

DEEP Hybrid-DataCloud

Intensive computing techniques for applications needing specialised hardware

Workshop: Creating Platform-Driven E-Infrastructure Innovation On EOSC 10 July 2019

<u>Mario David</u>, Cristina Duma, Valentin Kozlov and Alessandro Costantini

On behalf of all the partners of DEEP



DEEP HybridDataCloud



- Designing and Enabling E-Infrastructures for intensive data Processing in a Hybrid DataCloud
- Started as a spin-off project (together with eXtreme DataCloud - XDC) from INDIGO-DataCloud technologies
- H2020 project, EINFRA-21 call
- Runs November 1st 2017 April 2020
- 9 academic partners + 1 industrial partner:
 - CSIC, LIP, INFN, PSNC, KIT, UPV, CESNET, IISAS, HMGU, Atos



DEEP-HybridDataCloud



- Goal: prepare a new generation of e-Infrastructures that harness latest generation technologies, supporting deep learning and other intensive computing techniques to exploit very large data sources
- Global objective: promote the use of intensive computing services by different research communities and areas, and the support by the corresponding e-Infrastructure providers and open source projects
- **Ease and lower** the entry barrier for non-skilled scientists
 - **Transparent execution** on e-Infrastructures
 - Build ready to use modules and offer them through a catalog or marketplace
 - Implement common software development techniques also for scientist's applications (**DevOps**)



DEEP architecture





DEEP core components



Component	Workpackage
Kubernetes	WP4
Mesos/Marathon/ Chronos	WP4
OpenStack nova-lxd	WP4
udocker	WP4
PaaS Orchestrator	WP5
Orchent	WP5
Infrastructure Manager	WP5
CLUES-indigo	WP5
TOSCA types and templates	WP5, WP6
Monitoring System	WP5
CloudProviderRanker	WP5
Cloud Information Provider	WP5
INDIGO Virtual Router	WP5
Alien4Cloud - DEEP	WP6
DEEPaaS API	WP6

Documentation Configuration Deployment

Development & Upstream Contribution



IaaS components

Hybrid DataCloud

udocker: improve support of GPUs, improve support of low latency interconnects (Infiniband)

Openstack: nova-lxd - add/improve support of GPUs

Mesos/Marathon/Chronos: authn/authz with oidc/oauth, documentation/configuration support for GPUs and Infiniband, recipes for on-demand deployments.

Kubernetes: same as mesos

HPC Integration Tools: Provides the way to access HPC resources from the PaaS Orchestrator - containers in HPC through udocker



PaaS components



- Develop/Improve support of deployments requiring GPUs and Infiniband
- Better support for the hybrid multi-site deployments
- Provide Tosca recipes for all required deployments



DEEPaaS components

2



Alien4Cloud: TOSCA template composition & deployment

C BB S ⊕ localhost:8088/#/editor/application/TestAppl	/environment/47/cafbd-702e-4ce4-9b5f-d1783902043a/archiv ≥ ♥ ■ admin ● •
Environment Topology Editor (0.1.0-SNAPSHOT)	<>> <> <> <> <> <> <> <> <> <> <> <> <>
Compute	• Selected node × Compute · Type: Compute • · Properties = = =
	Capabilities Scalable Scalable O(O) Q, ⊕ ⊕ max_instances 1 O Q, ⊕ ⊕ default_instan 1 O

DEEP Open Catalog/Marketplace: provides the universal point of entry to all services offered by DEEP.

DEEP Open Catalog

Open Catalog!			
d is a project that aims to deliver a fra ed on your local laptop, on a producti	mework to easily develop Machine Learning and Deep Learning modules on top of e-Infrastructu ion server or on top of e-Infrastrucutres supporting the DEEP-Hybrid-DataCloud stack.	res. In the DEEP Open Catalog you can find ready to use modules in a variety of domain	
rt	2. Download & execute	Browse all modules	
ble modules and pick the one of	With Docker:		
amer is a good starting point to be (but it does not provide any	stocker search deepdic stocker rum -11 - p 5000/ deepdic/deep-oc-generic-container bocker possible? Try udscker instead stocker rum -3 6609/bid deepdic/deep-oc-generic-container wed GPU access? bee noda-docker together with Docker: stocker rum -11 - p 5000/bid@ deepdic/deep.oc.generic-container	DEEP OC Conus Classification services, library/tensorflow, library/lasagne, docker A trained Xception net on Tensorflow/Keras to classify conus maine snails. KNCW MOXE #	DEEP OC Massive Online Data services, docker Massive Online Data Streams analysis.
			DEEP OC Phytoplankton • services, library/tensorflow, library/lasagne, docker A trained Xception net on Tensorflow/Keras to classify phytoplankton.
		DEEP OC Plant Classification	DEEP OC Retinopathy

