

Design and implementation of a connectivity manager for virtual scalable network environments

Bachelorarbeit
von

Manuel Bergler

01. Dezember 2014 – 08. Februar 2015

Referent: Herr Prof. Dr. Thomas Baar
Betreuer: Herr Benjamin Reichel M. Sc.

Manuel Bergler
Urbanstr. 26
10967 Berlin

Hiermit versichere ich, dass ich die von mir vorgelegte Arbeit selbstständig verfasst habe, dass ich die verwendeten Quellen und Internet-Quellen vollständig angegeben habe und dass ich die Stellen der Arbeit – einschließlich Tabellen und Abbildungen –, die anderen Werken oder dem Internet im Wortlaut oder dem Sinn nach entnommen sind, auf jeden Fall unter Angabe der Quelle als Entlehnung kenntlich gemacht habe.

Berlin, den 08. Februar 2015

Manuel Bergler

Contents

List of figures	7
List of tables	8
List of algorithms	9
1 Introduction	10
1.1 Motivation	10
1.2 Network Architecture	10
1.3 Objective	11
1.4 Scope	11
1.5 Overview	11
2 Fundamentals and Related Work	13
2.1 Software-Defined Networking	13
2.1.1 Motivation	13
2.1.2 Software-Defined Networking Concept	15
2.1.3 SDN Architecture	15
2.1.4 OpenFlow	16
2.1.5 Open vSwitch	21
2.2 Cloud Computing Infrastructures	26
2.2.1 OpenStack	26
2.2.2 OpenStack Compute (Nova)	27
2.2.3 OpenStack Orchestration (Heat)	29
2.2.4 OpenStack Neutron	29
2.3 Conclusion	31
3 Requirements	32
3.1 Functional Requirements	32
3.1.1 Service-Level Agreement Enforcement	32
3.1.2 Optimal Virtual Machine Placement	32
3.1.3 Integration with Elastic Media Manager	32
3.2 Non-functional Requirements	33
3.2.1 Scalability	33
3.2.2 Modularity	33
3.2.3 Interoperability	33
4 State of the Art	34

4.1	Overview	34
4.2	OpenDaylight SDN Controller	34
4.3	Ryu SDN Controller	35
4.4	OpenStack Neutron - QoS Extension	36
4.5	Problem Statement	37
4.6	Conclusion	37
5	Design	38
5.1	Architecture Overview	38
5.2	Interface between Manager & Agent	39
5.3	Design of Connectivity Manager	40
5.3.1	Algorithm for Instance Placement	41
5.3.2	QoS Manager	42
5.4	Design of Connectivity Manager Agent	43
5.4.1	API	44
6	Implementation	45
6.1	Environment	45
6.2	Connectivity Manager Components and Operations	47
6.3	Connectivity Manager Agent Components and Operations	52
6.4	Conclusion	54
7	Evaluation	55
7.1	Feature Analysis	55
7.2	Connectivity Manager Integration with NUBOMEDIA project	55
7.3	Network Performance Analysis	55
7.3.1	Test-bed Configuration	55
7.3.2	Installation of Connectivity Manager Agent	56
7.3.3	Topology Definition	57
7.3.4	Scenario 1: Without Instance Placement Engine & QoS enabled . . .	58
7.3.5	Scenario 2: With Instance Placement Engine enabled, but without QoS	59
7.3.6	Scenario 3: Instance Placement Engine and QoS Manager enabled . .	60
7.4	Conclusion	63
8	Conclusion	65
8.1	Summary	65
8.2	Problems Encountered	65
8.3	Future Work	65

Bibliography	67
A List of acronyms	69

List of Figures

2.1	"Classical" switch components	14
2.2	Software-Defined Network architecture	15
2.3	OpenFlow Network Architecture	16
2.4	OpenFlow Switch components	17
2.5	OpenFlow pipeline processing	18
2.6	Packet flow through an OpenFlow switch	19
2.7	OpenFlow QoS as a meter	20
2.8	Architecture of Open vSwitch: kernelspace / userspace separation	22
2.9	Visualization of the interaction of the ovs-vsctl tool	23
2.10	QoS Queues attached to a Port in OVS	24
2.11	HTB: qdisc classes linked together as a tree	25
2.12	Interaction among OpenStack services	27
2.13	OpenStack Compute service	28
2.14	Neutron modular framework, including ML2 drivers	30
4.1	Architecture of OpenDaylight Virtualization edition	35
4.2	Ryu architecture	36
4.3	Neutron QoS Extension architecture	36
5.1	High-level architecture of the Connectivity Manager	38
5.2	Minimized architecture of the Connectivity Manager and its integrations	39
5.3	Workflow: Deployment of stack & Assignment of QoS policies	39
5.4	Class diagram: Metaclass (Interface) and implementation (Service) instantiated by Factory Agent	40
5.5	Class diagram: Connectivity Manager service and its helper classes	41
5.6	Deployment: Topology object from EMM	41
5.7	Instance Placement activities	42
5.8	Abstract design of the Connectivity Manager Agent	43
5.9	Class diagram: Connectivity Manager Agent - All packages	44
6.1	Implementation focus	45
6.2	Check resource utilization of hosts	47
6.3	Get total amount of required resources for topology	48
6.4	Deployment feasibility check	48
6.5	Single host deployment check	49
6.6	Set AZ per Unit	50
6.7	Method for setting QoS for all Units	50
6.8	Activity diagram: Get list of hosts	52

6.9	Activity diagram: Set QoS rates for all servers	53
6.10	Activity diagram: Get OVS Port ID for server	53
6.11	Activity diagram: Get QoS ID for OVS Port	54
7.1	Scenario 1: Placement of servers without Connectivity Manager	58
7.2	Scenario 1: Bandwidth comparison	59
7.3	Scenario 2: Placement of servers with Connectivity Manager	59
7.4	Scenario 2: Bandwidth comparison	60
7.5	Scenario 3: Placement of servers using the Connectivity Manager with QoS enabled	60
7.6	Scenario 3: Bandwidth comparison	63

List of Tables

2.1	Fields within a Flow table	18
2.2	Fields within a meter entry in the meter table	20
2.3	Fields within a meter band	20
2.4	Qdisc leaf configuration	25

List of Algorithms

Chapter 1

Introduction

1.1 Motivation

The demands on networks have changed dramatically in the past two decades, with an ever-growing number of people and devices relying on interconnected applications and services. The underlying infrastructure has been left mostly unchanged and is reaching its limits. In order to resolve this, Software Defined Networking (SDN) is going to be extending and replacing parts of traditional networking infrastructures. SDN separates the network into control and forwarding planes and therefore allows a more efficient orchestration and automation of network services.

The use of cloud-based services, with not only competitive pricing but also high-availability and fast network access, is taking over traditional purely hardware-based data centers. The ease of administration and deployment of new servers as Virtual Machines (VMs) on the fly make it possible to effortlessly create a topology of multiple servers.

Network services have different requirements, depending on the type of data packets and their characteristics. A classification and prioritization of network traffic can be achieved through Quality of Service (QoS). A new approach has to be made to enable and configure the use of QoS in virtualized cloud infrastructures like OpenStack through an API. The traffic flow should be controlled right from the deployment of Virtual Machines on.

1.2 Network Architecture

Today's traffic patterns, the rise of cloud computing and "big data" to only name a few examples, are exceeding the capacity of classic network architectures. With scalable computing and storage the common-place tree-structured networking infrastructure are not efficient and manageable enough.

The increasing complexity of problems that have to be faced in networks and the need to control network traffic through software, are only a selection of the reasons why the Open Networking Foundation (ONF) developed an approach called Software-Defined Networking (SDN).

SDN is a leading-edge approach where the network control is separated from the forwarding functions. The centralized network intelligence allows programming the network, without a need to access its underlying infrastructure. With the help of this architecture a shift of today's networks leading to more flexibility, programmability and scalability is taking place.

1.3 Objective

The primary objective of this work is the development of a network orchestrator which is able to apply Quality of Service to the network interfaces of Virtual Machines. These Virtual Machines are deployed on compute nodes that run the OpenStack Nova-service and are interconnected with Open vSwitches. Another task of the Connectivity Manager is to select which OpenStack compute node new servers should run on. This decision takes different runtime parameters into account. The CM should be applicable in environments at scale in terms of amount of servers.

1.4 Scope

The scope of this work includes a Connectivity Manager which is interfaced with the already existing implementation of the Elastic Media Manager. The second application is a Connectivity Manager Agent that runs on the cloud controller within the OpenStack infrastructure in order to provide access to the hosts of OpenStack Nova. These two components have to be implemented and integrated with the existing SDN software that OpenStack makes use of. The entire virtual network switching will utilize Open vSwitch, which is a fully-compliant OpenFlow switch. As a reference for a cloud infrastructure, multimedia communications like the Nubomedia project will be used as an evaluation. The deployment of this cloud topology is tested on different performance characteristics such as maximal network bandwidth, jitter, CPU utilization and memory usage.

In virtualized cloud infrastructure like OpenStack, the placement of Virtual Machines (VMs) on a particular compute node can be decided on by comparing different run-time parameters. The network connectivity between those VMs has to be prioritized and classified depending on the type of service that the server provides.

Currently there are a number of solutions for managing network connectivity between virtual servers. A comparison and their current limitations follows in the next section. The selected approach is the extension of existing network control and management services with Quality of Service (QoS) capabilities and the ability to choose a host for the deployment of the topology. In support of the thesis the Connectivity Manager will be implemented and the differences in network performance will be shown in one use-case.

1.5 Overview

Chapter 1 begins with the motivation for this thesis and gives a brief introduction into the objectives and the scope.

Chapter 2 gives an overview of traditional network concepts and a introduction to SDN and its components. Furthermore the different services that OpenStack consists of will be described.

Chapter 3 conceptualizes the state-of-the-art solutions that are currently available and evaluates their implementation and limitations.

Chapter 4 contains an analysis of requirements.

Chapter 5 gives an architectural overview of the Connectivity Manager. Moreover design aspects are introduced and illustrated in relation to their requirements.

Chapter 5 examines the implementation of the Connectivity Manager and Agent.

Chapter 6 evaluates the network performance tests on the basis of a particular use-case.

Chapter 7 summarizes the results of this work and gives an overview on possible future work.

Chapter 2

Fundamentals and Related Work

2.1 Software-Defined Networking

The origins of Software-Defined Networking (SDN) go back to as early as 1995, however the first implementations were only developed in 2001 and the promotion of SDN began with the foundation of the non-profit industry consortium Open Networking Foundation (ONF) in 2011. [1] The ONF is dedicated to push open standards like OpenFlow into the industry and adapt them. In this following section a brief overview of the SDN architecture and concepts, including the OpenFlow standard and Open vSwitch, is given.

2.1.1 Motivation

Today's internet is part of the modern society for private users, enterprises and vital infrastructure services. Networks are required to evolve in order to address the challenges that are entailed with new applications, services and a growing number of end-users.

With a more detailed view on the challenges of current network infrastructures one comes to see the following limitations: [2]

- **Inability to scale:** With the expansion of data centers, networks grow too. Configuring and managing these additional network devices comes at a high administrative effort. With the virtualization of data-centers network traffic patterns become more and more dynamic and unpredictable. With multi-tenancy an additional complication is introduced, because different end-users and services need require distinct network characteristics. Network management and scaling of the architecture cannot be achieved with a manual configuration of the underlying infrastructure.
- **Complexity:** In the past decades various new networking protocols have been adapted by the industry. To add or move a device, multiple of the already existing switches, routes and firewalls must be touched in order to manage and introduce new protocol-based mechanisms. With the virtualization of servers the amount of network interfaces that need connectivity is on the rise. The distribution of applications across a number of virtual machines are another demand that the current fairly static connections cannot adapt to.

- **Inconsistent policies:** For IT departments to apply a network- or data center-wide policy, a lot of devices and mechanisms may need to be reconfigured. Virtual Machines can be created and rebuilt within no time, but if access or security rules need to be updated then the benefits of this dynamic can't be utilized to its full extent.
- **Vendor dependence:** Standards are needed to match the requirements of the markets with the capabilities of networks and enable data center operators to customize the network to specific environments.

Traditionally decisions about traffic-flows through the network are made directly by each network device, because the control logic and forwarding hardware are tightly coupled.

2.1.1.1 Classical Switches & Routers

Packet forwarding (data plane) and routing decisions (control plane) in classical switching and routing are both within one device. The main components depicted in Figure 2.1 have the following functions:

1. The **forwarding path** typically handles data path operations for each packet. It generally consists of Application-Specific Integrated Circuits (ASIC), network-processors or general-purpose processors that forwards frames and packets at wire speed (line-rate). Their lookup functions can be further increased with memory resources like Content Addressable Memory (CAM) or Ternary Content Addressable Memory (TCAM) to contain the forwarding information.
2. The elements in the **control plane** are based on general-purpose processors that provide services like routing and signaling protocols, including ARP, MAC Learning and forwarding tables.

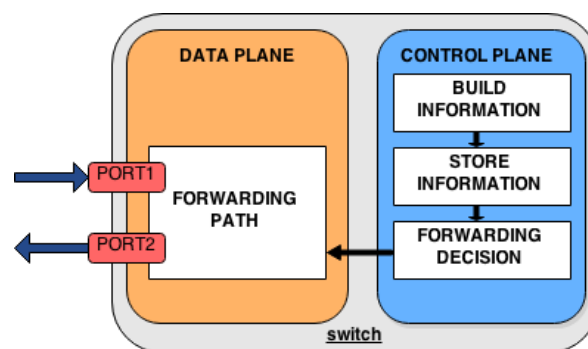


Figure 2.1: "Classical" switch components

A switch consists of multiple ports for incoming and outgoing data. Internal forwarding tables classify the packets and forward them to one specific or multiple ports. It does so by collecting MAC addresses and storing their corresponding port in tables. Layer 2 switches also support the segregation into virtual LANs (VLAN), which enables the network operator to logically isolate networks that share a single switch.

Routers forward packets on the Network layer (Layer 3) and routing-decisions are made based on IP addresses. They contain a routing table wherein paths to neighbouring networks are

stored, so that packets can be forwarded to their destination IP address. Other features that can be configured with routers are Quality of Service (QoS), Network Address Translation (NAT) and packet filtering.

The main differences between the classical architecture and SDN will be further described in the next section.

2.1.2 Software-Defined Networking Concept

SDN represents a novel dynamic, manageable, cost-effective and adaptable architecture [3] that is built to serve the dynamic infrastructures that are needed as a backbone for modern data centers. Opposed to the traditional approach, network control and forwarding functions are decoupled and thus can be programmed and divided into different applications and network services. The work of the Open Networking Foundation laid out the OpenFlow standard as the base for modern SDN solutions.

2.1.3 SDN Architecture

SDN separates the architecture into three distinct layers that communicate with each other through different APIs. This separation is shown in Figure 2.2.

- **Infrastructure Layer:** here all physical and virtual devices (e.g. switches and routers) that are capable of the OpenFlow Protocol provide forwarding mechanisms on different Network Layers.
- **Control Layer:** represents the 'network intelligence' and collects a global view of the network, by communicating with the switching elements through the so-called Southbound API.
- **Application Layer:** consists of business applications that allow the network operator to extend the SDN controller on an abstracted level, without being tied to the actual details of the implementation of the infrastructure. This communication with the Control Layer is also called Northbound API.

