# Alps
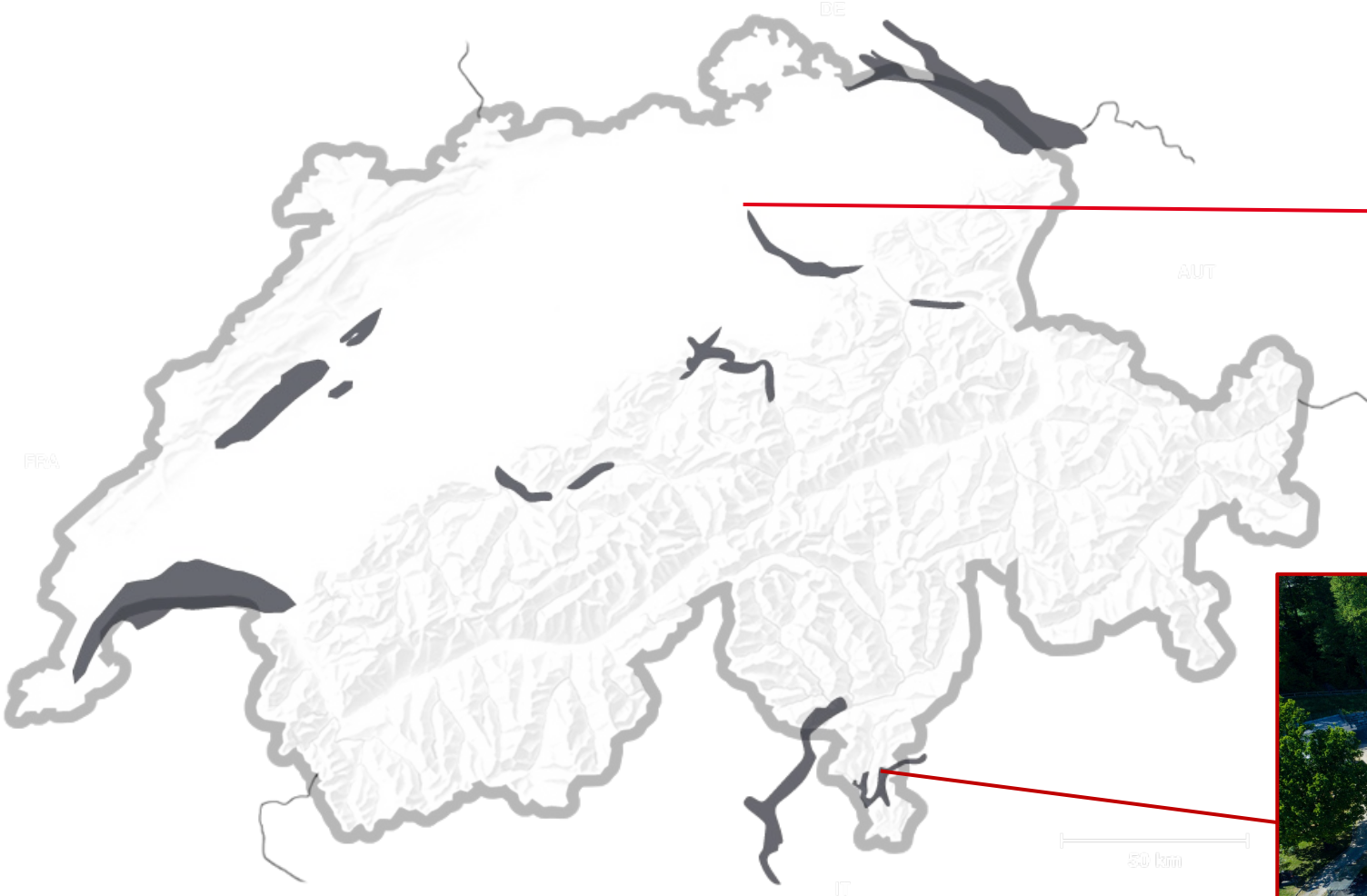# Cloud-native HPC at the Swiss National Supercomputing Centre

Dr. Riccardo Di Maria, ETH Zurich – CSCS

ISGC 2024 – ASGC, Taipei, Taiwan
March 26th, 2024

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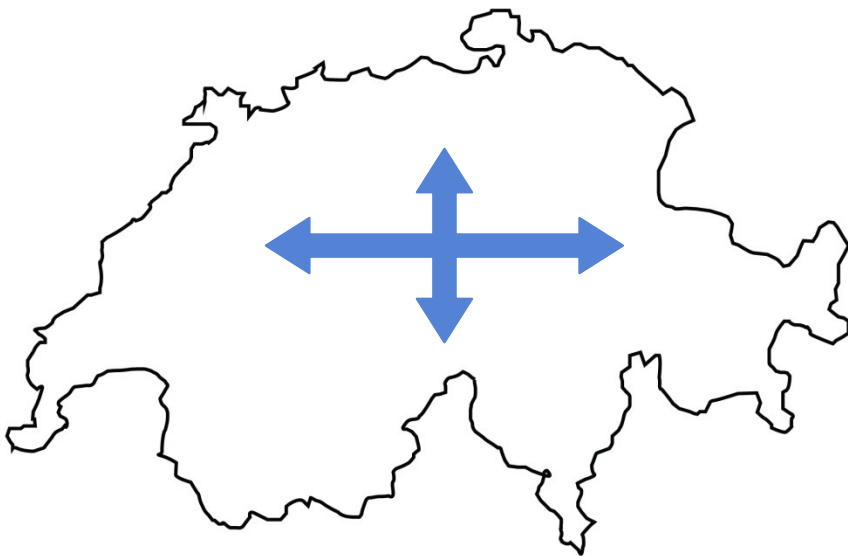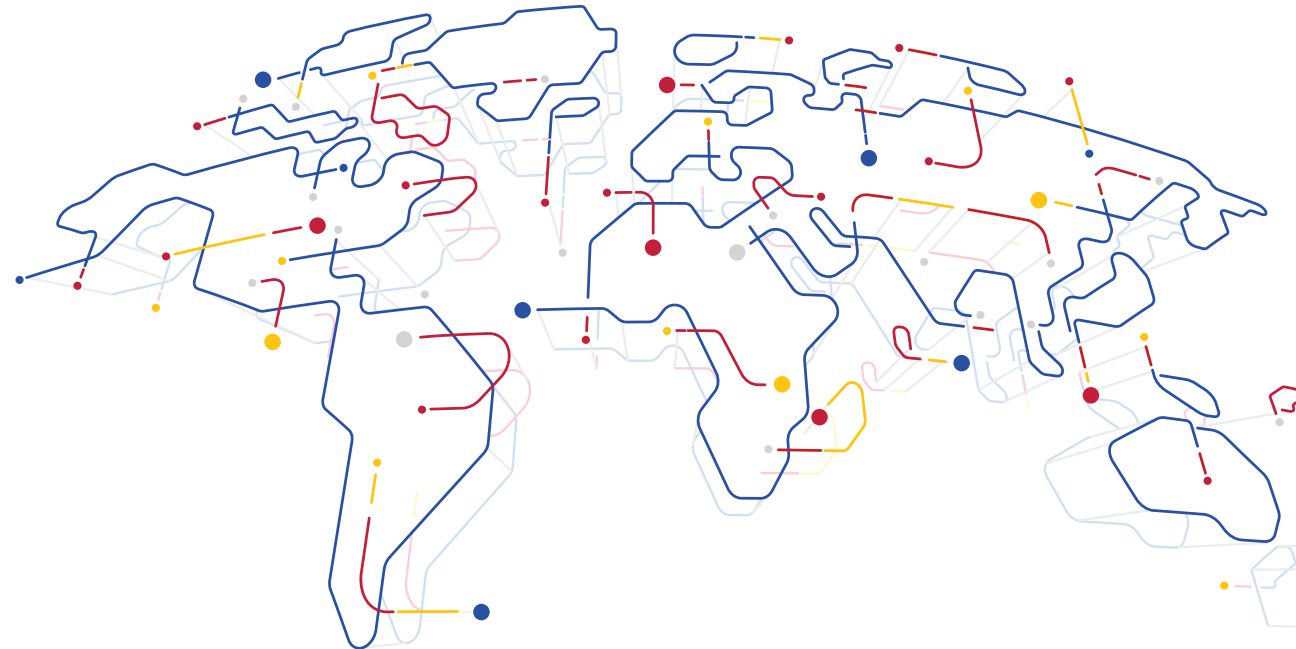


*ETH Zurich*



*CSCS Lugano*

# Mission

«We develop and operate a high-performance computing and data research infrastructure that supports world-class science in Switzerland»

- Located in Ticino since 1991

- National and international collaborations in the research of new technologies for HPC
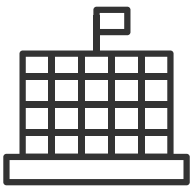
# Some numbers

## Staff

- 120 collaborators
- 25 nationalities
- Official language: English

## Building

- Lugano & Zurich offices
- 2000 m$^2$ machine room
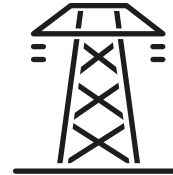- «Free cooling» with lake water

