

不均一分布データに対する 動的デステージング方式の有効性について

7Q-8

中山雅哉, 喜連川優, 高木幹雄
東京大学生産技術研究所

1. はじめに

我々は、クラスタリング技法 (Hash) に基づく高速結合演算処理方式の提案を行ってきた [1, 2]。本方式は、従来から用いられてきた Nest Loop 方式や, Sort Merge 方式に比べて, 特に大容量のデータベースに対して有効であることが, 文献 [3] 等にも述べられている。本稿では, Hybrid Hash 方式等, 他のクラスタリング技法による結合処理方式にはみられない, 分割パケットの動的な処理順序決定機構 (動的デステージング方式) を取り入れた場合の処理性能について考察を行っている。静的に処理順序を決定する方式では, 分割パケットのデータ分布により処理性能が大きく左右されていたが, 本手法を用いることで, パケットのデータ分布によらず, ほぼ一定の性能を得ることが確認した。

2. 本稿における演算処理方式と記号表記について

本稿では, ユーザから発行される一般の問合せ式のうち, 制約処理を施した 2 つのリレーション R と S ($R < S$ とする) に対して結合処理を施し, 結果をディスクに書き戻す操作のみを対象として取り扱っている。複雑な問合せ式の中間ノードとして働く場合に相当する, 結合結果に対する制約演算処理はここでは扱わない。

クラスタリング技法に基づく結合演算処理方式では, 対象リレーション R (及び S) が演算処理に利用できる主記憶サイズ M より大きい場合, 各リレーションを一旦 H 個の部分リレーション (パケットと呼ぶ) に分割しておく (ステージングフェイズ), 次に各パケットを順次取り出して Hash 操作を用いて結合演算処理を実行していく方法を取る (プロセッシングフェイズ)。更にパケット自身が M より大きい場合は, 再帰的にパケットを分割していくことで処理を実行するようにしている。

試作システムでは, ディスクとの入出力負荷を抑える為に, 主記憶空間を入出力単位である 2048 バイトの固定長ページに分割, 管理している。全体 $(N + 2)$ ページのうち, N ページをリレーション R のステージング処理に用いており, リレーション S の読み込みと, 結果リレーションの書き出し様に 1 ページずつ割り当てている (図 1)。

また, 性能評価の際に用いたリレーション R, S は, 共にタプル長 64 バイトで, 結合属性は, キー重複のない 4 バイト整数フィールドを用いている。

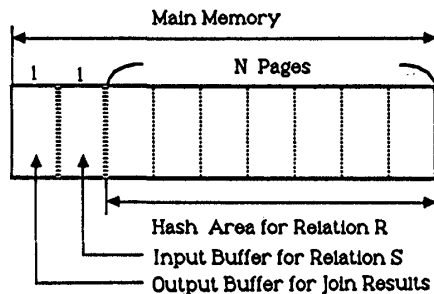


図 1 主記憶の利用構成

3. クラスタリング技法を用いた演算処理方式

Nest Loop 方式や, Sort Merge 方式との比較は, 他の文献 [2, 3] に譲り, 本稿では, 前章で概説したクラスタリング技法を用いた結合演算処理の方式について詳細に述べる。

3. 1. ステージングフェイズにおける処理

まず, ステージングフェイズにおける分割パケット数 H については, 以下の様な理由から最小になる様に選ばれることが多い。

- I. プロセッシングフェイズにおいては, 各パケットを処理する毎にハッシュテーブルの初期化が必要となる。
- II. H を大きくすると, 各パケットに配分されるタプル数が減少する為, 入出力バッファ内の無効領域が増大することが多い。
- III. 分割パケット数を増やすことで, ディスクへの書込みのための入出力バッファ数が増加し, 最初の処理に利用できるページ数が減少してしまう。

H を決定するにあたり, 従来から提案されてきた方法では, 各パケットのサイズがほぼ均一に N ページになると仮定するものが多い。例えば, Hybrid Hash 方式では,

