# Cloud & IPv6 – schön schaurig

SWITCH

Jens-Christian Fischer
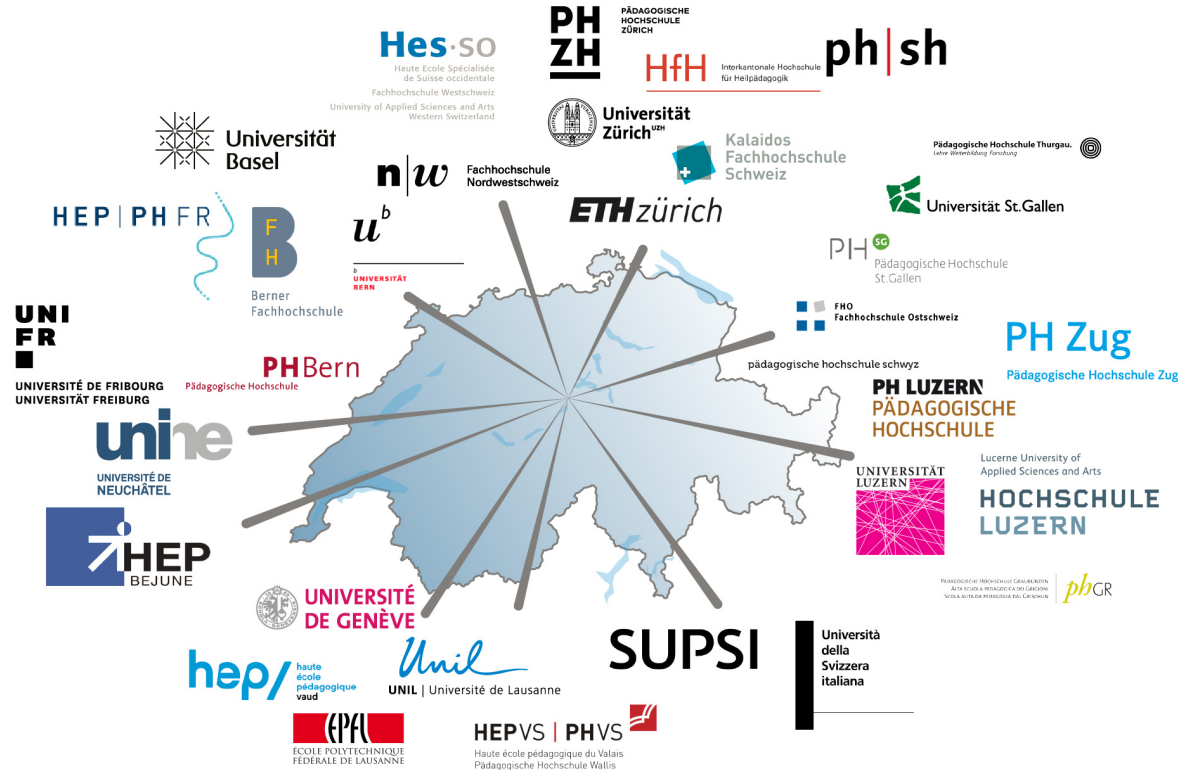jens-christian.fischer@switch.ch
@jcfischer / @switchengines

IPv6 Forum, 2019-07-01

# Academic community Switzerland

# SWITCHengines

Customer tailored computing and storage performance for universities, research and teaching – further developed in the SCALE-UP project mandated by swissuniversities.

## Customers

- Universities
- Research institutions
- eLearning Center
- University hospitals
- Spin-Offs

## Services

- SWITCHengines (IaaS)
- Virtual Private Cloud (VPC)
- SCALE-UP (academic project)

## Your benefits

- Your data in Switzerland
- Integrated network and security
- Support for academic use cases
- Simple administration and billing
- Created together with you

