

DevOps for Machine Learning in Academia

A personal view based on DEEP Hybrid DataCloud experience and HEAT project

Valentin Kozlov

STEINBUCH CENTRE FOR COMPUTING, KIT



GridKa School 2019 - The Art of Data

26-30 August 2019 KIT, Campus North, FTU



source: en.wikipedia.org

www.kit.edu

Outline



- What is DevOps?
- Why DevOps in Academia?
- DevOps techniques
- Examples:
 - DEEP Hybrid DataCloud
 - HeAT
- Summary

What is DevOps?





DevOps is a set of software development practices that combine software development (Dev) and information technology operations (Ops) to shorten the systems development life cycle while delivering features, fixes, and updates frequently in close alignment with business objectives.

... Academics and practitioners have not developed a unique definition for the term "DevOps"

ATLASSIAN

DevOps is a set of practices that **automates** the processes between software development and IT teams, in order that they can build, test, and release software **faster** and more **reliably**.

Keywords:

- set of practices
- release software faster and reliably
- *automation* of the process

Product pipeline







Product manager

Product pipeline

Developer

Tester

Karlsruhe Institute of Technology



Dev +Tests + Ops \rightarrow Work together





Iterative developments with a releasable product at every iteration

7 30-Aug-2019 V. Kozlov – DevOps for Machine Learning in Academia



Iterative developments with a releasable product *at every iteration*

8 30-Aug-2019 V. Kozlov – DevOps for Machine Learning in Academia



Iterative developments with a releasable product *at every iteration*

9 30-Aug-2019 V. Kozlov – DevOps for Machine Learning in Academia

Why DevOps for Scientific Development?



a set of practices

DevOps as

- for a *reliable* software to be delivered *faster*
- a 'single person skills' to maintain the product pipeline

this is ... what we often do/need in science:

DevOps		Science
Faster	\rightarrow	Art of experimentation
Reliability of software	\rightarrow	Reliability of software ©
'Single person' skills (supported by DevOps tools)	<i>></i>	 e.g. a master / PhD project Often a domain knowledge is required to develop and test the code Overcome bottlenecks where IT ops skills needed
Automation of tests (e.g. unit tests)	\rightarrow	Better structured and readable code, so that a student project can be continued
Code Quality checkers	\rightarrow	Maintainable code
Infrastructure as a Code, Containerization	\rightarrow	Scientific reproducibility

Product pipeline \rightarrow DevOps loop







source: https://marketplace.atlassian.com/categories/devops?utm_source=wac_marketplace_landing

DevOps: key techniques

	Technique	Examples
Plan	Tools to document your plans (Agile user stories) and design	Google docs, Wikipages
Build	Version Control Test stack (Test coverage)	GitHub
CI / CD	 Continuous Integration (CI) : Code merged in repository → Automated builds and code tests performed Continuous Delivery (CD) : CI success → Deploy in testing environments and perform integration tests Continuous Deployment (CD'): CD success → Deploy to production environment (e.g. as a stable release) 	Travis CI, Jenkins
Deploy	Cloud: Infrastructure as a code Docker containers	Puppet, Ansible udocker (HPC)
Ope- rate	Tools to monitor your application, infrastructure	Zabbix

Test stack pyramid

*) As a good first guess at Google: https://testing.googleblog.com/2015/04/just-say-no-to-more-end-to-end-tests.html

Unit testing in ML

- https://medium.com/@keeper6928/how-to-unit-test-machine-learning-code-57cf6fd81765
- Machine learning testing framework for TensorFlow: <u>https://github.com/Thenerdstation/mltest</u>

Eliminates:

- Sometimes variables are initialized but neural network layer is not (properly) connected, i.e. corresponding weights will not be updated
- Some layers (functions) have to be set into trainable state
- \rightarrow You check that all trainable variables are indeed updated after one iteration.

We may rephrase it as ...

