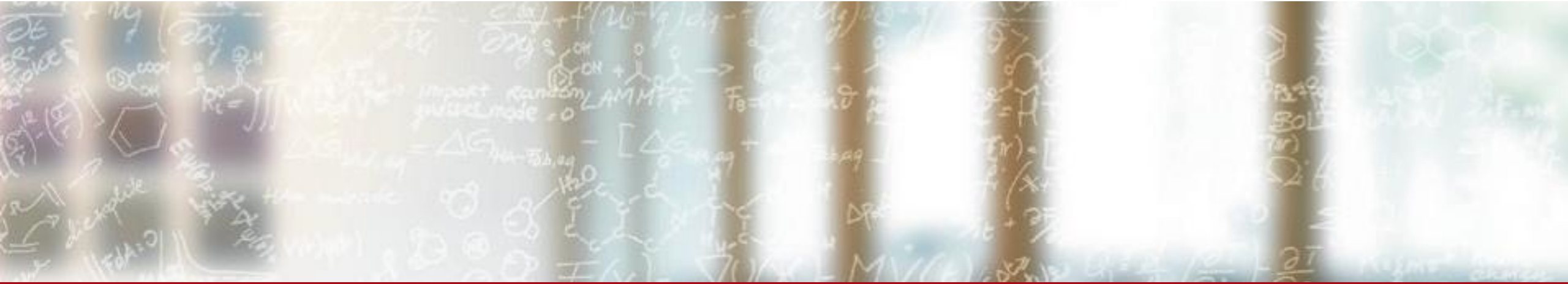




CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich



Versatile software-defined scientific clusters on CSCS' flagship infrastructure Alps

SKACH Spring Meeting 2024

Miguel Gila, CSCS

June 10, 2024



CSCS

ETH zürich



CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich

A bit of history about our supercomputers

Piz Daint

- Piz Daint is a Cray XC 40/50 with 7517 compute nodes
- It was commissioned in 2012 with a major upgrade/extension in 2016
- It's been our flagship system... since then, 8 years and counting!
- So far it has lived thru a lot of things:
 - +100M MC node-hours
 - +400M GPU node-hours
 - ~2800 users
 - ~60M user jobs

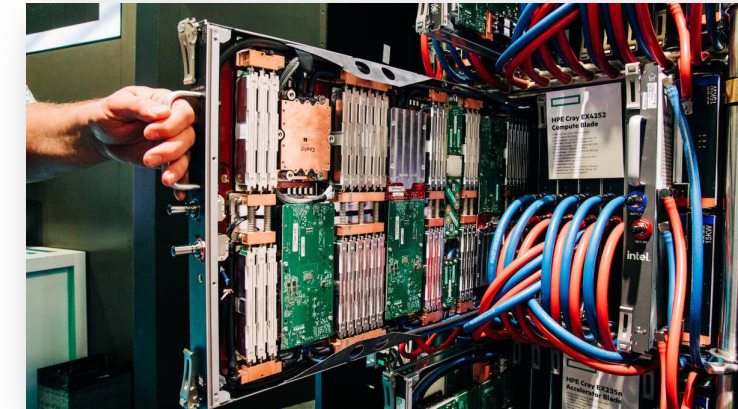
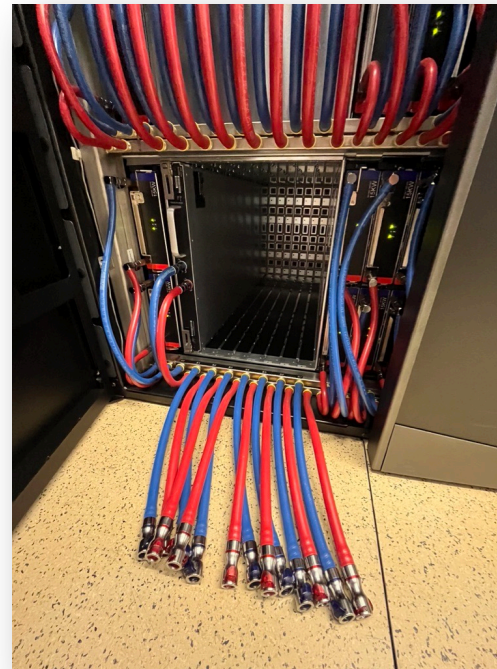
Specifications

