

SPUMONE 利用環境での Lock Holder Preemption の回避方法の検討と実装

三嶽 仁 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

1 はじめに

主にエンジニアリングコストの削減等を目的として、仮想化技術は組み込み機器にとっても重要な技術となりつつある。だが、仮想化環境は物理的なハードウェアを対象として設計、実装された既存のオペレーティングシステム的前提を常に提供可能とは限らず、Semantic Gap による予期せぬ弊害を発生させる危険性ももたらす。例えば、リアルタイム OS と汎用 OS を一つのハードウェア上で実行させた場合、汎用 OS 側のパフォーマンスが大きく低下することが先行研究で判明している。

本研究では、準仮想化技術 SPUMONE[3] を利用した環境での、弊害の代表的な 1 つである Lock Holder Preemption(LHP) の回避方法を提案する。RTOS には TOPPERS/JSP を、汎用 OS には Linux を利用した。

LHP は、ロックを獲得している Linux のスレッドを実行している仮想 CPU が RTOS の仮想 CPU に物理 CPU を奪われることで発生する。回避方法には、SPUMONE の仮想 CPU マイグレーション機能を用いる。仮想 CPU マイグレーション機能は、物理 CPU とその上で実行される仮想 CPU を、起動後も動的に変更可能にするものである。RTOS の仮想 CPU が物理 CPU にスケジューされる際、その物理 CPU 上で実行されている仮想 CPU がロックを獲得中のスレッドを実行中であれば、別の物理 CPU に RTOS の仮想 CPU をマイグレーションするという方式で実装を行った。

2 関連研究

Uhlig らは、準仮想化と完全仮想化の両方の方式で LHP を回避する方法を提案、実装している [1]。

この論文の提案する準仮想化環境での LHP の回避方法は、ゲスト OS を修正し、ロックを獲得していることを VM 側から知ることが出来るようにするというものである。VM がゲスト OS のスケジューリングを行う際、その情報に基づいて、プリエンプションを遅延させることにより LHP を防ぐ。

完全仮想化環境での LHP の回避方法は、カーネル空間ではほぼ全ての時間 1 つ以上のロックが獲得されているという経験則と、IA-32 プロセッサでは HLT 命令を実行した後は、割り込みを受け取るまではアイドルループを実行しているという性質を利用する。カーネ

ル空間とユーザ空間のスイッチや HLT 命令、割り込みは VM からモニタリング可能なので、それらの情報に基づいてゲスト OS がプリエンプト可能か不可能かを判断し、必要に応じてプリエンプトを遅延しスケジューリングを行う。

Wells らは、プロセッサを拡張し、ハードウェア的にスピロックの獲得を監視可能にし、その情報を LHP の回避に利用する方法を提案、実装している [2]。一定の時間内に実行された命令数が閾値以下であればゲスト OS はスピロックを獲得していると判断する、という機能をプロセッサに実装し、ゲスト OS の修正を全く行わないでロックの獲得状況を把握している。

3 提案方式

本研究では、仮想マシンモニタ SPUMONE から、Linux のスレッドがプリエンプト可能か否かを把握し、RTOS をディスパッチするコアを動的に変更することで可能な限り LHP を回避する方法を提案する。

仮想 CPU マイグレーション機能は、`rtos_cpu_affinity` というビットマスクとして利用される変数で RTOS の仮想 CPU が実行可能な物理 CPU を表現する。本研究では実験用に 4 コアのマシンを利用したので、下位 4 ビットのみが問題となる。全てのコアが RTOS を実行可能であれば、下位 4 ビットは 1111 となる。コア 0、コア 1 のみが実行可能であれば、0011 となる。Linux の全てのコアでロックが獲得されている場合は、0000 となる。

Linux カーネルはビルド時に `CONFIG_PREEMPT` オプションを有効にすると、`preempt_enable()` と `preempt_disable()` という 2 つのマクロが意味あるものとして定義される。これは、現在実行されているスレッドがプリエンプト可能か不可能かを識別するための情報を変更するマクロである。例えば、スピロックを獲得しているスレッドは、プリエンプトされると LHP のように無駄なビジーループを他のコアで発生させる原因になる。そのため、カーネル空間で実行中のスレッドのプリエンプトを許可する `CONFIG_PREEMPT` オプションが有効にされている場合でも、`preempt_count` が 1 以上の場合は、プリエンプトは発生しない。この `preempt_enable()` と `preempt_disable()` に変更を加え、`preempt_count` が 0 から 1 になった場合に `rtos_cpu_affinity` から実行中のコアのフ

