

Relatório Trabalho 2 - Programação Paralela

Pedro Amaral Chapelin

Pedro Willian Aguiar

Introdução

Para esse trabalho, foi pedido para que implementássemos o algoritmo dos **K vizinhos mais próximos**. Esse problema da programação consiste no seguinte: São duas matrizes de entrada, que representam pontos e suas coordenadas. Cada linha representa um ponto, e cada coluna os valores de suas dimensões (x, y, z, \dots), sendo a matriz **Q** de tamanho $nq \times d$ (nq = número de pontos em q , d = quantidade de dimensões dos pontos) e a matriz **P** de tamanho $np \times d$ (np = número de pontos em p) as matrizes de entrada, vale ressaltar que só as dimensões de cada matriz foram passadas como entradas, o conteúdo de cada uma é preenchido aleatoriamente. Além das dimensões de cada, o programa também recebe o número **k** como argumento, indicando quantos vizinhos mais próximos serão calculados para cada ponto.

A ideia principal dos **k vizinhos mais próximos** funciona da seguinte maneira: O programa deve devolver uma matriz **R** de tamanho $nq \times k$ que contenha as distâncias e os índices dos **k** pontos de **P** que são mais próximos de cada ponto em **Q**. Agora explicado como funciona o algoritmo, partimos para a implementação do mesmo. Basicamente, para cada ponto de **Q** foi necessário varrer todos os pontos de **P**, e para ir armazenando as menores distâncias calculadas utilizamos a estrutura **Heap** cedida pelo professor. Ao final de cada varredura por todos os pontos de **P** e com a Heap pronta, tudo que tínhamos que fazer era copiar o conteúdo dela para linha da matriz **R**.

Feita a implementação, agora o desafio era paralelizar esse processo todo usando a biblioteca *mpi.h*, para isso, toda a execução era baseada em nodos e processos que realizariam o trabalho dividido. Primeiramente, antes de entrar na função principal, definimos que o preenchimento das matrizes com números aleatórios seria realizado apenas no nodo 0, juntamente com a medição do tempo, utilizando a biblioteca do professor *chrono.h*. Logo após isso, realizamos a chamada da função principal:

```
// Execução principal do programa
calculaDistancias(Q, nq, P, npp, d, k, R, nq / nproc,
MPI_COMM_WORLD, processId);
```

Passamos todos os dados das matrizes juntamente com os dados necessários para paralelizar com o MPI. $nq / nproc$ significa o que cada nodo pegará da matriz **Q** para calcular, assim realizando uma divisão do processo todo pegando o número total de pontos em q (nq) pelo número total de processos passados na compilação ($nproc$). O MPI é responsável por dividir a matriz **Q** entre os nodos (função *MPI_Scatter*), enviar a matriz **P** completa para todos (função *MPI_Broadcast*), e no fim da função, juntar todos os resultados calculados paralelamente (*MPI_Gather*).

Depois de alcançado o passo do paralelismo, agora é hora de testar e anotar a diferença observada nos testes. Para testarmos, foi disponibilizado para nós um cluster de computadores, todos com as mesmas características a seguir de processamento:

Experiência:	Achar K menores com MPI								
	INFOS obtidas do programa lscpu sobre o hardware rodado no cluster								
Processador:	Intel(R) Xeon(R) CPU	E5462							
CPU MHz	2792,84								
L1d cache	256 KiB								
L1i cache	256 KiB								
L2 cache	24 MiB								

Para todas as execuções, os parâmetros internos do programa foram os mesmos:

knn-mpi 128(nq) 400000(np) 300(d) 1024(k)

Entretanto, foram realizadas 3 experiências diferentes com parâmetros do MPI e do

Cluster:

1- Rodar o programa para APENAS 1 processo MPI e medir o tempo da computação de knn

sbatch --exclusive -N 1 knn-mpi-slurm.sh

mpirun -np 1 knn-mpi 128(nq) 400000(np) 300(d) 1024(k)

2 - Rodar o programa para 4 processos MPI no mesmo host e medir o tempo da computação de knn

sbatch --exclusive -N 1 knn-mpi-slurm.sh

mpirun -np 4 knn-mpi 128(nq) 400000(np) 300(d) 1024(k)

```
mpirun -np 4 knn-mpi 128(nq) 400000(np) 300(d) 1024(k)
```

Testes	Tempo médio (s)	SpeedUp
Primeiro	~28	1.0
Segundo	~12	~2.2
Terceiro	~12	~2.2