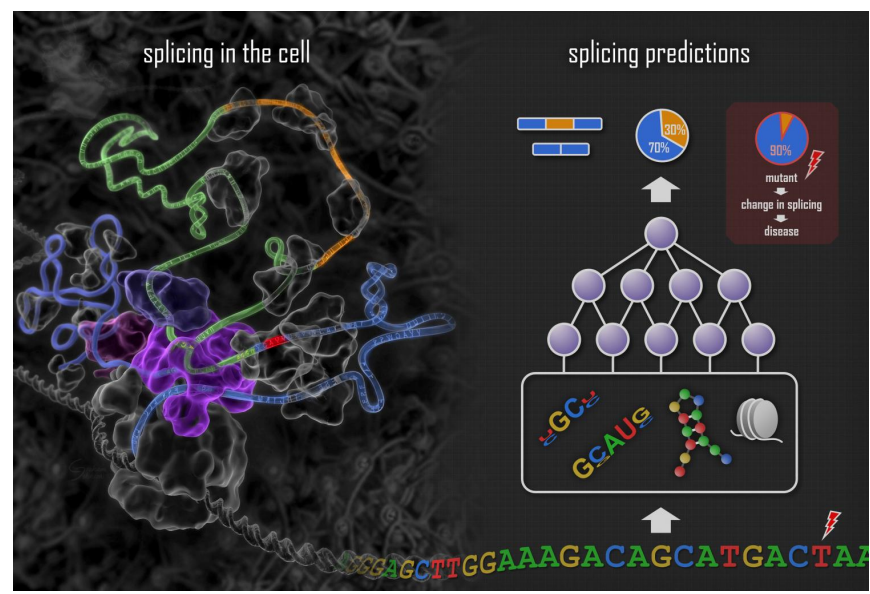


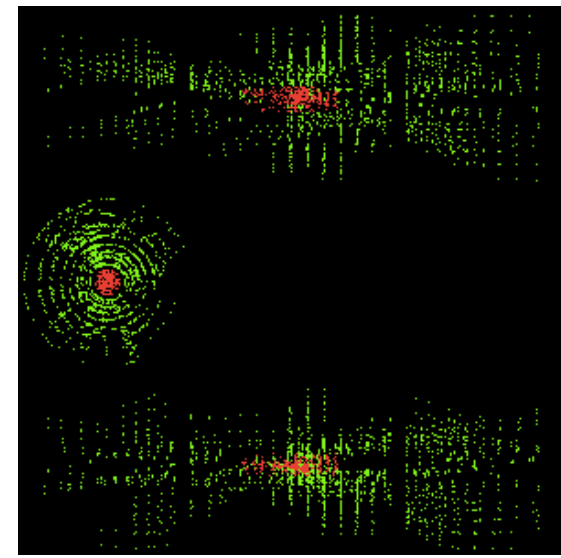
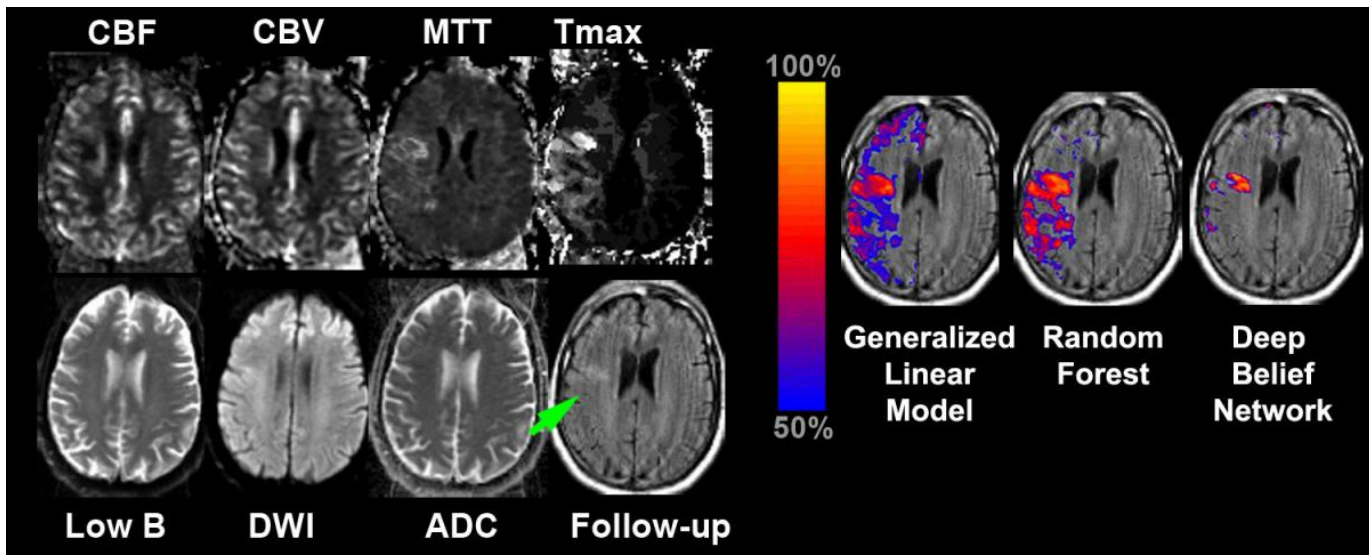
## Machine Learning

not-HEP

HEP  
towards High Luminosity LHC [1]



