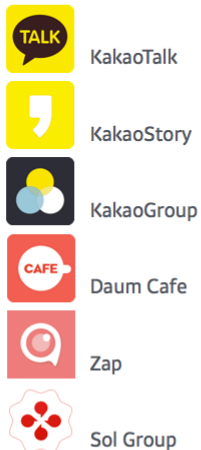L.T.H.

kakao

Scalable VM and Container Networking
using /32bit subnets and BGP routing

Andrew Yongjoon Kong
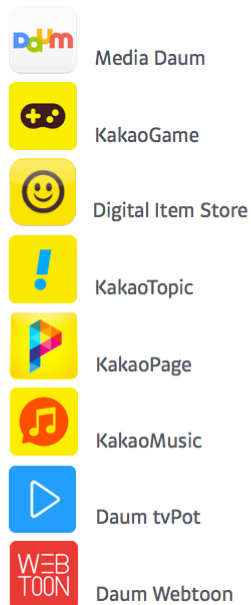
# DaumKakao

A Mobile Lifestyle Platform
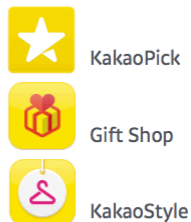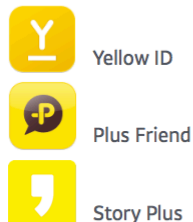
| Social Platform | Contents Platform | Commerce Platform | Marketing Platform | Local Platform | Personal Platform |
|---|---|---|---|---|---|
| KakaoTalk | Media Daum | KakaoPick | Yellow ID | Daum Map | KakaoHome |
| KakaoStory | KakaoGame | Gift Shop | Plus Friend | KakaoPlace | Sol calendar |
| KakaoGroup | Digital Item Store | KakaoStyle | Story Plus | | Sol Mail |
| Daum Cafe | KakaoTopic | | | | Daum Cluod |
| Zap | KakaoPage | | | | |
| Sol Group | KakaoMusic | | | | |
| | Daum tvPot | | | | |
| | Daum Webtoon | | | | |

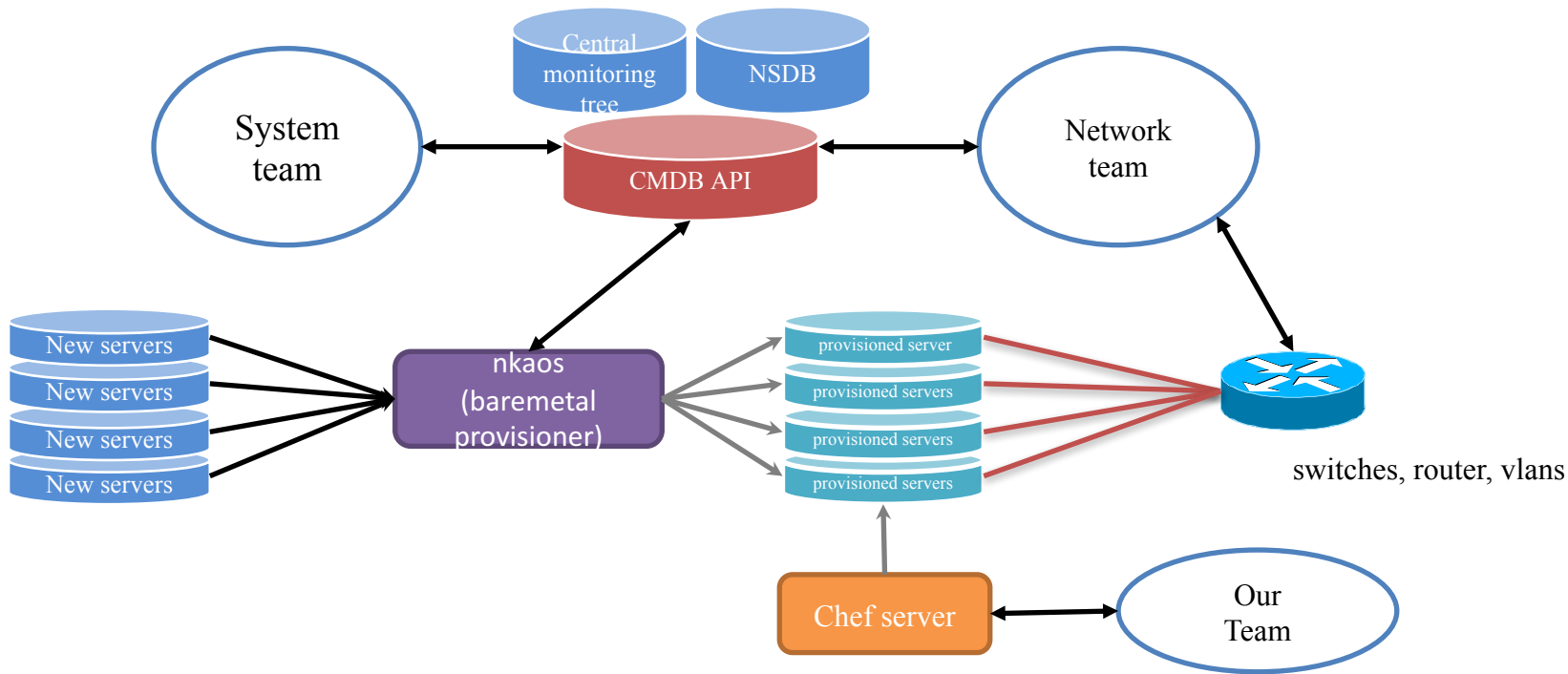96% of Korean smartphone users are using KakaoTalk messenger, 170 million users worldwide)

