

Protocoalele de transport UDP si TCP. Algoritmi de control al congestiei

Coordonator,

Conf. Dr. Ing. Ștefan Stăncescu

Lecu Radu Șerban

Tică Andra Maria

Vidrașcu Mihai

442A

Cuprins

A. Protocolul UDP

Vidrașcu Mihai

- a. Introducere
 - 1.1. Ce este UDP?
 - 1.2. Comparatie UDP-TCP
 - 1.3. Structura pachetelor
- b. RPC – Remote Procedure Call
- c. Protocole de transport de timp real
 - 3.1. RTCP
 - 3.1.1. Pachetele RTCP
 - 3.1.2. Intervalul de transmisie
- d. DCCP – Datagram Congestion Control Protocol
 - a. Introducere
 - b. Congestie
 - c. Caracteristici DCCP
 - d. Conexiunea DCCP
 - e. Pachetele DCCP
 - f. Header-ul
- e. TCP-Like Congestion Control
 - a. Relatia cu TCP
 - b. Exemplu pentru semiconexiune
- f. TFRC- TCP-Friendly Rate Control
 - 6.1. Mecanismele protocolului
 - 6.2. Ecuatia de transfer
- g. Bibliografie

B. Protocolul TCP

Lecu Radu Șerban

- 1. Introducere
- 2. Modelul serviciului TCP
- 3. Protocolul TCP
- 4. Header-ul TCP
- 5. Stabilirea conexiunii TCP
- 6. Intreruperea conexiunii TCP
- 7. Managementul conexiunii TCP
- 8. Politici de transmisie TCP
- 9. Managementul timerelor TCP
- 10. Algoritmi utilizati de TCP pentru eficientizarea transmisiei de pachete
 - 10.1. Slow start
 - 10.2. Congestion avoidance

- 10.3. Fast retransmit
- 10.4. Fast recovery
- 11. TCP tranzactional (T/TCP)
- 12. Bibliografie

C. Algoritmi de control al congestiei

Tică Andra Maria

- 1. Definirea problemei
- 2. Solutii posibile. Clasificarea Yang&Reddy a algoritmilor
- 3. Algoritmi de control ai congestiei
 - 3.1. Algoritmul RED
 - 3.1.1. Calcularea lungimii medii a cozii
 - 3.1.2. Calcularea probabilitatii de marcare a pachetelor
 - 3.1.3. Implementarea algoritmului RED
 - 3.2. Algoritmul GAIMD
 - 3.2.1. Modelarea ratei de transmisiune
 - 3.2.2. TCP-friendly GAIMD
 - 3.3. Algoritmii HTSCP si STCP
 - 3.3.1. HTSCP
 - 3.3.2. STCP
 - 3.4. Algoritmul binomial BI-TCP
 - 3.4.1. Binary search increase
 - 3.4.2. Additive increase
 - 3.4.3. Slow start
 - 3.4.4. Implementare
 - 3.5. Algoritmul SIMD
- 4. Bibliografie

Protocoalele de transport UDP si TCP. Algoritmi de control al congestiei

A. Protocolul UDP

Vidrașcu Mihai

1. Introducere

Internetul este probabil cea mai mare și mai complexă realizare a omului, este un sistem global de rețele de calculatoare interconectate, care pune la dispoziția miliardelor de utilizatori experiența și cunoștințele acumulate de alungul timpului de om. Este o rețea de rețele care interconectează milioane de calculatoare prin tehnologii de comunicare prin fir, radio sau optice.

Însă pentru a funcționa în armonie, toate calculatoarele trebuie "să vorbească aceeași limba". De aceea au fost create și standardizate o mulțime de protocoale. Protocolul de comunicare este un sistem de formate și reguli pentru schimbul de mesaje între sisteme de telecomunicații. Protocolul pentru internet este un protocol care asigură transmiterea datelor de la un sistem de calcul la altul, prin internet.

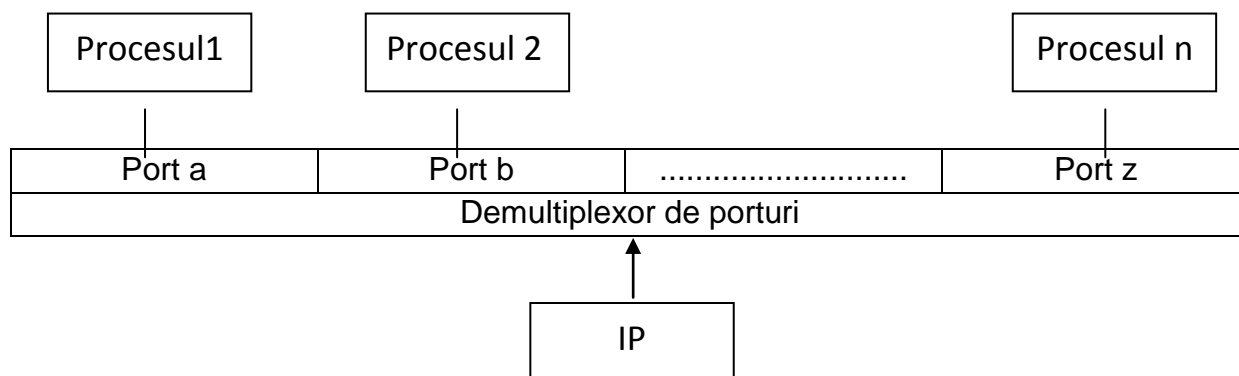
Protocoalele au fost standardizate și pot fi clasificate și după modelul stratificat OSI (Open Systems Interconnection). Există 7 nivele: fizic, legatura de date, rețea, transport, sesiune, prezentare și aplicație. Nivelul de transport oferă transparența transferului de date între utilizatori. Asigura fiabilitatea unei legături prin controlul vitezei, segmentare / desegmentare și controlul erorilor. [14]

În nivelul de transport există 2 protocoale importante, unul orientat pe conexiune și altul fără conexiune (sunt complementare unul altuia). Cel fără conexiune este UDP se ocupă în principal numai de transmisia pachetelor între aplicații. Cel orientat pe conexiune este TCP și este mult mai complex: realizează conexiunea, adaugă fiabilitate prin retransmisie, control pentru debit și congestie. [11]

1.1 Ce este UDP ?

UDP provine de la User Datagram Protocol și a fost proiectat în 1980 de David P. Reed. UDP oferă numai un serviciu minimal de transport (livrare ne-garantată de datagrame) și permite aplicațiilor acces direct la serviciul de datagrame al stratului de IP. UDP este folosit de aplicații care nu necesită servicii de nivelul TCP, sau care vor să folosească servicii de comunicare care nu sunt disponibile în TCP (ca multicast).

Singurele servicii pe care le oferă sunt verificarea datelor prin checksum și multiplexarea pe porturi. Deci o aplicație care folosește UDP trebuie să trateze direct problemele legate de comunicația E2E (end-to-end) pe care un protocol orientat pe conexiune le-ar fi soluționat, ca retransmisia pentru asigurarea fiabilității, segmentarea pe pachete și reasamblarea, controlul debitului, evitarea congestiei etc). [14]



Exemplificarea funcției de multiplexare pe porturi [16]

1.2 Comparatie UDP cu TCP

Se știe ca TCP (Transmission Control Protocol) este orientat pe conexiune, adică se folosește metoda de tip handshake pentru inițializarea comunicării între sisteme. Principalele caracteristici ale TCP sunt următoarele:

Siguranță: TCP gestionează confirmările de primire, retransmișiile și momentele de time-out. Se fac mai multe încercări de livrare a mesajului. Dacă se pierde ceva pe drum, serverul va cere din nou partea pierdută. Nu există cale de mijloc: fie nu există pachete lipsa, fie se întrerupe conexiunea.

Ordine: dacă se trimit două mesaje succesiv, ele vor fi recepționate în ordinea în care au fost trimse. În cazul în care pachetele ajung în alta ordine, TCP stochează datele dezordonate până ajung toate pachetele și apoi le reordonează și le livrează aplicației.

Streaming: datele sunt citite ca stream de octeți; nu există indicatori care să arate limitele unui segment.

TCP trimite 3 pachete pentru a inițializa o conexiune la un socket. Abia apoi poate începe să trimită date. TCP se ocupă și de fiabilitate și controlul congestiei.

UDP este un protocol mai simplu, fără conexiune. Protocoalele fără conexiune nu inițializează o conexiune dedicată între capete. Comunicarea se face prin transmiterea informației într-o singură direcție, fără a verifica starea sau disponibilitatea receptorului. Totuși, un puternic avantaj al UDP-ului este viteza, și este exploatată la maximum în cazul aplicațiilor VoIP (Voice over IP). O comunicare de tip handshake ar îngreuna transmiterea clară a vocii.

Caracteristicile UDP-ului sunt următoarele:

Nesigur: când un mesaj este trimis, nu se știe dacă va ajunge la destinație, se poate pierde pe drum. Nu se aplică conceptele de confirmare, retransmitere sau timeout.

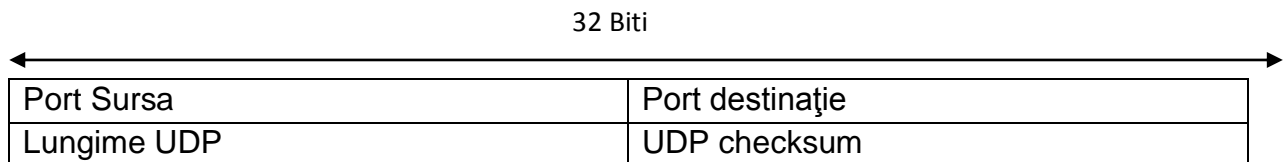
Fără ordine: dacă două mesaje sunt trimise succesiv către același receptor, nu se poate prezice ordinea în care vor ajunge.

Operarea bazată pe datagrame: pachetele sunt trimise individual și sunt verificate pentru integritate doar dacă sosesc la destinație. Pachetele au frontiere bine definite.

Nu există controlul congestiei: UDP nu evită congestiile de unul singur, și este posibil ca aplicațiile de viteză mare să ducă la blocaj dacă nu se implementează metode de control al congestiei la nivelul aplicației. [14]

1.3 Structura pachetelor

UDP transmite segmente formate dintr-un header de 8 Bytes, urmat de datele efective de transmis.



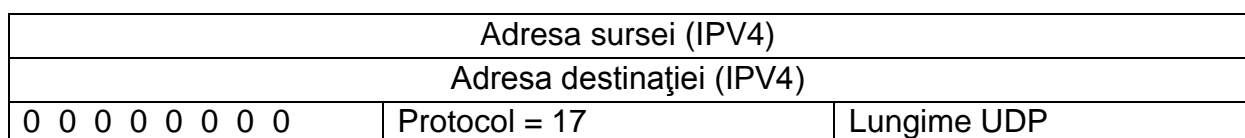
Cele două porturi identifică cele două sisteme de calcul care comunică, sursa și destinația. Când un pachet UDP sosește, încărcatura lui este preluată de procesul atașat portului destinație. Pe scurt, porturile asigură livrarea segmentelor aplicației corecte.

Portul sursă este necesar în cazul în care receptorul trebuie să răspundă emițătorului. Copiind portul sursă din pachetul abia sosit în câmpul de destinație al pachetului de răspuns, procesul care trimite răspunsul poate specifica ce proces de pe emițător îl primește.

Câmpul „Lungime UDP” include headerul de 8 Bytes și datele. Lungimea minimă este de 8 Bytes (pentru header). Lungimea maximă este de 65.515 bytes (dimensiune mai mică decât maximul reprezentabil pe 16 biți; limita este dată de dimensiunea maximă a pachetelor IP).

Suma de verificare este opțională și se folosește pentru creșterea fiabilității. Face verificarea header-ului, a datelor și un pseudoheader IP conceptual. Când se face calculul, câmpul checksum este setat la 0 și câmpul de date este bordat cu un byte adițional de 0 dacă lungimea lui este un număr impar. Algoritmul de verificare constă în sumarea tuturor cuvintelor de 16 biți în complement față de 1 și scoaterea complementului lui 1 din rezultat. Deci, când receptorul face calculul pe întregul segment, inclusiv pe câmpul de checksum, rezultatul ar trebui să dea 0. Dacă checksum-ul nu se calculează, atunci se va stoca cu valoarea 0.

În cazul IPv4, pseudoheader-ul arată ca în figura următoare:



Contine adresele IPv4 pe 32 de biți ale emițătorului și receptorului, numărul de protocol pentru UDP(17), și lungimea segmentului (inclusiv headerul). La IPv6 este similar, însă câmpul de checksum nu mai este opțional. Pseudoheader-ul va arăta astfel:

Adresa sursei (IPV6)	
Adresa destinației (IPV6)	
Lungimea pachetului din stratul superior	
0000...000 (24 de biti 0)	Urmatorul header

Includerea pseudoheader-ului în calculul checksum-ului ajuta la detectarea pachetelor livrate gresit, da includerea lui violeaza ierarhia protocolului, pentru ca adresele IP din el apartin stratului IP, nu celui UDP. TCP foloseste același pseudoheader pentru checksum-ul lui.

Trebuie tinut cont ca UDP nu are nici o metoda de control a congestiei și a debitului, și nici nu retransmite în cazul primirii unui segment eronat. Totul ramane pe seama proceselor user-ului. Înșã, oferã o interfață cãtre protocolul IP și poate și demultiplexa mai multe procese folosind porturi si, optional, detectie de erori la receptie. Este util în situatii de comunicare client-server. Deseori, clientul trimite o cerere scurta cãtre server și asteapta un rãapuns rapid. Dacă fie rãapunsul fie cererea s-au pierdut, clientul asteapta o perioada de time-out, și apoi incearca din nou. Pe langa faptul ca programarea e mai usoara, sunt necesare și putine mesaje (unul în fiecare directie) decat în cazul unui protocol care cere o setare initiala, ca TCP.

O aplicație care foloseste UDP în acest fel este DNS (Domain Name Server). Pe scurt, un program care trebuie să caute adresa IP a unui host oarecare poate trimite un pachet UDP care contine numele host-ului cãtre un server DNSS. Serverul raspunde cu un pachet UDP care contine adresa IP a hostului. Nu este nevoie de nici o setare initiala și eliberare ulterioara a conexiunii. [11]

UDP mai este folosit în aplicații care pun viteza inaintea calitatii, în aplicații de tip FINGER (FINGER este un protocol simplu pentru schimbul de status-uri și informații despre utilizator), si, în general, acolo unde cererile sunt rapide și necesita un singur pachet de rãapuns (exemplu: SNMP – Simple Network Management Protocol, RIP – Routing Information Protocol, DHCP – Dynamic Host Configuration Protocol). [6]

Traficul video și audio este în general transmis folosind UDP. Protocoalele de streaming audio-video au fost proiectate să suporte eventualele pierderi de pachete, deci va aparea doar o scadere slaba a calitatii, în locul unor intarzieri deranjante care ar aparea în cazul retransmisiei pachetelor pierdute. [7]

2. RPC – Remote Procedure Call

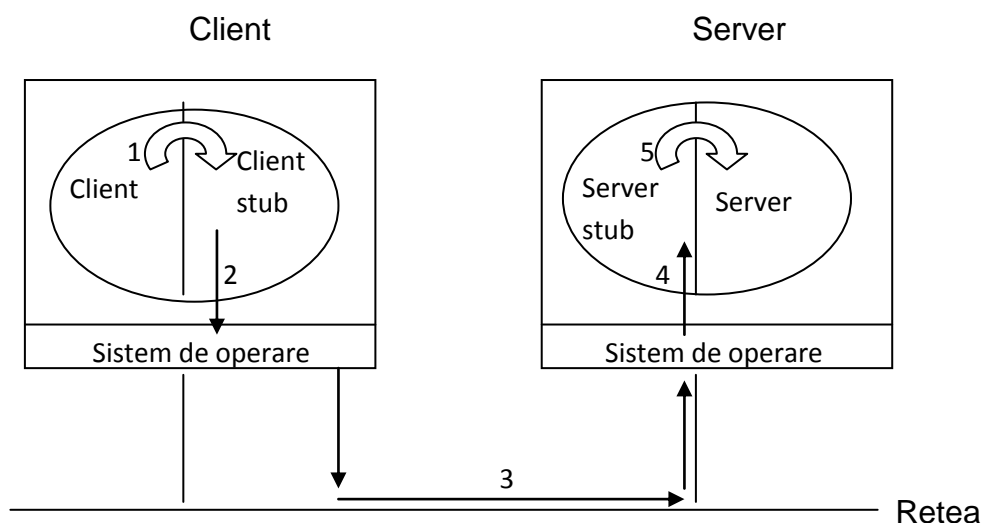
Intr-un fel, trimiterea unui mesaj cãtre un host și primirea unui rãapuns seamana cu apelarea unei funcții intr-un limbaj de programare. În ambele cazuri se porneste cu unul sau mai multi parametri și se returneaza un rezultat. Asemnarea aceasta se transpune în mediul de rețea sub forma apelurilor de procedura. Aceasta formulare usureaza mult programarea aplicatiilor. Exemplu: fie procedura `iaAdresaIP(numeHost)`, care funcționeaza prin trimiterea unui pachet UDP cãtre

serverul DNS și aștepta răspuns, cronometrând și încercând din nou, dacă acesta nu sosește în timp util. Astfel, toate detaliile rețelei sunt ascunse pentru programator.

Primii care au folosit acest concept au fost Nelson și Birrell, în 1984. Ei au venit cu ideea de a permite programelor să inițieze apeluri la proceduri aflate pe alte host-uri. Când un proces de pe masina A apelează o procedură de pe masina B, procesul apelant se suspendă și execuția procedurii are loc pe masina B. Informația este transportată de la apelant la apelat sub forma de parametrii, și se întoarce ca rezultat. Pentru programator, acest schimb de informație nu este vizibil. Această tehnică este cunoscută sub numele de RPC – Remote Call Procedure, și a devenit baza pentru multe din aplicațiile de lucru în rețea. Apelantul procedurii este denumit client; apelatul se numește server.

Scopul RPC este acela de a face un apel de procedură să pară ca unul local. În cea mai simplă formă, pentru a apela o procedură aflată în altă parte, programul client trebuie să dețină o bibliotecă, numită *client stub*, care reprezintă procedura serverului în spațiul de adresă al clientului. Similar, serverul are și el partea lui: *server stub*. Aceste proceduri ascund faptul că apelul de la client către server nu este local.

În figura următoare sunt prezentate etapele de creare a RPC. Pasul 1: clientul apelează stub-ul lui. Aici se apelează o procedură locală, cu parametrii introduși în stivă în modul obișnuit. În etapa a doua, stub-ul clientului comprimă parametrii într-un mesaj și face apel la sistemul de operare pentru a trimite pachetul. Această împachetare se numește „marshaling”.



În etapa a treia, sistemul de operare trimite mesajul de pe masina clientului către server. Pasul patru: sistemul de operare înmânează pachetul către stub-ul serverului, iar la etapa cinci, stub-ul serverului cheama procedura, după ce, în prealabil, a decomprimat parametrii. Răspunsul folosește aceleași etape, în sens invers.

De observat: procedura client, scrisa de programator, face un apel simplu de procedura, către stub-ul client, care are același nume ca procedura de pe server. Din moment ce procedura clientului și stub-ul lui se afla în același spațiu de adrese, parametrii sunt dati în modul obișnuit. În același mod, procedura de pe server este apelata de o procedura din același spațiu de adresa, cu parametrii așteptati. Astfel, în locul operațiilor de intrare-iesire facute pe socket-uri, comunicarea prin rețea se face prin mimarea unui apel de procedura simplu.

