

スーパーコンピュータにおける記憶階層について

阿部 仁 和田 英夫 石井 幸一 河辺 峻
(日立製作所)

スーパーコンピュータの内部計算能力は飛躍的に向上しており、これに伴い一段と高い入出力処理能力が必要である。しかし現在もその差は大きく、この問題を技術的に解決する必要がある。

スーパーコンピュータ HITAC S-820では、主記憶と外部記憶の間に「拡張記憶」と呼ぶ記憶階層を設けている。これは主記憶と外部記憶の中間に位置付けられる、高速・大容量の内部半導体記憶である。これによって、特にベクトルジョブの実行時間の短縮、OSの多重ベクトルジョブ実行環境の改善を行っている。

本文では、S-820拡張記憶の位置付けおよび性能改善の効果について報告する。

AN EFFECT OF THE EXTENDED STORAGE FOR SUPERCOMPUTER

Hitoshi ABE, Hideo WADA, Kouichi ISHII
and Shun KAWABE

Kanagawa Works, HITACHI, Ltd.

1 Horiyamashita, Hadano, Kanagawa, 259-13, Japan

In recent years a vector performance of a supercomputer is highly increasing, and according to it more higher I/O process performance is also necessary. But now a gap between them is apparently large. Therefore some technical improvements must be done.

We designed the extended storage for the supercomputer HITAC S-820, which is a highspeed and large capacity memory of semiconductor between a main storage and auxiliary magnetic disks. By an effective use of the extended storage, an execution time of vector job is shortened, and a direct time sharing execution of vector jobs is possible.

We report an effect of the extended storage for the supercomputer HITAC S-820.

1. はじめに

スーパーコンピュータの内部計算能力は飛躍的に向上しており、これに伴い一段と高い入出力処理能力が必要である。しかしその差は大きく、この問題を技術的に解決しなければならない。

スーパーコンピュータ HITAC S-820では、主記憶と外部記憶の間に拡張記憶と呼ぶ高速・大容量半導体記憶を設け、この問題を改善しようとしており、ベクトルジョブの実行時間の短縮、OSの多重ジョブ実行環境の改善に効果を発揮している。

本文では、S-820拡張記憶の位置付けおよび性能改善の効果について報告する。

2. 問題点

(1) 一般にスーパーコンピュータは、汎用のコンピュータに比して大容量の主記憶を実装できるが、それでも大規模ジョブでは必要なデータ及び途中結果が全て主記憶に入り切らないことがあり、そのためジョブ実行中にディスクへのI/O処理が必要となる。するとスーパーコンピュータの内部計算の速度向上に対応してCPU時間が飛躍的に短縮されるにもかかわらず、総実行時間はあまり短縮されないことになり、I/Oの多いジョブについてはスーパーコンピュータ利用の効果が減少してしまう。

(2) また上記のようにスーパーコンピュータのジョブは多量の主記憶を占有するため、多重ジョブ環境において、ジョブを一旦スワップアウト/インしようとするときディスクへのI/O処理が非常に多くなり、各ジョブへの影響とシステム全体の効率低下が生じる。したがって、スーパーコンピュータの大規模ジョブのスワップアウト/インは、多くの制約があり困難である。

HITACS-810/S-820では、この点を改善すべく、主記憶と外部記憶の間に位置する高速大容量の半導体記憶（拡張記憶）を設置している。特にS-820では、後者の問題にも対応できる拡張記憶

を開発した。

3. S-820概要

HITAC S-820は、最大ピーク性能3GFLOPSを実現した、ユニプロセッサでは世界最高速のスーパーコンピュータである。S-820はパイプライン方式のプロセッサであり、ベクトルプロセッサ（VP）、スカラプロセッサ（SP）、主記憶（MS）、拡張記憶（ES）、I/Oプロセッサ（IOP）、サービスプロセッサ（SVP）よりなるシステムである。図1にS-820モデル80の構成を示す。