- Trigger decision
- Tracking
- Detector quality monitoring
- Detector anomaly detection
- Simulation
- Computing resource optimisation [2]
- Physics analysis
- Particle Identification and particle properties

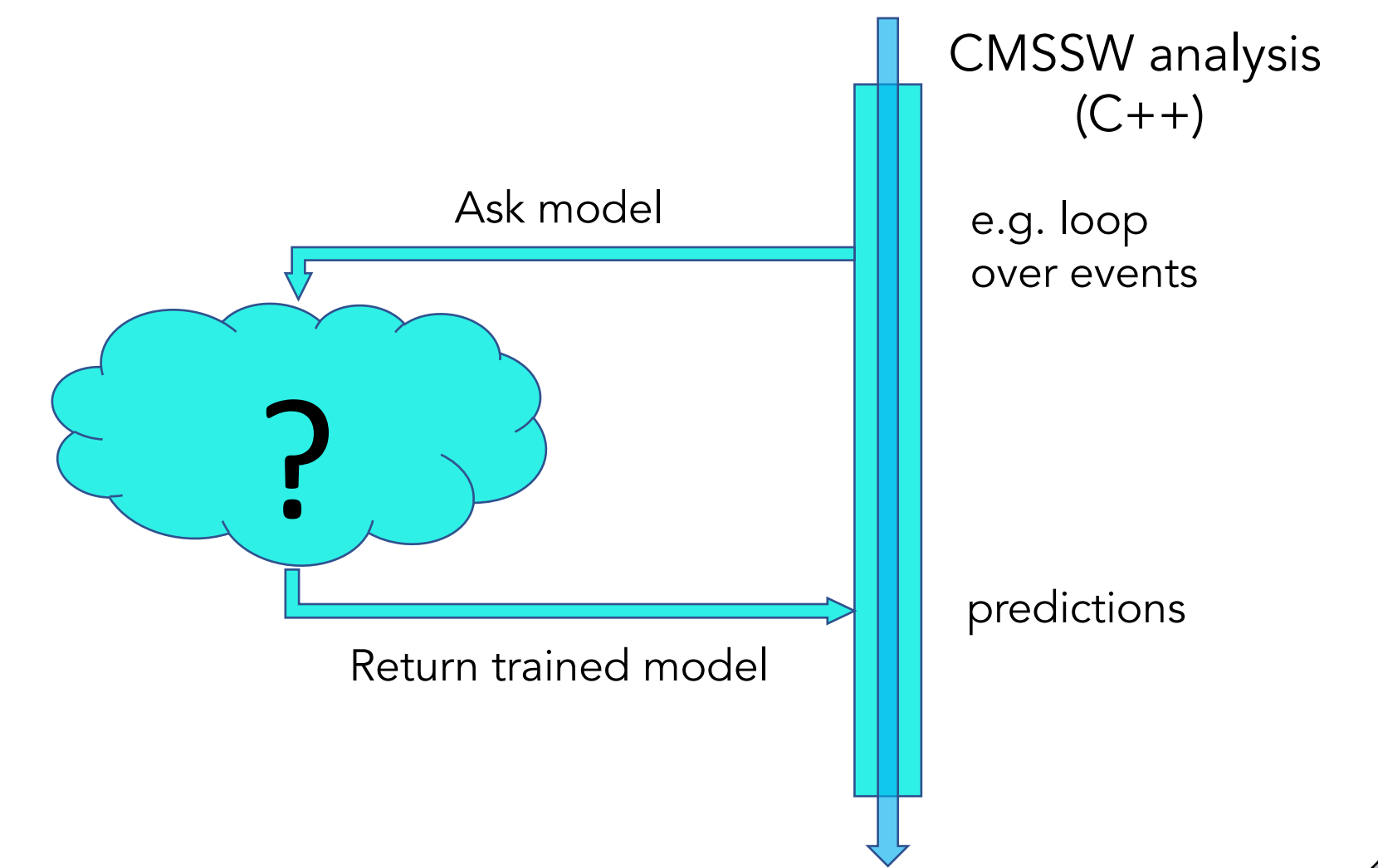


## Machine Learning "as a service" for CMS

The CMS experiment at CERN exploits various ML techniques due various physics and computing related projects. The construction and deployment of a ML project and its deployment for production use requires specific skills and it is a highly time-consuming task. There are no data science teams stably collaborating with CMS physicists and helping them to achieve their ML objectives. At the same time, the CMS physicists themselves rarely have specific data science skills to face such challenges alone. What is needed to design and run successfully a ML project is often not found in a basic CMS physicist expertise (or a HEP physicist, for what matters) whose primary competences are focussed on high energy physics, data analysis (including statistics), and whose ultimate goal is work towards a physics publication. Facing the need to improve a physics data analysis and understanding that ML might be an interesting exploration is hence just a first step towards actually embracing ML in an analysis.

