動的領域分割を用いた流体構造連成による サスペンション・フローの大規模 GPU 計算

都築 怜理^{1,a)} 青木 尊之^{1,b)}

概要:サスペンション・フローは流体中に多数の物体を含み,流体構造連成が支配的な複雑な流れである. 粒子法による計算を行うために物体を小さい粒子の集合で表現し,物体間の衝突には個別要素法(DEM)を 用い,流体計算には改良型 SPH 法を導入する.物体構成粒子は流体粒子や異なる物体の構成粒子間で相互 作用する.それらの受ける力とトルクを合算して物体の運動方程式の時間積分を行うが,多数の構成粒子 からの総和計算を伴うために計算はかなり複雑になる.GPU による計算において Linked-list を用いた近傍 粒子探索のデータアクセスを高速化するために,定期的な GPU のデバイス・メモリ上の粒子データのソー トが有効であることを確認した.また,力とトルクの総和計算を行う際,物体数に対して CUDA のスレッ ド並列を行う実装が高速であることも分かった.複数 GPU を用いるためにスライスグリッド法による動 的領域分割を行い,256 個の GPU により 8,743 万個の流体粒子と230 万個の物体構成粒子(物体としては 2,304 個)を用いたサスペンション・フローの大規模シミュレーションを実現した.

キーワード:流体-構造連成計算,粒子法,動的負荷分散,GPU コンピューティング

A Large-scale Suspension-flow Simulation directly solving fluid-structure interactions with Dynamic Load Balance on a GPU Supercomputer

Tsuzuki Satori^{1,a)} Aoki Takayuki^{1,b)}

Abstract: Suspension flows including fluid-structure interactions is one of challenging topics in fluid dynamics and many engineering applications. A large-scale particle-based method is a promising approach to carry out the numerical simulation for suspension flows. We have proposed an effective combination of the SPH method for fluids with the DEM (Discrete Element Method) collision for objects on a multi-GPU system. It is found that the sorting of particle data for the neighboring particle list using linked-list method improves the memory access greatly with a certain interval. We compare the two GPU implementations for the reduction of forces and torques applied to the particles constructing the objects. It is successful to apply the dynamic domain decomposition based on the 2-dimensional slice-grid method to our suspension flow simulation on multi-GPUs. A large-scale suspension flow of a tsunami interacting with a lot of objects with 87.3M particles including 2,304 cubic-shaped objects are successfully carried out on TSUBAME 2.5 at Tokyo Institute of Technology. The weak scalability is also measured from 4 GPUs to 256 GPUs.

Keywords: Fluid-Structure Interactions, Particle Method, Dynamic Load Balance, Multi-GPU computing

1. 序論

科学技術計算における粒子法とは、空間を任意に移動で

^{a)} tsuzuki@sim.gsic.titech.ac.jp

^{b)} taoki@gsic.titech.ac.jp

きる(ある程度の制限がある場合も含めて)点上で物理量を 時間積分する手法を指すことが多く,惑星軌道計算のよう な単一粒子計算よりも多数の粒子が相互作用しながら集団 運動する系の計算が中心となっている。粒子は原子や分子 と1対1対応する場合もあれば、プラズマ電子・イオンの 数万個を代表する場合もあるし、単に連続体を空間離散化

東京工業大学 Tokyo Institute of Technology, O-okayama 2-12-1, Meguro-ku, Tokyo 152–8550, Japan

するための代表点である場合もある.

ハイパフォーマンス・コンピューティングの分野では,重 力多体問題, 分子動力学計算, SPH(Smoothed Particle Hydrodynamics) 法による流体計算, 個別要素法 (DEM: Discrete Element Method) による粉体計算,メッシュフリー法によ る構造解析等があり、粒子間の相互作用や連結情報などは 大きく異なる。重力多体問題や分子動力学計算では相互作 用レンジが長く、N体問題に代表されるように相互作用の 計算負荷が支配的となるが、粉体計算では接触している粒 子間のみ反発や摩擦の相互作用が生じ、メモリ参照の比重 が高い. SPH 法などの粒子法による流体計算は、カーネル 半径と呼ばれる範囲内の 100 個程度の粒子と相互作用する ため、その中間的な計算負荷とメモリ参照と言える。 実際 の系は粒子数が非常に多く、数値計算ではできる限り多く の粒子を用いることで高い精度や計算できる現象が変わる ため, ACM ゴードンベル賞 [1][2] にも度々登場する大規 模計算が求められる分野である.近年,演算性能・メモリ 性能・電力比効率などの観点から、世界トップクラスのス パコンにおいて GPU がアクセラレータとして搭載されて いる. それらの GPU スパコンでのアプリケーション開発 の観点からは、CUDA プログラミングなどが必要であるば かりでなく、GPU ボード上のデバイス・メモリ量の制限や On-Chip の高速な共有メモリなどの階層的メモリ構造をう まく使う必要がある。また、GPUのデバイス・メモリ間の データ通信はホスト計算機を介する必要があり、通信時間 が大規模計算の大きなオーバヘッドになる可能性がある.

