

# Machine Learning as a Service for HEP on heterogeneous computing resources

**Luca Giommi**<sup>1</sup> (luca.giommi3@unibo.it), Valentin Kuznetsov<sup>2</sup>, Daniele Bonacorsi<sup>1</sup>, Daniele Spiga<sup>3</sup>

<sup>1</sup> University of Bologna and INFN Bologna, Italy

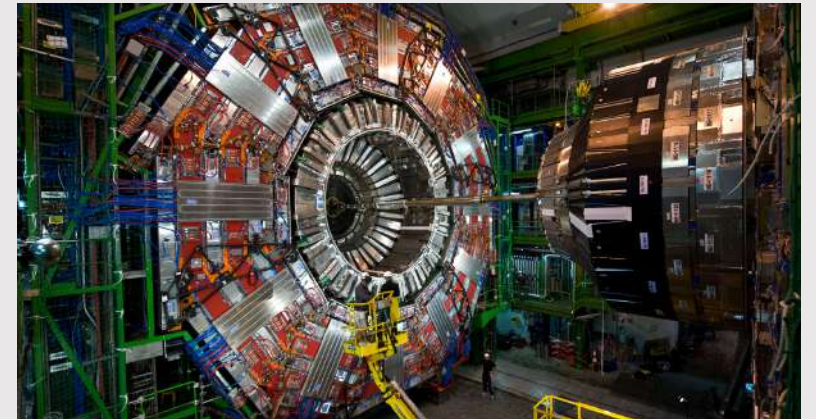
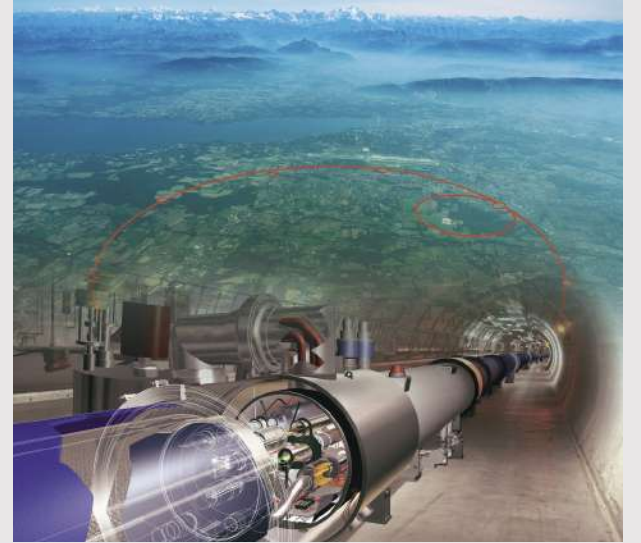
<sup>2</sup> Cornell University, USA

<sup>3</sup> INFN Perugia, Italy



# High Energy Physics at CERN

- The European Organization for Nuclear Research (CERN) was born in 1954 and it is based in the northwest suburb of Geneva on the Franco–Swiss border.
- At CERN is located the world's largest and most powerful particle accelerator called Large Hadron Collider (LHC).
- Inside the accelerator, two high-energy particle beams travel at close to the speed of light before they are made to collide at four locations, corresponding to the positions of four particle detectors: ATLAS, CMS, ALICE and LHCb.



# Why Machine Learning as a Service?

- Machine Learning techniques in the HEP domain are ubiquitous (e.g. detector simulation, object reconstruction, particles identification, Monte-Carlo generation) and will play a significant role also in the upcoming High-Luminosity LHC upgrade.
- One of the main obstacles in this scenario is the real gap among HEP physicists and ML experts, caused by the specificity of some parts of the HEP typical workflows and solutions.
- To close this gap and ease the physicists not ML-practitioners in the usage of ML techniques in their analyses, we propose a Machine Learning as a Service for HEP (MLaaS4HEP) solution as a product of R&D activities within the CMS experiment.

## Machine Learning in High Energy Physics Community White Paper

May 17, 2019

**Abstract:** Machine learning has been applied to several problems in particle physics research, beginning with applications to high-level physics analysis in the 1990s and 2000s, followed by an explosion of applications in particle and event identification and reconstruction in the 2010s. In this document we discuss promising future research and development areas for machine learning in particle physics. We detail a roadmap for their implementation, software and hardware resource requirements, collaborative initiatives with the data science community, academia and industry, and training the particle physics community in data science. The main objective of the document is to connect and motivate these areas of research and development with the physics drivers of the High-Luminosity Large Hadron Collider and future neutrino experiments and identify the resource needs for their implementation. Additionally we identify areas where collaboration with external communities will be of great benefit.

**Editors:** Sergei Gleyzer<sup>30</sup>, Paul Seyfert<sup>13</sup>, Steven Schramm<sup>32</sup>

**Contributors:** Kim Albertsson<sup>1</sup>, Piero Altoc<sup>2</sup>, Dustin Anderson<sup>3</sup>, John Anderson<sup>4</sup>, Michael Andrews<sup>5</sup>, Juan Pedro Araque Espinosa<sup>6</sup>, Adam Aurisano<sup>7</sup>, Laurent Basara<sup>8</sup>, Adrian Bevan<sup>9</sup>, Wahid Bhimji<sup>10</sup>, Daniele Bonacorsi<sup>11</sup>, Bjorn Burkle<sup>12</sup>, Paolo Calafiura<sup>10</sup>, Mario Campanelli<sup>9</sup>, Louis Capps<sup>2</sup>, Federico Carminati<sup>13</sup>, Stefano Carrazza<sup>13</sup>, Yi-Fan Chen<sup>4</sup>, Taylor Childers<sup>14</sup>, Yann Coadou<sup>15</sup>, Elias Coniavitis<sup>16</sup>, Kyle Cranmer<sup>17</sup>, Claire David<sup>18</sup>, Douglas Davis<sup>19</sup>, Andrea De Simone<sup>20</sup>, Javier Duarte<sup>21</sup>, Martin Erdmann<sup>22</sup>, Jonas Eschle<sup>23</sup>, Amir Farbin<sup>24</sup>, Matthew Feickert<sup>25</sup>, Nuno Filipe Castro<sup>6</sup>, Conor Fitzpatrick<sup>26</sup>, Michele Floris<sup>13</sup>, Alessandra Forti<sup>27</sup>, Jordi Garra-Tico<sup>28</sup>, Jochen Gemmler<sup>29</sup>, Maria Girone<sup>13</sup>, Paul Glayshe<sup>18</sup>, Sergei Gleyzer<sup>30</sup>, Vladimir Vava Gligorov<sup>31</sup>, Tobias Golling<sup>32</sup>, Jonas Graw<sup>2</sup>, Lindsey Gray<sup>21</sup>, Dick Greenwood<sup>33</sup>, Thomas Hacker<sup>34</sup>, John Harvey<sup>13</sup>, Benedikt Hegner<sup>13</sup>, Lukas Heinrich<sup>17</sup>, Ulrich Heintz<sup>12</sup>, Ben Hooberman<sup>35</sup>, Johannes Junggeburth<sup>36</sup>, Michael Kagan<sup>37</sup>, Meghan Kane<sup>38</sup>, Konstantin Kanishchev<sup>8</sup>, Przemysław Karpiński<sup>13</sup>, Zahari Kassabov<sup>39</sup>, Gaurav Kaul<sup>40</sup>, Dorian Keira<sup>3</sup>, Thomas Keck<sup>29</sup>, Alexei Klimentov<sup>41</sup>, Jim Kowalkowski<sup>21</sup>, Luke Kreczko<sup>42</sup>, Alexander Kurepin<sup>43</sup>, Rob Kutschke<sup>21</sup>, Valentin Kuznetsov<sup>44</sup>, Nicolas Köhler<sup>36</sup>, Igor Lakomov<sup>13</sup>, Kevin Lannon<sup>45</sup>, Mario Lassnig<sup>13</sup>, Antonio Limosani<sup>46</sup>, Gilles Louppe<sup>17</sup>, Aashrita Mangu<sup>47</sup>, Pere Mato<sup>13</sup>, Helge Meinhard<sup>13</sup>, Dario Menasce<sup>48</sup>, Lorenzo Moneta<sup>13</sup>, Seth Moortgat<sup>49</sup>, Meenakshi Narain<sup>12</sup>, Mark Neubauer<sup>35</sup>, Harvey Newman<sup>3</sup>, Sydney Otten<sup>50</sup>, Hans Pabst<sup>40</sup>, Michela Paganini<sup>51</sup>, Manfred Paulini<sup>5</sup>, Gabriel Perdue<sup>21</sup>, Uzziel Perez<sup>52</sup>, Attilio Picazio<sup>53</sup>, Jim Pivarski<sup>54</sup>, Harrison Prosper<sup>55</sup>, Fernanda Psihas<sup>56</sup>, Alexander Radovic<sup>57</sup>, Ryan Reece<sup>58</sup>, Aurelius Rinkevicius<sup>44</sup>, Eduardo Rodrigues<sup>7</sup>, Jamal Roric<sup>59</sup>, David Rousseau<sup>60</sup>, Aaron Sauers<sup>21</sup>, Steven Schramm<sup>32</sup>, Ariel Schwartzman<sup>37</sup>, Horst Severini<sup>61</sup>, Paul Seyfert<sup>13</sup>, Filip Siroky<sup>62</sup>, Konstantin Slazytkin<sup>43</sup>, Mike Sokoloff<sup>7</sup>, Graeme Stewart<sup>63</sup>, Bob Stienen<sup>64</sup>, Ian Stockdale<sup>65</sup>, Giles Strong<sup>6</sup>, Wei Sun<sup>4</sup>, Savannah Thais<sup>51</sup>, Karen Tomko<sup>66</sup>, Eli Upfal<sup>12</sup>, Emanuele Usai<sup>12</sup>, Andrey Ustyuzhanin<sup>67</sup>, Martin Vala<sup>68</sup>, Sofia Vallecorsa<sup>69</sup>, Justin Vasek<sup>56</sup>, Mauro Verzetti<sup>70</sup>, Xavier Vilasis-Cardona<sup>71</sup>, Jean-Roch Vlimant<sup>3</sup>, Ilija Vukotic<sup>72</sup>, Sean-Jiun Wang<sup>30</sup>, Gordon Watts<sup>73</sup>, Michael Williams<sup>74</sup>, Wenjing Wu<sup>75</sup>, Stefan Wunsch<sup>29</sup>, Kun Yang<sup>4</sup>, Omar Zapata<sup>76</sup>

