

Deep Learning - Challenge Problem

Gregor Kasieczka
(gregor.kasieczka@uni-hamburg.de)

1st Terascale ML School
2018-10-25



Universität Hamburg

DER FORSCHUNG | DER LEHRE | DER BILDUNG

Emmy
Noether-
Programm

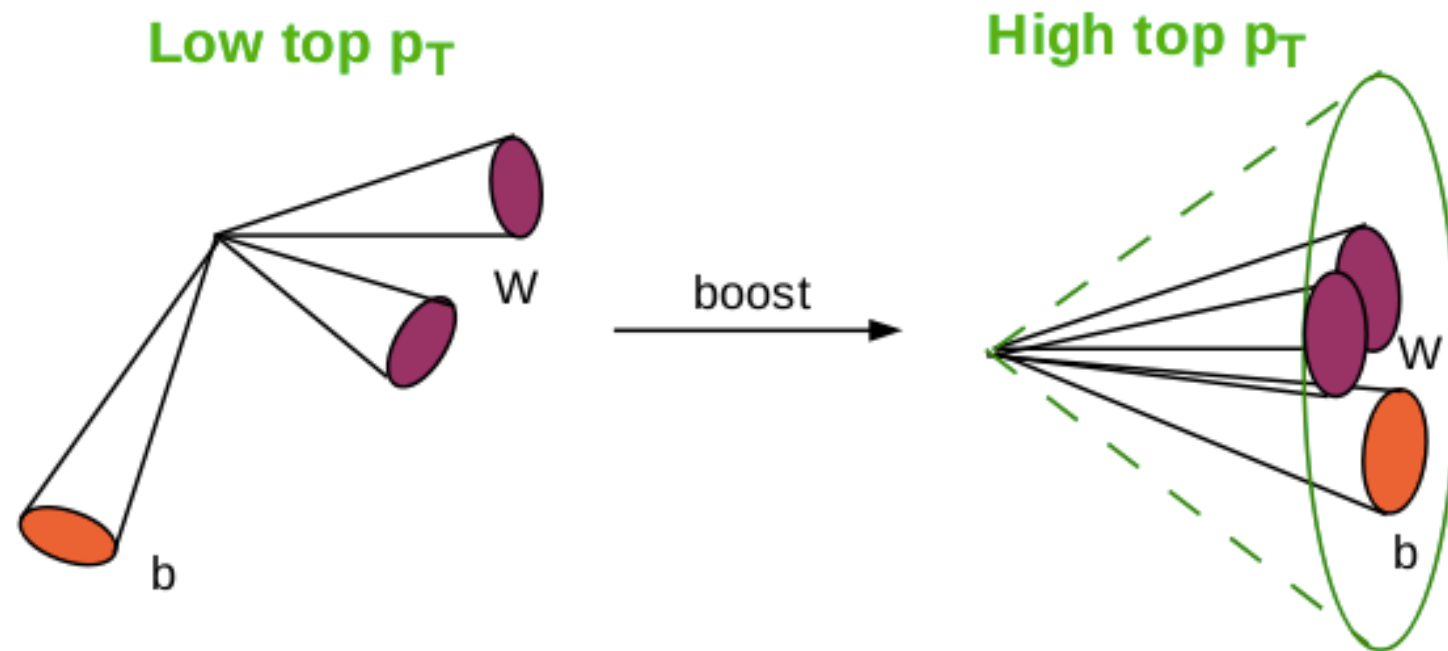
Deutsche
Forschungsgemeinschaft

DFG



Bundesministerium
für Bildung
und Forschung

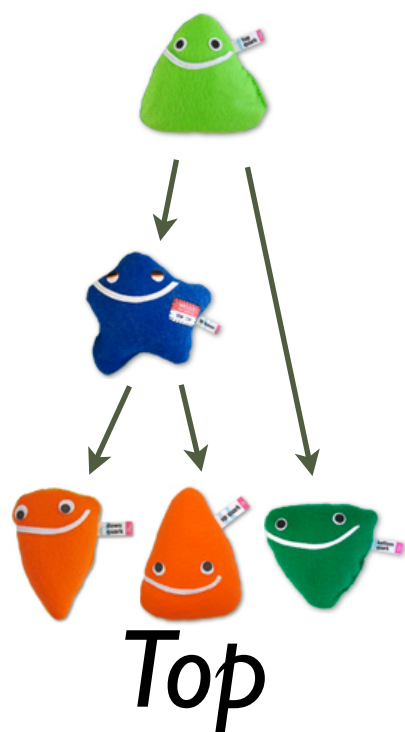
Heavy Resonance Tagging



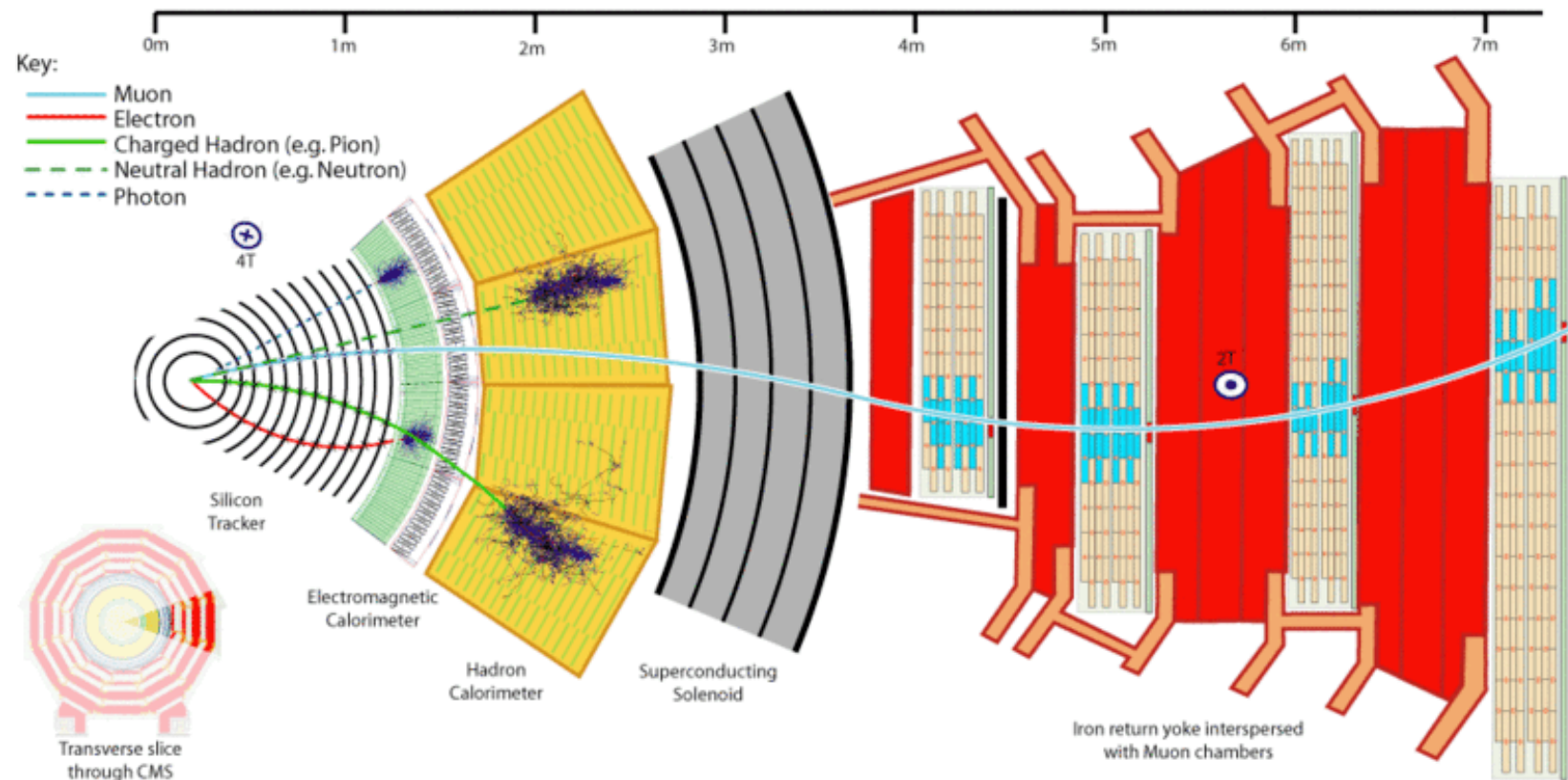
- Hadronically decaying top/Higgs/W/Z
- Contained in one (large-R) jet
- How to distinguish from light quark/gluon jets (and from each other)
- For new physics searches (and SM studies)

*Some Classical solutions:
(aka jet substructure)*

- Mass
Calculate using a grooming algorithm (eg mMDT/softdrop or pruning)
- Centers of hard radiation
n-subjettiness or energy correlation functions
- Flavour
b tagging of large-R jets or subjets
- Soft substructure
Color connection
- Inclusive reconstruction
HEPTopTagger V2, HOTVR
- Other substructure variables
Shower deconstruction, template tagger, ...



