# CALIPSOplus JRA2: Prototype of Data Analysis as a Service Platform

Aidan Campbell (ESRF)

Daniel Webster (PSI)

# Motivation

- Light sources are generators of big volumes of complex scientific data and their users need assistance in analysing the scientific data.

- Our aim is to provide a remote Data Analysis As a Service (DAAS) portal and platform for users to:

  - Access their experimental data

  - Use pre-packaged Data Analysis software available at each institute

  - Access onsite computer resources to assist with data reduction and processing

# CALIPSOPlus DAAS Portal

- Joint Project: ALBA/ESRF/PSI

- Written: Angular 7 / Django

- Common home page which connects user to institute portal

- Custom installation tested at: ESRF, ALBA, ELETTRA, PSI, DESY, SOLEIL, DLS

- Services:
  - Jupyter Notebooks
  - Containers
  - Virtual Machines

# CALIPSOPlus DAAS Portal Architecture*



See paper https://icalepcs2019.vrws.de/papers/wepha057.pdf

# Jupyter Notebooks (ESRF)

- ## Currently
  - 40+ users
  - Limited hardware resources

- ## Near Future
  - 100+ users
  - SLURM
    - Scalability
    - GPU Access
  - More GPUs available

# Jupyter Notebooks (DESY)

- **2018 - 2019**
  - 180+ users
  - Typically 60-120
  - concurrent sessions
  - SLURM
    - Scalability
    - GPU Access
  - 3 AMD dedicated nodes
    - Other partitions can
    - be used too

## Maxwell Jupyter Job Options

**Maxwell partitions** shared node on Jupyter partition

**Choice of GPU** none

**Note:** For partitions without GPUs (or choice of GPUs) the GPU selection will be set to 'none'

**Job duration** 1 hour

**Note:** on the shared Jupyter partition (jhub) the time limit is always 7 days!

| Current Status | | | | | |
|---|---|---|---|---|---|
| Partition | # nodes | # avail | # GPUs avail | # P100 avail | # V100 avail |
| jhub | 3 | 3 | 0 | 0 | 0 |
| maxwell | 61 | 46 | 0 | 0 | 0 |
| maxgpu | 19 | 10 | 10 | 5 | 5 |
| all | 327 | 186 | 0 | 0 | 0 |
| allgpu | 88 | 48 | 48 | 38 | 5 |

Spawn

# Compute Options: VM, Container, Bare Metal



**Virtual Server**

| APP1 | APP2 | APP3 |
|------|------|------|
| bin/lib | bin/lib | bin/lib |
| OS1 | OS2 | OS3 |
| Hypervisor | | |
| Host Operating System | | |
| Hardware | | |

**Container**

| APP | APP | APP |
|-----|-----|-----|
| bin/lib | bin/lib | bin/lib |
| Container Engine | | |
| Host Operating System | | |
| Hardware | | |

| APP1 | APP2 | APP3 |
|------|------|------|
| bin/lib | bin/lib | bin/lib |
| Host Operating System | | |
| Hardware | | |

**Bare Metal Server**

© 2016 Kumulus Technologies

Pre-installed software

Pre-installed software

Sudo privileges

CALIPSO plus

# Containers

- Computer that scientists and users can access with pre-installed software entirely in the web

- Scientists can:

  - Request a specific Linux OS (Ubuntu, CentOS, Debian etc)

  - Install their own software

  - Access their data (NFS)

  - Do analysis on site hardware from home/university

- Developers can:

  - Create/update container images with Github

  - Containers will update automatically on portal

# Virtual Machines

- Computer that scientists and users can access with pre-installed software entirely in the web

- Scientists can:

  - Request a specific OS (including Windows)

  - Can't install their own software (security risks)

  - Access their data

  - Do analysis on site hardware from home/university

- Use Cases tested:

  - PyMca (ESRF), pyFAI (ESRF), PtychoShelves (PSI), Savu (DLS), CrysFEL (DESY)

