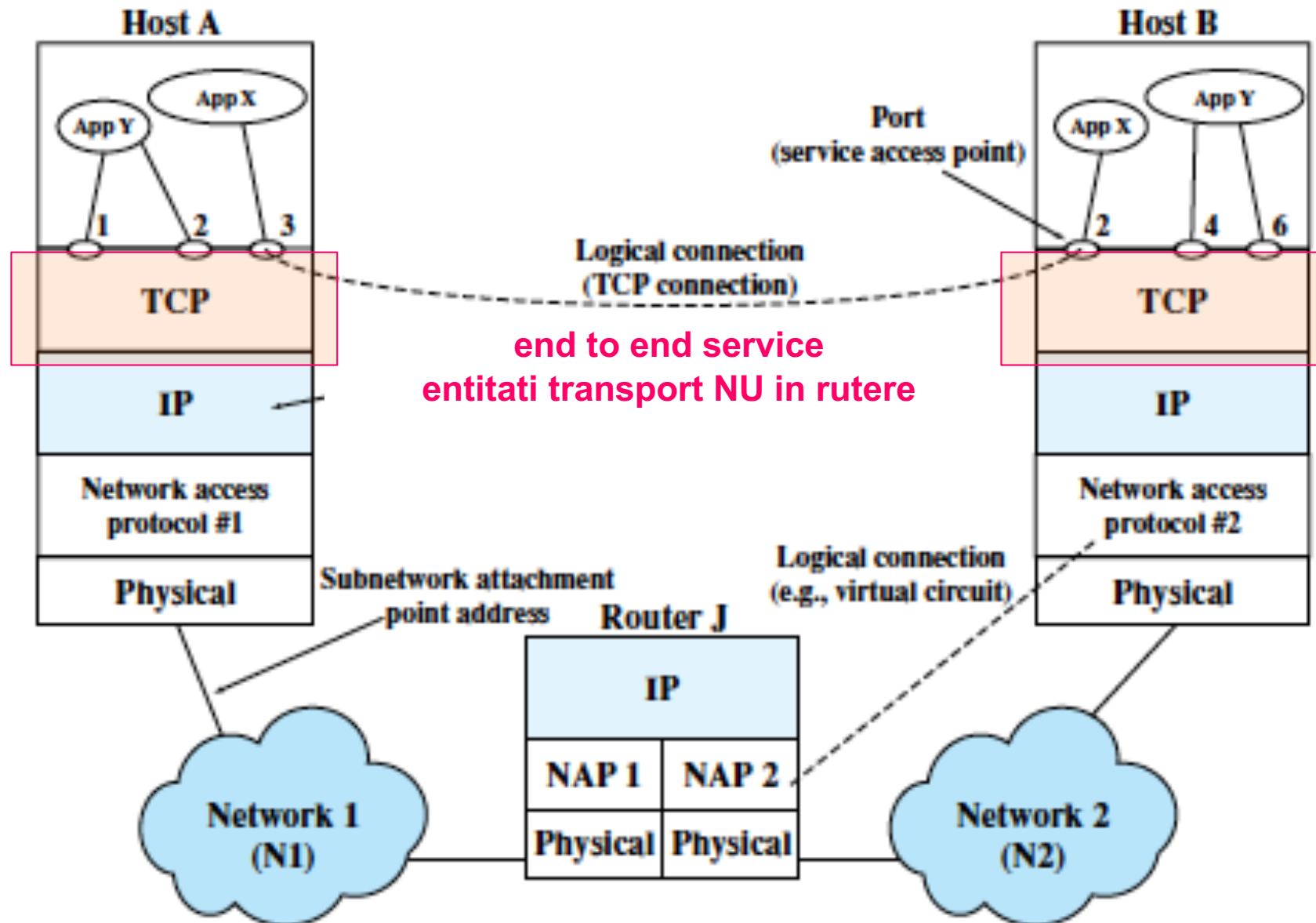




Nivelul transport

Nivelul transport in ierarhia TCP/IP

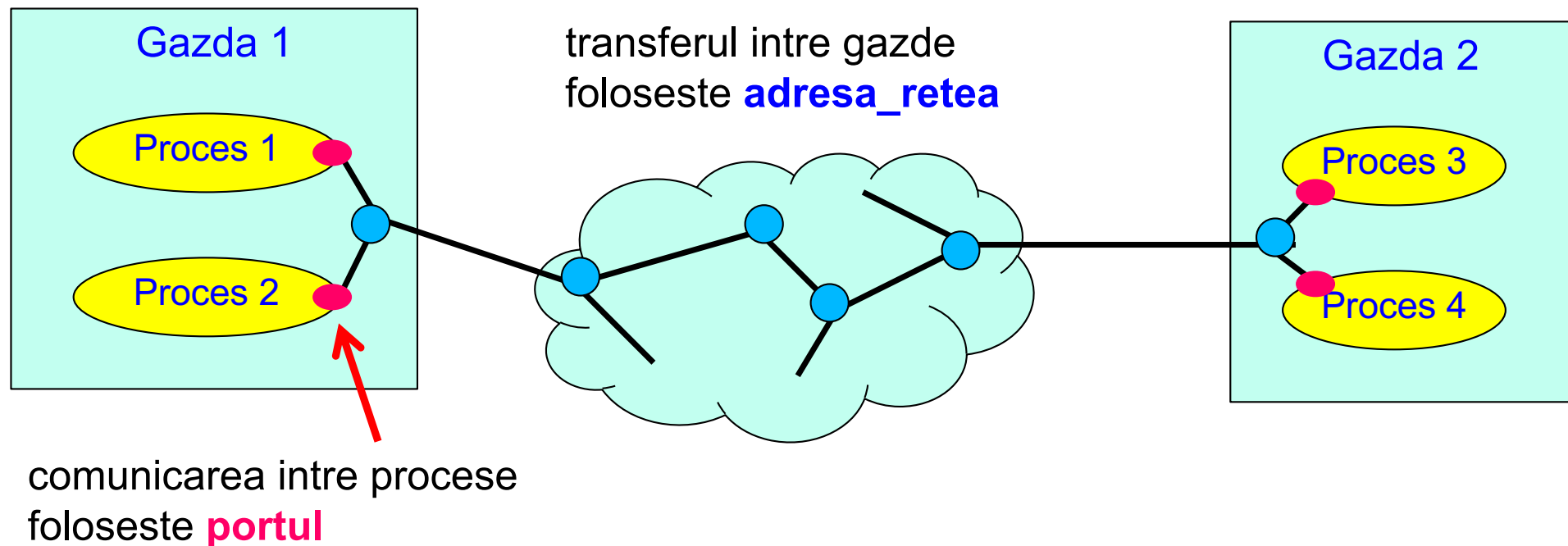


Comunicarea între aplicații

Reteaua asigură transmiterea **pachetelor** (datagrame) între **calculatoare** gazda

Transportul asigură comunicarea între aplicații

- legătura proceselor de aplicație cu **punctele** de acces la rețea
- **identificarea** unică a unui **punct de acces** prin **<adresa_retea, port>**





Servicii ale nivelului Transport

- **Servicii furnizate**

- transfer de date între procese de aplicație, folosind rețele de diverse tipuri
- interfața uniformă cu utilizatorii

- **Caracteristici**

- două tipuri de servicii:
 - orientate pe conexiune (connection oriented) - TCP
 - fără conexiune (connectionless) - UDP



Socket interface

- Serviciile nivelului transport sunt accesibile ca **API** – Application Programming Interface
- Oferită ca bibliotecă utilizator sau funcții OS
 - API descrie cum se apelează aceste funcții
- *Socket API*
 - Originară din Berkeley BSD UNIX
 - Disponibilă și pe Windows, Solaris, etc.
 - **socket** = punctul în care procesul de aplicație se atașază la rețea
 - identificat prin descriptor - **număr** (ca la fișiere)
- Nu e standard de jure ci *standard de facto*



Creare socket

`int socket(int family, int type, int protocol);`

`socket_descr = socket (protocol_family, comm_type, protocol)`

- deschide un socket
- intoarce `socket_descriptor` folosit in apelurile urmatoare

`protocol_family` selecteaza familia de protocoale

`PF_INET` - protocoale **Internet**

`PF_APPLETALK` - protocoale AppleTalk etc.

`comm_type` selecteaza tipul de comunicare

`SOCK_DGRAM` - fara conexiune - datagrama

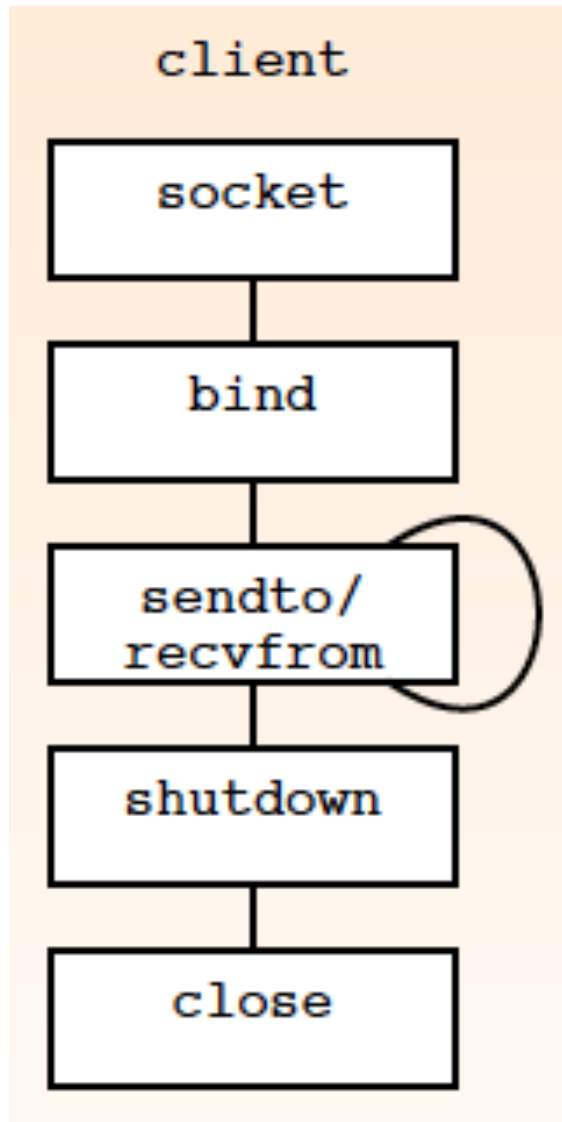
`SOCK_STREAM` - orientat pe conexiune – flux de octeti

`protocol` specifica protocolul

`IPPROTO_TCP` - TCP

`IPPROTO_UDP` - UDP

Serviciu fara conexiune

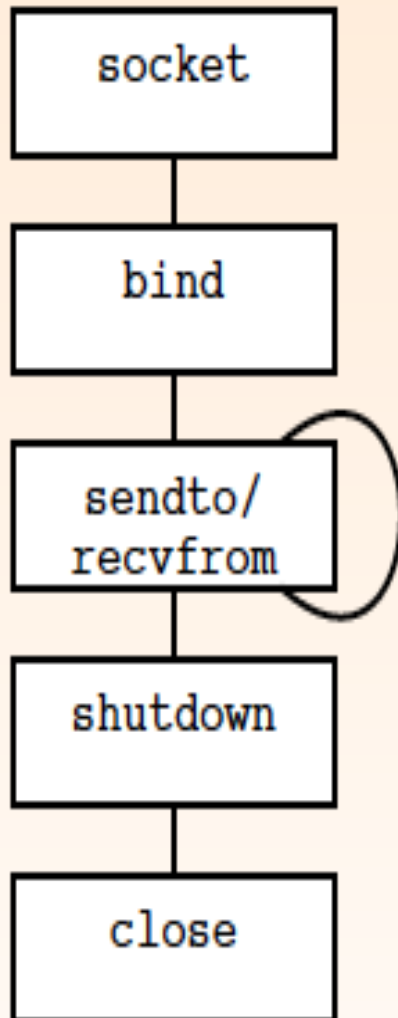


server este la fel!

bind - optional pentru initiator

Serviciu fara conexiune - server

server



Creaza socket si aloca-i resurse de sistem:

```
int s = socket (AF_INET, SOCK_DGRAM, IPPROTO_UDP);
```

Asociaza un socket cu <port, adresa_IP>:

```
int bind (int socket_descriptor, struct sockaddr* local_address, int  
address_length)
```

```
bind(s, &addr, sizeof(addr));           // addr locala - vezi slide urmator !
```

Primește mesaje de la un socket aflat la distanță:

```
int recvfrom (int socket_descriptor, char* buffer_address, int  
buffer_length, int flags, struct sockaddr* sender_address, unsigned int  
sendaddress_length)
```

```
recvfrom(s, buf, BUFLen, 0, NULL, NULL);
```

Oprește trimitere sau/și recepție de date

```
shutdown (s, SHUT_RD / SHUT_RDWR / SHUT_WR)
```

Inchide socket - termina utilizarea socket si elibereaza resurse alocate

```
close (s)
```




Setare adresa

Format TCP/IP:

```
struct sockaddr_in { u_char  sin_len;      /* total length of address */
                    u_char  sin_family;    /* family of the address */
                    u_short sin_port;      /* protocol port number */
                    struct  in_addr sin_addr; /* IP address */
                    char    sin_zero[8];    /* unused */
}
```

```
struct sockaddr_in serv_addr
memset ((char *) &serv_addr, 0, sizeof(serv_addr));
serv_addr.sin_family = AF_INET;           // adrese pentru Internet
serv_addr.sin_addr.s_addr = INADDR_ANY;
                                           // foloseste adresa IP a masinii
serv_addr.sin_port = htons(portno);
                                           // converteste de la host la network byte order
```

