

## ボランティアコンピューティング環境における動的クラスタ再構成法

菅原 雅也<sup>†1</sup> 福士 将<sup>†1</sup> 堀口 進<sup>†1</sup>

ボランティアコンピューティング (VC) 環境において計算機 (ノード) 同士の同期的な通信を含む並列計算を実行可能にするために、既存の計算モデルを拡張し、さらにシステムの性能管理手法としてノード群 (クラスタ) の動的再構成法を提案する。この手法では、ノードをネットワーク的な距離に基づいて階層的にクラスタリングすることで並列計算実行時の通信オーバーヘッドを削減し、クラスタ間で動的に所属ノードの移動を行うことで処理性能の調整を図る。シミュレーションによる評価から、大規模な並列計算を実行する VC において、本提案手法により最大で 16% の処理効率の向上が得られることを確認した。

### Dynamic Cluster Reconstruction for Volunteer Computing

MASAYA SUGAWARA,<sup>†1</sup> MASARU FUKUSHI<sup>†1</sup>  
and SUSUMU HORIGUCHI<sup>†1</sup>

This paper proposes a novel computation model and dynamic node-clustering method for Volunteer Computing (VC) systems, motivated by the need to handle parallel computation problems which include synchronous communications between computing nodes. The proposed method constructs several clusters of nodes based on their network distance to reduce the overhead of parallel communications, and allows node migration between clusters to control overall performance of VC systems. Simulation results indicate that the proposed method improves computational efficiency up to about 16 %.

#### 1. はじめに

近年、インターネット上の余剰計算機資源を活用して、大規模な並列分散処理環境を構築するボランティアコンピューティング (VC) が研究され、そして実際に運用されている。VC では、参加者がある計算プロジェクトに自主的に参加し、自身のパーソナルコンピュータ (PC) などの余剰計算資源を無償で提供する。スーパーコンピュータや高性能コンピュータを接続したようなグリッドシステムと比較すると、VC における個々の計算資源の性能は格段に低い。しかし、SETI@home<sup>1)</sup> や Folding@home<sup>2)</sup> の実用例に見られるように、インターネット上の膨大な数の PC を集めることで、それらのグリッドシステムを凌ぐ計算力をほぼ無償で入手できることが VC の大きな利点である。

これまで VC に適用されてきた計算は、高い並列性があり、ノード間通信をほとんど必要としない“単純並列処理”が可能な大規模並列計算である。これは、多数の PC に計算を分担させるために計算問題が分割可能であること、また通信が必要な場合、電源が切られるなどの理由で通信相手が不在になったときのペナルティが大きいこと、などが理由として挙げられる。しかし、科学技術計算などの全ての大規模並列計算がこれらの特性を持つとは限らないため、VC に適用可能な計算問題が限定

されることは望ましいことではない。

そこで本稿では、高性能な計算力を安価に提供できる VC 環境において、同期的な通信を含む大規模な並列計算問題を処理可能とすることを目的に、下記 2 点を提案する。

- (1) ノード間で同期的な通信を含む場合の計算モデル
- (2) ノードの動的な参加、離脱に対応するためのシステムの性能管理手法

(1) については、従来のマスタ・ワーカモデルを拡張し、同期的な通信を含むタスクをノード群 (クラスタ) に割り当てる計算モデルを提案する。また、通信オーバーヘッドを考慮するために、ノード間のネットワーク距離に応じた簡易通信コストのモデル化を行う。(2) については、ノードの動的な参加、離脱が発生する環境下で、クラスタ内でのノードのネットワーク距離とクラスタ間でのノード数を調整するための動的クラスタ再構成法を提案する。シミュレーションによる性能評価により、同期的な通信を含む大規模な並列計算問題を処理する VC 環境において、本提案手法により処理効率の向上が得られることを明らかにする。

#### 2. VC 環境とノードのクラスタリング

##### 2.1 既存の並列計算モデル

現在運用されているほとんどの VC では、計算モデルとして図 1 に示すマスタ・ワーカモデルを採用している。以下にその概要を示す。

- 処理対象である計算プロジェクトは複数の独立なタ

<sup>†1</sup> 東北大学 大学院情報科学研究科  
Graduate School of Information Sciences, Tohoku University

