

Kappa-P の並列問い合わせ処理

永沼和智¹, 佐藤裕幸¹, 河村元夫²

5 R - 4

1: 三菱電機 (株) 情報電子研究所 2: (財) 新世代コンピュータ技術開発機構

1 はじめに

第 5 世代コンピュータプロジェクトでは、大規模知識処理を目的として、並列推論マシン PIM [1] 及びその上で動作する並列知識処理システムの研究開発を行っている。Kappa-P [2] は、これらの知識処理システムに大量の複雑なデータの効率的な処理を提供するために開発された、並列非正規関係データベース管理システムである。Kappa-P は、PIM の各クラスに配置されたローカル DBMS と呼ばれる複数のデータベース管理システムから構成されており、クラス間での並列性を得るために、一つのテーブルを水平に分割し複数のローカル DBMS に配置することができる。そのテーブルへの問い合わせは、インタフェースプロセスにより構成テーブルへの並列な問い合わせに置き換えられ、各ローカル DBMS により並列に処理される。本論文では、この水平分割テーブル処理、また分割テーブルの欠点である通信コストの増大を防ぐための方法を、実測値と共に報告する。

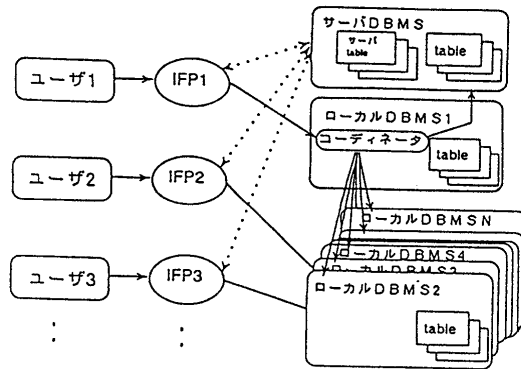


図 1: 問い合わせ処理

2 問い合わせ処理概要

2.1 処理の流れ

Kappa-P では、原始コマンドと呼ばれるプログラムインタフェースを提供している。原始コマンドとは、一つのテーブルに対して読みや更新の処理を行うコマンドで、インタフェースプロセスがユーザからコマンドを受ける。インタフェースプロセスは、コマンドを受けるとサーバと通信し、指定されたテーブルがどのローカル DBMS に存在するかと言う情報を得る。サーバは、データベースのグローバル情報を管理しているプロセスである。そして、インタフェースプロセスはそのローカル DBMS 情報をコマンドに付加し、コーディネータにその情報を送る。その情報によりコーディネータは、各ローカル DBMS へコマンドを送り、各ローカル DBMS でコマンドが処理される。(図 1 参照)

2.2 水平分割テーブル

水平分割テーブルとは、1つのテーブルをその属性値などにより複数の構成テーブルにレコード単位で分割し、複数のローカル DBMS に配置したものである。この分割方法には、属性値の範囲による分割、属性値のハッシュ値による分割、ラウンドロビンの 3 種類がある。属性の範囲による分割はその属性を条件に指定された時に有効であるが、構成テーブルのレコード数がまちまちになってしまうと言う欠点があり、これに対しラウンドロビン

は構成テーブルの大きさが均一になるので、各テーブルに対する処理時間が均一になりやすい。これらは、そのデータに対する検索を考慮して選択する必要がある。

水平分割テーブルに対する検索は、構成テーブルの存在する各ローカル DBMS で並列に処理されるが、ユーザは分割テーブルであることや構成テーブルの所在などを意識せずに一つのテーブルとして扱うことができる。ユーザから水平分割テーブルに対するコマンドが出されると、インタフェースプロセスはサーバとの通信により、そのテーブルの各構成テーブルの名前と存在するローカル DBMS を知る。そして、その情報によりコマンドは各ローカル DBMS に対する複数のコマンドに変換され、コーディネータ経由で各ローカル DBMS で並列に実行される。このときレコードの追加や削除のコマンドは、追加、削除すべき構成テーブルを分割情報から判断しその構成テーブルへのコマンドに変換される。

2.3 フィルタ付きレコード読み

テーブルを分割した場合にローカル DBMS 間の並列性は得られるが、その分ローカル DBMS 間の通信コストが増大すると言う欠点がある。特にユーザがデータを一度読み込み、そのデータに対して PIM の各クラスを使って並列に処理を行いたい場合には、一箇所にデータが集められてから再び分散させ、そのデータに対する処理を行った結果をもう一度集めるために、かなりの通信コストがかかる。この通信コストの増大を防ぐために、Kappa-P ではフィルタ付きレコード読みの機能を提供している。このフィルタ付きレコード読みでは、ユーザがその検索の結果に行う自分で定義した処理(フィルタ)を指定することにより、その処理が分割された構成テーブルのあるクラスで実行され、その処理の結果がユーザのいるクラスに集められる。このため、一度データを

Parallel query processing on Kappa-P
Kazutomo NAGANUMA¹, Hiroyuki SATO¹, Moto KAWAMURA²
1: Mitsubishi Electric Corporation 2: Institute for New Generation Computer Technology.

