

大容量高速記憶装置（LHS）の基本方式

西 直樹、大野直哉、妹尾義樹、中崎良成、藤原良二¹、実宝 昭²、増本健三²

日本電気株 C & C システム研究所、¹日本電気コンピュータシステム部、
²日本電気株コンピュータ技術本部

本稿では、大容量高速記憶装置 LHS (Large-capacity High-speed Storage) の記憶管理方式を報告する。LHSは階層記憶システムであり、半導体記憶装置とその他の比較的低速な記憶装置（磁気ディスク、光ディスク、磁気テープ等）から構成されている。LHSは一種の仮想記憶システムとなっており、この基本方式を「ワンレベル・ファイル・ストア」と称している。LHS実験システムをスーパーコンピュータの拡張記憶として構築し、大規模LU分解（3万2千元）プログラムを用いて評価した。全経過時間（10時間40分）に占めるLHSとのデータ転送時間比は21.7%であり、LHSが大規模数値計算環境で良好な性能を発揮することが確認されている。

The Architecture of LHS (Large-capacity High-speed Storage)

Naoki NISHI, Naoya OHNO, Yoshiki SEO, Ryousei NAKAZAKI,
Ryouji FUJIWARA¹, Akira JIPPOU², Kenzou MASUMOTO²

C&C Systems Research Laboratories, NEC Corporation

¹NEC Computer Systems Corporation

²Computer Engineering Division, NEC Corporation

4-1-1, Miyazaki, Miyamae-ku, Kawasaki, Kanagawa 213 JAPAN

This paper presents memory management techniques developed for the LHS (Large-capacity High-speed Storage). The LHS is a hierarchical file system with two kinds of file devices: a high-speed semiconductor storage and other relatively slower file devices (such as magnetic disks, optical disks and magnetic tapes). The LHS architecture is called "ONE LEVEL FILE STORE", which is a kind of virtual memory system. An experimental system was developed as an extended storage for supercomputers. An LU decomposition program was also designed for the LHS evaluation. This program solved a huge dense matrix of 32K²(8GB) with 10 hours and 40 minutes. The LHS data transfer time is 21.7% of all consuming time, which means the LHS is an effective file system for supercomputing environments.

1. はじめに

本報告は、通産省スーパーコンプロジェクトの一貫として研究開発された、大容量高速記憶装置 LHS (Large-capacity High-speed Storage) に関するものである*。LHS開発の目的は、「スーパーコンピュータに見合った大容量かつ高速な記憶装置」を研究することにあり、本目的の達成のため、LSIテクノロジや実装技術を含めたハードウェア装置開発、新しい入出力アーキテクチャを実現するための制御ソフトウェア開発、新装置／新方式の有効性実証のためのアプリケーション開発を行った。本稿では、このLHSでOS機能として重要な役割を果たしている、ワンレベル・ファイル・ストア方式に基づくファイル・システム統合アーキテクチャを中心に報告する。

2. ワンレベル・ファイル・ストアの狙い

計算機の記憶階層分類において最も一般的なのは、レジスタ、キャッシュ、メモリ、ディスク (I/O) といったハードウェアの物理的構成から見た記憶階層分類である。しかしながら、論理的な記憶階層分類は、ソフト／システム設計者の観点により大きく異なる。例えば、コンパイラ設計者はレジスタとメモリが主要な着目点であり、レジスタの最適利用のために骨身を削ることになる。リンク設計者はコンパイラとOSの間で最適かつ拡張性のある空間分割／配置に励む。

同様な意味で、OS設計者から見た記憶階層は、メモリ系とファイル系が大きな分かれ目であろう。ほとんどの計算機OSはメモリとファイルを異なるものとして利用者プログラムに提供している。また、OS内部においてもメモリ系とファイル系の独立性は高く、各々の目的／実現技術は異なったものとなる。他方、OSにおいてメモリとファイルを統合した環境として multics に端を発するワンレベル・ストアがあり、IBM社のいくつかのシステムで実際に実現されている[1, 2]。LHSのワンレベル・ファイル・ストアは、メモリとファイルは別物として捉えるが、ファイル系に関してより統合された環境をサポートすることを狙ったものである。LHSの基本構想を定めるにあたって認識した問題点／目標を以下に示す。

