



Programação Paralela

ARQUITETURAS E ESTRATÉGIAS DE ACELERAÇÃO

Aluno: **Gustavo de Macena Barreto**

Prof. Ricardo Menotti

Prof. Hélio Crestana Guardia

São Carlos, 9 de abril de 2023

1 Descrição do problema

Neste algoritmo, temos como objetivo realizar o cálculo de K-ésimos vizinhos próximos cujo intuito é identificar a categoria A qual este ponto se identifica. Esse tipo de algoritmo¹ é muito utilizado em redes neurais e sistemas de grafo, Nos quais é necessário realizar uma classificação genérica para esse tipo. Dessa forma, foi escolhido como base esse tipo de algoritmo por conta de sua aplicação. Seu funcionamento é relativamente simples, Tem-se como entrada: "n" pontos de classificação conhecida e por fim o ponto desejado classificar.

2 Estratégia de paralelização

Para desenvolvimento, foi inicialmente desenvolvido um algoritmo em OpenMP, no qual se tem maior facilidade e controle por meio de diretivas de compilação e critérios para separação de variáveis que podem ser utilizadas ao longo da paralelização do algoritmo. Além disso sua programação não é mais acessível e direta em relação à programação por meio de *kernels* (CUDA, openCL, SYCL e outras). Com isso, o resultado de programação não foi satisfatório tendo em vista as dependências durante a execução da ordenação(qsort), na qual custou muito tempo para a execução. Além disso, outro algoritmo que será mencionado em seguida, desenvolvido em CUDA, sofre do mesmo problema no quesito de ordenação.

A programação em openMP é baseada em pragmas, dos quais são característicos da linguagem. Esta linguagem oferece diversas funções para o gerenciamento de threads, que permite escalonar diversas operações simultâneas para uma dada tarefa, permitindo intensidade de processamento. Na imagem 1, será explicado a seguir a estratégia tomada para a paralelização.

¹O material pode ser acessado neste link: <https://github.com/MBGustav/WSCAD-2022/tree/main/A>

```

char knn_mp(Point *arr, int n, Point P, int k){
    int i, j, min_index;
    Point temp;
    char *result = (char*) malloc(sizeof(char)*k);

    #pragma omp parallel for private(i, j, min_index, temp) shared(arr, result)
    for (i = 0; i < k; i++) {
        min_index = i;
        for (j = i + 1; j < n; j++) {
            if(euclidean_distance_no_sqrt(arr[j], P) < euclidean_distance_no_sqrt(arr[min_index], P))
                min_index = j;
        }

        temp = arr[min_index];
        arr[min_index] = arr[i];
        arr[i] = temp;
    }

    //Write results --> smallest distances from "evaluate"
    #pragma omp barrier
    if (omp_get_thread_num() == 0) {
        for (i = 0; i < k; i++) {
            result[i] = arr[i].label;
        }
    }

    qsort(result, k, sizeof(char), compare_for_sort);

    char most_freq = most_frequent(result, k);
    return most_freq;
}

```

Figura 1: Algoritmo em openMP

É declarado o *parallel* for, no qual o loop principal será realizado de forma paralela, realizando o calculo de distancia entre o ponto P e o i -ésimo ponto com classificação. Vale ressaltar que nesta etapa não há condição de corrida entre as threads, visto que não há acesso concorrente entre as threads. Vale salientar que os vetores *arr* e *result* são elementos compartilhados, justamente para partilhar o acesso de dados simultaneamente entre as *threads*. Finalizada esta execução, é necessário aguardar a execução em sua totalidade a execução das distancias que preenchem os dados de *arr*, pois após este cálculos, haverá dependência entre dados para a transposição da categoria do ponto. Vale ressaltar que essa tranposição é feita justamente para obter-se qual a moda no vetor de char's result.

Nesta parte, a explicação será dividida em duas etapas: wrappers e declaração de kernels. Primeiramente, a declaração de wrappers tem como objetivo a facilitação do uso de kernels por meio da linguagem de paralisação CUDA, isso porque existe a necessidade de configuração previas do envio de kernel, tal como: alocação de memória dinâmica, total de blocos, e número de threads por bloco. Com isso, a distribuição e organização de dados se torna mais evidente ao longo da programação, facilitando assim eventual alteração ou melhorias de código , se necessário.

```

107 __global__ void count_labels(char* array, int k, int* counts) {
108     extern __shared__ int s_counts[];
109
110     int tid = threadIdx.x + blockIdx.x * blockDim.x;
111     int stride = blockDim.x * gridDim.x;
112
113     for (int i = tid; i < k; i += stride) {
114         atomicAdd(&s_counts[array[i]], 1);
115     }
116
117     __syncthreads();
118
119     for (int s = blockDim.x / 2; s > 0; s >>= 1) {
120         if (threadIdx.x < s) {
121             s_counts[threadIdx.x] += s_counts[threadIdx.x + s];
122         }
123         __syncthreads();
124     }
125
126     if (threadIdx.x == 0) {
127         atomicAdd(&counts[s_counts[0]], 1);
128     }
129 }

