

# 超並列 Teraflops マシン TS/1

4B-4

～分散共有メモリアーキテクチャ～

鈴木 真樹 田邊 昇 菅野 伸一 小柳 滋

RWCP<sup>†</sup> 超並列東芝研究室<sup>‡</sup>

## 1 はじめに

超並列 Teraflops マシン TS/1[1] では分散共有アクセス機構により共有メモリモデルによる並列プログラムを容易に作成することができ、また細粒度通信におけるオーバーヘッドを低減することができる。

PE に汎用マイクロプロセッサを用いる時、汎用マイクロプロセッサの物理アドレス空間よりも大きな分散共有メモリを高速にアクセスするためには工夫が必要である。また、分散共有メモリを常にワード単位でアクセスしていたのではメモリアクセスのたびに結合網遅延が発生し、効率が低下する。

本稿では PE に汎用マイクロプロセッサを用いる TS/1 における分散共有メモリを効率良くアクセスするためのアーキテクチャを発表する。

## 2 TS/1 における分散共有メモリ

分散共有メモリの導入の目的は、プロセッサ間チェイニング機構ではカバーできない不規則細粒度通信の通信オーバーヘッドの削減とプログラムに対して陽に送受信の記述をしなくても良い簡便な通信手段を提供することである。

### 2.1 マルチユーザ環境のサポート

図1は、TS/1における分散共有メモリのアドレス変移を示しており、仮想記憶の管理を容易に行なうために2レベルアドレス変換方式を採用している。しかし、単に2レベルアドレス変換方式をサポートしただけではマルチユーザ環境で使用する場合にアドレス空間の断片化やプロテクション違反の問題が発生する。そこで、TS/1ではユーザ単位に独立したグローバル仮想アドレスを用意する(複数の分散共有メモリ空間を用意する)ことで、マルチユーザ環境のサポートを行なう。2レベル変換において複

数のグローバル仮想アドレスを処理するためには、2レベル目の変換の際にユーザ毎に異なるアドレス変換テーブルを用いるだけで済むが、1レベル変換でユーザ毎に異なる分散共有メモリ空間を処理する為には、リモート PE でアクセスに対する正当性をチェックするための機構が新たに必要となる。

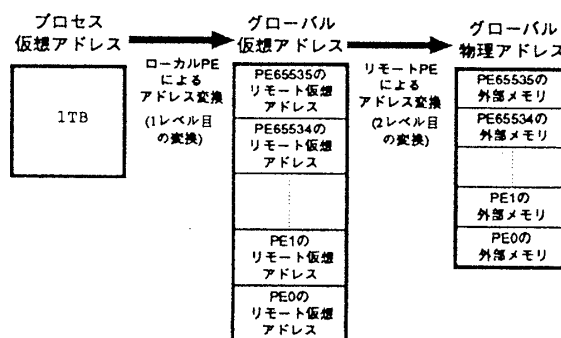


図1: TS/1の分散共有メモリのアドレス変移

### 2.2 分散共有メモリアクセス

汎用プロセッサではTS/1に実装されるようなテラバイトクラスの分散共有メモリアクセスするには物理アドレスが不足する。このため従来[2]は、分散共有メモリアクセスのとき、分散共有メモリに対するアクセスが発生する度にプロセッサに割込みをかけてソフトウェアにより処理を行っていた。しかしこれでは低オーバーヘッドの通信は不可能である。

TS/1では、上記の問題を解決するために、プロセッサ物理アドレスをセグメント方式を採用してグローバル仮想アドレスに拡張する機構を用いる。図2のようにR4400PCから出力される物理アドレス(プロセッサ物理アドレス)はTSC1のBANKレジスタの値を上位に連結されてグローバル仮想アドレスとなる。BANKレジスタを更新しないでプロセス仮想アドレスをプロセッサ物理アドレスに変換できるエントリのみをR4400PCのTLBに格納しておくことで、BANKレジスタを書き換えるときはR4400PCのTLB例外で発見できる。従って、BANKレジスタの管理はOSだけで処理され、ユーザからBANK

Massively Parallel Teraflops Machine "TS/1", - Distributed Shared-Memory Architecture - Masaki SUZUKI, Noboru TANABE, Shin-ichi KANNO, Shigeru OYANAGI

<sup>†</sup>Real World Computing Partnership(新情報処理開発機構)  
<sup>‡</sup>(株)東芝 研究開発センター 内

