

# Lexical Normalization for Code-switched Data and its Effect on POS Tagging

**Rob van der Goot**  
IT University of Copenhagen  
robv@itu.dk

**Özlem Çetinoğlu**  
IMS, University of Stuttgart  
ozlem@ims.uni-stuttgart.de

## Abstract

Lexical normalization, the translation of non-canonical data to standard language, has shown to improve the performance of many natural language processing tasks on social media. Yet, using multiple languages in one utterance, also called code-switching (CS), is frequently overlooked by these normalization systems, despite its common use in social media. In this paper, we propose three normalization models specifically designed to handle code-switched data which we evaluate for two language pairs: Indonesian-English (Id-En) and Turkish-German (Tr-De). For the latter, we introduce novel normalization layers and their corresponding language ID and POS tags for the dataset, and evaluate the downstream effect of normalization on POS tagging. Results show that our CS-tailored normalization models outperform Id-En state of the art and Tr-De monolingual models, and lead to 5.4% relative performance increase for POS tagging as compared to unnormalized input.<sup>1</sup>

## 1 Introduction

Social media provide an invaluable source of information for natural language processing (NLP) systems. Its informative and spontaneous nature leads to many interesting phenomena, like non-standard words, spelling errors and abbreviations. One particularly challenging and interesting phenomenon is the use of multiple languages within the same utterance, which is also called code-switching (CS) (Gumperz, 1982; Myers-Scotton, 1995; Toribio and Bullock, 2012).

Because most NLP models are designed to process canonical and monolingual data, their performance drops enormously when having to process

social media data (Eisenstein, 2013). One solution to this problem is lexical normalization: the translation of non-standard (e.g. social media) text to its canonical form (Han and Baldwin, 2011). Previous work has shown that by standardizing the data, we can improve the robustness of NLP systems (Derczynski et al., 2013; Zhang et al., 2013; van der Goot and van Noord, 2017). Nevertheless these systems overlook code-switching. (1) shows a code-switched tweet (upper) and its normalization annotation (lower), taken from an Indonesian-English CS corpus (Barik et al., 2019) (Indonesian in bold). This example demonstrates that CS complicates normalization, because it can be unclear in which language to normalize (e.g., *ak* is normalized to *aku* ‘I’ in Indonesian. English-only normalization systems would probably normalize it to *ok*).

(1) **ak** . luv u :( till die  
**aku** . love you :( till die  
*‘I love you till (I) die’*

Recently, there has been an increasing interest in the automatic processing of CS data, however, there has not been much work on its lexical normalization. To the best of our knowledge, only Adouane et al. (2019) focus entirely on lexical normalization for CS data in their work. For other works, normalization is a preprocessing step for downstream tasks: chunking (Sharma et al., 2016), parsing (Bhat et al., 2017, 2018), or machine translation (Barik et al., 2019). These CS normalizers are either rule-based and language-specific (Barik et al., 2019) or combine (Hindi) back-transliteration and normalization (Sharma et al., 2016; Bhat et al., 2017, 2018) thus, they are not directly applicable to other lexical normalization datasets. In this work,

- We are the first to present open-source normalization models specialized for CS lexical normalization without any language-specific components.

<sup>1</sup>Source code is available at: <https://bitbucket.org/robvanderg/csmonoise>. The Turkish-German data is available at: <https://github.com/ozlemcek/TrDeNormData>

- We provide a novel lexical normalization dataset by annotating a Turkish-German Twitter corpus (Çetinoğlu, 2016). We also align existing annotation layers – language IDs (LID) and part-of-speech (POS) tags – to normalization annotations.
- We evaluate three CS normalization models on two language pairs (Turkish-German (Tr-De), Indonesian-English (Id-En)). For both datasets, CS models reach performance in a similar range as monolingual models reach on monolingual datasets.
- Our CS-tailored normalization models outperform Id-En state of the art and set the state of the art for the Tr-De dataset.
- We show that our proposed normalization models improve the performance of POS taggers. For a broad perspective, we employ a variety of taggers (CRF, BiLSTM, BERT).

## 2 Related Work

**Lexical normalization** Traditionally, social media normalization approaches can be broadly divided into two types. The first stream of work uses techniques borrowed from machine translation (Aw et al., 2006; Pennell and Liu, 2011; Ljubešić et al., 2016). The second stream is based on a classic spelling correction framework (noisy-channel models) (Han, 2014). Here, they often apply three steps, detecting which words need to be replaced, generating candidates, and ranking these candidates. Later, it became evident that a two-step approach is sufficient (Jin, 2015; van der Goot, 2019), and the detection step was alleviated by considering the original word as a normalization candidate.

