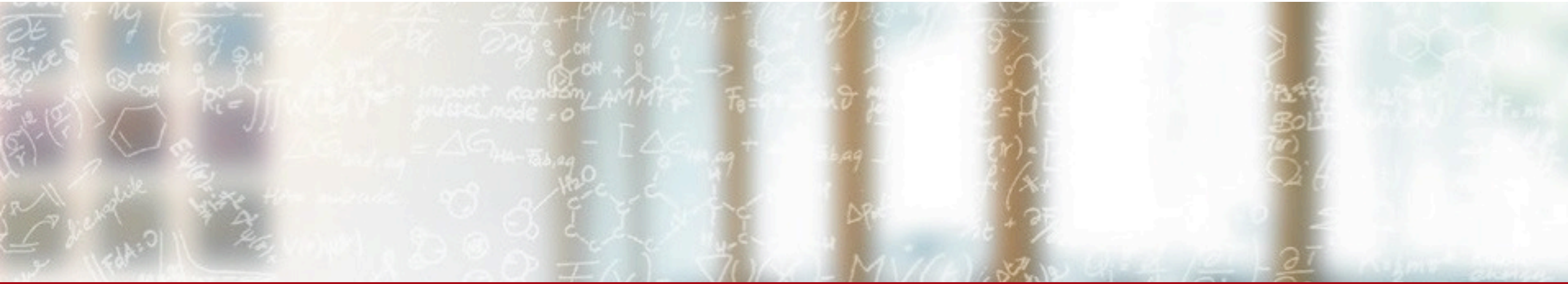




**CSCS**

Centro Svizzero di Calcolo Scientifico  
Swiss National Supercomputing Centre

**ETH** zürich



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

Dr. Riccardo Di Maria (ETH Zurich – CSCS)

Cray User Group 2023 – CSC IT Center for Science Ltd, Helsinki, Finland

May 9<sup>th</sup>, 2023

# Alps and Kubernetes at CSCS

## Disclaimer



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*



# 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*



# Alps

## Successor to Piz Daint



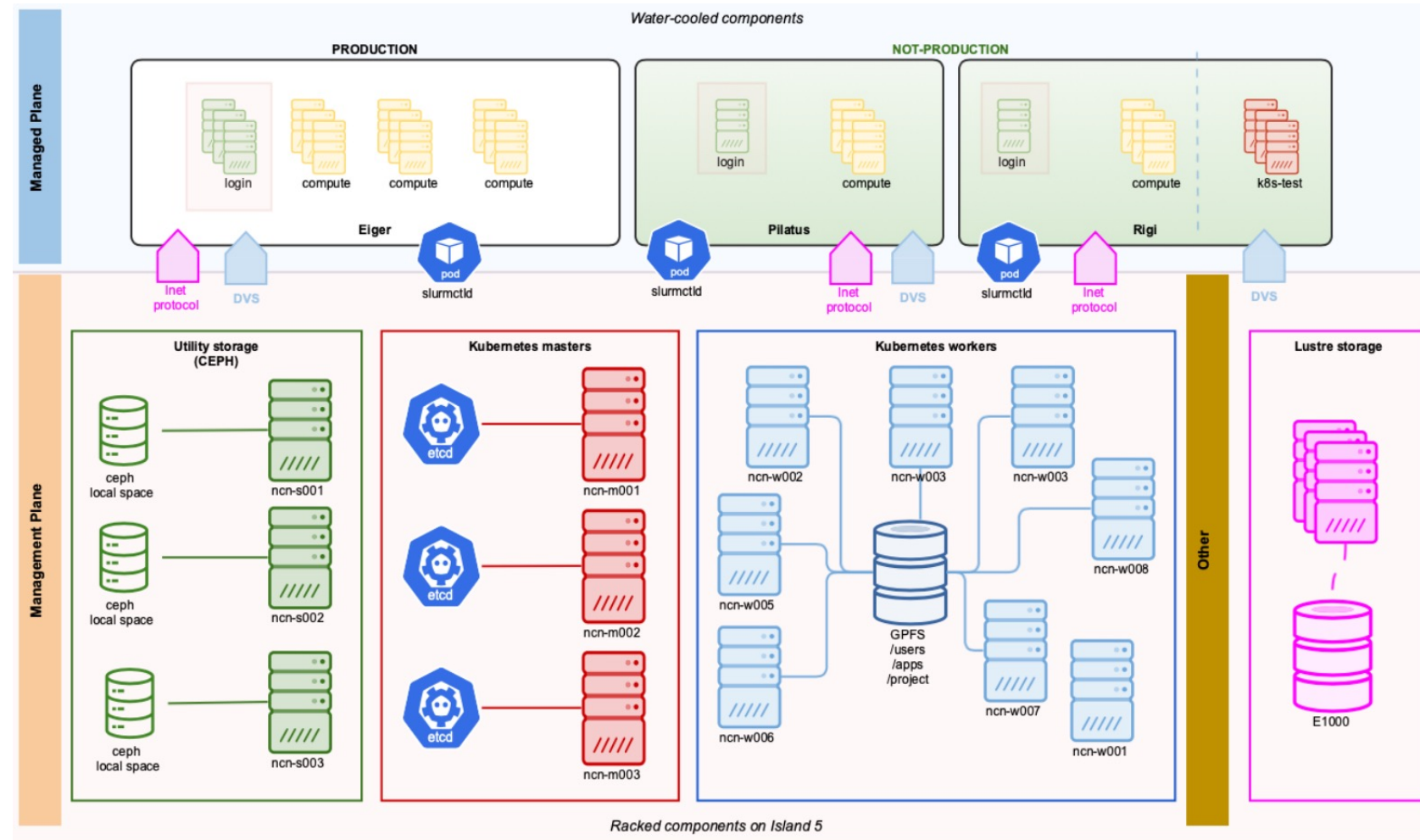
### ■ Alps at CSCS

- 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% April 2023 == ~1200 nodes)

