

TR-203

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、神谷茂雄、酒井 浩、柴山茂樹(東芝)
村上国男(NTT)、伊藤英則(ICOT)

September, 1986

©1986, ICOT

ICOT

Mita Kokusai Bldg. 21F
4-28 Mita 1-Chome
Minato-ku Tokyo 108 Japan

(03) 456-3191~5
Telex ICOT J32964

Institute for New Generation Computer Technology

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 (東芝), 神谷 茂雄 (東芝), 酒井 浩 (東芝)
柴山 茂樹 (東芝), 伊藤 英則 (ICOT), 村上 国男 (NTT)

Implementation and Evaluation of the Sorter in Relational Database Engine

Abstract

This paper describes the design considerations and the performance evaluation of a hardware sorter. Various kinds of sorters have been proposed. However, the sorters proposed so far aimed at the verification of hardware sorting algorithms and are not designed in detail. Our sorter is designed and implemented as a component of a relational database machine Delta developed by Japan's Fifth Generation Computer Project. The sorter is used not only for sorting but also as a preprocessor for relational database operations such as join. The sorter is based on the straight two-way merge-sort algorithm and implemented using currently available technology. The sorter consists of a linear array of 12 processing cells. The sorting operation is efficiently performed by pipeline processing synchronized with the data transfer rate of 3MB/sec. We discuss the functions of the sorter required for relational database processing and evaluate the performance theoretically and experimentally.

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。筆者等は、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⁽¹⁾⁽²⁾。本報告では、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⁽³⁾。

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⁽⁴⁾。その後、VLSI技術の進歩と関係データベースマシンの研究の興隆により、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⁽⁵⁾、 n 個のレコードを n 個のセルを用いてソートするバイトニックソータ⁽⁶⁾や並列計数ソータ⁽⁷⁾、 $\log n$ 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、ハードウェアが小形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⁽⁸⁾、ヒープソータ⁽⁹⁾、シストリックソータ⁽¹⁰⁾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、更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及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なため、汎用計算機の付加プロセッサとするのに適している。

マージソータについては、そのアルゴリズム⁽¹¹⁾及びシミュレーション⁽⁸⁾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。即ち、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、3MByte/secのデータ転送速度に同期して4096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

2. 1 アルゴリズム

アルゴリズムは2ウェイ・マージソート法を語（2バイト）単位で処理する方式にパイプライン化したものである。このアルゴリズムは、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。即ち、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 個のレコードよりなるストリングを次段のセルに出力する動作を繰返す。従って、ストリングの長さはセルで処理される毎に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3、レコード長Lが6バイトの時のパイプライン処理の様子を図1に示す。この図より、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。

$$T(S) = 2CLT + (2^a - C)LT - (L - 2a)T$$

ただし、 $a = \lceil \log_2 C \rceil$ で $\lceil \rceil$ は小数点以下を切上げた整数値、Tは1バイト当りの処理時間を示す。上式において、第2項はCの2の累乗からのずれによる遅れ時間、第3項はレコードの最初の2バイトが入力された時点から処理が開始されるためのレコードの終わりからみた先行時間を示す。

次に、セルiは原理的には $(2^{i-1} + 1) \times L$ のメモリ容量を持てばよいが、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、セルに後述するバス・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。

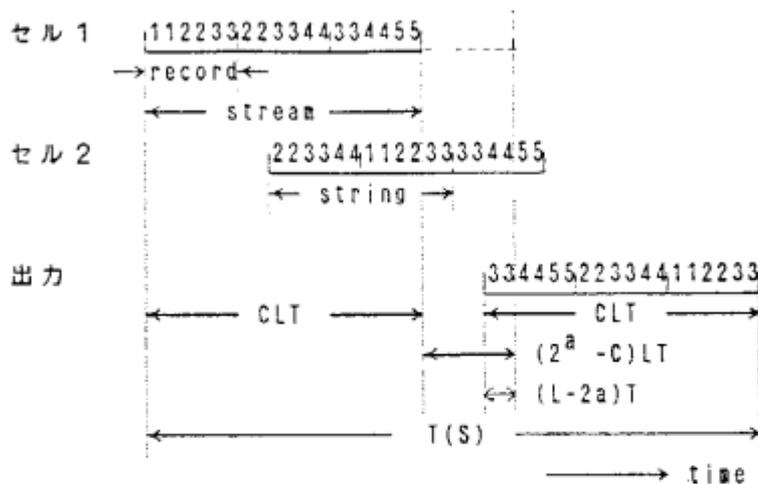


図1 パイプライン処理の様子

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、整数及び正規化された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、その後に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カ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。

