

## 関係データベース処理エンジンのソータの試作と評価†

岩田和秀<sup>††</sup> 神谷茂雄<sup>†††</sup> 酒井浩<sup>†††</sup>  
 柴山茂樹<sup>††††</sup> 伊藤英則<sup>††</sup> 村上国男<sup>††††</sup>

ソーティングは非数値処理分野における基本演算の1つであり、ソフトウェアの高速処理アルゴリズムの研究と同時にソータと呼ばれる専用ハードウェアによる高速処理の研究が活発に行われている。ソータの応用分野としては種々のものが考えられるが、その代表的なもの1つに関係データベース処理がある。しかし、これまでに提案されたソータは、アルゴリズムの検証や基本機能の実装が中心であり、ソータを実際に应用する場合の諸条件や処理時間の解析等が十分に検討されていない。そこで、筆者らはソータを関係データベースの処理に应用する場合に必要となる、レコード長、レコード数、キー長等のパラメータ、null値の扱い、重複レコードの検出等に柔軟な対応機能を持つソータの実現方式を考察した。次に、パイプライン化された2ウェイ・マージソート・アルゴリズムを採用して、上記の諸機能を盛り込んだセルの設計とゲートアレイによる実現を行い、これを用いてソータの試作と性能評価を行った。本ソータはセルを12個使用しており、3MByte/secのデータ転送速度に同期して4,096個の同一形式のレコードをソートする。本稿では、関係データベースの処理にソータを使用する場合にソータが具備すべき諸機能の考察、その機能を盛り込んだセルの設計例、ソート処理時間の詳細な解析および解析結果と実測値の比較評価について報告する。

### 1. ま え が き

第5世代コンピュータ・プロジェクトの前期(昭和57年度～昭和59年度)では、知識ベースマシン研究の第1歩として、Prologのファクトを格納する関係データベースマシンDeltaの開発が行われた<sup>1)</sup>。筆者らは、このプロジェクトに参加し、Deltaの特徴の1つである、関係代数演算等をパイプライン方式で実行する関係データベース処理エンジンの開発を行った<sup>2),3)</sup>。本稿では、このエンジンの基本構成要素の1つであるソータの特徴と試作結果について述べる。

ソーティングは非数値処理分野における基本演算の1つであり、古くから数多くのアルゴリズムが提案されて、目的に応じた使い分けが行われてきた<sup>4)</sup>。

1970年代に入ると、ホスト計算機のソート処理ルーチンを付加プロセッサのファームウェアで実行して、ホスト計算機の負荷を軽減する方式が提案された<sup>5)</sup>。その後、VLSI技術の進歩と関係データベースマシン研究の興隆により、ソータの提案が活発になった。

これまでに提案されたソータの主要なものには、磁気バブルメモリのループ構造を利用したソータ<sup>6)</sup>、 $n$ 個のレコードを $n$ 個のセルを用いてソートするバインニックソータ<sup>7)</sup>や並列計数ソータ<sup>8)</sup>、 $\log n$ 個のセルを用いて $n$ 個のレコードをソートするパイプライン方式のソータ等がある。これらを現状の技術を用いて試作するという観点からみると、 $\log n$ 個のセルを用いたパイプライン方式のソータが、ハードウェアが小型化できるので実用的と考えられる。

パイプライン方式のソータとしては、マージソータ<sup>9)</sup>、ヒープソータ<sup>10)</sup>、シストリックソータ<sup>11)</sup>等が発表されている。ヒープソータはレコード列の入力を終了してから結果が出力され始めるまでの出力遅れ時間がなく、かつメモリの使用効率が良いという優れた特徴を持っている。しかし、データの入出力端が同一のため、ソータの容量を単位とした連続的なパイプライン処理ができない。シストリックソータは入出力端を独立にしてヒープソータの問題点を解決し、さらに比較結果を移動歴として記憶することによりマルチウェイ・マージを可能とした大容量ソータである。しかし、マルチウェイ・マージを行う時のデータ転送制御方式等に未検討の課題が残されている。これらに対して、マージソータは若干の出力遅れ時間があることおよびメモリの使用効率が悪い欠点を持つものの、ソータの容量を単位とした連続的なパイプライン処理が可能で制御が簡単のため、汎用計算機の付加プロセッサとするのに適している。

† Implementation and Evaluation of the Sorter in a Relational Database Engine by KAZUhide IWATA (Institute for New Generation Computer Technology), SHIGEO KAMIYA, HIROSI SAKAI, SHIGEKI SHIBAYAMA (Toshiba Research and Development Center), HIDENORI ITOH (Institute for New Generation Computer Technology) and KUNIO MURAKAMI (NTT Communications and Information Processing Laboratories).

†† 新世代コンピュータ技術開発機構

††† (株)東芝総合研究所

†††† NTT 情報通信処理研究所

マージソータについては、そのアルゴリズム<sup>12)</sup>およびシミュレーション<sup>9)</sup>による動作確認等がすでに発表されている。しかし、これらの研究はアルゴリズムの検証が中心であり、ソータを実際に応用する場合の諸条件が十分に検討されていない。すなわち、ソータを関係データベースの処理に応用する場合には、レコード長、レコード数、キー長等のパラメータ、null 値の扱い、重複レコードの検出等に関する柔軟な対応が必要であるが、現段階ではこれらの実現方式を盛り込んだソータの設計は行われていない。

そこで、筆者らは関係データベースの処理にソータを使用する場合にソータで具備すべき諸機能を考察し、それらの機能を盛り込んだソータの設計を行った。次に、3 MByte/sec のデータ転送速度に同期して4,096 個の同一形式のレコードをソートできるソータ（固定長ソータ）の試作を行い、処理時間の解析結果と実測値の比較評価を行ったので、その概要を報告する。

