# The Czech Distributed Tier-2 Center

## Alexandr Mikula

## Institute of Physics
## of the
## Czech Academy of Sciences

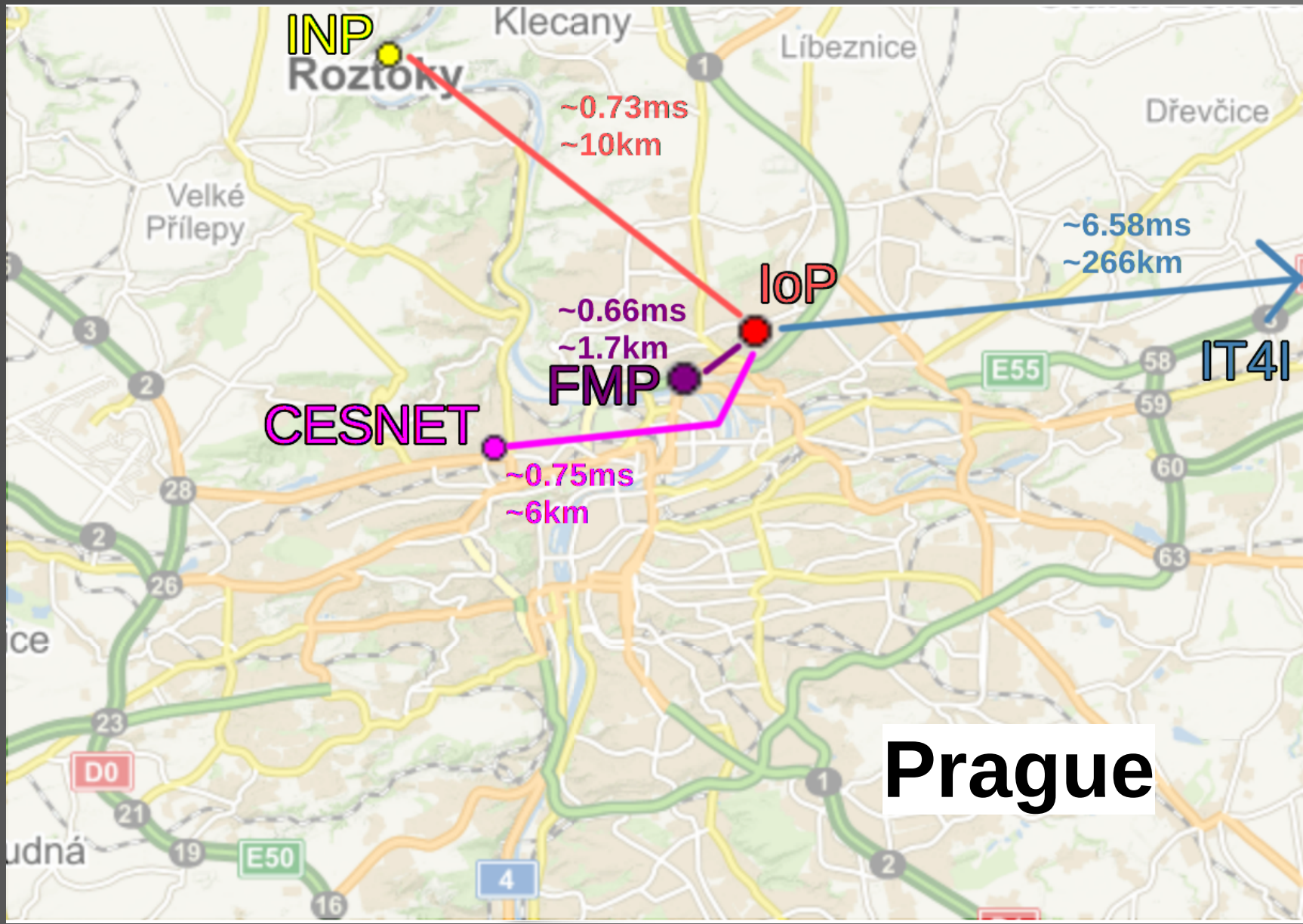Martin Adam, Dagmar Adamová, Jiří Chudoba, Petr Horák, Jana Uhlířová,  Petr Vokáč