În ciuda simplitatii, există probleme, una din cele mai mari fiind legata de folosirea pointer-ilor ca parametrii. În mod normal, un parametru pointer dat unei proceduri nu este o problema. Procedura apelata poate folosi pointerul în aceeași modalitate ca apelantul, pentru ca ambele există în același spațiu virtual de adresa. Însă, în cazul RPC, folosirea pointer-ilor ca parametrii este imposibila, pentru ca serverul și clientul se afla în spații total diferite.

Totuși, există situații în care ei pot fi utilizați. Fie primul parametru un pointer către o variabila întreaga "a". Stub-ul client poate impacheta pe a și trimite către server. Stub-ul server-ului creaza un pointer către a și îl da procedurii de pe server. Când procedura reda controlul stub-ului, acesta trimite a înapoi la client, unde noul a este copiat peste cel vechi, în cazul în care a fost modificat de server. Pe scurt, secvența de apel prin referință a fost înlocuita cu secvența de apelare prin copie și înlocuire. Din păcate, aceasta metoda nu funcționează întotdeauna. De exemplu, nu se poate aplica la un pointer care indica spre structuri complexe de date. De aceea se impun anumite restricții.

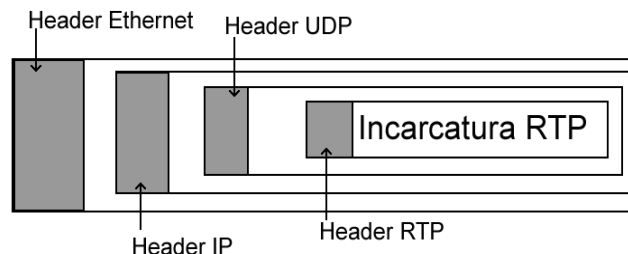
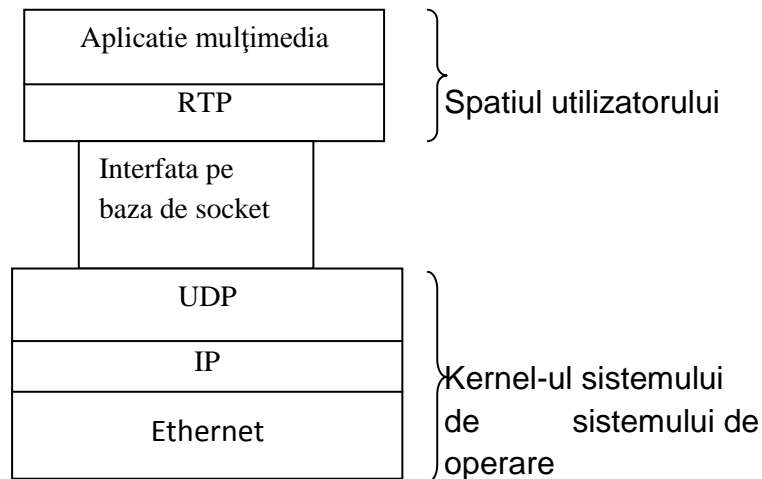
O alta problema este aceea ca în unele limbaje, ca C, este posibil și nu constituie o eroare, calculul produsului unor vectori, fără specificarea prealabilă a dimensiunii lor. Fiecare era terminat de o valoare specială cunoscută procedurilor. În cazurile acestea, stub-ul client nu poate comprima parametrii pentru ca nu știe exact cât de mari sunt.

Încă un obstacol: nu se poate deduce întotdeauna tipul de parametrii, nici după o specificare formală a lor în cod. Exemplu: *printf* care poate avea orice număr de parametrii, și care pot fi în orice ordine și de aproape orice tip.

A patra problema se referă la utilizarea variabilelor globale. În mod normal, procedurile ar putea comunica folosind variabile globale, dar dacă procedura apelată se afla la distanță, codul nu va reuși să le acceseze. [11]

3. Protocoale de transport de timp real

RPC între client și server este una din aplicațiile în care este foarte folosit UDP-ul. Alta utilizare: aplicații multimedia de timp real. Radio-ul prin internet, telefonica prin internet, programele *music-on-demand*, video-conferințele și altele au devenit extreme de populare și mai toate folosesc în principiu același protocol de transport în timp real. Acesta este RTP (Real Time Protocol) . Este descris în RFC3550 și este foarte răspândit printre aplicațiile multimedia. Transporta informația audio și video în pachete, și procesarea are loc în mare parte la receptor. Locul lor în stiva de protocoale se observă în figura următoare:



RTP funcționează în spațiul utilizatorului deasupra UDP-ului (în sistemul de operare). Aplicația multimedie lucrează cu mai multe stream-uri audio, video, text și poate și altele. Acestea sunt introduse apoi în biblioteca RTP, care se află în spațiul utilizatorului, împreună cu aplicația. Biblioteca multiplexează stream-urile și le codează în pachete RTP, pe care le trimite la un socket. Pe partea sistemului de operare, se generează pachete UDP care le include pe cele RTP și sunt trimise stratului IP pentru transmisie prin Ethernet. La receptor are loc procesul invers. Aplicația multimedie primește informația de la biblioteca RTP, și va reda informația media. Stiva de protocoale este în figura de mai sus, în partea stângă. În partea dreaptă este ilustrată imbricarea pachetelor. O consecință a acestei proiectări este dificultatea de a clasa RTP într-un anumit strat. Rulează în spațiul utilizatorului și este legat de aplicație, deci pare a fi un protocol de aplicație. Este totuși un protocol generic, independent de aplicație, care oferă doar facilități de transport, deci se încadrează în categoria protocoalelor de transport. O descriere mai bună: este un protocol de transport implementat în stratul de aplicație.

Funcția de bază a RTP este de a multiplexa mai multe stream-uri de timp real într-un singur stream de pachete UDP. Acesta din urmă poate fi trimis către o singură destinație (unicast) sau mai multe (multicasting). Pentru că RTP folosește UDP simplu, pachetele lui nu sunt tratate special de rutare, decât dacă niste caracteristici specifice de QoS (Quality Of Service) ale IP sunt activate. Nu există garanții speciale ale livrării, deci pachetele pot fi pierdute, corupte, întârziate etc. [11]

Totuși, RTP contine niste caracteristici specific care ajuta receptorii să proceseze informația multimediei. Fiecarui pachet trimis printr-un stream RTP i se aloca un număr cu 1 mai mare decat predecesorului sau. Aceasta număratoare permite masinii destinație să determine lipsa unuia sau mai multor pachete. Dacă unul lipseste, aplicatia decide ce se face. Poate omite un cadru video dacă acesta este tipul de informație transportat, sau poate interpola o valoare lipsa dacă se transporta informație audio. Retransmisia nu este o soluție practica, pentru ca, cel mai probabil, ar sosi prea tarziu pentru a mai fi utilizat. Consecinta acestui fapt: RTP nu are confirmare și nici un mechanism de cerere de retransmisie.

Un pachet RTP poate contine mai multe esantioane, care pot fi codate în orice fel (depend de aplicație). Pentru a putea conlucra (emițătorul cu receptorul), RTP defineste cateva profile (de exemplu: stream audio unic), și pentru fiecare din ele, se accepta mai multe formate de codare. Exemplu: un stream audio unic, poate fi codat PCM cu 8 biti (Pulse Code Modulation), esantionat la 8 KHz, codare predictiva, codare GSM, MP3 etc. RTP are un camp în header în care sursa poate specifica tipul de codare.

Alta facilitate: marcarea timpului (timestamping). În esenta, sursa asociaza un timestamp primului esantion din fiecare pachet. Timestamp-urile sunt relative la inceputul stream-ului, deci sunt mai importante diferentele între ele (valorile absolute sunt nesemnificative). Acest mechanism permite receptorului să plaseze esantioanele într-o memorie tampon (buffering) pentru a reda fiecare esantion un anumit interval de timp (atat cat trebuie pentru a putea fi perceput corect), indiferent dacă pachetul care trebuie nu a sosit inca.

Timestamping-ul nu numai ca reduce efectele variatiilor intarzierilor în rețea, dar permite și sincronizarea mai multor stream-uri între ele. Exemplu: un program de televiziune digitala, are un stream video și două stream-uri audio (pentru sunet stereo sau unul pentru coloana sonora originala și unul pentru dublare). Fiecare stream vine de la un dispozitiv fizic diferit, dar dacă sunt etichetate de un singur contor, pot fi redade sincronizat, chiar dacă sunt transmise dezorganizat.

În figura urmatoare este reprezentat headerul RTP. Contine 3 cuvinte de 32 de biti (si eventual niste extensii). Primul cuvânt contine câmpul pentru versiune (s-a ajuns la versiunea 2). Bitul P indica dacă pachetul a fost bordat pana la multiplu de 4 octeti. Ultimul octet de bordare arata cati octeti au fost adaugati. Bitul X indica existenta unei extensii. Formatul și semnificatia extensiei nu sunt definite. Se specifica doar lungimea, în primul cuvânt al extensiei.

Versiune	P	X	CC	M	Tipul de incarcatura	Numarul secventei
Timestamp						
Identificatorul sursei de sincronizare						
.....						
.....						
Identificatorii surselor contribuitoare						

Câmpul CC menționează câte surse contribuitoare există (între 0 și 15). Bitul M este un marker specific aplicației. Poate fi folosit pentru a marca începutul unui cadru video, începutul unui cuvânt într-un canal audio, sau orice altceva caracteristic aplicației.

Câmpul *Tip de încaratura* specifică ce algoritm de codare este folosit. Din moment ce fiecare pachet conține acest câmp, codarea se poate schimba în timpul transmisiei.

Numărul de secvență este un contor care se incrementează la fiecare pachet trimis. Este folosit pentru detectarea pierderii de pachete.

Timestamp-ul este creat de de sursa stream-ului pentru a reduce variația de întârziere (jitter) la recepție, prin decuplarea dintre redare și timpul de sosire al pachetului.

Identificatorul sursei de sincronizare spune cui stream aparține acel pachet. Este folosit pentru a multiplexa și demultiplexa mai multe stream-uri de date într-un / dintr-un singur stream de pachete UDP. Identificatorii surselor contribuitoare, dacă există, sunt folosiți când în studio-ul de înregistrare se folosesc mixere. În acest caz, mixer-ul este sursa de sincronizare, și stream-urile mixate sunt listate aici. [11]

3.1 RTCP

RTCP este un protocol care se ocupă numai de răspunsuri, sincronizări și interfață cu utilizatorul, dar nu transportă deloc date, și este strans legat de RTP. Prima funcție oferă feedback către surse în legătură cu întârzierile, lățimea de bandă, congestia și alte proprietăți ale rețelei. Informația furnizată de acest serviciu poate fi folosită de procesul de codare pentru a crește rata de date (și de a crește calitatea) atunci când rețeaua funcționează foarte bine, sau pentru a o reduce când există probleme în rețea. Primind răspunsuri continue, algoritmi de codare pot fi adaptați mereu pentru a oferi cea mai bună calitate posibilă, în limitele performanțelor curente ale rețelei. Exemplu: dacă lățimea de bandă crește (sau descrește) în timpul transmisiei, se poate folosi codare MP3, PCM pe 8 biți, sau codare Delta. Câmpul „tip de încaratura” este folosit pentru a indica mașinii destinație ce algoritm de codare s-a folosit pentru pachetul curent (de aceea se poate schimba codarea în timpul transmisiei).

O problemă a mecanismului de feedback este aceea că rapoartele se trimit către toți participanții la comunicare. Pentru o aplicație multicast, cu un număr mare de calculatoare, lățimea de bandă folosită de RTCP crește puternic. Pentru a preveni ocuparea bandei cu rapoarte, rata lor este micșorată, undeva sub 10% din banda media. Fiecare participant trebuie să știe dimensiunea benzii, pe care o poate calcula cu ajutorul emițătorului și a numărului de participanți, pe care îl poate determina ascultând porturile RTCP.

RTCP se ocupă și de sincronizarea între stream-uri. Diferite stream-uri pot folosi diferite surse de tact („ceasuri” diferite), cu granularități și derivate diferite. RTCP poate să le sincronizeze, în ciuda acestor impedimente. [11]

3.1.1 Pachetele RTCP

Sunt specificate mai multe tipuri de pachete pentru a transporta diverse informații de control:

- SR (*Sender Report*): contine statistici de transmisie și recepție de la participantii care emit activ.
- RR (*Receiver Report*): contine statistici de la participantii care nu emit activ
- SDES: descrie obiectele sursa, inclusiv CNAME
- BYE: indica sfârșitul participării unei mașini
- APP: pentru funcții specifice aplicației.

Fiecare pachet RTCP începe cu o secțiune fixă, similară pachetelor de date ale RTP. Urmează elemente structurate, care pot avea lungime variabilă, un funcție de tipul pachetului, însă trebuie să se limiteze pe o frontieră de 32 de biți. Această cerință de aliniere și câmpul „Lungime” din partea fixă a fiecărui pachet sunt necesare pentru a se putea „stivui” pachetele. Mai multe pachete RTCP pot fi concatenate, fără a fi nevoie de separatoare, pentru a forma un pachet compus care poate fi trimis într-un unic pachet al protocolului inferior (ca UDP). Nu se numără pachetele individuale RTCP din pachetul compus pentru că protocoalele din nivelele inferioare țin cont de lungimea totală și ele determină sfârșitul pachetului compus.

Fiecare pachet individual din pachetul compus poate fi procesat independent, fără a se impune o anumită ordine. Totuși, pentru a se putea folosi eficient funcțiile protocolului, se impun anumite constrângeri:

- Statisticile de recepție (în SR sau RR) ar trebui trimise cât mai des (în limita constrângerilor lățimii de bandă) pentru a maximiza rezoluția statisticilor, deci, periodic, unul din pachetele compuse trebuie să conțină și un pachet de raportare
- Participantii noi care primesc date trebuie să primească și CNAME-ul pentru o sursă cât mai repede posibil pentru a identifica sursa și pentru a asocia informația cu anumite scopuri (ca sincronizarea), deci fiecare pachet compus RTCP trebuie să includă și SDES CNAME.

În concluzie, toate pachetele RTCP trebuie trimise într-un pachet compus care conține cel puțin 2 pachete individuale, cu formatul următor:

- Prefix de criptare: dacă pachetul compus trebuie securizat, trebuie prefixat cu un număr pe 32 de biți, aleator, pentru fiecare pachet compus. Dacă este necesară bordarea pentru criptare, trebuie adăugată la ultimul pachet din ansamblu.
- SR sau RR: primul pachet din ansamblu trebuie să fie un pachet de raport pentru a facilita validarea headerului, chiar dacă nu există date de transferat (caz în care trebuie trimis un RR gol) și chiar dacă și celălalt pachet din ansamblu este de tip BYE.

- RR aditionale: dacă numărul de surse pentru care se raporteaza statistice depaseste 31 (adică nu mai incapa intr-un sungur pachet RR sau SR), atunci primul pachet de RR/SR trebuie urmat de un pachet aditional (sau mai multe, în funcție de numărul de surse).
- SDES: un pachet SDES care contine un obiect CNAME trebuie inclus în fiecare ansamblu RTCP. Optional pot fi incluse și alte surse dacă sunt cerute de o anumita aplicație, în limita constrangerilor legate de latimea de banda.
- BYE sau APP: alte tipuri de pachete RTCP pot aparea, în orice ordine, dar BYE trebuie să fie ultimul pachet trimis.

Un participant RTP trebuie să trimită doar un singur pachet compus RTCP pe intervalul de raportare, pentru ca estimarile privind banda pentru fiecare participant să fie corecte. Exceptie: când pachetul compus este impartit pentru criptare partiala. Dacă sunt prea multe surse și nu incap toate pachetele RR intr-un singur ansamblu fără depasirea MTU (maximum transmission unit), atunci doar subsetul care incapa intr-un MTU ar trebui inclus în fiecare interval. Subseturile trebuie selectate folosind metoda round-robin, astfel incat toate sursele să fie raportate. [4]

v=2	p	Nr RR	tip de pachet	Lungime mesaj
SSRC al raportului emitator				
Timestamp NTP(2 cuvinte de 32 biti)				
Timestamp RTP				
Contor cumulativ de pachete al emitatorului				
Contor cumulativ de octeti al emitatorului				
Blocul 1 al raportului de receptie				
Blocul 2 al raportului de receptie				
.....				

Pachet RTCP (cu header-ul inclus) [5].

3.1.2 Intervalul de transmisie

RTP este proiectat pentru a permite scalabilitate unei aplicatii, de la sesiuni de dimensiuni mici la sesiuni de mii de participanti. Exemplu: intr-o conferinta, traficul se limiteaza automat, pentru ca nu pot vorbi mai mult de 2-3 oameni deodata (nu ar intelege nimeni nimic). Deci, rata de date, la multicasting, este relativ constanta independent de numărul de participanti. Însă, traficul de control nu se limiteaza de la sine. Dacă rapoartele de receptie de la fiecare participant ar fi trimise cu o rata constanta, traficul de control ar creste liniar cu numărul de participanti. Deci rata trebuie scazuta și trebuie calculat dinamic un interval între pachetele transmise.

Presupunem ca există o limita a traficului de date, numit „latimea de banda a sesiunii”, care se divide între participanti. Aceasta se poate alege în funcție de un

cost sau de niste informații apriori despre latimea disponibila sesiunii. Deseori, latimea este egala cu suma latimilor nominale ale fiecarui emițător activ.

Latimea sesiunii o da de obicei o aplicație de management a sesiunii, atunci când invoca aplicatia media, insa aplicatiile media pot seta o valoare implicita, pe baza latimii pentru singur utilizator, pentru codarea selectata pentru sesiune. Aplicatia poate impune și limite la banda regulilo de multicast sau a altor criterii, insa toti participanti trebuie să aiba aceeași caloare pentru latimea sesiunii, astfel incat să se poata calcula un singur interval RTCP pentru toata lumea.

Calculul latimii de banda pentru traficul de date și de control include și protocoalele inferioare (ca IP și UDP) (ar trebui să le știe sistemul de gestiune al resurselor). Este de asteptat ca aplicatia să știe care din aceste protocoale este folosit. Header-ele de legatura nu sunt incluse în calcul, din moment ce pachetele vor fi incapsulate cu diferite header-e pe masura ce calatoresc.

Traficul de control trebuie limitat la o fractiune cunoscuta din latimea de banda a sesiunii (cat mai mica, dar suficienta pentru a nu impiedica sarcina de a transfera informații a protocolului de transport). De asemenea, trebuie să o cunoasca toti partiipantii pentru ca fiecare să își poata calcula independent partea lui de banda. Se recomanda ca fractiunea de latime de banda pentru RTCP să fie fixata la 5%. De asemenea, se recomanda ca o patrime din latimea pentru RTCP să fie dedicată participantilor care trimit date, astfel incat, în sesiunile cu număr mare de receptori, dar cu putini emițători, cei care vor să între în sesiune să poataprimi cat mai repede CNAME-ul pentru emițători. Desi valorile acestor constante nu sunt critice, este esential ca toti participantii la sesiune să folosească aceleasi valori, astfel incat să se calculeze același interval. Deci constantele acestea trebuie fixate pentru un anumit profil. Un profil poate să specifice ca latimea de banda asociata traficului de control să fie un parametru separat al sesiuni, și nu un procentaj strict din banda sesiunii. Folosirea unui parametru separat permite aplicațiilor cu rata adaptabila să sa seteze o latime care să corespunsa cu una tipica pentru datele trimise, și care e mai mica decat latimea maxima specificata de parametrul sesiunii.

