

Librăria online *ReadsProfiler*

Sebastian Ciobanu

Facultatea de Informatică,
Universitatea *Alexandru Ioan Cuza*, Strada Gen. Berthelot 16, Iași
`sebastian.ciobanu@info.uaic.ro`
`http://www.info.uaic.ro/`

Abstract. Combinând elemente de rețele de calculatoare (TCP implementat concurrent prin crearea cate unui fir de execuție pentru servirea fiecărui client și prin folosirea socketilor BSD), de baze de date (MySQL ca modalitate de stocare a datelor despre cărțile disponibile într-o librărie) și de algoritmică/statistică (algoritmi de recomandare de tip filtrare colaborativă bazată pe utilizator sau filtrare bazată pe conținut, similaritatea cosinus), aplicația *ReadsProfiler* de tip client/server se ocupă de accesul la o librărie online, oferind recomandări de cărți utilizatorilor.

Keywords: TCP, MySQL, cosine similarity, filtering, e-bookshop

1 Introducere

În ziua de astăzi, aproape orice magazin online are o modalitate sau mai multe de a își convinge proprii clienți să cumpere noi produse. Una din aceste metode este reprezentată de sistemul de recomandări.

Scopul unui astfel de sistem este de a genera recomandări semnificative unei colecții de utilizatori pentru obiecte sau produse care ar putea să îi intereseze. Să considerăm doar un exemplu minimal al unei librării online: utilizatorul Negrescu Vasile a căutat cărțile: *Ion*, de Liviu Rebreanu; *Răscoala*, de Liviu Rebreanu; *Pădurea spânzuraților*, de Liviu Rebreanu. Se poate intui că lui Vasile îi plac cărțile lui Liviu Rebreanu, așa că sistemul poate urma această regulă și îi va putea recomanda cartea: *Ciuleandra*, de același autor, carte de care Vasile, probabil, este interesat. Exemple din lumea reală ar include: sugestiile pentru cărți de pe Amazon, sugestiile de filme pe netflix etc.

Aplicația dezvoltată este axată pe 2 termeni cheie: librărie online și sistem de recomandări. Aceasta se numește *ReadsProfiler* și este o aplicație client-server prin care se permite accesul la o librărie online. Oferind o experiență tradițională prin **clasica interfață** cu utilizatorul, aplicația este **simplică de utilizat**, conține doar **cărți scrise sau traduse în limba română**, iar *tehnicile* utilizate pentru sistemul de recomandare sunt ideale pentru date relativ mici (față de numărul de produse de pe Amazon, spre exemplu) așa că **implementarea** poate fi de folos unui novice în crearea propriului sistem de recomandări. Un alt avantaj este faptul că programul este **gratis**, open-source.

2 Tehnologiile utilizate

În cadrul aplicațiilor în rețea **paradigmele** folosite cel mai frecvent sunt: modelul client/server, RPC și P2P. În cadrul primei paradigme, un server oferă servicii clienților. Serverul poate satisface cererile provenite de la clienți în mod **iterativ sau concurent**. De obicei, serverele concurente folosesc protocoale de comunicație orientate conexiune. În acest caz, cererile sunt procesate concurent.

În *ReadsProfiler*, s-au folosit:

- **paradigma: modelul client/server**
- **server: concurent**
- **protocol (la nivelul transport): TCP**

Serverul a fost ales să fie **concurent** din considerente practice: când se dorește a cumpăra un produs și clientul este pus să aștepte, rezultatul este dat de clientul care renunță la ideea de a mai cumpăra ceva.

TCP a fost ales pentru că este un protocol de transport orientat conexiune, fără pierdere de informații: dacă un client caută o carte și primește un rezultat neavând legătura cu acea căutare, acesta își pierde încrederea în magazinul online și astfel magazinul va avea de suferit; mai grav este momentul în care se descarcă o carte (pentru care poate s-a și plătit) și unele părți (sau mai multe) nu sunt ceea ce trebuie (citibile); deși UDP oferă o mai mare viteză, preferăm să aducem rezultate corecte într-un timp mai *lung* decât cel dat de folosirea UDP-ului (care în practică de multe ori nici nu este observabil); în astfel de situații (magazin online), doar TCP-ul poate fi adus în discuție. Ca **interfață de programare a aplicației**, am ales **socketii BSD**.