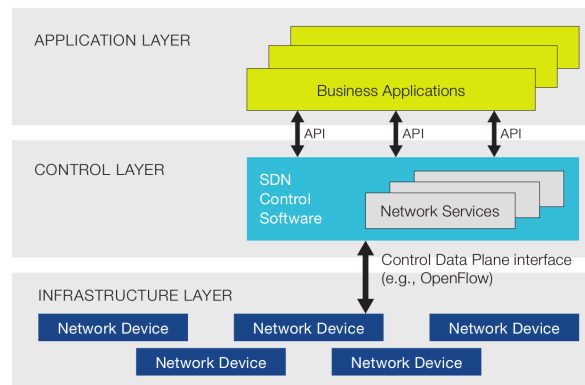


Figure 2.2: Software-Defined Network architecture

2.1.4 OpenFlow

With OpenFlow the Open Networking Foundation defined the first standard communications interface between the SDN architecture's control and forwarding layers. It enables manipulation and direct access to the forwarding plane of physical as well as virtual (hypervisor-based) network devices such as switches and routers. [2]

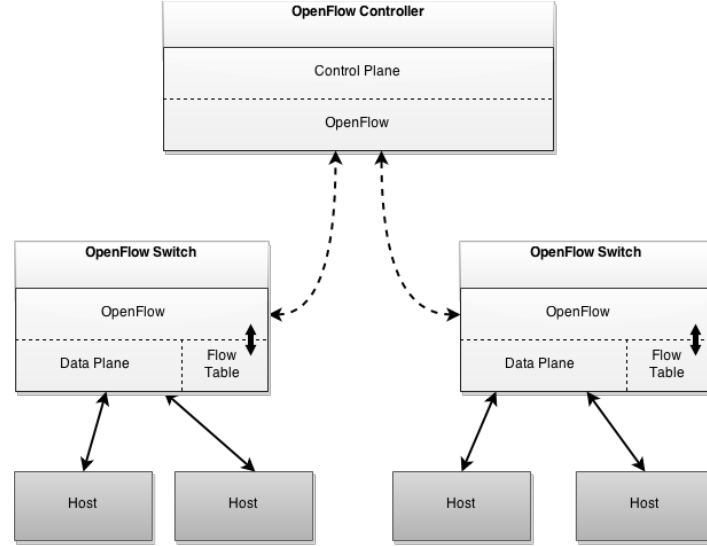


Figure 2.3: OpenFlow Network Architecture

OpenFlow first of all stands for the communications protocol that is used by SDN controllers to fetch information and configure switches. Additionally it is a switch specification with full support of OpenFlow.

Many of the OpenFlow-enabled switches and controllers up to today still solely support the OpenFlow version 1.0 (released in December 2009). The newest version at this date is 1.4, however the characteristics of OpenFlow in this work will focus on version 1.3, since it is the most recent specification which is supported by Open vSwitch.

The main features added since version 1.0 are, among others, support for VLANs, IPv6, tunnelling and per-flow traffic meters. [4]

Generally the switches are backwards-compatible to version 1.0. In the following description the focus lies on the required features of all OpenFlow capable devices, however it has to be mentioned that some of them are optional features.

2.1.4.1 OpenFlow Controller

The OpenFlow controller is segregated from the switch and has two interfaces. The north-bound interface is an API to the application layer for developing applications that control the network. The southbound interface connects with the underlying switches using the OpenFlow protocol.

2.1.4.2 OpenFlow Switch

There are two varieties of OpenFlow-compliant switches: [5]

- **OpenFlow-only:** in these switches all packets are processed by the OpenFlow pipeline and have no legacy features.
- **OpenFlow-hybrid:** support OpenFlow and normal Ethernet switching (including traditional L2 Ethernet switching, VLAN isolation, L3 routing, ACL and QoS). Most of the commercial switches that are available on the market today are of this type.

An OpenFlow switch includes one or more flow tables and a group table, which have the function of carrying out packet lookups and forwarding. Another component is the OpenFlow channel that is connected to the external controller [5].

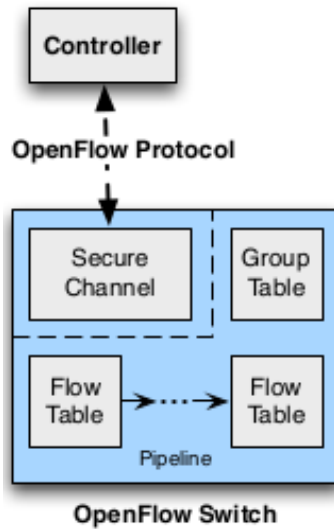


Figure 2.4: OpenFlow Switch components

Through the connection using the OpenFlow protocol, the controller is able to add, update and delete flow entries in Flow tables. This action can be performed either reactively or proactively. Sets of flow entries are stored in each flow table and each flow entry consists of *match fields*, *counters*, and a set of *instructions* used for matching packets.

Matching flow entries starts in the first flow table, however it may continue to additional flow tables. The first matching entry from each table is used and the instruction that is linked with that specific entry is performed. For packets without any matches a table-miss flow entry can be configured. Flow entries are usually forwarded to a physical port.

The instructions can either include actions or else modify the pipeline processing. Packet forwarding, packet modification and group table processing are available actions. Packets can be permitted to be sent to other tables with pipeline processing and additionally metadata can be exchanged between tables.

Packets can also be directed to a group, which contains a set of actions for flooding and more complex forwarding semantics (e.g. multipath, fast reroute and link aggregation).

2.1.4.3 OpenFlow Ports

OpenFlow ports are the network interfaces used for passing packets between OpenFlow processing and the rest of the network [5]. There are various types of ports that are supported by OpenFlow. This section will give a short overview about this port abstraction. Incoming OpenFlow packets enter the switch on an ingress port, are then processed by the OpenFlow pipeline and forwarded to an output port. (See OF Tables figure for processing).

There are three types of OpenFlow ports that must be supported by an OpenFlow switch:

- **Physical ports:** are hardware interfaces on a switch.
- **Logical ports:** don't directly interact with a hardware interface.
- **Reserved ports:** contain generic forwarding actions (e.g. sending to the controller, flooding or forwarding using traditional switch processing)

2.1.4.4 OpenFlow Tables

Pipeline Processing

The OpenFlow pipeline specifies how packets correspond with each of the flow tables [5].

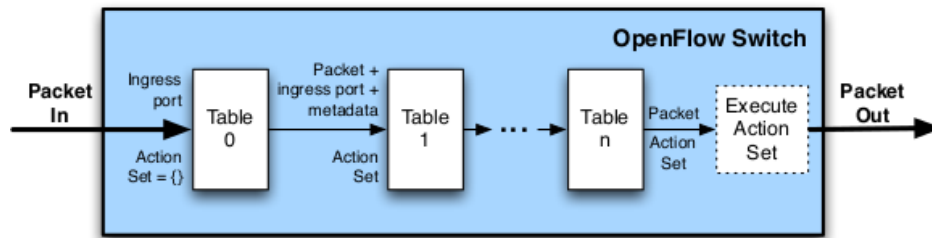


Figure 2.5: OpenFlow pipeline processing

As illustrated in the figure, each packet is matched against the flow entries starting at the first flow table, called flow table 0. The outcome of the match then decides if other of the sequentially numbered tables may be used. In the following sections the components of the Flow table, the matching procedures and different instructions will be described.

Flow Table

A flow table contains flow entries, each of which consist of the following fields [5]:

Match Fields	Priority	Counters	Instructions	Timeouts	Cookie
--------------	----------	----------	--------------	----------	--------

Table 2.1: Fields within a Flow table

- **match fields:** ingress port, packet headers and optionally metadata
- **priority:** sets the priority of the flow entry
- **counters:** is updated for matching packets

- **instructions:** to alter the action set or pipeline processing
- **timeouts:** set maximum amount of time or idle time before expiration of the flow
- **cookie:** is a opaque data value chosen by the controller

Each flow table entry is uniquely identifiable by its match fields and priority.

Packet Matching

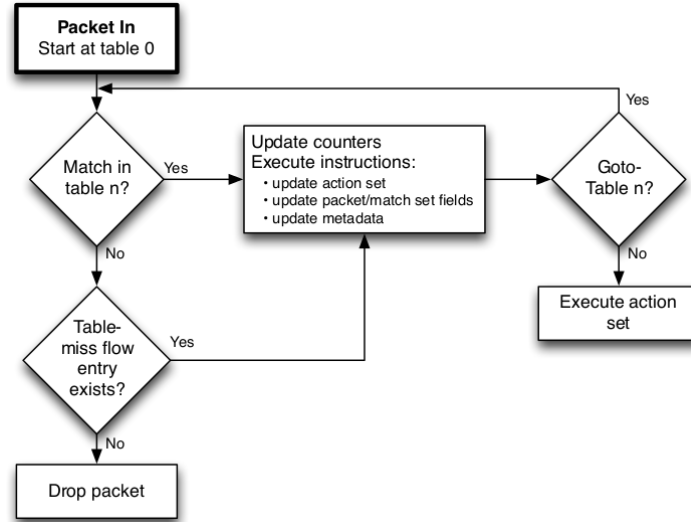


Figure 2.6: Packet flow through an OpenFlow switch

Upon a packet's arrival at the Flow Table, the packet match fields are extracted and used for the table lookup [5]. They include different packet header fields. Additionally the ingress port can be matched with metadata fields. If the values in the packet match fields are equal, then only the flow entry with the highest priority is selected. Furthermore the associated counters are updated and the instruction set is applied.

When the instruction set associated with a matching flow entry does not specify a next table, the pipeline processing stops. Only then the packet is processed with it's action set and in most cases forwarded as shown on the right-hand side of Figure 2.6. However, if the lookup phase does not match any of the entries, a table-miss event occurs.

Table-miss

Each flow table must support a table-miss flow entry which specifies how to process packets that are unmatched by other flow entries. The instructions associated with this entry are very alike to any other flow entries, packets are either forwarded to other controllers, dropped or it is continued with the next flow table. In case the table-miss flow entry is non-existent unmatched packets are dropped by default.

Group Tables

A group table consists of group entries and it provides a way to direct the same set of actions as part of action buckets to multiple flows. A flow entry is pointed to a group and enables additional methods of forwarding (e.g. broadcast or multicast).

Meter Tables

Meters are on a per-flow level and allow OpenFlow to implement various QoS operations, such as rate-limiting, but it can also be combined with per-port queues to implement more complex QoS like DiffServ.

The main components of a meter entry in the meter table are:

Meter Identifier	Meter Bands	Counters
------------------	-------------	----------

Table 2.2: Fields within a meter entry in the meter table

- **meter identifier:** a 32 bit unsigned integer uniquely identifying the meter
- **meter bands:** each meter band specifies the rate of the band and the action that is triggered by exceeding the limit
- **counters:** is updated when packets are processed by the meter

The rate of packets assigned to a meter are measured and controlled. Meters are directly attached to flow entries, as opposed to queues that are attached to ports. A meter is able to have one or more meter bands, each of which specifies the rate and the way packets should be handled. If the current measured meter rate reached the rate-limit, the band applies an action.

A meter band is identified by its rate and consists of the following fields:

Band Type	Rate	Counters	Type specific arguments
-----------	------	----------	-------------------------

Table 2.3: Fields within a meter band

- **band type:** defines how the packets are processed
- **rate:** selects the meter band for the meter and defines the lowest rate at which the band can apply
- **counters:** is updated when packets are processed by the meter
- **type specific arguments:** some band have optional arguments

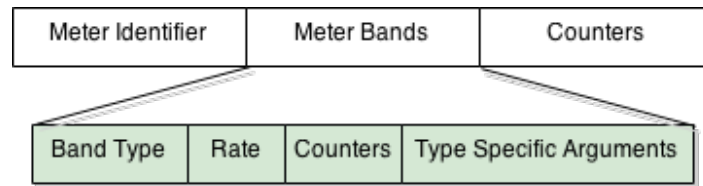


Figure 2.7: OpenFlow QoS as a meter

Instructions

Instructions are executed when a packet matches the flow entry [5]. Their result is either a change to the packet, action set and/or pipeline processing. There are different instruction types and some of them are required for an OpenFlow-enabled switch whereas others are optional:

- **Meter *meter-id*:** Directs packet to the specified meter. The packet may be discarded as the result of the metering.
- **Apply-Actions *action(s)*:** Applies the specific action(s) instantly, without any change to the Action Set.
- **Clear-Actions:** Immediately clears all the actions in the action set.
- **Write-Actions *action(s)*:** Merges the specified action(s) into the current action set.
- **Write-Metadata *metadata / mask*:** Writes the masked metadata value into the metadata field.
- **Goto-Table *next-table-id*:** Indicates the next table in the processing pipeline.

A maximum of one instruction of each type is associated with a flow entry and they are executed in the order as specified by the given list. Flow entries can also be rejected if the switch is not able to execute its instruction.

Action Set

An action set is associated with each packet, which is empty by default [5]. The action set can be modified using a *Write-Action* or a *Clear-Action* instruction. If there is no *Goto-Table* instruction within the instruction set of a flow entry the pipeline processing is halted and the actions in the action set of the packet are executed.

Actions

The following action types are available on OpenFlow-enabled switches [5]:

- **Output:** A packet is forwarded to a specified OpenFlow port.
- **Set-Queue:** Sets the Queue ID for a packet. This ID helps determining which Queue is used for scheduling and forwarding the packet. This forwarding behaviour allows to enable basic QoS support on a port.
- **Group:** Process the packet through the specified group.
- **Push-Tag/Pop-Tag:** The ability to push/pop tags such as VLAN.
- **Set-Field:** Modifies the values of header fields in a packet.
- **Change-TTL:** Set the values of IPv4 TTL, IPv6 Hop Limit or MPLS TTL in a packet.

2.1.5 Open vSwitch

2.1.5.1 Concept & Functionality

Open vSwitch (OVS) is an open source software switch that is used in virtualized server environments. It is able to forward traffic between virtual network interfaces (e.g. connected to Virtual Machines) and the physical network, as well as between different VMs on the same physical host. It can be controlled using OpenFlow and the OVSDB management protocol. It can run on any Linux-based virtualization platform i.e. KVM, VirtualBox, XEN, ESXi and is part of the mainline kernel as of Linux 3.3 but can be run on kernel 2.6.32 and newer [6].

OVS supports the following features [7]:

- 802.1Q VLAN model
- Link Aggregation Control Protocol (LACP)
- GRE and VXLAN tunneling
- fine-grained QoS control
- OpenFlow
- per VM interface traffic policing
- High-performance forwarding using a Linux kernel module

The internals of OVS are as follows [8]:

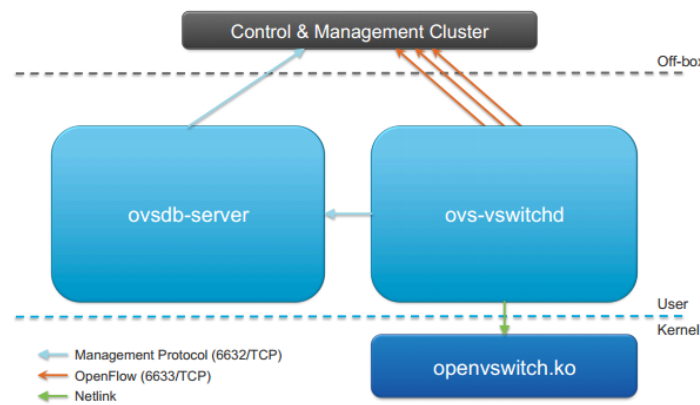


Figure 2.8: Architecture of Open vSwitch: kernelspace / userspace separation

The daemon which implements the switch is *ovs-vswitchd* and is shipped with an additional Linux kernel module for flow-based switching that it communicates with using the netlink protocol [9]. The configuration for the switch is queried from a lightweight database server named *ovsdb-server*. Generally the decision about how a packet is processed is made in userspace, yet all following packets hit the cached entry in the kernel.

It is also possible to run the switch entirely in userspace, but it decreases the performance drastically.

2.1.5.2 OVSDb

Each Open vSwitch daemon has a database (OVSDb) that holds its configuration [10]. The database is divided into multiple different tables with different purposes, with the ones related to this work outlined below:

- **Open_vSwitch:** Open vSwitch configuration
- **Port:** Port configuration
- **Interface:** A physical network device within a Port

- **QoS:** Quality of Service configuration
- **Queue:** Output queue for QoS
- **Controller:** OpenFlow controller configuration
- **Manager:** OVSDb management connection

2.1.5.3 Open vSwitch Management

Multiple different configuration utilities exist for OVS, however only `ovs-vsctl` is explained in this section. It is used for querying and updating the configuration of the switch through interaction with the `ovsdb-server`.

