

マルチグレイン並列処理におけるデータローカライゼーションのための近細粒度タスクスケジューリング

吉田 明正† 尾形 航† 岡本 雅巳† 合田 憲人† 笠原 博徳†
†早稲田大学 理工学部

1 はじめに

マルチプロセッサシステム上での Fortran プログラムの並列処理手法として、粗粒度並列処理(マクロデータフロー処理 [4]), 中粒度並列処理(ループ並列化), 近細粒度並列処理を階層的に組み合わせたマルチグレイン並列処理 [6] が提案されている。従来のマルチグレイン並列処理では、マクロタスク(粗粒度タスク)間で共有されるデータを集中型共有メモリ上に置き、マクロタスク間のデータ授受も集中型共有メモリを介して行なう方法がとられていた。しかし、このような方式では、集中型共有メモリを介したデータ転送オーバーヘッドが大きくなるという問題が生じる。

この問題点を解決するためには、プロセッサ内のローカルメモリを利用したマクロタスク間データ授受が重要となる。ローカルメモリの有効利用に関しては、単一ループを対象とした Array Privatization 法 [1] や、複数ループ間で PE 間通信の最小化を目指した静的なデータ分割・配置法 [2][3] が提案されている。しかし、これらの方法は、中粒度並列処理のみを対象としているため、実行時に処理とデータを動的に配置し複数粒度で並列処理を行なうマルチグレイン並列処理には適用できない。

一方、マクロデータフロー処理においては、ローカルメモリ経由でマクロタスク間データ授受を行う方法(データローカライゼーション手法) [5] が提案されている。しかし、この方法ではマクロタスク内部を複数プロセッサ上で並列処理するマルチグレイン並列処理に適用すると、MT 間で一部のデータしかローカルメモリを介して授受できないという問題点があった。

本稿では、このような問題点を解決するために、マクロタスク間データ転送を考慮した近細粒度タスクスケジューリング手法を提案する。

2 マルチグレイン並列処理

2.1 マクロデータフロー処理

マクロデータフローコンパイルーション手法 [4] [6] では、Fortran プログラムを、BPA(擬似代入文ブロック), RB(繰り返しブロック), SB(サブルーチンブロック)の3種類のマクロタスク(MT)に分割する。但し、粗粒度並列性を利用し、かつ、データローカリティを高めるため、マクロタスクの分割・融合を行い、マクロタスクの粒度の調整を行なう [4]。

コンパイラは、次に、BPA, RB, SB等のマクロタスク間のコントロールフローとデータフローを解析し、マクロフローグラフ(MFG)[6]を生成する。この後、コントロール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 [6] を解析する。この各マクロタスクの最早実行可能条件は、マクロタスクグラフ(MTG)[6]と呼ぶ無サイクル有向グラフで表すことができる。

最後に、コンパイラは、マクロタスクを実行時にプロセッサクラスタ(PC)に割り当てるためのダイナミックスケジューリングコードを生成する。

2.2 中粒度並列処理(ループ並列化)

前述のマクロデータフロー処理により PC に割り当てられたマクロタスクが、Doall 可能 RB である場合、PC 内部の PE によって中粒度すなわちイタレーションレベルで並列処理される。

2.3 近細粒度並列処理

ループ並列化不可の RB のボディ部に存在する基本ブロック、及び、ループ外部の BPA では、基本ブロック(BPA)がステートメントレベルの近細粒度タスクに分割され、PC 内の PE 上で並列処理される [6]。

近細粒度タスクの PC 内 PE への割り当てには、データ転送を考慮したヒューリスティックアルゴリズムである CP/DT/MISF 法, DT/CP 法 [6]、あるいは、提案する MT 間データ転送を考慮した近細粒度タスクスケジューリング手法を採用する。

3 MT 間データ転送を考慮した近細粒度タスクスケジューリング

本節では、データローカライゼーションの適用されるマクロタスク間で、PE 間データ転送を最小化するための、近細粒度タスクスケジューリング手法について述べる。

3.1 データ転送評価用擬似タスクの生成

本タスクスケジューリング手法では、データローカライゼーションを行なうマクロタスク集合から、データ転送評価用の擬似タスク(以後、単に擬似タスクと呼ぶ)を生成する。ここで、擬似タスクとは、MT 間データ転送を考慮して MT 内部の近細粒度・中粒度タスクをスケジューリングするために使用するものであり、実際にプロセッサ上で処理するものではない。ここでは、以下の3種類の擬似タスクを定義する。