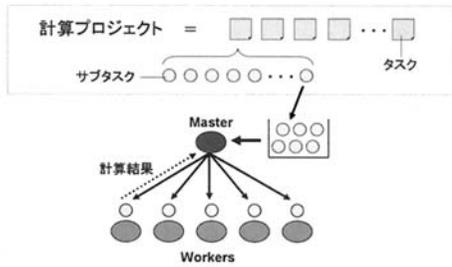


図1 既存のマスター・ワーカ型計算モデル

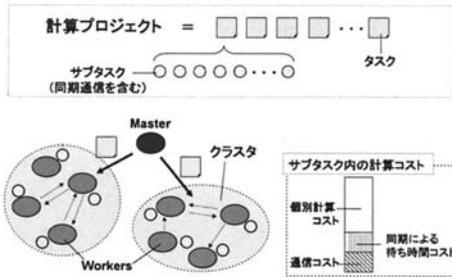


図2 タスク内に通信を含むマスター・ワーカ型計算モデル

スクに分割され、各タスクは複数の独立なサブタスクに分割される。

- VC システムはマスターと複数台のワーカから成る。
- マスタは、あるタスクの各サブタスクをワーカに対して割り当て、処理を依頼する。
- 各ワーカは、割り当てられたサブタスクを計算し、計算終了後に結果をマスターに返却する。

全てのサブタスクが処理された段階でそのタスクの計算が終了し、同様に全てのタスクが処理されると計算プロジェクトの終了となる。

## 2.2 同期通信を含む並列計算モデル

VC の応用範囲を広げるために、ノード（ワーカ）間の通信を考慮した並列計算モデルを新たに定義する。提案する計算モデルの概要を図2 および以下に示す。

- 計算プロジェクトは複数の独立なタスクに分割され、各タスクは互いに同期的な通信を行う複数のサブタスクに分割される。
- VC システムはマスターと複数のワーカから成り、各ワーカはタスクの処理単位として  $D$  個のクラスタに分類される。
- マスタは、あるタスクのサブタスクを、そのタスクを担当するクラスタ内のワーカに割り当てる。
- ワーカは、他のサブタスクを処理するワーカと通信を行いながら、割り当てられたサブタスクを計算し、結果をマスターに返却する。

本計算モデルでは、サブタスク間に通信を含むタスクが  $D$  個並列に計算される。この点が従来のモデルとは大きく異なる点である。

次に通信処理に要するコストをモデル化する。図2の

ように、各ノードにおいて、サブタスクの計算コスト（個別計算コスト）に加え、同期による待ち時間コスト（待ち時間コスト）と通信に要するコスト（通信コスト）を定義する。一般的な並列計算では、計算の途中で複数回の通信が発生する。本稿ではモデルの簡略化のため、それらの複数回の通信を1つにまとめた形で、1つの通信コストとして取り扱う。通信が開始可能になるまでの待ち時間コストも同様の取扱いをする。

各ノードは通信の待ち時間の間に新たなサブタスクをマスターから受け取り、その計算処理を行うことができるものとする。この場合、通信待ちのサブタスクの通信が開始できるようになったら、計算中のサブタスクの計算処理を一旦中断し、通信処理を行えるものとする。ワーカがVC環境から離脱した場合は、計算および通信の進行状況に関わらず、所持していたサブタスクの計算結果は全て破棄されるものとする。

## 3. 動的クラスタ再構成法

### 3.1 ノードのクラスタリング

2.2 節で定義した計算モデルでは、複数のタスクを並列に計算するために、タスクの処理単位として、ノードをクラスタリングする必要がある。この際、通信オーバーヘッドを削減するために、なるべく近いノード同士をクラスタリングするほうが望ましい。さらに、ノードの参加・離脱やタスクの粒度の変化に対応するために、任意数のノードから成る任意数のクラスタを生成する必要が生じる。

一般的に、多数のノードからなる大規模な分散処理環境では、スケーラビリティの観点から階層的なクラスタリングが適している。Banerjee ら<sup>3),4)</sup> や上田ら<sup>6)</sup> の手法では、通信用の物理ネットワークの上にノード管理用の論理的な階層型クラスタ構造を構築し、ノードの新規参加や離脱などの状況変化に応じて、これを動的に更新する機構を有している。しかしながら、ネットワーク距離に応じたクラスタリングは可能であるものの、任意数のクラスタの生成やクラスタ内ノード数の調整機能がないため、そのままの形で提案する計算モデルに適用することは困難である。