### Goal

The work presented in this poster contributes to build a ML "as a service" solution for CMS Physics needs, namely build an end-to-end data-service to serve ML trained model to the CMSSW framework. The basic idea is as simple as this: instead of asking each physicist who wants to exploit ML in their own task to just learn how to do it and do it themselves independently, each user would ultimately build a modified data analysis code where "calls" to an external service - as simple as calls to functions - would be added to return a trained ML model output that could be directly used in the analysis code (e.g. in loops over events) in a streamlined manner, thus hiding all the complexity related to the ML machinery via outsourcing this to an external service.



## Tensorflow "as a service" (TFaaS)

Choice of the Machine Learning framework

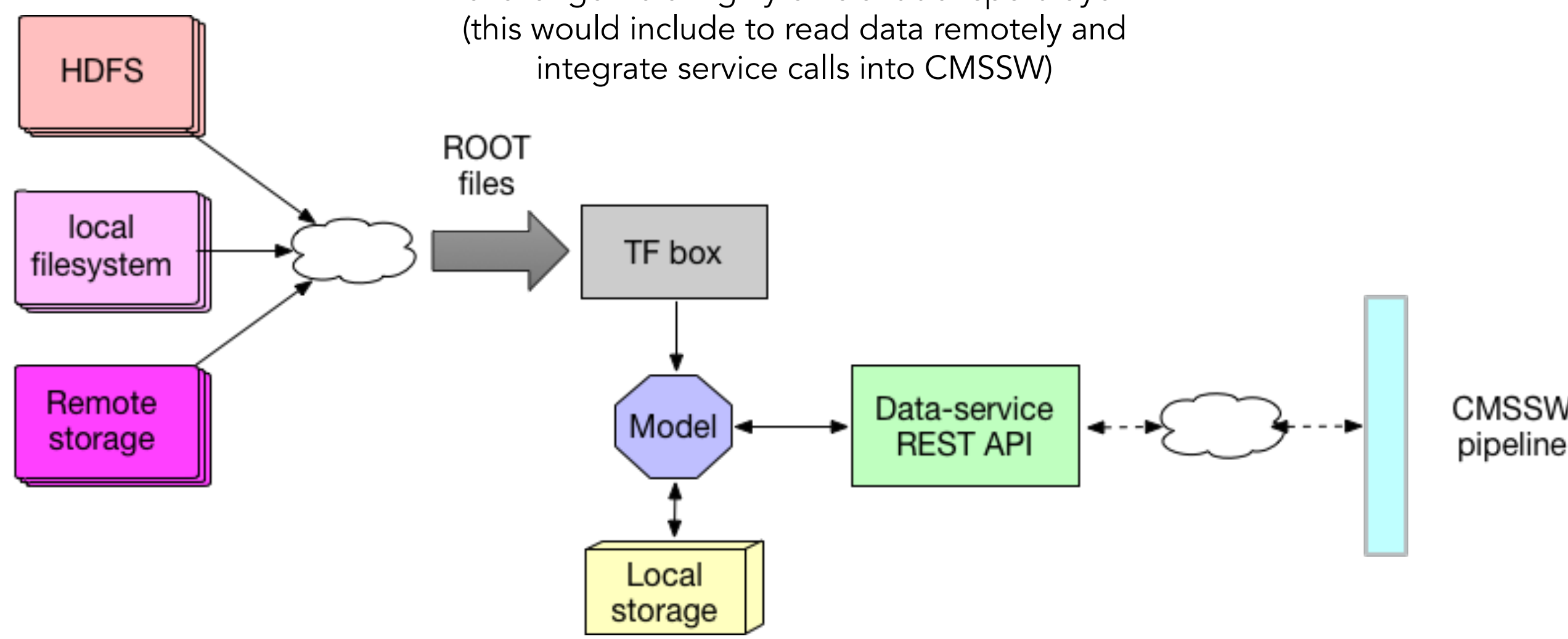


### Main functionalities of TFaaS [3]

read ROOT files and convert them into a suitable form as input to ML/DL systems

train the model on a TF box and serve the trained ML/DL model via REST API with data exchange via a highly efficient transport layer (this would include to read data remotely and integrate service calls into CMSSW)

deploy the service to the cloud (e.g. rent GPUs to train the model), and use it for predictions



