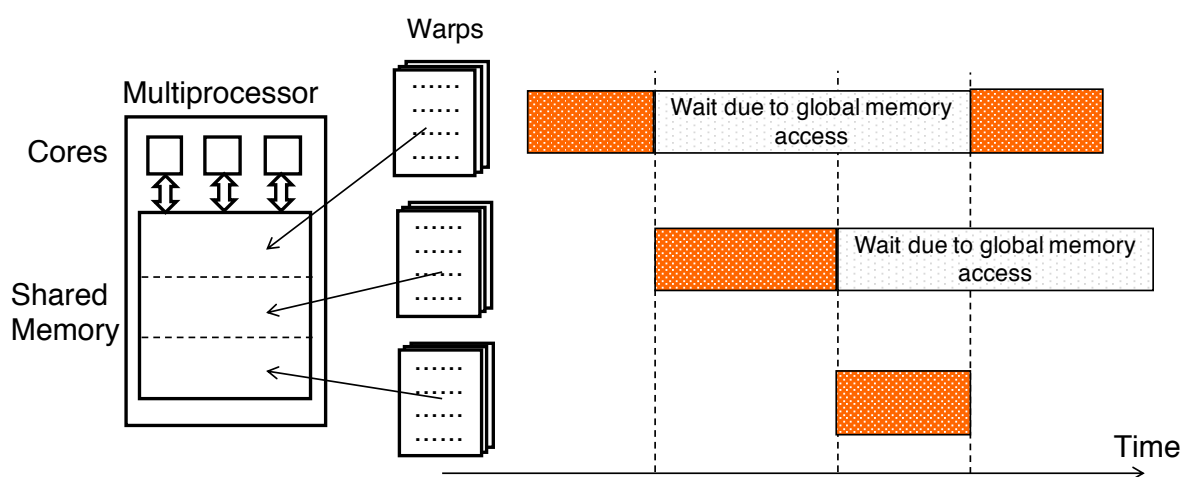

GPUプログラミング入門 講習会

GPUアーキテクチャとアルゴリズム設計

一関工業高等専門学校 未来創造工学科 情報・ソフトウェア系
小池 敦



目次

第 1 章 初日	3
1.1 はじめに	3
1.2 GPU と GPU コンピューティング	5
1.3 GPU プログラミングの情報源と GPU 環境の取得	7
1.4 GPU アーキテクチャ	10
1.4.1 概要	10
1.4.2 メモリ	16
1.4.3 マルチプロセッサ	27
1.5 GPU プログラミング	30
1.6 実習	33
第 2 章 二日目	37
2.1 はじめに	37
2.2 使用する GPU の仕様	39
2.3 環境設定	40
2.4 実習 1：最初のプログラム	40
2.4.1 並列化しないプログラム	41
2.4.2 CUDA による GPU プログラム（その 1：1 ブロックのみ使用）	43
2.4.3 CUDA による GPU プログラム（その 2：2 ブロック使用）	46
2.5 高速アルゴリズムの設計	51
2.6 実習 2：効率良く総和計算を行うプログラム	64
2.6.1 Cascading アルゴリズムによる総和計算	64
2.6.2 Thrust ライブラリとの比較	72
2.7 高速化トピックス	74
2.7.1 データに依存性がある場合	74
2.7.2 構造体の配列とコアレスアクセス	77
2.7.3 行列の共有メモリへの格納	78
2.7.4 ライブラリの活用	79
2.7.5 複数 GPU の活用	79
2.7.6 GPU ハードウェアの活用	80
付録：Linux 入門	82

第1章 初日

1.1 はじめに

数理工学講習会

GPUプログラミング入門

初日：GPUアーキテクチャとGPUプログラミング
一関高専 未来創造工学科 情報・ソフトウェア系
小池 敦



1

GPUプログラミング入門

- 目標
 - GPUアーキテクチャの特性を理解した上でGPUを用いた並列化ができる
- 対象者
 - C言語等でプログラミングをしたことがあり、GPUを用いた高速化に興味がある
 - GPUによる並列計算の概要を知りたい
- 必要な予備知識
 - C言語を触ったことがある

2

本講習会は、GPUを使って高速な並列プログラムを書きたい人のための二日間の入門講座です。目標は、単にプログラミングの文法を知るだけでなく、GPUのアーキテクチャを理解した上で、高速に動作するプログラムが書けるようになることです。対象者は、C言語等でプログラミングをしたことがあり、GPUを用いた高速化に興味がある人です。C言語については、簡単なコードが読めるくらいの知識があれば大丈夫です。また、単にGPUによる並列計算の概要を知りたいという人も対象としています。

GPUプログラミング入門

- 講習会の内容
 - 初日：GPUアーキテクチャとGPUプログラミング（座学中心）
 - 二日目：GPUプログラミング（実習中心）
- 必要な事前準備
 - 大阪大学サイバーメディアセンターの大型計算機試用アカウントを取得し、OCTOPUSにログインできるようにしておく
<http://www.hpc.cmc.osaka-u.ac.jp/service/intro/shiyo/>

3

講習会の内容はスライドの通りです。初日は、GPUのアーキテクチャを中心に話します。その中で、GPUのアーキテクチャを活かすためには、どのようにプログラミングすべきかという話もします。初日は座学中心で、実習については最後にサンプルプログラムを動かしてみる程度です。二日目は実際にGPUのコードを書きながら、プログラミングの文法と高速化のノウハウを学習します。

本テキストは大阪大学サイバーメディアセンターでの講習会で使うことを想定して作られています。もし他の環境で使用している場合は、適宜読み替えてください。サイバーメディアセンターを利用する場合、GPUが使用できるシステムはOCTOPUSです。事前にアカウントを取得し、OCTOPUSにログインできるようにしておいてください。

初日の内容

- GPUアーキテクチャ
 - GPUアーキテクチャの概要とその長所
 - GPUアーキテクチャは世代により大きく変わるが共通のもの+最新の動向を紹介する
- GPUプログラミング
 - GPUで汎用計算を行う方法を説明する（グラフィックの話はしない）
 - 1つのプログラムを複数GPUアーキテクチャで動作させる仕組み（CUDA）
- 計算機環境お試し

4

初日の内容の詳細はスライドのようになります。まず、GPUアーキテクチャの概要と長所について話すのですが、GPUアーキテクチャは世代によって、大きく変化しています。本講習会では主に世代で共通の部分を扱うこととし、最新の話については簡単に触れる程度とします。


次に、GPUを使ったプログラミングの仕方について概要を話します。後でも話しますが、今回の講習会では、GPUを使った並列計算の話をしていきます。GPUは本来グラフィック処理をするためのプロセッサですが、今回はグラフィックの話はしません。様々なGPUアーキテクチャ上で動作するプログラムを作る仕組みとしてCUDAがあります。CUDAにはプログラミング言語や開発環境が含まれており、それを説明します。

最後に、計算機環境のお試しとして、サンプルプログラムを実行してみます。もしうまくいかない人がいれば、二日目までに環境を確認しましょう。

1.2 GPUとGPUコンピューティング

GPU

- Graphics Processing Unit
- 元々はグラフィック処理の専用プロセッサ



Introducing The GeForce GTX 1080 Ti, The World's Fastest Gaming GPU
<https://www.geforce.com/whats-new/articles/nvidia-geforce-gtx-1080-ti>

5

それでは、本日の内容に入っていきます。まず、GPUについて説明します。GPUはGraphics Processing Unitの略称で、その名の通り、元々はグラフィック処理専用のプロセッサでした。

GPU

- 高い並列性能を持つため、グラフィック処理以外の汎用計算にも使われている
 - 最新モデル（GV100）では1デバイスで5120コア
 - 科学技術計算、深層学習などによく使われる
 - NVIDIA社はGPUコンピューティングと呼んでいる（元々はGPGPUと呼んでいた）

NVIDIA TESLA V100 GPUアーキテクチャ
<https://images.nvidia.com/content/pdf/tesla/Volta-Architecture-Whitepaper-v1.1-jp.pdf>

6

GPUにはとても多くのコアが搭載されています。コアというのは、プログラムを実行する装置のことで、ざっくり言えば、例えば4つのコアが搭載されていれば4つのプログラムを同時に並列に実行できるということです。コアでなくプロセッサと呼んでも良いのですが、複数のコアを一つのデバイスに搭載した装置もプロセッサと呼ぶので、それと区別するために通常はコアと呼ばれています。NVIDIAの最新GPUモデルであるGV100では、1デバイスに5120個のコアが搭載されています。CPUは4コアとか8コアとかのレベルなので、GPUには圧倒的な数のコアが搭載されているということがわかると思います。この圧倒的な数のコアを活用して、グラフィック処理以外の処理を行わせるということが盛んに行われています。シミュレーションなどの科学技術計算においては、既に多くの分野で活用されていますし、近年では深層学習での活用が進んでいます。GPUをグラフィック以外に活用することは、General Purpose GPU (GPGPU)と呼ばれていましたが、この分野を引っ張ってきたNVIDIA社は近年はGPGPUという用語を使用しなくなり、代わりにGPUコンピューティングと呼んでいるようです。この講習会ではこれのことを「GPUによる汎用計算」もしくは「GPUコンピューティング」と呼ぶことにします。また、本講習会でGPUプログラミングと言ったら、それはGPUコンピューティングのためのプログラミングを意味します。

GPU

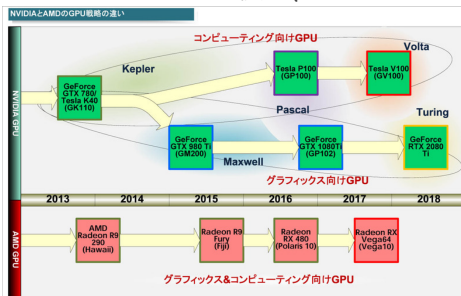
- 主な開発企業
 - NVIDIA
 - AMD
- 汎用計算においては、NVIDIAの1強
 - CUDAと呼ばれるNVIDIA社GPU向けのプログラミング環境（後述）が整備されていることが大きい
 - 高速化しやすい、そこそこプログラミングしやすい
 - 本講習会でもNVIDIA社GPUをターゲットとする

7

GPUは多くの半導体メーカーによって開発されていますが、代表的な企業としてNVIDIA社とAMD社を挙げることができます。しかし、GPUを用いた汎用計算に関しては、NVIDIA社の1強と言える状況です。これはNVIDIA社GPU向けにCUDAと呼ばれるプログラミング環境を整備していることが大きな要因ではないかと思えます。CUDAを使うと、それなりに簡単にプログラミングができ、頑張れば高速化のチューニングをすることも可能です。ということで、本講習会でもNVIDIA社のGPUを用いて、CUDA上でプログラミングを行います。

NVIDIA GPU

- 今回はTesla P100を使用(Pascalアーキテクチャ)



8

この図は後藤弘茂氏によるGPUモデルの変遷のまとめです。上の緑の行がNVIDIA社GPUで下の赤の行がAMD社GPUです。後藤氏のまとめによればNVIDIA社GPUは2014年頃グラフィック処置をメインとするGPUとGPUコンピューティングをメインにするGPUの2系統に別れたようです。私の感覚的にも、その頃深層学習でのGPU利用が流行ってきて2系統に別れたという印象です。そして、今回の講習会で使用するGPUは最上段の右から二番目にあるTesla P100です。これは最新モデルではありませんが、Pascalと呼ばれる高性能アーキテクチャを採用したモデルです。

1.3 GPUプログラミングの情報源とGPU環境の取得

GPUプログラミングの情報源

- NVIDIA社ドキュメント
 - CUDA C Programming Guide
<https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>
 - NVIDIAのWhite Paper (ハードウェア説明書)
<https://images.nvidia.com/content/pdf/tesla/Volta-Architecture-Whitepaper-v1.1-jp.pdf>
 - NVIDIA社 (日本) の各種資料 (検索する)
<https://www.nvidia.com/content/apac/gtc/ja/pdf/2017/1041.pdf>
<http://on-demand.gputechconf.com/gtc/2013/jp/sessions/8002.pdf>

9

GPUプログラミングの情報源

- NVIDIA Developer
 - CUDA Education & Training
<https://developer.nvidia.com/cuda-education-training>
 - GPU Accelerated Computing with C and C++
<https://developer.nvidia.com/how-to-cuda-c-cpp>
 - Training Material and Code Samples
<https://developer.nvidia.com/cuda-education>
- NVIDIA Developer Blog (一番詳細)
<https://devblogs.nvidia.com>
 - Mark Harris氏の記事 (下記は入門編)
<https://devblogs.nvidia.com/even-easier-introduction-cuda/>

10

GPUプログラミングの詳細説明に入る前に、これらについて知るための情報源について紹介したいと思います。この講習会で話す内容も基本的にはこれらの情報源を基にしています。この講習会は入門編ですが、さらに詳細な知識を得たいと思った際には、これらの情報源を参照すると良いと思います。

まずはNVIDIA社のドキュメントについて紹介します。最も基本となるドキュメントはCUDA C Programming Guideです。CUDAを使ってGPUプログラミングする際に必要となる情報が網羅的に書かれています。プログラミングをしていて詳細な情報が欲しくなった時には、基本的にはこの資料を見るのが良いと思います。次が、NVIDIA社のホワイトペーパーです。ホワイトペーパーという用語は、ソフトウェアの人には馴染みがないかもしれませんが、いわゆるハードウェアの詳細な説明書です。NVIDIAでは、GPUアーキテクチャごとにホワイトペーパーを作っており、そのGPUアーキテクチャの特徴や詳細な仕様を説明しています。なお、ホワイトペーパーについては、日本語版も作られているようです。また、日本のNVIDIA社は講習会などのために日本語の資料を作っており、検索すると見つかる場合もあります。

NVIDIA Developerのサイトには、GPUプログラミングを勉強するための各種資料が掲載されています。動画やサンプルコードもあり、これからGPUプログラミングを始めたいと考えている人には役立つと思います。

別の情報源として、NVIDIAエンジニアによるDeveloper Blogがあります。CUDA C Programming Guideなどでは説明されていない詳細仕様が書かれていることもあり、GPUアーキテクチャやCUDAについて詳細を知りたい時には重要な情報源です。特にMark Harris氏は多くの有益な記事を書いています。技術的な詳細の記事だけでなく、スライドに記載したような入門編の記事も書いています。

GPUプログラミングの情報源

- 後藤弘茂氏のニュース記事（幅広いネタ）
 - NVIDIAのマルチGPU戦略とインターコネクト帯域
<https://pc.watch.impress.co.jp/docs/column/kaigai/1122287.html>
 - AMD、7nmで最大64コアの「ZEN2」とNVIDIA Voltaを上回る「Radeon Instinct M60」
<https://pc.watch.impress.co.jp/docs/column/kaigai/1151995.html>
 - NVIDIAの巨大GPUを支えるTSMCのインタポーザ技術
<https://pc.watch.impress.co.jp/docs/column/kaigai/1064109.html>

11

GPUプログラミングの情報源

- CUDA関連の教科書
 - CUDA C プロフェッショナル プログラミング
https://www.amazon.co.jp/dp/B015R0M8TS/ref=dp-kindle-redirect?_encoding=UTF8&btkr=1
 - 詳しい。なんでも書いてある。初心者には詳しすぎるかも
 - GPUプログラミング入門—CUDA 5による実装
https://www.amazon.co.jp/dp/B00KCL6PFG/ref=dp-kindle-redirect?_encoding=UTF8&btkr=1
 - CUDAの説明はあっさり
 - 数値計算の例が多く載っているのでその分野の人はいいかも
 - はじめてのCUDAプログラミング
https://www.amazon.co.jp/dp/B000000000/ref=dp-kindle-redirect?_encoding=UTF8&btkr=1
 - よくまとまっているが、内容が古い

12

その他の情報源についてもいくつか挙げておきます。

まず、ニュース記事としては、後藤弘茂氏の記事が挙げられます。GPUプログラミング関連のニュース記事はNVIDIA社の発表をそのまま記事にしたものが多いですが、後藤弘茂氏の記事は、幅広い知識に基づいて深い洞察が行われているものが多いです。スライドのリンクのように、企業戦略の話だったり、半導体製造技術の記事もあります。

次にCUDAの教科書について、私が所有しているものについてコメントしておきます。「CUDA C プロフェッショナル プログラミング」はCUDAについての詳細が網羅的にまとめられた本です。CUDA C Programming Guideが難しすぎるとか、日本語で読みたいとか思っている人には良いと思います。ただし、分厚い本なので、初心者が一から全部読むのは辛いかもしれません。「GPUプログラミング入門—CUDA 5による実装」は、CUDAについての説明はあっさりしていて、これだけで理解するのは厳しいかもしれませんが、数値計算の例が多く載っているので、その分野の人はいいかもしれません。「はじめてのCUDAプログラミング」はCUDAについて簡潔にまとめられており、初学者に向いています。しかし、2009年に出版された本であり内容がかなり古いので、注意が必要です。最新のGPUやCUDAでは異なっているという内容が多くあります。

GPUプログラミングの情報源

- アーキテクチャの教科書
 - 「ヘネ&パタ」：ヘネシー&パターソン コンピュータアーキテクチャ 定量的アプローチ第5版 4.4節
<https://www.shoehisha.co.jp/book/detail/9784798126234>
 - 「パタ&ヘネ」：コンピュータの構成と設計 第5版
ダウンロード教材 Appendix C
<http://ec.nikkeibp.co.jp/nsp/dl/09842/index.shtml>
 - 無料で読めるが内容が古い

13

GPU アーキテクチャについて記載されている教科書についても紹介しておきたいと思います。一冊目が通称「ヘネ&パタ」と呼ばれているもので、正式名称はスライドの通りですが、この4.4節にGPU アーキテクチャの説明があります。GPU アーキテクチャの利点についてかなり詳しく説明されていると思います。二冊目は通称「パタ&ヘネ」で、正式名称は「コンピュータの構成と設計」ですが、この付録にGPU アーキテクチャの説明があります。これは素晴らしいことに、このリンクから無料でダウンロードできるのですが、内容がかなり古いので注意が必要です。

GPUの準備

- 汎用計算をするには、データセンター用GPUを使うのが一般的（画像用も使える）
 - 自分で買うと15万～100万以上
- クラウド環境
 - 阪大サイバーメディアセンターOCTOPUS（本ワークショップ）
 - Amazon Web Services (AWS) EC2
 - P3インスタンス（旧世代でよければ、P2, G3）
 - 1時間 3USD～（P2なら 0.9USD～）

Amazon EC2 P3 インスタンス
<https://aws.amazon.com/jp/ec2/instance-types/p3/>

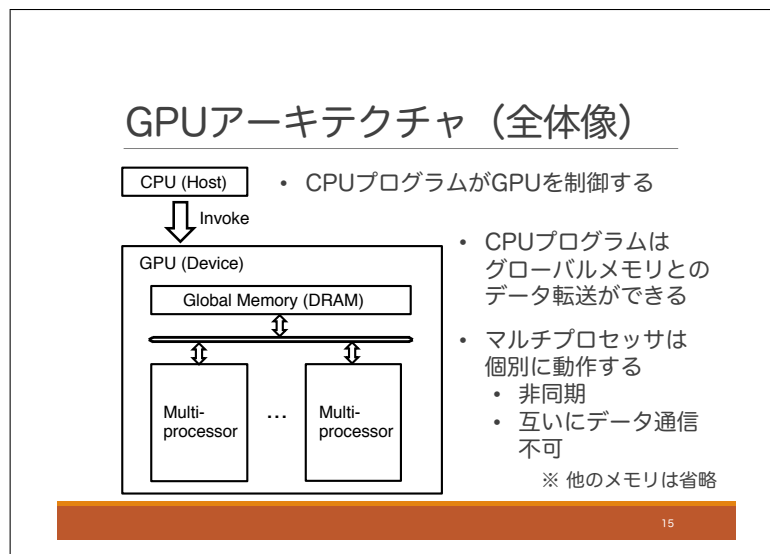
14

最後にGPU環境を準備する方法について紹介します。NVIDIAのGPUであれば、CUDAを用いたGPUプログラミングを行うことができます。スライド8で紹介した通り、NVIDIAのGPUはGPUコンピューティング向けとグラフィックス向けに分かれています。どちらもCUDAが使えるのですが、汎用計算における性能等を考えるとGPUコンピューティング向けを使う方がいいです。データセンター用GPUと説明されているものがGPUコンピューティング向けのGPUです。自分でGPU環境を構築しようとする、安くて15万円くらい、高いと100万円以上します。

お手軽にGPUプログラミングを始めるにはクラウド環境がおすすめです。本講習会で使用する大阪大学サイバーメディアセンターのOCTOPUSには試用期間があり、お試しで使ってみることもできます。また、Amazon Web Services、略してAWSのEC2というサービスでは、クラウド環境を時間単位の料金で使用することができます。P3インスタンスというものを選択すると最新GPUアーキテクチャのVoltaが1時間3USDドルで使用できますし、最新GPUでなくても良いなら、P2インスタンスを選択することで、1時間0.9USDドルで使用できます。

1.4 GPUアーキテクチャ

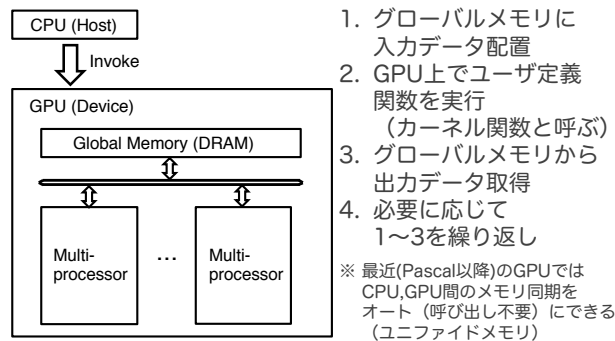
1.4.1 概要



それでは GPU アーキテクチャについて説明していきます。全体像はこのスライドのようになります。まず重要なこととして、GPU は単独では動作できず、常に CPU によって制御されます。CPU 上で動いているプログラムが GPU に対して処理を依頼すると GPU がそれを実行するという流れです。GPU プログラミングでは、CPU のことをホスト、GPU のことをデバイスと呼びます。例えば、ホストメモリと言った場合には CPU 上のメモリのことですし、デバイスメモリと言った場合は、GPU 上のメモリのことです。GPU は様々なモジュールで構成されますが、最も基本的なモジュールはグローバルメモリとマルチプロセッサです。本当はもっとたくさんの種類のメモリがあるのですが、ここでは省略しています。マルチプロセッサは、実際に処理を行う部分で、通常 1GPU デバイスの中に複数のマルチプロセッサが入っています。グローバルメモリはすべてのマルチプロセッサからアクセス可能なメモリです。CPU もグローバルメモリにアクセスできます。CPU プログラムは GPU に処理を依頼する前に、必要なデータをグローバルメモリに転送しておくことができます。少し細かいことを言うと、昔の GPU では CPU メモリ上にある必要なデータを必ず事前にグローバルメモリに転送しておかなくてはならなかったのですが、2016 年に発表された Pascal アーキテクチャ以降の GPU では CPU メモリ上のデータと GPU グローバルメモリ上のデータの自動同期の機能が搭載されたため、必ずしも事前に手動でデータ転送を行う必要はなくなりました。

次にマルチプロセッサについてですが、大きな特徴は、マルチプロセッサのそれぞれが個別に動作するということです。もう少し具体的に言うと、各マルチプロセッサは、プログラムを独自のペースで実行しており、他のマルチプロセッサとタイミングを合わせるといった機能は持っていません。他のマルチプロセッサとタイミングを合わせる命令は同期命令と呼ばれますが、GPU ではマルチプロセッサ間の同期命令は存在せず、各マルチプロセッサは非同期に実行されるということになります。また、マルチプロセッサ間でデータ通信を行うこともできません。これも GPU プログラミングを行う際の大きな制約になります。なお、CPU において、すべてのマルチプロセッサが処理を完了するのを待つということとはできません。あるマルチプロセッサでの処理結果を別のマルチプロセッサに渡したい時に、このことを活用することができます。

GPUプログラミングの処理概要



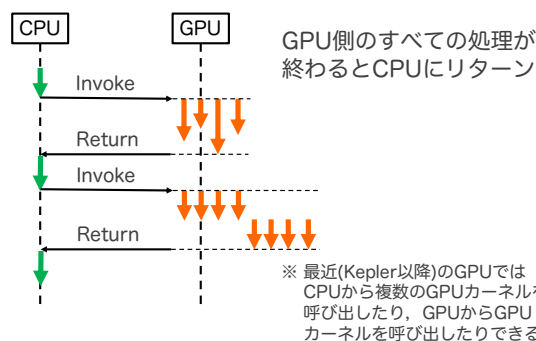
16

前のスライドと少し内容がかぶりますが、GPUプログラミングにおいて、プログラムがどのように処理を進めていくかということについて説明します。実際にはもっと様々な処理を行うことも可能なので、ここで説明するのは、基本的なパターンの一つです。