[arXiv:1807.02876v3](https://arxiv.org/abs/1807.02876v3) [physics.comp-ph]

[DOI: 10.1007/s41781-018-0018-8](https://doi.org/10.1007/s41781-018-0018-8)

# Machine Learning as a Service

## CLOUD MACHINE LEARNING SERVICES COMPARISON

	Amazon	Microsoft	Google	IBM
Automated and semi-automated ML services				
	Amazon ML	Microsoft Azure ML Studio	Google Prediction API	IBM Watson ML Model Builder
Classification	✓	✓	deprecated	✓
Regression	✓	✓		✓
Clustering	✓	✓		✗
Anomaly detection	✗	✓		✗
Recommendation	✗	✓		✗
Ranking	✗	✓		✗
Platforms for custom modeling				
	Amazon SageMaker	Azure ML Services	Google ML Engine	IBM Watson ML Studio
Built-in algorithms	✓	✗	✗	✓
Supported frameworks	TensorFlow, MXNet, Keras, Gluon, Pytorch, Caffe2, Chainer, Torch	TensorFlow, scikit-learn, Microsoft Cognitive Toolkit, Spark ML	TensorFlow, scikit-learn, XGBoost, Keras	TensorFlow, Spark MLlib, scikit-learn, XGBoost, PyTorch, IBM SPSS, PMML

- MLaaS is a set of tools and services including: data visualization, pre-processing, model training and evaluation, serving predictions, etc.
- Many of the world's leading cloud providers provide different types of MLaaS services, including Amazon, Microsoft, Google and IBM
- MLaaS service providers offer pre-defined models that can be used to cover standard use-cases, e.g. classification, regression, natural language processing, facial recognition, DeepLearning

# Issues with using existing solutions

---

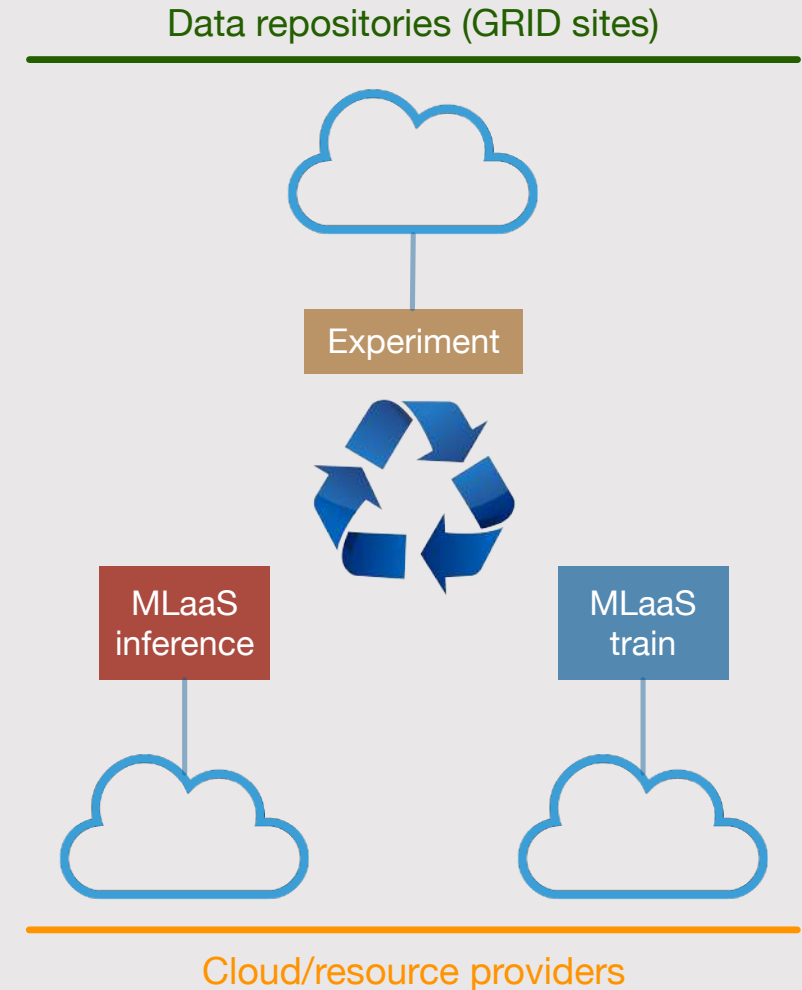
- Existing MLaaS services can't read HEP data directly in the **ROOT data-format**: most of the cases ML deal with either CSV or NumPy arrays representing tabular data
  - We don't use ROOT data directly in the ML framework, we need a conversion step
  - Pre-processing operations may be more complex than offered by service providers
- R&D for specialized solutions to speed-up inference on FPGAs, e.g. [HLS4ML](#)
  - These solutions are designed for optimization of the inference phase rather than targeting the whole ML pipeline from reading data, to training and serving predictions
- Custom solutions adopted in specific CMS analysis (e.g. [DOI: 10.1088/2632-2153/ab9023](#)) cannot easily generalized and do not represent "as a Service" solutions
- Recent solution with Spark platform for data processing and ML training ([DOI: 10.1007/s41781-020-00040-0](#)). Here data are read from the CERN EOS storage system, not allowing access to data stored in WLCG sites



# MLaaS for HEP

MLaaS for HEP aims at providing the following:

- **natively read HEP data**, e.g. be able to read ROOT files of arbitrary size from local or remote distributed data-sources via XrootD
- **use heterogeneous resources** both for training and inference, like local CPU, GPUs, farms, cloud resources, etc.
- **use different ML libs and frameworks** (Keras, TF, PyTorch, etc.)
- **serve pre-trained HEP models**, like a models repository, and access it easily from any place, any code, any framework.



