

Unikernel クラウド向けの VM グループ管理

2011SE252 田尻 翔太 2011SE269 遠山 恭平

指導教員：宮澤 元

1 はじめに

近年、ネットワーク経由でサービスを提供する、クラウドコンピューティング(クラウド)が普及している。クラウドでは、データセンタに設置された複数のサーバ計算機を多数のユーザが利用する。利用するユーザの多様さから、幅広い要求に対応することが求められる。そのため、CPU やメモリ、ストレージなどの物理的な資源を、その物理的構成にとらわれずに柔軟な形で提供したり、限りあるそれら物理資源を効率的に活用したりする必要がある。そこで、クラウドでは仮想化技術が取り入れられている。

仮想化技術を用いると、物理資源をその実構成にとらわれずに管理することができる。仮想化を実現するには、仮想マシンモニタ(VMM:Virtual Machine Monitor)が利用される。VMM 上では、仮想マシン(VM:Virtual Machine)が動作し、そこで通常の OS やアプリケーションを動かすことができる。

仮想化技術を用いる際の課題として、VM のパフォーマンスの問題がある。例えば、VMM が VM をスケジューリングし、VM がさらに内部で動いているプロセスのスケジューリングを行うといった、スケジューリング処理の多重化によるオーバーヘッドが生じる。また、ディスクやネットワークアクセス時に発生する I/O (Input/Output) エミュレーションでオーバーヘッドが生じる。こうしたオーバーヘッドは、VM のパフォーマンスに悪影響を及ぼす。

VM のパフォーマンスを改善するため、様々な研究が行われており、Unikernel[1] もそのような研究の 1 つである。Unikernel は、VM 上で動作する OS を Library OS を活用して軽量化するとともに、単一の VM で単一のサービスを提供するものである。これによって VM 内でプロセスのスケジューリングを行うために起こる、スケジューリングの多重化を解消できる。同時に、VM 内でスケジューリングを行うプリエンプティブスレッドは提供しない。

一方、大規模な並列処理を行うためにクラウドを用いることがある。例えば、クラウド上の VM で Hadoop[2] を用いて大規模計算を行うといったように、並列処理を利用することは少なくない。Unikernel では、並列処理が必要な場合、複数の VM を用いて VM 間通信を行っている。しかし、VM 間通信を行うとアドレス空間を頻繁にまたぐ必要があるため、余計なオーバーヘッドが生じてしまう。

本稿では、Unikernel のような単一のアプリケーションだけを VM で動作させるようなクラウドにおいて、大規模並列計算を支援するための VM のグループ管理について述べる。並列処理を行っている VM のグループ化を行うことにより、グループ化された VM のスケジューリングやメモリ割り当てを最適化可能である。

2 クラウドにおける仮想化環境

IaaS クラウドを構築するためにさまざまな仮想化環境が提案され、利用されている。本節ではこのような仮想化環境のうち Mirage OS(Mirage)、Xen、およびコンテナ型仮想化について述べる。

2.1 Mirage

Mirage は Unikernel の実装の 1 つで、モバイルプラットフォームや様々なクラウドコンピューティングで動作するセキュアでハイパフォーマンスな OS である [1, 3]。Mirage 向けのアプリケーションは関数型言語 OCaml で開発される。ソフトウェアスタック全体、言語ランタイムおよびシステムライブラリはハイパーバイザ上で直接動作するブート可能な VM イメージにコンパイルされる。図 1 のように Mirage は既存の VM と比較すると様々なソフトウェア層を省略していることがわかる。

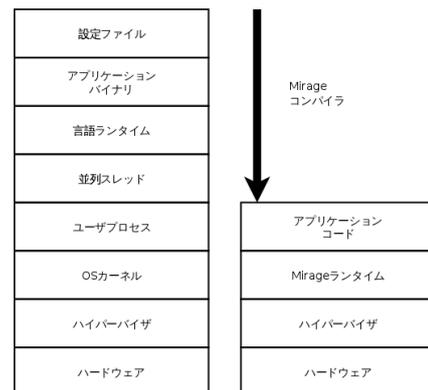


図 1 VM のソフトウェア層 (左:既存の VM. 右:Mirage)

