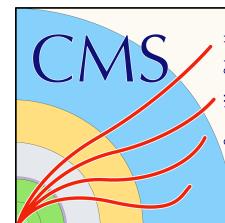


NTUA Top Tagger

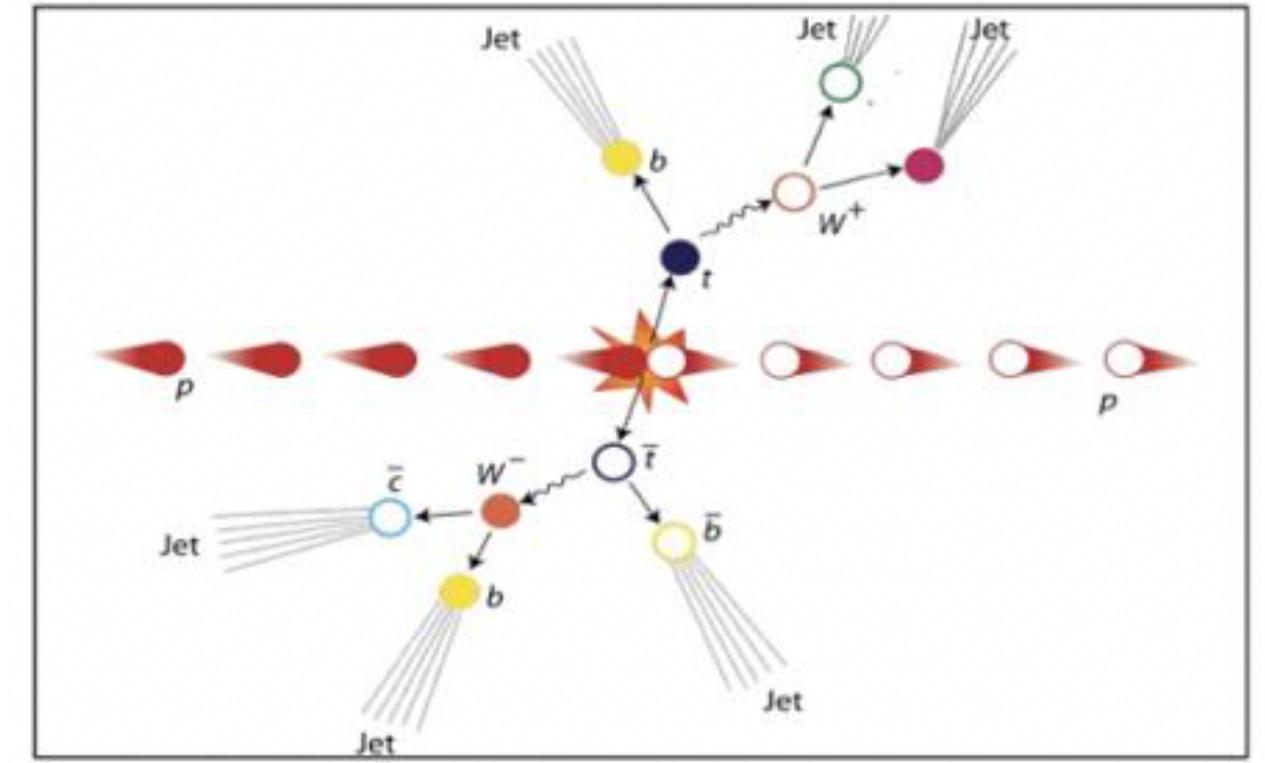
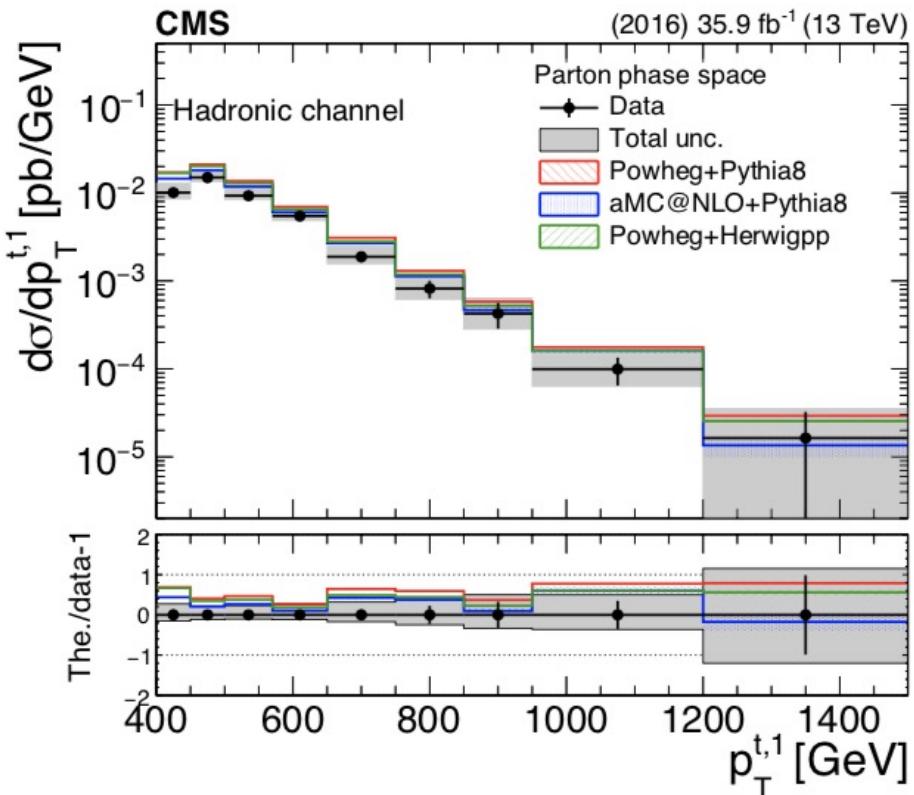
Tag & Probe methodology

G. Bakas, K. Kousouris, I. Papakrivopoulos, G. Tsipolitis



Analysis Overview

- Differential cross section for boosted ttbar pair fully hadronic final state
- Trying to identify two big jets that contain the products of the top/anti-top decay.



- A NN, for tagging ttbar events was used in TOP-18-013
- A BDT for tagging jets as tops is used in this analysis



Signal Selection

Variables	Selected Cut
pT leading jets	> 450 GeV
pT 2 nd leading jets	> 400 GeV
Njets	> 1
N leptons	= 0
eta (both leading jets)	< 2.4
mJJ	> 1000 GeV
jetMassSoftDrop (only for fit)	(50,300) GeV
Top Tagger	> 0.2
B tagging (2 btagged jets)	> Medium WP
Signal Trigger	

Control Region Selection

Variables	Selected Cut
pT leading jets	> 450 GeV
pT 2 nd leading jets	> 400 GeV
N leptons	= 0
eta (both leading jets)	< 2.4
mJJ	> 1000 GeV
jetMassSoftDrop (only for fit)	(50,300) GeV
Top Tagger	> 0.2
B tagging (0 btagged jets)	< Medium WP
Control Trigger	



Selection Overview

Region	Requirements
Signal Region (SR)	Baseline + topTagger + $m_{SD}^{jet1,2} \in (120,220)GeV + 2btags$
Control Region (CR)	Baseline + topTagger + $m_{SD}^{jet1,2} \in (120,220)GeV + 0btags$
Extended SR (SR_A) (QCD fit region)	Baseline + topTagger + $m_{SD}^{jet1,2} \in (50,300)GeV + 2btags$
Extended CR (CR_A) (QCD fit region)	Baseline + topTagger + $m_{SD}^{jet1,2} \in (50,300)GeV + 0btags$



Signal Extraction

$$S(x_{reco}) = D(x_{reco}) - C_{bkg}^{yield} N_{QCD}^{fit} C_{QCD}^{shape}(x_{reco}) Q(x_{reco}) - B(x_{reco})$$

Fiducial Yield

Transfer factor
from SR_A to SR

Measured dist
from data

Fitted number
of QCD events
in SR_A

QCD shape taken
from Data (CR)

QCD shape
correction factor

Subdominant bkg shape
and contribution (MC)



Signal Extraction

$$S(x_{reco}) = D(x_{reco}) - C_{bkg}^{yield} N_{QCD}^{fit} C_{QCD}^{shape}(x_{reco}) Q(x_{reco}) - B(x_{reco})$$

Fiducial Yield

Measured dist from data

QCD shape correction factor

QCD shape from Data

Number of QCD events in SR_A

Subdominant bkg contribution (MC)

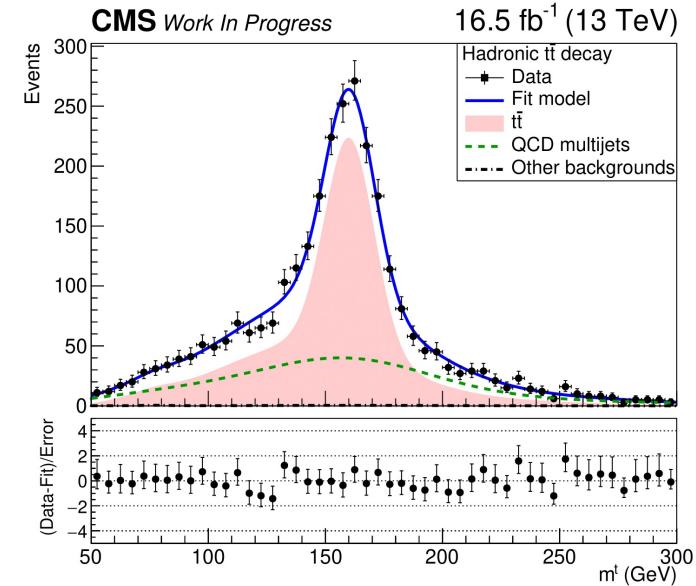
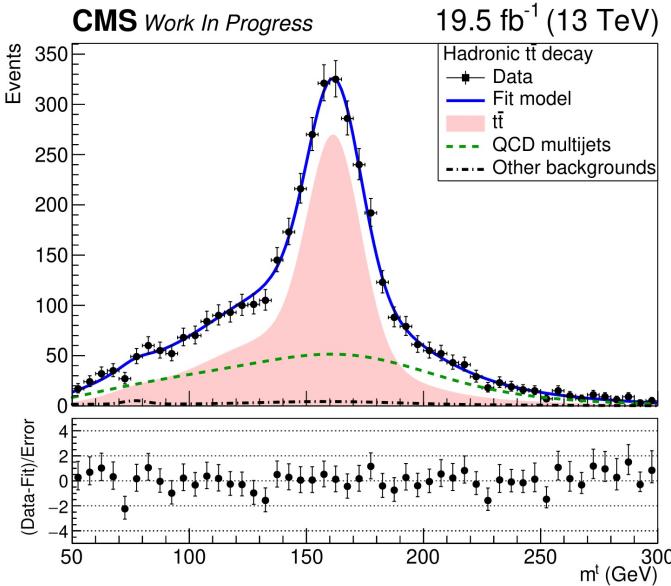
- We deploy a fit in the extended signal region (SR_A) to determine the Number of QCD events N_{QCD}^{fit}

$$D(m^t)^{(i)} = N_{tt}^{(i)} T^{(i)}(m^t, k_{MassScale}, k_{MassResolution}) + N_{bkg}^{(i)} B(m^t)(1 + k_1 x) + N_{sub}^{(i)} O^{(i)}(m^t)$$

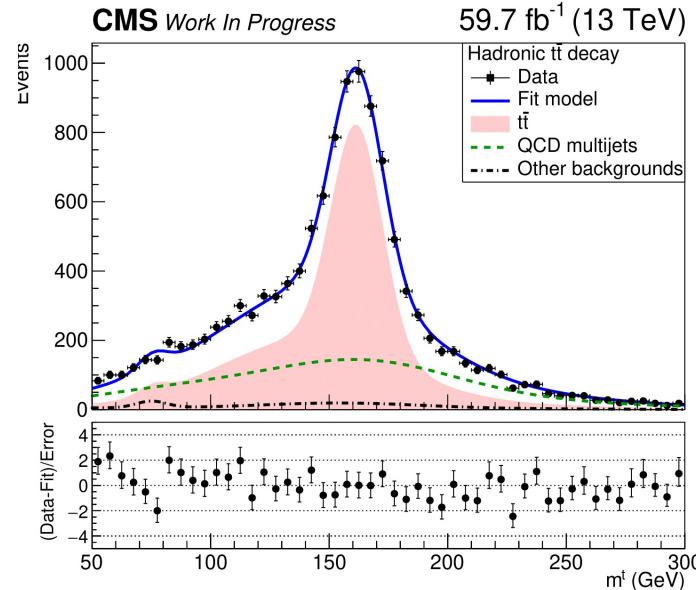
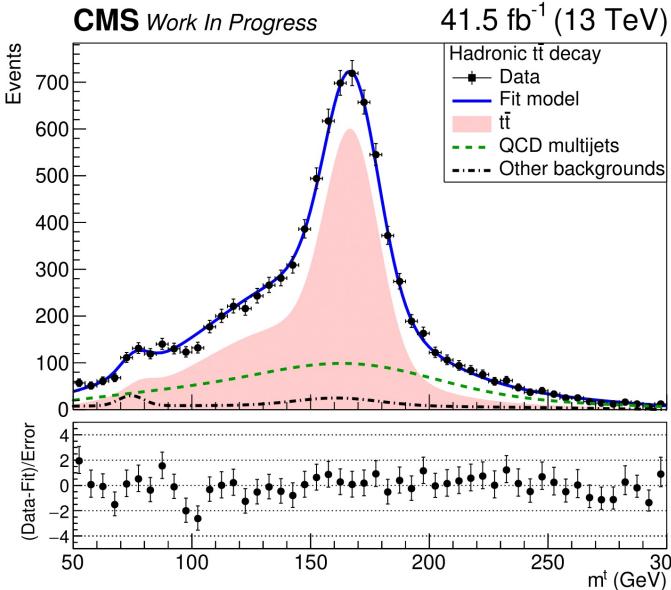


Mass Fit in SR_A

2016_preVFP
 $r = 0.691 \pm 0.028$



2016_postVFP
 $r = 0.640 \pm 0.029$



2018
 $r = 0.675 \pm 0.016$



Motivation

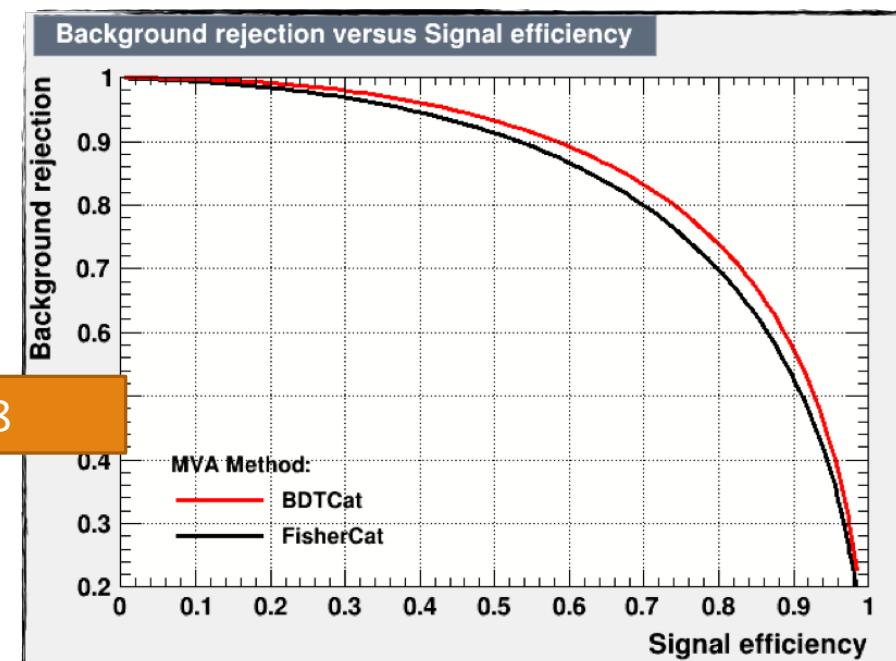
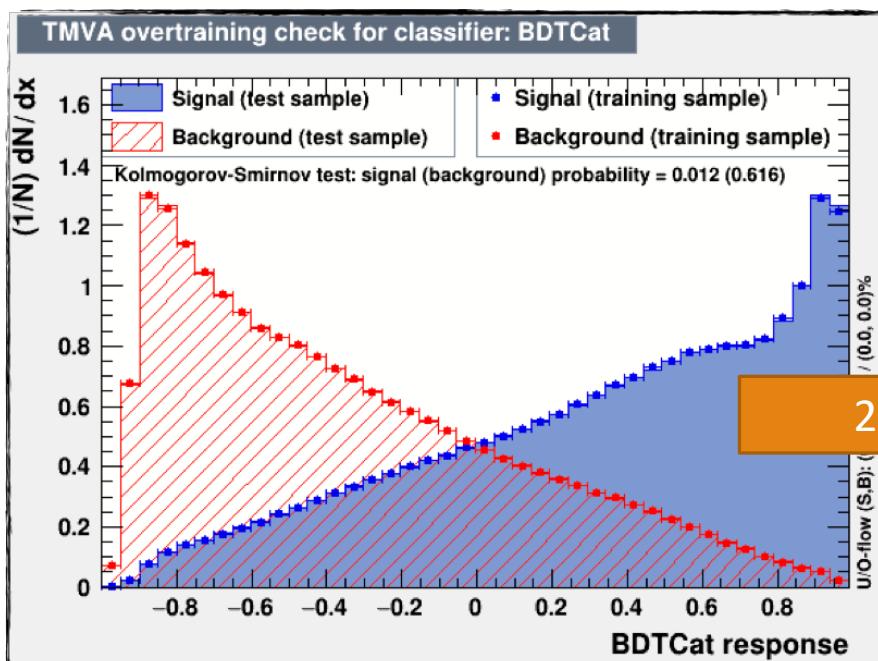
- The main background for this analysis is QCD
- A data driven method is used for subtracting it
- The method relies on the assumption that by inverting the b-tagging requirement in the signal region (SR) we can have the shape of the QCD contribution
- The tagger is required **not to** use b-tagging information and the full range of the subjet mass
 - An in-house BDT was developed to overcome this limitation
 - DeepAK8 uses b-tagging
 - Older top taggers use jet mass cuts

CMS top tagger $R = 0.8$	HEPTOPTAGGER $R = 1.5$	OptimalR $R = 0.5-1.5$
$\delta_p > 0.05$	$f_{\text{drop}} = 0.8$	$m_{23}/m_{123} > 0.35$
$A = 0.0004$	$m_{\text{cut}} = 30 \text{ GeV}$	$0.2 < \arctan \frac{m_{13}}{m_{12}} < 1.3$
$N_{\text{sub}} \geq 3$	$R_{\text{filt}}^{\max} = 0.3$	$f_W = 0.15$
$m_{\min} > 50 \text{ GeV}$	$N_{\text{filt}} = 5$	$140 < m_{123} < 220 \text{ GeV}$
$140 < m_{\text{jet}} < 220 \text{ GeV}$	$p_{T,\text{sub}} > 30 \text{ GeV}$	$p_{T,\text{sub}} > 30 \text{ GeV}$

BDT Output

- In house developed top tagger, for top candidate jets
 - BDT based
 - Input variables:
 - N-subjetness: τ_1, τ_2, τ_3
 - Energy correlation functions (ECF) ECFB1N2, ECFB1N3, ECFB2N2, ECFB2N3
 - Soft drop mass of the leading and subleading subjets
 - Fraction of the jet over the of all the jets in the event
 - **No b-tagging requirements**
- Phase space split in categories based on the pt of the jet:
 - [400, 600) GeV
 - [600, 800) GeV
 - [800, 1200) GeV
 - [1200, Inf) GeV
- Different training and working point for each year (Signal) and QCD (Bkg) samples used in the training

The use of DeepAK8 was investigated but it uses b-tagging so it is not applicable in our use case



Overview

- BDT Input and Output in the SR_B Region
 - SR_B: Baseline selection + tight Mass Cut (120,220) GeV, no TopTagger Selection
 - Leading + subleading in different pT regions:
 - [400,600], [600,800], [800, Inf)
 - [400,500], [500,600], [600, Inf)
 - Find Data vs MC Input and Output for UL our Analysis [here](#)
- Top Tagger Scale Factors
 - Data is subtracted QCD and Subdominant bkg (MC) so that the data sample is pure

$$\text{efficiency} = \frac{\text{Tight \& SR}}{\text{Tight \& Probe}} = \frac{\# (\text{1 jet pass baseline + Tight TopTagger Cut AND 1 jet pass SR})}{\# (\text{1 jet pass baseline + Tight TopTagger Cut AND 1 jet pass only baseline})}$$

- Implemented Randomization (check random jet) to fill histogram to avoid pT bias
- Divide the phase space into pT regions: [400-600]GeV, [600-800]GeV, [800-Inf]GeV
- For the QCD estimation, we perform a fit in both regions (Tight & Probe, Tight & SR):
 - Shape of QCD is estimated from Data while inverting btagging requirement
 - # QCD events in each region is calculated from fit using the Leading JetMassSoftDrop variable
 - To scale the ttbar → fit the Leading JetMassSoftDrop in each region and get the signal strength
 - For the evaluation of Signal distribution from data, we do the following:

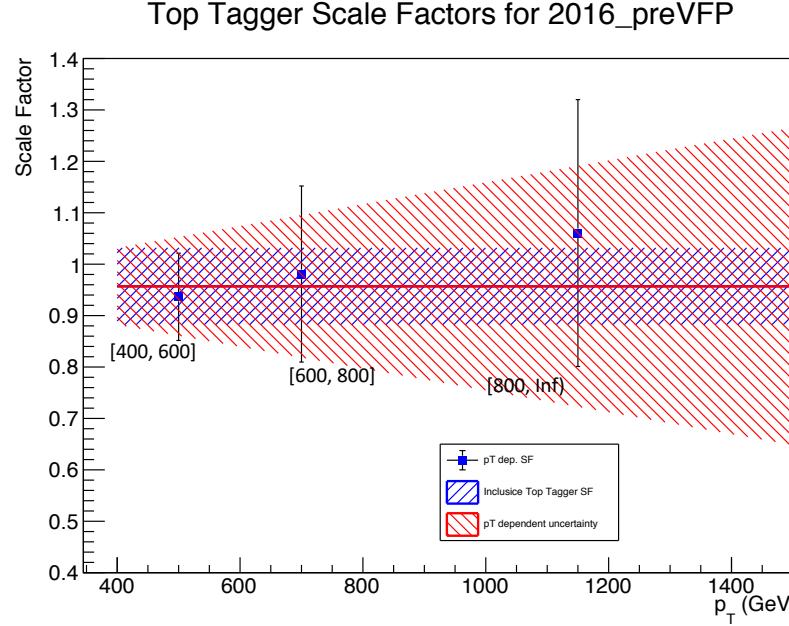
$$\forall \text{region}: S(x) = D(x) - N_{QCD} d_0(x) - \text{Sub. Bkg}(x)$$