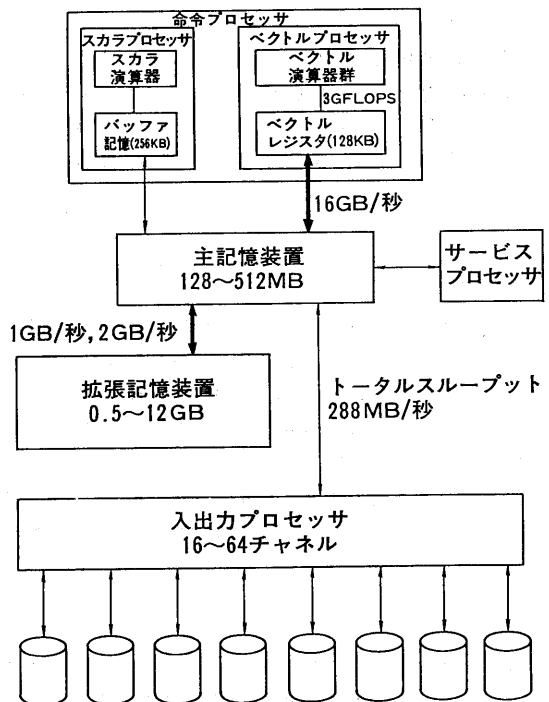


図1 S-820/80システム構成

スカラプロセッサは、通常のMシリーズの汎用命令とベクトルプロセッサ制御命令を実行できる。一方、ベクトルプロセッサは、専用の加算パイプライン、乗算パイプライン、除算パイプラインを持ち、ベクトル演算を高速にパイプライン処理できる。演算に必要なデータおよび演算結果のデー

タは超高速のベクトルレジスタに格納される。ベクトルレジスタのデータは、ロード／ストアパイプラインを通じて主記憶と転送される。(スカラプロセッサのようにキャッシュメモリは介さない。)主記憶とベクトルレジスタの転送を司るロード／ストアパイプラインは、合計して16Gバイト/秒の転送能力を持っている。

一方主記憶と入出力デバイスとのデータ転送は、I/Oプロセッサが行う。IOPは最大64チャンネルまで制御可能であり、チャンネルの下に各種のI/Oデバイスを接続する。その総転送能力は最大288Mバイト/秒あるが、チャンネルの転送能力は、各々最大3~6Mバイト/秒である。

拡張記憶(ES)は、主記憶(MS)と接続される半導体記憶装置であり、主記憶と外部記憶の間に位置する階層記憶をなしている。MSとデータ転送を行うことができ、その最大転送能力は2Gバイト/秒である。また、MSが最大容量512MB実装可能なのに対し、ESは最大容量12Gバイトまで実装可能である。このようにESは、MSに対して24倍の容量を持つことができ、チャンネルの転送能力に対し3桁近く高速である。

4. 拡張記憶のハード

半導体メモリ素子の集積度は年々進歩しており、多量の半導体記憶装置を安価に使用することが可能になってきている。S-820の拡張記憶も1MビットDRAMを使い、基板実装の技術水準を最高度に生かし、スペース効率及びコストパフォーマンスを高めている。

表1に拡張記憶の実装諸元を示す。1M DRAM 576個を、多層プリント基板の表裏両面に実装し、1枚当り64MBの記憶容量を実現している。本プリント基板を、96枚、電源装置などとともに90×86×172cmの筐体の実装する。この筐体1個で容量6Gバイト、2個で最大容量の12Gバ

トを実現している。

磁気ディスク装置に対してスペース効率も良く、コスト対性能比も良い。

ESはメカニカル部分を含まず、信頼性も高い。使用している半導体記憶素子の記憶障害に対しては、ECCコードによるメモリ障害の検出/訂正を行っており、2ビットの誤り訂正も可能である。

表1 S-820ESハード諸元

使用素子	1MビットDRAM
パッケージ	64MB/パッケージ 両面実装
転送速度	最大2Gバイト/秒
容量	最大12Gバイト
エラーチェック	2ビットエラー訂正

5. ベクトルジョブ性能の改善

前述のようにスーパーコンピュータのベクトル性能が増すにつれて、CPU時間は短縮されるが、I/O時間も短縮されないと、総実行時間は余り短縮されない。I/Oが多くベクトル化率が低いものは、スーパーコンピュータで実行する効果がないという結果になる。すなわち、ベクトルジョブのI/O時間の短縮が重要になってくる。S-820では、ディスク上ではなく、ES上にファイルを作成し、主記憶とのデータ転送を行えば、この処理時間が飛躍的に短縮できる。

S-810/S-820用OSは、ESのファイルをFORTRANの標準入出力としてサポートしている。OSは、ジョブのES使用宣言をJCL(ジョブ制御言語)から読み取り、自動的にこのファイルをES上に割り当てる。例えばJCLのファイルのDD文で、UNIT名=DASDならディスクを、UNIT名=ESならESを、そのファイルの実体として割り当てる。したがって、FORTRANソースプログラムを書き換えて、READ/WRITE文の

機器指定等を変更する必要もなく、非常に簡単にESを使用することができる。

I/Oファイルが多数あり全てES上に割当てると、ESの割当て上限値を越えてしまうような時は、そのうちI/O時間および回数が多く、ES化により最も効果が上がるファイルを選んで、ディスクからES使用に切り替える事でも良い。