Exemplu client

```
1  #include <sys/types.h>
2  #include <sys/socket.h>
3  #include <netinet/in.h>
4  #include <netdb.h>
5
6  #define SERVER PORT 12345 /* arbitrary, but client & server must agree */
7  #define BUF SIZE 4096 /* block transfer size */
8
9  int main(int argc, char **argv) {
10
11     int c, s, bytes;
12     char buf[BUF SIZE]; /* buffer for incoming file */
13     struct hostent *h; /* info about server */
14     struct sockaddr in channel; /* holds IP address */
15
16     if (argc != 3) fatal("Usage: client server-name file-name");
17     h = gethostbyname(argv[1]); /* look up host's IP address */
18     if (!h) fatal("gethostbyname failed");
19
20     s = socket(PF_INET, SOCK_STREAM, IPPROTO_TCP);
21     if (s < 0) fatal("socket");
22
23     memset(&channel, 0, sizeof(channel));
24     channel.sin family= AF_INET;
25     memcpy(&channel.sin addr.s addr, h->h addr, h->h length);
26     channel.sin port= htons(SERVER PORT);
27     c = connect(s, (struct sockaddr *) &channel, sizeof(channel));
28     if (c < 0) fatal("connect failed");
29
30     /* Connection is now established. Send file name including 0 byte at end. */
31     write(s, argv[2], strlen(argv[2])+1);
32     /* Go get the file and write it to standard output. */
33     while (1) {
34         bytes = read(s, buf, BUF SIZE); /* read from socket */
35         if (bytes <= 0) exit(0); /* check for end of file */
36         write(1, buf, bytes); /* write to standard output */
37     }
38 }
39
40 fatal(char *string) {
41     printf("%s\n", string);
42     exit(1);
43 }
```

Exemplu server

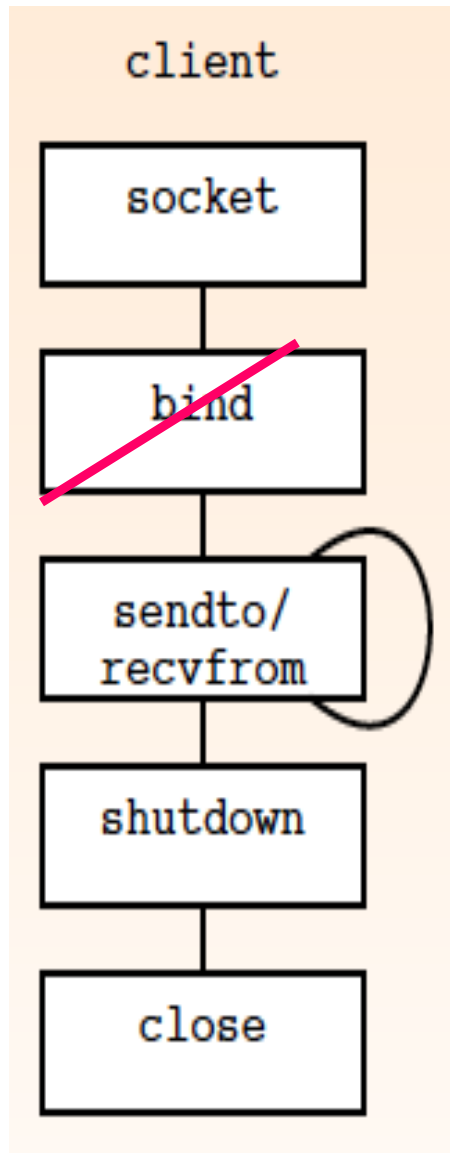
```

2  #include <sys/fcntl.h>
3  #include <sys/socket.h>
4  #include <netinet/in.h>
5  #include <netdb.h>
6
7  #define SERVER PORT 12345 /* arbitrary, but client & server must agree */
8  #define BUF SIZE 4096 /* block transfer size */
9  #define QUEUE SIZE 10
10
11 int main(int argc, char *argv[]) {
12
13     int s, b, l, fd, sa, bytes, on = 1;
14     char buf[BUF SIZE]; /* buffer for outgoing file */
15     struct sockaddr in channel; /* holds IP address */
16
17     /* Build address structure to bind to socket. */
18     memset(&channel, 0, sizeof(channel)); /* zero channel */
19     channel.sin family = AF_INET;
20     channel.sin addr.s addr = htonl(INADDR ANY);
21     channel.sin port = htons(SERVER PORT);
22     /* Passive open. Wait for connection. */
23     s = socket(AF_INET, SOCK_STREAM, IPPROTO TCP); /* create socket */
24     if (s < 0) fatal("socket failed");
25
26     setsockopt(s, SOL_SOCKET, SO_REUSEADDR, (char *) &on, sizeof(on));
27     b = bind(s, (struct sockaddr *) &channel, sizeof(channel));
28     if (b < 0) fatal("bind failed");
29
30     l = listen(s, QUEUE SIZE); /* specify queue size */
31     if (l < 0) fatal("listen failed");
32
33     /* Socket is now set up and bound. Wait for connection and process it. */
34     while (1) {
35         sa = accept(s, 0, 0); /* block for connection request */
36         if (sa < 0) fatal("accept failed");
37
38         read(sa, buf, BUF SIZE); /* read file name from socket */
39         /* Get and return the file. */
40         fd = open(buf, O_RDONLY); /* open the file to be sent back */
41         if (fd < 0) fatal("open failed");
42         while (1) {
43             bytes = read(fd, buf, BUF SIZE); /* read from file */
44             if (bytes <= 0) break; /* check for end of file */
45             write(sa, buf, bytes); /* write bytes to socket */
46         }
47         close(fd); /* close file */
48         close(sa); /* close connection */
49     }
50 }

```



Serviciu fara conexiune - client



Creaza socket și aloca-i resurse de sistem:

```
int s = socket (AF_INET, SOCK_DGRAM, IPPROTO_UDP);
```

Trimite mesaj la un socket aflat la distanță:

```
int sendto (int socket_descriptor, char* data_address, int data_length, int flags,  
struct sockaddr* dest_address, int destaddress_length)
```

```
sendto (s, buf, BUFLen, 0, &addr, sizeof(addr));
```

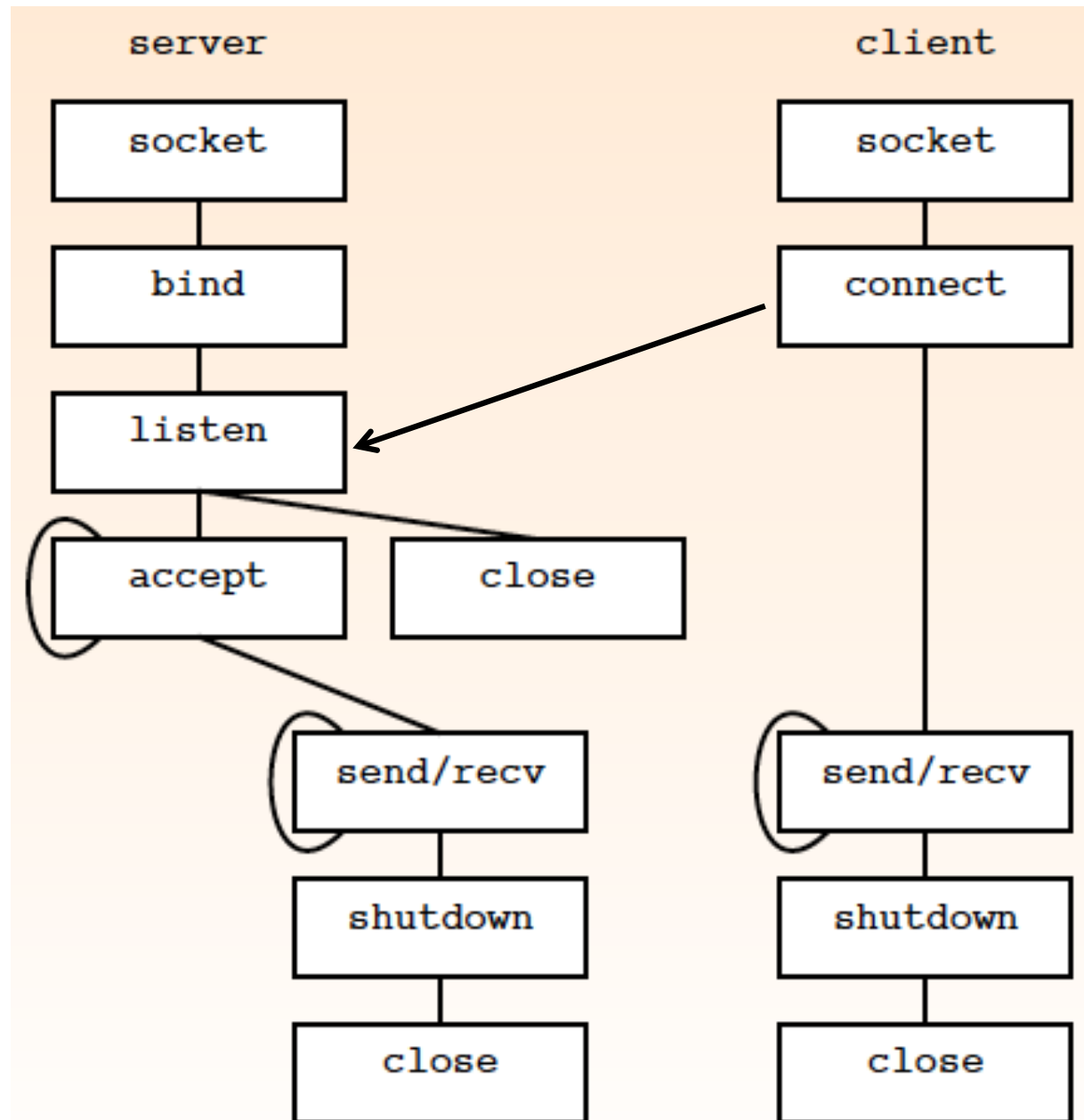
Opreste trimitere sau/și receptie de date

```
shutdown (s, SHUT_RD / SHUT_RDWR / SHUT_WR)
```

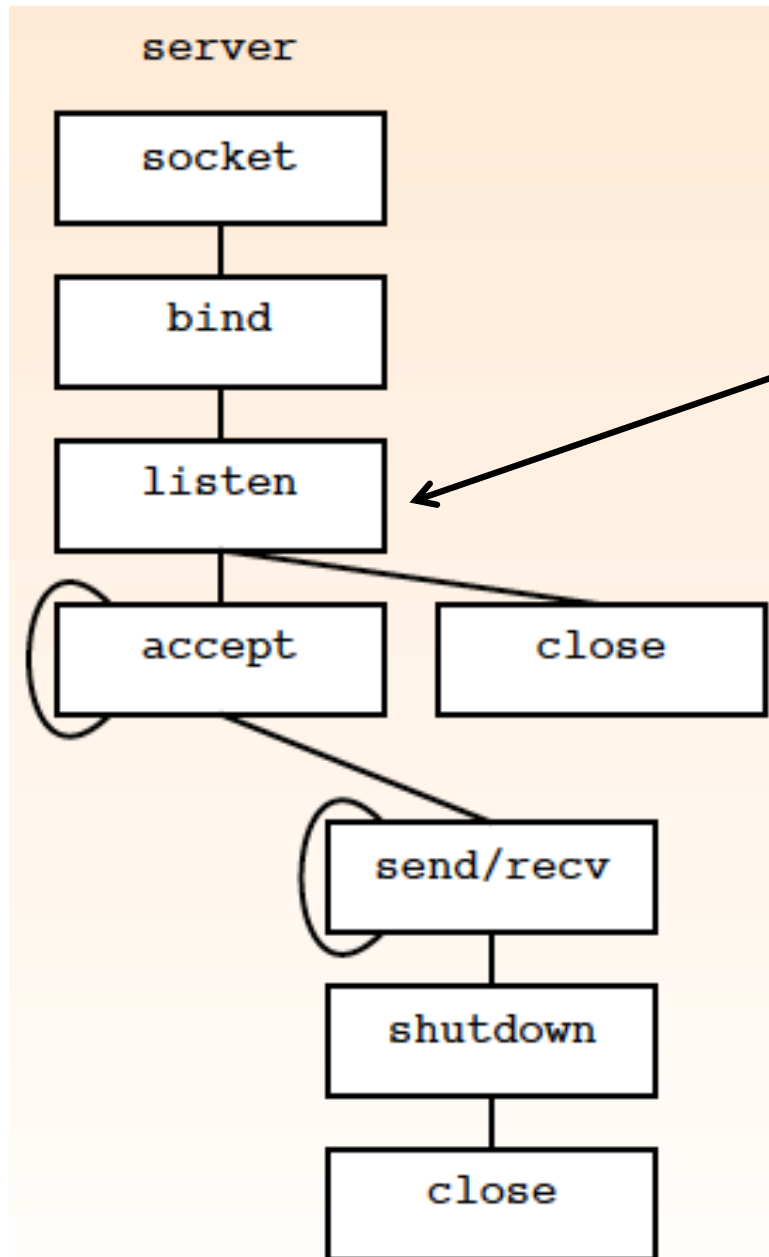
Termina utilizarea socket si elibereaza resursele alocate

```
close (s)
```

Serviciu orientat pe conexiune



Serviciu orientat pe conexiune - server



1. Creaza socket:

```
int ls = socket (AF_INET, SOCK_STREAM,
                 IPPROTO_TCP);
```

2. Asociaza un socket cu <port, adresa_IP>:

```
bind(ls, &addr, sizeof(addr));
```

3. Asteapta cereri de conectare la socket de ascultare (declara nr max cereri asteapta):

```
int listen (int socket_descriptor, int queue_size)
```

```
listen(ls, 5);
```

4. repetat - Accepta o cerere de conectare de la un client; creaza un nou socket pentru conexiune:

```
int accept (int listen_socket_descriptor, struct sockaddr*
            client_socket_addr, int* client_addrlen)
```

```
int s = accept(ls, NULL, NULL);
```

5. repetat - Trimite / primește pe s etc.

