

# Cuprins

<b>Capitolul 1</b>	<b>Introducere</b>	<b>1</b>
1.1	Descrierea temei . . . . .	1
1.2	Contextul proiectului . . . . .	1
1.3	Metoda abordată . . . . .	2
1.4	Motivație . . . . .	2
1.5	Structura lucrării . . . . .	3
<b>Capitolul 2</b>	<b>Obiective și specificații</b>	<b>4</b>
2.1	Obiectivele proiectului . . . . .	4
2.2	Specificația temei . . . . .	5
<b>Capitolul 3</b>	<b>Studiu bibliografic</b>	<b>7</b>
3.1	Algoritmii Fast Replica . . . . .	7
3.1.1	Fast Replica . . . . .	7
3.1.2	ALM-FR . . . . .	9
3.2	GridFTP . . . . .	10
3.2.1	Protocolul GridFTP . . . . .	10
3.2.2	Globus GridFTP . . . . .	11
3.3	Fast Parallel File Replication . . . . .	14
3.3.1	Abordare ierarhică . . . . .	14
3.3.2	I-BFS . . . . .	15
3.3.3	I-DFS . . . . .	16
3.3.4	I-SW . . . . .	16
<b>Capitolul 4</b>	<b>Fundamentare teoretică</b>	<b>18</b>
4.1	Rețele cu distribuție de conținut . . . . .	18
4.2	Arhitectura client-server . . . . .	19
4.3	Transferuri third-party . . . . .	19
4.4	TCP . . . . .	21
4.5	Tehnici de transfer . . . . .	21
4.5.1	Tehnica PUSH . . . . .	21
4.5.2	Tehnica PULL . . . . .	22

<b>Capitolul 5</b>	<b>Analiză și design</b>	<b>24</b>
5.1	Analiza problemei . . . . .	24
5.2	Distribuția serială . . . . .	25
5.3	Distribuția în doi pași . . . . .	26
5.4	MC-5 . . . . .	28
5.4.1	Replicarea pe două noduri . . . . .	28
5.4.2	Replicarea pe seturi reduse de noduri . . . . .	32
5.4.3	Analiza soluției propuse . . . . .	35
<b>Capitolul 6</b>	<b>Proiectare de detaliu și implementare</b>	<b>38</b>
6.1	Mediul de dezvoltare . . . . .	38
6.1.1	Python . . . . .	38
6.1.2	Twisted . . . . .	38
6.2	Implementare . . . . .	40
6.2.1	Implementare client . . . . .	40
6.2.2	Implementare server . . . . .	41
6.3	Protocoale de comunicare . . . . .	44
6.3.1	Client - Server . . . . .	44
6.3.2	Server - Server . . . . .	46
6.4	Calibrare parametrii . . . . .	47
<b>Capitolul 7</b>	<b>Testare și validare</b>	<b>50</b>
7.1	Testare funcțională . . . . .	50
7.2	Testare eficiență . . . . .	50
7.3	Comparare implementări . . . . .	53
<b>Capitolul 8</b>	<b>Manual de instalare și utilizare</b>	<b>55</b>
8.1	Instalare . . . . .	55
8.1.1	Resurse hardware . . . . .	55
8.1.2	Resurse software . . . . .	55
8.1.3	Instrucțiuni instalare . . . . .	56
8.2	Manual de utilizare . . . . .	56
8.3	Exemplu de utilizare . . . . .	57
<b>Capitolul 9</b>	<b>Concluzii</b>	<b>58</b>
9.1	Analiza rezultatelor obținute . . . . .	58
9.2	Dezvoltări ulterioare . . . . .	58
<b>Bibliografie</b>		<b>60</b>
<b>Capitolul 10</b>	<b>Anexe</b>	<b>62</b>

# Capitolul 1

## Introducere

### 1.1 Descrierea temei

În această eră a vitezei și a informației se pune foarte mult accentul pe acumularea și partajarea resurselor pentru a putea avea o putere computațională mare. Astfel au apărut sistemele distribuite, și odată cu ele și necesitatea de transfer de date între nodurile sistemului.

În această lucrare este analizată situația în care se dorește replicarea unor date pe mai multe noduri destinație din sistem. Deși pare o problemă banală la prima privire, ceea ce face această sarcină dificilă este faptul că numărul de noduri dintr-un sistem distribuit este de ordinul zecilor și chiar sutelor. De asemenea, volumul datelor de transferat poate fi foarte mare, de ordinul GB sau chiar TB, făcând din această problemă o adevărată provocare.

Aplicația propusă în lucrarea de față reușește să combată aceste probleme prin folosirea unor topologii ce implică divizarea nodurilor în grupuri mici și autonome de replicare și prin implementarea unor algoritmi eficienți pentru procesarea datelor ce împiedică apariția erorilor de memorie cauzate de volumele mari de date.

### 1.2 Contextul proiectului

Distribuirea rapidă și sigură a datelor de la o singură sursă la un număr mare de destinații localizate într-un sistem de mari dimensiuni a reprezentat un subiect intensiv cercetat pe parcursul ultimului deceniu. Această cercetare are aplicabilitate în contextul rețelelor cu distribuție de conținut (*en.* Content Delivery Networks), dar poate fi utilizată și în alte domenii. O rețea CDN reprezintă un sistem distribuit de servere de dimensiuni mari desfășurate în mai multe centre de date dispersate geografic. Scopul unui CDN este de a asigura un grad înalt de disponibilitate și de performanță a serviciilor oferite. Se folosește un set dedicat de servere pentru a distribui conținut clienților din partea serverului original pe care se afla datele cerute. În acest scop, conținutul este replicat pe aceste servere, mai

apropiate de utilizator, îmbunătățind astfel latența și reducând în același timp traficul din rețea.

Din punctul de vedere al utilizatorului, este irelevant de pe ce server îi sunt aduse datele, cu condiția ca acestea să fie prezentate în timp util. Rolul unui CDN devine astfel cel de a găsi cea mai bună soluție pentru fiecare utilizator, în funcție de locația acestuia. Pentru a putea oferi aceeași calitate a serviciilor tuturor utilizatorilor, este necesar ca datele să fie replicate pe cât mai multe servere dispersate geografic. În această lucrare vom vedea diferite abordări ale acestei probleme și soluțiile propuse pentru o replicare cât mai eficientă și adaptivă. Deși punctul de plecare al lucrării sunt rețelele cu distribuție de conținut, algoritmi propuși aici pot fi utilizați în orice tip de sistem distribuit ce necesită distribuirea unor volume mari de date la mai multe destinații.

### 1.3 Metoda abordată

Pentru scalarea aplicației la o multitudine de noduri destinație, abordarea folosită aici este de a le împărți în grupuri relativ mici, independente unul de celălalt. Comunicația se face doar în interiorul unui grup de replicare, pentru a menține numărul de conexiuni la un nivel acceptabil. Fiecare grup este responsabil pentru nodurile sale, iar protocoalele de comunicare asigură replicarea fiabilă a datelor.

Pentru a putea obține transferuri de dimensiuni mari, un algoritm de management al memoriei este responsabil pentru procesarea datelor la primirea lor. Deoarece una dintre cerințe este eficientizarea timpilor de transfer, aplicația se adaptează la specificațiile rețelei și la variațiile parametrilor.

### 1.4 Motivație

Pentru a motiva alegerea acestei teme, vom analiza un exemplu concret de utilizare. Să luăm în considerare existența unui sistem distribuit cu sute sau chiar mii de servere distribuite geografic. Presupunem că un utilizator dorește transferul unui fișier de dimensiuni mari de pe un server sursă pe mai multe servere destinație. Transferul acestui fișier folosind conexiuni punct-la-punct (*en.* P2P = Peer to Peer) între sursă și fiecare dintre destinații reprezintă o abordare naivă și prezintă limitări de performanță. După cum se precizează și în [1], o primă problemă este risipirea consumului de lățime de bandă. O altă problemă generată de această abordare este faptul că fiecare rata de transfer pentru fiecare destinație este limitată de caracteristicile conexiunii P2P.

Mai precis, dacă între serverul sursă  $S$  și una dintre destinații  $A$  există o cale cu viteza de transfer  $v_{SA}$ , este posibil să existe un alt server destinație  $B$ , astfel încât calea la serverul  $A$  ce trece prin  $B$  să aibă o viteză de transfer mai mică decât cea directă:  $v_{SA} > (v_{SB} + v_{BA})$ . În acest caz, este mai eficient ca transferul pe serverul  $A$  să se facă prin intermediul serverului  $B$ .

Chiar dacă am considera căile de la serverul sursă la destinații ca fiind cele mai eficiente, avem două noi probleme:

- În cazul în care sursa deschide conexiuni la toate destinațiile, poate să apară o problemă de suprasarcină, deoarece numărul de conexiuni suportate de un server este limitat. Aceasta poate rezulta în închiderea unor conexiuni, pentru a putea păstra altele deschise. Lățimea de bandă este împărțită între conexiuni, crescând timpul total de transfer.
- În cazul în care sursa deschide conexiuni la destinații în mod serial, adică fiecare nouă conexiune este deschisă după închiderea precedentei, fiecare server destinație trebuie să aștepte ca transferul la serverele care sunt înaintea lui în lista de transfer să primească fișierul. Lățimea de bandă poate să nu fie folosită în totalitate. Și în acest caz, timpul de transfer este mare, deoarece reprezintă suma tuturor timpilor de transfer către fiecare destinație.

## 1.5 Structura lucrării

Structura următoarelor capitole este:

- *Capitol 2: Obiective și specificații* | în acest capitol sunt descrise obiectivele propuse și sunt detaliate specificațiile funcționale ale problemei prezentate în lucrarea de față.
- *Capitol 3: Studiu bibliografic* | conține prezentarea și analiza altor lucrări ce tratează această problemă. Fiecare abordare este detaliată, pentru a se putea observa avantajele și dezavantajele fiecăreia.
- *Capitol 4: Fundamentare teoretică* | în acest capitol sunt descrise cunoștințele necesare pentru ca cititorul să înțeleagă această lucrare din punct de vedere teoretic.
- *Capitol 5: Analiză și design* | această parte prezintă schema propusă ca și soluție, descriind funcționalitatea fiecărui modul al acesteia.
- *Capitol 6: Proiectare de detaliu și implementare* | conține informații detaliate despre implementare.
- *Capitol 7: Testare și validare* | prezintă testele efectuate pentru a demonstra funcționalitatea aplicației și cele pentru a determina eficiența acesteia.
- *Capitol 8: Manual de utilizare* | conține detaliile legate de instalarea și utilizarea aplicației.
- *Capitol 9: Concluzii* | această secțiune conține concluziile lucrării, precum și sugestii pentru dezvoltări ulterioare.

# Capitolul 2

## Obiective și specificații

### 2.1 Obiectivele proiectului

Înainte de a defini specificațiile funcționale ale acestei lucrări, trebuie stabilit un set concret de obiective. Scopul este ca în final aplicația să îndeplinească toate obiectivele descrise în acest capitol. Obiectivele propuse sunt descrise în ordine cronologică, după cum urmează:

#### **Studierea domeniului**

Acest pas implică un studiu bibliografic intensiv în domeniul sistemelor distribuite, și mai precis al transferului volumelor mari de date. Principalele etape ale acestui pas sunt: studiul rețelelor cu distribuție de conținut, studiul transferurilor volumelor mari de date în sisteme distribuite, studiul diferitelor topologii și rețele, studierea și înțelegerea unei rețele de distribuție de tip multicast, studiul implementărilor existente și înțelegerea acestora.

#### **Analiza și compararea implementărilor existente**

După familiarizarea cu domeniul temei și înțelegerea diferitelor implementări existente, urmează o analiză a acestora, pentru a determina direcția de urmat. Fiecare lucrare existentă are avantaje și dezavantaje, iar scopul acestui pas este de a pune în balanță caracteristicile fiecăreia pentru a le alege pe cele mai bune în propria implementare.

#### **Dezvoltarea de topologii**

Analiza topologiilor existente rezultă în generarea unor noi structuri ce combină și balansează avantajele descrise mai sus. Aceste topologii sunt realizate între un nod sursă și minim două noduri destinație. Scopul este de a reuși generarea unor topologii eficiente pentru un număr mare de servere destinație, foarte dispersate geografic, și pentru transferul datelor de dimensiuni foarte mari. Pentru a diminua timpul total de transfer trebuie limitată comunicarea dintre noduri, dar în același timp trebuie ținut cont de faptul că adaptabilitatea se poate face doar prin comunicare.

**Dezvoltarea algoritmilor de transfer**

Generarea topologiilor este primul pas în obținerea unui transfer eficient, dar este foarte important și cum sunt transferate datele. Nodul sursă este foarte important în această etapă, deoarece reprezintă primul pas de distribuție.

**Implementarea diferiților algoritmi**

Acest pas este doar de implementare a algoritmilor găsiți anterior, urmărind topologiile studiate. Algoritmii pot fi modificați în cazul apariției unor limitări de implementare. De asemenea, în acest pas se pot folosi cunoștințe legate de mediul de implementare pentru a optimiza procesele de transfer.

**Testare și comparare**

În timpul analizei se fac diverse calcule pentru a determina timpii de transfer. În acest pas se face validarea acestor calcule, și se compară diferitele implementări, atât din punct de vedere al funcționalității, cât și al eficienței. Tot în acest pas trebuie determinate limitările pentru fiecare soluție: numărul maxim de noduri destinație, dimensiunea maximă a datelor ce pot fi transferate (fără a se lua în considerare capacitatea maximă a serverelor, care se presupune a fi infinită).

**Analiza rezultatelor**

Acest pas presupune o retrospectivă asupra muncii efectuate. Se analizează implementările făcute, cu avantajele și dezavantajele lor și se ajunge la un set de concluzii.

## 2.2 Specificația temei

Date de intrare:

**Fișier** Notăția generică de fișier se referă la un volum mare de date existent pe un sistem de calcul, numit sursă.

**Nod sursă** Noțiunea de nod se referă la o componentă a unui sistem distribuit, iar sursă la existența fișierului pe acesta.

**Noduri destinație** Tot componente ale sistemului distribuit, dar care nu posedă fișierul.

**Client** Componentă pasivă a sistemului distribuit, deoarece face parte din acesta doar pentru un scurt timp.

Cerințe funcționale:

- proiectarea unui sistem distribuit care conține toate elementele descrise mai sus.
- proiectarea unei arhitecturi de rețea pentru conectarea componentelor.
- dezvoltarea de protocoale de comunicare între componente.

- dezvoltarea de protocoale de transfer între componentele active ale sistemului.
- dezvoltarea de protocoale de replicare între componentele active ale sistemului.

Cerințe non-funcționale:

- scalabilitate – din punct de vedere al mărimii fișierului și al numărului de noduri destinație.
- eficiență – din punct de vedere al timpului de replicare.
- adaptabilitate – din punct de vedere al lățimii de bandă de pe conexiuni și al dimensiunii pachetelor de date.
- fiabilitate – din punct de vedere al nodurilor pe care se efectuează replicarea cu succes.



# Capitolul 3

## Studiu bibliografic

În acest capitol sunt analizate în detaliu alte lucrări, articole sau proiecte care implementează conceptele prezentate în lucrarea de față. Acestea au atât rolul de a informa în legătură cu alte soluții existente, cât și de a crea o bază teoretică și comparativă pentru implementarea proprie.

### 3.1 Algoritmii Fast Replica

#### 3.1.1 Fast Replica

În [2] este prezentat algoritmul *FastReplica*, un algoritm de distribuire prin tehnica de împingere a datelor (*en. push*). Aici se consideră, ca și în lucrarea de față, o rețea distribuită de mari dimensiuni și problema distribuirii de conținut în această rețea. În aceasta se găsesc:

- motivele pentru necesitatea unui asemenea algoritm – este vorba despre rețelele de distribuire de conținut CDN în care datele sunt replicate pentru a se evita supraîncărcarea unui server și blocajele în rețea.
- primele 3 cele mai populare metode de distribuire a conținutului în Internet: distribuire prin satelit, distribuire multicast, distribuire multicast la nivel de aplicație.
- descrierea algoritmului pentru două tipuri de rețele: mici (10-30 noduri) - *Fast Replica in the small* și mari (sute de noduri) - *Fast Replica in the large*
- extinderea algoritmului pentru a asigura transferuri de încredere (*en. reliable*).
- evaluarea performanțelor și compararea cu calculele precedente.

Există câteva idei de bază exploatate în [2]. Pentru a replica un fișier mare pe  $n$  noduri ( $n$  este în intervalul de 10-30 noduri), fișierul original este împărțit în  $n$  sub-fișiere

de dimensiuni egale și fiecare sub-fișier este transferat la un nod diferit din grup. După aceea, fiecare nod își propagă sub-fișierul celorlalte noduri din grup. Astfel, în loc de replicarea obișnuită a unui fișier întreg pe  $n$  noduri, cu ajutorul a  $n$  căi, conectând nodul original la întregul grup de replicare, FastReplica exploatează  $n * n$  căi din cadrul grupului de replicare unde fiecare cale este utilizată pentru transferul a  $1/n$  din fișier. Scopul este de a proiecta un algoritm scalabil și fiabil, care poate fi folosit pentru replicarea fișierelor de dimensiuni mari la un grup mare de noduri.

În figura 3.1 se pot observa cei doi pași ai algoritmului *FastReplica*: distribuție – 3.1(a) și colecție – 3.1(b).

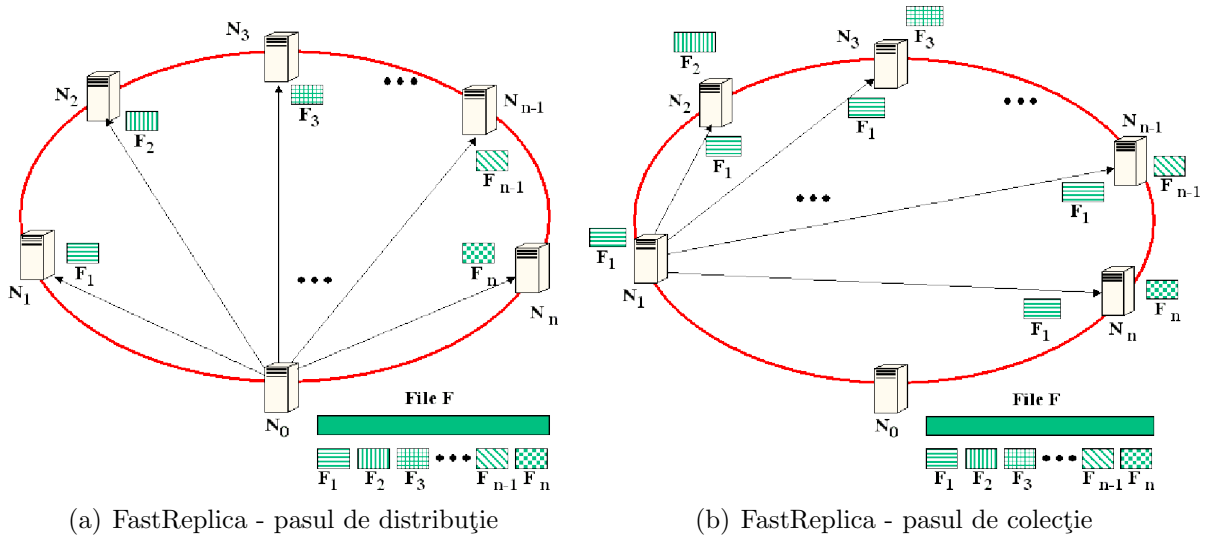


Figura 3.1: Ilustrare algoritm FastReplica

După cum se poate vedea, algoritmul *FastReplica* prezintă o multitudine de avantaje:

- timpul de replicare al fișierului este minimizat prin folosirea tuturor căilor disponibile în rețea:  $Timp_{mediu} = \frac{1}{n} \sum_{i=1}^n Timp^i(F)$  unde  $Timp^i(F)$  reprezintă timpul de transfer al fișierului  $F$  de la nodul sursă la nodul destinație  $N_i$ , iar  $n$  reprezintă numărul total de servere destinație.
- replicarea se face în doar doi pași: distribuție și colecție.
- modelul este simplu și ușor de înțeles, ceea ce îl face și ușor de implementat.
- algoritmul poate fi generalizat la un număr mult mai mare de noduri prin formarea așa-ziselor grupuri de replicare, utilizând o metodă iterativă.
- algoritmul este extins pentru a detecta și a corecta eventualele defecte ale nodurilor apărute atât înaintea începerii, cât și în timpul transferului.

