

LC-004

# マトリックス型超並列プロセッサのための処理分割手法の提案

## Task Partitioning Method for Massively Parallel Processor Based on Matrix Architecture

中村 仁† 田中 浩明† 坂主 圭史† 今井 正治† 武内 良典†

### 1. はじめに

近年、携帯電話やデジタルカメラなどの組み込みシステムにおいて、大量の演算を必要とするマルチメディアアプリケーションの実行が増加している。これに伴い、組み込みシステムに求められる処理能力への要求は日々高まってきている。一方、多様なマルチメディアアプリケーションの登場とともに、さまざまなマルチメディアデータの規格も次々と登場している。したがって、今後の組み込みシステムでは高い処理能力とともにさまざまな規格に対応可能な柔軟性を併せ持つ必要がある。従来はこういったさまざまなマルチメディア処理に対応した専用ハードウェアを組み込みシステムに搭載することによって処理能力を向上させていた。しかしながら、専用ハードウェアを用いる場合、特定の演算に対して高い処理性能を発揮することはできるが、新しい規格などへの柔軟性に欠ける。また、柔軟性の観点から、組み込みシステムに汎用プロセッサを用いてマルチメディアアプリケーションを実行する方法もあるが、これではあまり高い処理性能が期待できない。

そこで、高い処理性能と柔軟性を併せ持つプログラマブルデバイスとして、株式会社ルネサステクノロジ社よりマトリックス型超並列プロセッサ(MTA)が提案されている。MTAでは細粒度の2入力プロセッシングエレメント(PE)とその左右に配置された2つのメモリを1ラインとし、命令メモリに格納された命令に従って1,024ライン同時実行演算を行うことが可能である[1][3]。MTAは200MHz動作時にIntel Pentium 4 3.4GHz相当以上の処理能力を発揮することが報告されており[2]、同一処理の繰り返しが要求されるマルチメディアアプリケーションに非常に適している。

このMTAの処理能力を有効に使用するためには、並列性の高い繰り返し処理をアプリケーションプログラムから抽出し、適切にMTAで演算する必要がある。

そこで本研究では、アプリケーションプログラムに含まれている繰り返し処理を対象に、実行サイクル数を最小化するCPU/MTA分割手法について提案する。

本稿の構成は次のとおりである。まず第2章で関連研究について述べ、第3章でMTAの構成、実行モデルとMTA利用時の問題点を説明し、次の第4章で提案するシステムの概要、MTAを用いることによって削減されるサイクル数の算出手法、MTAで処理を行うループの選択手法について説明し、第5章で評価実験、第6章でまとめと今後の課題について述べる。

### 2. 関連研究

本研究と関連する研究としてR.Deepaらの手法[7]がある。[7]では、ヘテロジニアスなマルチプロセッサシステムにおけるタスクスケジューリング手法を提案しており、タス

クの依存関係と各タスクの各プロセッサごとの実行コストを元にタスクスケジューリングを行う。

本研究で対象とする、ホストCPUとMTAの2つのプロセッサを含むシステムはヘテロジニアスなマルチプロセッサシステムであるといえる。しかし、MTAはSIMD処理のみに特化したプロセッサであり、データメモリ量、命令メモリ量の一般的なプロセッサにはみられない制約を持つため、既存のマルチプロセッサシステムにおけるタスクスケジューリング手法をそのまま適用して処理の分割を行うことはできない。したがって本研究では、MTAの特徴を考慮し、MTAに適した処理分割手法を提案する。

### 3. MTA

本章ではまずMTAの構成について説明し、次にMTAの実行モデルについて説明する。そして最後にMTAの利用における問題を説明する。

#### 3.1 MTAの構成

MTAは2ビットの演算器を1,000を超える並列度で同時に動作させることが可能な、非常に細粒度のSIMDプロセッサであり、高い柔軟性と高い処理性能を両立している。本節では、MTAの構成について説明する。