DEEPaaS components

Github





9

SQA, Release, Maintenance, Support and Testbeds





DEEP core components: CI/SQA phase



"A set of Common Software Quality Assurance Baseline Criteria for Research Projects" http://hdl.handle.net/10261/160086



- 1. Code fetching
- 2. Code style check
- 3. Unit testing coverage
- 4. SLOC metrics gathering
- 5. Functional and integration testing
- 6. Code Review
- 7. Documentation
- 8. Automated Deployment
- 9. Security linter/scanner
- 10. Vulnerability check on dependencies
- 11. Delivery
- 12. Notifications





SQA controls

Java Google

Style

Python

PEP8

Code style standards used by DEEP-HDC products

Golang

Ansible

Style

YAML

specificatio

n

OpenStack

Style





Hide details

Revert



Products



Unit testing

Alien4Cloud-DEEP
 Cloud-Info-Provider
 DEEPaaS API
 IM
 PaaS Orchestrator
 Orchent
 TOSCA Types & Templates
 udocker
 Router



Technical documentation

Docs » Technical documentation

These pages contain technical notes software documentations, guides, tutorials, logbooks and similar documents produced with DEEP Hybrid DataCloud project

Mesos

- Introduction
- Testbed Setup
- Prepare the agent (slave) node
 Testing Chronos patch for GPU support
- Testing GPU support in Marathon
- Running tensorflow docker container
- References
- Enabling open-id connect authentication

Kubernetes

- DEEP : Installing and testing GPU Node in Kubernetes CentOS7
- Installing GPU node and adding it to Kubernetes cluster

OpenStack nova-lxd

- OpenStack nova-lxd installation via Ansible
- Deploying OpenStack environment with nova-lxd via DevStack
- Installing nova-lxd with Juju
- OpenStack nova-lxd testing configuration

uDocker

• uDocker new GPU implementation

Miscelaheous

· CBU charing with MDC

Bandit (Python): security check

Code style check

Python PEP8 (4 - 40%)

Golang (1 - 10%)

Java Google Style (2 - 20%)

OpenStack Style (1 - 10%)

YAML specification (1 - 10%)

Ansible Style (1 - 10%)

Back to v1.7.5 bandit

hardcoded_bind_all_interfaces: Possible binding to all interfaces. Test ID: B104 Severity: MEDIUM Confidence: MEDIUM File: <u>IM/config.py</u> More info: https://bandit.readthedocs.io/en/latest/plugins/b104_hardcoded_bind_all_interfaces.html

56 XMLRCP_PORT = 8899

- 57 XMLRCP_ADDRESS = "0.0.0.0"
- 58 ACTIVATE_REST = False

~	SonarCloud Code Quality check passed; 62.2% Est. post-merge coverage	Details
/	continuous-integration/jenkins/branch This commit looks good	Details
/	continuous-integration/jenkins/pr-merge This commit looks good	Details
/	continuous-integration/travis-ci/pr The Travis CI build passed	Details
1	continuous-integration/travis-ci/push The Travis CI build passed	Details
/	security/snyk - pom.xml (concept-reply-it) No new issues	Details
	t6pc-bot Test	Details

