

GPUゼミ (11/14)

参考

- CUDA By Example, by Jason Sanders + Edward Kandrot

(翻訳) インプレスジャパン

- GPU tutorial By 圭川

準備

手順:

- ssh to gpu computer (budou, melon, einstein, landsee.)
- Make "workspace" directory
- Do "git clone"

GPU ゼミ

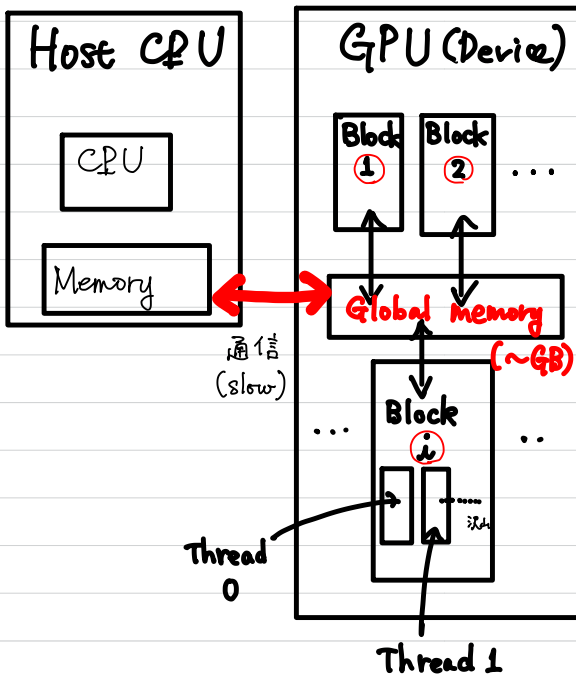
■ 基本事項

■ **CUDA** : NVIDIAが開発した
C/C++ ベースのGP GPU計算用
プログラミング言語

■ Host と Device

Host: CPU
Device: GPU

■ CUDAとGPU構造



- HostとDevice間の通信は、時間がかかる計算は、なるべく、Device内で主に行い、データの掃き出しに限り、Host (CPU) で処理することがコツ。

BlockとThreadは、3次元的に配置
上限: Block (x, y, z) =

(2147483647, 65535, 65535)
0(10⁹)

Thread: (x, y, z) = (1024, 1024, 64)

★ 2次元のみで十分

(注) このBlock数やThread数は、GPUの実際の構造を示しているわけではなく、CUDAが記載した数値

実際の構造

NVIDIA Tesla A100:

スレッド数 6812

(プロセッサ数)

2次元のみで十分

0(10⁹)程度 1024

↓ 単位 ↓

< block数 × thread数

[6812 スレッドごとに同時計算]

- GPU並列による高速化

```
for (int i=0; i<N; i++)
  for (int j=0; j<N; j++) {
    ....
  }
```

```
int i_global = blockIdx.x * blockDim.x
               + threadIdx.x;
if (i_global < N) {
  for (int j=0; j<N; j++) {
    ...
  }
```

このfor文を
おこなって
i_global (Max)
~ O(N)

- Device code (GPU 計算)

▶ Kernel

デバイス上でプログラムを走らせる関数内、
ホスト側で呼び出される関数

```
-- global -- void addVec ( ) {
    ....
}

int main ( ) {
    addVec <<< NB, NT >>> (変数1, 変数2)
    ....
}
```

↑ block ↑ thread

▶ Kernel 上で動く関数

- Kernel内でHost内の関数を呼び出すことは出来ない。
- Device上で呼び出してdevice上で動く関数が必要

```
-- device -- float square_dev ( ... ) {
    ....
}

-- global -- void kernel ( ... ) {
    ....
    X = square_dev ( );
    ....
}
```

■ CUDA programming

▶ CUDAの実行

hello.cu

```
#include <iostream>
using namespace std;
int main() {
    cout << "hello CUDA" << endl;
    return 0;
}
```

```
$ nvcc hello.cu
$ ./a.out
```

hello CUDA

▶ kernelの導入 何もいじり kernel

hello_kernel.cu

```
-- global -- void myKernel() {
    // nothing here.
}

int main() {
    cout << "Hello CUDA" << endl;
    myKernel <<< 1, 1 >>> ( );
    return 0;
}
```

山形関数構文:

関数名 <<< block数, thread数 >>> ();

Block p22

CudaMalloc
第1変数
(void **)

新たに割り当てられたアドレスを保持するポインタのポインタ。

pointer in device
↓
void **
↑
pointer in host

● GPUを用いた加算

add.cu

```
#include <iostream>
#include <cuda.h>
using namespace std;

-- global -- void add(int *c, int *a,
                      int *b) {
    *c = *a + *b;
}

int main() {
    int a, b, c; // values for the host
    int *a_dev, *b_dev, *c_dev;

    a_dev // values for the device: // *a_dev = ポインタの値
           is pointer

    // allocate memories on the device.
    cudaMalloc((void **)&a_dev, sizeof(int));
    cudaMalloc((void **)&b_dev, sizeof(int));
    cudaMalloc((void **)&c_dev, sizeof(int));

    a = 2;
    b = 7; // ポインタ ポインタ
           ↓ ↓
    cudaMemcpy(a_dev, &a, sizeof(int), cudaMemcpyHostToDevice);

    cudaMemcpy(b_dev, &b, sizeof(int), cudaMemcpyHostToDevice);

    add <<< 1, 1 >>> (c_dev, a_dev, b_dev);
                    pointer pointer pointer

    cudaMemcpy(&c, c_dev, sizeof(int), cudaMemcpyDeviceToHost);

    cout << "c: " << c << endl;

    cudaFree(a_dev);
    cudaFree(b_dev);
    cudaFree(c_dev);
    return 0;
}
```

```
$ nvcc add.cu
$ ./a.out
c: 9
```