複数 GPU を使う大規模な粉体計算は動的負荷分散の手 法により達成されている [3]. これにより高速化を目指す レーザープリンターのトナー粒子の挙動解析や,各種の化 学工学,錠剤製造プロセスにおいて,紛体シミュレーショ ンが可能になりつつある. SPH 法等の粒子法による流体計 算も,京コンピュータで ParMETIS を用いた領域分割によ る大規模計算 [4] や,複数 GPU の 1 次元の動的領域分割に よる大規模計算 [5] が報告されていて,本著者のグループ でも大規模な GPU 計算 [6] を行っている.

一方,流体と構造が共存する現象は実際に非常に多く, さまざまな分野で重要性が広く認識されている.東日本大 震災の際,無数の瓦礫を伴った津波が押し寄せてくる映像 は印象的であり,それらの衝突を考慮した津波被害シミュ レーションの重要性は容易に理解できる.地盤工学におけ る液状化現象も砂粒と流体の連成問題であるが,実スケー ルでシミュレーションされた例はない.土石流は水と土砂 を含んだ流れであり,斜面災害のより先進的なシミュレー ションを行うためには必須である.プラントの配管内の固 体を多数含んだ流れの解析は化学工学において大きなテー マである.これらはサスペンション・フローと呼ばれ,大 規模・高精度シミュレーションが殆ど行われていない.流 体に対する格子を用いた計算は精度の点で有利であるが, 移動する複雑形状の物体を扱うのは苦手である.一方、粒 子法は複雑形状や流体構造連成計算を容易に行うことがで きる利点を持っている.

本研究は多数の球形粒子の集合で表現する複雑形状の構造物と SPH の流体との大規模連成計算を GPU スパコンで 実現させることを目的とする.連結していない球形粒子を 用いた DEM 計算は近接相互作用だけであるが、多数の粒 子の剛体連結で構成する複雑形状の DEM 計算は個々の粒 子に働く力の総和計算が必要になり、そのような複雑形状 の物体が流体中に多数存在するサスペンション・フローで は通信コストが急激に増加する.複雑物体が分割領域間に 跨って存在する場合は通信が非常に煩雑になる.多数の複 雑形状の物体に対するメモリ管理、物体間の通信、領域間 の通信を効率的に制御する必要がある.

本論文では、動的な領域分割を行いつつ各計算領域での 流体粒子計算, DEM 粒子で構成される複雑形状物体の計算 の効率化と通信負荷の低減を行い, これまでに達成されて いない規模のサスペンション・フローの大規模計算を GPU スパコンで実現する.

2. 粒子法による流体構造連成計算

計算する系を構成する流体,壁,物体の全てを粒子により表現しているため,計算量域内には図1に示すように, 流体粒子,壁粒子,さらに各物体を構成する物体構成粒子 の3種類の粒子が存在する.通常は3種類の粒子のいずれ も同じデータ構造を持ち,これらの粒子に対して以下に示 す4種類の相互作用を計算することで時間積分を行い,粒 子法による流体構造連成計算が実現できる.

(1) 流体粒子-流体粒子間の相互作用 SPHや MPS (Moving Particle Simulation)[7] などの粒子法を用いてカー ネル半径内の粒子から受ける相互作用を計算し, 粒子 の速度や位置, 圧力や密度を時間発展させることがで きる.本研究ではこれまでの粒子法に比べて計算精度 を大幅に向上した改良型 SPH 法 [8] を導入する.

(2) 流体粒子-物体構成粒子・壁粒子との相互作用 流体

図1 粒子法による流体構造連成計算の模式図.

Fig. 1 The kinds of particles to compute fluid-structure interaction simulations using particle-based method.

粒子計算の中で壁粒子は静止した境界条件として,物 体構成粒子は移動境界条件として扱われる.