+



- Reconstruct energy with calorimeter
(improve resolution using tracker)
- Cluster energy deposits into jet
- Preprocess:

• center → rotate → flip (twice) → pixelate → crop → normalise

– center: centroid is at (0/0)

– rotate: principal axis is vertical

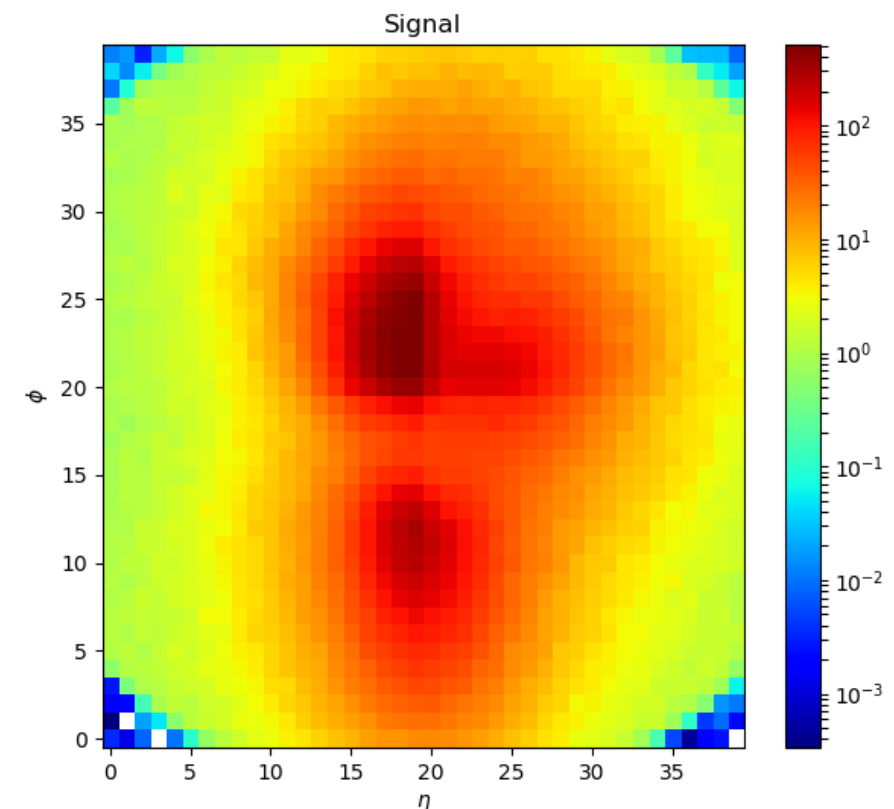
– flip: in $(x < 0, y > 0)$ -plane maximum intensity

– crop: to $n \times n$ images

– normalise: intensity of each pixel divided by total intensity

=

(Overlay of 100k images)



(jet images by Michel Luchmann)

Dataset Basics

- Goal: discriminate between hadronically decaying top quarks and QCD jets
 - Maximise area under ROC curve (AUC)
- Theory setting:
 - 14 TeV, hadronic tops for signal, qcd dijets background, delphes ATLAS detector card with Pythia
 - No MPI/pile-up included
 - We cluster particle-flow entries (produced by Delphes E-flow) into anti-kT 0.8 jets in the pT range [550,650]
 - All top jets are matched to a parton-level top within $\Delta R = 0.8$, and to all top decay partons within 0.8
 - We also require $|\eta|_{\text{jet}} < 2$

Version A: Constituents

- The leading 200 jet constituent four-momenta are stored, with zero-padding for jets with fewer than 200
- Constituents are sorted by pT, with the highest pT one first
- The truth top four-momentum is stored as truth_px etc.
- A flag (1 for top, 0 for QCD) is kept for each jet. It is called is_signal_new
- The variable "ttv" (= test/train/validation) is kept for each jet. It indicates to which dataset the jet belongs. It is redundant as the different sets are already distributed as different files.
- Either work with a fully connected architecture or try something else

Version B: Images

- center \rightarrow rotate \rightarrow flip (twice) \rightarrow pixelate \rightarrow crop \rightarrow normalise
 - center: centroid is at (0/0)
 - rotate: principal axis is vertical
 - flip: in $(x < 0, y > 0)$ -plane maximum intensity
 - crop: to $n \times n$ images
 - normalise: intensity of each pixel divided by total intensity

- Images after preprocessing of constituents
- Use a convolutional approach
- Provided by Michel Luchmann (Heidelberg)

Deep-learning Top Taggers or The End of QCD?
GK, Tilman Plehn, Michael Russell, Torben Schell
JHEP 05 (2017) 006
Pulling Out All the Tops with Computer Vision and Deep Learning
S Macaluso, D Shih, 1803.00107

