nem lassú
 - ▣ $cwnd$ exponenciálisan nő



Slow Start példa

56

- $cwnd$ gyorsan nő
- Lelassul, amikor...
 - ▣ $cwnd \geq ssthresh$
 - ▣ Vagy csomagvesztés történik



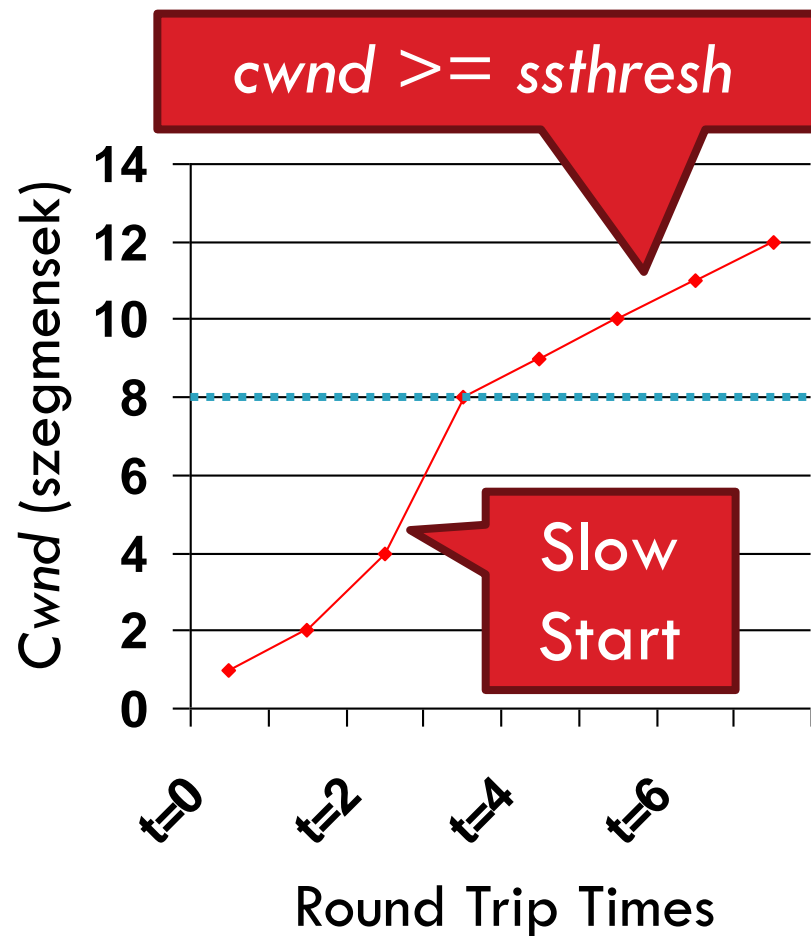
Torlódás elkerülés

57

- Additive Increase Multiplicative Decrease (AIMD) mód
- *ssthresh* valójában egy alsóbecslés a könyök pontra
- **Ha** $cwnd \geq ssthresh$ **akkor**
Minden nyugtázott szegmens alkalmával
növeljük a *cwnd* értékét $(1 / cwnd)$ -vel
(azaz $cwnd += 1 / cwnd$).
- Azaz a *cwnd* eggyel nő, ha minden csomag nyugtázva lett.

