

可変構造型並列計算機の並列オペレーティング・システム — プロセス管理とメモリ管理 —

恒富, 邦彦
九州大学大学院総合理工学研究科情報システム学専攻

草野, 和寛
九州大学大学院総合理工学研究科情報システム学専攻

福田, 晃
九州大学大学院総合理工学研究科情報システム学専攻

村上, 和彰
九州大学大学院総合理工学研究科情報システム学専攻

他

<https://doi.org/10.15017/17198>

出版情報 : 九州大学大学院総合理工学報告. 12 (2), pp.227-235, 1990-09-01. 九州大学大学院総合理工学研究科
バージョン :
権利関係 :

可変構造型並列計算機の並列オペレーティング・システム —プロセス管理とメモリ管理—

恒 富 邦 彦*・草 野 和 寛*・福 田 晃**
村 上 和 彰**・富 田 眞 治**
(平成2年5月31日 受理)

A Parallel Operating System for The Reconfigurable Parallel Processor —Process Management and Memory Management—

Kunihiko TUNEDOMI, Kazuhiro KUSANO, Akira FUKUDA,
Kazuaki MURAKAMI and Shinji TOMITA

A reconfigurable parallel processor under development at Kyushu University is a MIMD-type multipurpose multiprocessor system. This system can be tailored to a broad range of applications. Its operating system under development considers it as a tightly coupled multiprocessor. Its kernel is based on the message oriented method to implement data-parallel processing. The operating system provides users with functions of the light weight process (thread) as a parallel processing model for tightly coupled multiprocessor system. The overhead of the thread switching can be reduced by our scheduling scheme where processors are allocated and deallocated to each process. The spin lock mechanism with time-out facility is employed to realize more flexible scheduling than coscheduling. We evaluate two methods where address mapping table is copied and it is not copied. As a result, the latter is employed. We propose a method of lazy ceation of stack to have fast access to a stack.

1. はじめに

我々が現在開発中の「可変構造型並列計算機 (Reconfigurable Parallel Processor)」^{1),2)} は、汎用/多目的マルチプロセッサ・システムであり、128台のプロセッシング・エレメント (PE) が128×128のクロスバー網により接続されている。本システムは、種々の並列処理形態に対して計算機構成を適応させることを目的としており、相互結合網およびメモリにダイナミック・アーキテクチャを採用している。これによりシステムは Fig. 1 に示すアドレス空間を持ち、密結合型/疎結合型マルチプロセッサ双方をハードウェア・レベルで実現できる。このシステム上に構築する並列オペレーティング・システム (OS)^{3),4)} は、アドレス空間を制御管理するメモリ管理と、プロセス管理が必要となる。本稿では、このOSのプロセス管理とメモリ管理について述べる。

2. 並列オペレーティング・システムの設計方針

本システム開発の主な目的を以下に示す。

(1) ハードウェア・アーキテクチャと各種並列処理問題との親和性の検討

(2) 解くべき並列処理問題にハードウェアを動的に適合させることによるマルチプロセッサ・システムの応用分野の拡大。

(1), (2) を遂行するためには、並列/分散 OS の研究を含めた広範囲な並列処理の研究が不可欠である。

(2) を実現するためには、まず(1)の検討すなわち、ハードウェア・アーキテクチャと各種並列処理問題との親和性を定量的に検討する必要があると考える。そこで、我々は(2)を実現するOSの開発に先立って、(1)を対象としたOSを開発する。

