# The WLCG Journey at CSCS: from Piz Daint to Alps

Dr. Riccardo Di Maria (ETH Zurich – CSCS)

HEPiX Spring 2023 Workshop –  ASGC, Taipei, Taiwan
March 28th, 2023

# Alps and Kubernetes at CSCS

***Disclaimer:***



WORK IN PROGRESS

# The *Swiss National Supercomputing Centre,* located in Lugano, is a unit of the *Swiss Federal Institute of Technology* *in Zurich* (ETH Zurich)



*ETH Zurich*



*CSCS Lugano*

CSCS

| 3

ETH*zürich*

# Different infrastructure, different workloads, and different requirements

## The challenge of multiple customers

- **Different Infrastructure**
  - Flagship - CPU/GPU
  - Clusters - Customer Specific
    - WLCG
    - MeteoSwiss
    - CTA and SKA
    - …
  - OpenStack IaaS
  - Experimental Hardware

- **Different Workloads**
  - Classic HPC
    - SSH to login nodes
    - Submit jobs to Slurm
    - Wait for results
    - Repeat
  - Grid Computing
    - WLCG
  - Interactive Computing
    - Jupyter Notebooks
    - Remote Visualization
  - IaaS

*Piz Daint*



ETH zürich

# Alps

**Successor to Piz Daint**



Alps

- Alps at CSCS
    - HPE Cray EX (AMD Rome and Milan, ARM Grace, NVIDIA A100, etc.)
      → Shasta architecture and Slingshot
    - ***Infrastructure as Code***
      → designed from ground up for programmability of resources for workflows
      → multi-tenancy paradigm
      → Slurm/HPC and K8s/Cloud vClusters: persistent, on-demand, and/or elastic
    - Continued support for classic supercomputing use cases
    - Additional support for AI, ML and data-driven workflows
    - Phased installation/expansion (10-15% March 2023 == ~1200 nodes)

cscs

**ETH** *zürich*

# Virtual/Versatile/Volatile Cluster Configuration at CSCS

# WLCG @ CSCS

## Tier-2 for ATLAS, CMS, and LHCb under CHiPP Federation

**2022**
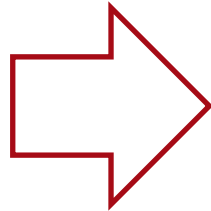
- ATLAS
  - 89 kHS06
  - 3.7 PB
- CMS
  - 77 kHS06
  - 2.8 PB
- LHCb
  - 56 kHS06
  - 2.5 PB

**2023**

- ATLAS
  - 112 kHS06
  - 4.4 PB
- CMS
  - 92 kHS06
  - 3.4 PB
- LHCb
  - 70 kHS06
  - 3.0 PB

- ❖ ~15 PB dCache on Ceph
- ❖ 100 AMD EPYC Rome nodes
  - 128 cores (256 CPUs), 256 GB RAM
  - "Mont Fort" cluster
  - 4 ARC-CEs
- ❖ +4 nodes for dev/tds instance
  - "Mont Gele" cluster, 1 ARC-CE
- ❖ Production CE
  - 300 TB shared CephFS NVMe
  - 4 TB local RBD NVMe per node
  - 64 GB CVMFS cache RBD NVMe per node

cscs

ETH zürich

# Kubernetes at CSCS (v1.0)

- ## Kubernetes Clusters at CSCS
  - shared internal CSCS-managed services (Fulen)
  - shared external user-managed (Combin)
  - dedicated for specific needs
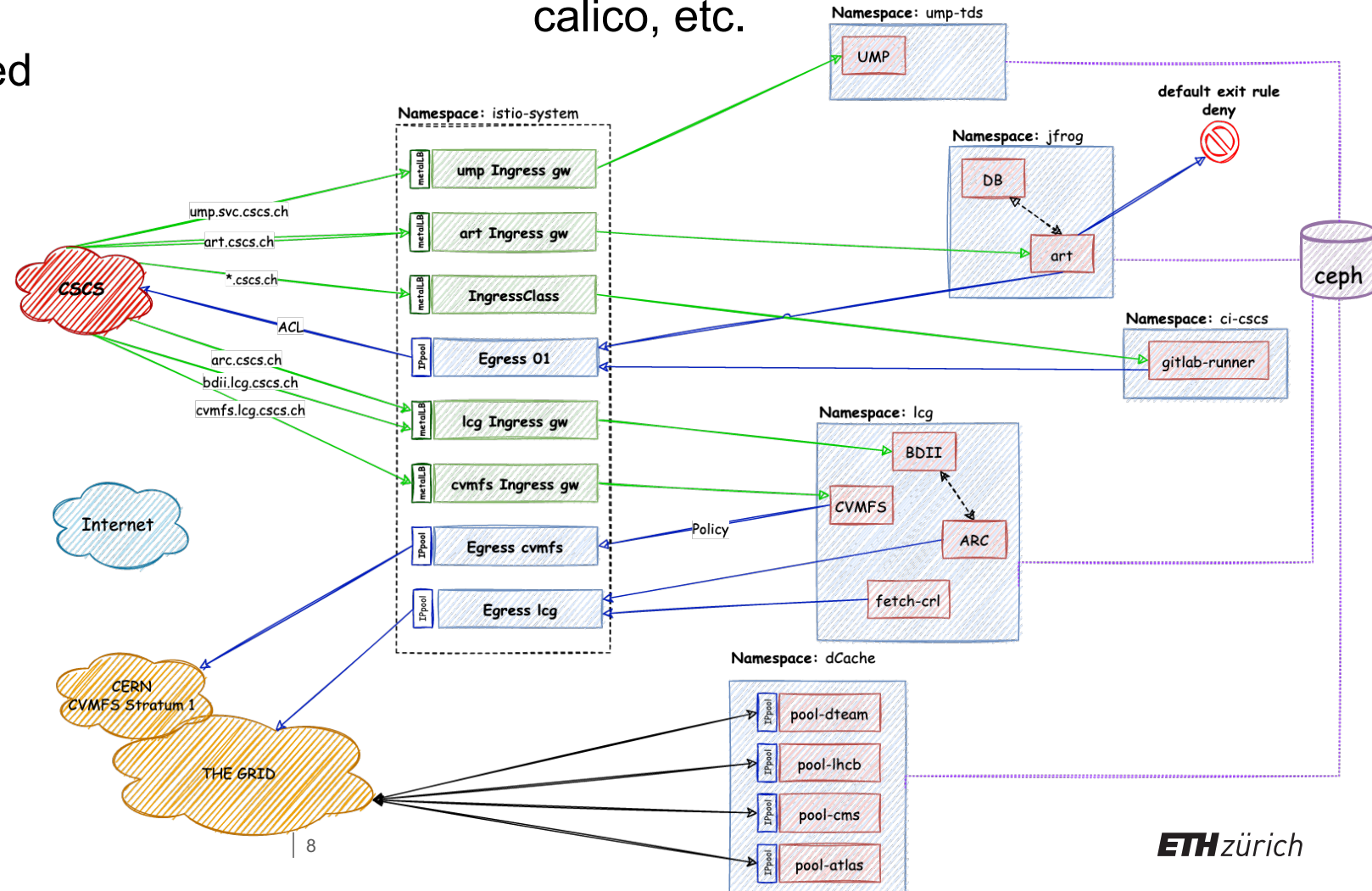
- ## Based on community "vanilla" Kubernetes

