



東京大学  
素粒子物理国際研究センター  
International Center for Elementary Particle Physics  
The University of Tokyo



# Status of ATLAS Tier2 Centre at Tokyo/ICEPP

3<sup>rd</sup> Sep. 2024

The 8<sup>th</sup> Asian Tier Center Forum (ATCF8)

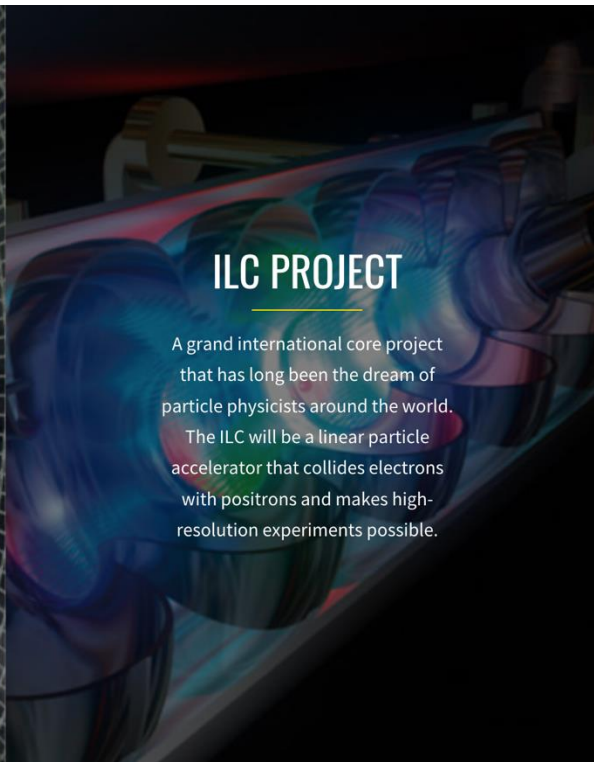
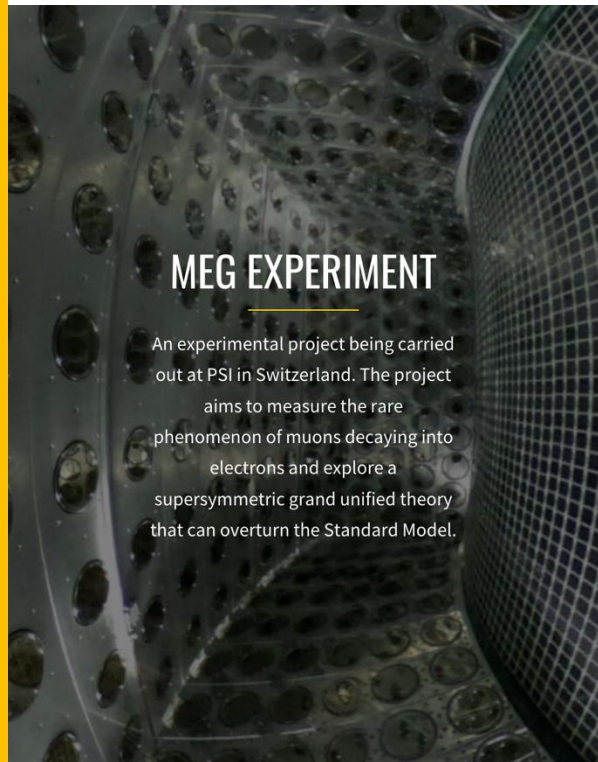
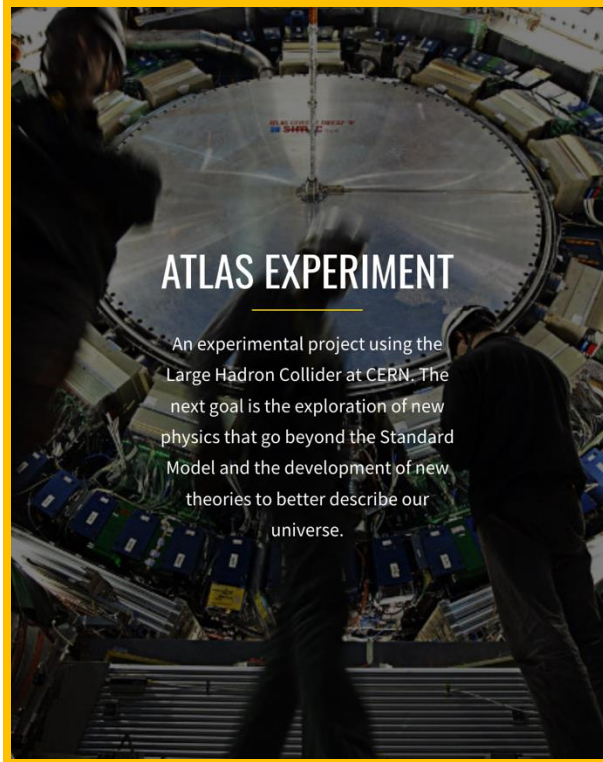
**Masahiko Saito**, on behalf of the operation team

ICEPP, The University of Tokyo

# International Center for Elementary Particle Physics (ICEPP)



## Main projects at ICEPP



## ATLAS-Japan group

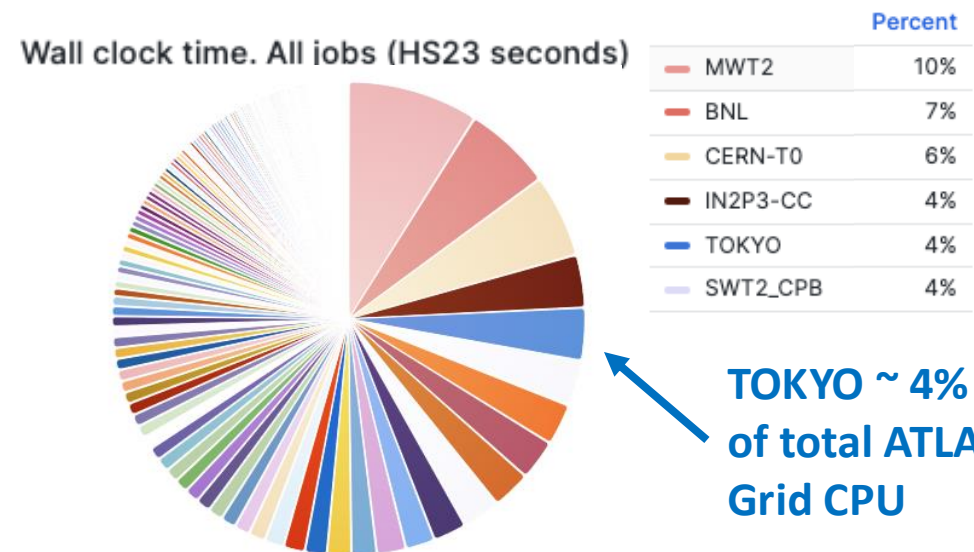
- 13 institutes and ~180 members (~40 members from ICEPP)
- Contributes to a wide area of the experiment
  - muon triggers, silicon tracker, **Tier2 operation**



➡ ICEPP operates **Tokyo Regional Analysis Center** for ATLAS/ATLAS-Japan

# Tokyo Regional Analysis Center

- Support ATLAS VO in WLCG (Tier2) and provide dedicated resources for ATLAS-Japan
- Tier2 (WLCG)
  - Worker nodes (ARC-CE/HTCondor): ~11k cores
    - ~4% of total ATLAS resources
  - Storage (dCache): ~13 PB
    - ~3% of total ATLAS resources
- Tier3 (ATLAS-Japan)
  - Interactive nodes: ~ 200 cores
  - Worker nodes (HTCondor): ~ 1.7 kcores
  - Storage (GPFS): 3 PB
  - GPU resources: V100, T4, A6000(x8)



