

## CHAPTER 6

### EVENT SELECTION

This chapter will describe the event selection for the signal regions of this analysis. The conceptual motivation for the categorization is discussed in Section 6.1. The details of the selection requirements are enumerated in Section 6.2. The optimization of the sensitivity is described in Section 6.3. A summary is provided in Section 6.4.

#### 6.1 Event selection categorization motivation

As introduced in Section 3.2, any process that is significantly impacted by the 26 WCs listed in Table 2.2 is considered to be a signal process for this analysis. The signal processes comprise  $t\bar{t}H$ ,  $t\bar{t}l\nu$ ,  $t\bar{t}l\bar{l}$ ,  $t\bar{t}lq$ ,  $tHq$ , and  $t\bar{t}t\bar{t}$ .

As outlined in Chapter 1, this analysis focuses on multilepton signatures of the  $t(\bar{t})X$  processes; events with two same-sign leptons are categorized as  $2\ell_{ss}$ , events with three leptons are categorized as  $3\ell$ , and events with four or more leptons are categorized as  $4\ell$ . All events are also required to contain jets, with one or more of them b-tagged. The events are further subdivided based on b-tag multiplicity, jet multiplicity, the lepton charge sum, and whether or not there is a same-flavor-opposite-sign pair of leptons with an invariant mass close to the Z mass. Aiming to isolate subsamples of events with distinct admixtures of each contribution, these subdivisions improve the sensitivity of the analysis by allowing the effects of the WCs (which impact each signal process differently) to be distinguished more distinctly.

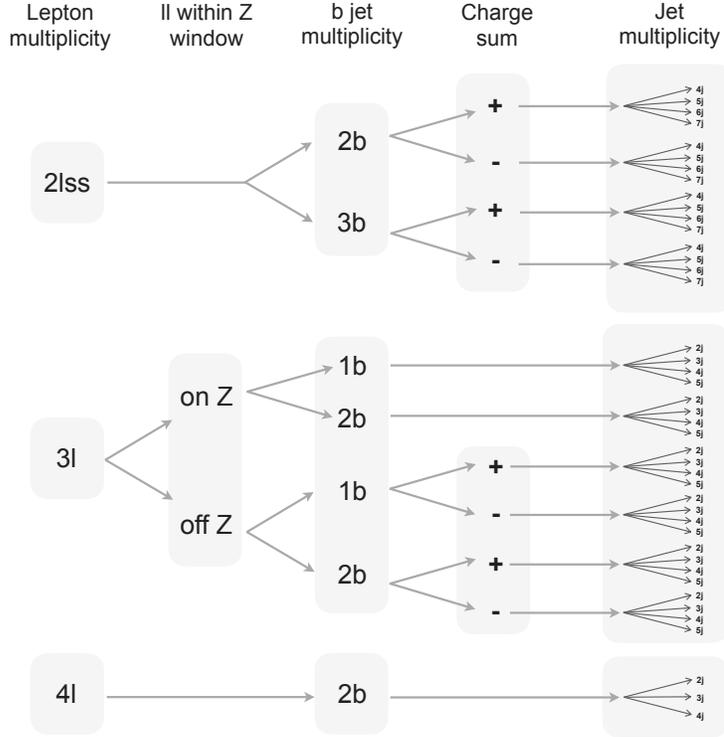


Figure 6.1. Summary of the event selection subdivisions. The details for the selection requirements are listed in Sections [6.2.1](#), [6.2.2](#), and [6.2.3](#).

The categorization scheme for this analysis is summarized in Figure [6.1](#) and the motivation for each subdivision is described below:

- Same-flavor-opposite-sign pair with an invariant mass close to the Z: This categorization helps to distinguish  $t\bar{t}l\bar{l}$  and  $t\bar{t}lq$  from the other processes that do not involve a Z. This distinction is also important for the identification of effects from the 2-quark-2-lepton WCs. These WCs are associated with vertices that directly produce a same-flavor-opposite-sign pair of leptons (without an intermediate Z), so the  $3\ell$  off-Z categories provide important sensitivity to these effects. The on-Z vs off-Z distinction is not applied for  $2\ell ss$  because  $t\bar{t}l\bar{l}$  and  $t\bar{t}lq$  do not naturally lead to  $2\ell ss$  final states.
- Multiplicity of b-tags: In the  $3\ell$  category, this distinction helps to separate single top processes from  $t\bar{t}X$  processes. In the  $2\ell ss$  category, we use a high b-tag multiplicity selection to isolate a subsample that is enriched in  $t\bar{t}t\bar{t}$ .