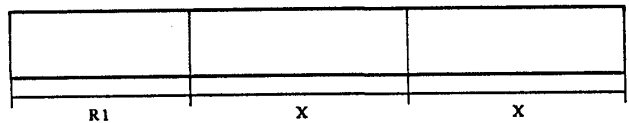
$$H - 1 = \left\lceil \frac{|R| - |M|}{|M| - 1} \right\rceil$$

$$|R_i| = |M| - H + 1$$

$$|R_i| = \frac{|R| - |R_1|}{H - 1} \quad (2 \leq i \leq H)$$

として H を決定している (図 2)。特にこの手法では, 最初に処理すべきパケット R_1 と, それ以外のパケットを区別し, R_1 の処理は, リレーション S のステージングフェイズとオーバーラップして実行することで, 入出力コストの削減を図っている。

Relation R



Available Memory

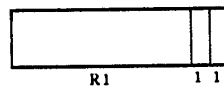


図 2 Hybrid Hash 方式におけるリレーションの分割方法

3. 2. プロセッシングフェイズにおける処理

プロセッシングフェイズでは, 各パケットに属するタプルを極力分割して, 結合演算処理のコストを削減する様にする。このため, ステージングフェイズでは, パケットをベ

ージ単位に連結した集りと考えるが、ここでは、タプル単位に連結した集りと考えて処理を行う。即ち、ステージングフェイズでは、 H 個の空間に分割して処理し、プロセッシングフェイズでは、 $32 * N (= 2048 * N / 64)$ 個の空間に分割して処理することになる。

Hybrid Hash 方式では、各バケット毎に結合演算処理が実行されている。

4. 動的デステージング方式による結合演算処理

動的デステージング方式 [2] は、従来とられてきた各手法と次の2点で大きく異なっている。

- ① 従来は、リレーションRのサイズをもとにして分割バケット数Hを最小とする様にしてきたが、本方式ではRによらず一定の値 ($N/2$) をとる様になっている。
- ② プロセッシングフェイズにおける各バケットの処理順序は演算処理を実行する前に決定されてきたが、これをデータの処理に伴い、動的にスケジューリングする方法をとっている。

前者は、前章で述べた様な欠点を持っているため、従来は取り入れられなかったものであるが、I. II. に関しては、ステージング処理の終了時に、各バケットのサイズをもとにスケジューリングを行い、小さなバケットはまとめてひとつの処理単位 (プロセッシングクラスタ) としてプロセッシングフェイズにおいて一度に処理する様にアルゴリズムの改良を行うことで解決している。III. に関しては、対処できる方法はないが、最初のプロセッシングクラスタの処理におけるメモリ利用効率のみが低下するだけなので、大容量のデータに対する処理全体の中で占める割合は低い。これに対し、より多くのバケットで分割することで、少々バケットサイズの変動に対しても、各々が主記憶容量を超えないで済むことに伴うオーバーヘッド軽減の効果は高いといえる。

後者は、従来の方式では考慮されていなかった各バケットのサイズが不均一となる場合に対応するための機構である。主記憶空間が足りなくなった時点で最もサイズの大きいバケットは、優先的にデステージングされて後で処理するバケットとして扱われる。この方式を用いることで、事前に予測したバケットの分布と実際のデータによるバケットの分布が異なる場合でも、最初に処理されるバケットや、その後の処理バケットのスケジューリングを行うことができるようになる。

5. 動的デステージング方式による演算処理の性能評価

2章でも述べた様に、動的デステージング方式をとる試作システムをワークステーション SUN3/260 上に作成し、不均一なバケット分布をとる2つのリレーションを結合した時の演算処理性能を測定した。本章はその結果についてまとめている。