Un profil mai poate specifica dacă latimea pentru traficul de control este divizata în doi parametri separati pentru participantii activi și pentru cei pasivi. Fie acesti parametrii S și R. Urmarind recomandara ca $\frac{1}{4}$ din banda să fie alocata emițătorilor, valorile recomandate pentru acesti parametrii sun de 1.25% , respectiv 3.75%.

Folosirea ambilor parametrii permite oprirea rapoartelor RTCP, prin setarea latimii de banda pentru cei care nu emit date la 0. Însă oprirea rapoartelor de receptie RTCP nu este recomandata, pentru ca sun necesare pentru feedback-ul calitatii și controlul congestiei. Totuși, poate fi avantajoasa pentru sistema care opereaza pe legături unidirectionale sau sesiuni care n-au nevoie de feedback-ul calitatii receptiei.[4]

4. DCCP – Datagram Congestion Control Protocol

4.1 Introducere

DCCP este un protocol al nivelului de transport orientat pe mesaje. Implementează setarea unei conexiuni sigure, închiderea ei, ECN (Explicit Congestion Notification), control de congestie, și negociere de caracteristici. DCCP oferă o care de accesare a mecanismelor de control al congestiei, fără a fi necesară implementarea lor în nivelul de aplicație. Permite fluxuri similare TCP, dar nu asigură și livrarea în ordinea transmisiei. Livrarea secvențială a mai multor fluxuri (ca în SCTP – Stream Control Transmission Protocol) nu este disponibilă în DCCP.

DCCP este utilizat în aplicații în care se impun constrângeri temporale asupra livrării pachetelor. Din această categorie fac parte jocurile online multiplayer, telefonie prin internet, streaming-ul de informații media (video, audio). Caracteristica esențială a acestor aplicații este aceea că mesajele vechi devin rapid expirate, își pierd utilitatea. Prioritate mai mare o au mesajele noi, de aceea nu este utilă încercarea de retransmitere a pachetelor (ar consuma timp și resurse de rețea inutile).

De asemenea, DCCP poate fi folosit ca mecanism general de control al congestiei pentru aplicații care au la bază protocolul UDP. Se poate adăuga și un mecanism pentru siguranță, eventual, unul pentru livrarea pachetelor în ordinea transmisiei. În acest context, DCCP permite utilizarea diferitelor mecanisme de control al congestiei, în general a celor TCP-friendly.

O conexiune DCCP conține atât trafic de confirmare, cât și trafic de date. Confirmările anunță emițătorul că pachetele lui au sosit la destinație sau dacă au fost marcate de ECN. Confirmările sunt transmise cu gradul de siguranță pe care îl cere mecanismul de control al congestiei. Este posibilă atingerea unui grad de 100% al siguranței.

4.2 Congestie

Am tot menționat cuvântul „congestie”, însă nu a fost definit clar.

Congestia rețelei este fenomenul de deteriorare a calitatii serviciilor (QoS – quality of service) produs de supraîncărcarea nodurilor rețelei, deci termenul este asociat mai ales rețelelor de dimensiuni mari, în care se vehiculează cantități importante de date.

Congestia are mai multe cauze: fie router-urile nu sunt suficient de rapide (procesurile lor sunt prea lente, și nu reușesc să golească cozile de așteptare în timp util, chiar dacă rețeaua este suficient de liberă), fie se transmit mai multe stream-uri care trebuie apoi trimise pe aceeași linie de ieșire (se adună iar cozi), fie buffer-urile nu sunt suficient de mari și se pierd pachete. În cazul unui trafic foarte mare, situația se poate agrava atât de tare încât nu mai este livrat nici un pachet. [7]

4.3 Caracteristici DCCP

DCCP are următoarele caracteristici:

- Un flux nesigur de datagrame, cu confirmare
- Negociere sigura de optiuni, inclusiv negocierea unui mecanism potrivit pentru controlul congestiei
- Protocol de tip handshake sigur pentru inițializarea și închiderea conexiunii
- Descoperirea unitatii maxime transmisibile pe calea aleasa (MTU = Maximum Transmission Unit)
- Mecanisme care permit serverelor să evite pastrarea de stari pentru încercări de realizare deconexiuni, neconfirmate, sau pentru conexiuni deja inchise.
- Control de congestie, inclusiv ECN și ECN NONCE (...)
- Mecanisme de confirmare care comunica pierderea pachetelor și informații ECN.
- Mecanisme optionale care comunica aplicației emitatoare, cu grad mare de siguranță, care pachete au ajuns la receptor și dacă ele au fost marcate de ECN, corupte sau eliminate în buffer-ul receptorului.
- Alegerea mecanismelor modulare de control al congestiei. În momentul de față, sunt specificate două metode: TCP-like Congestion Control (control de congestie asemanator TCP) și TCP-Friendly Congestion Control. [10]

Intregul scop al DCCP este acela de a oferi o metoda standard de a implementa controlul congestiei simplu și negocierea lui pentru aplicatiile de timp real. Un motiv pentru care DCCP este des folosit este acela ca permite utilizarea ECN, impreuna cu controlul congestiei end-to-end, pentru aplicații care altfel ar folosi UDP (care nu este orientat pe conexiune și nu are nici un mecanism de control de congestie, deci aplicatia care îl foloseste ar trebui să își implementeze propriul sistem de control). A fost proiectat astfel incat să genereze cat mai putine date de overhead, legat atat de dimensiunea antetelor pachetelor, cat și de stare și overhead CPU necesar la capatul de receptie.

Pe langa acestea, DCCP implementează și inițializarea unei conexiuni sigure, incheierea ei și negocierea de caracteristici.

Alta motivatie a folosirii DCCP e aceea ca a fost proiectat astfel incat să permita aplicațiilor să aleaga o alternativa la mecanismul curent de control al congestiei (TCP-Style Congestion Control) care injumatateste fereastra de congestie, ca răapuns la aparitia unei indicatii de congestie. DCCP permite aplicațiilor să aleaga între mai multe forme de control. Una din ele este controlul similar TCP (TCP-like Congestion Control), și injumatateste fereastra de congestieca răapuns la eliminarea sau marcarea unui pachet (ca în TCP). O alta alternativa este TFRC (TCP Friendly Rate Control). Este o forma de control bazată pe ecuatii, care minimizeaza schimbarile bruste ale ratei de trimitere.

4.4 Conexiunea DCCP

Orice conexiune DCCP ruleaza între două capete. Le voi numi DCCP A și DCCP B. Datele circula în ambele directii, folosind 4 tipuri de pachete:

1. pachete de la A la B
2. confirmări de la B la A
3. pachete de la B la A
4. confirmări de la A la B

Fiecare flux de pachete de mai sus reprezinta un set. Voi defini niste termeni pe care ii folosesc în continuare:

- sub-fluxuri = un subflux este un pachet fie de date, fie de confirmare, trimis într-o singura directie. Fiecare din cele 4 seturi de mai sus reprezinta un subflux (subfluxurile se pot intercala într-un anumit grad, deoarece confirmările pot fi trimise imediat după pachetele cu date).
- secvente = o secventa consta în totalitatea pachetelor trimise într-o singura directie, indiferent dacă sunt date sau confirmări. Seturile 1 cu 4, și 2 cu 3 de mai sus sunt secvente diferite. Fiecare pachet dintr-o secventa are un alt număr de secventa.
- Emitator/receptor pe jumătate de conexiune (half-connection) = în contextul unei conexiuni injumătate, emițătorul e cel care trimite date, iar receptorul le primește, trimițând înapoi confirmări. Exemplu: în semi-conexiunea A-B, DCCP A este emițătorul, iar DCCP B este receptorul.

Fiecare semiconexiune este gestionata de un mecanism de control al congestiei. Capetele negociaza aceste mecanisme la inițializarea conexiunii (cele două mecanisme nu trebuie să fie neaparat identice). Mecanismele conformante de control al congestiei corespund identificatorilor de byte, sau CCID-uri. CCID pentru o semiconexiune descrie cum emițătorul unei semiconexiuni limiteaza rata pachetelor de date, cum mentine parametrii necesari (ca ferestrele de congestie), cum receptorul trimite feedback-ul de congestie prin confirmări și cum gestionează rata confirmărilor.

Fiecare conexiune DCCP este initiata activ de unul din participanti, care se leaga la un socket DCCP aflat în stare pasiva de ascultare. Capatul activ e va numi client, iar cel pasiv: server. Majoritatea specificatiilor pentru DCCP sunt identice și pentru server și pentru client, inasa doar serverul poate cere inchiderea unei conexiuni. Deci clientul nu poate forta serverul să mentina status-ul unei conexiuni după ce aceasta a fost inchisa.

DCCP nu suporta deschidere simultana, ca TCP. Asta inseamna ca un host nu trebuie să raspundă unei cereri DCCP decat dacă portul destinație specificat în cerere corespunde unui socket local deschis pentru ascultare. DCCP nu poate lucra nici cu conexiuni semideschise (adică inchide ambele semiconexiuni ca pe o unitate întreaga).

DCCP foloseste mecanisme generice pentru negocierea proprietatilor conexiunii, ca CCID-urile active pe semiconexiuni. Un nume de proprietate, ca „CCID” este asociat, în general, la două caracteristici, una pentru fiecare semiconexiune. De exemplu, există 2 CCID-uri pentru o conexiune. Capatul care comanda o anumita caracteristica se numeste: locatia caracteristicii. Valorile caracteristicilor sunt negociate de optiunile: „Change”, „Confirm” și „Prefer”. „Change” (schimba) este trimis către locatia unei caracteristici pentru a cere schimbarea valorii pentru acea caracteristica. Locatia poate raspunde cu „Prefer” (preferinta), care cere schimbarea valorii de la capatul care a trimis cererea de schimbare, sau își poate schimba el insusi valoarea caracteristicii, raspunzand apoi cu „Confirm” (confirmare). Negocierea de caracteristici este una fiabila din cauza retransmisiilor.

4.5. Pachetele DCCP

DCCP foloseste noua tipuri de pachete:

1. DCCP-Cerere
2. DCCP-Raspuns
3. DCCP-Date
4. DCCP-Confirmare
5. DCCP-Confirmare de date
6. DCCP-Cerere de inchidere
7. DCCP-Inchidere
8. DCCP-Reset
9. DCCP-Transfer

Descrierea informaționala conexiunii:

1. Clientul trimite serverului un pachet DCCPCerere în care specifica porturile clientului și ale serverului, serviciul cerut și orice alte caracteristici care se negociaza, inclusiv CCID-ul pe care clientul vrea să îl folosească serverul. Clientul poate trimite în spatele acestui pachet și niste date, ca o cerere la nivelul aplicației, pe care serverul poate să o ignore.
2. Serverul trimite clientului un pachet DCCP-Raspuns, indicand ca va comunica cu clientul. Raspunsul indica optiunile și caracteristicile cu care serverul este de acord, incepe (sau continua) alte negocierii de caracteristici (dacă este necesar) si, optional, include și o procedura de inițializare de cookie-uri care impacheteaza toate aceste informații și care trebuie trimise inapoi la client pentru a completa conexiunea.
3. Clientul trimite serverului un pachet DCCP-Confirmare care confirma pachetul DCCP-Raspuns. Acesta confirma numărul initial de secventa al serverului și trimite inapoi inițializarea de cookie, dacă a existat vreuna un pachetul de răspuns. Poate include și negociere de caracteristici.
4. Pot urma (sau nu) mai multe pachete de confirmare pentru a finaliza negocierea de caracteristici. Clientul poate trimite și o cerere către aplicație în

spatele pachetului de confirmare, producand astfel un pachet DCCP-Confirmare de date.

5. În aceasta etapa, serverul și clientul fac schimb de pachete de date, pachete de confirmare, si, optional, pachete de confirmare de date, care contin date și confirmări. Dacă clientul nu are date de trimis, stunci serverul va trimite pachete DCCP-Date și DCCP-confirmare de date, în timp ce clientul va trimite numai pachete de confirmare.
6. Serverul trimite un pachet DCCP-cerere de inchidere (cerand inchiderea conexiunii). (numai serverul poate cere inchiderea conexiunii !).
7. Clientul trimite un pachet DCCP-Inchidere pentru a confirma inchiderea.
8. Serverul trimite un pachet DCCP-Reset și sterge starea conexiunii.
9. Clientul primește pachetul de RESET și pastreaza starea conexiunii un interval rezonabil de timp pentru a permite pachetelor ramase să soseasca înainte de a inchide conexiunea.

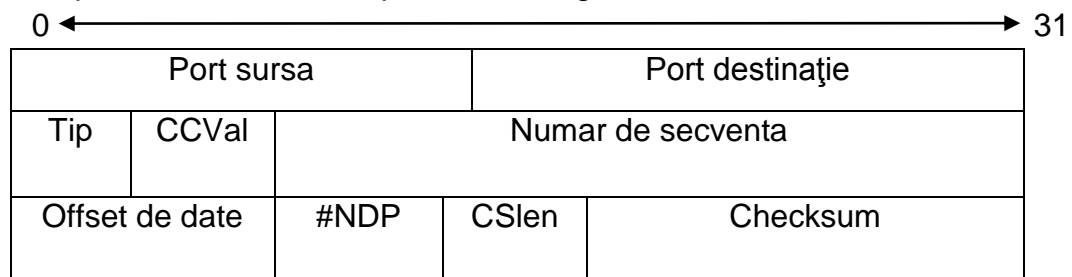
O alternativa la secventa de inchidere de mai sus este cea în care inchiderea este initiata de client:

6. Clientul trimite un pachet DCCP-Inchidere pentru a cere inchiderea conexiunii.
7. Serverul trimite un pachet DCCP-Reset cu câmpul „Scop” (campurile sunt tratate în paragrafele urmatoare) setat pe „Inchidere” și sterge starea conexiunii.
8. Clientul primește pachetul de RESET și pastreaza starea conexiunii un interval rezonabil de timp pentru a permite pachetelor ramase să soseasca înainte de a inchide conexiunea.

Aceasta metoda de comunicare cu handshake pentru initierea și inchiderea conexiunii permite serverului să refuze pastrarea oricarei stari pana când are loc handshake-ul cu clientul, ceea ce asigură ca clientul trebuie să tine starea de asteptare la inchiderea conexiunii.

4.6 Header-ul

Toate pachetele DCCP incep cu un antet generic:



Porturile sursa și destinație: reprezentate pe 16 biti fiecare, seamana cu campurile omoloage din TCP și UDP. Portul sursa reprezinta portul relevant la capatul care care trimite pachete, iar destinația este portul care asculta.

TIP: 4 biti, specifica tipul de mesaj DCCP. Sunt definite următoarele valori:

0. Pachet de cerere
1. Pachet de răspuns
2. Pachet de date
3. Pachet de confirmare(ACK)
4. Pachet de confirmare de date (Data Ack)
5. Pachet de cerere de închidere
6. Pachet de închidere
7. Pachet de reset
8. Pachet de transfer

CCVal: 4 biti, rezervat pentru trimiterea CCID-ului. CCID-ul emițător de la A la B, care este activ la DCCP A, poate trimite informații receptorului de la DCCP B, codând informația în CCVal.

Numarul de secventa: 24 biti. Acest camp este initializat de pachetul de cerere sau de răspuns și incrementat cu 1 la fiecare pachet trimis. Eceptorul foloseste aceasta informație pentru a determina dacă s-au pierdut pachete pe drum. Chiar și pachetele care nu contin date actualizeaza numărul de secventa. Aceste numere oferă și o protecție împotriva pachetelor vechi sau a celor trimise de programe malware. DCCP de rata mare poate necesita protecție împotriva numerelor de secventa impachetate. Exemplu: un flux de 10Gb/s cu pachete de 1500 octeti va trimite 2^{24} pachete în ~ 20 de secunde. În termenii calatoriei prin rețea, este un timp destul de mare, dar tot există riscuri. Numerele de secventa initiale ale celor două subfluxuri sunt setate de primele pachete de cerere, respectiv e răspuns trimise, și ar trebui alese ca pentru TCP. Un număr initial de secventa trebuie să includa o componenta aleatoare (sau pseudo-aleatoare) care să ingreuneze munca atacatorilor de a interveni în umerele de secventa. Numarul initial de secventa ales pentru un identificator de conexiune dat (adresa și portul sursa + adresa și portul destinație) ar trebui să creasca cu timpul, pentru a evita livrarea neadecvata a pachetelor vechi.

Offset-ul de date: 8 biti. Este offset-ul de la inceputul header-ului pana la inceputul incarcaturii utile a pachetului, masurata în cuvinte de 32 de biti.

#NDP (number of non-data packets): 4 biti, este numărul de pachete care nu contin date. DCCP seteaza acest camp cu numărul de pachete fără date trimise oana atunci în secventa să (modulo 16). Un pachet non-date este un pachet care nu contine informație de la utilizator (pachetele de confirmare, de închidere, de cerere de închidere și de reset sunt mereu fără date, în timp de pachetele de cerere, răspuns și de transfer pot avea data sau nu). Acest camp ajuta DCCP-ul receptor să decida dacă un pachet pierdut continea date de la utilizator. De exemplum pachetul 10 avea #NDP setat la 5. Pachetul 11 s-a pierdut, iad pachetul 12 are #NDP=5. Atunci receptorul își da seama ca pachetul 11 continea date, din moment ce nu s-a schimbat #NDP. Dacă acesta ar fi fost 6 pentru pachetul 12, pachetul nu ar fi continut deloc date.

CSLEN (Checksum length): 4 biti. Câmpul acesta specifică ce părți din pachet sunt acoperite de checksum (acesta acoperă cel puțin header-ul, opțiunile și un pseudoheader luat de la header-ul nivelului rețea). Dacă lungimea checksum-ului este 0, asta e tot ce acoperă. Dacă este 15, acoperă și încărcătura utilă a pachetului. Alte valori decât între 0 și 15 se consideră experimentale.

Checksum: 16 biti. DCCP folosește algoritmul de verificare de la TCP-IP. Câmpul de checksum este egal cu complementul lui 1 al complementului tuturor cuvintelor de 16 biti din header, opțiuni și din pseudoheader-ul luat de la headerul nivelului rețea, și, în funcție de lungimea specificată în CSLEN, o parte din încărcătura, poate chiar toată.

Pseudoheaderul este calculat tot ca pentru TCP. Pentru IPv4, are lungimea de 96 de biti și constă din adresele sursă și destinație, numărul de protocol IP pentru DCCP (bordat la stânga cu 8 biti de 0) și lungimea DCCP (pe 16 biti), cu tot cu header cu opțiuni + lungimea datelor.

Pachetele cu checksum invalid trebuie ignorate, și opțiunile lor nu trebuie procesate. [7]

5. TCP-Like Congestion Control

DCCP folosește identificatori de control al congestiei, sau CCID-uri (Congestion Control Identifiers) pentru a specifică ce mecanism de control al congestie se folosește pe o semiconexiune.