①入出力速度／容量の向上

I/O性能は、入出力装置自身の特性に依存す

*本研究は通商産業省工業技術院大型プロジェクトの一貫として新エネルギー・産業技術開発機構 (NEDO) から委託を受けて、実施したものである。

るところが大きい。半導体ディスク、磁気ディスク、光ディスク、磁気テープ等、転送速度と容量、記憶特性の異なる様々な記憶装置が存在する。この特性の差異を出来るだけ吸収し、より高速、かつ大容量のファイル・システムを実現することがLHSの目的である。

②ファイル・システム設計の階層化

一般的にOSが異なると、ファイルシステムも異なる。新しいOSには新しいファイルシステムを構築する必要があるが、ファイルシステムの下位層のプログラム機能は、異種OSであってもほぼ同等機能である場合が多い。

| 計算機 A | | | 計算機 B | | |
|------------|------|------|------------|------|------|
| ファイル／データ管理 | | | ファイル／データ管理 | | |
| トライバ | トライバ | トライバ | トライバ | トライバ | トライバ |
| 実装置 | 実装置 | 実装置 | 実装置 | 実装置 | 実装置 |

(a) 従来のファイル管理構築法

| 計算機 A | | | 計算機 B | | |
|------------------|------|------|------------------|------|------|
| ファイル／データ管理 | | | ファイル／データ管理 | | |
| ワンレベル・ファイル・ストア管理 | | | ワンレベル・ファイル・ストア管理 | | |
| トライバ | トライバ | トライバ | トライバ | トライバ | トライバ |
| 実装置 | 実装置 | 実装置 | 実装置 | 実装置 | 実装置 |

(b) ワンレベル・ファイル・ストア環境

図1 ファイルシステムの構築

図1にワンレベル・ファイル・ストアを用いたファイル管理法を示す。従来、ファイルシステムは実入出力装置に対するドライバの上位に構築され、基本的に2階層の構造となっている。ワンレベル・ファイル・ストアでは、これを3階層に広げ、ワンレベル・ファイル・ストア管理を中間層として設ける。ワンレベル・ファイル・ストアが実現する機能は、データ管理中の物理バッファ I/O 機能と、ファイル管理における名前と実体との対応を保持する機能である。これにより、異なるファイルシステムや装置ドライバの構築を標準化し、かつ実現を容易にすることが可能となる。大型計算機等では、ファイル／データ管理層として、通常のファイル、仮想記憶のBKST (Backing ST ore: 仮想記憶のページアウト用二次記憶装置)、データベース管理等、目的に応じた複数のファイル管理システムが同時に動作しており、これらの構築も容易になる。また、最上位のファイル／データ管理をネットワーク上で分散させることも可能である。実際、今回試作したLHSでは、ファイル管理／データ管理層が動作する計算機とワンレベル・ファイル・ストア以下のファイル記憶が

| | ディスクキャッシュ | 半導体ディスク | 拡張記憶 |
|-------------------------------|-------------------------|---|--|
| 記憶装置構成法 | | D R A M インタリープ | |
| ハード的なアクセス方法 | I / O 命令による | I / O 命令による | C P U 命令による |
| O S 内での捉え方 (ア'リケーションへの見せ方) | 基本的にはソフトから は意識する必要なし | ディスクと互換性をも たせるために、ドライ バで、シリンドやトラ ックをシミュレートしな いファイル（従来ファ イルと非互換）も実現 | 半導体ディスクと同じ 見せ方に加え、ディス クをシミュレートしな いファイル（従来ファ イルと非互換）も実現 |

表 1 半導体入出力装置の分類

動作する計算機を分離しており、両者の間は専用
ハードウェア・インターフェースで結合している。

③多種／多様な半導体記憶装置の一元化

現在、汎用大型計算機やスーパーコンピュータには3種類の半導体記憶装置が存在する（ディスク・キャッシュ、半導体ディスク、拡張記憶）。これらの装置本体は、比較的安価な DRAM をインターリープして使用して構成しているものの、ハードウェアの細部や、ソフト上での取扱いが異なり、ユーザからみた場合には別物になっている（表 1 参照）。L H S ではこれらを一元化することを図る。

