WIR SCHAFFEN WISSEN – HEUTE FÜR MORGEN

**Leonardo Sala :: AWI :: Paul Scherrer Institut**

# DARI, SLS and RA

**AWI Department meeting 2022.12.12 / PSI**

# Outline

- Who we are / What we do
- RA news and highlights
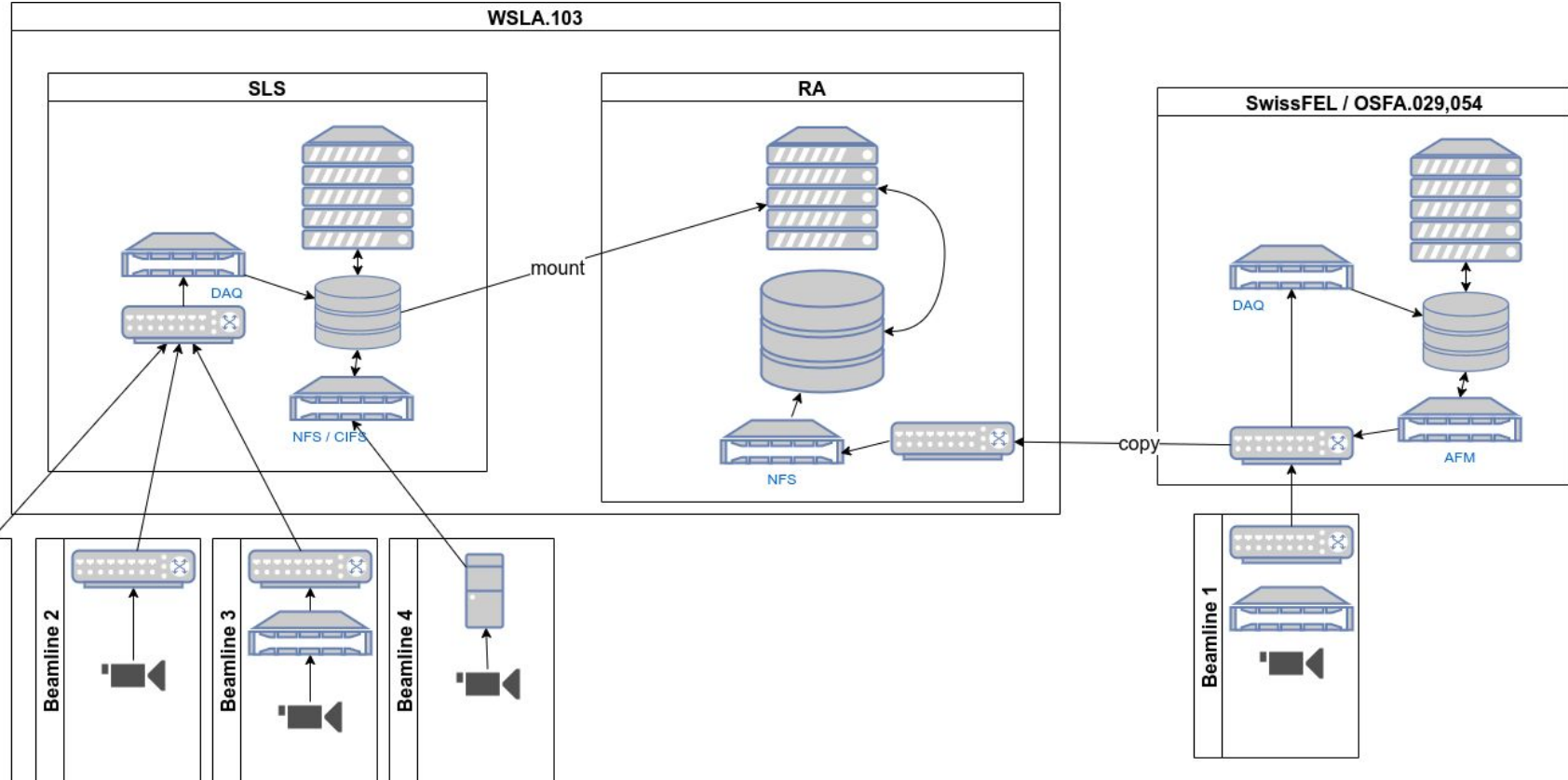- SLS news and highlights
- Next major projects

# Who / What

**Who** we are:
- **Ivano Talamo:** RA responsible and main admin
- **Alvise Dorigo:** SwissFEL responsible and main admin
- **Joshua Taylor:** SLS responsible and main admin
- **Krisztian Pozsa:** Expands / EOSCFuture projects (project)
- **Borys Sharapov:** SLS 2 DevOps / Sysadmin (project)
- **Leonardo Sala:** Group Lead

