# Machine Learning in ATLAS: activities and future challenges

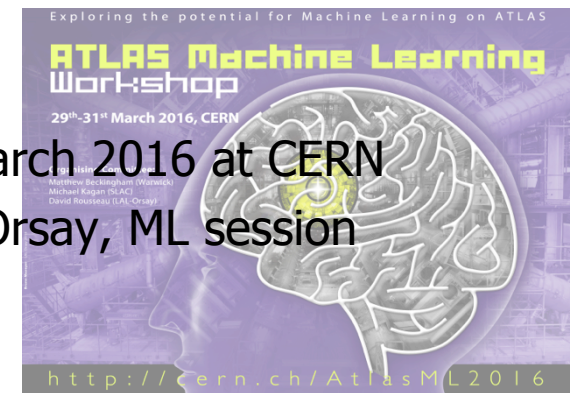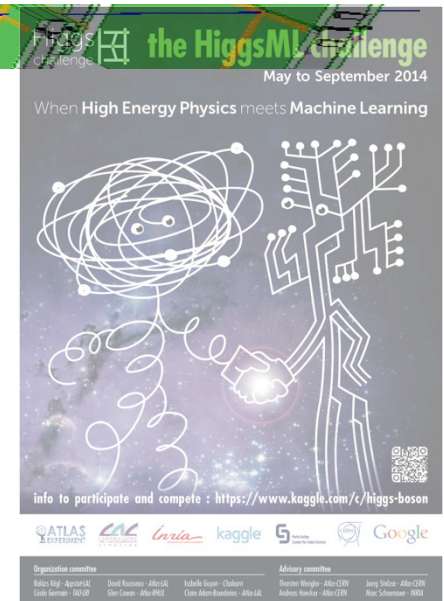**Matthew Beckingham, Michael Kagan, <u>David Rousseau</u>**

**for the ATLAS collaboration**

**OpenLab ML and Analytics workshop, 29th April 2016**

# ML events (with ATLAS participation)

- HiggsML Challenge, summer 2014
  - o ➜HEP ML NIPS satellite workshop, December 2014
- Connecting The Dots, Berkeley, January 2015
- DS@LHC workshop, 9-13 November 2015
  - o ➜future DS@HEP workshop
- LHC Interexperiment Machine Learning group
  - o Started informally September 2015, gaining speed
- Moscou/Dubna ML workshop 7-9th Dec 2015
- Heavy Flavour Data Mining workshop, 18-21 Feb 2016
- Connecting The Dots, Vienna, 22-24 February 2016
- (internal) ATLAS Machine Learning workshop 29-31 March 2016 at CERN
- Hep Software Foundation workshop 2-4 May 2016 at Orsay, ML session
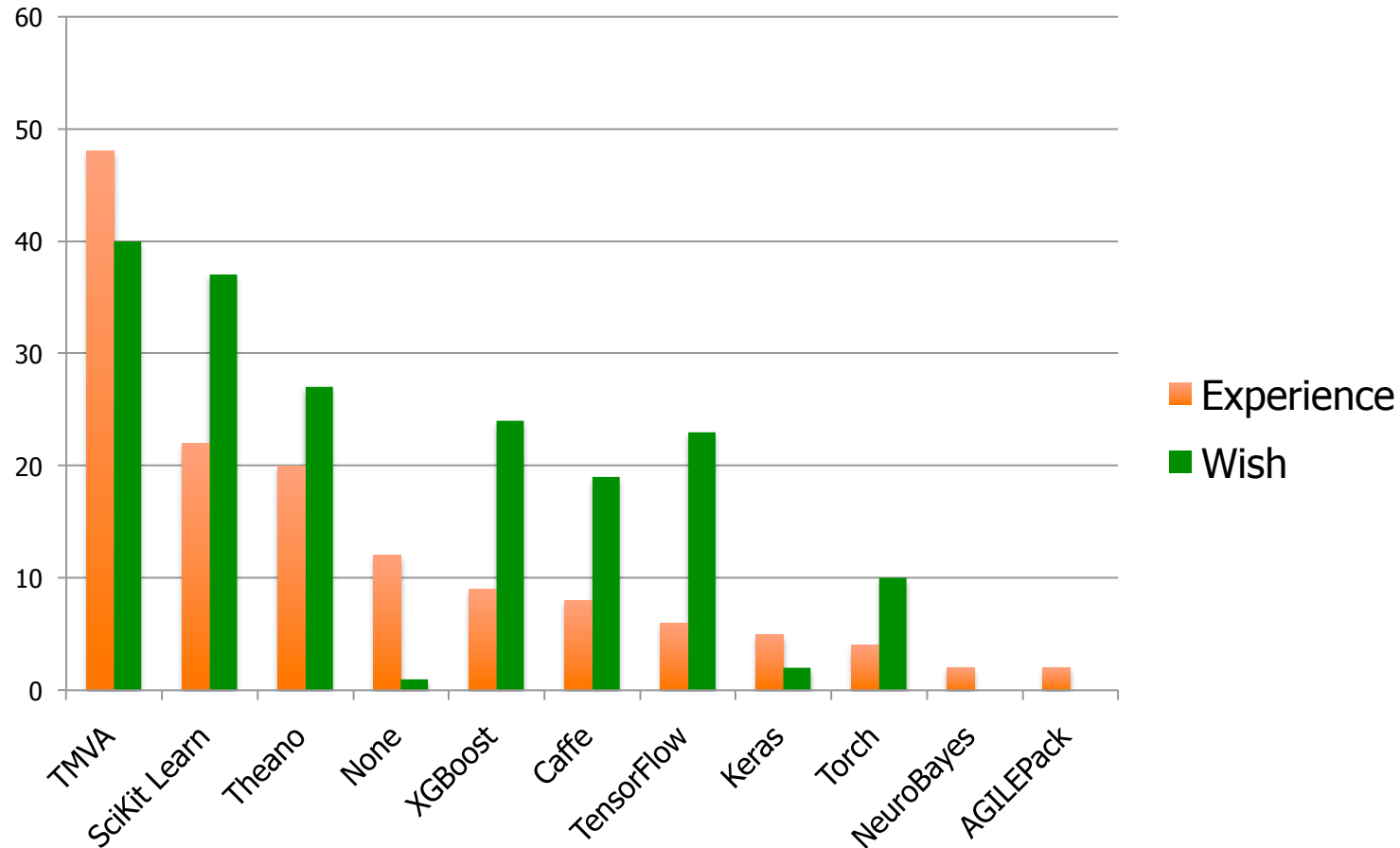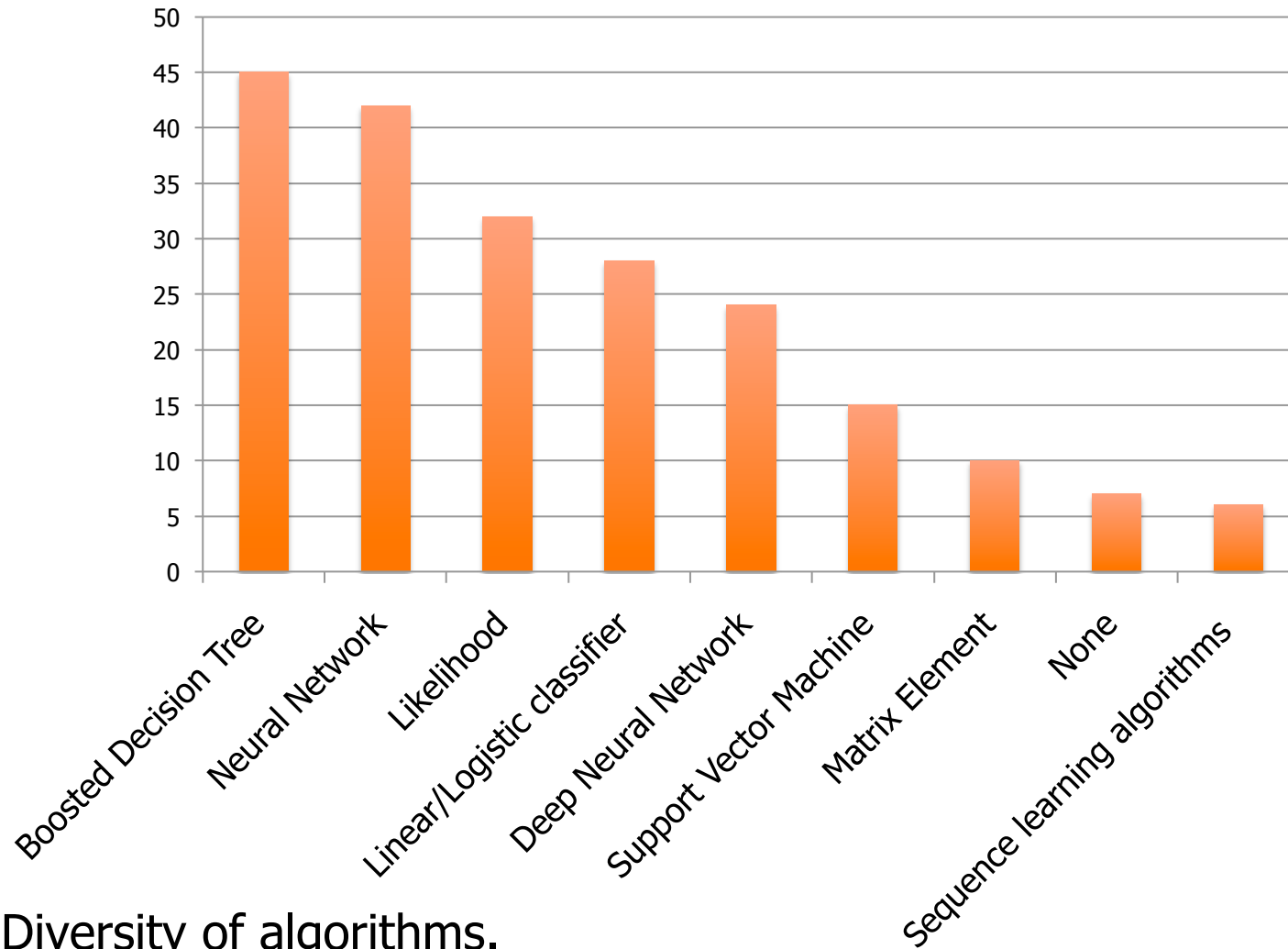- TrackML Challenge, fall 2016

