ワークスチーリング戦略のカスタマイズによる AMR法の高速化

中島潤^{1,a)} 河野瑛² 田浦健次朗^{1,b)}

概要:タスク並列処理によって記述されたアプリケーションを効率的に実行するためには、アプリケーショ ンやハードウェアに関する知識を考慮してタスクのスケジューリングを行うことが必要である.これを達 成するために、我々は以前にワークスチーリングの戦略をカスタマイズ可能なタスク並列処理の実行時シ ステムによって、プログラマがスケジューリング方針を最適化する手法を提案した.本稿では提案手法の 有効性を評価することを目的として、適合細分化格子法 (AMR)による移流計算に対して提案手法を適用 し、ワークスチーリング戦略をカスタマイズして評価を行なった結果について述べる.

キーワード:タスク並列処理,スケジューリング,適合細分化格子法

Accelerating Adaptive Mesh Refinement by Customizing Work Stealing Strategy

Jun Nakashima^{1,a)} Akira Kono² Kenjiro Taura^{1,b)}

Abstract: In order to execute task-parallel application efficiently, it is important to take characteristics of application and hardware into account. In order to address this issue, previously we proposed a runtime system which can customize the strategy of work stealing. In order to evaluate the effectiveness of our proposal, we applied the customization to an advection calculation using Adaptive Mesh Refinement. This paper describes the customized work stealing strategy and its evaluation results.

Keywords: Task Parallelism, Scheduling, Adaptive Mesh Refinement

1. はじめに

ノード内のコア数の増加,NUMA アーキテクチャによ るメモリの多階層化,アクセラレータの導入などによって 現在の並列計算機は非常に大規模かつ複雑なものとなって おり,さらに今後ますます複雑化していくものとみられて いる.それに対して現在広く普及しているプログラミング 処理系の多くは,並列計算機のハードウェアを単純に抽象 化した機能を提供しており,このような大規模並列計算機

1 東京大学

- Hitachi
- $^{\rm a)}$ nakashima@eidos.ic.i.u-tokyo.ac.jp

の性能を引き出すようにアプリケーションを記述すること は困難である.そのため,ハードウェアをより高度に抽象 化したプログラミングモデルを提供しており,平易な記述 で高い実行性能を達成できる生産性の高い並列プログラミ ング環境の構築が並列計算における重要な課題となってお り,これを解決するため Chapel[1] や X10[2] をはじめとす る様々なプログラミング処理系が提案され,活発に開発が 行われている.

これらの処理系の多くは,アプリケーションの並列性を 表現するための機能としてタスク並列処理を提供してい る.タスク並列処理では,処理全体が細粒度に分割して軽 量なスレッドに割り振られ (この軽量なスレッドは一般に タスクと呼ばれる.本稿でも以降タスクと呼ぶ),それらを 実行時システムがハードウェア (多くの場合計算ノード内

University of Tokyo

² 日立製作所

 $^{^{\}rm b)}$ tau@eidos.ic.i.u-tokyo.ac.jp

IPSJ SIG Technical Report

の CPU コアであるため,本稿でも以降 CPU コアである とする) に適切に割り振って実行する. +分に多くのタス クに処理が分割されていれば,処理系がタスクを CPU コ アに割り振ることによって動的な負荷分散が自動的に行わ れるため,タスク並列処理を用いることで再帰や多重ルー プをはじめとする,処理を単純に分割するだけでは効率的 な並列化ができないアプリケーションを平易に記述するこ とができる.

アプリケーション内で多くのタスクが生成,破棄される ため、タスク並列処理の実行時システムには、小さなオー バーヘッドでタスクを生成,破棄したうえで、それらを効 率的に各 CPU コアにスケジューリングすることが求めら れる.これを達成するためのスケジューリング手法として は、Lazy Task Creation[3] が広く知られているが、この手 法はタスク間の負荷のばらつきなどのアプリケーションの 特性やメモリ階層などのハードウェアの特性を考慮したス ケジューリングを行わない.

そのため、これらの特性を活かすことが特に重要になる 大規模な分散メモリ環境では、単に Lazy Task Creation を利用するだけでは効率的なスケジューリングを行えない と考えられる.今後並列計算機の複雑化にしたがって、こ の影響がますます大きくなることを考えると、アプリケー ションやハードウェアの特性について配慮したタスクのス ケジューリングを行うことは、タスク並列処理を用いて大 規模並列計算機上で高い実行性能を達成するうえでの重要 な課題であるといえる.

この課題を解決するために,我々はタスクのスケジュー リングの戦略をカスタマイズすることが可能な機能を備え たタスク並列処理系によって,プログラマがその方針をア プリケーションやハードウェアにあわせて最適化する手 法 [4] を提案しており,さらに単純なベンチマークを用い てその有用性を評価している.

我々はより実際のアプリケーションに近く,なおかつタ スク並列処理の利点を活かせるアプリケーションを用いて 提案手法の有用性を評価することを目的として,河野ら [5] による,タスク並列処理を用いて実装された適合細分化格 子法 (AMR)による移流計算を題材として,これに計算粒 度の不均衡と NUMA アーキテクチャを考慮するようなス ケジューリング戦略のカスタマイズを適用して評価を行 なった.本稿では,タスク並列 AMR に対して適用したス ケジューリング戦略とその評価について述べる.