3. ワンレベル・ファイル・ストア

前記の目的を考慮し、ワンレベル・ファイル・ストアは、高速な半導体記憶装置とその他の入出力装置を階層接続した構成を基本としている。また、論理的には上位ファイル管理／データ管理とのインターフェースを以下のように定めている。

①アドレス空間

図 2 にワンレベル・ファイル・ストアのアドレス空間を示す。ワンレベル・ファイル・ストアは 2^{68} B の論理記憶空間を有し、MB 単位でアドレス可能である。アドレス空間 (48bit) は、32ビットの空間 ID (SPID: SPace IDentifier) と 16 ビットの空間内ブロックアドレス (BID: Block IDentifier) からなる。ファイル管理層からの空間取得要求に際してアドレス空間を確保し、確保した 32 ビットのアドレス空間番号 (SPID) を引き渡す。

ファイル／データ管理層は、空間へのアクセスに先立ち、空間の使用をワンレベル・ファイル・ストアに通知する。OPEN 要求された空間に対して 16 ビットの一時的な空間番号 (CSPID: Current Space IDentifier) を割り当てる。ブロックの転送、ステージング等以降のアクセスは、この CSPID とブロック内アドレス BID の 32 ビット

のアドレスにより行われる。即ちファイル／データ管理層は、同時に最高 64 K 個迄の空間に対してアクセスすることが可能である。

このアドレッシング方式は、基本的にセグメンテーションを用いたものであり、セグメントの最大サイズは 64 GB に限定される。64 GB を越える大空間取得要求に対しては、ワンレベル・ファイル・ストア側が連続した SPID を割り当てることで、大空間の使用を可能にしている。



図 2 アドレス空間

②記憶管理用コマンド

上記のアドレス空間に対して行うことが可能な操作には以下のものがある。

- ・ 空間管理コマンド
 - 空間の取得／解放
 - 空間の OPEN/CLOSE
 - ステージイン／アウト
 - (+その他制御用コマンド)
- ・ 空間アクセス・コマンド
 - 空間ページへの READ/WRITE

ここで、ワンレベル・ファイル・ストアにおける空間の考え方について説明する。ワンレベル・ファイル・ストアにおいて、空間は属性を備える。属性には以下のようなものがある。

- ・性能属性
 - －低速、中速、高速
- ・記憶属性
 - －パーマネント／テンポラリ
- ・アクセス属性
 - －シーケンシャル、ランダム
- ・信頼性属性
 - －ジャーナル、二重化等

属性の指定は空間生成時に行う。この属性の指定／変更等は、空間の OPEN 時にも行えることが望ましいが、現段階では OPEN 時の指定は許していない。これは、空間情報を保持するための制御テーブル構造が空間属性によって異なるためであるが、将来的には OPEN 時の指定／変更を可能にしたいと考えている。

4. HPP-LHS 実験システム

図3にHPP-LHSシステムの構成を示す。HPP (High-speed Parallel Processors) は富士通㈱が制作された並列スーパーコンであり、4台のローカルメモリを有するベクトルプロセッサ (2 GFLOPS×3台、4 GFLOPS×1台) と、これらを接続する共有記憶装置から構成されている。この共有記憶装置にLHSも接続される。LHSの記憶装置は半導体記憶部 (4 GB) とディスク装置群から構成されている。HPP-LHSを結合しているデータ転送バスは1.5 GB/secの転送能力を有している。また、LHSのような仮想化ファイルシステムでは、ある論理アドレスに対応するデータが高速メモリ部に存在するのか、ディスク装置部に存在するのかを極めて高速に判定する必要がある。この目的のため、ハードウェア・ア