- ## The "Fulen case"
  - **dCache**
  - WLCG Services
    - **ARC-CE**, CVMFS, BDII, VO Boxes, etc
  - …

- ## Key features
  - metalLB, istio, cert manager, OIDC, calico, etc.
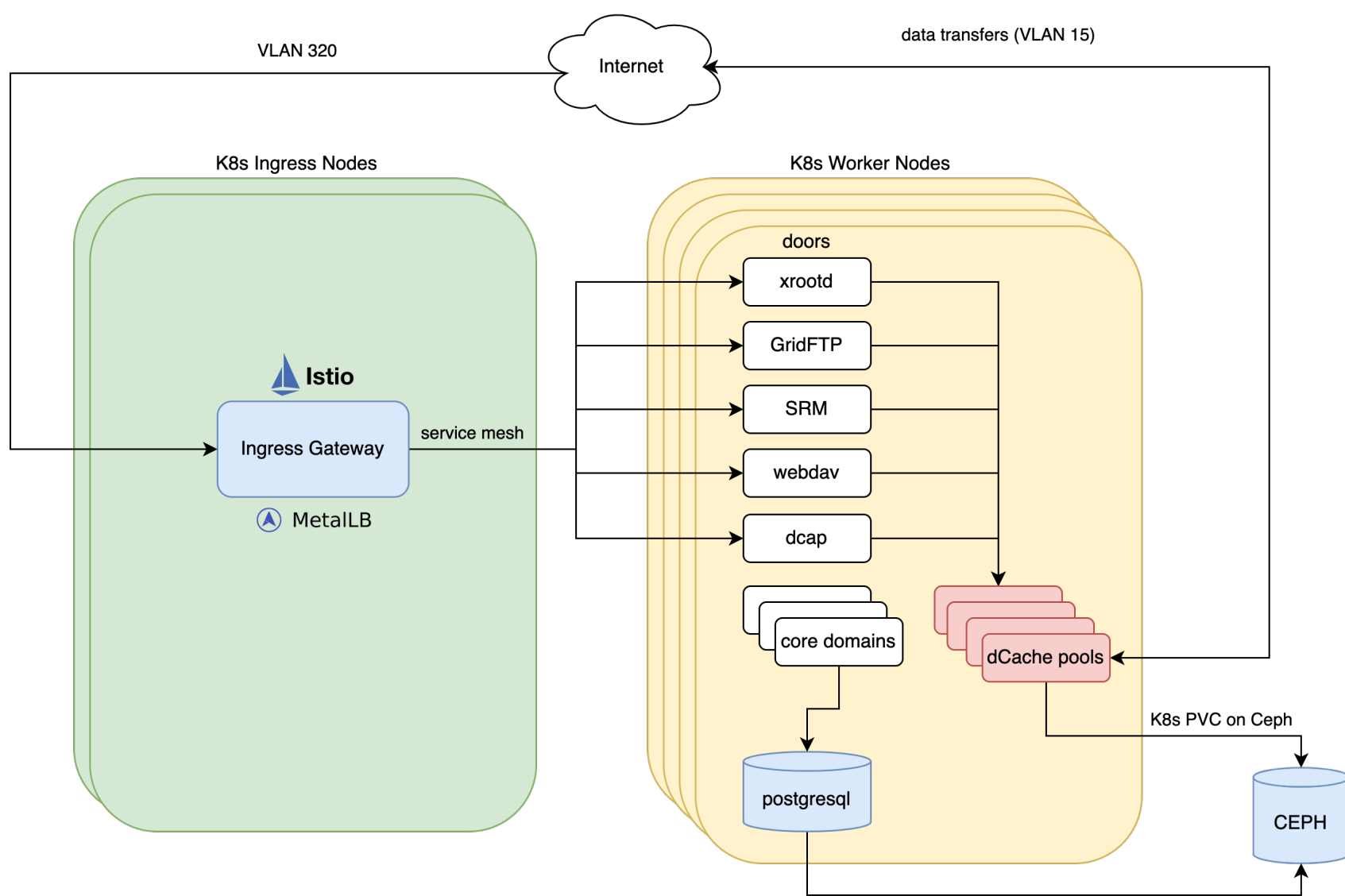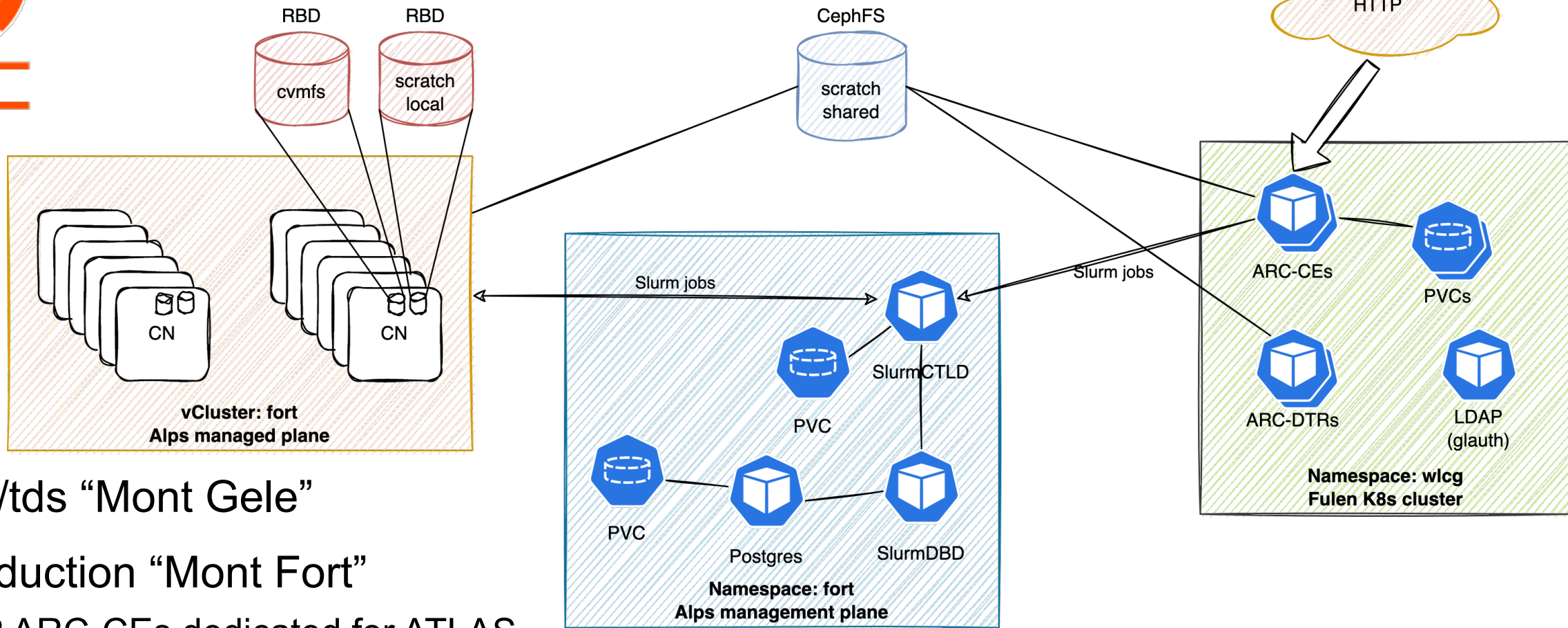
cscs