- (3)物体構成粒子-流体粒子の相互作用 物体構成粒子は 流体粒子との接触点において圧力およびせん断応力を 求める。接触点から重心を指す方向ベクトルと,それ に垂直な面への射影成分ベクトルに分解する。
- (4) 物体構成粒子-物体構成粒子の相互作用 複雑形状をした物体間の衝突の計算であり,接触している粒子間でDEM の衝突モデルによる反発と摩擦力を計算する.
 (3) と同じように接触点から重心を指す方向ベクトルと,それに垂直な面への射影成分ベクトルに分解する.

(1) と (2) の流体粒子の時間発展では通常の流体計算と同様に leap-frog 法, Runge-Kutta 法, 予測子-修正子法等の 2 次精度以上の時間積分法を用いる. ここでは予測子-修正子法を用いている. (3) と (4) の物体の時間発展では,相互作用計算のあとに以下の力とトルクの総和計算を行い,並進運動と回転運動について前進オイラー法による時間積分を行う.物体構成粒子 x_i に作用する力を f_i とし,各粒子の剛体重心からの相対座標を \tilde{r}_i として,物体の質量を M とすると,各物体の並進速度 v,角速度 ω ,角運動量 L は,以下のように求められる.

$$\mathbf{M}\frac{d\mathbf{v}}{dt} = \sum_{i \in solid} f_i \tag{1}$$

$$\frac{d\boldsymbol{L}}{dt} = \sum_{i \in solid} \tilde{\boldsymbol{r}}_i \times \boldsymbol{f}_i \qquad (2)$$
$$\boldsymbol{\omega} = \boldsymbol{I}^{-1} \boldsymbol{L} \qquad (3)$$

 I^{-1} は各時刻の慣性行列の逆行列であり初期の慣性行列の 逆行列 I^{-1} から回転行列により形式的に求められる.また, 物体の姿勢の管理にはクオータニオンを導入する.物体の 並進運動により更新された重心位置と,回転運動により更 新された姿勢から, x_i が更新される.

3. 粒子法の GPU 計算と高速化

3.1 Linked-list 法を用いた流体粒子の相互作用計算にお けるソートの有効性評価

各粒子の持つ速度や座標, 圧力などの時間更新される従 属変数は通常, 粒子構造体の中のメンバ変数として保持さ れ, GPU の Device メモリ (CUDA ではグローバルメモリ と呼ばれている) 上に確保される. GPU 上では1スレッド が1粒子を処理するスレッド並列計算を行う. 各スレッド では担当の粒子について粒子法の物理モデルに基いた演算 を行う.

GPU による演算は CUDA コアと呼ばれる最小演算ユニットでスレッド毎に行われるのに対し, 従属変数のデータはグローバルメモリ上にあるため粒子 *i* が粒子 *j* から受ける作用を計算する場合は粒子 *j* へのメモリアクセスの時



- 図2 Linked-list に合わせて定期的に実行されるセル番号をハッシュ 値とした粒子データのソート.
- Fig. 2 The sorting of particle data by cell ID executed from CPU at fixed interval.

間が問題となる. DEM では接触判定を行い, SPH や MPS などではカーネル半径内の粒子との相互作用を行う. 全粒 子との相互作用計算を行うことは非効率的であるため,計 算領域を空間格子(セル)に分割し自身と隣接するセルに属 する粒子とのみ相互作用計算を行うセル分割法を用いて計 算する.

流体計算では通常,精度維持のためにカーネル関数の影響半径を初期粒子間隔の 3~4 倍程度の大きさに設定する. セル分割法の一セルの辺の長さをカーネル関数の影響半径 以上にする必要があるため,少なくとも 27~64 個程度の粒 子が一セルに格納される.静的に空間格子のメモリを確保 し,すべての粒子番号をそれぞれの粒子が属するセルに登 録する通常のセル分割法ではメモリ消費量が多過ぎる [7]. 各セルにはセル内の一つの粒子の番号のみを登録し,各 粒子が同一セル内の粒子の番号を鎖状につなぎ保持する Linked-list を GPU 上で導入している.Linked-list の生成は 逐次的な処理を伴うため GPU では atomic 関数を用いて 逐次的に行われるが,全体の計算時間と比べて Linked-list の生成時間は無視できるほど小さい.

