



# **Cloud Networking (VITMMA02)**

## **Network Virtualization: Overlay Networks**

### **OpenStack Neutron Networking**

Markosz Maliosz PhD

Department of Telecommunications and Media Informatics  
Faculty of Electrical Engineering and Informatics  
Budapest University of Technology and Economics

Spring 2019



# OVERLAY NETWORKS

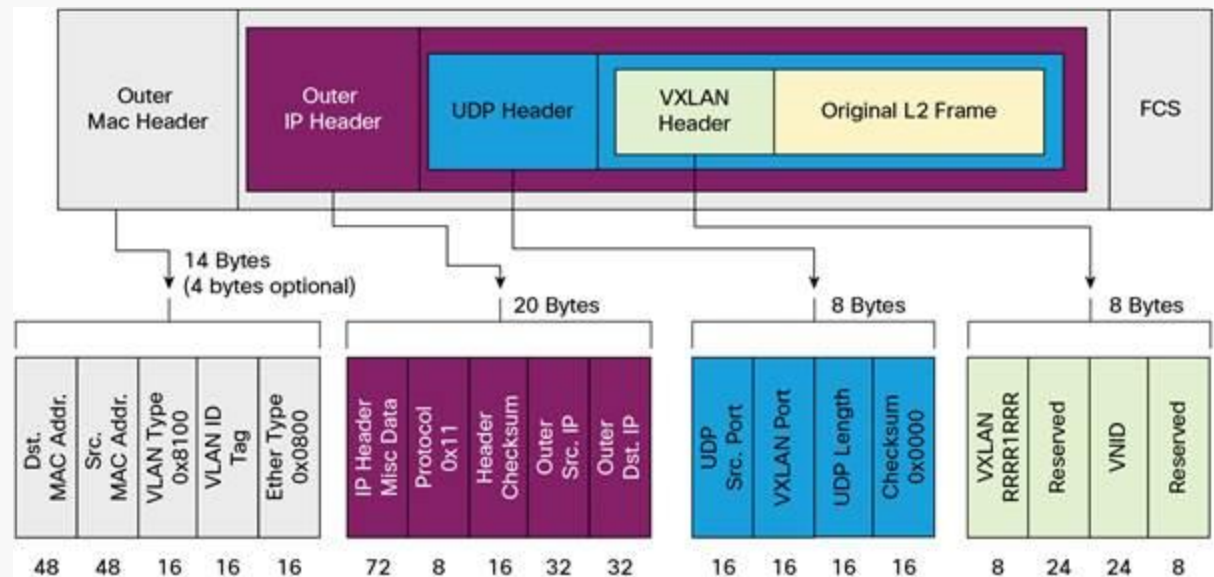


# Network Virtualization

- » Support for tenant separation
  - » Virtual Extensible LAN (VXLAN) – RFC 7348
    - » Cisco, VMware
    - » transport of virtual L2 traffic over physical L3 network
  - » Network Virtualization using Generic Routing Encapsulation (NVGRE)
    - » Microsoft, Intel, HP, Dell
  - » Generic Network Virtualization Encapsulation (GENEVE)
    - » superset of VXLAN and NVGRE
  - » Stateless Transport Tunneling (STT)
    - » Nicira ⇔ VMware

# VXLAN

- » Original L2 frame of tenant
  - » original MAC address and VLAN tag
- » MAC-in-UDP
- » VXLAN and UDP header
  - » VXLAN network ID (VNID) – identifying the tenant
    - » 24 bit  $\Rightarrow$  16 million tenant
- » physical network: IP routing (Layer3)

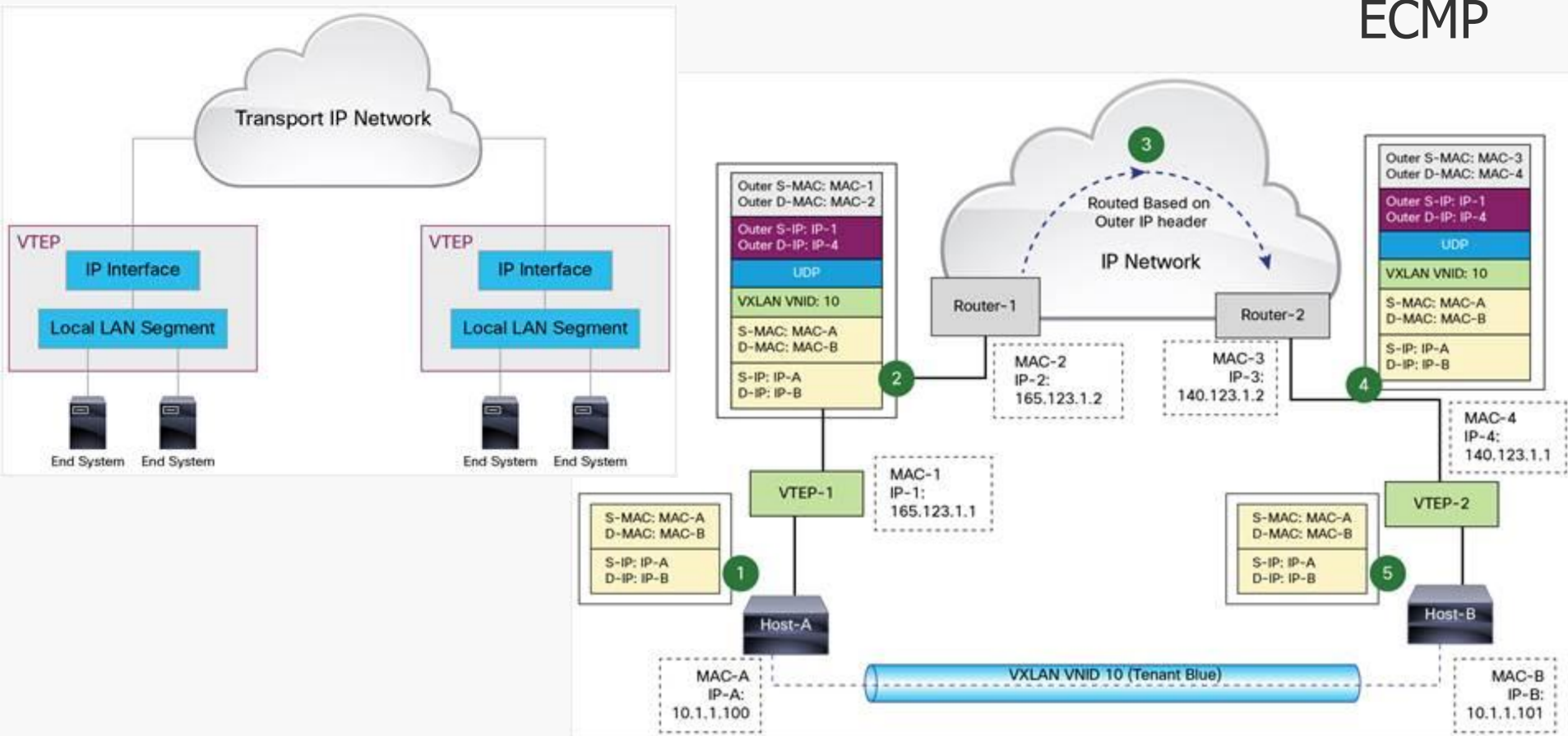




# VXLAN

- » VXLAN Tunnel End Point (VTEP)
- » MAC-to-VTEP tables by learning (IP multicast)
  - » all VTEP of a VNI in a multicast group