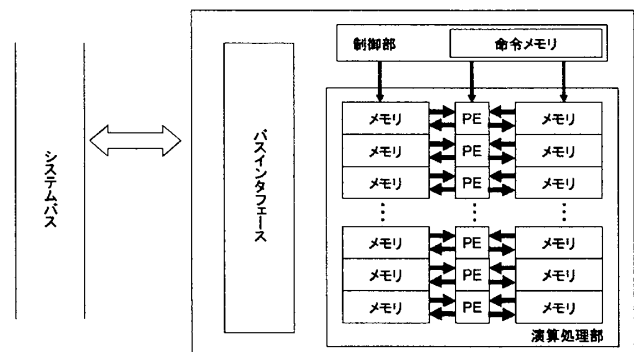


図1: MTAの構成

MTAは図1に示すように、演算処理部、制御部、バスインタフェースの部分からなる。

演算処理部はプロセッシングエレメント(PE)と呼ばれる小規模の演算器が縦に並び、その両翼にデータメモリとして用いられるメモリが配置されている。演算処理部では1,000を超える並列度で演算が実行される。

制御部はPEの処理内容を決定するための命令メモリを内蔵する。制御部の命令メモリの内容を変更することによって、演算処理部で実行される演算の内容を変化させることができる。

バスインタフェースは外部のバスと通信するためのインタフェースである。MTAはホストCPUのアクセラレータとしての利用を想定しており、このバスインタフェースを

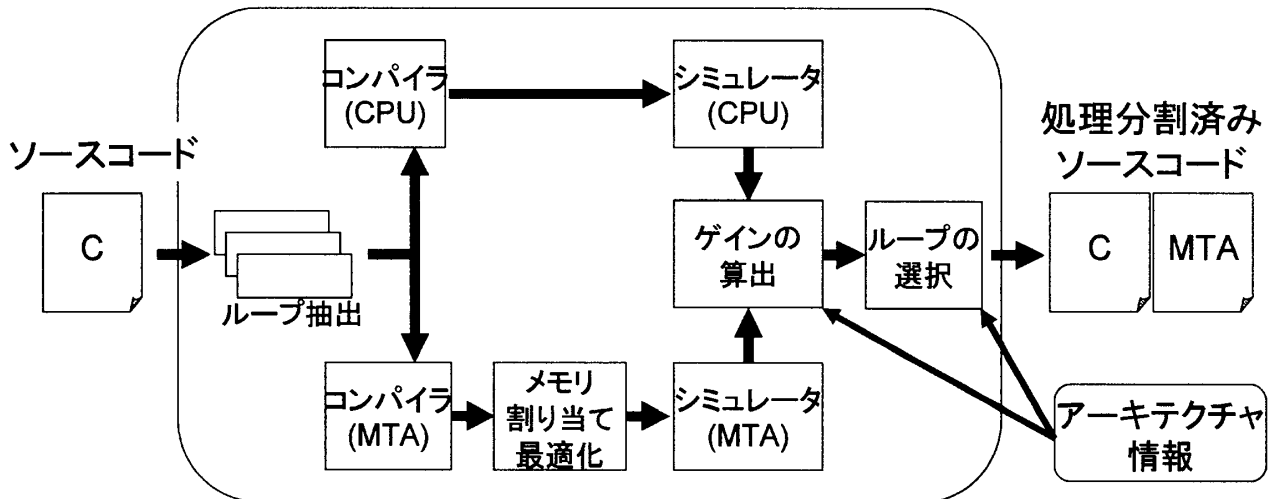


図2: CPU/MTA 分割手法の概要

通して処理対象データを送受信する。

また、PE数、データメモリ量、命令メモリ量などが異なるさまざまなバリエーションのMTAが存在する。

### 3.2 MTAの実行モデル

本節では、MTAの実行モデルについて述べる。

MTAではSIMD型演算器の特徴を生かして、並列に展開したループを処理する。そしてMTAで処理を行う場合、MTAを次の3つのフェーズで動作させる。

1. 入力フェーズ：メインメモリからMTAのデータメモリヘータを書き込む
2. 演算フェーズ：制御部の命令メモリの内容にしたがって処理を実行する
3. 出力フェーズ：MTAのデータメモリからメインメモリヘータを書き出す

例えば、MTA内のPE数を1,024とし、2,048回のループをMTAで処理する場合、まずループ1,024回分のデータをMTAのデータメモリに入力し、処理を実行し、処理されたデータを出力する。続けて、残りのループ1,024回分のデータについても同様に入力、処理、出力を行う。

MTAを用いて処理を高速化させるためには、データの入出力やPEの特性などを考慮して、どのような処理をホストCPUからMTAへ切り出すかが重要となる。以下にMTAで処理を切り出す場合に考慮しなければならない点について述べる。