2. タスク並列処理のスケジューリング

2.1 Lazy Task Creation

タスクの生成や合流処理を並列に行うため,タスク並列 処理によるアプリケーションは多くの場合,再帰的にタス クを生成する形で記述される.そのような形で記述され たタスク並列処理によるアプリケーションを効率的に実 呼称する)と、実行可能なタスクを格納するためのキュー (本稿ではランキューと呼称する)を配置し、それぞれの ワーカーがランキューからタスクを取りだし、実行するこ とでタスク並列処理を行う.

行するためのスケジューリング手法として, Cilk[6]をは

じめとする, 共有メモリ環境を対象とした多くのタスク

並列処理系において用いられている手法として Lazy Task

各ワーカーは以下の方針に基づいてタスクを実行する:

- タスクを作成する際は、それまでのタスクを中断して ランキューの先頭に格納して、新しく作成したタスク を即座に実行する
- タスクが終了あるいは中断した際はランキューの先頭 からタスクを取りだして実行する
- ランキューにタスクがない際には他のワーカーをランダムに選択して、その末尾からタスクを奪ってくる (ワークスチーリング)

ランキュー末尾からのワークスチーリングによって奪わ れるタスクは、タスクの親子関係の木構造の根に近く、将 来的に多くの子タスクを生成する、すなわち粒度が大きい タスクであるため、少ないスチール回数で効率的に負荷分 散をとることができる.また、ワークスチーリングが発生 した箇所を除けば、それぞれのタスクは逐次に再帰呼び出 しを行ったのとまったく同じ順序で実行されるため、再帰 の実行順を利用してメモリアクセスなどの最適化を行った 逐次プログラムをタスク並列処理によって並列化した際に も、まったく同様に最適化の効果が得られることが期待さ れる.

ワークスチーリングの際にタスクを奪う対象は実際に奪 われるタスクの粒度や、ハードウェアの構成などの要素を 考慮することなくランダムに選ばれる.そのため、特にこ れらの要素に配慮することが重要な場合に大きな性能低下 が発生しうる.例えば分散メモリ環境においては、データ にアクセスするオーバーヘッドがノード内とノード間で大 きく異なるため、ランダムなワークスチーリングではデー タの配置を意識したタスクのスケジューリングができない ことが原因で実行性能が大幅に低下すると考えられる.

2.2 ワークスチーリング戦略のカスタマイズ

様々なハードウェアやアプリケーションの特性を考慮した,効率的なタスクのスケジューリングを行うことを目的 として,我々はスケジューリング戦略をアプリケーション やハードウェアにあわせてプログラマがカスタマイズする 手法 [4] を提案して,それを MassiveThreads[7],[8] 上に実 装している.この手法では,タスクの実行順や,大きい粒度 のタスクのスチールを行えるといった Lazy Task Creation **IPSJ SIG Technical Report**

の利点を維持しつつスケジューリング戦略をカスタマイズ するために,カスタマイズ可能な範囲をワークスチーリン グの戦略に絞り,それに必要な機能をプログラマに提供し ている.

提案する手法を用いてワークスチーリング戦略のカスタ マイズを行う際の手順を以下に示す.

2.2.1 ヒントの割り当て

1	void user_task (void $*$ param){
2	
3	// 追加部分: ヒントを割り当てる
4	struct $\{\}$ hint;
5	set_task_hint(&hint, sizeof (hint));
6	// 再帰的にタスクを作成
$\overline{7}$	spawn_task(user_task,param1);
8	spawn_task(user_task,param2);
9	
10	}

まず,アプリケーションに変更を加え,ワークスチーリ ングのヒントとして用いられる任意のデータを各タスクに 割り当てる.ヒントの内容や割り当ては対応するタスクか ら任意のタイミングで変更,解除が可能である.ヒントと して割り当てるデータの種類は問われないが,確保,解放 の手間を簡略化できることを考えると,タスク内のローカ ル変数を割り当てることが最も望ましい*1.

Lazy Task Creation においては、タスクを作成する際に 新しく作成されたタスクが即座に実行され、親のタスクが ランキューに格納されて他のワーカーにスチールされる. したがって、ここで割り当てられるデータの内容は、タス クを作成した後に実行する処理に関するものとするべきで ある.例えば上の疑似コードでは、6行目でタスクが作成 されたことによって、このタスク他のワーカーにスチール された後は、そのワーカー上で7行目の param2 を引数と するタスクを作成するコードが実行される.したがって、 5 行目で割り当てられるヒントの内容は、7行目のタスク 作成に関するものであることが望ましい.

