

## カーネル・レベル・コードによるユーザ・レベル VMM の移植性の向上

榮樂 英樹<sup>†</sup>      新城 靖<sup>†‡</sup>      加藤 和彦<sup>†‡</sup>

<sup>†</sup> 筑波大学大学院システム情報工学研究科コンピュータサイエンス専攻

<sup>‡</sup> 科学技術振興機構 CREST

### 要旨

LilyVM は、静的なコード変換技術と準仮想化技術を用いて、実機用の OS をユーザ・プロセスとして動作させる仮想マシン・モニタである。従来の LilyVM では、多くの特殊な機能をホスト OS に要求しており、他の OS への移植性が低く、性能にも問題があった。この論文では、LilyVM においてカーネル・レベル・コードを用いることにより、他の OS への移植性を高めることについて述べている。カーネル・レベル・コードは、高い特権レベルで動作し、CPU 例外のリダイレクトを OS に依存せずに実行する。

## Improving Portability of a User-Level VMM by Kernel-Level Code

Hideki Eiraku<sup>†</sup>      Yasushi Shinjo<sup>†‡</sup>      Kazuhiko Kato<sup>†‡</sup>

<sup>†</sup>Department of Computer Science, University of Tsukuba

<sup>†‡</sup>Japan Science and Technology Agency CREST

### Abstract

LilyVM is a virtual machine monitor (VMM) which runs an operating system (OS) for a real machine as a user process using a static code translation technique and a paravirtualization technique. Former LilyVM required many special features of host OSes. Therefore, it had less portability to other OSes, and also had a performance problem. In this paper, we discuss about improving portability of LilyVM to other OSes by using a kernel-level code. The kernel-level code runs in a higher privilege level and redirects CPU exceptions without dependence on a host OS.

## 1 はじめに

1 台のコンピュータで複数のオペレーティング・システム (Operating Systems, 以下 OS と略す) を同時に動作させることが、一般的に行われるようになってきている。その目的としては、Unix と Windows のように、その上で動作する異なるタイプの応用プログラムを組み合わせる利用することや、Mac OS 9 と Mac OS X のようにオペレーティング・システムの進化を円滑に進めることなどがあげられる。この論文では、1 つの基盤となる OS を動作させ、そ

の上で別の OS を動作させることについて述べる。基盤となる OS を **ホスト OS**、ホスト OS 上で実行される OS を **ゲスト OS** と呼ぶ。

ソフトウェアによるエミュレーションで仮想化された実行環境を **仮想マシン** と呼ぶ。ゲスト OS による仮想的なハードウェアとの入出力をモニタし、エミュレートするものを、**仮想マシン・モニタ** (Virtual Machine Monitor, 以下 **VMM** と略す) と呼ぶ。一般に VMM は、実際のハードウェアをエミュレートすることにより、実機用の任意の OS を実行する

ことができる。このような仮想化技術を、**完全仮想化**と呼ぶ。完全仮想化に基づく VMM の欠点としては、CPU をほぼ完全にエミュレートしなければならないので、実現が難しいこと、高速化のためにハードウェアによる支援や機械語命令の動的な書き換えが必要になってしまうこと、および、動的な書き換えに伴うオーバーヘッドが存在することがある。完全仮想化に対して、仮想化に都合の良いアーキテクチャを定義し、ゲスト OS をそのアーキテクチャに移植することによって、仮想マシン上で OS を実行できるようにしたものがある。これは、**準仮想化** (paravirtualization) と呼ばれる。準仮想化の問題点は、実機用のゲスト OS を移植しなければならないことである。

我々は、IA-32 を対象として準仮想化技術を用いた VMM である LilyVM を開発している。LilyVM の特徴は、機械語命令の静的な変換により、ゲスト OS の移植を不要にしている点にある [4]。LilyVM では、ゲスト OS として NetBSD, FreeBSD および Linux が動作している。また、LilyVM は、ホスト OS として NetBSD, FreeBSD および Linux の上で動作している。従来の LilyVM は、ユーザ・レベルのコードのみを用いて実現されており、以下のようなホスト OS が提供している機能を使用していた。

