

B_010

PC クラスタ上での階層統合型粗粒度タスク並列処理の MPI 実装手法

MPI Implementation Scheme for Layer-Unified Coarse Grain Task Parallel Processing on a PC Cluster

吉田 明正†

Akimasa Yoshida

1 はじめに

マルチプロセッサシステム上での並列処理手法としては、従来よりループ並列化技術 [1] が広く用いられているが、最近ではループやサブルーチンレベルの粗粒度タスク並列処理 [2] が有効と考えられている。また、粗粒度タスク並列処理で用いられている階層型マクロタスクグラフを利用しつつ、対象プログラム中の異なる階層にまたがった粗粒度タスク間並列性を最大限に利用する階層統合型粗粒度タスク並列処理 [3] が提案されている。

本稿では、この階層統合型粗粒度タスク並列処理を、普及している PC クラスタにおいて MPI により実装する方法を提案する。また、データ転送オーバーヘッドを軽減するためのデータローカライゼーション技術も取り入れている。PC クラスタ上で行った性能評価の結果、提案手法の有効性が確認されている。

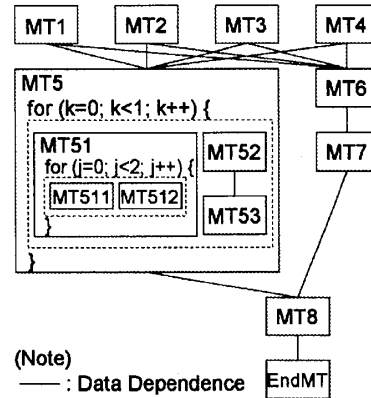
2 階層統合型実行制御を伴う粗粒度タスク並列処理

粗粒度タスク並列処理 [2] では、階層的にループやサブルーチン等の粗粒度タスク間の並列性を抽出し、粗粒度タスク (マクロタスク) をプロセッサあるいはプロセッサグループに割り当て並列処理する方式である。

粗粒度タスク並列処理による実行では、まず、プログラム (全体を第 0 階層マクロタスクとする) を第 1 階層マクロタスク (MT) に分割する。マクロタスクは、擬似代入文ブロック (基本ブロック)、繰り返しブロック (ループ)、あるいは、サブルーチンブロックの 3 種類から構成される。また、第 L 階層マクロタスク内部において、第 $(L+1)$ 階層マクロタスクを定義する。

マクロタスク生成後、各階層のマクロタスク間の制御フローとデータ依存を解析し、階層型マクロフローグラフ [2] を生成する。次に、制御依存とデータ依存を考慮したマクロタスク間並列性を最大限に引き出すために、各マクロタスクの最早実行可能条件 [2] を解析する。最早実行可能条件は、制御依存とデータ依存を考慮したマクロタスク間の並列性を最大限に表しており、マクロタスクの実行制御に用いられる。例えば、図 1 の MT5 の最早実行可能条件は、 $1 \wedge 2 \wedge 3 \wedge 4$ と求められ、MT5 は MT1~MT4 の実行終了後に実行可能となることを表している。階層統合型実行制御を伴う粗粒度タスク並列処理では、各階層ごとに求めた最早実行可能条件に、階層開始マクロタスクを導入し各種条件を変換する [3]。

階層統合型実行制御によるマクロタスクスケジューリングでは、全ての階層のマクロタスクが統一的に取り扱われ、それぞれの PE (グルーピングなし) に割り当てられて実行される。例えば、図 1 の階層型マクロタスクグラフの場合、図 2 に示すように、MT1~MT8 の第 1 階層マクロタスク、MT5 内部の MT51~MT53 の第 2 階層マクロタスク、MT51 内部の MT511~MT512 の第 3 階層マクロタスクを統一的に取り扱い、それらを各プロセッサに割り当てて実行する。この場合、第 1 階



(Note)

