



From terminal to platform



Challenges on building up a Cloud Science Platform in an HPC world at DZA

Overview

/2

In this presentation we want to show the technical challenges we are facing when building up a modern science platform for astrophysical projects and potential solutions from the perspective of computer science.

We will use the national German deSRC node from the SRCNet project as an example to illustrate hard challenges.

Agenda:

- SRCNet Intro and Status Quo
- HPC vs Cloud
- Options / solutions
- Summary

DZA – German centre for Astrophysics

13

- Currently being founded as a country-wide institute to support the German astrophysics research and collaborate with already existing established institutes
- Enable German participation in large astronomical research projects and facilities
- Have a focus on Radio-, Gravitational-Wave and Time-astronomy
- Consists of 3 closely interlinked science pillars:
 - Astro physics
 - Data science
 - Technology
- Mission to catalyze scientific, economical and sociological developments in the former coal mining area of Lusatia
 - “From coal to silicon”

DZA – Some example projects

14

- ESA Gaia DR4
 - ESA's billion star surveyor
 - https://www.esa.int/Science_Exploration/Space_Science/Gaia
 - DZA will aim to host DR4 as a partner institute
 - Hosting Type:
 - Catalogue data with various search types
 - Think: Storage + Web/TAP/IVOA access
- SKAO – SRCNet
 - <https://www.skao.int/en>
 - DZA will aim to support setup/installation of the official “deSRC” node
 - Hosting Type:
 - “Science Platform”
 - Think: Storage + Interactive/Batch Compute + Web/TAP/IVOA access

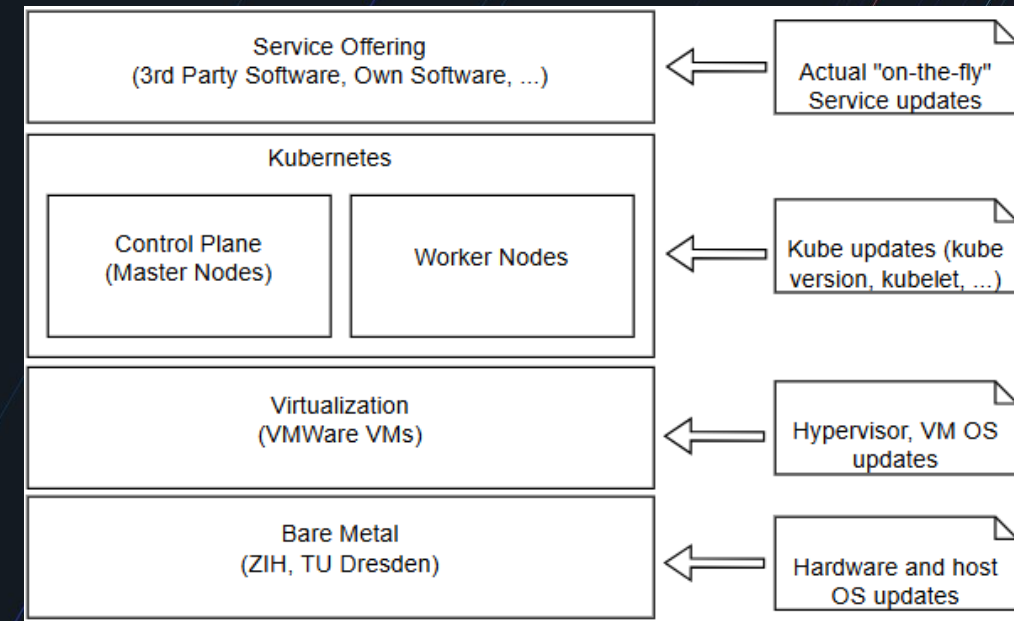
In the talk, we will focus on the technical/engineering challenges for such a “Science Platform” on the example of SRCNet

SRCNet Intro and Status Quo

SRCNet basic overview

/ 6

- Architecture:
 - Distributed node architecture to spread out redundant data storage and compute
 - Follow the example of CERN
- Software Stack:
 - Modern service like development program to follow the industry state of the art
 - Kubernetes (there is no official req. for Kube, but the developed artifacts (images/Helm charts/etc.) imply that req.
- Hardware Stack
 - Heterogenous IT infrastructure at each contributing country/institute



