

# Omni/SCASH における実行時性能評価に基づく 動的負荷分散拡張の実装と評価

栄 純明<sup>†1</sup> 松岡 聡<sup>†2</sup>  
佐藤 三久<sup>†3</sup> 原田 浩<sup>†4</sup>

コモディティクラスタ環境では様々な理由によりノード間に性能の heterogeneity が生じるケースが増加している。これはロードインバランスの原因の一つであり、何らかの動的負荷分散機能が必要である。本論文では Omni/SCASH に対して我々が現在行っている動的負荷分散機能拡張に関して報告する。動的負荷分散機能を用いることによって、アプリケーションプログラマが明示的にデータやタスクの配置を指定することなく、ランタイムシステムによって性能ヘテロなクラスタにおいてロードバランシングが行える。動的負荷分散拡張の一つであるループ再分割機能の評価の結果、ループ再分割機能自身の効果を確認した一方で、page migration 機能との協調動作が重要であることも確認した。

## Implementation and Evaluation of Dynamic Load Balancing Using Runtime Performance Monitoring on Omni/SCASH

YOSHIAKI SAKAE,<sup>†1</sup> SATOSHI MATSUOKA,<sup>†2</sup> MITSUHIISA SATO<sup>†3</sup>  
and HIROSHI HARADA<sup>†4</sup>

In the commodity cluster environment, there may be performance heterogeneity between nodes due to several reasons, from which the application programmer will suffer as load imbalances. To overcome these problems, some dynamic load balancing mechanisms are needed. In this paper, we report our ongoing work on dynamic load balancing extension to Omni/SCASH. Using our dynamic load balancing mechanisms, we expect that programmers can have load imbalances adjusted automatically by the runtime system without explicit definition of data and task placements in a commodity cluster environment with possibly heterogeneous performance nodes. The results of evaluation indicates that, our loop re-partitioning scheme which is one of the dynamic load balancing extension works well, and also it is important to combine loop re-partition with dynamic page migration.

### 1. はじめに

近年、HPC 向けの計算機の多くは SMP を多数結合したクラスタの形態を取る物が増加している<sup>9)</sup>。とりわけ汎用品を構成要素とするコモディティクラスタはコストパフォーマンスや管理の容易性などの面から、多くの研究機関や大学で使用されている。しかしながら、コモディティクラスタではプロセッサやネット

ワークテクノロジーの急速な進歩や費用面の問題からノードの段階的な追加や、プロセッサやメモリの段階的な増設が比較的頻繁に起こりノード間性能の heterogeneity が生じ、ロードインバランスの原因の一つとなっている。またノード間性能に heterogeneity が無いケースでもマルチユーザー環境がノード間の負荷バランスを崩し、結果的にロードインバランスにつながるケースもある。

環境およびアプリケーションごとにプログラマが明示的にロードバランシングを行うのは困難であり、下位のランタイムシステムによる自動負荷調整が必要とされている。この問題に対処するために、我々はコモディティクラスタ向けに OpenMP をプログラミングインターフェースとし、データの再配置および動的負荷分散を自動的に行うシステムを開発している。具体的には、性能的にヘテロな環境においての自動負荷分散機能として、実行時性能評価に基づくループ

<sup>†1</sup> 東京工業大学大学院情報理工学専攻 数理・計算科学専攻  
Tokyo Institute of Technology

<sup>†2</sup> 東京工業大学大学院情報理工学専攻 数理・計算科学専攻 / JST  
Tokyo Institute of Technology / JST

<sup>†3</sup> 筑波大学 電子・情報工学系 計算物理学計算センター  
Tsukuba University

<sup>†4</sup> ヒューレットパッカードジャパン  
Hewlett-Packard Japan, Ltd.

