

# Paradigma Map-Reduce (Programare Paralela)

---

## Structuri de Date

### Rezultatul Mapper-ilor

Fiecare thread Mapper produce un rezultat, acesta fiind un vector de perechi **word - fileID**. De vreme ce acelasi Mapper poate colecta acelasi cuvânt din fisiere diferite, **MapperResult** nu se poate implementa drept un dicționar **word - fileID**, ci ca un vector, în care aceste perechi să fie unice.

### WordList-ul

WordList-ul implementat de mine folosește două structuri de date (intermediare):

- Un vector de dicționare
- Un vector de liste

Cate un dicționar și cate o listă pentru fiecare literă din alfabet.

## Variabile Partajate

În main, se va crea o instanță a unei clase **SharedVariables** care reține toate valorile partajate între thread-uri. La crearea unui thread, **SharedVariables** de construcția argumentului thread-ului, fiecare thread primind referințe la aceste variabile partajate. În funcție de tipul thread-ului creat, se vor referenția variabile diferite (thread-urilor Mapper/Reducer li se vor da doar referințele de care au nevoie, nu mai mult).

Variabilele partajate între thread-uri conțin atât datele de interes (numele fisierele de intrare, rezultatele Mapper-ilor ...), cât și primitivile de sincronizare din **pthread.h** (mutex-uri, o variabilă condițională și o barieră).

Această clasă (**SharedVariables**) permite crearea cu ușurință a argumentelor thread-urilor.

În plus, tot **SharedVariables** se ocupă și de:

- Alocarea și dealocarea memoriei
- Inițializarea variabilelor partajate

## Programarea Paralela

### Problema Producatori - Consumatori

Paradigma Map-Reduce rezolvă problema producator-consumator, în ipostază mai mulți producatori - mai mulți consumatori:

- N-M: mai mulți producatori (fisiere de intrare), mai mulți consumatori (Mapperi)
- N-M: mai mulți producatori (rezultatele Mapper-ilor), mai mulți consumatori (cate un dicționar din WordList pentru fiecare literă din alfabet)
- N-M: mai mulți producatori (cate un dicționar din WordList pentru fiecare literă din alfabet), mai mulți consumatori (cele 26 de secțiuni din WordList - implicit și fisierele de ieșire corespondente - pentru

fiecare litera mica din alfabetul englez)

NOTA: Totusi, daca as fi implementat WordList-ul drept un dictionar/lista mare, si nu l-ar fi impartit pentru litera din alfabet, ultimele doua puncte din lista ar fi fost **many to one**, respectiv **one to many**.

## Thread-urile Mapper

Mapperi rezolva se rezolva problema producator-consumator, in ipostaza mai multi producatori (fisiere de intrare), mai multi consumatori (rezultatele mapper-ilor)

In main, thread-urile cu indicii in intervalul [0, argv[1]) vor reprezenta thread-urile Mapper.

Un thread Mapper:

- Va primi ca argument o referinta la clasa **MapperThread**
- Intr-o bucla infinita
  - Va pune lock pe mutex-ul cozii indicilor fisierelor de intrare
    - Daca coada nu mai contine niciun element, inseamna ca toate fisierele au fost parcurse si thread-ul curent nu mai are nimic de facut , astfel, da unlock la mutex si iese din bucla
    - Altfel, extrage un elemnent din coada, reprezentand indexul unui fisier de intrare si da imediat unlock la mutex-ul cozii. Astfel, niciun alt thread nu va avea acces la acel fisier afara de thread-ul curent, drept pentru care Mapper-ul isi poate incepe munca:
      - Parcurge fisierul, citind rand cu rand, caracter cu caracter, folosind o varabila auxiliara (**word**) pentru constructia cuvintelor. Atunci cand intalnim o litera (fie ea litera mare sau mica), transformam litera in litera mica si adaugam in **word**. In momentul in care intalnim un spatiu sau ajungem la capat de rand, adaugam **word**-ul construit intr-un **set** (**set**-ul va contine toate **cuvintele unice din fisier**)
      - Adaugam fiecare **cuvant** din **set**-ul precedent construit, alaturi de **ID**-ul fisierul parcurs, la rezultatul mapper-ului cu indicele egal cu ID-ul thread-ului mapper. De vreme ce inserarea in rezultatele mapper-ilor se face in paralel la indici diferiti (fiecare thread mapper insereaza doar in bucata da rezultat), nu mai este nevoie sa protejam aceasta variabila printr-un mutex
  - Cand a iesit din bucla (coada indicilor fisierelor de intrare este vid), asta inseamna ca thread-ul curent nu mai are nimic de facut, adica si-a terminat cu succes munca, drept pentru care incrementeaza o variabila care numara cati Mapperi si-au indeplinit rolul (variabila partajata atat intre toate thread-urile, atat Mapper, cat si Reducer, variabila la al cerui access este protejat printr-un mutex). In momentul in care valoarea acestei variabile pentru numarul de Mapperi finalizati este egala cu numarul total de Mapper, trimite thread-urilor Reduce un semnal de **broadcast** printr-o variabila conditionala (**pthread\_cond\_t**), **broadcast** care anunta ca toti Mapperi s-au terminat, iar Reducerii pot incepe

## Thread-urile Reducer

Reduceri rezolva problema producator-consumator, in ipostaza **mai multi producator - mai multi consumatori**:

- Mai intai: mai multi producatori (rezultatele Mapper-ilor) si mai multi consumatori (cate un dictionar in WordList pentru fiecare litera din alfabet)

- Apoi: mai multi producatori (cate un dictionar pentru fiecare litera din alfabet) si mai multi consumatori (pentru fiecare dictionar: cate un vector si un fisier de iesire).