本稿では、Banerjee ら<sup>3)</sup> の提案した階層型クラスタ構造を基にして、クラスタへの新規参加、サブクラスタの分割、リーダ選出などの手続きを明確に体系化する。またこれに加えて、クラスタ間で所属ノードの受け渡しを行う手続きを定義することで、クラスタ間における処理性能の調整を可能とする動的クラスタ再構成法を提案する。

### 3.2 階層型クラスタの基本構造

図3にBanerjee ら<sup>3)</sup> により提案されたノード管理のための階層型クラスタ構造を示す。図3に示されるように、全てのノードはネットワーク上で距離（ホップ数など）に応じて階層的にクラスタリングされる。各階層でサブクラスタが形成され、サブクラスタ内のノード（ネイバーノードと呼ぶ）は Heartbeat 信号と呼ばれる制御

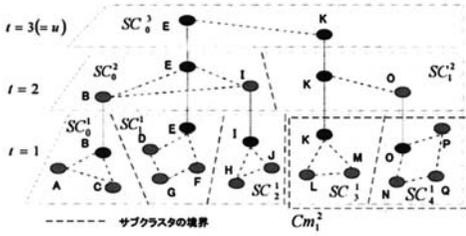


図3 階層型クラスタの基本構造

### Join ( $n_j$ )

- (1)  $t = u, n_r = L_i^u$  とする。
- (2) 新規参加ノード  $n_j$  は,  $n_r$  から  $SC_i^t$  のメンバリストを取得する。
- (3)  $n_j$  は取得したメンバリストに記されたノード (一つ下の階層のリーダー  $L_i^{t-1}$  に相当する) と自ノードとの距離測定を行い, 最も近いノードを新たな  $n_r$  とする。
- (4)  $t = 2$  の場合: 新規参加ノード  $n_j$  は  $n_r$  と同じ最下位層サブクラスタに所属する。  
 $t \neq 2$  の場合:  $t = t - 1$  として手順 2 へ戻る。

図4 新規ノード参加手続き

信号を互いに送信し合い, 距離情報などを交換する。以下, あるクラスタ  $C_i$  において  $t$  階層目のサブクラスタを  $SC_i^t$  で, また  $SC_i^t$  より下層の全てのノードの集合 (コミュニティと呼ぶ) を  $Cm_i^t$  で表す。ここで, 各クラスタにおける最下位層を第 1 層, 最上位層を第  $u$  層とし, 同一階層のコミュニティを区別するために便宜的な順序 ( $i$ ) を与える。各サブクラスタ  $SC_i^t$  にはリーダー  $L_i^t$  が存在し, 最上位層を除く全ての階層のサブクラスタにおいて, リーダは一つ上の層のサブクラスタのメンバになっている。図 3 に示されるように, VC に参加する全ノードは最下位層のいずれかのサブクラスタのメンバであり, リーダノードのみが  $t \geq 2$  の上位階層のメンバになっている。

ノードはネイバーノードの情報を各階層ごとにクラスタ表として保持し, 随時更新する。Heartbeat 信号の応答が一定時間返ってこない場合にそのノードが離脱したとみなし, ネイバーノードは該当するクラスタ表から離脱したノードの情報を削除する。

### 3.3 クラスタ再構成の手続き

#### 3.3.1 新規ノード参加アルゴリズム (Join)

VC に新規に参加するノード  $n_j$  は, あるクラスタの最下位層サブクラスタに所属する。 $n_j$  は, クラスタの最上位層から順に, 距離的に近いノード  $n_r$  を探索する動作を各階層で再帰的に実行することで, 所属先のサブクラスタを選択する。近いノード同士が同一のクラスタに所属することで, 以降に新規に参加するノードも同様に近いサブクラスタへ所属することができる。図 4 にアルゴリズムを示す。