Deși aceste avantaje fac *FastReplica* un model foarte folosit, studiul acestuia a scos la iveală și anumite dezavantaje:

- în timpul pasului de distribuție, nodul sursă  $N_0$  are  $n - 1$  conexiuni de ieșire, una către fiecare nod din grupul de replicare. Pentru valori mari ale lui  $n$ , numărul total de conexiuni poate cauza probleme.
- în timpul pasului de colecție, fiecare nod  $N_i$  are următoarele conexiuni:  $n - 1$  conexiuni de intrare și  $n - 1$  conexiuni de ieșire, totalizând  $2n - 2$  conexiuni pentru fiecare nod destinație. Pentru valori mari ale lui  $n$ , numărul total de conexiuni poate cauza probleme.
- conexiunile din pasul de distribuție nu sunt folosite la capacitatea lor maximă, deoarece după terminarea acestui pas, pe aceste conexiuni nu se mai transferă date utile către serverele destinație.
- nodul sursă trebuie să trimită fiecărui nod destinație adresele celorlalte noduri, pentru ca acesta să se poată conecta la ele. Astfel, comunicarea dintre noduri crește exponențial cu numărul total de destinații.

### 3.1.2 ALM-FR

În [3] ideile *FastReplica* sunt combinate cu tehnica nivelului de aplicație multicast (*en Application Level Multicast*) pentru optimizarea replicării fișierelor mari în cadrul CDN. Această metodă este denumită *ALM-FR*. Un transfer de fișiere în algoritmul clasic *FastReplica* urmează un mod de distribuire de tip stocare și înaintere, adică doar după ce un sub-fișier este primit de către un nod, acest nod începe să transmită sub-fișierul corespunzător nodurile ulterioare. În acest model, atunci când apare un nod defect, impactul se propagă doar la nodurile din grupul local de replicare. Algoritmul de recuperare este oarecum natural și simplu: aceasta implică doar nodurile din grupul replicare corespunzător.

Deoarece tehnica de multicast IP nu beneficiază de o implementare larg răspândită, multe eforturi industriale și de cercetare și-au redirecționat atenția către investigarea și implementarea multicastului la nivelul aplicației, în care nodurile distribuite pe internet acționează ca routere intermediare pentru a distribui eficient conținutul de-a lungul unui arbore predefinit. Mai detaliat, această lucrare propune împărțirea setului original de noduri în grupurile de replicare  $G_0, G_1, \dots, G_{k-1}$ , fiecare format din  $k$  noduri. Aceste grupuri sunt ordonate în arborii speciali de multicast  $\hat{M}, M^1, M^2, \dots, M^{m-1}$ . Fiecare dintre acești arbori este format din maxim  $m$  grupuri, unde  $\hat{M}$  reprezintă arborele multicast primar, iar  $M^1, M^2, \dots, M^{m-1}$  reprezintă arborii multicast secundari. Schema rezultată poate fi analizată în figura 3.2.

Fișierul este divizat, ca și în algoritmul original, în  $k$  sub-fișiere egale ca dimensiuni  $F_0, F_1, \dots, F_k$ . Diferența apare în pașii de execuție:



- performanța – acesta are suport inclus pentru utilizarea conexiunilor TCP paralele și a transferurilor multiple pentru a atinge rate de transfer mari.
- puncte de control – acestea sunt utile pentru a relua un transfer eșuat doar de la un anumit punct.
- transferuri de la terți – protocolul FTP, pe care se bazează GridFTP, are canale separate de control și de date, permițând transferuri de la terți, adică transferuri între două noduri, mediate de un al treilea ce acționează ca și client.
- securitate – oferă securitate pe ambele canale. Canalul de control este criptat în mod implicit, deoarece este canalul pe care se face prima dată comunicarea, iar canalul de date folosește autentificare, cu posibilitatea de a avea și protecția integrității și criptare.

### 3.2.2 Globus GridFTP

Globus Toolkit® este un *framework open source* pentru construirea de rețele de calcul, care oferă o implementare vastă și actuală a protocolului GridFTP. Acest *framework* oferă:

- o implementare pentru server (*globus-gridftp-server*).
- o implementare pentru client bazată pe scripturi ( *globus-url-copy* ).
- un set de biblioteci pentru dezvoltarea de clienți noi.

În [5] și [6] este prezentată munca cercetătorilor de la Argonne National Laboratory (ANL) și de la Ohio State University (OSU) legată de această temă și pusă în aplicare folosind Globus GridFTP. În [5] sunt prezentați algoritmi dinamici ce efectuează transferul simultan al bucăților de fișier din multiple replici ale acestuia, localizate în diverse surse. Algoritmul este adaptiv, atât în selectarea surselor, cât și în schimbarea dinamică a lățimii de bandă din rețea. Articolul [6] propune o altă modalitate de implementare a conceptelor GridFTP, prin generarea de rețele de acoperire (*en. overlay*). Beneficiile oferite de aceste tipuri de rețele sunt evidențiate printr-o arhitectură ce include două optimizări cheie:

- *multi-hop path splitting* – înlocuirea unei conexiuni directe între sursă și destinație cu o cale multi-hop, înălțuită prin noduri intermediare.
- *multi-pathing* – împărțirea datelor la sursă și trimiterea acestora pe multiple căi obținute folosind pasul anterior. Altfel spus, mai multe rute independente pot fi folosite pentru a transfera simultan bucățile de fișier la destinație.

### GridFTP Multicast

Implementarea Globus GridFTP oferă posibilitatea de a trimite pachete de date la mai multe destinații. Pentru aceasta, este necesar să se includă în stiva de drivere modulul de multicast. Atunci când un pachet ajunge la un server configurat în acest fel, el trece prima dată prin modulul de multicast, care îl transmite pe alte căi, la următoarele servere din structură (3.3(b), 3.3(c), 3.3(d)). Arhitectura GridFTP generează o topologie de tip arbore binar, după cum se poate vedea și în figura 3.3(a). Nodul sursă reprezintă rădăcina acestui arbore, iar nodurile destinații fii acestuia. [7]

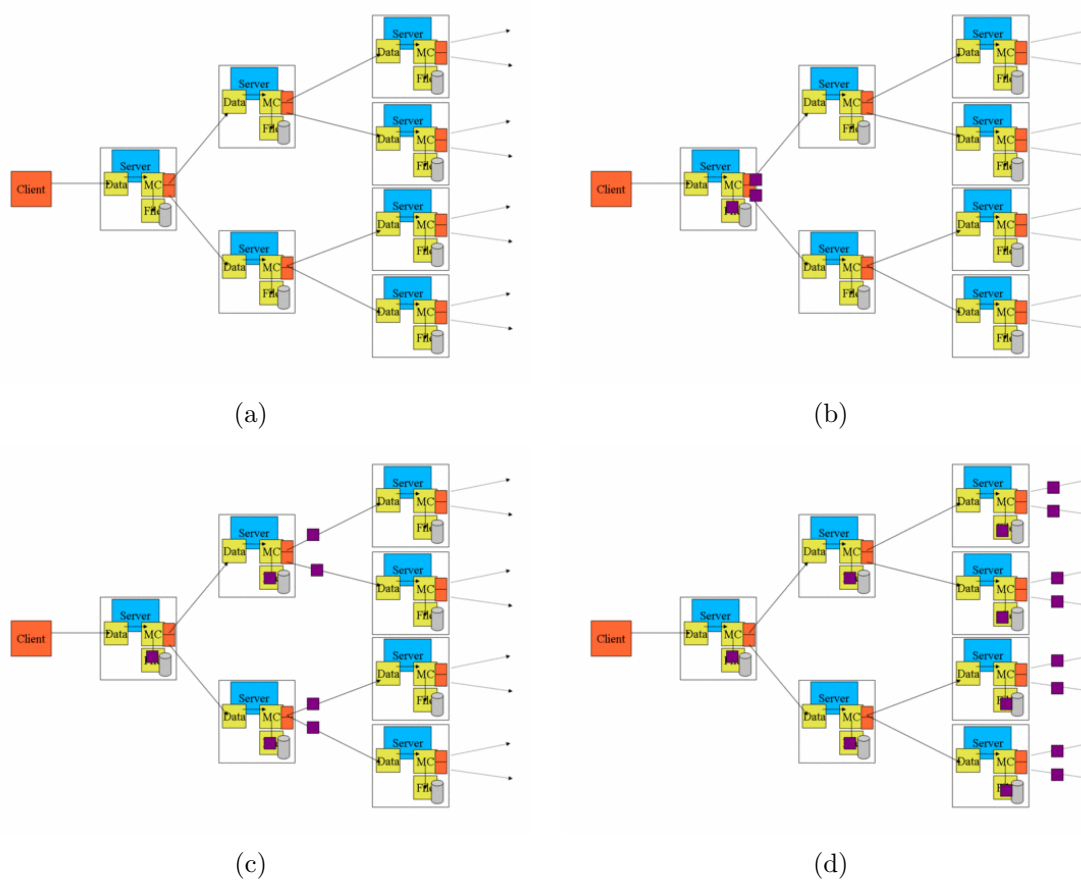


Figura 3.3: GridFTP multicast - schema generală pe pași [7]

Avantajele acestei abordări:

- Are o structură simplă, ușor de implementat și de utilizat. Structura de arbore binar se scalează bine pentru un număr mare de noduri destinație, deoarece numărul de conexiuni pentru fiecare nod rămâne maxim  $N$ .

- Are o arhitectură modularizată, ce permite atât utilizarea în combinație cu alte drivere existente, cât și extinderea prin implementarea de noi drivere.
- Timpul maxim de transfer este eficientizat prin faptul că fiecare pachet este trimis mai departe imediat ce este primit.
- Oferă posibilitatea de a crește numărul de fii  $N$  pentru a obține performanțe ridicate.

Dezavantaje:

- Pentru valori ale lui  $N$  egale cu numărul de noduri destinație, nu există nicio diferență între această topologie și o transmitere în serie a datelor.
- Volumul de date de control (adresele nodurilor, dimensiuni fișier) trimise pe legături este mare. De exemplu, primele două noduri din ierarhie trebuie să primească date despre  $N/2$  servere fiecare.
- Serverele aflate pe pozițiile frunze vor primi date care trec prin toți strămoșii acestora. În această trecere, există posibilitatea deteriorării sau a pierderii complete a datelor.
- Soluția nu este fiabilă: în cazul în care unul dintre serverele din ierarhie eșuează în timpul transferului, nici unul dintre nodurile din subarborele acestuia nu va mai primi pachetele necesare.

### GridFTP Striped

În cadrul proiectului Globus GridFTP a apărut și un *framework* numit GridFTP Striped. Ideile de bază ce au stat la apariția acestuia sunt descrise în [8]. Globus GridFTP Striped reprezintă un set de biblioteci pentru client și server, concepute pentru construirea de aplicații de date. Acesta folosește ideea că sursele de date sunt de multe ori clustere și organizează datele într-o structură adaptată. Astfel, datele sunt ținute pe mai multe servere care au un sistem comun de fișiere. Această organizare este comandată de faptul că uneori un server poate fi format din mai multe dispozitive fizice, care pot comunica între ele. Există două tipuri de noduri:

- *front end node* – acesta este nodul cu care se realizează conexiuni, cel care reprezintă celelalte noduri. Doar prin acesta se pot accesa datele, dar pe mașina fizică ce conține acest nod nu se află date, ci doar informații despre nodurile de date.
- *back end node* – pe aceste noduri se găsesc datele reale. Doar nodul front end se poate conecta la acestea, în cazul unui transfer atât de intrare, cât și de ieșire.

Arhitectura detaliată se poate vedea în figura 3.4.

Acest model are ca și caracteristici:

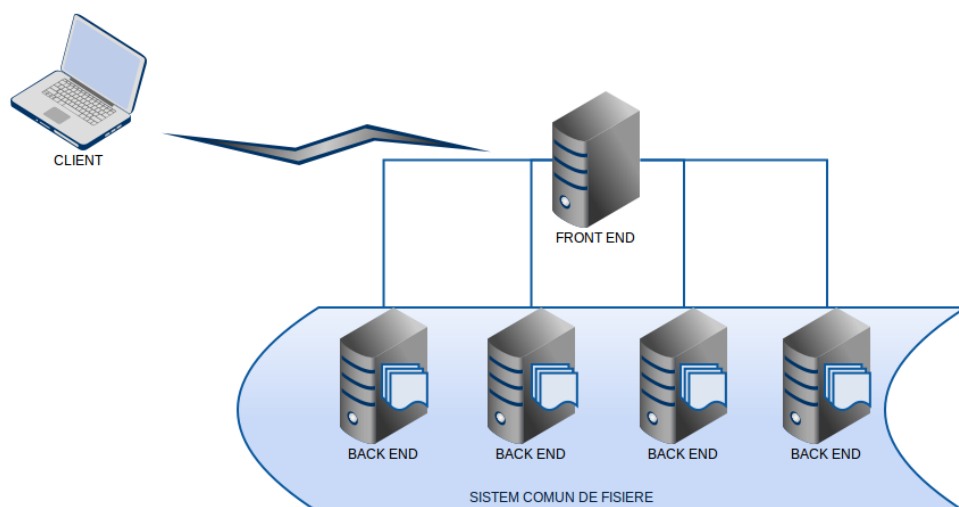


Figura 3.4: Arhitectura GridFTP striped

- modularitate: pentru utilizarea în diferite medii și cu diverse configurații fără nevoia de a modifica implementarea.
- eficiență: se preferă evitarea păstrării unor copii inutile de date, iar pentru că acestea sunt divizate pe mai multe locații fizice, crește performanța, cât și limitările existente.
- siguranță și securitate: deoarece doar nodul principal are acces la datele de pe nodurile secundare, sunt limitate accesurile neautorizate la posibile date confidențiale sau de importanță ridicată.

### 3.3 Fast Parallel File Replication

Replicarea paralelă a fișierelor mari simultan la mai multe locații este o parte integrantă a mediilor de date de tip grid, după cum se precizează și în [9]. În acest articol este analizat mediul de transfer GridFTP punct-la-punct și este conceput un mod de transfer punct-la-multipunct paralel pe baza acestuia. Această lucrare prezintă un motor de replicare rapidă paralelă a fișierelor, care creează mai mulți arbori de distribuție prin centralizarea sesiunilor GridFTP și organizarea acestora în structuri logice. În continuare sunt prezentate principalele arhitecturi analizate în [9]. Pentru fiecare dintre aceste arhitecturi se va ilustra arborele obținut prin aplicarea algoritmului pe figura 3.5.

#### 3.3.1 Abordare ierarhică

În această abordare, arborii sunt construiți prin plasarea nodurilor în manieră ierarhică, iar nodurile de pe nivelele superioare sunt utilizate pentru a transmite sub-fișiere



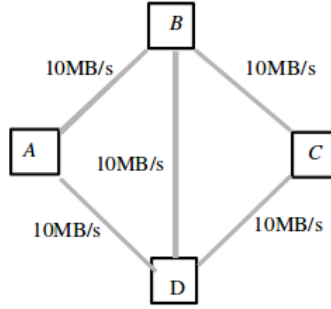


Figura 3.5: Exemplu de noduri și legături într-o rețea de replicare

nodurilor de pe nivelele inferioare. Această arhitectură se aseamănă foarte mult cu algoritmul *FastReplica*, după cum se specifică și în lucrare. Construcția rezultată formează două nivele ale ierarhiei; nu este informată de capacitatea legăturii sau a serverului și este independentă de topologia rețelei care stă la baza sa. Rețelele obținute sunt ilustrate în figurile 3.6(a), 3.6(b), 3.6(c). Sunt obținuți 3 arbori, dar performanța nu este cea așteptată.

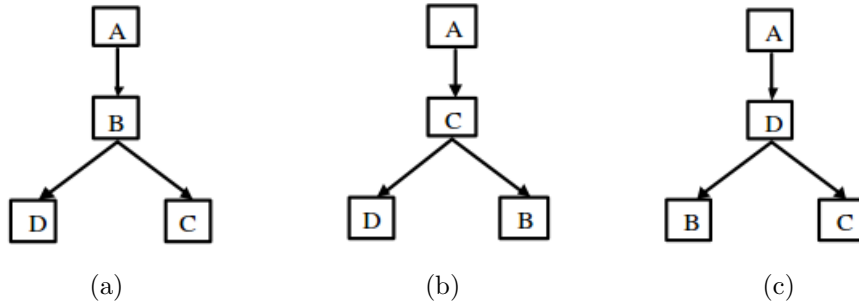


Figura 3.6: Arborii obținuți cu abordarea ierarhică

### 3.3.2 I-BFS

În această abordare, primul arbore este rezultatul unei căutări în lățime (*en. Breadth First Search*). Arborele rezultat are nodul sursă ca rădăcină, și se propagă la toate nodurile destinație, după cum se poate vedea în figura 3.3.2. Pentru arborele rezultat  $T_i$  se obține lățimea minimă de bandă  $B_i$ . Aceasta reprezintă un blocaj în rețea și este deci eliminată din setul de lățimi de bandă ale legăturilor din  $T_i$ . Algoritmul BFS este apelat din nou pentru a crea al doilea arbore, iar procesul este repetat până când nu se mai pot genera alți arbori valizi. Avantajul I-BFS, pe lângă timpul mic de convergență, este faptul că numărul de conexiuni concurente pentru un singur nod rămâne în limite acceptabile.

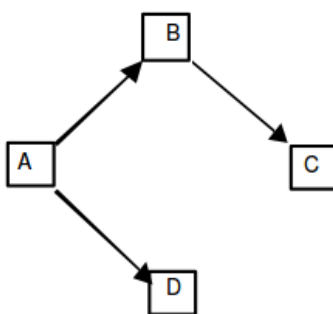


Figura 3.7: Arborele obținut cu algoritmul I-BFS

### 3.3.3 I-DFS

Această soluție este la fel ca I-BFS, cu diferența că algoritmul de căutare în adâncime (*en. Depth First Search*) este folosit pentru construirea arborilor. Rezultatul poate fi văzut în figurile 3.8(a) și 3.8(b). Trebuie remarcat faptul că deși numărul de arbori generați este mai mic decât numărul obținut cu abordarea ierarhică, acest algoritm are cel mai mic timp de convergență și se scalează cel mai bine pe acest exemplu.

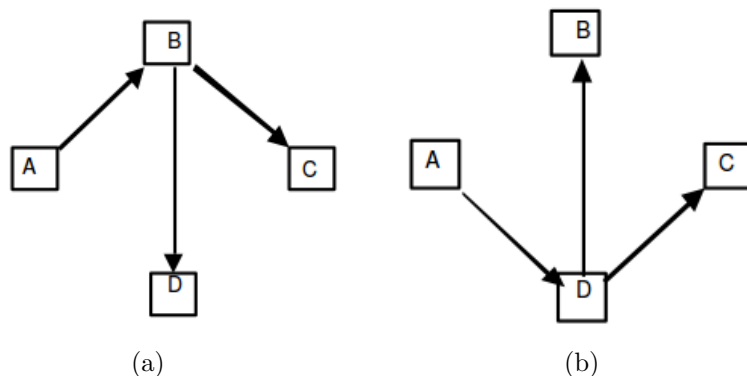


Figura 3.8: Arborii obținuți cu algoritmul I-DFS

### 3.3.4 I-SW

În această ultimă abordare, primul arbore este creat bazat pe găsirea celei mai scurte și mai late căi din rețea (*en. Shortest Widest*). Aceasta este obținută prin utilizarea algoritmului lui Dijkstra, în care metrica folosită pe muchii este lățimea de bandă. După crearea primului arbore, ceilalți sunt generați folosind aceeași manieră iterativă ca și la I-DFS și I-BFS. Figura 3.9 ilustrează structura generată cu acest algoritm.

În aceste ultime abordări (cele iterative) descrise în [9], topologia rețelei și lățimea de bandă disponibilă a legăturilor trebuie să fie cunoscute înaintea aplicării algoritmilor. Aceste informații pot fi achiziționate folosind un serviciu numit NWS (*en. Network Weather Service*) și detaliat în [10].

