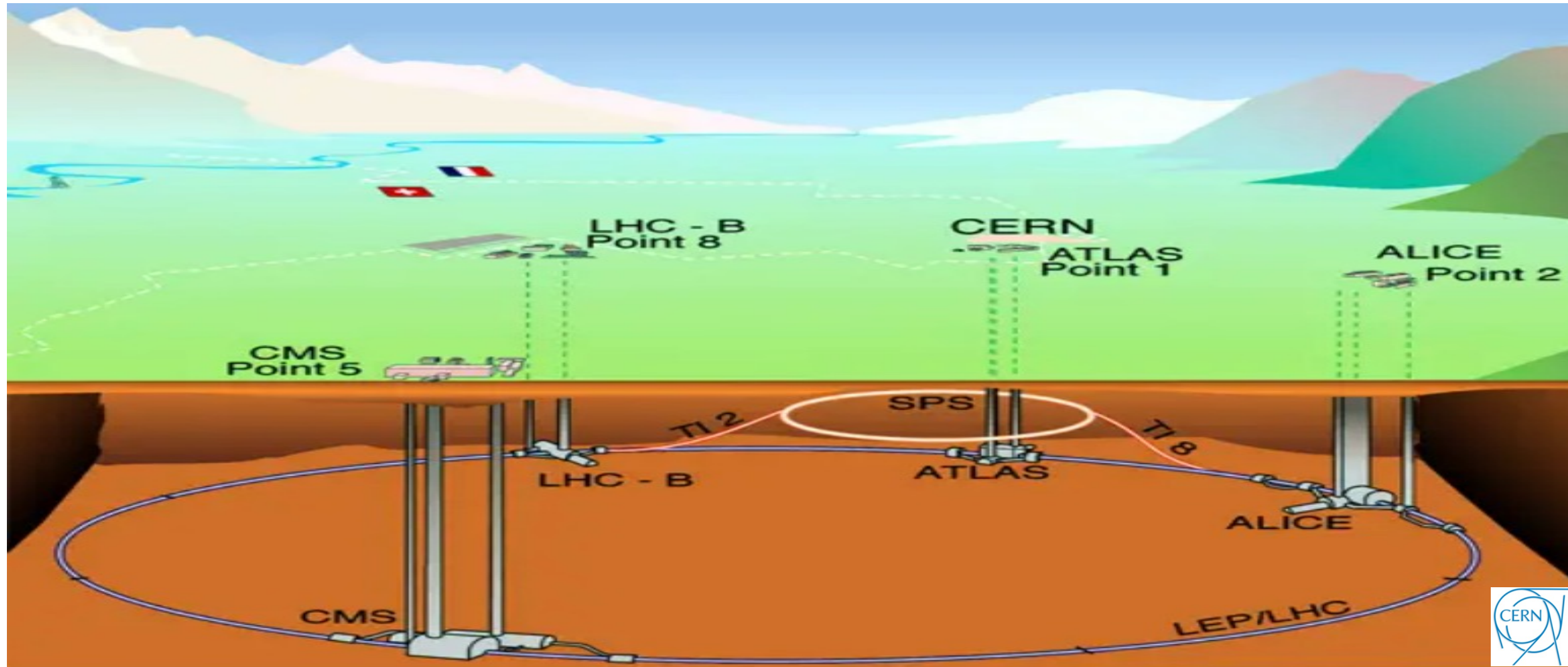# Moving towards IPv6-only in the German Tier-1 Data Center of the CERN Large Hadron Collider