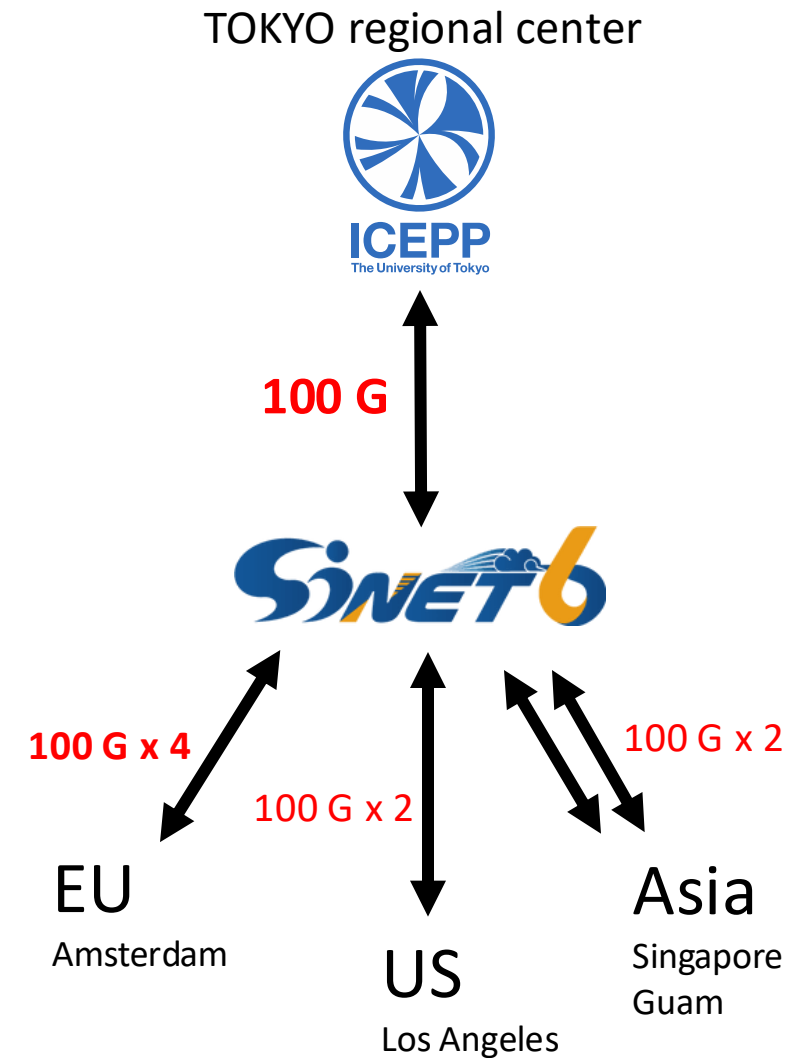
# Network

## Tokyo Tier2 RC ↔ SINET6

- Tokyo regional center is connected to SINET6.
- Bandwidth was **upgraded to 100 Gbps** in January this year.

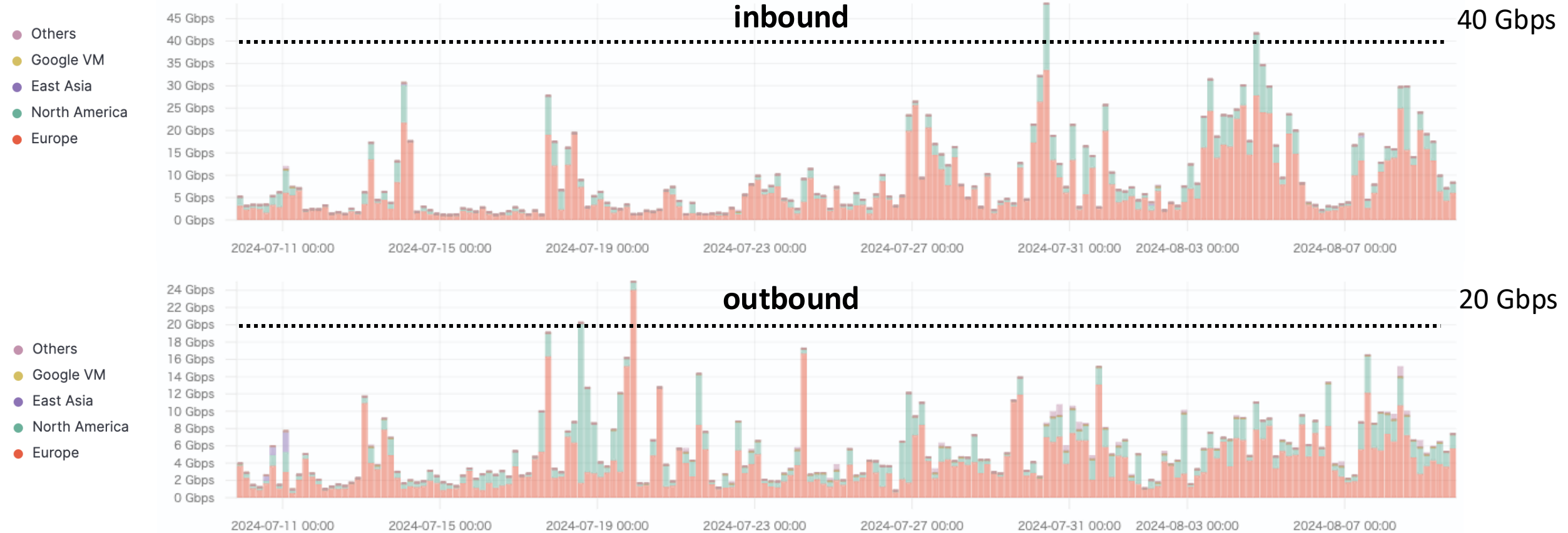
## SINET international connections

- Four >100G international lines to/from Tokyo
  - Tokyo - Amsterdam: 100G x 4
  - Tokyo - Los Angeles: 100G x 2
  - Tokyo - Singapore: 100G
  - Tokyo - Guam: 100G
- Geographical route of “Tokyo - Amsterdam” line was changed in April this year
  - **Larger bandwidth** and **higher latency** than before



# Network: Data transfers per regions

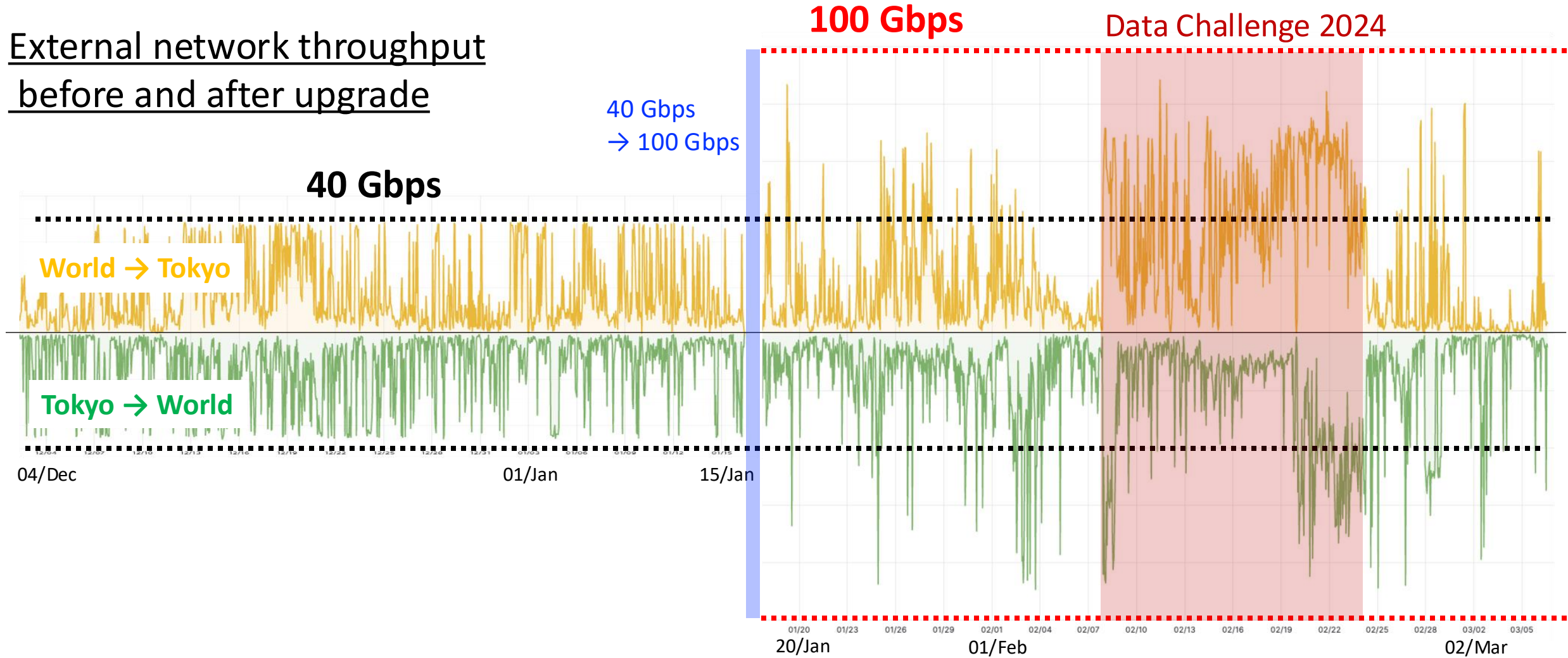
(dCache file servers ↔ LHCONE/Internet)



