

Tema de casă #2 – MPI

Responsabil de temă: Mihail Ionescu (mihai.ionescu@cs.pub.ro)

Data publicării: 07.11.2012

Data ultimei modificări a enunțului: 07.11.2012

Termenul de predare: 24.11.2012

Obiective

După realizarea acestei teme de casă studentul va fi capabil să:

- realizeze o aplicație bazată pe MPI.
- gestioneze eficient o structură de date mai complicată.

Cunoștințele necesare rezolvării acestei teme de casă:

- Programare în C, structuri..

Enunțul problemei

Realizați o aplicație de tip *Map/Reduce* bazată pe MPI. Într-o prezentare simplificată, Map/Reduce este o paradigmă generală de programare în sisteme distribuite, în care operația care trebuie executată să împartă în două faze. Faza 1, Map, în care un set de procese (numit mappers) împartă datele de intrare în părți egale. Fiecare mapper va trimite către un set de procese numit reducers datele pentru care este responsabil. Fiecare proces reduce implementează semantica aplicației și returnează datele procesului mapper care le-a apelat. Procesele mapper vor agrega datele și vor trimite înapoi către procesul principal datele agregate. La rândul lui, procesul principal agreghează toate aceste date și prezintă rezultatul final.

Concret, aplicația va trebui să numere cuvintele dintr-un text mare și să prezinte pentru fiecare cuvânt de câte ori apare, ordonat după frecvență. Procesul principal va citi un fișier de configurare care va specifica numărul de procese mapper disponibile și pentru fiecare mapper, câte procese reduce. Apoi, fiecare proces map va lua partea corespunzătoare a fișierului de care este responsabil, va crea în mod dinamic procesele reduce și va trimite fiecărui proces reduce datele pentru care este responsabil. Fiecare proces reduce va calcula frecvența cuvintelor pe partea sa de date și va returna aceste date, procesului mapper părinte. Acesta agreghează toate aceste date și le trimite înapoi la procesul părinte, care le agreghează din nou, le sortează după frecvență și le prezintă ca output într-un fișier.

Fișierul de input poate fi orice text. Testele vor trebui făcute cu un text mare, <http://www.gutenberg.org/files/2600/2600.txt>. Fișierul de output trebuie să fie de forma: cuvânt tab frequency pe fiecare linie, sortate după frequency.

Fisierul de configurare trebuie sa specifice numarul de mappers, numarul de reducers pentru fiecare mapper, numele fisierului de intrare si numele fisierului de iesire. Se presupune ca fiecare proces mapper are access la fisierul de intrare.

Toate testele vor fi facute pe un calculator Ubuntu 11.04 cu MPICH2. Tema va contine rezultatul masuratorilor (timpul de executie) variind numarul de mappers de la 2 la 4 din 1 in 1 si numarul de reducers de la 4 la 10 din 2 in 2.

IMPORTANT: Notarea se va face pe curba, unde programele care merg cel mai repede vor lua punctaj maxim, iar punctajul va scadea pe masura ce performantele scad. Timpii considerati vor fi aceia obtinuti pe masina de test.

In cazul in care programul nu merge (sau merge foarte prost) pe calculatorul de test, studentul poate arata ca merge pe calculator personal (laptop, etc). In acest caz, se vor scadea 10 puncte.

Tema va putea fi predată numai o saptamana dupa deadline, caz in care se vor scadea 5 puncte pe zi. Dupa o saptamana nu mai poate fi predată.

Fisierele care contribuie la rezolvarea temei trebuie **OBLIGATORIU** împachetate într-o arhiva de tip **'zip'**, cu numele 'Grupa_NumePrenume_TemaX.zip' (de exemplu, studenta Stan Sonia de la grupa 341C3 va trimite pentru tema 2 o arhiva cu numele 341C3_StanSonia_Tema1.zip). Aceasta arhiva trebuie **OBLIGATORIU** sa conțină și un fișier text README în care se explică soluția considerată în rezolvarea problemei.