# MLaaS4HEP R&D

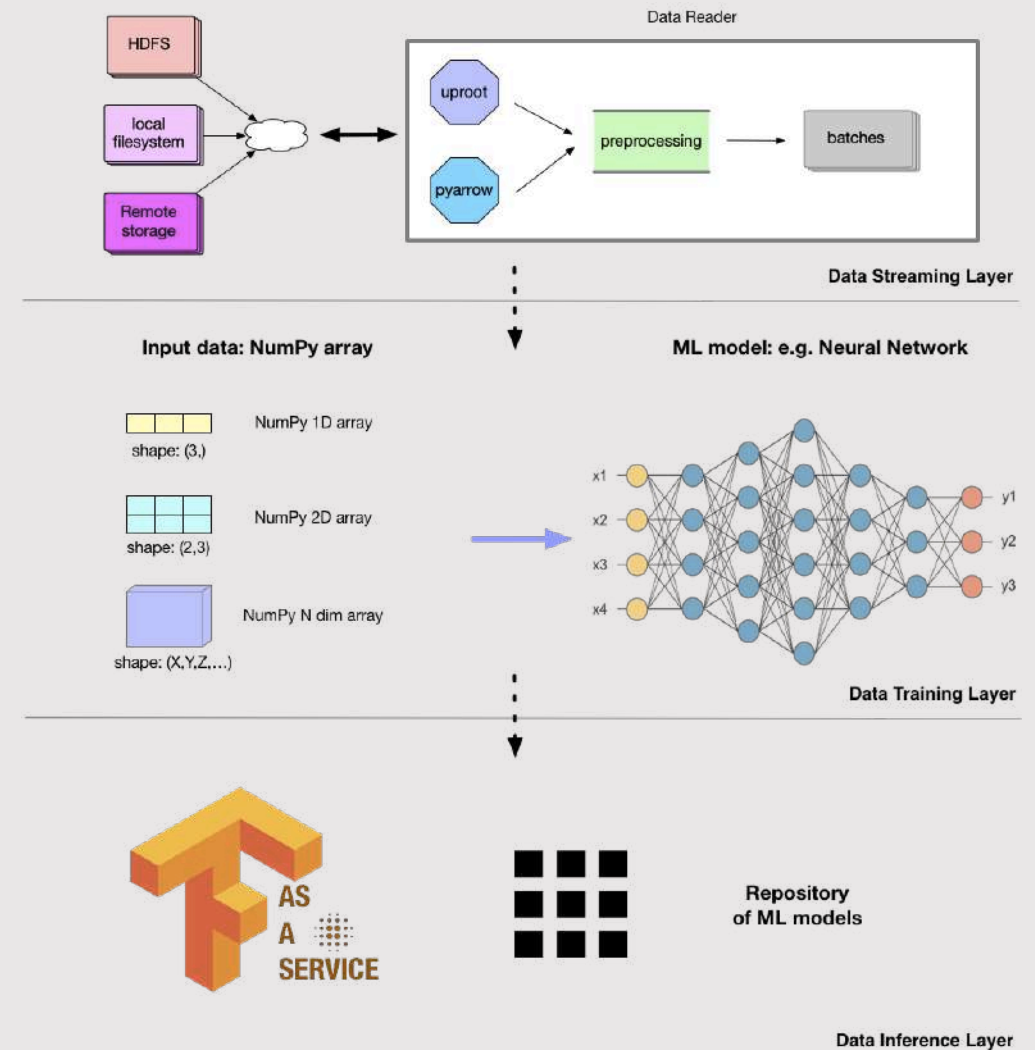
- **Data Streaming Layer** is responsible for local and remote data access of HEP ROOT files
- **Data Training Layer** is responsible for feeding HEP ROOT data into existing ML frameworks
- **Data Inference Layer** provides access to pre-trained HEP model for HEP users

All three layers are independent from each other and allow independent resource allocation

Data streaming and training tools: [github.com/vkuznet/MLaaS4HEP](https://github.com/vkuznet/MLaaS4HEP)

Data inference tool: [github.com/vkuznet/TFaaS](https://github.com/vkuznet/TFaaS)

Paper on arXiv ([arXiv:2007.14781v2](https://arxiv.org/abs/2007.14781v2) [hep-ex]) and submitted to the [CSBS journal](#)



# Data Streaming Layer

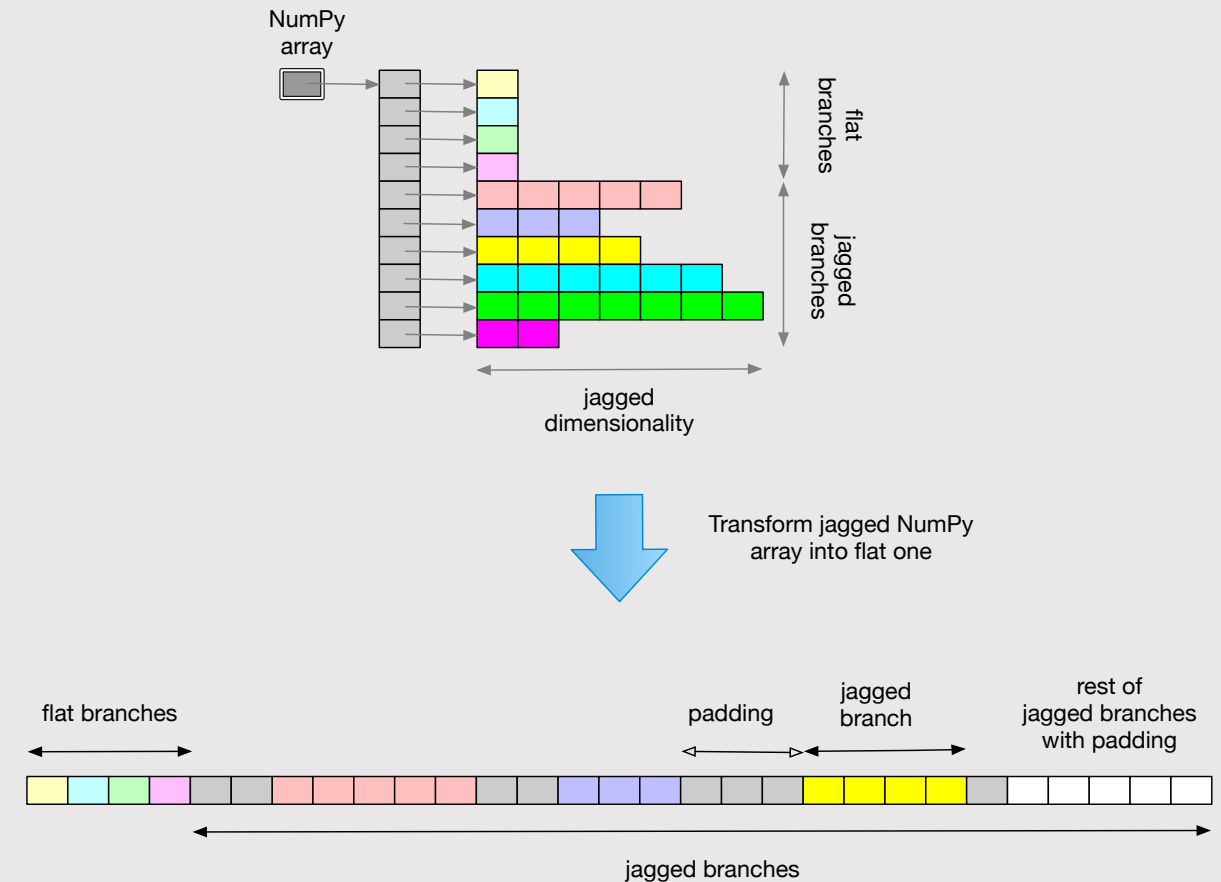
---

- The development of the DIANA-HEP **uproot** library provides the ability to read ROOT data in Python, access them as NumPy arrays, and implements XrootD access to read remote files
- MLaaS4HEP extends uproot library and provide APIs to feed data read from local and remote distributed ROOT files into existing ML frameworks
  - a Python Generator is created to read ROOT files and deliver them as chunks
    - such implementation provides efficient access to large datasets since it does not require loading the entire dataset into the RAM of the training node
    - read from multiple files is supported, taking the right proportion of data from each file
- The non-flat ROOT branches are read and represented by uproot as Jagged Arrays