Serviciu orientat pe conexiune - client

1. Creaza socket:

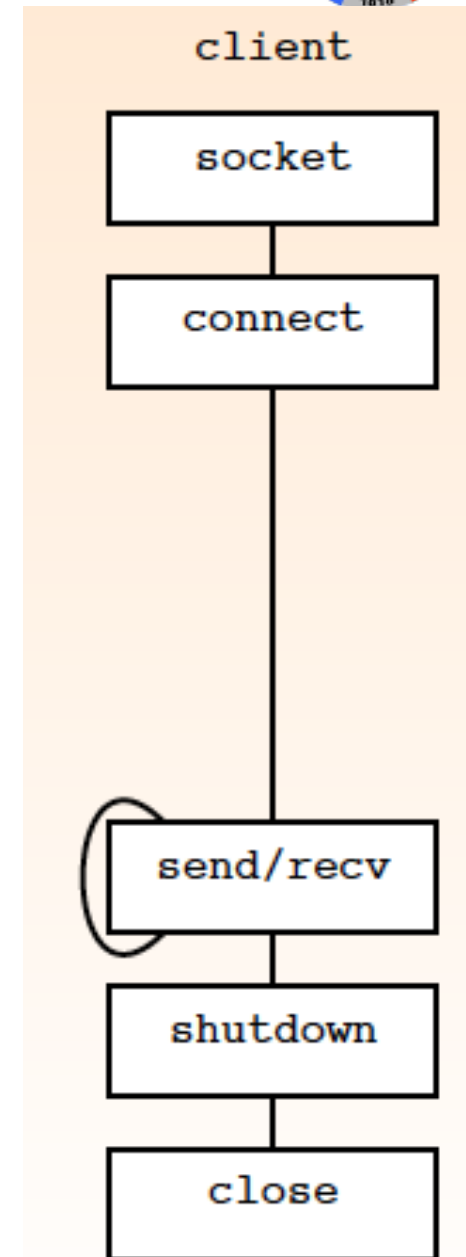
```
int s = socket (AF_INET, SOCK_STREAM,  
                IPPROTO_TCP);
```

2. Conecteaza un socket client (specificat prin descriptor) cu un socket server (precizat prin adresa):

```
int connect (int client_socket_descriptor, struct sockaddr*  
            server_socket_address, int  
            server_sockaddress_length)
```

```
connect(s, &addr, sizeof(addr));
```

3. repetat - Trimite / primește etc.





Transmisia de date cu TCP

send (s, buf, len)

`int send(int socket, const char* buf, int len, int flags);`

- Intoarce numarul de octeti trimisi
 - **Poate fi mai mic decat len !**

recv (s, buf, max_len)

`int recv(int socket, char* buf, int len, int flags);`

- Intoarce numarul de octeti primiti
 - **Poate fi mai mic decat max_len!**

`flags` indica optiuni speciale

MSG_OOB – trimite/primește date out-of-band

MSG_PEEK – livrează date primite, dar tratează ca necitite



Inchiderea conexiunii TCP

- Elibereaza resursele asociate conexiunii
- Informeaza capatul celalalt de inchiderea conexiunii
- API
 - **shutdown** (s, SHUT_RD/SHUT_RDWR/SHUT_WR)
`int shutdown (int socket, int how)`
 - opreste primirea – rejecteaza datele care sosesc
 - opreste transmiterea – ignora datele inca netrimise
 - ambele
 - **close** (s)
`int close (int socket);`
 - inchide socket (elibereaza structurile de date din kernel)



Protocoale de Transport

Doua potocoale majore

- UDP
- TCP

Aspecte discutate

- Formatul datelor
- Adresare
- Functionare



UDP - User Datagram Protocol

- UDP livrează datagrame utilizator - *user datagrams*
 - Livrare “Best effort” – datagramele pot fi pierdute, primite în altă ordine etc.
 - Sume de control pentru integritate
- Puncte de capăt UDP = *protocol ports* sau *ports*
- UDP identifică *adresa Internet* și *număr port* pentru sursă și destinație
- *Destination port* și *source port* pot diferi.

Antet UDP



checksum (calculat la fel ca la TCP) dar nefolosit

Nu control flux

Nu control erori

Nu retransmisie



Aplicatii care folosesc UDP

- DNS
- Voice over IP
- Online Games
- Se foloseste atunci cand:
 - Latența este foarte importantă
 - Livrarea tuturor datelor nu e necesara
 - Retransmisiile sunt implementate de aplicatii cand e nevoie (DNS)



TCP - Transmission Control Protocol

Livrare sigura pe retea nesigura (datagrame)

- ***Cel mai folosit protocol de transport***

Web

Email

SSH

Chat

Video streaming

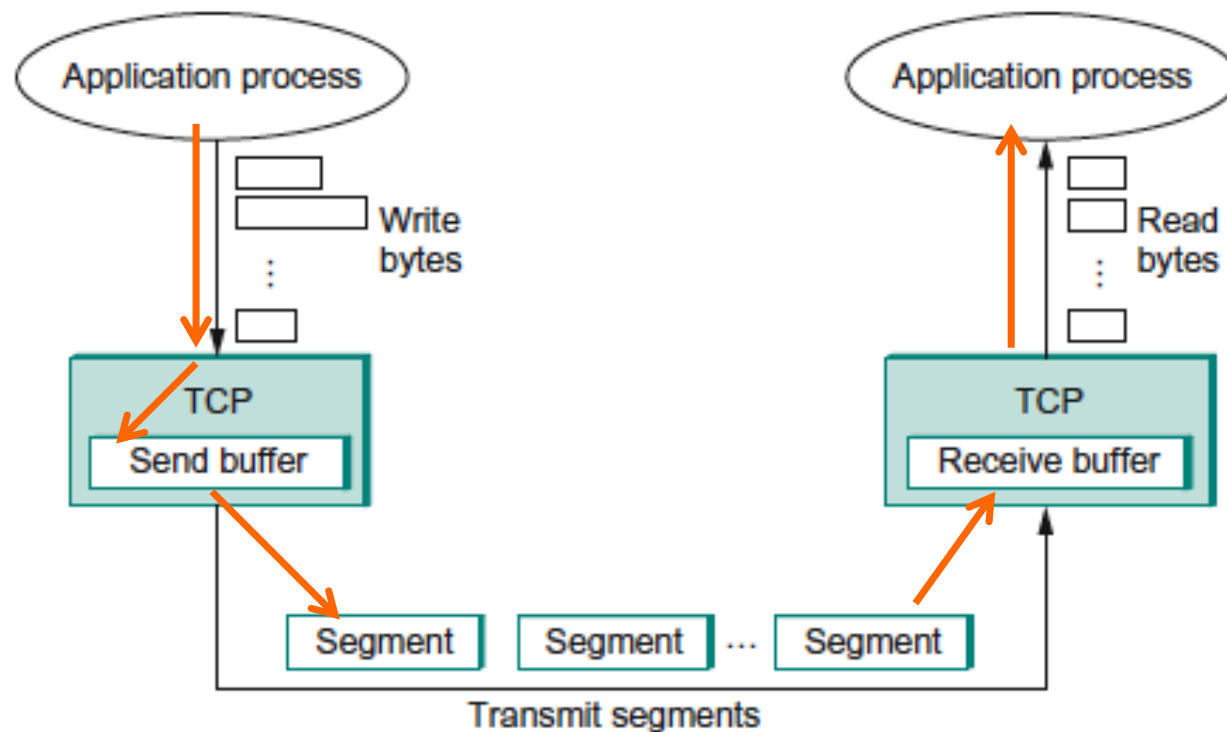
Peer-to-peer



Cateva porturi standard

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

TCP este orientat pe flux de octeți



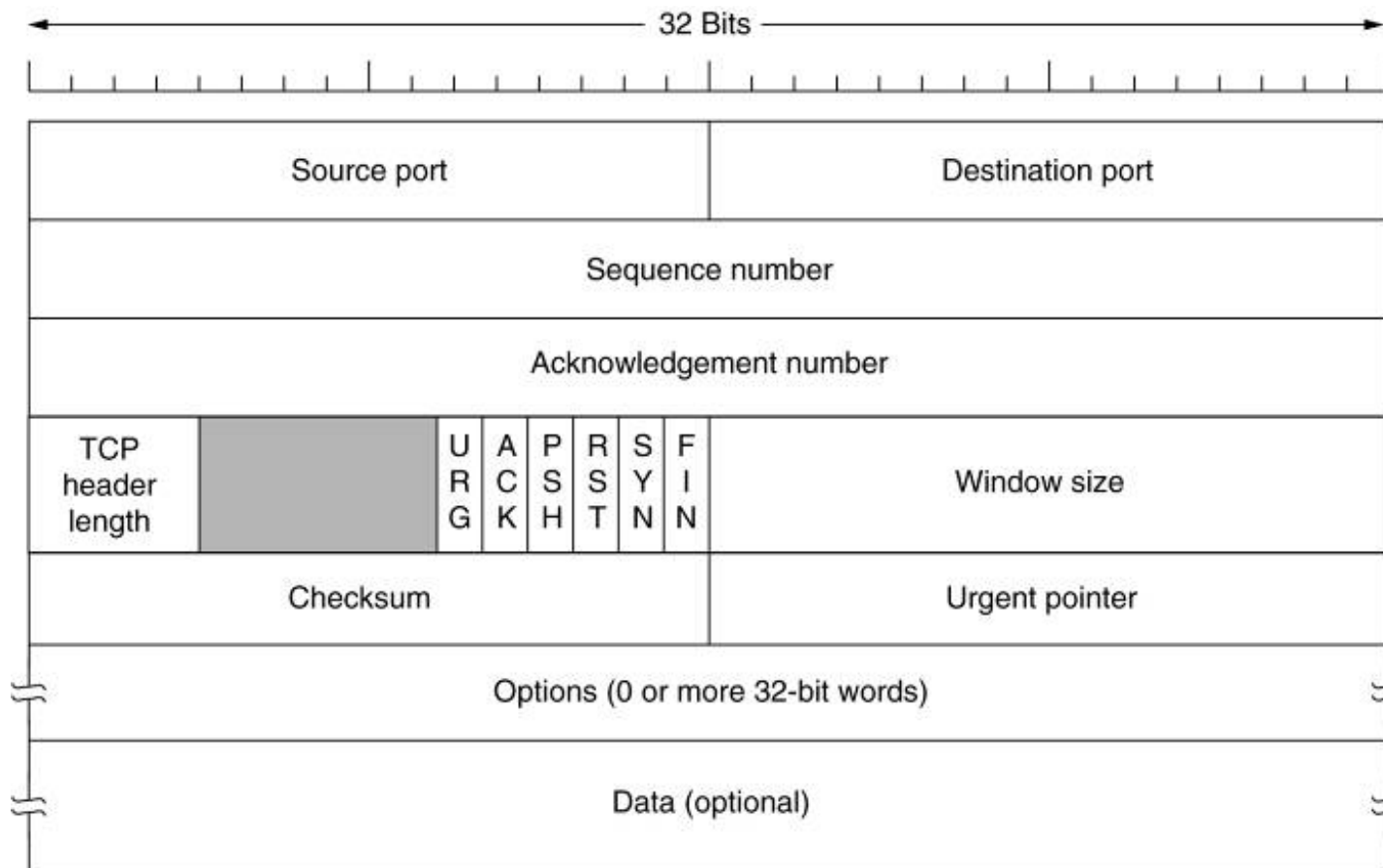
- Aplicația transmițătoare scrie octeți în conexiunea TCP
- TCP la sursă memorează octeții într-un buffer și transmite segmente
- TCP la destinație memorează segmentele într-un buffer
- Aplicația receptoare citește octeți (câți vrea!) din conexiunea TCP



Caracteristici

- Orientat pe conexiune
- Interfață **flux** (Stream)
 - transmisie si receptie siruri de octeti
- Face **controlul congestiei** adaptând viteza de transmisie la condițiile rețelei
- Garantează **transmisie în ordine și sigură** a datelor pe o conexiune
- **Full duplex**
- Stabilire sigură a conexiunii - three-way handshake
- Eliberare lină a conexiunii – fără pierdere de date

Antet segment TCP



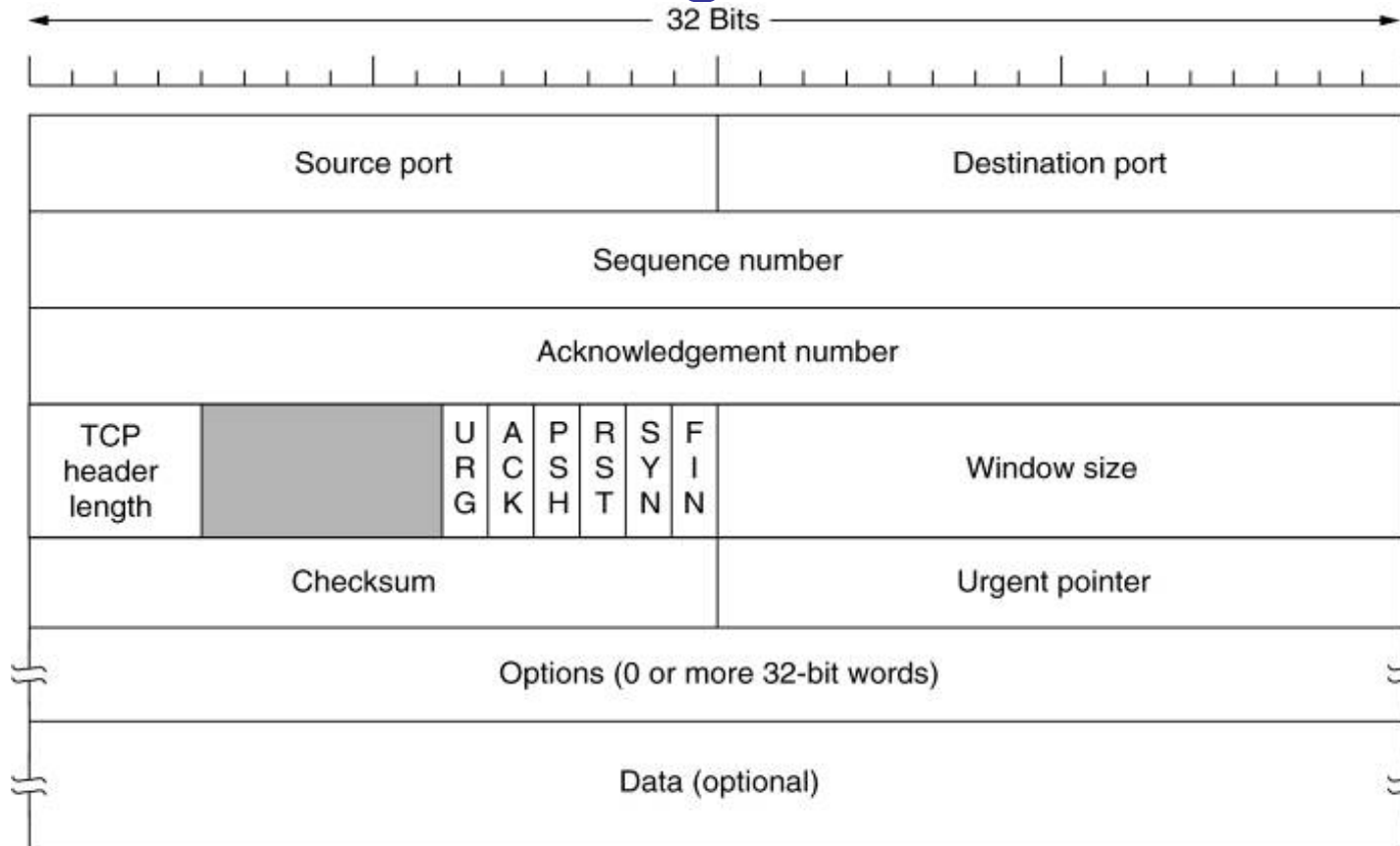
Sequence number – numărul **primului** octet din segment

