

L-014

VM間のトラフィック交流を考慮した仮想サーバの効率的な配置方法の提案 An Effective Arrangement of Virtual Machines Based on traffic matrix between VMs

朝倉 浩志†
Hiroshi ASAKURA

倉上 弘†
Hiroshi KURAKAMI

山田 博司†
Hiroshi YAMADA

1. まえがき

本稿ではデータセンタにおける仮想サーバ(VM; virtual machine)の効率的な配置方法を提案する。

サーバ仮想化技術の進展により仮想化されたサーバやネットワークインフラの提供を行う IaaS 事業が盛んである。クラウドプロバイダーと呼ばれる事業者が、「クラウド環境」を利用者に提供する。クラウド環境の特徴として利用者自身で自由にシステム構築を行うことができること、要求性能の変化に対し柔軟な変更が可能ながある。VMのCPU資源やメモリ容量の割り当てを増減させるスケールアップ/ダウン、VM数を増減させシステム全体としての処理能力を増減させるスケールアウト/インといった手法がある。物理サーバを更改、増減設していた時に比べ、要する時間が格段に減少したため一つの魅力となっている。

このような柔軟な変更要求に対応するためには、サーバ・ネットワーク双方の資源を統合的に管理し、自動的に資源を割り振っていく運用が必要である。

2. データセンタ内ネットワークの構成

以下にクラウドプロバイダーが構築するデータセンタの典型的なネットワーク構成を示す。通常、木構造構成が用いられる。クラウドの規模によって木構造の高さが異なるが、物理サーバを収容する ToR(Top of Rack; 通常サーバラックの上部に配置されることからそう呼ばれる)スイッチ、ToR スwitchを束ねる集約スイッチ、集約スイッチを束ねるコアスイッチ等から構成される。

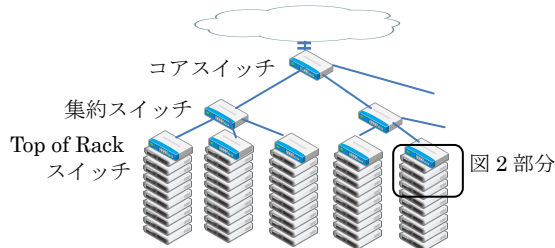


図1. ネットワーク構成例

物理サーバ上では複数のVMが動作し、それぞれ異なる利用者に貸し出されている。異なる利用者でブロードキャストドメインを分けるため tagged VLAN(IEEE802.1Q, 以下VLAN)を利用しネットワークを構成することが多い。また、同一の利用者が複数のL2セグメントを使用してシステム構築を行う場合がある。この場合も同様にVLANを用いて実現されることが多い。

このようなネットワーク構成において、新規契約やスケールアウトを目的としたVMの増設要求があると、特定のルールに従って物理サーバ上にVMが割り当てられる。ま

† 日本電信電話株式会社 NTT 情報流通プラットフォーム研究所, NTT Information Sharing Platform Laboratories.

た、解約やスケールインのための削除要求があると、VMが取り除かれる。

3. 課題

利用者によるVMの増減が繰り返されると、それぞれの物理サーバが提供できるCPU資源やメモリ資源は断片化される。VMに提供できる資源が断片化されると、クラウドプロバイダーが持つ資産を効率的に活用できなくなり問題である。この解決方法としてVMの効率的な配置方法を考える。既存の検討として、単一の物理サーバ資源を、収容したVMに対して等分(fair)に割り当てる方法については Ongaro[2]らが実装評価も含め研究を行っている。本研究は複数の物理サーバとネットワークに資源の対象を広げている点で異なる。また、本研究分野の製品として vmware 社の DRS(Distributed Resource Scheduler)[3]がある。CPU及びメモリの使用率についてルールを設定し、VMの再配置を自動化できるがトラフィック交流は考慮することができない。

このとき、CPU資源やメモリ資源だけに着目し、配置方法を決定していると別の問題が発生する。同じ利用者のVM、即ち通信が発生するVM同士が異なる物理サーバに収容される可能性がある。アプリケーションによっては、トラフィックが異なる物理サーバ間を往復することになる。また、異なる利用者のVM、即ち異なるトラフィック特性を持ったVMが同じ物理サーバに同居する場合がある。こういった状況では最悪の場合、物理サーバのNIC(Network Information Card)の帯域が不足しボトルネックになる場合があり[1]、高品質なクラウドサービスを提供できない。