- Data transfer volume: in(out)bound 40 (35) PB / year → **~200 TB / day**
  - 4PB data transfer as data challenge 2024.
- Dominant transfer region is Europe, followed by North America.
- 40 Gbps bottleneck is now gone with the upgrade.

# Network changes (1): bandwidth (40G to 100G)

External network throughput  
before and after upgrade



- Previously often saturated at 40Gbps. After upgrade, >40Gbps throughput is observed.
- Data Challenge 2024 was a good demonstrator (report later)

# Network changes (2): Tokyo-Amsterdam line

Before Mar 2024



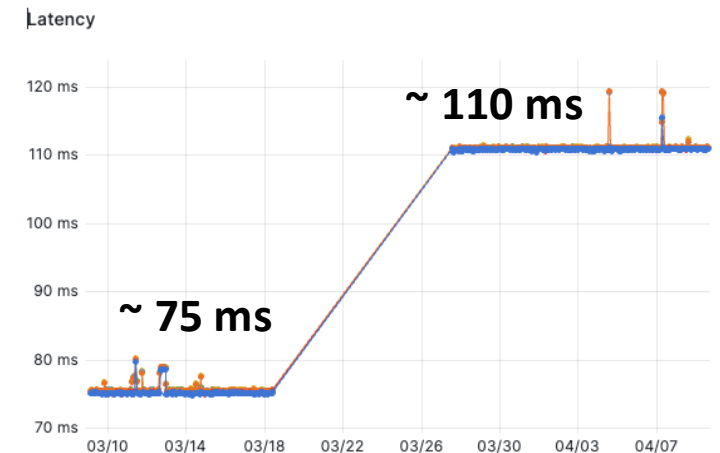
After Apr 2024



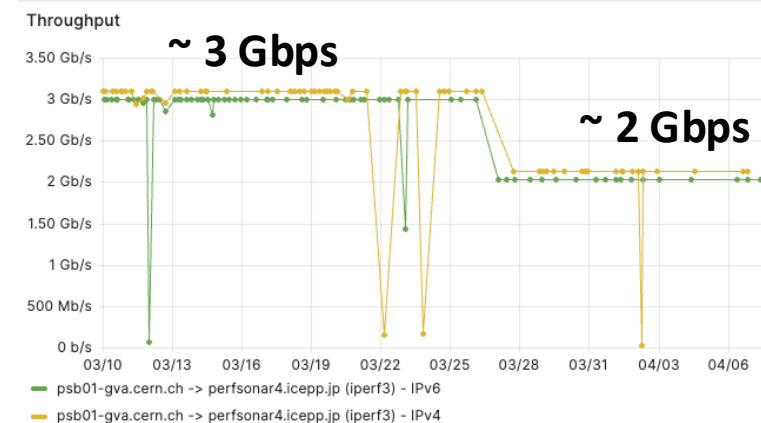
from [SINET6 webpage](#)

- Geographical cable routing has been changed
- Bandwidth improved: 100G to 100G x 4
- Latency increased: 150ms to 220ms
- No critical issues were observed with this change.

perfSONAR latency (Tokyo → Amsterdam)

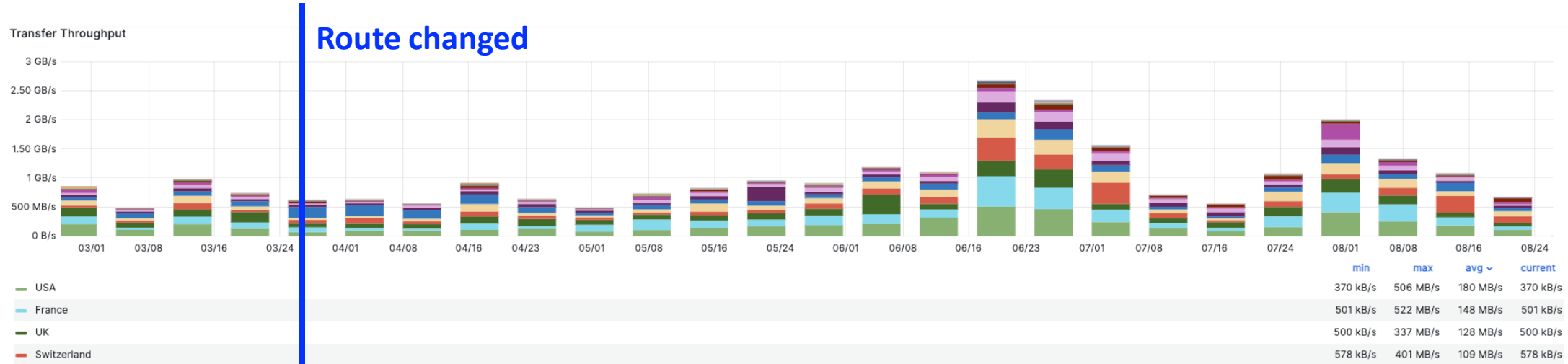


perfSONAR throughput (CERN → Tokyo)