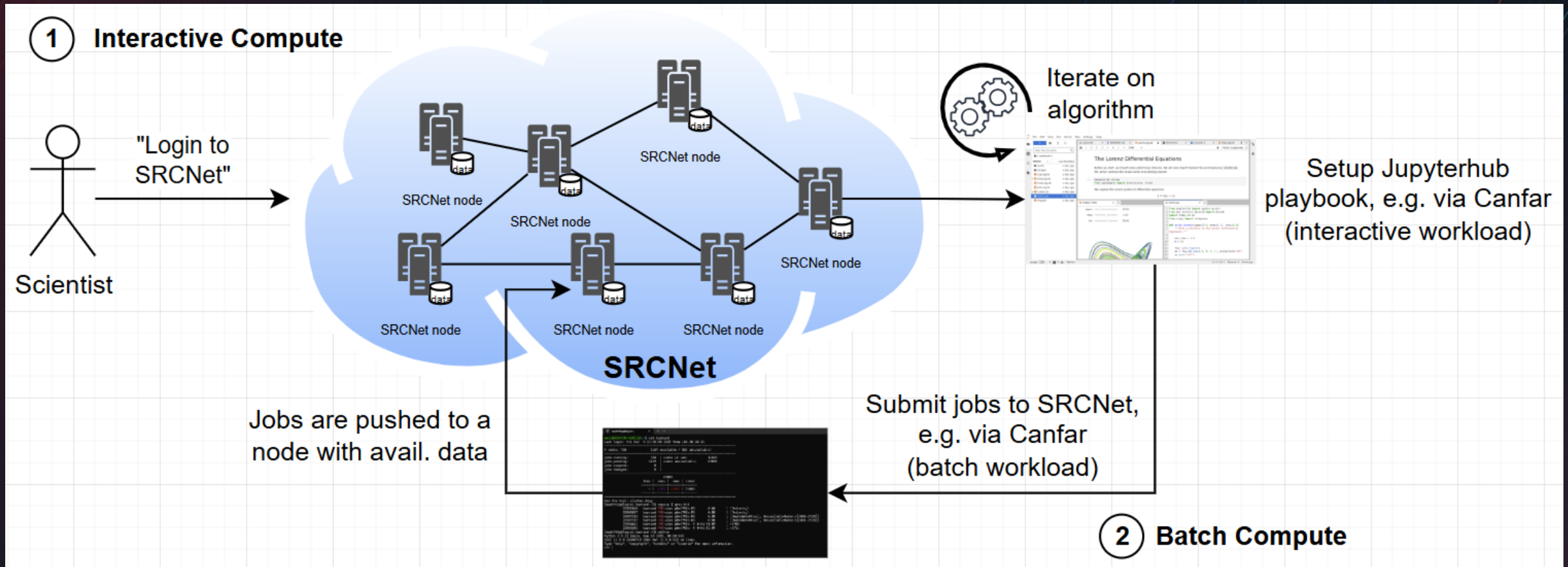
Opinion:

From the perspective of an Industry/Service background -> this is the right way to go

Assumed workflow on a science platform (HTTPS-based)

