

畳み込み実装5

shun sato

GPUを使ってみよう

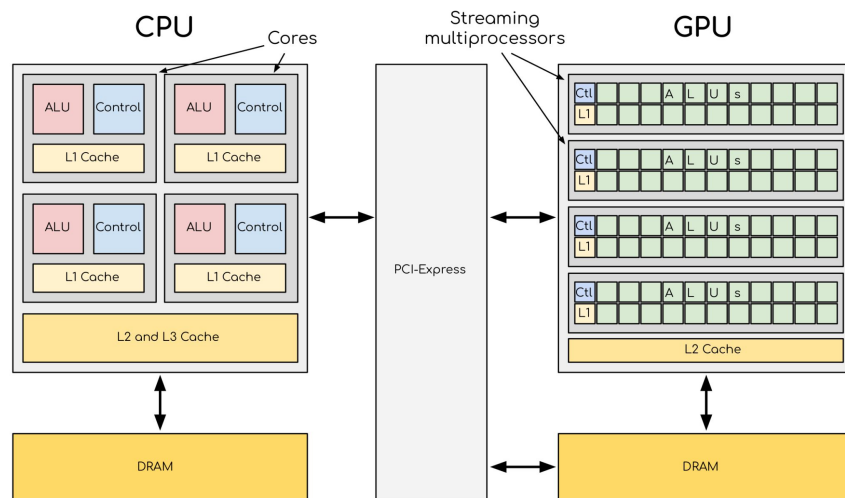
- 今や機械学習はGPUを使わないと**実行速度**が出せない
- GPUは何が得意ですか？



NVIDIAのGPU

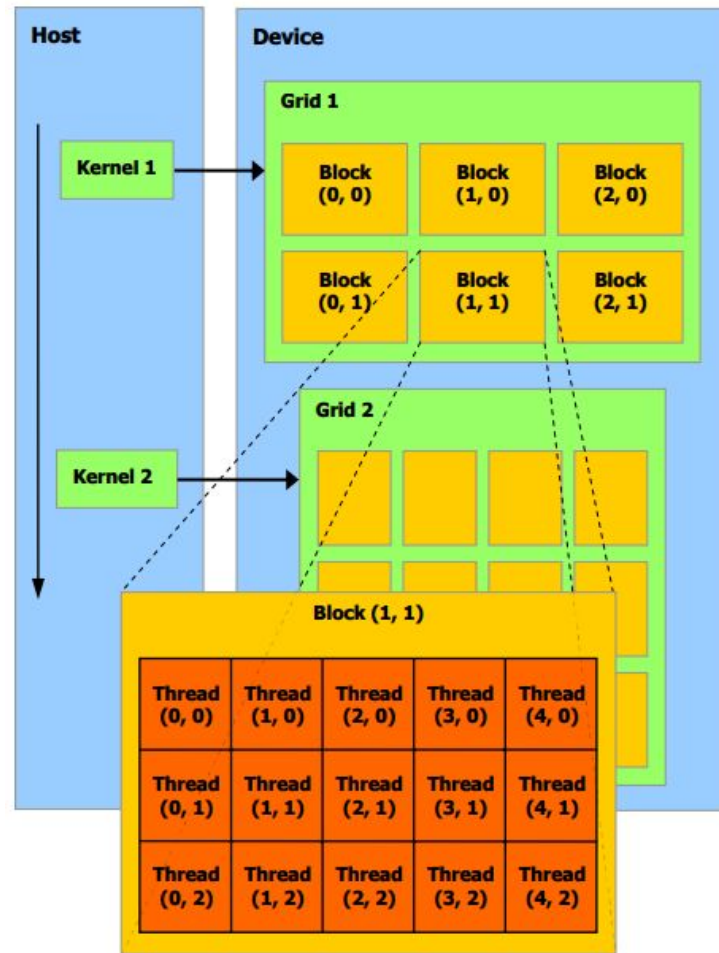
GPUとは

- GPUは**並列計算**が得意なハードウェア
- 大量の計算ユニットでまとめて計算
 - 深層学習モデル(行列計算)と相性がいい！
- RTX 4090
 - **16,384**個の並列計算ユニット
 - 24GBのメモリ
 - 2.5GHzで動作