CCID-ul similar TCP (il voi prescurta TCP-L) trimite date folosind o variantă apropiată de mecanismul de control folosit de TCP, incluzând mai multe confirmări selective: SACK (Selective Acknowledgements). Această variantă de control al congestiei se numește CCID2, și este potrivită pentru emițătorii care se pot adapta creșterilor bruste din fereastra de congestie, tipică metodei AIMD (Additive increase, Multiplicative Decrease) a lui TCP, și este utilă pentru cei care vor să profite de lățimea de bandă într-un mediu în care condițiile se schimbă rapid.

Există două tipuri de pachete: pachete de date, pe care le trimite emițătorii, și pachete de confirmare (ack), trimise de receptori.

CCID2 este potrivit pentru fluxurile DCCP care au nevoie de cât mai multă lățime de bandă, pe termen lung. Fluxul trebuie să tolereze variații mari ale ratei de transmisie, caracteristica AIMD, incluzând și înjumătățirea ferestrei de congestie ca răspuns pentru un eveniment de aglomerare. Poate fi folosit de aplicații care trebuie să transfere cât mai multe date în cel mai scurt timp posibil.

În contrast cu CCID2, CCID3 (sau TFRC = TCp Friendly Rate Control) este potrivit pentru fluxuri în care se preferă minimizarea schimbărilor bruste ale ratei de transmisie. De exemplu, CCID2 este mai potrivit decât CCID3 pentru streaming multimedial care folosește buffere mari la recepție înainte de redare, izolând într-un fel aplicația de schimbările rapide ale ratei de transmisie. Astfel de aplicații ar trebui să aleagă CCID2 și în locul TCP-ului, adăugând opțional și în forma de asigurare a fiabilității.

Alt avantaj al lui CCID 2 este ca mecanismele lui de control al congestiei similare TCP sunt clar intelese, iar dinamica traficului seamana cu cea de la TCP. [1]

5.1 Relatia cu TCP

Diferentele între CCID2 și controlul direct al congestiei prin mecanismul TCP sunt:

1. CCID 2 aplica controlul la confirmări, un mecanism care nu este standardizat inca pentru TCP.
2. DCCP este un protocol cu datagrame, deci mai multi parametrii ale caror unitati de masura sunt octetii în TCP (ca fereastra de control) se masoara în pachete la DCCP.
3. Fiind un protocol nefiabil, DCCP nu retransmite niciodata un pachet, deci mecanismele de control al congestiei care disting pachetele trimise initial de cele retransmise au fost reprojectate pentru contextul DCCP.

5.2 Exemplu pentru semiconexiune

Exemplul urmator arata evolutia unei semiconexiuni folosind controlul de congestie al CCID2:

1. Emitatorul trimite pachete DCCP de date, numărul lor fiind controlat de o fereastra de congestie, ca în TCP. Fiecare pachet DCCP-Data are un număr de secventa. Emitatorul trimite și o caracteristica AckRatio (raport de confirmări) în care specifica numărul de pachete de date care trebuie acoperite de un pachet de confirmare de la receptor. Implicit, AckRatio este 2. câmpul CCVal din headerul DCCP este setat la 0. Presupunand ca semiconexiunea are capacitate ECN, fiecare pachet DCCP-Data se trimite cu capacitate ECN.
2. Receptorul trimite un pachet DCCP-Ack, confirmand pachetele de date în funcție de setarea AckRatio. Fiecare pachet de confirmare are un număr de secventa și contine vectorul de confirmare (AckVector). Numarul de secventa confirmat într-un pachet DCCP-Ack este acela al pachetului de date cu cel mai mare număr (nu este ca la TCP, cu confirmări cumulative). Receptorul returneaza suma anunturilor ECN primite prin optiunile vectorului de confirmare permitand emițătorului să verifice probabilistic faptul ca receptorul primeste corect. Pachetele DCCP-Ack de la receptor au și ele capacitate ECN, din moment ce emițătorul va controla rata confirmărilor într-o maniera TCP-Friendly, folosind caracteristica AckRatio.
3. Emitatorul continua să emita pachete DCCP-Data, controlate de fereastra de congestie. Dupa primirea pachetelor de confirmare, emițătorul examineaza vectorii lor de confirmare pentru a afla dacă s-au pierdut sau distrus pachete, și regleaza fereastra de congestie. Fiind un transfer nesigur, emițătorul nu va retrimite pachetele eliminate.

4. Pentru ca pachetele de confirmare folosesc numere de secvență, emițătorul are niste informații despre pachetele pierdute sau marcate, și raspunde la pierderi sau marcare prin reglarea raportului de confirmare (AckRatio) trimis receptorului.
5. Emitatorul confirma confirmările receptorului cel puțin odată într-o fereastră de congestie. Dacă ambele semiconexiuni sunt active, confirmarea emițătorului despre confirmările receptorului este inclusă în confirmarea emițătorului despre pachetele de date ale receptorului. În cazul în care calea inversă a semiconexiunii este pasivă, emițătorul trimite cel puțin un pachet DCCP-DataAck (confirmare de date) pe fereastra de congestie.

Emitatorul estimează timpul de dus-întors, fie prin urmărirea timpilor confirmărilor (asa cum face TCP), fie prin opțiunea explicită de timestamping, și calculează valoarea de expirare (TimeOut) similar cum se calculează timpul de retransmisie din TCP. TimeOut-ul determină când un pachet de date nou poate fi transmis, atunci când emițătorul este limitat de fereastra de congestie și nu există nici un feedback de la receptor. [1]

6. TFRC – TCP-Friendly Rate Control

TFRC este un mecanism de control al congestiei pentru fluxuri unicast. A fost proiectat pentru aplicații care folosesc o dimensiune fixă a pachetului și poate rula alături de TCP cu pachet fix fără a consuma o lățime de bandă prea mare. Are însă o variație mai mică a timpului de traversare a rețelei, comparativ cu TCP, ceea ce îl face mai adecvat pentru aplicații ca media streaming sau telefonie, unde o rată de transmisie constantă este importantă.

Trade-off-ul pentru rată constantă de transmisie constă în faptul că răspunde mai greu la schimbările de bandă disponibilă. Deci TFRC ar trebui folosit numai când aplicația are nevoie explicită de transmisie „lină”, evitând înjumătățirea ratei de transfer (ca la TCP), ca răspuns la pierderea unui pachet. Pentru aplicații care trebuie să transfere cât mai multe date în cel mai scurt timp posibil, se recomandă TCP, însă dacă nu este necesară fiabilitatea, se poate folosi și schema de control al congestiei AIMD (Additive Increase, Multiplicative decrease).

TFRC a fost proiectat pentru a da performanțe maxime aplicațiilor care folosesc dimensiune fixă a segmentului, și variază rata de trimitere a pachetelor ca răspuns la congestie. [1]

6.1 Mecanismele protocolului

Pentru mecanismul de control, TFRC implementează direct o ecuație de transfer pentru rata de transmitere permisă, ca funcție de rata evenimentelor de pierdere și timpul de dus-întors. TFRC folosește direct ecuația de transfer a TCP, unde rata este o funcție de evenimentele de pierdere, timpul de dus-întors și dimensiunea segmentelor. Mecanismul de control funcționează astfel:

1. Receptorul masoara rata evenimentelor de pierdere și trimite aceasta informație emițătorului.
2. Emitatorul foloseste aceste mesaje pentru a determina timpul dus-intors pentru un pachet (RTT = round-trip time).
3. Rata evenimentelor de pierdere și RTT sunt introduse în ecuatia de transfer și rata de transmisie rezultata este limitata la cel mult dublul ratei de receptie, pentru a se putea transmite cu rata X (cea calculata).
4. emițătorul ajusteaza apoi propria rata de transmitere pentru a se potrivi ratei de transmisie X. [1]

6.2 Ecuatia de tranfer

Ecuatia de transfer este urmatoarea (o forma usor simplificata):

$$X[B/s] = \frac{s}{R \cdot \sqrt{\frac{2 \cdot b \cdot p}{3}} + (t_{RTO} \cdot (3 \cdot \sqrt{\frac{3 \cdot b \cdot p}{8}} \cdot p \cdot (1 + 32 \cdot p^2)))}$$

Unde:

- X[B/s] este media ratei de transmisie
- S este dimensiunea segmentului în octeti
- R = timpul dus-intors (RTT) în secunde
- p = rata evenimentelor de pierdere a pachetelor, cuplinsa între 0 și 1, și este fractiunea de pachete pierdute din totalul pachetelor emise
- t_{RTO} = timpul de expirare pentru retransmisie, exprimat în secunde
- b = numărul maxim de pachete care sunt confirmate cu un singur pachet de confirmare TCP. [1]

7 Bibliografie

- [1] RFC5348
- [2] RFC3448
- [3] RFC4828
- [4] RFC3550
- [5] <http://www.cs.odu.edu/~cs778/jeffay/Lecture6.pdf>
- [6] <http://www.networksorcery.com>
- [7] <http://opalsoft.net/qos/>
- [8] RFC4341
- [9] RFC 4336
- [10] RFC4340
- [11] Andrew S. Tanenbaum: Computer Networks
- [12] <http://www.scribd.com/doc/22272663/All-About-UDP>
- [13] <http://www.netfor2.com/udp.htm>
- [14] www.wikipedia.com
- [15] http://www.ardenstone.com/projects/seniorsem/reports/UDP_Protocol.html

- [16]www.unibuc.ro/prof/niculae_c_m/telecom/user_datagram_protocol_udp.htm

B. Protocolul TCP - Transmission Control Protocol

Lecu Radu Șerban

1) Introducere

Protocolul TCP (Transmission Control Protocol) este unul dintre protocoalele nivelului Internet al modelului TCP/IP

TCP a fost special conceput pentru a oferi un transport de biți de tip capăt-la-capăt, fără erori sau pierderi de date. El a fost definit în standardul RFC 793 (Request For Comment) scris în 1981.

O rețea de internet diferă de una globală deoarece diferite părți pot avea topologii, lățimi de bandă, timpi de întârziere diferiți, dimensiuni ale pachetelor sau alți parametri diferiți. Astfel TCP a fost astfel conceput să se adapteze la proprietățile diferite ale componentelor folosite în cadrul rețelei.

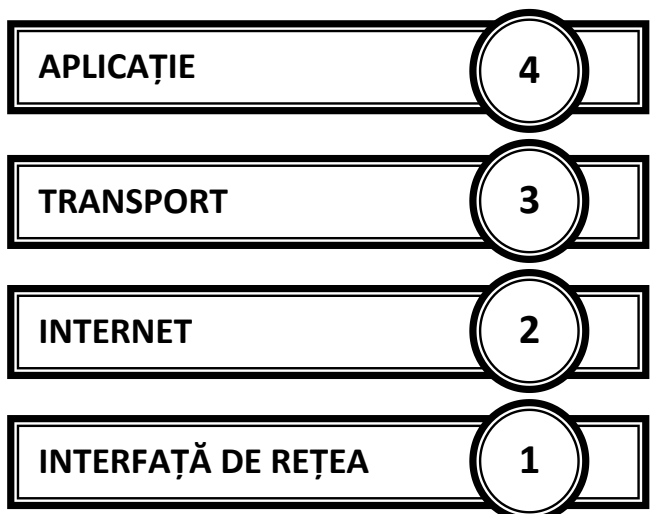
Fiecare mașină care suportă TCP este astfel concepută încât este capabilă să transmită informații și date prin transmisia unui șir de octeți. Un șir de octeți poate fi grupat în pachete de date de lungime de maxim 64KB incluzând și informațiile de transmisie (headerul). În cazul în care informația, cum ar fi cea conținută de un fișier, depășește această valoare, ea trebuie divizată în mai multe pachete de către emitor în așa manieră încât să poată fi recuperată și reconstruită la recepție. Pachetele sunt transmise independent, iar ele pot ajunge la destinație pe căi diferite și în altă ordine. [5]

Protocolul TCP reprezintă unul din componentele modelului TCP/IP (Protocol de control al transmisiei/ Protocolo Internet) care a fost creat de US DoD (Departamentul de apărare al Statelor Unite) din necesitatea unei rețele care ar putea supraviețui în orice condiții. Scopul rețelei TCP/IP era ca orice conexiune să nu se rupe în interiorul rețelei rețeaua în ansamblu să rămână intactă.

Astfel a trebuit concepută o arhitectură complexă, flexibilă și care avea în vedere aplicații divergente, printre care se transmit informații (cum ar fi vorbirea, fișiere) în timp real.

Pornind de la structura modelului OSI, cerințele de mai sus au condus la alegerea a 4 niveluri pentru modelul TCP/IP

Nivelul Transport este similar celui în cadrul stivei OSI, ocupându-se de probleme legate de siguranță, controlul fluxului și corecție de erori. El a fost astfel conceput încât să permită conversații între entitățile pereche sursă și destinație. Astfel au fost concepute 2 protocoale capăt-la-capăt: TCP și UDP.



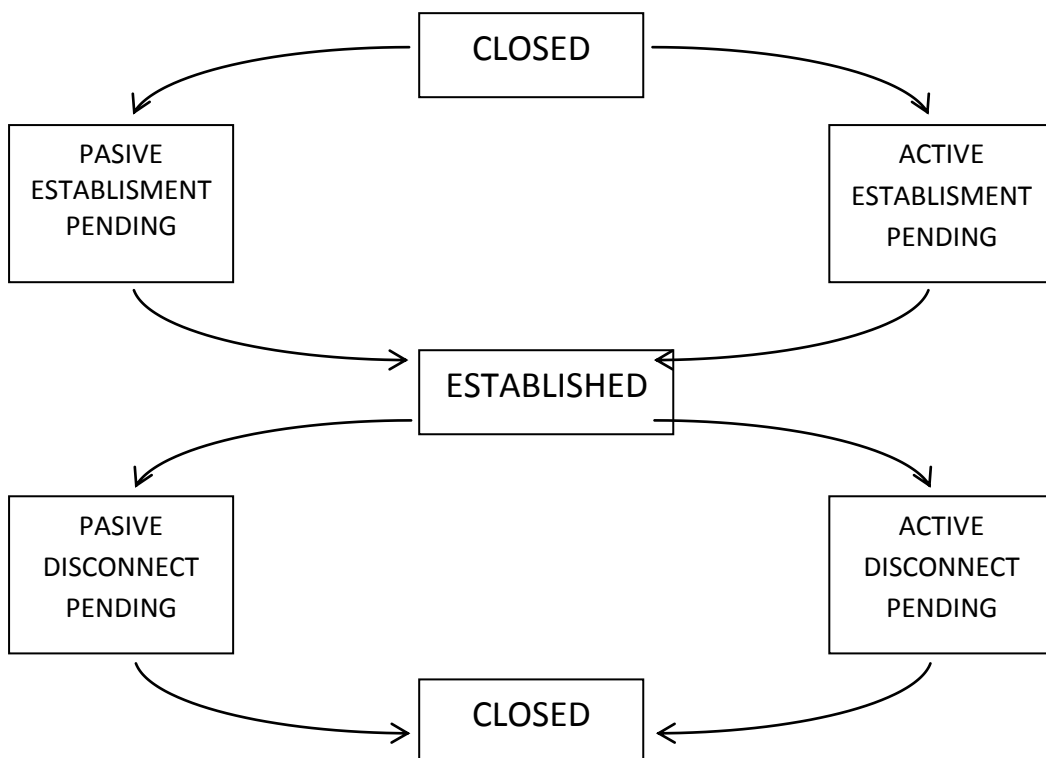
Protocolul TCP este un protocol orientat pe conexiune care permite ca un flux de octeți transmiși să ajunga la destinație ,fără erori, la orice altă mașină din cadrul aceleiași rețea. [4]

Modelul TCP este protocolul folosit de nivelul Transport al stivei TCP/IP.

2) Modelul Serviciului TCP

Serviciul TCP trebuie asigurat atat de emitator cat si de receptor deoarece se bazează pe principiul capăt-la-capăt și folosește socket-uri. Fiecarui socket îi corespunde un număr cu rol de adresă numit adresa IP a destinatarului și un alt număr alcătuit din 16 biți numit port. Pentru ca serviciul TCP să fie obținut trebuie realizată explicit o conexiune între mașina care transmite mesajul și cea care îl recepționează. Simplificat o conexiune poate fi reprezentată astfel:

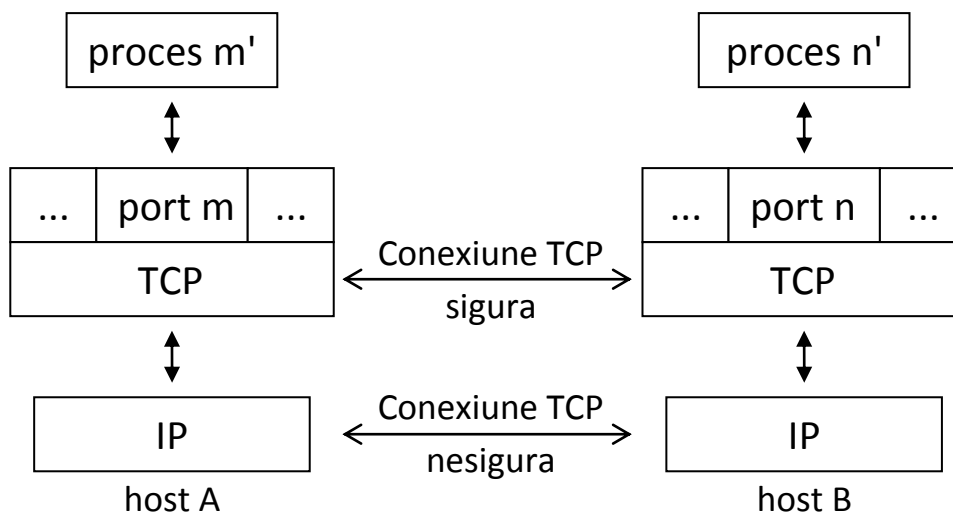
Astfel TCP poate fi comparat cu o conexiune telefonică din punct de vedere al modului în care utilizatorul vede această conexiune. Diferența de bază este dată de faptul că mesajele sunt transmise prin intermediul pachetelor fără a bloca o linie telefonică în cazul netransmiterii de informație.[1]



Astfel un anumit port se poate afla în diferite stări, cum ar fi IDLE, ESTABLISHED sau stări intermediare în cadrul cărora se realizează stabilirea conexiunii sau deconectarea. Conexiunea presupune 3 faze: realizarea conexiunii (IDLE->ESTABLISHED), transmiterea informației (starea rămâne ESTABLISHED) și întreruperea conexiunii (ESTABLISHED->IDLE).

Un socket poate fi folosit pentru mai multe conexiuni simultane. Astfel 2 sau mai multe conexiuni se pot termina la același socket. Acestea sunt identificate prin identificatori de socket aflați la ambele capete.

Conexiunea TCP poate fi reprezentată în felul următor:[6]



Din cele 65536 porturi (2^{16}), primele 1024 porturi sunt numite well-known-ports și sunt rezervate pentru servicii standard. Spre exemplu se dorește realizarea unui transfer de fișiere prin intermediul protocolului FTP (File Transfer Protocol). Aceasta se realizează prin conexiunea la portul 21 pentru a-și realiza conexiunea cu propriul demon. Demonul reprezintă o aplicație care rulează pe un anumit calculator fără control direct realizat de utilizator. Câteva exemple de porturi mai importante sunt cele din tabelul următor

Port	Protocol	Utilizare
21	FTP	File transfer
23	Telnet	Remote login
25	SMTP	E-mail
69	TFTP	Trivial file transfer protocol
79	Finger	Lookup information about a user
80	HTTP	World Wide Web
110	POP-3	Remote e-mail access
119	NNTP	USENET news

Ar fi posibil ca de exemplu demonului FTP să îi fie asociat portului 21 în momentul butării. Similar pentru celelalte porturi. Totuși acest lucru ar umple memoria cu demoni ce vor fi idle majoritatea timpului.