Acknowledgement number – numărul **următorului** octet așteptat

Window size - numărul de **octeți** care pot fi trimiși, începând cu octetul confirmat

Urgent Pointer – **deplasamentul**, față de *Sequence number*, ptr. info. urgentă

Antet segment TCP



URG Urgent pointer valid

ACK Acknowledge Number valid

PSH - push information to user

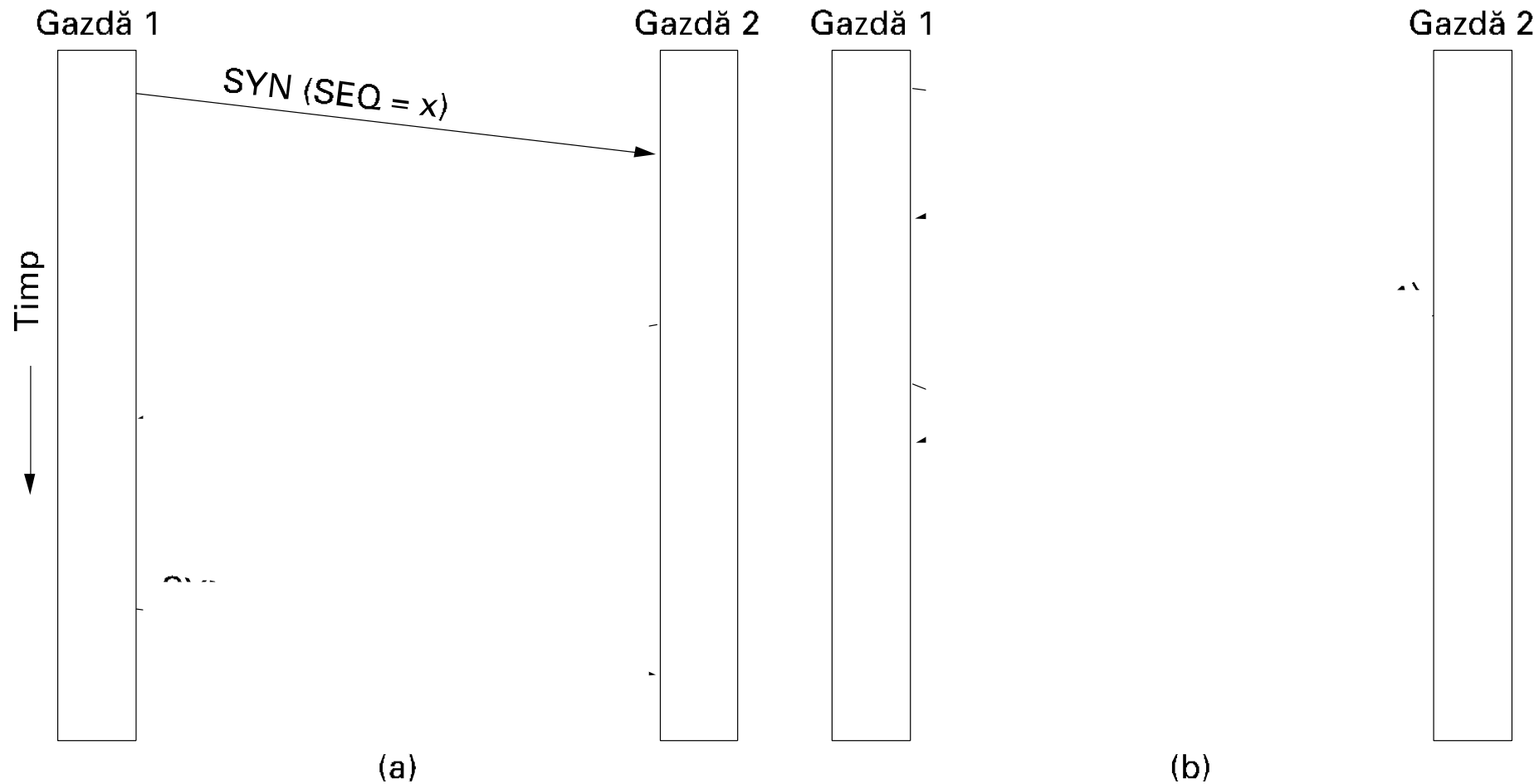
RST - close a connection due to an error

SYN - open connection

FIN - close a connection

Options: e.g. max TCP payload (implicit 536 octeti), selective repeat

Stabilirea conexiunii - Three way handshaking



(a) Cazul normal.

(b) Coliziune.

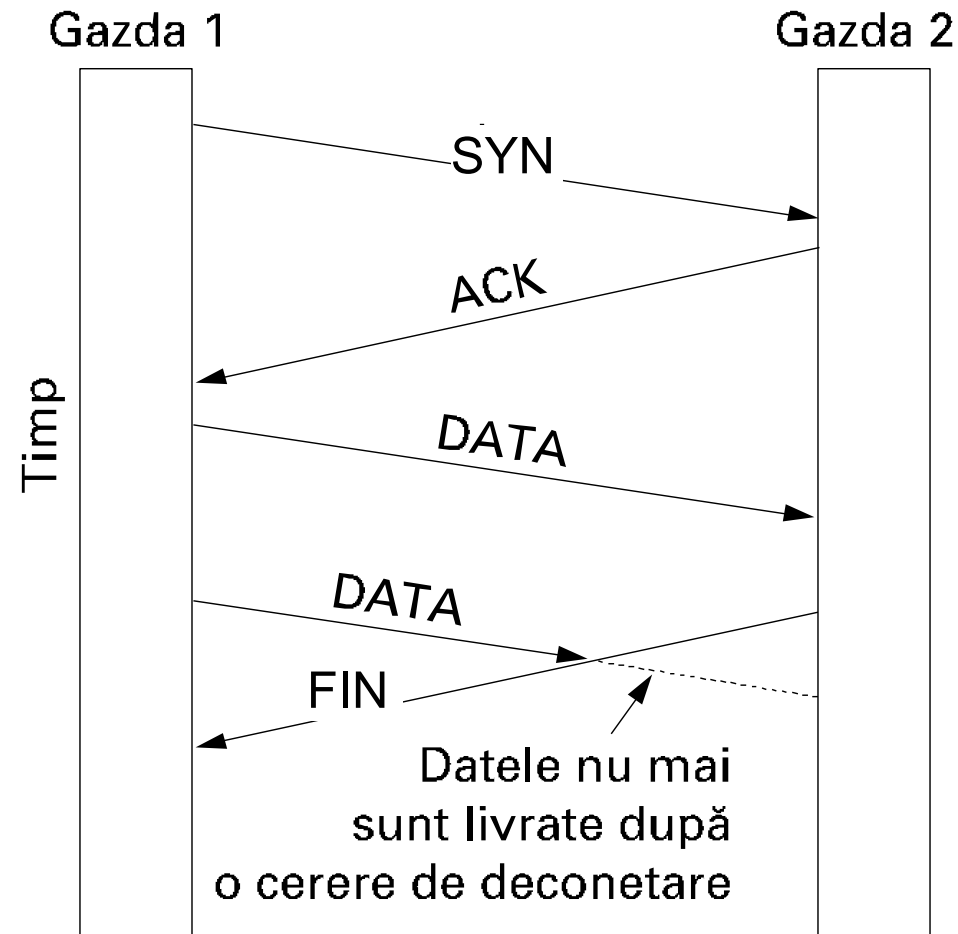
Rejectarea conexiunii



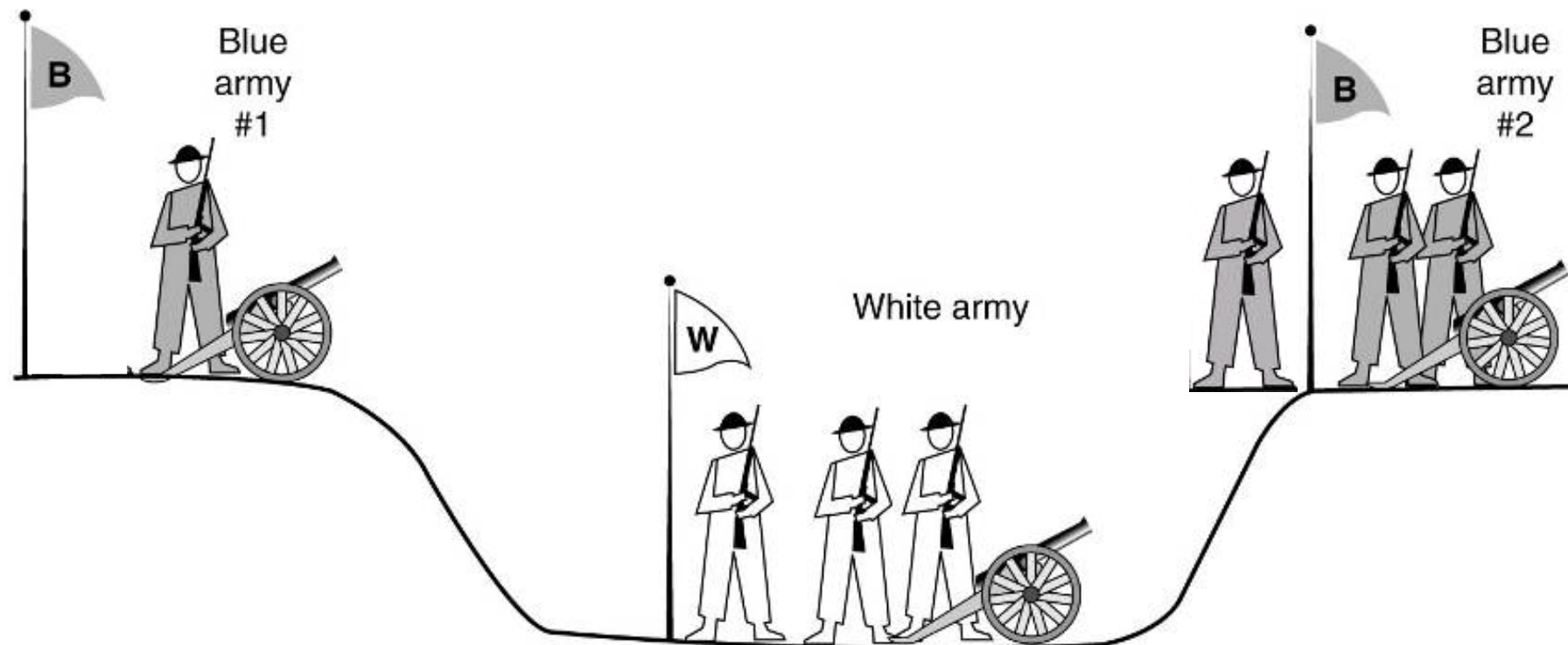
SYN intarziat
B accepta
A rejecteaza

A initiaza conexiune (SYN_i) + SYN_k intarziat
A refuza (RST)
B accepta SYN_i – raspunde cu SYN_j