# Disclaimer

RB2011UiAS-RM

POE    GIGABIT ETHERNET

USB

SFP

ETH1    ETH2    ETH3    ETH4    ETH5
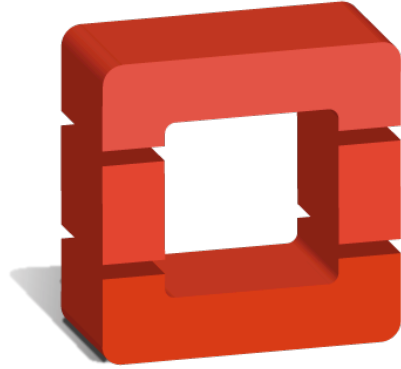
# We are building a cloud

## Kalender 2013

Kalenderpedia
Informationen zum Kalender

| Januar | Februar | März | April | Mai | Juni | Juli | August | September | Oktober | November | Dezember |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 Di Neujahr | 1 Fr | 1 Fr | 1 Mo Oster-montag | 1 Mi Tag der Arbeit | 1 Sa | 1 Mo | 1 Do | 1 So | 1 Di | 1 Fr | 1 So |
| 2 Mi | 2 Sa | 2 Sa | 2 Di | 2 Do | 2 So | 2 Di | 2 Fr | 2 Mo | 2 Mi | 2 Sa | 2 Mo |
| 3 Do | 3 So | 3 So | 3 Mi | 3 Fr | 3 Mo | 3 Mi | 3 Sa | 3 Di | 3 Do Tag der Dt. Einheit | 3 So | 3 Di |
| 4 Fr | 4 Mo | 4 Mo | 4 Do | 4 Sa | 4 Di | 4 Do | 4 So | 4 Mi | 4 Fr | 4 Mo | 4 Mi |
| 5 Sa | 5 Di | 5 Di | 5 Fr | 5 So | 5 Mi | 5 Fr | 5 Mo | 5 Do | 5 Sa | 5 Di | 5 Do |
| 6 So | 6 Mi | 6 Mi | 6 Sa | 6 Mo | 6 Do | 6 Sa | 6 Di | 6 Fr | 6 So | 6 Mi | 6 Fr |
| 7 Mo | 7 Do | 7 Do | 7 So | 7 Di | 7 Fr | 7 So | 7 Mi | 7 Sa | 7 Mo | 7 Do | 7 Sa |
| 8 Di | 8 Fr | 8 Fr | 8 Mo | 8 Mi | 8 Sa | 8 Mo | 8 Do | 8 So | 8 Di | 8 Fr | 8 So |
| 9 Mi | 9 Sa | 9 Sa | 9 Di | 9 Do Himmelfahrt (Vatertag) | 9 So | 9 Di | 9 Fr | 9 Mo | 9 Mi | 9 Sa | 9 Mo |
| 10 Do | 10 So | 10 So | 10 Mi | 10 Fr | 10 Mo | 10 Mi | 10 Sa | 10 Di | 10 Do | 10 So | 10 Di |
| 11 Fr | 11 Mo | 11 Mo | 11 Do | 11 Sa | 11 Di | 11 Do | 11 So | 11 Mi | 11 Fr | 11 Mo | 11 Mi |
| 12 Sa | 12 Di | 12 Di | 12 Fr | 12 So | 12 Mi | 12 Fr | 12 Mo | 12 Do | 12 Sa | 12 Di | 12 Do |
| 13 So | 13 Mi | 13 Mi | 13 Sa | 13 Mo | 13 Do | 13 Sa | 13 Di | 13 Fr | 13 So | 13 Mi | 13 Fr |
| 14 Mo | 14 Do | 14 Do | 14 So | 14 Di | 14 Fr | 14 So | 14 Mi | 14 Sa | 14 Mo | 14 Do | 14 Sa |