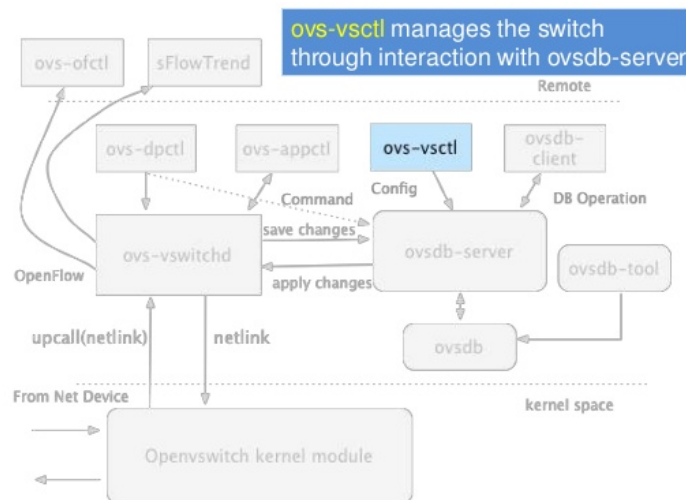


Figure 2.9: Visualization of the interaction of the `ovs-vsctl` tool

The tool also configures `ovs-vswitchd` and can be seen as a high-level interface for the database. The commands below are available for the OVS configuration that is needed to get a basic network setup running:

- `ovs-vsctl add-br %bridge%`
- `ovs-vsctl list-br`
- `ovs-vsctl add-port %bridge% %port%`
- `ovs-vsctl list-ports %bridge%`
- `ovs-vsctl get-manager %bridge%`
- `ovs-vsctl get-controller %bridge%`
- `ovs-vsctl list %table%`

2.1.5.4 QoS

With Open vSwitch QoS can be configured for ports that virtual machines are connected to. The minimal and maximal rate-limits can be defined in bit / s and applied to a queue. With the use of this feature the network traffic from the switches internal bridge to a virtual network interface can be shaped (egress). The QoS port assurances in OVS make use of the 'tc' implementation which is included in the Linux kernel.

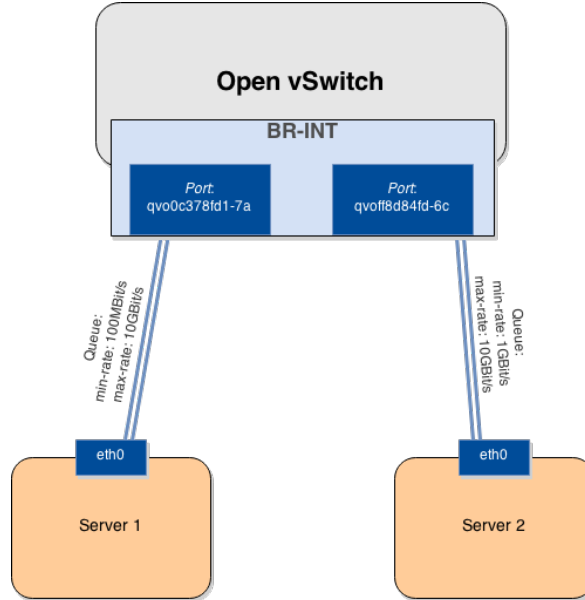


Figure 2.10: QoS Queues attached to a Port in OVS

Traffic control (tc) uses 'queuing discipline' (qdisc) for configuring the network interface [11]. They are the fundamental schedulers used under Linux. When a packet is sent, it is enqueued to the qdisc for the interface and shortly after the kernel is trying to get as many packets as it can from the qdisc, so they can be forwarded to the network adaptor driver. By default the 'pfifo_fast' qdisc is set in the kernel, which is a pure 'First In, First out' queue.

When setting QoS in OVS a classful qdisc named Hierarchy Token Bucket (HTB) is used. HTB is meant as a more understandable, intuitive and faster replacement for the Class Based Queuing (CBQ) qdisc in Linux [12]. It helps to control the use of the outbound bandwidth on a given link.

For further understanding the following graph shows a a HTB tree that contains on qdisc root and two classes (leaves) with different configurations:

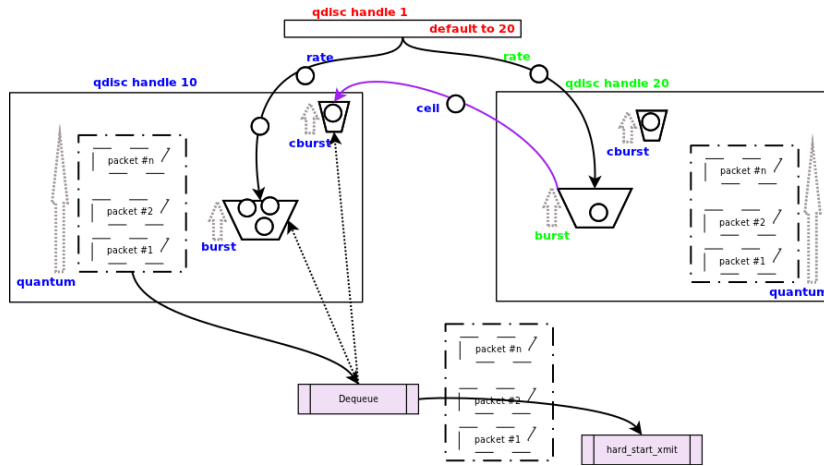


Figure 2.11: HTB: qdisc classes linked together as a tree
[13]

The terminology that is used for defining HTB classes is as follows:

- **burst:** the size of the token bucket
- **rate:** the speed at which tokens are generated and put in the bucket (in the leaf)
- **quantum:** the amount of bytes to serve from a single leaf at once
- **ceil:** the definite upper class rate
- **cburst:** the burst for ceil (gets computed)

This tree was created using the following commands:

```
$ tc qdisc add dev eth0 root handle 1: htb default 20
$ tc class add dev eth0 parent 1:0 classid 1:10 htb rate 200kbit
  ceil 400kbit prio 1 mtu 1500
$ tc class add dev eth0 parent 1:0 classid 1:20 htb rate 200kbit
  ceil 200kbit prio 2 mtu 1500
```

In figure 2.11 a root named 'qdisc handle 1' which applies filters to decide the direction of packets is depicted. It includes two leaves: 'qdisc handle 10' and set as default another leaf named 'qdisc handle 20'.

With the above listed commands the leaves were configured with the following parameters:

Leaf parameter	#10	#20
Rate value	200 KBit/s	200 KBit/s
Ceil value	400 KBit/s	200 KBit/s
Priority	1	2

Table 2.4: Qdisc leaf configuration

The difference between the rate and ceil value in leaf #10 means it can borrow 200 KBit/s more than it's actual rate, while leaf #20 cannot because its rate value equals the ceil value.

The HTB leaves always match one of the following three status [13]:

- **HTB_CAN_SEND:** the class can send using its own tokens
- **HTB_CANT_SEND:** the class can't send nor borrow and no packets are permitted to leave the class
- **HTB_MAY_BORROW:** the class can't send using its own tokens, but can try to borrow from another class

When a group of packets enters *tc* that are marked with the flag #10 and are directed to it, however the bucket does not contain enough tokens to let the first packets through then it will try to borrow tokens from leaf #20. The quantum value defines the maximal amount that leaf #10 will try to send at once. It is 1500 bytes, as that is the default MTU. If the first packet has a size of 1400 bytes and the bucket in leaf #10 has enough tokens for 1000 bytes it will try to borrow the remaining 400 bytes from the neighbor leaf. The quantum value can be used to privilege a leaf compared to others, because it specifies the tokens that it can borrow from other leaves in the tree.

By configuring *tc* the quantum value is not directly managed. It is automatically calculated based on an intermediary parameter named *r2q* with the following formula: $quantum = \frac{rate}{r2q}$. This value is set on the root qdisc and is used across all of its leaves. Its default value is set to 10.

2.1.5.5 GRE

Generic Routing Encapsulation (GRE) tunnels are used with Open vSwitch [14] to encapsulate isolated layer 2 network traffic between compute nodes. This helps enabling tenant networks in Neutron without the need configure a switch to trunk a range of VLANs.

2.2 Cloud Computing Infrastructures

2.2.1 OpenStack

The OpenStack project was founded by Rackspace Cloud and NASA, however up to this day more than 200 companies are contributing to it. With OpenStack it is possible to design, deploy and maintain a private or public cloud installation. It is a flexible, scalable and open-source approach that combines multiple technologies into one Infrastructure-as-a-Service (IaaS) [15]. All of the interrelated services include an API that offers administrators different ways of controlling the cloud, be it through a web-interface, a command-line client or a software development kit. All of the core components that are part of OpenStack are implemented in Python.

Conceptual architecture

The following graph shows the interaction between different OpenStack services that are involved in launching a virtual machine.

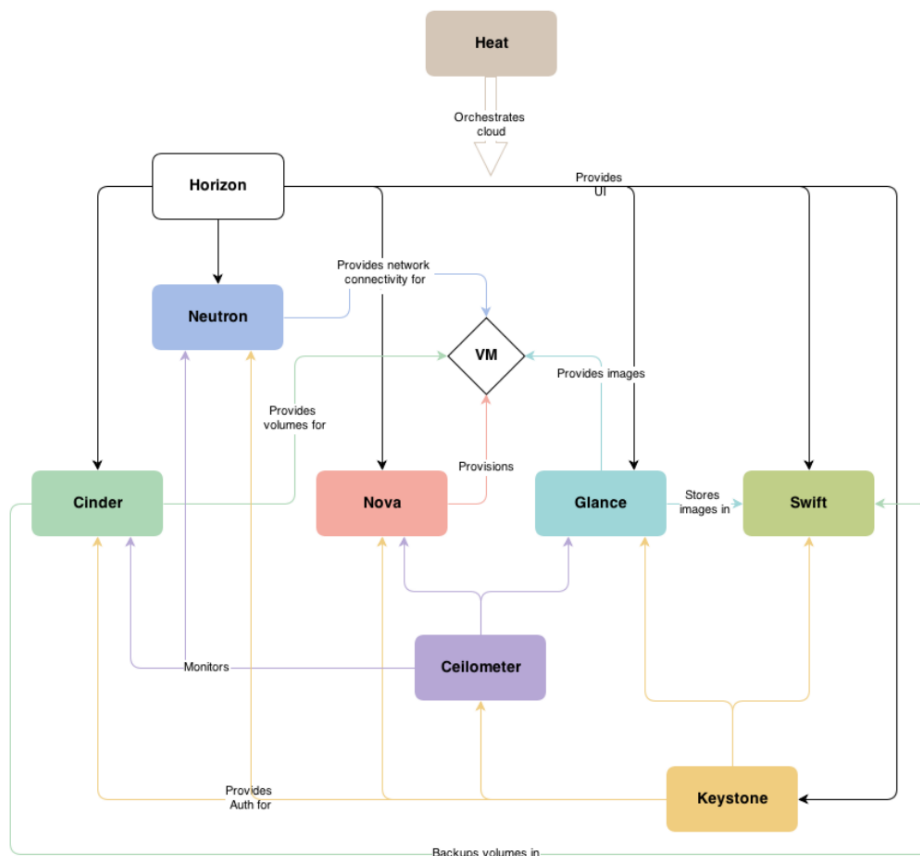


Figure 2.12: Interaction among OpenStack services
[17]

2.2.2 OpenStack Compute (Nova)

Nova is used to provision and manage large amounts of virtual machines . It supports different hypervisors and the number of physical hosts running the compute services can be scaled horizontally with no requirement of hardware resources from specific vendors. Hosts that provide Nova services are also called 'Compute Nodes'. Data center can be divided into so called tenants, which are isolated users with their own servers, security groups and externally reachable IP addresses (Floating IP addresses).

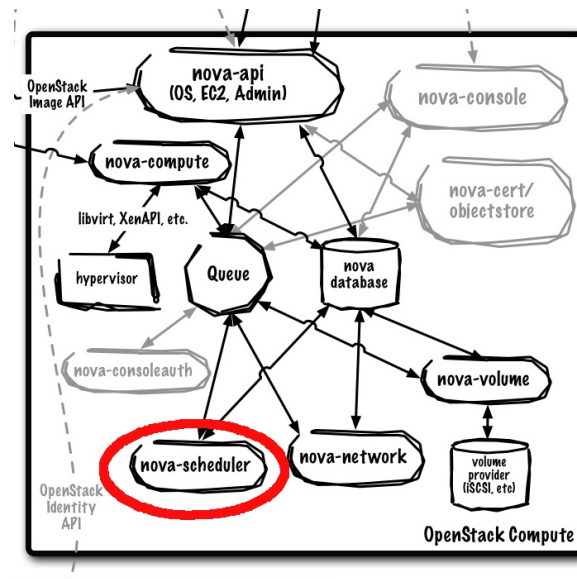


Figure 2.13: OpenStack Compute service

Compute Node segregation

An OpenStack cloud can be logically and physically grouped on different levels [15]:

- **Region:** A Region has its own full OpenStack deployment and can be physically at a different location. Regions share a set of Keystone and Horizon services to provide access control and the graphical management interface.
- **Availability Zone:** Inside of a Region, it is possible to logically group multiple compute nodes into Availability Zones. This zone can be specified when new servers or stacks (via Heat) are instantiated.
- **Host Aggregates:** Compute nodes can also be logically grouped into Host Aggregates by using meta-data to tag them. This feature can be used to separate nodes with certain hardware characteristics (e.g. with SSD drives) from others.

For zoning compute nodes availability zones will be used in the Connectivity Manager in order to achieve the best networking performance between individual servers.

Default VM Placement

Nova makes use of the nova-scheduler service in order to regulate on which host a new VM is launched. The nova-scheduler interacts with other components through the queue and central database repo. All of the compute nodes periodically send an update about their status, available resources and hardware capabilities to the nova-scheduler. It collects this data to make a decision about the placement based on a filter scheduler which takes the following criteria into account:

- are within the requested availability zone
- have sufficient RAM available
- are capable of servicing the request

By default Nova allows its RAM and vCPU resources to be overcommitted, which means that more than the physically available resources can be deployed. However this is to the detriment of their performance.

2.2.3 OpenStack Orchestration (Heat)

Heat provides a template-based orchestration service for creating and managing cloud resources. This means multiple OpenStack resource types (such as virtual machines, floating IP addresses, volumes, security groups and users) can be generated and also maintained with additional functionality like auto-scaling.

2.2.4 OpenStack Neutron

In the early versions of OpenStack, virtual networking was exclusively a sub-component of Nova called Nova-network. This service has its limitations, because it is closely coupled with networking abstractions and no API available. With Neutron the implementation is decoupled from the network abstraction and it provides a flexible management interface to administrators and users.

2.2.4.1 Networking Concepts

Neutron is responsible for defining network connectivity and addressing within OpenStack. In the main network abstraction the following components are defined: A network as a virtual layer 2 segment, a subnet as a layer 3 IP address space used within a network, a port as an interface to a network or subnet, a router that performs address translation and routing between subnets, a DHCP server responsible for IP address distribution, a security group for filtering rules acting as a cloud firewall and Floating IP addresses to give VMs external network access.

Neutron exposes an extensible set of APIs for creating and administering its components. Neutron consists of the following elements [18]:

- **neutron-server:** Provides the logic for SDN, however it does not contain any of its functionality in itself. It provides a generic API for the network operations, is modular and extended with the following agents.
- **L2 agent:** Plugin-specific agent that manages networking on a compute node. For more details, see Modular Layer 2 section.
- **DHCP agent:** Provides DHCP services to tenant networks through dnsmasq instances.
- **L3 agent:** Provides L3/NAT forwarding to allow external network access for VMs (virtual routers).
- **Metadata agent:** Acts as a proxy to the metadata service of Nova.

The different agents can interact with the neutron-server process through a message queue and the OpenStack Networking API. In most use-cases the neutron-server and the different

agents can run on the controller node or on a separate network controller node, with the plugin agent running on each compute node.

2.2.4.2 Modular Layer 2

The ML2 plugin is a framework that allows the simultaneous usage of multiple layer 2 networking technologies. It also helps vendors develop new agents, because an abstract method for their services is already specified.