- 指定したアドレスにメモリをマップする機能。
- システム・コールを検出する機能。
- 例外の詳細情報 (割り込み番号や、エラー・コードなど) を取得する機能。

これらは一般的な機能ではないため、移植性が低いという問題があった。また、以下のようなオーバーヘッドが大きかった。

- フォールトおよびトラップのエミュレーション。
- MMU のエミュレーション。
- FPU レジスタ切り替え。

この論文では、ホスト OS にカーネル・レベル・コードを追加することにより、これらの問題を解決する方法について述べる。カーネル・レベル・コードが、ホスト OS を介さずに CPU を制御することにより、ホスト OS への依存度を小さくする。これによって、従来の LilyVM で実現していなかった、Solaris への移植 [6] および Windows XP への移植を実現した。また、カーネル・レベル・コードを利用することで上に示したオーバーヘッドを削減する。フォールトおよびトラップは、カーネル・レベルの

VMM でエミュレーションを行うことで高速に処理する。MMU のエミュレーションは、カーネル・レベル・コードでページ・テーブルを操作することで、オーバーヘッドを削減する。FPU レジスタの切り替えは、ハイパバイザ・コールを用いて、FPU レジスタを使用している場合にのみ行われるようにする。これにより、Linux カーネルの構築にかかる時間を約半分に短縮している。

## 2 関連研究

VMware Workstation[8] は、IA-32, および, x86-64 (AMD64, Intel 64) を対象とした完全仮想化に基づく VMM である。ホスト OS として、Linux, Windows, および、Mac OS X 上で動作する。VMware Workstation では実際のハードウェアをエミュレートしており、実機用の様々な OS をゲスト OS としてそのまま実行することができる。VMware Workstation は、本方式と同じように、カーネル・レベルの VMM を追加しているが、本方式では、ゲスト OS とホスト OS が協調して動作するような機能を容易に追加できるようにしている点が異なる。例えば、LilyVM ではユーザ・レベル OS からホスト OS の TCP/IP スタックを利用することができる [5]。

coLinux[1] は、IA-32 を対象として、ゲスト OS として動作するように移植された Linux である。ホスト OS として、Windows, および、Linux 上で動作する。coLinux では、ゲスト OS をカーネル・レベルで動作させているため、ゲスト OS がクラッシュすると、ホスト OS を含めてシステム全体がクラッシュすることがある。本方式は、ゲスト OS に対して特権を与えていないため、システム全体がクラッシュすることがない。また、本方式では、ゲスト OS に対する修正を言語処理系で自動化している点が異なる。

Xen[3] は、IA-32, x86-64, および、IA-64 などを対象とした準仮想化に基づく VMM である。Xen はホスト OS 上で動作するのではなく、実機上で動作する。Xen 2.0 で実現される仮想マシンのハードウェアは、実際のハードウェアとは異なるため、ゲスト OS に対する修正が必要となる。Xen 3.0 では、CPU による仮想化機能に対応し、従来の準仮想化による仮想マシンに加えて、完全仮想化による仮想マシンも作成できるようになった [9]。ただし完全

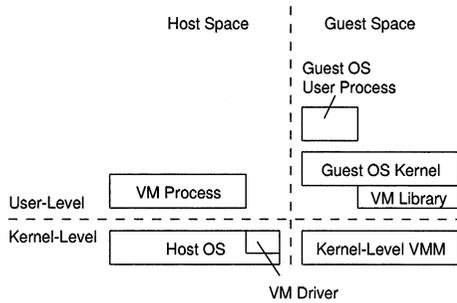


図 1: 本研究で提案する仮想マシン・モニタの構成

仮想化を用いると性能が低下する。本方式は、ホスト OS 上で動作するという点と、ゲスト OS に対する修正を言語処理系で自動化している点が異なる。

User-mode Linux[2] は、Linux を、他の Linux に移植したユーザ・レベル OS である。本方式は、ゲスト OS に対する修正を言語処理系で自動化している点が異なる。

