

SR2001 OSの開発コンセプト

2H-1

西門 隆*、岩崎 正明*、藤田 不二男**、山本 徹***、柴宮 実**

* (株)日立製作所 システム開発研究所、**同 ソフトウェア開発本部、

***日立ソフトウェアエンジニアリング(株)

1. はじめに

大規模科学技術計算では、これまで主にベクトル型スーパーコンピュータが使われてきた。しかし、メモリを共用するベクトル型スーパーコンピュータでは処理能力の向上に限界があり、より複雑な計算には対応しきれなくなってきた。

そこで、ローカルにメモリを持つプロセッサ(ノード)を高速な通信路で結び、各ノード上のプログラムがデータ交換をしながら並列に各ローカルメモリ上のデータを処理する分散メモリ型のアプローチが取られるようになってきた。SR2001は、この分散メモリ型の科学技術計算用の並列計算機である。

本稿では、SR2001について簡単に説明すると共に、それを制御するOSの開発コンセプトと概要について述べる。

2. SR2001の概要

SR2001は図1に示す通り、RISC型プロセッサを持つノードを2次元のクロスバ方式高速通信路で繋いだ構成を持つ。システムは機能分割されており、ディスクや外部ネットワーク装置等の周辺装置を持つI/Oノードと、周辺装置を持たず主に計算を担当する計算ノードにより構成する。I/Oノードの1つ(SIOU)にはコンソールが接続され、システム全体の管理を担当する。

各次元方向のクロスバ通信路は、各ノードのプロセッサとは独立したデータ交換装置(Exchanger)で繋がれており、任意のノード間には距離に関係なく、最大でも1つのExchangerを経由するだけで高速なデータ転送が行える。ノード間通信には以下の2パターンがあり、それぞれのパターンに対して、1対1通信だけでなく、複数ノードへのブロードキャスト通信が行える。

(1) バッファ型通信

受信側が用意したリングバッファにデータを転送する通信方式

(2) 直接メモリ転送型通信

受信側の格納場所を指定してデータを転送する通信方式

さらに、ノード間的高速同期のためにバリヤ同期機能を持つ他、システムメンテナンスやシステムデバッグ^[1]用にクロスバとは別の通信路を持つ。

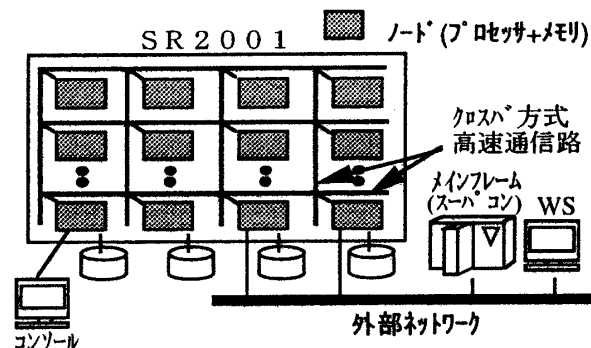


図1. SR2001のシステム構成

3. OSの開発コンセプト

科学技術計算用並列計算機では以下の2つが要求される。

(1) スケーラビリティ

並列計算機は、単に高い処理性能を実現するだけでなく、より高度な計算にも容易に対応可能とするため、要求に応じてノードや周辺装置を追加することで、プログラムを変更することなく、性能・機能が向上可能であることが必要である。すなわち、ユーザプログラム自身は、システムの構成を意識しなくても、OS側でシステムの増強を認識して、自動的に性能が向上することが必要である。

(2) マシン特性の最大限の活用

物理計算などでは、まだ解けていないグランドチャレンジ問題を解くことが並列機のメインの使用目的である。従って、(1)のOS機能によって自動的にスケラブルに性能が向上できるだけでなく、必要であれば部分的にユーザがマシンを意識してマシンに依存したチューニングを施し、そのマシンが持つ最大性能を引き出すことが可能であることが必要である。

そこで、SR2001では、これらの2点を考慮したOSの開発を目指した。

4. SR2001 OSの概要

前章で述べた開発コンセプトに基づき、SR2001では以下の特徴を持つOSを開発した。

(1) シングルUNIX+システムイメージの提供
ユーザプログラムにシステム構成を意識させずスケラブルに機能・性能の向上を可能とするため、システム全体を1つのUNIXシステムのように

Design Concept of SR2001 OS

Takashi NISHIKADO*, Masaaki IWASAKI*

Fujio FUJITA*, Toru YAMAMOTO**, Minoru SHIBAMIYA*

*Hitachi, Ltd., **Hitachi Software Engineering Co., Ltd.