ブ再分割機能を SDSM SCASH 上の OpenMP の実装である Omni/SCASH<sup>7)</sup> に実装した。本論文では Omni/SCASH のループ再分割機能拡張の実装と評価および、現在行っている SCASH レベルの page fault counting に基づく動的ページマイグレーション機能の概要に関して報告する。これらの機能を用いることによって、プログラマが明示的にデータやタスクの配置を指定すること無く、ランタイムシステムによって自動的にロードバランシングが実現できると期待している。

## 2. 背景

### 2.1 Omni OpenMP コンパイラ

Omni OpenMP コンパイラは OpenMP プログラムを入力とし、ランタイムライブラリ呼び出しを含むマルチスレッド C プログラムを生成する translator である。フロントエンドとして C-front および F-front があり、それぞれ C および Fortran コードを入力とし、中間コードである Xobject を出力する。Xobject の表現形式はデータタイプ情報を持った AST (Abstract Syntax Tree) で、その各ノードはソースコードの各要素を表現する Java のオブジェクトであり、Exc Java tools により容易に変換可能である。Exc Java tools はプログラムの Xobject の表現レベルでの解析や変換を実現するクラスやメソッドを提供する Java のクラスライブラリ群である。OpenMP プログラムのマルチスレッド C プログラムへの変換はこの Exc Java tools を用いて実現されている。生成されたマルチスレッド C プログラムはネイティブバックエンドコンパイラによってコンパイルされ、ランタイムライブラリとリンクされる。

### 2.2 SCASH

SCASH<sup>3)</sup> は低レイテンシかつ広バンド幅な通信ライブラリである PM<sup>8)</sup> と、OS の提供するメモリプロテクションなどのメモリ管理機能を用いて実現されているページベースの Software Distributed Shared Memory System である。SCASH のメモリモデルはマルチプルライタプロトコルを用いた Eager Release Consistency (ERC) であり、ユーザーレベルのランタイムライブラリとして実装されている。

ページ単位の一貫性維持プロトコルとしては invalidate および update の両方が実装されており、実行時に選択可能となっている。

SCASH では home ノードでページの最新データと、ページを共有しているノードの集合を保持するページディレクトリを管理し、base ノードが最新の home ノードを常に記録している。すべてのノードは、全ページの base ノードを記録しており、最新の home が不明になったノードは base に home を問い合わせることで最新の home を得ることができる。

### 2.3 Translation of OpenMP programs to SCASH

OpenMP のプログラミングモデルでは、特に指定しない場合グローバル変数は共有される。その一方で SCASH ではグローバルスコープで宣言された変数はプロセッサプライベートに割り当てられ、共有領域は実行時にメモリアロケーションプリミティブを用いて明示的に割り当てる必要がある。OpenMP プログラムを SCASH の “shemem メモリモデル” にコンパイルするために、コンパイラはグローバル変数を実行時に共有アドレスに確保するように変換する必要がある。

OpenMP のディレクティブは、SCASH のプリミティブを用いて同期やプロセッサ間通信を行うランタイムファンクションに変換される。

OpenMP の parallel ディレクティブでアノテートされた逐次プログラムを fork-join モデルの並列プログラムに変換する際、コンパイラは各 parallel リージョンを独立したサブファンクションとしてくり出す。マスターノードはランタイムファンクションを用いてスレーブスレッドを呼び出し、このサブファンクションを並列実行させる。各ノードのすべてのスレッドはプログラムの起動時に作成され、スレーブノードで fork が行われるまで待機している。また、Omni/SCASH では nested parallelism はサポートしていない。

SCASH ではすべての共有メモリ領域の一貫性はバリア同期の際に取られる。これは OpenMP のメモリモデルと対応している。OpenMP のロックと同期操作は特定のオブジェクト (e.g. ロック変数) に対する SCASH の一貫性保持プリミティブを明示的に用いて実現されている。

## 3. ループ再分割による動的負荷分散

1 節で述べたようにロードインバランスの原因は様々である: (1) 対象アプリケーションが元々ロードインバランスを持つ場合 (2) マルチユーザー環境の影響でノード間の負荷に差がある場合 (3) アプリケーションが性能ヘテロなクラスタで実行される場合、などが考えられる (2)(3) のケースではロードインバランスを事前に決定できないため、静的な負荷分散技法は不十分であり、実行時の性能情報に基づく動的負荷分散機能が必須である。