Torlódás elkerülés példa

58



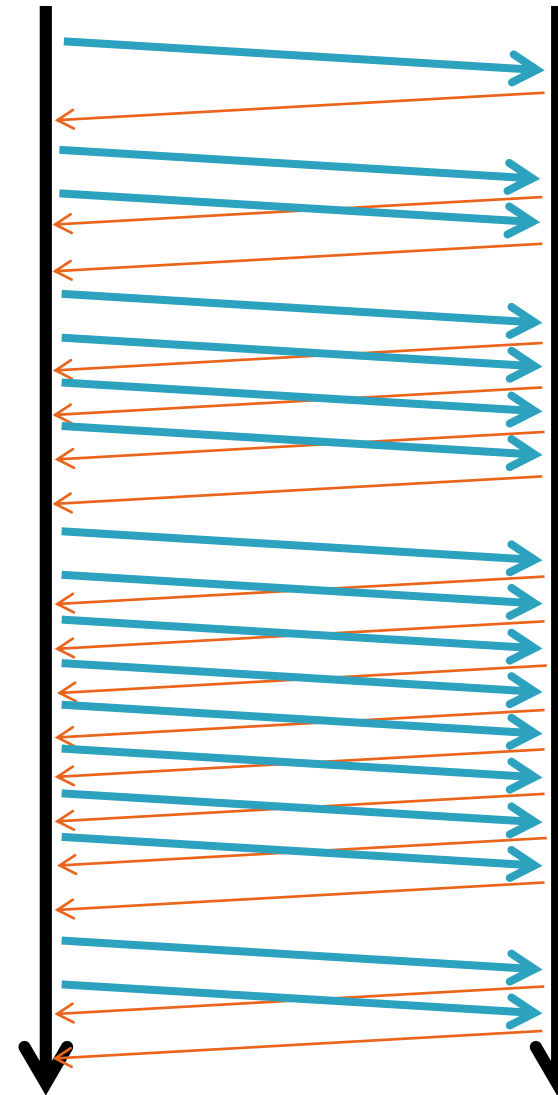
$cwnd = 1$

$cwnd = 2$

$cwnd = 4$

$cwnd = 8$

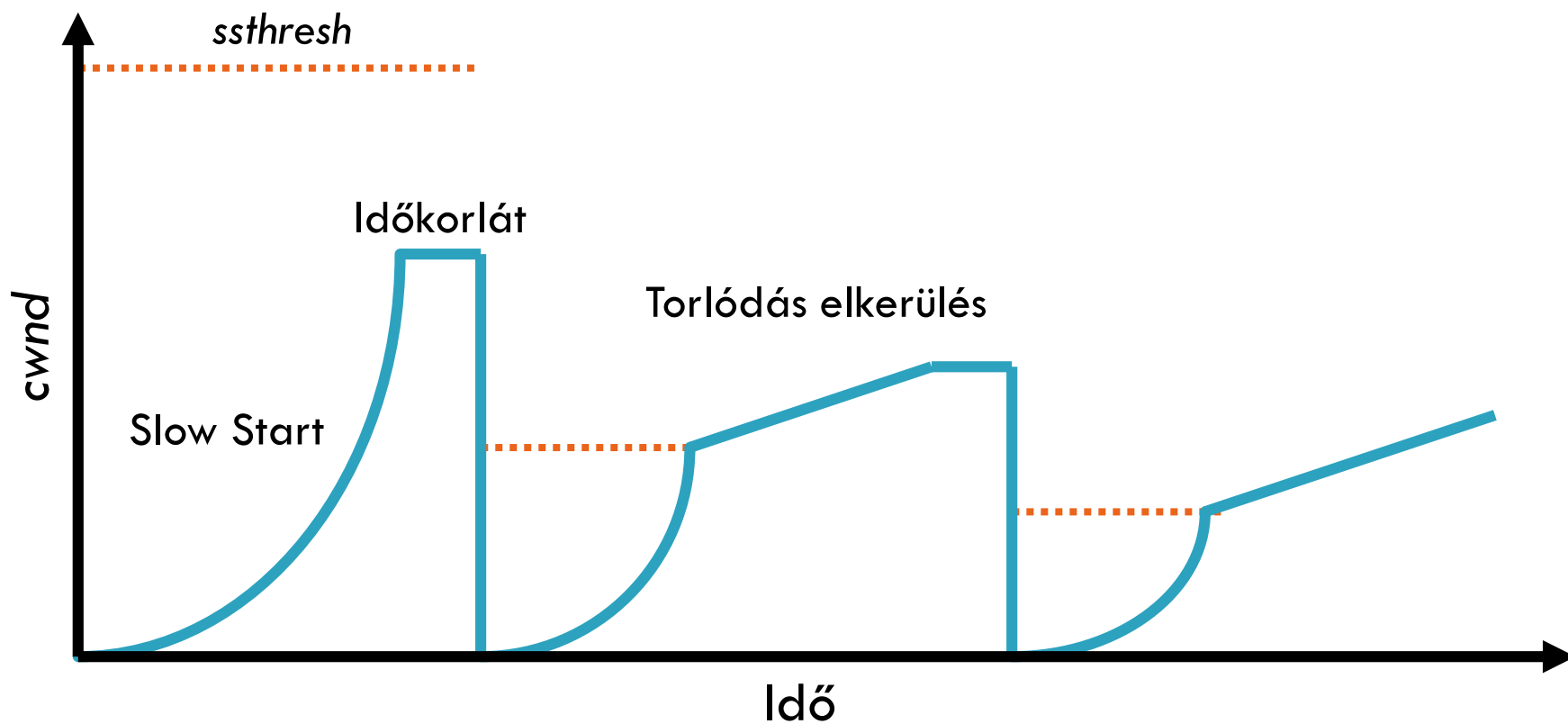
$cwnd = 9$



A teljes kép – TCP Tahoe

(az eredeti TCP)

59



Összefoglalás - TCP jellemzői

60

„A *TCP* egy kapcsolatorientált megbízható szolgáltatás kétirányú bájtfolyamokhoz.”

KAPCSOLATORIENTÁLT

- Két résztvevő, ahol egy résztvevőt egy *IP-cím* és egy *port* azonosít.
- A kapcsolat egyértelműen azonosított a résztvevő párral.
- Nincs se *multi-*, se *broadcast* üzenetküldés.
- A kapcsolatot fel kell építeni és le kell bontani.
- Egy kapcsolat a lezárásáig aktív.

Összefoglalás - TCP jellemzői

61

„A TCP egy kapcsolatorientált megbízható szolgáltatás kétirányú bájtfolyamokhoz.”

MEGBÍZHATÓSÁG

- ❑ Minden csomag megérkezése nyugtázásra kerül.
- ❑ A nem nyugtázott adatcsomagokat újraküldik.
- ❑ A fejléchez és a csomaghoz ellenőrzőösszeg van rendelve.
- ❑ A csomagokat számozza, és a fogadónál sorba rendezésre kerülnek a csomagok a sorszámuk alapján.
- ❑ Duplikátumokat törli.

Összefoglalás - TCP jellemzői

62

„A TCP egy kapcsolatorientált megbízható szolgáltatás kétirányú bájtfolyamokhoz.”

KÉTIRÁNYÚ BÁJTFOLYAM

- Az adatok két egymással ellentétes irányú bájtsorozatként kerülnek átvitelre.
- A tartalom nem interpretálódik.
- Az adatcsomagok időbeli viselkedése megváltozhat: átvitel sebessége növekedhet, csökkenhet, más késés, más sorrendben is megérkezhetnek.
- Megpróbálja az adatcsomagokat időben egymáshoz közel kiszállítani.
- Megpróbálja az átviteli közeget hatékonyan használni.

