

## ヘテロな OS の計算資源を活用するグリッド RPC の設計

上村佳史<sup>†</sup> 中島佳宏<sup>†,††</sup> 佐藤三久<sup>†</sup>

大規模グリッド環境では従来のクラスタに加え、研究室やオフィスなどで用いられている PC を利用することが考えられる。このようなグリッド環境では各々の PC には Windows や Linux といった様々な OS が用いられていると想定される。そのためヘテロな OS 環境に対応した仕組みが Grid RPC に必要となる。本論文では、ヘテロな OS 環境に対応する Grid RPC システムとして、Linux から Windows を計算資源として利用するためのシステムを提案する。Windows を Grid RPC のワーカに利用し、再コンパイルなしに Linux プログラムを実行する方法として、Linux プログラムを Windows 上でそのまま実行可能とするための Linux プログラム実行環境 BEE の設計・開発を行った。BEE を用いた評価実験によりシステムコールの実行性能は Windows のネイティブプログラムと同程度の性能が得られた。また、OmniRPC を用いた実アプリケーションを Linux と Windows の混在環境での実行を行い、Windows と Linux のヘテロな OS 環境で Windows を Grid RPC のワーカとして利用できることを示した。

### Design of Grid RPC system utilizing computing resources of heterogeneous OS

YOSHIFUMI UEMURA,<sup>†</sup> YOSHIHIRO NAKAJIMA<sup>†</sup> and MITSUHISA SATO<sup>†</sup>

To build a large scale grid environment, PCs in laboratory and offices are expected to be used as computational resources in addition to dedicated conventional clusters. PCs are running under various kinds of OS including Windows and Linux. To make use of such PCs as workers of Grid RPC system, we propose a mechanism of Grid RPC system that utilizes Windows PCs as computing resources from Linux PC. We have designed and implemented a system called BEE to allow a workers program exported from Linux to run under Windows, and it is integrated into OmniRPC system. BEE enables Linux binary program to run under Windows without recompilation. We have evaluated its performance of system calls and realistic applications in heterogeneous OS environment that contain Windows and Linux. The result shows that system calls on BEE can be executed at almost the same performance with a Windows native binary, and the Grid RPC system successfully utilizes PC of Windows OS as computing resources in heterogeneous OS environment.

#### 1. はじめに

近年、ネットワーク技術の進歩により広域に分散した計算機資源を統合し、共有することで大規模計算などを行うグリッドコンピューティングに関する技術が注目されている。このグリッドコンピューティング環境上に分散した計算資源を利用するための並列プログラミングモデルの一つとして Grid RPC が提案されている。Grid RPC では、広域に分散された計算機資源を遠隔手続き呼び出し (Remote Procedure Call) の形で利用することが出来る。

多くの計算機資源を利用することが必要となる大規模グリッド環境を考えた場合、従来の PC クラスタに加えて、オフィスや研究室などで利用されている PC

を利用することが考えられる。しかし、このような PC には Windows や Linux などの様々な OS が用いられていると想定される。このようなヘテロな OS 環境に対応し、オフィスなどで用いられている PC を計算機資源に利用するための方法が Grid RPC に必要となる。本稿では、このようなヘテロな OS 環境に対応した Grid RPC システムとして、Windows と Linux からなるヘテロな OS 環境において Linux から Windows を Grid RPC のワーカに利用するためのシステムを提案する。

Linux から Windows を Grid RPC のワーカに利用する方法として、次のような方法が考えられる。仮想マシンと Linux プログラムを Windows プログラムへの再コンパイルである。まず、VMware<sup>®</sup> や QEMU<sup>®</sup> などの仮想マシンを用いる方法の場合、この仮想マシン上に Linux 環境を構築することができる。これにより仮想マシン上の Linux を Grid RPC のワーカとして利用できる。

