

クラウド事業者から見た課題と今後 ～さくらのクラウド編～

さくらインターネット研究所

大久保 修一

ohkubo@sakura.ad.jp

さくらのクラウドとは？

IaaSの基本的なリソースを提供



サーバ

- 1コア/1GB～12コア/128GB
- 全42種類
- UNIX系OS各種、Windows



ネットワーク

- 共有グローバルセグメント
- 専用グローバルセグメント
- スタティックルート
- ローカルスイッチ
- ロードバランサ
- パケットフィルタ



ストレージ

- SSD 20GB, 100GB
- HDD 40GB～4TB
- アーカイブ、ISOイメージ領域

全てAPIで自由に操作可能

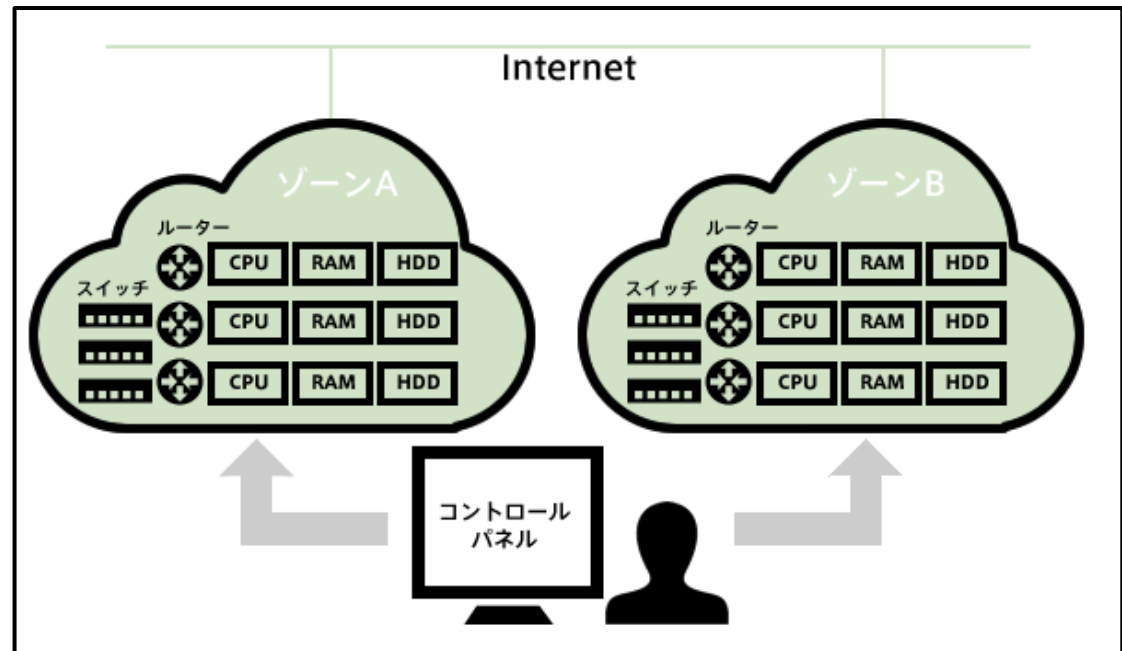
これらの組み合わせで
Software-Defined Data Centerを実現！

最近のアップデート

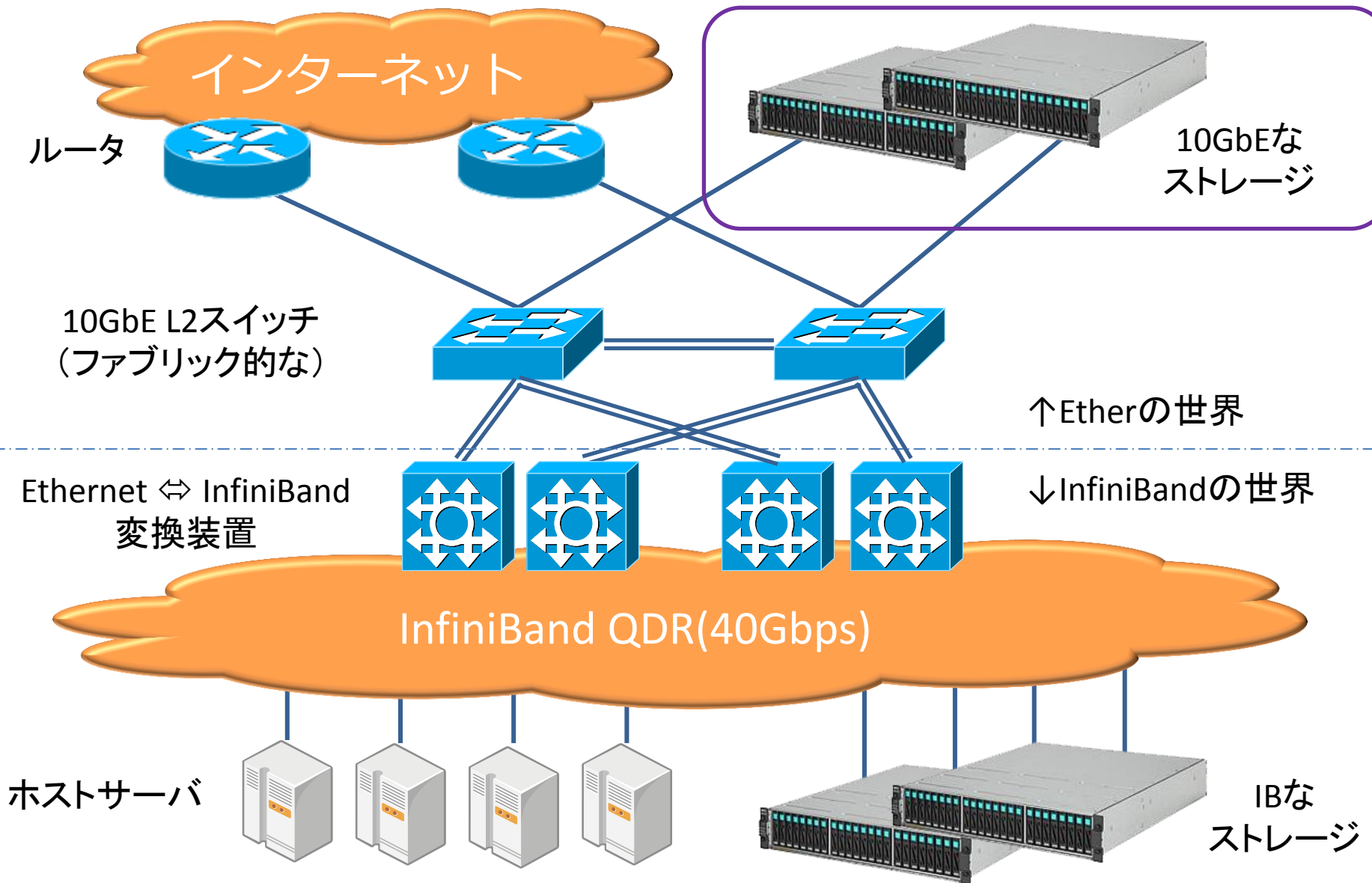
2013/10/8より
石狩第2ゾーン提供開始

ゾーンとは？

- AWSさんのアベイラビリティゾーン的な
- システムを完全に分離、障害が波及しない
- APIサーバ(クラウドコントローラ)とコンパネサーバも分離
- GSLB等と
組み合わせた
冗長システムの
構築が可能



以前のネットワーク(第1ゾーン)



新規のネットワーク(第2ゾーン)