The current state-of-the-art model for most languages is MoNoise (van der Goot, 2019), which is based on this two-step approach. A variety of modules are used for the generation of candidates. For the ranking, MoNoise complements features from the generation step with additional features, which are all combined in a random forest classifier that predicts the probability that a candidate is a ‘correct’ candidate. MoNoise is described in more detail in Section 4.2. More recently, sequence-to-sequence models (Lourentzou et al., 2019) and contextual embeddings (Muller et al., 2019) have been used for the lexical normalization task. These approaches have been shown to reach performances close to MoNoise on English benchmarks.

Like most NLP tasks, most research on normalization has been done on English datasets (Han and Baldwin, 2011; Baldwin et al., 2015). However, there has been some efforts on other languages, where usually only one language is considered, we refer to Sharf and Rahman (2017) and van der Goot (2019) for an overview of available resources.

### Processing of code-switched social media data

Early work on normalizing CS data focused on Hindi-English, as part of pipelines to achieve downstream tasks (Sharma et al., 2016; Bhat et al., 2017, 2018). As Hindi is Romanized in datasets and additional Hindi resources are in the Devanagari script, they include back-transliteration into the normalization step, thus defining the task beyond the scope of this paper. Nevertheless, all systems report a positive impact of normalization on their final task.

More recently, Barik et al. (2019) experiment on normalization for Indonesian-English. They use a rule-based approach supplemented by clusters derived from word embeddings, and show that normalization can be used to improve machine translation. Adouane et al. (2019) instead propose to use sequence-to-sequence models for normalizing Algerian Arabic data mixed with Modern Standard Arabic, French, Berber, and English. They show that their edit distance-based token-level aligner helps improve normalization.

When annotating the Tr-De dataset for normalization, we also adapted its POS tags (see Section 3.1). This gives us the opportunity to apply POS tagging as extrinsic evaluation. Besides research on Hindi-English that combines normalization and back-transliteration, most work either use normalization to improve tagging performance of monolingual social media data (Derczynski et al., 2013; van der Goot et al., 2017), or on POS tagging of CS data without normalization (AlGhamdi et al., 2016; Soto and Hirschberg, 2018). In this work, we combine these angles.

Because some of our proposed normalization models depend on language labels, we require a word-level language identification system. There is a wide variety of approaches used for this task, where early systems mostly used CRFs (Sequiera et al., 2015; Molina et al., 2016). More recently, neural networks based approaches have shown superior performance for this task (Zhang et al., 2018). We opt for three different architectures to observe the effect of the quality of language identification on normalization (Section 4.1).

	Seg+CS			Norm		Tok+Anon
Tokens	Semester§da	-yim		Semesterdayım		semesterdayım
LID	Mixed	TR	⇒	Mixed	⇒	Mixed
POS	NOUN	VERB		VERB		VERB

Figure 2: Mapping LID and POS tags from Seg+CS to Norm to Tok+Anon for the mixed word *Semesterdayım* ‘I am in semester’.

```
Raw: @Erkan1903 nerdee 3 semesterdayım dha.
Tok+Anon: @username|nerdee|3|semesterdayım|dha|.
Norm @username|Nerde|3|Semesterdayım|daha|.
Seg+CS: @username Nerde 3. Semester§da -yim daha .
```

Figure 1: Different annotation layers of a tweet from the Tr-De corpus, meaning ‘No way, I am still in the 3rd semester’. The German part is in bold. Raw: downloaded tweet; Tok+Anon: after tokenization and anonymization; Norm: after normalization; Seg+CS: after segmentation (e.g. the Turkish copular *-yim*), and CS boundaries (§) in Mixed tokens. The token alignment and normalization tasks are carried out on the Tok+Anon and Norm pairs.

### 3 Data

In this section we first describe the design decisions of the novel Turkish-German dataset, then we compare some basic statistics together with the existing Indonesian-English dataset (Barik et al., 2019).

#### 3.1 Turkish-German code-switched normalization corpus

We use the Turkish-German Twitter corpus from Çetinoğlu (2016) in our experiments. It consists of 17K tokens as 1,029 tweets. The raw tweets of the corpus have undergone three main steps of alternations after the collection: tokenization, normalization, and segmentation.<sup>2</sup> In addition, usernames and URLs are anonymized as `@username` and `[url]` respectively, and intra-word CS boundaries are marked in Mixed tokens with §. Each alternation layer is exemplified on a sentence from the corpus in Figure 1.