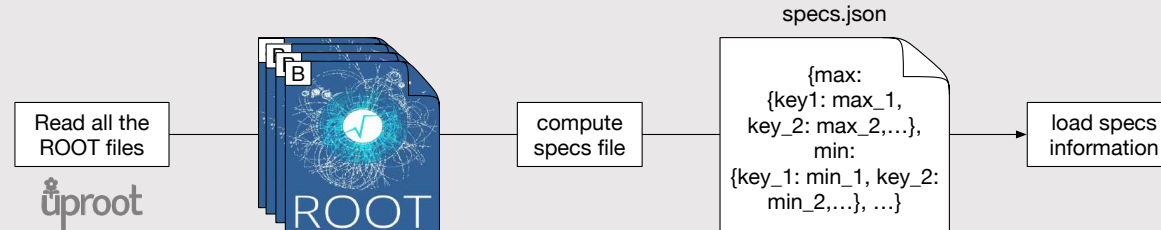


# Data Training Layer

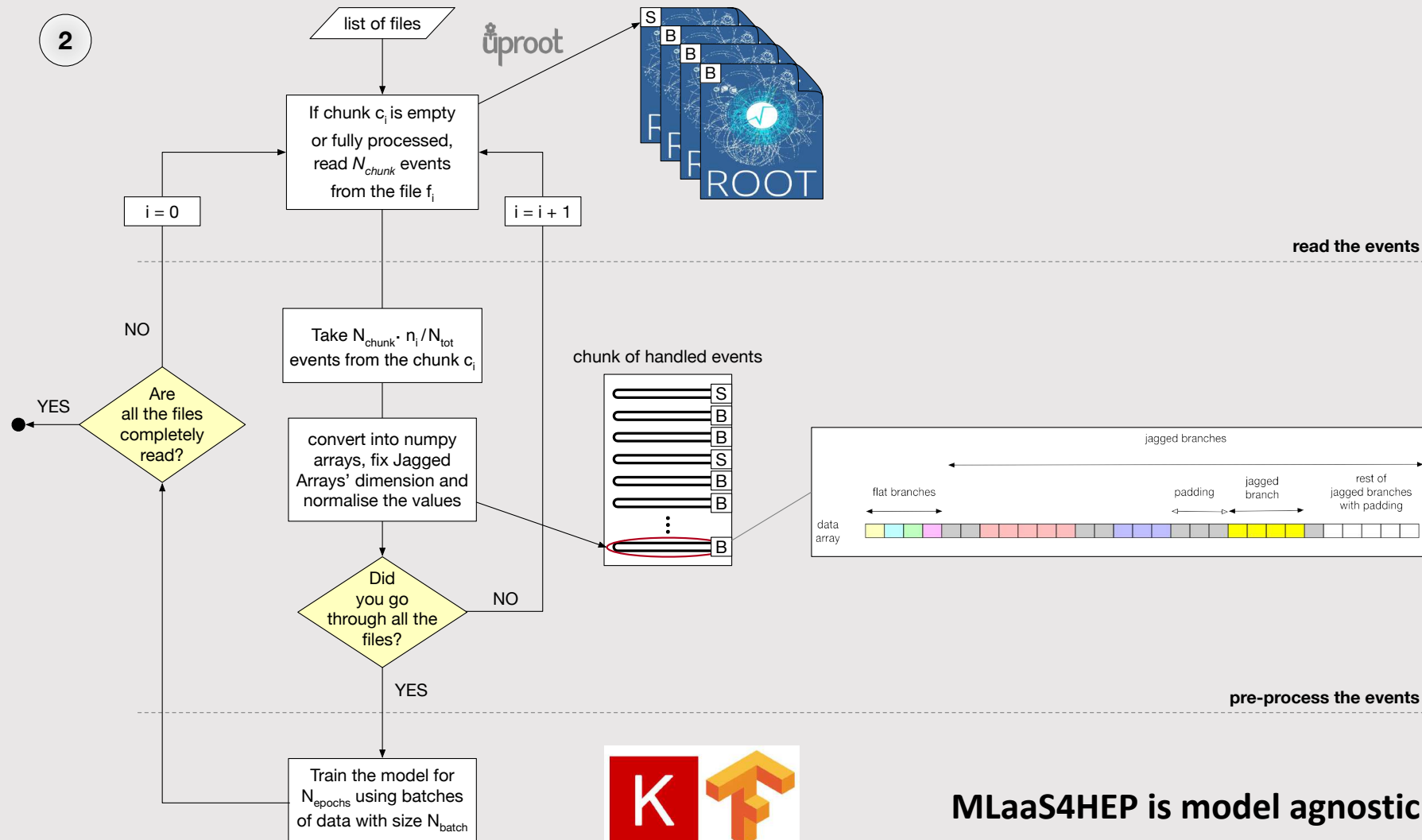
- Each event is a composition of flat and Jagged Arrays
  - Such data representation is not directly suitable for ML (dynamic dimension of Jagged Arrays across events)
- To feed these data into ML we need to resolve how to treat Jagged Arrays. We opted to flatten Jagged Arrays into fixed-size array with padding values through a two-step procedure:
  - know the dimensionality of every Jagged Array attribute;
  - update the dimension of jagged branches using padding values, which should be assigned as NaNs since all other numerical values can represent attribute spectrum.
- Keep the mask array with padding values location
- We provide a proper normalization of each attribute



1



2



MLaaS4HEP is model agnostic

train the ML model

```
./workflow.py --files=files.txt --labels=labels.txt --model=model.py --params=params.json
```

MLaaS  
workflow

Input  
ROOT files

Labels of  
ROOT files

User  
model

MLaaS  
parameters

### Keras model (model.py)

```
from tensorflow import keras
from keras.models import Sequential
from keras.layers import Dense, Dropout

def model(idim):
    "Simple Keras model for testing purposes"
    ml_model = Sequential([Dense(128,
                                activation='relu', input_shape=(idim,)),
                            Dropout(0.5),
                            Dense(64, activation='relu'),
                            Dropout(0.5),
                            Dense(1, activation='sigmoid')])
    ml_model.compile(optimizer=keras.optimizers.Adam(lr=1e-3),
                    loss=keras.losses.BinaryCrossentropy(),
                    keras.metrics.AUC(name='auc'))
```

### MLaaS parameters (params.json)

```
{
  "nevents": 30000,
  "shuffle": true,
  "chunk_size": 10000,
  "epochs": 5,
  "batch_size": 100,
  "identifier": ["runNo", "evtNo", "lumi"],
  "branch": "boostedAk8/events",
  "selected_branches": "",
  "exclude_branches": "",
  "hist": "pdfs",
  "redirector": "root://xrootd.ba.infn.it",
  "verbose": 1
}
```

### Input ROOT files (files.txt)

```
PATH/flatTree_ttHJetTobb_M125_13TeV_amcatnloFXFX_madspin_pythia8.root
PATH/flatTree_TT_TuneCUETP8M2T4_13TeV-powheg-pythia8.root
```

### Labels of ROOT files (labels.txt)

```
1
0
```

## MLaaS parameters

## Read remote root files

## Write and load the specs

```
./workflow.py --files=files.txt --labels=labels.txt --model=model.py --params=params.json
DataGenerator: <MLaaS4HEP.generator.RootDataGenerator object at 0x7f0cb58d7fd0> [29/Jun/2020:17:53:44] 1593445994.0
model parameters: {"nevents": 30000, "shuffle": true, "chunk_size": 10000, "epochs": 2, "batch_size": 500, "identifier": ["runNo", "evtNo", "lumi"],
"branch": "boostedAk8/events", "selected_branches": "", "exclude_branches": "", "hist": "pdfs", "root_director": "root://xrootd.ba.infn.it", "verbose": 1}
```

```
Reading root://xrootd.ba.infn.it//PATH_FILES/flatTree_ttHJetTobb_M125_13TeV_amcatnloFXFX_madspin_pythia8.root
# 10000 entries, 77 branches, 9.5222034454347 MB, 1.0169336795806885 sec, 9.36364252323795 MB/sec, 9.833482950553169 kHz
# 10000 entries, 77 branches, 9.53391551971455 MB, 1.2977769374847412 sec, 7.346343770133804 MB/sec, 7.705484441248654 kHz
# 10000 entries, 77 branches, 9.53866767883008 MB, 1.4104814529418945 sec, 6.7627033726234735 MB/sec, 7.089777734505208 kHz
--- first pass: 948348 events, (22-flat, 57-jagged) branches, 328 attrs
<MLaaS4HEP.reader.RootDataReader object at 0x7f840dbf4d50> init is complete in 4.852992534637411 sec
```

```
Reading root://xrootd.ba.infn.it//PATH_FILES/flatTree_TT_TuneCUETP8M2T4_13TeV-powheg-pythia8.root
# 10000 entries, 77 branches, 8.875920295715332 MB, 0.9596493244171143 sec, 9.24912889518941 MB/sec, 10.42047313071777 kHz
# 10000 entries, 77 branches, 8.868906021118164 MB, 1.2938923835754395 sec, 6.85443869497903 MB/sec, 7.728618026459661 kHz
# 10000 entries, 77 branches, 8.869449615478516 MB, 1.1267895698547363 sec, 7.87143389747739 MB/sec, 8.874771534572496 kHz
--- first pass: 1003980 events, (22-flat, 52-jagged) branches, 312 attrs
<MLaaS4HEP.reader.RootDataReader object at 0x7f8410e15f90> init is complete in 4.53512477847559 sec
```

```
write global-specs.json
load specs from global-specs.json for root://xrootd.ba.infn.it//$PATH_FILES/flatTree_ttHJetTobb_M125_13TeV_amcatnloFXFX_madspin_pythia8.root
load specs from global-specs.json for root://xrootd.ba.infn.it//$PATH_FILES/flatTree_TT_TuneCUETP8M2T4_13TeV-powheg-pythia8.root
init RootDataGenerator in 11.186564683914185 sec
```

```
label 1, file <flatTree_ttHJetTobb_M125_13TeV_amcatnloFXFX_madspin_pythia8.root>, going to read 4858 events
read chunk [0:4857] from /$PATH_FILES/flatTree_ttHJetTobb_M125_13TeV_amcatnloFXFX_madspin_pythia8.root
# 10000 entries, 77 branches, 9.52220344543457 MB, 1.3816642761230469 sec, 6.891835889507034 MB/sec, 7.237648228164387 kHz
total read 4858 evts from /$PATH_FILES/flatTree_ttHJetTobb_M125_13TeV_amcatnloFXFX_madspin_pythia8.root
```

```
label 0, file <flatTree_TT_TuneCUETP8M2T4_13TeV-powheg-pythia8.root>, going to read 5142 events
read chunk [4858:9999] from /$PATH_FILES/flatTree_TT_TuneCUETP8M2T4_13TeV-powheg-pythia8.root
# 10000 entries, 77 branches, 8.875920295715332 MB, 1.7170112133026123 sec, 5.169401473297779 MB/sec, 5.8240737873606205 kHz
total read 5142 evts from /$PATH_FILES/flatTree_TT_TuneCUETP8M2T4_13TeV-powheg-pythia8.root
```

