

CUDA を使うと並列演算が高速になる仕組み

鈴木 量三郎

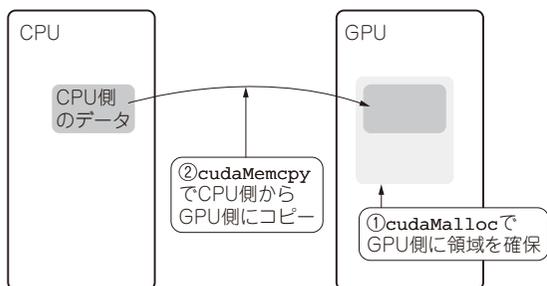


図1 GPUを使う大前提…GPUに計算を依頼するにはデータを転送する必要があります

● GPUに期待すること…並列演算による高速処理

ここでは、CUDAを使うと処理が高速になる仕組みをプログラマ視点で解説します。

特にエヌビディアのGPUを使うときは、PCI Express接続のグラフィックス・ボードを追加する場合があります。GPUには積和演算を並列にできる計算エンジンが多数搭載されており、CPUの逐次処理では時間のかかる計算を効率的に処理できます。GPUに計算を依頼することで、トータルとして高速に処理できることが期待できます。

● GPUを使う大前提…近くのDRAMへ転送が必要

CPUがGPUに計算を依頼するという事は、計算の元となるデータをGPUでも扱えるようにしなくてはなりません。そのためCPUは自分の手元にあるデータをいったんGPUの近くにあるDRAMに転送します。

▶データの転送

例えば、CUDAでよく使うAPIの組み合わせは、`cudaMalloc`と`cudaMemcpy`です。図1のように、`cudaMalloc`でGPU側のDRAMに領域を確保し、`cudaMemcpy`でCPU側のデータをコピーします。

▶プログラムの転送

GPUが実際に走らせるプログラム自身も、CPUからGPUへデータとして転送されます。

リスト1 CPUからGPUで実行する関数を呼び出すコード

```
vectorAdd<<<blocksPerGrid, threadsPerBlock>>>(
    d_A, d_B, d_C, numElements);
```

例えばリスト1のように、まるでCPUからGPUを呼び出しているように見えますが、実際には次の操作が裏で行われています。

- (1) CPU側でCUDAのカーネル・プログラムをコンパイル
- (2) GPU側のDRAMにプログラムを置く場所を確保
- (3) コンパイルされたデータをGPUのDRAMへ転送
- (4) GPUにプログラムの場所を教えて実行を指示

高速化の仕組み①…パイプラインとSIMD

● GPUでは意外と大変な足算

単純な $1234 + 5678$ の足し算を考えます。

```
1234
+ 5678
```

この処理は、本質的に並列化が困難です。手計算する場合も、一の位の足し算 $4 + 8$ から始めるでしょう。千の位の $1 + 5$ から始めようとする、百の位の足し算の結果によって繰り上がりが発生する可能性があるため、千の位だけ見て単純に結果を 6 と決定できないためです。少し考えると、このケースでは百の位は $2 + 6$ なので、繰り上がりが発生しないことが分かります。

幾つかの特殊なケースで並列化可能な箇所が見つかりますが、本質的に足し算は下の位から計算するという順序性があります。32ビットや52ビット (doubleの仮数部) の足し算をしようと思った場合、順序性があるため高速化するためには意外と工夫が必要です。

● 工夫1：パイプライン

足し算自身の並列化による高速化は困難ですが、ベルトコンベア方式で計算することで、多数の足し算の計算効率を上げることを考えます。