集めてから分散させると言うことがなくなりまた、一回の通信データもフィルタを通った後のものになるので、通信コストはかなり軽減される。

3 評価

実際に一つのテーブルを分割した場合に、どれだけの速度向上が見られるのか、64台のプロセッサを持つPIM/mを使って測定した。

3.1 測定方法

測定の対象としたデータは、公共のデータベースPIRの遺伝子データである。このテーブルのデータは約61メガあり、3万3千のレコードを持っている。このテーブルを16,32,64に分割し、それぞれ次の測定を行った。なお、ディスクの数が少ないため必要なデータはあらかじめ主記憶に読み込んでから測定を行ったが、分割しないテーブルは主記憶の容量が足りないために実行できなかった。

- (1) 属性sequenceがCHという文字列を含むレコードの検索。
この検索では、約4119件のレコードが結果として得られる。
- (2) フィルタ付きレコード読み。
指定したフィルタは、属性sequenceがある一定のパターンの文字列を含むレコードを検索するモチーフサーチである。この検索では、47件のレコードが結果として得られる。
- (3) フィルタ付きレコード読みを使わずに、一度レコードを読み込んでからその結果に対して(2)と同じフィルタを並列に実行する。

3.2 測定結果 / 考察

測定結果を表1に示す。

(1)の測定結果では、32分割の実行時間が16分割のテーブルの2分の1に近い数値となっている。これは、テーブルの大きさが2分の1となりそれがそのまま実行時間に反映している。実際は通信時間、前処理、後処理の時間は2分の1とはならないので、厳密には2分の1よりも多少大きな数値となっている。

しかし64分割の場合には、32分割の2倍とはならず約1.6倍となっている。これは、一つの分割テーブルに対する処理が非常に軽くなったため、通信にかかる時間が無視できなくなったと言うこと、各ローカルDBMSで行われる検索の前のジョブやトランザクションの生成、テーブルのロック、検索の後の終了処理などにかかる時間は分割テーブルであろうと普通のテーブルであろうと一定であるため、実際の検索時間が短くなるとこれらの時間の占める割合も無視できなくなったと考えられる。

(2)の測定では、32分割テーブルの実行時間は16分割の約2分の1となっており、64分割の実行時間は約4分の1となっている。これは、(1)のレコード検索と比べてここで指定したモチーフサーチの負荷がかなり高く実行時間もかかっているため、先に述べたような通信、前処理、後処理にかかる時間が無視できる範囲であったと考えられ、分割の効果が十分表れていると言える。

(3)の測定結果を(2)と比べてみるとそれぞれ約1.5倍

テーブル分割数	16	32	64
(1) レコード検索	33,331	18,158	11,575
(2) フィルタ付きレコード検索	132,447	68,892	37,664
(3) レコード検索 + フィルタ	206,549	101,173	66,190

表1: 各テーブルの実行時間 (msec)

から、2倍近い時間がかかっている。フィルタ付きレコード読みを使わずに同じ処理をする場合には、フィルタの処理に入る前のレコード読みの結果が一度ユーザのいるクラスタに集められてから、再び分散することになる。この通信コストがそのまま数字に表れていると考えられ、フィルタ付きレコード読みの有効性が確認された。

この測定結果から、(1)の前方一致、後方一致と言うような検索に対しては、このテーブルを64以上に分割してもそれほど効果が得られていない。これに対し(2)のモチーフ検索では、64分割でも十分な効果が得られており、更に分割した場合の効果も期待できる。このように、検索の負荷の大きさにより適した構成テーブルの大きさは異なり、そのデータの使用目的のより分割数を決定することが必要である。

4 おわりに

今回測定に使用した並列推論マシンPIMは64台のプロセッサを持っているが、現在ICOTでは256台版のPIMが動作している。また、Kappa-Pは大量で複雑なデータを扱うことを目的としており、今後は今回扱った以上に大量で複雑なデータに対して負荷の大きい処理が要求されると考えられる。そのような要求に対して水平分割テーブルは非常に有効な手段であることが今回の測定により確認できた。ただ、データの複雑さ、データに対する処理の負荷の大きさにより、最適な分割テーブル数が異なるため、今後は分割をする際の目安になるような解析をする必要があると思われる。

参考文献

- [1] 瀧: *Parallel Inference Machine PIM*, The International Conference on Fifth Generation Computer Systems'92, 1992.
- [2] 河村他: 並列データベース管理システムKappa-Pの概要, 第45回情報処全国大会, 5R-3, 1992, 10.
- [3] 横田他: *Overview of the Knowledge Base Management System (KAPPA)*, The International Conference on Fifth Generation Computer Systems'88, 1988.