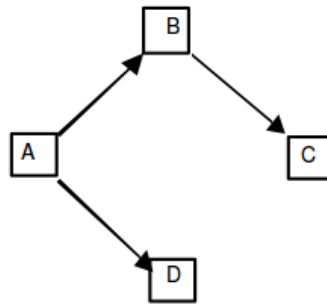


Figura 3.9: Arborele obținut cu algoritmul I-SW

# Capitolul 4

## Fundamentare teoretică

### 4.1 Rețele cu distribuție de conținut

O rețea cu distribuție de conținut este un sistem distribuit de mari dimensiuni, în care serverele sunt instalate în diverse centre de date dispersate pe Internet [11]. Scopul unei asemenea rețele este acela de a oferi conținutul cerut de către clienți în timp real, adică să aibă o performanță ridicată. Pentru a atinge această performanță, rețelele cu distribuție de conținut folosesc replicarea pe un grup de servere pentru a putea deservi cu un grad mare de disponibilitate toți clienții, indiferent de locația fizică a acestora.

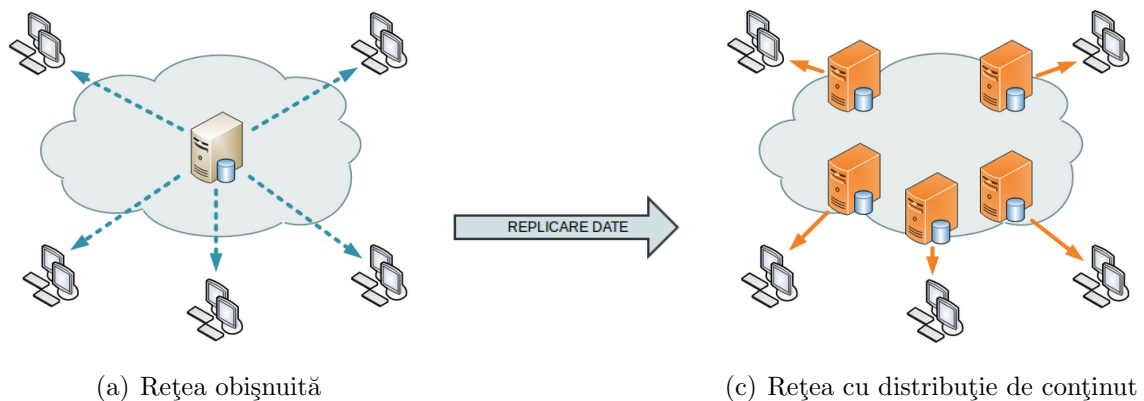


Figura 4.1: Crearea unei rețele de distribuție de conținut. Sursa [12]

Modalitatea de funcționare a unui CDN este descrisă în figura 4.1. În locul unui singur nod ce oferă servicii clienților, ca în figura 4.1(a), rețelele cu distribuție de conținut păstrează copii ale acestor servicii sau fișiere pe mai multe noduri, după cum se poate vedea în 4.1(c). Astfel, atunci când un număr mare de clienți vor să acceseze aceleași date sau servicii, va exista o balansare a cererilor și ofertelor. Se evită în acest mod supraîncărcarea unui singur server, în timp ce altele nu sunt utilizate. Deși au fost proiectate

inițial pentru fluxurile de date continue (*live streaming*), acum acestea sunt folosite pentru orice fel de date: media, obiecte descărcabile, aplicații, rețele sociale. Obiectivul este de a oferi un set de servicii stabile și fiabile, cu performanță ridicată.

## 4.2 Arhitectura client-server

În domeniul sistemelor distribuite, o arhitectură client-server este o structură formată din două părți ce funcționează după următoarele principii:

### Clientul

Este partea care cere de la server servicii sau date; clientul este cel care inițiază sesiunea de comunicare. Clientul este cel care are nevoie de resurse, el nu participă cu propriile resurse la această arhitectură.

### Serverul

Este partea care oferă serviciile și datele; serverul ascultă cererile clienților și le răspunde acestora. Serverele sunt deseori folosite pentru a da spre folosință resurse atât clienților, cât și altor servere. Ele pot restricționa sau limita accesul la aceste resurse, în funcție de cereri.

Pentru a putea comunica, clientul și serverul trebuie să vorbească același limbaj, deci să urmeze același protocol. După stabilirea unei conexiuni, clientul și serverul schimbă mesaje folosind acest protocol, după modelul cerere-răspuns. Acest protocol rulează la nivelul de aplicație. Comunicarea se realizează de obicei în cadrul unei rețele, dar este posibil ca atât clientul cât și serverul să se afle pe același dispozitiv fizic. Un client se poate conecta la mai multe servere, iar un server poate deservi mai mulți clienți în același timp. Clienții nu pot comunica între ei, dar două servere au posibilitatea de a comunica, dacă este necesar.

## 4.3 Transferuri third-party

După cum am precizat și în paragraful anterior, două servere pot comunica între ele. Dar chiar și pentru această comunicație este necesar un client care să o inițieze și să transmită unuia dintre servere datele celuilalt, pentru ca acesta să se poată conecta la cel ulterior. Primul server devine astfel un client pentru cel de-al doilea, dar clientul inițial rămâne totuși informat. Acest tip de comunicație poartă numele de *third party*, sau a treia parte, deoarece clientul este aici o a treia parte a comunicației. În cazul unui transfer de acest fel, clientul trebuie să se conecteze la servere pentru a le transmite cererea sa, iar apoi nodul care conține datele ce trebuie transferate se conectează la nodul care trebuie să primească datele. Această structură poate fi analizată în figura 4.2.

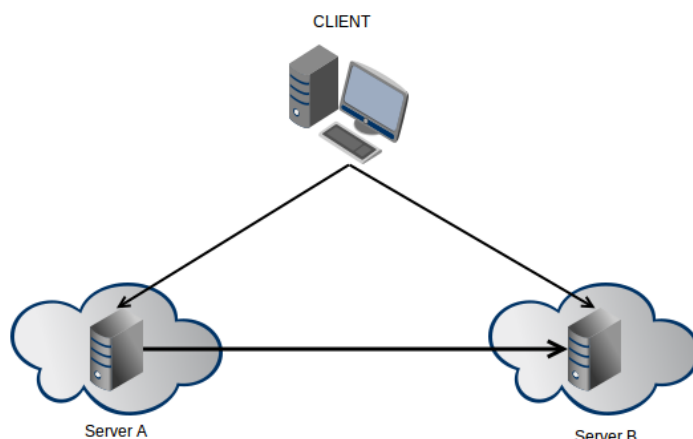


Figura 4.2: Transfer third-party

Canalele pe care clientul se conectează la serverele  $A$  și  $B$  se numesc canale de control, deoarece clientul este cel care controlează întreaga operațiune și care inițiază transferul. Între serverul  $A$  și serverul  $B$  se află canalul de date, deoarece pe această legătură se transferă datele. Protocolul de comunicație dintre client și server diferă de cel dintre servere, iar serverele trebuie să asculte pe două porturi: unul pentru clienți și unul pentru servere. Motivele unei asemenea abordări sunt variate:

- permite mutarea datelor între două locații aflate la distanță de client (prin distanță a se înțelege că serverele și clientul nu se află pe aceeași mașină fizică, ceea ce înseamnă că clientul nu are acces la datele de pe serverul  $A$  decât prin rețea).
- clientul are controlul asupra desfășurării transferului, el fiind cel care oferă serverelor informațiile necesare.
- la nivelul transferului de date există transparență pentru client, acesta nefiind implicat direct.
- este diminuat traficul pe o singură legătură, deoarece există legături diferite pentru date și pentru comenzi. Aceasta poate însemna și o scădere în traficul din rețea.
- serverul sursă preia responsabilitatea pentru efectuarea transferului de la client, devenind chiar el un client în momentul conectării la serverul destinație.
- abordarea third-party are aplicabilitate și în securitatea transferului.

## 4.4 TCP

*Transmission Control Protocol* este unul dintre cele mai importante protocoale de comunicare din suita IP. Acesta este descris în [13]. TCP a fost creat pentru a fi folosit ca protocol de transmisie între două noduri din rețelele bazate pe schimbarea de pachete și în sisteme interconectate. Ca și caracteristici, protocolul TCP este un protocol de mare încredere, fiabil și care oferă funcția de verificare și de corectare a erorilor. TCP aparține nivelului transport, deci este situat deasupra nivelului IP, căruia îi transmite segmente variabile de informație sub forma de datagrame. Această ierarhie se vede în figura 4.3, unde este simbolizată o parte din stiva de protocoale TCP-IP.

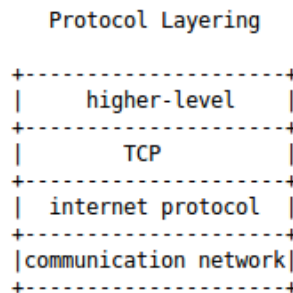


Figura 4.3: O parte din stiva de protocoale TCP-IP

Protocolul TCP este folosit în transferul de fișiere, deoarece oferă încrederea că datele vor ajunge la destinație, chiar dacă există erori în rețea. Este un protocol bazat pe flux de date (*en. data streaming*) care garantează că toți biții trimiși vor fi la fel cu cei primiți, în ordinea corectă. O alternativă la protocolul TCP este UDP (*User Datagram Protocol*), care nu oferă retransmisia datelor pentru a asigura corectitudinea, dar care este mult mai rapid. Pentru acesta este necesară o implementare proprie a anumitor caracteristici, în funcție de necesitățile cerute. Totuși, pentru transferul de fișiere în rețele cu comutare de pachete, cea mai bună alegere este TCP.

## 4.5 Tehnici de transfer

Pentru realizarea transferului de la un nod sursă la un nod destinație, se disting două tehnici de transfer diferite. Acestea sunt descrise în secțiunile următoare. Cele două tehnici sunt identificate din perspectiva nodului care inițiază distribuirea datelor.

### 4.5.1 Tehnica PUSH

Prima strategie de distribuire este cea de împingere a datelor de la sursă la destinație, după cum se poate vedea în figura 4.4. Această strategie este cel mai des folosită

în transferurile de fișiere, deoarece presupune distribuirea datelor cu inițiativa sursei (*en. sender-initiative*). Avantajul unei asemenea strategii este faptul că sursa poate controla transferul, deoarece știe ce date posedă și nu există problema ca datele să nu fie găsite sau să fie corupte.

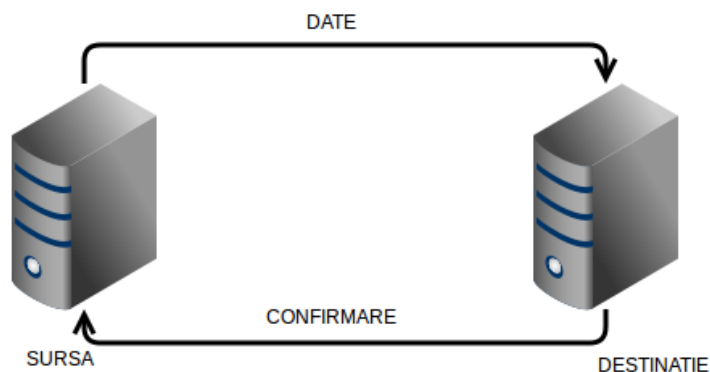


Figura 4.4: Transfer inițiat de nodul sursă

Comunicarea este realizată în doi pași:

- în primul pas, sursa se conectează la destinație și începe să îi trimită acestuia datele pe care le deține. Modalitatea de transmisie poate fi de două feluri: sincronă sau asincronă. În cazul unui transfer sincron, sursa trebuie să aștepte realizarea următorului pas înainte de a începe să trimită din nou date. În cazul unui transfer asincron, sursa poate să ignore confirmarea și să trimită datele în continuare, până la terminarea acestora. Ambele modalități au avantaje și dezavantaje.
- cel de-al doilea pas constă din răspunsul serverului destinație, care trebuie să confirme primirea datelor, în cazul în care transferul este acceptat. În caz contrar, destinația poate să refuze conexiunea sau chiar datele, trimițând înapoi la sursă un mesaj de eroare.

### 4.5.2 Tehnica PULL

Strategia de tragere a datelor de la sursă este opusă celei descrise anterior. În această strategie, cererea de date vine de la destinație, după cum este ilustrat și în figura 4.5. Deși nu așa des utilizată pentru transferul de fișiere, această tehnică oferă o mai mare flexibilitate, deoarece destinația poate să ceară exact datele de care are nevoie, și pe care nu le va refuza. Avantajul acestei tehnici este faptul că serverul destinație joacă rolul de client, el fiind capabil să ceară date de la sursă.

Și aici, comunicația se desfășoară în doi pași:



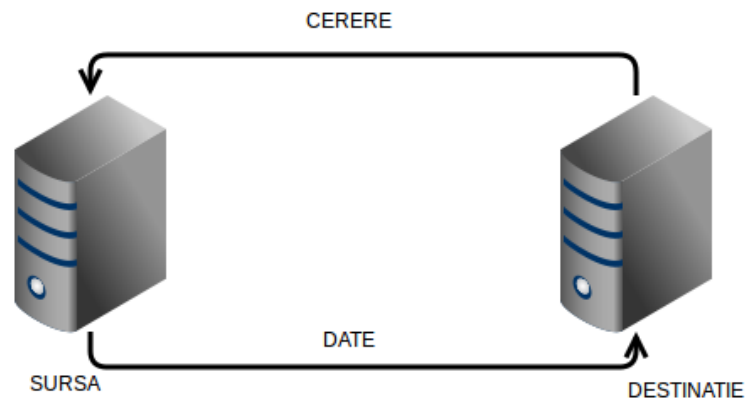


Figura 4.5: Transfer inițiat de nodul destinație

- inițiatorul transferului este destinația (*en. receiver-initiative*), care are nevoie de un anumit set de date. El trimite o cerere către serverul sursă prin care precizează datele de care are nevoie.
- al doilea pas constă din răspunsul serverului sursă, care va începe să trimită datele cerute către serverul destinație. În cazul în care sursa nu deține datele cerute, va răspunde cu un mesaj de eroare.

# Capitolul 5

## Analiză și design

În acest capitol sunt utilizate informațiile obținute în studiul bibliografic pentru a analiza problema ce trebuie rezolvată și pentru a putea propune diferite soluții la aceasta.

### 5.1 Analiza problemei

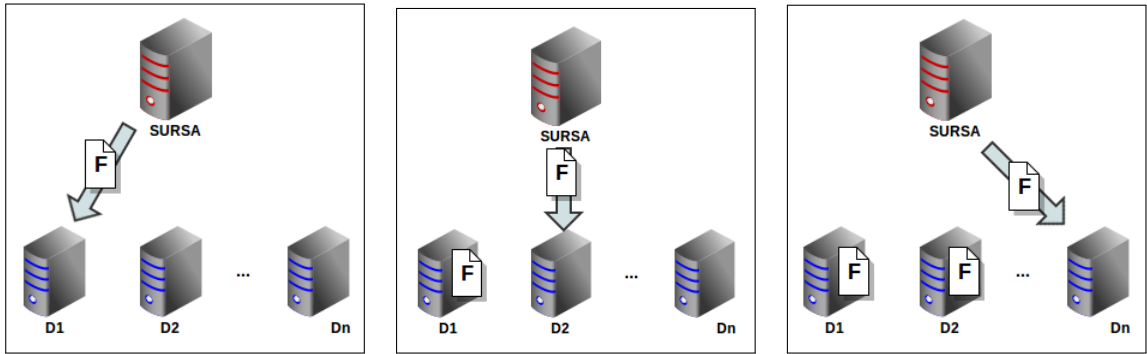
După cum am precizat și în primul capitol, problema abordată în lucrarea de față este a replicării unui set de date sub forma unui fișier pe un număr mare de noduri destinație în cadrul unui sistem distribuit. Scopul este acela de a determina numărul maxim de noduri pe care replicarea se poate face, și, de asemenea, dimensiunea maximă a datelor ce pot fi transferate. Desigur, algoritmul prezentat în această lucrare trebuie să îndeplinească următoarele cerințe:

- să permită transferul unui volum mare de date.
- să aibă o arhitectură scalabilă (*en. scalability*), ce permite replicarea pe un număr foarte mare de noduri destinație.
- să fie adaptiv (*en. adaptability*), deoarece nodurile pe care se face replicarea pot fi dispersate geografic, având lățimi de bandă diferite, de care trebuie ținut cont.
- să fie fiabil (*en. reliability*), adică să permită terminarea transferului chiar și în cazul apariției unor erori.
- să aibă un grad mare de disponibilitate (*en. availability*), ceea ce înseamnă că poate deservi un număr mare de clienți.
- să prezinte rezistență la erori și eșecuri (*en. fault and failure tolerance*); această caracteristică este opusul uneia dintre cele mai mari probleme din domeniul sistemelor distribuite: SPOF (*Single Point Of Failure*).

Până acum am prezentat mai mulți algoritmi rezultați din cercetările efectuate în acest domeniu, alături de avantajele și dezavantajele acestora. În continuare, vom prezenta trei abordări pentru această problemă: primele două sunt folosite pentru a realiza o comparație cu cea de-a treia, care reprezintă aportul personal în acest domeniu.

## 5.2 Distribuția serială

Prima soluție propusă pentru comparare este distribuirea serială a datelor. După cum se poate remarca și din denumire, aceasta presupune trimiterea fișierului la fiecare nod destinație de către nodul sursă. Pentru o mai bună înțelegere a acestei abordări, se poate consulta figura 5.1. Aici se poate vedea cum nodul sursă începe prin a se conecta la prima destinație, căreia îi trimite fișierul în totalitate (5.1(a)). După acest pas, indiferent dacă a fost realizat cu succes sau nu, nodul sursă începe conexiunea și respectiv transferul către următorul nod destinație (5.1(b)). Acest pas este repetat de  $n$  ori, unde  $n$  este numărul total de destinații. După ce fișierul este trimis și ultimului nod din setul de destinații (5.1(c)), transferul poate fi considerat încheiat. După acest ultim pas, în cazul unei rulări fără erori, fișierul se află pe fiecare dintre destinațiile vizate.



(a) Primul pas: primul nod destinație  $D_1$  primește fișierul.

(b) Al doilea pas: al doilea nod destinație  $D_2$  primește fișierul.

(c) Ultimul pas: ultimul nod destinație  $D_n$  primește fișierul.

Figura 5.1: Replicarea datelor pe  $n$  noduri folosind algoritmul de distribuție serială.

La o analiză mai atentă a acestei soluții, se poate vedea că timpul total de transfer este suma timpilor pentru fiecare destinație:

$$Timp_{transfer} = T_1(F) + T_1(F) + \dots + T_n(F) \quad (5.1)$$

unde  $T_i(F)$  reprezintă timpul de transfer al fișierului  $F$  pe nodul  $D_i$ . Dacă notăm cu  $BW_i$  lățimea de bandă a legăturii dintre sursă și nodul  $D_i$  și cu  $Size(F)$  dimensiunea

fișierului  $F$ , atunci relația de mai sus devine:

$$Timp_{transfer} = \frac{Size(F)}{BW_1} + \frac{Size(F)}{BW_2} + \dots + \frac{Size(F)}{BW_n} \quad (5.2)$$

Deși distribuția serială poate părea și chiar este o soluție inefficientă din punct de vedere al timpului, aceasta are totuși anumite avantaje:

- este scalabilă, deoarece indiferent de numărul de noduri destinație, în această arhitectură va exista la orice moment cel mult o conexiune (cea dintre nodul sursă și cel destinație curentă).
- este fiabil, deoarece în cazul în care transferul către unul dintre noduri nu poate fi efectuat, aceasta nu afectează celelalte noduri.
- este de încredere, deoarece sursa nu începe transferul către un alt nod până când nodul curent nu a primit întreg fișierul.

