# About me



## Wei Zhou

- Apache CloudStack community member since 2012

- Apache CloudStack committer since 2013.05

- Apache CloudStack PMC member since 2017.03

- Software Architect @ Shapeblue since 2021

- Member of Kubernetes org

- Member of OpenSDN.io

- Email:   weizhou@apache.org

- Github:  @weizhouapache

# Contents

# 01

## Why Routed mode

# Guest networks in CloudStack

Guest network types
- Shared
- Isolated
- L2

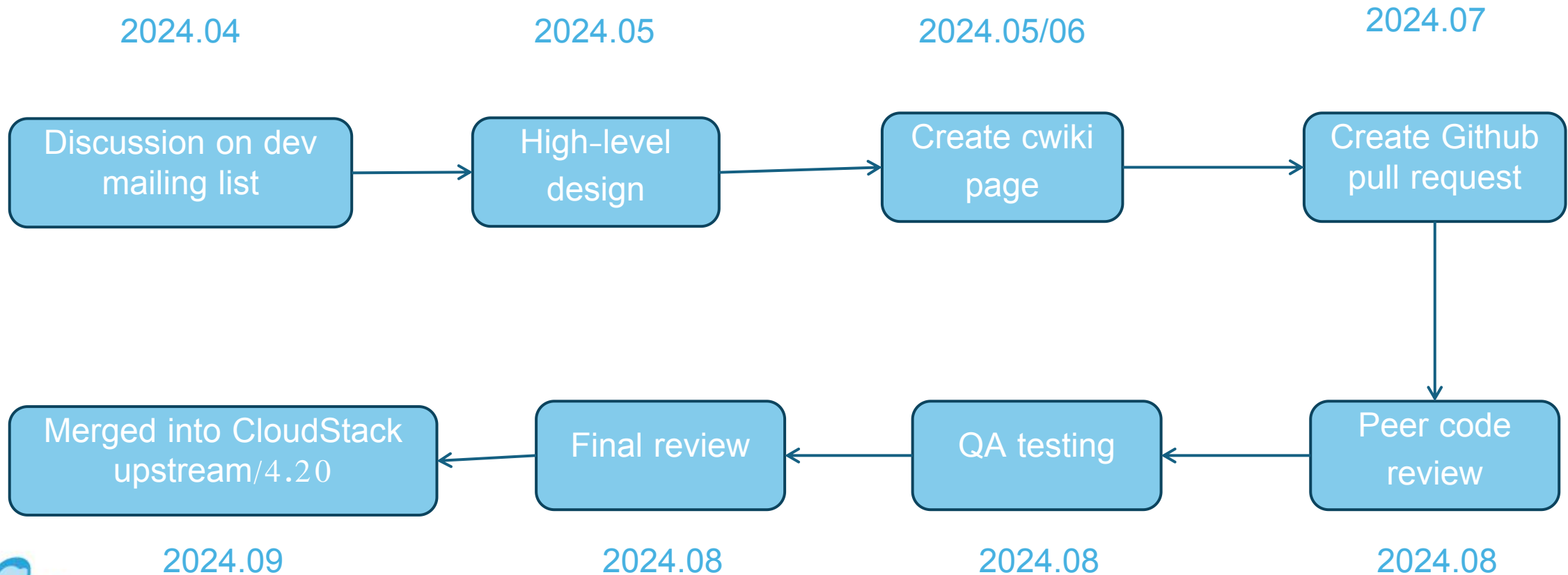Problems to address for Shared/Isolated network:
- IP assignment  (IPv4/IPv6 address, netmask, gateway. Dhcp/Dns)
- IP Routing       (Between VM instances and the Internet)
- Network Access Control   (For inbound/outbound traffic)

# Overview of network types

| | Shared network | Isolated network | New network type or mode ❓ |
|---|---|---|---|
| **IP assignment** | Direct IP<br>Publicly accessible | Private IP<br>No public access | Direct IP<br>Publicly accessible |
| **IP Routing** | CloudStack VR is not gateway<br><br>Requires operators:<br>- Configure gateway on the upstream router<br>- Create the network in ACS | CloudStack VR is Source NAT gateway<br><br>Ways to access VMs:<br>Static NAT, Load Balancing, Port Forwarding, VPN | No Source NAT.<br><br>Doesn't require operators' manual configuration.<br>Can be created by end users. |
| **Network Access control** | Security groups (KVM only*) | Egress rules<br>Firewall rules | Ingress/Egress Firewall rules |
| | | | |
| **Performance** | Good | Not as good as Shared network | Better performance than Isolated network |

# Routed mode: timeline

2024.04

2024.05

2024.05/06

2024.07

Discussion on dev mailing list → High-level design → Create cwiki page → Create Github pull request