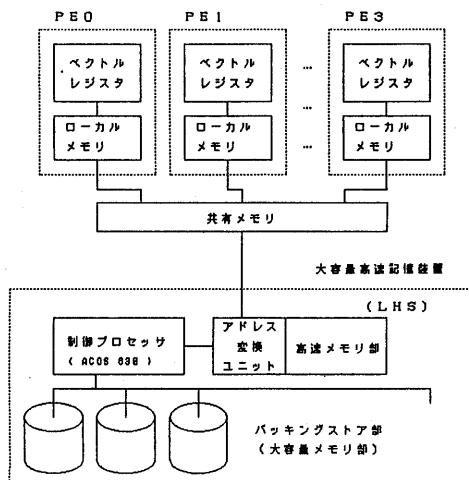


図3 HPP-LHSシステムの構成

ドレステーブルを備えて、高速メモリ部に HITするデータに対しては、制御プロセッサを介入させることなく、直接ハードウェアでデータ転送を行うことを可能としている。

アドレステーブルの論理構成を図4に示し、変換動作を説明する。HPPから転送要求とともに与えられた OPEN 中の空間アドレス (32bit) は、以下の手順で高速メモリ部に対象データが存在するかどうかがチェックされる。

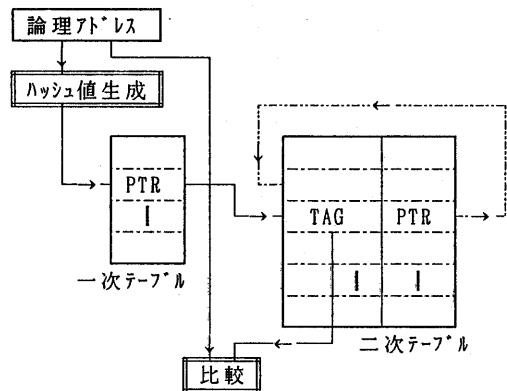


図4 アドレステーブルの構成

- ①：論理アドレスをハッシングし、一次テーブルへのポインタを求める。
 - ②：①で得られたポインタを用いて一次テーブルを参照し、二次テーブルへのポインタを求める。
 - ③：②で得られたポインタを用いて二次テーブルを参照する。二次テーブル内のエントリには、TAG フィールド (論理アドレス)、PTR フィールド (二次テーブルへのポインタ)、FLG フィールド (該当エントリの TAG/PTR の有効／無効を示す)
- の3つのフィールドが有り、これらを用いて次の動作を行う。

- ・TAG/PTR が無効であれば MISS HIT と判定
- ・TAG がチェック対象の論理アドレスと一致すれば、HIT と判定。
- ・TAG がチェック対象の論理アドレスと一致しない場合、PTR を再び二次テーブルへのポインタとして③の頭に戻る。

ここに示した変換方式は、いわゆる「逆引きテーブル」であり、二次テーブルのエントリ番号が、高速メモリ部の物理アドレスに対応する。即ち、二次テーブルのエントリ数は、LHSの高速メモ

リ部(4 GB)の全ページ数である4 Kエントリに等しい。また、本方式において、一次テーブルは高速メモリ部の特定ページに障害が発生した場合にページ単位での切り離しを可能とするために設けられている。このアドレス変換ユニットは、超高速のアドレス変換を可能するためにSi LSIに加えて3種のGaAs素子を使用し、これらを液浸冷却する方式が採用されている。

次に、HPP-LHSシステムの構成をファイルシステムの面から説明する(図5)。HPP-LHS間は、物理的に2種類のバスにより接続されている。第一のバスは空間管理系コマンドバスで、9,600 bps 無手順回線を用いている。通常、空間の作成／削除等のコマンドは、アプリケーション中の出現頻度が少なく、実験システムでは簡単な接続法を採用している。第二のバスは、空間アクセス系コマンド(READ/WRITE)の要求バスであり、転送要求は頻繁に発生、かつ処理遅延は性能低下を招くため、専用ハードウェアインターフェースを用いている。この第二のバスは1 MBブロックのデータ転送バスと共用しており、64 Bのデータバス幅を備えている。

