

解説



日本におけるオペレーティングシステム研究の動向

2.3 並列計算機 TOP-1 の OS†

河内谷 清久仁†

1. はじめに

TOP-1 (Tokyo research Parallel processor-1) プロジェクトは、日本アイ・ビー・エム(株)東京基礎研究所で、マルチプロセッサ・アーキテクチャの設計項目の評価および、実際の並列ハードウェアを用いての各種並列処理ソフトウェアの研究を目的として行われたプロジェクトである¹⁾。当プロジェクトは、ハードウェアの開発からオペレーティングシステム、言語サポートまでを含めたものであるが、本稿では TOP-1 上で動作する OS を中心に解説する。

図-1 に TOP-1 ハードウェアの外観を示す。そのシステム構成は図-2 のようになっている。TOP-1 は、10 台のプロセッシング・ユニット (PU) をスヌープ・キャッシュを介して共有バスで結合した、いわゆる UMA (Uniform Memory Access) 型の並列計算機である^{2)~4)}。各 PU には Intel i386DX と 128 K バイトのキャッシュが装備されており、キャッシュ間の一貫性はハードウェアによって保たれる。PU 間の同期には、バス・ロックとメッセージ割込みの機能が利用できる。TOP-1 から制御できる I/O デバイスとして 4 台の SCSI HD が用意されているが、この HD に直接アクセスできるのは、HD 制御 (HDC) カードに接続された特定の PU (図の PU 0) のみである。その他の I/O や外界へのアクセスは、フロントエンド・プロセッサとして接続されている IBM PS/55 を介して行われる。この PS/55 からは、TOP-1 上の共有メモリにアクセスしたり、各 PU との間でメッセージ割込みをかけることができる。

TOP-1 プロジェクトでは、このハードウェア上で様々な並列処理の実験を行うためのベースとして OS の開発を行った。開発にあたっては、以下の点を特に考慮している。

- 既存 OS との互換性から、Unix 系の OS を採用
- TOP-1 ハードウェアの評価機能のサポート
- OS 自身の評価機能を設け、OS の構成法を検討
- 並列プログラミングのためのサポート

次章以降では、OS 開発の際の問題点と、TOP-1 上で稼働している OS の構成について述べる。

2. OS 並列化の問題点

Unix 系の OS を並列マシンで動かす場合の最大の問題は、本来リエントラントになっていない OS カーネルをどうやって並列対応にするかという点にある。一般に Unix のカーネルは実行途中でプリエンプトされないことを前提として書かれているため、複数の PU 上で同時にカーネル・コードが実行されると、内部のデータ構造が破壊さ

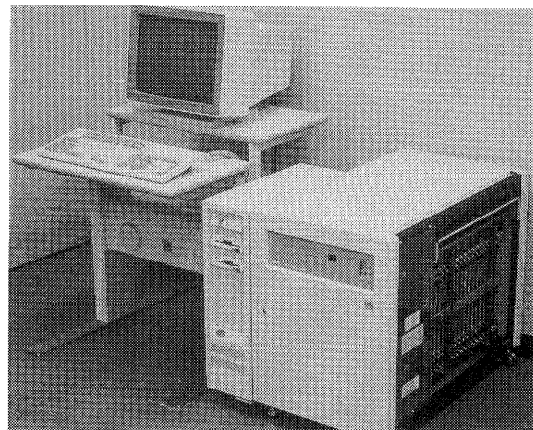


図-1

† Operating Systems on the TOP-1 Multiprocessor by Kiyokuni KAWACHIYA (IBM Research, Tokyo Research Laboratory).