Bruno Hoeft, Matthias Schnepf, Max Fischer, Andreas Petzold
**Karlsruhe Institute of Technology, Hermann-von-Helmholtz-Platz 1, 76344 Eggenstein-Leopoldshafen, {first-.familyname}@kit.edu**
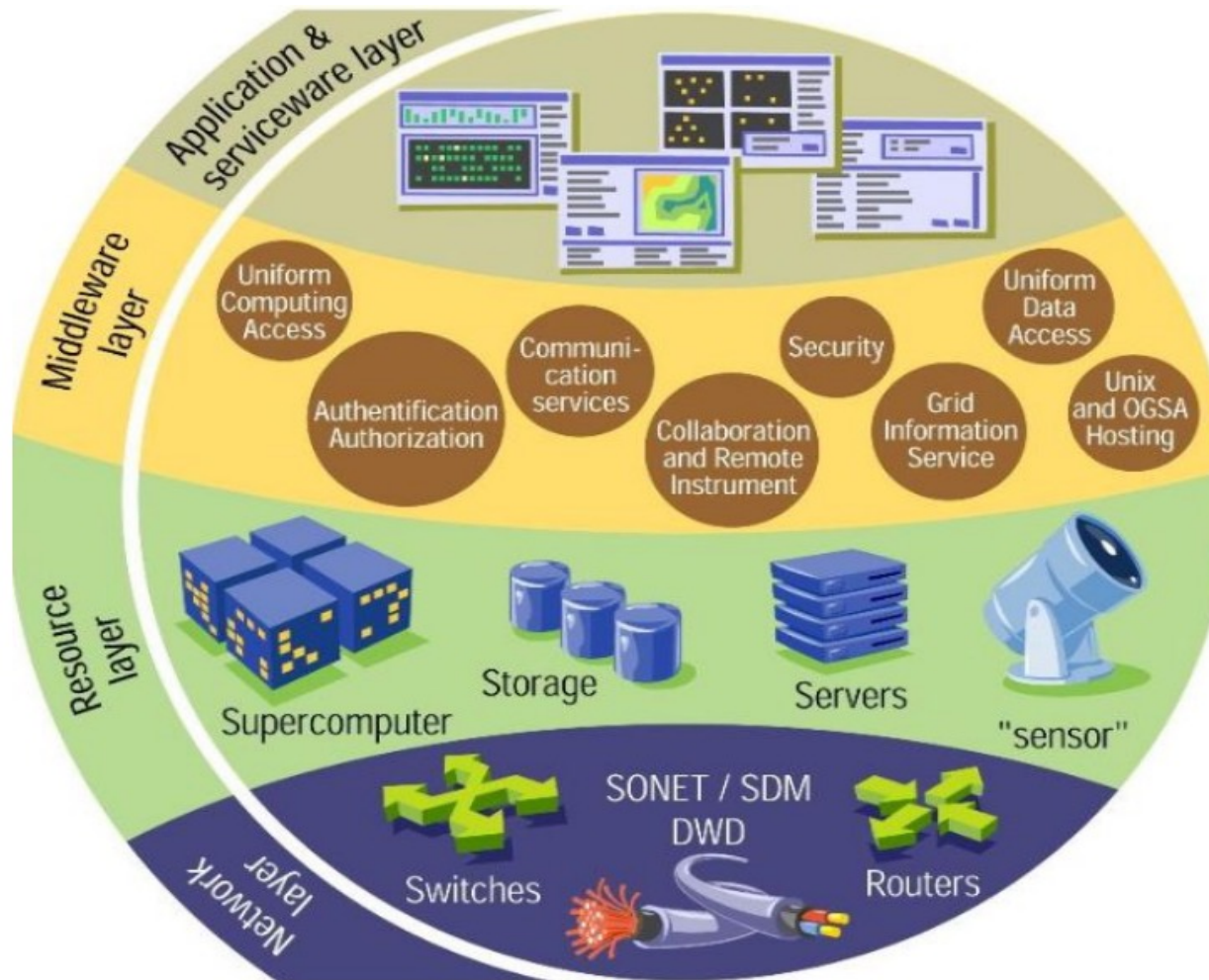
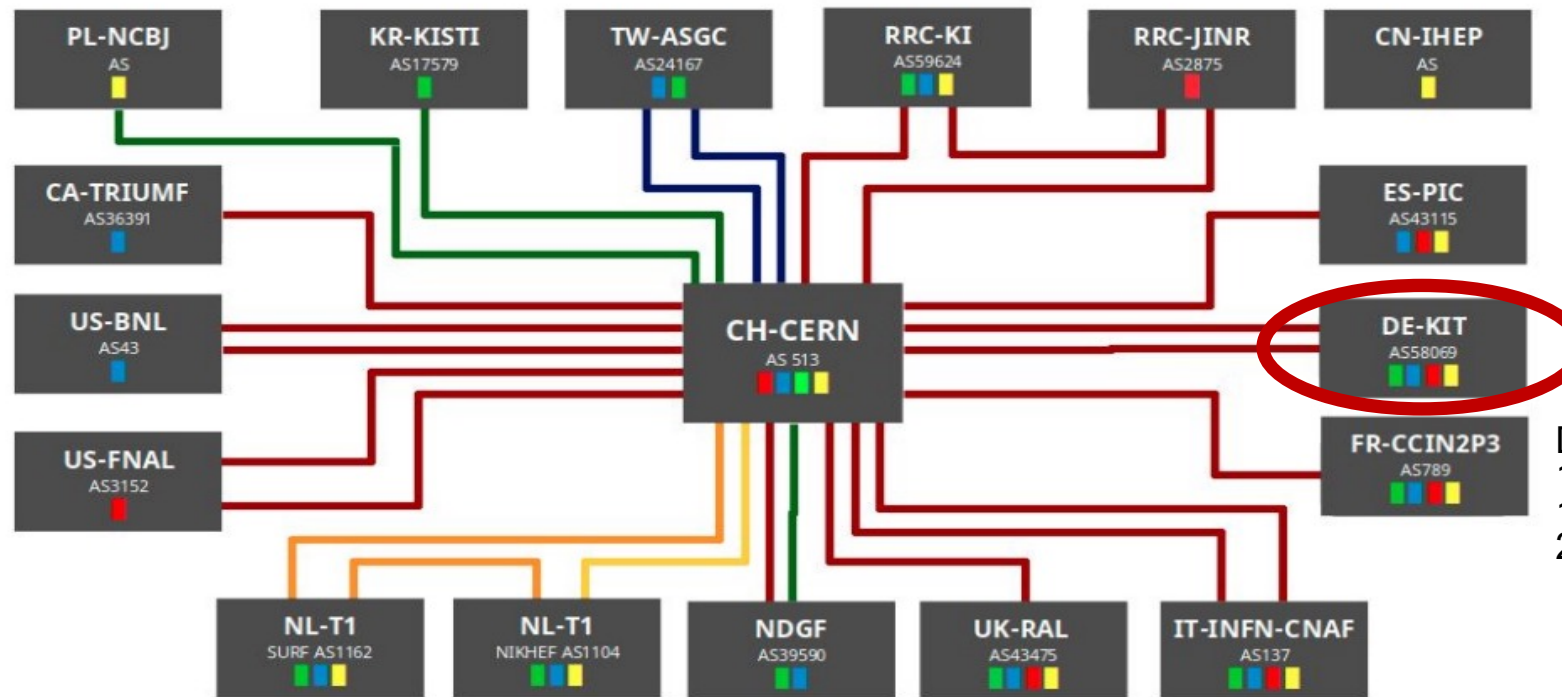# LHC accelerator and experiments

# WLCG



world wide distributed computing engine for the  Large Hadron Collider emerging data Arround 2000:

the memberstates decided for remote and distributed installations. The foundation for the **W**orldwide **L**arge hadron collider **C**omputing **G**rid (WLCG) was layed

# raw data calculation and tape storage centres



## Tier0 to Tier1s network

Dedicated 10/100/400Gbps links
16 Tier1s + 1 Tier0
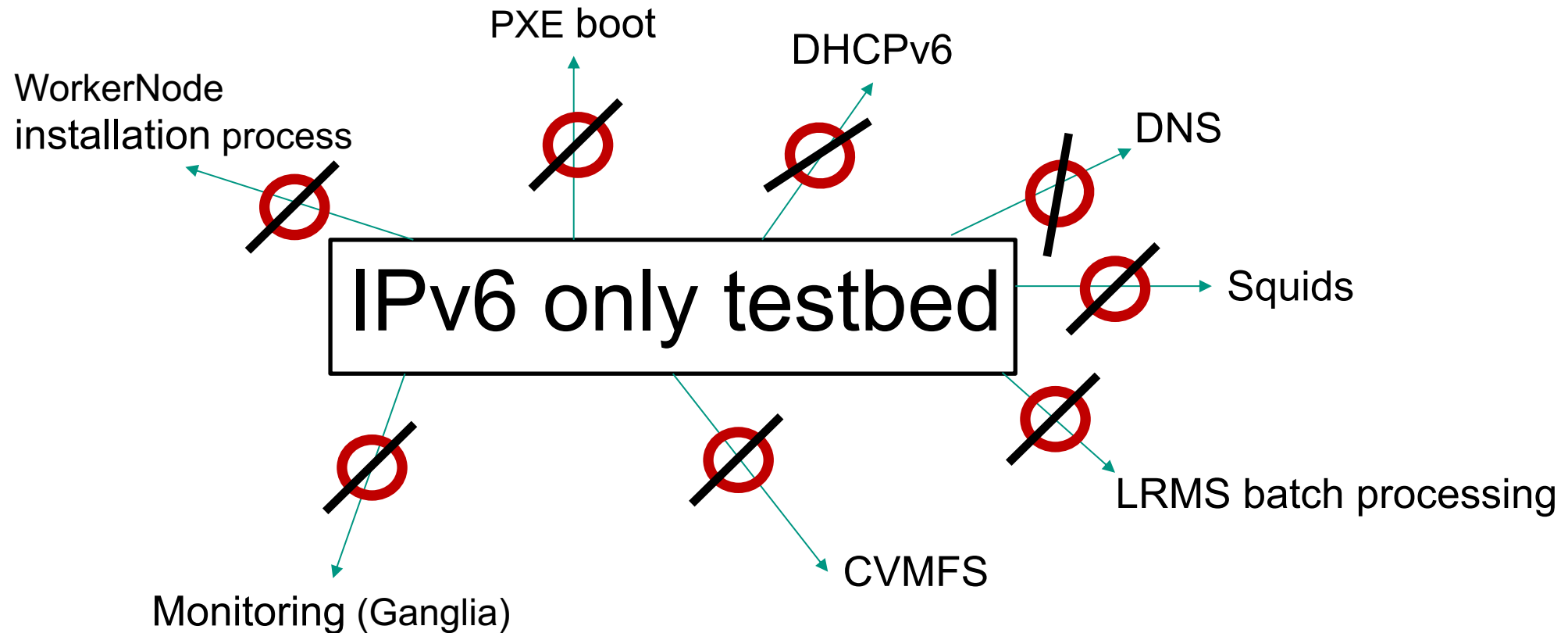12 countries in 3 continents
2.1 Tbps to the Tier0

