

高速ハードウェアソータを用いた SQLシステムの実現

伏見信也 岩崎浩文 安藤隆朗 佐藤重雄 小宮富士夫 樋口雅宏
三菱電機 コンピュータ製作所

データベースプロセッサGREOを用いたSQLシステムを試作した。本システムは、GREOの特長であるハードウェアによる高速ソート機能を利用し、SQLによる問合せをソートを中心とした実行形式に変換して高速に処理する。このために、SQLによる問合せからGREOのソートや結合機能の組み合わせへのコンパイル方式、及びその際の最適化技法を開発した。本システムは、オフィスコンピュータMELCOM80シリーズ上で試作され、現在評価中である。本稿では、本システムの構成、GREO用のSQLコンパイルの実現方式等について報告する。

Implementation of The SQL System Based on The High-Speed Sorter

Shinya Fushimi, Hirofumi Iwasaki, Takaaki Ando,
Shigeo Sato, Fujio Komiya, and Masahiro Higuchi

Computer Works, Mitsubishi Electric Corporation
325 Kamimachiya, Kamakura-Shi, Kanagawa 247, Japan

This paper reports the design and implementation of the new experimental SQL system. The system is unique in that it is based on the high-speed sorting function provided by the database processor called GREO. The new SQL compiler is developed to fully utilize the database primitives of GREO. The overall organization of the system is reported. Also, the details of the SQL compiler and the optimization techniques oriented to GREO are described.

1. まえがき

データベースソフトウェアGREO^[1]を用いたSQLシステムを試作した。本システムは、GREOの特徴であるハードウェアソフトウェアによる高速ソートを中心に実現された検索専用SQLシステムである。本稿ではそのシステム構成、SQLによる問合せのGREO命令群へのコンパイル技法等について報告する。

2. リレショナルデータベースソフトウェアGREO

2.1. ハードウェア構成

GREOは東京大学生産技術研究所との並列データベース処理に関する共同研究の成果を製品化したものである。図1にGREOのハードウェア構成を示す。GREOは当社のオフィスコンピュータMELCOM80シリーズの付加ソフトウェアとして使用される。GREOは、ハードウェアソフトウェア部とデータベース処理部により構成される。ハードウェアソフトウェア部は、LSIソフトウェア^[2]19石をハイライン接続することにより構成され、8MB/512Kワードまでのデータを一度に高速ソートすることができる。また、データベース処理部は、3つのマイクロプロセッサ(MC68020)、ローカル/共有メモリ、ハードウェアソフトウェアインタフェース等により構成されたマルチプロセッサアーキテクチャのデータベースエンジンである。これらマイクロプロセッサ群はソート処理の前後で選択、射影等の関係演算を実行し、また2ポート化されたソータのメモリを用いることにより、ソート後のデータのマージ、結合等の演算を高速に処理する。

