



Fake Estimation in the $t\bar{t}H(H \rightarrow b\bar{b})$ Analysis in the Single Lepton Channel with the ATLAS Experiment

Johannes Mellenthin
Supervised by Arnulf Quadt

II. Physikalisches Institut, Georg-August-Universität Göttingen

11th Annual Meeting of the Helmholtz Alliance "Physics at the Terascale"
28.11.2017

SPONSORED BY THE



Federal Ministry
of Education
and Research



BMBF-Forschungsschwerpunkt
ATLAS-EXPERIMENT

Physik bei höchsten Energien mit dem ATLAS-Experiment am LHC

FSP 103

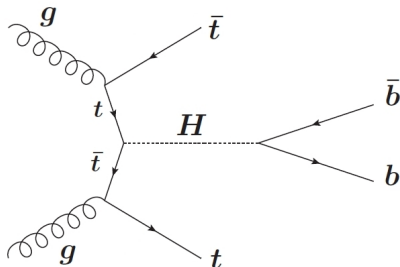
ATLAS

$t\bar{t}H(H \rightarrow b\bar{b})$

- Direct sensitivity to top and b -quark Yukawa coupling
- Study Higgs fermion couplings
- $H \rightarrow b\bar{b}$ has the highest BR (58%)

Results

- Full $t\bar{t}H(H \rightarrow b\bar{b})$ analysis with 36.1 fb^{-1} data gives $\mu > 2.0$ excluded at 95% C.L.
- Significance: 1.4σ (1.6σ exp.)
- *ATLAS-CONF-2017-076*

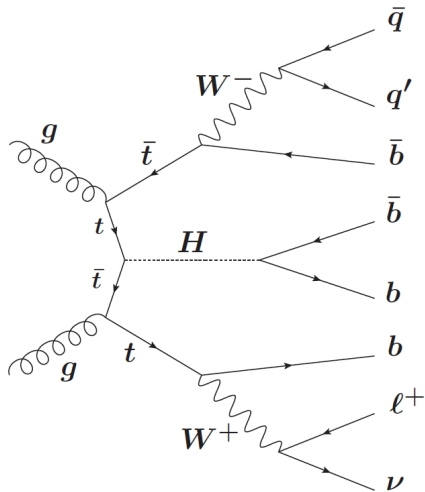


Single Lepton Channel

- 1 leptonic W decay
- 1 electron or 1 muon
- 6 jets with 4 b -jets

Very challenging analysis

- 4 b -jets in the signal region
 - Large background from $t\bar{t}$ + jets
- Strategy: Divide into different regions



Event Selection

- Events triggered by single lepton triggers
- At least 5 jets
- At least 2 b -tags at 60% working point (WP) or at least 3 b -tags at 77% WP

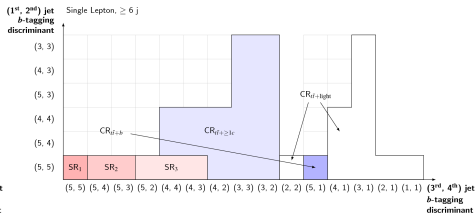
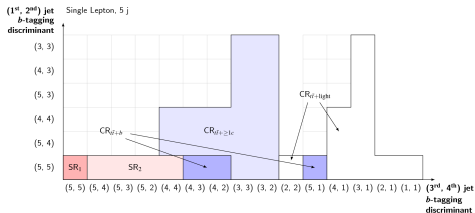
Categorisation

- Split in N jets
 - 5 jets
 - at least 6 jets
- Split in N b -tags
 - Use 4 working points for 5 ranges

	none (1)	loose (2)	medium (3)	tight (4)	very tight (5)
WP	[100%, 85%]	[85%, 77%]	[77%, 70%]	[70%, 60%]	[60%, 0%]
- Define 11 regions enriched in $t\bar{t}H$, $t\bar{t} + b$, $t\bar{t} + c$, and $t\bar{t} + \text{light jets}$

Analysis Regions

- 5 (ultra) pure **signal regions (SR)**; train BDT in each SR
- 6 **control regions (CR)** which are enriched either in $t\bar{t} + b$, $t\bar{t} + c$, and $t\bar{t} + \text{light jets}$ to constrain systematic uncertainties on the background