Am ales ca datele despre cărți (și nu numai) să fie stocate într-o bază de date **MySQL**. Față de SQLite, MySQL oferă **acces în rețea la baza de date și un grad mare de concurență** (parallelism) a interogărilor, lucruri chiar necesare unei aplicații în rețea precum este cea de față. Totodată, căutările frecvente ale utilizatorilor după diverse criterii (combinat) se pliază perfect în acest context al bazelor de date. Iarăși, **volumul de date** este un alt motiv pentru care am ales această opțiune.

3 Arhitectura aplicației

3.1 Adăugiri la interfața BSD

Pentru că în cadrul stivei de protocoale TCP/IP fiecare nivel oferă servicii nivelului superior, avem că nivelul transport (unde am folosit protocolul TCP) oferă servicii nivelului aplicație. Astfel, ca programatori, la nivelul aplicației vom ști că biții vor ajunge în cele din urmă la destinatar în ordine și în mod eficient (să nu mai vorbim de celelalte *încercări* pe care le-am fi avut dacă nivelurile inferioare nivelului transport n-ar fi existat).

Totuși, la nivel de aplicație, prin utilizarea interfeței socketilor BSD, nu putem fi siguri că un apel *read* sau *write* va duce sigur la citirea/scrierea numărului de bytes specificat. Iarăși, o altă problemă este *Dacă am format o conexiune între*

client și server și unul dintre cei doi o părăsește nepoliticos, cum poate cealaltă parte să-și dea seama că nu mai are cu cine comunica? Evident, problema se pune doar dacă cele două părți mai au ce comunica (*read*, *write*). Dacă nu, nu există nicio problemă. Problema se rezolvă ușor dacă partea rămasă se așteaptă la un mesaj de la cealaltă parte, pentru că apelul *read* își dă seama de acest lucru. Dacă însă partea rămasă transmite un mesaj, el va crede că mesajul a ajuns la destinatar cu succes (pentru că abia al doilea apel *write* își dă seama de *ruperea* conexiunii, nu și primul).

Aceste două probleme au fost rezolvate la nivelul aplicației prin implementarea funcțiilor *myRead*, *myWrite* și *isPeerClosed*. Mai multe detalii vom avea în secțiunea următoare.

3.2 Subniveluri ale nivelului aplicație

Luând în calcul complexitatea implementării unei astfel de aplicații, am considerat că ar fi maximă nevoie de metode precum *mySendLength*, *mySendMessage*, *mySendFile*, *mySendPage*, *sendTable*, *myRecvLength*, *myRecvMessage*, *myRecvFile*, *myRecvPage*, *recvTable* care să faciliteze comunicarea între client și server. Mai mult, în final *utilizatorul* funcției *main()* nu va avea acces decât la funcții de tipul *login*, *find*, atât la nivelul clientului, cât și la nivelul serverului (programatorul de la acest subnivel nu va fi nevoit să cunoască detalii de implementare ale acestor funcții).

3.3 Proiectarea unui protocol de comunicație între client și server cu scopul de a ușura comunicarea dintre cei doi

a) Definirea operațiilor permise clientului (comanda, informații de furnizat, semantica)

La nivelul cel mai de sus (al clientului utilizator de terminal), am definit următoarele operații:

1)

1. comanda: **signup**

informații de furnizat: *user*, *pass*

semantica: va înregistra un nou utilizator cu numele *user* și parola *pass*.

Dacă *user* este înregistrat deja, se va afișa un mesaj corespunzător.

```
@ReadsProfiler> signup
user: sebi
pass: parola
```

Fig. 1. Comanda **signup**

2. comanda: **login**

informații de furnizat: *user*, *pass*

semantica: va executa autentificarea utilizatorului. Dacă *user* nu este găsit în baza de date sau dacă parola nu corespunde, se va afișa un mesaj corespunzător. **Atentie!** Fără logare, majoritatea comenzilor sunt inaccesibile (doar **signup** e accesibilă).

```
@ReadsProfiler> login
user: sebi
pass: parola
```

Fig. 2. Comanda **login**

3. comanda: **logout**

informații de furnizat: -

semantica: va executa delogarea utilizatorului curent. În cazul în care nu este nimeni logat, va afișa un mesaj corespunzător.

```
sebi@ReadsProfiler> logout
Succes.
@ReadsProfiler>
```

Fig. 3. Comanda **logout**

4. comanda: **quit**

informații de furnizat: -

semantica: va închide programul în 3 secunde.

```
@ReadsProfiler> quit
Programul se va închide în 3 secunde
```

Fig. 4. Comanda **quit**

5. comanda: **find**

informații de furnizat: criteriile de căutare, valorile criteriilor de căutare

semantica: va duce la afișarea tuturor cărților ce îndeplinesc toate criteriile menționate. Dacă nu există astfel de cărți, se va afișa un mesaj corespunzător.