まず、手順1としてCPU上のプログラムはCPU上のメモリにあるデータをGPUのグローバルメモリに転送します。手順2では、CPU上のプログラムはユーザが作成した関数の実行をGPUに依頼します。すると、GPUはその関数を実行します。なお、GPU上で実行される関数はカーネル関数と呼ばれます。手順3では、CPU上のプログラムは、カーネル関数の実行終了後GPU上のグローバルメモリから必要なデータを取得します。その後、手順4として、必要に応じて手順1~3を繰り返します。

なお、前のスライドでも話した通り、手順1,3についてはCPUメモリとGPUメモリの自動同期機能を使えば不要です。これはユニファイドメモリと呼ばれており、Pascalアーキテクチャ以降のGPUでサポートされています。

基本的な処理の流れ

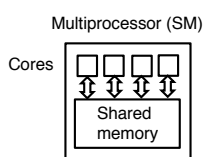


17

前のスライドの手順をシーケンスで表したのがこのスライドです。前のスライドの手順1~3を2回繰り返す場合の例になっています。繰り返しになりますが、これも基本的なパターンの一つで、常にこのようなシーケンスになるというわけではありません。CPUプログラムがカーネル関数を呼び出すと、GPU上の複数のスレッドが並列処理を行います。そしてすべての処理が終わると、処理がCPUに戻されます。

2013年に発表されたKepler以降のGPUでは、複数のカーネル関数を同時に実行したり、カーネル関数の中から別のカーネル関数を呼び出したりできるようになっています。また、カーネル関数実行中にCPUが他の処理をすることもできます。なので、実際にはもっと複雑なシーケンスになることが多いです。

マルチプロセッサ (Streaming Multiprocessor; SM)



- 複数のコアと共有メモリからなる
 - 32~192コア/ SM
- コアは（基本的には）同期しながら処理を行う
 - 複数のコアが同一命令を実行
 - 一部のコアのみが命令実行することも可能

※ 最新GPU(Volta) ではより柔軟なコアの制御ができる

18

次に、マルチプロセッサの中身について見ていきます。NVIDIA はこれを Streaming Multiprocessor と呼んでおり SM と略しているので、本講習会でもこの略称を使います。

マルチプロセッサは、主に複数のコアと共有メモリと呼ばれるメモリからなります。これまでと同様に、マルチプロセッサにはもっと様々なモジュールが搭載されており、それらが高速化に役立っているのですが、ここでは省略します。1 マルチプロセッサに搭載されるコアの数は、アーキテクチャによって異なるのですが、32 コアから192 コアの間くらいになっています。

マルチプロセッサ内のすべてのコアは、基本的には同期して処理を行います。つまり、マルチプロセッサ内のコアは常に同じ命令を実行します。命令は同じですが、対象とするデータはコアごとに変えることができます。また、指定した条件を満たすコアのみ命令を実行するということもできます。複数のコアのそれぞれが異なる命令を実行したい場合には、1つのコアのみが命令実行を行い、他のコアは休ませるという処理を繰り返すこととなります。つまり、命令を並列に実行することはできず、一つのコアずつ順に命令を実行していくこととなります。これでは効率の良い処理にならないので、プログラミングする際には極力このようなことがおきないようにアルゴリズムを設計する必要があります。

なお、2017年に発表された Volta アーキテクチャでは良い柔軟なコアの制御が可能になっています。それ以前のアーキテクチャでは、すべてのコアは与えられたプログラムを同期実行していたのですが、Volta では条件分岐が発生した際にコアごとに異なる実行パスを辿れるようになりました。とはいえ、異なる命令を並列に実行できるようになったわけではないので、条件分岐が発生すると効率が悪くなることには変わりはありません。

これから説明すること

- GPUアーキテクチャはなぜ優れているのか？
 - そのために並列/並行処理により処理が高速化される仕組みを説明する
- 並列/並行処理による高速化の仕組み（概要）
 - 処理を複数コアに割り当て、並列実行することで処理時間が短くなる
 - メモリアクセス時間が低減・隠蔽される

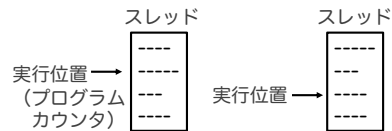
19

GPU アーキテクチャの概要を話しましたので、これからは GPU アーキテクチャはなぜ優れているのかについて説明していきたいと思います。それを説明するために、まず、そもそも並列処理や並行処理によりなぜ処理が高速化されるのかということについて説明したいと思います。高速化される要因は主に二つあります。一つ目は、当たり前ですが、処理を複数のコアに割り当てて並列実行することで、処理時間が短縮できるということです。一人でやると時間がかかる仕事も 10 人でやれば早く終わるという感じで、これはわかりやすいと思います。もう一つは、並列/並行処理によりメモリアクセス時間が低減・隠蔽されるということです。隠蔽という言葉は後で説明するので、ここでは説明しないでおきます。ともかく、メモリアクセスに必要な時間コストが小さくなるんだと思ってもらえば良いです。こちらはわかりづらいので、このあとは、メモリアクセスについて話していきたいと思います。

ちなみにこのスライドでは、並列/並行処理という言い方をしていますが、並列処理と並行処理には専門用語として使い分けがあります。並列処理と言った場合は、複数のコアを同時に動かすことで同時に複数の処理を行うことです。一方、並行処理は複数の処理を同時進行で実行することを意味し、この場合は、必ずしも複数のコアが同時に処理を実行する必要はありません。一つのコアが複数の処理を時分割で交互に処理する場合も並行処理となります。例えるなら、一人の人が複数の仕事をながら作業で同時こなしていくイメージです。10 人で 100 個の仕事をこなすような場合も並行処理と言えます。並列処置が高速化に繋がることはわかりやすいですが、並行処理自体が高速化に繋がることはわかりにくいかもしれません。しかし、並列処理を行わない並行処理、つまり一人で複数の処理をながら作業する場合でも一つ一つ個別に処理する場合に比べて処理が高速化されます。それは、メモリアクセス時間の隠蔽が大きな要因なのですが、それについて、これから説明していきたいと思います。

【用語】 スレッド

- プログラム実行の単位
 - 1CPU上で複数スレッドを実行可能
 - 複数スレッドは代わりばんこにCPUを利用する
⇒ 通常はOSがスレッドのスケジューリングを行う
 - GPUではハードウェアがスレッド管理を行う



20

詳細な説明に入る前にスレッドという用語を説明しておきたいと思います。スレッドというのはプログラム実行の単位の一つです。コンピュータを使用してプログラムを実行する際、基本的には、コンパイルされたコードを上から下に1行ずつ実行していきます。ただし、ジャンプ命令というのがあると指定行にジャンプします。コンピュータは今実行している行を常に保持しており、それをプログラムカウンタと言います。

ここで、複数のコードを並行処理する場合、プログラムカウンタなどの内部状態をコードごとに保持しておく必要があります。コンピュータでは、複数のコードの1つ1つをスレッドという単位で管理しており、1つのスレッドには一つのコードが割り当てられます。C言語の初学者が習うような基本的なプログラムは、一つのスレッドからなる場合が多いです。しかし、グラフィカルユーザインタフェースつまりGUIにより、マウスのクリックを検出するようなプログラムは基本的には複数のスレッドから構成されています。これは、マウスがクリックされた際、アプリを実行するためのコードとは別のコードが実行されるためです。

複数のスレッドの管理は通常はOSがCPUと連携して行うことが多いです。OS自体もいくつかのスレッドの集まりですが、OSが他のスレッドに関して、いつどのスレッドを実行するのかを管理しています。

後で詳しく話しますが、GPUはこれらとは異なったスレッド管理の仕組みを持っています。GPUではOSを用いずハードウェアがスレッドの管理を行っています。これにより高速なスレッドの切り替えが可能です。

高速化の課題

- コアを増やせばその分速くなる
というわけではない
- 速くならない主な要因
 - スレッド間の依存度が高い
(他のスレッドの処理完了待ち時間が長い)
⇒ 主にソフトウェアの問題 (二日目に扱う)
 - メモリアクセスに時間がかかる
(多くのコアが一斉にアクセスして渋滞する)
⇒ アーキテクチャ+ソフトによる適切な活用が重要

21

これまでの話の繰り返しになる部分もありますが、並列処理を用いて高速化を行う際の課題について見てみます。並列処理では、コア数を増やせばその分プログラムの実行が高速化されるというわけではありません。これはCPUのクロック周波数とは異なる状況です。クロック周波数については、周波数を大きくすれば、プログラムに何の変更を加えなくても実行速度が速くなりました。ところが並列処理については、そうではないということです。

コア数を大きくしても高速化されない原因の一つにスレッド間の依存度があります。例えば、100コアによる並列処理を行うためには、少なくとも100個のスレッドを作る必要がありますが、もし無理矢理100個のスレッドを作ったとしても、あるスレッドが別のスレッドの処理完了を待たないといけなような状況が多いと処理は高速化されません。これはソフトウェアを工夫することで改善される場合もありますが、これについては二日目に扱います。別の原因として、メモリアクセスに時間がかかるということがあります。各コアはメモリからデータを読み込んで処理をすることが多いので、コア数を増やしていくと、一般にメモリアクセスが渋滞します。これはGPUのように膨大な数のコアを搭載したデバイスでは特に深刻であり、GPUにはアーキテクチャとソフトウェアが協調して効率の良いメモリアクセスを実現するための仕組みがたくさん搭載されています。

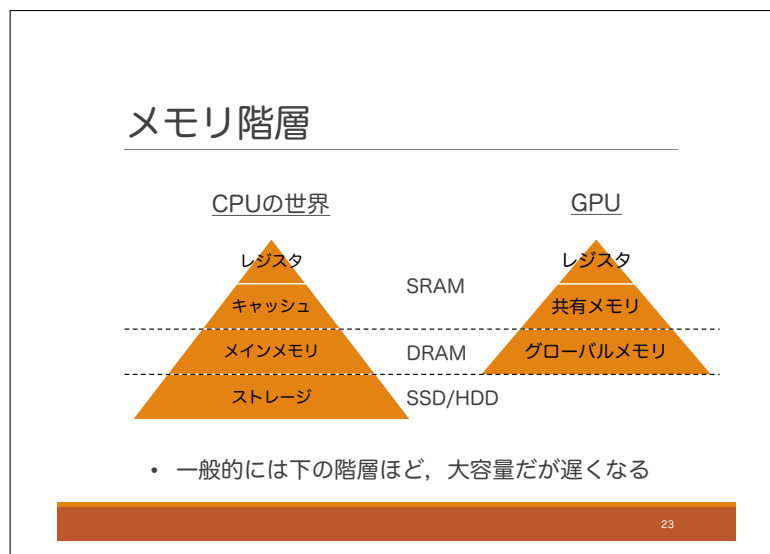
1.4.2 メモリ

メモリ

- メモリアクセスの効率化を議論する
 - メモリのハードウェアの特性を説明した後、各メモリへ効率的にアクセスする仕組みを説明する
- GPU
 - グローバルメモリ：DRAM
 - 共有メモリ：SRAM

22

これからしばらくの間 GPU が備えている効率の良いメモリアクセスの仕組みについて説明したいと思います。まずメモリのハードウェア的な特性について説明したのち、それを踏まえてメモリへの効率的なアクセス方法を説明します。GPU のグローバルメモリには DRAM というメモリが使われており、共有メモリには SRAM というメモリが使われていますので、DRAM と SRAM について説明します。



まず全体像の説明としてメモリ階層について見てみましょう。コンピューターアーキテクチャのメモリシステムでは、大容量化と高速化の両方の需要に対応するため、複数のメモリを組み合わせ使用します。まず、CPU においては、レジスタとキャッシュのために SRAM が使われています。これは高速アクセスができますが、低容量です。プログラム実行時により大きいデータを扱えるようにするため、別途メインメモリが搭載されますが、これには DRAM が使われます。DRAM は SRAM よりもアクセスに時間がかかります。さらに大容量なデータをコンピューターに保持しておくためにストレージが搭載されますが、これには SSD や HDD と呼ばれるデバイスが使われます。これは大容量ですが、DRAM よりもさらにアクセスに時間がかかります。

一方 GPU ですが、レジスタや共有メモリには高速アクセス可能な SRAM が使用されます。グローバルメモリにはメインメモリと同じく DRAM が使われています。

DRAM

- キャパシタ（コンデンサ）で情報を記憶する
- 格子状にキャパシタを配置する
 - 先頭行から順にアドレスを配置し、埋まったら二行目

		列						
行		1	0	1	1	0	1	アドレス小
		0	1	1	0	0	1	↓ アドレス大
		1	0	1	1	0	0	
		0	1	0	0	0	1	
		0	0	1	1	0	1	
		1	1	0	0	1	1	

24

まず、DRAM についてですが、DRAM はキャパシタを用いて情報を記録するメモリです。キャパシタというのは電荷を蓄えることができる素子で、つまりはコンデンサのことです。コンデンサに電荷が溜まっているかどうかで1ビットを記録できます。電荷が溜まっていれば1、溜まっていなければ0という感じです。DRAM ではこのキャパシタが格子状に配置されています。アドレスが1行目から順に配置され、行番号が大きくなるに従ってアドレスも大きくなっていきます。近年のDRAM ではこの格子が複数枚配置されていて、それぞれをバンクと呼んでいます。その際のアドレス配置について色々な方法が考えられますが、ここでは触れないでおきます。

DRAM

- データ読み込み後、電荷の再注入（プリチャージ）が必要となる
- 何もしなくても一定時間毎に電荷の再注入（リフレッシュ）が必要となる
- これらを行単位で管理する

25

DRAM のハードウェアとしての特徴はいくつかあり、一つ目はデータ読み込み処理の後、電荷の再注入が必要となることです。これはプリチャージと呼ばれています。読み込みの際、コンデンサに溜まっている電荷を吸い取ってしまうのでそれを元に戻してあげるのがプリチャージです。二つ目の特徴として、何もしない場合でも一定時間ごとに電荷の再注入が必要になります。これはリフレッシュと呼ばれています。コンデンサに溜まった電荷は何もしなくても少しずつ放電してしまうので、定期的に再チャージをするということです。DRAM ではプリチャージやリフレッシュを行単位で行います。DRAM のキャパシタが格子状に配置されていることにはこのようなメリットがあります。

DRAMのデータ読み込み手順

1. 行のアクティベート
 - 信号線を指定の行に接続する
2. 指定列の転送
 - アクティベートした行の中から必要な列を転送する
3. プリチャージ
 - 他の行にアクセスする前に電荷のチャージを行う (ディアクティベート)

データ転送以外にも時間がかかる

26

次にDRAMのデータ読み込み手順について見てみます。手順1が行のアクティベートで、信号線をDRAMの指定行に接続する処理です。手順2が指定列の転送で、アクティベートした行の中から、必要な列のデータを転送する処理です。手順3がプリチャージで、行の電荷を再チャージする処理です。これは他の行を読む前に行う必要があります。ディアクティベートとも呼ばれます。

ここで重要なことは、DRAMのデータを読む際、実際のデータ転送以外にも時間がかかるということです。

DRAMへの効率的アクセス

- 一度、行にアクセスしたら、できるだけたくさんのデータを転送した方が有利
- DRAMはバーストリード（行の連続するデータの高速転送）に対応している
 - ⇒ 連続するデータにアクセスした方が高速
 - ⇒ CPUではキャッシュメモリが有効

27

DRAMから高速にデータを転送するためには、一度行にアクセスしたら、できるだけたくさんのデータを転送した方が有利です。DRAMにはバーストリードという機能があり、行の中の連続するデータを高速に転送することができます。通常は読み込む列を指定するための制御信号をDRAMに送信するとそのデータが送られてくるのですが、バーストリードでは、制御信号をやり取りすることなく連続するデータを送信できるので高速なデータ転送が可能になります。

ということで、DRAMでは同じ行の連続するデータにアクセスすると効率的だということがわかります。CPUのキャッシュメモリはDRAMに対しこのようなアクセスをしており、とても効率的だと言えます。

ざっくりとした数字

- クロック周波数を1GHzとすると
クロック周期は1ns
- データアクセスに必要な総時間は10～
100nsのオーダー
 - グローバルメモリアccessのレイテンシは200～
400クロック (@CUDA Programming Guide)
- メモリバンド幅 (最大データ転送レート) は
672GB/sec (@Quadro RTX 6000)
 - 1nsあたり672Byte (4Byte変数168個分)

NVIDIA TURING GPU ARCHITECTURE White paper

https://www.nvidia.com/content/GTC/content_output千万/nvidia.com/content/technologies/turing-architecture/NVIDIA-Turing-Architecture-Whitepaper.pdf

28

ここで時間の数値に関する感覚を持ってもらうために時間に関するざっくりとした数字を見てみたいと思います。まずプロセッサのクロック周波数ですが、だいたい1GHzくらいのオーダーですので、クロック周期は1nsくらいということになります。ざっくりいうと、プロセッサは1nsに1回命令を実行できる能力を持っているということです。DRAM中のデータへのアクセスに必要な時間はざっくり10nsから100nsのオーダーです。CUDA C Programming Guideには、DRAMが使われているグローバルメモリへのアクセスの待ち時間(レイテンシ)は200～400クロックであると書かれています。ざっくり言うと、メモリアccessには他の単純命令の200命令分くらいの時間が必要になると言うことです。DRAMからデータを読み込む際の最大データ転送レートを見てみると例えば、Quadro RTX 6000と呼ばれるモデルでは、672GB/secで、1nsあたりに直すと672Byteです。これは4Byteの変数、168個分に相当します。これは結構速いのではないかと思います。つまりデータの転送自体は速いが、DRAMのaccessには色々な手順があるので、トータルの待ち時間としては遅くなるということです。

DRAMへの効率的なアクセス

1. キャッシュメモリ

- 多くのプログラムにおいて、一度アクセスしたアドレスとその周辺は近いうちにアクセスされる傾向がある（時間的局所性、空間的局所性）
- ある1ワードにアクセスした時、そのワードが含まれるブロック（行の一部）の全データをキャッシュにバースト転送する
- GPUもキャッシュメモリを持っている

29

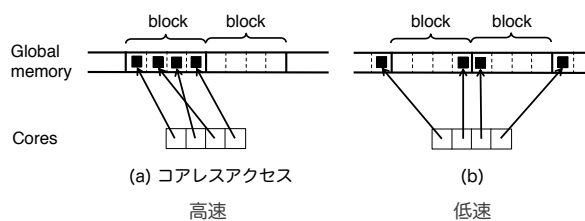
それでは、DRAMに効率的にアクセスするための具体的な方法についてみていきます。3つ紹介しますが、まず最初は、キャッシュメモリです。多くのプログラムでは一度アクセスしたアドレスやその周辺のアドレスには近いうちにアクセスされやすいということが経験的に知られています。これは時間的局所性、空間的局所性と呼ばれるものです。そこで、アドレスを複数のブロックに分割しておき、DRAMのあるアドレスにアクセスした際、そのアドレスを含むブロックをキャッシュに転送します。転送はDRAMのバースト転送を使って高速に行えます。また、キャッシュにはSRAMが使われるため、高速アクセスが可能です。

強調されることが少ないように思いますが、GPUもキャッシュを持っています。あまり強調されないのは、GPUはDRAMに効率的にアクセスするための別の手段を持っていてそれが強調されるためだと思えます。それについては、後で説明します。

DRAMへの効率的なアクセス

2. コアレスアクセス

- 複数のコアが連続するデータにアクセスする時、DRAMのバーストリードにより高速読込できる



30

2つ目は、並列アーキテクチャ限定の話になりますが、コアレスアクセスです。DRAMは複数のブロックに分かれています。が、(a)のように複数のコアがDRAMの同一ブロックにアクセスする時、バーストリードが使えるため、高速にアクセスできます。これをコアレスアクセスと呼んでいます。一方(b)のようにコアが別々のブロックにアクセスする時は、ブロックごとに別々にアクセスしないといけなくなるため低速になります。

例：SIMDアーキテクチャ

- 代表的な並列アーキテクチャの一つ
 - GPUアーキテクチャはSIMDの進化版と言える
- 1命令で複数のデータを並列処理する
 - 全データに対して同一の演算を行う (Single Instruction Multiple Data)
 - 演算器をデータ数分用意しておく
- 連続データを処理するのでコアレスアクセスが可能
 - 各コアが独立に動作するアーキテクチャより効率的

31

コアレスアクセスを活かせるアーキテクチャとして SIMD アーキテクチャがあります。SIMD アーキテクチャは代表的な並列アーキテクチャの一つです。GPU アーキテクチャは SIMD アーキテクチャの進化版であり、SIMD アーキテクチャの多くの特徴を引き継いでいます。SIMD アーキテクチャは複数のデータに対し同一命令の並列実行が可能なアーキテクチャです。そのために複数の演算器が搭載されています。SIMD とは Single Instruction Multiple Data の略です。SIMD アーキテクチャでは連続する領域のデータに対して、同一命令を実行するため、コアレスアクセスが可能です。各コアが独立に動作する並列アーキテクチャでは、アクセスするデータのアドレスはコアによりバラバラになるため、このようなコアレスアクセスはできません。SIMD アーキテクチャには様々な長所がありますが、効率の良いメモリアクセスも長所の一つと言えます。

DRAMへの効率的なアクセス

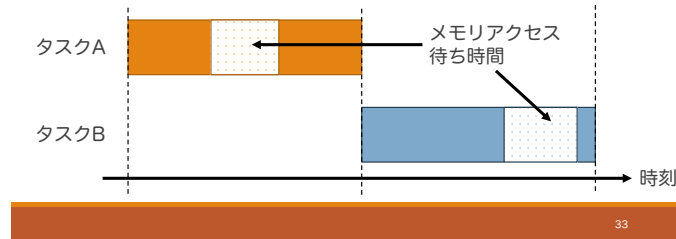
- ### 3. ハードウェアマルチスレッディング
- メモリアクセスによる待ち時間（レイテンシ）を隠蔽できる
 - 隠蔽？
 - 隠蔽したら速くなるの？
- ⇒ 具体例で説明する

32

3つめは、ハードウェアマルチスレッディングです。メモリアクセスによる待ち時間（レイテンシ）を隠蔽する技術で、これがGPUアーキテクチャの大きな特徴なのですが、「隠蔽」とは何か分かりづらく、また、隠蔽したからといって高速になるのかという疑問も湧いてくるところなので、レイテンシ隠蔽による高速化について、具体例を使って説明したいと思います。

レイテンシ隠蔽による高速化

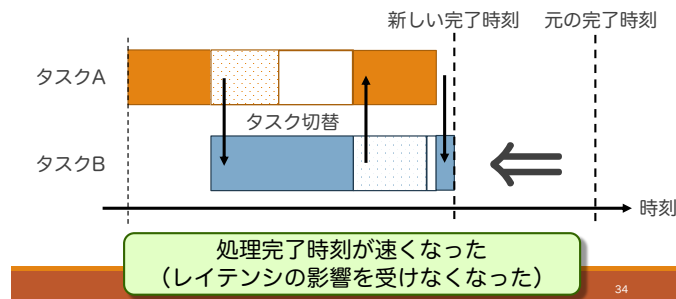
- 独立した二つのタスクA,Bを行うとする
- 1コアで、タスクAを終えた後にタスクBをやる場合、以下のようになったとする



今、独立した二つのタスク A,B を1コアで行うとします。まず一つ目のパターンとして、タスク A を終えた後にタスク B をやることにするとスライドのようになったとします。タスク A,B はともに途中でメモリアクセスが必要で、その際、データを要求してからデータが届くまで長時間待たされるものとします。水玉になっている領域は何もできず待たされている時間帯です。

レイテンシ隠蔽による高速化

- タスクA,Bをいつでも部分的に実行できるとすると・・・



この時、二つ目のパターンとして、コアはいつでもタスク A,B を部分的に実行できるものとします。すると、タスク A のメモリアクセス待ち時間にタスク B を実行できます。また、タスク B のメモリアクセス待ち時間にタスク A を実行できます。このようにすると、タスクの完了時間がメモリアクセス待ち時間の影響を受けなくなり、結果として、タスク A,B の順に実行する場合よりも早く処理が完了します。これがレイテンシの隠蔽と呼ばれるものです。