→ Combined profile likelihood fit to all regions is performed

Composition

- Signal process $t\bar{t}H$
 - Cross-section: ~ 0.5 pb
 - MadGraph5_aMC@NLO+Pythia8
 - Dominating background
 - $t\bar{t}$ + jets
 - Cross-section: ~ 832 pb
 - Powheg+Pythia8
 - Other backgrounds
 - $t\bar{t}V$, single top, W/Z + jets, diboson
 - **Multi-jet** (Fakes and non-prompt)
($\sim 4.5\%$ in CRs and $\sim 1.3\%$ in SRs)
- Good estimate of multi-jet background improves data/MC agreement and has impact on $t\bar{t}$ +HF normalisation

ATLAS Preliminary
 $\sqrt{s} = 13$ TeV
Single Lepton

$t\bar{t}$ + light
 $t\bar{t}$ + $\geq 1c$
 $t\bar{t}$ + $\geq 1b$
 $t\bar{t}$ + V
 Non- $t\bar{t}$

CR _{$t\bar{t}$ +light}^{5j}



CR _{$t\bar{t}$ + $\geq 1c$} ^{5j}



CR _{$t\bar{t}$ +b}^{5j}



SR₂^{5j}



SR₁^{5j}



SR^{boosted}



CR _{$t\bar{t}$ +light} ^{$\geq 6j$}



CR _{$t\bar{t}$ + $\geq 1c$} ^{$\geq 6j$}



CR _{$t\bar{t}$ +b} ^{$\geq 6j$}



SR₃ ^{$\geq 6j$}



SR₂ ^{$\geq 6j$}

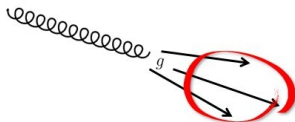
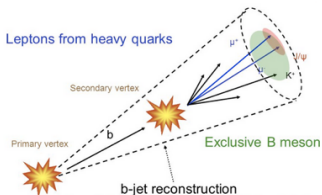


SR₁ ^{$\ge 6j$}



Fake and non-prompt leptons

- Non-prompt leptons
 - Semi-leptonic decays from c - and b -quarks
 - Photon conversion
 - Kaon decay (usually small)
 - Fake leptons
 - Jets can be misidentified as reconstructed lepton (mostly electrons)
 - Consider electrons and muons separately
 - Difficult to model this from MC, huge statistics needed
- Very important to measure this type of background from data

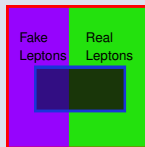


Selection

- Fully data-driven estimate using matrix method
- Introducing a **loose** lepton selection alongside the **tight**
- Measure behaviour of real and fake leptons in CRs
 - CRs contain events that are enriched in fakes and non prompt leptons
 - Use an extrapolation weight to relate these events to the background in the SR

$$N^{\text{loose}} = N_{\text{real}}^{\text{loose}} + N_{\text{fake}}^{\text{loose}}$$

$$N^{\text{tight}} = N_{\text{real}}^{\text{tight}} + N_{\text{fake}}^{\text{tight}}$$



- Define fake efficiency f and real efficiency r as


$$f = \frac{\text{small purple}}{\text{large purple}} \quad r = \frac{\text{small green}}{\text{large green}}$$

$$N^{\text{loose}} = N_{\text{real}}^{\text{loose}} + N_{\text{fake}}^{\text{loose}}$$
$$N^{\text{tight}} = r \cdot N_{\text{real}}^{\text{loose}} + f \cdot N_{\text{fake}}^{\text{loose}}$$

Lepton Selection

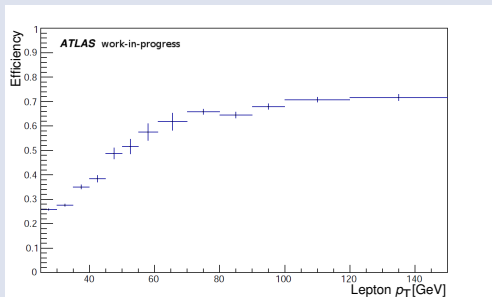
- Introducing a **loose** lepton selection alongside the **tight**

	Loose selection	Tight selection
Electron identification level	MediumLH	TightLH
Muon identification level	Medium	Medium
Lepton isolation requirement	None	Gradient

- Measure behavior of fakes in fake enriched regions
 - Electrons: single charged lepton + neutrino and low E_T^{miss}
 - Muons: decay from semileptonic b -quarks