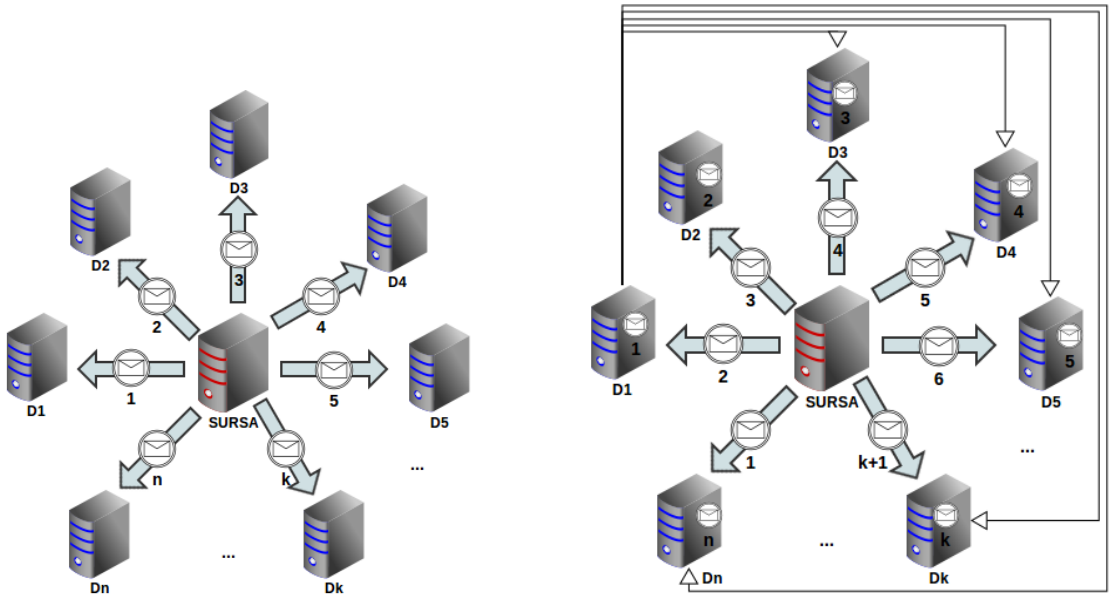
### 5.3 Distribuția în doi pași

O a doua abordare folosită pentru testarea algoritmilor ce urmează este o versiune modificată a *FastReplica*, prezentat în capitolul anterior. Aceasta presupune divizarea documentului în  $n$  părți egale ca dimensiune, numite subfișiere, și executarea a doi pași:

- în primul pas (numit în *FastReplica* pasul de distribuție), nodul sursă trimite fiecărei destinații  $D_k$  subfișierul  $k$ , după cum se poate observa și în figura 5.2(a). Acest pas este aici identic cu cel din *FastReplica*.
- în cel de-al doilea pas (numit în *FastReplica* pasul de colecție), nodul sursă trimite fiecărei destinații  $D_k$  subfișierul  $k+1$ , după cum se poate observa și în figura 5.2(b). Diferența dintre acest algoritm și *FastReplica* este faptul că în acest pas, fiecare nod destinație deschide doar  $n-2$  conexiuni concurente la celelalte noduri destinație.

Avantajele acestei versiuni față de *FastReplica* sunt următoarele:

- deoarece sursa trimite subfișiere și în al doilea pas, aceste legături sunt folosite la capacitatea lor maximă.
- fiecare nod primește două subfișiere de la sursă, ceea ce înseamnă că fiecare nod are  $n-2$  legături de intrare în al doilea pas.
- fiecare nod  $D_k$  trimite subfișierul primit în primul pas la nodurile  $D_{k+2}, D_{k+3} \dots D_{k-1}$  în cel de-al doilea pas. Deci vor exista  $n-2$  legături de ieșire în al doilea pas, totalizând  $2n-4$  conexiuni pentru fiecare destinație.

(a) Fiecare nod  $k$  primește pachetul  $k$ .(b) Fiecare nod  $k$  primește pachetul  $k + 1$ .Figura 5.2: Replicarea datelor pe  $n$  noduri folosind algoritmul în doi pași.

Folosind aceste informații se poate calcula timpul de transfer pentru această arhitectură. În ecuațiile ce urmează, considerăm  $n$  numărul de noduri destinație,  $F$  fișierul ce trebuie replicat,  $Size(F)$  dimensiunea acestui fișier și  $BW$  lățimea de bandă a legăturilor dintre oricare două noduri (se presupune că rețeaua este omogenă). Timpul total de transfer este suma timpilor execuției celor doi pași descriși mai sus:

$$Timp_{transfer} = Timp_{distributie} + Timp_{colectie} \quad (5.3)$$

Timpul de distribuție este egal cu timpul de transfer al unui subfișier pe unul dintre nodurile destinație:

$$Timp_{distributie} = \frac{Size(F)}{n \times BW} \quad (5.4)$$

Timpul de colectie este de fapt maximum dintre timpul de distribuție și timpul necesar unui nod destinație să trimită subfișierul pe care îl posedă pe toate cele  $n - 2$  conexiuni de ieșire:

$$Timp_{colectie} = \max\left(\frac{Size(F)}{n \times BW}, \frac{(n - 2) \times Size(F)}{n \times BW}\right) = \frac{(n - 2) \times Size(F)}{n \times BW} \quad (5.5)$$

Din ecuațiile 5.3, 5.4 și 5.5 rezultă că timpul total de replicare a fișierului pe toate destinațiile este:

$$Timp_{transfer} = \frac{Size(F)}{n \times BW} + \frac{(n-2) \times Size(F)}{n \times BW} = \frac{(n-1) \times Size(F)}{n \times BW} \quad (5.6)$$

## 5.4 MC-5

Soluția propusă pentru problema replicării unui fișier de dimensiuni mari pe multe noduri destinație a primit denumirea de **MC-5**. În următoarele paragrafe vom explica de ce a fost ales acest nume și cum funcționează această soluție. Vom vedea de asemenea avantajele și dezavantajele existente, plus o comparație cu ceilalți algoritmi deja prezentați.

### 5.4.1 Replicarea pe două noduri

Pentru a putea genera o arhitectură scalabilă pentru un număr mare de destinații, uneori este necesar să se pornească de la un număr mic de destinații. În această idee, prima parte a soluției prezentate aici este replicarea fișierului pe două noduri. Situația poate fi analizată în figura 5.3. Topologia este formată dintr-un nod sursă  $S$  și două noduri destinație  $D_1, D_2$ . Inițial, fișierul ce va fi replicat  $F$  se găsește doar pe nodul sursă. Acesta se conectează la ambele destinații în același timp și le transmite acestora informațiile necesare pentru a se putea realiza transferul:

- numele și calea fișierului  $F$ .
- dimensiunea acestuia, care în calculele ce urmează va fi notată cu  $Size(F)$ .
- dimensiunea unui pachet, deoarece pentru transfer, fișierul va fi trimis în pachete.
- numărul total de pachete.
- adresa IP a celui alt nod destinație (care de aici încolo va fi referit ca fratele nodului curent).
- portul pe care se poate face conectarea la nodul frate (acesta este necesar deoarece se dorește realizarea unei conexiuni bidirecționale între cele două destinații).

În primă fază, nodul sursă are datoria de a căuta informațiile corecte despre fișier: folosind calea fișierului, nodul sursă trebuie să verifice existența fișierului, drepturile de editare/copiere/modificare ale acestuia, cât și dimensiunea totală. Apoi, în funcție de dimensiune și de alți parametri, nodul sursă decide mărimea pachetelor de date și numărul acestora. Alegerea mărimii pachetelor va fi explicată mai detaliat în secțiunea de implementare. După trimiterea pachetelor de configurații către nodurile destinație, acestea vor răspunde nodului sursă cu o confirmare ce arată că au primit configurațiile și că au realizat conexiunea între ele. Acest pas înseamnă declanșarea transferului propriu-zis, deoarece

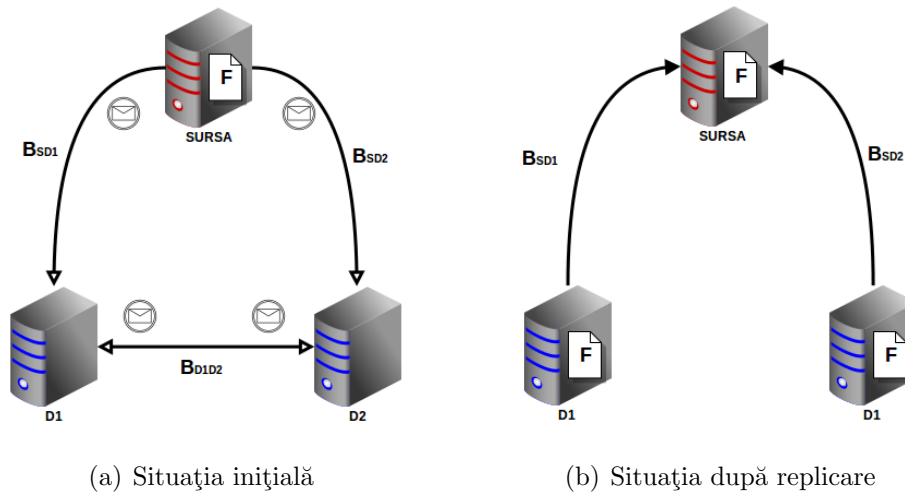


Figura 5.3: Replicarea datelor pe două noduri folosind toate legăturile disponibile.

Începând cu acest moment, nodul sursă va trimite fiecărei destinații câte un pachet, alternativ, începând din ambele capete ale acestuia, până ce se ajunge la un punct comun. Acest pas poate fi văzut în figura 5.3. La primirea unui pachet, fiecare nod destinație aplică o strategie de stocare și înaintare (*en. store-and-forward*) adică memorează pachetul primit și îl trimite mai departe nodului frate.

Folosind această strategie, prezentată în algoritmul 1, se observă că nodurile destinație vor primi date pe toate legăturile existente. Pentru o înțelegere mai bună a acestei strategii, poate fi consultată figura 5.4. Aici se pot vedea pachetele trimise de nodul sursă, cât și cele trimise de nodurile destinație între ele. Folosind acest algoritm, nodul sursă nu are nevoie de confirmarea primirii fiecărui pachet în parte de către destinații. După cum se poate vedea și în algoritmul 1, el va trimite pachetele de date concomitent către destinații începând cu ambele capete ale fișierului. Până nu se ajunge la un punct comun, nodul sursă va tot trimite pachetele, folosind tehnica de împingere a datelor (*en. push*), prezentată anterior. După ce se ajunge la acest punct comun, nodul sursă a trimis fiecare pachet disponibil o singură dată pe oricare dintre legături. Acest aspect este important deoarece se dorește evitarea trimiterii mai multor pachete pe aceeași legătură, la fel ca și primirea aceluiași pachet pe mai multe legături.

Evitarea cererii unor confirmări de primire pentru pachete este avantajoasă din punct de vedere al timpului, aceasta fiind o caracteristică importantă a tehnicii *push*. Trebuie precizat faptul că trimiterea unui pachet de către nodul sursă se referă la acțiunea de inițiere a transferului propriu-zis, și nu la momentul de primire a pachetului. Deci, după ce nodul sursă va termina de trimis primul set de pachete pentru fiecare nod, va exista o diferență de timp până când acestea se vor afla pe serverele destinație (în acest timp, zicem că pachetele se află în tranzit, adică în coada de mesaje TCP). În timp ce nodul sursă trimite aceste pachete, nodurile destinație comunică și ele folosind o legătură

bidirecțională, prin care trimit pachetele primite de la sursă. Folosind această tehnică, fiecare nod va primi într-un final toate pachetele (jumătate de la sursă, jumătate prin intermediul nodului frate). Dar această abordare ar însemna să nu fie utilizate toate legăturile existente la capacitatea lor maximă. De aceea, după ce sursa a trimis prima serie de pachete fiecărui nod, ea cere informații de la destinații ce o ajută să determine ce pachete mai sunt necesare pentru fiecare. Aceasta este tehnica de aducere a datelor (*en. pull*) și este folosită pentru a determina exact ce pachete lipsesc fiecărei destinații.

---

**Algorithm 1** Transferarea unui fișier pe două destinații

---

**Date intrare:**

- $D_1, D_2$  – nodurile destinație.
- $F$  – fișierul ce va fi transferat.
- $packsz$  – dimensiunea unui pachet.

**TRIMITE\_FISIER**( $D_1, D_2, F, packsiz$ )

deschide\_conexiune( $D_1$ )

deschide\_conexiune( $D_2$ )

$n \leftarrow \lceil size(F)/packsz \rceil$

**for**  $i = 1 \rightarrow j, j = n \rightarrow i$  **do**

    trimite( $D_1, F[i]$ )

    trimite( $D_2, F[j]$ )

$i \leftarrow i + 1$

$j \leftarrow j - 1$

**end for**

**while**  $True$  **do**

$ind1 \leftarrow preia\_ultimul\_index(D_1)$

$ind2 \leftarrow preia\_ultimul\_index(D_2)$

**if**  $ind1 \neq ind2$  **then**

        trimite( $D_1, F[i]$ )

$i \leftarrow i + 1$

        trimite( $D_2, F[j]$ )

$j \leftarrow j - 1$

**else**

**return**

**end if**

**end while**

---

Acești pași sunt repetați până când fiecare nod destinație trimite nodului sursă confirmarea că a primit toate pachetele și că fișierul este acum complet ( figura 5.3(b)). Conexiunile sunt oprite, întâi cele dintre destinații, iar după trimiterea confirmării, și cele



dintre sursă și destinații. Fișierul  $F$  se află acum pe toate nodurile implicate în transfer.

Algoritmul descris mai sus este simplu, dar puternic. Să analizăm câteva caracteristici ale acestuia, comparativ cu abordarea clasică de transfer serial:

- are un număr mai mare de conexiuni, deoarece există conexiuni și între destinații, dar acestea sunt utilizate la capacitatea lor maximă – sunt folosite mereu pentru a transfera informații.
- la fel ca și în abordarea serială, replicarea datelor pe un nod nu depinde de eșecul sau succesul fratelui acestuia – datorită modului în care se transferă fișierul, există certitudinea că toate pachetele de date vor ajunge la destinații, pe una dintre legăturile existente.
- spre deosebire de abordarea serială, aici sunt utilizate mai multe conexiuni, ceea ce se înseamnă fiabilitate ridicată, dar și o cantitate crescută de informații de configurare ce trebuie trimise de la sursă.
- timpul de transfer este evident mai mic decât în cazul unei soluții în care fișierul este transmis pe rând fiecărei destinații.
- comunicarea minimizată dintre noduri este un avantaj pentru minimizarea timpului de replicare.
- în cazul în care nu se poate realiza conexiunea între destinații, acest algoritm se va reduce la o schemă de distribuție serială.

Dacă realizăm o analiză a acestui algoritm, se observă că acesta seamănă cu o versiune cu două noduri a algoritmului în doi pași descris mai sus. Există totuși câteva diferențe importante:

- în această soluție fișierul nu este divizat în două subfișiere, ci în multiple pachete de dimensiuni considerabil mai mici decât ale subfișierelor. Aceasta permite realizarea transferului în mai mult de doi pași.
- deși numărul și tipul conexiunilor este același ca și la algoritmul în doi pași, aici felul în care sunt trimise datele diferă prin faptul că nodurile destinație au rolul de a retrimite orice pachete primesc de la sursă. Astfel, chiar dacă legătura dintre sursă și una dintre destinații nu este funcțională, fișierul poate fi trimis folosindu-se celălalt nod destinație ca intermediar.
- numărul relativ mare de pachete trimise în această abordare permite nodurilor să se adapteze la viteza legăturilor din rețea: pe legăturile mai bune vor fi transferate mai multe pachete.
- în situația în care nodul sursă decide divizarea fișierului în două pachete, cei doi algoritmi devin identici.

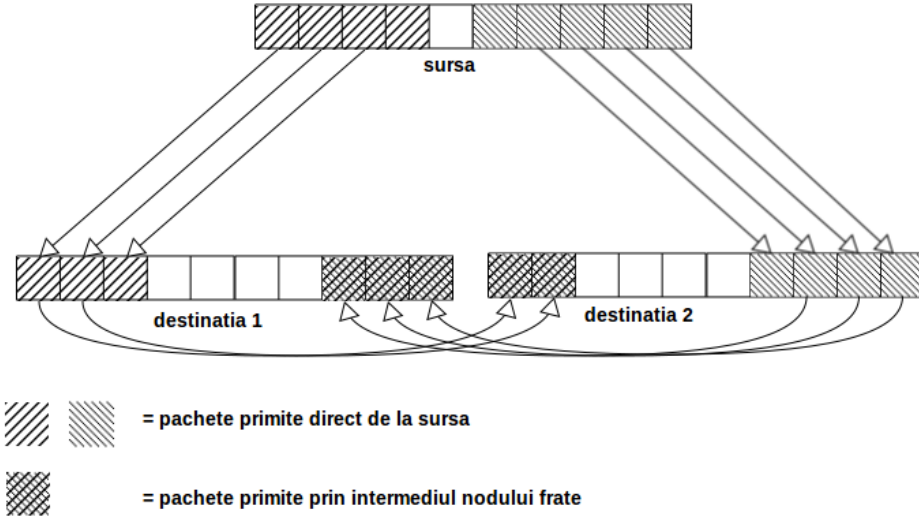


Figura 5.4: Divizarea documentului și trimiterea datelor pe toate căile posibile în arhitectura cu două noduri destinație.

Pentru calcularea timpului total de replicare în această situație, vom considera o rețea omogenă, cu lățimea de bandă  $BW$ .  $T_{D_1}$  și  $T_{D_2}$  sunt timpii de transfer pentru nodurile  $D_1$  respectiv  $D_2$ . Aceștia pot fi exprimați ca produsul dintre timpul de transfer pentru un pachet de date și numărul total de pachete trimise pe acea legătură. Dimensiunea unui pachet este raportul între dimensiunea totală a fișierului și numărul de pachete  $k$ .

$$T_{imp_{replicare}} = \max(T_{D_1}, T_{D_2}) = \frac{k}{2} \times \frac{\frac{Size(F)}{k}}{BW} = \frac{Size(F)}{2 \times BW} \quad (5.7)$$

#### 5.4.2 Replicarea pe seturi reduse de noduri

Pornind de la algoritmul descris anterior și de la o arhitectură minimală, cu doar două noduri, se poate generaliza și se poate găsi o arhitectură ce reușește să se scaleze pe un număr mare de noduri. Abordarea propusă în această lucrare pornește de la această idee, iar modalitatea de generare a acesteia este ilustrată în figura 5.5.

Ideea de bază a acestei abordări de genul împarte și cucerește (*lat. divide et impera*) este de a nu privi toate nodurile ca un mare set, ci de a le restructura în seturi mai mici, cu minim două și maxim cinci noduri fiecare. Aceste seturi se numesc grupuri de replicare și pot fi de patru feluri ( $G_1, G_2, G_3, G_4$ ), după cum se vede și în figura 5.5. Un grup este format din trei tipuri de noduri:

- nod sursă ( $S$ ) – acesta este unic și comun tuturor grupurilor de replicare.