#### 3.3.2 サブクラスタ分割アルゴリズム (Split)

第  $t$  層のサブクラスタ  $SC_i^t$  は, ネイバーノード数が  $k$  ノードを超えた場合に, 図 5 に示す K-means 法を横し

### Split ( $SC_i^t$ )

- (1)  $L_i^t$  が任意の 2 ノード  $n_a$  と  $n_b$  を分割後のサブクラスタ  $SC_a^t$  と  $SC_b^t$  のリーダー候補として選出する
- (2)  $L_i^t$  は  $SC_i^t$  内のネイバーノードに対し, 選出した  $n_a$  と  $n_b$  を通知する。
- (3)  $SC_i^t$  の各ノードは,  $n_a, n_b$  との距離を測定し, それぞれ  $SC_a^t, SC_b^t$  の近い方へ所属する。
- (4) 手順 3 で構成した  $SC_a^t, SC_b^t$  のそれぞれにおいて, 他ノードとの距離の最大値が最小のノードを選出し, この 2 つのノードを新たに  $n_a, n_b$  とする。
- (5) 手順 2~手順 4 を一定回数繰り返すことで  $SC_i^t$  を 2 つのサブクラスタ  $SC_a^t, SC_b^t$  に分割する。

図5 サブクラスタ分割手続き

### Merge ( $SC_i^t$ )

- (1)  $SC_i^t$  のリーダー  $L_i^t$  はネイバーノードに対し, 自身の所属する第  $t+1$  層のサブクラスタ  $SC^{t+1}$  のメンバリストを通知する。
- (2) 通知を受けた各ノードは取得したメンバリストの中から最も近いノードを探索し, そのノードが所属する第  $t$  層サブクラスタ  $SC_j^t$  ( $j \neq i$ ) のメンバとして所属先を変更する。
- (3)  $L_i^t$  も手順 2 と同様にして,  $SC^{t+1}$  のメンバの中から最も近いノードを探索し, 第  $t$  層のメンバとして所属先を変更する。またこの際に  $L_i^t$  は  $SC^{t+1}$  のメンバから削除される。

図6 サブクラスタ統合手続き

たアルゴリズムにより,  $SC_i^t$  を二つのサブクラスタ  $SC_a^t$  と  $SC_b^t$  に分割する。分割の際に  $SC_a^t, SC_b^t$  のリーダーとして新たに選出されたノード  $n_a$  と  $n_b$  は, 第  $t+1$  層のサブクラスタのメンバとなる。同時に分割前のリーダー  $L_i^t$  は第  $t+1$  層のメンバリストから削除される。また分割がクラスタの最上位層で行われると, クラスタの階層が一つだけ増え,  $n_a, n_b$  は新しい最上位層サブクラスタのメンバとなる。このためノード数の増大に応じて分割が発生することによりクラスタの階層が適宜増加していく。

### 3.3.3 サブクラスタ統合アルゴリズム (Merge)

Split アルゴリズムとは逆に,  $t$  階層目のサブクラスタ  $SC_i^t$  内のノード数が一定値  $l$  を下回った場合に, 図 6 に示す Merge アルゴリズムにより, サブクラスタを統合させる。Merge により,  $SC_i^t$  の所属メンバは,  $L_i^t$  が所属する第  $t+1$  層のサブクラスタを経由して, 別な第  $t$  層のサブクラスタへ所属し直す。また, 本稿では, 最上位層サブクラスタ  $SC^u$  内の所属ノード数が一台だけとなってしまった場合に,  $SC^u$  を削除しクラスタの階層を一階層分減らすこととした。第  $u-1$  層において Merge アルゴリズムが発生することにより, 所属ノード数の減少に応じてクラスタの階層数を減らすことが可能となる。

### 3.3.4 クラスタ間でのノード移動アルゴリズム (Transfer)

ある特定のクラスタの所属ノード台数が大きく減少すると, このクラスタに与えられたタスクの処理に大きな遅れが生じ, 結果的に VC 全体で著しく処理性能が低下する可能性がある。図 7 に示す Transfer アルゴリズムにより, クラスタ  $C_i$  からクラスタ  $C_j$  ( $j \neq i$ ) へ  $s$  台程