CUDAアーキテクチャ

- **CUDA**はNVIDIAのGPUが採用しているアーキテクチャ
 - 計算ユニットを**階層構造**でまとめている
1. Grid
 - カーネルごとに1つ発行
 2. Block
 - Threadをまとめる単位
 3. Thread
 - 命令を実行する単位



CUDAアーキテクチャ

色々な種類のメモリ

Registers

- **最速**: Threadごとに割り当て

Local Memory

- 遅い: Registerが足りないときに割り当て

Shared Memory

- 速い: Blockごとに割り当て(共有)

Constant Memory

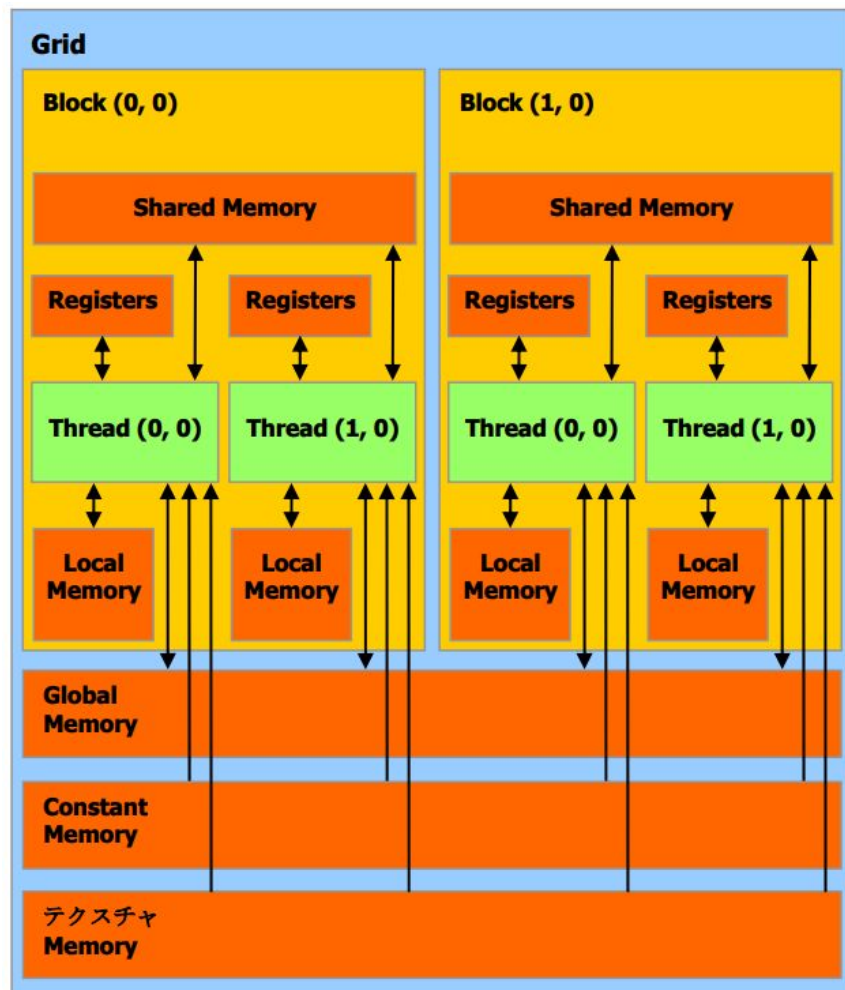
- 普通: 読み取り専用で全体共通

Texture Memory

- 普通: 2Dの参照に特化

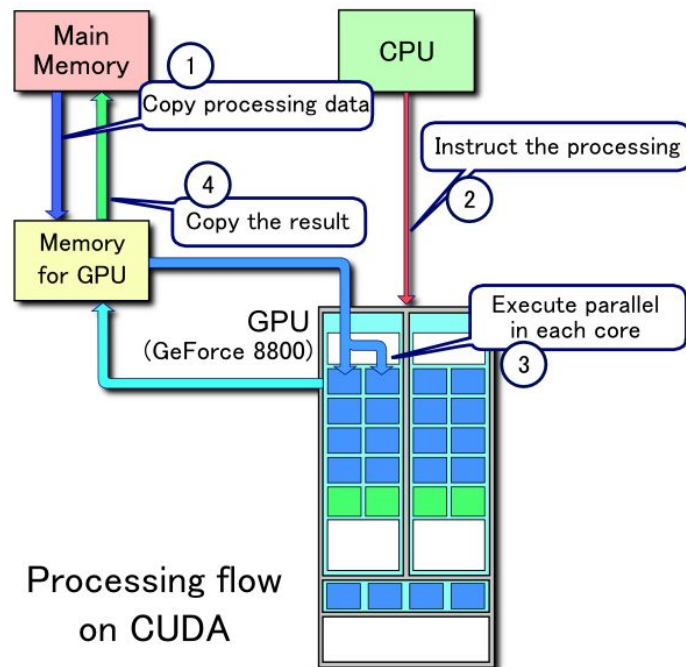
Global Memory

- **非常に遅い**: 全体でなんでも共通



GPUで命令を実行するまで

