

## データストリーム処理方式のデータベースコンピュータ

田中 譲

(北海道大学 工学部)

## 概 要

DBマシンの開発目標は、単にメインプロセッサの負荷の軽減を図るだけではなく、従来のデータベース処理においては基本演算とは見做されなかったジョインやディビジョンのような高度な集合演算を高速に実行するマシンの開発へと移りつつある。本稿では各々  $\log n$  個程度の比較器を持ち、パイプライン処理によりサーチソートを高速に処理するモジュールを提案し、これを用いたストリーム処理方式のデータベースコンピュータの概容を述べる。

## 1. 序 論

ファイルプロセッサを拡張し、ファイルのアクセスのみならず、データベース処理の一部を肩代りして、メインプロセッサの負荷を軽減することを目指したバックエンドプロセッサの研究は、70年代になって始められた。XDMS<sup>(1)</sup>等がこの範疇に入り、言わば第1世代のデータベースマシンと言わなければならない。その後、バックエンドプロセッサの研究は、70年代初めに新しく提案された関係データベースモデルの簡明さを取り込み、RAP<sup>(2)</sup>、CASSM<sup>(3-4)</sup>、RARE<sup>(5)</sup>等のデータベースマシンの提案を生んだ。これらのマシンは関係モデルの表形式をアーキテクチャに反映し、多数の回転メモリーデバイスと、数トラック毎に1個のプロセッサを持つ多重マイクロプロセッサの構成になっている。この型のDBマシンも多重トラックプロセッサ型と呼ぶことにする。これは言わば第2世代のDBマシンである。その後、CCDや磁気バブルデバイスの進歩により、回転メモリーとして、トラック毎にヘッドのある特殊なディスク装置のかわりに、これらの電子ディスクを用いた多重トラックプロセッサの提案がいくつかなされた<sup>(6)</sup>。これらは第2.5世代のDBマシンと言えらる。

一方、関係モデルの簡明さと普遍性は、知的データベース管理システムや知識データベースの研究に代表されるような高次の知的処理に関する基礎研究を著しく発展させた。これは、関係モデルが集合論に基づいており、従来のデータベースモデルでは基本的演算とは考えられなかった新しい型の集合演算を基本的演算として定義していることに負うところが大きい。特にジョインやディビジョンのように2つの関係に跨がる交差的演算は、データベースの処理に関する理論の発展に本質的な役割を果している。しかし、交差的演算を、ポインタメインフレームを駆使して、ソフトウェアで毎分の時間内に処理することは不可能である。

このような事情から、DBマシンの開発目的は、単にメインプロセッサの負荷を軽減することだけでなく、ノイマン型と呼ばれる現在のメインプロセッサでは毎分の時間内に処理することのできない集合演算を高速に処理するシステムの開発へと変遷しつつある。以前は、著者等は多重トラックプロセッサのエミュレータの開発<sup>(8-9)</sup>と解析的モデルの検討<sup>(10)</sup>により、この型のDBマシンは第2の開発目的に添うものはないとの結論を得た。

交差的集合演算等の集合演算の高速化には、ビットマップ、ハッシング関数、

サーチ/ソート等の技法が用いられ得る。既にこうした考えに基づく基本的な構成デバイスの提案もいくつか生まれている。この動向を、肩代りを開発目的とした第1, 第2世代のDBマシンから、集合演算の高速化を目指す第3世代のDBマシンの開発への推移の兆しと見る事ができる。前述の種々の技法の内、ビットマップ処理では、動的にビットマップを作成することによるオーバーヘッドと、交差的演算に要するビットマップが巨大になることに問題点がある。ハッシングは、交差的演算の前処理に有効であるが、交差的演算そのものを高速化するものではない。<sup>(10-11)</sup> 本稿で述べるシステムはサーチ/ソート技法を用いている。