# ML in Atlas

- ❑ Machine Learning (or rather Multi Variate Analysis as we used to call it) used almost since first data taking (2010) for reconstruction and analysis
- ❑ In most cases, Boosted Decision Tree with Root-TMVA, but recent explosion of usage and studies (see later)
- ❑ Recent Atlas ML workshop organised (by MB, MK, DR) to assess current usage, spot opportunities, and favour collaboration with ML experts
  - o 200 participants
  - o Survey of ML usage (see later)
  - o Most of the material shown today gathered there (but limited to published material)
  - o ATLAS ML forum being instantiated, will serve as a forum of discussions within ATLAS and with the outside world

# ATLAS ML Survey 1



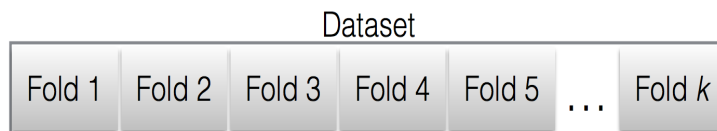☐ Already experience beyond TMVA. Plan to use new tools.

# ATLAS ML Survey 2



❑ Diversity of algorithms.

# Validation Techniques



- ❑ K-fold cross validation allows the estimation of the generalization error
  - ○ Not overly dependent on the exact training / testing split
  - ○ Average / RMS of k-fold errors gives estimate of true error rate
  - ○ Very standard in the non HEP world. Little used in HEP. Being integrated in Root-TMVA, good!
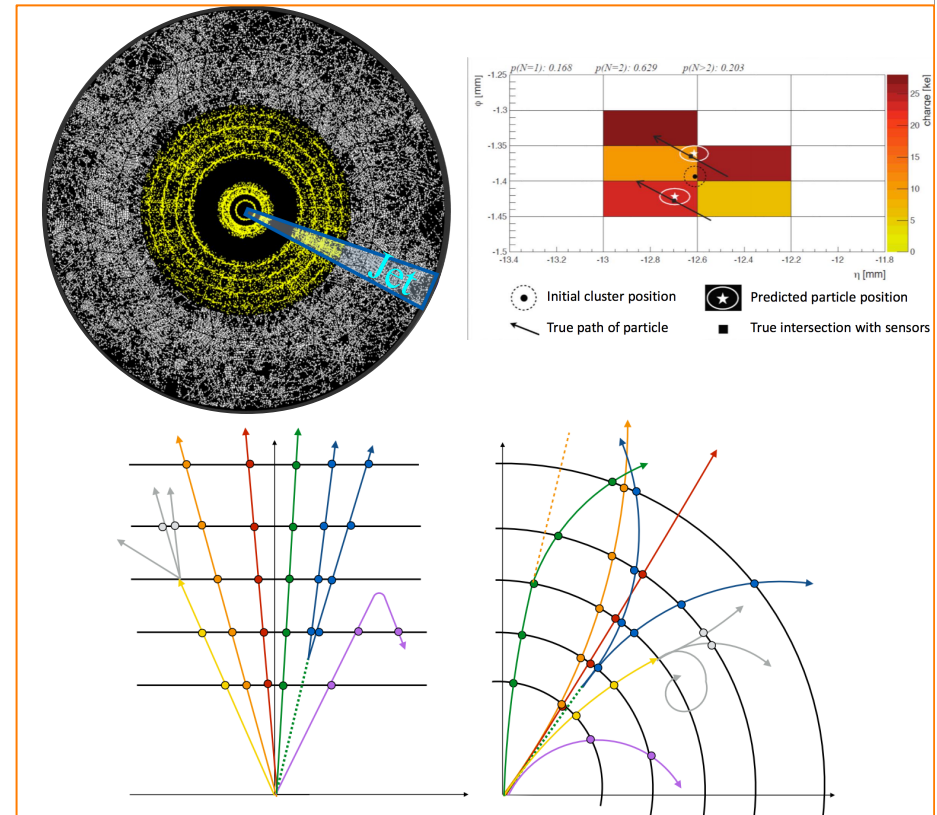


Dataset

| Fold 1 | Fold 2 | Fold 3 | Fold 4 | Fold 5 | … | Fold $k$ |

- ▸ Split the dataset into k randomly sampled independent subsets (folds).
- ▸ Train classifier with k-1 folds and test with remaining fold.
- ▸ Repeat k times.

