

単一システムイメージを提供するための仮想マシンモニタ

金田 憲 二^{†,††} 大山 恵 弘[†] 米 澤 明 憲[†]

我々は、クラスタ上に共有メモリ型マルチプロセッサマシンを仮想的に構築する仮想マシンモニタを設計・実装した。この仮想マシンの機能によって、クラスタを非常に簡便に利用することが可能となる。たとえば、共有メモリ型マルチプロセッサマシン用の並列アプリケーションを、無変更のままクラスタ上で実行することが可能になる。さらに、マルチプロセッサをサポートする OS (たとえば Linux) を、少量の変更で仮想マシンにインストールすることができる。我々は、8 台の物理マシン上に仮想的に 8-way のマルチプロセッサマシンを構築した。そして、その仮想マシン上で互いに独立な粗粒度タスクを並列に実行し、その実行時間を測定した。この実験の結果は、我々のアプローチで現実的な性能を達成できることを示している。

A Virtual Machine Monitor for Providing a Single System Image

KENJI KANEDA,^{†,††} YOSHIHIRO OYAMA[†] and AKINORI YONEZAWA[†]

We have designed and implemented a virtual machine monitor that virtualizes a shared-memory multi-processor machine on a commodity cluster. This functionality greatly simplifies utilization of commodity clusters. For example, it enables parallel applications for shared-memory multi-processor systems to run on clusters without any change of the applications. Moreover, commodity operating systems that support multi-processors (e.g., Linux) can be installed in a virtual machine with a small amount of modification. We built a virtual 8-way multi-processor machine on eight physical machines. We ran parallel coarse-grain tasks on Linux installed in the virtual machine and measured the execution time. The experimental result demonstrates the feasibility of our approach.

1. はじめに

コモディークラスタは、ネットワーク技術の進化や PC の価格低下にとともに、近年その重要度を増しつつある。たとえば、これまで並列計算などに携わることのなかった一般の自然科学の研究者が、個人・ワークグループ用として 4~32 ノード程度のクラスタを所有する、ということも現実にも可能となっている。

しかし、こうしたクラスタの利用に際しては、その利便性が低いという重大な問題が依然として存在する。たとえば、クラスタの各ノードに Linux などの通常の OS をインストールしただけでは、複数のマシンに分散した CPU やディスクなどの資源を大域的に管理する機構が欠如している。そのため(特に計算機科学を専門としない一般の人にとって)クラスタを効率的に

利用するのは難しいものとなっている。

そこで、本研究の目標として、こうしたクラスタ上に単一システムイメージ (Single System Image: SSI) を構築し、クラスタの利便性を向上させることを目指す。SSI の実現方法としては様々な既存研究 (たとえば SCORE¹⁾) が存在するが、本研究では特に、仮想マシンモニタ (Virtual Machine Monitor: VMM) を利用するというアプローチをとる。VMM というのは、実マシンと同等の処理 (たとえば OS の実行) が可能な仮想マシンを、実機上に構築するミドルウェアシステムである。CPU・メモリ・I/O デバイスなどのハードウェアをソフトウェアでエミュレーションすることで、これを実現している。有名な VMM としては、たとえば VMware Workstation²⁾ があげられる。

我々は、より具体的には、「ネットワークで結合された複数のマシン上に、共有メモリ型マルチプロセッサマシンを仮想的に構築する VMM」を設計・実装する。この VMM を、我々は *Virtual Multiprocessor* と呼ぶ。Virtual Multiprocessor は、たとえば N 台のシングルプロセッサマシン上に、 N プロセッサからな

[†] 東京大学大学院情報理工学系研究科
Graduate School of Information Science and Technology,
The University of Tokyo

^{††} 日本学術振興会
Japan Society for the Promotion of Science

る仮想マルチプロセッサマシンを構築することができる。クラスタを利用する際には、この仮想マシン上に Linux といったマルチプロセッサ対応の OS をインストールし、さらに、その OS 上で並列アプリケーションを記述・実行する。

SSI を提供するための既存研究と比較して、我々のアプローチは以下の3点で優れている。まず、1つ目の利点として、共有メモリ型マルチプロセッサマシンを仮定して記述された様々な並列アプリケーション（科学技術計算から Web サーバまで）を、無変更のまま分散環境上で実行できることがあげられる。特に、パラメータスイーブ型の並列アプリケーションなどを記述する際に、MPI³⁾ などの分散プログラミング言語を使用する代わりに、利用者が慣れ親しんだ既存の言語やツール（たとえば並列 make やシェルスクリプト）をそのまま使用することができる。

さらに2つ目の利点として、並列アプリケーションだけでなく、単に逐次プログラムを複数同時に走らせる場合にも、我々のアプローチが有用であることがあげられる。マルチプロセッサ用の OS カーネルも、少量の変更を加えるだけで仮想マシン上で動作させることができる。そのため、分散環境を管理するのに、たとえば Linux のプロセス空間やファイルシステムをそのまま用いることができる。仮想マシン上で走っている Linux から複数のプロセスを立ち上げると、それらのプロセスは、Linux のスケジューラによって、それぞれ異なる仮想プロセッサに対して割り当てられる。そして、最終的には、その仮想プロセッサが対応づけられた実プロセッサ上で実行されることになる。

最後に3つ目の利点として、VMM による資源のカプセル化 (resource encapsulation) を、安全性や信頼性の向上のために利用できることがあげられる。たとえば、クラスタをサーバーホスティングのために用いる際に、マシンを破壊する恐れのあるソフトウェアを安全に実行するためのサンドボックスとして、VMM を利用することができる^{4),5)}。また、仮想マシンの実行状態のスナップショットを取得・復元できるようにして、OS に何か不具合が生じたらそれ以前の状態に復旧する、といったことも可能になる^{2),6),7)}。

我々は、8 台の物理マシン上に仮想的に 8-way のマルチプロセッサマシンを構築した。その上でフィボナッチ数を計算するプロセスを 8 個並列に走らせたところ、仮想・物理シングルプロセッサマシン上で実行させた場合と比較して、最大約 6.6 倍の性能向上となった。この実験の結果から、VMM の頻繁な介入を必要としない（たとえば、I/O デバイスに頻繁にアク

セスしない）プログラムの場合、良い性能を達成できることが分かった。

本稿は以下のように構成される。2 章で本システムの概要について述べる。3 章でハードウェアの仮想化手法について説明する。4 章で共有メモリの一貫性制御について詳細に説明する。5 章で我々が行った予備実験について述べる。6 章で本システムが現在持つ限界とその解決策について議論する。7 章で関連研究との比較を行い、最後に 8 章でまとめと今後の課題について述べる。

2. Virtual Multiprocessor の概要

本章では、まず、Virtual Multiprocessor が提供する仮想マシンの特徴について説明する。次に、仮想マシンと実マシンの資源の対応づけについて述べる。

2.1 構築される仮想マシンの特徴

Virtual Multiprocessor によって構築される仮想マシンの特徴として、まず、命令セットアーキテクチャ (Instruction Set Architecture: ISA) のレベルで仮想化を行うことがあげられる。仮想マシン上では、その ISA を満たす OS を走らせることができる。現在の実装では IA-32 アーキテクチャを対象としている。