A TCP evolúciója

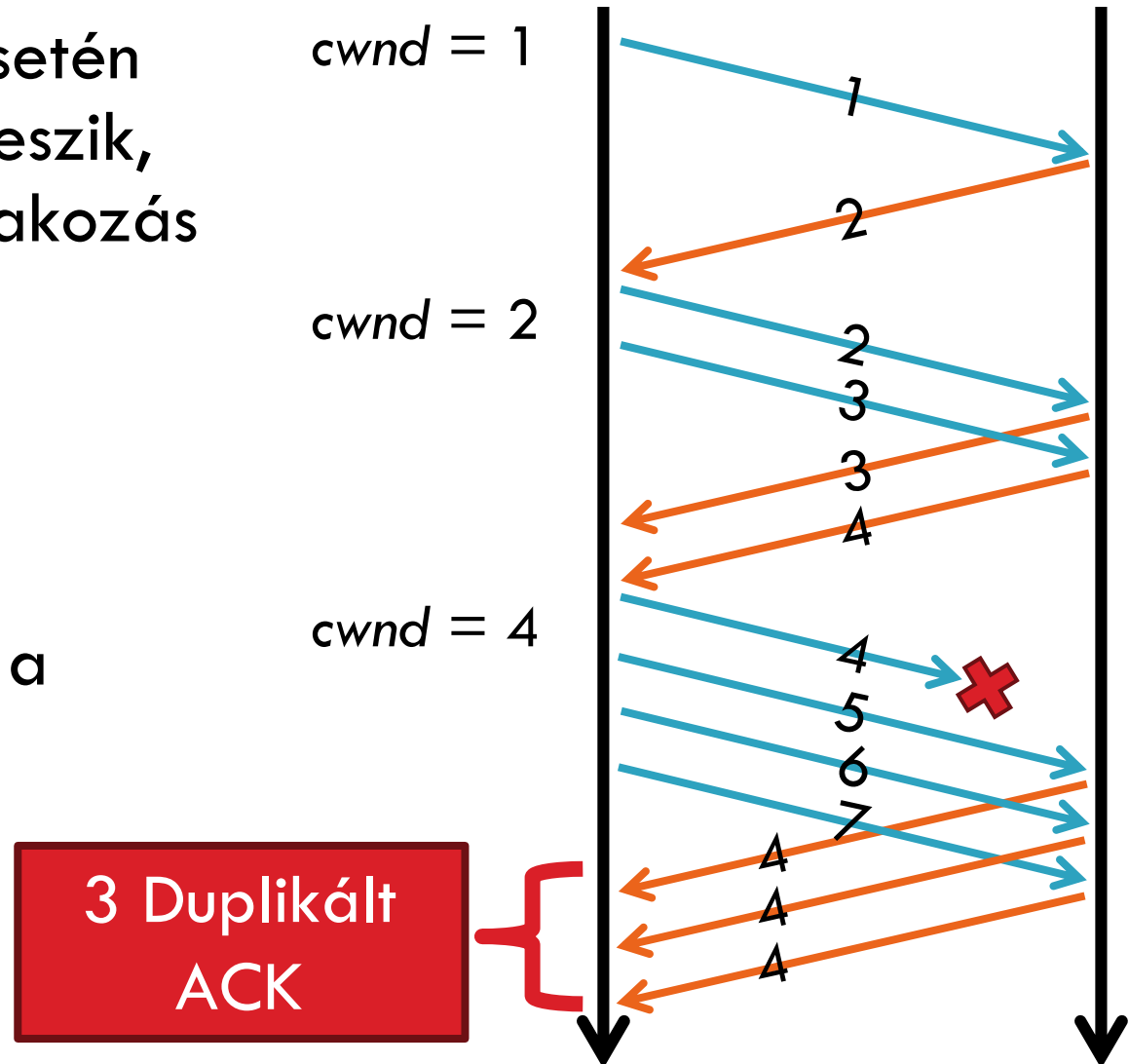
63

- Az eddigi megoldások a TCP Tahoe működéshez tartoztak
 - ▣ Eredeti TCP
- A TCP-t 1974-ben találták fel!
 - ▣ Napjainkba számos változata létezik
- Kezdeti népszerű változat: TCP Reno
 - ▣ Tahoe lehetőségei, plusz...
 - ▣ Gyors újraküldés (Fast retransmit)
 - 3 duplikált ACK? -> újraküldés (ne várjunk az RTO-ra)
 - ▣ Gyors helyreállítás (Fast recovery)
 - Csomagvesztés esetén:
 - $\text{set cwnd} = \text{cwnd} / 2$ (ssthresh = az új cwnd érték)

TCP Reno: Gyors újraküldés

64

- ❑ Probléma: Tahoe esetén ha egy csomag elveszik, akkor hosszú a várakozás az RTO-ig
- ❑ Reno: újraküldés 3 duplikált nyugta fogadása esetén
 - Explicit jele a csomagvesztésnek



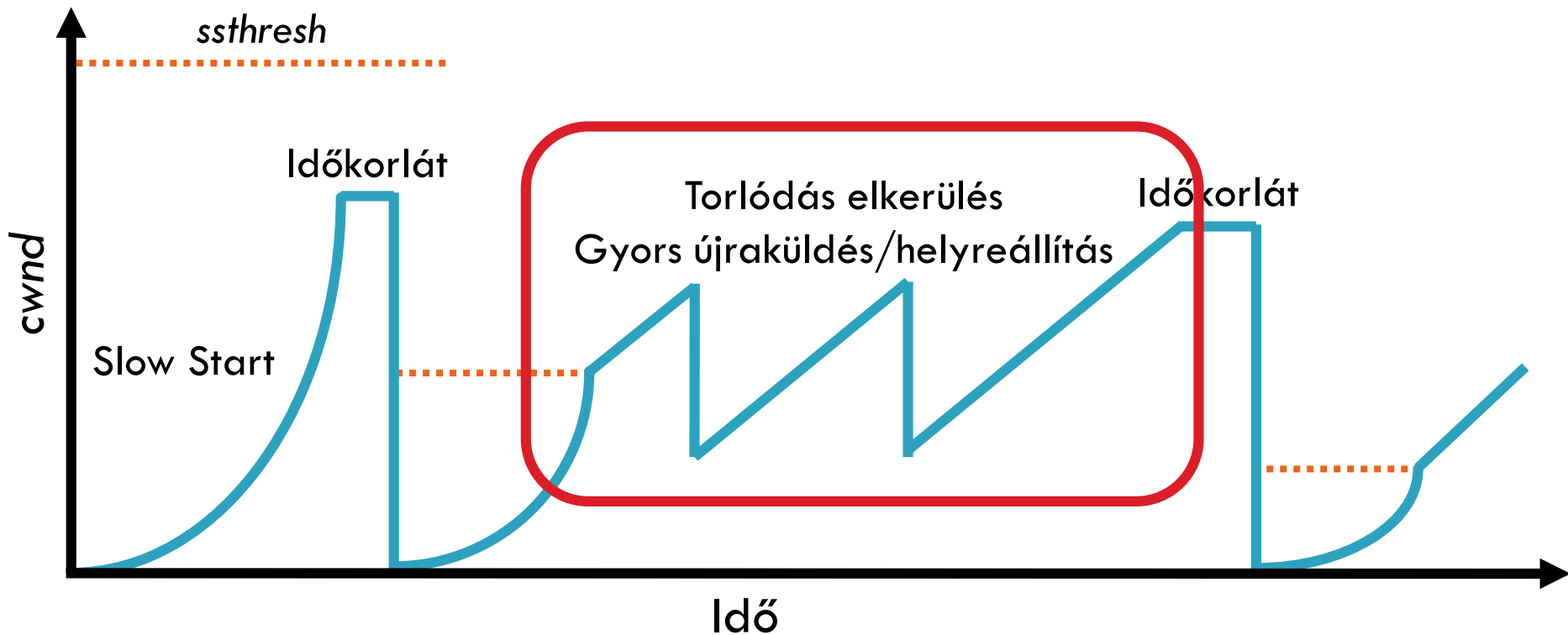
TCP Reno: Gyors helyreállítás

65

- ❑ Gyors újraküldés után módosítjuk a torlódási ablakot:
 - ▣ $cwnd := cwnd/2$ (valójában ez a *Multiplicative Decrease*)
 - ▣ $ssthresh :=$ az új $cwnd$
 - ▣ Azaz nem álltjuk vissza az eredeti 1-re a $cwnd$ -t!!!
 - ▣ Ezzel elkerüljük a felesleges slow start fázisokat!
 - ▣ Elkerüljük a költséges időkorlátokat
- ❑ Azonban ha az RTO lejár, továbbra is $cwnd = 1$
 - ▣ Visszatér a slow start fázishoz, hasonlóan a Tahoe-hoz
 - ▣ Olyan csomagokat jelez, melyeket egyáltalán nem szállítottunk le
 - ▣ A torlódás nagyon súlyos esetére figyelmeztet!!!

Példa: Gyors újraküldés/helyreállítás

66



- Stabil állapotban, a $cwnd$ az optimális ablakméret körül oszcillál
- TCP mindig csomagdobásokat kényszerít ki...

Számos TCP változat...