## 2. ハードウェア化の検討

ソータを用いて関係代数演算を効率よく処理するためには、ソート処理と関係代数演算が連続して行われるよう配慮する必要がある。このため、ソート・アルゴリズムには連続的なパイプライン処理が可能な2ウェイ・マージ法を採用し、入力レコードを語（2バイト）単位で処理する方式を検討した。本章では、アルゴリズムの概要と特徴、ソータを関係代数演算の前処理に使用するために必要となる機能とその実現方式について述べる。

### 2.1 アルゴリズム

2ウェイ・マージソート・アルゴリズムは、2つのソートされたレコード列（以下、ストリングと呼ぶ）のマージ操作を、繰り返し実行してソート処理を行うものである。本ソータは、このマージ操作を、1次元に12個配置したセルにより連続的に行う。すなわち、対象となる入力レコード列（以下、ストリームと呼ぶ）を1段目のセルから入力すると、 $i$  段目 ( $i=1, 2, \dots$ ) のセルは前段セルからの  $2^{i-1}$  個のレコードからなる2つのストリングをマージして、 $2^i$  個のレコードからなるストリングを次段のセルに出力する動作を繰り返す。したがって、ストリングの長さはセルで処理されるごとに2倍され、最終段セルからは  $2^{12}$  個のレコードからなるストリングが得られる。

レコード数  $C$  が5、レコード長  $L$  が2バイトのス

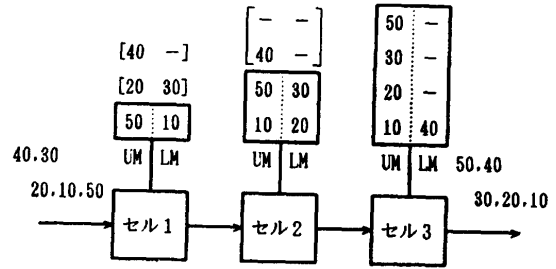


図1 2ウェイ・マージソータの原理  
Fig. 1 Principle of the 2-way merge sorter.

時刻	入力	セル 1 UM reg. LM reg.	セル 2 UM reg. LM reg.	セル 3 UM reg. LM reg.	出力	処理時間
1	50	50 —				CLT
2	10	50 (10) <sup>*</sup> —				
3	20	20 (50) —	—10 —			
4	30	— (20) <sup>*</sup> 30	5010 —			
5	40	40 (30) —	—50 (10) <sup>*</sup> —20			$(2^2-C)LT$
6		— (40)	—50 (20) <sup>*</sup> —30	—10 —		
7			4050 (30) <sup>*</sup> —	—2010 —		
8			—40 (50) —	—302010 —		$(a-L/2)2T$
9			— (40)	50302010 —		
10				—503020 (10) <sup>*</sup> —40		
11				—5030 (20) <sup>*</sup> —40	10	CLT
12				—50 (30) <sup>*</sup> —40	20	
13				—50 (40) <sup>*</sup> —	30	
14				— (50) —	40	
15					50	

図2 パイプライン処理の様子  
Fig. 2 Overview of pipeline processing.

トリーム (50, 10, 20, 30, 40 の順に入力されるとする) が入力された時の本ソータの動作原理を図1に、パイプライン処理の様子を図2に示す。ソータの各段は、2つのストリングを格納するメモリ (UM と LM) とセルで構成され、奇数番目の入力ストリングを UM に、偶数番目の入力ストリングを LM に格納する。ソータは、入力ストリームを長さ1のストリングが並んだものと見なし、これを入力順に2個ずつマージし、各段でストリング長を倍々にしていく。各段は、偶数番目のストリングの先頭データが入力されるごとに処理を開始し、図2に示すようなパイプライ

ン動作を行う。図2において、reg. は、2つのストリングのレコードの比較によって選択されたレコード(\*印)または無条件に選択されたストリングの残りのレコード(無印)を、次段へ出力するために保持するレジスタである。

図2より、ストリームが入力され始めた時点から結果が出力され終るまでのソート処理時間  $T(S)$  を求めると、次のようになる。

$$T(S) = 2CLT + (2^a - C)LT + (a - L/2)2T \quad (1)$$

ただし、 $a = \lceil \log_2 C \rceil$  で  $\lceil \cdot \rceil$  は小数点以下を切り上げた整数値、 $T$  は1バイト当りの処理時間である。上式において、第1項はストリームの入力時間と出力時間の和、第2項は  $C$  の2の累乗からのずれによる遅れ時間、第3項は入力ストリームがソータを通過するために必要な遅れ時間を示す。

本例の場合  $C=5$ ,  $L=2$ ,  $a=3$  であるから、 $T(S) = 30T$  となる。なお、本ソータの1ステップは2バイト単位の処理を行うので、 $T(S) = 30T$  をステップ数に換算すると、図2に示されているように、15ステップとなる。

次に、図2より明らかなように、各セルは偶数番目のストリングが入力され始めると、2バイト受け取ると同時に2バイト出力する動作を繰り返すので、セル  $i$  で必要になるメモリ容量は1ストリング分の  $2^{i-1} \times L$  である。しかし、メモリは2つのストリングを区別してしかもキュー機能を持つ必要があるので、本ソータでは  $2^i \times L$  の容量のメモリをセル  $i$  に実装してセルの制御回路を簡単にする方式を採用した。

## 2.2 レコード形式とデータタイプ

関係データベース処理では、レコードをキーとして扱う場合とレコードのあるフィールドをキーとして扱う場合がある。本ソータでは、セルに後述するパス・モード機能を持たせて、ストリーム内のレコードのレコード長とキー長をともに2バイト単位で最大4kバイトの長さまで指定できるようにした。

フィールドをキーとした演算を効率よくパイプライン処理するには、キーがレコードの先頭にある状態で入力される必要がある。そこで、ソータの入力側に1レコード分のメモリを用意して、レコードのフィールドを回転させ、出力側で元に戻す方式を採用した。

データベースの処理では様々なデータタイプが扱われるが、それをそのままソータで扱うことは困難である。そこで、ソータの入力側でキーを絶対値数に変換し、出力側で逆変換する方式を採用した。本ソータで

