## タスク間とスレッド間の通信を考慮した 可変並列度 Fork-Join タスクのスケジューリング

西川広記<sup>1</sup> 島田佳奈<sup>1</sup> 谷口一徹<sup>2</sup> 冨山宏之<sup>1</sup>

概要:本論文は、均質なマルチコア上における可変な並列度を有する Fork-Join タスクのスケジューリング手法を提案 する.可変な並列度を有するタスクは、各タスクが複数のスレッドに分割されて並列に実行されることを許し、その スレッド数はスケジューリングと同時に決定する.本研究では、DMA を用いたタスク間およびスレッド間の通信に 加えて Fork および Join の際に発生するオーバーヘッドを考慮したスケジューリング手法を提案する.

キーワード:タスクスケジューリング,整数計画法,並列タスク,マルチコア

#### 1. はじめに

近年,シングルコアにおける性能向上が鈍化しており, 組込みシステムにおけるマルチコア/メニーコア技術のさ らなる普及が期待されている.そのなかでも,マルチコア の効率的な利用に向け,マルチコアを用いてタスクの実行 順序を決定するタスクスケジューリングの研究が一層進め られている.文献[1]ではマルチコアを用いてシングルコア 上のみで各タスクが実行されることを想定したスケジュー リング手法が提案されている.しかしながら,現実世界の マルチメディア領域などにみられるアプリケーションでは, 多くのタスクが並列タスクと呼ばれ,各タスクは複数のス レッドへ分割され並列な実行が許容されたタスクである. そこで,タスクを複数のスレッドへ分割することを想定し た並列タスクのスケジューリング問題が近年盛んに研究さ れている [2-10].

並列タスクを想定したスケジューリング分野の研究領 域において,並列タスクは大きく二つの種類に分けられる. 一方は,各タスクの実行に用いられる並列度(スレッド数) がスケジューリングの事前に決定されている,固定並列度 タスクである [3]. 文献[3]では、リストスケジューリング に基づき並列タスクをスケジューリングする方法が提案さ れている.各タスクはそれぞれ固有の並列度を有しており, 複数のコア上で各タスクのスレッドが同時に実行されるよ うな方式を想定したスケジューリング問題を、整数計画法 に帰着させて定式化し、全体のスケジュール長を最短化す ることを目的にしている.他方は、各タスクの実行に用い られる並列度がスケジューリングと同時に決定される、可 変な並列度を有するタスクを想定したスケジューリングで ある [4-11]. 文献[4]は, 文献[3]とは異なり, 可変な並列度 を有するタスクを想定している. この研究ではデッドライ ンを満たしながらハードウェアコストの最小化を目指して いる. さらに文献[5-9]でも可変並列度タスクに関する研究 が行われている. 文献[6]はリアルタイムシステム向けの可

変並列度タスクのスケジューリングが研究されている.文 献[9]では、線形計画法で定式化された可変並列度タスクの スケジューリングに対し、制約プログラミングを用いて定 式化されたスケジューリングを比較しており、スケジュー ル長(メイクスパン)の最短化を目的としている.

しかしながら、現実のシステムにおいては複数のタスク の間でデータ転送などによる通信が必要な場合がある.例 えば、あるタスクを実行したコアから別のタスクを実行す るコアに対してデータ転送を行う場合には通信が必要であ る.また、並列タスクにおいては独立に異なるコアで実行 される複数のスレッド間に対してデータの通信が必要な場 合がある.こうしたタスク間の通信と、タスク内のスレッ ド間の通信の両方を考慮したタスクスケジューリング問題 は、我々の知る限りでは過去に多くは研究されていない. 文献[10]では、タスク間の通信を考慮した並列タスクのス ケジューリングを提案しているが、スレッド間の通信につ いては考慮していなかった.加えて、この研究ではタスク における全てのスレッドが同時に異なるコア上で実行され ることを想定していたが、文献[9]などにみられるように、 現実世界の並列タスクの多くは Fork-Join 型であり、分割さ れた複数のスレッドはそれぞれ独立にスケジュールされる.

本稿では、タスク間とスレッド間の通信を考慮した可変 並列度 Fork-Join タスクのスケジューリングを提案する.提 案するスケジューリング手法は整数計画法に基づいてスケ ジュール長を最短化することを目的とする.

本稿の構成は以下の通りである.まず第2章では、本研 究において提案されるタスク間とスレッド間の通信時間を 考慮した可変並列度 Fork-Join タスクのスケジューリング 問題について、例題を用いながら説明する.そしてそのス ケジューリング問題を整数計画法に基づいて定式化する. 第3章では実験における環境や比較手法について述べ、さ らにその結果について述べる.最後の第4章では、まとめ と今後の課題を述べる.

