



INSTITUTO SUPERIOR TÉCNICO

MESTRADO INTEGRADO EM ENGENHARIA ELECTROTÉCNICA E DE  
COMPUTADORES

ARQUITETURAS AVANÇADAS DE COMPUTADORES

**Paralelização e aceleração de um programa**

Guilherme ..... Maria Margarida Dias dos Reis      n.º 73099

Nuno Miguel Rodrigues Machado      n.º 74236

Lisboa, 3 de Maio de 2015

# Índice

<b>1</b>	<b>Introdução</b>	<b>1</b>
<b>2</b>	<b>Implementação no CPU</b>	<b>1</b>
<b>3</b>	<b>Implementação no GPU</b>	<b>3</b>
<b>4</b>	<b>Técnicas de aceleração e optimização</b>	<b>4</b>
<b>5</b>	<b>Secção de resultados</b>	<b>5</b>
<b>6</b>	<b>Conclusões</b>	<b>5</b>
<b>7</b>	<b>Anexos</b>	<b>5</b>

# 1 Introdução

## 2 Implementação no CPU

Inicialmente o algoritmo proposto foi implementado para correr só no CPU. Para isso, foi transcrito para o código C. O algoritmo divide-se em duas partes fundamentais:

- Criação do sinal amostrado;
- *Smoothing* do sinal amostrado;

### 2.1 Criação do sinal amostrado

O sinal a ser processado é composto pela soma de dois sinais sinusoidais mais um erro com uma amplitude máxima de 0,1. Todos os calculos são feitos em `float`, em primeiro lugar é realizado a alocação da memória de todas as variáveis necessárias para o calculo do sinal de entrada. De seguida está representado o código C em detalhe.

```
1 #define N 10000;
2 ...
3
4 int main() {
5
6     float *x, *y, *yest_cpu, *randomArray;
7     ...
8
9     /*Alocacao de memoria*/
10    x = (float *)malloc(N*sizeof(float));
11    y = (float *)malloc(N*sizeof(float));
12    yest_cpu = (float *)malloc(N*sizeof(float));
13    randomArray = (float *)malloc(N*sizeof(float));
14    ...
15    exit(0);
16 }
```

A implementação do algoritmo do sinal de entrada é composta por um ciclo que itera o numero de amostras do sinal pretendido. Em primeiro lugar é necessário gerar os valores ao sinal que vai ser processado, estes valores são gerados pela seguinte equação

$$X = i/10; \quad (2.1)$$

de seguida é gerado um valor aleatório entre 1 e -1, simulando o ruído resultante da amostragem do sinal. Este valor é gerado pela função `randn()`, o código da função está representado de seguida, é de salientar que o código foi obtido da Internet, onde este simula a função `randn()` do MatLab :

```
1 float randn()
2 {
3     float x1, x2, w, y1;
4     do
```

```

5  {
6      x1 = (float)(2.0 * rand() / RAND_MAX - 1.0);
7      x2 = (float)(2.0 * rand() / RAND_MAX - 1.0);
8      w = x1 * x1 + x2 * x2;
9  } while (w >= 1.0);
10
11  w = (float)sqrt((-2.0 * log(w)) / w);
12  y1 = x1 * w;
13  return y1;
14  }
15

```

Depois de obter os valores de  $X$  e do valor aleatório pode-se iterar os valores das amostras do sinal a ser processado. O código de seguida representa o ciclo que itera as amostras de  $X$ , do valor aleatório, `randomArray` e o sinal a ser processado,  $Y$ .

```

1  int main(){
2      for (int i = 0; i < N; ++i) {
3          x[i] = (float)i / 10;
4          randomArray[i] = randn();
5          y[i] = function((float)x[i], (float)randomArray[i]);
6      }
7      ...
8      exit(0);
9  }
10

```

## 2.2 *Smoothing* do sinal amostrado

O algoritmo de *Smooth* para o anulamento do ruído resultante da amostragem do sinal é aplicado segundo a expressão seguinte:

$$y_{est} = \sum_{i=0}^{N-1} \frac{\sum_{k=0}^{N-1} Kb(x_i, x_k) y_k}{\sum_{k=0}^{N-1} Kb(x_i, x_k)}; \quad (2.2)$$

$$Kb(x_i, x_k) = e^{-\frac{(x_i - x_k)^2}{2b^2}}; \quad (2.3)$$

A implementação do algoritmo em C baseia-se na utilização de um ciclo para o somatório exterior e outro ciclo para o somatório interior, as funções exponencial `expf` e potencia de base 2, `powf`, pertence á biblioteca, `math.h`. O código seguinte demonstra a utilização dos dois ciclos como também a das funções para o calculo do *smoothing*:

```

1  int main(){
2      float sumA, sumB;
3      ...
4
5      for (int i = 0; i < N; ++i) { //percorrer o yest
6          sumA = 0;

```

```

7     for (int j = 0; j < N; ++j) { //percorer o input dataset
8         sumA = sumA + ((expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2)))) * y[j
9     ]);
10    }
11    sumB = 0;
12    for (int j = 0; j < N; ++j) { //percorer o input dataset
13        sumB = sumB + expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2)));
14    }
15    yest_cpu[i] = sumA / sumB;
16 }
17 ...
18 exit(0);
19 }
20

```

### 3 Implementação no GPU

A implementação em GPU é dividida em 4 partes:

- Alocação da memória;
- Envio dos dados do CPU para o GPU;
- Iniciação do algoritmo em GPU;
- Envio dos dados do GPU para o CPU;

