# Data-driven estimation of fake $\tau$ background in Higgs searches in ATLAS

**Marzieh Bahmani, on behalf of the ATLAS Collaboration**[*]

*Institute of Nuclear Physics, Polish Academy of Sciences.*
*E-mail:* marzieh.bahmani@ifj.edu.pl

In analyses using reconstructed $\tau$ leptons, estimation of backgrounds arising for jets misidentified as hadronically decaying $\tau$s becomes a crucial issue. This paper presents two methods employed by the ATLAS experiment at the LHC, the fully data-driven fake factor method and the semi-data-driven fake rate method. Example of fake factor method in the background modelling of $H^{\pm} \to \tau\nu$ analysis and the fake rate method applied to the search for high mass resonances decaying to $\tau\tau$, in the $\tau_{\text{had}}\tau_{\text{had}}$ channel, are presented. The systematics associated with the methods are also discussed.

---

[*]Speaker.

## 1. Introduction and motivation

An important source of background for physics analyses using particle-level identification criteria, is misidentification. One of the primary motivation for estimating the fake $\tau$ background in the $\tau$ related analyses is that jets genuinely fake hadronically decaying $\tau$ leptons, and these analyses suffer from such a background. The fake $\tau$ background is not well modeled by Monte Carlo (MC) therefore, data driven techniques had to be developed. There are different approaches to the estimation of jet-to-$\tau$ misidentified contribution to a $\tau$ selection. In this paper, we focus on two main methods used in the ATLAS experiment [1]: the fake factor method and the fake rate method. The former method is fully data-driven, while the latter is semi-data-driven since the data-driven efficiency factor is applied to MC. In section 2 the general ideas behind the fake factor method are presented and the example analysis searching for $H^{\pm} \to \tau \nu$ [2] using the fake factor method is shown. Section 3 describes the fake rate method and an example of its application in the search for high mass resonances decaying to $\tau \tau$ [4], in the $\tau_{\text{had}} \tau_{\text{had}}$ channel is shown.

## 2. Fake factor data-driven method

The dominant background processes can be categorized based on the object that gives rise to reconstructed and identified hadronically decaying $\tau$ candidate [3]. These are mostly quark- or gluon-initiated jets fulfilling selection criteria of the signal region. These backgrounds are poorly modelled due to the statistical limitations in the sample of simulated events (e.g. multi-jet processes). Also the systematic uncertainties related to object misidentified as $\tau$ are not well known. Therefore, a data-driven approach is used to estimate this background. In the fake factor (FF) method, background processes where a quark- or gluon-initiated jet is reconstructed and identified as a $\tau$ candidate are estimated from data. For this purpose, an anti-$\tau$ selection is defined by requiring the $\tau$ candidate to fail the identification criteria of the nominal selection. The fake factor is defined as the ratio between the number of jets reconstructed as $\tau$ candidates and fulfilling the nominal $\tau$ identification criteria to the number of corresponding candidates failing the identification criteria (anti-$\tau$) and is measured in a dedicated control region (CR) enriched with fake $\tau$s:

$$\text{FF} = \frac{\text{N}^{\text{CR}}_{\tau-\text{ID}}(\text{data}) - \text{N}^{\text{CR}}_{\tau-\text{ID}}(\text{MC}, \tau \neq \text{j})}{\text{N}^{\text{CR}}_{\text{anti}-\tau-\text{ID}}(\text{data}) - \text{N}^{\text{CR}}_{\text{anti}-\tau-\text{ID}}(\text{MC}, \tau \neq \text{j})} \tag{2.1}$$

In order to obtain the $\text{N}^{\text{CR}}_{\text{anti}-\tau}$ and $\text{N}^{\text{CR}}_{\tau}$, contribution from true $\tau$ events in either categories are subtracted using simulation. The fake factors are usually measured in bins of $p_{\text{T}}$ or number of associated tracks in the $\tau$ hadronic decay (1-prong, 3-prong), they can also be measured in opposite- or same-sign regions, with or without b-jets, depending on the topology of interest for the analysis.