ECMP



# NVGRE

- » is very similar to VXLAN
- » basis: Generic Routing Encapsulation (GRE)
  - » generic header
  - » can encapsulate a wide variety of network layer protocols
  - » point-to-point links over an Internet Protocol network
- » NVGRE
  - » GRE header
    - »
 

```

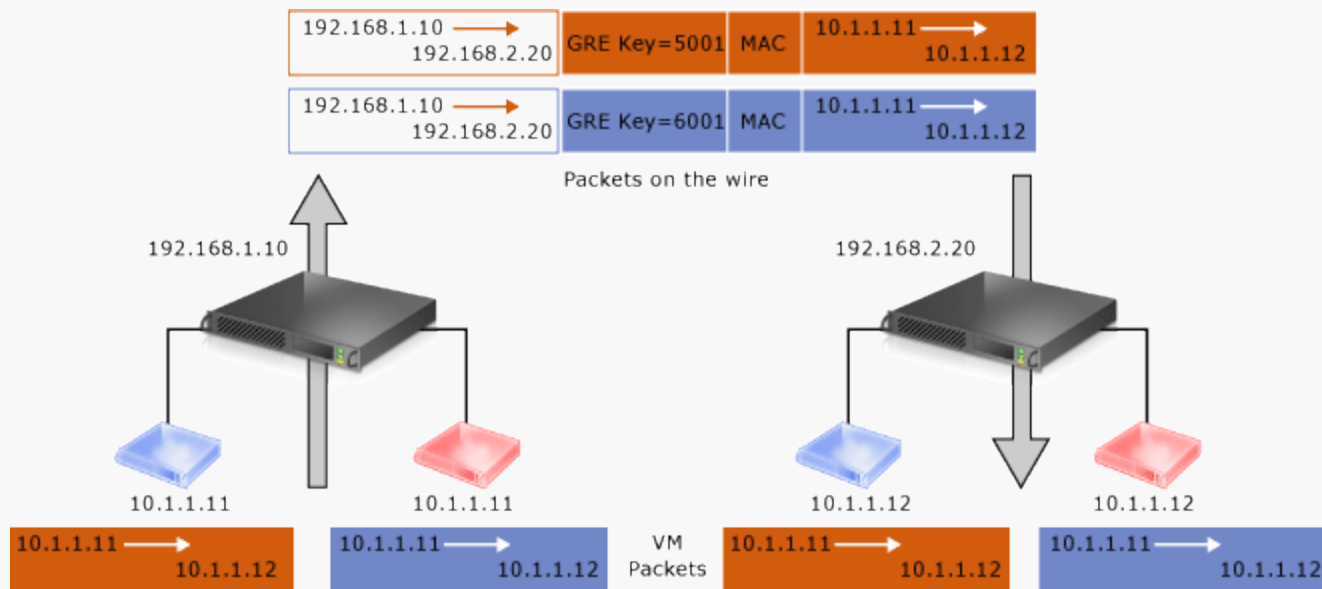
              +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
              |0| |1|0|   Reserved0   | Ver |   Protocol Type 0x6558   |
              +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
              |
              |               Virtual Subnet ID (VSID)               |   FlowID   |
              +-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+-----+
              
```

      - » Virtual Subnet Identifier (VSID) 24 bit ⇨ 16 million tenant
      - » FlowID: optional, unique flow identifier
        - » for ECMP hashing
    - » no inner VLAN tag (or it is removed)
      - » encoded into VSID or multiple VSIDs for the same tenant



# NVGRE

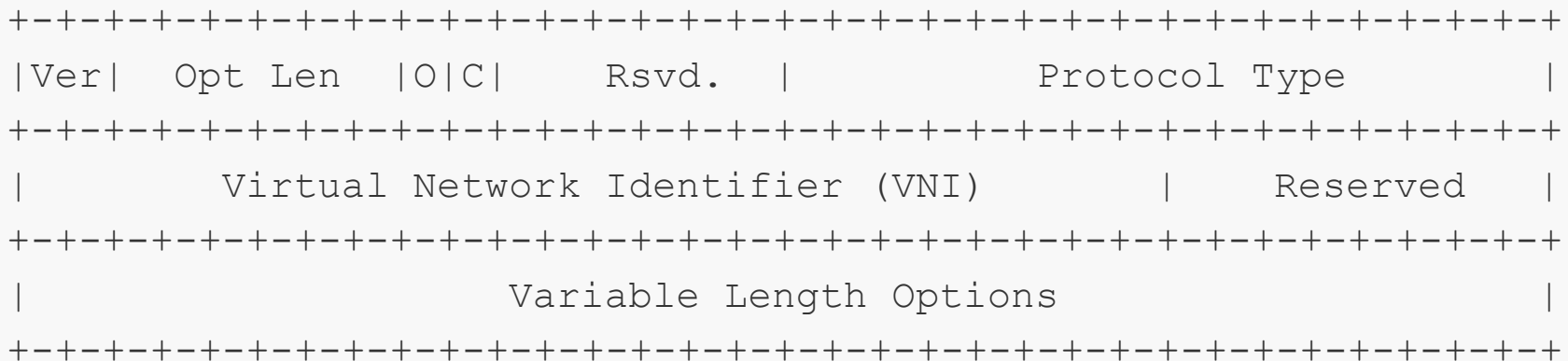
- » Network Virtual Endpoint (NVE)
  - » forwarding to the destination based on VSID and DMAC
- » not specified in the Internet draft
  - » dissemination of addressing information
  - » restoring VLAN information





# Generic Network Virtualization Encapsulation

- » MAC-in-UDP over IPv4/IPv6
- » universal, extensible
- » specifies only the encapsulation format
- » optional fields
  - » variable field lengths, flexibility
- » Geneve header:







# Location of Tunnel Endpoint

- » Inside hypervisor/vSwitch
  - » most common
  - » closest to the VMs
  - » consumes CPU resources
  - » TCP segmentation offload (TSO), checksum offload support for non-encapsulated packets
- » NIC
  - » offload support for tunneling
- » ToR switch
  - » VMs are unknown
  - » has to identify VM based on inner MAC to assign VNI/VSID





