

## E3Fabric Internals

Items	internals
Docs creating date	Feb 12,2018
Docs finishing date	
Initial author	jzheng.bjtu@hotmail.com
Description	This article mainly introduces internals of the e3fabric

E3fabric provides a highly available and highly performed virtual network for NFV workloads with the support of the COTS hardware by migrating virtual network into infrastructure network while retaining scalability. we achieve this by eliminating the overhead of virtual networking processing in compute entity<sup>1</sup>, thus preserving all the processor resource to VNF business logics<sup>2</sup>. While in order to scale datacenter-widely, we use dedicated servers to do virtual switching in the infrastructure network<sup>3</sup>. To construct the virtual network with maximized throughput and minimized latency and jitter, we take the centralized routing as the forwarding policy, generalized speaking, in a SDN<sup>4</sup> way. Note this will be working together with customized container<sup>5</sup> technology to realize a IaaS for NFV<sup>6</sup>.

This paper is focusing on networking virtualization. the subsequent chapters will cover the specific parts.

### Infrastructure Network Resource Introduction

As mentioned, E3fabric is a link level network virtualization solution, it provides virtual layer 2 network in datacenter in a very efficient way with little overhead. Let's clarify the infrastructure network resource first.

This section includes compute node network resource, infrastructure switch network resource and network node network resource.

Compute nodes host virtualized network function, the connectivity to virtual network is obtained through x86 based PCIe NIC<sup>7</sup>.the NIC MAY support SR-IOV based

---

<sup>1</sup> The universe compute wrapper often includes bare metal, container, virtual machine, we call the hosting server **compute node** all through the paper.

<sup>2</sup> the compute entity employs hardware assisted virtual switching to connect to e3fabric network, i.e. SR-IOV NIC switch

<sup>3</sup> this is viable to do switching in a large scale infrastructure with COTS large volume server and switch, by employing DPDK as the data plane acceleration, VPP shows us a good example of how to scale. It can handle traffic at the rate of potential hundreds of Gigabit per second. we call these hosting server **network node**.

<sup>4</sup> quite clearly, we do not use OpenFlow or other standard control plane protocol, it's proprietary one.

<sup>5</sup> We plan to employ containerd(<https://containerd.io>), since it's been a project of CNCF, and open sourced, it's long term developed and supported by the community.

<sup>6</sup> Usually we call it the NFVIaaS, this is a platform that deliver cloud native VNF services.

<sup>7</sup> Network Interface Card

physical function(PF) and virtual function(VF) partition in order to be able to support multiple networks. That's to say, **the traffic originated from VNF is mapped to virtual network by selecting an 802.1q VLAN tag for it.** and SR\_IOV switch is responsible for multiplexing or de-multiplexing traffic on the same NIC(link). The model is depicted below:

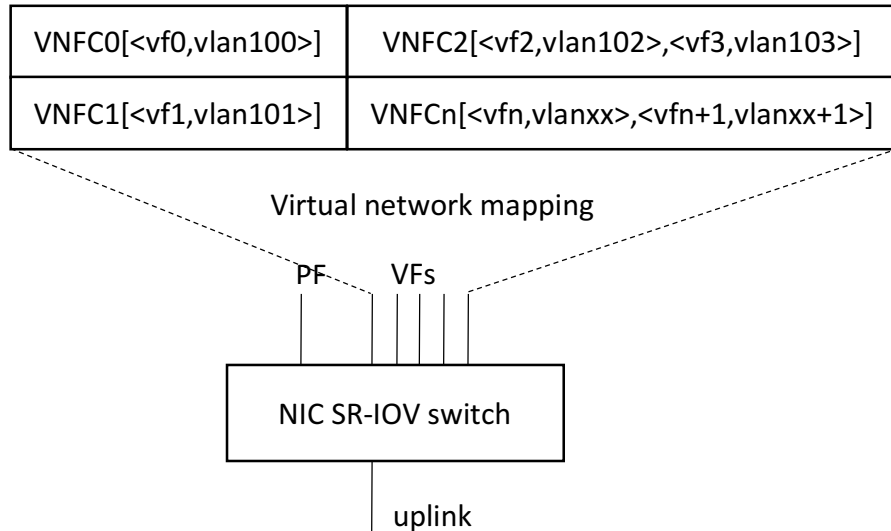


Figure 1 Host virtual network mapping

As Figure 1 illustrates, different VNFC<sup>8</sup>s map to separated VNFs, often every VF has chosen a VLAN ID for it. cases are common when VNFs need to attach more than one **virtual interface**<sup>9</sup> to the same virtual network, in such circumstance, these VFs share same VLAN ID<sup>10</sup>.

For administrative purpose, compute node should report these SR-IOV backed PFs and VFs to controller.