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



Alps

- Alps is an HPE Cray EX supercomputer meant to be our new flagship infrastructure
- Multi-phase installation started in 2020
- Some specs
 - 1024x MC nodes (AMD Rome 7742)
256/512GB RAM
 - 144x nVIDIA A100 GPU nodes
 - 32x AMD MI250x GPU nodes
 - ~10k of GraceHopper GPUs
- Slingshot network
- Two available zones (HA, non-HA)
- 100% liquid cooled (to the chips themselves)



Water cooled blades

Installation



Big delivery

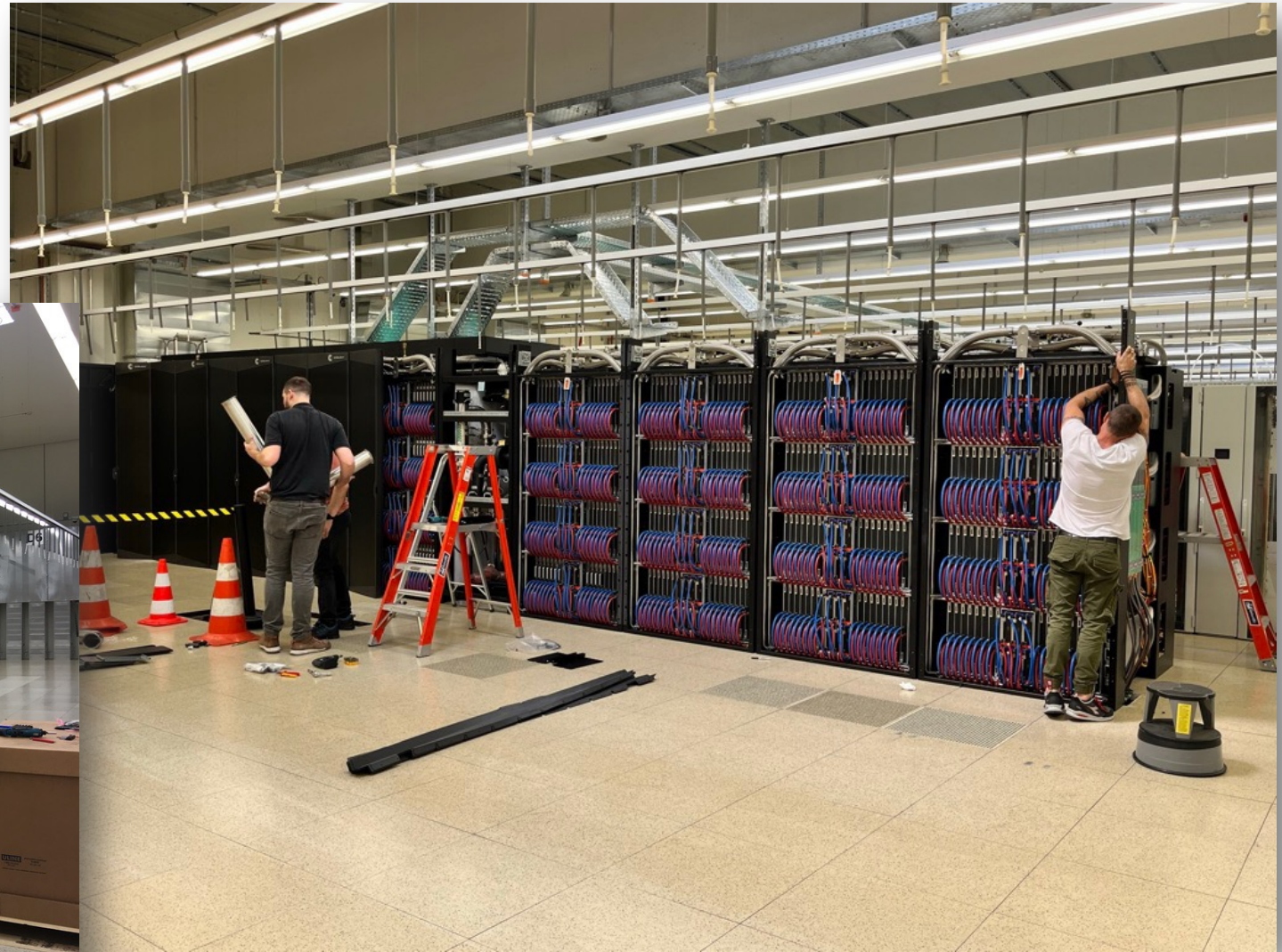


Moving the first racks

Installation



Some boxes



Front: water cooling

More pictures



Front: HA zone (compute + mgmt.)