粒子法の計算では、計算が進むにつれて粒子が移動する ために周囲の粒子のメモリ・アドレスはランダムになる. 例えば図2に示すように、123番のセル内に粒子番号が3 番と9998番と14000番の粒子があったとすると、隣接す るセル内のある粒子と123番のセル内の粒子との相互作用 を計算する際には、123番のセルに生成されたLinked-list をたどって粒子番号が3、9998、14000番の粒子に順にア クセスする.しかしながら、図2の上側のように粒子が並 んでいた場合、3番の粒子と9998番、14000番の粒子はメ モリの位置が大きく離れている.CUDAのwarpでは一度 に16連続アドレスの配列データをdata load するが、3番の 粒子を読みに行った際の連続した16アドレス内には9998 番や14000番の粒子のデータ存在しておらず、9998番と 14000番を含む16アドレスの data load がそれぞれ必要に なり、メモリ・アクセスに時間がかかってしまう.

図2の下側のように、粒子データをセル番号順に並び替 える操作を行うと、3番と、9998番、14000番の粒子はメ モリ上で連続するため、3番の粒子にアクセスする際に、 連続して隣接する 9998 番や 14000 番も data load され,高 速に Linked-list をたどることができる.本研究では GPU の Device メモリ上で粒子の従属変数が連続となるよう粒 子構造体の中に座標や速度や圧力など物理量のポインタを メンバ変数として保持し、それぞれのポインタに対して連 続的に配列を確保する SOA(Structure of Array)型のデータ 格納形式による実装を行う.通常良く行われるような物理 量をメンバ変数に持つ粒子構造体を配列として確保する AOS(Array of Structure)型のデータ格納形式の実装の場合 は、CUDA にも標準で付属する Thrust ライブラリの API を 用いるなどにより粒子の並び替えの操作は比較的容易に行 うことができる.しかし、SOA型のデータ格納形式の場合 は、並び替えのキー(ハッシュ値)となるセル番号を一時的 にバッファに保存するなど実装は複雑であり、また 20変 数以上に及ぶ物理量のそれぞれに対して並び替えの操作を 順次行うため効率的ではない。そこで簡単のため、セル番 号による粒子データの並び替えは GPU から CPU 側にデー タをコピーしてきた上で CPU 側で行う. セル番号による 並び替えのコストや、CPU と GPU の粒子データのコピー が計算の大きなオーバーヘッドとならないように回数を抑 えて定期的に実行する. Linked-list 法を用いない通常のセ ル分割法に対して粒子の並び替えを行う手法はこれまでに も多く報告されているが [9],メモリ制約の大きい GPU 上 で Linked-list 法を用いてメモリ容量を抑え, 合わせてセル 番号による粒子の並び替えを CPU 側から定期的に実行す る本手法の有効性は報告されておらず、新しい手法の提案 である。40万個の壁粒子と100万個の流体粒子を合わせ て140万個の粒子を用いた流体のダム崩壊計算に対して、 セル番号による並び替えを行う場合と行わない場合の1ス テップあたりの計算時間の変化を比較した結果を図3に示 す。セル番号による粒子の並び替えを行わない場合、計算 ステップが進むにつれて実行性能が低下して計算時間が増 加するのに対し、セル番号による並び替えを行う度に計算 性能が回復し計算時間を大幅に短縮できていることがわか り、5000 ステップ目の段階で並び替えを行わない場合と比 較して 8.6 倍の高速化を達成していることが分かる.

図3の計算では50ステップ毎に粒子の並び替えを実行 しており、1回の並べ替えに要する時間は350msecであっ た.50ステップで平均すると1ステップ当たり平均7msec であり、全体の計算時間と比較して十分に小さい。

3.2 物体粒子の時間積分に対する GPU 実装

3.1 節で提案する方法は,2章で述べた4種類の相互作用 計算における流体粒子の時間発展((1)と(2))において有効 である.一方,物体粒子の時間発展は,複雑形状の物体が





- 図3 GPU における Linked-list 法を用いた近傍粒子探索にセル番号 による粒子データのソートを行う場合と,行わない場合の計算 時間の比較.
- Fig. 3 The comparison of the computational elapsed time of the neighbor particle search using the linked-list with/without sorting by cell ID.

流体中に多数存在するサスペンション・フローの計算を前 提としているため、流体の時間積分の方法とは大きく異な る.計算領域内に全部で $N_{particle}$ 個の粒子があり、 N_{obj} 個 の物体があるとすると、物体の時間積分を行う疑似コード は Algorithm 1 のように書ける. Algorithm 1 中の条件分 岐は、粒子 j が物体 k の構成粒子であるかどうかを判定す るものである.

GPU 計算では、(A) 物体数についての外側のループをス レッド並列化する場合と、(B) 領域内の粒子数についての 内側のループをスレッド並列化する場合の2通りの実装が 考えられる.(A)を実装する場合、外側のループに CUDA のスレッド並列を割り当て、各スレッドが内側のループの 処理を行うように実装する.(B)を実装する場合、あらか じめ領域内の粒子数分の配列を用意しておき、GPU のス レッド並列で受ける力を書き込んでおく、外側のループ で k 番目の物体の総和計算部分には Thrust ライブラリの inclusive scan を実行する.