サーチ/ソートの高速並列処理アルゴリズムは既にいくつか知られているが<sup>(12-15)</sup> これらをDBマシンの基本モジュールとして採用するには次のような問題点がある。第1に、コストと開発可能性の観点から、メモリーを除く論理回路の複雑さや、処理の対象となるデータ数に比例するようなアーキテクチャを採用すべきでない。第2にデータベースの扱う膨大な量のデータが処理部に常駐すると仮定すると膨大な数の処理デバイスを必要とするので、処理部と記憶部は別れていて、両者の間、および処理部を構成するモジュール間でデータ転送がなされると仮定すべきである。したがって、データの処理部常駐を仮定したアーキテクチャや、転送に用いられるチャンネル容量を過大に仮定したアーキテクチャは現実的でない。一方、従来より、DBマシンの隘路はデータの転送に起因することは常識として認識されている。上述の第2点に、この認識と一見矛盾している。本稿で述べるシステムは、転送時間と処理時間を完全に重畳することにより、この矛盾を解決している。すなわち、著者等は、対象データはデータのストリームとして各種モジュール間を転送されるものと考え、データストリームの転送とサーチ/ソート等のデータストリームの処理とをパイプライン的に重畳させ、2つの過程を同時に進行させることのできるアーキテクチャの開発を目指している。本稿では、著者等の開発目標となっているデータストリーム処理方式のデータベースコンピュータの基本素子であるサーチエンジンとソートエンジンの2つのモジュールのアーキテクチャを明らかにし、データストリーム型データベースコンピュータ(以下DSDBCと略す)の概容を示す。<sup>(14)</sup>

## 2. サーチエンジン (SEE: Search Engine)

### 2.1. サーチエンジンの機能

DSDBCでのキー探索処理はSEEで行う。SEEは値の順に整列したn個の被探索キーからなる順序集合中で、別に与えられるm個の任意に並んだ探索キーを順々に探索し、値の等しい被探索キーが見つかった時はその被探索キーの順番を出力し、見つからなかった時はその旨を出力するモジュールである(図2.1)。

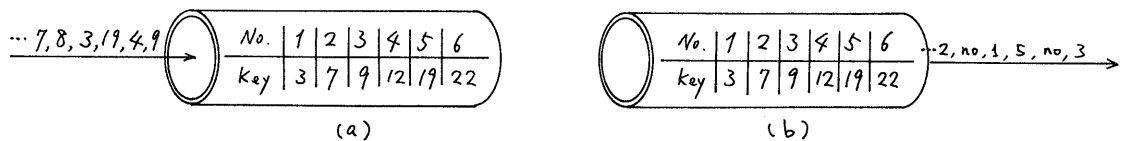


図2.1. SEEによるキー探索の一括処理。

SEEは、この処理過程における3つの過程、すなわち、探索キーの他のモジュールからの転送、探索処理、探索結果の他のモジュールへの転送をパイプライン処理により重畳させ、探索キーがデータストリームとしてSEEというパイプを通過する間に探索結果のデータストリームに変換することができる。

## 2.2. 左充填二分木とパイプライン探索

SEEの用いるデータ構造は図2.2に示すように、通常の二分木を左に詰められた形になっており、以下ではこれを左充填二分木と呼ぶ。小文字は整列した被探索キーが二分木のノードに格納される順序を示している。