2.2.2 ユーザー定義のワークスチーリング関数の実装

1	// ワークスチーリング関数: ワーカーがアイドル状態
2	// になった際に呼び出される
3	void user_defined_steal(int id)
4	{
5	task_handle t_stolen;
6	// <i>1</i> .スチール対象の候補を決める
7	$int candidates[] = get_random_workers();$
8	// 2.ヒントを読み取る
9	struct $\{\}$ hints = readydeque_peek(candidates);
10	/*
11	* 3.ヒントをもとに実際にスチールを行う
12	* 対象や順序を決める
13	*/
14	int targets[]=;
15	// 4.スチールを行う
[4]	で発表した際にはタスクの作成時にヒントの初期データとサ

*1 [4] で発表した際にはタスクの作成時にヒントの初期データとサイズを割り当て、それ以降は内容を変更することしかできなかったが、現在は改良が加えられ、上に説明したような動作となっている.また、タスク作成時に初期データを割り振る手法についても引き続き利用可能である.

for t in targets{ 16 $t_stolen = readydeque_trysteal(t,$ 17 steal_confirm, data); 18 if (t_stolen)break; 19 20} 21return t_stolen; 22} 23// 確認用関数: ワークスチーリングが成功することが 24// 確定した際呼び出される 25int steal_confirm(task_handle_t stolen, 2627void *userdata) 28{ 29int ok; 30 /* * 5.ここで実際にスチールするか判断 31 32 // 0以外を返すとスチール完了 33 // 0を返すと中断 34 35 return ok; 36 37 38 int main(int argc, char **argv) 3940// 6.ワークスチーリング関数を切り替え 41 42set_steal_func(user_defined_steal); 43 ... } 44

次に, ヒントを利用してワークスチーリングを行うユー ザー定義のワークスチーリング関数を実装して, 標準の ワークスチーリング関数と切り替える. ワークスチーリン グ関数はワーカーのランキューにタスクがなくなった際 に呼び出されるコールバック関数で, アイドル状態のワー カーの ID を引数としてとる. 典型的には, ワークスチー リング関数は, 以下のような手順でワークスチーリングを 行うように実装される.

- 候補選択 ワークスチーリング対象の候補となる複数の ワーカーを選択する.
- **ヒントの読み取り** それぞれの候補から盗むことができる タスクのヒントを読み取る.
- **スチール対象の決定** ヒントをもとにスチールを行う対象 のワーカーやそれらの間の順序を決定する
- **スチールの試行**対象のワーカーに対してワークスチーリ ングを試行する.この際,コールバック関数としてス チールされたタスクを確認するための関数を同時に指 定する.
- 結果の確認 ワークスチーリングの際にヒントを読み取ったものと同じタスクがスチールされるとは限らないため、実行するには望ましくないタスクがスチールされてしまう可能性がある。それに対処するため、スチールが成功することが確定した段階で、スチール試行時に指定したユーザー定義の確認用コールバック関数が呼び出される。この関数はスチールしたタスクをもとに、スチールを確定するか否かを判断する。
 - ユーザー定義のワークスチーリング関数,確認用関数,

およびワークスチーリング関数の切り替え部分の例を上の 擬似コードに示す.

3. 関連研究

ハードウェアやアプリケーションの特性に配慮した効率 的なタスクのスケジューリングを行うことを目的として, 様々なスケジューリング手法が提案されている.

例えば、Chen らが提案する CATS[9] は、プロファイリ ングによって得られたタスクが利用するデータサイズを推 測し、共有キャッシュからはみ出さないようにタスクの配 置を行う.また、Wimmerら [10] は一度に盗むタスクの数 や実行順などをアプリケーションの特性に応じてプログラ マが自由に設定できるスケジューラを提案している.

分散メモリ環境では、データにアクセスするオーバー ヘッドがノード内とノード間で大きく異なるため、分散メ モリ環境を対象としたタスク並列処理系の多くが、それに 配慮したスケジューリングを行っている.

例えば, Scioto[11], [12] では,タスクが作成される際, そのタスクがノード内で実行されることが望ましい場合 はランキューの先頭に追加するのに対し,別のノードで実 行されるのが望ましい場合は,ランキューの末尾に追加し て,他のノードに効率よくスチールされるようにしてい る.広域分散環境を対象としたタスク並列処理系である Satin[13], [14] は,広域ネットワーク上で効率的な負荷分散 を行うために,ネットワークの階層構造にも配慮してワー クスチーリングを行っている.

これら既存のスケジューリング手法の多くは,タスク並 列処理系に組み込まれており,実際の処理はプログラマか ら隠蔽されている.そのため,スケジューリング手法に対 して最適化を行う余地はあくまでタスク並列処理系が API として公開している範囲に限定されており,それ以上の高 度なチューニングを行うためには,タスク並列処理系の内 部的な実装に踏み込んで,スケジューラを改変することが 求められる.アルゴリズムのみが発表されており,実装が 入手できない手法をアプリケーションに対して適用する際 についても同様に,何らかの既存のタスク並列処理系を改 変してアルゴリズムを実装する必要がある.

それに対して我々が提案しているスケジューリングのカ スタマイズ手法は、タスク並列処理系の実装とユーザーが カスタマイズする部分とは分離されているため、処理系の 内部的な実装に踏み込むことなくスケジューリングの方針 をカスタマイズすることができる.しかもスケジューリン グ方針は C/C++言語の関数として実装されるため、カス タマイズの自由度が非常に高い.

さらに,既存のスケジューリング手法の中には,実行順 が逐次の再帰プログラムと異なるものがあるのに対して, 我々が提案しているカスタマイズ手法では,カスタマイズ の範囲がワークスチーリングの戦略に限局されていること



図 1 格子とそれに対応する四分木 ([15] より引用) Fig. 1 Mesh and Corresponding Quadtree (cited from [15])

により逐次再帰と同じ実行順が保証されているため,それ に由来する利点を失うことなくカスタマイズを行うことが できる.

