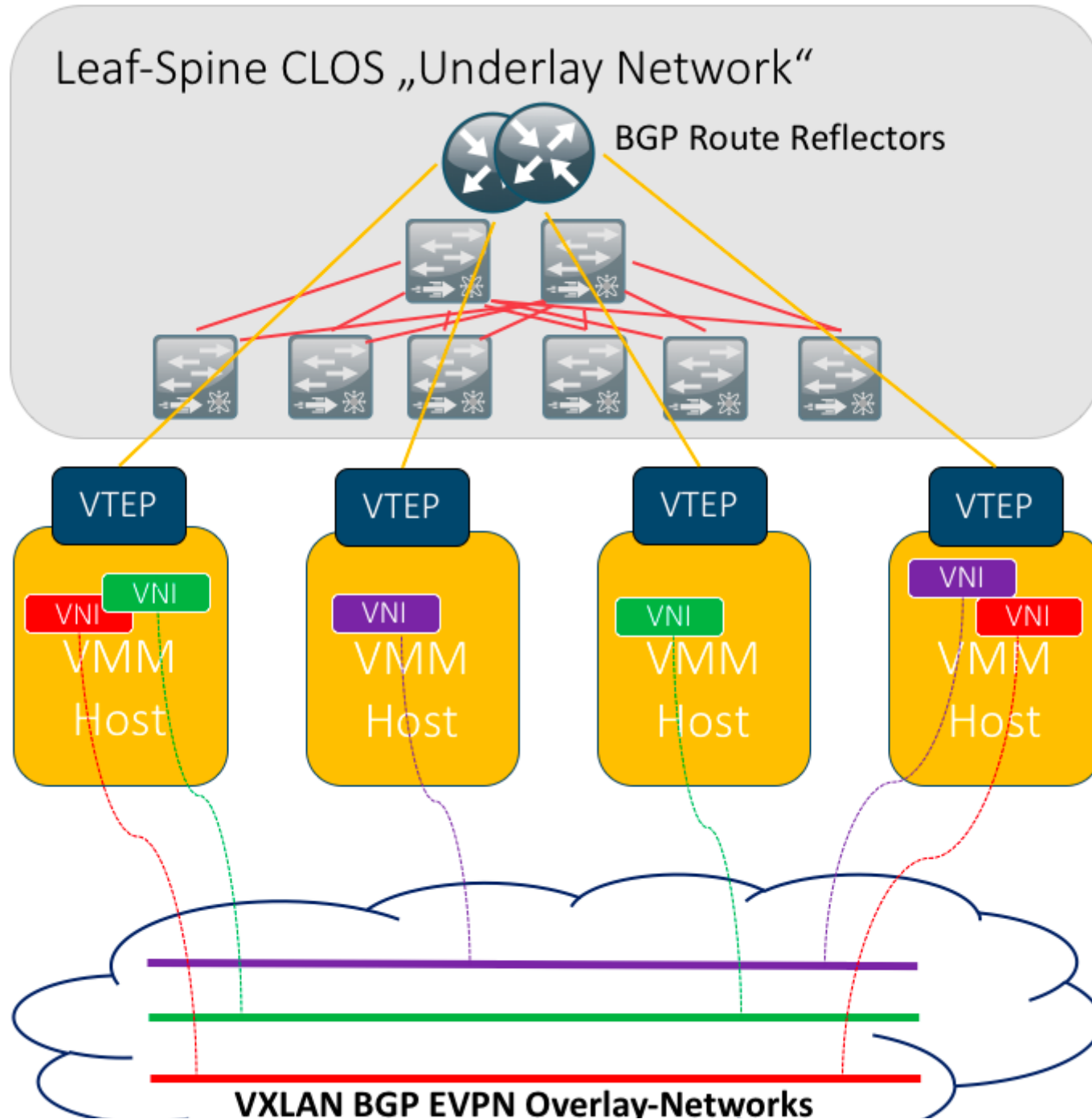


# **CLOUD BLOG – OPENSOURCE POWER @ENTERPRISE & DATACENTER**

Infrastructure-as-a-Service / Platform-as-a-Service / Datacenter 2.0



13. JULI 2017 VON SEBASTIAN

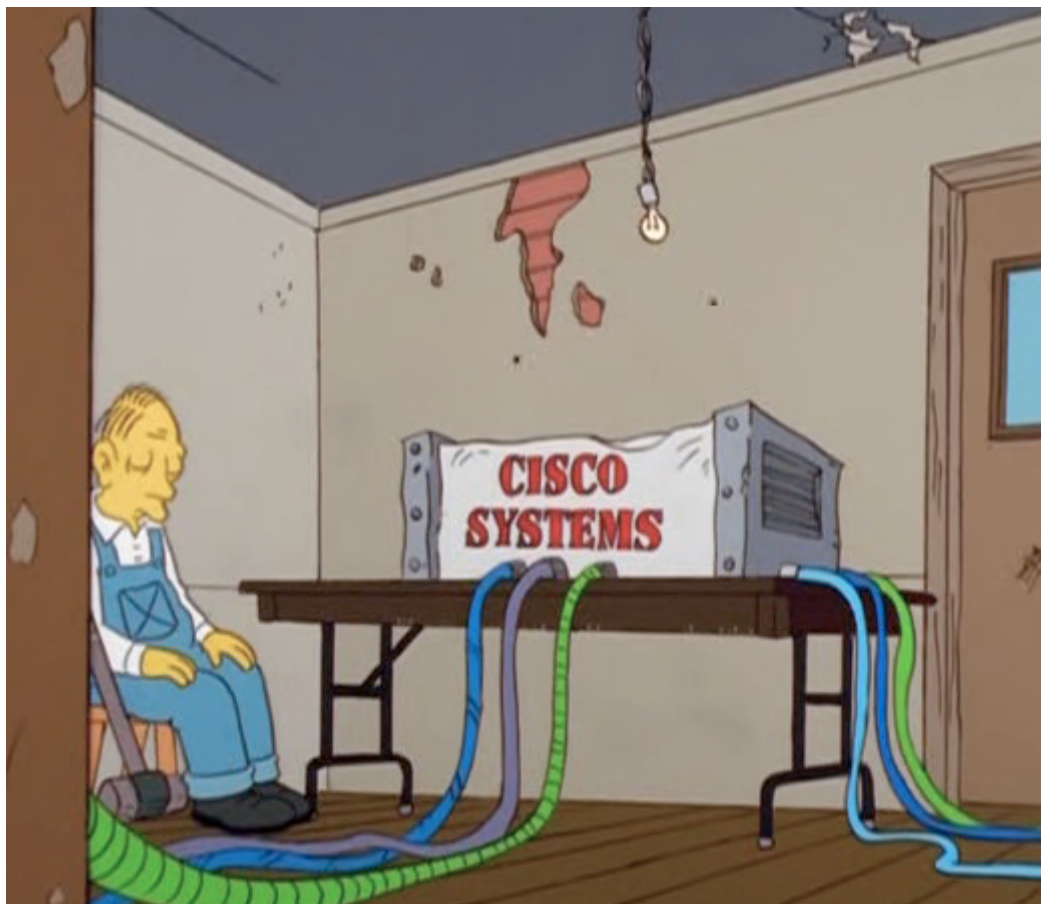
## BGP EVPN (MP-BGP + VXLAN) with OpenNebula

From my point of view a well designed datacenter fabric based on VXLAN with BGP EVPN as a kind of control plane, based on a Leaf-Spine (CLOS) topology is one of the secrets to gain highly scalable, resilient and performant multi-tenancy for datacenter and cloud networks.

Thanks to Mario's idea i started playing around with this fancy technology. And of course, my employer NTS is a networking company – and last but not least: it's SDN and that's cool 😊

In this blog post we try to build a lab setup with BGP-EVPN and OpenNebula. Perhaps i will go to develop an OpenNebula Add-On or VNM-Driver if it's needed.

Because it's inside our Lab i'm going to implement the BGP part with Docker and „Cumulus Quagga“, a software BGP implementation. In a real world scenario the Route Reflectors can be the Cisco ToR switches or any Leaf/Spine (Nexus) switch. Cause we use Standards, the implementation is open for any other vendor. But of course, there is only one vendor :



### Why BGP EVPN / MP BGP with VXLAN?

With BGP EVPN / MP-BGP as control plane for VXLAN we will work with proven and well-known standards – BGP empowers the whole internet. It's capable to handle thousands and thousands of routes with ease.

The following RFCs and Drafts are taken into account for this solution:

- VXLAN (RFC7348)
- BGP4 (RFC4271)