# GridKa



- **worker node farm**
  - 217 aktive hardware systems
  - 42500 compute cores

- **online-storage**
  - 99 PB effectiv storage capacity
  - 6824 HDDs
  - 100 Server

- **nearline-storage**
  - 85 PB saved on tapes
  - 135PB available capacity on tapes

- **wan network**
  - 2 x 100Gb/s direct to CERN (LHCOPN)
  - 2 x 100Gb/s to DFN (LHCONE overley)
  - 2 x 100Gb/s to Belwue

# Building IPv6 Testbed

## HEPiX- IPv6 working group asking for IPv6 only testbed

# DE-KIT – workernode migration towards IPv6



Pro-active IPv6 Monitoring at DE-KIT

**packet number decreased from monitoring in 2022 to 2023**
**- power budget depending workernodes were switched of (while still LHC MOU is full filled)**

# Detailed monitoring at DE-KIT (GridKa)

- Monitor all comunications between WorkerNodes and
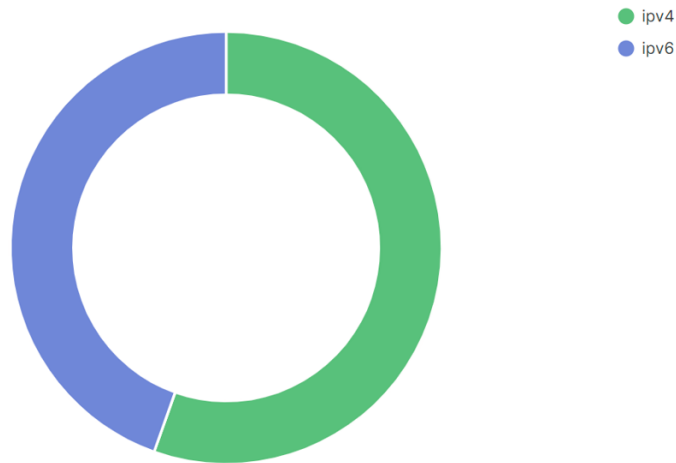
  - administration
  - job submission
  - Storage
  - ...

# Monitoring of process intercomunication at DE-KIT (GridKa)

- with packetbeat collecting network data

- logstach pushing the data to opensearch (former elastic search) for storing the data

- kibana for visualizing
   (no opensearch – only easy search requests)

   • started with a small set of workernodes (storing the data „longterm" → ~ 6 days)

   • while enlarging the set of workernodes graduately
      data keeping time had to be limited to less than one week only
      (for not exceeding the storage size of 0,5 Tbyte)
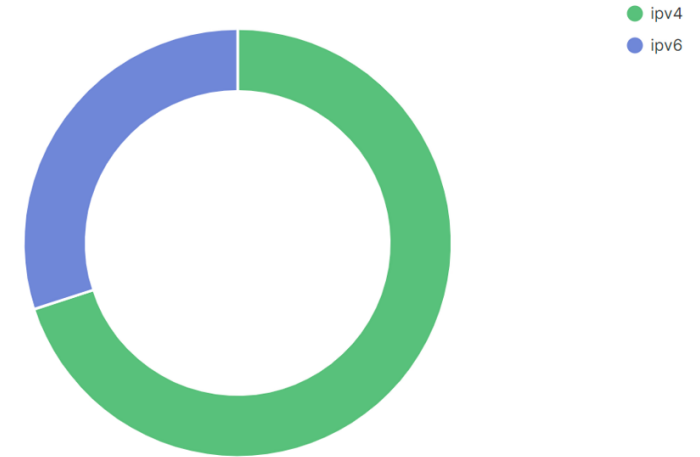
● Identify IPv4 protocol usage

