# Attention-based reconstruction for $t\bar{t}H(b\bar{b})$ in CMS

L. Benato<sup>1</sup>, A. Calandri<sup>2</sup>, M. Donega<sup>2</sup>, A. Gomez Espinosa<sup>2</sup>, G. Kasieczka<sup>1</sup>, <u>T. Lösche<sup>1</sup></u>, M. Meinhard<sup>2</sup>, C. Reissel<sup>2</sup>, D. Ruini<sup>2</sup>, R. Wallny<sup>2</sup>

"Physics at the Terascale" Annual Meeting, 26.11.2019

<sup>1</sup>IExp, Universität Hamburg; <sup>2</sup>IPA, ETH Zürich



**ETH** zürich



## Introduction - ttH motivation





- SM: fermion masses accounted for by Yukawa interactions between Higgs and fermion field
- + top Yukawa coupling essential probe of SM (expected:  $Y_{t}\approx$  1)
- coupling can be observed in Higgs production in association with top-quarks
- consider single-lepton channel

## Introduction - $t\bar{t}H(b\bar{b})$ challenges





Problems:

- high number of jets in final state
- irreducible background ( $t\bar{t} + b\bar{b}$ )  $\rightarrow$  cross section 8 times higher than signal
- often final state not completely reconstructed

 $\Rightarrow$  need to handle combinatorial assignment of jets to improve discrimination

#### Current state



- current state-of-the-art for SL channel: ANN output used as final discriminant [CMS-PAS-HIG-18-030]
- inputs: basic object kinematics and high-level variables (e.g. Matrix-Element-Method (MEM) discriminant)
- problem: MEM discriminant time consuming in computation (up to 10 min. per event on single CPU)
- solution: develop new DNN architecture with competitive discrimination power using simple kinematics as input

aivozritāt Hambur

#### **COBRA - attention based DNN**





A woman is throwing a <u>frisbee</u> in a park.



A little <u>girl</u> sitting on a bed with a teddy bear.

arXiv:1502.03044 [cs.LG]

- COBRA = <u>COmBinatoRics</u> based deep <u>Attention network</u>, developed and implemented at ETH Zurich
- attention originally developed by Google (arXiv:1706.03762 [cs.CL])
- allows network to focus on subset of inputs (e.g. certain parts of image)
- can detect specific features
- extensively used in image recognition and natural language processing
- idea for ttH: focus on jet-combinations to get a handle on combinatorics

## **COBRA** architecture & inputs





## input variables (Monte Carlo simulation):

- jets [8 floats/jet]:
  - $\eta$ ,  $\phi$ , Energy,  $p_X$ ,  $p_y$ ,  $p_z$ ,  $p_T$ , b-tag
- leptons [7 floats]:
  - +  $\eta$ ,  $\phi$ , Energy, p<sub>X</sub>, p<sub>y</sub>, p<sub>z</sub>, p<sub>T</sub>
- MET [5 floats]:
  - +  $\phi$ ,  $\sum E_{T}$ ,  $p_X$ ,  $p_y$ ,  $p_T$
- no prior categorization based on jet or b-tag multiplicity
- always one lepton and ten jets (zero-padded if less) per event  $[8 \cdot 10 + 7 + 5 = 92$  floats total]

#### Attention mechanism









- completely matched sample:
  - 3 jets matched to hadronically decaying top
  - 2 jets matched to Higgs
  - 1 jet matched to leptonically decaying top
- expectation: Higher attention weights for higher number of matched jets ightarrow not the case









- combinations with higher number of jets matched have higher attention
- accuracies (separately trained classifiers): trijet classifier 64.35%; Higgs classifier 38.7%; leptonic top classifier 71.86%





- use of pre-trained classifier does not enhance performance when trained on completely matched sample
- accuracy of 38.7% of dijet classifier probably too low to cause performance increase
- provide network with "optimal" dijet attention:
  - + no jet matched  $\rightarrow$  0
  - + 1 jet matched  $\rightarrow$  1
  - + 2 jets matched  $\rightarrow$  2
  - $\cdot\,$  normalise scores such that the sum is one