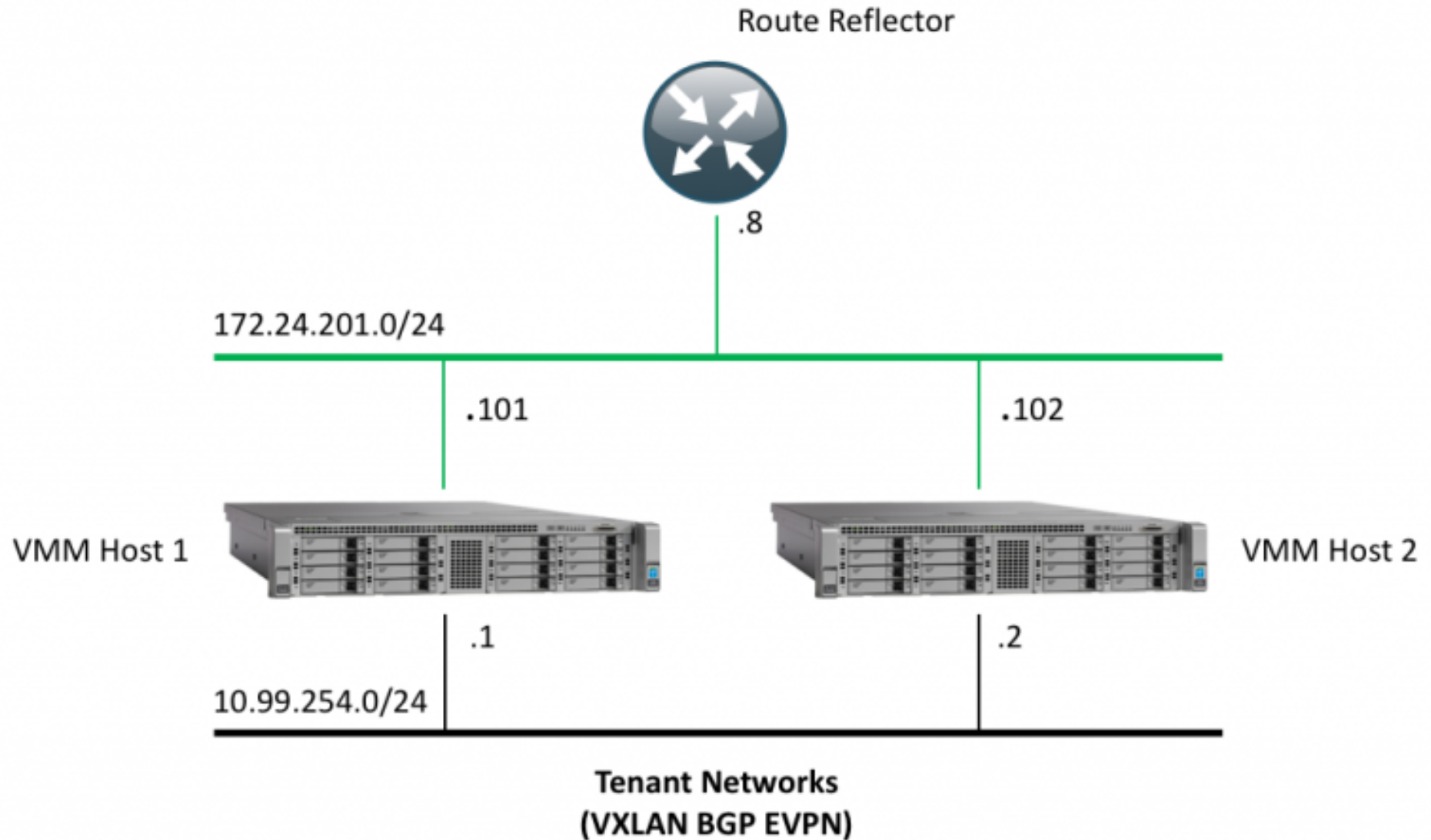
- BGP MPLS-Based Ethernet-VPN ([RFC 7432](#))
- [draft-ietf-bess-evpn-overlay-08](#)
- MP-BGP ([RFC4760](#))

VXLAN with BGP EVPN has several advantages over VXLAN with Multicast:

- Traffic engineering
- More scalability
- Better manageability
- Easier to handle than IP-Multicast in a large scale network

In our case BGP carries the reachability informations (f.e. the forwarding tables of our virtual switches/bridges) in Network Layer Reachability Information (NLRI) format for our EVPN. In particular it carries the MAC-Adresses and the information, where they are attached to (VTEP).

Our lab setup looks like this:

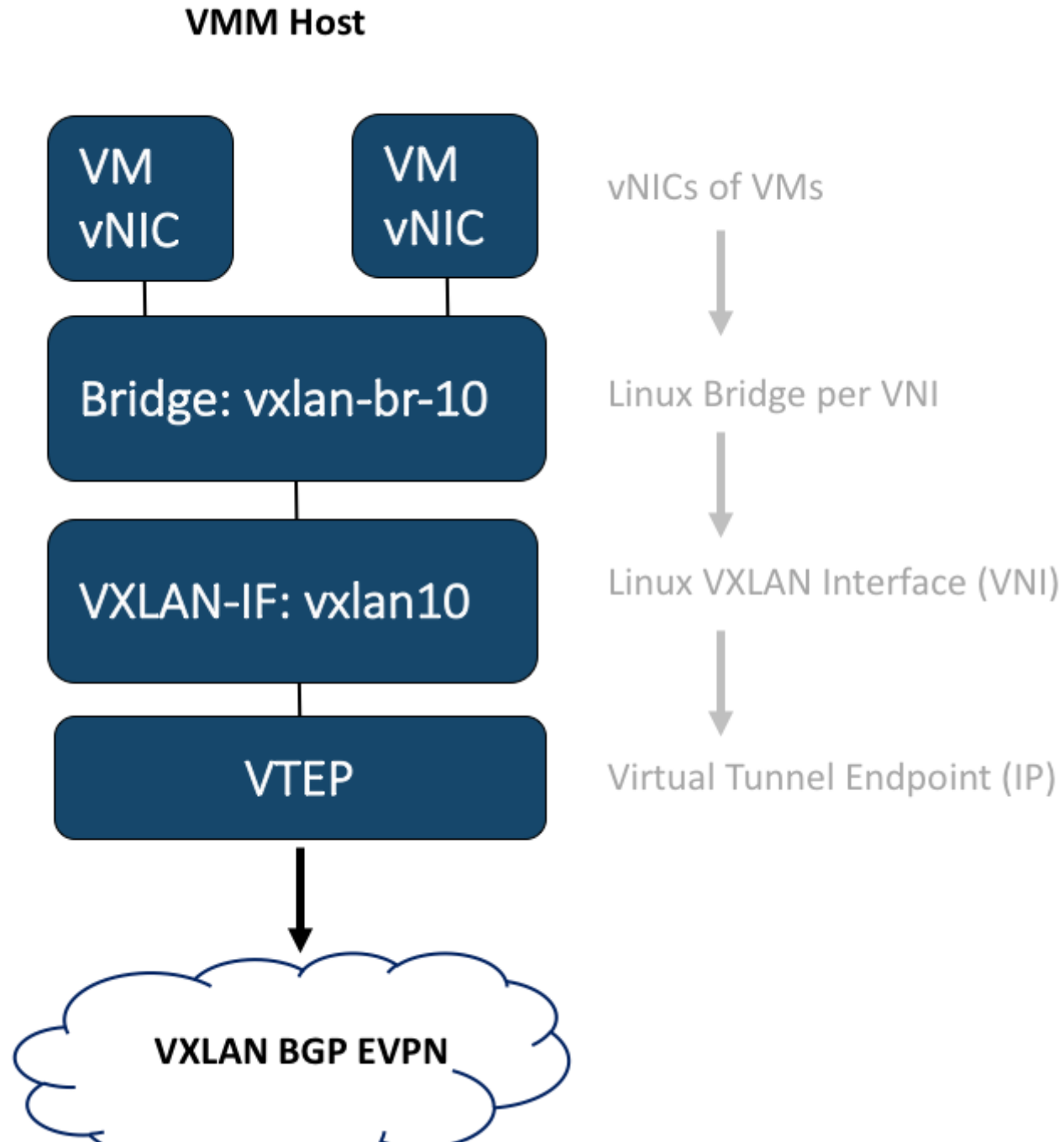


*VXLAN BGP EVPN Lab network*

Every VTEP (in our case the hypervisor host) exports two kinds of „routes“ via BGP:

- Type 2 Route: The local MAC addresses for the VNI
- Type 3 Route: The local VNIs with it`s IPv4-Endpoint – the VTEP

The networking stack inside the hypervisor hosts is based on proven and well known standards (vNICs of the VMs, attached to an Linux-Bridge and a VXLAN-Interface for each VNI enslaved to that bridge):







### Networking Stack inside a VMM Host

In our case we operate one Route-Reflector, also Docker based, on a VM in our lab. In a real world example we'll run at least two Route-Reflectors. The configuration of our Route-Reflector is quite simple:

- Create a peer-group called „vxlan“
- Create a dynamic membership for the VTEP network (172.24.201.0/24)
- Reflect all EVPN to all clients

First we had to pull the current Docker-Image for the Quagga implementation and start it:

```
1. docker pull cumulusnetworks/quagga
2.
3. # Quagga has to run in "privileged" mode to get the attached VNIs of the VMM hosts and to inject the FDBs
4. docker run -t -d --net=host --privileged --name BGP-EVPN-RR1 cumulusnetworks/quagga:latest
```

Now we are able to configure our RR:

```
1. docker exec -i -t BGP-EVPN-RR1 /usr/bin/vtysh
```

The CLI is like our well known Cisco CLI 😊

```
1. router bgp 64512
2.   bgp router-id 172.24.201.8
3.   bgp cluster-id 172.24.201.8
4.   bgp log-neighbor-changes
5.   no bgp default ipv4-unicast
6.   neighbor vxlan peer-group
7.   neighbor vxlan remote-as 64512
8.   neighbor vxlan capability extended-nexthop
```

```
9. neighbor vxlan update-source 172.24.201.8
10. bgp listen range 172.24.201.0/24 peer-group vxlan
11. !
12. address-family evpn
13.     neighbor vxlan activate
14.     neighbor vxlan route-reflector-client
15. exit-address-family
16. !
```