次に、Linux プログラムを Windows プログラムへ再

<sup>†</sup> 筑波大学大学院 システム情報工学研究科

Graduate School of Systems and Information Engineering, University of Tsukuba

<sup>††</sup> 日本学術振興会特別研究員

Research Fellow of the Japan Society for the Promotion of Science

コンパイルする方法である。これは Cygwin<sup>3)</sup> 環境を用いることで簡単に実現できる。Cygwin では Linux プログラムとソースコードでの互換性があるため、Cygwin 上で Linux プログラムを再コンパイルすることにより、Windows で利用できるプログラムを作成することができる。

我々は、Windows と Linux のヘテロな OS 環境において Linux から Windows を利用する方法として、Linux で作成されたプログラムを、そのまま Windows 上で実行するシステムとして Linux プログラム実行環境 BEE を設計した。

本稿の構成は以下のようになっている。まず、2 章で提案する Linux プログラムを Windows 上で実行する方式について述べる。3 章では提案システムである Linux プログラム実行環境 BEE の設計と実装について述べる。4 章では Grid RPC として用いた OmniRPC の概要と、BEE を用いた場合の OmniRPC の RPC システムについて述べる。5 章では性能評価について述べる。最後にまとめと今後の課題について述べる。

## 2. Linux プログラムの Windows 上での実行方式

ここでは、提案するシステムである Linux プログラムを Windows 上でそのまま実行するシステムの特徴を、仮想マシンや再コンパイルする方法と比較を通して述べる。

仮想マシンを用いることで、通常の Linux と同じように扱うことができる。これにより Grid RPC システムに Windows に対応する仕組みを加えることなく Windows を計算機資源として利用することが可能となる。しかし、Windows 上に仮想マシンを導入し、さらに Linux の OS を導入しなければならないという問題がある。さらに、仮想マシンは専用のメモリやディスク領域を多く使用する。そのためホスト OS が仮想マシンに占有されてしまう問題がある。また、Windows を Grid RPC のワーカとして利用するのに Linux 環境全てが必要なわけではない。

Cygwin を利用して Linux プログラムを Windows プログラムに再コンパイルする方法の場合、ソースコードでの互換性があるため、再コンパイルすることで Windows プログラムとすることができる。しかし、この場合、Linux に加えて Windows 環境を用意し、Windows プログラムを用意する必要がある。また、もとの Grid RPC に Windows に対応した仕組みを加えなければならない問題がある。

提案する Linux プログラムをそのまま Windows 上で実行するシステムの場合、Windows に仮想マシンほどの環境や Linux OS 環境は必要とせず済む。そのため、仮想マシンで使用されるメモリやディスクといった専用の領域が必要ではなくなる。また、Linux プロ

グラムを Windows 上で実行できるため、Windows 上のプログラムを必要としなくなる。

## 3. Linux プログラム実行環境 BEE の設計と実装

ヘテロな OS 環境において、Linux から Windows を Grid RPC のワーカとなる計算機資源として利用するために、Linux プログラム実行環境として BEE を開発した。ここでは、BEE の設計と実装について述べる。

Linux は様々なアーキテクチャに対応しているが、BEE では Linux プログラムをそのまま Windows 上で実行することを目的として実装しているため、IA-32 アーキテクチャ上で動作する Linux のプログラムを対象とする。

BEE のシステムは、2 つの部分から構成されている。1 つはプログラムローダである。Linux では ELF などの多種のプログラムのフォーマットが用いられているが、Windows では PE フォーマットなどがプログラムのフォーマットとなっている。そのため、Linux プログラムのフォーマットは Windows では対応していないため実行することができない。そこで、Linux プログラムのフォーマットを解析、実行するためのプログラムローダの実装をしている。

もう 1 つはシステムコールである。Windows と Linux で同じ IA-32 アーキテクチャを用いた場合、プログラム自体はそのまま実行することは可能である。しかし、システムコールの実装は OS 毎に異なっている。そのため、同じアーキテクチャであっても Windows 上で Linux プログラムをそのまま実行することが出来ない。そこで、Linux のシステムコールを実行するための機能を実装している。

