GPUクラスタにおける大規模都市気流シミュレーションの 最適化と性能モデル

高嵜 祐樹1 遠藤 敏夫1 松岡 聡1

概要:GPU 向けステンシル計算の規模は,通常 GPU のメモリ容量に制限されるが,テンポラルブロッキ ングと呼ばれる手法により性能劣化なく大規模化を実現可能である.本研究では,約 15,000 行のコード規 模を持つ GPU クラスタ向けアプリケーションである都市気流シミュレーションの大規模化・高性能維持 を実現する手法について述べる.プログラミングコストを抑えるために,自動で GPU メモリとホストメ モリ間でメモリスワップを行うランタイムライブラリである HHRT を用いる.その条件下で性能劣化を防 ぐための最適化を行うことにより,元プログラムの最大 85%の性能を達成した.さらに,性能予測モデル の構築により,性能に影響を与えるパラメータの絞り込みを可能にした.

Optimization and Performance Model for Large-scale Wind Simulation in Metropolitan on GPU Clusters

Yuki Takasaki¹ Toshio Endo¹ Satoshi Matsuoka¹

Abstract: The scale of stencil computation in GPU is usually limited by the memory size of GPU. Temporal blocking technique is used to break this limitation without performance degradation. Urban turbulence simulation code for GPU cluster, which consists of 15,000 lines, are modified to exploit the scale, keeping its performance. To reduce programming cost, HHRT runtime library is used to swap data between host memory and GPU memory automatically. With optimization to prevent performance degradation, we achieved up to 85% of the performance of the original program. We also revealed parameters affecting the performance by constructing performance model on memory swapping of HHRT.

1. はじめに

エクサスケール時代に向けた高性能計算システムにおけ る課題の一つにメモリウォール問題,つまりプロセッサの 計算速度に比べ,メモリのデータ転送速度と容量の増加 の進歩が遅いことがあげられる.近年,GPGPUやXeone Phiに代表されるメニーコアプロセッサの進歩により,ま すます対応の必要性は高まっており,その解決に向けてメ モリ階層を効率的に活用する技術が求められている.

本論文では, GPGPUを用いたスーパーコンピュータ上

東京工業大学 Tokyo Institute of Technology におけるステンシル計算を対象に取り上げ,問題規模の大 規模化と高性能性を目的とする.ステンシル計算は,流体 シミュレーションなどの分野の重要カーネルであり,一般 的には,シミュレート対象領域を規則格子で表し,各時間 ステップでは各格子点の値の更新を隣接する格子の値を 使って計算する.非常に並列性が高いこと,および必要と するメモリバンド幅性能が高いという性質から,GPU上の 実行に向いており,多くの成功例をもつ[1],[2].しかし, GPUに搭載されているデバイスメモリはホストメモリよ り小さく,GPU上のステンシル計算のほとんどにおいて, 対象とする問題規模はデバイスメモリサイズ内に限定され ていた.シミュレーションの分野において求められる,大 規模,高精度かつ高性能な計算を実現するには,デバイス メモリだけでなくホストメモリも加えたメモリ階層を効率 的に利用する必要がある.

本論文では,GPU 搭載スーパーコンピュータ向けに実 装された既存の流体アプリケーションを対象とし,MPI と CUDA を用いて記述された都市気流シミュレーション ソフトウェア [2] を取り上げ,大規模性・高性能性・高生 産性を実現する。第一段階として,対象ソフトウェアを HHRT と呼ばれるメモリスワップを行うランタイムライブ ラリ [4], [5] 上で動作させる.これにより,デバイスメモリ 容量を超えることができ,大規模性を達成する.次に性能 を高く維持するために,第二段階としてアルゴリズムのメ モリアクセス局所性を向上させる.そのためにテンポラル ブロッキング [6], [7], [8], [9] と呼ばれる手法を用いてソフ トウェアを書き換える.このとき,HHRT の実行モデルを 基にすることにより,書き換えコストを比較的小さくする ことができる.

我々はすでに,以上の手法が単純な7点ステンシルベン チマークにおいては適用可能であることを示してきた[5] が,大規模実アプリケーション(今回のアプリケーション は約15,000行から成る)に適用できるか明らかではなかっ た.そこで対象アプリケーションに我々の手法が適用可能 であるかの予備調査および,アプリケーションの性質を考 慮した最適化も重要であることを明らかにした.より具体 的には:

- 予備調査において,データ間の依存関係が隣接データ 間のみであり,各時間ステップにおけるリダクショ ンのような大域的な処理がないことの確認が必要で あった。
- 今回のアプリケーションは方向によって複数の境界条件の計算方法を用いており、それぞれが時間ブロッキングに対応可能か調査が必要であった。
- 今回のアプリケーションにおいて、プロセス間境界の MPI通信の方法がHHRT上での低速な動作を引き起こしており、修正が必要であった。

以上の手法により, GPU 搭載スーパーコンピュータで ある TSUBAME2.5 上での性能評価を通して,実際のステ ンシルアプリケーションに対して,大規模化,高性能化, 高生産性の3つを達成することを示す.