‡ 日本アイ・ビー・エム(株)東京基礎研究所

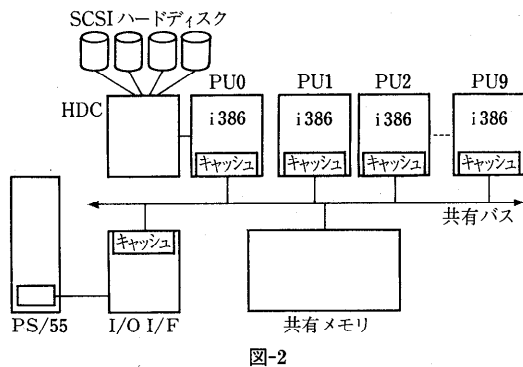


図-2

れてしまう。この問題を解決するには、各データへのアクセスに対して排他制御を行う必要がある。

Unix カーネルの並列化は、この排他制御のレベルによって、対称型とマスタ・スレーブ（非対称）型に大別できる。対称型では、カーネル内のデータ構造それぞれにロックを設け排他制御を行う。これに対し、マスタ・スレーブ型ではカーネル全体に対して排他制御を行う。複数のプロセスがカーネル・コードを実行しようとしたとき、そのうちの1つだけがカーネルに入ることを許可され、残りはカーネルの入口で待たされる。

マスタ・スレーブ型は実現が比較的容易である反面、カーネルの実行がシリアライズされるため、それがシステム全体のボトルネックになる危険性がある。対称型ではその心配はないが、複雑なカーネル・コード全体からロックすべきデータ構造を切り出し、しかもデッドロックなどが生じないように排他制御を行うには、かなり大がかりな変更が必要と考えられる。TOP-1 上では、まずマスタ・スレーブ方式の TOP-1 OS が開発された。続いて対称型マルチプロセッシングをサポートした OSF/1 TOP-1 が開発されている。

3. TOP-1 OS

TOP-1 上の OS 開発にあたってはまず、確実に動作する並列 OS を短時間で作成する必要から、マスタ・スレーブ方式で既存 Unix の並列化を行った。これが TOP 1 OS である⁵⁾。TOP-1 OS は、IBM の i386 版 Unix である AIX PS/2 をベースにしており、OS のカーネル・コードは特定の PU (KPU: Kernel PU, 図-2 の PU 0) でのみ実行される。残りの PU (UPU: User

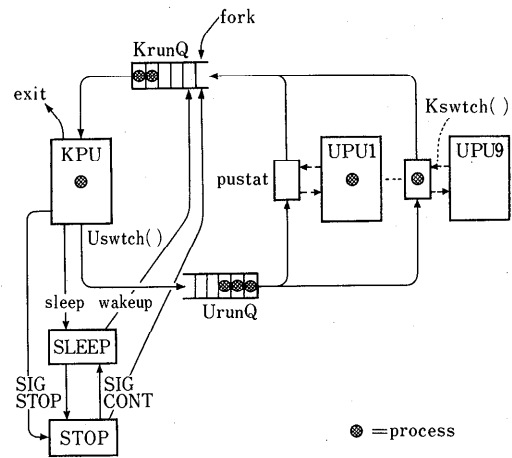


図-3

PU, 図-2 の PU 1~PU 9) ではプロセスのユーザ・コード部分が並列に実行される。

3.1 TOP-1 OS の開発

TOP-1 OS の最大の特徴はそのプロセス管理機構である。Unix プロセスは通常、ユーザ・コードを実行している状態（ユーザ・モード）とカーネル・コードを実行している状態（カーネル・モード）の二状態を行き来しながら動作している。TOP-1 OS では、プロセスのカーネル・モード部分は KPU で実行され、ユーザ・モード部分は複数ある UPU の一つで実行される。そのため、モードが切り替わるときに同時に PU も切り替える必要がある。このモード切替え用に、KrunQ, UrunQ という 2 つのプロセス・キューを用意してプロセス管理を行っている。これは一般の Unix の runQ (runnable process queue) を拡張したもので、KrunQ には KPU での実行を待っているプロセスが、UrunQ には UPU での実行を待っているプロセスが入られる。図-3 に、これらのキューの扱いも含めた TOP-1 OS におけるプロセスの状態遷移の様子を示す。PU 間のプロセスの受渡しは、図中の pustat (PU status) 構造体とメッセージ割込みを用いて行われる。

