

## SPUMONE 仮想環境下におけるリカバリ指向 OS 向け OS 間通信の実装

福島 拓† 石井正将† 巻島 一雄† 杵渕 雄樹† 中島 達夫†

早稲田大学理工学部コンピュータ・ネットワーク工学科 分散ユビキタスコンピューティング研究室‡

### 1 はじめに

我々の研究室では、新たな OS アーキテクチャを考案すると共に、その実装である SPUMONE と呼ばれるシステムを開発している [1, 2].

SPUMONE は組込みシステムにおける高パフォーマンス、信頼性、可用性、保守性、リアルタイム性を得るために、ハードウェア層と OS の間に仮想化層を設けて CPU 仮想レイヤーを提供する。SPUMONE は、単一の CPU を複数個の CPU があるように見せかけるもので、仮想 CPU をゲスト OS に提供することで、同時に複数の OS を動作させることを可能にしている。本研究では SPUMONE 上で複数のゲスト OS を協調して動作させる際に必要となる OS 間での通信 (IOC: InterOS Communication) について考察する。また、近年の高可用性を求める動きに合わせ我々の研究室では問題が起きたときに高速に再起動することで問題を解決するリカバリ指向 OS のテストを行っており、そのために SimpleOS という独自の OS を開発している。本研究では SPUMONE 上でゲスト OS として Linux と SimpleOS を動作させ、故障したと仮定して Linux 側から SimpleOS を再起動できるようにすることも目標としている。

### 2 関連研究

SPUMONE とは異なるハイパーバイザ型の仮想化技術のひとつに Xen[3] が挙げられる。Xen 上で動作する OS 間のデータ転送には TCP/IP を用いた方式と共有メモリを用いた方式が存在する。TCP/IP を用いた方式は一般的なインターフェースであるため多くの OS で利用されやすいが、TCP/IP の処理によるオーバーヘッドが大きく高パフォーマンスを得ることができない。一方共有メモリを利用した方式は高速であるが、通信に関する規約を自分達で定義し実装なくてはならない。そこでコントロール情報の転送のみに TCP/IP を採用し、データの転送には共有メモリを利用することで高速なデータ転送を実現した XenSocket[4] や Xwayl[5] が考案されている。さらに、Lamia Youseff らは Xen の shared page を改良してネイティブに共有メモリへアクセスで

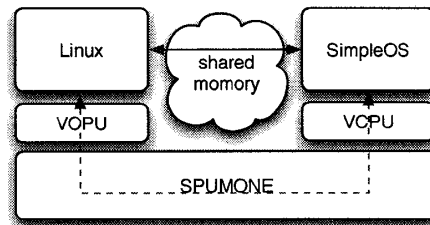


図 1: SPUMONE 上で動作する Linux と SimpleOS 間の IOC

きるようにすることによってハイパーコールの呼び出しを最小限に押さえてパフォーマンスを向上させることに成功している [6].

### 3 設計と実装

本研究では組込みシステム向けにマルチ OS 環境を用いることを想定しており、リアルタイム性が求められるから、IOC の通信方式として動作が高速な共有メモリ方式を用いることにした。また、データの一意性を保証する軽量なプロトコルを実装した。このように統一されたインターフェースを提供することでアプリケーション開発者はエンジニアリングコストを抑えることができる。

#### 3.1 SimpleOS 上での IOC の設計

データの転送を通知する手段として、あるゲスト OS が IOC を行なった際に SPUMONE が別のゲスト OS に仮想的な割り込みを発生させることを利用することにした。これにより IOC が行なわれるたびに割り込みが発生し、即座にあるゲスト OS から別のゲスト OS へと IOC が行われたことを通知することができる。通知された OS はあらかじめ割り込みハンドラを用意しておくことによって共有メモリに書き込まれた内容を読みとることができる。

#### 3.2 SimpleOS 上での IOC の実装

本研究では Linux が共有メモリに書き込みを行う際に SPUMONE の API を呼び出し、共有メモリへの書き込

†Taku FUKUSHIMA, Masaki ISHII, Kazuo MAKIJIMA, Yuki KINEBUCHI, Tatsuo NAKAJIMA

‡Department of Information and Computer Science, Faculty of Science and Engineering, Waseda University

みを IPI\*として SimpleOS へと伝える命令を SPUMONE に渡す。図 1 の破線がこの IPI 命令のパスに相当し、実線はデータのパスである。これを受けとった SPUMONE は SimpleOS に対し割込みを発生させる。すると SimpleOS は割込みハンドラへと処理を移し、Linux によって書き込みが行なわれた共有メモリを参照する。通常 Linux などで共有メモリの読み書きをする際には OS の標準ライブラリを利用するが、SimpleOS には標準ライブラリのようなものが存在しないのでそれらの必要な処理も実装した。