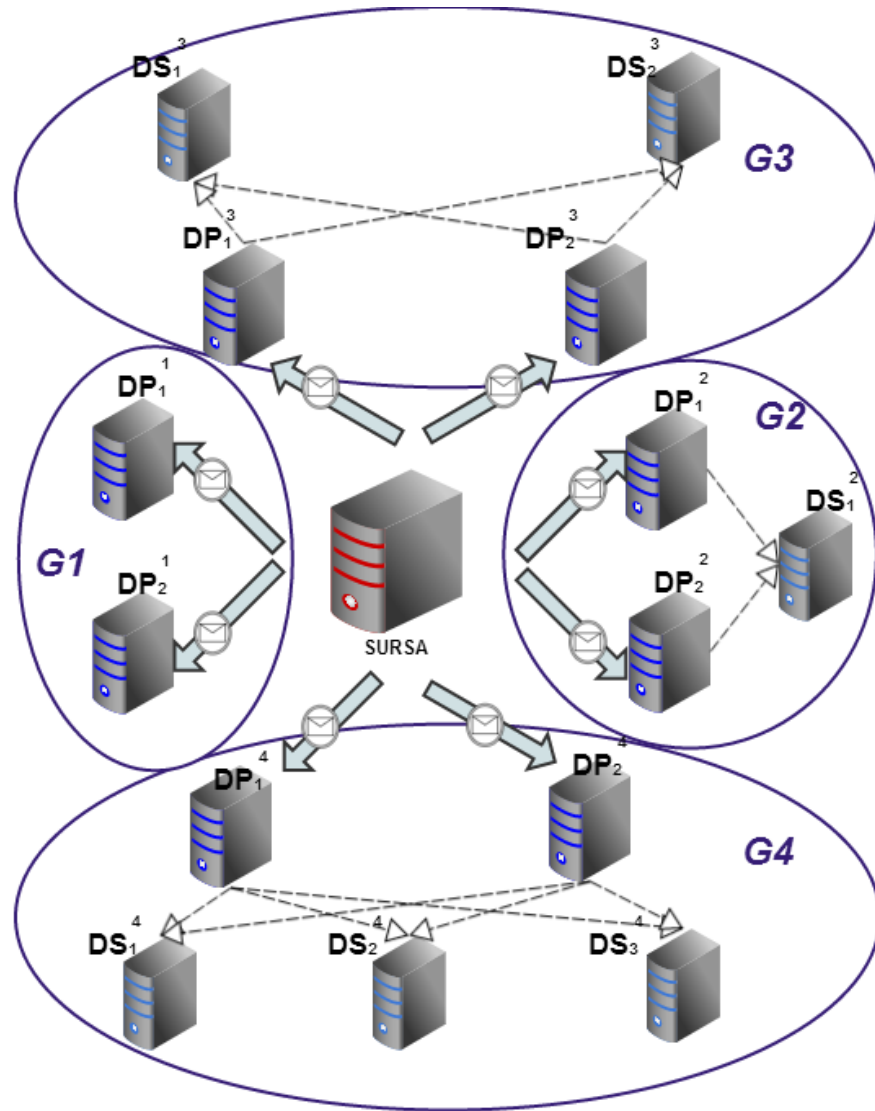


Figura 5.5: Arhitectura pentru replicarea datelor pe grupuri formate din maxim 5 noduri.

- noduri destinație principală ( $DP_i$ ) – două la număr, aceste noduri formează, împreună cu sursa, o structură de genul celei prezentate în subcapitolul anterior. Acestea se conectează direct la sursă, cât și între ele, și au rolul de a transmite mai departe pachetele primite.
- noduri destinație secundară ( $DS_i$ ) – aceste noduri pot fi în număr de zero, unu, doi sau trei și se conectează la destinațiile principale pentru a primi de la acestea fișierul. Ele nu comunică direct cu sursa și nu trimit pachetele primite mai departe.

În funcție de numărul total de destinații, nodul sursă decide schema de distribuție,

$n \pmod{5}$	Schema de distribuție		Conexiuni sursă
0		$k \times G_4$	$2 \times k$
1	$k = 0$	unicast	1
	$k > 0$	$2 \times G_2 + (k - 1) \times G_4$	$2 \times (k + 1)$
2	$k = 0$	$G_1$	2
	$k > 0$	$G_2 + G_3 + (k - 1) \times G_4$	$2 \times (k + 1)$
3	$k = 0$	$G_2$	2
	$k > 0$	$G_2 + (k - 1) \times G_4$	$2 \times (k + 1)$
4	$k = 0$	$G_3$	2
	$k > 0$	$G_3 + (k - 1) \times G_4$	$2 \times (k + 1)$

Tabela 5.1: Generarea schemei de distribuție din numărul total de destinații

adică ce fel de grupuri să genereze. În tabelul 5.1 se pot vedea schemele generate și numărul de conexiuni totale pentru nodul sursă pe baza acestor scheme. Aici,  $n$  reprezintă numărul de destinații, iar  $G_1, G_2, G_3, G_4$  tipurile de grupuri de replicare.

Din acest tabel se poate observa că numărul maxim de conexiuni pentru nodul sursă este  $2 \times (k + 1)$ , unde  $k = n \div 5$ . De aici rezultă maxim  $\frac{2 \times n}{5} + 2$  conexiuni folosind acest sistem. Acest rezultat face din topologia propusă una scalabilă pentru un număr mare de destinații. În algoritmul 2 sunt descrise acțiunile luate de un nod destinație la primirea unui pachet. Deoarece este posibil și chiar foarte probabil ca nodurile destinație să nu primească pachetele de date în ordinea naturală  $(1, 2, \dots, n)$  acestea nu le pot scrie direct în fișier. Rezultă ca va fi necesară o zonă de memorie pentru pachetele care nu sunt încă în ordine. Cea mai ușoară abordare, în acest caz, este de a memora toate pachetele primite într-o zonă tampon (*en. buffer*) și de a le scrie în fișier după ce sunt primite toate. Totuși, din cauza limitărilor de memorie, această abordare nu este mereu posibilă. Din această cauză, o soluție alternativă trebuie găsită. Pentru aceasta, algoritmul 2 a fost creat, pornind de la ideea de unei ferestre glisante (*en. sliding window*) de memorie. Fiecare nod sursă are un buffer de memorie de dimensiune  $r = \text{Size}(W)$  – aceasta este denumită dimensiunea ferestrei și este comunicată de către nodul sursă la începutul fiecărui transfer, deoarece depinde de mărimea fișierului. Deoarece nu se dorește memorarea tuturor pachetelor primite, ci doar a celor care nu pot fi scrise în fișier chiar la primire,  $r$  va fi mai mic decât numărul total de pachete. Să analizăm acum algoritmul pe pași:

- atunci când un pachet nou este primit, prima sarcină a nodului sursă este să verifice dacă aceasta poate fi scris direct în fișier.
- dacă da, nodul sursă trebuie să parcurgă fereastra și să găsească și pachetele următoare ce pot fi scrise.
- dacă nu, se caută un loc liber în fereastră, denumit slot, iar pachetul este memorat acolo până când va putea fi scris în fișier.

- în cazul în care nu mai există locuri libere în fereastră, pachetul cu cel mai mare index (deci cel care va fi scris cel mai târziu în fișier) va fi aruncat, iar noul pachet îi va lua locul.

---

**Algorithm 2** Procesarea unui pachet pe nodurile destinație.

---

**Date intrare:**

- $P$  – pachetul primit.
- $W$  – fereastra de memorie.
- $F$  – fișierul în care se vor scrie datele.

**PROCESARE\_PACHET**( $P, W, F$ )

```

 $urm \leftarrow \text{INDICE\_URM}(F)$ 
 $curent \leftarrow \text{INDICE\_PACHET}(P)$ 
if  $urm = curent$  then
  SCRIE( $F, P$ )
  while  $True$  do
     $urm \leftarrow \text{INDICE\_URM}(F)$ 
     $curent \leftarrow \text{INDICE\_MINIM}(W)$ 
    if  $urm = curent$  then
      SCRIE( $F, W[curent]$ )
    else
      break
    end if
  end while
else
  if ESTE_PLIN( $W$ ) then
     $slot \leftarrow \text{INDICE\_MAXIM}(W)$ 
  else
     $slot \leftarrow \text{GASESTE\_LIBER}(W)$ 
  end if
   $W[slot] \leftarrow P$ 
end if

```

---

### 5.4.3 Analiza soluției propuse

După detalierea arhitecturii propuse și a algoritmilor de transfer, vom vedea în continuare o analiză a acestei soluții din mai multe puncte de vedere:

#### Scalabilitate

Această caracteristică se referă la posibilitatea de a realiza topologii cu număr foarte mare de noduri, fără a cauza probleme de performanță. Deoarece nodurile destinație

sunt divizate în grupuri de replicare după cum a fost descris, numărul total maxim de conexiuni pentru fiecare dintre acestea este șase (pentru destinațiile principale, două conexiuni de intrare și patru de ieșire).

Problema se pune la nodul sursă care are două legături cu fiecare grup de replicare. Deci, pentru un număr de  $A$  grupuri, sursa va avea  $2 \times A$  conexiuni. Numărul de grupuri este dat de numărul de destinații împărțit la 5, deci pentru 100 destinații sursa va avea 40 conexiuni. Putem zice că soluția propusă este scalabilă, dar cu o anumită limită.

### **Adaptabilitate**

Acest punct se referă la modificarea anumitor parametrii de transfer sau adaptarea anumitor caracteristici pentru a obține o performanță ridicată. În cazul de față, trebuie remarcat că adaptabilitatea este realizată prin utilizarea tuturor legăturilor disponibile între noduri. Acestea sunt în permanență ocupate cu pachete de configurare sau de date, optimizând foarte mult transferul. De asemenea, se observă că pe legăturile mai bune, adică mai rapide, se vor trimite mai multe pachete, deoarece pe acestea vor veni mai multe cereri.

### **Fiabilitate**

Spunem despre un sistem că este fiabil sau de încredere dacă acesta reușește să îndeplinească scopul pentru care a fost proiectat de fiecare dată, indiferent de diversele medii în care e executat. Aici, fiabilitatea reprezintă garanția că întreg fișierul va ajunge pe nodurile destinație, chiar dacă anumite pachete sunt pierdute sau deteriorate.

Acest fapt poate fi demonstrat prin algoritmi prezentați anterior. Aceștia asigură trimiterea pachetelor de date prin faptul că sursa cere mereu confirmări de la destinații pentru primirea tuturor pachetelor. În cazul în care primirea nu e confirmată, se execută o retransmitere a pachetelor pierdute sau deteriorate.

### **Toleranță la erori**

Uneori, înainte sau chiar în timpul execuției, pot apărea erori pe serverele destinație. Toleranța la erori înseamnă că defectarea unui server nu ar trebui să afecteze în nici un fel replicarea fișierului pe celelalte servere. Să analizăm această problemă în contextul sistemului prezentat aici:

- în cazul în care serverul pe care apare eroarea este o destinație principală, destinațiile secundare care ar trebui să primească date de la acesta vor primi totuși pachetele prin intermediul fratelui acestuia, deci problema este rezolvată.
- în cazul în care serverul pe care apare eroarea este o destinație secundară, acesta nu este responsabil cu retransmiterea pachetelor către alte servere, deci nici aici nu se poate vorbi de propagarea erorilor.

În concluzie, folosind abordarea propusă, cu grupuri mici de replicare, putem fi siguri că fișierul va fi replicat pe toate serverele funcționale.

**Eficiență memorie**

Deoarece dorim să transferăm volume mari de date, managementul memoriei este o caracteristică necesară pentru acest sistem. Pentru nodul sursă, ideal este să nu se folosească zone inutile de memorie pentru a ține pachetele de date: aceasta se realizează prin citirea directă a datelor din fișier. Pentru nodurile destinație, este important să nu se stocheze pachetele primite în memoria internă a proceselor, deoarece aceasta este limitată și poate ridica excepții ce duc la terminarea prematură a transferului. De aceea, algoritmul 1 a fost creat, pentru a ține sub controlul nivelul memoriei folosite la un moment dat.

**Eficiență timp**

Scopul general al acestei lucrări este de a avea un sistem performant din punct de vedere al timpului de transfer, adică de a încerca minimizarea acestuia. După cum am văzut, pentru fiecare grup de replicare, sursa este considerată ca parte integrantă a acestuia, deoarece ea trimite pachetele concurent tuturor grupurilor. Această caracteristică este o importantă creștere a eficienței, la fel ca și folosirea tuturor legăturilor disponibile dintre noduri.

# Capitolul 6

## Proiectare de detaliu și implementare

### 6.1 Mediul de dezvoltare

Pentru implementarea aplicației din lucrarea de față a fost ales ca limbaj de programare Python, iar ca platformă de dezvoltare Twisted. Acestea sunt prezentate în secțiunile ce urmează.

#### 6.1.1 Python

**Python** este un limbaj de programare de nivel înalt apărut în anul 1991. Python oferă multiple paradigme de programare, precum programare orientată obiect, programare imperativă și funcțională. Acesta dispune de un sistem dinamic pentru tipurile de date și de management automat al memoriei, două puncte forte pentru aplicația prezentă. Este disponibil pentru toate marile tipuri de sisteme de operare și are licență *open source*. Permite atât interpretarea codului (rulare direct folosind fișierele sursă) cât și compilarea acestuia pentru a crește viteza de execuție. De asemenea, oferă posibilitatea de a structura codul și de a-l arhiva în module și librării simplu de utilizat.

Python a fost conceput cu o sintaxă ce permite exprimarea conceptelor în mai puține linii de cod decât în alte limbaje și cu un design ce conferă lizibilitate codului. Accentul este pus pe construirea programelor din punct de vedere funcțional prin evitarea scrierii de cod inutil. Acestea sunt câteva dintre motivele pentru care Python a fost ales pentru implementarea sistemului prezentat în lucrarea de față.

#### 6.1.2 Twisted

Twisted este o platformă pentru crearea de aplicații de rețea scrisă în Python și licențiată *open-source*. Aceasta este disponibilă la adresa <http://twistedmatrix.com>. Twisted este bazat pe paradigma de programare acționată de evenimente (*en. event-driven*). În figura 6.1 este prezentată structura unei aplicații realizate cu Twisted: este



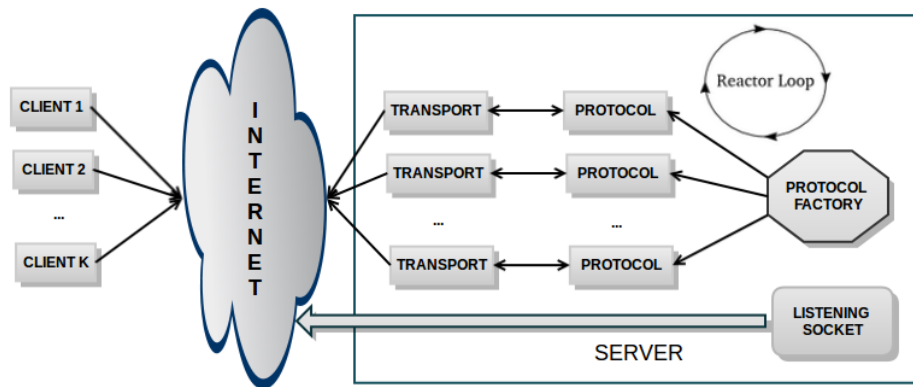


Figura 6.1: Funcționarea unui server implementat cu Twisted.

vorba despre un server ce așteaptă conexiuni de la clienți prin intermediul Internetului. În această imagine sunt câteva dintre cele mai importante elemente Twisted:

### Reactor

Această componentă este implementată conform modelului de proiectare reactor (*en. reactor design pattern*): o buclă care așteaptă apariția unor evenimente, pe care apoi le gestionează. Acest model este util pentru manevrarea cererilor concurente de la mai mulți clienți.

### Transport

Această componentă reprezintă o singură conexiune care poate trimite sau primi octeți. După cum sugerează și numele, aceasta reprezintă o abstractizare a transportului de date. Implementarea oferită de Twisted este non-blocantă și este recomandat să se folosească în locul unei implementări proprii.

### Protocol

Această componentă poate fi suprascrisă pentru a implementa un protocol specific de rețea. Fiecare instanță Protocol este utilizată pentru gestionarea unei singure conexiuni, deci aici vor fi memorate starea protocolului cât și datele parțial primite. Fiecărei instanțe Protocol îi corespunde o instanță Transport, care reprezintă conexiunea pe care protocolul acționează. Twisted oferă implementări pentru cele mai comune și mai des utilizate protocoale.

### Factory

Crearea instanțelor unui Protocol este datorată componentei Factory. Aceasta este implementată folosind modelul de proiectare fabrică (*en. factory design pattern*). De fiecare dată când o nouă cerere de conexiune apare, este de datorită acestei componente să instanțieze un nou Protocol.

## 6.2 Implementare



Figura 6.2: Structura aplicației pe module.

După cum a fost precizat și anterior, implementarea a fost făcută în Python, folosind platforma Twisted. Structura aplicației poate fi analizată în imaginea 6.2, din stânga. Aceasta este modularizată și structurată în două module ce sunt detaliate mai jos.

Când un modul este importat, interpretorul Python va îl va căuta în căile predefinite, și când îl va găsi, îl va adăuga în tabela locală de simboluri. Din motive de eficiență, fiecare module este importat o singură dată per sesiune. Ca o importantă accelerarea a timpului de pornire pentru programele scrise în Python care utilizează multe module, codul pre-compilat este utilizat (*en. byte-compiled code*). Modalitățile de lansare în execuție a clientului și a serverului vor fi descrise mai detaliat în manualul de instalare și utilizare, alături de resursele necesare.

Fiecare modul conține fișierul `__init__.py`, care este necesar pentru ca Python să trateze folderele în care se află ca și module; această abordare a fost aleasă pentru a preveni folderele cu nume comune de la a fi confundate cu adevăratele module Python (*e.g. string*). În cel mai simplu caz, `__init__.py` poate fi un doar un fișier gol, dar poate fi și folosit pentru a executa cod de inițializare pentru pachetul respectiv. În structura aceasta, fișierele `__init__.py` sunt goale, deci nu vor mai fi prezentate în paragrafele următoare.

### 6.2.1 Implementare client

Primul modul din imaginea 6.2 este denumit **client** și este format din următoarele fișiere sursă:

#### **mainClient.py**

Acesta este punctul de începere a execuției unui client, fișierul ce conține metoda principală de lansare (*en. main*). Deoarece lansarea se face din terminal primul pas este de a analiza opțiunile date ca argumente în linia de comandă:

**'-sa'/'-source-address'** Adresa IP a serverului sursă. Dacă nu este precizată, se va folosi adresa de buclă locală 127.0.0.1.

- '-sf'/'-source-filepath'\* Calea fișierului pentru serverul sursă. Acest argument este necesar.
- '-mc'/'-multicast'\* Calea fișierului de multicast în care sunt scrise, câte una pe linie, adresele IP ale serverlor destinație. Acestea trebuie să fie de forma  $X.X.X.X : [P]$ , unde  $X \in [0, 255]$  și  $P \in [20000, 25000]$ . Acest argument este necesar.
- '-sp'/'-source-port' Portul pe care se realizează conectarea la serverul sursă. Valoarea implicită este 22020.
- '-df'/'-dest-filepath'\* Calea fișierului pentru serverul sursă. Acest argument este necesar.
- '-dp'/'-dest-port' Portul pe care se va realiza conectarea la serverele destinație în cazul în care acesta nu este precizat în fișierul de multicast. Valoarea implicită este 22000.
- '-ps'/'-packet-size' Dimensiunea unui pachet exprimată în octeți. Valoarea implicită este de 67108864 octeți(64MB).
- '-ws'/'-window-size' Dimensiunea ferestrei glisante exprimată în număr de sloturi. Valoarea implicită este de 10 sloturi.
- '-lf'/'-logfile' Calea fișierului jurnal (*en. logfile*) în care se scriu diverse informații din timpul transferului și eventuale erori. Valoarea implicită este `/tmp/log-client.log`

După parsarea argumentelor din linia de comandă, se deschide fișierul de multicast pentru a se citi adresele serverlor destinație. Următorul pas este inițializarea protocolului de comunicare între client și sursă, prin pornirea reactorului și deschiderea canalelor de comunicație.

#### **clientFactory.py**

MyClientFactory este o clasă ce extinde ClientFactory, implementat în Twisted, și care are rolul de a instanția protocolul de comunicare.

#### **clientProtocol.py**

MyClientProtocol este clasa ce extinde LineReceiver și care conține protocolul de comunicare între client și serverul sursă. Metodele principale din această clasă sunt:

**connectionMade** apelată atunci când conexiunea s-a realizat.

**lineReceived** apelată la primirea unei linii de la celălalt capăt al comunicației.

### **6.2.2 Implementare server**

Pentru servere am ales o implementare separată pentru sursă, respectiv destinație. Fiecare din aceste instanțe ascultă pe un alt port, deci în cazul unui transfer nu mai este

necesară pornirea unei noi instanțe, făcând aplicația mai eficientă din punct de vedere al timpului. Modulul **server** este format din două submodule și fișierul de lansare în execuție a unui server – **mainServer.py**. Și aici, primul pas este parsarea argumentelor din linia de comandă:

'-sp'/'-source-port' Portul pe care va asculta serverul sursă. Valoarea implicită este 22020.

'-dp'/'-dest-port' Portul pe care va asculta serverul destinație. Valoarea implicită este 22000.

'-lf'/'-logfile' Calea fișierului jurnal (*en. logfile*) în care se scriu diverse informații din timpul transferului și eventuale erori. Valoarea implicită este `/tmp/log/server.log`

După parsarea argumentelor, sunt pornite cele două instanțe de protocol și este pornit reactorul. Din acest moment, serverul așteaptă conexiuni de la clienți.

