



JupyterHub on Kubernetes as a platform for developing secure shared environment for data analysis at MAX IV

Andrii Salnikov,
Zdenek Matej, Dmitrii Ermakov, Jason Brudvik

Interactive data analysis environment

- **Container images** with pre-defined and custom kernels
- **Kubernetes cluster**
 - as a **resource pool**:
 - Moderate CPU
 - Large RAM
 - V100/A100 GPUs
 - as a **deployment platform**:
 - review/prod/next lifecycle
 - CI testing of notebook images
 - as a **runtime environment**
- **Shared service** for staff and researchers
 - Remote-desktop style experience
 - Resources overcommit

The screenshot shows a JupyterLab interface with a terminal window open. The terminal displays the output of the `id` and `nvidia-smi` commands. The `id` command shows the user's group memberships, including MAX-Lab. The `nvidia-smi` command displays the NVIDIA-SMI version, driver version, CUDA version, and a table of GPU details.

GPU	Name	Fan	Temp	Perf	Persistence-M Pwr:Usage/Cap	Bus-Id	Disp.A Memory-Usage	Volatile GPU-Util	Uncorr. Compute M. MIG M.	ECC
0	NVIDIA A100	N/A	32C	P0	Off 124W / 300W	00000000:CA:00.0	Off N/A	N/A	On Default Enabled	

GPU ID	GI	CI	MIG Dev	Memory-Usage BAR1-Usage	SM	Vol Unc ECC	Shared CE ENC DEC OFA JPG
0	2	0	0	1256M1B / 40192M1B 2M1B / 65535M1B	42	0	3 0 2 0 0