### 3 基本設計

本研究で提案する VMM の構成を図 1 に示す。図において、ホスト空間 (Host Space) は、ホスト OS のカーネルとそのユーザ・プロセスが動作するアドレス空間を表し、ゲスト空間 (Guest Space) は、ゲスト OS のカーネル、そのユーザ・プロセスと、VMM が動作するアドレス空間を表している。LilyVM は以下の 4 つの要素から構成される。

- VM プロセス (VM process)。ホスト OS 上でユーザ・プロセスとして動作し、仮想マシン環境の資源を持つ。
- VM ドライバ (VM Driver)。VM プロセス、および、カーネル・レベル VMM との間の通信を行う。
- カーネル・レベル VMM (Kernel-level VMM)。ゲスト空間において高い特権レベルで動作し、例外、および、割り込み処理を行う。
- VM ライブラリ (VM library)。ゲスト OS カーネルの特権命令などのエミュレーションを行う。

ゲスト OS のカーネル (Guest OS Kernel) は、従来の LilyVM で用いているものと同様に、機械語命令の静的な変換を行ったものである。すなわち、特権命令や、センシティブな非特権命令が、VM ラ

イブラリ内の手続きを呼び出す命令に変換されている。また、カーネル・レベル VMM、および、VM ライブラリを同じアドレス空間に入れるため、アドレス空間の定数を変更している。また、ゲスト OS のカーネルは、入出力のための特別なデバイス・ドライバを含む。

また、ゲスト OS のユーザ・プロセス (Guest OS User Process) は、従来の LilyVM で用いているものと同様に、実機用のものをそのまま使用する。これにより、実機用のバイナリとして配布されている豊富なアプリケーションをそのまま活用することができる。

LilyVM の構成要素の中で、ホスト OS に依存している部分は、VM ドライバ、および、VM プロセスのメモリを確保する部分と、VM ドライバとの通信部分である。VM プロセスのその他の部分と、カーネル・レベル VMM、および、VM ライブラリは、ホスト OS に依存していない。このように、多くの構成要素をホスト OS に依存しないように作成することで、移植性を高めている。以下では、これらの構成要素の働きについて述べる。

#### 3.1 VM プロセス

VM プロセスは、以下のような処理を行うユーザ・プロセスである。

- 仮想マシン環境の資源 (CPU 時間やメモリ) を持つ。
- ゲスト OS を実行するための準備を行う。例えば、仮想マシン環境のためのメモリを確保したり、ゲスト OS カーネルを読み込んだりする。
- ホスト OS のシステム・コールを発行し、入出力を行う。
- 仮想マシンのハードウェア割り込みをエミュレートする。

VM プロセスは、ゲスト OS 実行の準備ができたら、VM ドライバに対してゲスト OS の実行を命令する。すると、ハードウェア割り込みが発生するか、ゲスト OS において停止命令が実行されるか、ディスクドライブ、ネットワーク等への入出力要求があるまでの間、ゲスト OS が実行される。VM プロセスは、ゲスト OS の実行が停止すると、仮想マシンのハードウェア割り込みをエミュレートするか、入出力の要求を処理し、VM ドライバに対してゲスト

OSの実行を続行させる。VMプロセスが入出力要求を処理することにより、入出力を標準的なAPIで実装できるため、移植性が向上している。

### 3.2 VMドライバ

VMドライバは、ホストOSのデバイス・ドライバとして高い特権レベルで動作する。VMドライバは、VMプロセス、および、カーネル・レベルVMMと通信し、割り込みベクタ・テーブルやセグメント等を切り替え、ゲスト空間に切り替える機能を持つ。