少し脱線しますが、日常的な例として、無駄な会議を考えてみましょう。無駄な会議はなくするのが一番ですが、様々な理由でなかなか無くせない場合もあります。そこで、別のアプローチとして、会議の前に雑用を溜め込んでおき、会議中に雑用を内職するという方法が考えられます。こうすれば無駄な会議であっても、時間を無駄にせずすみ、結果的に労働時間を減らすことができます。レイテンシの隠蔽とはそのような技術です。

レイテンシ隠蔽の実装

1. 複数のタスクをマルチスレッド化する
⇒ ソフトウェア側で頑張る
 2. スレッドを高速に切り替えられるようにする
⇒ ハードウェアがクロック単位で切り替える
- **ハードウェアマルチスレッディング**
 - ハードウェア管理によるマルチスレッディング
 - スレッドの高速な切り替えを可能にする

35

レイテンシ隠蔽を行うために必要なことを見てみます。まず、ソフトウェア側に必要な対応として、複数のタスクをマルチスレッド化する必要があります。先ほどの無駄な会議の例では、雑用を溜め込んでおくということに相当します。次にハードウェア側に必要なこととして、処理するスレッドを高速に切り替えられるようにする必要があります。スレッドの切り替えに時間がかかると、その度に待ち時間が発生してしまうので、クロックごとに実行するスレッドを切り替えられることが望ましいです。

それを実現する技術がハードウェアマルチスレッディングです。ハードウェア管理によるマルチスレッディングであり、実行スレッドの高速な切り替えを可能にします。CPUが搭載されたコンピューターでは、一般的にはOSがスレッドの切り替えを行います。しかしそのようにするとOSソフトウェアがスレッドを管理するための処理が待ち時間となり、頻繁にスレッドを切り替えると待ち時間が増えていきます。GPUでは、複数スレッドが実行可能状態で待機できるようにし、クロック毎にハードウェアが実行スレッドを選択できるようにすることで、スレッド切り替えの待ち時間が少なくなるようにしています。

SRAM

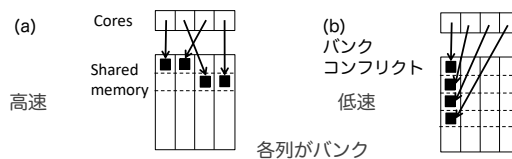
- 半導体（フリップフロップ回路）により情報を記憶する
- 高速アクセスが可能（1～数クロックで1ワード(通常64bit)を取得可能）
- 通常、複数のバンクで構成される
- マルチプロセッサの共有メモリがSRAM

36

それでは、次に SRAM に移りたいと思います。SRAM は DRAM よりも高速なので、計算時間に対するインパクトは少ないですが、それでも高速アクセスのためのコストがありますので、それを話したいと思います。まず、SRAM とは半導体のフリップフロップ回路を用いて情報を記録するメモリのことです。高速アクセスが可能であり、だいたい 1 クロックから数クロックで 1 ワードを取得できます。1 ワードというのは 64bit アーキテクチャなら 64bit のことです。並列アーキテクチャでは、SRAM は通常複数のバンクから構成されます。つまり、SRAM は物理的に分けられた複数のメモリブロックから構成され、そのそれぞれをバンクと呼びます。GPU では、マルチプロセッサ内の共有メモリに SRAM が使われています。

SRAM

- 各バンクに同時にアクセスできる
 - 複数のコアが別々のバンクにアクセスする時、同時に処理が可能
 - 逆に複数のコアが同一バンクにアクセスすると、待ちが発生する→バンクコンフリクト

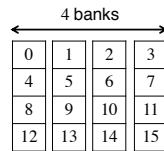


37

バンクのそれぞれは各コアと別々のバス（つまりデータ送信用の信号線）で繋がれており、各コアは異なるバンクに並列にアクセスすることができます。例としてコア数と共有メモリバンク数が共に 4 の場合で考えてみます。共有メモリの各列が 1 つのバンクだとします。(a) のように 4 つのコアがすべて別々のバンクにアクセスする時、並列にデータを取得できるため高速です。しかし、並列にアクセスできない場合もあります。複数のコアが同じバンクにアクセスすると、並列処理できません。(b) のように 4 つのコアがすべて同じバンクにアクセスすると並列処理は行われず、各コアにシーケンシャルにデータが転送されることとなります。このように複数のコアが同じバンクにアクセスすることをバンクコンフリクトと呼びます。つまり、共有メモリに高速にアクセスするには極力バンクコンフリクトを避ける必要があるということです。

SRAM上のアドレス配置

- 連続するアドレスが異なるバンクになるように配置される
⇒ 各コアが連続するワードにアクセスすることでバンクコンフリクトを回避できる



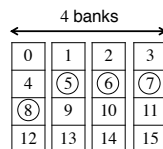
4バンクSRAMでのアドレス配置

38

バンクコンフリクトを避ける方法について考えるためにまず、SRAM上のアドレスの配置について説明します。SRAM上のアドレスは連続するアドレスが異なるバンクになるように配置されます。例えばバンク数が4つの場合はスライドのようになります。なので、例えば4つのコアが連続するアドレスにアクセスするような場合にはバンクコンフリクトは起きません。

例：SIMDとバンクコンフリクト

- 処理するデータを連続データにすることでバンクコンフリクトを回避できる



4並列SIMDでのバンクコンフリクト回避例
(○がアクセスするワード)

39

再び SIMD アーキテクチャの例を話すと、処理するデータを連続データにしている限りはバンクコンフリクトを回避できます。例えば、4並列の SIMD でアドレス 5,6,7,8 にアクセスする場合にはバンクコンフリクトは起きません。この意味でも SIMD アーキテクチャは効率的なメモリアクセスが可能であると言えます。

メモリアクセスまとめ

- DRAMもSRAMも連続データに同時アクセスすることで効率の良いデータ転送が行える
 - SIMDアーキテクチャによりこれを実現できる
- 効率的なDRAMへのアクセスの別のアプローチとして、ハードウェアマルチスレッディングがある
 - レイテンシ（待ち時間）を隠蔽できる

40

効率的なメモリアクセスの方法についてまとめてみます。まず、DRAMもSRAMも連続するアドレスに同時アクセスすることで、高速にデータ転送が行えます。SIMDアーキテクチャでは、これを活かして効率の良いメモリアクセスを実現しています。

DRAMに効率よくアクセスするための別の手段として、ハードウェアマルチスレッディングがあります。これはデータ転送自体を高速化するのではなく、レイテンシを隠蔽するというものです。繰り返しになりますが、レイテンシの隠蔽とは、待ち時間に他の仕事をする事で待ち時間を無駄にしないという技術です。

1.4.3 マルチプロセッサ

GPUマルチプロセッサの概要

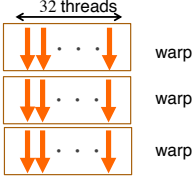
- NVIDIA社はSIMTと呼んでいる
- 特徴
 - SIMDと同様の効率的なメモリアクセス
 - ハードウェアマルチスレッディングによるレイテンシの隠蔽
- 「ヘネ&パタ」ではマルチスレッドSIMDプロセッサと呼ばれている

ヘネシー&パターソン コンピュータアーキテクチャ 定量的アプローチ第5版
<https://www.shoeshisha.co.jp/book/detail/9784798126234> 41

さて、それでは GPU のマルチプロセッサの説明に移りたいと思います。GPU マルチプロセッサの長所を理解することが今日の目的であり、今日のメインパートになります。GPU マルチプロセッサの方式ですが、NVIDIA 社は SIMT と呼んでいます。SIMT とは Single Instruction Multiple Threads の略です。その特徴ですが、SIMD アーキテクチャと同様の効率的なメモリアクセスとハードウェアマルチスレッディングによるレイテンシ隠蔽の機能の両方を備えていることです。コンピュータアーキテクチャの標準的な教科書の一つである「ヘネ&パタ」ではこれをマルチスレッド SIMD プロセッサと呼んでいます。つまりハードウェアによるマルチスレッディングが可能な SIMD プロセッサということです。

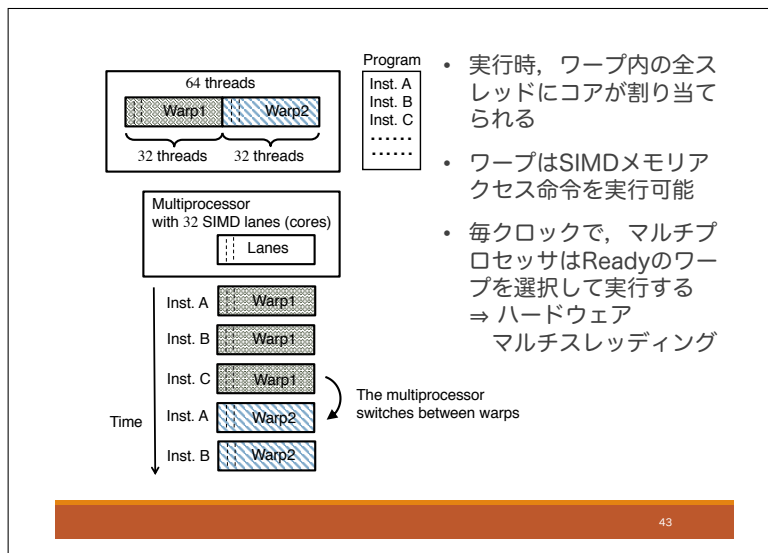
GPUマルチプロセッサの概要

- 各マルチプロセッサに大量のスレッドを割り当てる
- 全スレッドが異なるデータに対して同じ処理を実行する
- スレッドは32スレッドごとに1ワープとしてまとめられる
 - 32コアに1ワープ (32スレッド) を割り当て ⇒ SIMD命令



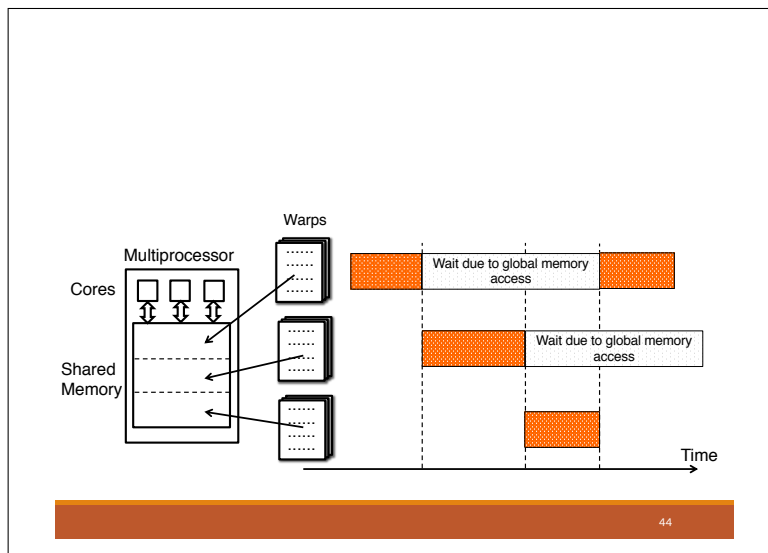
42

それでは、それをどのように実現するのかをみていきましょう。まず、各マルチプロセッサに大量にスレッドを割り当てます。もちろんソフトウェア側でたくさんスレッドを生成しておかないと大量のスレッドを割り当てることはできませんが、そのようにソフトウェアを設計したものとします。異なるデータに対して、これらの全スレッドが同じ処理を行います。SIMD ではすべてのコアが同じ処理を行います。今回の SIMT ではコアからスレッドに置き換わっています。マルチプロセッサに割り当てられたスレッドは 32 スレッドごとにまとめられます。この 32 スレッドのまとまりのことをワープと呼びます。そして基本的にはマルチプロセッサ内の 32 個のコアに 1 ワープ (32 スレッド) を割り当てます。ワープ内の全スレッドは常に同じ処理を行うので、マルチプロセッサはワープに対して、SIMD 命令を実行することになります。そして、この SIMD 命令をすべてのワープに対して実行することで、全スレッドに対して同じ処理を行うことができます。



次にワープのコアへの割り当てについて具体例を使って説明します。やや複雑ですが、大事な部分なのでじっくりみていきましょう。あるマルチプロセッサに64個のスレッドが割り当てられており、すべてのスレッドは命令A、命令B、命令C…という同じ処理を行うものとします。本当は1024個とか2048個とかくらい割り当てた方がいいのですが、ここではわかりやすくするため、64個とします。前のスライドで話した通り、マルチプロセッサへの割り当てスレッドは32スレッドずつワープとしてまとめられます。なので、この例では2つのワープがあることになります。これらのスレッドをマルチプロセッサ内の32個のコアで実行することを考えます。

この時、各ワープを交互に実行する必要はありません。ハードウェアが毎クロックにおいて、実行可能なワープの中からどのワープを実行するかを決定します。スライドの例では一つ目のワープの命令を3つ続けて実行したのち、二つ目のワープの命令を2つ続けて実行しています。これはSIMDマシンに対するハードウェアマルチスレッディングと言えます。通常のハードウェアマルチスレッディングではクロックごとに実行するスレッドを1つ選択していましたが、GPUではクロックごとに実行する32スレッドを選択します。このように32スレッド単位でガバッと入れ替えられるようにすることで、SIMDとハードウェアマルチスレッディングを両立しています。



GPU マルチプロセッサがハードウェアマルチスレッディングによりグローバルメモリアクセスのレイテンシを隠蔽する例をみてみましょう。一番上のワープがSIMD 命令によりグローバルメモリアクセスを行なったことで待ち時間に入ったとします。するとマルチプロセッサは他の実行可能なワープ、例えば真ん中のワープを実行します。そのワープもグローバルメモリへのSIMD アクセスにより待ち時間に入ったとするとマルチプロセッサは3つ目のワープを実行します。そして、一つ目のワープが待ちから復帰すると、再び一つ目のワープを実行します。このような感じで、グローバルメモリにSIMD 命令で効率よくアクセスしながらもレイテンシを隠蔽することができます。

GPUマルチプロセッサまとめ

- スレッドをワープ単位で管理することで
 - SIMDメモリアクセス命令が実行できる
 - ハードウェアマルチスレッディングが行える
- プログラム時の注意点
 - コア数よりも多くのスレッドを割り当てないとハードウェアマルチスレッディングが行えない
 - 1SMあたりのスレッド数は仕様上許される最大数にするのが良い (最大数に対するスレッド数の割合をオキュパンシと呼ぶ)
 - SIMDメモリアクセス命令においては、基本的に連続するデータを指定しないと効率的にならない

45

GPU マルチプロセッサについてまとめると、スレッドをワープ単位で管理することが特徴です。これにより、SIMD メモリアクセス命令に対応しつつも、ハードウェアマルチスレッディングも行えます。

これを踏まえてプログラミングする際の注意点は以下の通りです。まず、ハードウェアマルチスレッディングを行うためにコア数よりも多くのスレッドを生成しないとイケないということです。アーキテクチャごとにマルチプロセッサに割り当て可能なスレッドの最大数が決まっていますが、レイテンシ隠蔽の効果を大きくするためには、最大数のスレッドを割り当てるようにするのが良いです。最大数に対する実際の割り当てスレッド数の割合をオキュパンシと呼んでいます。なので、オキュパンシを100%にするのが良いということです。オキュパンシを大きくする際の障壁については二日目に話します。次にワープ内のスレッドは基本的にはメモリ内の連続する領域にアクセスするのが良いです。これはメモリのところで説明した通りです。

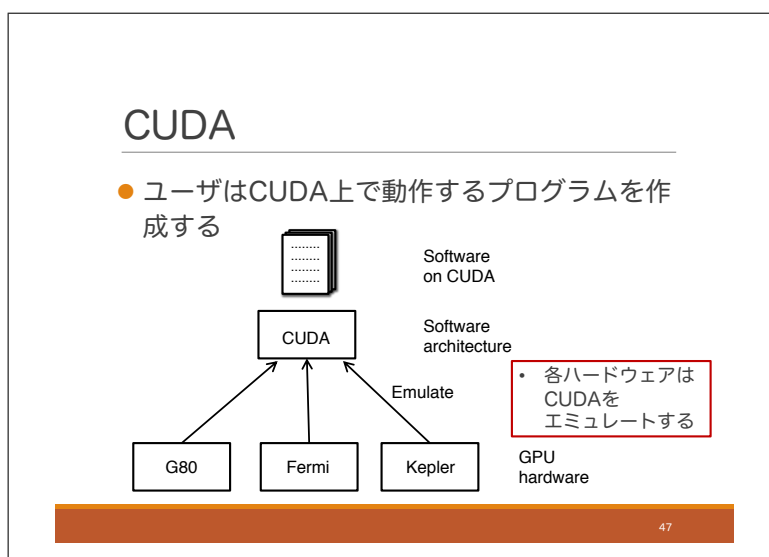
1.5 GPUプログラミング

GPUプログラミング

- NVIDIA GPUは世代やモデルごとに仕様が大きく変わる
- すべてのGPU上で動作するようなプログラムを書きたい
 - ⇒ CUDA (Compute Unified Device Architecture)
- NVIDIA GPU向けの計算モデル, 開発環境

46

さて、ここからはGPUプログラミングの説明をしていきます。まずは全体的な話をしますが、NVIDIA社のGPUは世代やモデルごとに仕様が大きく変わります。その際、異なるGPUを使うたびにプログラムを書き換えないといけないのでは大変ですので、すべてのGPUで動作するようなプログラムが書けることが望ましいです。そのようなことが可能になる仕組みがCUDAです。CUDAはCompute Unified Device Architectureの略ですが、アーキテクチャとなっていますが、それだけでなく、GPUプログラムの開発環境などもすべて含めてCUDAと呼ばれています。なお、CUDAはNVIDIA社GPUのためのものであり、他社GPUはサポートされていません。



CUDAのアーキテクチャの概要はこのようになります。まずCUDAはソフトウェアアーキテクチャとなっており、ソフトウェアはCUDA上で動作するように作成します。そして、GPUハードウェアはそれぞれCUDAをエミュレートします。このようにすることで、GPUハードウェアの仕様が異なっても同一のプログラムを実行できるようになっています。

CUDA

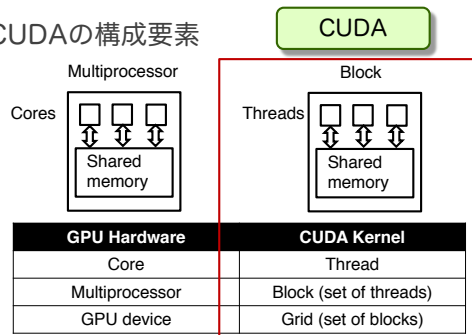
- Compute Capability
 - CUDAのバージョン番号であり、この値ごとに
対応する仕様（対応可能なGPU）が決まっている
 - Compute Capabilityを指定してコンパイルすること
で、最新の機能を使用することができる

48

どの GPU でも動かせるプログラムを書きたいという話をしましたが、実際にはそうでない場合もあり、例えば、最新 GPU の新しい機能を使って特別なことをやりたいということもあります。このような場合に対応するため、CUDA では対応する機能や仕様のセットに対してバージョン番号をつけており、それを Compute Capability と呼びます。各 GPU モデルに対し、そのモデルに付けられた Compute Capability を確認することで対応する機能や仕様を知ることができます。また、Compute Capability を指定してコンパイルすることで、一部の GPU でしか対応していない機能を使うことができます。

GPUデバイスとCUDA

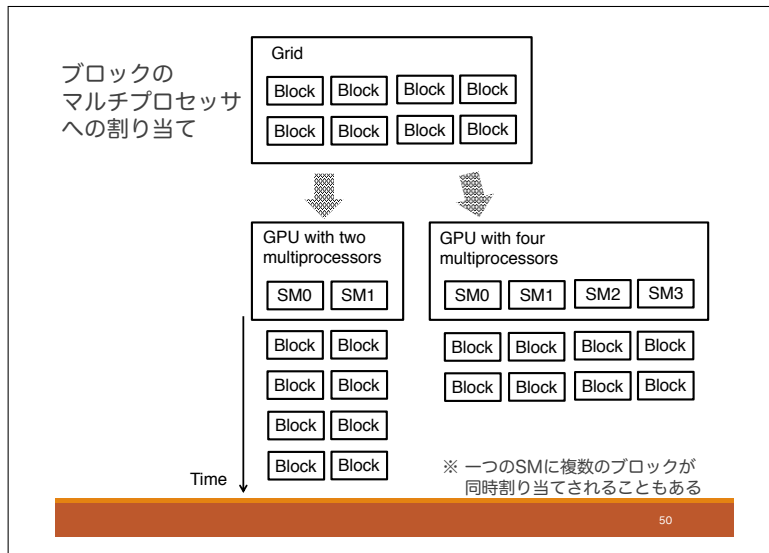
- CUDAの構成要素



49

CUDA アーキテクチャの中身について見てみましょう。CUDA では様々なスレッドの集合が定められています。GPU ハードウェアとの対応を見ながら見ていきましょう。

まず、コアに対応するのはスレッドです。実際には、一つのコア上で複数のスレッドが実行されます。次にマルチプロセッサに対応するのはブロックです。ブロックとは同一のマルチプロセッサで実行されるスレッドの集合です。実際には、一つマルチプロセッサ上で複数のブロックが実行されます。しかし、一つのブロックが複数のマルチプロセッサで実行されることはありません。最後に GPU デバイスに対応するのが、グリッドです。グリッドはブロックの集合です。



CUDA 上で開発されたプログラムが GPU 上でどのように実行されるかについて見てみましょう。例として、1 グリッド、8 ブロックからなる CUDA プログラムを考えます。まず、このプログラムを 2 つのマルチプロセッサからなる GPU で実行する場合には、左の矢印のように一つのマルチプロセッサが 4 つのブロックを担当することになります。もし 4 つのマルチプロセッサからなる GPU で実行する場合は、右の矢印のように一つのマルチプロセッサが 2 つのブロックを担当することになります。このように GPU の構成が変わっても同じプログラムを実行できます。

なお、このスライドでは各マルチプロセッサはブロックを一つずつ計算していますが、実際は一度に複数のブロックが割り当てられることが多いです。

1.6 実習

実習

- CUDAプログラムを実習環境で実行する
- 今日の目標
 - サンプルプログラムを実行できるようになる
- 二日目の目標
 - 基本的な文法を理解する
 - 高速化のノウハウを学ぶ

51

それでは実習に移りましょう。実習では CUDA 上で開発したプログラムを実習環境で実行します。今日の目標は、サンプルプログラムを実行できるようになることです。二日目は CUDA の基本的な文法を理解し、高速化のノウハウを身に着けることを目標とします。

実習のための基本事項

- 大規模並列計算システムOCTOPUSを使う
- 詳細は大阪大学サイバーメディアセンターのマニュアルを参照のこと
<http://www.hpc.cmc.osaka-u.ac.jp/system/manual/octopus-use/>
- 今回はコマンドベースでコード作成やプログラム実行を行う

52

実習のための基本事項をいくつか述べます。まず実習環境としては、大阪大学サーバーメディアセンターの大規模並列計算機システム OCTOPUS を使います。詳細はサイバーメディアセンター作成のマニュアルをご覧ください。また、今回はコマンドベースでコード作成やプログラム実行を行います。GUI が使える場合には Nsight と呼ばれる開発プラットフォームを使うのが便利ですが、今回は使用しません。

実習1：サンプルファイルを実行する

- 本計算機環境ではファイル実行の仕方が通常とは異なる
 - ジョブスケジューラにジョブを投げる
- 本実習ではビルド済みの実行ファイルを実行する
 - deviceQueryと呼ばれるプログラムを実行する
 - CUDAにおいて、最初に行われるサンプルの一つ
 - 使用するGPUの仕様を出力するプログラム

53

それでは、本日の実習としてサンプルファイルを実行してみます。ここで、注意事項ですが、本計算機環境ではファイル実行の仕方がパソコンで行う場合とは異なります。通常はコマンドラインから実行ファイル名を打ち込むだけですが、今回はそうではありません。本環境は複数のユーザが共同利用しているので、ジョブスケジューラと呼ばれる各ユーザのジョブを管理するソフトウェアが導入されています。プログラムを実行するにはジョブスケジューラにジョブを投げる必要があります。

この実習ではビルド済みの実行ファイルを実行してみます。実行するのは deviceQuery と呼ばれるプログラムで、これは CUDA において最初に実行されることが多いサンプルの一つです。使用する GPU の仕様を出力するプログラムになっており、この後のプログラミングにも役立ちます。

