

Tema de casa APD – map reduce

Iorga Elena, grupa 1411B

Cuprins

1.Prezentarea problemei.....	3
2.Modul de rezolvare	3
3.Prezentarea implementării.....	4
4.Concluzii	9
5.Bibliografie.....	9

1. Prezentarea problemei

Colecțiile sunt uneori atât de mari încât nu putem face construirea indexilor eficient pe o singură mașină, acesta fiind cazul și pentru World Wide Web care folosește indexarea distribuită pentru a face asta. Rezultatul procesului este distribuit pe mai multe mașini, fie după termen fie după document.

Astfel, trebuie găsită o metodă care să proceseze datele de intrare (care pot fi foarte mari) și a cărei performanță în timp să fie cât mai bună. O astfel de soluție este reprezentată de folosirea procesării distribuite care, deși este complexă, este necesară. Dacă unul dintre modulele folosite eșuează, atunci programul în sine va continua să lucreze, reprezentând încă un mare avantaj al acestei soluții. În această lucrare, vom rezolva problema indexării după termen.

2. Modul de rezolvare

Pentru a rezolva problema descrisă mai sus este o aplicație a algoritmul MapReduce. Acesta se bazează pe faptul că distribuim munca în bucăți mai mici pentru a fi mai ușor de rezolvat într-un timp mai scurt.

Vom avea o etapă de mapare unde un nod cu rolul de coordonator împarte problema „originală” în sub-probleme și le distribuie către workeri pentru procesare. Apoi, una de reducere în care rolul de coordonator colectează soluțiile sub-problemelor (ce au facut workerii) și le combină pentru a obține rezultatul final al procesării dorite.

Astfel, o figură care ar putea reprezenta vizual modul de rezolvare al problemei este:

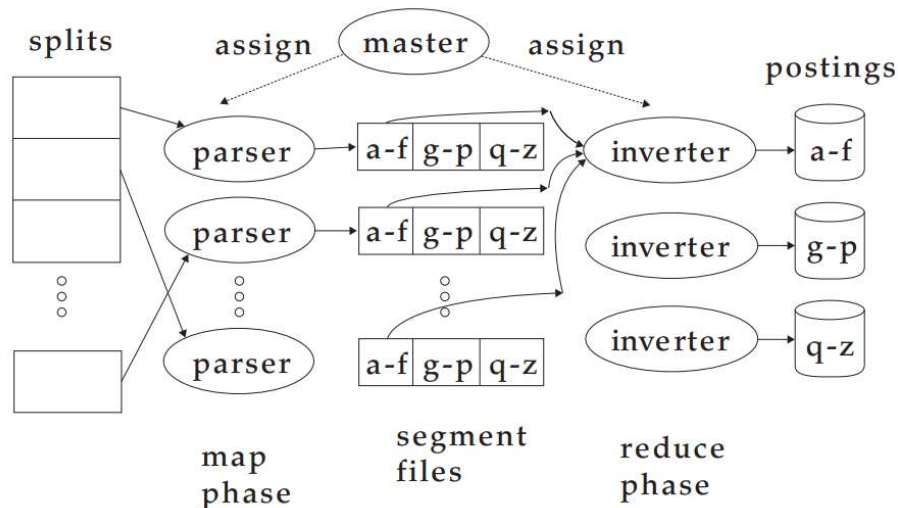


Fig. 2.1. Reprezentarea grafică a algoritmului

Pentru implementare, s-a folosit limbajul C++ cu MPI pentru a face procesarea paralelă. Programul a fost rulat pe un sistem de operare Windows, folosind Visual Studio Code cu extensia WSL pentru Debian.

3. Prezentarea implementării

În primul rând, s-a adăugat setul de date într-un director numit „date” în care sunt toate fișierele de intrare, 25 la număr.

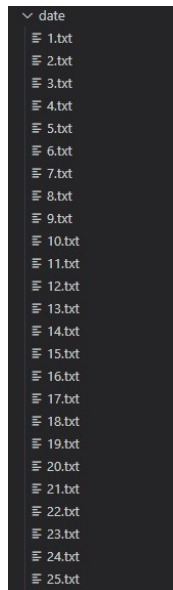


Fig. 3.1 Fișierul cu date de intrare

Apoi, din acest director au fost extrase care sunt numele tuturor fișierelor și adăugate într-un vector folosind această funcție:

```
std::vector<std::string> getFiles()
{
    std::vector<std::string> files;
    struct dirent *entry;
    DIR *dir = opendir("./date");

    if (dir == NULL)
    {
        return files;
    }
    string str1 = ".";
    string str2 = "..";
    while ((entry = readdir(dir)) != NULL)
    {
        cout << entry->d_name << endl;
        if (entry->d_name != str1 && entry->d_name != str2)
        {
            files.push_back(entry->d_name);
        }
    }
    closedir(dir);

    return files;
}
```