$$E = \frac{1}{k} \sum_{i=1}^{k} E_i.$$

Openlab ML workshop, ATLAS, 29th April 2016

6

# Reconstruction



❑ Clear upcoming challenges as we approach HL-LHC

❑ Generally, making everything robust to increased pileup, and resource usage will be vital

  o New techniques needed

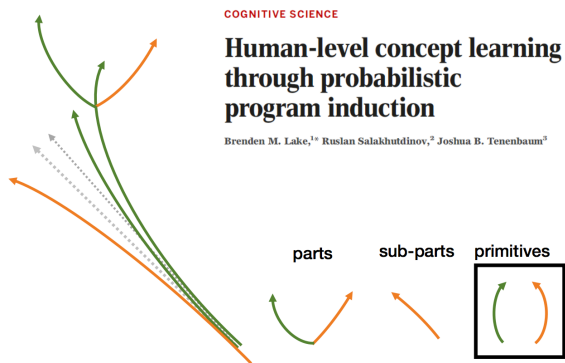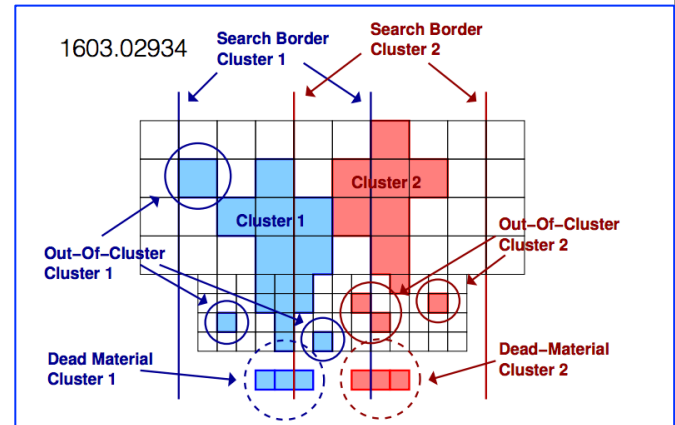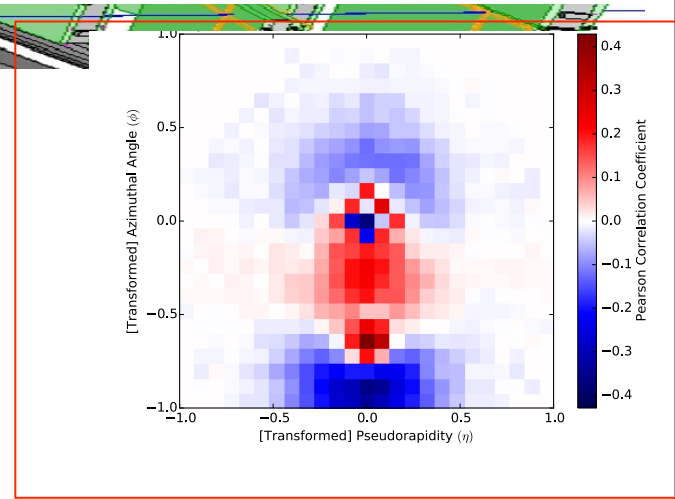  (e.g. TrackML challenge, end of this talk)

# Looking at Data in New ways, using low level info



- Look at data in new ways, potentially a lot of exciting tools available!  Just 2 examples:
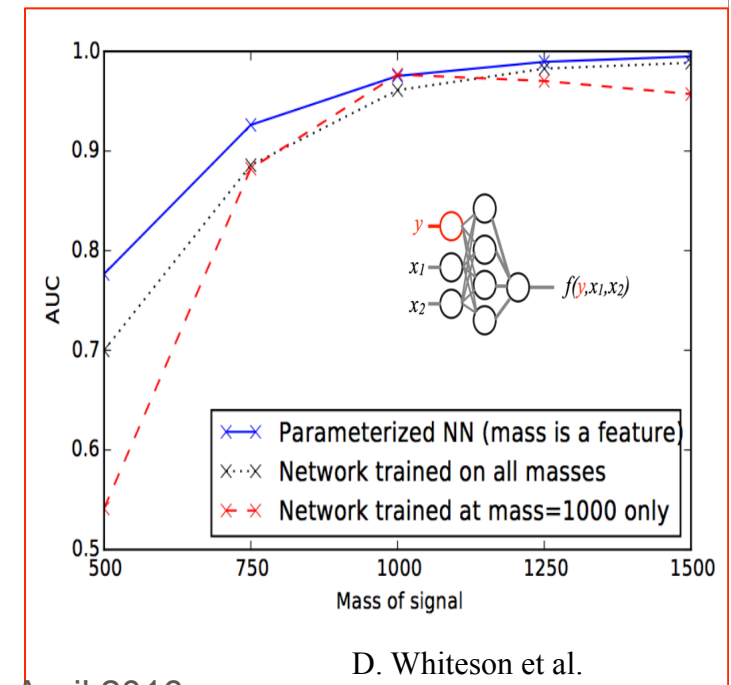  - Computer vision and Imaging
    - Already being studied for jet physics
    - Potential for use for e/gamma cluster ID, topocluster ID / calibration, tau clusters?
      - Similar ideas being investigated for ID in LArTPC's
  - Sequence learning / Machine translation
    - Process variable length sequences with recurrent neural networks
    - Fist studies for processing tracks for b-tagging

COGNITIVE SCIENCE

**Human-level concept learning through probabilistic program induction**

Brenden M. Lake,[1]* Ruslan Salakhutdinov,[2] Joshua B. Tenenbaum[3]

parts    sub-parts    primitives

1603.02934

Search Border Cluster 1
Search Border Cluster 2
Out-Of-Cluster Cluster 1
Cluster 1
Cluster 2
Out-Of-Cluster Cluster 2
Dead Material Cluster 1
Dead-Material Cluster 2

# ML in Analysis
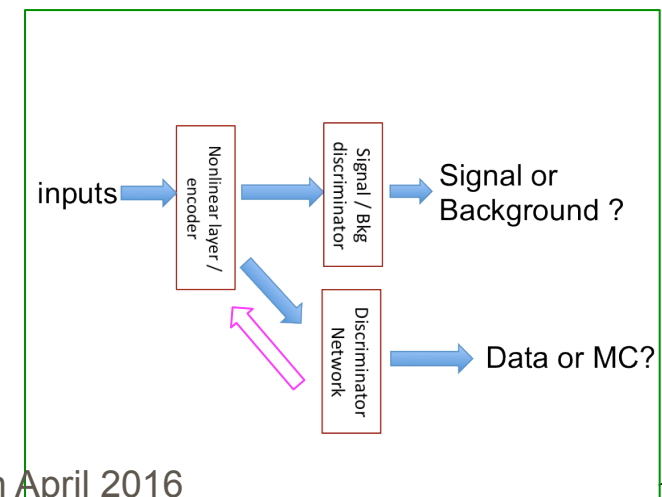
- ML is vital for some analysis
  (e.g. ttH, Single top, $B_s \rightarrow \mu\mu$)
- For most analyses, ML can help push sensitivity
  - We are very good at doing physics already!
- Tend to replace cut based selection by BDT
  - "We run a lot of BDT's" → break down complex analyses into classification problems
- New ideas to ease examination of many signal regions
- ML may help improve statistical analysis techniques
- Take step back: Investigate new ways to look at data,
  - Reconstructing complex final states
  - Interplay between ML and Matrix Element techniques
  - Different approaches to handle multiple classifications
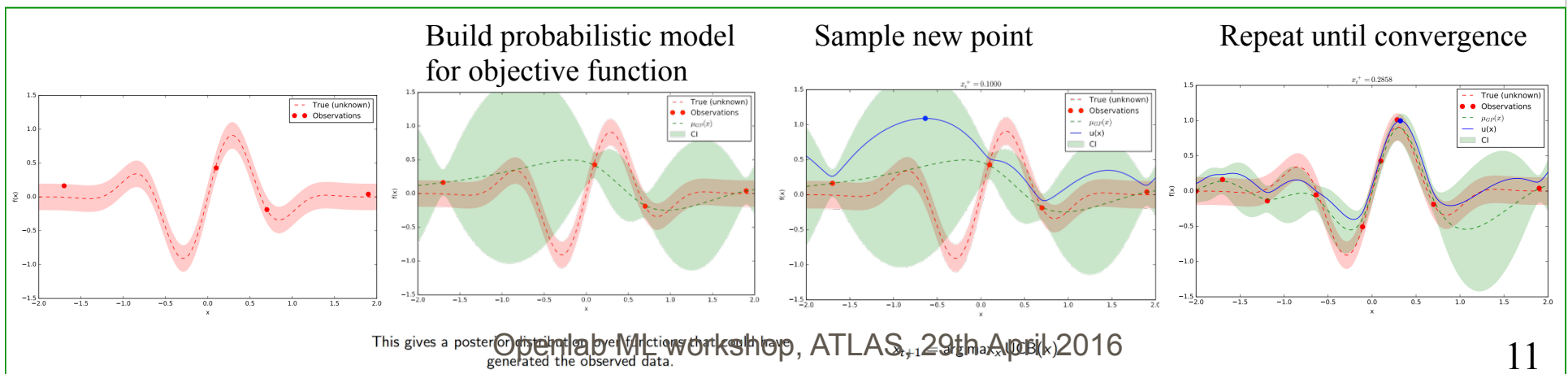  - Improved inputs to MVA can lead to significant gains