ラグをクリアしてプリエンプトの回避を SPUMONE に依頼し, preempt_count が 1 から 0 になった場合にフラグをセットしてプリエンプト可能であることを SPUMONE に伝える。

Linux の全ての仮想 CPU でプリエンプト不可能な状態のスレッドが実行されている場合は, RTOS が直前にスケジュールされたコアの Linux の仮想 CPU をプリエンプトする。この場合は LHP が発生してしまうが, RTOS のリアルタイム性を保証するための措置である。

4 評価と考察

評価環境として用いたハードウェアの仕様を表 1 に示す。

表 1: 評価環境

ターゲットボード	RP1 評価ボード
CPU アーキテクチャ	SH4A
CPU	600MHz × 4
メモリ	128MB

評価用アプリケーションには, Rusty Russell 氏の作成したスケジューラ用ベンチマークプログラム hackbench を, より現実的なワークロードを表現するように修正したものを利用した。変更点は, 組込み機器には過剰と思われる 1 グループ 40 個の IPC 用プロセスを 4 個に減らし, 無駄なメモリコピーをそれぞれのプロセスで行わせることでユーザランドでの実行時間を増加させたという 2 点である。

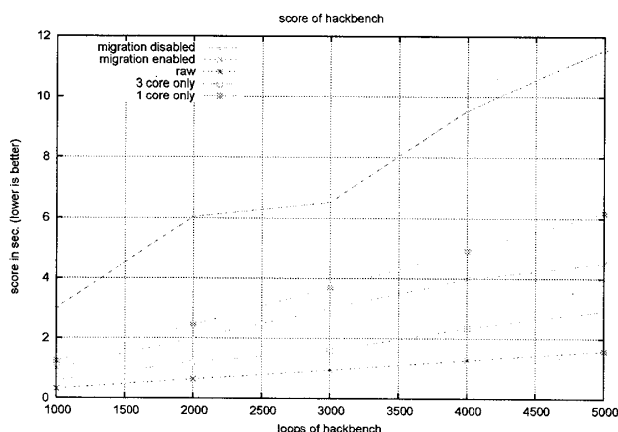


図 1: hackbench のスコア

評価結果を図 1 に示す。図中の migration disabled という折れ線は本研究による変更前のパフォーマンスのスコア

である。migration enabled という折れ線は本研究による変更後のスコアである。raw という折れ線は, RTOS が時間を消費しない状況で, hackbench が 4 コア全てを利用した状況でのスコアである。3 core only という折れ線は, RTOS が時間を消費しない状況で, sched_setaffinity() を利用して hackbench が利用するコアを 3 個に制限した状況でのスコアである。1 core only という折れ線は, 3 core only と同じ条件で hackbench が利用するコアを 1 個に制限した状況でのスコアである。変更前, 変更後では, RTOS が 90% の CPU 時間 (500ms 周期で 450ms) を消費している。

本研究による変更を加える前では, このように RTOS が高負荷な状況になると, Linux のパフォーマンスは 1 コアのみを利用した時よりもさらに低くなる。変更を加えることで, 2 倍以上にスコアは改善されるが, それでも依然として理想である 3.1 コア分のスコアには届いていない。

5 まとめ

本研究では, preempt_enable() と preempt_disable() に変更を加えることで, SPUMONE 利用環境での LHP を可能な限り回避し, Linux の性能の向上を達成した。

ベンチマークによるスコアは向上が見られたが, それでも依然として理想とするスコアとの間には差がある。これは, RTOS のリアルタイム性を保証するため, Linux の全ての仮想 CPU がプリエンプト不可能だった場合には LHP の発生を許可し, RTOS の実行を優先させているからだと考えられる。

参考文献

- [1] Volkmar Uhlig, Joshua LeVasseur, Espen Skoglund, and Uwe Dannowski. Towards Scalable Multiprocessor Virtual Machines. *In Proceedings of the 3rd Virtual Machine Research & Technology Symposium (VM'04)*
- [2] Philip M. Wells, Koushik Chakraborty, and Gurindar S. Sohi. Hardware Support for Spin Management in Overcommitted Virtual Machines. *15th International Conference on Parallel Architectures and Compilation Techniques (PACT '06)*
- [3] Yuki Kinebuchi, Takushi Morita, Kazuo Makijima, Midori Sugaya and Tatsuo Nakajima. Constructing a Multi-OS Platform with Minimal Engineering Cost *In Analysis, Architectures and Modelling of Embedded Systems: Third IFIP TC 10 International Embedded Systems Symposium (IESS 2009)*