The Seg+CS layer is annotated with language IDs and POS tags (Çetinoğlu and Çöltekin, 2016). The LID tag set consists of TR (Turkish), DE (German), Lang3 (third language), Mixed (intra-word CS), NE (named entity), Ambig (both Turkish and German and cannot be disambiguated in given context), Other (punctuation, numbers, URLs, emoticons, symbols). Additionally, named entities are tagged with their language label next to the NE tag,

<sup>2</sup>Morphosyntactic split of words into subwords, cf. (Çetinoğlu and Çöltekin, 2016) for details.

e.g. ‘Germany’ is annotated in the corpus as follows depending on the language: *Almanya* NE.TR, *Deutschland* NE.DE, *Germany* NE.Lang3. The POS annotation adopts the Universal Dependencies (UD) tag set (Nivre et al., 2016).

**Preprocessing for normalization** The original version of the corpus has only the Raw and Seg+CS layers and only tweet-level alignment between them. As our work focuses only on normalization we created the intermediate layers Tok+Anon and Norm that leave out other tasks. Since MoNoise requires word-aligned annotations, we also provided these alignments.

We anonymized and tokenized the raw tweets to achieve the Tok+Anon layer. For tokenization, we use a slightly modified version of `twookenize.py`<sup>3</sup> (O’Connor et al., 2010). To obtain the Norm layer, we merged back segmented tokens and removed CS boundaries on the Seg+CS layer.

After this stage, we aligned Tok+Anon and Norm on the token level automatically using Giza++ (Och and Ney, 2003). We parsed the resulting alignment files to align the actual tokens and corrected them manually. There are 15,715 1:1, 520 1:n, and 147 n:1 alignments.

**LID and POS alignment** The existing LID and POS tags are on the Seg+CS layer; since we base our experiments on the Tok+Anon layer, we need to map the annotations. This is done in two steps following the Seg+CS ⇒ Norm ⇒ Tok+Anon order. Due to segmentation merges in the first step, and 1:n and n:1 token alignments in the second step, there are non-trivial LID and POS alignments.

Figure 2 demonstrates a segmented word in the first column. The first segment *Semesterda* ‘in semester’ is Mixed with German *Semester* and Turkish locative case marker *da*. The second segment is the Turkish copular *-yim* ‘I am’. Their POS tags are NOUN and VERB, respectively. When segmentation is undone in the second column (Norm), their LID and POS are merged too. If two tokens

<sup>3</sup>[github.com/brendano/tweetmotif](https://github.com/brendano/tweetmotif)

	#words	%norm	% split	%merge	CMI
Id-En	18,758	14.13	1.33	0.17	28.20
Tr-De	13,217	25.97	3.01	1.04	22.44

Table 1: Descriptive normalization and code-switching statistics on the training split of the datasets. CMI is the code-mixing index (Das and Gambäck, 2014), averaged over all training sentences. %norm reflect the percentage of words which is normalized.

have the same LID, the merged token takes the same LID. If they are different, the resulting token is `Mixed`, as in the example.

POS tag merging rules can get more complicated, therefore, we used a heuristic that favors the POS tag of the second token in most cases.<sup>4</sup> When a `NOUN` segment is merged with a `VERB` segment, as in Figure 2 (`Seg+CS`  $\Rightarrow$  `Norm`), the merged token is assigned a `VERB` POS tag. For the `Norm`  $\Rightarrow$  `Tok+Anon` mapping, the alignment is 1:1, thus LID and POS are directly carried over.

### 3.2 Dataset characteristics

Besides the data described in the previous section, we use the Indonesian-English (Id-En) data from (Barik et al., 2019). The Id-En data is only annotated with language IDs and uses three labels: `ID`, `EN`, `UN` (Unspecified), whereas the Tr-De includes 12 labels (Section 3.1). To simplify the models and improve comparability, we map the language labels of the Tr-De dataset to `TR`, `DE` and `UN`. Named entities are mapped to their respective language tags, e.g. `NE.DE` to `DE`. `Mixed` tokens are mapped to `DE` as they are German words with Turkish inflection. `Lang3`, `Ambig` and `Other` are mapped to `UN`.

We divide both datasets into a train and test split (80-20%), and omit a development set due to small sizes. Since we want to leave test set out in analyses, we opt for 10-fold cross-validation on the training split of the data in experiments. Statistics of the training splits of the datasets are shown in Table 1. The datasets are relatively small, but a high ratio of words is normalized, including a high percentage of splits and merges. The percentage of in-vocabulary words is especially low in the Tr-De data, which is mainly due to the morphological richness of Turkish. The code-mixing index (CMI) (Das and Gambäck, 2014) indicates

<sup>4</sup>Turkish is agglutinative. Segmentation often happens by splitting derivational suffixes that bear the final POS tag.