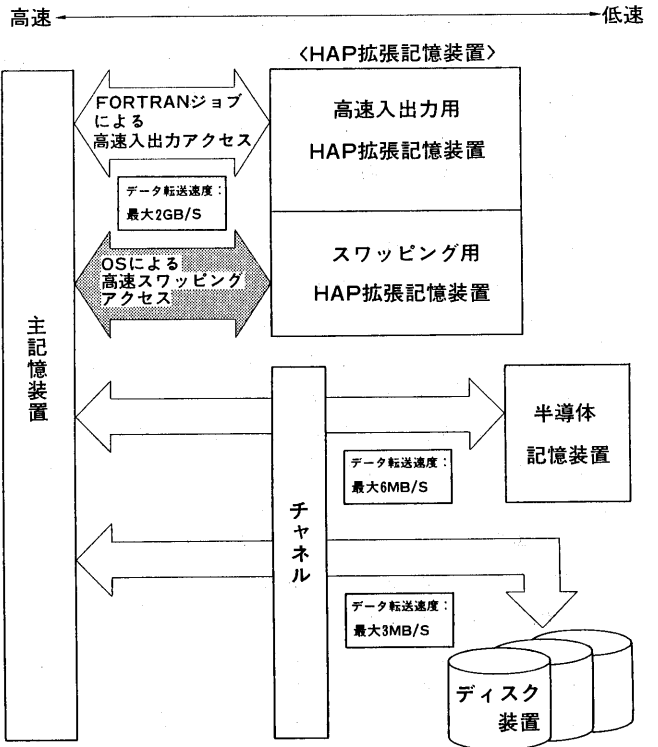
6. ベクトル多重ジョブ環境の改善

前述のように多量のMS領域を使用しているベクトルジョブは、実行途中でディスク装置にスワップアウト/インすることはI/Oネックを招き、性能上好ましくない。特に、複数のベクトルジョブの同時実行環境、例えばベクトルジョブのTSS下での直接実行などはその応答速度に問題が生じる。そこで、S-820では、ESをスワッピングファイルとしても用いることとした。これによりMS中の領域をESに高速にスワップアウト/インするため、マルチジョブ実行時のジョブ多重度の制約を軽減でき、ベクトルジョブのTSS直接実行環境の提供が可能となる。

このためにESをベクトルジョブのユーザファイル用とスワッピング用に分割できる。この分割は、各センタの運用状況に応じて、IPL時のシステムパラメータとして設定することが出来る。

7. 転送制御方式

データの転送方式は、2種類あり、それぞれを使いわけている。ベクトルジョブの中間ファイルとしてのMSとESとのデータ転送には、同期型転送が用いられる。一方、スワッピングのためのMSとESのデータ転送には非同期型転送が用いられる。同期型転送では、スカラプロセッサの命令制御部がこの転送を制御し、終了まで命令制御部は次の命令を実行しない。



効果 (1)ベクトルジョブスワッピングの高速化
(2)チャンネルビジー等の入出力ネックの緩和

図2 ESによる性能改善

非同期型転送では、スカラプロセッサの命令制御部はESに対して非同期転送の起動のみ行い、起動後は次の命令実行に移る。ESは命令制御部が与えた起動情報に従い、まずコマンド語(転送指令)をMSからフェッチし解釈したうえで、この指令によるデータ転送動作を開始する。即ち、I/O動作でのチャンネルの動作と同様のことを行う。S-820では複数のコマンド語があり、データ転送以外にコマンド語アドレスのジャンプ等も指定でき、コマンド語のチェーンができる。一連のコマンド語の実行が終了するとESは、スカラプロセッサに対して割り込みを起こす。これにより、非同期転送の終了と終了状態とを検知できる。

上記の形態を使いわけ理由は、ベクトルジョブ中のファイルアクセスでは一回当りのデータ転送量が必ずしも多くないので、同期型転送方式で転送終了まで待つ時間が非常に短く、非同期型転送方式により割り込み処理ルーチンを走ってからジョブに復帰するより、ジョブ性能に有利な場合が多いためである。これに対して、非同期型転送は、多量のデータ転送を伴うことが多いので、I/O処理と同様な使い方を採用している。

8. S-810拡張記憶との相違

S-820と前期機種S-810の拡張記憶の相違を表2に示す。

表2 S-810/820ESの相違点

	S-810	S-820
転送形態	同期転送	同期転送 非同期転送
転送速度 (GB/s)	1 or 0.5 max	2 or 1 max
容量	3GB max	12GB max

S-820では、ESの転送速度の向上およびジョブスワッピングのための非同期転送が追加されている。

9. 拡張記憶の効果

以下に拡張記憶の効果についての例を示す。

例1

(S-820使用、テストプログラム、シングルジョブ環境下、ESを64MB使用して測定)

単位：秒	ベクトル処理 +ディスク	ベクトル処理 +ES
CPU時間	0.82	0.08
I/O時間	139.07	-
総実行時間	139.89	0.08

このプログラムは、数値をES (or ディスク)へ書いて、それを読み出す動作を繰り返すもので、合計64MBの書き込みと64MBの読み出しを行っている。ESを使用すると、CPU時間=総実行時間(経過時間)となっており、I/O時間は0になった。