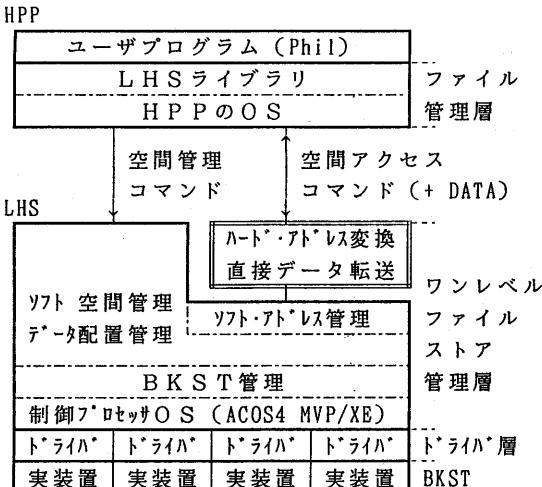


図5 LHSファイルシステム構造

ファイルシステムとして、最上位のファイル/データ管理層は、HPP上で動作する。アプリケーション中文字列で表現されるファイル名等をワンレベル・ファイルの空間番号に対応させることは、HPP側で行われる。また、本来この層で、ファイルの論理レコードと物理レコードの対応をとる必要があるが、今回の実験システムでは暫定的に論理レコードと物理レコードは等しいものとし、ユーザプログラムのデータ領域との直接デー

タ転送を行っている。基本的に今回のインプリメンテーションでは、データ管理にあたる機能は、HPP上の主力言語であるPhilに付属するI/O用ランタイム・サポート・ルーチンの形で提供されている。ファイルシステムの内、ワンレベル・ファイルを実現する層以下は、LHSの制御プロセッサ内で実現されている。

ここで、HPP上のアプリケーション・プログラムの動作が、システム全体としてどのように動作するかを説明しよう。

・空間の確保

ファイル名と記憶容量をパラメタとして空間を作成する。実験システムにおいては、空間属性はパーマネントなシーケンシャル空間として固定されている。LHS側では新規作成した空間番号を登録するとともに、この空間を管理するために必要となる最小限のページテーブル・エントリを作成する。

・空間の OPEN

ファイル名を用いて空間をOPENする。HPP上でファイル名は空間番号(SPID)に置き換えられる。LHSはOPENされたSPIDに対してCSPIDを与える。また、OPENされた空間に対応するページテーブルの一部を活性化し以降のアクセスに備える。

・空間へのアクセス

HPPから発せられたアクセス要求は、LHSのハードウェア・アドレス変換ユニットでチェックされ、HITしているならばハードウェアで直接データ転送を行う。MISS HITした場合には、制御プロセッサによりアドレス変換テーブルの再設定や、必要であればディスクとのデータ転送が行われる。新規アドレスに対するアクセスが発生した場合、この時点でページ・テーブル・エントリが割り当てられる。また、新規にデータがディスクに追い出される場合、BKSTの割り当てもこの時点で行われる。

・空間の CLOSE

HPPからCLOSE要求が発せられると、LHSは処理要求受付後、直ちにHPPに完了報告を通知する。CLOSE処理に伴って必要となる高速メモリ部からディスクへのデータ追い出しは遅延をともなって処理される。ディスクへの追い出しが完了しない間に対象空間が再度OPENされた場合、高速メモリ部に残存していたデータはそのまま再活用されることになる。

・空間の削除

HPP側ファイル管理ではファイル名が抹消される。LHS側においても、削除された空間に對して与えていた実資源（高速メモリ、BKST）を回収し、ページテーブル等の空間管理のために保持していた各種テーブルも解放する。

5. LHS制御ソフトウェアの構造

制御プログラムの構造を図5に示す。制御プログラムは大きく4つのI/O管理モジュールと、5つの内部処理モジュールから構成されている。

TIO : オペレータ端末 I/O
 FIO : 空間管理コマンド用 I/O
 XIO : 新装置管理
 DIO : 並列ディスク I/O
 CIM : 制御ソフト全体管理
 CIP : コマンド・インタプリタ
 ATM : アドレス変換管理
 SPM : 空間管理
 TIM : タイマ

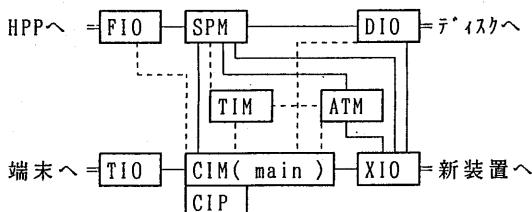


図5 LHS制御ソフトの構造

すべてのモジュールは並列動作可能であることを前提としており、合計18個のタスクで実現している。タスク間通信はメッセージ交換を基本としているが、ユーザ・データ・ブロック(1MB)等の大量データに関しては、バッファアドレスをメッセージとし、特に必要がない限りデータ自身のコピーは行わない。今回このような構造を採用したのは、LHS制御ソフトを分散環境でも効率的に動作可能としておくためである。