D. Whiteson et al.

# ML in Analysis : systematics

Q. Buat et al.

- ❑ Our experimental papers typically ends with
  - o measurement = m ± $\sigma$(stat) ± $\sigma$(syst)
  - o $\sigma$(syst) systematic uncertainty : known unknowns, unknown unknowns...
- ❑ Name of the game is to minimize quadratic sum of : $\sigma$(stat) ± $\sigma$(syst)
- ❑ ML techniques used so far to minimise $\sigma$(stat)
- ❑ Impact of ML on $\sigma$(syst) or even better global optimisation of $\sigma$(stat) ± $\sigma$(syst) is an open problem
- ❑ Worrying about $\sigma$(syst) untypical

of ML in industry

inputs → Nonlinear layer / encoder → Signal / Bkg discriminator → Signal or Background ?

→ Discriminator Network → Data or MC?

# ML in Simulation


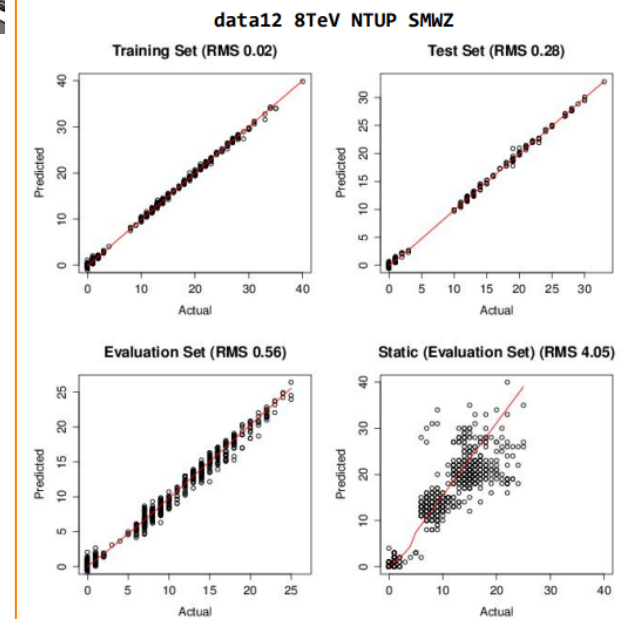
- ❑ We invest a lot of resources (CPU: ~100k cores *year, human) on very fine tuned simulations:
    - o so far very manual optimisation by super experts
    - o optimisation in many dimensions parameter space, with costly evaluation

- ❑ Now turning to more modern techniques e.g.:
    - o Bayesian Optimization and Gaussian Processes



Build probabilistic model for objective function

Sample new point

Repeat until convergence

This gives a posterior distribution over functions that could have generated the observed data.
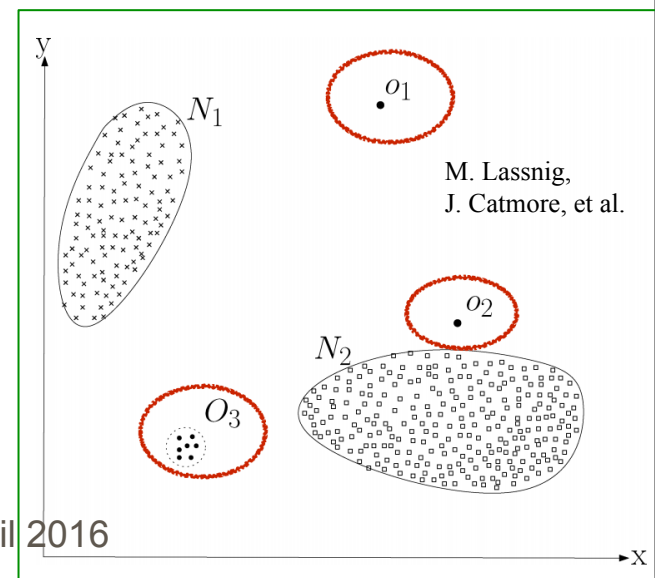
# ML in Software/Computing



- ❑ Workflow and request completion optimization
    - o Under dynamic task constraints

- ❑ Data brokering
    - o Limited CPU and world-wide distributed storage resources for an increasing volume of data
    - o How best to place data around different sites for optimal access? Where best to send jobs?

- ❑ Anomaly detection and prevention
    - o Sometimes, things don't work as expected
    - o Automated preventive measures
    - o Monitor computing infrastructure? Aid Data Quality? MC / Data production validation?

    - o Could also be useful for a generic search for new physics, when we don't know exactly what we are looking for?
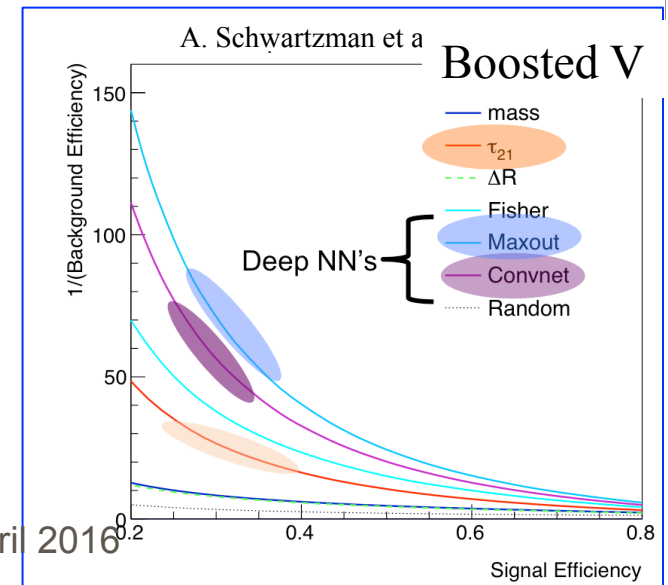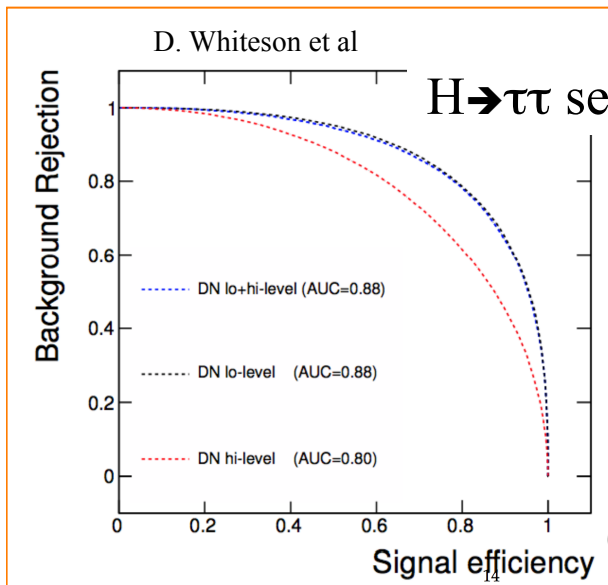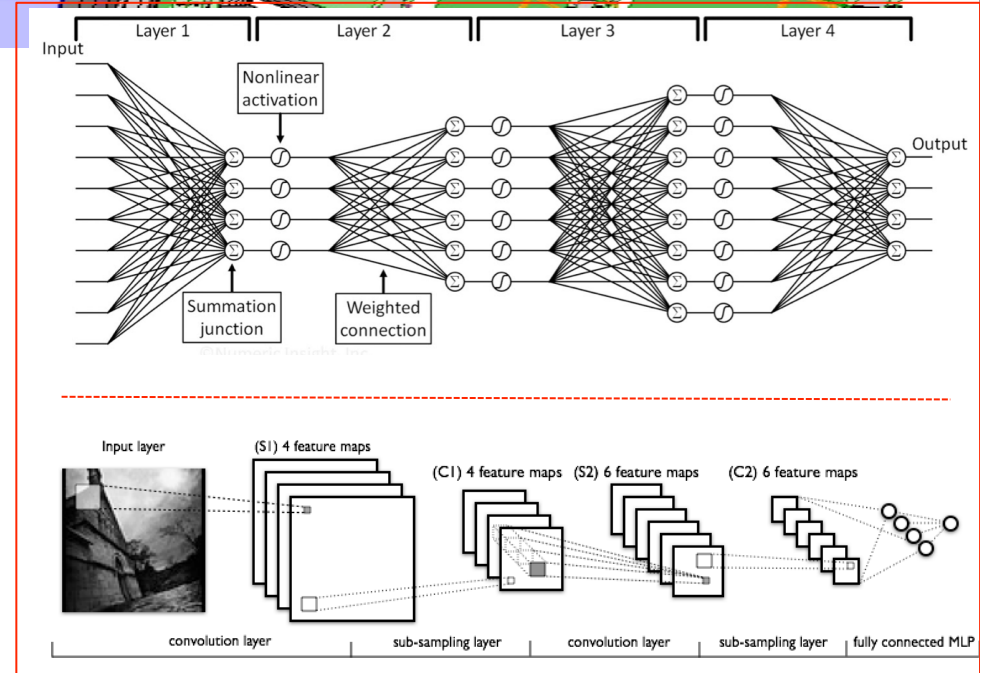