# Comparison

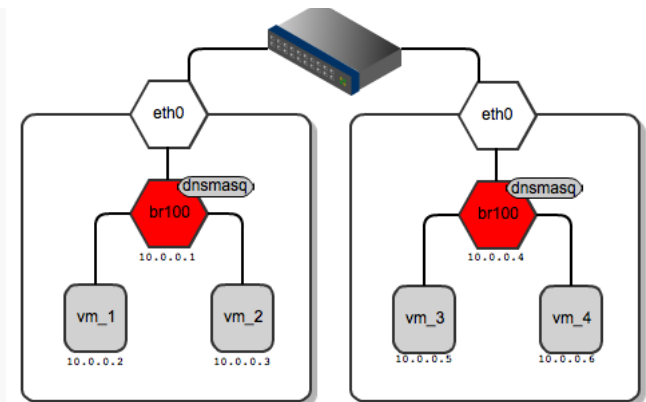
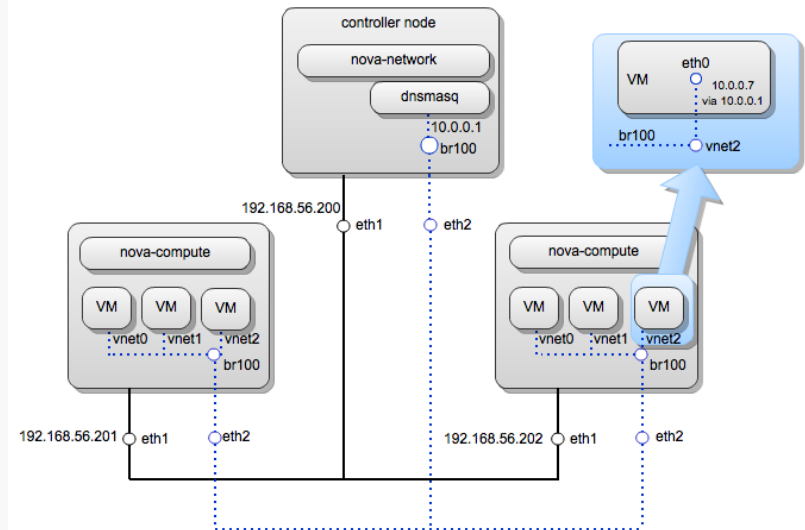
	VXLAN	NVGRE	STT
Extra bytes	50 (VLAN: +4)	42 (VLAN: +4)	First segment: 76 Others: 58 (VLAN: +4)
Protocol	UDP	GRE	TCP
Tenant separation	24 bit VNID	24 bit VSID	64 bit Context ID
ECMP flow differentiation (inner ⇒ outer flow)	Source UDP port	VSID + FlowID (8bit)	Source TCP port



# OPENSTACK NEUTRON

# OpenStack network architecture

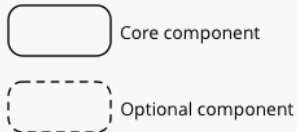
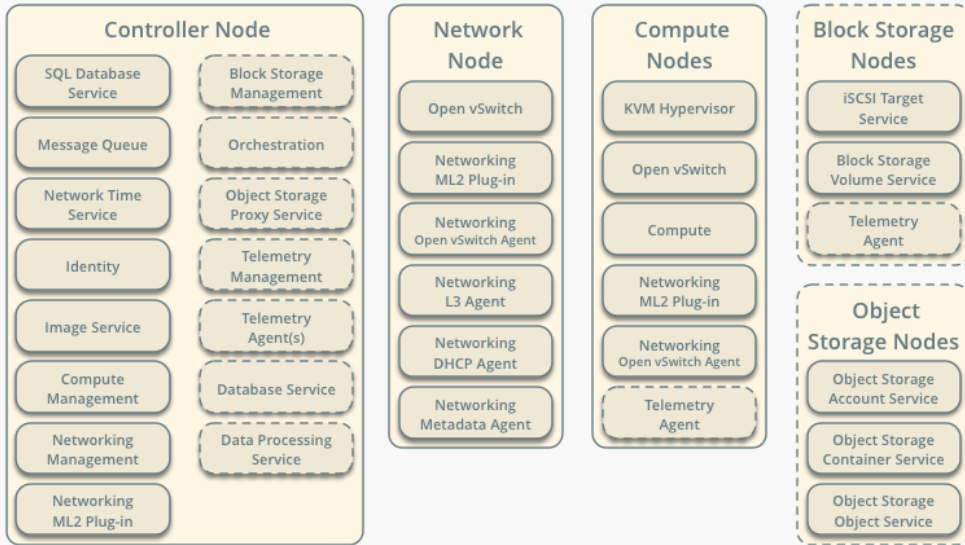
- » „Networking in OpenStack is a complex, multifaceted challenge.“ /OpenStack Operations Guide/
- » Network as a Service
- » functions
  - » IP addressing
    - » static, DHCP
    - » floating IP
  - » virtual networks
    - » flat, VLAN
  - » self-service
- » alternatives
  - » Nova networking / Neutron
  - » single-host / multi-host
- » Neutron
  - » plug-in architecture
  - » SDN/OpenFlow



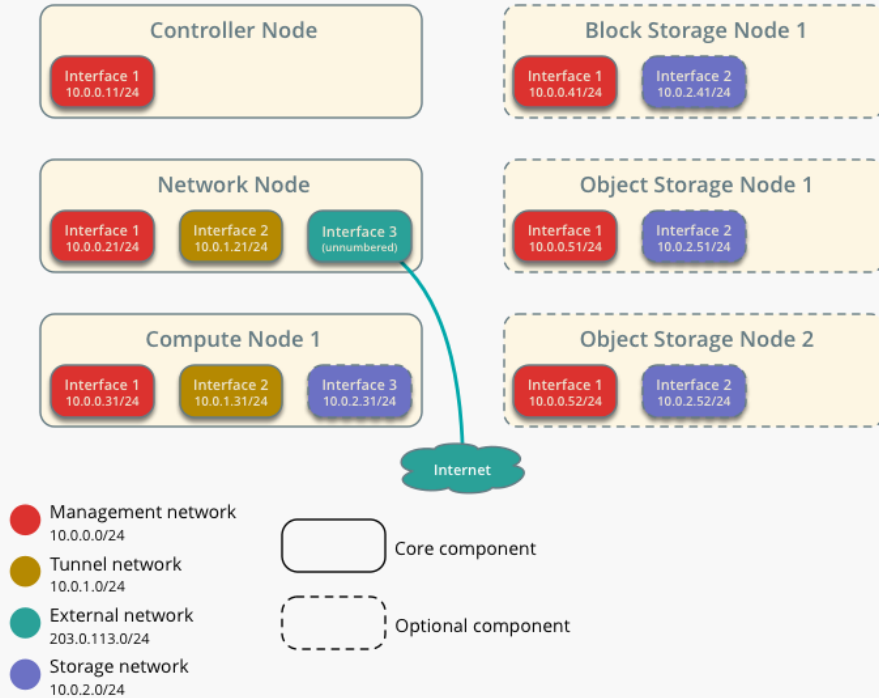


# Neutron network

Minimal Architecture Example - Service Layout  
OpenStack Networking (neutron)