Ca ordine a operațiilor, această comandă are sens doar dacă sunteți logat și nu ați accesat chiar înainte vreo pagină a vreunei cărți.

```
sebi@ReadsProfiler> find
Cautati dupa criteriile de cautare cunoscute. Doriti sa le afisam?[d/n]
d
Criteriile sunt:
0: <carteID>
1: <parte_din_titlu>
2: <nr_volum>
3: <nume_autor>
4: <prenume_autor>
5: <isbn>
6: <editura>
7: <anul_aparitiei>
8: <gen>
9: <subgen>
10: < rating<= >
11: < rating= >
12: < rating>= >
Scrieti nr. criteriilor dorite de cautare pe O SINGURA LINIE despartite prin spatii, apoi introduceti informatiile cerute.
1 2
Titlu: anna
Nr. volum: 1
```

Fig. 5. Comanda **find**

6. comanda: **list**
 informații de furnizat: -
 semantica: va afișa maxim 100 de cărțile ce sunt disponibile în cadrul librăriei.
 Ca ordine a operațiilor, această comandă are sens doar dacă sunteți logat și nu ați accesat chiar înainte vreo pagină a vreunei cărți.

```
sebi@ReadsProfiler> list
+-----+-----+-----+-----+-----+-----+
| # | Titlul cartii | # volum | Un autor | Editura | Un gen | Un subge |
+-----+-----+-----+-----+-----+-----+
| 1 | Notre-Dame de Paris | - | Hugo | Adevarul | Fictiune | Drama |
| 2 | Marile Sperante | 1 | Dickens | Adevarul | Fictiune | Drama |
| 3 | Marile Sperante | 2 | Dickens | Adevarul | Fictiune | Drama |
```

Fig. 6. Comanda **list**

7. comanda: **access**
 informații de furnizat: numărul cărții din ultimul afișaj (list/find) cu succes
 semantica: va accesa pagina aferentă cărții cu numărul dat (nr. a apărut în ultimul apărut în ultimul afișaj list/find). Dacă nr. nu este potrivit, se va afișun mesaj corespunzător. Ca ordine a operațiilor, această comandă are sens doar dacă sunteți logat, nu ați accesat chiar înainte vreo pagină a vreunei cărți și ați folosit măcar o dată list/find.
8. comanda: **download**
 informații de furnizat: numele fișierului salvat

```
sebi@ReadsProfiler> access
Introduceti numarul de carte din afisajul anterior.
4
Ai selectat cartea cu numarul: 4
Rating general: 1.0000
Titlu: Un yankeu la curtea regelui Arthur
# volum: -
Autori: Twain Mark
ISBN: 9783161484103
Editura: Adevarul
Anul aparitiei: 2008
Gen(uri): Fictiune
Subgen(uri): Istorie
Puteti descarca/da un rating cartii (daca nu ati facut-o deja).
sebi@ReadsProfiler> █
```

Fig. 7. Comanda **access**

semantica: va avea sens doar dacă ultima comandă folosită este **access**. Cartea accesată anterior va fi descărcată de utilizator la calea ../descarcari și va avea denumirea menționată.

```
sebi@ReadsProfiler> download
Introduceti numele fisierului ce va fi descarcat in "../descarcari". Daca numele exista deja, vechiul fisier va fi
yankeu
Succes.
sebi@ReadsProfiler> █
```

Fig. 8. Comanda **download**9. comanda: **rate**

informații de furnizat: ratingul dat cărții accesate ($0 \leq rating \leq 5$)

semantica: va avea sens doar dacă ultima comandă folosită este **access**. Cărții accesate anterior i se va atribui ratingul menționat doar dacă userul nu i-a mai dat vreun rating în trecut, caz în care se va afișa un mesaj corespunzător.

```
sebi@ReadsProfiler> rate
Dati o nota cartii de la 0 (de proasta calitate) la 5 (de buna calitate).
4
Ati dat deja un rating cartii: 1. Nu mai puteti vota inca o data.
sebi@ReadsProfiler> █
```