Fig. 3.2 Funcția folosită pentru a extrage fișierele

Folosind aceste fișiere, procesul cu rank-ul 0 care în cazul nostru va prelua rolul de master, va delega câte un fișier celorlalte procese (workerilor) pentru a putea face maparea.

```

if (rank == 0) // master
{
    int next_task = 0;
    while (next_task < 25)
    {
        int j = next_task + 1;
        const char *message = fls[next_task].c_str();
        MPI_Send(message, 9, MPI_CHAR, j, 1, MPI_COMM_WORLD);
        next_task++;
    }
    MPI_Recv(&mesaj, 1, MPI_INT, 25, 0, MPI_COMM_WORLD, &status);
    if (mesaj == 1)
    {
        map<std::string, vector<pair<string, int>>, less<string>> finallist;
        finallist = read_word_by_word_final("beforeReduce.txt", finallist);
        string stx = "\x02";
        string etx = "\x03";
        string eoh = "\x01";
        for (auto it = finallist.begin(); it != finallist.end(); it++)
        {
            ofstream outfile;
            outfile.open("final.txt", ios_base::app);
            outfile << "<" << it->first << ",{";
            for (auto pair : it->second)
            {
                if (pair.first != stx && pair.first != etx && pair.first != eoh)
                    outfile << pair.first << ":" << pair.second << ", ";
            }
            outfile << "}" << std::endl;
        }
    }
}

```

Fig. 3.3 Procesul master

Masterul ia fișierul din vector unde este sub forma de string și îl parsează în char* pentru a îl putea trimite la alte procese. După, așteaptă mesajul cum că un proces a terminat de scris în fișier (sunt șanse ca și celelalte să fi terminat), pentru a aștepta să nu existe posibilitatea mare de a parse fișierul înainte ca ceilalți să termine de scris în el și după se apucă de etapa de reducere.

Aceasta constă în a lua datele procesate de workeri care au fost scrise în fișierul *beforeReduce.txt* și a le adăuga într-un map format dintr-un string (cuvântul din fișier) și un vector de perechi care conține fișierul în care se află cuvântul și numărul de apariții. Pentru a prelua datele se folosește funcția *read_word_by_word_final* care procesează fișierul și returnează map-ul.

```

map<std::string, vector<pair<string, int>>> read_word_by_word_final(string filename, map<std::string, vector<pair<string, int>>, less<string>> finallist)
{
    ifstream file;
    string word;
    file.open(filename.c_str());
    while (file >> word)
    {
        string firstWord = word;
        string space = " ";
        file >> word;
        pair<string, int> pair;
        pair.first = word;
        file >> word;
        // cout << "incercare " << word << endl;
        pair.second = stoi(word);

        auto it = finallist.find(firstWord);
        if (it != finallist.end())
        {
            auto list = it->second;
            auto it_pair = std::find_if(list.begin(), list.end(), [&](const std::pair<std::string, int> &p)
            { return p.first != pair.first; });
            if (it_pair != list.end())
            {
                (*it_pair).second += 1;
            }
            else
            {
                vector<std::pair<string, int>> newValue;
                std::pair<string, int> par;
                par.first = pair.second;
                par.second = stoi(pair.first);
                newValue.push_back(par);
                finallist[firstWord] = newValue;
            }
        }
        else
        {
            vector<std::pair<string, int>> newValue;
            std::pair<string, int> par;
            par.first = pair.second;
            par.second = stoi(pair.first);
            newValue.push_back(par);
            finallist[firstWord] = newValue;
        }
        finallist[firstWord].push_back(pair);
    }
    file.close();
    return finallist;
}

```

Fig. 3.4 Funcția de citire din fișier a rezultatelor de la worker

Se extrag datele din fișier care știm că sunt puse în ordine de care avem nevoie de workeri, iar pe măsura ce se iau elementele de tip <cuvent, {fișier, numarAparitii}> se verifică dacă un cuvânt a fost adăugat deja în noua listă. Dacă da, dacă deja a fost adăugat și fișierul respective, se incrementează numărul de apariții. Dacă a fost adăugat cuvântul dar nu din fișierul curent, se adaugă și fișierul curent cu numărul de apariții la termen. Dacă nu a fost adăugat deloc până atunci, se adaugă simplu în listă și cuvântul cu fișierul și numărul de apariții.

