

2Q-7

パイプラインによる  
サイト間結合演算の並列処理

斎藤邦子、坂本明史、川上 英、疋田定幸

沖電気工業(株)

1. はじめに

質問処理において、結合演算はコストの高い処理を要する。分散データベースで、結合する表が異なるサイトにある場合、通信のオーバーヘッドも加わるため、効率化は重要である。分散では、複数の計算機による負荷の分担、並列処理の利点がある。従って、通信オーバーヘッドの削減と並列性のバランスを考慮した効率化が必要である。

並列処理の有効な方法に、パイプライン処理が提案されており<sup>[1]</sup>、結合演算処理では、パイプラインストール法が有効である<sup>[2]</sup>。パイプライン処理では、データがスムーズに流れ、処理を効率化するために、1回の通信で送るデータ件数(ブロッキングファクタ)に最適な値を設定することが重要である。本稿では、パイプラインストール法の処理において、ブロッキングファクタと並列性、通信オーバーヘッドの削減について考察し、最適なブロッキングファクタを示す。

なお、本研究は通商産業省工業技術院の大型プロジェクト「電子計算機相互運用データベースシステムの研究開発」の一環として研究開発を進めているものである。

2. 分散データベースモデル

分散データベースは、ローカルなDBMSを持つ複数のサイトから構成され、自サイト、他サイトのデータベースアクセスにはSQLの機能が提供されているとする。SQLは、データのフェッチを1タプル単位に行うものと定義している。しかし、サイト間の結合演算を行う場合、他サイトのデータのフェッチを1タプルずつ行っていたのでは、通信オーバーヘッドが大きい。通信回線の高速化により、転送量が通信時間に及ぼす影響は少ないので、他サイトのデータをnタプル(ブロッキングファクタ)ずつフェッチすれば、オーバーヘッドの減少が期待できる。以下では、パイプライン処理のためのプリフェッチの機能を示し、2サイト間、3サイト間の結合演算処理を説明する。

2.1 プリフェッチ

サイト1のデータをサイト2が処理する場合、パイプライン処理の目的は、サイト1から2へデータが流れ、サイト1、2の処理が止まらずに並列に行われることである。そのために、非同期通信プリミティブによるプリフェッチの機能が必要である。図1は、プリフェッチによる並列処理を示している。サイト2では、サイト1よりタプルを受け取ったならば、次のタプルの要求をだし(図1①)、応答を待たずに、受け取ったタプルの処理を行う。処理が終了し、次のタプルが

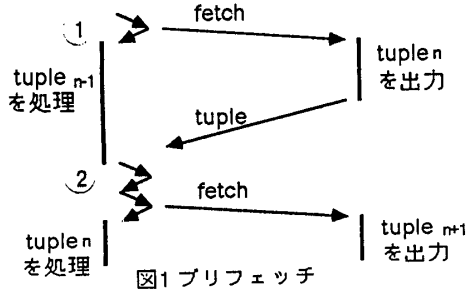


図1 プリフェッチ

必要になったならば、応答メッセージを受け取る(図1②)。

2.2 結合演算処理

まず、サイト1、2のリレーション  $R_1, R_2$  の結合の処理を説明する。 $B_n$ はサイトnから転送されるタプルのブロッキングファクタである。サイト1はサイト2の要求により、 $R_1$ から条件なしで  $B_1$ 個のタプルをフェッチし転送する。サイト2では、まず、 $R_2$ の検索のために、 $R_1$ の最初のタプルと結合するタプルを得るための条件を付ける。例えば、結合条件が  $R_1.A=R_2.A$  の場合、 $R_2$ の検索には、“ $R_2.A=パラメータ$ ”の条件が付けられており、パラメータに  $R_1$ のタプルのAの値をセットする。 $R_2$ のタプルをフェッチし、結合結果を得る。フェッチの結果、“該当タプルなし(End Of Tuple)”になったならば、 $R_1$ の次のタプルに対して、同様の処理を行う。サイト2は結合処理を行う前に、 $R_1$ の次の  $B_1$ 個のタプルを得るための要求をだしておく。次に、サイト1、2、3のリレーション  $R_1, R_2, R_3$  の結合の処理を説明する。サイト2は、サイト3からの要求により上記の方法で  $R_1, R_2$ の結合を行い、 $B_2$ 個のタプルを得たならば、転送する。サイト3は、 $R_1, R_2$ の結合結果と $R_3$ の結合を同様に行う。この場合も、次の  $R_1, R_2$ の結合結果の  $B_2$ 個のタプルのプリフェッチの要求を出しておく。

3. ブロッキングファクタと処理効率

3.1 最適なブロッキングファクタ

効率的にパイプライン処理を行うには、サイト2がサイト1のデータを待つ時間( $W_2$ )を短くしなければならない。図2はサイト1、2の処理を示したものであり、 $W_2$  は次の2つを合わせたものである。