- Charge sum of leptons: In categories that are naturally populated by  $t\bar{t}l\nu$  ( $2lss$  and  $3l$  off- $Z$ ), we distinguish events with a positive charge sum from events with a negative charge sum. This helps to distinguish  $t\bar{t}l\nu$  from  $t\bar{t}H$  (since the LHC is a pp collider,  $t\bar{t}l^+\nu$  will have a larger cross section than  $t\bar{t}l^-\nu$ , so  $t\bar{t}l\nu$  will contribute more significantly to the + categories while  $t\bar{t}H$  will populate the + and - categories symmetrically).
- Multiplicity of jets: This helps to distinguish processes which tend to produce more jets (e.g.  $t\bar{t}H$  and  $t\bar{t}t\bar{t}$ ) from processes which tend to produce fewer jets (e.g.  $t\bar{t}l\nu$ ,  $t\bar{t}l$ , and the single top processes).

## 6.2 Event selection category requirements

The details of the event selection requirements for the  $2lss$ ,  $3l$ , and  $4l$  categories will be enumerated in Sections [6.2.1](#), [6.2.2](#), and [6.2.3](#). In addition to the category-specific requirements, some requirements are applied to all categories:

- To remove background contributions from light resonances, events that contain a pair of leptons (passing the loose lepton requirements outlined in Section [5.3](#)) with an invariant mass less than or equal to 12 GeV are vetoed.
- Events with anomalously large  $p_T^{\text{miss}}$  (caused by e.g. detector noise) are removed with the CMS MET group  $p_T^{\text{miss}}$  Filters [46](#).
- To ensure that electrons are well measured, some requirements are implemented on top of the requirements outlined in Chapter [5](#): the number of missing hits (the `lostHits` property) is required to be 0, and the electrons are required to pass the conversion veto (the `convVeto` property must be `True`).
- For simulated events, the leptons in the event are further required to pass MC truth requirements to ensure that they are prompt (the `genPartFlav` property is required to be 1 or 15). This ensures that we do not include any MC contributions in the nonprompt estimation (which is estimated with a data-driven approach, as described in Chapter [8](#)).

### 6.2.1 The $2lss$ category

The  $2lss$  category requires at least two leptons to pass the fakeable lepton object selection requirements defined in Section [5.3](#). Ordered by cone- $p_T$ , the leading two leptons must also pass the tight selection requirements defined in Section [5.3](#) and

these must be the only leptons in the event that pass the tight requirements (i.e. there must not be more than two tight leptons in the event). The leptons are required to have the same charge. The  $p_T$  of the leading lepton is required to be greater than 25 GeV, and the  $p_T$  of the second lepton is required to be greater than 15 GeV.

For both electrons and muons, additional requirements are implemented to ensure that the charges are well measured. This helps to reduce the charge flip contribution. These requirements are applied on top of the object selection requirements outlined in Chapter 5. For muons, the `tightCharge` property is required to be greater than or equal to 1 (this requires that the ratio of the uncertainty on the  $p_T$  to the  $p_T$  is less than 0.2). For electrons, the `tightCharge` property is required to be greater than or equal to 2 (this requires that the multiple methods of calculating the sign of the electron charge [47] yield consistent results). Furthermore, events where the two tight leptons are electrons with an invariant within 10 GeV of the Z mass are vetoed; this also helps to reduce the contribution from charge flip events. For simulated samples, electrons are required to pass MC truth requirements to ensure that the charge has not been mismeasured (the electron's `matched_gen_gen_pdgId` is required to have the same sign as the electron's `pdgId`); this ensures that MC does not contribute to the charge flip background (which is estimated with a data-driven approach, as described in Chapter 8).

