

# アドレス空間の大きさに制限されない スレッド移動を実現する PGAS 処理系

原 健太郎<sup>†1</sup> 中 島 潤<sup>†1</sup> 田 浦 健次朗<sup>†1</sup>

利用可能な計算資源が動的に増減しうるクラウド環境で並列計算を実行させるためには、そのとき利用可能な計算資源に対応して並列計算の規模を動的に増減させる必要がある。しかし、有限要素法などの高性能並列科学技術計算のプログラムを、計算規模を動的に拡張/縮小できるように記述するのは難しい。そこで本研究では、「プログラムは十分な数のスレッドを生成するだけで良く、あとは処理系が透過的にそれら大量のスレッドをその時点で利用可能な計算資源に動的にマッピングする」モデルに基づく PGAS 処理系として DMI を提案する。特に、その要素技術として、アドレス空間の大きさに制限されることなく計算資源間でスレッド移動を実現するためのアドレス管理手法として random-address を提案し、その最適性を証明する。評価の結果、DMI を使うことで、既存の SPMD プログラムをわずかに変更するだけで計算規模を動的に拡張/縮小可能な並列プログラムを記述できることを確認した。また、利用可能な計算資源の増減に追隨して、実用的な有限要素法アプリケーションの並列度を効果的に増減できることを確認した。

## A PGAS Framework Achieving Thread Migration Unrestricted by the Address Space Size

KENTARO HARA,<sup>†1</sup> JUN NAKASHIMA<sup>†1</sup>  
and KENJIRO TAURA<sup>†1</sup>

In order to execute a parallel computation on a cloud environment in which the number of available resources can change dynamically, the scale of the parallel computation should expand or shrink dynamically in response to the dynamic increase or decrease of available resources. It is, however, difficult to describe most large-scale parallel scientific computings such as finite element methods so that they can scale up or down dynamically. Therefore, this paper proposes a PGAS framework named *DMI*, with which a programmer only has to create a sufficient number of threads because the framework transparently and dynamically schedules these threads on the available resources at the time. In particular, as an elemental technique for DMI, this paper pro-

poses *random-address*, which is a novel address management algorithm for live thread migration unrestricted by the address space size. This paper also gives the proof of the optimality of the random-address algorithm. We confirmed that in DMI we can easily develop parallel programs supporting dynamic scale-up and scale-down by slightly modifying existing SPMD programs. Furthermore, our evaluation showed that DMI can effectively change the parallelism of real-world scientific computings such as finite element methods in response to the dynamic increase or decrease of available resources.

### 1. 序 論

#### 1.1 背景と目的

有限要素法による応力解析や地震シミュレーション、粒子法による流体解析や分子シミュレーションなど、計算に長時間を要するような大規模な高性能並列科学技術計算への要請が高まっている。また、これらの大規模な高性能並列科学技術計算を実行するための計算基盤として、HPC クラウドに代表されるようなクラウドコンピューティング（以下、クラウド）環境<sup>(12), (42), (43)</sup> が広く注目されている。

クラウドでは、クラウドプロバイダと呼ばれる組織が大規模なデータセンタを構築し、インフラストラクチャ、プラットフォーム、ソフトウェアなどを整備して、それらをサービスとして利用者に提供する。そして利用者は、それらのサービスを従量制課金のもとで必要なときに必要な量だけ利用できる。このようなクラウドにおいては、煩雑なサーバ管理技術などが不要なうえ、急激な負荷変動にも柔軟に対応しやすいため、自前でデータセンタを管理するよりもコストパフォーマンスが優れる場合が多い。代表的なクラウドサービスとしては、仮想サーバやストレージなどの計算機資源を提供する IaaS (Infrastructure as a Service) としての Amazon EC2 や Amazon S3<sup>(1)</sup>、利用者が作成したアプリケーションの実行環境を提供する PaaS (Platform as a Service) としての Google App Engine<sup>(2)</sup> や Windows Azure<sup>(6)</sup>、エンドユーザのためのソフトウェアサービスを提供する SaaS (Software as a Service) としての Google Docs<sup>(3)</sup> や Salesforce.com の CRM<sup>(5)</sup> などがある。

このように、クラウドは多種多様なアプリケーション領域に対して効率的な実行環境を提供しているが、次節で考察するように、利用可能な計算資源が動的に増減しうるクラウド環

<sup>†1</sup> 東京大学大学院情報理工学系研究科

School of Information Science and Technology, The University of Tokyo

境において、長時間を要する大規模な高性能並列科学技術計算を効率的に実行できるような並列プログラムを記述するのは容易ではない。そこで本研究では、クラウド環境上での大規模な高性能並列科学技術計算を簡単に記述できるプログラミングモデルを提案したうえで、その処理系として DMI (Distributed Memory Interface) を提案して実装し、評価する。特に、その主要な要素技術としてスレッド移動について検討し、アドレス空間の大きさに制限されないスレッド移動を実現するアドレス管理のアルゴリズムとして random-address を提案する。

## 1.2 クラウド環境における並列計算

### 1.2.1 クラウドの特徴

第一に、クラウドの特徴を分析する。

クラウドでは、多数の計算資源を多数の利用者で利用することになる。よって、クラウドサービスを効率的に運用するためには、何らかの形で、有限個の計算資源が利用者間でスケジューリングされることになる。具体例で説明する。あるクラウドの中に利用者 A と利用者 B がいるとし、いま利用者 A の負荷が増大したとする。このときこのクラウドを利用している利用者が他にいなければ、利用者 A の計算規模を拡張することで負荷分散が図られる。やがて、別の利用者 B の負荷が利用者 A の負荷よりも増大したとすると、今度は利用者 A の計算規模を縮小するかわりに利用者 B の計算規模を拡張することで全体としての負荷分散が図られる、というような計算資源のスケジューリングがたとえば起きる。当然、どのような条件が成立したときに各利用者の計算規模がどのように拡張/縮小されるかは、各クラウドサービスの課金体系、対象とするアプリケーション、SLA (Service Level Agreement) などに依ってさまざまである。しかし、いずれにせよ、クラウドの本質は、全体の負荷状況に応じて各利用者の計算規模を拡張/縮小することで、有限の計算資源を多数の利用者間で効率的にスケジューリングするという点にある。このように、クラウド環境においては各利用者が利用可能な計算資源が動的に増減する。

このようなクラウド環境の代表例としては、Amazon EC2 Spot<sup>1)</sup> がある。Amazon EC2 Spot は、指示した数の仮想マシンを確実に利用することが可能な Amazon EC2 よりも、SLA の品質を落とすかわりに、より安価に計算資源を提供することを意図して作られたクラウドサービスである。Amazon EC2 Spot では計算資源を競り落として利用する。すなわち、需要と供給に応じて計算資源の時価が動的に変動しており、利用者が決めた入札価格がその時点での時価を上回っていれば、その時価で仮想マシンを利用することができる。たとえば、時価が \$0.03/時間のときに \$0.05/時間で入札すれば仮想マシンが起動し、やがて時価

が \$0.06/時間に上がったとすれば仮想マシンの電源が落とされる。よって、不要不急の計算であれば、時価よりも安い入札価格を設定することで「いつか実行されれば良いから安く実行」することができるし、急を要する計算であれば、時価よりも高い入札価格を設定することで「高くても良いからすぐに実行」することもできる。

### 1.2.2 要請される並列分散プログラミングモデル

第二に、このようなクラウド環境における並列計算に要請される並列分散プログラミングモデルについて考える。

まず、利用可能な計算資源が動的に増減しうるクラウド環境において長時間を要する並列計算はどう実行されるべきかを考えると、当然、その並列計算は、そのとき利用可能な計算資源に対応して計算規模を動的に拡張/縮小しながら実行される必要がある。たとえば、1000 ノードが利用可能なときにはその 1000 ノードを利用して実行され、やがて 10 ノードしか利用できなくなればその 10 ノードだけを利用して実行され、しばらくして 100 ノードが利用可能になればその 100 ノードを利用して実行されるよう、計算規模を動的に拡張/縮小しながら実行される必要があると言える。しかしながら、利用可能な計算資源の増減に対応して、計算規模を動的に拡張/縮小できるように並列科学技術計算を記述するのは明らかに難しい。言い換えると、外的な要因に対応して並列度を動的に増減させながら実行を継続できるような並列プログラムを記述するのは明らかに困難である。

以上の考察より、長時間を要する大規模な高性能並列科学技術計算をクラウド環境で実行するためには、計算規模を動的に拡張/縮小できるような並列プログラムを簡単に記述できる並列分散プログラミングモデルが必須であると言える。そこで本研究では、以下の並列分散プログラミングモデルを提案する：

- プログラマは、計算規模の拡張/縮小をいっさい考えることなく、単にアプリケーションの並列性をプログラムに記述するだけで良い。
- あとは処理系が透過的に、それらの並列性を物理的な計算資源にマッピングすることで、そのとき利用可能な計算資源に対応してプログラムの計算規模を動的に拡張/縮小してくれる。

要するに、プログラマは十分な並列性さえ並列プログラムに記述すれば良く、計算規模の動的な拡張/縮小は処理系が透過的に処理してくれるという並列分散プログラミングモデルである。この並列分散プログラミングモデルのもとでは、プログラマは並列度の変化を意識する必要がないため、高性能並列科学技術計算にとって代表的な SPMD 型で並列プログラムを記述することも可能である。

### 1.2.3 要請される並列分散プログラミング処理系

第三に、前述の並列分散プログラミングモデルを実現する並列分散プログラミング処理系について考える。

並列分散プログラミング処理系を設計するうえで重要となるのはデータ通信モデルである。一般に、並列分散プログラミング処理系のデータ通信モデルとしては、大きく分類して、メッセージパッシングモデルと共有メモリモデルが存在する。

メッセージパッシングモデルは、系内の各ノードに対して一意なランク（名前）が与えられ、ユーザプログラムではランクを用いたデータの送受信を send/recv 操作として明示的に記述する。メッセージパッシングモデルの利点は（1）ユーザプログラムにおける記述（send/recv）が下層ハードウェアで実際に発生する操作（send/recv）にそのまま対応するため、ユーザプログラムにとって本質的に必要とされる通信以外は発生せず通信に無駄が生じない点（2）データの所在管理や通信形態の決定に関してユーザプログラム側に自由度があるため、明示的なチューニングが行いやすく性能を引き出しやすい点などである。一方で、メッセージパッシングモデルの欠点は、ユーザプログラム側でデータの所在や通信形態を管理しなければならないため、動的で不規則なデータ構造を取り扱うような非定型な処理は非常に記述しづらい点などである。まとめると、メッセージパッシングモデルは、性能を引き出しやすい一方でプログラマにかかる負担が大きいデータ通信モデルと言える。メッセージパッシングモデルを採用する代表的な処理系には MPI<sup>(25), (31), (38), (44), (52), (53), (56)</sup> がある。

一方、共有メモリモデルは、物理的には分散した計算資源上に処理系が仮想的な共有メモリを構築することによって、ユーザプログラム側からは、あたかも通常の共有メモリ環境と同様の read/write 操作によってデータ通信を実現することができる。共有メモリモデルの利点は（1）共有メモリ環境上の並列プログラムと同様の read/write ベースの記述が可能のため、プログラマは各ノード上のデータ配置や煩雑なメッセージ通信を意識することなく並列アルゴリズムの開発に専念できる点、（2）既存のマルチコア並列プログラムとの概念的な乖離が小さい点などである。一方で、共有メモリモデルの欠点は（1）ユーザプログラム側の操作（read/write）と下層ハードウェアで実際に発生する動作（send/recv）が対応しないために、ユーザプログラムから処理系の挙動が把握しづらく明示的なチューニングを施しにくい点、（2）ユーザプログラムにとって本来必要な通信パターンが何なのかを処理系側で把握しづらいために、無駄なメッセージ通信が多量に発生する可能性が高い点などである。まとめると、共有メモリモデルは、プログラミングが容易な一方で性能を引き出しにくいデータ通信モデルと言える。共有メモリモデルを採用する代表的な処

プログラマ：並列性だけを記述すれば良い

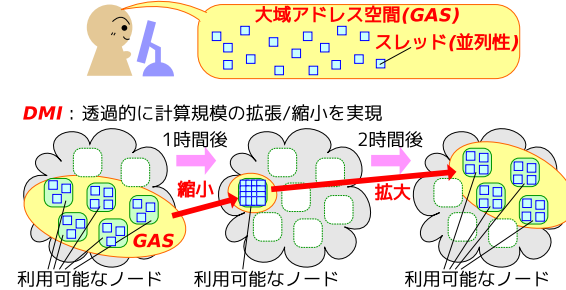


図1 DMIのコンセプト。

理系には、TreadMarks<sup>(8)</sup>、DSM-Threads<sup>(46), (47), (51)</sup>、SMS<sup>(68)</sup>、CRL<sup>(41)</sup>、UPC<sup>(11), (14), (18), (21)</sup>、Titanium<sup>(23), (30), (54), (55), (59), (60)</sup>、Co-Array Fortran<sup>(20), (21), (49), (58)</sup>、Global Arrays<sup>(27), (48)</sup> などがある。

以上の特徴をふまえて、クラウド環境における並列計算を簡単に記述させることを目標とする本研究では、データ通信モデルとしては、プログラミングが容易な共有メモリモデルを採用する。共有メモリモデルをさらに細かく分類すると、分散共有メモリモデル<sup>(8), (41), (46), (47), (51), (68)</sup>、ローカルビュー型のPGAS (Partitioned Global Address Space) モデル<sup>(20), (21), (23), (30), (49), (54), (55), (58)–(60)</sup>、グローバルビュー型のPGASモデル<sup>(11), (14), (18), (21), (27), (48)</sup> などさまざまなものがあるが、本研究では、read/writeに基づくプログラミングの容易さと性能をバランスさせるため、グローバルビュー型のPGASモデルを採用する。グローバルビュー型のPGASモデルでは、大域アドレス空間に対するread/write操作に基づく簡単なプログラミングが行えるうえ、ユーザプログラムからデータの分散配置を明示的に指示できるようになっており、リモートとローカルを明示的に区別できるため、性能の明示的なチューニングが行いやすい。グローバルビュー型のPGASモデルに基づく処理系としては、UPC、Global Arraysなどが代表的である。

以上の考察に基づき、本研究では、1.2.2節で提案した並列分散プログラミングモデルをグローバルビュー型のPGASモデルで実現する並列分散プログラミング処理系として、DMI (Distributed Memory Interface) を提案して実装し、評価する。DMIのコンセプトは以下のとおりである（図1）:

- プログラマは、計算規模の拡張/縮小を考えるとなく十分な数のスレッドを生成するように並列プログラムを記述するだけで良い。

- あとは処理系が透過的に、それら大量のスレッドをそのとき利用可能な計算資源にスケジューリングすることで、利用可能な計算資源の増減に対応して計算規模を動的に拡張/縮小してくれる。

このような DMI を実現するためには、以下の 3 つの要素技術が必須である：

- (1) 高性能な並列科学技術計算をサポートするためには、大域アドレス空間に対するアクセスが高性能に行えなければならない。よって、大域アドレス空間に対するアクセスを高性能化する技術。
- (2) 計算規模を拡張/縮小するためには、大域アドレス空間のコンシステンシがノードの動的な参加/脱退を越えて正しく維持されなければならない。よって、ノードの動的な参加/脱退を越えて大域アドレス空間のコンシステンシを維持する技術。
- (3) 大量のスレッドをそのとき利用可能な計算資源にスケジューリングするには、生きたスレッドをノード間で移動しなければならない。よって、生きたスレッドを安全に移動する技術。

このうち (1) と (2) に関しては著者らの研究<sup>29),65),66)</sup> を参照されたい。本稿では (3) のスレッド移動の要素技術について述べる。特に、2.3 節で説明するように、既存のスレッド移動技術<sup>9),10),38),45)</sup> では、生成できるスレッドの本数や各スレッドが使用できるメモリ量がアドレス空間の大きさに制限されるのに対して、本研究では、それらがアドレス空間の大きさに制限されないアドレス管理のアルゴリズムとして、random-address を提案する。

なお、本研究で提案するスレッド移動の手法は、ホモジニアスな環境を前提にする。実際のクラウド環境ではヘテロジニアスな環境が多いと思われるが、その場合には、仮想マシンを利用してホモジニアスな環境を作り出すことで DMI を利用できるようになる。

### 1.3 本稿の構成

第 2 章では、関連研究について述べ、特に既存のスレッド移動技術の問題点について言及する。第 3 章では、DMI のシステムを概観する。第 4 章では、クラウド環境における並列計算を簡単に記述できる DMI のプログラミングインタフェースについて説明する。第 5 章では、スレッド移動に伴うプログラミング制約について説明する。第 6 章では、random-address のアルゴリズムを提案し、random-address に基づいたアドレス管理を説明する。第 7 章では、DMI におけるスレッド移動の実装について説明する。第 8 章では、シミュレーションによる random-address の評価とアプリケーションベンチマークを行う。第 9 章では、結論と今後の課題を述べる。第 A.1 章では、random-address のアルゴリズムの最適性を証明する。

## 2. 関連研究

### 2.1 クラウドサービスにおける計算規模の拡張/縮小

まず、既存のクラウドサービスにおける計算規模の拡張/縮小について観察する。計算規模を拡張/縮小する際の粒度としては、粒度が大きい方から順に、仮想マシン、プロセス、スレッドが考えられる<sup>16),17)</sup>。

第一に、仮想マシンを粒度とする場合には、仮想マシンを起動/停止することで計算規模の拡張/縮小が実現される。これを実現するクラウドサービスとしては、Amazon EC2、Windows Azure などがあり、利用者は必要なときに必要な量だけ仮想マシンを利用することができる。これらのクラウドサービスの利点は、利用者に対して仮想マシンという汎用的な計算環境が提供されるため、利用者にとっての自由度が大きく、実行可能なアプリケーション領域が広いという点である。一方で、第一の欠点は、課金の粒度が CPU の使用時間などではなく仮想マシンの起動時間に基づいて行われてしまう点である。これは、仮想マシンは存在しているだけで無視できない量のメモリ資源を消費することに起因している。第二の欠点は、仮想マシンは 1 つの OS として動作するためメモリ消費量が多く、起動/停止には数分を要するので、計算規模の拡張/縮小の要求に対する応答性が悪い点である。たとえば、これらのクラウドサービスが仮想マシンのライブマイグレーション<sup>19),50)</sup> の技術をサポートすることにより、仮想マシンの移動時間を削減して応答性を改善することは可能かもしれないが、いずれにせよ、仮想マシンが使用しているメモリ全体を移動させる必要があることには変わりない。本研究が対象とするような並列科学技術計算にとっては、その並列科学技術計算が使用しているメモリだけが移動されれば十分な場合が多く、OS などを含めた仮想マシンの実行環境すべてが移動されることは要求されておらず、仮想マシン粒度での計算規模の拡張/縮小は不必要に重すぎる場合が多い。

第二に、プロセスを粒度とする場合には、プロセスを生成/破棄することで計算規模の拡張/縮小が実現される。これを実現するクラウドサービスとしては、Google App Engine などが代表的である。Google App Engine では、利用者が Java もしくは Python で記述した Web アプリケーションを登録しておく、その Web アプリケーションに対するクライアントからのリクエスト数の増減に応じて、計算規模が透過的に拡張/縮小され、利用者が何の意識を払わずとも負荷分散が図られる。Google App Engine の第一の利点は、計算規模の拡張/縮小の要求に対する応答性の良さである。たとえば、2010 年 7 月時点における無料コースでは、1 分間に最大 7400 個ものリクエストが処理可能とされている。この応答性

の良さは、プロセスが消費するメモリ量は仮想マシンより少なく、生成/破棄などの取扱いを高速に実現できることに起因している。第二の利点は、実際の CPU 使用時間に基づいた細粒度な課金を行える点である。これも、プロセスのメモリ消費量が少なく、取扱いが軽量であることに起因している。一方で、第一の欠点は、Google App Engine は Web アプリケーションに特化した作りになっており、高性能並列科学技術計算のようにプロセスどうしが密に結合して動作するアプリケーションには向いていない点である。たとえば、Google App Engine では、プロセス間のデータ共有は BigTable<sup>15)</sup> というデータベース経由で行われるが、このようなデータベースを利用して、高性能並列科学技術計算に要求されるようなプロセスどうしの密で複雑なデータ共有を効率的に実現することは難しいと考えられる。第二の欠点は、短時間で終了するアプリケーションしか実行できない点である。Google App Engine では主として Web アプリケーションを想定しているため、各プロセスは 30 秒以内に処理を終了させる必要がある。しかし、有限要素法などの並列科学技術計算を各計算部分が短時間で終了するように分割して記述することは困難である。

以上のように、既存のクラウドサービスには、並列科学技術計算に関して、ユーザプログラムから透過的かつ効率的に計算規模の拡張/縮小を行うには不十分な点が多いと言える。

## 2.2 並列科学技術計算におけるプロセス/スレッド移動

前節で指摘したように、長時間を要するような並列科学技術計算に対して計算規模の拡張/縮小を実現するには、プロセス/スレッドの粒度が適している。そこで本節では、クラウドサービスへの応用性は指摘されていないものの、それらに自然に応用可能だと思われる要素技術として、並列科学技術計算におけるプロセス/スレッド移動に関する既存研究について観察する。

まず、BLCR<sup>26)</sup> は (ある 1 つのノード上の) Linux のプロセスをチェックポイント/リスタートするためのカーネルモジュールである。BLCR を用いて、あるノードでチェックポイントしたプロセスを別のノードでリスタートさせることで、ノード間のプロセス移動を実現できる。ただし、プロセス移動を通じて IP アドレスなどの情報も透過的に移動させることは難しいため、BLCR は TCP/UDP ソケットのチェックポイント/リスタートに対応していない。そこで、研究<sup>52)</sup> では、MPI プロセスのチェックポイント時に on the fly な MPI メッセージをすべて回収し、リスタート時にそれらを復旧させる機能を BLCR に対して付け加えることで、MPI プロセスのプロセス移動を実現している。また、MPI-Mitten<sup>25)</sup> では、プロセス移動を越えて MPI における集合通信をサポートするための分散アルゴリズムが提案されている。さらに、Tern<sup>38)</sup> や Adaptive MPI<sup>31)</sup> では、MPI の各インスタンスを

プロセスではなくスレッドとして実装することで、プロセス移動よりも細粒度なスレッド移動を実現している。

このような MPI におけるプロセス/スレッド移動の要素技術は、さまざまな分野に応用されている。第一の応用は耐故障な並列計算の実現である。MPI プロセスのチェックポイント/リスタートを行うことで、MPI を用いた耐故障な並列計算を実現できる。さらに、研究<sup>56)</sup> では、故障しそうなノード上の MPI プロセスを、故障が起きる前に健康なノードにプロセス移動させることによって、チェックポイントの回数を削減している。第二の応用は動的な環境における並列計算の負荷分散である。クラウド環境にかぎらず、多数のユーザが共同利用する時分割方式のクラスタ環境などでは、利用可能な計算資源やその性能が動的に変化するため、その動的な変化に追従して並列計算を適応させることが重要になる。MPI Process Swapping<sup>53)</sup> では、 $n$  個のプロセッサを利用する MPI アプリケーションに対してあらかじめ  $n + m$  個のプロセッサを割り当てておき、計算資源の性能の変化に追従して MPI プロセスを移動させることにより、各時点でもっとも高性能な  $n$  個のプロセッサを選択して実行することを提案している。また、Adaptive MPI<sup>31)</sup> では、本研究と同様に、このように動的に計算資源の状況が変化する環境における並列計算をサポートするためには、「プログラマには十分な並列性のみを記述させることにし、あとは処理系が透過的にそれらの並列性を物理的な計算資源にマッピングする」プログラミングモデルが必要であることを指摘し、MPI のスレッド移動に基づいて負荷分散を行う処理系を実現している。

しかし、以上で述べたようなプロセス/スレッド移動およびその各種応用は、すべてメッセージパッシングモデルに基づくものである。著者らの知るかぎり、DMI のように、共有メモリモデルに基づいて、ユーザプログラムから透過的に計算規模の拡張/縮小を実現した研究は存在しない。

## 2.3 既存のスレッド移動と問題点

最後に、スレッド移動の要素技術について既存研究を観察する。

スレッド移動<sup>9),10),16),22),24),32)–37),39),45),57),61)–63)</sup> とは、あるプロセス内で実行しているスレッドを停止させ、そのスレッドのメモリを (特に別のノードの) 別のプロセスに移動させてから実行を復帰させることである。このとき、各スレッドのスタック領域やヒープ領域などのメモリ領域をプロセス間で移動させることになるが、これらのメモリ領域にはそのメモリ領域自身へのポインタが含まれている可能性がある。よって、移動先プロセスにおいて、単純に適当なアドレスにスレッドのメモリ領域を割り当ててしまうと、ポインタが無効化してしまい、プログラムの正しい実行を保証できなくなる。この問題に対しては、主に 2

つの解決策が提案されている。

第一の解決策は、移動元プロセスと移動先プロセスとで異なるアドレスにメモリ領域を配置することを許すかわりに、スレッド移動直後に、スレッドのメモリ領域に含まれるすべてのポインタを、移動先プロセスのアドレス領域に合わせて完全に正しく更新する手法<sup>22),33)–36)</sup>である。この手法では、スレッド移動時にどれがポインタなのかを処理系が完全に把握する必要があるので、どれがポインタなのかをプログラマに明示的に指定させたり、データフロー解析などのコンパイラ的手法を用いてポインタを自動的に発見したりする。しかし、前者の方法はプログラミングの負担を増大させるという問題があり、後者の方法は、本質的に C 言語は型安全な言語ではないのですべてのポインタを自動的に完全に発見することはできないという問題がある。

第二の解決策は、iso-address と呼ばれる方法で、アドレス空間全体をあらかじめいくつかに分割しておき、各スレッドが使用可能なアドレス空間を静的に決め打っておく方法<sup>9),10),45)</sup>である。これにより、あるスレッドが使用しているアドレス空間が他のいかなるスレッドによっても使用されていないことをつねに保証できるため、スレッド移動時には、移動元プロセスと移動先プロセスとでつねに同一アドレスにメモリを割り当てることができる。既存のスレッド移動の研究の多くは iso-address を用いている<sup>16),33)</sup>。しかし、iso-address は、計算規模がアドレス空間全体の大きさに制限されるという問題がある。アドレス空間全体の大きさを  $w$  バイト、スレッド数を  $n$ 、各スレッドが使用可能なメモリの大きさを  $s$  とすると、iso-address では  $ns = w$  が成立している必要がある。よって、32 ビットアーキテクチャであれば  $w = 2^{32}$  なので、たとえば  $n = 1024$  個のスレッドを生成するならば各スレッドが使用できるメモリ量はわずか  $s = 4\text{MB}$  であり、各スレッドが  $s = 2\text{GB}$  のメモリ量を使用するならばスレッドはわずか  $n = 2$  個しか生成できず、非現実的である。一方、近年の多くの 64 ビットアーキテクチャでは (CPU の実装に依存するが)  $w = 2^{47}$  のアドレス空間を利用できるため、これをもって iso-address の欠点は解消されたと見る向きもある<sup>32),39),57)</sup>が、これも楽観的である。なぜなら、 $w = 2^{47}$  であっても、 $n = 8192$  個ならば  $s = 64\text{GB}$ 、 $n = 1024$  個ならば  $s = 512\text{GB}$  であり、これらの数字は 2010 年 7 月現在のクラスタ規模や各ノードの搭載メモリ量から見れば十分に現実的な数字だからである。以上の考察より、今後ますます増大する計算規模に対応するためには、これら既存研究のアプローチでは不十分であり、計算規模がアドレス空間全体の大きさに制限されないスレッド移動の手法が要請されていると言える。当然、ハードウェアの進化にともなって  $w = 2^{47}$  という数字は今後増える可能性もあるが、そうであっても、計算規模がアドレス空間全体の大きさに制限され

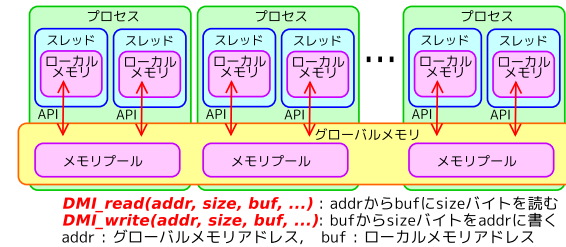


図 2 DMI のシステム構成。

ないスレッド移動の手法が存在することには価値がある。そこで DMI では、そのようなアドレス手法として random-address を提案する (6 節)。

### 3. 基本設計

本節では、DMI の基本設計について概観する。詳細に関しては論文<sup>29),65),66)</sup>を参照されたい。

#### 3.1 DMI の概要

DMI のシステム構成を図 2 に示す。DMI では各ノード上に任意個のプロセスを生成し、各プロセス内に任意個のスレッドを生成することができる。ただし、性能上は、1 ノードあたり 1 個のプロセスを、1 プロセッサあたり 1 個のスレッドを生成するのが望ましい。

DMI では、Co-Array Fortran や UPC などの GET/PUT 操作に基づく PGAS 処理系とは異なり、データのキャッシュをサポートするため、分散共有メモリと同様のコンシステンシ管理を行う。各プロセスは、メモリプールと呼ばれる一定量のメモリを DMI に提供する。このメモリプールの量は各プロセスを生成するときに明示的に指示できる。すると、DMI はこれらのメモリプールをメモリ資源として、ページテーブルなどのメモリ管理機構をユーザレベルで実装することによって、分散環境上に仮想的な共有メモリを構築する。この仮想的な共有メモリのことをグローバルメモリと呼ぶ。これにより、各スレッドは、グローバルメモリへの read/write を通じて、すべてのプロセスが提供するメモリプールに透過的にアクセスすることができる。このとき、アクセス対象のデータがそのプロセスのメモリプールに存在しない場合には、ページフォルトが発生して、他プロセスとの通信が行われてコンシステンシ管理が行われる。さらに、必要であれば、read/write したデータをメモリプールにキャッシュすることができる。このように DMI では、同一プロセス内の複数のスレッド