SEEによるパイプライン探索の概念図を図2.2の例を用いて図2.3に示す。

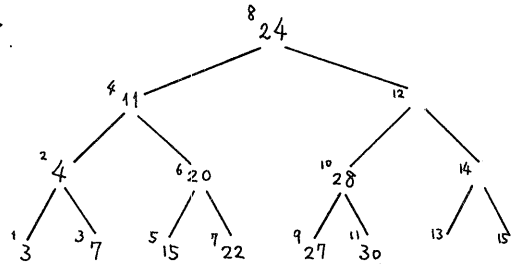


図2.2 左充填二分木の例。

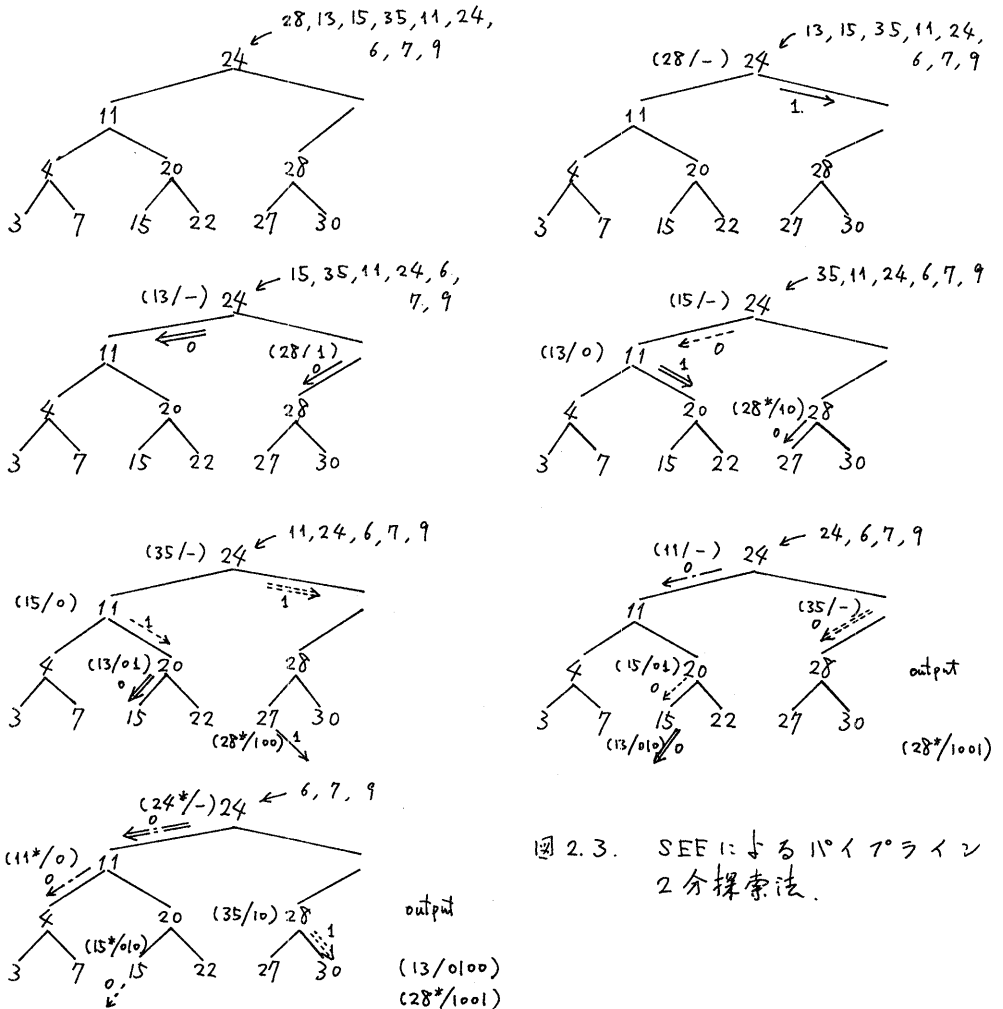
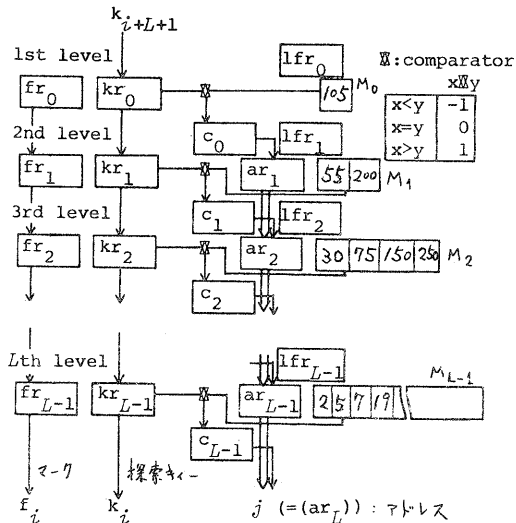