Ingredients

- Dense Layer (Fully connected layer)
`keras.layers.Dense(number_of_nodes, activation='relu')`
- Convolutional Layer (For image data, `kernel_size=(2,2)` corresponds to a 2x2 pixel large filter)
`keras.layers.Conv2D(number_of_filters, kernel_size, padding='same', activation='relu')`
- Flatten Layer (turn image data into a list of numbers. Needed to go from Conv to Dense layers)
`keras.layers.Flatten()`
- First layer:
 - Can be anything, but needs arguments `input_shape=(80,)` (constituents) or `input_shape=(40,40,1)` (images) as extra argument. Only the first layer should have the input shape argument.
- Last layer of your network:
 - `keras.layers.Dense(2, activation='softmax')`
 - (recommend to start with `activation='relu'` for all other layers)

Documentation at: <https://keras.io/>

Some ideas

- Both versions:
 - more data, learning rate, optimiser algorithm, epochs, dropout, early stopping, **Smart use of val?**
- Version A (constituents)
 - more constituents, calculate physics quantities (mass?), other preprocessing (normalisation?), more layers, more nodes, ...
- Version B (images)
 - More filters, filter size, more Conv layers, pooling, more dense layers/nodes, image processing, locally connected layers,

Getting Started

- Login to your AWS machine
 - `mv .c1ng challenge`
 - `cd challenge`
 - `jupyter notebook`
- Connect to the notebook server

Overview

<input type="checkbox"/>	0	▼	📁 /	
<input type="checkbox"/>	📄	do_constit.ipynb	Constituents	
<input type="checkbox"/>	📄	do_images.ipynb	Images	
<input type="checkbox"/>	📄	test_without_truth_100k.h5		
<input type="checkbox"/>	📄	test_without_truth_img_100k.h5		
<input type="checkbox"/>	📄	train.h5		
<input type="checkbox"/>	📄	train_img.h5		
<input type="checkbox"/>	📄	val.h5		
<input type="checkbox"/>	📄	val_img.h5		

Training

train.h5 and train_img.h5

1.2M labelled examples
Nominal sample

Validation

val.h5 and val_img.h5

20k labelled examples
Systematic applied: like “test”

