

TOP-1 オペレーティングシステム

穂積 元一

日本アイ・ビー・エム

高性能並列処理ワークステーションTOP-1の概要、および、スレッドを持ったそのOSについて解説する。TOP-1システムは80386マイクロプロセッサ10台を持つ密結合型共有メモリアイブのワークステーションで、そのオペレーティングシステムはAIXをベースとした並列処理用のOSである。TOP-1アーキテクチャの特徴、特にシステムソフトに関連する部分、および、我々のオペレーティングシステムの特徴であるスレッドの考え方に関して述べる。スレッドとは従来のプロセスに比べてスイッチングの軽いものであり、ある意味で実行コンテキストと考えられるものである。

TOP-1 Operating System

Motokazu Hozumi

IBM Japan Ltd. Tokyo Research Lab.

5-19 Sanban-cho Chiyoda-ku Tokyo Japan

This paper describes the overview of TOP-1 system which has 10 80386 micro-processors and its operating system. This operating system is based on AIX and has thread mechanism for supporting MP functions. Thread is a light weighted process and can be regarded as an executing context.

1. はじめに

我々は高速並列処理ワークステーション (TOP-1) の上にAIXをベースとした並列処理オペレーティングシステムを設計作成中である。第1版は既に稼働中である。

このオペレーティングシステムを解説するには、まずTOP-1のアーキテクチャの概要から示さなければならない。

TOP-1プロジェクトの目的とは、密結合共有メモリ型の並列処理ワークステーションを実現することで、汎用の高いコストパフォーマンスの実現と、それに伴う各種ソフトの研究にある。

2. TOP-1アーキテクチャの概要

TOP-1アーキテクチャの詳細な内容については、本年秋の情報処理学会全国大会に於て6件の発表が予定されており、そちらを参照されるのが良いが、現在まで新聞発表以外に詳しい情報がほぼ公表されていないので、オペレーティングシステムの概要を理解する上で必要な部分、および重要な特徴についてのみここでは述べる。

2.1 TOP-1システムの構成

図1に示されるとおり、TOP-1のシステム構成は、10台のインテル80386プロセッサを使用した、共有メモリ型マルチプロセッサシステムである。このシステムはローカルメモリを使用しておらず、すべてのCPUから完全に均質に共有された共有メモリのみからできている。従来はローカルメモリを使うことで、各ノードに固有なデータを割りつけ、これによって共有メモリとバスの競合を軽減するのが普通であった。しかし本システムはそのシステムソフトと相俟って汎用性を大きな特徴としているため、複雑なメモリ管理を要求するローカルメモリを使用せず、一切を共有メモリとしている。ただし共有メモリ型アーキテクチャの欠点であるバス競合を大幅に低下させるべき工夫としてスヌープキャッシュと高速バスを採用して、この問題に対処している。

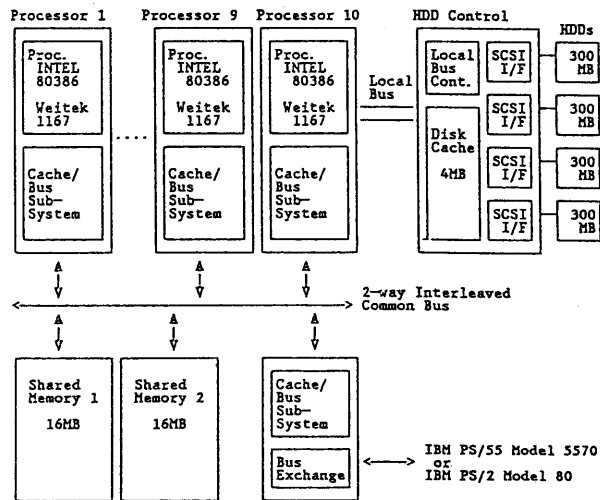


図1 TOP-1のハード構成

浮動小数点プロセッサとしては、80387より高速のWeitec1167を採用している。

又、I/Oに関しては、図1にも示されるとうりハードディスクとしてSCSIの300MBのものを4台実装しており、その他のI/Oについては弊社の製品IBM PS/55 Model 5570とI/Oカードを通じて直接接続することで行っている。