実習1：サンプルファイルを実行する

- 手順1：deviceQueryをユーザ領域にコピーする
 - `$ cd ~/`
 - `$ cp /octfs/apl/CUDA/cuda-9.0/samples/1_Uilities/deviceQuery/deviceQuery ~/`
- 手順2：ジョブスケジューラにジョブを登録するためのスクリプトを作る
 - `$ vi deviceQuery.nqs`
#!/bin/bash
#PBS -q LECTURE
#PBS -l elapstim_req=00:00:10,cpunum_job=1,gpunum_job=1
cd \$PBS_O_WORKDIR
./deviceQuery

54

それでは実際にやってみましょう。計算機システムにログインしたのち、スライドの手順1,2をやってみてください。手順1はサンプルファイルを自分のホームディレクトリにコピーする処理です。手順2ではジョブスケジューラにジョブを投げるためのスクリプトを作ります。viを立ち上げたのち、その下5行を打ち込んで保存してください。

なお、Linuxの操作に不安のある方は、付録のLinux入門の説明をチェックしてみてください。cd, cp コマンドやviについては後半の「基本的なLinuxコマンド」のところに説明があります。

実習1：サンプルファイルを実行する

```
#!/bin/bash      ← bashを使う（シバン）
#PBS -q LECTURE  ← 講習会用の待ち行列に並ぶ
#PBS -l elapstim_req=00:00:10,cpunum_job=1,
gpunum_job=1    ← 10秒間CPU1ノードGPU1台を使用する
cd $PBS_O_WORKDIR
./deviceQuery   ↑ ジョブ投入時のディレクトリに移動する
                ↑ 実行ファイルdeviceQueryを実行する
```

55

手順2で作成したスクリプトについて、少し中身を見てみましょう。1行目はbashのパスを指定しています。いわゆるシバンと呼ばれるもので、おまじないだと思って書いておきましょう。2行目はジョブの順番待ちのために並ぶ待ち行列を指定しています。この講習会ではLECTUREという特別な待ち行列を作っていますので、そこに並びます。3行目ではジョブ実行について設定しています。ここではCPU1ノードとGPU1台を10秒間使う設定になっています。下2行は、ファイルを実行するフォルダに移動して、実際にファイル実行する処理です。

実習1：サンプルファイルを実行する

- 手順3：スクリプトを実行しジョブを投入する
 - \$ qsub deviceQuery.nqs
 - 手順4：ジョブが終わるまで待つ
(コマンドは後述)
 - 実行結果を確認する
 - \$ less deviceQuery.nqs.o(リクエスト番号) ⇒標準出力
 - \$ less deviceQuery.nqs.e(リクエスト番号) ⇒標準エラー出力
- ⇒ 標準出力にGPU情報が記載されていれば成功

56

それでは実際にジョブを投入してみましょう。手順3によりジョブが投入されます。その後手順4として実行が終わるまで待ちましょう。ジョブの状態を確認するコマンドは次のスライドに書いてあります。最後に実行結果を確認します。実行結果ですが、端末に表示されるのではなく、ファイルに出力されます。プログラムの出力には、標準出力と標準エラー出力がありますが、スクリプト名の後ろに、ドット+o+リクエスト番号が書かれたものが標準出力です。また、スクリプト名の後ろに、ドット+e+リクエスト番号が書かれたものが標準エラー出力です。標準出力にGPUの仕様の情報が書かれていれば成功です。

実習1：サンプルファイルを実行する

- 参考：ジョブリクエスト関連のコマンド
 - \$ qstat ← リクエストの状態確認
 - \$ sstat ← 実行中ジョブの状態確認
 - \$ qdel 123456.oct ← リクエストの削除 (適宜数字変更)
 - \$ sstatall | less ← 投入されているすべてのリクエストの状態
- 詳細は大阪大学サイバーメディアセンターのマニュアルを参照のこと
<http://www.hpc.cmc.osaka-u.ac.jp/system/manual/octopus-use/scheduler/>
- OCTOPUS利用状況 (要ログイン)
<https://portal.hpc.cmc.osaka-u.ac.jp/secure/oct-smap/oct-smap.html>

57

最後に、リクエストしたジョブの状態を確認するコマンドを紹介します。まず、qstat は投入したジョブの状態を確認するコマンドです。sstat は実行中のジョブの詳細情報を確認するコマンドです。投入したジョブを削除するには qdel を使います。自分が投入した以外のジョブも含めてすべてのジョブのスケジューリング状態を調べるには sstatall を使います。詳細についてはサイバーメディアセンター作成のマニュアルをご覧ください。また、OCTOPUS の利用状況について、Web ページで確認できます。混雑状況の確認などにご利用ください。


第2章 二日目

2.1 はじめに

数理工学講習会

GPUプログラミング入門

二日目：GPUプログラミング（実習）
一関高専 未来創造工学科 情報・ソフトウェア系
小池 敦



1

GPUプログラミング入門

- 目標
 - GPUアーキテクチャの特性を理解した上でGPUを用いた並列化ができる
- 対象者
 - C言語等でプログラミングをしたことがあり、GPUを用いた高速化に興味がある
 - GPUによる並列計算の概要を知りたい
- 必要な予備知識
 - C言語を触ったことがある

2

二日目の講義を始めます。初日の繰り返しになりますが、本講習会は、GPUを使って高速な並列プログラムを書きたい人のための二日間の入門講座です。目標は、単にプログラミングの文法を知るだけでなく、GPUのアーキテクチャを理解した上で、高速に動作するプログラムが書けるようになることです。対象者は、C言語でプログラミングをしたことがあり、GPUを用いた高速化に興味がある人です。C言語については、簡単なコードが読めるくらいの知識があれば大丈夫です。また、単にGPUによる並列計算の概要を知りたいという人も対象としています。

GPUプログラミング入門

- 講習会の内容
 - 初日：GPUアーキテクチャとGPUプログラミング（座学中心）
 - 二日目：GPUプログラミング（実習中心）
- 必要な事前準備
 - 大阪大学サイバーメディアセンターの大型計算機試用アカウントを取得し、OCTOPUSにログインできるようにしておく
<http://www.hpc.cmc.osaka-u.ac.jp/service/intro/shiyo/>

3

講習会の内容はスライドの通りです。初日は、GPUのアーキテクチャを中心に話しました。初日は座学中心でしたが、最後にサンプルプログラムを動かしてみました。二日目は実際にGPUのコードを書きながら、プログラミングの文法と高速化のノウハウを学習します。

初日の繰り返しになりますが、本テキストは大阪大学サイバーメディアセンターでの講習会で使うことを想定して作られています。もし他の環境で使用している場合は、適宜読み替えてください。サイバーメディアセンターを利用する場合、GPUが使用できるシステムはOCTOPUSです。事前にアカウントを取得し、OCTOPUSにログインできるようにしておいてください。

二日目の内容： GPUプログラミング（実習）

- Hello world
 - 環境設定
 - CUDAの基本的な文法
 - CUDA C++とCUDA Fortran等があるが、今回はCUDA C++ (C++に対するCUDA拡張)を扱う
- 高速アルゴリズムの設計
 - 基本的な考え方
 - 例：リダクション
- 高速化関連トピックス

4

二日目の内容の詳細はスライドのようになります。

まず、ハローワールドということで、環境設定をしたのち、簡単なプログラムを書きながらCUDAの基本的な文法を学習します。CUDAのプログラミング言語は既存の言語をGPUが使えるように拡張する形で作られています。C++、Fortran、Pythonなどから拡張されていますが、今日はC++をGPU向けに拡張した言語であるCUDA C++を扱います。その後、高速なアルゴリズムを設計するための考え方をリダクションと呼ばれる問題を例にして説明します。最後に様々な高速化関連トピックを扱います。

2.2 使用する GPU の仕様

本日使用するGPU

- NVIDIA Tesla P100
- Pascalアーキテクチャ
- Whitepaper
<https://images.nvidia.com/content/pdf/tesla/whitepaper/pascal-architecture-whitepaper.pdf>
- 主要な仕様はdeviceQuery（初日）の結果に書かれている
 - `$ less deviceQuery.nqs.o`(リクエスト番号)

プログラミングを始める前に、本日使用する GPU の仕様をみておきましょう。大型計算機システム OCTOPUS には NVIDIA Tesla P100 という GPU が搭載されています。2016 年に発表された Pascal アーキテクチャを採用したモデルです。詳細はスライドに記載のホワイトペーパーをご覧ください。仕様については、初日に実行した deviceQuery と呼ばれるサンプルファイルを実行することでも取得することができます。初日にこのファイルを実行済みの場合には、その結果を見てみましょう。

Tesla P100の主な仕様

- トータルコア数：3584
 - 56マルチプロセッサ × 64(コア/マルチプロセッサ)
- グローバルメモリサイズ：16GB
- マルチプロセッサあたりの共有メモリサイズ：64KB
- マルチプロセッサあたりの最大スレッド数:2048（最大ワープ数:64）
- ブロックあたりの最大スレッド数：1024

主な仕様について、スライドに書いておきます。まず、1 デバイスあたりのコア数は 3584 です。内訳としては、56 個のマルチプロセッサが、それぞれ 64 個のコアを持っています。メモリについては、グローバルメモリのサイズが 16GB、マルチプロセッサあたりの共有メモリのサイズが 64KB です。

また、1 マルチプロセッサに最大 2048 個のスレッドを割り当てることができます。これをワープ数に直すと 64 個になります。また、1CUDA ブロックあたりの最大スレッド数は 1024 です。これについては後でできます。

2.3 環境設定

- CUDAファイルのコンパイルはnvccコマンドで行うため、そのコマンドへのパスを通す

```
$ echo 'export PATH=/octfs/apl/CUDA/cuda-9.0/bin:$PATH' >> ~/.bash_profile
$ source ~/.bash_profile
```

- 大阪大学サイバーメディアセンター 大規模計算機システムの「OCTOPUSの利用方法」を参照のこと
<http://www.hpc.cmc.osaka-u.ac.jp/system/manual/octopus-use/>
 - フロントエンド利用方法：
<http://www.hpc.cmc.osaka-u.ac.jp/system/manual/octopus-use/frontend/>

それでは、実習を始めましょう。まず環境設定をします。CUDA プログラムをコンパイルする際、nvcc というコマンドを使いますが、そのコマンドへのパスを通します。OCTOPUS 環境であれば、スライドのようにコマンドを打ってください。「.bash_profile」という設定ファイルにパス情報を追記して、その設定ファイルの再読み込みを行っています。

環境設定は以上ですが、OCTOPUS を使用する際はその前に関連ドキュメントを確認しておきましょう。その URL を書いておきます。

2.4 実習1：最初のプログラム

実習1. 最初のプログラム

- 問題
 - 16要素からなる整数配列のそれぞれの要素にindex番号（0～15）を書き込み表示しなさい

0
1
⋮
⋮
15

7

まず最初の実習として、スライドのような問題を考えます。16 要素の配列のそれぞれに index の値を書き込むという問題ですが、これを並列化しない場合とする場合の 2 通りでプログラミングしようと思います。

2.4.1 並列化しないプログラム

C++プログラム（並列化なし）による計算

ソースコード (ex1-c.cpp)

```
1  #include <stdio.h>
2
3  void WriteIndex(int* a, int n)
4  {
5      for (int i = 0; i < n; i++) {
6          a[i] = i;
7      }
8  }
9
10 int main()
11 {
12     // Init
13     int N = 16;
14     int* h_a = new int[N]; // on host (CPU)
15
16     // Main
17     WriteIndex(h_a, N);
18
19     // Display results
20     for (int i = 0; i < N; i++) {
21         printf("%d\n", h_a[i]);
22     }
23
24     // Free memory
25     delete [] h_a;
26
27     return 0;
28 }
```

まずは、並列化しない場合のソースコードを書いてみましょう。ファイル名を ex1-c.cpp とします。その後 CUDA プログラムを作ること考えて若干変なコードになっていますが、普通に読めると思います。

次に、このプログラムを実習環境で実行してみます。初日に話したように実習環境でプログラムを実行するには、ジョブを投入して待ち行列に並ぶ必要があります。

ジョブスクリプト (cuda.nqs)

- 10秒制限
- 実行ファイル名をa.outにする
- 講習会用の待ち行列 (LECTURE) に投入する
- 今回GPUは使用しないが気にしない🙄

```
1  #!/bin/bash
2  #PBS -q LECTURE
3  #PBS -l elapstim_req=00:00:10,cpunum_job=1,gpunum_job=1
4  cd $PBS_O_WORKDIR
5  ./a.out
```

まず、ジョブスクリプトを書きましょう。vi等のエディタで上記の内容を書いて保存しましょう。ファイル名はなんでもいいですが、ここでは cuda.nqs とします。制限時間を 10 秒とし、講習会用専用の待ち行列である LECTURE にジョブを投入することにします。

コンパイルとジョブ投入

- コンパイル(g++)により実行ファイルa.outが生成される

```
$ g++ ex1-c.cpp
$ qsub cuda.nqs
```

それではコンパイルして、ジョブを投入しましょう。今回は GPU を使用しない普通の c++プログラムなので、g++コマンドを使ってコンパイルすることにします。上記の g++コマンドを端末で実行すると a.out というファイルが出来上がります。次に qsub コマンドを用いて、ジョブを投入します。初日に反した通り、ジョブ投入後は qstat コマンドでリクエストの状態を確認できます。また、実行中のジョブに関しては、sstat で状態を確認できます。

2.4.2 CUDA による GPU プログラム (その 1:1 ブロックのみ使用)

ソースコード (ex1-cuda.cu)

```
1  #include <stdio.h>
2
3  __global__ void WriteIndex(int* a, int n)
4  {
5      int i = threadIdx.x;
6      if (i < n) {
7          a[i] = i;
8      }
9  }
10
11 int main()
12 {
13     // Init
14     int N = 16;
15     int* h_a = new int[N]; // on host (CPU)
16     int* d_a;              // on device (GPU)
17     cudaMalloc(&d_a, N * sizeof(int));
18
19     // Copy data from host to device
20     // Do nothing
21
22     // Main
23     int blocks = 1;
24     int threadsPerBlock = N;
25     WriteIndex<<<blocks, threadsPerBlock>>>(d_a, N);
26
27     // Copy data from device to host
28     cudaMemcpy(h_a, d_a, N * sizeof(int), cudaMemcpyDeviceToHost);
29
30     // Display results
31     for (int i = 0; i < N; i++) {
32         printf("%d\n", h_a[i]);
33     }
34
35     // Free memory
36     delete [] h_a;
37     cudaFree(d_a);
38
39     return 0;
40 }
```

CUDAプログラムによる計算（1ブロックのみ使用）

- 本プログラムでは同一ブロック内の16スレッドが、それぞれ1要素の書き込みを行う
- 同一ブロック内のスレッドは必ず同一マルチプロセッサ上で動作することに注意する

それでは CUDA による GPU プログラムを書きましょう。ファイル名を ex1-cuda.cu とします。CUDA プログラムでは拡張子を cu にすることが多いです。エディタで前ページのソースコードを入力してください。

ソースコード解説

- 3-9行目：デバイス(GPU)上で実行される関数
 - デバイス上で実行される関数のうち、ホスト(CPU)から呼ばれるものには `__global__` をつける
 - 参考：デバイス内で呼び出される関数には `__device__` をつける
 - ブロック内のスレッドには `0` から `blockDim.x - 1` (ブロック内のスレッド最大値 - 1)までのIDが割り当てられ、スレッドIDは `threadIdx.x` で参照できる。
- 16,17行目：デバイス上のグローバルメモリにメモリ確保するには `cudaMalloc` 関数を使用する
- 19,20行目：通常、GPUでの処理の前にホストのデータをデバイスに転送する処理を行うが、今回は不要
- 25行目、GPU関数（カーネル関数）を呼び出すには、関数名の後に三重ブラケットをつける。 `<<<blocks, threadsPerBlock>>>` と書いた時、ブロックの数が `blocks` になり、1ブロックあたりのスレッド数が `threadsPerBlock` になる。今回は1ブロック16スレッドにする。
- 27,28行目：GPU処理が完了したら、デバイスのデータをホストに転送する。そのために `cudaMemcpy` 関数を使用する。
 - 第二引数から第一引数へデータ転送される
 - 参考：ホストからデバイスにデータ転送するには `cudaMemcpyDeviceToHost` を `cudaMemcpyHostToDevice` にする
- 37行目：デバイス上のメモリを解放するためには、 `cudaFree` 関数を用いる

ソースコードを解説します。まず、3行目ですが、関数の戻り値指定の前に「`__global__`」がついています。これはCPU（ホスト）からGPU（デバイス）に対して呼び出される関数、つまりカーネル関数であることを示しています。ちなみにGPU（デバイス）の内部で使用される関数の場合は、関数の戻り値指定の前に「`__device__`」をつけます。GPU内の各スレッドがこの関数を個別に処理します。スレッドには0始まりのIDがついていて、今回はスレッド0からスレッド15までの16個のスレッドが作られます。この指定方法についてはすぐ後に説明します。このスレッドIDは `threadIdx.x` で参照できます。本関数ではスレッド i (i は0から15)が、個別に $a[i]$ に i を書き込みます。

次に `main` 関数を見てください。こちらはCPUで実行される処理です。まず、16、17行目について、デバイスのグローバルメモリにメモリ領域を確保する処理が「`cudaMalloc`」です。C言語のメモリ確保命令である `malloc` と違い、GPUデバイスの中にメモリを確保します。同様に、37行目の `cudaFree` はC言語の `free` と違い、GPUデバイス中のメモリを解放します。

初日にGPUプログラムの基本的なステップは、(1) 入力データのCPUからGPUへのコピー、(2) GPU上での関数実行、(3) 出力データのGPUからCPUへのコピーの3つであると述べました。今回は(1)はありません。25行目が(2)に相当します。GPUに関数を実行させるには、スレッドの数とブロック数を指定する必要があります。関数名の後に三重ブラケットを使って `<<<blocks, threadsPerBlock>>>` と書くと、ブロックの数が `blocks` になり、1ブロックあたりのスレッド数が `threadsPerBlock` になります。28行目は(3)のGPUからCPUへのデータコピーです。`cudaMemcpy` 関数を使うことで、CPU-GPU間でデータをコピーすることができます。

コンパイルとジョブ投入

- コンパイル(nvcc)により実行ファイルa.outが生成される

```
$ nvcc ex1-cuda.cu
$ qsub cuda.nqs
```

それでは、このソースコードをコンパイルして実行しましょう。ジョブスクリプトは「並列化なし」の時に作成した cuda.nqs を使用します。CUDA C++で書かれたコードをコンパイルするコマンドは nvcc です。コマンドオプションについては g++ と同様のものを持っていますが、それ以外にも GPU 向け専用のものもあります。後でいくつか出てきます。ジョブ投入コマンド qsub については並列化なしの時と同様です。ジョブを投入したら qstat, sstat コマンドで状態を確認しましょう。

プログラムの性能の考察

- 全スレッドがグローバルメモリにコアレスアクセスしている（連続する要素にアクセスしている）点については効率が良いと言える
- スレッド数が少ないため、ハードウェアマルチスレッディングは活用できない（もっと要素数が多いと活用できない）

補足

- ユニファイドメモリを使うとデバイス-ホスト間の手動でのデータ転送処理が不要になるが、今回は扱わない
- カーネル関数でブロック数、スレッド数を指定したが、それぞれ代わりに3次元座標で指定することもできる（例えば、1000と指定する代わりに(10, 10, 10) (10 × 10 × 10 = 1000である)と指定できる）。今回は省略する。

プログラムの性能について考察してみましょう。初日に GPU アーキテクチャの特性について話しましたが、アーキテクチャの特性を活かしたプログラムになっているでしょうか。7行目の配列 a へのアクセスに関して、この配列は cudaMalloc によって確保された d_a であるため、グローバルメモリ中に存在します。そして、CUDA においては、スレッドは 32 スレッドずつまとめられて1つのワープを形成しますが、スレッド 0 から順に 32 スレッドずつワープが作られます。よって、このプログラムで作成されるスレッド 0 から 15 は同じワープに属します。同一ワープに属するスレッドが $a[0]$ から $a[15]$ という連続領域にアクセスするので、このアクセスはコアレスアクセスであることがわかります。

また、今回はマルチプロセッサに割り当てられるワープは1つだけなので、ハードウェアマルチスレッディングは活用できません。もう少し要素数が多い場合でないと、活用することはできません。

最後にこのプログラムについて、いくつか補足しておきます。CUDA では CPU-GPU 間のデータ転送を自動化する仕組みとしてユニファイドメモリという機能を提供しています。しかし、この講習会では GPU の内部処理について理解したいため、この機能は使用しません。

ソースコード 25 行目のカーネル関数の呼び出しに関して、今回のプログラムではブロック数やブロックごとのスレッド数を1つの数値で指定しましたが、代わりに3次元座標形式で指定することもできます。例えば、1000と指定する代わりに(10,10,10)と指定することができます。こうするとブロックやスレッドは x 方向に 10 個、 y 方向に 10 個、 z 方向に 10 個の3次元配列として生成されるため、合計は 1000 個となります。この方式は空間のシミュレーションを行う場合などに便利です。

2.4.3 CUDA による GPU プログラム (その2:2ブロック使用)

CUDAプログラムによる計算 (複数ブロック使用)

- 本プログラムでは8スレッドからなるブロック2つを用いて計算を行う
 - 2つのブロックは同一マルチプロセッサ上で動作することもあれば、そうでないこともある (デバイスが自動でスケジューリングする)
 - スレッドIDはブロック毎に別々に割り当てられるということに注意する
- 計算時間を計測する処理も追加する

CUDA の文法を学ぶため、いま作成したプログラムを少し変更してみます。まず、ブロック数を2つにしてブロックあたりのスレッド数を8にします。これは勉強のためだけにやっていることに注意してください。こんなことをしてもまったく効率的になりません。また、各ブロックがどのマルチプロセッサで実行されるかについてはデバイスが自動でスケジューリングするため、ユーザは指定できないことにも注意しましょう。複数のブロックが同一のマルチプロセッサ上で動作することもあれば、そうでないこともあります。スレッドIDはブロックごとに個別に割り当てられますので、そのことにも注意しましょう。

ブロック数の変更と同時に、計算時間を計測する処理もプログラムに追加してみます。それでは、ソースコードを作成しましょう。レイアウトの関係上2つに分割しています。

ソースコード (ex1-cuda2.cu) (2分割のうち1つ目)

```
1  #include <stdio.h>
2
3  __global__ void WriteIndex(int* a, int n)
4  {
5      int i = blockDim.x * blockIdx.x + threadIdx.x;
6      if (i < n) {
7          a[i] = i;
8      }
9  }
10
11 int main()
12 {
13     // Init
14     int N = 16;
15     int* h_a = new int[N]; // on host (CPU)
16     int* d_a;              // on device (GPU)
17     cudaMalloc(&d_a, N * sizeof(int));
18
19     // Measure the elapsed time
20     cudaEvent_t start, stop;
21     float time_ms;
22     cudaEventCreate(&start);
23     cudaEventCreate(&stop);
24     cudaEventRecord(start, 0);
25
26     // Copy data from host to device
27     // Do nothing
```

ソースコード (ex1-cuda2.cu) (2分割のうち2つ目)

```
1
2 // Main
3 int threadsPerBlock = 8;
4 int blocks = (N + threadsPerBlock - 1) / threadsPerBlock; // round up
5 WriteIndex<<<blocks, threadsPerBlock>>>(d_a, N);
6
7 // Copy data from device to host
8 cudaMemcpy(h_a, d_a, N * sizeof(int), cudaMemcpyDeviceToHost);
9
10 // elapsed time
11 cudaEventRecord(stop, 0);
12 // Wait until cudaEventRecord is completed.
13 // (Note: cudaEventRecord is an asynchronous function)
14 cudaEventSynchronize(stop);
15
16 // Display results
17 for (int i = 0; i < N; i++) {
18     printf("%d\n", h_a[i]);
19 }
20 cudaEventElapsedTime(&time_ms, start, stop);
21 printf("exe time: %f ms\n", time_ms);
22
23 // Free memory
24 delete [] h_a;
25 cudaFree(d_a);
26 cudaEventDestroy(start);
27 cudaEventDestroy(stop);
28
29 return 0;
30 }
```

