

5 H-5

Symmetry S81 における GRACE HASH 方式の実装と評価

津高 新一郎 中野 美由紀 喜連川 優 高木 幹雄

東京大学 生産技術研究所

1はじめに

関係データベース処理の中で、結合演算は、選択演算などの他の関係演算に比べて処理負荷が重いことは良く知られており、その処理負荷を軽減すべく今まで種々の結合演算処理方式が提案されてきた[1,2,3,4,5]。なかでもハッシュ操作に基づく結合演算処理は、従来のソート処理に基づく方法に比べて高い性能が得られ、数々の研究成果が発表されている。

一方、データベースマシンの分野では、処理性能の向上を目指し、並列処理技術を取り入れたアーキテクチャが多数考案されている。当研究室では、共有メモリマルチプロセッサマシン上での関係演算処理方式として、最も処理負荷の重い結合演算の実装について検討している。32 MB の共有メモリ、18 台のプロセッサ、4 台のディスクから構成される Sequent 社の Symmetry を用い、この上で GRACE HASH 方式を用いた結合演算の実装を行なった。本稿ではその実装方式と評価について述べる。

2 GRACE HASH 結合演算処理方式

ハッシュ操作に基づいた結合演算処理とは、演算処理対象を独立した複数の空間に分割し、結合処理に必要なデータ検索空間を狭めることで、演算処理の高速化を図る方法である。以下、GRACE HASH 結合演算方式の処理の流れについて述べる。本方式は、次のような 2 つのフェイズから構成される。(結合演算の対象となるリレーションを R、S とする。)

分割フェイズ リレーション R の各タブルの結合属性にスプリット関数を適用してパケットに分割し、中間リレーション Ri に書き出す。続いて、リレーション S も同様にしてパケット分割する。

結合フェイズ 各パケット毎に結合演算処理を実行する。異なるスプリット値を持つパケット同士は結合される可能性が無いため、等しいスプリット値をとるパケット同士が主記憶上にステージングされて結合処理が施される。

3 GRACE HASH 方式の実装

当研究室の Symmetry S81 は、32 MB の共有メモリ、18 台のプロセッサ ($i_3 86$ 、16 MHz) 4 台のディスクから構成される。4 台のディスクは各々別のチャネルに接続され、並列にアクセスすることが可能である。

GRACE HASH 方式を本マシン上に実装する際には、共有メモリ上のデータに対するマルチプロセッサによる並列処理効果を考慮すると共に、入出力処理に対する並列処理方式についても留意しなければならない。そこで、今回の実装方式の特徴を以下に示す。

- 分割フェイズにおけるスプリット関数によるパケット分割、あるいは、結合フェイズ時の結合処理等のデータ処理を行なうプロセッサ群とディスクの入出力アクセス処理を行なうプロセッサ群を完全に分離し、各々の並列処理効果を最大限に得る。
- リレーションを複数台のディスクに分割して置き、並列アクセスによる入出力の高速化を図る。
- 入力用のプロセッサ、出力用のプロセッサとしてディスク 1 台につき各々 1 台のプロセッサを割り当て、入出力管理を専用プロセッサで集中して行なう。

- 入出力バッファをダブルバッファ化し、一方のバッファをディスクとのデータ転送用、もう一方をプロセッサとのデータ転送用とすることで、プロセッサの入出力待ち時間を減らす。
- 各パケット毎の結合処理では、ハッシュ関数を用いて小さなデータクラスタ群を生成し、プロセッサが並列に各データクラスタ毎に処理を行なう。

以下に、GRACE HASH 方式の実装について各フェイズ毎に詳述する。但し、結合演算の対象となるリレーションを R、S、入力用プロセッサを Pi、出力用プロセッサを Po、その他のプロセッサを Px とし、ディスクの台数を d 台とする。

分割フェイズ(図 1)

- ディスク上のリレーション R を d 個のプロセッサ Pi が読み出し、入力用ダブルバッファに格納する。
- n 個のプロセッサ Px は入力バッファ中のタブルを並列にローカルメモリにコピーした後、結合属性にスプリット関数を適用し、その値に応じた出力用ダブルバッファにコピーする。この時バッファがいっぱいになると、書き出し要求メッセージを出すと同時にバッファの切り替えを行う。
- d 個のプロセッサ Po は書き出し要求メッセージを受けて、出力用ダブルバッファの内容を部分リレーション Rj としてディスクに書き出す。
- リレーション S に関しては同様に 1~3. の処理を行なう。