2.2 Xen の概要

Xen[4] はオープンソースの VMM の 1 つである。Xen にはコントロールドメインと呼ばれる Domain-0 と、ゲストドメインと呼ばれる Domain-U の 2 種類の VM が存在する。それぞれが VMM によって仮想化され、資源の提供を受けて動作している(図 2)。

Xen の仮想化の方式には以下の 2 種類の方式がある
完全仮想化 ゲスト OS のコードを修正せず、Intel VT や AMD-V といった CPU の仮想化支援機能を利用し OS を動作させる。

準仮想化 Xen 上で動作させるためにゲスト OS のコードを修正し、ハードウェアエミュレーションのオーバーヘッドを減らす。

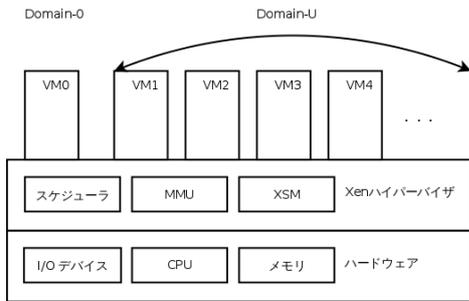


図2 Xenの構成要素

2.2.1 Xenの共有メモリ

Xenが従来持つ共有メモリを提供する機能として Grant Mechanism [5]がある。Grant Mechanismは、VMとVMMで共有される grant referenceを利用してメモリの共有を行う。grant referenceには、共有するメモリページのアドレスや、相手のVMのIDなどの情報が含まれている。相手のVMが共有されたページにアクセスするときに、ページ所有者が変更されることはないので、安全に処理が行われる。Grant Mechanismは、Domain-0とDomain-U間で処理が分かれる場合など、VM間で通信が必要な場合に、処理の効率化のため用いられる。

2.3 コンテナ型仮想化

VMMのオーバーヘッドのないLinux Containers(LXC)[6]やDocker[7]などのコンテナ型仮想化が開発されている。起動する全てのプロセスは物理マシンにインストールされたOS上で動作する。そのプロセスの一部を図3のようにグループ化し、他のグループやグループに属していないプロセスから隔離した空間で動作させる。コンテナ型仮想化ではOSは固定され、プロセスごとに様々なOSを利用することはできない。また、現状のLinux Kernelではプロセスのマイグレーションもサポートされておらず、ハードウェアをまたぐようなコンテナを利用することはできない。

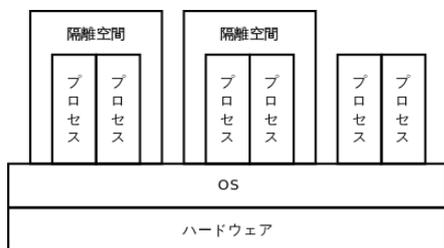


図3 コンテナの仕組み

3 Unikernelクラウド向けのVMグループ管理の設計

我々は、Unikernel型のクラウドのVMMにおいて大規模並列計算のサポートを行うためにVMをグループとし

て管理する機構を導入する。本節ではVMのグループ化、VMグループを考慮したVMスケジューリングとVMへのアドレス空間の割り当てについて述べる。

3.1 VMのグループ化

Unikernel型のクラウドで並列処理を行う場合、それぞれ単一スレッドを動作させる複数のVMが互いに通信を行う。このように関係の強いVMをグループとしてまとめて管理することで、VMのスケジューリングとVMへのアドレス空間の割り当てをVMグループを考慮して行うことができる。

図4の例では、VM0、VM1、VM2の3つのVMが起動しており、このうちVM1とVM2が互いに通信を行っている。VM1とVM2をグループとしてまとめ、1つのVMグループとしてVMのスケジューリングとアドレス空間を割り当てる。