Alps

# vCluster Configuration at CSCS

- vCluster
  - dedicated compute administered by namespaced K8s resources (K8s4CSM)
- Software-defined infrastructure (IaC) and CPE features
  - multiple Slurm instances
- WLCG context:
  - shared + tailored CFS layers
  - no login nodes
  - HTC workflows
    - single and multi-core jobs
    - no MPI
    - no hyper-threading
  - Slurm fine-tuned



# HPC and Kubernetes

- Full service on HPC
  - security challenges
    - VLANs should help → need testing
    - ad-hoc configurations between management and managed plane
  - inefficiency on costly resources
  - additional "virtualisation" layer → complexity (e.g. network)
- Front-end service on external K8s, compute on HPC
  - efficient use of HPC resources
  - necessity of workflow/job scheduler
    - impact on management plane
  - necessity of middleware/interface between customer and compute

# Moving to Kubernetes

- Main advantages
  - Decoupling from the infrastructure
    - Storage with CSI
  - Declarative configuration
  - Reusage of code
  - Load balancing
  - Automated rollouts and rollbacks
  - Self-healing
  - Secret management
  - Observability and traffic management
  - **Disaster recovery management and one-button deployment**
- Main challenges
  - Additional “moving parts” and complexity layers
    - Networking: Cilium vs. Calico, service mesh
  - Security
    - Additional configuration



# kubernetes



# Kubernetes Tools at CSCS

## ■ Rancher (SUSE)

- Kubernetes cluster orchestrator
  - multi-tenancy
  - role-based access control
  - monitoring
- Multi-cloud and bare-metal
  - Deployment process simplified
- Integration with Harvester and VMWare
- Cluster templating
- Security oriented
- K8s cluster using Cilium for CNI
  - leveraging extended Berkeley Packet Filter (eBPF) technology
  - offering transparent visibility and control of network traffic between services, enabling fine-grained policy enforcement and network segmentation

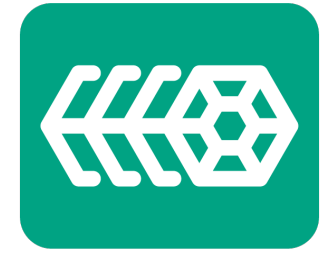
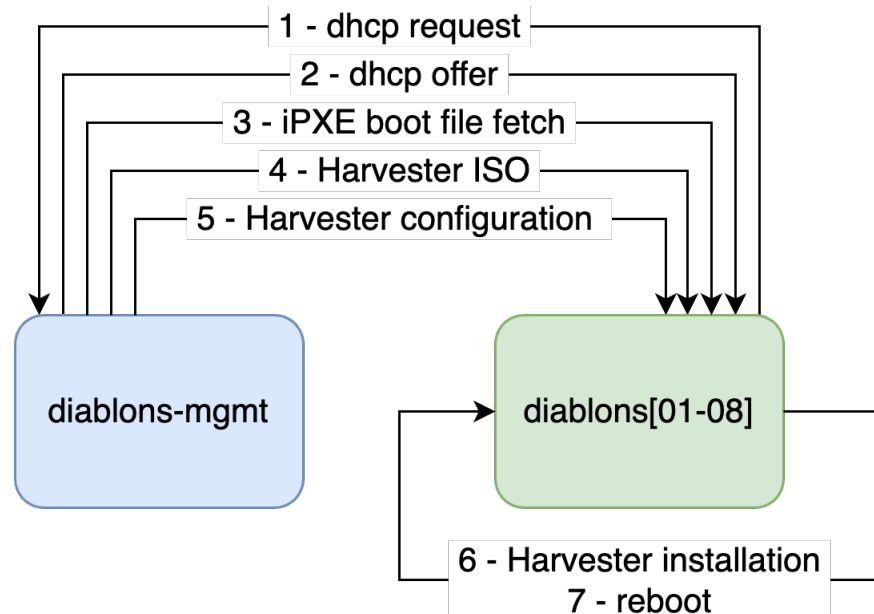


- 3 dedicated servers in HA
  - Intel dual-socket 12-core 128 GB RAM
  - provisioned with Metal-as-a-Service (MaaS) by Canonical
  - Rancher installed via RKE2 through Ansible

# Kubernetes Tools at CSCS

## ■ Harvester (SUSE)

- Hyperconverged Infrastructure (virtualization)
  - master/worker nodes of K8s clusters are VMs
- Network isolation (VLANs)
- Longhorn Storage
- Installed via iPXE boot through the network:



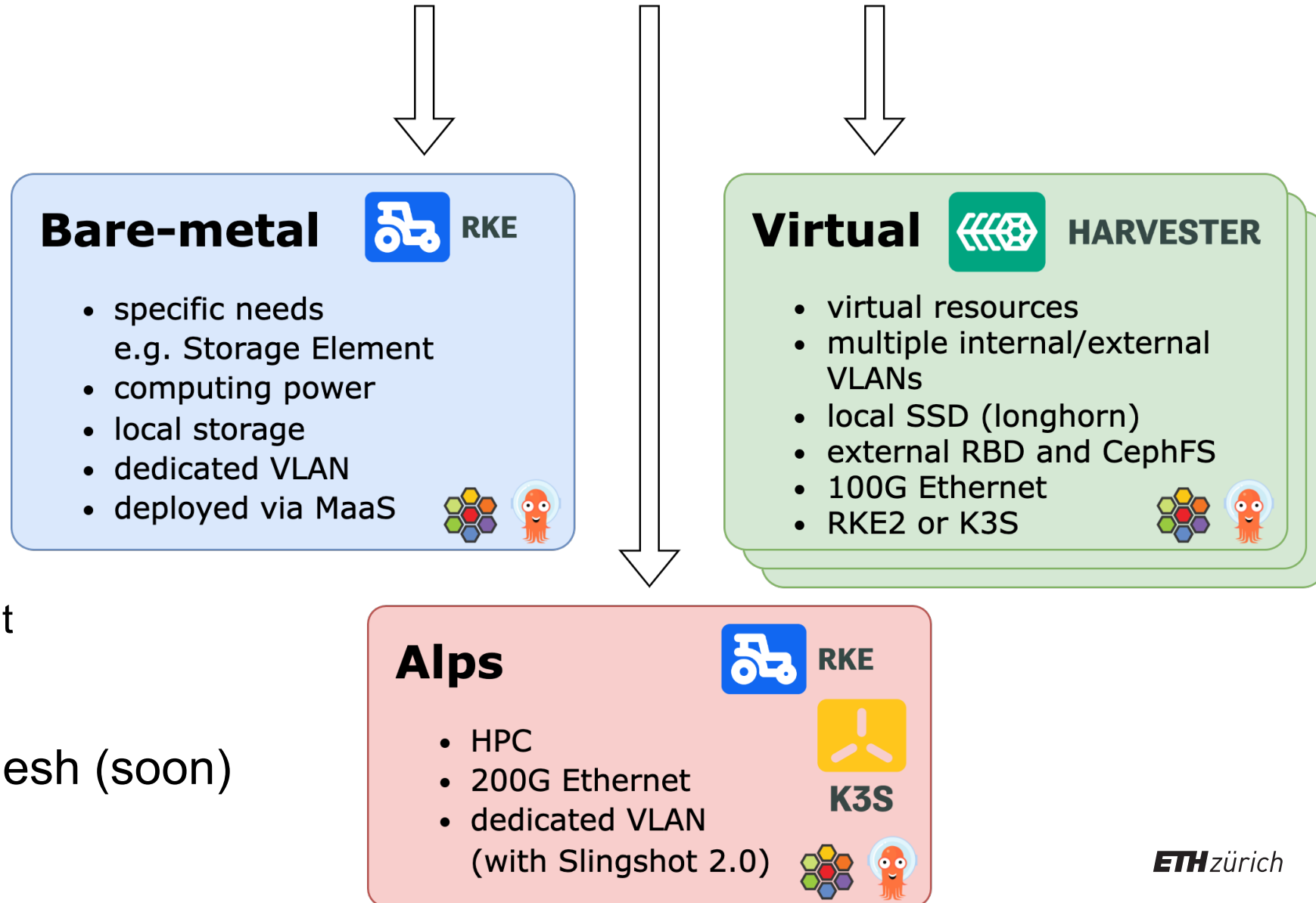
# HARVESTER

- 8 dedicated servers in HA (“Diablons”)
  - AMD EPYC 64-core 512 GB RAM  
8 TB NVMe local storage
  - 25 Gb/s (management network)  
100 Gb/s (VLAN network)  
HA mode, using LACP (in IEEE 802.3ad)
  - flexibility to scale up physical cluster

# Kubernetes at CSCS

## Scenarios

- On-demand clusters
  - different needs and requirements
- RKE2(/K3S) clusters
  - VLAN isolation
  - Rancher managed upgrades
- ArgoCD
  - cluster configuration
  - application deployment
- Cilium CNI
- Cilium/Istio Service Mesh (soon)



