

Thesis project and overview of the 1st year of PhD

1st seminar of the PhD school in Data Science and Computation

Luca Giommi (Univ. of Bologna and INFN) E-mail: *luca.giommi3@unibo.it* Supervisors: Daniele Bonacorsi and Claudio Grandi



• Who am I?

I am a PhD student in Data Science and Computation of the XXXIV cycle. I am associated with the National Institute for Nuclear Physics (INFN) and CERN. My research work is within the CMS-Bologna group.

• What is 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.

• What is CMS?

The Compact Muon Solenoid (CMS) is a general-purpose detector at the LHC. It has a broad physics programme ranging from studying the Standard Model (including the Higgs boson) to searching for extra dimensions and particles that could make up dark matter.

What is CNAF?

CNAF is the biggest computing facility of INFN. Located in Bologna, since 2003 CNAF has hosted the Italian Tier-1 data center of the World-wide LHC Computing Grid, providing the resources, support and services needed for data storage and distribution, data processing and analysis, and Monte Carlo production.





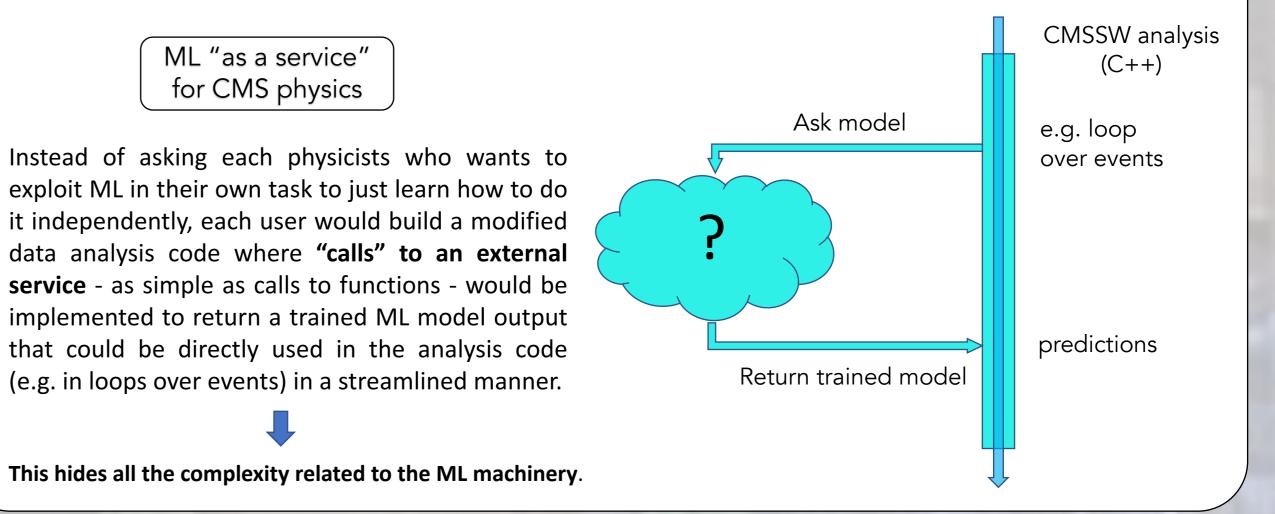


About my thesis project... where we were: master thesis

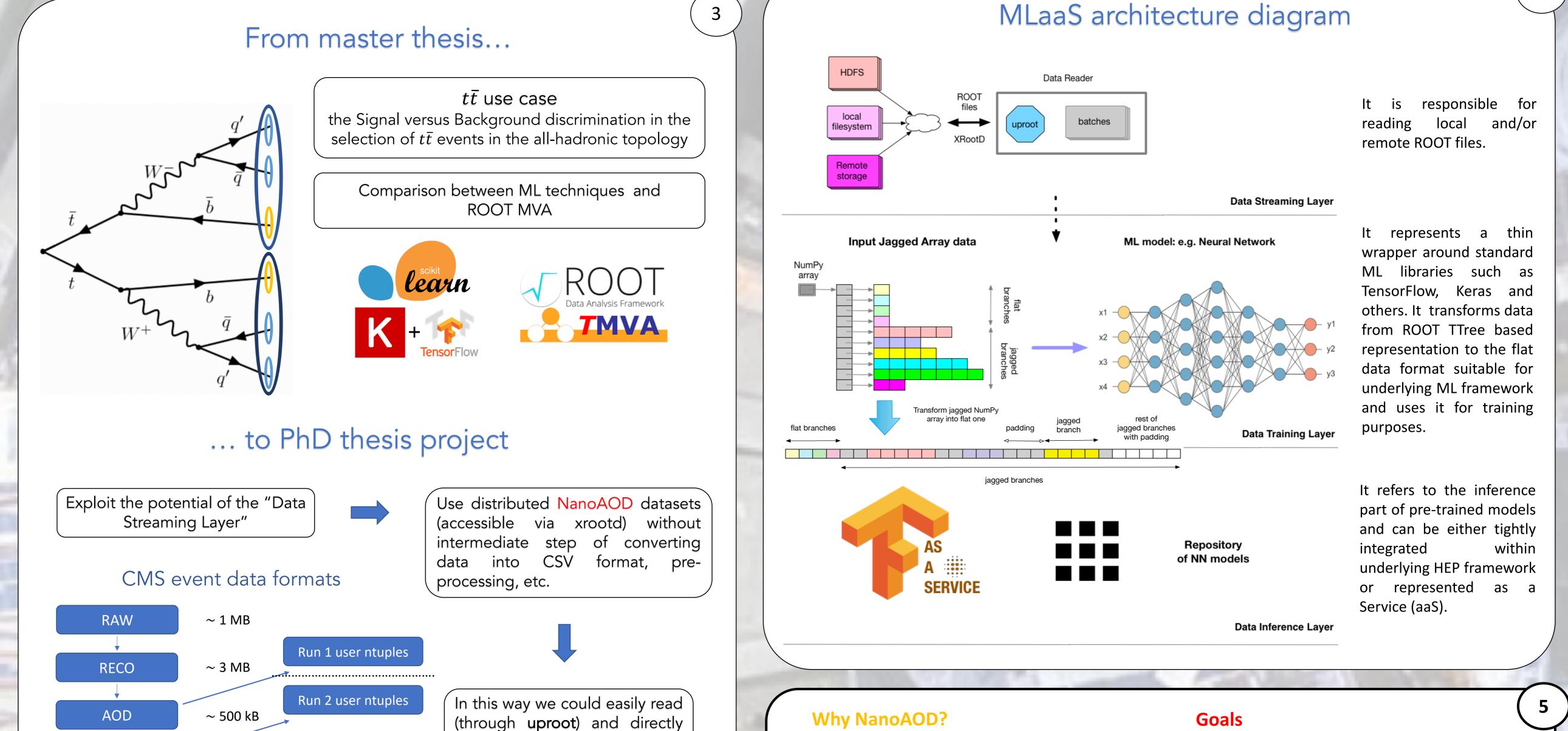
A full proof-of-concept demonstration of an end-to-end data service to provide trained Machine Learning (ML) models to the CMS software framework (CMSSW) and its usage in Signal/Background (S/B) discrimination in $t\bar{t}$ selection

ML "as a service"

for CMS physics



7



6



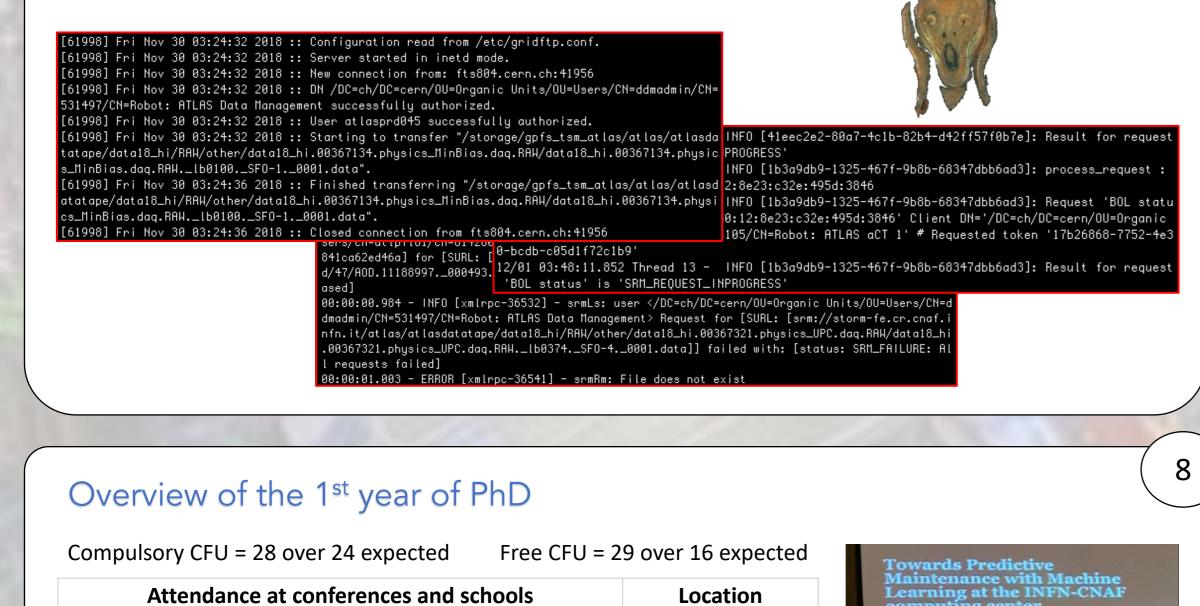
use O(10 TB) of datasets in the training of the ML models