次に、ワンレベル・ファイル・ストアの実現を図っている空間管理部について説明する。空間管理における特徴として、次の4点があげられる。

- ・論理BKSTの提供
- ・空間管理テーブルの仮想化
- ・空間属性に応じた空間管理テーブル構造の採用
- ・動的データ配置制御

以下、この4点について説明する。

①論理BKSTの提供

LHSに使用される物理的なBKST装置は、制御ソフト内で、論理的に1MBページからなる論理BKST装置として仮想化されて取り扱われている。本構成により記憶管理の上位層は、BKSTが並列ディスク化されていることや、二重化されていること等を意識しないで構造になっている。論理BKSTと物理BKSTとのマッピングはシステム構成定義によりユーザが規定することが可能である。

②空間管理テーブルの仮想化

ワンレベル・ファイル・ストアにおけるユーザアドレス空間の最大量は 2^{68} Bである。仮に、この超巨大な記憶空間(2^{48} ページ)を保持するためには必要となる管理用テーブルを1ページ当たり4Bとするならば、管理オーバヘッドは $3 \cdot 8^{-6}$ $(2^{50}/2^{68})$ となり、極めて効率的な記憶管理が実現可能である。しかし、ここで問題となるのは記憶管理テーブルの絶対量自身である。相対的には優秀な記憶効率が得られるとしても、記憶管理テーブルに 2^{50} Bが必要であり、記憶管理テーブル自身の記憶管理が最大の問題点となる。

この問題に関し、LHSでは記憶管理テーブルを仮想化することで対応している。図6に示すように、制御プロセッサメモリ上には、高速メモリ部とBKST間のデータ転送用バッファ(1MB×4枚)に加えて、記憶管理テーブル用のキャッシュ領域(4KBページサイズ)を用意している。このキャッシュを使用するのは、OPEN中の空間のページ・テーブル、及び、BKST管理用テーブルである。これらのテーブルは、4KB単位に制御プロセッサ・メモリ上に展開され、必要に応じてディスク上のテーブルと交換される。ページ枠は、リファレンス・カウンタとページ・ロック機構を併用し、LRUアルゴリズムで管理されている。

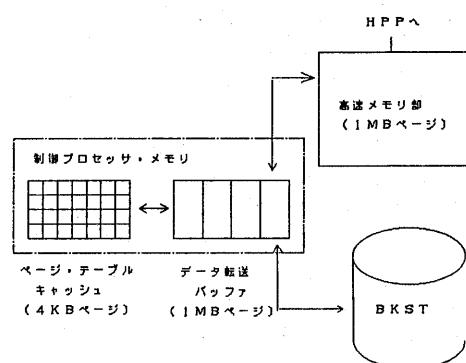


図6 ページテーブルの仮想化

③空間属性に応じた空間管理

テーブル構造の採用

前項で示したように、LHSにおいてはユーザデータ自身が仮想記憶化されるとともに、空間管理用のテーブルも仮想化されている。従って、空間管理テーブル自身のミスヒットを少なくすることも重要となる。そこで、LHSでは空間属性に応じて空間管理テーブルを異なる形式で表現している。

図7はシーケンシャル空間に対する空間管理テーブルを示している。シーケンシャル空間に対するページテーブルはBIDをオフセットとして用いる標準的なテーブル構成である。ページテーブルは4KB単位にディスクと交換されるが、シーケンシャル・アクセスされる限りページテーブル・ミスは多発しない。