Akademie věd
České republiky

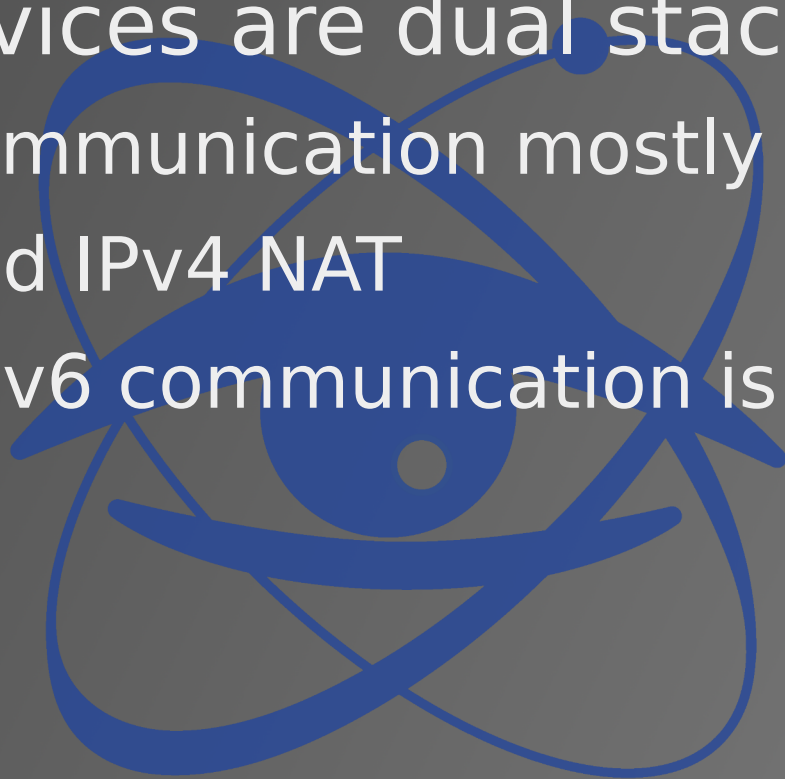FYZIKÁLNÍ ÚSTAV
Akademie věd ČR, v.v.i.