2. 背景

2.1 GPU 搭載マシンのメモリ階層

スーパーコンピュータの性能を向上させるために GPU アクセラレータや Intel Xeon Phi の利用が広がっており, 本論文では NVIDIA 社の GPU について述べる.このよう なシステムのメモリ階層は、従来のキャッシュ、ホストメモ リ、HDD からなるメモリ階層より、高階層となっている. ホストメモリは、数十から数百 GB の容量と数十 GB/s の メモリバンド幅であるのに対し、デバイスメモリは、容量 6GB~12GB とメモリバンド幅 250GB/s と、小容量で高 速なメモリである傾向にある.デバイスメモリ容量を超え るためには、ホストメモリのメモリバンド幅より更に小さ い8GB/s という PCI-Express バスを通したデータ転送を 効率よく行うことや転送量を削減することが重要となる.

2.2 ステンシル計算

ステンシル計算は、流体シミュレーション分野などで、 使われる一般的なカーネルである.シミュレートされる領 域を規則格子で表し、時間経過による各格子点の値の変化 を計算する.各格子点の値の更新には、隣接する格子の値 を使って計算される.各格子点の計算は、独立に計算する ことが出来るため、非常に並列性が高い計算となっている. また,ダブルバッファリング手法が良く使われ,本論文で もそれを仮定する.

GPU クラスタ上での 7 点ステンシル計算を例にとり,プロセス間の通信を MPI, GPU のプログラミングを CUDAで記述された実装の例を図 1(a) に示す.これは GPU メモリを超えられない単純な例である.GPU のメモリ容量を超えるデータに対応した実装例を図 1(b) に示す.大容量のデータを扱うために、領域を部分領域に分割し、部分領域を入れ替えながら計算して行く.MPI 通信に加え、2 回の CPU-GPU 間のデータ転送を全部分領域にわたって毎ステップ行わなければならないので、性能が大幅に低下してしまうという問題がある.



図 1 (a) 通常の MPI+CUDA、(b) 単純な大容量化に対応した アルゴリズム.

2.3 テンポラルブロッキング

テンポラルブロッキングは、ステンシル計算のメモリ局 所性を向上させるための手法である [6], [9].ここでは計 算対象の配列を小さい部分空間に区切り,その部分空間に ついて,時間ステップ計算を他と独立に複数回進める.こ れにより,時間ステップ毎に配列全体をスキャンする通 常の方法よりも局所性が良好である.もともとはキャッ シュヒット率を向上させるために提案された手法である が,GPUのデバイスメモリを超えるステンシル計算の隣 接通信で必要な CPU-GPU 間の通信回数を減らすために も用いられている [12],[13].

テンポラルブロッキングを用いて図 1(b) を改良すると図 2 のようなアルゴリズムとなる.部分領域ループの中に更 に内部時間ループが追加されていることがわかる.時間ブ ロッキングする時間幅をテンポラルブロッキングサイズ kとし,時間発展する回数を t とおくと,各部分領域が順番に 内側の時間ループで k 回計算するので,外側の時間ループ は t/k 回の計算となる.図 1(b) と比べると,PCI-Express 通信を 1/k 回に削減している.一回あたりの転送量はほぼ 変わらない(どちらも部分空間全体 + 袖領域)ので,計算 全体の PCI-Express 通信量もほぼ 1/k にできる.

しかしながら図 1(a) のような既存アプリケーションが 存在するときに,図2のように書き換えるのはプログラミ ングコストが高く,これを改良するために,次に述べるメ モリスワップランタイム HHRT を用いる.



図 2 大規模対応テンポラルブロッキングのアルゴリズム

2.4 HHRT

HHRT (Hybrid Hierarchical Runtime) [4], [5] は, メモ リ階層間のメモリスワップをサポートするランタイムライ ブラリで,メモリスワップを必要とするアプリケーション のプログラミングコストを抑えることを目的としている. HHRT は,現在,CUDA と MPI でコーディングされたア プリケーションを対象に,GPU-CPU 間のメモリスワップ をサポートするランタイムライブラリである.