Minimal Architecture Example - Network Layout  
OpenStack Networking (neutron)





# Nova and Neutron Network

## » Nova

- » basic networking functions
  - » network address translation (NAT), DHCP, DNS
- » only support L2 bridge networking
  - » allows virtual interfaces to connect to the outside network through the physical interface
- » limited scalability
  - » VLAN, DNS&DHCP (dnsmasq)

## » Neutron

- » network abstraction
- » L2/L3 network, self-service
  - » e.g. more LAN segments for a web service
- » Load Balancing, Virtual IP, VPN, firewall
- » overlay VLAN tunneling
- » Distributed Virtual Router (from Juno)



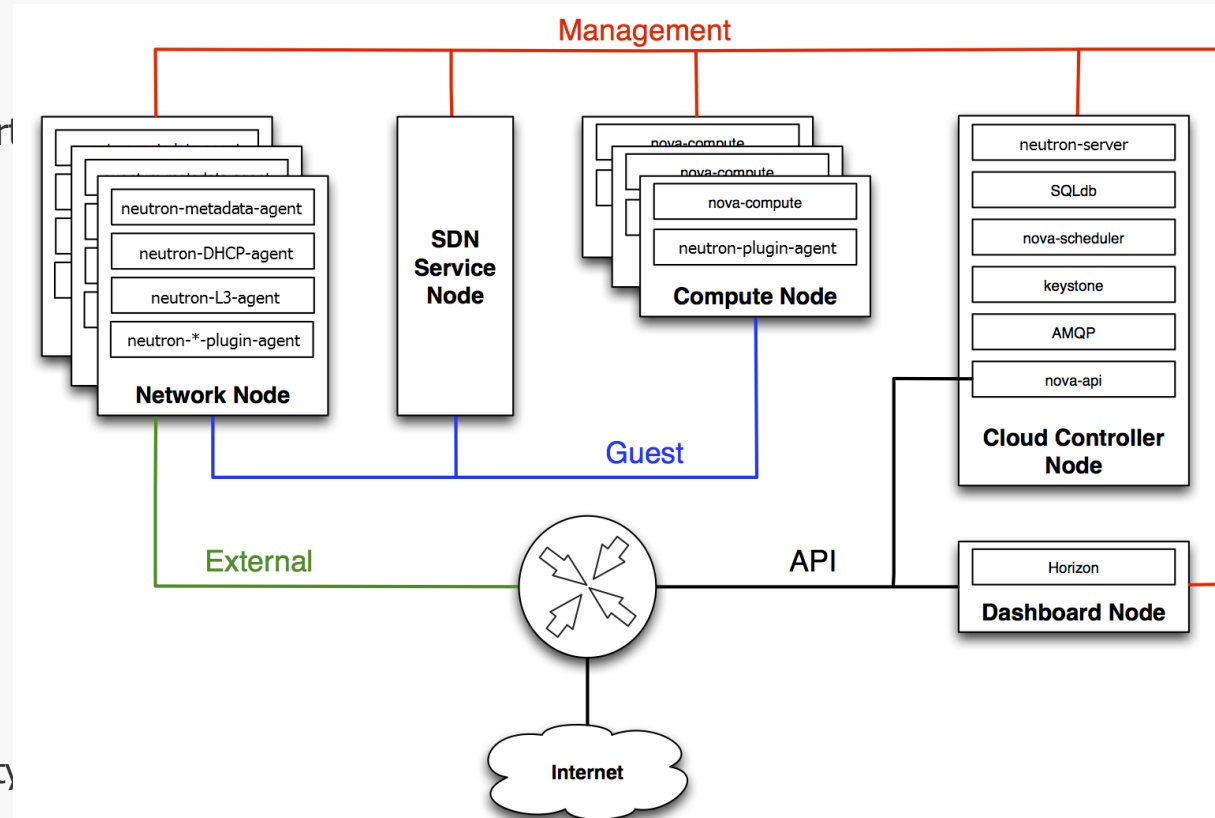
# Neutron network abstraction

- » External /physical/ network, e.g. Internet
- » Internal networks to connect VMs
  - » virtual: network, subnetwork, router
  - » floating IP from external network address space for reaching a VM from outside
- » Security groups
  - » firewall rules
  - » assigned to the VM
- » Open vSwitch
  - » core plugin
  - » br-int (integration bridge)
    - » connected to VMs
  - » br-ex
    - » connected to external network



# Neutron components

- » server + plugin + agent architecture
  - » neutron-server
    - » on controller node
    - » handling API requests
    - » network model and port IP address setup
  - » plugin – extensions: neutron-\*-plugin
    - » on network node
  - » plugin-agent: neutron-\*-agent
    - » on compute node
    - » managing the local vswitch
  - » general agents
    - » DHCP: neutron-dhcp-agent
    - » L3 agent: neutron-l3-agent
      - » L3/NAT functionality towards the external network
      - » implementation: Linux IP stack and iptables





# Modular Layer 2 (ML2) plugin

- » Managing different L2 network technologies in uniform way
- » Operates with openvswitch, linuxbridge and Hyper-V L2 agents
- » Type drivers for different network types
  - » Flat
  - » Local (DevStack single box)
  - » VLAN
  - » GRE
  - » VXLAN

