

知識ベースマシン Mu-X (1)

—並列制御方式—

酒井 浩, 仲瀬 明彦, 柴山茂樹
(株) 東芝 総合研究所

物井 秀俊, 伊藤 英則
(財) 新世代コンピュータ技術開発機構

1 はじめに

通産省第5世代コンピュータプロジェクトの一環として、知識ベースマシン Mu-X [1] の試作が進められている。この「知識ベースマシン」という用語は、未だ明確な定義がなされていないが、Mu-Xでは第5世代コンピュータらしい特徴を備えたデータベースマシンととらえ、複数の推論マシンと結合した形態の知識情報処理システムを目標としている。

Mu-Xはマルチプロセッサシステムであり、並列処理によりデータベース演算を高速化するため、種々の工夫がなされている。本稿では、Mu-Xにおける並列処理制御方式について述べる。

2 Mu-Xの目標と設計方針

2. 1 Mu-Xの目標

Mu-Xの目標は、第5世代コンピュータらしい特徴を備えたデータベースマシンの実現である。そのため情報の表現能力を推論マシンに近づける目的で、データや問合せ等は項で表現することにした。そして、それらに対する演算として項間の単一化やその他のPrologの評価可能述語を採用した。また、ホストマシンとしてPSIを想定し、マルチトランザクション環境の実現を目指した。問合せの傾向としては、一貫性のチェックなど処理負荷の重い演算と電子化辞書の見出し語検索など処理負荷の軽い演算が混在することを想定し、レスポンスタイムとともにスループットを重視することにした。

2. 2 Mu-Xの設計方針

Mu-Xの設計にあたっては、このように機能面で高い目標を実験機規模で実現するため、図1に示すように汎用マイクロプロセッサを処理要素(PE: Processing Element)とするマルチプロセッサ構成をとり、ソフトウェア指向のシステムとした。アーキテクチャの特徴は、各PEにディスク装置をつけてデータベースシステムの課題とされる主記憶とディスク間のI/Oボトルネック解消を図ったこと、PE間の通信等に必要な共有記憶を2種類のメモリ(ワードアクセス可能な通常の共有メモリと高速ページアクセス

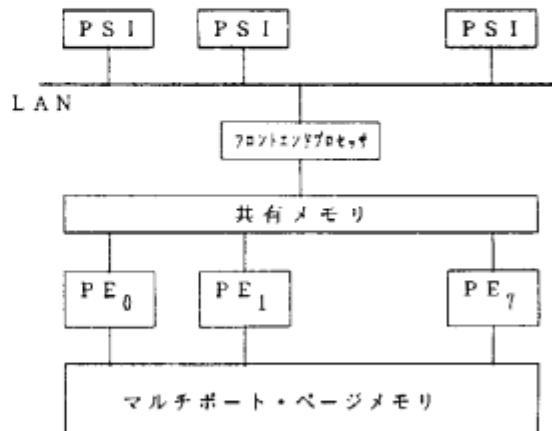


図1 Mu-Xのハードウェア構成

可能なマルチポート・ページメモリ)のハイブリッド構成とすることにより、データの参照特性に応じた使い分けを可能にしたことである。

3 PE間の機能分担と並列制御方式

3. 1 PE間の機能分担

一般にデータベースシステムで行なうべき処理機能は、ホストマシンから到着するキュエリの解析やデータ辞書の更新を行なう知識管理機能とリレーションに対する演算を処理する知識演算機能に分けることができる。従来のマルチプロセッサ構成のシステムでは、各々のPEにどちらか一方の機能を割当ることが多かった[2]。しかしながらMu-Xでは、各々のPEに両方の機能を分担させることにした。これは、2. 1で述べたように問合せの傾向として、知識管理機能に比べ知識演算機能が重い演算と軽い演算が混在するとの想定に基づいて、その比率が変動した場合でも安定した性能が発揮できるようにするためにある。

あるホストマシンから到着するキュエリの解析は、ある特定のPEで処理することにより、ホストマシンにローカルな情報は対応するPEの内部で保持することにした。これは、アクセスネットとなりやすい共有メモリに置くべき情報を減らすためである。一方、各PEの知識演算機能は

Knowledge Base Machine Mu-X Control Mechanisms for Parallel Processing
Hiroshi SAKAI¹, Akihiko NAKASE¹, Shigeaki SHIBAYAMA¹, Hidetoshi MONOI², Hidenori ITOH²
¹TOSHIBA Corporation, ²Institute for New Generation Computer Technology

すべてのホストマシンからのキュエリに応じた処理を担当させることにした。そのため、知識管理機能で並列処理用内部コマンドを共有メモリ上に作成し、各PEの知識演算機能ではその指示に基づいた処理をする方式を採用した。