```

(a)

```

77 //Smallest numbers from Point - Kernel
78 __global__ void k_smallest(CUPoint *arr, CUPoint to_eval, int n, int k, char *result) {
79     int tid = threadIdx.x + blockIdx.x * blockDim.x;
80     int stride = blockDim.x * gridDim.x;
81     int i, j, min_index;
82     CUPoint temp;
83
84     for (i = tid; i < k; i += stride) {
85         min_index = i;
86         for (j = i + 1; j < n; j++) {
87             if (distance_no_sqrt(arr[j], to_eval) < distance_no_sqrt(arr[min_index], to_eval))
88                 min_index = j;
89         }
90         temp = arr[min_index];
91         arr[min_index] = arr[i];
92         arr[i] = temp;
93     }
94
95     // __syncthreads();
96     //Write results -> smallest distances from "evaluate"
97     if (tid == 0) {
98         for (l = 0; l < k; l++) {
99             result[l] = arr[l].label;
100         }
101     }
102 }
103
104 __global void count_labels(char* array, int k, int* counts) {
105 }

```

(b)

Figura 2: Implementação em Cuda - Declaração de Kernels

Na imagem 2(b) e 2(a), é apresentada a declaração por meio da linguagem de programação CUDA. Nessa etapa é realizada a compilação direta para linguagem de programação voltada para placas de vídeos da NVidia. O primeiro algoritmo apresenta um contador de categorias cujo intuito é contabilizar a categoria que mais se repete, tendo em vista a quantidade total daquele elemento, isto é, a moda daquele grupo representado pelo vetor nomeado de "array". Já no segundo, é declarado o Kernel que tem o papel principal de classificar os pontos mais próximos dado o ponto P. Neste inquéria é realizado de maneira paralela o cálculo de distâncias sem o uso de raízes quadráticas pois isso reduz o tempo de execução e também não será necessário para esse caso. Após o cálculo de distância é preenchido o vetor result, cuja função é acumular os pontos mais próximos na ordem de k, no qual será utilizado pelo primeiro algoritmo mencionado anteriormente.

Conforme mencionado anteriormente, os wrappers tem como objetivo facilitar e evidenciar características na programação na placa de vídeo. É possível notar funções de transmissão de dados de host para device e vice-versa. Além disso, existem alocações

```

133 char most_frequent(char *d_array, int k) {
134     const int num_labels = 256; // assuming ASCII encoding
135     int* d_counts;
136     cudaMalloc(&d_counts, num_labels * sizeof(int));
137     cudaMemcpy(d_counts, 0, num_labels * sizeof(int));
138
139     // char* d_array;
140     // cudaMalloc(&d_array, k * sizeof(char));
141     // cudaMemcpy(d_array, array, k * sizeof(char), cudaMemcpyHostToDevice);
142
143     const int threads_per_block = 256;
144     const int num_blocks = (k + threads_per_block - 1) / threads_per_block;
145     const int shared_memory_size = num_labels * sizeof(int);
146
147     count_labels<<num_blocks, threads_per_block, shared_memory_size>>>(d_array, k, d_counts);
148
149     int* h_counts = new int[num_labels];
150     cudaMemcpy(h_counts, d_counts, num_labels * sizeof(int), cudaMemcpyDeviceToHost);
151
152     char most_freq = 0;
153     int most_freq_count = 0;
154     for (int i = 0; i < num_labels; i++) {
155         if (h_counts[i] > most_freq_count) {
156             most_freq_count = h_counts[i];
157             most_freq = (unsigned char) i;
158         }
159     }
160
161     delete[] h_counts;
162     cudaFree(d_counts);
163     cudaFree(d_array);
164
165     return most_freq;
166 }

```

(a)

(b)

Figura 3: Implementação em Cuda - Declaração de *Wrappers*

prévias com chamadas de prevenção de erros ao longo desses tipos de chamadas. Vale salientar que esses tipos de checagem de erros, feitos por meio de *checkCudaError*, visam evitar falha de transmissão entre dispositivos.

3 Análise da escalabilidade: esperada e obtida

Apesar da característica e desempenho notável para o cálculo de distâncias de maneira paralela, seja em openMP, seja em cuda, isso não foi o suficiente para obter um resultado satisfatório. A necessidade de ordenar os vetores e consequentemente a redução dos mesmos para um vetor exige a necessidade de um dataset relativamente grande, para que haja ganho de empenho nas plataformas. O principal gargalo observado nesse algoritmo está relacionado com a transmissão e quantidade de dados necessários. Uma possível melhoria para este algoritmo seria a locação de memória em local para melhor manejo dos dados presentes na placa de vídeo.

4 Conclusão

Por fim há de se concluir, e consequente redução dos dados em posições menores corroboraram para um desempenho não tão satisfatório com esperado. Possíveis melhorias podem ser feitas nesses códigos, das quais permitem agregar maior desempenho para o algoritmo, tal como melhor manejo e transmissão de dados, além de ordenar ou refaturar o atual modelo de ordenação desse algoritmo. Além disso seu atual funcionamento se encontra similar a CPU na qual seu funcionamento serial se tornou equivalente com acelerador.