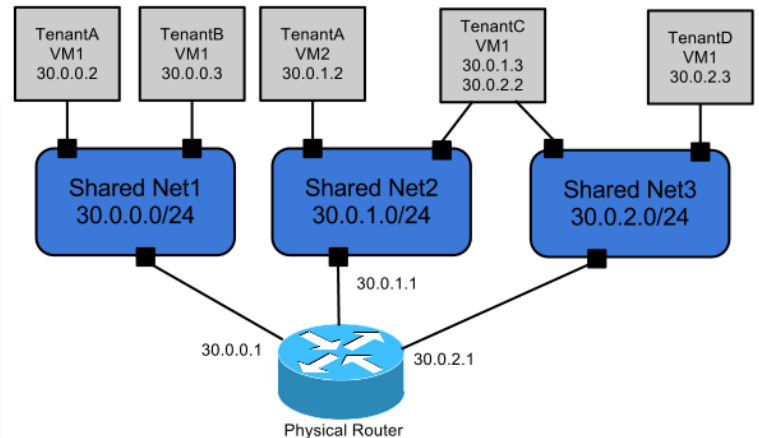
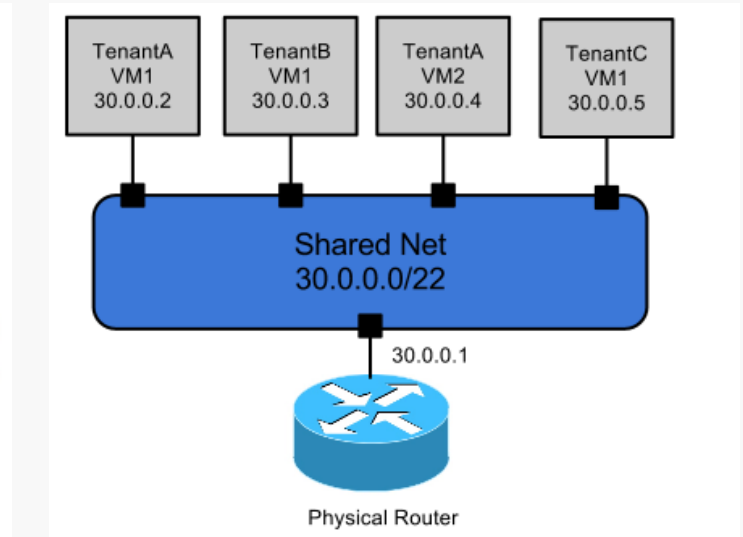
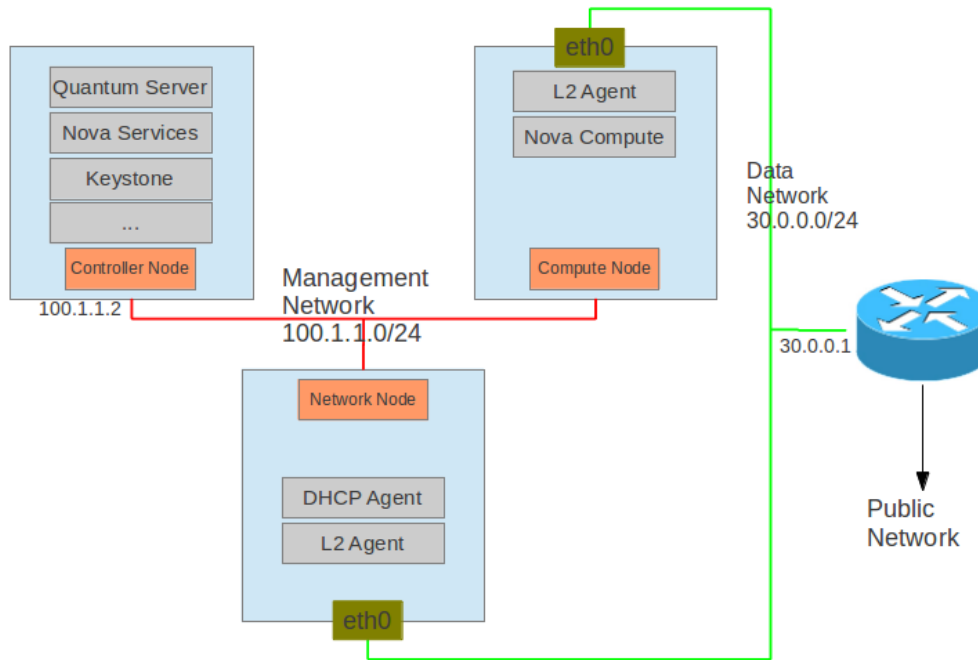


# Network namespaces

- » kernel level namespaces, not only for networking
  - » file system, process, user, etc.
- » isolated Layer2 networks with overlapping IP addresses
- » separating virtual interfaces, routers
- » e.g. dhcp-agent and l3-agent runs in different namespaces
- » In practice
  - » `ip netns`
    - » lists available network namespaces
  - » `ip netns exec <namespace> <command>`
    - » e.g. `ip netns exec qdhcp-e521f9d0-a1bd-4ff4-bc81-78a60dd88fe5 ip a`

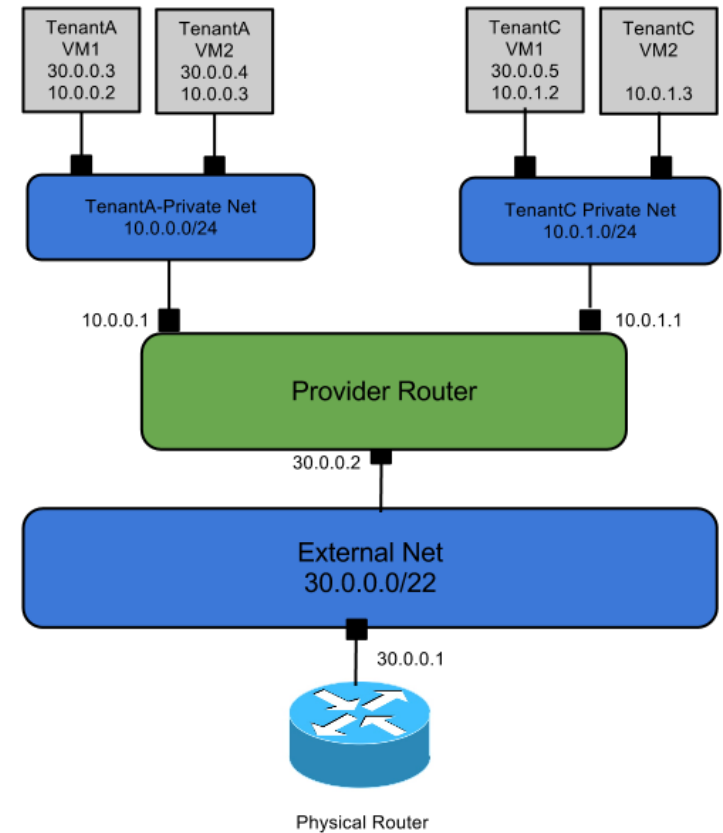
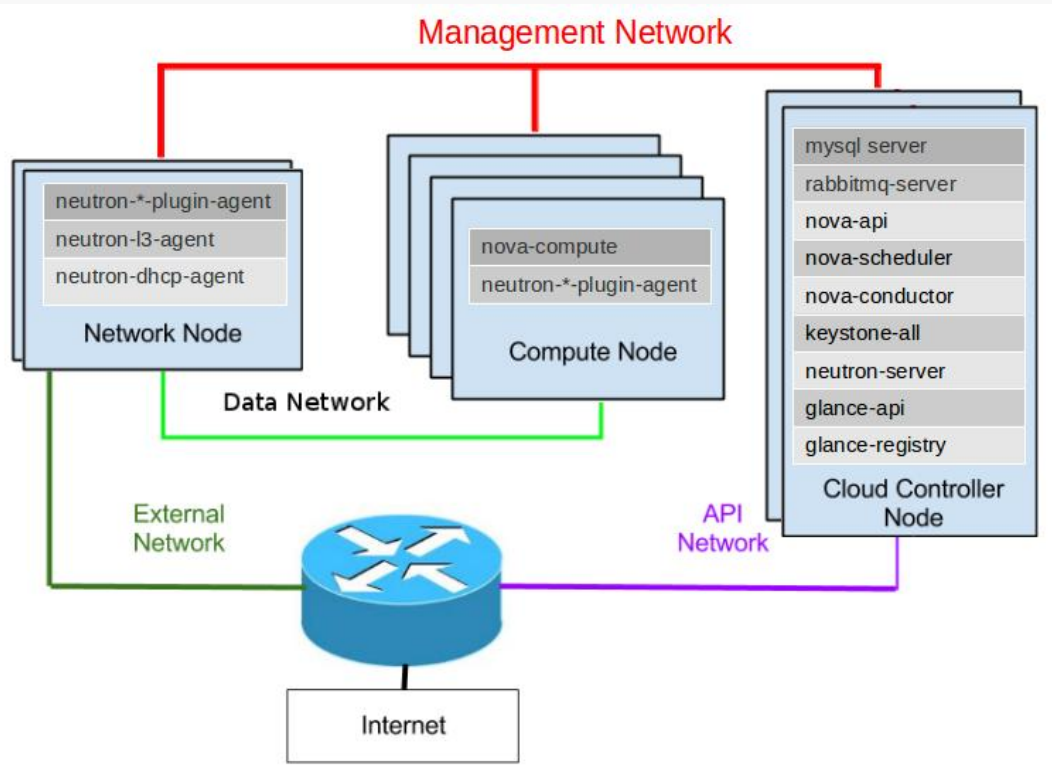


