

LINE Data Center Networking with SRv6

Hirofumi Ichihara
co-author: Toshiaki Tsuchiya
LINE corporation

About Me

- Hirofumi Ichihara
- LINE Corporation
 - Network Development Team
- Network Software Developer
 - SDN/NFV
 - OpenStack Neutron
 - Docker
 - Kubernetes

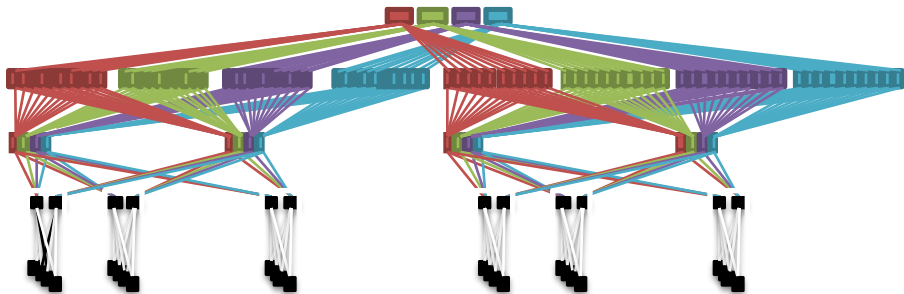
LINE Services and Networks

Full L3 CLOS Network*

- Single tenant network
- LINE message service and related services running



...



* Excitingly simple multi-path OpenStack networking: LAG-less, L2-less, yet fully redundant
<https://www.slideshare.net/linecorp/excitingly-simple-multi-path-openstack-networking-lag-less-l2-less-yet-fully-redundant>

Exclusive Network for Services

- Service with specific requirements running
- Building specific network for each service



LINE Pay

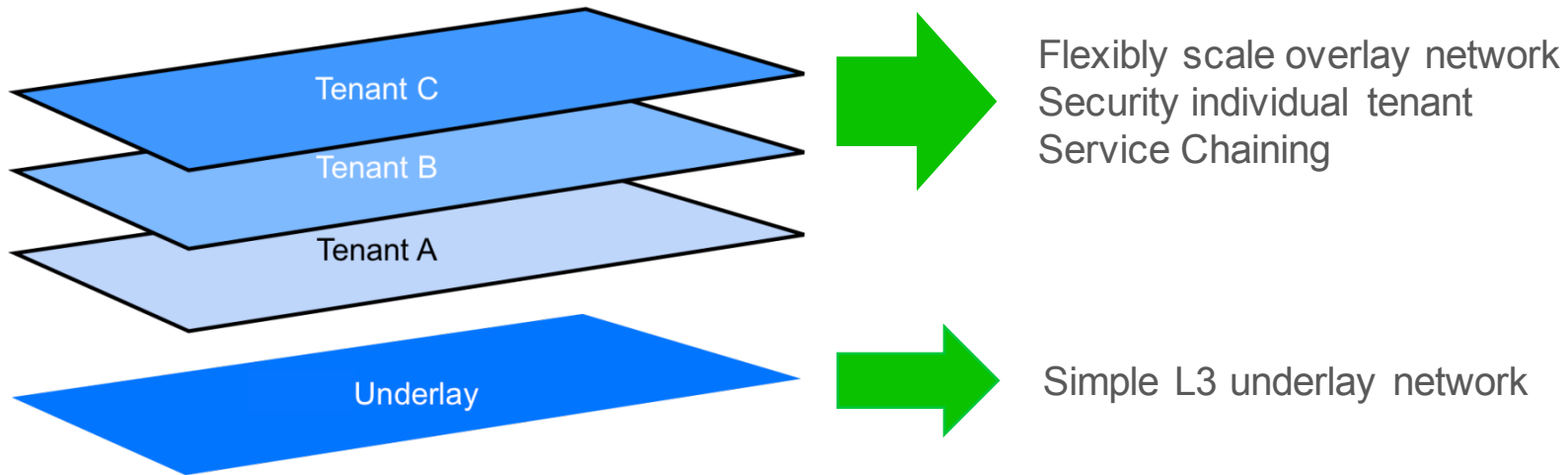
Other: Fintech Business



**Many fragment underlay networks
Many works to design and build
Management cost increases**

Multi tenant network

- Sharing underlay network decrease management cost
- Achieve policy for each service(tenant) on overly network



Multi tenancy

VXLAN

Pros

- More information
- Many network devices support

Cons

- Lose advances of full-L3
- Need additional protocol to achieve service

IPv6 Segment Routing (SRv6)

Pros

- IPv6 forwarding only on underlay
- Support segregation and service chaining with Segment ID

Cons

- No information about DC use case
- No network device support

+ SRv6 future

Adopted SRv6

SRv6

Segment ID (SID)

- 128bit number(IPv6 address)
- **Locator**: Information for routing to SRv6 node(parent node). It must be unique within a SR domain
- **Function**: Information to identify the action to be performed on the parent node



Segment Routing Header (SRH)

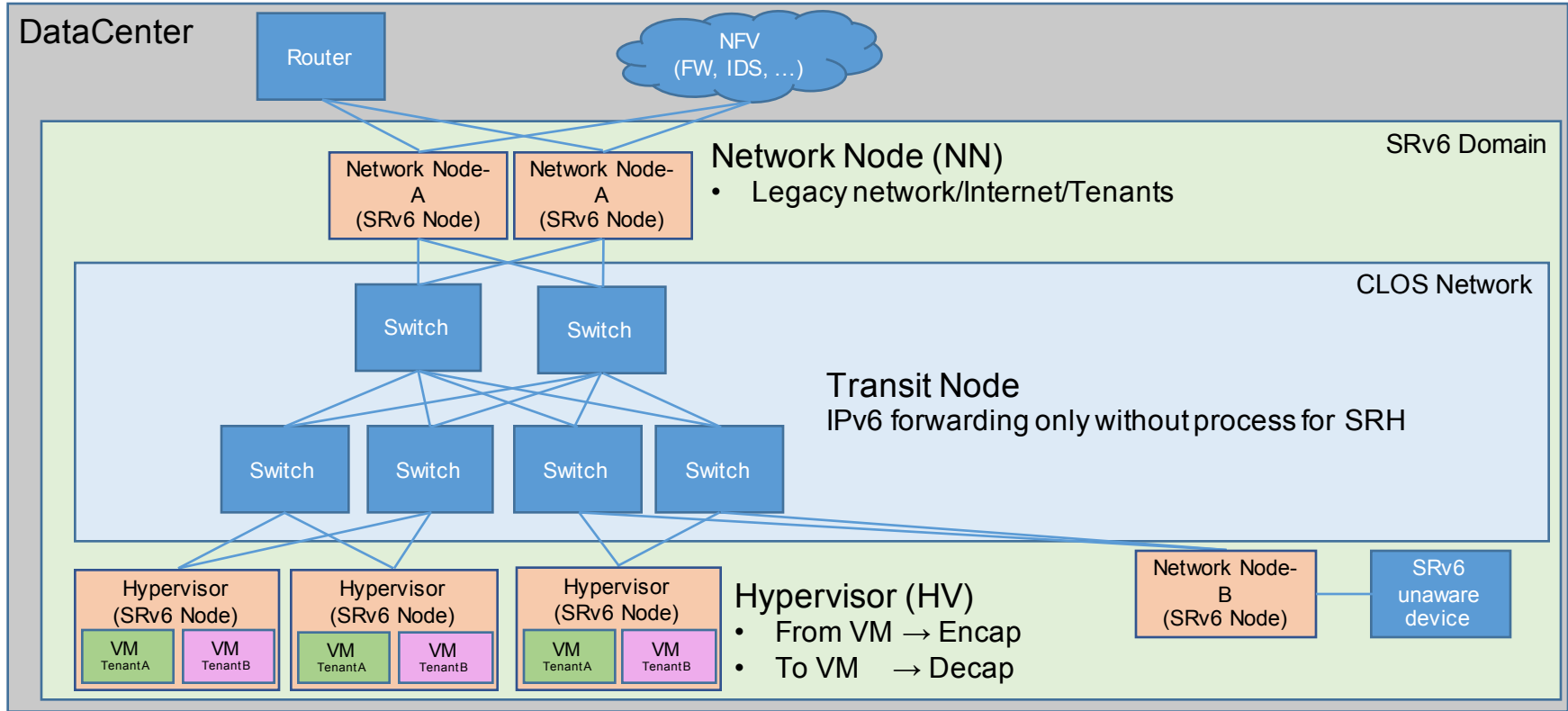
- IPv6 extension header
- Including a Segment List, Segment Left points out current point of Segment List and so on