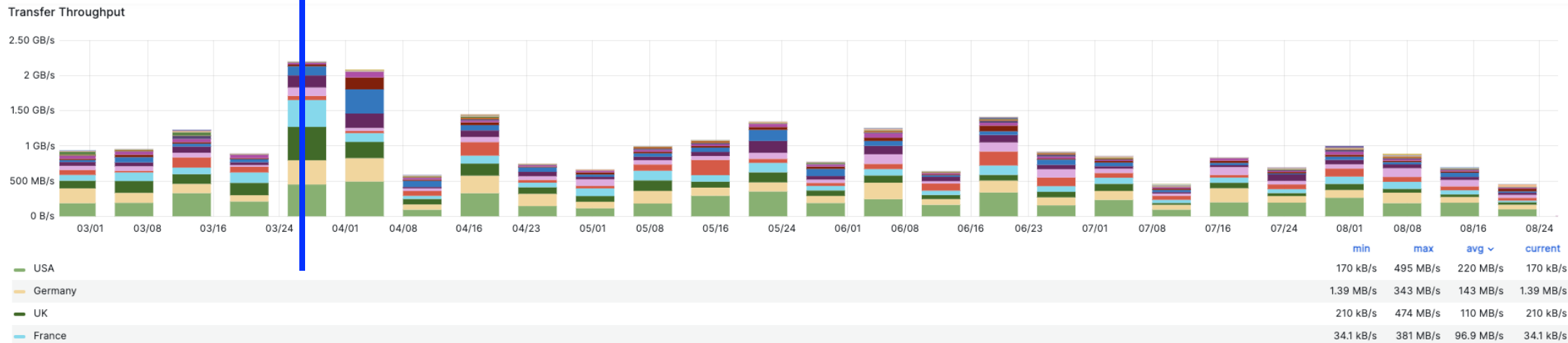


# Network: Throughput to/from Tokyo T2

World → Tokyo



Tokyo → World



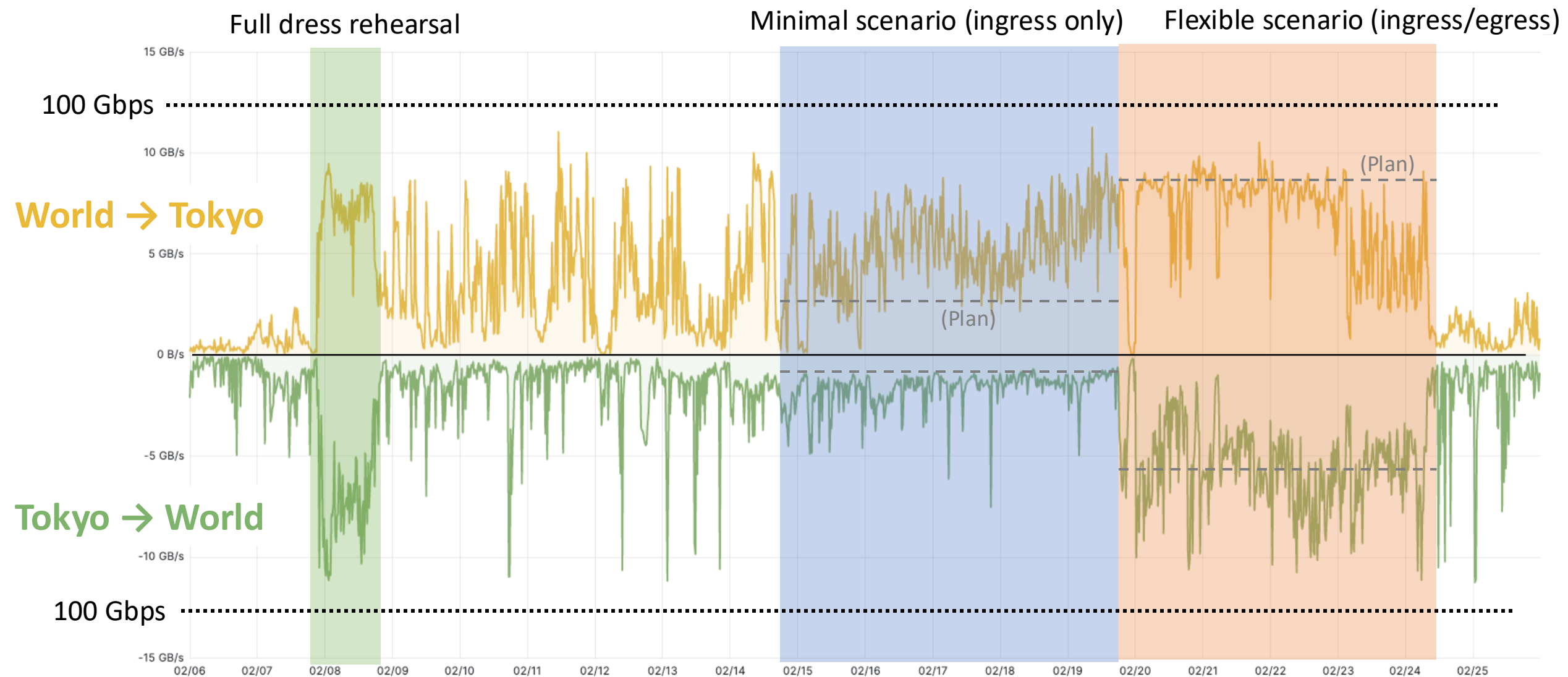
There was no visible reduction in throughput.



# Data Challenge 2024: overview

- Test of 25% of HL-LHC network transfer
  - Aim is to discover the bottleneck/performance/scalability of the network and SE.
- Two (+one) phases and (planned) transfer rates for TOKYO
  0. Full dress rehearsal (7 Feb)
  1. Minimal scenario (12 ~ 18 Feb): Hierarchical (T0 ↔ T1 ↔ T1 ↔ T2)
    - ingress: 21.8 Gbps, egress: 6.1 Gbps
  2. Flexible scenario (19 ~ 23 Feb): Mesh (T0 ↔ T1 ↔ T1 ↔ T2 ↔ T2 ↔ T0)
    - ingress: 64.0 Gbps, egress: 44.1 Gbps
- (These data transfer were executed in parallel with the production transfer)

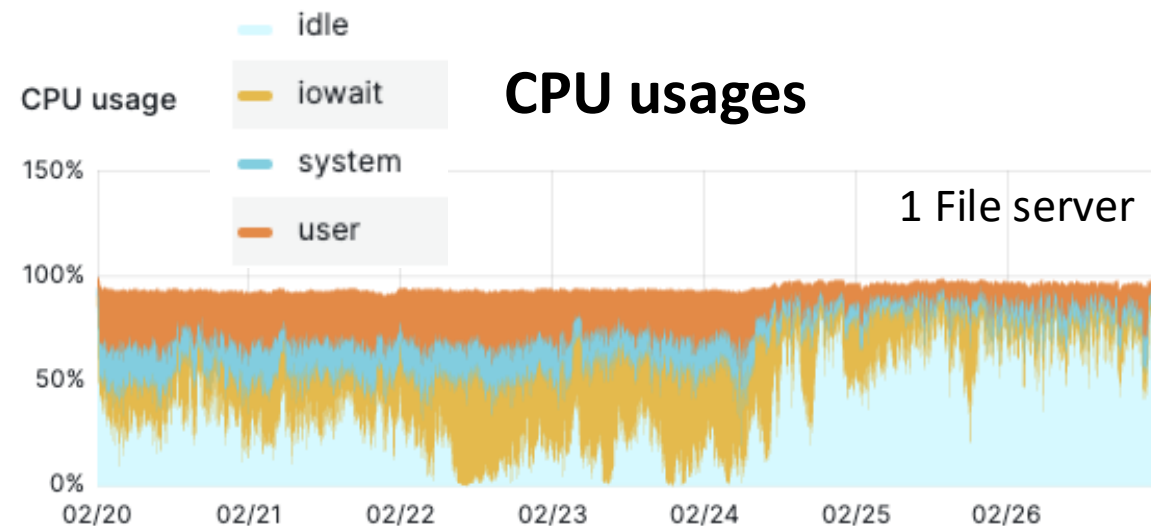
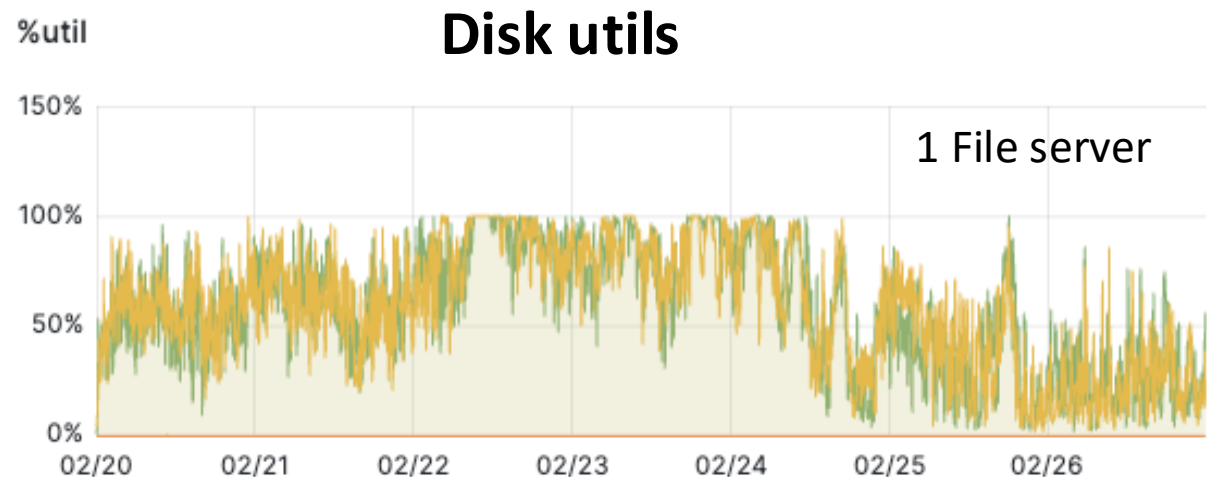
# Data Challenge 2024: Throughput at Tokyo RC



- Successfully operated our SE system at ~100Gbps for O(weeks).

