



INSTITUTO SUPERIOR TÉCNICO

MESTRADO INTEGRADO EM ENGENHARIA ELECTROTÉCNICA E DE
COMPUTADORES

ARQUITETURAS AVANÇADAS DE COMPUTADORES

Paralisação e aceleração de um programa

Guilherme Branco Teixeira	n.º 70214
Maria Margarida Dias dos Reis	n.º 73099
Nuno Miguel Rodrigues Machado	n.º 74236

Lisboa, 1 de Junho de 2015

Índice

1	Introdução	1
2	Implementação no CPU	1
3	Implementação no GPU	3
4	Técnicas de aceleração e optimização	4
5	Secção de resultados	8
6	Conclusões	12
7	Anexos	12

1 Introdução

Pretende-se paralisar e acelerar um algoritmo de *smoothing* usando um processador gráfico como ferramenta de computação paralela. Neste relatório está demonstrado o algoritmo em CPU, as optimizações necessárias e paralisações essenciais de forma a ter os melhores resultados no GPU.

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:

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

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

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 ...
19 exit(0);
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 Optimizaç 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 for (int i = 0; i < N; ++i) {
2   for (int j = 0; j < n ;j++){
3     sum = (expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2))));
4     sumA = sumA + sum* y[j];
5     sumB = sumB + sum;
6   }
7   yest[i] = sumA / sumB;
8 }
```

4.2 Paralisa  o dos ciclos

Analisando a optimiza  o do c digo anterior verificou-se que pode ser paralisado em duas situa  es:

- 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   for (int j = 0; j < n ;j++){/*Ciclo Interno*/
3     sum = (expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2))));
4     sumA = sumA + sum* y[j];
5     sumB = sumB + sum;
6   }
7   yest[i] = sumA / sumB;
8 }
```

4.2.1 Paralisa  o do ciclo externo

Analisando o ciclo externo verificou-se que n o existe depend ncias de dados sendo por isso a primeira implementa  o a ser realizada. Assim sendo cada itera  o do ciclo externo   executada pelas *threads* dos blocos chamados. Como o GPU utilizado   composto por 1024 *threads* por bloco, quando se pretende processar N itera  es   necess rio saber o n  mero de blocos a serem chamados para isso

divide-se o número máximo de iterações pelo número máximo de *threads* por blocos. De seguida está implementado o código do *kernel* que corre no GPU e a função de chamada pelo CPU.

```

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
6      if (i < n){
7          for (int j = 0; j < n ;j++){
8              sum = (expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2)))));
9              sumA = sumA + sum* y[j];
10             sumB = sumB + sum;
11         }
12         yest[i] = sumA / sumB;
13     }
14 }
15 int main() {
16     ...
17     dim3 dimBlock(BLOCK_SIZE, 1, 1);
18     dim3 dimGrid(N / BLOCK_SIZE + 1, 1, 1);
19     funtion_smooth <<< dimGrid, dimBlock >>>(d_x, d_y, d_yest, N);
20     ...
21     exit(0);
22 }

```

O *kernel* a ser executado é composto pelo ciclo interno, ou seja, cada *thread* executa o *smoothing* de cada amostra do sinal. É de salientar que a variável *i* identifica qual a iteração que o ciclo exterior executava, para isso seguiu-se o seguinte esquema,

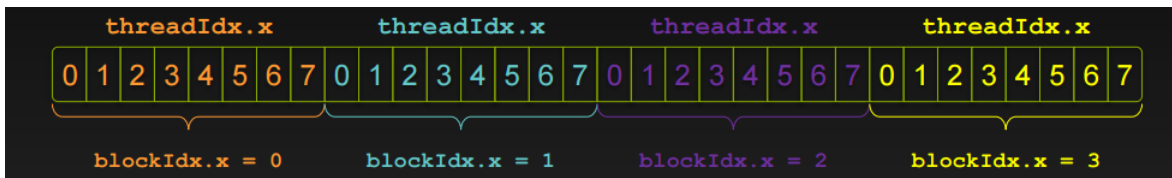


Figura 1: Esquema para indexação única.