Modul de realizare este următorul: Un singur demon, numit inetd (Internet demon) este atașat la rândul său la porturi multiple și așteaptă prima conexiune. Când se realizează acest lucru, inetd realizează un nou proces și execută demonul corespunzător, lăsându-l pe acesta să se ocupe de cerere. Cu toate acestea este posibilă și realizarea unor demoni utilizați permanent cum ar fi portul 80 (HTTP) și inetd în rest.

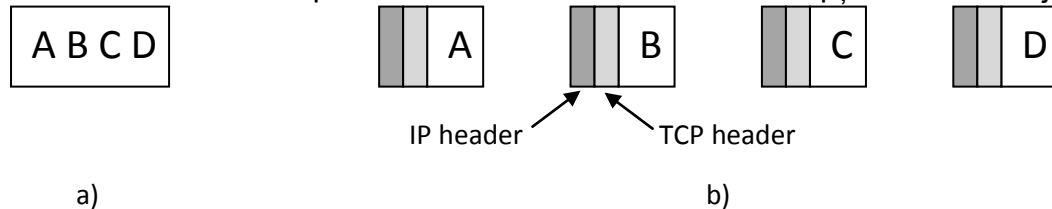
Toate conexiunile TCP sunt de tip full-duplex capăt-la-capăt. Full-duplex înseamnă că transmisia se poate face bidirecțional și simultan între cele 2 calculatoare. Capăt-la-capăt arată că nu pot exista mesaje de tip multicast sau broadcast.

O conexiune TCP reprezintă o înșiruire de biți (octeți), care are o structură bine definită și care este împărțită în diferite componente. Pe lângă mesaj (șirul de octeți ce reprezintă informația) la fiecare nivel se mai adaugă și alte informații numite

headere. Acestea conțin informații cum ar fi ip sursa, ip destinație, port destinație, informații de detecție a erorilor, informații de ordonare a pachetelor pentru reconstrucția întregului mesaj din fragmentele sale, etc. și sunt asociate fiecărui nivel în parte. În cazul lipsei unor astfel de informații pachetele pot ajunge la destinație dar sunt inutilizabile, sau pot chiar să nu fie recepționate.

Să presupunem că avem un mesaj împărțit în 4 fragmente (A,B,C,D). Fiecărui mesaj îi va corespunde un header ip și unul TCP ca în figurile următoare:

Protocolul TCP se ocupă numai de transmiterea sau recepționarea mesajului



nu și de interpretarea lui. La emisie mesajului i se adaugă informațiile suplimentare din headere. La recepție mesajul este separat de headere fiind transmis nivelelor superioare.

Important este de menționat că pentru realizarea conexiunii se folosesc anumite pachete cu ajutorul cărora emițătorul și receptorul comunică pentru realizarea conexiunii. Pachetele care sunt folosite pentru realizarea conexiunii au prioritate față de altele în care se transmite informație. Acestor pachete li se asociază un flag numit URGENT flag. Când un pachet urgent este recepționat acesta are prioritate față de celelalte fiind primul interpretat chiar dacă alte pachete anterioare au venit înaintea acestuia.[1]

3) Protocolul TCP

Receptorul și transmițătorul comunică prin intermediul unor segmente. Un segment TCP este alcătuit dintr-un header de 20 biți la care se adaugă mesajul. Dimensiunea header-ului plus cea a mesajului pot avea împreună până la 65515 octeți.

Protocolul de bază al TCP-ului este sliding window protocol, care se bazează pe principiul ferestrei glisante. În momentul în care emițătorul transmite un segment, acesta porneste și un timer. Când segmentul ajunge la destinație, receptorul transmite un segment (cu sau fără date) care conține un număr de confirmare egal cu următoarea secvență pe care o așteaptă. Dacă timerul emițătorului este depășit înainte de recepția mesajului de confirmare, mesajul este retransmis. [6]

Spre deosebire de UDP care transmite pachete între 2 hosturi fără a oferi siguranța că acestea ajung la destinație, TCP este protocol orientat pe conexiune. Astfel pentru fiecare pachet transmis de la sursă la destinație, cu ajutorul arhitecturii TCP, trebuie să parcurse mai multe etape:

- realizarea conexiunii dintre cele 2 hosturi
- schimbarea de informații între acestea
- întreruperea conexiunii

În cazul UDP, un pachet este transmis fără a avea siguranța că informația a ajuns de la sursă la destinație. Același pachet transmis folosind protocolul TCP, va oferi siguranța că pachetul a ajuns la destinație. Mai mult, în cazul în care locația dorită nu poate fi găsită, emițătorul va fi informat de acest lucru.

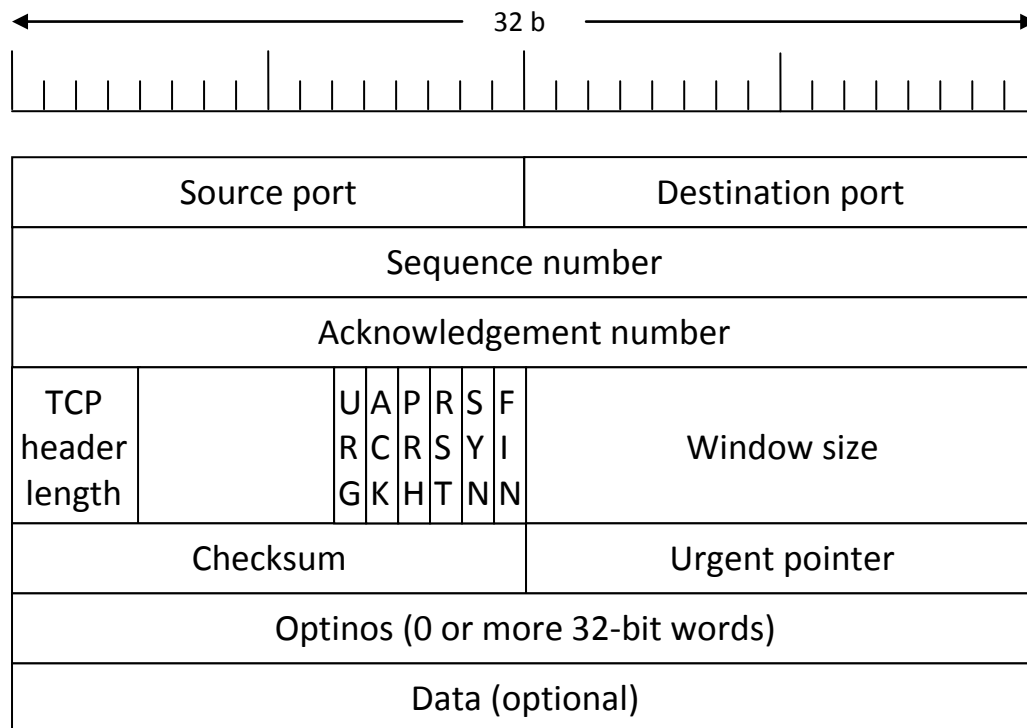
Există 2 cazuri în care nu se poate transmite un pachet.

a) Nu se poate realiza conexiunea între cele 2 hosturi. În momentul în care se dorește realizarea conexiunii locația nu poate fi găsită sau nu acceptă realizarea conexiunii.

b) Conexiunea a fost deja realizată anterior între cele 2 hosturi dar pachetul nu a putut fi transmis. Acest lucru se poate întâmpla fie din cauză că unul din hosturi s-a deconectat fără a anunța întreruperea conexiunii (ex: întreruperea curentului), fie fără ca cererea de închidere să ajungă la celălalt capăt al conexiunii (legătura dintre calculatoare a fost întreruptă și nu există altă rută alternativă).

4) Header-ul TCP

Fiecare segment începe cu un format fix de 20 B urmat de cuvinte duble de 32 biti optionale și de date. Câmpul de date poate să conțină până la 65495 octeți (65535 - 20 octeți pt header-ul ip -20 octeți ai header-ului TCP).



Portul sursă și portul destinație arată porturile de emisie și de recepție. Portul plus adresa ip formează o adresă unică de 48 biti.

Sequence number și Acknowledgement number realizează funcțiile de verificare dacă a ajuns pachetul la destinație sau nu. Ambele au 32 de biti.

TCP header length arată lungimea headerului în cuvinte de 32 biți. Informația este necesară deoarece headerul include și câmpul de opțiuni care este variabil.

Urmează 6 biți nefolosiți și alți 6 care reprezintă flaguri care au rol diferit

- URG - Urgent pointer - descris mai sus
- ACK - Acknowledge flag - arată că Acknowledgement number este valid adică pachetul are rol de confirmare a recepției mesajului
- PSC - Push Function - arată că mesajul trebuie transmis aplicației cât mai repede posibil. Astfel se oferă prioritate pachetului față de alte pachete.

- RST - Reset conexiun - arată apariția unei probleme ce impune resetarea conexiunii

- SYN - Sincronizare - este folosit pentru realizarea conexiunii.

- FIN - No more data from sender - este folosit pentru închiderea conexiunii.

Window size are rolul de a arăta lungimea câmpului de date în octeți.

Checksum reprezintă o sumă de verificare cu ajutorul căruia se determină dacă datele transmise sunt eronate sau nu. Aceiași operație este realizată în ambele capete iar dacă ele coincid atunci este foarte probabil ca informația obținută să fie corectă.

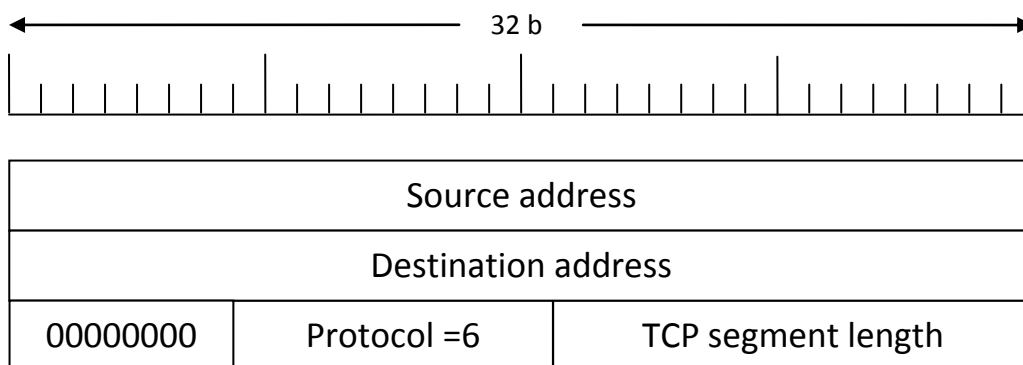
Urgent pointer este folosit când se dorește ca anumite date să fie procesate în cel mai scurt timp posibil. El indică unde se termină ultimul byte de date urgente.[5]

Options indică anumite informații opționale ce pot fi transmise în scopul obținerii anumitor funcționalități. Una din cele mai importante funcționalități este lungimea maximă a segmentului (MSS Maximum Segment Size). Aceasta are importanță în cazul în care avem o linie de transmisie cu multe erori. Dacă segmentul are lungimea mai scurtă, probabilitatea ca să avem cel puțin un bit eronat este mai mică iar în cazul în care se retransmite pachetul se retransmite o cantitate mai mică de informație. Astfel în anumite cazuri este de preferat să împărțim un mesaj în fragmente mai multe dar mai mici decât să avem pachete de dimensiune mare. În același timp număr mare de pachete duce la cantitate mare de informație datorată header-ilor dar inutilă pentru utilizator și care crește durata de transmisie a întregului mesaj. Valoarea predefinită a dimensiunii mesajului este de 536 B pentru date.

Câmpul de date reprezintă informația ce va fi transmisă către sau recepționată de la nivelul de aplicație pe portul cerut.

Pseudoheader-ul este folosit pentru calculul câmpului Checksum el nefiind transmis în cadrul pachetului. Se observă că acesta conține și informații care nu aparțin nivelului TCP deoarece IP-ul este o informație a nivelului IP, violând astfel ierarhia TCP/IP.

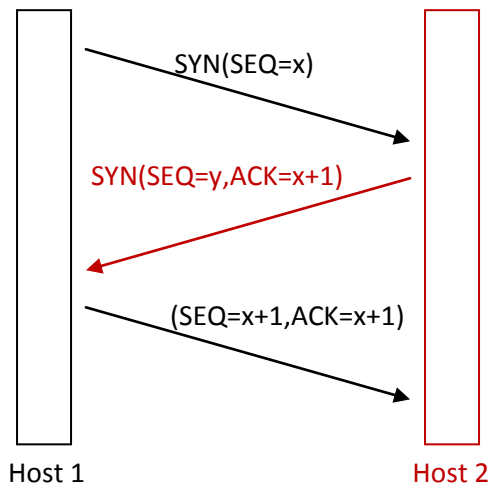
Următoarea figură arată structura pseudoheader-ului[6]:



5) Stabilirea conexiunii TCP

Realizarea conexiunii presupune o secvență de 3 pachete transmise între cele 2 hosturi cunoscută sub numele de three-way handshake. Pentru stabilirea conexiunii este necesar ca un host, numit Rețea, să aștepte pasiv realizarea conexiunii. În același timp celălalt host, numit Client, va transmite un pachet prin care cere realizarea conexiunii la adresa IP și la portul dorit. Pachetul transmis are setate flagul SYN =1 și ACK=0. Rețeaua răspunde prin transmiterea unui pachet având flagurile SYN=1 și ACK=1. În final Clientul retransmite un pachet având setați SYN=0

si ACK=1, moment în care s-a realizat conxiunea în cazul în care nu au aparut alte erori de-a lungul întregului proces. [1]



Exista posibilitatea ca simltan ambele hosturi să încerce realizarea conexiunii între aceiasi 2 socketi. În final este realizată numai o singura conexiune.[1]

6) Întreruperea conexiunii TCP

Conexiunea full-duplex dintre hosturi poate fi văzută si ca 2 conexiuni semi-duplex

Pentru a se realiza întreruperea conexiunii, oricare din cele 2 hosturi pot cere acest lucru cu ajutorul flagului FIN. Deoarece când unul dintre hosturi cere întreruperea conexiunii celălalt poate transmite date, întreruperea conexiunii se face simultan în ambele părți. Porcedeu este realizat similar realizarii conexiunii fiind realizat printr-o secventa de 3 pachete.

Există însă situația ca mesajul de încheiere a conexiunii să nu fie recepționat. Dacă în intervalul de timp maxim în care se poate transmite si receptiona un mesaj nu s-a primit un răspuns se încheie conexiunea numai într-o singura parte, ramanand ca celalalt host să observe că nimeni nu mai asculta mesajele sale. [1]

7) Managementul conexiunii

Conexiunea TCP poate fi reprezentată printr-un automat cu număr finit de stări. Automatul are 11 stări prezentate în tabelul următor.

Ca orice automat, în funcție de starea curentă nu se poate trece decât în anumite stări. Succesiunea acestori stări este reprezentata în graficul de mai jos. Tot în grafic seunt prezentate si primitivele (comenzile cu care utilizatorul face trecerea dintr-o stare în alta).

Stare	Descriere
CLOSED	Nici o conexiune nu este activă
LISTEN	Servărul așteaptă un apel
SYN RECEIVED	A fost primită o cerere de conexiune
SYN SENT	Aplicația a început deschiderea conexiunii
ESTABLISHED	Starea normală pentru transferul de date
FIN WAIT 1	Aplicația a zis că a terminat de transmis date
FIN WAIT 2	Cealaltă parte a acceptat închiderea conexiunii
TIME WAIT	Așteptare ca toate pachetele să fie transmise
CLOSING	Ambele părți încearcă închiderea conexiunii simultan
CLOSE WAIT	Cealaltă parte a inițiat închiderea conexiunii
LAST ACK	Așteaptă ca toate pachetele să fie transmise

Fiecare conexiune începe cu starea CLOSED. Această stare este părăsită numai când se face o deschidere pasivă (hostul are rol de rețea, adică se așteaptă cererea de conexiune de la alt host) sau când se face una activă (hostul are rol de client, el apelează un alt host cu rol de rețea, căruia îi cere realizarea conexiunii).

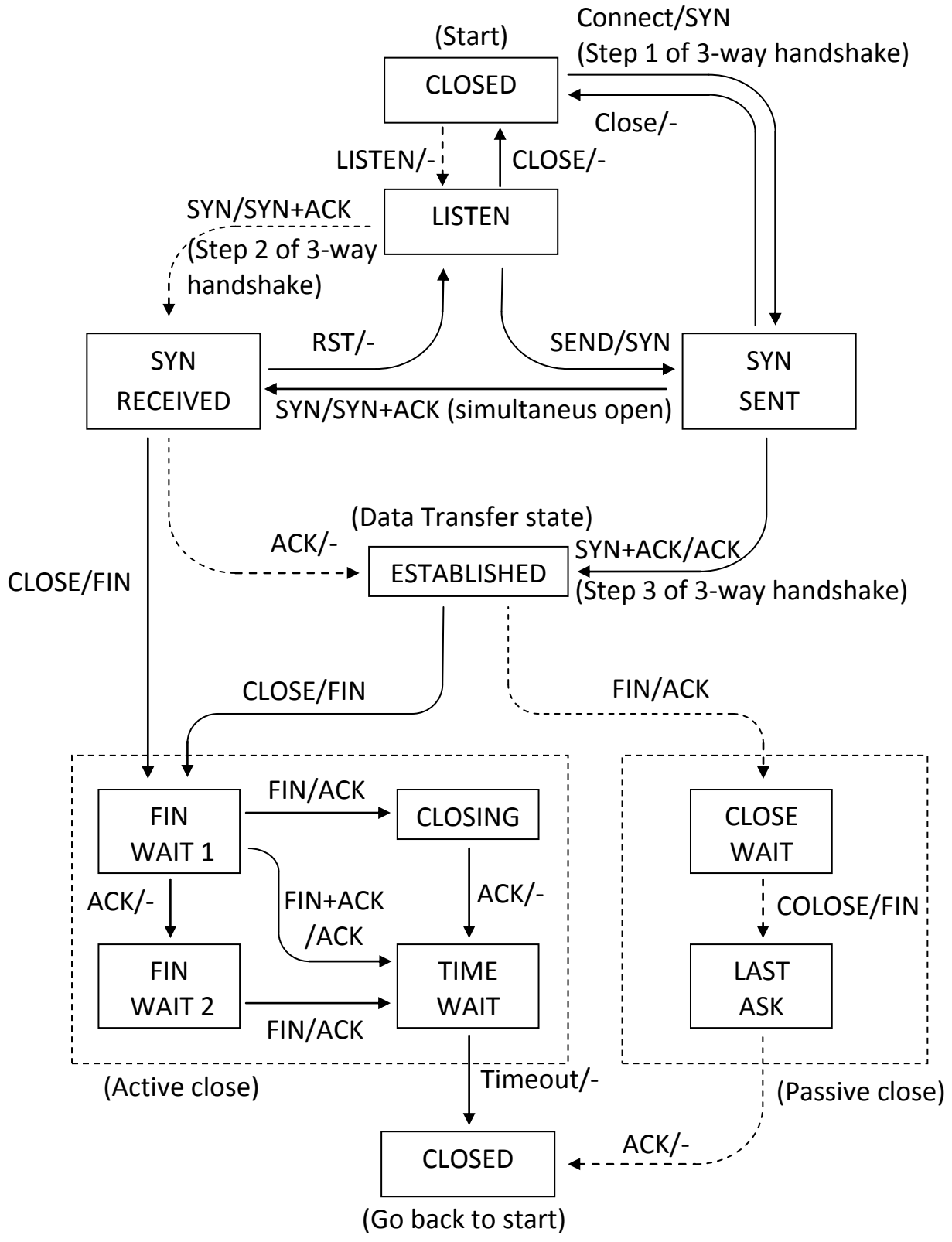
Rețeaua trece în starea LISTEN, stare în care așteaptă cererea de conexiune. Clientul trece în starea SYN SENT și transmite un segment de tip SYN. Dacă există un host de tip rețea la adresa cerută, care ascultă pe portul corespunzător, acesta recepționează mesajul și transmite un mesaj de tip SYN+ACK. În cazul în care clientul nu primește un astfel de mesaj după un anumit interval de timp, el trece înapoi în starea CLOSED. Altfel el transmite un semnal de tip ACK, moment în care se realizează conexiunea.