仮想マシンの他の特徴として、準仮想化 (para-virtualization) であることがあげられる。つまり、仮想マシンの ISA は、実ハードウェアの ISA とほぼ同一であるが、完全に同じではない。詳しくは後述するが、実マシンと仮想マシンの ISA の相違点は、コントロールレジスタや EFLAGS レジスタなどの特殊なハードウェアへのアクセスに関連したものである。よって、汎用レジスタのみアクセスする多くのユーザアプリケーションは、ソースコードの変更を必要としないで済む。特殊なハードウェアにアクセスする OS のカーネルやデバイスドライバは、ソースコード変更を必要とする。ただし、その変更も、独自のアセンブリ変換器を用いることで半自動的に行い、手動変更による手間は少ない¹¹⁾。詳しくは 3.2 節で述べる。

以上のように仮想マシンを設計した理由は、主に、少ない実装コストで、十分な性能を持つ仮想マシンを実現するためである。ちなみに、現在の実装では Linux をサポートしている。

2.2 仮想マシンと実マシンの資源の対応

Virtual Multiprocessor は、仮想マシンと実マシンの資源を以下のように対応づける（図 1 参照）。

User-mode Linux⁸⁾, Zap⁹⁾, SoftwarePot¹⁰⁾ などのような、Application Binary Interface (ABI) のレベルでの仮想化ではない。

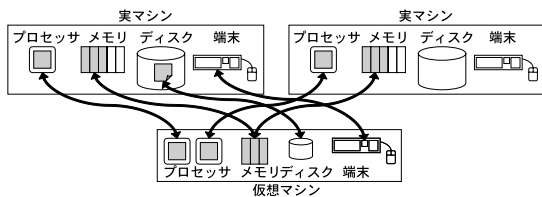


図 1 実マシンと仮想マシンの対応

Fig. 1 Mapping between physical machines and a virtual machine.

プロセッサ 基本的には、仮想マシンのプロセッサと実マシンのプロセッサは 1 対 1 に対応する．たとえば、 N 個のプロセッサからなる仮想マシンを構築するためには、複数の実マシンから総計 N 個のプロセッサを確保する必要がある．

メモリ 実マシンのメモリの一部を、仮想マシンのそれとして使用する．より具体的には、 M MB の共有メモリを仮想化するためには、各仮想プロセッサごとに M MB のメモリを実マシンから確保する必要がある．

I/O デバイス どれか 1 つの実マシンにあるデバイスを、仮想マシンのそれとして使用する．たとえばハードディスクの場合、実マシンのどれかに置かれたファイルをディスクイメージとして利用する．シリアル端末の場合も同様に、どれかの実マシンの仮想端末を、仮想マシンのシリアル端末として利用する．

3. ハードウェアの仮想化

本章では、プロセッサ・メモリ・I/O デバイスといった各種ハードウェアの仮想化処理の実装について述べる．メモリの仮想化は、アドレス空間と一貫性制御のそれぞれの仮想化を必要とする．プロセッサおよびメモリのアドレス空間の仮想化手法は、LilyVM¹¹⁾ とほぼ同様である．それ以外のメモリの一貫性制御と I/O デバイスの仮想化手法は、本システムに特有な実装となっている．

3.1 基本的な実装方針

Virtual Multiprocessor は、LilyVM¹¹⁾ や FAUmachine¹²⁾ と同様に、実マシン上で走る OS の上に位置する．この方式は、VMM がハードウェア上に直に置かれる方式^{4),5),13)} と比較してオーバーヘッドが大きい反面、少ない実装コストで IA-32 アーキテクチャの仮想化などを実現することができるといった利点を持つ¹⁴⁾．

VMM は、大半の命令を実プロセッサ上で直に実行し、実ハードウェアの状態と干渉する命令(たとえば、

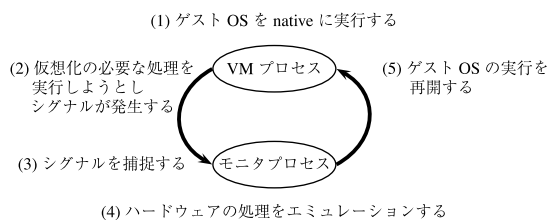


図 2 Virtual Multiprocessor の基本実行サイクル

Fig. 2 Basic execution cycle of Virtual Multiprocessor.

cr3 などのシステムレジスタを読み書きする命令)のみ、ソフトウェアでエミュレーションする(このエミュレーションを必要とする命令を、センシティブ命令と呼ぶ)．以上の動作を実現するために、1 つの仮想プロセッサごとに、VM プロセスとモニタプロセスという 2 つのユーザプロセスを用意する．

VM プロセス 仮想マシンのプロセッサのうちどれか 1 つを担当し、ゲスト OS を native に実行する．VM プロセスが仮想化を必要とする処理を行おうとすると、シグナルが発生する．たとえば、ユーザレベルでは実行することのできない特権命令を実行しようとする、SIGSEGV シグナルが発生する．

モニタプロセス VM プロセスの実行を監視し、必要に応じてセンシティブ命令の実行の仮想化など、ハードウェアの仮想化処理を行う．たとえば、VM プロセスがシグナルを発生させた場合、それを ptrace システムコールにより捕捉する．そして、捕捉したシグナルの種類と、その時点での VM プロセスのメモリやレジスタの状態を基に、必要とする仮想化処理を特定し、エミュレーションする．

この VM プロセスとモニタプロセスの実行サイクルを図 2 に示す．ちなみに、ゲスト OS 上で新たにプロセスが生成された場合、そのプロセスは VM プロセスの内部で仮想的に実行される．新たにホスト OS 上でプロセスが生成されるわけではない．

3.2 プロセッサの仮想化

プロセッサの仮想化は、大きく分けて、センシティブ命令の仮想化と、割込みや例外処理の仮想化からなる．まず、センシティブ命令の仮想化について述べる．センシティブ命令は、特権命令と非特権命令の 2 種類に分類される¹⁵⁾．特権命令(たとえば lgdt 命令)は、プロセッサの特権レベルが 0 より大きい場合、実行時に例外を発生させる．それに対して非特権命令(たとえば sgdt 命令)は、低い特権レベルで実行しても例外は発生しない．

センシティブ命令の実行を捕捉する方法は、その命

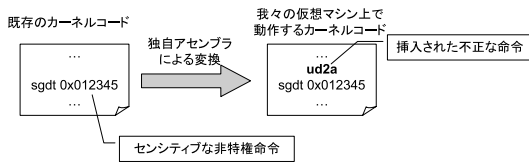


図 3 カーネルコードの変換

Fig. 3 Translation of kernel code.