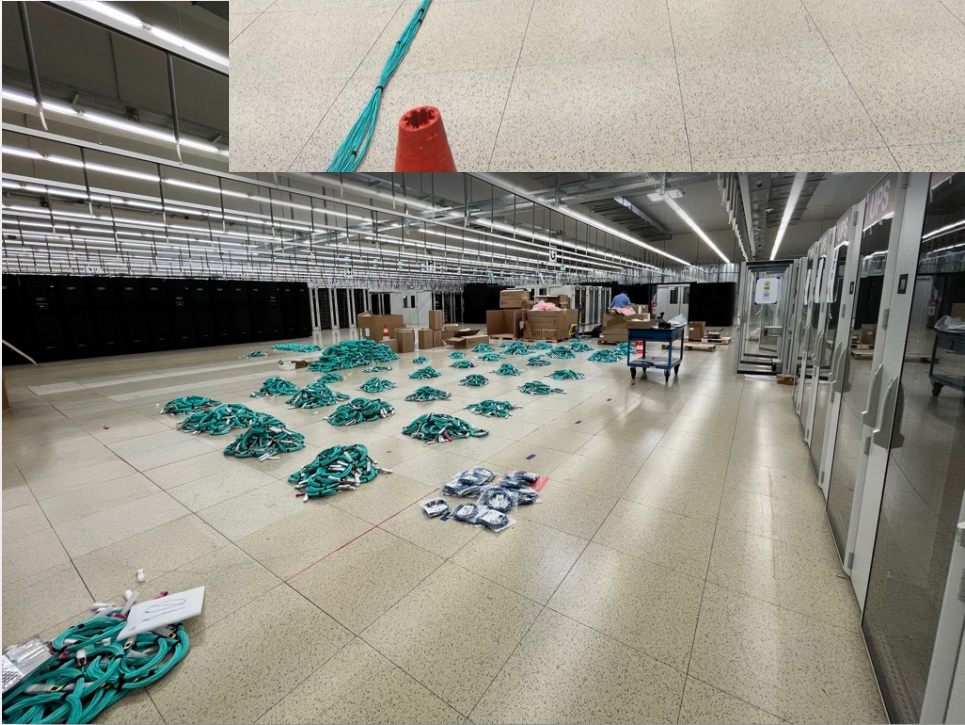


Back: Slingshot network

What's below, in front, and on top



Internal pipes and doors and panels



Some Slingshot fiber cables



Our mission

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.

- To achieve this, 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
- But, as the numbers of science domains and projects grows
 - 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)



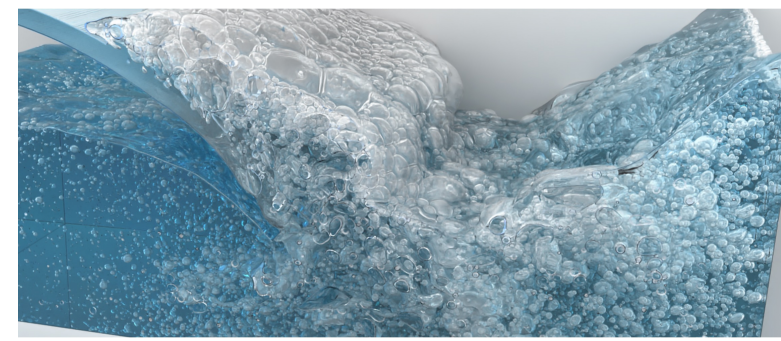
CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich

HPC and Cloud convergence

HPC design principle – towards science



- Resources – bare metal
 - Compute nodes, high-speed network, high-speed large storage optimized for write
 - Exclusive access to compute resources with accelerators
 - Shared network and storage
- Services
 - Resource manager – node allocations with queues – usually large jobs (nodes and time)
 - Fully vertically integrated software stack from the kernel to the high-level library
 - Pre-defined user environments with limited choices
- Access
 - Design for interactive shell with SSH jumping through a chain of login nodes
 - Posix IAM (limited authentication and authorisation)

High performance → fully integration of stack → limited set of services

Cloud design principle – towards enterprise

- Resources – virtualisation
 - vResources: vCPU, vGPU – uses virtual images/containers
 - Bare metal nodes runs many different vResources
 - Multi tenancy for compute, network and storage
 - Price model cheap for vResource to very expensive for bare metal
- Services
 - DIY infrastructure - define the configuration of your infrastructure
 - Bring your own service
 - Extensive offering of pre-configured services
- Access
 - UI in your browser, interactive shell in your browser or with SSH
 - Web IAM: complex protocol for authentication and authorisation