HHRT を利用した場合の特徴は, 複数の MPI プロセス が単一 GPU の計算リソースを共有していることであり, これらのプロセス間合計ではデバイスメモリ容量を超え るデータを扱うことができる.なお,合計利用メモリ量が デバイスメモリ容量を超える場合には,プロセス単位での (OSのようにページ単位ではない)スワッピングを行う. 詳細は過去の文献を参照されたい[4],[5].

この HHRT 単体はデータの追い出しを自動化する一方, テンポラルブロッキングのようなアルゴリズムの改良を行 うものではない.そのためその組み込みは依然ユーザの責 任であるものの,図2よりも単純に,図3のように記述す ることができる[5].



図 3 HHRT とテンポラルブロッキングを組み合わせた アルゴリズム

3. 都市気流シミュレーション

今回大規模化に使用する小野寺らが開発した都市気流 シミュレーション [2] について説明する.これは,東京の 10Km 四方のエリアを対象に,実際の建造物のデータを 使い 1m 間隔の格子解像度で気流シミュレーションを行 い,高層ビルの背後に発生する渦によるビル風や幹線道 路に沿って流れる風の道などを再現する実ステンシルア プリケーションである.格子ボルツマン法(LBM:Lattice boltzmann methd)とよばれる手法が使われている.LBM は,連続対として記述された流体に対して,離散化された 空間格子状において,並進と衝突をする仮想的な粒子の集 合を速度分布関数として仮定し,格子状の粒子の速度分布 関数について時間発展を解いて行く手法である.19方向も しくは 27 方向の速度を,それぞれ格子データで表現する のが特徴の一つである.

7 点ステンシルプログラムと都市気流シミュレーション の違いは,以下のような点が挙げられる.

7 点ステンシルは,格子点をダブルバッファリングのための2つの3次元配列で表すのに対し,LBMは3次元配列を19ないし27個,さらにダブルバッファリ

ングのために2倍の個数用いる.

- 都市気流シミュレーションはまた 16 個の物理量のための三次元配列を必要とする.
- 7 点ステンシルは1つのカーネル関数のみ用いるのに対し,都市気流シミュレーションは各ステップあたり6つのカーネルを順次計算する.
- 境界条件は、7点ステンシルは単純なディレクレ条件を用いる(境界値は固定)のに対し、都市気流シミュレーションは、X軸を周期的境界条件、YとZ軸はノイマン境界条件を使って計算する。

以上のように,計算に必要なデータ数,1ステップあたりの計算,境界条件が主な違いとなっており,コードの行数 も約15,000行となる大規模なものである.

 大規模・高性能・高生産性の実アプリケー ションにおける実現

都市気流シミュレーションの提案手法である HHRT と テンポラルブロッキングの導入について説明する.本アプ リケーションの既存実装は,前章に述べたようにベンチ マークより複雑であるものの,基本的には図1(a)のような 構成を持っている.これを,図3のようにソフトウェアを 改造し,HHRT 上で実行を行う.

4.1 HHRT の導入

第一段階として,既存実装をほぼそのまま HHRT 上で コンパイル・実行可能とするために2つの小さなコード変 更を行う.

- hhrt.h のインクルードの追加
- HH_initHeap() 関数による各プロセスの最大利用メモ リ量の宣言

上記の変更を行い, HHRT ライブラリとリンクすることに より、GPU メモリサイズを超えた問題サイズにおいても プログラムを実行することが可能となる.

HHRT を使わない場合には,大規模問題サイズの対応 のためには,配列の部分空間への分割および,明示的な CPU-GPU 間の転送の記述などのために,図1(b)のよう なプログラム全体に係るループ構造の変更が必要となって しまう.これに対し,HHRTの利用により,上記のような コード内の局所的・機械的な変更により,プログラミング コストを抑えつつ,同様の効果を得ることができる.

4.2 テンポラルブロッキングにむけた予備調査と実装

テンポラルブロッキングは,アクセス局所性を向上させ るアルゴリズムであるが,どんなステンシル計算にも適用 可能なわけではない.テンポラルブロッキングは,ある時 刻における点 x の計算が,「直前の時刻の,点 x の近傍の値 にのみ依存」するような計算の場合に適用可能である.こ れに反する例としては,共役勾配法のように,各ステップ において領域全体に係る内積が必要な場合が挙げられる. また,各時間ステップにおいて広範囲の平均操作などが必要な場合も障害となる.

今回の都市気流シミュレーションにおいて,上記の条件 を満たすアルゴリズムであるか否かについて,カーネル関 数を中心としたコード精査により調査した.その結果,各 点の計算は近傍にのみ依存することが確認でき,テンポラ ルブロッキングは適用可能であると判断した.