Create the chunk

Model: "sequential"

Layer (type)	Output Shape	Param #
dense (Dense)	(None, 128)	49152
dropout (Dropout)	(None, 128)	0
dense_1 (Dense)	(None, 64)	8256
dropout_1 (Dropout)	(None, 64)	0
dense_2 (Dense)	(None, 1)	65
Total params: 57,473		
Trainable params: 57,473		
Non-trainable params: 0		

Init the ML model

Perform training cycle

Train on 7000 samples, validate on 3000 samples

Epoch 1/2

7000/7000 [=====] - 2s 220us/sample - loss: 1.5275 - auc: 0.7845 - accuracy: 0.7307 - val\_loss: 2.5731e-04 - val\_auc: 1.0000 - val\_accuracy: 1.0000

Epoch 2/2

7000/7000 [=====] - 0s 20us/sample - loss: 0.1406 - auc: 0.9883 - accuracy: 0.9543 - val\_loss: 8.8477e-06 - val\_auc: 1.0000 - val\_accuracy: 1.0000

# Data Inference Layer

---

- The Data Inference Layer is implemented as TensorFlow as a Service ([TFaaS](#)), written in the **Go** programming language
  - The Go programming language natively supports concurrency and it is very well integrated with the TF library
- TFaaS is capable of serving any TensorFlow model and it can be used as a global repository of pre-trained HEP models (it is experiment agnostic and can work with any HTTP based client)
- Both Python and C++ clients were developed on top of the REST APIs (end-points) and other clients can be developed thanks to HTTP protocol used by the TFaaS Go RESTful implementation
  - [C++ client](#) library talks to TFaaS using ProtoBuffer data-format, all others use JSON (see [examples](#))
- TFaaS allows a rapid development or continuous training of TF models and their validation: clients can test multiple TF models at the same time
  - A demo server hosted by CERN is online: <https://cms-tfaas.cern.ch/>
- From R&D to a production service: expanding the team on code maintainance and ops support

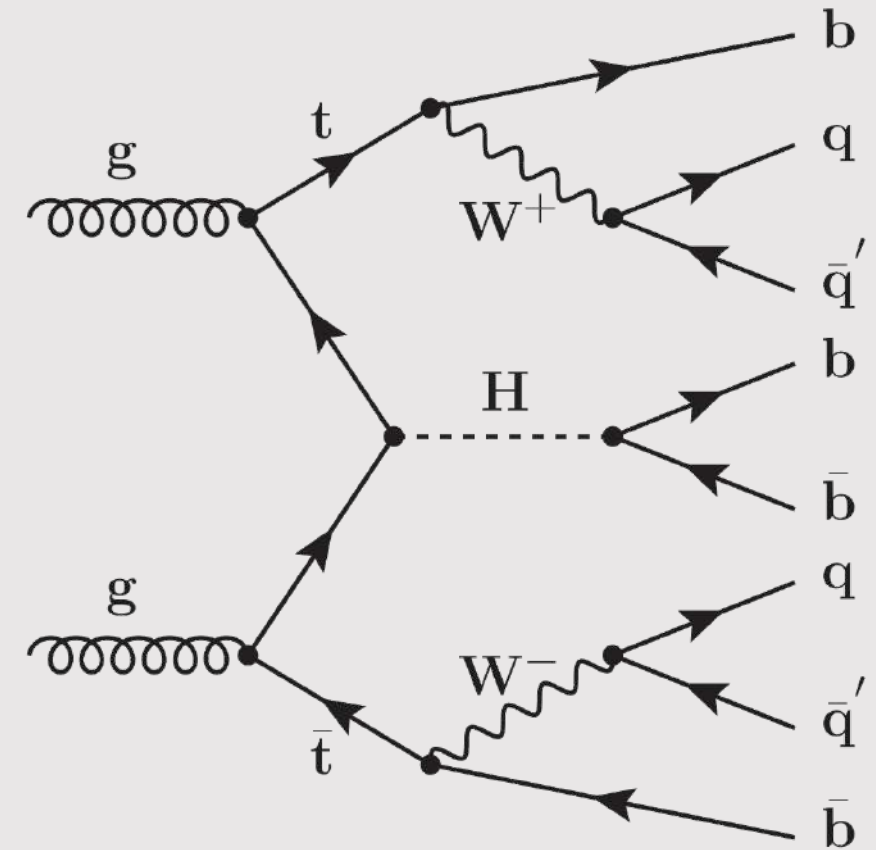


# Real case scenario:

## $t\bar{t}H(bb)$ analysis in the boosted, all-hadronic final states

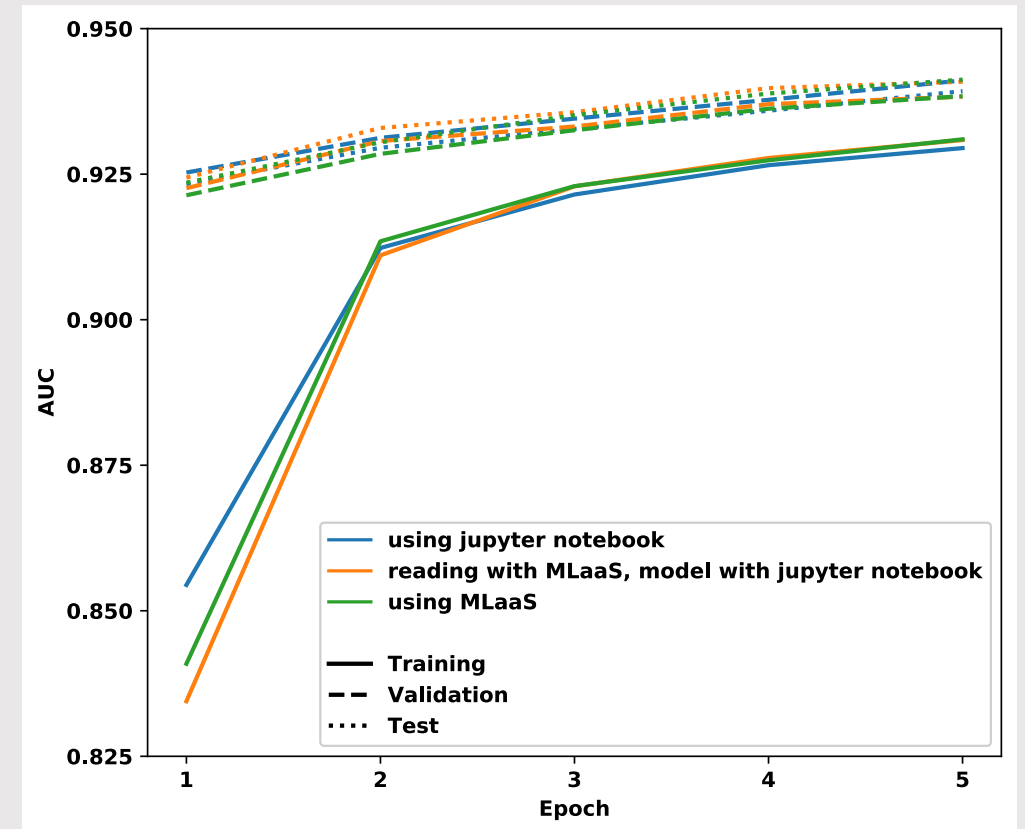
- We performed a proof-of-concept of the entire pipeline using CMS NANOAOD.
- We recently validated the MLaaS framework, namely tested the infrastructure on real physics use-case. We chose a signal vs background discrimination problem in a  $t\bar{t}H$  analysis. This allows us to:
  1. validate MLaaS results from the physics point of view
  2. test performances of MLaaS framework