At the VTEP side we also had to run a BGP implementation, which advertises all connected VNIs to the Route-Reflector:

```
1. router bgp 64512
2.     bgp router-id 172.24.201.102
3.     no bgp default ipv4-unicast
4.     neighbor vxlan peer-group
5.     neighbor vxlan remote-as 64512
6.     neighbor vxlan update-source vnet0
7.     neighbor vxlan capability extended-nexthop
8.     neighbor 172.24.201.8 peer-group vxlan
9.     !
10.    address-family evpn
11.        neighbor vxlan activate
12.        advertise-all-vni
13.    exit-address-family
14.    !
```

From now, all VNIs will be advertised from the VTEP to the RR. The RR will reflect all EVPN to all connected neighbours. Of course there are much more features (policies, communities, ...) – i cannot cover them, because i'm not the „networking guy“ – that's stuff for my brilliant networking colleagues like Mario.

To create some VNIs at our VMM hosts i crafted a little script to create the VNI, VTEP and a bridge to attach the VNICs of the VMs on it:

```
#!/bin/bash

# 2017 NTS / S.Mangelkrmaer
# BGP EVPN VXLAN Lab

for ((vni = 10 ; vni <= 20 ; vni += 1)); do

    echo "Creating VXLAN Network with VNI-ID: $vni"

    # VTEP IP of physical iface eth3
    vtep=`ip -4 addr show eth3 | grep -Po 'inet \K[\d.]+'`
    ip link add vxlan${vni} type vxlan id ${vni} dstport 4789 local $vtep nolearning

    # Create Bridge
    brctl addbr vxlan-br-${vni}

    # Add VTEP to Bridge
    brctl addif vxlan-br-${vni} vxlan${vni}

    # Disable STP
    brctl stp vxlan-br-${vni} off

    # Set links up
    ip link set up dev vxlan-br-${vni}
    ip link set up dev vxlan${vni}
done
```

Within our OpenNebula we simply create new „Bridged Networks“.

In this example we will use two VXLAN networks „vxlan10“ and „vxlan20“, create 4 virtual machines on two different VMM hosts and start our testing.

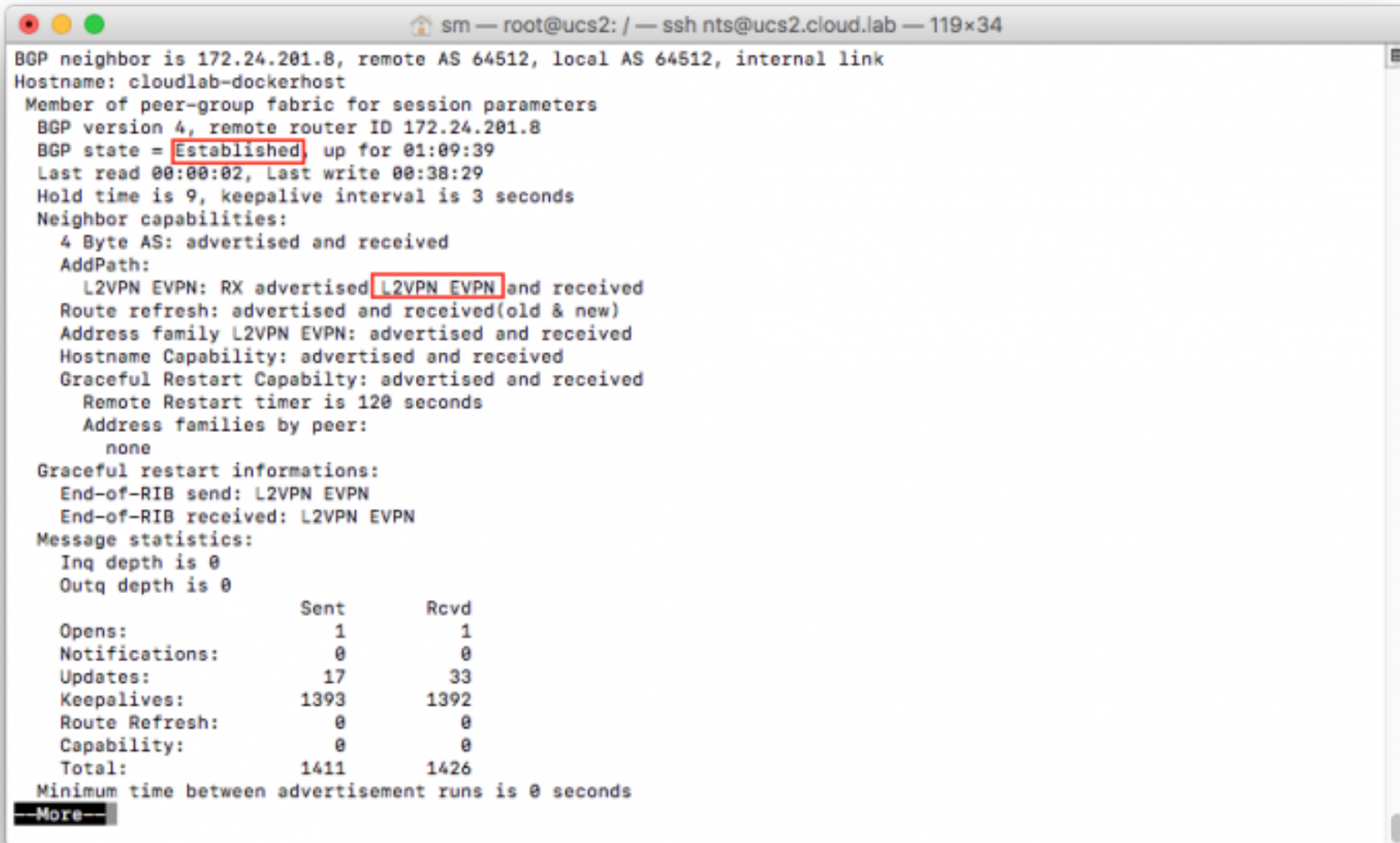
After instantiating the VMs we will check if we can see the VXLAN networks at the VTEP – this can be done inside our Quagga-Docker-Container with the `show interface $VNI` command:

```
sm — root@ucs2: / — ssh nts@ucs2.cloud.lab — 119x34
ucs2#
ucs2# show interface vxlan10
Interface vxlan10 is up, line protocol is up
  Link ups:      0    last: (never)
  Link downs:   0    last: (never)
  PTM status: disabled
  vrf: Default-IP-Routing-Table
  index 25 metric 0 mtu 1500
  flags: <UP,BROADCAST,RUNNING,MULTICAST>
  Type: Ethernet
  HWaddr: 4e:ab:46:26:f3:39
  inet6 fe80::4cab:46ff:fe26:f339/64
  Interface Type Vxlan
  VxLAN Id 10
  Access VLAN Id 1
  Master (bridge) ifindex 26 ifp 0x562399386b00
ucs2#
ucs2#
ucs2# show interface vxlan20
Interface vxlan20 is up, line protocol is up
  Link ups:      0    last: (never)
  Link downs:   0    last: (never)
  PTM status: disabled
  vrf: Default-IP-Routing-Table
  index 45 metric 0 mtu 1500
  flags: <UP,BROADCAST,RUNNING,MULTICAST>
  Type: Ethernet
  HWaddr: 62:83:fb:1b:2e:53
  inet6 fe80::6083:fbff:fe1b:2e53/64
  Interface Type Vxlan
  VxLAN Id 20
  Access VLAN Id 1
  Master (bridge) ifindex 46 ifp 0x56239938bba0
ucs2#
```