上記の判断のもと,都市気流シミュレーションについて, 図 1(a)の構成から図 3 の構成への変更を手作業で行った. ここで,HHRTを利用することにより,コード変更は比較 的局所的なもので済んだ,具体的には,下記のような変更 である.

- 時間ループの開始箇所と終了箇所を変更し,時間ループの二重化を行った.
- 内側時間ループの進行にしたがって,計算対象に含める袖領域が一段ずつ狭まっていくような,空間ループ 開始・終了点の調整を行った。
- MPIによる境界通信を,外側時間ループと内側時間 ループの間に移動し,さらに k 段の境界を一気に通信 するよう変更した.
- 境界条件については、ノイマン条件の計算は内側時間 ループの中にとどめたのに対して、周期境界条件のた めの MPI 通信を、前項目と同様に、外側時間ループ と内側時間ループの間で行った。

この結果,主な変更は各ループの開始箇所と終了箇所に 集中し,実アプリケーションにおいてコードを長くする主 要因であるループ内の計算ロジックについては,変更は不 要であった.さらに,コードの移動も上記3番目,4番目 の点に抑えられた.その結果,テンポラルブロッキング導 入におけるコード変更量は,全体コード量約15,000行に 対して,約300行に抑えられた(元プログラムソースとの diffに基づく計量).

 4.3 HHRT のスワップ動作に合わせた MPI 通信の最 適化

上記のステップまでで,基本的には大規模・高性能・高生 産性は実現されたが,ここではさらなる性能最適化につい て述べる.詳細は評価の章に譲るが,上記の基本的な実装 を GPU 搭載スーパーコンピュータで実行したところ,予 想よりも HHRT 内部のスワップ処理回数が3倍程度多く なってしまい,そのため性能が期待よりも大きく低下して いた.調査の結果,元プログラムの実装の性質と,HHRT の性質の相性が悪く,以下のような問題が起こっているこ とがわかった.

HHRT のスワップ処理は、ブロッキング通信命令が呼ばれたときに発生するため,複数回ブロッキング通信が呼ばれるアプリケーションでは,この複数回のスワップ処理

が性能低下の原因となってしまう.都市気流シミュレー ションでは、全体空間は3次元にプロセス間分割されてお り,19方向ないし27方向の近傍値を用いるためには,各 プロセスは隣接する26プロセスとの境界通信を必要とす る.元プログラムでは,X,Y,Z方向の通信を順番に行 う3回のブロッキング通信を行っており,MPI通信のた めだけに5回のスワップ処理を行うような実装となってい た.通常の(HHRTを使わない)実行ではこれは大きな問 題とはならなかったが,HHRTの「MPI通信がブロックす る箇所では,プロセスがスワップアウトしうる」という性 質により性能上の問題が発生した.つまり,1度スワップ アウトすれば十分なところ,3回のスワップ処理となり, PCI-Express通信量を5倍に押し上げてしまっている.

これを改善するために,MPI-Waitall1回にまとめるようなコード変更を行った.つまり,26個の全ての隣接プロ セスに対して同時にノンブロッキング通信を起動すること とした.各隣接プロセスに対して MPI-Isend と MPI-Irecv が起こるため,52個の通信となる.その後,すべての通信 を1回の MPI-Waitall で待つように変更を行った.この最 適化にともなうコード変更量は約800行であった.前節の テンポラルブロッキング自身の組み込みよりも大きくなっ たが,これは三次元の隣接通信を行う元プログラムの通信 部分がもともとコードがやや複雑だったためである.

4.4 HHRT のスワップデータ量削減による最適化

HHRT のスワップ処理は,スワップ処理時に GPU 上に 確保したデータを保護するために,ホストメモリヘデータ を退避する.しかし,ステンシル計算において,次ステッ プで全てのデータを保持する必要がないことが調査の結果 わかった.このため,スワップ処理では,無駄なデータ転 送を行っているため,無駄なデータ転送を削減することで, 通信時間を減らすことができる.

そこで,HHRTには,HH_madviseというAPIがあり, スワップイン・スワップアウト時のデータ転送において, 転送しないデータを選択することで,そのデータを転送し ないようにすることができる.

HH_madvise 関数は,配列の先頭アドレス,データサ イズ,データの状態を引数として渡すことで,そのデー タをスワップイン・スワップアウトするか選択すること ができる.データの状態はHHMADV_NORMALとHH-MADV_CANDISCARDの2状態があり,それぞれスワッ プイン・スワップアウト時に、配列データの値を保証する かしないか,つまり,スワップイン・スワップアウト時に データ転送をするかしないかを表している.このようにし て,値の保証が必要ないデータをスワップイン・スワップ アウト時に転送せず,転送時間を削減することができる. 転送不要なデータとして,ステンシル計算の場合,ダブ ル・バッファリングの一方は,値の保証の必要がないため, 一方の転送を無視することができる.都市気流シミュレー ションなどの実アプリケーションの場合,配列数が多く バッファの一方だけでは性能向上率は少ない.しかし,次 ステップで使われる値を格納しない配列があるアプリケー ションの場合,計算ステップ終了後は、値を保持する必要 がないため,スワップイン・スワップアウト時にデータ転 送の必要がなくなるため,転送時間の削減ができる.