At least four jets (with  $p_T > 30$  GeV and  $|\eta| < 2.4$ ) are required. Of these jets, at least two must pass the loose working point for the DeepJet algorithm, and at least 1 of these must also pass the medium working point for the DeepJet algorithm; i.e. there must be at least two loose b-tagged jets, at least one of which must also be a medium b-tagged jet. A subcategory is defined for events with at least three medium b-tagged jets; this allows us to isolate a collection of events that is relatively enriched in  $t\bar{t}t\bar{t}$ .

### 6.2.2 The $3\ell$ category

The  $3\ell$  category requires at least three leptons to pass the fakeable lepton object selection requirements outlined in Section 5.3. Ordered by cone- $p_T$ , the leading three leptons must also pass the tight requirements defined in Section 5.3 and these must be the only leptons in the event that pass the tight requirements (i.e. there must not be more than three tight leptons in the event). The  $p_T$  of the leading lepton is required to be greater than 25 GeV, and the  $p_T$  of the second lepton is required to be greater than 15 GeV. If the third lepton is an electron,  $p_T > 15\text{GeV}$  is required; if the third lepton is a muon,  $p_T > 10\text{GeV}$  is required.

At least 2 jets (with  $p_T > 30\text{GeV}$  and  $|\eta| < 2.4$ ) are required. Of these jets, at least one must pass the medium working point for the DeepJet algorithm; i.e. there must be at least one medium b-tagged jet. The events are subdivided based on whether there is exactly medium b-tagged jet, or more than one medium b-tagged jet. This helps to distinguish between the single top processes and the  $t\bar{t}X$  processes.

### 6.2.3 The $4\ell$ category

The  $4\ell$  category requires at least four leptons to pass the fakeable lepton object selection requirements defined in Section 5.3. Ordered by cone- $p_T$ , the leading four leptons must also pass the tight requirements defined in Section 5.3. The  $p_T$  of the leading lepton is required to be greater than 25 GeV, and the  $p_T$  of the second lepton is required to be greater than 15 GeV. For the trailing leptons, the requirements are  $p_T > 15\text{GeV}$  for electrons and  $p_T > 10\text{GeV}$  for muons.

At least four jets (with  $p_T > 30\text{GeV}$  and  $|\eta| < 2.4$ ) are required. Of these jets, at least two must pass the loose working point for the DeepJet algorithm, and at least 1 of these must also pass the medium working point for the DeepJet algorithm; i.e. there must be at least two loose b-tagged jets, at least one of which must also be a medium b-tagged jet.

### 6.3 Optimization studies

Based on the multiplicity of leptons, multiplicity of b-jets, sum of lepton charges, and the invariant mass of dilepton pairs, the binning described in Sections [6.2.1](#), [6.2.2](#), and [6.2.3](#) results in 11 independent categories. Further subdividing the categories by the jet multiplicity leads to 43 independent categories. We refer to this binning as inclusive  $N_{jet}$  binning, where “inclusive” refers to the fact that  $N_{jet}$  is the finest subcategorization. The inclusive  $N_{jet}$  binning provides good sensitivity to many of the WCs studied in this analysis, and is similar to the categories used in Ref. [\[10\]](#) (the predecessor to this analysis, which made use of only 2017 data). However, this analysis makes use of more than three times the data that was available for Ref. [\[10\]](#), and the additional statistics allow a more differential approach to be applied. In order to gain additional sensitivity to EFT effects, we bin the events in each of the 43 categories according to a kinematical variable.

In principle, a kinematic distribution could be fit for each category, resulting in  $43 \cdot n$  total bins, where  $n$  is the number of bins in each kinematic distribution, assuming  $n$  is the same for all categories. However, it is not necessary to use the same kinematic binning in every category, as the binning may be adjusted to account for varying statistics. In the limit where a single bin is used for the differential variable, the inclusive  $N_{jet}$  binning would be recovered for the given category. For the categories defined in this analysis, we have found that using 4 or 5 differential bins provides a good increase in sensitivity while maintaining reasonable statistics in each category.