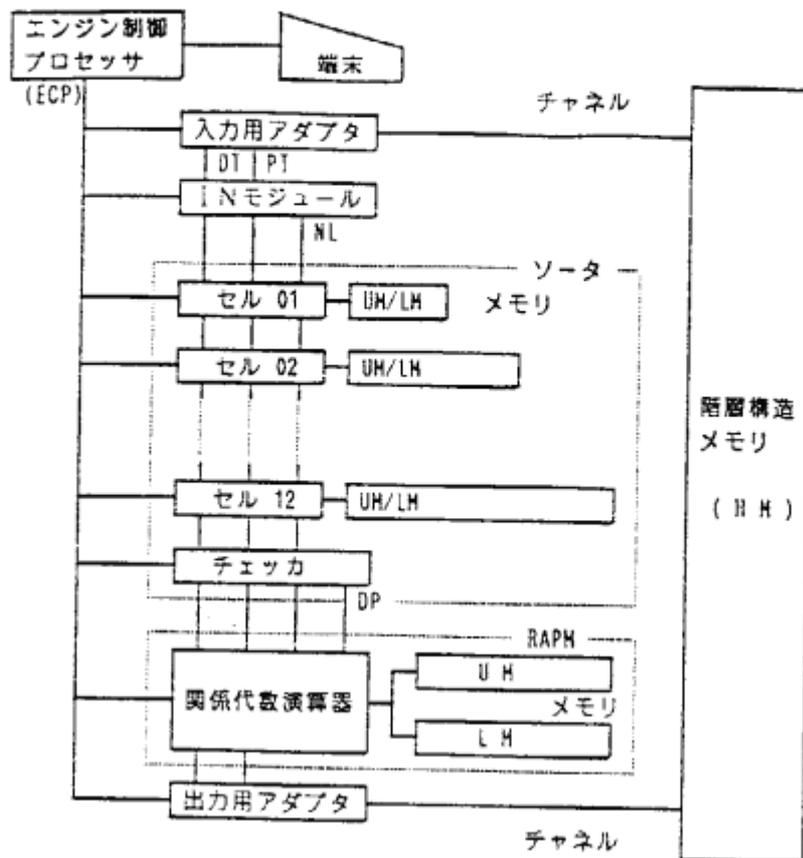


図 2 RDBEのシステム構成

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

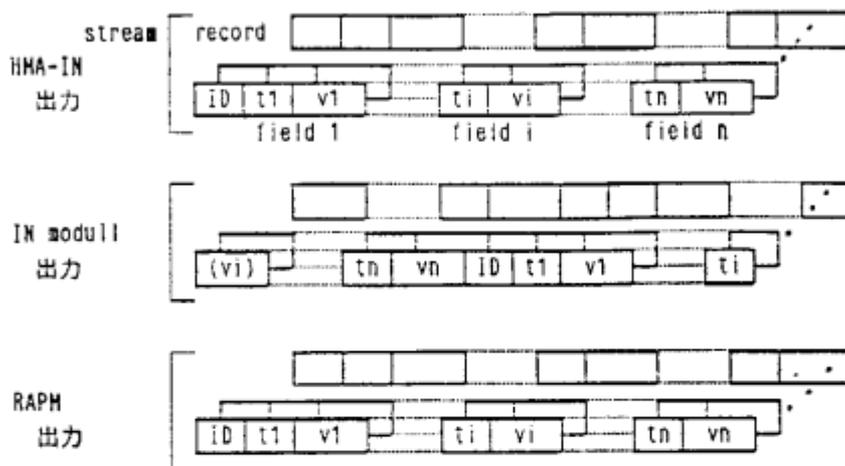


図3 レコード形式とフィールドの回転操作

生ずる操作等を行う。Null信号はタグを調べて1語毎に生成するもので、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84Kバイトメモリ(UH,LM)と関係代数演算器で構成され、関係代数演算、マージ演算及び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及びRAPMでは処理できない算術演算等を実行する。
- (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3MByte/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、RDBEに入力されるストリームのレコード形式とINモジュール及び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及び値、 (V_i) は絶対値を示す。

図2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値であること及び前のレコードのキー値と等しいことをそれぞれ示す。

3.3 ソフトウェア構成

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、本稿では前者の機能を実現する制御プログラムについて述べる。