令が特権命令か非特権命令であるかによって異なる．特権命令を捕捉するためには，実行時に発生した例外を捕捉すればよい（VM プロセスはユーザプロセスであり，特権レベル 3 で動作していることに注意）．それに対して，非特権命令を捕捉するためには，より複雑な処理を必要とする．具体的には，ゲスト OS のカーネル・デバイスドライバを静的に書き換え，すべてのセンシティブ命令の直前に不正な命令を挿入する（図 3 参照）．この不正な命令の挿入は，カーネルコンパイル時に，独自のアセンブリ変換器を用いて自動的に行われる．そして，モニタプロセスは，この不正な命令の実行によって発生した例外を捕捉することによって，非特権命令の実行直前に VM プロセスの実行を止める．

この静的なカーネルコードの変換は，小さな実装コストで非特権命令の仮想化が実現できるという利点を持つが，変換にソースコードを必要とするという欠点も持つ．たとえば，バイナリしかない Linux のカーネルモジュールや，Windows などのソースコードの公開されていない OS は，現在の実装では扱うことができない．

次に，割り込みや例外処理の仮想化について述べる．割り込みや例外処理の仮想化は，基本的には，以下の手順で実装される．まず，モニタプロセスが割り込みや例外の発生を検知する．これは，ptrace システムコールによってシグナルやシステムコール呼び出しを捕捉したり，APIC のマップされたメモリ領域への書き込みを捕捉したりすることによって実現する．次に，モニタプロセスは，その発生を検知した割り込み・例外を，適切な仮想プロセッサへと配送する．基本的には自分の担当する仮想プロセッサに対して割り込み・例外を配送するが，プロセッサ間割り込みが発生した場合は，TCP/IP で通信を行いながら指定された遠隔プロセッサに配送する．そして，割り込み・例外の配送された先の仮想プロセッサを担当するモニタプロセスは，割り込みディスクリプタテーブルなどを参照しながら，適切なハンドラへと VM プロセスの制御を移す．I/O デバイスの起こす外部割り込みについては，3.4 節で詳細に述べる．

3.3 共有メモリの仮想化

共有メモリの仮想化は大きく分けて，アドレス空間の仮想化と一貫性制御の実現の 2 つを必要とする．前述のとおり，アドレス空間の仮想化に関しては，LilyVM とほぼ同様の手法である．一貫性制御に関しては，本システムに特有な処理である．

まず，アドレス空間の仮想化について簡単に説明する．アドレス空間の仮想化を実現するためには，セグメント機構とページング機構のそれぞれの仮想化を実現する必要がある．セグメント機構の仮想化は，Linux を動作させるのに必要な処理のみ実装されている．具体的には，セグメントレジスタへの読み書きの仮想化は実装されているが，アドレス変換は未実装で，各セグメントの先頭アドレスはつねに 0 だと仮定している．

ページング機構の仮想化は以下のようにして実現している．まず，実マシン上に，仮想マシンのメモリイメージを保持するファイルを用意する．そして，仮想マシンのページディレクトリ・テーブルを参照することによって，VM プロセスの各ページをメモリイメージファイルへとマップする．より詳細には，VM プロセスがあるページにアクセスし SIGSEGV シグナルが発生した時点で，そのページを mmap システムコールによってイメージファイルへとマップする．また，ページディレクトリ・テーブルへの変更が有効になったとき（たとえば，cr3 レジスタの値が更新されたとき）に，ページディレクトリ・テーブル中にすでに存在しなくなったマッピングを munmap システムコールによって解放する．

また，仮想マシンがゲスト OS に提供する仮想アドレス空間の範囲を，0x00000000 以上 0xb0000000 未満に限定している（図 4 参照）．その理由は，0xc0000000 以上 0xffffffff 以下の領域はホスト OS のカーネルアドレス空間であり，0xb0000000 以上 0xc0000000 未満の領域はハードウェアの仮想化に必要な情報（たとえばシステムレジスタの値）を格納するために使用されるからである．ゲスト OS のカーネルコードを静的に変更し，カーネルアドレス空間の先頭アドレスを 0xa0000000 にしている．

次に，一貫性制御の仮想化について説明する．既存のプログラムを変更することなく仮想マシン上で走らせることを可能にするために，基本的にはソフトウェア分散共有メモリシステムなどと同様の手法をとる．具体的には，VM プロセスのページの読み書き権限を mprotect システムコールによって適宜制限することによって，一貫性制御を実装する．たとえば，あるページの情報が古くなった場合は，そのページへの読

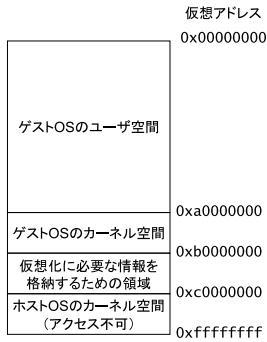


図 4 モニタプロセスのメモリアウト

Fig. 4 Memory layout of the monitor process.

み書きを禁止する。すると、そのページに VM プロセスがアクセスした際に SIGSEGV シグナルが発生する。モニタプロセスはこのシグナルを捕捉すると、他のモニタプロセスと TCP/IP で通信をしながら、ページの最新の状態を取得する。そして、そのページへの読み書きを許可し、VM プロセスの実行を再開させる。この共有メモリの一貫性制御の詳細については 4 章で述べる。

最後に、これまで述べてきたアドレス空間の仮想化と一貫性制御の仮想化でともに必要となる処理について述べる。まず、mmap, munmap, mprotect システムコールを呼び出す処理について説明する。これらのシステムコールは、それを呼び出したプロセスのページのマッピングやアクセス権限しか変更できない仕様になっている。そのため、モニタプロセスからでなく VM プロセスからシステムコールを呼び出すように、以下のように実装する必要がある。まず、VM プロセスのメモリ中に、システムコールを呼び出すためのコードを用意しておく。具体的には、ゲスト OS がアクセスしない領域 (0xb0000000 以上 0xc0000000 未満) の領域にコードを張り付けておく。そして、ページのマッピングやアクセス権限の変更が必要な際には、モニタプロセスは制御をいったん VM プロセスに移し、VM プロセスがその張り付けたコードを実行するようにする。

次に、SIGSEGV シグナルの扱いについて述べる。以上の記述から分かるように、様々な理由によって SIGSEGV シグナルが発生することに注意する必要がある。たとえば、モニタプロセスは SIGSEGV シグナルを捕捉した際に、それがページフォルトの発生によるものなのか、それとも共有メモリの一貫性制御によるものなのかを判別する必要がある。具体的には、図 5 に示されたフローチャートに従って、SIGSEGV シグナルの発生原因を特定し、それに応じた仮想化処理を行う。

3.4 I/O デバイスの仮想化

現在の実装は、ハードディスクやシリアル端末などの I/O デバイスのエミュレーションをサポートしている。また、I/O デバイスへアクセスする手段としては、主に、プログラム I/O 方式 (in/out 命令の発行)、Direct Memory Access (DMA) 方式、メモリマップド I/O 方式があげられるが、現在の実装では、基本的には、プログラム I/O 方式と DMA 方式のみをサポートしている。

具体的には I/O デバイスのエミュレーションのために、全デバイスの状態をソフトウェア上で管理する I/O サーバが 1 つ用意されている。モニタプロセスは、必要に応じてこのサーバと通信を行う。たとえば、ゲスト OS が in 命令によって p 番のポートから読み込みを行おうとする場合を考える。この場合、モニタプロセスは in 命令の実行を捕捉し、サーバに問合せをする。この問合せを受けたサーバは、 p 番のポートへ対応づけられたデバイスの動作をエミュレーションし、デバイスから読み出される値をモニタプロセスに返信する。

また、I/O サーバは各種デバイスの外部割り込みを以下のようにエミュレーションする。まず、デバイスの状態を定期的に (少なくともタイマ割り込みの間隔ごとに) 監視し、割り込みを発生させる必要のあるデバイスが存在するか調べる。そして、もし割り込みを発生させる必要がある場合は、基本的には、以下の手順で割り込みを仮想プロセッサに配送する。