インターネット

10GbEのみでシンプルに！
以前IBで組んでいたストレージは
10GbEで組み直した。

ルータ

多数のVLANを
あらかじめ設定

10GbE L2スイッチ

ToR

ToR

ToR



ホストサーバ

1ラック30台ずつ



ホストサーバ

1ラック30台ずつ

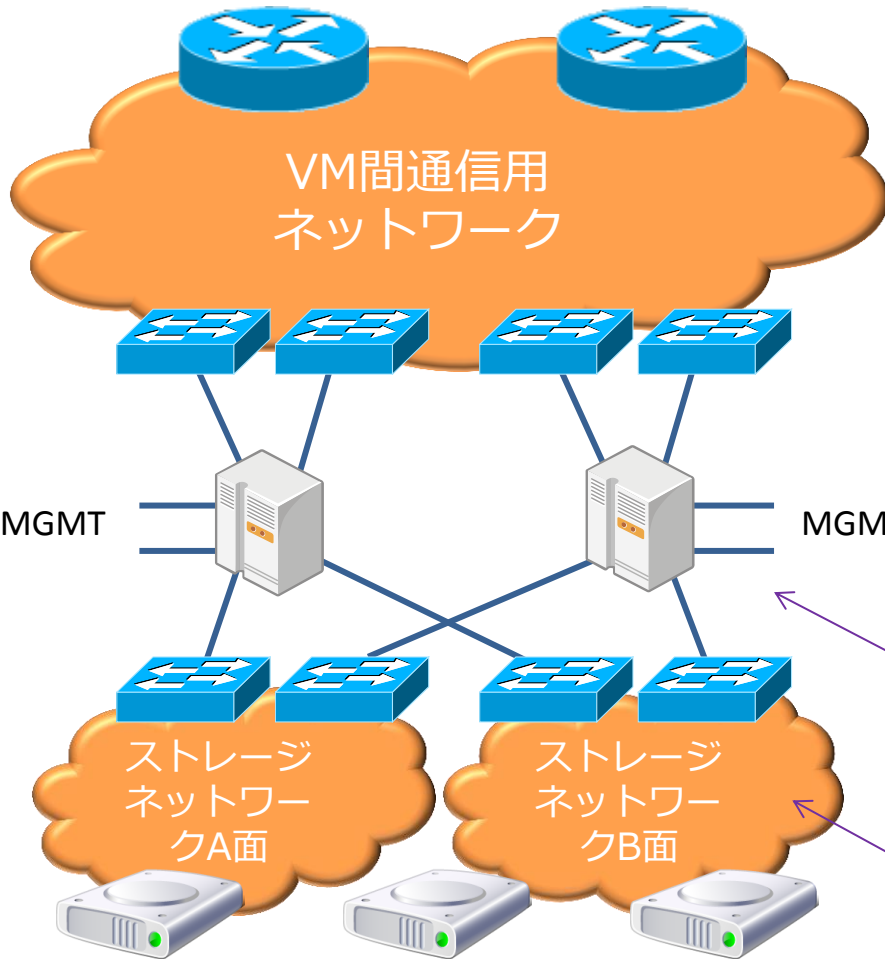


10GbEなストレージ

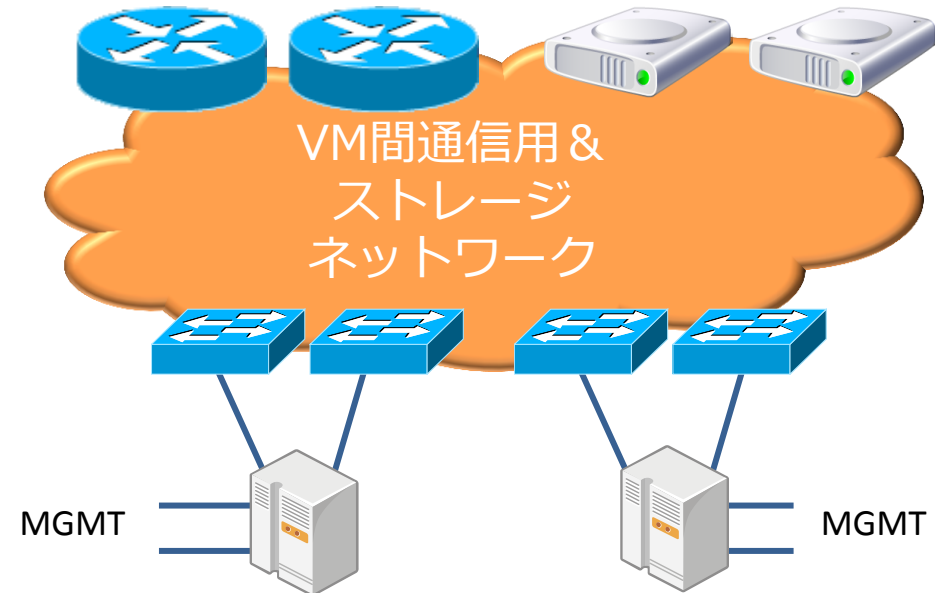
課題1 : ストレージネットワーク

ネットワークの統合

普通に安全に組むなら...