data12 8TeV NTUP SMWZ

M. Lassnig et al.

M. Lassnig, J. Catmore, et al.

# Potential for Deep Learning



- ❑ Variety of studies inside and outside ATLAS show large potential for using DNN's
    - o Improving reconstruction (particle identification)
    - o Improving Searches and Analyses
    - o ...

- ❑ These have been seen to be very powerful algorithms in ML, definitely worth more exploration on ATLAS



D. Whiteson et al

H➔ττ search

- DN lo+hi-level (AUC=0.88)
- DN lo-level (AUC=0.88)
- DN hi-level (AUC=0.80)



A. Schwartzman et al

Boosted V

- mass
- $\tau_{21}$
- $\Delta R$
- Fisher
- Maxout
- Convnet
- Random

Deep NN's

# Tools



- ❑ Root-TMVA heavily used, but it is limited
  - o Easy to use on lxplus and local batch systems, nicely integrated to Root
  - o But does not currently have many modern algorithms / validation techniques
  - o Can lead to large usage of memory
  - o We understand new development efforts underway, very welcome
- ❑ Many people trying new algorithms / techniques are using common data science tools (Scikit-learn, Xgboost, Theano, TensorFlow,…)
  - o Such tools not available on lxplus
  - o Quite some "plumbing" to integrate these tools in typical HEP workflow
  - o Common installations for popular tool would lower barrier to use in ATLAS

# ML platforms

- ❏ Training time can become prohibitive (days), especially Deep Learning, especially with large datasets
- ❏ With hyper-parameter optimisation, cross-validation, number of trainings for a particular application large ~100
- ❏ We're exploring ML platforms :
  - o Dedicated cluster (with GPUs)
  - o Relevant software preinstalled (VM)
  - o Possibility to load large datasets (GB to TB)

# Open Data



- ❑ Public dataset are essential to collaborate (beyond talking over beer/coffee) on new ML techniques with ML experts (or even physicists in other experiments)
    - o can share without experiments NDA
- ❑ Some collaborations built on just generator data (e.g. Pythia) or with simple detector simulation e.g. Delphes
    - o Good for a start, but inaccurate
- ❑ Effort to have better open simulation engine (e.g. ACTS for tracking, see later)
- ❑ Role of CERN Open Data portal:
    - o We (in ATLAS) initially saw its use for outreach purposes
    - o But after all ML collaboration is a kind of scientific outreach
    - o ➔We've uploaded there in 2015 the data from Higgs Machine Learning challenge (essentially 4-vectors from full G4 ATLAS simulation Higgs->tautau analysis)
    - o We consider releasing more datasets dedicated to ML studies
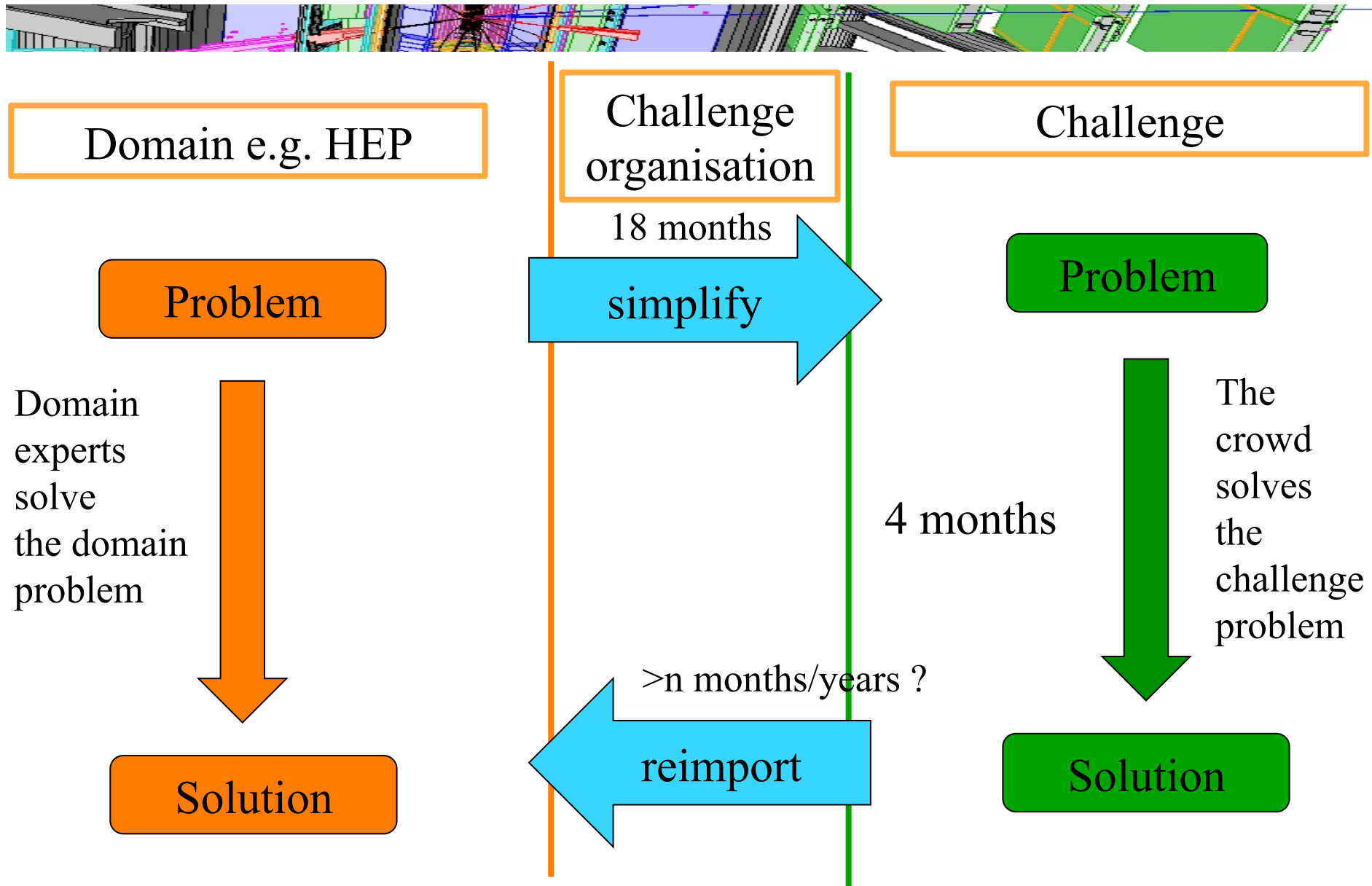
# Challenges (competition)