# CC IoP CAS

- Distributed site located in CZ
- Several institutions involved
  - Institute of Physics CAS (IoP)
  - Institute of Nuclear Physics CAS (INP)
  - Faculty of Mathematics and Physics CUNI (FMP)
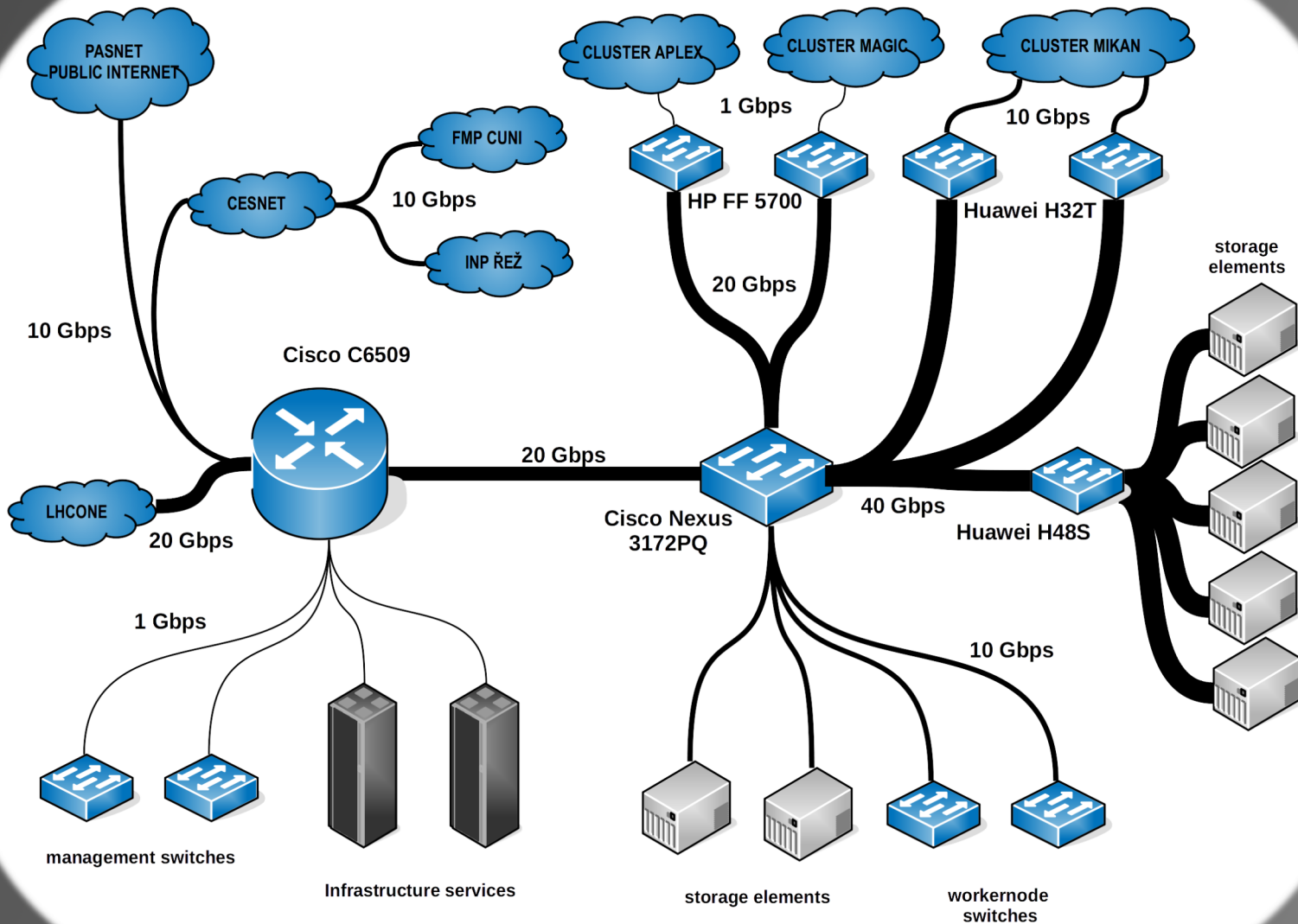  - CESNET
  - IT4Innovations