- ① プリフェッチによる待ち時間
- ② サイト1から最初のデータが届くまでの待ち時間

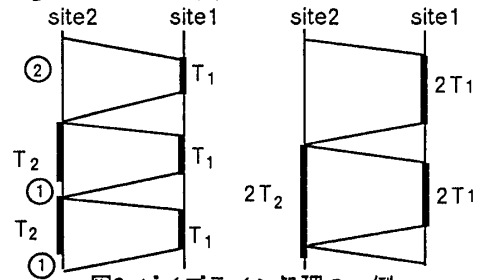


図2 パイプライン処理の一例

①は、サイト2の処理が終わり、次のタプルが必要となったときに、サイト1からタプルが届いていないために生じる待ち時間である。この待ち時間は、通信コストが転送量に依存しないとすると、ブロッキングファクタに最適な値を設定することで、削減できる。図2では、サイト1からサイト2に送るタプル数を2倍にしても、通信時間は変わらないので、1タプルあたりの通信オーバーヘッドは1/2になり、処理時間の差に吸収される。従って、次の式を満たす場合、①は生じない。

$$T_{i2}-T_{o1} \geq T/B_1 \dots (1)$$

$T_{o1}$ : サイト1が、1タプル出力するための時間

$T_{i2}$ : サイト2が、受信したデータ、1タプルを処理する時間

$B_1$ : サイト1が1度に転送するタプル数

$T$ : 要求とデータの転送のための通信時間(転送時間)

(1)を満たしていない場合、一回の通信で生じる待ち時間は、  
 $(T + B_1 * T_{01}) - B_1 * T_{i2} = T + B_1 * (T_{01} - T_{i2})$   
 であり、転送回数は、 $N_1/B_1$

$N_1$ : サイト1が転送する全タプル数  
 であるので①は、

$$(T + B_1 * (T_{01} - T_{i2})) * (N_1/B_1)$$

になる。

☒は、並列動作に入るまでの待ち時間であり、

$$T + B_1 * T_{01}$$

になる。

従って、ブロッキングファクタ( $B_1$ )と全体の待ち時間( $W_2$ )の関係は、

i) 式(1)を満たしている場合(①の待ち時間がない)

$$W_2 = T + B_1 * T_{01}$$

ii) 式(1)を満たしていない場合

$$W_2 = T + B_1 * T_{01} + (T + (B_1 * (T_{01} - T_{i2}))) * (N_1/B_1)$$

である。ブロッキングファクタが大きくなると、通信のオーバーヘッドは削減されるが、並列動作に入るまでの待ち時間が増える。上記の式より、 $W_2$ を最小にするブロッキングファクタ $Bo_1$ は、以下のようになる。

i)  $T_{i2} > T_{01}$  かつ  $T/(T_{i2}-T_{01}) < \sqrt{N_1} * T / T_{01}$  の場合

$$Bo_1 = T/(T_{i2}-T_{01})$$

ii) i) 以外の場合

$$Bo_1 = \sqrt{N_1} * T / T_{01}$$

i)では、通信オーバーヘッドは、最初のデータの転送のときだけである。ii)は、これ以上大きな値にすると、通信のオーバーヘッドは下がるが、最初の待ち時間が増え、並列性が下がる場合である。

### 3.2 ネストループ法でのブロッキングファクタ

パイプラインネストループ法を用いた2サイト、3サイト以上の結合での、最適なブロッキングファクタについて考察する。そのためには、データを送るサイトと、受け取って処理するサイトの処理時間を見積らなければならない。各サイトの処理時間は、(フェッチの回数)\*(1タプルのフェッチに要する時間)とみなすことができる。リレーションのタプル数、セレクトィビティ等を用いて、処理時間を見積もる。セレクトィビティは、結合結果のサイズを見積もるためのものであり、 $R_1, R_2$ のセレクトィビティを  $S_{12}$  とすると結果タプル数は  $N_1 * N_2 * S_{12}$  になる。

#### ① 2サイト結合

$R_1, R_2$ の結合で、サイト2は、 $R_1$ の各タプルについてBOTまで検索するので、フェッチの回数は、(結果タプル数 $N_1$ )となる。これを、 $N_1$ で割ったものが、 $R_1$  1タプルを処理するのに必要なサイト2のフェッチの平均回数であり、 $N_2 * S_{12} + 1$ となる。従って、 $T_{01}, T_{i2}$ は以下のような見積りになる。