Fake Efficiency

$$\frac{N(\text{Data} - \text{MC}_{\text{real}})_{\text{tight}}}{N(\text{Data} - \text{MC}_{\text{real}})_{\text{loose}}}$$

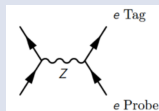


- Efficiencies vary as function of kinematic quantities such as lepton p_T , lepton η , leading jet p_T , E_T^{miss}

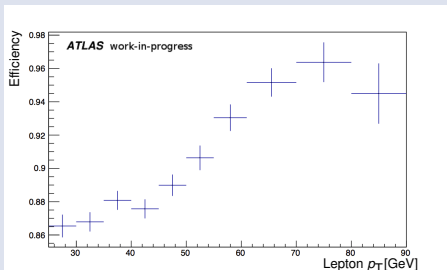
Real Efficiency

- Tag-and-Probe on $Z \rightarrow \ell\ell$ events

$$\frac{\text{Number of probes that pass tight}}{\text{Number of all probes}}$$



- Events with a pair of SFOS loose or tight leptons and at least one jet



Weight

- Calculate a per-event weight and apply it to loose selection data

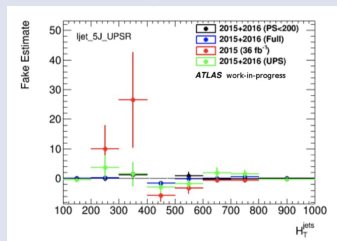
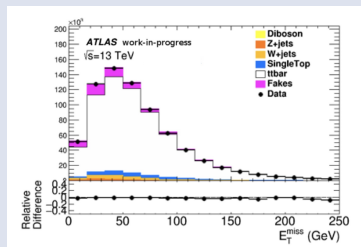
$$w = \frac{f}{r - f} (r - P)$$

where P is 1 if the loose event passes the tight selection and 0 otherwise

- Negative bin entries can occur if most of the loose events pass the tight selection
 - Need to go to a looser lepton definition
 - High fake rates
- Use pre-scaled (PS) triggers (only every n^{th} event passing a (loose) trigger is saved)

Limitations

- Good modelling in both electron and muon channels, but:
 - Only in low b -tagging multiplicity regions
 - “Spikes” in larger b -tagging multiplicity regions
 - Events from the electron channel with high trigger pre-scales
 - Bins with low statistics containing only tight (loose) events → negative (very large) weights



→ Use a **tag-rate-function (TRF)** to increase statistics

What is TRF?

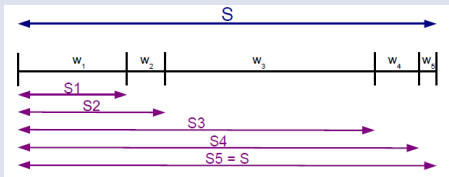
- Usually used to avoid having fluctuations in MC in regions with low statistics due to large b -tagging multiplicities (we have ≥ 4 b -tag jets in SR)
 - Idea: Instead of removing events if they do not pass the number of required b -tags, consider all events, reweighted based on the **probability of the event to contain n b -tags**
- Keep inclusive statistics
- Given $\epsilon(f, \eta, p_T)$ being the efficiency for a jet with η , p_T , and flavour f to be b -tagged, the probability for an event with N jets to contain 1 b -tag is

$$P_{=1} = \sum_{i=1}^N \left(\epsilon_i \prod_{i \neq j} (1 - \epsilon_j) \right)$$

- Inclusive b -tagging regions are computed with $P_{\text{incl}} = 1 - P_{=0}$

How to compute it

- To compute discriminating variables, b -tagged jets probability needs to be known
 - Compute sum of TRF weights S of all permutations corresponding to the desired number of b -tags (e.g. 1 b -tag among 5 jets)
 - Some of the weights w_i can be much larger than others
 - Throw a random number uniformly distributed and see in which w_i segment it falls
 - b -tag the jets based on the configuration of the i^{th} permutation
- Probability to pick a permutation i is proportional to its TRF weight