Function examples

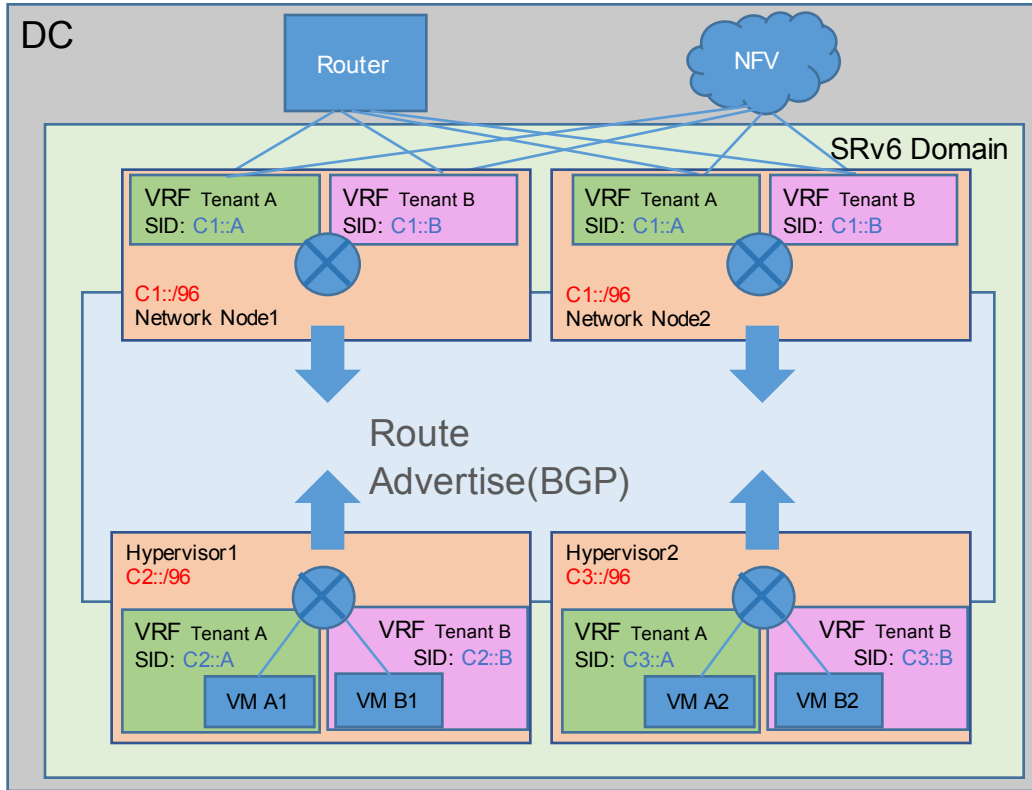
- **T.Encaps(Encap)**: Encapsulation packet with IPv6 header and SRH
- **End.DX4(Decap)**: Remove IPv6 header and SRH from packet and then forward next hop
- **End.DT4(Decap)**: Remove IPv6 header and SRH from packet and then lookup routing table and forward (DT4 is not implemented in Linux Kernel so we used DX4 although DT4 is better)

SRv6 Data Center Network Data Plane

Data Plane - Architecture



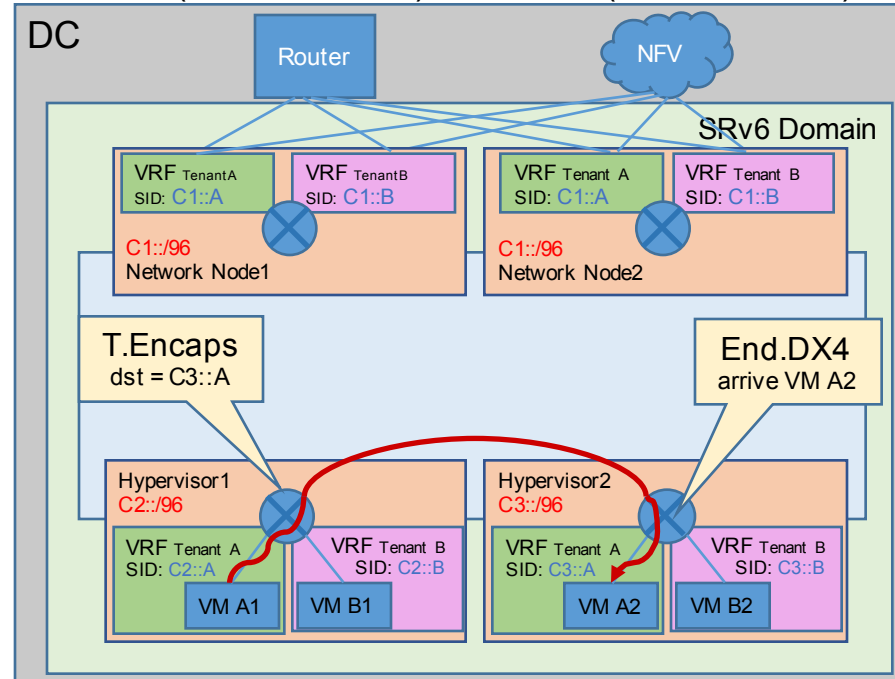
Data Plane - SID, Routing



- Create VRF (I3master device) for each tenant on NetworkNode, Hypervisor
- Assign IPv6 address /96 block (**Locator**) to nodes(NetworkNode, Hypervisor)
- Add identifier for each tenant to the Locator as **Function** (LINE uses specific address from 169.254.0.0/16 each tenant)
- Advertise /96 IPv6 address(Locator) via BGP

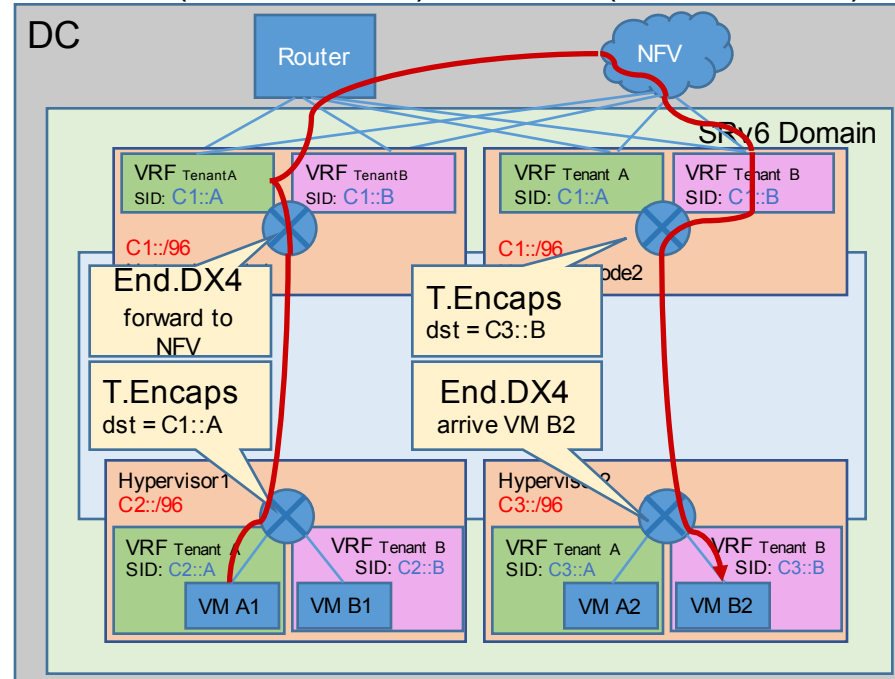
Data Plane - Packet flow in a tenant

VM A1 (HV1 TenantA) → VM A2 (HV2 TenantA)