(1) MT が Doall ループの場合、Doall ループのインデックス範囲を PC(プロセッサクラスタ)内の PE 数個に均等に分割し、この分割された各部分 Doall ループ(中粒度タスク集合)を擬似タスクと定義する。例えば、図 1(a)の部分プログラム(プログラム中から MT 間データローカライゼーション適用範囲として自動抽出される)中の Doall ループ RB3 は、PC が 2PE で構成される場合、図 1(b)に示すように、擬似タスク L3.1 と L3.2 に分割される。

(2) MT が基本ブロック(BPA)の場合には、図 1(b)の BPA1 のように、基本ブロック内の各ステートメントを擬似タスクと定義する。

(3) MT がシーケンシャルループの場合には、ループ内部の各ステートメント群(ボディ部の各ステートメントをループイタレーション回数分展開したもので、ループディストリビューションにより分割されたループと似たイメージ)を擬似タスクと定義する。但し、データローカライズされる MT 集合内に中粒度並列処理される MT が含まれている場合には、擬似タスク

ク生成前に、シーケンシャルループのインデックス範囲を、PC内のPE数個に均等に分割しておく。これは、シーケンシャルループの近細粒度タスクのスケジューリング(タスクのプロセッサへの割当て)をループ制御変数の値によって変更できるようにするためである。例えば、図1(a)のシーケンシャルループRB2は、インデックス範囲で図1(b)のRB2aとRB2bに分割され、その後、RB2aとRB2b内部で図1(b)の[]に示した擬似タスクに分割される。

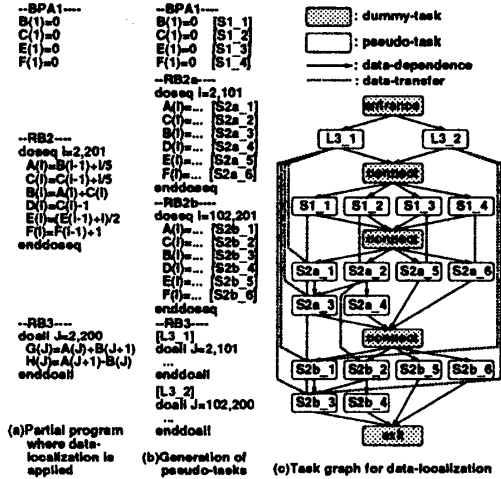


図1: データローカライゼーション用タスクグラフの例
3.2 データローカライゼーション用タスク

グラフの生成

次に、データローカライゼーションの対象となる複数MT間のデータ転送を表したデータローカライゼーション用タスクグラフを生成する。

具体的には、まず、各MTごとに、擬似タスク間のデータ依存(図1(c)において矢印付きの実線)を表したタスクグラフを作る。次に、データローカライゼーションを行なうMT集合中で、中粒度並列処理の適用されるMT、近細粒度並列処理の適用されるMTの順に、各MTのタスクグラフを接続する。このとき、各MTのタスクグラフの入口及び出口には、ダミーの入口ノード、接続ノード、出口ノード(図1(c)では *entrance*, *connect*, *exit*)を用いる。

その後、MT間及びシーケンシャルループ(MT)のイタレーション間において、擬似タスク間のデータ転送エッジ(図1(c)において太点線)を付加し、データローカライゼーション用タスクグラフを生成する。この際、データ依存エッジ及びデータ転送エッジは、データ転送時間の重み(1データの転送時間 t_{trans} × 転送データ数)をもたせる。例えば、図1(c)において、擬似タスク $S2a.1$ と $L3.1$ の間では、配列 $A(2) \sim A(101)$ のデータ転送が必要となるので、データ転送時間は $t_{trans} \times 100$ となる。なお、前述のように、このデータローカライゼーション用タスクグラフは、データローカライゼーションのためのMT内近細粒度タスクのスケジューリングのみ使用されるもので、MT間の実行順序制約は表していない。

