## Boosted Decision Trees and *b*-jet Trigger Calibration Studies for $t\bar{t}H(b\bar{b})$ Fully Hadronic Analysis

Bartolini Giovanni

CPPM University of Aix-Marseille

October 18, 2018







### Outline

#### Introduction

- The Standard Model of Particle Phisics
- The ATLAS experiment at LHC
- Higgs Boson and Top Yukawa coupling

#### 2 Fully Hadronic ttH Analysis Strategy

#### **3 BDT studies**

- Reconstruction BDT Step
- Classification BDT Step

#### 4 b-jet Trigger Calibration with tt Dilepton Events

#### 5 Conclusion

## Introduction

### The Standard Model of Particle Phisics

- The **Standard Model** (SM) of Particle Physics is a gauge theory that classifies all known elementary particle and describes 3 of the 4 known interaction forces: Strong, Weak and Electromagnetic
- Discovery of **Higgs** boson in 2012 completed the set of predicted elementary particle
- Very successful and predictive theory, but has still many shortcomings:
  - inclusion of gravity
  - neutrino masses and oscillations
  - matter/anti-matter asymmetry
  - evidence of dark matter existence
- Strong chase for New Physics
  - direct search of new particles
  - indirect evidence trough deviations in SM predictions

#### Standard Model of Elementary Particles



## The ATLAS experiment at LHC

- The ATLAS experiment is placed in one of the 4 interaction point of LHC
- Almost at the end of Run 2:
  - ▶ already collected ~130 fb<sup>-1</sup>







## Higgs Boson and Top Yukawa Coupling

- The discovery of the **Higgs** boson in 2012 started an effort on the precise measure of its properties
- The top quark is the heaviest elementary particle
  - has the highest Yukawa coupling: Y<sub>t</sub>
- Anomalous values for *Y<sub>t</sub>* are predicted by many *Beyond the Standard Model* (BSM) theories
- Associated production  $(t\bar{t}H)$  only way to directly measure  $Y_t$ 
  - many accessible final states: γγ, multi-lepton, lepton+jets, all hadronic
- ATLAS first *ttH* observation recently published using Run 1 and Run 2 data\*
  - with observed(expected) significance of 6.3(5.1)





arXiv:1806.00425 [hep-ex]

## Fully Hadronic ttH Analysis Strategy

## $t\bar{t}H$ Production in Fully Hadronic Final State

- Most abundant final state
  - ▶  $\sim$ 33% of total  $t\bar{t}H$  production
- No neutrino
  - full event reconstruction
- Ideal for differential analysis
  - explore the CP nature of Y<sub>t</sub>
- Challenging experimental signature:
  - ▶ 8 quarks, 4 *b*-quarks
  - Large QCD background
  - irreducible  $t\bar{t} + b\bar{b}$  background
- First analysis published in Run 1<sup>\*</sup>:  $\mu_{t\bar{t}H} = 1.6 \pm 2.6$ 
  - ► brought 10% improvement on the significance of the  $t\bar{t}H(H \rightarrow b\bar{b})$  combination





arXiv:1604.03812 [hep-ex]

### Jets as Quark objects

- The nature of the strong interaction does not allow to quarks to be available as free particles
  - How can we see quarks as final state particles?
- Quarks and gluons that are produced from a collision will *hadronize* producing a collimated flow of hadronic particles
  - a specific algorithm is used to reconstruct this flow as a single object, called JET
- Jets coming from *b* quarks have particolar properties that allow them to be distinguished from other jets
  - ▶ the identification of a jet as coming from a *b* quark is called *b*-tagging



### **Event Selection and Categorization**

- Multijet trigger
  - ▶ ≥4 jets with  $p_T$  > 100 GeV (120 GeV for 2017 data) and  $|\eta| < 2.5$
- Lepton veto
  - to avoid overlap with other channels
- Offline selection:
  - same as the trigger + ≥ 2 additional jets with p<sub>T</sub> > 25 GeV and |η| < 2.5</p>
  - ≥ 2 jets b-tagged
- Categorization in jet and *b*-jets multiplicity
  - ▶ 6, 7, 8 or ≥9 jets
  - ▶ 2, 3 or ≥4 *b*-jets