都市気流シミュレーションでは,このような配列が複数 あったため,HH_madviseによる最適化を行うことができ る.1つの配列データに対するこの最適化に伴うコード変 更は最低2行である.

5. 性能モデルの構築

提案手法を導入したアプリケーションを実行するために は、1GPU に処理させる問題サイズに対して、領域の分割 数、つまり 1GPU を共有する MPI プロセス数と最適なブ ロッキング段数を指定する必要がある.このため、最適な MPI プロセス数とブロッキング段数を推定するための性能 モデルを構築する.

MPI+CUDA のアプリケーションは,主に cudaMemcpy の HostToDevice と DeviceToHost, そして Kernel 計算と MPI 通信の4つから構成されている.今回のモデル構築 では,MPI 通信は,HHRT のメモリスワップ処理の実行 時間で隠蔽できると仮定して構築する.そのため,Host-ToDevice, DeviceToHost,Kernel 計算の3つからなる性 能モデルを構築する.

モデル構築にあたり各実行時間は次のようにあらかじめ 算出しておく必要がある.cudaMemcopyは,HostToDeviceとDeviceToHostのそれぞれの実行バンド幅を測定し, 通信するデータサイズにより実行時間を算出する.同様 に,カーネル計算も実行Flopsを測定し,そのflop数から 算出している.ただし,MPI通信のためのバッファコピー は,実行バンド幅とデータサイズから算出する.

提案手法と通信の最適化を行った都市気流シミュレーションの動作モデルを図4から図7に示す.各図は,1GPU を共有している MPI プロセス数 n,1GPU が同時にメモ リ資源を共有出来るプロセス数 m,ブロッキング段数 k と したときの各処理実行時間を表している,また,各図の長 方形が,各時間を表しており,それぞれの色はスワップイ ン(赤),スワップアウト(青),計算(オレンジ),実行待 ち(緑)を表し,P0~P3は MPI プロセスのランクを表し ている.

また,メモリ以外の各ハードウェア資源は複数プロセス で同時に占有することは出来ないため、競合が起こった 場合に実行待ちとなっている.ただし,スワップインとス ワップアウト処理は,同時に実行することが出来る.

図4から図7の場合分けに対応したモデル式を以下に

示す .

$$T_{Iter}(1,n,k) = n(T_{INk} + T_{CALk} + T_{OUTk})$$
(1)

$$T_{Iter}(2, n, k) = \frac{n}{2} (T_{INk} + T_{CALk} + T_{OUTk})$$
(2)

$$T_{Iter}(2, n, k) = n \max\{T_{INk}, T_{CALk}, T_{OUTk}\}$$
(3)

$$T_{Iter}(m, n, k) = n \max\{T_{INk}, T_{CALk}, T_{OUTk}\}$$
(4)

式 $T_{Iter}(m,n,k)$ は,各図の黒矢印間の時間を表している.また, T_{INk} , T_{OUTk} , T_{CALk} は,それぞれプロッキング段数 k の時のスワップイン,スワップアウトの転送時間とカーネル計算時間を表しており,スワップイン,スワップアウトの転送時間には,MPI通信に必要なバッファ領域の CPU-GPU 間の転送時間を,カーネル計算時間には,MPI 通信に必要なバッファと配列間のコピー時間を含んでいる.

各式はそれぞれ m によって場合分けされており,それそれ (1)m = 1, (2)(3)m = 2, $(4)m \ge 3$ の3つに場合分けされる.m = 2の時に,各実行時間の内に $T_a \ge T_b + T_c$ となるような T_a が「(2)存在しない,(3)存在する」という条件によって式が異なる.ここで, T_a , T_b , T_c には, T_{INk} , T_{OUTk} , T_{CALk} が入る.