Source	Indonesian	English	Turkish	German
Wikipedia	75	2,162	55	776
Twitter	510	5,018	203	89

Table 2: Size of raw data (in million words) from both data sources.

the (average of the) amount of words not written in the majority language for each sentence. The relatively high CMI for both datasets indicates a high frequency of code-switching occurs in the data.

In both datasets there are a small amount of sentences without normalization (8 and 76 for respectively Id-En and Tr-De), which might be desirable for evaluation of (over)normalization, as in a real-world setup one also does not know beforehand whether normalization is necessary. In more than half of the sentences the number of normalized words is larger than 3. Furthermore, there are some sentences (5-10 per dataset) with a very high normalization ratio (>70%), which are all in capitals.

### 3.3 Monolingual Data

Our baseline model (MoNoise) exploits monolingual data from both the source and the target domain (canonical data) to train word embeddings and estimate n-gram probabilities. To this end, we utilize Wikipedia dumps from 01-01-2020 and random tweets collected throughout 2012 and 2018 from the Twitter API, filtered by the FastText language classifier (Joulin et al., 2017). We tokenized this data based on whitespaces, and removed all duplicate sentences/tweets. The sizes of the collected raw datasets are shown in Table 2.

## 4 Models

In this section we describe the models used for word-level language identification (4.1), lexical normalization (4.2) and POS tagging (4.3).

### 4.1 Word-level language identification

We treat language identification as a sequence labeling task where the label of each word is a language ID. We evaluate three sequence labeling libraries: 1) MarMoT (Mueller et al., 2013), a higher-order conditional random fields tagger 2) Bilty (Plank et al., 2016), a BiLSTM tagger, also incorporating character level information 3) a BERT-based (Devlin et al., 2019) tagger named MaChAmp (van der Goot et al., 2021). For Bilty, we project polyglot embeddings (Al-Rfou et al., 2013) of each

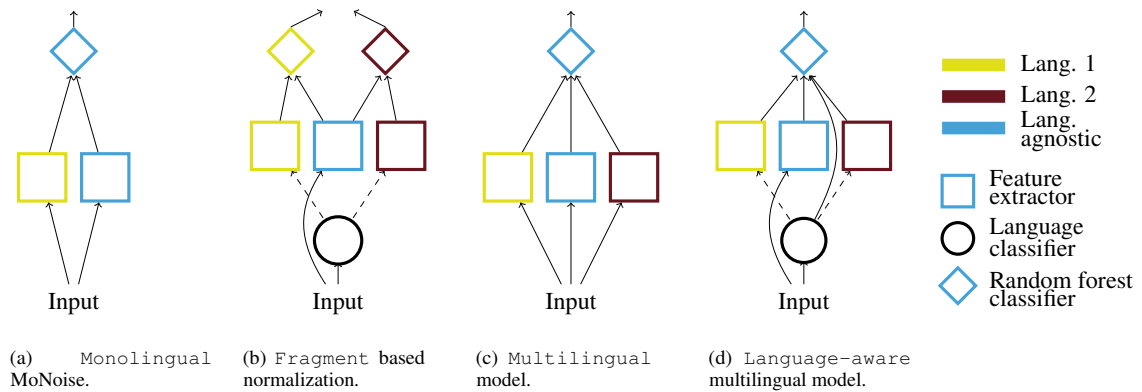


Figure 3: Overview of the different proposed variations of MoNoise. Dashed lines mean that only one of the two paths is taken, decided by the language identification. For model (a), there can be two versions, one with features from Lang. 1 (shown here) and one based on Lang. 2.

language of the language pairs to the same space using MUSE (Lample et al., 2018), whereas for MaChAmp, we use multilingual BERT.<sup>5</sup> We use the default settings for all toolkits.

## 4.2 Normalization

We choose to use MoNoise (van der Goot, 2019) as a baseline and starting point for our proposed models for two main reasons: 1) Normalization annotation for code-switched data is scarce, and MoNoise is specifically strong in low-resource setups because of its dependence on external resources (generated from raw data); 2) It is the only normalization model that has shown to be effective in multiple languages. Below we first introduce the standard monolingual MoNoise model, and then all the proposed extensions which are focused on code-switched data. A schematic overview of all models is shown in Figure 3.