17

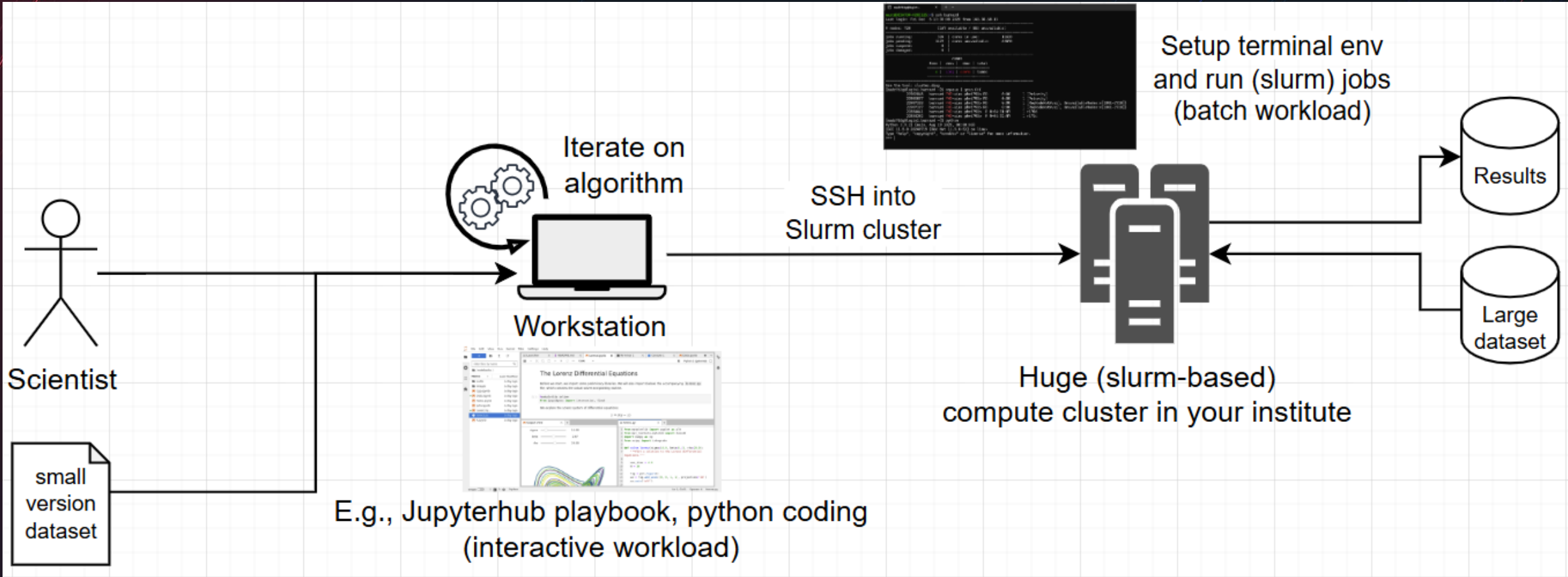
Or any variation of that



Status Quo: Workflow of a traditional science use-case (SSH-based)