## General process

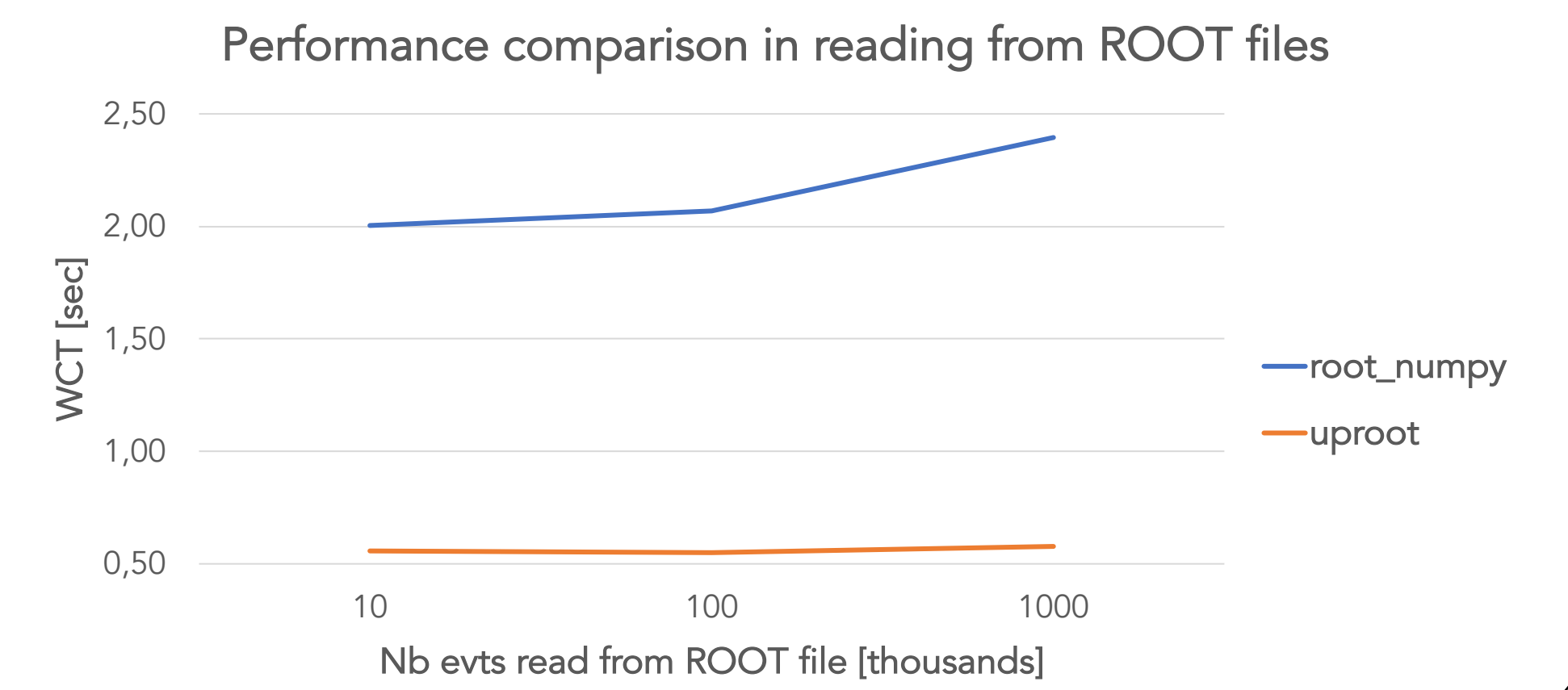
- Select use case
- Create a model

## In this work

- $t\bar{t}$  selection
- Scikitlearn-based
  1. Data Preparation  $\rightarrow$  hardest part to read ROOT file using a new tool ("uproot")
  2. Data Validation
  3. Algorithm selection
  4. Parameter tuning
- Create a keras-tensorflow model