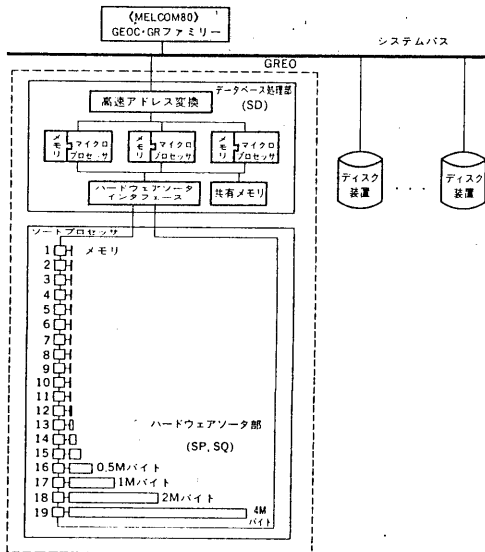


図1 GREOのハードウェア構成

2.2. ソフトウェア構成

GREOのソフトウェア構成を図2に示す。GREOの3台のマイクロプロセッサ上には専用のマルチプロセッサ/マルチプロセッサの制御OS(SY/OS)が搭載されている。GREOは固有のデータベース処理機能を持たず、ホスト計算機上に格納されているプログラムをダウンロードし、これをGREO内部の3つのCPU上のOSにより多重プロセッサとして動作させることによって所用の機能を実現する。このダウンロードプログラムを拡張機能と呼び、対応するプロセッサをフィルタと呼ぶ。このフィルタに対してディスク上のデータを連続的に入力することにより、データの選択、結合等を実行する。現在、ソート(sorting)、マージ(merger)、索引生成(index generation)、結合(join)等の機能を実現する拡張機能プログラムが用意されている。以下ではコンパイルでの慣例に従って、これらを各々ソート命令、結合命令等と呼ぶ。

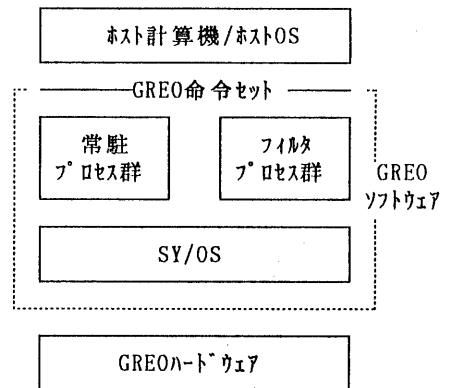


図2 GREOのソフトウェア構成

2.3. 命令セット

本試作では、GREOの命令の内、ソート命令、マージ命令、結合命令のみを用いた。図3にGREOのソート命令、結合命令の構造を示す。

2.3.1. ソート命令

ソート命令は、入力されてくるデータストリームを指定されたキーによりソートして出力する。ソート命令では、ソートの前処理としてデータストリームに対する選択、射影、ソートキーの抽出/生成機能、またソートの後処理としてソート結果に対するグループ化、グループ化により生成された各グループに対する集合関数(aggregation)計算、重複値処理、

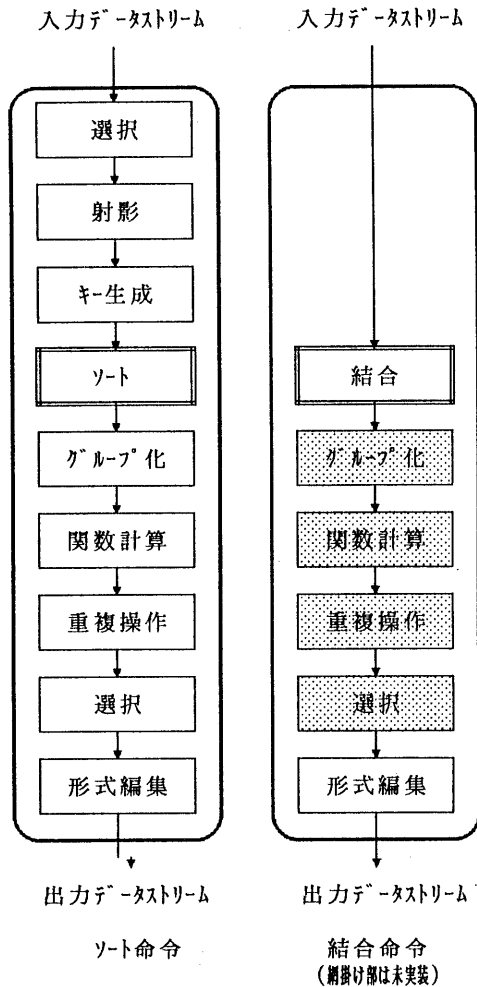


図3. GREOのソート/結合命令

グループ化後の値に対する選択機能、及びこれらにより得られた個々の値を結果レコードの形式に編集する機能が各々指定可能である。一度にソート可能なデータ容量は現在8MBであるが、これを超えるデータのソートの際には8MB単位の連(string)を連続的に生成し、後にこれらをマージ命令により併合することによりソートを完成させる。マージ命令に対しては、ソート命令と同様な機能が指定可能であり、同様の議論が成立するため、ここでは詳細を省略する。