/ 8

Or any variation of that

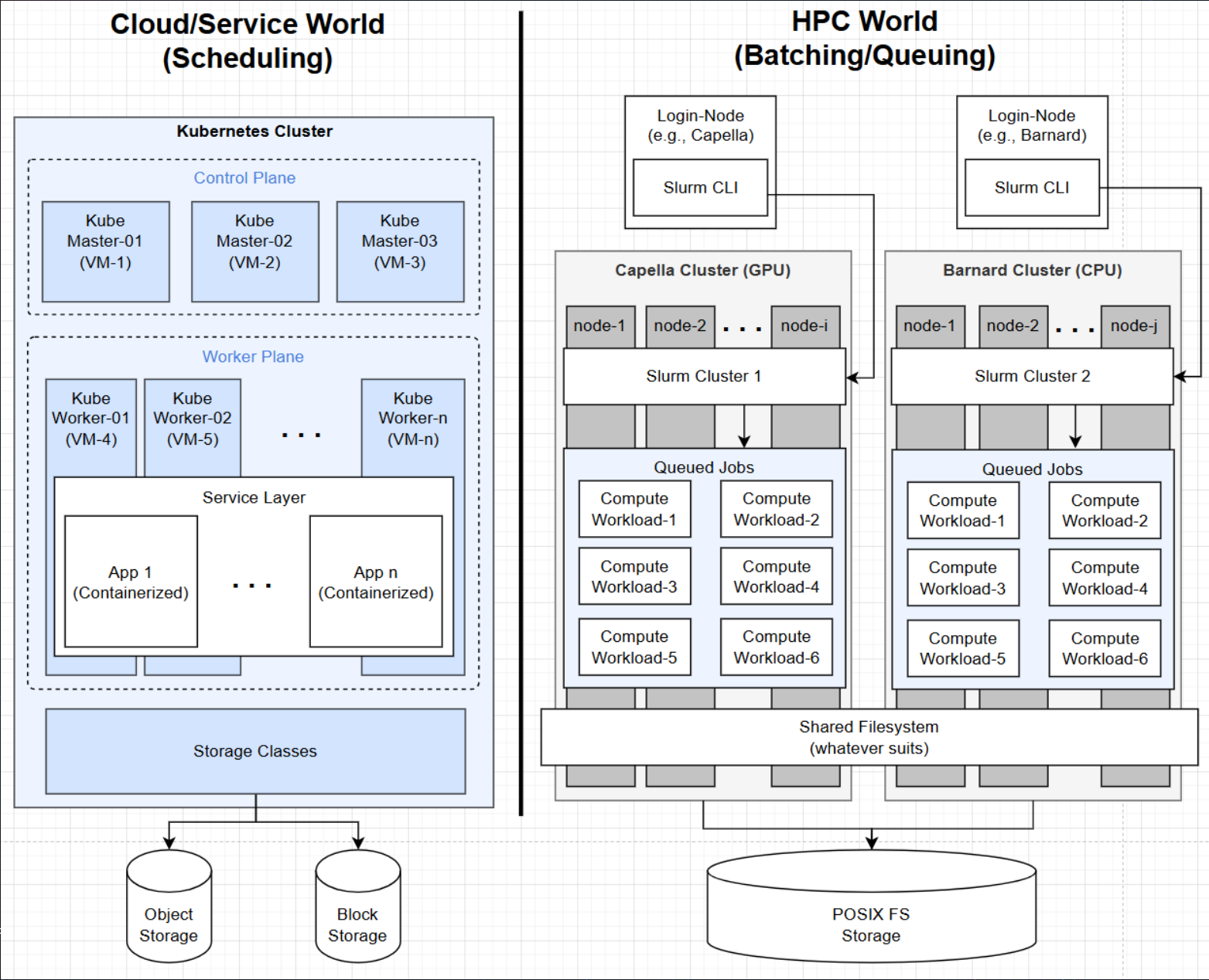


Opinion:

From an academic/batch-job perspective -> this has been an easy and reliable traditional approach

HPC vs Cloud

HPC vs Cloud - Comparison



HPC vs Cloud - Comparison

/ 11

Category	Slurm (HPC)	Kubernetes (Cloud)
Target use-case	Heavy, but “short-running” batch compute	Never failing, always available, always up-to-date end-user services
High Availability	Platform: not required Jobs: not required	Platform: required (9er rules) Services: required (9er rules)
Environment	Self-managed terminal env – scripting Singularity Containers	Entirely image/container-based
Automation	Mostly manual job triggers (can be scripted though)	GitOps oriented fully automated approach
Queuing/Batching	Native	Scheduling only (but Batching can be added, e.g. Kueue)
Backup/Recovery	Input/Result data: required Jobs: not required	Required for all services
Resource Sharing Model	Between all authorized parties (reservation possible)	Between all authorized parties
Main Adoption	Academia/Research	Industry

So what does that mean for “deSRC”