4. タスク並列処理による Tree-Based AMR

この章では,スケジューリングのカスタマイズの題材として用いた,河野ら [5] によるタスク並列 Tree-Based AMR を用いた二次元空間の移流計算について説明する.

4.1 Tree-Based AMR

数値シミュレーションによって偏微分方程式を解く際に は計算対象の空間を離散化して,離散化された空間上の格 子点に対する計算を行うことで解を得る.得られる解の精 度は格子の細かさによって決まるため,高い精度の解を得 るためには非常に多くの格子点が必要になり,それに伴っ て計算量も非常に大きくなる.しかし,実際には細かい格 子による計算が必要な領域は領域内の一部であることが多 い.そこで,高精度な計算が必要な部分には細かい格子を, 低い精度の計算でも十分な部分には粗い格子を動的に割り 当てることで,計算量の削減と高い精度の計算を両立する 手法が適合細分化格子法 (AMR)[16] である.

Tree-Based AMR は AMR の一種で,立方体 (二次元の 場合は正方形)の形状の直交格子を各座標軸に沿って再帰 的に二分割して細分化することで,適合的な格子を生成す るものである.計算領域を離散化する格子全体が四分木 (三次元の場合は八分木)で表現され,計算対象の格子点は 四分木のリーフノードに対応している (図 1).

4.2 計算の流れ

タスク並列 Tree-Based AMR による移流計算は, C++ 言語によって記述されており, MassiveThreads を用いて 並列化されている.

各タイムステップは以下に述べる calculate, check, remesh, balance, ghost の5つのフェーズからなり, い ずれのフェーズも再帰的にタスクを生成して四分木全体を 並列に処理するような形で記述されている.

calculate 一次精度風上差分法によって各格子点のタイ ムステップを進める.

2013/8/1

- check 各格子点について、細分化、集約化の必要がある かどうかの判断は Parthasarathy ら [17] による,隣接 格子点との物理量の差の統計量から判断する手法に よって行われる. このフェーズではその判断に必要と なる, 各格子点における隣接する格子点との物理量の 差の平均,および標準偏差を計算する.
- remesh check フェーズで計算した統計量と、隣接格子 点との物理量の差の最大値を比較して, その結果をも とに格子の細分化,および集約化を行う.具体的には, valmax を隣接格子点との物理量の差の最大値, µ を平 均, δ を標準偏差, K_{high},K_{low} を定数としたとき,細 分化、および集約化の判断は以下の式にしたがって行 われる.
 - $val_{max} \mu \ge K_{high}\delta$ のとき:細分化
 - $val_{max} \mu \leq K_{low}\delta$ のとき:集約化 細分化、および集約化の際には物理量の補間が必要だ が、この実装では簡単化のため、細分化する際には単 に粗い格子の物理量をコピーしている.また,集約化 を行う際には細かい格子の平均をとっている.
- balance 解の精度を保つためには、隣接する格子の細分 化レベル差が1以内に収まっている必要があるため, 隣接する格子間の細分化レベルを比較して、その制約 を満たすように細分化を行う.このフェーズは河野 ら [15] によって提案された分割統治法に基づくアルゴ リズムによって実装されている.
- ghost 細分化レベルが異なる境界の計算を正しく行うた めに,両者の境界にある粗い格子の一つ下のレベルに ダミーの格子 (ghost cell とよばれる) を配置して,粗 い格子の物理量をそれに補間する. この補間も簡単化 のため、粗い格子の物理量をコピーによって実装され ている.

calculate, remesh フェーズは四分木全体を並列にた どって処理を行うのみなのに対して, check, balance, ghost フェーズは、各タスクが再帰的に子タスクを生成し て全ノードに対する処理を行った後、子タスクの終了を待 ち、改めて子ノードどうしが接する部分に対して再帰的に タスクを生成して処理を行う,という複雑な制御構造を もっている.

なお、典型的な AMR の実装では毎タイムステップ実行 する必要があるのは calculate と ghost フェーズのみで, 他のフェーズは適宜間隔を空けて行われるものであるが, 本稿の評価では全フェーズを毎タイムステップごとに実行 している.

4.3 スケジューリングの課題

上に述べたタスク並列 Tree-Based AMR の特性から導 かれるスケジューリングの課題としては、以下の二点が挙 げられる.

- 不均衡な木の構造 AMR の格子を表す四分木は、細分化 が行われた一部の領域にノードが集中する不均衡なも のとなるため、ランダムワークスチーリングによる負 荷分散が前提としている、再帰の親子関係の根に近い ほど粒度が大きいという特徴が成り立たない. そのた めスチールの回数が増加し、それに伴うオーバーヘッ ドによって実行性能が低下すると考えられる.
- メモリへの負担 フェーズごとに木構造をたどる必要があ るため、メモリのバンド幅が AMR の性能に大きく影 響する.近年主流となっている NUMA 構成のマシン でメモリバンド幅を活かすためには、タスクとデータ が配置されるソケットをできるだけ一致させる必要が あるが、ランダムワークスチーリングはこれを考慮し ないため、メモリのバンド幅を十分に活かすことがで きない.

5. ワークスチーリング戦略のカスタマイズ

この章では、4.3節で述べた課題を解決するために実装 したワークスチーリング戦略のカスタマイズ内容について 説明する.

5.1 方針

ワークスチーリング戦略のカスタマイズは以下の二つの 方針に基づいて行なった.