github PR

alberto-brigandi merged commit 7002da1 into master on 12 Nov 2018





DEEP core components: Deployment Pilot Preview





DEEP core components: Software releases





DEEP-1/Genesis release



DEEP - 1 (Genesis) QC reports

Product	VCS	vcs	License	CodeStyle	Unit Testing (%)	Functional Testing	Integration Testing	Docs	Code Review	Automated Deployment	Security	Artefacts
	code	tag										
Alien4Cloud-"plugin"	0	1.1-r0	0	0		*	1.4	0	0	🚖 = manual	0	docker image: 🥑
Cloud-Info-Provider- DEEP	0	0.10.4	ø	0	✓ = 97 (lowest: 86)	0	-	0	0	★ = manual	0	rpm & debs: 🥑
DEEPaaS API	0	0.1.0	0	0	= 99 (lowest 95)	0	0	0	0	0	0	docker image:
IM	0	1.7.5	0	0	= 95 (lowest 75)	0	0	0	0	0	*	docker image: 🕑 rpm & debs: 🕑
PaaS Orchestrator	0	2.1.1 -FINAL	0	0	= 75 (lowest 71)	0	ø	0	0	★ = manual	0	docker image: 🥑
Orchent	0	1.2.2	0	*	*	0	0	0	0	0	0	rpm & debs: 🥑
TOSCA Types & Templates	0	3.0.0	0	N/A	N/A	*	*	*	0	0	-	tarballs:
udocker	0	1.1.3	0	0	= 92 (lowest 62)	0	0	0	0	0	0	tarballs:
vRouter	0	DEEPv1	0	1.00						0		Ansible playbook

Pinned Tweet

DEEP Hybrid-DataCloud @DEEP_eu - Jan 18 The first software release of our project, codenamed #DEEPGenesis, is out! We deliver a comprehensive architecture for #AI #DeepLearning #MachineLearning as #PaaS and #SaaS components that are orchestrated thanks to @indigodatacloud solutions

deep-hybrid-datacloud.eu/2019/01/18/dee...





Upstream contributions

Upstream Code	Link	Partner		
Tosca-parser	https://github.com/openstack/tosca-parser/commit /3af43cb9a88863ccb7f8991ca379163fe23f33e3 https://github.com/openstack/tosca-parser/commit /23cd991e884b67d41378a69ee30ef38b5d760e68 https://github.com/openstack/tosca-parser/commit /c08022d0b71ca7936e084bc40805af0b89b724ae	UPV		
Apache-libcloud	https://github.com/apache/libcloud/pull/1215 https://github.com/apache/libcloud/pull/1242 https://github.com/apache/libcloud/pull/1269			
Apache OpenWhisk	https://github.com/apache/incubator-openwhisk-devtools/pull/165 https://github.com/apache/incubator-openwhisk-devtools/pull/162 https://github.com/apache/incubator-openwhisk-devtools/pull/161	CSIC		
cloud-info-provider	https://github.com/EGI-Foundation/cloud-info-provider/pull/137 https://github.com/EGI-Foundation/cloud-info-provider/pull/126 https://github.com/EGI-Foundation/cloud-info-provider/pull/119	CSIC		

User communities and applications



Citizen Science: Plant classification Image Classification

For training and testing image classifiers (CNNs; TensorFlow). From this model the following services are derived:

- Plants (dataset: up to 1 TB)
- Conus marine snails
- Seeds
- Phytoplankton







Earth Observation: Satellite Imagery

Explore application of **Machine Learning** for satellite imagery (e.g. remote object detection, terrain segmentation, meteorological prediction).

Currently being developed is **super-resolution service** to upscale low resolution bands to high resolution with deep learning (e.g. DSen2; TensorFlow. Dataset: ca. 1 TB)





super-resolution



User communities and applications



Biological and Medical Science: Retinopathy

Diabetic retinopathy is a fast-growing cause of blindness worldwide. The use-case focuses on a deep learning approach (CNNs; TensorFlow) to automated classification of retinopathy based on color fundus retinal photography images (DR=0 (=healthy) .. 4 (blind)). Dataset: ca. 100 GB

Computing Security: Massive Online Data Streams: Online analysis of data streams

Intrusion detection systems: provide an architecture able to analyze massive on-line data streams, also with historical records, in order to generate alerts in real-time. Based on proactive time-series prediction adopting artificial neural networks (e.g. LSTM, GRU; TensorFlow). Dataset: 100 GB currently, then up to 2 TB /day

Physics: Post-processing

Of HPC simulations (Lattice QCD): analysis of a large number of configurations for Lattice QCD simulation. Move the configurations to long-term storage, perform checks and metadata operations. Requirements: Infiniband, data of 1 TB





DR=0





dataset, prediction (train), prediction (test)

DEEP serves different users' profiles





DEEP from a user point of view





DEEP CI/CD for user applications



Continuous Delivery



Development and integration in cloud resources.

