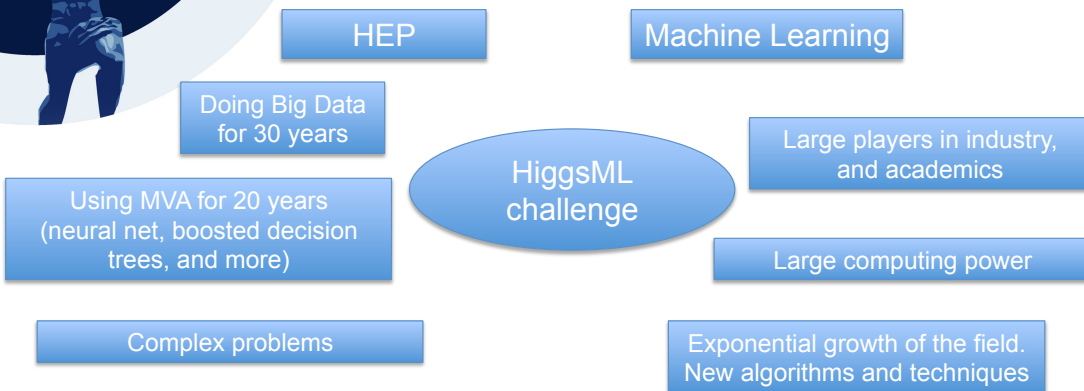
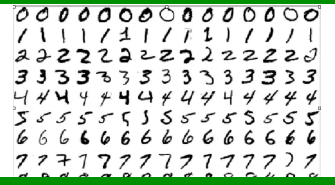


## The Higgs Boson Machine Learning Challenge

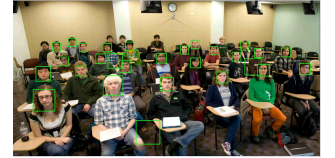
Will Davey (Universität Bonn, GE) for the ATLAS Collaboration



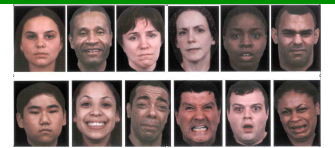
### Character recognition



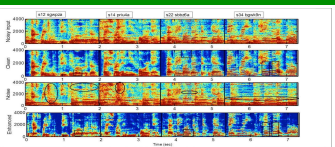
### Real time face detection



### Emotion recognition



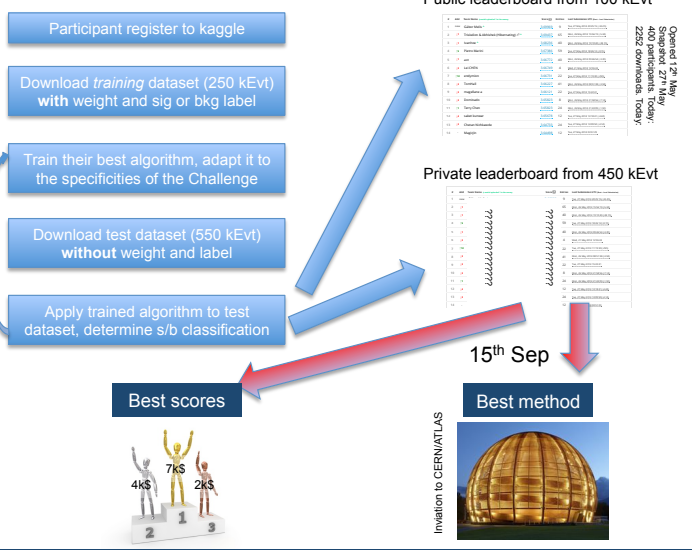
### Speech recognition



Let's put some ATLAS simulated data on the web and ask data scientists to develop the best machine learning algorithm to find the Higgs !

Jointly organised by ATLAS physicists and data scientists:

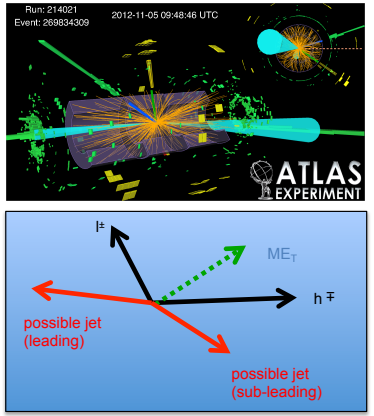
### How does it work ?



### Choice of ATLAS H to tau tau analysis

- "Evidence for Higgs Boson Decays to tau tau Final State with the ATLAS detector" ATLAS CONF-2013-108
- 4.1 sigma observed (3.2 sigma expected)
- Direct evidence of Higgs coupling to leptons
- Complex analysis (e.g. signal H mass ~ one sigma from dominant Z to tau tau background, VBF signature...)

### lepton-hadron topology



### Choice of figure of merit

Need one and only one robust estimator of the quality of the classification algorithm : Approximate Medium Significance **AMS**

Given  $s$  and  $b$  expected number of signal and background normalised to 2012 luminosity:  
 $s = \sum(\text{selected signal})$  weights,  
 $b = \sum(\text{selected bkg})$  weights,

Decided to use the "Asimov" formula (G. Cowan, K. Cranmer, E. Gross, and O. Vitells, "Asymptotic formulae for likelihood-based tests of new physics", EPJCC, vol. 71, pp. 1–19, 2011.) with « regularization » :  
 $b = b + 10$  (avoid large fluctuations in small regions)

$$AMS = \sqrt{2} \sqrt{(s+b)} \log(1 + s/b) - s$$

### Real analysis vs Challenge

Real analysis	Challenge
1. Systematics	1. No systematics
2. 2 categories x n BDT score bins	2. No categories, one signal region
3. Background estimated from data (embedded, anti tau, control region) and some MC	3. Straight use of ATLAS Geant4 MonteCarlo : signal Higgs, backgrounds Z, W, top
4. Weights include all corrections. Some negative weights.	4. Weights include normalisation and generator weight. Negative weight events rejected.
5. Potentially use any information from all 2012 data and MonteCarlo events	5. Only use variables and events preselected by the real analysis
6. Few variables fed in two BDT	6. All BDT variables + categorisation variables + primitives 3-vector
7. Significance from complete fit with Nuisance Parameters, Control Regions, etc...	7. Significance from "regularised Asimov"
8. MVA with TMVA BDT	8. MVA "no-limit"

### Variables provided

Weight and signal/background label	PRimitive 3-vectors allowing to compute the DER variables (mass neglected)	
weight	PRI_tau_pt	
label	PRI_tau_eta	
Conference note DER ived variables used for categorization or BDT (VBF and Boosted categories):	PRI_tau_phi	
	PRI_jet_num	
	DER_mass_MMC	PRI_jet_leading_pt
	DER_mass_trans_met_jet	PRI_jet_leading_eta
	DER_mass_vis	PRI_jet_leading_phi
	DER_pt_h	PRI_jet_subleading_pt
	DER_deltaeta_jet_jet	PRI_jet_subleading_eta
	DER_mass_jet_jet	PRI_jet_subleading_phi
	DER_prodelta_jet_jet	PRI_jet_subleading_pt
		PRI_jet_all_pt

Simpler, but not simple!

<https://www.kaggle.com/c/higgs-boson>

Contact : higgsml@lal.in2p3.fr