現在開発中の、密結合型メモリ・モデルを基本とした並列OSの設計方針を以下に示す。

(1) 密結合方式に適した並列処理モデルをユーザに提供する。

*情報システム学専攻修士課程

**情報システム学専攻

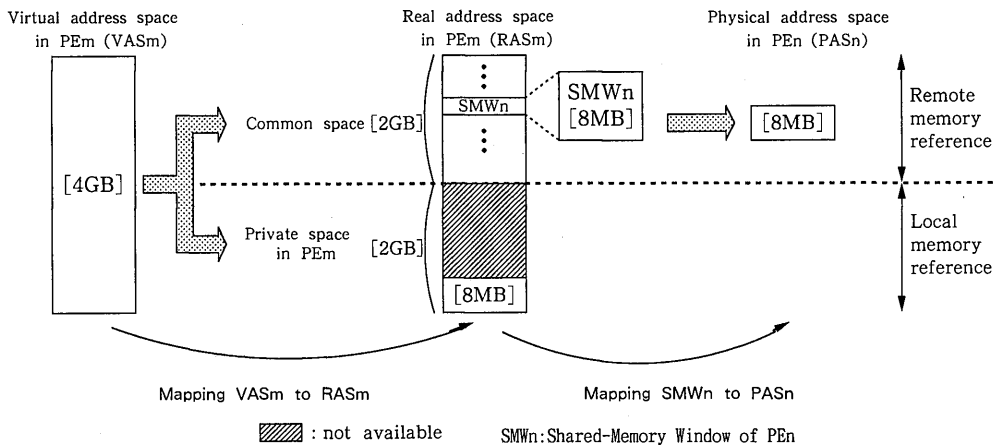


Fig. 1 Addressing scheme.

ユーザプログラムをマルチプロセッサ上で中粒度の並列処理（スレッド）に分割して並列実行させる機能を提供する。詳しくは3章で述べる。

(2) カーネル内部を可能な限りデータパラレル化する（負荷分散）。

ある1つのシステム・コールの実行軌跡である制御フローに着眼した場合、従来の密結合型マルチプロセッサ用OSでは以下の方式が採られている。

(a) 制御フローが1つ、あるいは複数のプロセッサで実行されるにしても、いずれの場合でも1つの逐次的な制御フローとして実行される。

(b) (a)の方式をさらに進めて1つのシステム・コールを機能分割し、並列実行できるところは並列実行する（コントロールパラレル化）。

我々は、これらの方式をさらに進めて複数のスレッドの同時生成などデータパラレル化できるところはデータパラレル化する。すなわち、オーバーヘッドの問題に対処する必要があるが、密結合型マルチプロセッサ用OSの内部処理のデータパラレル化の可能性を追求する。

これを実現する方式として、メッセージ指向のOSを構築する。すなわち、Fig. 2に示すようにデータごとの処理に各機能を細分化し、実行要求をメッセージとして各処理のキューにつなぐメッセージ方式により、カーネルの内部処理を並列化する。

(3) 多様なシステム・コールを提供する。

ユーザが多様な並列処理問題を扱えるようにするため、多様なシステム・コールを提供する。また、ソフ

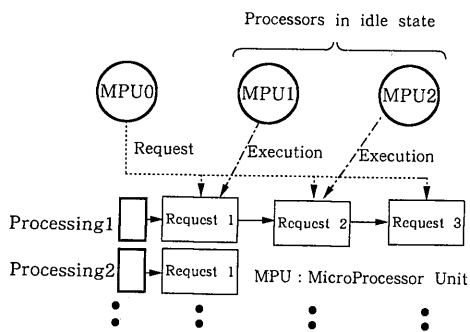


Fig. 2 Data parallel in kernel.

トウェア資産の継承という観点から、一般的なOSであるUNIXと互換性をもたせる。

3. プロセス管理

3.1 並列処理モデル

UNIXをはじめとする従来のシングルプロセッサ上のOSでは、実行の単位としてプロセスを考え、同時に複数のプロセスが並行に実行できる環境を提供し、またスループットを向上するために時分割スケジューリングを行っている。従ってマルチプロセッサ上では、これらのプロセスの並列実行が可能となり、各プロセスがお互いに独立である場合には処理が高速になる。しかし、ある仕事内が並列に実行可能であり、互いに協調動作を行う処理が存在する場合に、これをプロセスに分割して実行させるとプロセス生成、通信および同期のオーバーヘッドが大きくなるなどの問題が存在し、

プロセスモデルによる記述は並列処理には適しているとはいえない。

我々は、ユーザのプロセス内部を並列実行可能な処理の流れ（スレッド）に分割して実行することによりこの問題を解決する。同一プロセス内の複数のスレッドは仮想アドレス空間を共有しており、通信や同期処理を1つのユーザ空間のみで処理することが可能となる。プロセスは、1つの仮想アドレス空間とその空間に割り当てられた資源および複数のスレッドからなる。スレッドは、プログラム・カウンタ、レジスタ、スタック領域および実行に関する情報からなる。

3.2 メッセージ指向によるプロセス管理構築

各スレッドにはスレッド・コントロール・ブロック（THCB）とスタック領域が割り当てられ、実行時の情報などをTHCBに格納している。従って、システム・コール `thread-1fork(n)` によって一度に多量のスレッドを生成すると、THCBも多量に必要となり、これを初期化する処理が逐次的になされる限りボトルネックとなる。