サスペンション・フローの計算を前提とするために 140 万粒子のダム崩壊問題に対して計算領域内に一つあたり 1000 個の物体構成粒子からなる 100~400 個の立方体を加



図4 実装(A)を用いた場合の計算時間の内訳.

Fig. 4 The breakdown of the computational elapsed time with the implementation of (A).







えて配置し, (A) と (B) のそれぞれの実装に対して, 1 ス テップあたりの計算時間を測定した結果をそれぞれ図4と 図5に示す.

GPU上で物体番号について並列化を行う実装(A)は物 体数が増えても緩やかに計算時間が増加するのに対し,物 体数について並列化されていない実装(B)を用いた場合, 実装(A)と比べ計算時間の増加が急である。物体数が300 個を超えると実装(B)の方が計算時間を要しており,数千 個に及ぶ多数の物体を扱う本研究のような場合には,実装 (A)を用いることにより,物体数が増えても計算時間の増 加を抑えることができる。

3.3 単体 GPU での実行性能

単体 GPU の計算について,NVIDIA 社の提供する性能 解析ツールである Visual Profiler を用いてカーネル関数 の 実行時間や FLOPS(Floating-point Operations Per Second) を 測定した.140万個の流体粒子に対して,1000個の粒子か ら構成される立方体を 400 個配置したダム崩壊問題に対す る実装 (A) を用いた場合の実行性能の内訳を図6に示す.





- Fig. 6 Elapsed time of each GPU Kernel (Numeric values on each bar show performance (Units: GFlops)).
- 表1 セル番号によるソートを毎ステップ実行した場合の実行性能の 理想値。
- Table 1
 The ideal performance for the case with the execution of sorting in every time step.

計算項目	単精度	倍精度
予測子:	102.6 (2.59)	88.1 (6.72)
粒子粘性:	94.3 (2.38)	82.5 (6.29)
粒子圧力:	104.7 (2.65)	122.1 (9.32)
座標の修正子:	67.4 (1.70)	91.1 (6.95)
圧力補間:	90.6 (2.29)	106.8 (8.15)
速度の修正子:	67.2 (1.70)	91.7 (7.00)

単位: GFlops (括弧内は理論ピーク性能に対する割合 (%))

浮動小数点演算は流体計算部分に対して単精度計算を行 い,数値誤差の蓄積が生じ易い総和計算については倍精度 計算を行っている.左から順に流体の時間発展計算(図4 における水色)6種類のカーネル関数と,物体の時間発展の ための力とトルクの総和計算を行うカーネル関数,トルク の総和計算に用いる重心(式(2)の*r*_i)を決めるための構成 粒子の位置座標の総和計算を行うカーネル関数となってい る.それぞれの実行性能(GFlops 値)も記載してある.

流体の時間発展を計算する各カーネル関数の実行性 能は図6の場合では最大で約65.6 GFlops となり, NVIDIA Tesla K20X の単精度の理論ピーク性能である 3.95 TFlops に対して約1.66%と低い. セル番号による粒子のソートを 毎ステップで実行した理想的な場合の各カーネル関数だけ の実行性能を表1に示す. 括弧内は理論ピーク性能に対す る実行性能の割合を示しており,単精度で2.65%,倍精度 で9.32%となっている. 図6における流体の時間発展の 計算はそれに対して約62%程度であり, ソートによるオー バーヘッドがかかることを考慮すれば,妥当な値である.

物体の時間発展のための総和計算を行うカーネル関数は メモリアクセスが支配的であるため、実行性能は非常に低 い. 一つあたりの構成粒子数が少ない物体が数万個から数 十万個あるような粒度の細かいサスペンション・フロー計



図7 2 次元のスライスグリッド法を用いた動的負荷分散の概要. Fig. 7 Outline of dynamic load balancing based on 2-dimensional slicegrid method.

算には実装 (A) はより有効となる. 実装 (B) は総和計算に Thrust ライブラリの inclusive_scan を用いており計算効 率は良いが,物体数について並列化できないため,物体数 が数個以下の場合に有効になる.

本研究では計算コードの最適化は十分とは言えず,実行 性能の上限値について詳細な議論を行うためには,Roofline Model[10] などを用いた検討が必要になる.