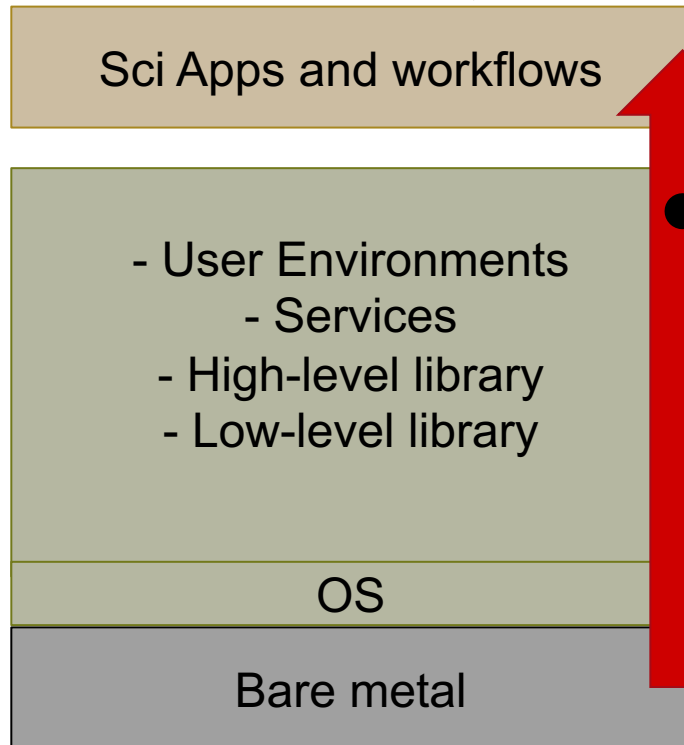
Virtualization at scale → high flexibility → limited performance