この問題を解決するために、設計方針で述べた様にメッセージ方式によりプロセス管理を構築し、処理を並列化する。スレッド生成のシステム・コールを実行しTHCBの初期設定を行うときには、THCBの初期化処理サーバに、生成されるスレッドごとに要求をつなぐ。アイドル状態のプロセッサはキュー検索を行うループを実行しており、キューから要求を1つ取り出し実行する。複数のアイドル状態のプロセッサが存在すれば、各々のスレッドのTHCBが並列に初期化される。このように複数スレッド生成の並列化を実行できるが、このとき、メッセージ操作のオーバーヘッドが問題となる。すなわち、並列化とそれに伴うオーバーヘッドのトレードオフの問題である。これに関しては、複数スレッド生成のシステム・コールを受け付けたプロセッサでまとまった数のスレッドを生成し、残りのスレッド生成の依頼をメッセージにして他のプロセスで処理する方法などが考えられる。

3.3 スケジューリング

3.3.1 スケジューリングの要件

マルチプロセッサ上のプロセス管理で重要なことは、プロセッサのスケジューリングである。スケジューリングによってシステム全体のスループットが左右され、また通信や同期のためのオーバーヘッドを極力回避することができるからである。スケジューリングには、仕

事のプロセッサへの時間的、空間的マッピングを実行以前に行う静的スケジューリングと、システムの状況により実行時に行う動的スケジューリングがある。我々は両方について研究しているが、ここでは、動的スケジューリングについて述べる。通信パターン、通信コストを考慮した静的スケジューリングについては文献5)を参照されたい。

今 n 個のプロセッサが結合されているマルチプロセッサを考える。UNIXのマルチプロセスの時分割スケジューリングをそのままマルチプロセッサに適用すると、実行可能なプロセスのキュー（プロセス・キュー）を用意し、キューを検索して優先度の高いものから n 個を選んで各プロセッサに割り付け、プロセッサはそれぞれの優先度に対し割り当てられる実行時間（タイム・スライス）だけ処理を行うこととなる。タイム・スライスを使いきると、プロセスの優先度を再計算し、最高の優先度を持つ実行可能なプロセスを選択して再びユーザ・プロセスを実行する。

スレッドは実行の単位であるから、上記のスケジューリングにおいてプロセスをスレッドに置き換えてそのまま適用すると、実行可能なスレッドのキュー（スレッド・キュー）がシステム全体で共有され、そのキューの優先度の高いスレッドが n 個実行されるスケジューリングとなる。しかし、スレッドが同一プロセス内処理の並列性の記述のために用いられ、複数のスレッドが特に密接な関係にあることから、上記の方式は次の問題が生じる。

(1) 同期によるオーバーヘッドの増大

同一プロセス内のスレッド間の同期としては、ユーザ空間において共有のロック変数を持ちこれがアンロックされるまで繰り返しリードするスピン・ロックが有効である。これによりスレッド間の同期がカーネルの処理を必要としないために高速に行える可能性がある。しかし、同一プロセス内の複数のスレッドが生産者と消費者の関係にある場合には、両者が共に実行されていることが保証されていなければ能率が上がらない。生産者があるデータの生産を行う以前に消費者が励起されて実行を開始しても、データへの排他制御のためのスピン・ロックを繰り返すのみでタイム・スライスを消費し、資源の無駄使いとなる。上記のスケジューリングでは、同一プロセス内のスレッドがばらばらに実行されることになるので、この問題が生じる。

(2) 並列性を実現する上でのオーバーヘッドの増大

上記のスケジューリングでは、スレッドのタイムアウト毎にスレッドの優先度を各プロセッサが再計算することになるので、システム内に存在するプロセッサの台数より遙かに大きい並列性をスレッドを用いて実行する場合、スレッドの切り換え回数が増大する。

(3) スレッド切り換えによるオーバーヘッドの増大

各プロセッサのキャッシュ上には、そのスレッドが用いた情報がキャッシングされている。従って、過去に実行されたスレッドが再び励起される場合およびスレッドを新たに実行する場合には、そのスレッドと同一プロセス内のスレッドの処理を実行していたプロセッサに割り当てられると、キャッシュ内のデータの再利用が可能となり処理が高速となる。上記の方式では、プロセス単位のプロセッサの割当てを考慮していないので、新たに実行されるスレッドが異なるプロセス内のスレッドとなる可能性が大である。このとき、キャッシュのミスヒットが頻発すると共に仮想空間の切り換えのオーバーヘッドが生ずる。