本研究では、データ並列アプリケーションの実行時負荷分散機能として、実行時性能計測に基づくループ再分割の手法を用いる。OpenMP の *dynamic* や *guided* などのスケジューリングオプションはある程度の負荷分散を達成できるが、双方ともスケジューリングのたびに集中管理されている chunk queue へアクセスしアップデートする必要があるため同期操作を伴い、特にクラスタのような分散メモリ環境では大きなオーバーヘッドとなる。このオーバーヘッドを避け

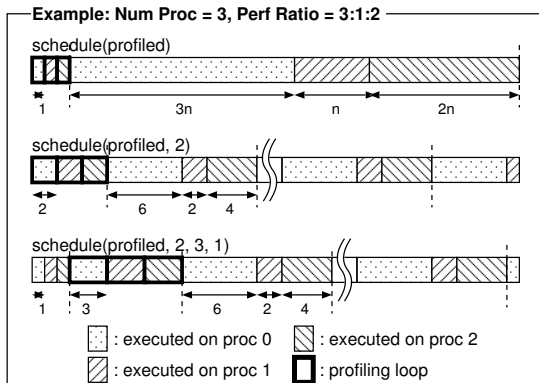


図 1 Example of profiled scheduling

るため、我々の提案する profiled スケジューリングではターゲットループの初期イテレーションの実行時性能を計測し、各プロセッサが計測した性能差に基づきプロセッサローカルに chunk サイズを調整することで負荷分散を達成する。

profiled スケジューリングは OpenMP の *schedule* clause で以下のように指定する。

```
schedule(profiled[, chunk_size[,  

    eval_size[, eval_skip]])
```

profiled スケジューリングが指定されると、まず eval\_skip で指定した回数のイテレーションを各スレッドで通常実行し、続く eval\_size で指定した回数のイテレーションの実行時性能計測を行う。その後、各スレッドは実行時性能計測の結果に基づき、残りのイテレーション空間を chunk\_size を基底値として性能比にしたがってサイクリック分割する。

ただし、eval\_skip, eval\_size, chunk\_size はこの順に省略可能であり、eval\_skip を省略もしくは 0 とした場合には先頭のイテレーションから実行時性能計測を行う。eval\_size を省略した場合には eval\_size = 1 として実行時性能計測を行う。chunk\_size を省略もしくは 0 とした場合には残りのイテレーション空間を実行時性能比に基づいてブロック分割する。

図 1 に profiled スケジューリングの例を示す。ここではプロセッサ数を 3 としその性能比を 3:1:2 とする。

アプリケーションプログラマが並列ループに対して profiled スケジューリングを指定すると、Omni はそのパラレルリージョンをサブファンクションとしてくりだし、パラレルリージョンの実行に参加しているスレッドによって並列実行される。それに加えて、Omni によりイテレーションの実行に要する時間を正確に計測するためループの前後に計測コードが挿入され、スレッドごとのターゲットループの実行性能の差に基づきループ再分割を行う (図 2)。実装では正確な性能計測のために、計測コードでは PAPI<sup>2)</sup> を用い

```
#pragma parallel omp for schedule(profiled)
for (i = 0; i < N; i++) {
    LOOP_BODY;
}

static void __ompc_func(void **__ompc_args){
    int i;
    {
        int lb, ub, step;
        double start_time = 0.0, stop_time = 0.0;
        lb = 0, ub = N, step = 1;
        _ompc_profiled_sched_init(lb, ub, step, 1);
        while (_ompc_profiled_sched_next(&lb, &ub, start_time, stop_time)) {
            _ompc_profiled_get_time(&start_time);
            for (i = lb; i < ub; i += step) {
                LOOP_BODY;
            }
            _ompc_profiled_get_time(&stop_time);
        }
    }
}
```

図 2 profiled スケジューリング時のコード変換

て、CPU がサポートするハードウェアリアルタイムカウンタを利用している。

各スレッドに割り当てられるサブループのインデックス、上限下限などの情報は各スレッドごとに管理され、\_ompc\_profiled\_sched\_init() で初期化される。各スレッドに割り当てられるサブループのイテレーション空間の計算は\_ompc\_profiled\_sched\_next() によって計算される。\_ompc\_profiled\_sched\_next() では\_ompc\_profiled\_get\_time() によって計測される実行時性能を元に、イテレーション空間をスレッド間で調整することでループ再分割を実現している。