ATLAS (s = 13 TeV allhad	Internal	[]fi+lig][fi+V []fi+≥][Jii+≥1b []Single[]CβCD	ATLAS 15 = 13 TeV, 36.1 f alhad	Internal b <sup>-1</sup>
fije Stor	6je 40		6 is 3be 5 5 - 0.1%	6 (e 4b) 50 = 0.3% 57 = 0.5
7/e 500	Tye 4b		() () () () () () () () () ()	7)e 4bi sn = 0.6%
Eje Ste	51+ 40		Bie 3be se = 0.5% 0 0.5	Bije 4bi 55 = 1.0%
59.50e	59-46i		(m) $\overline{0}$ 0.5	9(j 4b) 50 = 1.5% 50 0.5

data 15+16	8j,3b	8j,≥4b	≥9j,3b	≥9j,≥4b
ttH	$50.9 \pm 0.3$	$13.9 \pm 0.2$	$76.4 \pm 0.4$	$24.5 \pm 0.2$
ttb	$1033.7 \pm 25.4$	$149.3 \pm 9.6$	$1502.6 \pm 30.8$	$281.9 \pm 12.8$
ttc	$360.4 \pm 15.8$	$10.1 \pm 2.7$	$437.9 \pm 17.5$	$17.6 \pm 3.6$
ttl	$505.0 \pm 19.0$	$8.0 \pm 2.5$	$431.3 \pm 17.6$	$4.3 \pm 1.8$
ttV	$30.9 \pm 1.0$	$5.2 \pm 0.4$	$45.2 \pm 1.5$	$10.2 \pm 0.7$
singletop	63.5 ± 4.2	$6.5 \pm 1.4$	57.3 ± 3.9	$6.3 \pm 1.3$
QCD	$8946.1 \pm 35.6$	$1203.9 \pm 5.1$	$7238.9 \pm 34.9$	$1196.9 \pm 5.9$
Total Background	$10940 \pm 50$	$1373 \pm 12$	$9713 \pm 53$	$1518 \pm 15$

### **Event Selection and Categorization**

- Multijet trigger
  - ▶ ≥4 jets with  $p_T$  > 100 GeV (120 GeV for 2017 data) and  $|\eta| < 2.5$
- Lepton veto
  - to avoid overlap with other channels
- Offline selection:
  - same as the trigger + ≥ 2 additional jets with p<sub>T</sub> > 25 GeV and |η| < 2.5</p>
  - ≥ 2 jets b-tagged
- Categorization in jet and *b*-jets multiplicity
  - ▶ 6, 7, 8 or ≥9 jets
  - ▶ 2, 3 or ≥4 *b*-jets

ATLAS (s = 13 TeV alltad	Internal	[]tī+ko]]tī+∀ []tī+≥[]tī+≥1b []Singko]]CβCD	ATLAS 15 = 13 TeV, 36.1 fb althad	Internal -1
fije Stor	6je 45		(g) 0:5-0.1%	(0) (0) (0) (0) (0) (0) (0) (0)
7/e 3be	7/0 451		(2) 300 50 = 0.3% 6 0.5	(0) (0) (0) (0) (0) (0) (0) (0)
Sie Ste	× 44		(g) (g) (g) (g) (g) (g) (g) (g) (g) (g)	Bio 4bi 56 = 1.0% 0 0.5
9/34	9i 4bi		(g) $\overline{0}$ 0.5	99 461 50 - 1.6%