以上述べた問題点(1)～(3)を解決するスケジューリングを考える。

3.3.2 コスケジューリング

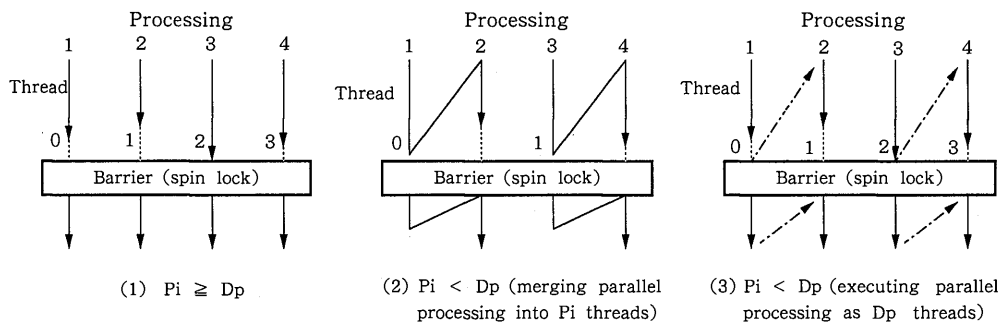
上記の問題を解決するために Medusa⁹⁾ をはじめ多くの並列 OS がコスケジューリングを行っている。コスケジューリングとは、プロセス・キューにプロセス内に存在するレディ状態のスレッド数すなわち並列処理数(並列度)をもたせ、このプロセスを実行に移すときは並列度だけプロセッサを割り当てて処理が並列に実行されることを保証する方式である。これによって、(1)の問題は解決できる。しかし、並列度がシス

テム内のプロセッサ数を越える場合、また割当て可能なプロセッサ数が並列度より少ない場合などは柔軟に対処することができない。

これを回避する方法として、ユーザによって指定された並列度よりアイドル状態プロセッサ数が少ない場合に、並列度をアイドル状態のプロセッサ数へ再分割する方法が挙げられる (Fig. 3 (2))。すなわち、システム・コールにより現在のアイドル状態のプロセッサ数を確認し、アイドル・プロセッサ数だけに並列度をマージし、それらをスレッドとして実行する。しかしこの方式では、一度並列度をアイドル・プロセッサ数分だけにマージさせてしまうと、新しくアイドル状態のプロセッサが発生した場合でも、並列度を再び上げることができない。

3.3.3 本 OS におけるスケジューリング

上記のようにコスケジューリングは動的に変化するシステムのアイドル状態のプロセッサ数を反映できない欠点をもつ。これを解決するスケジューリング方式として、本 OS では、2段階スケジューリングを実現する。すなわち、第1段階目のスケジュールは、UNIX と同様にプロセッサの実行時間からプロセス優先度を決定し、優先度の高いプロセスからレディ状態スレッド数に応じてプロセッサを割り当てる。コスケジューリングと違い、割当て可能なプロセッサ数がスレッド数より少ない場合でもプロセッサを割当てる。第2段階目のスケジュールは、プロセスに割当てられたプロセッサをプロセス内のスレッドに対し FCFS (First Come First Served) 方式により割り当てる。これにより高優先度のプロセスは全スレッドが同時に



P_i : the number of idle processors
 D_p : Degree of parallel processing

Fig. 3 Parallel processing scheme.

処理を行うが、低優先度のプロセスも、一部のスレッドが処理を進めることができる。この処理は、優先順位を持つプロセス・キューをシステム全体で1つ共有し、各プロセスごとにスレッドキューを割り当てることにより実現される。この手順を以下に示す。

1) プロセスのスケジューリング

プロセスを実行していない（アイドル状態の）プロセッサはプロセス・キューを検索し、最高優先度のプロセスを選択する。キューの要素には、プロセス id とそのプロセス内のレディ状態のスレッド数が記録されており、プロセッサが割り当てられる毎にこのスレッド数を減ずる。プロセスの処理中に新たなスレッド生成が発生する場合（Fig. 4 (1)）およびスレッドがサスペンド状態から励起されて実行待ち状態になる場合には、スレッド数に加算する。スレッド数が0となれば、その要素をキューから外す。以上の処理を行い、プロセスを選択したプロセッサは2)のスレッドのスケジューリングに移る。