<sup>1</sup> 立命館大学

Ritsumeikan University

<sup>2</sup> 大阪大学 Osaka University





1 pre

(6)

1\_3 (12)

(6)

1\_2 (12)

1\_1 (12) S

2

1\_4 (12)

	スレッド実行時間			スレッド間通信	
#並列度	前処理	ボディ	後処理	前通信	後通信
		スレッド			
1	0	40	0	0	0
2	2	21	2	2	2
3	2	14	4	2	2
4	6	12	6	1	1

(b) タスク1におけるスレッド実行時間・スレッド間通信時間



(c) 並列化されたタスクによるタスクグラフ

E

図 1 タスクスケジューリング問題の例

<sub>(5)</sub>(5

### 2. タスク間とスレッド間の通信を考慮した可 変並列度 Fork-Join タスクのスケジューリング

本章では、本稿で提案するタスク間とスレッド間の通信 を考慮した可変並列度 Fork-Join タスクのスケジューリン グを、2.1 節において例題を用いてそのスケジューリング の概要を説明し、また 2.2 節においてスケジューリング問 題を定式化する.

#### 2.1 スケジューリング例

本節では、本稿で提案するタスク間とスレッドの通信を 考慮した可変並列度 Fork-Join タスクのスケジューリング がどのようなものかを例を用いて説明する.図1がそのス ケジューリング問題の例である.図1(a)に示される重み付 き有向非循環グラフは一般的なアプリケーションをグラフ 化したものであり、本稿ではこれをタスクグラフと呼ぶ. まず、各ノードはアプリケーションにおけるタスクを示し ている. "S"と"E"でラベル付けされたタスクはそれぞれ開 始タスクと終了タスクを示しており、本研究ではこれらを ダミータスクと呼び、コアなどの資源を用いた演算処理を 要さないタスクとして定義する.そのほかのタスクにおい て、ノード内に数が記されたタスクはタスクグラフにおい てコアなどの資源を用いた演算処理を要するタスクとして (d) スケジューリング結果

定義される.ノード内における数はそれぞれタスクの ID を示している.各タスクの横にある()内における数は、そ のタスクがシングルコアで実行されたときの実行時間を示 しており、例えばタスク1をシングルコアで実行した場合 の実行時間は40単位時間である.ただし、開始タスクと終 了タスクはダミータスクであるため資源を用いた演算処理 を要さない.したがって、これらの実行時間は0である. 次に、タスクグラフにおける各エッジはタスク間の依存関 係を示す.各エッジ上の数は二つのタスク間で通信に要す る時間を示す.ここで、簡単化のために開始タスクと終了 タスクは資源による演算処理を要さないタスクであるため、 開始タスクから出ているエッジおよび終了タスクに向かっ ているエッジにおける通信時間は0としても一般性を失わ ない.

本稿におけるスケジューリング問題において、タスクは 可変並列度 Fork-Join タスクを想定しており、このタスクは 複数のスレッドに並列化可能である. 並列化されたタスク はマルチコア上において独立に、かつ、並列に実行できる. 例えば、タスク 1 が可変並列度 Fork-Join タスクであり 4 スレッドで実行される場合について考える. このとき、タ スク1に対して、タスクをスレッドへ分割するための前処 理を行うスレッド、計算を要する4つのボディスレッド、

それら4つのスレッドの計算結果を集約するための後処理 を行うスレッドの計6つのスレッドが生成される.図1(b) はタスク1を例とした、各並列度における各スレッドの実 行時間とスレッド間における通信時間の表を示す. タスク 1が4並列で実行される場合,前処理の実行時間は6単位 時間であり、4 つのボディスレッドの実行時間はそれぞれ 12時間単位であり,後処理の実行時間は6単位時間である. また、スレッド間通信は前処理からボディスレッドへ向け て行う通信を行う前通信とボディスレッドから後処理へ向 けて行う後通信に分けられ、タスク1が4並列で実行され る場合の通信時間はそれぞれ1である.また、この実行時 間および通信時間の表は各タスクに対して与えられる. 図 1(c)は図1(a)で示したタスクグラフのタスクのうち、タス ク1,3,4 が並列化された場合のタスクグラフである.図1 (d)は図1(c)のスケジューリング結果である. ここでタスク 1に注目すると、タスク1の前処理が時刻0から時刻6に かけてコア2で実行されている.1番目と3番目のボディ スレッドは、同じコア2で実行されている.一方、2番目 と4番目のボディスレッドはそれぞれコア1とコア3で実 行されている. これらのスレッドは前処理とは異なるコア で実行されているため、前通信が発生している. そのため 2番目のボディスレッドは時刻7から実行が開始されてい る. 次にタスク1の後処理はコア3で実行されている. こ のため、コア3以外で実行された1番目、2番目3番目の ボディスレッドから後処理に向けて通信が発生する.次に, 図1(d) におけるタスク3の後処理に注目する.ここでは, タスク3の後処理からタスク5に向けて2単位時間分の通 信が発生している.これは、タスク3の後処理とタスク5 の実行がそれぞれ異なるコア上でマッピングされているた めにタスク間の通信が必要となるためである. 簡単化のた め、本稿では各コアはどこへ向けても通信が可能であるよ うに十分に接続されているものとし、かつ、通信の競合が 発生しないような通信路を想定しているが、その場合もス ケジューリング問題の一般性は失われない. ただ、本研究 を拡張させて通信路の制約や通信の競合などの共有リソー スに対する制約を容易に付与することができる.