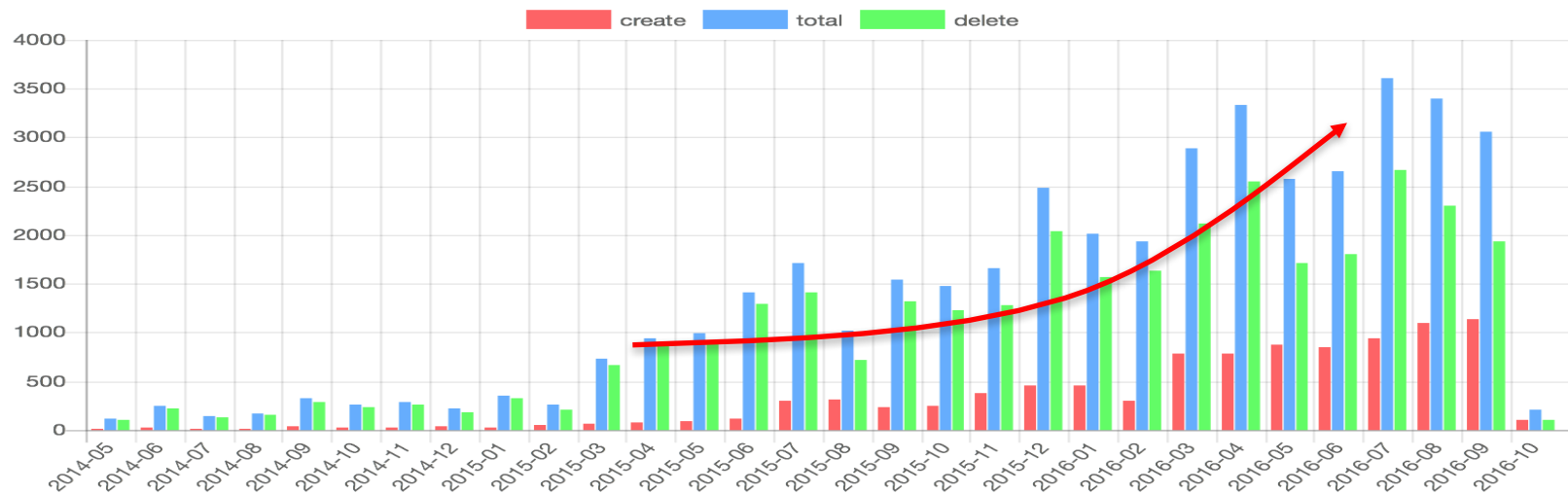2nd largest search and portal

# The Peaceful operation

When we're running out of resources ( cpu, memory, disk ),
Just add new( or additional ) resources to existing one.



switches, router, vlans

# The Growth(I)

kakao

VM creation speed is accelerating

# The Growth(II)

Spend more than 45M krane ( $45,000) per month
– this also increased.



krane

|  | 2015-01 | 2015-02 | 2015-03 | 2015-04 | 2015-05 | 2015-06 | 2015-07 | 2015-08 | 2015-09 | 2015-10 | 2015-11 | 2015-12 |