マスタ・スレーブ型の並列化では、大部分のカーネル・データに関しては排他制御は不要であるが、UPU でのプロセス実行時に使用される一部のデータに関しては排他制御などの対応が必要になる。そのため、TOP-1 OS ではメモリ管理機構、クロック処理、シグナル処理などの部分にも

変更が行われている。クロック処理の部分にはさらに、TOP-1 キャッシュの統計情報収集のための機能が追加されている。

3.2 並列プログラミングのサポート

TOP-1 OS における最もプリミティブな並列実行単位はプロセスである。1つのアプリケーションが多数のプロセスを fork し、それらが UPU 上で同時実行されることで並列性が生まれる。プロセス間のデータ共有には System V 系の Unix で用意されている共用メモリ・セグメントが利用できるが、その他に新しく追加した sharedfork というシステムコールを利用することもできる。sharedfork では、プロセスのデータ空間を共有したまま fork が行われる。スタック空間は共有されず、通常の fork と同様プロセスごとに別のものが用意される。プロセス間の同期に関しては、System V 系 Unix のセマフォの他、共有メモリ上に用意した同期変数に対して Busy Wait 方式の同期をとる方法も利用できる。

これらの機能は非常に低レベルのものであり、アプリケーションを書く際に使いやすいものとはいえない。そのため TOP-1 プロジェクトでは、並列プログラミング環境の研究も同時に行われている。C 言語でのプログラミングには、前述のプリミティブを直接利用することもできるが、より軽い同期や排他制御をめざして、ユーザレベルで動作する並列 LWP (Light Weight Process) パッケージも用意されている⁶⁾。このパッケージでは、まず sharedfork を使ってメモリを共有した複数のプロセスを生成し、その上でユーザレベルのコンテキストスイッチングを行っている。並列サポート言語処理系としては、TOP-1 Common Lisp と Parallel Fortran 処理系の研究が行われている。TOP-1 Common Lisp は、Kyoto Common Lisp をベースに future などの並列プリミティブの追加と並列ガーベージコレクション等の拡張を行ったもので、実装には前述の並列 LWP パッケージが用いられている^{7,8)}。TOP-1 上の Parallel Fortran は、IBM Parallel Fortran のサブセットで、PARALLEL LOOP 文や PARALLEL CASES 文などが利用できる。実装は、プリプロセッサと並列サポートライブラリにより行われている^{9,10)}。

3.3 ハードウェアと OS の評価

TOP-1 プロジェクトの目標の1つは、並列ハードウェアおよび OS 自身の構成法の評価にある。そのため、TOP-1 OS 上には各種データのモニタ機能が用意されている。代表的なものとして、TOP-1 OS モニタ、TOP-1 トレーサ、キャッシュ情報収集ツールなどがある。

TOP-1 OS モニタは主に OS のデバッグのために作られたもので、共有メモリの内容表示などを、OS の走行中にリアルタイムで行うことができる⁵⁾。パスロックをかけてメモリ内容をフリーズしたり、メモリ内容の変更をすることも可能である。このモニタは、TOP-1 に接続された PS/55 から共有メモリにアクセスすることで実現されている。さらに、共有メモリ上に用意された PU ごとの「仮想端末エリア」をブラウズすることで、各 PU における OS コードのトレースやプロセスの実行状況の表示も行える。

実際の並列アプリケーション実行時のモニタ機能としては、TOP-1 トレーサとキャッシュ情報収集ツールが利用できる¹¹⁾。TOP-1 トレーサは、アプリケーションのメモリアクセス状況をトレースするためのツールで、i386 の Single-Step Trap 機能を使って実現されている。キャッシュ情報収集ツールは、TOP-1 キャッシュに内蔵された統計情報収集機能を利用したもので、キャッシュのヒット率やデータ共有の状況などをリアルタイムで知ることができる。TOP-1 プロジェクトでは、これらの機能を利用して、ハードウェア、OS、アプリケーションの評価が行われている^{12,13)}。