#### 2.2 整数計画法に基づくスケジューリング問題の定式化

本節では、均質マルチコアにおける可変並列度 Fork-Join タスクのスケジューリングを整数計画法に基づき定式化す る.本研究で開発したスケジューリング手法では、各タス クの並列度の決定、マルチコア上へのマッピング、タスク スケジューリングを同時に行う.本稿ではこのスケジュー リング手法は整数計画法に基づいて定式化されるが、この 定式化にいくつかの数式を新たに書き足して線形計画法と して定式化することも容易に行うことができる.

*par<sub>i,k</sub>* はタスク*i* が*k*個のボディスレッドに分割される 場合に1になる 0-1 決定変数である. すなわち*k*コアで実行 されるときに1の値をとり,それ以外の場合には0をとる.

$$\forall i, \qquad \qquad \sum_k par_{i,k} = 1 \tag{1}$$

Time\_pre<sub>i,k</sub>, Time\_body<sub>i,k</sub>, Time\_post<sub>i,k</sub>, はそれぞれ, タスクがk個のボディスレッドに分割された場合における, 前処理, ボディや後処理のスレッドの実行時間である. ま た, これらの値は与条件として与えられる. したがって, 前処理, ボディ, 後処理のスレッドの実行時間を表す time pre<sub>i</sub>, time post<sub>i</sub>, time\_body<sub>i</sub>は次の通りになる.

 $\forall i, \qquad time\_pre_i = \sum_k Time\_pre_{i,k} \times par_{i,k} \qquad (2)$ 

 $\forall i, \qquad time_post_i = \sum_k Time_post_{i,k} \times par_{i,k} \tag{3}$ 

$$\forall i, j, \quad time\_body_{i,j} = \sum_k Time\_body_{i,j,k} \times par_{i,k}$$
(4)

 $Comm_intra_pre_{i,k}$ および $Comm_intra_post_{i,k}$ は、タスクi がkコアで実行されるとき、それぞれ前処理からボディに 対する通信時間とボディから後処理に対する通信時間を表 す.  $Comm_intra_pre_{i,k}$ および $Comm_intra_post_{i,k}$ はいずれ も与条件として与えられる.  $pre_intra_{i,j}$ は、前処理とボデ ィとの間に通信が必要な場合に1の値をとる 0-1 決定変数 である. 同様に $post_intra_{i,j}$ は、ボディと後処理との間に通 信が必要な場合に1の値をとる 0-1 決定変数である.

- $\forall i, j, \ comm_intra_pre_{i,j} = \sum_k (Comm_intra_pre_{i,k} \times par_{i,k}) \times pre_intra_{i,j}$ (5)
- $\forall i, j, \ comm_intra_post_{i,j} = \sum_k (Comm_intra_post_{i,k} \times par_{i,k}) \times post_intra_{i,j}$ (6)

pre\_intra<sub>i,j</sub>はタスクiの前処理およびj番目のボディスレ ッドが異なるコアにマップされるときに1の値をとる.同 様にpost\_intra<sub>i,j</sub>はタスクiのj番目のボディスレッドおよ び後処理が異なるコアにマップされるときに1の値をとる. map\_pre<sub>i,k</sub>およびmap\_post<sub>i,k</sub>はタスクiの前処理と後処理 がそれぞれk番目のコアにマップされた場合に1の値をと る 0-1 決定変数である. map\_body<sub>i,j,k</sub>は,タスクiのj番目の ボディスレッドがk番目のコアにマップされたときに1の 値をとる 0-1 決定変数である.