### Monolingual (Figure 3a)

MoNoise consists of two parts, a candidate generation step and a candidate ranking step. For the generation of candidates, a spelling correction system (Aspell),<sup>6</sup> word embeddings and a dictionary based on the training data are used. Features from these modules are then supplemented with n-gram probabilities based on Wikipedia and Twitter data and other features indicating whether a word is present in the Aspell dictionary, whether it contains an alphabetical character, the length of a candidate compared to the original word, and whether it starts with a capital. For the novel proposed models, we will split up the features based on whether they

require *language-specific* resources (spelling correction, word embeddings and n-grams features; yellow and red in Figure 3), or whether they are *language-agnostic* (all other features; blue in Figure 3). For the ranking of the candidates a random forest classifier (Breiman, 2001) is used, which predicts the probability whether a candidate is correct. An obvious disadvantage when applying monolingual MoNoise on CS-data is that many features are language-specific (e.g. spelling correction, word embeddings, n-grams), which is sub-optimal for tokens from another language. Since our datasets and evaluation include capitals, we use the version of MoNoise including capitalization handling (van der Goot et al., 2020).

### Fragments (Figure 3b)

The baseline model has the deficiency that it has the language-specific features only for one language, while normalizing texts for two languages. An intuitive way of improving this model would be to split up the input data into monolingual fragments, and train two separate monolingual models. The fragments are created by splitting the data on every CS point, where words with the UN label are converted to the label of the previous word. This setup has the advantage that the normalization model itself does not need any adaptation, and it can thus be used with any normalization model. The disadvantages are that it is dependent on a language label, two separate classifiers have to be trained and the context is interrupted.

### Multilingual (Figure 3c)

Instead of using two separate random forest classifiers, we can exploit both feature sets simultaneously in one classifier. This means that for every

<sup>5</sup>multi.cased.L-12\_H-768\_A-12

<sup>6</sup>www.aspell.net

language-specific feature, we now have two features. In this setup, the model is not explicitly informed about the language of input words, however, some of the features (especially n-gram probabilities) will have a very high correlation with this information. This model has the advantage that only one classifier has to be trained, and no language labels are necessary. It has the disadvantage that it uses more features for the classifier compared to the `Monolingual` and `Fragments` models, which increases the complexity of the classification.

### Language-aware (Figure 3d)

Some of the language-specific features of the `Multilingual` model will be rather superfluous for words in the other language. For example, it will search for Turkish words in German word embeddings, and also use n-gram counts based on the German Wikipedia. To avoid this, we can use only one copy of each language-specific feature, and generate them based on the language label (the same language labels as in the `Fragments` model are used). More concretely, this means that for a German word, we will generate uni-gram probabilities based on German data, whereas for Turkish we will use Turkish data; these are then modeled as one feature in the model. On top of this, we also add a feature that indicates which language a word belongs to. There might be some mismatches in the importance of features because different data sources and languages are used. Because the language label is known, and a random forest classifier can model feature interactions intrinsically (Breiman, 2001), these mismatches should not be problematic. This model has the advantage that the number of features stays almost the same as in the `Monolingual` model (+1, the language ID), but a disadvantage is that it requires language labels.

### 4.3 POS tagging

For POS tagging, we examine the same three sequence labeling systems as used for language identification (Section 4.1): `MarMoT`, `Bilty` and `MaChAmp`. For each normalization setting, we normalize the input data, and use this normalized text as input for the POS tagger, which is trained on canonical data.

## 5 Evaluation

In this section we evaluate each of the three sub-tasks (LID, normalization, POS), where for the lat-

Model	Id-En	Tr-De
MarMoT	92.71	92.91
Bilty	*93.81	*94.31
MaChAmp	*95.17	*95.67

Table 3: Word level accuracies for language identification (10-fold).

ter two we also examine the effect of exploiting the prediction of the previous tasks. Unless mentioned otherwise, we report the results of 10-fold cross-validation on the training split of the data. For all experiments, we use a paired bootstrap test on the sentence level with 1,000 samples to test significance. For all results, we order the models by the complexity of the implementation as compared to `MoNoise` (first fragments, as the original model can be used as a black box, then multilingual because it does not need a language classifier, and finally the language-aware model). An \* next to results denotes a significant difference for  $p < 0.05$ , of a model always as compared to the previous model (corresponding to the previous column in Table 6, the previous row in other tables) for the same data.