Astfel, de obține un fișier final denumit *final.txt* unde vor fi datele finale sub forma <term, {docId1:count1, docId2:count2, ...}>.

Pentru procesele care sunt diferite de master, acestea este asignat rolul de mapare.

```

else
{
    char file_name[9];
    MPI_Recv(&file_name, 9, MPI_CHAR, 0, 1, MPI_COMM_WORLD, &status);
    cout << "Process " << rank << " received file name " << file_name << endl;
    map<std::string, int> numbers;

    numbers = read_word_by_word(file_name, numbers);

    pair<string, map<std::string, int, less<string>>> pair1;
    pair1.first = file_name;
    pair1.second = numbers;

    list<pair<string, pair<string, int>>> firstList;
    for (auto it = numbers.begin(); it != numbers.end(); it++)
    {
        pair<string, pair<string, int>> newElement;
        newElement.first = pair1.first;
        newElement.second.first = it->first;
        newElement.second.second = it->second;
        firstList.push_back(newElement);
    }
    list<pair<string, pair<string, int>>> secondList;
    for (auto it = numbers.begin(); it != numbers.end(); it++)
    {
        pair<string, pair<string, int>> newElement;
        newElement.first = it->first;
        newElement.second.first = pair1.first;
        newElement.second.second = it->second;
        secondList.push_back(newElement);
        ofstream outfile;
        outfile.open("beforeReduce.txt", ios_base::app);
        string space = " ";
        string nothing = "";
        if (it->first != space && it->first != nothing && pair1.first != space)
            outfile << it->first << " " << pair1.first << " " << it->second << " ";
    }
    int mesaj = 1;
    MPI_Send(&mesaj, 1, MPI_INT, 0, 0, MPI_COMM_WORLD);
}

```

Fig. 3.5 Codul pentru workeri

Fiecare worker practic citește din fișierul de intrare care i-a fost asignat cu funcția `read_word_by_word` :


```

map<std::string, int, less<string>> read_word_by_word(string filename, map<std::string, int, less<string>> numbers)
{
    ifstream file;
    string word;
    string path = "./date/";
    cout << path + filename.c_str();
    file.open(path + filename.c_str());
    while (file >> word)
    {
        word.erase(remove(word.begin(), word.end(), '.'), word.end());
        word.erase(remove(word.begin(), word.end(), ','), word.end());
        word.erase(remove(word.begin(), word.end(), '('), word.end());
        word.erase(remove(word.begin(), word.end(), '\\'), word.end());
        word.erase(remove(word.begin(), word.end(), '\"'), word.end());
        word.erase(remove(word.begin(), word.end(), '\"'), word.end());
        word.erase(remove(word.begin(), word.end(), ';'), word.end());
        transform(word.begin(), word.end(), word.begin(), ::tolower);
        if (numbers[word] == 0)
            numbers[word] = 1;
        else
        {
            numbers[word] = numbers[word] + 1;
        }
    }
    file.close();
    return numbers;
}

```

Fig. 3.6 Funcția de citire din fișierele inițiale

Astfel, se iau toate cuvintele din fișier și se numără. Pentru a avea un rezultat cât mai bun, cuvintele sunt luate fără anumite elemente de punctuație pentru a nu parsă de exemplu cuvânt apple ca fiind diferit de apple, . Se returnează un map de string și int care este ordonat după string.