4. OSF/1 TOP-1

TOP-1 OS は、マスタ・スレーブ型の構成をとったため、カーネル処理の頻度や比率が高いプログラムにおいては、KPU での処理がボトルネックになり、システムの並列性が十分に生かされない場合があった。ユーザ・モードでの処理がほとんどの、数値計算のようなアプリケーションでは、並列部分について 8 プロセッサで 7 倍程度の性能向上が得られる¹⁰⁾が、スケジューラや I/O の比率が高い処理では、2 倍程度の性能向上で頭打ちになるものもみられた¹⁴⁾。また、構造が固定的で OS 自身の構成法の研究には使用しにくい面

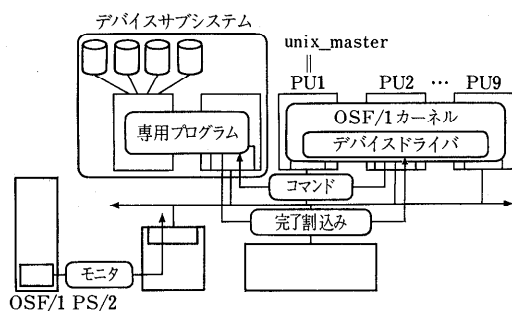


図-4

もあった。そのため、第二ステップとして、OSF/1 TOP-1の開発を行った。これは、OSF/1 1.0をTOP-1に移植したものである。OSF/1 1.0は、CMUで開発されたMach 2.5¹⁷⁾をベースにしたOSで、対称型のマルチプロセッシングをサポートしている。並列化は、カーネル内の各データ構造にロックを追加することで実現されており、ロックプリミティブとして、spin, read/write, mutexの3種類が使用されている¹⁸⁾。

TOP-1への移植にあたっては、これらの並列化用のプリミティブの実装と、機種依存のデバイス制御部などの開発を行った。図-4にOSF/1 TOP-1の構成を示す。OSF/1 TOP-1では、PU0でデバイス制御のための専用プログラムを走らせ、「プログラマブルなデバイスサブシステム」と見せせるようにしている。これにより、残りの9台のPUからは同等にデバイスにアクセスできるようになり、OS構成の自由度が大きくなっている。また、OSF/1 TOP-1上にも3.3で述べたものと同様の評価ツールが実現されている。

TOP-1プロジェクトにおいて、OSF/1 TOP-1は主に並列OS自身の構成法の研究のために用いられた。主なものとしては、デバイスサブシステムを利用したデバイスアクセス手法の比較¹⁵⁾や、ページ管理機構やキャッシュプロトコルの違いによる性能差の測定¹⁶⁾などがあげられる。前者では、デバイスアクセスの並列化により5~20%の性能向上が得られることが示されている。また後者では、マルチユーザ環境にはupdate型の、シングルユーザ環境にはinvalidate型のキャッシュ・プロトコルがより適していることなどが示されている。

5. おわりに

本稿では、並列計算機TOP-1について、OSを中心に解説した。一般にOSの研究には、「結果が見えにくい」「いじりにくい」という問題があるように思う^{*}。前者の問題を解決するには、OSカーネルだけでなくその上のサーバや言語、ミドルウェアも含めたトータルな研究を行っていく必要があるだろう。この点では、TOP-1プロジェクトは一定の成果をあげられたと考えている。後者の問題点に関しては、TOP-1プロジェクトでは根本的な解決は図られなかったように思うが、今後の方向としては、マイクロカーネルに代表される新しいOS構成法を考えていく必要があると感じている。