# Neutron: single/multiple flat network



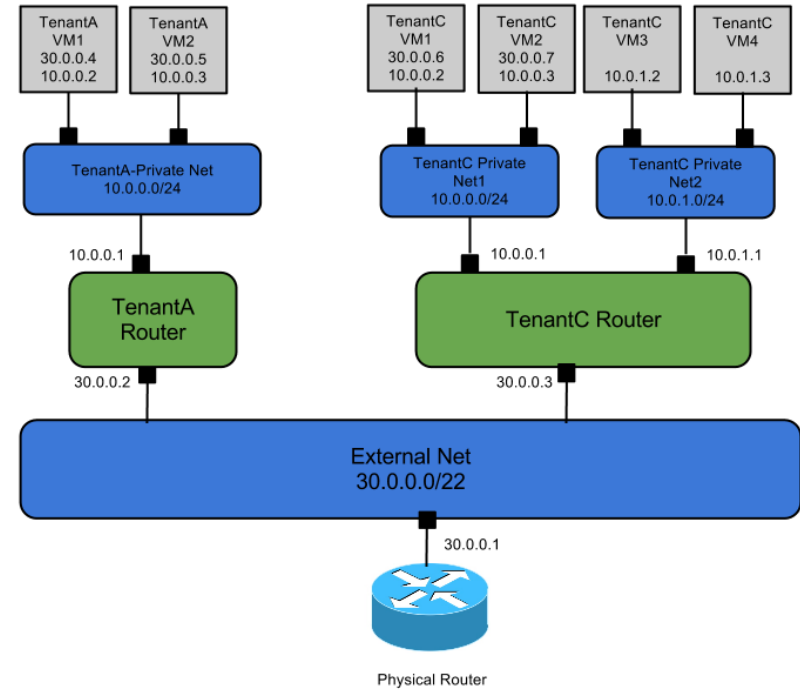
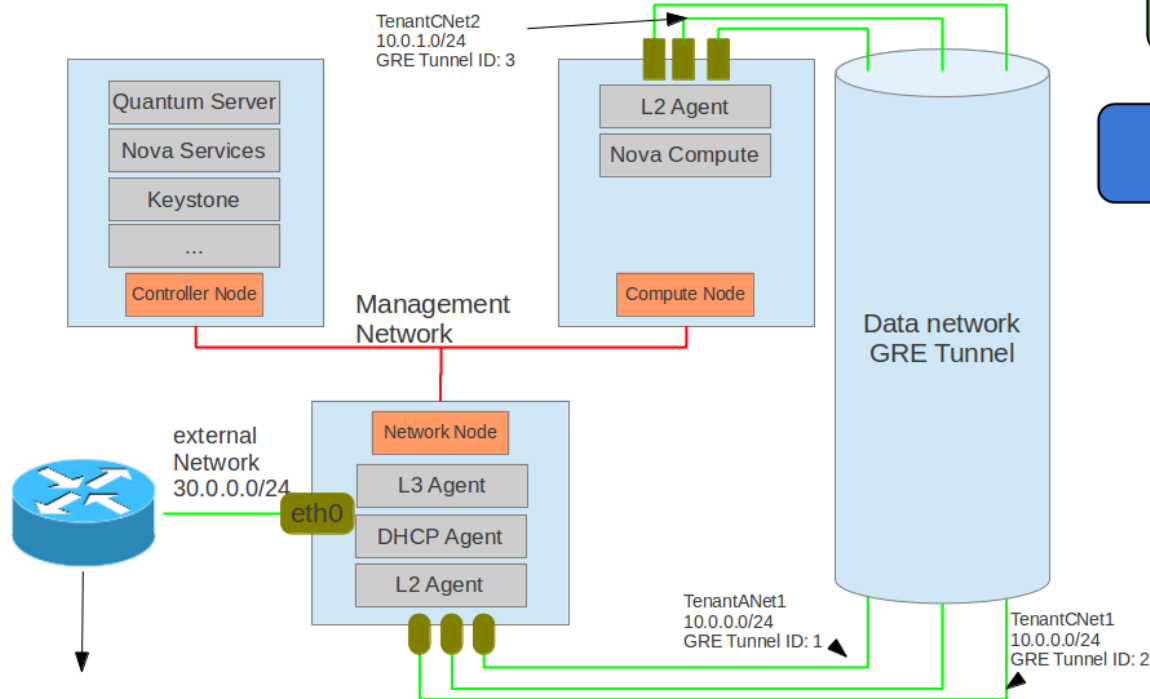