Fraction of events used in the cross section measurement that are also used in the Top Tagger SF measurement is of the order of 35%

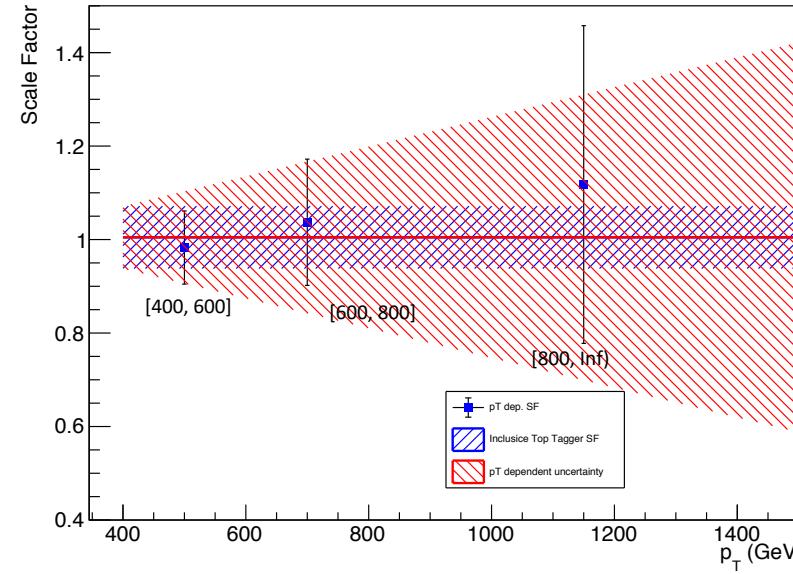


Scale Factors

2016 preVFP

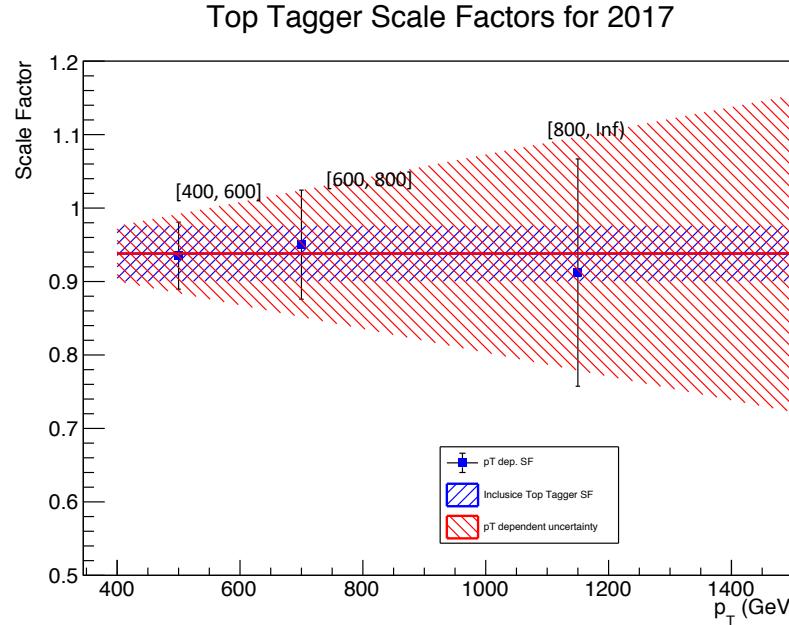


Top Tagger Scale Factors for 2016_postVFP

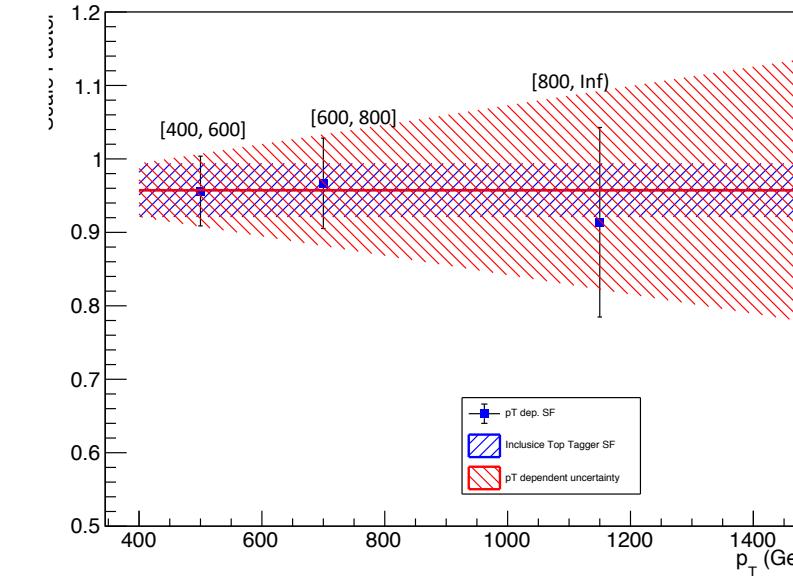


2016 postVFP

2017



Top Tagger Scale Factors for 2018

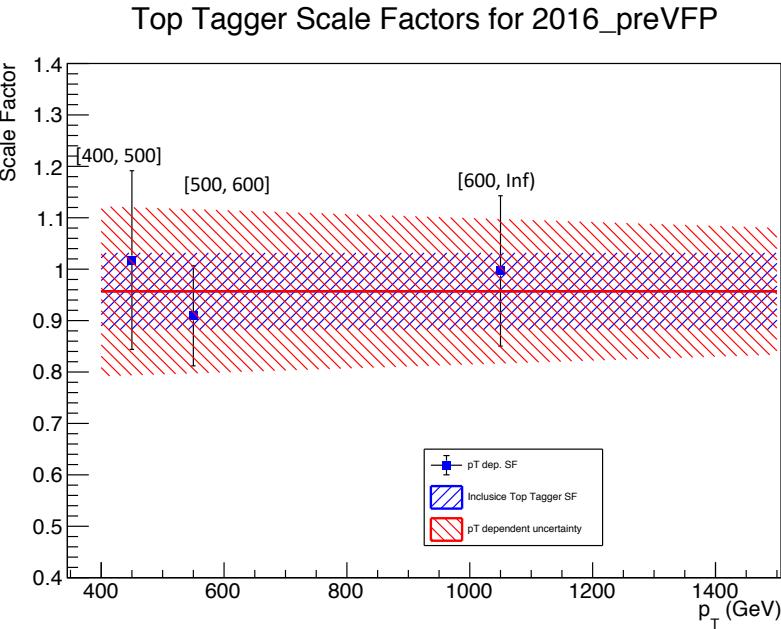


2018

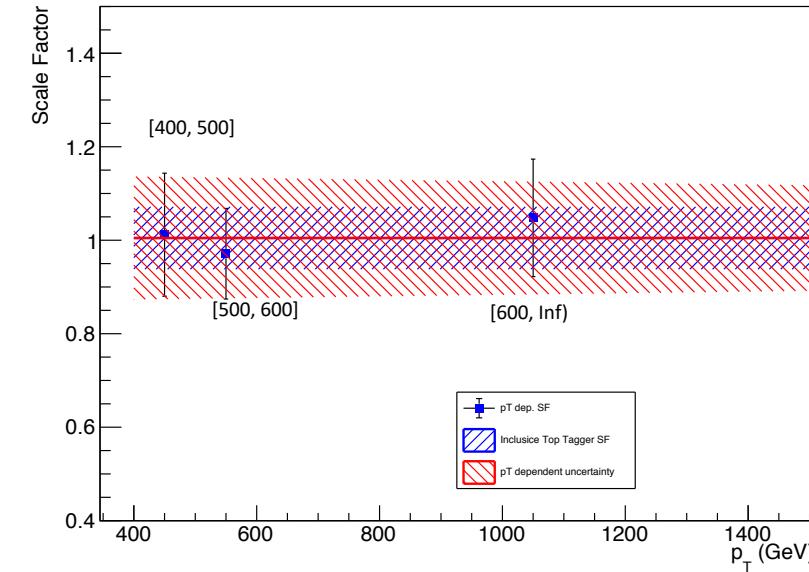


Scale Factors

2016 preVFP

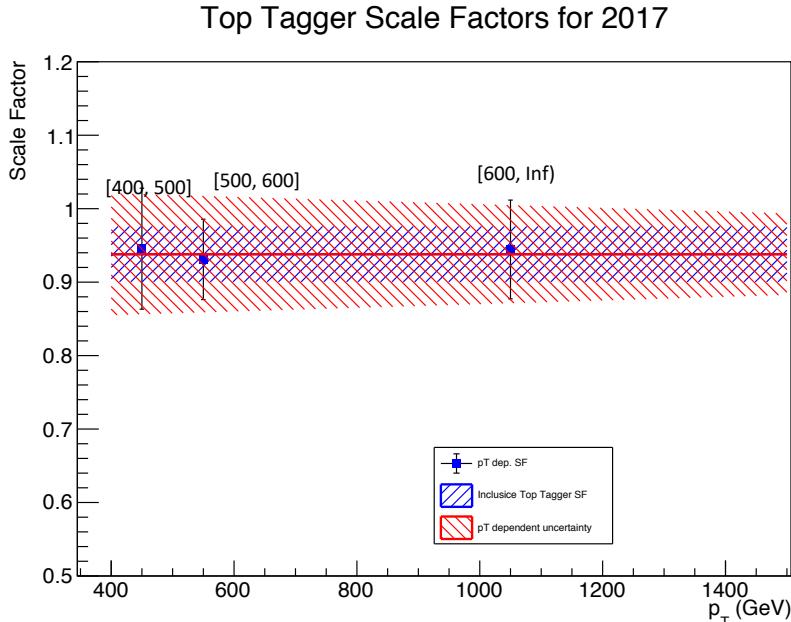


Top Tagger Scale Factors for 2016_postVFP

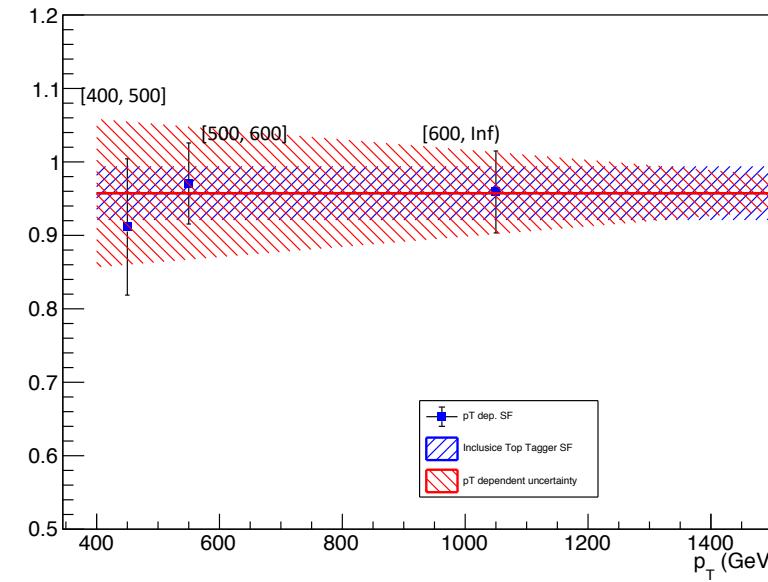


2016 postVFP

2017



Top Tagger Scale Factors for 2018



2018



Systematic uncertainties list

- Jet energy scale
- Jet energy resolution
- B-tagging
- Luminosity
- PDF variations
- ISR, FSR
- Renormalization and factorization scales
- Parton shower (α_s)
- Parton shower (hdamp)
- Tune variations, cp5up, cp5down

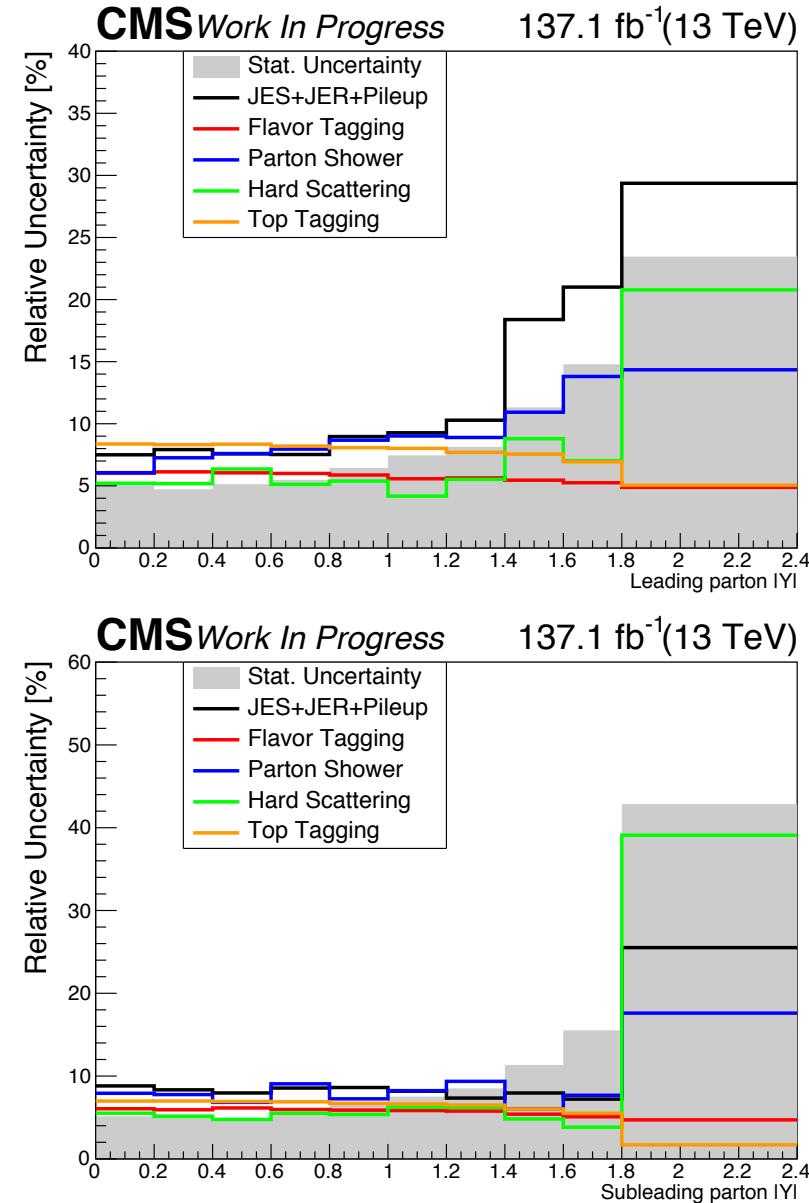
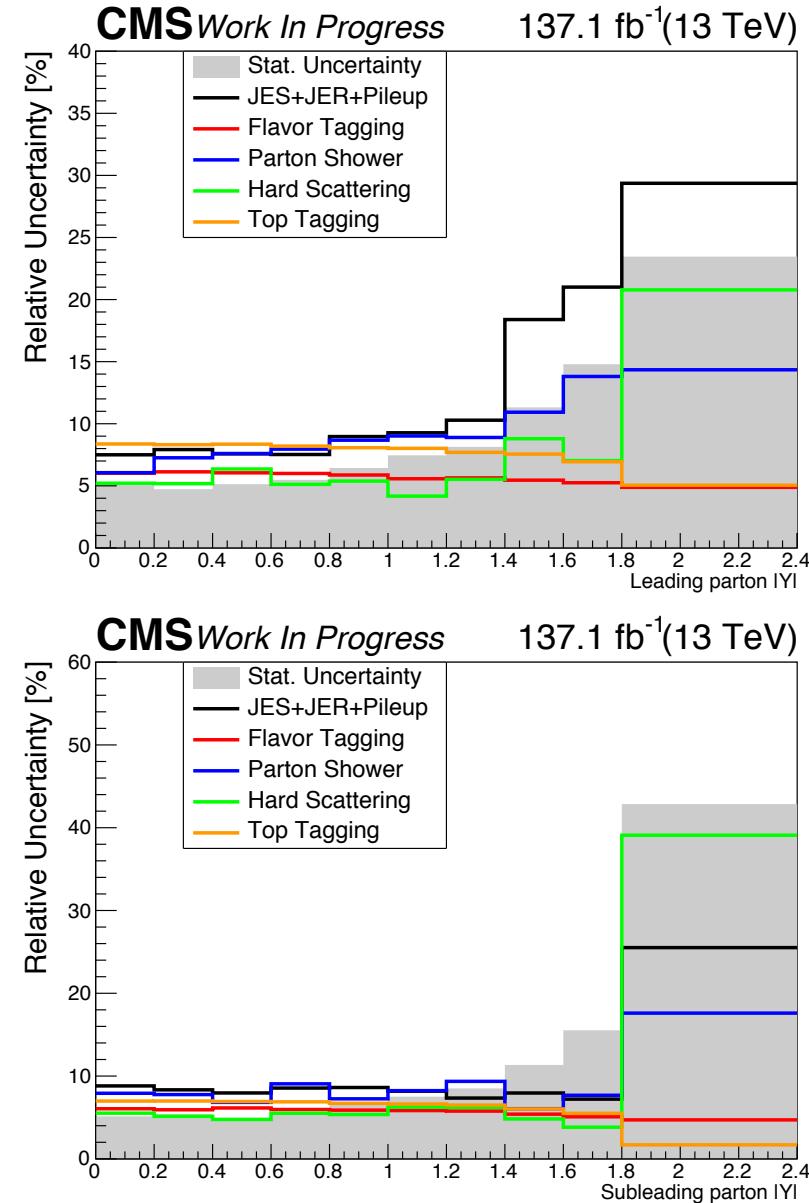
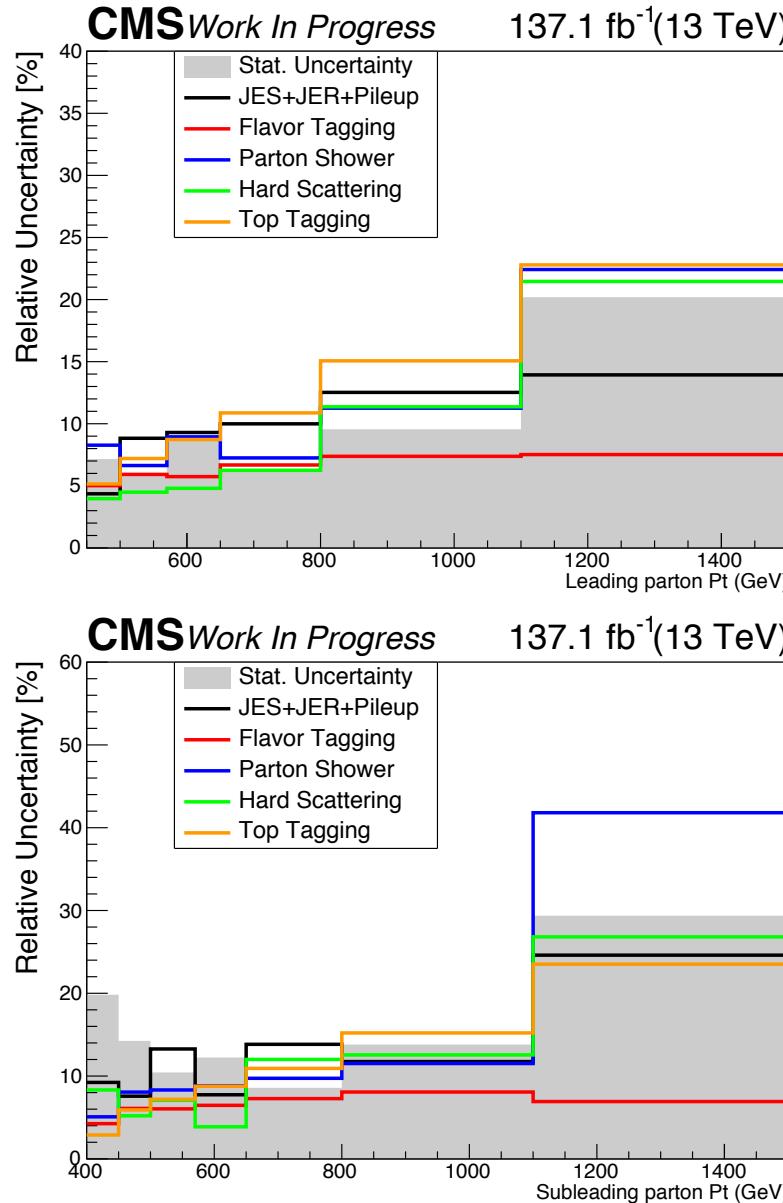


Combination

- Perform the signal extraction for each individual year
- The signal extraction is performed for the nominal samples and for each variation separately
- For the nominal sample and each variation, add the separate years to form a common signal
- Correlations used in the addition of the uncertainties based on the corresponding recommendations from each group
- Calculate combined response matrix
- Calculate combined efficiency and acceptance
- Unfold and calculate cross section for the nominal and each variation separately
- Calculate the total systematic uncertainty to produce the final measurement