It is better to have a simple DevOps pipeline than no pipeline

Advance your pipeline iteratively

- Analyze your current way of software development
- Think what improvement would bring you the most benefit
- Implement

Example #1

https://deep-hybrid-datacloud.eu/

https://marketplace.deep-hybrid-datacloud.eu/

https://docs.deep-hybrid-datacloud.eu

https://www.youtube.com/playlist?list=PLJ9x9Zk1O-J_UZfNO2uWp2pFMmbwLvzXa

DEEP Hybrid DataCloud: context

Designing and Enabling E-Infrastructures for intensive data Processing in a Hybrid DataCloud

- Started as a spin-off project (together with XDC) from INDIGO-DataCloud technologies
- H2020 project, EINFRA-21 call
- Runs November 1st 2017 April 2020
- March 27, 2019 : mid-term Review (Luxembourg)
- 9 academic partners + 1 industrial partner: CSIC, LIP, INFN, PSNC, KIT, UPV, CESNET, IISAS, HMGU; Atos

Ease and lower the entry barrier for scientists

in using intensive computing techniques, e.g. deep learning, to exploit very large data sources

- Build ready to use modules and offer them through an open catalog or marketplace (Docker images)
- Implement common software development techniques also for scientist's applications (DevOps)
- Transparent execution on e-Infrastructures (pilot testbed, orchestration)

DEEP pilot use-cases

Image Classification

Plants, Conus marine snails, Seeds, Phytoplankton (CNN, TF)

Satellite Imagery

Super-resolution service (e.g. DSen2; TF)

Karlsruhe Institute of Technology

Retinopathy

DL to analyze color fundus retinal photography images (CNN, TF)

DEEP Leading Interaction: dog's breed detection

to test the infrastructure and provide working examples

Massive Online Data Streams: Online analysis of data streams Intrusion detection systems (LSTM/GRU, TF)

DevOps tools from a DEEP user perspective

Data Science Template

- 1. \$ pip install cookiecutter
- 2. \$ cookiecutter https://github.com/indigo-dc/cookiecutter-data-science
- 3. answer questions
- Two directories created: <user_project> and <DEEP-OC-user_project>
 - This are also git repositories.
 - Both have 'master' and 'test' branches.
 - Dockerfile, Jenkinsfile, and python files necessary for DEEPaaS API integration are pre-populated
 - .stestr.conf, tox.ini for testing in Python

<user_project></user_project>	<deep-oc- user_project></deep-oc-
data	Dockerfile
docker	Jenkinsfile
docs	
dogs_breed_det	
models	
notebooks	E) motodoto ison
references	
reports	
.dockerignore	
juitignore	
stestr.conf	
Jenkinsfile	
README.md	
requirements-dev.txt	
requirements.txt	
Setup.cfg	
Setup.py	
test-requirements.txt	
test_environment.py	
🖹 tox.ini	SCC, KIT

Data Science Template

- 1. \$ pip install cookiecutter
- 2. \$ cookiecutter https://github.com/indigo-dc/cookiecutter-data-science
- 3. answer questions
- Two directories created: <user_project> and <DEEP-OC-user_project>
 - This are also git repositories.
 - Both have 'master' and 'test' branches.
 - Dockerfile, Jenkinsfile, and python files necessary for DEEPaaS API integration are pre-populated
 - .stestr.conf, tox.ini for testing in Python

<user_project></user_project>	<deep-oc- user_project></deep-oc-
🖿 data	Dockerfile
docker	Jenkinsfile
docs	
dogs_breed_det	
models	
notebooks	docker-compose.yml
references	metadata.json
reports	
Jockerignore	
.gitignore	
stestr.conf	
Jenkinsfile	
	-
README.md	
requirements-dev.txt	
requirements.txt]
🖹 setup.cfg	-
a setup.py	
test-requirements.txt	
test_environment.py	