Fig. 9. Comanda **rate**10. comanda: **back**

informații de furnizat: -

semantica: va avea sens doar dacă ultima comandă folosită este **access**. Ne va întoarce la rezultatul ultimei comenzi **find/list**.

```
sebi@ReadsProfiler> back
```

#	Titlul cartii	# volum	Un autor	Editura	Un gen
1	Notre-Dame de Paris	-	Hugo	Adevarul	Fictiun
2	Marile Sperante	1	Dickens	Adevarul	Fictiun

 Fig. 10. Comanda **back**

 11. comanda: **recommend**

informații de furnizat: -

semantica: va da 1 sau mai multe recomandări de cărți userului curent. Ca ordine a operațiilor, această comandă are sens doar dacă sunteți logat, nu ați accesat chiar înainte vreo pagină a vreunei cărți și ați folosit măcar o dată list/find.

```
sebi@ReadsProfiler> recommend
Va recomandam cartile numerotate cu 2 si 8. Folositi comanda "find" cu optiunea 0 pentru a a
sebi@ReadsProfiler> █
```

 Fig. 11. Comanda **recommend**

 12. comanda: **help**

informații de furnizat: -

semantica: va da informații despre comenzi și reguli de aplicare a lor (inf. asemănătoare cu cele de aici).

b) Raportarea erorilor

 Pe lângă erorile interne aferente fiecărui participant la conexiune, mai apar erori *interesante* și pentru partener. Eroarea cea mai distructivă dacă o putem denumi așa este cea menționată în subsecțiunea 3.1. Celelalte erori ce vor fi comunicate pe rețea vor fi:

- signup: User deja existent.
- login: Ați greșit userul/parola.
- find: Nu există astfel de cărți.
- rate: Ați dat deja un rating cărții. Nu mai puteți vota încă o dată.

Celelalte erori (cum ar fi neîncadrarea ratingului între 0 și 5) vor fi tratate doar local, neimplicând și partenerul. Pentru fiecare eroare vom avea perechea (cod de eroare, mesaj explicativ).

3.4 Stocarea informațiilor

În cadrul bazei de date MySQL, este folosită o serie de tabele pentru stocarea datelor. Vom arunca doar o privire de ansamblu asupra lor:

- opere(opID,titlu,nr_volum)

- carti(isbn,editura,calea,an_aparitie,rating,opID)
- autori(autID,nume,prenume)
- opere_autori(opID,autID)
- ierarhie(subgen,gen)
- abordari(opID,subgen)
- history(carteID,user,accessed,downloaded,rating)
- recommendedTo(carteID,user)
- view-ul avgRating(carteID,avgrate)

3.5 Algoritmul de recomandare

Conform cerinței, utilizatorilor li se oferă recomandări de cărți. Vom folosi 2 metode de recomandare:

- **filtrare bazată pe conținut**; astfel vom folosi faptul că avem cărțile sortate pe genuri și că autorii abordează mai multe genuri în operele lor

- **filtrare colaborativă bazată pe utilizatori**; astfel, sistemul de recomandări se va putea îmbunătăți în acuratețe, luând în considerare factori precum gusturile unui utilizator, preferințele altor utilizatori cu gusturi similare, rating-ul pe care l-au primit cărțile de la clienții cărora le-au fost recomandate.

Sisteme de recomandare bazate pe conținut se recomandă articole al căror conținut e similar cu cel al altor itemi agreeți de user în trecut sau care se potrivește cu atributele userului.

Sisteme cu filtrare colaborativă - unui utilizator i se recomandă itemi pe baza voturilor tuturor userilor din trecut.

Mai ales pentru filtrarea cea din urma va trebui să calculăm similaritatea dintre 2 utilizatori. Există numeroase idei de calcul (distanța euclidiană, Manhattan etc.), dar pentru datele noastre, în care nu fiecare user a dat rating la aproape toate cărțile (cel mai probabil, niciunul nu a făcut aceasta), vom folosi similaritatea cosinus care se descurcă bine și în caz de sistem bazat pe utilizator (cazul nostru). Pentru a înțelege mai bine, trebuie să ne imaginăm fiecare utilizator ca pe un vector n -dimensional, unde n =numărul de cărți din baza de date. Fiecare componentă a vectorului reprezintă ratingul dat de utilizator cărții *de pe poziția respectivă*.

$$\text{similarity} = \cos(\theta) = \frac{\mathbf{A} \cdot \mathbf{B}}{\|\mathbf{A}\| \|\mathbf{B}\|} = \frac{\sum_{i=1}^n A_i B_i}{\sqrt{\sum_{i=1}^n A_i^2} \sqrt{\sum_{i=1}^n B_i^2}}$$

Fig. 12. Similaritatea cosinus [9]

Pentru a recomanda o carte lui X , calculăm similaritățile dintre X și ceilalți. Îl alegem pe cel ce are similaritatea dintre el și X cea mai mare. Fie acesta Y .