Totals:
- 2015-01: 37 811 340
- 2015-02: 35 245 330
- 2015-03: 42 441 940
- 2015-04: 46 086 450
- 2015-05: 59 857 170
- 2015-06: 67 636 380
- 2015-07: 83 622 260
- 2015-08: 103 393 480
- 2015-09: 119 856 680
- 2015-10: 135 056 980
- 2015-11: 131 972 260
- 2015-12: 149 897 920

Segment values:
- 2015-01: 3 154 220 / 2 172 720
- 2015-02: 3 065 280 / 1 965 000
- 2015-03: 3 616 580 / 2 456 400
- 2015-04: 3 721 140
- 2015-05: 4 497 200 / 12 627 040 / 21 333 840
- 2015-06: 5 361 660 / 12 776 760 / 24 679 120
- 2015-07: 7 902 420 / 15 810 240 / 37 236 080 / 3 430 880
- 2015-08: 9 267 040 / 20 424 040 / 52 519 680 / 3 677 600
- 2015-09: 10 170 840 / 24 971 080 / 64 787 360 / 3 725 440
- 2015-10: 10 974 420 / 27 127 440 / 77 257 280 / 3 518 560
- 2015-11: 12 103 540 / 26 791 760 / 75 356 800 / 3 340 800
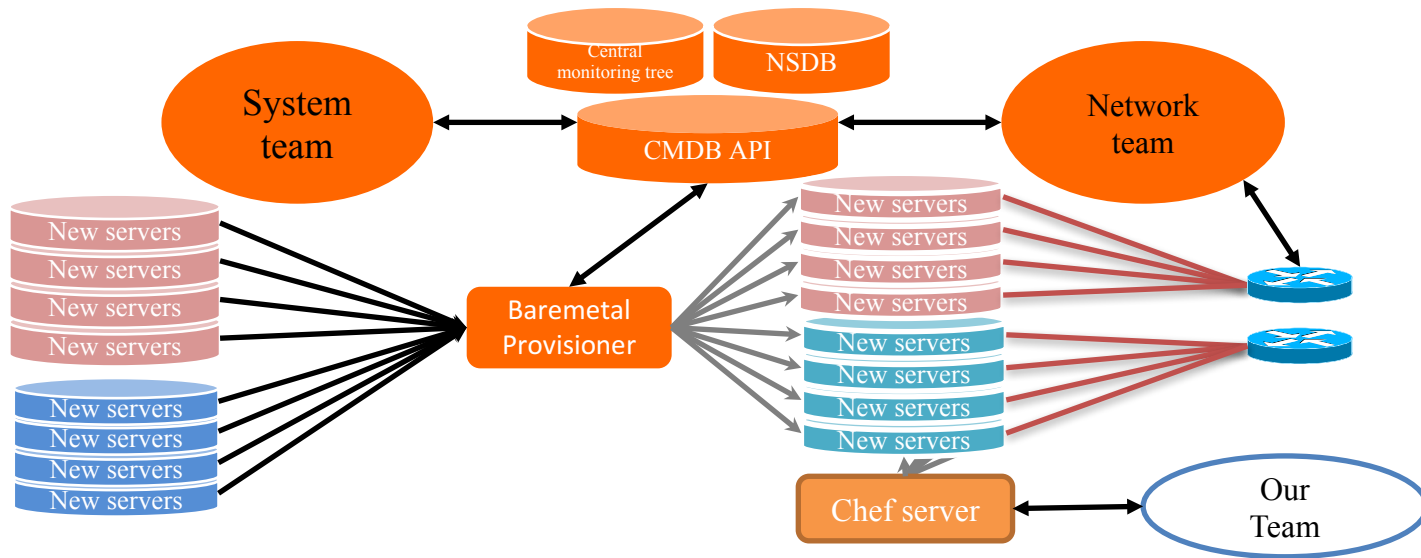- 2015-12: 12 939 920 / 31 016 720 / 3 728 160 / 87 933 920 / 3 361 440

1 krane = 1 Won ( $0.001)