ETH zürich

# on Kubernetes



- K8s came after WLCG and CTA requirements were set

- ~1 year in production

- dCache pool services run as K8s pods

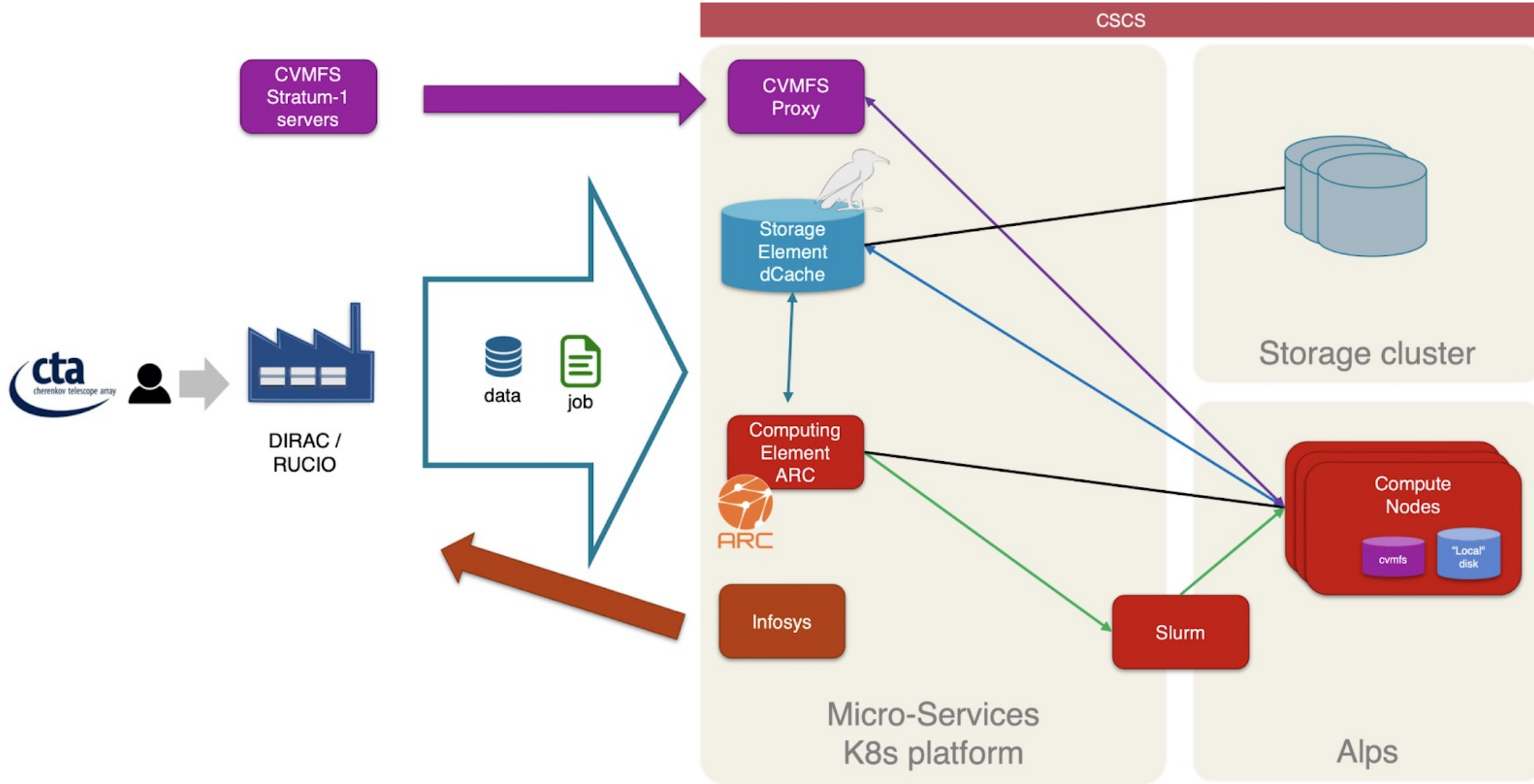- Pods