While the same kinematical variable may be used across all categories, it would also be possible to use a different variable in every category; any case in between these two extremes may also be implemented. However, it is important to keep in mind that the WCs cannot be fully isolated or associated with a single category, so it is not possible to choose a particular variable for each WC. Nevertheless, since some WCs

impact certain categories more strongly than others, it is possible to target WCs by choosing specific variables for categories that may be particularly sensitive to the given WCs. The goal is thus to find a variable that provides good sensitivity to all WCs, or to find a combination of different variables (to use in different categories) that may improve the sensitivity to the target WCs without degrading the sensitivity to other WCs. To assess the sensitivity provided by a given variable, we compare the limits obtained from the differential fit against the limits obtained from the more inclusive  $N_{jet}$  fits. For these optimization studies, asimov data (i.e. simulated data that is equal to the SM prediction) was used; signals, backgrounds, and systematic uncertainties were included in the fit.

Since the contributions of many EFT vertices scale with energy, a variable that is related to the center of mass energy of the collision may provide generally good sensitivity for many WCs. For this reason, we investigated several variables related to the overall energy of the event. For example, we studied the  $S_T$  variable, which is defined as the scalar sum of the  $p_T$  of all of the leptons and jets in the event. It is also interesting to consider variables associated with the highest  $p_T$  object in an event, as it is possible that these objects could be associated with the EFT vertex in the process. For example, we can consider the leading lepton  $p_T$ , or the  $p_T$  of the leading object (lepton or jet). It is also interesting to consider variables that combine multiple high- $p_T$  objects. For example, we may consider the  $p_T$  of the leading pair of objects from the collection of leptons and jets in the event, a variable we refer to as  $p_T(lj)_0$ . Testing our sensitivity to these and other similar variables, we found an improvement of approximately 50% for most of our WCs (compared to a the more inclusive approach of fitting to the  $N_{jet}$  distributions).

For the on-Z categories, it is also interesting to consider the  $p_T$  of the same flavor opposite sign lepton pair. Using this variable (referred to as  $p_T(Z)$ ) for all on-Z categories while using one of the other variables (e.g.  $p_T(lj)_0$ ) for all other categories,

we observe improvements for some WCs (most notably WCs from among the “two-heavy” category); however, we observed a significant decrease in sensitivity to  $c_{Qq}^{31}$  and  $c_{Qq}^{38}$ . Members of the “two-light-two-heavy” category of WCs, these two WCs are unique in that they give rise to q-q'-t-b vertices. These WCs can thus contribute to  $3\ell$  on-Z 2b 2j (and 3j) final states. In cases where  $c_{Qq}^{31}$  and  $c_{Qq}^{38}$  contribute to these signatures, the Z boson is not part of the EFT vertex, so it is not optimal to use  $p_T(Z)$  as the differential variable, thus explaining the loss in sensitivity to these WCs when using  $p_T(Z)$  for all onZ categories. To mitigate this effect, we tested the scenario where  $p_T(Z)$  was used for all on-Z categories except the on-Z 2b 2j/3j final states. This modification indeed mitigated the degradation of the sensitivity to  $c_{Qq}^{31}$  and  $c_{Qq}^{38}$ , while still providing improvements in sensitivity for the “two-heavy” WCs. Figure 6.2 summarizes the sensitivity observed for several of the differential variables considered during the optimization studies.