Como o número de identificação de cada *thread* só identifica a *thread* dentro de cada bloco é necessário obter uma indexação única e igual à do ciclo exterior paralisado. Assim sendo, e analisando a figura anterior obteve-se a seguinte expressão para o calculo do de *i*, $i = \text{blockIdx.x} * \text{blockDim.x} + \text{threadIdx.x}$. Onde o *blockIdx.x* é o identificador do bloco, *blockDim.x* dimensão de cada bloco e *threadIdx.x* é o identificador de cada *thread*.

4.2.2 Paralisação do ciclo interno

Ao contrário do ciclo externo o ciclo interno tem dependência de dados dificultando assim a paralisação desse mesmo. Para a paralisação interna optou-se por dividir o número máximo de iterações

em quatro partes, reduzindo assim o tempo de processamento em cada *thread*. O código seguinte representa a implementação referida:

```

1  __global__ void funtion_smooth(float *x, float *y, float *yest, int n){
2      int i = blockIdx.x* blockDim.x + threadIdx.x;
3      int j = 0;
4      float sumA=0.0, sumB=0.0, sumA_Total=0.0, sumB_Total=0.0, sum = 0.0;
5      __shared__ float sumA_partial[BLOCK_SIZE];
6      __shared__ float sumB_partial[BLOCK_SIZE];
7
8      if ((i/MULTHREADS) < n){
9          if (!(threadIdx.x&0x02)){
10             if (!(threadIdx.x&0x01)){
11                 for (j = 0; j < (n/MULTHREADS); j++){
12                     sum = (expf(-powf((x[i]>>2] - x[j]), 2) / (2 * powf(SMOOTH, 2))));
13                     sumA = sumA + sum* y[j];
14                     sumB = sumB + sum;
15                 }
16             }else{
17                 for (j = (n/MULTHREADS); j < 2*(n/MULTHREADS); j++){
18                     sum = (expf(-powf((x[i]>>2] - x[j]), 2) / (2 * powf(SMOOTH, 2))));
19                     sumA = sumA + sum* y[j];
20                     sumB = sumB + sum;
21                 }
22             }
23         }else{
24             if (!(threadIdx.x&0x01)){
25                 for (j = 2*(n/MULTHREADS); j < 3*(n/MULTHREADS); j++){
26                     sum = (expf(-powf((x[i]>>2] - x[j]), 2) / (2 * powf(SMOOTH, 2))));
27                     sumA = sumA + sum* y[j];
28                     sumB = sumB + sum;
29                 }
30             }else{
31                 for (j = 3*(n/MULTHREADS); j < 4*(n/MULTHREADS); j++){
32                     sum = (expf(-powf((x[i]>>2] - x[j]), 2) / (2 * powf(SMOOTH, 2))));
33                     sumA = sumA + sum* y[j];
34                     sumB = sumB + sum;
35                 }
36             }
37         }
38         sumA_partial[threadIdx.x]= sumA;
39         sumB_partial[threadIdx.x]= sumB;
40         ...
41     }
42 }

```

Como se pode verificar no código anterior, existe quatro ciclos com um quarto do número máximo de iterações. Ao dividir o processamento em quatro *threads* é necessário ter uma memória partilhada:

```

1 __global__ void funtion_smooth(float *x, float *y, float *yest, int n){
2 ...
3     __shared__ float sumA_partial[BLOCK_SIZE];
4     __shared__ float sumB_partial[BLOCK_SIZE];
5 ...
6 }

```

Para que a organização de cada *thread* seja executada de forma correcta é necessário organizar as *threads* de forma coerente, assim sendo é verificado os últimos dois *bits* menos significativos da variável `thread.idx` de forma a executar o ciclo correcto, por exemplo, se a variável `thread.idx` for igual a 1 então executa o 2 ciclo.

Com os dados guardados na memória partilhada é necessário sincronizar todo de forma a obter o resultado acumulado final. Para a sincronização é utilizado a função de sistema `__syncthreads()`, que sincroniza todas as *threads* do mesmo bloco. Depois da sincronização pode-se acumular os resultados de todas as quatro *threads* para isso é utilizado a *thread* que tem os dois *bits* menos significativos a zero, ou seja, a primeira *thread* de cada grupo de quatro. De seguida está demonstrado o código:

```

1 ...
2 if(!(threadIdx.x&0x03)){
3     for(int l=0 ; l<MULTHREADS;l++){
4         sumA_Total += sumA_partial[threadIdx.x + l];
5         sumB_Total += sumB_partial[threadIdx.x + l];
6     }
7     yest[i>>2] = sumA_Total / sumB_Total;
8 }