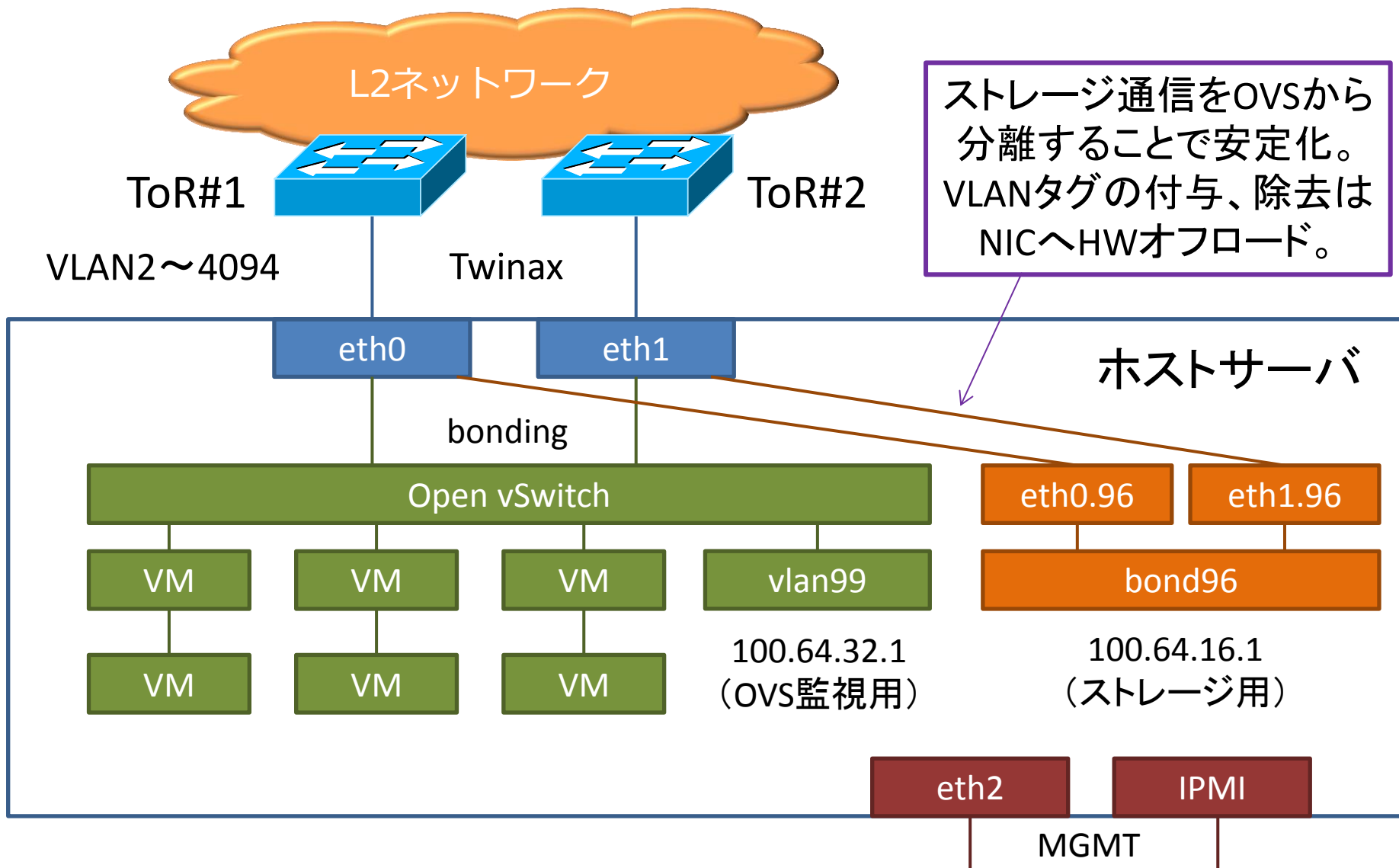
できれば一緒にしたい



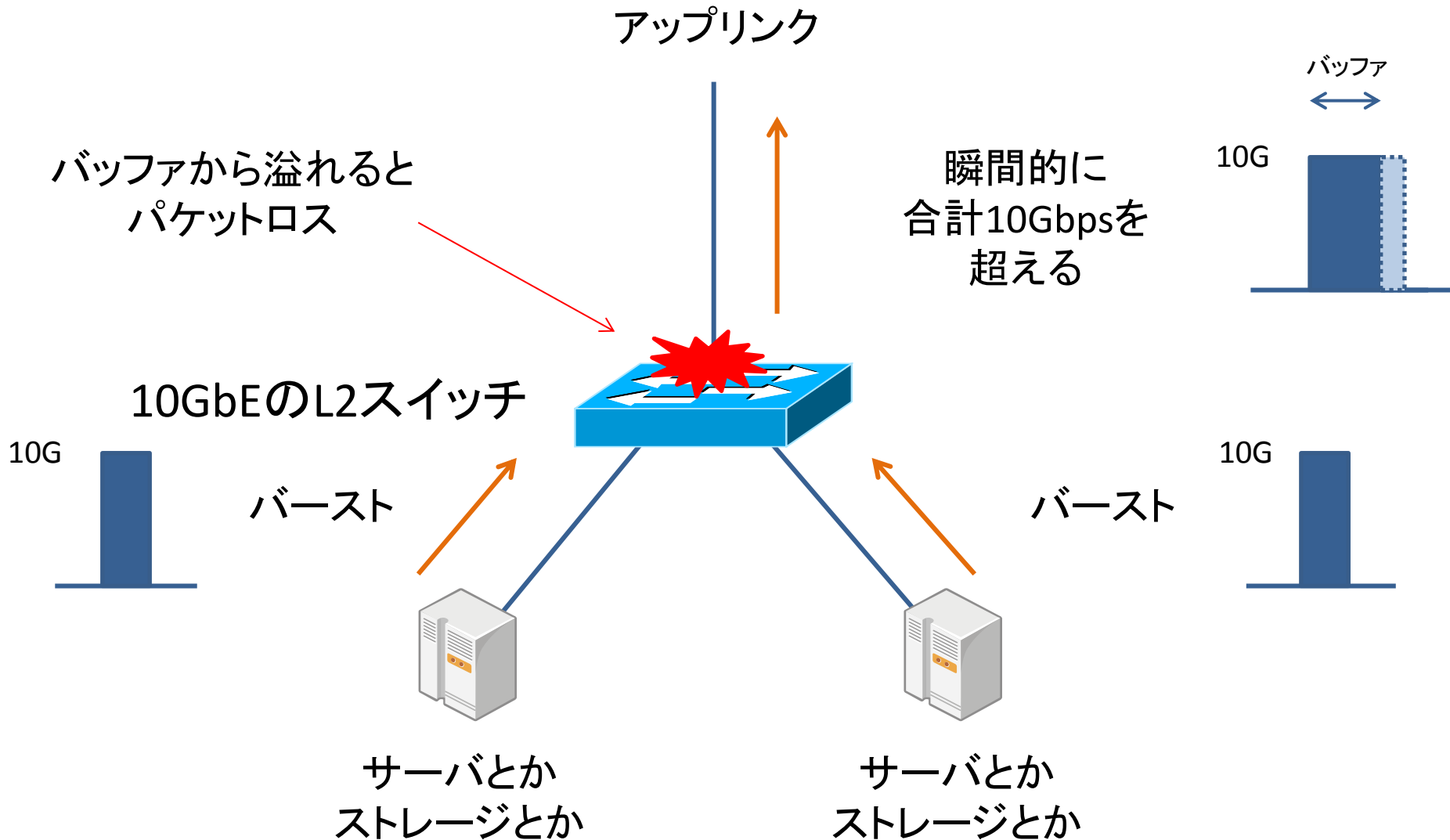
配線もスイッチも
たくさん必要

文化の違う
ネットワーク

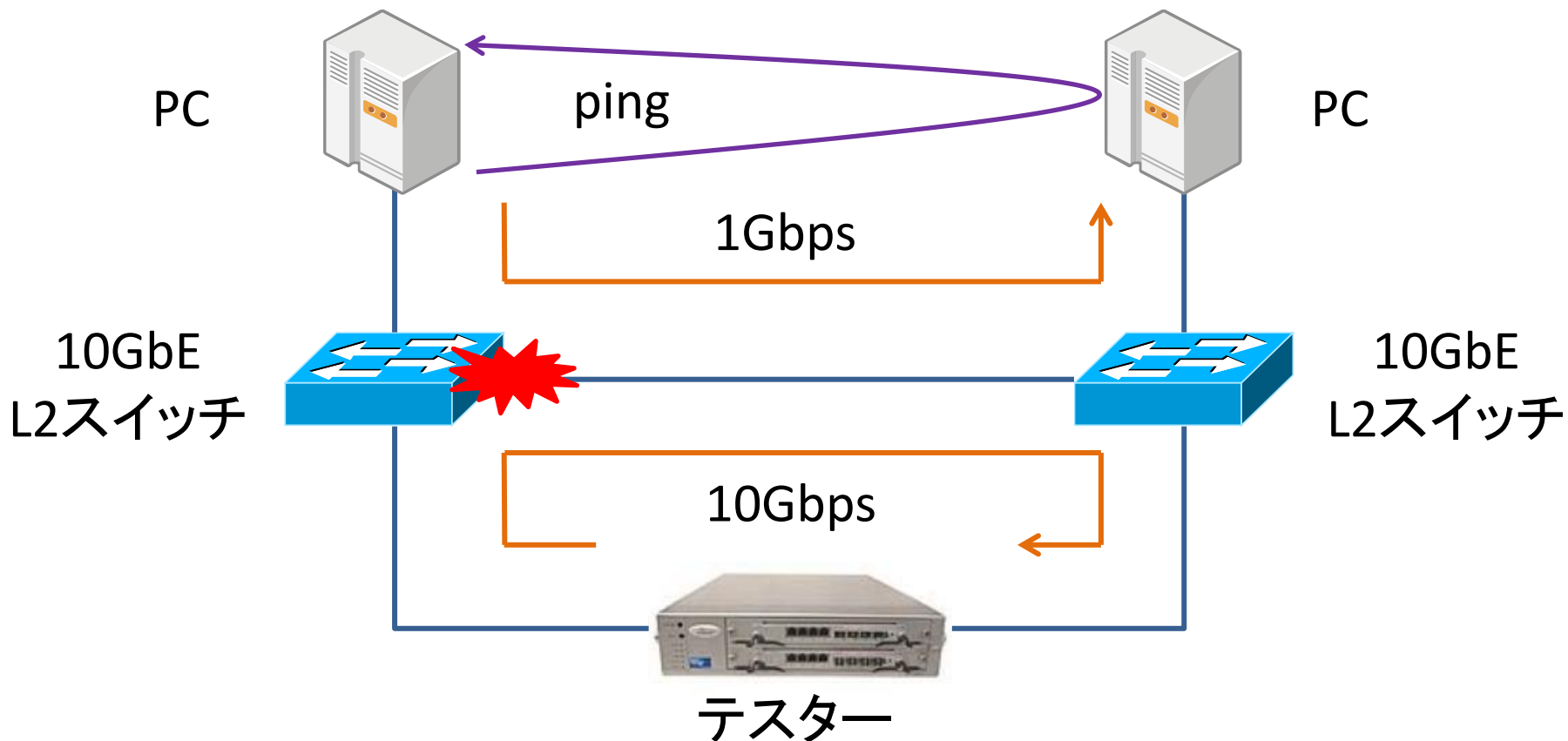
HV内でのストレージ通信の分離



バーストによるパケットロスのケア

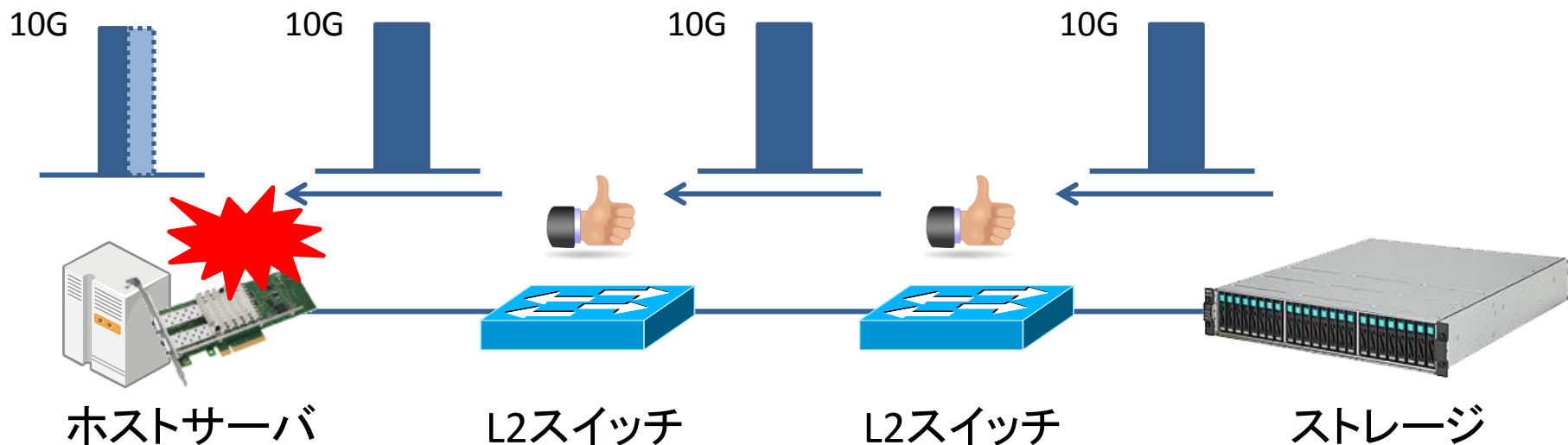


Deep BufferなL2スイッチを導入



ポートを溢れさせ、何 μ secの packets を貯められるかを ping のレイテンシで測定し、バッファサイズを逆算する

ホストの10GbE NICも・・・



現状、iSCSIイニシエータのWindowサイズで調整

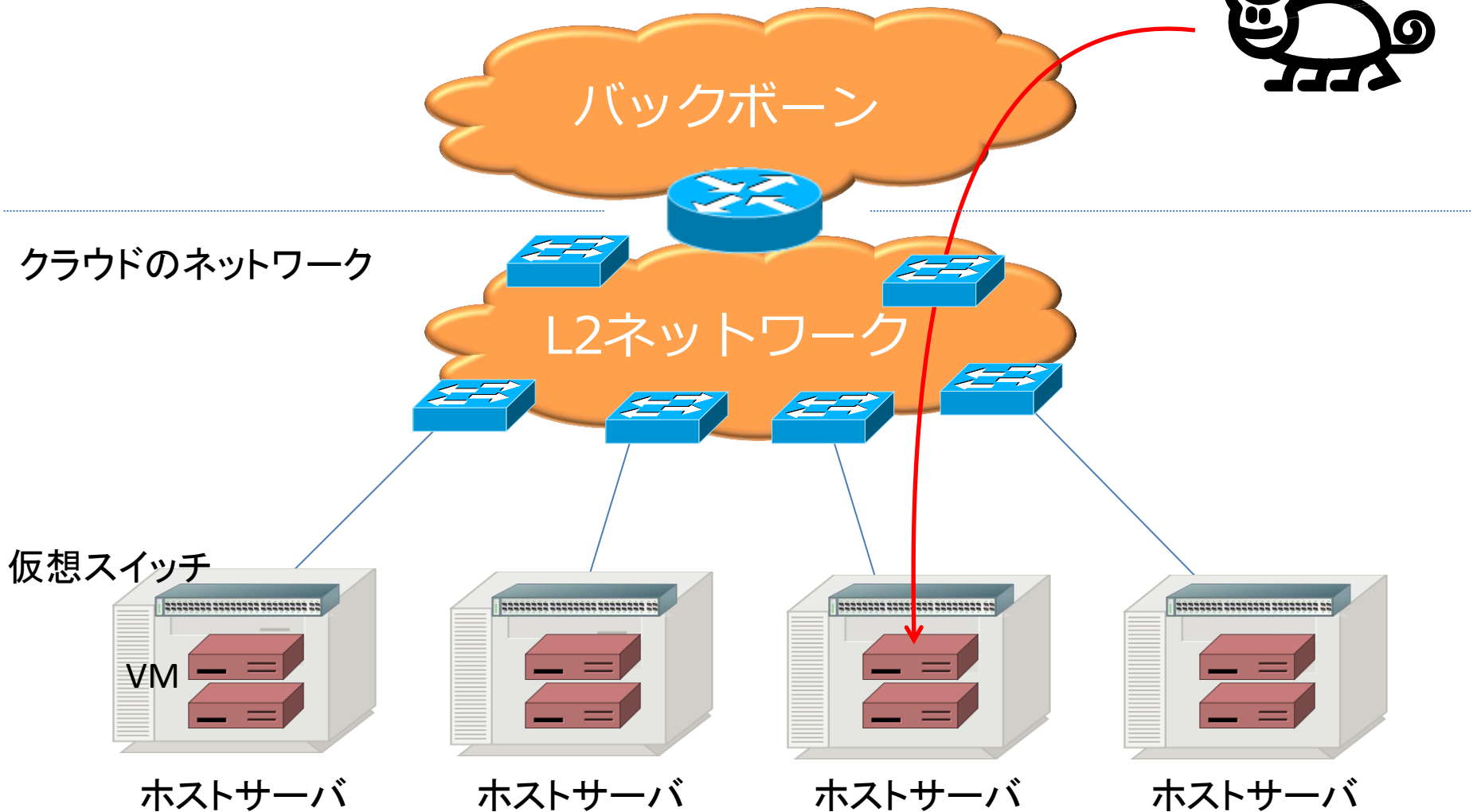