図2.3. SEEによるパイプライン二分探索法。

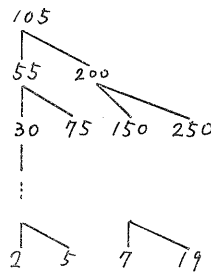
2分探索原理に従い、各探索キーはルートを起点に下方へと探索が進められるが、この過程で各レベル毎に"0"か"1"の1ビットずつの情報として付加される。探索が左下へ進むときは、それまで得られた情報の最後尾に"0"が1ビット付加され、右下へ進むときは"1"が1ビット付加される。探索が空ノードに遭遇したときは左下へ進むものとし、探索キーに等しい被探索キーを持つノードに遭遇したときは探索キーにマークを付け左下へ進む。図ではマークを"\*"で表わしている。左充填2分木において、各レベル毎のキー比較の処理は独立であり、図2.4に示すように、各レベル毎に独立のメモリーと比較器を備えたアーキテクチャにより、図2.3に示すようなバイナリ探索処理が可能である。(4) 図2.3において、 $(28*/1001)$ の1001 (=9)は探索キー"28"が10 (=9+1)番目の被探索キー $K_9$ と等しいことを表わしており、 $(13/0100)$ の13に見つからなかった探索キーに対しては、 $0100 (=4)$ は $K_{4-1} < 13 < K_4$ であることを表わす。

被探索キーに値の重複がある場合には最小番号のキーを探索する。探索法と被探索キーの充填法を少し変更することにより、最大番号のキーを探索することもできる。この2つの探索モードを用いることにより、探索キーに等しい被探索キーが何番目から何番目までの間に連続しているかを2個のSEEで同時に調べることができる。



fr<sub>l</sub>: flag register, kr<sub>l</sub>: key register,  
lfr<sub>l</sub>: loading factor register,  
ar<sub>l</sub>: address register, M<sub>l</sub>: memory

2.4 (a)



2.4 (b)

図2.4. SEEのアーキテクチャの概略

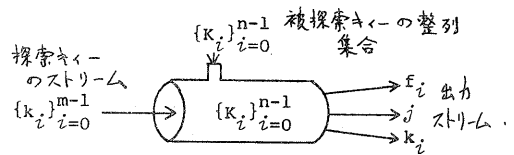


図2.5. SEEのハードウェア・シンボル

図2.4のマークというのは、図2.3中の $(x*/y)$ における"\*"をさし、アドレスというのはjをさす。被探索キーは図のように各レベルのメモリーに格納される。格納順の規則性により、格納は簡単なハードウェアで実現できる。各キーの適当な位置への格納は、キーの他デバイスからの転送に渡れることなく処理され、転送過程と整列集合から左充填2分木への変換過程は完全に重畳される。

図2.5にSEEのハードウェア記号を示しておく。

被探索キーの各  $K_i$  に対し、対応するレコード部  $R_i$  が存在し、順番い々  $R_i$  を求めることが必要の場合には、SEE の出力ストリームのアドレス部  $i$  を各々  $R_i$  に変換する変換器を SEE の後に置けばよい。これは通常の RAM に化ならず。

SEE は几个の整列した被探索キーと  $M$  個の探索キーの一括探索を  $O(m + \log n)$  の時間で処理することが出来る。試作システムのレベル数と語長は 12 レベル / 2 バイトである。SEE はチップ化に最適している。

### 3. ソートエンジン (SOE: Sort Engine)

5, 7, 10, 25, 18, 9, 20

#### 3.1. ソートエンジンの機能

DSDBC のソート処理は SOE で行う。SOE は勝手な順序で並んだキーを、キーの値の順序に並べかえる (図 3.1)。

SOE はキー集合の他デバイスからの転送入力過程と整列処理の一部を重畳させ、整列処理の残りと整列結果の他デバイスへの転送出力過程とも重畳させ、最後のキーが入力され終るやいなや、整列結果の最初のキーの出力を行うことが出来る。