- Using similar pricing with AWS EC2
- Network/Disk usage not included
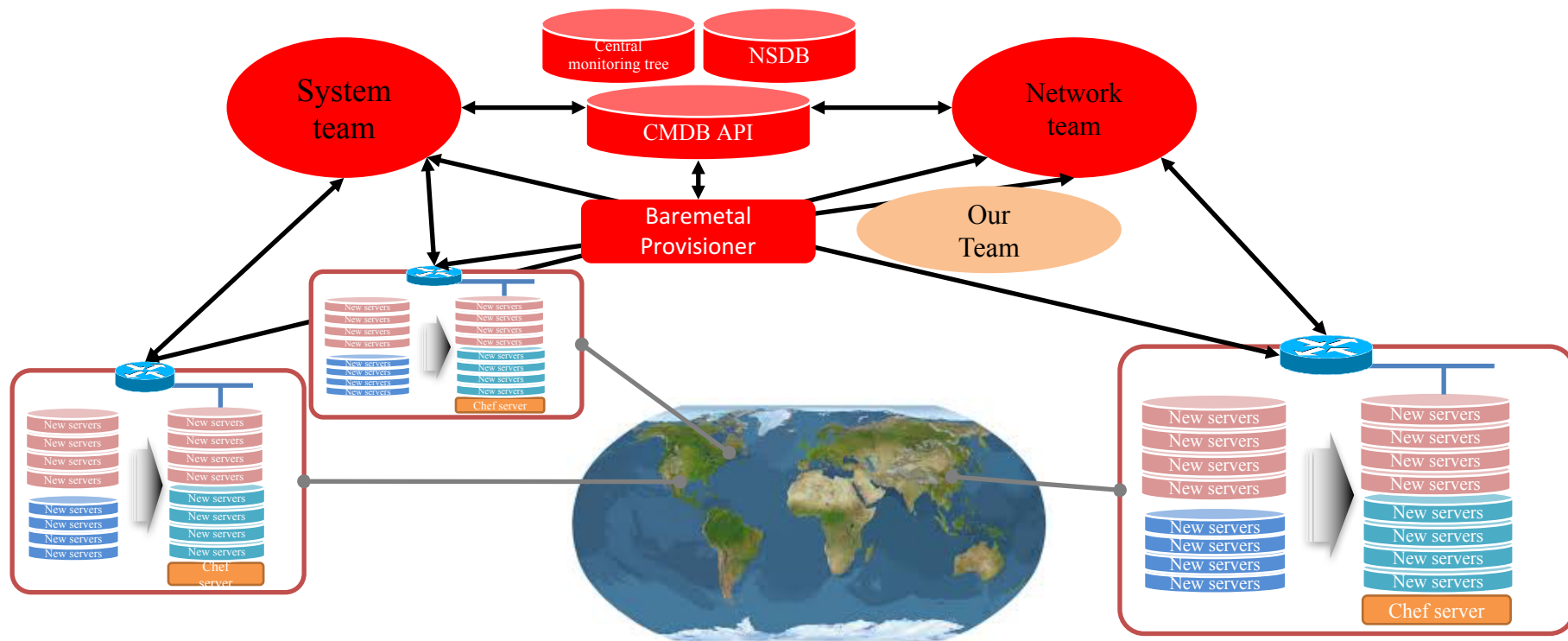
kakao

kakao

Growth is accelerating

- No. of Engineer is growing

- New Pilot services or experiments are growing.

- The resources depletion speed is accelerating → this simply make more work to resource management teams

# The Growth(IV)

Scale, The only driving force disrupt everything.

# The Growth – Lesson learned

kakao

Growth doesn't come alone

- Infra growth includes scale-up , scale-out as well

- Scale-up includes these

  - Add Server, Storage, Switches
  - Add more power facility to supply juice fluently
  - This is not that difficult.

- Scale-out include these
  - Add New Datacenters, New Availability Zones
  - This is nightmare!

This leads radical changes over everything

- The way of preparing, provisioning
- The way of monitoring, logging, developing

Some Numbers

**1021** tenants

**662** pull request since 2014.9

**136** VMs are created/deleted per day

Some information about kakao Openstack

kakao

openstack upgraded from grizzly to **Liberty**

total **4Region**

additional service **Heat/Trove/Sahara/Octavia**

kakao

Resources for Openstack finally comes to be exhausted

- CPU, Memory, Storage always experience shortages.

- They have skewness.

- Sometimes, CPU depleted. Sometimes, Storage depleted.
  - All resources are able to be re-balanced.
  - you can migration clients' VM ( image , volume )

- IP is also Resources.
  - Very limited than our expectations
    - No of IP counts is limited.
    - Location of IP also is limited.
  - Managing these Resources is getting tougher issue.

# OpenStack Neutron Network

We've been using Provider Network (VLAN)

– ML2 plugin

– From OVS → LinuxBridge.