mount Ceph RBD volumes through Kubernetes CSI
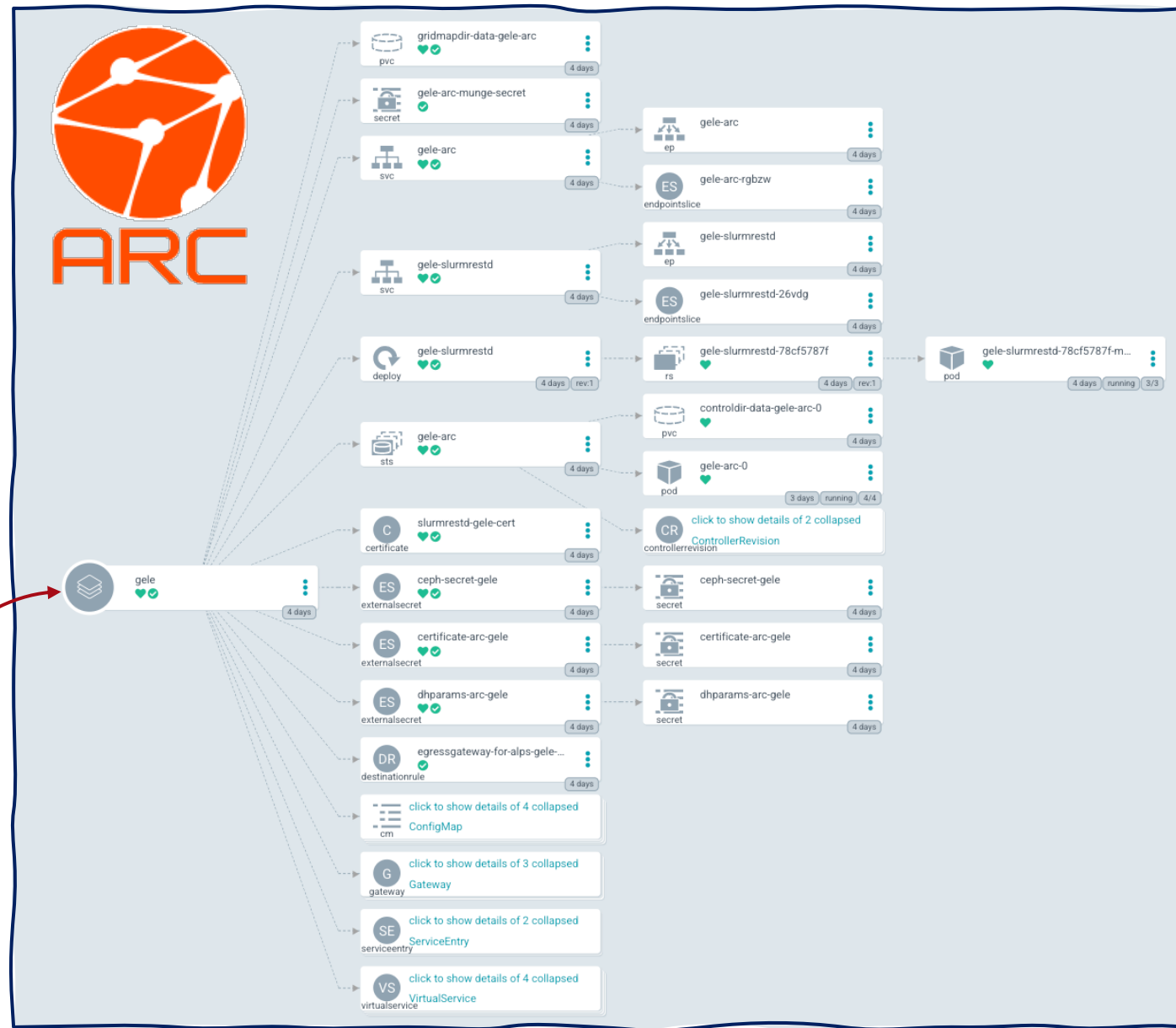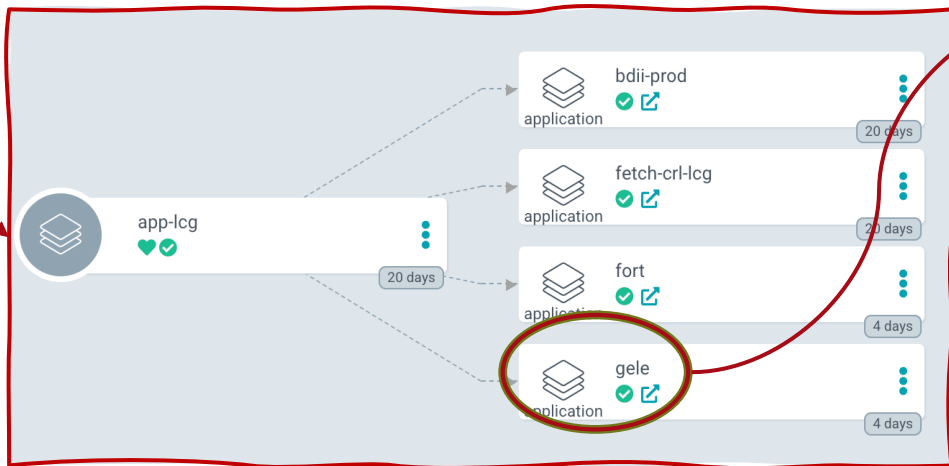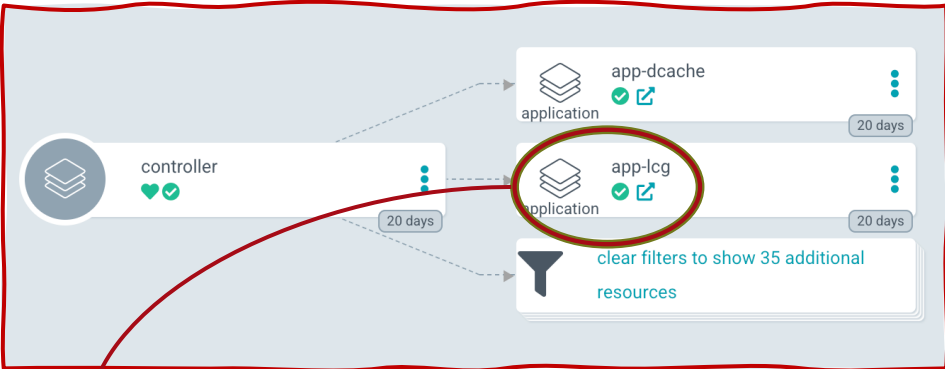
# on Kubernetes