がメモリプールを共有キャッシュ的に利用する構成となっており、同一プロセス内のスレッド間のデータ共有は通常の共有メモリ経由で実現される。すなわち、DMI は、分散レベルの並列性と同時に、マルチコアレベルの並列性を統合的に活用したハイブリッドプログラミングを透過的に実現している。また、多数のプロセスを利用することで巨大なグローバルメモリを構築し、遠隔スワップシステムとして動作させることもできる。ページングを繰り返すうちにメモリプールの使用量が指定量を超えてしまう場合があるが、その場合には、ページ置換アルゴリズムに基づいて他プロセスへのページアウトが行われる。

DMI は C 言語の共有ライブラリとして実装されており、コンパイラや OS には一切手を加えていないため移植性が高い。DMI の処理系は約 22000 行からなる C 言語で実装されており、86 個の API を提供している。

### 3.2 大域アドレス空間に対するアクセス

DMI では、各スレッドが malloc などによって確保するローカルメモリと、グローバルメモリを明確に分離している。ローカルメモリは通常の共有メモリであり、malloc/free を通じて確保/解放し、通常の変数参照や配列参照などによって read/write できる。一方、グローバルメモリは、DMI\_mmap() 関数/DMI\_munmap() 関数によって確保/解放し、DMI\_read() 関数/DMI\_write() 関数によって read/write する。以下、これらの API について詳しく述べる。

DMI では、region-based<sup>41)</sup> な分散共有メモリと同様に、ユーザプログラムの振る舞いに合致した任意のコンシステンシ粒度を指定してグローバルメモリを確保することができる。DMI では、コンシステンシ粒度のことをページ、そのサイズをページサイズと呼ぶ。具体的には、DMI\_mmap(int64\_t \*addr, int64\_t page\_size, int64\_t page\_num, ...) 関数を呼び出すことによって、ページサイズが page\_size のページを page\_num 個持つグローバルメモリを確保できる。つまり、PGAS 処理系に即した表現で言えば、任意のブロックサイズに基づくブロックサイクリックなデータ分散を指定することができる。たとえば、巨大な行列行列積をブロック分割によって並列に行いたい場合には、各行列ブロックのサイズをページサイズに指定してグローバルメモリを割り当てれば良い。

グローバルメモリに対して、read/write するには、DMI\_read(int64\_t addr, int64\_t size, void \*buf, ...) 関数/ DMI\_write(int64\_t addr, int64\_t size, void \*buf, ...) 関数を使う。DMI\_read(int64\_t addr, int64\_t size, void \*buf, ...) 関数は、グローバルメモリアドレス addr から size バイトをローカルメモリアドレス buf に read する。一方で、DMI\_write(int64\_t addr, int64\_t size, void \*buf, ...) 関数は、ローカルメモリアド

レス buf から size バイトをグローバルメモリアドレス addr に write する。DMI では、アドレス領域 [addr, addr + size) が 1 ページに収まるような DMI\_read() 関数/DMI\_write() 関数に対する Sequential Consistency が保証されており、直観的に理解しやすいコンシステンシモデルのもとで並列プログラムを記述できる。複数のページにまたがって DMI\_read() 関数/DMI\_write() 関数を呼び出すことも可能であるが、この場合には、要求されたアドレス領域全体がページ単位のアドレス領域に内部で分割され、その分割された各アドレス領域に対して独立に DMI\_read() 関数/DMI\_write() 関数が呼び出されたのと同様の結果になる。

このように DMI では、API 呼び出しを通じてしかグローバルメモリにアクセスできないが、API 呼び出しの形式にすることの利点は、OS のメモリ保護違反に頼ることなく、高度に柔軟なユーザレベルでのコンシステンシ管理が可能になる点である。第一に、ユーザレベルでコンシステンシを管理することで、OS のページサイズに制限されない任意粒度でのコンシステンシ管理が可能になっている。これにより、ページサイズが固定されている page-based な分散共有メモリと比較すると、ページフォルトの回数を大幅に抑制できるうえ、データ転送の単位をユーザプログラムにとって必要十分なサイズまで巨大化させられるため、効率的な通信を実現できる。第二に、DMI\_read() 関数/DMI\_write() 関数の引数としてさまざまな情報を与えることができるため、各 read/write の粒度で、明示的で強力な最適化手段を提供することができる。具体的には、マルチモード read/write という機能によって、(1) read/write したときに各ページをキャッシュするかいなか、(2) さらにキャッシュする場合にはそのページを invalidate 型としてキャッシュするのか update 型としてキャッシュするのか、(3) そのページのキャッシュ管理のオーナー権を read/write したプロセスに移動させるかいか、を各 read/write の粒度で明示的に指定することができる。さらに、非同期に DMI\_read() 関数/DMI\_write() 関数を発行できる。この非同期操作を用いることで、アプリケーションが要求するコンシステンシ強度に応じて、DMI によって保証されている Sequential Consistency を明示的に緩和させることもできる。

一方で、API 呼び出しの形式にすることの欠点はプログラム記述の煩雑さである。しかし、作業的には面倒であるが論理的に難解なわけではない。もう 1 つの欠点は、ユーザレベルでコンシステンシ管理を行うため、ページフォルトを引き起こさないアクセスに対しても逐一ユーザレベルの検査が入るうえ、グローバルメモリとローカルメモリの間でのメモリコピーが必要となるため、オーバーヘッドが大きい点である。しかし、多くのアプリケーションでは、ここで生じるオーバーヘッドによる損失よりも、前述したような柔軟なコンシステンシ管理によって得られる利得の方が大きいと考えられる。以上のような API の特徴を

ふまえると、高性能な DMI プログラムを記述するためには、できるかぎり大きなグローバルメモリ領域に対して DMI\_read() 関数/DMI\_write() 関数を発行することで、DMI\_read() 関数/DMI\_write() 関数が呼ばれる回数を少なくすることが重要である。

### 3.3 動的なプロセスの参加/脱退

DMI では、動的なプロセスの参加/脱退を越えてグローバルメモリのコンシステンスを維持するプロトコルが実装されており、並列計算を実行中に動的にプロセス（ノード）を参加/脱退させることができる。プログラミングインタフェースとしては、参加/脱退しようとしているプロセスをポーリングする API、それらの参加/脱退を承認する API、スレッドを生成/破棄する API などが提供されており、参加プロセスに対してスレッドを生成したり脱退プロセスからスレッドを破棄したりすることで、動的なノードの参加/脱退に対応して計算規模を動的に拡張/縮小させる並列プログラムを記述することができる。言い換えると、DMI では、計算規模を拡張/縮小させるためには、プログラマがスレッドの生成/破棄を通じて明示的にスレッド数を増減させる必要がある。

一般に、タスクパラレルなアルゴリズムでは、並列計算を実行中に明示的にスレッド数を増減させるのが比較的容易である。なぜなら、タスクパラレルなアルゴリズムでは、アルゴリズム上の論理的な並列度と実際の物理的な並列度とが概念的に分離されているため、アルゴリズム上はいつ何個のスレッドが参加していても良く、各タスクの粒度において自由なタイミングでスレッドを増減させられるからである。一方で、SPMD 型のアルゴリズムでは、並列計算を実行中に明示的にスレッド数を増減させるのが困難である。なぜなら、SPMD 型のアルゴリズムでは、アルゴリズム上の論理的な並列度と実際の物理的な並列度が一致しているからである。SPMD 型のアルゴリズムでは、すべてのスレッドが同期的に動作するよう記述されるため、並列計算の途中でスレッド数を増減させるためには (1) すべてのスレッドが同期をとり (2) スレッド数を増減させ (3) 新たなスレッド数において各スレッドの担当範囲を分割し直す、という手順を行う必要があるからである。これはプログラミング上の負担を増大させるだけでなく、担当範囲の再分割に時間がかかる場合には性能上の問題も引き起こす。たとえば、8.3.1 節に示す有限要素法においては、担当範囲の分割は物体全体の領域分割に相当するが、この領域分割は計算時間を要する処理である。

以上をふまえて、本研究では、以上のような基本設計を持つ DMI に対してスレッド移動の技術を付け加え、1.2.3 節で述べたコンセプトに基づき、計算規模を動的に拡張/縮小できる並列プログラムをより簡単に記述できるプログラミングインタフェースを整備する。

```
typedef struct arg_t {
    ...; /* 各スレッドに渡す引数 */
}arg_t;

void DMI_main(int argc, char **argv)
{
    arg_t arg;
    int rank, pnun;
    int64_t sched_addr, arg_addr, handle[THREAD_MAX];
    ...;
    DMI_mmap(&sched_addr, sizeof(DMI_scheduler_t), 1, ...); /* スレッドスケジューラを管理するグローバルメモリを確保 */
    DMI_mmap(&arg_addr, sizeof(arg_t) * pnun, 1, ...); /* 各スレッドへの引数を格納するグローバルメモリを確保 */
    for(rank = 0; rank < pnun; rank++) {
        arg = ...; /* 各スレッドに渡す引数を設定 */
        DMI_write(arg_addr + rank * sizeof(arg_t), sizeof(arg_t), &arg, ...); /* グローバルメモリに書き込む */
    }
    DMI_scheduler_init(sched_addr); /* スレッドスケジューラを初期化 */
    for(rank = 0; rank < pnun; rank++) { /* スレッドを生成 */
        DMI_scheduler_create(sched_addr, &handle[rank], arg_addr + rank * sizeof(arg_t));
    }
    for(rank = 0; rank < pnun; rank++) { /* スレッドを回収 */
        DMI_scheduler_join(sched_addr, handle[rank], ...);
    }
    DMI_scheduler_destroy(sched_addr); /* スレッドスケジューラを破棄 */
    DMI_munmap(sched_addr, ...);
    DMI_munmap(arg_addr, ...);
    ...;
}

int64_t DMI_thread(int64_t arg_addr) /* 各スレッドの処理 */
{
    arg_t arg;
    int iter;
    DMI_read(arg_addr, sizeof(arg_t), &arg, ...); /* グローバルメモリから引数を読み込む */
    for(iter = 0; iter < ITER_MAX; iter++) {
        DMI_yield(); /* スレッドスケジューラにスケジューリングのチャンスを与える */
        ...; /* 各イテレーションの処理 */
    }
}
```

図 3 DMI のプログラミングインタフェース。

## 4. プログラミングインタフェース

1.2.3 節で述べたように、本研究では、プログラマは十分な数のスレッドを生成するだけで良く、あとは処理系が透過的に、それらのスレッドを、そのとき利用可能な計算資源にマッピングするような並列プログラミング処理系を提案する。そのプログラミングインタフェースを図 3 に示す。図 3 において、DMI が起動されると DMI\_main() 関数が実行される。そして、DMI\_scheduler\_init() 関数を呼び出して DMI のスレッドスケジューラを初



期化したあと、`DMI_scheduler_create()` 関数を呼び出すことで任意個のスレッドを生成できる。生成されたスレッドは `DMI_thread()` 関数として実行され始める。また、終了したスレッドを `DMI_scheduler_join()` 関数で回収したり、`DMI_scheduler_detach()` 関数でデタッチしたりできる。

ここで、スレッドスケジューリングは DMI によって透過的に行われるが、DMI におけるスレッド移動は、プリエンティブではなく協調的に行われる。具体的には、DMI によってスレッド移動が必要であると判断された任意のタイミングでスレッド移動が起きるわけではなく、それ以降で、移動対象のスレッドが初めて `DMI_yield()` 関数を呼び出した時点で、その `DMI_yield()` 関数の内部で協調的なスレッド移動が起きる。この `DMI_yield()` 関数は、DMI によってスレッド移動の指示が届いていなければ何も行わずにすぐに返り、スレッド移動の指示が届いていれば内部でスレッド移動を行って移動先のプロセスで返る。したがって、外的な計算資源の変化に追従してスレッド移動を応答性良く行うためには、プログラマは、移動される可能性のあるスレッドがある程度短い間隔で `DMI_yield()` 関数を呼び出すようにプログラムを記述しておく必要がある。反復計算を行うプログラムであれば、たとえば図 3 のように、各イテレーションの先頭に `DMI_yield()` 関数を記述すれば良い。

このように、DMI では、通常の共有メモリ環境におけるスレッドプログラミングと同様のプログラミングインタフェースにより、計算規模の拡張/縮小を意識することなく並列プログラムを記述できる。

## 5. プログラミング制約

### 5.1 制約が必要になる理由

DMI のように、1 個のプロセス内に多数のスレッドが存在するマルチスレッド型のモデルにおいて、各スレッドを粒度としたスレッド移動を行うためには、各スレッドの使用アドレス領域が独立している必要がある。たとえば、プロセス  $p$  の中にスレッド  $i$  とスレッド  $j$  があり、スレッド  $i$  はスレッド  $j$  のスタック領域のどこかへのポインタ  $d$  を使用しているとすると、このとき、スレッド  $j$  だけを別のプロセス  $q$  にスレッド移動させたとすると、スレッド  $i$  が使用しているポインタ  $d$  が無効化されてしまい、実行を正しく継続できなくなる。別の例として、プロセス  $p$  内のスレッド  $i$  とスレッド  $j$  が、プロセス  $p$  のグローバル変数  $g$  を使用している状態で、スレッド  $j$  だけを別のプロセス  $q$  に移動させることを考える。この場合には、グローバル変数  $g$  をプロセス  $q$  に移動させてもささなくても、スレッド  $i$  とスレッド  $j$  のいずれか一方の実行を正しく継続できなくなる。また、そもそも、プロセ

ス  $q$  にもグローバル変数  $g$  がすでに存在しているとすれば、プロセス  $p$  のグローバル変数  $g$  をプロセス  $q$  へ移動することによって、プロセス  $q$  のグローバル変数  $g$  の値が書きつぶされてしまう。また、ローカルなファイル IO などはスレッド移動を越えてサポートできない。以上からわかるように、各スレッドを粒度としたスレッド移動を安全に行うためには、C 言語で記述できることすべてをサポートできるわけではなく、アドレス領域の使用を何らかに制約する必要がある。よって本章では、まず 5.2 節でアドレス領域をモデル化したうえで、そのモデルに基づいて 5.3 節でプログラミング制約を記述する。さらに、5.5 節でそのプログラミング制約を緩和する。

### 5.2 アドレス領域のモデル化

まず、アドレス領域をモデル化する。DMI では、アドレス領域全体を以下の 6 種類に分類して扱う (図 4):

$register_i^p$  プロセス  $p$  内のスレッド  $i$  のレジスタ領域。 $register_i^p$  はマシンによってスレッドローカルに管理される。

$stack_i^p$  プロセス  $p$  内のスレッド  $i$  のスタック領域。 $stack_i^p$  はマシンによってスレッドローカルに管理される。

$threadheap_i^p$  プロセス  $p$  内のスレッド  $i$  のヒープ領域。 $threadheap_i^p$  の定義は、プロセス  $p$  内のスレッド  $i$  が `DMI_thread_mmap()` 関数/`DMI_thread_mremap()` 関数/`DMI_thread_munmap()` 関数を呼び出すことで確保/解放されるアドレス領域である。 $threadheap_i^p$  は DMI によってスレッドローカルに管理される。`DMI_thread_mmap()` 関数および `DMI_thread_mremap()` 関数は、返り値として通常のポインタを返すので、このアドレス領域は通常のメモリアクセスによって使用できる<sup>\*1</sup>。

$static^p$  プロセス  $p$  の静的変数領域。 $static^p$  はマシンによってプロセスローカルに管理される。

$processheap^p$  プロセス  $p$  のヒープ領域。 $processheap^p$  の定義は、プロセス  $p$  に含まれるいずれかのスレッドが (`malloc()` 関数/`free()` 関数などを經由して) システムコールの `mmap()` 関数/`mremap()` 関数/`munmap()` 関数を呼び出すことで確保/解放されるアドレス領域である。 $processheap^p$  はマシンによってプロセスローカルに管理される。

\*1 なお、プロセス  $p$  内のスレッド  $i$  が `DMI_thread_mmap()` 関数によって確保したアドレス領域を、プロセス  $p$  内の別のスレッド  $j$  が `DMI_thread_munmap()` 関数で解放することはできない。また、プログラミングの便宜のため、`DMI_thread_mmap()` 関数/`DMI_thread_munmap()` 関数/`DMI_thread_mremap()` 関数の上位関数として、`DMI_thread_malloc()` 関数/`DMI_thread_free()` 関数/`DMI_thread_realloc()` 関数を提供している。

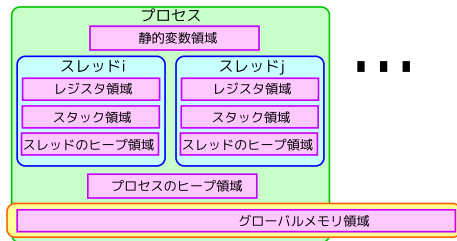


図 4 DMI におけるアドレス領域のモデル化.

```
static char *buf;
static pthread_mutex_t mutex = PTHREAD_MUTEX_INITIALIZER;

int printf(format, ...)
{
    pthread_mutex_lock(&mutex);
    if(buf == NULL) {
        buf = malloc(4096);
    }
    ...; /* buf を使って文字列 format を値で埋める */
    write(1, buf, strlen(buf));
    pthread_mutex_unlock(&mutex);
}

```

図 5 printf() 関数の実装例.

*dmi* DMI が実現するグローバルメモリ領域。*dmi* の定義は、`DMI_mmap()` 関数/`DMI_munmap()` 関数を呼び出すことで確保/解放されるアドレス領域である。*dmi* は DMI によって全プロセスで共有されるように管理される。

### 5.3 プログラミング制約

以上で述べたアドレス領域のモデルに基づき、プログラミング制約について述べる。第 4 章で述べたように、DMI では、プリエンティブなスレッド移動ではなく、ユーザプログラムに `DMI_yield()` 関数を呼び出してもらうことで協調的なスレッド移動を行う。このとき、この `DMI_yield()` 関数に関して以下のプログラミング制約が守られる必要がある：

**プログラミング制約** プロセス  $p$  内のスレッド  $i$  が `DMI_yield()` 関数を呼び出す時点において、そのスレッド  $i$  の実行を正しく継続するために必要なすべてのデータは、 $register_i^p$ 、 $stack_i^p$ 、 $threadheap_i^p$ 、*dmi* のいずれかのアドレス領域に含まれている\*1。

したがって、`DMI_yield()` 関数を呼び出す時点で、グローバル変数 ( $static_i^p$ ) を使用していたり、`malloc()` 関数で確保したアドレス領域 ( $processheap_i^p$ ) を使用していたりすると、スレッド移動後の実行の正しさは保証されない。

ここで注意すべき点は、このプログラミング制約は、`DMI_yield()` 関数を呼び出す時点についてしか言及していないという点である。よって、`DMI_yield()` 関数を呼び出す時点でプログラミング制約が守られてさえいれば、`DMI_yield()` 関数を呼び出していない時点においては、 $static^p$ 、 $processheap^p$  のアドレス領域を使用しても問題はない。たとえば (1)  $p=malloc(\dots)$  関数によりアドレス領域  $p$  を確保し (2) アドレス領域  $p$  を使用して計算を行い (3) `free(p)` 関数によりアドレス領域  $p$  を解放する、という処理を行う関数として  $f()$

関数を考える。このとき、`DMI_yield()` 関数を呼び出す前に  $f()$  関数を呼び出したとしても、プログラミング制約には違反しない。別の例としては、`DMI_read()` 関数や `DMI_write()` 関数などの DMI 関数は、内部的には `malloc()` 関数などを使用しているが、すべての DMI 関数は、その DMI 関数の実行が終了した時点で  $static^p$ 、 $processheap^p$  にはスレッドの実行を継続するために必要なデータを残さないように実装されているため、`DMI_yield()` 関数を呼び出す前に DMI 関数を呼び出しても、当然プログラミング制約には違反しない\*2。

以上の観察をまとめると、スレッド移動を行うユーザプログラムに対しては一定のプログラミング制約が課せられるものの、このプログラミング制約のもとでは以下の記述が許されているため、8.2 節で評価するような並列科学技術計算を記述できるだけの記述力は保たれていると言える：

- `DMI_thread_mmap()` 関数/`DMI_thread_munmap()` 関数/`DMI_thread_mremap()` 関数によるスレッドローカルなアドレス領域の確保/解放と、そのアドレス領域に対する通常のメモリアクセス。
- `DMI_mmap()` 関数/`DMI_munmap()` 関数/`DMI_mremap()` 関数によるグローバルメモリの確保/解放と、そのアドレス領域に対する `DMI_read()` 関数/`DMI_write()` 関数などによるメモリアクセス。
- グローバルメモリに対する同期操作、プロセスの参加/脱退操作などの各種 DMI 関数の呼び出し。

### 5.4 スレッド移動の手順

前節で述べたプログラミング制約のもとでは、以下の手順により、プロセス  $p$  内のスレッド  $i$  をプロセス  $q$  へと移動させることができる：

- (1) スレッド  $i$  が `DMI_yield()` 関数を呼び出したとき、スレッド  $i$  の移動が指示されれば、スレッド  $i$  を停止させる。
- (2) プロセス  $p$  における  $register_i^p$ 、 $stack_i^p$ 、 $threadheap_i^p$  のアドレス領域を、プロセス  $q$  に対して送信する。
- (3) プロセス  $q$  は、受信した  $register_i^p$ 、 $stack_i^p$ 、 $threadheap_i^p$  のアドレス領域を、プロセス  $p$  で使用されていたアドレス領域とまったく同一のアドレス領域に割り当てる。
- (4) プロセス  $q$  はスレッド  $i$  を復旧させ、`DMI_yield()` 関数を返す。

\*1 ただし、非同期 DMI 関数に関しては、その非同期操作が完了するまでは、スレッドの実行を継続するために必要なデータが  $static^p$ 、 $processheap^p$  に残されているため、非同期操作の完了を回収するまでは `DMI_yield()` 関数を呼び出してはいけない。

\*1 なお、このプログラミング制約は 5.5 節において少し緩和して言い直される。

$dmi$  のアドレス領域に関しては、スレッド移動に伴って何もする必要はない。なぜなら、グローバルメモリに関しては、DMI のコンシステンシプロトコルによって、どの時点でどのスレッドからアクセスされても問題がないようにコンシステンシが維持されているためである。なお、上記の (3) において、プロセス  $p$  において  $register_i^p$ ,  $stack_i^p$ ,  $threadheap_i^p$  が使用していたアドレス領域がプロセス  $q$  で使用されていない保証はないため、まったく同一のアドレス領域に割り当てることができる保証はない。この対策に関しては第 6 章で述べる。

### 5.5 プログラミング制約の緩和

5.3 節において、DMI におけるスレッド移動時のプログラミング制約は、並列科学技術計算を記述するための記述力を保っていると述べた。しかし、グローバル変数を使用できないことはプログラマに対して一定の不便を強いると考えられる。なぜなら、グローバル変数を使用できないということは、グローバル変数を使用しうるすべてのライブラリ関数をユーザプログラムから呼び出せないことを意味するからである。したがって、(実装に依存するが) `printf()` 関数、`malloc()` 関数などの、内部でグローバル変数を使用する `libc` 共有ライブラリの多くのライブラリ関数は呼び出せないことになる。ところが、実は、5.3 節で述べたプログラミング制約は少し緩和させることができ、特定の条件を満たすライブラリ関数であれば、内部でグローバル変数を使用していても安全に呼び出すことができる。以下では、グローバル変数を使用する `printf()` 関数を例にとり、プログラミング制約がどう緩和できるかを議論する。

図 5 に示すような実装の `printf()` 関数を考える。`libc` 共有ライブラリの `printf()` 関数の実装では、任意長のフォーマット文字列を許可したり、出力のバッファリングなどを行っていると思われるが、簡単のため省略する。この `printf()` 関数では、1 回目の呼び出しでグローバル変数 `buf` にメモリを確保し、2 回目以降の呼び出しでは 1 回目の呼び出し時に確保したメモリを再利用する仕様になっている。言い換えると、グローバル変数 `buf` にスレッドの実行を継続するために必要なデータを格納したまま、関数が返る仕様になっている。よって、プロセス  $p$  内のスレッド  $i$  を別のプロセス  $q$  に移動させることを考えたとき、スレッド  $i$  が `DMI_yield()` 関数を呼び出す前に一度でも `printf()` 関数を呼び出しているならば、`DMI_yield()` 関数を呼び出す時点で、スレッドの実行を継続するために必要なデータが `static^p` に格納されていることになり、プログラミング制約に違反してしまう。具体的には、DMI ではスレッド移動の際に `static^p` を移動させないため、スレッド  $i$  が移動先プロセス  $q$  で最初に `printf()` 関数を呼び出した時点でグローバル変数の不一致が起き、問題が

生じるように思われる。ところが、実際には何の問題も生じない。なぜなら、図 5 における `printf()` 関数の実装を注意深く観察するとわかるように、スレッド  $i$  が移動先プロセス  $q$  で呼び出した `printf()` 関数は、移動元プロセス  $p$  のグローバル変数 `buf` の値とは無関係に、その時点での移動先プロセス  $q$  のグローバル変数 `buf` の値に基づいて正しく実行されるからである。すなわち、この `printf()` 関数に関しては、グローバル変数を使用しているにもかかわらず、スレッド移動に伴ってグローバル変数を移動させなくても問題は生じない。

以上のような現象が生じる理由は、この `printf()` 関数は、実際にはグローバル変数を使用しているものの、任意のスレッドによって任意の順序で呼び出されても正しく実行されるようなセマンティクスでグローバル変数を使用しているためである。言い換えると、セマンティクスとしては、スレッドの実行を継続するために必要なデータがグローバル変数に入っていないと見なすことができるためである。以上の議論より、先ほど定義したプログラミング制約は以下のように緩和することができる：

**プログラミング制約** プロセス  $p$  内のスレッド  $i$  が `DMI_yield()` 関数を呼び出す時点において、そのスレッド  $i$  の実行を正しく継続するために必要なすべてのデータは、 $register_i^p$ ,  $stack_i^p$ ,  $threadheap_i^p$ ,  $dmi$  のいずれかの領域に含まれているセマンティクスになっている。

ここで問題なのは、各ライブラリ関数が上記のプログラミング制約を満たすかどうかは、通常は仕様として規定されていないため、逐一ライブラリ関数の実装を調べなければならない点である。しかし、本節で主張すべきことは、ライブラリ関数の実装を注意深く検討しさえすれば、仮にそのライブラリ関数が内部で `static^p`, `processheap^p` のアドレス領域を使用しているとしても、スレッド移動の安全性に影響を与えないようにそれらのライブラリ関数を呼び出すことができる場合がある、ということである。実際、スレッドセーフな `libc` 共有ライブラリの関数の中には安全に呼び出せるものも多い。

ここまでグローバル変数に関連する問題を議論したが、既存研究においても、グローバル変数の取扱いはスレッド移動を行ううえで深刻な問題とされてきた。第一に、`PM2`<sup>(9),(10)</sup>, `Adaptive MPI`<sup>(31)</sup>, `MigThread`<sup>(34)–(36)</sup>, `Arachne`<sup>(24)</sup> など多くの既存研究はそもそもグローバル変数の使用を禁止している。第二に、Windows 環境においてスレッド移動を実現する `Tern`<sup>(38)</sup> では、スレッド移動に伴うスレッドローカルストレージの移動をサポートしているため、プログラマは「各スレッドにとってグローバルな」変数を利用できる。しかし、Windows 環境とは異なり、DMI が想定している Linux 環境では、ユーザレベルからランタイムにスレッドローカルストレージを抽出するのは難しい。第三に、Java においてスレッ

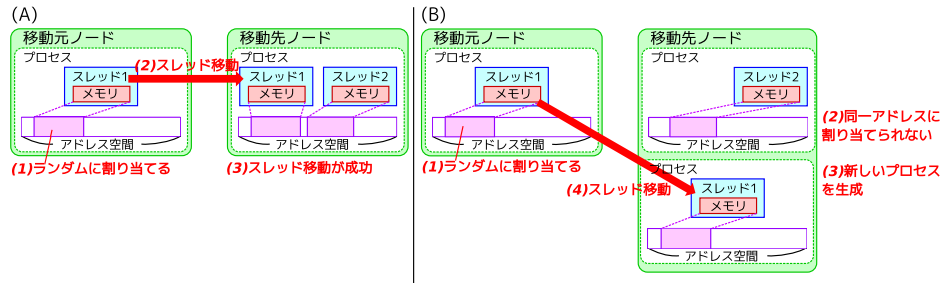


図 6 random-address のアルゴリズム (A) アドレスが衝突しない場合 (B) アドレスが衝突する場合

ド移動を実現している JESSICA<sup>2(37),62,63)</sup> では, Delta Execution というマスターカ型の手法を用いて, グローバル変数のサポートを実現している. Delta Execution では, 系内に存在するすべてのスレッドは, マスタノードに「親スレッド」を持つ. 各スレッドはマスタノードおよび各ワーカノードを自由にスレッド移動できるが, あるスレッド  $i$  がワーカノードにおいて, グローバル変数へのアクセスやファイルアクセスなどプロセス依存な操作を行おうとした場合には, プログラムの実行をマスタノードに存在するスレッド  $i$  の「親スレッド」に引き渡し, マスタノード上でそのプロセス依存な操作を実行する. そして, それらのプロセス依存な操作が完了したあと, 再びプログラムの実行をワーカノードに引き戻す. これにより, プロセス依存な操作はすべてマスタノードで実行されることになるため, ユーザプログラムからグローバル変数などを使用しても差し支えない. しかし, この Delta Execution は Java のバイトコードに手を加えることで実現されており, DMI のような C 言語における処理系に応用させるのは難しい.