- MTAで実行する処理の内容やデータ量によって削減されるサイクル数は大きく異なるため、データの入出力も考慮しながら、サイクル数削減効果の大きい処理をホストCPUからMTAへ切り出さなければならない。また、MTAで処理を実行している間は命令メモリを書き換えることができない。したがって、MTA内の命令メモリ容量内でサイクル数削減効果の大きい処理をホストCPUからMTAへ切り出さなければならない。
- MTA内のPEは2ビットで処理を実行する。したがって、単体のPEはホストCPU内の演算器と比

較すると性能が劣るため、並列度の高い処理をMTAで実行する必要がある。また、MTAはその構造上、異なる処理を並列に実行できないため、並列度の高い同一の処理をホストCPUからMTAへ切り出さなければならない。

- MTAのPE数、データメモリ量、命令メモリ量、そしてMTAと外部メモリ間のデータの転送速度は利用するMTAおよびMTAを含むシステムによって異なる。そのため、利用するシステムによって、サイクル数削減効果やMTAで処理可能な命令数が変化する。したがって、選択したMTAおよびMTAを含むシステムに応じて、削減されるサイクル数を最大化するように、処理をホストCPUからMTAへ切り出さなければならない。

以上のことから本稿では、対象アプリケーションを解析し、実行サイクル数を最小化するようなホストCPUからMTAへ処理を切り出す手法(CPU/MTA分割手法)を提案する。

## 4. CPU/MTA 分割手法

図2に提案するCPU/MTA分割手法の概要を示す。提案手法は、MTAを含む組み込みシステム上で実行する対象のアプリケーションのC言語ソースコードを入力とし、ホストCPU実行部分のC言語プログラムとMTA実行部分の命令のソースコードを出力する。

### 4.1. CPU/MTA 分割手法の概要

以下にCPU/MTA分割手法の処理の概要を示す。

1. 入力されたC言語のソースコードからループを抽出する。
2. 抽出したループをホストCPUで処理する場合に必要なサイクル数を見積もる。
  - I. ループをホストCPUのコンパイラでコンパイルする。
  - II. 回路シミュレータを用いて、ホストCPUで処理を行う場合に必要となるサイクル数を見積もる。

3. 抽出したループを MTA で処理する場合に必要なとなるデータ量, 命令メモリ量, 演算フェーズで必要となるサイクル数を見積もる.
  - I. ループを MTA のコンパイラでコンパイルし, 初期 MTA コードを生成する.
  - II. MTA メモリ最適化アルゴリズム[4]を用いて MTA コードのメモリ最適化を行い, 最適化された MTA コードを得る.
  - III. 最適化された MTA コードを MTA のシミュレータを用いて, MTA で処理を行う場合に必要となるデータ量, 命令メモリ量, 演算フェーズで必要となるサイクル数を見積もる.
4. 2, 3 で得た情報とアーキテクチャ情報により, 各ループを MTA で実行する際に削減されるサイクル数(削減量)をそれぞれ計算する.
5. 命令メモリの範囲内で削減量を最大化するように MTA で処理するループを選択する.
6. MTA で処理するループ部分を書き換えたホスト CPU の C 言語のソースコードと MTA に切り出されたループ部分の MTA コードを出力する.

提案手法中の, ホスト CPU のコンパイラと回路シミュレータ, MTA のコンパイラと MTA のシミュレータは既存のものを使用する.

MTA はその演算モデルの特徴により, 演算に使用する変数が, PE を挟んで両翼に配置されているメモリのどちらに配置されているによって, 演算にかかる実行サイクル数が異なってくる[4]. 提案手法では, 我々が別途開発した, メモリ割り当て最適化手法[4]を用いて, MTA コンパイラが主力したコードを最適化して, MTA での実行サイクル数を算出している.

次に, 提案手法中の削減量の算出手法とループの選択手法を説明する.

#### 4.2 削減量の計算

提案手法では以下の情報を用いて, ループを MTA で実行した場合の実効サイクル数の削減量を求める. ここで, 入力されたソースコードに含まれるループの集合を  $L = \{l_i | i=1, \dots, n\}$  とする.