🖹 tox.ini

DEEPaaS API

get_metadata()	Return metadata from the exposed model	predict_file(path, **kwargs)	Perform a prediction from a file in the local filesystem
get_train_args(*args)	Function to expose to the API what are the typical training parameters	predict_data(*args, **kwargs)	Perform a prediction from the data passed in the arguments.
train(*args)	Function to train your model.	predict_url(*args)	Perform a prediction from a remote UR

DEEP as a Service API endpoint [Base URL: /] http://0.0.0.0:5000/swagger.json	
DEEP as a Service (DEEPaaS) API endpoint.	
models Model information, inference and training operations	\sim
GET /models/ Return loaded models and its information	
GET /models/{model_name} Return <model_name> models metadata</model_name>	
POST /models/{model_name}/predict Make a prediction given the input data	
PUT /models/{model_name}/train Retrain model with available data	

Jenkins CI/CD pipeline

25

Code testing & Quality Control

flake8 / PEP8

unit tests

Metrics gathering

Bandit security scanner

Docker build & push to registry

docker hub :cpu-test, :gpu-test :latest, :cpu, :gpu DEEP Open Catalog

github.com/deephdc

test branch

master branch

Jenkins CI/CD reports

Recent Changes

Q Full Stage View

- coverage.py report
- Bandit report
- Open Blue Ocean
- 🔑 GitHub

#16

- PyLint Warnings
- Coverage Report
- SLOCCount Results

Embeddable Build Status

🠢 Bu	ild History	trend —
find		×
#30	Aug 11, 2019 10:32 PM	
#29	Aug 11, 2019 10:24 PM	
#28	Aug 11, 2019 10:07 PM	
#27	Aug 11, 2019 9:47 PM	
#26	Aug 11, 2019 9:39 PM	
#25	Aug 11, 2019 9:36 PM	
#24	Aug 11, 2019 9:28 PM	
#23	Aug 11, 2019 5:31 PM	
#22	Aug 11, 2019 6:48 AM	
#21	Aug 10, 2019 10:23 AM	
#20	Aug 9, 2019 10:55 PM	
#19	Aug 8, 2019 9:20 PM	
🥥 <u>#18</u>	Aug 6, 2019 9:03 AM	
#17	Aug 5, 2019 10:18 PM	

Aug 5, 2019 9:19 PM

	(Averag	e <u>full</u> run time: ~2	811
Œ	30	_	- 2
	Aug 12	1	
	00:32	commit	

Stage View

Declarative: Checkout SCM	Code fetching	Style analysis: PEP8	Unit testing coverage	Metrics gathering	Security scanner	Re-build Docker images	Declarative: Post Actions
7s	1s	2min 42s	47s	6s	28s	2min 21s	647ms
10s	3s	3min 24s	3min 38s	45s	3min 16s	16min 26s	570ms
7s	1s	3min 23s	33s	349ms	81ms	67ms	551ms

250

https://marketplace.deep-hybrid-datacloud.eu

https://marketplace.deep-hybrid-datacloud.eu

Browse all modules

DEEP Open Catalog

SCC, KIT

DEEP OC Massive Online Data Streams

G GET THIS MODULE	DEEPHDC/DEEP-OC-MODS	
y DEEP-Hybrid-DataClou	ud Consortium Created: Tue 19 Februa	ary 2019 - Updated: Wed 14 August 2019
sonicos dockor		
SEIVICES, UUCKEI		

build passing

This use case analyzes online data streams in order to generate alerts with time-bounded constrains and in real-time. The main study is focused on building additional intelligent module using NN and DL techniques in co-function with underlying Intrusion Detection Systems (IDS) supervising traffic networks of compute centers. Preserving old data for historical purposes, security analysts will be able to supervise generated alerts and to enhance cyber security [1, 2] for such centers when large IT infrastructures and devices products a huge amount of data streaming continuously and dynamically.