```

5 Secção de resultados

Depois da implementação anteriormente referida foi analisado os tempos de execução

5.1 Tempos

Para o cálculo dos tempos foi usado a função `clock_gettime()` e a função `timeDiff()`. Todas as estruturas e funções não representadas estão na biblioteca `sys/time.h`.

```

1 double timeDiff(struct timespec tStart, struct timespec tEnd)
2 {
3     struct timespec diff;
4
5     diff.tv_sec = tEnd.tv_sec - tStart.tv_sec - (tEnd.tv_nsec<tStart.tv_nsec ? 1 : 0);
6     diff.tv_nsec = tEnd.tv_nsec - tStart.tv_nsec + (tEnd.tv_nsec<tStart.tv_nsec ?
7         1000000000 : 0);
8
9     return ((double)diff.tv_sec) + ((double)diff.tv_nsec) / 1e9;
10 }

```

Para o calculo do tempo de execucao em CPU foi calculado o intervalo de tempo desde a primeira iteracao do ciclo exterior até á ultima tiracao.

```

1 int main(){
2     ...
3     struct timespec timeVect[7];
4     double timeCPU;
5     ...
6     clock_gettime(CLOCK_REALTIME, &timeVect[0]);/*Inicio da contagem do tempo*/
7
8     for (int i = 0; i < N; ++i) { //percorrer o yest
9         sumA = 0;
10        for (int j = 0; j < N; ++j) { //percorer o input dataset
11            sumA = sumA + ((expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2)))) * y[j]);
12        }
13        sumB = 0;
14        for (int j = 0; j < N; ++j) { //percorer o input dataset
15            sumB = sumB + expf(-powf((x[i] - x[j]), 2) / (2 * powf(SMOOTH, 2)));
16        }
17        yest_cpu[i] = sumA / sumB;
18    }
19    clock_gettime(CLOCK_REALTIME, &timeVect[1]);/*Fim da contagem do tempo*/
20    timeCPU = timeDiff(timeVect[0], timeVect[1]);
21    printf("    ... execution took %.6f seconds\n", timeCPU);
22    ...
23    exit(0);
24 }

```

O calculo do tempo de execucao do GPU referido foi obtido somando o tempo de alocao de memoria, o tempo de transferencia de dados do CPU para o GPU, o tempo de execucao do *kernel* e a transferencia dos resultados do GPU para o CPU.

```

1 int main(){
2     ...
3     struct timespec timeVect[7];
4     double timeCPU, timeGPU[7];
5     ...
6     clock_gettime(CLOCK_REALTIME, &timeVect[1]);/*Inicializacao do contador*/
7
8     cudaMalloc(&d_x, N*sizeof(float));
9     cudaMalloc(&d_y, N*sizeof(float));
10    cudaMalloc(&d_yest, N *sizeof(float));
11    clock_gettime(CLOCK_REALTIME, &timeVect[2]);/*Tempo de alocao de memoria*/
12
13    cudaMemcpy(d_x, x, N*sizeof(float), cudaMemcpyHostToDevice);
14    cudaMemcpy(d_y, y, N*sizeof(float), cudaMemcpyHostToDevice);
15    clock_gettime(CLOCK_REALTIME, &timeVect[3]);/*Tempo de transferencia de dados CPU/
    GPU*/
16

```

```

17  dim3 dimBlock(BLOCK_SIZE, 1, 1);
18  dim3 dimGrid(N / BLOCK_SIZE + 1, 1, 1);
19
20  funtion_smooth <<< dimGrid, dimBlock >>>(d_x, d_y, d_yest, N);
21  clock_gettime(CLOCK_REALTIME, &timeVect[4]);/*Tempo de execucao do kernell*/
22
23  cudaMemcpy(yest_gpu, d_yest, N*sizeof(float), cudaMemcpyDeviceToHost);
24  clock_gettime(CLOCK_REALTIME, &timeVect[5]);/*Tempo de transferencia de dados GPU/
    CPU*/
25  ...
26  exit(0);
27 }

```

5.1.1 Paralisação do ciclo externo

De seguida está implementado uma tabela com os resultados dos tempos de execução da paralisação do ciclo externo em função do número de amostras. Também se encontra representado na tabela o valor do *speed-up* obtido com a implementação do GPU. O *speed-up* é calculado pela seguinte equação, tmp_{cpu} tempo de execução no CPU e tmp_{gpu} tempo de execução no GPU.

$$speedup = \frac{tmp_{cpu}}{tmp_{gpu}}; \quad (5.1)$$

Tabela 1: Valores dos tempos de execução no CPU e GPU para diferentes valores de amostras.