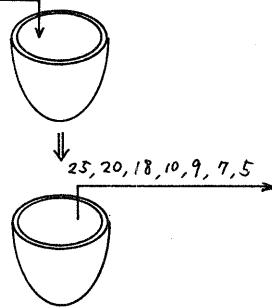


図 3.1. SOE によるキー集合の整列処理。

#### 3.2. ヒープとパイプライン・ヒープソート

SOE で用いるデータ構造は、ヒープソートで用いられるヒープである (図 3.2)。SOE で用いるアルゴリズムは基本的にはヒープソートであるが、ヒープソートそのものは並列処理やパイプライン処理に不向きである。著者等は、ヒープソートとは異なる方法で、ヒープソートと同じ効果をもたらす、しかもパイプライン処理に適したアルゴリズムを新たに開発した。パイプライン・ヒープソートは通常のヒープソートと同じようにヒープ構成とヒープ出力の 2 つのフェーズを持つ。順序を変えて、先にヒープ出力を説明する。

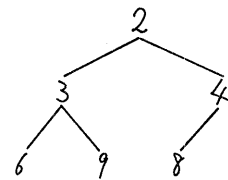


図 3.2. ヒープの例。

図 3.3 の (0)~(7) に図 3.2 に示すヒープの出力過程を示す。図 3.2 のヒープは、レベル数 3 の SOE では図 3.3 の (0) のように格納され、キーの格納されないノードには "∞" が入る。"∞" は SOE の語長で表現できる最大数を表わし、キーがこの値をとることはないとして仮定する。

ヒープ出力の動作では、奇状態と偶状態が交互に繰り返される。四中、丸で囲ったキーは、そのキーが既に親ノードへ転送されており、キー値はまだそのノードに残っているが、論理的にはそのノードがキーを持たない "空孔" 状態になっていることを示す。奇状態 (偶状態) では、ルートから数えて奇 (偶) 数番目のレベルに 1 個の空孔が存在し、偶 (奇) 数番目のレベルでは、空孔の下に位置する 1 対のキーの間の比較が行われ、小さい方のキーが上の空孔へと転送される。最下位レベルに空孔が生じると、次のステップに移る前に "∞" が入

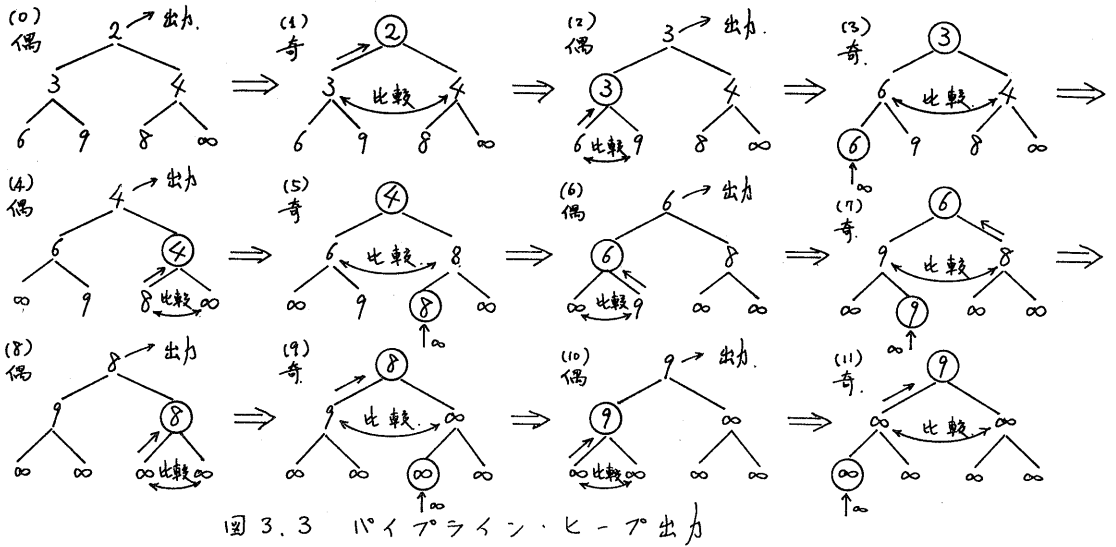


図 3.3 パイプライン・ヒープ出力

れられる。偶状態ヤルートに位置するキーが出力されるが、出力を繰り返す内にルートにいつか必ず " $\infty$ " が転送されてくる。このとき、ヒープ出力は終了する。終了時には、ヒープの各ノードがすべて " $\infty$ " となることに注意されたい。ヒープ出力に要する時間は、キー数を  $n$  として  $O(n)$  である。