– Network Team plan/setup networks (the VLAN, IP[subnet], Gateways)

– Mapping availability zone / Neutron Network to that Physical networks

# Resource Imbalance

## After Running multiple Available Zones

- Experiencing resource imbalance between zones, naturally
- Filter Scheduling won't helpful.
- Migration is a proper solution. ( add extra resource is better If possible )

openstack
scheduler

openstack™

X

Hey Openstack,
Create 1 VM ( 1cpu, 1 IP, 1 Storage)

VLAN.1

VLAN.2

VLAN.3

Zone1

1 CPU
1 storage
No IP
Left

Zone2

No CPU
No
Storage
1 IP

Zone3

# Resource Imbalance & Remedies

kakao

## Develop Network Count filter

- Check Remaining IP count for each zone, treat ip count as resource
- Select the zone which have more ip count
- but e



| | re... | Zone | m1_m... | m1_db | m1_la... | m1_s... | m1_xl... | m1_2... | IP Total | IP Used | IP Avail |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | 1 | Dev_Zone | 41(79) | 8(8) | 16(41) | 57(94) | 8(17) | 8(8) | 1485 | 954 | 531 |
| | 2 | DB_KR3_Prod_Zo... | 37(0) | 9(0) | 22(0) | 37(0) | 14(0) | 9(0) | 495 | 3 | 492 |
| | 3 | DB_KR3_Prod_Zo... | 37(0) | 9(0) | 22(0) | 37(0) | 14(0) | 9(0) | 495 | 3 | 492 |
| | 4 | DB_Zone00 | 20(0) | 5(0) | 11(0) | 24(0) | 8(0) | 5(0) | 495 | 12 | 483 |
| | 5 | DB_Zone01 | 24(0) | 5(0) | 12(0) | 3...(...) | ...(...) | | 495 | 5 | 490 |
| | 6 | Prod_Zone | 24(5) | 5(3) | 10(5) | 33(5...) | ...(5) | ...(5) | 495 | 236 | 259 |
| | 7 | DB_Alpha_Zone01 | 15(0) | 4(0) | 9(0) | | | | 495 | 5 | 490 |
| | | | | | | | | ...(0) | 495 | 9 | 486 |

1-8 of 8

kakao

## Rethinking about Connectivity

# Rationale

## Rethinking about Connectivity (Overlay)

- it solve remote link layer separation issue.
- Still have issue with IP management. and Gateway ( Packet Forwarding)

we need to thinks of those requirement

- IP movement inter-rack, inter-zone, inter-dc(?)
- IP resource imbalance
- Fault Resilience
- Dynamically check status of network
- Simple IP Resource Planning and Management

We thinks Router as best candidate

- It dynamically detects and exchanges changes. (via dynamic routring protocol)
- It is highly distributed.
- It have HA ( e.g. VRRP)
- the issue is that most of time routing is done in ranges (a.k.a Subnet)
  - Because of Memory and CPU issue

Finally, Come to route only IP

Generally, Known as /32 network.

10.0.0.1 / 32  or
IP 10.0.0.1  netmask 255.255.255.255

- No L2 (link) consideration needed anymore ( no subnet )
- With Dynamic Routing Protocol,  it move every where.
- Simple IP planning ( Just think of IP ranges )
- It's very Atomic Resource, it keeps its IP after migration through zones

# How it setup

1.  install nova/neutron agent.
2.  create neutron network ( name: freenet, subnet: 10.10.100.0/24)

# How it setup

1. install nova/neutron agent.
2. create neutron network ( name: freenet, subnet: 10.10.100.0/24)
3. user create VM

# How it works

1. install nova/neutron agent.
2. create neutron network ( name: freenet, subnet: 10.10.100.0/24)
3. user create VM
4. update Routing(with Dynamic routing protocol)

| Routing Table | |
|---|---|
| 1 | 10.100.10.2/32  via 192.1.1.201 |

192.1.1.202

192.1.1.201

**Compute node**

dhcp-server process

| Routing Table | |
|---|---|
| Default GW | 192.168.1.1 eth1 |
| Host Route | dest 10.10.100.2/32 to 10.10.100.1 |

10.10.100.1

eth1

advertising:
via Dynamic Routing Protocol

linux bridge

IP:10.10.100.2/32
GW: 10.10.100.1

vm

neutron-dhcp-agent

neutron-linuxbridge-agent

nova-compute

eth0

Controller

Phase 1

kakao

Use RIP and OSPF

- Heterogeneous setting will be burden
- Using Default GW as eth1 even for compute node.
  Management and service network mixed.

# kakao

## Use BGP and switch namespace

– Isolating vm's traffic using switch namespace.