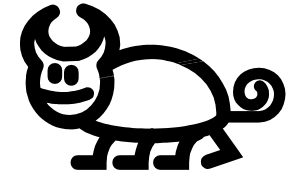
```
# iscsiadm -m node -o update ¥  
-n node.conn[0].tcp.window_size -v 32768
```

ストレージの課題について

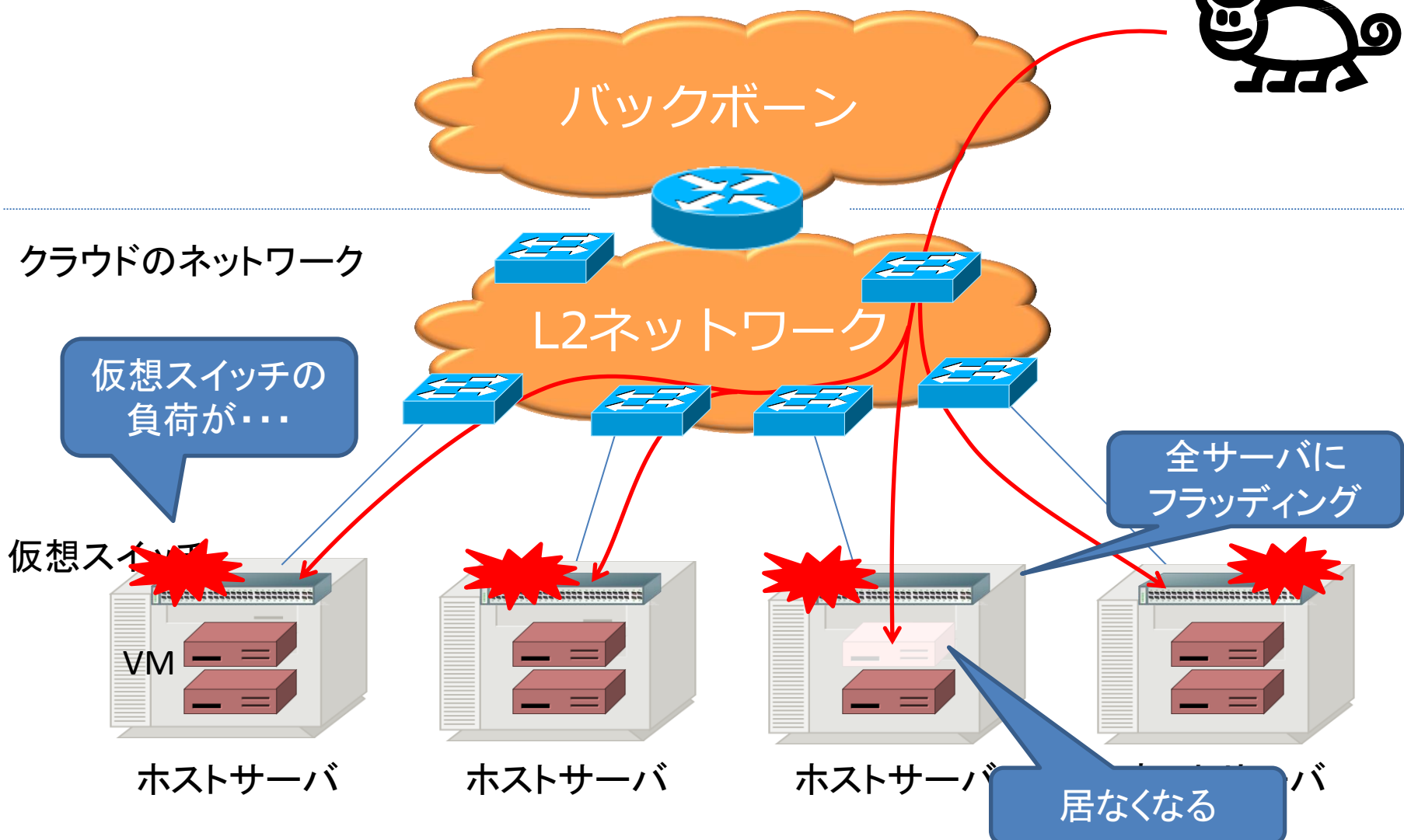
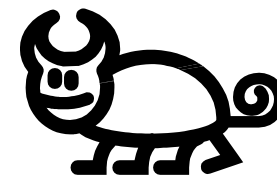
1. ストレージとVM間通信（インターネット通信）を1面のネットワークで干渉しないようにしたい
2. バーストに対して、バッファで吸収するのも限界
 - End to Endのフロー制御機構が欲しい
 - DCB / FCoE が一つの解か？
 - ファブリックに接続する全機器の対応が必要
 - NIC(CNA)、L2スイッチ、ストレージ(ある?)
 - ルータ、Firewall、LoadBalancer(ない?)