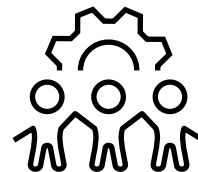## User Lab

- 2300 users
- 130 projects

## Budget

- CHF 30 Mio. operating budget
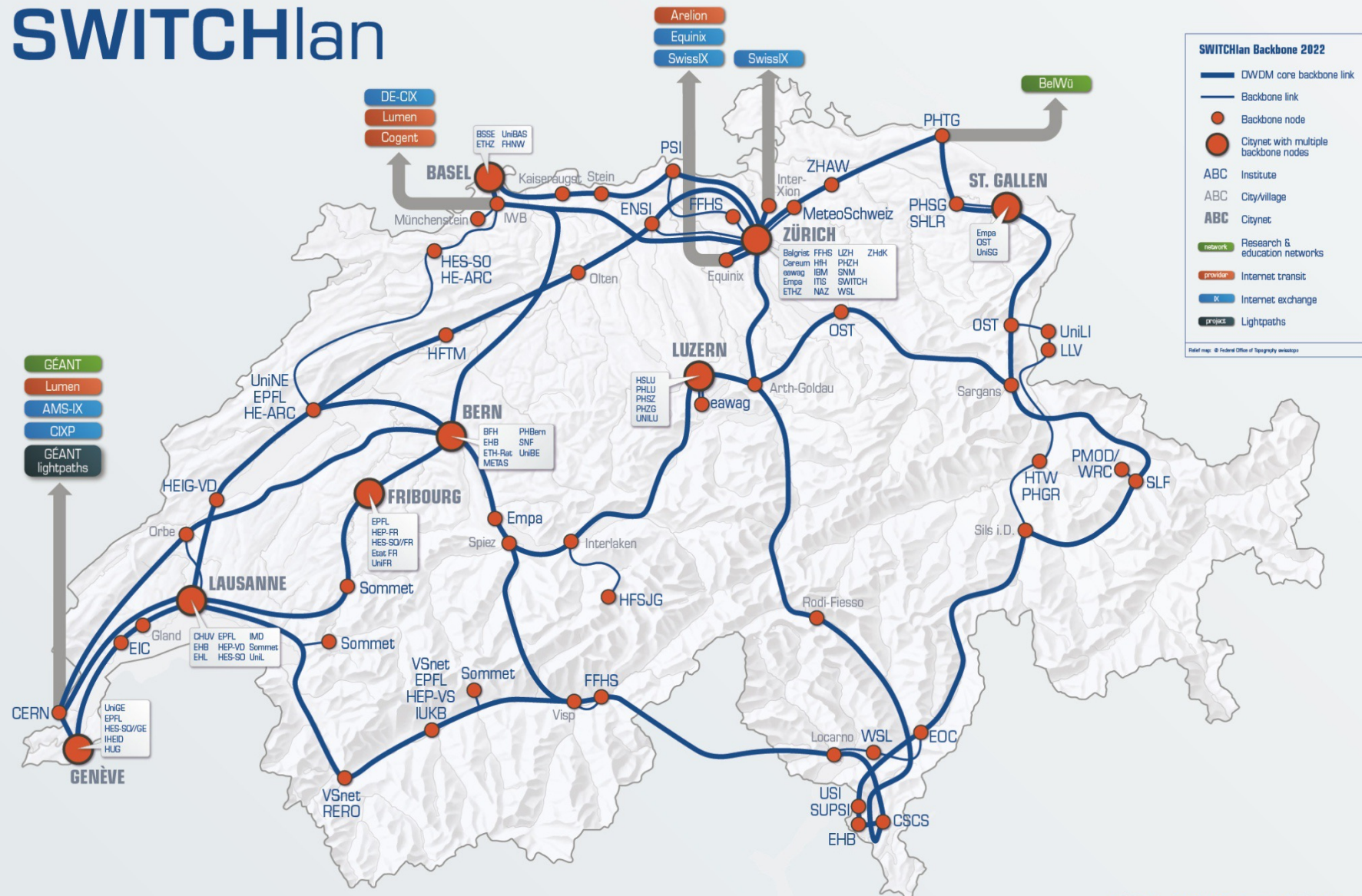- CHF 20 Mio. IT investment

## Electricity

- Currently 11 MW
- Possible extension to 25 MW
- 100% hydroelectrical source

## Third-party

- MeteoSwiss, NCCR Marvel, PSI, CHiPP, Empa, ETH Zurich, SDSC, USI, UZH, BlueBrain ...

# User Lab Program, PASC, and Partnerships

The **User Lab Program** and **PASC** (Platform for Advanced Scientific Computing) initiative provide access to resources and knowhow based on a *peer review process* and are the funded through the HPCN initiative by the ETH Rat

**Partnerships** provide access to services at CSCS based on *external funding*, and involve collaboration, exchange of knowhow, provisioning of resources

- MeteoCH
- PSI
- CHiPP
- CTA / SKA
- HBP
- UZH
- EMPA
- Euler
- BBP
- …



Publication output chart — Number of Publications by year (2012–2021), categories: Chemistry & Materials, Physics, Earth & Environmental Science, Mechanics & Engineering, Life Science, Computer Science

# Case Studies

# MeteoSwiss - Numerical Weather Forecasts

- MeteoSwiss computes its daily weather forecasts on two CSCS supercomputers (**Tsa and Arolla**)
  - future weather conditions and possible natural hazards

- Essential information for air traffic control services and disaster mitigation (e.g. radioactive leak)
  - 18 nodes with 8xV100 GPUs each + 20 post-proc nodes

- Currently migrating into Alps
  - production and development vClusters
  - AI/ML based on Kubernetes

# CHiPP - Analysis of Data from the Large Hadron Collider (LHC)





- Swiss particle physics community in the context of the **LHC** at CERN
  - help understanding the building blocks of our Universe by particle collisions

- On behalf of CHiPP, CSCS operates a mid-size Tier-2 grid site for three of the four experiments: ATLAS, CMS, and LHCb

- Grid site fully running on HPC resources
  - migrating into Alps
  - main services "Kubernetes-ised"

# The Paul Scherrer Institute (PSI)





- The **Paul Scherrer Institute** is the largest research centre for natural and engineering sciences within Switzerland
  - world-class research in three main subject areas:
    - matter and material
    - energy and environment
    - human health
- PSI operates large scale research facilities (e.g. SwissFEL) and conducts research that produces a large amount of data
  - >10PB of archived data at CSCS

- Currently migrating their dedicated HPC cluster to Alps
  - around 100 nodes currently at PSI
  - including Kubernetes-based management of clusters
  - working on geo-redundancy of the CSCS archive

PAUL SCHERRER INSTITUT



cscs

**ETH** *zürich*

# Shared HPC resources for Universities and Research Institutes

- The University of Zurich moved their local HPC service into CSCS' shared resources back in 2014 on Piz Daint
  - migrated into Alps in 2021
- Resources are shared among 15 different groups inside UZH
  - UZH local Scientific IT service (S3IT) distributes the resources and supports their scientists with their research
- Classic HPC
  - ~200 multi-core nodes

**University of Zurich** UZH

*Similar model followed by **Empa** (Swiss Federal Laboratories for Materials Science and Technology), **Eawag** (Swiss Federal Institute of Aquatic Science and Technology), **USI** (Università della Svizzera Italiana), **FHGR Chur** (University of Applied Sciences of the Grisons), **SDSC** (Swiss Data Science Center) and others*

cscs

**ETH**zürich

# Observatories – Square Kilometer Array & Cherenkov Telescope Array



- Shaping the future of Astronomy at Petabyte-to-Exabyte scales is of strategic importance for Switzerland to support research, education and innovation

- **Cherenkov Telescope Array (CTA)** and the **Square Kilometer Array (SKA)** Observatories are at the forefront of this new Big Data astronomy
  - ~1 Exabyte per year, to be distributed around the globe

cscs

ETH zürich

# MARVEL – Compute and Long-Term Storage Resources

- **MARVEL** (National Centre of Competence in Research - NCCR) aims to radically transform and accelerate the design and discovery of novel materials

- CSCS supports MARVEL by providing compute and long-term storage resources, and the two organizations work closely together in defining the **Materials Platform** of the future

- Marvel is already running on Alps; Materials Cloud uses OpenStack and is now moving into Kubernetes

# EXCLAIM : Cloud-Resolving Weather and Climate

## High-level implementation plan



Figure 2: High-level road-map of the EXCLAIM project showing the five key development threads

A community united around an ambitious roadmap that jointly decides on the milestone simulations, shares data for analysis, and develops a platform with dedicated functionality in a specialized platform. Model for how we see communities develop their HPC roadmaps in the future.



**EXCLAIM**

Extreme scale computing and data platform for cloud-resolving weather and climate modeling

https://exclaim.ethz.ch

CSCS

ETH zürich

# Architecture: a Change in Paradigm

# CSCS Goals

Founded in 1991, CSCS, the Swiss National Supercomputing Centre, develops and provides the key supercomputing capabilities required to solve important problems to science and/or society. The centre enables world-class research with a scientific user lab that is available to domestic and international researchers through a transparent, peer-reviewed allocation process. CSCS's resources are open to academia, and are available as well to users from industry and the business sector. The centre is operated by ETH Zurich and is located in Lugano with additional offices in Zurich.