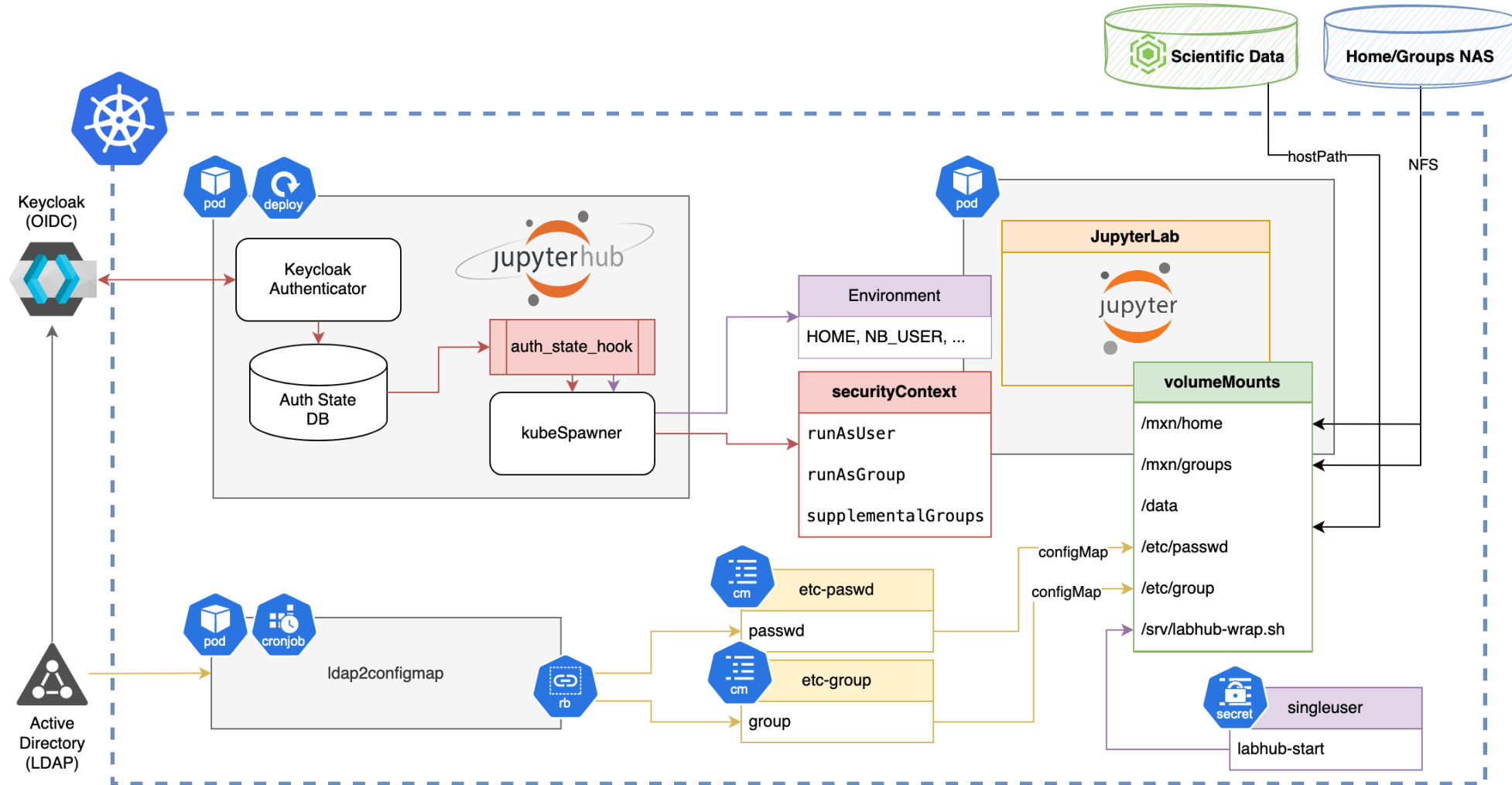


Goals and technical requirements

- **Key objective:** fully unprivileged container environment that operates seamlessly with existing **LDAP user credentials**
- **Functional Requirements:**
 - Integration with **MAX IV storage systems** (home, group, data)
 - Run **any notebook images** without modifications
 - Ensure available **resources visibility**
 - Efficient **sharing of available GPU** resources between users
 - Observability of **usage metrics**
- **Operation Requirements:**
 - [Zero to JupyterHub with Kubernetes](#) Helm Chart **without** modifications
 - Just custom hooks and proper `values.yaml`