**What** we do
- Manage **IT infrastructure** for Photon experiments (storage, compute, network, …)
- Manage special **services** (archiving, opendcim, jupyterhub, …)
- Enable **DevOps** and best practices to enable scientists / staff to do their job

# Big picture

# Some numbers

- **~9000 cores, ~20 PiB, ~70 TB ram, ~30 managed switches (Infiniband, Ethernet)**
- managed by **Puppet and ansible**
  - infrastructure as code backed up by Gitlab
  - Puppet for basic / standard OS installation
  - Ansible for special setups, pipelines and operations
- monitored by **Icinga and InfluxDB / Grafana**
  - for more details see Alvise's talk
  - looking into AIT ELK
- We even have a test Openshift k8s cluster
  - used for gitlab runners and tests

# Highlight: server installation

Our server installation is mostly automated:
- physical server installation
- register default admin credentials and required variables
- run playbook that:
    - configure BIOS based on profiles
    - register system in linux inventory
    - configure RAID, boot device, …
    - boot up server
- based in industry-standard Redfish API

Next steps:
- automatic filling of our Data Center management system (opendcim)
- this is possible now as we recently installed a version with a RESTful API

# RA updates

**Slurm resources management:**
- migrated away from full node allocation towards resources allocation
- this allows us to:
  - efficiently manage CPU and GPU resources
  - limit resources usage with cgroups -> less interference
  - have similar setup to Merlin -> can bother Marc even more :D

**New Jupyterhub interface**
- more dynamic (python + javascript) - depending on queue, shows different options
- more improvements to come, like run time checks

**New storage**
- project-funded storage is more than 5 years old now
- replacing 2 x 2.2 PiB systems with 1 x 6 PiB system
- delivery this week (Xmas present)

# Highlight: tape retrieve

Thanks to a joint effort by AWI (Ivano Talamo, Krisztian Pozsa, Stephan Egli, Carlo Minotti) and AIT (Peter Huesser, Michael Kallmeier), **simple one-click tape retrieve from CSCS Petabyte Archive to RA storage is available now**. Fixing now some bugs



A similar mechanism will allow retrieval from tape to CSCS Object Storage (possibly early 2023)

# SLS updates

**Automatic quota warning system**
- overcomes some icinga limitations
- automatic warning emails sent to beamline scientists
- scientists can self manage thresholds and address list

**ACLs**
- most data writers run as root -> not good
- explore ACLs usage to write data without root privileges
- prototype with MX successful, plan to propagate to SLS and SwissFEL together with the developers
- A tool to verify ACL policies is being implemented

**DAQ support**
- support MX Jungfraujoch efforts
- migration away from IBM Power architecture (Filip Leonarski)

# SLS updates / II

**Migration from Samba wide-links**

- current data access over samba mounts e-account home directory, access data directory over symlink
- this is not supported anymore due to security concerns
  - it also creates quota issues when mounted over Windows
- new mountpoints will be created during the long Shutdown

**Plan:**

- Enable new mountpoints in our prod_2 cluster
- Migrate beamlines from prod_1 to prod_2
- Keep prod_1 running as-is in case of issues
- Reinstall prod_1 and upgrade Spectrum Scale versions

# Highlight: services deployment

Quite some work has been put in the past to make services deployment to DAQ nodes reproducible and automatic -> this is the way SLS DAQ nodes are mostly managed since some time

New effort to improve the system and support beamlines-managed services, e.g. MX analysis pipeline

**Requirements:**
- separate code from data (config files)
- restrict control over code and pipeline definition
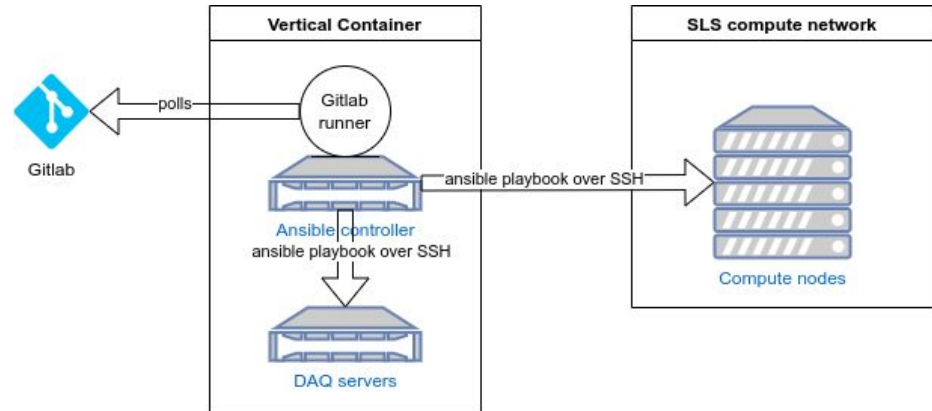- have a simple and intuitive interface