- (1) I/O サーバが、割り込み配送先の仮想プロセッサを決定する。現在の実装では、外部割り込みの配送先は基本的には固定されており、I/O サーバと同一の実マシン上にある仮想プロセッサに対して送られる。
- (2) I/O サーバは、その仮想プロセッサを担当する VM プロセスに対して、適当なシグナル (たとえば SIGUSR1) を送る。
- (3) モニタプロセスが送信されたシグナルを捕捉する。
- (4) モニタプロセスは、割り込みディスクリプタテーブルなどを参照しながら、適切なハンドラへと VM プロセスの制御を移す。

4. 共有メモリの一貫性制御

Virtual Multiprocessor は IA-32 アーキテクチャの

優先度の低いプロセッサに割り込みを配送する、I/O APIC による割り込みの動的な分配は、まだ実装されていない。

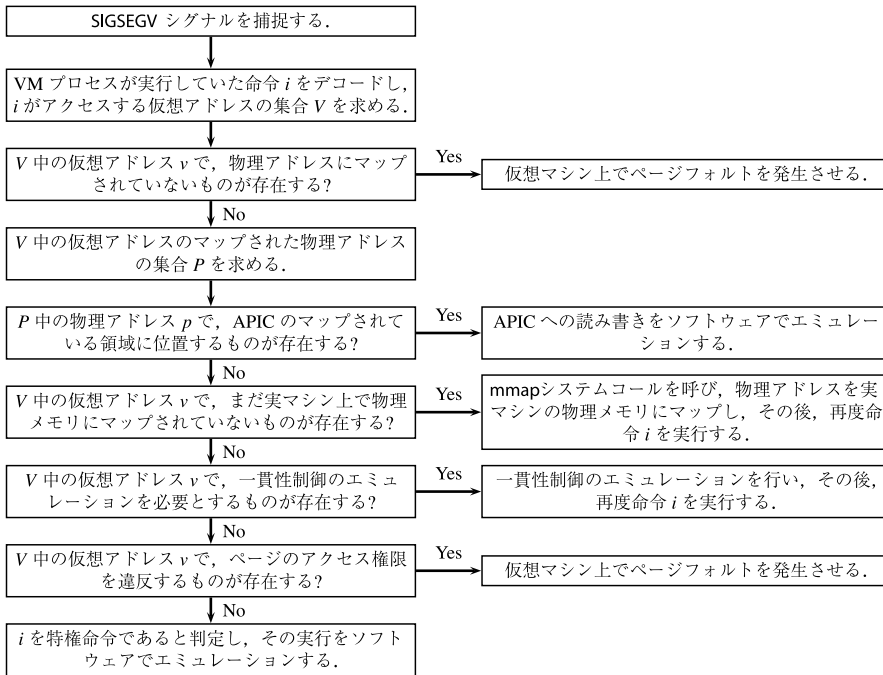


図 5 SIGSEGV シグナルの扱い
Fig. 5 Handling of SIGSEGV signals.

仮想化を目指しており、共有メモリの一貫性制御も、その IA-32 のメモリモデルを満たしている必要がある。多くの既存のソフトウェア分散共有メモリシステムに関する研究は release consistency^{16),17)} など他のメモリモデルを対象としていることから、この IA-32 のメモリモデル上での一貫性制御は我々のシステムに特有のものであるといえる。

本章では、まず、IA-32 のメモリモデルについて概説する。次に、我々が実装した simple なメモリ一貫性アルゴリズムについて述べる。

4.1 IA-32 のメモリモデルの概要

IA-32 のメモリモデルは、「あるプロセッサが行ったメモリアクセスが、自分および他のプロセッサにどういった順序で反映されるか」を定める、より具体的には、IA-32 のメモリモデルはプロセッサ順メモリモデル (Processor-ordered Memory Model) であり、以下の規則を満たす¹⁸⁾。

- あるプロセッサ x が行った読み書きについては、プログラムの文面どおりの順で x に反映される。
- あるプロセッサ x が行った書き込みが x に反映される順序と、他のプロセッサに対して反映される順序は同一である。
- 異なるプロセッサが行った書き込みの間では、それらが反映される順序に制約はない。

$proc_i$: 仮想プロセッサ i
$pages_i^n$: 仮想プロセッサ i の n 番ページ
$p.state$: ページ p の状態 (invalid, read_only, or read_write)
$p.c$: ページ p の内容
$p.own$: ページ p のオーナー
$p.copysset$: ページ p の複製を持つプロセッサの集合
$p.busy$: ページ p を処理中かどうかを示すフラグ (true or false)

図 6 アルゴリズムの記述のための変数
Fig. 6 Variables for algorithm description.

また、IA-32 は、上記のメモリモデルを強めるための命令を提供している。たとえば、mfence などの直列化命令 (serializing instruction) は、「直列化命令 i を実行した時点で、プログラムの文面上で i より前に位置する全メモリアクセス命令が、すでに全プロセッサに反映済みである」ことを保障する。

4.2 Simple なアルゴリズム

我々は、simple な一貫性制御アルゴリズムを実装した。このアルゴリズムは、ページ単位でメモリの共有・非共有が管理される write invalidate プロトコル¹⁹⁾ の一種である。同一ページに対して、読み込みは複数のプロセッサが同時に行うことができるが、書き込みは同時に 1 つのプロセッサしかできない。

図 6 と図 7 は、上記のアルゴリズムのより詳細な記

```

access nth page with a when violation(pagesin, a)
⇒ stop the execution of the VM process;
   send (fetch, n, a, i) to manager(n);

receive(fetch, n, a, s) when pagesin.busy = false
⇒ let p be pagesin;
   p.busy := true;
   match a with
   read ⇒ send (inv, n, a, s, p.own) to p.own;
   write ⇒ forall x ∈ S(p, s) do
       send (inv, n, a, s, p.own) to x;

receive (inv, n, a, s, o)
⇒ let p be pagesin;
   match a with
   read ⇒ p.state := read_only;
   write ⇒ p.state := invalid;
   if o = proci then send (ack, n, a, p.c) to procs;

receive (ack, n, a, c)
⇒ let p be pagesin;
   p.c := c;
   match a with
   read ⇒ p.state := read_only;
   write ⇒ p.state := read_write;
   send (finish, n, a, i) to manager(n);
   restart the execution of the VM process;

receive(finish, n, a, s)
⇒ let p be pagesin;
   match a with
   read ⇒ p.copyset := p.copyset ∪ { procs };
   write ⇒ p.copyset := { procs };
           p.own := procs;
   p.busy := false;

ただし
manager(n) : n ページのマネージャー
             (たとえば, プロセッサ数を N とするとき
              manager(n) = procn mod N)
violation(p, a) ≡ p.state = invalid ∨
                (p.state = read_only ∧ a = write)
S(p, s) = {x ∈ p.copyset : x ≠ procs ∨ x = p.own}

```

図 7 Simple なメモリ一貫性アルゴリズム (仮想プロセッサ i における)

Fig. 7 Simple memory consistency algorithm (for virtual processor i).