# Future Developments

- Integrate experiment data within the portal

    - Custom plug-ins required

- Building tailored Notebooks for beamlines

- Create containers/virtual machines without needing an experiment

- Write report on needs of Calipsoplus community wrt European Open Science Cloud

- Integrate results with PaNOSC and ExPaNDS DAAS portal

# PSI Deployment

- This project represents an ideal opportunity to prove emergent technologies

  - Containerisation and orchestration thereof

  - Microservice Architectures

  - "DevOps" methodologies

- PSI chose to deploy the portal on Red Hat OpenShift

  - Enterprise-grade Kubernetes distribution

- CALIPSOplus portal had to be deployed as a Microservice Architecture

  - A wealth of experience gained for PSI in this methodology

# PSI Use Case

- ## We engaged our cSAXS Beamline for our CALIPSOplus use case

  - **cSAXS** – Coherent Small-Angle X-ray Scattering - uses Ptychography (among other techniques) for image reconstruction

  - **Ptychography**: computationally generate images by processing coherent interference patterns

  - The application we are using for this is **PtychoShelves**

    - **PtychoShelves**, a versatile high-level framework for high-performance analysis of ptychographic data

      

    - Paper can be found here: **https://scripts.iucr.org/cgi-bin/paper?zy5001**

      - *(Will paste into chat window)*

# PSI Architecture
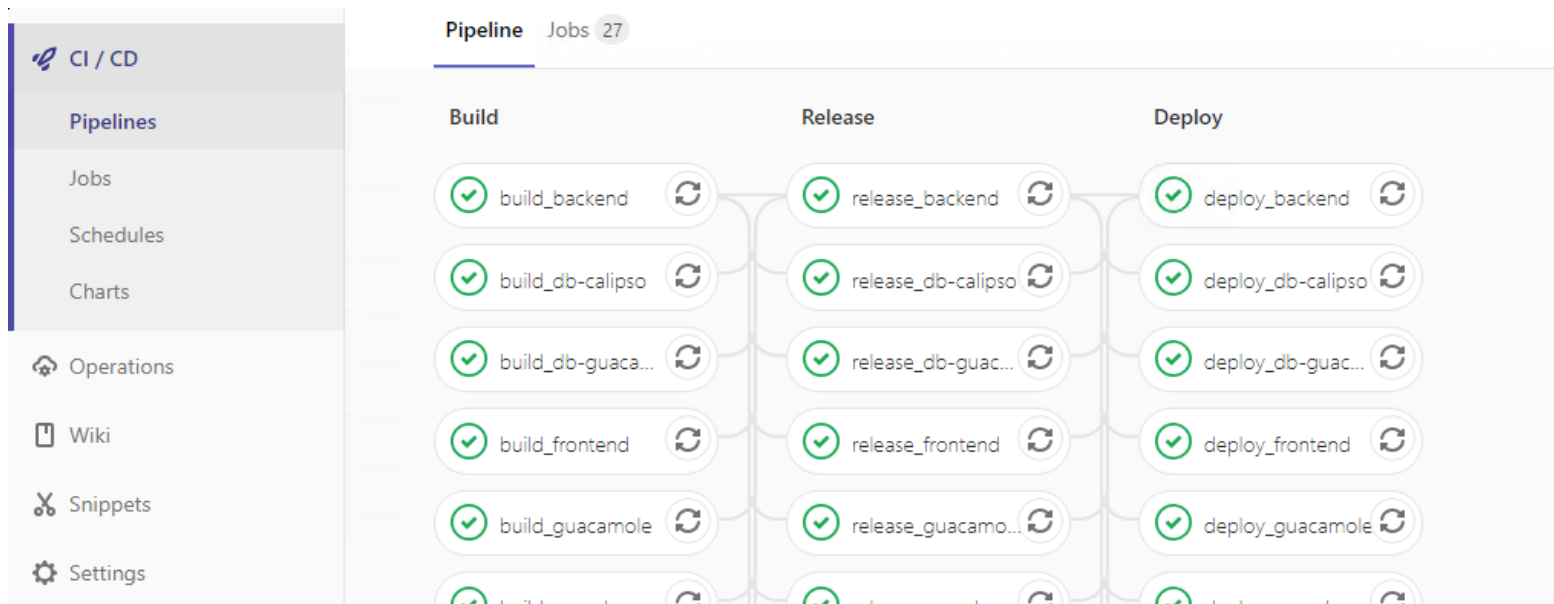
# Deploying the Application

- We deploy CALIPSOplus as a Microservice Architecture via "Continuous Integration/Continuous Deployment" (CI/CD) from GitLab:



- A new build is triggered, released, and deployed automatically to our OpenShift instance, upon committing new code

- This is "**DevOps**" in action

# Demonstration - Screenshots

- We select the proposal we are interested in, launch our chosen container on top of this data, and connect to it via RDP or VNC



- A Linux "system" will now be at our disposal, and our scientific application is already mapped-in and ready to run

# Demonstration – Screenshots

- We can then run our MATLAB package, and get our results:

# Live Demo – please standby..

- We will now show a live demo of the portal in action

# Conclusions

- CALIPSOplus JRA has been very useful in bringing together sites to collaborate and share a prototype portal for providing Data Analysis as a Service