# Neutron: provider router



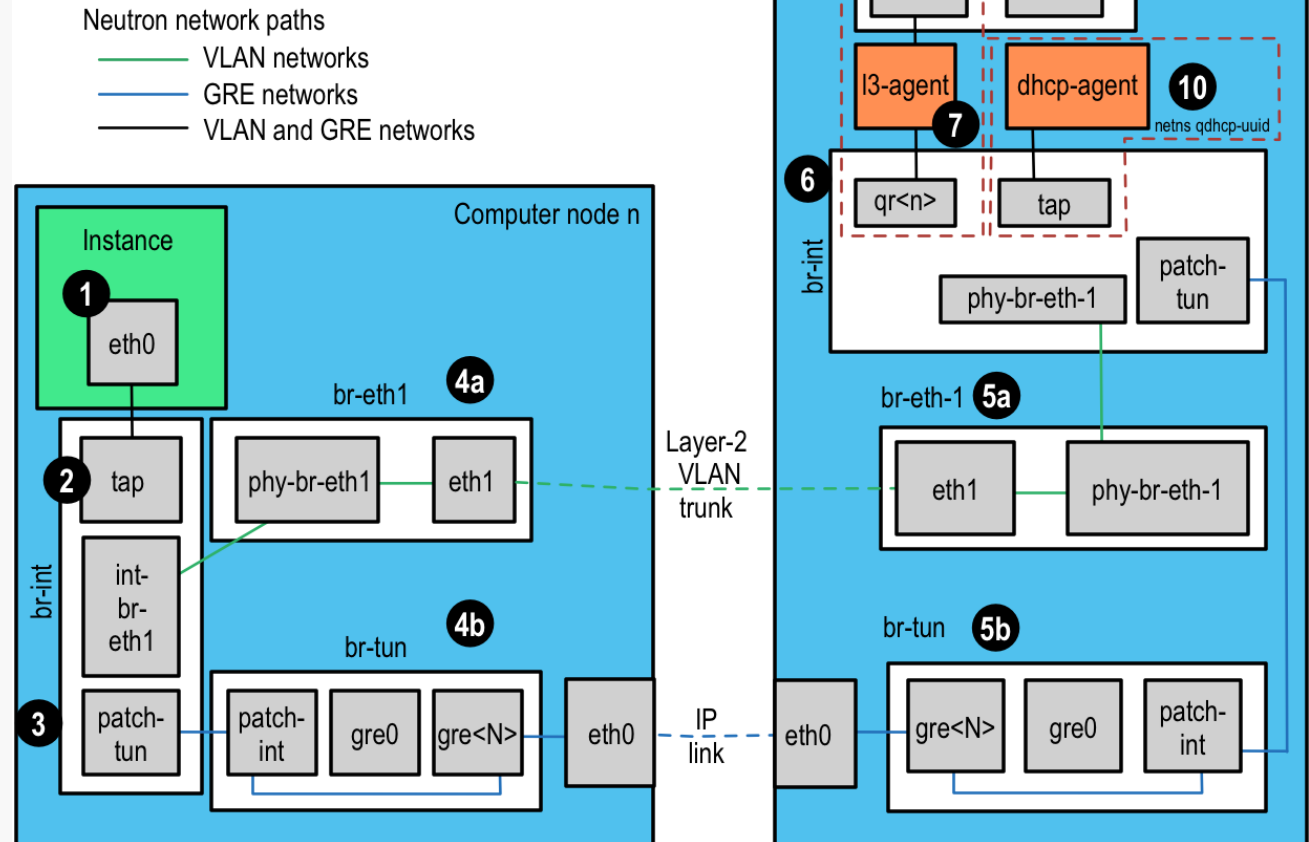


# Neutron: tenant routers



# Path of a packet

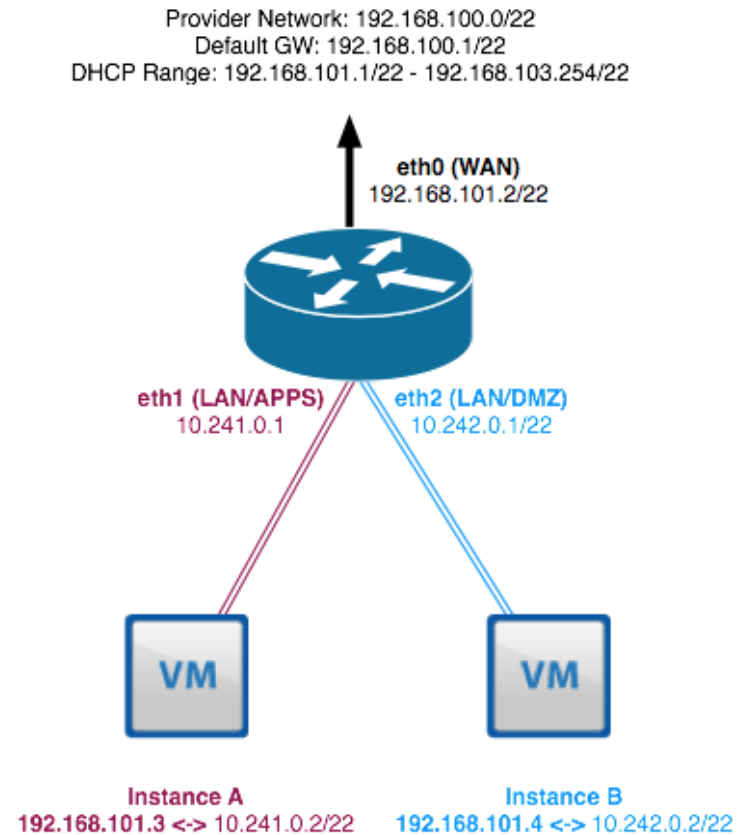
- » Test Access Point (TAP) device
- » int-br: integration bridge
- » br-eth1: VLAN internal/external tag translation
- » veth: between int-br-eth1 and phy-br-eth1



# Floating IP

- » Neutron router
  - » gateway for VMs
  - » iptables/NAT rules in the namespace of router
    - » nova network: in hypervisor
  - » floating IP addresses allocated from the public network address range

**Diagram 1.1 - Logical Neutron Router**



- Diagram 1.1 -

**eth0** is connected to a PROVIDER network.  
**eth1** is connected to a TENANT network.  
**eth2** is connected to a TENANT network.