Options / Solutions

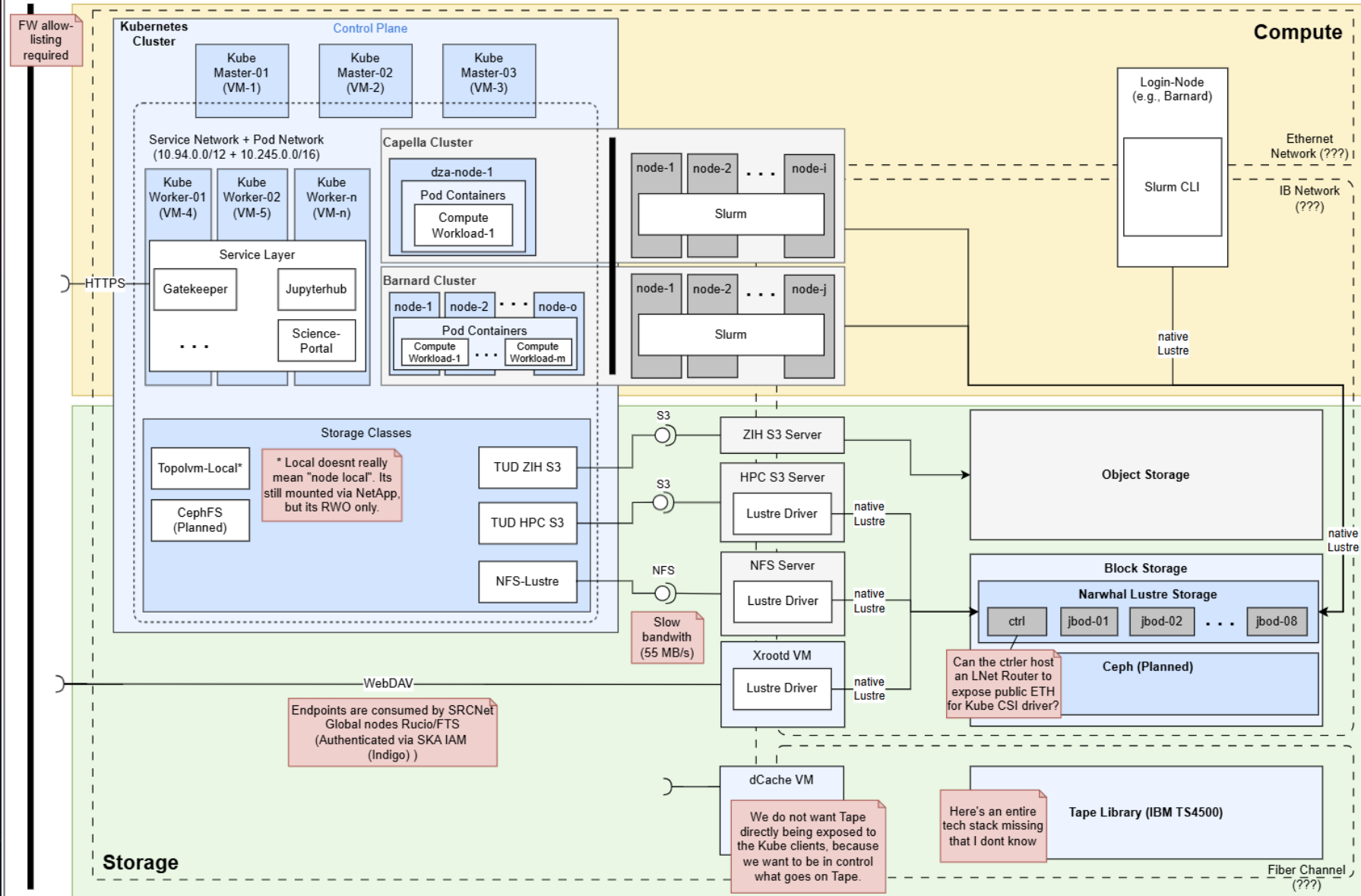
deSRC - constraints

/13

- DZA will host the “deSRC node” from the datacenter at TU Dresden
- All our resources are added into the existing Slurm clusters Barnard/Capella
 - The initial assumption was:
 - This is yet another traditional academic intense compute research project
 - Nobody expected Kubernetes 😊
- We have to somehow make it work!