| 15 Di | 15 Fr | 15 Fr | 15 Mo | 15 Mi | 15 Sa | 15 Mo | 15 Do | 15 So | 15 Di | 15 Fr | 15 So |
| 16 Mi | 16 Sa | 16 Sa | 16 Di | 16 Do | 16 So | 16 Di | 16 Fr | 16 Mo | 16 Mi | 16 Sa | 16 Mo |
| 17 Do | 17 So | 17 So | 17 Mi | 17 Fr | 17 Mo | 17 Mi | 17 Sa | 17 Di | 17 Do | 17 So | 17 Di |
| 18 Fr | 18 Mo | 18 Mo | 18 Do | 18 Sa | 18 Di | 18 Do | 18 So | 18 Mi | 18 Fr | 18 Mo | 18 Mi |
| 19 Sa | 19 Di | 19 Di | 19 Fr | 19 So | 19 Mi | 19 Fr | 19 Mo | 19 Do | 19 Sa | 19 Di | 19 Do |
| 20 So | 20 Mi | 20 Mi | 20 Sa | 20 Mo Pfingst-montag | 20 Do | 20 Sa | 20 Di | 20 Fr | 20 So | 20 Mi | 20 Fr |
| 21 Mo | 21 Do | 21 Do | 21 So | 21 Di | 21 Fr | 21 So | 21 Mi | 21 Sa | 21 Mo | 21 Do | 21 Sa |
| 22 Di | 22 Fr | 22 Fr | 22 Mo | 22 Mi | 22 Sa | 22 Mo | 22 Do | 22 So | 22 Di | 22 Fr | 22 So |
| 23 Mi | 23 Sa | 23 Sa | 23 Di | 23 Do | 23 So | 23 Di | 23 Fr | 23 Mo | 23 Mi | 23 Sa | 23 Mo |
| 24 Do | 24 So | 24 So | 24 Mi | 24 Fr | 24 Mo | 24 Mi | 24 Sa | 24 Di | 24 Do | 24 So | 24 Di |
| 25 Fr | 25 Mo | 25 Mo | 25 Do | 25 Sa | 25 Di | 25 Do | 25 So | 25 Mi | 25 Fr | 25 Mo | 25 Mi 1. Weih-nachtstag |
| 26 Sa | 26 Di | 26 Di | 26 Fr | 26 So | 26 Mi | 26 Fr | 26 Mo | 26 Do | 26 Sa | 26 Di | 26 Do 2. Weih-nachtstag |
| 27 So | 27 Mi | 27 Mi | 27 Sa | 27 Mo | 27 Do | 27 Sa | 27 Di | 27 Fr | 27 So | 27 Mi | 27 Fr |
| 28 Mo | 28 Do | 28 Do | 28 So | 28 Di | 28 Fr | 28 So | 28 Mi | 28 Sa | 28 Mo | 28 Do | 28 Sa |
| 29 Di | | 29 Fr Karfreitag | 29 Mo | 29 Mi | 29 Sa | 29 Mo | 29 Do | 29 So | 29 Di | 29 Fr | 29 So |
| 30 Mi | | 30 Sa | 30 Di | 30 Do | 30 So | 30 Di | 30 Fr | 30 Mo | 30 Mi | 30 Sa | 30 Mo |
| 31 Do | | 31 So | | 31 Fr | | 31 Mi | 31 Sa | | 31 Do | | 31 Di |

