

5N-9

PIE64 のデータベースマシンとのインタフェースについての考察

吉田 実, 田中英彦  
 東京大学工学部\*

1 はじめに

我々の研究室では、並列推論マシン PIE64 の開発をすすめている。[1] PIE64 は、データベースマシンをつなぐことにより、大規模なデータを扱うアプリケーションを走らすことができるようになる。ここでは、PIE64 とデータベースマシンをつなぐ場合の問題点について考察する。

2 データベース呼び出し

PIE64 では、コミットドチョイス型言語 FLENG[2] を走るため、FLENG からデータベースにアクセスできるようにする必要がある。通常、データベースマシンではデータベースの整合性を保つトランザクション単位でアクセスするのが常套手段である。よって、次のようにトランザクションを生成する述語を設ける。トランザクションの中間状態は、ユーザプログラムからは見えない。

```
transaction(Result, Transaction)
```

述語 transaction は、Transaction で与えられるコマンドを処理し、その処理結果の状態を Result に返す。Transaction は、例えば、以下のようなリストになっている。

```
[replace([name], select((city = "Tokyo"), address),  
         Tuples, WriteTuples)]
```

Transaction は、[command1, command2, ...] というリストになっていて、各コマンドは、commandName(arg1, arg2, ...) という形式である。replace というコマンドなら、replace(AttributeListToReplace, Relation, Tuples, WriteTuples) という形になる。AttributeListToReplace には、書き換える必要のあるアトリビュートのリストをかく。アトリビュートには、必要ならリレーションの名前をプリフィックスとしてつける。リレーションには関係代数演算を書くことができる。同じ名前のリレーションが2度以上でてきてあいまいになるときは、alias(RelationName, AliasedName) を使って、リレーションの名前を変える。例えば、

```
[alias(family, family1),  
 read(join((family.father = family1.son),  
         family, family1), Ts)]
```

\*PIE64 - Interface to the database machine, YOSHIDA Minoru, TANAKA Hidehiko, University of Tokyo

```
<Transaction> ::= [<Command> {, <Command>}]  
  
<Command> ::= replace(<AttributeList>, <Relation>,  
                    <Tuples>, <Tuples>)  
                | delete(<OrgRelation>, <Relation>)  
                | insert(<OrgRelation>, <Tuples>)  
                | alias(<Relation>, <Relation>)  
                | read(<Relation>, <Tuples>)  
  
<Relation> ::= select(<Condition>, <Relation>)  
                | join(<Condition>, <Relation>,  
                    <Relation>)  
                | union(<Relation>, <Relation>)  
                | project(<AttributeList>, <Relation>)  
                | cartesianProduct(<Relation>,  
                                    <Relation>)  
                | setIntersection(<Relation>,  
                                    <Relation>)  
                | thetaJoin(<Condition>, <Relation>,  
                            <Relation>)  
                | naturalJoin(<Relation>, <Relation>)  
                | division(<Relation>, <Relation>)  
                | <OrgRelation>
```

表 1: Transaction Commands

のように使う。Tuples K は、結果のリレーションの各タプルが返ってくる。WriteTuples は、AttributeListToReplace でかかれた書き換えるべきリレーションのアトリビュートの値をいれる。そうすると、Tuples K に対応する部分のタプルが書きかわる。コマンドとしては、その他に read(RelationName, Tuples) などを用意する。トランザクションの内部では、ユーザから見て、逐次性が保たれる。例えば、

```
[read(address, Tuples), delete(address, Tuples)]
```

を 処 理 す る  
 場合には、read が終了したあと、一括して delete が行なわれるように見える。トランザクションのコマンドの一部を表 1 にあげる。例での Tuples, WriteTuples のような入出力用のストリームの管理は、他のプロセスがストリームの枠組をつくり、transaction が作ることはしない。これは、巨大なリレーションを作ってしまったとき、全てを PIE64 のメモリ上に書き出してしまい、メモリがなくなること防ぐためである。