Run locally on your computer

Using Docker

You can run this module directly on your computer, assuming that you have Docker installed, by following those stops:

The principle of the solution is proactive time-series prediction [5] adopting NNs as well as

Deployment of artefacts

OASIS N TOSCA

Topology and Orchestration Specification for Cloud Applications TOSCA template \rightarrow Simple Profile in YAML.

Alien 4 Cloud

GUI to compose TOSCA templates, deployment in a cloud

orchent – CLI to perform deployments using TOSCA templates

PaaS Orchestrator Dashboard – Web GUI for deployments using TOSCA templates

HPC environment:

a basic user tool to execute simple docker containers in user space *without* requiring root privileges.

- a docker like CLI
- provides GPU access
- MPI / Infiniband can be used

HEAT - Helmholtz Analytics Toolkit

https://github.com/helmholtz-analytics/heat/

HeAT

31

Scientific data analytics library for HPC systems build on top of **PyTorch**

- Operates on heterogeneous hardware like GPU/CPU systems
- Allows computation on *distributed* systems
- **Distributed tensor data object**: operations like basic scalar functions, linear algebra algorithms, slicing or broadcasting operations

HeAT Use cases

Earth System Modelling

Research with Photons

Aeronautics and Aerodynamics

Structural Biology

Neuroscience

 $\alpha = 0^{\circ}$

 $\alpha = 2^{\circ}$

 $\alpha = 6$

ICIT

Karlsruhe Institute of Technology

CONTINUOUS SEEDERCH NUNICATION CONTINUOUS OPERATE BUILD Version control, issue code review tracking **GitHub GitHub**

PLAN

Travis Cl docker

INTEGRATION

DEPLOY

codecov.io

Mattermost.

sprint

GitHub

planning

for discussions

COMMUN

HeAT CI/CD pipeline

34

Summary

DevOps is a set of practices to release software faster and reliably, and automate this process as much as possible.

Key techniques are:

- *iterative development*
- version control
- test stack

- Cl/{Delivery | Deployment}
- Pipeline-as-a-Code
- Infrastructure-as-a-Code
- Monitoring

This set of practices and available tools help *a single person* to maintain the *full product pipeline* \rightarrow this is often the case in science!

While in Machine Learning it is not an easy task to run System/Performance tests, one should start with simpler (unit)tests, which already help a lot.

Credits and Links

- Agile: <u>https://en.wikipedia.org/wiki/Agile_software_development</u>
- Coursera course "Continuous Delivery & DevOps" by University of Virginia
- Online: Jez Humble: <u>https://continuousdelivery.com</u> (there're many, of course!) at medium.com: <u>https://medium.com/search?q=DevOps</u>
- Atlassian Marketplace: <u>marketplace.atlassian.com/categories/devops?utm_source=wac_marketplace_landing</u>
- GitHub Marketplace: <u>https://github.com/marketplace</u>
- Cookiecutter: <u>https://github.com/cookiecutter/cookiecutter</u>
- How to (unit) test machine learning code:
 - https://medium.com/@keeper6928/how-to-unit-test-machine-learning-code-57cf6fd81765
 - Machine learning testing framework for TensorFlow: <u>https://github.com/Thenerdstation/mltest</u>
 - TensorFlow Benchmarks: <u>https://github.com/tensorflow/benchmarks/tree/master/scripts/tf_cnn_benchmarks</u>
- Bandit security scanner: <u>https://github.com/PyCQA/bandit</u>, <u>https://pypi.org/project/bandit/</u>
- OASIS TOSCA: <u>https://www.oasis-open.org/committees/tc_home.php?wg_abbrev=tosca</u>
- udocker: <u>https://github.com/indigo-dc/udocker</u>
- Deep Learning online books: <u>livebook.manning.com/book/deep-learning-with-python</u> (F.Chollet); <u>www.deeplearning.ai/machine-learning-yearning</u> (A. Ng); <u>www.deeplearningbook.org</u> (I. Goodfellow);