Jenkins pipeline for user applications: CI/SQA: **Code style**, **security scan**.

CD:

Immediate availability of application.Automatic building (Docker images).Automatic publishing (Docker Hub).Notification (email to developers).



DEEP user applications





Achievements: users perspective



- Encapsulation and isolation of different environments (using container)
 - Enable operativity in different infrastructure
- **DEEPaaS** as an entry points with **flexible design**
 - Allow different training arguments for each use-case
- DEEP DS template and DEEPaaS API
 - ready-made template that facilitates standardisation and (semi)automatic creation of necessary files and codes



- DEEP OC software automation DevOps: CI/CD pipelines for user applications
- DEEP Leading Interaction
 - A good way to provide technical support for users as well as accompanying "HowTo" documentation with details
- Udocker extended support
 - Older Linux kernels are also supported and does not require root privileges
 - Functionalities were critically valuable in test running docker images at an old local cluster.
- Orchestrator Dashboard: WebUI where users can submit TOSCA templates to the orchestrator

Use cases status



	Problem	Goal	DEEP services	Status
Plant Classification	Automatically identify plant species from images using Deep Learning	Perform image classification on different datasets by performing the so called transfer learning	Tensorflow, Keras via DEEP PaaS Orchestrator, OIDC-agent	Deployed through orchestrator and oidc-agent in the DEEP testbed
Satellite Imagery	Explore possible applications of machine learning techniques to satellite imagery from different sources	Support a super -resolution service to upscale low resolution bands to high resolution with Deep Neural Networks	Tensorflow, Keras, Pytorch via DEEP Paas Orchestrator, OIDC-agent, Alien4Cloud	Only DEEPaas API available Will be implemented in the DEEP testbed via A4C
Retinopathy	Automated classification of retinopathy based on color fundus retinal photography images	Perform deep learning approach for image classification	Tensorflow via DEEP PaaS Orchestrator, OIDC-agent	Deployed through orchestrator and OIDC-agent in the DEEP testbed
MODS	Intrusion detection systems: provide an architecture able to analyze massive on-line data streams, also with historical records, in order	Generate alerts in real-time using ML and DEEP learning approaches.	Tensorflow, Keras via DEEP PaaS Orchestrator, OIDC-agent, Alien4Cloud	Deployed through orchestrator and OIDC-agent in the DEEP testbed. To be implemented via A4C



DEEP-HybridDataCloud highlights

DEEP vision & work on Software Quality Assurance



• Vision:

- We support the HLEG vision on **delivering quality software for the EOSC**
- We produced "A set of common software quality assurance baseline criteria for research projects"
 - http://hdl.handle.net/10261/160086
 - Done together with the eXtreme DataCloud and INDIGO projects
 - On 2019: Open for collaborations, collaborative document: <u>https://indigo-dc.github.io/sqa-baseline/</u>
- **Objective:** *align baseline criteria within different projects*

• Work:

- Current EOSC synergies (eXtreme-DataCloud)
 - SQA baseline
 - Automation: continuous integration and delivery for core products
 - Common library for CI/CD pipeline functionalities
 - Agile software development: jump-started from WP2 requirements
- **DEEP goes beyond**: automation techniques supporting user communities
 - Continuous integration and delivery pipelines in place: Docker Hub images re-creation triggered by changes in i) DEEPaaS and ii) application itself
 - Initial continuous deployment prototype: readiness/provision of training and inference as a service
 - Rendering and generation of the marketplace portal: leveraging (JSON) schema- validated metadata descriptions