# Snapshot of a dashbord

at 08.09.22
all worker nodes already dual-stack
deployed
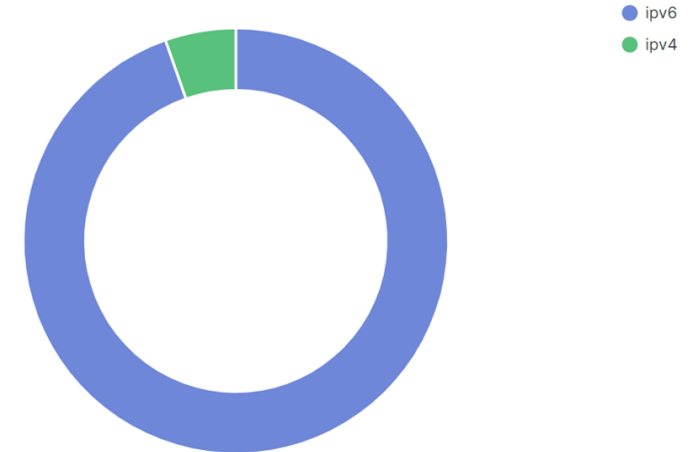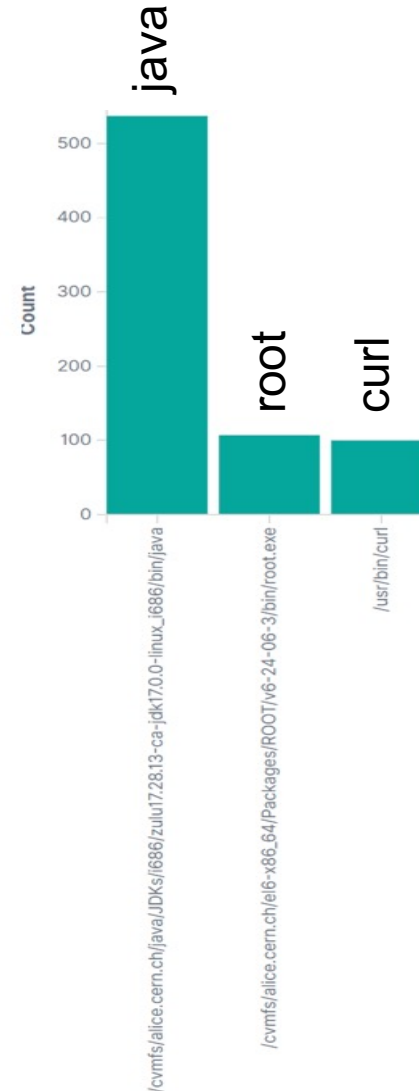
**IPv4/IPv6 incoming traffic**
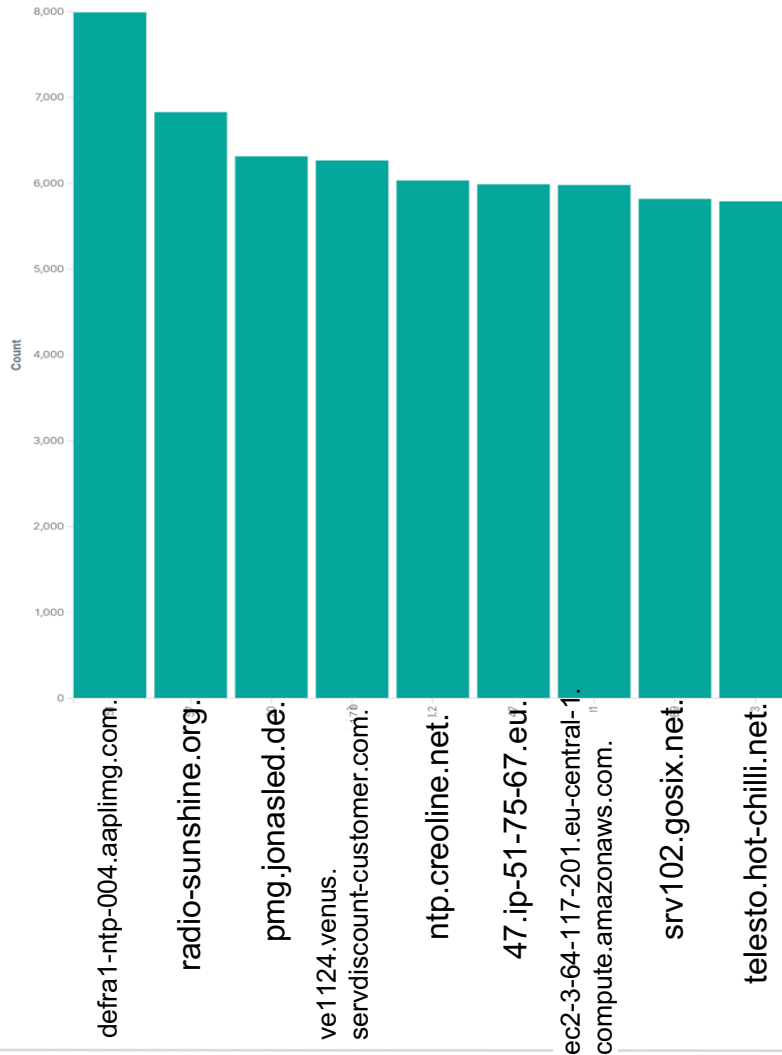
- ipv4
- ipv6

**IPv4/IPv6 Packages**

- ipv4
- ipv6

**IPv4/IPv6 outgoing traffic**

- ipv6
- ipv4

# NTP ?



- Many NTP / port 123 connections
  - During 24 hours approx. 210.000
  - NTP → IPv4 only (depending on dualstack enabling of rack-manager (40.000 internal))
  - Monitoring was first pointing especially 10.1.12 and 10.1.18 → checking later showed that much more racks running ntp check via private addr. (NAT)
  - 160.000 external communications → some of the destination server have quite dubious „names"
- process-tracking
  - The numbers of NTP communication process and matched process is not matching yet

# S O L V E D

- NTP.ORG
  → returns sometimes funny addresses
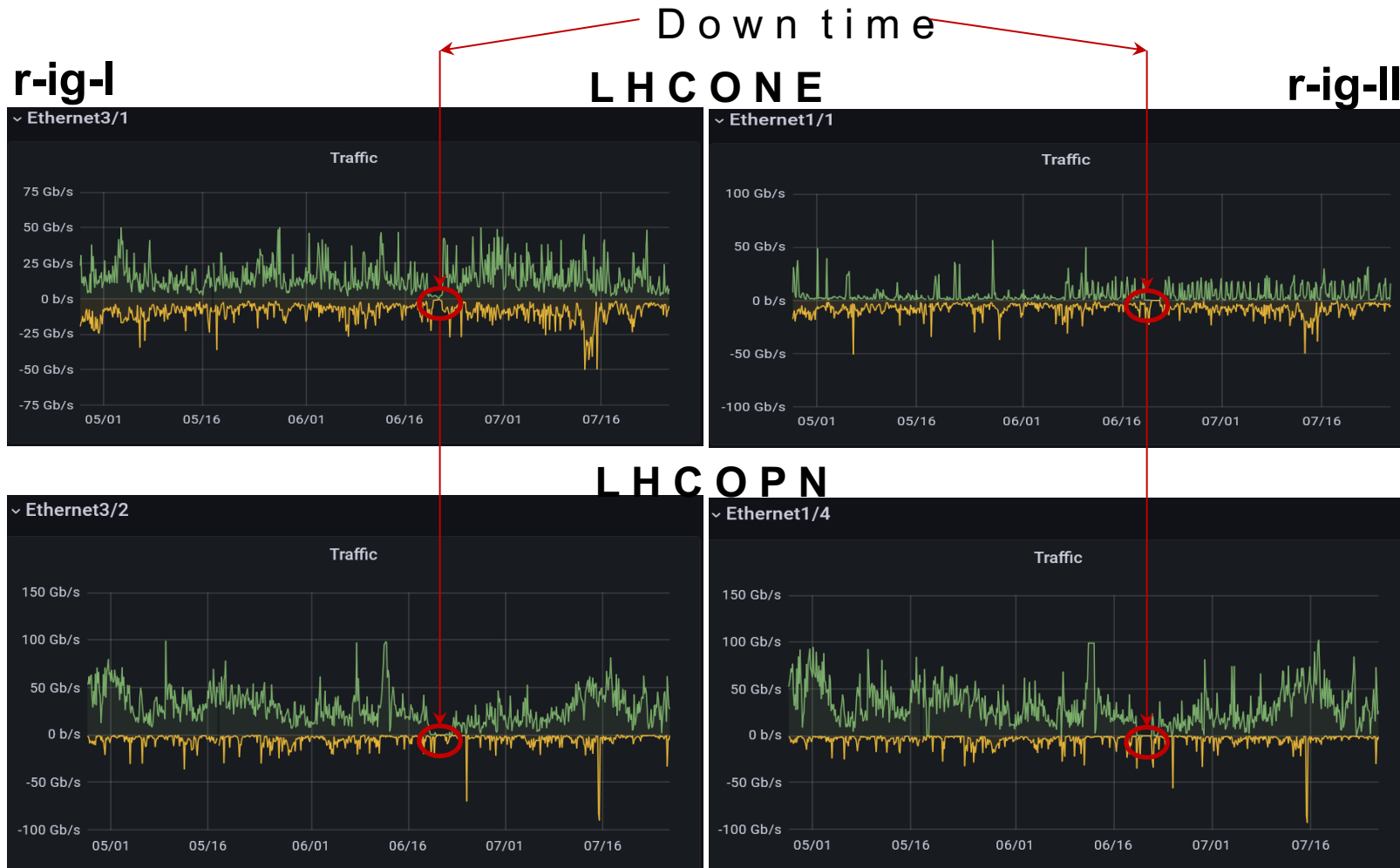
# dCache upgrade to 7.2.15

