

4 A C - 7

並列論理型言語 KLIC による 並列データベースの I/O 処理の高速化

宮崎 純

横田 治夫

北陸先端科学技術大学院大学 情報科学研究科

1 はじめに

我々は、並列論理型言語 KLIC[1] によるアクティブデータベースシステム [2] のプロトタイプ Parade[3] を提案している。KLIC の論理変数によるプロセス間のメッセージ通信は、無共有並列計算機上でデータベース処理を実現する際に非常に適合性が良い。これまでの研究により、データベース処理を KLIC のプロセス間のメッセージ通信を利用することにより実現し、汎用ワークステーションならびに超並列計算機 nCUBE2 上で動作させ、その記述能力と移植性の高さを実証してきた。しかしながら、データベース処理において非常に大きな鍵を握る二次記憶および PE 間のデータ転送などの I/O 処理が、システムの性能のボトルネックとなっていることが分かってきた。これは、KLIC 標準の I/O が性能を重視して設計されたものではないことに起因する。

本稿では、KLIC による並列データベース処理を実用的な速度で動作させるために、KLIC のジェネリックオブジェクトを用いて低レベルの I/O 性能の向上を目指し、その初期評価を並列計算機 nCUBE2 を用いて行なった。

2 I/O 処理の高速化

2.1 ページオブジェクト

KLIC 標準の I/O は、外部表現と内部表現を交換するために parse/unparse 機構を必要とし、かつそのデータアクセス単位は小さい。そのため、ソフトウェアオーバーヘッドの大きいものとなっており、しかもディスク I/O の面から見ても、スループットを高くすることは困難である。また、PE 間のデータ転送は基本的に 1 レベルごとの小さなメッセージで行なわれるため、メッセージハンドリングオーバーヘッドの大

きな並列計算機では、そのコストは重い(図 1 参照)。

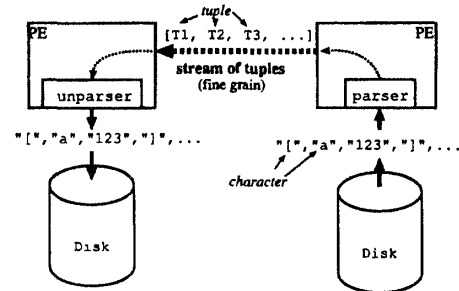


図 1: KLIC 標準 I/O

一般的にデータベース処理では、I/O が全体のボトルネックとなることが多く、このボトルネックを解決することがデータベース全体の処理速度の向上につながる。通常のデータベースでは、ディスクをページ単位といった大きなデータ単位でアクセスしてディスク I/O のスループットを高くすることが行なわれている。これを KLIC で行なうためには、KLIC のデータ型を拡張するためのジェネリックオブジェクトを用いるのが、最も簡単である。そこで、新しいデータ型としてページ単位でディスクをアクセスできるよう、ジェネリックオブジェクトを用いてページオブジェクトを定義し、ディスク I/O のスループットの向上を図った。ページ本体は、KLIC のヒープとは別空間に取り、ゴミ集めによる無駄なデータの移動を減らすことも試みている。このページオブジェクトによるページ単位のディスクアクセスは、B-Tree などのアクセスパス構造との親和性の向上にも貢献する。また、このページオブジェクト単位で PE 間のデータ転送を行なうことにより、通信のスループットの向上も同時に達成することができる。一般に無共有並列計算機は、メッセージのセットアップコストがメッセージバッファの獲得やロックのため、通信のスループットに比較してかなり大きい。しかし、KLIC はプロセッサ間通信の際に、通常 1 レベルずつの細かなメッセージが転送される。KLIC はマルチレベルでのメッセージ通信もサポートしているが、これはメッセージの encode/decode というソフトウェアの

オーバーヘッドが伴う。しかしながら、前述のページオブジェクトを用いれば、連続空間上の大きなデータ単位でメッセージを送受信することが可能となり、メッセージのセットアップコストを相対的に小さくでき、さらにデータの encode/decode といったソフトウェアのオーバーヘッドも削減できる。以上によりトータルな通信オーバーヘッドを大きく改善することができる(図2参照)。

