Istanziazione e utilizzo di batch system on demand su infrastrutture Cloud. L'esempio pilota dell'esperimento AMS Welcome and Introduction

### Daniele Spiga spiga@infn.it



**EOSC-hub** 

Dissemination level: Public/Confidential *If confidential, please define:* Disclosing Party: (those disclosing confidential information) Recipient Party: (to whom this information is disclosed, default: project consortium)



EOSC-hub receives funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No. 777536.





Intro:

some history and the preamble of EU project Quick overview of AMS02 and its computing

Integration activities motivations and objectives what we did

**Results and future** 





The early stage of the presented activity was in the context of the project INDIGO-DataCloud.

- Just thoughts, discussions and ideas

We started real activity in the context of EOSC-hub project (2018)

- in the framework of DODAS Thematic Service (see later)
  - one of the main objective was (and is) the exploitation towards new communities

## **Thematic Services of EOSC-hub**

Thematic services of EOSC-hub are service which come from external partners willing to collaborate with the project.

**EOSC-hub** 

 The project has the objective of involving early adopters, stakeholder etc... aiming at reaching out to new user groups and service providers.

### DODAS is one of the TS

marketplace@e	osc-hub.eu		USC	-nub		<b>Istituto Nazionale di Fisio</b> Log in Contact						
<b>E</b> O	SC-hut	0				CART 1 Product						
сом	IPUTE DATA	& STORAGE	PLATFORMS	THEMATIC SERVICES	IDENTITY & SECURITY	PROFESSIONAL SERVICES						
🕈 🔵 Thematic se	ervices Dynamic	On Demand An	alysis Service									
DYNAMIC ON ANALYSIS SEF	I DEMAND RVICE	Dyna	mic On	Demand Ana	lysis Service							
Component Metadata Infrastructure		Simpli	Simplify the access and management of computing resources									
		DODAS a clouds to	DODAS acts as cloud enabler designed for scientists seeking to easily exploit distributed and heterogeneous clouds to process data. Aiming to reduce the learning curve as well as the operational cost of managing									
DARIAH Dynamic On Den Service	nand Analysis	commun provisioni	ty specific servic ng, creating, ma	ces running on distributed naging and accessing a p	l cloud, DODAS completel ool of heterogeneous com	y automates the process of aputing and storage resources.						
ENES Climate An	alvtics Service	DODAS p	rovides:									
LifeWatch ERIC Plants Identification App		<ul> <li>A comp centers resource</li> </ul>	<ul> <li>A comprehensive approach to opportunistic computing, with the possibility of orchestrate multiple centers (e.g. campus facilities, public or private clouds, to gather all available computing and storage resources).</li> </ul>									
OPENCoastS		<ul> <li>A simple absorb</li> </ul>	e solution for ela peaks of usage.	stic computing site exten	sions, e.g. extension of allo	cated resources in order to						
WeNMR Suite for Structural Biology		<ul> <li>An easy and controlled procedure to dynamically instantiate a spot 'Data Analysis Facility', for example a mission specific site. This is meant as the generation of an ephemeral WLCG-Tier as a Service to share computing and data resources with collaborators.</li> </ul>										
VIEWED PRO	DUCTS	<ul> <li>The sup cloud re</li> </ul>	sources.	r condor batch systems ar	id BigData Platform on de	mand over multi-backend laas						
	WeNMR Portal	How to a	ccess the servic	<b>e</b> vice, users are first require	d to register online.							
	WeNMR is a Virtual	See the o	nline manual:									
	Community	https://do	das-ts.github.io/	dodas-doc/								

#### Perugia, 25-28 Novembre 2019

a Nuclear

# EOSC-hub Reminder: What is DODAS in a nutshell











courtesy of Matteo Duranti

### installed on the International Space Station, ISS, on May 19, 2011

- operations 24h/day, 365d/year, since the installation
- 300k readout channels + 1500 temperature sensors
- acquisition rate up to 2kHz