❑ Challenges are essentially a way to create a buzz around an open dataset dressed with a benchmark

- o HiggsML (ATLAS) 2014
- o FlavourML (LHCb) 2015
- o future TrackML (ATLAS+CMS) 2016?

❑ Buzz in non-HEP world to get the attention of ML specialists

❑ More on this now

# HiggsML in a nutshell



- ❑ Why not put some ATLAS simulated data on the web and ask data scientists to find the best machine learning algorithm to find the Higgs ?
  - o Instead of HEP people browsing machine learning papers, coding or downloading possibly interesting algorithm, trying and seeing whether it can work for our problems
- ❑ Challenge for us : make a full ATLAS Higgs analysis simple for non physicists, but not too simple so that it remains useful
- ❑ Also try to foster long term collaborations between HEP and ML
- ❑ Do not underestimate the time to learn common languages (e.g. hand waving explanation of S/sqrt(B) not enough)
- ❑ Do not underestimate the percolation time :
  - o 1) New ML ideas ➔ 2) Demo on toy data set ➔ 3) Demo in real ATLAS analysis➔4) published ATLAS analysis ==> we're still between 1 and 2 for most new ideas

# From domain to challenge and back

Domain e.g. HEP

Challenge organisation

Challenge

18 months

simplify

Problem

Problem

Domain experts solve the domain problem

4 months

The crowd solves the challenge problem

>n months/years ?

reimport

Solution

Solution

# Higgs Machine learning challenge



- ❑ See talk DR CTD2015 Berkeley
- ❑ An ATLAS Higgs signal vs background classification problem, optimising statistical significance
- ❑ Ran in summer 2014
- ❑ 2000 participants (largest on Kaggle at that time)
- ❑ Outcome
  - o Best significance 20% than with Root-TMVA
  - o BDT algorithm of choice in this case where number variables and number of training events limited (NN very slightly better but much more difficult to tune)
  - o XGBoost best BDT on the market (quite wide spread nowadays)
  - o Wealth of ideas, documented in JMLR proceedings v42
  - o Still working on what works in real life what does not
  - o Raised awareness about ML in HEP
- ❑ Also:
  - o Winner Gabor Melis hired by DeepMind
  - o Tong He, co-developper of XGBoost, winner of special "HEP meets ML" price got a PhD grant and US visa

# Future Tracking Machine Learning challenge

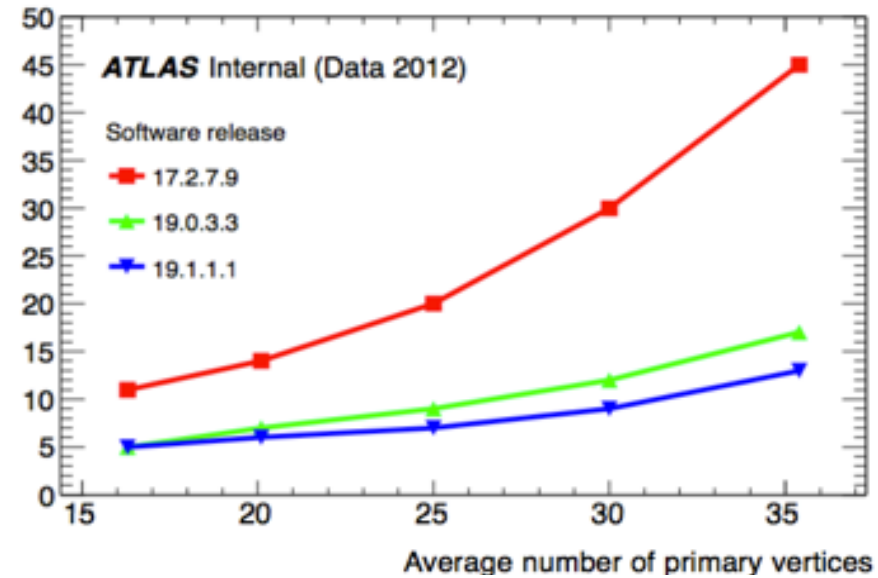**A collaboration between ATLAS and CMS physicists, and Machine Learners**
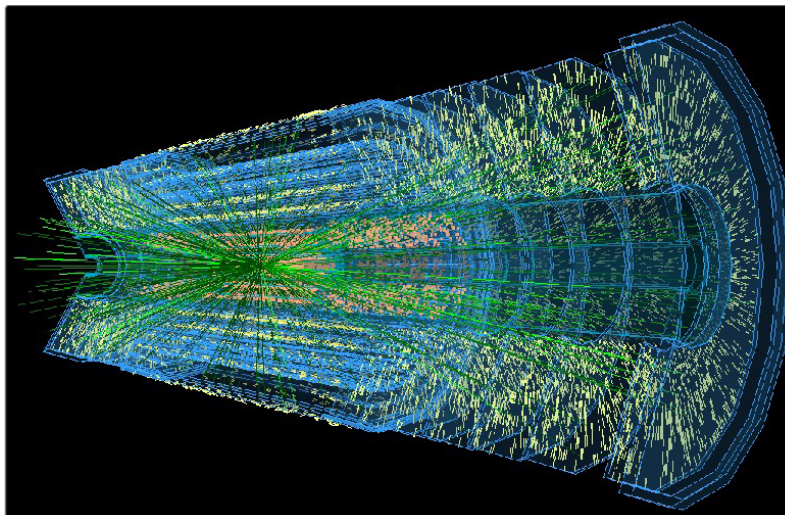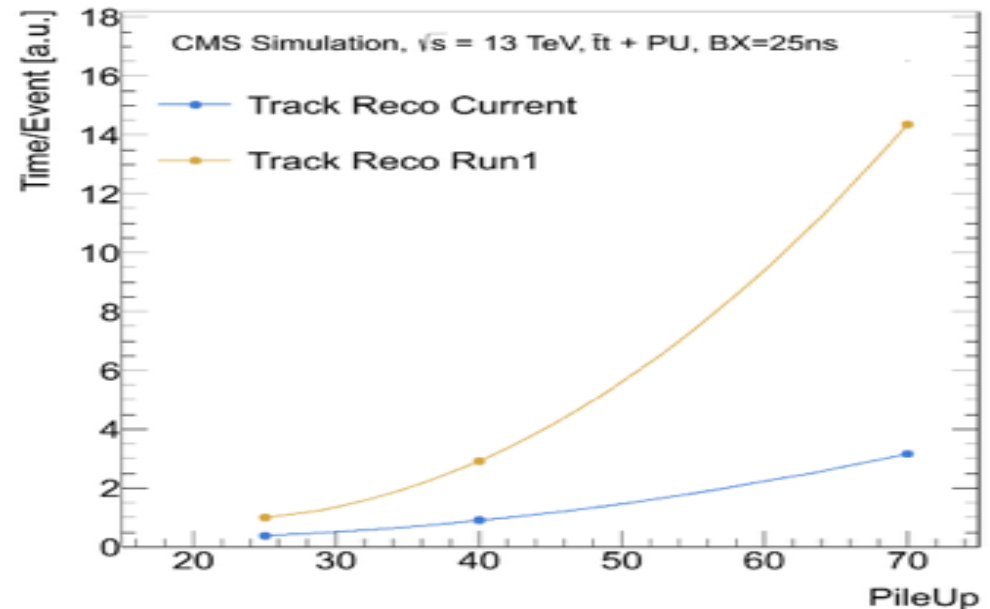
# TrackML : Motivation 1