Angaben ohne Gewähr

© www.kalenderpedia.de

# Software

# But wait

# Datacenters in Zurich and Lausanne

**Zurich University of the Arts, Toni Areal**

**Université de Lausanne, Géopolis**

Photo Fabrice Ducrest © UNIL

# Now

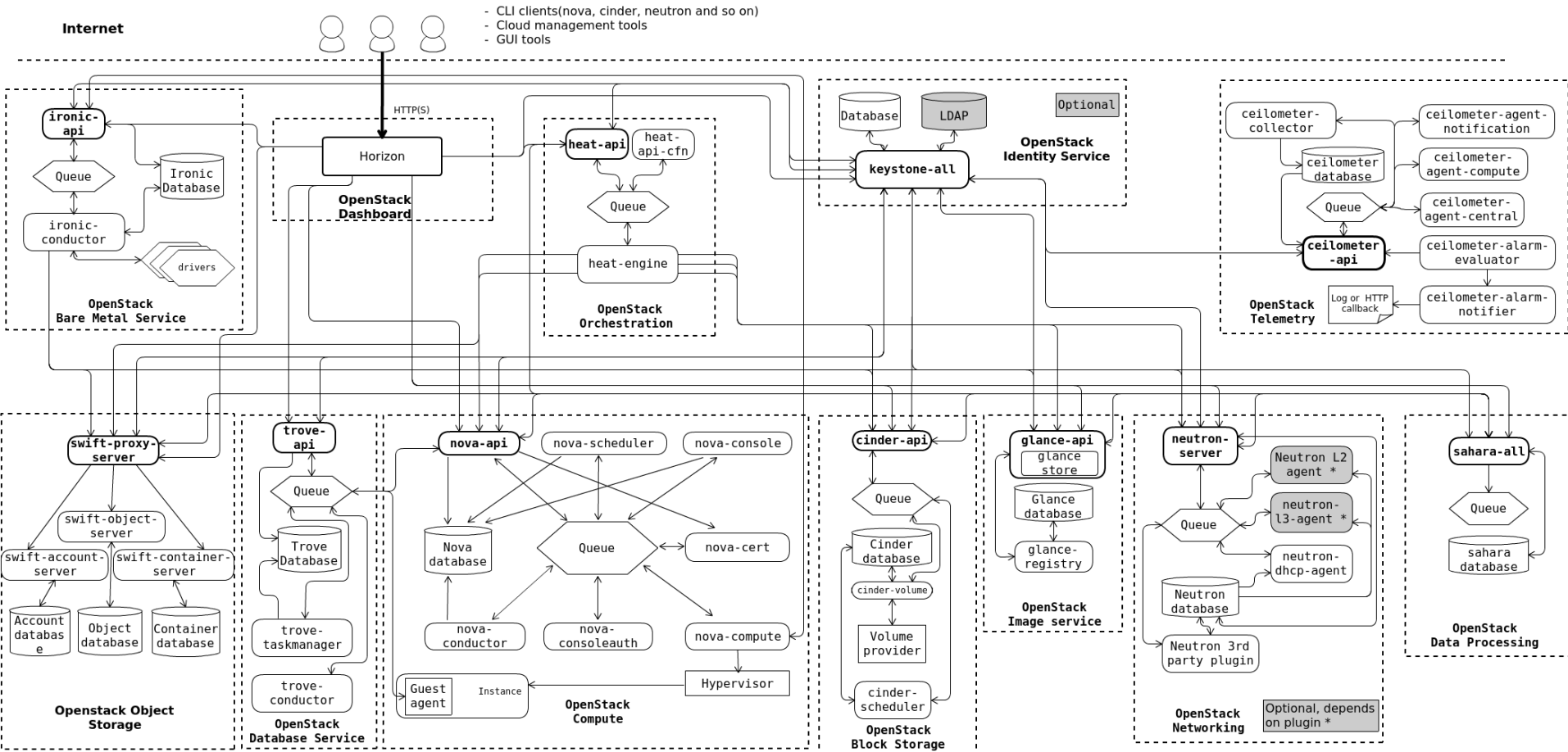| | |
|---|---|
| Physical servers: | ~ **400** |
| CPU cores: | **4076** (physical cores) |
| Memory: | ~ **32** TB |
| Storage: | ~ **10** PB (Ceph SATA) / ~ 1400 Disks |
| | ~ **100** TB (Ceph SSD) / 50 NVMe |
| GPU: |  **8** Titan XP |
| | **16** T4 |
| | **34** P100 |
| Network: | **26** Cumulus Linux 40 & 100Gbs switches |
| | Dual 10 Gbs; upgrading to 100 Gbs (Q2 2019) |
| | L2 tunnel to campus networks (VPC) |

SWITCH

Internet

- CLI clients(nova, cinder, neutron and so on)
- Cloud management tools
- GUI tools

HTTP(S)

