

3P-3

# 並列処理用OS SKY-1の スケジューリング方式

上脇 正 山口 伸一郎 齊藤 雅彦

日立製作所 日立研究所

## 1. はじめに

密結合マルチプロセッサ(TCMP)の多くは、各プロセッサがキャッシュメモリを持っており、このキャッシュメモリが、システムの性能上非常に重要な役割を果たす。

しかし、スレッドのスケジューリングの方法によっては、キャッシュメモリの性能が大きく低下してしまう問題がある。例えば、スレッドがプロセッサを移動すると、今まで蓄えられたキャッシュデータが無駄になり、キャッシュミス連続しておこす<sup>1)</sup>。

また、キャッシュメモリがバスをスヌープして自分が記憶しているデータに対する書き込みを検出すると、そのデータを無効化するTCMPでは、共有データを持つスレッド同志を異なるプロセッサで実行すると、キャッシュメモリの無効化が頻発する問題がある<sup>2)</sup>。

そこで、我々が試作しているTCMP用OS・SKY-1(System Kernel for You)<sup>3)</sup>では、これらの問題を解決する手法をスレッドスケジューリング方式に組み込んでいる。

## 2. SKY-1スケジューリング

SKY-1スケジューリングの戦略は3つある。

### 2.1 同一プロセッサ優先

同一プロセッサ優先は、タイムスライスの終了、イベント待ちの発生等により一度中断されたスレッドを再開するときには、可能な限り前回実行されていたプロセッサで実行する方式である。この方式によれば、スレッドの不必要なプロセッサ間移動がなくなり、移動によるキャッシュミスが減少する。

### 2.2 スレッドのグループ化

スレッドのグループ化は、図1のようにスレッドをグループとして管理し、同一グループに属するスレッドを同一のプロセッサで実行するものである。スレッドT<sub>0</sub>、T<sub>1</sub>はグループG<sub>0</sub>に属し、プロセッサP<sub>0</sub>で実行される。

先に述べたようにデータを多く共有するスレッドを異なるプロセッサで実行すると、キャッシュ一致化のための無効化が頻発する。共有データの多いスレッド同志をグループにして、同一プロセッサで実行すれば、共有デー

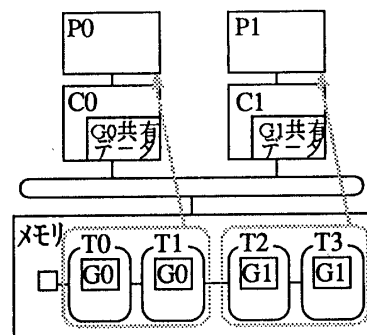


図1 スレッドグループ化

タが複数のキャッシュメモリに入ることが少なくなる。図1ではグループG<sub>0</sub>のデータはキャッシュC<sub>0</sub>に集中する。するとキャッシュ一致化のための無効化が減少、キャッシュヒット率が向上する。

スレッドグループ化には、他の使いかたもある。同一タスクに属するスレッドをグループ化すると、スレッドスイッチのオーバーヘッドが減少する。つまり、タスクの異なるスレッド同志のスイッチは、同一タスクのスレッドのものに比べて遅い。同一タスクに属するスレッドをグループ化すると、スレッドスイッチが同一タスクのもの同志であることが多くなり、スレッドスイッチのオーバーヘッドが減少する。

また、資源を共有するスレッドをグループ化すると、排他制御を簡単化できる。ある資源を使用する全てのスレッドが一つのプロセッサのみで実行されることが判っている場合、スレッドスイッチと割り込みを禁止すれば、資源を使用するときにロックする必要がない。ロック処理は時間がかかるため、処理速度向上に役立つ。また、資源の排他制御を割り込み禁止のみで行っているシングルプロセッサ用ソフトウェアも、この機能によれば、改造なしにマルチプロセッサで実行可能となる。

### 2.3 負荷分散

同一プロセッサ優先とスレッドのグループ化を行っていると、あるプロセッサは処理するスレッドが沢山あるにも拘らず、アイドルとなるプロセッサが出来ることがある。そこで、プロセッサ間での負荷を平均化するため

に、一定時間実行されていないスレッドについては、プロセッサの移動を認めた（移動を禁止することも可能）。

実行されていないスレッドのデータは、徐々に他のスレッドのデータによりキャッシュメモリから追い出される。長く実行されていないスレッドは、キャッシュメモリにデータが少なく、移動による性能への影響が少ない。