参考文献

- 1) 鈴木他：共有記憶型並列システムの実際，コロナ社並列処理シリーズ16 (1993)。
- 2) 清水他：高性能マルチプロセッサ・ワークステーションTOP-1，並列処理シンポジウムJSPP'89, pp.155-162 (1989)。
- 3) Oba, N. et al.: TOP-1: A Snoop-Cache-Based Multiprocessor, Proc. 9th. Annual IEEE Int. Phoenix Conf. on Computers and Communications, pp.101-108 (1990)。
- 4) Shimizu, S. et al.: Design and Evaluation of Snoop-Cache-Based Multiprocessor, TOP-1, Proc. Int. Symp. on Shared Memory Multiprocessing, pp.209-217 (1991)。
- 5) 河内谷他：TOP-1オペレーティング・システムの構造，情報処理学会研究報告90-OS-48, pp.25-34 (1990)。
- 6) 渦原，森山：マルチプロセッサシステムにおけるライトウエイトプロセス機構，日本ソフトウェア科学会第6回大会，pp.353-356 (1989)。
- 7) Tanaka, T. and Uzuhara, S.: Multiprocessor Common Lisp on TOP-1, Proc. 2nd. IEEE Symp. on Parallel and Distributed Processing, pp.617-622 (1990)。
- 8) 渦原：共有メモリ型マルチプロセッサにおける並列ガーベージコレクション，情報処理学会研究報告90-ARC-83 (SWoPP '90), pp.205-210 (1990)。
- 9) 大澤：TOP-1における流体問題解析，情報処理学会研究報告90-ARC-83 (SWoPP '90), pp.31-36 (1990)。
- 10) 大澤：TOP-1における流体問題解析，電子情報通信学会論文誌D-I, Vol. J75-D-I, No. 8, pp.757-764 (1992)。

^{*} 商用化という点で考えるとさらに、「従来OSとの互換性」や「はやらないとだめ」などの問題もある。

- 11) 大庭他: 並列処理ワークステーション TOP-1 の性能評価環境, 情報処理学会研究報告 90-ARC-83 (SWoPP '90), pp. 139-144 (1990).
- 12) Yamanouchi, N.: Performance Effects of Program Structures on a Snoop-Cached Multiprocessor System, Proc. InfoJapan '90, IPSJ, pp. 339-346 (1990).
- 13) Horiguchi, S. and Nakada, T.: Performance Evaluation of Parallel Fast Fourier Transform on a Multiprocessor Workstation, J. Parallel and Distributed Computing, Vol. 12, No. 2, pp. 158-163 (1991).
- 14) Shiratori, T. et al.: OSF/1 on the TOP-1 MP Workstation, IBM/TRL Tech. Rep. RT5026 (1992).
- 15) 河内谷他: MP UNIX におけるデバイスアクセス手法の比較, 情報処理学会研究報告 92-OS-56 (SWoPP '92), pp. 89-96 (1992).
- 16) 山崎: マルチユーザ・オペレーティングシステムの基でのスヌープキャッシュの動作効率の評価, 並列処理シンポジウム JSPP '93, pp. 371-378 (1993).
- 17) Accetta, M. et al.: Mach: A New Kernel Foundation for UNIX Development, Proc. USENIX 1986 Summer Conf., pp. 93-112 (1986).
- 18) OSF: Symmetrical Multiprocessing in the OSF/1 Operating System, OSF/1 White Paper, OSF-OSMP-WP13-1190-1 (1990).
(平成 6 年 7 月 1 日受付)



河内谷清久仁 (正会員)

1963 年生。1985 年東京大学理学部情報科学科卒業。1987 年同大学大学院理学系研究科情報科学専門課程修士課程修了。同年日本アイ・ビー・エム(株)入社。以来、同社東京基礎研究所にてオペレーティングシステムやマルチメディアシステムをはじめとするシステムソフトウェアの研究に従事。現在、同研究所副主任研究員。1994 年本会全国大会奨励賞受賞。