ヒープ構成は、ヒープのレベル数に比例するような連想的論理回路によって処理する。図 3.4 の (a) のようなヒープにキー "5" を挿入するには、図の次経路上のデータを同時に読みだし、これに "5" を挿入し、再び次経路に同時に書き込めばよい。このように得られる木はヒープの性質を常に満足する。ヒープ構成は  $O(n)$  で実行できる。

図 3.5 に SOE のアーキテクチャの概略を示す。試作中のものは 12 レベル / 2 バイトである。SOE がチップ化に適したアーキテクチャを持っている。

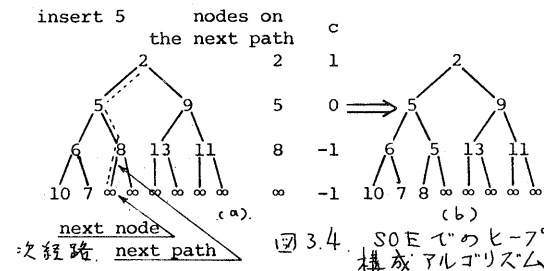
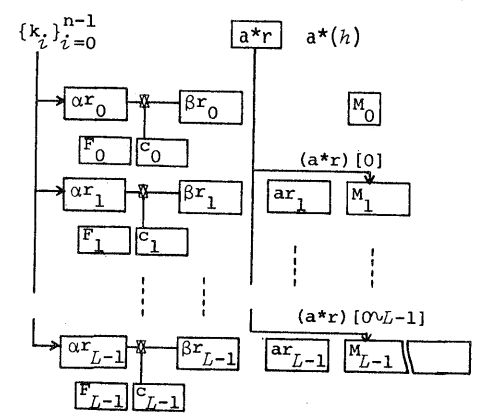


図 3.4 SOE でのヒープ構成アルゴリズム

$a^*(0)$	0	0	0	0	0	0	...	0
$a^*(1)$	0	1	0	0	0	0	...	0
$a^*(2)$	1	1	0	0	0	0	...	0
$a^*(3)$	0	0	1	0	0	0	...	0
$a^*(4)$	0	1	1	0	0	0	...	0
$a^*(5)$	1	1	0	0	0	0	...	0
$a^*(6)$	1	1	1	0	0	0	...	0
$a^*(7)$	0	0	0	1	0	0	...	0
$a^*(8)$	0	0	0	1	1	0	...	0
...								
$a^*(14)$	1	1	1	1	0	0	...	0
$a^*(15)$	0	0	0	0	1	0	...	0
$a^*(16)$	0	0	0	1	1	0	...	0
...								

3.5 (a)  
 $a^*r$  の発生するアドレスのシーケンス



3.5 (b)  
 $\alpha r_i, \beta r_i$ : registers,  $F_i$ : even/odd flag,  $M_i$ : memory,  $ar_i$ : address register,  $a^*r$ : address generator.  $(a^*r)[0:L-1]$ : the leftmost  $L$  bits of the register  $a^*r$ .

図 3.5 SOE のアーキテクチャの概略

以上の説明では、キーのみを考慮してきたが、キーとレコードの対応をキーに因りて整列しなおすためには SOE と殆ど同じアーキテクチャを持つ SOE シミュレータ (SOES) や SOE 中の対応するキーの動きをレコード部でシミュレートすればよい。これに必要な SOE から SOES へ送られる信号線は、halt 信号と、各レベルの比較結果の出力線で、(レベル数) + 1 本である。