Deconectare abruptă cu pierdere de date

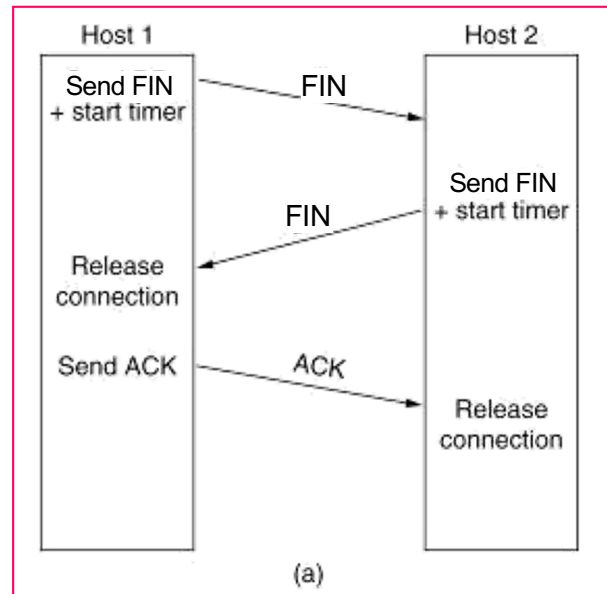


Problema celor doua armate

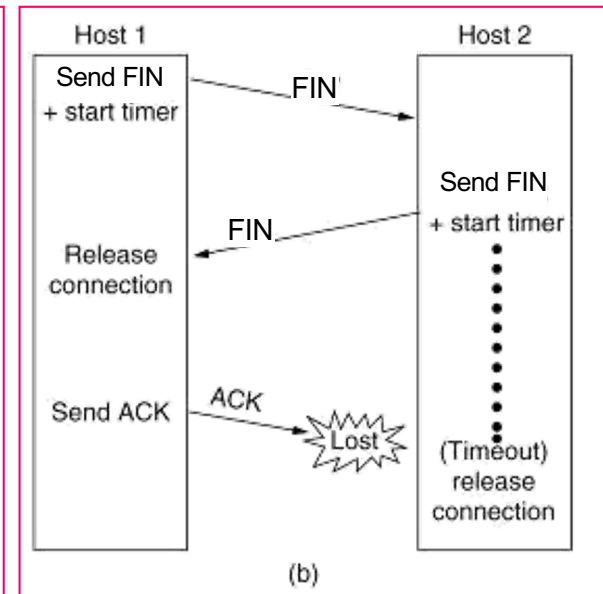


Patru scenarii de eliberarea conexiunii

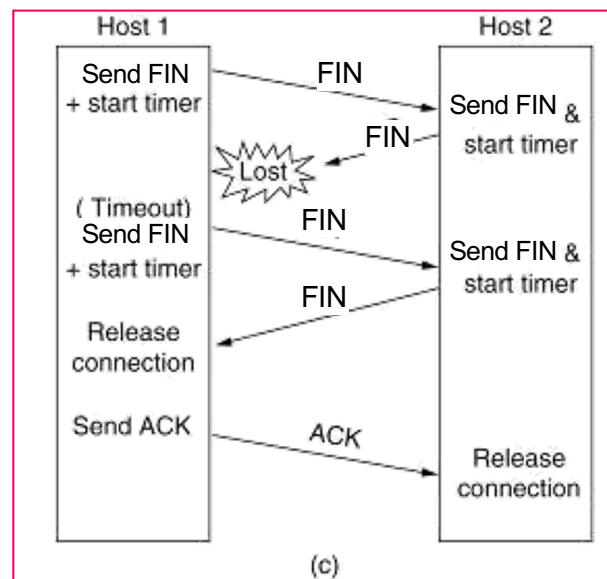
normal



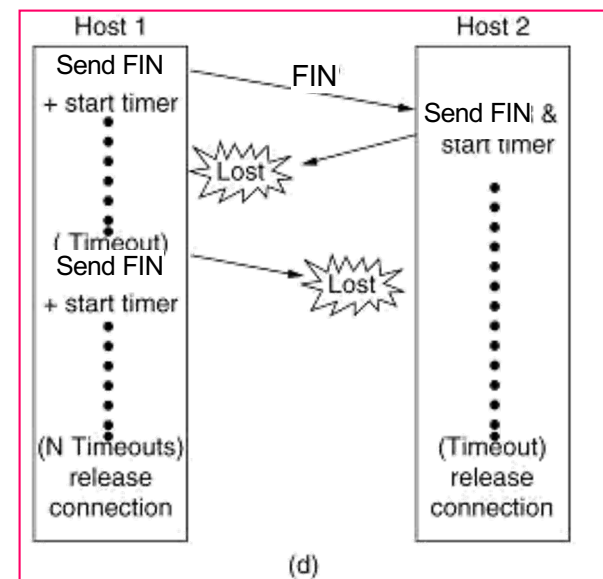
ACK final
pierdut



raspuns
pierdut

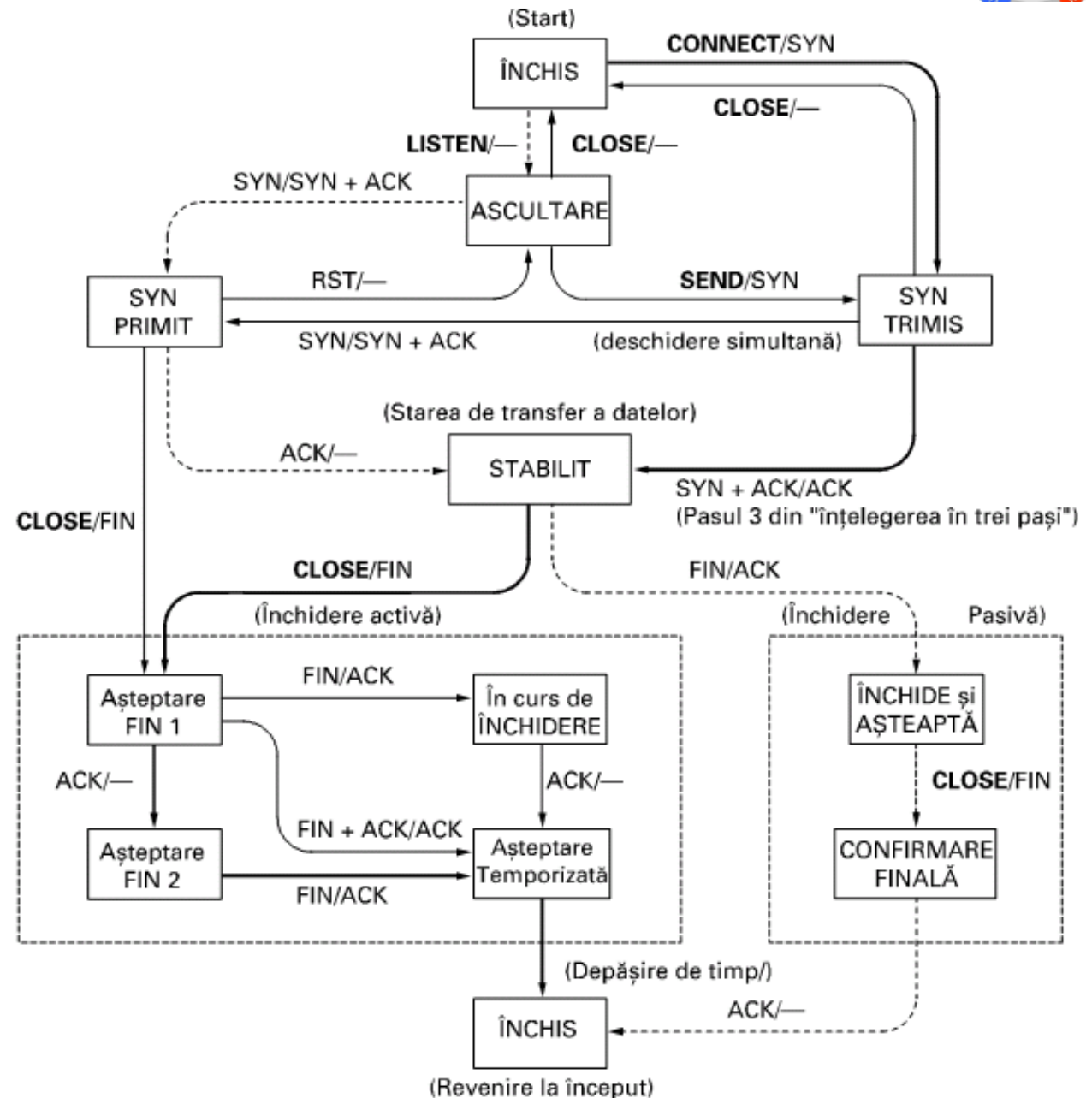


Raspuns si
FIN pierdute



Gestiunea Conexiunilor TCP (mașina de stări)

- client (**activ**) cere conectarea
- server (**pasiv**) raspunde
- Linie groasa continua = cale normala pentru client
- Linie intrerupta = cale normala server
- Linii subtiri = evenimente exceptionale
- Fiecare tranzitie are eticheta **eveniment / actiune**

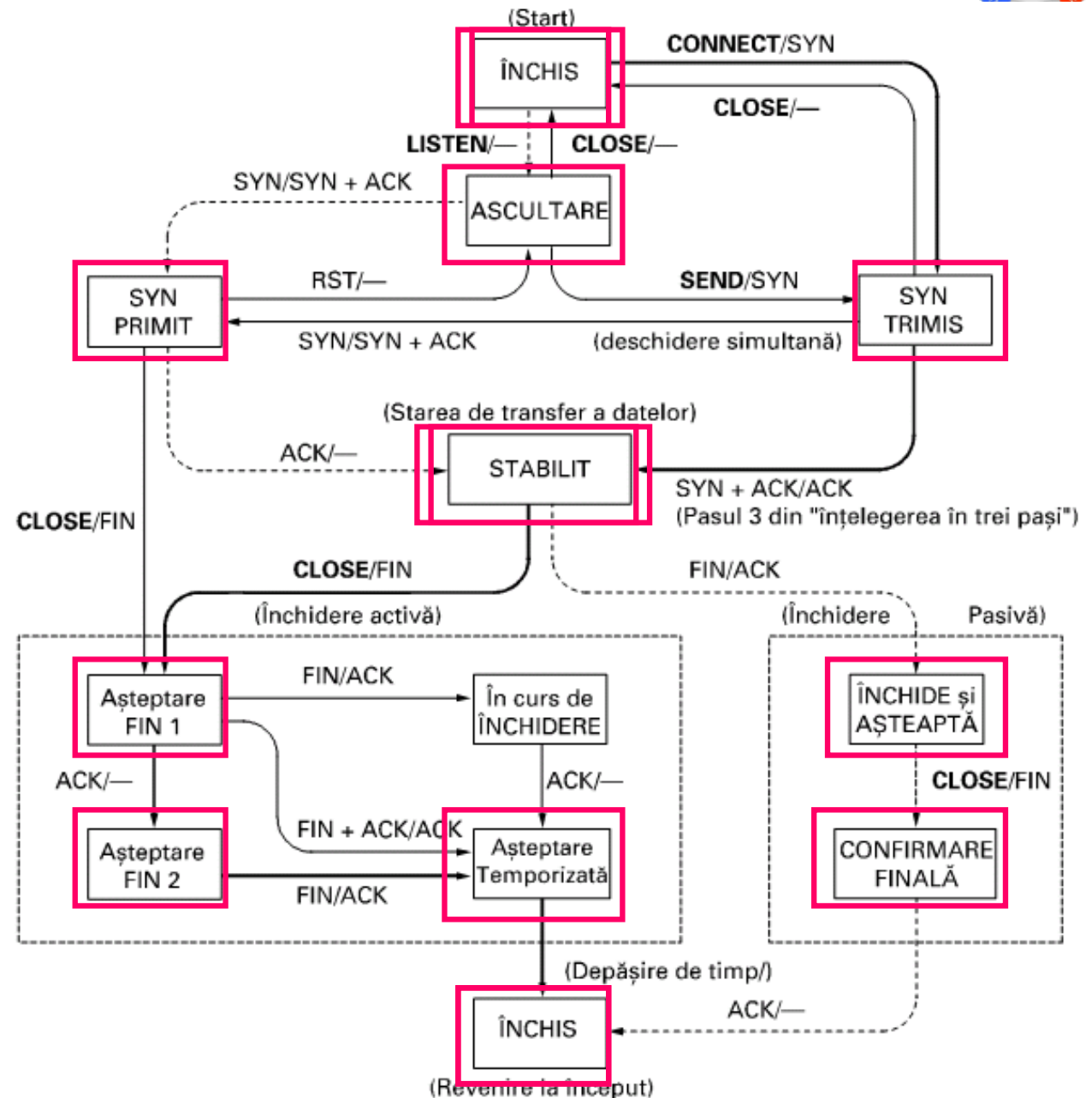


Client

Linie groasa
continua

Server

Linie intrerupta





Corectitudinea segmentelor TCP

1 - **Suma de control (checksum)** din antet are **16 biti** si include

1. antet
2. incarcatura segment TCP (date)
3. un pseudo-antet

adresele IP sursa si destinatie

protocolul (6 pentru TCP)

lungime segment TCP (include antetul)

Algorithm: la **transmisie**

aduna cuvinte de 16 biti in complement fata de 1

complementeaza rezultatul

scrie rezultatul in antet

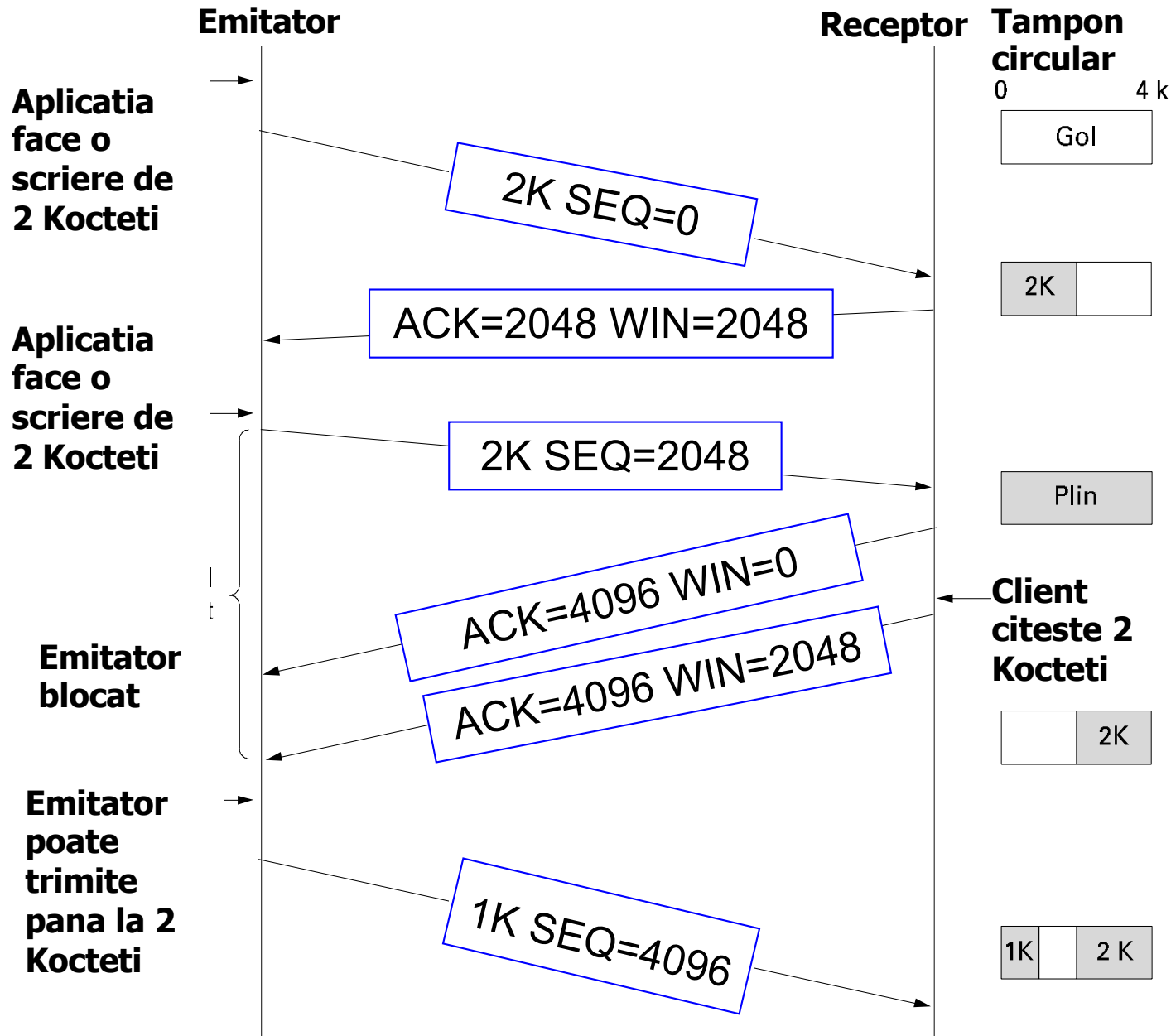
la **receptie**

aduna cuvinte de 16 biti – rezultatul trebuie sa fie zero

2 - **Acknowledgement** number

Corectia se face prin retransmisie

Controlul fluxului de date



Receptorul
specifica
fereastra de
receptie
disponibila

Un emitator blocat
poate trimite:

- date urgente
- un segment de 1 octet ptr a afla fereastra (daca anuntul precedent al receptorului s-a pierdut)



Probleme dimensiuni câmpuri antet

Numere secvență de 32-biti

- Durata ciclu de numarare – depinde de viteza transmisie
 - 1 saptamana pentru 56kbps
 - 57 min pentru 10Mbps
 - 34 sec pentru 1Gbps (**sub 120 sec care este timp viata maxim in Internet**)

Problema: pot apare segmente diferite, cu acelasi numar de secventa?