TRF with Matrix Method

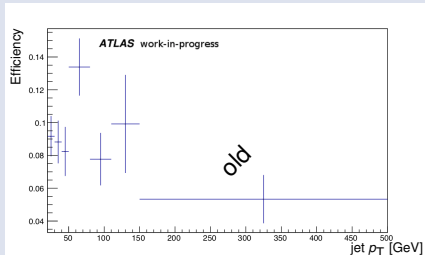
- Typically applied to MC samples where the truth flavour of a jet is known
- This information is not available for the matrix method (data-driven)
- But b -tagging can be used to label the jets as “ b ” or “not- b ”
- Efficiencies are flavour-blind

$$\epsilon(X|N_{\text{jets}}) = \frac{X_{b\text{-tagged}}}{X_{\text{all}}}$$

- **Hybrid TRF approach:** Reduce possible modeling discrepancies by extrapolating from inclusive sample with m b -tags to desired n b -tags ($n > m$)
 - Smaller statistics gain compared to 0 inclusive b -tags
- Use **matrix method** to estimate fakes in a b -tag inclusive region
- Apply **TRF** weight to estimate on desired high b -tagging multiplicity region

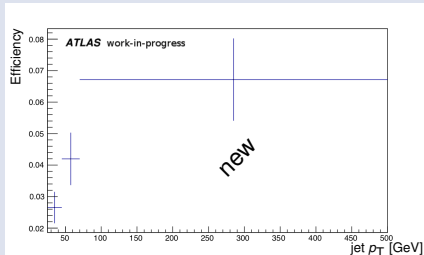
$t\bar{t}H$ Analysis

- Extrapolate from at least 1 b -tagged jet with a WP of 85%
- Lepton η , leading jet p_T and Delta R as variables for the parametrisation of the fakes (MM)
- Without PS trigger for the low p_T region in 2016 data
- Computation of TRF efficiencies
 - For WPs of 60%, 70%, 77%, 85%
 - Hybrid approach: ≥ 1 b -jet
 - 4 jet inclusive
 - Rebinning
 - Decrease uncertainty
- Different b -tagging efficiency parametrisations for TRF, but focus on jet p_T (most significant)
 - 00001: jet p_T
 - 00011: jet p_T and jet η
 - 01001: Delta $R(\text{jet}, \ell)$ and jet p_T
 - 10001: Delta $\phi(\text{jet}, E_T^{\text{miss}})$ and jet p_T



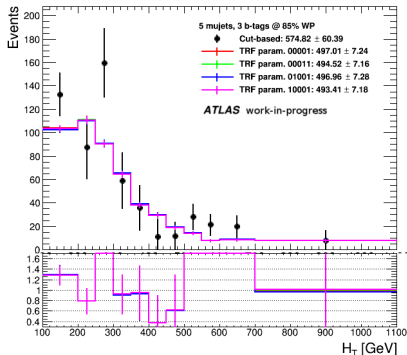
$t\bar{t}H$ Analysis

- Extrapolate from at least 1 b -tagged jet with a WP of 85%
- Lepton η , leading jet p_T and Delta R as variables for the parametrisation of the fakes (MM)
- Without PS trigger for the low p_T region in 2016 data
- Computation of TRF efficiencies
 - For WPs of 60%, 70%, 77%, 85%
 - Hybrid approach: ≥ 1 b -jet
 - 4 jet inclusive
 - Rebinning
 - Decrease uncertainty
- Different b -tagging efficiency parametrisations for TRF, but focus on jet p_T (most significant)
 - 00001: jet p_T
 - 00011: jet p_T and jet η
 - 01001: Delta $R(\text{jet}, \ell)$ and jet p_T
 - 10001: Delta $\phi(\text{jet}, E_T^{\text{miss}})$ and jet p_T



Results

- H_T^{had} : scalar sum of the jets' p_T
- Comparison of cut-based with TRF
- Agreement of both methods for electrons as well as muons
- Reduction of the statistical uncertainties with TRF
- Very good agreement between the different parametrisations
→ use jet p_T
- Removal of spikes
- Smoothens distribution



Fakes in the $t\bar{t}H$ Analysis

- Estimating fakes and non-prompt leptons is important for a good $t\bar{t}H$ measurement
 - Uncertainties for TRF efficiencies are small
 - Different TRF parameterisations very similar → use jet p_T
 - TRF smoothens distributions
- Implementation for continuous b -tagging

The image features a large, light blue, stylized logo consisting of the letters 'G' and 'A' intertwined. The 'G' is on the left and the 'A' is on the right. Below the 'A' are the numbers '1', '7', '3', and '7' in a light blue, serif font. The text 'Thank you for your attention!' is centered over the logo in a bold, black, sans-serif font.

Thank you for your attention!