HPC and Cloud concepts to enable Science

Aim at Science as a Service

SSH / POSIX IAM

Defined by Business service
Cloud native IAM

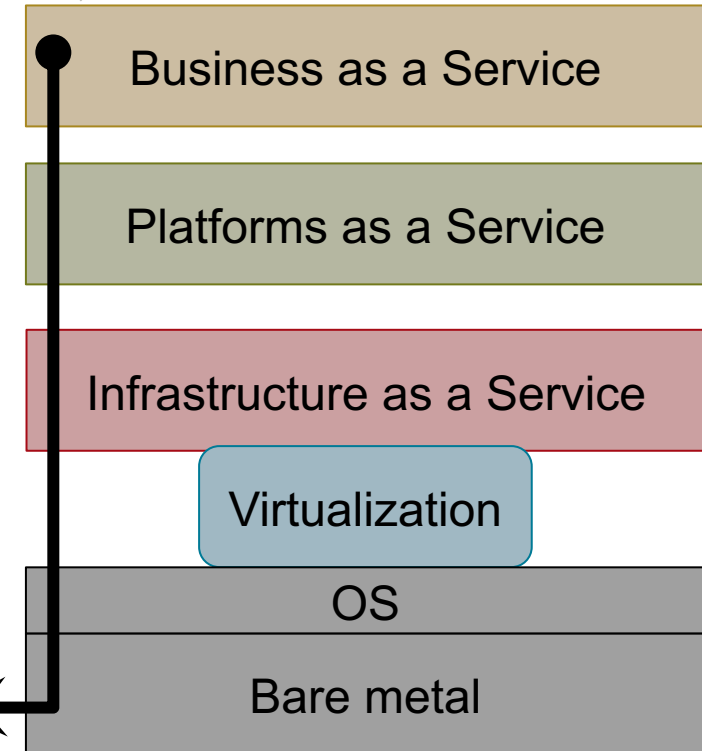


HPC

1. Performance
 1. OS drivers
 2. Low-level library
 3. High-level library
 4. User Env for HPC

2. Layer flexibility
 1. Infra as Code
 2. Multi-tenancy: QoS and isolation

3. Software as a Service
 1. User Env and Apps
 2. Workflows and APIs
 3. Access controls



Cloud



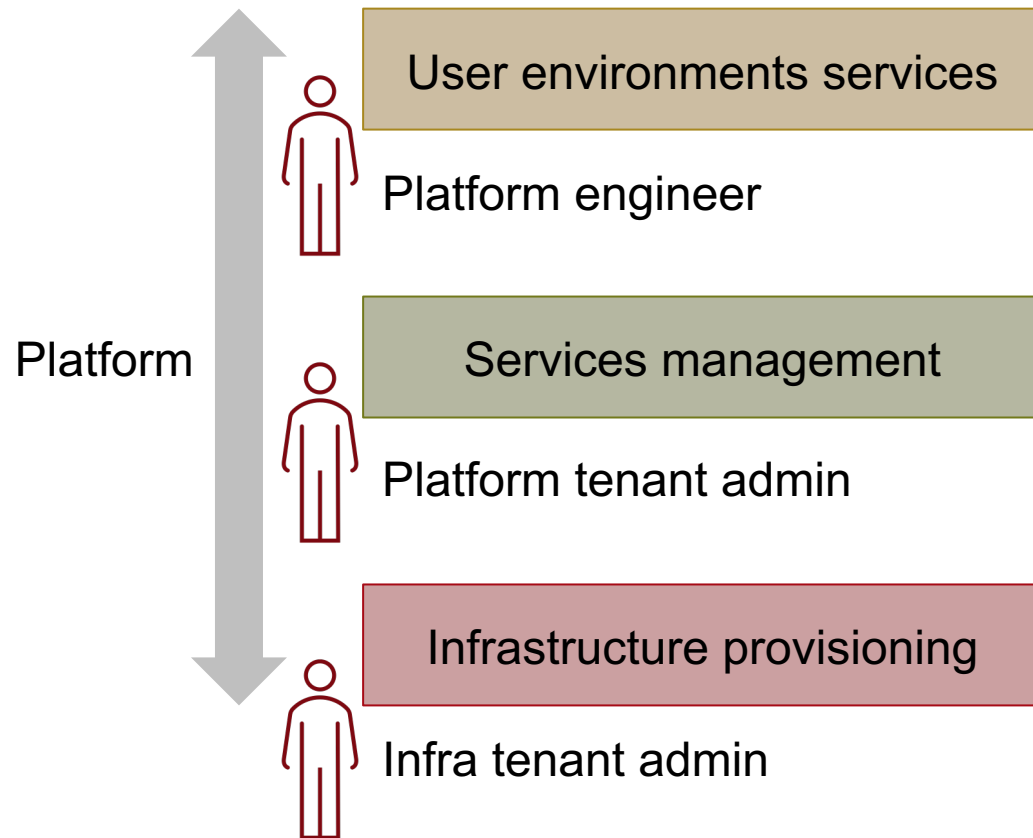
CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich

Versatile software-defined cluster (vCluster)

