

Chordy

Iulian Baltag

Faculty of Computer Science Iasi
<https://www.info.uaic.ro>

Abstract. Implementarea unui design DHT bazat pe Chord si structurat in doua aplicatii server si client.

Keywords: DHT · Chord · TCP.

1 Introducere

O tabela hash distribuita (DHT) functioneaza prin asocierea cheilor cu nodurile dintr-un sistem distribuit. Fiecare nod va stoca valorile corespunzatoare cheilor distribuite acestuia. Un exemplu bine cunoscut de design pentru DHT este chord. Acesta este un algoritm si protocol pentru o tabela hash distribuita *peer to peer*. Chord specifica cum o cheie este atribuita unui nod si cum un nod poate descoperi o valoare gasind mai intai nodul responsabil de cheia pentru acea valoare. In aplicatia prezentata ulterior un client poate trimite comenzi unui server(nod) pentru adaugarea, stergerea si gasirea unei chei.

2 Tehnologii Aplicate

- Aplicatia este implementata in limbajul de programare C.
- Am implementat un protocolul de comunicare pe baza TCP, oferind o conexiune fiabila si orientata pe flux intre server si clienti.
 - TCP asigura transmiterea mesajelor in ordine si fara erori, lucru necesar pentru functionarea unui DHT precum Chord. Rutarea cheilor, comunicarea intre noduri, analizarea succesorilor necesita transmiterea de mesaje fara erori.
 - TCP permite conexiunea permanenta a nodurilor. Acest lucru este necesar deoarece nodurile comunica frecvent pentru a mentine structura retelei. Oferă un canal bidirectional.
- Pentru gestionarea conexiunilor multiple am folosit threaduri. Acestea asigura concurenta si sunt mai eficiente ca procesele in privinta resurselor si a timpului.
- Pentru stocarea perechilor *cheie-valoare* am folosit o baza de date *sqlite*. Am creat cate un tabel pentru fiecare nod.
- Pentru functia de hash-ing am folosit libraria *openssl* utilizand SHA1.

3 Structura aplicației

Aplicația Chord este structurată în server și client. Acestea comunică între ele prin socketuri. Utilizatorii vor avea la dispoziție următoarele comenzi:

- **ADD** $\langle key \rangle \langle value \rangle$ Adaugă o pereche cheie-valoare în sistemul distribuit. Perechea cheie-valoare este trimisă nodului responsabil de stocarea acesteia (după hashing și eventualele rutări).
- **FIND** $\langle key \rangle$ Găsește valoarea asociată unei chei. Nodul curent verifică dacă cheia îi aparține, dacă nu folosește mecanismul de rutare pentru a găsi nodul responsabil. Valoarea este returnată clientului.
- **REMOVE** $\langle key \rangle$ Ștergerea perechei cheie-valoare asociată cu cheia respectivă. Similar ca mai sus, nodul responsabil de cheie o șterge pe aceasta din memoria sa locală.

Următoarea diagramă ilustrează funcționalitatea proiectului:

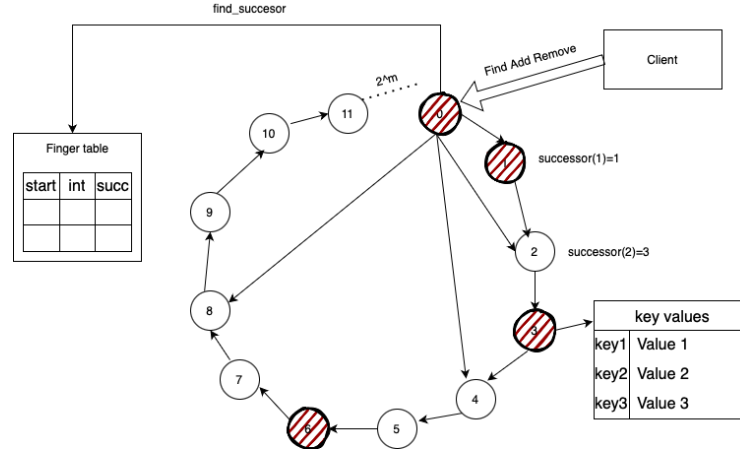


Fig. 1. Diagramă Chord

Se poate utiliza tabela *Finger table* pentru optimizarea căutărilor. Fiecare nod stochează local perechile cheie-valoare asociate acestuia precum și nodurile succesor și predecesor. Nodurile hash-urate cu roșu sunt nodurile "online" din rețea. Un astfel de nod este responsabil de toate valorile care îl preced pe inel începând de după ultimul nod online. În rețea vom avea 2^m noduri și fiecare nod va avea un ID format din primii m biți din hash.