不均一なバケット分布のリレーションとは、図3に示す様に、第kバケットのサイズ R_k が、第1バケットのサイズ R_1 のk倍になる (三角分布) データ分布をとるリレーションを指している。ここで、各バケットのサイズ R は、 $\sum R_k = R$ をもとにして計算される。

まず、Hybrid Hash 方式と、動的デステージング方式について、バケット分布が均一の場合と不均一の場合の性能を比較すると、図4に示す様になる。Hybrid Hash 方式では、バケット分布により大きく性能が変動するのに対し、動的デステージング方式では、不均一なバケット分布に対しても殆ど性能低下が見られないことが分る。これは、先に述べた様に、分割バケット数を大きくして、動的にバケットの処理順序をスケジューリングすることで、プロセッシングクラスタのサイズをほぼNにすることが可能となるためである。これは、図5の様に $H=500$ に固定した時、主記憶サイズNを変化させてもそれほど性能低下を招かないことから明らかである。

この他にも、バケット分布を図6の様にした場合も、ほぼ同様の結果を得ており、データ分布の変動に対して柔軟に対応できることが明らかとなった。

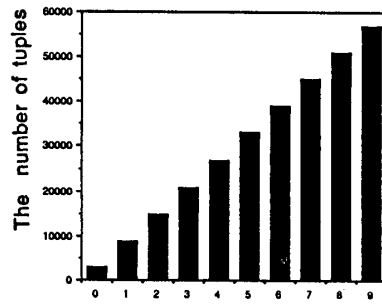


図3 不均一なデータ分布例1

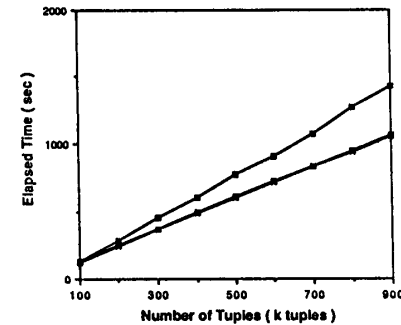


図4 各方式の性能評価結果

- Hybrid Hash 方式
- 動的デステージング方式
- 均一分散の時

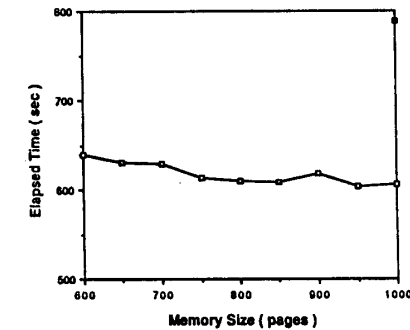


図5 主記憶サイズを変動させた時の性能

- Hybrid Hash 方式
- 動的デステージング方式

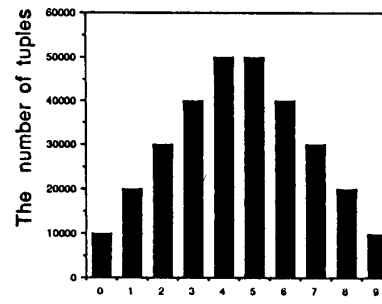


図6 不均一なデータ分布例2

6. まとめ

本稿では、動的デステージング方式を用いた結合演算処理方式について、従来から提案されてきたHybrid Hash 方式との比較を行い、提案する方式の柔軟性について述べてきた。また、ワークステーション上に実装した試作システムを用いて、分割バケットのデータ分布が不均一な場合の演算処理性能について、測定結果をまとめている。この結果、データの分布特性によらずほぼ同程度の性能が得られることが示されている。

参考文献

[1] 中山他, 「クラスタリング技法に基づく大規模リレーションの結合演算処理方式とその評価」, 36回情報大, 3F-5, 1988
 [2] 中山他, 「動的クラスタリング技法を用いた結合演算処理の性能評価」, 信学技法, DE87-21, 1988
 [3] D. DeWitt, et al., 「Multiprocessor Hash-Based Join Algorithms」, Proc. of Conf. on VLDB 1985