Atât clientul, cât și rețeaua după ce primesc pachetul ACK, trec în starea ESTABLISHED. În această stare se realizează comunicarea bidirecțională între cele 2 hosturi.

Metoda prezentată mai sus poartă numele de 3-way Handshake.

În orice moment, oricare din cele 2 hosturi poate realiza întreruperea conexiunii. Aceasta se realizează prin metoda prezentată mai sus.

Există 2 moduri de întrerupere a conexiunii: activă, în cazul celui care a inițializat întreruperea conexiunii, sau pasivă în cazul celui care este anunțat de întreruperea conexiunii.[1]



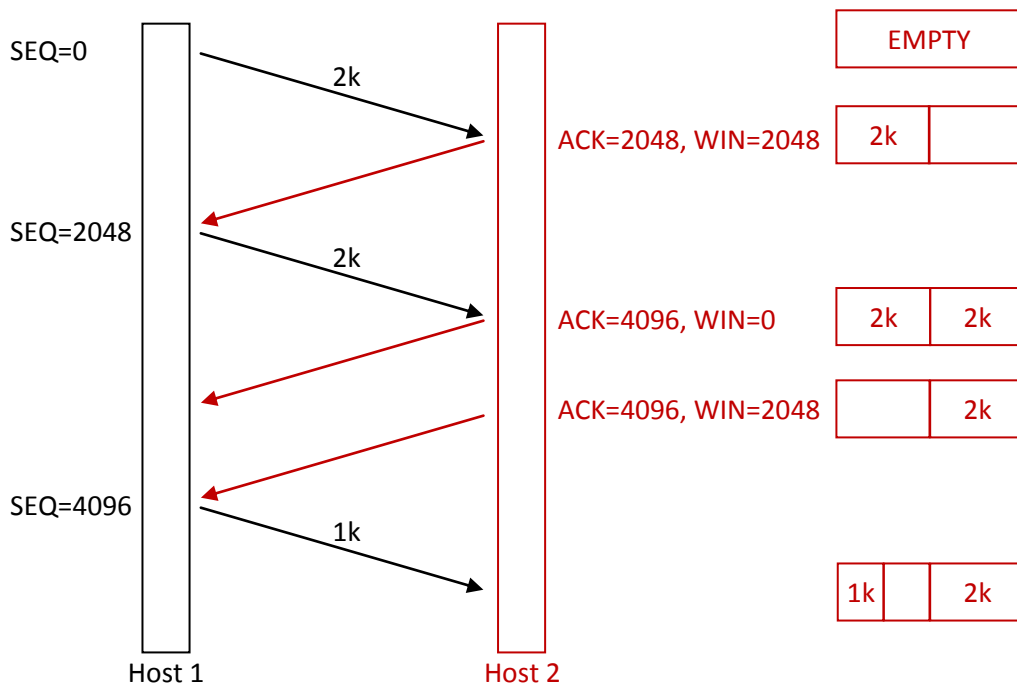
8) Politici de transmisie TCP

În capitolele anterioare s-a discutat despre modul în care se realizează conexiunea dintre 2 entități diferite, cum se întrerupe conexiunea dintre acestea, structura segmentelor și a headerelor acestora. Toate acestea au ca scop final oferirea posibilității de a transmite informații între hosturi fără pierdere de informație.

Transmisia de informații între cele 2 hosturi se face numai când se presupune că ambele capete ale conexiunii sunt în starea ESTABLISHED. Astfel se porneste de la ideea că a fost realizată anterior conexiunea.

Fiecare host are un buffer în care odată pachetele recepționate ele sunt stocate într-un buffer. Acest buffer are o dimensiune limitată. Astfel dacă are loc transmiterea unui pachet care nu poate fi stocat atunci acesta va fi pierdut fiind necesară retransmisia sa.

Pentru a minimiza cantitatea de informație pierdută s-au introdus protocoale de management ale ferestrei. Modul de funcționare este prezentat în figura următoare:



După cum se observă în imagine, receptorul, în momentul în care trimite mesaje de ACK, trimite și informații suplimentare pentru a informa emițătorul despre spațiului liber din buffer, prin utilizarea câmpului WIN. Valoarea câmpului WIN arată spațiul liber disponibil al bufferului. Astfel, dacă bufferul este plin se va opri transmisia. În momentul în care se eliberează o zonă din buffer, receptorul trimite un mesaj cu dimensiunea zonei eliberate din interiorul bufferului adică cantitatea maximă de informație pe care emițătorul o poate transmite. [1]

9) Managementul timerelor TCP

Pentru transmisia datelor, se pun următoarele probleme:

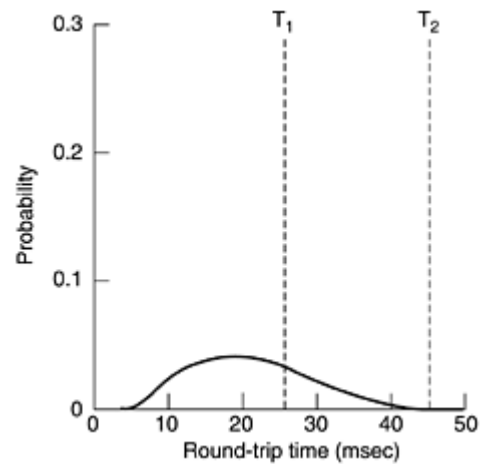
- fiecare segment trebuie să ajungă la destinație;
- un segment ajuns la destinație nu trebuie să aibe biți eronați.

A doua condiție este verificată cu ajutorul checksum din headerul TCP. Prima condiție presupune existența unui timer numit timerul de retransmisie.

TCP-ul folosește multiple timere. Cel mai important este cel de retransmisie. Atunci când un segment este transmis, este pornit un timer. Dacă se recepționează un ACK pentru acel segment înainte ca timerul său să expire atunci se consideră că mesajul a ajuns la destinație. În caz contrar se retransmite segmentul iar timerul se repornește. Se pune problema alegerii dimensiunii timerului.

În cazul în care timerul este mic, există o foarte mare probabilitate ca pachetul să fie retransmis înainte de recepționarea ACK-ului. Astfel se poate provoca o congestie datorită introducerii de pachete inutile în rețea. Dacă timerul este foarte mare atunci intervalul de timp după care are loc retransmisia este foarte mare, astfel că se va aștepta un interval mare de timp recepționarea segmentului dacă acesta a nu a fost recepționat de prima dată.

Timerul trebuie astfel ales încât să se obțină optimul pentru cele 2 cazuri.[1]



10) Algoritmi utilizați de TCP pentru eficientizarea transmisiei de pachete

Pentru a rezolva diferite probleme ale TCP s-au adoptat 4 algoritmi folosiți de TCP a căror documentație se găsește în RFC 2001 publicat în 1997. Aceștia sunt:

- Slow start
- Congestion avoidance
- Fast retransmit
- Fast recovery algorithms

10.1) Slow start

Versiunile inițiale de TCP ar începe o conexiune cu clientul injectând numeroase segmente în rețea de până la lungimea maximă a ferestrei. Acest lucru nu reprezintă o problemă în cazul în care ambele hosturi se află în același LAN. Problemele apar atunci când există rutere între hosturi. Unele rutere intermediare trebuie să păstreze pachetele într-o memorie. Drept urmare vom avea o încărcare exagerată a memoriei și deci la umplerea ei.

Algoritmul pentru evitarea acestor situații este Slow start. Acesta operează prin a observa că rata maximă la care noile pachete ar trebui injectate în rețea este egală cu rata cu care răspunsurile se întorc de la celălalt capăt.

Pentru realizarea algoritmului se adaugă o nouă fereastră emitorului numită fereastră de congestie (cwnd - congestion window), când o nouă conexiune este realizată cu un host din altă rețea, fereastra de congestie este inițializată la un segment (dimensiunea segmentului cerută de la celălalt capăt sau o valoare default, în general egală cu 536 sau 512). De fiecare dată când un nou ACK este recepționat, cwnd este crescută cu încă un segment. Emițătorul poate să transmită până la minimul dintre lungimea ferestrei de congestie și cea de advertisement a receptorului. Fereastra de congestie este controlată de emițător, iar cea de advertisement este controlată de receptor.

Principiul de funcționare este următorul: emițătorul începe prin a transmite un segment. Inițial $cwin=1$. Când receptorul primește mesajul el va trimite un răspuns de tip ACK. Când emițătorul îl recepționează $cwin$ devine 2. De data aceasta se pot

transmite până la 2 pachete. Drept răspuns vor fi 2 ACK-uri deci $cwin=4=2^2$. În continuare $cwin$ devine $2^3, 2^4, 2^5$ și așa mai departe astfel că vom avea o creștere exponențială. Este posibil ca, de exemplu când $cwin$ devine 16, emițătorul să nu mai dorească să mai transmită 16 pachete, ci numai 5. Ca rezultat $cwin$ va deveni 21 nu 32, caz în care nu vom mai avea o creștere exponențială. Dacă continuă creșterea numărului de pachete, la un moment dat capacitatea rețelei poate fi atinsă iar un router intermediar va începe să renunțe la pachete, acestea fiind pierdute. Pierdere lor determină micșorarea ferestrei.

Deși inițial Slow start a fost aplicat numai pentru entități aflate în rețele diferite în prezent este folosită și pentru rețelele locale.[7][9][12]

10.2) Congestion avoidance (Evitarea congestiei)

Congestia poate să apară atunci când, pachetele care intră într-un router și care urmează să iasă toate în același LAN, depășesc debitul maxim suportat de LAN. Congestion avoidance este un procedeu prin care se pot evita astfel de situații pentru a minimiza numărul de pachete pierdute.

Presupunerea algoritmului este că pachetele pierdute prin eroare sunt foarte mici (maxim 1%), deci pierderea unui pachet semnează existența unei congestii între sursă și destinație. Există 2 indicatori ai pierderii unui pachet: expirarea timerului sau recepționarea a 2 ACK-uri.

Congestion avoidance și Slow start sunt algoritmi diferiți cu obiective diferite. Totuși când congestia apare TCP trebuie să încetinească rata de transmisie a pachetelor în rețea, și apoi să invoce Slow start pentru a crește din nou rata de transmisie. De aceea ele sunt în practică implementate împreună.

Cei 2 algoritmi necesită 2 variabile pentru fiecare conexiune: $cwnd$ (congestion window) și $ssthresh$ (Slow start threshold).

Algoritmul combinat funcționează după cum urmează:

I) Inițializarea conexiunii date. Se setează $cwnd$ la 1 segment și $ssthresh$ la 65535.

II) Rutina de transmisie TCP nu transmite mai mult decât minimul dintre $cwnd$ și dimensiunea ferestrei de advertisement a receptorului.

III) Când apare congestia indicată de expirarea timerului sau de recepționarea de ACK dublu, o jumătate din dimensiunea ferestrei curente (minimul dintre $cwnd$ și dimensiunea bufferului receptorului dar cel puțin egală cu 2) este salvată în $ssthresh$. Suplimentar dacă congestia este determinată datorită expirării timpului $cwnd$ este setat la 1 segment (Slow start).

IV) Când este recepționat un nou ACK de la celălalt capăt, $cwnd$ este crescut cu 1 dar acesta crește în funcție de ce realizează TCP: Slow start sau Congestion avoidance.

Dacă $cwnd$ este mai mic sau egal cu $ssthresh$, TCP se aplică Slow start, altfel este realizat Congestion avoidance. Astfel TCP se află în Slow start până în momentul în care $cwnd$ atinge $ssthresh$ după care se realizează Congestion avoidance.

În cazul slow start $cwnd$ crește cu 1 pentru fiecare ACK primit. În cazul Congestion avoidance $cwnd$ este crescut cu $\frac{segsize^2}{cwnd}$ pentru fiecare ACK primit, unde $segsize$ și $cwnd$ au ca unitate de măsură bytes. Astfel vom avea o creștere liniară în cazul Congestion avoidance comparativ cu cea exponențială pentru cazul slow start. [7][9][12]

10.3) Fast retransmit (Retransmitere rapidă)

Fast retransmit reprezintă o modalta prin care TCP reduce timpul prin care emițătorul așteaptă pînă cînd retransmite segmentul pierdut.

Emițătorul TCP ,în varianta sa inițială, folosește un timer pentru a detecta pierderea de pachete. Dacă un ACK nu este recepționat pentru un anumit pachet pînă la expirarea timerului, emițătorul va presupune că pachetul a fost pierdut deci va trebui retransmis. De multe ori acest timp duce la întârzieri nedorite care de multe ori pot fi eliminate.

ACK dublu stă la baza algoritmului.

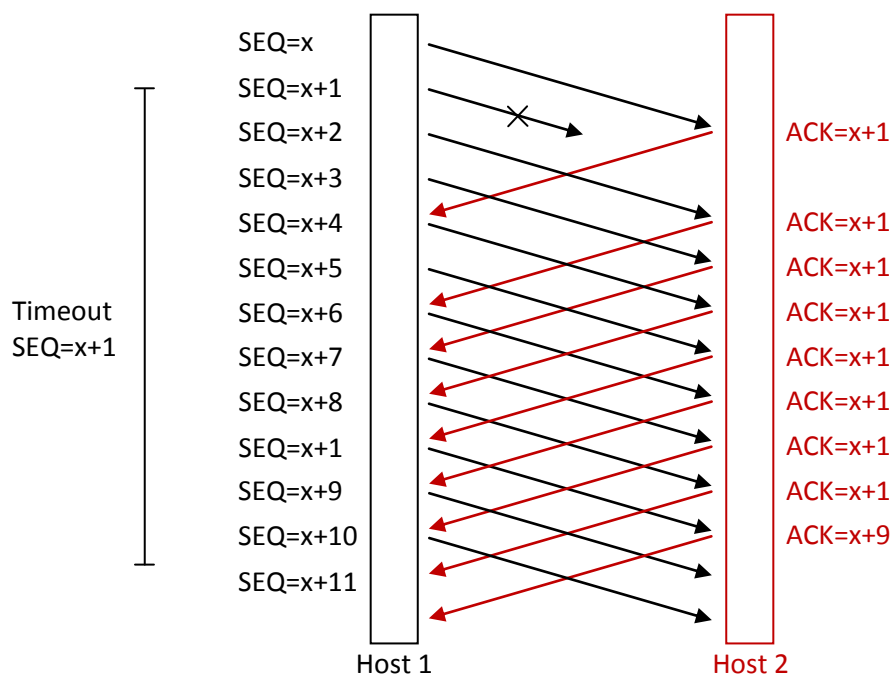
Un ACK dublu poate fi obținut din 2 motive: pierdere de pachet, sau inversare de pachete.

După recepționarea unui pachet cu $SN=x$ (Sequence number), receptorul transmite un mesaj cu $ACK=x+1$ (Acknowledgement number) care înseamnă că acesta a primit pachetul x și așteaptă pachetul $x+1$. În urma acestui procedeu emițătorul este informat despre pachetele care au ajuns la destinație.

Să presupunem că dorim să transmitem un șir de pachete, iar pachetul $SN=x+1$ nu ajunge la destinație, el fiind pierdut. $SN=x$ ajunge deci emițătorul recepționează $ACK=x+1$. Când $SN=x+2$ ajunge la destinație emițătorul recepționează $ACK=x+1$ în loc de $x+3$ deoarece se așteaptă încă pachetul $x+1$. Similar la $SN=x+2$, $SN=x+3$, $SN=x+4$, etc. emițătorul primește $ACK=x+1$.

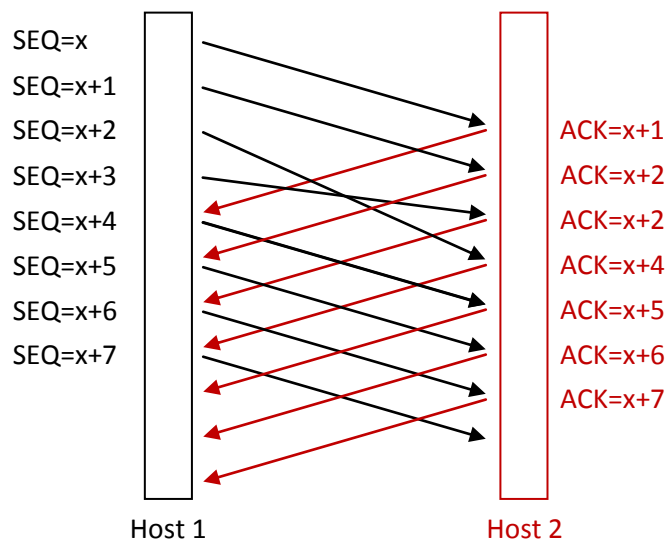
Dacă se primesc o succesiune de 4 ACK-uri egale atunci pachetul respectiv se retransmite deoarece este considerat pierdut.

Se observă că deși timerul nu a expirat încă pentru pachetul $SEQ=x+1$ el este retransmis mai repede. Există posibilitatea ca pachetul să fie retransmis cu mult înaintea timerului și astfel se poate ajunge la întârzieri mult mai mici.



Pentru cazul cu inversare de pachete(pachetele ajung în alta ordine) algoritmul funcționează conform figurii de mai jos.