**ironic-api**

Queue

Ironic Database

ironic-conductor

drivers

**OpenStack Bare Metal Service**

Horizon

**OpenStack Dashboard**

**heat-api**    heat-api-cfn

Queue

heat-engine

**OpenStack Orchestration**

Database    LDAP    Optional

**OpenStack Identity Service**

**keystone-all**

ceilometer-collector    ceilometer-agent-notification

ceilometer database    ceilometer-agent-compute

Queue    ceilometer-agent-central

**ceilometer-api**    ceilometer-alarm-evaluator

Log or HTTP callback    ceilometer-alarm-notifier

**OpenStack Telemetry**

**swift-proxy-server**

swift-object-server

swift-account-server    swift-container-server

Account database    Object database    Container database

**Openstack Object Storage**

**trove-api**

Queue

Trove Database

trove-taskmanager

trove-conductor

**OpenStack Database Service**

**nova-api**    nova-scheduler    nova-console

Nova database    Queue    nova-cert

nova-conductor    nova-consoleauth    nova-compute

Guest agent    Instance    Hypervisor

**OpenStack Compute**

**cinder-api**

Queue

Cinder database

cinder-volume

Volume provider

cinder-scheduler

**OpenStack Block Storage**

**glance-api**

glance store

Glance database

glance-registry

**OpenStack Image service**

**neutron-server**

Neutron L2 agent *

Queue    neutron-l3-agent *

Neutron database    neutron-dhcp-agent

Neutron 3rd party plugin

**OpenStack Networking**

Optional, depends on plugin *

**sahara-all**

Queue

sahara database

**OpenStack Data Processing**

© 2019 SWITCH | 12

## added logic to make the creation of networks (IPv4 only) validation a...

… bit smarter:

 – detects if the cidr is already in use
 – detects if any existing smaller networks are within the range of requested cidr(s)
 – detects if splitting a supernet into # of num_networks && network_size will fit
 – detects if requested cidr(s) are within range of already existing supernet (larger cidr).

IPv6 logic remains intact yet had not been improved by this code.

master    pike-em    ...    12.0.0a0

**John Tran** authored and **Tarmac** committed on Aug 14, 2011

Showing **2 changed files** with **237 additions** and **5 deletions**.

Overview    Code    **Bugs**    Blueprints    Translations    Answers

# [RFE] Support metadata service with IPv6-only tenant network

Bug #1460177 reported by 👤 Baodong (Robert) Li on 2015-05-29

This bug affects 13 people                                                    🔥 76

| Affects | Status | Importance | Assigned to | Milestone |
|---|---|---|---|---|
| ▷        🔳 neutron | Triaged 📝 | Wishlist | Unassigned | |

➕ Also affects project ❓    ➕ Also affects distribution/package    🔄 Nominate for series

## Bug Description

EC2 metatdata service is supported by nova metadata service that is running
in the management network. Cloud-init running in the instance normally
accesses the service at 169.254.169.254. Cloud-init can be configured with
metadata_urls other than the default http://169.254.169.254 to access the
service. But such configuration is not currently supported by openstack. In
order for the instance to access the nova metadata service, neutron
provides proxy service that terminates http://169.254.169.254 and forwards
the request to the nova metadata service, and responds back to the
instance. Apparently, this works only when IPv4 is available in the tenant
network. For an IPv6-only tenant work, to continue the support of this
service, the instance has to access it at an IPv6 address. This requires
enhancement in Neutron to support it.

A few options have been discussed so far:
    -- define a well-known ipv6 link-local address to access the metadata
service.
    -- enhance IPv6 RA to advertise the metadata service endpoint to
instances. This would require standards work and enhance cloud-init to
support it.
    -- define a well-known name for the metadata service and configure
metadata_urls to use the name. The name will be resolved to a datacenter
specific IP address. The corresponding DNS record should be pre-provisioned
in the datacenter DNS server for the instance to resolve the name.

👤 **Brian Haley (brian-haley)** wrote on 2018-06-01:

And I have another question. Say neutron and cloud-init can be upgraded to
support an IPv6-only metadata request. Are there additional changes
required to the API? For example, there is a local-ipv4 value today, is a
local-ipv6 needed? I'm confused by the wording of the 'network/interfaces/
macs/mac/ipv6s' field - "The IPv6 addresses associated with the interface.
Returned only for instances launched into a VPC." - local-ipv4 doesn't
mention VPC.