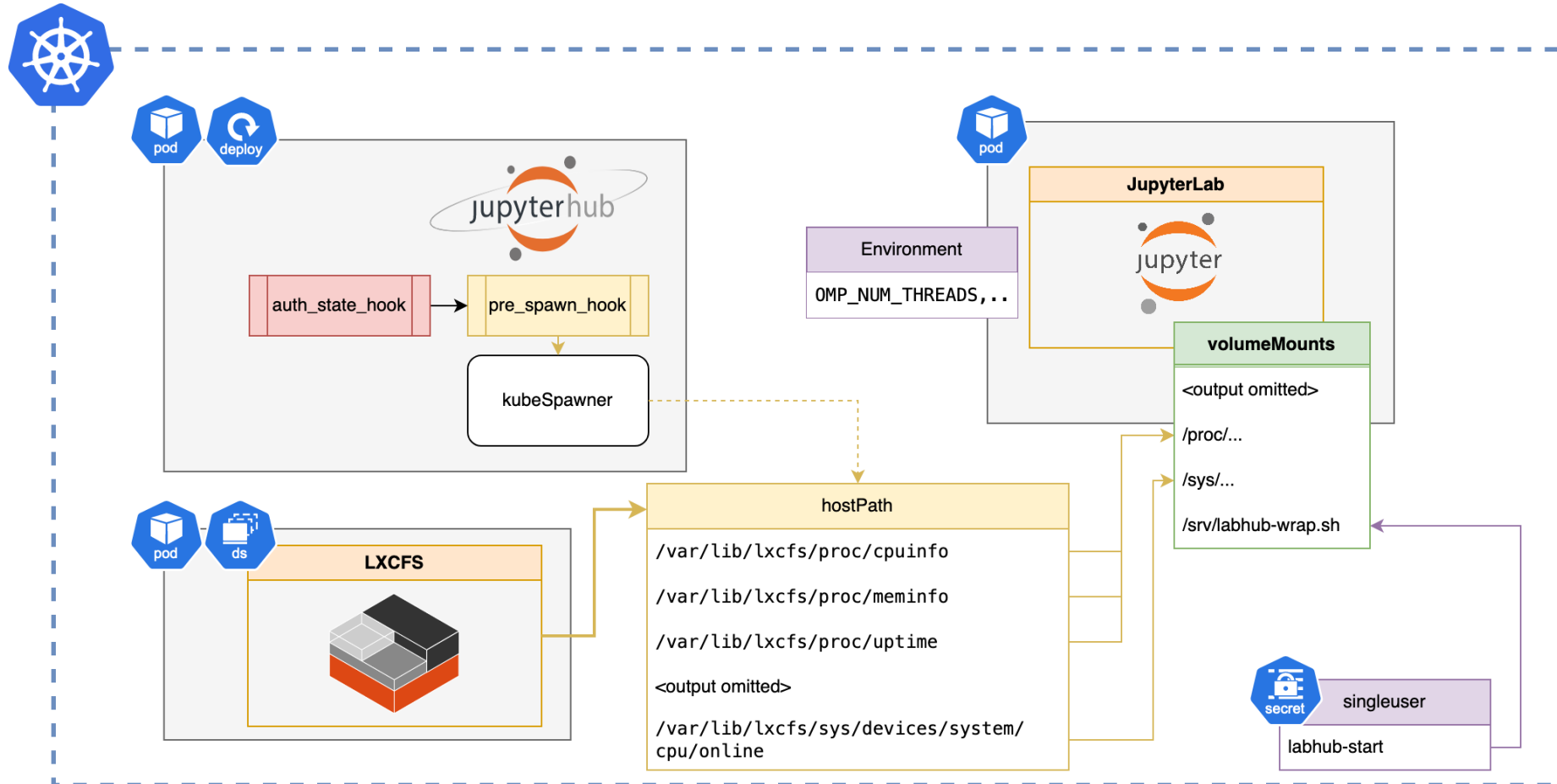
Existing LDAP credentials

- **UID/GIDs** from Token to securityContext
- **NSS data sync from LDAP** to configMap to mount inside container
- **Environment variables** to define HOME directory, etc
- Wrapper **startup script** to bootstrap environment
- **Storage mounts** are simply defined in the values.



LXCFS: Resources visibility

- [LXCFS](#) is a FUSE filesystem offering **overlay files** for `cpuinfo`, `meminfo`, `uptime`, etc
- **Deployed as DaemonSet** on Kubernetes level
- Visible CPU and RAM **container limits**
- Mounted to `/proc` and `/sys` in **pre_spawn hook**
- Defining additional environment variables in startup scripts