結合フェイズ(図 2、3)

- ディスク上の部分リレーション Rj を Pi が読み出し、入力用ダブルバッファに格納する。
- Px は入力バッファ中のタブルをローカルメモリにコピーしたのち、結合属性にハッシュ関数を適用し、その値に応じた共有メモリ領域にコピーする。
- 1~2. の処理を Rj のタブル全てについて行なう。
- 次に、Rj に応する部分リレーション Sj を Pi が読み出し、入力用ダブルバッファに格納する。
- 各 Px は入力バッファ中のタブルを並列にローカルメモリにコピーしたのち、結合属性にハッシュ関数を適用し、その値に応じた共有メモリ領域を参照する。結合属性値が一致すると、両方のタブルは出力用ダブルバッファにコピーされる。
- Po は出力バッファを監視しており、いっぱいになった時点で結果リレーションとしてディスクに書き出す。
- 1~6. の処理を全ての部分リレーション Rj、Sj について行なう。

4 性能評価

タブル長 208 バイト、結合属性長 4 バイト、属性値がユニークで順番がランダムであるようなリレーションに対し、4 kB バイトのページ 256 個からなる 1 M バイトのメインメモリをステージングバッファとして用いてコストの測定を行なった。

まず、入出力以外に用いるプロセッサの数が 1 台の時に、ディスクの台数とリレーションのサイズの変化させてコストを測定した。その結果を図 4. に示す。処理コストはリレーションの件数に比例して増加しており、Symmetry 上で我々の提案した GRACE HASH 実装方式が有効であることが確認できる。また、ディスクの並列アクセスに関しては台数を増やしてもほぼ処理コストはリレーションの件数に比例して増加しており、GRACE HASH 方式の性能として理想的な結果が得られている。

⁰The Implementation and Evaluation of the GRACE HASH Algorithm on Symmetry S81
S.Tsudaka, M.Nakano, M.Kitsuregawa, M.Takagi
Institute of Industrial Science, University of Tokyo

また、リレーションの大きさを2万件とし、ディスクの台数と入出力以外の処理に用いるプロセッサの数を変化させた。その結果を図5に示す。この図から、ディスクによる台数効果がみられ、ディスクの台数の増加に従って処理性能はほぼ線形に伸びている。しかし、プロセッサが9台の時、ディスク2台では約1.6倍、ディスク4台では約3倍しか処理性能が向上していないことがわかる。これは、リレーションをディスク毎に分割した結果、各ディスクをアクセスする入出力プロセッサ間で同期をとる必要が生じ、オーバーヘッドが生じるためであると思われる。また、プロセッサが1台の場合、9台の時と比べてディスクの台数効果はあまり顕著でない。これは、ディスクの複数化によって入出力コストが軽減されると、入出力以外のコストが支配的になってくるためと思われる。

さらに、共有メモリ上のデータ処理に対するマルチプロセッサの並列処理効果を明確にするため、リレーションを1万件とし、ディスクとのデータ転送を行なわない場合のハッシュを用いた結合処理フェーズの性能測定を行なった。その結果を図6に示す。プロセッサの増加に従い、処理性能はほぼ線形に伸びていることが確認でき、Symmetry上での並列処理が結合処理コストの軽減に効果的なことがわかる。

5 おわりに

GRACE HASH方式のSymmetry上における実装方法とその性能について報告し、ハッシュを用いた処理方式が有効であることが確認された。関係データベース演算の中でも最も負荷の重い結合演算に対して、ディスクの台数効果、プロセッサの台数効果とともに明確に認められ、Symmetryが関係データベースシステムを構築するうえで極めて高い性能を持っていることが確認された。

現時点ではディスクとのデータの転送にシステムコールを使っているため、入出力はOSの介在によりバッファリング等の管理を受け、そのコストを詳細に評価するのは困難である。従って、OSによる管理を受けないローデバイスを使用した実装方式を現在考慮中である。また、結合演算だけでなく、集計演算等のハッシュを用いた処理方式が有効な関係データベース演算についてもSymmetry上で検討を行なう予定である。