– adopting same dynamic routing scheme to compute node

| Routing Table | |
|---|---|
| 1 | 10.100.10.2/32  via 192.1.1.201 |

192.1.1.202

eBGP

**Compute node**

Switch Namespace    dhcp-server process

| Routing Table | |
|---|---|
| Default GW | 192.168.1.1 eth1 |
| Host Route | dest 10.10.100.2/32 to 10.10.100.1 |

192.1.1.201

iBGP

eth1

10.10.100.1

linux bridge

IP:10.10.100.2/32

neutron-dhcp-agent

neutron-linuxbridge-agent

vm

| Routing Table | |
|---|---|
| Default GW | x.x.x.x eth0 |

nova-compute

eth0

Controller

global name space

kakao



| Routing Table | |
|---|---|
| 1 | 10.100.10.2/32 via t |

| Routing Table | |
|---|---|
| 1 | 10.100.10.2/32 via tor2 |

| Routing Table | |
|---|---|
| 1 | 10.100.10.2/32 via RT1 |

| Routing Table | |
|---|---|
| 1 | 10.100.10.2/32 via RT3 |

tor1
tor2
tor3

**Compute node1**
Switch Namespace
linux bridge
neutron-dhcp-agent
neutron-linuxbridge-agent
nova-compute
rt1
AZ1

10.10.100.2/32
vm

**Compute node2**
Switch Namespace
linux bridge
neutron-dhcp-agent
neutron-linuxbridge-agent
nova-compute
rt2
global name space

**Compute node1**
Switch Namespace
linux bridge
neutron-dhcp-agent
neutron-linuxbridge-agent
nova-compute
rt3
AZ2
global name space

**Compute node2**
Switch Namespace
linux bridge
neutron-dhcp-agent
neutron-linuxbridge-agent
nova-compute
rt4
global name space

**Compute node1**
Switch Namespace
linux bridge
neutron-dhcp-agent
neutron-linuxbridge-agent
nova-compute
rt5
AZ3
global name space

**Compute node2**
Switch Namespace
linux bridge
neutron-dhcp-agent
neutron-linuxbridge-agent
nova-compute
rt6
global name space

# What we solve?

kakao

Simple IP planning

– only IP ranges matter. (no more VLAN, IP subnet, Router planning)

Resource imbalancing

– No chance of IP imbalancing.

Fault Resilience

– If one router gone, it propagated by Dynamic routing protocol to other router

Distributed

– deciding routing path is very distributed. No single point of failure.

– scale out nature.
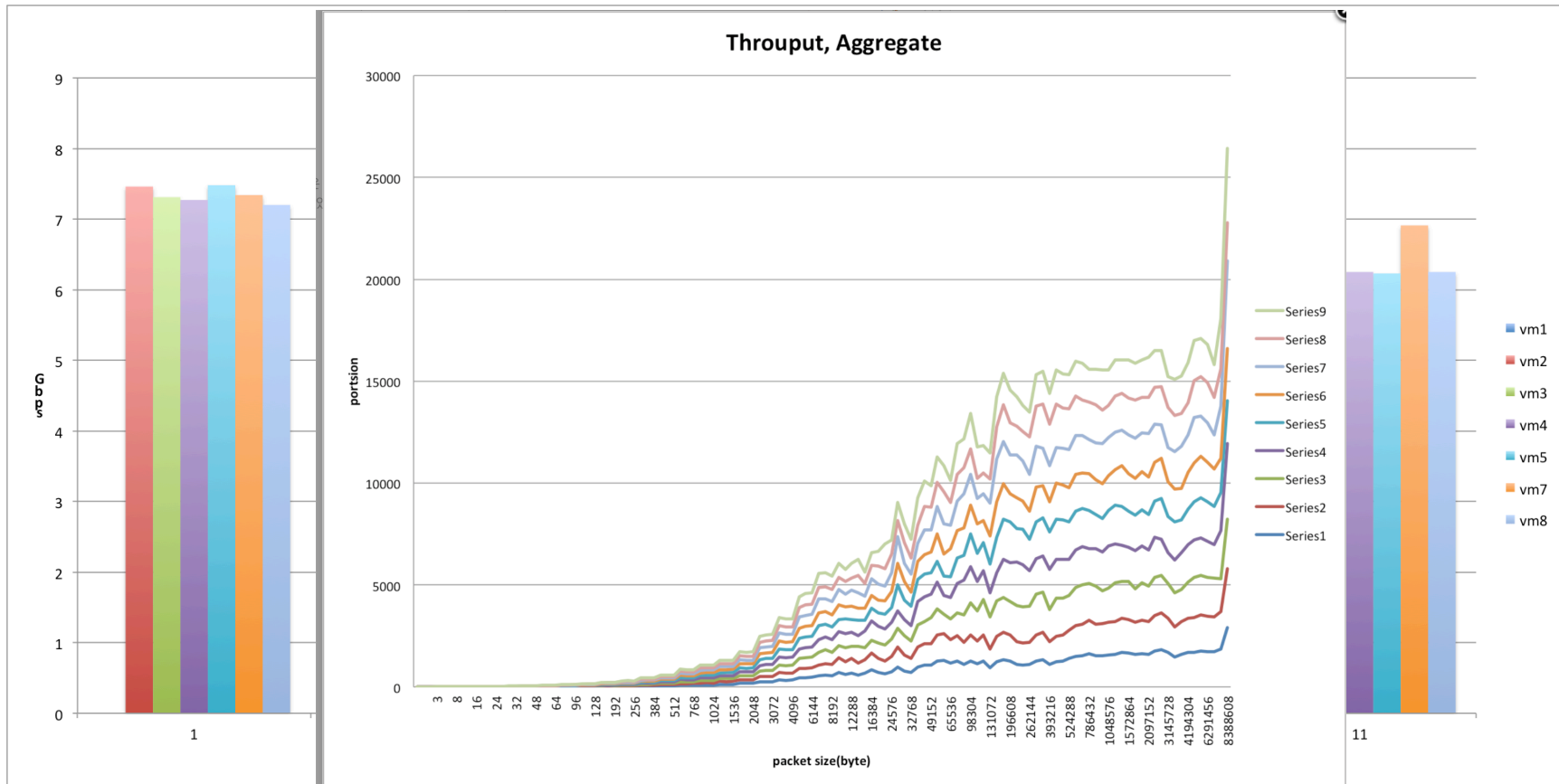
What we still have to solve?