- スチール粒度の向上 不均衡な木構造においてもできるだ け粒度の大きなワークスチーリングを行うために、各 タスクに粒度を表す指標を付加して、それをもとに粒 度の大きなタスクを選択してスチールする. 各フェー ズの処理は木を再帰的にたどるものであるため, 粒度 の指標としては、タスクが処理を行う四分木を根とし た部分木のノード数とする.スチール粒度の向上に よって,スチールの回数ならびにオーバーヘッドが減 少するだけでなく、後述するハードウェアプリフェッ チ機構に配慮したメモリ再配置がより効果的に機能す ることが期待される.
- NUMA 構成への配慮 四分木の各ノードがどのワーカー で処理されたかを記録しておき、それをもとに、ある ワーカーがワークスチーリングを試みる際は、できる だけ直前のフェーズにそのワーカー上で処理された ノードを利用するタスクをスチールするようにする. さらに, 定期的に木構造の各ノードをワーカーに対 応する CPU ソケット上のメモリに再配置すること で、ソケットをまたぐメモリアクセスの頻度を削減し、 NUMA 構成のマシンにおけるメモリアクセスの効率 を高める. ハードウェアプリフェッチ機構を活用する ために,メモリの再配置の際には、タスクの実行順に 沿ってノードのメモリアドレスができるだけ連続にな るように留意する.本稿執筆には、このメモリ再配置

機能の実装が間に合わなかったため,これは今後行う 予定である.

あるワーカーがワークスチーリングを試みた際, 粒度が 大きいが直前に違うワーカーで実行されていたタスクと, 粒度が小さいが直前にそのワーカーで実行されたタスクが 同時にスチール可能になる場合など,これらの方針が競合 することも考えられるが,その際は後者を優先する.タイ ムステップごとの四分木の構造変化はそれほど大きくない ため,直前に行われたスケジューリングと同じ方針をとる ことによる性能劣化は少ないことが期待されるためである.

5.2 追加のデータ構造

1	struct Tree {
2	// タスク並列 <i>Tree–Based AMR</i> のデータ構造
3	double T[2]; // 物理量
4	struct Tree *child; // 子ノードへのポインタ
5	double highest_diff; // 隣接格子点との差分の最大値
6	int level; // 細分化レベル
$\overline{7}$	bool is_leaf; // 葉ノード <i>(=</i> 格子点)であれば <i>true</i>
8	// ワークスチーリング戦略をカスタマイズする
9	// ために追加するデータ
10	int n_nodes; // 自身を根とした部分木のノード数
11	int last_worker; // 最後に実行したワーカー
12	$- attribute_{-}((aligned(64)));$
13	
14	// 各タスクに割り当てるヒント
15	struct ws_hint{
16	int n_nodes;
17	int last_worker;
18	};

まず、粒度の指標となる部分木のノード数と、最後にそのタスクに関する処理を実行したワーカーを格納するために、四分木のノードを表す構造体にそれらを表す整数 n_nodes と last_worker を追加する.また、各タスクに 割り当てるヒントも同様に、n_nodes と last_worker か らなる構造体として定義している.

上にデータ構造の宣言部の疑似コードを示す. なお,構造体のメンバ数自体は増加しているが, false sharing を避けるために構造体のサイズはキャッシュラインのサイズで切り上げられているため,追加する前後で実際に確保されるデータの大きさは変化していない.

5.3 スケジューリングヒントの割り当て

1	void traverse_tree (Tree *t){
2	// ヒントを確保, 初期化してタスクに割り当てる
3	struct ws_hint hint;
4	hint.n_nodes = $t \rightarrow n_nodes;$
5	$set_task_hint(\&hint, sizeof(hint));$
6	// last_worker を更新
7	$t \rightarrow last_worker = get_worker_num();$
8	// 子ノードを並列にたどるタスクを作成する
9	for $(i=0;i<4;i++)$ {
10	// <i>child[i+1]</i> が盗まれた後最初に実行されるので
11	// それが実行されたワーカーをヒントに指定
12	if (i<3)
13	$hint.last_worker = t - > child[i+1].last_worker;$
14	$spawn_task(traverse_tree, t->child[i]);$

15} _____終了待ち 16 17sync_tasks; // ヒントを解除 18 set_task_hint(NULL,0); 1920 /* * check, balance, ghost フェーズでは, 21* ここでさらに再帰的にタスクを作成する 22 23*/ 2425// 最後にノード数を更新 (Calc フェーズのみ) 26 $t \rightarrow n_n des = 1 + t \rightarrow child[0] + t \rightarrow child[1]$ +t -> child[2] + t -> child[3];2728

各フェーズで再帰的にタスクを作成して四分木をたどる 部分の擬似コードを上に示す.子タスクを作成する直前に ローカル変数としてスケジューリングヒントを確保して, 対応する四分木のノード数をヒントに書き込んだ後タスク に割り当てる.次に,ノード構造体のlast_workerメン バに現在タスクを実行しているワーカーの ID を書き込む.

子タスクが生成される際には、現在のタスクが他のワー カーにスチールされて、別のワーカーで次の子タスクが生 成される.したがって、ヒントのlast_workerメンバに はタスクが作成される直前に、次に作成されるタスクがも つ子ノードのlast_workerを書き込む.