# Data Challenge 2024: File server performance

during DC2024  
read: ~ 200 MB/s, 2500 io/s

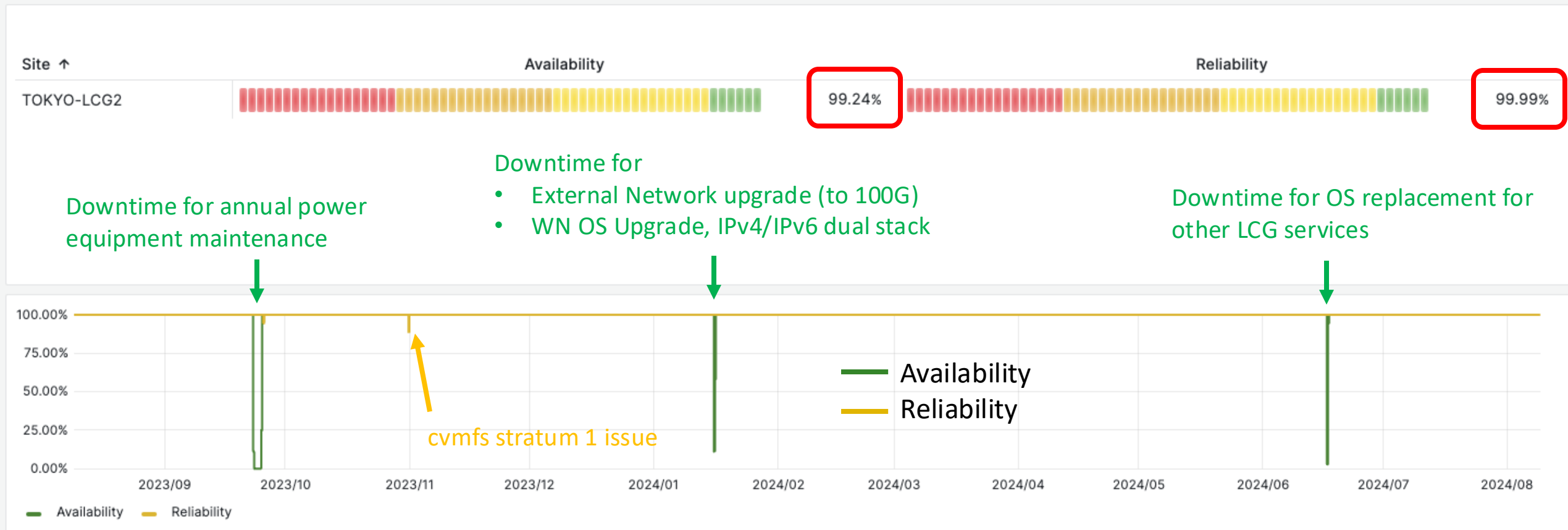


- File server performance is reaching its limit.
- Disk utils reached 100% and a large fraction of IO wait was observed during DC2024
  - due to IO intensive user jobs, as well as DC2024 data transfer
- Improving file server performance is a top priority for HL-LHC.

# **Tier2 operation and Grid services**

# Tie2 operation: stable and reliable

Availability & Reliability (last 1 year)



2023-09

2024-08

**High availability (~99%) and reliability (~99.99%) operation**

# Grid services: overview



**Data**

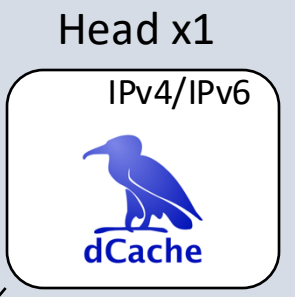
APEL

**Job**



## TOKYO site

### Storage Element (dCache)



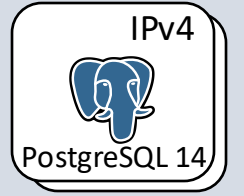
FS x24



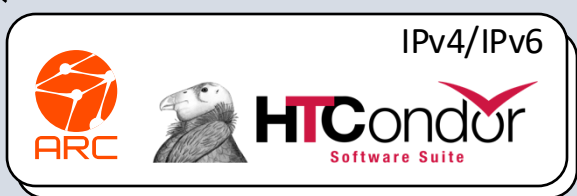
Disk array x48



Database x2

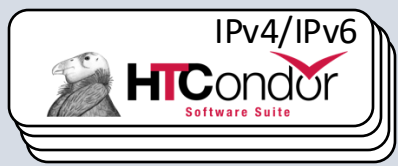


Head x2

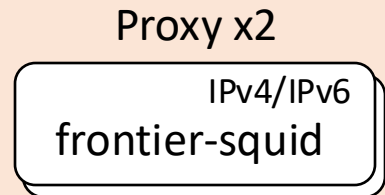


### Computing Element (ARC-CE, HTCondor)

Worker x224



cvmfs & DB



Perfsonar x2



- ~10 head nodes
- 224 worker nodes
- 24 file servers
- 48 disk array

# Grid services: overview (recent changes)



Data

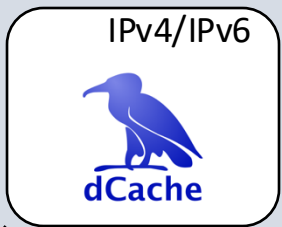
Job



## TOKYO site

Storage Element (dCache)

Head x1



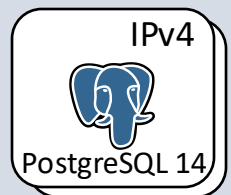
FS x24



Disk array x48



Database x2



APEL

IPv4 BDII

IPv4 Argus

Head x2



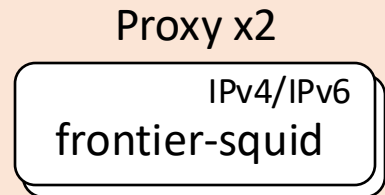
Worker x224



Computing Element (ARC-CE, HTCondor)

- CentOS7 to Alma9 for all nodes!
- IPv4/IPv6 dual stack for CE/WN
- Decommission of argus and top-bdii

cvmfs & DB



Perfsonar x2



# Grid services: recent changes

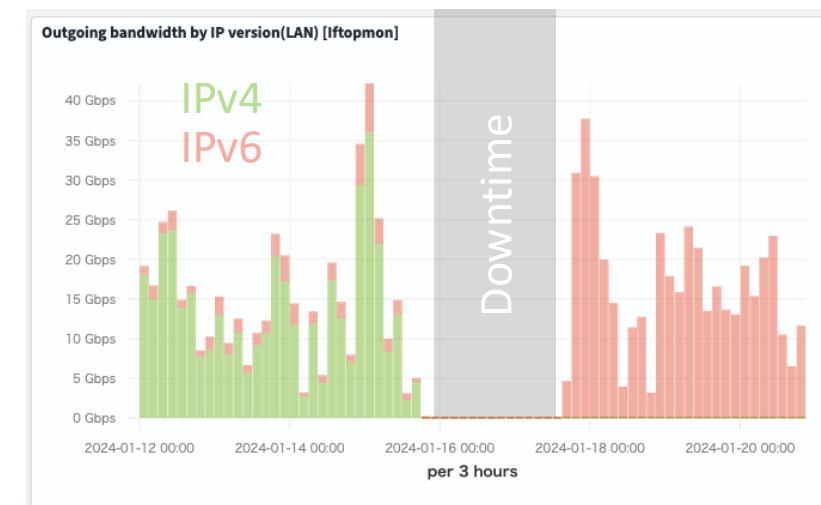
## OS migration (CentOS7 to Alma9)

- CentOS7 reached EOL in June 2024.
- Tokyo site has already completed the OS migration. All grid services have been moved to **Alma 9**.
  - Half day downtime in January 2024 for worker node
  - Half day downtime in June 2024 for other services (CE, SE, BDII, etc.)
  - Cleaned up legacy services/configurations:
    - Argus and top-BDII have been decommissioned.
  - No major troubles during migration.