- dev/tds "Mont Gele"

- Production "Mont Fort"
  - 2 ARC-CEs dedicated for ATLAS
  - 2 ARC-Ces dedicated for CMS and LHCb
  - Common Slurm queue

# WLCG and CTA Workflows at CSCS

# The Fulen Cluster and ArgoCD

## CI/CD for ARC-CE on Kubernetes



*(mmh, not quite!)* | 12

# Kubernetes at CSCS (v2.0)

**RANCHER®**
BY SUSE

- On-demand K8s clusters for clients and customers with different needs and requirements

- K3S/RKE2 spawned clusters with tagged VLAN isolation
  → improved istio management

- ArgoCD for cluster configuration and/or application deployment

- Cilium as K8s CNI

## Baremetal — RKE
- specific needs e.g. SE, monitoring
- computing power
- local storage
- dedicated VLAN
- deployed via MaaS

## Virtual — HARVESTER
- virtual resources
- multiple internal/external VLANs
- RKE2 or K3S
- 100G Ethernet
- local SSD (longhorn)
- external RBD and CephFS

## Alps — RKE
- HPC
- 400G Ethernet
- dedicated VLAN
  (after Slingshot 2.0 upgrade)

K3S

cscs

**ETH** zürich

# Kubernetes at CSCS (v2.0)

- Baremetal
  - e.g. monitoring/ECK → Dino Conciatore *"Dynamic Deployment of Data Collection and Analysis Stacks at CSCS"*, HEPiX 2023
  - on-going WLCG and CTA dCache instances migration
- Alps
  - challenges:
    - cluster persistency and CI/CD
    - admin privileges for customers → Slingshot 2.0 upgrade on-going → dedicated VLANs to be tested
  - PoC/MVP for PSI
- Virtual
  - *quite a few…*

cscs

**ETH** *zürich*

# Kubernetes Multi-Cluster Design

- ## Cluster for client:
  - ### etcd cluster S3-backup
  - ### CSI CephFS and RBD
  - ### velero
  - ### beats
  - ### ingress nginx
  - ### metalLB
  - ### external-DNS
  - ### cert-manager

- ## External-secrets

- ## Vault

- ## ArgoCD

Internet        Firewall        CSCS Network

CSCS shared external

CSCS TDS shared external

**IAM**
SSH-service
Keycloak
Waldur
...

CSCS shared internal

CSCS TDS shared internal

**IAM TDS**
SSH-service
Keycloak
Waldur
...

DWDI dev

vCEF dev

SOLE

**Critical DevOps services**
Vault
JFrog
GitLab
ArgoCD

**WLCG/CTA**
dCache
ARC-CE
CVMFS
...

xyz dev

PSI TDS

Rosa (JupyterHub and BuildFarm)