以下にTOP-1の重要な特徴であるスヌープキャッシュとMP-OSを作成する上で基本的なプロセッサ間非同期通信機構について述べる。

2.2 スヌープキャッシュ

バス及びメモリの競合を低減するために、上に述べたようにTOP-1では各CPUにキャッシュを用いている。これはデータと命令を区別するものでなく、大きさは128Kバイトのものである。

キャッシュのサイズは大きければ、それだけヒット率が高くなるが、逆にコストも高くなるわけで、このかめあいで決められるものであるが、我々は高速SRAMを用いることで128Kバイトのキャッシュを実現している。

マルチプロセッサでのキャッシュの最大の問題は同一のメモリロケーションに対して複数のコピーがキャッシュに存在してしまう点で、いかにしてこれら複数のコピーの矛盾をなくし、一致性を保証するかということである。この問題を解決するために、TOP-1ではスヌープキャッシュ方式を採用している。

我々は4つの状態を導入することで、できる限り各キャッシュ内でローカルに処理できる比率を高く保ち、共有バスに対する要求回数を低く抑える事を計っている。

2.3 プロセッサ間非同期通信機構

マルチプロセッサにおいて効率の良い並列処理を行うためには、プロセッサ間の同期機構が大切である。TOP-1システムは共有メモリ方式であるので、共有メモリ上に共有変数を置くことで、プロセッサ間の同期を計る事ができる。しかし共有変数だけでは非同期に発生するイベントをプロセッサ間で効率よく伝達できない。このためにハードとしてのプロセッサ間非同期通信機構が必要になる。

もちろん通常のセマファは80386のLOCK#信号を用いる事でTest-and-Setが実現できプロセッサ間の排他制御を行っている。

イベントドリブな通信を行うために、TOP-1では図2に示すような、Everybody Broadcast, Anybody Broadcastという2種類のプロセッサ間非同期通信機構を用意している。

Everybody Broadcastとは指定したすべてのプロセッサがメッセージを受信できたとき、成功と見なされ、一つでも受取に失敗するとこのメッセージパッシングは失敗となるものである。

Anybody Broadcastとは指定した宛先のうち一つ以上のプロセッサが受信できれば成功と見なされ、全ての宛先が受取に失敗したときに、失敗と見なされるものである。例えば、Everybody BroadcastはTLBのコンシステンシーのプログラムで用いられ、Anybody BroadcastはI/O割りこみに用いられる。

3. TOP-1オペレーティングシステム

上に述べたTOP-1システム上にAIXをベースとしてMP-OS TOP-1オペレーティングシステムを作成中であるが、まず並列処理オペレーティングシステム作成のアプローチを述べ、次にどのような研究課題があり、更にそれらをどのように解決しているかについて述べる。

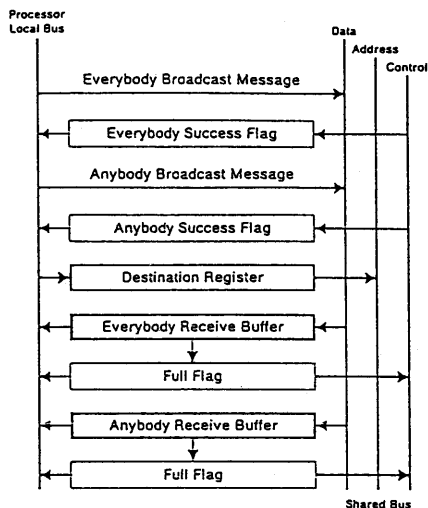


図 2: Message Passing Hardware

1. プロセッサがある I/O アドレスに出力命令を開始する。
2. キャッシュ・バス・コントローラがバスを獲得する。
3. 宛先の受信バッファが空のとき、プロセッサの出力データがバッファに書かれる。
4. バスが解放され、プロセッサの I/O 命令が終了する。

3. 1 TOP-1 オペレーティングシステムのアプローチ

TOP-1 オペレーティングシステムはシングルプロセッサ用の OS である AIX をベースとして MP-OS を作成しようとしている。そこで次のようなアプローチが考えられる。