data 17	8j,3b	8j,≥4b	≥9j,3b	≥9j,≥4b
ttH	$36.0 \pm 0.4$	$9.4 \pm 0.2$	$63.0 \pm 0.5$	$19.8 \pm 0.3$
ttb	$594.9 \pm 22.1$	$75.8 \pm 6.6$	$1023.2 \pm 29.7$	$178.8 \pm 11.6$
ttc	$197.2 \pm 12.9$	$9.2 \pm 4.1$	$321.1 \pm 17.3$	$7.0 \pm 2.0$
ttl	$278.1 \pm 16.6$	$5.4 \pm 2.7$	$278.1 \pm 16.7$	$4.7 \pm 1.9$
ttV	$20.6 \pm 1.1$	$3.6 \pm 0.5$	$32.6 \pm 1.6$	$5.8 \pm 0.7$
singletop	$46.2 \pm 4.2$	$1.7 \pm 0.6$	$45.7 \pm 4.1$	$4.7 \pm 1.0$
QCD	5748.7 ± 28.5	$763.6 \pm 4.0$	5279.2 ± 30.2	$859.1 \pm 5.1$
Total Background	6886 ± 42	$859 \pm 9$	$6980 \pm 49$	$1060 \pm 13$

## **Analysis Strategy**

- QCD multijet is main background, but really hard to reproduce in simulation
  - estimation with data-driven method: Tag Rate Function Multijet (TRF<sub>MJ</sub>)
- Signal vs background discrimination
  - with Multivariate Analysis (MVA) method called Boosted Decision Tree (BDT)
- Maximum likelihood fit
  - to extract the final results









## **BDT** studies

### Multivariate Analysis for Fully Hadronic $t\bar{t}H$



### **Reconstruction BDT Step**

- **Goal**: find the best association between jets reconstructed in the detector and the final state partons
- $\bullet$  Large mutliplicites  $\rightarrow$  large combinatorics
  - For 36 up to thousands possible ways to reconstruct the  $t\bar{t}H$  system



#### **Reconstruction BDT Step**

- 2 different BDTs using reconstructed resonances and angular correlations between jets
  - recoBDT: tries to reconstruct only tī system
    - $\rightarrow$  no bias on the Higgs candidate mass
  - recoBDT\_withHiggs: full tTH system reconstruction
    - $\rightarrow$  higher reconstruction efficiency
- Trained using  $t\bar{t}H$  simulation
  - signal: correct quarks-jets association
  - background: all other possible associations



#### **Reconstruction BDT Step**



 Higgs boson candidate properly reconstructed 57% and 75% of times and full event properly assigned 41% and 53% respectively in events with full matching available

## **Classification BDT Step**

- Goal: perform signal vs background discrimination
- Combines reconstruction results from previuos step with global event kinematics
- Optimization against 2 different main background sources:
  - irreducible  $tar{t} + bar{b} 
    ightarrow$  need more simulation
  - QCD multijet, using data-driven simulation from TRF<sub>MJ</sub>



### **Classification BDT Step**

- Variables optimization performed separately in each signal region with a recursive method
  - From a preselected set of variables takes the one with the highest s/b separation
  - Adds recursively the variable that brings the biggest improvement
  - stops when the improvement from the addiction of a new variable is less than 1%



### **Classification BDT Step**

- Variables optimization performed separately in each signal region with a recursive method
  - From a preselected set of variables takes the one with the highest s/b separation
  - Adds recursively the variable that brings the biggest improvement
  - stops when the improvement from the addiction of a new variable is less than 1%



# b-jet Trigger Calibration with tt Dilepton Events

## **b-jet Trigger Calibration**

- trigger-level *b*-tagging improved with respect to Run 1, now close in performance to offline algorithms
- ullet Use of b-jet trigger results on an increase of signal efficiency of a factor  ${\sim}3$



- Data/MC Scale factors are used to calibrate online algorithms
  - SFs are derived for different jet flavors
- b-tagging used both at trigger level and at recostruction level
  - combined online + offline calibration
  - Geometrical association between online and offline jet objects

## *b*-jet Trigger Calibration with $t\bar{t}$ Events

- Event selection:
  - activate one Lepton + boffperf trigger
  - exactly 2 leptons with  $p_T > 28$  GeV and opposite charge
  - exactly 2 jets with  $p_T > 35$  GeV
  - eµ channel
  - *m*<sub>lj</sub> cuts
  - both jets matched



#### Lepton + boffperf prescaled triggers

b-tagging algorithm will run for each jet without taking decisions

#### *m<sub>lj</sub>* cuts

- Idea: improve bb-purity by finding Jet + lepton combinations which corresponds to the top quarks.
- For *b*-Jets the invariant mass of the combination should be smaller then the top mass.
- The combination we found which seems to be the most promising in reducing background is the one which minimizes the sum of the squared invariant mass of both possible "Ij-combinations" in the event.