### 5.1 Language identification

Results for the language identification task are reported in Table 3. Unsurprisingly, the performances are in line with the chronological order of the introduction of the systems, and their computational complexity. It should be noted that for `MaChAmp` we used pre-trained embeddings which were trained on the largest amount of external data. When inspecting the performance per language label, we saw that the ‘UNspecified’ is by far the most difficult. Even though this class contains punctuation, it also contains many harder cases, where a word belongs to any language other than `Lang1` and `Lang2`, or when the annotator is uncertain. Barik et al. (2019) use a conditional random fields classifier with a variety of features for this task, and report 90.11 accuracy for the full Id-En dataset in a 5-fold cross-validation setting. Which, despite differences in data splits, confirms that our results are competitive.

### 5.2 Normalization

For lexical normalization, a wide variety of evaluation metrics is used in the literature, ranging from accuracy (Han and Baldwin, 2011), F1 score (Baldwin et al., 2015) and precision over out-

Model	Id-En	Tr-De
LAI	73.24	74.03
MFR	*88.35	*78.57
Monolingual-lang1 (Tr/Id)	*94.76	*79.81
Monolingual-lang2 (De/En)	*94.31	80.58
Fragments	94.73	*81.24
Multilingual	94.84	*81.74
Language-aware	94.79	81.68

Table 4: Normalization performance of the baselines and the proposed models (10-fold accuracy). For the models dependent on language labels, we used the labels predicted by MaChAmp.

of-vocabulary words (Alegria et al., 2013), to CER and BLUE score (Ljubešić et al., 2016). Because the word order is fixed in our task, and to ease interpretation of the results, we opt to use simple accuracy on the word level, where we consider all words (i.e., also the unnormalized words).

To interpret the scores, we include three baselines: 1) leave-as-is (LAI), which always outputs the original word, i.e. its accuracy is equivalent to the percentage of words that are not normalized 2) most-frequent-replacement (MFR), which uses the most frequent replacement from the training data for each word 3) monolingual MoNoise, which can be trained on either of the languages within a language pair (two models).

Results for the different models are compared in Table 4. For the Id-En dataset, the differences between all proposed models are small and not significant. Even the monolingual models perform remarkably well, and only small gains are observable when using the multilingual model. We also compared our results to Barik et al. (2019), using their evaluation metric as their model/output was not available. The metric is non-deterministic, as it uses accuracy over unique OOV words.<sup>7</sup> Nevertheless, our average estimated result for Multilingual is 69.83 for this metric, outperforming their score of 68.50.

For the Tr-De dataset, the scores are generally lower, indicating that this dataset (and perhaps language pair) is more difficult. Especially now, we can observe that the code-switched adaptations lead to substantially higher scores. To our surprise, Multilingual and Language-aware

<sup>7</sup>Which can be normalized differently dependent on context, we confirmed this with the authors.

Model	Id-En	Tr-De
Fragments (MarMoT)	94.66	80.77
Fragments (Bilty)	*94.71	*80.89
Fragments (MaChAmp)	94.73	*81.24
Fragments (Gold)	*94.81	81.71
Language-aware (MarMoT)	94.74	81.24
Language-aware (Bilty)	94.76	81.57
Language-aware (MaChAmp)	94.79	81.68
Language-aware (Gold)	94.90	82.18

Table 5: Effect of different language predictions on normalization models (10-fold accuracy).

perform on par, even though the multilingual model does not rely on language labels. Fragments performs significantly worse. This leads to the conclusion that language labels are not directly beneficial for lexical normalization (in this setup). In general, the performances are in a similar range as for monolingual datasets (van der Goot, 2019).<sup>8</sup>

**Model behavior** Besides the metrics reported in the table, we also examined precision and recall. Precision is generally much higher (1.1 to 3 times, see Appendix B) than recall especially for Tr-De, which is in line with previous observations (van der Goot, 2019). This means that the model is conservative and only replaces cases for which it is rather certain, which arguably is a desirable behavior.

**Effect of language predictions** To evaluate the effect of the language predictions, we run both the Fragments and the Language-aware models with all language predictions from Section 5.1 as well as the gold language labels. The results (Table 5) show that the performance of the language identification has a positive effect on the normalization performance. Although it is not significant in most cases, it should be noted that significance is only tested compared to the previous model.

**Language labels** Looking at the normalization performance breakdown on language labels shows that the gains of our proposed models are consistently smaller on Indonesian and Turkish compared to respectively English and German (see Appendix A for full results). This was to be expected, as for these languages the model has less external

<sup>8</sup>van der Goot (2019) used error reduction rate as main evaluation metric, for which the multilingual model would score 80.72 (Id-En) and 30.42 (Tr-De). The reported scores on monolingual datasets are 77.09 for En and 28.94 for Tr.

