Outline

- Objectives and user stories
- An interTwin k8s virtual node
- State of the art
- First pilot experiences
- Current status and plans



interTwin

Enabling transparent access to heterogeneous computing providers via interLink

Diego Ciangottini -







- Federate a set of highly heterogeneous and disparate providers
 - enabling a "transparent payload offloading"
- Flexible extension of existing computing clusters to remote resources
 - E.g. HPC, specialized hardware...





User stories: Scheduling K8s POD on HPC

This is the most generic one:

- I have a container and I want it executed on a node with N GPUs and M GB of RAM
- I want to create a simple container to be scheduled on a remote Slurm batch on an HPC center





For development of ML pipelines, Jupyter notebooks are a common solutions

- I want access to a hub where I can ask for the spawning of a JupyterLab instance on a node with GPUs
- In other words I want to schedule a JupyterLab container on a remote slurm cluster
 - Similar to the previous one, with the exception of bidirectional connection needed b/w Lab and HUB that is a bit tricky



User story: pipelines/workflows

- Frameworks for DAG/workflow managements are usually well integrated with K8s APIs
 - Airflow/Kubeflow pipelines/Argo wokflows/MLFlow are no exceptions
- Pilot exploitation possible thanks to interTwin AI Workflow management activities
 - CERN group in particular has technically produced a demo container image for this





- Executing payload in response of an external trigger
 - Being it a storage event or a web server call
- Pilot exploitation with <u>dCNiOS</u> (dCache + Nifi + OSCAR) devs for a first round of validation with payload executed on a virtual node



Building an "offloading" prototype

- Extend the container orchestration de-facto standard (K8s) to support offloading under the hood
 - With little to no knowledge required from the end user perspective
- What about making a Kubernetes cluster send your containers/pod to a "virtual" node that seamlessly takes care of your application lifecycle on a remote server/hpc batch queue?
 - While exposing the very SAME experience of running a pod on the cloud resources
 - Meaning that the API layer exposed is the "normal" set of K8s apis

N.B. We aim to use Kubernetes as the workhorse for the "offloading", NOT as the user interface though



Kelsey Hightower 🥏 @kelseyhightower

The problem is we asked developers to do all that. Kubernetes is not a tool for developers. They can use it, but we have to be honest, Kubernetes is low level infrastructure and works best when people don't know it's there.

7

The enabling technology

Virtual kubelet (VK):

"Open-source Kubernetes kubelet implementation that masquerades as a kubelet. This allows Kubernetes nodes to be backed by Virtual Kubelet providers"

Knoc project has been the trigger of our interest.



All in all the flow we want to achieve



Currently working on

<u>interLink</u>

We extended the VK solutions with a first draft of a generic API layer for delegating pod execution on ANY remote backend

k8s POD requests are digested through the API layer (e.g. deployed on an HPC edge) into batch job execution of a container.

Custom plugins



Building the first prototype

- We have a first implementation ready for a first round of validation
 - A standalone "remote docker run" implementation for testing and development
 - A SLURM quick start installation of interLink components
 - Validated already on different pilot systems (see later)

A development guide for creating custom plugin (like the slurm one) is available.

Dick Start

• Fastest way to start using interlink, is by deploying a VK in Kubernetes using the prebuilt image:

kubectl create ns vk
kubectl kustomize ./kustomizations
kubectl apply -n vk -k ./kustomizations

• Then, use Docker Compose to create and start up containers:

docker compose -f docker-compose.yaml up -d

- You are now running:
 - A Virtual Kubelet
 - · The InterLink convice
 - A D Quick-start: InterLink
- Submit a Install binaries

kubec curl -sfL https://cloud-pg.github.io/interLink/itwinctl.sh | sh -s - install

Start daemons

curl -sfL https://cloud-pg.github.io/interLink/itwinctl.sh | sh -s - start

Restart daemons

curl -sfL https://cloud-pg.github.io/interLink/itwinctl.sh | sh -s - restart

Stop daemons

curl -sfL https://cloud-pg.github.io/interLink/itwinctl.sh | sh -s - stop

What it takes from user/framework POV

- Virtual nodes are available as normal / nodes BUT
 - They are not included into the unconditioned scheduling
- Pods created by either users or frameworks have to EXPLICITLY indicate a k8s node affinity and toleration
 - This is the intended behavior since those are not general purpose resources, but rather computing specific ones
- NO shared FS are available:
 - Configmaps, secrets and empty dirs are the only available volume kind