Merged into CloudStack upstream/$4.20$ ← Final review ← QA testing ← Peer code review
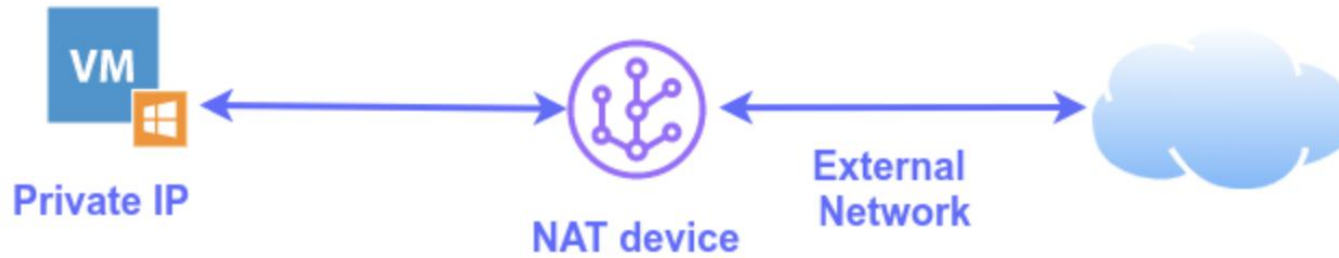
2024.09

2024.08
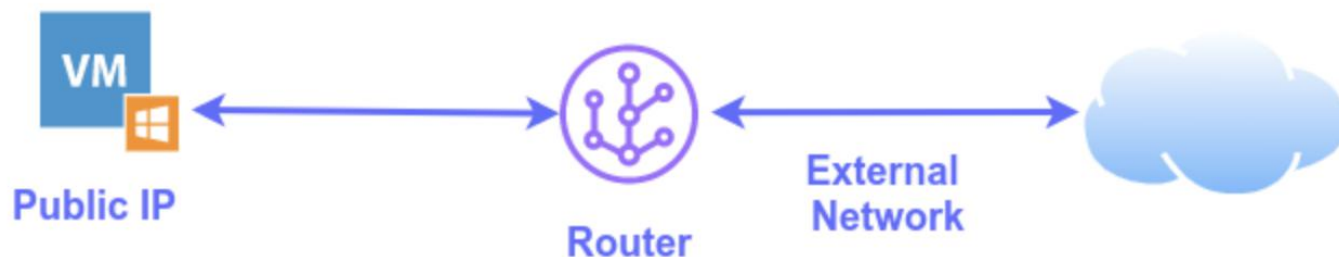
2024.08

2024.08

## 02 What is Routed mode

# New concept in ACS 4.20: Network mode

- NATTED mode



  - Default network mode for Isolated networks
  - Virtual Router (VR) as Source NAT (Network Address Translation) gateway

- ROUTED mode:



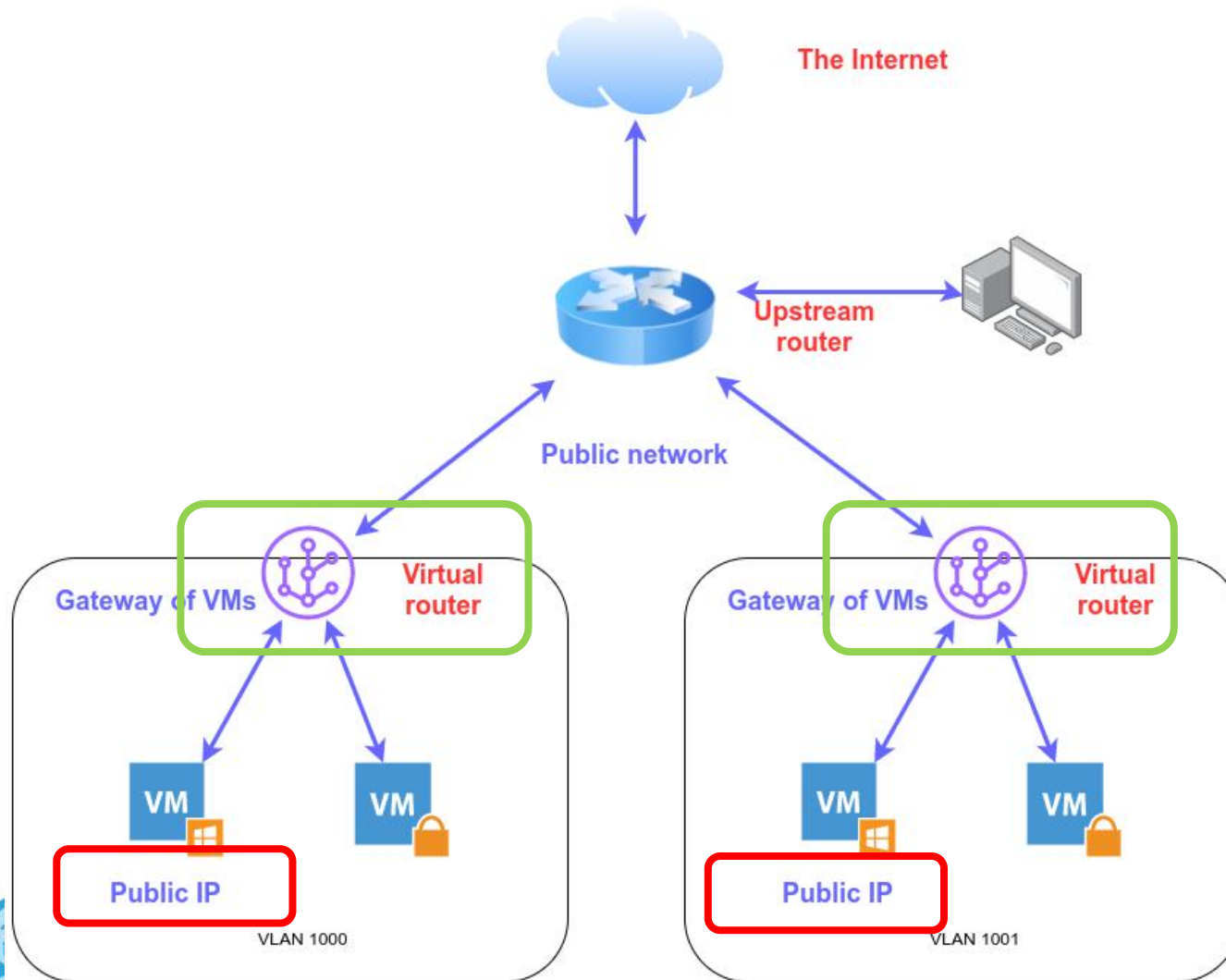  - Virtual Router (VR) as Gateway

**?**

Routed mode has already been used in CloudStack.
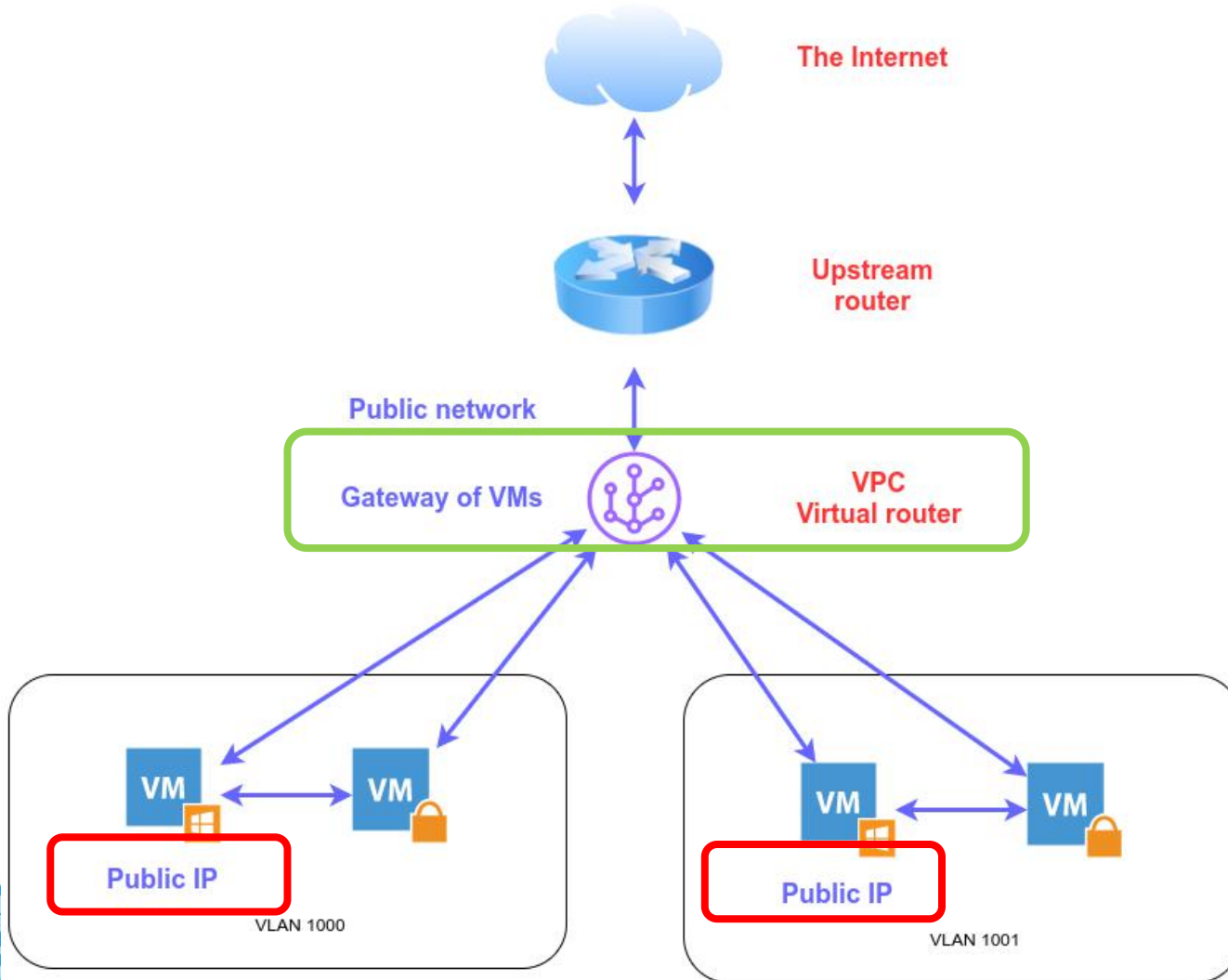
Which feature uses it ?