Model	LAI	Multiling.	Lang.-aware	Gold
MarMoT-POS	61.92	*65.50	65.47	*69.14
Bilty-POS	*65.23	*67.99	*68.26	*72.04
MaChAmp-POS	65.60	*68.25	68.13	*71.27

Table 6: Accuracies for Tr-De POS tagging, using a variety of normalization strategies.

data (Section 3.3) and while the model was originally not evaluated for Indonesian, Turkish had the lowest performance in van der Goot (2019).

**Qualitative analysis** Both Multilingual and Language-aware correct most frequent normalization mistakes well. This means for Id-En, abbreviations ( $yg \mapsto yang$  ‘which is’ in ID), slang words ( $gw \mapsto saya$ ,  $gue \mapsto saya$  ‘I’ in ID), phonetic spelling ( $kalo \mapsto kalau$  ‘if’ in ID); and for Tr-De emoticons, restoring Turkish-specific characters, restoring vowels ( $cnm \mapsto canım$  ‘my dear’ in TR), and punctuation replacements. On the Id-En data, however, there is a higher number of these frequent replacements compared to the Tr-De dataset, which explains the high scores and small variability for Id-En in Table 4 and 5.

For the Tr-De dataset, the most common mistakes include: not correcting capitalization in the beginning of a sentence, merging of words, monolingual ambiguous cases depending on context ( $mi \mapsto [mi, mi]$ , question clitics in TR), and tokenization and punctuation mistakes ( $?:D \mapsto ? :D$ ). In comparison, for the Id-En dataset, the models make rather different errors: in-vocabulary words which should be normalized are left as is ( $kaya \mapsto seperti$ ,  $usah \mapsto perlu$ ), normalizations which are lexically very distant are not found ( $lw \mapsto kamu$ ), and English contractions are often not replaced ( $isnt \mapsto is not$ ). Error analysis on the Id-En dataset revealed that correction of capitalization was annotated inconsistently. However, because in most cases the normalization was lowercased, this did not have a large effect on performance.

Interestingly, Language-aware is better in correcting words that exist in both languages. For instance, *ne* is the informal form of *eine* ‘a/one’ in German, and also means ‘what’ in Turkish. The dataset annotations expect the  $ne \mapsto eine$  normalization. While Multilingual fails to do so, Language-aware corrects them. We believe language IDs play a positive role here in defining the context, and although in general both models perform on par, if a dataset contains many such

ambiguous words, Language-aware could be preferable.

### 5.3 POS tagging

For POS tagging, we only look at Tr-De as Id-En is not annotated with POS tags. We employ a pipeline approach; we first normalize our training data in a 10-fold setting, and then apply the tagger on this normalized data. The taggers are trained on a shuffled concatenation of the Turkish-IMST (Sulubacak et al., 2016) and German-GSD (McDonald et al., 2013) datasets of UD version 2.5 (Nivre et al., 2020). Now that none of the CS data is used during training, 10-fold cross-validation is not necessary. We directly apply the taggers on the full training data. This way the exact same data split is used for evaluation as in the 10-fold setting in the previous sections. Even though we have POS tags available for the gold normalization (Section 3.1), we do not have gold tags for predicted normalization, and to keep the comparison fair we evaluate using the Tok+Anon POS tags. When a word is split or merged, we use the alignment and check whether the correct tag is present. In other words: we select one tag based on an oracle selection.<sup>9</sup>

Results in Table 6 show that, surprisingly, Bilty performs competitive to MaChAmp across most settings. Considering the differences between the normalization models, the Multilingual model and the Language-aware model perform on par, but there is still a marginal gap compared to the gold normalization.

We also analyzed the confusion matrices of the POS tagger, the full analysis can be found in Appendix C, we will shortly summarize findings here. 1) Bilty is mainly outperforming MaChAmp in gold due to better recognition of symbols (emojis), 2) Bilty is more sensitive to different normalization strategies, whereas for MaChAmp the differ-

<sup>9</sup>It should be noted that this makes splitting beneficial, and this metric can easily be tricked by splitting every token so it should be used with caution. However, our proposed normalization models have a low rate of splitting (114 versus 398 in gold) and merging is not handled at all.



Model	Normalization		POS
	Id-En	Tr-De	Tr-De
LAI	74.03	67.02	60.77
Monolingual (Id/De)	*94.62	76.33	*63.47
Multilingual	94.27	*78.28	*64.06
Language-aware	94.32	77.83	*63.92
Gold	*100.00	*100.00	*67.75

Table 7: Normalization and POS tagging accuracies on test data. The POS tagging model is the same MaChAmp model for all results, only the normalization strategy for the input changes.

ences between them are minimal, 3) Performance on nouns improves a lot after normalization, especially for German (due to corrected capitalization of nouns), 4) The second POS tag which improved most are verbs, investigation showed that this is mainly because Turkish-specific characters are replaced by their ASCII counterparts, which helps the tagger assign the correct POS.