System Overview - Bare Metal nodes

HPC World



Compute Gateway options - Interlink

/15

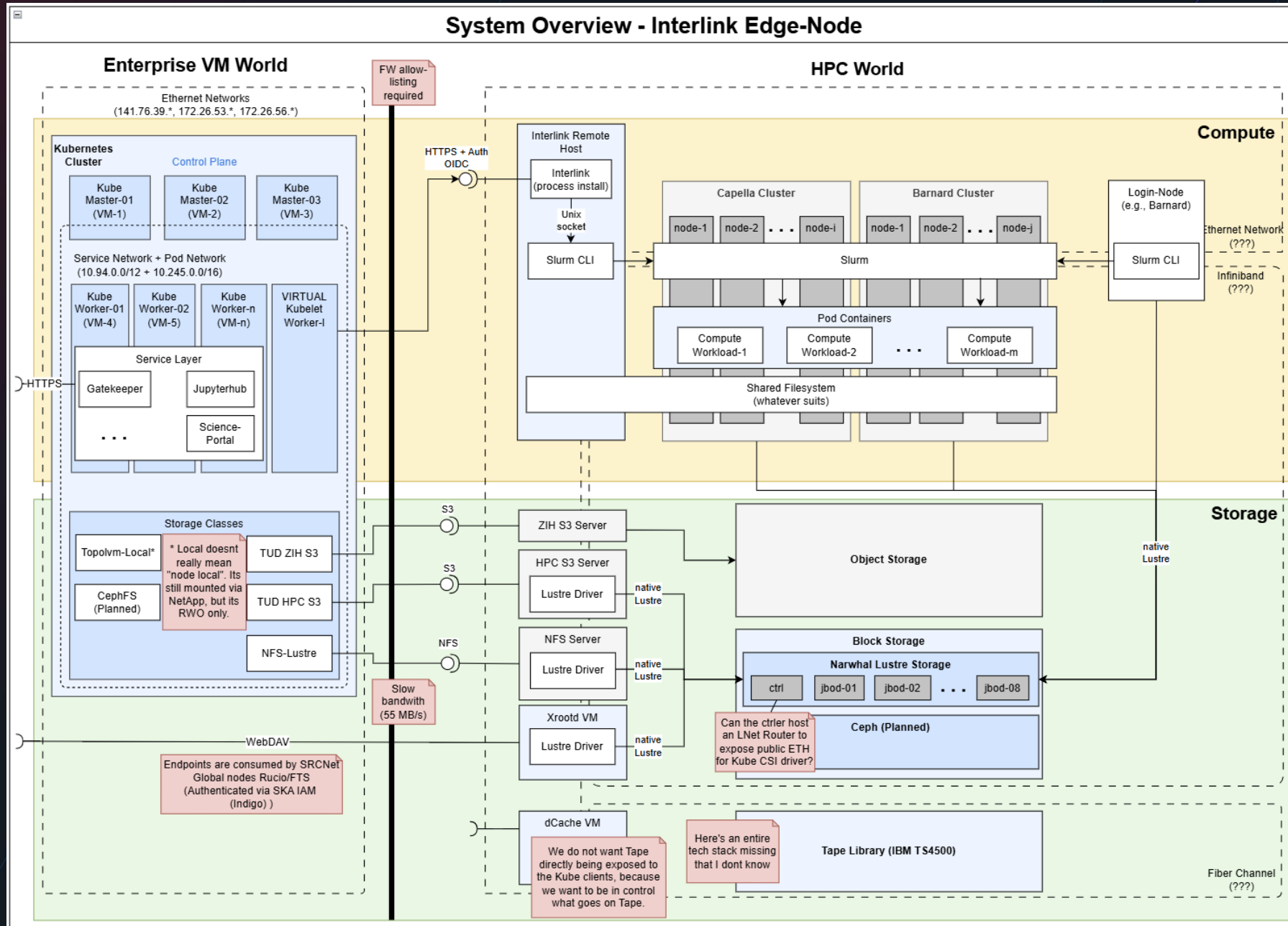
- Comes out of European Project “InterTwin” (<https://www.intertwin.eu/>)
- Allows to connect together the Kube Control Plane with the `slurm` CLI
- Exposes the HPC resources (cluster) as a new virtual Kubelet (= a node inside the Kube cluster)
- On its way to become a CNCF project
- Overview: <https://interlink-project.dev/>
- Code: <https://github.com/interlink-hq/interLink>