他方、ランダムアクセス空間に対しては、図8に示すテーブル形式を採用している。この形式は基本的にハードウェア・アドレス変換テーブルと

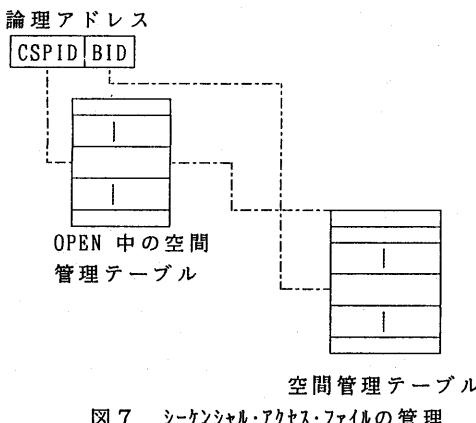


図7 シーケンシャル・アクセス・ファイルの管理

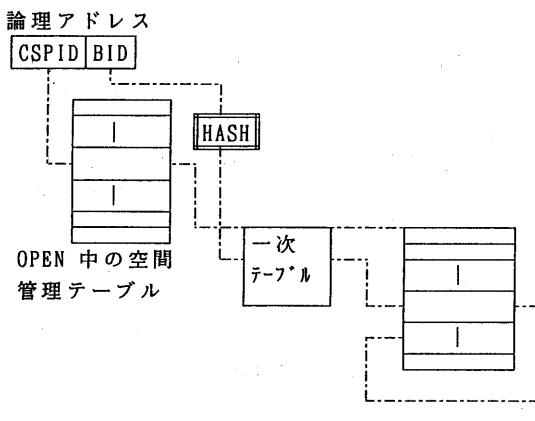


図8 ランダム・アクセス・ファイルの管理

同じ形式であり、ハッシングとポインタを用いて検索を行う。このテーブル形式はランダム空間中のすべてのアドレスが実際に使用された場合には効率が悪いが、スペース空間として使用された場合に効果を發揮する。即ち、実際に書き込みが発生している空間アドレスに対してのみ、ページテーブル・エントリが割り当てられる。

④動的データ配置制御

高速メモリ部とBKS部の最適な記憶配置を目的とし、動的データ配置制御を行うために2種類の機構の取り込みを図っている。

- ・アプリケーションからのステージイン／アウト要求に基づいた記憶配置制御

LHSに対する記憶管理用コマンドの一種として、ステージイン／アウト・コマンドが用意されている。このコマンドは、アプリケーションプログラムから記憶配置制御を行う場合に利用する。また、OSの内部からこの機能を使うと有効な場合もある。例えばJOB管理において、JOB実行開始に先だって、ユーザデータセットをプリロードして、処理高速化を図ることが可能となる。

- ・LHSに対する記憶参照履歴を用いた、自動データ配置制御

LHSは、HPPから参照されたアドレスの履歴をバッファ中に保持している。このバッファは、 FIFOとして管理されており、バッファ更新は、高速メモリ部へのミスヒット発生時、あるいは高速メモリ部へのヒットが1K回程度発生した時点を契機として行われる。また、このバッファの各エントリは、OPEN中の空間ごとにリンクされている。基本的にはこの参照履歴情報を用いて自動データ配置制御を行うことが可能である（現状では未完の機能）。一般にシーケンシャル空間に対するプリフェッチ等が考えられるが、より多くの参照パターンに対して自動配置アルゴリズムを検討して順次機能拡張をしていく必要がある。この機構の性質上、空間単位での最適化と全体としての最適化のバランスの取り方にも工夫する必要がある。

6. 性能評価

HPPとLHSを結合したシステムとして評価した結果を示す。評価システムは平成元年11月に富士通㈱沼津工場において実際に結合し、システム検査を実施後、各種評価を実施した。評価は、基本性能特性評価とアプリケーション時性能評価の2段階にわけて行っている。表2に基本性能を

示す。この結果は、HPP上の性能測定プログラムを用いて経過時間ベースで測定したものである。従って、この結果は、HPP上のI/Oサポートルーチン、LHS制御ソフト、HPPとLHS双方のハードウェア等のすべてのオーバヘッドが含まれた実効性能である。

| 性能測定ケース | 処理時間 (msec) | 処理速度 (MB/sec) |
|------------|----------------|------------------|
| 高速メモリ登録時 | 268 | 3.73 |
| リード・ヒット時 | 0.964 | 1.037 |
| ライト・ヒット時 | 0.973 | 1.028 |
| ディスク書き込み時 | 898 | 1.13 |
| ディスク読み込み時 | 856 | 1.17 |
| ディスク・スワップ時 | 1,358(2MB) | 1.47 |