Soluție

- Folosire opțiuni TCP (RFC 1323)
 - TCP Timestamps
 - asociaza o amprenta de timp fiecarui segment
 - rezolva numere de secventa duplicate



Probleme dimensiuni câmpuri antet (2)

Fereastra receptor (camp **Window size** de 16 biti – echivalent 64 KB)

Transmitere **500 Kb** pe legatura **1 Gbps** ocupa 500 μ sec

La intarziere de **20 ms** pe sens confirmarea se primeste dupa **40 ms** =>
ocupare canal pe un ciclu complet este mică - **1.25%**

Ocupare completă in ambele direcții: produs **bandwidth*delay** =
= 1 Gbps * 40 ms = 40 milioane biti

Condiții - fereastra receptor ar trebui sa fie \geq **bandwidth*delay**

- camp **Window size** nu se poate mari

Soluție

- Folosire **opțiuni** TCP (RFC 1323)
 - **Window Scale** – factor de scalare a câmpului **Window size** de până la **2^{14}** ori
 - ➔ ferestre de până la **2^{30}** octeți



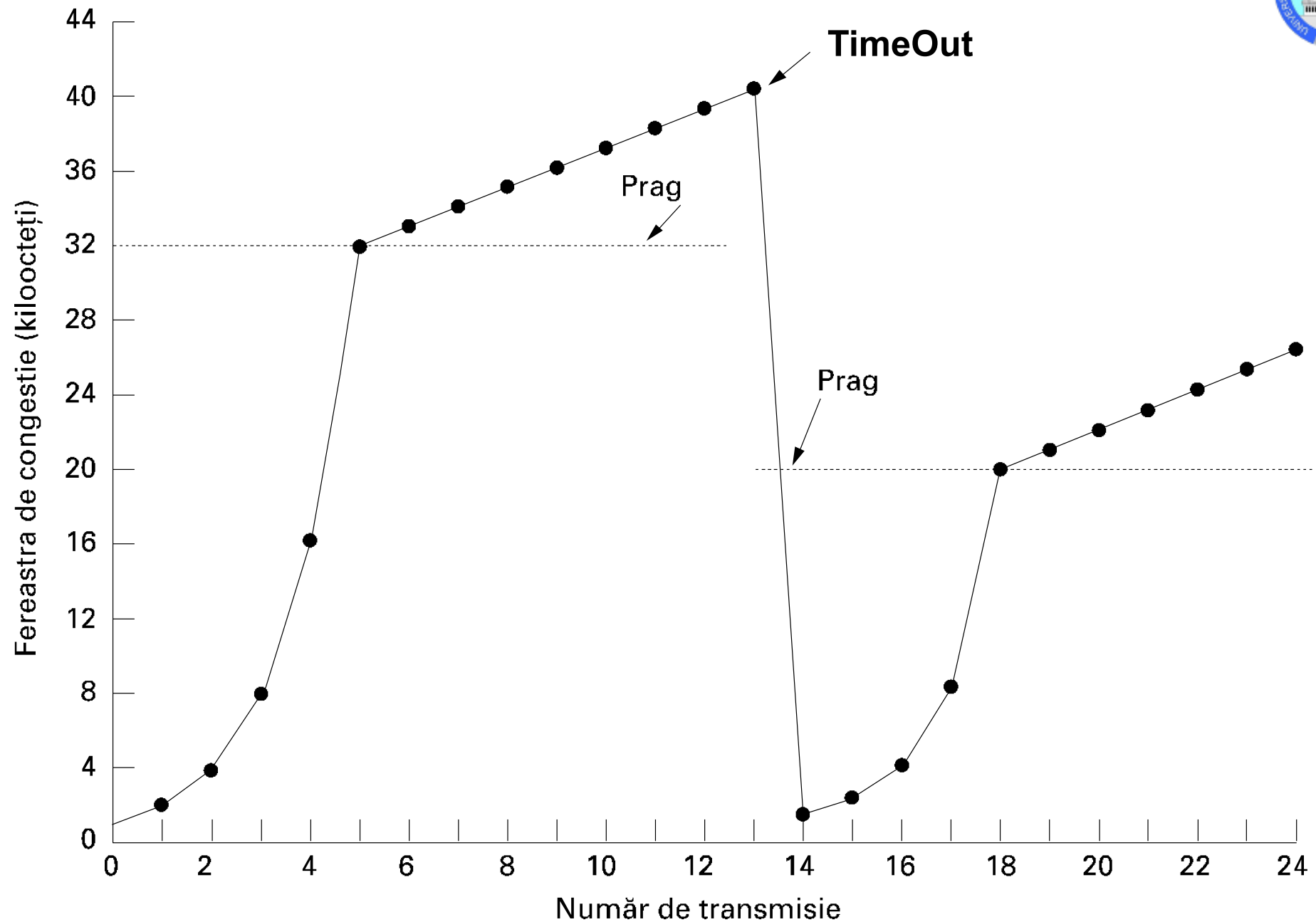
Controlul congestiei

- Fluxul de date transmis pe o conex. TCP limitat de minimul dintre:
 - dimensiunea fereastrei receptorului
 - capacitatea rețelei (**fereastra de congestie**)
- Algoritm de **stabilire fereastra de congestie**
 - transmite un segment de **dimensiune maximă** pe conexiunea stabilită
 - dubleaza volumul de date – rafală de segmente - (creștere exponențială) la fiecare transmisie confirmată la timp
 - la **primul timeout** opreste procedeul si fereastra ramane la valoarea ultimei transmisii confirmate la timp (fara timeout)

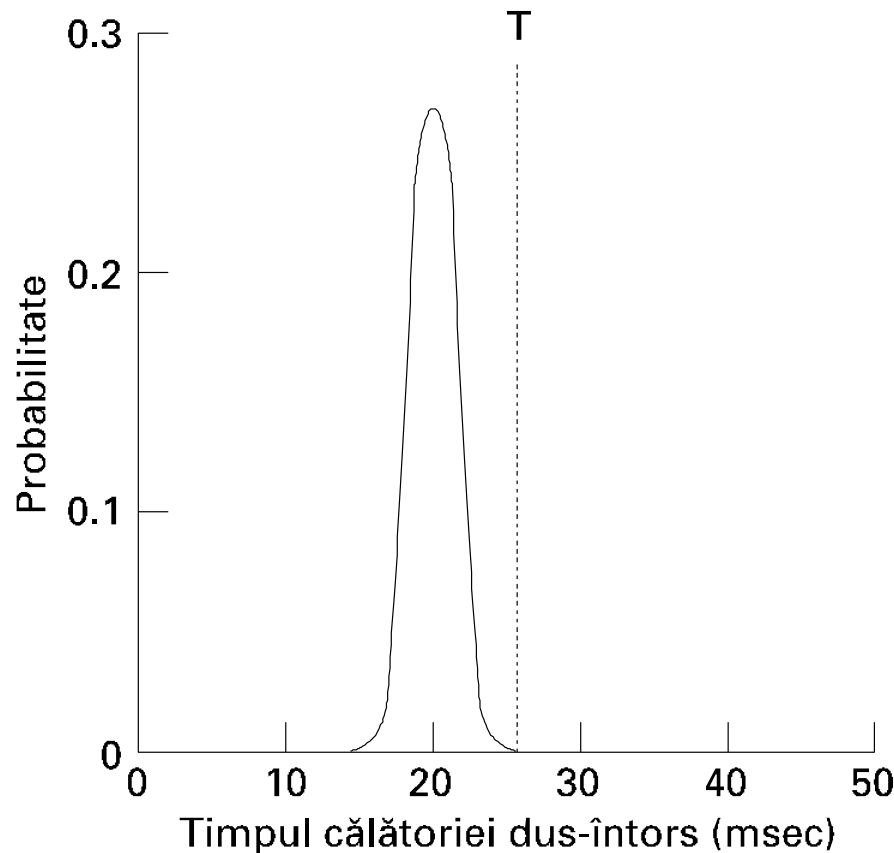


Controlul congestiei

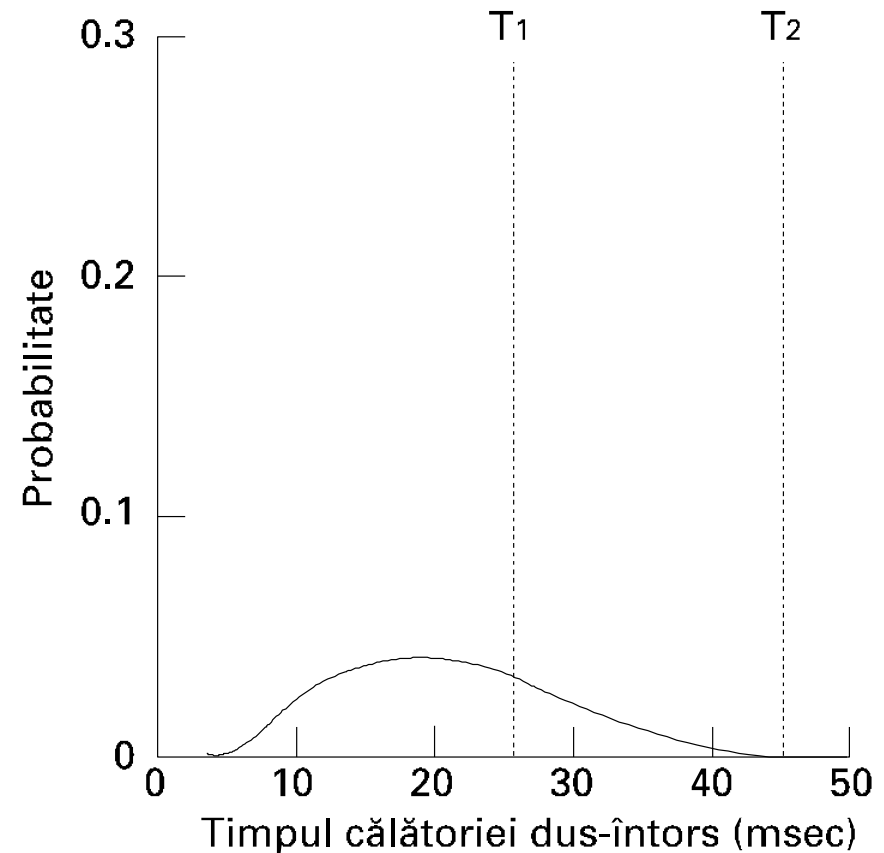
- Algoritmul de **control al congestiei**
 - folosește un **prag** (threshold)
 - la un timeout pragul setat la jumătate din fereastra de congestie
 - se aplica procedeul de creștere (exponentială) a ferestrei de congestie până se atinge pragul
 - peste prag se aplica o creștere liniară (cu câte un segment de dimensiune maximă o dată)



Gestiunea ceasurilor in TCP



(a)



(b)

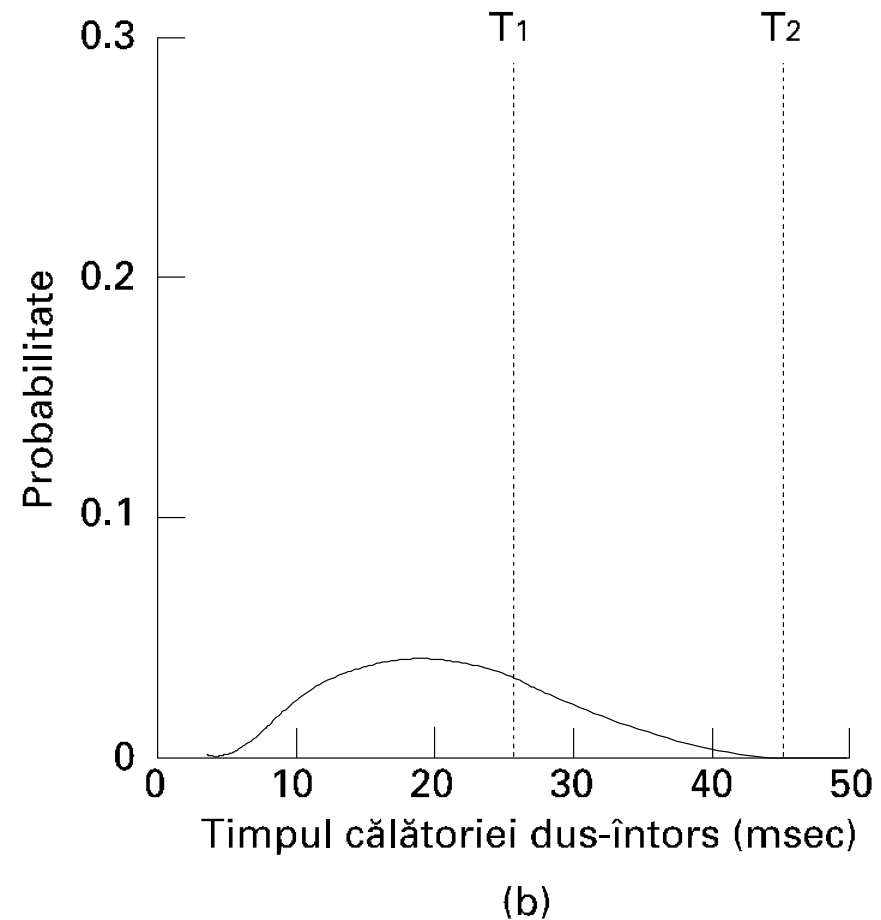
(a) Densitatea de probabilitate a timpilor de sosire ACK in nivelul **legatura de date**.