Data Plane - Packet flow between tenants

VM A1 (HV1 TenantA) → VM B2 (HV2 Tenant B)



Data Plane - Real config on Network Node

Encap

```
[NetworkNode]# ip route show table 12
```

		Segment List	
10.122.12.113	encap seg6 mode encap segs 1	[2400:dcc0::a7a:4d8d:a9fe:108]	dev vrf5c0594737b87 scope link
10.122.12.114	encap seg6 mode encap segs 1	[2400:dcc0::a7a:4d8e:a9fe:108]	dev vrf5c0594737b87 scope link
10.122.12.115	encap seg6 mode encap segs 1	[2400:dcc0::a7a:4d8f:a9fe:108]	dev vrf5c0594737b87 scope link

Destination IPv4 address of VM Locator(HV Address) Function (IPv4 address to identify each tenant)

Decap

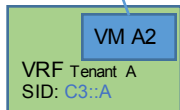
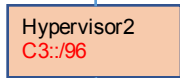
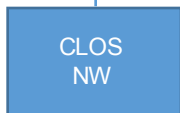
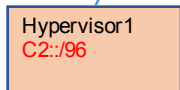
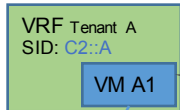
```
[NetworkNode]# ip -6 route show table local
```

local	2400:dcc0::a7a:4d87:a9fe:102	encap seg6local	action End.DX4	nh4	169.254.1.2	dev vrf01b1db9dd10f	metric 1024	pref medium	
local	2400:dcc0::a7a:4d87:a9fe:104	encap seg6local	action End.DX4	nh4	169.254.1.4	dev vrf01b1db7f5d2b	metric 1024	pref medium	
local	2400:dcc0::a7a:4d87:a9fe:108	encap seg6local	action End.DX4	nh4	169.254.1.8	dev vrf5c0594737b87	metric 1024	pref medium	

... Locator(HV Address) Function (Tenant identifier) IPv4 address to identify each tenant. They are assigned to VRF IF (That is magic to lookup VRF with End.DX4)

Data Plane - Real behavior

IP: 10.122.12.36



IP: 10.122.12.35

```
[VM-A1]$ ping 10.122.12.35 -c 10
PING 10.122.12.35 (10.122.12.35) 56(84) bytes of data.
64 bytes from 10.122.12.35: icmp_seq=1 ttl=63 time=0.356 ms
64 bytes from 10.122.12.35: icmp_seq=2 ttl=63 time=0.461 ms
...
64 bytes from 10.122.12.35: icmp_seq=10 ttl=63 time=0.415 ms
--- 10.122.12.35 ping statistics ---
10 packets transmitted, 10 received, 0% packet loss, time 9000ms
```

HV2: Decap

Remove IPv6, SR header

```
▶ Frame 5: 98 bytes on wire (784 bits), 98 bytes captured (784 bits)
▶ Ethernet II, Src: 7e:55:34:6d:37:3a (7e:55:34:6d:37:3a), Dst: fa:16:
▼ Internet Protocol Version 4, Src: 10.122.12.36, Dst: 10.122.12.35
  0100 .... = Version: 4
  .... 0101 = Header Length: 20 bytes (5)
▶ Differentiated Services Field: 0x00 (DSCP: CS0, ECN: Not-ECT)
  Total Length: 84
  Identification: 0x0000 (0)
▶ Flags: 0x4000, Don't fragment
  Time to live: 63
  Protocol: ICMP (1)
  Header checksum: 0x0e6f [validation disabled]
  [Header checksum status: Unverified]
  Source: 10.122.12.36
  Destination: 10.122.12.35
▶ Internet Control Message Protocol
```

HV1: Encap

Insert IPv6, SR header

```
▼ Internet Protocol Version 6, Src: 2400:dcc0::a7a:4d8e:a9fe:102, Dst:
  0110 .... = Version: 6
  ▶ .... 0000 0000 .... = Traffic Class: 0x00 (D
  .... 0000 0000 0000 0000 = Flow Label: 0x0000
  Payload Length: 108
  Next Header: Routing Header for IPv6 (43)
  Hop Limit: 63
  Source: 2400:dcc0::a7a:4d8e:a9fe:102
  Destination: 2400:dcc0::a7a:4d8f:a9fe:102
▼ Routing Header for IPv6 (Segment Routing)
  Next Header: IPIP (4)
  Length: 2
  [Length: 24 bytes]
  Type: Segment Routing (4)
  Segments Left: 0
  First segment: 0
  ▶ Flags: 0x00
  Reserved: 0000
  Address[0]: 2400:dcc0::a7a:4d8f:a9fe:102
  ▼ [Segments in Traversal Order]
    Address[0]: 2400:dcc0::a7a:4d8f:a9fe:102
  ▶ Internet Protocol Version 4, Src: 10.122.12.36, Dst: 10.122.12.35
  ▶ Internet Control Message Protocol
```

SRv6 Data Center Network Control Plane

SRv6 Control Plane Choices

- ISIS
- OSPF
- BGP
- SDN Controller

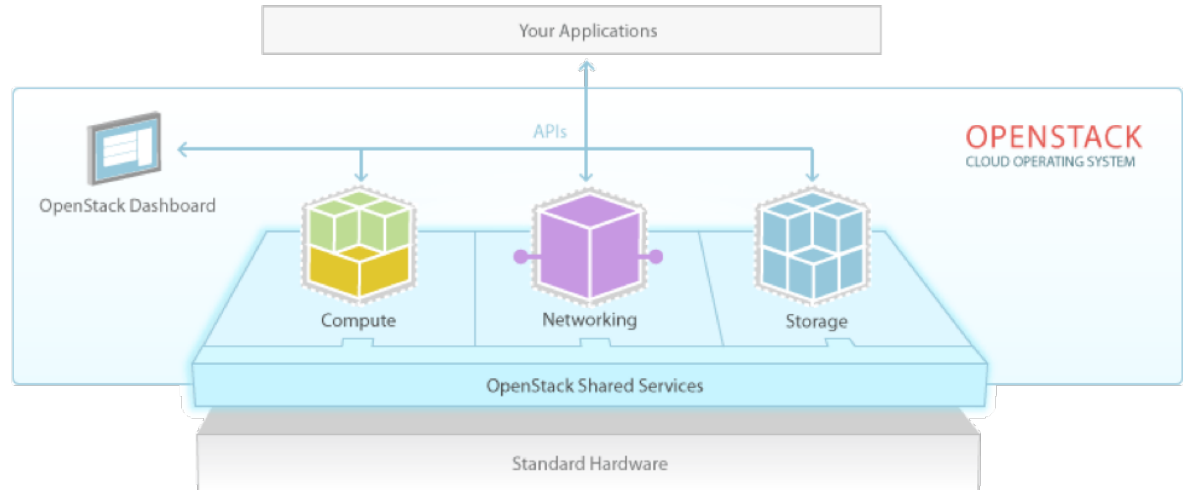
LINE uses OpenStack as Private
Cloud Controller so adopted
SDN Controller

OpenStack

- Cloud Operating system
- Support Multi Hypervisor
- Support various SDN controllers and Storage appliances