扱う入力データのタイプは、絶対値数、整数および正規化されたIBMフォーマットの浮動小数点数である。

## 2.3 Null値と重複キーの取り扱い

データベースでは、フィールドの値が不明または定義されていない時、その値はnull値と呼ばれ、特別な扱いがされる。本ソータでは、入力ストリームにnull値が含まれている場合、まず正常なキー値を持つレコードを出力し、その後null値のキーを持つレコードを入力順に出力するようにした。実現方式としては、セルの制御回路を簡単にするため、ソータの入力側で後述するnull信号を発生させる方式を考案した。

等しいキーを持つレコードを検出することは、レコードをキーとして扱う場合はunique演算等で、フィールドをキーとして扱う場合はjoin演算等で重要となる。等しいキーの検出はセルで容易にできるが、タグの操作がセルの制御回路を複雑にする。本ソータでは、最終セルの出力にチェックと呼ぶ専用回路を付加して、ソート結果のチェックと重複信号の発生を行うようにした。なお、重複したキーを持つレコードは、入力された順に出力してマルチ・キー・ソートに対応できるようにした。

## 3. システム構成の概要

### 3.1 システム構成

本ソータは関係データベース処理エンジン(以下、RDBEと略記する)の基本構成要素の1つとして開発されたので、ソータの性能評価は図3に示すRDBEのシステム構成で行った。RDBEは入力用アダプタ(HMA-IN)、INモジュール、ソータ、関係代数演算モジュール(RAPM)、出力用アダプタ(HMA-OUT)および制御用マイクロプロセッサ(ECP)で構成され、チャンネルを介して階層構造メモリ(HM)に接続されている。

RDBEはECPの制御により、HMからのデータをHMA-IN、INモジュール、ソータ、RAPM、HMA-OUTを介してパイプライン処理し、結果をHMに転送する動作を行う。

### 3.2 各モジュールの機能

- (1) HMA-INとHMA-OUT: RDBEとHM間のインタフェースの制御を行う。
- (2) INモジュール: フィールドの回転操作、データ・タイプの変換操作およびnull信号を発生す

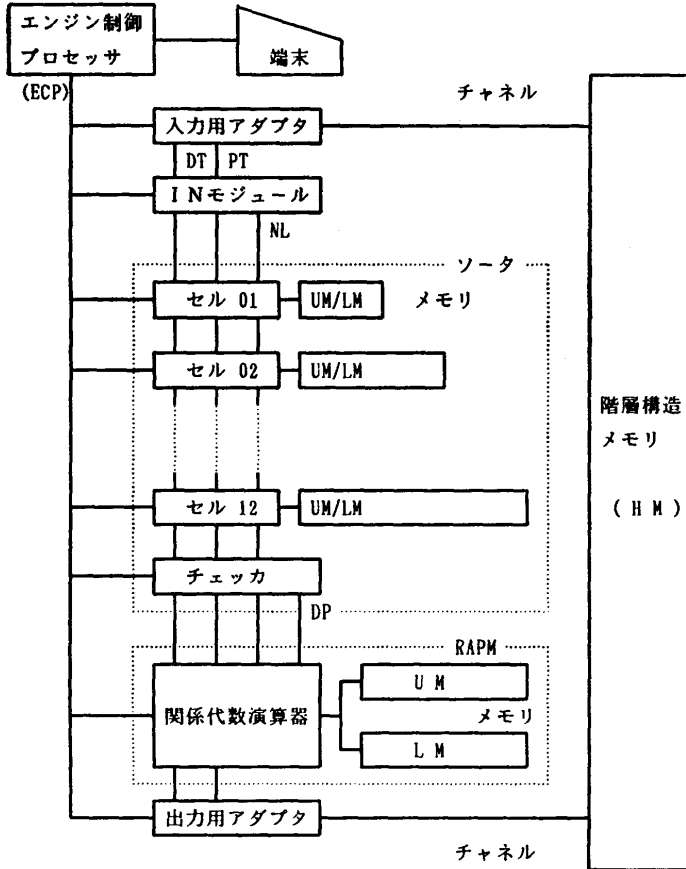


図3 RDBEのシステム構成  
Fig. 3 RDBE system configuration.

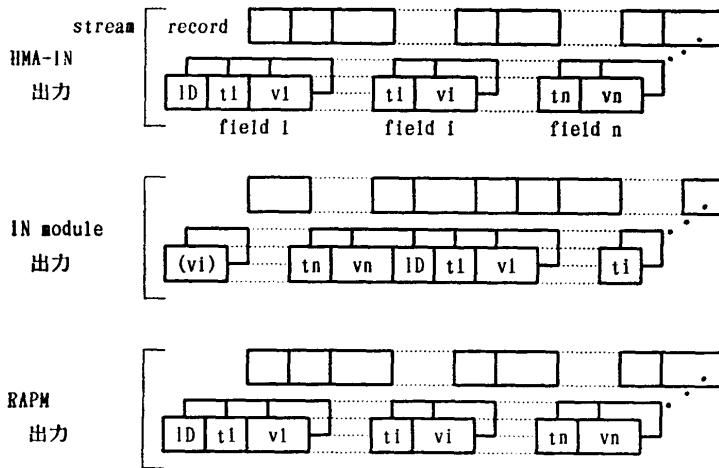


図4 レコード形式とフィールドの回転操作  
Fig. 4 Record format and field ordering.