Se observă că nu există probleme aduse algoritmului în cazul în care 2 pachete sunt inversate. Dacă însă vom avea un pachet întârziat foarte mult, există cel mult riscul ca să fie retransmis, iar la destinație va ajunge de 2 ori. Singura



problemă care se pune este aceea de detecție a pachetului pentru eliminarea sa. [9][12][13]

10.4) Fast recovery algorithms

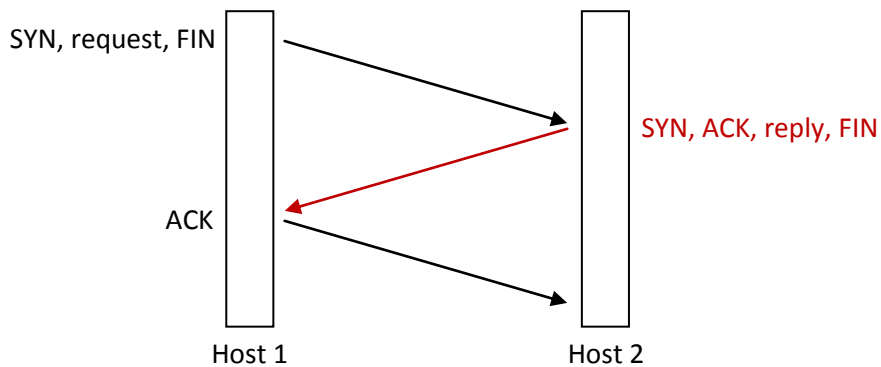
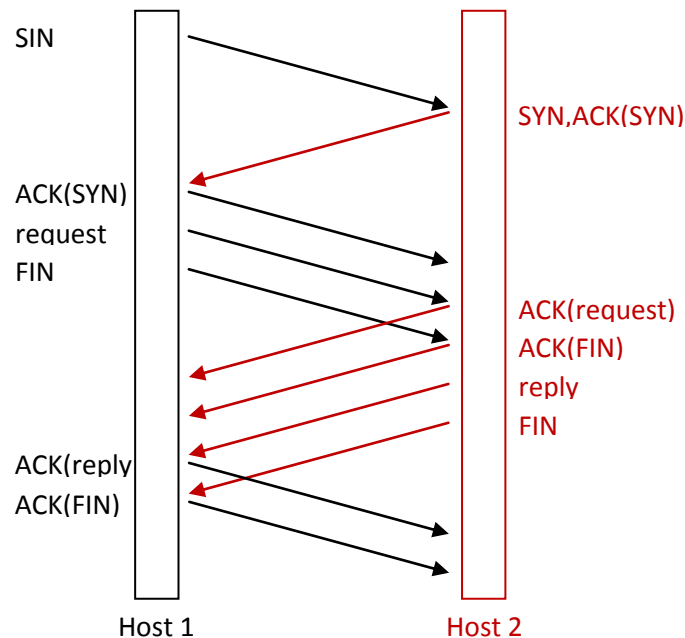
În cazul Fast Recovery Algorithm, după ce fast retransmis cea ce pare a fi un segment pierdut, este aplicat Congestion Avoidance în loc de Slow start. Aceasta este o îmbunătățire care permite debite mari cu un nivel de congestie moderat, în special pentru ferestre mari.

Motivul pentru care nu se aplică Slow start este că recepționarea ACK-urilor duble informează TCP-ul mai mult decât prin faptul că pachetul a fost pierdut. Deoarece receptorul nu poate genera decât ACK-uri duble când un alt segment este recepționat, acel segment a părăsit rețeaua și se găsește în interiorul bufferului receptorului. Există în continuare schimb de date între cele 2 capete, iar TCP-ul nu vrea să reducă abrupt fluxul de date prin trecerea în Slow start. [7][12]

11) TCP Transactional (T/TCP)

Protocolul TCP fiind un protocol orientat pe conexiune presupune realizarea unei conexiuni indiferent de numărul de pachete transmise între cele 2 hosturi. După cum se observă în imaginea următoare pentru transmiterea unui singur pachet de cerere a unei informații (request) și a unuia de răspuns (reply) este necesară transmiterea unui număr mare de pachete între cele 2 hosturi care duc la creșterea timpului de transmisie a datelor utile.

În cazul folosirii protocolului TCP Tranzacțional mai multe pachete sunt cumulate într-unul singur astfel fiind necesară transmiterea numai a 3 pachete între hosturi comparativ cu 11 din varianta inițială. [1]



12) Bibliografie

- [1] Andrew S. Tanenbaum, Computer Networks
- [2] http://en.wikipedia.org/wiki/Transmission_Control_Protocol
- [3] http://en.wikipedia.org/wiki/File:Tcp_state_diagram_fixed_new.svg
- [4] <http://ro.wikipedia.org/wiki/TCP/IP>
- [5] <http://condor.depaul.edu/jkristof/technotes/tcp.html>
- [6] http://www.unibuc.ro/prof/niculai_c_m/telecom/transm_ctrl_prot_tcp.htm
- [7] http://www.cisco.com/web/about/ac123/ac147/ac174/ac195/about_cisco_ipj_archive_article09186a00800c83f8.html
- [8] <http://www3.gdin.edu.cn/jpkc/dzxnw/jsikj/chapter3/35.htm>
- [9] <http://www.faqs.org/rfcs/rfc2001.html>
- [10] <http://www.cs.umd.edu/~shankar/417-F01/Slides/chapter3b/sld015.htm>
- [11] http://nptel.iitm.ac.in/courses/IIT-MADRAS/Computer_Networks/pdf/Lecture37_TCPConnectionMgmt.pdf
- [12] <http://en.wikipedia.org/wiki/Slow-start>
- [13] <http://www.fcoe.ru/english/data-transfer-protocols-overview/25-fcip>

C. Algoritmi de control al congestiei

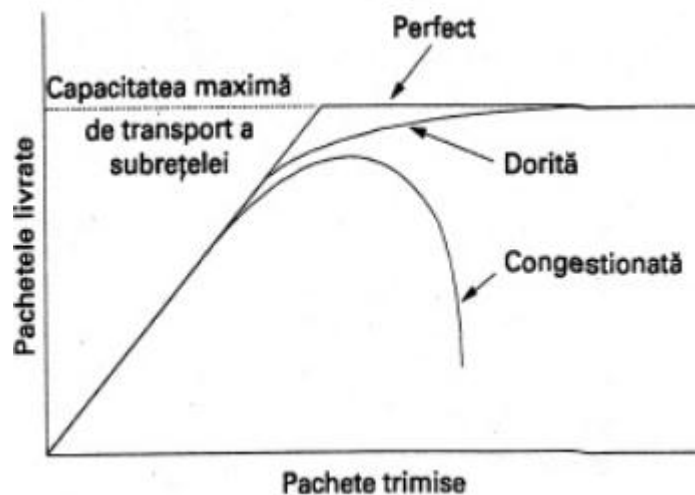
Tică Andra Maria

1. Definirea problemei [5]

Congestia apare în rețelele de calculatoare atunci când încărcarea depășește posibilitățile acesteia de a transfera datele. Încărcarea unei rețele la un moment dat nu depinde numai de capacitatea ei de transmitere ci și de erorile care apar în mediul de transmisie, de viteza de prelucrare în noduri și de mecanismele de confirmare folosite de receptor.

Controlul acestui fenomen intră în responsabilitatea nivelelor de transport și de rețea. Apare ca rezultat al traficului intens la nivelul transport, propagându-se la nivelul rețea, performanțele totale ale sistemului degradându-se prin întârzierea și pierderea pachetelor de date.

Această situație este prezentată în graficul următor:



[5]

Acesta prezintă cele trei situații posibile: transportul optim, transportul dorit și cazul în care într-o rețea sunt prezente foarte multe pachete și performanțele acesteia se degradează apărând fenomenul de congestie. Putem spune că atunci când numărul de pachete emise în rețea nu depășește capacitatea de transport, acestea sunt livrate integral, excepție făcând cele care prezintă erori de transmisie. În acest caz numărul pachetelor livrate este proporțional cu numărul celor emise. Atunci când traficul crește prea mult, ruterele încep să nu mai facă față și să piardă pachete. Această situație se poate deteriora complet la un trafic intens, ducând sistemul în imposibilitatea de a mai livra pachete.

Congestia apare ca rezultat al mai multor factori.

În situația în care la sosirea unui număr mare de pachete provenind de pe mai multe linii de intrare într-o singură linie de ieșire, atunci se va forma o coadă. Pentru a

le păstra pe toate, sistemul necesită memorie, iar creșterea acestei capacități de memorare duce la înrăutățirea congestiei și nu la ameliorarea ei.

Un alt factor care cauzează congestia este reprezentat de viteza procesoarelor. Dacă unitatea centrală a rutelor este lentă în execuția funcțiilor sale, cozile pot crește. Și liniile cu lățime de bandă scăzută pot provoca congestia. Schimbarea liniilor cu unele mai performante și păstrarea aceluiași procesor sau invers ajută puțin, deoarece problema ține de incompatibilitatea între părțile sistemului și aceasta va persista până la aducerea la echilibru a tuturor componentelor.

Controlul congestiei trebuie să asigure capabilitatea rețelei de a transporta întreg traficul implicat. Este o problemă globală care implică comportamentul tuturor calculatoarelor gazdă și al rutelor.

2. Soluții posibile. Clasificarea Yang&Reddy a algoritmilor.

Prezența congestiei înseamnă că încărcarea momentană a sistemului este mai mare decât capacitatea resurselor care pot fi gestionate de sistem la acel moment de timp. Sunt posibile două soluții dar eficacitatea acestora nu oferă o rezolvare definitivă a acestei probleme: sporirea resurselor sau reducerea încărcării.

Sporirea resurselor constă în creșterea temporară a lățimii de bandă între anumite puncte ale rețelei. Dar uneori nu este posibilă creșterea capacității de transfer sau aceasta și-a atins deja limitele. În acest moment singura cale de a rezolva congestia este reducerea încărcării prin următoarele metode: interzicerea unor servicii către anumiți utilizatori, degradarea serviciilor pentru o parte sau pentru toți utilizatorii sau planificarea cererilor utilizatorilor într-o manieră mai previzibilă.

Această abordare conduce la împărțirea soluțiilor în două grupe: în buclă deschisă și în buclă închisă.

Soluțiile în buclă deschisă încearcă să rezolve problema printr-o proiectare atentă, asigurând evitarea apariției acesteia. După ce sistemul se pornește și se verifică funcționarea acestuia, nu se mai fac niciun fel de corecții. Se iau decizii fără a se ține cont de starea curentă a rețelei: instrumentele de realizare a controlului în buclă deschisă decid când să se accepte trafic nou, când să se distrugă pachete și care dintre acestea, realizând planificarea deciziilor în diferite puncte din rețea.

Soluțiile în buclă închisă se bazează pe conceptul de reacție inversă, această abordare având trei părți:

- Monitorizarea sistemului pentru a detecta când și unde se produce congestia.
- Trimiterea acestor informații către locurile unde se pot executa acțiuni.
- Corectarea problemei prin ajustarea funcționării sistemului.

Se cunosc numeroși algoritmi pentru controlul congestiei, iar Yang și Reddy au realizat o clasificare a acestora, studiind în principal cele patru clase de algoritmi centrați pe host-uri:

- ✓ Algoritmi centrați pe host-uri
 - Algoritmi în buclă deschisă
 - Algoritmi care acționează asupra sursei
 - Algoritmi care acționează asupra destinației
 - Algoritmi în buclă închisă
 - Algoritmi cu feedback implicit – sursa deduce existența congestiei din observații locale, cum ar fi timpul necesar pentru întoarcerea confirmărilor.
 - Algoritmi cu feedback explicit – pachetele sunt trimise înapoi de la punctul unde s-a produs congestia către sursa, pentru a o avertiza.
- ✓ Algoritmi centrați pe rutere
 - Ruterile înlătură congestia prin renunțarea la pachete.
 - Ruterile semnalează congestia și nu participă activ la înlăturarea ei.

3. Algoritmi de control a congestiei

3.1. Algoritmul RED – Random Early Detection [1] – face parte din clasa algoritmilor de evitare a congestiei centrați pe rutere, unde ruterile au un rol activ în înlăturarea congestiei prin marcarea pachetelor. Această abordare încearcă să prevină formarea de cozi la intrarea în rețea prin renunțarea la pachete înainte de producerea propriu-zisă a fenomenului de congestie. Congestia incipientă este detectată de gateway-uri prin calculul dimensiunii medii a cozii. Acestea pot anunța congestiunea unor conexiuni fie prin renunțarea la pachetele care ajung în dreptul gateway-ului respectiv sau prin trimiterea unui bit cu rol de flag în header-urile pachetelor. Când dimensiunea medie a cozii depășește un anumit prag prestat, gateway-urile renunță sau marchează fiecare pachet care ajunge cu o anumită probabilitate în funcție de valoarea mediei calculate. În acest mod, algoritmul RED păstrează dimensiunea medie a cozii scăzute dar permite acumularea ocazională în coadă a unui număr mare de pachete.

În rețelele de mare viteză se dorește conceperea de gateway-uri care să suporte cozi de dimensiuni cât mai mari pentru a evita congestiunea acestora. Este din ce în ce mai important să se găsească mecanisme pentru o livrare optimă a

pachetelor din punctul de vedere al întârzierilor, dar în același timp să se păstreze o dimensiune medie a cozilor cât mai scăzută.

Cel mai efektiv mod de detectare al congestiei are loc chiar la gateway-urile rețelei, deoarece numai aici se poate vizualiza comportamentul cozilor în timp. Metoda de monitorizare a dimensiunii medii a cozii la gateway-ul rețelei și de a semnaliza conexiuni care se află în faza incipientă de congestie este bazată pe presupunerea că va fi util să avem cozi la gateway-uri unde traficul mai multor conexiuni este multiplexat după algoritmul FIFO. FIFO este util aici la reducerea întârzierii pentru o anumită conexiune atunci când ea este foarte încărcată, prin împărțirea întârzierii totale între toate conexiunile.

Algoritmul RED este conceput pentru rețelele unde marcarea unui pachet este suficientă pentru a semnaliza prezența congestiei. Acest mecanism de control monitorizează dimensiunea medie pentru fiecare coadă de ieșire și alege într-un mod aleatoriu ce conexiune să fie notificată de apariția congestiei. Congestiile tranzitorii sunt eliminate prin creșterea temporară a dimensiunii cozii. Congestiile de lungă durată sunt sesizate prin creșterea dimensiunii medii a cozii, având ca rezultat transmiterea aleatorie de răspunsuri către unele conexiuni pentru a-și micșora ferestrele. Probabilitatea ca o conexiune să fie notificată de apariția congestiei este proporțională cu proporția de participare a conexiunii respective la încărcare.

Algoritmul RED calculează dimensiunea medie a cozii folosind un filtru trecejos de tip EWMA – Exponential/Weighted Moving Average. Dimensiunea medie este comparată cu două praguri, unul de minim și celălalt de maxim. Când media este mai mică decât pragul minim niciun pachet nu este marcat. Când media depășește pragul maxim, toate pachetele care ajung sunt marcate. Dacă media se află între cele două praguri, fiecare pachet care ajunge este marcat cu probabilitatea p_a , unde $p_a = f(\text{media avg})$. Astfel, de fiecare dată când un pachet este marcat, probabilitatea ca un pachet să fie marcat ca venind de la o anumită conexiune este proporțională cu proporția de participare a acelei conexiuni la încărcare.

În cele ce urmează este dat algoritmul general RED:

Pentru fiecare pachet care ajunge

Se calculează dimensiunea medie a cozii avg

Dacă $min \leq avg \leq max$

Se calculează probabilitatea p_a

Se marchează pachetul cu probabilitatea p_a

Altfel

Dacă $max \leq avg$ se marchează pachetul.

Dupa cum se vede algoritmul RED constă in calcularea dimensiunii medii a cozii și a probabilității de marcarea pachetelor. Prin determinarea mediei se determină gradul de încărcare a rețelei. Probabilitatea de marcarea pachetelor arată frecvența cu care pachetele sunt marcate știind că se cunoaște nivelul curent de congestie.

Algoritmul detaliat este următorul:

Inițializare: $avg=0$, $count=-1$

Pentru fiecare pachet care ajunge

Se calculează dimensiunea medie a cozii avg

*Dacă coada nu e goală $avg=(1-w)avg+w*q$*

Altfel $m=f(time-q_time)$ și $avg=(1-w)^m avg$

Dacă $min \leq avg \leq max$

$Count=count+1$

Se calculează probabilitatea pa : $pb=exp(avg-min)/(max-min)$

*$pa=pb/(1-count*pb)$*

Se marchează pachetul probabilitatea pa și $count=0$

Altfel

Dacă $max \leq avg$ se marchează pachetul și $count=0$

Altfel $count=-1$

Când coada devine vidă $q_time=time$.

Variabile salvate: avg , q_time =timpul in care coada e goala, $count$ =numărul de pachete ajunse de la ultimul pachet marcat.

Parametrii ficși: w = „greutatea” cozii, min , max , $maxp=max(pb)$

Alți parametrii: pa , q =dimesiunea cozii curente, $time$ =timpul curent, $f(t)=o$ funcție liniară de t .

Calculul dimensiunii medii ține cont de perioada in care coada e goală, estimând numărul m de mici pachete care ar fi putut fi transmise in acest timp. După perioada in care coada a fost goală se calculează dimensiunea medie ca și cum m pachete ar fi ajuns in coada in acea perioadă. Dupa cum avg variază de la min la max probabilitatea de marcarea pachetelor variază liniar de la 0 la $maxp$: $pb=exp(avg-min)/(max-min)$.

Probabilitatea finală de marcarea a pachetelor pe crește încet în timp ce count crește începând de la marcarea ultimului pachet: $pa=pb/(1-count*pb)$. Se asigură astfel diminuarea timpului de așteptare între marcarea pachetelor.

Există și opțiunea de a măsura coada în bytes și nu în pachete. În acest caz dimensiunea medie reflectă întârzierea medie la gateway. Când se folosește această opțiune, algoritmul trebuie modificat pentru a asigura că probabilitatea ca un pachet să fie marcat să fie proporțională cu dimensiunea pachetului în bytes:

$$pb= \max p(\text{avg}-\text{min})/(\text{max}-\text{min})$$

$$pb=pb * \text{dim_pachet}/\text{dim_max_pachet}$$

$$pa=pb/(1-count*pb)$$

„Greutatea” cozii (w) este determinată de dimensiunea și durata încărcărilor bruște ale cozii care sunt admise la gateway. Pragurile min și max sunt determinate de dimensiunea medie dorită a cozii, care depinde la rândul ei de caracteristicile rețelei.

3.1.1 Calcularea lungimii medii a cozii

Se folosește un filtru trece-jos de tip EWMA – Exponential/Weighted Moving Average:

$$avg=(1-w)avg+w*q$$

„Greutatea” w determină constanta de timp a filtrului trece-jos.

✓ Limita superioară a lui w

Dacă w este prea mare, la gateway nu vor putea fi eliminate congestiile cauzate de încărcări tranzitorii.

Presupunem că inițial coada este goală, cu o medie 0, lungimea acestora crescând de la 0 la L pachete. După ce și ultimul L pachet a ajuns în coadă, o să avem:

$$avgL = \sum_{i=1}^L iw(1-w)^{L-i} = w(1-w)^L \sum_{i=1}^L i \left(\frac{1}{1-w}\right)^i = L+1 + \frac{(1-w)^{L+1} - 1}{w}$$

Mai sus s-a folosit următoarea identitate:

$$\sum_{i=1}^L ix^i = \frac{x + (Lx - L - 1)x^{L+1}}{(1-x)^2}$$

Fiind dat pragul minim \min și știind că admitem încărcări de L pachete la gateway, atunci w trebuie să satisfacă următoarea inecuație:

$$L+1 + \frac{(1-w)^{L+1} - 1}{w} < \min$$

✓ **Limita inferioară a lui w**

Dacă w are setată o valoare prea mică atunci modificările din coadă se vor sesiza greu, deoarece valorile lui avg vor înregistra fluctuații foarte reduse. Astfel stadiile inițiale de congestie vor fi mai greu de depistat.