- CSCS has been running supercomputers and clusters for years, using logical abstractions (projects, Slurm queues, POSIX permissions, etc.) to partition and distribute computing power to the different groups of users

- As the numbers of science domains and projects grow
  - provisioning dedicated clusters becomes expensive (cost, manpower, management, etc.)
  - integrating completely different workloads on a single system is complex (e.g., Slurm on Piz Daint, WLCG HTC vs. UL HPC queues)

# HPC and Cloud Convergence

- Science and engineering requires more and more computer-assisted experiments
  - simulation of physical phenomena and behaviours
  - Digital Twins
  - design engineering products
  - AI/ML statistic solutions

1. HPC offers high-performance compute and data access
   - improves Time to Solution
   - manage efficiently data to compute
   - bare-metal performance, fixed amount of resources

   1. + 2. = ?

2. Cloud offers high flexibility for business needs
   - XaaS – business logic as a service
   - economy of scale – oversubscription of resources
   - virtualized resources, scalable to infinity (and beyond)

To infinity and beyond...

cscs

ETH zürich

# Achieve Best of HPC and Cloud

- Performance and flexibility
  - container as a virtualization layer with OCI hooks
    - keep OS near bare metal – accelerators and high-speed network drivers
    - bring low-level libraries in the container
  - own user environment
    - decouple HPC programming environments from underlying layers
    - UE can potentially become "just" an artifact mounted in the container

- Separation of concerns with layers
  - Platforms
    - provisioning of services with Kubernetes and/or Nomad
    - container as an abstraction layer for compute nodes
  - **Infrastructure as code**
    - APIs and configuration management
    - multi-tenancy: exclusive <u>compute</u>, <u>network</u>, and <u>storage</u> segregation
    - update components independently and minimise downtime across tenants

- HPC business logic
  - web-facing API to access HPC resources (submit jobs, move data)
  - web gateway

**vCluster**
versatile software-defined cluster

# Infrastructure Today and Tomorrow: the Systems







- top-left: **Piz Daint**
  current Cray HPC platform of CSCS
  (Intel CPU + NVIDIA P100)

- top-right: **LUMI**
  flagship European HPC system
  ETHZ/CSCS is a partner (AMD CPU+GPU)

- bottom-left: **Alps**
  future Cray Shasta HPC platform of CSCS
  (NVIDIA GraceHopper, …)

# The Research Infrastructure: Alps

# Alps

- Alps is an HPE Cray EX supercomputer meant to be our new flagship <u>infrastructure</u>

- Multi-phase installation started in 2020

- Multiple **geo-distributed** infrastructure

- Specification
  - 1024x MC nodes (AMD Milan and Rome 7742) 256/512GB RAM
  - 144x NVIDIA A100 GPU nodes
  - 32x AMD MI250X GPU nodes
  - 128 AMD MI300 GPU nodes
  - *thousands* of GraceHopper nodes

  - Slingshot network
  - 2 availability zones (HA, non-HA)
  - 100% liquid cooled

water cooled blades

# Installation



big delivery



moving the first racks

# Installation



some boxes



water cooling

# Installation



front: HA zone (compute + management)



back: Slingshot network

# Installation



internal pipes, doors and panels





some Slingshot fiber cables

# Alps Phase II:
# the Grace-Hopper "Super Chip"



Each node will have 4 Grace-Hopper modules

- 1 Grace CPU socket and
  1 Hopper GPU per module

- all-to-all cache-coherent memory NVLINK
  between all host and device memory



- one NIC per module

# Network, Slingshot, and Isolation

- Fully redundant connectivity to internet

- 2x 400 Gbps CSCS – Internet links

- Fully redundant HSN – ethernet connectivity

- Dedicated VLANs and Public IP space available

- Slingshot networks are based on Ethernet

- HPC niceties such as traffic congestion, low latency

- Ethernet niceties such as VLAN isolation and native TCP/IP stack (ARP as well)

- 8x 100Gbps uplinks from Alps' Slingshot to CSCS network
  - more in the future



Internet

Ext. FS

CSCS Ethernet Network

Lustre

CN

CN

CN

CN

VLAN

CN

VLAN

CN

VLAN

CN

VLAN

Alps Slingshot
High Speed Network

cscs

**ETH** *zürich*

# Storage in the Alps Era

- Consolidate backends and technologies

- <u>Most</u> storage areas based on Lustre

- Filesystems/spaces available
  - capstor
  - iopsstor
  - purpose, dedicated filesystem

- Metadata ops for the different areas hitting different servers

- Evaluating alternatives for multi-tenancy

- CEPH is available



**other storage**
blazing fast
or dedicated

**facility**

**capacity storage**
persistent and fast

**slow storage**
long-term
dataset archive

Eiger vCluster

Lustre appliance
/scratch/iopsstor
(iopsstor)

Clariden vCluster

Lustre appliance
/scratch/eiger
/users
/project
/store
(capstor)

tape library
(backup)

Lustre appliance
/scratch/mchstor1
(dedicated)

Tasna vCluster

CSCS

**ETH** *zürich*

# Storage in the Alps Era: ClusterStors

- capstor
    - 100 PiB HPE ClusterStor E1000D
    - spinning disks
    - raw performance:
        - **~1TB/s**
        - 300k write iops | 1.5M read iops
    - 8480 spinning disks (16 TB each)
    - 6 metadata servers
    - 11 full racks
    - Slingshot 11

- iopsstor
    - 3.2 PiB HPE ClusterStor E1000F
    - all flash
    - raw performance:
        - 240GB/s write | 600 GB/s read
        - **13.5M write iops | 18.4M read iops**
    - 240 NVMe devices (30TB each)
    - 2 metadata servers
    - 1 rack
    - Slingshot 11

cscs

ETH zürich

# Storage Pics



capstor: 126 disks per enclosure!



back of iopsstor



back of Piz Daint's Sonexion 3000

# High Level View of Monitoring

- CSCS runs a large instance of ElasticSearch and Kafka on Kubernetes
  - 112 nodes, 2TB NVMe/SSD per node
  - 100 Gb outgoing connectivity
  - 45000 docs processed per second
  - 230 billion documents online

- Overall monitoring using a combination of
  - Nagios
  - ElasticSearch
  - Grafana
  - Kibana

- Storage and network metrics across the centre

- Dino Conciatore *"Dynamic Deployment of Data Collection and Analysis Stacks at CSCS"*, HEPiX 2023, Taipei, Taiwan

- vCluster Compute Node telemetry
  - in-band (/proc/cray/counters + metricbeat, pushed to ElasticSearch)
  - out-of-band (based on LDMS, pushed to ElasticSearch)

- Dashboards and metrics can be tailored to specific needs

# vCluster: versatile software-defined cluster

# Alps and vClusters

- Alps is partitioned in logical units (vClusters) aggregated into platforms

- Each of these platforms or vClusters serve a group of user communities with common interests or needs
  - HPC platform: 3 vClusters (prod, test, dev)
  - Grid-like platform: 3 vClusters (WLCG prod, WLCG test, CTA prod)
  - etc.

- vClusters are instantiated on specific hardware and areas depending on the use-cases
  - network
    - VLAN enabled for potential network separation
  - services, storage, etc.
  - users and configuration/policies
  - flexibility in terms of hardware and software choices

- Resource elasticity and scalability: growing or shrinking vClusters is easier

# Cloud-native Supercomputing and Big Data

**vCluster: versatile software-defined cluster**
- custom user environments
- manage platform services, provisioning of clusters
- possibility of network and storage isolation



| WLCG | SKA | CTA | Materials Cloud | …. | ICON-22 |

**HPC Platform**

**Weather & Climate**

**ALPS Infrastructure**

| User environments management | Platforms and services management | Infrastructure as Code |

A layered, versatile research infrastructure, providing the environment for various communities to excel and innovate in their fields

cscs

ETH zürich

# Cloud-native Supercomputing and Big Data



Piz Daint
User Lab and WLCG

MeteoSwiss
Arolla/Tsa
COSMO/ModInterim

# vCluster Access Methods and Configuration

- **Dedicated login nodes**
  - intended for workflow preparation
  - non-CPU intensive pre/post processing activities

- **Basic access via SSH**
  - shell interaction: preparing and submitting batch jobs
  - transferring configuration files and input/output files
  - MFA enabled

- **Alternative access methods and services**
  - FirecREST - https://products.cscs.ch/firecrest/
    - RESTful services gateway and interface for managing HPC resources

- **Container runtime**
  - Sarus - https://products.cscs.ch/sarus/
    - user-friendly way to instantiate feature-rich containers from Docker images