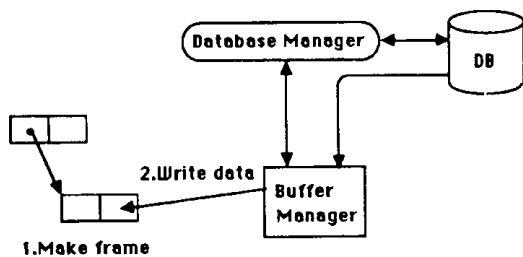


図 1: a Structure of reading data from a database machine

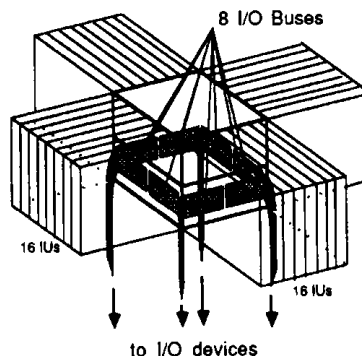


図 2: I/O interface of PIE64

### 3 処理系

図1のように、データベースマネージャはバッファマネージャに有限長バッファ通信[3]を行なわせ、例えば、データ読みだし時にはリストセルが作られたらすぐに値を代入できるようにする。有限長バッファ通信をユーザのストリーム上でやるという方法もあるが、差分リストを管理しなければならないので、ユーザにシステムプログラムほどのことをさせるのはやや負担になると思われる。しかし、そのことがユーザに見えなければ構わないので、FLENG++のなどで、ユーザから見えないようにする機構がつけば直接差分リストを管理する方法も使える。一方、PIE64とデータベース間の通信では、プロトコルを決めて通信することになるが、ポインタのようなものを渡しておき、サスペンドしたらバインドフックをかけておき、リストセルが作られたら、書き込むというような、FLENG 処理系の機構はデータベースマシン上には持ち込まない。変数の具体化の監視はPIE上で行ない、インタフェースプロセスが、上限値のあるバッファリングをして出力ストリームの具体化を待っているようにする。データベースとの間では、PIEのプロセッサ間通信のようにメモリ結合がないので、ポインタを渡すことは二度手間になるため、本当に必要なデータのみを値で渡すようにする。上の例なら、*Tuples*をprojectに必要なデータだけをPIE64にもってくるようにする。また、データベースユニット内で行えることはなるべく内部ですましてしまい、PIEはコマンドを送るだけで実際のデータはできる限り取りこまないようにする。基本的なデータベース操作はかなりの部分、データベースユニット内で行うことができるはずである。

### 4 ハードウェアのインタフェース

ハードウェアのインタフェースは、実際にどのような通信が起きるかを見て決定する必要があるが、PIE64のIU間ネットワークにはハードウェアの制約上、新たにデータベースインタフェース用のネットワークを挿入するのはほとんど困難なので、図3のように、別にネットワークを構築する。物理的には、図2のようになる。

### 5 まとめ

PIE64とデータベースマシンのインタフェースについて考察した。ハードウェア的には、それほど高速なインタフェース

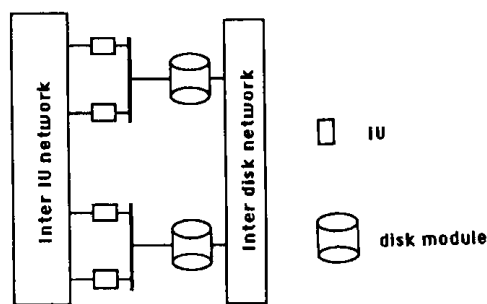


図 3: Network

は、普通のアプリケーションに対しては要求されない。ソフトウェア的には、トランザクション単位での処理を行うことによって、複数のプロセスからデータベースを整合性を保ってアクセスすることができる。

### 参考文献

- [1] 小池, 田中, “並列推論エンジン PIE64 の全体構成”, 本大会, (1988)
- [2] Nilsson M. and Tanaka H., “FLENG Prolog - Turning supercomputers into Prolog machines”, Proc. Logic Prog. Conf. '86, Tokyo, June, 1986
- [3] Takeuchi, A. and Furukawa, K., “Bounded Buffer Communication in Concurrent Prolog”, New Generation Computing, Vol.3, No.2, pp.145-155, 1985