制御プログラムの構成を図4に示す。各部の機能の概要は次の通りである。カーネル部

は割り込み処理や特権命令の処理、入出力制御部はRDBEの各モジュールの制御、関係代数演算処理部はRDBEコマンドの解釈と実行、トップ・レベルはRDBEコマンド処理のスケジューリングをそれぞれ行う。なお、テスト・評価プログラム部には、データベースのサイズの指定と作成機能及びRDBEコマンドの種類をメニューから選択できる機能を持たせ、端末よりソータの性能測定が行えるようにした。

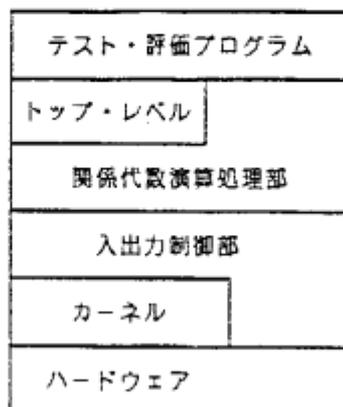


図4 RDBE制御プログラムの構成

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通りである。

- ソート・レコード数 : 4096個
- レコード長L : $2 \leq L \leq 4096$ バイト
- キー長K : $2 \leq K \leq 4096$ バイト
- セル数N : $N = 12$
- メモリ素子 : $8 \times 8k$ 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 : $M = 64k$ バイト
- 処理速度T : $T = 330ns$ / バイト

本ソータはレコード長が16バイト以下の時、4096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・モードの他に、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バス・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バス・モードを用いることにした。従って、本ソータの実効ソート容量Eは、次式のようになる。

$$E = L \times 2^e$$

ただし、 e は実効動作セル数で $e = \min([\log_2 M/L] \cdot N)$ 、 N は実装セル数、 M は最終段セルの実装メモリ容量、 $[\]$ はガウス記号を示す。

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。

$$J = \min([\log_2 M/L] \cdot N, [\log_2 2(C-1)])$$

即ち、1 から $(N-J)$ 番目のセルがバス・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番目のセルがバス・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。入力データを直接RAPHに転送する時は、全セルをバス・モードで動作させる。

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、アドレス生成部、制御部及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トア(SEL)及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,RIAC)、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・フロップ及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Uメモリ(UM)とLメモリ(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。