次に、プログラムローダとシステムコールのそれぞれの実装について述べる。

### 3.1 プログラムローダ

Linux のプログラムを実行するには、そのフォーマットに対応したプログラムローダを実装することで実行可能となる。Linux では、プログラムのフォーマットは多くの種類が存在している。BEE では一般的に普及しているフォーマットである ELF<sup>1)</sup> を対象としてプログラムローダの実装を行った。また、ELF には動的リンクと静的リンクがある。BEE ではプログラム単体で実行することが可能となる静的リンクした実行プログラムのみを対象とした。これは、動的リンクの場合、共有ライブラリを Windows が持たなければならないためと、共有ライブラリのバージョンが異なっていた場合に実行することができなくなってしまうという問題が起きるためである。プログラムローダは Linux プログラムのフォーマットの解析を行い、ローダプログラム上にプログラムを展開し実行する。このときの仮想メモリの状態を図 1 に示す。

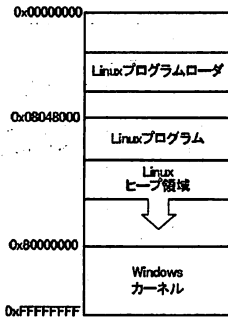


図1 プログラムローダによる Linux プログラムの展開

### 3.2 システムコール

入出力などの処理はシステムコールを発行することで実行できる。このシステムコールは OS により異なる実装がなされている。そのため Linux プログラムを Windows 上で実行する場合、システムコールの実装が異なるため実行することができない。Windows でのシステムコールは、NT/2000 の場合は *int 2E*、XP などの場合では *sysenter* 命令を用いて、ソフトウェア割り込みを発生させることで実行している。Linux のシステムコールは *int 0x80* により、ソフトウェア割り込みを発生させ実行している。このように Windows では *int 0x80* に対応する割り込み処理は無いためである。

Windows では Linux のシステムコールに用いられている *int 0x80* に対応する割り込み処理は特に指定されていない。そこで、BEE では *int 0x80* に対応するソフトウェア割り込み処理を、新たに追加することでシステムコールの実行を可能にしている。新たにソフトウェア割り込みの処理を追加する方法として BEE ではデバイスドライバを利用している。

割り込み処理を追加には、割り込みディスクリプタテーブル (Interrupt Descriptor Table)<sup>5)</sup> を操作する必要があり、割り込みディスクリプタテーブルに指定する割り込み処理はカーネル上に用意しておかなければならない。通常のユーザプログラムではユーザ空間上にプログラムが展開されてしまうため、これを割り込み処理として割り込みディスクリプタテーブルに指定することができない。しかし、デバイスドライバを用いた場合は可能となる。デバイスドライバは OS の起動時に OS に組み込まれカーネル上に展開される。そのためデバイスドライバを利用して、カーネル上に新たな割り込みの処理を追加することができる。これにより、Linux のシステムコールに対応した処理などの任意のソフトウェア割り込み処理などを実装することができる。デバイスドライバを利用して割り込みディスクリプタテーブルに割り込み処理を追加したときの実行の流れを図 2 に示す。

システムコールのソフトウェア割り込み処理はデバイスドライバを利用することで可能となるが、実際の

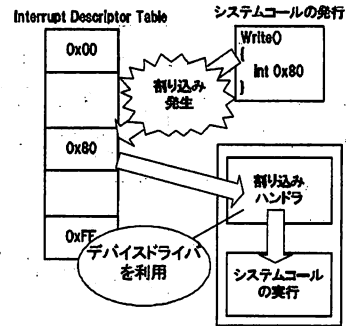


図2 並列固有値計算プログラムによる性能評価

システムコールの処理を実装する必要がある。Linux ではシステムコールは多くの種類が存在する。しかし、Windows を Grid RPC のワーカとして利用するのに、このような数多くのシステムコールを実装する必要はない。そこで、BEE ではワーカとして利用するのに必要最低限のシステムコールのみを実装する。現在実装しているシステムコールは *write/read* などの基本的なファイル I/O と、マスタとのデータのやり取りに必要なネットワーク I/O のみを実装している。