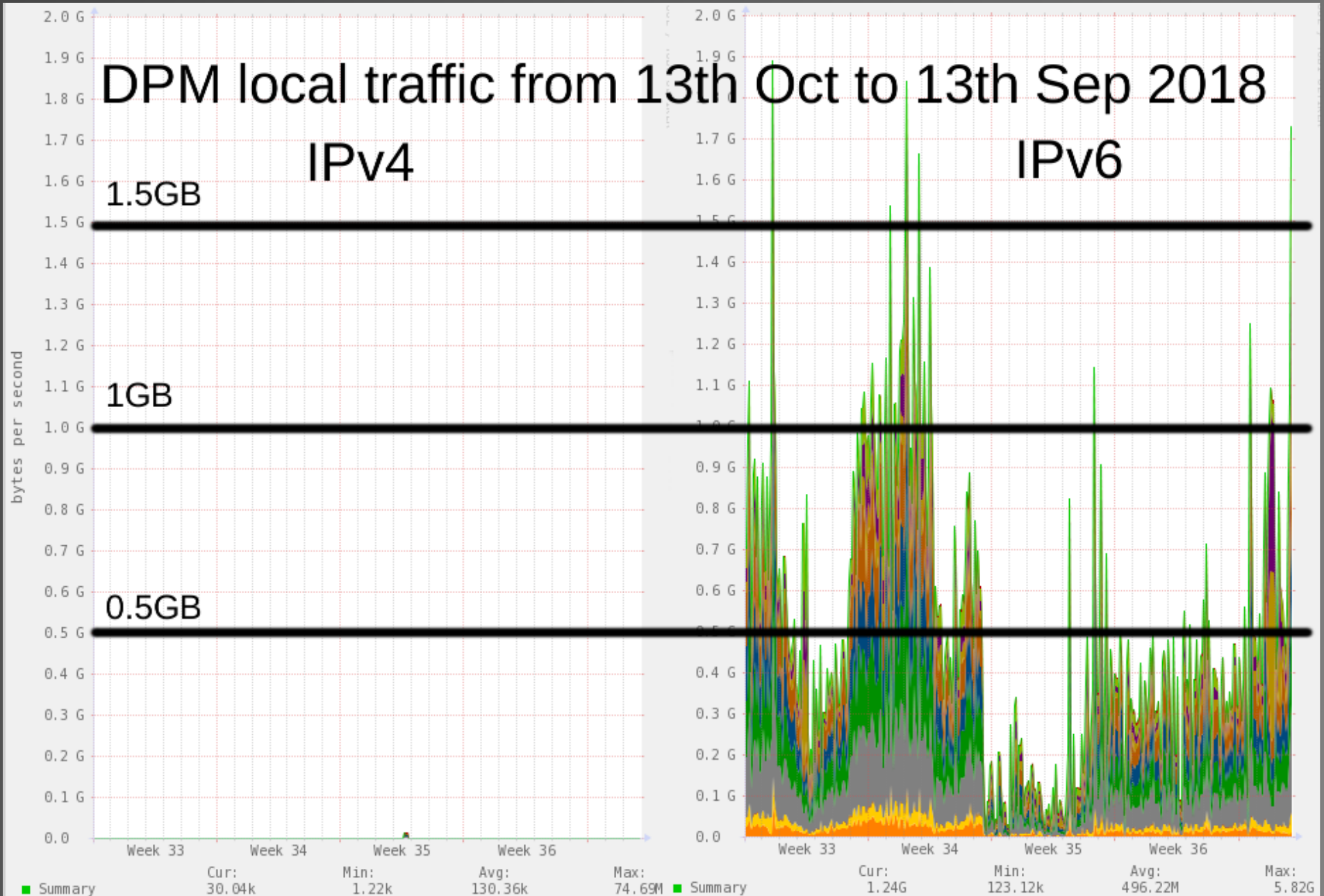# CC IoP CAS

# Network Connectivity

- Most of services are dual stack IPv4/IPv6
  - Internal communication mostly IPv6
  - WNs behind IPv4 NAT
  - External IPv6 communication is growing
