

# 階層記憶型並列処理の制御ソフトウェア

3 P-5

加藤 喜郎, 橋本 伸, 本間 義夫, 岡田 信, 神谷 幸男  
富士通 (株)

## 1. はじめに

科学技術計算分野への適用を目的とする、階層記憶型の並列処理システム (HPP) <sup>(1)</sup> を開発した。共用記憶装置を用いた階層記憶型のシステムでは、主記憶共用の密結合型と異なり、並列処理制御及びデータ転送のオーバーヘッドが大きくなる。本稿では、HPPシステムの制御ソフトウェアについて、効率のよい並列処理方式と共用記憶管理方式、及びユーザ支援ツールとの連携 <sup>(2)</sup> について報告する。

## 2. 制御ソフトウェアの機能

本稿で述べる制御ソフトウェアは、図1に示すHPPシステム上で、科学技術計算を高速に行うことを目的としており、以下の機能を持つ。

- 並列実行制御機能  
並列実行単位 (プロセス) の管理、同期機構など。
- 階層記憶管理機能  
CSU空間の割り当て、PES-CSUデータ転送など。
- 支援ツール連携機能  
ユーザにチューニング情報などの提供を行う、ユーザ支援ツールのためのトレース情報 (ログ) 出力機能。

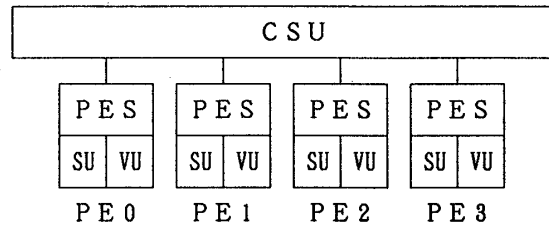


図1 HPPシステム構成図

## 3. 並列処理方式

HPPシステムでは、PE上で動作する“タスク”が仮想プロセッサとしての働きをする。

タスクは、並列実行の単位であるプロセスを実行する。Phil言語でクラスタとして記述された手続きの実行がプロセスである。同一のクラスタが2度実行されたとき、制御ソフトウェアはそれらを2個のプロセスとして管理する。

タスクとプロセスとの割り当て処理方式は、キューに付なされたプロセスを実行待ち状態のタスクがキューから外して実行する、ダイナミックなスケジューリング方式である。この方式を採用すると、PEの性能が均等でない場合でも実行状況に合わせて適当な負荷分散となることが実証できた。図2に、子プロセスの生成と実行開始の処理を表す。

### 1) プロセス生成

プロセスを新たに生成するプロセスを親プロセス、生成されるプロセスを子プロセスと呼ぶ。また、同時に生成された子プロセスを兄弟プロセスと呼ぶ。

子プロセスを生成するときは、子プロセスの入口点、パラメタなどの情報をプロセス制御ブロック (PCB) に格納する。PCBをCSU上のキューに接続することによって、その子プロセスは実行待ち状態となる。

親プロセスは、子プロセスを生成すると全ての子プロセスが終了するまで待ち状態となる。親プロセスを実行していたタスクは、他のプロセスを実行することができる。

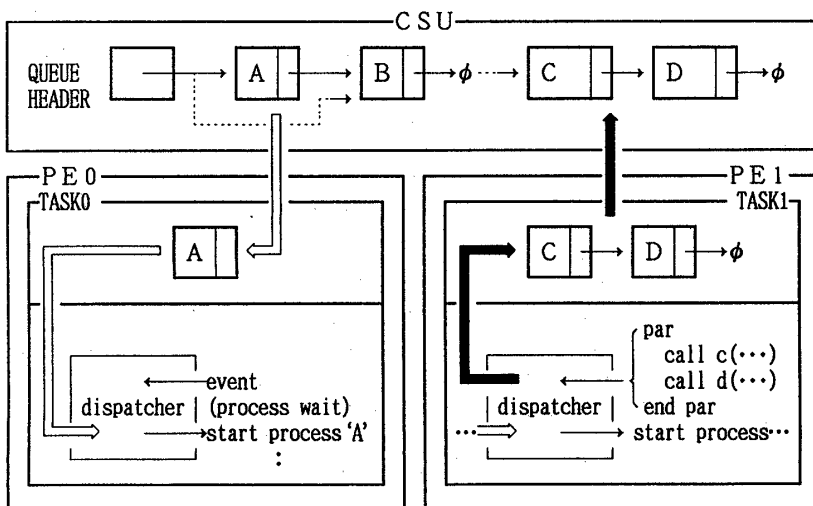


図2 プロセスの生成と実行開始

## 2) 同期機構

Phil言語で提供される, post文, wait文, lock文, barrier 文に対応したライブラリで同期機構を実現した。

post/wait及びlockでは, ユーザがCSU上に宣言したイベント変数及びロック変数をキューヘッダとして, 排他制御と待ち状態となったプロセスのキューイングを行う。

barrier では, 制御ソフトウェアがbarrier 変数をCSU上に確保する。barrier 変数にはキューヘッダの他に子プロセス数を格納するフィールドを持ち, 兄弟プロセスの待ち合わせに使用する。