(b) Densitatea de probabilitate a timpilor de sosire ACK pentru **TCP**.

Gestiunea ceasurilor in TCP (2)

Configurare incorecta – performante slabe:

- Prea lung (T_2) – transmitatorul asteapta mult ptr retransmisie
- Prea scurt (T_1) – trafic inutil generat de transmitator





Stabilire time-out

- **Timeout** diferit la fiecare conexiune - setat dinamic
- Se folosesc metode empirice
- Transmitatorul alege *Retransmission TimeOut* (**RTO**) pe baza *Round Trip Time* (**RTT**)

M este timpul masurat pana la primirea ack

$$\mathbf{RTT} = \alpha * \mathbf{RTT} + (1 - \alpha) * \mathbf{M} \quad \text{cu } \alpha = 7/8$$

$$\mathbf{RTO} = \beta * \mathbf{RTT} \quad \text{cu } \beta = 2$$

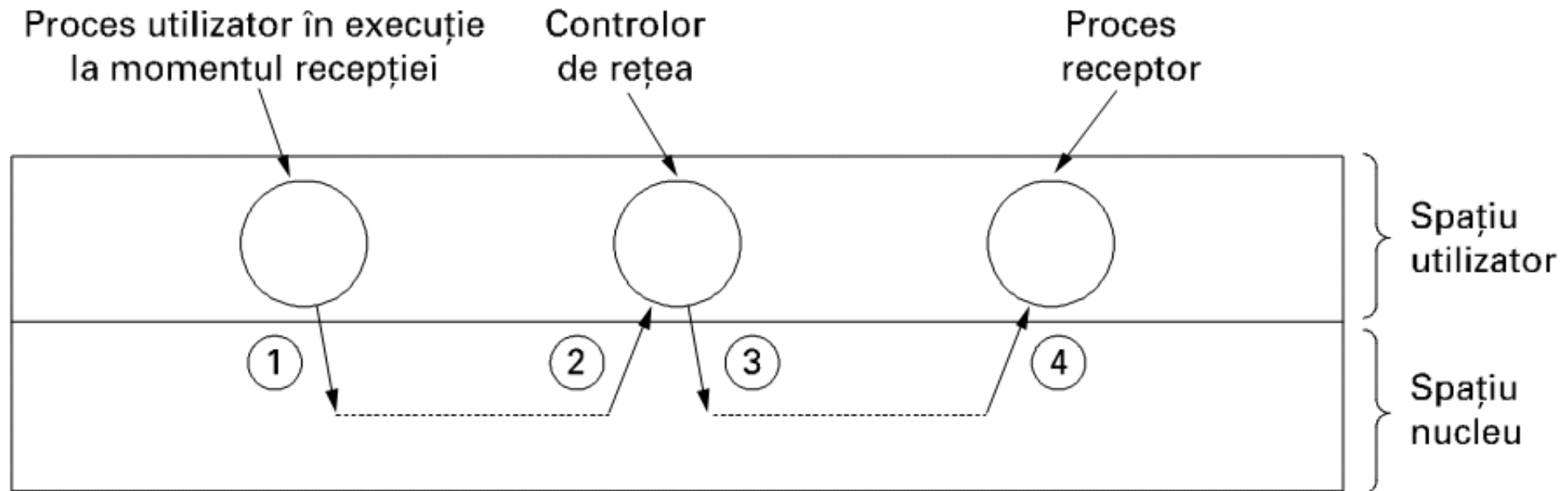
- Alegere dupa *deviatia standard* (*DS*);

D aproximeaza DS

$$\mathbf{D} = \alpha * \mathbf{D} + (1 - \alpha) * |\mathbf{RTT} - \mathbf{M}|$$

$$\mathbf{RTO} = \mathbf{RTT} + 4 * \mathbf{D}$$

Proiectarea pentru performanță



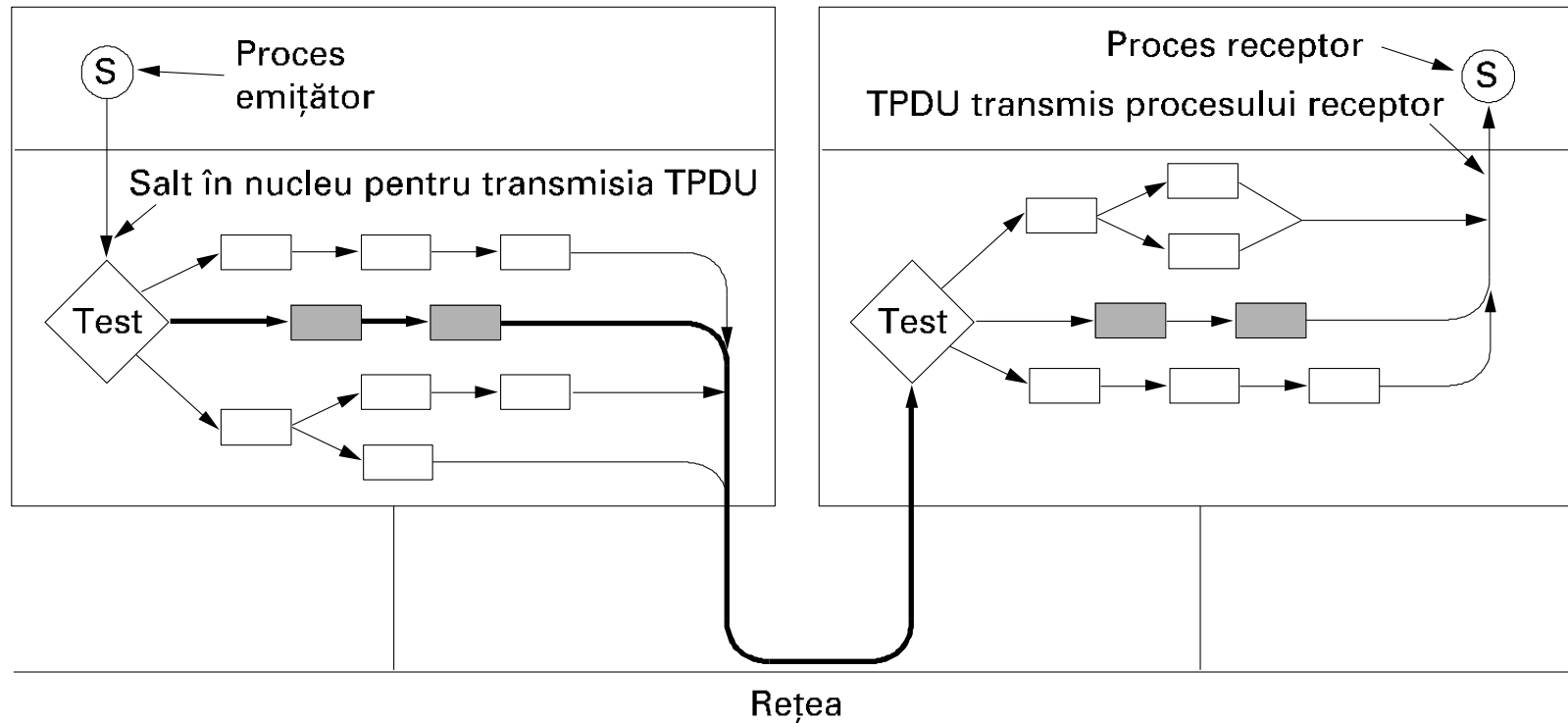
Problema: Multe **schimbari de context user – kernel** pentru a manevra pachete cu un Controlor de rețea în spațiul utilizatorului.

Controlorul assemblează segmentele pentru a le furniza apoi procesului receptor

Soluția: **reducerea** numărului de schimbări de context - acumularea mesajelor sosite, în memorie tampon și livrarea în grup către utilizator.

Similar, gruparea mesajelor de transmis, în memorie tampon și trimiterea lor grupată

Prelucrare rapida TPDU



Calea rapida intre transmitator si receptor este cu line groasa. Pasii sunt reprezentati cu gri.

Test caz normal:

starea = ESTABLISHED

nu se incearca inchiderea conexiunii,

TPDU normal (nu out-of-band),

suficient spatiu la receptor

} → alege fast path

Prelucrare rapida TPDU (2)

| | | | |
|------------------------|--------|------------------|-------------|
| Source port | | Destination port | |
| Sequence number | | | |
| Acknowledgement number | | | |
| Len | Unused | | Window size |
| Checksum | | Urgent pointer | |

(a)

(a) Antet TCP.

| | | | | | | |
|---------------------|-----|----------|--------------|-----------------|--|-----------------|
| VER. | IHL | TOS | Total length | | | |
| Identification | | | | | | Fragment offset |
| TTL | | Protocol | | Header checksum | | |
| Source address | | | | | | |
| Destination address | | | | | | |

(b)

(b) Antet IP.

Transmitator

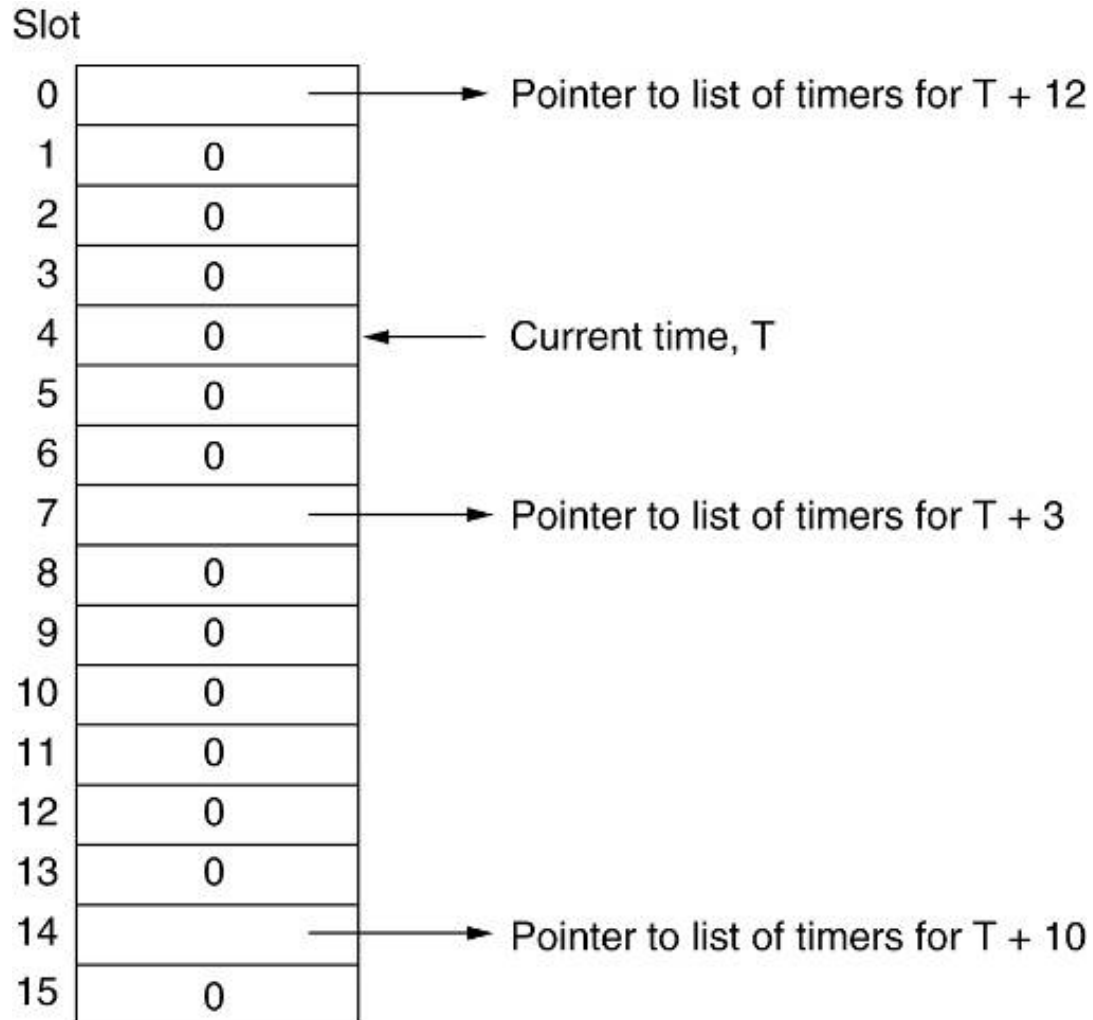
Pastreaza un **prototip** de mesaj in entitatea de transport – campuri nemodificate in unitati de date consecutive; la fel pentru pachet IP

Campurile gri sunt luate din **prototip** fara modificari. Celelalte se calculeaza pentru fiecare segment

Receptor

- Localizeaza inregistrarea conexiunii din TPDU intr-o **tabela** hash
- Testeaza pentru cazul normal (similar cu transmisia)
- Actualizeaza inregistrarea conexiunii (starea curenta)
- Copiaza datele la utilizator si calculeaza suma de control
- Transmite confirmarea

Prelucrare Fast TPDU (timer management)



“Timing wheel.”