Still many issue
- – Apply this to physical server
- – Making Router setup by API ( REST, RPC) using seed BGP( only advertising)
- – ACL propagation using API ( e.g. Flowspec)
- – Shared storage base service

# Performance Test VMs to VMs

kakao



Throuput, Aggregate

kakao



```
krane-prod-md2-48# sh ip ro sum
Route Source          Routes              FIB
kernel                28                  28
connected             11                  11
ebgp                  0                   0
ibgp                  1                   0
------
Totals                40                  39
krane-prod-md2-48#  sh bgp mem
84 RIB nodes, using 9408 bytes of memory
32 BGP routes, using 2048 bytes of memory
31 Adj-Out entries, using 1240 bytes of memory
1 Nexthop cache entries, using 24 bytes of memory
5 BGP attributes, using 280 bytes of memory
5 BGP extra attributes, using 440 bytes of memory
2 BGP AS-PATH entries, using 64 bytes of memory
1 BGP AS-PATH segments, using 24 bytes of memory
2 peers, using 9120 bytes of memory
24 hash tables, using 960 bytes of memory
36 hash buckets, using 864 bytes of memory
```

# Application of /32bit network: /32bit route + DNAT
→ 1:1 NAT (A.K.A FloatingIP )

| Routing Table | |
| --- | --- |
| 1 | 10.10.100.2/32 via 192.1.1.201 |
| 2 | 10.10.100.3/32 via 192.168.1.202 |
| 3 | 192.168.100.2/32 via 192.168.1.201 |

192.1.1.202

Compute node1

| IPTable | |
| --- | --- |
| DNAT | Dest 192.168.100.2 is forwarded to 10.10.100.2 |

eth1

192.1.1.201
Compute Node Router

Switch Namespace

| Routing Table | |
| --- | --- |
| Default GW | 192.168.1.1 eth1 |
| Host Route | dest 10.10.100.2/32 to 10.10.100.1 |
| connected | dest 192.168.100.2 |

linux bridge

IP:10.10.100.2/32

vm

global name space

# Application of /32bit network: ECMP + DNAT
→ Scalable Loadbalancer

VIP: 192.168.100.2 is ECMPed

Aggregation

TOR1    192.1.1.201

TOR2

**Compute node1**

| IPTable | |
|---|---|
| DNAT | Dest 192.168.100.2 is forwarded to 10.10.100.2 |

eth1

eth1

192.1.1.201

Compute Node Router

**192.1.1.202**

Compute Node Router

Switch Namespace

| Routing Table | |
|---|---|
| Default GW | 192.168.1.1 eth1 |
| Host Route | dest 10.10.100.2/32 to 10.10.100.1 |
| connected | dest 192.168.100.2 |