述である。図 6 には、アルゴリズム中に現れる変数が示されている。図 7 には、仮想プロセッサ i の動作が記述されている。各動作は $G \Rightarrow A$ というシンタックスで記述され、 G が満たされると A がアトミックに実行される。また、メッセージの送信先と送信元が同一プロセッサの場合、実際の実装ではメッセージは送信されず、ローカルに処理される。

5. 予備実験

我々は、Virtual Multiprocessor のプロトタイプを実装し、予備実験を行った。このプロトタイプでは、

最大 8 台の物理マシン上に 8-way の仮想マルチプロセッサマシンを構築し、その仮想マシン上で Linux を走らせることができる。また、我々は、仮想マシン上で gcc や make などの既存のプログラムが動作することも確認した。

予備実験として、システムの基本性能を評価するため、まず、共有メモリの一貫性制御以外の仮想化処理（たとえば、センシティブ命令の実行や I/O デバイスへの読み書きのエミュレーション）にかかるオーバヘッドを測定するために、仮想シングルプロセッサマシン上でいくつかのベンチマークプログラムを実行した。そして次に、共有メモリの一貫性制御のオーバヘッドを測定するために、仮想マルチプロセッサマシン上で互いに独立なタスクを並列に実行した。

すべての実験において、Intel Xeon 2.4 GHz, 2 GB RAM, 1 Gigabit Ethernet NIC からなるマシンを用いた。また、ホスト OS とゲスト OS とともに Linux 2.4 を用いた。

5.1 仮想シングルプロセッサマシンの性能評価

共有メモリの一貫性制御以外の仮想化処理にかかるオーバヘッドを測定するために、物理および仮想シングルプロセッサマシン上で、いくつかの逐次プログラムを実行した。

表 1 は、実験に使用したベンチマークプログラムの説明と、物理・仮想シングルプロセッサマシンにおけるそのプログラムの実行時間を示している。この実験結果から、getpid や gcc の実行において、システムコール呼び出しや I/O デバイスの仮想化にかかるオーバヘッドが非常に大きいことが分かる。我々は、既存の VMM（たとえば Xen¹³⁾ や CoVirt²⁰⁾）によって開発された最適化手法を適応することで、このオーバヘッドは削減可能であると考えられる。

5.2 仮想マルチプロセッサマシンの性能評価

1 ~ 8 台の物理マシン上にそれぞれ 1-way ~ 8-way の仮想マルチプロセッサマシンを構築し、その上でフィボナッチ数を計算するプロセスを 8 個並列に走らせた。このプログラムを実行した結果、図 8 に示される性能向上を得た。fib(n) は、 n 番目のフィボナッチ数を計算する際の性能向上を表している。この図から、実行されるタスクの粒度が大きければ台数効果がでていくことが分かる。たとえば、fib(46) を動作させたときは、8 プロセッサで約 6.6 倍の性能向上となった。オーバヘッドの要因を調べるため、fib(44) を実行

取得可能な BIOS の制限などのため、現在の我々の実装では 8 プロセッサまでしか扱っていない。

表 1 逐次ベンチマークプログラムの説明と、そのプログラムの物理・仮想シングルプロセッサマシン上での実行時間 (単位: 秒)

Table 1 Description of Sequential benchmark programs and their execution time on a physical and a virtual single-processor machine (unit: seconds).

プログラム名	プログラムの説明	物理マシン上での 実行時間 (P)	仮想マシン上での 実行時間 (V)	オーバーヘッド (V/P)
fib	フィボナッチ数を計算する	22.6	22.1	0.97
getpid	getpid を 100,000 回発行する	0.05	18.1	354
ls	数百のファイルの情報を表示する	0.03	6.64	255
gcc	C プログラムをコンパイルする	0.14	0.98	6.81

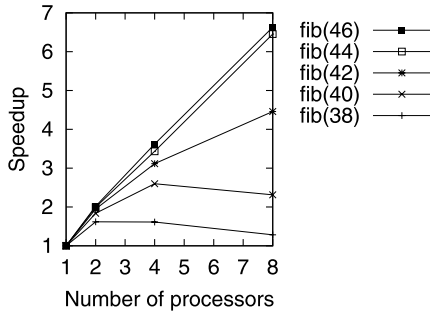


図 8 並列フィボナッチの性能向上
Fig. 8 Speedup of parallel fibonacci.

表 2 fib(44) の実行時間の内訳 (単位: 秒)

Table 2 Breakdown of execution time of fib(44) (unit: seconds).

プロセッサ数	Total	Native	Mem	Misc	Idle
1	180.0	177.8	0.0	2.2	0.0
2	90.3	87.9	1.0	1.1	0.3
4	52.4	43.7	3.0	0.4	5.3
8	27.9	22.1	3.7	0.1	2.0

Total : 全実行時間
 Native : ゲスト OS が native に実行されていた時間
 Mem : 一貫性制御のエミュレーションにかかる時間
 Misc : センシティブ命令などのエミュレーションにかかる時間
 Idle : 仮想マシンが halt 命令を実行していた時間

した際の実行時間の内訳を計測した。表 2 はその結果を示しており、それぞれの値は、各プロセッサごとにかかった時間の平均値を示している。Total は、8つのプロセスすべてが実行を終えるまでにかかった全実行時間を表す。Native は、ゲスト OS が実機上で走っていた時間である。フィボナッチ数を計算するのにかかった時間を主に表す。Mem は、モニタプロセスが行う仮想化処理のうち、共有メモリの一貫性制御の仮想化にかかった時間である。遠隔マシンとの通信にかかった時間を主に表す。Misc は、モニタプロセスが行う仮想化処理のうち、共有メモリの一貫性制御以外の処理 (たとえば、センシティブ命令の実行や I/O デバイスへの読み書きのエミュレーション) にかかった時間を示す。Idle は、実行可能なプロセスが存在せず仮想プロセッサが halt 命令を実行していた時間を表す。表 2 から、プロセッサ数が増加するにつれて、共有メモリの一貫性制御の仮想化のオーバーヘッドが増大していることが分かる。

そこで、共有メモリの一貫性制御のオーバーヘッドについて、より詳しい分析を行った。具体的には、fib(44) 実行時における、ページのフェッチが発生した仮想アドレスの分布 (図 9 参照) と、1 回のページのフェッチにかかる時間の分布 (図 10 参照) とを測定した。図 9 では、横軸がプログラム開始からの経過時間を、縦軸がページのフェッチが発生した仮想アドレスを表

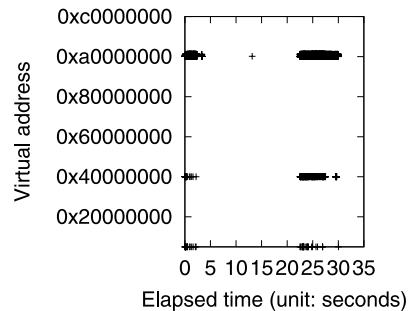


図 9 メモリ共有のためにモニタプロセスによってフェッチされた仮想アドレスの分布 (fib(44) における)
 Fig. 9 Distribution of virtual addresses fetched by the monitor processes for memory sharing (for fib(44)).

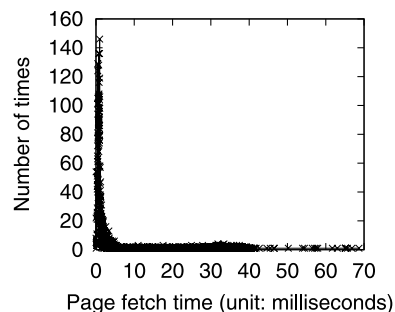


図 10 個々のページフェッチが完了するのにかかった時間の分布 (fib(44) における)
 Fig. 10 Distribution of times which individual page fetch requests took to complete (for fib(44)).