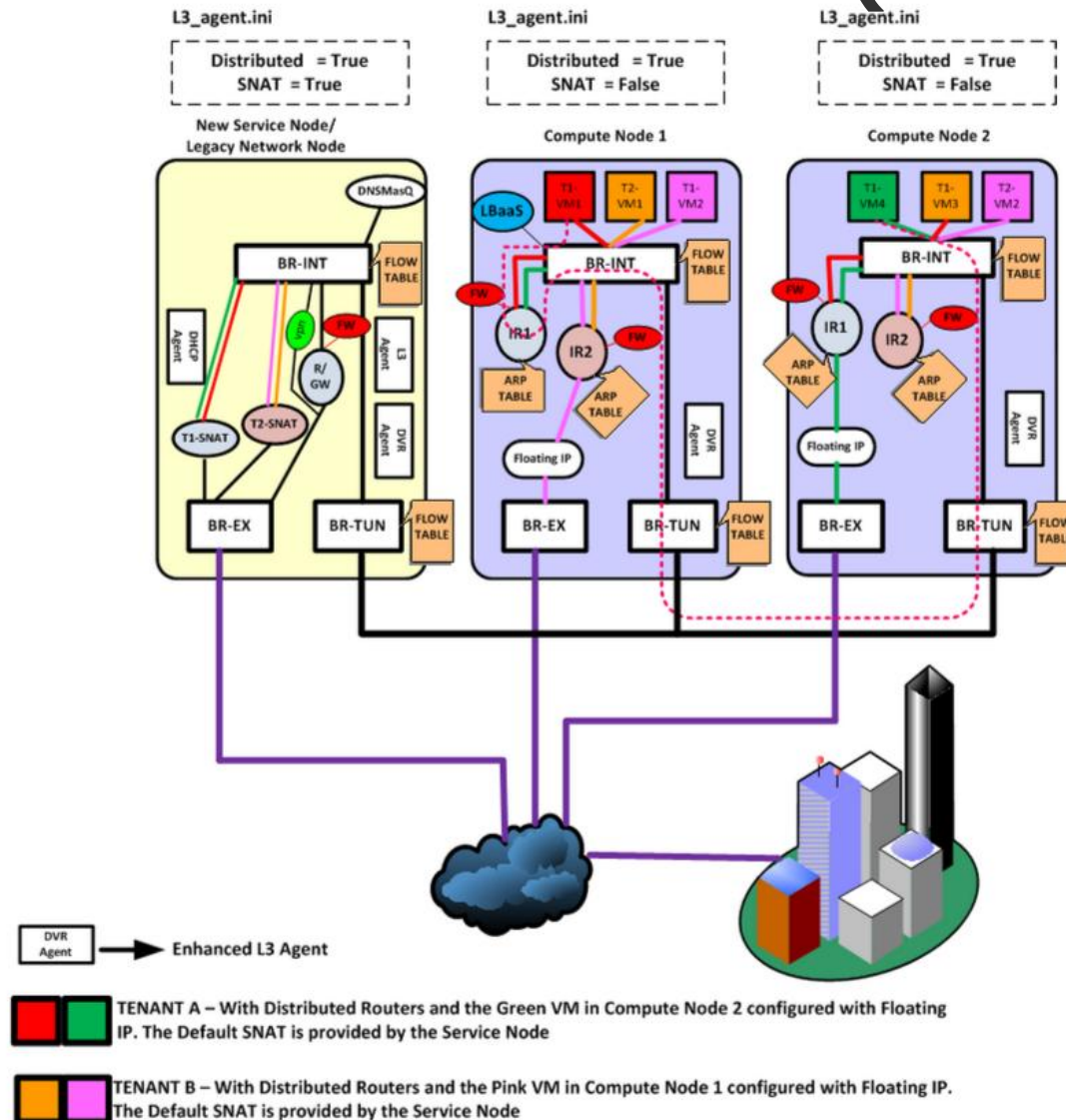
**Floating IPs** are assigned from the DHCP range of the PROVIDER network:

DHCP Range: 192.168.101.1/22 - 192.168.103.254/22

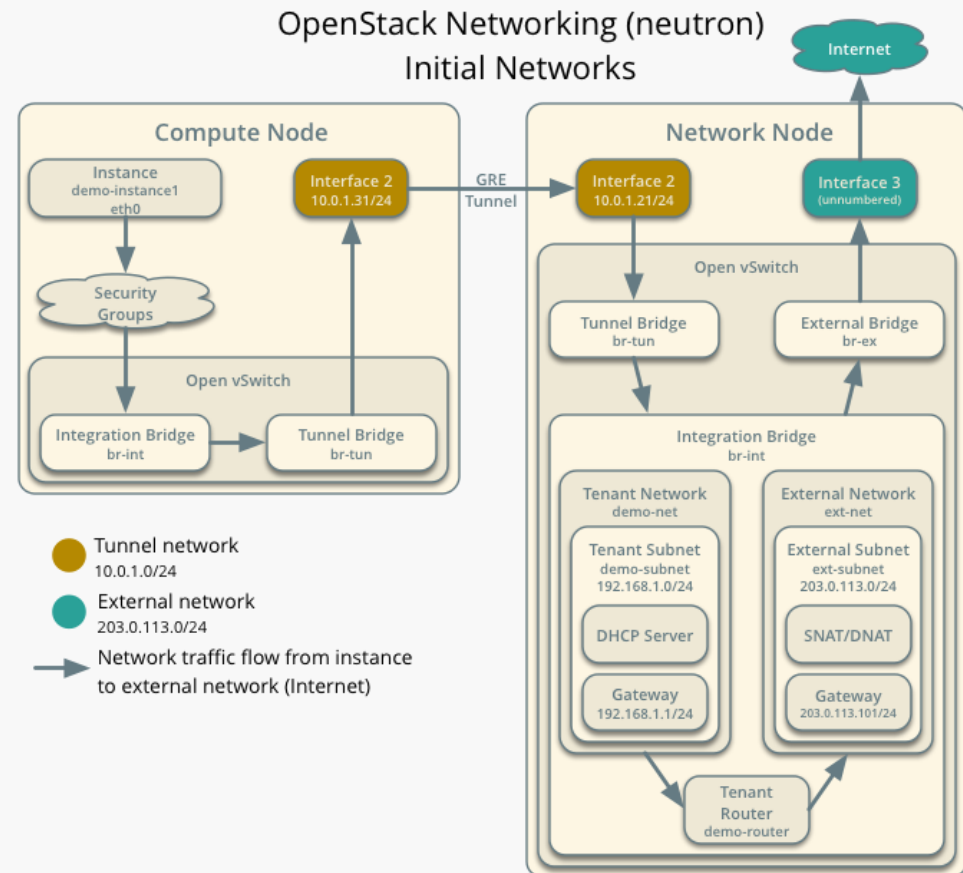




# Distributed Virtual Router (DVR)



# Virtual network configuration



## » Open vSwitch

» setup by ovs-dpctl / OpenFlow

» e.g. mapping VM MAC address and hypervisor transport IP address



# References

- » Overlay Virtual Networking Explained, Ivan Pepelnjak, NIL Data Communications, 2011.
- » <http://docs.openstack.org>
- » <https://developer.rackspace.com/blog/neutron-networking-l3-agent/>
- » [https://www.rdoproject.org/Networking\\_in\\_too\\_much\\_detail](https://www.rdoproject.org/Networking_in_too_much_detail)