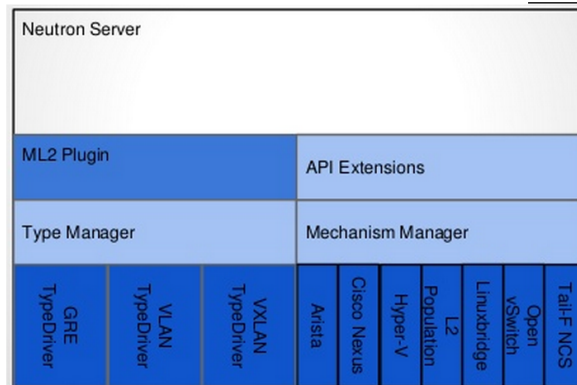


Figure 2.14: Neutron modular framework, including ML2 drivers

The plugin interfaces with the type driver and the mechanism driver. The type driver defines the network types that can be declared when a new network is created and currently includes: local, flat, vlan, gre and vxlan. The mechanism driver enables the mechanism for accessing these network types, i.e. Open vSwitch, Linux Bridge or other vendor-specific solutions.

ML2: Open vSwitch

Open vSwitch is the ML2 mechanism driver that is set as default when installing OpenStack with Neutron enabled using the devstack script. It is also the most commonly deployed agent.

2.2.4.3 Network Distinction

Neutron is connected to different networks. For internet routable connections the external network is used. The management network is created by the network operator and is mapped to pre-existing networks within the datacenter, which is used to connect the different hosts. The tenants within OpenStack have their own isolated self provisioned private networks. Those can optionally be connected to other tenant or external networks. The abstraction of those tenant networks is possible through network namespaces. This allows overlapping IP addresses within the datacenter.

2.2.4.4 Neutron Workflow

From the time when Neutron is started as a process within the OpenStack controller node to booting a virtual machine it passes through the following steps:

1. Start Neutron-Server
2. Start Open vSwitch Agent
3. Start L3-Agent
4. Start DHCP-Agent
5. Start Metadata-Agent
6. Create Networks
7. Create Routers
8. Boot VMs

Once the VM is started it has full network connectivity within the tenants network and if a floating IP address was attached to it, then it can also be reached from the internet.

2.3 Conclusion

In this chapter the technological fundamentals for this thesis were laid out. Starting with the motivation for the use of SDN and the difference to the traditional networking approach to describing the OpenFlow standard which forms the basis of the work of the Open Networking Foundation. The Open vSwitch makes full use of the OpenFlow standards and protocol and forms the basis for the network switching infrastructure used in this thesis. Furthermore it provides the service for enabling network traffic shaping which is described on a deep technical level. All necessary components that are used in connection to this work are introduced.

Chapter 3

Requirements

3.1 Functional Requirements

This section identifies the functional requirements of the Connectivity Manager, specifically as needed for the NUBOMEDIA use-case.

3.1.1 Service-Level Agreement Enforcement

One of the key objectives of the Connectivity Manager is to grant different Service-Level Agreements (SLA) to the links between Virtual Machines. The agreement is set as Quality of Service assurances with the minimal and maximum bandwidth rate set. Network performance problems can provide a negative experience for the end-user, as well as productivity and economic loss. This is why some services need to have an ensured premium traffic.

3.1.2 Optimal Virtual Machine Placement

The placement of Virtual Machines makes a tremendous difference in terms of their network and overall resource performance. It needs to be evaluated which placement makes for the best-available network bandwidth between VMs within the internal network. The current utilization of hosts need to be taken into account as well. A stack should only be deployed if the resources are available at the time of deployment without overcommitting any hardware resources.

3.1.3 Integration with Elastic Media Manager

The Connectivity Manager needs to be integrated with the Elastic Media Manager (EMM) which is used for deploying a topology of resources within a cloud infrastructure. Furthermore it provisions the instances and manages them during their runtime through services such as upscaling the amount of instances after certain utilization alarms are triggered. The CM communicates with the EMM in order to enable the two previously-mentioned requirements for the overall platform.

3.2 Non-functional Requirements

Non-functional requirements generally specify criteria to do with the operation of a system and not with its behavior. Thus the Connectivity Manager should also fit the following characteristics.

3.2.1 Scalability

Today's data-centers can grow in a fast-pace, especially in connection with automated up-scaling of compute resources at a certain level of utilization. This is why the underlying virtualized network software needs to be scalable too.

3.2.2 Modularity

Building modular software not only simplifies further development for a third-party, but also makes it easier to exchange certain parts of the software for improvements or maintenance at a later date. The separation into two different components with a defined API makes it more flexible.

3.2.3 Interoperability

In the case of the use of Open vSwitch, interoperability is given because of the availability for various architectures. The integration into the Linux Kernel and the use of standardized protocols such as OpenFlow are a significant factor.

Chapter 4

State of the Art

4.1 Overview

Three existing solutions for extending Neutron with additional SDN features have been tested. They were selected based on the requirements given in section 3.

4.2 OpenDaylight SDN Controller

OpenDaylight is fully implemented in Java. The Controller platform has multiple Northbound & Southbound interfaces. OpenDaylight exposes a single common OpenStack Service Northbound API which exactly matches the Neutron API. The OpenDaylight OpenStack Neutron Plugin simply passes through and therefore pushes complexity to OpenDaylight and simplifies the OpenStack plugin. The ML2 mechanism driver in Neutron has to be set to the OpenDaylight ML2 plugin with the ODL agent running on the Compute Nodes. The OpenDaylight controller can be run on the Control Node or on a separate VM [19]. The Open vSwitch database (OVSDB) Plugin component for OpenDaylight implements the OVSDB management protocol that allows the southbound configuration of vSwitches. The OpenDaylight controller uses the native OVSDB implementation to manipulate the Open vSwitch database. The component comprises a library and various plugin usages. The OVSDB protocol uses JSON/RPC calls to manipulate a physical or virtual switch that has OVSDB attached to it [21].

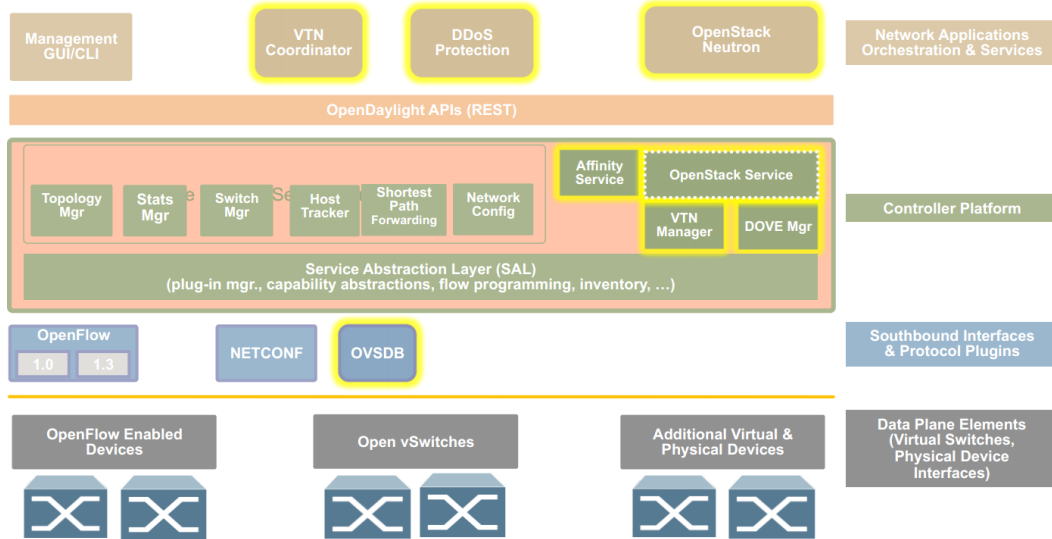


Figure 4.1: Architecture of OpenDaylight Virtualization edition

The OVSDB component is accessible through a Northbound ReST API, which enables the operator to connect to the OpenFlow controller and modify various OVSDB tables. Through this API QoS rules can be deployed. Because it connects directly to the OpenVSwitch tables, all the QoS types that come with OpenVSwitch can be deployed (DSCP marking, setting priority, min-/max-rate for switch ports & OpenFlow Queues). In the local testbed we were able to successfully deploy QoS rules on the ports of Virtual Machines.

4.3 Ryu SDN Controller

Ryu is a component-based software defined networking framework which fully supports OpenFlow 1.0, 1.2, 1.3 and 1.4 switches and is fully written in Python [22]. Ryu is a full featured OpenFlow controller that supports GRE and VLAN tunnelling. The OpenFlow controller that is embedded in the agent sets Flows on the switch by sending OpenFlow messages to the switch [23]. It includes a set of apps which build the base of the SDN controller like L2 switch, ReST interface, topology viewer and tunnel modules. Ryu also includes an app that allows to set QoS rules through a ReST interface which uses a OVSDB interaction library to apply those. The QoS rules can be either applied to a specific Queue within a VLAN or a Switch port. It supports DSCP tagging and setting the min-rate and max-rate of an interface.

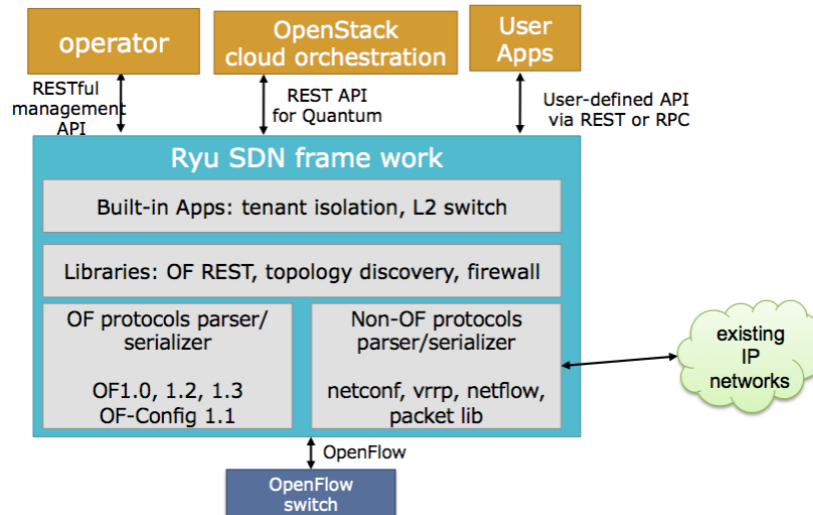


Figure 4.2: Ryu architecture

As of OpenStack IceHouse Ryu has been renamed to OFagent and is included in the Neutron repository. In order to use it as the SDN framework for OpenStack Neutron, OFagent has to be set as both the ML2 mechanism driver (running on the Control / Network node) and the Neutron agent (running on the Compute node).

4.4 OpenStack Neutron - QoS Extension

A Neutron extension has been partially implemented for OpenStack IceHouse which includes an API for setting and retrieving QoS on a per-tenant and per-port basis [24].

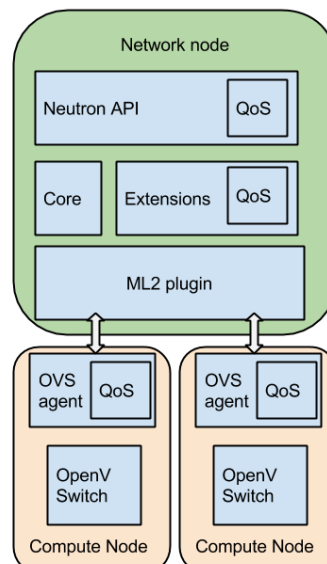


Figure 4.3: Neutron QoS Extension architecture

4.5 Problem Statement

This section lists the restrictions that have been discovered with the previously mentioned solutions.

Problems encountered with OpenDaylight:

The local testbed used for the integration of OpenStack Juno and OpenDaylight Helium consisted of 2 hosts, one running the OpenStack control node and OpenDaylight Controller and a OpenStack Compute Node on the second host. During the tests it was not possible to get the public network access for the Virtual Machines working, thus the L3 routing did not work. This and the fact that it ODL is very complex to debug and understand all underlying processes led us to the decision not to use OpenDaylight.

Problems encountered with Ryu:

The test of Ryu was unsuccessful due to a number of errors while stacking the test environment using Devstack. It was not possible to launch instances and test the QoS features. The lack of proper documentation for the interaction with OpenStack Neutron led us to look more into other SDN controllers for our particular use case. Currently Ryu doesn't support the Distributed Virtual Routing feature that has been introduced with OpenStack Juno.

Problems encountered with Neutron QoS Extension:

The implementation has not been finished and merged into Neutron, however the basic deployment of QoS seem to have been tested successfully.

At the moment it is not clear if the OpenStack community will keep working on this, according to the whiteboard it was deferred to Juno, but it's not included in the current release and no active development is stated in the code / documentation platforms of the OpenStack community.

The patch consists of an extension to the Neutron API which allows setting QoS rules through the Neutron Python client, the actual Neutron extension with the QoS, QoS Driver in the OpenVSwitch agent and an addition to the Neutron Database that includes QoS.

4.6 Conclusion

The listed problems further strengthen the motivation for the implementation of the Connectivity Manager. A feature comparison and analysis is given.

Tool/Solution	QoS support	OF Controller	OpenStack (Juno) integration
Ryu	X	X	No (tests in Juno failed)
OpenDaylight	X	X	No (tests in Juno failed)
Neutron QoS extension	not fully implemented	No	No (IceHouse patch not ported to Juno)

Chapter 5

Design

5.1 Architecture Overview

The Connectivity Manager is logically located between the EMM and the cloud infrastructure and provides the following two functionalities:

- **Optimal Instance Placement:** During the deployment of a stack an algorithm chooses where individual instances are placed within the cloud infrastructure.
- **Service-Level-Agreement enforcement:** Depending on the services that an instance provides to the rest of the stack, certain requirements for its network performance need to be fulfilled.

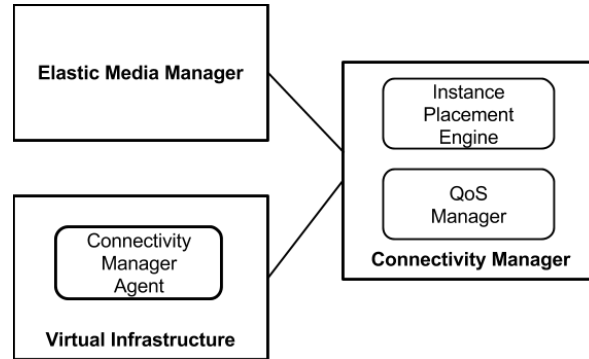


Figure 5.1: High-level architecture of the Connectivity Manager

The *Instance Placement Engine* determines if and where the instances should be deployed. It does so by comparing the current utilization and capacity of the available compute nodes within the availability zone.

The *QoS Manager* enforces different QoS policies based on the type of service that the instance is grouped in. A guaranteed and maximum bit-rate for the network port of an instance can be set. This way a certain network performance can be insured.

5.2 Interface between Manager & Agent

The Connectivity Manager and Agent are two separate applications that communicate using a ReST API.

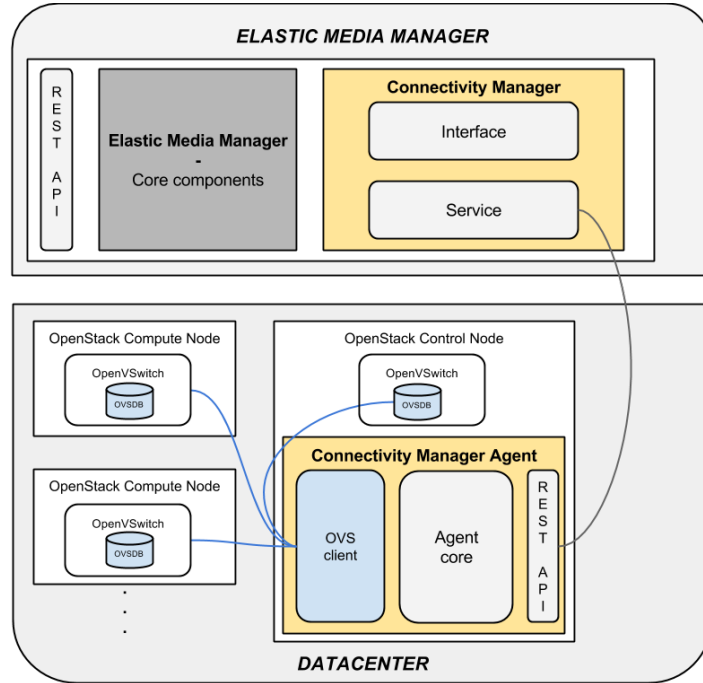


Figure 5.2: Minimized architecture of the Connectivity Manager and its integrations

This design was chosen first of all because the Connectivity Manager is integrated in the EMM, which is required to be placed anywhere outside of the data center. Second of all the Connectivity Manager Agent needs to the OVSDb on the compute nodes and consequently needs to be within the internal management network of the OpenStack infrastructure.