### 2.1 Considering quark-gluon jet composition

The fake factors are usually extracted in control regions enriched in either gluon-initiated or quark-initiated jets, as the probability for a hadronic jet to fake a $\tau$ depends on its origin. Depending on the analysis, there can be one or several control regions where the fake factors are measured. In case there is only one control region, one must ensure that the origin of fake $\tau$ composition is

| Multi-jet CR | W+jets CR |
|---|---|
| number of jet at least 2 | one electron or muon |
| $E_T^{miss} < 80$ GeV | at least one reconstructed $\tau_{\text{had}-\text{vis}}$ candidate |
| bjets veto, electron and muon veto | bjets veto |
| $p_T$ of $\tau > 30$ GeV | $p_T$ of electron and muon > 30 Gev |
| $m_T(\tau , E_T^{miss})$ >50 GeV | $60 < m_T(\ell, E_T^{miss}) < 160$ Gev |
| $\tau$ identification score $> 0.02$ | $\tau$ identification score $> 0.02$ |

**Table 1:** Control regions for fake factor measurement in $H^{\pm} \to \tau\nu$ analysis [2].

close to the one in the signal region. When two (or more [6]) control regions are used, and one is enriched in gluon-initiated jet, the FF for each bin is calculated as follows:

$$\text{FF} = \alpha_{\text{g}} \times \text{FF(g)} + [1 - \alpha_{\text{g}}] \times \text{FF(other(s))} \qquad (2.2)$$

where FF(g) and FF(other(s)) are the fake factor for gluon-initiated jet and other control regions, and $\alpha$ is the fraction of reflecting composition in the signal like anti-$\tau$ region. Therefore, one needs to estimate the $\alpha$.

In order to estimate the yield of fake $\tau$ background in the signal region, an anti-$\tau$ region is defined identical to the signal region but where $\tau$ candidate fails the $\tau$ identification requirement, instead of fulfilling it. Then in a bin i, the number of events with a jet misidentified as $\tau$ is given by :

$$\text{N}_{\text{fakes}}^{\tau}(\text{i}) = \text{N}_{\text{fakes}}^{\text{anti}-\tau-\text{ID}}(\text{i}) \times \text{FF(i)}, \qquad (2.3)$$

### 2.2 Fake factor method in $H^{\pm} \to \tau\nu$

In the $H^{\pm} \to \tau\nu$ analysis [2], in order to account for different sources of misidentified hadronically decaying $\tau$ lepton ($\tau_{\text{had}-\text{vis}}$) [3] in the signal region, fake factors are measured in two control regions of the data with different fractions of quark- and gluon-initiated jets, and then they are combined. The first control region with a significant fraction of gluon-initiated jets (multi-jet CR) is defined as shown in Table 1(left), Such events are collected using a combination of multi-jet triggers. The other control region enriched in quark-initiated jets (W+jets CR) is defined as shown in Table 1(right), and using single lepton trigger. The transverse mass $m_T$ of the $\tau$ candidate is obtained, as a function of the missing transverse energy $E_T^{\text{miss}}$ and the reconstructed $\tau$ momentum by eq. 2.4:

$$m_T = \sqrt{2 p_T^{\tau} E_T^{miss} (1 - cos\Delta\phi_{\tau,miss})} \qquad (2.4)$$

In the second control region the transverse mass variable $m_T(\ell, E_T^{miss})$ is computed as in the previous case using the lepton $p_T$ and separation in azimuthal angle from the missing transverse momentum of the event. The fake factors measured in these two control regions are shown in Figure 1 left plot. In the anti-$\tau_{\text{had}-\text{vis}}$ regions corresponding to the nominal event selections, the fraction of quark- and gluon-initiated jets misidentified as $\tau_{\text{had}-\text{vis}}$ candidates are then measured using a template-fit approach, based on variables that are sensitive to the difference in quark- and gluon-induced jets. For 3-prong $\tau_{\text{had}-\text{vis}}$ candidates, the $\tau$ identification score (based on the multivariate BDT approach) is used as a template. For 1- prong $\tau_{\text{had}-\text{vis}}$ candidates, the $\tau_{\text{had}-\text{vis}}$ jet width is used which is defined