EOSC-hub

- more than 600 microprocessors to reduce the rate from 7 Gb/s to 10 Mb/s
- 4 Science Runs (DAQ start/stop + calibration) per orbit: I Science Run = ~ 23 minutes of data taking
- On May 2019, ~135 billion triggers acquired
- 35 TB/year of raw data



**The AMS Experiment and AMS-02 detector** 







- Data collected by AMS (~ 35 TB/year) are transmitted from the International Space Station, through the NASA infrastructure, to CERN
  - Copied at INFN-CNAF (master copy on tape)
- Reconstruction step, called std, produces ~ 100 TB/year
  - About a reconstruction step/year due to improved condition (which in turn is due to improved detector understanding)
    - example: last reprocessing required > 400 CPU-Years, to process > 5 years of AMS data taking
- Data reconstruction performed on regional centers distributed among: Europe, China and USA
  - workflow is adapted, by hand, to the various regional centers
  - job submission is managed by the LSF or PBS batch systems





- Analysis groups skim the reconstructed data producing ntuples
  - typically a production step takes about 10 CPU-Years
- AMS Scientists generate also huge amount of Monte Carlo datasets
  - as example in 2015 AMS used about 8000 CPU-Years, in 2016 about 11000
- Moreover the expectation is that these numbers will increase as far as the experiment will collect data
  - improving the understanding of the detector newer and updated Monte
     Carlo are required
    - statistics must follow the increased statistics of the real data





- Stable, massive, resources:
  - CERN
  - CNAF
  - etc...
- Additional "stable" resources:
  - ASI
  - Small campus facilities (e.g "farm di Sezione")
- Temporary "free" resources: Opportunistic Computing
  - Any grant e.g. User grant and/or Funding Agency grant
  - Computing center @ China
  - Any provider from any partner of the collaboration





Credits Matteo Duranti

- AMS computing model has not originally designed to cope with cloud computing paradigm
- AMS software is not automatically portable
- Costs for deploying AMS Environment on new resources might be to high
  - setting-up environment requires huge effort wrt obtained gain
  - keep in mind the limited effort of personnelle fully committed on computing
- it does not result trivial to exploit modern computing infrastructures etc
- All this appears to be a limiting factor in the scope of the data analyst daily work
  - the value is not in setting up computing...





Need a technology to enable the following:

- Generation of on-demand batch systems for data processing
  - User friendly solution to exploit "any cloud platform"
- Exploitation of opportunistic computing
  - A solution to manage stateless sites
    - Resources w/o permanently dedicated support for a specific experiment and/or activity;
- Elastic extension of existing facilities/batches
  - e.g To absorb peaks of resource usage;





We realized that we had the technology to start filling the gaps...

As anticipated: we decided to try to use DODAS Thematic Service originally implemented for CMS and adapt to AMS needs

- the underlying technology (the building blocks) has been originally developed to be experiment independent
- we first performed a requirement analysis (see later)
- Toward the end of 2017 we started prototyping the integration as feasibility study











### Data Input:

Read from eosams or eospublic@CERN and from gpfs@CNAF

- the "official" AMS files
- the ancillary files

### Compute

Need a Batch System. No special requirements here

### Working Station/User interface

Need a place where the user keeps his working directories (As it is now for Italian users at CNAF)

### Data Output

Write to eospublic / eosams@CERN or to gpfs@CNAF via XRootD

### Software

Read the "static" libraries needed by the executable, common to all the users Read the "user" libraries needed by the executable





- Data:
  - Data I/O is performed using XRootD. This technology enable, by construction, the possibility to implement federation, a key to use multiple data sources
     Data caching, a key to avoid latency issues while reading data from remote location
- Batch System:
  - HTCondor is the most suitable solution (many facilities including CNAF migrating, AMS already uses HTCondor at CERN etc)
    - moreover, we consider this a good choice also because offer many handles to implement federation.
    - we used it as overlay technology (see later)