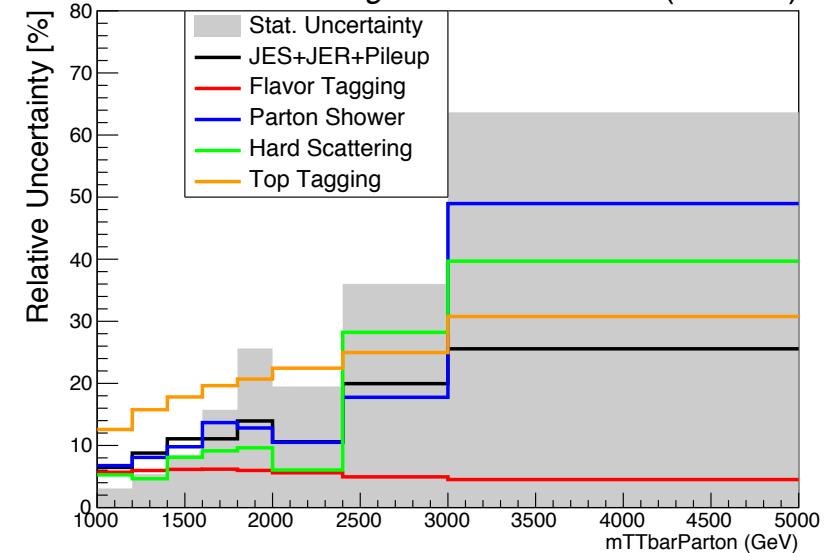
Systematic Uncertainties Breakdown



Systematic Uncertainties Breakdown

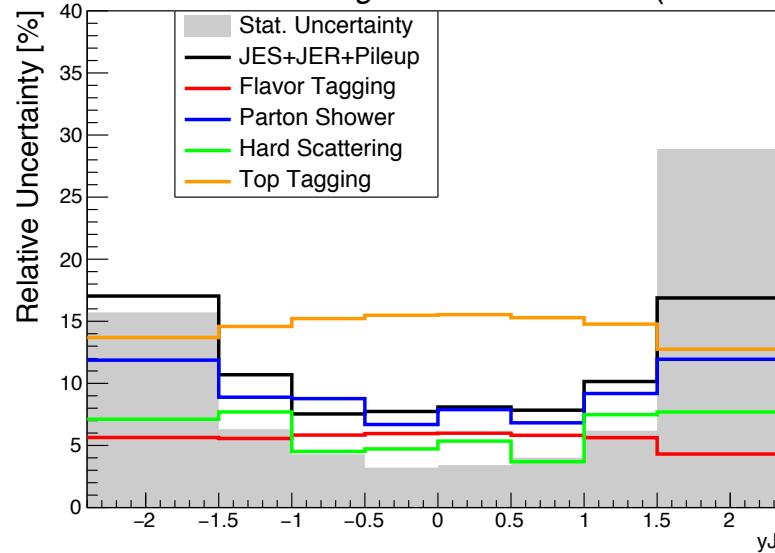
CMS Work In Progress

137.1 fb^{-1} (13 TeV)



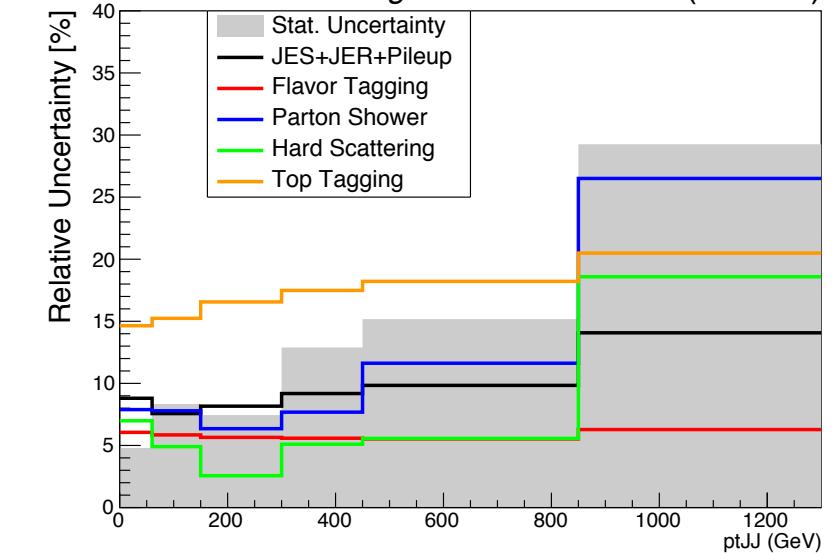
CMS Work In Progress

137.1 fb^{-1} (13 TeV)



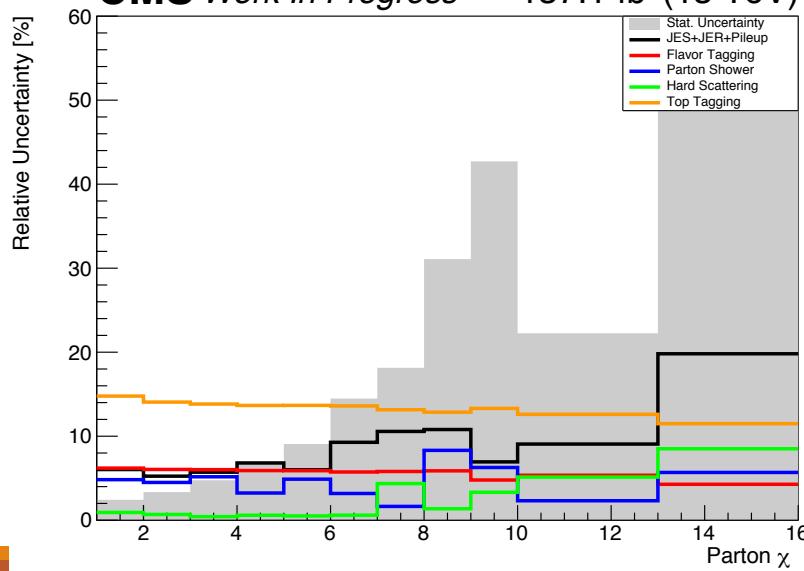
CMS Work In Progress

137.1 fb^{-1} (13 TeV)



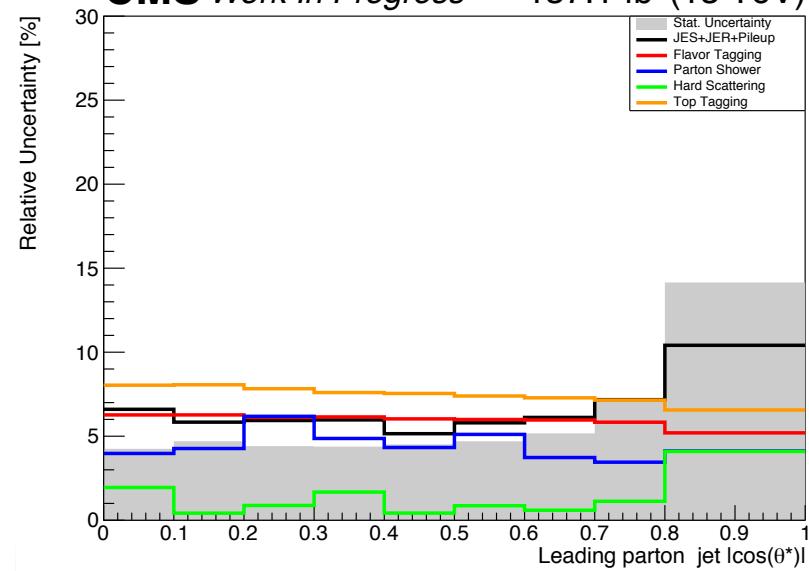
CMS Work In Progress

137.1 fb^{-1} (13 TeV)



CMS Work In Progress

137.1 fb^{-1} (13 TeV)

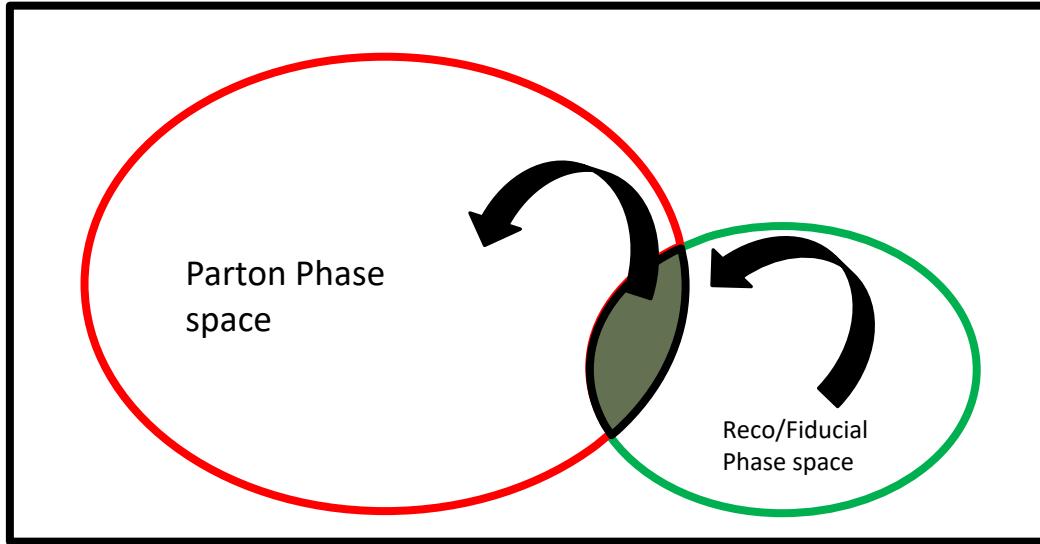


Parton & Particle levels

Particle level

Parton

Observable	Requirement
p_T^1	> 450 GeV
p_T^2	> 400 GeV
$ \eta^{t,\bar{t}} $	< 2.4
$m_{t\bar{t}}$	> 1000 GeV



Observable	Requirement
N_{jets}	>1
p_T^{jet1}	> 450 GeV
p_T^{jet2}	> 450 GeV
$ \eta^{jet1,2} $	< 2.4
$m_{SD}^{jet1,2}$	(120, 220) GeV
m_{jj}	> 1000 GeV

$$\frac{d\sigma_i^{\text{unf}}}{dx} = \frac{1}{\mathcal{L} \cdot \Delta x_i} \cdot \frac{1}{f_{2,i}} \cdot \sum_j \left(R_{ij}^{-1} \cdot f_{1,j} \cdot S_j \right)$$

true efficiency of the
reco+true selection

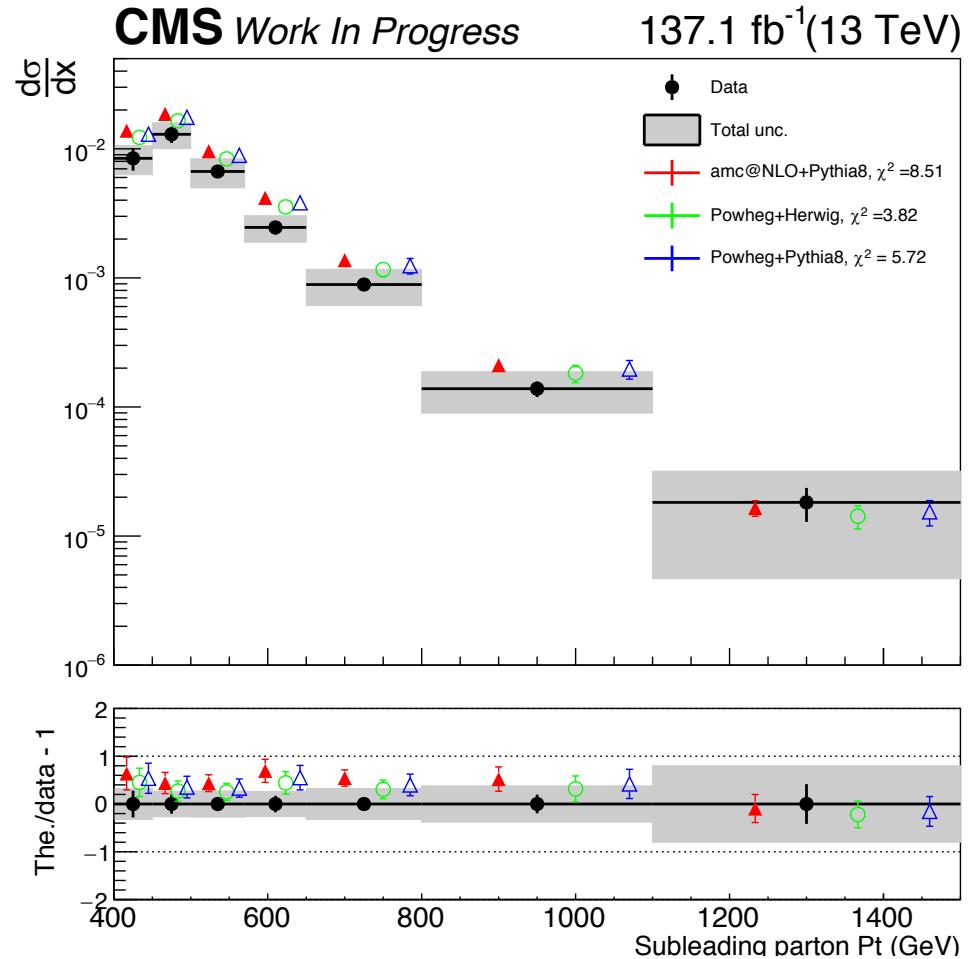
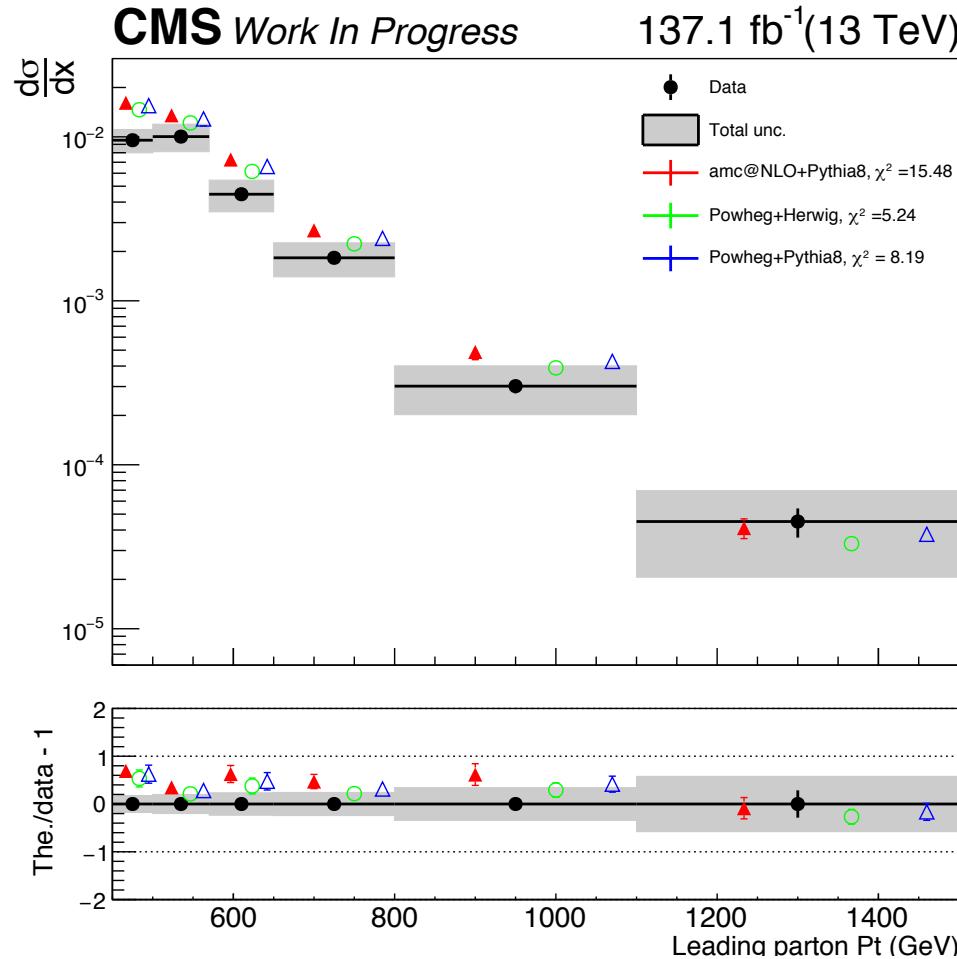
migration matrix

reco efficiency of the
reco+true selection

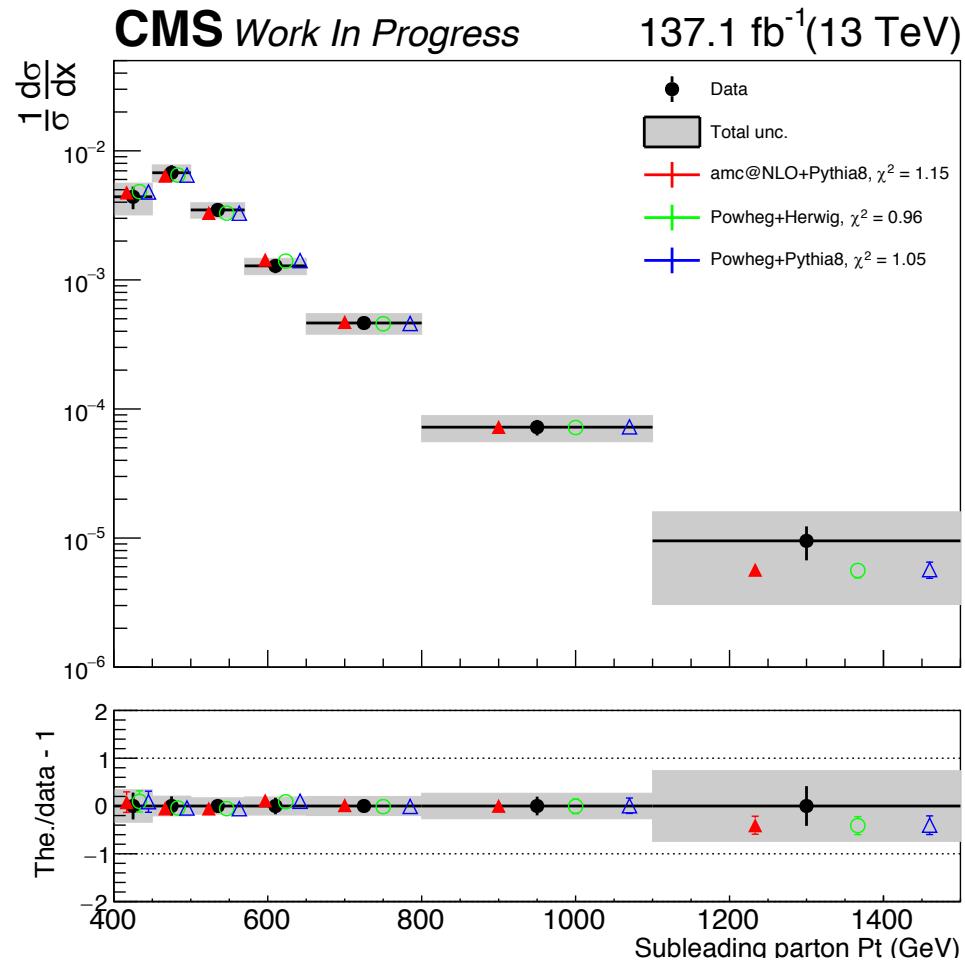
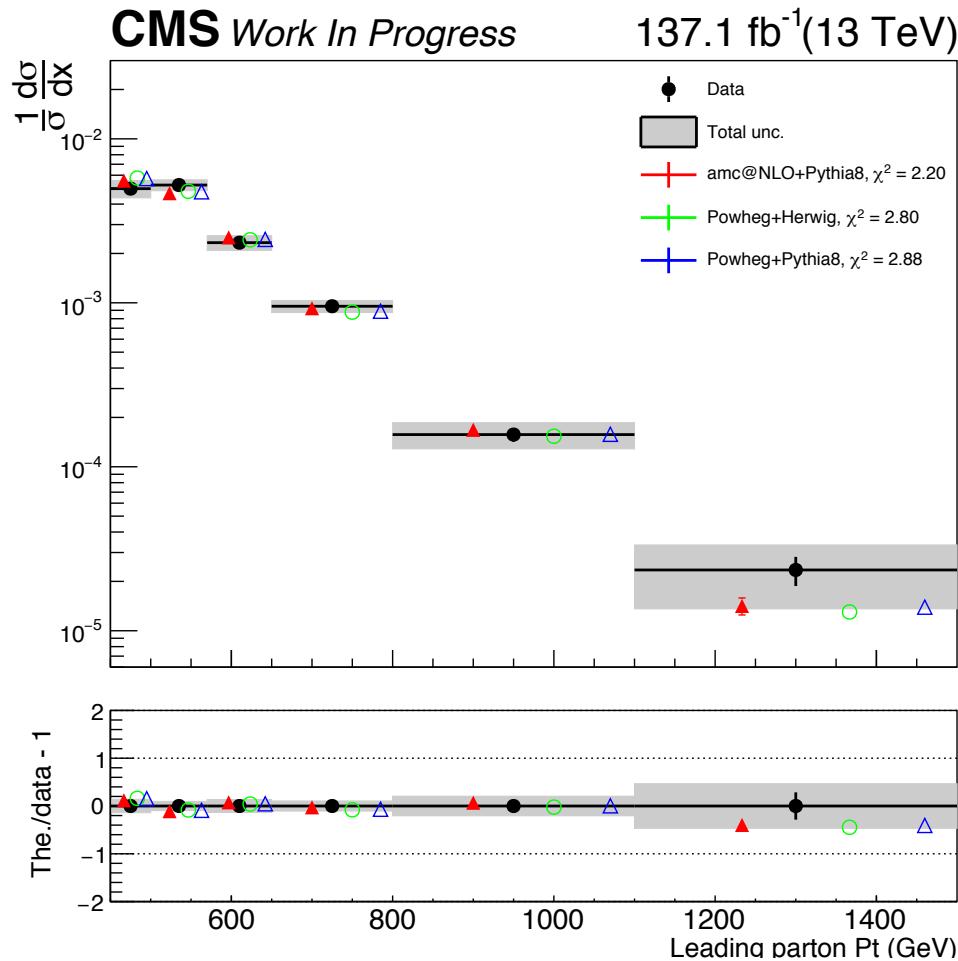
Unfolding: simple response matrix inversion w/o regularisation



Final Results (absolute)



Final Results (normalized)