さらに,タスクが終了する直前に,子のノード数をもと に,自身を親とする部分木のノード数を更新する.これは どのフェーズでも行うことができるが,各フェーズにこの 処理を追加してオーバーヘッドを評価した際,最も性能低 下の度合いが小さかった calc フェーズでのみ行うものと した.

なお, check, balance, ghost の3つのフェーズでは, 子に対して再帰的にタスクを作成して処理を行なった後, ノード間の接する部分を対象にしたタスク並列処理を改め て行う.ここにもヒントを割り当てることが可能だが,こ の部分の再帰はノード間の接する部分のみをたどるもので あるため,部分木全体のノード数は正確な粒度の指標には ならないと考えられる.そのため,現在の実装ではこれら のタスクを作成する前にヒントの割り当てを解除してい る.今後の課題として,この部分にも適切な指標を割り当 てることについても検討を行っている.

5.4 ワークスチーリング関数

1	int n_candidates;// 候補ワーカー数
2	
3	// ワークスチーリング関数: ワーカーがアイドル状態
4	// になった際に呼び出される
5	void mysteal(int id)
6	{
7	task_handle t_stolen;
8	// スチール対象の候補を決める
9	$int candidates[n_candidates]$
10	$= get_random_workers();$
11	// ヒントを読み取る
12	$\mathbf{struct} $ ws_hint hints[n_candidates]

情報処理学会研究報告

IPSJ SIG Technical Report

13	= readydeque peek(candidates):	
14	for $(int i=0; i < n candidates; i++)$	
15	int target:	
16	// ヒントとワーカー <i>ID</i> が一致する	
17	// タスクの中で最もノード数の	
18	// 大きいものを選択	
19	$target = pick_the_best_candidate(id,$	
20	candidates, hints);	
21	// ワークスチーリングを試行	E
22	$t_stolen = readydeque_trysteal(candidates[target]),$	
23	steal_confirm, &hints[target]);	
24	$hints[target].n_nodes=-1;$	
25	if (t_stolen)break;	\mathbf{C}
26	}	最適
27	return t_stolen;	
28	}	numa
29		
30	// 確認用関数: ワークスチーリングが成功することが	
31	// 確定した際呼び出される	
32	int mysteal_confirm(task_handle_t stolen,	
33	void *userdata)	
34	{	
35	struct ws_hint *hint_read = userdata;	
36	struct ws_hint *hint_stolen	
37	$= get_{hint_data(stolen)};$	
38	// ノート数か読み取つたものより	
39	$\frac{1}{16}$ (hint) = n and $\frac{1}{16}$ (hint) = n and $\frac{1}{16}$	
40	II (nint->n_nodes > nint_stolen->n_nodes)	
41	$\mathbf{return}(0;$	
42	$(get_worker_num()) == mit - >iast_worker){// ii \pm o p - h - in > -m + z$	
43	// 坑山の / ハ ID C 取りる // タマカを次く だ退合け	木の
44	// Jast worker も比較	
40	\mathbf{if} (hint->last worker	
40	l hint stolen -> last worker)	
48	return 0:	
49	}	
50	return 1;	
51	}	
	,	

ワークスチーリングを行う手順を上の擬似コードに示す. まず,複数のワーカーからワークスチーリング対象をラン ダムに選択して,それらからヒントを読み取る.次に,そ れらのヒントを比較して,ワーカーの ID とヒントに保存 されている ID が一致するタスクに対して,粒度が大きい 順にスチールを試行する.ワーカーの ID とヒントに保存 されている ID が一致するタスクがない場合は,単に粒度 が大きい順にスチールを試行する.ノード間の接する部分 を再帰的にたどるタスクはヒントの情報をもたないため, これらのタスクは粒度が0であるものとして,最も低い優 先度でスチールする.

ワークスチーリングの確認関数では、実際にスチールし たタスクと以前に読みだしたヒントのタスクの粒度を比較 して、想定よりも小さい粒度のタスクを盗んでしまった場 合はこれをキャンセルする.また、ワーカーの ID とヒン トに保存されているワーカーの ID が一致する場合は、これ についても確認して、一致しないならばキャンセルする.

6. 評価

表1に示す共有メモリ環境を用いて、5章で述べたワーク スチーリング戦略のカスタマイズに関する評価を行なった.

表 1 評価に用いた環境 Table 1 Experimental Setup

CPU	Opteron 6172 (2.1GHz) 4 ソケット	
010	$12 \ \exists \ \times 4 = 48 \ \exists \ $	
メモリ	128GB (32GB × 4 ソケット)	
	L1D: 64KB/コブ	
キャッシュ	L2: 512KB/コア	
	L3: 12MB/コア	
OS	Linux $2.6.32$ (Debian)	
C コンパイラ	GCC 4.7.2	
最適化オプション	-O3	
numactl パラメータ	interleave=all	

表 2 移流計算の問題設定

oblem Setting
r

空間の大きさ	(1.0, 1.0)
速度場	(0.1, 0.1)
境界条件	周期境界条件
CFL 数	0.4 (最も細かい格子)
タイムステップ数	500
細分化レベル	最小 8 (256×256) 最大 19
細分化閾値 K _{high}	0.1
集約化閾值 Klow	0.05
木のリーフあたりの格子点数	1
細分化を行う周期	毎タイムステップ
逐次再帰に移行する閾値	2048 ノード



- 図 2 細分化レベルを最小 5, 最大 8 とした場合の初期条件. 赤い 部分は物理量が大きい部分,青い部分は物理量が小さい部分で ある.
- Fig. 2 Initial Condition (refinement level=[5,8]). Red parts have high value and blue parts have low value.