- ❑ Tracking (in particular pattern recognition) dominates reconstruction CPU time at LHC
- ❑ HL-LHC (phase 2) perspective : increased pileup :
  - o Run 1 (2012): <>~20
  - o Run 2 (2015): <>~30
  - o Phase 2 (2025): <>~150
- ❑ CPU time quadratic/exponential extrapolation (difficult to quote any number)



150



Openlab ML work

# TrackML : Motivation 2



- ❏ LHC experiments future computing budget flat (at best)
- ❏ Installed CPU power per $==€==CHF expected increase factor ~10 in 10 years
- ❏ Experiments plan on increase of data taking rate ~10 as well (~1kHz to 10kHz)
- ❏ ➔HL reconstruction at mu=150 need to be as fast as Run1 reconstruction at mu=20
- ❏ ➔requires very significant software improvement, factor 10-100
- ❏ Large effort within HEP to optimise software and tackle micro and macro parallelism. Sufficient gains for Run 2 but still a long way for HL-LHC.
- ❏ >20 years of LHC tracking development. Everything has been tried?
  - o Maybe yes, but maybe algorithm slower at low lumi but with a better scaling have been dismissed ?
  - o Maybe no, brand new ideas from ML (i.e. Convolutional NN)
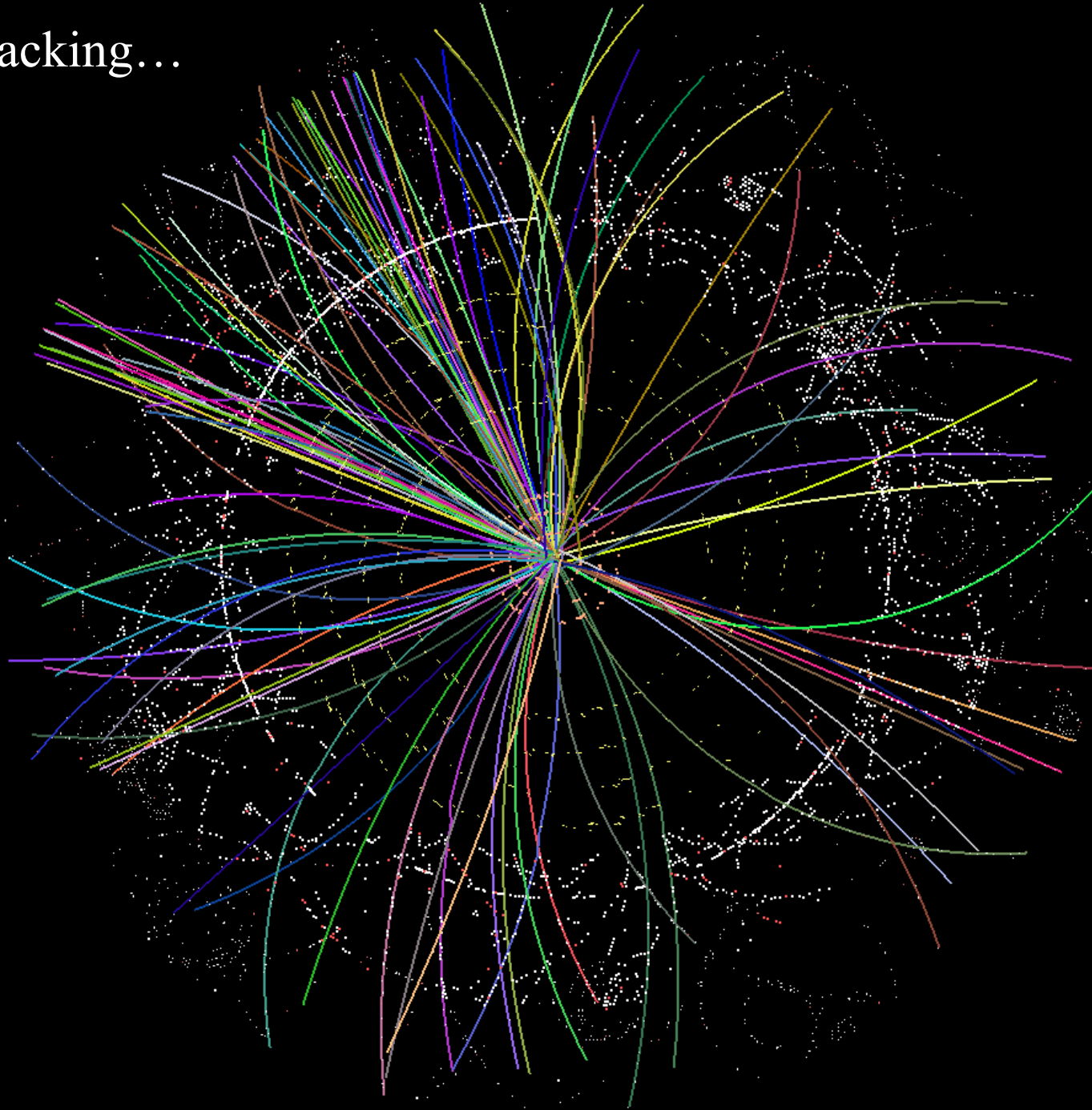- ❏ Need to engage a wide community to tackle this problem
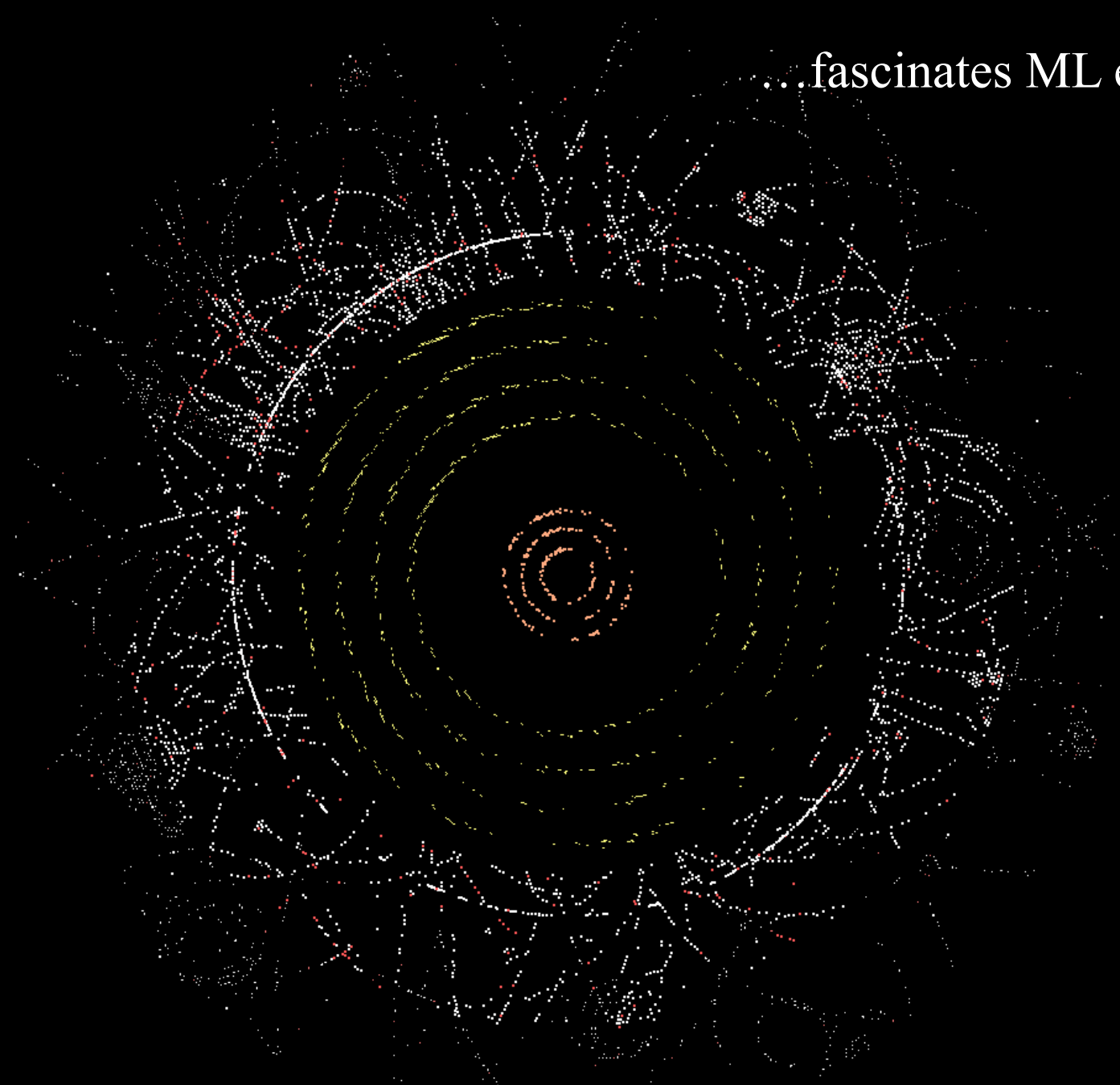
# TrackML : engaging Machine Learners