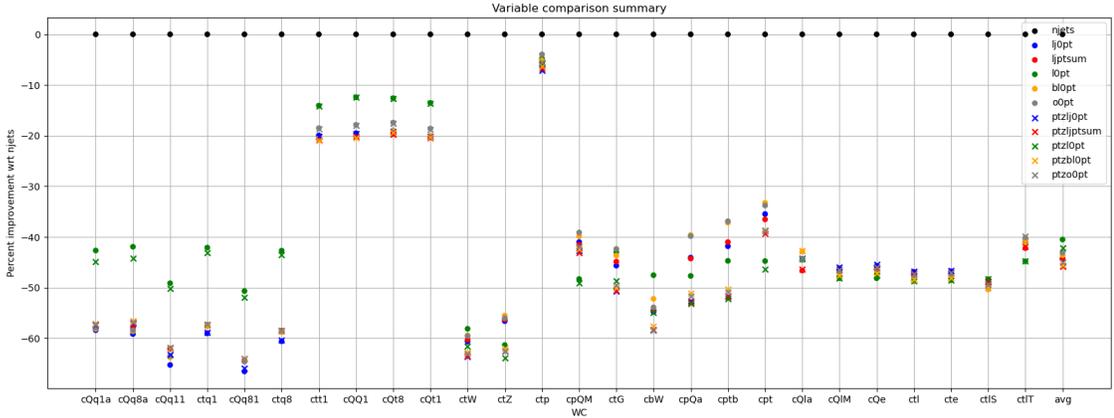


Figure 6.2. Summary of the sensitivity provided by fits to various differential distributions. The  $y$  axis represents the percent improvement with respect to the inclusive  $N_{jet}$  fit (based on the widths of the  $2\sigma$  confidence intervals from fits to asimov data).

Most of the differential variables in Figure 6.2 performed similarly to each other, providing a significant improvement in sensitivity compared to the inclusive  $N_{jet}$  fit. In general, the fits in which we use  $p_T(Z)$  for the selected on-Z categories (indicated with “x” shaped markers in Figure 6.2) provide better sensitivity than the cases where the same variable is used in all categories (shown with circular markers in Figure 6.2). The case where  $p_T(Z)$  is used for the selected on-Z categories and  $p_T(lj)_0$  is used for all other categories (i.e. “ptzlj0pt” in the plot, denoted with blue “x” shaped markers) shows consistently good sensitivity across all categories, as well as providing the best sensitivity to the “two-light-two-heavy” WCs. For these reasons, we choose to use the  $p_T(Z)$ - $p_T(lj)_0$  distributions as the kinematic distributions for this analysis. For the fit, we use 4 bins in  $p_T(lj)_0$  and 5 bins in  $p_T(Z)$ , resulting in 178 total bins.

#### 6.4 Event selection summary

Targeting the multilepton signatures of  $t(\bar{t})X$  processes, the event selection categories in this analysis constitute  $2lss$ ,  $3l$ , and  $4l$ . The events are further subdivided into 43 unique categories designed to differentiate as much as possible between the different  $t(\bar{t})X$  contributions. To gain additional sensitivity, the events in each of the 43 categories are binned according to a differential kinematical distribution, resulting in 178 total bins. The  $p_T(Z)$  variable is used for all of the on-shell Z categories, except for the 2 and 3 jet categories with 2 b-tagged jets; the  $p_T(Z)$  variable is thus used in 6 total categories. In the remaining 37 categories, the  $p_T(lj)_0$  variable is used. Binning the 43 analysis categories in terms of the  $p_T(lj)_0$  and  $p_T(Z)$  variables provides an improvement in sensitivity of a factor of about 2 (compared to the case where the 43 analysis bins are not further subdivided). Table 6.1 summarizes the selection requirements for each of the 43 categories in this analysis. Requirements separated by commas indicate a division into subcategories. The differential kinematical variable that is used in the category is also listed.

TABLE 6.1

## SUMMARY OF EVENT SELECTION CATEGORIES.

Category	Leptons	$m_{\ell\ell}$	b-tags	Lepton charge sum	Jets	Differential variable
2 $\ell$ ss 2b	2	No requirement	2	> 0, <0	4,5,6, $\geq$ 7	$p_T(lj)_0$
2 $\ell$ ss 3b	2	No requirement	$\geq$ 3	> 0, <0	4,5,6, $\geq$ 7	$p_T(lj)_0$
3 $\ell$ off-Z 1b	3	$ m_Z - m_{\ell\ell}  > 10\text{GeV}$	1	> 0, <0	2,3,4, $\geq$ 5	$p_T(lj)_0$
3 $\ell$ off-Z 2b	3	$ m_Z - m_{\ell\ell}  > 10\text{GeV}$	$\geq$ 2	> 0, <0	2,3,4, $\geq$ 5	$p_T(lj)_0$
3 $\ell$ on-Z 1b	3	$ m_Z - m_{\ell\ell}  \leq 10\text{GeV}$	1	No requirement	2,3,4, $\geq$ 5	$p_T(Z)$
3 $\ell$ on-Z 2b	3	$ m_Z - m_{\ell\ell}  \leq 10\text{GeV}$	$\geq$ 2	No requirement	2,3	$p_T(lj)_0$
3 $\ell$ on-Z 2b	3	$ m_Z - m_{\ell\ell}  \leq 10\text{GeV}$	$\geq$ 2	No requirement	4, $\geq$ 5	$p_T(Z)$
4 $\ell$	$\geq$ 4	No requirement	$\geq$ 2	No requirement	2,3, $\geq$ 4	$p_T(lj)_0$