2) スレッドのスケジューリング

スレッド・キューはプロセス対応に割り付けられる。プロセッサは実行するプロセスのスレッド・キューを検索し、キュー先頭のスレッドを実行する。実行中のスレッドがI/O待ちや、ページフォルトの処理などでサスペンドされた場合、再び同一プロセスのスレッドキューからレディ状態のスレッドを選択し処理する。Fig. 4 (2) に示すようにレディ状態のスレッドが存在しない場合にはじめてプロセスから解放され、1)のプロセス・スケジューリングに戻る。

この方式では、プロセッサのプロセス切り換えは、

(a) タイムアウト後に優先度の高いプロセスが発生した場合、(b) 割り込みにより優先度の高いプロセスにより横取り (preempt) される場合、および(c) プロセス内のレディ状態スレッドが存在しない場合に発生する。Fig. 5 に示すように、最初にプロセスを実行するプロセッサにタイマを設定し、タイムアウトになると自プロセッサはプロセス・キューの優先度を再計算する。優先度は、プロセッサの実行時間、初期優先度等により決定する。このとき、実行状態プロセスよりも高優先度のレディ状態プロセスが発生すると、低優先度の実行状態プロセスのプロセッサに割り込みをかけ、再スケジューリングを行わせる。これが(a) (b)の切り換え処理である。このようにタイムスライスを終了判定および優先度の再計算は、プロセスが起動されたときに最初に実行されるスレッドを実行していたプロセッサが行い、各プロセッサが行う必要がない。また、ディスパッチの対象とならないプロセスのスレッド処理はタイマを実行するプロセッサを除いては実行中

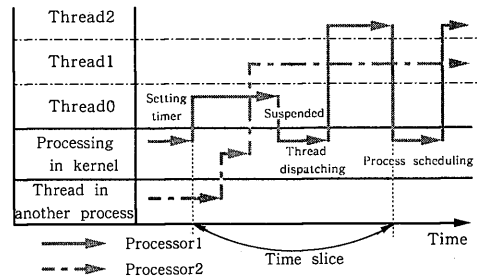


Fig. 5 Processor's behavior in two-level scheduling.

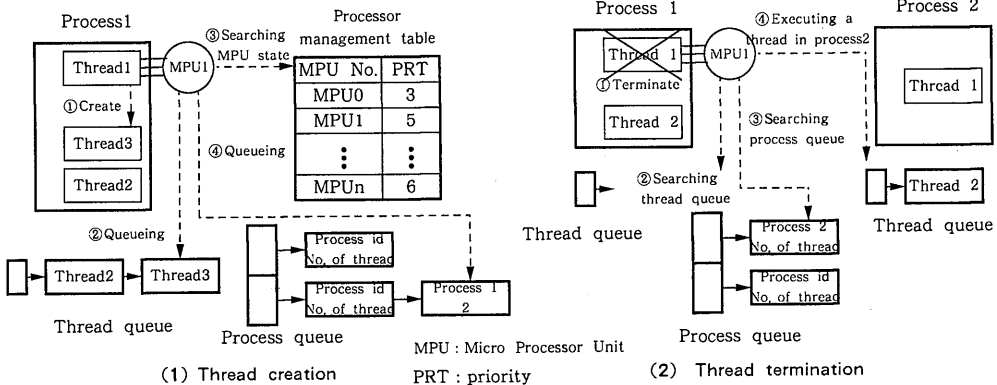


Fig. 4 Thread creation and termination.

断されることがなく、スレッドのプロセッサ間の移動も発生しない。従って問題点(2)が解決できる。

2段階のスケジューリングによりプロセッサはプロセスごとに割当てられ、タイムスライス内でスレッドの切り換えが発生した場合、同一プロセスのスレッドを優先して実行させる。これにより、仮想空間の切り換えやキャッシュヒット率の低下を抑えることができ問題点(3)が解決できる。

プロセスの優先度が高いものは、レディ状態のスレッド数だけプロセッサを割り当てられるので、スレッド数がアイドル状態のプロセッサ数よりも少ないときには全スレッドが同時に実行され、問題点(1)は回避することができる。しかし、優先度の低いものは割当てられたプロセッサ数が少なく、全スレッドが同時には実行されないことが考えられる。またスレッドのスケジューリングにおいては、実行中のスレッドがサスペンドまたは終了するまでスレッドの切り換えは発生しない。従って優先度の低いプロセスにおいては、スレッドがスピン・ロックを実行する場合には問題(1)が発生する。これを解決するためにスピン・ロックに時間制限(タイムアウト)をつけたものを用意する。スピン・ロックをある一定期間繰り返した後、ロックが解放されない場合にはスレッドをサスペンドさせ、スレッド・キューから新たにスレッドを選択し実行する (Fig. 3 (3))。