## 6. アドレス空間の大きさに制限されないスレッド移動

### 6.1 基本アイデア

DMI では, 計算規模がアドレス空間の大きさに制限されないスレッド移動の手法として, random-address を提案する. random-address のアルゴリズムの基本アイデアは以下のとおりである:

- (1) 各スレッドは, 自分以外のスレッドがどのアドレスにローカルメモリを割り当てているかに関する知識を持たない. 各スレッドは, 乱数を使って, ローカルメモリを割り当てるアドレスを決定する (図 6 (A)). よって, 各スレッドがローカルメモリを割り当てる

操作 (DMI\_thread\_mmap() 関数/DMI\_thread\_munmap() 関数/ DMI\_thread\_mremap()

関数/) は, 他のスレッドとの通信をいっさい必要とせず, 完全に独立に実行できる.

- (2) いま, ノード  $P$  上のプロセス  $p$  に存在するスレッド  $i$  を, ノード  $Q$  上のプロセス  $q$  へと移動させるとする. このとき, 「運が良ければ」, 移動元プロセス  $p$  においてスレッド  $i$  が使用しているアドレスは, 移動先プロセス  $q$  では使用されていない. この場合には, 移動先プロセス  $q$  において, スレッド  $i$  のローカルメモリを移動元プロセス  $p$  と同一のアドレスに割り当てることで, スレッド移動を完了させる (図 6 (A)).
- (3) スレッド  $i$  の移動時に, 「運が悪ければ」, スレッド  $i$  が移動元プロセス  $p$  において使用しているアドレスが, すでに移動先プロセス  $q$  でも使用されている. この場合には, 当然, 移動先プロセス  $q$  において, スレッド  $i$  のローカルメモリを移動元プロセス  $p$  と同一のアドレスに割り当てることができない. そこで, 移動先プロセス  $q$  が存在するノード  $Q$  上に新しいプロセス  $q'$  (要するに新しいアドレス空間) を生成し, 新しいプロセス  $q'$  の中にスレッド  $i$  を移動させる (図 6 (B)).

このように, random-address ではアドレスが衝突した場合にプロセス数が増えるが, スレッドスケジューリングを適切に行ってスレッドが存在しなくなったプロセスを破棄するようにすれば, スレッド移動を繰り返してもプロセス数が増え続けることはない. たとえば, ノード  $P$  上に存在するスレッド  $i$  とスレッド  $j$  に関して, ある時刻  $t$  において, スレッド  $i$  が使用しているアドレス集合とスレッド  $j$  が使用しているアドレス集合に重なりがあり, スレッド  $i$  はノード  $P$  上のプロセス  $p$  にスレッド  $j$  はノード  $P$  上のプロセス  $a'$  に入っている状況を考える. このとき, 仮に, スレッド  $j$  が使用しているアドレス集合と, 別のノード  $Q$  上にすでに存在しているプロセス  $q$  が使用しているアドレス集合に重なりがなければ, スレッド  $j$  をプロセス  $q$  の中へ移動させることで, プロセス  $a'$  を破棄することができる. あるいは, 時刻  $t + \Delta t$  において, スレッド  $i$  またはスレッド  $j$  が使用しているローカルメモリが変化して, スレッド  $i$  が使用しているアドレス集合とスレッド  $j$  が使用しているアドレス集合に重なりがなくなったとすれば, スレッド  $j$  をプロセス  $p$  の中にスレッド移動させることで, プロセス  $a'$  を破棄することができる. 理論的には,  $m$  個のスレッド  $x_0, x_1, \dots, x_{m-1}$  に関して, ある時刻  $t$  において, これら各スレッド  $x_i$  が使用しているアドレス集合を  $S_{x_0}^t, S_{x_1}^t, \dots, S_{x_{m-1}}^t$  とするとき, これら  $m$  個のスレッドは, 最小  $f(S_{x_0}^t, S_{x_1}^t, \dots, S_{x_{m-1}}^t)$  個のプロセスに格納することができる. ここで  $f(S_{x_0}^t, S_{x_1}^t, \dots, S_{x_{m-1}}^t)$  とは, 以下の条件を満たす  $F$  のうち最小の値とする:

条件  $m$  個の集合  $S_{x_0}^t, S_{x_1}^t, \dots, S_{x_{m-1}}^t$  を,  $F$  個のグループ  $G_0, G_1, \dots, G_{F-1}$  に分

類したとする．つまり，

$$\forall i (0 \leq i \leq m-1), \exists j (0 \leq j \leq F-1), \forall k (0 \leq k \leq F-1 \wedge k \neq j) : \\ S_i^t \in G_j \wedge S_i^t \notin G_k$$

となるように各  $S_i^t$  を分類したとする．このとき，

$$\forall i (0 \leq i \leq F-1), \forall S_j^t (S_j^t \in G_i), \forall S_k^t (S_k^t \in G_i \wedge k \neq j) : S_j^t \cap S_k^t = \emptyset$$

が成り立つ．

しかし，当然ながら，任意の時刻  $t$  において， $m$  個のスレッドを  $f(S_{x_0}^t, S_{x_1}^t, \dots, S_{x_{m-1}}^t)$  個のプロセスに格納するようにスレッドを管理するのは現実的ではない．実際には，ノード間のスレッドの負荷バランス，スレッド移動に要する時間，プロセス数を増やすことによるオーバーヘッド，各スレッド間でのデータ共有の度合いなどの要素を総合的に考慮して，スレッドスケジューリングを最適化する必要がある．ただし，本研究ではスレッドスケジューリングの最適化は考察の対象外とし， $m$  個のスレッドを  $n$  個のプロセスにスケジューリングする場合には，単純に，スレッド番号が若い方から順に，各プロセスに  $m/n$  個（または  $m/n+1$  個）ずつスレッドを割り当てるようなスレッドスケジューリングを行うものとする．

以上からわかるように，random-address においては，動的にプロセスを生成したり破棄したりする必要があり，これを実現するためには，当然，そのプロセスを動的にグローバルメモリに参加させたり脱退させたりする必要がある．つまり，この random-address は，DMI がプロセスの動的な参加/脱退に対応しているからこそ可能な手法であると言える．

## 6.2 アドレス衝突確率の最小化

前節で述べたように，random-address では，スレッド移動時に移動先プロセスでアドレスが衝突した場合には，そのプロセスが存在するノード上に新しいプロセスを生成することによってスレッドを移動させる．しかし，一般論として，スレッド間のデータ共有の方がプロセス間のデータ共有よりも高速であり，協調動作するインスタンスはプロセスとして実装するよりもスレッドとして実装の方が性能上望ましいことをふまえると，アドレスの衝突を理由として同一ノード内に多数のプロセスを生成することは性能上不利である．

DMI に即して言えば，3.1 節で述べたように，DMI ではプロセスを単位としてグローバルメモリのコンシステンシ管理を行っているため，1 ノード内のプロセス数が増えると性能が劣化してしまう．具体的には，第一に，各プロセスにつき，他ノードからのメッセージを受信する receiver スレッド，それらのメッセージを処理する handler スレッド，ページの追い出しを担当する sweeper スレッドなどの多数の管理用スレッドが存在している．よって，1 ノード内のプロセス数を増やせば，1 ノード内に存在する管理用スレッドも増えてしまい，

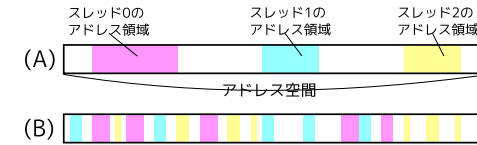


図7 アドレス領域の連続的な使用と離散的な使用 (A) 連続的な使用 (B) 離散的な使用．

計算本体を行うスレッドの性能が劣化してしまう．第二に，スレッド  $i$  とスレッド  $j$  が同一のプロセスに属していればスレッド  $i$  とスレッド  $j$  とでページのキャッシュを共有できるのに対して，別のプロセスに属している場合にはページのキャッシュを共有できないためである．たとえば，スレッド  $i \rightarrow$  スレッド  $j$  の順序でページ  $p$  をキャッシュしようとする場合，スレッド  $i$  とスレッド  $j$  が同一のプロセスに属している場合には，ページフォルトは1回で済むのに対して，別のプロセスに属している場合には，ページフォルトが2回発生してしまう．このように，1 ノード内のプロセス数が増えると性能が劣化するため，あるノード内に  $m$  個のスレッドを生成する場合，1 個のプロセス内に  $m$  個すべてのスレッドを格納するのが性能上もっとも望ましい．

以上の観察より，同一ノード内のプロセス数を増やさないようにするために，できるだけアドレス衝突を起きにくくする手法が必要だと言える．すなわち，random-address においては，乱数を使うとはいえ，本当にランダムにアドレスを割り当てるのではなく，スレッド移動時にアドレスが衝突する確率を最小化するための工夫が必要である．ここで考えるべき問題はおよそ以下である（問題の正確な定義は第 A.1 章で与える）：

問題  $m$  個のスレッド  $x_0, x_1, \dots, x_{m-1}$  を考える．各スレッド  $x_i$  は，自分以外のスレッドがどのアドレスにメモリを割り当てているか知らないとする．このとき，「どの2つの異なるスレッド  $x_i$  とスレッド  $x_j$  に対しても，スレッド  $x_i$  が使用するアドレス集合とスレッド  $x_j$  が使用するアドレス集合が共通部分を持たない確率」を最大にするためには，各スレッドがどのようなアドレス割り当ての戦略を採用すれば良いか？

そして，上記の問題に対する最適な戦略の1つは以下であることが証明できる（証明は第 A.1 章で与える）：

最適な戦略（の1つ） 各スレッド  $x_i$  は，できるかぎりアドレスを連続的に使用する

上記の意味は，「各スレッドが離散的にランダムにアドレスを使用するよりも，連続的にアドレスを使用する方がスレッド移動時のアドレス衝突確率が小さい」ということである．たとえば，図7(A)のようにアドレスを使用する方が，図7(B)のようにアドレスを使



$Z_i$  : Set of region

```

when thread  $i$  is created:
   $Z_i \leftarrow \emptyset$ 
  return

when thread  $i$  mmapsize size bytes:
  foreach  $(ptr, length) \in Z_i$  do
     $p \leftarrow \text{fixed\_mmap}(ptr + length, size)$ 
    if  $p \neq \text{MAP\_FAILED}$  then
       $Z_i \leftarrow Z_i \setminus \{(ptr, length)\}$ 
       $\cup \{(ptr, length + size)\}$ 
      reduction( $Z_i$ )
      return
    endif
  endforeach
  while 1 do
     $addr \leftarrow \text{rand}(inf, sup)$ 
     $p \leftarrow \text{fixed\_mmap}(addr, size)$ 
    if  $p \neq \text{MAP\_FAILED}$  then
       $Z_i \leftarrow Z_i \cup \{(ptr, size)\}$ 
      reduction( $Z_i$ )
      return
    endif
  endwhile
  return

when thread  $i$  munmaps region  $(p, size)$ :
  munmap( $p, size$ )
  foreach  $(ptr, length) \in Z_i$  do
    if  $ptr == p$  and  $p + size == ptr + length$  then
       $Z_i \leftarrow Z_i \setminus \{(ptr, length)\}$ 
    else if  $ptr == p$  and  $p + size < ptr + length$  then
       $Z_i \leftarrow Z_i \setminus \{(ptr, length)\} \cup \{(p + size, length - size)\}$ 
    else if  $ptr < p$  and  $p + size == ptr + length$  then
       $Z_i \leftarrow Z_i \setminus \{(ptr, length)\} \cup \{(ptr, length - size)\}$ 
    else if  $ptr < p$  and  $p + size < ptr + length$  then
       $Z_i \leftarrow Z_i \setminus \{(ptr, length)\} \cup \{(ptr, p - ptr)\}$ 
       $\cup \{(p + size, ptr + length - p - size)\}$ 
    endif
  endforeach
  return

when thread  $i$  is destroyed:
  foreach  $(ptr, length) \in Z_i$  do
    munmap( $ptr, length$ )
  endforeach
   $Z_i \leftarrow \emptyset$ 
  return

```

図 8 random-address における各プロセスのアドレス管理のアルゴリズム。(ptr, length) は、アドレス ptr から始まる length バイトの連続領域を表す。inf と sup は、それぞれ、 $stack_i^p$  と  $threadheap_i^p$  を割り当てるために使用することのできるアドレス空間の下限値と上限値を表す。また、fixed\_mmap(addr, size) 関数は、「アドレス addr から size バイトが未使用であれば mmap し、使用中であれば MAP\_FAILED を返す」関数とする(7.3 節を参照)。reduction( $Z_i$ ) は、領域集合  $Z_i$  内の任意の 2 個以上の連続領域に関して、連続領域としてまとめられるものを 1 個の連続領域にまとめる関数とする。

用するよりスレッド移動時のアドレス衝突確率が小さい。

以上の観察に基づき、random-address では、各スレッドが使用するアドレスができるかぎり連続的になるように、図 8 に示すアルゴリズムに従って、各スレッドが使用するアドレス領域 ( $stack_i^p$  と  $threadheap_i^p$ ) を管理する。

### 6.3 random-address の改善

なお、各スレッドが使用するメモリ量を予測できるならば、さらに random-address を改善し、スレッド移動時のアドレス衝突確率を下げるができる。

前節の議論で導いた、「各スレッド  $x_i$  はアドレスを連続的に使用する」という戦略においては、各スレッド  $x_i$  は任意のアドレスから始めて連続的なアドレス領域を使用することができる。言い換えると、この戦略では、各スレッド  $x_i$  が使用するアドレス領域は 1 の整数

倍にしかアラインされない。ところが、この戦略に加えて、「各スレッド  $x_i$  が使用するアドレス領域は、align の整数倍にアラインされなければならない」という制約を導入し、align の値を適切に調整することによって、さらにアドレス衝突確率を下げるができる。この理由は、直観的には、各スレッド  $x_i$  が使用するアドレス領域をある程度大きな数の整数倍にアラインさせることによって、2 つのスレッドのアドレス領域のごく一部だけが衝突することに起因するアドレス衝突が起きにくくなるためである。

ここで問題となるのは、align の最適値である。align の値を 1 から増やしていく場合のアドレス衝突確率を定性的に考えると、ある一定の値まではアドレス衝突確率は下がるが、align の値が各スレッド  $x_i$  の使用するメモリ量よりも十分大きくなってしまうと、逆にアドレス衝突確率は上がってしまうことがわかる。すなわち、align には、各スレッド  $x_i$  が使用するメモリ量に依存した最適値が存在する。証明は省略するが、一般に、「すべてのスレッド  $x_i$  が  $b$  バイトを使用するならば、align =  $b$  のときにアドレス衝突確率が最小になる」ことが導ける。しかし、実際のプログラムにおいては、各スレッドが使用するメモリ量はスレッドごとに異なるうえ、そもそも各スレッドが使用するメモリ量を事前に知ることは難しい。したがって、実際には、各スレッドが使用するメモリ量を予測し、その予測に基づいて経験的に align の値を設定することが必要になる。

### 6.4 アドレス領域の管理

random-address では、アドレスが衝突した場合には新しいプロセスを生成して、その新しいプロセスの中にスレッドを移動させるが、当然ながら、このとき新しいプロセスの中へのスレッド移動は確実に成功しなければならない。言い換えると、生成された直後のプロセスが使用しているアドレス領域(以下、初期領域と呼ぶ)が、移動したいスレッドが使用しているアドレス領域(以下、移動領域と呼ぶ)と重なっていることは許されない。なぜなら、仮に重なってしまったとすると、その新たに生成されたプロセスが使用できないことになるため、さらにもう 1 個新しいプロセスを生成する必要があるため、この作業を繰り返すうちに無限個のプロセスを生成してしまう可能性があるためである。また、将来的に、スレッド移動を拡張して分散チェックポイント/リスタートを実装しようとした場合に、いくつ新しいプロセスを生成してもスレッドをリスタートできないという事態も起こりうる。以上をふまえて、本節では、初期領域と移動領域が重ならないことを保証できるようなアドレス領域の管理について考える。

第 5 章で述べたように、DMI では、アドレス領域全体を  $register_i^p$ ,  $stack_i^p$ ,  $threadheap_i^p$ ,  $static^p$ ,  $processheap^p$ ,  $dmi$  の 6 種類のアドレス領域に分類して管理する。そして、5.2 節



と 5.4 節の説明より、このうち移動領域に含まれる可能性があるのは、 $register_i^p$ 、 $stack_i^p$ 、 $threadheap_i^p$  の 3 つであり、初期領域に含まれる可能性があるのは、 $static^p$ 、 $processheap^p$ 、 $dmi$  の 3 つである。したがって、初期領域と移動領域が重ならないようにするためには、 $register_i^p$ 、 $stack_i^p$ 、 $threadheap_i^p$  のアドレス領域と  $static^p$ 、 $processheap^p$ 、 $dmi$  のアドレス領域が重ならないように管理すれば良い。そこで、DMI では、利用可能なアドレス領域全体（たとえば 64 ビットアーキテクチャであれば  $2^{47}$  バイト）を前半部分と後半部分の 2 つに分割し、 $register_i^p$ 、 $stack_i^p$ 、 $threadheap_i^p$  のアドレス領域に関しては前半部分に割り当て、 $processheap^p$ 、 $dmi$  のアドレス領域に関しては後半部分に割り当てるよう、アドレス領域を管理する。特に、前半部分に関しては、図 8 に示すアルゴリズムに従って、各スレッドの使用アドレス領域ができるかぎり連続的になるようアドレス領域を管理する。 $static^p$  に関しては、静的に確保されるアドレス領域であるため、そもそも DMI の処理系がアドレス領域を制御できる自由度はないが、 $static^p$  はどのアドレスに配置されたとしても移動領域と重なることはないので、問題にならない。なぜなら、同一アーキテクチャの実行環境において同一の（再配置可能でない）実行バイナリを実行するならば、 $static^p$  のアドレス領域の位置はすべてのプロセスで同一になるため、移動領域が  $static^p$  に重なることはありえないからである。なお、ここで考慮したアドレス領域以外にも、コードが配置されるアドレス領域などがあるが、それらのアドレス領域が配置される位置に関しても、同一の（再配置可能でない）実行バイナリを同一アーキテクチャの実行環境で実行するならば同一になるため、移動領域と重なることはない。

なお、前半部分と後半部分の大きさをどう配分するかに関しては最適化の余地があるが、現在の実装では単純にアドレス空間全体を 2 等分するよう実装している。以上で述べたアドレス管理の実装について第 7 章で述べる。

## 7. 実装

本節では、第 6 章で述べたスレッド移動および random-address によるアドレス管理を、ユーザレベルで実装する方法を説明する。

### 7.1 スレッドのチェックポイント/リスタート

DMI では、`ucontext_t`<sup>4)</sup> と呼ばれる、Linux においてユーザレベルでコンテキストスイッチを行う機構を利用して、スレッドのチェックポイント/リスタートを実装している。`ucontext_t` では、`makecontext(ucontext_t *ucp, void *func, ...)` 関数を呼び出すことで、新しいコンテキスト `ucp` を作成できる。このとき、コンテキスト `ucp` が初めてコン

テキストスイッチされたときに実行される関数を `func` として指定できるほか、コンテキスト `ucp` が使用するスタック領域も明示的に指定できる。また、`swapcontext(ucontext_t *oucp, ucontext_t *ucp)` 関数を呼び出すことで、この関数を呼び出した現在のコンテキストを `oucp` に保存し、別のコンテキスト `ucp` にコンテキストスイッチすることができる。スレッドのチェックポイント/リスタートの手順を示す：

- (1) `DMI_scheduler_create()` 関数が呼ばれることにより、あるプロセス  $p$  でスレッドが起動される。この時点におけるこのスレッドのコンテキストを  $c_0$  と名付ける。
- (2) コンテキスト  $c_0$  は、`makecontext(c1, DMI_thread, ...)` 関数を呼び出すことで、スタック領域を明示的に指定して新しいコンテキスト  $c_1$  を作る。この新しいコンテキストを  $c_1$  と名付ける。
- (3) コンテキスト  $c_0$  が、`swapcontext(c0, c1)` 関数を呼び出す。その結果、現在のコンテキスト  $c_0$  が  $c_0$  に保存され、コンテキスト  $c_1$  へのコンテキストスイッチが起きる。そして、ユーザプログラムに定義されている `DMI_thread()` 関数がコンテキスト  $c_1$  で実行される。
- (4) やがて、このスレッドに対するスレッド移動の必要が生じ、コンテキスト  $c_1$  で実行されているユーザプログラムから `DMI_yield()` 関数が呼び出されたとする。このとき、`DMI_yield()` 関数は内部で (3) で保存した  $c_0$  を指定して `swapcontext(c1, c0)` 関数を呼び出す。その結果、コンテキスト  $c_1$  が  $c_1$  に保存されたあと、コンテキスト  $c_0$  へのコンテキストスイッチが起き (3) においてコンテキスト  $c_0$  が呼び出した `swapcontext()` 関数が返る。
- (5) この時点で、コンテキスト  $c_1$  つまり `DMI_thread()` 関数のコンテキストが使用しているレジスタ領域およびスタック領域は  $c_1$  に保存されている。よって、コンテキスト  $c_0$  は、 $c_1$  そのものと、コンテキスト  $c_1$  のスタック領域、および（別に管理している）スレッドのヒープ領域の 3 つを、移動先のプロセス  $q$  に対して送信する。
- (6) レジスタ領域、スタック領域、スレッドのヒープ領域を受信した移動先プロセス  $q$  は、それら 3 つの領域を移動元プロセス  $p$  と同一のアドレス領域に割り当て可能かどうかを調べ、不可能ならば移動元プロセス  $p$  に対して失敗通知を返す。可能ならば、それら 3 つの領域をまったく同一のアドレス領域に割り当て、移動元プロセス  $p$  に対して成功通知を返す。
- (7) 移動先プロセス  $q$  は、移動元プロセス  $p$  から受信した  $c_1$  をそのまま指定して `swapcontext(c1, c0)` 関数を呼び出すことで (4) において移動元プロセス  $p$  がコン

テキスト  $c_1$  で呼び出した `swapcontext()` 関数が移動先プロセス  $q$  で返る。その結果、移動元プロセス  $p$  とまったく同一のコンテキスト  $c_1$  で、`DMI_yield()` 関数を返すことができ、リスタートが完了する。

- (8) 移動元プロセス  $p$  は、6 において移動先プロセス  $q$  が返す成功/失敗通知を待機する。成功通知が受信されれば、そのままスレッドを終了させる。失敗通知が受信された場合、移動先プロセス  $q$  が存在するノードに対して新しいプロセス  $q'$  を生成したあと、プロセス  $q'$  に対して再度スレッド移動を試みる。

## 7.2 システムコールのハイジャック

### 7.2.1 ハイジャックの必要性

6.4 節で述べたアドレス管理においては、利用可能なアドレス領域全体を前半部分と後半部分の 2 つに分割し、 $register_i^p$ 、 $stack_i^p$ 、 $threadheap_i^p$  のアドレス領域は前半部分に割り当て、 $processheap^p$ 、 $dmi$  のアドレス領域は後半部分に割り当てるとようなアドレス管理を行う必要がある。そのためには、DMI の処理系が、各アドレス領域の確保/解放の処理をハイジャックして、そのアドレス領域の確保/解放の処理を具体的にどのアドレスに対して行うかを明示的に制御する必要がある。

ここで、各アドレス領域に関して、アドレスの明示的な制御が可能かどうかを確認する。第一に、7.1 節で述べた実装では、 $register_i^p$  は `ucontext_t` 型の変数として扱い、 $stack_i^p$  は `makecontext()` 関数を呼び出すときに明示的に指定するので、DMI の処理系はこれらのアドレス領域を明示的に制御することができる。第二に、 $threadheap_i^p$  は、定義より、`DMI_thread_mmap()` 関数/`DMI_thread_munmap()` 関数/`DMI_thread_mremap()` 関数の呼び出しによって確保/解放されるアドレス領域であるから、DMI の処理系がそのアドレス領域を明示的に制御することができる。第三に、 $processheap^p$  と  $dmi$  は (`mmap()` 関数などを經由して) システムコールの `mmap()` 関数/`mremap()` 関数/`munmap()` 関数によって確保/解放されるアドレス領域である。よって、それらのシステムコールの関数が確保/解放するアドレス領域を明示的に制御するためには、何らかのレイヤーにおいてシステムコールをハイジャックする必要がある。具体的には、アドレス領域の確保/解放に関連するシステムコールである、`mmap()` 関数/`mremap()` 関数/`munmap()` 関数/`mprotect()` 関数/`brk()` 関数をハイジャックする。

### 7.2.2 どのレイヤーでハイジャックすべきか

どのレイヤーにおいてシステムコールをハイジャックすべきかを、`mmap()` 関数を例にして議論する。以降では、システムコールの `mmap()` 関数を `sys_mmap()` 関数、`libc` 共有ライ

ブラリの `mmap()` 関数を `libc_mmap()` 関数、ハイジャック後に実行したい `mmap()` 関数を `hijack_mmap()` 関数と表記する。

Linux カーネル 2.6 においては、C 言語の実行バイナリから `libc_mmap()` 関数が呼び出されたとき以下の動作が起きる<sup>67)</sup>：

- (1) 実行バイナリから呼び出された `libc_mmap()` 関数が動的リンクされており、かつその `libc_mmap()` 関数の呼び出しが初回ならば、`libc_mmap()` 関数のアドレスが検索される。
- (2) アドレスが求まったあとで、`libc_mmap()` 関数が呼び出される。
- (3) `libc_mmap()` 関数がシステムコールの `sys_mmap()` 関数を呼び出し、カーネルレベルに制御が移る。
- (4) システムコールテーブルが検索され、`sys_mmap()` 関数のアドレスが求められる。
- (5) このプロセスが `ptrace` のデバッグプロセスになっていれば、デバッガプロセスに対してシグナルを送り、システムコールが発行されたことを通知する。
- (6) `sys_mmap()` 関数の本体が実行される。
- (7) このプロセスが `ptrace` のデバッグプロセスになっていれば、デバッガプロセスに対してシグナルを送り、システムコールが完了したことを通知する。
- (8) ユーザレベルに制御が戻り、`libc_mmap()` 関数が返る。

上記の手順を観察すると、`sys_mmap()` 関数をハイジャックして代わりに `hijack_mmap()` 関数を実行させられると思われる場所は (1)(3)(4)(5) の 4 箇所ある。以下ではこの 4 箇所におけるハイジャックの概要と問題点を議論する。

- (1) におけるハイジャック 環境変数 `LD_PRELOAD` を使用することで共有ライブラリの検索順序を変更する手法である。これにより (1) の手順において `libc_mmap()` 関数が発見される前に `hijack_mmap()` 関数を発見させることができ、`hijack_mmap()` 関数を実行させることができる。しかし、この手法は `libc_mmap()` 関数が静的リンクされている場合には使えない。たとえば、`libc` 共有ライブラリの `malloc()` 関数は、内部で `libc_mmap()` 関数を呼び出すが、`libc` 共有ライブラリ自体は静的リンクされているならば、`libc` 共有ライブラリの `malloc()` 関数の内部で呼び出される `libc_mmap()` 関数を検知することはできず、`hijack_mmap()` 関数を実行させることができない。

- (3) におけるハイジャック `libc_mmap()` 関数のコードを書き換えることで、`sys_mmap()` 関数の代わりに `hijack_mmap()` 関数を実行させる手法である。静的に行う手法と動的に行う手法がある。静的に行う手法では、`libc_mmap()` 関数のソースコードを書き換

えてコンパイルした改造 libc 共有ライブラリを用意し、DMI を実行するときには、通常の libc 共有ライブラリではなく、改造 libc 共有ライブラリを使用するようにすれば良い。しかし、この手法は、libc 共有ライブラリが多様なアーキテクチャに対応して実装されているためにソースコードの変更箇所が広範囲に及び点、各実行環境ごとに改造 libc ライブラリを用意しなければならず移植性が低い点などの問題を抱える。そこで、DMI では、実装の容易化と移植性の向上のため、動的に libc 共有ライブラリのコードを書き換える手法を採用した。この手法の詳細と得失は次節で議論する。

- (4) におけるハイジャック カーネルモジュールを使用して、システムコールテーブル内の `sys_mmap()` 関数のアドレスを `hijack_mmap()` 関数のアドレスに書き換えることで、`hijack_mmap()` 関数を実行させる手法である。しかし、Linux カーネル 2.6 以降では、セキュリティ上の理由により、システムコールテーブルの先頭アドレスを示す変数 `sys_call_table` が `extern` されなくなったため、カーネルモジュールからシステムコールテーブルのエントリを操作することができない。そのため、この手法を使うためには、変数 `sys_call_table` を `extern` するようにカーネルを変更する必要があるが、カーネルの変更は移植性を落としてしまう問題がある。また、カーネルモジュールの実行には `root` 権限が必要な点も問題である。
- (5) におけるハイジャック `ptrace` を使用して該当プロセスをデバッグプロセスとして登録し、デバッグプロセスで発行されるすべてのシステムコールをデバッガプロセスから監視する手法である。デバッグプロセスで発行された `sys_mmap()` 関数をデバッガプロセスでフックしたあと、デバッガプロセスから `PTRACE_POKEUSER` を使ってデバッグプロセスのコード領域を書き換えたり、システムコールの引数レジスタを書き換えたりすることで、デバッグプロセスに `hijack_mmap()` 関数を実行させることができる。しかし、本来 `ptrace` はプロセスの監視用に作られており、DMI が利用している `pthread` を監視するには不十分な点が多いという問題がある。たとえば、`ptrace` によってデバッガプロセスでシステムコールをフックした時点で、どの `pthread` によって呼び出されたシステムコールかを容易に知る手段が存在しない。