#### 3.1 Alocação da memória

No início da implementação da paralisação em CUDA é necessário alocar a memória total a ser enviada do CPU para o GPU. Neste caso também foi alocado a memória total necessária para guardar o resultado do *smoothing*. De seguida apresenta-se o código para alocação da memória do *device* que com tem o GPU:

```

1  int main(){
2      float *d_x, *d_y, *d_yest;
3      ...
4      cudaMalloc(&d_x, N*sizeof(float));
5      cudaMalloc(&d_y, N*sizeof(float));
6      cudaMalloc(&d_yest, N *sizeof(float));
7      ...
8      exit(0);
9  }
10

```

### 3.2 Envio dos dados do CPU/GPU ou GPU/CPU

Com a memória alocada, o passo seguinte é transferir os dados obtidos na secção 2.1,  $X$  e  $Y$ , para o *device*. É utilizado a função `cudaMemcpy`, que pertence à biblioteca `cuda.h`, recebe o ponteiro de destino e o ponteiro onde está a memória a ser transferida, é necessário definir a dimensão de dados a ser transferidos e por fim é necessário definir o sentido da transferência usando as seguintes mascaras, `cudaMemcpyHostToDevice` sentido do CPU para o GPU e `cudaMemcpyDeviceToHost` sentido GPU para o CPU. No código seguinte está implementado a função descrita:

```
1  int main(){
2      float *x, *y, *yest_cpu,*yest_gpu, *randomArray;
3      float *d_x, *d_y, *d_yest;
4      ...
5      /*Envio de dados para o device*/
6      cudaMemcpy(d_x, x, N*sizeof(float), cudaMemcpyHostToDevice);
7      cudaMemcpy(d_y, y, N*sizeof(float), cudaMemcpyHostToDevice);
8      ...
9      /*Envio de dados para o CPU*/
10     cudaMemcpy(yest_gpu, d_yest, N*sizeof(float), cudaMemcpyDeviceToHost);
11     ...
12     exit(0);
13 }
14
```

### 3.3 Iniciação do algoritmo em GPU

Depois de enviar todos os dados para o GPU o *kernel* está pronto para ser invocado. O *kernel* representa o código que vai ser executado pela GPU, este é definido pela declaração `__global__` antes da função C. Ver código seguinte,

```
1  __global__ void funtion_smooth(float *x, float *y, float *yest, int n){
2      ...
3  }
4
```

Com o *kernel* definido este é executado usando a seguinte configuração,

$$NomeDaFuncao <<< NB, NT >>> \quad (3.1)$$

Onde NB é o numero de blocos a ser lançados no GPU e NT o numero de *threads* por bloco.

## 4 Técnicas de aceleração e optimização

O código a ser paralisado é referente à segunda secção, *Smoothing* do sinal amostrado, do capítulo, Implementação no CPU. Analisou-se a estrutura do algoritmo e verificou-se a possibilidade de optimização e paralisação dos ciclos.

## 4.1 Optimiza  o

Analisando o algoritmo proposto identificou-se duas situa  es principais de optimiza  o, n  mero de ciclos e acesso    mem  ria. Come  ou-se ent  o por reduzir o n  mero de ciclos do algoritmo, os dois ciclos interiores podem ser reduzidos a um s  , com esta altera  o verificou-se que se podia reduzir o acesso    mem  ria. Isto   , como as vari  veis `sumA` e `sumB` calculam-se da mesma forma tirando a diferen  a de `sumA` ser multiplicada por `y[i]`. Fez-se a seguinte altera  o, a parte comum    calculada em primeira inst  ncia e o resultado    guardado numa vari  vel auxiliar, `sum`. Assim para o calculo de `sumB`    s   necess  rio aceder a cache e obter os valores `sumB` e `sum` e para o calculo de `sumA` acede de igual forma, vai    cache retirar os valores de `sumA` e `sum` e um acesso    mem  ria global, `y[j]`. O c  digo de seguida demonstra a explica  o feita anteriormente:

```
1 __global__ void funtion_smooth(float *x, float *y, float *yest, int n){
2     int i = blockIdx.x* blockDim.x + threadIdx.x;
3     float sumA=0.0, sumB=0.0, sum=0.0;
4
5     if (i < n){
6         for (int j = 0; j < n ;j++){
7             sum = (expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2)))));
8             sumA = sumA + sum* y[j];
9             sumB = sumB + sum;
10        }
11        yest[i] = sumA / sumB;
12    }
```

## 4.2 Paralisa  o dos ciclos

- Paralisa  o do ciclo externo;
- Paralisa  o do ciclo externo e interno;

Para melhor compreens  o no c  digo seguinte est   representado qual o ciclo externo e qual o ciclo interno.

```
1 for (int i = 0; i < N; ++i) { /*Ciclo Externo*/
2     sumA = 0;
3     for (int j = 0; j < N; ++j) { /*Ciclo interno*/
4         sumA = sumA + ((expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2)))) * y[j]);
5     }
6     sumB = 0;
7     for (int j = 0; j < N; ++j) { /*Ciclo interno*/
8         sumB = sumB + expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2)));
9     }
10    yest_cpu[i] = sumA / sumB;
11 }
```

4.2.1 Paralisação do ciclo externo

4.2.2 Paralisação do ciclo interno

## 5 Secção de resultados

5.1 Tempos

5.2 Resultados do *Smooth*

## 6 Conclusões

## 7 Anexos