Show VXLAN interface at VTEP

We see the VNI and the correct detected linux bridge (where the vNICs are attached to). This looks good.

Next we check if the VTEPs correctly has established their BGP session to the RR with the `show bgp neighbours` command:



```
sm — root@ucs2: / — ssh nts@ucs2.cloud.lab — 119x34
BGP neighbor is 172.24.201.8, remote AS 64512, local AS 64512, internal link
Hostname: cloudlab-dockerhost
Member of peer-group fabric for session parameters
BGP version 4, remote router ID 172.24.201.8
BGP state = Established up for 01:09:39
Last read 00:00:02, Last write 00:38:29
Hold time is 9, keepalive interval is 3 seconds
Neighbor capabilities:
 4 Byte AS: advertised and received
AddPath:
  L2VPN EVPN: RX advertised L2VPN EVPN and received
Route refresh: advertised and received(old & new)
Address family L2VPN EVPN: advertised and received
Hostname Capability: advertised and received
Graceful Restart Capabilty: advertised and received
  Remote Restart timer is 120 seconds
  Address families by peer:
    none
Graceful restart informations:
End-of-RIB send: L2VPN EVPN
End-of-RIB received: L2VPN EVPN
Message statistics:
Inq depth is 0
Outq depth is 0

      Sent      Rcvd
Opens:          1         1
Notifications:  0         0
Updates:        17        33
Keepalives:    1393      1392
Route Refresh:  0         0
Capability:     0         0
Total:         1411      1426
Minimum time between advertisement runs is 0 seconds
More
```