**Transfer** ( $C_i, C_j, s$ )

- (1)  $C_j$  の最上位層サブクラスタ  $SC_j^u$  のリーダー  $L_j^u$  は  $C_i$  に対し Join アルゴリズムを実行し、最も距離の近い最下位層サブクラスタ  $SC_i^l$  を探索する（ここでは探索のみを行い所属動作は行わない）。
- (2) 手順1で到達した  $SC_i^l$  から上位層に向かって、所属ノード数が  $s$  以上のコミュニティ  $Cm_i^h$  を順次探索する。
- (3)  $SC_i^h$  の所属ノード  $L^{h-1}$  が  $C_j$  に対してそれぞれ Join アルゴリズムを実行する。
- (4)  $Cm_i^h$  の他のノードは、第  $h-1$  層のリーダーが手順3で到達した最下位層のサブクラスタへ移動する。

図7 クラスタ間でのノード移動手続き

**Leader-election** ( $SC_i^l$ )

- (1)  $L_i^l$  の離脱を感知した段階で、ノードは自身のクラスタ表を参照し、 $SC_i^l$  内の他ノードとの最大距離値が最も小さいノードを新規リーダー候補とみなして、このノードに対してクエリする。
- (2) 他ノードからのクエリを一定数受けたノードは、新しいリーダーノード  $L_n^l$  として他のノードへ通知を行う。

図8 新規リーダー選出手続き

度のノードを移動させる。

### 3.3.5 新規リーダー選出アルゴリズム

(Leader-election)

離脱や故障などにより、サブクラスタ  $SC_i^l$  のリーダーノード  $L_i^l$  が喪失した場合、図8に示すアルゴリズムにより  $SC_i^l$  内に残された他のノードが自律的に新規リーダーノード  $L_n^l$  を選出する。なお本稿では、 $L_n^l$  の選出に必要なクエリ数を、ネイバーノード数の半数とした。

### 3.4 動的クラスタ再構成法

前節で述べた手続きを用いて、VC環境上の性能管理手法として、動的クラスタ再構成法を提案する。提案する手法は、ワーカの参加・離脱、または、タスクの粒度の変化に対応するために、クラスタ間でノードを移動し、動的に再構成を行う。ここで、全参加ノード数を  $n$ 、クラスタ数を  $D$  として、手法の説明をする。また、簡略化のために、各タスクを構成するサブタスク数は一律に  $S$  とし、1個のサブタスクは  $R$  台のワーカで冗長に計算される冗長計算を仮定する。

#### (1) 階層型クラスタ管理構造への登録

計算プロジェクトに参加している全ノードを前節で説明した階層型クラスタ管理構造<sup>3)</sup>に登録する。

#### (2) クラスタ数 $D$ の決定

次式により、クラスタ数（同時に処理可能なタスク数）を決定する。

$$D = \left\lfloor \frac{n}{S \times R} \right\rfloor \quad (1)$$

#### (3) クラスタ再構成

定期的に、最も所属ノード数の多いクラスタ  $C_i$  から最も所属ノード数の少ないクラスタ  $C_j$  に対して Transfer アルゴリズムを実行する。Transfer アルゴリズムの引数である移動台数  $s$  は、以下の式によって決定する。

$$s = \min\left\{\frac{1}{2}(x - S \cdot R), 2(S \cdot R - y)\right\} \quad (2)$$

ただし、 $x, y$  はそれぞれクラスタ  $C_i, C_j$  の所属ノード数である。

### 3.5 新規参加ノードの所属クラスタの決定

VCに新規に参加するノード  $n_j$  は、3.3.1節で示した Join アルゴリズムにより、あるクラスタ内において距離的に近いサブクラスタに所属するが、まず最初に、所属するクラスタを決定する必要がある。ここで、各クラスタに Join アルゴリズムを適用し、最下位層のリーダーとの距離が最も近いクラスタを選択する方法が考えられる。しかし、クラスタ内の全ノードとの距離を考慮しているわけではないため、クラスタの物理的なネットワークポロジによっては、最下位層のリーダーとの距離が最短であっても、距離的に遠いノードが存在する可能性もある。本稿では、 $n_j$  との距離が最も近い最下位層のリーダー  $n_r$  と最も遠い最下位層のリーダー  $n_f$  の2つのノードを考慮することで、この問題に対処する。 $n_j$  が所属するクラスタを決定する際に、Join アルゴリズムで  $n_r$  を、Join アルゴリズムと同様のアルゴリズムで  $n_f$  を探索し、次式で定義される近接度が最も小さいクラスタに所属する方法をとる。

$$dis(n_j, n_r) + \alpha \cdot dis(n_j, n_f) \quad (3)$$

ただし、 $dis(n_j, n_r)$  はノード  $n_j$  と  $n_r$  間の距離であり、 $\alpha$  は遠いノードの存在をどの程度考慮するかを決定するパラメータである（本稿では  $\alpha = 1.0$  としている）。