- a. OS のカーネルをできるだけエントラント化する。
- b. OS のカーネルはどれか一つの特定の CPU でのみ実行させ、ユーザプロセスを他の CPU で並列動作させる。
- c. OS のカーネルはシリアリユーザブルのまま、任意の CPU で動作させるようにする。

b は概念が明確である分作成は容易であるが、よりよいパフォーマンスを求めると、どうしても a を実現しなければならない。ただし a は b に比べて作成は困難であり、又幾通りかの実現方法があり、将来的な並列処理に耐えられるカーネルの構造をも実現しようとする、どのようなアプローチでも良いというわけではない。

更に作成上のことを考慮に入れると、並列処理 OS のデバッグはかなり困難なものであり、初期の設計が余程完全であるか、ステップバイステップに作成するかのいずれかと思われる。我々は既存の OS の並列化を行っているので、一からの完全な設計は望めない。そこで b → c → a のアプローチで作成することとし、最終段階の a を実現する段で十分な将来的な機能分割を達成する事を目指している。

3.2 TOP-1オペレーティングシステムの研究課題

このような並列処理オペレーティングシステムを作成する上で、どのような研究課題を考えて行って要るかを述べる。

1. OSのカーネルの機能分割
2. ユーザに対する並列処理のより良いモデルの提供
3. 並列処理特有の問題の解決

1については、今後の並列オペレーティングシステムを考える上で本質的な問題であり、現時点では多くを論文にできないが、概念のみ説明する。

2については、我々は一般にスレッドと呼ばれているモデルを導入しているが、これについては幾通りかの考え方があり、そのうちの一つについて以下に説明する。

3については、TLBフラッシュの問題を含めた並列処理特有の資源管理の問題解決の事であり、知られている事も多いのでここには述べない。

3.2.1 カーネルの機能分割

カーネルのできるだけ多くの部分をリエントラントにして、並列度を上げ性能向上を計ることは、並列処理オペレーティングシステムの作成上の一つの眼目である。もちろん必要なクリティカルセクションをきちんと認識し、ロックを取る事によってこれらは達成される。しかしいかにしてデッドロックを回避するか、更には将来的な展望を考えた場合望まれるカーネルの構造との調和が問題として残る。

そこで我々はマルチサーバモデルというものを採用しており各サーバは一つのみまとまった機能単位として、サーバ単位の並列性の実現を考えている。これによってサーバの入れ替えや追加更には機能分割型のアーキテクチャにも柔軟に対処できることを目指している。

3.2.2 スレッドモデル

一般にスレッドと呼ばれているモデルのインプリメントには幾通りか存在し、もちろんそれぞれに意味が異なっている。例えばSUNのlight-weight-processはある意味のコルーチンとして実現されている。又CMUで開発されたオペレーティングシステムMACHにもスレッドモデルは採用されている。我々はスレッドを共有メモリを持つ並列計算のモデルとして実現している。この共有メモリのみならず実行環境をプロセスといい、その中の実行文脈(コンテキスト)をスレッドと言っている。

この両者の関係を図3に示す。

我々がこのようなスレッドを導入した意味は大きく言って次の2点から成る。

1. もしコルーチンのままだと、マルチプロセッサ上でスレッドによる実時間並列処理ができない。メモリを共有する'軽いプロセス'としてはコルーチンも意味があるが真の並列処理が実現できない。
2. UNIXのSystem Vにはメモリセグメントの概念があり、複数のプロセスでメモリを共有できる。ただしこの場合はプログラムコードやデータ全般が共有されるわけではなく、従ってより密な相互作用が不可能になっている。更に我々のスレッドでは全空間を図4のように共有しているので、同一プロセス内のスレッドスイッチの場合にTLBフラッシュの必要がなく効率的である。System Vのプロセスの場合は当然これをMP化するとプロセススイッチ時にTLBフラッシュが必要になる。