*Show BGP neighbours*

At our RR there has to be a Type 3 route for every VNI and a Type 2 route for every MAC inside the VNI of a VMM host:

```

BGP table version is 0, local router ID is 172.24.201.8
Status codes: s suppressed, d damped, h history, * valid, > best, i - internal
Origin codes: i - IGP, e - EGP, ? - incomplete
EVPN type-2 prefix: [2]:[ESI]:[EthTag]:[MAClen]:[MAC]
EVPN type-3 prefix: [3]:[EthTag]:[IPlen]:[OrigIP]

  Network          Next Hop          Metric LocPrf Weight Path
Route Distinguisher: 172.24.201.101:1
*>i[2]:[0]:[0]:[48]:[02:00:0a:0a:0a:07]
  172.24.201.101          100      0 i
*>i[3]:[0]:[32]:[172.24.201.101]
  172.24.201.101          100      0 i
Route Distinguisher: 172.24.201.101:2
*>i[3]:[0]:[32]:[10.99.254.1]
  10.99.254.1             100      0 i
Route Distinguisher: 172.24.201.101:3
*>i[3]:[0]:[32]:[10.99.254.1]
  10.99.254.1             100      0 i
Route Distinguisher: 172.24.201.101:4
*>i[3]:[0]:[32]:[10.99.254.1]
  10.99.254.1             100      0 i
Route Distinguisher: 172.24.201.101:5
*>i[3]:[0]:[32]:[10.99.254.1]
  10.99.254.1             100      0 i
Route Distinguisher: 172.24.201.101:6
*>i[3]:[0]:[32]:[10.99.254.1]
  10.99.254.1             100      0 i
Route Distinguisher: 172.24.201.101:7
*>i[3]:[0]:[32]:[10.99.254.1]
  10.99.254.1             100      0 i
--More--

```

Type 2 Route = MAC

Type 3 Route = IP of VTEP










Show BGP EVPN routes at RR

The syntax of the Type 2 and 3 route is like this:

- **Type 2 Route:** [2]:[ESI]:[EthTag]:[MACLen]:[MAC]
- **Type 3 Route:** [3]:[EthTag]:[IPLen]:[OrigIP]

Now we're able to run some checks. First we start a „PING from the first VM (running on the first VMM Host) to the second VM (running on the secondary VMM Host). Then we will check the Forwarding Tables of our „virtual switches“ – the Linux Bridged with their VXLAN interfaces attached to.

We will ping from SRC\_IP: 10.22.10.5 to DST\_IP: 10.22.10.7 – see the details of the VMs below:

<input type="checkbox"/>	ID	Owner	Group	Name	Status	Host	IPs	
<input type="checkbox"/>	187	oneadmin	oneadmin	vxlan-bgp-evpn-test-vni-10-3	RUNNING	ucs1.cloud.lab	10.22.10.7	
<input type="checkbox"/>	186	oneadmin	oneadmin	vxlan-bgp-evpn-test-vni-10-1	RUNNING	ucs2.cloud.lab	10.22.10.5	
<input type="checkbox"/>	185	oneadmin	oneadmin	vxlan-bgp-evpn-test-vni-10-2	RUNNING	ucs1.cloud.lab	10.22.10.6	
<input type="checkbox"/>	184	oneadmin	oneadmin	vxlan-bgp-evpn-test-vni-10-0	RUNNING	ucs2.cloud.lab	10.22.10.4	
<input type="checkbox"/>	183	oneadmin	oneadmin	vxlan-bgp-evpn-test-vni-20-1	RUNNING	ucs1.cloud.lab	10.22.20.4	
<input type="checkbox"/>	182	oneadmin	oneadmin	vxlan-bgp-evpn-test-vni-20-2	RUNNING	ucs2.cloud.lab	10.22.20.5	
<input type="checkbox"/>	181	oneadmin	oneadmin	vxlan-bgp-evpn-test-vni-20-0	RUNNING	ucs1.cloud.lab	10.22.20.6	
<input type="checkbox"/>	180	oneadmin	oneadmin	vxlan-bgp-evpn-test-vni-20-3	RUNNING	ucs2.cloud.lab	10.22.20.7	

*VM Details*

The result looks good – as expected:

## VNC Connected (unencrypted) to: QEMU (one-186)

```
localhost:~# ifconfig eth0
eth0      Link encap:Ethernet  HWaddr 02:00:0A:16:0A:05
          inet addr:10.22.10.5  Bcast:10.22.10.255  Mask:255.255.255.0
          UP BROADCAST RUNNING MULTICAST  MTU:1500  Metric:1
          RX packets:20 errors:0 dropped:7 overruns:0 frame:0
          TX packets:5 errors:0 dropped:0 overruns:0 carrier:0
          collisions:0 txqueuelen:1000
          RX bytes:1328 (1.2 KiB)  TX bytes:322 (322.0 B)

localhost:~#
localhost:~# ping -c 4 10.22.10.7
PING 10.22.10.7 (10.22.10.7): 56 data bytes
64 bytes from 10.22.10.7: seq=0 ttl=64 time=1.300 ms
64 bytes from 10.22.10.7: seq=1 ttl=64 time=0.635 ms
64 bytes from 10.22.10.7: seq=2 ttl=64 time=0.642 ms
64 bytes from 10.22.10.7: seq=3 ttl=64 time=0.619 ms

--- 10.22.10.7 ping statistics ---
4 packets transmitted, 4 packets received, 0% packet loss
round-trip min/avg/max = 0.619/0.799/1.300 ms
localhost:~# _
```

