

コンテナクラスタ用ポータブルロードバランサの検討

高橋公俊¹ 合田憲人^{2,1}

1. はじめに

Linux コンテナのクラスタからなる WEB システムが注目を集めている。Linux コンテナは可搬性が高く、BCP(Business continuity planning)や DR(Disaster Recovery)の目的で WEB システムを異なるクラウドプロバイダのクラウド基盤間でマイグレーションすることが容易であるという利点がある。

しかし既存のコンテナ管理システム Kubernetes[1]を調査したところ、異なるクラウドプロバイダ間でマイグレーションを実現するためには、ロードバランサの仕組みに課題があることがわかった。具体的には、Kubernetes は外部のロードバランサに依存しており、コンテナクラスタの起動と同時にクラウドプロバイダが提供する API を通してロードバランサのセットアップを行う。ところが、ロードバランサの仕組みや API はクラウドプロバイダによって異なり、今の所一部の主要なクラウドプロバイダでしかロードバランサを使うことができない。また、オンプレミスデータセンター用のロードバランサ製品は多岐に渡り、その多くを Kubernetes から API で制御することは非常に困難である。本稿ではこの課題を解決するためにロードバランサ自体をコンテナとして実装する方法について検討する。

2. Kubernetes におけるロードバランサの課題

Fig.1 に Kubernetes システムの構成図を示す。Kubernetes システムは Master と複数の Node で構成される。Master が kubectl コマンドによるリクエストを受けると、各ノードで Pod が作成される。Pod は IP アドレスなどのリソースを共有する複数のコンテナをまとめたもので、Kubernetes では Pod 単位でコンテナが実行される。

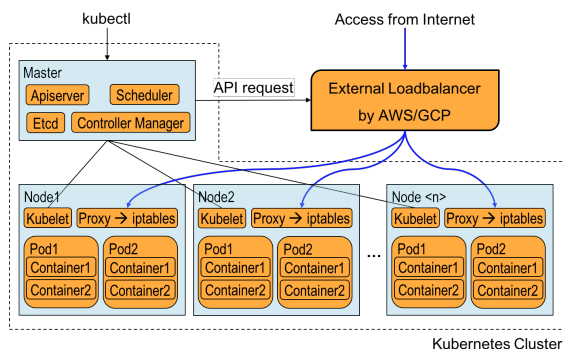


Fig.1 Kubernetes システム構成図

Kubernetes でのロードバランサのセットアップは、Master からの API リクエストを通して行われる。セットアップされたロードバランサはインターネットからのアクセスを全ての Node に均等に振り分け、各 Node でそれらのアクセスをさらに適切な Pod に DNAT[2]を用いて接続している。以上が Kubernetes における外部ロードバランサの利用形態であるが、この方法には次のような問題がある。(1)ロードバランサの仕組みや API はクラウドプロバイダによって異なるため、Kubernetes にサポートされた非常に限られたクラウドプロバイダ以外はこの

仕組みを使うことができない。(2)オンプレミスデータセンターで使われるロードバランサ製品は多岐に渡り、そもそもそれらの多くをサポートすることは非常に困難である。(3)Pod への DNAT は iptables を利用しているため、同時にデプロイされる WEB システムの数が増えると、設定が複雑になりデバッグが困難になる。(4)一旦全ての Node にアクセスを均等に分散し、そこから目的の Pod に DNAT を用いて接続しているため、アクセス経路が複雑になり障害ポイントの把握が困難である。

3. 本研究での提案手法

Fig.2 に本稿が提案するロードバランサを用いた場合のシステム構成を示す。

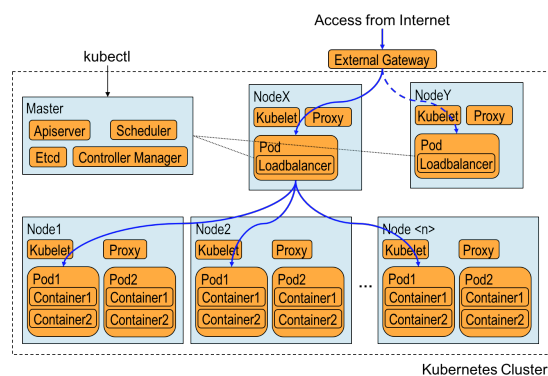


Fig.2 ポータブルロードバランサシステム構成図

提案のロードバランサは次のような特徴を持っている。(1)ロードバランサ自体が Kubernetes の Pod として起動される。(2)ロードバランサ用コンテナは冗長化されており、Active-Backup 構成をとる。(3)負荷分散先の Pod の情報を Master から取得し、動的に分散先に加えたり分散先から外したりすることができる。

このような構成をとることで、既存の Kubernetes システムでのロードバランサの問題を解決することができる。すなわち、ロードバランサ自体をコンテナ化しているため、Kubernetes による API サポートの如何に関わらず、オンプレミス環境でもクラウドプロバイダでも同じ仕組みのロードバランサを利用することができる。また iptables DNAT に依存しないためデバッグの複雑さが無く、外部からのアクセス経路も単純で、通信障害発生時にも、障害箇所を容易に特定することができる。

4. 実装検討

今回我々は、提案方式の実現可能性を明らかにするために、オンプレミス KVM 仮想マシン[3]と AWS EC2 上に Fig. 2 の Kubernetes システムを構築し、ロードバランサ実装方法について検討を行った。