👤 **YAMAMOTO Takashi (yamamoto)** wrote on 2018-06-22:

@Brian

just because non VPC (EC2-Classic?) doesn't support ipv6?
https://docs.aws.amazon.com/AWSEC2/latest/UserGuide/using-vpc.html

👤 **Miguel Lavalle (minsel)** on 2018-11-07:

**tags**:added: rfe-postponed
      removed: rfe-triaged

To post a comment you must log in.

# Neutron Networking

- Took us 2 years to get a default IPv6 address to a newly started VM

- Still some manual work required for IPv6 routed internal private networks

# VMs with global routed IPv6

Displaying 6 items

| | Instance Name | Image Name | IP Address | Flavor |
|---|---|---|---|---|
| ☐ | t4e | Ubuntu Xenial with GPU support (SWITCHengines) | 10.0.2.86<br>2001:620:5ca1:1f0:f816:3█████████<br><br>**Floating IPs:**<br>86.119.38█████ | g1.c16r176-4t4 |
| ☐ | t4b | Ubuntu Xenial with GPU support (SWITCHengines) | 10.0.2.41<br>2001:620:5ca1:1f0:f816█████████<br><br>**Floating IPs:**<br>86.119█████ | g1.c16r176-4t4 |
| ☐ | inv_rescue | Ubuntu Bionic 18.04 (SWITCHengines) | 10.0.0.99<br>2001:620:5ca1:1f0:f816:3█████████ | m1.small |
| ☐ | engines-admin -runner | - | 10.0.3.153<br>2001:620:5ca1:1f0:f816:3eff█████████<br><br>**Floating IPs:**<br>86.119.3█████ | m1.medium |

# Complex SDN Setups

2001:620:5ca1:1f0::/64, 10.0.0.0/16

10.0.23.0/24, 2001:620:5ca1:1ee::/64

2001:620:5ca1:131::/64, 10.1.49.0/24

SWITCH
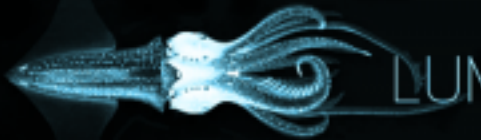
BOBTAIL

DUMPLING

FIREFLY

GIANT

HAMMER

JEWEL

ceph
LUMINOUS

NAUTILUS

# Ceph – Software Defined Storage

- IPv6 from the beginning

- No *) problems

*) almost

# Brocade L2 Switches

- Clear IPv6 neighbour caches (manually)
- Otherwise machines would loose connectivity

Which is bad in a storage cluster

# Hard problems

- Random (huge) performance problems on random VMs (with big amount of IOPS)

Nothing seems to work

Attract "Management Attention"

Blame the running Bluestore migrations

Ask the people at CERN

Spend 2 Weeks trying to reproduce it

Start sniffing the network

# Reproducibility

We found that the TCP traffic from the writing VM (the sender) to an OSD (the receiver) was limited to **one** 512-byte TCP segment **every 200 ms**

**Install new Kernel**

# KVM / libvirt

- Talks IPv6 to Ceph Storage Cluster

# Kubernetes Warms Up to IPv6

25 Feb 2019 11:55am, by Mary Branscombe



There's a finite number of public IPv4 addresses and the IPv6 address space was specified to solve this problem some 20 years ago, long before Kubernetes was conceived of. But because it was originally developed inside Google and it's only relatively recently that cloud services like Google and AWS have started to support IPv6 at all, Kubernetes started out with only IPv4 support.