- Feedback from users was positive and demonstrated the need for such services

- Jupyter service is being generalised at most sites

- Other services (containers+VM) are under test

- Main difficulty encountered in providing DAAS services in production is the lack of data analysis policy at all the sites (a survey has been prepared to get feedback from sites)

- Future developments will be in EOSC with PaNOSC+ExPaNDS

# Acknowledgments

- **ESRF:**

- A.Campbell, A.Götz, A.Rau, , J.Kieffer, T.Vincent, A.Sole, A.de Maria, M.Retegan, V.Rene-Nicolas, B.Roussel

- **ALBA:**

- D.Salvat, A.Camps, D.Sanchez

- **Elettra:**

- G.Kourousias, I.Ardian, D.Palmisano

- **Paul Scherrer Institut:**

- M.van Daalen, SEgli, D.Webster, AAshton

- **DESY:**

- J. Reppin, F.Schluenzin

- **Diamond LightSource**

- T.Schoonjans

- **SOLEIL:**

- M.Oursy, G.Viguier

- **Helmholtz Zentrum Dresden Rossendorf:**

- B.Schramm, M.Grobosch

# Thank You

*...questions?*

# Backup slides

- Extra slides which can be used in case there are questions on certain details

# Jupyter Notebook Prototypes

- SLURM (DESY, ESRF)

  - Interactive notebook server created using a scheduler

  - GPU support

  - 12 hour session

  - Cannot install new software (except by pip install --user)

- SudoSpawner

  - Single machine with GPU support

  - Very limited hardware resources

- Kubernetes (ALBA, PSI)

  - Load balanced notebooks (more computing resources)

  - Unlimited session time

  - Can install new software (apt-get, pip, etc) temporarily

  - Notebook is culled after X hours of inactivity

  - Multiple custom notebook images

CALIPSO<sup>plus</sup>

# Virtual Machines

- Virtual machines can be created using multiple systems used by each site.

- ESRF:
  - KVM

- Others:
  - OpenStack
  - KVM
  - Citrix
  - Vsphere / VMWare

- Creating a container can take time
  - Maintain "bank" of virtual machines available at all times

# Containers

- Previously

  - Limited to Docker

  - Limited to a single machine (not scalable)

  - Apache Guacamole

  - NFS for data access

- Now

  - Kubernetes

    - Scalability

    - Orchestration

  - Apache Guacamole but in the same browser

  - NFS for data access

# Administration

- View resource usage for Containers and Virtual Machines

- Upload new container and vm images

- Manage all containers and virtual machines

- Manage all users