■ ベクトル計算

- 逐次処理
(sequential processing)

timer.cuh

↳ measureTime()
は2回呼ぶ間の時間をms
で返す。

- CPUを用いたベクトル計算

add_vec_serial_cpu.cu

```
#include <iostream>
#include <cuda.h>
#include "timer.cuh"
// current directory
// に置く。
```

```
using namespace std;
const int Nv = 1e+6;
void setup_vec(int *vec, int a){
    for(int i=0; i<Nv; i++)
        vec[i] = i*a;
}
```

```
void add_vec(int *c, int *a, int *b){
    for(i=0; i<Nv; i++)
        c[i] = a[i] + b[i];
}
```

```
int main() {
    int *a, *b, *c;
    a = (int *) malloc (Nv * sizeof(int));
    b = (int *) malloc (Nv * sizeof(int));
    c = (int *) malloc (Nv * sizeof(int));
```

```
    setup_vec(a, 1);
    setup_vec(b, 2);
```

```
    double ms;
    measureTime();
    for(i=0; i<100000; i++)
        add_vec(c, a, b);
```

```
    ms = measureTime();
    cout << "Time" << ms/1000
         << "ms" << endl;
```

```
// Free
    free(a);
    free(b);
    free(c);
    return 0;
}
```

和の計算を100000回繰り返した際の
計算時間を計測。

(3分くらいかかると)

add_vec_serial_gpu.cu

```
--global-- void add_vec (int *c, int *a, int *b){  
    for (int i=0; i<Nv; i++)  
        c[i] = a[i] + b[i];  
}
```

```
int main ( ) {  
    int *a, *b, *c;  
    int *a-dev, *b-dev, *c-dev;  
  
    a = (int *) malloc (Nv * sizeof (int));  
    b = (int *) malloc (Nv * sizeof (int));  
    c = (int *) malloc (Nv * sizeof (int));  
    cudaMalloc ((void **)&a-dev, Nv * sizeof (int));  
    cudaMalloc ((void **)&b-dev, Nv * sizeof (int));  
    cudaMalloc ((void **)&c-dev, Nv * sizeof (int));  
  
    setup_vec (a, 1);  
    setup_vec (b, 1);  
  
    cudaMemcpy (a-dev, a, Nv * sizeof (int), cudaMemcpy Host To Device );  
    cudaMemcpy (b-dev, b, Nv * sizeof (int), cudaMemcpy Host To Device );  
  
    double ms;  
    measureTime ( );  
    for (int i=0; i<1000001000; i++)  
        add_vec << 1, 1 >> (c-dev, a-dev, b-dev);  
    cudaMemcpy (c, c-dev, // gpu Nv * sizeof (int), cudaMemcpy Device To Host );  
    ms = measureTime ( );  
    cout << "Time" << ms/1000 << "ms" << endl;  
  
    free(a);  
    free(b);  
    free(c);  
    cudaFree (a-dev);  
    cudaFree (b-dev);  
    cudaFree (c-dev);  
  
    return 0;  
}
```

- 並列計算
 - ▶ thread の活用

add-vec-parallel.cu

```
#include <iostream>
#include <cuda.h>
#include "timer.cuh"
using namespace std;
const int Nv = 1000000;
const int NT = 512;
const int NB = (Nv + NT - 1) / NT;

.....

__global__ void add_vec(int *c, int *a, int *b) {
    int i_global = blockIdx.x * blockDim.x + threadIdx.x;

    if (i_global < Nv)
        c[i_global] = a[i_global] + b[i_global];
}

int main() {
    .....
    double ms;
    measureTime();
    for (int i = 0; i < 1000; i++)
        add_vec <<< NB, NT >>> (c_dev, a_dev, b_dev);

    cudaMemcpy(c, c_dev, Nv * sizeof(int), cudaMemcpyDeviceToHost);
    ms = measureTime();
    cout << "Time:" << ms / 1000 << "ms" << endl;
    ....
    return 0;
}
```

block ID block 数 256 個
 thread ID 1024

Lecture2 総和計算

並列計算において、総和計算は注意が必要

悪い例

```
const int Narr = 10000000
const int NT = 512
const int NB = (Narr + NT - 1) / NT
-- global -- gpu_sum(int *sum_dev, int *arr_dev){
    int i_global = blockIdx.x * blockDim.x + threadIdx.x;
    if(i_global < Narr)
        *sum_dev += arr_dev[i_global];
}

....
int sum(int *arr){
    int sum;
    int *sum_dev, *arr_dev;
    cudaMalloc((void**)&sum_dev, sizeof(int));
    cudaMalloc((void**)&arr_dev, Narr * sizeof(int));
    cudaMemcpy(arr_dev, arr, Narr * sizeof(int), cudaMemcpyHostToDevice);
    gpu_sum << NB, NT >>> (sum_dev, arr_dev);
    cudaMemcpy(sum, sum_dev, sizeof(int), cudaMemcpyDeviceToHost);
    cudaFree(sum_dev);
    cudaFree(arr_dev);
    return sum;
}

// main
```

問題点:

- ▶ *sum_dev に 複数の スレッド が 同時に access する。
- ▶ 計算結果は 正しい 結果 になる
- ▶ コンパイラは、この問題を 教えてくれない。

• Atomic 演算

sum_atomic.cu

```
....  
__global__ gpu_sum(int *sum_dev, int *arr_dev){  
    int i_global = blockIdx.x * blockDim.x + threadIdx.x;  
    if(i_global < Narr)  
        atomicAdd(sum_dev, arr_dev[i_global]);  
}  
....
```

atomicAdd を用いることで、access の衝突は起さず済む。
ただし、sum_dev にある thread A が書き込んでいる間、他の thread は待機
している為、逐次処理と同等。