レジスタの存在を意識すること無く分散共有メモリをアクセスすることが可能である。

また、ユーザからは分散共有メモリのアクセスもローカルメモリのアクセスと変わらず、また、BANKレジスタの使用に関して若干のオーバーヘッドが発生するが、従来のように分散共有メモリのアクセスの度にプロセッサに割込みが発生する方式に比べてオーバーヘッドは遥かに少ないと思われるので低オーバーヘッドの通信を行なうことができる。本方式において小規模な構成でありアドレスが足りる場合は、BANKレジスタの値は変化しないのでオーバーヘッドがほとんどなく、常に割込みがかかる方式より有利である。

プロセッサ物理アドレス、グローバル仮想アドレスやBANKレジスタの値は、PE番号(X,Y,Z)やリモート仮想アドレスの項目より構成されており、プロセッサ物理アドレスを拡張するときに項目のビット境界を静的に変更できる。そのため、実装されているPE構成に見合った最適なアドレス拡張を行なうことができる。

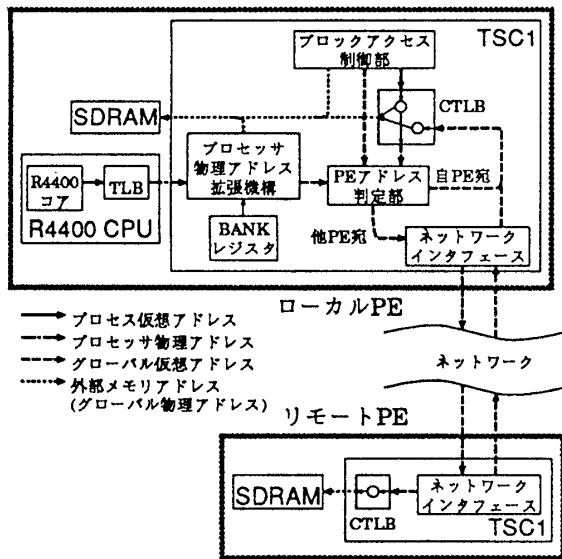


図 2: プロセッサ物理アドレス拡張機構

### 2.3 ブロックアクセス

前述の分散共有メモリアクセス機構を用いることにより分散共有メモリのアクセスにかかる通信オーバーヘッドは削減されるが、リモートの分散共有メモリのアクセスにはキャッシングしない空間をアクセスするため常にネットワークによる遅延が生じるので、ローカルメモリのアクセスに比べるとアクセス時間が遅い。また、分散共有メモリをプロセッサから出力されるプロセッサ物理アドレスでアクセスする場合、リード命令のときにプロセッサがロックされるのでプロセッサの効率が悪くなる。ネットワーク

の観点で見ても、分散共有メモリのアクセスはネットワークに細かいメッセージが多く発生するのでアドレス等のヘッダ情報の転送に通信バンド幅がとられるため、ネットワークの効率も悪くなる。

以上の問題を解決するためにTS/1ではTSC1にDMAコントローラを内蔵し、分散共有メモリのブロック単位のアクセスをサポートする。これにより、プロセッサをロックすることなく分散共有メモリのアクセスが可能となることからプロセッサの効率が向上する。また、ネットワークにはブロック状のメッセージが発生するので、ワード単位のアクセスに比べてメッセージヘッダの情報が少なくなり、ネットワークの効率も向上する。

ブロックアクセスの起動は、通信オーバーヘッドの削減のため、ユーザモードからブロックアクセスのコマンドをメモリに記憶しておき、システムプログラムにより予めプロセス仮想アドレスにマッピングされたDMAコントローラ起動レジスタに対してブロックアクセスのコマンドの先頭アドレスを書き込むことにより行なわれる。これにより、ブロックアクセスの起動は、システムプログラムの介在なしにユーザモードだけで処理を行なうことができ、低オーバーヘッドの通信が行なわれる。

また、通常のメッセージパッシングでは、受信側でメッセージの書き込みアドレスを決定するために割込みが発生し、オーバーヘッドの一要因となっているが、我々のブロックアクセスでは、書き込みアドレスもメッセージ中に指定されており、受信側での割込みが不要であるためオーバーヘッドを少なくすることができる。

### 3 おわりに

本稿では、低通信オーバーヘッドを目指した超並列TeraFlopsマシンTS/1における分散共有メモリアーキテクチャについて述べた。市販の汎用プロセッサを用いた場合でも、テラバイトクラスの分散共有メモリを低通信オーバーヘッドでアクセスすることができる方法について述べた。

今後は、詳細設計を進めながら評価を行なう予定である。

### 参考文献

- [1] 田邊, 菅野, 鈴木, 小柳: 「超並列テラフロップスマシンTS/1の構想」, 情報処理学会研究報告 93-ARC-101-6, pp.41-48, 1993.
- [2] 加納, 中田, 奥村, 大竹, 中村: 「並列マシンCenju2上の有限要素法による非線形変換解析」, JSPP '93, pp.379-386, 1993.