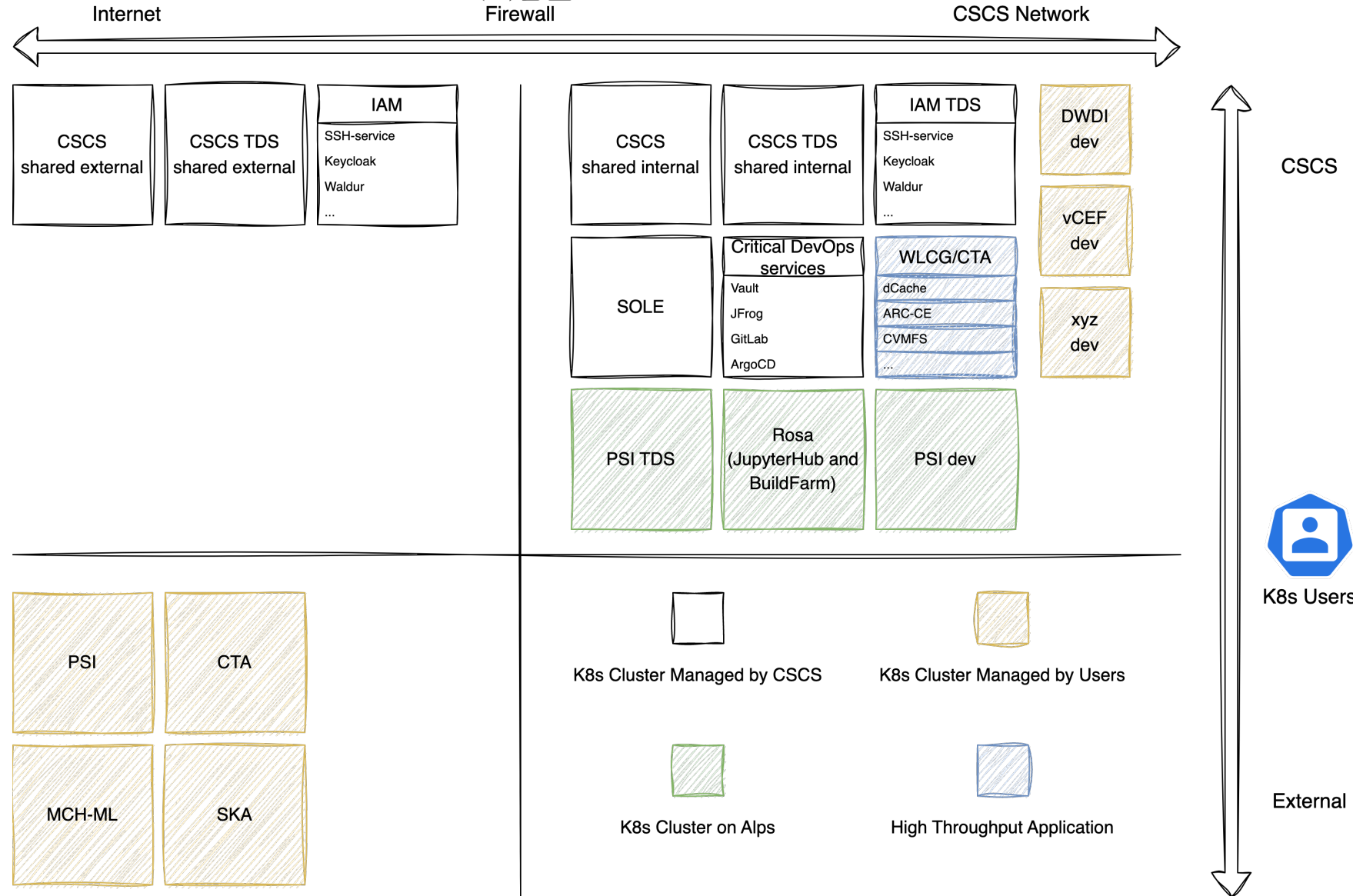
# Kubernetes at CSCS

- Baremetal
  - e.g. monitoring/ECK → Dino Conciatore “[\*Dynamic Deployment of Data Collection and Analysis Stacks at CSCS\*](#)”, HEPIX 2023, Taipei, Taiwan
  - dCache instances (Storage Element for GRID-like Workloads)
- Alps
  - Challenges
    - cluster persistence
    - networking and security
      - 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...*

# 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



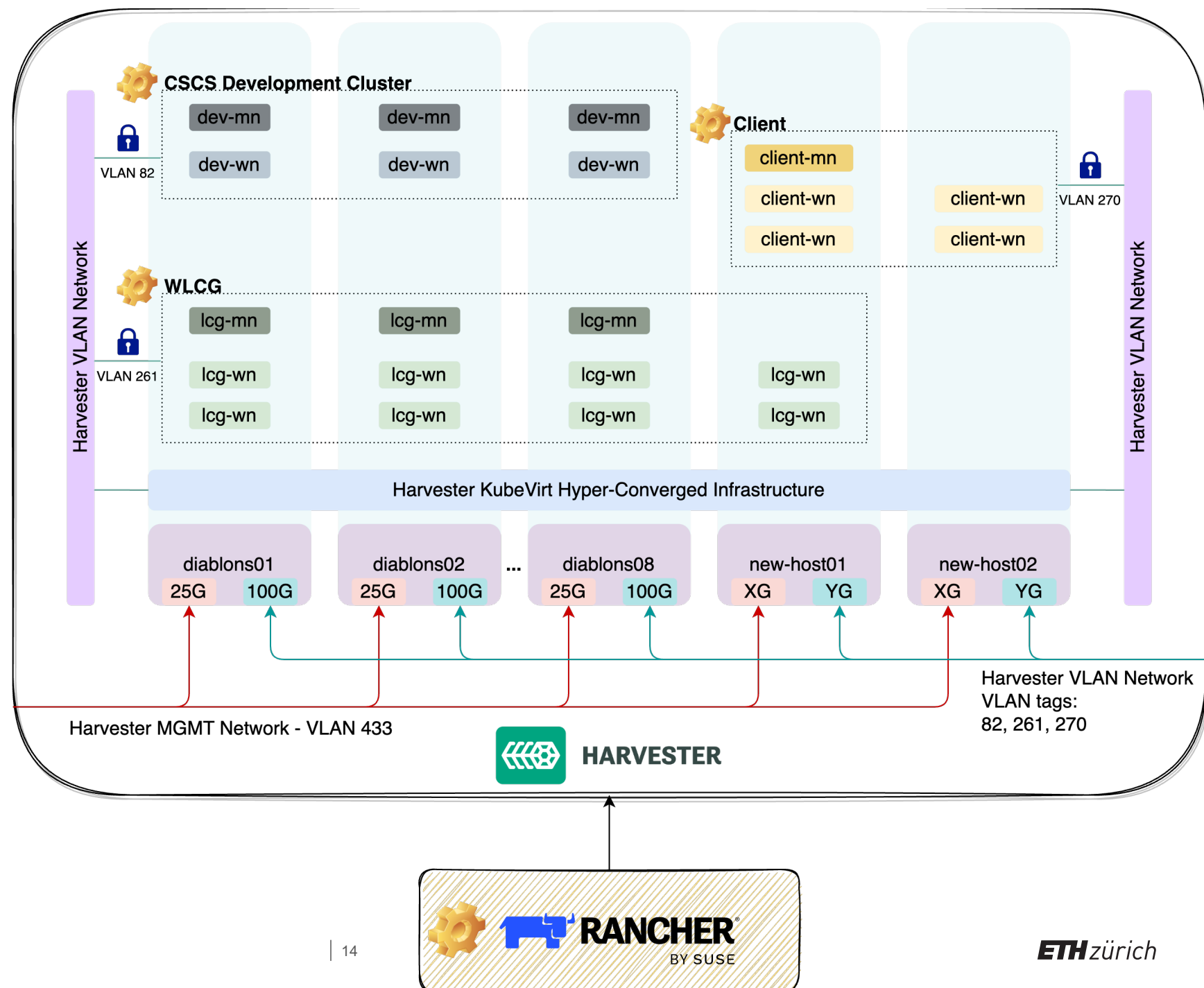
# Harvester at CSCS

## Harvester Nodes:

- physical servers
- KubeVirt cluster
- MGMT network
- VLAN network

## Harvester Cluster:

- iPXE boot
- fetch configuration
- image based install
- cloud-init provisioning
- VLAN network