3. 性能評価

SKY-1スケジューリング方式の性能を評価するために、図2のモデルを用いて解析的シミュレーションを行った。モデルには、68030(33MHz)を4台接続したマルチプロセッサを用いた。スレッドは、平均250msに1個生成され、発生したスレッドは実行待ち行列に入り、スケジューリングされるのを待つ。スケジューラがスケジューリングアルゴリズムに従い、実行待ち行列のスレッドをプロセッサに割り振る。タイムスライスが1sで、それを使い終わったスレッドは次のスレッドとスイッチして実行待ち行列の最後に戻される。I/Oアクセスが発生したスレッドもスイッチを行う。キャッシュミス：ヒットのアクセス時間比は10：1とした。

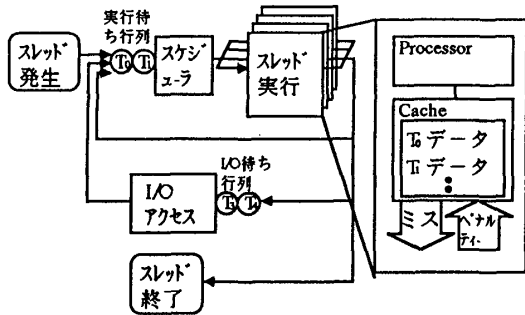


図2 シミュレーションモデル

SKY-1方式を、スレッドスイッチの必要となったプロセッサが、単純に実行待ち行列の先頭から、スレッドを取り出して実行する従来方式と比較したところ以下のような結果が得られた。

(1) キャッシュミス率のキャッシュ容量依存性

図3はキャッシュメモリの容量を変化させたときのキャッシュミス率の変化を表している。キャッシュ容量が128kB時、SKY-1方式は、キャッシュミス率が45%（同一プロセッサ優先により15%、スレッドグループ化により30%）低減する。これは13%のプロセッサ性能向上に相当する。また、キャッシュ容量が2MBと多くなると、キャッシュミス率の低減も65%と大きくなる。

(2) キャッシュミス率の共有データ量依存性

図4は、スレッド間の共有データ量（あるスレッドのライトアクセスが、他のスレッドのキャッシュデータを無効化する確率）とキャッシュミス率の関係のグラフで

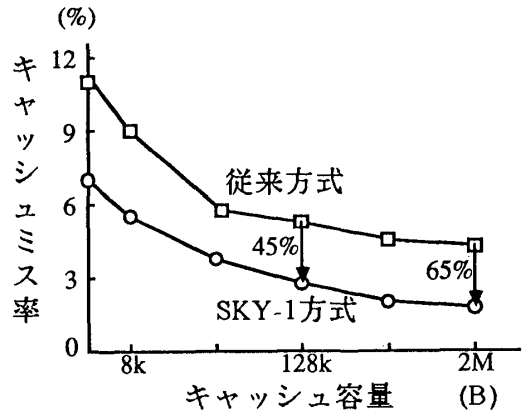


図3 キャッシュ容量依存性

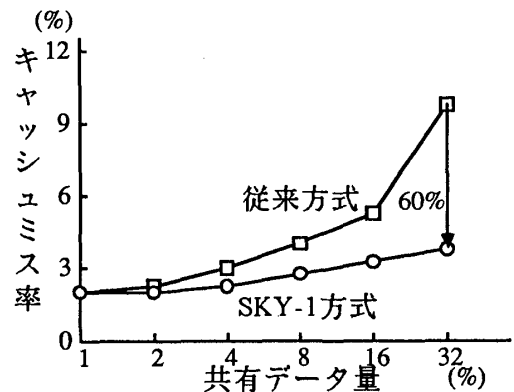


図4 共有データ量依存性

ある。共有データ量が増加すると、従来方式はキャッシュミス率が急増するが、SKY-1方式ではほとんど変化がない。共有データ量が32%の時、SKY-1方式によってミス率が60%減少し、プロセッサ性能が2.6%向上する。プロセスに代わりスレッドを使用するようになると、共有データ量は必然的に多くなるため、本方式が重要となってくる。

4. おわりに

並列処理OS・SKY-1のスケジューリング方式について述べた。本方式の特徴は、同一プロセッサ優先とスレッドグループ化によりキャッシュメモリの性能を十分に引き出すことにある。

5. 参考文献

- 1) Sites, R. L., et al.: Multiprocessor Cache Analysis Using ATUM: Proc. ISCA (1988)
- 2) Eggers, S. J., et al.: The Effect of Sharing on the Cache and Bus Performance of Parallel Programs: Proc. ASPLOS III (1989)
- 3) 山口 他: 並列処理用OS SKY-1 - 開発構想 - : 情報処理学会第39回全国大会 (1989)
- 4) 齊藤 他: 並列処理用OS SKY-1のメモリ管理方式: 情報処理学会第39回全国大会 (1989)