Nº Amostras	Tempo no CPU (segundos)	Tempo no GPU (segundos)	SpeedUp
1000	0,058	0,0014	40,2
10000	3,94	0,0197	200,4
20000	15,184	0,0753	201,7
30000	33,757	0,1496	225,7
40000	59,644	0,2167	275,3

De seguida também está implementado o gráfico do *speed-up* em função do número de amostras.

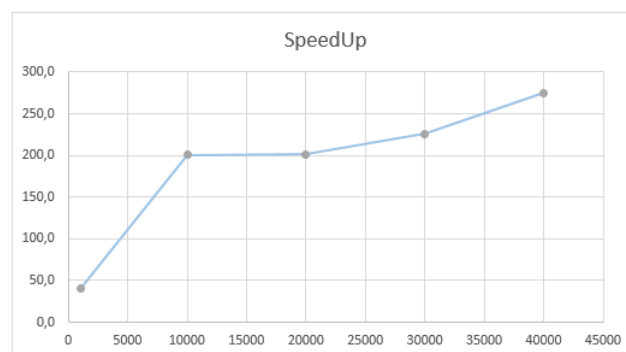


Figura 2: Gráfico do *speed-up* em função das amostragem.

5.1.2 Paralisação do ciclo externo e interno

Está representado os resultados dos tempos de execução referentes á paralisação do ciclo externo e interno. Também se encontra representado na tabela o valor do *speed-up* obtido com a implementação do GPU. O *speed-up* é calculado pela equação anterior 5.1.

Tabela 2: Valores dos tempos de execução no CPU e GPU para diferentes valores de amostras.

Nº Amostras	Tempo no CPU (segundos)	Tempo no GPU (segundos)	SpeedUp
1000	0,061	0,0015	41,7
10000	3,94	0,0580	67,9
20000	15,19	0,2287	66,4
30000	33,731	0,4822	70,0
40000	59,621	0,8121	73,4

De seguida também está implementado o gráfico do *speed-up* em função do número de amostras.

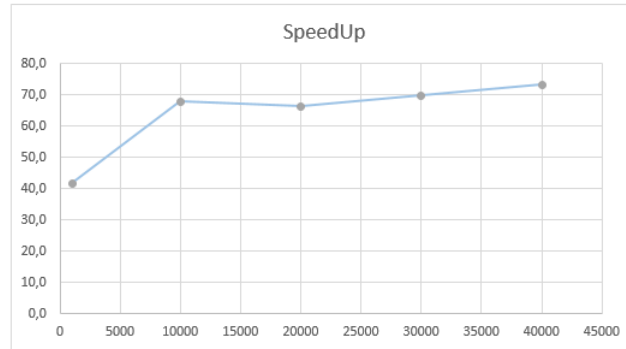


Figura 3: Gráfico do *speed-up* em função das amostragem.

5.2 Resultados do *Smooth*

Está representado nos gráficos seguintes os resultados das amostras processadas pelo CPU e GPU com 20000 amostras.

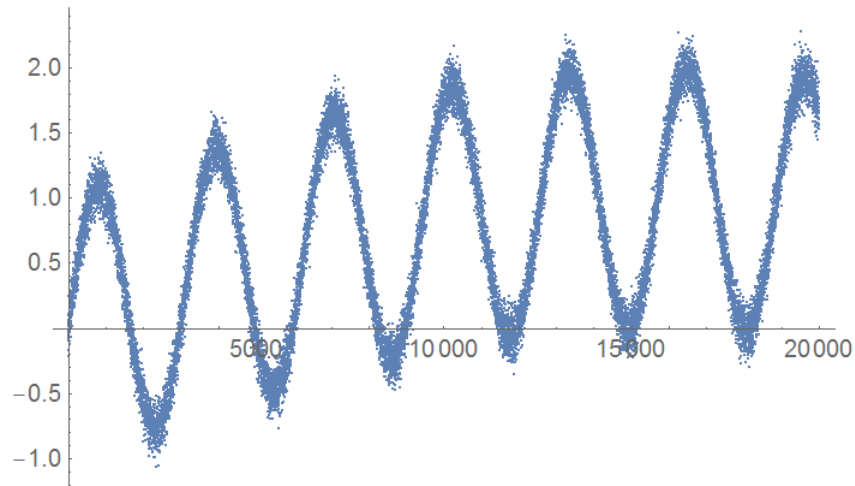


Figura 4: Gráfico do sinal amostrado com ruído.