図 3 に\_ompc\_profiled\_sched\_next() によるループ再分割アルゴリズムの概要を示す。

この操作で他のスレッドとの通信をとまなうのは速度データ交換の操作だけであり、スレッド間の性能差が、アプリケーション自体のロードインバランスによらず、ノード性能差や定常的な負荷によるものである場合には、初回のみ通信をとまなうループ再分割が行われ、それ以降は計算済みの割り当て比率によって、イテレーション割り当てが行われる。そのため、スケジューリングごとに中央管理された chunk queue にアクセスするために通信の発生する dynamic, guided などのスケジューリングに比べて、速度の向上が期待できる。

#### 4. 評価

profiled スケジューリングの評価に NAS Parallel Benchmarks (NPB-2.3) の EP, CG の OpenMP パージョンを利用した。オリジナルの Fortran パージョンから OpenMP への移植は RWCP によって行われたものである。

EP EP はモンテカルロ法などで使用される乱数の生成を行うカーネルベンチマークプログラムである。種となる乱数列に対して隣り合った要素のみに依存する計算が行われ、 $2n$  個の浮動小数点の疑似乱数が生成される。データ依存はごく少なく

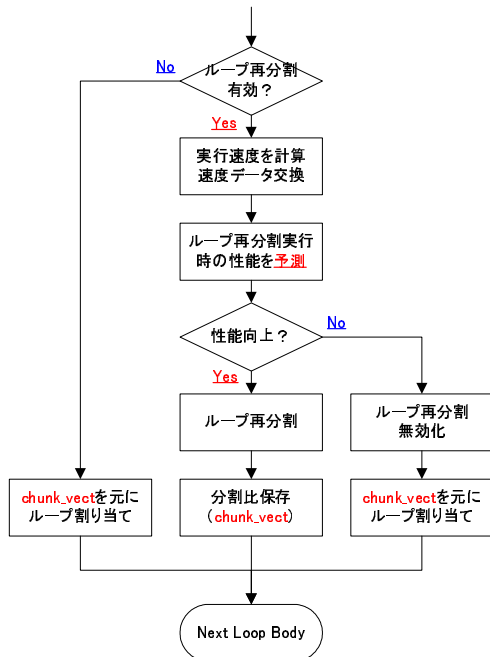


図 3 ループ再分割アルゴリズムの概要

表 1 評価環境: Performance Heterogeneous Cluster

	Fast nodes	Slow node
CPU	Pentium III 500MHz	Celeron 300MHz
Cache	512KB	128KB
Chipset	Intel 440BX	Same as "Fast"
Memory	SDRAM 512MB	Same as "Fast"
NIC	Myrinet M2M-PCI32C	Same as "Fast"

この問題は Embarassingly Parallel と呼ばれている。こういった特徴から、profiled スケジューリングのデータのローカリティに左右されない純粋な性能を評価するため EP を用いた。

**CG** CG は大規模で正値対称な粗行列の最小固有値の近似解を Conjugate Gradient (CG) 法を用いて求める問題である。行列積と不規則な通信を含む unstructured grid computation の典型であり、CG のカーネルループでは共有配列に対する頻繁なアクセスが生じる。そのためデータのローカリティが性能に大きな影響をおよぼす。その結果、profiled スケジューリングで動的負荷分散を行ったことによる性能向上が相殺される可能性がある。

#### 4.1 評価環境

評価環境を表 1 に示す。評価に使用した性能ヘテロなクラスタは CPU 性能のみヘテロな設定となっており Pentium III 500MHz のノードと Celeron 300MHz のノードから構成され、それぞれ Fast, Slow と表記することにする。

OS には RedHat 7.2 を、クラスタシステムソフト

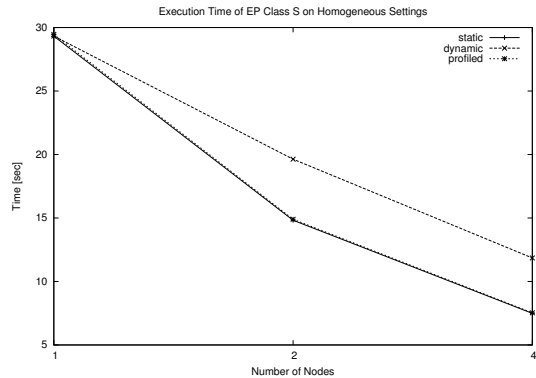


図 4 Execution Time of EP Class S on Homogeneous Settings (Pentium III nodes only)