- **Job scheduler**
  - Slurm as workflow manager
    - accounting via dedicated SlurmDB
  - default scheduling policies and configuration
    - single job queue
    - exclusive node usage
  - additional tailoring/fine-tuning enabled
    - compute node usage quotas
    - additional queues
    - advanced priority management
    - job usage reporting

cscs

ETH zürich

# vCluster
# Levels of Interaction

- Cray System Management (CSM) / Shasta
  - micro service architecture control plane management



Infrastructure Admin

vCluster Admin

vCluster User

Infrastructure Admin Control (SSH to Management plane nodes)

ClusterAdmin Role API access to own vCluster

User Access Deployed Services Managed Plane (SSH, etc)

Manta (Node power actions, image building, node health) https://github.com/eth-cscs/manta

Admin Access Deployed Services Managed Plane

User Access Deployed Services Kubernetes

API Acces Only

**Management Plane (HPE CSM)**

CSM API

Node Management Network

CSM HA Microservices

High Speed Network (Future removal)

CSM Utility Storage

**Alps Managed Plane (Slingshot HSN Network)**

vCluster (Optional VLAN/VNI Isolation)

CN

LN

Optional Extension to Rancher Cluster

Subdirectory Isolated Mounts

Subdirectory Isolated Mounts

CAPSTOR (Lustre)

IOPSTOR (Lustre)

Slingshot

Edge Router Ethernet/Slingshot

Ethernet

**Support Services**

vCluster Support Services (Rancher/Harvester) (VLAN Isolation)

Kubernetes API

Kubernetes Services

VLAN Isolation

CSI, RBD, CephFS

VAST (NFS)

Ceph

cscs

# HPC and Cloud-Native

# Service Orchestration

a. **Full service on HPC**
  - security challenges
    - VLANs help
    - ad-hoc configurations between management and managed plane
  - inefficiency on costly resources
  - additional "virtualisation" layer
    → complexity (e.g. network)

b. **Adjacent/front-end services orchestrated within**
  - compute nodes of Alps using NOMAD
  - dedicated Kubernetes clusters
    - efficient use of HPC resources
    - necessity of middleware/interface between user and compute

- **Orchestration and maintenance driven by use-case needs**

- **Advantages**
  - decoupling from the infrastructure
  - 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**
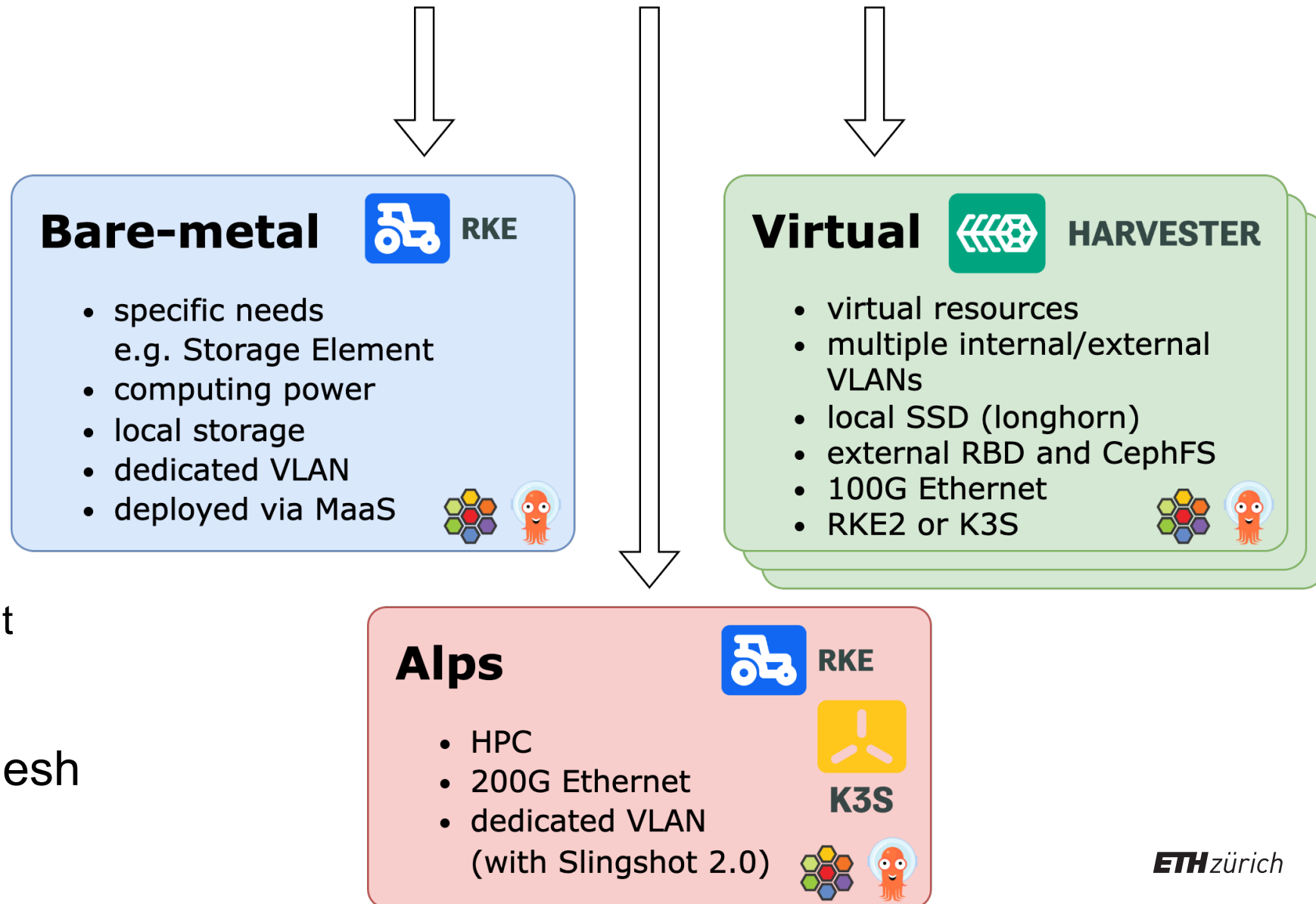
- **Challenges**
  - additional "moving parts" and complexity layers
    - networking: Cilium vs. Calico, service mesh
  - security
    - additional configuration

- Riccardo Di Maria *"The WLCG Journey at CSCS: from Piz Daint to Alps"*, HEPiX 2023, Taipei, Taiwan

# Kubernetes at CSCS
## Scenarios

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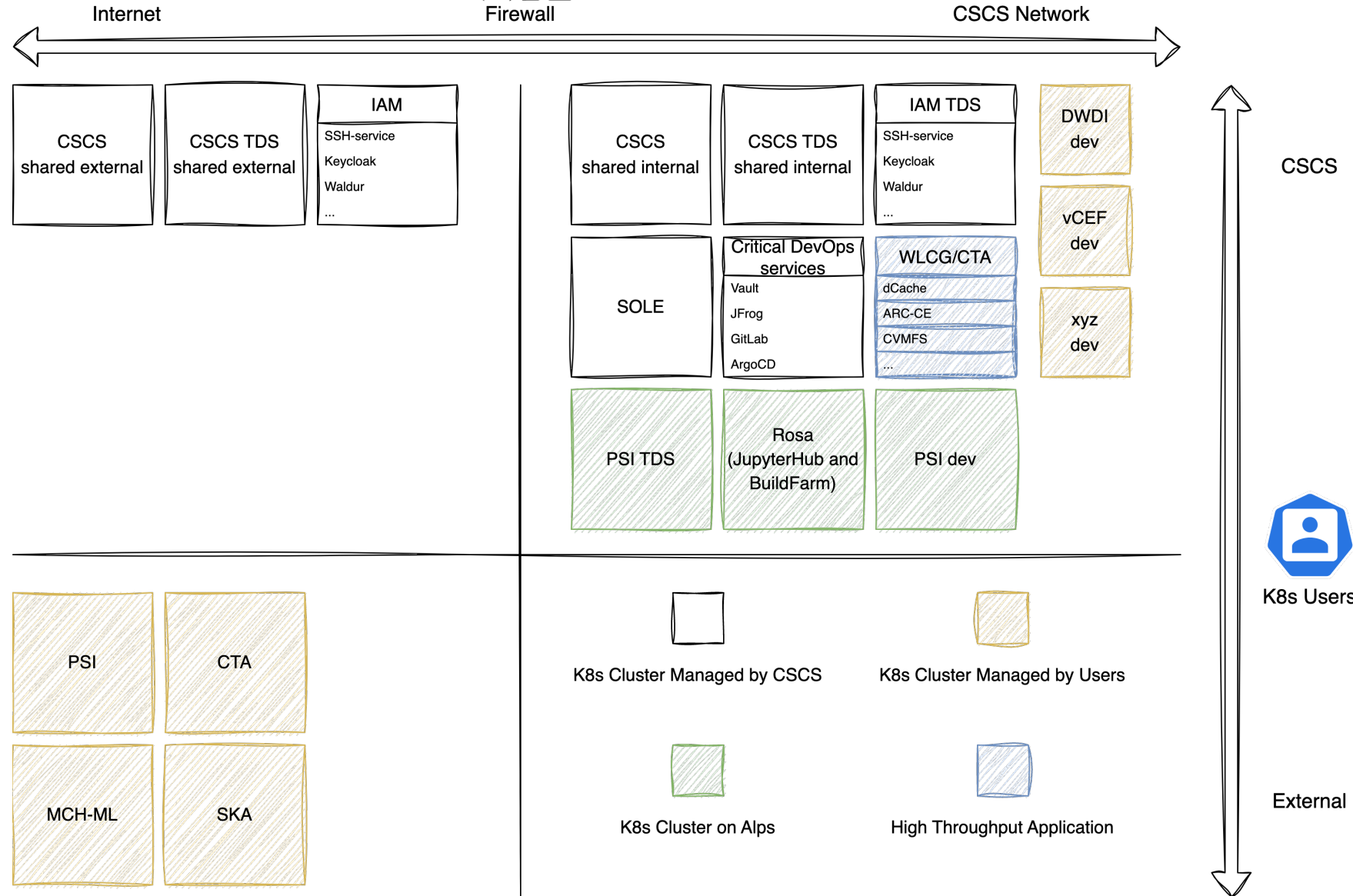- Cilium CNI
- Cilium/Istio Service Mesh

