

5 Q-4 密結合型マルチプロセッサ上での自律的な動的負荷分散制御
 Automatic and Dynamic Load Balancing Scheme on Tightly Coupled MultiProcessors

中山善太郎 (Yoshitaro Nakayama) 大川善邦 (Yoshikuni Okawa)
 NEC 交換移動通信事業本部* 日本大学工学部情報工学科†

1 はじめに

本稿では、プロセッサ同士が局所的に結合する MNC 並列計算機において、各プロセッサが自身の状態に合わせて負荷分散と問題処理をバランス良く行なうための自律分散アルゴリズムを提案し、実験とシミュレータによる解析を行なった。

2 MNC 並列計算機

本稿ではプロセッサ (Processor Element : PE) が 4 ポート共有メモリ (Quad Port RAM : QPR) により局所的に結合する MNC(Maximum Neighbor Connection: 最大隣接結合) 並列計算機を用いる。MNC 並列計算機は、ある PE から見て1つの QPR を介して通信可能なプロセッサ数が最大になる結合方式の 1 つである。したがって、自律分散制御を行なう場合、より多くの隣接 PE から情報を得ることができる有利な構造である。

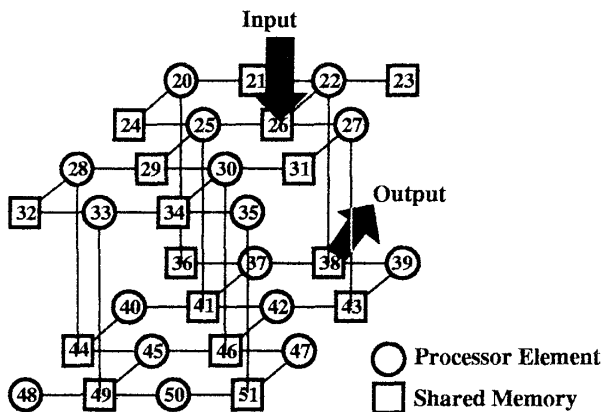


図 1: 実験に用いた MNC 並列計算機の結合

3 処理問題

データ処理時に PE 間で依存関係のない $y = f(x)$ 型の問題を扱うことにした。各 PE はデータ処理のための関数 $f()$ をあらかじめ持つ。ソースから並列計算機

*Switching and Mobile Systems Operations Unit, NEC Ltd.
 †Dep. of Information engineering, Faculty of Engineering, Nihon University

に入力される大量のデータ $x_i(i = 1, 2, \dots)$ (データ X と呼ぶ) は、各プロセッサで並列的に y_i (データ Y と呼ぶ) に変換され、並列計算機からシンクに出力される。 x_i と y_i は 1 対 1 で対応しており、すべてのデータについて、データ X から Y への変換時間は同一であることを想定している。

4 自律分散アルゴリズム

自律分散制御では並列計算機内部の各 PE が共有メモリを介して得られた局所情報から自律的に自身の状態を判断し、処理を行なう。各 PE が行なう処理は、隣接する PE との情報交換、取得した情報に基づく行動計画、行動の実行であり、各 PE がこれらを繰り返し実行することで全体として処理が進行する。PE が処理を行なうときには、

- 制御に用いる情報
- PE が行なう行動の種類と内容
- 行動の計画方法

が重要になってくる。

4.1 参照情報

各プロセッサ PE_i が制御に用いる情報を、次の状態ベクトル X とした。

$$X = [X_i, Y_i, RX_{ij}^B, RY_{ij}^B] \quad (1)$$

表 1 に、式 (1) に出てくる状態量の意味を示す。

表 1: 制御の状態量

X_i, Y_i	PE _i が保持している X・Y の量
X_{ij}, Y_{ij}	PE _{ij} が保持している X・Y の量
SX_{ij}^B, SY_{ij}^B	PE _i から PE _{ij} 宛の X・Y の量
RX_{ij}^B, RY_{ij}^B	PE _{ij} から PE _i 宛の X・Y の量

4.2 行動の種類

PE_i が行なう行動は X を Y に変換する問題処理と負荷分散のための転送処理の 2 種類である。データには X

と Y の 2 種類があり, MNC 並列計算機では QPR を介して PE 間のデータ転送を行なうため, PE から QPR への転送を PUT, その逆を GET とすると 4 つ転送処理が考えられる.

- 問題処理のための行動 (X から Y の変換)
 - PROC : $f()$ により X_i を Y_i に変換
- 負荷分散のための行動 (データ転送)
 - XPUT, YPUT : X_i, Y_i をローカルメモリから QPR に転送
 - XGET, YGET : RX_{ij}, RY_{ij} を QPR からローカルメモリに転送

効率的な処理を行うためには, これら 5 つの行動を組み合わせて実行しなければならない.