ソースコード解説

- 1つ目のファイルの5行目：全ブロックのスレッドに一意的な通し番号を振る処理。
 - `blockIdx.x` はブロックIDで、各ブロックに対し0から「ブロック数-1」までの番号が振られる。
 - `blockDim.x` はブロックあたりのスレッド数、`threadIdx.x` はブロック内でのスレッド番号。
- 2つ目のファイルの3,4行目：ブロック数とブロックあたりのスレッド数を決める処理。
 - 先に「ブロックあたりのスレッド数」を決めてから、その後必要なブロック数を計算している。
 - 「配列サイズ/1ブロックあたりのスレッド数」（今回は $16/8 = 2$ ）がブロック数になるが、この値が整数にならない時は切り上げる（そのため式が複雑になっている）。
- 計算時間計測処理
 - 1つ目のファイルの19-24行目：計測用のイベント作成とstartイベントの時刻記録
 - 2つ目のファイルの10-14行目：stopイベントの時刻記録
 - 2つ目のファイルの20,21行目：startイベントからstopイベントの経過時間算出。時間はミリ秒の単位で得られ、0.5マイクロ秒の精度を持つ（ベストプラクティスガイド）。
 - 2つ目のファイルの26,27行目：イベント削除

ソースファイルの解説をします。

まず、1つ目のファイルの5行目は全ブロックのスレッドに一意的な通し番号を振る処理です。スレッドIDはブロック毎に個別に割り当てられるIDですが、全スレッドの通し番号があったほうが便利なが多いため、これは定番の処理です。「`blockIdx.x`」はブロックIDで、各ブロックに対し0から「ブロック数-1」までの番号が振られます。また「`blockDim.x`」はブロックあたりのスレッド数、「`threadIdx.x`」はブロック内でのスレッド番号です。

次に、2つ目のファイルの3,4行目はブロック数とブロックあたりのスレッド数を決める処理です。ここでは「ブロックあたりのスレッド数」を先に決めてから、その後に必要なブロック数を決めています。ブロックあたりのスレッド数が8の時、必要なブロック数は明らかに $16/8 = 2$ なのですが、ここでは少し複雑な書き方をしています。基本的には「配列サイズ/1ブロックあたりのスレッド数」がブロック数になりますが、この値が整数にならない時は整数値に切り上げる必要があります。今回は切り上げ処理は不要ですが、定番処理なのであえて切り上げ処理を含んだ式を記載しています。

残りの変更箇所は計算時間計測に関連する処理です。計算時間を計測するには処理開始時刻と処理終了時刻の差分をとれば良いですが、CUDAでは時刻を取得するための `cudaEventRecord` というAPIと時刻を保持するための `cudaEvent_t` という構造体が用意されています。まず、1つ目のファイルの19-24行目では計測用のイベントを作成しstartイベントの時刻記録を行なっています。2つ目のファイルの10-14行目ではstopイベントの時刻記録を行なっています。注意事項として、時刻を記録するAPIである `cudaEventRecord` は非同期関数です。つまり、stopイベントへの時刻の書き込みが終わっていてもリターンします。そこで、`cudaEventSynchronize` というAPIを呼ぶことで書き込みが完了するまで待つ必要があります。2つ目のファイルの20,21行目はstartイベントからstopイベントの経過時間を算出して表示する処理です。時間はミリ秒の単位で得られます。CUDAのベストプラクティスガイドによれば、精度は0.5マイクロ秒とのこと。最後に2つ目のファイルの26,27行目はイベントの削除処理です。

コンパイルとジョブ投入

- コンパイル(nvcc)により実行ファイルa.outが生成される
- 最適化O3オプションとアーキテクチャを指定することで効率的なコードを生成する
 - アーキテクチャ指定についてはdeviceQueryの結果およびnvccのドキュメントを参照のこと
<https://docs.nvidia.com/cuda/cuda-compiler-driver-nvcc/index.html#options-for-steering-gpu-code-generation>

```
1 | $ nvcc -O3 -arch=sm_60 ex1-cuda2.cu
2 | $ qsub cuda.nqs
```

先ほど話した通りコンパイルには nvcc コマンドを使用しますが、今回はオプションをつけてみます。まず、O3 オプションをつけると最適化が行われます。これは g++ と同様です。次に arch オプションをつけると、アーキテクチャを考慮した最適化が行われます。今回しようとしている P100 では sm_60 を指定します。この値については、サンプルファイルの deviceQuery の実行結果や nvcc のドキュメントを参照してください。

ジョブの投げ方はこれまでと同様です。

プロファイリング (パフォーマンスの解析)

- CUDAプログラムのためのコマンドベースのプロファイリングツールとしてnvprofがある
- 詳細はNVIDIAのProfiler User's Guideを参照のこと
<https://docs.nvidia.com/cuda/profiler-users-guide/index.html#nvprof-overview>
 - オプションを指定することでより詳細な情報を取得できる
- GUIが使える時は総合開発環境であるNsightを使ってプロファイリングするのが便利
- nvprofを利用するためにはジョブスクリプト(prof.nqs)を以下のようにする

```
1 | #!/bin/bash
2 | #PBS -q LECTURE
3 | #PBS -l elapstim_req=00:00:10,cpunum_job=1,gpunum_job=1
4 | cd $PBS_O_WORKDIR
5 | nvprof ./a.out
```

この実習の最後に、プロファイリング、すなわちパフォーマンス等の詳細解析のやり方について話しておきたいと思います。計算時間が想定より遅い場合に、プロファイリングをすることにより、どの処理が遅いのかを解析することができます。コマンドベースのプロファイリングツールとして、nvprofがあります。詳細はNVIDIAのProfiler User's Guideをご覧ください。オプションを指定することで様々な情報を取得することができます。一部については今日の後半でやります。

なお、GUIが使える場合には、Nsight という総合開発環境を使ってプロファイリングするのが便利ですが、今日は端末上で実習しているため、nvprofを使います。

nvprofを使うには、nvprofコマンドの後に実行ファイル名を指定します。本演習環境では、ジョブスクリプトをこのように書き換えて、prof.nqs という名前で保存します。

- ジョブの投入

```
$ qsub prof.nqs
```

- 結果は標準エラー出力に出力される（以下は具体例）

```
-bash-4.2$ cat prof.nqs.e389728
==92462== NVPROF is profiling process 92462, command: ./a.out
==92462== Profiling application: ./a.out
==92462== Profiling result:
   Type  Time(%)   Time     Calls   Avg       Min       Max  Name
GPU activities: 60.75%  2.0800us      1  2.0800us  2.0800us  2.0800us  WriteIndex(int*, int)
                39.25%  1.3440us      1  1.3440us  1.3440us  1.3440us  [CUDA memcpy DtoH]
API calls:      99.75%  520.14ms      1  520.14ms  520.14ms  520.14ms  cudaMalloc
                0.13%   682.87us     94  7.2640us   114ns   282.80us  cuDeviceGetAttribute
                0.04%   214.33us      1  214.33us  214.33us  214.33us  cudaFree
                0.04%   209.93us      1  209.93us  209.93us  209.93us  cuDeviceTotalMem
                0.01%   56.766us      1  56.766us  56.766us  56.766us  cuDeviceGetName
                0.01%   45.203us      1  45.203us  45.203us  45.203us  cudaLaunch
                0.01%   30.483us      1  30.483us  30.483us  30.483us  cudaMemcpy
                0.00%   11.443us      2  5.7210us  2.2800us  9.1630us  cudaEventRecord
                0.00%   8.3070us      2  4.1530us    682ns   7.6250us  cudaEventCreate
                0.00%   6.3850us      1  6.3850us  6.3850us  6.3850us  cudaEventSynchronize
                0.00%   6.1000us      2  3.0500us    193ns   5.9070us  cudaSetupArgument
                0.00%   5.0610us      1  5.0610us  5.0610us  5.0610us  cudaEventElapsedTime
                0.00%   3.0570us      3  1.0190us    185ns   2.6390us  cuDeviceGetCount
                0.00%   1.9980us      2    999ns     517ns   1.4810us  cudaEventDestroy
                0.00%   1.2720us      1  1.2720us  1.2720us  1.2720us  cudaConfigureCall
                0.00%    595ns       2    297ns    146ns    449ns   cuDeviceGet
```

ジョブを投入すると、処理完了後、標準エラー出力にプロファイリング結果が出力されます。これはその一例です。

2.5 高速アルゴリズムの設計

高速アルゴリズムの設計

- カーネル関数（GPU起動関数）を繰り返し呼ぶことで処理を行う
- 各カーネル関数では、以下の2つの階層を考えると良い
 - 全体: スレッド生成と各スレッドへのタスク振り分け
 - Local: ブロック内でのスレッドの処理

8

ここからは高速アルゴリズムを設計する方法について解説します。まず、繰り返しになりますが、GPU プログラムは、CPU から GPU に対し繰り返しカーネル関数を呼び出すことで処理を行います。各カーネル関数を設計する際は、2つの階層を考慮すると良いと思います。一つはスレッドをいくつ生成し、タスク全体をどのようにスレッドに割り振るかという点で、もう一つはブロック内のスレッドがどのように効率よく処理を行うかという点です。ここでは一つ目の階層を全体、二つ目の階層をローカルと呼ぶことにします。これらは独立して考えられるものではないですが、ここでは一つずつ見ていくことにします。

全体のアルゴリズム設計

- 気をつけること
 - スレッド数を十分に大きくする
 - ⇒ スレッド数を多くしないと
 - ハードウェアマルチスレッディングを活かせない
 - ⇒ すべてのマルチプロセッサに最大数のスレッドを割り当てておくことを目指す
- オキュパンシ（占有率）
 - 割り当て可能なスレッド数に対する実際の割り当てスレッド数の比 ⇒ 100%にすることが望ましい

9

まず、アルゴリズムの全体設計において考慮すべき点について説明します。一つ目は、スレッド数を十分に大きくするという点です。スレッド数をコア数より大きくする必要があるということはすぐに分かると思いますが、初日に話した通り、GPU のハードウェアマルチスレッディングを活用するためにスレッド数をコア数よりも大幅に大きくする必要があります。マルチプロセッサには同時に割り当て可能なスレッド数の最大値が決まっていますが、各マルチプロセッサに最大数のスレッドを割り当てられるようにする必要があります。

初日の復習になりますが、GPU デバイスへの最大割り当て可能なスレッド数に対する実際の割り当てスレッド数をオキュパンシと呼びます。GPU デバイスの特性を活用するためには、オキュパンシが 100%になることが望ましいです。

オキュパンシの増やし方

1. ブロック数を大きくする

- 1つのマルチプロセッサに複数のブロックが割り当てられる

2. ブロック内のスレッド数を大きくする

- 仕様上の最大値は1024
- 32の倍数にすると効率的（∵1ワーブ32スレッド）

⇒ どちらでも良いができるだけ2を活用した方が効率的なアルゴリズムを設計しやすい

10

それではオキュパンシを大きくするための方法を見ていきます。単にスレッド数を多くすれば、オキュパンシが大きくなるという訳ではありません。それ以外にもメモリ使用量を考慮する必要があります。ですが、まずはスレッド数を大きくする方法についてみていきます。

総スレッド数を大きくするには、ブロック数を大きくする方法とブロックごとのスレッド数を大きくする方法が考えられます。まず一つ目ですが、一つのマルチプロセッサに複数のブロックを同時割り当て可能ですので、ブロック数を大きくすることはスレッド数を増やすのに有効です。ただし、1マルチプロセッサに同時割り当て可能なブロック数の上限がアーキテクチャごとに決まっていますので注意が必要です。今回使用する P100 では 32 個が上限となっています。他のアーキテクチャについては CUDA C Programming Guide の Compute Capability の章に記載されています。

スレッド数を増やす二つ目の方法はブロック内のスレッド数を大きくすることです。ただし、ブロック内のスレッド数についても上限が決まっており、今の所はすべてのアーキテクチャで 1024 です。なお、1 ワーブあたりのスレッド数は 32 なので、ブロック内のスレッド数を 32 の倍数にしておくことで効率の良いアルゴリズムを設計しやすいです。

これら 2 つのどちらの方法でもスレッド数を増やせるのですが、出来るだけ後者を活用した方がアルゴリズムを設計しやすいことが多いです。これを説明するためにブロックについて復習と補足説明をします。

ブロック

- ブロック内の全スレッドは必ず同一のマルチプロセッサで実行される
 - ブロック内のスレッドは共有メモリを用いてデータのやり取りをすることができる
- 複数のブロックが同時に一つのマルチプロセッサに割り当てられることもある
 - オキュパンシ計算時にはこのことも考慮する
 - ブロック間のデータ交換はできない

11

ブロック内の全スレッドは必ず同一のマルチプロセッサを用いて実行されます。そして、同じブロック内のスレッドとは共有メモリなどを介してデータのやり取りができます。一方で、異なるブロックのスレッドとはデータのやり取りができません。なので、出来るだけ同じブロックに多くのスレッドがいた方がアルゴリズムを設計しやすくなります。

また、複数のブロックが同時に1つのマルチプロセッサに割り当てられることもあります。ただしその場合であってもブロック間のデータやり取りはできません。ちなみにP100では1マルチプロセッサに同時割り当て可能な最大スレッド数は2048です。1ブロックあたりの最大スレッド数は1024なので、マルチプロセッサに最大数のスレッドを割り当てると、少なくとも2つのブロックがマルチプロセッサに割り当てられていることとなります。

ワーブ

- ブロック内のスレッドは32スレッドごとに1ワーブにまとめられる
- マルチプロセッサのコア数が64の時、1クロックにつき2ワーブをコアに割り当てる
⇒ $32 \times 2 = 64$
- ワーブ内のスレッドは基本的に、常に同じ命令を実行する (Voltaは別)
- ブロック内の各ワーブは異なるタイミングで命令を実行する

12

ワーブについても復習と補足説明をしておきましょう。初日に話した通り、ブロック内のスレッドは32スレッドごとにワーブとしてまとめられます。そしてワーブ内の全スレッドは基本的には常に同じ命令を実行するのでした。2017年に発表されたVoltaアーキテクチャでは少し異なりますがここでは置いておきます。P100ではマルチプロセッサあたりのコア数は64です。この時、各クロックにおいて2つのワーブがマルチプロセッサに割り当てられます。どのワーブが割り当てられるかはハードウェアにより計算されます。なので、ブロック内の各ワーブは異なるタイミングで命令列を実行していくこととなります。このタイミングを合わせる命令もありますが、それについては後で説明します。

オキュパンシを増やす際の障壁

- マルチプロセッサ内のスレッドは共有メモリとレジスタを共有する
 - ⇒ スレッド数が多いと1スレッドあたりのメモリ割り当て量が小さくなる
 - ⇒ メモリ使用量の大きい手法は使えない
- 一般に1スレッドに割り当てられるデータ量が大きいほど、アルゴリズムを設計しやすい(次の例参照)

13

それでは本線に戻って、オキュパンシの増やし方について説明します。単にスレッド数を多くすればオキュパンシが大きくなるという訳ではありません。オキュパンシを大きくする際の障壁としてメモリ使用量があります。それについて説明します。

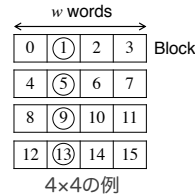
マルチプロセッサに複数のスレッドを割り当てるとき、それらのスレッドはハードウェア資源を共有します。特に共有メモリについては、決められた容量をみんなで使うことになります。CUDAではブロックごとに共有メモリを確保するのですが、マルチプロセッサにたくさんのブロックを割り当てようとする、1ブロックあたりで使えるメモリ量は少なくなります。逆にブロックのメモリ使用量が多い場合、1つのマルチプロセッサに割り当て可能なブロック数は小さくなります。

GPUではレジスタもスレッドごとに個別に割り当てられます。もし複数のスレッドが同じレジスタを利用する場合、スレッドを切り替えるたびにレジスタ内容の切り替え処理が必要になり時間がかかりますが、GPUではそのような処理が不要です。これはハードウェアマルチスレッディングでクロックごとに実行スレッドを高速切り替える際に有利です。しかし、このような仕様によりマルチプロセッサ内のスレッド数を増やすと、1スレッドあたりが使用できるレジスタ数は少なくなります。逆にスレッドのレジスタ使用量が多い場合、1つのマルチプロセッサに割り当て可能なスレッド数は小さくなります。

ということで、オキュパンシのことを考えると共有メモリやレジスタの使用量は十分に小さくする必要があります。しかし一般的にはメモリ使用量が多い方が効率的なアルゴリズムを設計しやすくなるので、これらはトレードオフの関係になっています。このトレードオフについて一つ例を挙げてみます。

よくあるトレードオフ (オキュパンシ vs メモリアクセス)

- 32 × 32 行列がグローバルメモリに格納されており、各列にアクセスしたい
- 32個の各スレッドが1要素にアクセスするとコアレスアクセスではなくなってしまう
- そこで一旦全要素を共有メモリにコピーすることにする
⇒ コアレスアクセスが可能になる



14

よく遭遇するトレードオフのシチュエーションとして、行列計算があります。行列計算では、行列の各列にアクセスしたいことがよくあります。行列積を計算する場合がその典型例です。グローバルメモリに格納された行列の各列にアクセスすることを考えます。この例では4コアで4×4行列の各列にアクセスすることを考えます。つまり最初のアクセスで4つのコアは0,4,8,12にアクセスし、二回目のアクセスで1,5,9,13にアクセスします。以下同様です。この場合、各コアはグローバルメモリの不連続なアドレスにアクセスすることになりコアレスアクセスではなくなります。

この時、メモリ使用量を大きくすることで、コアレスアクセスを可能にする方法があります。それは、まず最初に行列の全データをグローバルメモリから共有メモリにコピーするという方法です。その際、1回目で0,1,2,3をコピーし、2回目で4,5,6,7をコピーするという感じで行ごとにデータをコピーすることで全てのデータコピーをコアレスアクセスにより行うことができます。1回データを共有メモリにコピーしてしまえば、もはやコアレスアクセスについて気にする必要はありません。ここでは、行列の各列にアクセスしたいわけですので、コピーした全てのデータが使われることになり、無駄がありません。また、共有メモリのバンクコンフリクトについて気になる人がいるかもしれませんが、うまく回避する方法があります。これは後でやります。

ということで、この方法はすごくいいような気がするのですが、問題はメモリ使用量が大きいうことです。これによりオキュパンシが小さくなってしまいます。

よくあるトレードオフ (オキュパンシ vs メモリアクセス) (続き)

- メモリ使用量の解析
 - 各要素が4バイト変数とする32×32行列のサイズは $32 \times 32 \times 4 = 4096B = 4KB$
 - P100の共有メモリ：64KB
 - よって最大でも $64/4=16$ ワープしか配置できない
 - P100の最大ワープ数/SM：64
 - オキュパンシは $25\%(=16/64)$
- 扱う行列サイズを16×16にするとオキュパンシは100%になるが、メモリアクセスは多くなる

15

今回の実習で使用する P100 を例として、具体的な値を計算してみます。グローバルメモリ内に4バイト変数を要素とする 32×32 行列があり、ワープ内のスレッドが各列にアクセスするものとします。この時、行列の全要素を共有メモリにコピーすることを考えると、そのサイズは $32 \times 32 \times 4B = 4KB$ になります。一方で P100 マルチプロセッサの共有メモリ容量は64KBです。1つのワープが一つの行列を扱うとすると、1マルチプロセッサに配置可能なワープ数は $64/4 = 16$ となります。P100 では1マルチプロセッサあたり64ワープを配置可能という仕様でした。最大64ワープに対して、16ワープしか配置できないので、オキュパンシは $16/64$ で25%となります。

ちなみにこの全データコピーのアイデアを活かしたまま、オキュパンシを大きくする方法があります。それは扱う行列サイズを 16×16 にすることです。すると行列サイズは1KBになるので、オキュパンシは100%となります。しかし、このようにするとメモリアクセス回数は大きくなります。

以上のように、行列の列アクセスではオキュパンシとメモリアクセス回数のトレードオフが発生します。

全体のアルゴリズム設計まとめ

- オキュパンシ100%を目指す
 - 作成するスレッド数 (=ブロック数×ブロックあたりのスレッド数) をデバイスで同時割り当て可能なスレッド数以上にする
 - 例：P100はSM数：56, SMあたりの最大ワープ数：64
⇒ 作成ワープ数を $3584 (=56 \times 64)$ 以上にする
 - ただし、オキュパンシはローカルで使用するアルゴリズムの共有メモリ使用量 (とレジスタ使用量) に依存する

16

全体のアルゴリズム設計についてまとめておきます。目指すのはオキュパンシを100%にすることです。そのためには作成するスレッド数を大きくすることと、使用するメモリ、レジスタ量が小さくなるアルゴリズムを設計することが必要です。

Localのアルゴリズム設計

- ブロック内のスレッドの処理の設計法を述べる
- 32スレッドごとに1ワーブになっていることに注意する
- ブロック内のワーブは同一の命令列を実行するがタイミングはバラバラである
- 全ワーブを同期する命令も用意されている
 - `_syncthreads()`
 - 遅いワーブが追いつくのを待つ

17

続いてローカルのアルゴリズム設計について説明します。つまりブロックの中の各スレッドに対するアルゴリズムの設計ということです。

まず繰り返しになりますが、ブロックの中のスレッドは32ワーブずつ1つのワーブを構成していることに注意します。また、ブロック内のすべてのワーブが同じ命令列を実行しますが、ワーブにより実行されるタイミングはバラバラです。あるワーブは非常に速く計算が進み、別のワーブは非常に遅く計算が進むかもしれません。すべてのワーブを強制的に同じタイミングにするコマンドとして`_syncthreads`関数があります。これは、この関数が書かれた地点で一番遅いワーブが追いつくのを待つという命令です。共有メモリを用いてワーブ間でデータのやり取りをしたい場合などに便利です。

Localのアルゴリズム設計

- ワーブ内のスレッドはSIMD的に動作することに注意する（Voltaは別）
 - グローバルメモリへのコアレスアクセス
 - 共有メモリへのバンクコンフリクトを避ける
 - 条件文では、非該当のコアは待ち状態になる
- オキュパンシを大きくするために共有メモリ使用量を極力小さくする

18

次に、ワーブ内のスレッドはSIMD的に動作することに注意が必要です。Voltaアーキテクチャはちょっと異なりますがここでは置いておきましょう。このことを活用するとグローバルメモリにコアレスアクセスできたり、共有メモリのバンクコンフリクトを避けたりすることができます。また条件文において、条件に非該当のスレッドは待ち状態になります。そのため条件文を大量に使うとコアの稼働率を下げることになるので注意が必要です。

また、全体アルゴリズムでも述べましたが、オキュパンシを大きくするために共有メモリ使用量とレジスタの使用量を極力小さくする必要があります。

アルゴリズム設計の具体例： リダクション

- リダクションとは？
 - 総和の一般化
- $\sum_{i=0}^3 a_i = a_0 + a_1 + a_2 + a_3$
- $a_i \uparrow a_j \stackrel{\text{def}}{=} \max(a_i, a_j)$ と定義すると
 $\max(a_0, a_1, a_2, a_3) = a_0 \uparrow a_1 \uparrow a_2 \uparrow a_3$
- ここで、 $a_0 \oplus a_1 \oplus a_2 \oplus a_3$ を考えると、
 \oplus を+にすれば総和、 \uparrow にすればmaxになる

19

アルゴリズム設計をする時に全般的に考慮しないといけないことは、これまで話した通りです。個々の高速化テクニックについては後半で簡単に紹介しようと思います。ここからは、これまで話してきた考慮すべき項目についてどのようにアルゴリズム設計に適用するのかということについて説明したいと思います。実際にアルゴリズムを設計してみるのが良いと思いますので、ここでは具体例としてリダクションと呼ばれる操作のためのアルゴリズムを設計してみます。

リダクションというのは配列の総和を求める計算を一般化したものです。総和というのは配列の全ての要素に対し、加算という二項演算を繰り返し適用し、1つの値を返す処理ですが、リダクションとは二項演算子について加算だけでなく任意の演算子を取れるようにしたものです。やや細かい話をしておく、演算子は結合則を満たすものとして、また、本日は可換である場合を扱います。

色々細かい話をしていますが、本日の内容に関しては、配列の総和を高速計算するアルゴリズムを設計する思っただけで問題ありません。ここで言いたいのは、そのアルゴリズムが他の様々な計算にもそのまま使えますよということだけです。他の計算の例として、配列の最大値計算が挙げられます。今、2要素のうち大きい方を返す演算を上矢印で書くことにすると、配列の要素の最大値は、配列に関してこの上矢印でリダクションをすることにより得られます。

アルゴリズム設計の具体例： 総和計算

- 問題
 - n要素からなる配列の総和を求めよ
- これから2つのやり方を紹介する
 - Tree-basedアルゴリズム
 - Cascadingアルゴリズム
- 今回は配列の総和を考えるが、
max, minなども同じ方法で計算できる