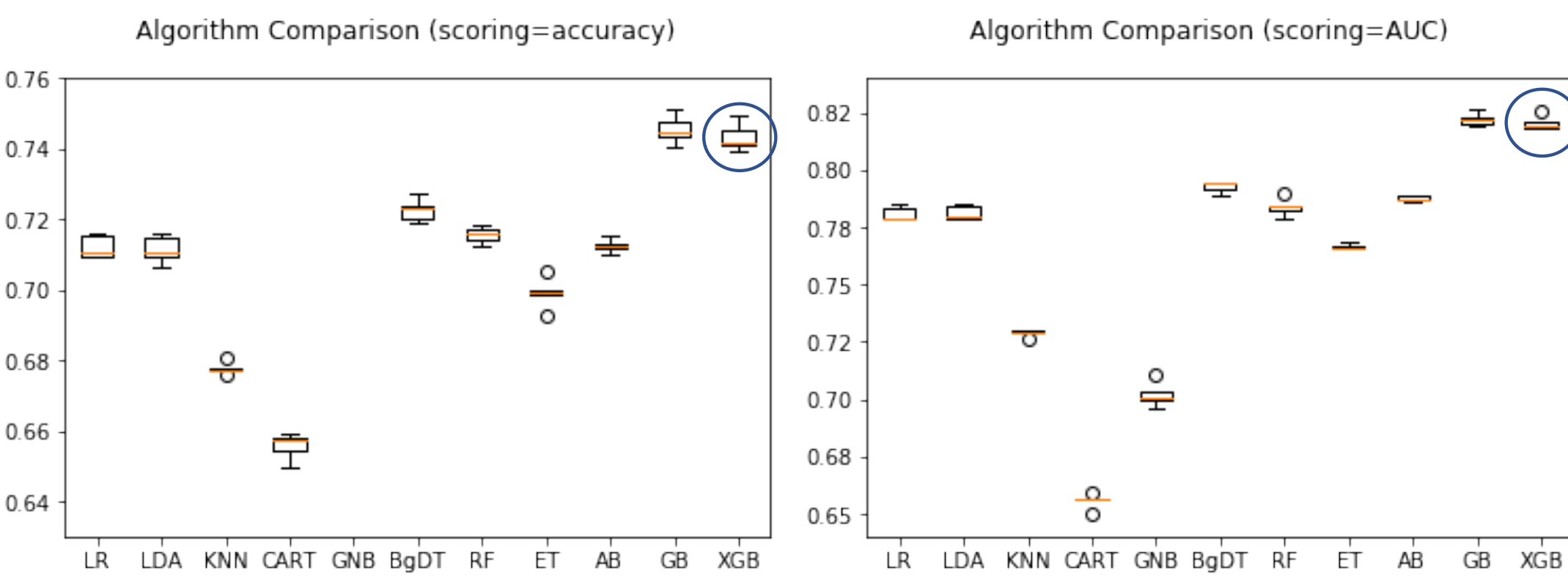
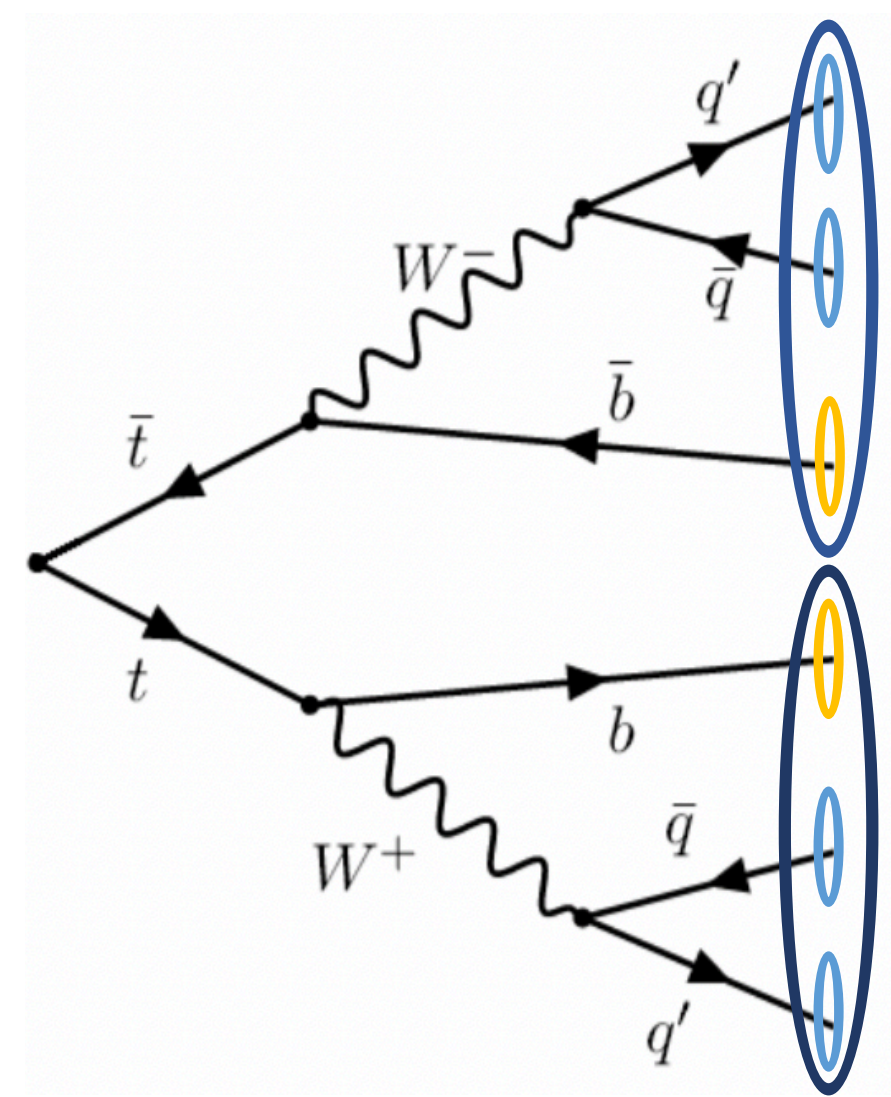
### Interface with the prototype

Disclaimer: this is a direct comparison among two approaches, and not a performance benchmark.