VMドライバは、VMプロセスからゲストOSの実行を命令されると、アドレス空間や割り込みベクタ・テーブルなどを、ゲスト空間に切り替える。実行がゲスト空間から戻ると、VMドライバは、ゲストOSが停止した原因を調べる。もし、ハードウェア割り込みが原因であった場合は、本来のホストOSの割り込み処理ルーチン呼び出す。

VMドライバは、ホストOSに依存している。実際に必要な機能の一覧を表1に示す。表には、Linuxにおける関数名を示した。これらの機能は、一般的なOSのカーネルで提供されているものであり、ホストOSへの依存度は小さくなっている。

### 3.3 カーネル・レベルVMM

カーネル・レベルVMMは、ゲストOSを実行している間に発生した例外や、ゲストOSユーザ・プロセスが発行するシステム・コールを、ゲストOSヘリダイレクトする。この処理は、ホストOSの例外ハンドラに依存しないので、移植性が向上している。

ゲストOSを実行している間にハードウェア割り込みが発生すると、カーネル・レベルVMMの割り込みハンドラに制御が移る。ハードウェア割り込み

表 1: VMドライバが使用しているホストOSカーネルの機能

ホストOS (Linux) の機能	説明
virt_to_phys(), phys_to_virt()	アドレスの変換
copy_from_user(), copy_to_user(), memcpy()	メモリのコピー
kmalloc(), kfree() register_chrdev(), unregister_chrdev()	メモリ管理 デバイス・ドライバの登録

を処理するのはホストOSであるため、割り込みハンドラは、アドレス空間や割り込みベクタ・テーブルなどをホスト空間に切り替える。その後、3.2節で述べたように、VMドライバによってホストOSカーネルの割り込みハンドラが呼び出される。

カーネル・レベルVMMは、VMライブラリにハイパバイザ・コールを提供する。ハイパバイザ・コールについて詳しくは4.3節で述べる。

### 3.4 VMライブラリ

VMライブラリは、従来のLilyVMと同様に、ゲストOSと同じアドレス空間にロードされ、特権命令や、センシティブな非特権命令のエミュレーションを行う。MMUの操作などは、カーネル・レベルVMMを呼び出す。

VMライブラリは、ゲストOSカーネルからホストOSの機能呼び出すための手続きも持っている。従来のLilyVMでは、それらの手続きの中で、VMライブラリから直接ホストOSのシステム・コールを発行していた。本方式では、入出力を行うハイパバイザ・コールは、VMプロセスで処理される。このため、オーバーヘッドが大きくなっているが、移植性は維持される。

## 4 カーネル・レベルVMMとVMライブラリの実装

ここでは、ホストOSに依存していないカーネル・レベルVMMと、カーネル・レベルVMMを呼び出すVMライブラリの実装について述べる。

### 4.1 MMUのエミュレーション

従来のLilyVMでは、ホストOSのmmap()システム・コールを用いて、ホストOSが持つページ・テーブルへの反映を行っていた。mmap()システム・コールでは、以下に示すような、非標準の機能を使わなければならないかかったり、ホストOSによって制限がある場合があった。

- 引数に指定したアドレスへの割り当てを強制するMAP\_FIXEDフラグ。
- MAP\_FIXEDフラグを使用した0番地への割り当て。

- 割り当て可能な領域の数の制限。たとえば Linux では、1 つのプロセスが割り当てられる領域は約 6 万個までに制限されている。

カーネル・レベル VMM では、ホスト OS の機能を使わずに、MMU を直接操作することにより、これらの問題を解決できるため、移植性が向上する。

ゲスト OS は、内部に自分自身で管理しているページ・テーブルを持っている。このページ・テーブルは、仮想マシンの仮想アドレスから、仮想マシンの物理アドレスへのマッピングを保持している。プログラムを実行するときに必要なのは、仮想マシンの仮想アドレスから実マシンの物理アドレスへのマッピングであり、カーネル・レベル VMM は、このページ・テーブルを持つ。