## Sursă

Componentele pentru implementarea serverului sursă pot fi văzute în modulul **server.sender**. Acestea sunt:

### **senderFactory.py**

MySenderFactory extinde ServerFactory, și are rolul de a crea instanțe MySender-Protocol.

### **senderProtocol.py**

MySenderProtocol extinde LineReceiver și conține protocolul de comunicare între client și serverul sursă. Metodele importante din această clasă, pe lângă cele moștenite de la Protocol, sunt:

**\_incFileDone** apelată de instanța de trimitere a fișierului atunci când destinația cu care comunica a primit toate pachetele. Utilizată pentru a contoriza numărul de destinații care au primit întreg fișierul, pentru a putea notifica clientul la sfârșitul transferului.

**\_getPacketById** apelată de instanța de trimitere a fișierului pentru a primi un nou pachet. În această metodă se citesc **packetSize** octeți din fișier și se creează un pachet de date, folosind formatul din figura 6.3. Aici, antetul este format din două câmpuri: **zero**, care reprezintă identificatorul unui pachet de date și este format din patru octeți cu valoarea 0; și **seqNr**, care reprezintă numărul de secvență al pachetului, adică pe ce poziție va fi scris în fișier.

**\_getConfigurations** apelată de ambele destinații principale; doar prima dintre ele va primi configurațiile, pe care le va trimite mai apoi nodului frate și nodurilor fii.

**fileSenderFactory.py**

MyFileSenderFactory extinde ServerFactory, și are rolul de a crea instanțe MyFileSenderProtocol.

**fileSenderProtocol.py**

MyFileSenderProtocol extinde LineReceiver și conține protocolul de comunicare între serverul sursă și un server destinație (acesta este descris în secțiunea de protocoale).

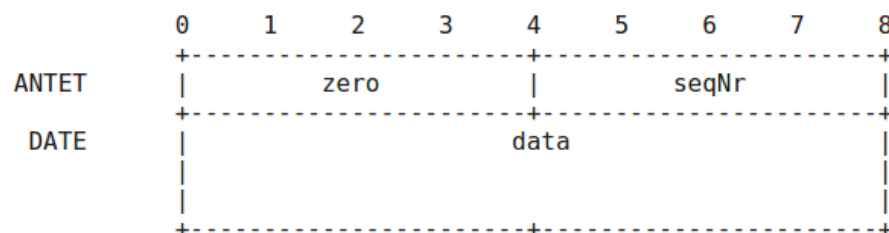


Figura 6.3: Structura unui pachet de date.

**Destinație**

Fișierele sursă ce formează implementarea serverului destinație se află în modulul **server.receiver**. Acestea sunt:

**receiverFactory.py**

MyReceiverFactory extinde ServerFactory, și are rolul de a crea instanțe MyReceiverProtocol.

**receiverProtocol.py**

MyReceiverProtocol extinde LineReceiver și conține protocolul de comunicare între serverul sursă și serverul destinație principală. În această clasă este implementată procesarea pachetelor de date primite folosind algoritmul 2. Metode importante:

**\_newConnection** metodă apelată la realizarea unei conexiuni și folosită pentru a contoriza toate conexiunile existente la un moment dat. Aceasta este necesară deoarece destinația principală are rolul de a trimite mai departe pachetele de date pe care le primește de la sursă, deci conexiunile trebuie realizate înainte de primirea datelor.

**\_processData** această metodă este apelată când se detectează primirea unui pachet de date. Aici se realizează despachetarea acestuia și se decide dacă pachetul va fi scris direct în fișier sau va fi stocat în fereastra de memorie. Tot aici se hotărăște ce pachete sunt trimise nodurilor fii.

**multiEchoFactory.py**

MultiEchoFactory extinde Factory, și este clasa ce creează instanțe MultiEchoProtocol, câte una pentru fiecare destinație secundară pe care se dorește retransmiterea pachetelor primite de la sursă.

**multiEchoProtocol.py**

MultiEchoProtocol extinde LineReceiver și este folosit pentru a realiza conexiunile cu destinațiile secundare. Aici, o mare importanță o au metodele moștenite de la LineReceiver:

**connectionMade** la stabilirea unei conexiuni, aceasta este adăugată în lista de noduri fii.

**connectionLost** la pierderea conexiunii, aceasta este scoasă din listă.

## 6.3 Protocoale de comunicare

În figura 6.4 se poate vedea diagrama de clase a aplicației în raport cu clasele utilizate din Twisted. După cum se poate observa, toate protocoalele de comunicare moștenesc clasa LineReceiver. Acesta oferă posibilitatea de a trimite linii sau date neprelucrate, care sunt terminate de un delimitator. Lungimea maximă a unei linii ce poate fi trimisă este de 16384 octeți, dar aceasta poate fi modificată pentru a se mapa necesităților aplicației. Principalele metode implementate de această clasă:

**lineReceived** apelată la primirea unei linii întregi.

**rawDataReceived** apelată la primirea unui set de date neprelucrate.

**sendLine** apelată pentru a trimite o linie la celălalt capăt al conexiunii.

### 6.3.1 Client - Server

Deoarece în această aplicație este vorba de un transfer *third-party*, în primul rând vom avea comunicare între client și servere. Datorită numărului mare de servere destinație, clientul va comunica doar cu sursa, trimițându-i acesteia informațiile necesare pentru transfer. Deoarece nu există transfer de date între client și sursă, acestea vor comunica doar prin comenzi, după cum este descris mai jos. Trebuie remarcat faptul că nu este exclus ca clientul și sursa să se afle pe aceeași mașină fizică. Mai jos sunt descrise comenzile care stau la baza acestei comunicații. Comenzile ce au săgeata spre dreapta ( $\rightarrow$ ) sunt trimise de la client la sursă, iar cele ce au săgeata spre stânga ( $\leftarrow$ ) de la sursă la client. Descrierea este făcută din perspectiva clientului, deoarece el este inițiatorul conversației:

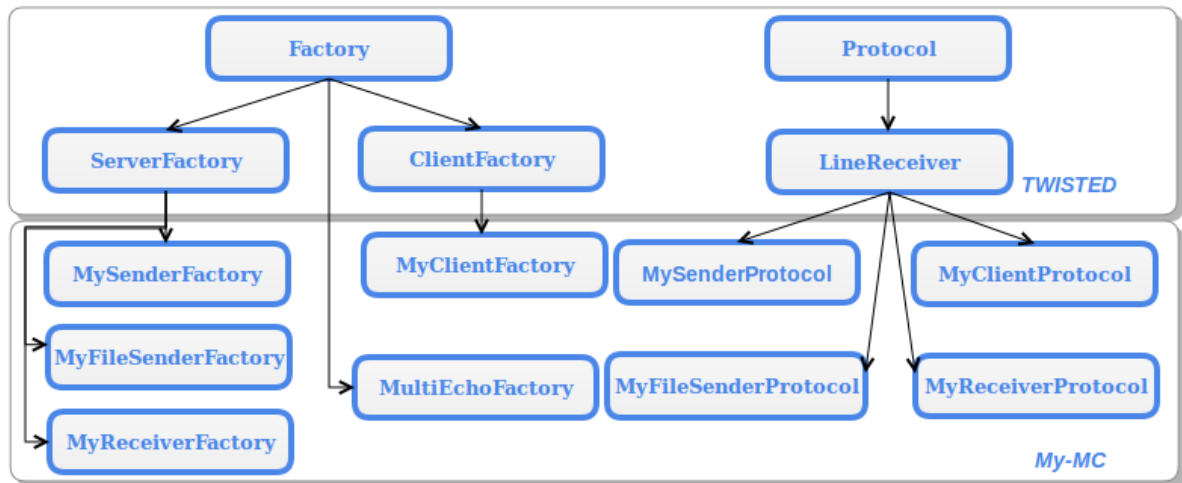


Figura 6.4: Diagrama de clase a aplicației.

**SRFI cale\_fisier\_sursa** → trimite calea fișierului de pe serverul sursă. La primire, sursa trebuie să verifice existența fișierului și drepturile acestuia de citire, copiere și modificare.

**DSFI cale\_fisier\_destinatie** → trimite calea fișierului pentru serverele destinație. Acestea vor trebui să verifice drepturile de scriere pentru calea respectivă.

← **FISZ dimensiune\_fisier** primește dimensiunea fișierului în octeți, ceea ce înseamnă că fișierul a fost găsit pe serverul sursă, iar drepturile acestuia permit citirea datelor.

**NRDS numar\_destinatii** → trimite numărul total de destinații către sursă.

**PKSZ dimensiune\_pachet** → trimite dimensiunea unui pachet în octeți.

**WDSZ dimensiune\_fereastră** → trimite dimensiunea ferestrei în număr de sloturi.

← **DSOK** primește confirmarea de la sursă că s-a alocat suficient spațiu pentru a memora toate destinațiile, pentru o ulterioară conectare la acestea.

**DSIP adresa\_IP port** → trimite adresa IP și portul unei destinații. Această comandă este trimisă pentru fiecare destinație în parte.

← **TRST** primește confirmarea începerii transferului. În acest moment, sursa a reușit să se conecteze la destinații și urmează să trimită fișierul.

- ← **STAT** primește confirmarea terminării transferului. După acest moment, clientul se va deconecta de la serverul sursă.
- ← **ERRR** această linie este primită doar în cazul apariției unei erori pe server: fișierul nu este găsit, comanda primită de server a fost alterată sau eronată, fișierul nu a putut fi deschis pentru citire sau a apărut o eroare internă.

### 6.3.2 Server - Server

După primirea instrucțiunilor de transfer de la client, serverul sursă devine el însuși client, conectându-se la serverele destinație principală, pentru a începe transferul. Protocolul folosit în această sesiune de comunicare este diferit de cel cu clientul și este descris în paragraful ce urmează. Aici, comenzile cu săgeata orientată spre dreapta (→) sunt trimise de la sursă la destinația principală, iar cele cu săgeata orientată spre stânga (←) de la destinație la sursă. Narațiunea este făcută din perspectiva sursei, deoarece ea este inițiatoarea conversației. De remarcat că sursa comunică doar cu destinațiile principale:

**SCHM numar\_schema** → trimite numărul schemei de distribuție către destinație. Această informație este utilă pentru ca destinația să știe câte noduri fii va avea.

**SIBL adresa\_IP port** → trimite adresa IP și portul destinației frate.

**CHLD adresa\_IP port** → trimite adresa IP și portul unei destinații cu rol de nod fiu; această comandă va fi trimisă pentru fiecare nod fiu în parte, în funcție de schema de distribuție.

- ← **CONN** primește confirmarea de conectare la nodul sursă și la nodurile fii. Acest pas înseamnă că serverul sursă poate să înceapă să trimită comenzile de configurare a transferului.

**FPTH cale\_fisier** → trimite calea fișierului ce va fi trimis. Această informație este trimisă mai departe către nodurile fii și către nodul frate.

**PACK nr\_pachete dimens\_pachet** → trimite numărul de pachete și dimensiunea unui pachet de date. Această informație face parte din setul de configurații, deci va fi trimisă mai departe.

**WDSZ dimensiune\_fereastră** → trimite dimensiunea ferestrei de memorie. Și această comandă este trimisă nodului frate și nodurilor fii, aceasta fiind ultima configurație necesară pentru începerea transferului.



- ← **SEND** primește confirmarea procesării setărilor de configurație, ceea ce înseamnă ca transferul poate începe. Acum nodul sursă începe să trimită pachete de date folosind tehnicile descrise în algoritmi.
- REQC** → trimite o cerere de confirmare de primire a tuturor pachetelor de date. Destinația poate să răspundă acestei cereri cu una dintre următoarele comenzi.
- ← **NREC id\_pachet** primește de la destinație o cerere pentru pachetul cu numărul de secvență `id_pachet`, indexul primului pachet lipsă de pe acest server.
- ← **RECA** primește de la destinație confirmarea primirii tuturor pachetelor de date. Fișierul este complet.

## 6.4 Calibrare parametrii

Pentru a putea vorbi de eficiență, trebuie stabilit nivelul la care se face retransmiterea datelor și care va determina gradul de paralelism în comunicarea dintre diferite noduri. De această granularitate depinde adaptabilitatea protocolului prezentat aici. Cu cât e mai mare dimensiunea unei unități de transfer, cu atât mai mică este și adaptabilitatea. Totuși, cu cât e mai mică unitatea respectivă, cu atât mai mult se complică mecanismul de evidență a pachetelor și operația de transfer. Ideea de bază pentru a rezolva această problemă este că dimensiunea unui pachet de date poate fi crescută/redușă în funcție de anumiți factori.

În implementarea propusă există posibilitatea de a lansa un transfer personalizat, prin setarea parametrilor din linia de comandă. Este vorba despre dimensiunea pachetelor și dimensiunea ferestrei de memorie. Acești doi parametri influențează timpul de transfer total, iar alegerea unor valori potrivite pentru aceștia este esențială atunci când se dorește performanță.

În acest scop, o serie de transferuri de test au fost realizate. Acestea sunt detaliate în tabelele de mai jos.

În tabelul 6.1 sunt prezentate transferurile realizate pentru determinarea dimensiunii unui pachet. Pentru acestea s-a folosit un fișier de 1GB și două noduri destinație. După cum se poate vedea în tabel, cel mai bun rezultat a fost obținut pentru valoarea de 1.75MB. Totuși, trebuie ținut cont de faptul că sunt mai mulți factori ce influențează evoluția timpului de transfer.

Următorul pas pentru găsirea unei valori corecte a fost varierea dimensiunii fișierului (aceasta a luat valori între 50MB și 1.6GB) și a numărului de noduri destinație (cu valori între două și opt destinații). După realizarea acestor teste s-a ajuns la concluzia că cei mai buni timpi de transfer sunt dependenți de dimensiunea pachetelor raportat la dimensiunea totală a fișierului, adică la numărul de pachete. S-a constatat că rezultatele optime se

Nr. rulare	Dimensiune pachet	Număr pachete	Timp de transfer
1	10.0 MB	100	4m 26s
2	5.0 MB	200	2m 56s
3	2.5 MB	400	2m 16s
4	1.75 MB	800	2m 9s
5	1.0 M	1000	2m 18s
6	1.5 M	666	2m 18s
7	2.0 M	500	2m 16s

Tabela 6.1: Determinarea dimensiunii optime a unui pachet de date pentru un fișier de 1000 MB.

Nr. rulare	Dimensiune fereastră	Număr pachete	Timp de transfer
1	200	800	2m 9s
2	300	800	2m 29s
3	100	800	2m 26s
4	150	800	2m 7s
5	250	800	2m 20s
6	220	800	2m 12s
7	210	800	2m 10s

Tabela 6.2: Determinarea dimensiunii optime a ferestrei de memorie pentru pachete de 1.75 MB și fișier de 1000 MB.

obțin atunci când fișierul este divizat în aproximativ 350-450 pachete de date, indiferent de dimensiunea acestora.

Desigur, pentru fișiere foarte mici, nu este recomandat să se folosească această tactică, deoarece se poate ajunge la pachete de dimensiuni mult prea mici. Deși timpul de transfer în cazul acestor fișiere este oricum destul de scăzut, este recomandabil ca pentru fișierele care au mai puțin de 50MB să nu se coboare dimensiunea unui pachet sub 102400 octeți (100 KB).

După alegerea unei valori potrivite pentru mărimea pachetelor de date, urmează setarea dimensiunii ferestrei de memorie. În cazul în care aceasta este setată să fie egală cu numărul de pachete, algoritmul de procesare a datelor își pierde utilitatea. Pentru dimensiuni prea mari ale acestei ferestre, apare aceeași problemă ca și în cazul lipsei ei: posibilitatea unei erori de memorie cauzată de volumul mare de date. Pentru valori prea mici ale acestui parametru, fereastra se va umple foarte repede, și foarte multe pachete vor fi aruncate, pentru ca apoi ele să fie retrimise de către sursă. Consecința acestui comportament este creșterea timpului de transfer direct proporțional cu scăderea ferestrei.

Din nou, ne confruntăm cu problema găsirii unei situații de mijloc, una care să rezolve ambele probleme descrise mai sus. În acest scop, tabelul 6.2 prezintă încercările

efectuate pentru a găsi o valoare optimă pentru fereastra de memorie. Testele au fost rulate pentru un fișier de 1GB și pachete de 1.75MB. Prin varierea numărului de destinații și a dimensiunii fișierului, s-a ajuns la concluzia că dimensiunea ferestrei depinde de numărul total de pachete. Mai precis, aceasta trebuie să fie aproximativ un sfert din numărul de pachete, pentru a nu cădea în nici una dintre extreme.

Desigur, aceste valori sunt general valabile, dar depind foarte mult de caracteristicile mașinii fizice pe care rulează serverul și de memoria pe care aceasta o are la dispoziție. De aceea parametrul pentru dimensiunea ferestrei și pentru dimensiunea pachetelor sunt încă modificabili, pentru a permite utilizatorului să aleagă cea mai bună soluție pentru fiecare transfer în parte.

# Capitolul 7

## Testare și validare

Secțiunile anterioare au detaliat proiectarea și implementarea unui sistem de replicare a volumelor mari de date în sisteme distribuite. Totuși, pentru utilizarea acestui sistem, trebuie demonstrată eficiența acestuia. Demonstrația se face printr două serii de teste.

### 7.1 Testare funcțională

Primul set de teste are rolul de a demonstra funcționalitatea sistemului și de a afla limitele acestuia în funcție de:

- numărul de destinații.
- volumul de date transferate.

Este ușor de observat că prin creșterea acestor parametri crește și timpul de transfer. De aceea, trebuie determinate limitele superioare. Testele efectuate în vederea găsirii volumului maxim de date au funcționat pentru fișiere de până la 1.7GB, replicate pe trei noduri. Trebuie remarcat faptul că între numărul de destinații și volumul datelor există o relație de inversă proporționalitate: creșterea valorii volumului datelor determină scăderea numărului maxim de destinații. Astfel, fișierele relativ mici pot fi replicate pe un număr mare de noduri. Din teste s-a obținut că replicarea se poate face pe 10 noduri cu fișiere de până la 500MB.

### 7.2 Testare eficiență

După ce am aflat limitele operaționale ale sistemului, al doilea set de teste are rolul de a vedea eficiența acestuia.

Pentru a realiza testarea eficienței, este necesară o privire comparativă asupra celor trei metode descrise în această lucrare. Prima dată vom vedea cum se comportă soluția

intitulată ‘distribuție serială’. Rezultatele testelor rulate folosind această metodă pot fi văzute în tabelul 7.1 și în figura 7.1.

		Număr destinații						
		2	3	4	5	6	7	8
Dimens. fișier	1MB	1.243	1.125	0.944	0.957	2.145	0.963	1.849
	10MB	2.646	3.222	3.2	3.58	6.367	5.025	9.194
	50MB	10.22	15.27	17.29	21.001	33.379	34.65	39.787
	100MB	18.183	28.365	32.252	39.618	55.078	40.584	56.662
	250MB	89.251	114.928	151.751	190.531	206.137	235.763	276.821
	500MB	89.251	114.928	151.751	190.531	206.137	175.763	276.821
	1GB	177.007	225.607	292.612	374.827	340.272	416.349	475.435

Tabela 7.1: Rezultate algoritm distribuție serială

Tabelul prezintă timpii de transfer (până ce fișierul e replicat pe toate nodurile) pentru două până la opt destinații și pentru fișiere de 1, 10, 50, 100, 250, 500 respectiv 1000 MB. Graficul 7.1 a fost generat folosind acest tabel. De aici se poate observa că timpul de transfer are o dependență liniară față de numărul de noduri.

