

## OpenCLを用いたGPGPUの数値計算高速化の検討

A study on accelerating numerical calculation of GPGPU using OpenCL

南 昇吾\*

Shogo Minami

魚谷 智志\*

Mikuro Yamamoto

祢津 佑\*

Keisuke Usukura

増田 信之\*

Nobuyuki Masuda

## 1. まえがき

近年,PCの画像処理用に使用されるGPU(Graphics Processing Unit)を数値計算に用いるGPGPU(General purpose computing on GPU)が一般的に使用されるようになってきた。GPGPUは,基本的に演算を並列化することで計算の高速化を図る。そのため,演算の中に条件分岐処理が含まれる場合,計算プロセッサ同士の待ち時間が生じるため,パフォーマンスが落ちてしまう[1]。そこで本研究では,ヘテロジニアス並列コンピューティング環境のためのフレームワークであるOpenCLを用いて,GPGPUにおける条件分岐処理が多く含まれる数値計算の高速化について検討を行った。計算題材として,FDTD法(Finite Difference Time Domain method)を使用した。

具体的には,条件分岐を含む領域と含まない領域とで計算領域を分割し,それぞれを別のカーネル関数として実行させる。結果として,計算量が少ない場合では分割数を少なく,計算量が多い場合では分割数を多くすることで演算が高速となることが分かった。

## 2. FDTD法

FDTD法は,電磁波シミュレーションなどに用いられる電磁界解析法の一つである。マクスウェル方程式を差分化し,時間領域で解くことから,最終的な結果だけでなく,結果に至るまでの電磁界の変化を解析することができる。本稿では,境界条件としてMurの吸収境界条件を用いた二次元FDTD法について解析を行った。

## 2.1 FDTDの計算

今回は電解が1方向成分のみの場合であるTE波を考える。マクスウェル方程式を時間,空間に対して中心差分し,以下の二次元FDTD法の式を得る。

$$H_x^{n+1/2}(i, j+1/2) = H_x^{n-1/2}(i, j+1/2) - \Delta t / \mu \Delta y \{ E_z^n(i, j+1) - E_z^n(i, j) \} \quad (1)$$

$$H_y^{n+1/2}(i+1/2, j) = H_y^{n-1/2}(i+1/2, j) - \Delta t / \mu \Delta x \{ E_z^n(i+1, j) - E_z^n(i, j) \} \quad (2)$$

$$E_z^{n+1}(i, j) = E_z^n(i, j) - \Delta t / \epsilon \Delta y \{ H_x^{n+1/2}(i, j+1/2) - H_x^{n+1/2}(i, j-1/2) \} + \Delta t / \epsilon \Delta x \{ H_y^{n+1/2}(i+1/2, j) - H_y^{n+1/2}(i-1/2, j) \} \quad (3)$$

ここで,空間離散間隔を $\Delta x, \Delta y$ とし,時間離散間隔を $\Delta t$ とする。また, $n$ は時間ステップ, $E$ は電界, $H$ は磁界を表す。また,パラメーター $\epsilon$ と $\mu$ はそれぞれ誘電率と透磁率を表す。Murの吸収境界条件については別紙に譲る[2]。

\*東京理科大学基礎工学部

## 2.2 計算手順

計算手順は,ある時間ステップにおける磁界と電界の値をGPUで並列計算し,計算結果をグローバルメモリ上に保存させる。ステップを進めて再び並列計算を行い,求めたい時間ステップ分繰り返し計算を行ったら演算を終了する。本稿では,磁界計算部分に条件分岐が発生しないため,吸収境界条件を含む電界計算部分についてのみカーネルの分割を検討した。

## 3. カーネル分割方法

本稿では,電界計算部分のカーネルを7通りの分割方法について比較した。カーネルの分割方法について次の図1に示した。

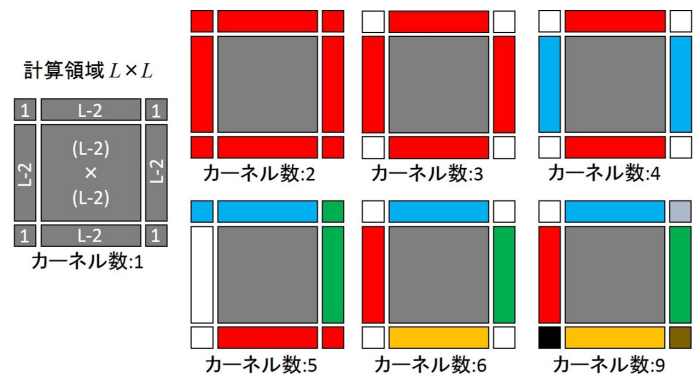


図1: 分割方法

図1は,同じ色で塗られた領域が同一のカーネルとなるように表記してある。また,吸収境界条件に当たる領域は計算量が中央部分に比べて比較的少ない周辺部分が相当する。

## 4. 実験結果

## 4.1 実験環境

ここでは,磁界計算と電界計算をGPUによって行い,電界部分についての計算時間の比較を行う。電界部分におけるカーネルの分割方法は,前章で述べたように7通りについて比較する。計算時間の計測にはOpenCLのプロファイリングAPIを用い,CPU-GPU間のメモリアクセス時間はグローバルメモリを用いるためどの分割方法も等しいとする。カーネル起動時間に関しては今回に限り無視し,純粋な計算時間だけの比較を行った。時間ステップ数を850,計算領域を $L \times L$ ,thread数を $x$ 方向 $y$ 方向共に16とし,初期条件として磁界は0,電界は半分の領域を0もう半分を1とした。計測時間は850ステップの平均を1ステップ分として算出した。また,ソフトウェア開発環境を表1に示した。また,計算に使用したGPUについて概要を表2に示した[3]。