**RANCHER®** BY SUSE

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

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

**Alps** RKE K3S
- HPC
- 200G Ethernet
- dedicated VLAN (with Slingshot 2.0)

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
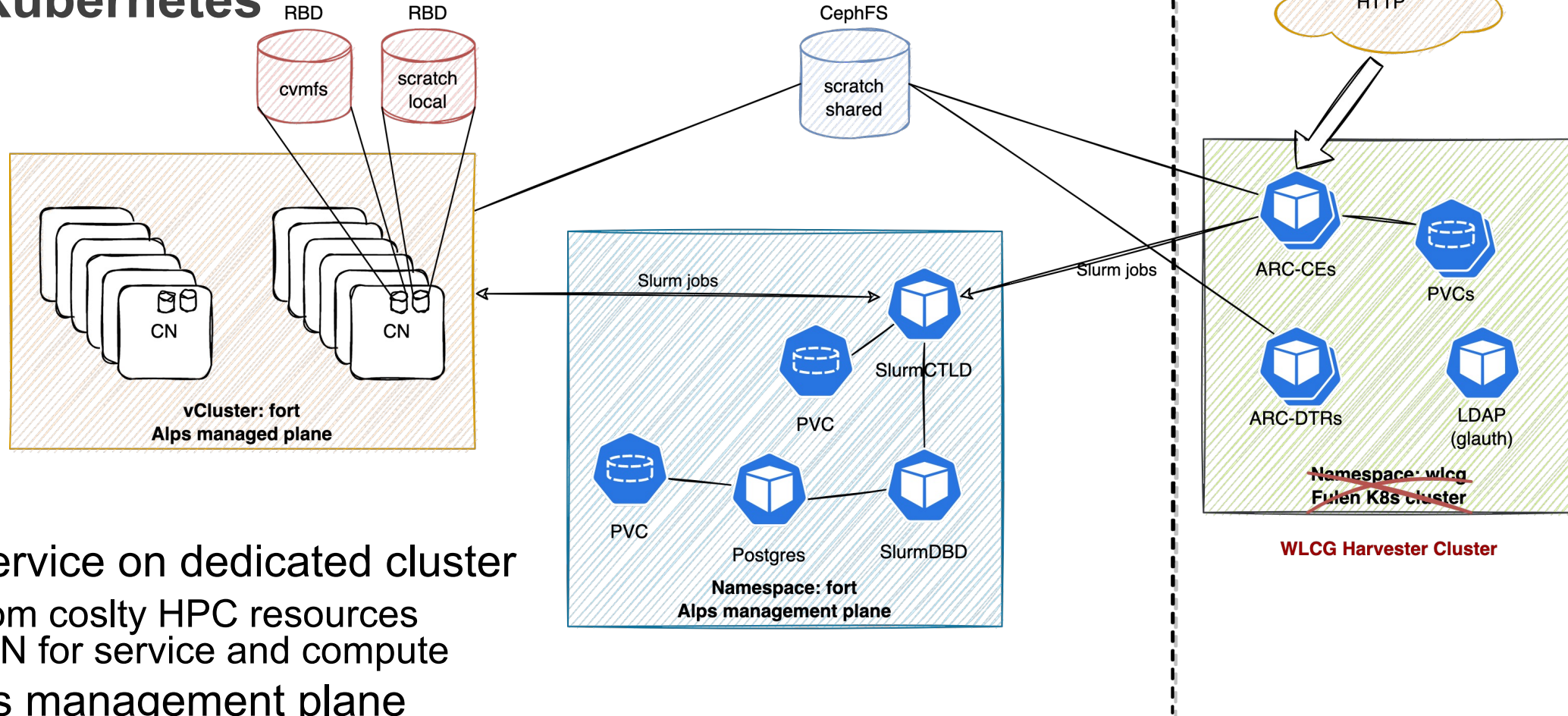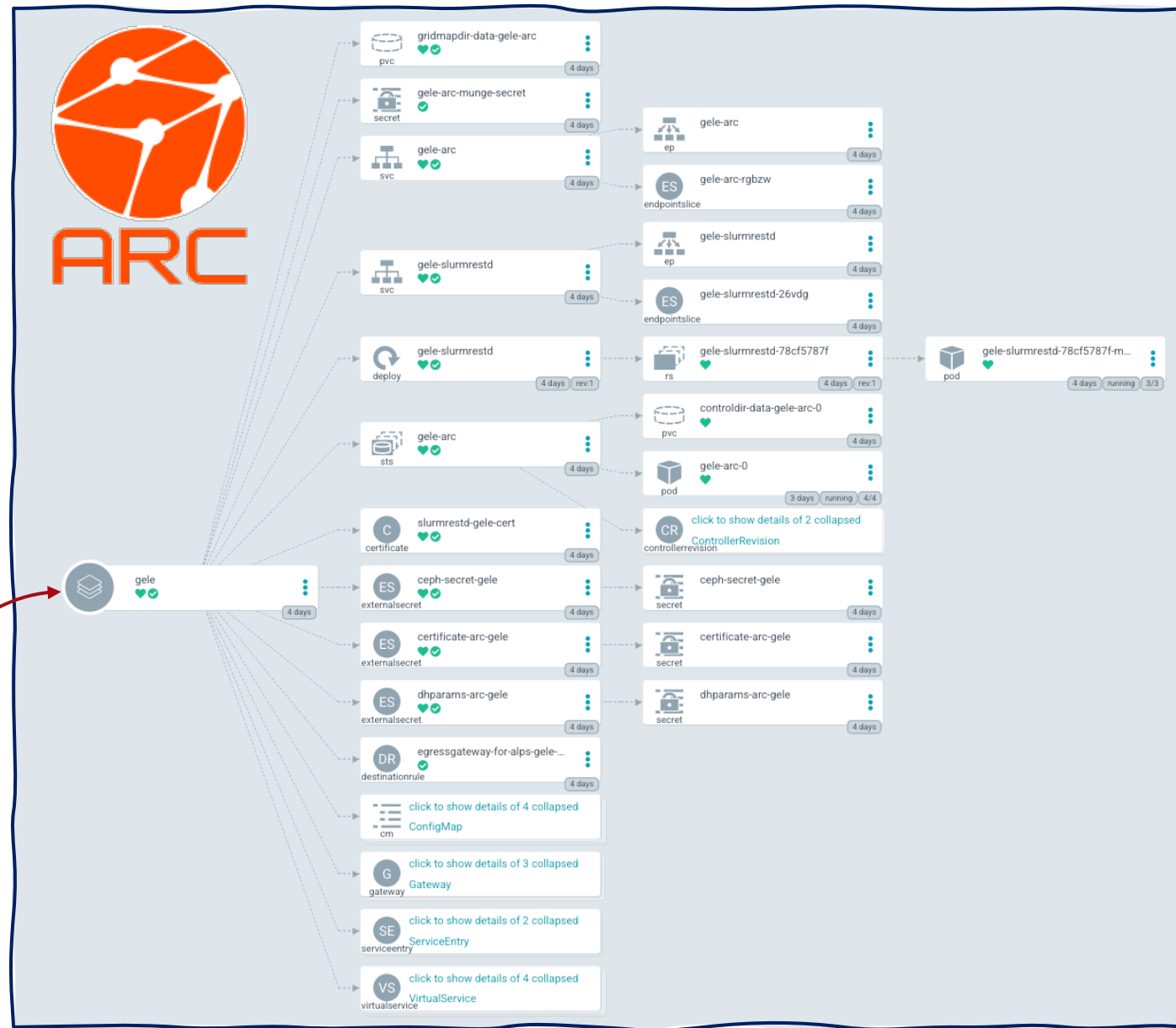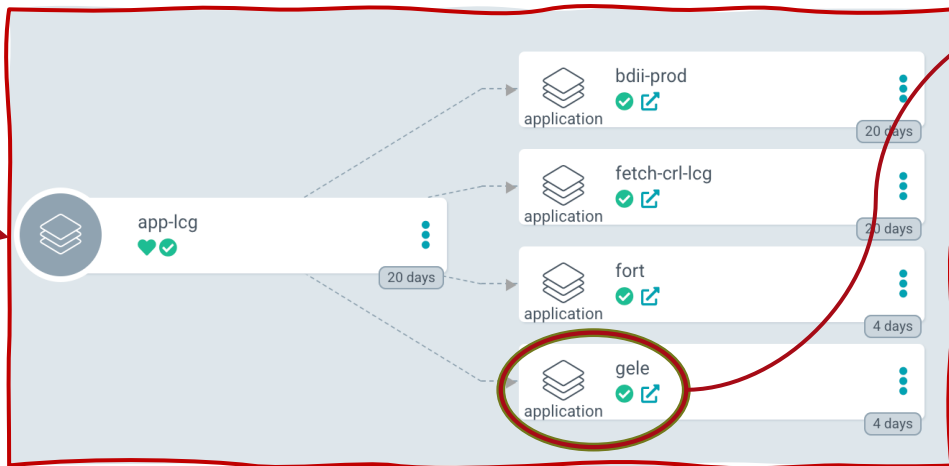
- ## Vault