67

- ❑ Tahoe: az eredeti
 - ❑ Slow start és AIMD
 - ❑ Dinamikus RTO, RTT becsléssel
- ❑ Reno:
 - ❑ fast retransmit (3 dupACKs)
 - ❑ fast recovery ($cwnd = cwnd/2$ vesztes esetén)
- ❑ NewReno: javított gyors újraküldés
 - ❑ Minden egyes duplikált ACK újraküldést vált ki
 - ❑ Probléma: >3 hibás sorrendben fogadott csomag is újraküldést okoz (hibásan!!!)...
- ❑ Vegas: késleltetés alapú torlódás elkerülés
- ❑ ...

TCP a valóságban

68

- Mi a legnépszerűbb variáns napjainkban?
 - ▣ Probléma: TCP rosszul teljesít nagy késleltetés-sávszélesség szorzattal rendelkező hálózatokban (a modern Internet ilyen)
 - ▣ Compound TCP (Windows)
 - Reno alapú
 - Két torlódási ablak: késleltetés alapú és vesztes alapú
 - Azaz egy összetett torlódás vezérlést alkalmaz
 - ▣ TCP CUBIC (Linux)
 - Fejlettebb BIC (Binary Increase Congestion Control) változat
 - Az ablakméretet egy harmadfokú egyenlet határozza meg
 - A legutolsó csomagvesztéstől eltelt T idővel paraméterezett

Nagy késleltetés-sávszélesség szorzat (Delay-bandwidth product)

69

- ❑ Probléma: A TCP nem teljesít jól ha
 - ❑ A hálózat kapacitása (sávszélessége) nagy
 - ❑ A késleltetés (RTT) nagy
 - ❑ Vagy ezek szorzata nagy
 - $b * d =$ maximális szállítás alatt levő adatmennyiség
 - Ezt nevezzük késleltetés-sávszélesség szorzatnak
- ❑ Miért teljesít ekkor gyengén a TCP?
 - ❑ A slow start és az additive increase csak lassan konvergál
 - ❑ A TCP ACK ütemezett (azaz csak minden ACK esetén történik esemény)
 - A nyugták beérkezési gyorsasága határozza meg, hogy milyen gyorsan tud reagálni
 - Nagy RTT → késleltetett nyugták → a TCP csak lassan reagál a megváltozott viszonyokra

Célok

70

- ❑ A TCP ablak gyorsabb növelése
 - ▣ A slow start és az additive increase túl lassú, ha nagy a sávszélesség
 - ▣ Sokkal gyorsabb konvergencia kell
- ❑ Fairség biztosítása más TCP változatokkal szemben
 - ▣ Az ablak növelése nem lehet túl agresszív
- ❑ Javított RTT fairség
 - ▣ A TCP Tahoe/Reno folyamatok nem adnak fair erőforrás-megosztást nagyon eltérő RTT-k esetén
- ❑ Egyszerű implementáció

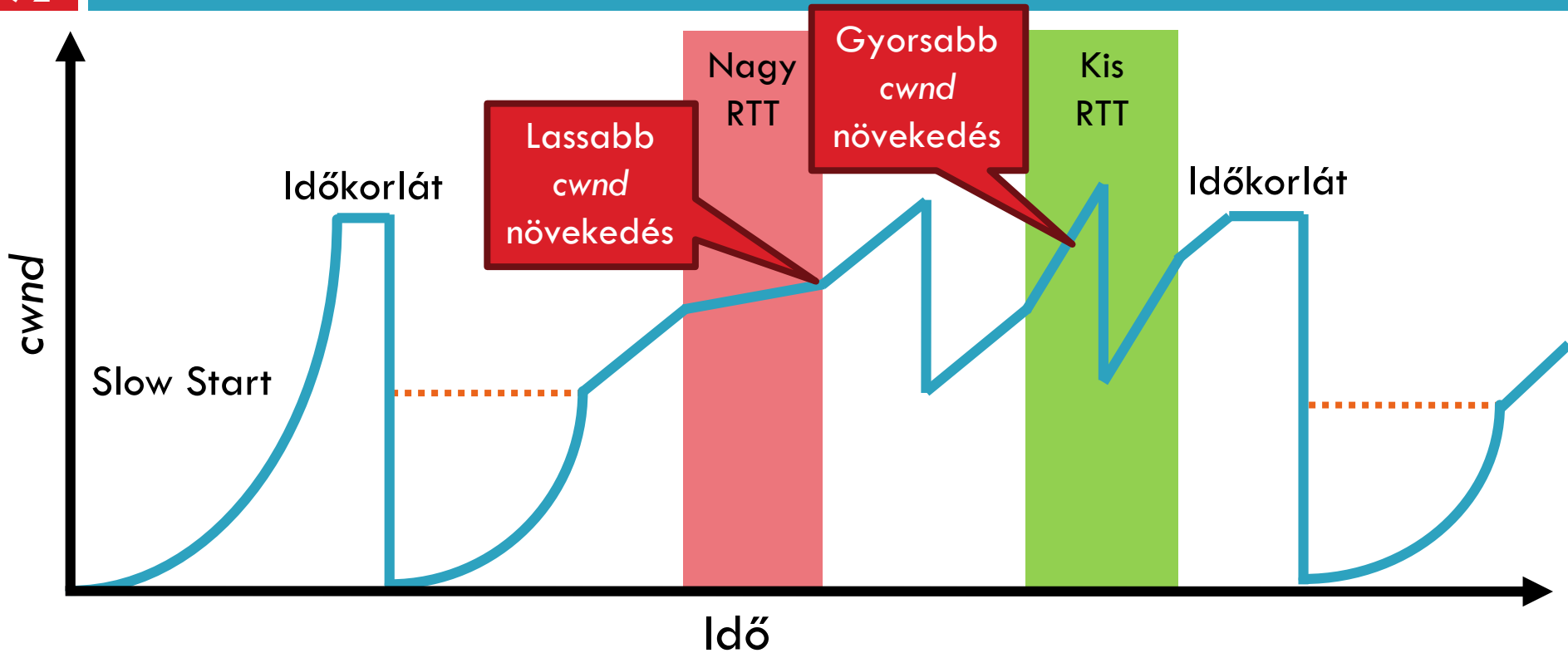
Compound TCP

71

- Alap TCP implementáció Windows rendszereken
- Ötlet: osszuk a *torlódási ablakot* két különálló ablakba
 - ▣ Hagyományos, vesztes alapú ablak
 - ▣ Új, késleltetés alapú ablak
- $wnd = \min(cwnd + dwnd, adv_wnd)$
 - ▣ $cwnd$ -t az AIMD vezérli AIMD
 - ▣ $dwnd$ a késleltetés alapú ablak
- A $dwnd$ beállítása:
 - ▣ Ha nő az RTT, csökken a $dwnd$ ($dwnd \geq 0$)
 - ▣ Ha csökken az RTT, nő a $dwnd$
 - ▣ A növekedés/csökkenés arányos a változás mértékével