図 5 (2)m = 2の実行モデル



図 0 (3)m = 2 の美行モデル



6. 性能評価

6.1 評価環境

今回の性能評価のために,東京工業大学学術国際情報 センターのスーパーコンピュータ TSUBAME2.5 の Thin ノードを使って測定した.Thin ノードの構成は,Intel Xeon 5670 2.93GHz (6 cores)プロセッサを 2 ソケット, NVIDIA Tesla K20X を 3 基搭載しており,ホストメモリ の容量 54GB (一部 96GB)となっている.各計算ノード 間は,Infiniband QDRx4 によって接続されている.今回 の評価実験では,1ノードあたり1GPUを使用し,1ノー ドでの実験はホストメモリ 96GB のノードで実験し,複数 ノードを使用した実験はホストメモリ 54GB のノードで 実験を行った.また,開発環境は,OS が SUSE Linux 11 sp1,コンパイラが gcc 4.3.4,MPI が OpenMPI 1.6.5,そ して CUDA は 6.0 を使用した.

6.2 1GPU での評価

GPU のデバイスメモリを超えるデータ処理の性能を評価するために,ホストメモリ容量 96GB の Thin ノードを使用した.比較対象のプログラムの詳細を以下に示す.

- NORMAL:提案手法を一切使わない元プログラム
- HH:NORMAL に HHRT を導入したプログラム
- HH_TB:HH にテンポラル・ブロッキングを導入した プログラム
- HHTB_MPI:HH_TB を MPI 通信の最適化したプログ ラム
- HHTBMPI_OPT:HHTB_MPI に HH_madvise の最適 化したプログラム

測定対象となる各プログラムを,問題サイズを変化させた時の性能の推移を図8に示す.図にプロットされた性能は,最適なプロッキング段数と分割数の時の性能を使用した.

実験で利用した GPU K20X は,デバイスメモリ容量が 6GB であるため,NORMAL は,6GB 以上の問題サイズ を実行することができないことがわかる.

NORMAL に HHRT を導入した HH は,デバイスメモ リ容量である 6GB の制限を超え,約 48GB の問題サイズ までの実行を可能にしている.しかし,6GB を超えると NORMAL の 5%程度まで性能が低下している.これは、1 ステップ毎にスワップ処理が発生してしまうためである.

HH のスワップコストを削減するためにテンポラルブ ロッキングを導入した HH_TB は, HH に比べ最大 4.9 倍 の性能向上し, NORMAL の最大 27%の性能となった.し かし, HH_TB のスワップ回数を調査したところ,1 ステッ プあたり3回のスワップを行っていることがわかった. HH_TB のコードは、1 ステップあたり3回の MPI_Wait 命令がよばれており,このため, HHRT のスワップ条件に よって,3回のスワップを行い性能が低下していた.

そこで、HH_TB の MPI のブロッキング通信を1回にま とめる最適化した HHTB_MPI は, HH_TB に比べ最大 2.4 倍の性能向上し, NORMAL の最大 64%を達成した.この 結果から, 複数回の MPI 通信を行うアプリケーションに 対して, MPI 通信をまとめる最適化が必要であることがわ かった.

HHRT の実行モデルでは, MPI プロセスのスワップ処理 を行う必要があるため, 計算より速度の遅い PCI-Express 通信が実行時間の大半を占める.そこで, PCI-Express の 通信時間を減らすために, 通信データ量を削減する最適化 をした HHTBMPI_OPT は, HHTB_MPI に比べ最大 1.4 倍性能向上し, NORMAL に対して最大 85%の性能を達成 した.

容量の増加と伴に性能が低下している原因は,HHRT の MPI プロセスの管理コストとホストメモリ容量である. HHRT の MPI プロセスの管理コストは,複数の MPI プロ セスが 1GPU を共有するのに必要なコストで,1MPI プロ セスあたり約 72MB の管理領域を GPU メモリ上に確保す るため,分割数が多くなると GPU メモリを圧迫してしま う.ホストメモリ容量による制限は,HHRT は,ホストメ モリの半分の容量までの問題サイズしか扱えないためであ る.どちらの場合も,ブロッキング段数を多くするとメモ リ容量が大きくなるため,ブロッキング段数を増やすこと による性能向上ができないことが原因である.

以上の結果から,提案手法を導入したアプリケーションは,MPI通信の最適化とHint機能による最適化によって, 性能劣化を最大15%程度にまで防ぐことが可能であること がわかる.



図 8 1GPU での実行結果

6.3 ブロッキング段数の評価

図 9 は、ブロッキング段数を変化させたときの HHTB_MPIと HHTBMPI_OPT の性能の変化を表した グラフである.HHTBMPI_OPT は HH_madvise によって スワップ処理時のスワップデータサイズが削減されたこと で、最適なブロッキング段数が HHTB_MPIと異なり、ま た、HHTB_MPIより少ないブロッキング段数で、最適値 となっていることがわかる.また,最適値以降性能が低下 している理由は,ブロッキング段数を増やすことによる冗 長計算の増加によって,計算コストが,スワップ時間など の通信コストを上回ったためである.