- ## ArgoCD

# on Kubernetes



- Front-end service on dedicated cluster
  - off-load from coslty HPC resources
  - same VLAN for service and compute
- 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



*(DIY solution)*

GitOps

| 43

# User Experience

# User Experience | On-Boarding, Training, and Support

- Documentation

  - CSCS Knowledge Base https://docs.cscs.ch

  - user portal https://user.cscs.ch (frozen; content is transitioning to KB)

  - tutorials https://www.cscs.ch/publications/tutorials

  - video portal https://www.cscs.ch/publications/video-portal


- User management

  - partners can self-manage users https://account.cscs.ch

- User support portal https://support.cscs.ch

cscs

ETH zürich

# User Experience | CPE - Cray Programming Environment

- HPE vendor provides a comprehensive software stack of compilers and libraries to cover a wide range of user's needs in building, debugging and profiling applications

- Cray Programming Environment (CPE) includes the following compilers and MPI libraries for C/C++/Fortran

  - Cray compiler collection based on Clang
  - GNU
  - NVCC
  - Intel and AMD (AOCC) for AMD Zen2/3 nodes

# User Experience | UENV - Programming User Environment

- CSCS programming user environment (UENV) based on Spack
  - software stacks built from a simple recipe
  - stored and versioned in GIT
  - each UENV is a single SquashFS image
  - reproducible builds

- Benefits of UENV approach
  - users can define, update and deploy their own stacks
  - software stacks can be added and updated with no intervention from system admins
  - smaller stacks targeted towards specific user communities reduce maintenance and support

- CSCS provides
  - support and maintenance of the tooling used to define and build the stacks and their integration on Alps
  - CI/CD pipelines for automated testing and deployment of end-user UENV
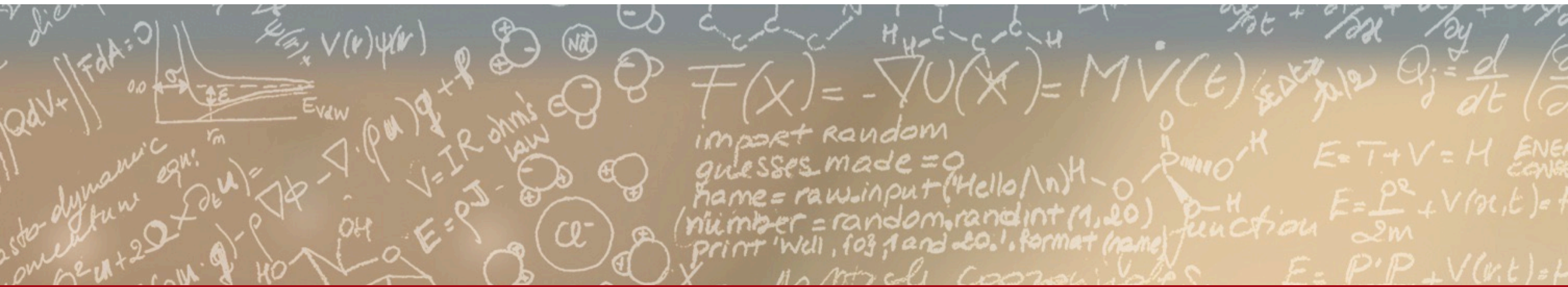
cscs

ETH zürich

# User Experience | Visualisation

- CSCS provides support for installation and usage of visualisation software tools
  - ParaView
  - Visit

- Extensive documentation provided on CSCS user portal


- The support for remote access to such tools is a cornerstone of our visualization service
  - usage of ParaView in client-server mode tested from Canada, Australia, Northern European countries and more
- Highly reliable service
  - tools have built-in support for adapting performance based on available bandwidth

# Summary and Conclusions

- Alps towards improving standard HPC through cloud-native approach
  - VLANs as lifeblood
  - multi-tenancy for an increase variety case studies

- Multiple geo-distributed infrastructure

- Heterogeneous infrastructure
  → NVIDIA and AMD GPU, x86, ARM, GH200

- Managed by a micro service architecture control plane
  → CSM/Shasta

- Slingshot → performance and zone segregation

- Common infrastructure, IaaS and PaaS + DevOps technologies reduces time to solution

- Platform per group of user communities
  - vCluster technology as logical unit
    → convergence of Cloud and HPC
  - e.g. CSCS Tier-2 Grid Site

- **Infrastructure as Code** based implementation of Alps vClusters and of Rancher-managed K8s-clusters
  - scale the infrastructure dynamically and according to the changing requirements of the partners

- Science as a Service concept with innovative resource access

- 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

# Thank you for your attention!

*Contact: riccardo.dimaria@cscs.ch*

cscsch

# Backup Slides

# CSCS Facility – Water Cooling System



**Lugano Lake**

-13m
25°C (max)

~75kW regain (3x)

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

**LD Islands Building (7MW)**

-45m
6°C

~250kW (3x)

9°C    17°C    19°C    27°C

CSCS

- pipeline length: 2.8 km
- height: 30 m
- max. flow rate: 760 l/s

CSCS

ETH *zürich*

# Evolution



SVILUPPO DELLA POTENZA DI CALCOLO PRESSO L'ETH ZURIGO E IL CSCS

# Piz Daint

- Piz Daint is a Cray XC 40/50 with 7517 compute nodes

- Commissioned in 2012
  - major upgrade/extension in 2016

- CSCS flagship system… since then, 8 years and counting!

- Lived through
  - >111M MC node-hours
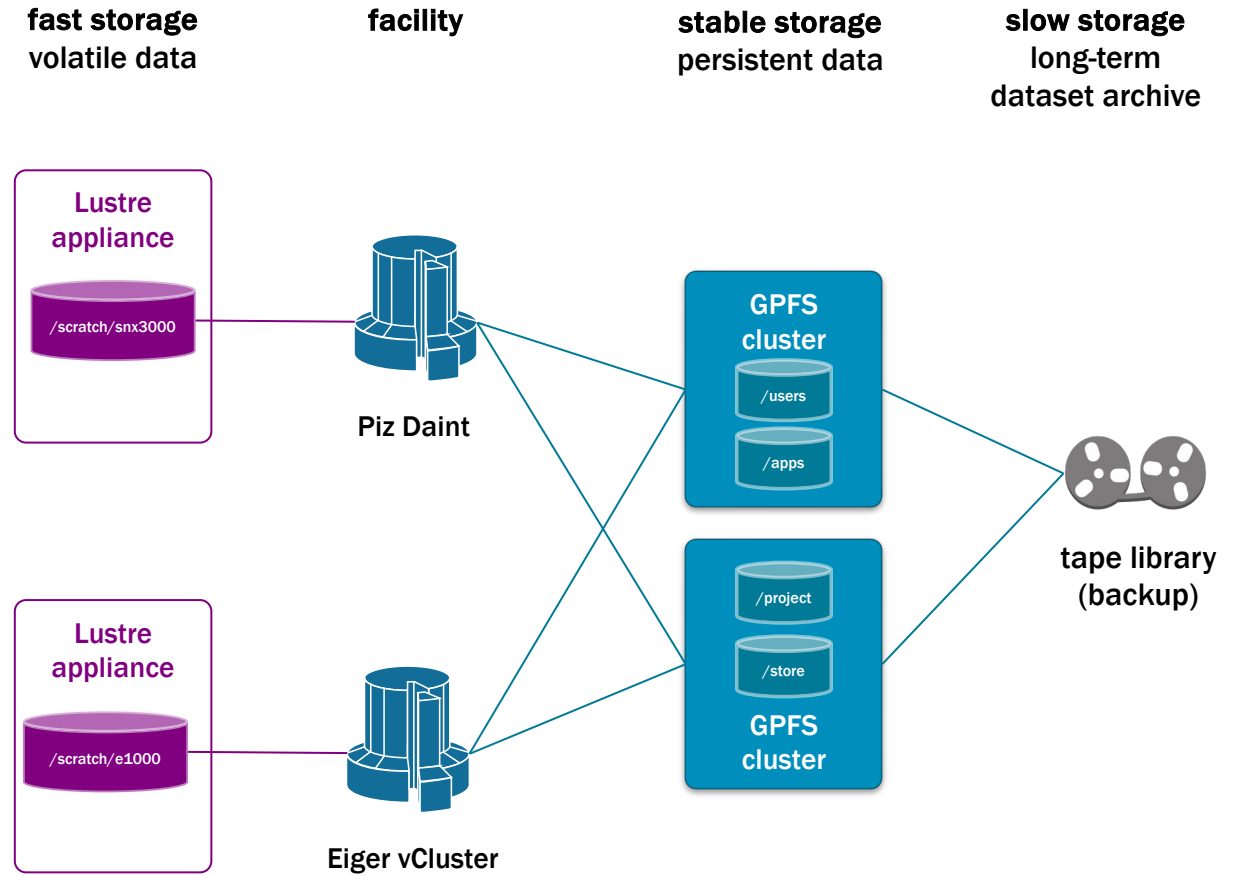  - >349M GPU node-hours
  - ~2800 users
  - ~50M user jobs