6.1 問題設定

計算に用いる問題には、二次元空間の移流方程式を表 2 に示す設定のもとで解くものを用いた.初期条件には、 円盤状に物理量が大きい部分が存在している、図 2 のよ うな物理量の分布を用いた.これは Langer ら [18] によ る Charm++[19] によって記述された AMR の実装の評価 に用いられているデータセットを参考にして設定したも





の*2で、物理量が大きい部分と小さい部分の境界部分のみに細分化が行われるため四分木が不均衡なものとなり、さらに境界部分が時間発展によって移動するため細分化や集約化が頻発するような設定となっている.

6.2 計算の粒度

ワークスチーリングのカスタマイズによるタスクの粒度 向上の効果を確認するために,ワーカー数を48に固定し て,ワークスチーリング対象の候補数を変化させた際に, 計算がどの程度の粒度で行われているかを大まかに計測し た.この計測は,calculateフェーズにおいてスチールさ れることなく連続で実行された部分木のノード数を2のべ き乗間隔に集約したヒストグラムを作成して,それに対し て区間の最小値をかけ,それぞれの粒度が全体に占める割 合を大まかに計算することで行なっている.

得られた結果を図 3 に示す (粒度が1ノードから 4095 ノードまでの区間は統合してある).カスタマイズを行わな い場合 (グラフ内の"Original"以降のグラフにおいても同 様)と比較して、カスタマイズを行った場合は大きい粒度の 計算が占める割合が増加していることがわかる.calculate フェーズにおいては候補数の変化による粒度の変化はほと んどみられなかったが、他のフェーズにおいて粒度の傾向 が同様になるとは限らないため、今後他のフェーズについ ても同様な調査を行う予定である.

6.3 直前のタイムステップとのワーカーの一致度

ワーカー数を48に固定して、ワークスチーリング対象 の候補数を変化させた際の、木の各ノードの処理が直前の フェーズと同じワーカーで実行される確率を図3に示す. ただし、check、balance、ghostフェーズにおける、ノー ド間の接する部分のみを再帰的にたどる部分に関してはこ



図 4 直前のフェーズと同じワーカーでノードの処理が実行される 確率

Fig. 4 Worker Matching Ratio vs. # of Candidates



れに含まれない.カスタマイズを行わない場合 (グラフ内 の"Original") は直前のフェーズと同じワーカーでノード の処理が行われる確率はおよそ 2%であるのに対して,カ スタマイズを行った場合はその確率が 30%以上に増加して いる.

6.4 フェーズごとの実行時間

ワーカー数を 48 に固定して, ワークスチーリング対象 の候補数を変化させた際のフェーズごとの実行時間を図 5 に示す.カスタマイズを行わない場合と比較して,カスタ マイズを行なった場合は特に balance, ghost フェーズの 実行時間が減少している.候補数の変化による実行時間の 変化はランダムワークスチーリングとの差と比較してわず か (候補数 8 ワーカーの場合と 48 ワーカーの場合との差が 0.5 秒程度) ではあるが,候補数を増やすほど実行時間は短 くなっている.

^{*2} 細分化や集約化の方針,粗い格子から細かい格子への補間方法, 四分木の葉あたりの格子点数,および細分化を行う周期が Langer らの評価と異なっている.

35



図 6 ワーカー数を変化させた際の実行性能 Fig. 6 Performance vs. # of Workers

6.5 実行性能

ワークスチーリング対象の候補数を [(ワーカー数) - 1] として、ワーカー数を変化させた際の実行性能を図6に示 す.縦軸の単位は一秒あたりのタイムステップである.カ スタマイズを行なったもの (グラフ内の "Customized") は カスタマイズを行わないものと比較してワーカー数が多い 部分で性能が向上しており、性能向上の割合が最も大きい 48 ワーカーの場合は単位時間あたりのタイムステップ数が 9.8%向上している.

NUMA に配慮したメモリ再配置機構がまだ実装されて いないことから、この性能向上の理由は主にタスクの粒度 が向上したことによると推測される. 今後より詳細な解析 を行い,性能向上の理由を明らかにする予定である.

7. まとめと今後の課題

タスク並列処理を用いた Tree-Based AMR による移流 計算の性能を向上させるために、部分木のノード数を指標 として、スチールの粒度を向上させると同時に、NUMA 環 境に配慮して直前に実行されたワーカーにできるだけタス クが配置されるようにワークスチーリング戦略のカスタマ イズを行なった.

評価を行った結果、カスタマイズによってタスクの粒度、 直前に実行されたワーカーに配置される確率の両方が改善 していることが確認できた.また,48コア使用時の実行性 能も 9.8%改善された.

今後は以下の実装や評価を予定しており、これらを通じ てタスク並列処理による Tree-Based AMR ならびに本稿 で述べたカスタマイズ手法の有効性を明らかにする予定で ある.