3.3.4 スケジューリングの比較

コスケジューリングと本 OS のスケジューリングとの比較を問題点(1)～(3)について行う。

問題点(1)については、コスケジューリングは、いかなるプロセスも、その内部処理を同時実行させることにより完全に解決している。一方、本 OS のスケジューリングは、優先度が高いプロセスのみ内部処理を同時に実行させる。並列度より少ないプロセッサで実行される低優先度のプロセスについては問題点(1)は解決されないで、タイムアウト付きのスピン・ロックを実現する。これにより、タイムスライスを全てスピン・ロックに費やす可能性は低くなり、問題点(1)を回避できるが、同期によるオーバーヘッドがスピン・ロックの制限時間によって左右される。制限時間が短い場合にはスレッドの切り換えが頻繁に発生し、スピン・ロックの同期の軽さをいかにすることができない。制限時間の長い場合は、全スレッドが同時に実行されていないときは、早急にスレッドの切り換えが必要である

が、タイムスライスのほとんどをスピン・ロックに費やしてしまい同期のオーバーヘッドが大きくなる。また、コスケジューリングのように同一プロセスの全スレッドの処理進行状況をそろえることは保証されないで、スピン・ロックそのものに費やされる時間も増加する。スピン・ロックの制限時間の検討は今後の課題である。

一般的に、ロックを実現する方法としてスピン・ロックとサスペンド・ロックがある。このいずれを用いるかはロックが解除されるまでの時間とコンテキスト切換えに要する時間とのトレードオフで決まる。従って、上記のようにタイムアウト付きのスピンロックで実現する方法の他に、ロック変数に関するスレッドが全て実行されている場合にはスピンロックを用い、そうでない時にはサスペンドロックを用いる方法も考えられる。この検討は今後の課題である。

次に問題点(2)について考察する。コスケジューリング、本 OS のスケジューリングとも優先度はプロセスごとに存在する。いま、実行されていたプロセスの優先度よりも高い優先度のプロセスが発生し、このプロセスにプロセッサを割当てたため、実行されていたプロセスのプロセッサ数が減少したとする。この場合、コスケジューリングでは、並列度の少ない別のプロセスを起動しなくてはならないのでキャッシュ内のデータを生かすことができない。これに対し、本 OS のスケジューリングではプロセッサ数が減っても実行を続けるのでキャッシュ内のデータを使うことができる。また、コスケジューリングよりもスレッドの切り換え回数が減少する。

問題点(3)については、両方のスケジューリングとも、プロセスが実行を続ける限り、プロセッサはプロセス間を移動しないため、キャッシュの性能を生かすことができ、解決できる。しかしコスケジューリングは、上記のようにプロセッサ数が不足するとプロセスは実行を停止するので、プロセスが実行を続ける可能性が本 OS のスケジューリングよりも低い。また、同一プロセス内のスレッドは仮想空間を共有しており、このプロセスのスレッドを実行したプロセッサがキャッシングしたデータを他のスレッドが利用できる可能性がある。従って、本 OS のスケジューリングでは、プロセッサ数が並列度より少ないときには、プロセッサが同一プロセッサ内のスレッドを優先して実行するのでキャッシュの内容を再利用することができる。コスケジューリングでは、このようなキャッシュ内の

データの再利用はできない。

4. メモリ管理

4.1 概要

本システムでは、メモリ構成のダイナミック・アーキテクチャを実現するために、ローカル/リモート・アーキテクチャを採用している。このために、Fig. 1に示したように3つのアドレス空間（仮想アドレス空間、実アドレス空間、物理アドレス空間）を導入して、2レベルのアドレス変換を行う。実アドレス空間は、全PEに共通のコモン空間と各PEの私有領域であるプライベート空間から構成されている。メモリ管理では、これら3つのアドレス空間と、2レベルのアドレス変換に用いられるアドレス変換テーブルの管理を行う必要がある。本システムに特徴的な実アドレス空間の管理において問題となるのは、仮想アドレス空間からコモン空間へのマッピングを無効化するのに大きなオーバーヘッドを伴うことである。これに関しては、文献7)を参照されたい。本節ではアドレス変換テーブルとスタックの管理について述べる。