### 4. OmniRPC の概要と BEE を用いた OmniRPC の RPC システム

Linux プログラム実行環境 BEE で用いる Grid RPC システムとして、OmniRPC<sup>7)6)</sup> を使用した。OmniRPC は、ローカルなクラスタ環境から広域ネットワークで構成されたグリッド環境までのシームレスな並列プログラミングを可能にするマスタ/ワーカ型の Grid RPC システムである。

OmniRPC のシステムを図 3 に示す。OmniRPC の Linux での実装では、マスタからの遠隔手続き呼び出しにより、リモート実行プログラムを実行する。このとき、直接ワーカのリモート実行プログラムを実行はせずに、ssh などにより認証を行った後に Agent をリモートノード上に起動する。その後、マスタからの RPC コールは Agent を経由し、リモート実行プログラムの実行は Agent が行う。

BEE を用いた場合の Windows での OmniRPC システムを図 4 に示す。Windows はワーカとしての利用をするため、マスタとして利用するのは Linux のみである。OmniRPC は Agent を ssh などによる認証を行った後に起動するが、Windows の場合、このような ssh による認証などは通常はできない。そのため Agent は常に起動した状態で動作するように変更を加える必要がある。また、リモート実行プログラムは Linux プログラムであるので、Agent は Linux プログラム実行環境である BEE を経由してプログラムを実行するように変更する。後は Linux での OmniRPC と同じように

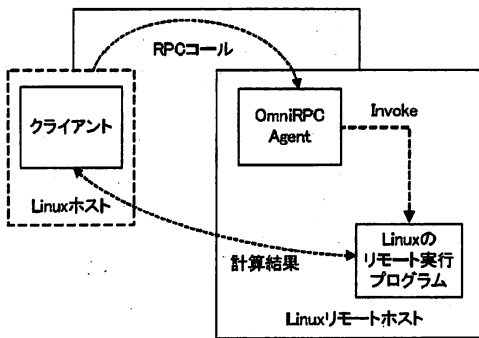


図3 OmniRPCのRPCシステム

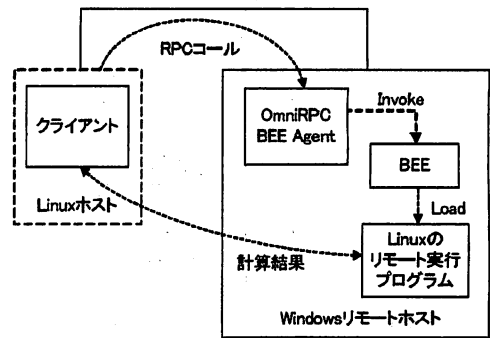


図4 BEEを用いたOmniRPCのRPCシステム

BEEにより実行されたLinuxプログラムがマスタであるLinuxホストと通信を行う。

また、OmniRPCで想定しているグリッド環境では、インターネット上で複数のクラスタ環境が相互接続されている環境である。そのような環境では、それぞれのクラスタ環境内では通常ファイルの共有がなされている。しかし、大規模グリッド環境を想定した遊休PCを利用する場合には、ノードは単一のPCとして存在するためファイルの共有をすることができない。そこで、単一のPCを簡便に利用できるようにするために、ワーカが必要なリモート実行プログラムを所持していない場合には、マスタが自動的に必要とするリモート実行プログラムを配布する機能を実装した。これにより、WindowsやLinuxにおいてリモート実行プログラムなどをあらかじめ取得しておくといった必要がなくなった。

## 5. 性能評価

提案システムであるBEEとの性能比較には、VMwareとCygwinを用いて行った。性能評価には、それぞれの環境におけるシステムコールの実行性能と通信バンド幅の性能、実アプリケーションによる性能としてOmniRPCを用いた並列固有値計算プログラムにより評価する。また、ヘテロなOS環境で実行としてWindowsとLinuxの混在環境で動作確認を行った。