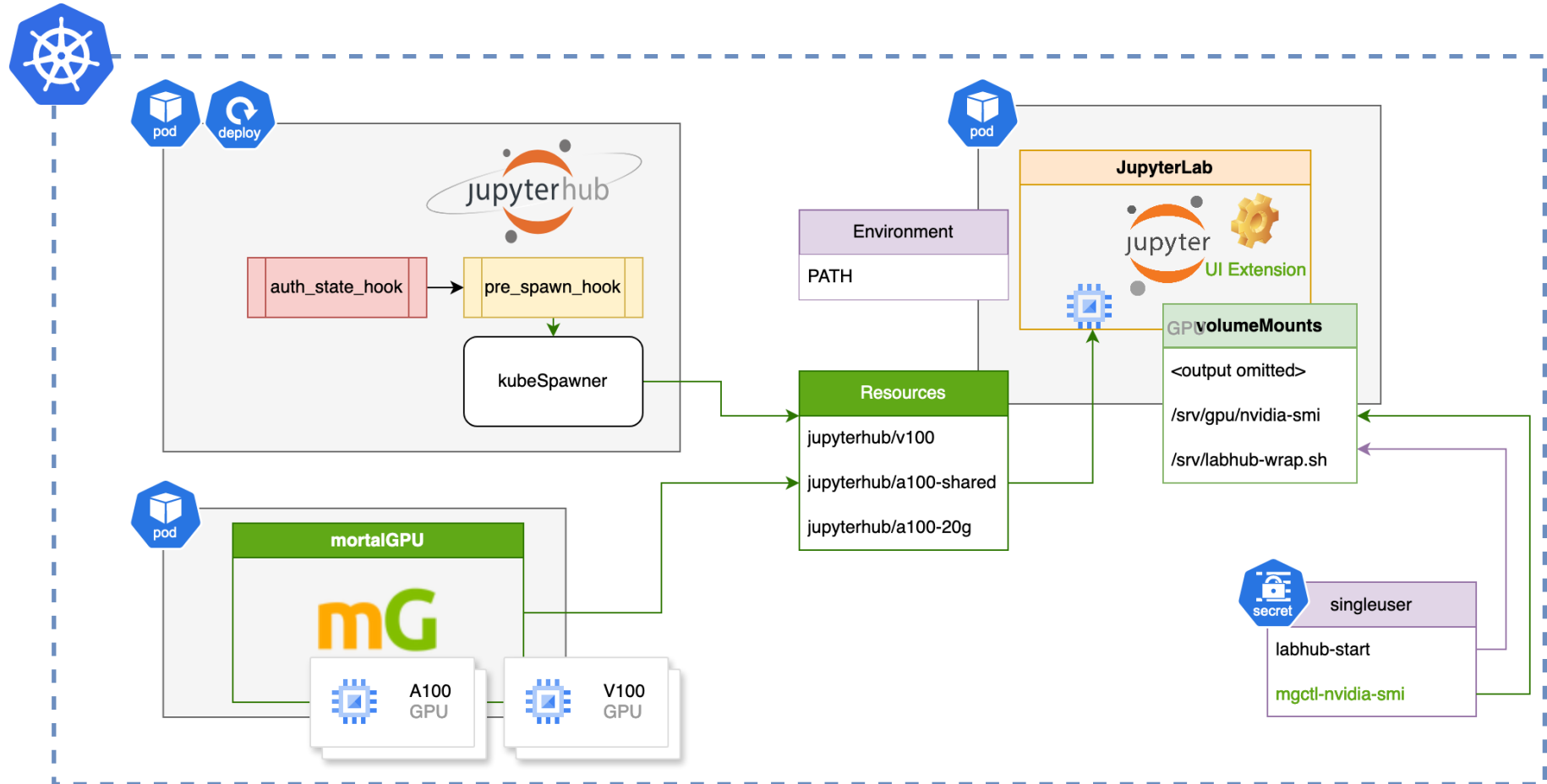
GPU sharing: MortalGPU development

- Kubernetes **device plugin** for **GPU memory overcommit**, while maintaining **allocation limit per GPU workload** - the approach used for sharing RAM on Kubernetes.
- Fork of MetaGPU with development focus on **interactive workloads** run by mortals (with operations support by mortal admins)
- Provides:
 - **Device Plugin**: represent GPU (or MIG partition) with configurable number of meta-devices (e.g. 320 of mortalgpu/v100)
 - **Memory enforcement** based on the usage monitoring data
 - **Kubernetes-aware observability** in general and container-scoped resource usage in particular:
 - mgctl tool and Prometheus exporter



Jupyterhub with MortalGPU

- Kubernetes DaemonSet
- GPU RAM resource requests and limits, defined the same way as RAM
- Multiple MortalGPU resources available (different GPUs and partitions)
- Wrapper over mgctl to provide `nvidia-smi` output for container processes only