## Upgrade from dCache version 6.2.34 to 7.2.15

Two day downtime at June 20th and 21st 2022

- HTTP-TPC transfers now prefer IPv6 address, if both endpoints support it.

- fixed handling of Storage Resource Reporting (SRR) requests over IPv6

- Handle IPv6 address when running HTTP(s) Third Party Copy (TPC) with gridsite delegation

- Storage Resource Manager (SRM) : Fix IPV6 logging for SRM
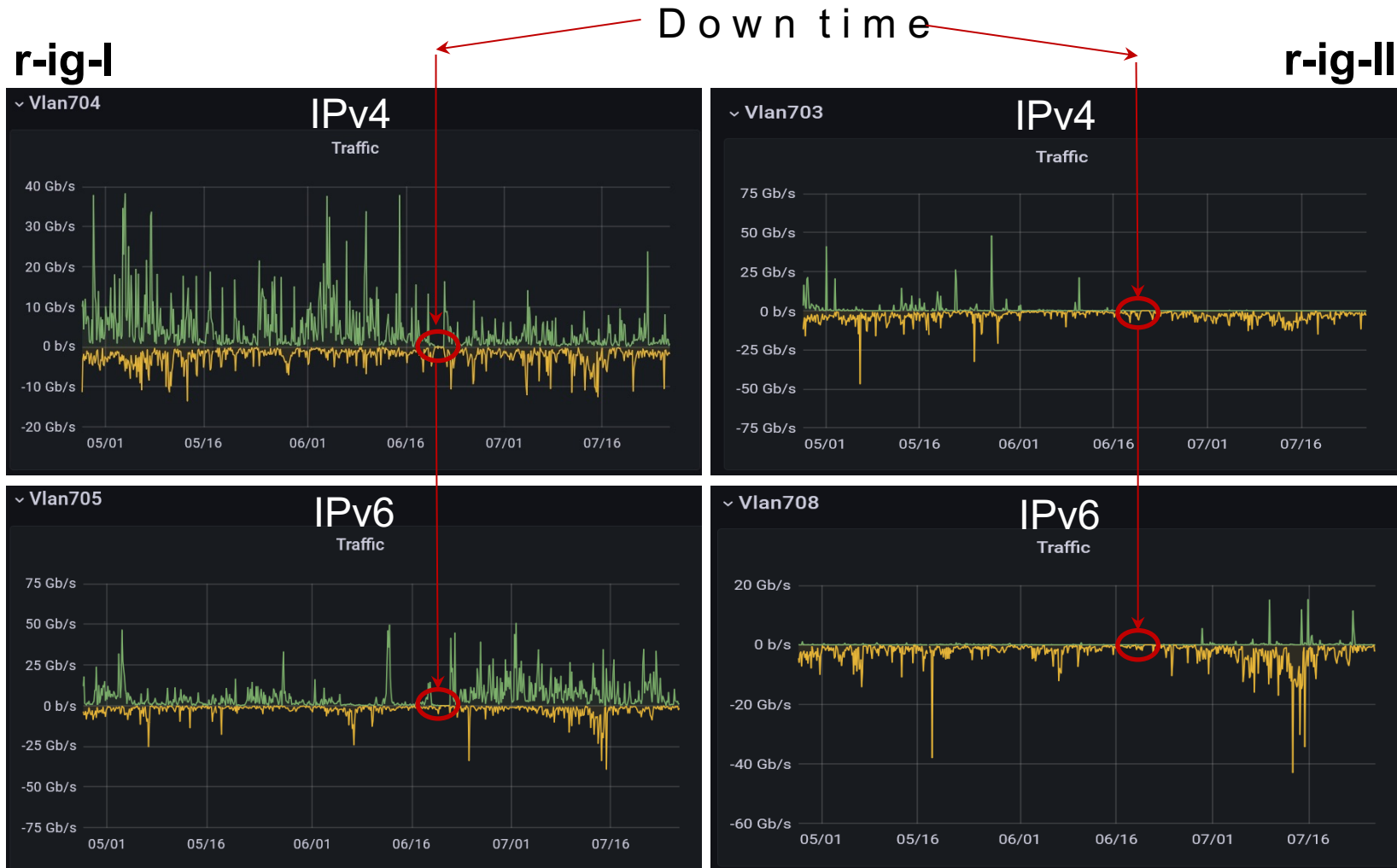
# WAN Interfaces