Căutam o carte C căreia Y i-a dat rating *pozitiv* și căreia X nu i-a dat rating. C va fi cartea de recomandat lui X . Procedeu anterior se va perfecționa prin utilizarea primilor k (k fiind știut în prealabil) utilizatori similari cu X .

În cazul nostru, am folosit un algoritm de recomandare bazat pe conținut și unul bazat pe useri. Dacă prin astfel nu găsim nicio recomandare luăm **topul** cărților după rating și recomandăm prima carte căreia userul curent nu i-a dat rate.

Algoritm bazat pe conținut sună astfel:

1. U =user curent
2. extragem cărțile din istoricul lui U cu ratingul ≥ 3
3. extragem autorii acelor cărți
4. extragem subgenurile abordate de acei autori
5. extragem cărțile care abordează acele subgenuri și cărora U nu le-a dat nici rate, nici download și care au fost recomandate userilor
6. calculăm pentru fiecare carte media dată de notele acordate de userii cărora le-au fost recomandate
7. ordonăm cărțile după această medie
8. *returnăm* prima carte.

Algoritm bazat pe useri sună astfel:

1. U =user curent= (r_1, \dots, r_n) , n =nr. de cărți, r_i =ratingul dat cărții i de către U
2. formăm și ceilalți useri astfel
3. calculăm similaritatea dintre U și fiecare dintre ceilalți useri prin similaritatea cosinus
4. ordonăm userii după această similitudine
5. alegem primul user = V
6. căutăm cărți ale lui V neratingate de U
7. *returnăm* prima găsită.

4 Detalii de implementare

În cadrul protocolului de comunicare client-server, la nivel de implementare, putem adauga:

- Formatul mesajelor:
 - de obicei, mesajele vor fi de tipul lungimea_mesajului + mesaj
 - în cazul diferențierii erorilor de mesaje de succes (acolo unde trebuie trimise erori perechii de comunicare), lungimea va fi și ea prefixată cu '0'-mesaj de succes, '1'-mesaj de eroare
 - un alt tip de mesaj este cel în care am trimis 0/1/2 recomandări clientului și am prefixat astfel întreg mesajul cu '0'=0 rec., '1'=1 rec. 2=2 rec.
- la *download*: vom deschide fișierul în format binar ("rb", "wb") și apoi îl vom trimite; dacă îl deschideam în format normal ("r", "w") era posibil să găsim un EOF pe la mijlocul fișierului și transferul ar fi fost deci doar pe jumătate;