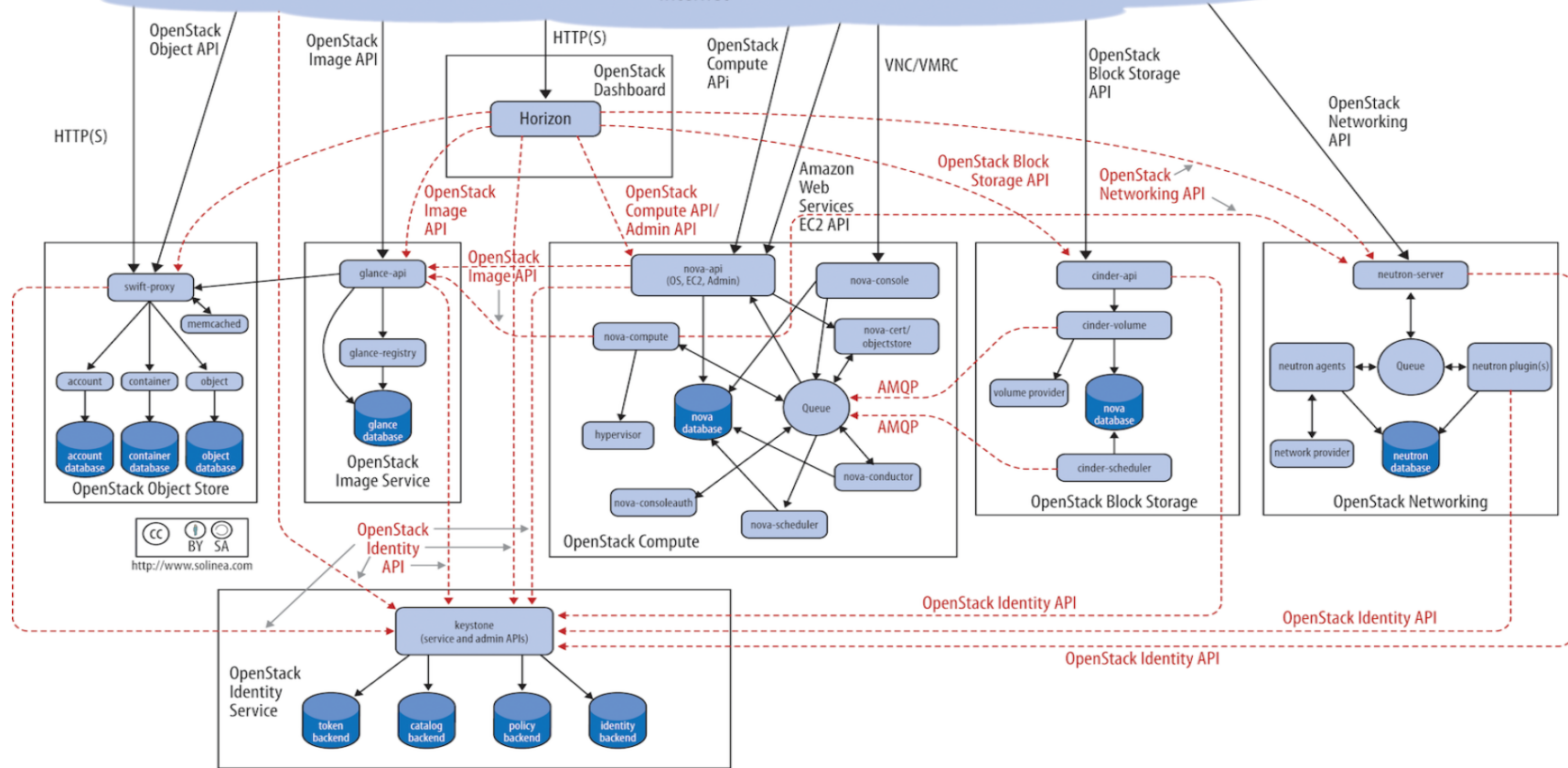
openstack®





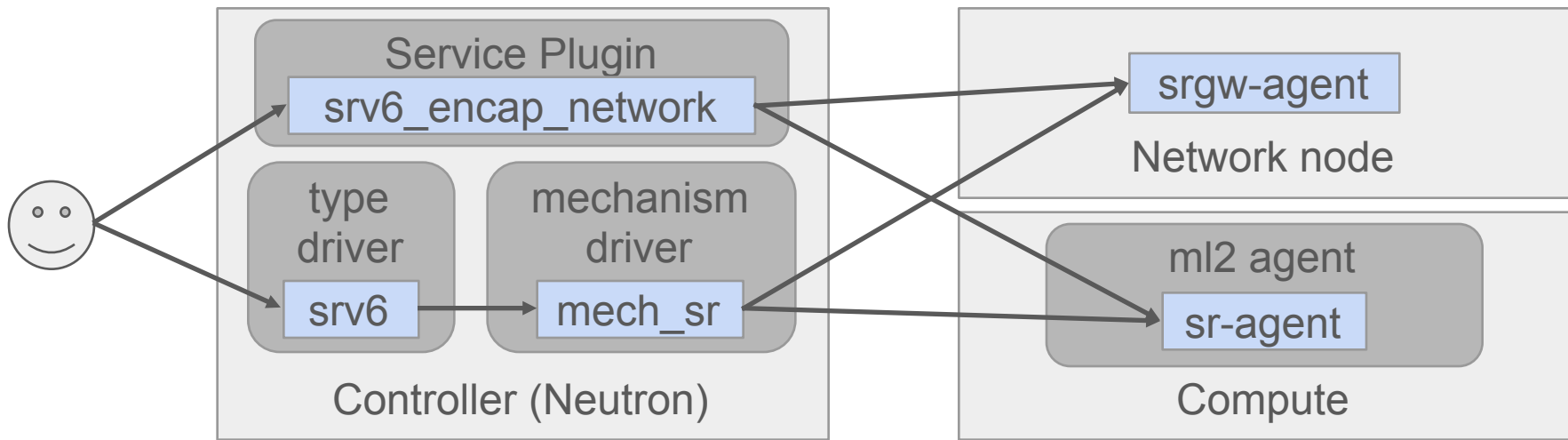
- Command-line interfaces (nova, neutron, swift, etc)
- Cloud Management Tools (Rightscale, Enstratus, etc)
- GUI tools (Dashboard, Cyberduck, iPhone client, etc)