## ICMP PING

Let's dig a bit deeper down the layers. First let's have a look to the FDBs of the VMM hosts:

## FDB VNI 10 @VMM-HOST "UCS1"

```
1. # bridge fdb show dev vxlan10
2.
3. 12:37:89:6b:a3:49 master vxlan-br-10 permanent
4. 02:00:0a:16:0a:05 master vxlan-br-10
5. 02:00:0a:16:0a:05 vlan 1 master vxlan-br-10
6. 12:37:89:6b:a3:49 vlan 1 master vxlan-br-10 permanent
```



```
7. 00:00:00:00:00:00 dst 10.99.254.2 self permanent
8. 02:00:0a:16:0a:05 dst 10.99.254.2 self offload
```

#### FDB VNI 10 @VMM-HOST "UCS2"

```
1. # bridge fdb show dev vxlan10
2.
3. 02:00:0a:16:0a:07 master vxlan-br-10
4. 4e:ab:46:26:f3:39 vlan 1 master vxlan-br-10 permanent
5. 4e:ab:46:26:f3:39 master vxlan-br-10 permanent
6. 02:00:0a:16:0a:07 vlan 1 master vxlan-br-10
7. 00:00:00:00:00:00 dst 10.99.254.1 self permanent
8. 02:00:0a:16:0a:07 dst 10.99.254.1 self offload
```

Our SRC-VM has the MAC „02:00:0a:16:0a:05“ and runs on VMM host „UCS2“, which has the VTEP-IP „10.99.254.2“. The Forwarding Database of our VXLAN bridge at the VMM host „UCS1“ shows a entry for this MAC address pointing to the VTEP IP of the VMM host „UCS2“:

```
1. 02:00:0a:16:0a:05 dst 10.99.254.2 self offload
```

The DST-VM ha the MAC „02:00:0a:16:0a:07“ and runs on VMM host „UCS1“, which has the VTEP-IP „10.99.254.1“. The entry in the FDB of VMM host „UCS2“ looks like this:

```
1. 02:00:0a:16:0a:07 dst 10.99.254.1 self offload
```

The „00:00:00:00:00:00 “ MAC is the „Default-Route“ for any BUM-Traffic.

Looks good!

This information (Type 2 and Type 3 Routes) are distributed by our BGP EVPN. Let's check this layer.

First on our VMM host „UCS1“.

Display some information about our VNI 10:

```
1. ucs1# sh evpn vni 10
2.
3. VNI: 10
4. VxLAN interface: vxlan10 ifIndex: 20 VTEP IP: 10.99.254.1
5. Remote VTEPs for this VNI:
6. 10.99.254.2
7. Number of MACs (local and remote) known for this VNI: 2
8. Number of ARPs (IPv4 and IPv6, local and remote) known for this VNI: 0
```

The interesting information is our local VTEP IP and a list of the remote VTEPs for this VNI. Also there is some information about the MAC count inside the VNI.

To display more details we use the command „sh evpn mac vni 10“:

```
1. ucs1# sh evpn mac vni 10
2.
3. Number of MACs (local and remote) known for this VNI: 2
4. MAC                Type    Intf/Remote VTEP    VLAN
5. 02:00:0a:16:0a:05 remote 10.99.254.2
6. 02:00:0a:16:0a:07 local  one-187-0
```

Now we can close the circle to our FDBs 😊

And – of course – it's BGP we also took a look to the routing information:

```
1. ucs1# sh bgp evpn route vni 10
```

```
2.
3. BGP table version is 0, local router ID is 172.24.201.101
4. Status codes: s suppressed, d damped, h history, * valid, > best, i - internal
5. Origin codes: i - IGP, e - EGP, ? - incomplete
6. EVPN type-2 prefix: [2]:[ESI]:[EthTag]:[MAClen]:[MAC]
7. EVPN type-3 prefix: [3]:[EthTag]:[IPlen]:[OrigIP]
8.
9.      Network          Next Hop          Metric LocPrf Weight Path
10. *>i[2]:[0]:[0]:[48]:[02:00:0a:16:0a:05]
11.                10.99.254.2          0    100     0 i
12. *> [2]:[0]:[0]:[48]:[02:00:0a:16:0a:07]
13.                10.99.254.1                32768 i
14. *> [3]:[0]:[32]:[10.99.254.1]
15.                10.99.254.1                32768 i
16. *>i[3]:[0]:[32]:[10.99.254.2]
17.                10.99.254.2          0    100     0 i
18.
19. Displayed 4 prefixes (4 paths)
```

This is only a short, quick and dirty blog post about this impressive technology. Perhaps there are some projects for implementations like this. But that's a part, i'm not allowed to write it down here.

 **ALLGEMEIN, OPENNEBULA**