

4 L-10

Symmetry S81 における結合演算の並列処理の性能評価

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してきた。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[1, 2, 3, 4]。

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、並列処理技法を取り入れたアーキテクチャが多数考案されている。我々は、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、8台のディスクから構成される Sequent 社の Symmetry を用い、この上に並列処理技法を取り入れた結合演算処理方式を実装した[5, 6]。本稿では、その実装方式と性能評価について述べる。

2 実装の特徴

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 386, 1.6 MHz) 8台のディスクから構成される。8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。

実装する結合演算処理方式としては、Grace Hash 方式を採用する。これは大規模なリレーション同士の結合演算に適したアルゴリズムであり、ハッシュ操作に基づいて演算処理対象を独立した複数の空間に分割し、結合処理時の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、というのがその基本的な考え方となっている (図 1, 2, 3)。

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。

プロセッサ割り当て プロセッサを2種類、すなわちディスクの入出力アクセス処理 (図 1, 2, 3 で濃い矢印) を行なうプロセッサ群 (入出力プロセッサ) と、その他のメモリ上での比較、転送などのデータ処理 (図 1, 2, 3 で薄い矢印) を行なうプロセッサ群 (ジョインプロセッサ) に分離する。入出力プロセッサは各ディスクに1台ずつ与えられ、入出力の集中管理を行なう。ジョインプロセッサと入出力プロセッサはリード/ライトキューによってデータの通信を行なう。これらによりディスクのノンストップ化と入出力処理とオンメモリ処理の両方における高度な並列化を実現する (図 4)。

ファイル構成 水平分割によりリレーションを複数台のディスク上へ分割して配置し (ストライピング)、並列にアクセスを行なって入出力を高速化する。また、入出力コストを詳細に評価するため、ローデバイスを使用した入出力を採用する。

3 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、64 k バイトのページ 256 個からなる 16 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、ディスク台数、プロ

セッサ数、リレーションサイズをさまざまに変化させて結合演算を行ない、そのコストの測定を行なった。

3.1 リレーションの大きさと総コスト

まず、ディスクが 8 台のとき、10万件から 100 万件の大きさを持つリレーションに対し結合演算を行ない、そのコストを測定した。そのコストと、入出力にかかるコストを図 5 に示す。処理にかかる総コストはリレーションの件数にほぼ比例して増加し、また、そのうちのほとんどを入出力コストが占めており、オーバーラップ処理がほぼ完全に行なわれていることから、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。

3.2 ディスクの台数効果

次に、ディスクの数を 1 台から 8 台まで変化させた場合に、10 万件のリレーションの結合演算を行ない、その性能比を求めた (図 6)。ディスク 4 台まではほぼ線形に性能向上が見られるが、それ以上の台数では急激に性能が低下し、ディスク 8 台では 6 度程度である。その主な原因としては、入出力の単位を 6.4 k バイトと比較的大きくとてあるため、ディスク間でのリレーションの分割が不均等になるという点や、ディスク 1 台当たりのメモリが少くなり、入出力バッファを圧迫するということ、また、ディスクの台数の増加により入出力時間にばらつきが生じ、一番遅いディスクがボトルネックとなってコストが上昇することなどが考えられる。

3.3 ジョインプロセッサの影響

総コストはディスクの台数、すなわち入出力プロセッサの数だけでなく、ジョインプロセッサの数にも依存する。10 万件のリレーションを用い、ディスクの数とジョインプロセッサの数を変えて結合演算処理を行ない、そのコストを測定した (図 7)。ディスクが 8 台のときはジョインプロセッサが 4 台以上、4 台のときは 2 台以上でコストがほぼ一定となっている。ディスクが 2 台、1 台の時はコストはジョインプロセッサの数にはほとんど依存しない。このことから、ディスクの性能を十分に生かすためにはディスク (入出力プロセッサ) 1 台あたりおよそ 0.5 台のジョインプロセッサが必要であることがわかる。

4 おわりに

Grace Hash 方式の Symmetry 上における実装方法を述べるとともにその性能の評価を行ない、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が明確に認められ、共有メモリ型マルチプロセッサーアーキテクチャが関係データベースシステムを構築するうえで極めて有効であることが確認された。また入出力を管理するプロセッサとその他の処理をするプロセッサの数をさまざまに変化させて測定し、最大の効率を与えるためのそれらの適正な比を求めることが出来た。

現在、Symmetry の上で Dynamic Grace Hash の実装が進んでいる。これは、ステージングバッファをさらに有効に利用することによって入出力コストを低下させることを目的とした、Grace Hash アルゴリズムの改良版である。また、類似のアルゴリズムとして、Hybrid Hash アルゴリズムがある。これらの実装を完成させ、リレーションのサイズによってどのようなアルゴリズムが最適なのか、またその時のプロセッサの割