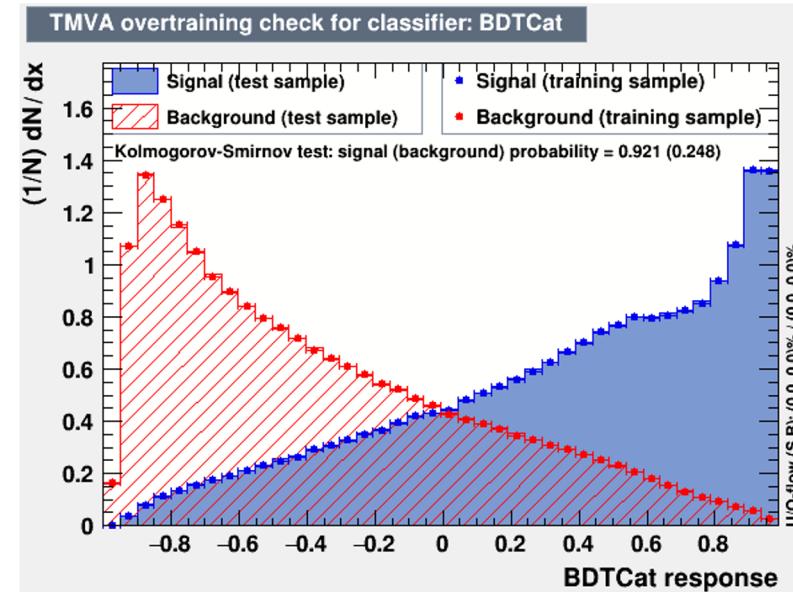
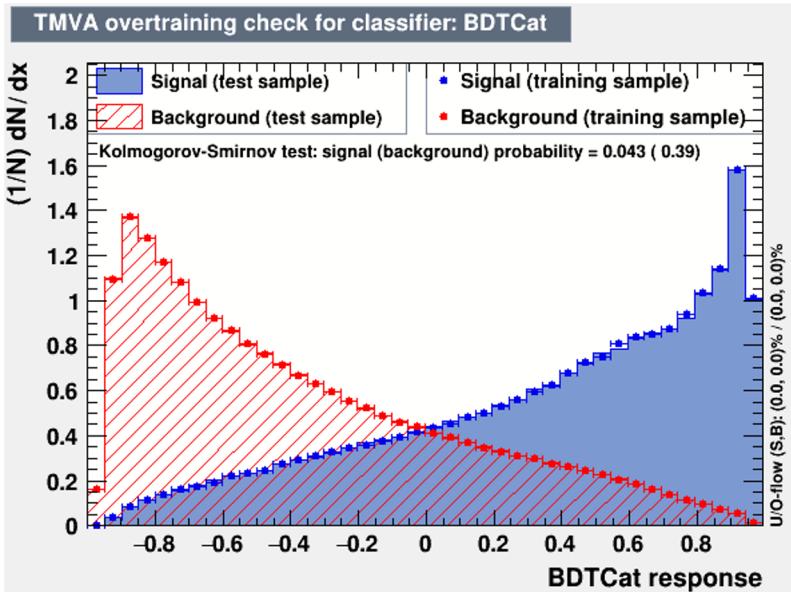


Backup



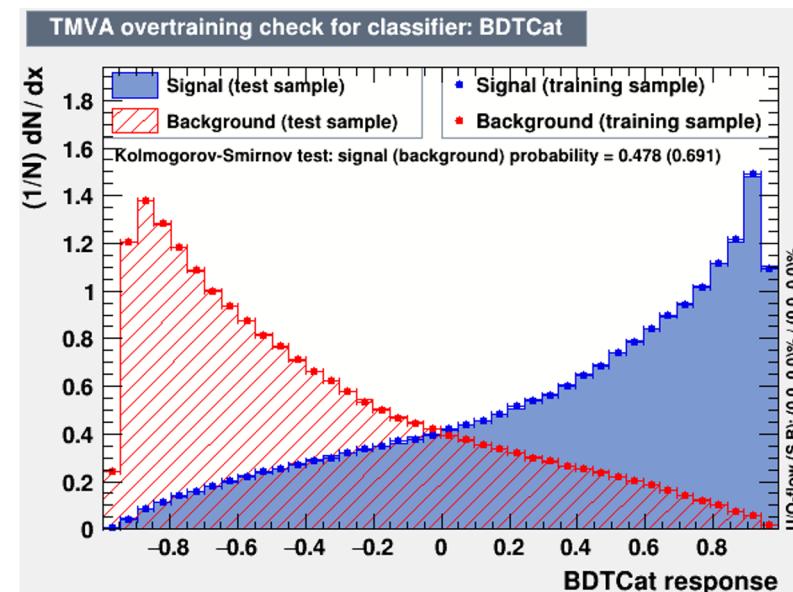
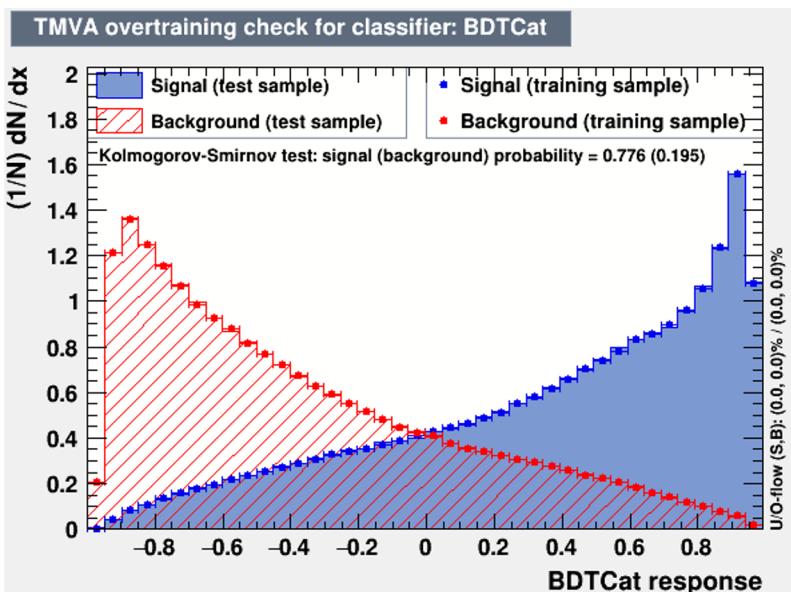
Top Tagger

2016_preVFP



2016_postVFP

2017

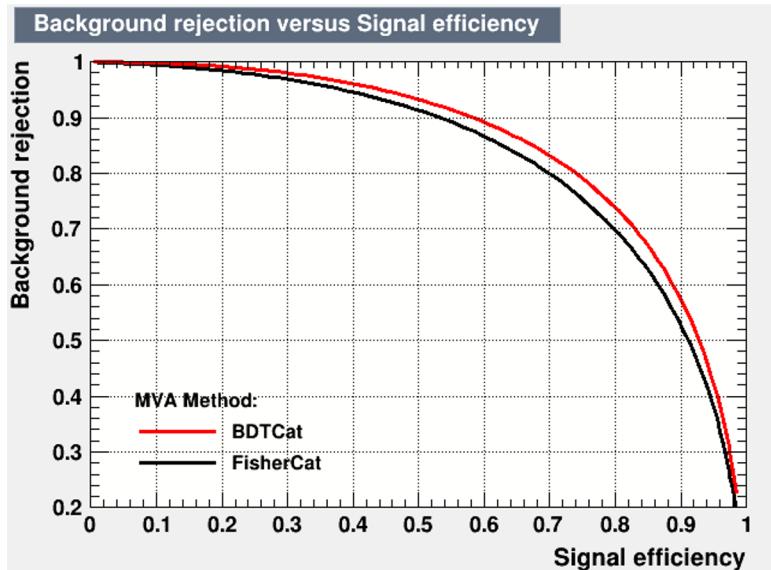


2018

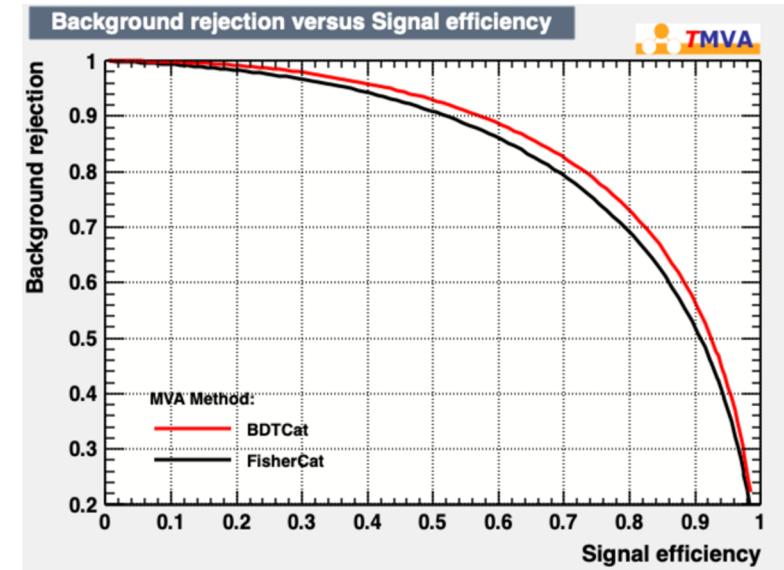


Top Tagger

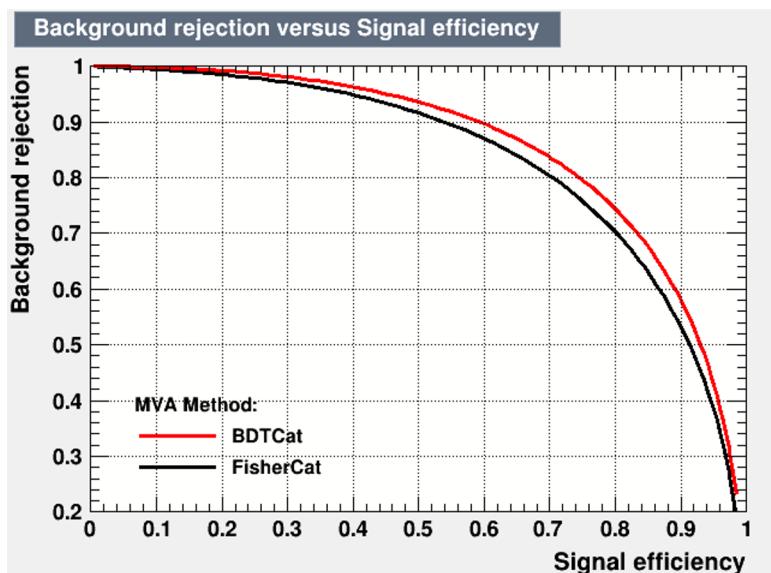
2016_preVFP



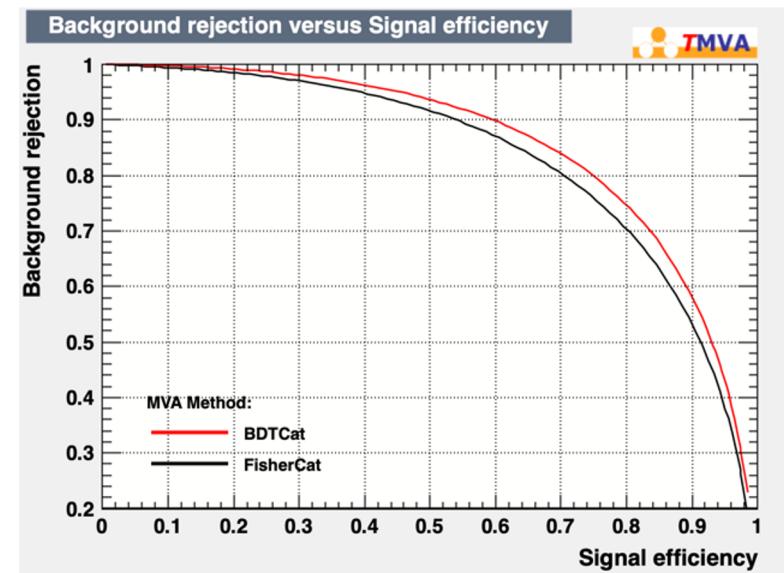
2016_postVFP



2017

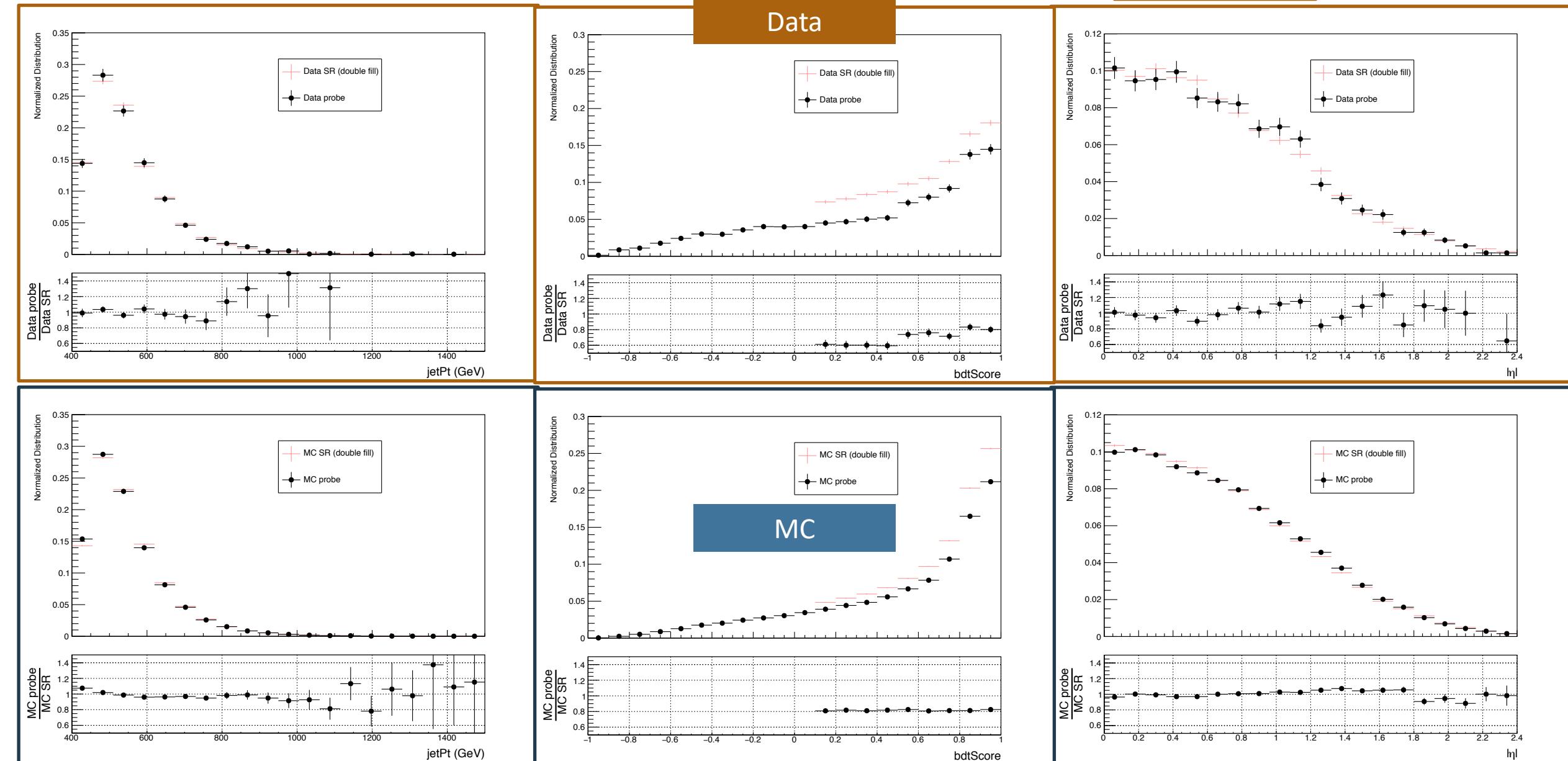


2018



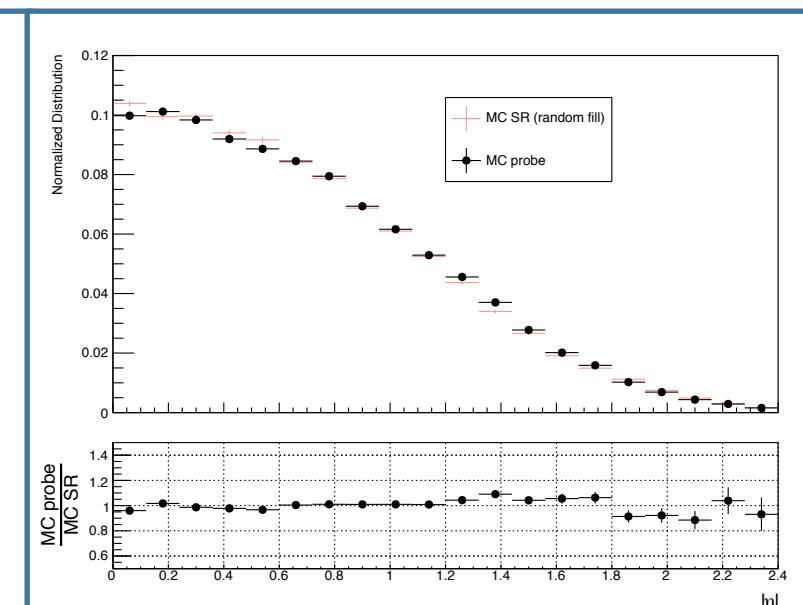
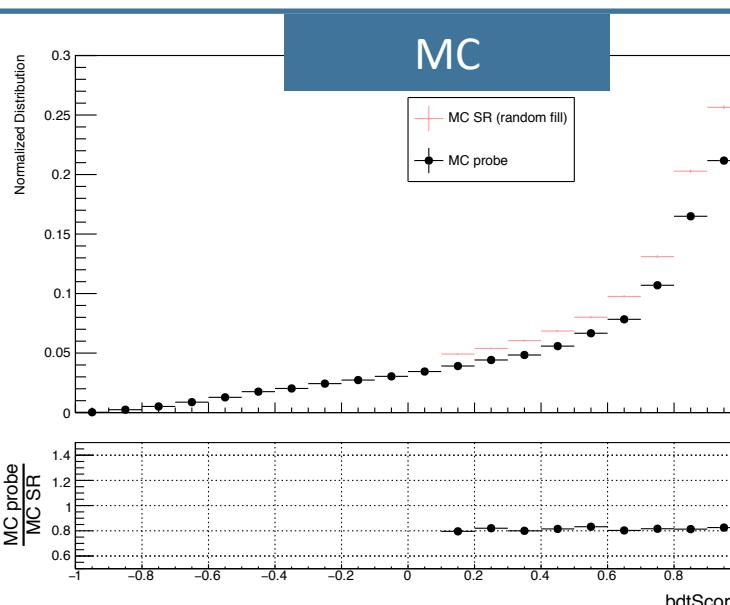
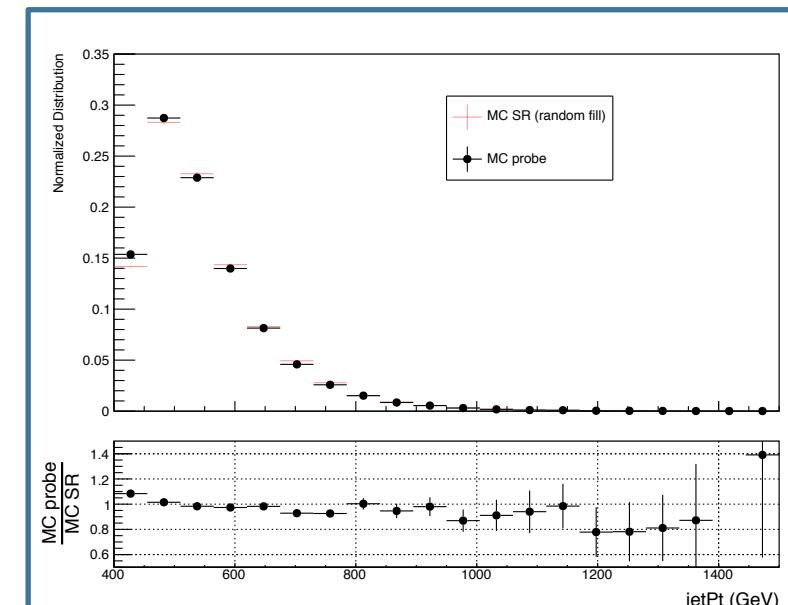
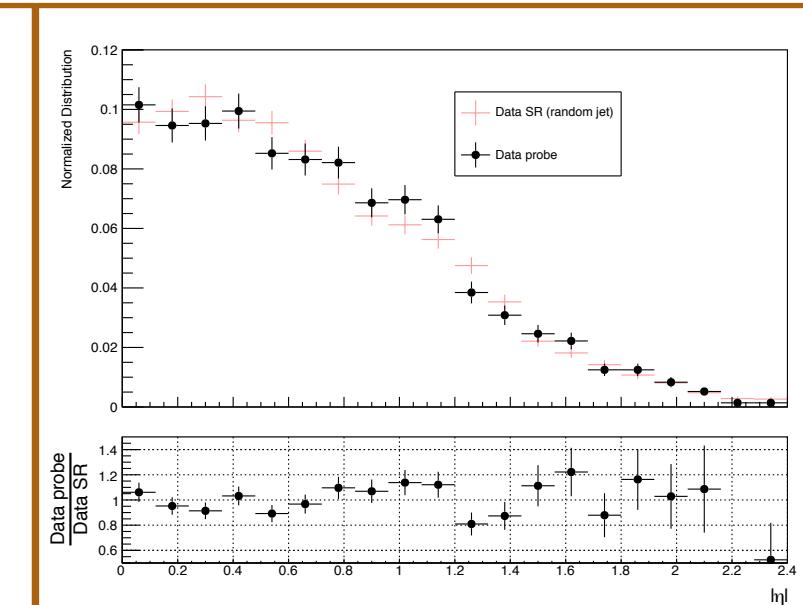
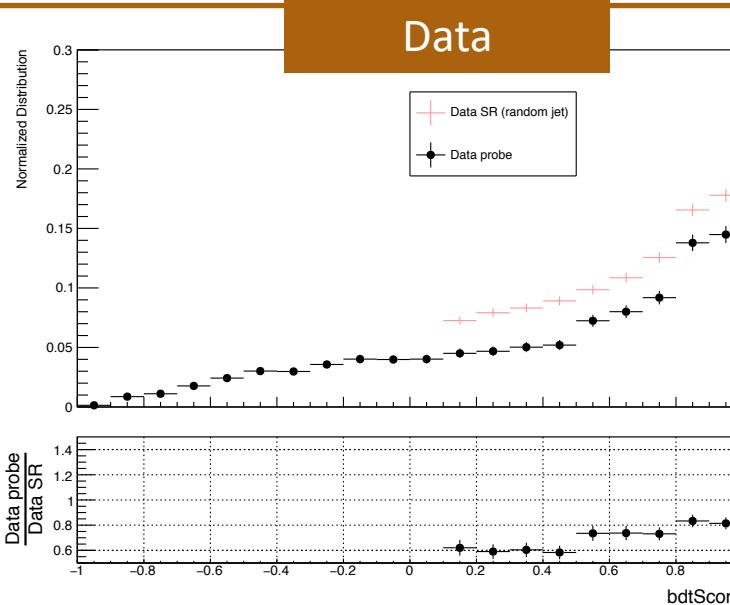
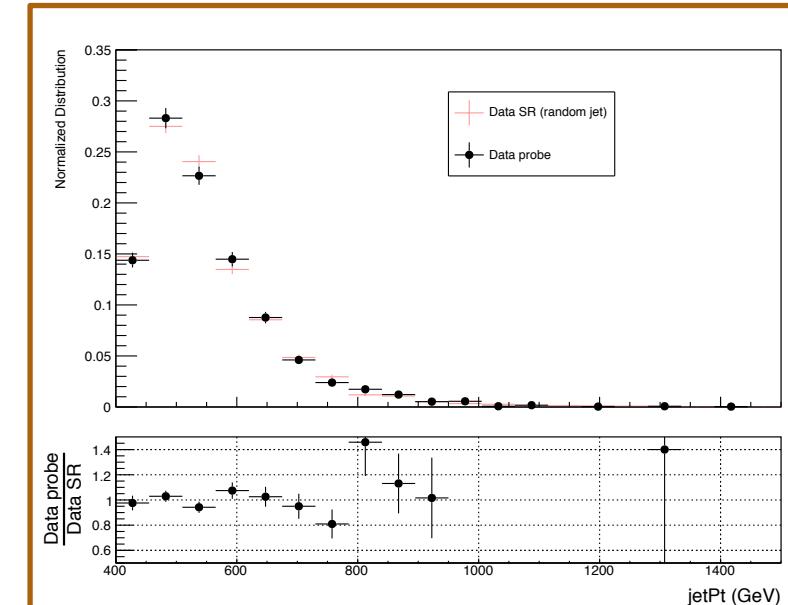
Shape Comparison of the probe jet and the jets used for the measurement (double Filled)