ウェアには SCore をコンパイラには gcc-2.96 -O をそれぞれ用いている。

評価結果では profiled スケジューリングに関して、もっとも良い性能の得られた chunk\_size, eval\_size, eval\_skip いずれも指定しなかった場合の結果のみを示す。

#### 4.2 性能ホモな設定における結果

まずはじめに、profiled スケジューリング自身のオーバーヘッドを確認するために、データ配置が性能に影響を与えない EP を用いて性能的にホモな設定でベンチマークを行った。

図 4 から分かるように、dynamic スケジューリングでは chunk 割り当てごとに中央管理されたループインデックスにアクセスするためのリモートアクセスがオーバーヘッドとして現れており、性能低下が見られるのに対し、profiled スケジューリングではオーバーヘッドがほとんどなく static スケジューリングと同等の性能が得られている。

#### 4.3 性能ヘテロな設定における結果

図 5 に static, dynamic, profiled の各スケジューリングポリシーで EP class S を実行した際の性能を示す。

点線に四角でプロットされているグラフはノードとしてすべて "Fast" を利用する設定で static スケジューリングを用いた場合の EP の性能で、期待される性能の上限値を示している。実線にプラス記号でプロットされているグラフはヘテロな設定で static スケジューリングを用いた場合の EP の性能で、性能改善の基準値と見なすことができる。

EP のターゲットループのイテレーション空間は比較的小さく、その結果 chunk queue へのアクセス回数が少なくすむため、dynamic スケジューリングのオーバーヘッドは小さく static スケジューリングよりも高い性能を示している。

profiled スケジューリングでは、EP のターゲットループのループボディサイズが比較的大きく正確に実

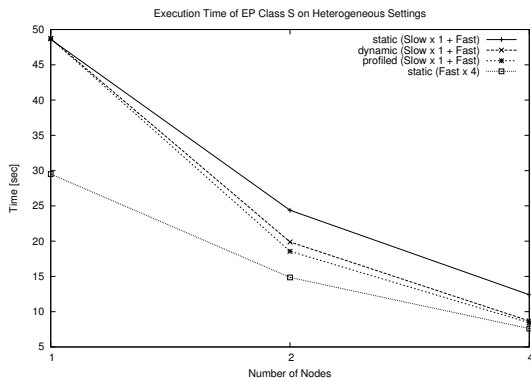


図 5 Execution Time of EP Class S on Heterogeneous Settings (one Celeron node + Pentium III nodes)

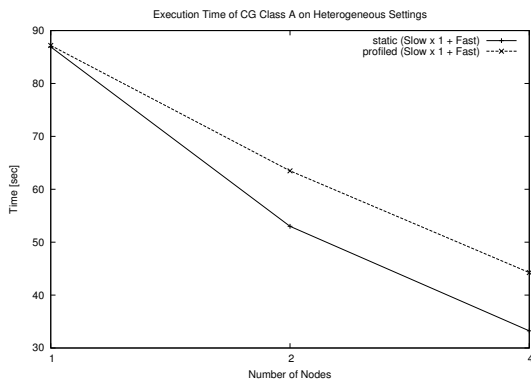


図 6 Execution Time of CG Class A on Heterogeneous Settings (one Celeron node + Pentium III nodes)

行時性能の計測が行えているため、ノード間性能比に沿って chunk\_vector が計算されていることを確認した。その結果、性能ヘテロな環境での EP は profiled スケジューリングでもっとも良い性能が得られている。