4 Aspecte de implementare

Pentru comunicare vom folosi *socketuri*. În secțiunile de cod de mai jos este ilustrat protocolul prin care clientul comunică cu serverul urmând pașii *connect*, *send*, *recv* și inițierea structurii *sockaddr* necesare comunicării.

Inițiere:

```

1 int s = socket(AF_INET, SOCK_STREAM, 0);
2   if (s == -1)
3       perror("[client] Eroare la socket");
4
5   struct sockaddr_in server;
6   server.sin_family = AF_INET;
7   server.sin_port = htons(port);
8   server.sin_addr.s_addr = inet_addr(ip);

```

Conectarea la server:

```

1   if (connect(s, (struct sockaddr *) &server, sizeof (
2       struct sockaddr)) == -1)
3       perror("[client] Eroare la connect().\n");

```

Trimiterea comenzilor:

```

1   printf("Se trimite comanda %s\n", comanda);
2   send(s, comanda, strlen(comanda), 0);

```

Primirea răspunsului:

```

1   int bytes_received = recv(s, buffer, BUFFER_SIZE, 0);
2   if (bytes_received == -1)
3       perror("[client] Eroare la recv() de la server.\n");
4   else
5   {
6       buffer[bytes_received] = '\0';
7       printf("Răspuns primit de la server: %s\n", buffer);
8   }

```

Acestea sunt urmate de instrucțiunile *bind*, *listen*, *accept* din server. În cadrul serverului ne folosim de threaduri prin intermediul bibliotecii *pthread.h*. Astfel se va crea un nou thread pentru fiecare comandă dată de useri.

```

1   pthread_t thread;
2   if(pthread_create(&thread, NULL, handle_request, (
3       void*)new_socket) != 0) {
4       perror("[server] Thread creation failed");
5       close(client_socket);
6       free(new_socket);
7       continue;
8   }
9   pthread_detach(thread);

```

Fiecare comandă primită din server va fi parsată de către handler **ȘI FORWARD-ATĂ DACĂ ESTE CAZUL** (la succesorul imediat sau utilizând finger table), comenzile transmise de clienții din afara rețelei vor fi *FIND*, *ADD*, *REMOVE*. Vom gestiona secțiunile de cod critice cu ajutorul mutexurilor din biblioteca *pthread.h*.

```

1  if(strncmp(buffer, "FIND_", 5) == 0) {
2      //implemenatre FIND
3      pthread_mutex_lock(&mutex_chord);
4      /* sectiune critica */
5      pthread_mutex_unlock(&mutex_chord);
6  }
7  else if(strncmp(buffer, "ADD_", 4) == 0) {
8      //implementare ADD
9      pthread_mutex_lock(&mutex_chord);
10     /* sectiune critica */
11     pthread_mutex_unlock(&mutex_chord);
12 }
13 else if(strncmp(buffer, "REMOVE_", 7) == 0) {
14     //implementare REMOVE
15     pthread_mutex_lock(&mutex_chord);
16     /* sectiune critica */
17 }
18
19 .
20 . etc
21 .
22
23 else {
24     snprintf(raspuns, BUFFER_SIZE-1, "Please use:_(ADD_<
25     key>_<value>_;FIND_<key>_;REMOVE_<key>_;exit_
        :_");
    }

```

Alte comenzi utilizate pentru menținerea rețelei vor fi *UPDATE_P*, *UPDATE_S*, *REQUEST_KEYS*, *TRANSFER_KEY*, *PING* etc.

Nodurile vor fi reprezentate printr-o structură. Nodul inițial va avea succesorul și predecesorul el însuși urmând să fie actualizate la reuniunea cu alte noduri.

```

1  typedef struct {
2      int id;                // ID nod (hash-ul portului)
3      int port;              // Port nod
4      int successor;         // Port succ
5      int predecessor;       // Port pred
6      int finger_table[M];
7      int start[M];
8      int backup_successor;   // Al doilea succes (succesorul
9                              // succesoriului)
10 }Node;
11
12 Node node;

```

Criptarea cheilor se va face folosind SHA1 din librăria *openssl*.

```

1 int hash_sha1(int key) {
2     //Convertim key in string
3     char string_cheie[64];
4     snprintf(string_cheie, sizeof(string_cheie), "%d", key);
5
6     unsigned char digest[SHA_DIGEST_LENGTH]; //
7     SHA_DIGEST_LENGTH = 20
8
9     SHA1((unsigned char*)string_cheie, strlen(string_cheie),
10         digest);
11
12     //Interpretam primii 4 bytes din SHA1 ca un numar (32 de
13     //biti)
14     uint32_t hash_val = 0;
15     for (int i = 0; i < 4; i++) {
16         hash_val = (hash_val << 8) | digest[i];
17     }
18
19     //Facem modulo 2^m = RING_SIZE
20     return (hash_val % 256);
21 }