IPv6 support

# Routed Isolated network: Topology



- Very similar to "Isolated network"

- Virtual router (VR) as Gateway of VMs

- Virtual router (VR) provides services
  - Dhcp/Dns/Userdata
  - Firewall

- Differences from Isolated network
  - Public IP vs Private IP
  - No StaticNat/Lb/PF/VPN support

# Routed VPC: Topology



- Very similar to "VPC"

- VPC Virtual router (VPC VR) as Gateway of VMs

- VPC Virtual router (VPC VR) provides services
  - Dhcp/Dns/Userdata
  - Network ACL

- Differences from VPC
  - Public IP v.s. Private IP
  - No StaticNat/Lb/PF/VPN support

# IP Routing in Routed mode

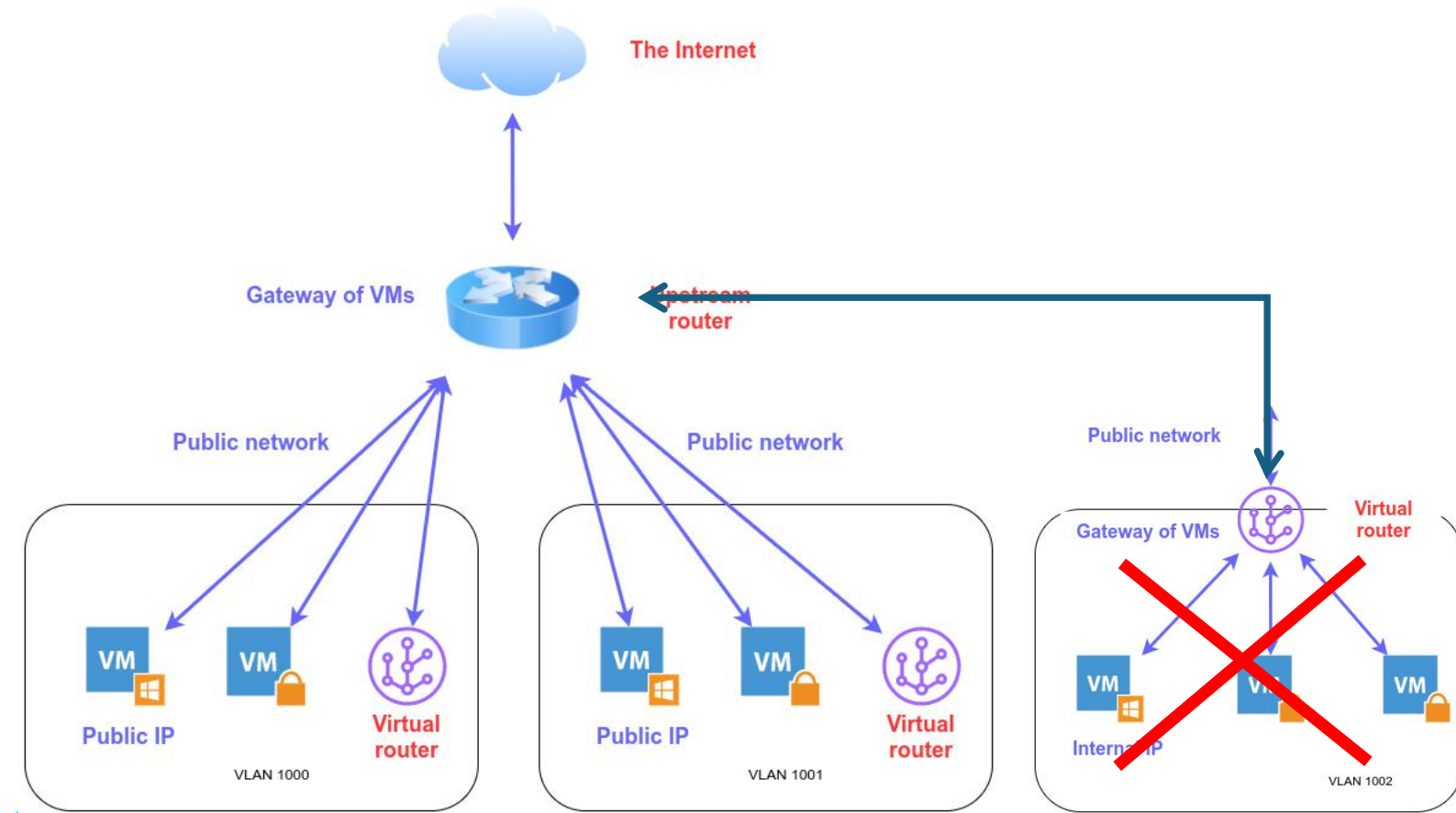| Routing mode | What the operators need to do |
|---|---|
| **Static routing** | Operators have to manually add **static routes** for each Routed network in the upstream router.<br><br>Tips: IPv6 implementation already supports it. |
| **Dynamic routing** | Operators configure **Dynamic BGP** in the upstream router<br><br>The routes for guest networks will be automatically advertised to (upstream and virtual) routers via dynamic routing protocol (BGP). |