For the phase of validation we used 9 ROOT files, 8 of background and 1 of signal. Each file has 27 branches, with 350 thousand events for the whole pool of files and a total size of almost 28 MB. The ratio between signal and background is 10.8%.



# MLaaS4HEP validation

- Validate the MLaaS4HEP approach by comparing it with alternative methods on the reference use-case
  - We used a simple NN with Keras in all methods
- Validation successful: physics results are not impacted
- The AUC score is also comparable with the BDT-based analysis, performed within the TMVA framework by a subgroup of the CMS HIG PAG



AUC score

# MLaaS4HEP performance

---

- For this phase we used all available ROOT files without any physics cut. This gave us a dataset with 28.5M events with 74 branches (22 flat and 52 Jagged), and a total size of about 10.1 GB.
- We performed all the tests running MLaaS framework on
  - macOS, 2.2 GHz Intel Core i7 dual-core, 8 GB of RAM
  - CentOS 7 Linux, 4 VCPU Intel Core Processor Haswell 2.4 GHz, 7.3 GB of RAM CERN Virtual Machine
- The average available bandwidth was approximately 129 Mbit/s and 639 Mbit/s using macOS and CERN VM, respectively.
- The ROOT files are read from local file-systems (SSD storages) and remotely from the Grid sites. In particular, we read files remotely from three different data-centers located at
  - Bologna (BO)
  - Pisa (PI)
  - Bari (BA)

# MLaaS4HEP performance results

---

- Based on the resource we used and if the ROOT files were local or remote, we obtained:
  - ❖ **specs computing phase (chunk size = 100k events)**
    - Event throughput: 8.4k – 13.7k evts/s
    - Total time using all the 28.5M events: 35 – 57 min
  - ❖ **chunks creation in the training phase (chunk size = 100k events)**
    - Event throughput: 1.1k – 1.2k evts/s
    - Total time using all the 28.5M events: 6.5 – 7.5 hrs
- The time to train the ML model is not included in the performance. It is independent from the MLaaS4HEP framework but depends on the underlying ML framework, the complexity of the used ML model, and the available hardware resources.
- We estimate that projecting these results for datasets at the **TB scale** and using the same hardware resources, the specs computing phase will take  $O(100)$  hours and the training phase will take  $O(1k)$  hours (plus the time required to train the ML model).
  - Further optimization of the MLaaS4HEP pipeline will be required to process TB or PB scale datasets and it may involve parallelization of I/O, distributed ML training, etc.

PRELIMINARY

# TFaaS performance

---

- We did TFaaS benchmarks on CentOS 7 Linux, 16 cores, 30 GB of RAM in two modes:
  - using 1k calls with 100 concurrent clients,
  - using 5k calls with 200 concurrent clients.
- We tested both JSON and ProtoBuffer data formats while sending and fetching the data to/from the TFaaS server.
- In both cases, we achieved a throughput of **500 req/sec**. These numbers were obtained by serving a mid-size pre-trained NN model with 27 features and 1024x1024 hidden layers used in the physics analysis discussed. Similar performance was found for image classification datasets (MNIST).
- The actual performance of TFaaS will depend on the complexity of served ML model and available hardware resources.
- Even though a single TFaaS server may not be as efficient as an integrated solution, it can be horizontally scaled, e.g. using Kubernetes or other cluster orchestrated solutions, and may provide the desired throughput for concurrent clients.

# Towards MLaaS4HEP cloudification

---

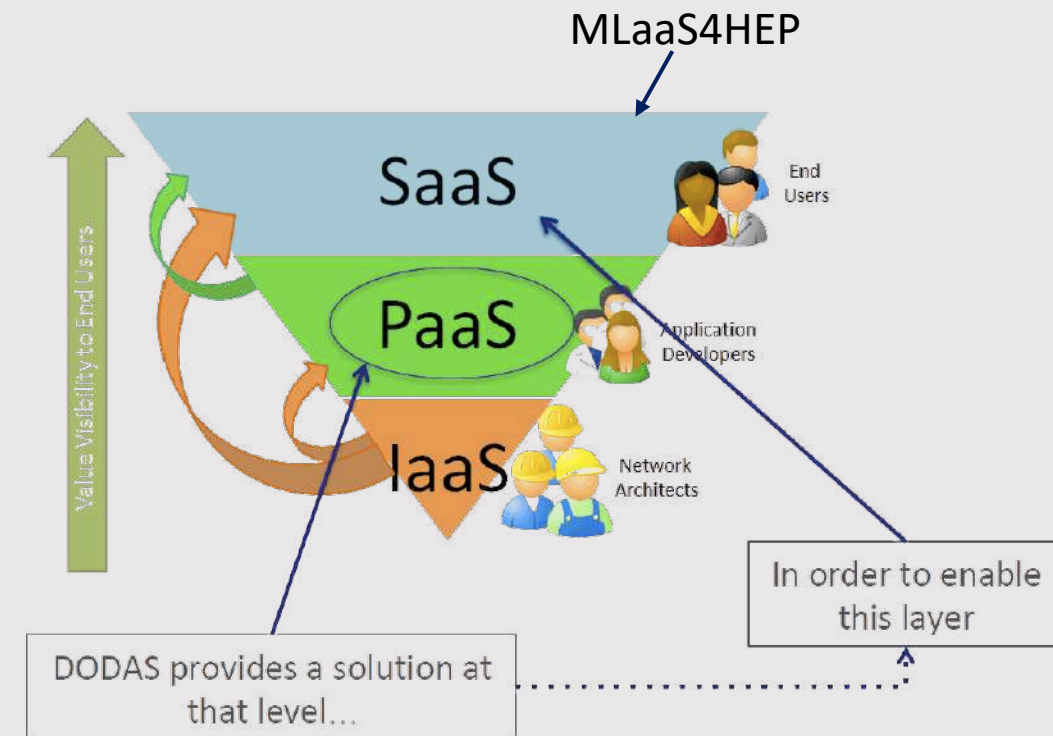
- We created a service performing the ML pipeline using local and remote ROOT files
  - The performance strictly depends on the available hardware resources
- How to improve the performance?
  - Adopt new solutions in the code
  - Invest in better and more expensive on-premise resources
  - Move to the cloud
- The operation of cloudification has two benefits:
  - opens us to potentially more performing resources
  - provides a real “as a Service” solution for the user
- We started to work for a MLaaS4HEP cloudification using DODAS





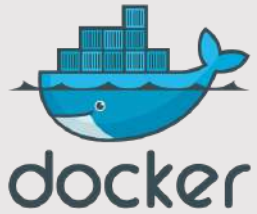
Dynamic On Demand Analysis Service (DODAS) is a Platform as a Service tool for generating over cloud resources and on-demand, container based solution.

- DODAS completely automates the process of provisioning, creating, managing and accessing a pool of distributed and heterogeneous computing and storage resources.
- DODAS has a high level of modularity, a key to a generic applicability.
  - Being modular, the architecture provides the ability to easily customize the workflow depending on the community computational requirements.
  - Implements services composition model based on templates
- Both HTCondor batch system and platform for the Big Data analysis based on Spark, Hadoop etc, can be deployed using “any cloud provider” with almost zero effort.