- MTA に固有のパラメタ
  - $M$ : MTA の命令メモリ量 (byte)
- MTA を含むシステムに固有のパラメタ
  - $R$ : メインメモリ, MTA 内のデータメモリ間のデータの転送速度(byte/cycle)
- ループごとのパラメタ
  - $C_{mta,i}$ : ループ  $l_i$  を MTA で処理するために必要なサイクル数 (cycle)
  - $C_{cpu,i}$ : ループ  $l_i$  をホスト CPU で処理するために必要なサイクル数 (cycle)
  - $E_i$ : ループ  $l_i$  の演算フェーズに必要なサイクル数 (cycle)
  - $D_{in,i}$ : ループ  $l_i$  で処理されるためにホスト CPU から MTA へ転送されるデータ量 (byte)
  - $D_{out,i}$ : ループ  $l_i$  で処理し, MTA からホスト CPU へ転送されるデータ量 (byte)

- $l_i$ : ループ  $l_i$  を MTA で処理するために必要な命令メモリ量 (byte)

MTA の命令メモリ量  $M$  は使用する MTA が決定されることで決まる. また, メインメモリと MTA 間のデータ転送速度も, MTA を使用するシステムのシステムバスの性能によって決定される.

ループ  $l_i$  をホスト CPU で処理した場合に必要なサイクル数  $C_{cpu,i}$  は, ループ  $l_i$  をホスト CPU のコンパイラでコンパイルし, 回路シミュレータを用いることで求まる.

ループ  $l_i$  で処理されるためにホスト CPU から MTA へ転送されるデータ量  $D_{in,i}$  と, ループ  $l_i$  で処理して MTA からホスト CPU へ転送されるデータ量  $D_{out,i}$  は, ループ  $l_i$  を MTA コンパイラでコンパイルすることで求められる.

また, ループ  $l_i$  を MTA で処理するために必要な命令メモリ量  $l_i$  も MTA コンパイラでコンパイルすることで求められる.

3.2 節で述べたように, ループ  $l_i$  を MTA で処理するために必要なサイクル数  $C_{mta,i}$  は, 入力フェーズ, 演算フェーズ, 出力フェーズに必要なサイクル数の和となる. 演算フェーズに必要なサイクル数  $E_i$  は MTA のシミュレータを用いることで求められる. 入力フェーズ, 出力フェーズに必要なサイクル数はメインメモリ, MTA のデータメモリ間のデータの入出力量をそれぞれ転送速度で割ることによって求められる. したがってループ  $l_i$  を MTA で処理する場合に必要な総サイクル数  $C_{mta,i}$  は以下のように計算される.

$$C_{mta,i} = \frac{D_{in,i}}{R} + E_i + \frac{D_{out,i}}{R}$$

#### 4.3 ループの選択

4.2 節で求めた各ループの  $C_{cpu,i}$  と  $C_{mta,i}$  を用いて, どのループを MTA で処理するかを決定する. 特にループを MTA で処理した場合に必要な命令メモリ量によっては, 全てのループを MTA で処理出来ないことがある. そこで本節では, 命令メモリを制約とし, 実行サイクル数を最小にするようなループの MTA への切り出し問題(CPU/MTA 分割問題)を定式化する.

ループをホスト CPU で処理するか MTA で処理するかを表す変数  $x_i$  を定義する. ループ  $l_i$  をホスト CPU で処理する場合  $x_i=0$  とし,  $l_i$  を MTA で処理する場合  $x_i=1$  とする.

前述したように, MTA の命令メモリ量以上のループは処理出来ないため, 全てのループを MTA で処理できないこともある. したがって, MTA で処理するループの命令コード量の総和が命令メモリ量を超えないことが制約条件となり, 以下の式となる.

$$\text{制約: } \sum_{i=1}^n l_i \cdot x_i \leq M$$

アプリケーションの実行サイクル数を最小にするように, 上記制約を満たす範囲で, どのループを MTA で処理するかを決定する必要がある. つまり, 上記制約を満たす範囲で, MTA で処理することで削減されるサイクル数を最大化する必要がある. ここで, ループを MTA で処理することで削減されるサイクル数は,  $l_i$  をホスト CPU で処理した場合に必要なサイクル数  $C_{cpu,i}$  と  $l_i$  を MTA で処理した場合に必要なサ