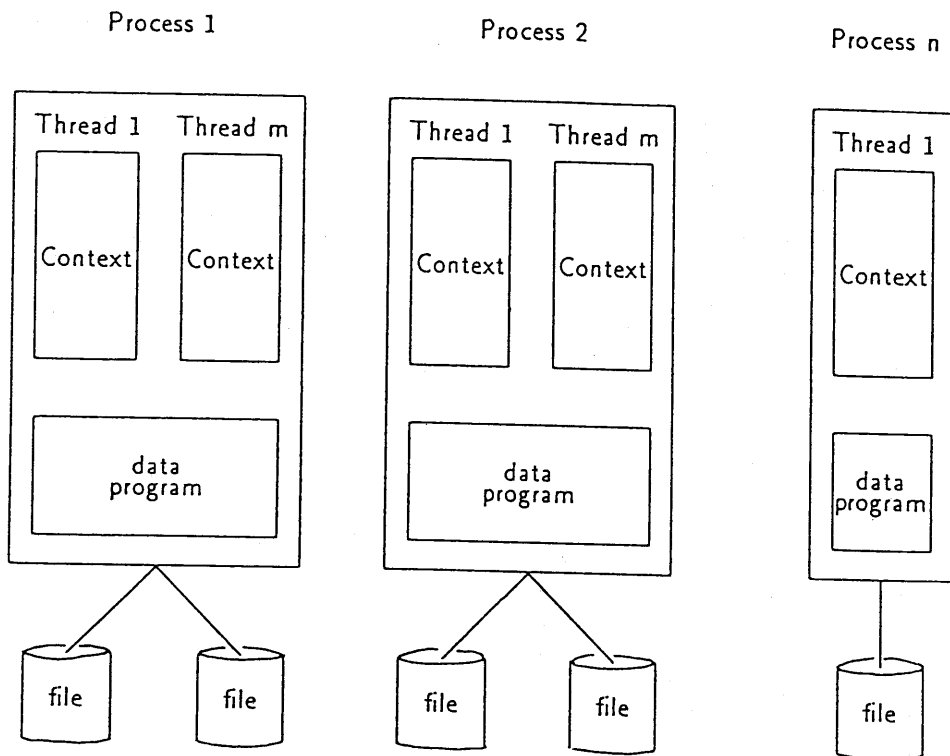


図3 プロセスとスレッド

次に同一プロセス内の各スレッドのスタックについて述べる。これは言語処理系をスレッドを用いてインプリメントする場合のスコールールと深く関係する問題である。

基本的にスタックのインプリメントには2通りの考え方がある。

- a. スタック領域は同一プロセス内の全てのスレッドで共有し、それらの間のチェインニングや保護は何もサポートしない。すべて言語処理系に任せる。
- b. ハードのセグメント機構等を用いて各フレーム毎の保護をサポートする。更に親子関係といったチェインニングもOSで行う。

aの方法は言語処理系が自由にスタックの使用法を定義できる利点があるが、保護や空フレームの回収等も自分でやらなければならない欠点がある。

- Virtual memory

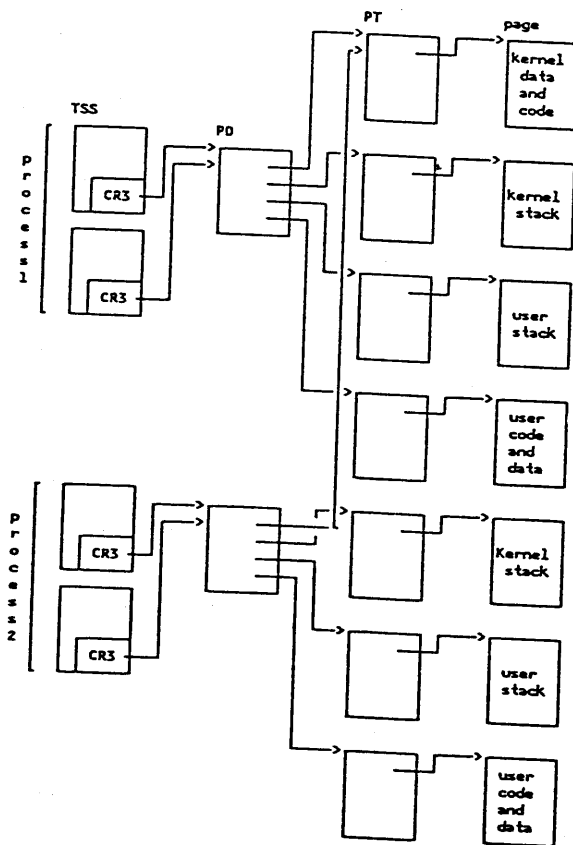


図4 プロセスとスレッドの仮想記憶機構

これに対してbでは、OSやハードである程度の事をしてくれる代わりに、自由なスコープルールや親子関係は定義できなくなる。LISPやCのような比較的自由的な言語ではaのほうが向いているといえ、Adaのような厳密な言語ではそれにあつたスタックインプリメントがbの方法で求められるだろう。将来的な互換性や言語の移植性といった事まで含めると軽々には決定できない問題といえる。