以上で検討した手法はいずれも、システムコールをハイジャックできない場合が存在するか、または移植性が低いという点で問題がある。

### 7.2.3 動的に共有ライブラリのコードを書き換える

以上の観察をふまえて、DMI では、ほぼすべての場合にシステムコールをハイジャックでき、かつ移植性の高い手法として、libc 共有ライブラリのコードを動的に変更する手法

を提案する。この手法では、`libc_mmap()` 関数が呼び出されたときに `hijack_mmap()` 関数が呼び出されるよう、DMI のプログラムの実行が開始された直後に libc 共有ライブラリのコード領域を書き換える。

第一に、基本アイデアを説明する。まず、`libc_mmap()` 関数の先頭に、`hijack_mmap()` 関数のアドレスへのジャンプ命令を挿入する。これによって、`libc_mmap()` 関数が呼び出されたとき、つねに `hijack_mmap()` 関数を実行できるようになる。しかし、これだけでは不十分で、`hijack_mmap()` 関数が DMI のアドレス管理のための一連の処理を行ったあと、実際にアドレス領域を割り当てようとして `libc_mmap()` 関数を呼び出すと、再度 `hijack_mmap()` 関数が呼び出され、無限に再帰してしまう。そこで、本来の `libc_mmap()` 関数も呼び出せるようにするため、本来の `libc_mmap()` 関数を呼び出すためのエントリポイントを `true_mmap()` 関数という名前で作成しておく。これにより、`hijack_mmap()` 関数は、`true_mmap()` 関数を呼び出すことで実際にアドレス領域を確保できるようになる。まとめると、

- (1) C 言語の実行バイナリが `libc_mmap()` 関数を呼び出す。
- (2) `libc_mmap()` 関数はすぐに `hijack_mmap()` 関数へジャンプする。
- (3) `hijack_mmap()` 関数がアドレス領域を確保するときは `true_mmap()` 関数を呼び出す。
- (4) `true_mmap()` 関数は本来の `libc_mmap()` 関数のエントリポイントであり、この内部で `sys_mmap()` 関数が呼び出されることで、実際のアドレス領域の確保が実現される。という順序で関数が呼び出されるように、libc 共有ライブラリのコード領域を書き換える。

第二に、実装について説明する：

- (1) `libc_mmap()` 関数と `hijack_mmap()` 関数の先頭アドレスを、それぞれ `libc_mmap`、`hijack_mmap` とおく (図 9 (A))。
- (2) 「 $x \geq injection$  であり、かつ `libc_mmap` から  $x$  バイト先のアドレスがちょうど命令境界になっている」条件を満たす  $x$  のうち最小の  $x$  を  $x_0$  とする。ここで、`injection` は (4) においてこの位置に挿入したいアセンブリコードのバイト数であり、x86-64 アーキテクチャの場合 12 バイトである。`libc_mmap` から何バイト目が命令境界になっているかは、`objdump` コマンドなどを使用して libc 共有ライブラリを逆アセンブルすることで調べる。たとえば、libc-2.3.6 では  $x_0 = 12$  であり、libc-2.7 では  $x_0 = 14$  である。空いている適当なアドレス領域に  $x_0 + injection2$  バイトを確保し、その先頭アドレスを `true_mmap` とする。ここで、`injection2` は (5) においてこの位置に挿入したいアセンブリコードのバイト数であり、x86-64 アーキテクチャの場合 13 バイトである。最終的には、この位置に `true_mmap()` 関数のエントリポイントを作成す

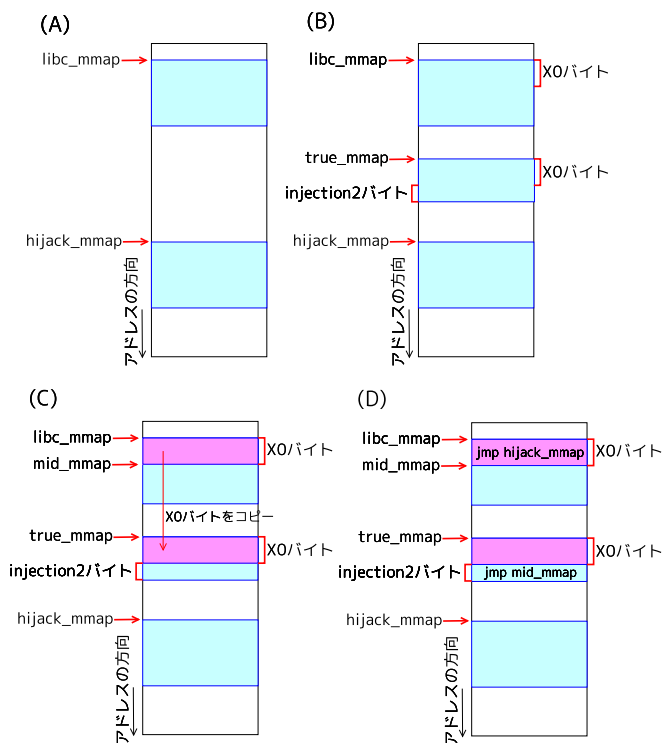


図 9 共有ライブラリのコード領域を動的に書き換える手順 .

る ( 図 9 ( B ) ) .

- (3) アドレス領域  $[libcc\_mmap, libcc\_mmap + x_0)$  を , アドレス領域  $[true\_mmap, true\_mmap + x_0)$  にコピーする . ここで , アドレス  $libcc\_mmap + x_0$  を `mid\_mmap` とおく ( 図 9 ( C ) ) .
- (4) `libcc\_mmap()` 関数が呼ばれた直後に `hijack\_mmap()` 関数に処理を飛ばすため , 以下のアセンブリコード ( サイズを `injection` バイトとする ) をアドレス領域  $[libcc\_mmap, libcc\_mmap + injection)$  に挿入する :

```
mov    hijack_mmap, %rax
jmpq   *%rax
```

これにより , `libcc\_mmap()` 関数が呼び出されると , そのままの引数で `hijack\_mmap()`

関数が呼び出されるようになる .

- (5) アドレス `true\_mmap` の位置に , 本来存在していた `libcc\_mmap()` へのエントリポイントを作るために , 以下のアセンブリコード ( サイズを `injection2` バイトとする ) をアドレス領域  $[true\_mmap + x_0, true\_mmap + x_0 + injection2)$  に挿入する ( 図 9 ( D ) ) :

```
mov    mid_mmap, %r11
jmpq   *%r11
```

これにより , `true\_mmap()` 関数が呼び出されると , もともとアドレス領域  $[libcc\_mmap, libcc\_mmap + x_0)$  の位置に配置されていたアセンブリコードが実行されたあと , アドレス `mid\_mmap` へとジャンプし , アドレス  $libcc\_mmap + x_0$  からアセンブリコードが実行されることになる . すなわち , `true\_mmap()` 関数を呼び出すことで , 本来存在していた `libcc\_mmap()` 関数を呼び出すことができる . なお , レジスタ `rax` や `r11` を使い分けている理由は , システムコール発行時のコーリングコンベンションに基づき , その時点で使用されていないレジスタを使用するためである .

以上の処理を , `main()` 関数が実行される前に実行することにより , プログラム内で発行されるすべての `libcc\_mmap()` をハイジャックでき , `processheapp` と `dmi` のアドレス領域を明示的に制御することができる .

なお , この手法では , `libc` 共有ライブラリをハイジャックしているだけであって , システムコールそのものをハイジャックしているわけではない . したがって , アセンブリコードで直接 `sys_mmap()` 関数を呼び出したり , `libc` 共有ライブラリの `syscall()` 関数から直接 `sys_mmap()` 関数を呼び出す場合など , `libcc\_mmap()` 関数を經由せずに呼び出される `sys_mmap()` 関数はハイジャックできないという欠点がある .

### 7.3 指定したアドレス領域の mmap

ここまでの議論では , `registerip` , `stackip` , `threadheapip` , `staticp` , `processheapp` , `dmi` の 6 種類の各アドレス領域の確保/解放の処理をハイジャックする手法を述べた . よって , あとは DMI のアドレス管理に基づいてアドレス領域を明示的に制御すれば良いが , これを実装するには , 当然「指定したアドレスに安全にアドレス領域を割り当てる」関数が必要である . 具体的には , アドレス `addr` とサイズ `size` を指定したとき「アドレス領域  $[addr, addr + size)$  が使用されていないならばアドレス領域  $[addr, addr + size)$  を割り当て , すでに使用されているならば `MAP_FAILED` を返す」仕様の関数 ( 図 8 における `fixed_mmap()` 関数 ) が必要である . よって , 本節では , カーネルを変更することなくユーザレベルでそのような関数を実装するアルゴリズムを説明する .

```

00: void* fixed_mmap(void *start, size_t length, int prot, int flags)
01: {
02:     void *ptr;
03:     int32_t ret;
04:     pthread_mutex_t mutex = PTHREAD_MUTEX_INITIALIZER;
05:
06:     pthread_mutex_lock(&mutex); /* ロック */
07:     errno = 0;
08:     ptr = true_mremap(start, PAGESIZE, PAGESIZE * 2, 0); /* start からの 1 ページを, 2 ページへと拡張しようとする */
09:     if(ptr == MAP_FAILED && errno == EFAULT) {
10:         ptr = true_mmap(start, PAGESIZE, prot, flags | MAP_FIXED, -1, 0); /* 必ず成功する */
11:         assert(ptr != MAP_FAILED);
12:         ptr = true_mremap(start, PAGESIZE, length, 0); /* start からの 1 ページを, length バイトへと拡張しようとする */
13:     }
14:     if(ptr == MAP_FAILED) { /* 失敗したら */
15:         ret = true_munmap(start, PAGESIZE); /* 後片付け */
16:         assert(ret == 0);
17:     }
18:     } else if(ptr != MAP_FAILED) { /* 最初の mremap に成功してしまったら */
19:         ret = true_munmap(start + PAGESIZE, PAGESIZE); /* 後片付け */
20:         assert(ret == 0);
21:         ptr = MAP_FAILED; /* fixed_mmap() は失敗 */
22:     }
23:     pthread_mutex_unlock(&mutex);
24:     return ptr;
25: }

```

図 10 fixed\_mmap() 関数のアルゴリズム .

まず, libc 共有ライブラリの仕様<sup>4)</sup> に基づけば, libc\_mmap(void \*start, size\_t length, int prot, int flags, int fd, off\_t offset) 関数だけでは, 意図するアドレスに安全にアドレス領域を割り当てることはできないことを確認する. 第一の方法として, 確保したいアドレス addr を第一引数 start に指定する方法があるが, libc\_mmap() 関数の仕様によれば, start はあくまでもヒントとして使用されるだけであり, 仮に指定したアドレス領域 [start, start + length) が使用されていなくとも, アドレス領域 [start, start + length) が確保できる保証はない. 実際, 8.3 節の環境では, これが起きることを確認している. 第二の方法として, libc\_mmap() 関数の第四引数の flags に MAP\_FIXED オプションを指定する方法を使えば, つねに指定したアドレス領域 [start, start + length) を確保することができる. しかし, アドレス領域 [start, start + length) がすでに使用されている場合には上書き確保されてしまう仕様であるため, この方法も安全ではない.

以上の動機をふまえて, DMI では, 図 10 に示すアルゴリズムで fixed\_mmap() 関数を実装している. 図 10 では, まず 8 行目で, true\_mremap() 関数を呼び出し, アドレス start から始まる 1 ページを 2 ページへと拡張しようとする. このとき, アドレス start から始

まる 2 ページに関して (1 ページ目の状態, 2 ページ目の状態) の組合せとしては以下の 5 通りの場合が考えられるが, libc\_mremap() 関数の仕様<sup>4)</sup> によれば, この各場合に対して, libc\_mremap() 関数は以下の値を返す:

(未使用, 未使用) 戻り値は MAP\_FAILED, errno は EFAULT .

(未使用, 使用中) 戻り値は MAP\_FAILED, errno は ENOMEM .

(使用中, 未使用) 戻り値は start, errno は設定されない .

(使用中, 1 ページ目と同じ属性で使用中) 戻り値は MAP\_FAILED, errno は ENOMEM .

(使用中, 1 ページ目とは異なる属性で使用中) 戻り値は MAP\_FAILED, errno は ENOMEM .

したがって, 9 行目の if 文の条件を満たすのは (未使用, 未使用) の場合のみである. つづいて 10 行目では, 1 ページ目に対して true\_mmap() 関数が MAP\_FIXED オプション付きで発行され, 1 ページ目のアドレス領域が確実に割り当てられる. いまは 1 ページ目が未使用であることが 9 行目の if 文により保証されているため, 10 行目の true\_mmap() 関数によって既存のアドレス領域が破壊されることはない. ここまで (使用中, 未使用) の状態に変化する. 次に 12 行目の true\_mremap() 関数により, いま割り当てた 1 ページ目のアドレス領域を length バイトにリサイズすることを試みる. libc\_mremap() 関数の仕様によれば, これが成功するのは, アドレス領域 [start, start + length) が使用されていない場合のみである. そして, これはいま実現すべき fixed\_mmap() 関数の仕様に他ならない.

## 8. 性能評価

### 8.1 シミュレーションによる random-address の評価

random-address のアルゴリズムを評価するため, シミュレーションによって, さまざまなパラメータに対するスレッド移動時のアドレス衝突確率を調べた. シミュレーションを用いた理由は, 実際にスレッド移動を行って評価すると時間がかかりすぎるうえ, 評価できるスレッド数やメモリ量の規模が実際のマシンの資源量に制限されてしまうためである. 本シミュレーションでは, 乱数として, 原始多項式  $x^{521} + x^{32} + 1$  に基づく 64 ビットの M 系列乱数を使用した.

ここで, 次のような状況を考える. 利用可能なアドレス空間の大きさを  $2^{\text{address}}$  バイトとし, 系内に process 個のプロセスが存在するとする. また, 各プロセスは合計 memory バイトのローカルメモリを使用しているとする. 言い換えると, 各プロセス  $p$  に関して, そのプロセス  $p$  内に存在するすべてのスレッド  $i$  が使用している  $\text{register}_i^p$  と  $\text{stack}_i^p$  と  $\text{threadheap}_i^p$  のメモリ量の総和が memory バイトであるとする. さらに, 各プロセスは合計 memory バ

イトのローカルメモリを割り当てるために、大きさが等しい  $chunk$  個の離散化されたアドレス領域を使用しているとする。言い換えると、各プロセスが使用するアドレス領域は、アドレス空間上で  $chunk$  個の小アドレス領域に分かれており、この各小アドレス領域の大きさはすべて  $memory/chunk$  バイトであるとする。 $chunk = 1$  の場合が、アドレス領域を連続的に使用する場合に相当する。そして、この  $chunk$  個の小アドレス領域の先頭アドレスは、 $align$  の整数倍の値の中からランダムに選ばれるとする。

本シミュレーションでは、以上のような状況において、すべてのプロセスに含まれるすべてのスレッドを、ある 1 個のノード  $P$  の中へとスレッド移動させるとき、ノード  $P$  の中に何個のプロセスが生成されるかを、さまざまな  $address, process, memory, chunk, align$  の値に対して測定した。以下では、このとき生成されるプロセス数を、 $address, process, memory, chunk, align$  の関数として、 $N(address, process, memory, chunk, align)$  と表す。当然、 $N(address, process, memory, chunk, align)$  が小さいほどアドレス衝突確率が小さいことを意味し、 $N(address, process, memory, chunk, align)$  が大きいほどアドレス衝突確率が大きいことを意味する。なお、すべてのプロセスに含まれるすべてのスレッドを、ノード  $P$  の中へとスレッド移動させるとき、これらのスレッドをどのような順番でノード  $P$  へ移動させるかに応じて、ノード  $P$  に生成されるプロセス数は変化しうるが、本シミュレーションではランダムな順序ですべてのスレッドを移動させた。このようなシミュレーションを、 $(address, process, memory, chunk, align)$  の各組合せに対して 10 回行い、測定された 10 個の値の最小値、最大値、平均値を算出した。この最小値、最大値、平均値をそれぞれ、 $N_{\min}(address, process, memory, chunk, align)$ 、 $N_{\max}(address, process, memory, chunk, align)$ 、 $N_{\text{avg}}(address, process, memory, chunk, align)$  と表す。

図 11 から図 20 までに、シミュレーション結果のグラフを示す。これらのすべてのグラフでは  $align = 1$  としている。また、各グラフが 1 個の  $chunk$  の値に対応しており、各グラフ中の 1 個の折れ線が 1 個の  $process$  の値に対応している\*1。各折れ線中の 1 個の点は、 $N_{\text{avg}}(address, process, memory, chunk, align)$  をプロットしており、その点に対するエラーバーの下限が  $N_{\min}(address, process, memory, chunk, align)$  を、エラーバーの上限が  $N_{\max}(address, process, memory, chunk, align)$  をプロットしている。具体例を挙げると、図 11 のグラフにおいてもっとも右下の赤い点は、「 $P$  アドレス空間全体が  $2^{32}$  バイトの環境

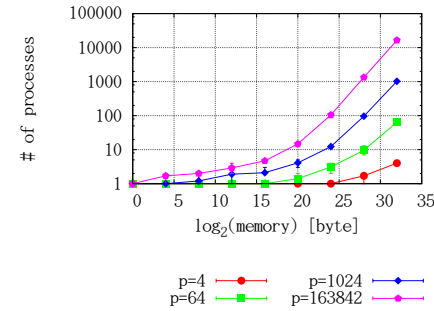


図 11  $N(32, process, memory, 1, 1)$  .

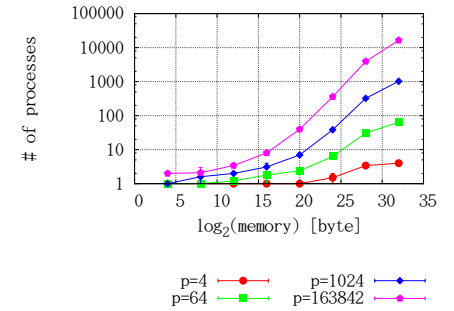


図 12  $N(32, process, memory, 16, 1)$  .

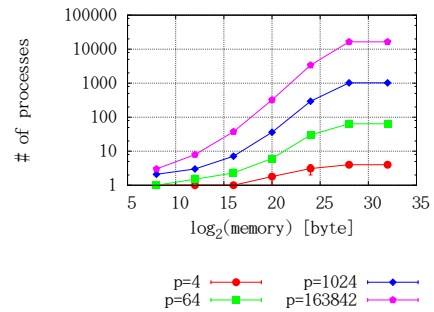


図 13  $N(32, process, memory, 256, 1)$  .

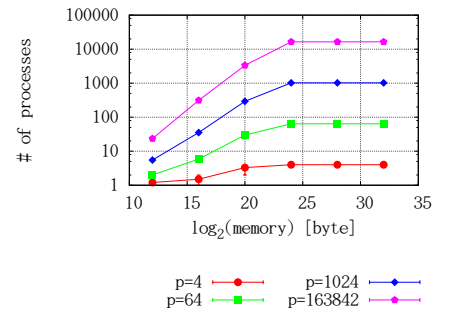


図 14  $N(32, process, memory, 4096, 1)$  .

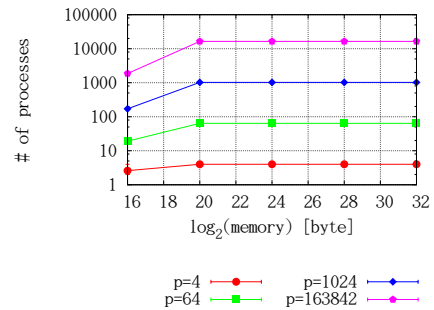


図 15  $N(32, process, memory, 65536, 1)$  .

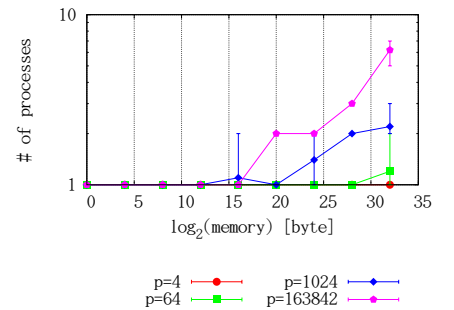


図 16  $N(47, process, memory, 1, 1)$  .

\*1 グラフ中では  $process$  を  $p$  と略記している。



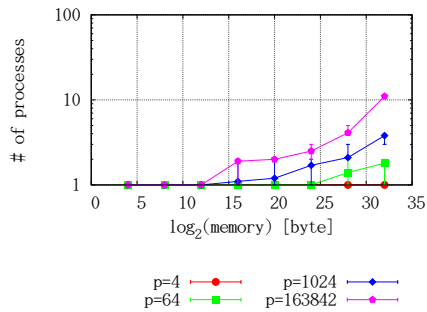


図 17  $N(47, process, memory, 16, 1)$ .

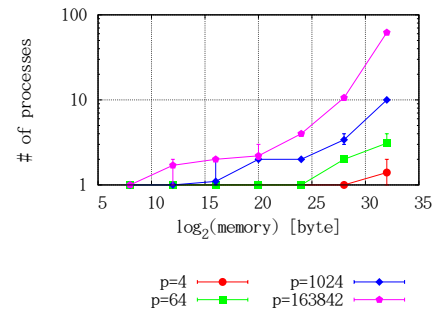


図 18  $N(47, process, memory, 256, 1)$ .

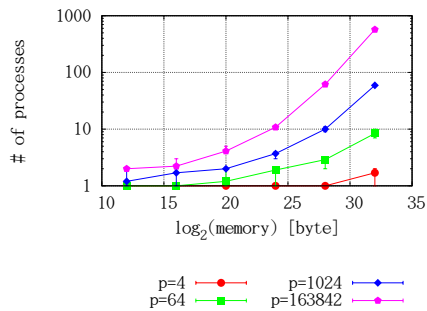


図 19  $N(47, process, memory, 4096, 1)$ .

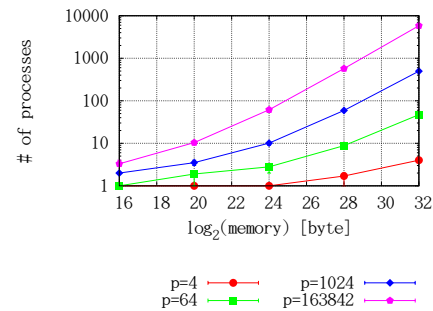


図 20  $N(47, process, memory, 65536, 1)$ .

に 4 個のプロセスがあり、各プロセスが  $2^{32}$  バイトを使用している。また、各プロセスはその  $2^{32}$  バイトを ( $chunk = 1$  なので) 1 個の連続的なアドレス領域として確保している。さらに、その連続的なアドレス領域の先頭アドレスはランダムに決まっている』という状況において、これら 4 個のプロセスに存在するすべてのスレッドを、1 個のノードにランダムな順序で詰め込むとき、そのノードに生成されたプロセス数の平均値を表している。ただし、いまの場合「 $2^{32}$  バイトのアドレス空間から、 $2^{32}$  バイトの 1 個の連続的なアドレス領域を確保する方法」はアドレス領域  $[0, 2^{32})$  を確保する方法の 1 通りしか存在せず、ランダム性はない。そして、アドレス領域  $[0, 2^{32})$  を使用している 4 個のプロセスに関して、これらのプロセスに存在するすべてのスレッドを  $2^{32}$  バイトのアドレス空間を持つノードに

移動させたとすれば、当然、生成されるプロセス数は必ず 4 個になる。したがって、図 11 のグラフにおいてもっとも右下の赤い点の値は 4 であり、エラーバーの下限も上限も 4 になっている。なお、6.4 節で述べたように、各プロセス内では各スレッドが使用するアドレス領域が重ならないようなアドレス管理が行われるため、各プロセス内に何個のスレッドが存在するかは問題にならない。また、 $chunk$  の値は  $memory$  の値以下である必要があるため、図 11 から図 20 のグラフでは、 $chunk$  の値が大きくなるにつれて、折れ線の左側が存在しなくなっている。 $addr$  として  $2^{32}$  と  $2^{47}$  を使用しているのは、それぞれ、既存の多くの 32 ビットアーキテクチャと 64 ビットアーキテクチャで使用可能なアドレス空間をシミュレートするためである。

このシミュレーション結果より、以下の事実が読み取れる：

- (1)  $align$  と  $memory$  と  $chunk$  と  $process$  を固定したとき、 $addr$  が増加するほどアドレス衝突確率が小さい。
- (2)  $address$  と  $align$  と  $chunk$  と  $process$  を固定したとき、 $memory$  が増加するほどアドレス衝突確率が大きい。
- (3)  $address$  と  $align$  と  $chunk$  と  $memory$  を固定したとき、 $process$  が増加するほどアドレス衝突確率が大きい。
- (4)  $address$  と  $align$  と  $memory$  と  $process$  を固定したとき、 $chunk$  が増加するほどアドレス衝突確率が大きい。

ここで (1) と (2) と (3) は定性的に考えて明らかである。一方 (4) は、「各プロセスのアドレス領域が連続的に使用されるとき、アドレス衝突確率がもっとも小さく、各プロセスが使用するアドレス領域が離散化するにしたがって、アドレス衝突確率が大きくなる」ことを述べている。すなわち、この結果は、6.2 節で述べた random-address の戦略が確かに最適であることを裏付けている。なお、random-address の最適性の正確な証明は第 A.1 章で行う。

また、図 20 より、以下の定量的な事実が読み取れる：

- $2^{47}$  バイトのアドレス空間において random-address を用いれば、65536 個のプロセスがそれぞれ  $2^{32}$  バイトのローカルメモリを割り当てたととしても、これらすべてのプロセスはわずか平均 4 個のアドレス空間に詰め込むことができる。実際のスレッドスケジューリングでは、利用可能な計算資源が急激に減少しないかぎり、すべてのスレッドを同一プロセスに詰め込むことは起きえないことをふまえると、 $2^{47}$  バイトのアドレス空間で 65536 個のプロセスを生成するような状況では、アドレス衝突はほぼ起きない

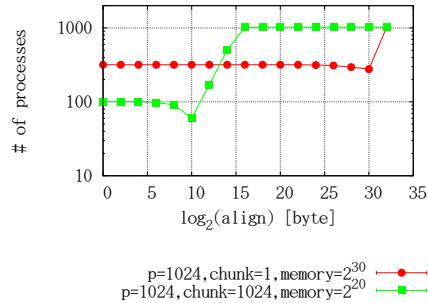


図 21  $N(32, 1024, 2^{30}, 1, align)$  と  $N(32, 1024, 2^{20}, 1024, align)$ .