4. 動的負荷分散を用いた複数 GPU 計算

著者らの研究グループでは、これまで個別要素法による 粉体シミュレーションの複数 GPU 計算および SPH 法によ る流体の複数 GPU 計算に対して、Fig.7 に示すように計算 領域をスライスグリッド法により動的に 2 次元領域分割 し、分割された各小領域内の粒子数を一定に保つ負荷分散 を導入してきた [3][6][11][12]. ここでは、その手法を流体 粒子と多数の物体が存在するサスペンション・フローの計 算に図8のように適用する.

小領域境界の移動は小領域を微少幅 △l で分割して各分 割内の粒子を数え上げ,粒子数の均一化に必要な境界領域 に含まれる粒子を GPU 上でパッキングし,隣接小領域の GPU に転送する方法を取っている.異なるノード上に搭載 された GPU 間では直接のデータ通信はできないため,ホ



- 図8 サスペンション・フローの流体構造連成問題の複数 GPU 計算 に対する 2 次元のスライスグリッド法を用いた動的負荷分散.
- Fig. 8 Dynamic load balance based on 2-dimensional slice-grid method for the suspension flow problems with multiple GPUs.

スト計算機を介して通信を行っている。パッキングされた 粒子を CPU 側に転送し, CPU 間で MPI ライブラリ [13] に おける MPI_Isend と MPI_Irecv を用いた非同期通信を行 い、受信した粒子をもとに境界線を変更して領域の再分割 が完了する。粒子データの GPU 間通信が発生するのは以 下の3種、粒子が運動したことにより小領域間を移動する 場合,小領域の境界を移動したことにより粒子の所属する 小領域が異なった場合、カーネル半径内の粒子と相互作用 するためにハロー領域(=小領域の境界から、カーネル半径 の長さだけ内側の領域)に存在する粒子を隣接する小領域 にコピーする場合である.これらの粒子移動により、GPU のグローバルメモリ上で粒子が移動前に占めていたメモリ 領域は以降の計算で参照されなくなり、計算が進行するに つれてこのような未参照メモリがグローバルメモリ上に離 散的に発生し、使用するアドレスが断片化していく. そこ で、粒子データを先頭から詰め直して未参照メモリを除外 するため、3.1節で述べたセル番号によるソートとは別に 粒子の再整列を実行する.

流体-構造連成の複数 GPU 計算では、一つの物体の構成 粒子が同一小領域内にある場合は少なく、領域境界を跨り 複数の GPU のグローバルメモリ上に構成粒子のデータが 分散する。従って物体の時間積分は力とトルクの総和計算 においてノード間で MPI ライブラリによる通信が必要にな る.まず自身の小領域内で総和計算を行い MPI_Gather を 用いて結果をルートに送信する. ルートでは,受信した物 理量を足し合わせて物体の時間積分した結果を MPI_Bcast を用いて各小領域にブロードキャストする.現在は1対多 型の MPI による集団通信となっていて,これが並列化効 率を低下させる。物体の時間積分は重心が存在する小領域 のホストで行うことにより通信量を低減させることができ るが、その物体の構成粒子が存在する小領域との関係を毎 ステップ更新するための効率的な手法を開発する必要があ るため本論文では行っていない。物体の時間積分に関する 動的負荷分散は導入されていないが、サスペンション・フ ローであっても流体粒子に比べれば物体の数は圧倒的に少 なく計算負荷は小さい。

5. TSUBAME2.5 における大規模計算

東京工業大学学術国際情報センターの GPU スパコン TSUBAME2.5 は, 1 / ードに GPU (NVIDIA KeplerK20X) が 3 枚装着され,全体で 4,224 枚搭載されている.TSUB-AME2.5 の複数 GPU を用いた粒子法によるサスペンショ ン・フローの大規模計算を行い,合わせてその実行性能を 測定する.

5.1 複数 GPU によるサスペンション・フローの大規模 GPU 計算

最大規模のサスペンション・フローとして, 計算領域の

縦横高さを 144 m×160 m×60 m とし,深さ 1.6 m の静止 した水を張り,そこに合計 2,304 個の浮遊する物体を配置 する.流体の挙動を計算するために合計で 8743 万個の流 体粒子を用い,物体を構成する粒子は 1,000 個に固定して いる. 256 個の GPU を用い,スライスグリッド法による 2 次元の領域分割を行う.初期状態は物体が規則正しく並ん でいて,1 GPU 当たりの計算領域は 9 m×10 m×10 m と なり,6 個の物体を含んでいる.津波が押し寄せてくる状 況を想定し,右側に高さ 4.8 m の水柱を設定している.