Compound TCP példa

72



- Agresszívan reagál az RTT változására
- Előnyök: Gyors felfutás, sokkal fairebb viselkedés más folyamatokkal szemben eltérő RTT esetén
- Hátrányok: folyamatos RTT becslés

TCP CUBIC

73

- Alap TCP implementáció Linux rendszereken
- Az AIMD helyettesítése egy „kübös” (CUBIC) függvénnnyel

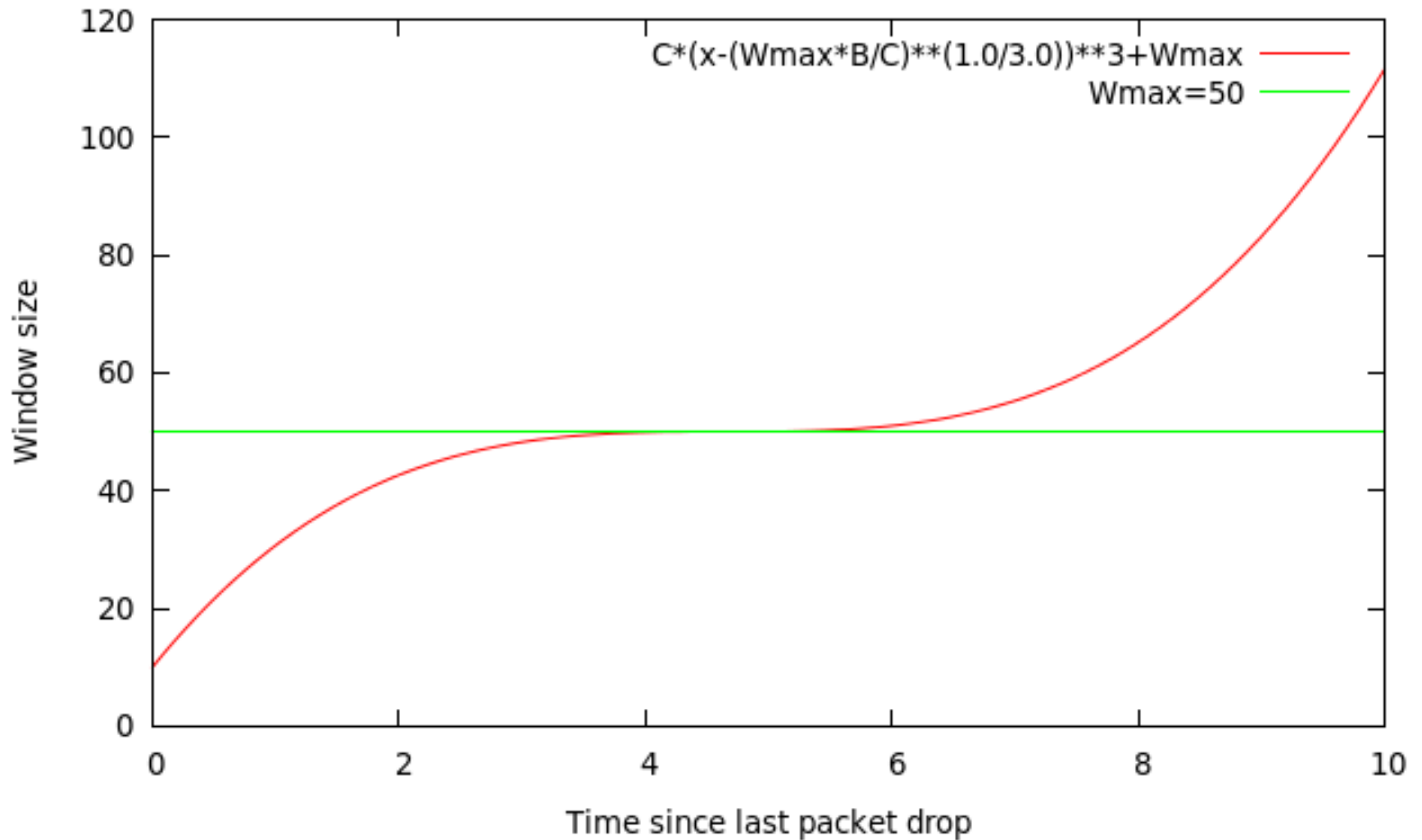
$$W_{cubic} = C(T - K)^3 + W_{max} \quad (1)$$

C is a scaling constant, and $K = \sqrt[3]{\frac{W_{max}\beta}{C}}$

- $B \rightarrow$ egy konstans a multiplicative increase fázishoz
- $T \rightarrow$ eltelt idő a legutóbbi csomagvesztés óta
- $W_{max} \rightarrow$ cwnd a legutolsó csomagvesztés idején