現在のカーネル・レベル VMM では、Xen と同様に、ゲスト OS 内にあるページ・テーブルのシャドウを VMM 内に持っている。ゲスト OS はプロセスごとにページ・テーブルを作成し、それらを切り替えながら動作する。ゲスト OS がページ・テーブルを切り替える時に、カーネル・レベル VMM はシャドウを切り替える。

MMU のエミュレーションのために、ページ・フォールトを利用してしている部分がある。これについては 4.2 節で述べる。

## 4.2 ハードウェア割り込み、例外処理とシステム・コール

従来の LilyVM では、例外処理とシステム・コールのリダイレクションを、`ptrace()` システム・コールとシグナルにより実現していた。これらの機能は OS に依存していた。例えば、ホスト OS が Linux の場合、`ptrace()` を用いてシステム・コールの発行を検出し、そのシステム・コールを無効化することができたが、NetBSD および FreeBSD ではそのようなことができなかった。

本方式では、ゲスト OS の実行中に発生した例外やシステム・コールは、カーネル・レベル VMM が受け取り、ゲスト OS にリダイレクトする。具体的には、まず、例外やシステム・コールのスタック・フレームをゲスト OS のメモリ上に作成し、スタック・ポインタを変更する。次に、プログラム・カウンタをゲスト OS カーネルのハンドラに変更する。

カーネル・レベル VMM は、基本的に例外をその

ままゲスト OS カーネルへリダイレクトする。しかし、下で述べる一般保護例外、浮動小数点演算に関する例外、および、ページ・フォールトについては、カーネル・レベル VMM で例外の内容を調べ、その他の処理をする場合がある。

IA-32 では、256 種類の割り込みがある [7]。そのうち、割り込み番号 0 から 31 までは、CPU の例外などに使用されるため予約されている。残りの割り込み番号 32 から 255 まだが、ハードウェア割り込みやトラップに使用される。しかし、使用される割り込み番号は、OS によって異なる。カーネル・レベル VMM は、ホスト OS に依存しないようにするため、IA-32 において予約されていない割り込み番号 32 から 255 までの割り込みを、ハードウェア割り込みとして扱う必要がある。しかし、Linux のユーザ・プロセスがシステム・コールを発行する時、トラップ命令により割り込み番号 128 の割り込みが生成される。したがって、すべての割り込みを無条件にハードウェア割り込みとして扱うことはできない。

この問題を解決するために、割り込み番号 32 から 255 の割り込みを生成するトラップ命令をユーザ・レベルから実行できないように特権レベルを設定し、一般保護例外を通じて、システム・コールを検出するようにした。また、現在、ハイパバイザ・コールのために割り込み番号 255 を使用しており、その検出にも一般保護例外を使用している。

一般保護例外は、上で述べたように、トラップの検出に使用している。トラップでない場合や、ゲスト OS によりユーザ・プロセスからの生成が禁止された割り込みである場合は、一般に、ユーザ・プロセスが不正な処理を行ったことを意味する。従って、例外をゲスト OS に処理させるために、ゲスト OS カーネルへのリダイレクションを行う。

ページ・フォールトは、MMU のエミュレーションにおいて重要なイベントのひとつである。ページ・フォールトが発生したときに、ゲスト OS のページ・テーブルをアクセスし、必要に応じてシャドウへの反映を行う。また、ページ・フォールトによって、ゲスト OS カーネルがページ・テーブルを変更しようとしていることも検出し、シャドウへの反映を行う。ゲスト OS のページ・テーブルの情報を調べた結果、ページが不在であったり、アクセス権のないページへのアクセスであった場合は、ゲスト OS カーネルにページ・フォールトを通知する。

浮動小数点演算に関する例外は、FPU レジスタの

保存と読み込みと、ゲスト OS の遅延評価機能（制御レジスタ CR0 の TS ビット）のために使用する。従来の LilyVM では、FPU を使わない場合でも、ゲスト OS のコンテキスト切り替え時に FPU レジスタの内容を入れ替える必要があった。その理由は、IA-32 の遅延評価機能をユーザ・プロセスから利用することができないためである。カーネル・レベル VMM では、ゲスト OS に対して遅延評価機能を提供することにより、この問題を解決している。