A : Qu	set of Co ality As	ommon Sor surance Ba	ftware
Chi	ena ior	Abstract	entures and best practices
The purpose of to conform a 5 ecosystem relation of the purpose of	of this document is to def fortware Quality Assuran ated projects for the adeq at Log	ce plan to serve as a reference wit uate development and timely deliv	hin the European research very of software products
The purpose of to conform a Seconystem relation of the purpose of	of this document is to def fortware Quality Assuran and projects for the adeq at Log Date	ce plan to serve as a reference wit suare development and timely deliv Comment	hin the European research very of software products
The purpose of to conform a 5 ecosystem relation of the purpose of	of this document is to delioftware Quality Assurant ated projects for the adequated projects for the adequated by the second strength of	ce plan to serve as a reference wit juate development and timely deliv Comment First draft versi Updated criter	hin the European research very of software products.
The purpose of to conform a 5 ecosystem relation of the secosystem relation of the secosystem relation of the secosystem relation of the second of the secon	of this document is to dely Assuran ated projects for the adeq the log base of the adeq base of the log base	co plan to serve as a reference wit suare development and timely delix Comment First draft versi Updated criter	hin the European research very of software products ion ta
The purpose to conform a 5 ecosystem ref Document Source V1.0 V2.9	If this document is to the document is a to the document is a star and projects for the adequent the document of the document	co plan to serve as a reference wit uate development and timely delix Comment First draft versi Updated criter	hin the European researcher very of software products non Sa

DEEP-Genesis: 1st platform and release



- First software release and prototype platform released Jan. 2018
- More than 12 software components, 4 different services, several upstream contributions, more than 10 models in marketplace



- DEEP-Genesis: initial service catalog
 - All services are OIDC-ready and follow the AARC and AARC2 blueprint recommendations

Service	Functionalities	Preview endpoint
Visual application topology composition and deployment	 Graphical composition of complex application topologies Deployment through PaaS orchestrator 	https://a4c.ncg.ingrid.pt
ML/DL training facility as a service	 Provide continuous training and retraining of developed models 	On demand
DEEP as a Service	 Deployment of DEEP Open Catalog components as server-less functions Provide inference/prediction endpoints 	(beta, internal preview only) https://deep.cloud.ifca.es/
DEEP Open Catalog	 Ready-to-use machine learning and deep learning applications, including: Machine learning frameworks + JupyterLab ML/DL ready to use models BigData analytic tools 	https://marketplace.deep-hybrid-datacloud.eu



- Collaboration with **EINFRA-21** projects:
 - **eXtreme-DataCloud**: Integration of data management solutions (XDC) and computing solutions (DEEP), exploiting event driven executions. Work on software quality.
 - **DARE**: Provide ML/DL services to integrate into workflows
- Collaboration with **other** initiatives:
 - **EOSC-Hub**: integration of developments into production tools (cloud-info-provider, TOSCA-Parser).
 - **EGI.eu**: Improved support for accelerators in Cloud resources
 - Developments merged upstream
- Collaborations with **external** user communities:
 - Royal Botanical Garden of Madrid, LifeWatch ERIC, Mouse Motor Lab (Rowland Institute Harvard), Centre for Automatic and Robotics (CSIC)

Selected DEEP early results

Hybrid DataCloud



- B. M. Nguyen, H. Phan, D. Q. Ha, G. Nguyen. "An Information-centric Approach for Slice Monitoring from Edge Devices to Clouds", Procedia Computer Science Volume 130, 2018, Pages 326-335. DOI: 10.1016/j.procs.2018.04.046
- Published articles by user communities not in the project, exploiting DEEP-HybridDataCloud software components:
 - I. Heredia Cacha. Application of a Convolutional Neural Network for image classification to the analysis of collisions in High Energy Physics. CHEP 2018 Conference, Sofia, Bulgaria. Oral Contribution
 - L. Lloret; I. Heredia; F. Aguilar; E. Debusschere; K. Deneudt; F. Hernández. Convolutional Neural Networks for Phytoplankton identification and classification. Biodiversity Information Science and Standards. 2018. Oral Contribution
 - F. Pando; I. Heredia; C. Aedo; M. Velayos; L. Lloret; J. Calvo. Deep learning for weed identification based on seed images. Biodiversity Information Science and Standards. 2018. Oral Contribution

Sector Sector



Thank You



- Serverless framework planned to be developed
 - based on OpenWisk platform
- Further development of FLAAT:
 - FLAsk support for handling OIDC Access Tokens
- Extend the CI/CD pipeline:
 - Include deployment, testing of produced application Docker images from its DockerHub repository.
- General improvements:
 - Improve documentation based on feedback from users.
 - Perform training actions: through dedicated video conferences and webinars
 - Individual support through 1-to-1 TeleConferences
- Development Docker Image (DDI)