しかしTOP-1は研究用のマシンなので、言語処理系にとって自由度の高い方法aを採用している。もちろんこのほうが言語処理系の研究にとって都合が良いからである。

次にスレッドスタックと言語の変数スコープの関係について、図5に従つてもう少し詳しく眺めてみる。左がわのツリーはスレッドの生成関係を示しており、右側のボックスは言語の構造（プロシージャ、ファンクション）による変数のスコープを示している。つまり、左側でいうとスレッドAがスレッドBとCを生成し、更にBがDとEをCがFとGを生成したことを、しめしている。右側でいうと、構造BとCからは構造Aの変数は参照可能である事を、BとCでは互いに参照不可能の事を示している。又右と左で同一の名前は構造AがスレッドAで実行される事を示している。この様な場合のスレッドスタックのインプリメントは参照可能方向のチェーンと参照不可能間の保護を同時にサポートしなければならない。これには動的関係つまりスレッドの生成関係と静的関係つまり言語構造の入れ子関係の二つが必要になる。OSだけでこれを全て処理すると、一義的なインプリメントとなつて、微妙な言語仕様（LISPのSpecial Variables等）のインプリメントを困難なものとしてしまう。

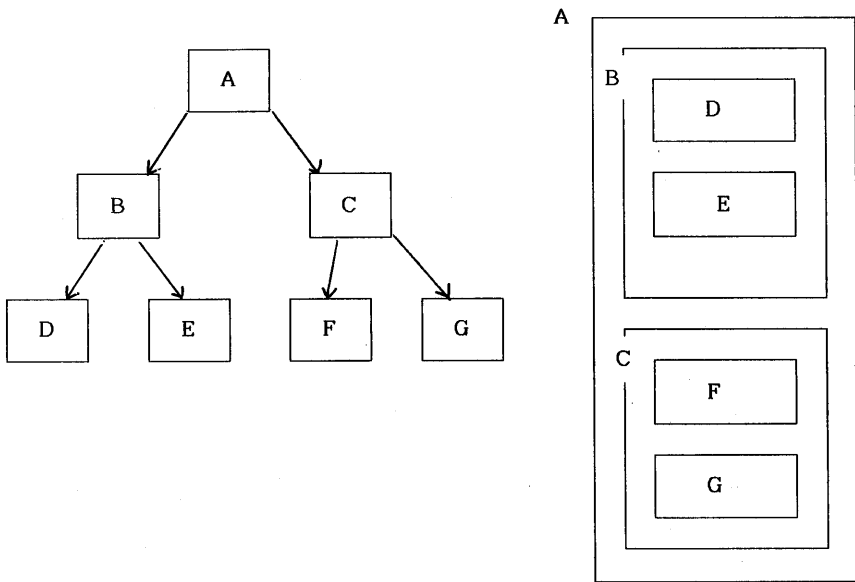


図5 スタックとフレームの関係

4. おわりに

現在TOP-1は稼働中であり、我々は新しいバージョンのMP-OSを作成中であり、同時にCommon-LISPの並列化にも取り組んでいる。いくつかの並列アプリケーションを走行させる予定である。もちろん我々のOSはAIXの完全上方互換であるので、AIXのアプリケーションはそのまま並列に走行させる事ができる。今後これら並列アプリケーションから有用なデータの収集が期待される。

5. 参考文献

- "System Service Overview" Sun Microsystem Inc.
- "MACH Kernel Interface Manual" Robert V baron, David Black et al Carnegie-Mellon Univ. 1987
- "MACH: A Basis for Future UNIX Development" Avadis Teranian Jr. Richard F. Rashid
Carnegie-Mellon Univ. 1987
- "MACH Threads and UNIX Kernel" Avadis Teranian Jr. Richard F. Rashid et al.
Carnegie-Mellon Univ. 1987