Compute Gateway options - Slinky

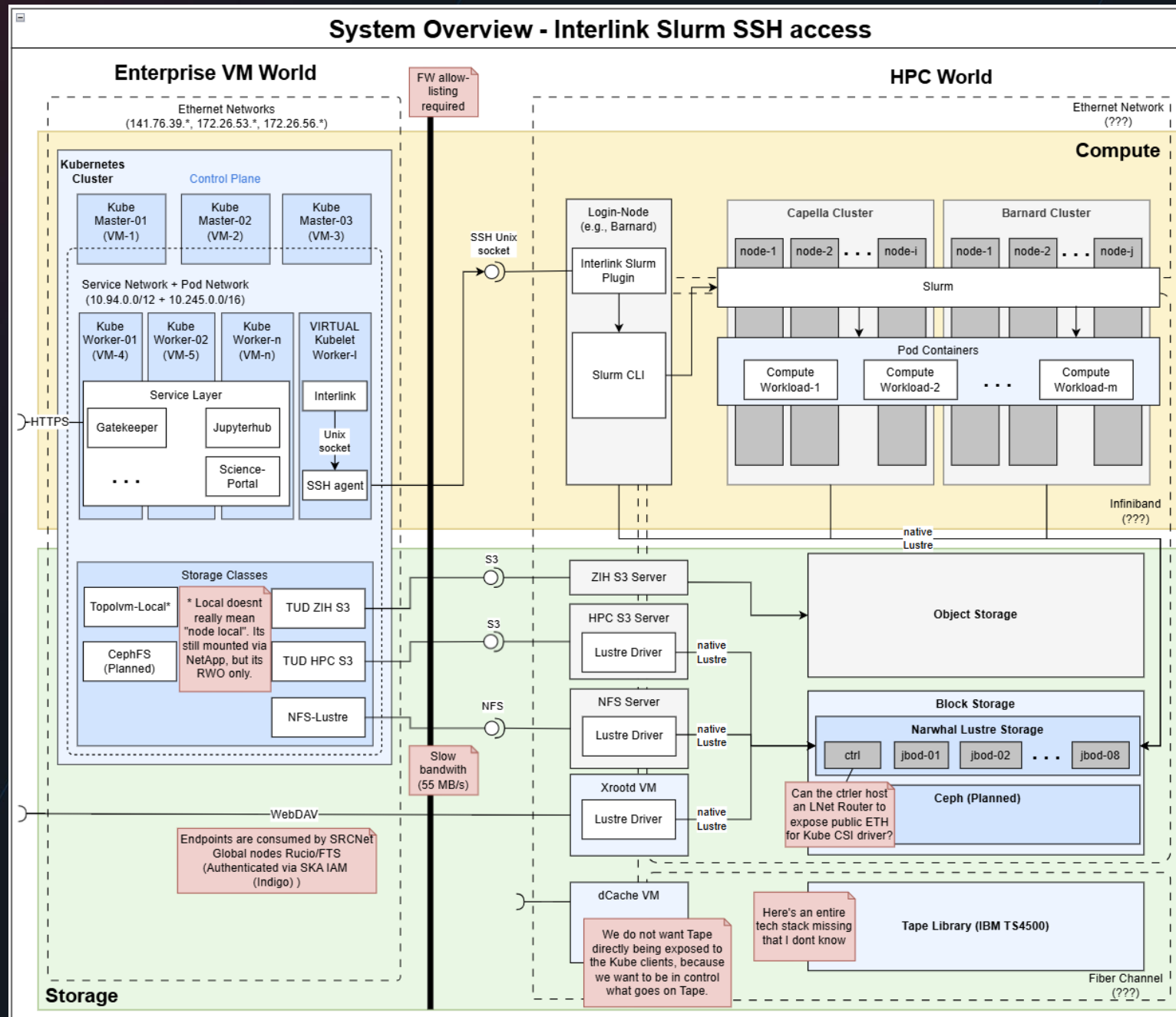
/16

- Developed by SchedMD (the Devs of Slurm)
- Allows to connect together the Kube Control Plane with the `slurmctld` daemon
- Postulates 3 different operation scenarios:
 - Over
 - Adjacent
 - Under
- See Talk by Tim Wickberg at SC 2022
 - <https://slurm.schedmd.com/SC22/Slurm-and-or-vs-Kubernetes.pdf>
- Overview: <https://www.schedmd.com/slinky/why-slinky/>
- Code: <https://github.com/slinkyproject>

System Overview - Interlink Edge-Node



System Overview - Interlink Slurm SSH access



Summary

Looks like a gap between HPC – Cloud?

Summary - Differences

/ 20

1. The end-user workflows / processes to “do science” look different
 1. We move from “use a terminal” -> “use a webportal”
 2. No need for workstation data copy anymore
 3. “Change management” and “expectation management” for end-users is required
 1. Do they need to learn “kubectl” to interact with Kube or a “webfrontend” to interact with the Science Platform?
2. The underlying technical foundation changes
 1. We move from “use Slurm Job processing” -> “use Kubernetes service layer”
 1. (But do we really? Or will both enter in a marriage?)

Summary – Tech Challenges

/ 21

1. Is it better to have one single beefy Kube cluster or schedule workload outside?
2. Can we reliably and consistently schedule workload via Kube on an existing Slurm cluster?
3. What is the best solution to make the same data accessible inside Kube for a service/app and inside Slurm for heavy computation (without copying the data)?
 1. Blockstorage vs ObjectStorage vs NFS vs etc.
 2. Can we avoid data copies between “user areas” vs “raw data areas”
4. What will be the “winning” orchestration strategy of the astro algorithm/workload?
 1. Push via packaged container as “standalone” workload in Kube?
 2. Push via “direct code” through Software platform (e.g., Jupyterhub, Canfar, etc.)?
5. And plenty plenty more ...