### 5.1 実験環境

実験は本システムを用いる環境である、マスタにLinuxを使用しワーカにWindowsを使用する環境で行った。表1にそれぞれのマシン環境とVMware環境、使用したノード数を示す。VMwareとCygwin環境はWindowsと同じマシン環境上に構築した。

### 5.2 システムコールの性能評価

システムコールの性能評価は、基本的なファイルへのWrite/Readシステムコールにより行った。測定の方法は、1万回のシステムコールを発行し、バッファサイズを変化していったときの実行時間を比較する

Master	Xeon 2.4GHz, メモリ 1GByte, Gigabit Ethernet Linux kernel 2.6.9, gcc 3.4.5, 1ノード
Worker (Windows)	Xeon 3.2GHz, メモリ 2GByte, Gigabit Ethernet Windows Server 2003 Enterprise Edition SP1 4ノード
Worker (Linux)	Xeon 2.4GHz, メモリ 1GByte, Gigabit Ethernet Linux kernel 2.6.9, 4ノード
VMware	仮想メモリ 1GByte, ネットワーク接続ブリッジ Linux kernel 2.6.11, gcc 4.0.0

ことにより評価をした。また、性能比較の基準としてWindowsのネイティブプログラムの測定も行っている。

ファイルへのwriteシステムコールを行ったときの測定結果を図5、readシステムコールを行ったときの測定結果を図6に示す。図の縦軸が、1万回のシステムコールを発行するのに要した時間であり、10回測定した平均値である。横軸がそのときのバッファサイズである。

BEEはWindowsのネイティブプログラムやCygwinと同程度の性能が得られているが、1回のシステムコールの発行に平均で、Windowsのネイティブプログラムと0.6μsec、Cygwinと0.5μsecほど遅くなっている。これは、BEEの実装ではLinuxプログラムによるシステムコール発行時とWindows APIを呼び出す時の2度の割り込みの発生してしまうため、通常のプログラムよりも割り込みが多く発生することによるオーバーヘッドのためであると考えられる。また、VMwareではWrite、Readともにバッファサイズが10KByteあたりから実行性能が大きく低下している。

### 5.3 通信バンド幅の性能評価

通信バンド幅の性能測定はLinuxとそれぞれの環境で行った。図7に結果を示す。図の横軸がメッセージサイズである。このメッセージサイズを用いて、各ノード間でメッセージのping-pongを行い、この送受信に要した時間から転送バンド幅を算出した。