る操作等を行う。Null 信号はタグを調べて1語ごとに生成するもので、パリティ・ビットと同様に専用線で転送される。

(3) ソータ: 12 個のセルと1個のチェッカにより構成され、入力ストリームをソートした後、その結果のチェックとキーの重複を検出して重複信号を生成する。

(4) RAPM: 2つの 64k バイトメモリ (UM, LM) と関係代数演算器で構成され、関係代数演算、マージ演算および IN モジュールで行った各種変換の逆変換等を行う。

(5) ECP: 上記のモジュールの制御、RAPM では処理できない算術演算の実行およびソータの性能測定用データの収集等を行う。

(6) HM: 汎用コンピュータで実現された大容量メモリで、RDBE とは 3 MByte/sec のデータ転送速度を持つ2つのチャンネルで結合されている。

図3において、DT, PT, NL, DP はそれぞれデータ線 (16 ビット)、パリティ・ビット線 (2ビット)、null ビット線 (1ビット)、重複ビット線 (1ビット) を示す。Null ビット線と重複ビット線は、データ線上のデータの属するレコードのキー値が、null 値であることおよび前のレコードのキー値と等しいことをそれぞれ示す。

図4に、RDBE に入力されるストリームのレコード形式と IN モジュールおよび RAPM におけるキー・フィールド  $i$  の回転操作の様子を示す。同図において、ID はレコード識別子、 $t_i$  と  $v_i$  はそれぞれ  $i$  番目のフィールドのタグおよび値、 $(v_i)$  は  $v_i$  の絶対値数表現を示す。

RDBE のソフトウェアは、ハードウェアを制御する機能とハードウェアでは実行できない算術演算等を行う機能を持つが、詳細は文献<sup>1)</sup>で述べた。

#### 4. ソータのハードウェア構成

##### 4.1 ソータの仕様

試作したソータの概略仕様は次のとおりである。

ソート・レコード数 : 4,096 個

レコード長  $L$  :  $2 \leq L \leq 4,096$  バイト  
 キー長  $K$  :  $2 \leq K \leq 4,096$  バイト  
 セル数  $N$  :  $N=12$   
 メモリ素子 :  $8 \times 8$  kビット SRAM  
 セル  $i$  のメモリ容量:  $2^{i+4}$  バイト  
 セル 12 のメモリ容量  $M$   
 :  $M=64$  kバイト  
 処理速度  $T$  :  $T=330$  nsec/バイト  
 本ソータはレコード長が 16 バイト以下の時、  
 4,096 個のレコードをソートする能力を持つ  
 が、セルのメモリ容量が固定されているので、  
 この値はレコード長により変化する。そこで、  
 セルには 2 つの入力ストリングをマージする  
 ソート・モードのほか、入力されたデータを  
 メモリを介さずにそのまま出力するパス・モ  
 ードを用意し、レコード長が大きい時はそれを格  
 納できるメモリ容量を持つセルまでパス・モ  
 ードを用いることにした。したがって、本ソータ  
 の実効ソート容量  $E$  は、次式のようになる。

$$E = L \times 2^e \quad (2)$$

ただし、 $e$  は実効動作セル数で  $e = \min([\log_2 M/L], N)$ ,  $N$  は実装セル数,  $M$  は最終段セルの実装メモリ容量,  $[\ ]$  は小数点以下を切り捨てた整数値を示す。

4.2 セルの制御方式

セルの 2 つの動作モードは、入力ストリームのレコード数  $C$  とレコード長  $L$  により、次式に従って使い分ける。

$$J = \min([\log_2 M/L], N, [\log_2 (C-1)]) \quad (3)$$

すなわち、1 から  $(N-J)$  番目のセルがパス・モードで、 $(N-J+1)$  から 12 番目のセルがソート・モードで動作するよう制御する。例えば、 $L=100$  バイトで  $C=200$  の時、 $J=8$  となるので、1 から 4 番目のセルがパス・モードで、5 から 12 番目のセルがソート・モードで動作する。なお、入力データを直接 RAPM に転送する時は、全セルをパス・モードで動作させる。

ソート・モードの動作は、1 語単位の繰り返し処理、レコード・レベルの繰り返し処理、ストリング/ストリーム・レベルの繰り返し処理を行う 3 階層構造で制御される。一方、パス・モードの動作は、1 語単位の繰り返し処理を行う階層だけで制御される<sup>3)</sup>。

4.3 セルの構造

試作したセルの構造を図 5 に示す。セルは演算部、

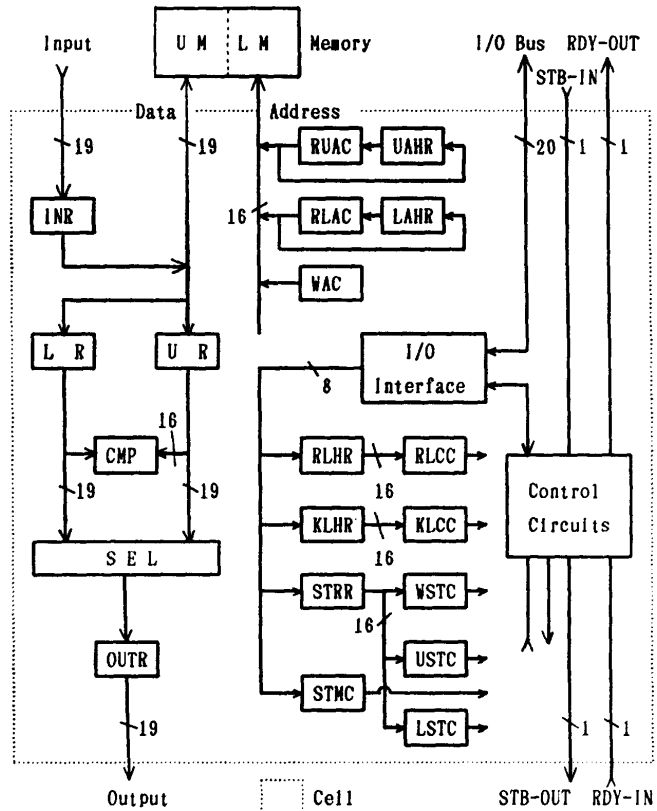


図 5 セルの構造  
 Fig. 5 Structure of the cell.