Applying this selection to the data and simulated samples described in Chapter 3, Table 6.2 shows the resulting event yield in each category (summed over jet bins) for the data and for the SM prediction. The observed event yields are generally larger than the predicted event yields across all of the categories; overall, the observed yield (3927 events) is about 14% higher than the prediction (3440.0 events). However, it should be noted that there are significant systematic effects (described in Chapter 9) that can influence many bins in a correlated way. For this reason, the agreement between the prediction and the observation should not be judged until after a likelihood fit incorporating the systematic uncertainties (as described in Chapter 10) has been performed.

TABLE 6.2

EXPECTED SM YIELDS AND OBSERVATIONS IN THE ANALYSIS  
 CATEGORIES (SUMMED OVER JET CATEGORIES).

	$2\ell ss\ 3b\ -$	$2\ell ss\ 3b\ +$	$2\ell ss\ 2b\ -$	$2\ell ss\ 2b\ +$	$3\ell\ 1b\ -$	$3\ell\ 1b\ +$	$3\ell\ 2b\ -$	$3\ell\ 2b\ +$	$3\ell\ \text{on-Z}\ 1b$	$3\ell\ \text{on-Z}\ 2b$	$4\ell\ 2b$
tWZ	0.46	0.47	6.66	6.7	4.78	4.77	1.69	1.69	63.34	20.3	2.72
Diboson	0.1	0.25	12.08	15.79	30.45	29.64	2.02	3.19	338.24	34.35	4.81
Triboson	0.04	0.07	2.16	3.15	0.95	1.33	0.1	0.17	16.02	2.61	0.45
Charge flips	1.62	1.57	17.5	17.4	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Nonprompt	6.72	8.94	112.61	120.56	56.95	56.5	11.86	10.47	55.31	10.06	0.0
Conversions	0.96	0.84	11.45	9.57	3.37	2.93	3.25	3.07	0.78	0.68	0.01
Sum bkg	9.9	12.14	162.45	173.17	96.52	95.16	18.92	18.59	473.69	68.0	7.99
$t\bar{t}l\nu$	12.54	23.7	144.18	272.76	25.69	47.38	27.26	50.67	10.18	11.42	0.03
$t\bar{t}l\bar{l}$	12.29	12.31	119.02	119.64	51.58	51.04	48.44	49.78	320.72	295.81	40.22
$t\bar{t}H$	9.5	9.48	83.3	83.48	23.51	23.24	22.92	22.71	9.62	9.69	3.5
$t\bar{t}lq$	0.47	0.87	6.51	11.75	5.46	9.54	2.4	4.23	111.51	48.11	0.01
$tHq$	0.12	0.23	1.42	2.61	0.47	0.83	0.35	0.61	0.35	0.23	0.03
$t\bar{t}t\bar{t}$	9.61	9.53	7.58	7.46	0.87	0.85	4.88	4.92	0.21	1.3	0.55
Sum sig	44.53	56.12	362.01	497.71	107.59	132.88	106.23	132.93	452.59	366.56	44.35
Sum expected	$54 \pm 6$	$68 \pm 7$	$524 \pm 50$	$671 \pm 63$	$204 \pm 23$	$228 \pm 24$	$125 \pm 11$	$152 \pm 13$	$926 \pm 132$	$435 \pm 46$	$52 \pm 6$
Observation	71.0	68.0	608.0	781.0	233.0	270.0	148.0	158.0	1074.0	466.0	50.0