This solution requires no special functions for infrastructure switch except that it MAY support 802.1q based forwarding mechanism to optimize underlay traffic. as TRILL<sup>11</sup> emphasizes, the underlay infrastructure switch provides 802.1D(802.1Q) Layer2 connectivity, these switches can employ their own loop free technologies, So does E3Fabric.

As default deployment, all ports of a switch should work in TRUNK mode to allow all configured VLAN traffic to pass. Since these VLAN is locally scoped, we can configure them<sup>12</sup> in management plane.

Network node use PCIe NIC as with compute node does, it still employs SR-IOV to partition virtual functions in order to balance load to different CPUs in case NIC is of 25Gbps/40Gbps/100Gbps rate when a single CPU can not handle all the traffic in a burst traffic pattern. Unlike RSS<sup>13</sup> or Flow Director, EoMPLS<sup>14</sup> is not supported to share traffic among CPUs for most high volume NICs, SR-IOV partitioning works it around.

<sup>8</sup> VNFC: virtual network function component.

<sup>9</sup> From VNF's view, they would rather call the networking attach point virtual interface.

<sup>10</sup> Control Plane guarantees it.

<sup>11</sup> TRILL: Transparent Interconnect of Lots of Links

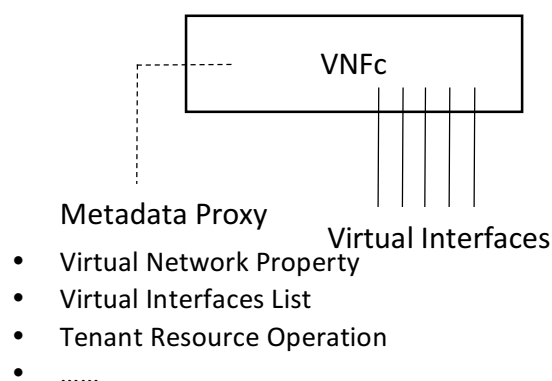
<sup>12</sup> allowed VLAN range is associated to LAN\_ZONE object.

<sup>13</sup> RSS: Receive Side Scaling.

<sup>14</sup> EoMPLS: Ethernet over MPLS, RFC 4448

## Virtual Network Resource introduction

For Cloud Native user/developer, what the VNF wants to know is what virtual network(s) it's accessing, it knows the exact virtual interface. NFVlaaS MUST expose these virtual interfaces along with their metadata to VNFs.



*Figure 2 virtual network resource overview*

More detail about metadata will be described in NFVlaaS documents. these virtual interfaces will be delivered to VNF container with the requested format.

## Virtual Network Orchestration Scenario

This chapter includes the several cases of Virtual Network Service Orchestration. Here we define two terms: **LANZONE affinity** and **anti-LANZONE affinity**.

- **LANZONE affinity:** the policy that compute nodes which connected to the Ethernet service must reside in the same LANZONE, the Ethernet Service involves no network nodes to do forwarding, instead, they share the same VLAN ID.
- **Anti-LANZONE affinity:** the policy that compute nodes MUST span across more than one LANZONE, and network nodes act as intermediate proxy to forward traffic from one LANZONE to another.

### In-Host Ethernet service orchestration

In-Host Ethernet Service is the service case where two or more VNFCs reside in the same compute node, multiple virtual interfaces share the same VLAN ID and exchange packets via SR-IOV NIC switch<sup>15</sup>. each single virtual interface in the host can join and leave the Ethernet service without involving network nodes.

### In-LANZONE Ethernet service orchestration

In-LANZONE Ethernet service is the service case where all the virtual interfaces of VNFCs are connected to the same LANZONE, comparing with In-Host Ethernet Service, the difference is these interfaces MAY span more than one compute node, but they are still in

---

<sup>15</sup> this is less optimized because packet flows through PCIe bus to NIC switch, then echoes back, it's still hairpin pattern. to improve performance, **we MAY use in-memory channel to replace it if there is no other VNFCs joining the virtual Ethernet service.**

the same LANZONE, even no Ethernet service instance is created at network nodes. In-Host Ethernet service is the special case of In-LANZONE Ethernet service.

#### Cross-LANZONE Ethernet service orchestration

Cross-LANZONE Ethernet service is the service case where the virtual interfaces MAY be connected to more than one LANZONE, the virtual interfaces which are connected to different LANZONES MAY have different VLAN IDs<sup>16</sup>, each of the virtual interfaces can join and leave the Ethernet service at any time. Once a LANZONE joins an Ethernet service, the VLAN ID is determined, subsequent virtual interfaces will share the same VLAN ID. the VLAN ID must be within the defined VLAN range.

---

<sup>16</sup> The VLAN ID is locally determined, scaling is not a problem.