Internet



Neutron SRv6 Plugin - networking-sr

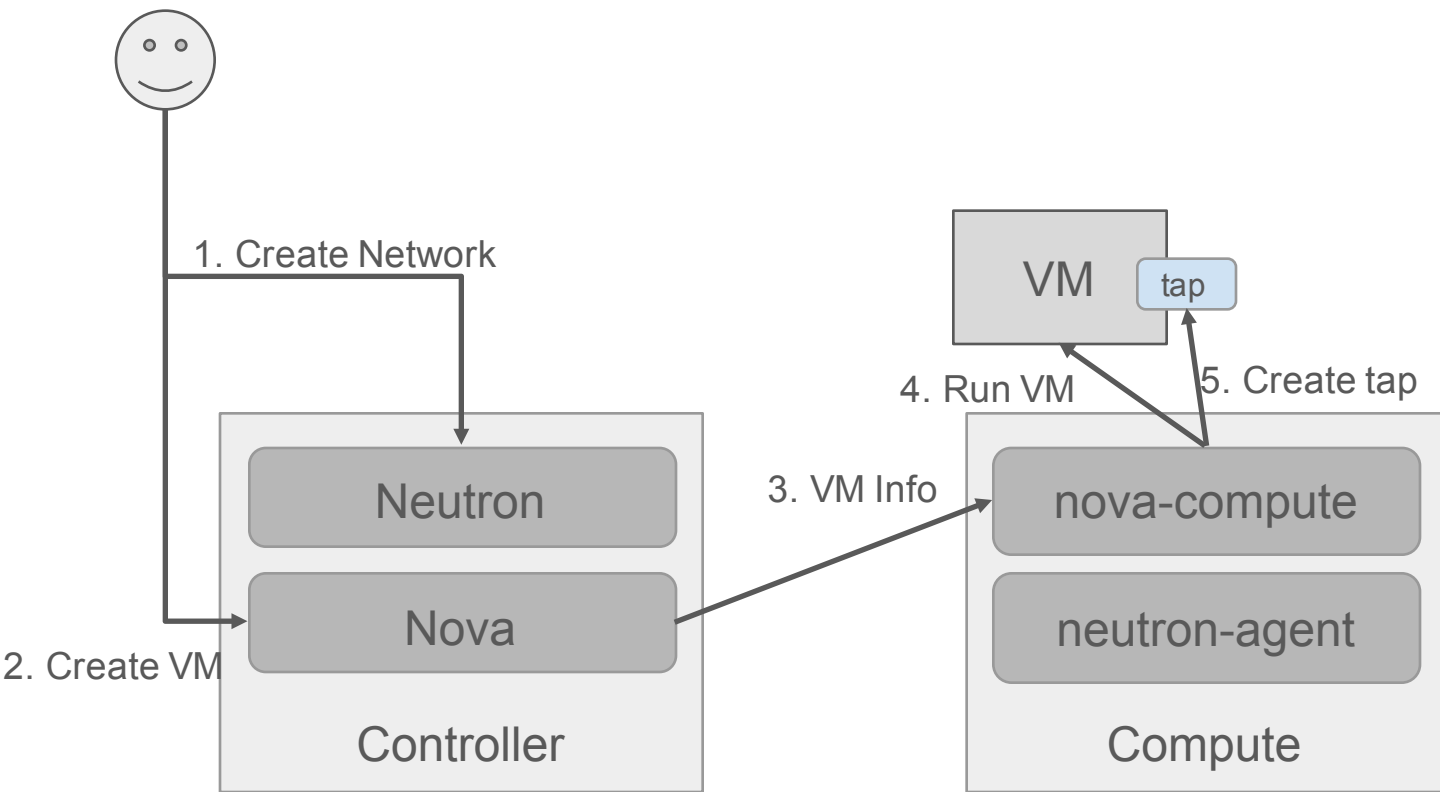
- ML2 mechanism/type driver and agent
- Gateway agent on network nodes
- Service plugin for new API to add SRv6 encap rule