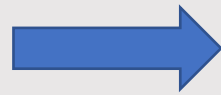


[PoS\(ISGC 2018 & FCDD\)024](#)

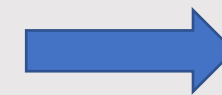
# MLaaS4HEP cloudification with DODAS



Creation of a [docker image](#) able to run the workflow.py script



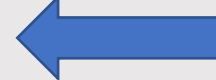
Create an Ansible playbook to automatize the configuration and deployment of the container with dependencies



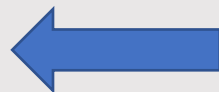
Convert the Ansible playbook into an [Ansible role](#)



Creation of a [Tosca template](#) to define the resource requirements and the input parameters for the creation of the docker container



Create the deployment from command line



Run workflow.py interactively or with jupyterhub

```
dodas create lgiommi-template.yml  
dodas login <infID> <vmID>
```



# Functional test

- We tested MLaaS4HEP using an infrastructure created with DODAS
  - Ubuntu 18 Linux, 8 AMD Opteron 62xx class CPU 2.6 GHz, 16 GB RAM VM
  - We used both local and remote (stored in the Pisa data-center) ROOT files

		MacOS and CERN VM	DODAS infrastructure
specs computing phase	Event throughput	8.4k – 13.7k evts/s	5.8k – 7.3k evts/s
	Total time using all the 28.5M events	35 – 57 min	66 – 82 min
chunks creation in the training phase	Event throughput	1.1k – 1.2k evts/s	0.5k evts/s
	Total time using all the 28.5M events	6.5 – 7.5 hrs	16 hrs

RESULTS OF THE  
FUNCTIONAL TEST

# MLaaS4HEP using Jupyterhub

- We provide a SaaS solution for a sharable jupyter notebook
- Token-based access to the jupyterhub, with the support for a customizable environment



**DEMO**

- Integrate cloud storage for managing the required files (ROOT files, ML model, etc.)

```
# . ./shared/setup_local
(base) # cd /workarea/shared/folder_test
(base) # ../../workarea/MLaaS4HEP/src/python/MLaaS4HEP/workflow.py --files=files_test.txt --labels=labels_test.txt --
model=keras_model.py --params=params_test.json
model parameters: {"nevt": -1, "shuffle": true, "chunk_size": 10000, "epochs": 5, "batch_size": 100, "identifier": ["runNo", "evtNo",
"lumi"], "branch": "events", "selected_branches": "", "exclude_branches": "", "hist": "pdfs", "redirector": "root://gridftp-storm-
t3.cr.cnaf.infn.it:1095", "verbose": 1}
Reading ttH_signal.root
# 10000 entries, 29 branches, 1.10626220703125 MB, 0.034181833267211914 sec, 32.364039645948566 MB/sec, 292.5530623775014 kHz
# 10000 entries, 29 branches, 1.10626220703125 MB, 0.022344589233398438 sec, 49.50917626973965 MB/sec, 447.53563807084936 kHz
```

# Summary

---

We built a MLaaS solution for HEP where:

- local and remote ROOT files can be directly read
- complexity of data transformation from ROOT I/O to ML frameworks is hidden from the user
- resources can be used dynamically and independently for training and inference layers
- Python based ML framework of user choice can be used (e.g. TensorFlow, Keras, Pytorch)
- customization via JSON configuration: total number of events to read; chunk size of data to read; select or exclude branches to read, choice of XrootD redirector
- we validated the MLaaS framework with a physics use-case and achieved similar results as BDT analysis
- we did performance tests for the streaming and training layers

And ...

- we moved towards the MLaaS4HEP cloudification using DODAS
  - provides interactive and jupyterhub access for an user-friendly platform where run MLaaS4HEP
  - provides compliance with the INFN-Cloud portfolio of services

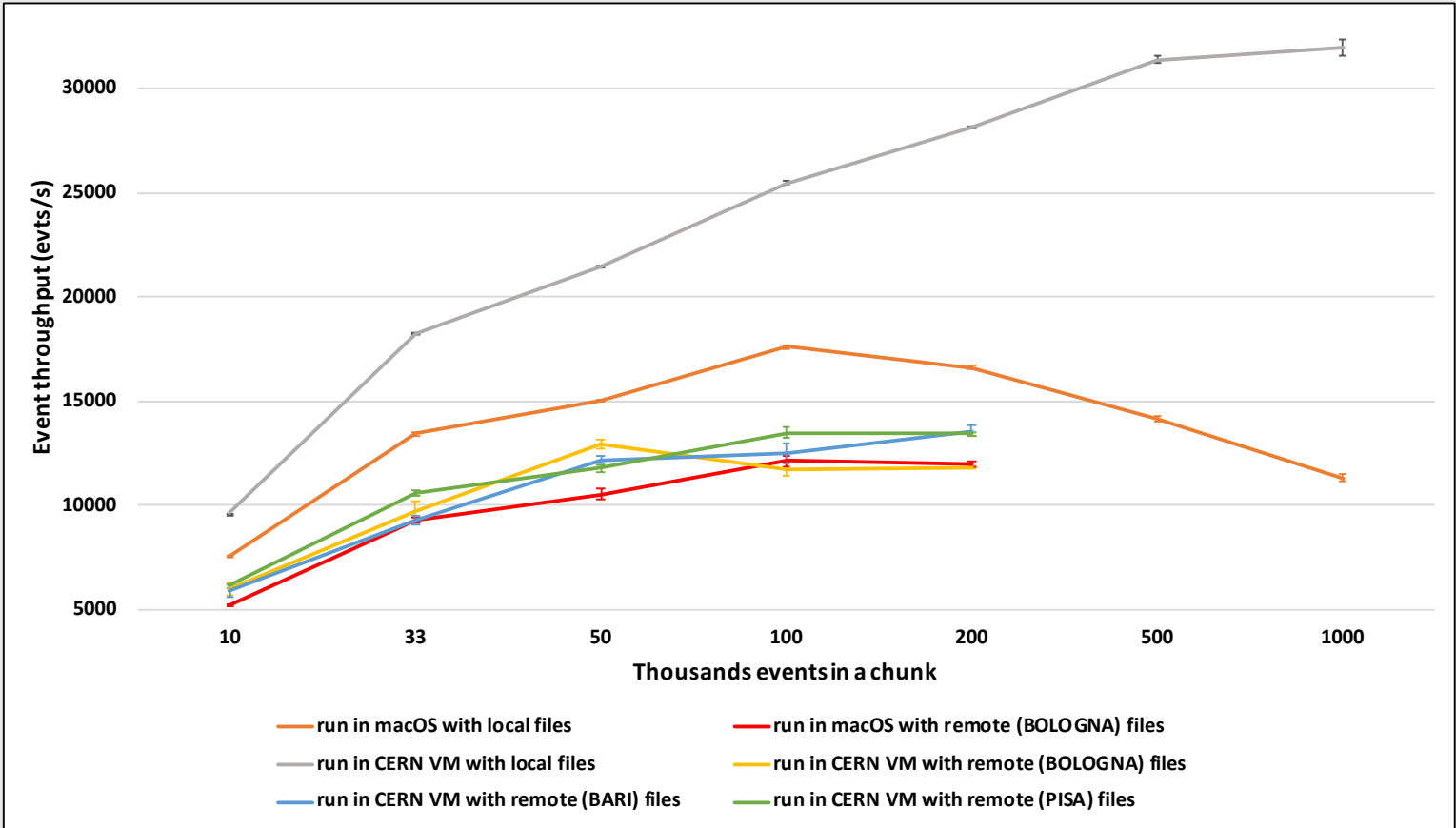
**Planning:** dynamically load user-based pre-processing functions; try FPGAs to speed up the inference phase, I/O parallelization and distributed ML training, ...