SOE のハードウェア記号は図 3.6 に示すものを用いる。

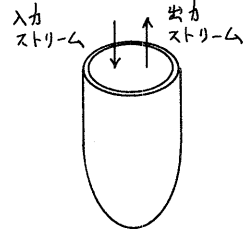


図 3.6. SOE のハードウェアシンボル

#### 4. DSDBC の概要

SOE, SEE は適当な結合により、処理可能なデータ数とデータの語長に関して機能を拡張することが可能である。<sup>(4)</sup>

DSDBC は、SOE, SEE の多数のモジュールと、インテリジェント 2 次記憶装置をコミュニケーションネットワークで結合した DB マシンである。図 4.1 に現時点での DSDBC の構想を示す。DSDBC はデータストリームを許すように拡張したデータフロー計算機の 1 種と考えることもできる。

DSDBC のデータフロー言語と、これを用いた関係演算式に対するいくつかのプログラム例<sup>(5)</sup>を参照されたい。これらの研究により、交差演算であるジョインとデイクジョンの所要時間が  $O(n)$  になることがわかってゐる。

例として、全システムの転送レートが  $5 \times 10^5$  項目/sec であるとき、 $10^6$  個の組を持つ 2 つの関係をジョインし、結果として、やはり  $10^6$  程度の組を持つ関係を生成するには、10 秒程度しか要しない。

SOE の採用により、転置ファイルは動的に高速に作成することができ、補助的なファイルは一切必要としない。

このように、DSDBC は、時間と空間の効率にあって優れた DB マシンと考えられる。

謝辞. 貴重な御助言を頂く東大宇宙航空研究所大須賀節雄助教授と電総研古川康一代に深謝します。

#### 参考文献

- (1) Canaday et al., CACM Vol. 17, pp 575-582, 1974
- (2) Ozkarahan, E.A., Proc. AFIPS, Vol. 44, pp 379-387, 1975
- (3) Copeland, G.P. et al., Proc. 1st Annual Symp. Computer Architecture, 1973.

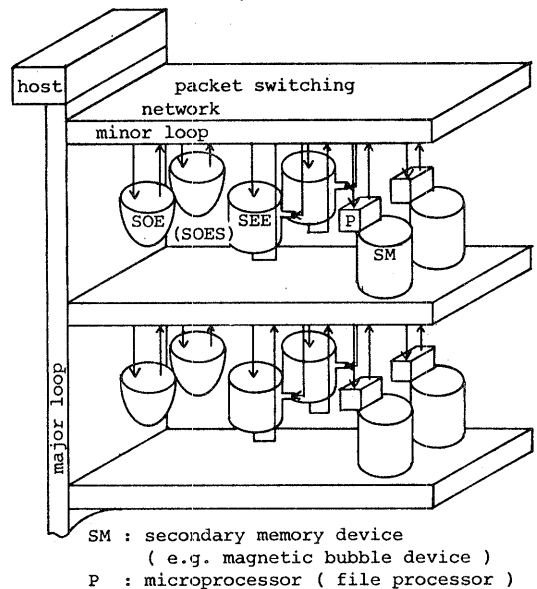


図 4.1. DSDBC のアーキテクチャの概略

- (4) Su, S.Y.W. et al., ACM TODS, Vol. 3 No. 1 pp 57-91, 1978.
- (5) Lin, C.S. et al., ACM TODS, Vol. 1, No. 1, 1976.
- (6) Schuster et al., IEEE Trans. Computers Vol. C-28, No. 6
- (7) Ozkarahan et al., Proc. 4th VLDB, pp 300-311, 1978
- (8) 田中護他, 信学研叢, EC 96-77, 1976.
- (9) 田中護他, 四学会連合学会論文集, 1977.
- (10) Bobb, E., ACM TODS, Vol. 4, No. 1, pp 1-29, 1979
- (11) McGregor, D.R. et al., Proc. 2nd VLDB, pp 103-116, 1976.
- (12) Batchier, K.E., 1968 SJCC, pp 307-314.
- (13) Nassimi, D. et al., IEEE Trans. Computers, Vol. C-28, No. 1
- (14) Tanaka, Y. et al., Proc. IFIP '80, 1980 (to appear).
- (15) 田中護他, 昭和 55 年度 情報処理学会全国大会論文集 pp. 493-494, 1980.