図2で具体例を示す。典型的な3層構造のWebシステムをクラウド環境で実現している。VMで提供されたロードバランサ(LB)、アプリケーションサーバ(App)とWebサーバ(Web)が異なる物理サーバに収容されている。またLBと同じ物理サーバに他利用者が使う動画配信サーバ(Streaming)が収容されている。

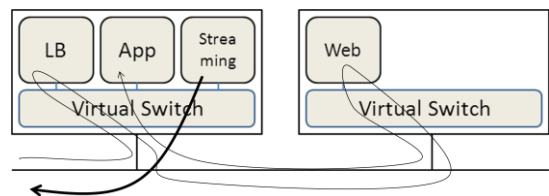


図2. NICでボトルネックが発生する可能性がある例

LBがリクエストを受け付け、Webにリクエストを送出する。WebはAppにリクエストを要求する。レスポンスは逆を辿る。最悪の場合、リクエストだけで物理サーバにあるNICを3回経由する。現在、クラウド環境で使われる物理サーバのNICは、コスト面からGbEが中心であり実測値の上限は800~900Mbpsである。リクエストが物理サーバ間を往復することで使用できる帯域は減少する。Streamingが

使う帯域を考慮すると利用可能帯域はさらに下がる。NICの性能を上げて対処することは可能であるが、問題の本質はVM間トラフィック交流を考慮せずVMの割り当てを決めているところにある。このため、トラフィックを考慮したVMの効率的な配置方法を課題とする。

4. 提案するアルゴリズム

VM間トラフィック交流に着目し、トラフィック量の大きいVMのペアを優先的に同じ物理サーバに割り当てて行くことで課題を解決するアルゴリズムを提案する。

N 個のVMが存在するとする。

b_{ij} = Traffic volume from v_i to v_j , $t_{ij} = (\max(b_{ij}), v_i, v_j)$

・基本アルゴリズム

入力: トラフィックの集合 $T = \{t_{ij}\}$ ($i = 1, \dots, N, j = i, \dots, N$)

出力: アサインされた物理サーバの集合 S_{new}

手順0: VM間のトラフィック交流 T について各要素 t_{ij} の $\max(b_{ij})$ の降順でソートしておく。

手順1: 手順0でソートされた順に、 t_{ij} 各々のトラフィック $t_{ij} = (\max(b_{ij}), v_i, v_j)$ について a), b) の条件で評価し、2つのVM, v_i と v_j を新しい物理サーバ集合 S_{new} のいずれかに割り当てる(図3)。

- v_i が新しい物理サーバに割り当てられていない場合、未使用の物理サーバ $s \in S_{new}$ に割り当てる。
- v_j が新しい物理サーバに割り当てられていない場合、 v_i が割り当てられている物理サーバ s_j に割り当てる。不可能な場合は、未使用の物理サーバ $s \in S_{new}$ に割り当てる。

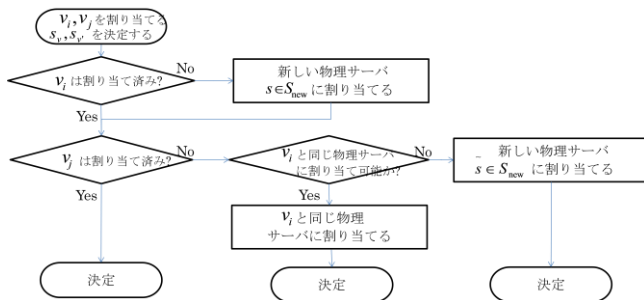


図3. 基本アルゴリズム手順1

ここで割り当て可能かどうかについては以下のようなアルゴリズムで判断する。

・割当判断アルゴリズム

与えられたVMが要求するCPUを c_{new} , メモリを m_{new} , 割当たときに物理サーバ外へ出るトラフィックを \tilde{t}_{new} とする。また、割り当て対象となる物理サーバが提供できるCPU資源を C , メモリ資源を M , 物理NICの処理能力を \tilde{T} , 収容している仮想サーバが各々使用しているCPUを c_i , メモリを m_i , 物理サーバ外へのトラフィック量を \tilde{t}_i とする。