2.3.2. 結合命令

結合命令は、入力されてくる複数のソート済みのデータストリームを指定されたキーにより結合す

る。結合命令は、ソート命令と同様な後処理が指定可能となるよう設計されているが、集合関数計算等、これら後処理の機能の一部は結合キーとグループ化処理のキーが同一の場合のみ有効なものであり、使用頻度が低いとの判断から本試作では使用されていない。

3. システムの概要

図4にシステムの構成を示す。本システムは、ワイルドビューAMELCOM80上に実現されている。以下、その主要構成要素について概略を述べる。

3.1. SQLのオブジェクトコード

SQLのコンパイル結果を記述するオブジェクトコードをプラン(plan)と呼ぶ。プランは、(一般に複数の)オペレータ、親言語の変数値、オペレータのプラン内での実行順序などの制御情報からなる。プランは、テキスト部とデータ部からなり、再入可能に構成されている。主としてオペレータ群とその実行順序に関する情報がプランテキスト部におかれ、またコンパイル時には確定しない親言語の変数値記述や、検索した結果を表示するためのフォーマット情報がプランのデータ部に保持される。

オペレータ(operator)は、GREOやSQL中核による処理の基本単位である。この内、ソート及び結合オペレータは、主としてGREOに対するソート命令、結合命令を記述する。また、中間ファイルの生成、削除、レコード単位のフェッチ等のオペレータが定義されている。オペレータは、更新機能等、将来のSQLの完全仕様の実装を考慮して設計されている。

3.2. SQLコンパイラ

SQLコンパイラは、SQL文を解析し、これをプランに変換する。SQLコンパイラは、与えられたSQL文から必要なオペレータ群を生成し、実行順序や同時実行可能な演算を解析して、結果をプランテキストに格納する。SQL文と、それをコンパイルして得られたプランの構造の例を図5に示す。ここでは、後述する選択述語の分配による最適化が仮定されている。これらコンパイラの実現方式の詳細は4節で述べる。