$n_j$  が VC に参加する際、存在する全てのクラスタに対して近接度を計算すると、メッセージ交換によるコストが増大する。そこで本稿では、 $n_j$  が近接度を計算するクラスタの候補を全クラスタのうち所属台数の少ない下位  $\lfloor D \times H \rfloor$  個に制限する。ここで、 $H(0.0 \leq H \leq 1.0)$  はどの程度候補を限定させるかを決定するパラメータである。この方法により、新規参加時のメッセージ交換コストを減少させるだけでなく、参加ノードを処理性能（＝ノード台数）の小さなクラスタに優先的に所属させ、特定のクラスタが肥大化することを防ぐことが可能となる。すなわち、ノードの新規参加時に Transfer アルゴリズムを適用するのと似たような効果が得られる。

## 4. 評価と考察

### 4.1 シミュレーション実験の概要

本提案手法の有効性を検証するために、VCのシミュレーション実験を行い、計算ターン数による評価を行った。物理ネットワークとして、1000台、5000台規模の2種類の Transit-Stub 型のインターネットトポロジをトポロジジェネレータ GT-ITM<sup>5)</sup> によって作成した。

シミュレーション実験で用いたパラメータを表1に示す。本研究では全てのタスクは同一数の  $S$  個のサブタスクに分割されるとし、VCに参加するノードの処理性能は均一であると仮定する。

シミュレーションの流れを以下に示す。

- (1) トポロジ内に存在する全ノードのうち70%のノードをランダムに参加させ、3.4節で示したクラス

表 1 実験パラメータ

トポロジ内ノード数 $N$	1000, 5000
参加離脱率 $P$	0.005, 0.01, 0.02
タスク数 $T$	500, 1000, 1500
サブタスク数 $S$	10 ~ 50 ( $N = 1000$ ) 60 ~ 100 ( $N = 5000$ )
サブタスクの計算コスト	20 ターン
サブタスク間の通信コスト	所属クラスタ内の 他ノードとの距離 (ホップ数)の最大値
冗長度 $R$	2
サブクラスタ内ノード数の上限 $k$	6
サブクラスタ内ノード数の下限 $l$	2
Transfer 手続きの実行頻度 $TI$	20 or 実行せず
新規参加時の所属先クラスタ制限 $H$	0.5, 1.0

タ再構成法に従って一つの初期クラスタを構成する。なおこの際に、クラスタ数が  $D$  個となるまでクラスタの分割、もしくは統合を繰り返すが、この分割・統合操作を含めて 100 ターン経過させた状態をタスク配布前の初期状態とする。

- (2) 手順 (1) で構成した初期状態の各クラスタに対して個別のタスクを割り当て処理を開始する。
- (3) タスクの処理が完了したクラスタに対しては未割当のタスクを新たに配布し処理を継続させる。  $T$  個のタスク全てが処理完了した段階で計算プロジェクトの完了とし、シミュレーションを終了する。

なお、タスク配布の前後に関わらず、毎ターンごとにノードの新規参加と離脱をそれぞれ  $N \times P$  台ずつ発生させる。また、毎ターンごとに必要に応じて提案手法により動的にクラスタ内の構成を更新する。

#### 4.2 離脱の頻度と処理時間の関係

$N = 1000, 5000$  の場合における、参加離脱率と処理時間の関係をそれぞれ図 9, 図 10 に示す。なおサブタスク数はそれぞれ  $S = 20, 60$  としている。それぞれの図において、No-Clustering はクラスタリングを全く行わず、 $D$  個のタスクおよびタスク内のサブタスクをラウンドロビンで配布する手法であり、Clustering は提案手法を示している。提案手法において、 $H = 1.0, 0.5$  は、3.5 節で説明した所属先クラスタの候補数を制限するパラメータ  $H$  をそれぞれ 1.0, 0.5 に設定することを意味しており、Transfer は  $TI$  ターン毎に Transfer アルゴリズムを実行することを意味している。