# Network Access Control in Routed mode

- Routed networks
  - Egress rules (improved)
  - IPv4 Routing firewall (new!)
  - IPv6 firewall
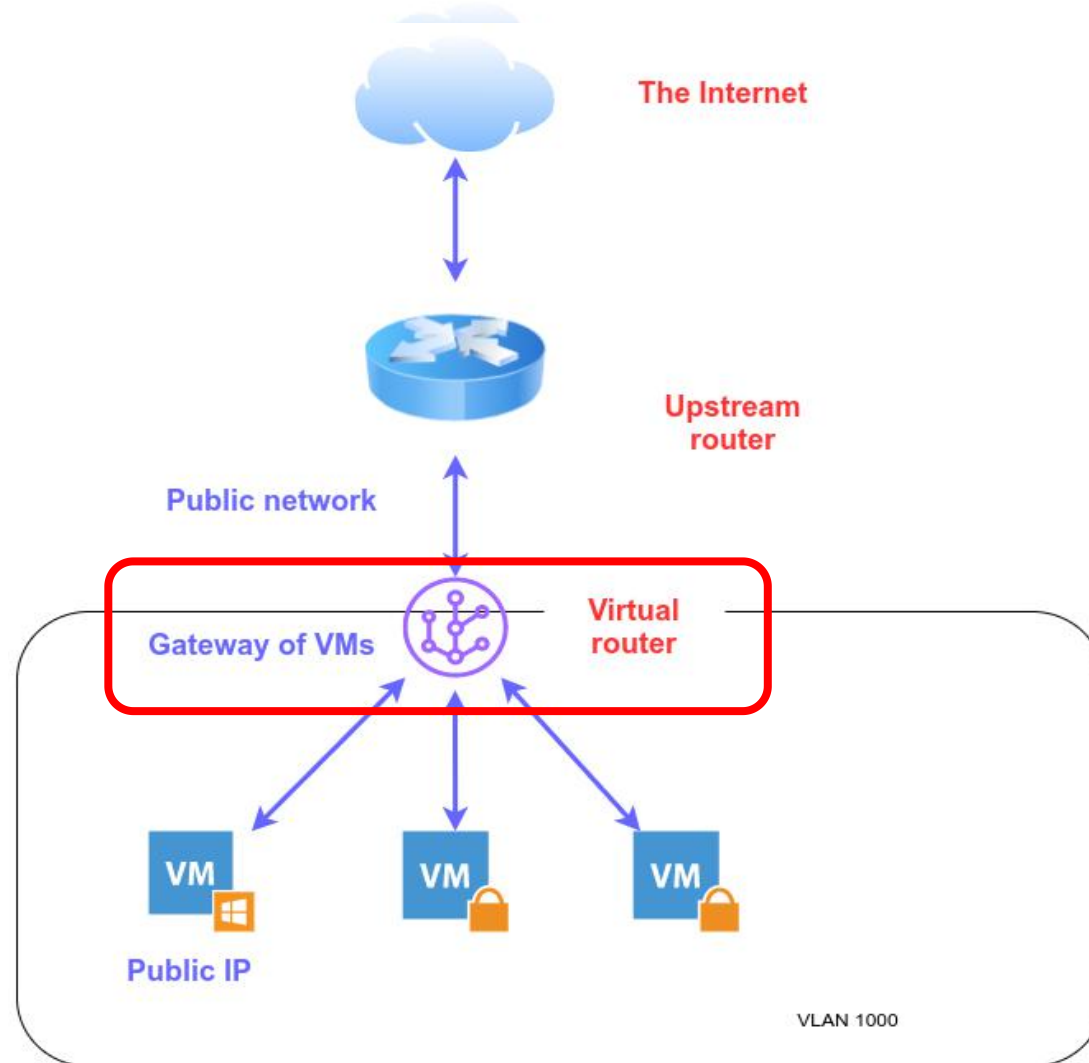
- Routed VPC
  - Network ACL (improved)