## $t\bar{t}$ use case

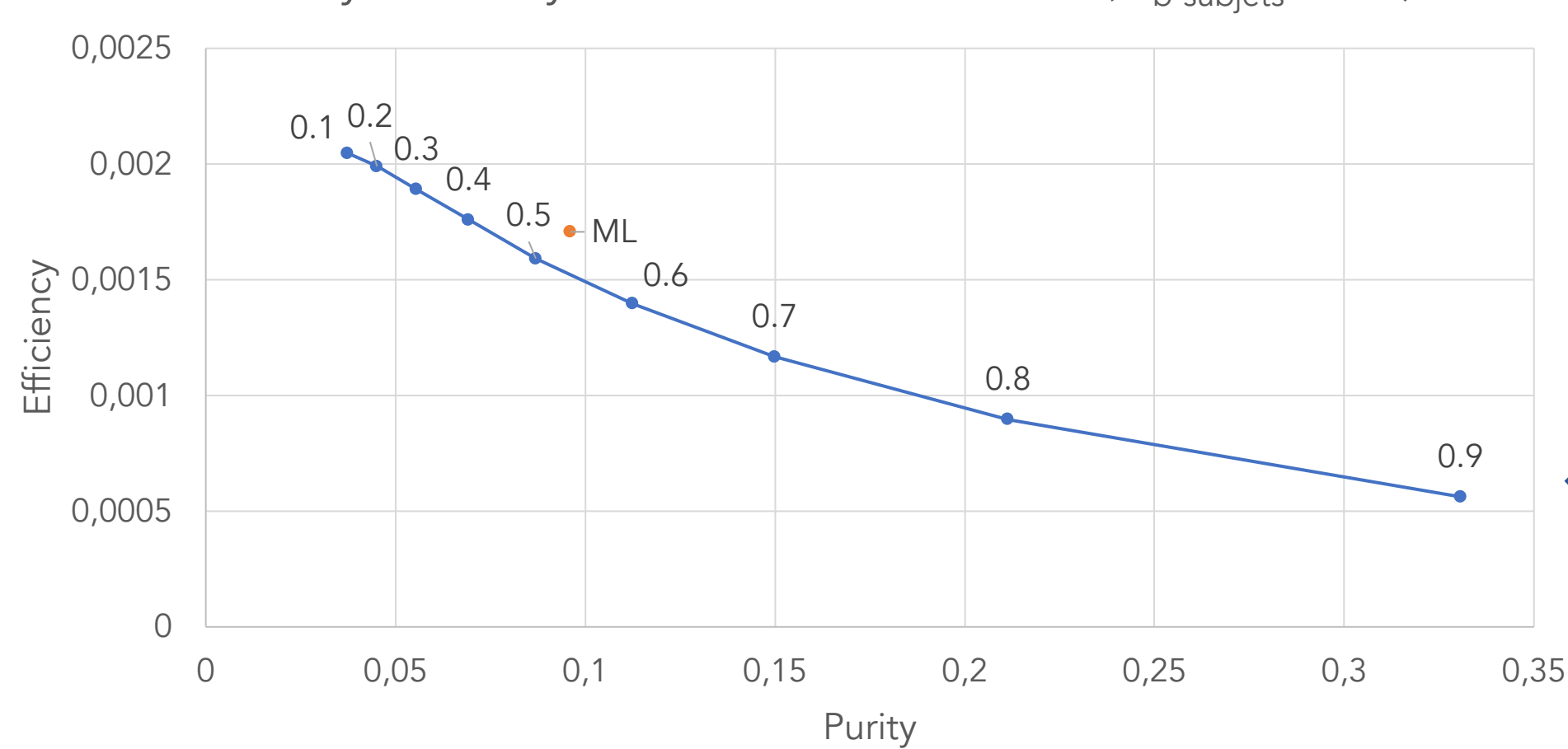
the Signal versus Background discrimination in the selection of  $t\bar{t}$  events in the all-hadronic topology



XGBoost (eXtreme Gradient Boosting) is an advanced implementation of gradient boosting algorithm. It offers many improvements with respect to GBM, among which the most relevant are: i) regularization, that helps to reduce overfitting; ii) parallel processing, hence it delivers performances that are demonstrably better than GBM; iii) high flexibility.