課題2 : DoSアタック対策

ある日発生した事象

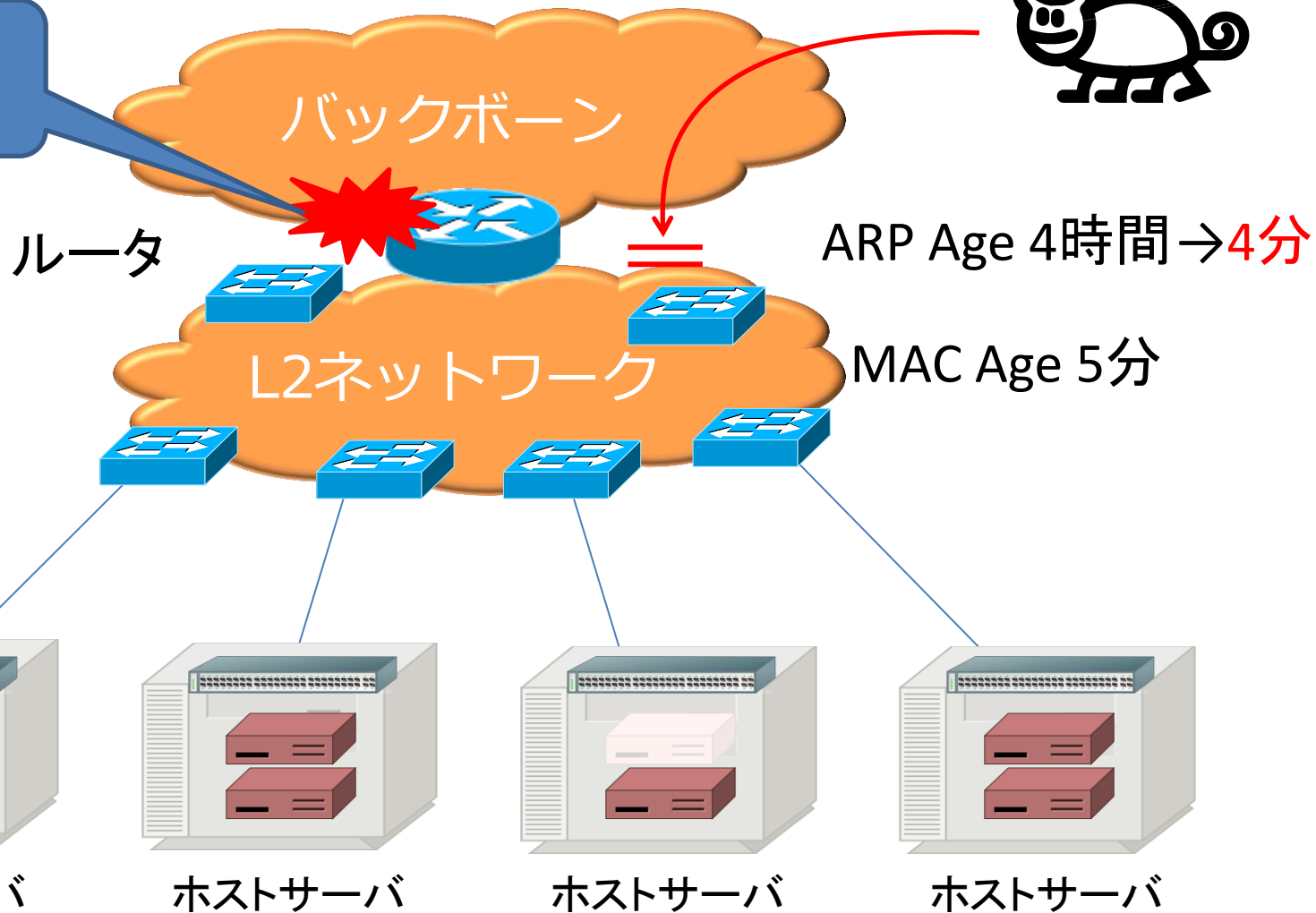
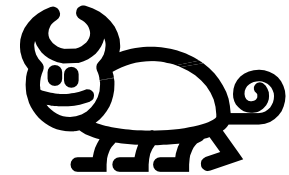


ある日発生した事象



ルータのARPタイムを調整

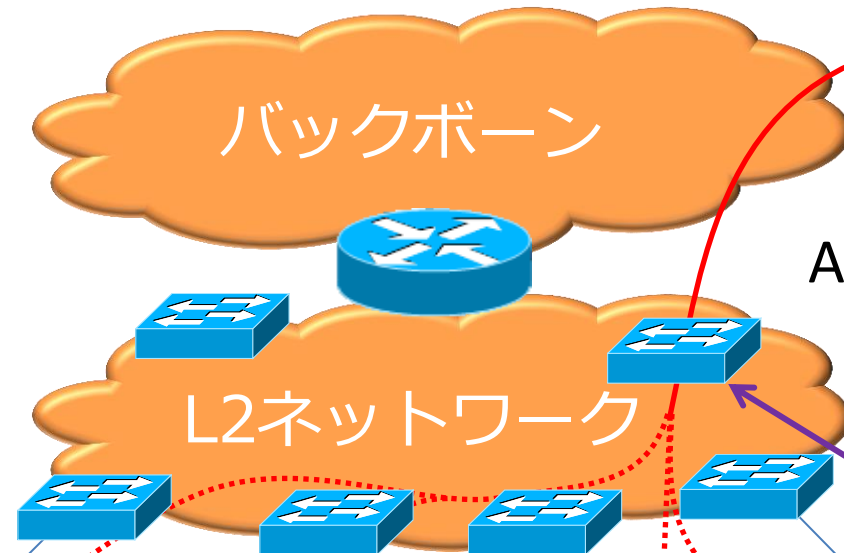
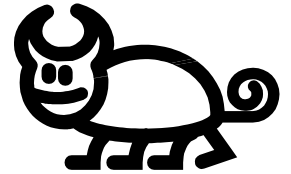
ルータの
CPU負荷が...



解決策は？

1. 超高速にARPのリフレッシュができるルータ
2. L2ネットワークのAgingタイムを延長
 - VMのマイグレーション時や
 - bonding failoverの切り替わりに懸念が...
3. unknown unicastのフラッディング制限を行う