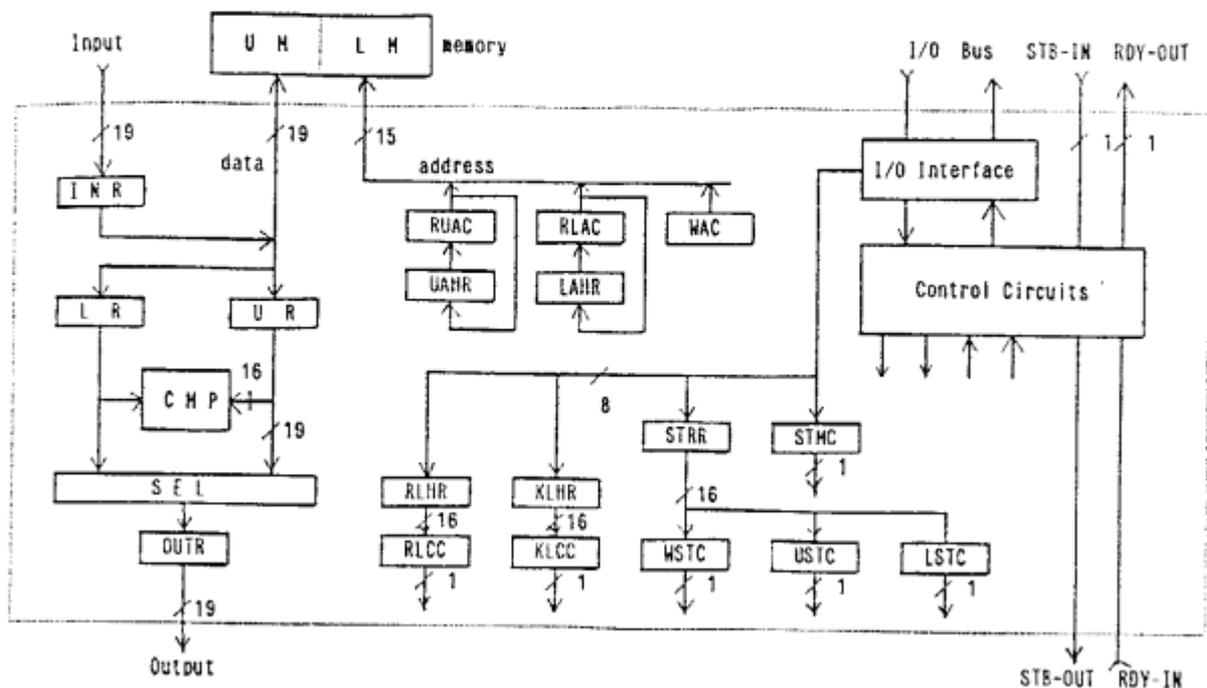


図5 セルの構造

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、ソータの入力側及び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3サイクルで1語の処理を行う。処理タイミング・チャートを図6に示す。サイクルタイムは220 nsである。

第1サイクルは、セル間の制御信号STB-IN(OUT)とRDU-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、及び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の比較及びINRデータのメモリへの書き込み動作を行う。

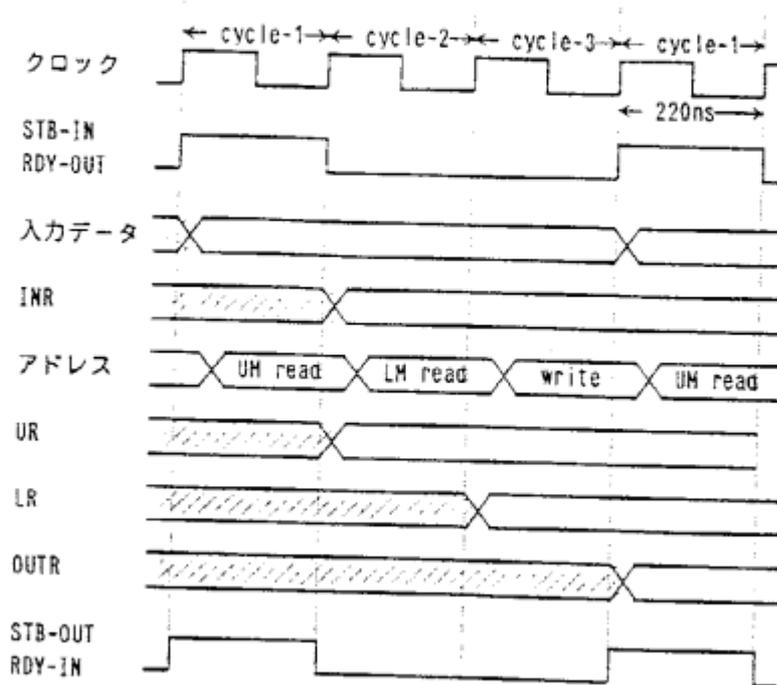


図6 処理タイミング・チャート

5. セルの状態遷移

5. 1 ソート・モードの状態遷移

このモードの動作は、次に示すような3階層構造で制御される。各階層での状態遷移を図7に示す。

- (1) 下位階層：1語単位の繰返し処理を制御する。URD、LRD、WT/Cの3つの状態が、それぞれ図6の第1、第2、第3サイクルに対応した動作を行う。
- (2) 中位階層：レコード・レベルの繰返し処理を制御する。CMP状態は等号が成立して比較結果がまだ決定していない状態で、URのデータをOUTRに転送する。比較結果が決まると、UH(LM)のレコードが選択された場合はSELU(SELL)状態へ遷移する。なお、この状態の始めにnull信号がチェックされ、一方がnull値であれば正常なキーのレコードが、両方ともnull値であればLMのレコードが選択される。SELU(SELL)状態はUH(LM)のレコードの残りの語を出力する。
- (3) 上位階層：ストリング/ストリーム・レベルの繰返し処理を制御する。S01状態は第1ストリングをUHに格納する。S02,S03状態は第2ストリングをLMに格納しつつ、両メモリのレコードの比較を行い、LMへの格納が終ると、両メモリにレコードが有る時はS04へ、UH(LM)だけに有る時はS05(S06)へ遷移する。S04状態は第3ストリングをUHに格納しつつ、