している．図 10 では，横軸が 1 回のフェッチにかかる処理時間を表しており，縦軸がフェッチの発生数を表す．両図ともに，全プロセッサの総計を求めている．図 9 から，プログラムの実行直後と終了直前に，ユーザ空間とカーネル空間の両方で，ページのフェッチが大量に発生していることが分かる（ゲスト OS のカーネル空間が 0xa0000000 から始まっていることに注意）．また，図 10 から，大半のページのフェッチの処理は数ミリ秒で完了するが，一部のフェッチは処理するのに数十ミリ秒必要としていることが分かる．以上のことから，false sharing などのために同一ページへのフェッチが集中した際に，オーバーヘッドが大きくなっていることが考えられる．

6. 議 論

この章では，本システムが現在持つ限界とその解決策について議論する．

6.1 共有メモリの一貫性制御アルゴリズム

4.2 節で述べた simple なアルゴリズムは，十分効率の良いものとはいえない．そこで，我々は，IA-32 のメモリモデルにより適した，効率の良い一貫性制御アルゴリズムの案を示す．

具体的には，1 つのページに対して複数のプロセッサが同時に書き込むことを可能にすることで，最適化を図る．このために，IA-32 のメモリ順序モデルが，「同期命令を実行するまで，自プロセッサの行った書き込みが遠隔プロセッサに反映されるとは限らない」という性質を持つことを，以下のように利用する．

- 同期命令が実行されるまでは，自プロセッサが行った書き込みを他のプロセッサに反映させない．
- 同期命令が実行された際に，直前の同期命令から今までの間に行われた書き込みを他のプロセッサに反映させる．文献 16) にあるように，メモリ状態の差分を計算することで，これを実現する．

上記のような最適化を通してアルゴリズムの性能を改善することが可能であると我々は考えるが，それには，いくつか解決しなければならない問題が残っている．その 1 つは，同期命令などの実行を捕捉するのにかかるオーバーヘッドが大きくなってしまおうという問題である．そこで，我々は，同期命令の捕捉を実現するのに，現在の ptrace システムコールを用いる方式の代わりに，ホスト OS の一部を改変するという方式²⁰⁾をとることで，コンテキストスイッチの回数を減らしオーバーヘッドを削減する．

6.2 メモリの一貫性制御以外の仮想化処理のオーバーヘッド

現在では，ホスト OS のカーネルは改良せず，また，ゲスト OS に加える改良も最小限のものにとどめるという方針で，VMM は実装されている．この方式は，VMM の実装コストが少なく済むという利点があるものの，5.1 節の実験結果が示すように，システムコール呼び出しや I/O デバイスの仮想化にかかるオーバーヘッドが大きくなってしまおうという欠点がある．

これらのオーバーヘッドを削減する手段としては，まず，既存の VMM（たとえば Xen¹³⁾ や CoVirt²⁰⁾）によって開発された技術を我々の VMM に適応することが考えられる．たとえば，Xen などのように VMM をハードウェアの上に直に置く方式をとることで，ptrace を用いることによって発生するコンテキストスイッチの数を削減することができる．また，仮想化に適した設計のアーキテクチャ^{21),22)} 上での実装に変更することも考えられる．

6.3 不均等メモリアクセスの考慮

本システムにおいては，仮想マシンのメモリは複数の物理マシンにまたがって共有されるため，通常の NUMA アーキテクチャ同様メモリアクセスが不均等になる．今回実験で用いた Linux カーネル 2.4 は，NUMA アーキテクチャに対応していないため，たとえば，プロセスがスケジューリングによってプロセッサ間を頻繁に移送されると，遠隔メモリへのアクセスによってオーバーヘッドが発生してしまう．実際，表 2 の実験結果で共有メモリのエミュレーションのオーバーヘッドが大きくなっている理由の 1 つとして，カーネルが NUMA 非対応であることが考えられる．

そこで，こうしたオーバーヘッドを回避するために，NUMA アーキテクチャに対応した Linux カーネル 2.6 をゲスト OS として用いる．カーネル 2.6 においては，たとえば，メモリの局所性を考慮したプロセススケジューリングが提供されている．各プロセスは，原則的には特定のプロセッサに割り当てられて，プロセッサ間の移送は起こらない．そのため，プロセスの移送時に発生するメモリ転送によるオーバーヘッドを回避することができる．

ちなみに，まだ本システムの実装が十分ではなく，カーネル 2.6 が安定して動作をしなかったため，今回の実験はカーネル 2.4 を用いて行った．

6.4 物理マシン上に確保する必要のあるメモリ量

現方式では， M MB の共有メモリを仮想化するためには，各仮想プロセッサごとに M MB のメモリを実マシンから確保する必要がある．これに対して，仮

想マシンの物理メモリを各仮想プロセッサが分担して確保することで、必要なメモリ量を削減する手法が考えられる。たとえば、仮想プロセッサ数を N 、キャッシュサイズを C とするとき、各仮想プロセッサごとに確保するメモリ量が $M/N + C$ で収まるようにする。

この分割方式と比較して、現方式は以下のような利点・欠点を持つ。現方式の利点は、以下のとおりである。分割方式は、確保しているページがキャッシュサイズより大きくなった場合に、一部のページを他の仮想プロセッサに移譲する仕組みが必要になる。そのため、共有メモリの一貫性制御が複雑になり、より大きなオーバーヘッドを必要とする。現方式では、そうした処理によるオーバーヘッドを回避することができる。

現方式の欠点は、実マシン上に確保する必要のあるメモリの総量が、仮想プロセッサ数に比例して増加するという点である。

6.5 動的な資源の変化への対応

現在の実装では、仮想マシンの資源と実マシンの資源の対応づけは静的に決定され、利用する物理マシンを動的に変更することはできない。そのため、複数の利用者が存在し、利用可能な物理マシンが動的に変化する環境では、本システムを用いることができない。

そこで今後は、より柔軟に仮想マシンの資源と実マシンの資源の対応づけを行うことで、利用可能な物理マシンが動的に変化する環境下でも、本システムを動作可能にする。具体的には、以下の方式によって、実マシンの構成によらず、つねに一定数の仮想プロセッサを提供可能にする。まず、各実プロセッサに1つまたは複数の仮想プロセッサを割り当てることで、実マシンの増加に耐えられる十分な仮想プロセッサをあらかじめ用意する。そして、マシンの増減に応じて、仮想プロセッサの実プロセッサへの割当てを動的に変更させる。たとえば、実マシンが削除される場合には、そのマシンに割り当てられた仮想プロセッサを他の実マシンに移譲する。また、各物理プロセッサに割り当てられる仮想プロセッサ数が不均等になる場合には、time ballooning²³⁾などの技術を利用して負荷を調整する。

6.6 耐故障性の導入

現在の方式では、物理マシンのうちどれか1つでも故障してしまうと、システム全体が停止してしまう。そのため、クラスタなど頻繁に故障が発生する環境では、本システムを用いることができない。