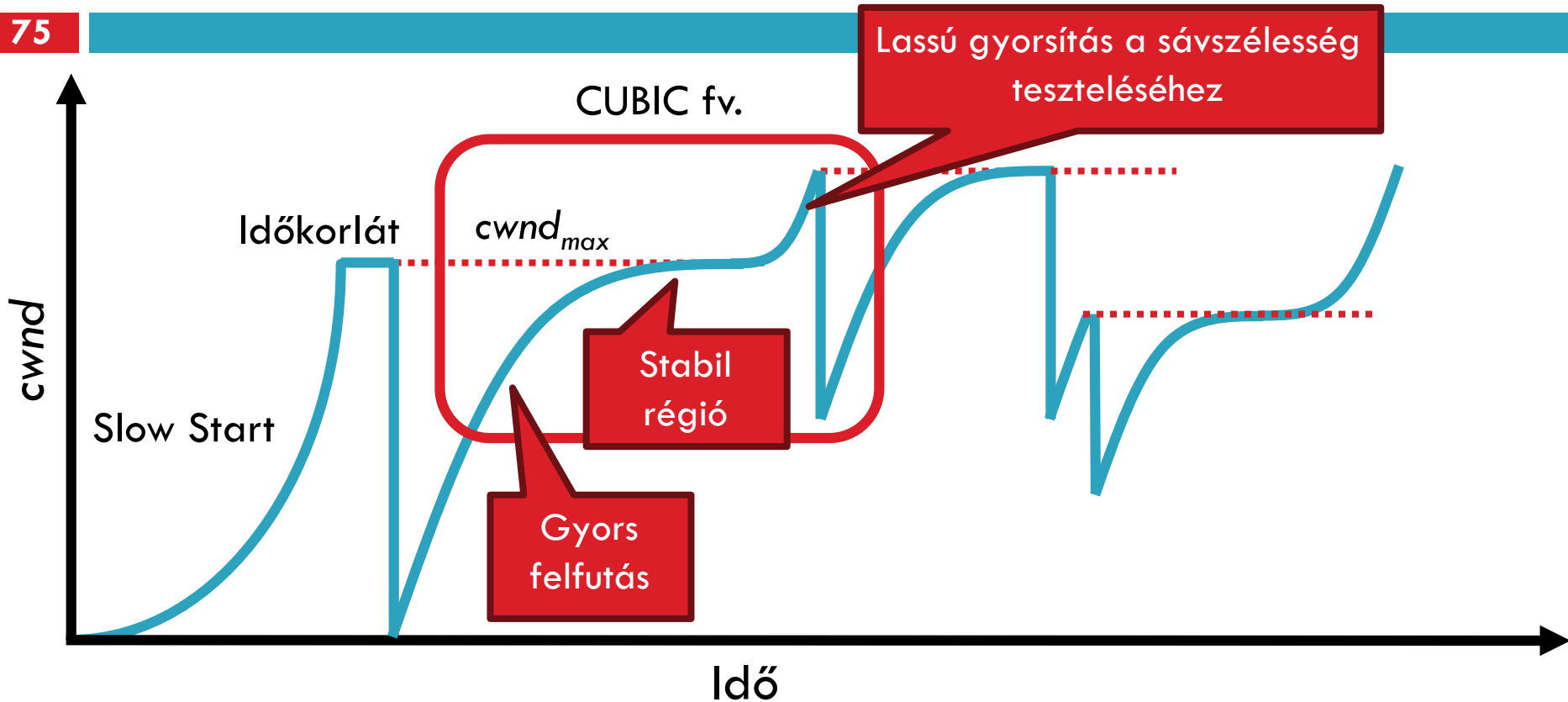
TCP CUBIC

74



TCP CUBIC példa

75



- Kevésbé pazarolja a sávszélességet a gyors felfutások miatt
- A stabil régió és a lassú gyorsítás segít a fairness biztosításában
 - ▣ A gyors felfutás sokkal agresszívabb, mint az additive increase
 - ▣ A Tahoe/Reno variánsokkal szembeni fairnesshez a CUBIC-nak nem szabad ennyire agresszívnak lennie

Problémák a TCP-vel

76

- Az Internetes forgalom jelentős része TCP
- Azonban számos probléma okozója is egyben
 - ▣ Gyenge teljesítmény kis folyamok esetén
 - ▣ Gyenge teljesítmény wireless hálózatokban
 - ▣ DoS támadási felület

Kis folyamok (flows)

77

- ❑ Probléma: kis folyamok esetén torz viselkedés
 - ❑ 1 RTT szükséges a kapcsolat felépítésére (SYN, SYN/ACK)
 - pazarló
 - ❑ *cwnd* mindig 1-gyel indul
 - Nincs lehetőség felgyorsulni a kevés adat miatt
- ❑ Az Internetes forgalom nagy része kis folyam
 - ❑ Többnyire HTTP átvitel, <100KB
 - ❑ A legtöbb TCP folyam el se hagyja a slow start fázist!!!
- ❑ Lehetséges megoldás (Google javaslat):
 - ❑ Kezdeti *cwnd* megnövelése 10-re
 - ❑ TCP Fast Open: kriptográfiai hashek használata a fogadó azonosítására, a három-utas kézfogás elhagyható helyette hash (cookie) küldése a syn csomagban

Wireless hálózatok

78

- ❑ Probléma: A Tahoe és Reno esetén csomagvesztés = torlódás
 - ▣ WAN esetén ez helyes, ritka bit hibák
 - ▣ Azonban hamis vezeték nélküli hálózatokban, gyakori interferenciák
- ❑ TCP átvitel $\sim 1/\sqrt{\text{vesztési ráta}}$
 - ▣ Már néhány interferencia miatti csomagvesztés elég a teljesítmény drasztikus csökkenéséhez
- ❑ Lehetséges megoldások:
 - ▣ Réteg modell megsértése, adatkapcsolati információ a TCP-be
 - ▣ Késleltetés alapú torlódás vezérlés használata (pl. TCP Vegas)
 - ▣ Explicit torlódás jelzés - Explicit congestion notification (ECN)

Szolgáltatás megtagadása

Denial of Service (DoS)

79

- ❑ Probléma: a TCP kapcsolatok állapottal rendelkeznek
 - ▣ A SYN csomagok erőforrásokat foglalnak az serveren
 - ▣ Az állapot legalább néhány percre fennmarad (RTO)
- ❑ SYN flood: elég sok SYN csomag küldése a servernek ahhoz, hogy elfogyjon a memória és összeomoljon a kernel
- ❑ Megoldás: SYN cookie-k
 - ▣ Ötlet: ne tároljunk kezdeti állapotot a serveren
 - ▣ Illesszük az állapotot a SYN/ACK csomagokba (a sorszám mezőbe (sequence number mező))
 - ▣ A kliensnek vissza kell tükrözni az állapotot...