3.3 データ転送評価用擬似タスクのスケジューリング

最後に、データローカライゼーション用タスクグラフ上の擬似タスクをPC内のPEにスケジューリングする。このとき、近細粒度タスクスケジューリングにおけるデータ割当ての基準となる部分 *Doall* ループで構成

された擬似タスクは、ループ制御変数の初期値の小さい部分 *Doall* ループから順に、PC内の番号の小さいPEに割り当てられる。一方、ステートメント(群)で構成された擬似タスクは、DT/CP(Data Transfer/Critical Path)法のスケジューリングアルゴリズムにより、データローカライゼーション用タスクグラフ上の擬似タスク間データ転送時間を考慮して、PC内のPEに割り当てる。このような擬似タスクのスケジューリングにより、各MT内の擬似タスクに対応する近細粒度タスクまたは中粒度タスクのスケジューリングが求められる。

4 OSCAR 上での性能評価

本性能評価では、図1(a)のプログラムに、提案手法を適用し、OSCARシミュレータ上で実行した結果について述べる。OSCARは、ローカルメモリ(LM)と分散型共有メモリ(DSM)を持つ32ビットRISCプロセッサ(PE)を、集中型共有メモリ(CSM)に3本のバスで接続した分散/集中共有メモリ型マルチプロセッサシステムである。

図1(a)のプログラムを1PE上でシーケンシャル実行すると、処理時間は10.04[ms]であった。次に、4PEを用いて従来の *Doall* 処理により実行すると、シーケンシャルループ等を並列処理できないため、処理時間は8.16[ms]でありほとんど短縮されない。それに対して、4PEからなる1PCを用いて中粒度並列処理と近細粒度並列処理を適用すると、処理時間は4.36[ms]に短縮される。

さらに、提案する近細粒度タスクスケジューリングを行ない、データローカライゼーション手法を適用した場合には、中粒度並列処理を適用する *Doall* ループと近細粒度並列処理を適用するシーケンシャルループとの間でLM経由データ授受が行われ、1PC(4PE)上での処理時間は3.39[ms]に短縮された。この場合、処理時間はCM経由データ授受の場合に比べて22.2%短縮されている。

5 まとめ

本稿では、マルチグレイン並列処理において、各粒度の並列性を最大限に利用しつつ、複数MT間でのデータローカライゼーションを効果的に実現するための、近細粒度タスクスケジューリング手法を提案した。また、OSCARシミュレータ上での性能評価の結果より、提案手法の有効性が確かめられた。

本研究の一部は、文部省科学研究費(特別研究員奨励費 053760, 一般研究 (b)05452354, 一般研究 (c)05680284)により行なわれた。

参考文献

- [1] P.Tu, D.Padua: "Automatic Array Privatization", 6th Annual Workshop on Languages and Compilers for Parallel Computing (1993).
- [2] M.Gupta, P.Banerjee: "Demonstration of Automatic Data Partitioning Techniques for Parallelizing Compilers on Multicomputers", IEEE Trans. on Parallel and Distributed Systems, Vol.3, No.2, pp.179-193 (1992).
- [3] J.M.Anderson, M.S.Lam: "Global Optimizations for Parallelism and Locality on Scalable Parallel Machines", Proc. of the SIGPLAN '93 Conference on Programming Language Design and Implementation, pp.112-125 (1993).
- [4] 笠原, 合田, 吉田, 岡本, 本多: "Fortran マクロデータフロー処理のマクロタスク生成手法", 信学論(D-I), Vol. J75-D-I, No.8 (1992-08).
- [5] 吉田, 前田, 尾形, 笠原: "マクロデータフロー処理におけるデータローカライゼーション手法", 情報学論, Vol.35, No.9 (1994-09).
- [6] H.Kasahara, H.Honda, A.Mogi, A.Ogura, K.Fujiwara, S.Narita: "Multi-Grain Parallelizing Compilation Scheme for OSCAR", 4th Workshop on Language and Compilers for Parallel Computing (1991).