+ UNIX は、X/Open Company Limitedがライセンスしている米国ならびに他の国における登録商標です。

見せる機能を提供する。この機能は、図2に示す通り、機能別のOSサーバが連携して1つのOSを構成するマイクロカーネル方式により実現する。すなわち、メモリ管理やプロセス間通信機能等の基本機能だけを持つマイクロカーネルを各ノードに登載し、各ノードの機能に応じてさらにファイルサーバやネットワークサーバを動作させる構成をとる。

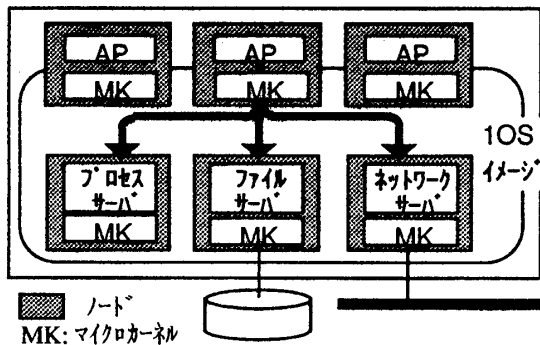


図2. SR2001 OSの構成

例えば、周辺装置を持たない計算ノード上のプログラムからのI/O要求は、マイクロカーネルを通して、対応するサーバにメッセージ通信の形で送られ処理される。各サーバへの要求の振り分けは、ノードの構成とは独立した概念の、システム内でユニークに管理されているファイル名、プロセスid等使って行うため、ノード追加に対応してOS側で振り分け先を自動的に増やすことで機能・性能を向上させることが可能である。

並列機を制御する方式としては、本方式以外にも各ノードにUNIXを登載してシステム全体をクラスタとして管理するクラスタ制御方式があるが、本方式は以下の点で並列機制御に適している。

- (a) 各ノード間でファイルシステムを互いにマウントし合う必要がなく、ノード数を増しても管理が複雑にならない。
- (b) 各ノードに外部ネットワークアドレスを割り当てる必要がなく、各サイトに割り当てられた数少ないネットワークアドレスの使用を最低限にすることが可能。
- (c) 計算ノードでのOSによるメモリ占有量が少ないため、ユーザは多くの実メモリを使用可能。

(2) 用途別によるノード割り当て機能

並列計算機上で動作するプログラムには、複数ノード使って計算をするプログラムもあれば、1つのノード上で動作してユーザとの会話処理を行うプログラムもある。後者のタイプのプログラムは、1つのノードを他と共用してTSS的に実行しても問題ないが、典型的な並列プログラムでは、各ノードのプログラムが同期して動くことを前提としているため、一部のノードを共用してノード間の処理時間がばらつく、プログラム間で待ちが発生するという問題

が生じる。そこで本システムでは、図3に示す通り、用途別にノードを分割し、その分割単位に設定されたノード共用・非共用の属性に従い、各ノードをプログラムに割り当てる機能を提供する。

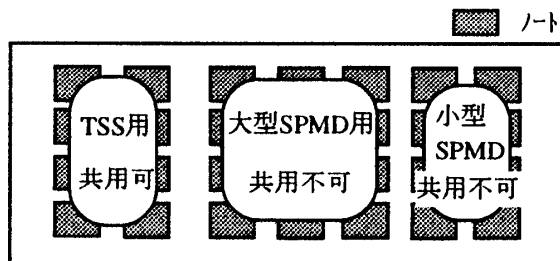


図3. 用途別ノード割り当て機能

(3) 低オーバーヘッドノード間通信機能^[2]

従来UNIXのTCP/IP等を使った通信では、OS内でのプロトコル処理のため、送受信時に数百マイクロ秒のオーバーヘッドが発生する。そこで図4に示す直接メモリ転送型通信を使用し、受信側が予め指定したユーザ空間領域に送信側が直接データを転送する機能を提供する。これにより、送受信時のOSのプロトコル処理をなくし、必要最小限の起動オーバーヘッドでノード間通信を行うことを可能としている。

また、通信状況等のモニタリング機能を提供し、マシンに依存したチューニングを可能としている。^[3]

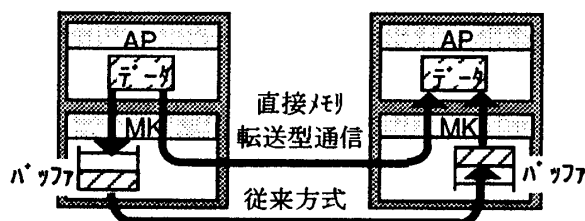


図4. 低オーバーヘッドノード間通信

5. 終わりに

科学技術計算用並列計算機SR2001用OSの開発コンセプト及び概要について述べた。本OSは、マイクロカーネル方式を用いてOS側でスケーラブル性を実現するとともに、ユーザがマシン特性を最大限活用したプログラミングをすることも可能にした。

参考文献

- [1] 山本、他：SR2001におけるカーネルデバッグ 情報処理学会第50回全国大会(1995)
- [2] 藪田、他：SR2001における高速プロセッサ間通信機能 情報処理学会第50回全国大会(1995)
- [3] 吉松、他：SR2001における並列トレーサ/モニタ機能 情報処理学会第50回全国大会(1995)