#### 5.4 Test data

On the test data we take both the ‘no normalization’ and the best baseline (which are monolingual Indonesian for Id-En and monolingual German for Tr-De), and compare these to our best two proposed normalization models. The results in Table 7 show that, parallel to 10-fold cross-validation results (Table 4), `Multilingual` and `Language-aware` scores are similar and their difference is insignificant for both datasets. This leads to the conclusion that `Multilingual` is the most elegant model, as it is not dependent on language labels. On the Tr-De dataset the proposed models are clearly outperforming the baselines. However, on the Id-En dataset the differences are small (and not significant) between the monolingual model and both of our proposed models.

For Tr-De, we take the test set normalized by systems in the second column of Table 7 and apply MaChAmp for POS tagging. The results in the third column show that the POS tagger follows the trend in normalization scores, and performs slightly better when using the multilingual model, beating the LAI baseline (i.e. not using normalization) with 5.4% relative improvement.

## 6 Conclusion

Code-switching provides many challenges for NLP systems. In this work we attempt to overcome

some of these challenges by normalizing the data, and evaluating the downstream effect of this for POS tagging. For evaluation we use an Indonesian-English dataset (Barik et al., 2019) as well as a German-Turkish dataset (Çetinoğlu, 2016), for which we provided novel normalization layers and adapted existing LID and POS annotation.

We proposed three different models to normalize CS data. The two best-performing models are `Language-aware` and `Multilingual`. The first model exploits language labels, to identify for which language to generate features, whereas the second model combines features for both languages. The differences in performance between these two systems was not significant for any of the 10-fold experiments nor on the test data, so in most cases the multilingual model would be preferable, as it has no dependence on language labels.

We showed that normalizing the input before POS tagging results in significantly higher POS accuracies for CS data. Gold normalization experiments showed that there is still room for improvement for normalization models to help POS tagging.

An interesting property of the proposed model is that it does not have to be trained on intrasentential CS data. In fact, it can be trained on a mix of two monolingual datasets, thereby handling many more language pairs. We hope to evaluate this setup if resources (i.e., normalization test data for a CS language pair, and monolingual normalization training data for both languages) become available.

## Acknowledgements

We would like to thank Barbara Plank, Alan Ramponi, Marija Stepanovic, and Agnieszka Falenska, for feedback on early drafts. We also thank Manuel Mager and Sevde Ceylan with their help on Tr-De data alignment and Anab Maulana Barik for sharing the Id-En data. The first author is partially funded by an Amazon Research Award. The second author is funded by DFG via project CE 326/1-1 ‘‘Computational Structural Analysis of German-Turkish Code-Switching’’ (SAGT).

## References

Wafia Adouane, Jean-Philippe Bernardy, and Simon Dobnik. 2019. [Normalising non-standardised orthography in Algerian code-switched user-generated data](#). In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 131–140.

- Hong Kong, China. Association for Computational Linguistics.
- Rami Al-Rfou, Bryan Perozzi, and Steven Skiena. 2013. [Polyglot: Distributed word representations for multilingual nlp](#). In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 183–192, Sofia, Bulgaria. Association for Computational Linguistics.
- Inaki Alegria, Nora Aranberri, Víctor Fresno, Pablo Gamallo, Lluís Padró, Inaki San Vicente, Jordi Turmo, and Arkaitz Zubiaga. 2013. [Introducción a la tarea compartida Tweet-Norm 2013: Normalización léxica de tuits en español](#). In *Tweet-Norm@SEPLN*, pages 1–9.
- Fahad AlGhamdi, Giovanni Molina, Mona Diab, Thamar Solorio, Abdelati Hawwari, Victor Soto, and Julia Hirschberg. 2016. [Part of speech tagging for code switched data](#). In *Proceedings of the Second Workshop on Computational Approaches to Code Switching*, pages 98–107, Austin, Texas. Association for Computational Linguistics.
- AiTi Aw, Min Zhang, Juan Xiao, and Jian Su. 2006. [A phrase-based statistical model for SMS text normalization](#). In *Proceedings of the COLING/ACL 2006 Main Conference Poster Sessions*, pages 33–40, Sydney, Australia. Association for Computational Linguistics.
- Timothy Baldwin, Marie Catherine de Marneffe, