

Raport: Mnożenie macierzy

Heorhii Lopatin (hl456366)

2 maja 2023

1 Wstęp

Celem tego zadania jest stworzenie, przetestowanie oraz porównanie wydajności różnych wersji kernela CUDA, które służą do obliczania iloczynu macierzy.

2 Opis implementacji

Zaimplementowano 5 wersji kernela. Pierwsza - przeniesienie kodu CPU na GPU. Pozostałe wersje dzielą macierz na bloki o rozmiarach K na K i liczą każdy blok za pomocą jednego wątku.

Kernel 1

```
1  __global__ void mult_gpu1(real* A, real* B, real* C, int N) {
2      int column = blockIdx.x * blockDim.x + threadIdx.x;
3      int row = blockIdx.y * blockDim.y + threadIdx.y;
4      int locA, locB;
5      int k;
6      real c = 0.0;
7      locA = N * row;
8      locB = N * column;
9      if (column < N && row < N) {
10         for (k = 0; k < N; k++) {
11             c += A[row * N + k] * B[k*N + column];
12         }
13         C[row * N + column] = c;
14     }
15 }
16
```

Kernel 2

```
1  template <uint32_t dim>
2  __global__ void mult_gpu2(real* A, real* B, real* C, int N) {
3      int start_x = blockIdx.x * dim;
4      int start_y = blockIdx.y * dim;
5
6      int iter = blockIdx.x < gridDim.x ? DIM : DIM - (DIM * gridDim.x - N);
7
8      __shared__ real OUT[dim][dim];
9      __shared__ real rowB[dim];
10
11     int i = threadIdx.x;
12
13     for (int n = 0; n < dim; n++)
14         OUT[n][i] = 0;
15     __syncthreads();
16
17     int row = start_y + i;
18     for (int k = 0; k < N; ++k) {
19         if (i < iter)
20             rowB[i] = B[k * N + start_x + i];
21         __syncthreads();
22
23         real colA = row < N ? A[row * N + k] : 0;
24         __syncthreads();
25
26         for (int n = 0; n < dim; n++) {
27             OUT[i][n] += colA * rowB[n];
28         }
29         __syncthreads();
30     }
31     __syncthreads();
32     if (row < N) {
33         for (int n = 0; n < dim; n++) {
34             C[row * N + start_x + n] = OUT[i][n];
35         }
36     }
37     __syncthreads();
38 }
39
```

Kernel 3

```
1  template <uint32_t dim>
2  __global__ void mult_gpu3(real* A, real* B, real* C, int N) {
3      int start_x = blockIdx.x * dim;
4      int start_y = blockIdx.y * dim;
5
6      int iter = blockIdx.x < gridDim.x ? DIM : DIM - (DIM * gridDim.x - N);
7
8      __shared__ real OUT[dim][dim + 1];
9      __shared__ real rowB[dim];
10
11     int i = threadIdx.x;
12
13     for (int n = 0; n < dim; n++)
14         OUT[n][i] = 0;
15     __syncthreads();
16
17     int row = start_y + i;
18     for (int k = 0; k < N; ++k) {
19         if (i < iter)
20             rowB[i] = B[k * N + start_x + i];
21         __syncthreads();
22
23         real colA = row < N ? A[row * N + k] : 0;
24         __syncthreads();
25
26         for (int n = 0; n < dim; n++) {
27             OUT[i][n] += colA * rowB[n];
28         }
29         __syncthreads();
30     }
31     __syncthreads();
32     if (row < N) {
33         for (int n = 0; n < dim; n++) {
34             C[row * N + start_x + n] = OUT[i][n];
35         }
36     }
37     __syncthreads();
38 }
39
```

Kernel 4

```
1
2  template <uint32_t dim>
3  __global__ void mult_gpu4(real* A, real* B, real* C, int N) {
4      int start_x = blockIdx.x * dim;
5      int start_y = blockIdx.y * dim;
6
7      int iter = blockIdx.x < gridDim.x ? DIM : DIM - (DIM * gridDim.x - N);
8
9      __shared__ real OUT[dim][dim + 1];
10     real rowB;
11
12     int i = threadIdx.x;
13
14     for (int n = 0; n < dim; n++)
15         OUT[n][i] = 0;
16     __syncthreads();
17
18     int row = start_y + i;
19     for (int k = 0; k < N; ++k) {
20         if (i < iter)
21             rowB = B[k * N + start_x + i];
22         __syncthreads();
23
24         real colA = row < N ? A[row * N + k] : 0;
25         __syncthreads();
26
27         for (int n = 0; n < dim; n++) {
28             OUT[i][n] += colA * __shfl_sync(0xffffffff, rowB, n);
29         }
30         __syncthreads();
31     }
32     __syncthreads();
33     if (row < N) {
34         for (int n = 0; n < dim; n++) {
35             C[row * N + start_x + n] = OUT[i][n];
36         }
37     }
38     __syncthreads();
39 }
40
41
```

Kernel 5

```
1
2 template <uint32_t dim>
3 __global__ void mult_gpu5(real* A, real* B, real* C, int N) {
4     int start_x = blockIdx.x * dim;
5     int start_y = blockIdx.y * dim;
6
7     int iter = blockIdx.x < gridDim.x ? DIM : DIM - (DIM * gridDim.x - N);
8
9     real OUT[dim];
10    real rowB;
11
12    int i = threadIdx.x;
13
14    for (int n = 0; n < dim; n++)
15        OUT[n] = 0;
16    __syncthreads();
17
18    int row = start_y + i;
19    for (int k = 0; k < N; ++k) {
20        if (i < iter)
21            rowB = B[k * N + start_x + i];
22        __syncthreads();
23
24        real colA = row < N ? A[row * N + k] : 0;
25        __syncthreads();
26
27        for (int n = 0; n < dim; n++) {
28            OUT[n] += colA * __shfl_sync(0xffffffff, rowB, n);
29        }
30        __syncthreads();
31    }
32    __syncthreads();
33    if (row < N) {
34        for (int n = 0; n < dim; n++) {
35            C[row * N + start_x + n] = OUT[n];
36        }
37    }
38    __syncthreads();
39 }
40
```

3 Wyniki

Kod przetestowano na macierzach 1000x1000. Wyniki przeprowadzonego eksperymentu przedstawiono w tabeli:

3.1 Tabela 1. Podsumowanie wyników

Wersja kodu	Konfiguracja kernela	Średni czas wykonania (ms)	Przyspieszenie względem CPU
CPU	-	4356.31±0.61	1.00
Kernel #1	8 x 1	48.27±0.01	90.25±0.02
Kernel #1	16 x 1	33.17±0.02	131.31±0.08
Kernel #1	32 x 1	26.80±0.01	162.55±0.06
Kernel #2	32 x 1	120.51±6.95	36.15±2.22
Kernel #2	64 x 1	257.80±0.45	16.90±0.03
Kernel #2	96 x 1	321.30±0.07	13.56±0.01
Kernel #2	128 x 1	321.38±0.09	13.55±0.01
Kernel #3	32 x 1	19.48±0.01	223.68±0.09
Kernel #3	64 x 1	23.30±0.05	186.96±0.46
Kernel #4	32 x 1	17.24±0.01	252.71±0.07
Kernel #5	32 x 1	10.54±0.01	413.15±0.10

4 Podsumowanie

Patrząc na wyniki, możemy wyciągnąć kilka wniosków:

- Należy dbać o możliwych konfliktach banków pamięci.
- Jeśli można uniknąć dostępu do/utworzenia obiektów w pamięci, należy to zrobić

Najszybszy był kernel w 5 wersji. Wydaje się, że przez małą liczbę zmiennych użytych w funkcji, mogą być przechowywane w rejestrach, a to oznacza, że dostęp do zawartości oraz wymiana zmiennych są szybkie. Dodatkowo, wyniki przechowywane w pamięci lokalnej, co też ułatwia dostęp.

Kernele w wersji 3,4 też były w miarę szybkie, jednak wolniejsze od 5 wersji przez dostęp do pamięci współdzielonej w każdym obrocie pętli.