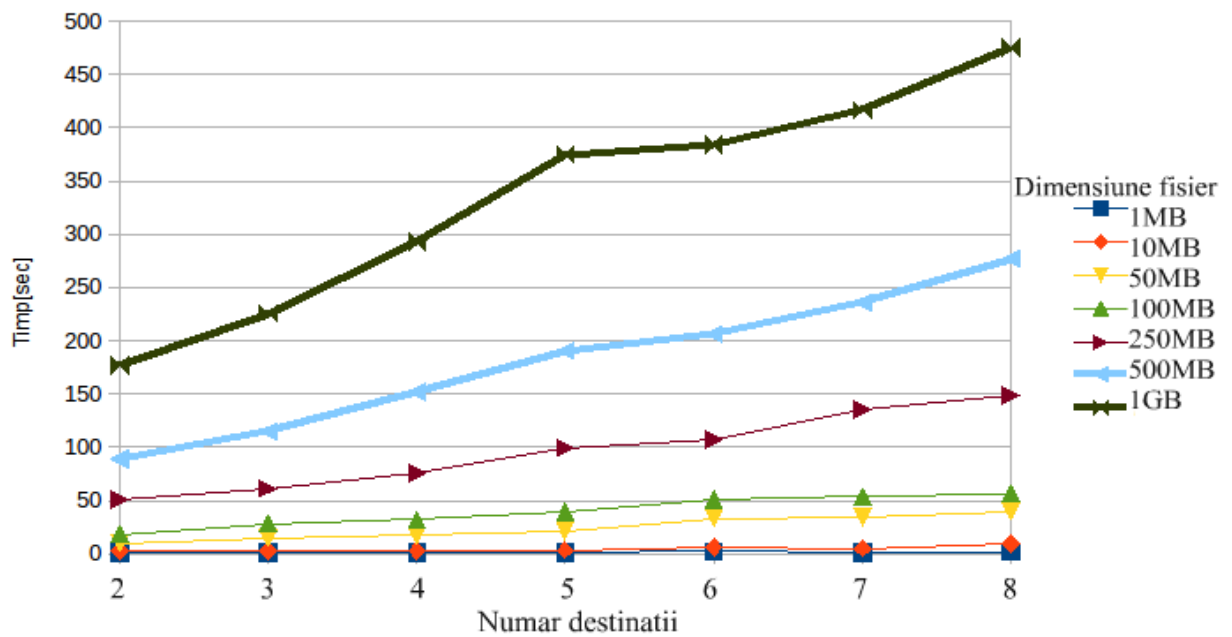


Figura 7.1: Rezultate grafice algoritm distribuție serială

Cel de-al doilea set de teste este aplicat metodei ‘în doi pași’, care este o versiune ușor modificată a *FastReplica*. Pentru aceste teste au fost folosite între trei și zece destinații, iar fișierele transferate au avut dimensiunile de 1, 10, 50, 100, 250 respectiv 500 MB.

		Număr destinații							
		3	4	5	6	7	8	9	10
Dimens. fișier	1MB	1.54	1.41	1.31	1.21	1.13	1.10	1.22	1.48
	10MB	5.34	4.25	4.27	3.88	3.68	3.49	3.36	3.56
	50MB	22.94	17.92	18.2	17.66	16.749	15.74	17.22	18.27
	100MB	72.17	55.94	47.61	42.26	39.37	36.82	47.23	56.012
	250MB	442.82	326.67	265.22	209.9	177.8	162.99	292.88	382.69
	500MB	1467.248	1326.25	1140.05	898.77	762.46	669.8	702.6	956.78

Tabela 7.2: Rezultate algoritm în doi pași

Rezultatele acestor teste pot fi studiate în tabelul 7.2, respectiv în figura 7.2. Timpii de transfer sunt prezentați în tabel, iar pentru o mai bună înțelegere a eficienței se poate consulta graficul generat de aceștia.

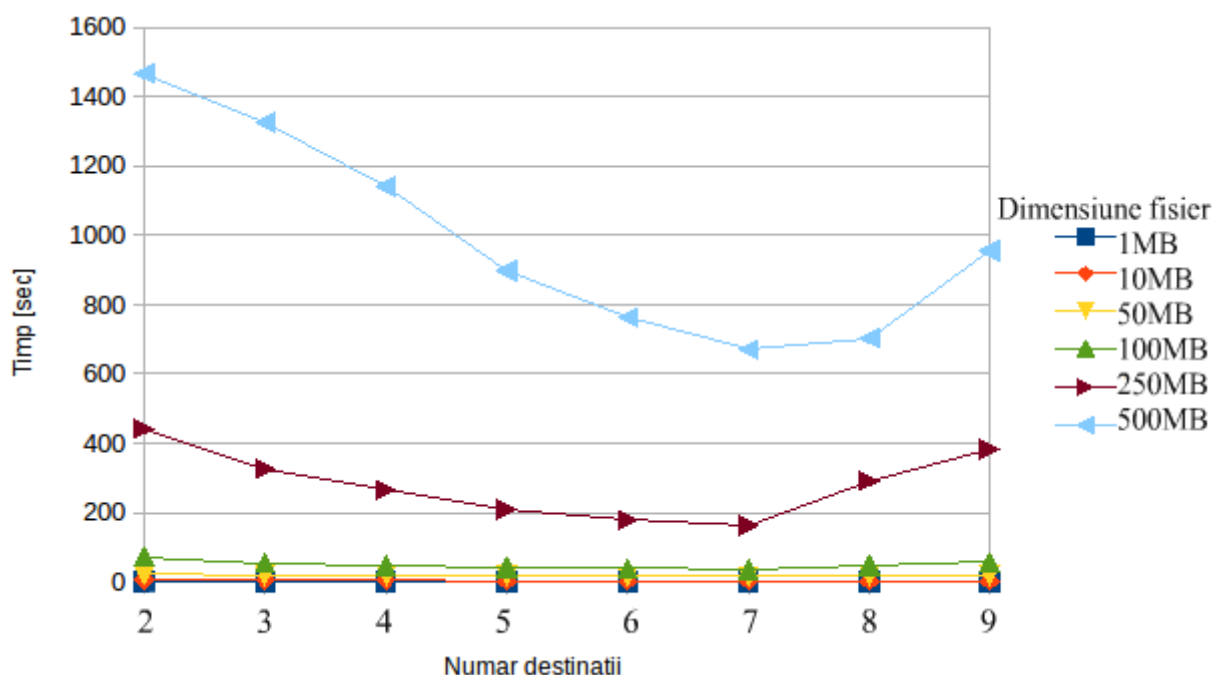


Figura 7.2: Rezultate grafice algoritm în doi pași

În grafic se poate vedea că timpii de transfer sunt relativ mari pentru un număr mic de noduri, scăzând până la un anumit punct, pentru ca apoi să crească din nou. Acest comportament se explică prin faptul că în această metodă fișierul este împărțit într-un număr de subfișiere egal cu numărul de destinații. De aceea, pachetul de date, care aici se numesc subfișiere, au dimensiuni mari, și necesită un timp de transfer individual foarte mare. Totuși, pentru valori potrivite ale numărului de noduri și ale volumului de date, această metodă este eficientă.

		Număr destinații							
		2	3	4	5	6	7	8	9
Dimens. fișier	50MB	5.958	8.919	9.38	9.284	12.055	11.635	11.935	34.473
	100MB	12.095	20.443	18.524	21.777	26.799	29.875	38.17	79.463
	250MB	32.051	38.427	38.816	38.796	85.775	58.06	71.024	203.35
	500MB	62.994	80.404	91.775	101.202	118.335	116.21	188.571	341.92
	1GB	129.542	174.297	188.29	215.158	338.338	320.54	336.817	666.54

Tabela 7.3: Rezultate algoritm MC-5.

Ultima soluție testată este MC-5. Pentru aceasta s-au folosit între două și nouă noduri destinație, iar pentru volumul de date 50, 100, 250, 500 respectiv 1000 MB. Rezultatele pentru această metodă sunt ilustrate în tabelul 7.3, respectiv în figura 7.3.

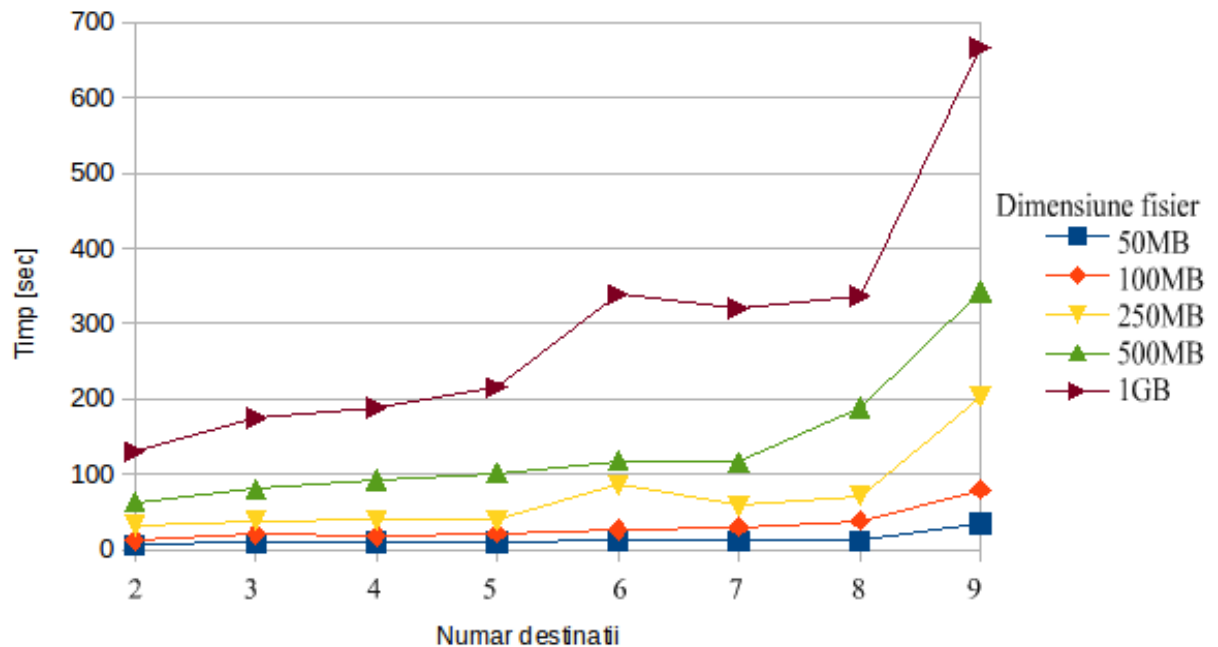
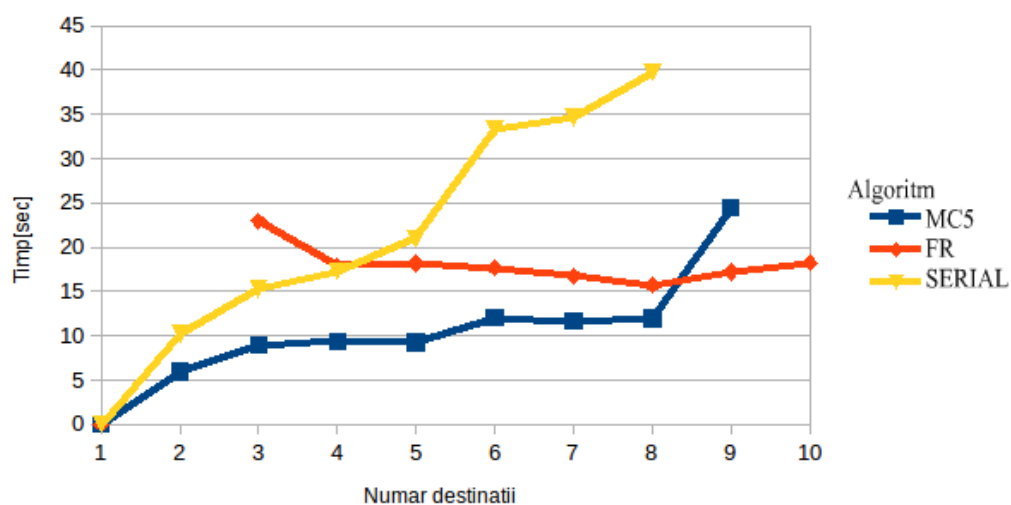


Figura 7.3: Rezultate grafice algoritm MC-5.

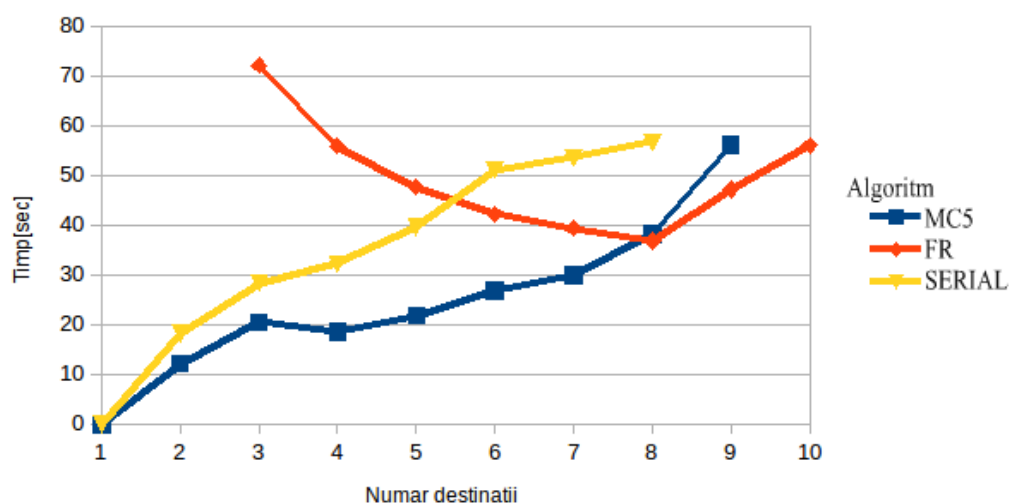
Tabelul conține, ca și în celelalte cazuri, timpii de transfer pentru valorile specificate. Analiza graficului arată că soluția aceasta înregistrează o creștere foarte mică a timpilor de transfer până ce se ajunge la șase destinații. Acest comportament apare din cauza structurii grupurilor de replicare pe baza cărora funcționează MC-5.

## 7.3 Comparare implementări

În figurile 7.4(a) și 7.4(b) se pot vedea comparativ performanțele fiecărei soluții descrise.



(a) Comparație pentru fișier de 50MB.



(b) Comparație pentru fișier de 100MB.

Figura 7.4: Compararea algoritmilor MC-5, FR și SERIAL pentru diferite dimensiuni de fișiere.

Soluția propusă în această lucrare este denumită MC-5, cea în doi pași este denumită FR, iar cea de distribuție serială SERIAL. Este ușor de observat din aceste grafice că MC-5 reprezintă o soluție bună din punct de vedere al timpilor de transfer. Comparativ cu SERIAL, această abordare are timpi de transfer cu până la 70% mai buni, iar în medie, transferurile sunt mai rapide folosind MC-5. Față de FR, soluția MC-5 nu este mereu mai bună, dar totuși obține timpi de transfer mai mici în aproximativ 80% din cazuri.



# Capitolul 8

## Manual de instalare și utilizare

Acest capitol conține instrucțiunile de instalare și utilizare ale aplicației prezentate. De asemenea, pentru un caz de utilizare, sunt prezentate și rezultatele așteptate în cazul unei execuții fără erori.

### 8.1 Instalare

#### 8.1.1 Resurse hardware

Ca și resurse hardware, este necesar un sistem de calcul care are minim următoarele configurații:

- procesor dual-core
- 2GB de memorie RAM
- 50MB spațiu de stocare pe disc
- placă de rețea

#### 8.1.2 Resurse software

Pentru a putea rula aplicația prezentată în lucrarea de față, următoarele resurse software sunt necesare:

- sistem de operare Ubuntu 12.04 LTS pe 32 biți
- Python 2.7.4

```
adelaneacsu@ubuntu:~$ sudo apt-get update  
adelaneacsu@ubuntu:~$ sudo apt-get install -y python2.7
```

- Twisted 13.0.0

```
adelaneacsu@ubuntu:~$ sudo apt-get install python-twisted
```

### 8.1.3 Instrucțiuni instalare

Codul sursă al algoritmilor prezentați aici se află în arhiva ‘my-protocols.tar.gz’, ce poate fi dezarhivată folosind comanda următoare:

```
adelaneacsu@ubuntu:~$ tar -xzf my-protocols.tar.gz
```

În interiorul acestei arhive se află alte trei arhive, câte una pentru fiecare soluție descrisă aici. Acestea pot fi dezarhivate în același mod ca mai sus. Punctul de interes este ‘my-mc.tar.gz’, în care se află soluția propusă în acest document.

## 8.2 Manual de utilizare

În primul rând trebuie precizat că pentru realizarea unui transfer folosind această aplicație sunt necesare cel puțin trei mașini virtuale: una pentru lansarea clientului și a sursei și două pentru serverele destinație. Nu este obligatoriu ca acestea să fie pe mașini fizice diferite, dar acest lucru este de preferat ținând cont că se dorește ca nodurile să fie dispersate.

Pentru pornirea unui server, indiferent dacă acesta va juca rol de sursă sau de destinație, se folosește:

```
adelaneacsu@ubuntu:~$ python my-mc/server/mainServer.py
```

Înainte de a lansa un client în execuție, trebuie creat fișierul de multicast, ce conține adresele IP ale serverelor destinație. Pentru pornirea clientului, se folosește următoarea comandă minimală (argumentele ce apar aici sunt necesare):

```
adelaneacsu@ubuntu:~$ python my-mc/client/mainClient.py  
-sf \textit{path_to_source_file}  
-df \textit{path_to_destination_file}  
-mc \textit{path_to_multicast_file}
```

Dacă se dorește măsurarea timpului de rulare a procesului client, se poate folosi comanda *time* din Linux:

```
adelaneacsu@ubuntu:~$ time python my-mc/client/mainClient.py .
```

Comenzile de mai sus pot primi argumente suplimentare pentru a eficientiza și pentru a personaliza transferul bazat pe cunoștințele utilizatorului despre rețea. De exemplu, la realizarea fișierului de multicast, adresele IP ale nodurilor destinație care se află în

aceeași rețea locală (LAN) pot fi scrise succesiv, pentru a fi conectate în același grup de replicare.

De asemenea, alegerea dimensiunii pachetelor joacă un rol important în eficientizarea transferului, după cum a fost demonstrat în capitolele anterioare. Există și alte argumente ce pot fi folosite după placul utilizatorului; important este să se înțeleagă conceptele care stau la baza acestuia.

## 8.3 Exemplu de utilizare

Pentru acest exemplu, avem la dispoziție mașinile virtuale din tabelul 8.1.

ID	Adresă IP	Nume	Locație
ADE	193.226.6.229	-	Cluj, România
CJ	193.226.5.102	gt-ige.utcluj.ro	Cluj, România
MA	147.96.25.28	gridway.fdi.ucm.es	Madrid, Spania
LY	134.158.75.91	vm-91.lal.stratuslab.eu	Lyon, Franța
AT	62.217.122.13	vm013.grnet.stratuslab.eu	Athena, Grecia

Tabela 8.1: Mașini virtuale disponibile pentru teste.

Considerăm mașina locală ADE, de pe care va fi lansat clientul. Tot aceasta va constitui și sursa, deci aici se află fișierul ce dorim să-l replicăm, numit ‘replica.file’. Primul pas este crearea fișierului de multicast numit aici ‘mc.file’. Acesta poate avea forma:

```
193.226.5.102
147.96.25.28:22000
134.158.75.91
62.217.122.13:22000
```

Pe mașinile virtuale CJ, MA, LY și AT se pornesc serverele, folosind instrucțiunile de mai sus. Comanda pentru lansarea clientului:

```
adelaneacsu@ubuntu:~$ time python my-mc/client/mainClient.py
-sf replica.file -df copy.file -mc mc.file
```

Răspunsul așteptat este în stilul acesta:

```
Connection made: 127.0.0.1 on port 22020
The file has 524288000 Bytes of data. Starting transfer ...
Transfer fully completed. Connection ended.

real    7m2.950s
user    0m0.140s
sys     0m0.036s
```

# Capitolul 9

## Concluzii

### 9.1 Analiza rezultatelor obținute

Un sistem distribuit a fost dezvoltat conform specificațiilor inițiale ce permite replicarea volumelor mari de date de pe un nod sursă pe un set de noduri destinație. Acesta se bazează pe o abordare de tip *divide et impera* și presupune împărțirea setului de destinații în seturi mai mici, care sunt mai ușor de analizat și de manevrat.