## IPv4/IPv6 dual stack

- Added IPv6 addresses to CE and WN in January this year.
- Completed IPv4/IPv6 dual stacks for all public T2 services
  - except BDII

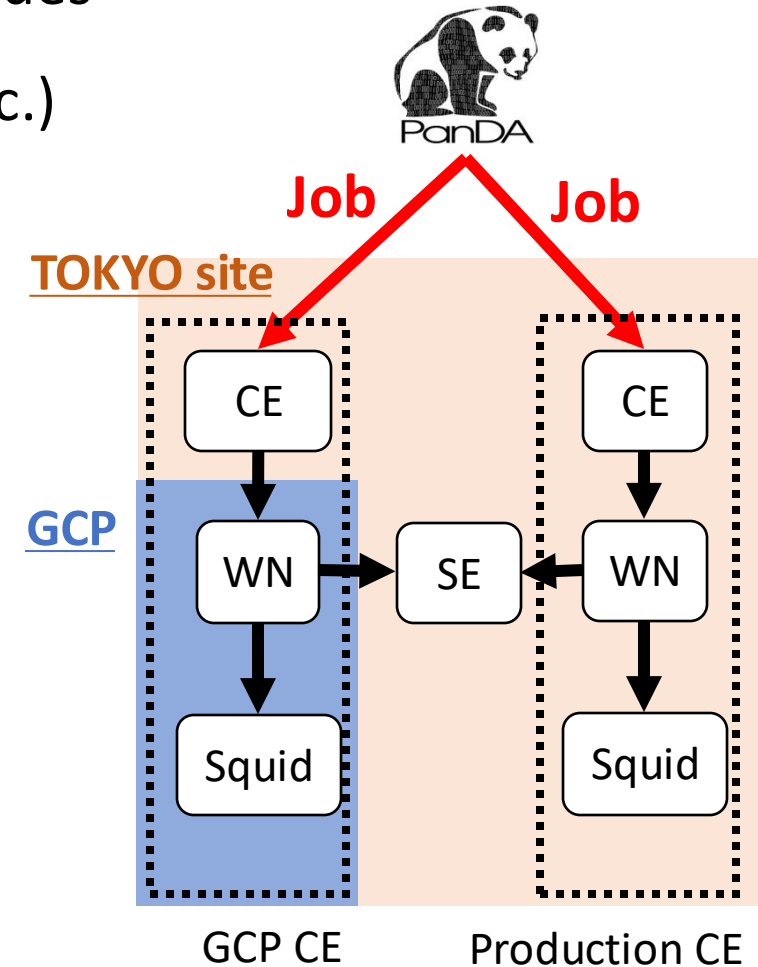
IP versions to file server from local nodes





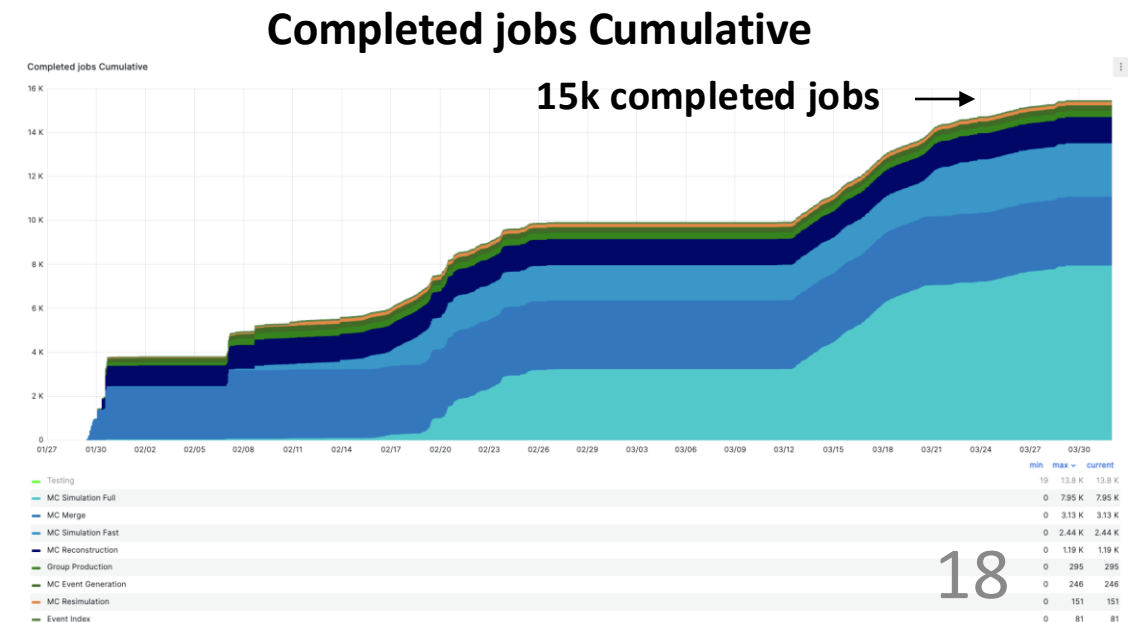
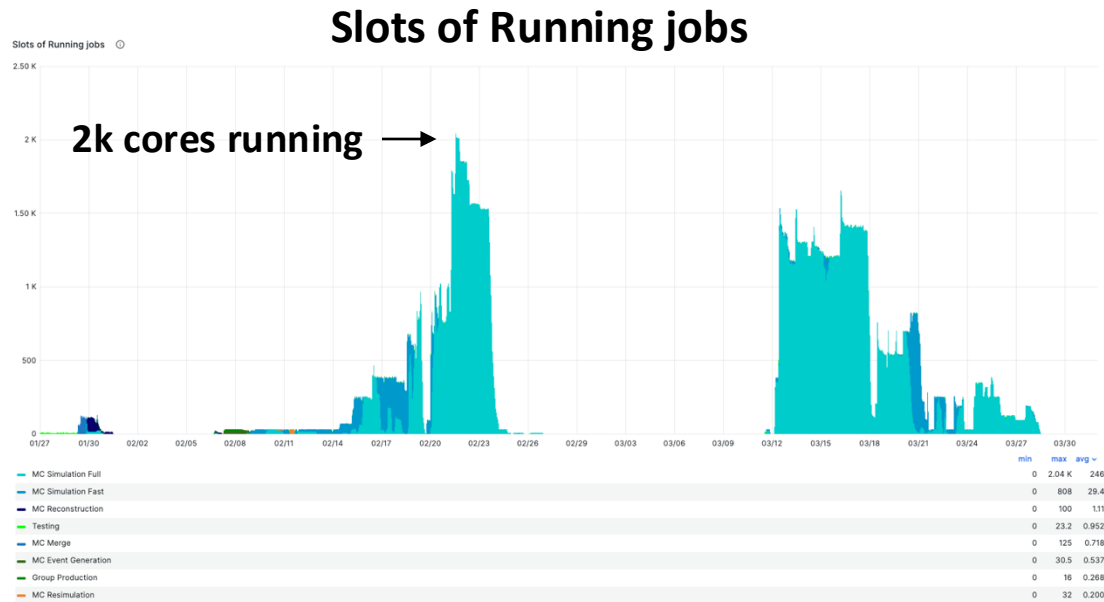
# Misc: Google Cloud Platform

- Started R&D to use Google Cloud Platform (GCP) for worker nodes
  - to use custom/specialized nodes (high mem, ARM, GPU, etc.)
- Configuration
  - A new CE node was created in Tokyo RC domain.
  - WNs and frontier-squid were created in GCP Tokyo region
    - Worker nodes were created as preemptible nodes, while squids were created as non-preemptible nodes.
    - GCP nodes can be easily cloned using a disk template.
  - GCP WNs access to Tokyo T2 SE
    - Network cost of reading the SE from GCP is free, but cost of writing to the SE from GCP is very high.
    - Simulation jobs are a good for GCP CE, because they are CPU intensive.