そこで今後は、耐故障性を導入し、物理マシンが数台故障した中でもVMMが動作し続けることを可能にする。分散チェックポイント技術²⁴⁾や仮想マシンの複

製技術²⁵⁾を利用することで、これを実現する。

7. 関連研究

7.1 仮想マシンモニタ

vNUMA²⁶⁾は、我々のVMMと同様、分散環境上に共有メモリ型マルチプロセッサマシンを仮想的に構築するシステムである。このvNUMAとVirtual Multiprocessorの相違点として、まず、vNUMAがItaniumアーキテクチャを対象としているのに対して、Virtual MultiprocessorはIA-32アーキテクチャを対象としていることがあげられる。

また、別の相違点として、共有メモリの一貫性制御アルゴリズムがあげられる。vNUMAが元に行っているIvy¹⁹⁾では、ページのinvalidationが同期的に行われる(メッセージ送信後、アックを受信するまで待機する)のに対して、Virtual Multiprocessorでは非同期に行われる。それによって、1回のページのフェッチにおけるメッセージの送信回数と、クリティカルパス長(ゲストOSが実行を再開するまでに最低限送信の必要なメッセージ数)は、各システムで表3のようになる。ただし、フェッチのリクエスト元のプロセッサを除く、ページの複製を持つプロセッサの数を c とする。この表から、read時のメッセージ送信回数はvNUMAの方が優れているが、write時のメッセージ送信回数とクリティカル長については、Virtual MultiprocessorとvNUMAは同性能、または、Virtual Multiprocessorの方が優れていることが分かる。これらの差は定数倍であるので、実アプリケーションで有意な差が生じるかどうかを、今後実験を通して検証する。

また、vNUMAの論文中にはccNUMAを仮想的に

表3 IvyとVirtual Multiprocessorのメモリー一貫性アルゴリズムの比較

Table 3 Comparison of memory consistency algorithm between Ivy and Virtual Multiprocessor.

	アクセス	マネージャ	メッセージ数	クリティカルパス長
Ivy	read	リクエスト元 オーナー それ以外	2 2 3	2 2 3
	write	リクエスト元 オーナー それ以外	2 + 2c 2 + 2c 3 + 2c	4 4 5
Virtual Multi- processor	read	リクエスト元 オーナー それ以外	2 3 4	2 2 3
	write	リクエスト元 オーナー それ以外	2 + c 3 + c 4 + c	2 3 4

c: フェッチのリクエスト元のプロセッサを除く、ページの複製を持つプロセッサの数

構築すると述べられているが、共有メモリを仮想化するために実マシン上に確保するメモリ量については、明示的には触れられていない。そのため、メモリ使用量の点で、Virtual Multiprocessor と vNUMA とを比較できていない。

vNUMA のほかには、以下のような VMM が関連研究としてあげられる。Virtual Iron²⁷⁾ も、分散環境上に共有メモリ型マルチプロセッサマシンを仮想的に構築するシステムだが、詳細が未公開であるため十分な比較を行えていない(2005年12月現在)。Disco²⁸⁾ と VMware ESX Server⁵⁾ は、仮想マルチプロセッサマシンを構築する VMM であるが、これらのシステムは単一のマルチプロセッサマシン上でしか動作しない。そのため、分散環境上で SSI を提供するためにこれらのシステムを利用するのは難しい。

7.2 クラスタ用ミドルウェアシステムと OS

SSI を実現するクラスタ用ミドルウェアシステムについて、数多くの研究がなされてきた。例として、SCore¹⁾ や Condor²⁹⁾ などがあげられる。これらのミドルウェアシステムを用いると、遠隔マシンへの並列ジョブ投入などを行うことが可能になる。

しかし、これら既存のクラスタ用ミドルウェアシステムが提供する SSI は、Virtual Multiprocessor の提供する SSI と比較して、機能が制限されたものとなっている。たとえば、共有メモリを仮定して記述された既存の並列プログラムをそのまま分散環境上で実行することはできない。また、CPU などの資源を大域的に管理する機構が存在するが、それぞれのシステムを持つ特有の操作・処理に習熟する必要がある。

また、Linux カーネルに一部変更を加え、分散環境上で動作させるという研究(たとえば MOSIX³⁰⁾ や Kerrighed³¹⁾) も存在する。これらの研究で提案されたシステムには、カーネルの改変に多大な手間を必要とするという問題がある。

7.3 ソフトウェア分散共有メモリシステム

共有メモリを仮定して記述された既存の並列プログラムを分散環境上で実行可能にするソフトウェア分散共有メモリシステムの例として、Shasta³²⁾ があげられる。Shasta は、コンパイラによって実行ファイルを変換し、ロード・ストアの直前にチェックコードを挿入することによって、共有メモリを実現する。

Shasta がユーザアプリケーションしか対象としていないのに対して、我々のシステムは、ユーザアプリケーションだけでなく OS カーネルも分散環境上で実行可能にする。また、対象としているメモリモデルも異なり、Shasta は MIPS アーキテクチャを対象とし

ているが、Virtual Multiprocessor は IA-32 アーキテクチャを対象としている。

8. おわりに

本稿では、ネットワークで結合された複数のマシン上に共有メモリ型マルチプロセッサマシンを仮想的に構築するシステム Virtual Multiprocessor について述べた。本システムによって、クラスタなどの分散環境を簡便に利用することが可能になる。予備実験によって、粗粒度で共有メモリへのアクセスが少ないアプリケーションであれば、良い性能が得ることを実証した。こうした性質を持つアプリケーションには Parameter sweep や並列 make など多くの有用な並列プログラムが含まれるため、我々のシステムが現実的に役に立つものとする。本システムのプロトタイプは、<http://www.yl.is.s.u-tokyo.ac.jp/~kaneda/vmp> から取得できる。

今後の課題は、6章で述べた問題の解決である。また、SPLASH-2 や Apache といった、より現実的なアプリケーションを仮想マシン上で走らせ性能評価を行う。それと同時に、よりプロセッサ数を増やした環境での実験も行う。

謝辞 本研究の一部は、科学技術振興機構戦略的創造事業の支援を受けた。

参考文献

- 1) SCore Cluster System Software.
<http://www.pcluster.org/>
- 2) VMware Inc. <http://www.vmware.com/>
- 3) Message Passing Interface (MPI) standard.
<http://www-unix.mcs.anl.gov/mpl/>
- 4) Whitaker, A., Shaw, M. and Gribble, S.D.: Scale and Performance in the Denali Isolation Kernel, *Proc. 5th Symposium on Operating Systems Design and Implementation (OSDI '02)*, pp.195-209 (2002).
- 5) Waldspurger, C.A.: Memory Resource Management in VMware ESX Server, *Proc. 5th Symposium on Operating Systems Design and Implementation (OSDI '02)*, pp.181-194 (2002).
- 6) Whitaker, A., Cox, R.S. and Gribble, S.D.: Configuration Debugging as Search: Finding the Needle in the Haystack, *Proc. 6th Symposium on Operating Systems Design and Implementation (OSDI '04)*, pp.77-90 (2004).
- 7) Sato, O., Potter, R., Yamamoto, M. and Hagiya, M.: UML Scrapbook and Realization of Snapshot Programming Environment, *Proc.*