2018



Shape Comparison of the probe jet and the jets used for the measurement

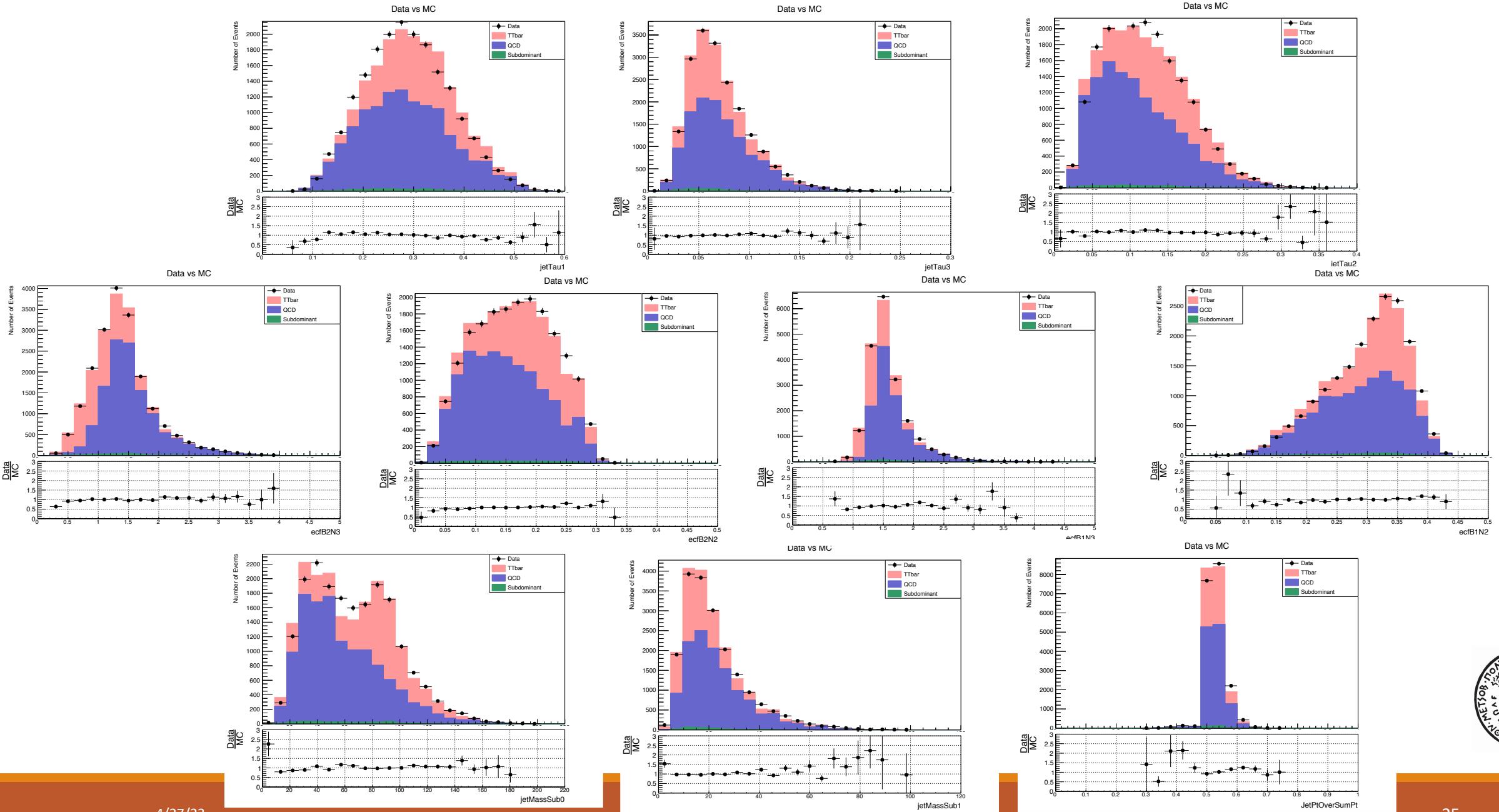
2018



Top Tagger Input Variables

2018

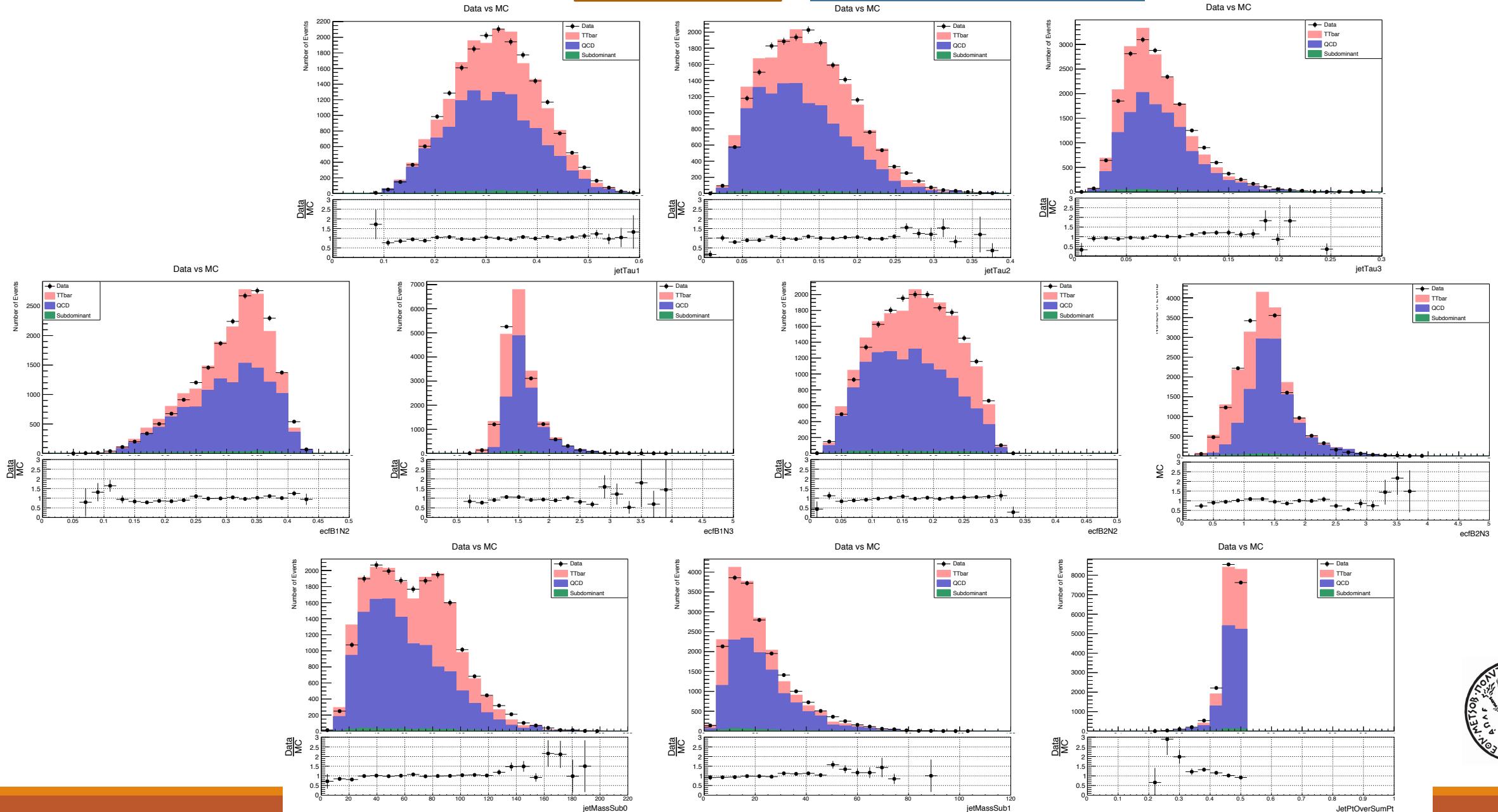
Leading jet



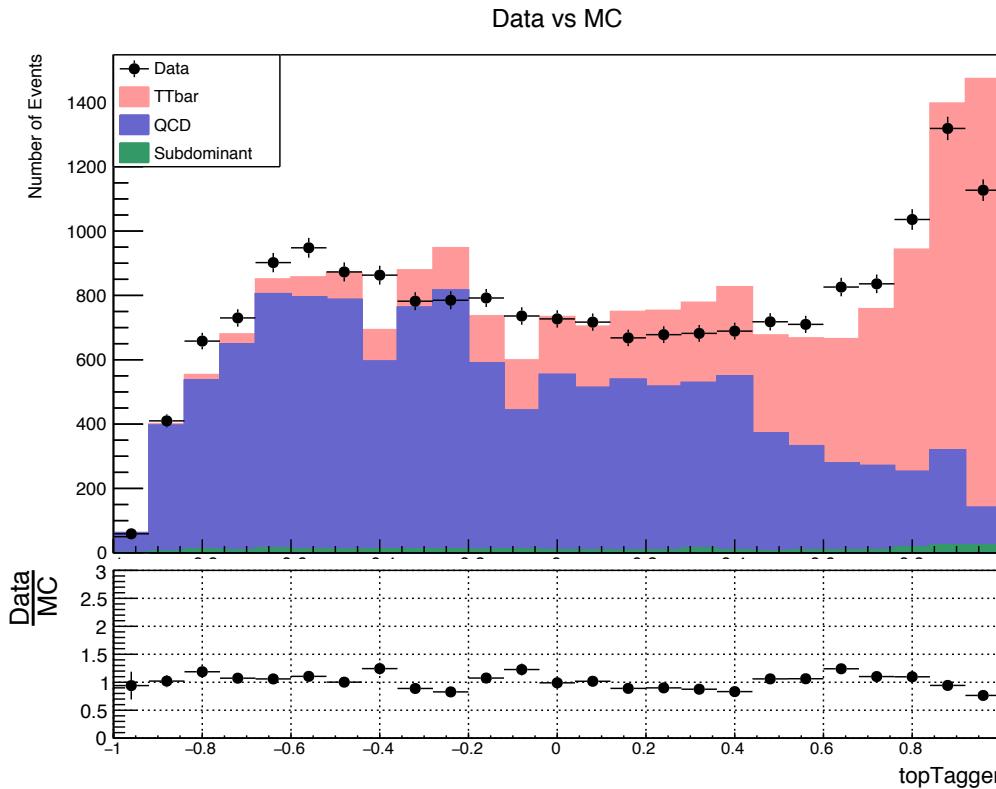
Top Tagger Input Variables

2018

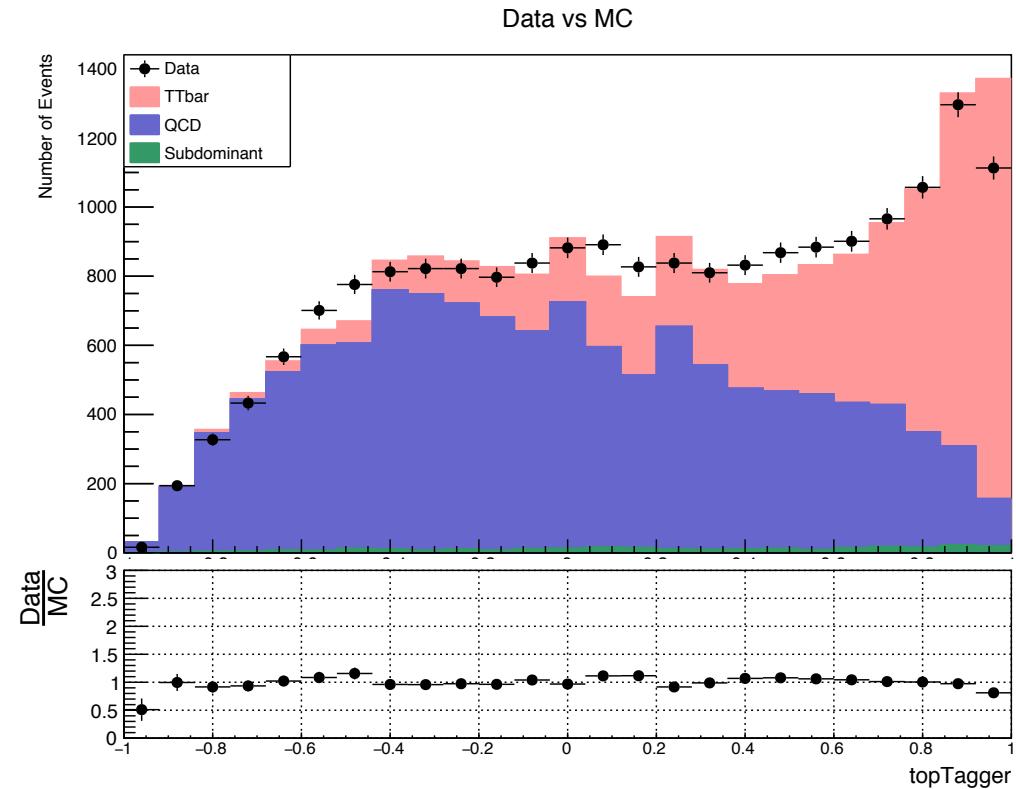
Second Leading Jet



Leading Jet



Second Leading Jet



Top Tagger Efficiencies

Table 27: Top Tagger efficiency Values for 2016 preVFP.

Eff. Type	Eff. Data (stat)	Eff. $t\bar{t}$ (stat)	Eff. $t\bar{t}$ (stat + systematic)
Inclusive	0.757 ± 0.058	0.791 ± 0.009	0.791 ± 0.01
p_T [400, 600]GeV	0.742 ± 0.067	0.793 ± 0.011	0.793 ± 0.021
p_T [600, 800]GeV	0.774 ± 0.134	0.79 ± 0.016	0.79 ± 0.025
p_T [800, Inf)GeV	0.824 ± 0.198	0.777 ± 0.037	0.777 ± 0.047

Table 29: Top Tagger efficiency Values for 2016 postVFP.

Eff. Type	Eff. Data (stat)	Eff. $t\bar{t}$ (stat)	Eff. $t\bar{t}$ (stat + systematic)
Inclusive	0.79 ± 0.052	0.786 ± 0.008	0.786 ± 0.011
p_T [400, 600]GeV	0.776 ± 0.061	0.79 ± 0.01	0.79 ± 0.021
p_T [600, 800]GeV	0.81 ± 0.104	0.781 ± 0.015	0.781 ± 0.024
p_T [800, Inf)GeV	0.861 ± 0.259	0.77 ± 0.035	0.77 ± 0.046

Table 31: Top Tagger efficiency Values for 2017.

Eff. Type	Eff. Data (stat)	Eff. $t\bar{t}$ (stat)	Eff. $t\bar{t}$ (stat + systematic)
Inclusive	0.814 ± 0.032	0.868 ± 0.006	0.868 ± 0.009
p_T [400, 600]GeV	0.81 ± 0.04	0.867 ± 0.008	0.867 ± 0.017
p_T [600, 800]GeV	0.827 ± 0.063	0.871 ± 0.012	0.871 ± 0.021
p_T [800, Inf)GeV	0.793 ± 0.132	0.869 ± 0.029	0.869 ± 0.037

Table 33: Top Tagger efficiency Values for 2018.

Eff. Type	Eff. Data (stat)	Eff. $t\bar{t}$ (stat)	Eff. $t\bar{t}$ (stat + systematic)
Inclusive	0.792 ± 0.03	0.827 ± 0.005	0.827 ± 0.008
p_T [400, 600]GeV	0.789 ± 0.039	0.825 ± 0.006	0.825 ± 0.014
p_T [600, 800]GeV	0.805 ± 0.051	0.833 ± 0.01	0.833 ± 0.02
p_T [800, Inf)GeV	0.752 ± 0.104	0.822 ± 0.024	0.822 ± 0.037

Table 28: Top Tagger efficiency Values for 2016 preVFP using JMAR proposed p_T regions.

Eff. Type	Eff. Data (stat)	Eff. $t\bar{t}$ (stat)	Eff. $t\bar{t}$ (stat + systematic)
Inclusive	0.757 ± 0.058	0.791 ± 0.009	0.791 ± 0.01
p_T [400, 500]GeV	0.806 ± 0.136	0.792 ± 0.021	0.792 ± 0.031
p_T [500, 600]GeV	0.721 ± 0.076	0.793 ± 0.013	0.793 ± 0.022
p_T [600, Inf)GeV	0.785 ± 0.114	0.787 ± 0.014	0.787 ± 0.024

Table 30: Top Tagger efficiency Values for 2016 postVFP using JMAR proposed p_T regions.

Eff. Type	Eff. Data (stat)	Eff. $t\bar{t}$ (stat)	Eff. $t\bar{t}$ (stat + systematic)
Inclusive	0.79 ± 0.052	0.786 ± 0.008	0.786 ± 0.011
p_T [400, 500]GeV	0.782 ± 0.1	0.773 ± 0.018	0.773 ± 0.029
p_T [500, 600]GeV	0.774 ± 0.076	0.8 ± 0.012	0.8 ± 0.02
p_T [600, Inf)GeV	0.817 ± 0.097	0.779 ± 0.013	0.779 ± 0.025

Table 32: Top Tagger efficiency Values for 2017 using JMAR proposed p_T regions.

Eff. Type	Eff. Data (stat)	Eff. $t\bar{t}$ (stat)	Eff. $t\bar{t}$ (stat + systematic)
Inclusive	0.814 ± 0.032	0.868 ± 0.006	0.868 ± 0.009
p_T [400, 500]GeV	0.808 ± 0.069	0.854 ± 0.014	0.854 ± 0.023
p_T [500, 600]GeV	0.812 ± 0.047	0.872 ± 0.009	0.872 ± 0.018
p_T [600, Inf)GeV	0.822 ± 0.058	0.870 ± 0.011	0.870 ± 0.019

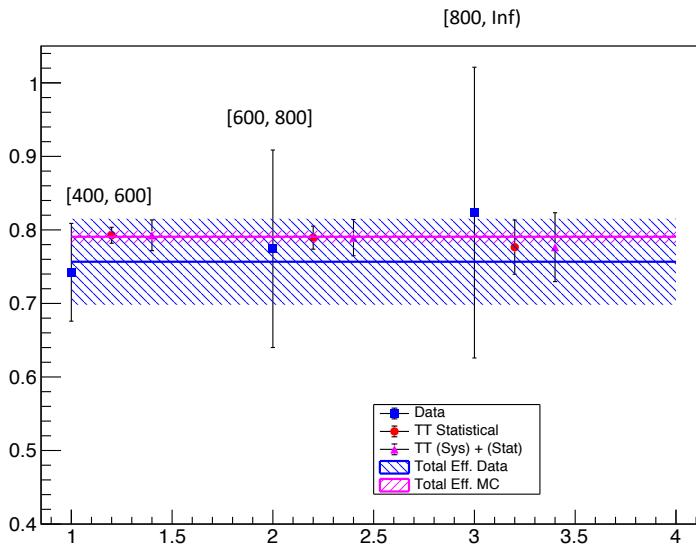
Table 34: Top Tagger efficiency Values for 2018 using JMAR proposed p_T regions.