as follow:

$$w_\tau = \frac{\Sigma[p_T^{track} \times \Delta R(\tau_{\text{had-vis}}, track)]}{\Sigma p_T^{track}} \tag{2.5}$$

where the sum runs over the tracks satisfying $\Delta R(\tau_{\text{had-vis}}, track) < 0.4$. In order to account for unknown gluon- and quark-initiated jets composition in the signal region, a linear combination of the two templates is defined as :

$$f(x|\alpha_{\text{MJ}}) = \alpha_{\text{MJ}} \times f_{\text{multi-jet}}(x) + (1 - \alpha_{\text{MJ}})f_{\text{W+jets}}(x) \tag{2.6}$$

with a free parameter $\alpha_{MJ}$ and the f(x) is the $\tau_{\text{had-vis}}$ jet width or the $\tau$ identification score. $f_{\text{multi-jet}}$ and $f_{\text{W+jets}}$ are two binned templates obtained in the multi-jet and W+jets control regions defined above, respectively. This linear combination is fitted to the normalized distribution measured in the signal region, by varying the $\alpha_{MJ}$ and in every bin of $p_T$ minimizing the $\chi^2$ distribution for each channel separately. Finally, from the best fit values of $\alpha_{MJ}$, combined fake factors are obtained by:

$$\text{FF}^{\text{comb}}(i) = \alpha_{\text{MJ}}(i) \times \text{FF}^{\text{multi-jet}}(i) + (1 - \alpha_{\text{MJ}}) \times \text{FF}^{\text{W+jets}}(i) \tag{2.7}$$

where i refers to each bin in the parametrization of fake factor. The combined fake factors, used in the $\tau_{\text{had-vis}}$+jets and $\tau_{\text{had-vis}}$+lepton signal regions are shown in Figure 1 right plot.



**Figure 1:** Fake factors parameterized as a function of $p_T^\tau$ and number of tracks. The left plot shows the fake factor in the multi-jet and w+jet CRs. Errors represent the statistical uncertainties. The right plot shows fake factors after reweighting by $\alpha_{MJ}$ in the $\tau_{\text{had-vis}}$+jets and $\tau_{\text{had-vis}}$+lepton channel. [2].

The dominant sources of systematic uncertainty of fake factor method are coming from a) the range of the $\tau_{\text{had-vis}}$ identification score in the anti-$\tau_{\text{had-vis}}$ definition of control samples, which modifies the corresponding fraction of quark- and gluon-initiated jets, as well as the event topology. b) the contamination of true $\tau_{\text{had-vis}}$ candidates fulfilling the anti-$\tau_{\text{had-vis}}$ selection. c) the statistical uncertainty of the control sample. d) the statistical error on the best fit value of $\alpha_{MJ}$. The impact of systematic uncertainty is different according to the $H^+$ mass. For low and intermediate mass of $H^+$ the dominant source of systematic uncertainty is caused by fake factor method, while in the high mass range it is caused by signal modelling as the contribution of background is smaller.

## 3. Fake rate method

Fake rates are defined as ratios of event yields with identified $\tau$s to the yields of all $\tau$ candidates without identification applied. They are applied to non-true $\tau$ objects in a signal-like region in MC. This is a semi-data-driven method, since fake rates are applied to simulated events. Fake rates are measured in dedicated control region as:

$$FR = \frac{N_{\tau-ID}(data) - N_{\tau-ID}(MC, \tau \neq j)}{N_{noID-\tau}(data) - N_{noID-\tau}(MC, \tau \neq j)} \quad (3.1)$$

where MC events in which a reconstructed $\tau$ is associated with a true $\tau$ at the generator level are subtracted. Fake rates are usually parameterized in bins of number of tracks, $p_T$ and $\eta$. Examples of analysis using fake rate method are $A/H/Z' \rightarrow \tau\tau$(had had) [4], hh$\rightarrow bb\tau\tau$ (had had) [6] and LFV $Z' \rightarrow l\tau$ [7].