図 9 ブロッキング段数による影響

6.4 複数ノードでの評価

図 10 は、1 ノードあたり 1GPU を使用して,ノード数 を変化させていった場合のウィークスケーリング性能を表 している.問題設定は,NORMAL が 5.5GB で,OPT は HHTBMPLOPT と同じプログラムで,OPT_1,2,3 は問題サ イズの違いを表しており,それぞれ 5.5GB,10.8GB,16.2GB である.NORMAL と OPT_1 を比較すると,HHRT や TB の 2 重ループのコストによって OPT_1 は,数%の性能劣 化しているが,ほほ同等の性能で,スケーラビリティの劣 化が無いことがわかる.

全体で比較すると, OPT_2, OPT_3は, ほぼ同等の性能 で, スケーラビリティも高く, NORMALの約80%程度の 性能を維持していることから,提案手法は, 複数ノードを 使用した大規模環境にも対応していることがわかる.



図 10 複数ノードでの実行結果

6.5 性能モデルの評価

提案手法を適応したアプリケーションの最適な性能を 予測する性能モデルの評価を行った.対象プログラムは、 HHTBMPI_HINT で,問題サイズは,10.8GB に設定し, 分割数とブロッキング段数を変化させたときの性能モデ ルの予測値と HHTBMPI_HINT の実測値を比較し,その 結果を図 11 から図 14 に示す.各図はそれぞれプロセス 数 4,8,16,32 の時の 1 ステップあたりの実行時間を表 しており, 横軸はブロッキング段数、縦軸に実行時間を とっている.性能モデルの 1 ステップあたりの予測値は, $T_{step}(m,n,k) = T_{Iter}(m,n,k)/k$ として算出した.予測値 と実測値の差はあるものの,ブロッキング段数の変化に伴 う実行時間の変化をかなり表現できていると考えられる. ただし,4プロセスにおいて,ブロッキング段数 1 と 2 で, 大きくずれているが,これは,MPI 通信時間を考慮してい ないことが,原因であると考えられる.

実測値の最適解は, (n,k) = (16,8)のときで,これは予 測値の最適解に一致している.以上の結果から,構築した 性能モデルは,パラメータの絞り込みが可能であると考え られる.













図 14 32 プロセスの場合

7. 関連研究

ステンシル計算は連続体のシミュレーションにおいて重要なカーネルであり,多数のアプリケーションがスパコン 上で実装されてきた[1],[2],[3].

ステンシル計算の最適化技法についても,空間ブロッキ ング,テンポラル・ブロッキングなど数多く提案・評価され てきたが,以降では後者に注目する.時間方向にまたがっ た計算順序の変更により,メモリアクセス局所性を改良する テンポラルブロッキング手法は,当初は主に CPU のキャッ シュ利用効率の向上のために提案された[6],[7],[8],[9]. また GPU においてデバイスメモリとオンチップ共有メモ リ間のデータ移動を削減する研究もなされている[10].本 研究と同様に,GPU デバイスメモリの容量制限を超える ためにデバイスメモリとホストメモリ間のデータ移動を削 減する研究としては Mattes らの FDTD の実装[11]が挙げ られる.さらに我々の研究グループでは,多数 GPU・多 数ホスト環境におけるテンポラルブロッキングの採用によ り,大規模・高性能ステンシル計算が可能であることを実 証した[12],[13].

しかしながら上記の研究においては,プログラミングコ ストの軽減についての議論が不足していると考えられる. プログラマが最適化を行う際には,時間ループ・空間ルー プ双方の大幅な書き変えが必要であり,それは特に複雑な 構造を持つ実アプリケーションにおいて困難となると考え る.それに対して我々は,MPIおよび CUDA で記述され た実アプリケーションがすでに存在するという前提のも と,原則的に時間ループ部分の書き換えと,HHRT ライブ ラリ上での実行により,デバイスメモリ容量超えが可能で あることを実証した.

ステンシル計算の記述を容易にすることを目指したアプ ローチとして, Physis[14] などのドメイン固有言語 (DSL) を用いることが考えられる.このアプローチでは,アプリ ケーションが DSL で記述されていることを前提にしてお り,すでに並列化されたアプリケーションを仮定する本研 究とは異なる.一方で,DSL によるアプローチでは、テン ポラルブロッキングを含めた最適化をシステム側で行い, ユーザから隠ぺい可能であると期待される.