と考えられる．

- $2^{47}$  バイトのアドレス空間において 65536 個のプロセスがそれぞれ  $2^{32}$  バイトのローカルメモリを割り当てるとき、各プロセスがそのローカルメモリを 16384 個の離散的な小アドレス領域として割り当てると、これらすべてのプロセスを詰め込むに必要なアドレス空間の数は平均 5828 個にもなる（図 20 のグラフにおけるもっとも右上の点）．この結果より、アドレスを連続的に割り当てることで、アドレス衝突確率を大幅に低減できることがわかる．

次に、6.3 節で述べた、 $align$  を大きくすることによってアドレス衝突確率が下がることの効果を調べる．図 21 には、以下の 2 つの折れ線をプロットしている：

- $align$  の値を  $2^0$  バイトから  $2^{32}$  バイトまで変化させたときの、 $N_{\min}(32, 1024, 2^{30}, 1, align)$  と  $N_{\max}(32, 1024, 2^{30}, 1, align)$  と  $N_{\text{avg}}(32, 1024, 2^{30}, 1, align)$  ．
- $align$  の値を  $2^0$  バイトから  $2^{32}$  バイトまで変化させたときの、 $N_{\min}(32, 1024, 2^{20}, 1024, align)$  と  $N_{\max}(32, 1024, 2^{20}, 1024, align)$  と  $N_{\text{avg}}(32, 1024, 2^{20}, 1024, align)$  ．

図 21 より、 $N_{\text{avg}}(32, 1024, 2^{30}, 1, align)$  は  $align = 2^{30}$  において、 $N_{\text{avg}}(32, 1024, 2^{20}, 1024, align)$  は  $align = 2^{10}$  において極小値をとることがわかる．これは、 $N_{\text{avg}}(32, 1024, 2^{30}, 1, align)$  は  $align = 2^{30}$  では各小アドレス領域の大きさが  $memory/chunk = 2^{30}$  であり、 $N_{\text{avg}}(32, 1024, 2^{20}, 1024, align)$  では  $memory/chunk = 2^{10}$  であるためである．以上の結果は、6.3 節における考察に合致する．

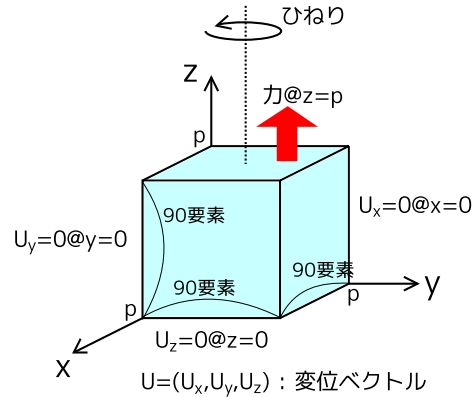


図 22 有限要素法による応力解析．

図 8 に示した random-address におけるアドレス管理のアルゴリズムでは、本シミュレーションにおける各小アドレス領域が、ほぼ、1 個のプロセス内に存在する各スレッドが使用する連続的なアドレス領域に相当する．したがって、各スレッドが使用するアドレス領域の大きさを事前に見積もれるならば、 $align$  の値をその値に設定するのが望ましい．

## 8.2 アプリケーションベンチマーク

### 8.3 実験環境

実験環境としては、Intel Xeon E5530 2.40GHz (4 コア)  $\times$  2 の CPU、24GB のメモリ、カーネル 2.6.26-2-amd64 の Linux で構成されるマシン 16 ノードを、10Gbit イーサネット でネットワーク接続した、合計 128 プロセッサのクラスタ環境を用いた．以降の実験では、DMI/MPI を  $n$  プロセッサで実行する際には、8 個の (DMI の) スレッド/MPI プロセスを  $\lfloor n/8 \rfloor$  台のノードに立て、残りの  $n - 8 \times \lfloor n/8 \rfloor$  個の (DMI の) スレッド/MPI プロセスを別の 1 つのノードに立てるプロセス構成とした．コンパイラには gcc 4.3.2、MPI には OpenMPI 1.4.2 と mpich2-1.2.1p1、最適化オプションには -O3 を使用した．

#### 8.3.1 有限要素法

3 次元立方体物体に対して、図 22 に示すような境界条件を課したときの応力解析を有限要素法で行った．この有限要素法では、3 次元立方体が  $90^3$  個の要素に分割されており、各要素に対しては Sequential Gauss Algorithm に基づいて  $z$  軸回りに最大 200 度のひねりが加えられている．この有限要素法は、第 2 回クラスタシステム上の並列プログラミングコンテスト<sup>7)</sup> の題材として使用された、実際の工学に基づく実用的な並列科学技術計算である．非常に収束させづらい問題であるため、さまざまな実用的な工学的手法が必要となる．

アルゴリズムの概要は以下のとおりである．詳細は著者らの資料<sup>64)</sup> を参考にされたい：

- (1) 立方体物体全体をプロセッサ数個の領域に分割する．このとき、立方体領域を単純に直方体領域の集合として分割するのではなく、領域間オーバーラップやフィルインを考慮した非定型的な領域分割を行う．また、収束性を改善させるため、修正 RCM オーダリング<sup>40)</sup> によって領域内の節点 (各要素における 8 個のコーナー点のこと) を並び替える．
- (2) 与えられている節点間の結合関係を連立一次元方程式  $Ax = b$  として表現する．ここで  $A$  は疎行列になる．
- (3) 連立方程式  $Ax = b$  を、BiCGSafe 法<sup>28)</sup> と呼ばれる反復法を用いて解  $x$  が収束するまで反復計算する．各イテレーションの先頭では、収束性を改善するために、フィルインレベル 3 のブロック不完全 LU 分解による前処理と、領域間オーバーラップ 2 に

よる Restricted Additive Schwarz 法<sup>13)</sup> に基づく前処理を適用する。

(4) さらに、各イテレーションの先頭において、DMI\_yield() 関数を記述する。

スレッド移動に対応した DMI プログラムは、スレッド移動に非対応の通常の SPMD 型の DMI プログラムを変更することで実装したが、この際の変更点は (1) 各イテレーションの先頭で DMI\_yield() 関数を記述すること、(2) DMI\_yield() 関数をまたいで使用されるヒープ領域の確保/解放を DMI\_thread\_malloc() 関数/ DMI\_thread\_realloc() 関数/ DMI\_thread\_free() 関数に置き換えること、の 2 点のみだった。すなわち、既存の SPMD 型の並列科学技術計算に対して、わずかな変更を加えるだけで、計算規模を拡張/縮小可能なプログラムが得られることを確認できた。

まず、図 23 図 24 には、スレッド移動を行わない場合において、実行するプロセッサ数を変化させたときの DMI, mpich2, OpenMPI の実行時間とウィークスケラビリティを示す。なお、mpich2 と OpenMPI のプログラムとしては、第 2 回クラスタシステム上の並列プログラミングコンテストにおける優勝プログラムを使用した。図 23 図 24 より、DMI は mpich2 と同等で、OpenMPI よりも高い性能を達成していることがわかる。OpenMPI の性能の低さは、OpenMPI における send/recv の遅さに起因していることがわかっている。

次に、DMI において、128 個のスレッドを生成し (1) 最初はノード 0~ ノード 7 で実行し (合計 8 ノード, 64 プロセッサ) (2) 第 50 回目の反復処理終了直後にノード 8~ ノード 15 の 8 ノードを参加させ (合計 16 ノード, 128 プロセッサ) (3) 第 101 回目の反復処理終了直後にノード 0~ ノード 11 の 12 ノードを脱退させる (合計 4 ノード, 32 プロセッサ)、というように利用可能な計算資源を動的に増減させた。このときの各イテレーションの実行時間を図 25 に示す。なお、この実験では、各イテレーションの実行時間をある程度大きくするため、有限要素法の要素数を  $90^3$  ではなく  $150^3$  とした。図 25 より、利用可能な計算資源の増減に対応して計算規模を動的に拡張/縮小できていることが読み取れる。なお、合計 4 ノードで実行した場合の実行時間が揺れているが、これは、各プロセッサあたり 4 個のカーネルスレッドを割り当てているためカーネルによるスレッドスケジューリングがばらつくためである。測定の結果、この有限要素法では、各スレッドが 510MB~520MB のローカルメモリを消費し、グローバルメモリは 335MB 消費されていた。また (2) における 8 ノード参加時には、8 ノードの参加処理と 120 スレッドのスレッド移動 (=57.6GB のメモリ移動) が発生して 17.3 秒を要した。さらに (3) における 12 ノード脱退時には、12 ノードの脱退処理と 120 スレッドのスレッド移動 (=57.6GB のメモリ移動) が発生し 30.9 秒を要した。なお、8 ノードしか参加させていないにもかかわらず 120 スレッドもの

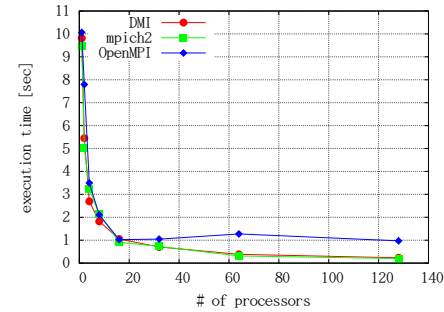


図 23 プロセッサ数を変化させた場合の、有限要素法の実行時間の変化。

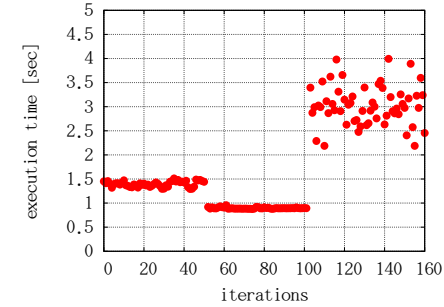


図 25 有限要素法において、計算規模を動的に拡張/縮小した場合における各イテレーションの実行時間。

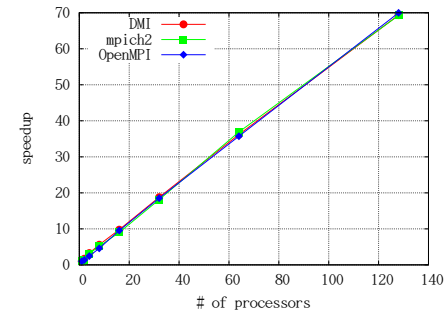


図 27 N 体問題のウィークスケラビリティ。

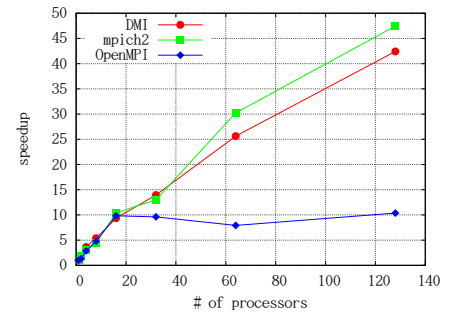


図 24 有限要素法のウィークスケラビリティ。

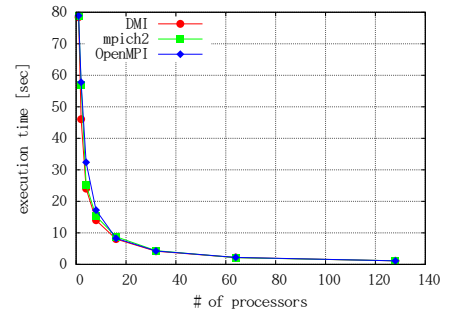


図 26 プロセッサ数を変化させた場合の、N 体問題の実行時間の変化。

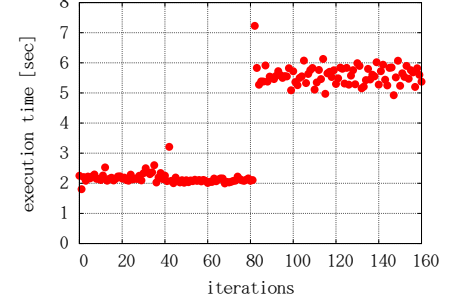


図 28 N 体問題において、計算規模を動的に拡張/縮小した場合における各イテレーションの実行時間。

スレッド移動が発生するのは、6.1 節で述べたように、現実装では、計算規模が変化するたびに、スレッド番号が若い方から順に、各プロセスに  $m/n$  個（または  $m/n + 1$  個）ずつスレッドを割り当てるようなスレッドスケジューリングを行っているためである。また、この実験では  $align = 1$  としたが、これらのスレッド移動に伴うアドレス衝突は発生しなかった。以上の結果より、DMI が、実用的な並列科学技術計算に対して、利用可能な計算資源の増減に対応して効果的に計算規模を動的に拡張/縮小できることを確認できた。

### 8.3.2 N 体問題

N 体問題を題材にして性能評価を行う。N 体問題のアルゴリズムの概要は以下のとおりである：

- (1) 3次元格子状に並んだ  $l_1 \times l_2 \times l_3$  個の粒子を考える。各粒子に位置と速度の情報を持たせ、初期的に適当な位置と初速度を与える。
- (2)  $n$  個のプロセッサで実行する場合、 $l_1 \times l_2 \times l_3$  個の直方体状に並んだ粒子を  $l_1$  の方向に  $n$  等分し、 $i$  番目のプロセッサには  $i$  番目の領域を担当させる。つまり、各プロセッサは  $l_1 \times l_2 \times l_3/n$  個の粒子を担当する。
- (3) 以下の処理を反復する：
  - (a) 各プロセッサが、そのプロセッサの担当範囲の粒子の位置の情報をすべてのプロセッサに対して送信する Allgather 型の通信を行うことで、すべてのプロセッサが全粒子の位置情報を把握する。
  - (b) 各プロセッサは、全粒子の位置情報に基づいて、そのプロセッサの担当範囲の粒子の位置と速度の情報を更新する。

まず、図 26 図 27 には、スレッド移動を行わない場合において、実行するプロセッサ数を変化させたときの DMI, mpich2, OpenMPI の実行時間とウィークスケールビリティを示す。この実験では、 $l_1 = l_2 = l_3 = 24$  とした。図 26 図 27 より、DMI も mpich2 も OpenMPI も同等のスケールビリティを達成していることがわかる。

次に、DMI において、128 個のスレッドを生成し、前節と同様にして (1) 最初はノード 0~ノード 7 で実行し (2) 第 41 回目の反復処理終了直後にノード 8~ノード 15 の 8 ノードを参加させ (3) 第 82 回目の反復処理終了直後にノード 0~ノード 11 の 12 ノードを脱退させる、というように利用可能な計算資源を動的に増減させたときの各イテレーションの実行時間を図 28 に示す。なお、この実験では、 $l_1 = 24$ ,  $l_2 = l_3 = 28$  とした。この N 体問題では、各スレッドが 9.7MB のローカルメモリを消費し、グローバルメモリは 441KB 消費されていた (2) における 8 ノード参加時には、8 ノードの参加処理と 120 スレッドの

スレッド移動 (=1.19GB のメモリ移動) が発生して 2.22 秒を要した。さらに (3) における 12 ノード脱退時には、12 ノードの脱退処理と 120 スレッドのスレッド移動 (=1.19GB のメモリ移動) が発生し 5.51 秒を要した。また、この実験では  $align = 1$  としたが、これらのスレッド移動に伴うアドレス衝突は発生しなかった。

## 9. 結 論

### 9.1 ま と め

本研究の貢献は以下のとおりである：

- 利用可能な計算資源が増減しうるクラウド環境において長時間を要する並列計算を実行するためには、計算規模を動的に拡張/縮小するような並列計算が必要であることを指摘した。また、そのような並列計算を簡単に記述するためには、プログラマには十分な並列性だけを記述させておき、あとは処理系が透過的に計算規模の拡張/縮小を行うような並列分散プログラミングモデルが必要であることを提案した。
- そのような並列分散プログラミングモデルを実現する、グローバルビュー型の PGAS モデルに基づいた並列分散プログラミング処理系として、DMI を提案して実装し、評価した。評価の結果、有限要素法を用いた応力解析という実用的な並列科学技術計算に対して (1) 既存の SPMD 型の並列プログラムに対する変更点がわずかであること、(2) 計算資源の増減に伴って動的に並列度を変化させられることを確認した。
- DMI の主要な要素技術として、スレッド移動技術を検討し、今後ますます増大する計算規模の拡張に伴い、従来のスレッド移動技術では限界に達する可能性を指摘したうえで、アドレス空間の大きさに制限されないスレッド移動を実現するためのアドレス管理のアルゴリズムとして random-address を提案した。また、random-address の最適性について数学的な証明を与えるとともに、シミュレーションによってその最適性を確認した。random-address は、DMI がプロセスの動的な参加/脱退に対応しているからこそ実現できるアドレス管理のアルゴリズムである。

著者らの知るかぎり、グローバルビュー型の PGAS モデルに基づいて、ユーザプログラムから透過的に計算規模の拡張/縮小を実現した研究は存在せず、DMI の全体的な設計および実装には新規性があると言える。また、今後の計算規模の拡張に伴って従来のスレッド移動技術では限界に達する可能性を指摘し、アドレス空間の大きさに制限されないスレッド移動技術を新たに提案している点にも新規性がある。

## 9.2 今後の課題

現時点では、6.1 節に示すような単純なスレッドスケジューリングしか行っていないが、DMI のように透過的にスレッドをスケジューリングする処理系においては、スレッドスケジューリングの最適化が性能上きわめて重要である。DMI においてスレッドスケジューリングを最適化するには、ノード間のスレッドの負荷バランス、スレッド移動に要する時間、プロセス数を増やすことによるオーバーヘッド、各スレッド間でのデータ共有の度合いなどの要素を総合的に考慮する必要がある。DMI では、API を通じてのみグローバルメモリにアクセスすることになっているため、各スレッド間のグローバルメモリの共有度合いを容易に把握できるほか、メモリ確保/解放関連のシステムコールをハイジャックしているため、各スレッドのローカルメモリの使用状況も把握できるなど、スレッドスケジューリングにとって必要な多くの情報を容易に取得できる。したがって、これらの情報に基づいた効率的なスレッドスケジューリングを検討する必要がある。

## 参 考 文 献

- 1) Amazon EC2 [Online]. <http://aws.amazon.com/ec2/>.
- 2) Google App Engine [Online]. <http://code.google.com/intl/appengine/>.
- 3) Google Docs [Online]. <http://docs.google.com/>.
- 4) Linux Manpages [Online]. <http://linuxmanpages.com/>.
- 5) Salesforce.com [Online]. <http://www.salesforce.com/>.
- 6) Windows Azure [Online]. <http://www.microsoft.com/windowsazure/>.
- 7) 第2回クラスタシステム上のプログラミングコンテスト [Online]. <https://www2.cc.u-tokyo.ac.jp/procon2009-2/>.
- 8) Christiana Amza, Alan L.Cox, Hya Dwarkadas, Pete Keleher, Honghui Lu, Weimin Yu Ramakrishnan Rajamony, and Willy Zwaenepoel. TreadMarks: Shared Memory Computing on Networks of Workstations. *IEEE Computer*, Vol.29, No.2, pp. 18–28, Feb 1996.
- 9) Gabriel Antoniu, Luc Bouge, and Raymond Namyst. An Efficient and Transparent Thread Migration Scheme in the PM2 Runtime System. *Proceedings of the 11 IPPS/SPDP'99 Workshops Held in Conjunction with the 13th International Parallel Processing Symposium and 10th Symposium on Parallel and Distributed Processing*, pp. 496–510, 1999.
- 10) Gabriel Antoniu and Christian Perez. Using Preemptive Thread Migration to Load-Balance Data-Parallel Applications. *Proceedings of the 5th International Euro-Par Conference on Parallel Processing*, pp. 117–124, 1999.
- 11) Rafik A.Salama and Ahmed Sameh. *Potential Performance Improvement of Collective Operations in UPC*. John von Neumann Institute for Computing, 2007.
- 12) Rajkumar Buyya, CheeShin Yeo, Srikumar Venugopal, James Broberg, and Ivona Brandic. Cloud Computing and Emerging IT Platforms: Vision, Hype, and Reality for Delivering Computing as the 5th Utility. *Future Generation Computer Systems*, Vol.25, pp. 599–616, 12 2008.
- 13) Xiao-Chuan Cai and Marcus Sarkis. A Restricted Additive Schwarz Preconditioner for General Sparse Linear Systems. *SIAM Journal on Scientific Computing*, Vol.21, No.2, pp. 792–797, 9 1999.
- 14) William Carlson, Thomas Sterling, Katherine Yelick, and Tarek El-Ghazawi. *UPC Distributed Shared Memory Programming*. WILEY INTER-SCIENCE, Jun 2005.
- 15) Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C.Hsieh, Deborah A.Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E.Gruber. Bigtable: A Distributed Storage System for Structured Data. *ACM Transactions on Computer Systems*, Vol.26, , Jun 2008.
- 16) V.Chaudhary and H.Jiang. Techniques for Migrating Computations on the Grid. *Engineering the Grid: Status and Perspective*, pp. 399–415, Jan 2006.
- 17) Po-Cheng Chen, Cheng-I Lin, Sheng-Wei Huang, Jyh-Biau Chang, Ce-Kuen Shieh, and Tyng-Yeu Liang. A Performance Study of Virtual Machine Migration vs. Thread Migration for Grid Systems. *International Conference on Advanced Information Networking and Applications*, Mar 2008.
- 18) Barton Christopher, Cascaval Clin, Almasi George, Zheng Yili, Farreras Montse, Chatterje Siddhartha, and AmaralJose Nelson. Shared memory programming for large scale machines. *ACM SIGPLAN Notices*, Vol.41, No.6, pp. 108–117, Jun 2006.
- 19) Christopher Clark, Keir Fraser, Steven Hand, Jacob Gorm Hansenf, Eric Julf, Christian Limpach, Ian Pratt, and Andrew Warfield. Live migration of virtual machines. *Proceedings of the 2nd conference on Symposium on Networked Systems Design and Implementation*, Vol.2, pp. 273–286, 2005.
- 20) Cristian Coarfa, Yuri Dotsenko, Jason Eckhardt, and John Mellor-crummey. Co-Array Fortran Performance and Potential: An NPB Experimental Study. *16th Internatitnal Workshop on Languages and Compilers for Parallel Computing*, Vol. 2958, pp. 177–193, 2004.
- 21) Cristian Coarfa, Yuri Dotsenko, John Mellor-Crummey, Francois Cantonnet, Tarek El-Ghazawi, Ashrujit Mohanty, YiYi Yao, , and Daniel Chavarria-Miranda. An Evaluation of Global Address Space Languages: Co-Array Fortran and Unified Parallel C. *10th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, pp. 36–47, 2005.
- 22) David Cronk, Matthew Haines, and Piyush Mehrotra. Thread Migration in the Presence of Pointers. *Proceedings of the 30th Hawaii International Conference on*

- System Sciences: Software Technology and Architecture*, Vol.1, pp. 292–302, 1997.
- 23) Kaushik Datta<sup>1</sup>, Dan Bonachea<sup>1</sup>, and Katherine Yelick<sup>1</sup>. Titanium Performance and Potential: An NPB Experimental Study. *18th International Workshop on Languages and Compilers for Parallel Computing*, Vol. 4339, pp. 200–214, 2006.
  - 24) Bozhidar Dimitrov and Vernon Reg. Arachne: a portable threads system supporting migrant threads on heterogeneous network farms. *IEEE Transactions on Parallel and Distributed Systems*, Vol.9, pp. 459–469, May 1998.
  - 25) Cong Du and Xian-HeSun Sun. MPI-Mitten: Enabling Migration Technology in MPI. *IEEE International Symposium on Cluster Computing and the Grid*, pp. 11–18, May 2006.
  - 26) Jason Duell. The design and implementation of Berkeley Lab’s linuxcheckpoint/restart. Technical report, Ernest Orlando Lawrence Berkeley National Laboratory, Apr 2005.
  - 27) Message Passing Interface Forum. MPI: A Message-Passing Interface Standard Version 2.2. Technical report, Message Passing Interface Forum, Sep 2009.
  - 28) Seiji Fujino, Maki Fujiwara, and Masahiro Yoshida. BiCGSafe Method Based on Minimization of Associate Residual (in Japanese). *Transactions of the Japan Society for Computational Engineering and Science*, Vol.8, pp. 145–152, 2006.
  - 29) Kentaro Hara and Kenjiro Taura. A Global Address Space Framework for Irregular Applications (accepted, short paper). *High Performance Distributed Computing*, Jun 2010.
  - 30) Paul Hilfinger, Dan Bonachea, David Gay, Susan Graham, Ben Liblit, Geoff Pike, and Katherine Yelick. Titanium Language Reference Manual. Technical report, University of California at Berkeley, Aug 2006.
  - 31) Chao Huang, Orion Lawlor, and L.V.Kale. Adaptive MPI. *International workshop on languages and compilers for parallel computing*, Vol. 2958, pp. 306–322, Oct 2003.
  - 32) Ayal Itzkovitz, Assaf Schuster, and Lea Shalev. Thread Migration and its Applications in Distributed Shared Memory Systems. *Journal of Systems and Software*, Vol.42, No.1, pp. 71–87, Jul 1998.
  - 33) Hai Jiang and Vipin Chaudhary. Compile/Run-Time Support for Thread Migration. *Proceedings of the 16th International Parallel and Distributed Processing Symposium*, pp. 58–66, 2002.
  - 34) Hai Jiang and Vipin Chaudhary. MigThread: Thread Migration in DSM Systems. *International Conference on Parallel Processing*, p. 581, 2002.
  - 35) Hai Jiang and Vipin Chaudhary. On Improving Thread Migration: Safety and Performance. *Proceedings of the 9th International Conference on High Performance Computing*, pp. 474–484, 2002.
  - 36) Hai Jiang and Vipin Chaudhary. Thread Migration/Checkpointing for Type-Unsafe C Programs. *International conference on high performance computing*, Vol. 2913, pp. 469–479, Nov 2003.
  - 37) Matchy J.M.Ma, Cho-Li Wang, and Francis C.M.Lau. Delta Execution: A preemptive Java thread migration mechanism. *Cluster Computing*, Vol.3, pp. 83–94, 2000.
  - 38) Jian Ke and Evan Speight. Tern: Thread Migration in an MPI Runtime Environment. Technical report, Cornell, Nov 2001.
  - 39) K.Thitikamol and P.Keleher. Thread migration and communication minimization in DSM systems. *Proceedings of the IEEE, Special Issue on Distributed Shared Memory*, Vol.87, pp. 487–497, 3 1999.
  - 40) Wai-Hung Liu and Andrew H.Sherman. Comparative Analysis of the Cuthill-McKee and the Reverse Cuthill-McKee Ordering Algorithms for Sparse Matrices. *SIAM Journal on Numerical Analysis*, Vol.13, No.2, pp. 198–213, Apr 1976.
  - 41) Kirk L.Johnson, M.Frans Kaashoek, and Deborah A.Wallach. CRL: High-Performance All-Software Distributed Shared Memory. *Proceedings of the Fifteenth Symposium on Operating Systems Principles*, Vol.29, No.5, pp. 213–228, Mar 1995.
  - 42) L.Vaquero, L.Rodero-Marino, J.Caceres, and M.Lindner. A Break in the Clouds : Towards a Cloud Definition. *SIGCOMM Computer Communication Review*, pp. 137–150, 2009.
  - 43) Armbrust M., A.Fox, R.Griffith, A.D.Joseph, R.Katz, A.Konwinski, G.Lee, D.A.Patterson, A.Rabkin, I.Stoica, and M.Zaharia. Above the Clouds: A Berkeley View of Cloud Computing. Technical report, UC Berkeley Reliable Adaptive Distributed Systems Laboratory, 2 2009.
  - 44) KaoutarEl Maghraoui, Boleslaw K.Szymanski, and Carlos Varela. An Architecture for Reconfigurable Iterative MPI Applications in Dynamic Environments. *International Conference on Parallel Processing and Applied Mathematics*, Vol. 3911, pp. 258–271, Sep 2006.
  - 45) Scott Milton. Thread Migration in Distributed Memory Multicomputers. Technical report, Australia National University, 1998.
  - 46) Frank Mueller. Distributed Shared-Memory Threads: DSM-Threads. *Workshop on Run-Time Systems for Parallel Programming*, pp. 31–40, Apr 1997.
  - 47) Frank Mueller. On the Design and Implementation of DSM-Threads. *Conference on Parallel and Distributed Processing Techniques and Applications*, pp. 315–324, Jun 1997.
  - 48) Jarek Nieplocha, Manojkumar Krishnan, Bruce Palmer, and Vinod Tipparaju. The Global Arrays User’s Manual. Technical report, Pacific Northwest National Laboratory, Jul 2009.



- 49) Robert W. Numrich and John Reid. Co-array Fortran for parallel programming. *ACM SIGPLAN Fortran Forum*, Vol.17, No.2, pp. 1–31, 1998.
- 50) Michael R.Hines and Kartik Gopalan. Post-copy based live virtual machine migration using adaptive pre-paging and dynamic self-ballooning. *Proceedings of the 2009 ACM SIGPLAN/SIGOPS international conference on Virtual execution environments*, pp. 51–60, 2009.
- 51) Thomas Roblitz and Frank Mueller. Combining Multi-Threading with Asynchronous Communication: A Case Study with DSM-Threads using Myrinet via BIP and Madeleine. *Myrinet User Group Conference*, pp. 131–138, Sep 2000.
- 52) Sriram Sankaran, Jeffrey M. Squyres, Brian Barrett, Vishal Sahay, and Andrew Lumsdaine. The Lam/Mpi Checkpoint/Restart Framework: System-Initiated Checkpointing. *International Journal of High Performance Computing Applications*, Vol.19, pp. 479–493, 2005.
- 53) Otto Sievert and Henri Casanova. A Simple MPI Process Swapping Architecture for Iterative Applications. *International Journal of High Performance Computing Applications*, Vol.18, pp. 341–352, Aug 2004.
- 54) Jimmy Su, Tong Wen, and Katherine Yelick. Compiler and Runtime Support for Scaling Adaptive Mesh Refinement Computations in Titanium. Technical report, Electrical Engineering and Computer Sciences University of California at Berkeley, Jun 2006.
- 55) Jimmy Su and Katherine Yelick. Automatic Support for Irregular Computations in a High-Level Language. *19th IEEE International Parallel and Distributed Processing Symposium*, Vol.1, p. 53b, 2005.
- 56) Chao Wang, Frank Mueller, Christian Engelmann, and Stephen L.Scott. Proactive process-level live migration in HPC environments. *Proceedings of the 2008 ACM/IEEE conference on Supercomputing*, pp. 1–12, Nov 2008.
- 57) Boris Weissman, Benedict Gomes, Jurgen W.Quittek, and Michael Holtkamp. Efficient Fine-grain Thread Migration with Active Threads. *Proceedings of the 12th. International Parallel Processing Symposium on International Parallel Processing Symposium*, p. 410, 1998.
- 58) Robert W.Numrich, John Reid, and Kieun Kim. Writing a Multigrid Solver Using Co-array Fortran. *4th International Workshop on Applied Parallel Computing, Large Scale Scientific and Industrial Problems*, Vol. 1541, pp. 390–399, 1998.
- 59) Katherine Yelick, Paul Hilfinger, Susan Graham, Dan Bonachea, Jimmy Su, Amir Kamil, Kaushik Datta, Phillip Colella, and Tong Wen. Parallel Languages and Compilers: Perspective from the Titanium Experience. *Journal of High Performance Computing Applications*, Vol.21, No.3, pp. 266–290, 2007.
- 60) Kathy Yelick, Luigi Semenzato, Geoff Pike, Carleton Miyamoto, Ben Liblit, Arvind Krishnamurthy, Paul Hilfinger, Susan Graham, David Gay, Phil Colella, and Alex Aiken. Titanium: A High-Performance Java Dialect. *ACM 1998 Workshop on Java for High-Performance Network Computing*, Vol.10, No. 11–13, pp. 825–836, Feb 1998.
- 61) Wenzhang Zhu, Cho-Li Wang, and Francis C.M.Lau. Lightweight Transparent Java Thread Migration for Distributed JVM. *International Conference on Parallel Processing*, p. 465, Oct 2003.
- 62) Wenzhang Zhu, Cho-Li Wang, and Lau F.C.M. JESSICA2: a distributed Java Virtual Machine with transparent thread migration support. *Fourth IEEE International Conference on Cluster Computing*, pp. 381–388, 2002.
- 63) Wenzhang Zhu and Submitted Wenzhang Zhu. JESSICA2: A distributed Java virtual machine with transparent thread migration support. *IEEE International Conference on Cluster Computing*, Sep 2002.
- 64) 原健太郎. 有限要素法における連立方程式ソルバの並列化. 第9回PCクラスタシンポジウム, Dec 2009.
- 65) 原健太郎, 田浦健次郎, 近山隆. DMI: 計算資源の動的な参加/脱退をサポートする大規模分散共有メモリアンタフェース. *SWoPP2009*, Aug 2009.
- 66) 原健太郎, 田浦健次郎, 近山隆. DMI: 計算資源の動的な参加/脱退をサポートする大規模分散共有メモリアンタフェース. *情報処理学会論文誌 (プログラミング)*, Vol.3, No.1, pp. 1–40, Mar 2010.
- 67) 高橋浩和, 小田逸郎, 山幡為佐久. Linux カーネル 2.6 解説室. ソフトバンククリエイティブ, Nov 2006.
- 68) 緑川博子, 飯塚肇. ユーザレベル・ソフトウェア分散共有メモリ SMS の設計と実装. *情報処理学会論文誌ハイパフォーマンスコンピューティングシステム*, Vol.42, No. SIG9, pp. 170–190, Aug 2001.