さらに、共有メモリの番地や大きさは SPUMONE に定義されており、IOC 用のディスクリプタテーブルが複数の `ioc_mem_t` という構造体で構成されている。それぞれの `ioc_mem_t` は、読み込み用と書き込み用の `ioc_mem_fifo` という配列で実装された円環状の FIFO 型データ構造を含んでいる。図 2 に共有メモリ構成の概略を示す。この FIFO 型のデータ構造への読み込みおよび書き込みは非同期に行なわれ、したがってそのままではデータの一貫性が保証されない。例えば、あるプロセス A があるメモリ領域を参照したところでコンテキストスイッチが起こり、別のプロセス B がその値を変更してしまった後で再びプロセス A へ制御が移った場合はそのメモリ領域の値はプロセス A が想定していないものになってしまう。

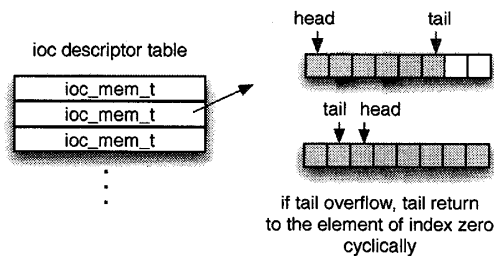


図 2: IOC に使用される共有メモリ構成の概略

#### 4 考察

前章で提示した共有メモリを用いた非同期通信時に発生する問題を ABA 問題という (図 3)。ABA 問題を回避し、データの一貫性を保つためにはセマフォを利用した排他制御や Lock-Free Queue [7] などを利用する方法がある。Lock-Free Queue は Fetch&Add (FAA) や Compare&Swap (CSW) などのアトミックに行なわれる Read-Modify-Write を利用して実装され、ロックを利用したものよりも高いパフォーマンスを実現することが

できる。アトミックな加算や減算はハードウェアによってサポートされている命令を使用し、SPUMONE が動作する SH4A アーキテクチャ上ではアセンブラの `movli` 命令および `movco` 命令で実現される。

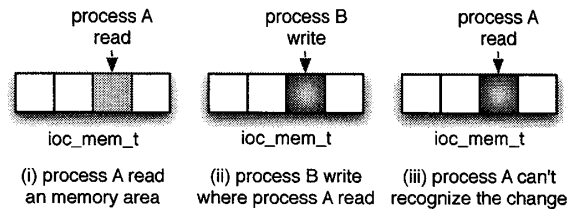


図 3: ABA 問題の例

#### 5 まとめと今後の課題

SPUMONE 上で動作する独自 OS である SimpleOS 上に IOC を実装することに成功し、Linux と SimpleOS の間で通信をすることができるようになった。これにより Linux から SimpleOS を再起動させることが可能になり、例えば異常を検知したときに再起動させるモニタリングサービスのテストができるようになった。そういったモニタリングサービスは本研究の範囲外となるので別の機会に譲りたいと思う。また、IOC は IPC と非常に似ているため、IPC で行なわれている研究を応用しやすいと考えられる。例えばメッセージキューやメールボックス型の通信を行うことでより強固で柔軟性のある通信ができると考えられる。

#### 参考文献

- [1] Tatsuo Nakajima, Hiroo Ishikawa, Yuki Kinebuchi, Midori Sugaya, Sun Lei, Alexandre Courbot, Andrej van der Zee, Aleksi Aalto, and Know Ki Duk: An Operating System Architecture for Future Information Appliances (2008)
- [2] Wataru Kanda, Yu Yumura, Yuki Kinebuchi, Kazuo Makijima: SPUMONE: Lightweight CPU Virtualization Layer for Embedded Systems
- [3] Paul Barham, Boris Dragovic, Keir Fraser, Steven Hand, Tim Harris, Alex Ho, Rolf Neugebauer, Ian Pratt, Andrew Warfield: Xen and the Art of Virtualization (2003)
- [4] Xiaolan Zhang, Suzanne McIntosh, Pankaj Rohatgi, and John Linwood Griffin: XenSocket: A High-Throughput Interdomain Transport for Virtual Machines
- [5] Xway: <http://sourceforge.net/projects/xway/>
- [6] Lamia Youseff, Dmitrii Zagorodnov, and Rich Wolski: Inter-OS Communication on Highly Parallel Multi-Core Architecture
- [7] Jhon D. Valois: Implementing Lock-Free Queues (1994)

\*IPI: InterProcessor Interruption