3.3. SQL中核

SQL中核は、SQLコンパイラによって生成されたプランを制御、実行する。SQL中核は、プランの

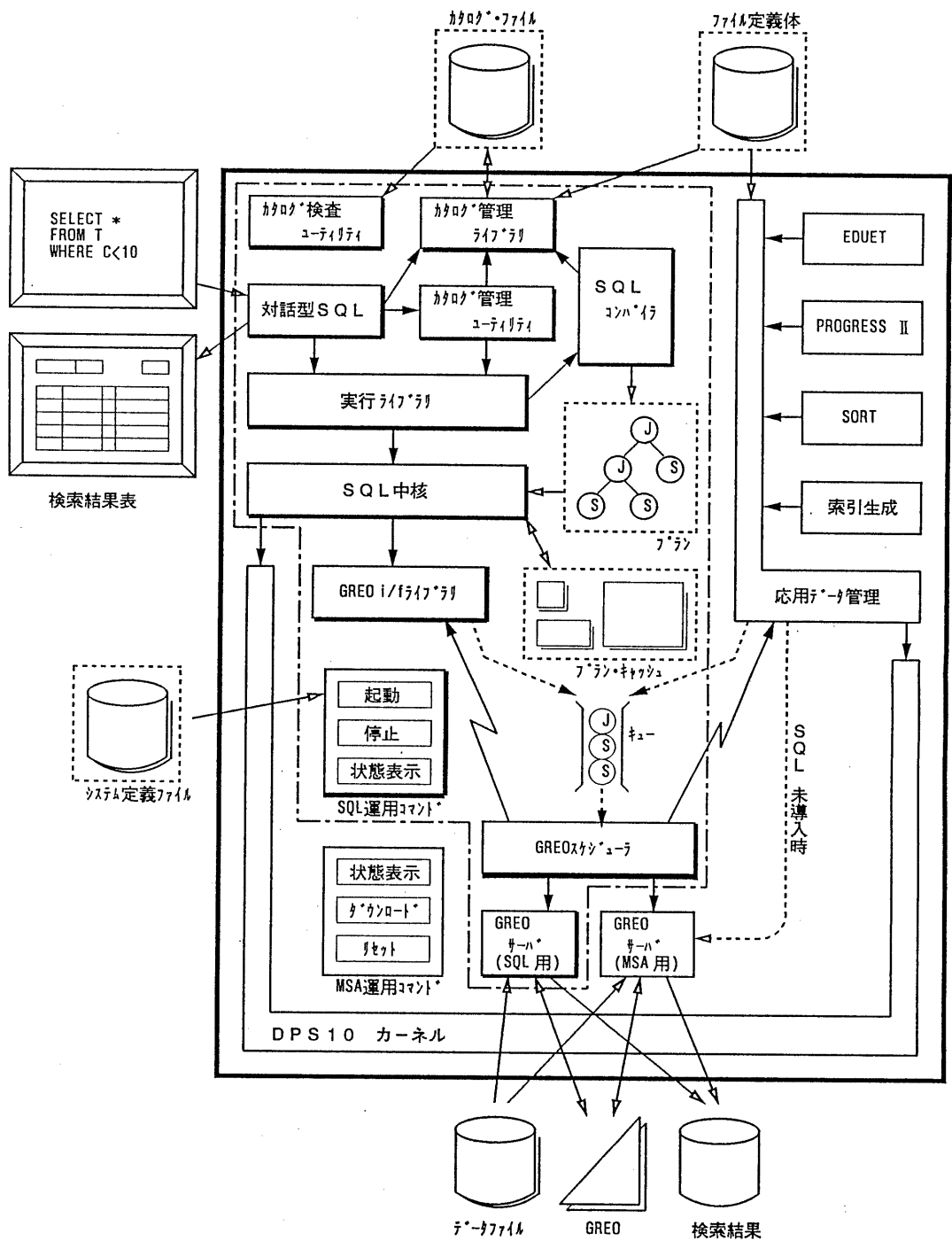


図4. SQLシステムの構成

○問合せ：

```
SELECT * FROM T, U
WHERE T.A = U.B AND U.B > 100
      AND T.A+U.A = 10
ORDER BY T.C ----- ①
```

○プランの構造：

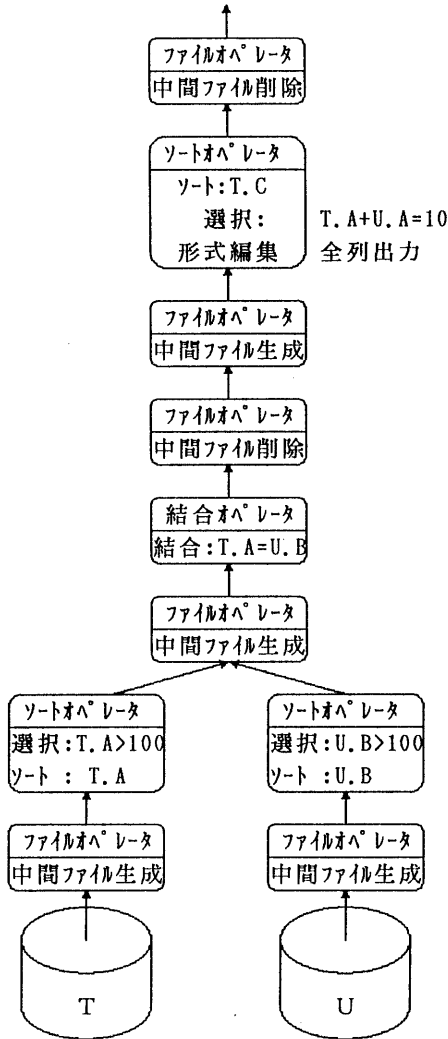


図5. SQL文とプラン

実行を開始する前に、まずそのテキスト部をプランキャッシュ(図4)にロードする。さらに、データ部をアプリケーション毎に割付けられた領域に格納する。テキスト部を共通空間中のキャッシュに配置することにより、複数のユーザが同一のプランを実行する