t	est-vk	NotReady	agent	6m39s	
1	ega-new-vk _	Ready	agent	17d	
	piVersion: v1				
ŀ	ind: Pod				
n	etadata:				
	name: test-pod-	-cfg			
	namespace: osco	ar			
	annotations:				
	slurm-job.knd	oc.io/flags: "	job-name=test	-pod-cfg -t 2800ntasks=8	nodes=1
5	pec:	-			
	restartPolicy:	OnFailure			
	containers:				
	- image: docker	://busybox:late	st		
	volumeMounts:	:			
	- name: foo				
	mountPath:	"/etc/foo"			
	readOnly: t	true			
	env:				
	- name: BOO	OKKEEPING			
	value: ke	eepmehere			
	command:				
	- sleep				
	- infinity				
	imagePullPoli	icy: Always			
	name: busyech	10			
	dnsPolicy: Clus	sterFirst			
A	nodeSelector:				
	kubernetes.ic	/hostname: vega	I-new-vk		
	tolerations:				
	- key: virtual-	node.interlink/	'no-schedule		
	operator: Exi	ists			
	volumes:				
	- name: foo				

Status of pilot implementations

Infrastructure: the following sites were federated under a common K8s cluster deployed on cloud resources

- Vega EuroHPC
 - Deployed interLink Slurm layer on a VM on the edge of the HPC
- Juelich
 - Deployed "remote docker run" interLink on a login node

(b NA te	pase) → ~ } ME est-pod-cfg-	<pre>cubectl get p -x-oscar-x-os</pre>	ood -n oscar READY scar 1/1	test-pod-c: STATUS Running	fg-x-c RESI 1	oscar-x-oscar CARTS AGE 24h	-o wide IP 127.0.0.1	NODE vega-new-vk
ciangottin 1e Sep 26	id@inter 09:41:12 JOBID 4086045	twin ~]\$ s 2023 PARTITION cpu	NAME test-pod	ne -i5 USER ciangott	ST R	TIME 1-00:41:14	NODES 1	NODELIST(REASON cn0130

Applications

- Oscar at Vega
 - Serverless execution of a pod on Vega has been functionally tested
 - A lot of useful feedback already \rightarrow looking forward to the next round of iteration
- MLFlow at Juelich
 - A ML container has been successfully delegated to docker on the login node at Juelich



- A first version of the interTwin offloading system has been prototyped and the fruitful collaboration with Vega and Juelich represented a key to the success for this early implementation
 - Along with the tests of real workflows
- Lessons learnt (technical):
 - Managing runtime environment in a transparent way is possible... but a challenge i.e.
 - volume and Apptainer/Singularity
 - Site resources access policies
 - FS permission etc..
- Next steps and priorities
 - Further integration and enhancements with Data-Lake (i.e. Rucio ecosystem)
 - Offloading integration with INDIGO-PaaS integration
 - Extend testbeds and pilots within interTwin communities

BACKUP





- The initial idea is a model where each k8s cluster admin is responsible for authorizing access to the remote resources
 - On the remote side the HPC can allow connection only for "trusted" admins/VOs
 - Additional custom policies can be integrated at the plugin level
 - e.g. contacting IdP for additional user info
 - Requesting additional information from pod annotation (e.g. the username of the payload owner)





- Current implementation support generic OAuth2 access token via Authorization header
 - Thanks to <u>oauth2 proxy</u> instantiated automatically on the edge service
 - it can be configured to accept connection only from services presenting a token with a valid group claim (configurable) and audience



Orchestrating deployment on cloud

Compute Resource orchestration is a critical feature to **implement infrastructure federation functionality**

- to automate the deployment and setup of user-defined services
- make decisions regarding provider selection

Identification of the best Cloud provider by utilising a variety of features based on both static and dynamic metrics.

- The static metrics can be defined in advance
- Dynamic information collected at runtime and continuously updated. Intertwin enhancement





How the cloud dynamic orchestration and federation integrates with the offloading mechanism?

- With interTwin we plan to enhance the orchestration system for the workload offloading mechanism
 - on-demand deployed services will be automatically configured to offload payloads.

For the system to successfully implement the workload offloading mechanism, it is essential to include and rely on a detailed information system that collects and publishes information about the different heterogeneous providers in the federation.

E.g. we plan to use the interLink APIs layer to send telemetry/traces

The overall configuration should be fully transparent to the end users. Keeping the possibility to go "advanced mode" and tweak manually the desired parameters.



The design has been discussed and agreed (see the deliverable)

The plan is to start with a semi-manual setup

- We deploy services via PaaS orchestrator
- We configure the services (by hands initially) to offload

This is a procedure to be adopted on the pilots (see later)

A key is that WP6 and users start defining requirements (i.e. TOSCA templates/Helm charts/whatever for the services to be deployed)