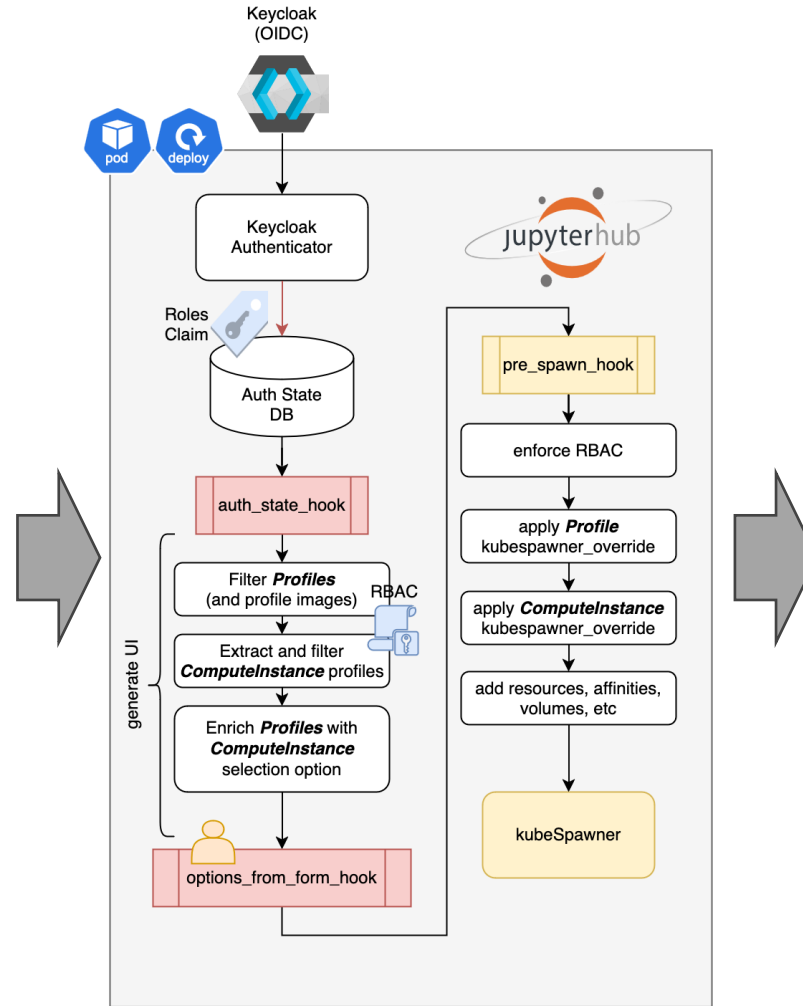


Compute Instance profiles and RBAC

```

1  profileList:
2  # Compute instance resource profiles
3  - compute_instances:
4    a100_shared:
5      display_name: "Shared A100 (Limits: 100GB RAM, 8 CPU, 8GB
6      cpu_guarantee: 0.1
7      mem_guarantee: 500M
8      cpu_limit: 8
9      mem_limit: 100G
10     gpu_resource_name: jupyterhub/a100-shared
11     gpu_guarantee: 4
12     gpu_limit: 80
13     node_type: a100
14     a100_20g:
15       display_name: "Dedicated A100 20G Partition (Limits: 100GB
16       required_role: a100-20g
17     # ...
18   - display_name: "Beamline-specific analysis"
19     description: "Notebooks tailored for specific beamline needs"
20     slug: "beamline"
21     allowed_compute_instances:
22       - a100_shared
23       - a100_10g_quick
24       - a100_10g
25       - a100_20g
26       - v100
27     profile_options:
28       image:
29         display_name: Image
30         choices:
31           bloch:
32             display_name: "Bloch (h5py, igor, ipympl, ipywidgets,
33             kubernetes_override:
34             image: harbor.maxiv.lu.se/jupyterhub/bloch-notebook

```



MAX IV JupyterHub

Common analysis

Generic commonly used analysis tools collection

Image: HDF5 simple analysis (bohrium, h5py, hdf5plugin, matplotlib, numpy, p...

Compute Profile:

- Shared A100 (Limits: 100GB RAM, 8 CPU, 8GB VRAM)
- Dedicated A100 10G Partition (Limits: 50GB RAM, 8 CPU, 10GB VRAM, 30 min)
- Dedicated A100 10G Partition (Limits: 100GB RAM, 8 CPU, 10GB VRAM, 12 ho
- Dedicated A100 20G Partition (Limits: 100GB RAM, 8 CPU, 20GB VRAM, 12 ho**
- Shared V100 (Limits: 50GB RAM, 8 CPU, 8GB VRAM)

Beamline

Notebooks tailored for specific beamline needs

Image: Bloch (h5py, igor, ipympl, ipywidgets, lmfit, matplotlib, numpy, scipy, s...