# Networking in K8s

- Pods support IPv4 & IPv6 – it just works
- Internal K8s Services only work with IPv4 (even though it is claimed that IPv6 is supported)

=> Run everything in IPv4

# The day we shut IPv4 down

- The day we shut of IPv4 outbound connectivity in one of our clusters

- And everything continued to work

- For hours and hours

- Until Kubernetes evicted a number of pods (Because k8s can be stupid)

- And tried to rebuild the images

- And we discovered:

```
                              3. fischer@macjcf: ~ (zsh)
Last login: Mon Jul  1 09:25:14 on ttys012
~  dig aaaa dockerhub.io

; <<>> DiG 9.10.6 <<>> aaaa dockerhub.io
;; global options: +cmd
;; Got answer:
;; ->>HEADER<<- opcode: QUERY, status: NOERROR, id: 43504
;; flags: qr rd ra; QUERY: 1, ANSWER: 0, AUTHORITY: 1, ADDITIONAL: 1

;; OPT PSEUDOSECTION:
; EDNS: version: 0, flags:; udp: 4096
;; QUESTION SECTION:
;dockerhub.io.                    IN      AAAA

;; AUTHORITY SECTION:
dockerhub.io.           10744   IN      SOA     a.dns.gandi.net. hostmaster.gandi.net. 1484668036 108
00 3600 604800 10800

;; Query time: 82 msec
;; SERVER: 130.59.31.248#53(130.59.31.248)
;; WHEN: Mon Jul 01 15:51:21 CEST 2019
;; MSG SIZE  rcvd: 112

~  
```
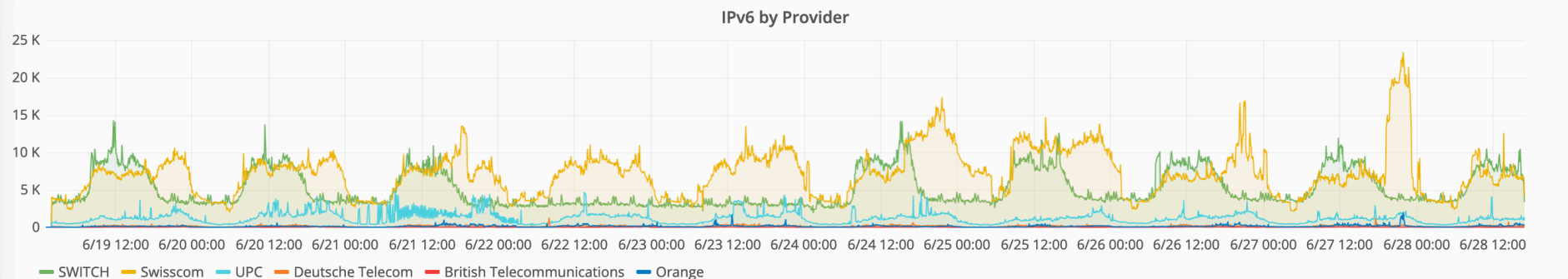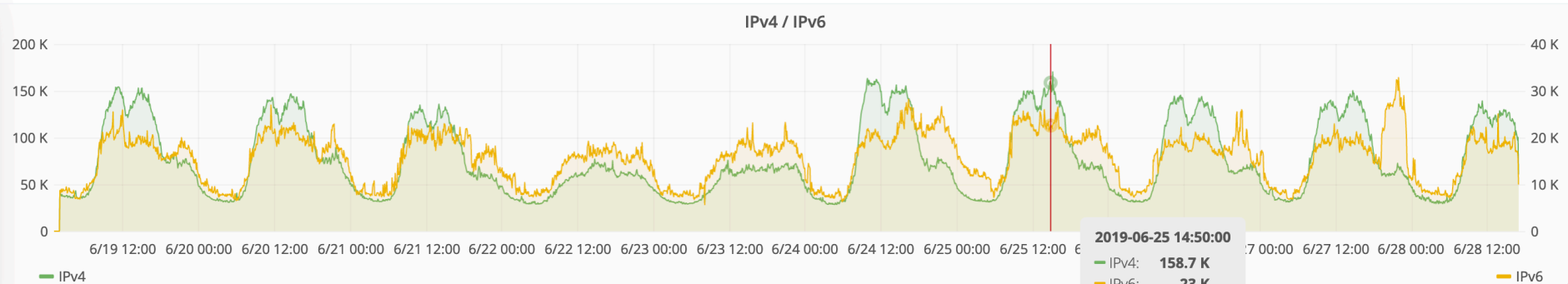
# What about the users

# Enduser IPv6