図9の最下図は時刻6.8 secの計算結果(上側8枚の右列 上から2番目)の拡大図である.サスペンション・フロー の計算により,物体が浮遊していることで津波が砕波する 先端がクリアでなくなっていることが分かる.物体を搬送 するために波のエネルギーを奪われるために流れは緩やか になり,物体同士が衝突することで非常に不規則な流れに なっている.大規模計算では多数の浮遊する物体と流れと の相互作用を計算することができ,物体が集団的に流され る速度などの津波被害予測への貢献が期待できる.

物理時間で 6.8 sec に相当する計算を時間刻み 5×10⁻⁴ sec で 13600 ステップ計算するのに必要な計算時間は 43.1 時 間であった. 100 ステップ毎に計算結果のバックアップを 行い,合計が 1.1 TB の計算結果を GPU 側から CPU 側に 転送する時間と計算結果を出力するファイル I/O の時間ま で含めた全計算時間は 96.1 時間であった.

5.2 TSUBAME2.5 における弱スケーリング

GPU の台数に比例させて、粒子を敷き詰める水平方向の 面積を拡大させる.1ステップの実行時間をT,全粒子数 を N_{particle} として Performance = N_{particle}/T により弱スケー リングを評価する.

図 10 において赤は 1 ステップの計算全体の実行性能, 青は流体の時間発展部分の実行性能,緑は総和計算 (座標, 力,トルク)部分の実行性能を表す.黒の実線は 4GPU の 全体の実行性能を GPU の台数に比例させ理想値であり,点 線は 4GPU の流体の時間発展の実行性能を GPU の台数に 比例させた理想値である.

流体計算部分のみの実行性能(青)は,理想値(点線)に対 して、32GPUと256GPUを用いた場合にそれぞれ96.7%、 52.8%であった.計算時間全体の実行性能(赤)は、32GPU と256GPUを用いた場合に理想値(実線)に対してそれぞ れ73.4%、27.3%であった.流体計算のみの場合と比べて、 総和計算が加わることにより並列化効率が悪化することを 示している.総和計算部分では、座標、力、トルクの3種 類について各々 Reduction 操作を行っており、ノード内で Reduction を行った後、全プロセス間でMPI_Gather を用い て全体の Reduction を行う.図10における1回の Reduction 操作に要する計算時間の内訳を図11に示す.弱スケー リングが保たれていることによりノード内の Reduction 時







Fig. 9 A simulation result of suspension flow with using 87,430,000 particles with 2304 objects with 256 GPUs.

間は一定に保たれており,GPU 数の増加に伴う実行性能の 低下はノード間の Reduction(一対多型の MPI による集団通 信) が原因だと分かる.







図 11 図 10 における Reduction の内訳. Fig. 11 The break down of the data reductions in fig.10.

6. 結論と今後の課題

サスペンション・フローの GPU コードを開発し,最大 で 256 GPU を用いて 144 m×160 m×60 mの計算領域に 8,743 万個の流体粒子と 2,304 個の浮遊する物体を含んだ 大規模計算を実行することができた.また,以下のことが 明らかになった.(a) Linked-list を用いた近傍粒子探索に対 して定期的にソートすることで 8.6 倍に高速化することが できた.(b) 物体構成粒子同士の相互作用により受ける力 とトルクを合算して物体の時間積分を行う部分では,物体 数に対して CUDA のスレッド並列を行うことにより,物 体数が増えても計算時間の増加を抑えることができる.(c) これまで紛体のみの計算や流体のみの粒子計算に適用して きたスライスグリッド法による動的領域分割が流体と物体 を含むサスペンション・フローの流体構造連成計算に適用 することができた.

東京工業大学学術国際情報センターのスパコン TSUB-

HPCS2015 2015/5/20

AME2.5 において弱スケーリングを調べ,4GPUを用いた 場合に対して,256 GPUを用いた場合には流体計算部分の みで 52.8%,連成計算を含む場合に 27.3% に並列化効率が 低下する結果を得た.