## Comparison between ML and ROOT MVA

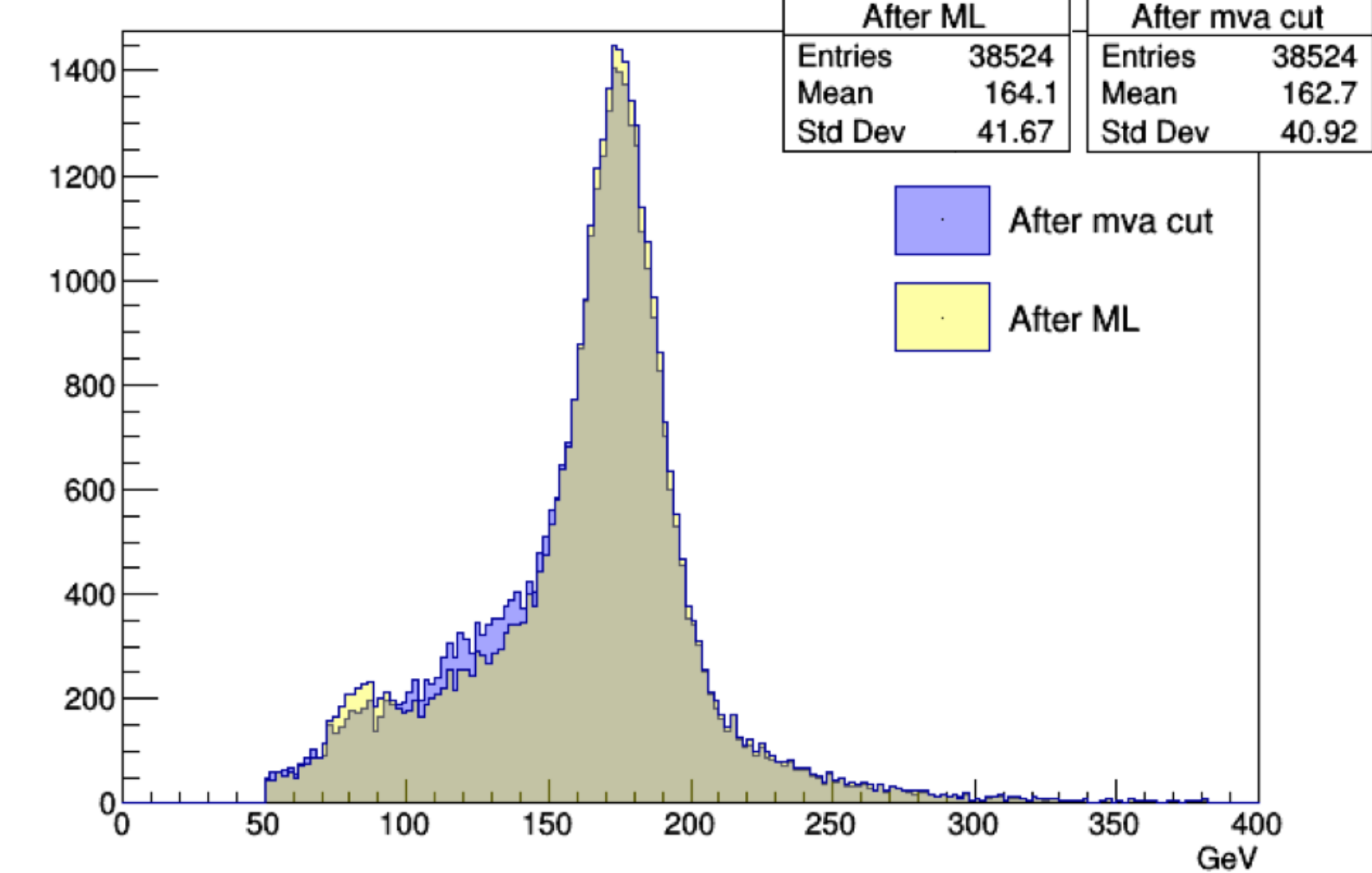
Efficiency VS Purity for different "mva" cuts ( $N_{b-subjects} \geq 1$ )



Efficiency and purity for different mva cuts (the values quoted above the dots), compared to the same obtained by ML, in the case of  $N_{b-subjects} \geq 1$ .

ML is not worst respect to MVA.

"jetMassSoftDrop\_0" from MC sample ( $N_{b-subjects}=2$ )



Invariant mass of the leading jet after the "soft drop declustering" algorithm for the MC data with  $N_{b-subjects} = 2$ .