### 4.3 VM ライブラリとハイパバイザ・コール

カーネル・レベル VMM は、VM ライブラリに対してハイパバイザ・コールを提供する。実装したハイパバイザ・コールの一覧を表 2 に示す。これらは、従来の LilyVM で使用していたシステム・コールよりも、CPU の機能に近い、すなわち、抽象度の低い機能としている。これによって、カーネル・レベル VMM が多くの処理を行うことができ、性能が向上する。

VM ライブラリは、従来の LilyVM で用いていたものをベースとし、一部を修正して使用した。そのまま利用したのは、命令のエミュレーションを行う部分の一部である。たとえば、セグメント・レジスタやフラグ・レジスタの値を読み出す命令のエミュレーションを行う部分は、ホスト OS に依存していないため、そのまま利用した。修正した部分は、以下の点である。

- MMU のエミュレーション。
- セグメント機能のエミュレーション。
- シグナル処理。
- 時計、および、タイマーのエミュレーション。
- ディスク、ネットワーク、および、シリアル・ポートの入出力処理。

MMU のエミュレーションについては、性能を上

げるため、そのほとんどをカーネル・レベル VMM に実装している。そのため、MMU を制御する命令のエミュレーションにおいて、ゲスト OS のページ・テーブルにアクセスし、ホスト OS の `mmap()` システム・コールを使用していた部分をすべて削除し、ハイパバイザ・コールに置き換えた。

セグメント機能のエミュレーションについては、従来の LilyVM では、VM ライブラリからシステム・コールを発行していたために、次のような移植性を低下させる問題があった。

- FreeBSD や Linux などは、LDT (Local Descriptor Table) を変更する機能を持っているが、OS によって API が大きく異なる。
- Linux 2.6 系列には、GDT (Global Descriptor Table) の一部を変更するための機能である `set_thread_area()` システム・コールが追加されているが、Linux 2.4 以前や、FreeBSD, NetBSD などにそのような機能はない。

VM ライブラリでハイパバイザ・コールを使用し、カーネル・レベル VMM で GDT と LDT を変更することにより、これらの問題を解決した。ハイパバイザ・コールはホスト OS に依存しないように設計しているため、VM ライブラリの移植性が向上している。

従来の LilyVM では、例外を処理するために、VM ライブラリにおいてホスト OS に依存したシグナル機能を利用していた。本方式では、カーネル・レベル VMM が例外をすべて処理するようにした。

VM ライブラリにおいて、時計、および、タイマーのエミュレーションは、システム・コールを発行していた部分を、ハイパバイザ・コールに置き換えている。ハイパバイザ・コールの処理においてカーネル・レベル VMM は、実際のタイマーを制御してカウンタの値を取得し、その値と、ホスト空間からゲスト空間に切り替える時に得た時刻情報を元に、時刻を計算し、VM ライブラリに結果を返す。

表 2: ハイパバイザ・コールの一覧

ハイパバイザ・コール	説明
<code>set_cr3()</code>	ページ・テーブル (CR3 レジスタ) をセットする
<code>invalidate_page()</code>	指定されたページの TLB をフラッシュしページをマップする
<code>set_gdt()</code>	GDT のエントリを変更する
<code>set_ldt()</code>	LDT のエントリを変更する
<code>run_user()</code>	ゲスト OS のユーザ・プロセスに制御を移す
<code>call_host()</code>	ホスト OS 上の VM プロセスを呼び出す
<code>gettimeofday()</code>	時刻情報を得る

ディスク、ネットワーク、および、シリアル・ポートの入出力処理についても、システム・コールを発行していた部分を、ハイパバイザ・コールに置き換えている。カーネル・レベル VMM は、VM プロセスに入出力要求を伝える。