**Specifications**

| | |
|---|---|
| Model | Cray XC40/XC50 |
| XC50 Compute Nodes | Intel® Xeon® E5-2690 v3 @ 2.60GHz (12 cores, 64GB RAM) and NVIDIA® Tesla® P100 16GB – 5704 Nodes |
| XC40 Compute Nodes | Two Intel® Xeon® E5-2695 v4 @ 2.10GHz (2 x 18 cores, 64/128 GB RAM) – 1813 Nodes |
| Login Nodes | Intel® Xeon® CPU E5-2650 v3 @ 2.30GHz (10 cores, 256 GB RAM) |
| Interconnect Configuration | Aries routing and communications ASIC, and Dragonfly network topology |
| Scratch capacity | 8.8 PB |



cscs

ETH zürich

# Storage in the Piz Daint Era

- Different storage backends depending on
  - use-case
  - performance envelope
  - technology constraints
  - historical reasons

- Persistent data and tape access stored on IBM Spectrum Scale (GPFS) clusters and a myriad of backends (SAS or FC disks, SAN, JBODs, etc.)

- Volatile data with fast access patterns on Lustre appliances



**fast storage**
volatile data

**facility**

**stable storage**
persistent data

**slow storage**
long-term
dataset archive

Lustre appliance
/scratch/snx3000

Piz Daint

Lustre appliance
/scratch/e1000

Eiger vCluster

GPFS cluster
/users
/apps

/project
/store

GPFS cluster

tape library
(backup)

cscs

**ETH** zürich

# Storage in the Piz Daint Era

- Sonexion 3000 (Daint's scratch)
  - 8.8 PiB Cray Sonexion 3000
  - spinning disks
  - raw performance:
    - 112 GB/s write | 125 GB/s read
  - 1640 HDDs (8TB each)
  - 2 metadata servers (20x 800GB SSD each)
  - 3 full racks
  - InfiniBand FDR

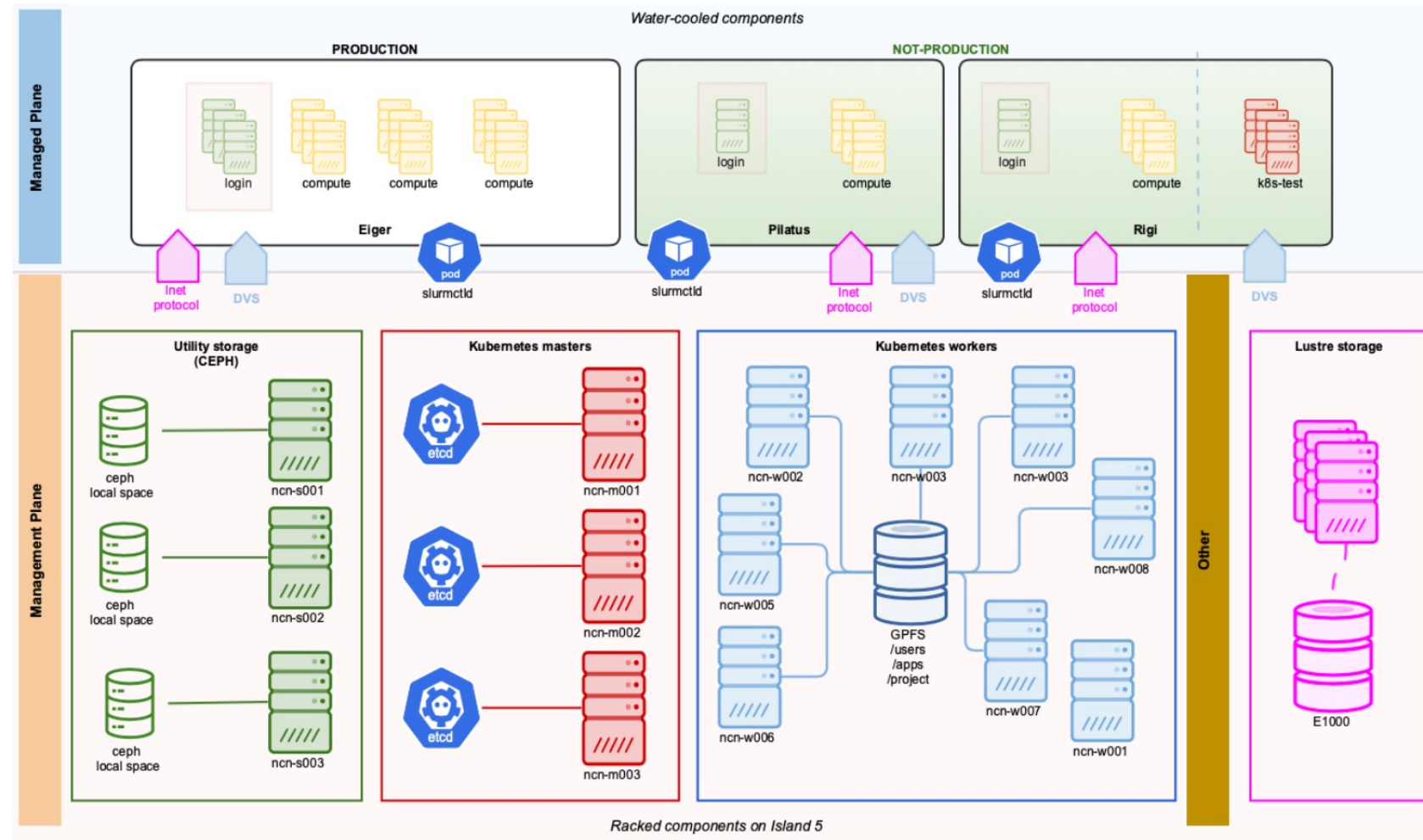comparable performance

- Alpstor
  - 10 PiB HPE ClusterStor E1000
  - spinning disks
  - raw performance:
    - 120 GB/s
  - 240 HDDs (16TB each)
  - 2 metadata servers
  - ~2 half racks
  - Slingshot 10

# 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

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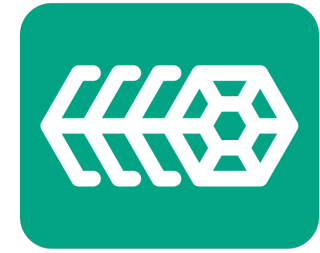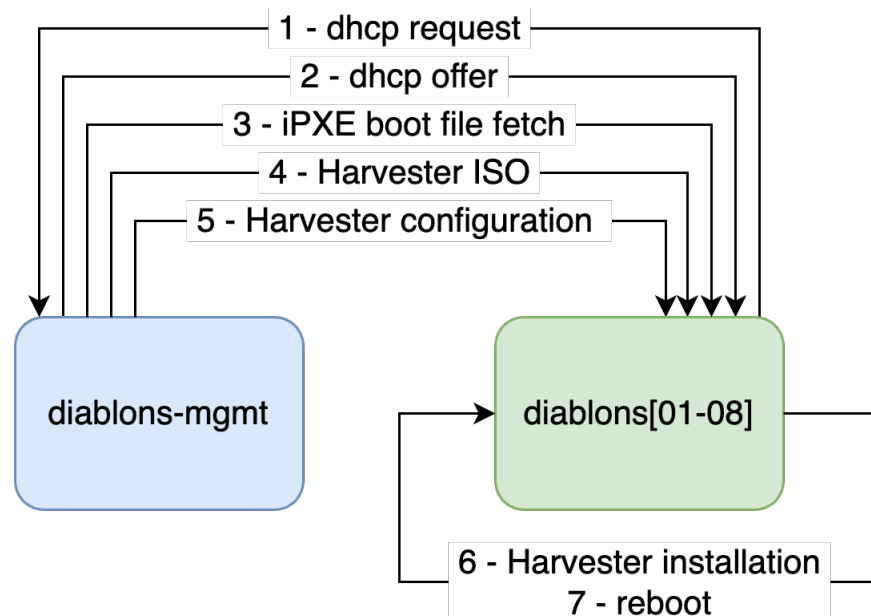
- 3 dedicated servers in HA
  - Intel dual-socket 12-core 128 GB RAM
  - provisioned with Metal-as-a-Service (MaaS) by Canonical
  - Rancher installed via RKE2 through Ansible