# Worldwide LHC Computing Grid @ 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

- ❖ Ceph on commodity hardware

- ❖ 51 storage servers delivering 530 TiB and 22 PiB of usable NVMe and HDD capacity, respectively
- ❖ ~15 PB through dCache for WLCG

- ❖ 100 AMD EPYC Rome nodes

- 128 cores (256 CPUs), 256 GB RAM
- “Mont Fort” cluster
- **4 ARC-CEs**

- ❖ +4 nodes for dev/tds instance

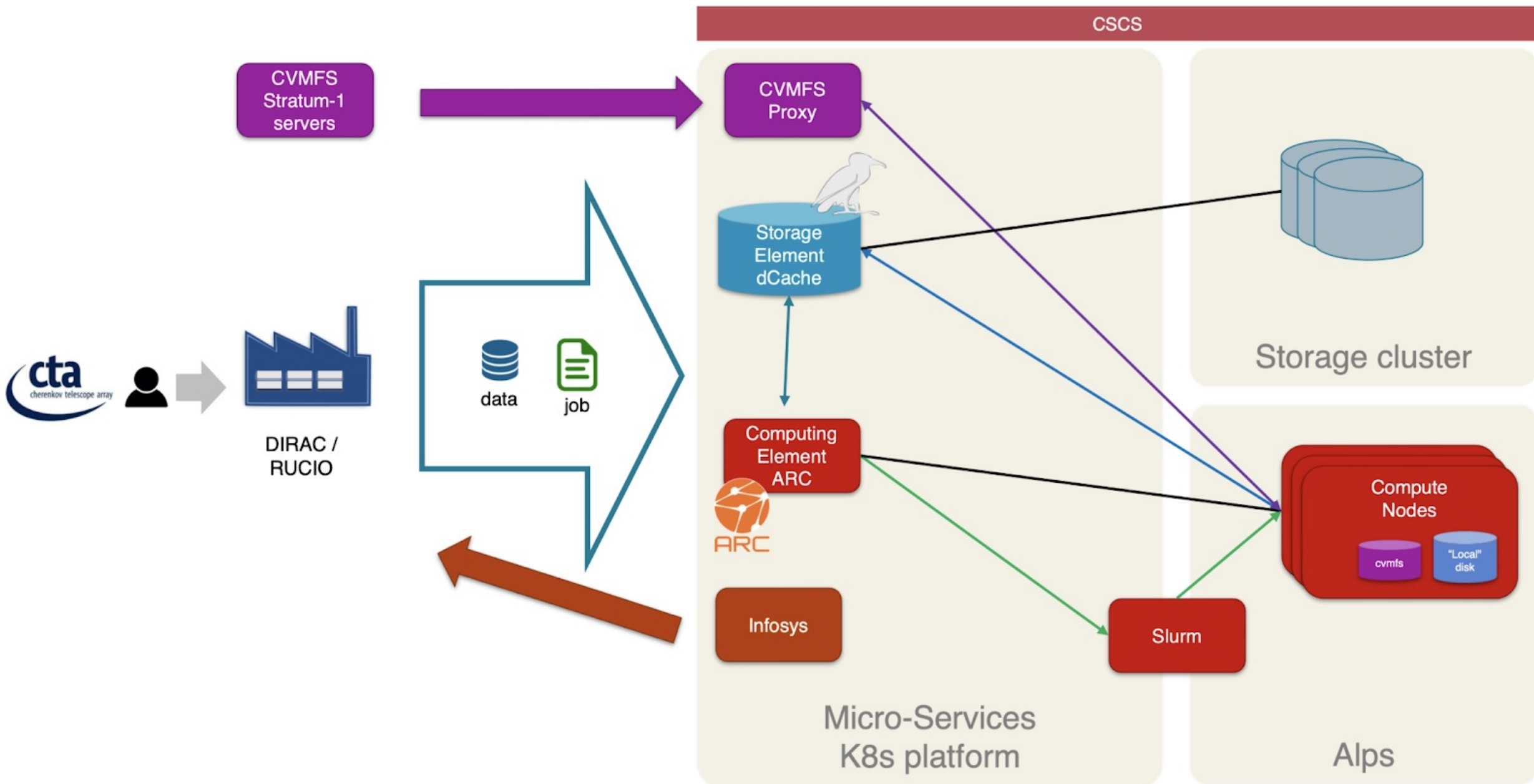
- “Mont Gele” cluster, 1 ARC-CE

- ❖ Production CE

- 300 TB shared CephFS NVMe
- 4 TB local RBD NVMe per node
- 64 GB CVMFS cache RBD NVMe per node