Test

test_without_truth_100k.h5 and
test_without_truth_img_100k.h5

100k *unlabelled* examples
Systematic applied: like “real” data

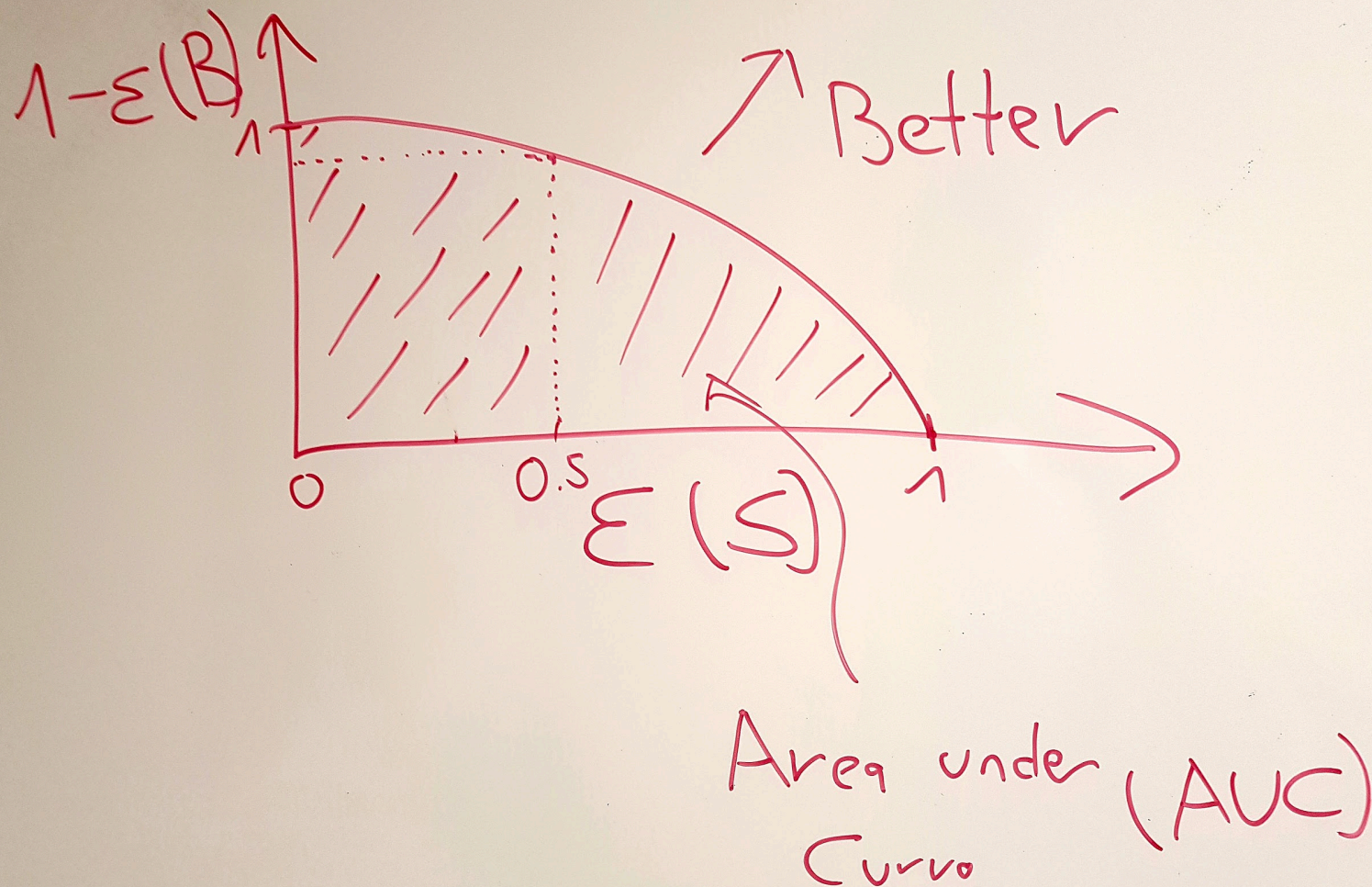
Data

- Based on
 - <https://goo.gl/XGYju3>
- But with a systematic twist (-;
(Thanks to Sven!)

Your task

- Calculate probability to be signal for each jet in the test sample
- Make sure to *keep the ordering!*
- Use images or constituents as input - both are fine.
- Produce a zipped .npy file and upload to challenge page
- Winner is best area-under-curve on subset of the test sample*
(*we keep the right to veto solutions)
- Submissions close **today at 23:00**
(Hamburg time)

Performance Measures



ROC:
Receiver operation characteristic

AUC:
Area under curve

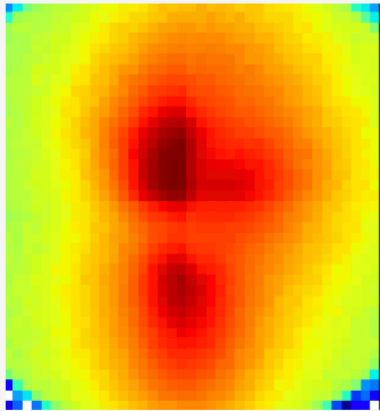
Other caveats

- There will be a small symbolic price for the best submission
- We ask the three best submissions to give a brief presentation tomorrow
- Team work is encouraged, individual submissions are ok as well
- Fair-play: Don't mess with other peoples machines, etc.
- However, I think there is one clever “cheaty” (that does not involve hacking, etc.) way for a perfect score

Way to submission

- Train network
- Apply to test sample
- Download result.zip file
- Go to:
 - <https://competitions.codalab.org/competitions/20432>
- (you might need to register at Codalab)
- Submit your result!

Score in Leaderbord is based on random 10% of data!
Use full sample for final ranking..



Top Quark Tagging Challenge

Secret url: https://competitions.codalab.org/competitions/20432?secret_key=e7c978d6-683e-459f-95e7-720dd4e0f086

Organized by gregor.kasieczka - Current server time: Oct. 24, 2018, 5:08 p.m. UTC

► Current

Challenge Phase

June 30, 2013, midnight UTC

End

Competition Ends

Oct. 25, 2018, 11 p.m. UTC

Learn the Details

Phases

Participate

Results

Overview

Evaluation

Welcome!

A Challenge of tagging top quarks with systematic uncertainties. Organised as part of the [1st Terascale School of ML](#).

[Learn the Details](#)[Phases](#)[Participate](#)[Results](#)[Get Data](#)[Files](#)[Submit / View Results](#)**Challenge Phase**

Phase description

None

Max submissions per day: 20**Max submissions total: 20**

Click the Submit button to upload a new submission.

Optionally add more information about this submission

Submit

Here are your submissions to date (✓ indicates submission on leaderboard):

#	SCORE	FILENAME	SUBMISSION DATE	STATUS	✓	
1	0.9043234788	result (1).zip	10/24/2018 16:57:00	Finished	✓	+