The sequence diagram below displays the work-flow that the CM passes during the run-time.

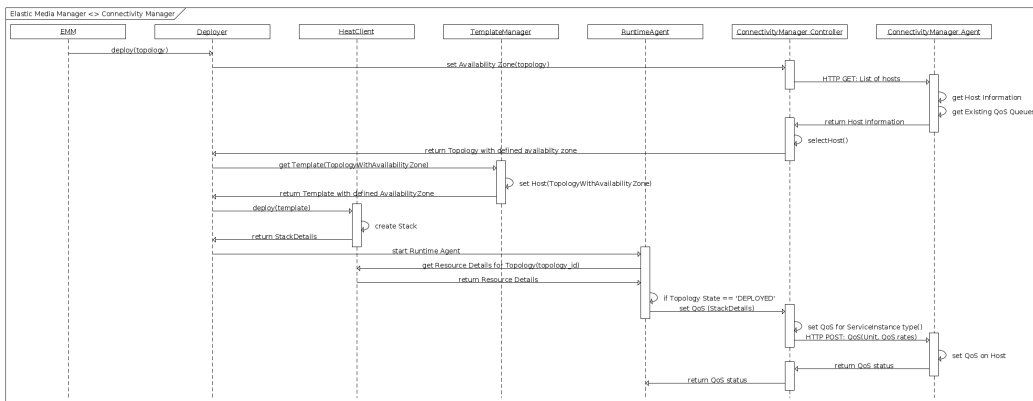


Figure 5.3: Workflow: Deployment of stack & Assignment of QoS policies

As visible in the above figure, the Connectivity Manager receives the Topology object that contains a description of the configuration and specifications of the whole cloud. For the placement decision the CM needs to get the information about the current state of the infrastructure. This exchange with the CM Agent occurs through the given API. Upon reception of that data, the placement algorithm sets the availability zone for each instance within the topology. The topology is then converted into a Heat template by the Template Manager. Once the template got deployed by the Heat Client a runtime agent starts. The purpose of the runtime agent is to continuously check the state of the stack. Once the stack has reached the 'DEPLOYED' state, the runtime agent requests the CM to set the QoS policies according to previously configured values. This configuration is subsequently transmitted to the CM Agent whose task is then to enable it on the according ports of the instances within the Open vSwitch.

5.3 Design of Connectivity Manager

The design of the Connectivity Manager is based on the framework that already exists in the Elastic Media Manager.

It consists of a highly dynamic Factory, abstract Interfaces and the actual implementation as a Service. The Factory Agent reads a configuration file in which the class name of the service for the according interface is defined. It then instantiates a instance from that given class. This means it is easy to replace the actual implementation while the interface to the other components remains identical.

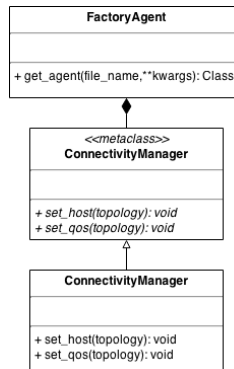


Figure 5.4: Class diagram: Metaclass (Interface) and implementation (Service) instantiated by Factory Agent

The Connectivity Manager interface and service only contain the methods that receive calls from the Elastic Media Manager. However additional helper classes are needed in order to provide the appropriate configuration and the communication with the Connectivity Manager Agent, as shown in the next figure.

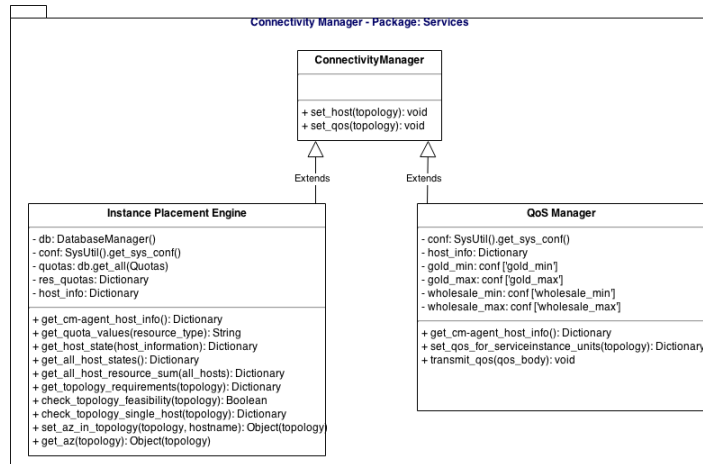


Figure 5.5: Class diagram: Connectivity Manager service and its helper classes

The topology that is created by the EMM contains the following required resources for the stack:

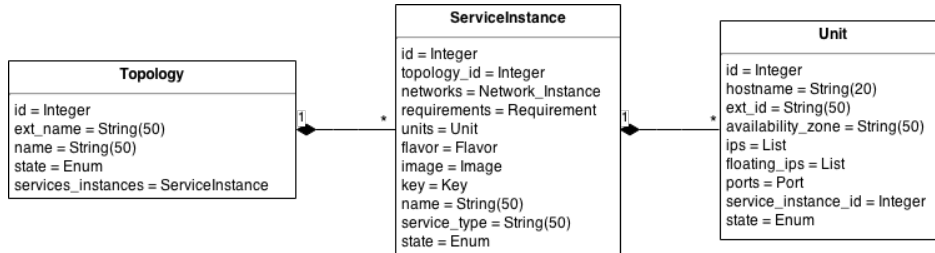


Figure 5.6: Deployment: Topology object from EMM

This object is handed over to the CM during the deployment phase. The topology can contain multiple Service Instances, which defines general parameters such as which networks the underlying instances need to be connected to, the flavor wherein the resources (e.g. amount of RAM & CPU) are specified [25], the image that contains the operating system and additional software packages and the key that is needed to establish a console connection via SSH.

Each Service Instance can contain multiple Units (instances). Among others it contains parameters like IP addresses, the availability zone that it gets deployed on and the external ID which is unique across the whole infrastructure.

5.3.1 Algorithm for Instance Placement

The following section describes the algorithm that is used for setting the availability zone for the instances, and thus placing them in the most advantageous way based on the project requirements.

Each tenant has a set of quotas which limit the resource usage within the OpenStack infrastructure. Limitations can be made based on the following resource types (excerpt): Amount

of instances, vCPU's, RAM, Floating IP addresses and fixed IP addresses [25].

Considering the requirement to have optimal connectivity between individual instances, it is preferable to position them on the same compute node, given that its resources allow this.

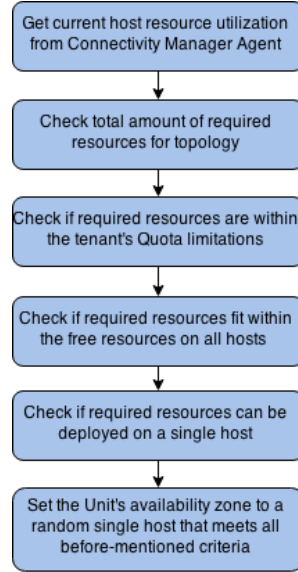


Figure 5.7: Instance Placement activities

5.3.2 QoS Manager

The QoS is deployed on the Units port. This design was chosen, because it needs to be differentiated between different services. The bandwidth rates are required to be changeable and therefore can't be hard-coded. By default the following rates and classes will be set:

Name	Minimum rate	Maximum rate
Wholesale	100 MBit/s	1 GBit/s
Gold	100 MBit/s	10 GBit/s

All Units that are part of the 'Media Server' service instance will be associated with the Gold class and all units of other types in the topology will be part of the Wholesale class.

5.4 Design of Connectivity Manager Agent

The CM Agent is separated into the WSGI Application wherein the routes for the ReST API and the corresponding methods that will be called are defined. The Agent class calls the different Clients to get the resources status and make changes to their configuration.

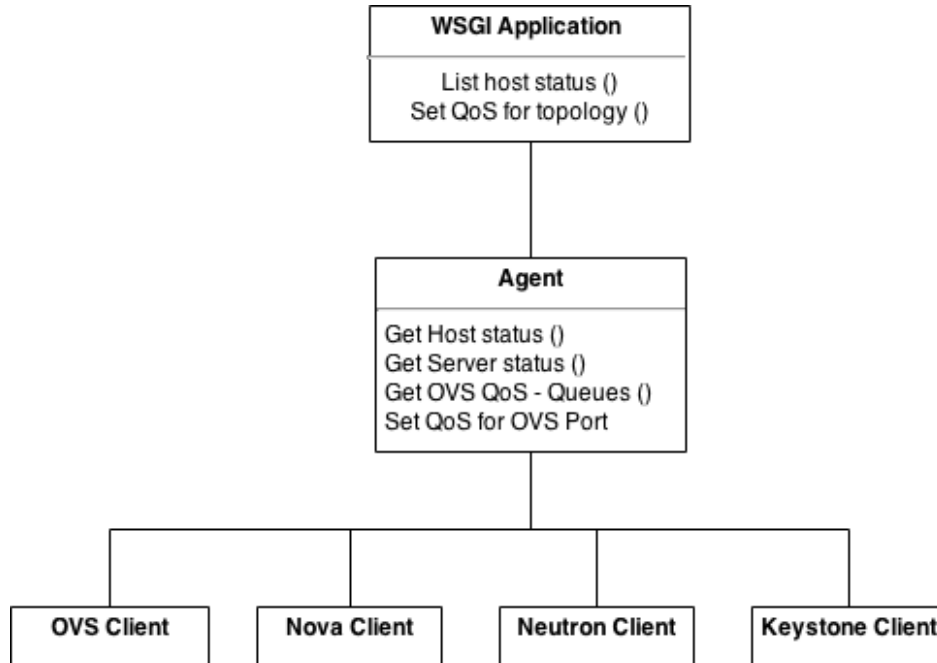


Figure 5.8: Abstract design of the Connectivity Manager Agent

For connecting to the OVSDb on each host the OVS Client is needed. It contains methods for listing all of its ports, interfaces, qos's and queues as well as for applying actions to a port, creating queues and creating new qos's.

The OpenStack Identity API is what the Keystone Client is connecting to in order to get the authentication endpoint and token for Neutron.

The Nova Client binds to the OpenStack Nova API and is required for getting the status of all Compute Nodes and the servers running on it.

Each server can have multiple network ports, which are managed by Neutron. Its Client can retrieve a list of all ports for the servers within one tenant.

Finally the OVS Client makes use of various calls to the OVS Control tool (ovs-vsctl) for reading various tables of the OVSDb and making changes for enabling QoS on a per-port basis.

During the design process the following package structure was chosen:

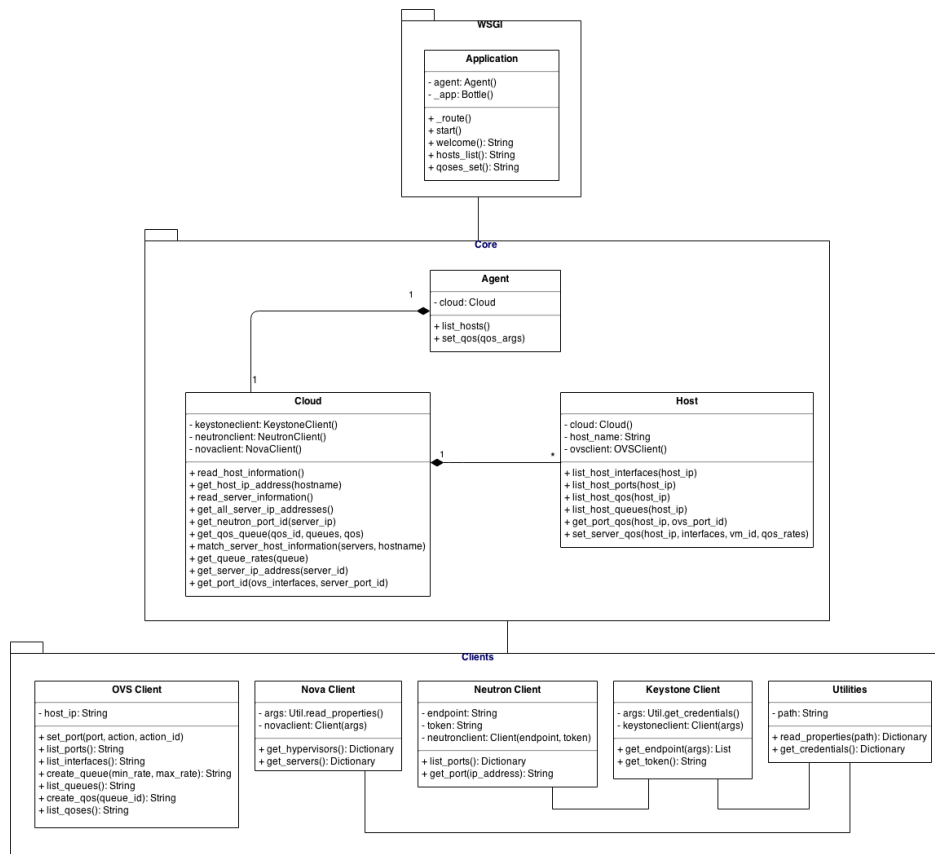


Figure 5.9: Class diagram: Connectivity Manager Agent - All packages

5.4.1 API

In the following table the ReST API that the Connectivity Manager Agent exposes to the Connectivity Manager is described:

HTTP method	Path	Method body	Description
GET	/hosts	-	List all available hosts and their resources
POST	/qoses	JSON: {'Hostname': {'Server ID': {'QoS rates'}}}	Set QoS rates for Servers.

Chapter 6

Implementation

6.1 Environment

The software developed in this thesis is completely realized using the Python programming language. This choice was made because OpenStack offers clients that connect to their API and the Elastic Media Manager is also programmed in Python.

PyCharm was selected as the integrated development environment (IDE) to simplify the programming and testing lifecycles. The code is under revision control using git and the repository that contains both the code for the CM and CM Agent consists of two main branches: master and develop. The develop branch holds the latest changes and upon successful testing those were merged back into master, which is always in a production-ready state.

6.1.0.1 Project Structure

The code is separated into two different projects in order to allow testing of the integration concurrently. The following graphic outlines the focus within the Elastic Media Manager (EMM), whilst the CM Agent is separated and has its own structure.

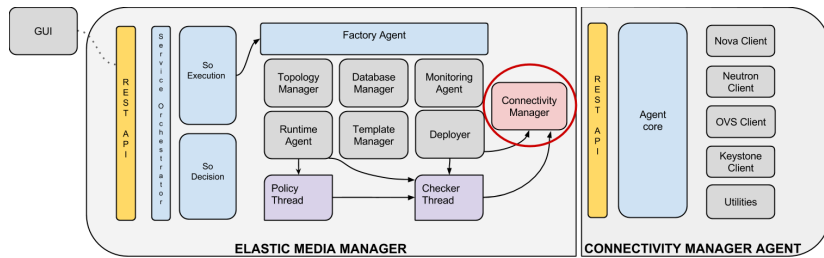


Figure 6.1: Implementation focus

The structure for the CM is dictated by the already existing implementation of the EMM and the scope of this thesis includes solely its extension with a Connectivity Manager interface and service plus the needed changes in the other interfaces to make use of the methods within the CM.