∀i, j, pre\_intra<sub>i,j</sub>

$$= \begin{cases} 0 & \text{if } map\_pre_{i,k} = map\_body_{i,j,k} \text{ for any } k \\ 1 & \text{otherwise} \end{cases}$$
(7)

∀i, j, post\_intra<sub>i, j</sub>

$$= \begin{cases} 0 & \text{if } map\_body_{i,j,k} = map\_post_{i,k} \text{ for any } k \\ 1 & \text{otherwise} \end{cases}$$
(8)

*start\_pre*<sub>i</sub>および *finish\_pre*<sub>i</sub>はそれぞれ, タスク *i* の前処 理における開始時刻および終了時刻を示す. 同様に,





図 2 タスク間とスレッド間通信を考慮した MFJ タスクのスケジューリング結果

start\_body<sub>i,j</sub> および finish\_body<sub>i,j</sub>, start\_post<sub>i</sub> および finish\_post<sub>i</sub>はそれぞれ,ボディスレッドの開始時刻と終了 時刻,後処理の開始時刻と終了時刻を示す.

 $\forall i, j, \qquad finish\_pre_i = start\_pre_i + time\_pre_i \qquad (9)$ 

 $\forall i, j, \quad finish\_body_{i,j} = start\_body_{i,j} + time\_body_i \quad (10)$ 

 $\forall i, j, \quad finish\_post_i = start\_post_i + time\_post_i \quad (11)$ 

その際に、スレッド間の通信を考慮すると以下の制約式 が必要となる.

 $\forall i, j, finish\_pre_i + comm\_intra\_pre_{i,j} \leq start\_body_{i,j}$  (12)

 $\forall i, j, \quad finish\_body_{i,j} + comm\_intra\_post_{i,j} \leq start\_post_i$ (13)

複数のボディスレッドはマルチコア上で独立に実行す ることができる.それら複数のボディスレッドは並列,に 異なるコア上で実行されるか,あるいは,同じコア上で実 行される.後者の場合,同一時刻において複数のスレッド を同一のコア上で実行することが出来ない.この資源制約 は次のように表される.

$$\forall i, j1, j2, k \ (j1 \neq j2),$$
 (14)

# $$\begin{split} map\_body_{i,j1,k} &= map\_body_{i,j2,k} \\ \rightarrow \left(finish\_body_{i,j1} \leq start\_body_{i,j2}\right) \\ & \lor \left(finish\_body_{i,j2} \leq start\_body_{i,j1}\right) \end{split}$$

ここまでスレッド間における資源制約および順序制約 に注目してきた.タスク間についても同様に資源制約と順 序制約を満たす必要がある.しかし,本論文では割愛する. 本研究におけるスケジューリング問題の目的は,スケジ ュール長 (メイクスパン)の最短化である.すなわち,本 問題における目的関数は次の通りになる.

#### Minimize:

max<sub>i</sub>{finish\_post<sub>i</sub>}

(15)

本スケジューリング問題は整数計画問題として定式化 される.上式のいくつかは線形式ではないが、これらは容 易に線形化できる.したがって、本定式化は汎用の ILP ソ ルバー上で解くことができる.

#### 3. 実験

#### 3.1 実験手法

本節では、本研究で提案するタスク間・スレッド間の通信を考慮した可変並列度 Fork-Join タスクのスケジューリング手法と実験に用いたベンチマーク及び比較手法について説明する. ベンチマークには TGFF [11]を基にして生成した 14 個のタスクグラフを用いた. ターゲットシステム

は4コアおよび8コアの2種類を想定している. 求解ソル バーには ILOG CP Optimizer 12.8 を使用している. 本研究 で開発したスケジューリング問題は非常に複雑であり,解 空間が広いため実用的な時間内に最適解を求めることは困 難である. そのため,本研究では求解時間を5時間として 既定し,その時点における許容解をその手法における解と して評価に用いる. なお,実験環境は Core i9 7980XE (2.6GHz)プロセッサを使用しており,主記憶は 128GB であ る.

本研究では、比較手法として以下の4手法を用いる.

*Max*: 各タスクは全てのコア上で実行される. 言い換え ると, それぞれのタスクは *N* コアのターゲットシステムに おいて *N* 個のボディスレッドに分割される. そしてそれら のタスクは逐次的に実行される.

MS: 文献[10]で提案されている可変並列度タスクのスケ ジューリング手法である. この手法は本稿の手法とは異な り、タスクにおける全てのボディスレッドは同時に開始さ れることを想定している.

*MFJ*:本稿で提案している,可変並列度 Fork-Join 型タス クのスケジューリング手法である.

*MS-FJ*: 2 段階のヒューリスティックな可変並列度 Fork-Join タスクのスケジューリング手法である. 始めの 1 時間で MS 手法を用いて各タスクの並列度を決定し, 残り の4時間でスレッドのスケジューリングを行う手法である.