—: Data Dependence

図 1 階層型マクロタスクグラフ。

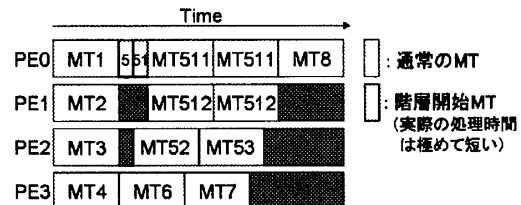


図 2 階層統合型粗粒度タスク並列処理の実行イメージ。

層から第 3 階層までの並列性 (例えば MT511, MT512, MT52, MT6 の間の並列性) が同時に利用されており、実行時間が大幅に短縮されることが分かる。

3 階層統合型粗粒度タスク並列処理の MPI 実装方法

本章では、PC クラスタ上で階層統合型粗粒度タスク並列処理を MPI により実装する方法を述べる。本手法では集中型ダイナミックスケジューリング方式を採用しているため、PC クラスタ上の 1PE をスケジューリング処理用とし、他 PE をマクロタスク実行用とする。

スケジューリング処理用 PE におけるスケジューリングの基本手順を以下に示す。

- (i) 階層統合型レディマクロタスクキューから絶対 CP 長 [3] の大きい MT_y を取り出し、アイドル状態の PE_j に対してマクロタスク番号 y を送信 (MPI_Send) する。
- (ii) MT_y で参照される共有データを、スケジューリング処理用 PE から、 MT_y を実行する PE_j に送信 (MPI_Send) する。
- (iii) あるマクロタスク MT_z の実行が PE_k で終了した場合、そのマクロタスクを実行した PE の番号 k を受信 (MPI_Recv) する。
- (iv) MT_z で定義された共有データを PE_k から受信 (MPI_Recv) する。
- (v) MT_z の終了状態を設定する。
- (vi) 新たに実行可能になるマクロタスクを階層統合型レディマクロタスクキューに投入する。

†東邦大学理学部情報科学科

Department of Information Science, Toho University

(vii) EndMTが終了してない間は (i) に戻る。EndMTが終了した際には、マクロタスク実行用 PE に終了信号を送信 (MPI_Send) する。

次に、マクロタスク実行用 PE (PE_j) におけるマクロタスク実行の基本手順を以下に示す。

- (i) スケジューリング処理用 PE から送信されたマクロタスク番号 y を受信 (MPI_Recv) する。
- (ii) MT_y で参照される共有データを、スケジューリング処理用 PE から受信 (MPI_Recv) する。
- (iii) MT_y の実行コードを実行する。
- (iv) MT_y の実行を終了した PE の番号 j を、スケジューリング処理用 PE に送信 (MPI_Send) する。
- (v) MT_y で定義された共有データを、スケジューリング用 PE に送信 (MPI_Send) する。
- (vi) EndMT の終了信号を受信 (MPI_Recv) していない間は (i) に戻る。

なお、関連研究としては、PC クラスタ上でマクロデータフロー処理を実現するために、データ到達条件 [4] を用いる方法が提案されているが、階層統合型実行制御やデータローカライゼーションは対象とされていない。

4 データローカライゼーションによる PE 間データ転送最小化

本手法では、ダイナミックスケジューリングによりマクロタスクを割り当てており、マクロタスク間共有データは、スケジューリング処理用 PE において管理する方式をとる。このため、マクロタスク間データ転送は、スケジューリング処理用 PE を経由して行われることになる。

本手法では、このようなマクロタスク間データ転送を最小化するために、パーシャルスタティック割当てを用いたデータローカライゼーション手法 [5] を導入し、同一階層のマクロタスク間データ転送のみならず、異なる階層のマクロタスク間データ転送に対しても、できるだけ PE 上のローカルメモリを介して行う。

例えば、図 3 のプログラムにおけるデータローカライゼーション適用部分は、(MT1, MT711), (MT2, MT712), (MT3, MT713), (MT4, MT72, MT75), (MT5, MT73, MT76), (MT6, MT74, MT77) となる。これらの各適用部分のマクロタスクは、パーシャルスタティック割当てを伴うダイナミックスケジューリングにより同一 PE に割り当てられ、スケジューリング処理用 PE を介したデータ転送は軽減される。