4.3 各行動の実行判断

XGET, YGET, YPUT については, 目標状態が一意であるため, 隣接 PE 間の平均データ量に基づいて判断した. XPUT と PROC については, PE_i の負荷の適正值 \widetilde{X}_i よりも保持データが多いときには XPUT, 少ないときには PROC を行なうようにした.

4.3.1 プロセッサ PE_i の負荷の適正值 \widetilde{X}_i

各時刻における, 自 PE を基準にして隣接 PE の処理能力を相対的に表す状態量 PA (Processor Ability) を用いて, PE_i の負荷の適正值 \widetilde{X}_i を,

$$\widetilde{X}_i = \frac{\sum_j (PA_{ij} X_{ij}) + X_i}{\sum_j d_{ij} e_j + 1} \quad (2)$$

とした. ただし, PA_{ij} は以下の式で表される.

$$PA_{ij}(t) = PA_{ij}(t - t_{pa}) + \frac{G_{pa}}{t_{pa}} \int_{t-t_{pa}}^t \frac{(X_i(s) - X_{ij}(s)) |X_i(s) - X_{ij}(s)|}{\{X_i(s)\}^2} ds \quad (3)$$

式 (3) の記号は G_{pa} :ゲイン, t_{pa} :PA 更新時間間隔, X_{ij} : PE_i から見た PE_{ij} の保持データ量を示す.

PA_{ij} は PE_i を基準とした PE_{ij} への PA である. 式 (3) より, PA は, 将来の自 PE と隣接 PE の保持データ量の比を 1 に近づけるように作用することがわかる.

5 実験

16 台の PE と 16 台の QPR から構成される図 1 の MNC 並列計算機に制御アルゴリズムを実装し, 実験を行なった. 処理問題は様々な重さの処理を実現できるように, C 言語の for 文による空ループとした. データ X は for 文のループ数を含むサイズ 8[byte] とし, 1,000,000[個] の X をソースの QPR26 から計算機に流

し込んだ. データ Y のサイズも 8[byte] とし, シンクの QPR38 から回収した.

PE 数 16 までの部分は実機を用い, それ以降はシミュレータを用いて, PE 台数を 1~128 の間で変化させたときに, 処理効率, 総転送データ数を測定した結果を図 2, 3 に示す. 処理問題は 100[loop/data], 1000[loop/data] の 2 種類を取り扱った.

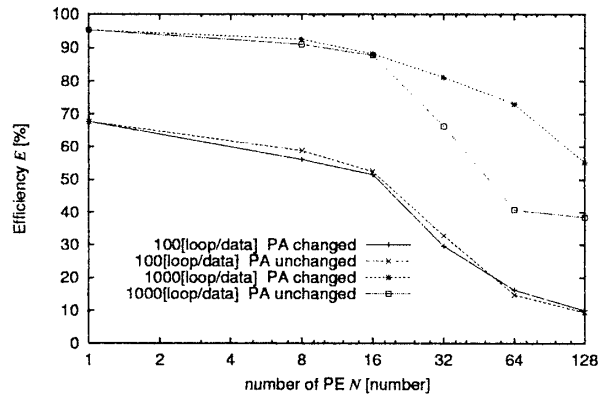


図 2: シミュレータによる処理効率の測定

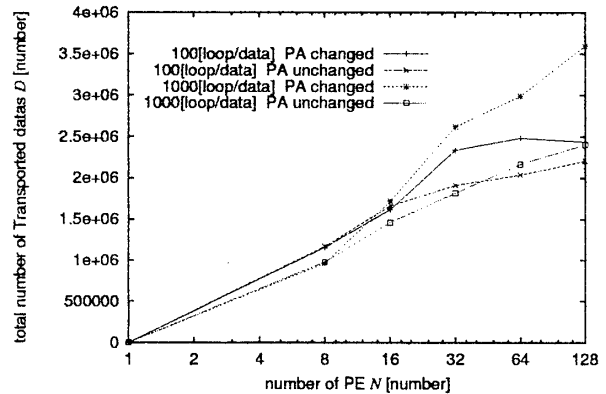


図 3: シミュレータによる総転送データ数の測定

図 2, 3 より, 負荷の軽い 100[loop/data] 問題に対しては, データ転送を抑制して転送オーバーヘッドを削減でき, 重い 1000[loop/data] 処理問題に対しては, 積極的なデータ転送による負荷分散を行なって, 計算機の資源を有効に利用できることが確認できた.

6 おわりに

本稿ではプロセッサ同士が共有メモリにより局所的に結合する MNC 並列計算機における自律分散アルゴリズムを提案した. このアルゴリズムは, ある処理問題に対して, それが必要とする計算機の規模を動的に推定して, 計算機の資源を有効に利用する制御が行なえることを確認した.