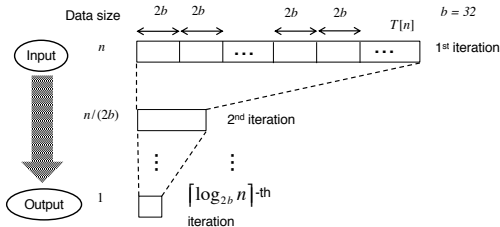
Optimizing Parallel Reduction in CUDA, Mark Harris
<https://developer.download.nvidia.com/assets/cuda/files/reduction.pdf>

20

ということで配列が与えられた際、その総和を求める問題を考えます。配列はGPUのグローバルメモリに格納されているものとして、この問題に対して、二つのアルゴリズムが知られています。Tree-basedアルゴリズムとCascadingアルゴリズムです。後者の方が高速なのですが、これが良いアルゴリズムだということを理解するのが今日の目標です。

Tree-basedアルゴリズム

- 全体のアルゴリズム設計
 - 各ワープに64要素を割り当て1要素にリダクション
 - 要素数が1になるまで繰り返す



21

まず、Tree-based アルゴリズムを見てみましょう。まず全体のアルゴリズム設計についてです。このアルゴリズムでは、まず入力データを 64 要素からなる小ブロックに分割します。そして、1つの小ブロックを一つの CUDA ブロックに割り当てることにします。各 CUDA ブロックでは 64 要素の和を計算して、グローバルメモリの所定の場所へ出力します。入力データが n 個だとすると、この処理の終了後、データサイズはだいたい $n/64$ 個ぐらいになります。この処理をデータサイズが 1 になるまで繰り返すと最終的な答えを得ることができます。

Tree-basedアルゴリズム

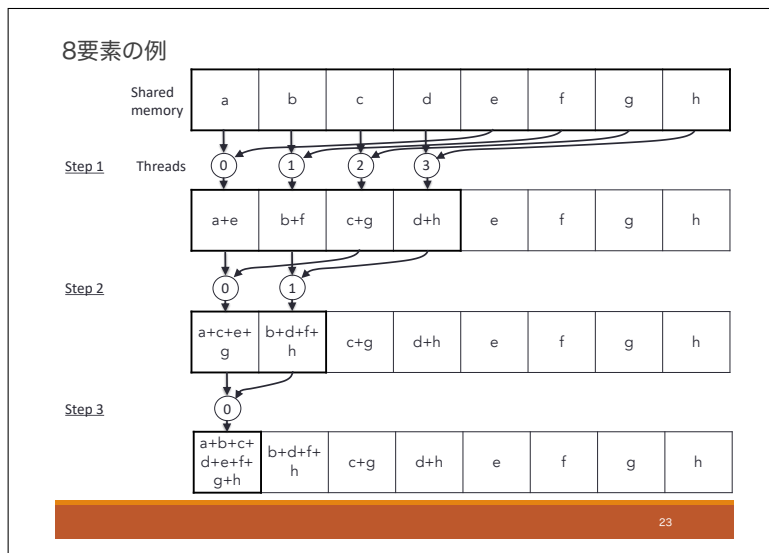
- ローカルのアルゴリズム設計
 - 1ステップごとに要素を半分にする
 - 要素数が1になるまで繰り返す

22

次にローカルのアルゴリズム設計について説明します。ここでは各ブロックは 64 要素の和を計算します。まずブロックの構成についてですが、各ブロックが 64 個のスレッドで構成されるようにします。32 スレッドごとに 1 つのワープとなることに注意してください。

まず、ブロック内の 64 スレッドは、入力データをグローバルメモリから共有メモリにコピーします。次に 32 スレッドの各スレッドが 2 要素ごとの和を求めて共有メモリに書き込むと、32 個の出力が得られます。次に 16 スレッドが 2 要素の和を求めると 16 個の和が得られ、これを繰り返すと 64 要素全ての和が求まります。

なお、ブロックあたりのスレッド数をもっと大きくすると、もう少し効率的なアルゴリズムができるのですが、ここでは深入りしないことにします。ブロックあたりのスレッド数を増やすと小ブロックのサイズを 64 より大きくできて、全体アルゴリズムにおける繰り返し処理の回数が少なくなります。



これは8要素の和を求める例です。四角が共有メモリで、丸がスレッドです。8スレッドが8要素を共有メモリにコピーした後に、共有メモリの内容がどのように書き換わるかを時系列で示しています。各ステップでデータサイズが半分になる様子を確認してください。この例では、1ステップ後に有効データ数が4になり、2ステップ後に2、3ステップ後に1となっています。ローカルアルゴリズムにおける和計算のステップ数は、要素数の対数をとったものになります。データサイズが64要素の場合は $\log_2 64 = 6$ ステップです。

Tree-basedアルゴリズム： ローカルのアルゴリズム

- 共有メモリ使用量の解析
 - 1ブロックあたり64変数分の共有メモリ使用。
1変数4バイトとすると256Byte使用
⇒ 32ブロック (64ワーブ) でも8KByte
⇒ オキュパンシに影響を与えない
- 計算量
 - スレッド数をbとすると $O(\log b)$
 - 待ち状態のコア (スレッド) が多いので非効率

24

ローカルアルゴリズムのメモリ使用量や計算量を解析してみましょう。

まず、共有メモリの使用量ですが、1ブロックあたり64変数分の共有メモリを使用します。もし1変数のサイズが4バイトなら、ブロック全体では、 $4 \times 64 = 256$ バイトです。1ブロックは2ワーブで構成されており、1マルチプロセッサあたりの最大ワーブ割り当て数は64です。よって、最大で32ブロックをマルチプロセッサに割り当て可能です。もし32ブロックを割り当てたとしてもメモリ使用量は $32 \times 256\text{Byte} = 8\text{KB}$ です。P100マルチプロセッサの共有メモリサイズは64KBですので、メモリサイズが原因でオキュパンシが低下することはなさそうです。

次にローカルアルゴリズムの計算量ですが、ブロックのスレッド数をbとすると、 $O(\log b)$ となります。このアルゴリズムの欠点として、コアの稼働率が低くなるということが挙げられます。ワーブ内の32スレッドのそれぞれにコアが割り当てられますが、和計算の2ステップ以降は16個以下のスレッドしか使われません。この時、未使用スレッドに割り当てられたコアは待ち状態となります。Tree-basedアルゴリズムでは、後半のステップで待ち状態のコアが多くなります。

Cascadingアルゴリズム

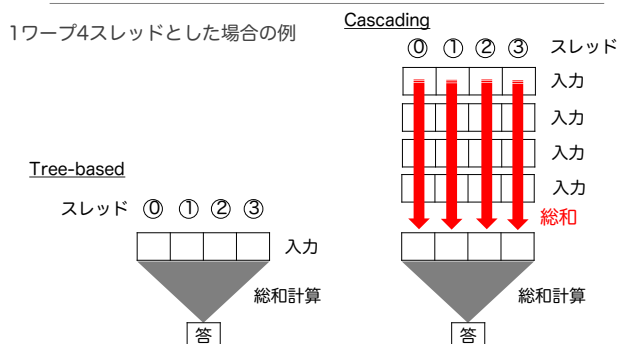
- Tree-basedアルゴリズムでは、
ローカルリダクションが非効率だった
⇒ 待ち状態のコアが多くもったいない
- 各ワープへの割り当てデータ量を増やすこと
で効率的にすることができる
⇒ Cascadingアルゴリズム

25

それでは次に Cascading アルゴリズムを見てみましょう。前のスライドで見たように、Tree-based アルゴリズムでは、ローカルリダクションが非効率的でした。待ち状態のコアが多いため、それらをもっと活用できれば、高速化できそうです。

実は各ワープへの割り当てデータ量を増やすことで、効率的にすることができます。これを Cascading アルゴリズムと呼びます。

Cascadingアルゴリズム： ローカルのアルゴリズム

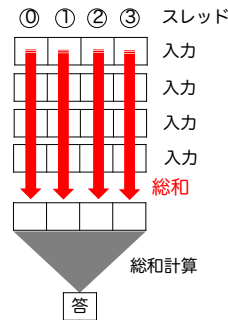


26

Cascading アルゴリズムの概要を Tree-based アルゴリズムと対比しながら説明します。ここでは、入力データは各ワープに分割して割り当てられるものとし、各ワープが割り当てデータの総和を求めるためのアルゴリズムを説明します。図を見やすくするため、図では1ワープ4スレッドとして描かれています。本当は1ワープは32スレッドです。Tree-based アルゴリズムでは、スライドの左側のように読み込んだ入力データをスレッドが協調しながら計算します。一方、Cascading アルゴリズムでは、その処理の前に各スレッドが個別に総和を計算するフェーズが入ります。ワープ内のスレッドはコアレスアクセスによりグローバルメモリから入力を複数回読み込み、各スレッドは決まった位置 (index) に対する総和を求めます。その後その結果を共有メモリに書き込み、後は Tree-based アルゴリズムと同様の総和計算を行います。

Cascadingアルゴリズム： ローカルアルゴリズム

- 共有メモリ使用量の解析
 - 共有メモリは各列の結果の格納のみに使用
⇒ オキュパンシに影響を与えない
- 計算量
 - 読み込む入力ブロック数が多ければ ($\log b$ 以上なら b : スレッド数)
最後の Tree-based 部分の計算時間は無視できる
 - コアがほぼフル稼働する



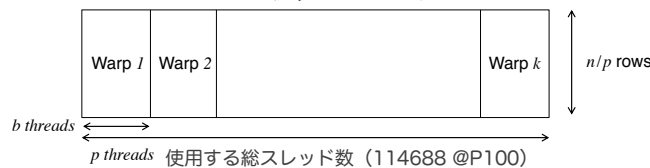
27

ローカルアルゴリズムについて、メモリ使用量と計算量の解析をしてみましょう。まず共有メモリの使用量ですが、各入力ブロックについては、スレッドが有するレジスタに格納すればよく共有メモリを使う必要はありませんので、メモリ使用量は Tree-based アルゴリズムと同様になります。すなわち、メモリ使用量が原因でオキュパンシが低下することはありません。

次に計算量ですが、読み込む入力ブロック数が十分に多ければ、最後の Tree-based 部分の計算量は漸近解析としては無視できます。各ワープに十分に多い n_w 個の入力データが割り当てられる時、1 ワープのスレッド数を b とすると計算量は $O(n_w/b)$ となります。このアルゴリズムはコアがほぼフル稼働するので、効率が良いと言えます。

Cascadingアルゴリズム： 全体のアルゴリズム

- 先に使用するワープ数 k を決める
 - デバイスへの最大割り当て可能数にする (3584 @ P100)
- 各ワープに均等にデータを割り振る
- 各ワープの処理後、残った k 個をリダクション



28

次に全体のアルゴリズムについて、説明します。Cascading アルゴリズムでは、各ワープにできるだけたくさん入力データを割り当てたほうが効率が良くなります。しかし、それを大きくしようとすると生成するワープの数が小さくなり、オキュパンシが低下してしまいます。そこで、オキュパンシを 100% に保てるためのワープ数を先に計算して、各ワープに均等に入力データを割り振ることにします。P100 では 56 個のマルチプロセッサが搭載されており、各マルチプロセッサでは最大 64 個のワープを同時割り当て可能です。よって、ワープを $56 \times 64 = 3584$ 個生成すればオキュパンシを 100% にできることがわかります。これはスレッド数に直すと $3584 \times 32 = 114688$ 個です。

ということで、入力データを 3584 個に分割し、各ワープがそれらを並列に総和計算します。最後に 3584 個のデータを一つのワープを用いて総和計算します。ここもさらに高速化できますが、そもそもそれほど時間のかかる処理ではないため、ここでは扱いません。

アルゴリズム設計まとめ

- アルゴリズム設計の際は以下に注意する
 - オキュパンシ100%を目指す
(P100では総ワーブ数を3584以上, ワーブあたりの共有メモリ使用量を1KB以下にする, レジスタの使用量も要注意だが今回は省略)
 - グローバルメモリへのコアレスアクセス
 - 共有メモリのバンクコンフリクト回避 (今回省略)
 - 待ち状態のコアを少なくする

29

最後に繰り返しになりますが, アルゴリズム設計時の注意点をまとめておきます. まず, オキュパンシを100%にすることを目指します. 具体的に言うと, P100では総ワーブ数を3584個以上にし, ワーブあたりの共有メモリ使用量を1KB以下にすることが必要です. レジスタ使用量が多すぎる場合もオキュパンシが低下しますが, ここでは省略します.

そしてメモリアクセスの仕方として, グローバルメモリへのコアレスアクセスと共有メモリへのバンクコンフリクト回避が有効です. 後者についてはこの例では出てきませんでしたが, 本日後半の高速化トピックの中でバンクコンフリクト回避のためのテクニックが少しだけ出てきます.

あとは, 待ち状態のコアを少なくすることが重要です. マルチプロセッサ内のコアはSIMD命令を実行しているという感覚を持ちながらアルゴリズムを設計すると良いと思います.

アルゴリズム設計まとめ

- 計算量の解析
 - KoikeらのGPU計算モデル (AGPUモデル) では, コア数 p とマルチプロセッサ内のコア数 b を用いて, 以下の3つを計算し, アルゴリズムを評価する
 - マルチプロセッサごとのSIMD命令回数
 - グローバルメモリへのアクセス回数
(コアレスアクセスでない時はその分多くする)
 - 共有メモリ使用量

Koike and Sadakane, A Novel Computational Model for GPUs with Applications to Efficient Algorithms
<http://www.iijnc.org/index.php/iijnc/article/view/96>

30

最後にやや宣伝になりますが, 発表者らはGPUアルゴリズムを評価するための計算モデルを提案しています. ここでは, コア数 p とマルチプロセッサ内のコア数 b を用いて, マルチプロセッサごとのSIMD命令回数, グローバルメモリへのアクセス回数, 共有メモリ使用量を算出し, アルゴリズムを評価します. このモデルでは1マルチプロセッサでは1ワーブのみが実行されると仮定しますが, オキュパンシについては共有メモリ使用量から別途評価します.

2.6 実習 2：効率良く総和計算を行うプログラム


実習2. 効率良く総和計算を行うプログラム

- 問題
 - 2^{25} 要素からなる配列にはランダムに0または1が書き込まれている。この時、この配列の総和を求めなさい

31

これまで解説してきた総和のアルゴリズムを実際に動かしてみましょう。スライドに書いてある通り、2の25乗個の要素からなる配列の総和を求めましょう。各要素は0または1です。これは桁あふれを考慮しなくいいようにするための仕様です。

2.6.1 Cascading アルゴリズムによる総和計算

 **Cascadingアルゴリズムによる総和計算**

- 各ブロックの処理の流れ
 1. ワープ内のスレッドが個別に担当データの総和を求める
 2. ワープ内の各スレッドが求めた結果に対する総和計算を行う
 3. ブロック内の各ワーブが求めた結果に対する総和計算を行う

各ブロックの処理の流れは以下の3つです。まず、ワーブ内の各スレッドが個別に担当データの総和を計算します。これはこれまで説明したように、ワーブ内のスレッドがコアレスアクセスにより32個のデータを取得したのち個々のスレッドが1要素分の和の計算を行うのでした。この処理がすべて終わったら次にワーブ内の各スレッドが持っている結果について、総和計算を行います。これはこれまで話してきた通りです。最後に、これはこれまでの説明で話してないですが、ちょっとした工夫として、ブロック内の各ワーブの結果について、さらに総和計算を行います。これで、あるブロックで担当したすべてのデータの総和が計算できます。

それではソースコードを見ていきましょう。長いので4分割して表示します。最初の2つがカーネル関数、つまりGPUで動作する関数で、残りの2つがCPU側の処理です。

ソースコード (ex2.cu) (4分割のうち1つ目)

```
1 #include <stdio.h>
2 #include <math.h>
3 // #include <thrust/reduce.h>
4 // #include <thrust/device_ptr.h>
5
6 // Device parameter
7 const int NUM_SM = 56;
8 const int WARPS_PER_SM = 64;
9 const int WARPS_PER_BLOCK = 32;
10 const int NUM_BLOCKS = NUM_SM * WARPS_PER_SM / WARPS_PER_BLOCK;
11 const int THREADS_PER_WARP = 32;
12 const int THREADS_PER_BLOCK = WARPS_PER_BLOCK * THREADS_PER_WARP;
13
14 __global__ void ReductionCascading(
15     int* inArray, unsigned int numElements, int* halfwayResult)
16 {
17     // thread id in each block
18     const unsigned int tid = threadIdx.x;
19     // warp id in each block
20     const unsigned int wid = tid / THREADS_PER_WARP;
21     // lane id = thread id in each warp
22     const unsigned int lid = tid % THREADS_PER_WARP;
23     // sequential thread id in the kernel
24     const unsigned int sid = blockIdx.x * blockDim.x + threadIdx.x;
25     const unsigned int numTotalThreads = gridDim.x * blockDim.x;
26
27     int item1 = 0;
28     __shared__ int x[WARPS_PER_BLOCK * THREADS_PER_WARP];
29     __shared__ int y[WARPS_PER_BLOCK];
30
31     const int numIter = 1 + (numElements - 1) / numTotalThreads; //round up
32     // Main
33     for (int i = 0; i < numIter - 1; i++) {
34         item1 += inArray[numTotalThreads * i + sid];
35     }
36     // Last
37     int idx = numTotalThreads * (numIter - 1) + sid;
38     if (idx < numElements) {
39         item1 += inArray[idx];
40     }
41     x[tid] = item1;
42     __syncwarp();
```

ソースコード (ex2.cu) (4分割のうち2つ目)

```
1
2 // Reduction for a warp
3 x[tid] += x[tid + 16];
4 __syncwarp();
5 x[tid] += x[tid + 8];
6 __syncwarp();
7 x[tid] += x[tid + 4];
8 __syncwarp();
9 x[tid] += x[tid + 2];
10 __syncwarp();
11 x[tid] += x[tid + 1];
12 if( 0 == lid ){
13     y[wid] = x[tid];
14 }
15 __syncthreads();
16
17 // Reduction for a block
18 if (tid < 32) {
19     y[tid] += y[tid + 16];
20     __syncwarp();
21     y[tid] += y[tid + 8];
22     __syncwarp();
23     y[tid] += y[tid + 4];
24     __syncwarp();
25     y[tid] += y[tid + 2];
26     __syncwarp();
27     y[tid] += y[tid + 1];
28 }
29 if( 0 == tid ){
30     halfwayResult[blockIdx.x] = y[tid];
31 }
32
33 return;
34 }
35
```

ソースコード (ex2.cu) (4分割のうち3つ目)

```
1 int main()
2 {
3     int logNumElements = 25; // input_size = 2^25
4     //int logNumElements = 10; // input_size = 2^25
5     int numElements = 1 << logNumElements;
6     int* h_input = new int[numElements]; // on host (CPU)
7     int* h_halfwayResult = new int[NUM_BLOCKS];
8     int h_output = 0;
9     int* d_input; // on device (GPU)
10    int* d_halfwayResult;
11    cudaMalloc(&d_input, numElements * sizeof(int));
12    cudaMalloc(&d_halfwayResult, NUM_BLOCKS * sizeof(int));
13    srand(0);
14    for( unsigned int i = 0; i < numElements; ++i) {
15        h_input[i] = (int)(rand() % 2);
16    }
17    int answer = 0;
18    for (int i = 0; i < numElements; i++) {
19        answer += h_input[i];
20    }
21
22    // Measure the elapsed time
23    cudaEvent_t start, stop;
24    float time_ms;
25    cudaEventCreate(&start);
26    cudaEventCreate(&stop);
27
28    // Copy data from host to device
29    cudaMemcpy(d_input, h_input, numElements * sizeof(int),
30              cudaMemcpyHostToDevice);
31
32    // Main
33    printf("numElements = %d (2^%d)\n", numElements, logNumElements);
34    cudaEventRecord(start, 0);
35    ReductionCascading<<<NUM_BLOCKS, THREADS_PER_BLOCK>>>(
36        d_input, numElements, d_halfwayResult);
37
38    // Copy data from device to host
39    cudaMemcpy(h_halfwayResult, d_halfwayResult, NUM_BLOCKS * sizeof(int),
40              cudaMemcpyDeviceToHost);
```

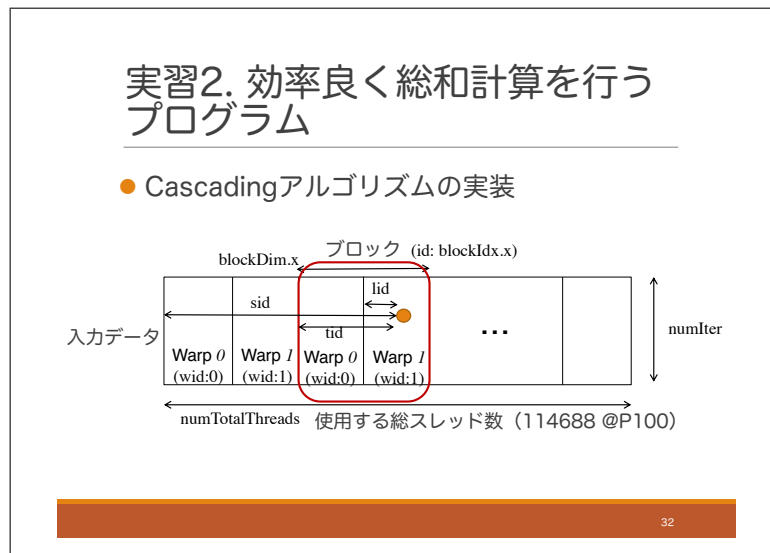
ソースコード (ex2.cu) (4分割のうち4つ目)

```
1
2     h_output = 0;
3     for (int i = 0; i < NUM_BLOCKS; i++) {
4         h_output += h_halfwayResult[i];
5     }
6
7     // elapsed time
8     cudaEventRecord(stop, 0);
9     // Wait until cudaEventRecord is completed.
10    // (Note: cudaEventRecord is an asynchronous function)
11    cudaEventSynchronize(stop);
12
13    // Display results
14    printf("Our result = %d, answer = %d\n", h_output, answer);
15    cudaEventElapsedTime(&time_ms, start, stop);
16    printf(" exe time: %f ms\n", time_ms);
17
18    // Thrust Library
19    // int h_outputThrust;
20    // thrust::device_ptr<int> d_input_ptr = thrust::device_pointer_cast(d_input);
21    // cudaEventRecord(start, 0);
22    // h_outputThrust = thrust::reduce(d_input_ptr, d_input_ptr + numElements);
23    // cudaEventRecord(stop, 0);
24    // // Wait until cudaEventRecord is completed.
25    // // (Note: cudaEventRecord is an asynchronous function)
26    // cudaEventSynchronize(stop);
27    // printf("Thrust result = %d, answer = %d\n", h_outputThrust, answer);
28    // cudaEventElapsedTime(&time_ms, start, stop);
29    // printf(" exe time: %f ms\n", time_ms);
30
31    // Free memory
32    delete [] h_input;
33    delete [] h_halfwayResult;
34    cudaFree(d_input);
35    cudaFree(d_halfwayResult);
36    cudaEventDestroy(start);
37    cudaEventDestroy(stop);
38
39    return 0;
40 }
41
```

ソースコード解説 (4分割のうち1つ目)

- 6-10行目は各種パラメータの設定. 内容は変数名の通り
- 14行目から2つ目のソースコードの最後まではカーネル関数
- 17行目から25行目は各種ID等の定義. 詳細は後述する
- 27行目の変数 `item1` に各スレッドの総和の途中経過を保存する. このように書くと通常はレジスタに変数が確保される
- 28,29行目: 共有メモリ上に変数を定義する. `__shared__` とつけると共有メモリ上にメモリが確保される. `x` はワープ内の各スレッドの総和計算結果を格納する配列. `y` は各ワープの総和計算結果を格納する配列.
- 31-42行目は各スレッドが担当データの総和を計算する処理. 33-35行目がメイン処理. 各スレッドに割り当てる要素数について, 入力データ数によっては一つ多く割り当てられるスレッドとそうでないスレッドが出てくるので, イテレーションの最後の1回にはその判定処理を入れる.
- 42行目の `__syncwarp()` は, ワープ内のスレッドで同期をとる処理. P100では不要だが, Voltaアーキテクチャではスレッド内のスレッドが同期処理されないことがあるので, 入れてある

1つ目のソースコードの説明は上にまとめた通りです. 変数の意味については以下のスライドにまとめます.