vCluster layers and tenant concept

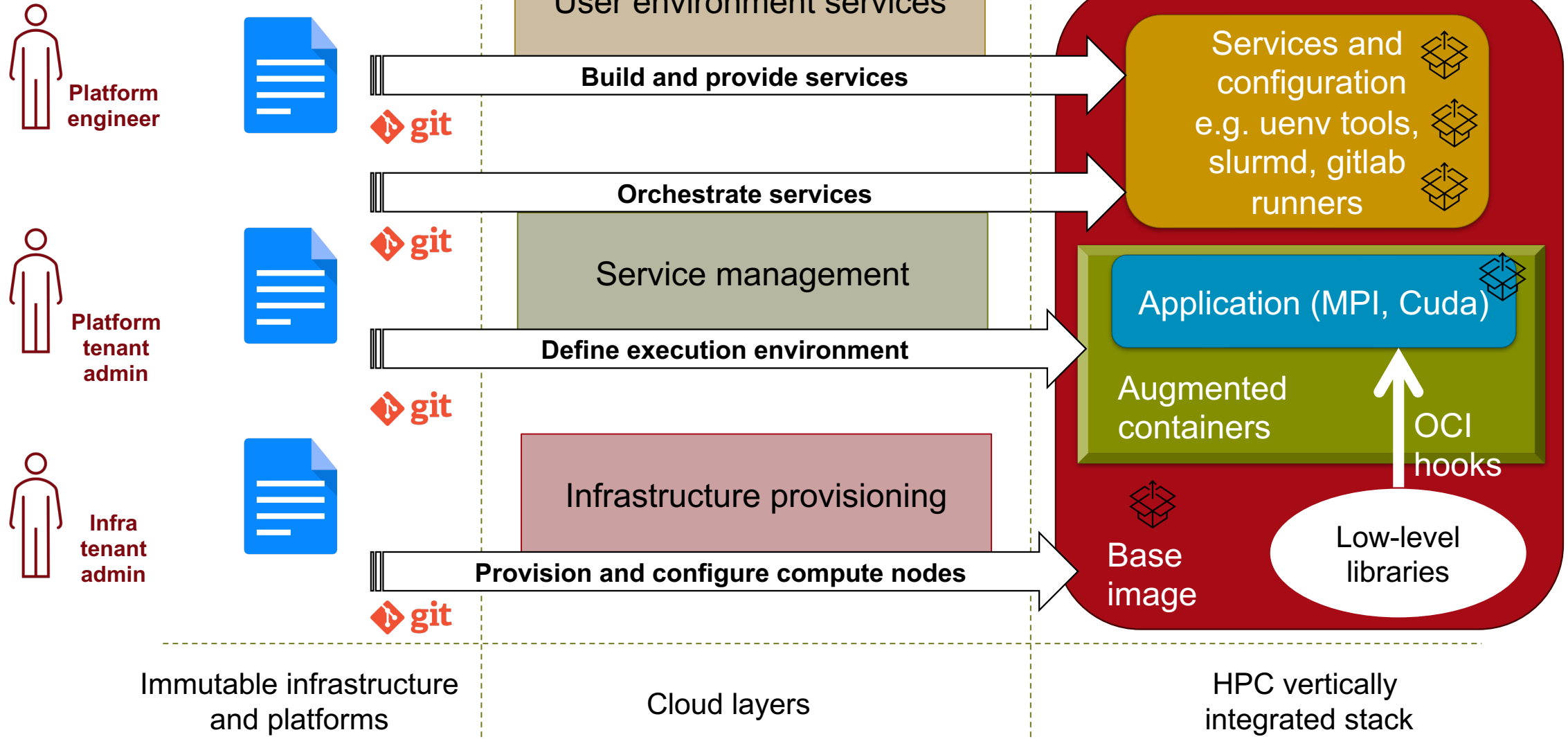


- User tailored environments
- Programmable resource access
- Scientific application build services

- Orchestration of platform services
- Execution environments
- Soft tenant (resource labelling)

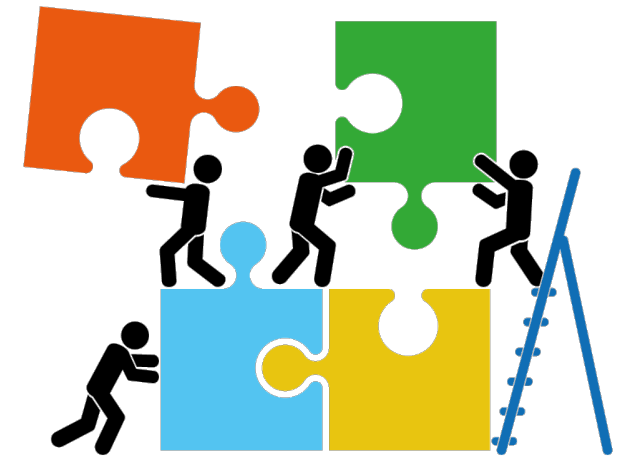
- Interface to the management plane
- Hard tenant (network segregation)