図 6 に CG class A を static および profiled スケジューリングを用いて実行した場合の性能を示す。

CG kernel ベンチマークは共有配列へのアクセス数が多いため、リモートアクセスのコストが高いクラスタのような分散メモリ環境では特に、タスクとデータを同じノードに割り当てるようにすることが性能の観点から重要である。さらに、たとえばデータ初期化フェーズと計算フェーズのように、“parallel for”リージョンごとにループ再分割が行われるため、タスクとデータの affinity が崩れる可能性がある。その結果、CG の場合には static スケジューリングに対して profiled スケジューリングの性能が低くなっている。この profiled スケジューリングがデータローカリティに対しておよぼす影響を次節で調査する。

#### 4.4 profiled スケジューリングがデータローカリティに与える影響

表 2 に “Slow node” 1 ノードと “Fast node” 3 ノード

表 2 Memory Behavior of the CG Class A on Celeron×1 + Pentium III×3 settings for static/profiled scheduling

	static	profiled
L2 miss ratio	29.6%	31.1%
Page Fault at SCASH	16456	27201
Barrier	5088	8006

の計 4 ノードを用いて、CG Class A を static および profiled スケジューリングで実行した場合の、L2 cache miss ratio、バリア同期、SCASH レベルの page fault の回数を示す。なお、L2 cache miss ratio は PAPI を用いて以下のように計算している：

$$\text{L2 cache miss ratio} = \frac{\text{PAPIL2\_TCM}}{\text{PAPIL1\_DCM}} \times 100$$

L2 cache miss ratio に関しては static、profiled とともに同程度の値であるが、SCASH レベルの page fault に関しては profiled スケジューリングの方が static スケジューリングより多数生じている。これは profiled スケジューリングによってループ再分割が行われるたびにタスクの割り当てが変更され、タスクとデータの affinity が崩れるためである。その結果、より多くのリモートメモリアccessを伴い、性能低下につながっている。

profiled スケジューリングにおけるバリア同期回数の増加は以下のような理由による：まず、profiled スケジューリングではループ再分割を決定するために実行時性能のデータをすべてのノードで交換している。また、現状ではリモートメモリ参照の際に page migration が行われていないため、ループ再分割によってデータのローカリティが変化した場合、性能の予測値がループ再分割後の実際の性能と一致しない可能性が高い。この不一致がループ再分割の不必要な繰り返しにつながり、結果として余計なバリア同期増加により性能の低下の原因となっている。

この結果、ループ再分割後にデータのローカリティを回復するために page migration を組み合わせることが必須であると言え、現在その作業を行っている。また、動的ループ再分割はデータのローカリティに、page migration は実行時性能にそれぞれ影響を与えるため、両手法を協調動作させる際に、現在行っているように実行時性能比に沿ったループ再分割を即座に行うのではなく、漸次的に適用することを考えている。これによって不必要なループ再分割の繰り返しを避け、すばやく stable なタスクおよびデータ配置に安定することが期待できる。

## 5. Ongoing Work: Page Fault Conting に基づく Page Migration

CG の評価結果からもわかるように、SDSM のような分散メモリ環境で高性能を達成するにはデータのローカリティを考慮することが特に重要である。しかし

ながら、本来共有メモリ環境を想定している OpenMP にはデータの配置をプログラマーが指定する手段はない。ハードウェアが提供する page reference counter を元に頻繁に参照するノードに、データをページ単位でマイグレーションする手法<sup>6)</sup> や静的なデータ配置手法<sup>1),4),5)</sup> が提案されているが、SDSM 環境ではページの参照回数を正確に数えることは非常に大きなオーバーヘッドをとともう、動的なループ再分割を行うため静的なデータ配置手法は不十分、などの理由によりいずれの手法も利用できない。

そこで、本研究では SDSM レベルの page fault つまり、メモリ参照がローカルで解決できなかった回数を数え、あるページに対するリモート参照のもっとも多いノードにそのページをマイグレーションする。直接的にすべてのページ参照をカウントする手法に対し、この方法ではローカル参照がカウントに入らないため、ページマイグレーションによってかえってリモート参照が増大してしまうという可能性を避けられない。しかし、我々は SPMD スタイルのアプリケーションをターゲットとしており、並列ループにおけるメモリアクセスパターンがイテレーション間で変化しないという前提においては、数回のページマイグレーションで最適なページ配置に収束することが期待できる。