PSI dev

PSI

CTA

MCH-ML

SKA

CSCS

K8s Users

External

K8s Cluster Managed by CSCS

K8s Cluster Managed by Users
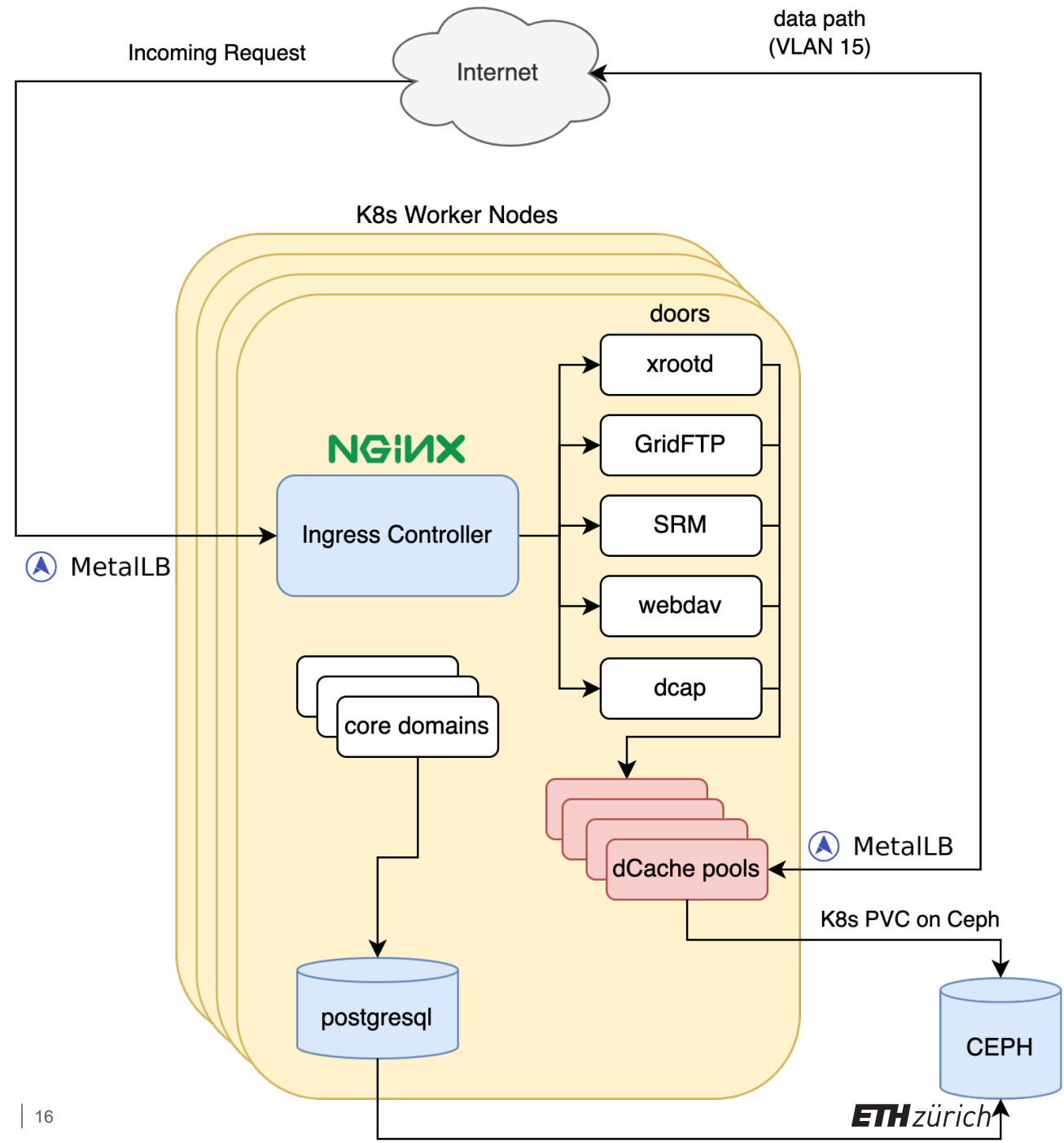
K8s Cluster on Alps