一定の資源 (α および β) を残し, c_{new} , m_{new} を C , M から提供できる場合、収容できると判断する。

$$\begin{aligned} \sum m_i + m_{new} &< M - \alpha \\ \sum c_i + c_{new} &< C - \beta \\ \sum \tilde{t}_i + \tilde{t}_{new} &< \tilde{T} * \gamma \end{aligned}$$

($\alpha, \beta > 0$ の場合は資源留保ありの収容判断, γ はトラフィック量の総和/uplink = oversubscription の係数)

α および β は定数や、収容を判断するVMのスケールアップ実績に応じた値を変数としてとる場合が考えられる。

5. 実験

$M=2000$ 台の仮想サーバを仮定し、複数台の物理サーバ(10CPU, 32GBメモリ)に割り振る。それぞれの仮想サーバが要求するCPU量(最大3CPU), メモリ量(最大10GB)をランダムに割り振る。トラフィック1VMあたり最大5交流あるとし、交流一つあたり最大1Gbpsの範囲でランダムに割り振った。結果は、乱数発生種の種を変えつつ10回行いその平均を示す。

上記の条件で(1)トラフィックを考慮せずVMが要求するCPUとメモリのみを満たせるよう物理サーバに順番に割り当てた場合、(2)提案アルゴリズム(内部留保なし; $\alpha = \beta = 0, \gamma = 2.4$), を比較した。結果は表1の通り。サーバ外通信は物理サーバの外に流れるトラフィック量の総和である。

	物理サーバ	サーバ外通信/台	サーバ外通信
(1)	483.2 台	7.03Gbps	3089.2Gbps
(2)	1592.7 台	1.54Gbps	2233.4Gbps

表1. 実験結果の比較

提案アルゴリズムは2000台のVMを約20%集約し、物理NICのボトルネックになる可能性のある物理サーバ外通信は平均1.54Gbpsに抑えられた。本値は各VMの最大値の合計で計算しており、GbE NICを装備した物理サーバと仮定すると、約1.5:1程度のoversubscriptionである。これに対し、従来の方法である(1)の場合、約7:1にもなるため、ボトルネックが発生しサービスレスポンスへの悪影響を及ぼす可能性がある。また物理サーバ外に抜けるトラフィック量については(1)の手法と比較し、約27%削減できた。提案手法では、ネットワーク全体のトラフィックを削減できることも分かった。一方、提案アルゴリズムでは一つのVMしか収容していない物理サーバが1068.6台あることも分かった。一層の物理サーバ台数削減が可能であると考えられる。

7. まとめ

本稿ではサーバ仮想化が利用者にもたらす特徴的な効果を述べ、それにより発生する物理サーバ資源の断片化という問題を挙げた。解決のため、VMの効率的な配置方法を提案した。このとき、VM間のトラフィック交流を考慮したことが従来と異なる。実験評価を行った結果、CPU及びメモリ資源のみを考慮したときに比べ、必要な物理サーバは増えるが、NICがボトルネックになる可能性は少なくなった。また、必要な物理サーバについては、まだ削減可能であると考えられる。

今回はVM間の片方向トラフィックを評価に使ったが、NICは通常全二重で使われるため双方向のトラフィックを対象とすれば高効率な配置ができる可能性がある。また、得られた最適配置を実現するためにマイグレーションが必要となるが、VMを入れ替える順番や移動回数も考える必要がある。これらを考慮した、収容アルゴリズムへ拡張を図る予定である。

[1] “すべてわかる仮想化大全2011”, 日経BP, 2011

[2] Diego Ongaro, Alan L. Cox and Scott Rixner, Scheduling I/O in Virtual Machine Monitors, In Proceedings of the ACM SIGPLAN/SIGOPS International Conference on Virtual Execution environments, 2008

[3] DRSの機能, 仮想インフラストラクチャ内の仮想マシンのリソース, vmware社, available at http://www.vmware.com/jp/products/vi/vc/drs_features.html