```

Pentru funcționalitatea proiectului am implementat următoarele funcții:

```

1 void send_join_message(int port_conectare, int port);
2 void send_update_message(int port, const char *message);
3 int send_message(int port, const char *message, char*
4     response);
5 void show_node();
6 int send_find_successor(int target_port, int id);
7 int closest_preceding_node(int id);
8 int find_successor(int id);
9 void update_finger_table();
10 void show_finger_table();
11 void start_table();
12 void *stabilizer();
13 int preceding_finger(int id);
14 void transfer_keys_to_new_node(int new_node_port, int
15     new_node_pred);
16 int check_node_alive(int port);
17 void is_succesor_alive();

```

Inițializarea bazei de date se face în modul următor:
Ne folosim de o tabelă pentru fiecare nod din rețea

```

1 void init_database(int id_nod) {
2     char sql[256];
3     char *err_msg = NULL;
4
5     if (sqlite3_open("bd_centralizata.db", &db)) {

```

```

6         perror("Eroare la deschiderea bazei de date");
7         exit(1);
8     }
9     printf("Successfully opened 'bd_centralizata.db'.\n");
10
11     //Setam un busy timeout si bd in modul wall pentru a
12     ajuta concurenta
13     sqlite3_busy_timeout(db, 1000);
14     if (sqlite3_exec(db, "PRAGMA journal_mode=WAL;", NULL,
15         NULL, &err_msg) != SQLITE_OK) {
16         fprintf(stderr, "Error setting WAL mode: %s\n",
17             err_msg);
18         sqlite3_free(err_msg);
19     }
20
21     // Pornim cu un table gol
22     snprintf(sql, sizeof(sql), "DELETE FROM nod_%d;", id_nod)
23     ;
24     if (sqlite3_exec(db, sql, 0, 0, &err_msg) != SQLITE_OK) {
25         if (strstr(err_msg, "no such table") == NULL) {
26             printf("SQL delete error");
27             sqlite3_free(err_msg);
28             sqlite3_close(db);
29             exit(1);
30         }
31         sqlite3_free(err_msg);
32     } else {
33         printf("Existing data from 'nod_%d' has been cleared
34             .\n", id_nod);
35     }
36
37     //Creeaza tabela
38     snprintf(sql, sizeof(sql), "CREATE TABLE IF NOT EXISTS
39         nod_%d (key INTEGER PRIMARY KEY, value TEXT);",
40         id_nod);
41     if (sqlite3_exec(db, sql, 0, 0, &err_msg) != SQLITE_OK) {
42         printf("SQL table error");
43         sqlite3_free(err_msg);
44         sqlite3_close(db);
45         exit(1);
46     }
47
48     printf("Table 'nod_%d' has been created.\n", id_nod);
49 }

```

Tabelele finger table se actualizează automat la nodurile afectate de intrarea sau ieșirea din rețea a unui alt nod. De asemenea acestea se actualizează periodic pentru toate nodurile active. Pentru menținerea acestor tabele am utilizat funcțiile:

```

1 void update_finger_table() {
2     for (int i = 0; i < M; i++) {
3
4         int successor = find_successor(node.start[i]);
5         if (successor == node.port || successor <= 0) {
6             printf("[upd_finger_table] Entry %d remains the
7                 same (port %d).\n", i, node.port);
8             successor = node.port;
9         }
10
11         node.finger_table[i] = successor;
12         printf("[upd_finger_table] Finger table [%d]: %d\n", i
13             , successor);
14     }
15 }
16
17 void fix_fingers() {
18     static int next = 0;
19     int start = (node.id + (1 << next)) % (1 << M);
20
21     int successor = find_successor(start);
22     if (node.finger_table[next] != successor) {
23         printf("[fix_fingers] Updating finger [%d] from %d to
24             %d\n", next, node.finger_table[next], successor);
25         node.finger_table[next] = successor;
26     }
27
28     next = (next + 1) % M;
29 }

```

Un punct important în rețeaua chord este transferul de chei în urma părăsirii nodurilor sau după intrarea altor noduri în rețea. Pentru părăsirea voluntară am implementat următorul protocol: Nodul care urmează să părăsească rețeaua trimite UPDATE_S la predecesor și UPDATE_P la succesori și transferă cheile succesoriului. La JOINUL unui nod în rețea următorul protocol este respectat: Nodul care va fi predecesorul nodului care intră în rețea va trimite mesaje de tip UPDATE_S, UPDATE_P noului nod. Noul nod va trimite o cerere REQUEST_KEYS noului său succesori și va primi răspunsuri de tipul TRANSFER cheie valoare. Transferul cheilor se realizează în modul următor:

```

1 void transfer_keys_to_new_node(int new_node_port, int
2     new_node_pred) {
3     printf("[transfer_keys] Node %d transferring keys to port
4         %d\n", node.id, new_node_port);
5
6     int new_node_id = hash(new_node_port);
7     int pred_new_node = hash(new_node_pred);