結局どうしたか？

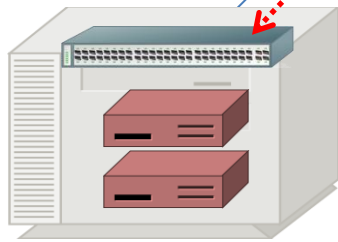


ARP Age 4分 → 4時間

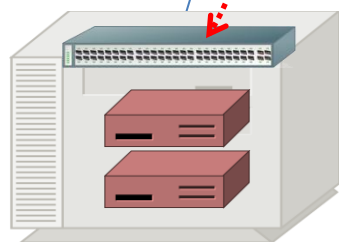
MAC Age 5分

フラッディングするけど、問題ないレベル

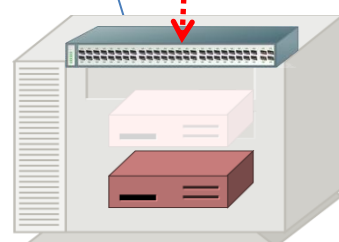
unknown unicastのフラッディング制限



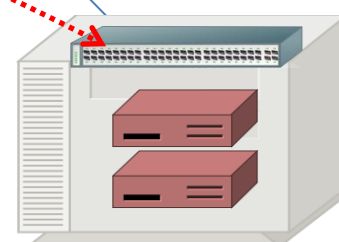
ホストサーバ



ホストサーバ



ホストサーバ



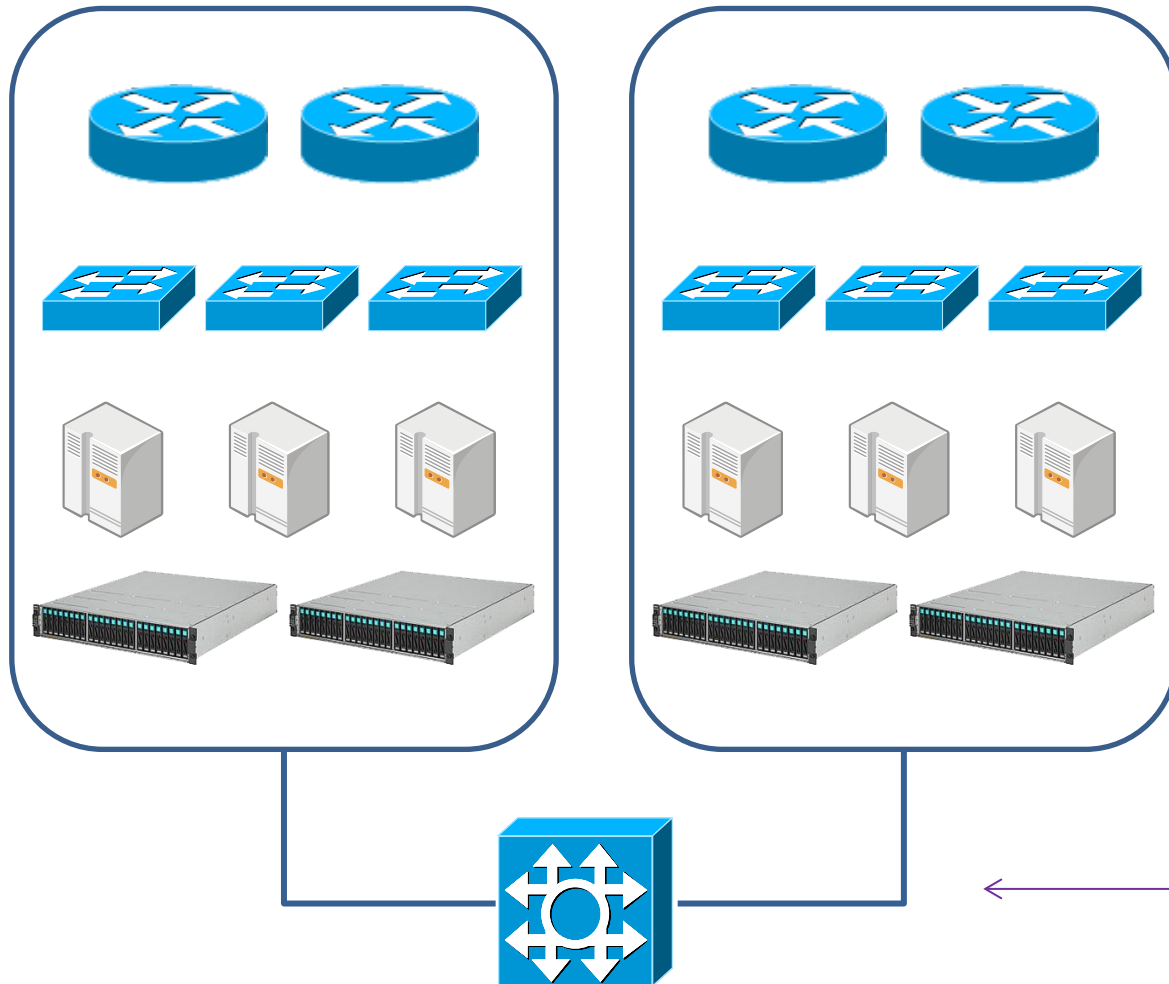
ホストサーバ

課題3 : L2網のスケーラビリティ

ゾーン単位の分割とゾーン間接続

第1ゾーン

第2ゾーン