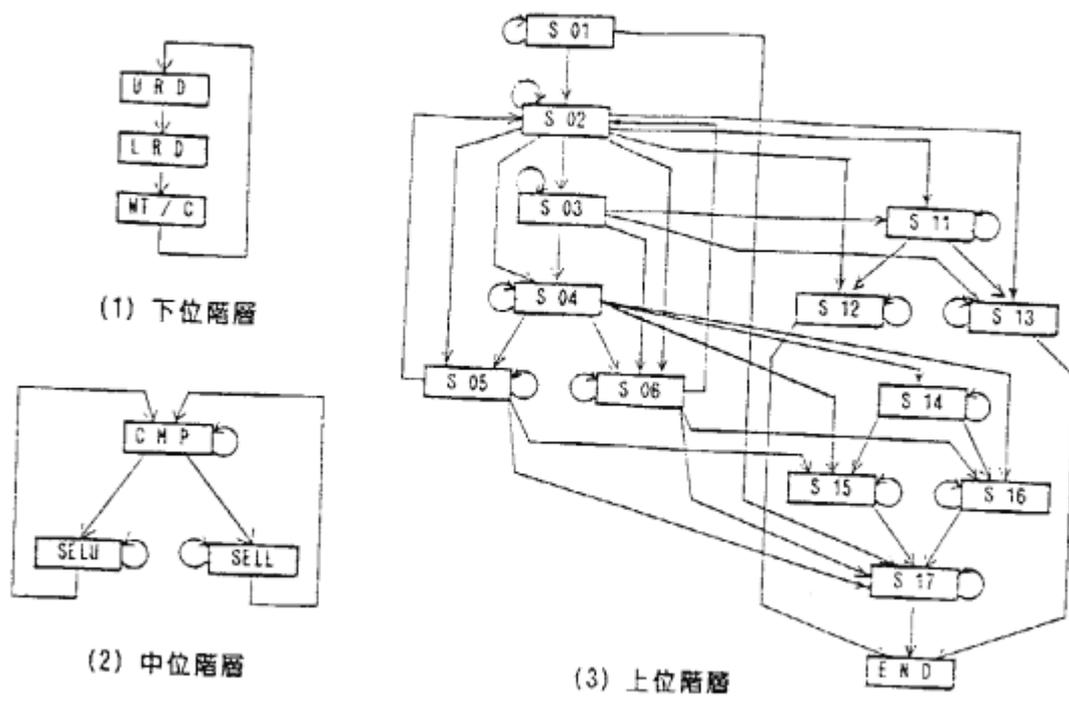


図7 状態遷移

両メモリのレコードの比較を行い、LM(LM)の比較中のストリングがすべて出力されるとS06(S05)へ遷移する。S05(S06)状態はLM(LM)のレコードの出力と入力レコードのLMへの書き込みを行い、S04 状態へ遷移する。

以下のS11 ~ S17 状態は、上記S02 ~ S06 の状態でストリームが入力され終わった時の処理を行う。

S11 状態は両メモリのレコードの比較を行い、LM(LM)にレコードが無くなるとS12(S13)へ遷移する。S12(S13)状態はLM(LM)のストリングを出力してEND へ遷移する。S14 状態は両メモリのレコードの比較を行い、LM(LM)にレコードが無くなるとS15(S16)へ遷移する。S15(S16)状態はLM(LM)の残ったストリングを出力してS17 へ遷移する。S17 状態はLMの最後のストリングを出力してEND へ遷移する。

5. 2 バス・モードの状態遷移

このモードは、上記下位階層の3つの状態遷移を次のように単純化して1階層で制御される。URD 状態で入力データのINR へのセットとOUTRデータの出力動作を、LRD 状態でINR のデータのLRへの転送を、WT/C状態でLRのデータのOUTRへの転送を行う。

従って、1語は660ns でセルを通過する。

6. 性能評価

6.1 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 $F = \lfloor CL/E \rfloor$ と表わすと、 $F = 1$ 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変わるので、処理時間の一般式を導くのは複雑になる。そこで、以下の解析では、ストリームの大きさが64k バイトの整数倍の場合を考え、 $E = 64k$ バイトとする。

(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。従って、 $T(1)$ は次式で表わされる。

$$T(1) = t_0 + t(S) + 3LT + t_2$$

ここで、 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 2.1節で求めた $T(S)$ で $CL \cdot E$ としたソータでの処理時間、 $3LT$ 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8に示す。同図において、 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。従って、 $T(2)$ は次のようになる。