The Connectivity Manager contains the ReST API, the Agent core, clients and other helper classes.

6.1.0.2 Local OpenStack Test Environment

In order to test the Connectivity Manager Agent and the use of its OpenStack API clients a testbed was installed. This test-bed was set up using Vagrant, as it allows to start virtual machines from the command-line and can easily be provisioned and managed. In order to test the software across multiple compute nodes a set up with 2 VM's was installed.

For the installation of OpenStack the devstack script was used, which takes care of not only the deployment of the different components but also their configuration. The configuration parameters are set in the 'local.conf' file. For the OpenStack cluster controller the following configuration was used:

```
[[ local | localrc ]]
ADMIN_PASSWORD=pass
DATABASE_PASSWORD=pass
RABBIT_PASSWORD=pass
SERVICE_PASSWORD=pass
SERVICE_TOKEN=a682f596-76f3-11e3-b3b2-e716f9080d50

HOST_IP=192.168.120.15
OVS_PHYSICAL_BRIDGE=br-ex
MULTIHOST=1

# Enable Logging
LOGFILE=/opt/stack/logs/stack.sh.log
VERBOSE=True
OFFLINE=True
RECLONE=no
LOG_COLOR=True
SCREEN_LOGDIR=/opt/stack/logs

# Neutron
disable_service n-net
enable_service q-svc
enable_service q-agt
enable_service q-dhcp
enable_service q-l3
enable_service q-meta

# OpenStack API paths
MYSQL_HOST=192.168.120.15
RABBIT_HOST=192.168.120.15
GLANCE_HOSTPORT=192.168.120.15:9292
KEYSTONE_AUTH_HOST=192.168.120.15
KEYSTONE_SERVICE_HOST=192.168.120.15
```

```
IMAGE_URLS="$IMAGE_URLS, http://cloud-images.ubuntu.com/releases/trusty/
release/ubuntu-14.04-server-cloudimg-amd64-disk1.img"
```

The configuration file for the second node, which solely runs Nova, the Open vSwitch agent and the Rabbit MQ is identical except for its enabled services and host IP address:

```
HOST_IP=192.168.120.16
ENABLED_SERVICES=n-cpu, rabbit, neutron, q-agt, q-l3
```

During the implementation and testing phase the Connectivity Manager Agent was installed on the controller node while the EMM was executed from the local machine.

6.2 Connectivity Manager Components and Operations

6.2.0.3 Selection of Best-Performing Hypervisor

Step 1: Retrieve current host utilization from CM Agent

In order to decide on the placement the Connectivity Manager first of all needs to retrieve the host information from the Agent. It then sums up the different resource utilizations: amount of servers currently running on it, the total amount of used RAM & vCPUs.

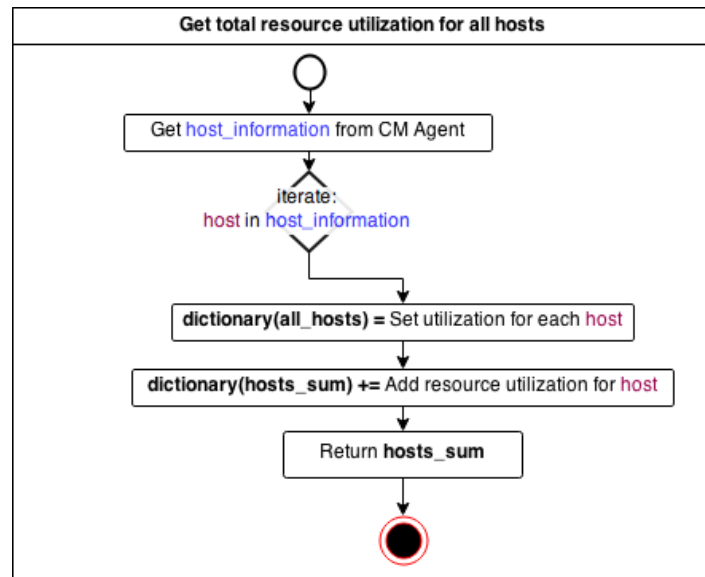


Figure 6.2: Check resource utilization of hosts

The amount of resources that are required in total to deploy the topology on the tenant needs to be calculated by adding up the amount of resources that are needed for each Unit. This can easily be done by checking its flavor.

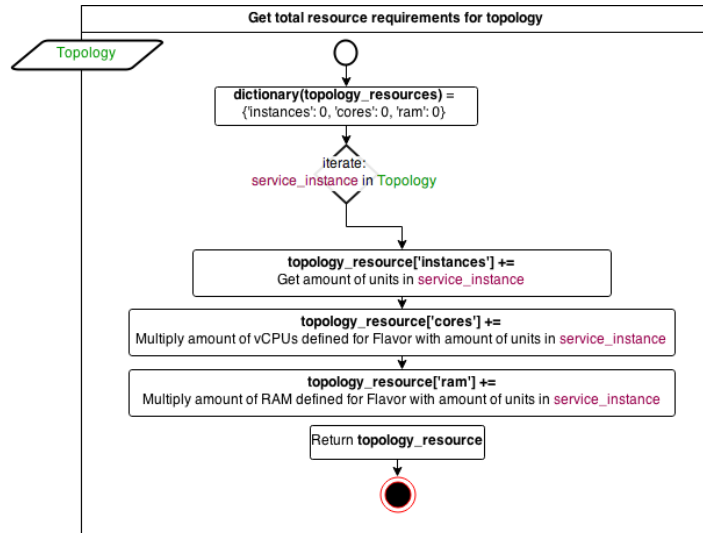


Figure 6.3: Get total amount of required resources for topology

Step 2: Check if Topology is within the limitations of the Quota and currently available resources on the tenants hosts.

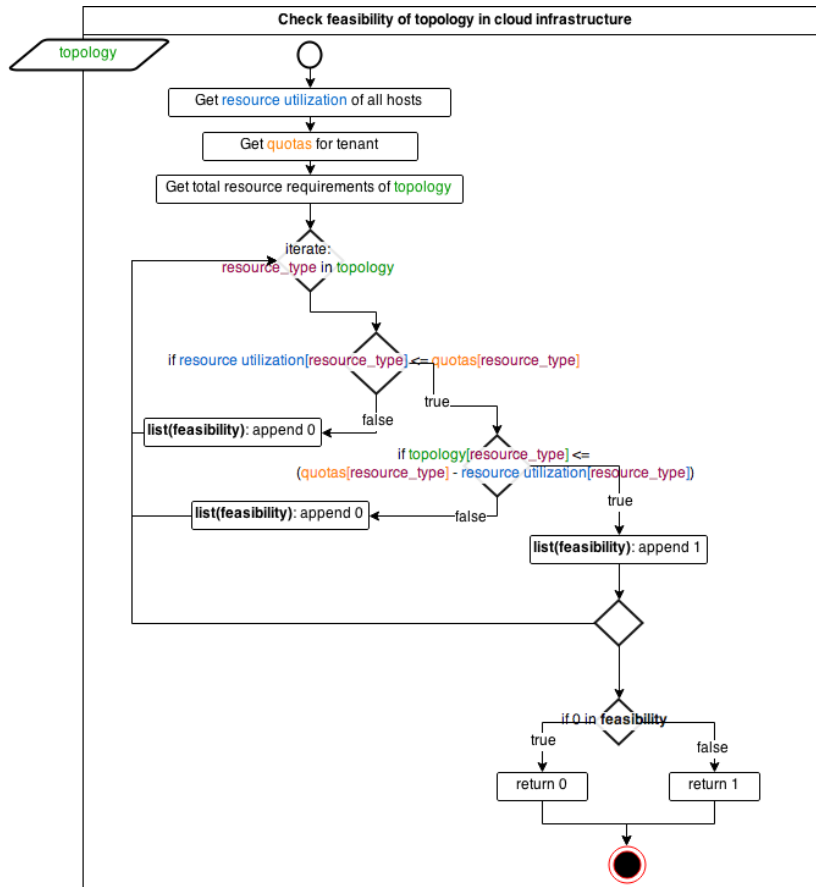


Figure 6.4: Deployment feasibility check

Step 3: Check whether the Topology can be deployed on a single host.

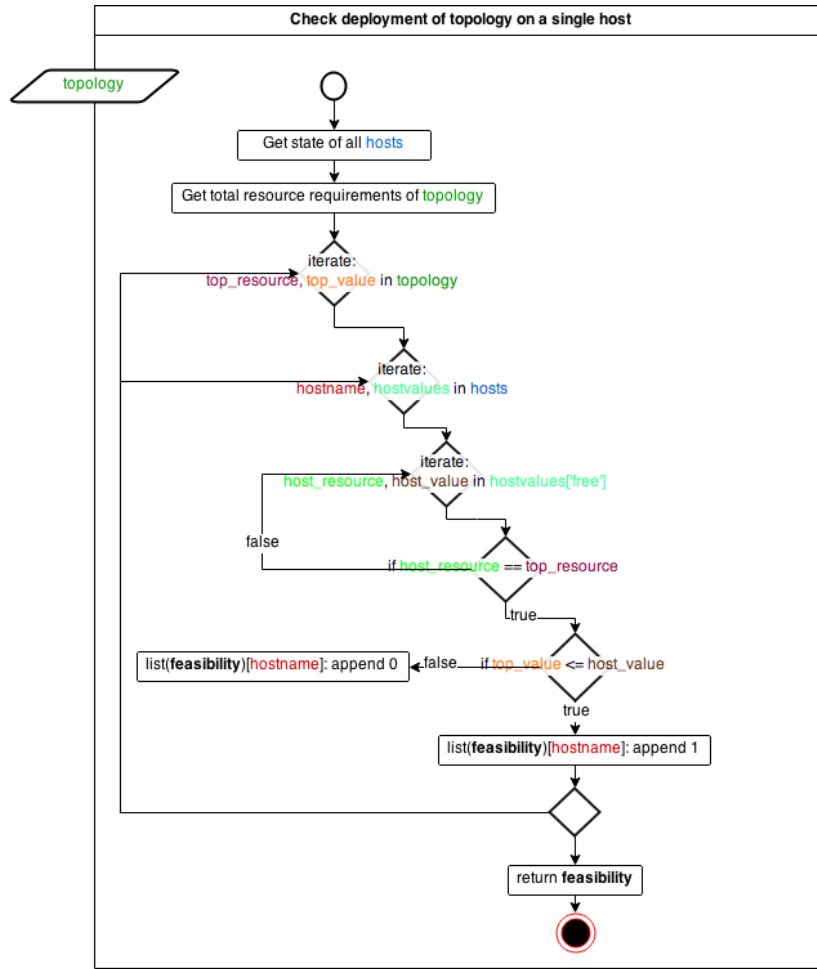


Figure 6.5: Single host deployment check

Step 4: Set selected host as availability zone for each Unit.

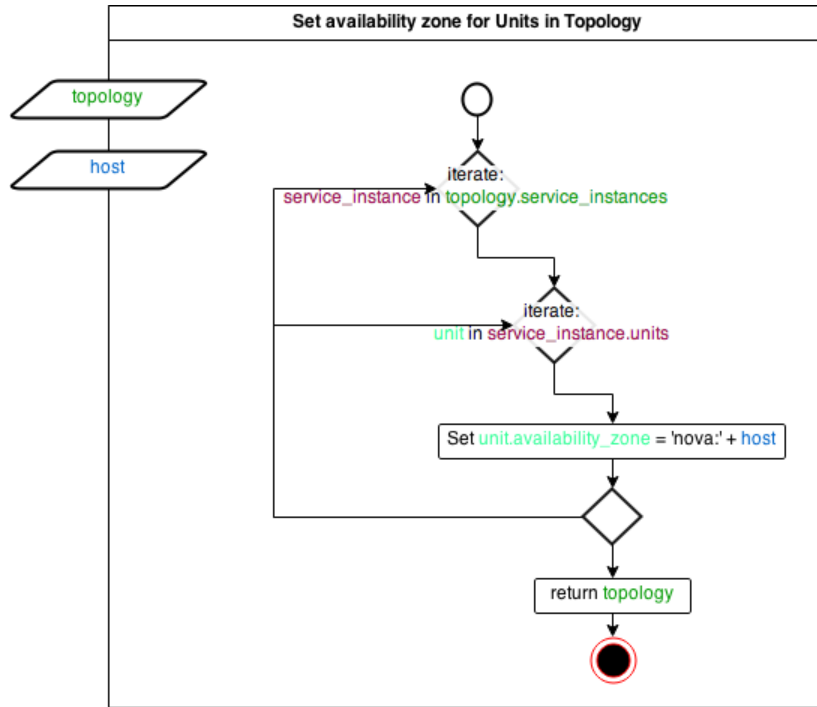


Figure 6.6: Set AZ per Unit

6.2.0.4 Enabling QoS for servers

.. TODO: Mention configuration file!!