### 3.1 Fake rate in high mass resonances decaying to $\tau\tau$

In the search for high mass resonances decaying to $\tau\tau$ [4], in the $\tau_{had}\tau_{had}$ channel, background originating from quark- and gluon-initiated jets that are misidentified as hadronic $\tau$ decays in processes other than multi-jet production are estimated using simulation. Simulated jets misidentified as hadronically decaying $\tau$ are weighted by the fake rates. This not only ensures the correct fake yield, but enhances the statistical precision of the estimate, since the events which were failing the identification are not removed. The fake rate FR$_{tau-ID}$, for both leading and sub-leading $\tau$ candidate, is defined as the ratio of number of $\tau$ candidates that pass a $\tau$ identification score cut, $N^{pass,tau-ID}$, to the total number of $\tau$ candidates, $N^{total}$, the only difference is that leading $\tau$ candidates should also pass the single $\tau$ trigger requirement. It is measured in W+jets events as:

$$FR_{tau-ID}(p_T, N_{track}) = \frac{N^{pass,tau-ID}(p_T, N_{tarck})}{N^{total}(p_T, N_{tarck})}|_{W+jets} \quad (3.2)$$

Events in the W+jets control region are selected by a single-muon trigger with a $p_T$ threshold of 36 GeV. They are required to contain one isolated muon with $p_T > 40$ GeV matched to the object that passed the trigger. There must be no additional muons or electrons and at least one $\tau$ candidate with opposite charge to the muon. Later all simulated events for background processes other that multi-jet events, processed by the search analysis, are assigned a weight given by:

$$w_{MC} = \prod_{i \in \{lead, sub-lead\}} (1 - \delta^i[1 - FR^i_{tau-ID}(p^i_T, N^i_{track}]) \quad (3.3)$$

where $\delta^i$ is 1 if the $\tau$ candidate originates from a jet and 0 otherwise.

The uncertainty in the fake rates used to weight simulated non-multi-jet events in the $\tau_{had}\tau_{had}$ channel is dominated by the limited statistics of the fake regions and can reach 40%. The requirement of opposite charge between muon and $\tau$ candidate enhance the quark-initiated jet composition. To evaluate the systematic uncertainty from applying these fake rates to simulate samples with different jet origin, the fake rate are also calculated for the same sign events which have higher fraction of gluon initiated jets, resulting in lower fake rates as shown in Figure 2. A relative uncertainty of 60% is assigned to cover the range of the measured fake rates for events with opposite- or same-sign $\tau$ candidates. The uncertainty is omitted for W+jets events as they are expected to have the same jet composition as events in the control region.

**Figure 2:** Tau-ID fake-rate measured in W($\mu\nu$)+jets data events, shown separately for muons of the same sign and opposite sign to the reconstructed $\tau$ candidate. Opposite-sign events are depicted by black circles and same-sign events by blue stars. The systematic uncertainty covers differences due to jet composition and is added to the statistical uncertainty [5].

## 4. Conclusion

The most commonly used approaches for the estimation of the contamination from misidentified hadronic $\tau$ decays in ATLAS analyses are the fake factor and fake rate methods. The fake factor method is universal and precise. It estimates entire background from all sources, however in this method the relative quark/gluon composition of jets in control and signal regions needs to be known. The fake rate method is a semi-data-driven approach and is applied to Monte Carlo samples, hence it can be used for estimating backgrounds which are modelled by MC. In this method the statistical precision of the estimate is enhanced, since the events failing the $\tau$ identification are not discarded. The choice of the optimal strategy for the determination of the background, among the two described above, depends on the specific analysis and it can even be a combination of them.

## References

[1] ATLAS Collaboration, 2008 JINST 3 S08003.

[2] ATLAS Collaboration, JHEP 1809 (2018) 139.

[3] ATLAS Collaboration, Eur.Phys.J. C75 (2015) no.7, 303

[4] ATLAS Collaboration, JHEP 1801 (2018) 055.

[5] ATLAS Collaboration, JHEP 1507 (2015) 157.

[6] ATLAS Collaboration, Phys. Rev. Lett. 121, 191801 (2018).

[7] ATLAS Collaboration, Phys. Rev. D 98, 092010 (2018).