Compute Profile: Shared A100 (Limits: 100GB RAM, 8 CPU, 8GB VRAM)

Nordugrid ARCV7 Client

Development and testing of Nordugrid ARC

Compute Profile: CPU-small (Limits: 10G RAM, 2 CPU)

Development

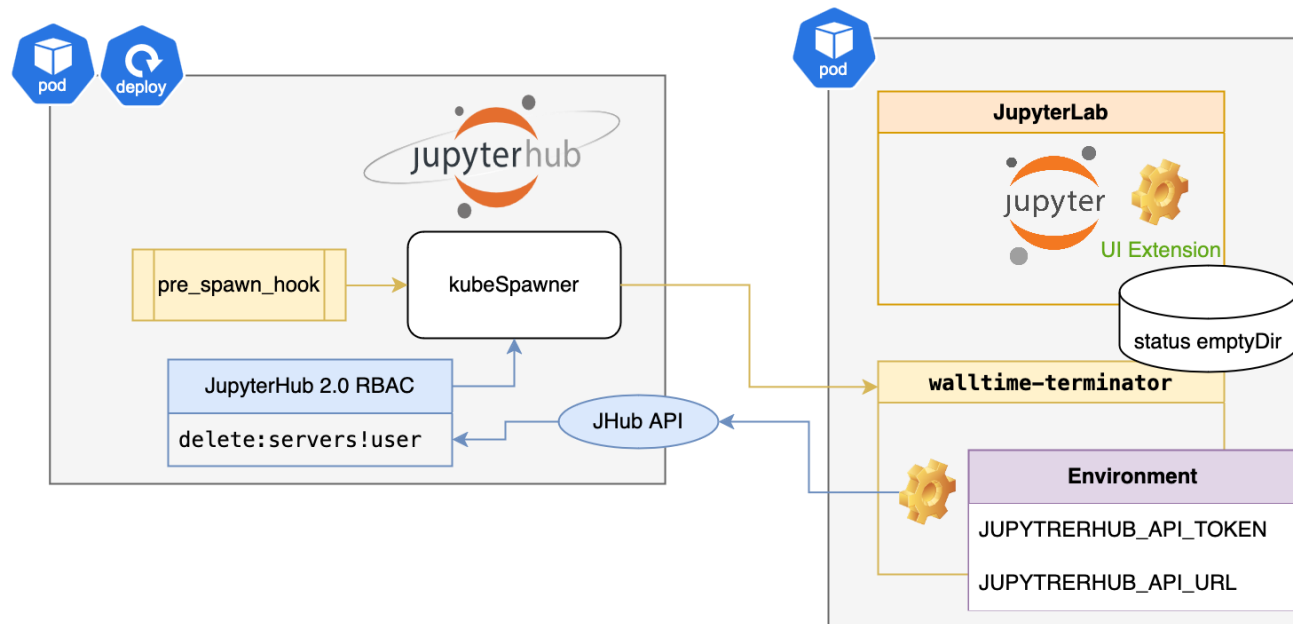
Development and testing of notebooks (not for general use)

Compute Profile: Dedicated A100 10G Partition (Limits: 100GB RAM, 8 CPU, 10GB VRA...