Presupunem că dimensiunea cozii se schimbă de la 0 la un pachet, că odata cu sosirea și plecarea pachetelor dimensiunea ei de un pachet nu se modifică și că dimensiunea medie inițială a fost 0. În acest caz este nevoie să ajungă $-1/\ln(1-w)$ pachete până când avg ajunge la valoarea $0.63=1-1/e$.

✓ **Setarea pragurilor de min și max**

Valorile optime pentru min și max depind de lungimea medie dorită a cozii. Valoarea optimă pentru max depinde mai ales de întârzierea maximă admisă de gateway. Se folosește în general o valoare a lui max astfel încât să fie dublul lui min.

3.1.2 Calcularea probabilității de marcarea a pachetelor

Probabilitatea inițială de marcarea a pachetelor pb este calculată ca funcție liniară a dimensiunii medii a cozii:

$$pb = \frac{avg - \min}{\max - \min}$$

✓ **Metoda 1. Variabile aleatoare geometrice.**

Fie pb - probabilitatea de marcarea a fiecărui pachet și X - numărul de pachete care ajung între două marcări succesive. X este o variabilă aleatoare geometrică cu parametrul pb și $E[X]=1/pb$.

Deoarece fiecare pachet este marcat cu probabilitatea pb avem:

$$P(X = n) = (1 - pb)^{n-1} pb$$

Cu o dimensiune medie a cozii constantă, scopul este să marcăm pachete la intervale regulate. Nu este de dori să avem pachete marcate nici prea des, dar nici să înregistrăm intervale lungi de timp între marcări succesive. Se realizează astfel o sincronizare globală, anumite conexiuni modificându-și ferestrele în același timp.

✓ **Metoda 2. Variabile aleatoare uniforme.**

In acest caz X este o variabilă aleatoare distribuită uniform in mulțimea $\{1,2,\dots,1/pb\}$. Acest lucru se întâmplă dacă probabilitatea de marcarea a fiecărui pachet este $pb/(1-\text{count} \cdot pb)$, unde count =numărul de pachete nemarcate între două marcări succesive.

In acest caz:

$$P(X = n) = \frac{pb}{1 - (n-1)pb} \prod_{i=0}^{n-2} \left(1 - \frac{pb}{1 - i pb}\right) = pb \text{ pentru } 1 \leq n \leq 1/pb$$

$$P(X = n) = 0 \text{ pentru } n > 1/pb$$

Pentru această metoda $E[X]=1/(2pb)+1/2$.

3.1.3 Implementarea algoritmului RED

Initializare: avg=0 si count=-1

Pentru fiecare pachet care ajunge

Se calculeaza noua dimensiune medie avg

Daca coada nu este goala atunci $avg=avg+w(q-avg)$

Altfel $avg = (1-w)^{(time-q_time)s} avg$

Daca $min \leq avg \leq max$ atunci

$count=count+1$ si $pb=C1 \cdot avg - C2$

Daca $count > 0$ si $count \geq \text{round}(R/pb)$ atunci se marcheaza pachetul care va ajunge si $count=0$

Daca $count=0$ atunci $R=\text{random}[0,1]$

Altfel

Daca $max \leq avg$ atunci se marcheaza pachetul care va ajunge si

$count=-1$

Altfel $count=-1$

Cand coada devine vida $q_time=time$

R – un numar aleator, s – timp de transmisiune tipic.

Concluzie

RED este un mecanism eficient pentru evitarea congestiei la gateway-uri in cooperare cu protocoalele de transport. Este un algoritm relativ simplu care poate fi implementat in rețelele de mare viteză.

Exista încă aspecte de rezolvat in legatură cu acest algoritm: determinarea dimensiunii medii optime a cozii pentru a maximiza numărul de pachete transmise corect și a minimiza întârzierile pentru diverse configurații de rețele, implementarea RED cu alte protocoale de transport inafară de TCP.

3.2 Algoritmul GAIMD – General Additive Increase Multiplicative Decrease [2] este o strategie de evitare a congestiei prin ajustarea dimensiunii ferestrei. Astfel, in starea de evitare a congestiei dimensiunea ferestrei va fi crescută cu un parametru α , iar la apariția propriu-zisă a congestiei dimensiunea acesteia va fi micșorata multiplicativ cu β .

Această tehnică a fost pentru prima data considerată de Chiu&Jain. Ei au demonstrat că dacă α și β respectă următoarele relații, atunci controlul congestiei GAIMD are stabilitate și se realizează corect: $\alpha > 0$ și $0 < \beta < 1$. Studiul lor a considerat numai cazul când toate fluxurile de transmisiuni din rețea folosesc aceleași valori pentru α și β și nu au studiat efectele acestor parametrii asupra performanței. Yang&Lam au continuat studiul, explorând in detaliu relația dintre performanța și diferite valori ale celor doi parametrii. In particular, ei au *modelat rata de transmisiune* in funcție de α , β , a ratei de pierderi, a mediei de timp tur-retur, media timpului de pauză și a numărului de pachete pe care fiecare ACK le confirmă. In continuare au adaptat acest rezultat, găsind o relație simplă intre α și β pentru ca fluxul de date să fie de tip TCP-friendly.

3.2.1. Modelarea ratei de transmisiune

O sesiune GAIMD incepe in starea *slowstart* când dimensiunea ferestrei de congestie este dublată pentru fiecare fereastra de pachete confirmate ca primite. La prima indicare a congestiei, fereastra de congestie este tăiata in două și sesiunea intră in starea de evitare a congestiei. In aceasta stare, fereastra de congestie este mărită cu α/W pentru fiecare ACK primit, unde W este dimensiunea ferestrei de congestie curentă. Vom spune astfel că dimensiunea ferestrei va fi crescută cu α la fiecare tur-retur. GAIMD schimbă dimensiunea ferestrei scăzând din aceasta βW când congestia a fost detectată printr-o confirmare de tip triplu-dublu ACK. Dacă detectarea s-a făcut prin detectarea unui timeout, dimensiunea ferestrei va fi setată la valoarea 1.

$W = W + \alpha/W$, $\alpha > 0$ – evitarea congestiei

$W=W - \beta W$, $0 < \beta < 1$ – congestie detectată printr-o confirmare de tip triplu-dublu ACK

$W=1$ - detectarea unui timeout

La modelarea ratei de transmisiune s-au făcut următoarele presupuneri:

- ✓ Sursa mereu are date de transmis.
- ✓ Rata de transmisiune este un proces aleatoriu.
- ✓ Se ignoră impactul etapei de slowstart, punându-se accentul pe mecanismul de evitare a congestiei.
- ✓ Mecanismul GAIMD de evitare a congestiei lucrează cu timpi de tur-retur.
- ✓ Pierderile într-un tur-retur sunt independente: când un pachet este pierdut înseamnă că toate pachetele care îi urmează în acel tur-retur sunt de asemenea pierdute. Se notează cu p probabilitatea ca un pachet să fie pierdut.

Formula ratei de transmisiune este:

$$T_{\alpha,\beta}(p, RTT, T_0, b) = \frac{1}{RTT \sqrt{\frac{2b(1-\beta)}{\alpha(1+\beta)}} p + T_0 \min(1, 3\sqrt{\frac{(1-\beta^2)b}{2\alpha}} p) p(1+32p^2)}$$

Observăm că numitorul formulei de mai sus este alcătuit din doi termeni pe care îi vom nota astfel:

$$TD_{\alpha,\beta}(p, RTT, b) = RTT \sqrt{\frac{2b(1-\beta)}{\alpha(1+\beta)}} p$$

$$TO_{\alpha,\beta}(p, T_0, b) = T_0 \min(1, 3\sqrt{\frac{(1-\beta^2)b}{2\alpha}} p) p(1+32p^2)$$

Numitorul formulei conține numai TD dacă congestia a fost identificată prin triplu-dublu ACK. Al doilea termen TO este adăugat când congestia poate fi indicată atât prin triplu-dublu ACK cât și prin timeout.

Termenul $Q = \min(1, 3\sqrt{\frac{(1-\beta^2)b}{2\alpha}} p)$ aproximează probabilitatea de timeout.

3.2.2 TCP-friendly GAIMD

Din formula anterioară se observă că este posibil să controlăm perechile (α , β) astfel încât să obținem rata de transmisiune dorită. Putem selecta parametrii α , β astfel încât:

$$T_{\alpha,\beta}(p, RTT, T_0, b) = T_{1, \frac{1}{2}}(p, RTT, T_0, b)$$

Aceste perechi (α, β) care satisfac egalitatea de mai sus formează curba TCP-friendly.

➤ **TD TCP-friendly**

In continuare păstrăm numai primul termen al numitorului, renunțând la variabilele p , RTT și b din ambii membri ai ecuației:

$$TD_{\alpha,\beta}(p, RTT, b) = TD_{1, \frac{1}{2}}(p, RTT, b)$$

$$\frac{1-\beta}{\alpha*(1+\beta)} = \frac{1-0.5}{1*(1+0.5)}$$

In urma calculelor obținem:

$$\alpha = \frac{3(1-\beta)}{1+\beta}$$

➤ **TO TCP-friendly**

Acum păstrăm cel de-al doilea termen al numitorului și renunțăm la variabilele p, T_0 și b din ambii termeni ai ecuației:

$$TO_{\alpha,\beta}(p, T_0, b) = TO_{1, \frac{1}{2}}(p, T_0, b)$$

$$\sqrt{\frac{(1-\beta^2)}{\alpha}} = \sqrt{\frac{(1-0.5^2)}{1}}$$

In urma calculelor obținem:

$$\alpha = \frac{4(1-\beta^2)}{3}$$

3.3 Algoritmii HSTCP – High Speed TCP si STCP – Scalable TCP [3]

3.3.1 HSTCP – High Speed TCP este un algoritm de tip additive increase multiplicative decrease (AIMD), cu deosebirea că cei doi factori α și β depind de dimensiunea ferestrei.

Considerăm w_i ca fiind dimensiunea ferestrei chiar inaintea unei pierderi și RTT_i , timpul de tur-retur, pentru două ferestre $i=1,2$. Fie t , intervalul dintre două pierderi consecutive.

$$w_1 = w_1(1 - \beta(w_1)) + \alpha(w_1) \frac{t}{RTT_1}$$

$$w_2 = w_2(1 - \beta(w_2)) + \alpha(w_2) \frac{t}{RTT_2}$$

Substituind $\alpha(w) = \frac{2w^2\beta(w)p(w)}{2 - \beta(w)}$ și $w = \frac{0.15}{p(w)^{0.82}}$ obținem:

$$\frac{w_1}{w_2} = \left(\frac{2 - \beta(w_2)}{2 - \beta(w_1)} \frac{RTT_2}{RTT_1} \right)^{4.56}$$

3.3.2 STCP – Scalable TCP este un algoritm de tip multiplicative increase multiplicative decrease (MIMD) cu factorii α și β :

$$w_1 = w_1(1 - \beta) + (1 + \alpha) \frac{t}{RTT_1}$$

$$w_2 = w_2(1 - \beta) + (1 + \alpha) \frac{t}{RTT_2}$$

3.6 Algoritmul binomial BITCP – Binary Increase TCP [3] este un algoritm de control al congestiei care adaptează controlul ferestrei în funcție de dimensiunea acesteia. Acesta constă în două părți: căutare prin creștere binară (binary search increase) și creșterea aditivă a dimensiunii ferestrei (additive increase).

3.6.1 Binary search increase

În această etapă privim controlul congestiei ca pe o problemă de căutare în care sistemul poate răspunde prin da sau nu atunci când vrem să stim dacă rata de transmisiune curentă depășește sau nu capacitatea rețelei. Fereastra minimă curentă poate fi estimată ca dimensiunea ferestrei atunci când nu avem pierderi.

Dacă se știe dimensiunea maximă a ferestrei, putem să aplicăm o tehnică de căutare binară pentru a seta dimensiunea ferestrei curente la dimensiunea medie dintre minimul estimat și maximul cunoscut. Dacă această dimensiune dă pierderi, noua fereastră poate fi tratată ca un nou maxim, iar dimensiunea ferestrei după pierderea pachetelor devine noul minim. Media acestor valori devine noul optim.

Acest proces de recalculare a minimului și maximului are loc până când diferența dintre dimensiunea maximă și cea minimă scade sub un nivel prestabilit numit incrementare minimă S_{min} .

3.6.2 Additive increase

Pentru a asigura o convergență rapidă a algoritmului și corectitudinea, se combină căutarea binară cu o strategie de creștere aditivă. Când distanța dintre dimensiunea medie și minimul curent este foarte mare, creșterea dimensiunii ferestrei direct la cea medie este un factor de risc în rețea. Astfel, când distanța dintre dimensiunea ferestrei curente și cea medie este mai mare decât o valoare presetată, numită increment maxim S_{max} , în loc să creștem dimensiunea ferestrei direct la acea valoare medie în următorul RTT, o să o creștem cu S_{max} până când distanța va deveni mai mică decât S_{max} , moment în care dimensiunea va fi setată direct la valoarea medie.

3.6.3 Slow start

Atunci când dimensiunea ferestrei curente devine mai mare decât valoarea maximă, nu cunoaștem noul maxim. În acest moment se va aplica tehnica de căutare binară care va seta un maxim egal cu o valoare presetată și dimensiunea ferestrei curente la o valoare minimă. Vom folosi o strategie de tip „slow start” pentru a căuta o valoare a maximului mai mică decât S_{max} . Dacă considerăm $cwnd$ dimensiunea curentă a ferestrei și incrementul maxim S_{max} , la fiecare RTT dimensiunea ferestrei va crește astfel: $cwnd+1$, $cwnd+2$... $cwnd+S_{max}$. În acest mod se caută banda disponibilă până când se asigură faptul că o creștere cu S_{max} este nedăunătoare. După această etapă se va crește dimensiunea la fiecare RTT cu S_{max} .

3.6.4 Implementare

Se folosesc următorii parametri presetăți:

- low_window – acest algoritm este activat în momentul în care dimensiunea ferestrei curente depășește această valoare
- S_{max} – increment maxim
- S_{min} – increment minim
- β – factorul de scădere multiplicativă a dimensiunii ferestrei
- $default_max_win$ – maximul presetat

Se folosesc următoarele variabile:

- max_win – dimensiunea maximă a ferestrei care la început a valoarea presetată
- min_win – dimensiunea minimă a ferestrei
- $prev_win$ – dimensiunea maximă chiar înainte de setarea noului maxim
- $target_win$ – media dintre minim și maxim

- `cwnd` – dimensiunea ferestrei de congestie
- `is_BITCP_ss` – variabila booleana care indica daca suntem sau nu in etapa slow start. Se initializeaza cu false.
- `ss_cwnd` – o variabila care stocheaza cu cat se creste `cwnd` curent in etapa slow start
- `ss_target` – stocheaza fiecare valoare a `cwnd` dupa fiecare RTT in etapa slow start

Algoritmul BI-TCP:

Când se intră in fast recovery:

```

if(low_window<=cwnd){
    prev_max=max_win;
    max_win=cwnd;
    cwnd=cwnd*(1-β);
    min_win=cwnd;
    if(prev_max>max_win) max_win=(max_win+min_win)/2;
    target_win=(max_win+min_win)/2;
} else { cwnd=cwnd*0.5; }

```

Cand nu se află in fast recovery și ajunge un ACK pentru un nou pachet:

```

if(low_window>cwnd) { cwnd=cwnd+1/cwnd; return;}
if(is_BITCP_ss is false ){
    if(target_win-cwnd<Smax) cwnd+=(target_win-cwnd)/cwnd;
    else cwnd+=Smax/cwnd;
    if(max_win>cwnd){
        min_win=cwnd;
        target_win= (max_win+min_win)/2;
    } esle { is_BITCP_ss=true;
    ss_cwnd=1;
    ss_target=cwnd+1;
}

```

```

        max_win=default_max_win;
    }
} else {
    cwnd+=ss_cwnd/cwnd;
    if(cwnd>=ss_target){
        ss_cwnd=2*ss_cwnd;
        ss_target+=ss_cwnd;
    }
    if(ss_cwnd>=Smax) is_BITCP_ss=false;
}

```

3.7 Algoritmul SIMD – Square Increase Multiplicative Decrease [4] este un mecanism de control al congestiei cu memorie. Folosește micșorarea multiplicativă precum AIMD dar creșterea dimensiunii ferestrei se face cu rădăcina timpului scurs de la detectarea ultimei pierderi. SIMD este TCP-friendly și TCP-compatibil. Are un comportament convergent mult mai bun decât TCP-friendly GAIMD și algoritmul binomial descrise în subcapitolele anterioare.

SIMD folosește pe lângă dimensiunea ferestrei curente și informații memorate referitoare la conexiunile rețelei. Această informație este reprezentată de dimensiunea ferestrei la momentul detectării ultimei pierderi. Fereastra este micșorată multiplicativ dar creșterea se realizează cu rădăcina timpului scurs de la detectarea ultimei pierderi.

Algoritmul SIMD

Creștere: $W = W + \alpha\sqrt{W - W_0}$, $\alpha > 0$ unde W_0 este dimensiunea ferestrei după ultima micșorare a acesteia.

Micșorare: $W = W - \beta W$, $0 < \beta < 1$

Fie $w(t)$ aproximarea continuă a dimensiunii ferestrei la momentul t scurs din momentul începerii creșterii ferestrei și $w = w(0)$.

Folosind interpolarea liniară și aproximarea continuă în formula de creștere, obținem:

$$\frac{dw(t)}{dt} = \alpha\sqrt{w(t) - w_0}, \text{ de unde rezultă: } \frac{1}{\sqrt{w(t) - w_0}} dw(t) = \alpha dt .$$

După integrare se obține: $2\sqrt{w(t)-w_0} = \alpha t + C$. Pentru $t=0$ avem $w(t)=w(0)$ de unde rezultă $C=0$.

Se obține: $w(t) = w_0 + \frac{\alpha^2}{4} t^2$. Se poate astfel spune că SIMD devine „agresiv” cu timpul adică poate să pună la dispoziție lățime de bandă în plus când dispune de aceasta.

SIMD TCP-friendly – definirea lui α atunci când se cunosc β și variabila de stare w_{\max} .

Definim $\alpha = \frac{3\sqrt{\beta}}{(1-2\beta/3)\sqrt{2w_{\max}}}$ și o înlocuim în formula lui $w(t)$:

$$w(t) = w_0 + \frac{9\beta}{8(1-2\beta/3)^2 w_{\max}} t^2.$$

Din ultima formulă putem spune că SIMD poate fi privit ca un AIMD al cărui parametru α variază.

Concluzie

SIMD este un algoritm cu memorie care converge mai rapid decât algoritmi fără memorie AIMD și binomial. De exemplu, dacă există două fluxuri de tip SIMD competitive atunci cel cu dimensiunea ferestrei mai mică este mai „agresiv”. Această proprietate duce la obținerea unui comportament convergent mai bun.

3.9 Bibliografie

- [1] „Random Early Detection Gateways for Congestion Avoidance”, Sally Floyd, Van Jacobson
- [2] „General AIMD Congestion Control”, Yang Richard Yang, Simon S. Lam
- [3] „Binary Increase Congestion Control for Fast, Long Distance Networks”, Lisong Xu, Khaled Harfoush, Injong Rhee
- [4] „TCP - friendly SIMD congestion Control and Its Convergence Behavior”, Shudong Jin, Liang Guo, Ibrahim Matta, Azer Bestavros
- [5] „Computer networks”, Andrew S. Tanenbaum, David J. Wetherall, Fifth Edition, 2011