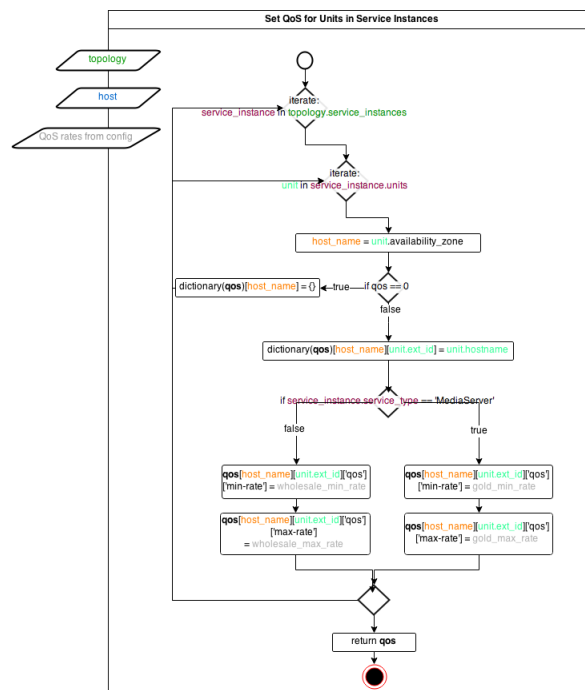


Figure 6.7: Method for setting QoS for all Units

6.3 Connectivity Manager Agent Components and Operations

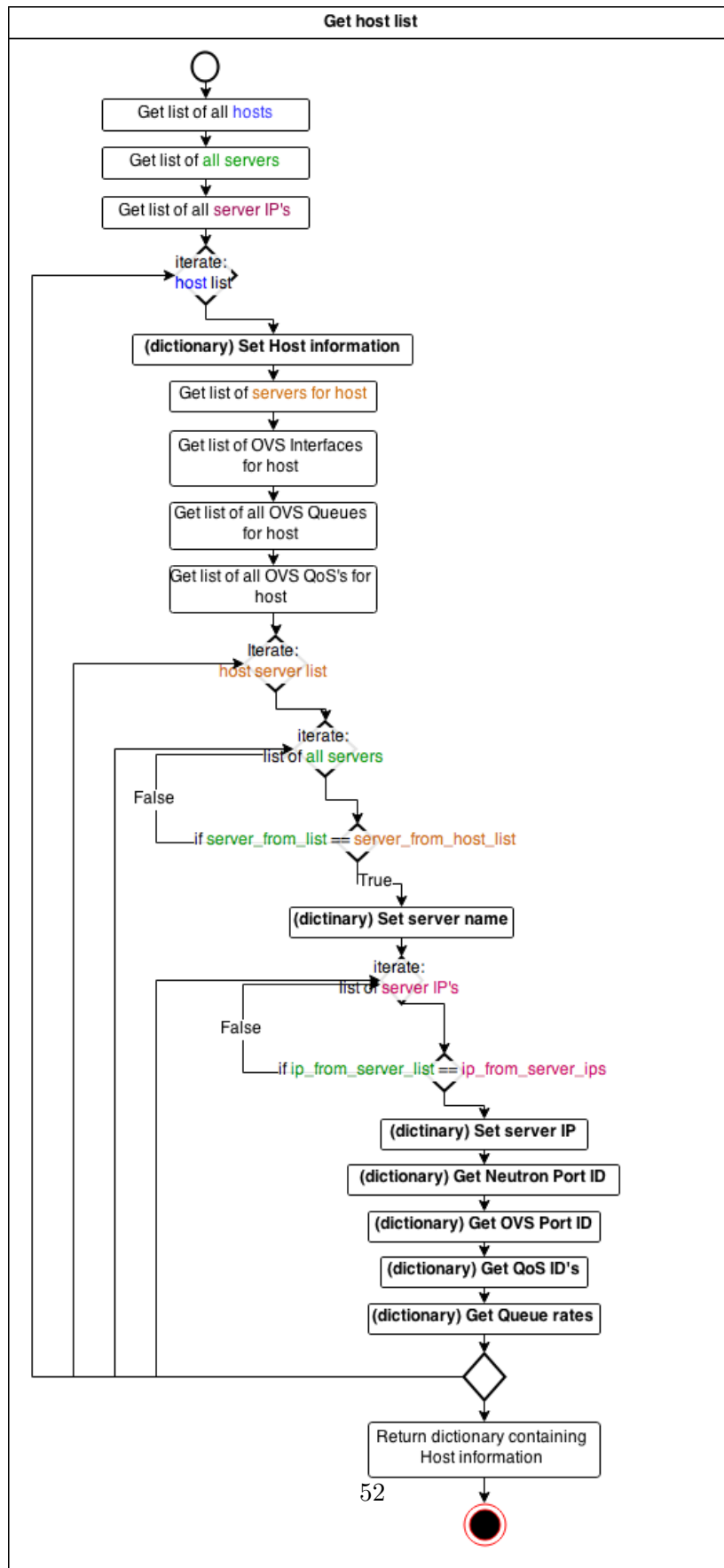


Figure 6.8: Activity diagram: Get list of hosts

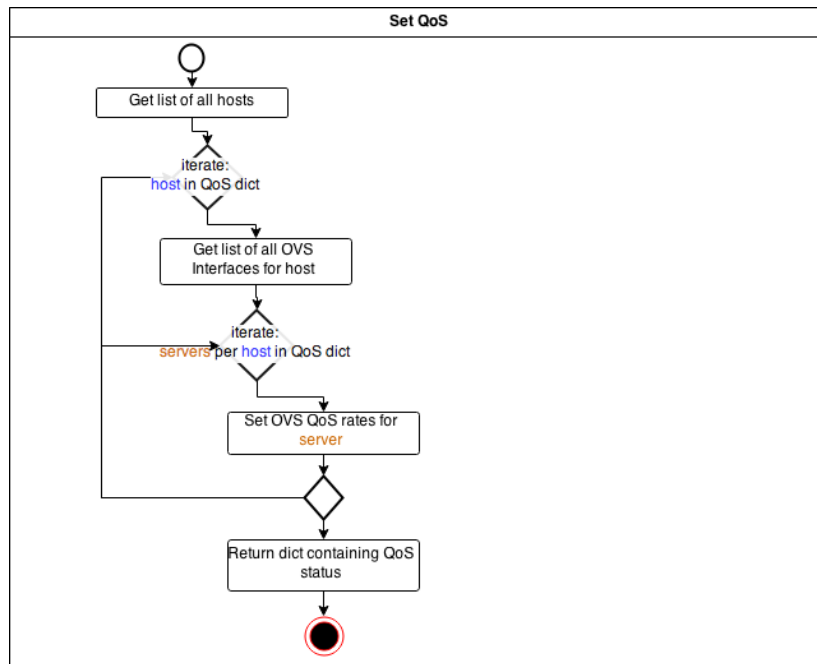


Figure 6.9: Activity diagram: Set QoS rates for all servers

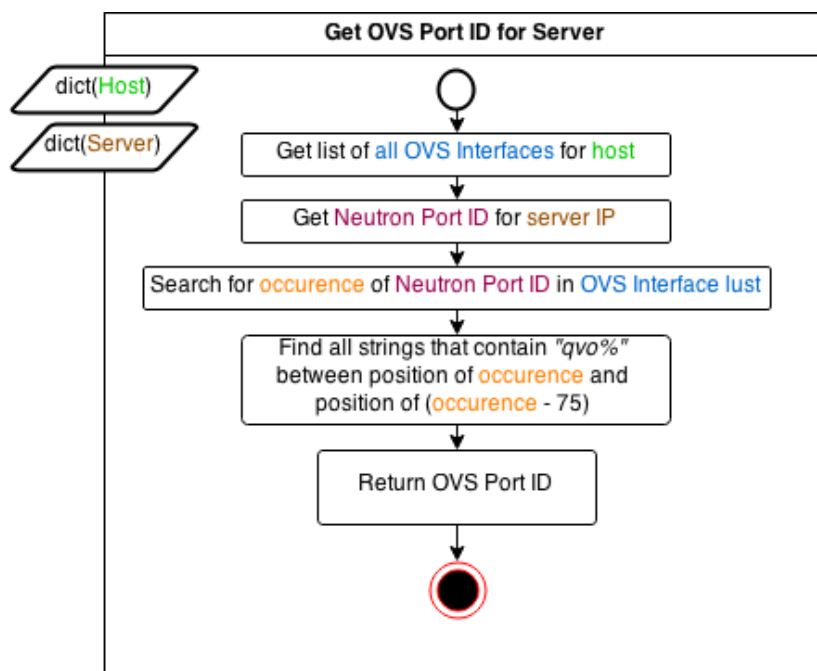


Figure 6.10: Activity diagram: Get OVS Port ID for server

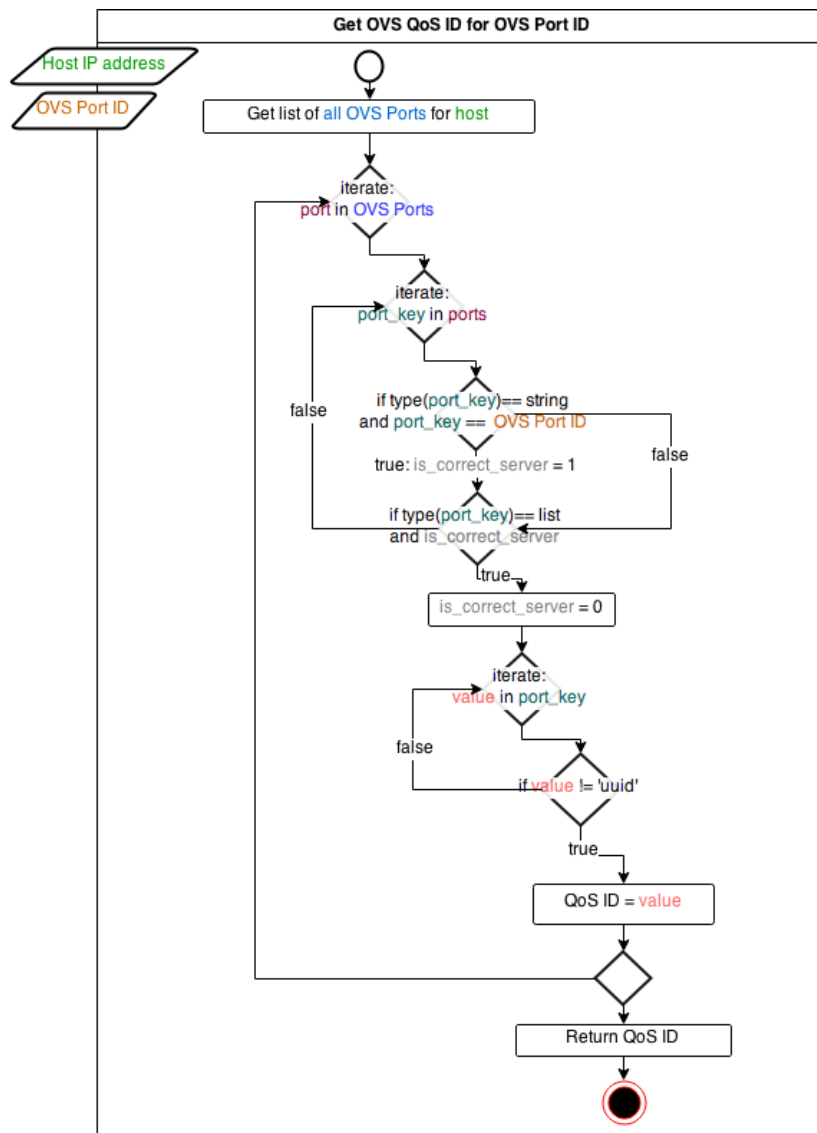


Figure 6.11: Activity diagram: Get QoS ID for OVS Port

Tests

6.4 Conclusion

Chapter 7

Evaluation

7.1 Feature Analysis

7.2 Connectivity Manager Integration with NUBOMEDIA project

The Connectivity Manager (CM) is part of the NUBOMEDIA platform and is placed between the virtual network resource management of the cloud infrastructure and the multimedia application. The main focus of the CM is related to management and control of network functions of the virtual network infrastructure provided by OpenStack.

Nubomedia is an elastic Platform-as-a-Service (PaaS) cloud for interactive social multimedia [26]. Its architecture is based on media pipelines: chains of elements providing media capabilities such as encryption, transcoding, augmented reality or video content analysis. These chains allow building arbitrarily complex media processing for applications. As a unique feature, from the point of view of the pipelines, the NUBOMEDIA cloud infrastructure behaves as a single virtual super-computer encompassing all the available resources of the underlying physical network.

7.3 Network Performance Analysis

The following section contains information about the used configuration and then different scenarios that were used to evaluate the effectiveness of the Connectivity Manager. All scenarios used the same topology and were run on a tenant without any other deployed servers or stacks. The bandwidth tests were performed using the iperf tool.

7.3.1 Test-bed Configuration

The testbed consists of two nodes with the following hardware characteristics:

Name	Controller node	Compute node
Hostname	datacenter-4	dc4-comp
OS	Ubuntu 14.04.1 LTS	Ubuntu 14.04.1 LTS
RAM	12 GB	8 GB
CPU	8-core Intel Core i7-4765T CPU @ 2.00GHz	4-core Intel(R) Core(TM) i3-2120T CPU @ 2.60GHz
Ethernet card	Intel Corporation Gigabit Ethernet Connection I217-LM	Intel Corporation 82579LM Gigabit Network Connection

The two nodes are connected to a Gigabit-Ethernet switch. The installation of OpenStack was performed using the devstack script as outlined in section ... (Implementation / Devstack) .

7.3.2 Installation of Connectivity Manager Agent

A setup script exists in order to make it easier to get the CM Agent running. It builds installs all the necessary Python packages in a virtual environment, in order to have all packages isolated from the already existing Python set-up. This ensures that all packages are in the required version and don't interfere with the ones that are needed OpenStack or other applications.

First of all the git repository needs to be cloned from the remote git server. For the installation the cm-agent.sh script needs to be executed with the 'install' option.

```
stack@datacenter-4:~/nubomedia$ ./cm-agent.sh
Usage: cm-agent.sh option
options:
  install    - install the server
  update     - updates the server
  start      - start the server
  uninstall  - uninstall the server
  clean      - remove build files
```

The installation process includes setting up the virtual environment, installing all required Python packages and copying the configuration file to the /etc/nubomedia folder.

The configuration file needs to be customized, so it contains the IP address of the controller node and the correct OpenStack credentials:

```
stack@datacenter-4:~$ cat /etc/nubomedia/cm-agent.properties
os_username=admin
os_password=pass
os_auth_url=http://192.168.41.45:5000/v2.0
os_tenant=demo
```

Lastly it can be run in a screen session using the following command:

```
$ venv/bin/python cm-agent/wsgi/application.py
```


7.3.3 Topology Definition

The topology that is used for the evaluation contains the following services instances:

data/json_file/topologies/topology-local.json:

```
{
  "name": "local_nm_template_minified",
  "service_instances": [
    {
      "name": "Controller",
      "service_type": "Controller"
    },
    {
      "name": "Broker",
      "service_type": "Broker"
    },
    {
      "name": "MediaServer",
      "service_type": "MediaServer"
    }
  ]
}
```

It can be deployed using a test application which performs a HTTP POST to the EMM API at the */topologies* path.

The service types are further defined in another JSON file, which includes their configuration, networks and other parameters that are needed for provisioning. As one example the Media Server service is given below:

data/json_file/services/MediaService.json

```
{
  "service_type": "MediaServer",
  "version": "1",
  "image": "trusty-iperf",
  "flavor": "ml.mini",
  "key": "nubomedia",
  "configuration": {
  },
  "size": {
    "min": 1,
    "def": 3,
    "max": 5
  },
  "networks": [
    {
      "name": "Network-1",
      "private_net": "8048fd67-70a6-447d-a779-8a86f9eeb35d",
      "private_subnet": "0df3f54c-d1af-4b82-8376-18baa11d0e98",
      "public_net": "62024eab-23c2-4a81-a996-87af4d252282",

```

```

        "security_groups": [
            "SecurityGroup-MediaServer"
        ]
    },
    "requirements": [
        {
            "name": "$BROKER_IP",
            "parameter": "private_ip",
            "source": "Broker",
            "obj_name": "Network-1"
        }
    ]
}

```

7.3.4 Scenario 1: Without Instance Placement Engine & QoS enabled

In this first scenario the deployment without any influence of the Connectivity Manager is shown. Here Nova randomly decides about the placement of the servers. In this case 4 servers were placed on the control node and one server on the compute node.

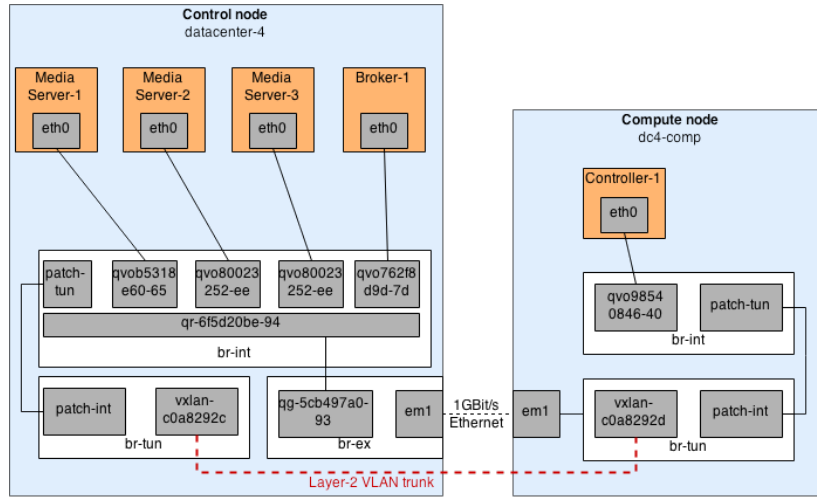


Figure 7.1: Scenario 1: Placement of servers without Connectivity Manager

For testing the bandwidth the MediaServer-2 was used as a TCP server using iperf. All other servers connected to it in client mode sending and retrieving TCP packets in a timeframe of 10 seconds. The following graph shows the bandwidth usage of the servers.

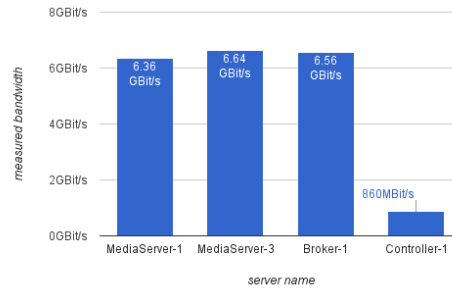


Figure 7.2: Scenario 1: Bandwidth comparison

As visible the network performance of the server on the separate node performs much worse, which is also due to the fact that there is only a Gigabit-Ethernet connection between the nodes.

7.3.5 Scenario 2: With Instance Placement Engine enabled, but without QoS

In this next scenario, the Instance Placement Engine was enabled and therefore the availability zone in the topology was successfully set to a single host.