Extra containers = extra features

Walltime enforcement

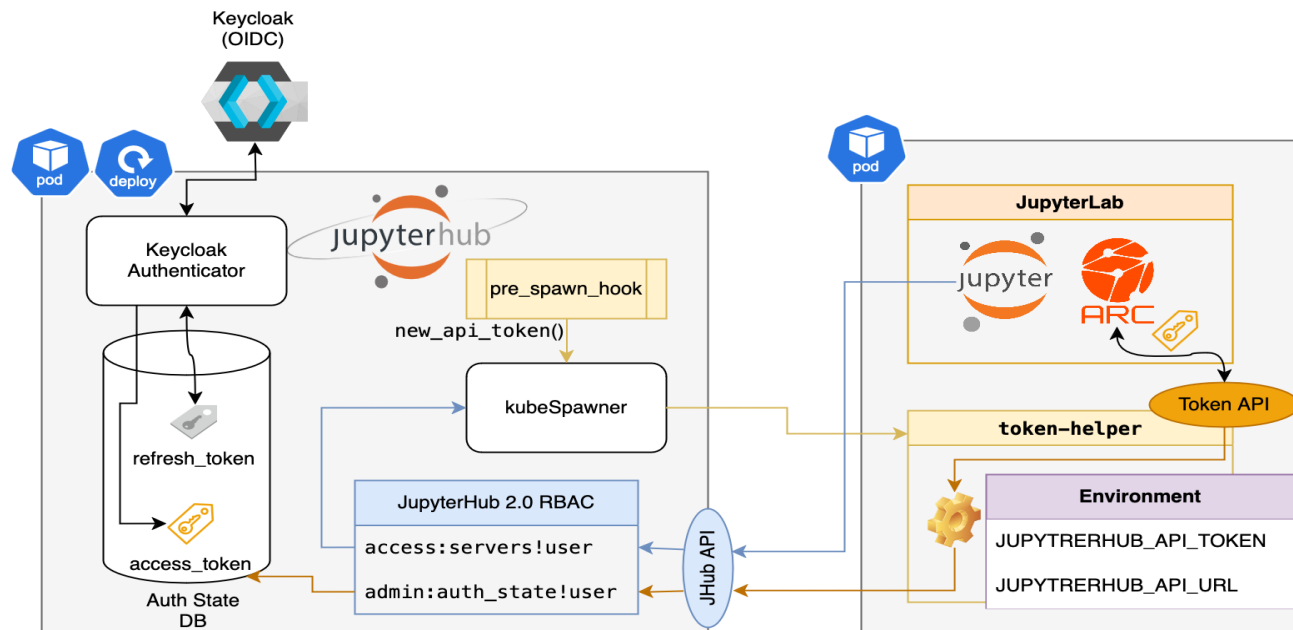
Mem: 450 / 51200 MB | GPU: 1164 / 10201 MB | Walltime: 28 / 30 min



- KubeSpawner is capable of **running additional containers** in the user Pod
- Isolated walltime countdown container **terminating user server** via JupyterHub API
 - Using the JupyterHub [RBAC](#) feature
- **Developed UI extension** to show values to end-user

Extra containers = extra features

”OIDC-agent” for ARC



- [KeyCloak Authenticator](#) to refresh access tokens
- Isolated `token-helper` container with privileges to read `auth_state`
 - Using the Jupyterhub [RBAC](#) feature
- API to provide **only Access Tokens** to JupyterLab container
 - wrapper to use in ARC CLI transparently