Kitekintés

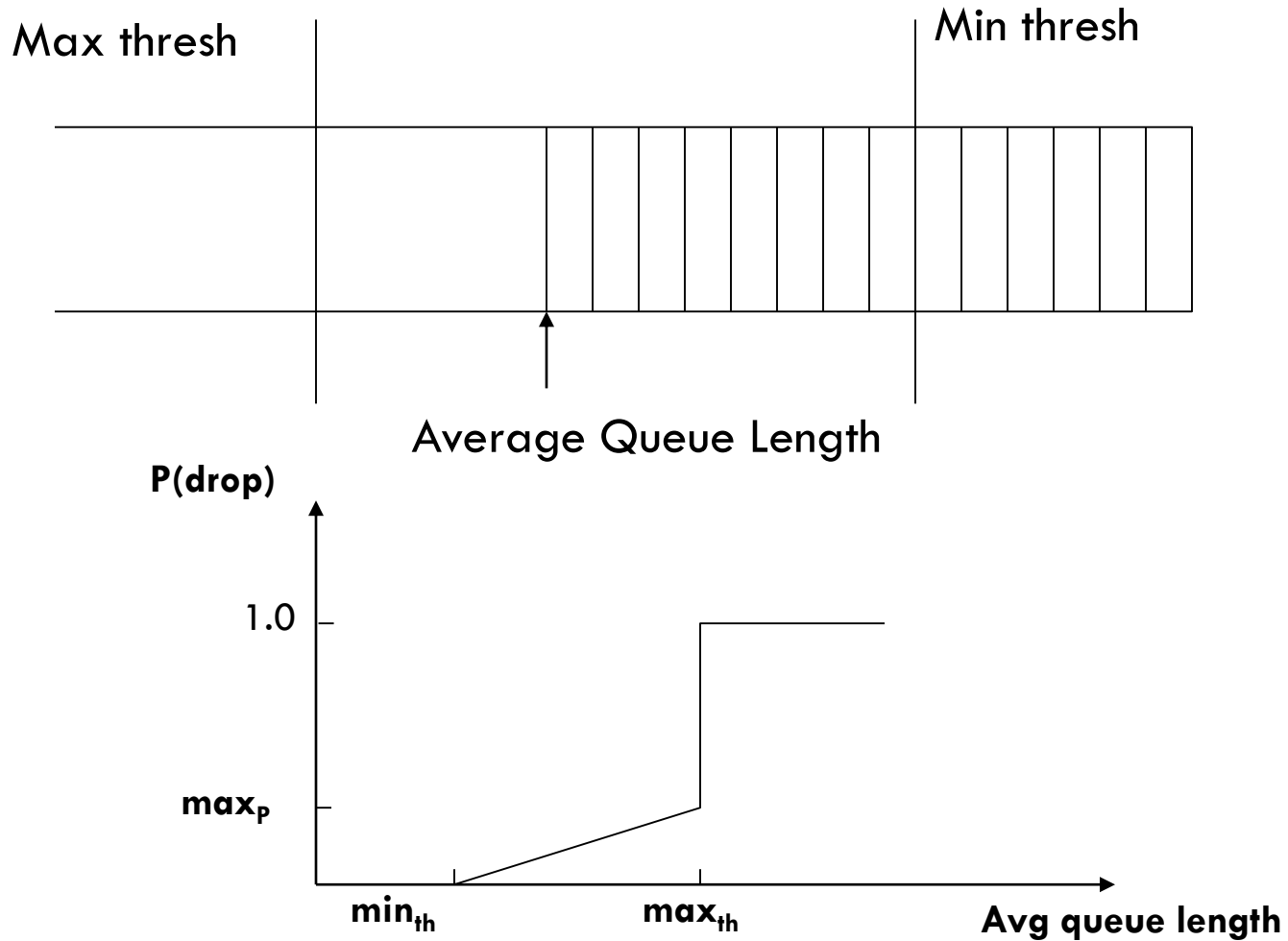
Typical Internet Queuing

- ❑ FIFO + drop-tail
 - ▣ Simplest choice
 - ▣ Used widely in the Internet
- ❑ FIFO (first-in-first-out)
 - ▣ Implies single class of traffic
- ❑ Drop-tail
 - ▣ Arriving packets get dropped when queue is full regardless of flow or importance
- ❑ Important distinction:
 - ▣ FIFO: scheduling discipline
 - ▣ Drop-tail: drop policy

RED Algorithm

- ❑ Maintain running average of queue length
- ❑ If $\text{avgq} < \text{min}_{\text{th}}$ do nothing
 - ▣ Low queuing, send packets through
- ❑ If $\text{avgq} > \text{max}_{\text{th}}$, drop packet
 - ▣ Protection from misbehaving sources
- ❑ Else mark packet in a manner proportional to queue length
 - ▣ Notify sources of incipient congestion
 - ▣ E.g. by ECN IP field or dropping packets with a given probability

RED Operation



RED Algorithm

- Maintain running average of queue length
- For each packet arrival
 - ▣ Calculate average queue size (avg)
 - ▣ If $\min_{th} \leq avg < \max_{th}$
 - Calculate probability P_a
 - With probability P_a
 - ▣ Mark the arriving packet: drop or set-up ECN
 - Else if $\max_{th} \leq avg$
 - ▣ Mark the arriving packet: drop, ECN

Csomag dobás vagy ECN jelölés

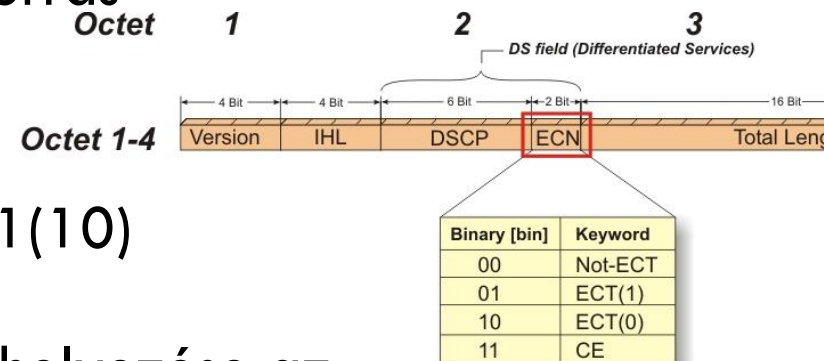
85

❑ Csomag dobás

- ▣ Újra küldés szükséges
- ▣ Egyszerűbb megvalósítás
- ▣ Timeout lejártá után tud reagálni a forrás

❑ ECN jelölés

- ▣ Végpont támogatás szükséges
- ▣ Az IP csomag ECT-0 (01) vagy ECT-1 (10) jelöléssel
- ▣ Dobás helyett -> ECN CE (11) jel elhelyezése az IP fejlécben
- ▣ A fogadó érzékeli a CE jelet, majd a visszamenő TCP nyugtába bebillent egy ECE flaget, mely jelzi a forrásnak a torlódást
- ▣ Hagyományos TCP (CUBIC, RENO, stb.) források az ECE flaget csomagvesztésnek értelmezik, de újra küldés nem szükséges.