Apoi, workerii iau și mapează cuvintele acestea specificând și fișierul din care au parsat datele, rezultând o pereche dintre fișier și map-ul de cuvinte cu aparițiile lor sub forma de <docId, {term1:count}>. Pe acesta trebuie să îl inverseze sub forma <term, {docId:count}> și face acestea reordonând datele într-o listă de perechi de forma cerută. Acestea sunt adăugate într-un fișier numit *beforeReduce.txt* pentru a putea fi folosite după. În acest fel, fișierul va avea date sub această formă:

```

obscurely 21.txt 1 guarantee 5.txt 4 obscurities 21.txt 1 hidden 8.txt 7 furious 12.txt 1 other 23.txt 90 districts 1.txt 26 morally 14.txt 5 elements 16.txt 17 held 17.txt 4 distrusts 1.txt 1
whisper 4.txt 1 concluding 25.txt 6 discussion 18.txt 1 disembark 18.txt 1 coats 10.txt 2 rushes 9.txt 3 hearkened 24.txt 1 automat 3.txt 1 crops 15.txt 1 furiously 12.txt 1 feet--i 11.txt 1
bumpkin 19.txt 1 disembodied 18.txt 1 passed 22.txt 4 coax 10.txt 1 whispered 4.txt 1 others 23.txt 38 lesser 6.txt 1 bumpy 19.txt 1 for--eternity? 20.txt 1 disturb 1.txt 1 disgorged 18.txt 1
morals 14.txt 33 heart 24.txt 12 others: 23.txt 3 more 14.txt 49 conclusio 25.txt 1 obscurity 21.txt 2 moreover 14.txt 6 eyed 13.txt 1 discomfort 2.txt 1 morose 14.txt 1 passes 22.txt 6
disconcerted 2.txt 1 automatic 3.txt 5 lesson 6.txt 1 discontent 2.txt 2 cock 10.txt 2 hell 17.txt 2 hide 8.txt 5 for? 20.txt 1 cross 15.txt 6 conclusion 25.txt 55 forbade 20.txt 1 discover 2.
txt 1 otherwise 23.txt 3 cockcrow 10.txt 2 heavenly 24.txt 2 fell 11.txt 19 discovered 2.txt 3 disturbance 1.txt 1 bun 19.txt 1 hello 17.txt 4 discreet 2.txt 1 furred 12.txt 1 automatics 3.txt
2 conclusion* 25.txt 1 ought 23.txt 19 force-fed 20.txt 1 our 23.txt 75 sad 9.txt 1 forced 20.txt 3 extravagant 7.txt 1 ours 23.txt 6 forcing 20.txt 1 passing 22.txt 7 codes 10.txt 1 hideous 8.
txt 1 automats 3.txt 1 extreme 7.txt 2 cross- 15.txt 1 automobile 3.txt 2 whistling 4.txt 1 furnace 12.txt 3 helmet 17.txt 3 safe 9.txt 3 guaranteed 5.txt 1 disturbing 1.txt 2 disguises 18.txt
1 coevals 10.txt 1 eyelashes 13.txt 1 automobiles 3.txt 1 most 14.txt 24 hides 8.txt 5 guaranteeing 5.txt 1 furnish 12.txt 1 observe 21.txt 1 said 9.txt 11 heavens 24.txt 1 discuss 2.txt 3
disgust 18.txt 1 elevation 16.txt 1 observed 21.txt 4 disgusts 18.txt 1 bunch 19.txt 3 coffee 10.txt 3 help 17.txt 22 observes 21.txt 1 eleventh 16.txt 3 obstinately 21.txt 1 guard 5.txt 1 eyes
13.txt 86 firearms 20.txt 1 fellow 11.txt 4 guarded 5.txt 1 cross-examination 15.txt 4 disturbingly 1.txt 1 obtain 21.txt 2 motion 14.txt 1 foreign 20.txt 2 helpful 17.txt 1 passions 22.txt 2
eloquence 16.txt 7 dish 18.txt 1 hiding 8.txt 1 height 24.txt 1 obvious 21.txt 2 let 6.txt 4 auxiliary 3.txt 1 disturbs 1.txt 1 discuss? 2.txt 1 cross-examined 15.txt 1 hie 8.txt 1 conclusion-
25.txt 1 dishes 18.txt 3 white 4.txt 1 eyes---- 13.txt 1 furniture 12.txt 2 disheveled 18.txt 1 ourselves 23.txt 18 saint 9.txt 1 further 12.txt 6 out 23.txt 20 motive 14.txt 12 held 24.txt 2
fasclock 20.txt 1 help 24.txt 9 hope 24.txt 52 outwired 23.txt 1 hopeformed 24.txt 1 happy 24.txt 1 hope 24.txt 57 hope? 24.txt 1 forest 20.txt 3 guidance 5.txt 2 eyes--a 13.txt 1 follow--for

```

Fig. 3.7 Captură din fișierul *beforeReduce.txt*

În ordinea pe care o așteaptă coordonatorul : cuvânt, fișier, numărul de apariții.

