

データベースプロセッサ RINDA の データベースアクセス方式

5Q-9

板倉 一郎、 中村仁之輔、 井上 潮

NTT情報通信処理研究所

1. はじめに

リレーショナルデータベース処理の性能向上を図るため、サーチ・ソートを高速に処理するデータベースプロセッサ RINDA (1) を開発した。本稿では、RINDA を制御する DBMS のデータベースアクセス方式について述べる。

2. 検索処理方式

RINDA は CSP (内容検索プロセッサ) と ROP (関係演算プロセッサ) から構成される。

CSP はディスクに格納された表を指定された条件でサーチし、条件に合致する行に対し射影を行い CPU に転送する。ROP は CPU から転送された表の行のソート、ふるい落とし、重複除去等を行い、結果を CPU に転送する。

CSP の出力と ROP の入出力が複数行単位であるため、検索処理は一時表 (検索の途中結果を格納する一時的な表でディスク/バッファ上に存在し、検索終了後削除される) を単位として、以下のような手順で実行される (図 1)。

- ① CSP で表を検索し結果を T1 (一時表を意味し、以下の説明でも同様) に格納する。
- ② T1 に対しソフトウェアで処理する条件式があればその判別を CPU で行い、結果を T2 に格納する。
- ③ ソート、結合、グループ化が指定されている場合は、T2 を ROP でソートし結果を T3 に格納する。
- ④ 結合、副問合せが指定されている場合は、①~③ をそれぞれの表に対して行い結果を一時表に格納し、CPU で結合、副問合せを行う。

なお、表をサーチしその結果をソートするという単純な問合せについては、ディスク/バッファ間の I/O 時間を

削減するため、以下の手順で処理する CSP/ROP 連動方式を実現している。

- ① CSP のサーチ結果をバッファに転送する。
- ② バッファのデータを ROP に転送する。
- ③ CSP のサーチを再開し、①、② を繰り返す。
- ④ ROP からのソート結果を一時表に格納する。

問合せが CSP/ROP 連動で処理可能であるかは言語処理部で検索条件式の解析により判断する。

3. バッファ管理方式

3.1 要求条件

CSP・ROP を有効に利用し、問合せを処理するため以下の処理が効率良く実行されなければならない。

- ① CSP のサーチ・ROP のソートのためのデータ転送
- ② 一時表への連続ページアクセス

一般的に使用されているページ単位のバッファ管理方式では、次の問題があり効率良く検索を実行することができない。

- ① CSP・ROP の起動がページ単位となり、CSP・ROP 起動回数が多くなる。
- ② 処理の単位が一時表であるため、連続的にページをアクセスすればよいにもかかわらず、ページ単位の入出力ではディスクのシーク・サーチ動作が多くなる。

3.2 RINDA 用バッファ管理方式

(1) 実現機能

上記の問題点を解決するため、RINDA 用のバッファ管理として以下の機能を実現した。

- ① CSP/ROP 起動回数を減らすため、バッファを複数ページ単位 (トラック単位) で一時表に割当てする。
- ② 一時表のページは基本的にはシーケンシャルにアクセスされることを考慮し、シーク・サーチ動作回数を削減するため、バッファ/ディスク間の入出力をトラック単位で行う。

また、バッファを有効に利用するため、次の機能も実現した。

- ③ 一時表にアクセスする時点でバッファを確保し割当てる。ただし、結合処理時には 3 個の一時表 (入力 2 個、出力 1 個) が必要となることを考慮し、最大 3 個までバッファを動的に確保する。

(2) 動作概要

(1) の機能を実現したバッファ管理方式の動作概要を以下に示す。

- (A) 一時表アクセス時のバッファの割り当て
3 個の一時表まで動的にバッファを確保し割当てる。

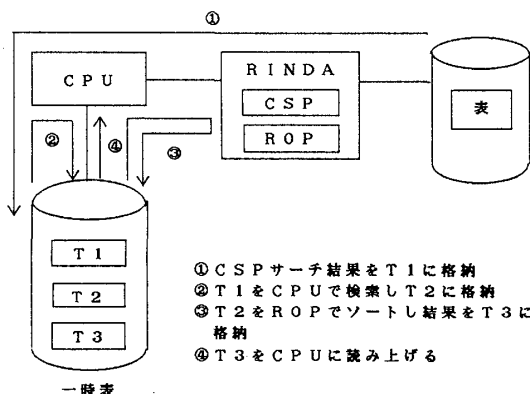


図1 検索処理の流れ

RINDA - Relational Database Processor: Database Access Method

Ichirou ITAKURA, Jinnosuke NAKAMURA, Ushio INOUE

NTT Communications and Information Processing Laboratories

3個使用中であればFIFOによりバッファを割当てる。

(B) ページアクセス時のバッファの割当て

バッファ上にアクセスするページが存在するか調べる。存在しないならば、アクセスするページを含むトラック全体をバッファに読み込む。ただし、既にバッファ上に追い出す必要のあるページが存在する場合には、トラック単位でディスクに書き込む。