- 2nd International Symposium on Software Security (ISSS)*, pp.281–295 (2003).
- 8) User-mode Linux.
<http://user-mode-linux.sourceforge.net/>
 - 9) Osman, S., Subhraveti, D., Su, G. and Nieh, J.: The Design and Implementation of Zap: A System for Migrating Computing Environments, *Proc. 5th Symposium on Operating Systems Design and Implementation (OSDI '02)*, pp.361–376 (2002).
 - 10) Kato, K. and Oyama, Y.: SoftwarePot: An Encapsulated Transferable File System for Secure Software Circulation, *Proc. 1st International Symposium on Software Security (ISSS)*, pp.112–132 (2002).
 - 11) Eiraku, H. and Shinjo, Y.: Running BSD Kernels as User Processes by Partial Emulation and Rewriting of Machine Instructions, *Proc. BSDCon 2003*, pp.91–102 (2003).
 - 12) Höxer, H.-J., Buchacker, K. and Sieh, V.: Implementing a User-Mode Linux with Minimal Changes from Original Kernel, *Proc. Linux-Kongress 2002*, pp.72–82 (2002).
 - 13) Barham, P., Dragovic, B., Fraser, K., Hand, S., Harris, T., Ho, A., Neugebauer, R., Pratt, I. and Warfield, A.: Xen and the Art of Virtualization, *Proc. 19th ACM Symposium on Operating Systems Principles (SOSP-19)*, pp.164–177 (2003).
 - 14) Sugeran, J., Venkitachalam, G. and Lim, B.-H.: Virtualizing I/O Devices on VMware Workstation's Hosted Virtual Machine Monitor, *Proc. USENIX Annual Technical Conference (USENIX '01)*, pp.1–14 (June 2001).
 - 15) Robin, J.S. and Irvine, C.E.: Analysis of the Intel Pentium's Ability to Support a Secure Virtual Machine Monitor, *Proc. 9th USENIX Security Symposium*, pp.129–144 (2000).
 - 16) Keleher, P., Cox, A.L., Dwarkadas, S. and Zwaenepoel, W.: TreadMarks: Distributed Shared Memory on Standard Workstations and Operating Systems, *Proc. USENIX Winter 1994 Technical Conference*, pp.115–131 (1994).
 - 17) Bennett, J.K., Carter, J.B. and Zwaenepoel, W.: Munin: Distributed Shared Memory Based on Type-specific Memory Coherence, *Proc. 2nd ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming (PPoPP '90)*, pp.168–176 (1990).
 - 18) Intel Corporation: *IA-32 Intel Architecture Software Developer's Manual Volume 3: System Programming Guide* (2003).
 - 19) Li, K. and Hudak, P.: Memory Coherence in Shared Virtual Memory Systems, *ACM Trans. Comput. Syst. (TOCS)*, Vol.7, No.4, pp.321–359 (1989).
 - 20) King, S.T., Dunlap, G.W. and Chen, P.M.: Operating System Support for Virtual Machines, *Proc. USENIX Annual Technical Conference (USENIX '03)*, pp.71–84 (2003).
 - 21) Intel Corporation: *Intel Virtualization Technology Specification for the IA-32 Intel Architecture* (2005).
 - 22) Advanced Micro Devices: *AMD64 Virtualization Codenamed "Pacifica" Technology Secure Virtual Machine Architecture Reference Manual* (2005).
 - 23) Uhlig, V., LeVasseur, J., Skoglund, E. and Dannowski, U.: Towards Scalable Multiprocessor Virtual Machines, *Proc. 3rd Virtual Machine Research and Technology Symposium (VM '04)*, pp.43–56 (2004).
 - 24) Elnozahy, E.N.M., Alvisi, L., Wang, Y.-M. and Johnson, D.B.: A Survey of Rollback-recovery Protocols in Message-passing Systems, *ACM Computing Surveys (CSUR)*, Vol.34, No.3, pp.375–408 (2002).
 - 25) Bressoud, T.C. and Schneider, F.B.: Hypervisor-based fault tolerance, *ACM Trans. Comput. Syst. (TOCS)*, Vol.14, No.1, pp.80–107 (1996).
 - 26) Chapman, M. and Heiser, G.: Implementing Transparent Shared Memory on Clusters Using Virtual Machines, *Proc. USENIX Annual Technical Conference (USENIX '05)*, pp.383–386 (2005).
 - 27) Virtual Iron Software.
<http://www.virtualiron.com/>
 - 28) Bugnion, E., Devine, S. and Rosenblum, M.: Disco: Running Commodity Operating Systems on Scalable Multiprocessors, *Proc. 16th ACM Symposium on Operating Systems Principles (SOSP-16)*, pp.143–156 (1997).
 - 29) Litzkow, M., Livny, M. and Mutka, M.: Condor — A Hunter of Idle Workstations, *Proc. 8th International Conference of Distributed Computing Systems (ICDCS '88)*, pp.104–111 (1988).
 - 30) Barak, A. and La'adan, O.: The MOSIX Multicomputer Operating System for High Performance Cluster Computing, *Journal of Future Generation Computer Systems (FGCS)*, Vol.13, No.4–5, pp.361–372 (1998).
 - 31) Morin, C., Lottiaux, R., Vallée, G., Gallard, P., Utard, G., Badrinath, R. and Rilling, L.: Kerrighed: A Single System Image Cluster Operating System for High Performance Computing, *Proc. 9th International European Confer-*

ence on Parallel Processing (Euro-Par '03), pp.1291-1294 (2003).

- 32) Scales, D.J., Gharachorloo, K. and Thekkath, C.A.: Shasta: A Low Overhead, Software-Only Approach for Supporting Fine-Grain Shared Memory, *Proc. 7th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS-VII)*, pp.174-184 (1996).

(平成 17 年 7 月 29 日受付)

(平成 17 年 11 月 14 日採録)



金田 憲二

1979 年生。2001 年東京大学理学部情報科学科卒業，2003 年東京大学大学院情報理工学系研究科コンピュータ科学専攻修士課程修了。現在，東京大学大学院情報理工学系研究科コンピュータ科学専攻博士課程に在学。クラスタ・グリッド計算と分散システムに興味を持つ。ACM 会員。



大山 恵弘

1973 年生。2001 年東京大学大学院理学系研究科情報科学専攻修了。博士（理学）。2001 年から 2003 年まで科学技術振興事業団研究員として筑波大学に勤務。2003 年より東京大学大学院情報理工学系研究科助手。情報処理学会平成 13 年度，16 年度論文賞受賞。日本ソフトウェア科学会 2003 年度論文賞受賞。興味はセキュリティ，システムソフトウェア，プログラミング言語，並列分散処理。



米澤 明憲

1947 年生。MIT 計算機科学科博士課程修了（Ph.D. 計算機科学）。MIT 計算機科学研究所および人工知能研究所において並列・分散計算モデルの研究に従事し，「並列オブジェクト」概念の構築に寄与した。帰国後，東京工業大学を経て，1988 年に東京大学大学院コンピュータ科学専攻教授となり現在に至る。日本ソフトウェア科学会理事長，ドイツ国立情報科学技術研究所（GMD）科学顧問を歴任。2004 年 3 月より 3 年間内閣府総合規制改革会議委員，同教育分野主査（独）産業技術総合研究所情報セキュリティ研究センター（東京秋葉原在）副センター長を兼務。日本ソフトウェア科学会フェロー，米国計算機学会終身フェロー（ACM Fellow）。