1. **ホスト**(CPU側)から**デバイス**(GPU側)に必要なデータを転送
2. ホストが**カーネル**(GPUの命令)を発行 ➡ 転送
3. デバイスがカーネルを実行
4. デバイスからホストに結果を転送(完了)



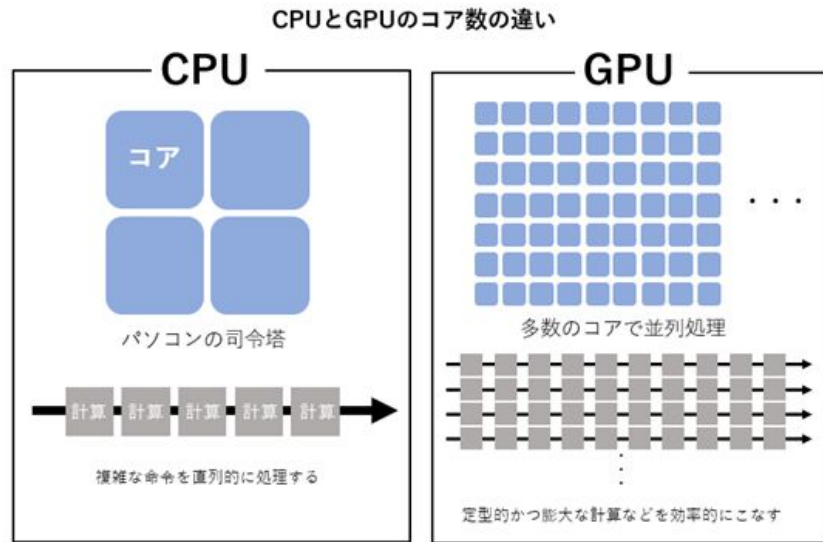
CPU vs GPU

CPU

- 計算は**超高速**
- 複雑な命令も高速に実行
- コアが**4~24**個

GPU

- 各コアの計算は遅い
- シンプルな命令の実行に特化
- コアが**1000**個以上
- GPUが得意な計算
 - 計算する**回数**が決まっている
 - **同時に**計算ができる➔ **画像処理**や**行列演算**が得意



<https://allai.jp/what-is-gpu/>

畳み込みをGPUで処理する

畳み込み演算

- 各ピクセルは同時に計算してもいい
 - GPUと相性抜群
- CPUはループで頑張ってた
 - $1024 \times 1024 = 1,048,576$
- GPUのボトルネック
 - 画像の転送の方が時間かかる
 - 処理は一瞬