これらの結果から、参加離脱の頻度が大きくなるにつれて計算プロジェクトにかかる処理時間が増大していることが分かる。参加離脱の頻度が小さい場合には、クラスタリングを行わない手法に比べて提案手法がより短い処理時間で計算を完了できることが確認された。しかしながら、提案手法において Transfer 手続きや新規参加時の所属先クラスタ候補の制限を行わない場合は、参加離脱率が 2.0% 程度になると処理効率の低下が見られた。これはノードの頻繁な離脱により、所属ノード数の少ないクラスタ内でのサブタスク計算において同期による待ち時間コストが増大し、著しい処理性能の低下が起こっ

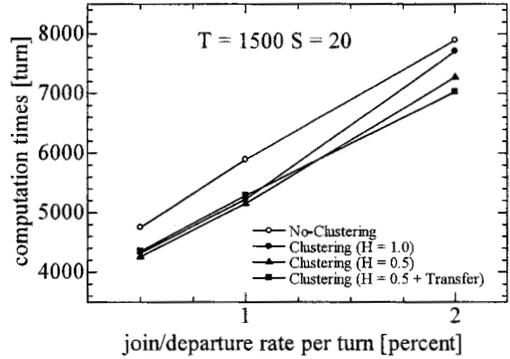


図 9 計算時間の変化 ( $N = 1000$ )

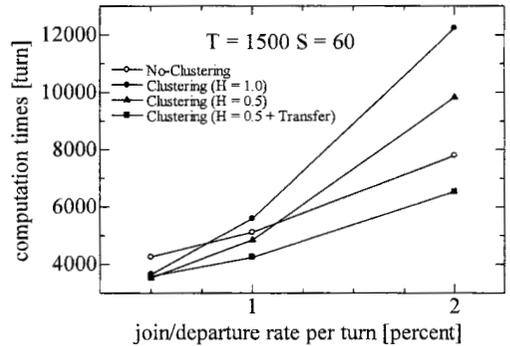


図 10 計算時間の変化 ( $N = 5000$ )

たものと考えられる。一方で、Transfer 手続きや新規参加時の所属先クラスタ候補の制限を行った場合は、離脱頻度が大きい場合にも、比較手法に比べて計算時間が短縮されていることが分かる。

#### 4.3 タスク数と処理時間の関係

タスク数の増加に対する処理時間の関係を図 11, および図 12 に示す。なお縦軸は提案手法とラウンドロビン手法における計算時間の比をプロットしている。この値が 1 より小さいほど、クラスタリングを行わないラウンドロビン手法に比べて提案手法の計算時間が短いことを表している。前項に示した結果と同様に、クラスタ間で処理性能を調整することにより、処理効率の向上が可能となる。また、タスク数が増大しても処理効率の向上を確認することができた。

#### 4.4 タスクの粒度と処理時間の関係

次に、 $N = 1000, 5000$  の場合において、サブタスク数、すなわちタスクの粒度と処理時間の関係を調べた。なお今回の実験では  $T = 1500, P = 2.0$  とした。得られた結果を図 13, 図 14 に示す。タスクの粒度を大きくすると、全参加ノードを多数のクラスタに分割することができないためクラスタサイズが大きくなり、通信コスト削減の効果が期待できない。図 13, 図 14 の結果からも、

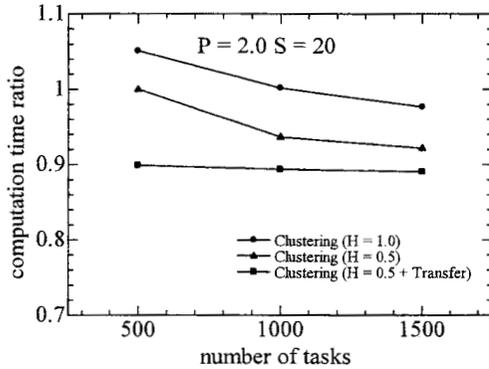


図 11 計算時間比 ( $N = 1000$ )