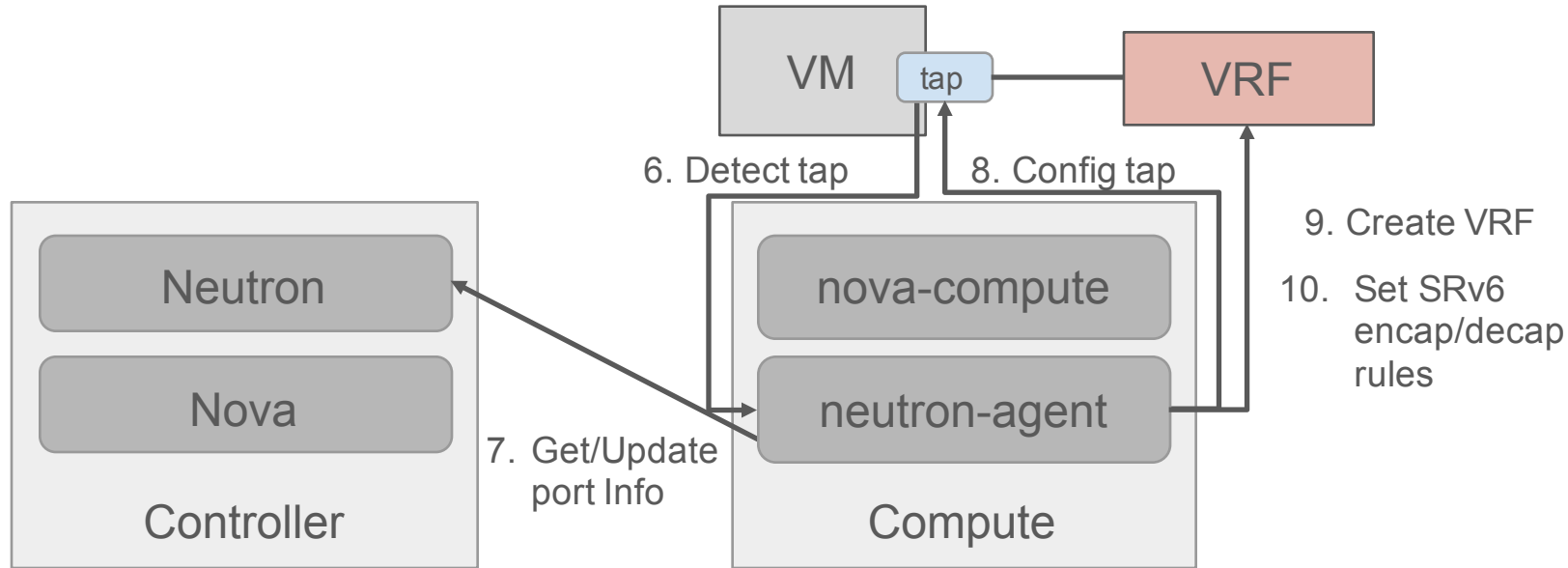
SRv6 Data Center Network Control Plane

ML2 mechanism/type driver and agent

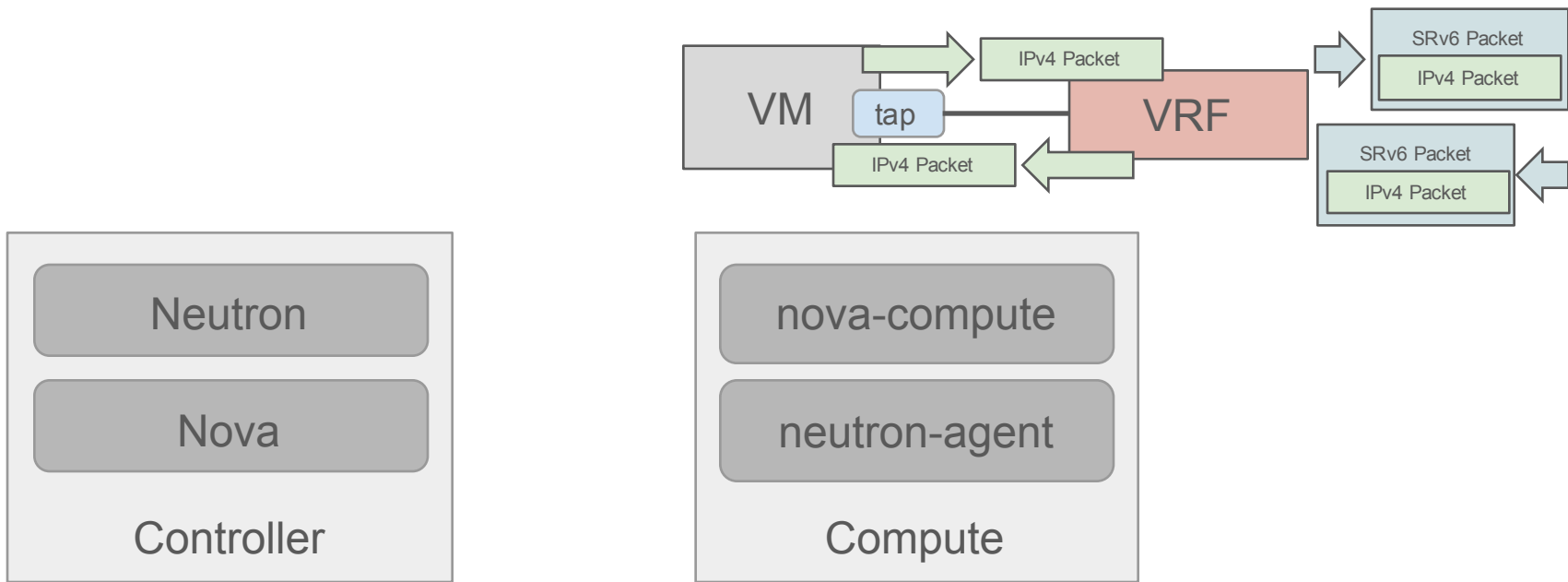
Nova, Neutron Behavior - VM create



Nova, Neutron Behavior - Network configuration



Packets for VM encap/decap on VRF



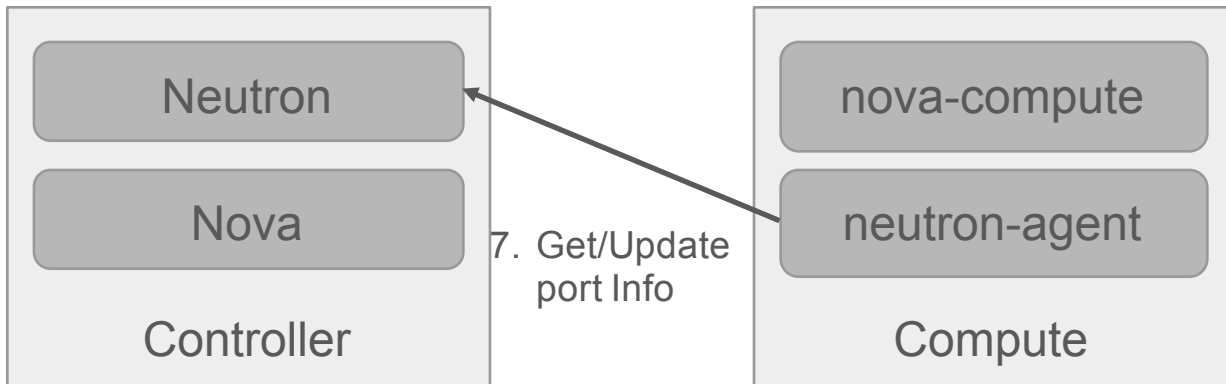
How does sr-agent get VRF info?

Virtual Machine Configuration

1. Create network
2. Create VM
3. Notify VM info
4. Run VM
5. Create tap

Network Configuration

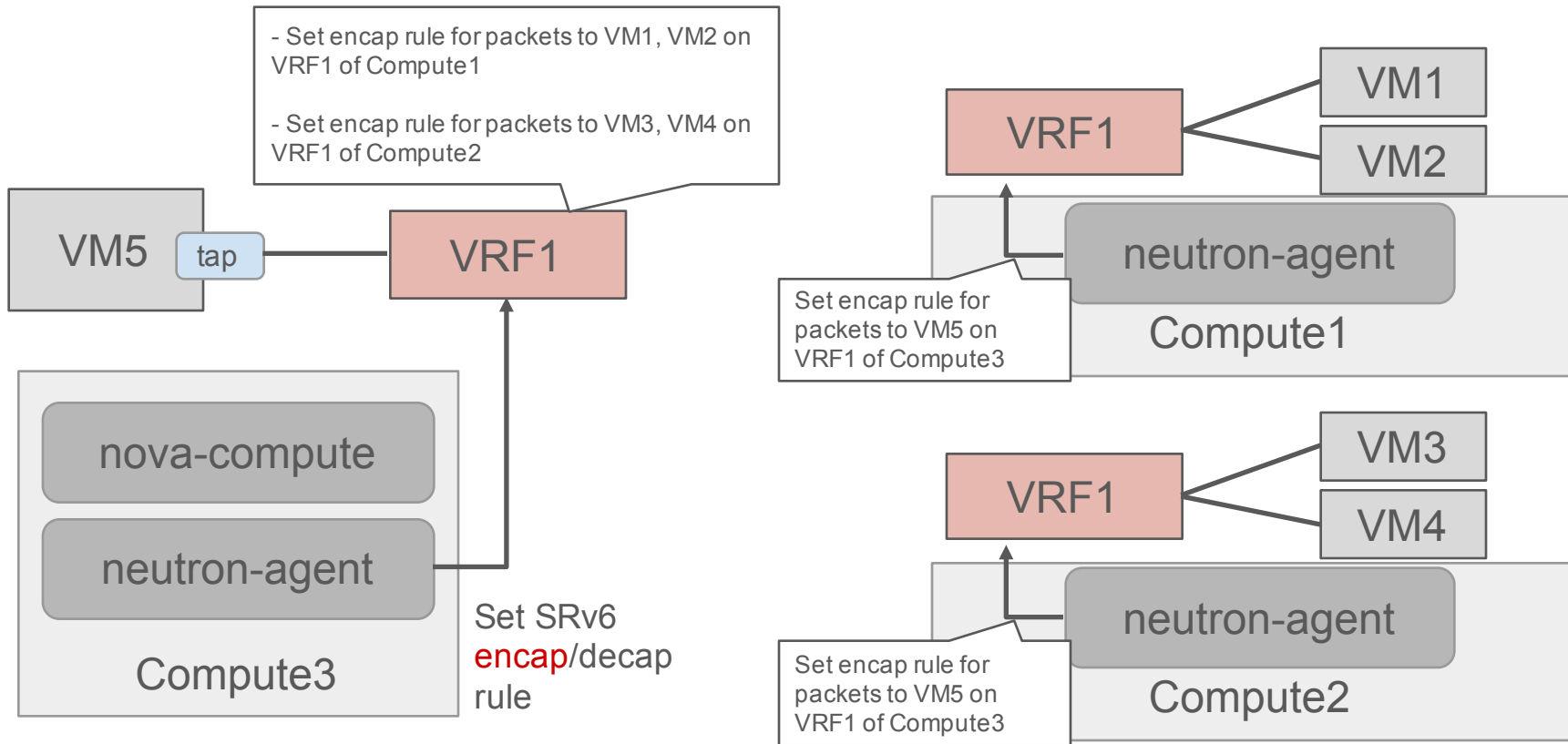
6. Detect tap
7. Update/Get port info
8. Config tap
9. Create VRF
10. Set SRv6 encap/decap rules



VRF info in Port binding:profile

```
{  
  "port":{  
    "binding:profile":{  
      "segment_node_id": "2400:dcc0::a7a:4d8e", # Locator(Hypervisor address) where VM with the port running  
      "vrf": "vrf644606a29039", # VRF IF name for the port. The name is combined by "vrf" + tenant_id + network_id  
      "vrf_cidr": "169.254.1.0/24", # IP CIDR of VRF for the port  
      "vrf_ip": "169.254.1.44" # IP Address of VRF for the port  
    }  
  }  
}
```