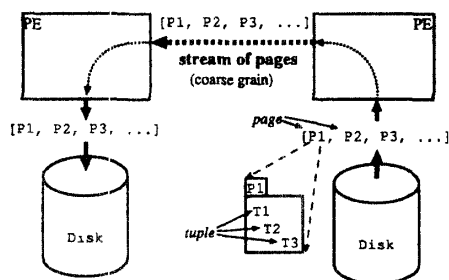


図2: ページオブジェクト I/O

2.2 ページオブジェクトの拡張

このページオブジェクトに対して様々なメソッドを定義し、リストなどの KLIC の基本データ構造に変換することなく直接データを扱えるようにすることができる。例えば、選択、射影、結合演算などを直接メソッドとして定義し、ページ単位で効率よく処理することが可能である。

3 ページオブジェクトの I/O 性能

4KB 単位でディスクアクセスとメッセージ通信ができるよう、ジェネリックオブジェクトを用いてページオブジェクトを定義し、ページオブジェクトと KLIC 標準の I/O の性能の比較を nCUBE2 上で行なった。実験に用いた nCUBE2 は、表1で示す I/O 性能を持つ。

表1: nCUBE2 の I/O 性能

Disk Read	1.1 MB/s
Disk Write	0.46 MB/s
PE 間通信	2.2 MB/s

表2は1タプル当たり180Bである3000タプルのリレーションにおける、ディスクからの読出し、書込み、PE間通信のI/O性能の比較を行なったものである。表中の termio は KLIC 標準の I/O, codedio は parser/unparser を専用の符合化を行ない最適化を行なったもの、pobjio は今回定義したページオブジェクトを用いた I/O の場合である。

この結果からも分かるように、ページオブジェクト

表2: I/O 性能の比較 (単位: msec)

	termio	codedio	pobjio
Disk Read	2.56×10^4	934	907
Disk Write	2.35×10^5	2130	2030
PE 間通信	4170	4180	165

を用いた I/O は、KLIC 標準の I/O に比較して著しい I/O 性能の向上が得られることが明らかである。これは、ページオブジェクトが、KLIC 内部表現と外部表現の変換を行なう parser/unparser のオーバーヘッドがなく、かつ大きなサイズのデータ単位で I/O を行なうため、I/O スループットが向上したからである。専用の符合化を行なう codedio の parser/unparser は、このオーバーヘッドが小さく、ページオブジェクトに比較してディスク I/O がわずかながら劣る程度であるが、PE 間通信は1レベルずつであるため、通信性能の改善は得られない。

4 おわりに

並列論理型言語 KLIC を用いてデータベース処理を行なう場合、ディスクや通信の I/O がボトルネックとなる。これは、KLIC の I/O の粒度が小さく、また、I/O の際のソフトウェアのオーバーヘッドが大きいためである。本稿では、KLIC のジェネリックオブジェクトを用いて新しくページオブジェクト型を定義し、これを用いて I/O をページ単位で効率よく行なう方法について述べた。さらに nCUBE2 上での実験を通して、このページオブジェクトを用いた I/O が、従来の I/O の性能に比較して大きな性能の改善が得られることを示した。

今後、このページオブジェクトに様々なメソッドを定義することにより、ページ単位で効率の良いデータベース処理を行なうことができるように拡張を行なう予定である。

参考文献

- [1] T. Fujise, T. Chikayama, K. Rokusawa and A. Nakase: "KLIC: A Portable Implementation of KL1", Proc. of FGCS '94, pp. 66-79 (1994).
- [2] J. Widom and S. Ceri: "Active Database Systems: Triggers and Rules For Advanced Database Processing", Morgan Kaufmann Publishers, INC., San Francisco, CA (1996).
- [3] 横田, 官崎, 土屋: "並列アクティブデータベースエンジン Parade の並列処理", 文部省科学研究費重点領域研究『高度データベース』松江ワークショップ講演論文集, 第2巻, pp. 426-433 (1996).