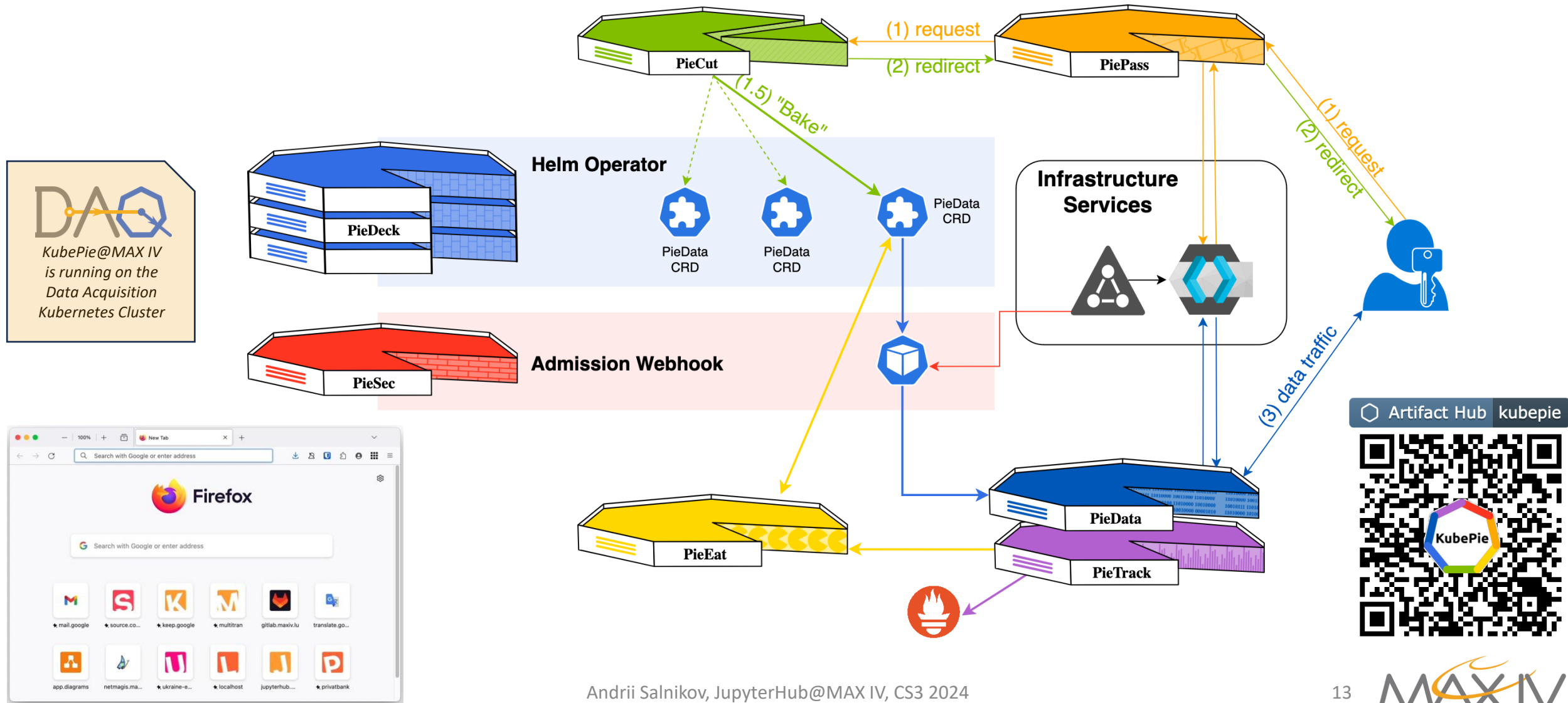
KubePie: sharing existing data over https

- **Idea:** "own" web server for each user with the correct UID/GIDs
 - Sounds crazy? But we do run such Pods for each user in JupyterHub!
- **KubePie** is harnessing Kubernetes' scalability and deployment capabilities by running, managing and securing web servers for every user
- **KubePie** is strictly relying on *OpenID Connect* flow or *OAuth2 bearer tokens* when it comes to the user identification
 - OAuth2 used in ARC PoC for data transfers
 - Claims-based user-mapping during Pod instantiation (admission)
 - Other auth credentials accessible via OIDC:
 - WebDAV with S3-like credentials is implemented as an example

CHALLENGE ACCEPTED



KubePie: Baking process



Conclusions

- **Extensibility of both JupyterHub and Kubernetes** *allows* to build data analysis platforms, matching organization needs in functionality and security.
- **LXCFS on the Kubernetes** *brings* allocated resources visibility to both interactive and batch containerized workloads.
- **Flexible and observable GPUs sharing with MortalGPU** *enriches* the interactive shared environments with CUDA capabilities.
- **Compute Instance profiles and RBAC** *extends* the usage patterns of the shared platform, improving the end-user experience.
- **Additional containers** in the running Pod *open a way* to securely add features beyond the usual JupyterHub capabilities.

Thank you for attention!



Source code and deployment configuration can be found on gitlab.com



We are working towards establishing similar deployment for providing EOSC service as Open Data analysis platform

Mail To: andrii.salnikov@maxiv.lu.se

Andrii Salnikov, JupyterHub@MAX IV, CS3 2024

15