り当てはどうあるべきか、といったことについて今後研究を進めていく予定である。

参考文献

- [1] 喜連川優,他.「動的処理バ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」, 情報学会論文誌第30巻, 1989.
- [2] M.Kitsuregawa,et al."The Effect of Bucket Size Tuning in the Dynamic Hybrid GRACE Hash Join Method", VLDB 89,1989.
- [3] D.J.Dewitt,R.Gerber."Multiprocessor Hash-Based Join Algorithms", VLDB 85,1985.
- [4] L.D.Shapiro."Join Processing in Database Systems With Large Memories", ACM TODS, Vol.11, No.3, 1986.
- [5] 津高、中野、喜連川、高木。「Symmetry S81におけるGRACE HASH方式の実装と評価」、情報処理学会第40回全国大会、1990。
- [6] 津高、中野、喜連川、高木。「Symmetry S81における結合演算の並列処理に対する考察」、情報処理学会第41回全国大会、1990。

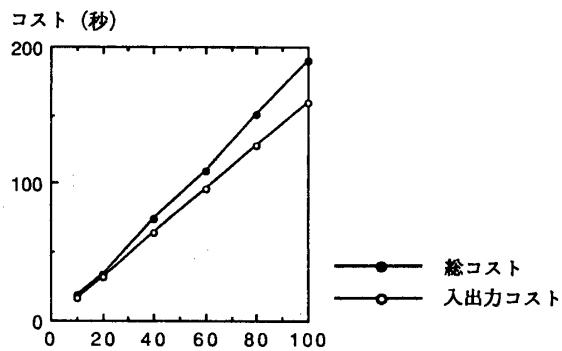


図 5. 測定結果 1

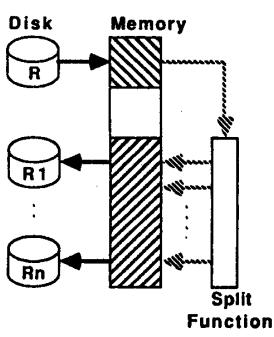


図 1. Split Phase

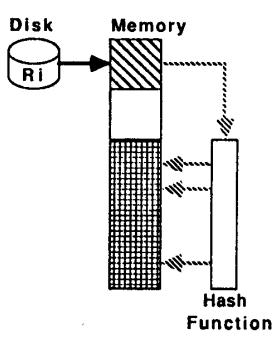


図 2. Build Phase

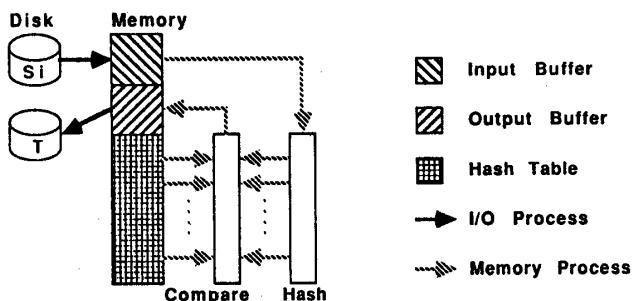


図 3. Probe Phase

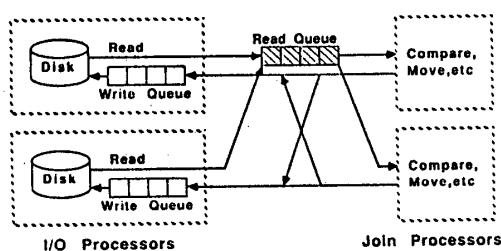


図 4.

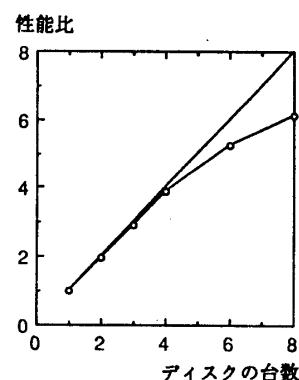


図 6. 測定結果 2

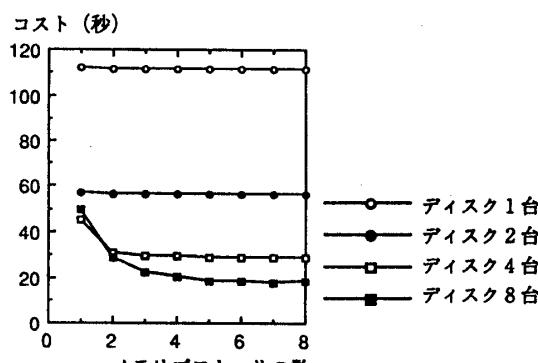


図 7. 測定結果 3