## 3) 入出力機能

入出力は, PEOで動作する特定のタスク (TASK-0) で全て動作する。この機能のためPE台数が増えても, 入出力装置が1組あれば並列ジョブが動作することができる。

入出力に関する情報は, 要求元からCSU経由でTASK-0の入出力処理ルーチンに渡る。

## 4. 階層記憶管理

制御ソフトウェアは, CSUのアクセスを極力減らし, オーバヘッドが小さくなるよう最適化を行った。Phil言語で記述されたCSU-PES転送は, 実行時に可能であればPES上でバッファリングして一括転送を行った。制御情報は, 以下のような方法でアクセス回数の削減を行った。

- 初期化時に値が定まるテーブルはPESにコピーして参照する。
- 連続した領域になるよう割り付ける。

各プロセスが使用するCSU領域は, 初期化時に各タスクに分配した領域を, 各タスク専用の管理テーブルを使って管理した。管理テーブルはPESに置くことができるので, CSU領域の獲得のためにCSUをアクセスする必要がなく, 排他制御のオーバヘッドも削減できた。

## 5. ツール連携

制御ソフトウェアには, 実行時にログを出力する機能がある。ログは, 子プロセス生成, PES-CSUデータ転送開始・終了などの事象ごとに, 事象の発生時刻, 転送データ量などを記録した情報である。ログはログバッファ一杯になったとき, 及びジョブの終了時にログ用ファイルに出力する。ログファイルはユーザ支援ツールの入力データとなる。また, 簡易なツールを用いて, 実行時ライブラリ自身のデバッグにも使用した。

ログを出力すべき事象は, PES-CSUの転送最適化など, 目的により異なる。また, デバッグが目的の場合は異常終了する直前の数レコードが得られればよいこともある。この問題解決には, チューニングツール又はデバッグツールとの, より有機的な連携が必要である。

## 6. 成果

HPPシステムの性能評価のために, 以下の評価用プログラムの並列実行を行った。

- SPIN: 磁場解析, 転送行列法
  - 円周率計算: 並列FFTによる多倍長計算
  - NIEL: 橋の強度解析
  - 帯行列解法: 疎行列係数の連立一次方程式の解法
- なお, SPIN及び帯行列解法は, 高水準並列処理記述言語PARAGRAPH, 他はPhil言語で記述されたプログラムである。

表1は, 評価プログラムの実行時間(経過時間)及び台数効果である。表2は, 制御ソフトウェアのオーバヘッド及びグラニュラリティの代表的な値である。

表1 評価プログラムの実行性能  
経過時間 単位: 秒

プログラム	1タスク	4タスク	
	経過時間	経過時間	台数効果
帯行列解法	284	100	2.8
SPIN	1811	479	3.8
円周率計算	3521	1375	2.6
NIEL	228	108	2.1

表2 オーバヘッドとグラニュラリティ

子プロセス生成	1 ms
lock文 (成功)	0.7 ms
lock文 (失敗, 別プロセス起動)	0.8 ms
評価プログラムのグラニュラリティ	0.8~7.8 sec

## 7. おわりに

実用プログラムの並列実行に成功し, HPPシステムの有効性が実証できた。評価に用いた大規模な科学技術計算プログラムでは, 制御ソフトウェアのオーバヘッドは, グラニュラリティに比べて十分小さいことが分かった。

ツール連携のためのログ出力機能は, ユーザ支援のみでなく, 制御ソフトウェア自身のデバッグにも有効であった。

**謝辞** 本研究は, 通商産業省工業技術院大型プロジェクト「科学技術用高速計算システムの研究開発」の一環として, 新エネルギー・産業技術総合開発機構(NEDO)から委託を受けて, 実施したものである。

御協力頂いた, 日本原子力研究所物理部 別役主任研究員, 日本原子力研究所東海研究所計算センター 横川氏, 東京大学大型計算機センター 金田助教授, 京都大学工学部土木工学教室 渡邊教授に感謝する。

## 【参考文献】

- (1) 橋本, 他: 階層記憶型並列処理, 情報処理学会第41回全国大会, 1990
- (2) 渡辺, 他: 並列処理におけるプログラム開発支援システム, 情報処理学会第41回全国大会, 1990