Below a simple use of TFaaS is shown: after the TFaaS server is launched and the model loaded, the prediction for an event in terms of probability is produced.

```
str957-135:Demo luca.giommi$ cd TFaaS/src/Go/
str957-135:Go luca.giommi$ ./tfaaS -config config.json &
[1] 8518
str957-135:Go luca.giommi$ [INFO][8080] <Config port=8083 dir=/Users/luca.giommi/Demo/keras_to_tensorflow base= auth=false config
Print verbose=0 log=text cnt=key>
[INFO][8080] Starting HTTP server
      Addr=":8083"

str957-135:Go luca.giommi$ curl -L -k -i -X POST http://localhost:8083/upload -F 'name=model1' -F 'params=@Users/luca.giommi/D
emo/TFaaS/src/Go/params1.json' -F 'model=@Users/luca.giommi/Demo/keras_to_tensorflow/model1/11_04_model.h5.pb' -F 'labels=@Users/luc
a.giommi/Demo/keras_to_tensorflow/sb_labels.csv'
HTTP/1.1 100 Continue

[INFO][8160] store as params1.json      fileName=params1.json
[INFO][8160] Uploaded                  File=/Users/luca.giommi/Demo/keras_to_tensorflow/model1/params1.json
[INFO][8160] Uploaded                  File=/Users/luca.giommi/Demo/keras_to_tensorflow/model1/11_04_model.h5
.pb
[INFO][8160] Uploaded                  File=/Users/luca.giommi/Demo/keras_to_tensorflow/model1/sb_labels.csv
HTTP/1.1 200 OK
Date: Fri, 25 May 2018 19:05:28 GMT
Content-Length: 0
Content-Type: text/plain; charset=utf-8

str957-135:Go luca.giommi$ curl -s -L -k --key -H "Content-type: application/json" -d '{"keys": ["nJets", "nLeptons", "jetEta_0
", "jetEta_1", "jetEta_2", "jetEta_3", "jetEta_4", "jetEta_0", "jetEta_1", "jetEta_2", "jetEta_3", "jetEta_4", "jetEtaSoft
Drop_0", "jetEtaSoftDrop_1", "jetEtaSoftDrop_2", "jetEtaSoftDrop_3", "jetEtaSoftDrop_4", "jetPhi_0", "jetPhi_1", "jetPhi_2
", "jetPhi_3", "jetPhi_4", "jetPt_0", "jetPt_1", "jetPt_2", "jetPt_3", "jetPt_4", "jetTau_0", "jetTau_1", "jetTau_2", "jetTa
u_3", "jetTau_4", "jetTau_2_1", "jetTau_2_2", "jetTau_2_3", "jetTau_2_4", "jetTau_3_0", "jetTau_3_1", "jetTau_3_2", "jetTa
u_3_3", "jetTau_3_4"], "values": [2.0, 0.0, 0.9228423833649999, -1.1428750753399999, 0.0, 0.0, 0.0, 155.239425659, 142.7896099
85, 0.0, 0.0, 0.35365982056, 128.549587141, 0.0, 0.0, 1.9385502176299998, -1.17742347717, 0.0, 0.0, 481.41979980
5, 449.04394531199995, 0.0, 0.0, 0.296788358391, 0.286615312089, 0.0, 0.0, 0.164555286895, 0.19625715911408002, 0.0,
0.0, 0.0, 0.11722382675, 0.155224091644, 0.0, 0.0, 0.8]}' https://localhost:8083/json
[INFO][8193] load to cache               Model=model1
2018-05-25 21:06:01.763295: I tensorflow/core/platform/cpu_feature_guard.cc:140] Your CPU supports instructions that this Tens
oFlow binary was not compiled to use: SSE4.2 AVX AVX2 FMA
[INFO][8193] load TF model              Label=@Users/luca.giommi/Demo/keras_to_tensorflow/model1/sb_labels.csv
v Model=/Users/luca.giommi/Demo/keras_to_tensorflow/model1/11_04_model.h5.pb
[0.6046947]
```

## Conclusions

An end-to-end data service has been developed to provide trained ML models to the CMS software framework and in particular the proof-of-concept has been demonstrated in the s/b discrimination in the all-hadronic channel in the  $t\bar{t}$  decay. A simple demo [4] has been created that shows how a common user can use TFaaS to make predictions.

Next steps: review and improve all the steps done and move the model creation and training on cloud.

## References

- [1] The HEP Software Foundation (HSF), A Roadmap for HEP Software and Computing R&D for the 2020s. 2017. arXiv:1712.06982
- [2] V. Kuznetsov, T. Li, L. Giommi, D. Bonacorsi, T.Wildish, Predicting dataset popularity for the CMS Experiment. Journal of Physics: Conference Series 762.1 (2016), arXiv:1602.07226.
- [3] TensorFlow as a Service (TFaaS). URL: <https://github.com/vkuznet/TFaaS>
- [4] URL: <https://drive.google.com/file/d/11pwt9dOJC9EN3miYkHExel6dd4ba0/view>

## Contact

luca.giommi2@studio.unibo.it