Data Center TCP: DCTCP

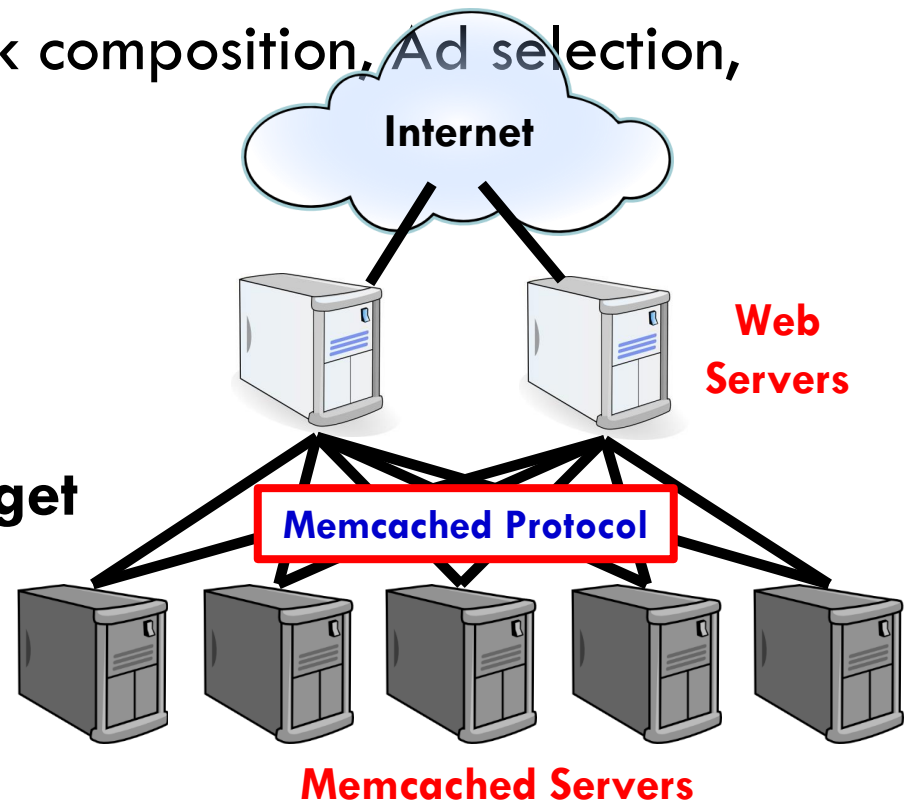
Generality of Partition/Aggregate

- The foundation for many large-scale web applications.
 - ▣ Web search, Social network composition, Ad selection, etc.

- Example: **Facebook**

Partition/Aggregate ~ Multiget

- ▣ Aggregators: **Web Servers**
- ▣ Workers: **Memcached Servers**



Workloads

88

□ Partition/Aggregate
(Query)



Delay-sensitive



□ Short messages [50KB-1MB]
(Coordination, Control state)



Delay-sensitive



□ Large flows [1MB-50MB]
(Data update)



Throughput-sensitive



Impairments

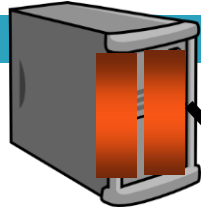
89

- ❑ Incast
- ❑ Queue Buildup
- ❑ Buffer Pressure

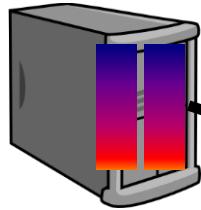
Incast

90

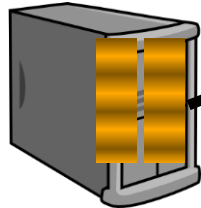
Worker 1



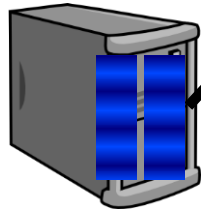
Worker 2



Worker 3



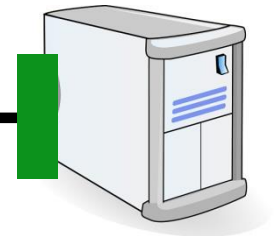
Worker 4



- Synchronized mice collide.

➤ **Caused by Partition/Aggregate.**

Aggregator



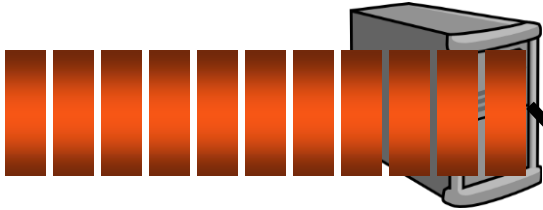
$RTO_{min} = 300 \text{ ms}$

← **TCP timeout**



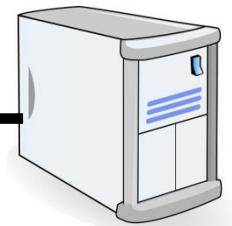
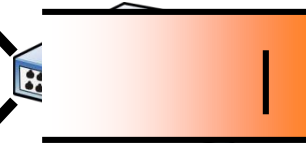
Queue Buildup

Sender 1

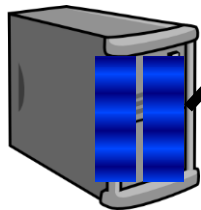


- Big flows buildup queues.
 - Increased latency for short flows.

Receiver



Sender 2



- Measurements in Bing cluster
 - For 90% packets: $RTT < 1ms$
 - For 10% packets: $1ms < RTT < 15ms$

Data Center Transport Requirements

92

1. High Burst Tolerance

- Incast due to Partition/Aggregate is common.

2. Low Latency

- Short flows, queries

3. High Throughput

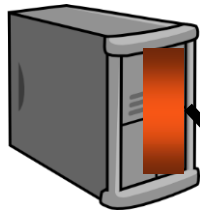
- Continuous data updates, large file transfers

The challenge is to achieve these three together.

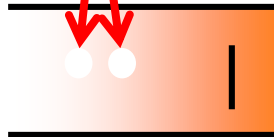
DCTCP: The TCP/ECN Control Loop

Sender 1

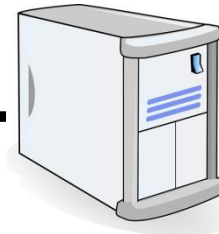
ECN = Explicit Congestion Notification



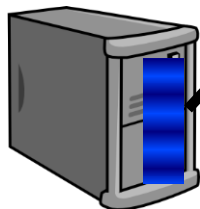
ECN Mark (1 bit)



Receiver



Sender 2



DCTCP: Two Key Ideas

18

1. React in proportion to the **extent** of congestion, not its **presence**.
 - ✓ Reduces **variance** in sending rates, lowering queuing requirements.

ECN Marks	TCP	DCTCP
1 0 1 1 1 1 0 1 1 1	Cut window by 50%	Cut window by 40%
0 0 0 0 0 0 0 0 0 1	Cut window by 50%	Cut window by 5%

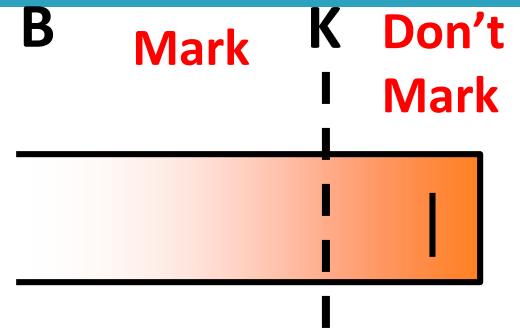
2. Mark based on **instantaneous** queue length.
 - ✓ Fast feedback to better deal with bursts.

Data Center TCP Algorithm

19

Switch side:

- Mark packets when **Queue Length > K**.



Sender side:

- Maintain running average of ***fraction*** of packets marked (α).

In each RTT:

$$F = \frac{\# \text{ of marked ACKs}}{\text{Total \# of ACKs}}$$

$$\alpha \leftarrow (1 - g)\alpha + gF$$

- **Adaptive window decreases:** $cwnd \leftarrow (1 - \frac{\alpha}{2})cwnd$

- Note: decrease factor between 1 and 2.

Köszönöm a figyelmet!