# Use case 1: Scalable private cloud



- Private cloud for large organization

- Uses Shared networks with Internal IPs

- New or scale out application
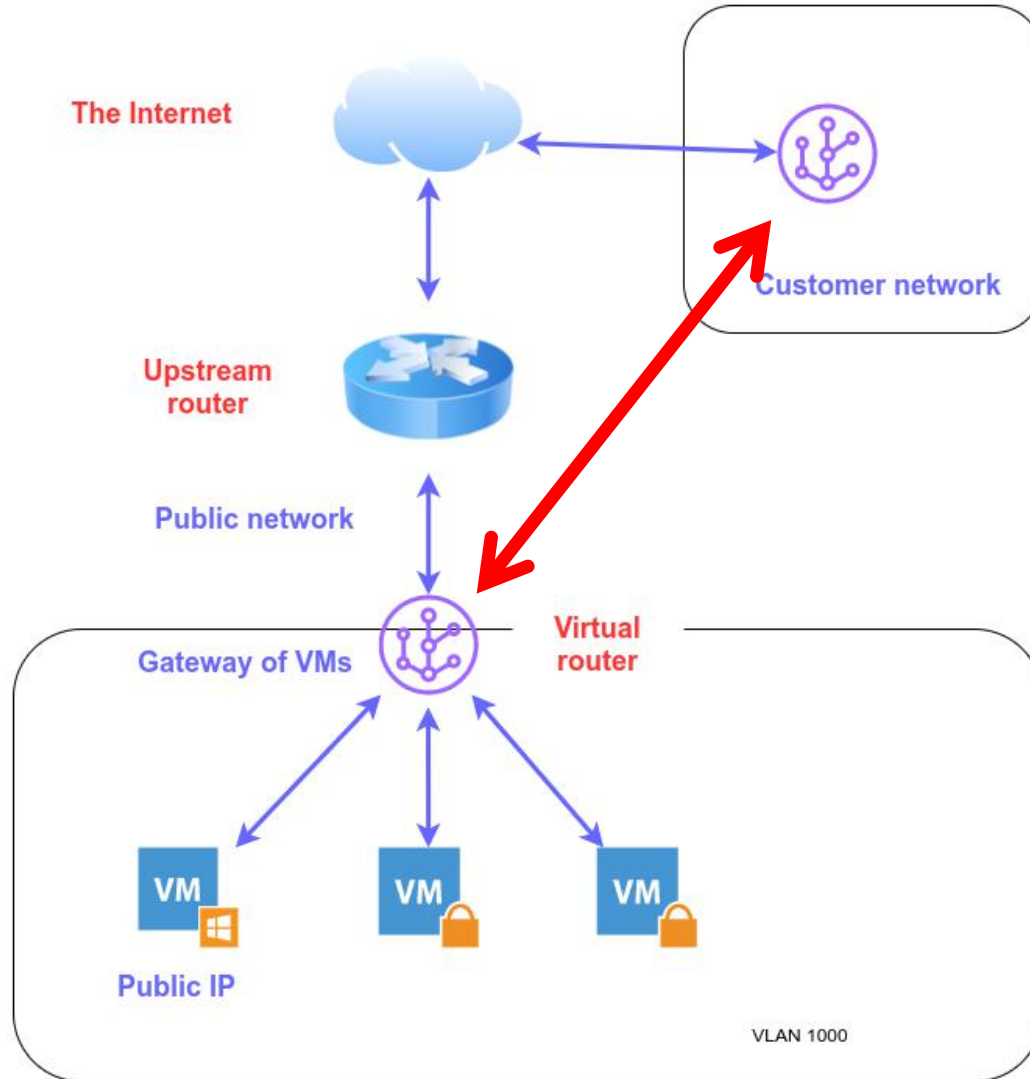
- More flexible and scalable

# Use case 2: VM with Access control on VMware



- Shared network on VMware (or XenServer/Xcp-ng)

- No SG support as Security group rules are applied on hypervisors (KVM only)

- Option: Routed network with access control support

# Use case 3: Bring Your own IPs



- Connects Customer network to Routed network

- via Dynamic routing

- Easy to setup a hybrid cloud

03 How to configure a Routed network

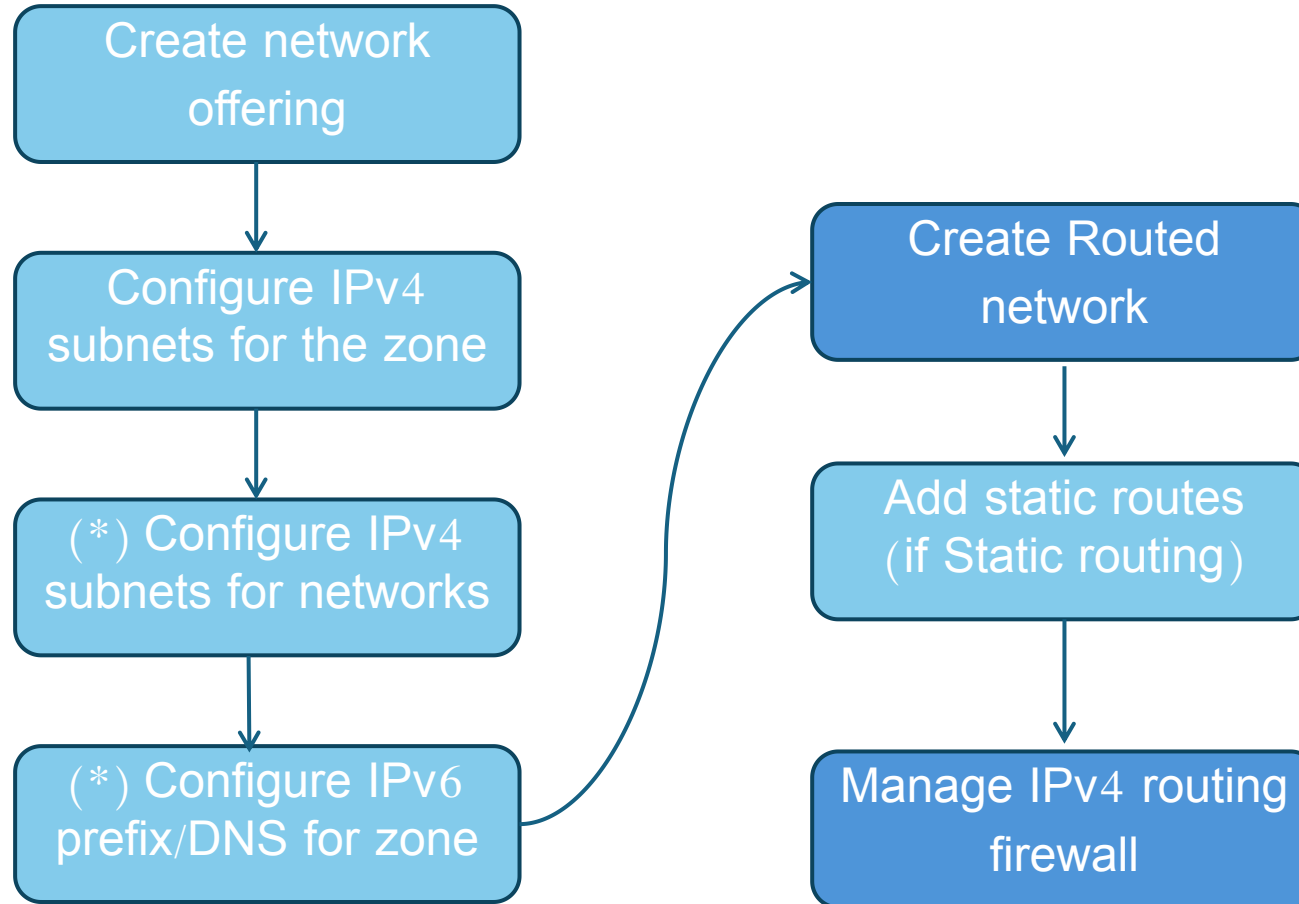Network Mode (i)