通信バンド幅においてもBEEはWindowsのネイティブプログラムと同程度の93.9MByte/secの性能が

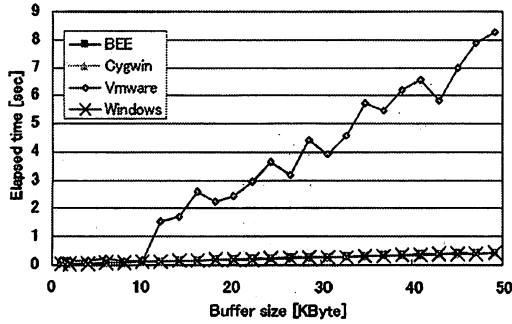


図5 ファイルへの write システムコールの実行時間

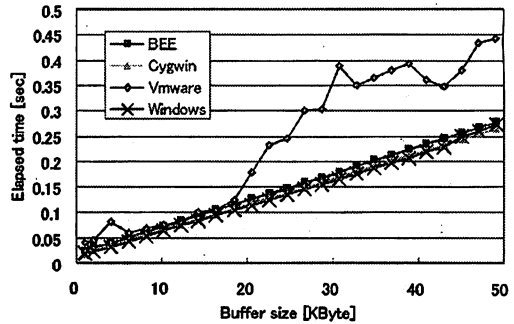


図6 ファイルへの read システムコールの実行時間

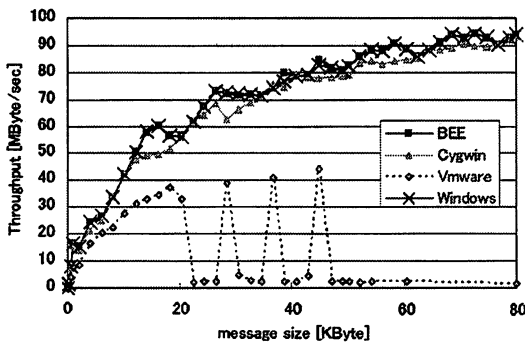


図7 環境によるスループットの比較

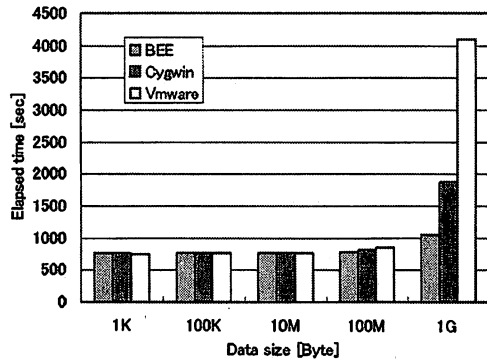


図8 OmniRPCを用いたデータサイズによる実行速度の比較

得られた。しかし、Cygwinでは93.6MByte/secと多少性能が低下している。これはCygwinではメッセージ転送の際にメモリコピーが発生しているためである。また、VMwareではバッファサイズが20KByteあたりから急激な性能低下が発生している。

#### 5.4 OmniRPCでの実行性能

予備評価として、Grid RPCを用いて1ジョブに30秒を必要とし、これが100ジョブある場合の性能を示す。Grid RPCシステムとしてOmniRPCを用い、それぞれの環境でWindowsを4ノード用いた。図8に結果を示す。横軸が1ジョブを実行するときのワークに転送するデータサイズである。

この結果より、データサイズが100MByteまでは実行性能にほとんど差は見られない。データサイズが1GByteになるとデータ転送に時間がかかってしまい、BEEやCygwinでは実際に使用されたのは3ノード、VMwareではネットワーク速度などの性能が悪いため2ノードであった。また、逐次での実行性能に比べ、4ノードで4倍の実行性能が出ている。

#### 5.5 並列固有値計算プログラムを用いた性能評価

実アプリケーションとしてOmniRPCを用いた並列固有値計算プログラム<sup>9)</sup>による性能評価を行った。このプログラムは、大規模な一般固有値問題のための解法であり、複素空間上の円周上の点に対応する方程式を並列に解くことにより、その円周内にある固有値を効率的に求めるものである。図9に、BEE、Cygwin、VMwareのそれぞれの環境で、ワークを2ノード、4ノードとしたときの実行結果を示す。図の縦軸は5回実行したときの実行時間の平均値である。

この結果より、BEE環境がCygwin環境よりも若干遅くなっていることがわかる。このような結果となったのは、まずOmniRPCの通信は1KByte単位で送受信を行っている。そのため、Cygwinと比較した場合に通信性能で、BEEの割り込みによるオーバーヘッドが現れるためと考えられる。また、並列固有値計算プログラムでは通信に加えてmmapシステムコールが多く回数実行されている。しかし、現在のBEEでのmmapシステムコールの実装は、あまり効率的な方法とはなっていない。そのため、mmapシステムコール

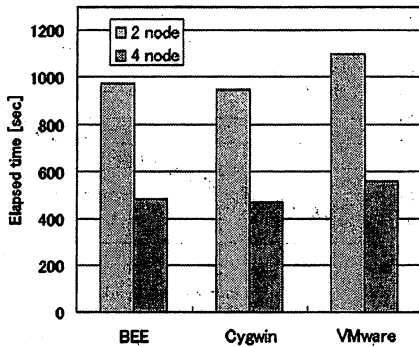


図9 並列固有値計算による性能評価

が多くの回数実行されるようなプログラムでは、かなりのオーバーヘッドとなってしまふと考えられる。

### 5.6 Windows と Linux の混在環境での実行

並列固有値計算プログラムを用いて、Windows と Linux の混在環境として、Windows と Linux を2ノードずつと4ノードずつの場合で測定を行った。図10に結果を示す。ここでWindows2ノードとLinux2ノードの実行時間がWindows4ノードよりも遅くなっているのは、Linuxマシンの性能がWindowsマシンの性能と異なるためである。この結果より、Windows と Linux が混在した場合でも利用できることが確認できた。