アドレス生成部、制御部およびインタフェース部より構成され、メモリに接続される。

演算部は入力データ用レジスタ (INR)、データ保持用の 2 つのレジスタ (UR, LR)、比較器 (CMP)、出力データを選択するセレクタ (SEL) および出力データ保持用レジスタ (OUTR) より構成され、2 つのデータを比較して条件に合致したデータを出力する動作を行う。

アドレス生成部は書込みアドレスカウンタ (WAC)、2 つの読み出しアドレスカウンタ (RUAC, RLAC)、2 つのアドレス保持用レジスタ (UAHR, LAHR) より構成され、入力データの書込みアドレスと比較データの読み出しアドレスを生成する。UAHR (LAHR) は比較中のレコードの先頭アドレスを保持するもので、選択されなかったレコードを次のサイクルで読み出す時に用いる。

制御部は入力ストリームのレコードを計数するカウンタ (STMC)、レコード長を保持するレジスタ (RLHR)、レコード長制御カウンタ (RLCC)、キー長保持レジスタ (KLHR)、キー長制御カウンタ (KLCC)、

ストリングのレコード数保持レジスタ (STRR), 入力ストリングのレコードを計数するカウンタ (WSTC), メモリに格納されたストリングの未処理レコードを計数するカウンタ (USTC, LSTC), 状態遷移を制御する 15 個のフリップ・フロップおよび外部インタフェースより構成され, セル全体の制御を行う。

メモリは入力ストリングを区別するため UM と LM に論理的に分割されており, 入力ストリングを UM, LM の順に交互に格納する。

#### 4.4 セルの外部インタフェース

セルは, セル間のデータ転送と ECP の入出力バスに関するインタフェース機能を持つ。

セル間のデータ転送は, 転送準備完了信号 STB-IN (OUT) と転送要求信号 RDY-IN (OUT) により前後のセルの準備が完了した時に起動され, 一方で準備が完了していなければ待ち状態に入るよう制御される。このため, ソータの入力側または出力側でデータ転送が一時的に中断されてもパイプライン動作は乱れない。

ECP からセルにセットされる制御情報には, レコード長, キー長, ストリームとストリングのレコード数, セルの動作モード指定がある。

#### 4.5 セルの処理タイミング

本セルは, サイクルタイム 220 nsec で動作し, 3 サイクルで 1 語の処理を行う。したがって, 処理速度  $T$  は 330 nsec/バイトである。ソート・モードの処理タイミング・チャートを図 6 に示す。第 1 サイクルは, セル間の制御信号 STB-IN (OUT) と RDY-IN (OUT) の 4 つが共に ON の時に開始され, 入力データの INR へのセット, UM のデータを UR に読み出す動作, および OUTR データの出力動作を行う。第 2 サイクルは LM または INR のデータを LR に読み出す動作を, 第 3 サイクルは UR と LR の比較および INR データのメモリへの書き込む動作を行う。

なお, パス・モードでは, 第 1 サイクルで入力データの INR へのセットと OUTR データの出力動作が, 第 2 サイクルで INR データの LR への転送が, 第 3 サイクルで LR データの OUTR への転送が行われる。

#### 4.6 セルのゲートアレイ化

セルの回路規模は 2 入力 NAND 回路換算で約 5,000 ゲートであり, 汎用の IC 回路で実装するとソータの 1 段分で 30 cm×30 cm の基板 1 枚を必要と

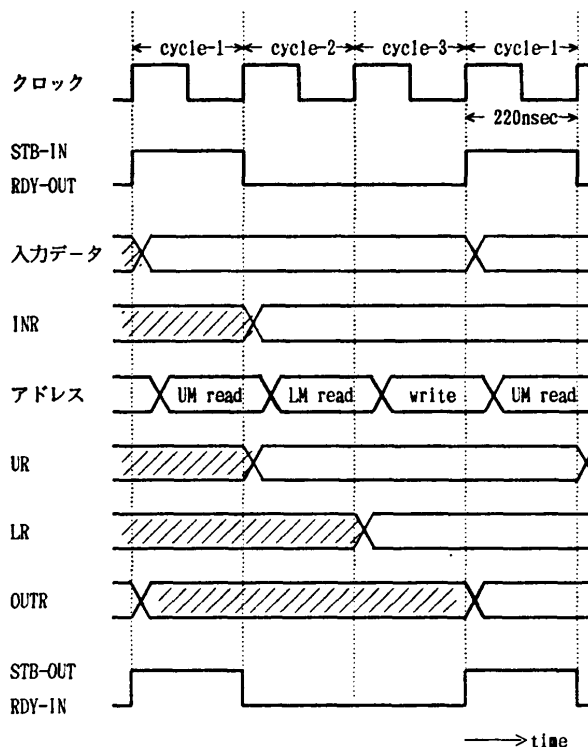


図 6 タイミング・チャート  
Fig. 6 Timing chart.

する。そこで, ソータの小型化と高信頼化を計るため, セルは CMOS ゲートアレイで実現し, 上記の基板に 4 段分を実装した<sup>1)</sup>。ゲートアレイの大きさは, セル内のレジスタやカウンタのテストを容易化する回路の追加, ゲート当りの遅延時間を 2 nsec 以内にするための配線長の短縮が行えるよう 8,000 ゲートのものを使用した。パッケージは, メモリとの信号線が 40 本, 前段セルからの入力信号線と後段セルへの出力信号線が各々 21 本, セルの制御用信号線が 20 本必要となることから, 179 ピンのものを使用した。開発期間は, 論理設計が 5 か月, 論理シミュレーションと遅延時間解析が 2 か月, 製造が 3 か月, テストが 1 か月で, 合計 11 か月である。

## 5. 性能評価

### 5.1 処理時間の解析

入力ストリームは, 実効ソート容量  $E$  以下のサブストリームに分割して処理される。この時, サブストリームの個数により処理方式が異なる。サブストリームの個数  $F$  を  $F = \langle CL/E \rangle$  と表すと,  $F=1$  の時はソータを,  $F=2$  の時はソータと RAPM を,  $F>2$  の時は  $F=2$  の処理と RAPM を用いた処理方式と