場合にテキストを共有することが可能となる。定型業務等では、同一のSQL文が複数のアプリケーションから実行されることが多く、このような場合には一旦ロードしたプランをキャッシュに保持し、共有することによって実行時のオーバーヘッドを軽減することが可能となる。

SQL中核はデータ駆動の原理によりプラン中のオペレータ群の実行を制御する。即ち、まずプラン中の実行可能な葉(leaf)オペレータを検索し、これに対してGREGOスケジューラに処理要求を発行する。その後、これら要求に対するGREGOによる処理が終了し、対応するオペレータの親オペレータが実行可能となると、親オペレータの処理を起動する。最終的に全てのオペレータの実行が完了すると、目的とする検索結果が得られる。

3.4. GREGOスケジューラとGREGOサーバ

GREGOスケジューラは、SQL中核からの要求に対しGREGOの割付けを制御するGREGOの資源管理スケジューラである。また、GREGOサーバは、GREGOの提供する各命令毎にその実行を制御するデータベース機能単位のサーバである。

GREGOスケジューラは、SQL中核より起動要求を受け取ると、使用可能なGREGOを探し、要求されている機能に対応したGREGOサーバを起動する。複数の要求に対しては優先順位により実行順序を決定する。GREGOサーバは、検索対象ファイルからGREGOへのデータ転送、GREGOからの検索結果のファイルへの転送等を制御する。処理が終了すると、これをGREGOスケジューラに通知する。GREGOスケジューラは当該処理の完了をSQL中核に対して伝えると共に、起動要求が待ち行列にあれば優先順位に従って対応するGREGOサーバを再び起動する。

本GREGOスケジューラにより、GREGOが複数台搭載されているシステムでは、複数のオペレータを同時に処理することが可能となる。例えば図5に示されたプランを単独で実行した場合には葉にあたる二つのソートオペレータは同時に処理される。また、GREGOの一部が故障した場合には、GREGOスケジューラより縮退運転が実現される。

4. SQLコンパイルの実現方式

4.1. 概要

SQLコンパイルは、構文解析、実行木の生成、最適化、コード生成の各部により構成されている。この内、構文解析等は通常のコンパイルと同

様であり、ここでは記述を省略する。

本コソバ[®] には以下の特徴がある。

- (1) 全てのSQLによる検索要求を上記の2命令の組み合わせにコソバ[®] 化する。特に、結合演算の実現アルゴリズムとしては、ソート/マージ結合のみを対象としている。
- (2) 可能な全ての等価変形に対し、問合せ処理途中で表同志の直積が必要となる表の集合を検出し、これが存在すればI²ラ[®]とする。これは通常のアプリケーションでは直積を生成することは稀であり、また直積生成のために長時間GREQOを使用することは資源管理の点で有利ではないことによる。
- (3) 等結合により結合される表の集合に対し、可能ならばその集合中の表に指定されている選択述語を等結合される全ての表に分配し、なるべく多数の表に対して予め入力データの絞り込みを行う。
- (4) 複数の2-way結合を単一のn-way結合、nキー結合に変換し、GREQOへの一度のデータストリーム転送により結合演算を実行可能とする。
尚、本システムで対象としたSQL言語は、JISのSQL仕様^[4]を基本とし、一部SQL2^[5]の仕様を追加した。また、入れ子にわたる列の参照を含む副問合せに関しては、プリプロセッサによる標準化を仮定しており、ここでは議論しない。以下ではSQLの仕様に従い、表(table)、列(column)、行(row)等の用語を用いる。