Eff. Type	Eff. Data (stat)	Eff. $t\bar{t}$ (stat)	Eff. $t\bar{t}$ (stat + systematic)
Inclusive	0.792 ± 0.03	0.827 ± 0.005	0.827 ± 0.008
p_T [400, 500]GeV	0.739 ± 0.074	0.811 ± 0.011	0.811 ± 0.019
p_T [500, 600]GeV	0.807 ± 0.045	0.832 ± 0.007	0.832 ± 0.018
p_T [600, Inf)GeV	0.797 ± 0.046	0.832 ± 0.009	0.832 ± 0.021

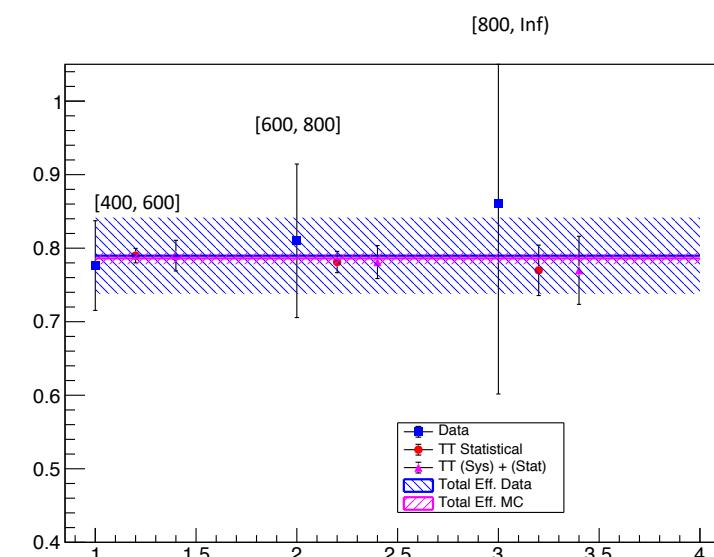


TagAndProbe Efficiency per Pt region

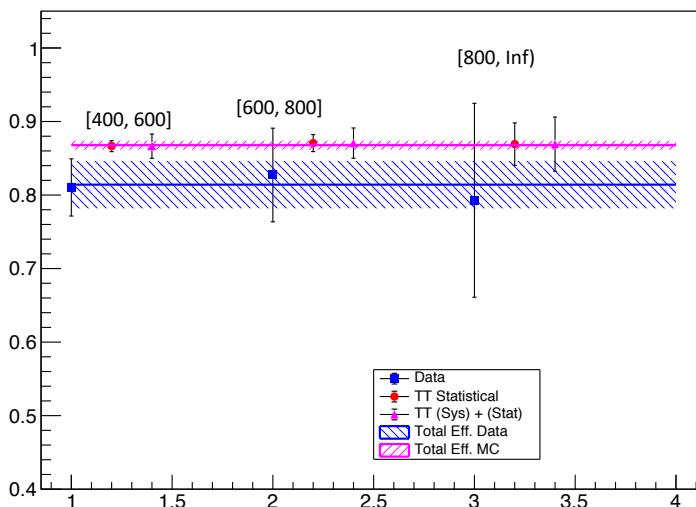
2016 preVFP



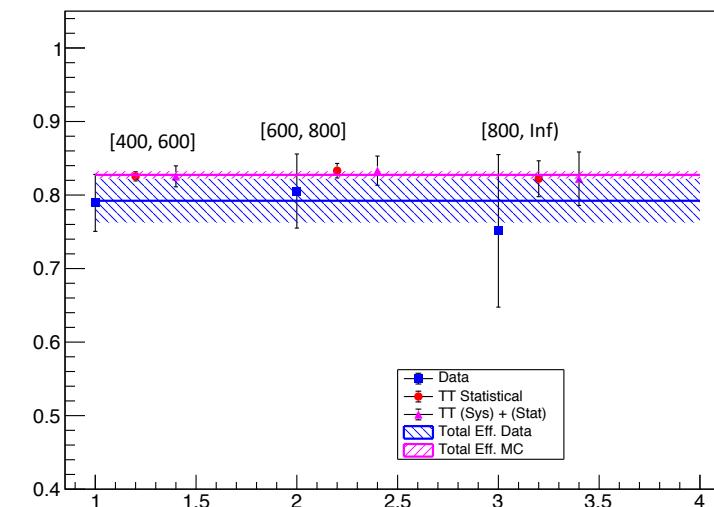
2016 postVFP



2017

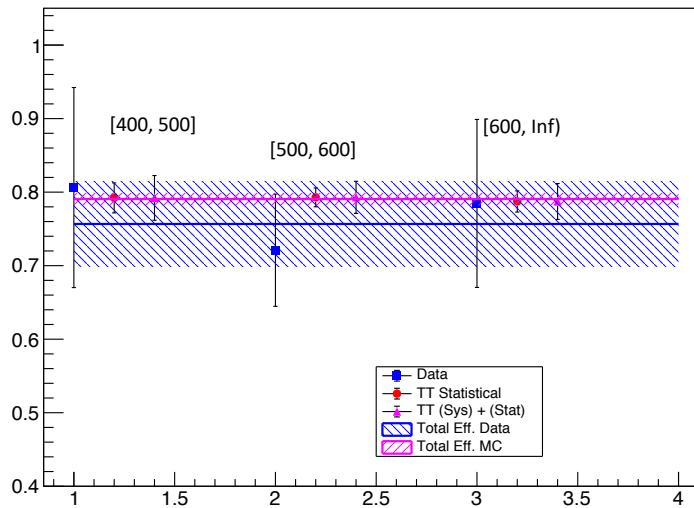


2018

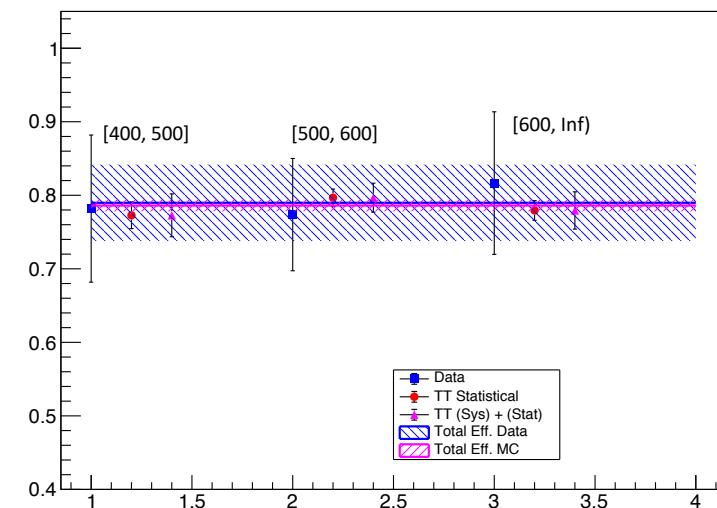


TagAndProbe Efficiency per Pt region (JMAR regions)

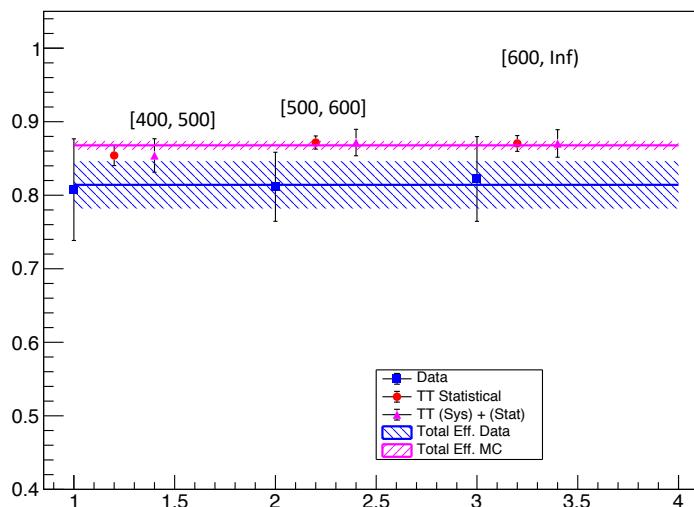
2016 preVFP



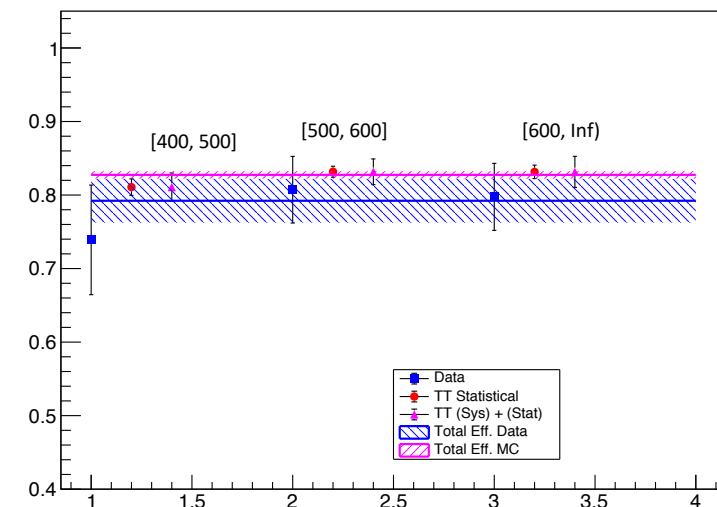
2016 postVFP



2017



2018



Scale Factor Values

Table 35: Top Tagger SF Values for 2016 preVFP.

SF Type	Value \pm error
Inclusive	0.957 \pm 0.074
p_T [400, 600]GeV	0.937 \pm 0.085
p_T [600, 800]GeV	0.981 \pm 0.17
p_T [800, Inf)GeV	1.06 \pm 0.26

Table 36: Top Tagger SF Values for 2016 preVFP using JMAR proposed p_T regions.

SF Type	Value \pm error
Inclusive	0.957 \pm 0.074
p_T [400, 500]GeV	1.02 \pm 0.173
p_T [500, 600]GeV	0.91 \pm 0.097
p_T [600, Inf)GeV	0.997 \pm 0.15

Table 37: Top Tagger SF Values for 2016 postVFP.

SF Type	Value \pm error
Inclusive	1.01 \pm 0.067
p_T [400, 600]GeV	0.983 \pm 0.078
p_T [600, 800]GeV	1.04 \pm 0.135
p_T [800, Inf)GeV	1.12 \pm 0.34

Table 38: Top Tagger SF Values for 2016 postVFP using JMAR proposed p_T regions.

SF Type	Value \pm error
Inclusive	1.01 \pm 0.067
p_T [400, 500]GeV	1.01 \pm 0.132
p_T [500, 600]GeV	0.971 \pm 0.097
p_T [600, Inf)GeV	1.05 \pm 0.13

Table 39: Top Tagger SF Values for 2017.

SF Type	Value \pm error
Inclusive	0.938 \pm 0.038
p_T [400, 600]GeV	0.935 \pm 0.046
p_T [600, 800]GeV	0.95 \pm 0.059
p_T [800, Inf)GeV	0.912 \pm 0.155

Table 40: Top Tagger SF Values for 2017 using JMAR proposed p_T regions.

SF Type	Value \pm error
Inclusive	0.938 \pm 0.038
p_T [400, 500]GeV	0.946 \pm 0.082
p_T [500, 600]GeV	0.931 \pm 0.055
p_T [600, Inf)GeV	0.945 \pm 0.068

Table 41: Top Tagger SF Values for 2018.

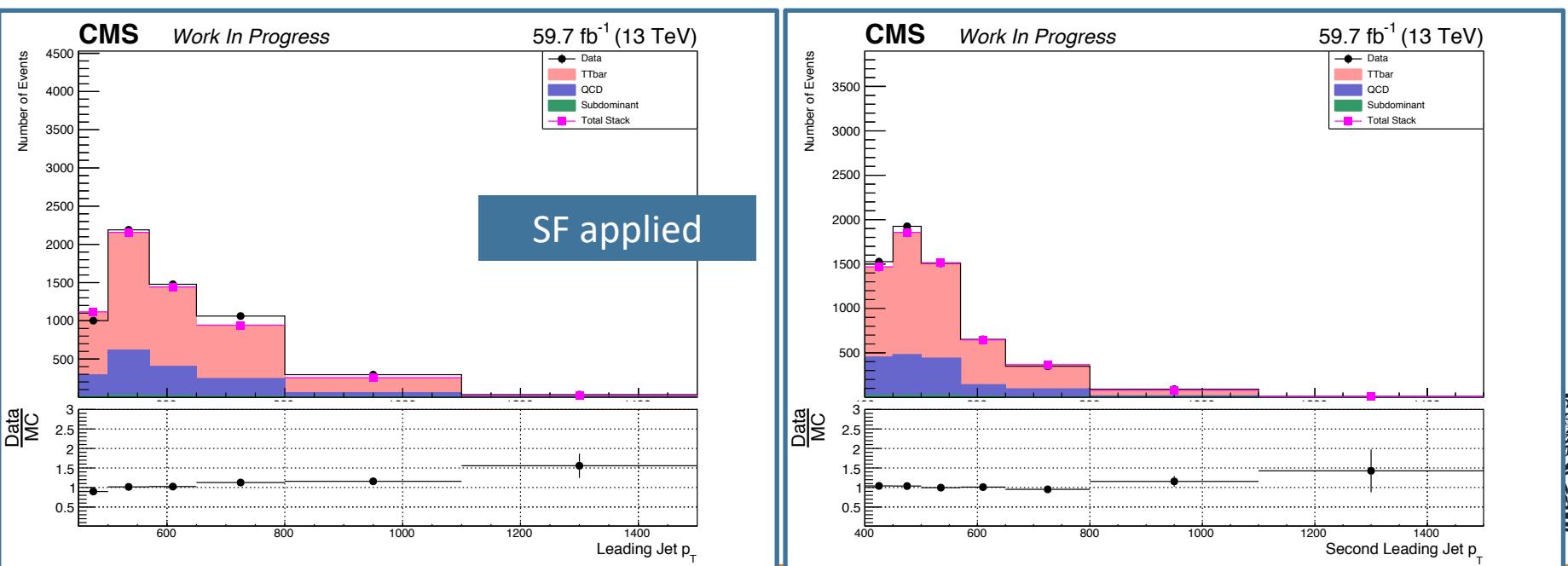
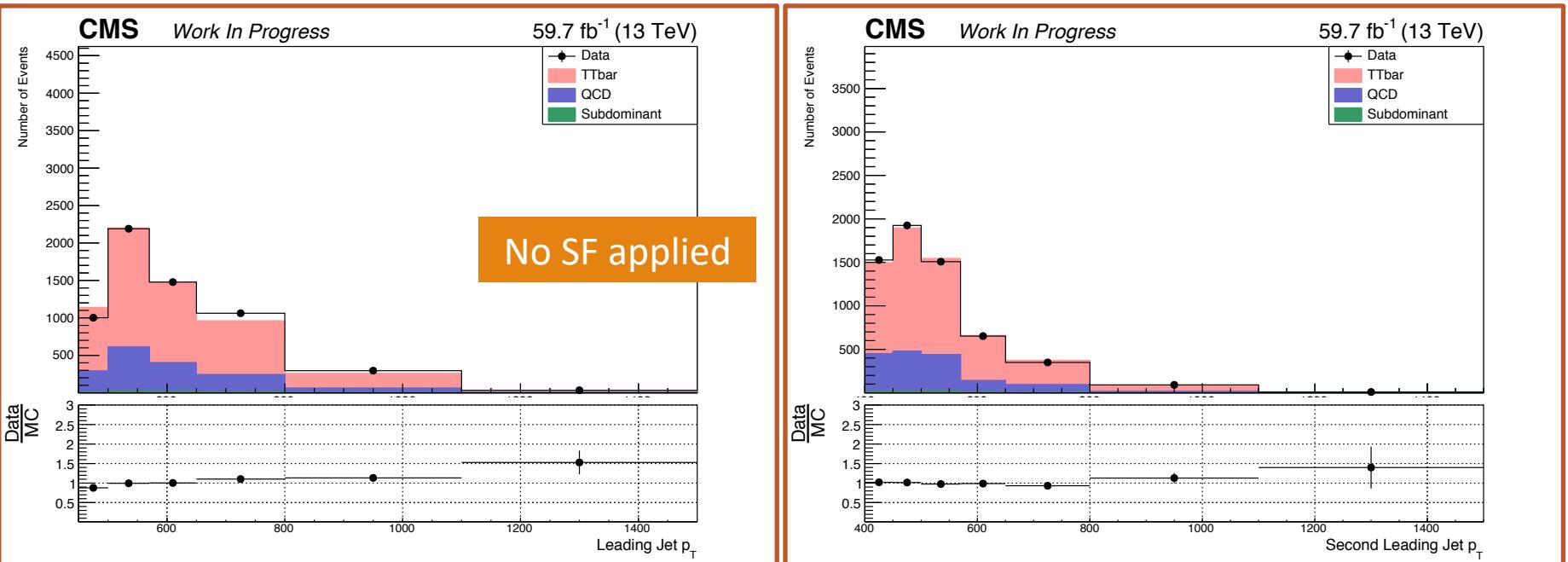
SF Type	Value \pm error
Inclusive	0.958 \pm 0.037
p_T [400, 600]GeV	0.956 \pm 0.048
p_T [600, 800]GeV	0.967 \pm 0.062
p_T [800, Inf)GeV	0.914 \pm 0.13

Table 42: Top Tagger SF Values for 2018 using JMAR proposed p_T regions.

SF Type	Value \pm error
Inclusive	0.958 \pm 0.037
p_T [400, 500]GeV	0.912 \pm 0.093
p_T [500, 600]GeV	0.971 \pm 0.055
p_T [600, Inf)GeV	0.959 \pm 0.056

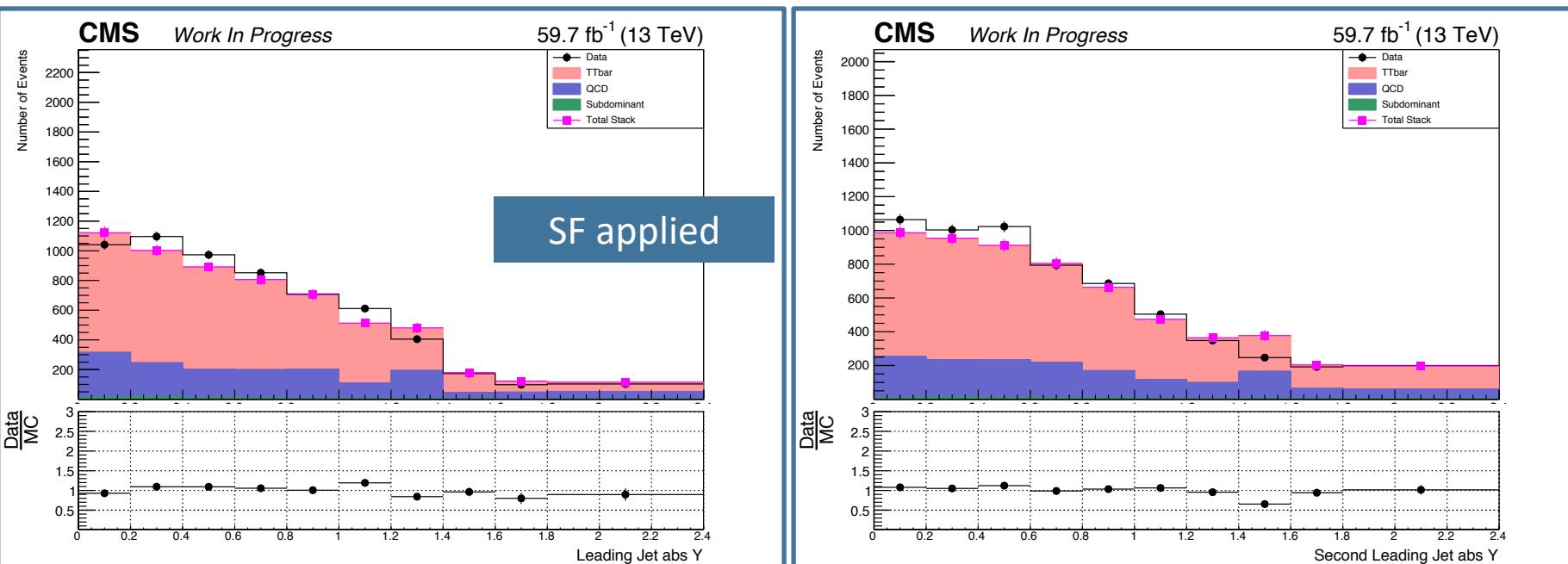
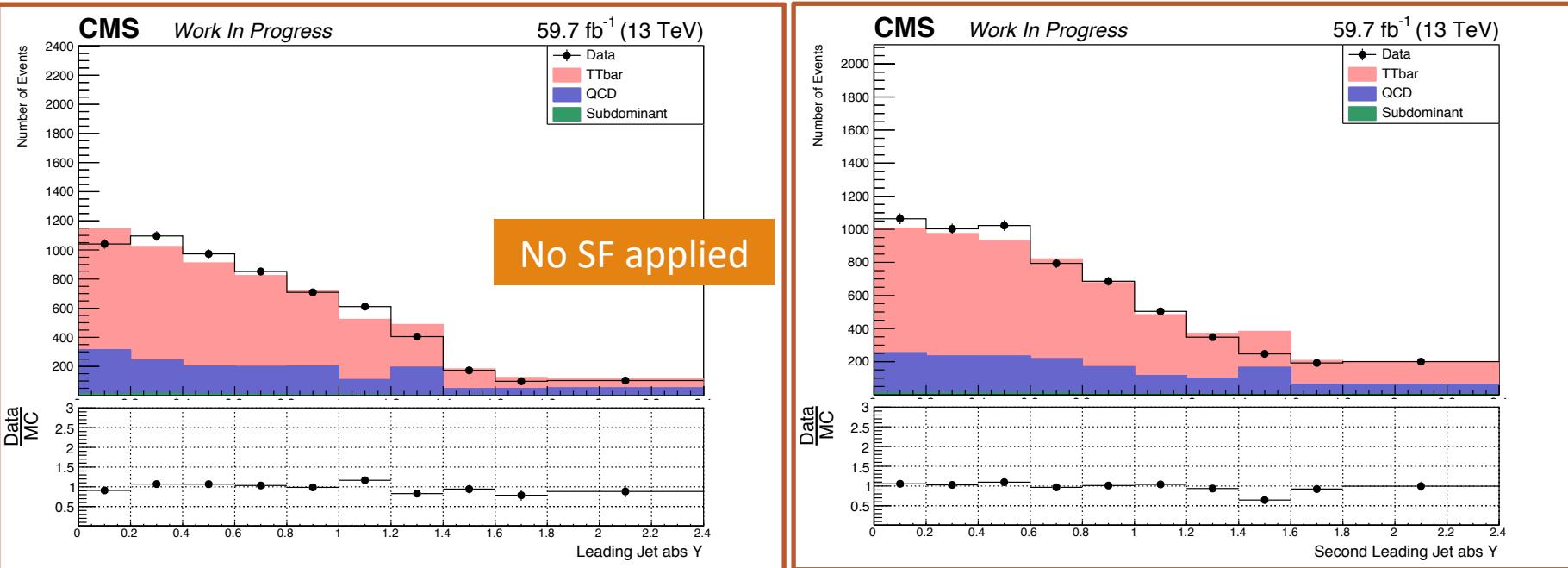
Data vs MC