1 slot = 1 clock tick

Time current $T=4$

Programare time-out peste 7 tick-uri
→ insereaza eveniment in lista de la slot 11

Anulare -> cauta in lista de la slot corespunzator si elimina eveniment

La fiecare **Clock tick**, un pointer avaneaza cu un slot, circular

Daca slot nevid, proceseaza toate evenimentele



The Real-Time Transport Protocol

Folosit pentru aplicatii multimedia de timp real

- muzica sau video la cerere, videoconferinte etc.

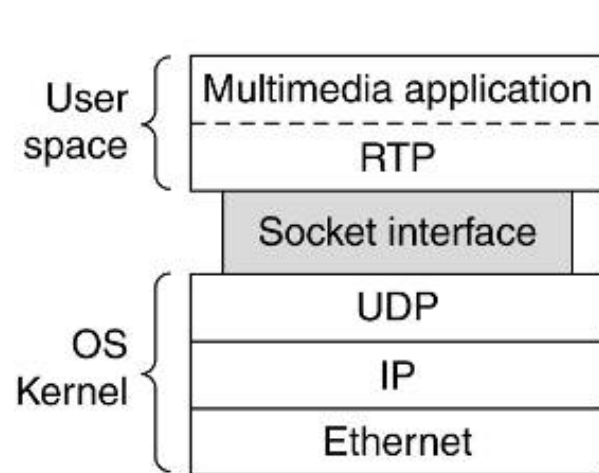
Funcție principală: multiplexare fluxuri RTP și transmiterea lor ca un **singur șir** de pachete **UDP**

La recepție, RTP livrează aplicației datele multimedia

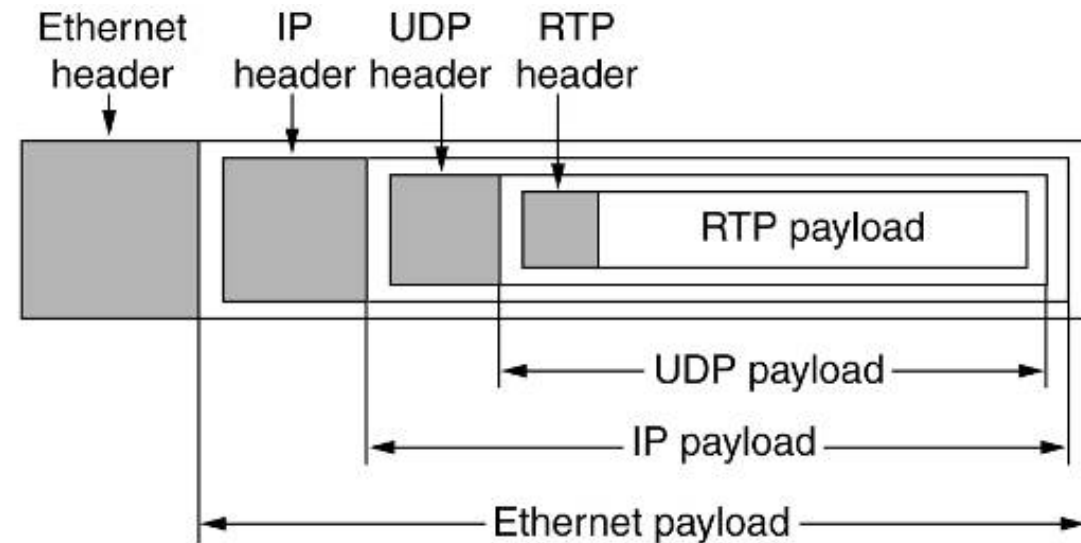
Fără retransmiterea pachetelor

- UDP nu asigură corectitudinea transmisiei
- RTP atribuie un număr de **secvență** fiecărui pachet
- **tratarea** unui pachet pierdut se face **la aplicație**: receptorul poate “interpola” pachetul absent sau ignora eroarea

The Real-Time Transport Protocol (2)



(a)



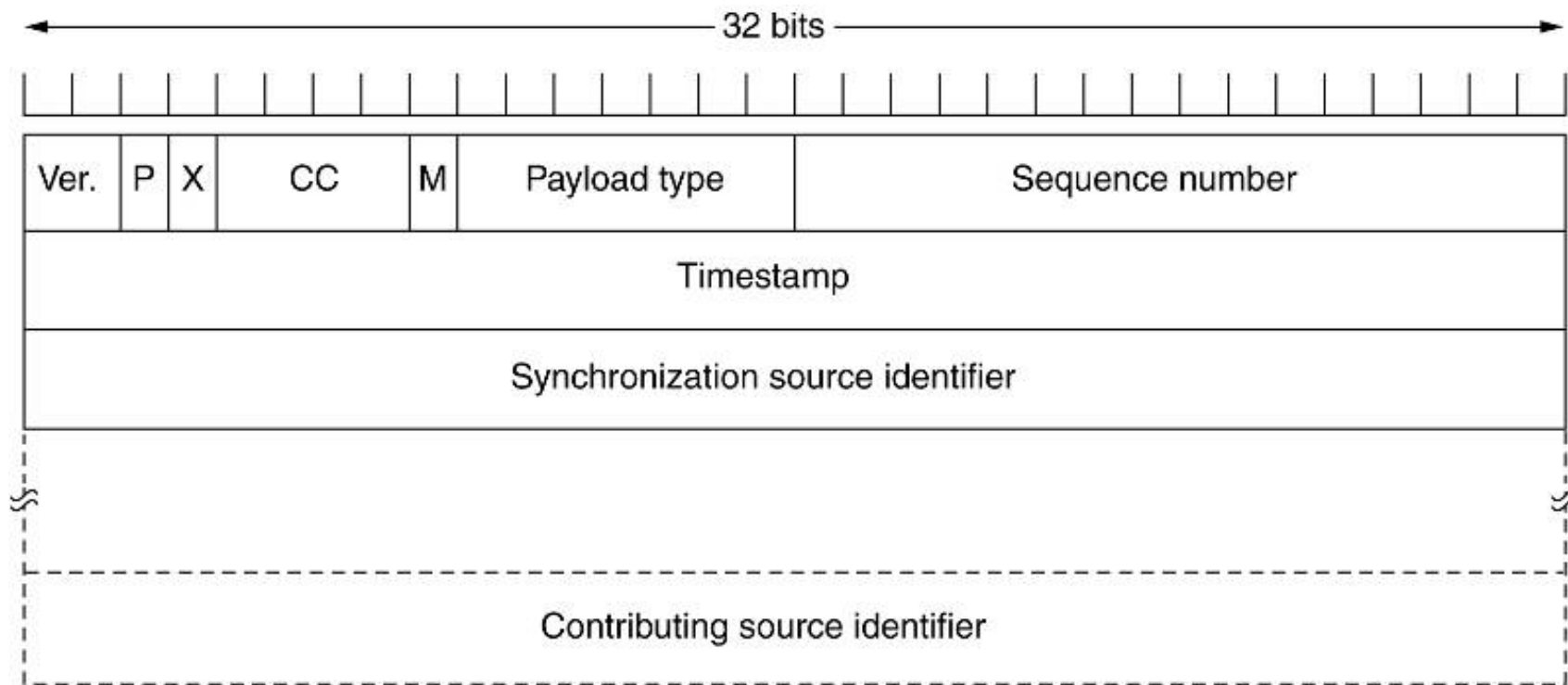
(b)

(a) Pozitia RTP in stiva de protocoale.

(b) Format pachet.

Formatul datelor din "RTP payload" este specific aplicatiilor

RTP permite definirea unor **profile** (ex. single audio stream) si, pentru fiecare profil, a mai multor **formate**



Antet RTP

Timestamp – amprenta de timp relativa la inceputul fluxului

- reduce efectul intarzierilor variabile
- permite sincronizarea intre stream-uri

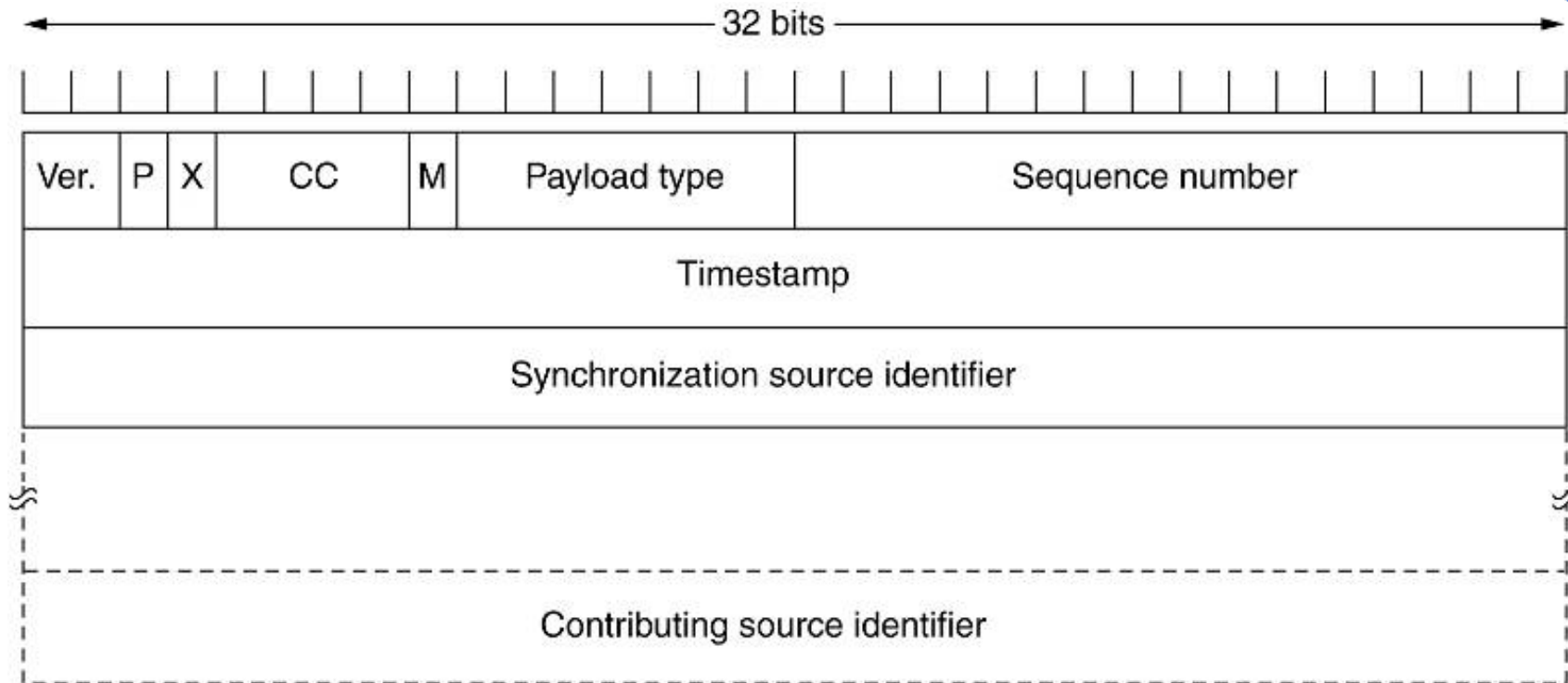
Synchro – identifica sursa de sincronizare (ex. microfon, camera video)

- receptorul grupeaza pachetele dupa sursa, pentru redare

Contrib – lista surselor care au contribuit la continutul curent; folosit pentru mixere;

- ex. la audio-conf, lista vorbitorilor ale caror discursuri sunt in pachet curent

Antet RTP



P – padding - pachet extins la multiplu de 4 octeti

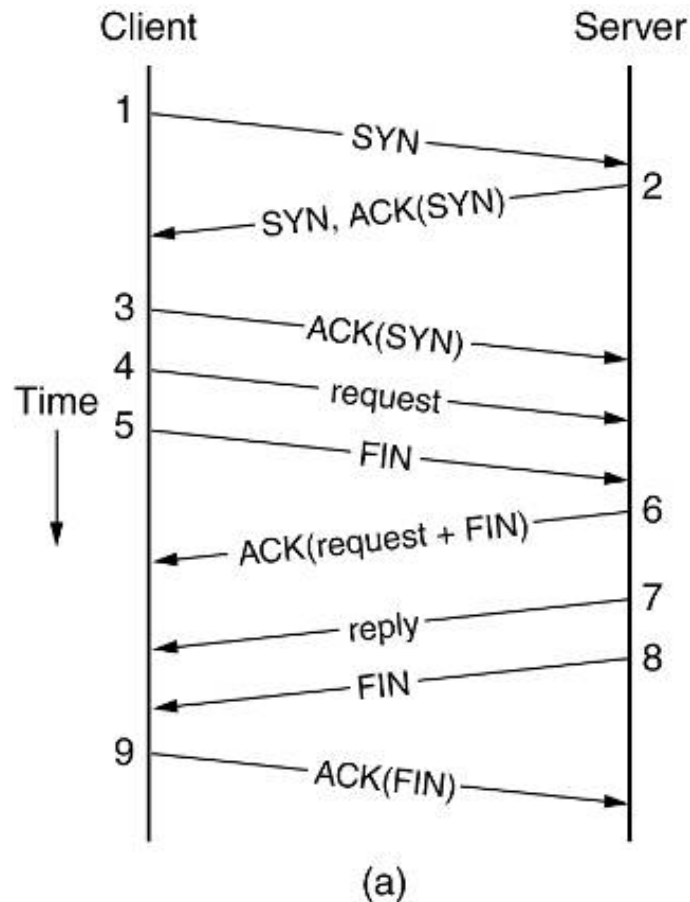
X – extension - antet extins (primul cuvant da lungimea)

CC – no. antete surse contributoare

M – mark (specific aplicatiei. Ex. start **video frame**)

Payload type – e.g. MP3

TCP Tranzactional



(a) RPC – Remote Procedure Call - folosind TPC normal.

(b) RPC folosind T/T

Permite transfer de date odata cu stabilirea conexiunii