ROUTED

Routing mode (i)

Static | Dynamic

Create network offering

↓

Configure IPv4 subnets for the zone

↓

(*) Configure IPv4 subnets for networks

↓

(*) Configure IPv6 prefix/DNS for zone

Performed by operators
Performed by end users

Create Routed network

CIDR size (i)

28

↓

Add static routes (if Static routing)

↓

Manage IPv4 routing firewall

https://docs.cloudstack.apache.org/en/latest/adminguide/networking/dynamic_static_routing.html

(*) This step is **optional**

November 20 - 22, 2024 | Madrid, Spain

**04** How to configure Dynamic routing

# The Internet is a network
# of autonomous systems (AS).

Source

Destination

## Routing:

the process of selecting a
path for a packet

# Border Gateway Protocol (B

- **BGP** is a networking protocol designed to exchange routing data among autonomous systems (AS).

- autonomous system number (ASN):
  - 16 bits or 32 bits
  - public ASN
  - private ASN
    - 64512 to 65534
    - 4200000000 to 4294967294

Network Mode ⓘ

ROUTED

Routing mode ⓘ

Static | **Dynamic**

**Create network offering**

↓

**Configure IPv4 subnets for the zone**

↓

**(*) Configure IPv4 subnets for networks**

↓

**(*) Configure IPv6 prefix for zone**

↓

**Create BGP Peers for zone**

**Create AS number range for zone**

Create AS Range

65000 | 65100

Add BGP peer                                    ✕

* AS Number

64999

IP Address

10.200.0.1

IPv6 IP address

fc00:2024:9:7::1

Password

•••••••••                                        ⦰

Set reservation

⬤○

Cancel | **OK**

(*) This step is **optional**

November 20 - 22, 2024 | Madrid, Spain

# Routing tables in VMs and routers

# What's more: functionalities

- Routed VPC
  - Static and Dynamic routing (*)
  - Network ACL support

- Supports DualStack
  - IPv4
  - IPv4 and IPv6

- Supports Kubernetes Cluster on Routed network and VPC

- **Routed mode is supported by VMware with NSX**

15:40–16:10

# What's more: Performance test

- VM template: ubuntu 24.04 (noble) cloud image
- tool: iperf (TCP port 5001, 10 seconds, max of 5 times)

| | To upstream router | To Shared network | To Isolated network B with PF | To Isolated network B with DNAT | To Isolated network B with Lb | To Routed network D |
|---|---|---|---|---|---|---|
| From upstream router | – | 8.23 Gbps | 5.94 Gbps | 6.56 Gbps | 2.10 Gbps | 4.58 Gbps |
| From Shared network | 7.98 Gbps | 7.31 Gbps | 4.59 Gbps | 5.62 Gbps | 1.97 Gbps | 4.61 Gbps |
| From Isolated network (without Static NAT) | 2.42 Gbps | 4.18 Gbps | 2.57 Gbps | 3.20 Gbps | 1.81 Gbps | 2.24 Gbps |
| From Isolated network (with Static NAT) | 3.08 Gbps | 3.94 Gbps | 2.80 Gbps | 3.22 Gbps | 2.29 Gbps | 1.21 Gbps |
| From Routed network | 3.76 Gbps | 5.04 Gbps | 2.65 Gbps | 4.02 Gbps | 2.37 Gbps | 3.07 Gbps |

Testing, suggestions and ideas are very welcome !

https://github.com/apache/cloudstack/discussions

05

Deep dive: How it works

# Virtual Router (aka VR) for Routed mode

- **Interfaces for Isolated network VR (*)**
  - eth0: guest
  - eth1: control
  - eth2: public/external

- **Enable IP forwarding**

*sysctl net.ipv4.ip_forward=1*

*sysctl net.ipv6.conf.all.forwarding=1*

- **VR routing table**

```
root@r-6-VM:~# route -n
Kernel IP routing table
Destination     Gateway         Genmask         Flags Metric Ref    Use Iface
0.0.0.0         10.200.0.1      0.0.0.0         UG    0      0        0 eth2
10.200.0.0      0.0.0.0         255.255.255.0   U     0      0        0 eth2
169.254.0.0     0.0.0.0         255.255.0.0     U     0      0        0 eth1
202.38.80.16    0.0.0.0         255.255.255.240 U     0      0        0 eth0
202.38.80.32    10.200.0.102    255.255.255.240 UG    20     0        0 eth2
202.38.80.64    10.200.0.103    255.255.255.240 UG    20     0        0 eth2
```



**Publi**

**eth2**

**eth1**

**Gateway of VMs**

**eth0**

**Virtual router**

**VM**

**VM**

# Dynamic routing: FRR

- **FRRouting (FRR)** is a free and open source Internet routing protocol suite for Linux and Unix platforms.