#### lj-Combination, invariant Mass



- Veto events with one  $m_{
m lj}$  > 175GeV (pprox top-mass), or constrain flavor fractions

### *b*-jet Trigger Calibration with $t\bar{t}$ Events

• b-jet content extracted using likelihood method

$$\mathcal{L}_{\rm E}(\boldsymbol{\rho}_{\rm T,1}, \boldsymbol{\rho}_{\rm T,2}, w_1, w_2 | \mathcal{P}_{\rm b}(w | \boldsymbol{\rho}_{\rm T})) = [f_{\rm bb}(\boldsymbol{\rho}_{\rm T,1}, \boldsymbol{\rho}_{\rm T,2}) \mathcal{P}_{\rm b}(w_1 | \boldsymbol{\rho}_{\rm T,1}) \mathcal{P}_{\rm b}(w_2 | \boldsymbol{\rho}_{\rm T,2}) \\ + f_{\rm bl}(\boldsymbol{\rho}_{\rm T,1}, \boldsymbol{\rho}_{\rm T,2}) \mathcal{P}_{\rm b}(w_1 | \boldsymbol{\rho}_{\rm T,1}) \mathcal{P}_{\rm l}(w_2 | \boldsymbol{\rho}_{\rm T,2}) \\ + f_{\rm lb}(\boldsymbol{\rho}_{\rm T,1}, \boldsymbol{\rho}_{\rm T,2}) \mathcal{P}_{\rm l}(w_1 | \boldsymbol{\rho}_{\rm T,1}) \mathcal{P}_{\rm b}(w_2 | \boldsymbol{\rho}_{\rm T,2}) \\ + f_{\rm ll}(\boldsymbol{\rho}_{\rm T,1}, \boldsymbol{\rho}_{\rm T,2}) \mathcal{P}_{\rm l}(w_1 | \boldsymbol{\rho}_{\rm T,1}) \mathcal{P}_{\rm l}(w_2 | \boldsymbol{\rho}_{\rm T,2})$$
(1)

$$\begin{split} &f_{f_1,f_2}(\rho_{\Gamma,1},\rho_{\Gamma,2}) = \text{fraction of flavour combination } [f_1,f_2]. \text{ (Extracted from Simulation)} \\ &\mathcal{P}_f(w_1|\rho_{\Gamma,1}) = \text{pdf for a b-tagging weight } w \text{ of jet with flavour } f \text{ and a given } \rho T, 1. \end{split}$$

$$\mathcal{L}(\mathcal{P}_{\mathrm{b}}(w|p_{\mathrm{T}})) = \prod_{\mathrm{data}} \mathcal{L}_{\mathrm{E}}(\mathrm{data}|\mathcal{P}_{\mathrm{b}}(w|p_{\mathrm{T}}))$$
(2)

$$\epsilon_{\rm b}(\boldsymbol{p}_{\rm T}) = \int_{w_{\rm cut}}^{\infty} dw' \mathcal{P}_{\rm b}(w_{\rm f}|\boldsymbol{p}_{\rm T}) \tag{3}$$

#### **Calibration Results: Online only**



#### **Calibration Results: Combined**



## **SF Stability Studies**

- Studied SF dependency to different possible source of sistematic uncertainties:
  - $\eta$ , pile-up and data period

85% WP

• SF seems to be robust to all of them







## Conclusion

#### Conclusion

- Fully hadronic  $t\bar{t}H$  channel:
  - Large BR, but dominated by QCD multi-jet production
  - Large statistic available and event fully reconstructable



- My contribution to improve the analysis:
  - Calibration of trigger b-tagging efficiency
    - ★ replace multijet trigger with *b*-jet trigger to increase by a factor ~3 the signal selection efficiency
  - Implementing a 2 steps strategy for MVA based signal/background discrimination
    - \* Reconstruction step to resolve combinatorics
    - \* Classification step with optimization for QCD background discrimination

# **END**

## Thanks for the attention!