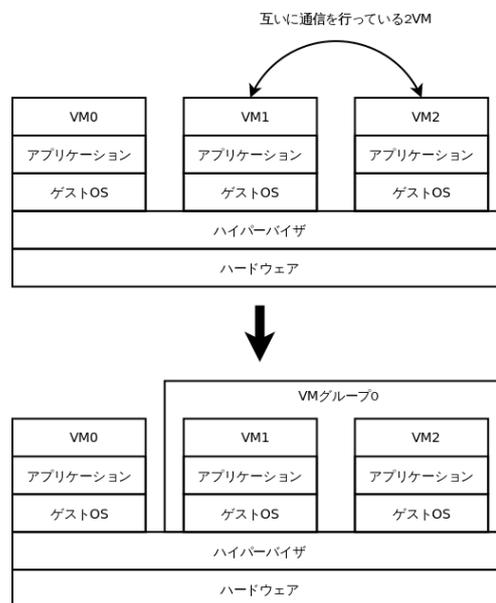


図4 関係の強いVMをグループ化する

3.2 VMグループを考慮したVMスケジューリング

Unikernel型のクラウドでは1つのVMが1スレッドと考えられるが、各VMの仮想CPUに物理CPUがVMスケジューラによってバラバラに割り当てられる。VMをグループ化した場合、図5のように、VMのグループを複数CPUが割り当てられた1つのVMと見なし、VMスケジューリングを行う。

3.3 共有アドレス空間の提供

Unikernel型のクラウドにおける現状のメモリ割り当ては、VMをプロセスのように扱っているとも考えられる。VMはそれぞれ独自のアドレス空間を所有し、その上で処理を行う。この方法では、並列処理を行う際に、アドレス空間の切り替えに伴うオーバーヘッドが生じてしまう。図6は、通常のVMへのメモリ割り当てを表している。各VM

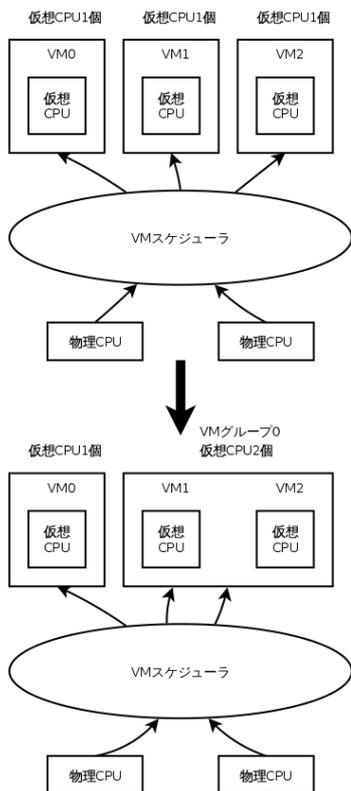


図5 VMグループを考慮したVMスケジューリング

のアドレス空間には、それぞれ別の物理メモリが割り当てられる。XenのGrant Mechanismを使えば、メモリの一部を他のVMに見せることができるが、アドレス空間をまるごと共有することは考えられていない。

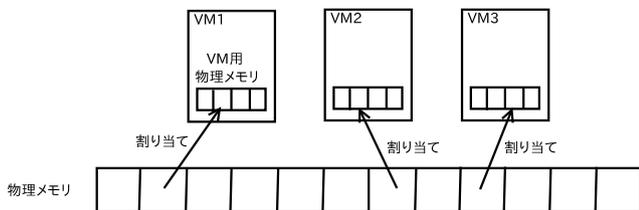


図6 VMグループのない場合のアドレス空間の割り当て

我々は、VMグループを考慮したメモリの割り当てを行う。同一のVMグループに属するVM同士でアドレス空間を共有することにより、各VMをスレッドのように扱うことができる。アドレス空間が共有されるので、アドレス空間の切り替えによるオーバーヘッドを削減できる。図7は、VMグループを考慮した場合のメモリ割り当ての例である。同じグループに属するVMは、アドレス空間が共有される。Unikernelのように単一VMで単一サービスを提供する仕様と合わせると、VMがスレッドのように扱われることになる。

4 VMグループの実装

3.1節のVMのグループ化と3.2節のVMグループを考慮したVMスケジューリングをXenのバージョン4.4.1

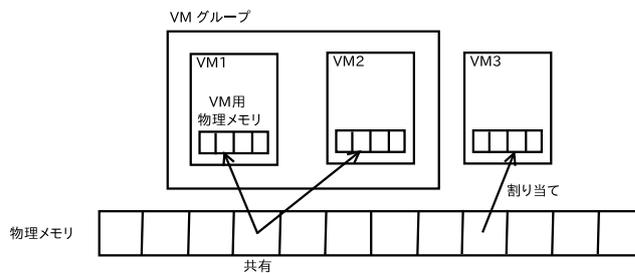


図7 VMグループを考慮したアドレス空間の割り当て

をベースに行った。なお、3.3節の共有アドレス空間は未実装である。

4.1 VMのグループ化

VMのグループ化を実装するにあたり、XenのVMの作成処理とVMの削除処理を拡張した。具体的な実装を以下に示す

- Xenハイパーバイザに新たにVMグループ構造体を追加した。VMグループは複数ありそれぞれVMグループIDを持っている。
- 作成するVMの情報に追加するVMグループのIDを追加し、VMグループIDを指定したVMの作成を可能にした。実装はxl createコマンドとハイパーコールの処理を拡張する形で行った。
- VMの削除処理を拡張し、VMをグループから外す処理を追加した。VMグループにVMが存在しなくなる場合にはそのVMグループを削除する。

4.2 VMグループを考慮したスケジューリング

デフォルトのドメインスケジューラであるCreditスケジューラに変更を加えた。Unikernel型のIaaSクラウドで並列処理を行う場合、各VMが並列処理のスレッドとして動作し、VM同士が互いに通信を行っている。グループ内のVMのどれかが動作した後は、VM間の通信の待ち時間を減らすために、同一グループ内の別のVMが優先して動作するようにスケジューラを拡張した。

5 実験

VMグループを考慮したスケジューラの効果を確認するために実験を行った。具体的には、Hadoopを使ったモンテカルロ法により円周率の近似値を算出するアプリケーション[8]を動作させ、その計算時間を測定する。

5.1 実験環境

表1に示されるPC1台を実験に使用した。並列計算処理を行わせるVMをこのPCの上に3つ作成する。作成した3つのVMをVM1, VM2, VM3とした。作成したVMの仕様を表2に示す。Hadoop環境の構築にはCloudera社のCloudera Manager[9]バージョン5.2.1を使用し、VM1をCloudera Managerサーバ、VM2とVM3を並列計算処理を行わせるホストとした。

表 1 PC の仕様

CPU	Intel®Core™i7-2600
クロック周波数	3.4GHz
メモリ	8GB
HDD 容量	500GB
コア数	4 コア 8 スレッド
OS	Ubuntu Srver 14.04 64bit

表 2 VM の仕様

	VM1	VM2, VM3
CPU コア数	2	
メモリ	2GB	1.5GB
HDD 容量	32GB	
OS	Ubuntu Srver 14.04 64bit	
Java バージョン	1.6.0.31	

5.2 実験結果

計算プログラムのマップ数と 1 マップ当たりのサンプル数をそれぞれ変化させ、計算時間を測定した。それぞれ 10 回ずつ計測を行い計算時間の平均値を求めた。マップ数が 10 の場合の実験結果を図 8 に、100 の場合を図 9 に示す。