4.2. 実行木の生成

実行木(Execution Strategy Tree)とは、問合せが参照している表の集合を結合演算により結合し、一つの出力表とするための結合順序を決定すると共に、これら各結合演算に付随して実行される選択演算等を確定するための木構造である。木の各ノードは、演算結果としての表を表し、またその表を生成する結合述語、結合結果に対する選択述語を持つ。ここではノードが表現する表、または結合結果の表の集合をノードにラベルとして付加し、ノードの結合述語、選択述語を各々「J:」、「S:」により表す。本コソバ[®] における実行木の生成は、通常のコソバ[®] の意味解析に対応する。

実行木生成の基本アルゴリズムは、良く知られている関係代数木生成アルゴリズム^[6]を拡張したものである。まず、WHERE句が比較述語等の基本述語のAND結合のみから構成される場合を考える。図5のSQL文①を例に説明する。

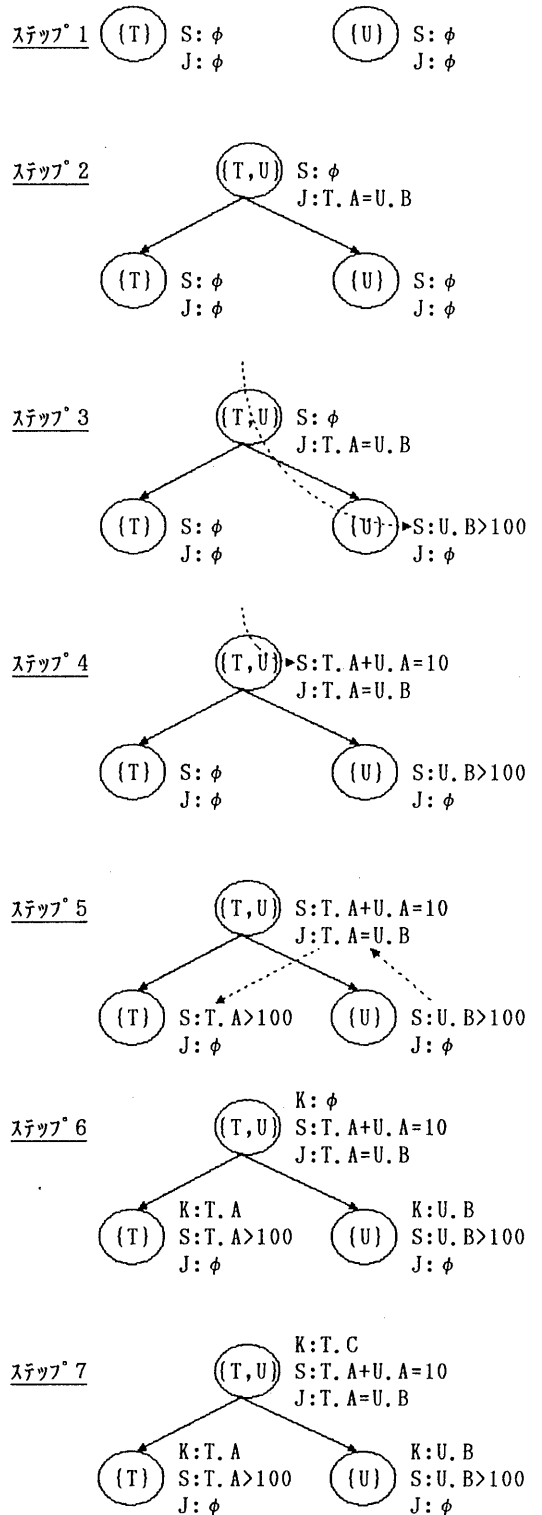


図6. SQL文のコソバ[®] 化

最初にFROM句内の各表に対応してノードを生成し、これを初期化する。この例では図6ステップ1に示すノードの森(forest)が得られる。

次に、WHERE句中の各述語を次々と取り出し、述語の各々の辺が参照する表(の集合)をこの森から探す。この際、ANDで接続された基本述語の集合は実行順序に関して可換であり、適用順序に依らず同一の結果が生成されることが基本となる。この場合、 $T.A=U.B$ を例にとると、この述語の左辺、右辺が参照する表の集合は各々{T}、{U}であり、各々の表集合を含むノードが森に存在する。これは $T.A=U.B$ がこれらノードが表現する表を結合することを意味する。このような述語が存在すれば新たにノードを作る。このノードは結合演算の結果を表し、結合されるノードを新しいノードの子に設定し、そのJ:フィールドに当該述語を設定する。この例ではTとUの結合結果を表現する新しいノードが生成され、TとUに対応するノードが子に設定される。また、新しいノードのJ:に述語 $T.A=U.B$ が付加される。このノードは、TとUの全ての列を含むと仮定される。これにより実行木は図6ステップ2に示すものとなる。

次に、 $U.B>100$ を考える。この述語は{T,U}上の述語であるが、ノード{T,U}の右の子(U)の列のみを参照している。これから実行木上でこの述語を{T,U}から{U}に移動することができる。(U)は葉であり、これ以上述語は移動できず、この述語はノード(U)の選択述語としてS:に設定される。これは、選択述語を先に実行し、結合演算実行以前にデータの絞り込みを行うことを意味する(図6ステップ3)。