# Thanks for the attention



# Backup slides

1



	reading time (s)	specs comp. time (s)	time to complete step ① (s)	event throughput for reading + specs comp. (evts/s)
macOS with local files	1633 (9)	958 (2)	2599 (11)	11055 (49)
macOS with remote files (BO)	2365 (49)	974 (10)	3353 (57)	8585 (149)
VM with local files	1131 (3)	963 (2)	2102 (5)	13690 (34)
VM with remote files (BO)	2455 (68)	959 (2)	3427 (67)	8396 (158)
VM with remote files (BA)	2304 (88)	961 (2)	3279 (89)	8801 (241)
VM with remote files (PI)	2129 (41)	1044 (78)	3186 (83)	9047 (228)

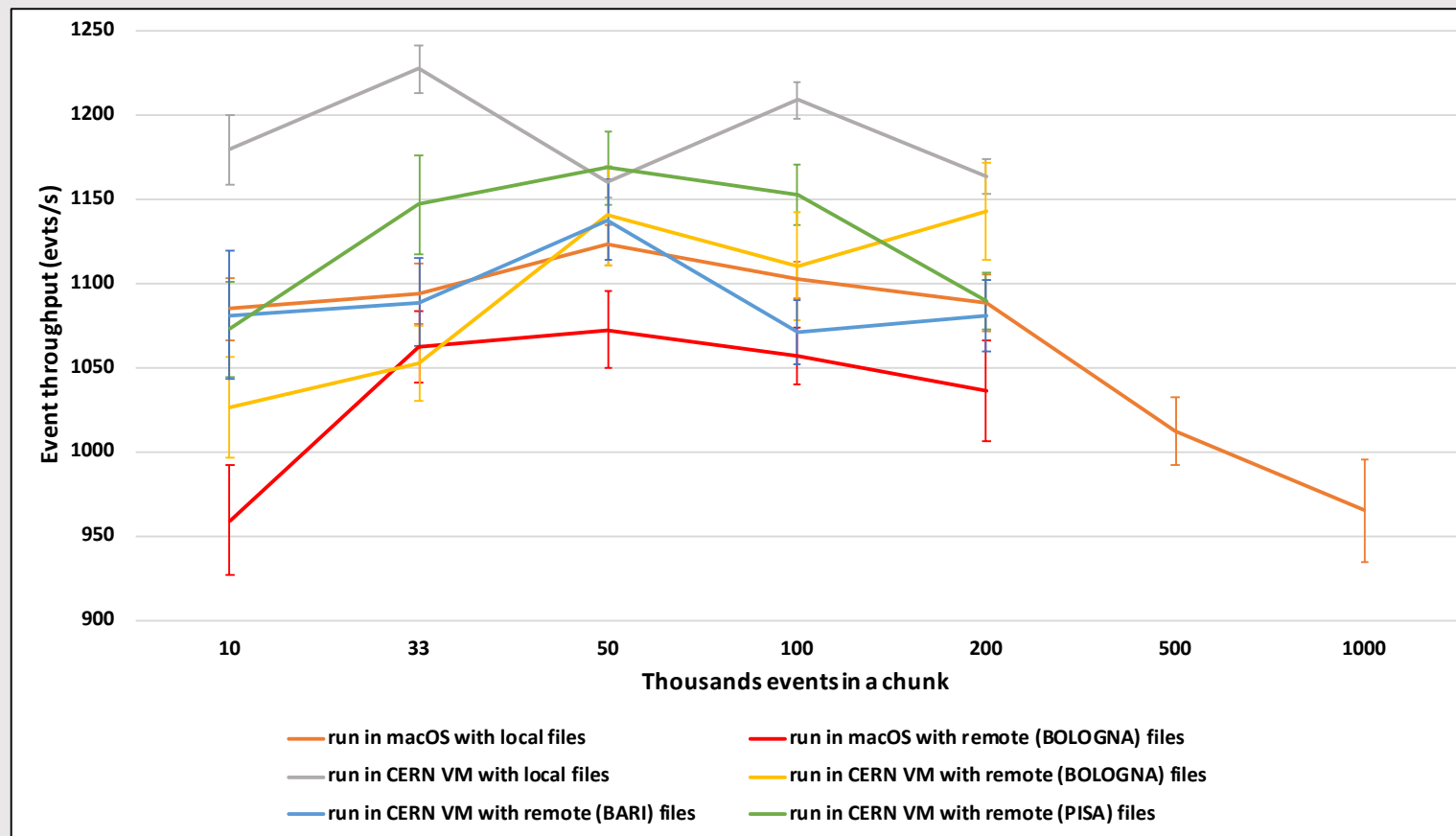
~ 19-41 min

~ 16-17 min

~ 35-57 min

~ 8.4k – 13.7k evts/s

Values for chunk size fixed to 100 thousands events



	event throughput for creating a chunk (evts/s)	event throughput for pre-processing a chunk (evts/s)
macOS with local files	1102 (11)	1157 (7)
macOS with remote files (BO)	1057 (17)	1138 (4)
VM with local files	1209 (11)	1247 (2)
VM with remote files (BO)	1110 (32)	1243 (5)
VM with remote files (BA)	1071 (19)	1153 (4)
VM with remote files (PI)	1152 (18)	1234 (5)

~ 1.1k – 1.2k evts/s

Values for chunk size fixed to 100 thousands events

# More on MLaaS4HEP performance

---

- The training step takes about 7 hours for both MacOS and CERN VM, using all the 28.5M events (plus the time required to train the ML model), where:
  - 36% is used to extract and convert each event in a list of NumPy arrays;
  - 27% is used for fixing the Jagged Arrays' dimension;
  - 26% is used for the normalization step;
  - 6% is used for creating the masking vectors.
- We estimate that, projecting these results for datasets at the **TB scale** and using the same hardware resources, the specs computing phase will take  $O(100)$  hours and the training step will take  $O(1k)$  hours (plus the time required to train the ML model)
  - Further optimization of the MLaaS4HEP pipeline will be required to process TB or PB scale datasets and it may involve parallelization of I/O, distributed ML training, etc.
- The time to train the ML model is not included in the performance shown. This time is independent from the MLaaS4HEP framework since it is determined by usage of the underlying ML framework, the complexity of used ML model and the available hardware resources.
  - In our case the training time for a chunk of 100k events is about 11 seconds and 13 for MacOS and CERN VM, respectively.

# How to use TFaaS (1)

- Follow the installation instructions [here](#)
- Setup url to point to your TFaaS server (e.g url=http://localhost:8083 or url=<https://cms-tfaas.cern.ch/>)
- If necessary convert the model into a TensorFlow one (e.g using [this](#) solution)

- create upload json file, which should include:

- fully qualified model file name
- fully qualified labels file name
- model name you want to assign to your model file
- fully qualified model parameters json file name

**Example of an upload.json file:**

```
{ "model": "/path/model.pb",  
  "labels": "/path/labels.txt",  
  "name": "model_name",  
  "params": "/path/params.json" }
```

**Example of a params.json file:**

```
{ "name": "model_name", "model":  
  "model.pb", "description": "my model  
  description", "labels": "labels.txt",  
  "inputNode": "dense_1_input",  
  "outputNode": "output_node0" }
```

- the model parameters json file is used on a TFaaS server side. It context should be the following:
  - model name (how your model will be named in TFaaS server, e.g. the one in the upload.json file)
  - model file name (name of your model file will be used in TFaaS server, e.g. the one in the upload.json file)
  - model labels file name (similar to model file name but used for labels file)
  - description string (provide details about your model)
  - inputNode of your TF model (can be found by inspecting ptxt)
  - outputNode of your TF model (can be found by inspecting ptxt)

# How to use TFaaS (2)

---

- We provide pure python [client](#) to perform all necessary actions against TFaaS server. Here is short description of available APIs
- upload a model to the server
  - `tfaas_client.py --url=$url --upload=upload.json`
- list existing models in TFaaS server
  - `tfaas_client.py --url=$url --models`
- delete given model in TFaaS server
  - `tfaas_client.py --url=$url --delete=model_name`
- prepare input json file for querying model predictions
  - e.g. `{"keys":["attribute1", "attribute2"], values: [1.0, -2.0]}`
- get predictions from TFaaS server
  - `tfaas_client.py --url=$url --predict=input.json`
- get image predictions from TFaaS server
  - `tfaas_client.py --url=$url --image=/path/file.png --model=ImageModel`



- INFN is offering to its users a comprehensive and integrated set of Cloud services through its dedicated INFN Cloud infrastructure
- The current operational state of INFN Cloud is pre-production but already serving several INFN experiments and collaborations, with full production state and general availability expected in 2021
  - The access to the INFN Cloud services is currently reserved to INFN personnel or personnel with whom INFN has established formal collaborations, such as research associates
- The INFN Cloud portfolio, available by an easy to use web interface but also exploitable via command line interfaces, is defined upon clear user requirements
- The INFN Cloud services are based on modular components and span the IaaS, PaaS and SaaS
  - All services are described by TOSCA templates (which can refer internally to other components such as Ansible playbooks)
- For more information visit the official [web page](#)