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

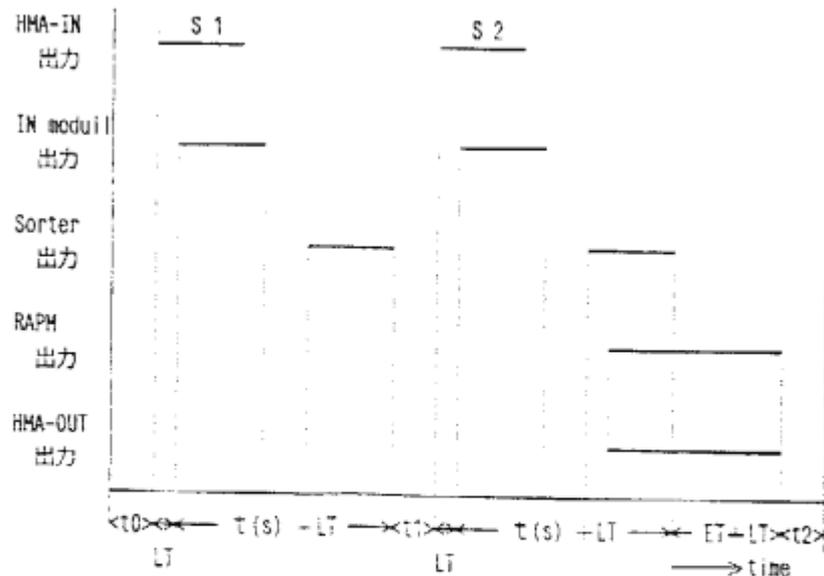


図8 各モジュールにおける処理時間

(8) $F > 2$ 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、これを $F = 2$ の方式で処理する。次に、 $S1, S2$ を E の単位に分割、即ち ($S11, S12$) 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。従って、 $T(4)$ は次のようになる。

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

ただし、 $A = (t_1 + 2t_1(s) + ET + 5LT)$ 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最悪処理時間を示す。即ち、 $S11$ 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} 回繰返し、その後上記のマージ操作を $(2^i E / 2E) \log(2^i E / 2E)$ 回繰返すことにより実行する。従って、 $T(2^i)$ は次のようになる。

$$T(2^i) = t_0 + t_2 + A 2^{i-1} + B(i-1) 2^{i-1}$$

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

6. 2 測定結果

測定は図 2 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。レコードの長さはすべて 16 バイトで、レコード全体をキーとして扱った。

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
131,072	6,159	8,160
262,144	14,511	18,390
524,288	33,423	45,860

表 1 ソート処理時間

6. 3 考察

本ソータの処理時間の計算値と測定値を表1に示す。t0 ~ t2 の値としては、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。これは、図2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t0 ~ t2 以外にRDBEとHM間で多数の内部コマンドが発行されるためのオーバーヘッドである。従って、このオーバーヘッドを減すためには、RDBEをHMの入出力機器としてHMから制御することが有効と考えられる。

7. むすび

関係データベースの処理を行う場合に必要となるレコード数、レコード長、キー長等のパラメータ、null値の取り扱い、重複レコードの取り扱い等の諸機能を考慮したソータの設計と試作結果について述べた。本セルの回路規模は、約5000ゲートであり、8000ゲートのゲートアレイで実現することができた。

本ソータは、単体で64K バイト、RAPMを連動させた時は128Kバイトのストリームを3MByte/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。従って、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 $(2^{i-1} + 1) \times L$ のメモリ容量を持たばよいので、ポインタを用いて削減することも可能である⁽¹⁰⁾。しかし、セルの制御回路が複雑になること及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。

参考文献

- (1) 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).
- (2) 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).
- (3) Knuth,D.E.:The Art of Computer Programming,Vol.3.Sorting and Searching, Reading,Addison Wesley,pp. 11-388(1973).
- (4) Barsamian,H.:Firmware Sort Processor System,U.S.Patent,3.713.107(1973).
- (5) 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,pp553-563(1980).
- (6) Kassini,D. and Sahni,S.:Bitonic Sort on a Mesh Connected Parallel Computer, IEEE Trans.Comput.,Vol.C-27,No.1,pp.2-7(1979).
- (7) 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.J65-D,No.2,pp.179-186(1982).
- (8) 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.J66-D,No.3,pp.332-339(1983).
- (9) 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).
- (10) 上肥: 大容量ファイルを整列するシストリック・ソータ, 信学論(D),Vol.J67-D,No.3,pp.281-288(1984).
- (11) Toda,S.:Algorithm and Hardware for a Merge Sort Using Multiple Processors,IBM J.RES.DEVELOP.,Vol.22,No.5,pp.509-517(1978)