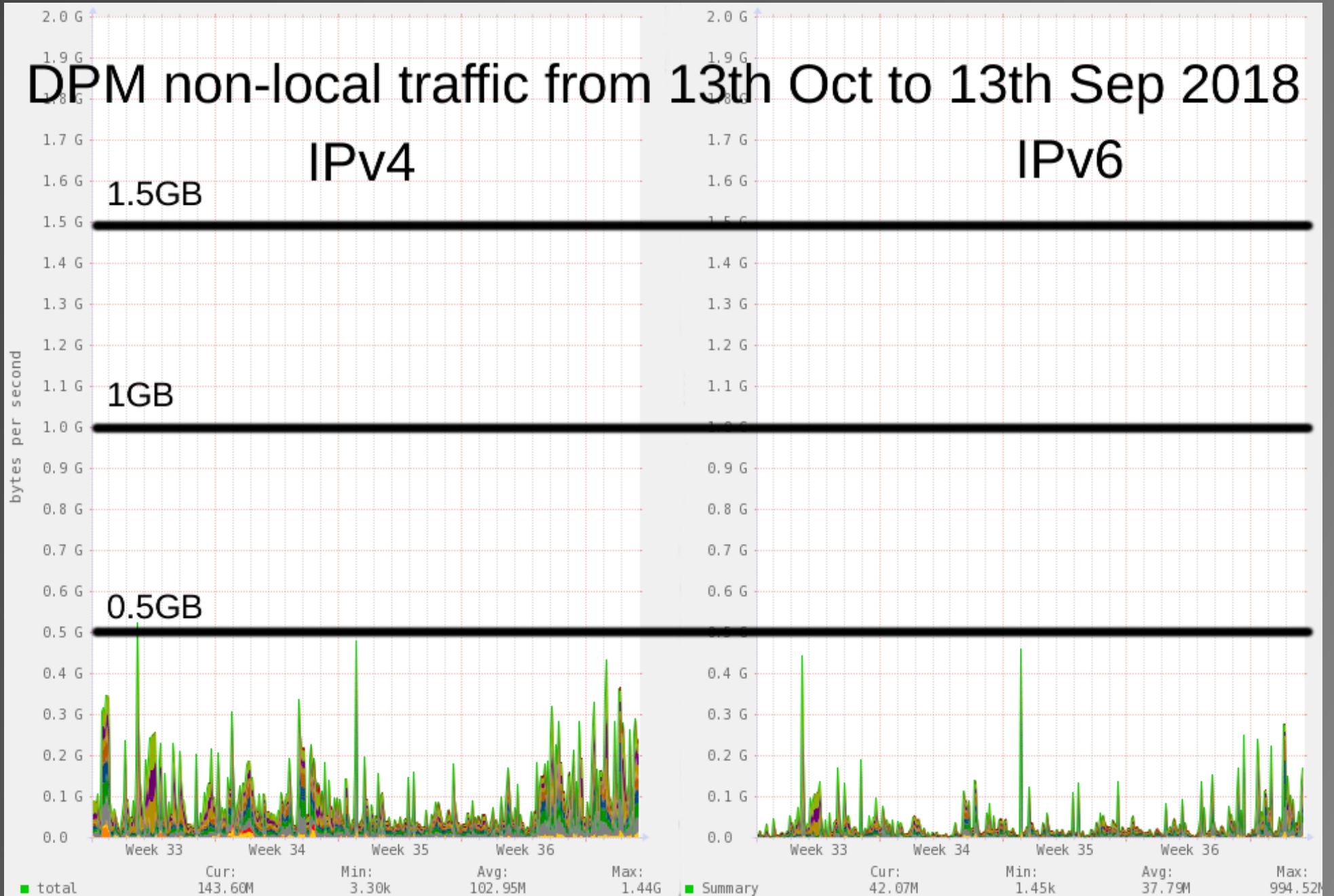
# Network connectivity

DPM local traffic from 13th Oct to 13th Sep 2018

IPv4

IPv6

DPM non-local traffic from 13th Oct to 13th Sep 2018
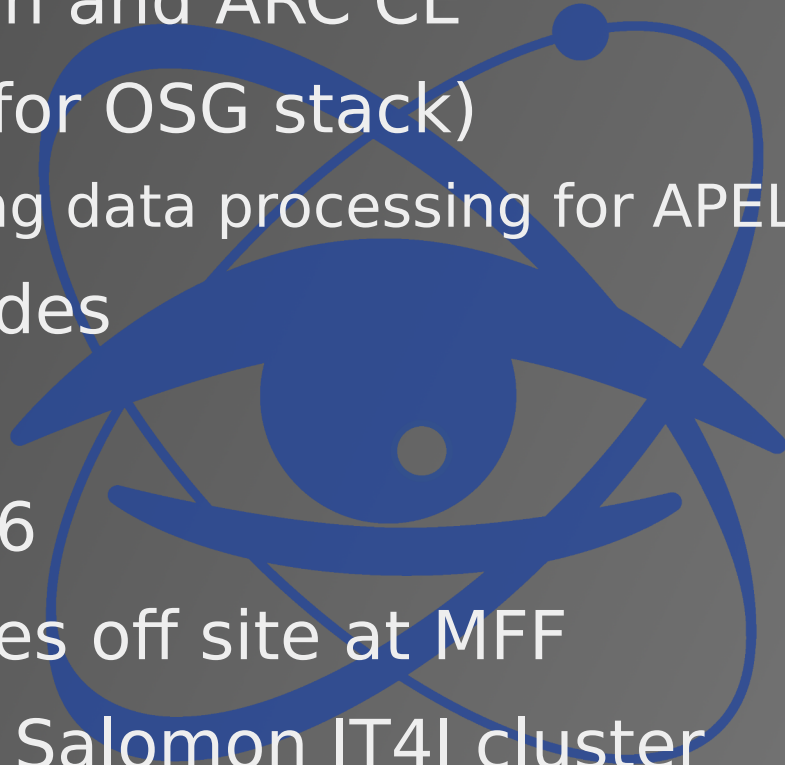
IPv4