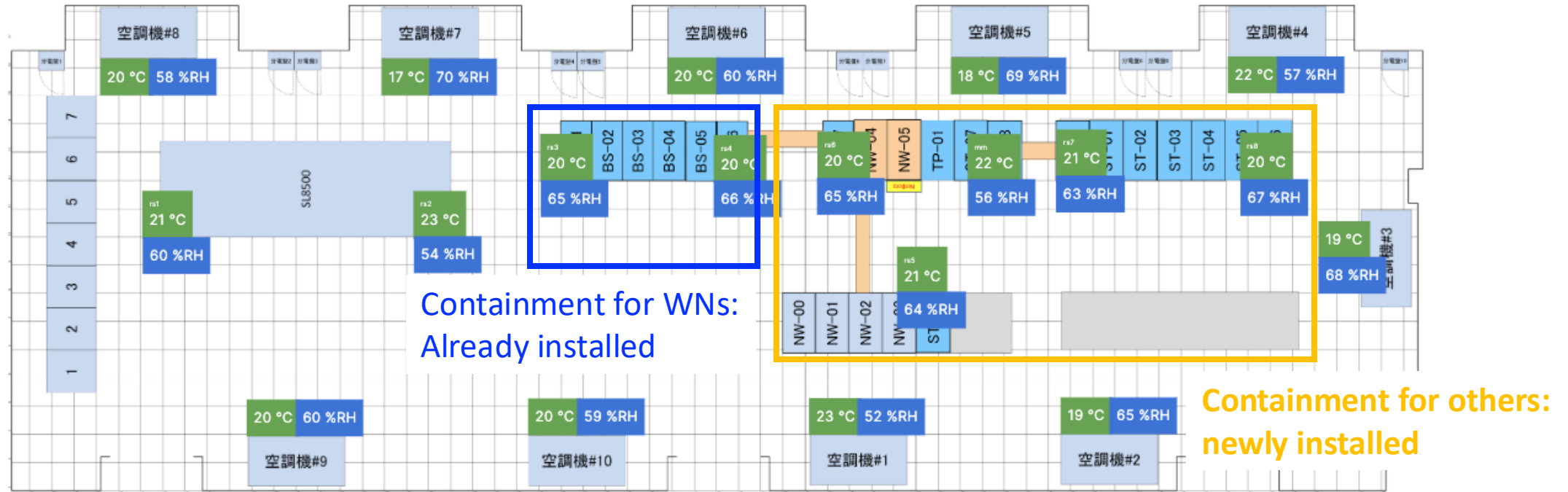
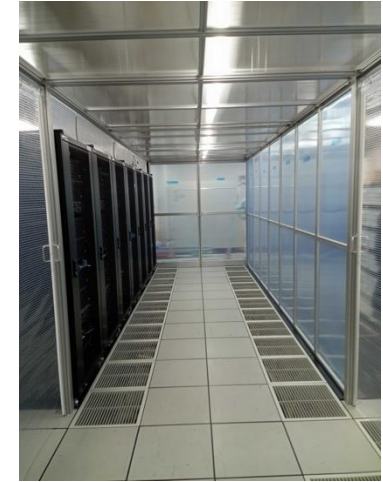
# Misc: Google Cloud Platform: Observation

- Result:
  - Wall time (success jobs): 1.5 Billion sec
  - The number of completed jobs (production jobs): 15 K ~O(10<sup>-3</sup>) of the Tokyo T2 1y resources
- Scalability test up to 2000 worker nodes.
  - Preemptible nodes are terminated very often depending on region/zone and time.
  - We distributed to 3 zones in the same region (Tokyo). But need to distribute more, if using more worker nodes and avoid sudden termination of all worker nodes.



# Misc: Room temperature control

- Built a cold aisle containment for storage/network racks.
- The temperature in Tokyo is getting higher, and our air conditioners are approaching their performance/life limit. It's planned to replace them gradually.



# Summary and plan

- ICEPP regional analysis center is operating stably.
- Contributes to ~4% CPU and ~3% Disk of ATLAS sites
- Network upgraded
  - 100G external network
  - SINET's international network route changes: higher bandwidth, higher latency
- OS of all Tier2 nodes has been successfully migrated to Alma9.
- Next hardware replacement
  - Tokyo site's hardware has been leased and replaced every 3 years so far. The next one was scheduled for Jan 2025.
  - Due to recent increases in delivery times, we will change our previous strategy. The next hardware replacement will be at least one year later. Until then we will continue to use the current hardware.