参考文献

- [1] 喜連川優、他。「動的処理パケット選択方式による結合演算処理の詳細評価」、情処学会論文誌第30巻、1989。
- [2] 喜連川優、他。「動的処理パケット選択手法に基づくハッシュ結合処理方式とその性能評価」、情処学会論文誌第30巻、1989。
- [3] M.Kitsuregawa, et al. "The Effect of Bucket Size Tuning in the Dynamic Hybrid GRACE Hash Join Method", VLDB 89, 1989.
- [4] D.J.Dewitt, R.Gerber. "Multiprocessor Hash-Based Join Algorithms", VLDB 85, 1985.
- [5] L.D.Shapiro. "Join Processing in Database Systems With Large Memories", ACM TODS, Vol.11, No.3, 1986.

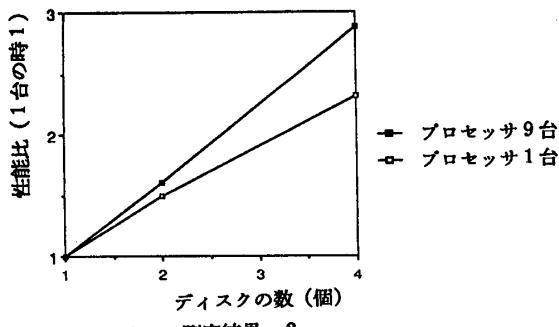


図5: 測定結果-2

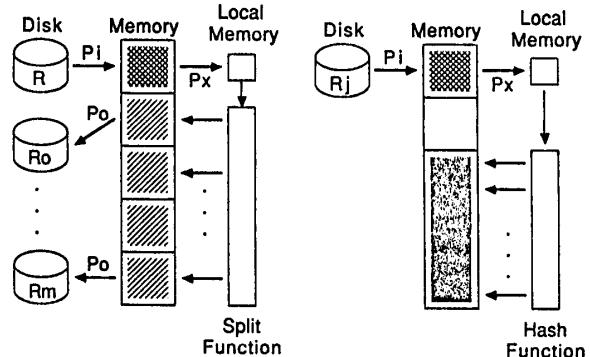


図1: 分割フェーズ

図2: 結合フェーズ-1

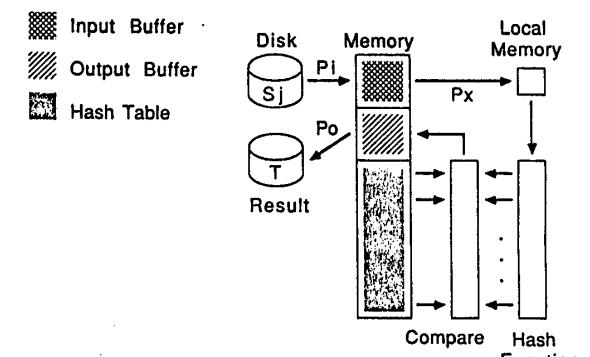


図3: 結合フェーズ-2

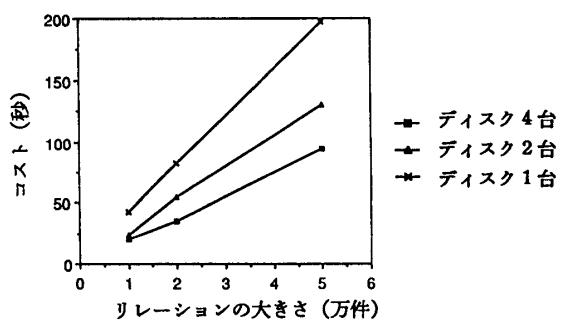


図4: 測定結果-1

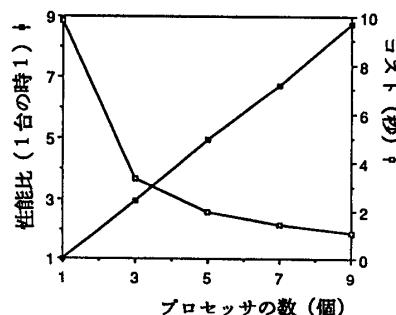


図6: 測定結果-3