- NUMA環境に配慮したメモリの再配置を行う機構の 実装および評価
- 今回得られた結果のより詳細な解析
- Langer らによる実装と同一条件での性能比較 •

謝辞 本研究は、JST CREST「高性能・高生産性アプ リケーションフレームワークによるポストペタスケール高 性能計算の実現」および JSPS 特別研究員奨励費 (課題番 号 248391) の助成を受けて行われたものである.

参考文献

- [1]Crav. The chapel parallel programming language. http: //chapel.cray.com.
- Philippe Charles, Christian Grothoff, Vijay Saraswat, [2]Christopher Donawa, Allan Kielstra, Kemal Ebcioglu, Christoph von Praun, and Vivek Sarkar. X10: An Object-Oriented Approach to Non-Uniform Cluster Computing. In OOPSLA '05: Proceedings of the 20th annual ACM SIGPLAN conference on Object-oriented programming, systems, languages, and applications, pp. 519–538, New York, NY, USA, 2005. ACM.
- [3]E. Mohr, D. A. Kranz, and R. H. Halstead, Jr. Lazy Task Creation: A Technique for Increasing the Granularity of Parallel Programs. IEEE Trans. Parallel Distrib. Syst., Vol. 2, No. 3, pp. 264–280, 1991.
- [4]Jun Nakashima, Sho Nakatani, and Kenjiro Taura. Design and Implementation of a Customizable Work Stealing Scheduler. In Proceedings of the 3rd International Workshop on Runtime and Operating Systems for Supercomputers, ROSS '13, pp. 9:1-9:8, Eugene, Oregon, 2013. ACM.
- 河野瑛,田浦健次朗. 分割統治法を用いたタスク並列 Tree-[5]based AMR アルゴリズム. Master's thesis, 東京大学, February 2013. http://hdl.handle.net/2261/54203.
- [6]Robert D. Blumofe, Christopher F. Joerg, Bradley C. Kuszmaul, Charles E. Leiserson, Keith H. Randall, and Yuli Zhou. Cilk: An Efficient Multithreaded Runtime System. SIGPLAN Not., Vol. 30, No. 8, pp. 207-216, 1995.
- 中島潤,田浦健次朗. 高効率な I/O と軽量性を両立させる [7]マルチスレッド処理系. 情報処理学会 論文誌 プログラミ ング (PRO), Vol. 4 No. 1, pp. 13–26, March 2011.
- Jun Nakashima. MassiveThreads: A Lightweight [8] Thread Library for High Productivity Languages. http://code.google.com/p/massivethreads/.
- [9]Quan Chen, Minyi Guo, and Zhiyi Huang. CATS: Cache Aware Task-Stealing based on Online Profiling in Multisocket Multi-core Architectures. In Proceedings of the 26th ACM international conference on Supercomputing, ICS '12, pp. 163–172, New York, NY, USA, 2012. ACM.
- Martin Wimmer, Daniel Cederman, Jesper Larsson [10]Träff, and Philippas Tsigas. Work-stealing with Configurable Scheduling Strategies. In Proc. of PPoPP '13, pp. 315-316, 2013.
- James Dinan, Sriram Krishnamoorthy, D. Brian Larkins, [11] Jarek Nieplocha, and P. Sadayappan. Scioto: A framework for global-view task parallelism. In Proc. of ICPP '08, pp. 586-593, 2008.
- [12]James Dinan, D. Brian Larkins, P. Sadayappan, Sriram Krishnamoorthy, and Jarek Nieplocha. Scalable Work Stealing. In Proc. of SC '09, pp. 53:1-53:11, 2009.
- [13]Rob V. van Nieuwpoort, Thilo Kielmann, and Henri E. Bal. Satin: Efficient Parallel Divide-and-Conquer. In Proc. from the 6th International Euro-Par Conference on Parallel Processing, pp. 690-699, 2000.
- [14] Rob V. van Nieuwpoort, Thilo Kielmann, and Henri E. Bal. Efficient Load Balancing for Wide-Area Divide-and-Conquer Applications. In Proc. of PPoPP '01, pp. 34-

IPSJ SIG Technical Report

43, 2001.

- [15] 河野瑛,田浦健次朗. タスク並列モデルを用いた Treebased AMR の評価. Technical Report 35,東京大学大学 院情報理工学系研究科,東京大学大学院情報理工学系研究 科, jul 2012.
- [16] Marsha J. Berger and Joseph E. Oliger. Adaptive Mesh Refinement for Hyperbolic Partial Differential Equations. Technical report, Stanford University, 1983.
- [17] V. Parthasarathy, Y. Kallinderis, and Kengo Nakajima. Hybrid Adaptation Method and Directional Viscous Multigrid with Prismatic-Tetrahedral Meshes. *AIAA pa*per 95-067, 1995.
- [18] Akhil Langer, Jonathan Lifflander, Phil Miller, Kuo-Chuan Pan, , Laxmikant V. Kale, and Paul Ricker. Scalable Algorithms for Distributed-Memory Adaptive Mesh Refinement. In Proceedings of the 24th International Symposium on Computer Architecture and High Performance Computing (SBAC-PAD 2012), New York, USA, October 2012.
- [19] L.V. Kalé and S. Krishnan. CHARM++: A Portable Concurrent Object Oriented System Based on C++. In A. Paepcke, editor, *Proceedings of OOPSLA'93*, pp. 91– 108. ACM Press, September 1993.