- Software
  - Several options have been evaluated and we concluded to proceed with CVMFS as the most suitable technical solution
    - experiment software already managed through CVMFS
    - user/group "custom" libraries stable enough (changes are >> hourly)
    - however this requires a piece more in the infrastructure and a action more from user perspectives
- Working Station:
  - we decided to rely on CNAF for the early integration
  - No needs for third solution although we could easily implement a automated solution to create a dedicated UI

## **EOSC-hub** The first end-to-end integration





## EOSC-hub And looking under the hood





#### spiga@infn.it

Perugia, 25-28 Novembre 2019

# **EOSC-hub** A first scale test during 2018



- A first integration finished during August
  - Remote submission of analysis jobs from CNAF AMS-User Interfaces: condor submit –pool XXX
    - Very same UI used to submit to LOCAL LSF
  - Input Data read from AMS EOS @ CERN
  - Jobs run@TSystem (Open Telekom Cloud)
  - Produced output copied on a temporary storage@TSystem (XRootD)
    - Third party copy from XRootD@OTC to srm://storm-fe-ams.cr.cnaf.infn.it:8444
- A scale test during Aug.-Sept. performed with real workflow (Matteo Duranti
  - 30+TB of produced data (ntuple production)
  - About 500K jobs





## **Implemented scenarios**



- Batch system on demand
  - opportunistic resources
  - exploiting our own (INFN) resources
    - ReCaS@Bari and Cloud@CNAF
- Exploitation of ASI resources
  - see later for details
- Moreover we worked on federation of distributed resources through overlay HTCondor
  - ReCaS@Bari and OpenStack@Perugia

EOSC-hub	AMS resources landscape	tetato Radenale di Facia Radene
<ul> <li>Stable, mas</li> <li>CERN</li> <li>CNAF</li> <li>etc</li> <li>Additional s</li> <li>ASI</li> <li>Small camp</li> <li>Temporary</li> <li>Any grant e</li> <li>Any provide</li> </ul>	essive, resources: stable resources: bus facilities (e.g "farm di Sezione") "free" resources: Opportunistic Computing e.g. User grant and/or Funding Agency grant er from any partner of the collaboration	

## EOSC-hub Computing resources@ASI



22

At the Space Scientific Data Center (SSDC) of the Italian Space Agency (ASI) we have an AMS farm: 384 cores, 90TB

- there is no dedicated support to AMS, researchers take care of sysadmin roles
- access to resources only via local batch
  - data must be copied manually etc.

AMS wanted to minimise the computing effort... so the idea was to extend resources management via Openstack (as a zone) and use DODAS to implement a batch system



## The underlying implementation



Il modello proposto si basa fondamentalmente su un overlay network per Il **management dei servizi di OpenStack** geograficamente distrubuiti

EOSC-hub

- Sicurezza:
  - Comunicazione remota dei servizi OS.
     Crittografia forte
  - Applicazioni utente finale: tutto può rimanere come oggi. Stesse policy autorizzative
    - Ovviamente il setup proposto offre altri gradi di libertà
- Bandwidth:
  - L'overhead rispetto al setup attuale è minimo (solo traffico di funzionamento OS)

Nota: L'obiettivo del proposal è che il traffico sul canale INFN-PG 
SSDC serva **solamente** al management di OpenStack







Everything is transparent. Just see more worker nodes

- We deployed (automatically) a cluster where HTCondor services run on Openstack @PG and WorkerNodes (startd) are in ASI