Version controlled textual description



Composability → Versatility

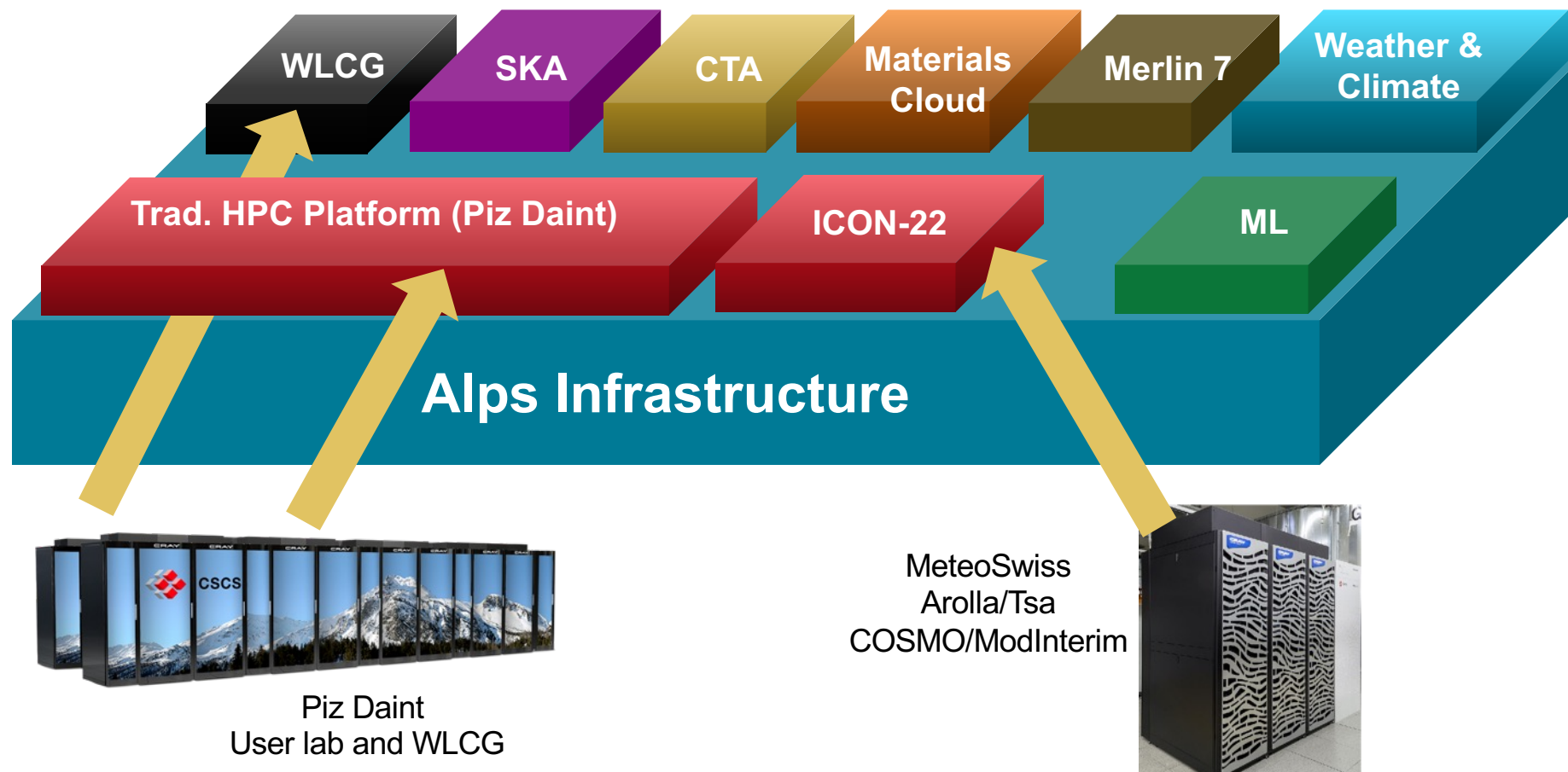
- This permits to build an HPC cluster using components developed and maintained by potentially different groups of people
- Each admin can choose which products to put where, and when
- This **composability** is the key to **versatility**
- Bridge gap between traditional HPC and cloud technologies
- Rendering Alps more adaptable and flexible