5 MPI 実装による階層統合型粗粒度タスク並列処理の性能評価

性能評価に用いた PC クラスタは、7 台の PC を 1000BASE-T のイーサネットに接続した構成となっている。各 PC の仕様は、PentiumIII: 933MHz, メモリ: 256MB, OS: VineLinux3.2 (Kernel-2.4.31), MPI: MPICH2 (Version-1.03) となっている。

階層統合型粗粒度タスク並列処理を、図 3 の 3 階層の評価用プログラムに適用し、PC クラスタ上で性能評価を行う。このプログラムの各マクロタスクは、基本的に、`for (i=2; i<10000; i++) { b[i]=b[i-1]+b[i-2]+a[i]; }` のような形状をしており、各マクロタスク間のデータ転送数は 10000 個となっている。

本性能評価では、PC クラスタの 1PE をスケジューリング処理用とし、他 PE をマクロタスク実行用とした並列プログラムを MPI を用いて作成した。PC クラスタ上での実行結果は図 4 の通りであり、6PE 実行において、データローカライゼーション無 (集中データ配置) の場

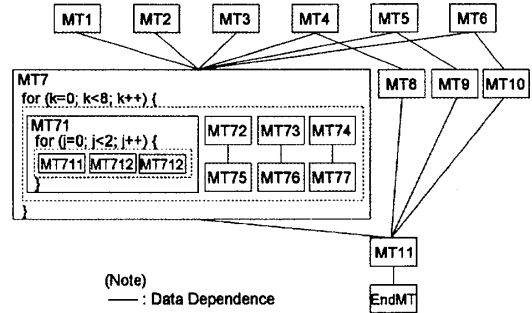


図 3 評価用プログラムのマクロタスクグラフ。

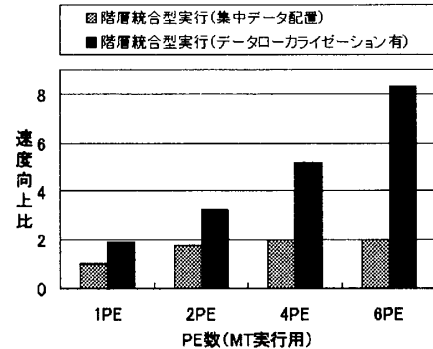


図 4 PC クラスタ上での階層統合型実行。

合 2.00 倍、データローカライゼーション有の場合 8.33 倍の速度向上が得られている。

6 おわりに

本稿では、PC クラスタ上で階層統合型粗粒度タスク並列処理を実現する方法を提案した。本手法により、ダイナミックスケジューリングを用いた階層統合型粗粒度タスク並列処理を MPI 実装できるようになり、全階層にまたがった粗粒度タスク間並列性を利用することが可能となる。

また、データローカライゼーション技術を導入することにより、データ転送オーバーヘッドの最小化も実現している。PC クラスタ上での性能評価の結果からも、提案する MPI 実装手法の有効性が確認されている。今後の課題としては、ベンチマークプログラムを用いて性能評価を行うことがあげられる。

参考文献

- [1] M. Wolfe. High performance compilers for parallel computing. Addison-Wesley Publishing Company, 1996.
- [2] 笠原博徳, 小幡元樹, 石坂一久. 共有メモリマルチプロセッサシステム上での粗粒度タスク並列処理. 情報処理学会論文誌, Vol. 42, No. 4, 2001.
- [3] 吉田明正. 粗粒度タスク並列処理のための階層統合型実行制御手法. 情報処理学会論文誌, Vol. 45, No. 12, 2004.
- [4] 本多弘樹, 上田哲平, 深川保, 弓場敏嗣. 分散メモリシステム上でのマクロデータフロー処理のためのデータ到達条件. 情報処理学会論文誌, Vol. 43, No. SIG 6(HPS 5), 2002.
- [5] 吉田明正, 越塚健一, 岡本雅巳, 笠原博徳. 階層型粗粒度並列処理における同一階層内ループ間データローカライゼーション手法. 情報処理学会論文誌, Vol. 40, No. 5, 1999.