なる。

4.1 節で述べたように、本ソータはストリームの大きさによりセルの動作個数が増えるので、 $F=1$  の場合にその詳細な解析を行い、それ以外はストリームの大きさが 64k バイトの整数倍で  $E=64k$  バイトとして解析する。

(1)  $F=1$  の場合の処理時間  $T(1)$ : ストリームは HMA-IN, IN モジュール, ソータ, RAPM, HMA-OUT を介して HM に送られる。したがって、 $T(1)$  は次式で表される。

$$T(1) = t_0 + t(S) + (N-J)2T + 3LT + t_2 \quad (4)$$

ただし、 $t_0$  はコマンドの解釈とエンジンを起動するための時間、 $t(S)$  は(1)式の  $T(S)$  で  $a=J$  としたソータでの処理時間、 $(N-J)2T$  はパス・モードで動作するセルを通過する時間、 $3LT$  は IN モジュールとチェッカと RAPM における 1レコード分のバッファリング遅れ時間の和、 $t_2$  は HM への終了報告時間を示す。

(2)  $F=2$  の場合の処理時間  $T(2)$ : ストリーム  $S$  は 2つのサブストリーム ( $S1, S2$ ) に分割されてソートされ、RAPM でマージされる。この場合の各モジュールにおける処理時間を図 7 に示す。同図において、 $t_0, t_2$  および  $LT$  は  $F=1$  の場合と同様であり、 $t_1$  は  $S1$  から  $S2$  への切替えとエンジンにパラメータをセットする時間を示す。したがって、 $T(2)$  は次のようになる。

$$T(2) = t_0 + t_2 + (t_1 + 2t(s) + ET + 5LT) \quad (5)$$

(3)  $F > 2$  の場合の処理時間  $T(2^i)$ :  $F=4$  の例を検討して  $F=2^i$  の場合に拡張する。 $F=4$  の時、まずストリーム  $S$  を 2つのサブストリーム ( $S1, S2$ ) に分割し、これらを  $F=2$  の方式で処理して、ソートされたサブストリーム ( $S10, S20$ ) を求める。次に  $S10, S20$  を  $E$  の単位に分割、すなわち ( $S11, S12$ ) と ( $S21, S22$ ) にして、これら 4つのサブストリームのマージ操作を RAPM で行う。したがって、 $T(4)$  は次のようになる。

$$T(4) = t_0 + t_2 + 2A + 2(2t_1 + 3LT + 2NT + 2T + 3ET) \quad (6)$$

ただし、 $A = (t_1 + 2t(s) + ET + 5LT)$  である。

上式の右辺において、 $2A$  は  $F=2$  の処理を 2 回行う時間を、括弧内は  $S11$  と  $S21$  をマージする場合の

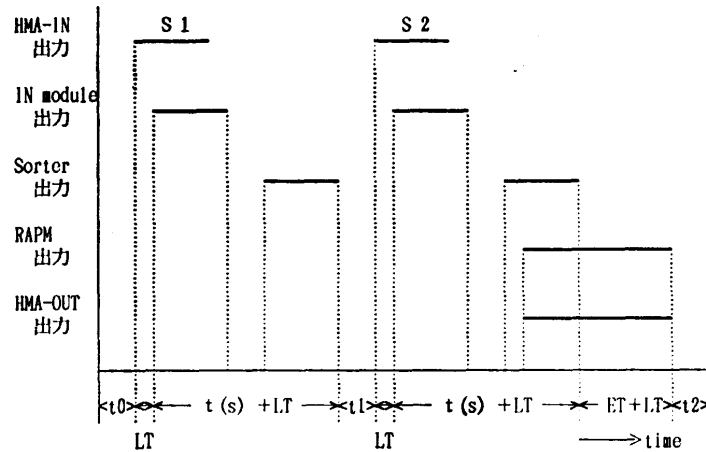


図 7 各モジュールにおける処理時間  
Fig. 7 Processing time in each module.

最悪処理時間を示す。すなわち、 $S11$  をソータをパス (遅れ時間  $(N+1)T$ ) して RAPM の UM に入力 (入力時間  $ET$ ) し、次に  $S21$  を同じくソータをパスして RAPM の LM に入力しつつ、UM と LM のレコードのマージ (処理時間  $2ET$ ) を行う時間を示す。次に、 $F=2^i$  の場合の処理は、まずストリームを  $2E$  単位でソーティングするため  $F=2$  の処理を  $2^{i-1}$  回繰り返す、その後上記のマージ操作を  $(2^i E / 2E) \log(2^i E / 2E)$  回繰り返すことにより実行する。したがって、 $T(2^i)$  は次のようになる。

$$T(2^i) = t_0 + t_2 + A2^{i-1} + B(i-1)2^{i-1} \quad (7)$$

ただし、 $B = (2t_1 + 3LT + 2NT + 2T + 3ET)$  である。

## 5.2 測定結果

測定は図 3 のハードウェア構成で、端末より RDBE コマンドを入力し、ECP 内の評価プログラムにより、コマンドの解釈からレスポンスの作成までの時間を計測した。なお、評価の対象としたストリームは、ECP で乱数を発生させ、それをすべて HM のバッファ上に格納して行った。 $F=1$  の場合の本ソータの処理時間の計算値と測定値 (■) を図 8 に、それ以外の場合を表 1 に示す。

## 5.3 考察

ソート処理時間の計算値には、 $t_0 \sim t_2$  の値として、それぞれ実測平均値 10ms, 3ms, 5ms を用いたが、測定値は計算値よりデータ量が増加するにつれて 12% から 35% ほど大きくなっている。この原因は、図 3 の評価システムが RDBE から HM にデータ要求を出して処理する方式のため、上記の  $t_0 \sim t_2$  以外に HM が RE からの内部コマンドを解釈してバッ