4.1 KVM 仮想環境

Fig.3 に KVM 環境でのロードバランサが稼働する Node のネットワーク構成図を示す。

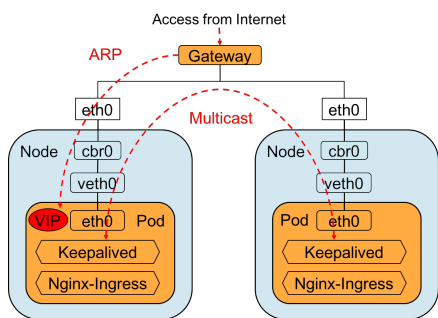


Fig.3 KVM 環境ロードバランサネットワーク構成

外部からゲートウェイに到達したVIP宛パケットは、MACアドレスをARP解決したのちに、イーサネットフレームとしてロードバランサ用Podに送られる。ロードバランサ用PodではNginx-Ingress[4]コンテナが稼働しており、そこから複数のPodにHTTPパケットを送信することで負荷分散を実現している。ロードバランサ用Podでは、keepalived[5]も稼働しており、他のロードバランサPodとの間でVRRPプロトコル[6]によりVIPを共有している。この構成では、ロードバランサ用Podとゲートウェイとの間でARP通信ができることが必要で、また他のPodとマルチキャスト通信ができることも必要である。従って、Fig. 3中のcbr0とeth0はブリッジ接続される必要がありPodのIPアドレスはゲートウェイと同じブロードキャストドメインになければならない。

KVM環境では、このような構成をとることで、負荷分散機能、フェイルオーバーとも問題なく動作した。

4.2 AWS EC2

Fig.4にAWS EC2上でのロードバランサの構成図を示す。外部からパブリックな固定IP(AWSではEIP[7]と呼ばれる)に向けて送られたIPパケットは、ゲートウェイ上で任意のプライベートIP(VIPとする)宛のパケットにDNATを用いて変換される。AWS EC2環境でEIP宛のパケットのDNAT変換先はAPIを通して変更可能である。またVIPのアドレスアサインも同様に変更可能である。AWS EC2環境ではARPが使えないが以上の仕組みを利用して、外部からのIPパケットをPodへ到達させることが可能である。

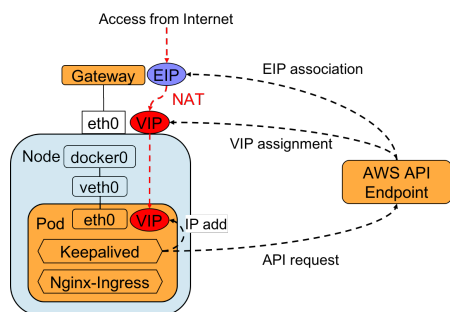


Fig.4 AWS EC2 環境ロードバランサネットワーク構成

ロードバランサの冗長化に関しては、AWS EC2上ではマルチキャストを使うことができないのでVRRPを用いることができない。VRRPにはユニキャストモードもあるが、あらかじめ相手のIPアドレスが決まっている必要があり、Kubernetesコンテナ環境ではIPを指定してPod起動ができないので利用

不可能である。従って、現状で冗長化は難しく、さらなる検討が必要である。

KubernetesではPod起動時に動的にIPアドレスが付与されるが、AWS EC2環境上では、APIやWEBコンソールを介さず付与されたIPアドレスを用いると通信が制限されてしまう。そのためPod間の通信にはトンネリング型のオーバーレイネットワークが必要で、今回はflannel(udp,vxlan)[8]を利用した。

以上のような構成によりロードバランサの負荷分散機能の動作は確認できた。しかしながら、上述のようにフェイルオーバーについてはさらなる検討が必要であることがわかった。

5. 結論

本稿では、異なるクラウドプロバイダ間でKubernetesシステムをマイグレーションするためのロードバランサの実装方法を提案した。提案方式の有効性をオンプレミスのKVM環境およびAWS EC2上で検証したところ、KVM環境においては負荷分散機能、フェイルオーバーとも問題なく動作した。AWS EC2においては、負荷分散機能は問題なかったが、フェイルオーバーについてはさらなる検討が必要であることが明らかになった。本問題を解決するために、今後、冗長化のメンバを知るためのメンバシップとメンバ間でのリーダーエレクトションが可能な別のプロトコルについて検討する予定である。

6. 謝辞

本研究の一部は、JSPS 科研費 24240006 の助成を受けている。

参考文献

- [1] Marmol, Victor, Rohit Jnagal, and Tim Hockin. "Networking in Containers and Container Clusters." *Proceedings of netdev 0.1, February* (2015).
- [2] Linux-ip.net. (2016). *5.5.Destination NAT with netfilter (DNAT)*. [online] Available at: <http://linux-ip.net/html/nat-dnat.html> [Accessed 15 Nov. 2016].
- [3] Kivity, Avi, et al. "kvm: the Linux virtual machine monitor." *Proceedings of the Linux symposium*. Vol. 1. 2007.
- [4] GitHub. (2016). *nginxinc/kubernetes-ingress*. [online] Available at: <https://github.com/nginxinc/kubernetes-ingress> [Accessed 14 Nov. 2016].
- [5] Keepalived.org. (2016). *Keepalived for Linux*. [online] Available at: <http://www.keepalived.org/> [Accessed 14 Nov. 2016].
- [6] RFC5798 | Hinden, R., Ed., "Virtual Router Redundancy Protocol (VRRP)", RFC 3768, DOI 10.17487/RFC3768, April 2004, <<http://www.rfc-editor.org/info/rfc3768>>.
- [7] Docs.aws.amazon.com. (2016). *Elastic IP Addresses - Amazon Elastic Compute Cloud*. [online] Available at: <http://docs.aws.amazon.com/AWSEC2/latest/UserGuide/elastic-ip-addresses-eip.html> [Accessed 15 Nov. 2016].
- [8] GitHub. (2016). *coreos/flannel*. [online] Available at: <https://github.com/coreos/flannel> [Accessed 14 Nov. 2016].