(C) 動作例

一時表T1のP1~P15を5ページ単位でROPに転送し、P6'~P10'をバッファに転送する直前のバッファの状態を示す(図2)。

4 ディスク並列アクセス方式

RINDAを用いた検索処理では表全体のサーチを行うため、表のデータ量に比例した処理時間がかかる。表のサ

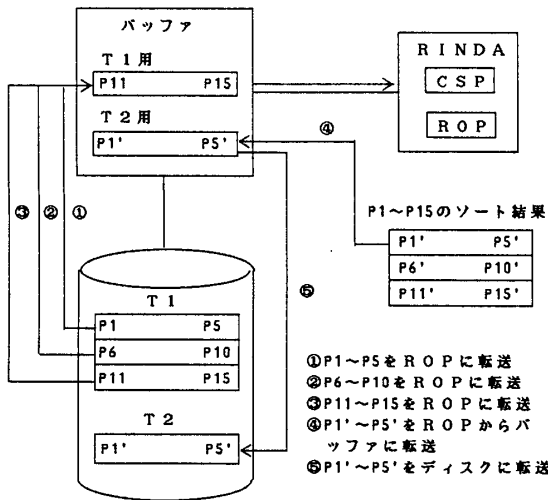


図2 バッファ管理の動作例

ーチ時間を短縮するため、表を複数のボリュームに分散格納し、複数台のCSPで並列にサーチするディスク並列アクセスを実現した(図3)。

4.1 分散格納

表の複数ディスクへの分散格納実現に当たり以下の事項を考慮した。

①各ディスクに均等にデータを格納する。

②格納ディスク決定のためのオーバーヘッドが少ない。

特定の列の値に基づきハッシング、クラスタリング等により格納ページ・格納ディスクを決定する方法は、列の値の分布による偏り、列の値の更新に伴う行の移動が発生する可能性があり、必ずしも良い方法とは言えない。このため、行は到着順に格納し、ページ単位で周期的にディスクに割当てる方法を実現した。

4.2 並列検索

ディスク並列アクセス時には以下のようにバッファを使用し、並列に検索を実行する。

①一時表に割当てられたバッファを並列にサーチを行うCSPに均等に割当てる。

②複数台のCSPを起動し、並列に検索を行う。

③CSPの検索動作は、指定された検索範囲のサーチが終了、あるいは、検索範囲のサーチは終了していないが検索結果を格納するバッファが満杯になったときに終了する。前者で終了し未使用バッファが存在する場合は、バッファ不足で検索が中断しているCSPにバッファを再度割当て、検索を再開する。

5. おわりに

本稿では、RINDAのデータベースアクセス方式として検索処理、バッファ管理、ディスク並列アクセスについて述べた。一時表を単位とした処理方式に着目したバッファ管理方式により、一時表を格納するディスクとのI/O時間短縮を図った。

6. 参考文献

- [1] 井上他「データベースプロセッサRINDAのアーキテクチャ」、情処第37回全国大会

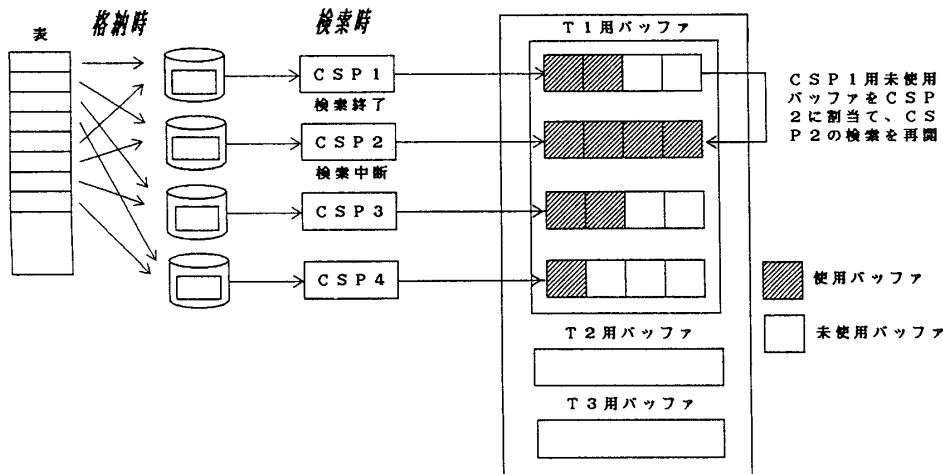


図3 ディスク並列アクセス方式の概念