今後の課題として、ノード内及びノード間の Reduction は 多くの改善の余地がある。ノード間の Reduction について、 今回は座標、力、トルクについて別々に Reduction を行っ たが、力とトルクのノード間 Reduction についてはパッキ ングすることにより通信回数を半分に抑制できる。物体を マスターノードで管理しているが、重心を含むノードが管 理することにより、通信負荷の分散を行うことができる。 現在のノード内の Reduction では、物体が変形する場合も 想定して毎回重心を求め直すことをしているが、剛体の問 題に特化すれば、重心については総和計算が必要なくなる。 また、物体の内部まで粒子を充填する必要もなくなるため、 ノード内の Reduction による計算コストは軽減される。

本研究ではこれまで殆ど成されてこなかったサスペン ション・フローの大規模計算を GPU スパコンで実現でき たことが最も大きな意義であり,物体ごとに形状が異なる ようなより現実的な場合も含め,効率化と高速化を進めて 行く必要がある.

謝辞

本研究の一部は,科学研究費補助金・基盤研究(S) 課題番号 26220002「ものづくり HPC アプリケーションのエクサ スケールへの進化」,科学技術振興機構 CREST「ポストペ タスケール高性能計算に資するシステムソフトウェア技術 の創出」,学際大規模情報基盤共同利用・共同研究拠点,特 別研究員奨励費・課題番号 14J12004「超大規模粒子シミュ レーションの演算加速器型スパコンにおけるフレームワー クの構築」,および革新的ハイパフォーマンス・コンピュー ティング・インフラから支援を頂いた.記して謝意を表す.

参考文献

- [1] Hamada, T., Yokota, R., Nitadori, K., Narumi, T., Yasuoka, K. and Taiji, M.: 42 TFlops hierarchical N-body simulations on GPUs with applications in both astrophysics and turbulence, in *Proceedings of the International Conference* on High Performance Computing, Networking, Storage and Analysis, SC'09, pp. 1–12 (2009).
- [2] Ishiyama, T., Nitadori, K. and Makino, J.: 4.45 Pflops Astrophysical N-body Simulation on K Computer: The Gravitational Trillion-body Problem, in *Proceedings of the International Conference on High Performance Computing, Networking, Storage and Analysis*, SC'12, pp. 1–10 (2012).
- [3] Tsuzuki, S. and Aoki, T.: Large-scale granular simulations using Dynamic load balance on a GPU supercomputer, in *Poster at the 26th IEEE/ACM International Conference on High Performance Computing, Networking, Storage and Analysis*, SC'14 (2014).
- [4] Murotani, K., Koshizuka, S., Ogino, M., Shioya, R. and Nakabayashi, Y.: Development of distributed parallel ex-

plicit Moving Particle Simulation (MPS) method and zoom up tsunami analysis on urban areas, in *Poster at the* 26th IEEE/ACM International Conference on High Performance Computing, Networking, Storage and Analysis, SC'14 (2014).

- [5] Domnguez, J., Crespo, A., Valdez-Balderas, D., Rogers, B. and Gmez-Gesteira, M.: New multi-GPU implementation for smoothed particle hydrodynamics on heterogeneous clusters, *Computer Physics Communications*, Vol. 184, No. 8, pp. 1848 – 1860 (2013).
- [6] 都築怜理, 青木尊之, 井元佑介, 田上大助: GPU スパコン における動的負荷分散を用いた粒子法による大規模流体 シミュレーション, 第 28 回数値流体力学シンポジウム講 演予稿集 (2014).
- [7] 越塚誠一,柴田和也,室谷浩平:粒子法入門,丸善(2014).
- [8] 井元佑介,田上大助:ある一般化粒子法に用いる近似作用素の打ち切り誤差解析とその応用,第27回計算力学講演会講演論文集(2014).
- [9] Bader, M., Bode, A., Bungartz, H.-J., Gerndt, M., Joubert, G. R. and Peters, F. eds.: *Parallel Computing: Accelerating Computational Science and Engineering (CSE)*, Vol. 25 of *Advances in Parallel Computing* (2014).
- [10] Williams, S., Waterman, A. and Patterson, D.: Roofline: An Insightful Visual Performance Model for Multicore Architectures, *Commun. ACM*, Vol. 52, No. 4, pp. 65–76 (2009).
- [11] 都築怜理, 青木尊之, 下川辺隆史: GPU スパコンにおける 1 億個のスカラー粒子計算の強スケーリングと動的負荷 分散, 情報処理学会論文誌. コンピューティングシステム, Vol. 6, No. 3, pp. 82–93 (2013).
- [12] Tsuzuki, S. and Aoki, T., : Large-scale Particle Simulations using Dynamic Load Balance on TSUBAME 2.0 Supercomputer, in *Proceedings of the 5th Asia Pacific Congress on Computational Mechanics and and 4th International Symposium on Computational* (2013).
- [13] MPI Forum, : Message Passing Interface (MPI) Forum Home Page, http://www.mpi-forum.org/ (2009).