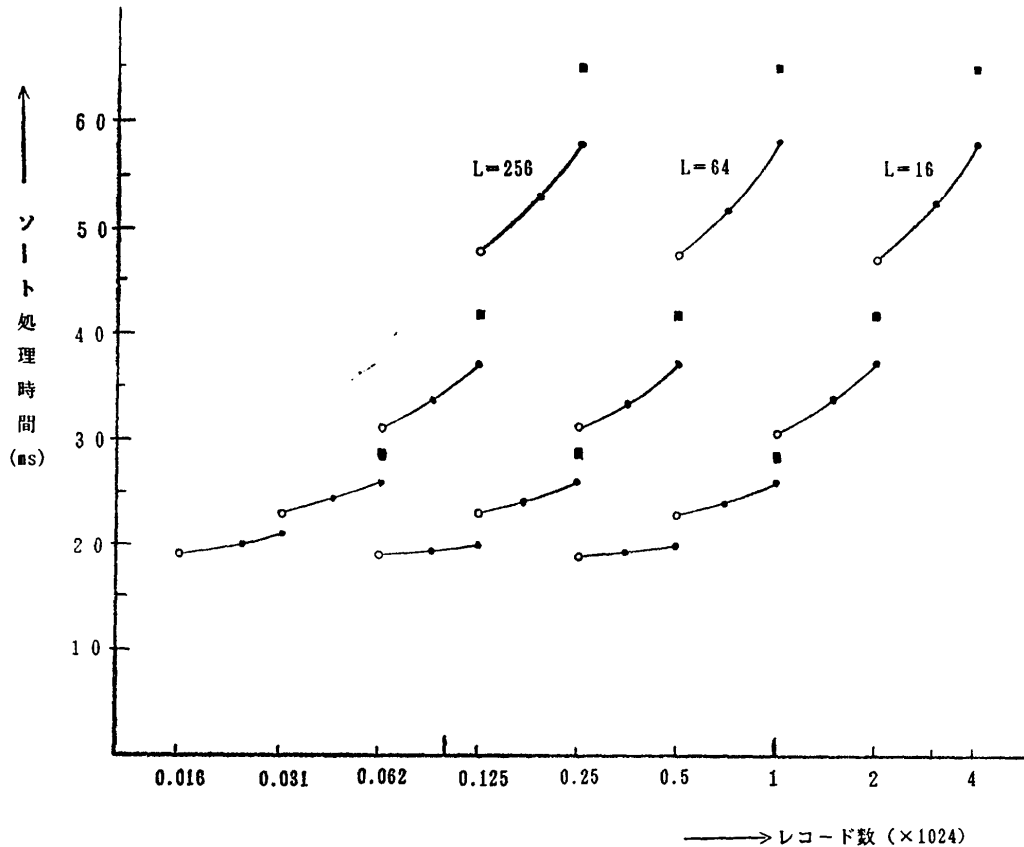


図 8 ソート処理時間 (F=1)  
Fig. 8 Sorting time (F=1).

表 1 ソート処理時間 (F≥2)  
Table 1 Sorting time (F≥2).

ストリームの レコード数	ソート処理時間 (msec)	
	計算値	実測値
1	15	16
4,096	57	65
8,192	123	141
16,384	369	460
32,768	999	1,340
65,536	2,535	3,330
132,072	6,159	8,160
262,144	14,511	18,390
524,288	33,423	45,860

(レコード長: 16 バイト)

ファ管理等を行うオーバーヘッドが存在するためである。したがって、このオーバーヘッドを減らすためには、RDBE を HM の入出力機器として HM から制御するか、または RDBE に HM の機能を含めてしまうことが有効であるので、現在後者の検討を進めている。

図 8 より、本ソータはレコード長が変化した場合、1 回にソートできるレコード数は減少するが、ソート容量 (レコード数×レコード長) からみた処理時間は一定であり、セルをパス・モードで動作させることによる遅れは無視できることがわかる。しかし、レコード数の 2 の累乗値に変曲点があることから明らかなように、レコード数が 2 の累乗値を少し越えた場合の遅れ時間が大きくなる欠点がある。しかし、平均値で見れば 10% 程度の損失であるので、アルゴリズムの改良による効果とセルの複雑度の増加とのトレード・オフを検討中である。

### 6. むすび

関係データベースの処理に必要なレコード数、レコード長、キー長等のパラメータの取り扱い、null 値の取り扱い、重複レコードの検出等の諸機能を考慮したソータの設計と試作機の評価を行った。その結果、unique 演算や join 演算等の重負荷となる演算が本ソータでの前処理と RAPM の演算により、 $O(n)$



で処理できることが確認できた。

本ソータは、単体で 64 k バイト、RAPM を連動させた時は 128 k バイトのストリームを 3 MByte/sec の処理速度でソートする。この処理速度は、現在のチャンネルの最高転送速度に相当する。したがって、本ソータは汎用計算機の付加プロセッサとして、磁気ディスクから主記憶へのデータ転送時に、ソート処理や関係代数演算を行う等の応用に適している。後者への応用については別途報告する予定である。

なお、本ソータは関係代数演算の前処理に用いることを主目的に開発したので、

- (1) 大量レコードのソートでは性能が低下する。
- (2) セルのメモリ使用効率が悪い。

等の欠点を持つ。(1)については、RAPM にマルチウェイ・マージ機能を付加する、汎用計算機の付加プロセッサとして用いる場合は汎用計算機でマルチウェイ・マージを行う等の方法が考えられる。しかし、マルチウェイ・マージを効率よく行うためには磁気ディスクを含めた検討が必要であり、今後の課題である。