**AMD EPYC Rome → HS06/CPU = 22.46**

# WLCG and CTA Workflows at CSCS






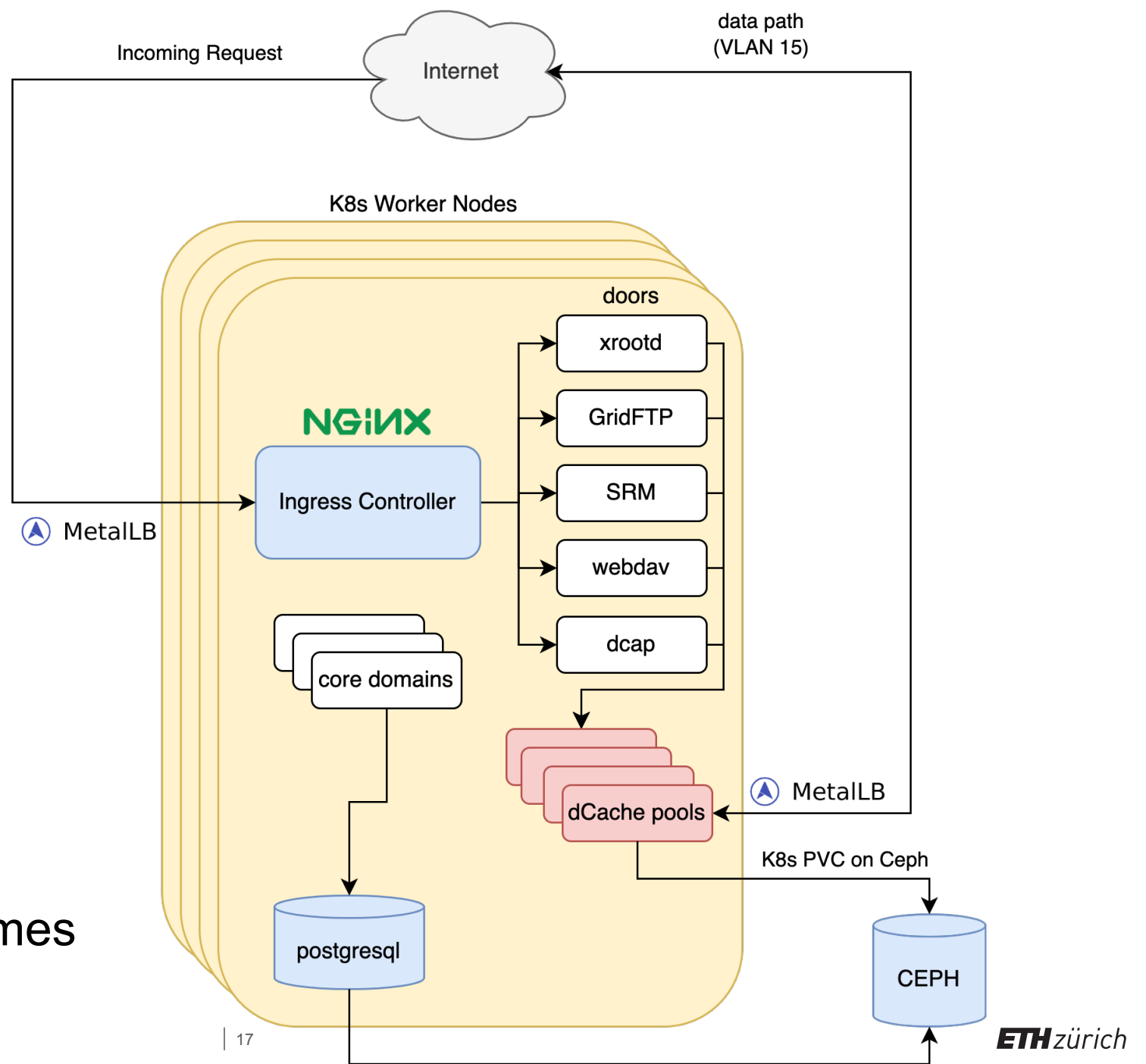


# on Kubernetes

Bare-metal  RKE

- specific needs e.g. Storage Element
- computing power
- local storage
- dedicated VLAN
- deployed via MaaS 