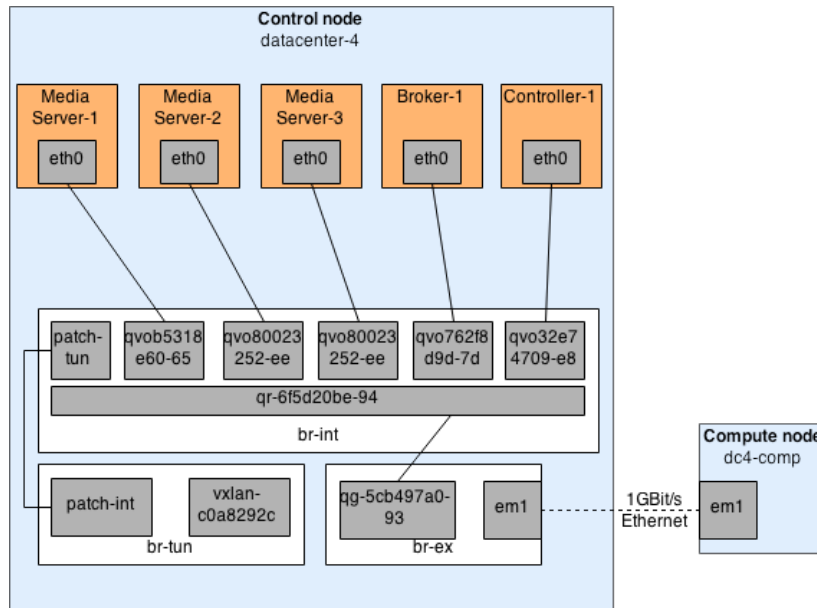


Figure 7.3: Scenario 2: Placement of servers with Connectivity Manager

The following graph shows that the available bandwidth is now evenly distributed for all servers. However the traffic of the MediaServers should be prioritized, which is why QoS is needed to further improve the connectivity according to the requirements.

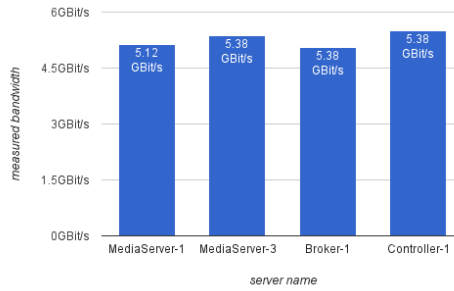


Figure 7.4: Scenario 2: Bandwidth comparison

7.3.6 Scenario 3: Instance Placement Engine and QoS Manager enabled

The servers are again placed on a single host and the Quality of Service configuration that was previously set in the configuration is applied. All media servers have a guaranteed bandwidth rate of 100 MBit/s and a maximum rate of 10 GBit/s, while all servers of other instance types have a rate between 100 MBit/s and 1 GBit/s.

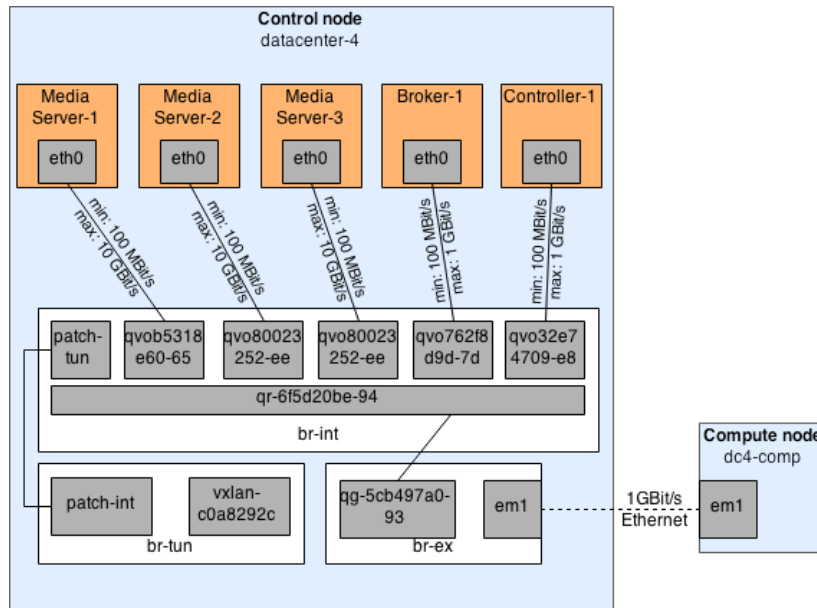


Figure 7.5: Scenario 3: Placement of servers using the Connectivity Manager with QoS enabled

For this last scenario two servers were set as iperf-servers and three servers as iperf-clients. The 'Controller-1' server acted as the server for a connection with the 'MediaServer-2' while the 'MediaServer-3' received packets from 'MediaServer-1' and 'Broker-1'.

The topology with it's configured QoS Queues is also visible by performing a HTTP GET to the Connectivity Manager Agent at the /hosts path:


```

        "uuid":"9b3794a1-7046-4f9d-afb6-0e13a4bbd2d5"
    },
    "name":"MediaServer-3",
    "neutron_port":"58137c93-bc56-4ddd-8330-8af161bc3e3d"
},
"45aecdb2-75ee-4876-af69-4793a7d4c242":{
    "ip":"10.0.0.17",
    "ovs_port_id":"qvo80023252-ee",
    "qos":{
        "queues":{
            "0":{
                "rates":{
                    "max-rate":"10000000000",
                    "min-rate":"100000000"
                },
                "uuid":"e8e6fdf1-62c3-42ca-a354-6437bf7eb9f2"
            }
        },
        "uuid":"241d9954-2048-4b00-bdfc-14c9c5191151"
    },
    "name":"MediaServer-2",
    "neutron_port":"80023252-ee83-4205-9503-7e86a16d7286"
},
"dedd78db-2be4-4655-b08a-8cac40087823":{
    "ip":"10.0.0.16",
    "ovs_port_id":"qvo762f8d9d-7d",
    "qos":{
        "queues":{
            "0":{
                "rates":{
                    "max-rate":"10000000000",
                    "min-rate":"100000000"
                },
                "uuid":"8b0f7fbb-e511-46ee-b20d-eb7b0f0a5c18"
            }
        },
        "uuid":"edcfcbb5-1bfc-4797-a101-845ceb55bc1c"
    },
    "name":"Broker-1",
    "neutron_port":"762f8d9d-7d7f-45fb-836b-eea3ebaad71b"
},
"1eb23d4c-fb5a-42eb-9943-2fb2585b8c04":{
    "ip":"10.0.0.14",
    "ovs_port_id":"qvo32e74709-e8",
    "qos":{
        "queues":{
            "0":{
                "rates":{

```

```

        "max-rate": "1000000000",
        "min-rate": "100000000"
    },
    "uuid": "81d2a604-bc78-41fd-a3ac-71c550251286"
},
{
    "uuid": "db182f5b-afba-4163-9ca5-51501aed1959"
},
{
    "name": "Controller-1",
    "neutron_port": "32e74709-e812-4f14-bbd9-2e9926fa1377"
}
},
{
    "instances": 5,
    "ram_total": 11813,
    "ram_used": 3072,
    "cpu_used": 5,
    "id": 1
}
}

```

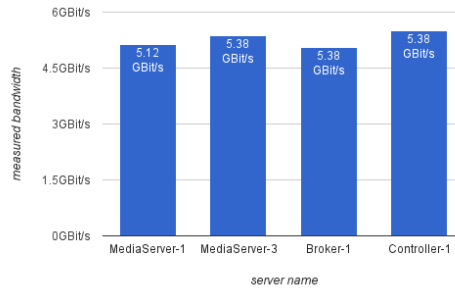


Figure 7.6: Scenario 3: Bandwidth comparison

The first two bars show that the egress port on the Open vSwitch that the *'MediaServer-3'* is connected to is limited to a combined bandwidth of 10 GBit/s. This is why the two connections share the available bandwidth and are nearly equal. It has to be remarked that the network performance is the average of a 10-second bandwidth test.

In the other QoS class the *'MediaServer-2'* which is connecting to *'Controller-1'* uses almost the full bandwidth of the available 1 GBit/s.

7.4 Conclusion

The three test scenarios show an improvement of the network performance, which is given through placing the servers on a single host. Additionally with enabling Quality of Service, the network bandwidth is shaped according to the requirements. It has been shown that the

connection of MediaServers is prioritized to the connection of servers of other service instance types. The Service-Level-Agreement of the Gold and Wholesale classes are fulfilled.

Chapter 8

Conclusion

8.1 Summary

With the use of SDN and OpenFlow a framework for creating programmable networks already exist. However the integration and use of its full extent within OpenStack is not given today. The bandwidth supplied to different services can not be influenced and topologies that are deployed using the Orchestration tool Heat can not be placed on specific hosts. With this work the first step for those main objectives is made. With the help of the Open vSwitch configuration tool it is possible to get information about the current state of the switched network and enable QoS using Queues. The needed configuration and information retrieval can now be accessed through a ReST API. The Elastic Media Manager is able to access this information and make decisions about where the topology has the best-performing network connectivity. The Service-Level-Agreements can be easily changed and extended. The evaluation shows step-by-step the improvements in bandwidth rates with the use of the Connectivity Manager.

8.2 Problems Encountered

One big problem during the development process was getting a stable testbed running with the correct configuration. Devstack was a big help in this process, because not many manual corrections need to be made, other than adjusting the configuration file. However some of the features that are in the 'stable' version of OpenStack Juno still need further bug-fixing and can't be used in a production environment. During the test of the state-of-the-art solution a lot of time was spent on testing different versions, which was unsuccessful. Due to this fact, comparisons to other solutions for enabling Quality of Service can't be made.

8.3 Future Work

For future development the algorithm for selecting the best-performing host could take more dynamic factors into account. One possibility is to check the network-speed of all the host's

network interfaces. It would be beneficial to implement QoS as a Neutron extension and make use of the already existing Neutron API.

Integrate with Neutron API as Extension / Plugin QoS Flow matching

Bibliography

- [1] N. Feamster, J. Rexford, E. Zegura: „The Road to SDN“, ACM Queue, Volume 11, Issue 12, December 30, 2013.
- [2] Open Networking Foundation: „Software-Defined Networking: The New Norm For Networks “, ONF White Paper, April 13, 2012.
- [3] Open Networking Foundation: „Software-Defined Networking (SDN) Definition “, <https://www.opennetworking.org/sdn-resources/sdn-definition>
- [4] S. Natarajan, SDN Hub: „OpenFlow version 1.3 tutorial “, <http://sdnhub.org/tutorials/openflow-1-3/>
- [5] Open Networking Foundation: „OpenFlow Switch Specification “, Version 1.3.0, June 25, 2012.
- [6] Open vSwitch: „Open vSwitch - Frequently Asked Questions “, http://git.openvswitch.org/cgi-bin/gitweb.cgi?p=openvswitch;a=blob_plain;f=FAQ;hb=HEAD
- [7] Open vSwitch: „Open vSwitch - Readme “, <https://github.com/openvswitch/ovs/blob/master/README.md>
- [8] J. Gross, VMware: „Programmable Networking with Open vSwitch “, LinuxCon, September, 2013. <https://events.linuxfoundation.org/sites/events/files/slides/OVS-LinuxCon%202013.pdf>
- [9] J. Pettit, E. Lopez, VMware: „OpenStack: OVS Deep Dive “, 07 November, 2013.
- [10] Open vSwitch: „Open vSwitch Manual“, <http://openvswitch.org/ovs-vswitchd.conf.db.5.pdf>
- [11] B. Hubert: „tc - Linux man page“, <http://linux.die.net/man/8/tc>
- [12] M. Devera, D. Cohen: „HTB Linux queuing discipline manual“, <http://luxik.cdi.cz/~devik/qos/htb/manual/userg.htm>
- [13] J. Vehent: „Journey to the Center of the Linux Kernel: Traffic Control, Shaping and QoS“, http://wiki.linuxwall.info/doku.php/en:ressources:dossiers:networking:traffic_control
- [14] RDO: „Using GRE tenant networks “, 2014, https://openstack.redhat.com/Using_GRE_Tenant_Networks

- [15] OpenStack Foundation: „OpenStack Operations Guide“, February 01, 2015, <http://docs.openstack.org/openstack-ops/openstack-ops-manual.pdf>
- [15] OpenStack Foundation: „OpenStack Configuration Reference“, Compute-Scheduler section, February 01, 2015, http://docs.openstack.org/juno/config-reference/content/section_compute-scheduler.html
- [16] OpenStack Foundation: „OpenStack Compute“, 2015, <https://www.openstack.org/software/openstack-compute/>
- [17] OpenStack Foundation: „OpenStack Installation Guide“, January 29, 2015, <http://docs.openstack.org/juno/install-guide/install/apt/content/>
- [18] OpenStack Foundation: „OpenStack Training Guide“, Chapter 7, January 31, 2015, <http://docs.openstack.org/training-guides/content/operator-network-node.html>
- [19] K. Mestery, D. Meyer, Linux Foundation: „Introduction to OpenDaylight and Hydrogen“, 2014, <https://www.openstack.org/assets/presentation-media/osodlat1.pdf>
- [21] C. Dixon, OpenDaylight: „OVSDB Integration“, October 1, 2014, <https://github.com/opendaylight/docs/blob/master/manuals/developers-guide/src/main/asciidoc/ovsdb.adoc>
- [21] D. Pemberton, A. Linton, S. Russell, University of Oregon: „RYU OpenFlow Controller“, 2014, <https://nsrc.org/workshops/2014/nznog-sdn/raw-attachment/wiki/WikiStart/Ryu.pdf>
- [22] D. Pemberton, A. Linton, S. Russell, University of Oregon: „RYU OpenFlow Controller“, 2014, <https://nsrc.org/workshops/2014/nznog-sdn/raw-attachment/wiki/WikiStart/Ryu.pdf>
- [23] Yamamoto, OpenStack: „Neutron OFAgent - Comparison with OVS“, December 5, 2014, <https://wiki.openstack.org/w/index.php?title=Neutron/OFAgent/ComparisonWithOVS&oldid=69704>
- [24] S. Collins, OpenStack: „Neutron QoS“, October 25, 2013, <https://wiki.openstack.org/w/index.php?title=Neutron/QoS&oldid=34054>
- [25] OpenStack Foundation: „OpenStack Admin User Guide“, February 01, 2015, <http://docs.openstack.org/user-guide-admin/user-guide-admin.pdf>
- [26] L. Lopéz, NUBOMEDIA: „What’s NUBOMEDIA?“, <http://nubomedia.eu/page/whats-nubomedia>

Appendix A

List of acronyms

ACL	→ Access Control List
API	→ Application Programming Interface
AZ	→ Availability Zone
CBQ	→ Class Based Queueing
CM	→ Connectivity Manager
CPU	→ Central Processing Unit
DHCP	→ Dynamic Host Configuration Protocol
DiffServ	→ Differentiated Services
DSCP	→ Differentiated Services Code Point
EMM	→ Elastic Media Manager
GBit/s	→ Gigabit per second
GRE	→ Generic Routing Encapsulation
HTB	→ Hierarchy Token Bucket
HTTP	→ Hypertext Transfer Protocol
IaaS	→ Infrastructure-as-a-Service
IP	→ Internet Protocol
JSON	→ JavaScript Object Notation
KVM	→ Kernel-based Virtual Machine
MBit/s	→ Megabit per second
ML2	→ Modular Layer 2
MPLS	→ Multiprotocol Label Switching
MQ	→ Message Queue

MTU → Maximum Transmission Unit
NAT → Network Address Translation
OF → Open Flow
ONF → Open Networking Foundation
OVS → Open vSwitch
OVSDB → Open vSwitch Database
PaaS → Platform-as-a-Service
qdisc → Queueing Discipline
QoS → Quality of Service
RAM → Random-Access Memory
ReST → Representational State Transfer
RPC → Remote Procedure Call
SDN → Software-Defined Networking
SLA → Service-Level Agreement
SSH → Secure Shell
TC → Traffic Control
TCP → Transmission Control Protocol
TTL → Time to live
vCPU → Virtual Central Processing Unit
VLAN → Virtual Local Area Network
VXLAN → Virtual Extensible Local Area Network
VM → Virtual Machine
WSGI → Web Server Gateway Interface