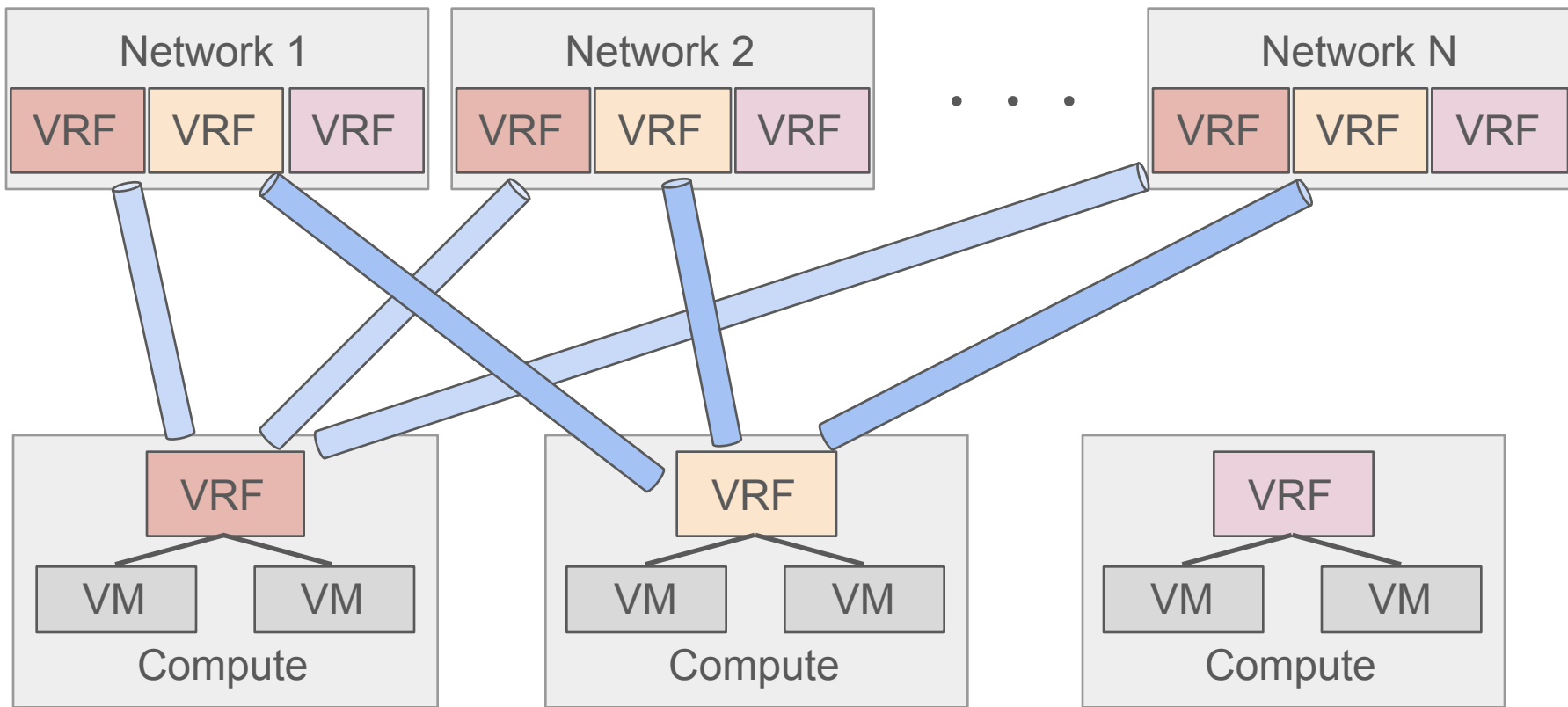

Set encap rule from Port info of each VM



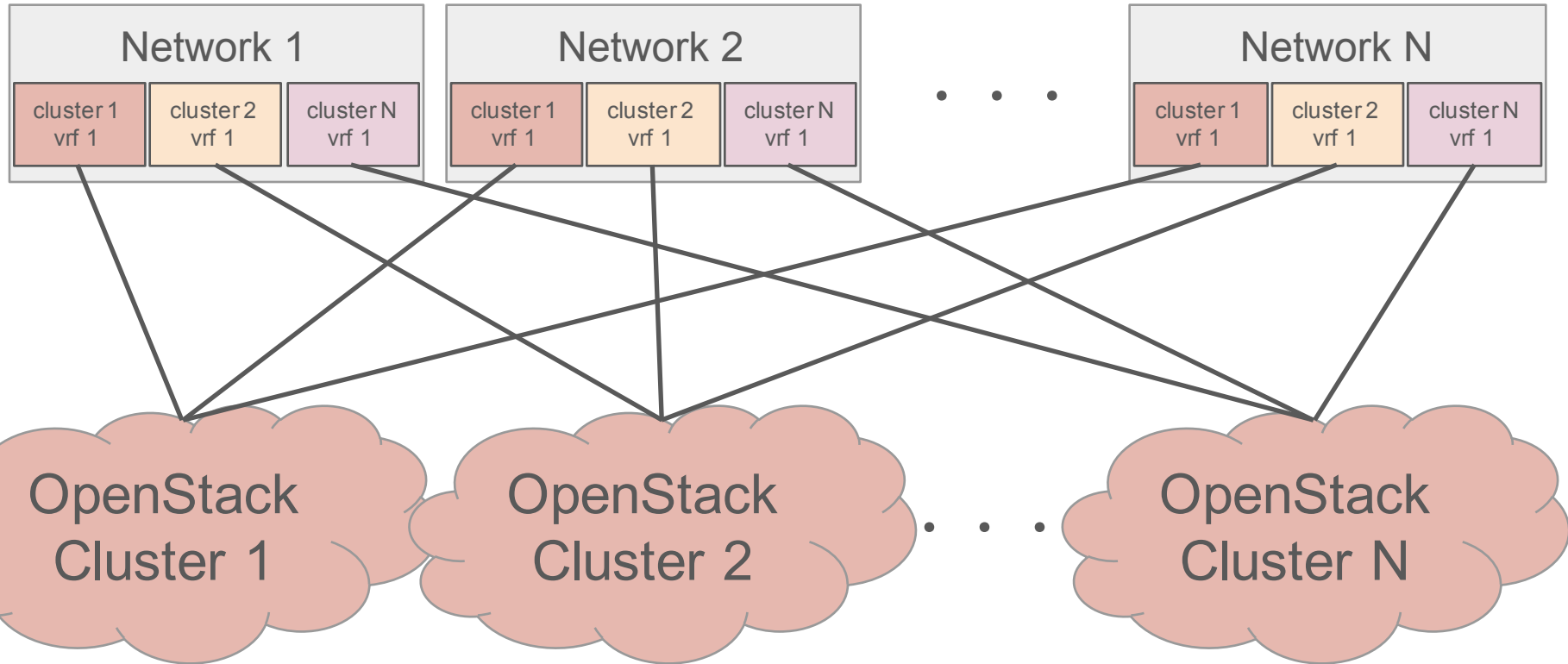
SRv6 Data Center Network Control Plane

Gateway agent on network nodes

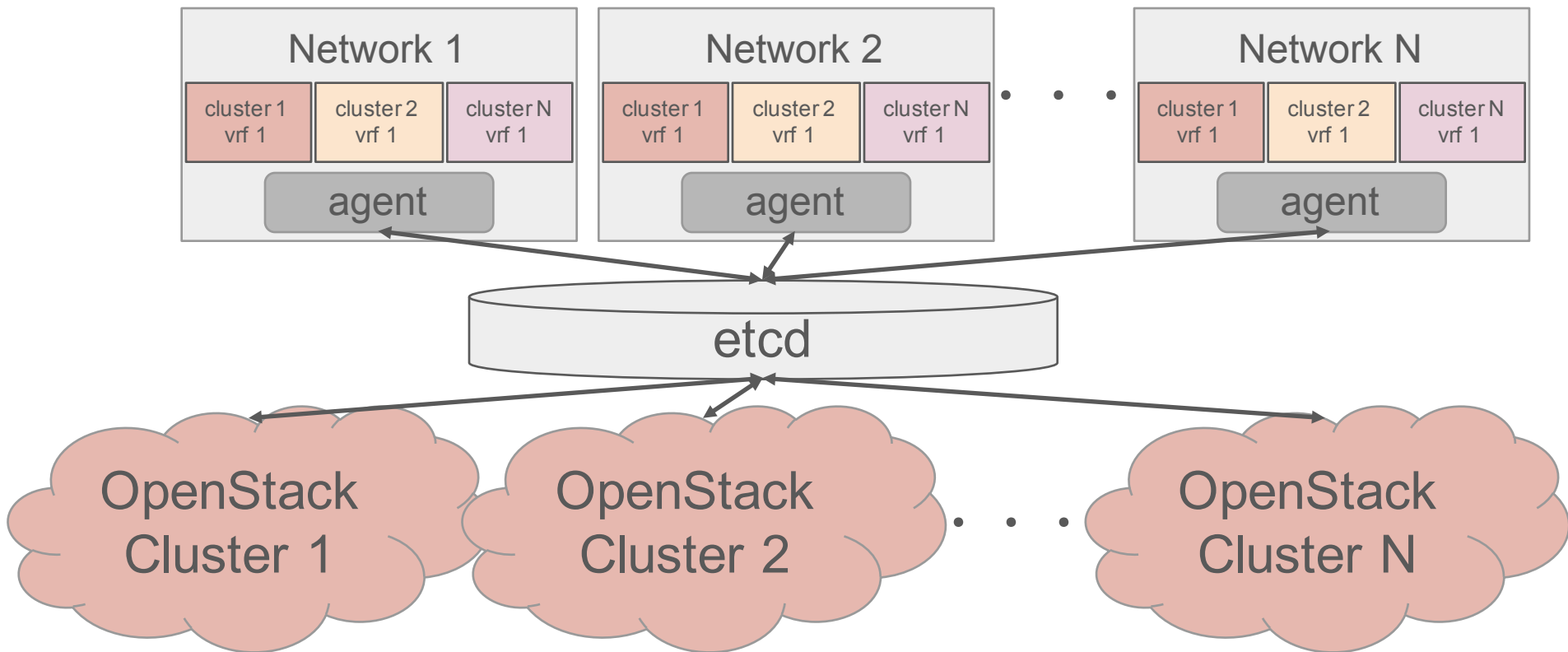
Network Node Requirements: Scale



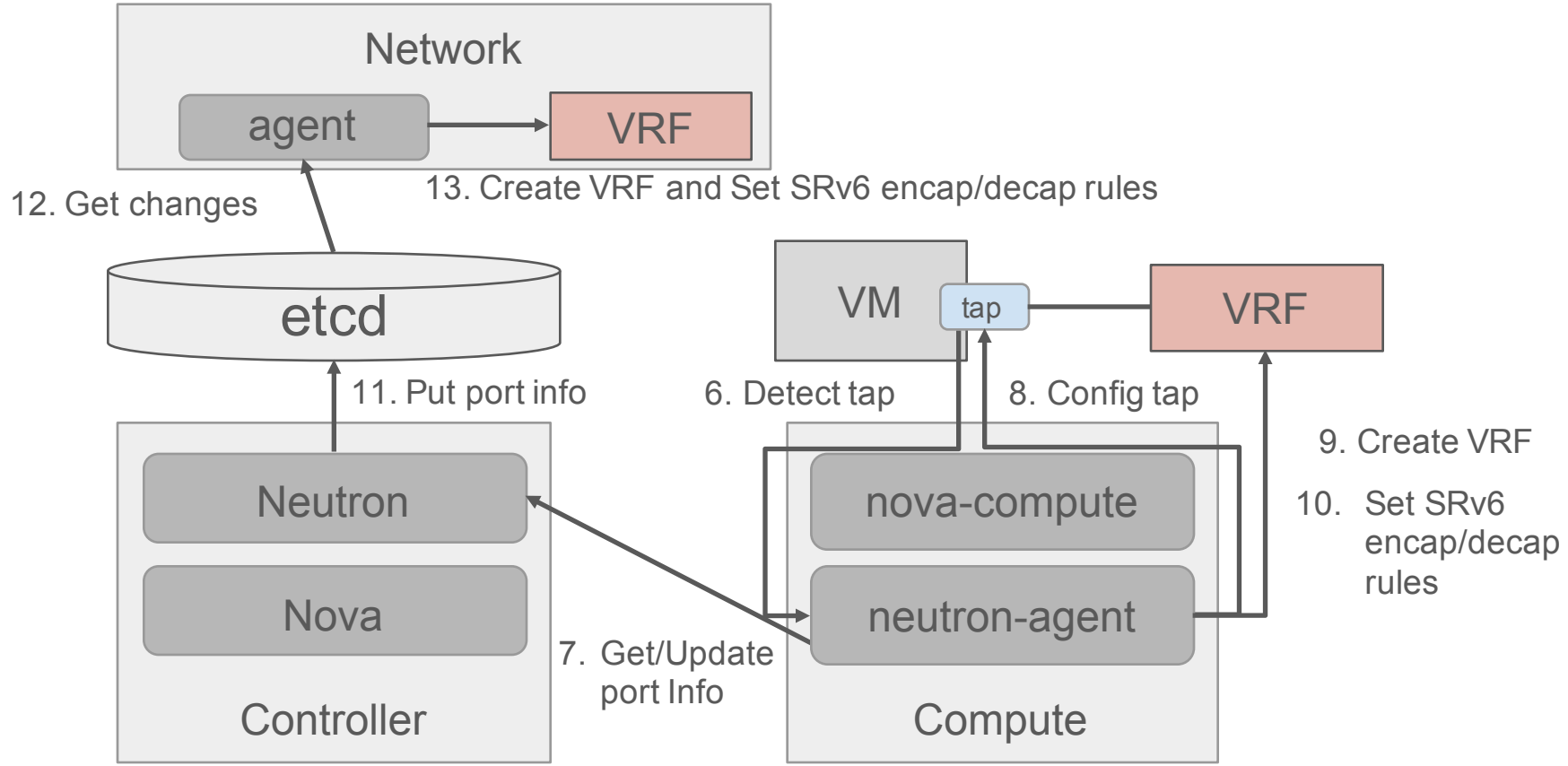
Network Node Requirements: Multi clusters



Etcd + Agent Model



Notify New Encap/Decap Rule via Etcd



SRv6 Data Center Network Control Plane

**Service plugin for new API to add
SRv6 encap rule**

srv6_encap_network API

SRv6 Encap Network

SRv6 Encap Network

The **srv6_encap_network** extension lists, creates, shows information for, and updates srv6_encap_network resource.

GET

/v2.0/srv6_encap_networks
List srv6 encap networks

POST

/v2.0/srv6_encap_networks
Create srv6 encap network

GET

/v2.0/srv6_encap_networks/ {srv6_encap_network_id}
Show srv6 encap network

PUT

/v2.0/srv6_encap_networks/ {srv6_encap_network_id}
Update srv6 encap network

DELETE

/v2.0/srv6_encap_networks/ {srv6_encap_network_id}
Delete srv6 encap network

```
{
  "srv6_encap_networks": [
    {
      "network_id": "fbc5f08e-0cb0-4b5c-a5ce-ac7032f50c7b",
      "tenant_id": "d988b205c6e142669a290dd80010587c",
      "project_id": "d988b205c6e142669a290dd80010587c",
      "encap_rules": [
        {
          "nexthop": "fc00:17::a00:fa",
          "destination": "10.0.201.200"
        }
      ],
      "id": "43938fab-ce22-442f-b537-24f2768de773"
    },
    {