- 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

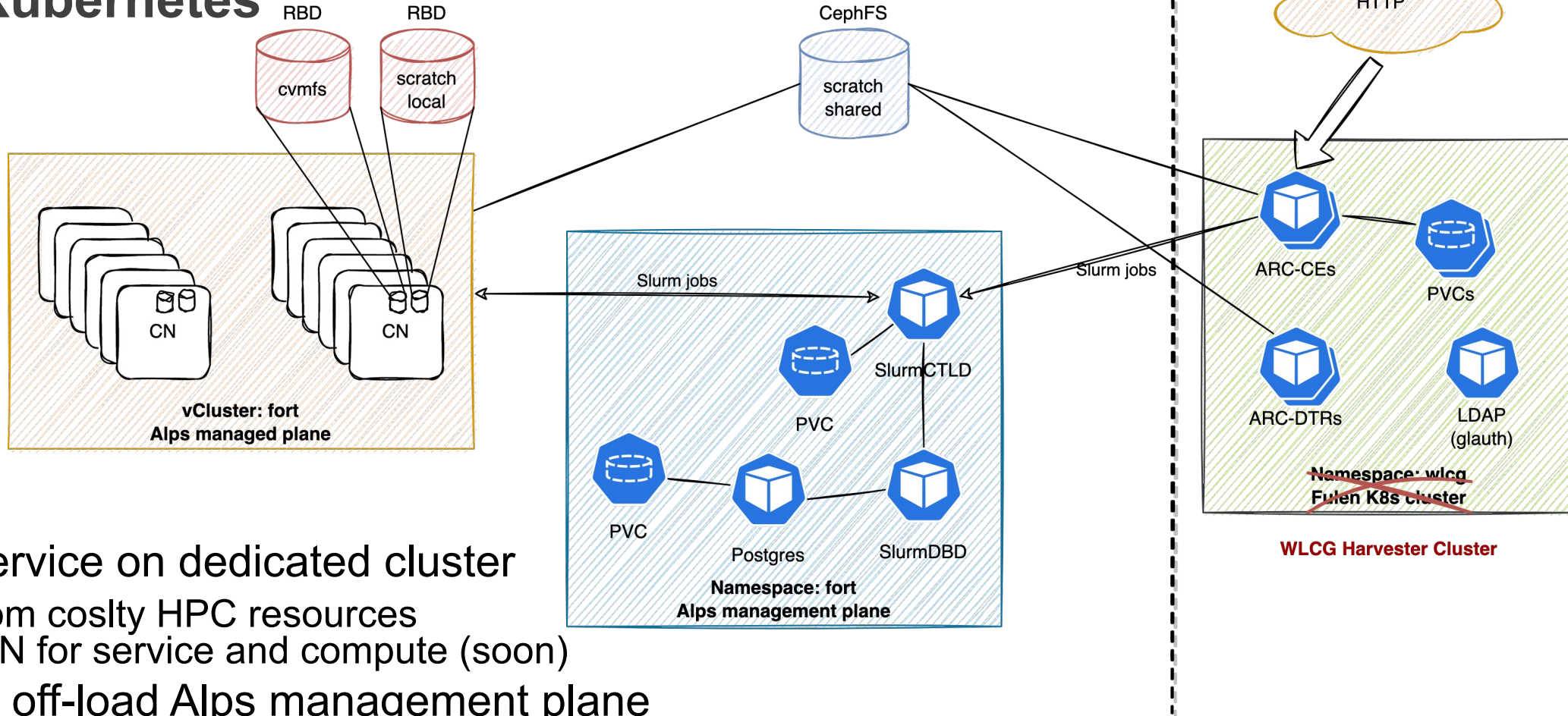




# on Kubernetes

**Virtual HARVESTER**

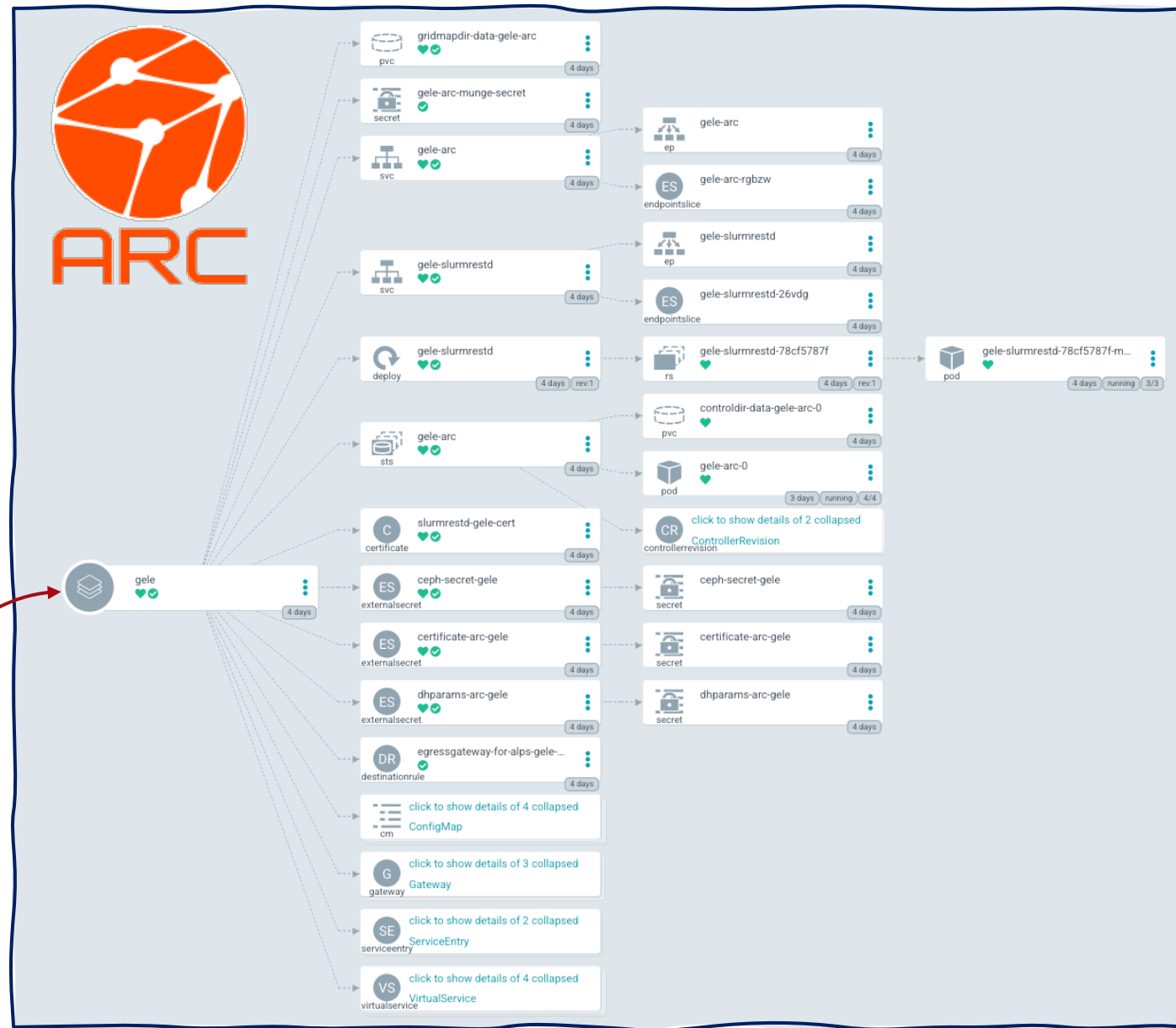
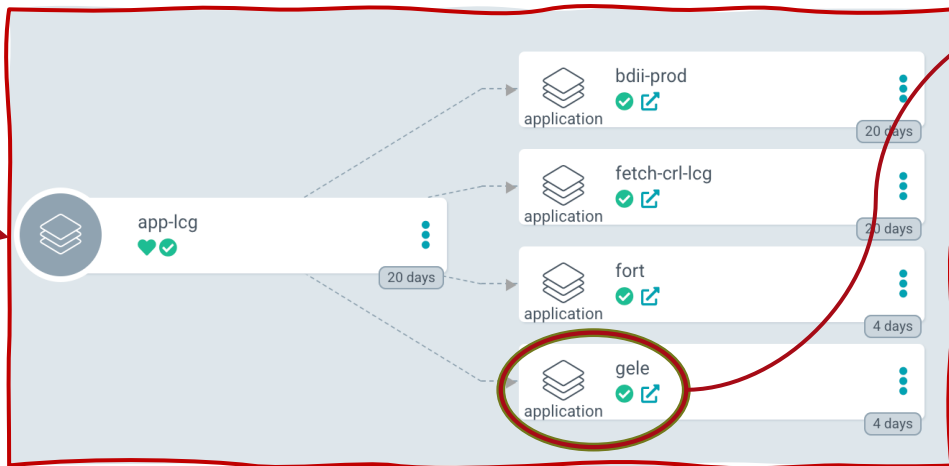
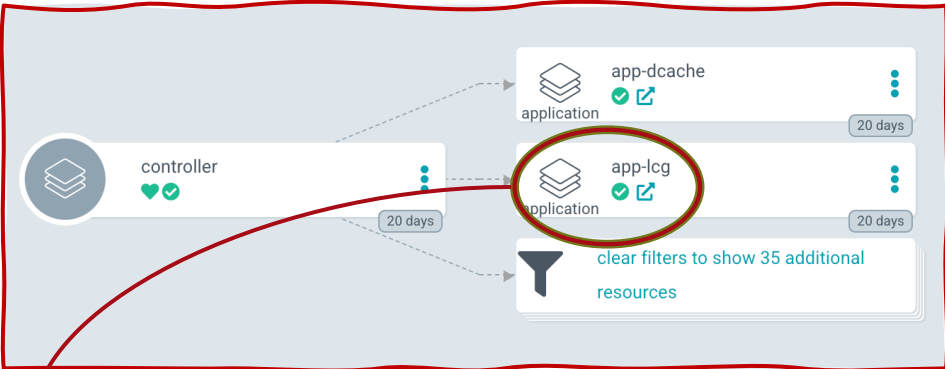
- virtual resources
- multiple internal/external VLANs
- local SSD (longhorn)
- external RBD and CephFS
- 100G Ethernet
- RKE2 or K3S