**Compute node2**

| IPTable | |
|---|---|
| DNAT | Dest 192.168.100.2 is forwarded to 10.10.100.3 |

| Routing Table | |
|---|---|
| Default GW | 192.168.1.1 eth1 |
| Host Route | dest 10.10.100.3/32 to 10.10.100.1 |
| connected | dest 192.168.100.2 |

Switch Namespace

linux bridge

IP:10.10.100.2/32

LB

global name space

linux bridge

IP:10.10.100.3/32

LB

global name space
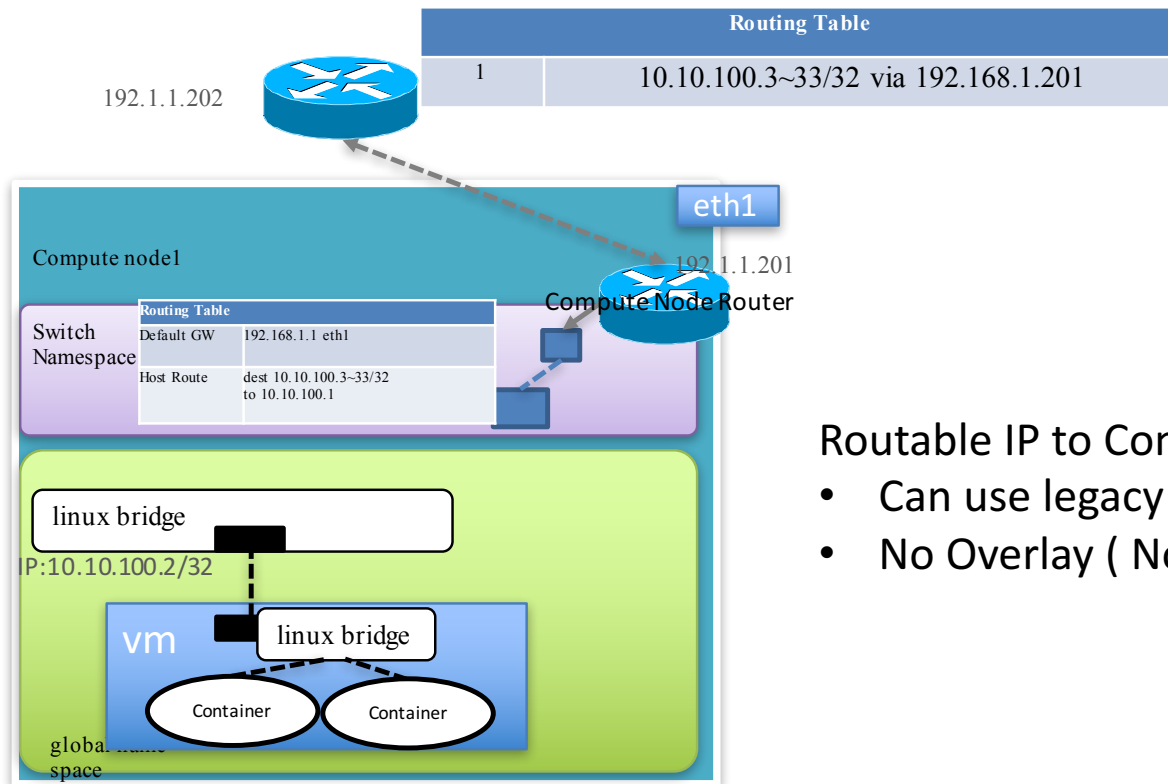
**Application of /32bit network**:

Multiple Routing Entry ( AKA, Fixed IPs) + Container Bridge Network

→ Scalable Container Network

kakao

192.1.1.202

| Routing Table | | |
|---|---|---|
| 1 | 10.10.100.3~33/32 via 192.168.1.201 | |

eth1

Compute node1

192.1.1.201

Compute Node Router

Switch Namespace

| Routing Table | | |
|---|---|---|
| Default GW | 192.168.1.1 eth1 | |
| Host Route | dest 10.10.100.3~33/32 to 10.10.100.1 | |

linux bridge

IP:10.10.100.2/32

vm

linux bridge

Container

Container

global name space

Routable IP to Container:

- Can use legacy IP base Monitoring
- No Overlay ( No complexity )

# Q&A

Thanks