イクル数 $C_{mta,i}$ との差である。以上より、目的関数は次の式となる。

$$\text{目的関数：最大化：} \sum_{i=1}^n (C_{cpu,i} - C_{mta,i}) \cdot x_i$$

以上の議論をまとめ、本研究であつかう問題をCPU/MTA分割問題として以下の様に定式化する。

CPU/MTA分割問題

- 入力
  - ・ ループ集合  $L = \{l_i; i = 1, \dots, n\}$
  - ・ ループ $l_i$ をMTAで処理するために必要なサイクル数 $C_{mta,i}$
  - ・ ループ $l_i$ をホストCPUで処理するために必要なサイクル数 $C_{cpu,i}$
  - ・ MTAの命令メモリ量 $M$
- 出力
  - ・  $x_i$
- 制約
  - ・  $\sum_{i=1}^n I_i \cdot x_i \leq M$
- 目的
  - ・ 最大化：  $\sum_{i=1}^n (C_{cpu,i} - C_{mta,i}) \cdot x_i$

ここで、CPU/MTA分割問題の、命令メモリ量 $M$ をナップザックの大きさ、各ループの命令メモリ量を大きさ、各ループの実行サイクルの削減量を価値とみなすと、CPU/MTA分割問題はナップザック問題として扱うことができる。

ナップザック問題はNP-完全であることが知られており[5]、色々な解法が提案されているが、本稿では、分枝限定法を用いてこの問題を解く。

分枝限定法は全探索の1種で、これ以上の探索が不要であると判断された時点で枝刈りを行う手法である。

## 5. 評価実験

提案手法の有効性を確認するために評価実験を行った。実験には、ホストCPUとしてエイシップ・ソリューションズ社のBrownie\_std32プロセッサを採用したシステムを想定し、ホストCPU、MTA、バスの動作周波数はともに同じ動作周波数とした。また、対象アプリケーションとしてOpenJPEG library[6]の5-3ウェーブレット変換のソースコードから、48種のループを含むソースコードを用い、解の導出には、分枝限定法によりCPU/MTA分割問題を解くプログラムを作成して用いた。

ループ部分に関してサイクル数を比較した結果を表1に示す。

	サイクル数	削減率
CPUのみ使用	184,728	-
提案手法	61,484	67%

表1：サイクル数の比較

提案手法を用いてMTAを利用すると、MTAを用いない場合と比較してサイクル数が約67%削減され、この例題の場合、約3倍高速に実行可能であることが示された。また解は48ループ中18ループをMTAに割り当てるものとなり、探索時間はPentium4 3.0GHz、1GB RAMの環境で1秒未満であった。

以上より、提案手法の有効性を確認した。

## 6. まとめと今後の課題

本稿では、CPU/MTA分割手法を提案し、その有効性を示した。

今後の課題としては、データの依存関係を考慮することにより、MTA内のメモリにデータを入出力するサイクル数を削減することが挙げられる。これは入力されたソースコードにおいて、データの依存関係を考慮して処理の順序を適切に入れ替えることによって実現できると考えられる。しかしながら、処理の連続する異なるループ間で、同じ変数を利用している場合においても、ただデータの入出力処理を削除するのではなく、データメモリの配置を考慮して、配置が異なる場合にはデータメモリを再配置するといった処理が必要となる。

### 文 献

- [1] M. Nakajima et.al., "A 40GOPS 250mW Massively Parallel Processor Based on Matrix Architecture," Proceedings of ISSCC, pp. 410-411, 2006.
- [2] 大野隆行, 山崎博之, 飯田全広, 久我守弘, 末吉敏則, "Matrix Processing Engine のメディア処理アプリケーションによる性能評価," 信学技報, Vol. 106, No. 49, RECONF2006-6, pp. 31-36, 2006.
- [3] 中田 清, 中島雅美, 野田英行, 谷崎哲志, 行天隆幸, "40GOPS 250mW マトリックス型超並列プロセッサの開発 モバイル向け SoC にも組み込み可能な超高速プロセッサ," 信学技報, vol. 106, no. 71, ICD2006-25, pp. 19-23, 2006.
- [4] 小橋 晶, 谷口 一徹, 坂主 圭史, 武内 良典, 今井 正治, 中田 清, "マトリックス型超並列プロセッサにおける変数のメモリ割り当て最適化手法," 信学技報, 2007.
- [5] Michael R. Garey, David S. Johnson, "Computers and intractability. A guide to the theory of NP-completeness", San Francisco, W. H. Freeman, 1979
- [6] OpenJPEG library ([www.openjpeg.org](http://www.openjpeg.org))
- [7] R. Deepa, T. Srinivasan, D. Doreen, H. Miriam, "An Efficient Task Scheduling Technique in Heterogeneous Systems Using Self-Adaptive Selection-Based Genetic Algorithm", PARELEC, p343-348, 2006.