例2

(S-820使用、数値計算実ジョブ、シングルジョブ環境下、ESを31MB使用して測定)

	ベクトル処理 +ディスク	ベクトル処理 +ES
CPU時間	1111秒	1095秒
I/O時間	1981秒	306秒
総実行時間	3092秒	1401秒

上記の例では、全てのI/OがES化されていないが、I/O時間は相当に短縮されている。また、拡張記憶の制御はI/O制御に比べて簡単なので、I/OをESに切り換えると、CPU時間として計測されている時間も多少減少している。

例3

(S-820、マルチジョブでのスワッピング改善例)

図3に測定結果を示す。各ジョブは、コンパイル、リンケージ、実行の3ステップからなる。ジョブ多重度を3とした時の必要な総記憶容量は、実主記憶容量の1.5倍である。

ディスク使用時は、各ジョブの実行時間は、シングルラン時の約8倍を要した。これに対し、ES使用時は3ジョブの終了時まで、シングルラン時の2.5倍で済んだ。

また、ディスク使用時のCPU使用率の低下に比べて、ES使用時にはCPUをほ

ば100%使い切っている。

(図3ではジョブのシングルラン時の実行時間を1として、相対値で実行時間を示している。)

て、現在のESの技術をさらに発展させることが考えられ、今後の検討課題である。

1.1. おわりに

スーパーコンピュータHITAC S-820の記憶階層、特に拡張記憶装置(ES)の位置付けおよびその性能に対する効果について報告した。

S-820では、主記憶と外部記憶の間に拡張記憶を設け、ベクトルジョブの実行時間の短縮、OSの多重ジョブ実行環境の改善を計っている。

スーパーコンピュータは、今後も適用範囲がますます拡大する方向にあり、さらに効率的なコンピュータシステムの開発を推進していきたいと考えている。

参考文献

- 1) 河辺、他「シングルプロセッサで最大性能2 GFLOPSのS-820」日経エレクトロニクス、1987年12月27日号、N0437, pp111-125
- 2) 大房、他、「スーパーコンピュータ HITAC S-820のオペレーティングシステム」、日立評論、1987年12月号、Vol169, pp21-26
- 3) 長瀬、他、「FORTRANベクトルコンパイラとチューニングソフトウェア」、日立評論、1987年12月号、Vol169, pp27-34

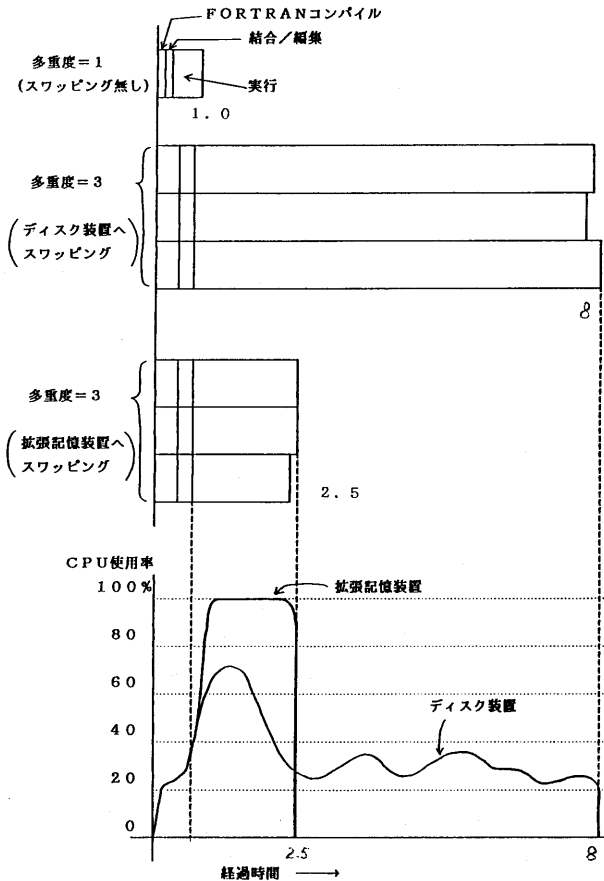


図3 ESによるスワッピング効果例

1.0. 拡張記憶装置の今後

拡張記憶はスーパーコンピュータの階層記憶の一つとして、今後とも重要な意味を持つ。その主記憶に対する大容量性とI/Oに対する高速性という特長から、I/Oとは異なる高速インタフェースを外部に対して持つことを考えることができる。これにより、大量のデータを高速に転送する。

また今後増えるマルチプロセッサ型のスーパーコンピュータの高速の共用記憶装置とし