High Throughput Application

cscs

| 15

ETH zürich
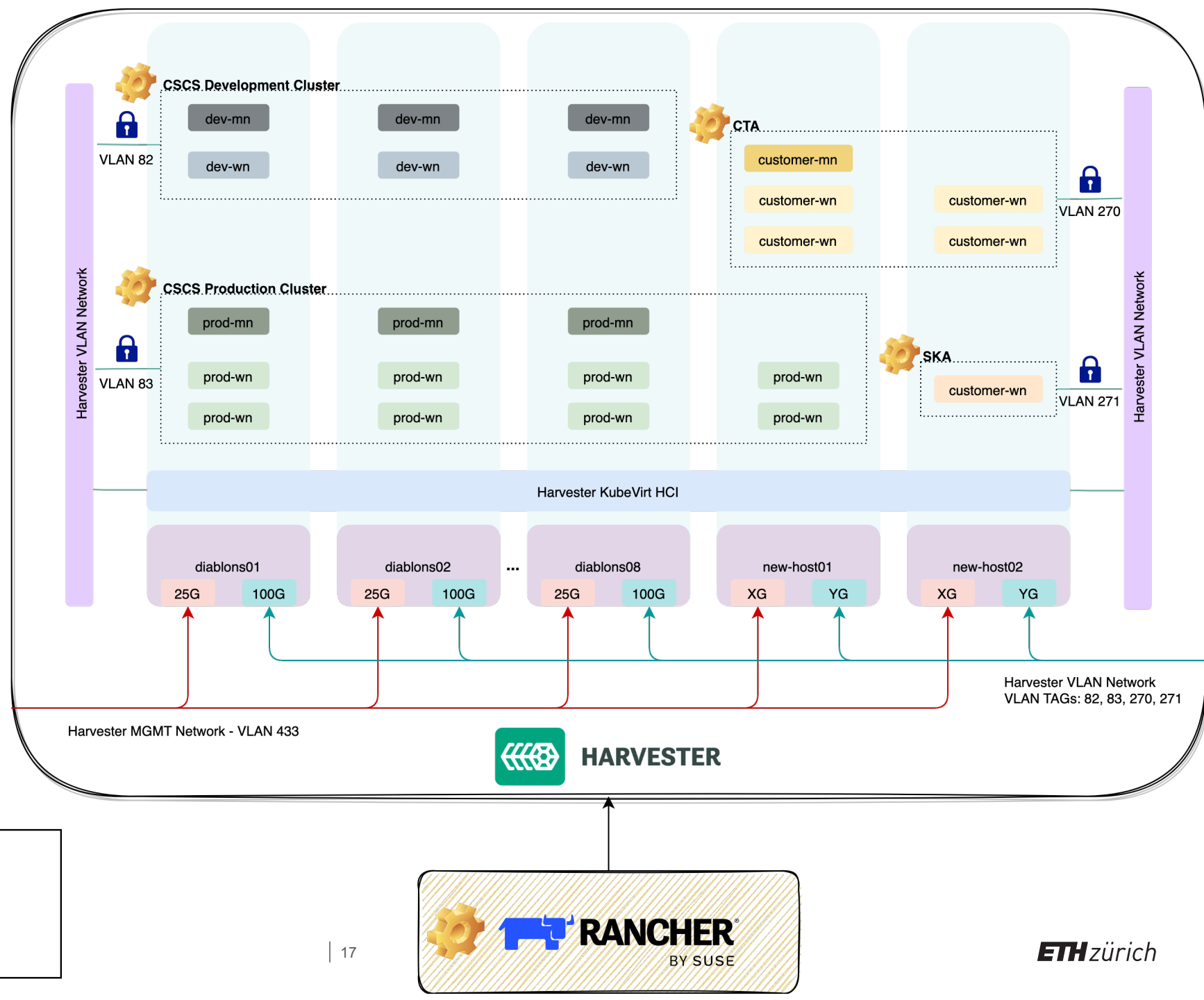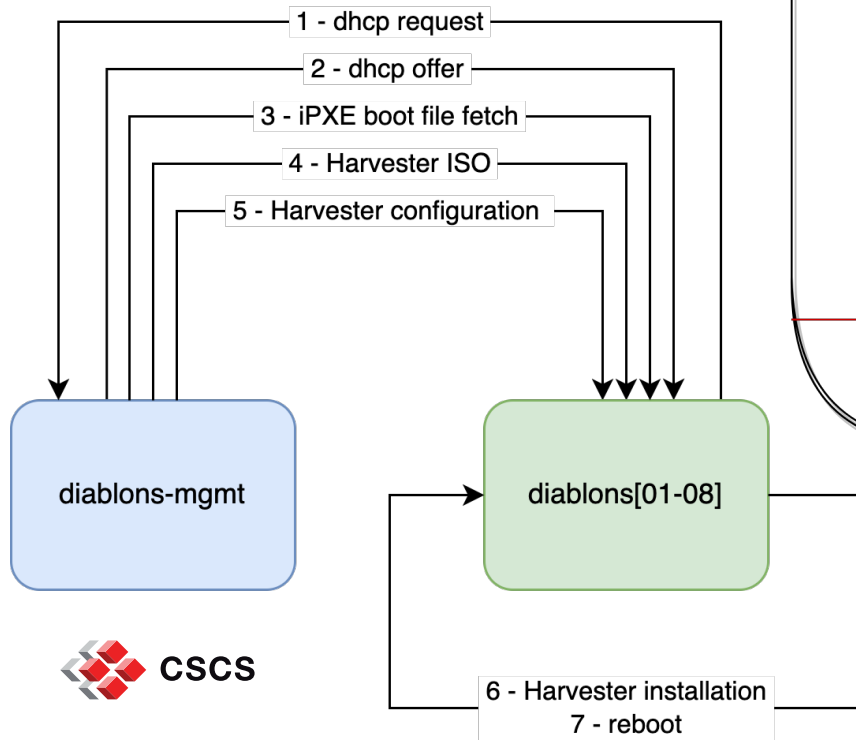
# Harvester at CSCS

- Cluster deployment →

- Harvester deployment:

# Why are we moving services to Kubernetes?

## What's the point of using Kubernetes?

- Main advantages
  - Load balancing
  - Storage orchestration
  - Automated rollouts and rollbacks
  - Automatic bin packing
  - Self-healing
  - Secret and configuration management
  - Observability and traffic management
  - **Disaster recovery management and one-button deployment**
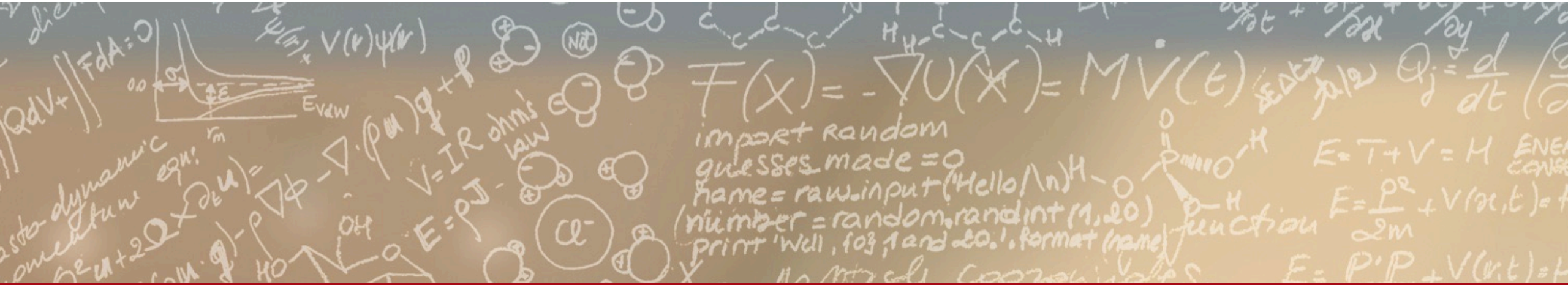
- Main challenges
  - Additional "moving parts" and complexity layer
    - Networking: Cilium vs. Calico, and service mesh
  - Security
    - Additional configuration and additional MAC (mandatory access control) configuration
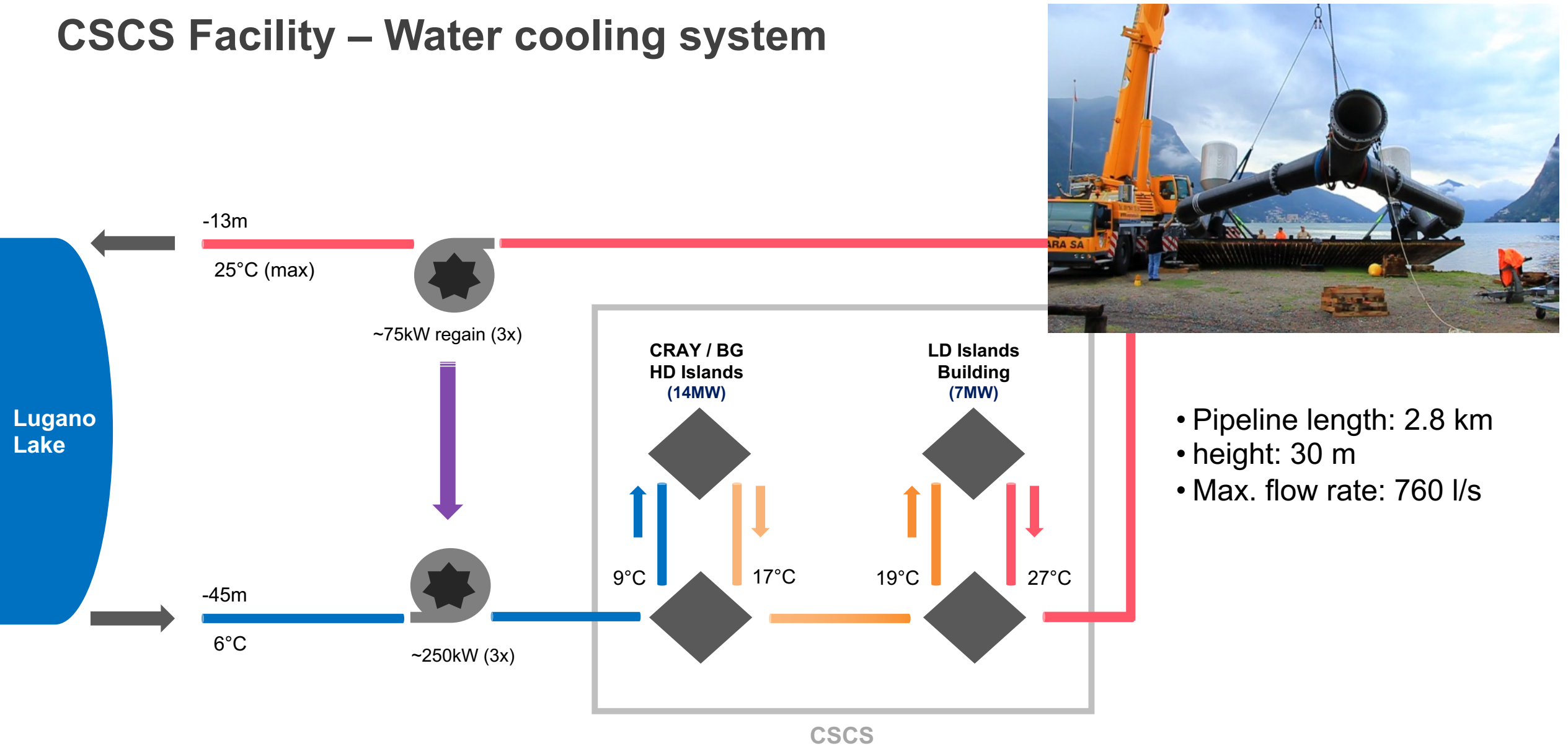
**Thank you for your attention.**

**Questions?**

*Contact: riccardo.dimaria@cscs.ch*

# CSCS Facility – Water cooling system



**Lugano Lake**

-13m
25°C (max)

~75kW regain (3x)

-45m
6°C

~250kW (3x)

**CRAY / BG HD Islands (14MW)**

9°C    17°C

**LD Islands Building (7MW)**

19°C    27°C

CSCS

- Pipeline length: 2.8 km
- height: 30 m
- Max. flow rate: 760 l/s

CSCS

ETH *zürich*

# Ceph at CSCS

- **Existing implementation**
  - NVMe (~300TB)
  - HDD (~11PB)

- **Expansion phase ahead**

- **On-going:**
  - Rucio backend integration with S3



**APP** → **LIBRADOS**

A library allowing apps to directly access RADOS, with support for C, C++, Java, Python, Ruby, and PHP

**APP** → **RADOSGW**

A bucket-based REST gateway, compatible with S3 and Swift

**HOST/VM** → **RBD**

A reliable and fully-distributed block device, with a Linux kernel client and a QEMU/KVM driver

**CLIENT** → **CEPH FS**

A POSIX-compliant distributed file system, with a Linux kernel client and support for FUSE

**RADOS**

A reliable, autonomous, distributed object store comprised of self-healing, self-managing, intelligent storage nodes