最後に、 $T.A+U.A=10$ は{T,U}上の述語であるが、その左辺は左右両方の子の列を参照しており、これ以上実行木上を移動できない。また{T,U}のJ:には既に結合述語が設定されているため、この述語は選択述語としてそのS:に付加される。これは、この述語がJ:に設定されている述語による結合結果に対する選択演算であることを意味する(図6ステップ4)。

WHERE句にNOTを含む場合には、De Morganの法則によりこれを変形し、基本述語の論理を反転することで全体をNOTの無い形に変形する。WHERE句にORを含む場合は、基本述語をORで接続した和項(disjunctive)をANDで接続した標準形(disjunctive normal form)に変形し、これらと和項群に対して上記の基本アルゴリズムを適用する。ここで和項により等結合

がなされる場合があることに注意する。例えば、①のWHERE句が

$$「T.A+1 = U.A \text{ OR } T.B*2 = U.B/U.C」$$

の場合、この述語全体により表TとUが結合される。これらと和項による結合は、和項中の各述語による結合を個々に実行し、これらの結果をUNIONすることにより実現される。これから判るように、和項内の全ての基本述語の左辺、右辺が各々同一のノードを参照するように和項内の基本述語の両辺を交換することができる場合に限って和項は結合述語となる。例えば、上の例のWHERE句が

$$「T.A+1 = U.A \text{ OR } U.B/U.C = T.B*2」$$

である場合、コンパイルは、和項内の2番目の基本述語の左辺と右辺を交換することにより、和項内の全ての基本述語の左辺がTを、また全ての述語の右辺がUのみを参照していることを検出する。

更に、SQLでは、これらORによる結合に対して種々の表現が可能なる点に注意する必要がある。例えば、①のWHERE句として

$$「T.A+1 \text{ IN } (U.A+1, U.B-U.C)」$$

と

$$「T.A+1 = U.A+1 \text{ OR } T.A+1 = U.B-U.C」$$

は等価であり、共にTとUを結合する結合述語となる。また、同様に

$$「T.A - U.A = 0」 \text{ は } 「T.A = U.A」$$

と等価であり、結合述語と見なせるが、実用上の観点から本システムでは結合述語判定の等価変形の対象としていない。

4.3. 最適化

本コンパイルは生成された実行木に対して以下の最適化を行う。

- (1) 実行木中の各ノードに対し、その結合述語が和項でない場合、その任意の子のS:フィールドに設定されている選択述語の内、元のノードのJ:にある結合述語のいずれかの辺を対象としているものがあればこれを結合に関与