## 5 VM ドライバの実装

VM ドライバは、3 章で述べたように、カーネル・レベル・コードを用いる LilyVM において、ホスト OS に存在する部分である。本研究では、ホスト OS として Linux, Solaris, および、Windows XP に対して VM ドライバの実装を行った。ここでは、それぞれのホスト OS に対する実装について述べる。

### 5.1 Linux

まず Linux をホスト OS として VM ドライバを実装した。VM ドライバは、VM プロセスから `ioctl()` システム・コールを使用して呼び出される。VM ドライバは、ホスト OS の機能として、表 1 に示す機能を使用している。

### 5.2 Solaris

Linux の次に、Solaris をホスト OS として VM ドライバを実装した。Solaris においても、VM ドライバは、VM プロセスから `ioctl()` システム・コールを使用して呼び出される。そのため、VM プロセスは Linux のものと同一である。

VM ドライバは、表 1 に示す機能として、Linux の `kmalloc()` に対応する `kmem_alloc()` 関数、`kfree()` に対応する `kmem_free()` 関数を使用している。また、`virt_to_phys()` および `phys_to_virt()` に直接対応する関数はないが、`vmem_alloc()`、`hat_devload()`、`hat_unload()`、`vmem_free()` 関数を用いて対応する機能を実装した。

### 5.3 Windows XP

次に、Windows XP をホスト OS として VM ドライバを実装した。Windows XP においては、Windows XP の API は Unix のものと大きく異なるため、Unix と互換性のある API を実装している Cygwin

を使用した。しかし、Cygwin においても、`ioctl()` によるデバイス・ドライバの呼び出しができない。そこで、VM プロセスと VM ドライバの間の通信には `ioctl()` ではなく、Windows XP の `DeviceIoControl()` API を使用した。そのため、VM プロセスは、`ioctl()` の部分を除いて Linux 版と同一である。

VM ドライバは、表 1 に示す機能として、Linux の `virt_to_phys()` に対応する `MmGetPhysicalAddress()` 関数、`kmalloc()` に対応する `MmAllocateContiguousMemory()` 関数、`kfree()` に対応する `MmFreeContiguousMemory()` 関数を使用している。また、`phys_to_virt()` に直接対応する関数はないが、`MmMapIoSpace()`、`MmUnmapIoSpace()` 関数を用いて対応する機能を実装した。

Windows XP への実装において問題となるのは、アドレス空間である。全体で 4GB (1GB=2<sup>30</sup> バイト) あるアドレス空間のうち、カーネル用に 1GB を確保している Linux、0.5GB を確保している Solaris と異なり、Windows XP では 2GB が確保されている。しかし、ゲスト OS を実行するには残りの 2GB では足りないため、カーネル空間を残したままゲスト OS を実行することができない。そこで、ゲスト OS を実行する時は、アドレス空間全体を切り替え、ホスト OS のカーネルがアドレス空間に残らないようにした。

## 6 実験

カーネル・レベル VMM の性能を測定し従来の LilyVM との比較を行った。実験で用いたハードウェア、および、OS は以下の通りである。

**CPU** Pentium 4 3.00GHz

**ホスト OS のメイン・メモリ** 1GB

**ゲスト OS のメイン・メモリ** 64MB

**ホスト OS カーネル** Linux 2.6.15.4

**ゲスト OS カーネル** Linux 2.6.15.4

設定として、SMP サポートを無効とし、ホスト OS とゲスト OS のタイマー割り込みの頻度を 100Hz とした。また、VM プロセスに割り当てられたメモリがページ・アウトされないようにした。

実験では、実機 (ホスト OS)、従来方式、および、本方式において、標準設定の Linux 2.6.16 のコンパイルを行った。そのときの `make` コマンドの実行時間を測定した。結果を図 2 に示す。図 2 では、