## 6. まとめと今後の課題

本稿ではヘテロなOSを計算機資源とするGrid RPCシステムとして、Linux と Windows のヘテロなOS環境におけるLinuxからのWindowsのワーカ利用の方法として、Windows上でLinuxプログラムを実行するための実行環境BEEを開発・評価を行った。

BEEではプログラムローダと、システムコールのための割り込み処理をWindowsに追加することで、Windows上でLinuxプログラムを実行させることを可能にした。BEEにより、仮想化環境やWindows用のプログラムを必要とせずにGrid RPCのワーカとして利用できるようになった。

性能評価の結果、システムコールの実行性能はWindowsネイティブプログラムと同程度の性能が得られた。また、Linux と Windows の混在環境においても実行を確認し、ヘテロなOS環境でWindowsをGrid RPCのワーカとして利用することが可能となった。

今後の課題として、Linux kernelのバージョンが2.6になってから、スレッド関連の実装<sup>4)</sup>が変化している。また、新しいスレッドモデルの実装はglibcなどのバージョンなどとも関連している。そのためプログラムの実行方法も変化しているため、現在のBEEでは実行することが不可能である。このような問題があるため、

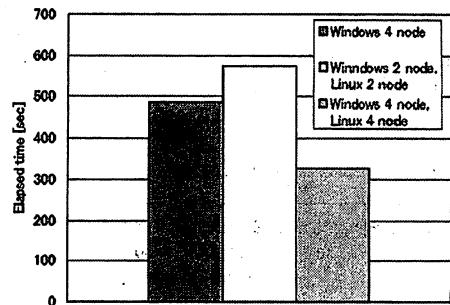


図10 並列固有値計算のWindowsとLinuxの混在環境での実行

動作できるようにBEEを変更していくことを考えている。また、Linux と Windows 間でのプロセスマイグレーションの実装を検討している。

謝辞 本研究の一部は、文部科学省科学研究費補助金 基盤研究A 課題番号172002, 特別研究員奨励費 課題番号17・7324, および、日仏共同研究プログラム(SAKURA)による。

## 参考文献

- 1) TIS Committee: Tool Interface Standard(TIS) Executable and Linking Format(ELF) Specification version 1.2 (1995).
- 2) Bellard, F.: QEMU, a Fast and Portable Dynamic Translator, *Proceedings of USENIX 2005 Annual Technical Conference*, pp. 41-46 (2005).
- 3) Cygwin: <http://sourceware.org/cygwin/>.
- 4) Drepper, U. and Molnar, I.: The Native POSIX Thread Library for Linux, <http://people.redhat.com/drepper/nptl-design.pdf>.
- 5) Intel(R): IA-32 IntelR Architecture Software Developer's Manuals (2004).
- 6) Nakajima, Y., Sato, M., Boku, T., Takahashi, D. and Gotoh, H.: Performance Evaluation of OmniRPC in a Grid Environment., *SAINT Workshops*, pp. 658-665 (2004).
- 7) OmniRPC: <http://www.omni.hpcc.jp/OmniRPC/>.
- 8) VMware: <http://www.vmware.com/>.
- 9) 櫻井鉄也, 多田野寛人, 早川賢太郎, 佐藤三久, 高橋大介, 長嶋雲兵, 稲富雄一, 渡邊寿雄梅田宏明: 大規模固有値問題の master-worker 型並列解法, 情報処理学会 ACS 論文誌, Vol. 46, No. 10, pp. 1-8 (2005).