      "network_id": "d76c20be-5c2a-40c5-bbd5-0b192fa3ff9c",
      "tenant_id": "aac15739c8034f60b2e8278e84563919",
      "project_id": "aac15739c8034f60b2e8278e84563919",
      "encap_rules": [
        {
          "nexthop": "fc00:17::a00:fa",
          "destination": "10.0.201.201"
        }
      ],
      "id": "70989e81-eae4-490c-b016-f665e6bc872f"
    }
  ]
}
```


srv6_encap_network resource

- id: Identifier for resource
- tenant_id/project_id: Identifier for project/tenant of resource
- network_id: Identifier of network which resource is assigned
- encap_rules: SRv6 encap rule list
 - destination: IPv4 address for specific destination of packet
 - nexthop: SID packets should be encapsulated

NFV(LBaaS) and networking-sr with new API



1. Create VIP

LBaaS

Network

VRF1

agent

4. Set SRv6
encap rule

tenant_id: Tenant User belongs
network_id: Network VM connects
encap_rules: destination is VIP, nexthop is SID of
VRF1 on Network node

3. Notify encap rule

Neutron
Controller

2. Add encap rule for VIP by
srv6_encap_netowrk API

VIP encap seg6 mode encap segs NetworkNode_VRF1_SID

VM1

tap

VRF1

nova-compute

neutron-agent

Compute

Summary

- SRv6 network for data center use case
 - Multi tenant networks
- Data plane architecture
 - SRv6 Encap/Decap support on Hypervisors and Network nodes
 - End.DX4 + Routing to VRF (Kernel doesn't have End.DT4)
- Control plane architecture
 - OpenStack Neutron SRv6 plugin networking-sr
 - Gateway agent with etcd for large scale
 - New API to add SRv6 encap rule