する全てのノードに分配する。例えば、上の例①では結合の右辺はU, Bであり、またこれを対象とした述語U, B>100がその右の子に存在するから、これをその左の子に分配することができる(図6ステップ5)。

- (2)実行木中の各ノードにおいて、J:とS:が共に設定されている場合には、これらの述語群を比較し、入れ換えを行って、最も結合による絞り込み効果が高いと判断される述語を結合述語としてJ:に再設定する。例えばJ:, S:が

J:T.A=U.A OR T.B = U.B
S:T.C=U.C

であれば、J:に設定されている結合述語にはUNIONが必要とされるため、S:にある述語の方が絞り込み効果が高いと判断でき、結合述語との入れ換えを行う。

- (3)ノード群のJ:述語を単一のn-way結合に併合する。例えば、J0:T.A=U.A, J1:U.A=V.AをJ:T.A=U.A=V.Aに変換する。尚、この変形は本試作では実行木生成の前処理として実装されている。
- (4)各ノードに対し、J:が設定されている場合には、これをそのS:中の述語群と結合述語を併合し、単一のnキ結合に変換することを試みる。例えば、

J:T.A=U.A → J:(T.A, T.B)=(U.A, U.B)
S:T.B=U.B S:φ

なる変形が可能である。

4.4. コード生成

以上により生成された実行木が、コード生成によりGREOのソート、結合命令に変換される。上述のように、これら命令には種々の補助機能が指定可能であり、コード生成ではこれら命令中の補助機能をなるべく有効に利用し、生成する命令の数の減少させることを試みる。例えば、実行木中のあるノードの結合演算の実行のためには、結合キーに対するソートが必要であるが、このソートの際に、ソート命令の補助機能を利用して、子ノードの選択演算が実行可能である。このように、コード生成では、命令で指定可能な補助機能群に実行木に設定されてい

る演算を最密充填することが目標となる。

このために、実行木のノードに対して、ソートキーフィールド(K:)を追加し、各ノードでは、結合→選択→ソートの順に指定された演算を実行するものとする。例えば上の例に対する実行木は、結合の前処理としてのソートを考慮すると図6ステップ6のように変形される。また、ORDER-BYやGROUP-BY、HAVING等が指定されていた場合にはこれらを実行木に追加する。この例ではORDER-BYに指定された列を根ノードのK:に設定する(図6ステップ7)。

この実行木に対して、木の根から始めて再帰的にノードを辿り、各ノードに対するコードを生成する。例えばK:とS:が共に設定されていれば、S:に設定されている選択演算をK:に相当するソート命令の補助機能として実行可能で、この場合には一つのソート命令でK:とS:の処理が可能となる。一方、K:がなく、S:が指定されている場合には新たにソート命令を生成する必要がある。図6に示された実行木に対して得られたコードが図5に示されたプランである。

本システムでは、これらK:, S:, J:の全ての可能な組み合わせに対して生成される命令のプランを予めコンパイルがデタベ-ス化して記憶しており、table drivenな方法によりこれらコード生成を実現している。また、このデタベ-スを変更することによりGREOの命令セットの拡張にも容易に対処することが可能である。

5. おわりに

GREOを用いた試作SQLシステムの概要と、そのコンパイルの実現方式について報告した。本システムは、当社のMELCOM80シリーズでの実用化に向けて更に評価、改良を行う予定である。

- [1]安藤他「リリショナルデタベ-スプロセッサGREOの構成」, 電子情報通信学会, DE89-37 1989
- [2]伏見他「LSIソートプロセッサ」, 電子情報通信学会, DE88-2 1988
- [3]井上他「デタベ-スプロセッサRINDAの設計と実現」, 情報処理学会論文誌, 31, 3 1990
- [4]JIS デタベ-ス言語SQL/X3005, 1990
- [5]ISO/ANSI, 「Database Language SQL2 and SQL3」 X3H2-90-001, 1990
- [6]Ullmann, J.D. 「Principles of Database Systems」, Computer Science Press 1982