Sistemul respectă atât cerințele funcționale, cât și pe cele legate de eficiență, scalabilitate, adaptabilitate și fiabilitate. Prin experimente repetate, s-a demonstrat că abordarea folosită aici produce rezultate bune comparativ cu alte soluții existente. Performanța este atinsă prin folosirea unor algoritmi ce permit replicarea chiar și în cazul volumelor foarte mari de date.

În concluzie, obiectivele descrise la începutul lucrării au fost îndeplinite cu succes. Analiza rezultatelor testării dovedește că abordarea propusă are un potențial mare de dezvoltare și de extindere, ceea ce face din soluția aceasta o bază solidă pentru crearea unui sistem mai complex, îmbunătățit.

### 9.2 Dezvoltări ulterioare

Deși soluția propusă și-a îndeplinit scopul pentru care a fost creată, în acest domeniu există mereu loc de îmbunătățiri. De aceea, în această secțiune sunt propuse posibilități de dezvoltări ulterioare bazate pe sistemul prezentat în lucrarea de față:

- extinderea protocoalelor de comunicare pentru a permite transferul mai multor fișiere simultan.
- implementarea unui sistem complementar ce permite aducerea unui fișier de la mai multe surse concurrent.
- protejarea sistemului împotriva atacurilor malițioase prin metode de autentificare și prin securizarea conexiunilor.

- folosirea mai multor conexiuni TCP paralele între două noduri pentru a crește rata de date trimise.
- dezvoltarea unor algoritmi ce permit calcularea valorilor optime pentru pachetele de date și pentru fereastra de memorie.

Structura modulară a sistemului prezentat aici permite realizarea unora dintre aceste propuneri fără modificarea codului existent. Totuși, pentru anumite puncte, este necesară o nouă analiză a problemelor și a soluțiilor posibile.

# Bibliografie

- [1] J. Byers, J. Considine, M. Mitzenmacher, and S. Rost, “Informed content delivery across adaptive overlay networks,” in *SIGCOMM '02: Proceedings of the 2002 Conference on Applications, technologies, architectures, and protocols for computer communication*. New York, NY, USA: acm, 2002, pp. 455–465.
- [2] J. Lee and L. Cherkasova, “Fastreplica: Efficient large file distribution within content delivery networks,” in *In 4th USENIX Symposium on Internet Technologies and Systems*. USENIX Association, 2003, p. 7.
- [3] L. Cherkasova, “Optimizing the reliable distribution of large files within cdns,” in *ISCC '05: Proceedings of the 10th IEEE Symposium on Computers and Communications*. IEEE Computer Society, 2005, p. 7.
- [4] W. Allcock, “Gridftp protocol specification (global grid forum recommendation gfd.20),” *International Journal of High Performance Computing Applications: Special issue on Grid Computing: Infrastructure and Applications*, Mar. 2003. [Online]. Available: <http://www.globus.org/alliance/publications/papers/GFD-R.0201.pdf>
- [5] G. Khanna, U. Catalyurek, T. Kurc, R. Kettimuthu, P. Sadayappan, and J. Saltz, “A dynamic scheduling approach for coordinated wide-area data transfers using gridftp,” in *Proceedings of the 2008 IEEE international parallel and distributed processing symposium*. IEEE Computer Society, 14-18 April 2008, pp. 1 – 12.
- [6] T. Kurc, R. Kettimuthu, P. Sadayappan, J. Saltz, G. Khanna, and U. Catalyurek, “Using overlays for efficient data transfer over shared wide-area networks,” in *The International Conference for High Performance Computing, Networking, Storage, and Analysis*. IEEE Computer Society, 15-21 Nov. 2008, pp. 1 – 12.
- [7] A. N. Laboratory. (2013) Gridftp multicast. Internet page last accessed on June 2013. [Online]. Available: <http://www.mcs.anl.gov/bresnaha/mc/>
- [8] W. Allcock, J. B. R. Kettimuthu, M. Link, C. Dumitrescu, I. Raicu, and I. Foster, “The globus striped gridftp framework and server,” in *In SC '05: Proceedings of the 2005 ACM/IEEE conference on Supercomputing*, 2005.

- [9] R. Izmailov, S. Ganguly, and N. Tu, “Fast parallel file replication in data grid,” in *In Future of Grid Data Environments workshop (GGF-10)*, 2004.
- [10] R. Wolski, N. T. Spring, and J. Hayes, “The network weather service: A distributed resource performance forecasting service for metacomputing,” *Journal of Future Generation Computing Systems*, vol. 15, pp. 757–768, 1999.
- [11] A. Vakali and G. Pallis, “Content delivery networks: status and trends,” *Internet Computing, IEEE*, pp. 68 – 74, Nov.-Dec. 2003.
- [12] W. F. Inc. (2013) Content delivery network. Internet page last modified on 22 June 2013. [Online]. Available: [http://en.wikipedia.org/wiki/Content\\_delivery\\_network](http://en.wikipedia.org/wiki/Content_delivery_network)
- [13] I. S. I. . U. of Southern California, “TRANSMISSION CONTROL PROTOCOL,” RFC 793, Internet Engineering Task Force, Sep. 1981. [Online]. Available: <http://www.ietf.org/rfc/rfc793.txt>

# Capitolul 10

## Anexe

```
import logging

from twisted.protocols.basic import LineReceiver

class MyClientProtocol(LineReceiver):

    def __init__(self, factory):
        self.factory = factory

    def connectionMade(self):
        msg = 'Connection made: %s' % self.transport.getPeer()
        logging.info(msg)
        print 'Connection made: %s on port %d' % (self.transport.
getPeer().host, self.transport.getPeer().port)

        line = 'SRFI ' + self.factory.sourceFile
        self.sendLine(line)
        line = 'DSFI ' + self.factory.destFile
        self.sendLine(line)

    def lineReceived(self, line):
        print line
        data = line.strip().split(' ')
        if data[0] == 'FISZ':
            # file size
            logging.info('File ' + self.factory.sourceFile + ' has '
+ data[1] + ' bytes.')
            size = int(data[1])
            if size == 0:
                logging.info('Transfer stopped.')
```



```

        print 'The file does NOT exist on the source server
or it does not have reading rights.'
        print 'Communication with the source server will now
stop.'

        self.factory.callback._stop()
    else:
        print 'The file has %d Bytes of data.' % size
        sLine = 'NRDS ' + str(len(self.factory.destinations))
        self.sendLine(sLine)
        sLine = 'PKSZ ' + str(self.factory.packetSize)
        self.sendLine(sLine)
        sLine = 'WDSZ ' + str(self.factory.windowSize)
        self.sendLine(sLine)

    elif data[0] == 'DSOK':
        # received number of destinations
        for dst in self.factory.destinations:
            sLine = 'DSIP %s %d' % (dst[0], dst[1])
            self.sendLine(sLine)

    elif data[0] == 'TRST':
        # transfer started
        msg = 'Transfer started ...'
        logging.info(msg)
        print msg

    elif data[0] == 'STAT':
        # statistics
        msg = 'Transfer fully completed. Connection ended.'
        print msg
        logging.info(msg)
        self.factory.callback._stop()

    elif data[0] == 'ERRR':
        msg = 'An error ocurred. Connection ended.'
        print msg
        logging.info(msg)
        self.factory.callback._stop()

```

clientProtocol.py

```

import logging
import struct
import time

```

```

from twisted.protocols.basic import LineReceiver
from twisted.internet.protocol import ServerFactory
from twisted.internet.endpoints import TCP4ClientEndpoint
from twisted.internet import reactor, defer

from multiEchoFactory import *

class MyReceiverProtocol(LineReceiver):

    def __init__(self, factory):
        self.factory = factory
        self.MAXLENGTH = 524288010

    def connectionMade(self):
        if self.factory.source is None:
            #print 'source is %d' % self.transport.getPeer().port
            self.factory.source = self
        logging.info('Connection made: %s' % self.transport.getPeer())

    def connectionLost(self, reason):
        if self.factory.source == self:
            self.factory.source = None
            for echoer in self.factory.echoFactory.echoers:
                echoer.transportloseConnection()
        logging.info('Connection lost: %s' % self.transport.getPeer())

    def lineReceived(self, line):
        if line[0:4] == '0000':
            # this is a packet – store and forward
            if self.factory.source == self:
                for echoer in self.factory.echoFactory.echoers:
                    echoer.sendLine(line)
            return self._processData(line[4:])
        else:
            print 'rc = %s %s' % (line, self.transport.getPeer().host)

            data = line.strip().split(' ')
            if data[0] == 'SCHM':
                # SCHM number of schema (0, 1, 2 or 3)
                try:
                    self.factory.nrConnections = int(data[1]) + 1

```

```

        except IndexError:
            self.sendLine('ERRR')
            return

    elif data[0] == 'SIBL' or data[0] == 'CHLD':
        # SIBL/CHLD IP port
        try:
            IP = data[1]
            port = int(data[2])
            endpoint = TCP4ClientEndpoint(reactor, IP, port)
            endpoint.connect(self.factory.echoFactory).
addCallback(self._newConnection)
        except IndexError:
            self.sendLine('ERRR')
            return

    elif data[0] == 'FPTH':
        # FPTH filepath
        try:
            self.factory.fileObj = open(data[1], 'wb')
            self.factory.nextToWrite = 0
            logging.info('File %s opened for writing. Raw
mode set.' % data[1])
        except IndexError:
            self.sendLine('ERRR')
            return

        # start statistics
        self.factory.startTime = time.time()
        self.factory.minTime = 1000000
        self.factory.maxTime = 0
        self.factory.avgTime = 0
        self.factory.nrPacksRec = 0
        # forward configurations
        if self.factory.source == self:
            for echoer in self.factory.echoFactory.echoers:
                echoer.sendLine(line)

    elif data[0] == 'PACK':
        # PACK nrPacks sizePack
        try:
            self.factory.nrPackets = int(data[1])
            self.factory.packetSize = int(data[2])
            self.factory.packetsReceived = [False] * self.
factory.nrPackets

```

```

        except IndexError:
            self.sendLine('ERRR')
            return
        if self.factory.source == self:
            for echoer in self.factory.echoFactory.echoers:
                echoer.sendLine(line)

    elif data[0] == 'WDSZ':
        # WDSZ windSz
        try:
            self.factory.windowSize = int(data[1])
            self.factory.window = [""] * self.factory.
windowSize
            # -1 if the slot is free, packet index otherwise
            self.factory.slotBusy = [-1] * self.factory.
windowSize
        except IndexError:
            self.sendLine('ERRR')
            return
        if self.factory.source == self:
            for echoer in self.factory.echoFactory.echoers:
                echoer.sendLine(line)
            self.factory.source.sendLine('SEND')

    elif data[0] == 'REQC':
        # REQC — check of all packets arrived
        allSent = True
        for bitIndex in range(len(self.factory.
packetsReceived)):
            if not self.factory.packetsReceived[bitIndex]:
                allSent = False
                self.sendLine('NREC ' + str(bitIndex))
                break
        if allSent:
            self.sendLine('RECA')

def _processData(self, rawData):
    logging.info('Received %d bytes.' % (len(rawData) - 4))
    # split : header + data
    currIndex = int(rawData[0:4])
    print 'PACK %d' % currIndex
    # statistics
    self.factory.nrPacksRec += 1

```

```

currTime = time.time()
delta = currTime - self.factory.startTime
if delta < self.factory.minTime:
    self.factory.minTime = delta
elif delta > self.factory.maxTime:
    self.factory.maxTime = delta
self.factory.avgTime += delta
self.factory.startTime = currTime

if self.factory.packetsReceived[currIndex] == False:
    # only keep first copy of each packet, since they might
    arrive on more than one links
    self.factory.packetsReceived[currIndex] = True
    if self.factory.nextToWrite == currIndex:
        # current packet must be written directly to file
        self.factory.fileObj.write(rawData[4:])
        self.factory.nextToWrite += 1
        countdown = self.factory.windowSize
        done = False
        # also write all packets that come after
        while not done:
            for slotIndex in range(self.factory.windowSize):
                countdown -= 1
                if self.factory.slotBusy[slotIndex] == self.
factory.nextToWrite:
                    self.factory.fileObj.write(self.factory.
window[slotIndex])
                    self.factory.slotBusy[slotIndex] = -1
                    self.factory.nextToWrite += 1
                    countdown = self.factory.windowSize
                if countdown == 0:
                    done = True
            else:
                # find an empty slot to store data
                for slotIndex in range(self.factory.windowSize):
                    if self.factory.slotBusy[slotIndex] == -1:
                        self.factory.window[slotIndex] = rawData[4:]
                        self.factory.slotBusy[slotIndex] = currIndex
                        break

if self.factory.nextToWrite == self.factory.nrPackets:
    # file was completely written
    self.factory.fileObj.close()
    if len(self.factory.echoFactory.echoers) == 0:

```

```

        self.sendLine('CHRA')

    def _newConnection(self, deferred):
        if len(self.factory.echoFactory.echoers) == self.factory.
nrConnections:
            # all nodes are connected
            self.sendLine('CONN')

```

receiverProtocol.py

```

import os
import logging

from multiprocessing.connection import Client
from twisted.internet.protocol import ServerFactory
from twisted.protocols.basic import LineReceiver

class MyFileSenderProtocol(LineReceiver):

    def __init__(self, factory):
        self.factory = factory

    def connectionMade(self):
        self.sendLine('SCHM ' + str(len(self.factory.children)))
        self.sendLine('SIBL ' + self.factory.sibling[0] + ' ' + str(
self.factory.sibling[1]))
        for child in self.factory.children:
            self.sendLine('CHLD ' + child[0] + ' ' + str(child[1]))
            logging.info('Connection made: %s' % self.transport.getPeer())
        )

    def connectionLost(self, reason):
        logging.info('Connection lost: %s' % self.transport.getPeer())
        )

    def lineReceived(self, line):
        print 'fs = %s %s' % (line, self.transport.getPeer().host)
        data = line.strip().split(' ')
        if data[0] == 'CONN':
            # get configuration data
            if self.factory.parent._getConfigurations(self.factory.
groupNr) == 1:
                # first connected
                self.sendLine('FPTH ' + self.factory.parent.dspath)

```

```

        nrPackets = self.factory.parent.fileSize / self.
factory.parent.packetSize
        if self.factory.parent.fileSize % self.factory.parent
.packetSize != 0:
            nrPackets += 1
            self.sendLine('PACK ' + str(nrPackets) + ' ' + str(
self.factory.parent.packetSize))
            self.sendLine('WDSZ ' + str(self.factory.parent.
windowSize))

    elif data[0] == 'SEND':
        # server received configurations, now is asking for data
        packet = self.factory.parent._getPacketByld(self.factory.
parity)
        k = 2
        while packet != -1:
            self.sendLine(packet)
            packet = self.factory.parent._getPacketByld(self.
factory.parity + k)
            k += 2
            self.sendLine('REQC')

    elif data[0] == 'NREC':
        # at least one packet did not arrive to the destination,
        resend it
        packet = self.factory.parent._getPacketByld(int(data[1]))
        self.sendLine(packet)
        self.sendLine('REQC')

    elif data[0] == 'RECA':
        # all file was received
        print '%s %s' % (line, self.transport.getPeer().host)
        self.factory.parent._incFileDone(self)

```

fileSenderProtocol.py

```

import os
import logging

from twisted.protocols.basic import LineReceiver
from twisted.internet.endpoints import import TCP4ClientEndpoint,
TCP4ServerEndpoint
from twisted.internet import import reactor

from fileSenderFactory import *

```

```
class MySenderProtocol(LineReceiver):

    def __init__(self):
        self.directPeers = []
        logging.info('Server started.')

    def connectionMade(self):
        self.fileDoneCnt = 0
        self.destinations = []
        self.client = self.transport.getPeer()
        logging.info('Connection made: %s' % self.client)

    def connectionLost(self, reason):
        if hasattr(self, 'client'):
            logging.info('Connection ended from: %s' % self.client)
        else:
            logging.error('Connection ended unexpectedly. Reason: %s'
                           % reason)

    def lineReceived(self, line):
        #print 'sd = %s %d' % (line, self.transport.getPeer().port)
        data = line.strip().split(' ')
        if data[0] == 'SRFI':
            try:
                self.filepath = data[1]
            except IndexError:
                self.sendLine('ERRR')
                return
            if not os.path.exists(self.filepath):
                self.sendLine('FISZ 0')
                return
            self.fileSize = os.stat(self.filepath).st_size
            self.sendLine('FISZ ' + str(self.fileSize))
            try:
                self.fileObj = open(self.filepath, 'rb')
            except IOError:
                self.sendLine('ERRR')
                return

        elif data[0] == 'DSFI':
            try:
                self.dstpath = data[1]
            except IndexError:
```



```

        self.sendLine('ERRR')
        return

elif data[0] == 'NRDS':
    try:
        self.nrDestinations = int(data[1])
    except IndexError:
        self.sendLine('ERRR')
        return
    # decide distribution schema
    self.schemaIndices = [0, 0, 0, 0] # 2, 3, 4, 5
    if self.nrDestinations % 4 == 0:
        self.schemaIndices[2] = self.nrDestinations / 4
    elif self.nrDestinations % 4 == 1:
        self.schemaIndices[2] = self.nrDestinations / 4 - 1
        self.schemaIndices[3] = 1
    elif self.nrDestinations % 4 == 2:
        k = self.nrDestinations / 4
        if k == 0:
            self.schemaIndices[0] = 1
        elif k == 1:
            self.schemaIndices[1] = 2
        else:
            self.schemaIndices[2] = k - 2
            self.schemaIndices[3] = 2
    elif self.nrDestinations % 4 == 3:
        k = self.nrDestinations / 4
        if k == 0:
            self.schemaIndices[1] = 1
        elif k == 1:
            self.schemaIndices[1] = 1
            self.schemaIndices[2] = 1
        elif k == 2:
            self.schemaIndices[1] = 1
            self.schemaIndices[2] = 2
        else:
            self.schemaIndices[2] = k - 3
            self.schemaIndices[3] = 3
    self.sendLine('DSOK')

elif data[0] == 'PKSZ':
    try:
        self.packetSize = int(data[1])
    except IndexError:

```

```

        self.sendLine('ERRR')
        return

    elif data[0] == 'WDSZ':
        try:
            self.windowSize = int(data[1])
        except IndexError:
            self.sendLine('ERRR')
            return

    elif data[0] == 'DSIP':
        try:
            destIP = data[1]
            destPort = int(data[2])
        except IndexError:
            self.sendLine('ERRR')
            return
        dst = [destIP, destPort]
        index = len(self.destinations)
        self.destinations.append(dst)

        if index == self.nrDestinations - 1:
            start = 0
            group = 0
            for i in range(len(self.schemaIndices)):
                for k in range(0, self.schemaIndices[i]):
                    lc = TCP4ClientEndpoint(reactor, self.
destinations[start][0], self.destinations[start][1])
                    rc = TCP4ClientEndpoint(reactor, self.
destinations[start+1][0], self.destinations[start+1][1])
                    lc.connect(MyFileSenderFactory(self, self.
destinations[start+1], self.destinations[start+2:start+2+i], group
, 0))
                    rc.connect(MyFileSenderFactory(self, self.
destinations[start], self.destinations[start+2:start+2+i], group,
1))
                    start += (2 + i)
                    group += 1
            self.configs = [0] * group

    def _incFileDone(self, peer):
        self.fileDoneCnt += 1
        logging.info('Transfer to %s completed.' % peer)
        if self.fileDoneCnt == self.nrDestinations:

```

```

        finishMsg = 'All destinations received the file. Transfer
fully completed.'
        logging.info(finishMsg)
        print finishMsg
        for cPeer in self.directPeers:
            cPeer.transportloseConnection()
        self.sendLine('STAT')
    else:
        self.directPeers.append(peer)

def _getPacketById(self, packetId):
    offset = packetId * self.packetSize
    chunk = ''
    if self.fileSize > offset:
        # there is still data to read from file
        currOffset = self.fileObj.tell()
        if currOffset != offset:
            self.fileObj.seek(offset)
        chunk = '0000' + str(packetId).zfill(4) + self.fileObj.
read(self.packetSize)
        return chunk
    else:
        return -1

def _getConfigurations(self, groupNumber):
    if self.configs[groupNumber] == 0:
        self.configs[groupNumber] = 1
        return 1
    else:
        return -1

```

senderProtocol.py