VLAN数、MAC数の
上限でシステムを分割

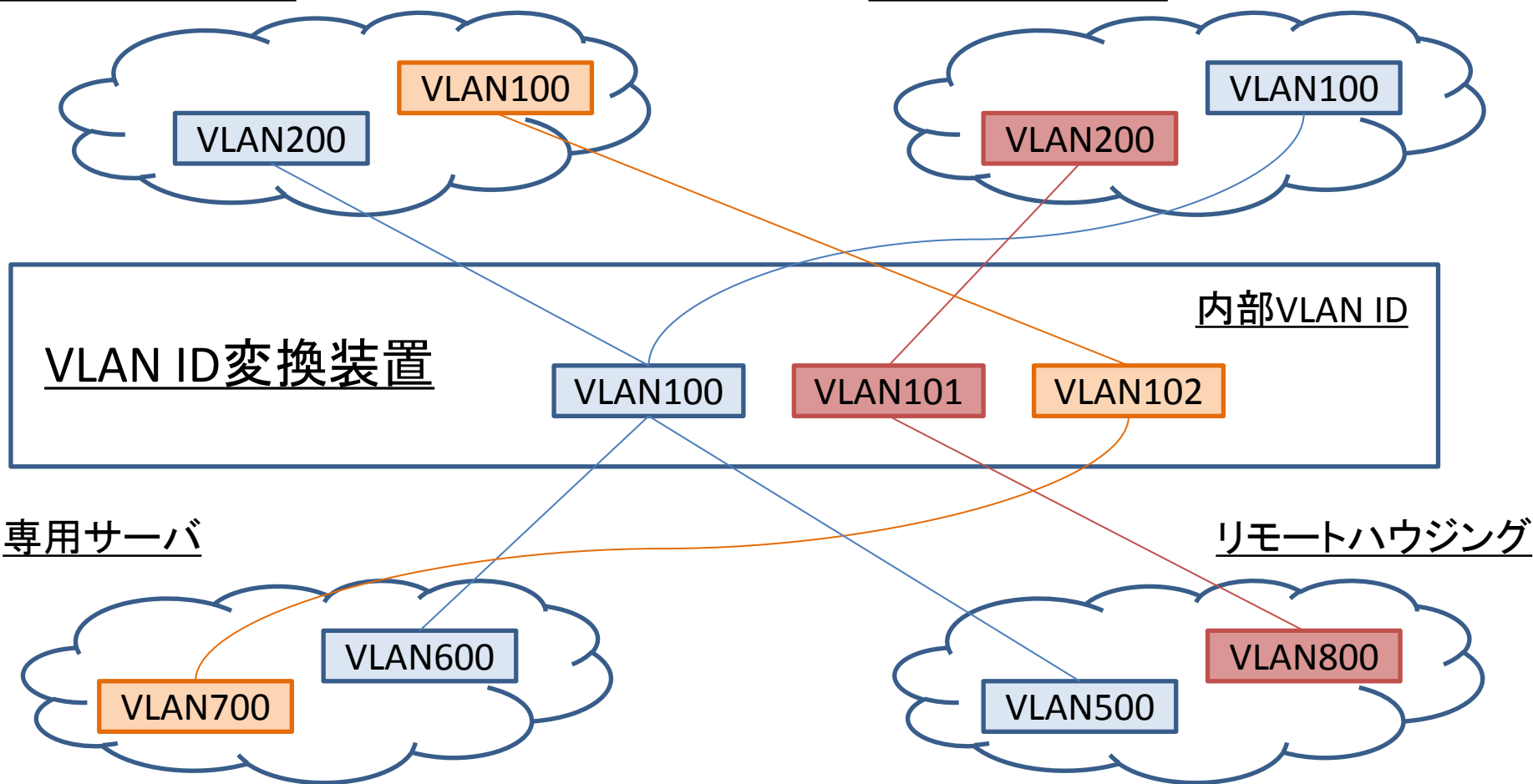


VLAN ID変換の
仕組みで相互接続

ネットワークが継ぎ接ぎになる

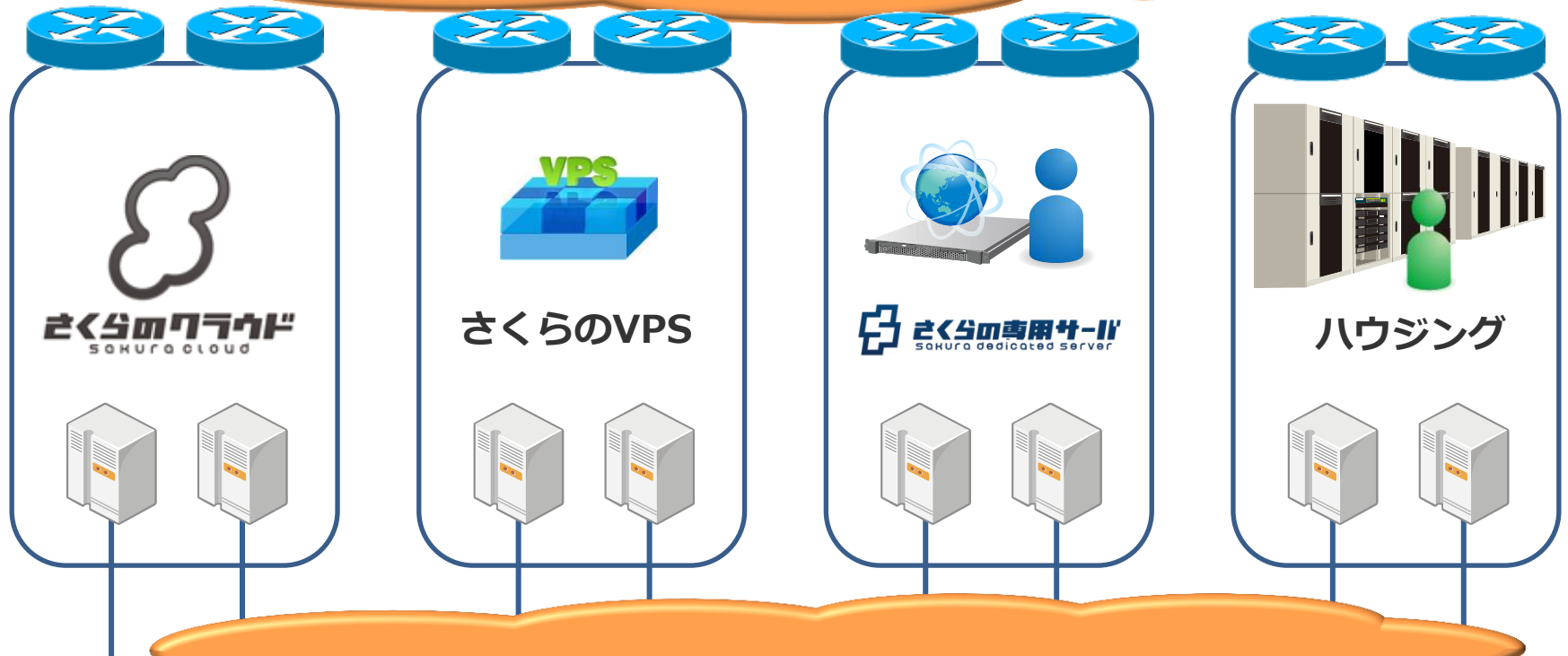
クラウド第1ゾーン

クラウド第2ゾーン



本当はこうしたい

インターネット



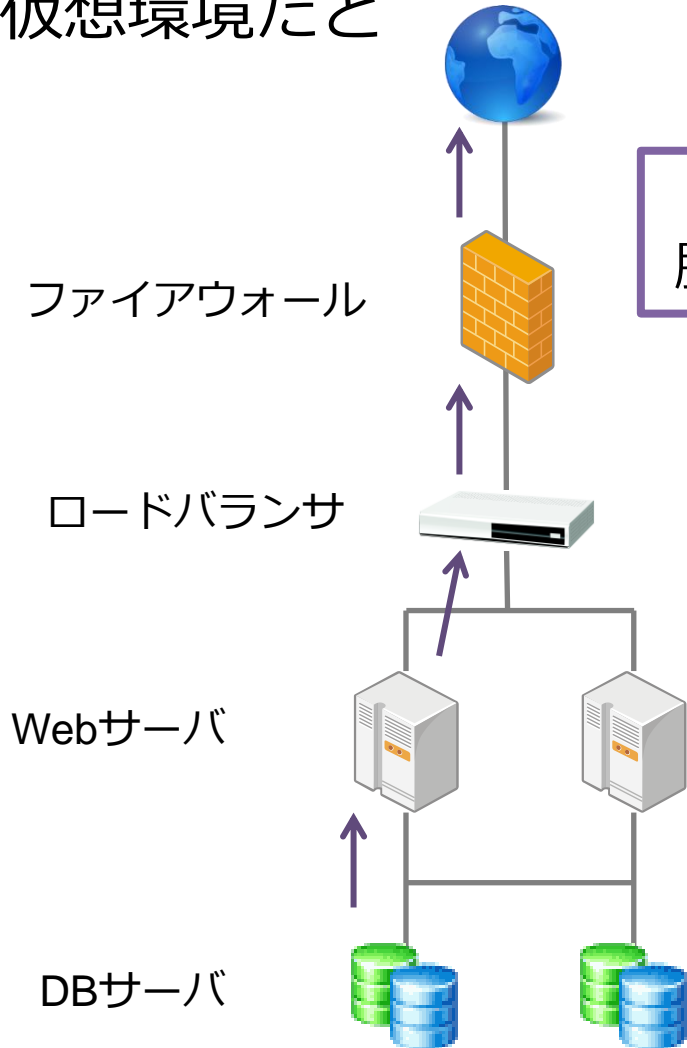
自由に接続できるVPN

そこで・・・SDN!?

- 2010年頃～
クラウドサービス開始前より
オーバーレイ方式の実装方法を検討
- 2011年頃～
プロダクトがいくつか出てきた
某N社さんのソフトウェアの検証を実施
- 2012年頃～
ミドクラさんと一緒に共同研究を実施
現在、具体的な商用導入に向けた取り組み