- Front-end service on dedicated cluster
  - off-load from costly HPC resources
  - same VLAN for service and compute (soon)
- Necessity to off-load Alps management plane
- Challenges from HTC workflow: storage and data-staging
- CVMFS exploited to fetch images (lightweight in comparison with HPC-standard) then executed in nested containers on Alps compute nodes

# GitOps at CSCS (ArgoCD)

## Configuration Management



# Summary and Conclusions

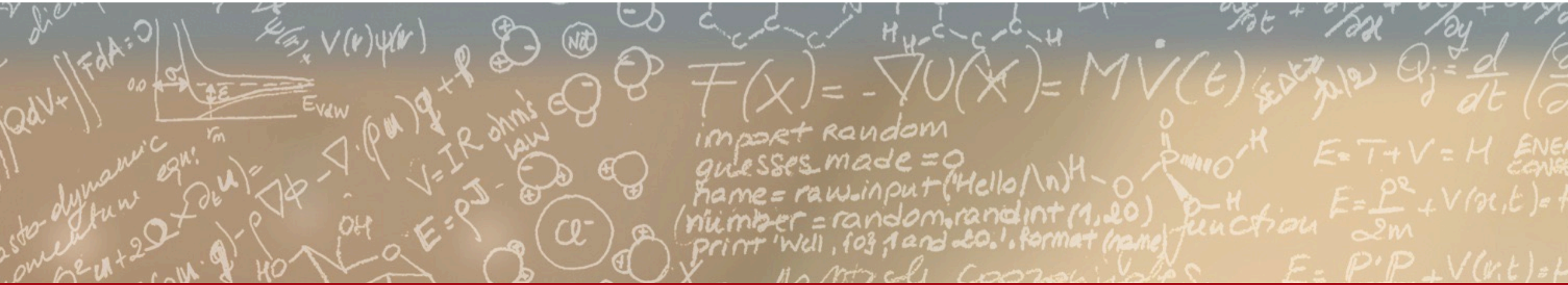
- Alps towards improving standard HPC through Cloud
  - VLANs as lifeblood
  - multi-tenancy for an increase variety of workloads, hence customers and clients
- IaC-based implementation of Alps vClusters and of Rancher-managed K8s-clusters
  - scale the infrastructure dynamically and according to the changing requirements of the customers
- Rancher/Harvester supporting management of clusters and off-load from HPC
  - central management of external and internal clusters
  - facilitating handling of micro-services
- ArgoCD eases deployment of services and configuration management
  - improved disaster recovery and CI/CD
  - potential deployment on external Clouds
- CSCS Tier-2 Grid Site as daily benchmarking exercise
  - challenging HTC workflows
  - pioneering K8s-isation of core components



**CSCS**

Centro Svizzero di Calcolo Scientifico  
Swiss National Supercomputing Centre

**ETH** zürich



**Thank you for your attention.**

**Questions?**

Contact: [riccardo.dimaria@cscs.ch](mailto:riccardo.dimaria@cscs.ch)