**r-ig-I** (DE-KIT Border Router):
left two Interfaces
- Ethernet 3/1 (Internet + LHCONE) +
- Ethernet 3/2 (LHCOPN)

**r-ig-II** (DE-KIT second Border Router):
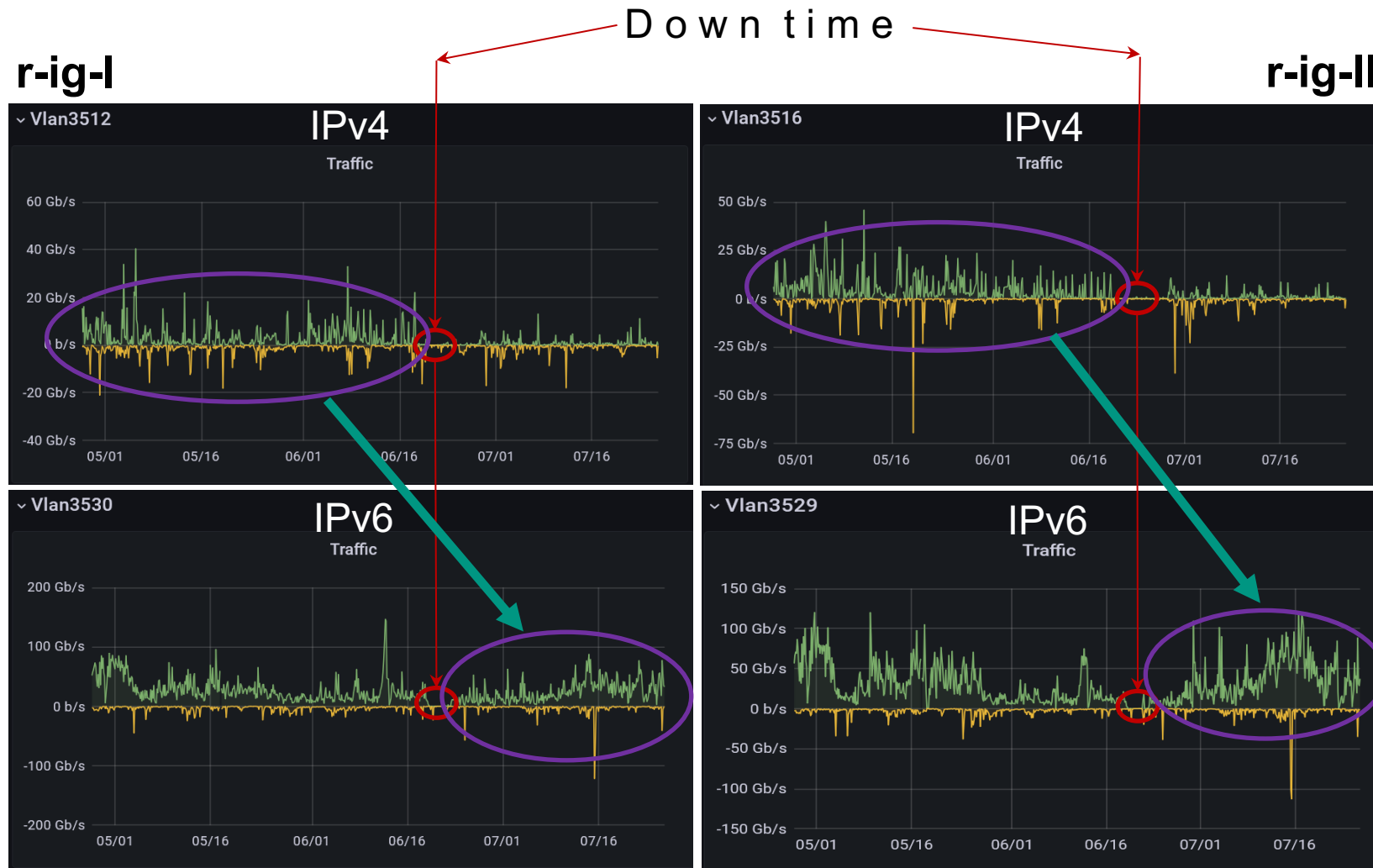right two Interfaces
- Ethernet 1/1  (Internet + LHCONE) +
- Ethernet 1/4  (LHCOPN)

# LHCONE IPv4 / IPv6
## transfer pattern after downtime



Graph over 90 days
Traffic of LHCONE
moved partioly from the IPv4 vlans
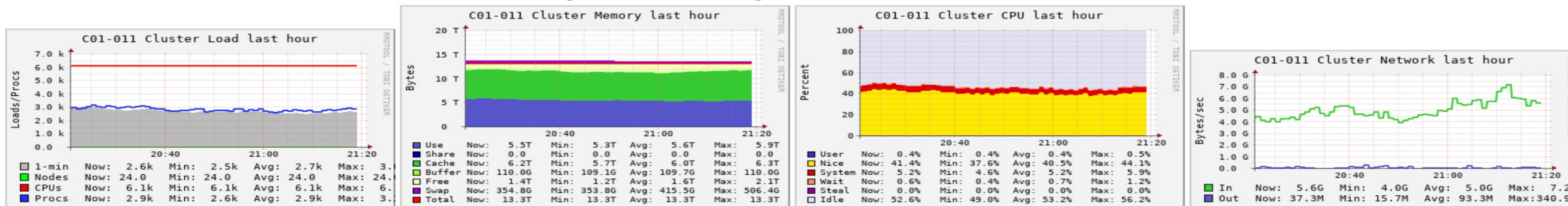after the downtime to the IPv6 Vlans

# LHCOPN IPv4 / IPv6
## transfer pattern after downtime



Graph over 90 days
Traffic of LHCOPN
moved from the IPv4 vlans
after the downtime to the IPv6 Vlans

# Monitoring

## G A N G L I A



- Migration of Ganglia to IPv6 will not persuit
- Ganglia will be replaced by opensearch, kibana and grafana

# Logstach → is now IPv6

Logstach (port 5047) → dual-stack deeployed

statistic:

28-07-2022 → IPv4 385k – IPv6 1,41M

23-10-2022 → IPv4 476k – IPv6 1,39M

23-12-2022 → IPv4 227k – IPv6 450k

**30-10-2023** → IPv4 906k – IPv6 864k

# Closer look at DNS



**IPv4** **IPv6**
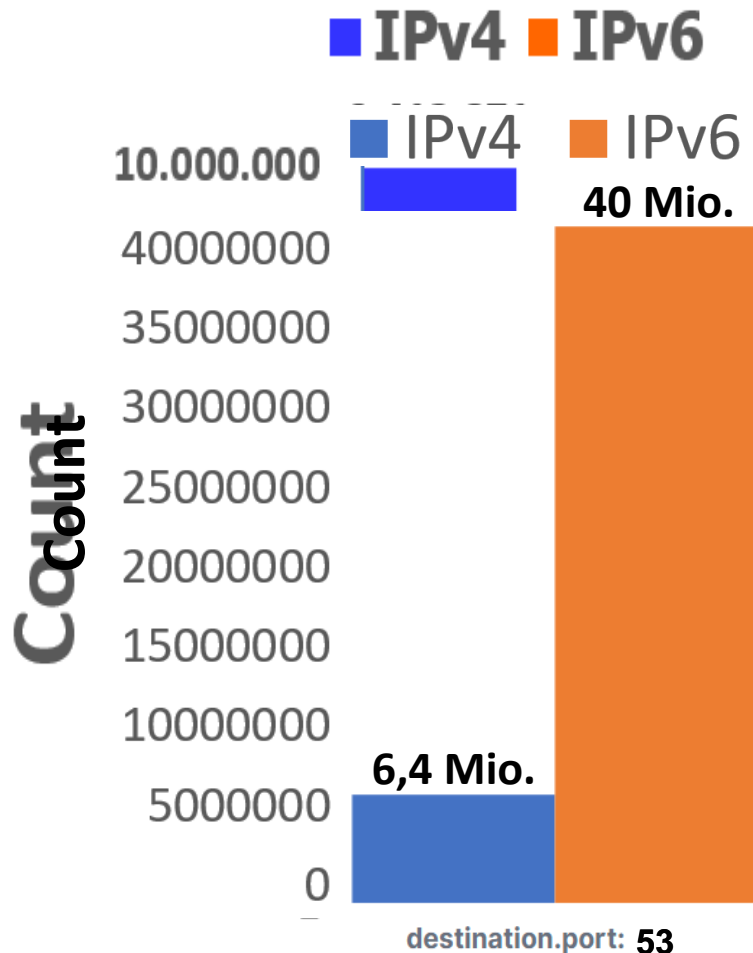
- GridKa DNS:
  - IPv4 only count : 9,412,871 (24 hours)
  - DNS (Bind) Server and WN is already dual-stack
  - at WN resolve.conf first lines IPv4
    - Make sure IPv6 DNS server addresses listed and
    - place it before IPv4
    - every new deployed host:
      the first lines are IPv6 resolver addresses
      of the **resolve.conf** file **followed by the IPv4 addresses**
      - `nameserver 2a00:139c:address`
      - `nameserver 2a00:139c:address`
      - `nameserver 10.privat-address`
      - `nameserver 10.privat-address`

  → **Resolve.conf update: reprovisioning required**

# Administrative Services

- at each rack is a Rack Manager deployed:
    - Starting in 2001 with private IPv4 only
    - Migration process initiated (but still in progress)
      → enable dual-stack (AAAA)
        - NTP
        - rsyslog  (→ migration → still pending (port 514))
        - Monitoring (GmonD → Ganglia Client)
        - DHCP  (→ migration to DHCPv6 pending)

# WN – deployment process

- Redhat Satellite Server (foreman)
  - Used for management of most GridKa hosts:
    - Manages redhat Subscriptions
    - Controlls kickstart installations (DHCP / PXE)
    - Provides yum repos
    - Provides CA (certificate authority) and ENC (encryptor) functionalities for puppet
  - Uses modular architecture. Additional functionalities can be added via so called capsules
    - TFTP server (IPv6 ready - dual-stack)
    - Puppetmaster (IPv6 ready - dual-stack)
    - Pulp  (software repository management (IPv6 ready - dual-stack))
    - DNS (IPv6 ready - dual-stack)
    - DHCP (currently DHCPv6 capsule not available)

# Details of Squid

- SQUIDS (Proxyserver and Web-Cache):
  - some SQUIDS still IPv4 only (**migration to dualstack in proccess**)
  - Significant part of connections via public IPv4
  - => to check:  if CVMFS can prefer IPv6?
    **(CVMFS → CernVM-File-System)**
    - CVMFS sending via http request to squid
    - CVMFS has DN configuriert that needs to be resolved
      → default chooses IPv4 address
    - **Solution** => cvmfs_ipfamily_prefer=6 → **not tested yet** `(end of 2022)`

# SQUIDS migrated all to dual-stack

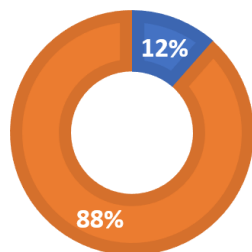During the second half of 2022 all SQUIDS migrated to dual-stack deployment

CVMFS now
- manly IPv6 but:
- on WorkerNodes uses IPv6 (with deployed flag: CVMFS_IPFAMILY_PREFER=6 )
- CVMFS frontier uses still IPv4 even while both systems dual-stack
- but switching of IPv4 → froniters will operate over IPv6
- the CMS CVMFS frontiers offers in site-local-config.xml the Option:

**<frontier-connect>**
...
   **<prefer ipfamily="6"/>**
...
**</frontier-connect>**

| **26-07-2022** | **23-10-2022** | **25-10-2023** |
|:---:|:---:|:---:|
| IPv4 : 1,25 mio. IPv6: 9,6 mio. | IPv4 : 4,44 mio. IPv6: 18 mio. | IPv4 :   64 k   IPv6: 22 mio. |

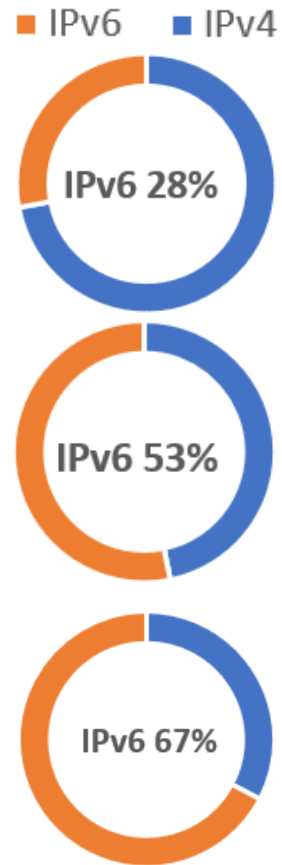# Batch-Processing -- LRMS (HT-Condor) all dual-stack

- LRMS (**Local Resource Management System**) HTCondor at GridKa (all dual-stack and set to **prefer** the protocoll **IPv6** (Port 9618/9)
  - 4080 – HTCondor (rooster-deamon) → migrated all towards IPv6 (HTCondor → startd)
  - percentage increased toward IPv6 at 28-06-2022→ IPv4: **895k** to IPv6: **255k**
  - 1,2% IPv4    28-07-2022 → IPv4: 27k,  IPv6: 2,17 mio.
  - **11%**  IPv4    02-01-2023 → IPv4: 287k,  IPv6: 2,28 mio.
  - **18%**  IPv4    31-10-2023 → IPv4: 2,68 mio.,  IPv6: 11,7 mio.

**Less then 1% (0,0049%) of IPv4 is internal traffic**

**(communication with home → the LRMS demons uses protocol of Home-Institution)**
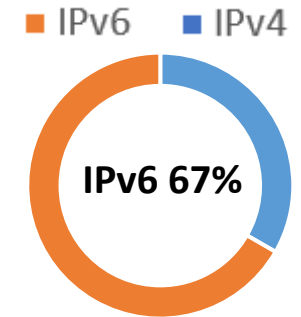
# A view statistics



- 15-04-2022:
  - IPv4: → 80 mio.
  - IPv6: → 31 mio.