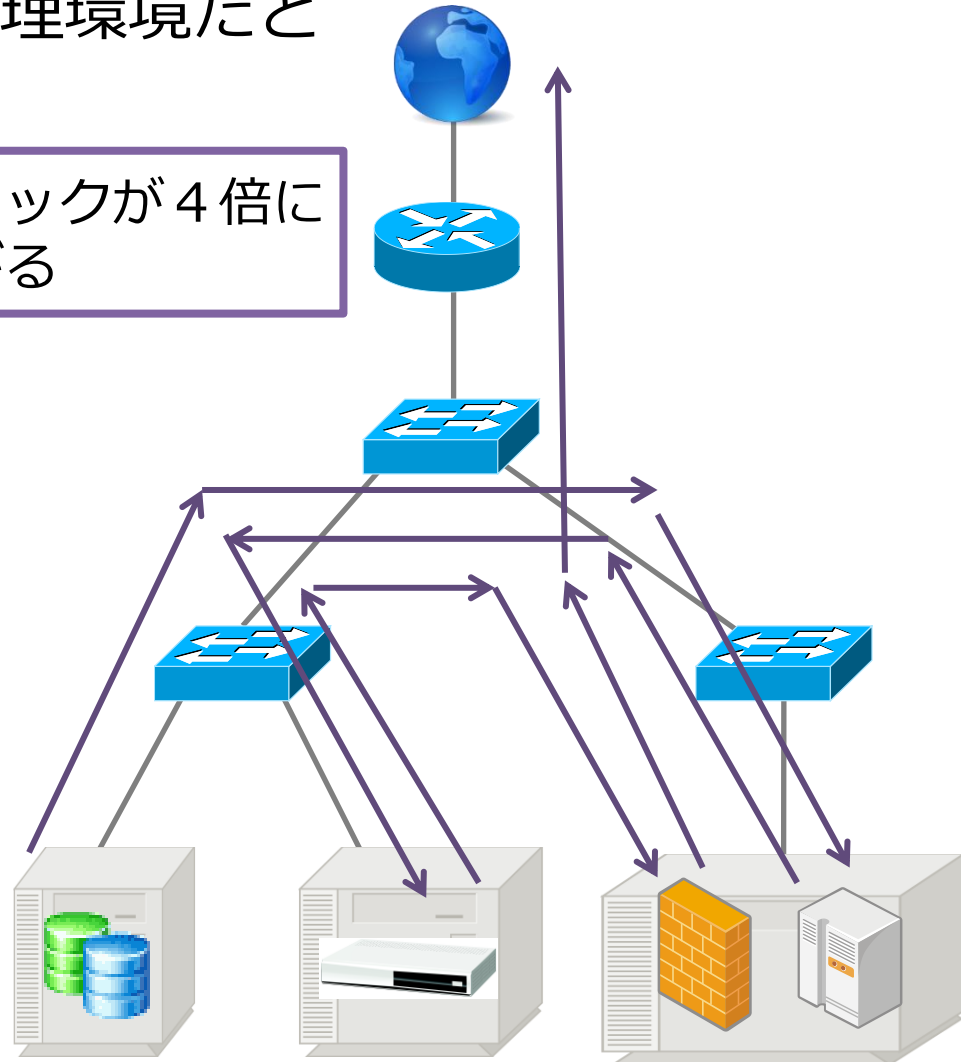
East-Westトラフィックの問題

仮想環境だと



物理環境だと

トラフィックが4倍に膨れ上がる

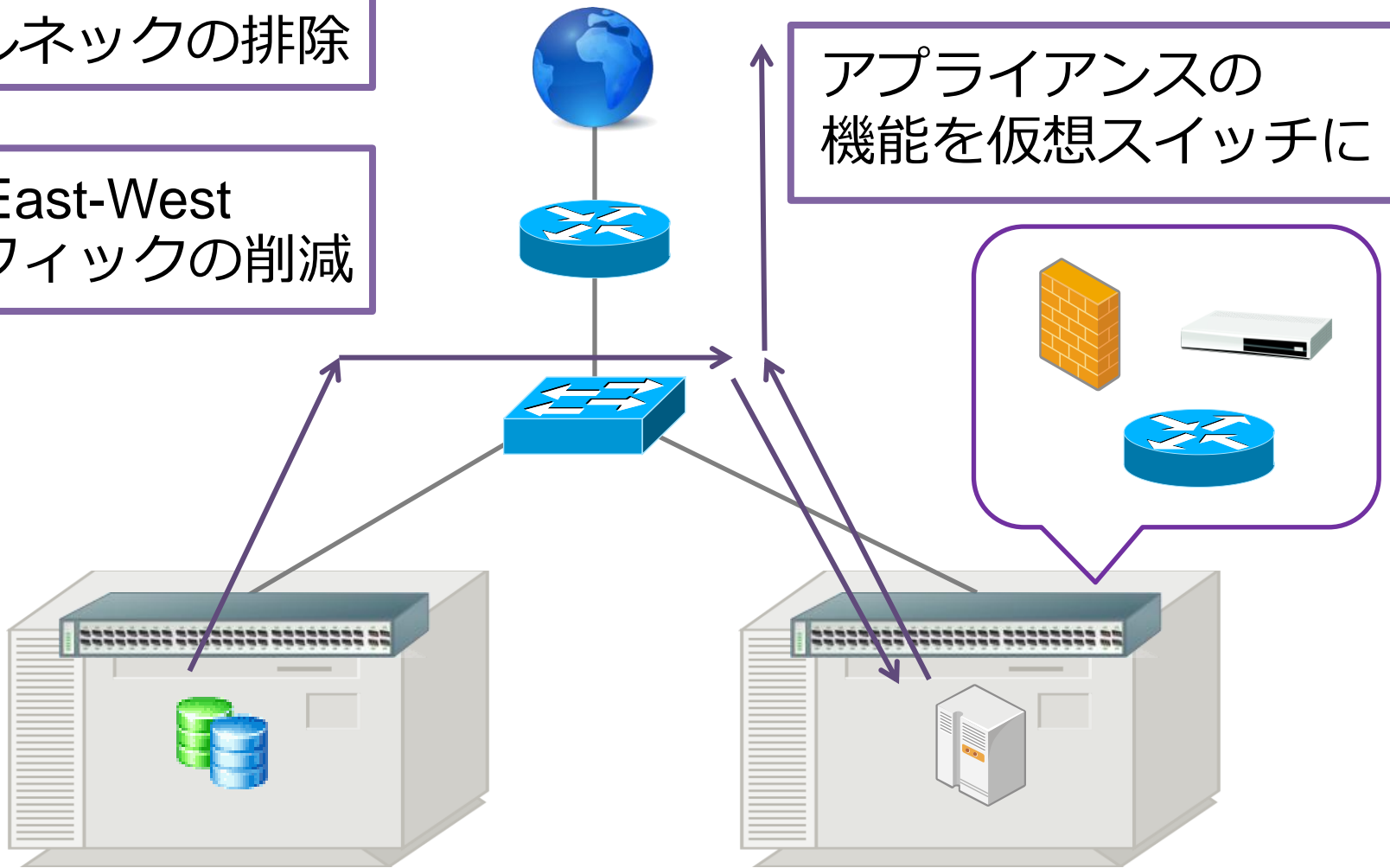


仮想スイッチのインテリジェント化

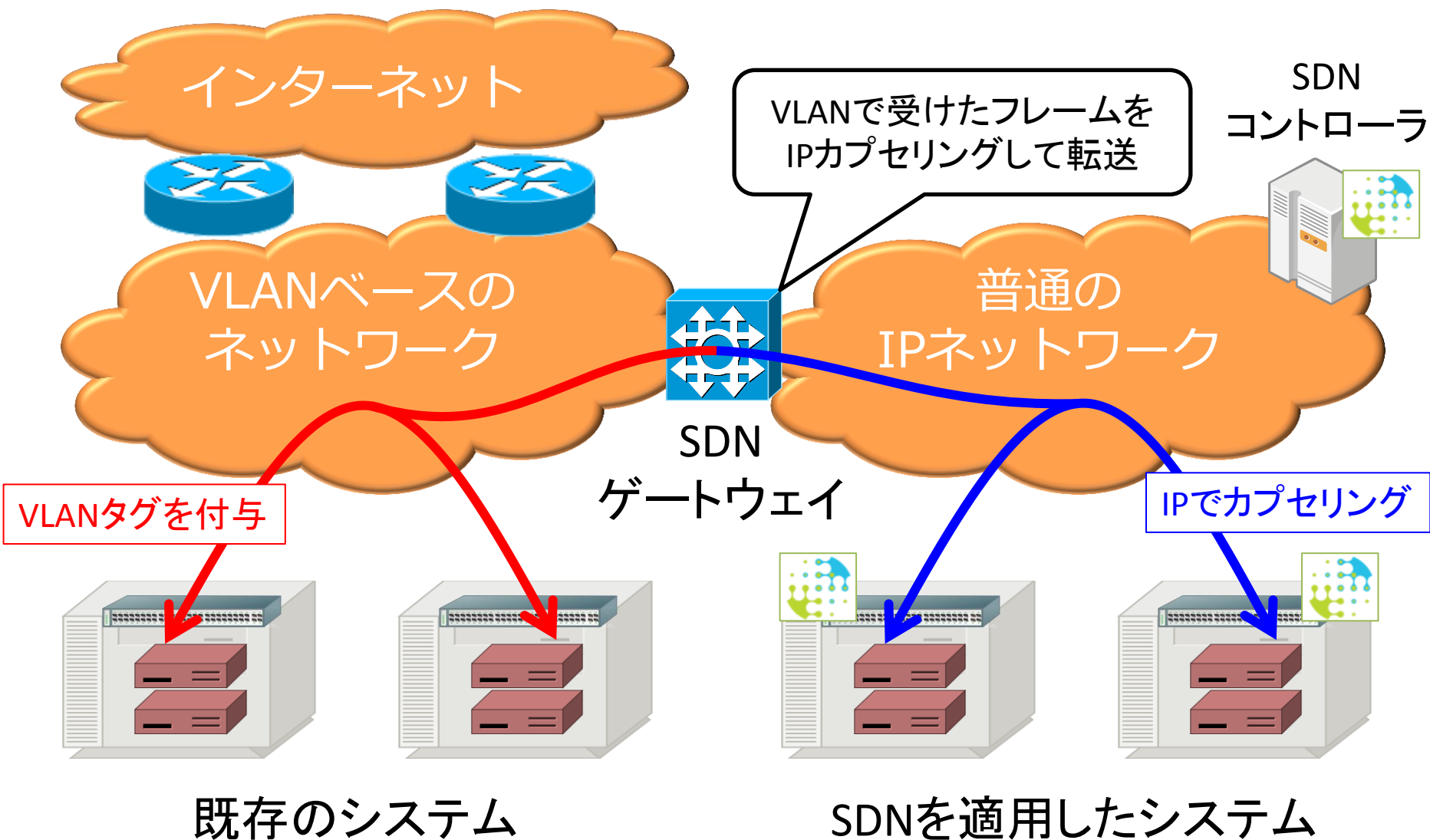
ボトルネックの排除

East-West
トラフィックの削減

アプライアンスの
機能を仮想スイッチに



SDNへの段階的移行



ご清聴ありがとうございました