cscs

**ETH**zürich

# 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:



- 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)
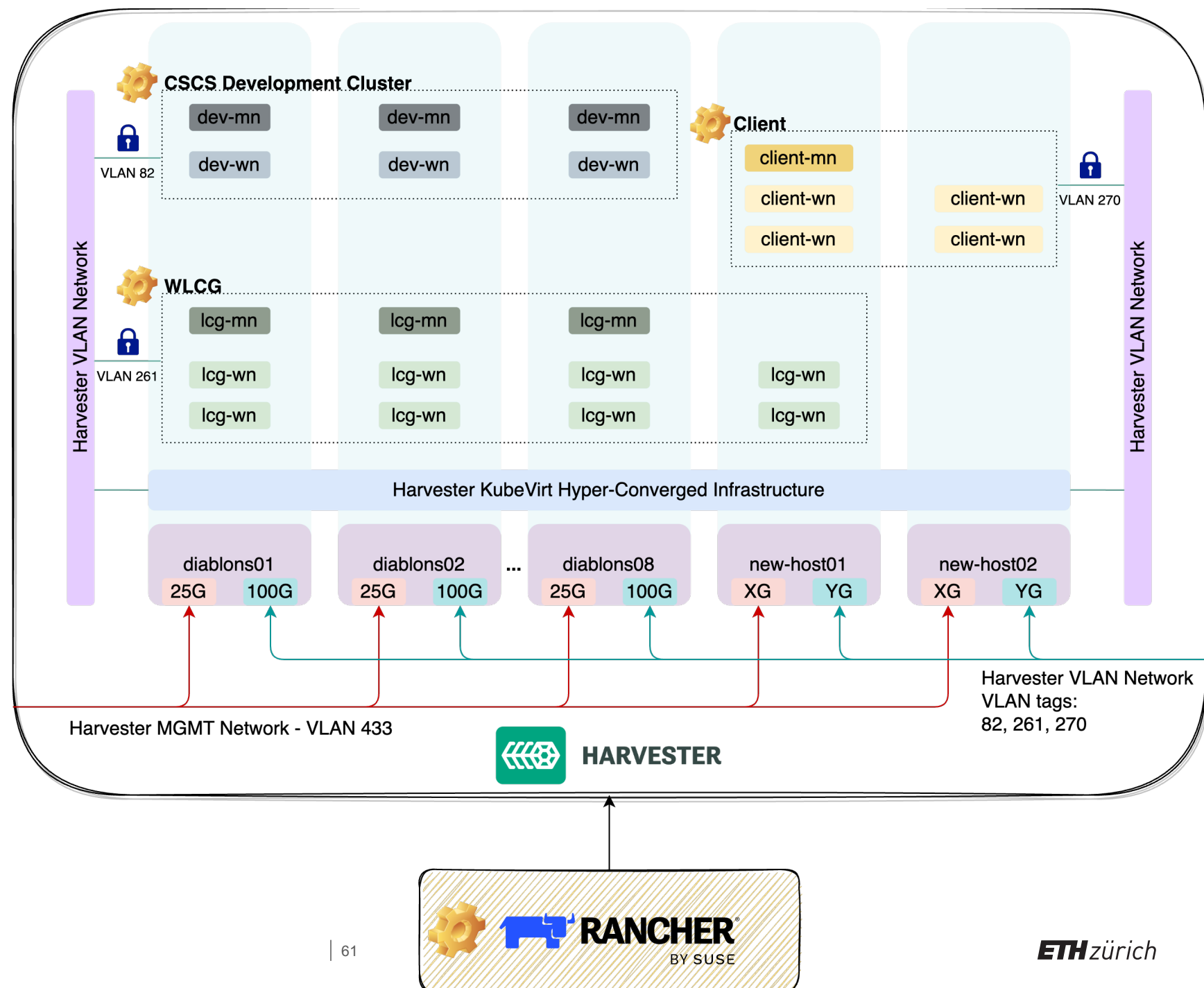  - flexibility to scale up physical cluster

# 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

**+20/25%**

### 2023

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

**AMD EPYC Rome → HS06/CPU = 22.46**

- ❖ 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
  - o 128 cores (256 CPUs), 256 GB RAM
  - o "Mont Fort" cluster
  - o **4 ARC-CEs**
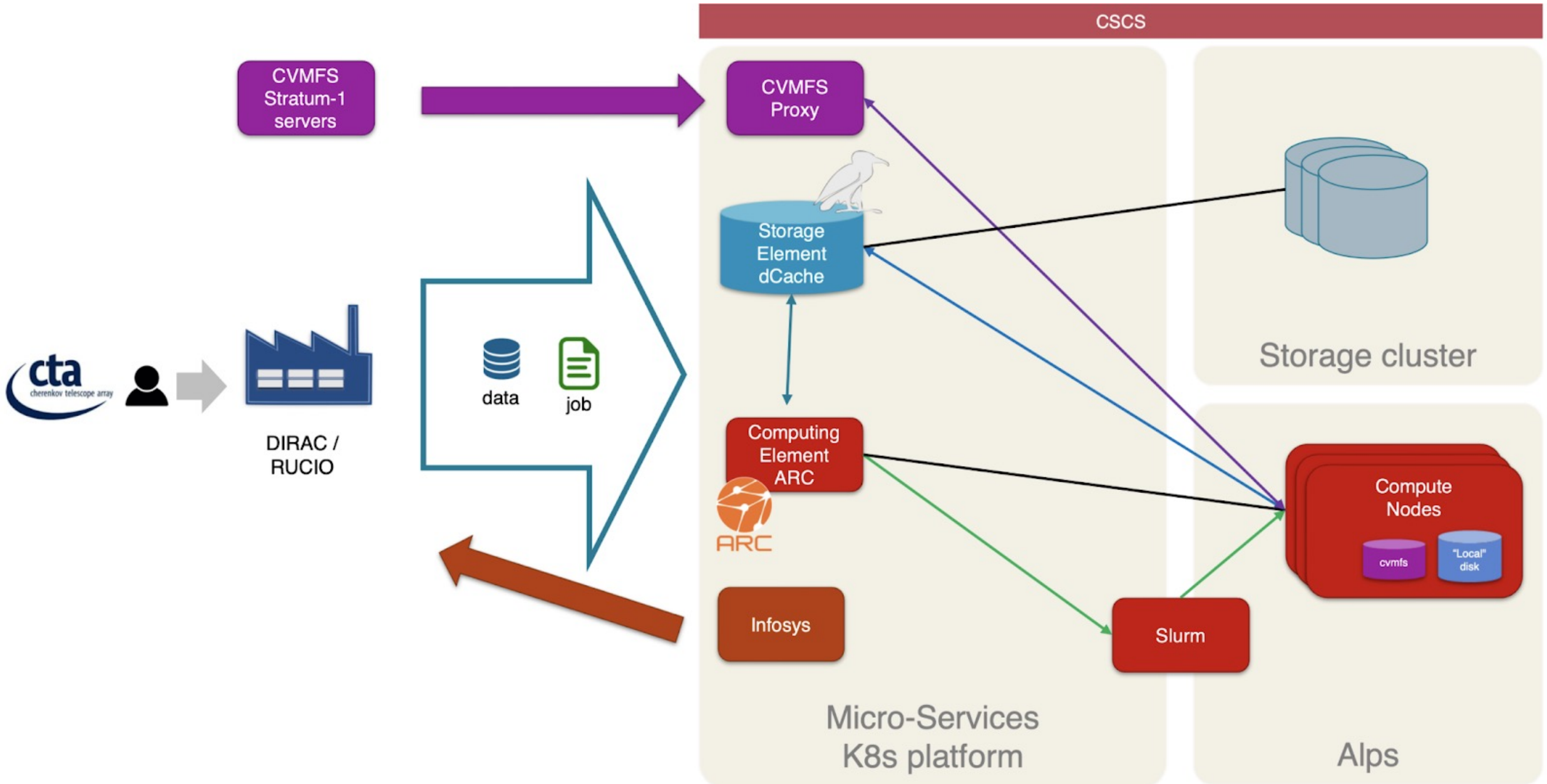- ❖ +4 nodes for dev/tds instance
  - o "Mont Gele" cluster, 1 ARC-CE
- ❖ Production CE
  - o 300 TB shared CephFS NVMe
  - o 4 TB local RBD NVMe per node
  - o 64 GB CVMFS cache RBD NVMe per node

# 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

- ~2 year in production

- dCache pool services run as K8s pods

- Pods mount Ceph RBD volumes through Kubernetes CSI

cscs

data path
(VLAN 15)

Incoming Request

Internet

K8s Worker Nodes

doors

xrootd

GridFTP

SRM

webdav

dcap

NGINX

Ingress Controller

MetalLB

core domains

dCache pools

MetalLB

K8s PVC on Ceph

postgresql

CEPH

| 64

ETH zürich