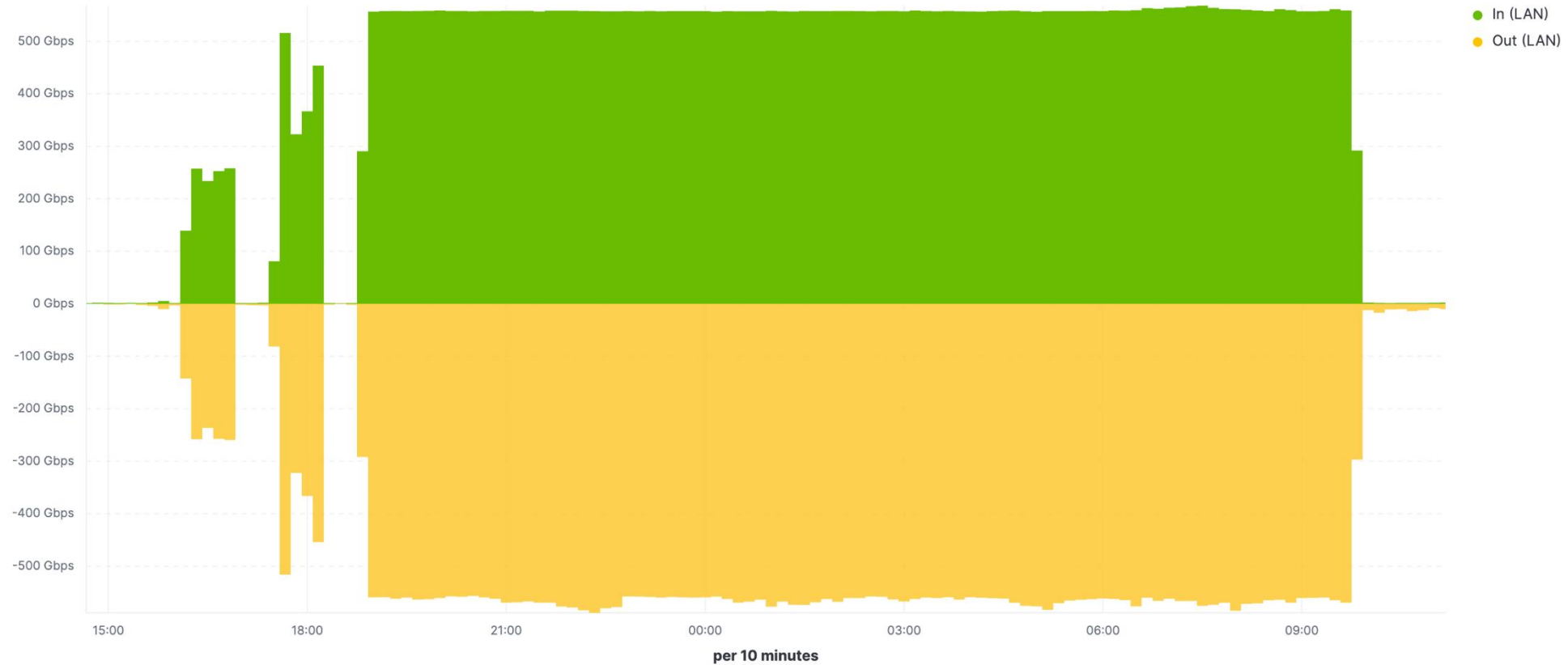
# Backup

# The 5<sup>th</sup> system vs the 6<sup>th</sup> system

		Total	For Tier2
CPU	5 <sup>th</sup> system	336 nodes, 10752 cores (16 cores / CPU) Intel Xeon Gold 6130 2.10 GHz (Skylake) 204 kHS06 1.2 TB HDD x2 / node	240 nodes, 7680 cores 18.97 HS06 / core 3.0 GB RAM / core
	6 <sup>th</sup> system	304 nodes, 15808 cores (26 cores / CPU) Intel Xeon Gold 5320 2.2 GHz (Icelake) 337 kHS06 1.92 TB SSD / node	224 nodes, 11648 cores 21.34 HS06 / core 2.5 GB RAM / core
Disk storage	5 <sup>th</sup> system	72 disk arrays, RAID6 15,840 TB (10TB / HDD)	48 disk arrays, RAID6 10,560 TB (10TB / HDD)
	6 <sup>th</sup> system	72 disk arrays, RAID6 22,176 TB (14 TB / HDD)	48 disk arrays, RAID6 14,784 TB (14 TB / HDD)

# Network (LAN)

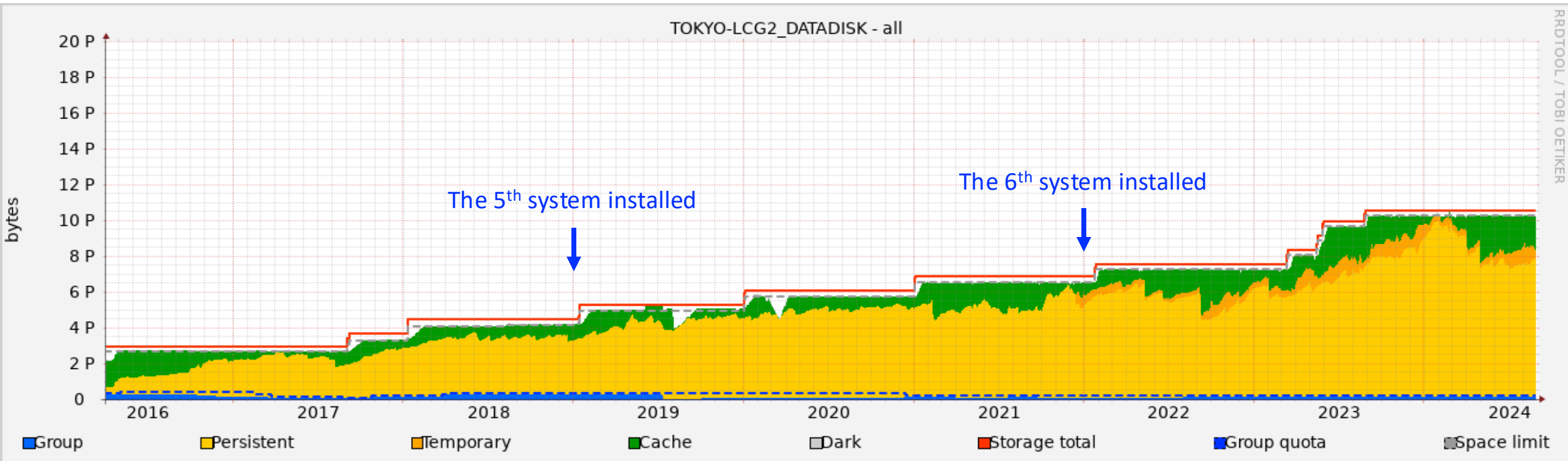
File servers ↔ File servers



- Performed network stress test between file servers during system migration phase
- Observed: ~560 Gbps. (Ideal bandwidth: 25 Gbps x 24 (file servers) = 600 Gbps)

# Storage element (SE)

## Storage volume provided for ATLAS DATADISK



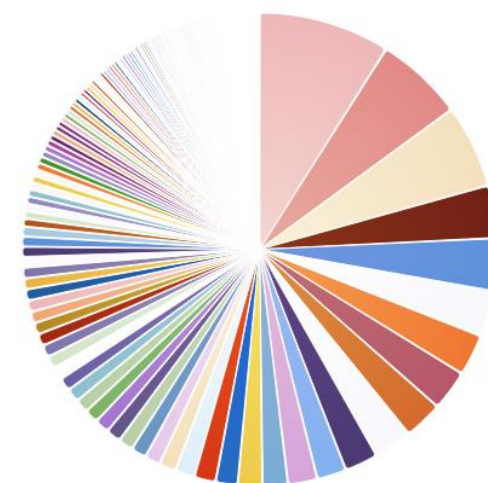
- Increase provided storage volume year by year
- Almost all of the quota are used



# Provided resources for ATLAS as Tier2

- One of the biggest Tier2 sites
  - CPU: ~4% of total ATLAS resources
  - Disk: ~3% of total ATLAS resources
  - cf. ATLAS-Japan member ratio to author list is ~ 3%

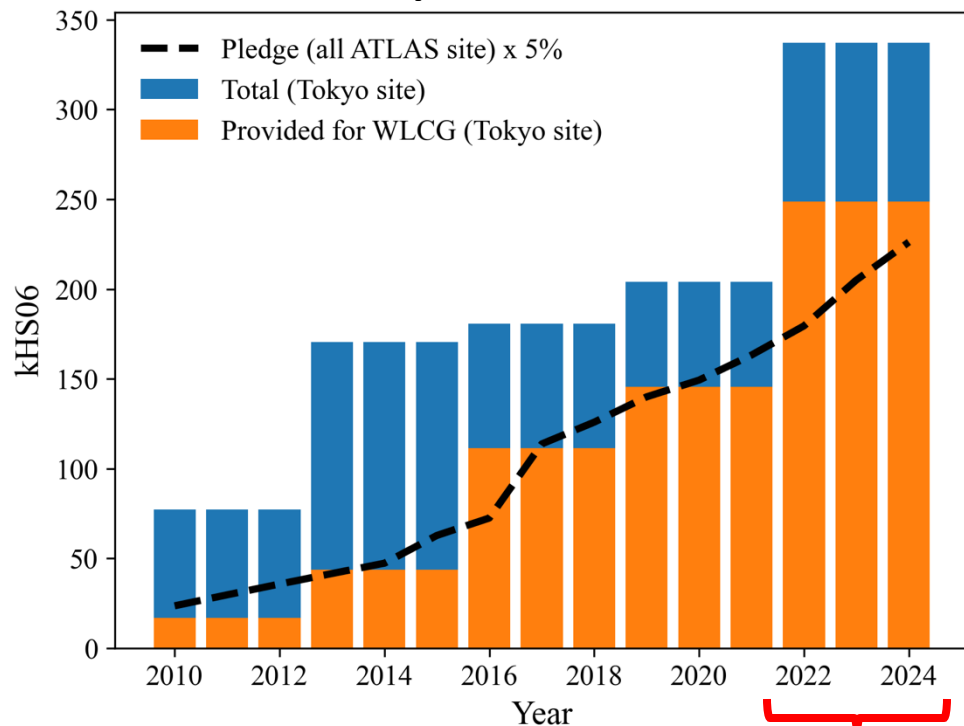
Wall clock time. All jobs (HS23 seconds)



**TOKYO ~ 4% of total ATLAS Grid CPU**

	Percent
MWT2	10%
BNL	7%
CERN-T0	6%
IN2P3-CC	4%
TOKYO	4%
SWT2_CPB	4%

### Compute resources



6<sup>th</sup> system: 2022 - 2024

### Storage resources