## 付 録

### A.1 random-address の最適性の証明

#### A.1.1 証明すべき定理

各スレッドは、自分以外のスレッドがどのアドレスをどれくらい使用しているかに関する知識を持たないとする。このとき、スレッド移動時のアドレス衝突確率を最小化する戦略のひとつが、各スレッドがアドレスを連続的に使用する戦略であることを証明する。ただし、各スレッドが、自分以外のスレッドがアドレスをどう使用しているかに関する知識を完全に持たないと仮定してしまうと、そのスレッドはどのような戦略に基づいてアドレスを使用したとしても、アドレス衝突確率が変わらないという結論になってしまう。したがって、各ス

レッドは、自分以外のレッドがどのアドレスをどれくらい使用しているかに関する知識を持たないが、自分以外のレッドがどのような戦略に基づいてアドレスを使用しているかは知っている」ことを仮定する。以下の議論では、まず「戦略」とは何かを（直観的な意味と合致するよう）合理的に定義し、問題を数学的に定式化したうえで、証明すべき定理を導く。

$m$  個のレッド  $x_0, x_1, \dots, x_{m-1}$  を考え、これらの各レッド  $x_i$  ( $0 \leq i \leq m-1$ ) が使用するアドレス空間を  $A = \{0, 1, \dots, n-1\}$  とする。このとき、各レッド  $x_i$  は、アドレス集合  $A = \{0, 1, \dots, n-1\}$  に含まれる  $n$  個のアドレスを何らかの順序で使うことになるが、ここでは、「各レッド  $x_i$  は置換  $\sigma_i$  に従って一様にアドレスを使用する」ともと仮定する。ここで、「レッド  $x_i$  が置換  $\sigma_i$  に従って一様にアドレスを使用する」とは以下の意味である：

定義 1 順列  $\{0, 1, \dots, n-1\}$  の置換  $\sigma_i = \{\sigma_i(0), \sigma_i(1), \dots, \sigma_i(n-1)\}$  を考え、さらにこの置換  $\sigma_i$  に対して以下の集合  $C^{\sigma_i}(A)$  を考える：

$$C^{\sigma_i}(A) = \{\{\sigma_i(a_i), \sigma_i(a_i+1), \dots, \sigma_i(a_i+b_i-1)\} \mid 0 \leq a_i \leq n-1, 0 \leq b_i \leq n\} \quad (1)$$

$C^{\sigma_i}(A)$  はアドレス集合の集合である。このとき、レッド  $x_i$  が使用するアドレス集合が、つねに  $C^{\sigma_i}(A)$  の要素集合であるとき、「レッド  $x_i$  は戦略  $\sigma_i$  に従ってアドレスを使用する」と定義する。さらに、レッド  $x_i$  が使用するアドレス集合が、 $C^{\sigma_i}(A)$  の各要素集合を等確率でとる\*1とき、「レッド  $x_i$  は戦略  $\sigma_i$  に従って一様にアドレスを使用する」と定義する。

具体例として、 $n = 4$  とし、順列  $\{0, 1, 2, 3\}$  の置換  $\sigma_i = \{2, 1, 3, 0\}$  を考える。このとき、 $C^{\sigma_i}(A)$  は、

$$C^{\sigma_i}(A) = \{\{\}, \{2\}, \{1\}, \{3\}, \{0\}, \{2, 1\}, \{1, 3\}, \{3, 0\}, \{0, 2\}, \\ \{2, 1, 3\}, \{1, 3, 0\}, \{3, 0, 2\}, \{0, 2, 1\}, \{2, 1, 3, 0\}\}$$

となる。よって、「レッド  $x_i$  が置換  $\sigma_i$  に従ってアドレスを使用する」とは、レッド  $x_i$  が使用するアドレス集合は、つねに、 $\{\}, \{2\}, \{1\}, \{3\}, \{0\}, \{2, 1\}, \{1, 3\}, \{3, 0\}, \{0, 2\}, \{2, 1, 3\}, \{1, 3, 0\}, \{3, 0, 2\}, \{0, 2, 1\}, \{2, 1, 3, 0\}$  のいずれかであるという意味である。また、「レッド  $x_i$  が置換  $\sigma_i$  に従って一様にアドレスを使用する」とは、レッド  $x_i$  は、アドレ

\*1 具体的な値は重要ではないが、具体的には  $1/|C^{\sigma_i}(A)| = 1/(n^2 - n + 2)$  である。なぜなら、集合  $C^{\sigma_i}(A)$  の定義式 (1) において、 $b_i = 0$  の場合にはすべての  $a_i$  に対してアドレス集合  $\{\}$  が生成されること、 $b_i = n$  の場合にはすべての  $a_i$  に対してアドレス集合  $\{\sigma_i(0), \dots, \sigma_i(n-1)\}$  が生成されること、 $1 \leq b_i \leq n-1$  の場合には各  $a_i$  に対して  $n$  個の異なるアドレス集合  $\{\sigma_i(a_i), \sigma_i(a_i+1), \dots, \sigma_i(a_i+b_i-1)\}$  が生成されることから、結局、 $|C^{\sigma_i}(A)| = 1 + 1 + (n-1)n = n^2 - n + 2$  となるからである。

ス集合として  $\{\}, \{2\}, \{1\}, \{3\}, \{0\}, \{2, 1\}, \{1, 3\}, \{3, 0\}, \{0, 2\}, \{2, 1, 3\}, \{1, 3, 0\}, \{3, 0, 2\}, \{0, 2, 1\}, \{2, 1, 3, 0\}$  を等確率で使用するという意味である。

なお、置換  $\sigma_j$  が置換  $\sigma_i$  の巡回置換であるとき、「レッド  $x_i$  が置換  $\sigma_i$  に従って（一様に）アドレスを使用する」と「レッド  $x_i$  が置換  $\sigma_j$  に従って（一様に）アドレスを使用する」ことはまったく等価である。よって、以降の議論では、置換  $\sigma_j$  が置換  $\sigma_i$  の巡回置換であるとき、 $\sigma_i = \sigma_j$  と表記することにする。

ここで、「各レッド  $x_i$  が置換  $\sigma_i$  に従って一様にアドレスを使用する」という仮定 A は、本節冒頭で述べた「各レッドは、自分以外のレッドがどのアドレスをどれくらい使用しているかに関する知識を持たないが、自分以外のレッドがどのような戦略に基づいてアドレスを使用しているかは知っている」という仮定 B の直観的な意味と合致しており、かつ十分に一般的であることを強調しておく。なぜなら、仮定 A のもとでは、各レッド  $x_i$  に対して置換  $\sigma_i$  を適切に選ぶことによって「各レッド  $x_i$  がアドレス集合  $A = \{0, 1, \dots, n-1\}$  に含まれる  $n$  個のアドレスを任意の順序で使用する」戦略すべてを表現できるからである。言い換えると、仮定 A のもとでは、仮定 B が直観的に意味するであろうすべての戦略を表現できる。また、仮定 A における「一様に」という条件は、仮定 B における「各レッドは、自分以外のレッドがどのアドレスをどれくらい使用しているかに関する知識を持たない」という状況を反映している。

以上の考察より、「レッド移動時のアドレス衝突確率を最小化する戦略のひとつは、各レッドがアドレスを連続的に使用する戦略である」ことを証明するために証明すべき定理として以下が得られる。

定理 1  $m$  個のレッド  $x_0, x_1, \dots, x_{m-1}$  を考え、各レッド  $x_i$  は置換  $\sigma_i$  に従って一様にアドレスを使用するとする。このとき、 $\sigma_0, \sigma_1, \dots, \sigma_{m-1}$  がとりうるすべての組合せ  $(\sigma_0, \sigma_1, \dots, \sigma_{m-1})$  のうち（各  $\sigma_i$  の選び方は  $n!$  通り存在するから、組合せは全部で  $n!^m$  通り存在する）、「どの 2 つの異なるレッド  $x_i$  とレッド  $x_j$  に対しても、レッド  $x_i$  が使用するアドレス集合とレッド  $x_j$  が使用するアドレス集合が共通部分を持たない確率」が最大になるのは、 $\sigma_0 = \sigma_1 = \dots = \sigma_{m-1} = \epsilon$  の場合である。ただしここで  $\epsilon$  は恒等置換を表す。

さらに、定理 1 を一般化して次の定理を考える：

定理 2  $m$  個のレッド  $x_0, x_1, \dots, x_{m-1}$  を考え、各レッド  $x_i$  は置換  $\sigma_i$  に従って一様にアドレスを使用するとする。このとき、 $\sigma_0, \sigma_1, \dots, \sigma_{m-1}$  がとりうるすべての組合せ  $(\sigma_0, \sigma_1, \dots, \sigma_{m-1})$  のうち、「どの 2 つの異なるレッド  $x_i$  とレッド  $x_j$  に対しても、

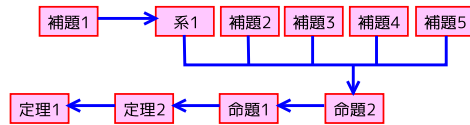


図 29 命題や補題の論理関係.  $A \rightarrow B$  は, 証明において  $A$  から  $B$  を導くことを意味する.

スレッド  $x_i$  が使用するアドレス集合とスレッド  $x_j$  が使用するアドレス集合が共通部分を持たない確率が最大になるのは,  $\sigma_0 = \sigma_1 = \dots = \sigma_{m-1}$  の場合であり, かつその場合に限られる.

明らかに, 定理 2 は定理 1 の拡張になっており, 定理 2 が証明されれば, 定理 1 も証明されたことになる. 次節では, 定理 2 の証明を行う.

#### A.1.2 証明

定理 2 を証明する. 以降, 命題や補題をいくつか立てて証明を進めるが, これらの導出関係を図 29 に示しておく.

まず, 以下の命題を考える. これは定理 2 において  $m = 2$  とした場合に相当する:

命題 1 スレッド  $x$  とスレッド  $y$  を考え, スレッド  $x$  は置換  $\sigma_x$  に従って一様にアドレスを使用し, スレッド  $y$  は置換  $\sigma_y$  に従って一様にアドレスを使用するとする. このとき, スレッド  $x$  が使用するアドレス集合とスレッド  $y$  が使用するアドレス集合が共通部分を持つ確率が最小になるのは,  $\sigma_x = \sigma_y$  の場合であり, かつその場合に限られる.

命題 1 の証明に入る前に, 以降の議論で使用する記号をいくつか導入する:

- アドレス集合  $A$  の冪集合を  $P(A)$  と表す. スレッド  $x$  が使用しているアドレス集合を  $S_x$ , スレッド  $y$  が使用しているアドレス集合を  $S_y$  とすれば, 明らかに  $S_x \in P(A)$ ,  $S_y \in P(A)$  が成り立つ. さらに,  $S_x \in C^{\sigma_x}(A)$ ,  $S_y \in C^{\sigma_y}(A)$  も成り立つ.
- アドレス集合  $A$  の冪集合  $P(A)$  の中で, 大きさが  $i$  であるようなアドレス集合の集合を  $P_i(A)$  と表す. すなわち, 集合  $P_i(A)$  の任意の要素は大きさ  $i$  の集合であり,  $|P_i(A)| = {}_n C_i$ ,  $P_i(A) \cap P_j(A) = \emptyset$  ( $i \neq j$ ),  $P(A) = \bigsqcup_{i=0}^n P_i(A)$  が成り立つ. 同様に, アドレス集合の集合  $C^{\sigma_y}(A)$  の中で, 大きさが  $i$  であるようなアドレス集合の集合を  $C_i^{\sigma_y}(A)$  と表す.
- 任意のアドレス集合の集合  $C$  と任意のアドレス集合  $S_0, S_1, \dots$  に関して, 集合  $C$  に属するすべてのアドレス集合  $S \in C$  のうち,  $(S_0 \cap S \neq \emptyset) \wedge (S_1 \cap S \neq \emptyset) \wedge \dots$  を満足する  $S$  の個数を,  $M(C; S_0, S_1, \dots)$  と表す. たとえば,  $M(C^{\sigma_y}(A); S_x)$  は,  $C^{\sigma_y}(A)$  に属するすべてのアドレス集合  $S_y \in C^{\sigma_y}(A)$  のうち,  $S_y \cap S_x \neq \emptyset$  を満足するような

$S_y$  の個数を表す. 以上のような定義から, 明らかに, 任意のアドレス集合  $S_i, S_j$  に対して,

$$M(C; \dots, S_i, \dots, S_j, \dots) = M(C; \dots, S_j, \dots, S_i, \dots)$$

が成り立つ. さらに, 任意のアドレス集合  $S_i, S_j$  に対して,  $S_i \cup S_j = S_i + S_j - S_i \cap S_j$  が成り立つので,

$$M(C; \dots, S_i \cup S_j, \dots) = M(C; \dots, S_i, \dots) + M(C; \dots, S_j, \dots) - M(C; \dots, S_i, S_j, \dots) \quad (2)$$

が成り立つ.

- 以降では  $n$  を法とする剰余環での議論が多くなるため, 任意の整数  $i$  に対して  $\text{mod}(i, n)$  を  $\underline{i}$  と略記することにする. よって,  $0 \leq i < n$  ならば,

$$\underline{i} = i \quad (3)$$

が成り立つ. また, 任意の整数  $i, j$  に対して,

$$\underline{i + j} = \underline{i} + \underline{j} \quad (4)$$

が成り立つ.

- 虚数単位  $j$  と整数  $i_0, i_1, \dots, i_{k-1}$  に対して, 複素平面上の  $k$  個の点  $e^{2\pi i_0 j/n}$ ,  $e^{2\pi i_1 j/n}$ ,  $\dots$ ,  $e^{2\pi i_{k-1} j/n}$  を考える. 偏角を  $[0, 2\pi)$  の範囲で考えるとき, ある整数  $\alpha$  ( $0 \leq \alpha \leq k-1$ ) が存在して,  $0 \leq \arg(e^{2\pi i \text{mod}(\alpha, k) j/n}) \leq \arg(e^{2\pi i \text{mod}(\alpha+1, k) j/n}) \leq \dots \leq \arg(e^{2\pi i \text{mod}(\alpha+k-1, k) j/n}) < 2\pi$  の関係が成り立つとき,  $i_0 \preceq i_1 \preceq \dots \preceq i_{k-1} \preceq$  と表す. 特に, 整数  $\alpha'$  ( $0 \leq \alpha' \leq k-1$ ) に対して,  $\arg(e^{2\pi i \text{mod}(\alpha', k) j/n}) < \arg(e^{2\pi i \text{mod}(\alpha'+1, k) j/n})$  が成り立つとき,  $i_0 \preceq i_1 \preceq \dots \preceq i_{\alpha'} \prec i_{\alpha'+1} \preceq \dots \preceq i_{k-1} \preceq$  と表す. 図形的に言えば,  $i_0 \preceq i_1 \prec i_2 \prec i_3 \preceq$  は, 4 個の点  $e^{2\pi i_0 j/n}$ ,  $e^{2\pi i_1 j/n}$ ,  $e^{2\pi i_2 j/n}$ ,  $e^{2\pi i_3 j/n}$  がこの順序で円周上に反時計回りに配置されており, かつ,  $e^{2\pi i_1 j/n}$  と  $e^{2\pi i_2 j/n}$  および  $e^{2\pi i_2 j/n}$  と  $e^{2\pi i_3 j/n}$  は異なる点であることを意味する (図 30 (A)).

以上の定義より, 任意の整数  $i_0, i_1, \dots, i_{k-1}$  に対して,

$$i_0 \preceq i_1 \preceq \dots \preceq i_{k-1} \preceq \\ \iff \underline{i_1 - i_0} + \underline{i_2 - i_1} + \dots + \underline{i_{k-1} - i_{k-2}} + \underline{i_0 - i_{k-1}} = n \quad (5)$$

が成り立つ. また,  $i_0 \preceq i_1 \preceq i_2 \preceq$  ならば,

$$\underline{i_1 - i_0} + \underline{i_2 - i_1} = \underline{i_2 - i_0} \quad (6)$$

が成り立つ.

ここで, 大きさが  $b_x$  の任意のアドレス集合  $S \in P_{b_x}(A)$  は, 任意の置換  $\sigma_y$  に対して,

$i_0 < i_1 < \dots < i_{b_x-1}$  なる整数  $i_0, i_1, \dots, i_{b_x-1}$  を用いて,

$$S = \{\sigma_y(i_0), \sigma_y(i_1), \dots, \sigma_y(i_{b_x-1})\}$$

と表現できることに注意する．たとえば,  $n = 8, b_x = 4$  として,  $S = \{1, 2, 5, 7\}$  は, 置換  $\sigma_y = \{5, 2, 1, 0, 7, 3, 6, 4\}$  を用いて,  $S\{\sigma_y(0), \sigma_y(1), \sigma_y(2), \sigma_y(3)\}$  と表現できる．よって, 以降では証明の直観的な理解を助けるため, アドレス集合の様子を図 30 のような円周上の点集合として図示することにする．

いま導入した記号を用いて, 命題 1 を定式化し, より証明しやすい形式の命題 2 に言い換える．

命題 1 においては, スレッド  $x$  とスレッド  $y$  は, それぞれ置換  $\sigma_x$  と置換  $\sigma_y$  に従って「一様に」アドレスを使用することを仮定しているので, 「 $\sigma_x = \sigma_y$  のときに, スレッド  $x$  が使用するアドレス集合とスレッド  $y$  が使用するアドレス集合が共通部分を持つ確率が最小になる」という命題 1 の題意は, 「 $\sigma_x = \sigma_y$  のときに, スレッド  $x$  が使用するアドレス集合  $S_x \in C^{\sigma_x}(A)$  とスレッド  $y$  が使用するアドレス集合  $S_y \in C^{\sigma_y}(A)$  のすべての組合せ  $(S_x, S_y)$  (全部で  $|C^{\sigma_x}(A)| \times |C^{\sigma_y}(A)|$  個ある) のなかで,  $S_x \cap S_y \neq \emptyset$  を満たすような組合せ  $(S_x, S_y)$  の個数が最小になる」ということと等価である．そして, これをさらに言い換えると, 「置換  $\sigma_y$  が与えられたとき, すべての置換  $\sigma_x$  のなかで, スレッド  $x$  が使用するアドレス集合  $S_x \in C^{\sigma_x}(A)$  とスレッド  $y$  が使用するアドレス集合  $S_y \in C^{\sigma_y}(A)$  のすべての組合せ  $(S_x, S_y)$  (全部で  $|C^{\sigma_x}(A)| \times |C^{\sigma_y}(A)|$  個ある) のなかで,  $S_x \cap S_y \neq \emptyset$  を満たすような組合せ  $(S_x, S_y)$  の個数が最小になる」ような置換  $\sigma_x$  とは, 置換  $\sigma_y$  である」ということと等価である．したがって, 命題 1 を言い換えると, 「 $P(A)$  に属するすべてのアドレス集合  $S_x \in P(A)$  のなかで, 『すべてのアドレス集合  $S_y \in C^{\sigma_y}(A)$  のうち,  $S_y \cap S_x \neq \emptyset$  を満足するような  $S_y$  の個数』が最小になるような  $S_x$  の集合は  $C^{\sigma_y}(A)$  である」と言い換えることができる．さらに, これを  $S_x$  の大きさによって分解して言い換えると, 「 $0 \leq b_x \leq n$  とする.  $P_{b_x}(A)$  に属するすべてのアドレス集合  $S_x \in P_{b_x}(A)$  のなかで, 『すべてのアドレス集合  $S_y \in C^{\sigma_y}(A)$  のうち,  $S_y \cap S_x \neq \emptyset$  を満足するような  $S_y$  の個数』が最小になるような  $S_x$  の集合は  $C_{b_x}^{\sigma_y}(A)$  である」と言い換えることができる．さらに, これを先ほど導入した記号を用いて言い換えると, 「 $0 \leq b_x \leq n$  とする.  $P_{b_x}(A)$  に属するすべてのアドレス集合  $S_x \in P_{b_x}(A)$  のなかで,  $M(C^{\sigma_y}(A); S_x)$  が最小になるような  $S_x$  の集合は  $C_{b_x}^{\sigma_y}(A)$  である」と言い換えることができる．

以上の議論により, 命題 1 と等価な命題 2 が得られる:

命題 2  $0 \leq b_x \leq n$  とする.  $P_{b_x}(A)$  に属するすべてのアドレス集合  $S_x \in P_{b_x}(A)$  のう

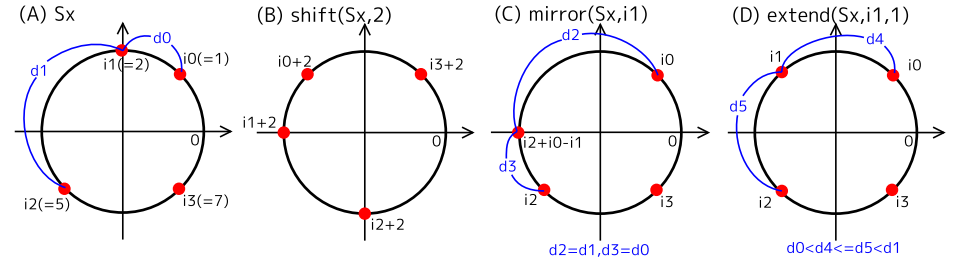


図 30 アドレス集合  $S_x$  と各写像の具体例 (A)  $S_x$  (B)  $\text{shift}(S_x, s)$  (C)  $\text{mirror}(S_x, i_\alpha)$  (D)  $\text{extend}(S_x, i_\alpha, s)$  .

ち,  $M(C^{\sigma_y}(A); S_x)$  を最小にする  $S_x$  の集合は  $C_{b_x}^{\sigma_y}(A)$  である．

命題 2 をわかりやすく書き下すと,  $b_x$  の値に応じて次のようになる． $P_{b_x}(A)$  に属するすべてのアドレス集合  $S_x \in P_{b_x}(A)$  のうち,  $M(C^{\sigma_y}(A); S_x)$  を最小化する  $S_x$  は,

- $b_x = 0$  のとき,  $\{\}$  の 1 通りのみである．
- $b_x = n$  のとき,  $\{0, 1, \dots, n-1\}$  の 1 通りのみである．
- $1 \leq b_x \leq n-1$  のとき, 以下の  $T^0, T^1, \dots, T^{n-1}$  の合計  $n$  通りのみである:

$$T^0 = \{\sigma_y(0), \sigma_y(1), \dots, \sigma_y(b_x-2), \sigma_y(b_x-1)\},$$

$$T^1 = \{\sigma_y(1), \sigma_y(2), \dots, \sigma_y(b_x-1), \sigma_y(b_x)\},$$

⋮

$$T^{n-1} = \{\sigma_y(n-b_x+1), \sigma_y(n-b_x+2), \dots, \sigma_y(n-1), \sigma_y(0)\}.$$

まず  $b_x = 0$  のときには,  $P_0(A) = \{\{\}\}$  であるから, 命題 2 は明らかである． $b_x = 1$  のときには,  $P_1(A) = \{\{0\}, \{1\}, \dots, \{n-1\}\}$  であるから, 対称性より命題 2 は明らかである．また  $b_x = n$  のときには,  $P_n(A) = \{\{0, 1, \dots, n-1\}\}$  であるから, 命題 2 は明らかである．したがって, 以降の議論においては,  $2 \leq b_x \leq n-1$  の場合に命題 2 が成り立つことを証明する．命題 2 の証明が終わるまでの間,  $2 \leq b_x \leq n-1$  を仮定して議論する．

さて,  $S_x \in P_{b_x}(A)$  なる任意の  $S_x$  は,  $i_0 < i_1 < \dots < i_{b_x-1}$  なる整数  $i_0, i_1, \dots, i_{b_x-1}$  を用いて,

$$S_x = \{\sigma_y(i_0), \sigma_y(i_1), \dots, \sigma_y(i_{b_x-1})\}$$

と表すことができる (図 30 (A)) . なお, 以降の議論では  $i$  の添字は  $b_x$  を法とする剰余環上で計算するものとする．つまり,  $i_{\text{mod}(j, b_x)}$  を  $i_j$  と略記する．

ここで、任意の  $S_x \in P_{b_x}(A)$  に対して、3つの写像  $shift$ ,  $mirror$ ,  $extend$  を以下のように定義する：

**shift** 任意の整数  $s$  に対して、

$$shift(S_x, s) = \{\sigma_y(i_0 + s), \sigma_y(i_1 + s), \dots, \sigma_y(i_{b_x-1} + s)\}.$$

**mirror**  $\sigma_y(i_\alpha) \in S_x$  を満たす任意の  $i_\alpha$  に対して、

$$mirror(S_x, i_\alpha) = S_x \setminus \{\sigma_y(i_\alpha)\} \cup \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}.$$

**extend**  $(\sigma_y(i_\alpha) \in S_x) \wedge (i_\alpha \prec i_\alpha + s \prec i_{\alpha+1} \prec) \wedge ((i_\alpha + s) - i_{\alpha-1} \leq i_{\alpha+1} - (i_\alpha + s))$  を満たす任意の整数  $s$  と  $i_\alpha$  に対して、

$$extend(S_x, i_\alpha, s) = S_x \setminus \{\sigma_y(i_\alpha)\} \cup \{\sigma_y(i_\alpha + s)\}.$$

なお、いまは  $2 \leq b_x \leq n-1$  を仮定しているので、 $i_{\alpha-1} = i_{\alpha+1}$  の可能性はあるが ( $b_x = 2$  のとき)、 $i_{\alpha-1} \neq i_\alpha$  かつ  $i_\alpha \neq i_{\alpha+1}$  が成り立つことに注意する。

$n = 8$ ,  $b_x = 4$  とした場合の、3つの写像  $shift$ ,  $mirror$ ,  $extend$  の例を、それぞれ図 30 (B) (C) (D) に示す。直観的には、図 30 に示すような円周上の距離で考えたとき、 $shift(S_x, s)$  の図形的意味は「 $S_x$  全体を距離  $s$  だけ移動させる」、 $mirror(S_x, i_\alpha)$  の図形的意味は「点  $i_\alpha$  を、点  $i_{\alpha-1}$  と点  $i_{\alpha+1}$  の中点に関して対称な位置に移動させる」、 $extend(S_x, i_\alpha, s)$  の図形的意味は「点  $i_\alpha$  を距離  $s$  だけ移動させる。ただし、このとき円周上で  $i_{\alpha-1}$ ,  $i_\alpha$ ,  $i_\alpha + s$ ,  $i_{\alpha+1}$  の順に反時計回りに点が並んでおり、かつ、 $i_\alpha + s$  と  $i_{\alpha+1}$  の距離は  $i_{\alpha-1}$  と  $i_\alpha + s$  の距離以上になっていなければならない」という意味である。

ここで、3つの写像  $shift$ ,  $mirror$ ,  $extend$  に関して、以下の3つの補題を考え、証明する：

**補題 1** 任意のアドレス集合  $S_x^0, S_x^1, \dots \in P_{b_x}(A)$  と任意の整数  $s$  に対して、

$$M(C^{\sigma_y}(A); S_x^0, S_x^1, \dots) = M(C^{\sigma_y}(A); shift(S_x^0, s), shift(S_x^1, s), \dots)$$