- 26-07-2022:
  - Ipv4 → 44 mio.
  - Ipv6 → 50 mio.

- 23-10-2022:
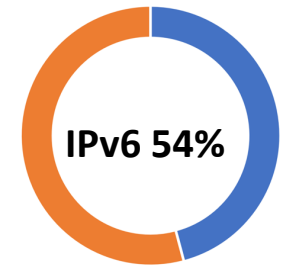  - IPv4 →  69 mio.
  - IPv6 → 142 mio.

- 20-12-2022:
  - IPv4: → 42 mio.
  - IPv6: → 86 mio.

- 31-10-2023:
  - IPv4: → 122 mio.
  - IPv6: → 144 mio.

(packets in 24 hours)
# of WorkerNodes included in the statistic expanded

# Next steps

- migration of Rackmanager – work in progress
- Narrow down the still IPv4 communication
  - packet monitoring configured
    - to list all unhandled IPv4 packets
      - 4080 – Condor rooster Monitor deamon → solved
      - 8884 – Alice: operation report
      - 2049 – NFS
      - 8649 – Ganglia gmond
      - 1094 – XrootD
      - 961[89] – LRMS (less than 1% only internal to WN-Farm)
  - PXE – Boot + DHCPv6 (first boot addr. Distribution)
- Identify the next service for IPv6 migration tasks

**Ports**

**IPv4 Adresses**

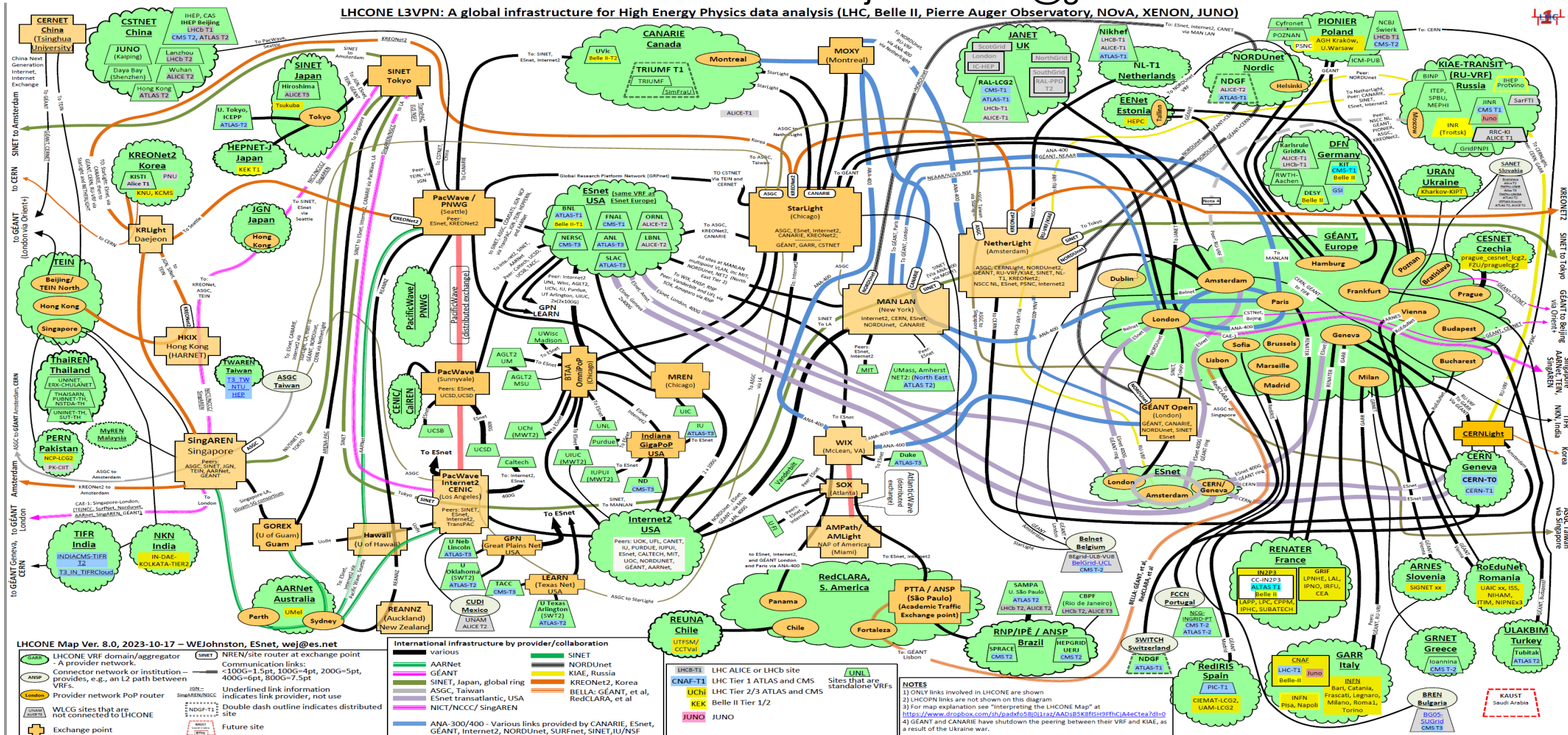# Thx for your attention

# backup Slides

# LHCONE

William E Johnston <johnstonwe@gmail.com>

LHCONE L3VPN: A global infrastructure for High Energy Physics data analysis (LHC, Belle II, Pierre Auger Observatory, NOvA, XENON, JUNO)

# Details of Alice VOBoxes:

- ALICE VOBoxes:
  - Client to VOBox prefers IPv4 (ALICE Monitoring (UDP))
  - => to check the possibility of IPv6 migration with ALICE (still ongoing)
    - dual-stack enabling works and
    - if Preference towards IPv6 is possible
    - ALICE is constrained by IPv6 unavailability on other sites
  - → advice of Alice : switch of IPv4 at VO-BOX (the none monitoring VO-BOX)
    - Timing still under discussion
  - Monitoring (port 8884 / IPv4 only) → 11 Mio. (/24 hours)

- XRootD:
  - via public IPv4 (ALICE)
  - All ALICE XRootD SE are dual-stack deployed
  - older version of XRootD → upgrade to current XRootD should improve, is still pending
  - → advice of Alice : get IPv6 ready – but wait for switching it on till complete Alice is IPv6 ready

- Dest port 1094 –Ipv4/ipv6 → XRootD (alice, belle2, atlas, cms)