Diagrama aplicației - nivel conceptual

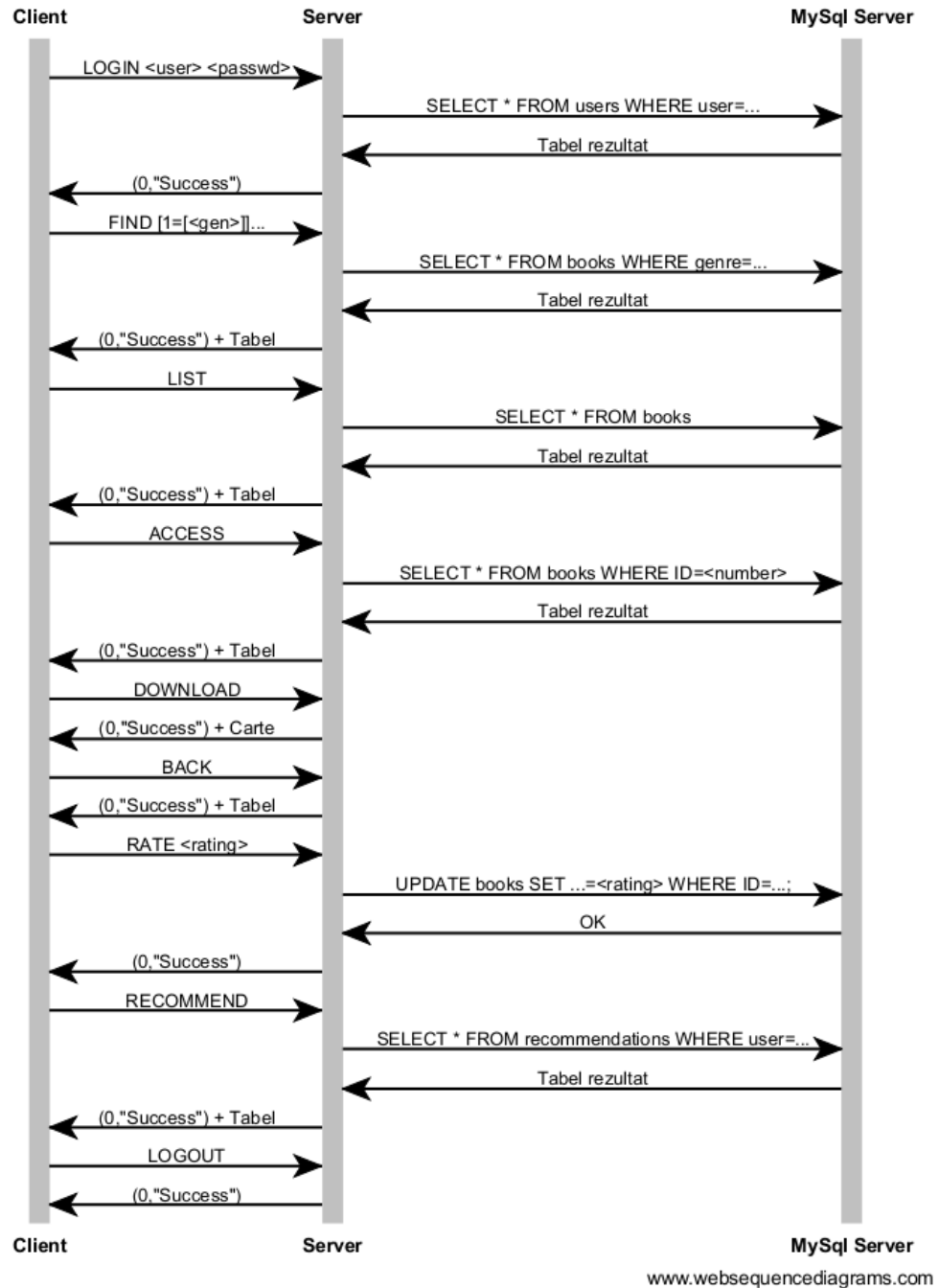


Fig. 13. Diagrama aplicației - nivel conceptual

- trimiterea de mesaje face apel la următoarele primitive care asigură citirea/scrierea numărului de bytes specificați:

Primitive citire/scriere implementate de utilizator:

```
ssize_t myWrite(int fd, const void *buf, size_t count)
{
    size_t left=count;
    ssize_t written=0;
    const char *charBuffer=buf;
    while(left>0)
    {
        if(isPeerClosed(fd)==1)
        {
            return -2;
        }
        written=write(fd,charBuffer,left);
        if(written==-1)
        {
            return -1;
        }
        left-=written;
        charBuffer+=written;
    }
    return count;
}
```

```
ssize_t myRead(int fd, void *buf, size_t count)
{
    size_t left=count;
    ssize_t beenRead=0;
    char *charBuffer=buf;
    while(left>0)
    {
        beenRead=read(fd,charBuffer,left);
        if(beenRead==-1)
        {
            return -1;
        }
        else if(beenRead==0)
        {
            return -2;
        }
        left-=beenRead;
        charBuffer+=beenRead;
    }
}
```

```

    return count;
}

```

- Interfața de programare: bazată pe socketBSD După cum am menționat anterior, serverul va fi concurrent. La nivel de implementare, vom folosi Pthread pentru UNIX.