$$T_{01}: C_s$$

$$T_{i2}: C * (N_2 * S_{12} + 1)$$

$C_s$ : 条件なしのシーケンシャルなフェッチに要する時間

$C$ : 条件付のフェッチに要する時間

#### ② 3サイト以上の結合

$R_1, R_2, R_3$ を結合する場合、サイト2からサイト3へ送るタプルの最適ブロッキングファクタを求めるために、サイト2、3の処理時間の見積りについて述べる。サイト2で、 $R_1, R_2$ を結合するためのフェッチの回数は①の通りである。これを結果タプル数で割ったものが、 $R_1, R_2$ の結合結果1タプルを得るために必要なフェッチの平均回数であり、 $1 + 1/(N_2 * S_{12})$ となる。サイト3が  $R_1, R_2$ の結合結果1タプルを処理するために必要な時間は、①の  $T_{i2}$ と同様に見積もることができる。従って、 $T_{02}, T_{i3}$ は以下のようになる。 $n$ サイトの結合の場合も同様に見積もることができる。

$$T_{02}: C * (1 + 1/(N_2 * S_{12}))$$

$$T_{i3}: C * (N_3 * S_{23} + 1)$$

以上の見積りで得た最適ブロッキングファクタを、転送時間、タプル数等に適当な値を設定した例で示す。

例) データベースアクセス時間を  $C_s=20msec, C=40msec$ 、リレーションのタプル数を  $N_1=10000, N_2=12000, N_3=11000$ 、セレクトィビティを  $S_{12}=1/12000, S_{23}=1/10000$  とする。転送時間( $T$ )が、80msec, 160msec, 800msec の場合について示す。800msec は、衛星通信のように通信時間がかかる場合を想定している。

$R_1, R_2$ の結合では、 $T_{01}=20msec, T_{i2}=80msec$  となり、

$T=80msec$  の場合  $Bo_1=2$ 、

$T=160msec$  の場合  $Bo_1=3$

$T=800msec$  の場合  $Bo_1=14$  である。

$R_1, R_2, R_3$ の結合では、 $T_{02}=80msec, T_{i3}=84msec$  となり、

$T=80msec$  の場合  $Bo_2=20$ 、

$T=160msec$  の場合  $Bo_2=40$ 、

$T=800msec$  の場合  $Bo_2=200$  である。

以下のグラフは、サイト2から転送するタプルのブロッキングファクタとサイト3の待ち時間の関係を示したものである。 $B_2 \leq Bo_2$ の場合、 $B_2$ の増加による通信のオーバーヘッドの減少が大きく、 $W_2$ も減少する。 $B_2 > Bo_2$ では、通信のオーバーヘッドはなくなるが、 $B_2$ の増加により最初の待ち時間が増えるため、 $W_2$ は増加する。衛星通信のように、通信オーバーヘッドが大きい場合、ブロッキングは有効である。

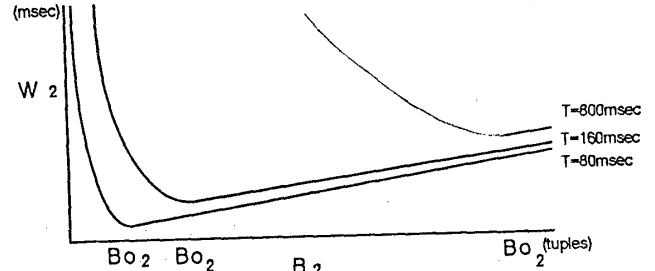


図3 ブロッキングファクタと待ち時間

### 4. おわりに

本稿では、パイプラインネストループ法におけるブロッキングファクタと処理効率について述べた。並列性を下げる要因に(1)並列動作に入るまでの待ち時間、(2)通信オーバーヘッドがある。これらは、ブロッキングファクタに左右されることがわかった。(1)、(2)は相反するものであり、ブロッキングファクタを増やすことで、(2)は削減できるが、(1)が大きくなる。さらに、考察を深め、最適なブロッキングファクタを得ることができた。

質問処理では、処理方法、処理プラン(結合の順番等)が重要である。従来、通信のオーバーヘッドを考慮し、転送量の削減を目的としたプランの決定についての研究が多かった。今後、並列性を考慮したプランの決定も必要である。また、以下の2点を考慮した処理時間の見積りについても、検討していきたい。

①値のばらつきによる、結合のためのフェッチ回数のばらつき

②複数利用者による、1サイトへの負荷の集中によるフェッチ時間の変化

[1] Kiyoki, Y., Kato, K., Masuda, T., "A Stream-Oriented Approach to Distributed Query Processing in a Local Area Network", Proc. 1986 ACM SIGSMALL/PC Symp. Small Syst., pp. 146-155, May, 1986.

[2] Hikita, S., Kawakami, S., Haniuda, H., "A Stepwise Approach to Distributed Database Systems by Database Machines", Proc. 1985 ACM SIGSMALL Symp. Small Syst., pp. 18-24, Dec. 1986.

[3] Selinger, P. G., Astrahan, M. M., Chamberlin, D. D., Lorie, R. A., Price, T. G., "Access Path Selection in a Relational Database Management System", Proc. ACM SIGMOD 1979 Int. Conf. Management of Data, pp. 23-34, 1979