IPv6

1.5GB

1GB

0.5GB

# Supported Projects

- Tier 2 site
  - ATLAS and ALICE
- Other projects
  - Pierre Auger observatory
  - Cherenkov telescope array
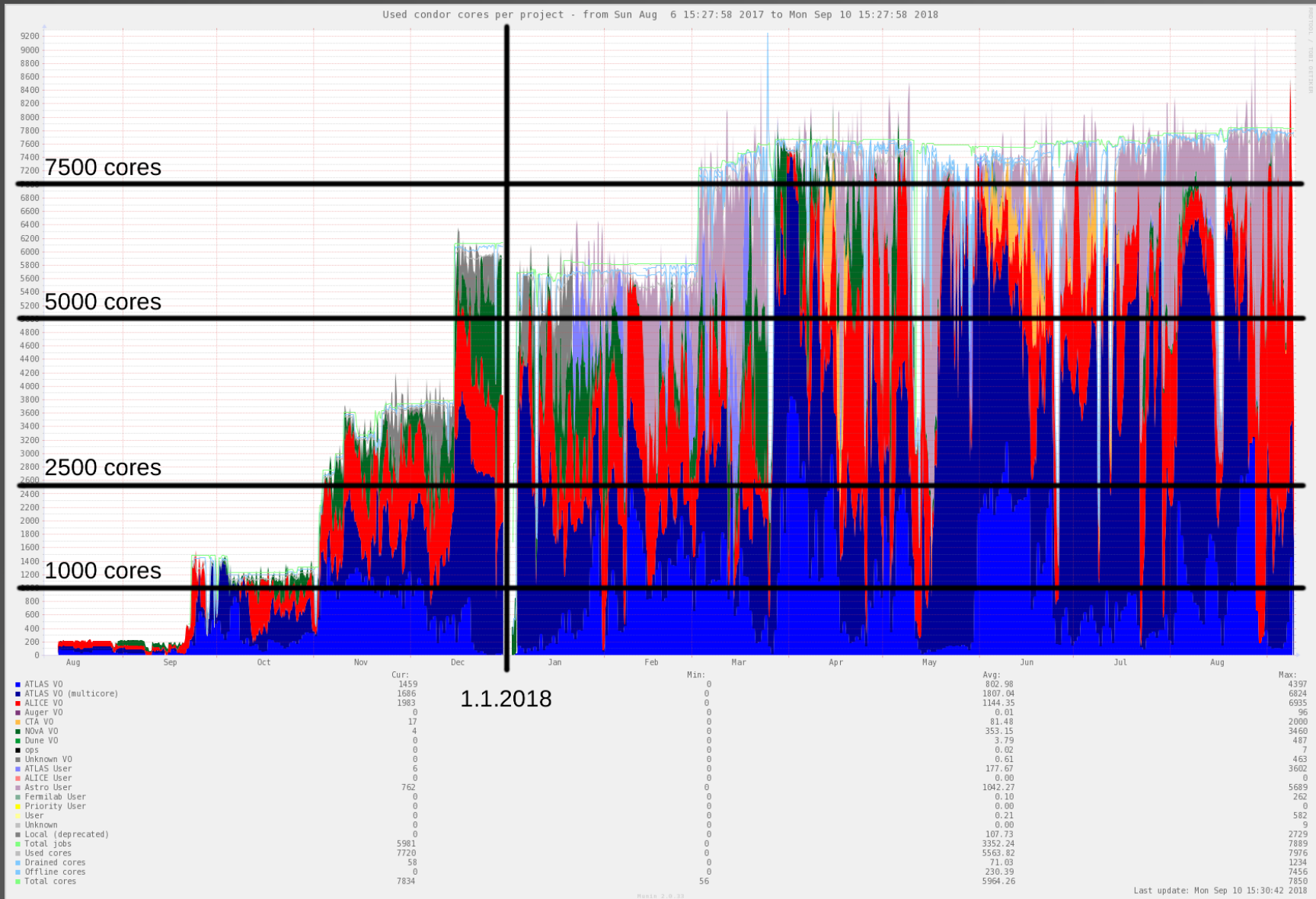  - NOvA
  - DUNE
- Local users