(2)については、原理的にはセル  $i$  は  $2^{i-1} \times L$  のメモリ容量を持たねばよいので、ポインタを用いて削減することも可能である<sup>9)</sup>。しかし、セルの制御回路が複雑になることおよびセルのメモリ容量が2倍ずつ増えていく特殊構成のため、本ソータの規模では  $2^i \times L$  のメモリ容量を実装する方がハードウェアを小型化できる。現在、本ソータの開発経験を基に、可変長キーを処理できるソータの設計を進めている。

**謝辞** 本研究は新世代コンピュータ技術開発機構からの委託研究の一環として実施したものであり、研究に参加された同機構第3研究室、(株)日立製作所、(株)東芝の関係者に深く感謝いたします。

## 参 考 文 献

- 1) 角田, 宮崎, 柴山, 横田, 伊藤, 村上: 関係代数演算専用エンジンを備えた関係データベース・マシン Delta, 日経エレクトロニクス, No. 378, pp. 235-280 (1985).
- 2) Sakai, H., Iwata, K., Kamiya, S., Abe, M., Tanaka, A., Shibayama, S. and Murakami, K.: Design and Implementation of the Relational Database Engine, *FGCS '84*, pp. 419-426 (1984).
- 3) Kamiya, S., Iwata, K., Sakai, H., Matsuda, S., Shibayama, S. and Murakami, K.: A Hardware Pipeline Algorithm for Relational Database Operation and Its Implementation Using Dedicated Hardware, *IEEE 12th ISCA*, pp. 250-

257 (1985).

- 4) Knuth, D.E.: *The Art of Computer Programming*, Vol. 3, *Sorting and Searching*, pp. 11-388, Addison Wesley, Reading (1973).
- 5) Barsamian, H.: Firmware Sort Processor System, U.S. Patent, 3.713.107 (1973).
- 6) Chung, K.M., Luccio, F. and Wong, C.K.: On the Complexity of Sorting in Magnetic Bubble Memory Systems, *IEEE Trans. Comput.*, Vol. C-29, No. 7, pp. 553-563 (1980).
- 7) Nassini, D. and Sahni, S.: Bitonic Sort on a Mesh Connected Parallel Computer, *IEEE Trans. Comput.*, Vol. C-27, No. 1, pp. 2-7 (1979).
- 8) 安浦, 高木: 並列計数法による高速ソーティング回路, 信学論 (D), Vol. J 65-D, No. 2, pp. 179-186 (1982).
- 9) 喜連川, 伏見, 桑原, 田中, 元岡: パイプラインマージソータの構成, 信学論 (D), Vol. J 66-D, No. 3, pp. 332-339 (1983).
- 10) Tanaka, Y., Nozaka, Y. and Masuyama, A.: Pipeline Searching and Sorting Modules as Components of a Data Flow Database Computer, *IFIP 80*, pp. 427-432 (1980).
- 11) 土肥: 大容量ファイルを整理するシストリック・ソータ, 信学論 (D), Vol. J 67-D, No. 3, pp. 281-288 (1984).
- 12) Todd, S.: Algorithm and Hardware for a Merge Sort Using Multiple Processors, *IBM J. Res. Develop.*, Vol. 22, No. 5, pp. 509-517 (1978).

(昭和 61 年 9 月 19 日受付)

(昭和 62 年 4 月 15 日採録)



岩田 和秀 (正会員)

昭和 18 年生。昭和 43 年名古屋大学工学部電気工学科卒業。昭和 48 年同大学院博士課程修了。同年東京芝浦電気(株)(現(株)東芝)入社、総合研究所に勤務。電力変換装置、産業用ロボット、専用プロセッサなどの研究・開発に従事。昭和 62 年(財)新世代コンピュータ技術開発機構に出向。現在、同機構研究計画部にて海外交流を担当。電気学会会員。

**神谷 茂雄** (正会員)

昭和24年生。昭和48年早稲田大学理工学部電子通信学科卒業。昭和50年同大学院修士課程修了。同年東京芝浦電気(株)(現(株)東芝)入社。総合研究所情報システム研究所において、マルチプロセッサ、データベースマシンの研究開発を行う。現在、同社半導体技術研究所にてマイクロプロセッサの研究・開発に従事。電子情報通信学会会員。

**酒井 浩** (正会員)

昭和30年生。昭和53年東京大学工学部計数工学科卒業。昭和55年同大学院修士課程修了。同年東京芝浦電気(株)(現(株)東芝)入社。総合研究所情報システム研究所において、知識ベースマシンの研究・開発に従事。人工知能学会会員。

**柴山 茂樹** (正会員)

昭和28年生。昭和50年東京大学工学部電子工学科卒業。同年東京芝浦電気(株)(現(株)東芝)入社。総合研究所情報システム研究所において、アレイプロセッサ、CT用画像再構成プロセッサ、データベースマシン、知識ベースマシンの研究・開発に従事。この間、昭和57年～60年(財)新世代コンピュータ技術開発機構に出向。電子情報通信学会会員。

**伊藤 英則** (正会員)

昭和21年生。昭和44年3月福井大学工学部卒業。昭和49年3月名古屋大学大学院工学系研究科博士課程電気及電子専攻修了。工学博士。昭和49年4月、NTT通信研究所入社。現在、(財)新世代コンピュータ技術開発機構に勤務。研究室長。これまで、オートマトンと数理論語理論の研究と大型計算機の研究開発に従事。現在、知識情報処理システムの研究開発に従事。電子情報通信学会、人工知能学会各会員。

**村上 国男** (正会員)

昭和15年生。昭和39年茨城大学工学部電気工学科卒業。同年日本電信電話公社入社。現在、NTT情報通信処理研究所知能処理研究部総括担当主席研究員。工学博士。この間、昭和57年～60年(財)新世代コンピュータ技術開発機構第一研究室長として出向。コンパイラ、ファイル管理方式、DIPS用OS、機能分散型コンピュータ・アーキテクチャ、知識ベースマシンなどの研究・開発に従事。ACM、電子情報通信学会、日本ソフトウェア科学会、人工知能学会各会員。