この図の見方ですが, 入力データを2次元に配置し, 各スレッドが1列の総和を計算します. その際の自スレッドの担当列のindexがsidです. この図では各ブロックは2つのワープを持つものとして書かれていますが, 実際には32個のワープを持っています. あるスレッドが所属するブロックの通し番号は, `blockIdx.x` で取得でき, ブロック内でのスレッドの通し番号は `threadIdx.x` で取得できます. また, 各ブロックに割り当てられているスレッドの数は `blockDim.x` で取得できます.

スレッドに対して, 様々なIDのつけ方が考えられるので, 先に計算しておくことにします. ワープ内でのスレッドID(`lid`: Lane ID), ブロック内でのスレッドID(`tid`: Thread ID), カーネル関数全体でのスレッドID(`sid`: Sequential Thread ID)などはよく使われるので先に計算しておくこととアルゴリズムを書く際に便利です. ブロック内でのワープのID (`wid`: Warp ID) もよく使います. ちなみにこれらの変数名は必ずしも標準的とは言えませんが, Lane ID という名称については, 比較的広く使われている印象です.

ソースコード解説 (4分割のうち2つ目)

- 2-15行目はワーブ内の各スレッドの計算結果に対して総和をとる処理。 `__syncwarp()` はVoltaアーキテクチャ用に入れてある。3-11行目では処理を行う必要のないスレッドも処理を行うことになるが、余計な条件文を入れたくないので気にしないことにする。
- 15行目の `__syncthread()` はブロック内の全スレッドの同期をとる処理。すべてのスレッドが共有メモリに結果を書き込み終わった後に次の処理に進む必要があるため、この処理が必要。
- 17-28行目は各ワーブの計算結果に対して総和をとる処理。18行目で条件を指定することにより1ワーブのみを使って計算を行なっている。
- 29-31行目はブロックでの計算結果をグローバルメモリに書き込む処理。これによりグローバルメモリには各ブロックの結果が書き込まれる。

ソースコード解説 (4分割のうち3つ目)

- これより先はCPUで処理されるmain関数
- 3-20行目は初期化処理で、必要なメモリ領域の確保や、入力の準備、正しい答えの準備などを行なっている
- 22行目から40行目はGPUを使用する時の定番処理

ソースコード解説 (4分割のうち4つ目)

- 2-5行目は各ブロックの総和計算結果を足し合わせて最終的な出力を得る処理。GPUで行なっても良いが、あまり時間がかからないと思われるので、今回はCPUで処理する。

2つ目から4つ目のソースコードの説明は上にまとめた通りです。
それでは、このプログラムを実習環境で実行しましょう。

ジョブスクリプト (cuda.nqs)

- 10秒制限
- 実行ファイル名をa.outにする
- 講習会用の待ち行列 (LECTURE) に投入する
- 実習1と同内容

```
1 | #!/bin/bash
2 | #PBS -q LECTURE
3 | #PBS -l elapstim_req=00:00:10,cpunum_job=1,gpunum_job=1
4 | cd $PBS_O_WORKDIR
5 | ./a.out
```

まず、上記のように実習1と同様のジョブスクリプトを書きます。

コンパイルとジョブ投入

- コンパイル(nvcc)により実行ファイルa.outが生成される

```
$ nvcc -O3 -arch=sm_60 ex2.cu
$ qsub cuda.nqs
```

上記のようにコンパイルして、ジョブを投入しましょう。ジョブを投入したら qstat や sstat コマンドで状況確認しながらジョブが終了するのを待ちましょう。ジョブが完了したら標準出力を確認し、正しく実行できたことを確認しましょう。

オキュパンシのプロファイリング

- オキュパンシをプロファイリングするためのジョブスクリプト (occupancy.nqs) を作成する

```
1  #!/bin/bash
2  #PBS -q LECTURE
3  #PBS -l elapstim_req=00:00:10,cpunum_job=1,gpunum_job=1
4  cd $PBS_O_WORKDIR
5  nvprof --metrics achieved_occupancy ./a.out
```

これまでオキュパンシにこだわってアルゴリズムを設計してきましたが、実際にオキュパンシが高くなっていることを確認しましょう。実習 1 とは異なるプロファイリングとして、実際に達成したオキュパンシの値を計測してみます。そのためには、nvprof コマンドのオプションとして、`-metrics achieved_occupancy` とつけます。詳細については、NVIDIA の Achieved Occupancy の説明ページを参照してください。

(<https://docs.nvidia.com/gameworks/content/developertools/desktop/analysis/report/cudaexperiments/kernellevel/achievedoccupancy.htm>)

ジョブ投入

```
$ qsub occupancy.nqs
```

それでは、ジョブを投入し、処理が終わるのを待って標準エラー出力を確認しましょう。高いオキュパンシが実現できているのを確認できるはずです。

2.6.2 Thrust ライブラリとの比較

Thrustライブラリとの比較

ソースコード (ex2.cuを変更)

- ex2.cuのThrustに関する行 (1つ目のソースコードの3,4行, 4つ目のソースコードの19-29行) のコメントアウトを取り除く

ジョブスクリプト (cuda.nqs)

- 前節と同じものを使う

コンパイルとジョブ投入

- コンパイル(nvcc)により実行ファイルa.outが生成される

```
$ nvcc -O3 -arch=sm_60 ex2.cu
$ qsub cuda.nqs
```

- ジョブ完了後, 標準出力ファイルを開き, 自作とライブラリの処理時間をチェックする

最後に CUDA のライブラリである Thrust を使った場合と処理時間を比較してみましょう。ソースコードの Thrust に関する部分のコメントアウトを取り除いてコンパイルした後, 再びジョブを投入してみましょう。そしてジョブ終了後, 出力を確認し, 自作のアルゴリズムとライブラリで処理時間を比較してみましょう。自作アルゴリズムの方が処理時間が短くなっているはずです。

処理時間プロファイリング用のジョブスクリプト (prof.nqs)

- 10秒制限
- 実行ファイル名をa.outにする
- 講習会用の待ち行列 (LECTURE) に投入する
- 実習1と同内容

```
1 | #!/bin/bash
2 | #PBS -q LECTURE
3 | #PBS -l elapstim_req=00:00:10, cpunum_job=1, gpunum_job=1
4 | cd $PBS_O_WORKDIR
5 | nvprof ./a.out
```

処理時間の詳細を調べるためにプロファイリングをしてみましょう。まずそのためのジョブスクリプトを作成します。これは実習1と同じものです。

処理時間のプロファイリング

- ジョブ投入

```
$ qsub prof.nqs
```

- 自作とライブラリでメイン処理の処理時間がほとんど変わらないことを確認する

それではジョブを投入し、処理終了後、標準エラー出力を確認してみましょう。総和の計算部分では、自作アルゴリズムとライブラリはほぼ同じ計算時間になっているはずです。

2.7 高速化トピックス

高速化トピックス

- いくつかの観点ごとに高速化手法を列挙する
 - 設計とデータ格納法
 - データに依存性がある場合
 - 構造体の配列とコアレスアクセス
 - 行列データの共有メモリへの配置方法
 - ライブラリの活用
 - 複数GPUの活用
 - GPUハードウェア(CUDA機能)の活用

33

ここからは高速化トピックスということで、高速化のための技法を列挙していきます。まずはアルゴリズムの設計とデータ格納法に関していくつかのトピックを紹介し、その後、ライブラリの活用、複数GPUの活用、GPUハードウェアの活用などについてトピックを紹介し、

2.7.1 データに依存性がある場合

データ依存性がある場合

- データ依存性があると並列化は難しいが、単純な依存性であれば解決できることもある
- 例1：Prefix sum

	0	1	2	3
入力	①	②	③	④
出力	0	①	①+②	①+②+③

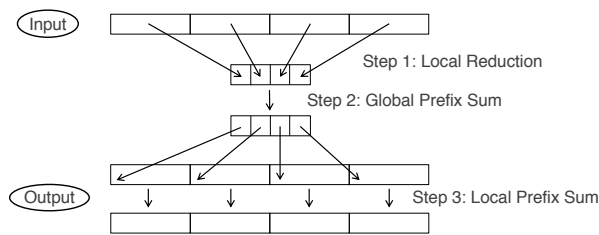
一見前から順に計算するしかないようにも見えるが・・・

34

アルゴリズムの設計とデータ格納法に関する最初のトピックはデータに依存性がある場合の処理についてです。データに依存性があるというのは、ここでは、ある出力データを計算するために別の出力データが計算されるのを待たなければならないという場合を指しています。このような依存性があると、並列化は難しくなるのですが、簡単な依存性であれば解決できることもあります。

その例を2つ挙げたいと思いますが、一つ目はプレフィックスサムです。プレフィックスサムというのは図のような出力を計算する処理です。ざっくりいうと各indexにおいて自分より左にいる入力配列の和を求める計算です。例えば出力配列のindex 2 は入力配列のindex 0 と 1 の和となります。これを計算する場合、どのindexも一つ前のindexに依存性があるため、一見、前から順に計算していくしかないようにも見えます。しかし、このような場合に並列化を行うための定番のやり方があります。

Prefix sumの並列アルゴリズム



先に各ブロックの先頭の値を決めてしまえば
あとは、並列に計算できる

35

プレフィックスサムを複数のマルチプロセッサを用いて並列に計算することを考えます。ここでは、4つのマルチプロセッサを使うものとして説明します。その場合、入力配列を4分割して、各ブロックをマルチプロセッサの一つずつ分配する方法が考えられます。しかし、前のスライドで説明したように、このままでは1つ目のブロックから順に計算していくしかありません。各ブロックのが並列に計算を行えるようにするにはどうすればよいでしょうか。少し考えてみると、実は、各ブロックの先頭の出力値さえ分かればあとは各ブロックで並列に計算が進められるということがわかります。

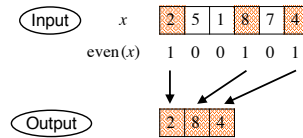
そこで、先に各ブロックの先頭の出力値を求めることを考えます。そのためには各ブロックの総和の値が必要になるので、まずそれを並列で計算します。それが図のステップ1です。

次にステップ1の結果に対してプレフィックスサムを計算します。それが図のステップ2です。実はこれが各ブロックの先頭の出力値になります。例えば、2つ目のブロックの先頭の出力値は、1つ目のブロックの総和になりますし、3つ目のブロックの先頭の出力値は、1つ目のブロックと2つ目のブロックの総和になります。このプレフィックスサムの計算は並列計算せずシーケンシャルに計算することにします。入力要素数に比べてブロック数が十分に小さい場合は、このように計算したとしてもステップ2の計算量は漸近的には無視できます。

ということでめでたく各ブロックの先頭の出力値が計算できたので、あとは、各ブロックにおいて並列にプレフィックスサムを計算します。これがステップ3です。各ブロックのプレフィックスサムの計算についても色々工夫できそうですが、ここでは深入りしないでおきます。

データ依存性がある場合

- 例2: 条件に合うデータをフィルタリング
 - 配列から偶数だけを取り出すとする
 - ⇒データの書き込み先が他のデータに依存する



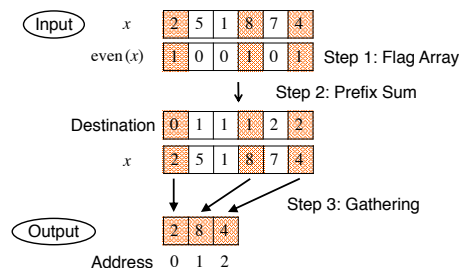
36

データに依存性がある場合の処理についての二つ目の例として、データのフィルタリングを考えてみます。膨大な入力データの中から条件を満たす要素のみを取り出したいという状況です。ここでは、入力配列の中から偶数のみを取り出して、配列としてまとめて出力することを考えましょう。

各要素が偶数かどうかを判定することは並列に行えるのですが、問題は出力の配列を作る部分です。偶数が見つかったとして、それを出力配列のどこに書き込めばよいのかわかりません。それは他の出力に依存します。このような場合への対処についても定番のやり方があります。

データフィルタリングのアルゴリズム

- 先にデータ書き込み先を求める(Prefix sum)



37

このような場合への対処法ですが、データの出力先がわからないことが問題なので、先にそれを求めることにします。実はその計算にはプレフィックスサムを活用できます。

まず、ステップ1として、入力の各要素が偶数かどうかを判定し、結果をフラグ配列に書き込みます。偶数なら1、そうでなければ0を書くことにします。

そして次がポイントなのですが、ステップ2として、このフラグ配列のプレフィックスサムを計算します。そうすると実はプレフィックスサムの値が出力の書き込み先 index になっています。プレフィックスサムにより、自分より左に条件を満たす要素が何個あるかを計算できていることを確認してください。

最後に条件を満たす要素について、並列に出力配列への書き込みを行います。書き込み先についてはステップ2の結果を使います。これがステップ3です。これをGPUで実行する場合、少し工夫する余地がありそうですが、ここでは深入りしないでおきます。

2.7.2 構造体の配列とコアレスアクセス

構造体の配列と コアレスアクセス

- 構造体の配列にアクセスしようとする
とコアレスアクセスが難しい

obj[0] obj[1] obj[2] obj[3]

a	b	c	d	a	b	c	d	a	b	c	d	a	b	c	d
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

各コアがメンバーaにアクセスすると
コアレスアクセスでなくなる

⇒ メンバーごとに配列を作ると
コアレスアクセス可能になる

a	a	a	a	b	b	b	b	c	c	c	c	d	d	d	d
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

38

それでは次のトピックに移ります。次は構造体の配列を作る際の注意事項です。

グローバルメモリにデータを保持する際、構造体の配列を作ると、コアレスアクセスが難しくなることがあります。スライドの上の図では各コアが構造体のメンバー *a* にアクセスしようとしています。飛び飛びのアドレスにアクセスしているため、コアレスアクセスではなくなっています。これへの対処法として、下の図のようにメンバーごとに配列を作ることになればコアレスアクセスが可能になります。

上の図のような構成は AoS (Array of Structures)、下の図のような構成は SoA (Structures of Array) と呼ばれています。GPU プログラミングでは、SoA を活用することで効率的なメモリアクセスができるようになります。

2.7.3 行列の共有メモリへの格納

行列の共有メモリへの格納

- 前述の通り、オキュパンシを考慮する
- 格納位置についても注意が必要
 - 行アクセス、列アクセスでバンクコンフリクトが起こらないようにする ((b)のようにする)

(a) 列アクセスでバンクコンフリクト

(b) 列アクセスでもバンクコンフリクトが発生しない

39

次は行列などの2次元データを扱う際の注意事項です。本日の前半で話した通り、2次元データの列に対して並列アクセスする場合は、工夫が必要です。まず、前半に話した通り、オキュパンシを考慮して、あまり大きいサイズのデータを一気に共有メモリにコピーしないほうがいいです。

さらにコピーする際のデータの配置についても工夫が必要です。バンク数4の共有メモリにマルチプロセッサの4コアが並列アクセスする場合を考えます。図の各列が1つのバンクになっているとします。まず4×4の16データをグローバルメモリから共有メモリにコピーし、左のような順番で並べたとします。①から⑯はデータの順番です。最初の4要素をコアレスアクセスによりグローバルメモリから取得し、1行目にコピーし、次の4要素は2行目にコピーという感じで、16要素を4要素ずつ共有メモリにコピーします。その後、列の各要素に並列アクセスしようとする、全てのコアが同じバンクにアクセスすることになり、バンクコンフリクトが発生します。

これに対し、データの格納方法を工夫することで、バンクコンフリクトを回避することができます。グローバルメモリからコピーした4要素を共有メモリに貼り付ける際、貼り付け先のアドレスを変更し、右の図のようにします。このようにすると、列に並列アクセスする際も各コアが別々のバンクにアクセスすることになり、バンクコンフリクトが発生しません。

2.7.4 ライブラリの活用

ライブラリの活用

- CUDAでは総和計算ですら効率的なコードを書くのは一苦労
⇒ 極力ライブラリに頼るのが良い
- リダクション, ソートなど: Thrust
- FFT: cuFFT
- ライブラリの使い方など:
<https://www.slideshare.net/NVIDIAJapan/1072-cuda>

40

これまで話してきた通り, GPU では総和計算ですら効率的なコードを書くのは一苦労です. 様々な問題に対し, GPU 上で動くプログラムを書きたいという場合には, 極力ライブラリに頼るのが良いと思います. リダクションやソートなどの基本的な処理については, Thrust というライブラリがありますし, FFT についても cuFFT があります. 他にもたくさんのライブラリがありますので, 例えばスライドの URL を参考に見てみてください.

2.7.5 複数 GPU の活用

複数GPUの活用

- 1ノードに複数のGPUを設置して高速化
- 以前はGPU間の通信にはCPUを介す必要があり低速だった
- 現在はGPU間が直接通信できるため高速 (NVLink)

41

一つのマシンに複数の GPU を取り付けて高速化することができます. 以前は GPU 間の通信は一旦 CPU を介す必要があり低速でしたが, 現在は GPU 間で直接通信ができるため高速です. 初日に紹介した CUDA C プロフェッショナルプログラミングという本には複数 GPU を使用したプログラミングについての詳細な解説があります. 興味のある方はチェックしてみてください.

2.7.6 GPUハードウェアの活用

GPUハードウェア (CUDA機能) の活用

- GPUハードウェア
 - コンスタントメモリとテクスチャメモリ
 - リードオンリーメモリを活用してアクセスを高速化する
 - ユニファイドメモリ
 - ホストとデバイスのデータ同期を自動化
 - 必要なデータのみが同期されることで効率化できる
 - テンソルコア
 - 行列計算をハードウェアで行うことで高速化する

42

最後に GPU のハードウェアを活用した高速化について簡単に紹介したいと思います。

まずはコンスタントメモリとテクスチャメモリです。どちらもリードオンリーのデータのためのメモリです。コンスタントメモリは GPU デバイス内に存在する定数格納用の小容量メモリですが、各マルチプロセッサ内のコンスタントキャッシュと呼ばれる部分に値がキャッシュされるため、高速アクセスが可能です。メモリサイズは今の所どのモデルでも 64KB になっています。全スレッドが同じ値を読み込む必要があるような場合に高速にデータ読み込みを行えます。テクスチャメモリとはグラフィック処理におけるテクスチャ、つまりオブジェクト表面に貼り付けられる画像を格納するためのメモリで、2次元データを扱う際、アクセスした場所の近傍のデータへのアクセスが高速化されるように設計されています。アプリケーションによっては、グローバルメモリにアクセスするよりも高速なアクセスが可能です。

次はユニファイドメモリです。これは、CPU と GPU 間のデータ同期を自動化する技術です。ユニファイドメモリを使わない場合、CPU と GPU の間のデータ転送はすべてプログラムで記述しなければいけません。これを使う場合は、プログラムにデータコピーを記述しなくても、必要に応じて自動で同期が行われます。不要なデータコピーがなくなることで、プログラムが高速化されることもあるようです。

テンソルコアは 4×4 の行列演算をハードウェアで行うモジュールで、行列計算を高速化できます。近年 GPU が深層学習に多く活用されていますが、そこでは行列計算を高速化することが重要です。テンソルコアを活用することでその処理が高速化されます。

GPUハードウェア (CUDA機能) の活用

- CUDA機能
 - マルチストリーム
 - GPUに異なる複数のタスクを行わせたり, GPU処理中にCPUで別の処理を行ったりできる
 - ダイナミックパラレルリズム
 - GPUから別のGPUカーネル関数を呼ぶことができる
 - CPUを介さずに動的にタスクを追加できる

43

ここからはハードウェアというよりは CUDA の機能と言った方がいいかもしれませんが、それをいくつか紹介します。

まず、マルチストリームですが、複数のタスクを GPU に同時に行わせたり、GPU の処理中に CPU で別の処理を行ったりできる機能です。

次に、ダイナミックパラレルリズムですが、GPU の処理中に別の GPU カーネル関数を呼ぶことができる機能です。これにより、CPU を介さずに動的なタスクの追加などができるようになります。

GPUハードウェア (CUDA機能) の活用

- CUDA機能
 - Warpシャッフル
 - 他のスレッドに割り当てられたレジスタを読むことができる
 - 独立型スレッドスケジューリング (Volta)
 - 命令をワープ内のスレッドごとに異なるタイミングで実行できる

44

次がワープシャッフルです。これはスレッドがワープ内の他のスレッドのレジスタを読めるようにする機能です。これによりワープ内のスレッドとのデータのやり取りが高速化されます。これまではワープ内のスレッド間でのデータのやり取りは、共有メモリを介して行う必要がありましたが、ワープシャッフルにより、他のスレッドに直接アクセスできるようになりました。

最後が独立型スレッドスケジューリングです。これは 2017 年に発表された Volta アーキテクチャに搭載された機能なので、今回使用している P100 では使用できません。これはワープ内のスレッドごとに異なるタイミングで命令を実行できるようにする機能です。この講習会では、ワープ内のスレッドは常に同じ命令を同じタイミングで実行すると説明してきましたが、Volta アーキテクチャではこれは当てはまりません。この機能によりアプリケーションによってはプログラミングがしやすくなりますが、ワープ内スレッドの同期実行を前提に作成されたプログラムは正常に動作しなくなることもあります。

付録：Linux 入門

Linux入門

一関高専 未来創造工学科 情報・ソフトウェア系
小池 敦

1

今日の内容

- Linuxの基本
 - 端末（ターミナル）の基本操作
 - ユーザーとグループ、スーパーユーザー
 - パーミッション
 - パイプラインとリダイレクト
 - 環境変数と設定ファイル
 - Linuxディストリビューション
- 基本的なLinuxコマンド（別資料）
 - ネットワーク関連コマンド

2

端末（ターミナル）の基本操作

- \$ (プロンプト)の後ろにコマンドを打つ
 - \$ ifconfig
- 基本的な操作
 - 矢印左右：カーソルの移動
 - 矢印上下：履歴の表示
 - Tabキー：入力の補完 (\$ ifc[Tabキー] ⇒ifconfig)
 - Ctrl-c：ジョブ（コマンド実行）の強制終了
 - Ctrl-a,e：カーソルを先頭（最後）に移動する
 - Ctrl-d: カーソル位置の文字（前方文字）を消去

3

練習

- ping www.google.co.jp をした後、
中断して、tracert www.google.co.jp
を実行してください

実行例

```
$ ping www.google.co.jp
```

```
Ctrl-c
```

```
上矢印 Ctrl-a Ctrl-d(4回) tracer[Tab]
```

```
$ traceroute www.google.co.jp
```

4

端末（ターミナル）の基本操作

- 端末はどこかのディレクトリ（フォルダ）で作業している
 - 端末起動時はホームディレクトリ (/home/ユーザー名)にいる
- 現在の作業ディレクトリの表示
 - `$ pwd`
- 特別なディレクトリ名
 - `~` : ホームディレクトリ (/home/ユーザー名)
 - `.` : 作業ディレクトリ
 - `..` : 作業ディレクトリの一つ上のディレクトリ

5

端末（ターミナル）の基本操作

- 作業ディレクトリの変更
 - `$ cd /etc` ⇒ /etc ディレクトリに移動
 - `$ cd ..` ⇒ 一つ上のディレクトリに移動
 - `$ cd` ⇒ ホームディレクトリに移動
- ディレクトリ内のファイル一覧表示
 - `$ ls`
 - `$ ls -l` ⇒ 各ファイルの詳細情報を表示
 - `$ ls -a` ⇒ 隠しファイルを含む全ファイルを表示
 - `$ ls -la` ⇒ 全ファイルを詳細表示

6

ユーザーとグループ

- ユーザーは必ずどこかのグループに所属
 - 複数のグループに属することも可能
- 各ファイルにも所有者とグループが設定される
 - `$ touch test.txt` ⇒ test.txtファイルを作成
 - `$ ls -l` ⇒ 所有者とグループを表示
- ファイルの所有者情報を変更するには
chownコマンドを使う

7

スーパーユーザー

- スーパーユーザー(rootユーザー)
 - 管理者権限を持つ
 - 強力な権限を持つのでスーパーユーザーで作業する時は作業ミスに注意する
- 一般ユーザーの一部はスーパーユーザー権限でのコマンド実行を許可されている

8

スーパーユーザー権限での コマンド実行

- コマンドの前にsudo をつける
 - 例：\$ sudo less /var/log/secure
⇒ 管理者権限でセキュリティ関連ログをチェック
- スーパーユーザーに変身
 - \$ su -
⇒ スーパーユーザーに変身（戻るにはexit）

9

ユーザーとグループ、スーパー ユーザー

- sudoを実行可能なユーザー
 - /etc/sudoersに記載されている
(本ファイルを編集する時はvisudoを使うこと)
- 一部の慣習として、wheelグループに属するユーザーにsudo権限を与えることがある
 - ユーザーにsudo権限を与えるには、
そのユーザーをwheelグループに所属させればよい