表 3: 倍精度演算において OpenCL1.2 環境で GTX 980ti を使用したときにおける計算時間

計算領域 [ $L \times L$ ]	計算時間 [ $\mu s$ ]						
	kernel1	kernel2	kernel3	kernel4	kernel5	kernel6	kernel9
18×18	<b>12.922</b>	18.089	22.843	32.878	34.180	38.622	54.980
34×34	<b>14.414</b>	21.478	26.612	24.584	36.658	34.451	51.521
66×66	<b>14.428</b>	18.218	23.524	25.328	32.544	35.528	57.641
130×130	23.584	<b>21.553</b>	25.793	26.385	39.178	37.029	55.630
258×258	78.729	<b>35.650</b>	39.734	39.131	47.249	48.119	60.559
514×514	1,289.254	104.412	<b>100.708</b>	101.861	108.121	111.741	126.605
1026×1026	761.568	363.622	<b>362.473</b>	366.654	370.002	373.985	388.912
2050×2050	3,004.890	1,665.115	1,477.734	<b>1,462.744</b>	1,472.051	1,476.197	1,518.152
4098×4098	15,805.632	8,214.094	<b>7,573.825</b>	7,578.008	7,599.301	7,602.571	7,608.656
8194×8194	335,262.834	39,161.256	<b>35,991.476</b>	35,995.554	36,044.913	36,695.552	36,032.242

表 1: ソフトウェア開発環境

CPU	Intel Core i7-5820K 3.30GHz ×2
メモリ	64 GB (32GB×2)
GPU	NVIDIA GeForce GTX 980ti
OS	Ubuntu 16.04LTS
開発環境	OpenCL 1.2

表 2: GTX 980ti 概要

計算コア	2,816
ベースクロック	1,000MHz
メモリ量	6GB

#### 4.2 計測結果と考察

表 3 に各々のカーネル分割時における 1 ステップ分の電界計算にかかる計測時間を示した。表 3 より、カーネル 1 に対する、分割を行った際の高速化比は計算領域の大きさによってまばらな値となったが、計算領域 8194\*8194 では分割した場合の高速化比は 9.3 倍を超える高いものとなった。ここで高速比とは単独カーネルの時間を分割カーネルで最も速かった時間で割った比である。

また、表 3 において赤字で示した数字は同一計算領域で最速の時間である。この結果から計算量がある程度少ない場合にはカーネルを分割しないほうが高速となっていることがわかる。この原因として、GPGPU には Warp という概念があり、32 個のスレッドを 1 つのグループとして並列化する [4][5]。そのため計算量が少ない場合、カーネルを分割すると境界条件部分の計算量が少ないため、GPU のリソースを活かすことができずパフォーマンスが低下する。

一方、計算量が多い場合、カーネルを分割する方が計算が高速となっている事がわかる。ここで、カーネルを分割することによってプログラムが高速となったことについて考える。カーネルを分割すると、条件分岐 (if 文) がカーネル関数内に存在しなくなるため、ワーブダイバージェンスが発生しなくなった事が要因であると考え

られる。しかし計算量によってはカーネルの数が多すぎてもプログラムが高速にはならないことがわかった。

#### 5. まとめと今後の課題

本論文では、OpenCL を用いて GPGPU が苦手とする条件分岐処理を含む数値計算における高速化の検討を行った。GPU 上で起動するカーネル関数を複数用いることにより、1 つのカーネル関数内における条件分岐処理の影響を減らし、計算時間の低減を行えると仮定した。その結果、計算量が比較的多い場合では分割数を増やすことで計算時間を低減することができた。しかし、計算量が少ない場合にはカーネルを分割しないほうが高速となることがわかった。

今回、計算量に応じて最適なカーネル数が異なることがわかった。そのため今後の展望として、計算量に応じて最適なカーネルの分割を選択するといったアルゴリズムの開発を検討する。また、OpenCL を用いた GPGPU の特性をより研究し、更なる高速化を図っていきたい。

#### 謝辞

本研究は JSPS 科研費 15K00175, の助成を受けたものです。

#### 参考文献

- [1] John Cheng, Max Grossman, Ty McKercher, "CUDA C プロフェッショナルプログラミング", 株式会社インプレス, 2015
- [2] 宇野 亨, "FDTD 法による電磁界およびアンテナ解析", コロナ社, 1998
- [3] "NVIDIA Home Page", <http://www.nvidia.co.jp/object/geforce-gtx-980-ti-jp.html>
- [4] 伊藤 智義, "GPU プログラミング入門", 講談社, 2013
- [5] 株式会社フィクスターズ, "OpenCL 入門", インプレスジャパン, 2012