❏ Suppose we want to improve the tracking of our experiment

❏ We read the literature, go to workshops, hear/read about an interesting technique (e.g. ConvNets, MCTS…). Then:

- o Try to figure by ourself what can work, and start coding➔traditional way
- o Find an expert of the new technique, have regular coffee/beer, get confirmation that the new technique might work, and get implementation tips➔better

❏ …repeat with each technique…

❏ Much much better:

- o Release a data set, with a benchmark, and have the expert do the coding him/herself
- o ➔ he has the software and the know-how so he'll be (much) faster even if he does not know anything about our domain at the beginning
- o ➔engage multiple techniques and experts simultaneously (e.g. 2000 people participated to the Higgs Machine Learning challenge) in a comparable way
- o ➔even better if people can collaborate
- o ➔a challenge is a dataset with a benchmark and a buzz
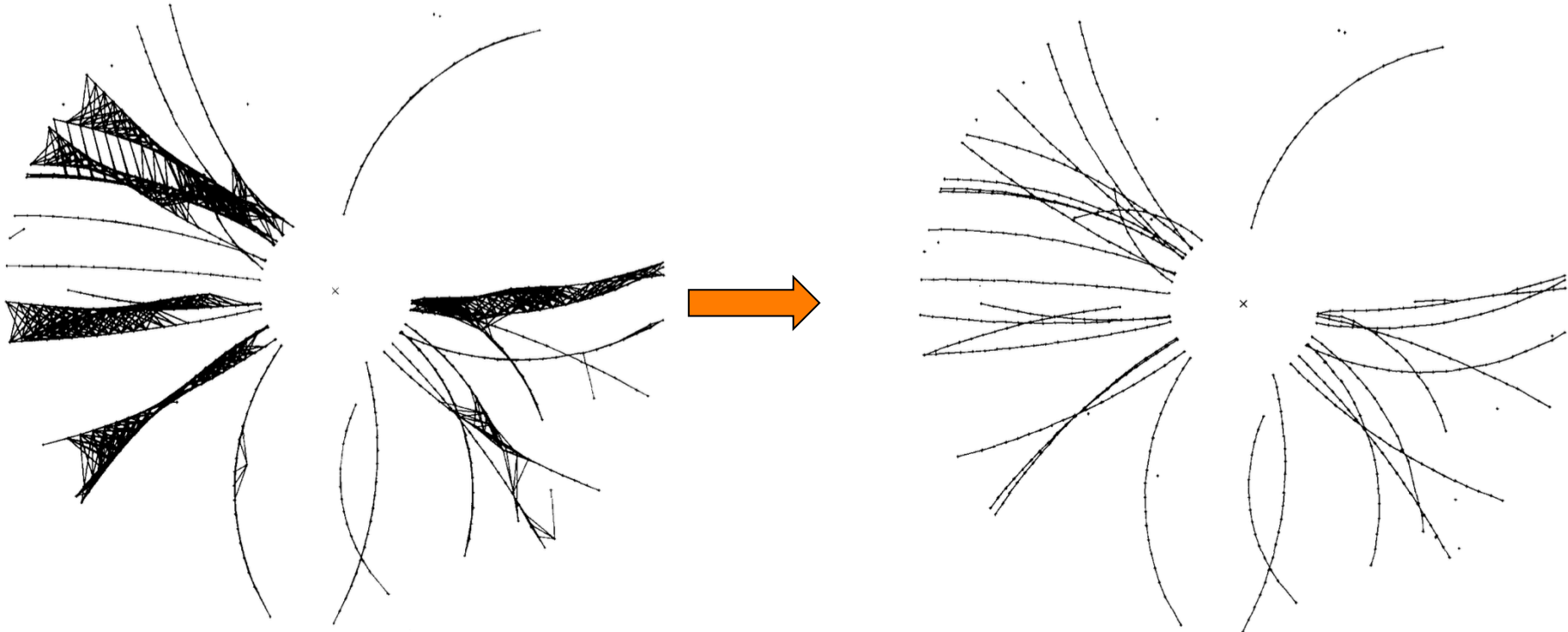- o Looking for long lasting collaborations beyond the challenge

HEP tracking…



25

...fascinates ML experts

Openlab ML workshop, ATLAS, 29th April 2016

# TrackML : An early attempt



- ❑ Stimpfl-Abele and Garrido (1990) (ALEPH)
- ❑ All posssible neighbor connections are built, the correct ones selected by the NN (not used in production)
- ❑ Also PhD Vicens Gaitan 1993, winner of Flavour of Physics challenge

# TrackML :CPU measurement

- Contrary to HiggsML or Flavour of Physic challenge need to evaluate CPU time
  - We know already how to solve the problem, but not quick enough (by factors)
  - CPU time to find the tracks
  - Cap on memory used (e.g. one x86-64 core with 2GB)
  - Training time unlimited
- Some platforms (see AutoML, Codalab, Topcoder) now allow to automatically upload, compile and run software
  - ➔well defined hardware (CPU and memory available)
  - ➔uniform comparison
  - Could also use an Amazon instance
- Positive side-effect : limit diversity of software languages and libraries
- Also the training dataset is very large (~1TB), better left on the platform
- We're more interested in the detailed algorithm (as it would be explained in a technical paper) rather than the software itself (but we do want to see the software)
- We're more interested in new approaches than in super-optimised version of old approaches
- We're looking for industry involvement there : a powerful platform to handle 100-1000 training per day on 1TB dataset

# Conclusion

- ❑ Machine Learning techniques widely used on ATLAS
- ❑ Recent explosion of novel (for HEP) ML techniques, novel applications for Analysis, Reconstruction, Simulation, Trigger, and Computing
- ❑ ATLAS Machine Learning group being set up by June 2016
- ❑ Building collaborations with ML scientists on a variety of issues
- ❑ Looking for powerful hardware+software platforms for:
  - o Deep NN training
  - o Tracking ML challenge