$\leftrightarrow$ $\rightarrow$ $\bigcirc$ $\bigcirc$ $\bigcirc$ Non sicuro   ope	enstack.fis	ica.unipg.it/horizor	n/project/instances/					\$	Y 💏	🐵 💽 🗢	<b>b</b> 🙆 🚳 🐻 (	ə 🕶 🔥 🖲 👹 :
🛗 App 🔇 Condividi su Face 🛅 Tab -	aperte 🗎	🛾 Conferenze 🗎 /	ams 🛅 dampe 🗎	HERD 🛅 Explot	ech 🛅 SW	CERN 🛅 INFN 🛅 Unive	rsità 🛅 F	isica 🛅 CNAF		otizie 🛅 Bar	nche 🔇 Ultime notizie	» 🛅 Altri Preferiti
FisGeo & INFN Perugia	■ ams -	i										🕹 duranti 👻
Progetto			Nome istanza = 👻				FILTRO	AVVIA ISTANZA		ELIMINA ISTANZE MORE ACTIONS •		
COMPUTE		Nome Istanza	Nome dell'Immagine	Indirizzo IP	Dimensione	Coppia di chiavi	Stato	Zona di Disponibilità	Task	Stato attivazione	Tempo a partire dalla creazione	Actions
Panoramica		userimage- 155939428316	ubuntu-16.04- multinet2	192.168.0.167	N/A	im-d2a20cec-846d-11e9-9238- 0242ac120002	Attivo	asi	None	In esecuzione	3 giorni, 21 ore	CREA ISTANTANEA +
Istanze Volumi		userimage 155939428316	ubuntu-16.04- multinet2	192.168.0.165	N/A	im-d2a20cec-846d-11e9-9238- 0242ac120002	Attivo	asi	None	In esecuzione	3 giorni, 21 ore	CREA ISTANTANEA +
Immagini Accesso e sicurezza RETE		userimage- 155939432068	ubuntu-16.04- multinet2	ams-net <sup>192.168.0.166</sup> infn- farm	m1.medium	im-e8ff6782-846d-11e9-96bf- 0242ac120002	Attivo	nova	None	In esecuzione	3 giorni, 21 ore	CREA ISTANTANEA •
ORCHESTRAZIONE				193.204.89.79								

## **Federation through HTCondor Flocking**

AKA a native HTCondor federation solution

EOSC-hub

- easy to implement from the HTCondor perspectives
  - few lines of config file (flock from, flock to)
  - plus AuthN/Z
- It is transparent from the user point of view
  - just keep submitting jobs as before





**Batch** 





Openstack cloud provider in Bari and In Perugia

PaaS federation layer allow to build container based HTCondor overylay







Remote submission and/or federation of distributed pools requires authentication

- So far we rely GSI which impy X509 certificates
  - we rely on INDIGO-IAM as a building block in this context
    - this start from a OIDC Token and has capability to deliver also X509

### This is a key point because:

- allows integration with legacy system
- does not require any extra knowledge to the user
- it provide the status of art, which means we are ready for integrating services and systems supporting OIDC (Json Web Token).
  - we will see in practice in the next days. We'll discuss a bit in detail on Thursday



## In Summary, what happened so far...



- **Resources exploited** (rough extimation):
  - 700k jobs (2-3 hours each 1000 cores \* 2 months ≃ 15-20% of resources available at CNAF for AMS)
  - 30TB generated data (copied to CNAF)
- Used Infrastructures:
  - HelixNebula Science Cloud (TSystem)
  - Google Cloud
  - Cloud@CNAF
  - ReCaS@Bari
  - OpenStack@PG
  - OpenStack@PG-ASI (not yet at scale)







The main gap (this might be a bit INFN biased) is to improve the federation strategy

- Both in term of policies and also from an architectural point of view
  - global view of harvested resources
- Monitoring and accounting still to be improved
- better documentation....

The main objective for the future is to consolidate the current work and possibly to extend it to future experiments and collaboration

- this is under discussion (very preliminary)

# Thank you for your attention!

*Questions*?



Seosc-hub.eu ♥ @EOSC\_eu



This material by Parties of the EOSC-hub Consortium is licensed under a Creative Commons Attribution 4.0 International License.

Contact