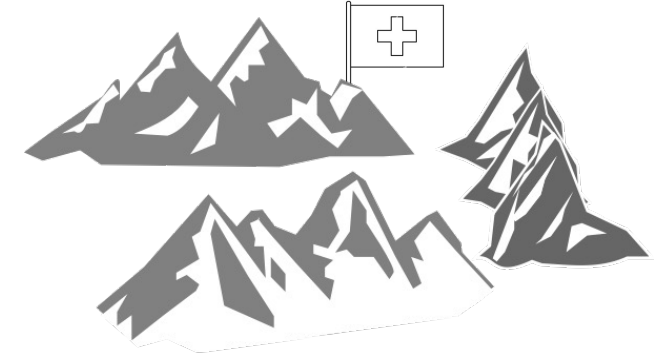
vClusters on Alps

Software-defined infrastructure

Consolidation of platforms on a shared infrastructure



Some vClusters in operation on Alps

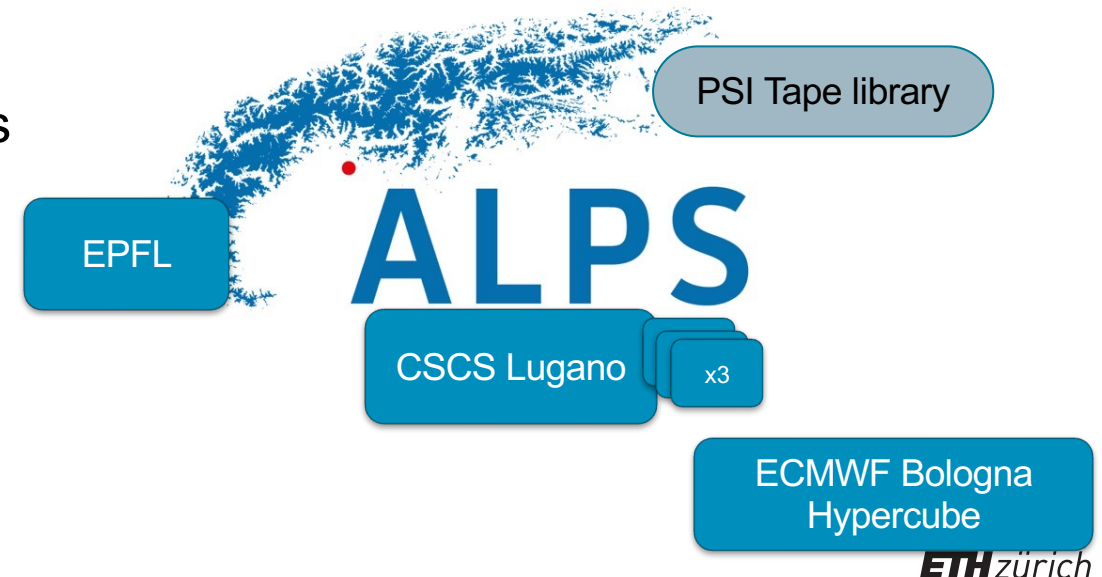


Platform	vCluster name	Scope
MCH	Balfrin	R&D
	Tasna	Production
HPC Platform	Eiger	Production
	Pilatus	Staging
	Rigi	R&D
WLCG Platform	Fort	Production
	Gele	R&D
	Noir	Production (CTA)
AI/ML Platform	Clariden	AI/ML
Testbeds	Rosa	K8s testing
	Adula	HW validation
	Zinal	Internal development
	Bar	Internal development
	Hohgant	Internal development
PSI	PSI-dev	R&D
	PSI-tds	Pre-production

vCluster Mobility

- On Cloud providers for company business scalability
 - Need to understand infrastructure provisioning
 - Common abstraction layer
- On redundant infrastructure at the same site for ensuring service availability during maintenance
- On redundant infrastructure at different sites for geo-redundancy service
 - Synchronisation of control planes across sites
 - Synchronisation of managed planes across sites
 - Synchronisation of data

AlpsE at EPFL



Wrap-up

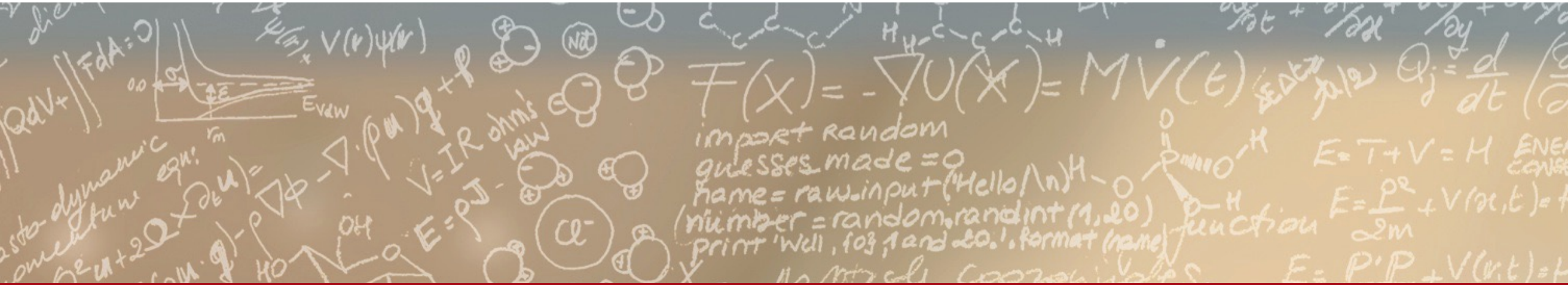
- vCluster is a set of technologies
 - User environments
 - Service management
 - Infrastructure provisioning
- Enable abstractions with performance
 - Soft and Hard tenancy
 - Augmented containers using OCI Hooks
- Offers multi-tenancy on a single large HPC infrastructure
 - HPC and Cloud convergence



CSCS

Centro Svizzero di Calcolo Scientifico
Swiss National Supercomputing Centre

ETH zürich



Thank you for your attention.