3.2 並列処理制御方式に関する設計方針

先に述べたように、Mu-Xでは内部コマンドを用いて演算の並列処理を実現することにしたが、細部の設計ではキュエリ到着からレスポンス作成までの総処理量を小さくすることを目標に、次のような点に留意した。

- (1) ある演算を実現に要するコマンド数を減らす。
- (2) 並列処理アルゴリズムはPE間の同期処理のコスト等も考慮して取捨選択する。
- (3) 複数PEで同じ処理の重複が起こらないようとする。

4 並列処理用内部コマンドの概要

4.1 コマンドの意味

並列処理コマンドの意味として次の2方式を取り上げ、その得失を検討した。

- (1) ひとつのコマンドで、PE1台がなすべき演算を規定する。
- (2) ひとつのコマンドで、PEのグループがなすべき演算を規定する。

その結果、(1)の方式を採用すると、少なくとも処理にかかるPE台数だけコマンドを生成する必要があり、生成する側の処理負荷が大きくなること、各コマンドの処理の所要時間に違いがあるとPE間の負荷分散がうまく図れないという短所があることがわかり、そのような欠点のない(2)の方式を採用することにした。一般に、メッセージ通信によりPE間の同期をとるシステムでは(1)の方法をとらざるを得ないが、Mu-Xのように共有メモリを用いるシステムでは(2)の方法をとることができ、その点でMu-Xは共有メモリ方式の利点を活かしているといえる。

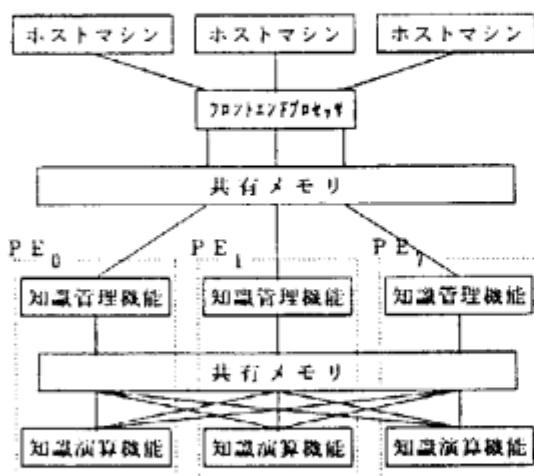


図2 Mu-XのPE間の機能分担

4.2 PE指定コマンドとPE無指定コマンド

Mu-Xでは並列処理の対象となる関係は、水平分割によりIEのディスクに分散格納されるか、マルチポート・ページメモリ上に格納されるかのいずれかとした。そのため、演算対象が前者の場合はディスク装置を有するPEに処理を委ねる必要があるが、後者の場合にはどのPEでも処理可能である。そこで、内部コマンドでは、処理すべきPEをビット列形式で指定するPE指定コマンドと、同時に実行可能なPE台数の上限だけを指定するPE無指定コマンドを用意した。

また、ひとつの内部コマンドを複数PEで処理する場合、処理の開始/終了のタイミングが各々のPEでまちまちである。そこで内部コマンドに関わるすべての処理の終了を効果的に判定する手段が必要となる。そのためコマンドごとに各PEによる処理の開始/終了を示すビット列を用意した。そして、PE指定コマンドでは要求ビット列と終了ビット列の一一致により、またPE無指定コマンドでは開始ビット列と終了ビット列の一一致によりコマンドとしての終了を判定するようにした。

4.3 並列処理の粒度とオブジェクトカウンタ

PE無指定コマンドでは、4.1で述べたように各PEがどのオブジェクトの処理をすべきか前もって決めないで、あるオブジェクトの処理が終了したものから順に次のオブジェクトの処理を開始するようにした。この処理の粒度は、選択演算ではページとした。また、結合演算ではネストドループで行なう場合はページとリレーション全体の結合演算、動的クラスタリングを前処理として行なった場合は対応するパケット間の結合演算とした。

また、Mu-Xの内部コマンドでは各PEが動的に処理対象を獲得する方式であるため、その処理を僅かなオーバヘッドで実現することが望まれる。そこでコマンドのパラメタとして次に処理すべき対象を示すカウンタを用意し、各PEが処理対象を獲得するごとにカウンタの値をその分だけ更新する方式を採用した。

5 おわりに

Mu-Xでは、各PEに知識管理機能と知識演算機能を持たせ、両者の負荷バランスが変動する場合でも安定した性能を目指している。また、両機能の間で行われるデータベース演算の並列処理の制御は、共有メモリ上に作成する内部コマンドを用いて実現している。今後、システム全体の評価を通じて方式の有効性を確認する予定である。

参考文献

- [1] 酒井 他：“知識ベースマシン Mu-X (3) 並列処理のための基本機能”，第36回情報処理全国大会予稿
- [2] Su, S. Y. W.: "Database Computer Principles, Architectures & Techniques", McGraw-Hill, 1988.