2018



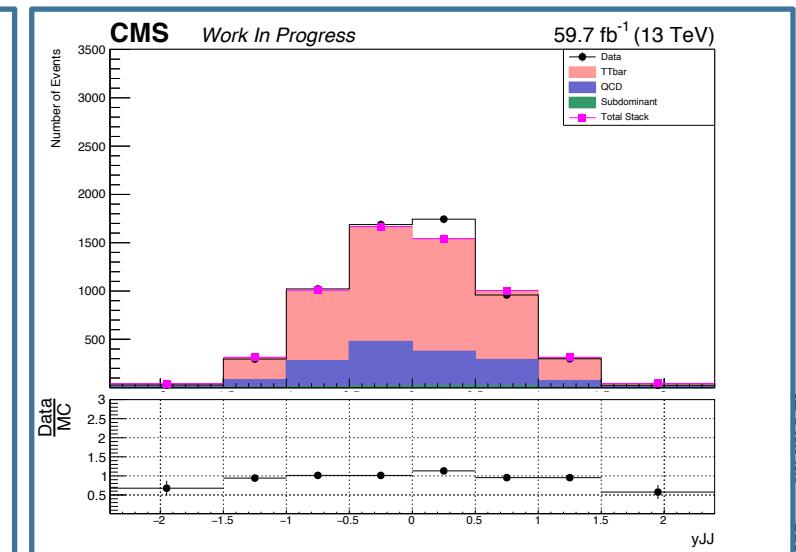
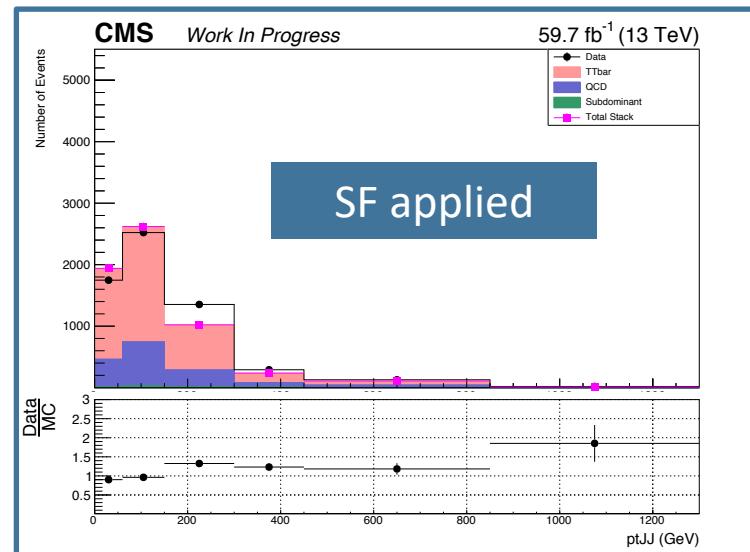
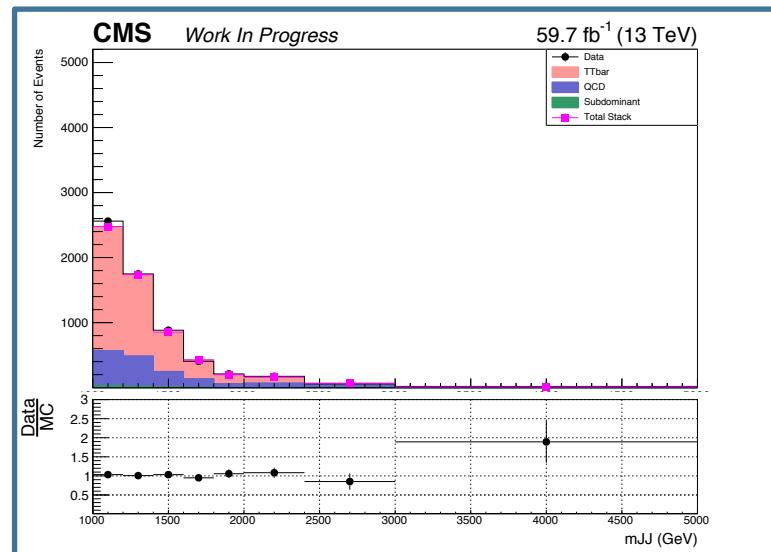
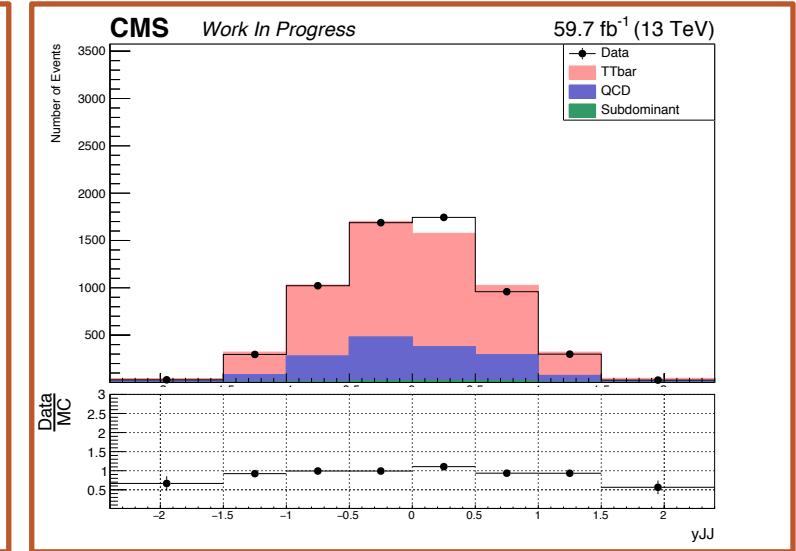
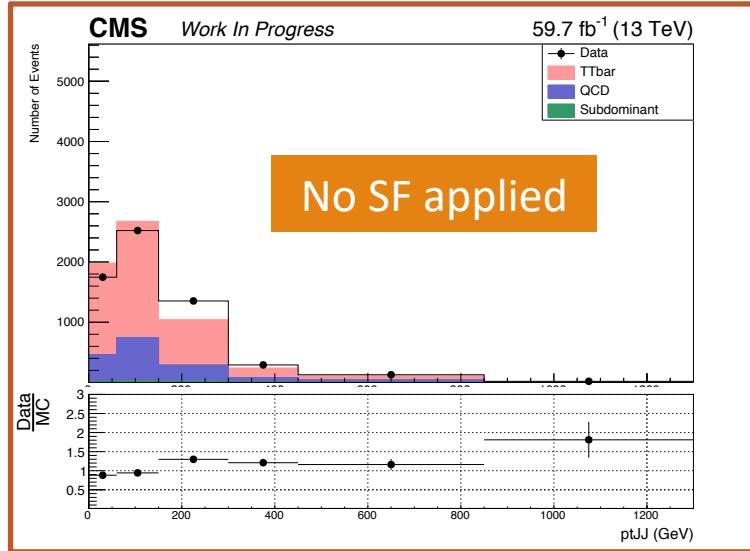
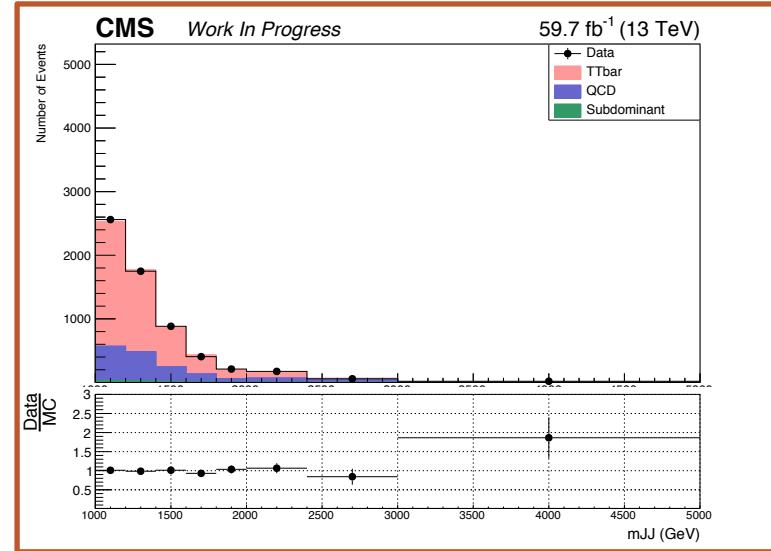
Data vs MC plots

2018



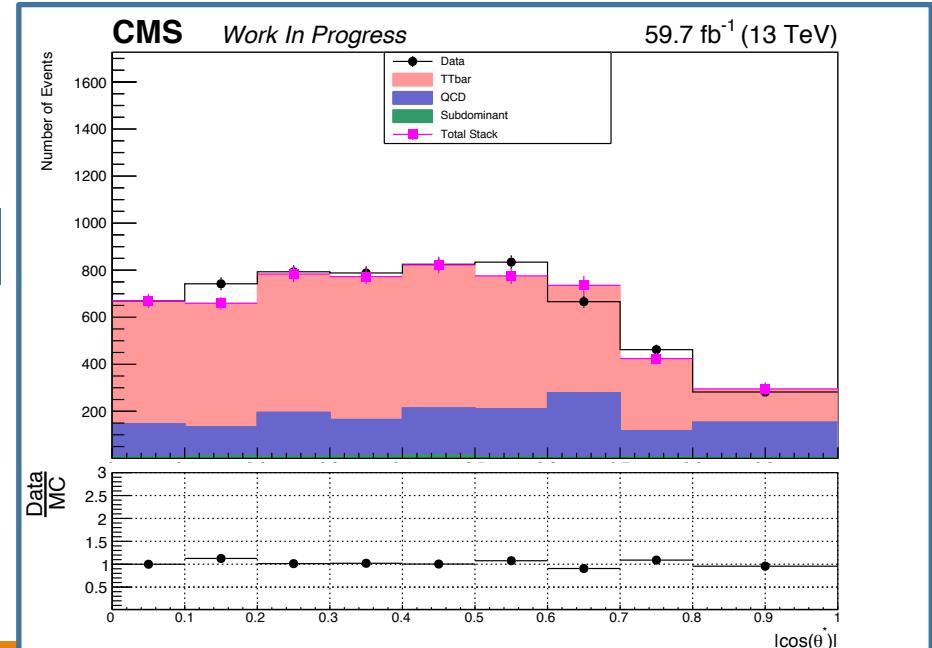
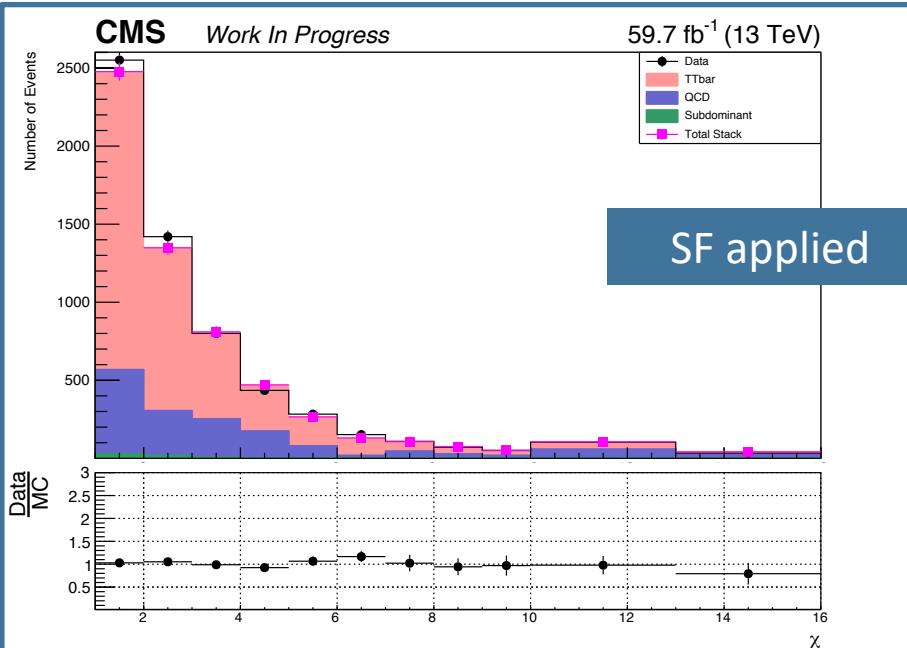
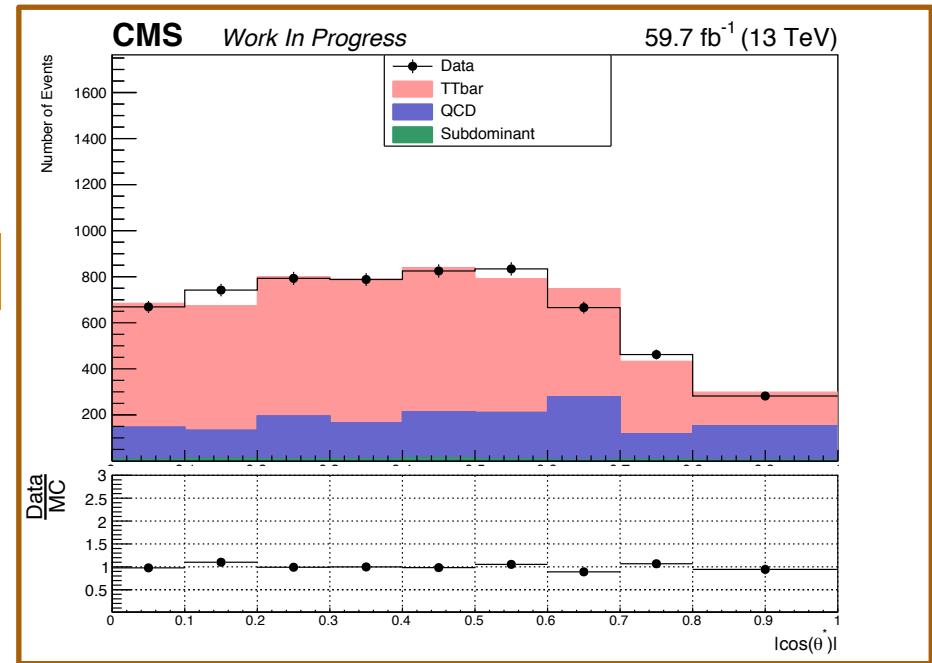
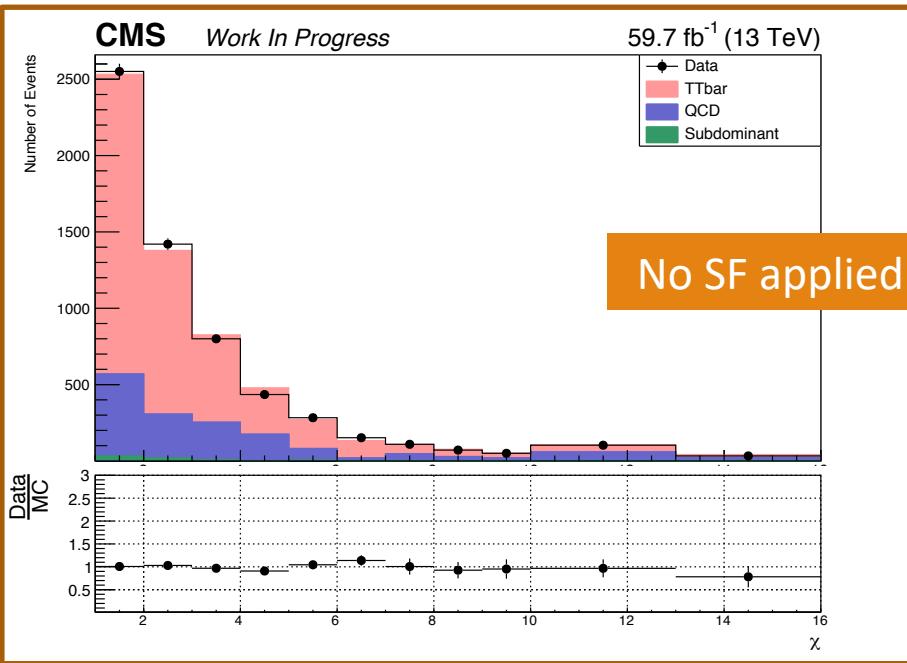
Data vs MC plots