が成り立つ。

**補題 2** 任意のアドレス集合  $S_x \in P_{b_x}(A)$  と  $\sigma_y(i_\alpha) \in S_x$  を満たす任意の  $i_\alpha$  に対して、

$$M(C^{\sigma_y}(A); S_x) = M(C^{\sigma_y}(A); mirror(S_x, i_\alpha))$$

が成り立つ。

**補題 3** 任意のアドレス集合  $S_x \in P_{b_x}(A)$ 、および  $(\sigma_y(i_\alpha) \in S_x) \wedge (i_\alpha \prec i_\alpha + s \prec i_{\alpha+1} \prec) \wedge ((i_\alpha + s) - i_{\alpha-1} \leq i_{\alpha+1} - (i_\alpha + s))$  を満たす任意の整数  $s$  と整数  $i_\alpha$  に関して、

$$M(C^{\sigma_y}(A); S_x) < M(C^{\sigma_y}(A); extend(S_x, i_\alpha, s))$$

が成り立つ。

補題 1 を示す。まず、 $C^{\sigma_y}(A)$  の定義は、式 (1) より、

$C^{\sigma_y}(A) = \{\{\sigma_y(a_y), \sigma_y(a_y + 1), \dots, \sigma_y(a_y + b_y - 1)\} \mid 0 \leq a_y \leq n-1, 0 \leq b_y \leq n\}$  である。ここで、任意の  $S_y = \{\sigma_y(a_y), \sigma_y(a_y + 1), \dots, \sigma_y(a_y + b_y - 1)\} \in C^{\sigma_y}(A)$  ( $0 \leq a_y \leq n-1, 0 \leq b_y \leq n$ ) に対して、

$$(S_x^0 \cap \{\sigma_y(a_y), \sigma_y(a_y + 1), \dots, \sigma_y(a_y + b_y - 1)\} \neq \emptyset) \wedge$$

$$(S_x^1 \cap \{\sigma_y(a_y), \sigma_y(a_y + 1), \dots, \sigma_y(a_y + b_y - 1)\} \neq \emptyset) \wedge \dots$$

$$\iff (shift(S_x^0, s) \cap \{\sigma_y(a_y + s), \sigma_y(a_y + 1 + s), \dots, \sigma_y(a_y + b_y - 1 + s)\} \neq \emptyset) \wedge$$

$$(shift(S_x^1, s) \cap \{\sigma_y(a_y + s), \sigma_y(a_y + 1 + s), \dots, \sigma_y(a_y + b_y - 1 + s)\} \neq \emptyset) \wedge \dots$$

が成り立つ。すなわち、アドレス集合  $S_x^0, S_x^1, \dots$  のすべてがアドレス集合  $S_y \in C^{\sigma_y}(A)$  と共通部分を持たないならば、そのような  $S_x^0, S_x^1, \dots$  に対して、アドレス集合  $shift(S_x^0, s)$ ,  $shift(S_x^1, s)$ ,  $\dots$  のすべてがアドレス集合  $S'_y$  と共通部分を持たないようなアドレス集合  $S'_y \in C^{\sigma_y}(A)$  がちょうど1つ存在する。したがって、「すべての  $S_y \in C^{\sigma_y}(A)$  のうち、 $(S_x^0 \cap S_y \neq \emptyset) \wedge (S_x^1 \cap S_y \neq \emptyset) \wedge \dots$  を満たす  $S_y$  の個数」と「すべての  $S_y \in C^{\sigma_y}(A)$  のうち、 $(shift(S_x^0, s) \cap S_y \neq \emptyset) \wedge (shift(S_x^1, s) \cap S_y \neq \emptyset) \wedge \dots$  を満たす  $S_y$  の個数」は等しい。よって、

$$M(C^{\sigma_y}(A); S_x^0, S_x^1, \dots) = M(C^{\sigma_y}(A); shift(S_x^0, s), shift(S_x^1, s), \dots)$$

が成り立つ。以上より、補題 1 が示された。■

系 1 任意のアドレス集合  $S_x \in P_{b_x}(A)$  と任意の整数  $s$  に対して、

$$M(C^{\sigma_y}(A); S_x) = M(C^{\sigma_y}(A); shift(S_x, s)).$$

補題 1 より明らかである。■

補題 2 を示す。左辺と右辺をそれぞれ計算し、両者が一致することを示す。

まず、式 (2) を用いて補題 2 の左辺を計算すると、

$$\begin{aligned} M(C^{\sigma_y}(A); S_x) &= M(C^{\sigma_y}(A); (S_x \setminus \{\sigma_y(i_\alpha)\}) \cup \{\sigma_y(i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}) \\ &\quad - M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_\alpha)\}) \end{aligned} \quad (7)$$

となる。上式の第1項、第2項、第3項をそれぞれ  $D_1$ ,  $D_2$ ,  $D_3$  とおく。ここで、 $D_3 = M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_\alpha)\})$  は「すべての  $S_y \in C^{\sigma_y}(A)$  のうち、 $((S_x \setminus \{\sigma_y(i_\alpha)\}) \cap S_y \neq \emptyset) \wedge (\{\sigma_y(i_\alpha)\} \cap S_y \neq \emptyset)$  を満たす  $S_y$  の個数」を意味しているが、これは「すべての  $S_y \in C^{\sigma_y}(A)$  のうち、 $(\{\sigma_y(i_{\alpha-1}), \sigma_y(i_{\alpha+1})\} \cap S_y \neq \emptyset) \wedge (\{\sigma_y(i_\alpha)\} \cap S_y \neq \emptyset)$  を満たす  $S_y$  の個数」に等しい。なぜなら、 $i_0 \prec i_1 \prec \dots \prec i_{\alpha-1} \prec i_\alpha \prec i_{\alpha+1} \prec \dots \prec i_{b_x-1} \prec$  なので、 $((S_x \setminus \{\sigma_y(i_\alpha)\}) \cap S_y \neq \emptyset) \wedge (\{\sigma_y(i_\alpha)\} \cap S_y \neq \emptyset)$  を満たすような  $S_y$  は、必ず  $\sigma_y(i_{\alpha-1})$  ま

たは  $\sigma_y(i_{\alpha+1})$  を含むからである．すなわち，

$$\begin{aligned} D_3 &= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1}), \sigma_y(i_{\alpha+1})\}, \{\sigma_y(i_\alpha)\}) \end{aligned}$$

が成り立つ．さらに  $D_3$  の計算を進めると，

$$\begin{aligned} D_3 &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1}), \sigma_y(i_{\alpha+1})\}, \{\sigma_y(i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\} \cup \{\sigma_y(i_{\alpha+1})\}, \{\sigma_y(i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \\ &\quad - M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \quad (\because \text{式 (2)}) \end{aligned} \quad (8)$$

が得られる．上式の第 1 項，第 2 項，第 3 項をそれぞれ  $D_4$ ， $D_5$ ， $D_6$  とおく．以上をまとめると，

$$M(C^{\sigma_y}(A); S_x) = D_1 + D_2 - (D_4 + D_5 - D_6) \quad (9)$$

となる．

同様にして，式 (2) を用いて補題 2 の右辺を計算すると，

$$\begin{aligned} &M(C^{\sigma_y}(A); \text{mirror}(S_x, i_\alpha)) \\ &= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\} \cup \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \\ &\quad - M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \end{aligned}$$

となるので，この第 1 項，第 2 項，第 3 項をそれぞれ  $D'_1$ ， $D'_2$ ， $D'_3$  とおく．先ほどと同様にして  $D'_3$  を計算すると，

$$\begin{aligned} D'_3 &= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1}), \sigma_y(i_{\alpha+1})\}, \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\} \cup \{\sigma_y(i_{\alpha+1})\}, \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \\ &\quad + M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \\ &\quad - M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \quad (\because \text{式 (2)}) \end{aligned}$$

が得られる．上式の第 1 項，第 2 項，第 3 項をそれぞれ  $D'_4$ ， $D'_5$ ， $D'_6$  とおく．以上をまとめると，

$$M(C^{\sigma_y}(A); \text{mirror}(S_x, i_\alpha)) = D'_1 + D'_2 - (D'_4 + D'_5 - D'_6) \quad (10)$$

となる．

$D_1$ ， $D_2$ ， $D_4$ ， $D_5$  と  $D'_1$ ， $D'_2$ ， $D'_4$ ， $D'_5$  の大小を比較すると，

$$D_1 = M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) = D'_1, \quad (11)$$

$$\begin{aligned} D_2 &= M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); \text{shift}(\{\sigma_y(i_\alpha)\}, i_{\alpha-1} + i_{\alpha+1} - 2i_\alpha)) \quad (\because \text{補題 1}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \\ &= D'_2, \end{aligned} \quad (12)$$

$$\begin{aligned} D_4 &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}) \\ &= M(C^{\sigma_y}(A); \text{shift}(\{\sigma_y(i_{\alpha-1})\}, i_{\alpha+1} - i_\alpha), \text{shift}(\{\sigma_y(i_\alpha)\}, i_{\alpha+1} - i_\alpha)) \quad (\because \text{補題 1}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \\ &= D'_5, \end{aligned} \quad (13)$$

$$\begin{aligned} D_5 &= M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \\ &= M(C^{\sigma_y}(A); \text{shift}(\{\sigma_y(i_\alpha)\}, i_{\alpha-1} - i_\alpha), \text{shift}(\{\sigma_y(i_{\alpha+1})\}, i_{\alpha-1} - i_\alpha)) \quad (\because \text{補題 1}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_{\alpha-1} + i_{\alpha+1} - i_\alpha)\}) \\ &= D'_4 \end{aligned} \quad (14)$$

となる．次に， $D_6$  と  $D'_6$  の大小を考える．まず，

$$\begin{aligned} &(\underline{i_{\alpha-1} + i_{\alpha+1} - i_\alpha} - i_{\alpha-1} + i_{\alpha+1} - (\underline{i_{\alpha-1} + i_{\alpha+1} - i_\alpha} + i_{\alpha-1} - i_{\alpha+1})) \\ &= \underline{i_{\alpha+1} - i_\alpha + i_\alpha - i_{\alpha-1} + i_{\alpha-1} - i_{\alpha+1}} \quad (\because \text{式 (4)}) \end{aligned}$$

$$= n \quad (\because \text{仮定より } i_{\alpha-1} \leq i_\alpha \leq i_{\alpha+1} \leq \text{と式 (5)})$$

であるから，式 (5) より， $i_{\alpha-1} \leq \underline{i_{\alpha-1} + i_{\alpha+1} - i_\alpha} \leq i_{\alpha+1} \leq$  が成り立つ．このことと，仮定より  $i_{\alpha-1} \leq i_\alpha \leq i_{\alpha+1} \leq$  であることを考慮すると，ある  $S_x$  が  $\sigma_y(i_{\alpha-1})$  と  $\sigma_y(i_\alpha)$  と  $\sigma_y(i_{\alpha+1})$  の 3 要素を含むならば， $S_x$  は  $\sigma_y(i_{\alpha-1})$  と  $\sigma_y(\underline{i_{\alpha-1} + i_{\alpha+1} - i_\alpha})$  と  $\sigma_y(i_{\alpha+1})$  の 3 要素も含むことがわかる．したがって，「すべての  $S_y \in C^{\sigma_y}(A)$  のうち， $(\{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \cap S_y \neq \emptyset$  を満たす  $S_y$  の個数」は，「すべての  $S_y \in C^{\sigma_y}(A)$  のうち， $(\{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(\underline{i_{\alpha-1} + i_{\alpha+1} - i_\alpha})\}, \{\sigma_y(i_{\alpha+1})\}) \cap S_y \neq \emptyset$  を満たす  $S_y$  の個数」に等しい．よって，

$$\begin{aligned} D_6 &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \\ &= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(\underline{i_{\alpha-1} + i_{\alpha+1} - i_\alpha})\}, \{\sigma_y(i_{\alpha+1})\}) \\ &= D'_6 \end{aligned} \quad (15)$$

となる．

式 (9)(10)(11)(12)(13)(14)(15) より，

$$M(C^{\sigma_y}(A); S_x) = M(C^{\sigma_y}(A); \text{mirror}(S_x, i_\alpha))$$

が成り立つ．以上より，補題 2 が示された． ■



補題 3 を示す。まず、補題 3 の左辺を計算すると、式 (7)(8) より、

$$\begin{aligned}
& M(C^{\sigma_y}(A); S_x) \\
&= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}) \\
&\quad - M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_\alpha)\}) \\
&= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}) \\
&\quad - (M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\})) \\
&\quad - M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_{\alpha+1})\}) \tag{16}
\end{aligned}$$

となる。上式の第 1 項、第 2 項、第 3 項、第 4 項、第 5 項を、それぞれ  $D_1, D_2, D_3, D_4, D_5$  とおく。

また、補題 2 のときの議論と同様にして、補題 3 の右辺を計算すると、

$$\begin{aligned}
& M(C^{\sigma_y}(A); \text{extend}(S_x, i_\alpha, s)) \\
&= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\} \cup \{\sigma_y(i_\alpha + s)\}) \\
&= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}) \\
&\quad - M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_\alpha + s)\}) \quad (\because \text{式 (2)}) \\
&= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}) \\
&\quad - M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1}), \sigma_y(i_{\alpha+1})\}, \{\sigma_y(i_\alpha + s)\}) \\
&\quad (\because i_{\alpha-1} \preceq i_\alpha + s \preceq i_{\alpha+1} \text{ かつ } i_{\alpha-1} \preceq i_\alpha \preceq i_{\alpha+1} \preceq) \\
&= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}) \\
&\quad - M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\} \cup \{\sigma_y(i_{\alpha+1})\}, \{\sigma_y(i_\alpha + s)\}) \\
&= M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}) \\
&\quad - (M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha + s)\})) \\
&\quad + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}, \{\sigma_y(i_{\alpha+1})\}) \\
&\quad - M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_{\alpha+1})\}) \quad (\because \text{式 (2)}) \tag{17}
\end{aligned}$$

となる。上式の第 1 項、第 2 項、第 3 項、第 4 項、第 5 項を、それぞれ  $D'_1, D'_2, D'_3, D'_4, D'_5$  とおく。

以降では、 $D_1, D_2, D_3, D_4, D_5$  と  $D'_1, D'_2, D'_3, D'_4, D'_5$  の大小を比較する。まず、 $D_1, D_2, D_5$  と  $D'_1, D'_2, D'_5$  を比較すると、

$$D_1 = M(C^{\sigma_y}(A); S_x \setminus \{\sigma_y(i_\alpha)\}) = D'_1, \tag{18}$$

$$\begin{aligned}
D_2 &= M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}) = M(C^{\sigma_y}(A); \text{shift}(\{\sigma_y(i_\alpha)\}, s)) \\
&= M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}) = D'_2, \tag{19}
\end{aligned}$$

$$D_5 = M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_{\alpha+1})\}) = D'_5 \tag{20}$$

が成り立つ。

次に、 $D_3 + D_4$  と  $D'_3 + D'_4$  を比較する。いま、 $b_y \neq b'_y$  ならば  $C_{b_y}^{\sigma_y}(A) \cap C_{b'_y}^{\sigma_y}(A) = \emptyset$  であり、 $C^{\sigma_y}(A) = \bigsqcup_{b_y=0}^n C_{b_y}^{\sigma_y}(A)$  が成り立つから、任意の集合  $S_x$  に対して、

$$M(C^{\sigma_y}(A); S_x) = \sum_{b_y=0}^n M(C_{b_y}^{\sigma_y}(A); S_x)$$

が成り立つことに着目する。したがって、

$$\begin{aligned}
& D_3 + D_4 \\
&= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \\
&= \sum_{b_y=0}^n M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}) + \sum_{b_y=0}^n M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \\
&= \sum_{b_y=0}^n \left( M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}) + M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}) \right), \\
& D'_3 + D'_4 \\
&= M(C^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha + s)\}) + M(C^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}, \{\sigma_y(i_{\alpha+1})\}) \\
&= \sum_{b_y=0}^n M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha + s)\}) \\
&\quad + \sum_{b_y=0}^n M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}, \{\sigma_y(i_{\alpha+1})\}) \\
&= \sum_{b_y=0}^n \left( M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha + s)\}) \right. \\
&\quad \left. + M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}, \{\sigma_y(i_{\alpha+1})\}) \right)
\end{aligned}$$

と展開できる。ここで、

$$d_3(b_y) = M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha)\}),$$

$$d_4(b_y) = M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_\alpha)\}, \{\sigma_y(i_{\alpha+1})\}),$$

$$d'_3(b_y) = M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_{\alpha-1})\}, \{\sigma_y(i_\alpha + s)\}),$$

$$d'_4(b_y) = M(C_{b_y}^{\sigma_y}(A); \{\sigma_y(i_\alpha + s)\}, \{\sigma_y(i_{\alpha+1})\})$$

とおくと,  $D_3 + D_4$  と  $D'_3 + D'_4$  は,

$$D_3 + D_4 = \sum_{b_y=0}^n (d_3(b_y) + d_4(b_y)), \quad (21)$$

$$D'_3 + D'_4 = \sum_{b_y=0}^n (d'_3(b_y) + d'_4(b_y)) \quad (22)$$

と表すことができる. よって, 以下では, 各  $b_y$  ( $0 \leq b_y \leq n$ ) の値に応じて,  $d_3(b_y) + d_4(b_y)$  と  $d'_3(b_y) + d'_4(b_y)$  がどのような値をとるか調べる.

まず,  $d_3(b_y)$  について考える.

(i)  $0 \leq b_y \leq \underline{i_\alpha - i_{\alpha-1}}$  のとき. アドレス集合  $S_y = \{\sigma_y(a_y), \sigma_y(a_y + 1), \dots, \sigma_y(a_y + b_y - 1)\} \in C_{b_y}^{\sigma_y}(A)$  ( $0 \leq a_y \leq n - 1$ ) が,  $\sigma_y(i_{\alpha-1})$  と  $\sigma_y(i_\alpha)$  の両方の要素を含むことはありえない. よって,  $d_3(b_y) = 0$  である.

(ii)  $\underline{i_\alpha - i_{\alpha-1}} < b_y \leq n - 1$  のとき. アドレス集合  $S_y = \{\sigma_y(a_y), \sigma_y(a_y + 1), \dots, \sigma_y(a_y + b_y - 1)\} \in C_{b_y}^{\sigma_y}(A)$  ( $0 \leq a_y \leq n - 1$ ) が,  $\sigma_y(i_{\alpha-1})$  と  $\sigma_y(i_\alpha)$  の両方の要素を含むのは,  $a_y$  が,  $i_\alpha - b_y + 1, i_\alpha - b_y + 2, \dots, i_{\alpha-1} - 1, i_{\alpha-1}$  を満たす場合である. よって,

$$d_3(b_y) = \underline{i_{\alpha-1} - (i_\alpha - b_y + 1) + 1} = \underline{b_y - (i_\alpha - i_{\alpha-1})}$$

である. ここで,  $0 \leq b_y - \underline{i_\alpha - i_{\alpha-1}} < n$  であることに注意すると,

$$d_3(b_y) = \underline{b_y - (i_\alpha - i_{\alpha-1})} = b_y - \underline{i_\alpha - i_{\alpha-1}} \quad (\because \text{式 (3)(4)})$$

となる.

(iii)  $b_y = n$  のとき. 明らかに  $d_3(b_y) = 1$  である.

同様にして,  $d_4(b_y)$ ,  $d'_3(b_y)$ ,  $d'_4(b_y)$  についても計算し, 結果をまとめると以下のようになる:

$$d_3(b_y) = \begin{cases} 0 & \text{if } 0 \leq b_y \leq \underline{i_\alpha - i_{\alpha-1}} \\ b_y - \underline{(i_\alpha - i_{\alpha-1})} & \text{if } \underline{i_\alpha - i_{\alpha-1}} < b_y \leq n - 1 \\ 1 & \text{if } b_y = n \end{cases} \quad (23)$$

$$d_4(b_y) = \begin{cases} 0 & \text{if } 0 \leq b_y \leq \underline{i_{\alpha+1} - i_\alpha} \\ b_y - \underline{(i_{\alpha+1} - i_\alpha)} & \text{if } \underline{i_{\alpha+1} - i_\alpha} < b_y \leq n - 1 \\ 1 & \text{if } b_y = n \end{cases} \quad (24)$$

$$d'_3(b_y) = \begin{cases} 0 & \text{if } 0 \leq b_y \leq \underline{(i_\alpha + s) - i_{\alpha-1}} \\ b_y - \underline{((i_\alpha + s) - i_{\alpha-1})} & \text{if } \underline{(i_\alpha + s) - i_{\alpha-1}} < b_y \leq n - 1 \\ 1 & \text{if } b_y = n \end{cases} \quad (25)$$

$$d'_4(b_y) = \begin{cases} 0 & \text{if } 0 \leq b_y \leq \underline{i_{\alpha+1} - (i_\alpha + s)} \\ b_y - \underline{(i_{\alpha+1} - (i_\alpha + s))} & \text{if } \underline{i_{\alpha+1} - (i_\alpha + s)} < b_y \leq n - 1 \\ 1 & \text{if } b_y = n \end{cases} \quad (26)$$

となる.

ここで, 式 (23)(24)(25)(26) において, 場合分けの境界値になっている  $b_y$  の値たち  $0, \underline{i_\alpha - i_{\alpha-1}}, \underline{i_{\alpha+1} - i_\alpha}, \underline{(i_\alpha + s) - i_{\alpha-1}}, \underline{i_{\alpha+1} - (i_\alpha + s)}, n$  に関して, その大小関係を調べると,

$$\begin{aligned} & 0 < \underline{i_\alpha - i_{\alpha-1}} \quad (\because \text{仮定より } i_{\alpha-1} < i_\alpha <) \\ & < \underline{(i_\alpha + s) - i_{\alpha-1}} \quad (\because \text{仮定より } i_{\alpha-1} < i_\alpha < \underline{i_\alpha + s} <) \\ & \leq \underline{i_{\alpha+1} - (i_\alpha + s)} \quad (\because \text{仮定そのまま}) \\ & < \underline{i_{\alpha+1} - i_\alpha} \quad (\because \text{仮定より } i_\alpha < \underline{i_\alpha + s} < i_{\alpha+1} <) \\ & < n \quad (\because \text{明らか}) \end{aligned} \quad (27)$$

が成り立つ. 図 30 (D) を見ると, この大小関係が図形的に理解できる.

以上を踏まえて,  $d_3(b_y) + d_4(b_y)$  と  $d'_3(b_y) + d'_4(b_y)$  の大小を,  $b_y$  の境界値に応じてまとめると以下のようになる.

(i)  $0 \leq b_y \leq \underline{i_\alpha - i_{\alpha-1}}$  のとき.

$$(d_3(b_y) + d_4(b_y)) - (d'_3(b_y) + d'_4(b_y)) = (0 + 0) - (0 + 0) = 0 \quad (28)$$

である.

(ii)  $\underline{i_\alpha - i_{\alpha-1}} < b_y \leq \underline{(i_\alpha + s) - i_{\alpha-1}}$  のとき.

$$(d_3(b_y) + d_4(b_y)) - (d'_3(b_y) + d'_4(b_y)) = ((b_y - \underline{(i_\alpha - i_{\alpha-1})}) + 0) - (0 + 0) > 0 \quad (29)$$

である.

(iii)  $\underline{(i_\alpha + s) - i_{\alpha-1}} < b_y \leq \underline{i_{\alpha+1} - (i_\alpha + s)}$  のとき.

$$\begin{aligned}
& (d_3(b_y) + d_4(b_y)) - (d'_3(b_y) + d'_4(b_y)) \\
&= ((b_y - (\underline{i_\alpha - i_{\alpha-1}})) + 0) - ((b_y - ((i_\alpha + s) - i_{\alpha-1})) + 0) \\
&= ((i_\alpha + s) - i_{\alpha-1}) - (i_\alpha - i_{\alpha-1}) \\
&= ((i_\alpha + s) - i_\alpha) + (\underline{i_\alpha - i_{\alpha-1}}) - (\underline{i_\alpha - i_{\alpha-1}}) \quad (\because i_{\alpha-1} \preceq i_\alpha \preceq \underline{i_\alpha + s} \preceq \text{と式 (6)}) \\
&= \underline{i_\alpha + s} - i_\alpha \\
&= s \\
&> 0
\end{aligned} \tag{30}$$

である。

(iv)  $i_{\alpha+1} - (i_\alpha + s) < b_y \leq i_{\alpha+1} - i_\alpha$  のとき。

$$\begin{aligned}
& (d_3(b_y) + d_4(b_y)) - (d'_3(b_y) + d'_4(b_y)) \\
&= ((b_y - (\underline{i_\alpha - i_{\alpha-1}})) + 0) - ((b_y - ((i_\alpha + s) - i_{\alpha-1})) + (b_y - (\underline{i_{\alpha+1} - (i_\alpha + s)}))) \\
&= -b_y + (\underline{i_{\alpha+1} - (i_\alpha + s)}) + ((i_\alpha + s) - i_{\alpha-1}) - (i_\alpha - i_{\alpha-1}) \\
&= -b_y + (\underline{i_{\alpha+1} - i_{\alpha-1}}) - (\underline{i_\alpha - i_{\alpha-1}}) \quad (\because i_{\alpha-1} \preceq \underline{i_\alpha + s} \preceq i_{\alpha+1} \preceq \text{と式 (6)}) \\
&= -b_y + (\underline{i_{\alpha+1} - i_\alpha}) + (\underline{i_\alpha - i_{\alpha-1}}) - (\underline{i_\alpha - i_{\alpha-1}}) \quad (\because i_{\alpha-1} \preceq i_\alpha \preceq i_{\alpha+1} \preceq \text{と式 (6)}) \\
&= -b_y + (\underline{i_{\alpha+1} - i_\alpha}) \\
&\geq 0
\end{aligned} \tag{31}$$

である。

(v)  $i_{\alpha+1} - i_\alpha < b_y \leq n - 1$  のとき。

$$\begin{aligned}
& (d_3(b_y) + d_4(b_y)) - (d'_3(b_y) + d'_4(b_y)) \\
&= ((b_y - (\underline{i_\alpha - i_{\alpha-1}})) + (b_y - (\underline{i_{\alpha+1} - i_\alpha})) \\
&\quad - ((b_y - ((i_\alpha + s) - i_{\alpha-1})) + (b_y - (\underline{i_{\alpha+1} - (i_\alpha + s)})))) \\
&= ((i_{\alpha+1} - (i_\alpha + s)) + ((i_\alpha + s) - i_{\alpha-1})) - ((i_{\alpha+1} - i_\alpha) + (\underline{i_\alpha - i_{\alpha-1}})) \\
&= (\underline{i_{\alpha+1} - i_{\alpha-1}}) - (\underline{i_{\alpha+1} - i_{\alpha-1}}) \quad (\because i_{\alpha-1} \preceq i_\alpha \preceq \underline{i_\alpha + s} \preceq i_{\alpha+1} \preceq \text{と式 (6)}) \\
&= 0
\end{aligned} \tag{32}$$

である。

(vi)  $b_y = n$  のとき。

$$(d_3(b_y) + d_4(b_y)) - (d'_3(b_y) + d'_4(b_y)) = (1 + 1) - (1 + 1) = 0 \tag{33}$$

である。

以上の式 (21)(22)(28)(29)(30)(31)(32)(33) より，

$$D_3 + D_4 = \sum_{b_y=0}^n (d_3(b_y) + d_4(b_y)) > \sum_{b_y=0}^n (d'_3(b_y) + d'_4(b_y)) = D'_3 + D'_4 \tag{34}$$

が成り立つ。なお，式 (34) における大小関係が  $\geq$  ではなく  $>$  になるのは，式 (29) において  $>$  が入っており，かつ，式 (27) より，式 (29) を満たす  $b_y$  が少なくとも 1 個存在するためである\*1。

式 (16)(17)(18)(19)(20)(34) より，

$$\begin{aligned}
M(C^{\sigma_y}(A); S_x) &= D_1 + D_2 - (D_3 + D_4 - D_5) \\
&< D'_1 + D'_2 - (D'_3 + D'_4 - D'_5) \\
&= M(C^{\sigma_y}(A); \text{extend}(S_x, i_\alpha, s))
\end{aligned}$$

が成り立つ。以上より，補題 3 が示された。■

以上によって，補題 1，補題 2，補題 3 が示された。次に，以下の補題を考えて証明する：  
補題 4  $T^0 = \{\sigma_y(0), \sigma_y(1), \dots, \sigma_y(b_x - 1)\}$  とする。このとき，アドレス集合  $T^0$  に対して，*shift* または *mirror* または *extend* の写像を有限回適用することによって， $P_{b_x}(A)$  に属する任意のアドレス集合  $S_x \in P_{b_x}(A)$  を構成することができる。

補題 4 を示すために，いくつか準備を行う。 $S_x \in P_{b_x}(A)$  なる任意の  $S_x$  は， $i_0 < i_1 < \dots < i_{b_x-1} <$  なる整数  $i_0, i_1, \dots, i_{b_x-1}$  を用いて，

$$S_x = \{\sigma_y(i_0), \sigma_y(i_1), \dots, \sigma_y(i_{b_x-1})\}$$

と表すことができる (図 30(A))。このとき，すべての  $j$  ( $0 \leq j \leq b_x - 1$ ) に関して， $\underline{i_k - i_{k-1}} \geq \underline{i_j - i_{j-1}}$  を満たすような  $k$  ( $0 \leq k \leq b_x - 1$ ) が少なくとも 1 個存在するので，そのような  $k$  のうちの 1 個を  $\alpha$  とおく。すなわち， $\alpha$  は，

$$\forall j (0 \leq j \leq b_x - 1) : \underline{i_\alpha - i_{\alpha-1}} \geq \underline{i_j - i_{j-1}} \tag{35}$$

を満たす。図形的には，点の間の距離が最大になるような 2 点を  $i_\alpha$  と  $i_{\alpha-1}$  とおいている。

さらに，以下の操作  $O$  を考える：

$$v \leftarrow b_x - 1$$

$$T \leftarrow T^0$$

$$k \leftarrow 0$$

**while**  $k < b_x - 1$  **do**

$$T \leftarrow \text{mirror}(T, k) \quad /* \text{step1} */$$

\*1 式 (30) でも  $>$  が入っているが，式 (27) より，式 (30) を満たす  $b_y$  は存在しない可能性がある。

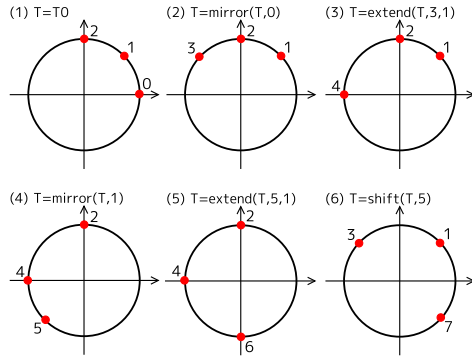


図 31  $T^0 = \{\sigma_y(0), \sigma_y(1), \sigma_y(2)\}$  に対して操作  $O$  を適用することで,  $S_x = \{\sigma_y(1), \sigma_y(3), \sigma_y(7)\}$  を得るまでの手続き.

```

v ← v + 1      /* step2 */
if  $i_{\alpha+k+1} - i_{\alpha+k} - 1 \neq 0$  then
  T ← extend(T, v,  $i_{\alpha+k+1} - i_{\alpha+k} - 1$ )  /* step3 */
endif
v ←  $v + i_{\alpha+k+1} - i_{\alpha+k} - 1$   /* step4 */
k ← k + 1
endwhile
T ← shift(T,  $i_\alpha - (b_x - 1)$ )  /* step5 */

```

操作  $O$  の図形的意味を確認するために, たとえば,  $n = 4, b_x = 3$  とし,  $T_0 = \{\sigma_y(0), \sigma_y(1), \sigma_y(2)\}$  に対して操作  $O$  を適用することで,  $S_x = \{\sigma_y(1), \sigma_y(3), \sigma_y(7)\}$  を得るまでの手続きを図 31 に示す. いまの場合,  $i_0 = 1, i_1 = 3, i_2 = 7$  に関して,  $i_0 - i_2 \geq i_2 - i_1, i_0 - i_2 \geq i_1 - i_0$  であるから,  $i_\alpha = i_0$  であることに注意したい.