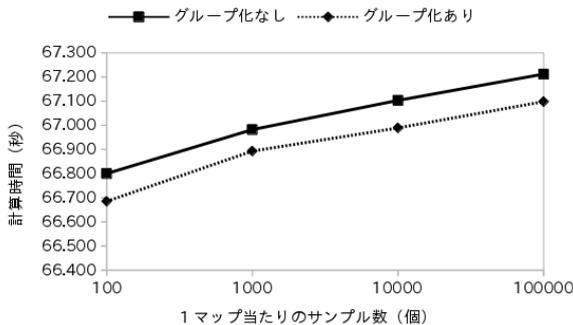


図 8 マップ数 10 の場合の実験結果

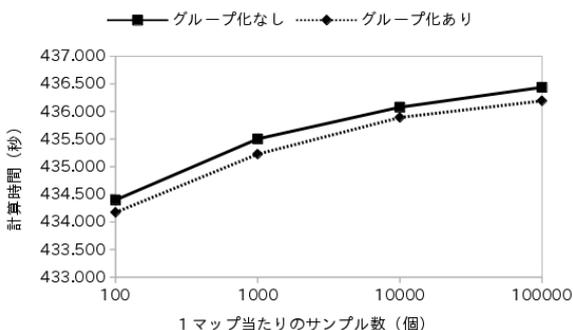


図 9 マップ数 100 の場合の実験結果

マップ数や、1 マップ当たりのサンプル数にかかわらず、3 つの VM が同一のグループに属して処理を行う場合の処理時間が短い。VM グループを考慮した VM スケジューラの効果が期待できることがわかる。

6 まとめと今後の課題

我々は、Unikernel 型の IaaS クラウドにおける大規模並列計算を支援するために、VM のグループ化を提案した。VM をグループ化することで、VM スケジューリングとアドレス空間の割り当てを VM グループを考慮した形で行うことができる。我々が提案する機能を Xen ハイパーバイザに追加することで実装を行った。VM グループを考慮したスケジューラを用いた実験により、スケジューラの効果が確認できた。

今後は、VM グループを考慮した共有アドレス空間を実装すること、実装した機能の有用性を確認するために更に実験を行うことが必要である。また、VM グループの機能をネットワーク環境に対応させるための検討を行う。

参考文献

- [1] A. Madhavapeddy, R. Mortier, C. Rotsos, D. Scott, B. Singh, T. Gazagnaire, S. Smith, S. Hand, and J. Crowcroft, “Unikernels: Library Operating Systems for the Cloud,” *SIGPLAN Not.*, vol. 48, no. 4, pp. 461–472, Mar. 2013. [Online]. Available: <http://doi.acm.org/10.1145/2499368.2451167>
- [2] The Apache Software Foundation, “Welcome to Apache™Hadoop®!” [Online]. Available: <http://hadoop.apache.org/>, 2015 年 1 月 10 日アクセス
- [3] “Mirage OS.” [Online]. Available: <http://openmirage.org/>, 2015 年 1 月 10 日アクセス
- [4] “The Xen Project, the powerful open source industry standard for virtualization.” [Online]. Available: <http://www.xenproject.org/>, 2015 年 1 月 10 日アクセス
- [5] K. K. Ram, J. R. Santos, and Y. Turner, “Redesigning xen’s memory sharing mechanism for safe and efficient I/O virtualization,” in *WIOV’10: Proceedings of the 2nd conference on I/O virtualization*. Berkeley, CA, USA: USENIX, 2010.
- [6] “Linux Containers.” [Online]. Available: <https://linuxcontainers.org/>, 2015 年 1 月 10 日アクセス
- [7] “Docker - Build, Ship, and Run Any App, Anywhere.” [Online]. Available: <https://www.docker.com/>, 2015 年 1 月 10 日アクセス
- [8] *Hadoop 徹底入門 第 2 版*. 翔泳社, 2013. [Online]. Available: <https://books.google.co.jp/books?id=4AgfBAAQBAJ>
- [9] Cloudera, “Cloudera Manager 5.” [Online]. Available: http://www.cloudera.com/content/cloudera/en/documentation/core/v5-2-x/topics/cm_intro_primer.html, 2015 年 1 月 10 日アクセス