- AUC increases with provided attention accuracy
- AUC of pretrained COBRA lower than value at  $\approx$ 39%
- further research in enhancing the ability of COBRA to find the Higgs necessary



- tTH important for better understanding mass generation in SM
- major challenge: irreducible background ( $t\bar{t} + b\bar{b}$ ) 8 times higher than signal
- we investigated the use of attention to improve combinatorial assignment of jets in  $\ensuremath{t\bar{t}}\xspace H$
- when used for signal/background classification, attention distribution independent of number of matched jets
- if trained as classifier  $\rightarrow$  attention increases for higher number of matched jets
- proof-of-principle: performance increases with attention accuracy
- pretrained classifiers show similar perfromance as  $\text{FCN} \rightarrow \text{accuracy of attention}$  network needs to be improved

# Backup

#### Background



- tt + b-jets: extra b-jets from (overlapping)
   b-hadrons
- tt + cc: at least one extra charm jet from one or more overlapping hadrons
- tt + light flavour (LF): events which do not fit in other categories
- minor backgrounds: tt-production in association with massive vector bosons or jets, production of W, Z, or γ with jets, diboson processes and single top quark production

### MEM method

- Problem: difficult to discriminate between  $t\bar{t}H(b\bar{b})$  and  $t\bar{t} + b\bar{b}$  $\Rightarrow$  Matrix-Element-Method (MEM)
- uses both experimental and theoretical information by evaluating scattering amplitudes
- for every event, two associated weights are computed according to:

$$\mathsf{P}(\mathbf{x}|lpha) \propto rac{1}{\sigma_lpha} \int \mathrm{d} \mathbf{\Phi}(\mathbf{y}) \quad |\mathcal{M}_lpha|^2 \left(\mathbf{y}
ight) \quad \mathcal{W}(\mathbf{x},\mathbf{y}) \quad \mathcal{W}(\mathbf{x},\mathbf{y})$$

- +  $\alpha$  can stand for signal or background-only hypothesis
- $|M_{\alpha}|^2$  (y) is LO matrix element,  $W(\mathbf{x}, \mathbf{y})$  transfer function simulating detector response with  $\mathbf{x}$  being measured and  $\mathbf{y}$  the true parameters
- discriminant is likelihood ratio:  $M_i = \frac{P(\mathbf{x}_i | sig)}{P(\mathbf{x}_i | sig) + k \cdot P(\mathbf{x}_i | bkg)}$ ; k = optimised parameter

#### Performance comparison



- sample:  $\geq$  4 jets,  $\geq$  3 btags
- no significant difference in performance visible
- other architectures (e.g. with pretrained classifiers) show no performance increase
   → rate at which correct Higgs combination is found still too low

## Comparison of b-tagging algorithms



Full sample CSV v2



- AUC performance decrease for CSV v2
- no significant change in the difference of performance between FCN and COBRA

DNN model	COBRA w/ att.	COBRA w/o att.	FCN comb.	simple FCN	Current CMS
AUC	0.80	0.80	0.79	0.80	N/A
Upper stat. uncertainty	0.194	0.206	0.203	0.194	0.22
Lower stat. uncertainty	0.194	0.206	0.203	0.194	0.21

## Object selection and samples

Object selection	SL		
Application of all MET filters	$\checkmark$		
Electrons: tight ID, isolation cuts	$ m p_T > 30$ GeV, $ \eta  < 2.1$		
Muons: tight ID, isolation cuts	$p_{\mathrm{T}} > 26\mathrm{GeV}$ , $ \eta  < 2.1$		
Jets: AK4PFchs,	$p_{ m T} > 30~{ m GeV},  \eta  < 2.4$		
Loose ID with lepton removal and PU jet ID			
b-tagging via deepCSV			
MET	> 20 GeV		

Samples:

- process: ttH(bb) SL; tag: 12Apr2018
- process: ttjets SL; tag: 12Apr2018\_new\_pmx
- number of events used in training  $\approx$  12.5 Mio.
- no further event selection used in training