4.2 アドレス変換テーブル

密結合型マルチプロセッサにおいて、複数プロセッサ上で同一プロセスの異なるスレッドが実行されることを考えると、同一プロセス内のスレッドは仮想アドレス空間を共有するので、アドレス変換テーブル（仮想アドレス空間から実アドレス空間への変換テーブル）へのアクセス競合が発生することが考えられる。本システムは、分散メモリ構成のローカル/リモート・アーキテクチャを採用している。このことから、スレッドの使用するアドレス変換テーブルの管理方法について、以下の2つが考えられる。

(1) 全てのスレッドが同一のアドレス変換テーブルを使用する。

(2) アクセス競合を避けるために、スレッドが実行されるプロセッサ毎に、または実行されるプロセッサの n 個毎にアドレス変換テーブルのコピーを持たせる。

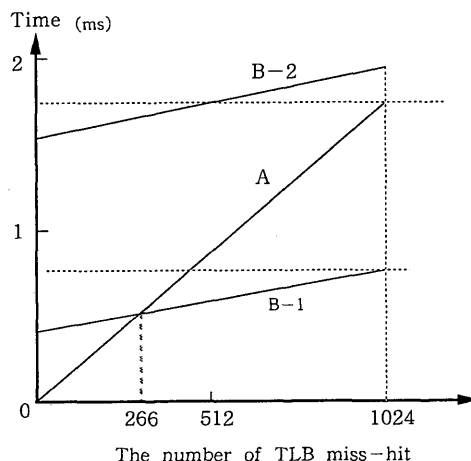
以下にこの2つの方式について考察する。

(1)の方法を用いると、アクセス競合が発生する。しかし、(2)の方法のようにアドレス変換テーブル間のコヒーレンス問題は発生しない。この時に、アドレス変換テーブルと異なるプロセッサ上で実行されているスレッドはリモート・アクセスを行うことになるが、

そのオーバーヘッドはTLB (Translation Look-aside Buffer) やキャッシュにより、かなり軽減されると考えられる。

(2)の方法を用いると、アクセスの競合は少なくなる、またはなくなると言えるが、アドレス変換テーブル間のコヒーレンスを保つことが必要になる。つまり、複数のアドレス変換テーブルの内容が一致していない状態が発生するのを防ぐ必要がある。特に、ページの入れ替え時のコヒーレンスが問題となる。例えば、あるスレッドがページフォルトを起こして、ページをメモリに読み込んで、他のアドレス変換テーブルを変更するまでの間に、それを用いているスレッドが同じページフォルトを発生させることも考えられる。この場合には、メモリ上に同じページが複数個存在するとなると共に、アドレス変換テーブル間のコヒーレンスが保てなくなる。また、アドレス変換テーブルのコピーの手間もローカル/リモート・アーキテクチャをとる本システムにおいてはかなりのオーバーヘッドとなる。

以上に述べたように、この問題は、リモート・アクセスと、アドレス変換テーブルのコピー及びコヒーレンスの処理との間のトレードオフによりどちらが良いかが決まることになる。



A : scheme (1)
B-1 : scheme (2) (one page transfer by DMA)
B-2 : scheme (2) (one page transfer by SPARC)

Fig. 6 Access times to the address mapping table.

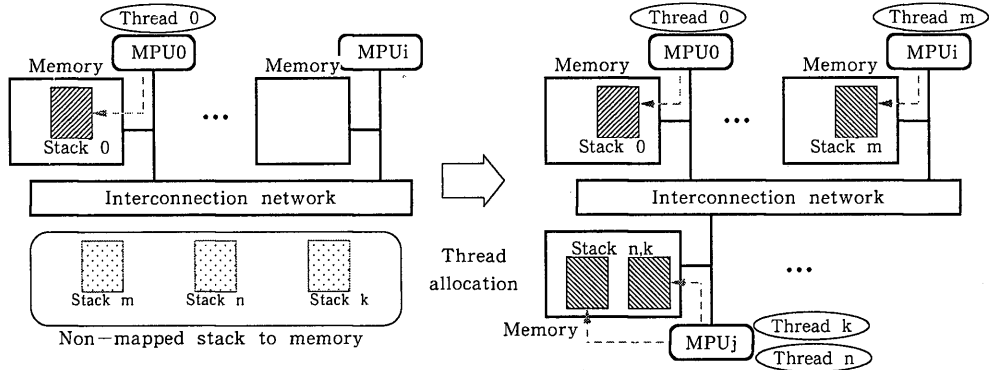


Fig. 7 Lazy creation of stack.