## 6. おわりに

本論文では我々が Software Distributed Shared Memory, SCASH 上の OpenMP の実装である Omni/SCASH に対して現在行っている、動的負荷分散拡張に関して報告した。我々の目標は、ノード間性能の heterogeneity やマルチユーザ環境などに起因するアプリケーションのロードインバランスを半自動的に解決する手法の実現である。このような問題に対して静的なアプローチは不十分であり、我々はターゲットループの実行時性能セルフプロファイリングに基づくループ再分割を実装し、またターゲットループにおける page fault counting に基づくページマイグレーションを用いた動的負荷分散技法の開発を現在行っている。これらの手法を用いることにより、ユーザプログラムが明示的にデータとタスクの配置を記述することなく、ラインタイムシステムによって最適なロードバランシングが行える。

評価の結果、データローカリティが性能に影響を与えない EP ベンチマークでは、実行時性能に基づきループ再分割を行う profiled スケジューリングは、static スケジューリングと比較してもオーバーヘッドなく、性能へテロな環境で十分な負荷分散が行え、もっとも良い性能を示した。しかしながら、ループ再分割によってデータのローカリティが損なわれるため、データのローカリティが性能に強く影響する CG ベンチマークでは SDSM のリモートメモリ参照の増加、および

ループ再分割を繰り返してしまうことによるバリア同期の増加により、性能低下が見られた。これは、データのローカリティ回復のため profiled スケジューリングを動的ページマイグレーション機能と協調動作させる必要があることを示しており、ハードウェアサポートなしに SDSM レベルの page fault counting を行うことで一定の精度で軽量なページマイグレーションを実現する方法を提案した。現在は、提案するページマイグレーション機能を Omni/SCASH に実装中であるが、早急に profiled スケジューリングとページマイグレーションの協調動作の効果を確認する評価を行う必要がある。

## 参考文献

- 1) J. Bircsak, P. Craig, R. Crowell, J. Harris, C. A. Nelson, and C. D. Offner. Extending OpenMP for NUMA Machines: The Language. In *Proceedings of Workshop on OpenMP Applications and Tool (WOMPAT'2000)*, July 2000. San Diego, USA.
- 2) S. Browne, J. Dongarra, N. Garner, G. Ho, and P. Mucci. A Portable Programming Interface for Performance Evaluation on Modern Processors. *The International Journal of High Performance Computing Applications*, 14(3):189–204, Fall 2000.
- 3) H. Harada, H. Tezuka, A. Hori, S. Sumimoto, T. Takahashi, and Y. Ishikawa. SCASH: Software DSM using High Performance Network on Commodity Hardware and Software. In *Proceedings of Eighth Workshop on Scalable Shared-memory Multiprocessors*, pages 26–27. ACM, May 1999.
- 4) A. Hasegawa, M. Sato, Y. Ishikawa, and H. Harada. Optimization and Performance Evaluation of NPB on Omni OpenMP Compiler for SCASH, Software Distributed Memory System (in Japanese). In *IPSI SIG Notes*, 2001-ARC-142, 2001-HPC-85, pages 181–186, Mar. 2001.
- 5) J. Merlin. Distributed OpenMP: Extensions to OpenMP for SMP Clusters. In *Invited Talk. Second European Workshop on OpenMP (EWOMP'00)*, Oct. 2000, Edinburgh, Scotland.
- 6) D. S. Nikolopoulos, T. S. Papatheodorou, C. D. Polychronopoulos, J. Labarta, and E. Ayguadé. Is Data Distribution Necessary in OpenMP? In *Proc. of Supercomputing 2000*, Nov. 2000. Dallas, TX.
- 7) M. Sato, H. Harada, and Y. Ishikawa. OpenMP compiler for Software Distributed Shared Memory System SCASH. In *Proceedings of Workshop on OpenMP Applications and Tool (WOMPAT'2000)*, July 2000. San Diego, USA.
- 8) H. Tezuka, A. Hori, Y. Ishikawa, and M. Sato. PM: An Operating System Coordinated High Performance Communication Library. In P. Sloot and B. Hertzberger, editors, *High-Performance Computing and Networking '97*, volume 1225, pages 708–717. Lecture Notes in Computer Science, Apr. 1997.
- 9) <http://www.top500.org/>.