表2 基本性能特性

結果を解析する。まず、ヒット時性能であるが、HPP-LHS間のデータ転送バス性能1.5GB/secに対して、実効性能は1.03GB/secとなっている。実時間にして300μsec程度の損失が発生している。これは、アドレス変換ユニット等の各種付加ハードウェア・オーバヘッドと、HPP側の性能測定プログラム(I/Oサポート・ルーチンを含む)のオーバヘッドの影響と考えられる。次に、LHS側ディスクにデータ・アクセスが発生した場合であるが、この場合の実効性能は1MB強にまで低下しており、ヒット時性能と比較すると1000倍近い性能差を生じている。ディスク自身は3MB/secチャネル4本を用い、4スピンドルに対し、16KB単位の並列アクセスを行っている。アドレス変換テーブルの再設定と、制御ソフトのオーバヘッドで400~500msecをロスしていると考えられる。今後の改良の第一課題は、ミスヒット時性能向上を図ることである。

次に、アプリケーション実行時性能として、世界最大規模の密行列に対するLU分解を実行させた場合の性能を表3に示す。このLU分解プログラムは、HPPの4台の並列プロセッサと4階層のメモリ(ローカル、共有、LHS高速部、LHSディスク)を最適に利用するために新たに開発したアルゴリズムを採用している。3万2千元の行列サイズの場合、HPP-LHSシステムでの全経過時間(10時間40分)に対し、LHSとのデータ転送所用時間は21.7%である。残りはHPPのCPU経過時間とHPP内のローカル-共有メモリ間データ転送によって消費されている。LHSとの転送時間の大半は、高速メモリ部とディスク間の転送であり、ここでもディスク性能の向上が重用な課題であることが示されている。

| | | |
|------------------|---------------------|---------------------|
| 行列サイズ* | 16,348 ² | 32,768 ² |
| 経過時間(sec) | 7,586 | 38,387 |
| HPPのローカル | 転送量(GB) | 624 |
| メモリと共有 | 時間(sec) | 3,761 |
| メモリ間転送 | 速度(MB/sec) | 166 |
| LHS | 転送量(GB) | 69 |
| HIT時 | 時間(sec) | 66 |
| の転送 | 速度(MB/sec) | 1,040 |
| LHS | 転送量(GB) | 1.9 |
| MISS HIT | 時間(sec) | 1,741 |
| 時の転送 | 速度(MB/sec) | 1.1 |
| LHS高速メモリ部登録(sec) | 274 | 1,097 |

表3 アプリケーション(LU分解)性能

7. おわりに

本稿では、大容量高速記憶装置LHSの記憶管理方式とその基本構想であるワンレベル・ファイル・ストアに関して報告した。ワンレベル・ファイル・ストアは、ファイル・システム構築の統合を狙った方式であり、半導体記憶とその他の低速記憶装置を階層接続することにより、大容量かつ実効的に高速な記憶装置の実現を目指している。HPPとLHSから構成される実験システムにおいて評価を実施した。LHSの大容量性を生かし、世界最大規模のLU分解の実行を行った。全経過時間に占めるLHSとのデータ転送時間比は21.7%であり、良好な性能が得られている。

今後の課題として、LHSの二次記憶装置のデータ転送速度を向上させることができられる。また、ワンレベル・ファイル・ストア方式として考案したすべての機構が既に実現出来たわけではない。属性空間とその特性を生かした効率的な記憶管理法の検討を進めていくことが重要である。

謝辞

本研究の実施にあたりご指導頂いた、C&Cシステム研究所 石黒所長、コンピュータ・システム研究部 小池部長に感謝します。LU分解プログラムの開発を行って頂いた日本電気技術情報システム開発部の白戸氏に感謝します。また、本プロジェクト遂行にあたり様々なご協力を頂いた関係諸機関各位に深く御礼申し上げます。

参考文献

- [1] E. I. Organick, "The Multics System: An Examination of Its Structure," Cambridge, the MIT Press, 1972.
- [2] G. G. Henry et al., "IBM System 38 Technical Developments," General System Division, IBM Corp., 1987.