8. まとめと今後の課題

本研究は,ステンシル計算を対象に,高性能,大容量, 低プログラミングコストの3つを実現するための提案手法 として,テンポラルブロッキングとHHRTを適用すること 提案し,また,その最適化を行った.提案手法の導入によ り,GPUのメモリ容量を超えるデータを処理することを 可能にし,元プログラムの最大85%の実行を可能にした. これにより,本提案手法は,GPU対象の実ステンシルア プリケーションの大規模化のための最適手法の一つである ことを示した.また,性能モデルの構築により,提案手法 の導入によって現れるパラメータの最適解の絞り込みを可 能にした.

今後の課題として,問題サイズ増加に伴う性能劣化を防 ぐために,HHRTのMPIプロセスの管理方法の改善と,更 なる大規模化のためのSSD,HDDへの対応が必要である.

謝辞 本研究は JST-CREST の研究課題「ポストペタス ケール時代のメモリ階層の深化に対応するソフトウェア技 術の深化に対応するソフトウェア技術」および科学研究費 助成事業(基盤研究(S)23220003)の支援によります.

参考文献

- [1] Takashi Shimokawabe, Takayuki Aoki, Tomohiro Takaki, Akinori Yamanaka, Akira Nukada, Toshio Endo, Naoya Maruyama, Satoshi Matsuoka: Peta-scale Phase-Field Simulation for Dendritic Solidification on the TSUB-AME 2.0 Supercomputer. In Proceedings of IEEE/ACM Supercomputing (SC11), 11pages (2011).
- [2] 小野寺直幸,青木尊之,下川辺隆史,小林宏充:格子ボル ツマン法による 1m 格子を用いた都市部 10km 四方の大 規模 LES 気流シミュレーション.情報処理学会ハイパ フォーマンスコンピューティングと計算科学シンポジウ ム (HPCS2013), pp. 123–131 (2013).
- [3] Ryoji Takaki, Kazuomi Yamamoto, Takashi Yamane, Shunji Enomoto, Junichi Mukai: The Development of the UPACS CFD Environment. In proceedings of ISHPC 2003, pp. 307-319 (2003).
- [4] 遠藤敏夫: 並列プログラムをメモリ階層利用可能とするランタイム.情報処理学会研究報告 2013-HPC-140 No 43, 8pages (2013).
- [5] T. Endo and G. Jin: Software Technologies Coping with Memory Hierarchy of GPGPU Clusters for Stencil Computations. In Proceedings of IEEE Cluster Computing (CLUSTER2014), 8pages (2014).
- [6] Michael E. Wolf and Monica S. Lam: A Data Locality Optimizing Algorithm. In proceedings of ACM PLDI 91, pp. 30–44 (1991).
- [7] David Wonnacott: Using Time Skewing to Eliminate Idle Time due to Memory Bandwidth and Network Limitations. In proceedings of IEEE IPDPS 2000, pp. 171–180 (2000).
- [8] M. Wittmann, G. Hager, and G. Wellein: Multicoreaware parallel temporal blocking of stencil codes for shared and distributed memory. Workshop on Large-

Scale Parallel Processing (LSPP10), in conjunction with IEEE IPDPS2010, 7pages (2010).

- [9] T. Minami, M. Hibino, T. Hiraishi, T. Iwashita and H. Nakashima: Automatic Parameter Tuning of Three-Dimensional Tiled FDTD Kernel. In Proceedings of The Ninth International Workshop on Automatic Performance Tuning (iWAPT2014), 8pages (2014).
- [10] Anthony Nguyen, Nadathur Satish, Jatin Chhugani, Changkyu Kim, and Pradeep Dubey: 3.5-D blocking optimization for stencil computations on modern CPUs and GPUs. In Proceedings of IEEE/IEEE Supercomputing (SC10), 13pages (2010).
- [11] L. Mattes and S. Kofuji: Overcoming the GPU memory limitation on FDTD through the use of overlapping subgrids. International Conference on Microwave and Millimeter Wave Technology (ICMMT), pp.1536– 1539 (2010).
- [12] G. Jin, T. Endo and S. Matsuoka: A Multi-level Optimization Method for Stencil Computation on the Domain that is Bigger than Memory Capacity of GPU. AsHES workshop, held with IEEE IPDPS2013, 8pages (2013).
- [13] G. Jin, T. Endo and S. Matsuoka: A Parallel Optimization Method for Stencil Computation on the Domain that is Bigger than Memory Capacity of GPUs. In Proceedings of IEEE Cluster Computing (CLUSTER2013), 8pages (2013).
- [14] N. Maruyama, T. Nomura, K. Sato, and S. Matsuoka: Physis: An Implicitly Parallel Programming Model for Stencil Computations on Large-Scale GPU-Accelerated Supercomputers, IEEE/ACM Supercomputing (SC11), 12pages, Seattle (2011).