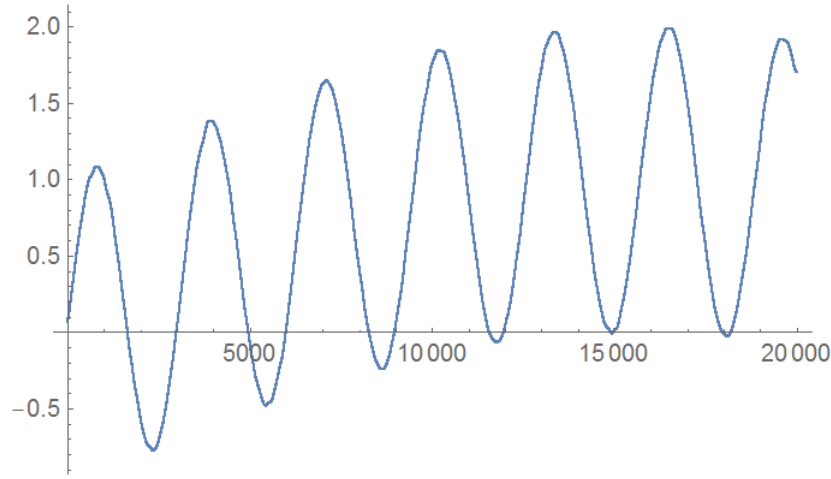


Figura 5: Gráfico do sinal amostrado com *smoothing* do ruído implementado em CPU.

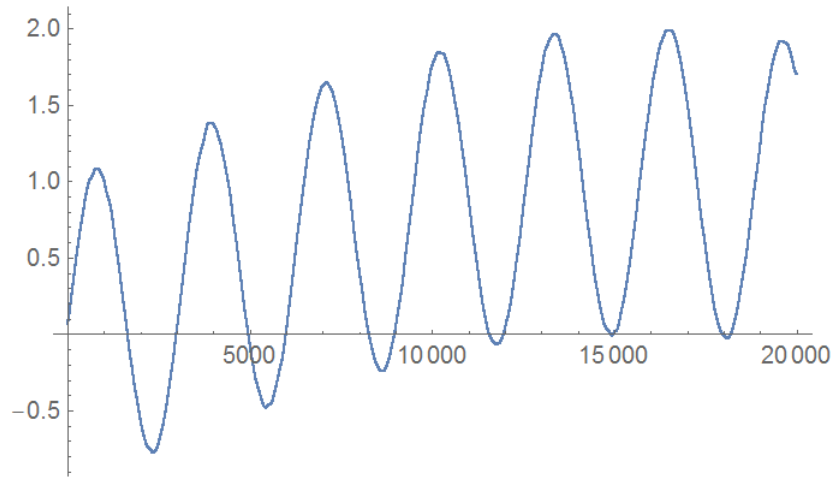


Figura 6: Gráfico do sinal amostrado com *smoothing* do ruído implementado em GPU.

6 Conclusões

A otimização escolhida para este algoritmo foi a implementação da paralisação só no ciclo externo, esta escolha é referente aos resultados dos *speed-ups* anteriormente referidos, pois existe um *speed-up* três vezes superior ao da paralisação do ciclo externo e interno. Estes resultado era de esperar devido ao *overhead* que existe na sincronização das *threads* e acesso á memória partilhada. Mesmo com os maus resultados obtidos foi importante implementar esta segunda paralisação devido à complexidade e a utilização de diversas funções do CUDA.

7 Anexos

Descrição dos ficheiros em anexo: **Smooth_1v.cu** ficheiro que contem o código da primeira paralisação, paralisação do ciclo externo. **Smooth_2V.cu** ficheiro que contem o código da segunda para-

lisação, paralisação do ciclo externo e interno.

Na pasta **Datasets** estão os resultados obtidos para 1000, 10000 e 20000 amostras.