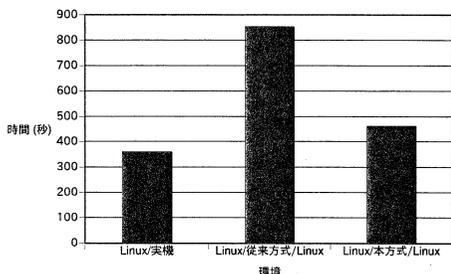


図 2: アプリケーションベンチマークの結果 (Linux 2.6.16 のコンパイル時間)

左から順に実機、従来方式、そして、本方式における実行時間を示している。このように、実行時間が従来の LilyVM と比べて約半分に短縮している。これは、カーネル・レベル VMM が例外のリダイレクションと MMU のエミュレーションを行うことにより、オーバヘッドを削減したためである。

また、本方式は実機の 1.3 倍程度遅くなっている。これは、ゲスト空間とホスト空間の間でアドレス空間や割り込みベクタ・テーブルなどを切り替えることによるオーバヘッドが大きいためである。

## 7 おわりに

この論文では、ホスト OS にカーネル・レベル・コードを追加することにより、ユーザ・レベル VMM の移植性を向上させる方法について述べた。その特徴は、カーネル・レベル・コードが、ホスト OS を介さずに CPU を制御することにより、ホスト OS への依存度を小さくしている点にある。LilyVM は、この方法により、新たに Solaris、および、Windows XP 上で動作するようになっている。

また、カーネル・レベル・コードを利用することで、フォールトおよびトラップのエミュレーションや、MMU のエミュレーションなどのオーバヘッドを削減している。これにより、Linux カーネルの構築にかかる時間を約半分に短縮している。

今後の課題は、MMU のエミュレーションにおいて、ホスト OS の VM プロセスに割り当てられたメモリがページ・アウトされた場合に対応することである。また、FreeBSD、NetBSD、および、Mac OS X といった OS への移植も実現したい。

## 謝辞

LilyVM の開発の一部は、萩谷プロジェクト・マネージャのもと、情報処理振興事業協会 (IPA) 平成 15 年度未踏ソフトウェア創造事業の援助を受けて行いました。萩谷プロジェクト・マネージャには、さまざまなご指導をいただきました。心より感謝いたします。

## 参考文献

- [1] Aloni, D.: Cooperative linux, *Proceedings of the Linux Symposium*, pp. 23–31 (July 2004).
- [2] Dike, J.: A user-mode port of the Linux kernel, *the 4th Annual Linux Showcase & Conference* (2000).
- [3] Dragovic, B., Fraser, K., Hand, S., Harris, T., Ho, A., Pratt, I., Warfield, A., Barham, P. and Neugebauer, R.: Xen and the Art of Virtualization (2003).
- [4] Eiraku, H. and Shinjo, Y.: Running BSD Kernels as User Processes by Partial Emulation and Rewriting of Machine Instructions, *USENIX BSDCon 2003 Conference (BSDCon'03)* (September 2003).
- [5] 榮樂英樹, 新城靖, 加藤和彦: ユーザレベル OS のためのユーザレベルネットワーク機能, 第 3 回情報科学技術フォーラム (FIT2004) (2004 年 9 月).
- [6] 榮樂英樹, 新城靖, 加藤和彦: 複合ハイパバイザによる仮想計算機の高速度化, 第 17 回コンピュータシステム・シンポジウム (ComSys2005) ポスター & デモ・セッション (2005 年 11 月).
- [7] Intel Corporation: *IA-32 Intel Architecture Software Developer's Manual* (2006).
- [8] Sugerman, J., Venkitachalam, G. and Lim, B.-H.: Virtualizing I/O Devices on VMware Workstation's Hosted Virtual Machine Monitor, *Proceedings of the General Track: 2002 USENIX Annual Technical Conference*, Berkeley, CA, USA, USENIX Association, pp. 1–14 (2001).
- [9] University of Cambridge, UK et al.: *Xen 3.0 Users' Manual*, <http://www.cl.cam.ac.uk/research/srg/netos/xen/readmes/user/user.html> (2006).