# Gitlab pipelines - architecture

**Solution:**
- IaaC with Ansible playbooks
- Gitlab as interface, deployment jobs through pipelines
- access control based on repositories

**Advantages:**
- fine grained control
- reproducible and versionable
- web interface to output

# Pipelines

SLS Online IT > MX > configs > **Repository**

master ⌄    configs / **spotter.yaml**

**Updated yaml files for all workers**
wojdyla_j authored 8 months ago

📄 **spotter.yaml** 📋 380 Bytes

```
 1  cn_x06da:
 2    - name:  spotter
 3      beamline: x06da
 4      version: stable
 5      workers_n: 12
 6  cn_x06sa:
 7    - name:  spotter
 8      beamline: x06sa
 9      version: stable
10      workers_n: 18
11  cn_x10sa:
12    - name:  spotter
13      beamline: x10sa
14      version: stable
15      workers_n: 22
16  cn_vagrant:
17    - name: spotter
18      beamline: x06sa
19      version: stable
20      workers_n: 10
```

SLS Online IT > MX > configs > **Pipelines**

**All** 71    Finished    Branches    Tags

Filter pipelines

| Status | Pipeline | Triggerer | Commit | Stages |
|--------|----------|-----------|--------|--------|
| ⊘ passed | #16496 | | ⑂ **winter-shut...** ⦾ bf231570  Remove x06da-cn* and mx-cn*. Add r... | ✓ ✓ |
| ⊘ passed | #14418 latest | | ⑂ **master** ⦾ 2d2ef2d8  add jobworker on ra-c-017 | ✓ ✓ |
| ⊘ passed | #14417 | | ⑂ **master** ⦾ 4a7903b7  Merge branch 'master' of git.psi.ch:sl... | ✓ ✓ |

```
3184  x10sa-cn-122.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3185  x10sa-cn-123.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3186  x10sa-cn-124.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3187  x10sa-cn-125.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3188  x10sa-cn-126.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3189  x10sa-cn-127.psi.ch        : ok=22    changed=1    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3190  x10sa-cn-128.psi.ch        : ok=22    changed=1    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3191  x10sa-cn-129.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3192  x10sa-cn-130.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3193  x10sa-cn-131.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3194  x10sa-cn-132.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3195  x10sa-cn-133.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3196  x10sa-cn-134.psi.ch        : ok=20    changed=0    unreachable=0    failed=0    skipped=22    rescued=0    ignored=0
3197  xbl-daq-37.psi.ch          : ok=7     changed=0    unreachable=0    failed=0    skipped=5     rescued=0    ignored=0
3198  Thursday 01 December 2022  13:47:39 +0100 (0:00:02.449)       0:05:09.419 *****
3199  ===============================================================================
3200  psi.adp --------------------------------------------------------------- 143.52s
3201  psi.spotter ------------------------------------------------------------ 46.49s
3202  stat ------------------------------------------------------------------- 31.57s
3203  include_role ----------------------------------------------------------- 27.14s
3204  psi.adm ---------------------------------------------------------------- 17.55s
3205  psi.dimmer ------------------------------------------------------------- 9.59s
3206  psi.jfjoch_writer ------------------------------------------------------ 5.76s
3207  include_vars ----------------------------------------------------------- 5.20s
3208  file ------------------------------------------------------------------- 4.25s
3209  systemd ---------------------------------------------------------------- 4.18s
3210  ansible.builtin.service_facts ------------------------------------------ 3.99s
3211  set_fact --------------------------------------------------------------- 3.94s
3212  gather_facts ----------------------------------------------------------- 3.86s
3213  include_tasks ---------------------------------------------------------- 2.32s
3214  ~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~
3215  total ----------------------------------------------------------------- 309.38s
3216  Playbook run took 0 days, 0 hours, 5 minutes, 9 seconds
3217  Cleaning up file based variables                                          00:00
3219  Job succeeded
```

14

# Next big projects

- **Compute node merge**
  - manage SLS and RA compute nodes as a unique pool of resources
  - prototype for feasible switching between different configurations
- **Storage WTO** - we need a new one to replace existing systems and procure SLS2 resources
- **Resources API**
  - Allowing users/beamline to self manage resources with no or minimal intervention from admins e.g. compute node reservation, quota extensions, …
    - possible by exposing resources operations and workflows via APIs
  - early discussions about DUO integration
- **Data lifecycle**: write - read - archive - delete - retrieve
  - streamline policies and workflows
  - including paid storage for projects / grants
- **SLS2**
- **Documentation and dashboards**

# Highlight: WHGA server room

Migration away from SLS server room during dark period.

Started using WHGA server room for RA:
- 4 compute nodes
- Infiniband switch (2x100G, additional links to be added in the next weeks)
- 6 PB Storage in January (Lenovo DSS-G260)

Plan to gradually migrate compute nodes, phase out storage based on lifecycle. No downtimes.

# Questions?