- It is a flat ntuple format with only standard data types (e.g. int, float, vectors), so it is simple to export to modern machine learning frameworks
- Many CMS analysis are switching to the new 1kb/event format
- We don't aim to reproduce an entire analysis with all details, but to demonstrate its feasibility using NanoAOD ROOT files and the MLaaS4HEP framework.
- Performance benchmarks (CPU vs GPU vs TPU, and various versions) for the training phase

Project at CNAF

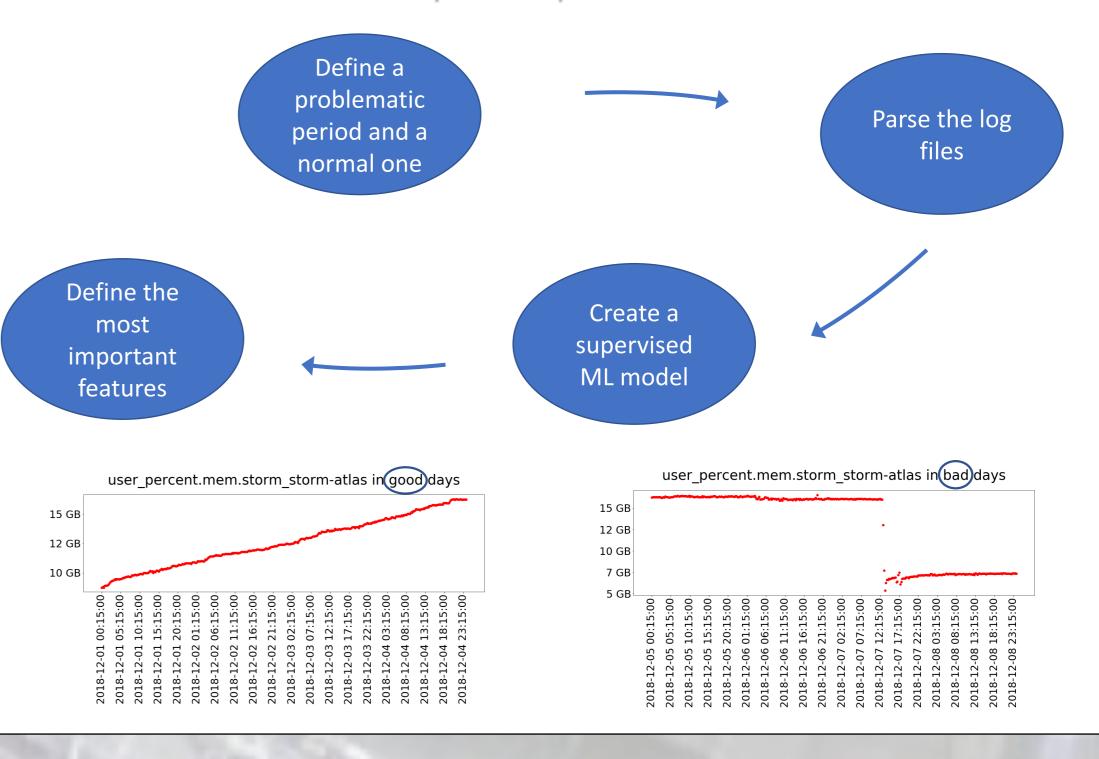
In order to increase efficiency and to remain competitive in the long run, CNAF is launching various activities aiming at implementing a global predictive maintenance solution for the site. Because of efficient storage systems are one of the key ingredients of Tier-1 operations, at CNAF an exploratory work started by investigating logs from the StoRM service.

Information about the status and the progress of the requests managed by the service is stored in log files, in a usually complex format

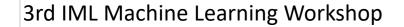


Trento (Italy)

Work presented at the ISGC International Conference (Taipei, April 2019)



References





LHCb/CERN & Microsoft Azure OpenHack

Third international School on Open Science Cloud (SOSC) Bologna (Italy)

International Symposium on Grids & Clouds (ISGC) 2019 Taipei (Taiwan)



[1] L. Giommi, Prototype of ML "as a Service" for CMS Physics in Signal vs Background discrimination, Master thesis and PoS LHCP2018 (2018) 093

[3] V. Kuznetsov, *Machine Learning as a Service for HEP*, arXiv:1811.04492 [hep-ex]

[4] MLaaS github page, https://github.com/vkuznet/MLaaS4HEP

[5] L. Giommi et al, Towards Predictive Maintenance with ML at the INFN-CNAF computing centre. Submitted to the proceedings of ISGC 2019, Taipei