#### 3.2 実験結果

図2(a), (b)は実験結果を示す. 横軸はタスクグラフを示 し, 括弧内の数は各タスクグラフに内包されるタスク数を 表している.縦軸はタスクグラフのスケジュール長を表し, 結果は全て Max 手法によって正規化した.

図2(a)は4コアを用いたときのスケジューリング結果である. MS 手法, MFJ 手法や MS-FJ 手法はそれぞれ Max 手法に比べると,平均で6.2%,16.0%,12.6% スケジュール長を短くすることに成功した.ほとんどの場合で,MFJ 手法が4つの中では最も良い結果を得た.

図2(b)は8コアを用いたときのスケジューリング結果で ある. MS 手法, MFJ 手法, MS-FJ 手法はそれぞれ平均で 7.7%, 8.1%, 23.2%スケジュール長の縮小に成功した. 4 コアの場合と異なる点として, MS-FJ 手法が4つの中で最 も良い結果を得ていることが挙げられる. これは, 8 コア を用いた場合には解空間が巨大化し,計算量が膨大になる ために5時間という実行時間内にMFJ 手法は良いスケジュ ール長を得ることが困難であるからである. また, いくつ かのケースでは MFJ 手法は MS 手法よりも悪い結果を示し たが, その一方で, MFJ 手法のヒューリスティックな手法 である MS-FJ 手法は MFJ 手法よりも効率的に解空間の探 索を行うことができた. したがって MS-FJ 手法が最も良い スケジュール長を得ることができたと考えられる.

#### 4. おわりに

本稿では、タスク間とスレッド間の通信を考慮した均質 なマルチコアによる可変並列度 Fork-Join タスクのスケジ ューリング手法を提案した.提案手法は、各タスクが複数 のスレッドに分割されて並列に実行されることを許し、並 列度の決定、マルチコア上へのマッピング、タスクスケジ ューリングを同時に行う.実験において、小規模なタスク グラフにおいては本研究で提案した手法が有用性を示し、 タスクグラフの規模が大きくなるにつれ、提案手法をさら に発展させたヒューリスティック手法が有用であることが 示された. 今後はさらに現実的なアプリケーションを用い たタスクグラフに対して、より高速なヒューリスティック 手法の提案を計画している.

**謝辞** 本研究は一部,キオクシア株式会社(旧社名 東 芝メモリ株式会社)の支援による.

#### 参考文献

- [1] M. Drozdowski, "Scheduling multiprocessor tasks: An overview," *European Journal of Operational Research*, 1996.
- [2] P. F. Dutot, G. Mounie, and D. Trystram, "Scheduling parallel tasks approximation algorithms," *Handbook of Scheduling: Algorithms, Models and Performance Analysis*, 2004.
- [3] Y. Liu, L. Meng, I. Taniguchi and H. Tomiyama, "Novel list scheduling strategies for data parallelism task graphs," *International Journal on Networking and Computing*, 2014.
- [4] H. Yang and S. Ha, "ILP based data parallel multi-task mapping/scheduling technique for MPSoC," *International SoC Design Conference*, 2008.
- [5] H. Yang and S. Ha, "Pipelined data parallel task mapping/scheduling technique for MPSoC," *Design Automation and Test in Europe (DATE)*, pp.69-74, 2009.
- [6] J. Sun, N. Guan, Y. Wang, Q. Deng, P. Zeng and W. Yi, "Feasibility of fork-join real-time task graph models: Hardness and algorithms," ACM Transactions on Embedded Computing Systems (TECS), 2016.
- [7] K. Lakshmanan, S. Kato, and R. Rajkumar, "Scheduling parallel real-time tasks on multi-core processors." *IEEE Real-Time Systems Symposium*, 2010.
- [8] C. Chen and C. Chu, "A 3.42-Approximation algorithm for scheduling malleable tasks under precedence constraints," *IEEE Transactions on Parallel and Distributed Systems*, 2013.
- [9] H. Nishikawa, K.Shimada, I. Taniguchi, and H.Tomiyama, "Scheduling of malleable fork-join tasks with constraint programming" in *Proc of International Synposium on Computing* and Networking (CANDAR), 2018.
- [10] K. Shimada, I. Taniguchi, and H. Tomiyama, "Communication-aware scheduling for malleable tasks," in *Proc.* of International Conference on Platcon Technology and Service, 2019.
- [11] R. P. Dick, D. L. Rhodes, and W. Wolf, "TGFF: task graph for free," *International Workshop on Hardware/Software Codesign*, 1998.