10

パーミッション

- すべてのファイルにアクセス制限（パーミッション）が設定されている
 - 誰に？
 - ファイル所有者
 - グループ
 - それ以外
 - 制限内容（これ以外もあるが省略）
 - 読み込み
 - 書き込み
 - 実行
- 2進数で指定する
例：1 1 0 (=6₁₀) ⇒ 読み込み書き込み可、実行不可
- 読み込み 書き込み 実行

11

パーミッション

- パーMISSIONの指定
 - ファイル所有者、グループ、それ以外の順に制限内容を10進数で書く
 - 例：6 4 4
 - 所有者 : 6 (=110₂)
 - グループ : 4 (=100₂)
 - それ以外 : 4 (=100₂)
- パーMISSIONを変更するにはchmodを使う
 - `$ chmod 600 test.txt`
⇒ test.txtのパーMISSIONを600に変更する

12

パイプライン

- あるコマンドの出力を次のコマンドの入力にする
 - 例 `$ ps aux` ⇒ 実行中の全プロセスを表示する
 - `$ ps aux | grep ssh`
 - ⇒ 実行中の全プロセスから sshが含まれる行を検索する
- (パイプライン利用してps auxの出力を grep ssh の入力にしている)

13

リダイレクト

- あるコマンドの出力をファイルに保存する
 - 例 `$ ps aux > test.txt`
 - ⇒ ps auxの実行結果をtest.txtに保存する
 - `$ ifconfig >> test.txt`
 - ⇒ ifconfigの実行結果をtest.txtに追記する (既存の内容の後ろに追記する)
 - `$ cat test.txt` ⇒ test.txtを表示する
- エラーメッセージのみを別ファイルに保存可能
 - 例 `$ aaa 2> error.txt` (追記したければ `2>>` とする)
 - `$ cat error.txt`

14

環境変数と設定ファイル

- 環境変数：設定情報を保持する変数
 - 例：PATH
⇒ パスなしでファイルを指定した際に探すディレクトリ
- 環境変数一覧表示： `$ printenv`
- 環境変数の編集
 - `~/.bash_profile`もしくは`~/.bashrc`に設定を書く
 - どちらに書いてもそんなに変わらない
(`.bash_profile`に書く人が多い?)
 - 設定後は再読み込みを行う `$ source ~/.bash_profile`

15

Linuxのディストリビューション

- Linux: UNIXを参考にヘルシンキ大学の学生だったリーナス・トーバルズ氏が開発
- オープンソースライセンスのため誰でも自由に改変して配布できる（ディストリビューション）
 - Debian系とRed Hat系が主流
- Linuxディストリビューションの例
 - CentOS: Red Hat系。企業サーバー等でよく使われる。
本演習環境はCentOS7。
 - Ubuntu: Debian系。デスクトップOSとして使われる。
 - Amazon Linux: Red Hat系。
Amazonのクラウドサービス (AWS) で利用できる。

16

基本的なLinuxコマンド

基本	ファイル操作	ファイル閲覧,編集等	管理者	ネットワーク	その他
man	cp	less	sudo	ifconfig	ps
pwd	mv	grep	su	ping	kill
ls	rm	find	shutdown	traceroute	df
cd	mkdir	cat	useradd	netstat	tar
	touch	diff	passwd	nslookup	
	chmod			dig	
	chown	vi		arp	
		emacs		ip	
				ss	
				curl	
				ssh	

基本

man	コマンドのマニュアル表示
コマンドや設定ファイル等のマニュアル表示。 表示後の操作は less を参照。	
[オプション]	
[使用例]	
\$ man ls	⇒ lsコマンドのマニュアルを表示
\$ man sshd_config	⇒ SSH設定ファイルsshd_configの設定マニュアルを表示

[このページのtopへ](#)

pwd	作業ディレクトリの表示
現在作業しているディレクトリ（フォルダ）の絶対パスを表示する。	
[オプション]	
[使用例]	
\$ pwd	⇒ 作業ディレクトリを表示

[このページのtopへ](#)

ls	ディレクトリ内のファイル一覧表示
作業ディレクトリ（フォルダ）に存在するファイル、ディレクトリの一覧を表示する。	
[オプション]	
-l: 詳細(サイズ,所有者,最終更新,パーミッション)表示。 -a: すべてのファイル,ディレクトリを表示(隠しファイル(ドットで始まるファイル)も表示)	

ls	ディレクトリ内のファイル一覧表示
[使用例]	
\$ ls	⇒ 作業ディレクトリ内のファイルを表示
\$ ls -a	⇒ 作業ディレクトリ内のすべてのファイルを表示
\$ ls -la	⇒ 作業ディレクトリ内のすべてのファイルの詳細を表示
\$ ls -la less	⇒ 作業ディレクトリ内のすべてのファイルの詳細をlessを使って表示

[このページのtopへ](#)

cd	ディレクトリの移動
作業ディレクトリ（フォルダ）を変更する。	
[オプション]	
[使用例]	
\$ cd /etc/	⇒ /etc ディレクトリに移動
\$ cd	⇒ ホームディレクトリ ~ (/home/ユーザの名前/) に移動
\$ cd ~/html/	⇒ ~(ホームディレクトリ) の下の html ディレクトリに移動
\$ cd tmp/	⇒ 作業ディレクトリ(.) の下の tmp ディレクトリに移動
\$ cd ..	⇒ 作業ディレクトリ(.) の一つ親のディレクトリ(..) に移動

[このページのtopへ](#)

ファイル操作

cp	ファイルのコピー
ファイルやディレクトリをコピーする。	
[オプション]	
-R: サブディレクトリも含めてすべてコピーする	
[使用例]	
\$ cp readme.txt readme.backup	⇒ readme.txt と同内容のファイルreadme.backupを作る
\$ cp ~/Downloads/Anaconda3-2.3.0-Linux-x86_64.sh .	⇒ Anaconda3-2.3.0-Linux-x86_64.shを作業ディレクトリ(.)にコピー
\$ cp -R src_dir target_dir	⇒ src_dirを丸ごとtarget_dirにコピーする
\$ cp -R src_dir/ target_dir	⇒ src_dirの中身を丸ごとtarget_dirにコピーする (src_dir自身はコピーされない)

[このページのtopへ](#)

mv	ファイル移動、名称変更
ファイルやディレクトリの移動や名称変更を行う。	
[オプション]	
[使用例]	
\$ mv main.c src/	⇒ main.cを作業ディレクトリ直下のsrcディレクトリに移動する
\$ mv old_name.txt new_name.txt	⇒ ファイルold_name.txtの名称をnew_name.txtに変更する

[このページのtopへ](#)

rm	ファイル削除
ファイルやディレクトリを削除する。	
[オプション]	
-r: サブディレクトリも含めてすべて削除する (-Rと同じ)	
-f: 確認メッセージを表示せずに即座に削除する	
[使用例]	

rm	ファイル削除
\$ rm index2.html	⇒ ファイルindex2.htmlを消去する
\$ rm -rf tmp	⇒ tmpフォルダを消去確認なしに完全に削除する

[このページのtopへ](#)

mkdir	ディレクトリの作成
ディレクトリを作成する。	
[オプション]	
[使用例]	
\$ mkdir src	⇒ srcディレクトリを作成する

[このページのtopへ](#)

touch	ファイルの新規作成、タイムスタンプ変更
ファイルの新規作成やファイルのタイムスタンプ（最終アクセス時刻、最終変更時刻）変更を行う。	
[オプション]	
-t: ファイルのアクセス時刻、変更時刻を修正する	
[使用例]	
\$ touch readme.txt	⇒ ファイルreadme.txtを作成する（中身は空）
\$ touch -t 201805311300 readme.txt	⇒ readme.txtのタイムスタンプを2018/5/31/ 13:00に変更する

[このページのtopへ](#)

chmod	ファイルのパーミッションの変更
ファイルのパーミッションを変更する。 "所有者(u)", "所属グループ(g)", "そのほか(o)"のそれぞれに対し "読み出し(r)", "書き込み(w)", "実行(x)" の許可(1)不許可(0)を設定。 rwxとも許可なら111で2進数から10進数に直し7が設定値となる。	
[オプション]	
-R: サブディレクトリも含めてすべてファイルのパーミッションを変更する	
[使用例]	
\$ chmod 755 a.out	⇒ a.outのパーミッションを755に設定する（所有者のみ全操作が可能で、それ以外は読み出し、実行が可能）
\$ chmod 644 readme.txt	⇒ readme.txtのパーミッションを644に設定する（所有者のみ読み込み、書き込みが可能で、それ以外は読み出しのみ可）
\$ chmod 600 id_rsa	⇒ id_rsa（SSH用の秘密鍵）のパーミッションを600に設定する（所有者のみ読み込み、書き込みが可能で、それ以外は読み出し書き込み不可）

[このページのtopへ](#)

chown	ファイル所有者の変更
ファイルの所有者やグループを変更する。	
[オプション]	
-R: サブディレクトリも含めてすべてファイルのパーミッションを変更する	
[使用例]	
\$ chown alice index.html	⇒ index.htmlの所有者をaliceに変更する
\$ chown alice:group1 index.html	⇒ index.htmlの所有者をaliceにグループをgroup1に変更する

[このページのtopへ](#)

ファイル閲覧・編集等

less	テキストファイルの内容表示
テキストファイルの内容を表示する。文書内の移動や単語の検索が可能である。	
[less起動中の操作]	
q : less を終了	
h : help を表示	
矢印キー上下 : 上下に1行分スクロール	
u,d : 上下に半ページ分スクロール	
スペースキー : 下に1ページ分スクロール	
b : 上に1ページ分スクロール	
g,G : ファイルの先頭,最後に移動	
/hoge : hoge を検索。n,Nで前後のhogeにジャンプ	
[使用例]	
\$ less readme.txt ⇒ readme.txtをlessを起動して閲覧する	
\$ ls -l /etc/ less ⇒ /etc/ ディレクトリの中身をlessを使って閲覧する	

[このページのtopへ](#)

grep	文字列の検索
指定した入力ファイルの中から条件を満たす行を検索する。正規表現を利用可能。	
[オプション]	
-i : 大文字と小文字の区別をしない	
[使用例]	
\$ grep DirectoryIndex /etc/httpd/conf/httpd.conf ⇒ httpd.confの中からDirectoryIndexを含む行を検索する	
\$ ps aux grep ssh ⇒ ps (プロセス状態表示) の実行結果からsshという文字を含む行を検索する	

[このページのtopへ](#)

find	ファイル検索
条件を満たすファイルを検索する。	
[オプション]	
-name : 検索するファイル名の指定 (部分一致条件を指定するには、"*hoge*"のようにワイルドカードを使う)	
-iname : 検索するファイル名の指定 (大文字小文字を区別しない)	
-type : ファイルの種類を指定する (-type f とすると通常ファイルのみを検索する (ディレクトリ等は排除される))	
[使用例]	
\$ find ~ -name "*bash*" ⇒ ホームディレクトリの中のすべてのディレクトリの中からbashという名前を含むファイルを検索する	
\$ find / -name network -type f 2>/dev/null ⇒ コンピュータ全体からnetworkという名前のファイルを検索する (2>/dev/null はエラーメッセージを捨てる処理)	

[このページのtopへ](#)

cat	ファイル内容の表示
ファイルの中身を表示する。lessのように文書内を移動するような機能はなく、ただ出力するだけである。	
[オプション]	
[使用例]	
\$ cat ~/.bashrc ⇒ .bashrcの中身を表示する	

[このページのtopへ](#)

diff	2ファイルの差分表示
------	------------

diff	2ファイルの差分表示
2つのファイルの差分を探して表示する。	
[オプション]	
-u : unified形式 (一つのコンテンツの中に両ファイルの内容を表示) で差分を表示する	
-y : 2つのファイルを左右に分けて表示する	
[使用例]	
\$ diff test.backup test.txt ⇒ test.backupとtest.txtの差分を表示する	

[このページのtopへ](#)

vi	テキストエディタviの起動
テキストエディタviを起動する。	
[vi起動中の操作]	
viには2つのモードがあるので、状況に応じて使い分ける	
<ul style="list-style-type: none"> • コマンドモード：検索、移動、セーブ、終了などのコマンドが使えるが文字入力ができない。iを入力すると編集モードに移行する • 編集モード：文字を入力するためのモード。Escキーを入力するとコマンドモードに戻る。 	
コマンドモードでの主なコマンド	
<ul style="list-style-type: none"> • i ⇒ 編集モードに移行し、現在のカーソル位置に文字入力を行えるようにする (Escキーでコマンドモードに戻る) • :wq ⇒ セーブしてviを終了する • :q! ⇒ セーブせずにviを終了する • :w ⇒ セーブする • dd ⇒ 1行削除 (切り取り) • u ⇒ 直前の操作の取り消し (アンドゥー) • /hoge ⇒ 文字列hogeを検索する。n,Nで前後のhogeに移動する。 • Ctrl-d ⇒ 半ページ下に移動する • Ctrl-u ⇒ 半ページ上に移動する • G ⇒ 最終行に移動する • 行数G ⇒ 「行数」で指定した行に移動する。先頭行なら1Gと入力する。 	
[使用例]	
\$ vi ~/.bash_profile ⇒ .bash_profileをviで編集する	

[このページのtopへ](#)

emacs	テキストエディタemacsの起動
テキストエディタemacsを起動する。環境によってはインストールされていないので注意。	
[emacs起動中の操作]	
C-x C-c ⇒ 終了する(CはCtrlキー；Ctrlを押しながらxとcを押す)	
C-x C-s ⇒ 保存する	
C-x C-w ⇒ 別名保存する	
C-x C-f ⇒ ファイルを開くまたは新規作成	
C-g ⇒ キャンセル (コマンド中断)	
C-@ ⇒ 範囲指定開始	
C-w ⇒ 指定した範囲をカット(キル)	
M-w ⇒ 指定した範囲をコピー(Mは通常はEscキー；Escを押した後wを押す)	
C-k ⇒ カーソルから行末までをカット	
C-y ⇒ 張り付け(ヤंक)	
C-x 1 ⇒ ウィンドウの分割をやめる (C-x 2 とするとウィンドウを縦に分割)	
C-s ⇒ 前向き検索	
C-r ⇒ 後ろ向き検索	
M-% ⇒ 問い合わせ置換(Mは通常はEscキー；Escを押した後%を押す)	
C-/ ⇒ もとに戻す(Undo)	
[使用例]	

emacs

テキストエディタemacsの起動

\$ emacs ~/.bash_profile ⇒ ~/.bash_profileをemacsで編集する

[このページのtopへ](#)

管理者

sudo

管理者としてコマンドを実行する

スーパーユーザー（管理者）としてコマンドを実行する。sudoを許可されたユーザーしか本コマンドを実行できない。一部のLinuxではwheelグループに属しているユーザーに権限を与えるという慣習もある。

[オプション]

[使用例]

\$ sudo less /var/log/secure ⇒ 管理者権限でセキュリティ関連ログを見る

[このページのtopへ](#)

su

他ユーザーに切り替わる

他のユーザーに切り替わる。元のユーザーに戻るにはexitと入力する。

[オプション]

- : 環境変数を新ユーザーに引き継がない（新ユーザーの環境を利用する）

[使用例]

\$ su - alice ⇒ ユーザーaliceに切り替わる

\$ su - ⇒ スーパーユーザー（管理者）に切り替わる

[このページのtopへ](#)

shutdown

シャットダウン、再起動

マシンのシャットダウンや再起動を行う。管理者権限が必要。

[オプション]

-h : シャットダウンする

-r : 再起動する

[使用例]

\$ sudo shutdown -h now ⇒ すぐにマシンをシャットダウンする

\$ sudo shutdown -r now ⇒ すぐにマシンを再起動する

[このページのtopへ](#)

useradd

ユーザーの追加

システムにユーザーを追加する。

[オプション]

-G : ユーザーが所属する補助グループを指定する

[使用例]

\$ sudo useradd alice ⇒ システムにユーザーaliceを追加する

\$ sudo useradd -G wheel alice ⇒ システムにユーザーaliceを追加し、wheelグループに加える（sudo実行権限を与える）

[このページのtopへ](#)

passwd

パスワード変更

ユーザーのパスワードを変更する

[オプション]

[使用例]

passwd

パスワード変更

\$ sudo passwd alice ⇒ ユーザーaliceのパスワードを変更する

[このページのtopへ](#)

ネットワーク

ifconfig

ネットワーク設定情報出力

ネットワーク設定情報の表示や変更を行う。

[オプション]

-a: 停止しているものも含めすべてのネットワークインターフェースの状態を表示する

[使用例]

\$ ifconfig ⇒ ネットワーク設定情報を表示する

[このページのtopへ](#)

ping

ネットワーク疎通確認

指定ノードへのネットワークの到達性を確認する。ネットワークに問題がなくてもpingに回答しないサーバーもあるので注意。Ctrl-c で中断する。

[オプション]

-c: 指定回数のみパケットを送付する

[使用例]

\$ ping www.google.co.jp ⇒ www.google.co.jpとの疎通確認を行う

\$ ping -c 5 www.google.co.jp ⇒ www.google.co.jpとの疎通確認を5回行う

[このページのtopへ](#)

traceroute

ネットワークの経路表示

指定ノードへのネットワーク経路を表示する。ネットワーク上のどこで通信が失敗しているかやボトルネックがどこかなどの解析に利用できる。ネットワークに問題がなくても応答しないサーバーもあるので注意。

[オプション]

[使用例]

\$ traceroute www.google.co.jp ⇒ www.google.co.jpへの経路を表示する

\$ traceroute www.tohoku.ac.jp ⇒ www.tohoku.ac.jpへの経路を表示する (サーバーが応答しない例)

[このページのtopへ](#)

netstat

ソケット一覧表示

ノードに設定されているソケット (他マシンとの通信路設定) の一覧を表示する。

[オプション]

-t: tcpソケットのみを表示する (--tcpとしても同じ)

-r: OSが保持するルーティングテーブルを表示する

-i: ネットワークインターフェースの統計状態を表示する

[使用例]

\$ netstat -t ⇒ TCPソケット一覧を表示する

\$ netstat -r ⇒ ルーティングテーブルを表示する

\$ netstat -i ⇒ インターフェースごとの統計情報を表示する

[このページのtopへ](#)

nslookup

ドメイン名からIPアドレスを調べる

ドメイン名から対応するIPアドレスを調べる。

nslookup	ドメイン名からIPアドレスを調べる
[オプション]	
-type: レコード (サーバーの種類等) の指定 (AAAA:IPv6アドレス、MX:メール、NS:ネームサーバー)	
[使用例]	
\$ nslookup www.ichinoseki.ac.jp ⇒ 一関高専WebサーバーのIPアドレスを調べる	
\$ nslookup -type=AAAA www.google.co.jp ⇒ Google検索サイトのIPv6アドレスを調べる	

[このページのtopへ](#)

arp	ARPテーブルの表示
OSで管理しているARPテーブル (IPアドレスとMACアドレスの対応表) の表示や設定を行う。	
[オプション]	
-a: ARPテーブルを表示する。	
-d: 指定したホストのARPエントリを廃棄する (ネットワーク設定が変更になった場合等に用いることができる)	
[使用例]	
\$ arp -a ⇒ ARPテーブルのすべてのエントリを表示する	
\$ sudo arp -d 192.168.0.1 ⇒ ARPテーブルから192.168.0.1に関するエントリを削除する	

[このページのtopへ](#)

dig	ネームサーバーへのドメイン名問い合わせ
ネームサーバーにドメイン名を問い合わせて結果を表示する。nslookupに比べ問い合わせ結果があまり加工されずに表示されるので、ネットワーク管理者にとっては扱いやすい (らしい)。	
[オプション]	
[使用例]	
\$ dig www.ichinoseki.ac.jp ⇒ 一関高専WebサーバーのIPアドレスを調べる	
\$ dig gmail.com mx ⇒ gmail.comのメールサーバーのIPアドレスを調べる	

[このページのtopへ](#)

ip	ネットワーク設定の表示
ネットワーク設定情報の表示や変更を行う。現在ifconfig等のコマンドはメンテナンスがされていないため非推奨となっており、それらの代替となるコマンドがipである。機能が多岐にわたるため、「ip addr」や「ip link」のようにサブコマンドを後ろにつけて使用する。情報表示だけでなく変更もできるがこの表では省略する。	
[サブコマンド]	
addr: 自身のIPアドレス等のネットワーク情報を表示する (aとしても同じ)	
route: ルーティングテーブルの情報を表示する (rとしても同じ)	
neighbor: ARPテーブルの情報を表示する (neighやnとしても同じ)	
help: ヘルプを表示する (nとしても同じ)	
[使用例]	
\$ ip a ⇒ 自身のIPアドレス情報を表示する (レガシーコマンドの ifconfig に相当; aはaddrの省略形)	
\$ ip r ⇒ ルーティングテーブルの情報を表示する (レガシーコマンドの netstat -r に相当; rはrouteの省略形)	
\$ ip n ⇒ ARPテーブル情報を表示する (レガシーコマンドの arp -a に相当; nはneighborの省略形)	
\$ ip -s link ⇒ インターフェースごとの統計情報を表示する (レガシーコマンドの netstat -i に相当)	

[このページのtopへ](#)

88

ソケット情報表示

ノードに設定されているソケット（他マシンとの通信路設定）の情報を表示する。現在 netstat コマンドはメンテナンスがされていないため非推奨となっており、その代替となるコマンドが ss である。

[オプション]

-t : tcpソケットのみを表示する (--tcpとしても同じ)

[使用例]

\$ ss -t ⇒ TCPソケット一覧を表示する（シグナーコマンドの netstat -t に相当）

[このページのtopへ](#)

curl

URLからのデータ取得

指定したURLからデータを取得する

[オプション]

-L : リダイレクトされた場合でもデータを取得する

-O : 取得したデータをファイル保存する（このオプションを指定しないときは標準出力に出力される）

[使用例]

\$ curl -L "http://zip.cgis.biz/xml/zip.php?zn=0218511" ⇒ 郵便番号検索APIで021-8511(一関高専)を検索し、結果を表示する

\$ curl -OL "http://weather.livedoor.com/forecast/rss/area/030010.xml" ⇒ Livedoorから盛岡の天気情報をダウンロードし、ファイルに保存する

[このページのtopへ](#)

ssh

SSHクライアントプログラムの起動

SSHクライアントプログラムを起動する

[オプション]

-p : ポート番号を指定する（指定しないときは22）

[使用例]

\$ ssh -p 10022 alice@198.51.100.14 ⇒ ユーザーaliceで192.51.100.14のポート10022にSSHログインする

[このページのtopへ](#)

その他

ps

実行中プロセス一覧表示

実行中のプロセスの一覧を表示する。

[オプション]

a : 自ユーザー以外のプロセスも表示する

u : 読みやすさを重視したフォーマットで出力する

x : 端末を持たないプロセス（OSが利用するプロセス等）も表示する

f : 階層構造（呼び出し関係）が分かるように出力する

[使用例]

\$ ps uf ⇒ 自ユーザーが起動したプロセスの一覧をみる

\$ ps auxf | less ⇒ すべてのプロセスの一覧をlessを使って閲覧する

[このページのtopへ](#)

kill

プロセスの強制終了

指定したプロセスID（PID）のプロセスを終了させる。PIDはpsコマンドを用いて調べられる。

kill	プロセスの強制終了
[オプション]	
[使用例]	
\$ kill 24567 ⇒ PIDが24567のプロセスを強制終了する	

[このページのtopへ](#)

df	ディスクの空き容量表示
ディスクの空き容量を表示する	
[オプション]	
-h: 人間に読みやすいフォーマットで表示する (容量をG,Mなどのサイズ文字を用いて表示する)	
[使用例]	
\$ df -h ⇒ ディスクの空き容量を表示する	

[このページのtopへ](#)

tar	ファイルの解凍・圧縮
ファイルの解凍・圧縮を行う。	
[オプション]	
-x: 解凍	
-c: 圧縮	
-v: 状況表示	
-f: 圧縮ファイル名	
[使用例]	
\$ tar -xvf hoge.tar ⇒ hoge.tarを解凍する(gzファイルやbz2ファイルの解凍も同じコマンド)	
\$ zip -r hoge.zip foo ⇒ fooフォルダをhoge.zipという名前でzip圧縮する (参考)	
\$ unzip hoge.zip ⇒ hoge.zipを解凍する (参考)	

[このページのtopへ](#)