În final, după ce master-ul face etapa de reducere și adaugă cuvintele, rezultă fișierul *final.txt* care va avea rezultatul programului:


```

15012 <drest,{16.txt:1,}
15013 <drew,{4.txt:2,8.txt:2,12.txt:9,20.txt:4,11.txt:12,13.txt:9,7.txt:30,1.txt:1,18.txt:2,2.txt:1,15.txt:12,}
15014 <drew,{8.txt:1,}
15015 <dreyfus,{5.txt:1,}
15016 <dribble,{9.txt:1,18.txt:1,}
15017 <dried,{17.txt:2,20.txt:1,11.txt:1,13.txt:1,7.txt:1,16.txt:2,1.txt:16,18.txt:3,2.txt:7,15.txt:2,}
15018 <dries,{5.txt:1,1.txt:1,18.txt:1,2.txt:2,15.txt:1,}
15019 <drift,{8.txt:1,12.txt:2,20.txt:2,13.txt:3,16.txt:1,1.txt:6,18.txt:1,2.txt:4,}
15020 <drift-wood,{7.txt:1,}
15021 <drifted,{12.txt:2,20.txt:4,13.txt:5,7.txt:2,1.txt:1,2.txt:2,15.txt:1,}
15022 <drifting,{12.txt:2,20.txt:1,11.txt:1,13.txt:3,1.txt:1,18.txt:3,2.txt:1,15.txt:1,}
15023 <drifts,{1.txt:1,18.txt:2,2.txt:4,15.txt:1,}
15024 <driftwood,{7.txt:2,2.txt:1,}
15025 <drill,{6.txt:1,13.txt:1,18.txt:6,2.txt:2,}
15026 <drilled,{2.txt:2,}
15027 <drilling,{18.txt:1,}
15028 <drills,{16.txt:1,18.txt:2,}
15029 <drily,{1.txt:1,}
15030 <drink,{4.txt:1,9.txt:2,8.txt:1,21.txt:8,6.txt:4,17.txt:1,12.txt:4,20.txt:8,13.txt:11,7.txt:10,16.txt:1,1.txt:16,18.txt:4,2.txt:18,15.txt:11,}
15031 <drink!,{20.txt:1,}
15032 <drink-,{20.txt:1,}
15033 <drink},{20.txt:1,13.txt:1,}
15034 <drinkable,{2.txt:1,}
15035 <drinker,{6.txt:1,15.txt:1,}
15036 <drinkers,{9.txt:1,8.txt:1,15.txt:1,}
15037 <drinkin,{15.txt:1,}
15038 <drinking,{4.txt:1,9.txt:2,8.txt:1,21.txt:4,6.txt:1,5.txt:1,12.txt:2,20.txt:5,13.txt:2,7.txt:5,16.txt:2,1.txt:6,18.txt:2,2.txt:7,15.txt:2,}
15039 <drinking!,{20.txt:1,}
15040 <drinking?,{4.txt:1,15.txt:1,}
15041 <drinks,{9.txt:2,12.txt:1,16.txt:1,1.txt:1,18.txt:2,2.txt:2,15.txt:2,}
15042 <drip,{20.txt:1,13.txt:1,1.txt:1,18.txt:3,}
15043 <dripping,{12.txt:1,20.txt:2,13.txt:1,7.txt:1,1.txt:1,18.txt:2,2.txt:9,}
15044 <dripping-wet,{2.txt:1,}
15045 <drippings,{1.txt:1,}
15046 <drips,{9.txt:2,1.txt:5,}
15047 <drive,{4.txt:1,8.txt:2,6.txt:2,12.txt:2,20.txt:4,11.txt:6,13.txt:20,7.txt:5,16.txt:1,1.txt:5,18.txt:3,2.txt:1,15.txt:5,}
15048 <drive--and,{13.txt:1,}
15049 <driven,{9.txt:1,22.txt:1,23.txt:1,8.txt:4,5.txt:2,12.txt:3,24.txt:2,20.txt:1,11.txt:2,7.txt:6,13.txt:2,16.txt:4,1.txt:1,18.txt:3,2.txt:3,15.txt:1,}

```

Fig. 3.8 Captură din fișierul final.txt

Se observă că datele sunt în ordine alfabetică întrucât s-au folosit map-uri ordonate după string-uri.

4. Concluzii

În concluzie, programul a reușit să facă ceea ce și-a propus folosind instrumentele precizate în modul de rezolvare. Acesta a rezultat un timp mult mai bun de rulare folosind procesele din MPI, un proces care delegă acțiunile și returnează datele cerute și restul care le îndeplinesc, decât ar fi rezultat un program rulat secvențial

5. Bibliografie

[1] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. *An Introduction to Information Retrieval*. Cambridge University Press, Cambridge, England, Online c 2009 Cambridge UP edition, 2009.

[2] *Proiect – MapReduce*, disciplina Algoritmi paraleli și distribuți

[3] ***, https://www.geeksforgeeks.org/map-vs-unordered_map-c/, ultima vizitare : 31/01/2023

[4] ***, <https://linuxhint.com/map-sort-key-c/>, ultima vizitare : 31/01/2023

[5] *Laborator nr. 2 – Introducere în MPI(1)*, disciplina Algoritmi paraleli și distribuți

[6] Michael Kleber. *The MapReduce paradigm*.

<https://sites.google.com/site/mriap2008/lectures>, January 2008, ultima vizitare : 31/01/2023