明らかに, 操作  $O$  は, *shift* または *mirror* または *extend* の写像を有限回適用することで終了する. したがって, 補題 4 を示すためには, 操作  $O$  が終了したとき  $T$  が  $S_x$  に一致することを示せば良い. そこで, 以下の補題を考える:

補題 5 操作  $O$  における  $k$  回目のループの先頭において, 以下のループ不変条件が成立する:

$$v = b_x - 1 + i_{\alpha+k} - i_\alpha, \quad (36)$$

$$T = \{\sigma_y(k), \sigma_y(k+1), \dots, \sigma_y(b_x - 1), \sigma_y(b_x - 1 + i_{\alpha+1} - i_\alpha), \dots, \sigma_y(b_x - 1 + i_{\alpha+k} - i_\alpha)\}. \quad (37)$$

補題 5 を数学的帰納法で示す. まず,  $k = 0$  の場合には,

$$v = b_x - 1,$$

$$T = T^0 = \{\sigma_y(0), \sigma_y(1), \dots, \sigma_y(b_x - 1)\}$$

であるから, 明らかにループ不変条件 (36)(37) が成り立つ.

次に,  $k$  の場合にループ不変条件 (36)(37) が成り立つことを仮定して,  $k + 1$  の場合にもループ不変条件 (36)(37) が成り立つことを言う.

第一に, 式 (36) について考える.  $k$  回目のループの先頭において,  $v = b_x - 1 + i_{\alpha+k} - i_\alpha$  が成り立つとすれば, step2 の操作によって,  $v$  は,

$$v = b_x - 1 + i_{\alpha+k} - i_\alpha + 1 = b_x + i_{\alpha+k} - i_\alpha \quad (\because \text{式 (4)}) \quad (38)$$

に変化する. さらに step4 の操作によって,  $v$  は,

$$v = b_x + i_{\alpha+k} - i_\alpha + i_{\alpha+k+1} - i_{\alpha+k} - 1 = b_x - 1 + i_{\alpha+k+1} - i_\alpha \quad (\because \text{式 (4)})$$

に変化し, これが  $k + 1$  回目のループの先頭で成り立つ. したがって, 式 (36) がループ不変条件であることが示された.

第二に, 式 (37) について考える.  $k$  回目のループの先頭において,

$$T = \{\sigma_y(k), \sigma_y(k+1), \dots, \sigma_y(b_x - 1), \sigma_y(b_x - 1 + i_{\alpha+1} - i_\alpha), \dots, \sigma_y(b_x - 1 + i_{\alpha+k} - i_\alpha)\}$$

が成り立つことを仮定する. step1 の操作によって,  $T$  は,

$$\begin{aligned} T &= \{\sigma_y(k+1), \sigma_y(k+2), \dots, \sigma_y(b_x - 1), \sigma_y(b_x - 1 + i_{\alpha+1} - i_\alpha), \dots, \\ &\quad \sigma_y(b_x - 1 + i_{\alpha+k} - i_\alpha), \sigma_y(b_x - 1 + i_{\alpha+k} - i_\alpha + (k+1) - k)\} \\ &= \{\sigma_y(k+1), \sigma_y(k+2), \dots, \sigma_y(b_x - 1), \sigma_y(b_x - 1 + i_{\alpha+1} - i_\alpha), \dots, \\ &\quad \sigma_y(b_x - 1 + i_{\alpha+k} - i_\alpha), \sigma_y(b_x + i_{\alpha+k} - i_\alpha)\} \end{aligned} \quad (39)$$

に変化する. 次にこの  $T$  に対して step3 の操作を適用することを考えるが, step3 の操作を適用する前に, この時点で写像 *extend* を適用できるための条件が成立していることを確認しなければならない. 写像 *extend* の定義により,

$$T = \{\sigma_y(\tilde{i}_0), \sigma_y(\tilde{i}_1), \dots, \sigma_y(\tilde{i}_{b_x-1})\} \quad (\tilde{i}_0 < \tilde{i}_1 < \dots < \tilde{i}_{b_x-1} <)$$

に対して，写像  $extend(T, \tilde{i}_\alpha, s)$  を適用するためには，

$$\sigma_y(\tilde{i}_\alpha) \in T, \quad (40)$$

$$\tilde{i}_\alpha \prec \tilde{i}_\alpha + s \prec \tilde{i}_{\alpha+1} \prec, \quad (41)$$

$$\underline{(\tilde{i}_\alpha + s) - \tilde{i}_{\alpha-1}} \leq \underline{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)} \quad (42)$$

の3条件が満足されねばならない．そこで，式(39)に対してstep3の操作を適用する時点で確かに式(40)(41)(42)が満足されていることを確認する．

式(39)に対してstep3の操作を適用する時点では，式(38)より  $v = \underline{b_x + i_{\alpha+k} - i_\alpha}$  になっているから，この時点における  $T, \tilde{i}_{\alpha-1}, \tilde{i}_\alpha, \tilde{i}_{\alpha+1}, s$  の値は，それぞれ，

$$T = \{\sigma_y(k+1), \sigma_y(k+2), \dots, \sigma_y(b_x-1), \sigma_y(\underline{b_x-1+i_{\alpha+1}-i_\alpha}), \dots, \sigma_y(\underline{b_x-1+i_{\alpha+k}-i_\alpha}), \sigma_y(\underline{b_x+i_{\alpha+k}-i_\alpha})\},$$

$$\tilde{i}_{\alpha-1} = \underline{b_x-1+i_{\alpha+k}-i_\alpha},$$

$$\tilde{i}_\alpha = v = \underline{b_x+i_{\alpha+k}-i_\alpha},$$

$$\tilde{i}_{\alpha+1} = k+1,$$

$$s = \underline{i_{\alpha+k+1}-i_{\alpha+k}-1}$$

になっていることに注意する．

(I) 式(40)が成り立つことを確認する． $\sigma_y(v) \in T$  が成り立つので，式(40)は成り立つ．

(II) 式(41)が成り立つことを確認する．そのためにまず， $\underline{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)}$  と  $\underline{\tilde{i}_{\alpha+1} - \tilde{i}_\alpha}$  を計算する．

$\underline{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)}$  については，

$$\begin{aligned} \underline{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)} &= \underline{(k+1) - ((b_x + i_{\alpha+k} - i_\alpha) + (i_{\alpha+k+1} - i_{\alpha+k} - 1))} \\ &= \underline{(k+1) - (b_x - 1) - (i_{\alpha+k+1} - i_{\alpha+k}) + (i_{\alpha+k} - i_\alpha)} \quad (\because \text{式(4)}) \\ &= \underline{(k+1) - (b_x - 1) - (i_{\alpha+k+1} - i_\alpha)} \\ &\quad (\because i_\alpha \preceq i_{\alpha+k} \preceq i_{\alpha+k+1} \preceq \text{と式(6)}) \\ &= \underline{n + (k+1) - (b_x - 1) - (i_{\alpha+k+1} - i_\alpha)} \quad (43) \end{aligned}$$

が得られる．ここで式(43)がとりうる値の範囲を考えると，

$$\begin{aligned} n &> n + (k+1) - (b_x - 1) - (i_{\alpha+k+1} - i_\alpha) \quad (\because k \leq b_x - 2 \text{ および } i_{\alpha+k+1} \neq i_\alpha) \\ &\geq n - \underline{(i_{\alpha+b_x-1} - i_{\alpha+k+1})} - \underline{(i_{\alpha+k+1} - i_\alpha)} \\ &\quad (\because i_{\alpha+k+1} \prec i_{\alpha+k+2} \prec \dots \prec i_{\alpha+b_x-1} \prec \text{より}) \\ \underline{i_{\alpha+b_x-1} - i_{\alpha+k+1}} &\geq (b_x - 1) - (k+1) \end{aligned}$$

$$= n - \underline{(i_{\alpha+b_x-1} - i_\alpha)} \quad (\because i_\alpha \preceq i_{\alpha+k+1} \preceq i_{\alpha+b_x-1} \preceq \text{と式(6)})$$

$$= n - \underline{(i_{\alpha-1} - i_\alpha)} \quad (\because \text{mod}(\alpha + b_x - 1, b_x) = \text{mod}(\alpha - 1, b_x))$$

$$= \underline{i_\alpha - i_{\alpha-1}} \quad (\because i_{\alpha-1} \preceq i_\alpha \preceq \text{と式(5)})$$

$$> 0 \quad (\because i_\alpha \neq i_{\alpha-1}) \quad (44)$$

が成り立つ．よって，式(3)(43)(44)より，

$$\underline{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)} = n + (k+1) - (b_x - 1) - \underline{(i_{\alpha+k+1} - i_\alpha)} \quad (45)$$

が成り立つ．

$\underline{\tilde{i}_{\alpha+1} - \tilde{i}_\alpha}$  については，

$$\begin{aligned} \underline{\tilde{i}_{\alpha+1} - \tilde{i}_\alpha} &= \underline{(k+1) - (b_x + i_{\alpha+k} - i_\alpha)} \\ &= \underline{k - (b_x - 1) - (i_{\alpha+k} - i_\alpha)} \quad (\because \text{式(4)}) \\ &= \underline{n + k - (b_x - 1) - (i_{\alpha+k} - i_\alpha)} \quad (46) \end{aligned}$$

となる．この式(46)は，式(43)における  $k+1$  を  $k$  に置き換えたものであるから，式(44)を導いたときの議論と同様にして，

$$n > n + k - (b_x - 1) - \underline{(i_{\alpha+k} - i_\alpha)} > 0 \quad (47)$$

と言える．よって，式(3)(46)(47)より，

$$\underline{\tilde{i}_{\alpha+1} - \tilde{i}_\alpha} = n + k - (b_x - 1) - \underline{(i_{\alpha+k} - i_\alpha)} \quad (48)$$

が成り立つ．

以上の結果をもとにして， $\underline{(\tilde{i}_\alpha + s) - \tilde{i}_\alpha} + \underline{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)} + \underline{\tilde{i}_\alpha - \tilde{i}_{\alpha+1}}$  を計算すると，

$$\begin{aligned} &\underline{(\tilde{i}_\alpha + s) - \tilde{i}_\alpha} + \underline{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)} + \underline{\tilde{i}_\alpha - \tilde{i}_{\alpha+1}} \\ &= \underline{s + \tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)} + \underline{(n - \tilde{i}_{\alpha+1} - \tilde{i}_\alpha)} \quad (\because \text{式(4)}, i_\alpha \preceq i_{\alpha+1} \preceq \text{と式(5)}) \\ &= \underline{(i_{\alpha+k+1} - i_{\alpha+k} - 1)} + \underline{(n + (k+1) - (b_x - 1) - (i_{\alpha+k+1} - i_\alpha))} \\ &\quad + \underline{(n - (n + k - (b_x - 1) - (i_{\alpha+k} - i_\alpha)))} \quad (\because \text{式(45)(48)}) \\ &= \underline{((i_{\alpha+k+1} - i_{\alpha+k} - 1) + 1)} + \underline{(n - (i_{\alpha+k+1} - i_\alpha))} + \underline{(i_{\alpha+k} - i_\alpha)} \\ &= \underline{((i_{\alpha+k+1} - i_{\alpha+k}) - 1 + 1)} + \underline{(i_\alpha - i_{\alpha+k+1})} + \underline{(i_{\alpha+k} - i_\alpha)} \\ &\quad (\because i_{\alpha+k+1} \neq i_{\alpha+k}, i_\alpha \preceq i_{\alpha+k+1} \preceq \text{と式(5)}) \\ &= n \quad (\because i_\alpha \preceq i_{\alpha+k} \preceq i_{\alpha+k+1} \preceq \text{と式(5)}) \end{aligned}$$

となる．したがって，式(5)より， $\tilde{i}_\alpha \preceq \underline{\tilde{i}_\alpha + s} \preceq \tilde{i}_{\alpha+1} \preceq$  が成り立つ．さらに，式(44)(45)

より  $\tilde{i}_{\alpha+1} \neq \tilde{i}_\alpha + s$  であり, 式 (47)(48) より  $\tilde{i}_{\alpha+1} \neq \tilde{i}_\alpha$  である. また, 操作  $O$  の定義より step3 を適用できるのは  $s = \frac{i_{\alpha+k+1} - i_{\alpha+k} - 1}{1} \neq 0$  のときであるから,  $\tilde{i}_\alpha \neq \tilde{i}_\alpha + s$  である. 結局,  $\tilde{i}_\alpha < \tilde{i}_\alpha + s < \tilde{i}_{\alpha+1} <$  が成り立つ. つまり, 式 (41) が成り立つことが確認できた.

(Ⅲ) 式 (42) が成り立つことを確認する. ここで, 式 (42) の左辺である  $\frac{\tilde{i}_\alpha + s}{1} - \tilde{i}_{\alpha-1}$  を計算すると,

$$\begin{aligned} \frac{\tilde{i}_\alpha + s}{1} - \tilde{i}_{\alpha-1} &= \frac{(b_x + i_{\alpha+k} - i_\alpha) + (i_{\alpha+k+1} - i_{\alpha+k} - 1) - (b_x - 1 + i_{\alpha+k} - i_\alpha)}{i_{\alpha+k+1} - i_{\alpha+k}} \quad (\because \text{式 (4)}) \\ &= \frac{i_{\alpha+k+1} - i_{\alpha+k}}{i_{\alpha+k+1} - i_{\alpha+k}} \end{aligned}$$

が得られる. 一方で, 右辺の  $\frac{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)}{1}$  に関しては, 式 (44)(45) より,

$$\frac{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)}{1} = n + (k+1) - (b_x - 1) - (i_{\alpha+k+1} - i_\alpha) \geq i_\alpha - i_{\alpha-1} \quad (50)$$

が成り立つ. よって, 式 (35)(49)(50) より,

$$\frac{\tilde{i}_\alpha + s}{1} - \tilde{i}_{\alpha-1} \leq \frac{\tilde{i}_{\alpha+1} - (\tilde{i}_\alpha + s)}{1}$$

が成り立つ. つまり, 式 (42) が成り立つことが確認できた.

以上の (Ⅰ)(Ⅱ)(Ⅲ) より, 式 (39) に対して step3 の操作を適用する時点で, 確かに式 (40)(41)(42) が満足されていることを確認できた.

そこで, 式 (39) に対して step3 の操作を適用すると,  $T$  は,

$$\begin{aligned} T &= \{\sigma_y(k+1), \sigma_y(k+2), \dots, \sigma_y(b_x - 1), \sigma_y(\frac{b_x - 1 + (i_{\alpha+1} - i_\alpha)}{1}), \dots, \\ &\quad \sigma_y(\frac{b_x - 1 + i_{\alpha+k} - i_\alpha}{1}), \sigma_y(\frac{b_x + i_{\alpha+k} - i_\alpha + i_{\alpha+k+1} - i_{\alpha+k} - 1}{1})\} \\ &= \{\sigma_y(k+1), \sigma_y(k+2), \dots, \sigma_y(b_x - 1), \sigma_y(\frac{b_x - 1 + (i_{\alpha+1} - i_\alpha)}{1}), \dots, \\ &\quad \sigma_y(\frac{b_x - 1 + i_{\alpha+k} - i_\alpha}{1}), \sigma_y(\frac{b_x - 1 + i_{\alpha+k+1} - i_\alpha}{1})\} \quad (\because \text{式 (4)}) \end{aligned}$$

に変化し, これが  $k+1$  回目のループの先頭で成り立つ. 以上により, 式 (37) がループ不変条件であることが示された. つまり, 補題 5 が示された. ■

補題 4 を示す. 補題 5 より, 操作  $O$  におけるループを抜けた直後の  $T$  は, ループ不変条件において  $k = b_x - 2$  としたものの, すなわち,

$$T = \{\sigma_y(b_x - 1), \sigma_y(\frac{b_x - 1 + i_{\alpha+1} - i_\alpha}{1}), \dots, \sigma_y(\frac{b_x - 1 + i_{\alpha+b_x-1} - i_\alpha}{1})\}$$

となっている. よって, この  $T$  に対して step5 の操作を適用すると,  $T$  は,

$$\begin{aligned} T &= \{\sigma_y(b_x - 1 + i_\alpha - (b_x - 1)), \sigma_y(\frac{b_x - 1 + i_{\alpha+1} - i_\alpha + i_\alpha - (b_x - 1)}{1}), \dots, \\ &\quad \sigma_y(\frac{b_x - 1 + i_{\alpha+b_x-1} - i_\alpha + i_\alpha - (b_x - 1)}{1})\} \\ &= \{\sigma_y(i_\alpha), \sigma_y(i_{\alpha+1}), \dots, \sigma_y(\frac{i_{\alpha+b_x-1}}{1})\} \quad (\because \text{式 (4)}) \\ &= \{\sigma_y(i_0), \sigma_y(i_1), \dots, \sigma_y(i_{b_x-1})\} \\ &= S_x \end{aligned}$$

に変化する. 以上によって, 操作  $O$  が終了したとき  $T$  が  $S_x$  に一致することが示された. つまり, 補題 4 が示された. ■

以上の系 1, 補題 2, 補題 3, 補題 4 を用いて, 命題 2 を示す. 系 1, 補題 2 より, 任意のアドレス集合  $S_x \in P_{b_x}(A)$  に対して写像 *shift* または写像 *mirror* を 1 回以上の任意回適用したアドレス集合を  $S'_x$  とすると,  $M(C^{\sigma_y}(A); S'_x) = M(C^{\sigma_y}(A); S_x)$  が成り立つ. また, 補題 3 より, 任意のアドレス集合  $S_x \in P_{b_x}(A)$  に対して写像 *extend* を 1 回以上の任意回適用したアドレス集合を  $S'_x$  とすると,  $M(C^{\sigma_y}(A); S'_x) > M(C^{\sigma_y}(A); S_x)$  が成り立つ. さらに, 補題 4 より, 任意のアドレス集合  $S_x \in P_{b_x}(A)$  は,  $T^0$  に対して, 写像 *shift* または写像 *mirror* または写像 *extend* を有限回適用することによって得られる. これらの事実を総合すると,  $M(C^{\sigma_y}(A); S_x)$  を最小化する  $S_x$  とは,  $T^0$  に対して, 写像 *shift* または写像 *mirror* を有限回適用して得られるアドレス集合のみであると言える. そのようなアドレス集合とは, 具体的には,

$$\begin{aligned} T^0 &= \{\sigma_y(0), \sigma_y(1), \dots, \sigma_y(b_x - 2), \sigma_y(b_x - 1)\}, \\ T^1 &= \{\sigma_y(1), \sigma_y(2), \dots, \sigma_y(b_x - 1), \sigma_y(b_x)\}, \\ &\quad \vdots \\ T^{n-1} &= \{\sigma_y(n - b_x + 1), \sigma_y(n - b_x + 2), \dots, \sigma_y(n - 1), \sigma_y(0)\}. \end{aligned}$$

のことであり, これは  $C_{b_x}^{\sigma_y}(A)$  にほかならない. 以上によって,  $P_{b_x}(A)$  に属するすべてのアドレス集合  $S_x \in P_{b_x}(A)$  のうち,  $M(C^{\sigma_y}(A); S_x)$  を最小化する  $S_x$  の集合は  $C_{b_x}^{\sigma_y}(A)$  であることが示された. つまり, 命題 2 が示された. ■

命題 2 が示されたので, その言い換えである命題 1 も成り立つ. よって, 命題 1 が示された. ■

命題 1 が成り立つことを用いて, 定理 2 を証明する.  $m$  個のスレッド  $x_0, x_1, \dots, x_{m-1}$  を考え, 各スレッド  $x_i$  は置換  $\sigma_i$  に従って一様にアドレスを使用するとする. また, スレッド  $x_{i_0}, x_{i_1}, \dots, x_{i_{k-1}}$  に関して, 「どの 2 つの異なるスレッド  $x_i$  とスレッド



$x_j$  ( $x_i, x_j \in \{x_{i_0}, x_{i_1}, \dots, x_{i_{k-1}}\}$ ) に対しても, スレッド  $x_i$  が使用するアドレス集合とスレッド  $x_j$  が使用するアドレス集合が共通部分を持たない事象」を  $E(x_{i_0}, x_{i_1}, \dots, x_{i_{k-1}})$  と表す. また, 事象  $E(x_{i_0}, x_{i_1}, \dots, x_{i_{k-1}})$  が起きる確率を  $p(E(x_{i_0}, x_{i_1}, \dots, x_{i_{k-1}}))$  と表す. たとえば,  $p(E(x_0, x_1))$  は, スレッド  $x_0$  が使用するアドレス集合とスレッド  $x_1$  が使用するアドレス集合が共通部分を持たない確率を意味する.  $p(E(x_0, x_1, x_2)) = p(E(x_0, x_1) \wedge E(x_1, x_2) \wedge E(x_2, x_0))$  などが成り立つ.

このとき, 定理 2 が成り立つことを数学的帰納法で示す.

まず,  $m = 1$  の場合には明らかに定理 2 は成り立つ. また,  $m = 2$  の場合には, 命題 1 より定理 2 は成り立つ.

次に,  $m$  ( $m \geq 2$ ) のときに定理 2 が成り立つことを仮定して,  $m + 1$  のときにも定理 2 が成り立つことを言う.  $p(E(x_0, x_1, \dots, x_{m-1}, x_m))$  を, 条件付き確率を用いて分解すると,

$$p(E(x_0, x_1, \dots, x_{m-1}, x_m)) = p(E(x_0, x_1, \dots, x_{m-1}))p(E(x_0, x_1, \dots, x_{m-1}) | (E(x_m, x_0) \wedge E(x_m, x_1) \wedge \dots \wedge E(x_m, x_{m-1}))) \quad (51)$$

$$(E(x_m, x_0) \wedge E(x_m, x_1) \wedge \dots \wedge E(x_m, x_{m-1}))) \quad (52)$$

となる. いま, 各スレッド  $x_i$  の置換の選び方は独立であることを用いて式 (52) を計算すると,

$$\begin{aligned} & p(E(x_0, x_1, \dots, x_{m-1}, x_m)) \\ &= p(E(x_0, x_1, \dots, x_{m-1}))p(E(x_m, x_0) \wedge E(x_m, x_1) \wedge \dots \wedge E(x_m, x_{m-1})) \\ &= p(E(x_0, x_1, \dots, x_{m-1}))p(E(x_m, x_0))p(E(x_m, x_1)) \dots p(E(x_m, x_{m-1})) \end{aligned} \quad (53)$$

となる. 式 (53) において,  $p(E(x_0, x_1, \dots, x_{m-1}))$  が最大になるのは, 数学的帰納法の仮定より  $\sigma_0 = \sigma_1 = \dots = \sigma_{m-1}$  のときにかぎられる. また, 式 (53) における各  $p(E(x_m, x_i))$  ( $0 \leq i \leq m-1$ ) が最大になるのは, 命題 1 より  $\sigma_m = \sigma_i$  のときにかぎられる. すなわち, 式 (53) が最大になるのは,  $\sigma_0 = \sigma_1 = \dots = \sigma_{m-1} = \sigma_m$  のときにかぎられる. したがって, 定理 2 が成り立つ.

以上より, 定理 2 が成り立つことが証明できた. ■

### A.1.3 アドレス衝突確率の定量的な評価

以上の証明より, 2 つのアドレス集合  $S_x$  と  $S'_x$  が与えられたとき, どちらの方がどれくらいアドレスが衝突しやすいのかを定量的に計算することができる. スレッド  $x$  とスレッド  $y$  を考え, スレッド  $y$  は置換  $\sigma_y$  に従って一様にアドレス集合を使用しているとす. また, スレッド  $x$  は  $b_x$  個のアドレスを割り当てようとしているとする. このとき, スレッド  $x$  が,  $b_x$  個のアドレスを割り当てるために, アドレス集合  $S_x^0 \in P_{b_x}(A)$  を使

用する場合とアドレス集合  $S_x^1 \in P_{b_x}(A)$  を使用する場合とでは, 前者の方が後者より,  $M(C^{\sigma_y}(A); S_x^0) - M(C^{\sigma_y}(A); S_x^1)$  だけ, スレッド  $y$  が使用するアドレス集合とアドレスが衝突しやすい\*1. 以下では,  $M(C^{\sigma_y}(A); S_x^0) - M(C^{\sigma_y}(A); S_x^1)$  の値を計算する.

まず, 任意のアドレス集合  $S_x = \{\sigma_y(i_0), \sigma_y(i_1), \dots, \sigma_y(i_{b_x-1})\}$  ( $i_0 < i_1 < \dots < i_{b_x-1}$ ) に対して,  $M(C^{\sigma_y}(A); \text{extend}(S_x, i_\alpha, s)) - M(C^{\sigma_y}(A); S_x)$  の値を計算すると, 式 (16)(17)(18)(19)(20) (21)(22)(28)(29) (30)(31)(32)(33) より,

$$\begin{aligned} & M(C^{\sigma_y}(A); \text{extend}(S_x, i_\alpha, s)) - M(C^{\sigma_y}(A); S_x) \\ &= \sum_{b_y=i_\alpha-i_{\alpha-1}+1}^{(i_\alpha+s)-i_{\alpha-1}} (b_y - (i_\alpha - i_{\alpha-1})) + \sum_{b_y=(i_\alpha+s)-i_{\alpha-1}+1}^{i_{\alpha+1}-(i_\alpha+s)} \quad (s) \\ &+ \sum_{b_y=i_{\alpha+1}-(i_\alpha+s)+1}^{i_{\alpha+1}-i_\alpha} (-b_y + (i_{\alpha+1} - i_\alpha)) \end{aligned} \quad (54)$$

となる.

ここで, アドレス集合  $T^0$  に対して操作  $O$  を適用することによってアドレス集合  $S_x^0 = \{\sigma_y(i_0^0), \sigma_y(i_1^0), \dots, \sigma_y(i_{b_x-1}^0)\}$  ( $i_0^0 < i_1^0 < \dots < i_{b_x-1}^0$ ) および  $S_x^1 = \{\sigma_y(i_0^1), \sigma_y(i_1^1), \dots, \sigma_y(i_{b_x-1}^1)\}$  ( $i_0^1 < i_1^1 < \dots < i_{b_x-1}^1$ ) を構成することを考える. すると,  $M(C^{\sigma_y}(A); S_x^0) - M(C^{\sigma_y}(A); S_x^1)$  は,

$$\begin{aligned} & M(C^{\sigma_y}(A); S_x^0) - M(C^{\sigma_y}(A); S_x^1) \\ &= (M(C^{\sigma_y}(A); S_x^0) - M(C^{\sigma_y}(A); T^0)) - (M(C^{\sigma_y}(A); S_x^1) - M(C^{\sigma_y}(A); T^0)) \end{aligned} \quad (55)$$

と表せる. 式 (55) において,  $M(C^{\sigma_y}(A); S_x^0) - M(C^{\sigma_y}(A); T^0)$  は,  $T^0$  から操作  $O$  によって  $S_x^0$  を構成するときに, 操作  $O$  の step3 を適用するたびに式 (54) を計算し, その総和を求めることで得られる. 同様に,  $M(C^{\sigma_y}(A); S_x^1) - M(C^{\sigma_y}(A); T^0)$  は,  $T^0$  から操作  $O$  によって  $S_x^1$  を構成するときに, 操作  $O$  の step3 を適用するたびに式 (54) を計算し, その総和を求めることで得られる. このように,  $M(C^{\sigma_y}(A); S_x^0) - M(C^{\sigma_y}(A); S_x^1)$  は単純な式にはならないが, 実際の数値を代入することで計算可能である.

\*1 ただし,  $M(C^{\sigma_y}(A); S_x^0) - M(C^{\sigma_y}(A); S_x^1)$  の次元は確率ではない.