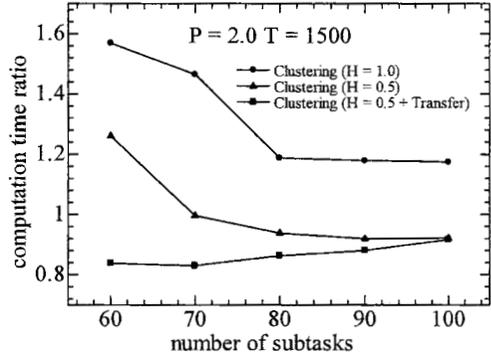


図 14 計算時間比 ( $N = 5000$ )

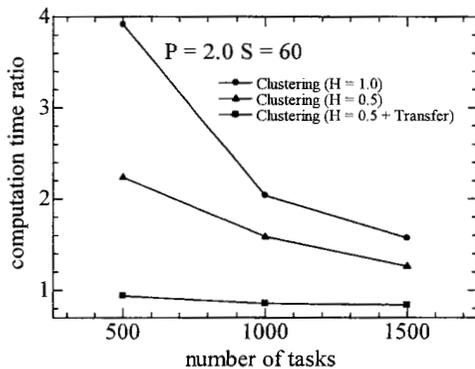


図 12 計算時間比 ( $N = 5000$ )

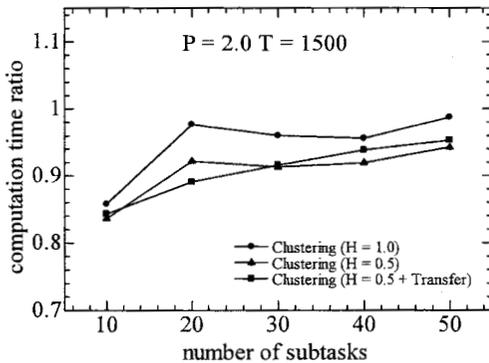


図 13 計算時間比 ( $N = 1000$ )

サブタスク数が多い場合に提案手法における計算時間がクラスタリングを行わない手法の値に近づいていることが分かる。また、これまでの実験結果と同様にクラスタ間の性能調整を行うことで、タスクの粒度によらず処理効率を向上させることができている。

## 5. まとめと今後の課題

本研究では、ノード同士に同期的な通信が発生するような大規模並列計算を VC で実行させることを目的に、

既存の並列計算モデルを拡張した。また、VC 環境において提供された計算機資源を管理するための構造として、階層型クラスタ管理構造を構築し、これを動的かつ自律分散的に再構成する手法を提案した。シミュレーションによる性能評価から、ノードの離脱頻度が大きく、局地的な処理性能の低下が起こるような場合でも、クラスタ間で処理性能の調整を行うことで、最大で約 16% の処理効率の向上が得られることを確認した。

今後の課題としては、ノードごとの性能の変動を考慮した場合の動的クラスタ再構成法の改良が挙げられる。また、計算時間を最小にするための冗長数の最適値の導出も挙げられる。

**謝辞** 本研究の一部は総務省戦略的情報通信研究開発推進制度 SCOPE 特定領域重点型研究開発助成 (061102002) によって行われた。関係各位に感謝いたします。

## 参考文献

- 1) SETI@home, <http://setiathome.ssl.berkeley.edu/>.
- 2) Folding@home, <http://folding.stanford.edu/>.
- 3) S. Banerjee, S. Parthasarathy, and B. Bhattacharjee. "A protocol for scalable application layer multicast", Technical Report CS-TR-4278, Department of Computer Science, University of Maryland College Park, USA, 2001.
- 4) S. Banerjee, C. Kommareddy, and B. Bhattacharjee. "Scalable peer finding on the internet", In proceedings of GLOBECOM'02, vol. 3, pp. 2205-2209, 2002.
- 5) E. W. Zegura, K. L. Calvert, and S. Bhattacharjee. "How to model an internetwork", In proceedings of IEEE Infocom, vol. 2, pp. 594-602, March 1996.
- 6) 上田達也, 安部広多, 石橋勇人, 松浦敏雄, "P2P 手法によるインターネットノードの階層的クラスタリング", 情報処理学会論文誌, Vol. 47, No.4, pp. 1063-1076, 2006.