方式(1), (2)におけるアドレス変換テーブルへのアクセス時間を Fig. 6 に示す. Fig. 6 から分かるように, TLB のミスヒット回数が数百回までは, 方式(1)の方がアクセス時間が短い. 一般にスレッドの大きさは TLB のミスヒット回数が数百回以下となる大きさと考えられること, 及び方式(2)におけるコヒーレンス処理のオーバーヘッドを考慮すると方式(1)が優れていると言える. 以上に述べたことにより, 本システムではアドレス変換テーブルは全スレッドで同一のものを使用する. つまり, アドレス変換テーブルのコピーは作らない.

4.3 スタック領域の遅延生成

スタックにアクセスするのはそれを所有しているスレッドのみであるから, ネットワークのトラフィックを減らすためにも, スレッドが実行しているプロセッサのローカル・メモリ上にスタックを置く方がよい. しかし, スレッドが実行されるプロセッサは, スケジューリングにもよるが, スレッド生成時には判明していないことが多い. このため, スレッドが実行を開始するまで実際のスタックの生成を遅らせることができれば, スレッドが実行されるプロセッサのローカル・メモリ上にスタックをとることができると考えられる (Fig. 7). もちろん, 仮想空間上にはスレッドのスタック領域は予め確保しておく必要がある. 上記のスタック領域の遅延生成は, 実行時に発生するページフォールトにより実行プロセッサを認識し, そのプロセッサ上で実際のスタック用の物理メモリのページ割り当て, 及びスタックへの生成時に必要な情報の書き込み処理等を行うことにより実現できると考えられる. この方法を実現できれば, 1 命令で複数のスレッドを

生成する場合にも, スタックの生成にはあまりオーバーヘッドを伴わないと考えられる. この方法によらない場合は, 複数のスレッドのスタックをできるだけ異なるプロセッサ上に配置して生成処理を行う方式が, 複数を単一プロセッサ上に生成する場合よりメモリへのアクセス競合が軽減される. これはコモン空間のエントリの仮想アドレス空間への割り当て方法, つまり割り当てアルゴリズムにより行うことができる. 本システムでは, 上記のスタック領域の遅延生成の手法を用いる. この詳細な実現は現在検討中である.

5. おわりに

以上, 開発中の並列 OS のプロセス管理およびメモリ管理について述べた. 現在 Sun-4 上にマルチスレッドが実現できる環境を構築しており, 今後マルチプロセス環境を実現する予定である.

謝 辞

我々と共に開発を進めている森, 蒲池 (現(株)日本電気), 廣谷 (現(株)日本電気), 福澤 (現(株)ソニー), 岩田, 甲斐, 上野, 杉山の各氏, および日頃ご討論頂く富田研究室の皆様へ感謝致します.

参 考 文 献

- 1) K. Murakami, S. Mori, A. Fukuda, T. Sueyoshi and S. Tomita: "The Kyushu University Reconfigurable Parallel Processor-Design Philosophy and Architecture-". Proc. of IFIP 11th World Computer Congress, pp. 995-1000 (1989).
- 2) K. Murakami, S. Mori, A. Fukuda, T. Sueyoshi and S. Tomita: "The Kyushu University Reconfigurable Parallel Processor-Design of Memory and Intercommunication Architectures-". Proc. of ACM SIGARCH Int'l Conf. on Su-

- percomputing, pp. 351-360(1989).
- 3) 福田, 福澤, 廣谷, 村上, 末吉, 富田: “可変構造型並列計算機の並列/分散オペレーティング・システム”, 情報処理学会オペレーティング・システム研究会, 89-OS-43-8 (1989).
 - 4) 福澤, 草野, 恒富, 福田, 村上, 富田: “可変構造型並列計算機のオペレーティング・システム・スレッドの実現”, 情報処理学会オペレーティング・システム研究会, 89-OS-45-6(1989).
 - 5) 蒲池, 森, 村上, 福田, 富田: “統合型並列化コンパイラ・システム・ネットワークシンセシサー”, 情報処理学会第40回全国大会講演論文集, 1G-2, pp. 653-654(1990).
 - 6) John K. Ousterhout, Donald A. Scelza and Pradeep S. Sindhu: “Medusa: An Experiment in Distributed Operating System Structure”, Communications of the ACM, Vol. 25, No. 2, pp. 397-409(1980).
 - 7) 草野, 福澤, 福田, 村上, 富田: “可変構造型並列計算機のオペレーティング・システム・メモリ管理”, 情報処理学会第40回全国大会講演論文集, 6G-1, pp. 740-741(1990).