# Computing power

- HTCondor Batch and ARC CE
- (HTCondor CE for OSG stack)
  - Own accounting data processing for APEL
- 285 Worker nodes
- 7800 cores
- 77k HEPSPEC06
- 47 Worker nodes off site at MFF
- ARC CE for the Salomon IT4I cluster
  - Opportunistic resource for ATLAS

# Computing power

# Other resources

- NFS servers
  - 5 servers
  - 212 TB available space
- 40TB Hitachi SAN array
- KVM Hypervisors
  - 3 AMD based severs with 128GB RAM
  - 2 Intel based servers with 256GB RAM
  - FC connection to SAN array
  - 42 virtual machines

# Monitoring

- Nagios with Check-MK
- Munin
- Observium
- Ganglia
- Puppet Explorer
- On site developed scripts

# Monitoring

# Software management

- Puppet configuration management
  - Git history control
- Spacewalk package management
- Migration to the Foreman

# Storage

- 12 DPM Servers
  - 3.8 PB available storage
  - ATLAS, Auger, CTA, …
- XrootD
  - ALICE experiment
  - 1.6 PB available storage
  - 3 servers on site
    - 1.3 PB
  - 4 servers off site at INP in Řež
    - 340.6 TB

# Software testing

- ZFS as an alternative for HW raid
  - Mixed results based on HW (type and age)
  - Still in testing
  - Some features are not marked stable
  - CDDL license

# Software testing

- DPM
  - Used version 1.10.3
  - DOME flavour core
  - Several bugs encountered and reported

# DPM short intro

- Storage management
- Easy to setup and maintain
- Multiple VO support
- X509 authentication
- Multi protocol support
- Great user/admin support
- http://lcgdm.web.cern.ch/dpm

# DPM DOME migration experience

- Easy to migrate

- Downtime needed

- In production for ATLAS on several sites besides our
  - Gaining some experience with it before recommending migration
  - Several bugs found, reported and fixed 1.10.4
    - Now in EPEL testing
    - Great work of Petr Vokáč
  - No major breakdown

- Own fast scripts for detecting dark and lost data

- Developers plan to drop SRM support
  - Active in WLCG DOMA group to bring TPC for xrootd and http
  - Built-in new non-SRM storage reporting

# DPM performance comparison

| command | protocol | legacy (s) | DOME (s) |
|---|---|---|---|
| ls | SRM | 0.48 | 0.46 |
| | XROOT | 0.28 | 0.22 |
| | DAVS | 0.24 | 0.9 |
| xattr | SRM | 0.59 | 2.85 |
| | XROOT | 0.91 | 0.33 |
| | DAVS | N/A | N/A |
| copy 1M file => DPM | SRM | 27.10 | 17.11 |
| | XROOT | 17.10 | 3.69 |
| | DAVS | 17.22 | 3.88 |
| copy 1M DPM => file | SRM | 23.51 | 18.36 |
| | XROOT | 4.65 | 5.01 |
| | DAVS | 5.43 | 6.06 |
| rm | SRM | 9.68 | 9.32 |
| | XROOT | 7.32 | 2.70 |
| | DAVS | 8.20 | 3.26 |

# DPM DOME bugs in 1.10.3

- Host not known error on el7
- SRM/DOME reported space discrepancy
  - Own script for fixing the difference
- Many dome_dirlist requests exhaust the MySQL handles
- dpm-xrootd does not create dir structure in the disk server
- HTTP authentication on dual stack
- The periodic cache cleanup may cause crashes

# Thank you for your attention