```

```

7   printf("[transfer_keys]_new_node_id=%d,_pred_new_node_id=
   _%d\n",new_node_id, pred_new_node);
8
9   char select_cmd[256];
10  snprintf(select_cmd, sizeof(select_cmd),"SELECT_key,_
   value_FROM_nod_%d;", node.id);
11
12  char *error = NULL;
13  char **rezultat = NULL;
14  int rows, cols;
15
16  pthread_mutex_lock(&mutex_db);
17  if (sqlite3_get_table(db, select_cmd, &rezultat, &rows, &
   cols, &error) != SQLITE_OK) {
18      fprintf(stderr, "[transfer_keys_to_new_node]_SELECT_
   error:_%s\n", error);
19      sqlite3_free(error);
20      return;
21  }
22
23  int index = cols;
24  for (int i = 0; i < rows; i++) {
25      int k = atoi(rezultat[index]);          // key
26      char *val = rezultat[index + 1];        // value
27      index += cols;
28
29      int key_hash = k; //pentru simplitate folosim insusi
   cheia ca hash
30
31      if (
32          // Caz normal(pred_new_node < new_node_id)
33          ((pred_new_node < new_node_id) &&
34           (key_hash > pred_new_node && key_hash <=
35            new_node_id))
36          ||
37          // Caz wrap-around(pred_new_node > new_node_id)
38          ((pred_new_node > new_node_id) &&
39           (key_hash > pred_new_node || key_hash <=
40            new_node_id))
41          ||
42          pred_new_node == new_node_id // Caz unde ramanem
   cu un singur nod
43      ) {
44          char msg[BUFFER_SIZE];
45          snprintf(msg, sizeof(msg),"TRANSFER_KEY_%d_%s_%d"
   , k, val, new_node_port);
46
47          send_update_message(new_node_port, msg);
48          printf("[transfer_keys]_Transferred_key_%d->_
   node_%d_OK.\n",k, new_node_port);

```



```

47
48     char delete_cmd[256];
49     snprintf(delete_cmd, sizeof(delete_cmd), "DELETE_
        FROM_nod_%d_WHERE_key_=_%d;", node.id, k);
50
51     char *delete_error = NULL;
52     if (sqlite3_exec(db, delete_cmd, 0, 0, &
        delete_error) != SQLITE_OK) {
53         fprintf(stderr, "Error_deleting_key_%d:_%s\n"
            , k, delete_error);
54         sqlite3_free(delete_error);
55     }
56 }
57 }
58 sqlite3_free_table(rezultat);
59 pthread_mutex_unlock(&mutex_db);
60 }

```

Am implementat în acest proiect:

- mecanismul de stabilire a succesorului și a predecesorului. Voi trimite mesajul JOIN către nodurile existente. Acestea își vor actualiza succesorii și predecesorii.
- mecanismul de rutare a cheilor. Dacă cheia nu aparține nodului curent se transmite succesorilor până la găsirea nodului responsabil de aceasta.
- mecanismul de stabilizare. Nodurile trebuie să verifice în permanență dacă succesorii și predecesorii sunt valizi.
- implementarea funcționalităților DHT. Adăugarea, ștergerea și găsirea cheilor. Inclusiv mecanismul de transfer a cheilor.
- optimizarea cu finger table (o tabelă cu m intrări care ne ajută să facem salturi mai mari de la un nod la altul până când un nod află că o cheie este stocată în nodul imediat succesor).
- utilizarea bazelor de date pentru stocarea perechilor cheie-valoare.

5 Concluzii

Pot fi aduse îmbunătățiri aplicației precum:

- Reținerea unei liste de noduri de backup în cazul căderii simultane a mai multor noduri care se află "în linie"
- Salvarea cheilor în mai multe locații pentru restituirea acestora în cazul deconectării neașteptate a nodurilor

Sunt multe scenarii de utilizare în viața reală a DHT bazat pe Chord precum rețelele peer-to-peer, sisteme de indexare și stocare(BitTorrent), implementări blockchain etc.

6 Bibliografie

1. <https://edu.info.uaic.ro/computer-networks/cursulaboratorul.php>
2. <https://edu.info.uaic.ro/computer-networks/files/NetEx/S12/ServerConcThread/servTcpConcTh2.c>
3. <https://courses.cs.duke.edu/cps212/spring15/15-744/S07/papers/chord.pdf>
4. <https://app.diagrams.net/>
5. https://www.overleaf.com/learn/latex/Code_listing
6. [https://en.wikipedia.org/wiki/Chord_\(peer-to-peer\)](https://en.wikipedia.org/wiki/Chord_(peer-to-peer))
7. <https://docs.openssl.org/master/>
8. <https://www.sqlite.org/docs.html>