- It implements **BGP**, OSPF, RIP, IS-IS, PIM, LDP, BFD, Babel, PBR, OpenFabric and VRRP, with alpha support for EIGRP and NHRP.

- FRR version for CloudStack 4.20.0: **8.4.4-1.1~deb12u1**

```
frr version 6.0
frr defaults traditional
hostname r-6-VM
service integrated-vtysh-config
ip nht resolve-via-default
router bgp 65053
 bgp router-id 10.200.0.101
 bgp default ipv6-unicast
neighbor 10.200.0.1 remote-as 64999
neighbor 10.200.0.1 password password3
neighbor fc00:2024:9:7::1 remote-as 64999
neighbor fc00:2024:9:7::1 password password3
address-family ipv4 unicast
  network 202.38.80.16/28
exit-address-family
address-family ipv6 unicast
  network 2a02:1810:248b:3b0a::/64
exit-address-family
line vty
```

# Dynamic routing: FRR status

- "vtysh"

```
r-6-VM# show bgp summary

IPv4 Unicast Summary (VRF default):
BGP router identifier 10.200.0.101, local AS number 65053 vrf-id 0
BGP table version 2
RIB entries 3, using 576 bytes of memory
Peers 2, using 1448 KiB of memory

Neighbor         V          AS   MsgRcvd   MsgSent   TblVer   InQ OutQ  Up/Down State/PfxRcd   PfxSnt Desc
10.200.0.1       4       64999        11        11        0     0    0 00:02:18           1        2 N/A
fc00:2024:9:7::1 4       64999         0         0        0     0    0    never     Connect        0 N/A

Total number of neighbors 2

IPv6 Unicast Summary (VRF default):
BGP router identifier 10.200.0.101, local AS number 65053 vrf-id 0
BGP table version 3
RIB entries 5, using 960 bytes of memory
Peers 2, using 1448 KiB of memory

Neighbor         V          AS   MsgRcvd   MsgSent   TblVer   InQ OutQ  Up/Down State/PfxRcd   PfxSnt Desc
10.200.0.1       4       64999        11        11        0     0    0 00:02:18           2        3 N/A
fc00:2024:9:7::1 4       64999         0         0        0     0    0    never     Connect        0 N/A

Total number of neighbors 2
r-6-VM# 
```

# IPv4 Routing firewall

- Implemented via "**nftables**"

| IP family | Table | Chain | Notes |
|-----------|-------|-------|-------|
| ip | ip4_firewall | INPUT<br>FORWARD<br>OUTPUT<br>fw_chain_egress<br>fw_chain_ingress | IPv4 routing firewall for Isolated networks |
| | ip4_acl | INPUT<br>FORWARD<br>OUTPUT<br>eth2_egress_policy<br>eth2_ingress_policy<br>eth3_egress_policy<br>eth3_ingress_policy | Network ACL for VPC tier 001<br><br>Network ACL for VPC tier 002 |

# IPv6 firewall

- Implemented via "**nftables**"

| IP family | Table | Chain | Notes |
|-----------|-------|-------|-------|
| ip6 | ip6_firewall | fw_input<br>fw_forward<br>fw_chain_egress<br>fw_chain_ingress | IPv6 firewall for Isolated networks |
| | ip6_acl | acl_input<br>acl_forward<br>eth2_egress_policy<br>eth2_ingress_policy<br>eth3_egress_policy<br>eth3_ingress_policy | Network ACL for VPC tier 001<br><br>Network ACL for VPC tier 002 |

# 06

Summary and Future work

# Summary

- New network mode: ROUTED

- Routing modes: Static and Dynamic

- Routed Network and Routed VPC with access control

# Future work

- **Static Routes**
  - Existing feature in ACS
  - Currently supports only VPC private gateway

  - Extend to support Isolated networks and VPC
    - NATTED and ROUTED mode
  - Next hop could be a router, a vm, or an IP in public/guest/private gateway network

  - Coming soon...

# Future work

- Support more FRR customizations
  - EBGP_MultiHop is supported

- Redundant VRs for Routed networks

- Routed Shared network

- Internal LB on Routed and Shared networks

# Acknowledgement

**Many thanks to**

- **Alex Mattioli**
  - proposal and high-level design

- **Nicolas Vazquez** and **Pearl D'silva**
  - AS number management
  - Routed mode Implementation in NSX integration

- **Kiran Chavala**
  - QA testing

# Q & A

#CSCollab24
@CloudStack

# Thank you!

#CSCollab24
@CloudStack