2018



Data vs MC plots

2018

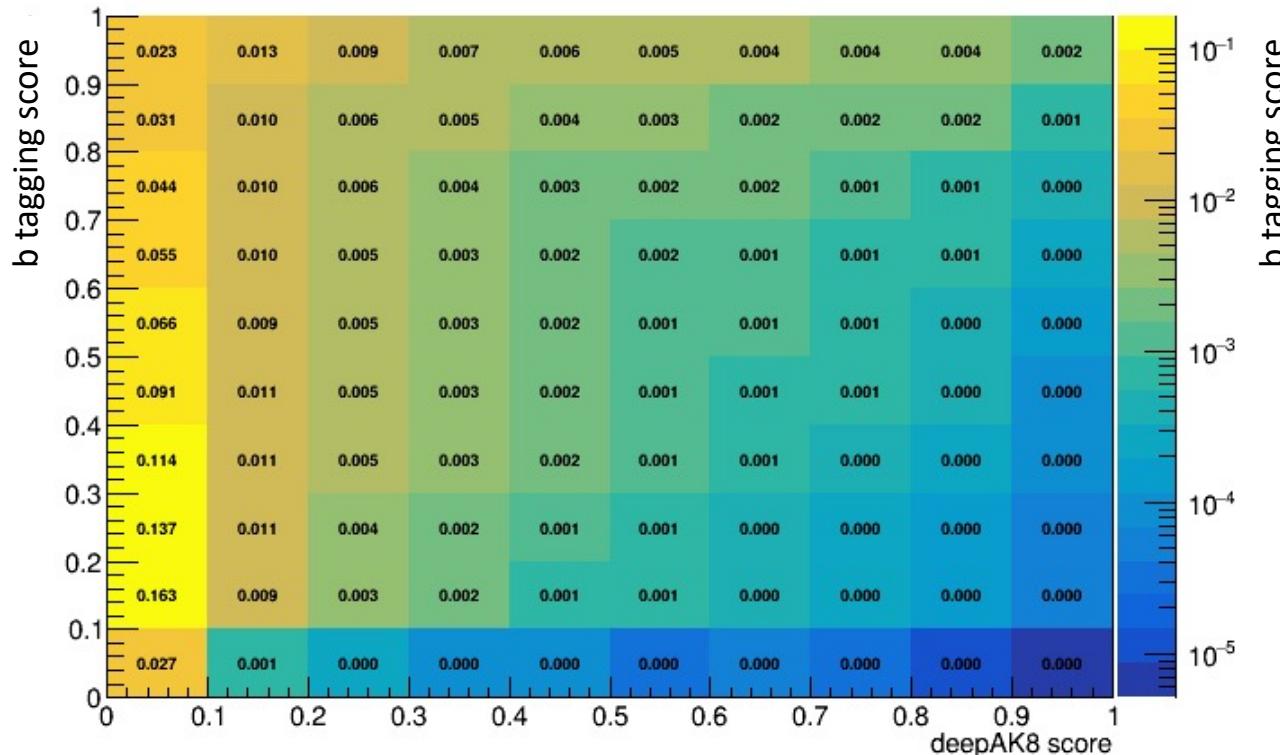


DeepAK8 backup

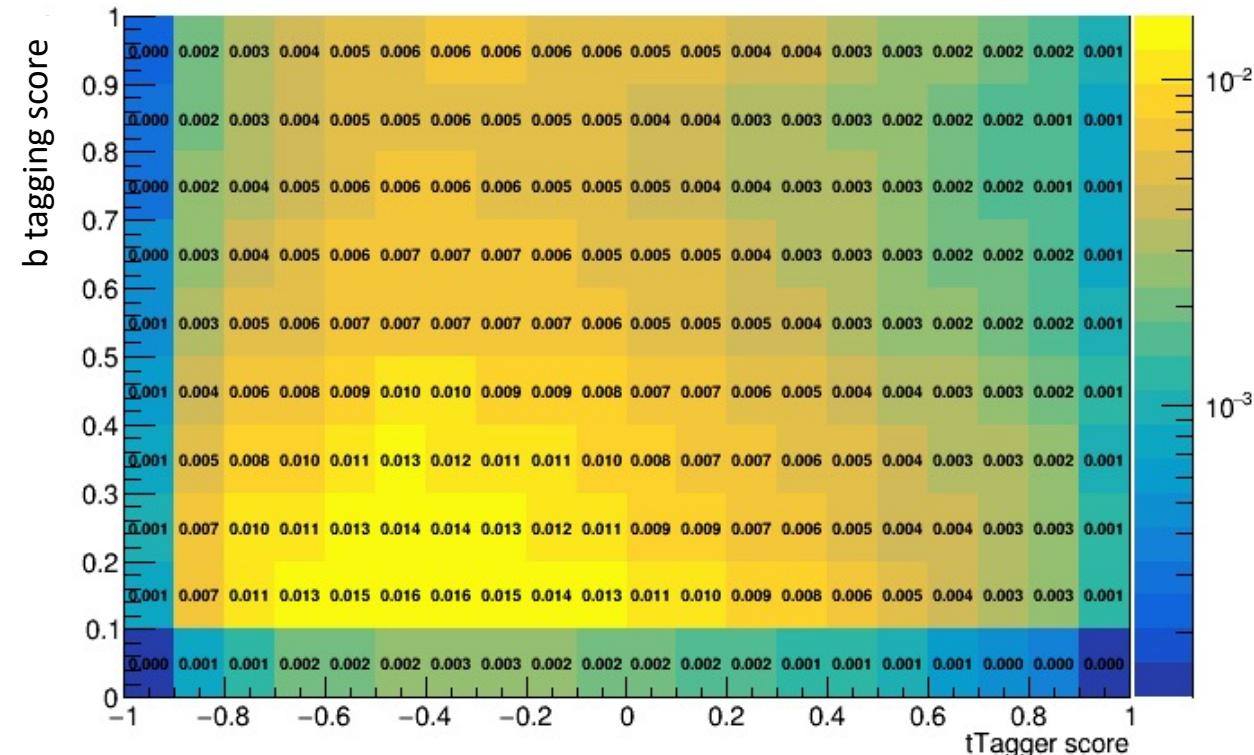


Btagging Correlation to DeepAK8 or TopTagger

deepAK8 b tagging correlation bkg sample

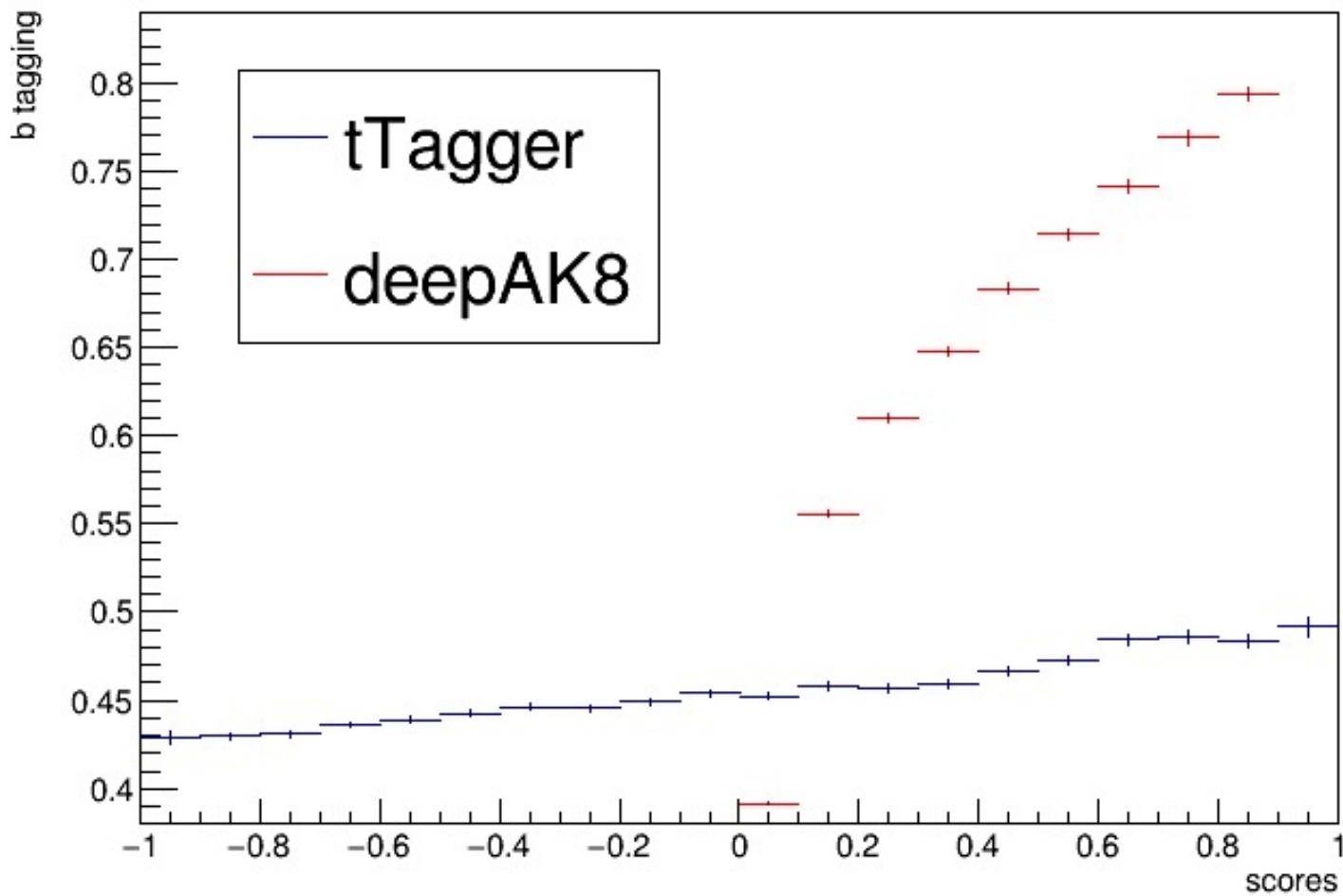


tTagger b tagging correlation bkg sample



Btagging Correlation to DeepAK8 or TopTagger profiles

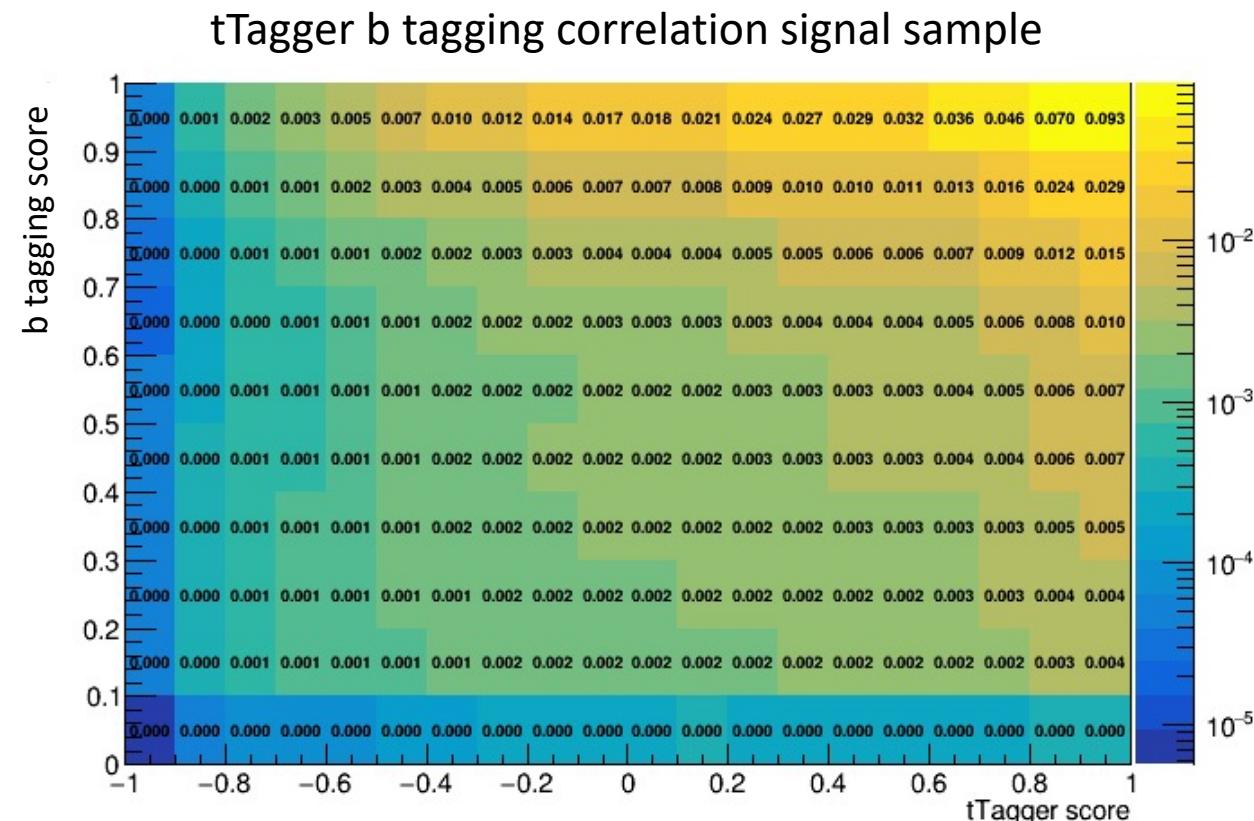
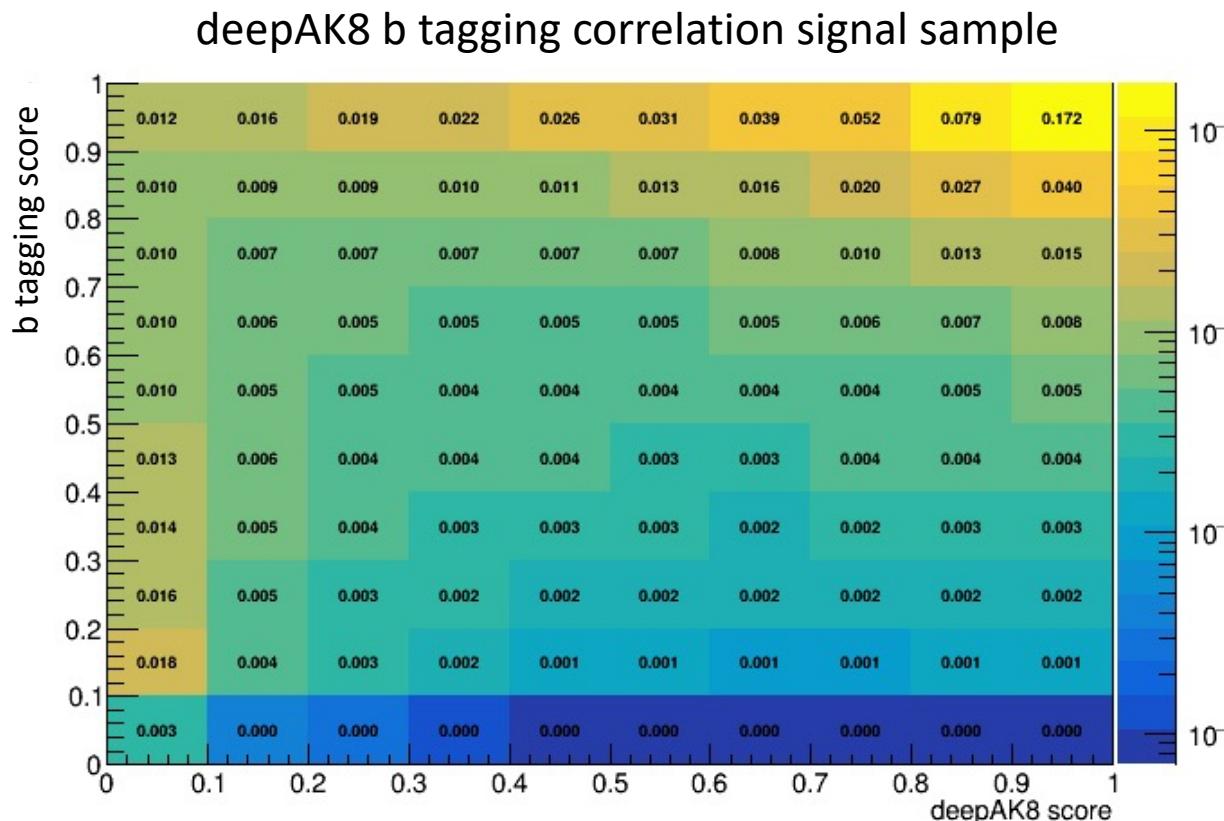
Tagger profiles bkg



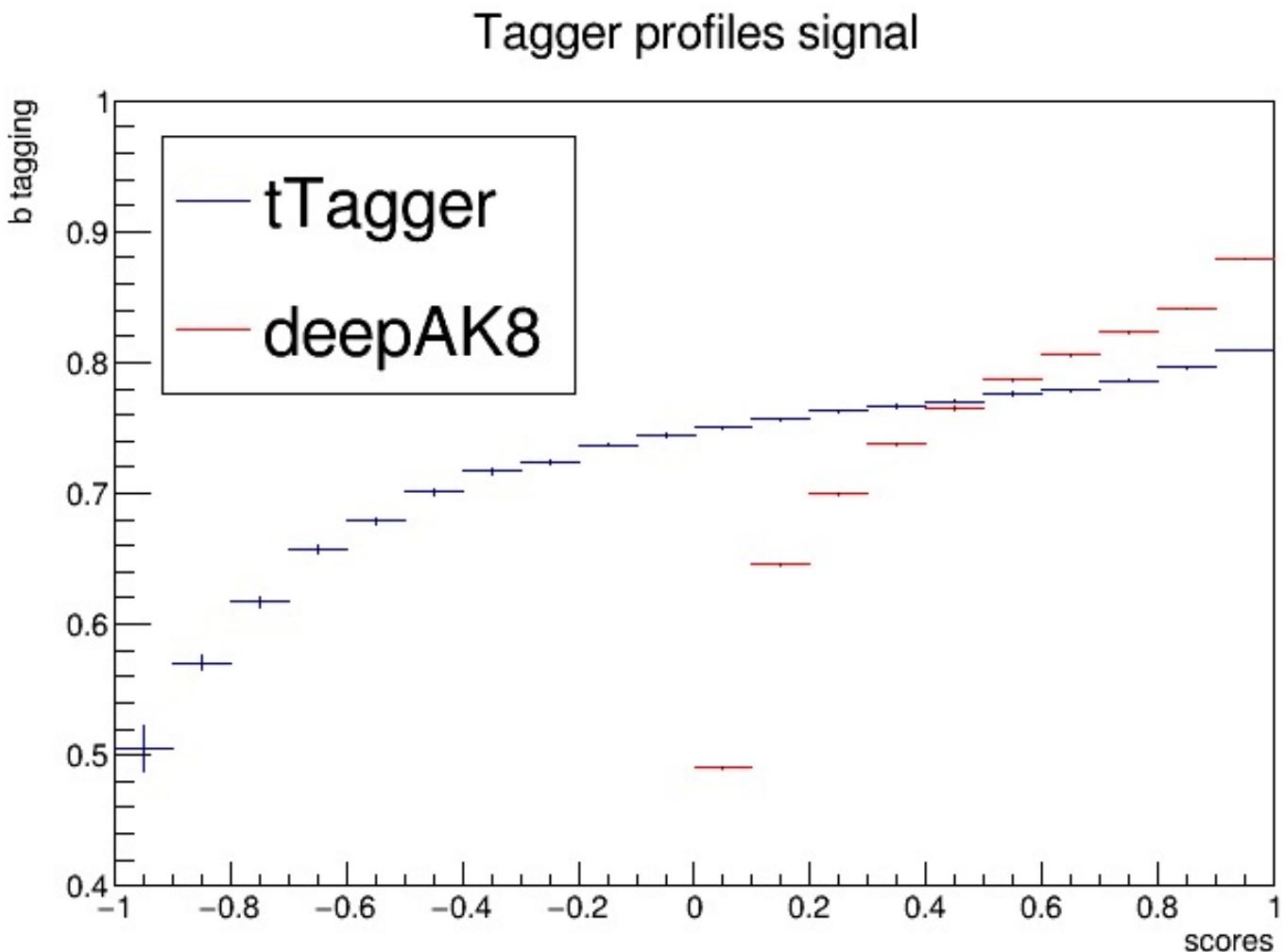
b-tagging mean value
vs
top tagging values



Btagging Correlation to DeepAK8 or TopTagger



Btagging Correlation to DeepAK8 or TopTagger profiles

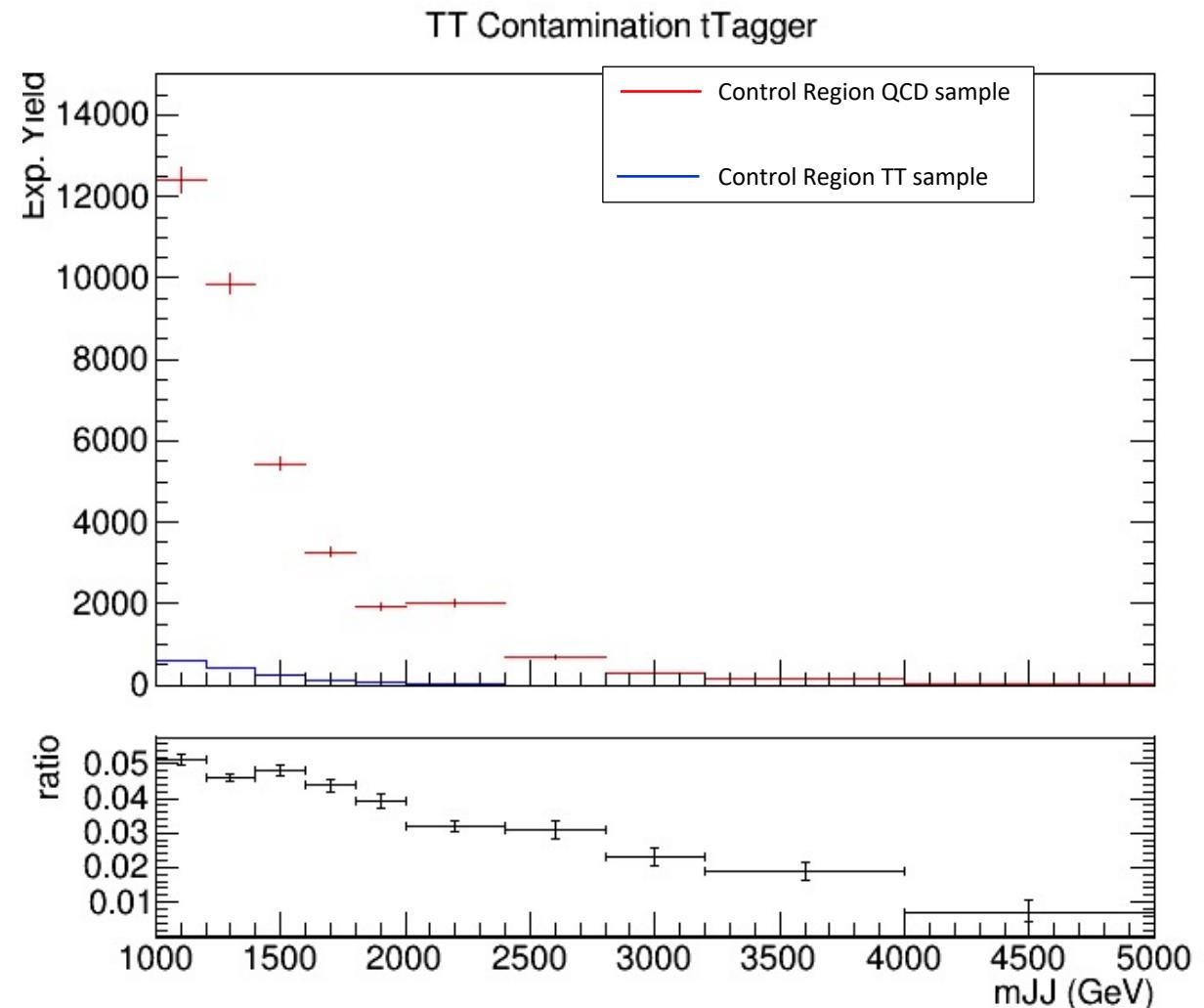
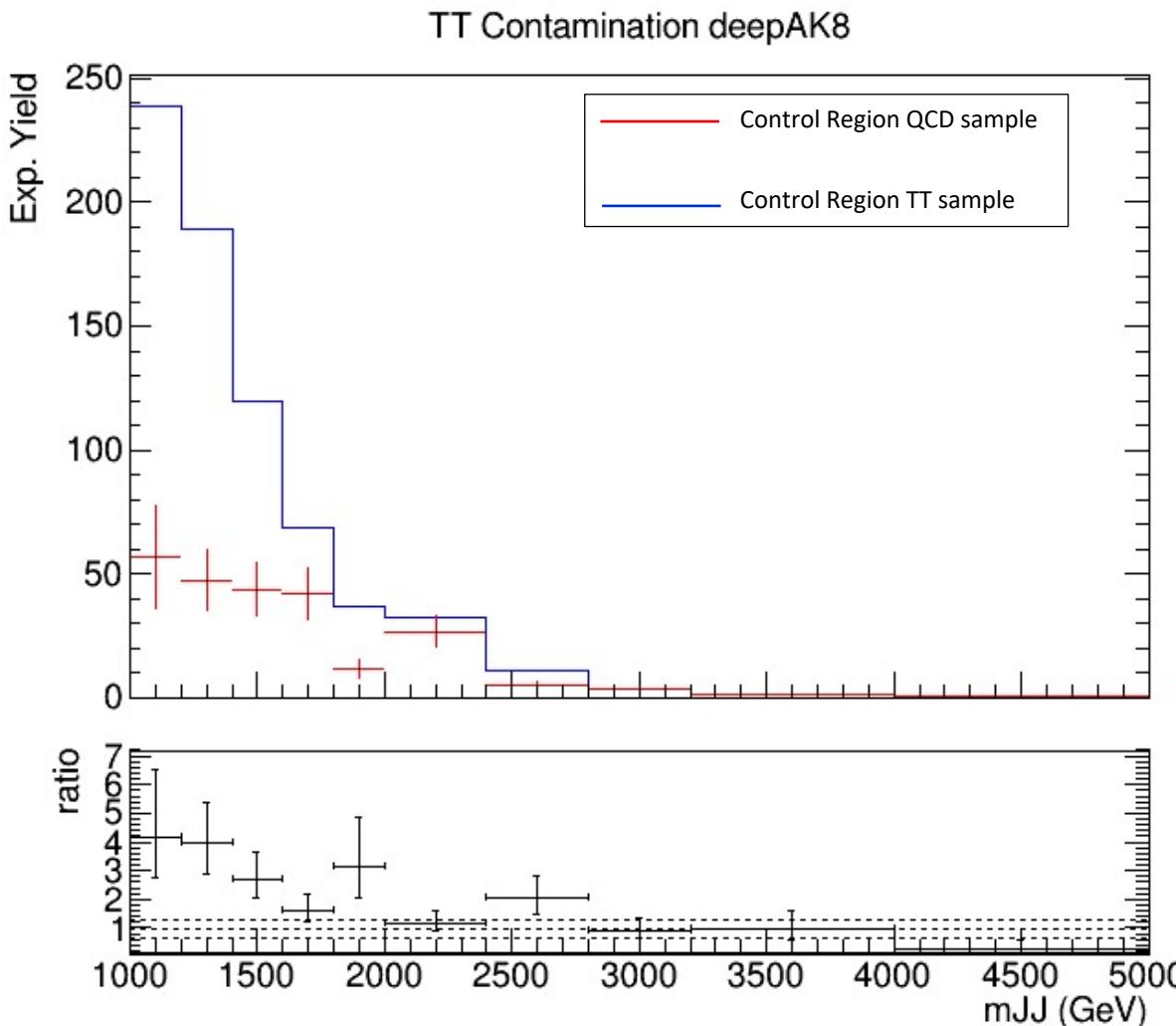


b-tagging mean value
vs
top tagging values



Control Region Contamination

- Expected yield from QCD Bkg samples and TT Signal sample in the CR ($b\text{tag} == 0$) vs m_{JJ}
- The QCD contribution is used to get the QCD shape



Control Region Contamination

- Expected yield from QCD Bkg samples and TT Signal sample in the CR ($b\text{tag} == 0$) vs jetPt
- The QCD contribution is used to get the QCD shape