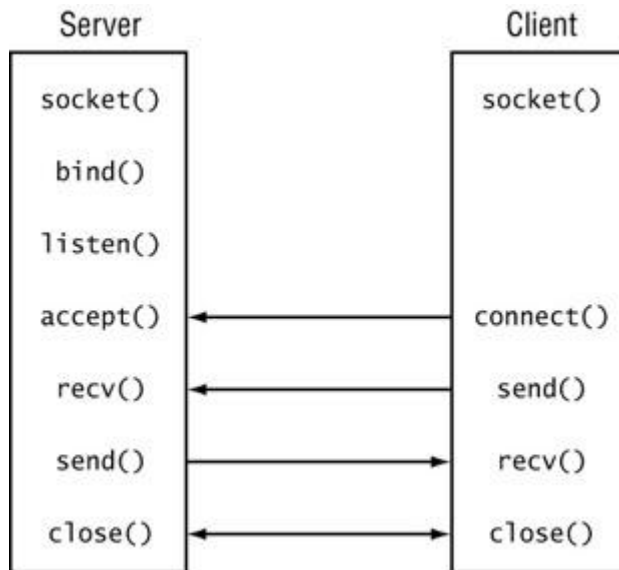


Fig. 14. Server concurrent TCP [14]

login server vs login client:

```

void login(int fd,MYSQL *con)// server
{
    int length=myRecvLength(fd);
    char *user=malloc(length);
    myRecvMessage(fd,length,user);
    user[length]='\0';
    length=myRecvLength(fd);
    char *pass=malloc(length);
    myRecvMessage(fd,length,pass);
    pass[length]='\0';
    char aux[1000];
    sprintf(aux,"select * from login where user='%s' and pass='%s'",user,pass);
    free(user);
}

```

```

free(pass);
MYSQL_RES * res=myQuery(con,aux,fd);
if(mysql_num_rows(res)==0)
{
    mySendType(fd,'1');
    char *message="Ati gresit userul/parola.\n";
    mySendLength(fd,strlen(message));
    mySendMessage(fd,strlen(message),message);
}
else
{
    mySendType(fd,'0');
    char *message="Succes.\n";
    mySendLength(fd,strlen(message));
    mySendMessage(fd,strlen(message),message);
}
mysql_free_result(res);
}

void login(int fd)// client
{
    if(flagLOGIN!=0)
    {
        printf("Eroare. Sunteti logat deja.\n");
        return;
    }

    mySendHeaderedMessage(fd,5,"login");

    char user[100],pass[100];
    printf("user: ");
    fflush(stdout);
    scanf("%s",user);
    printf("pass: ");
    fflush(stdout);
    scanf("%s",pass);
    mySendHeaderedMessage(fd,strlen(user),user);
    mySendHeaderedMessage(fd,strlen(pass),pass);

    char type=myRecvType(fd);
    int length=myRecvLength(fd);
    char *message=malloc(length+1);
    myRecvMessage(fd,length,message);
    message[length]='\0';
    printf("%s",message);
}

```

```

        fflush(stdout);
        free(message);

        if(type=='0')
        {
            flagLOGIN=1;
            strcpy(userGlobal,user);
        }
    }
}

```

5 Concluzii

În concluzie, soluția propusă oferă funcționalitatea cerută. Desigur, aceasta poate fi îmbunătățită prin adăugarea unui coș de cumpărături, combinarea mai multor metode de recomandare (folosirea filtrării bazate pe produs dacă baza de date devine prea mare), realizarea unei interfețe grafice.

References

1. Sabin Buraga și Gabriel Ciobanu, Atelier de programare în rețele de calculatoare, Polirom, Iași, 2001.
2. Cursurile disciplinei *Rețele de calculatoare* din cadrul FII, UAIC (predate de Conf. dr. Alboiaie Lenuța)
3. Stevens, Fenner, Rudoff, UNIX Network Programming Volume 1, Third Edition: The Sockets Networking API
4. <http://www.slideserve.com/nyoko/sisteme-de-recomandare-cu-filtrare-colaborativa>
5. <http://www.greenend.org.uk/rjk/tech/smtpreplies.html#connect>
6. <http://stackoverflow.com/questions/4813890/sqlite-or-mysql-how-to-decide>
7. <http://guidetodatamining.com/chapter2/>
8. <http://www.slideshare.net/lmsasu/curs-2-data-mining>
9. https://en.wikipedia.org/wiki/Cosine_similarity
10. <http://files.grouplens.org/papers/FnT%20CF%20Recsys%20Survey.pdf>
11. <http://stackoverflow.com/questions/26849705/how-to-transfer-pdf-mp3-or-mp4-files-using-socket>
12. http://liuj.fcu.edu.tw/net_pg/socket.html
13. <http://zetcode.com/db/mysqlc/>
14. <http://www.codeproject.com/Articles/12893/TCP-IP-Chat-Application-Using-C>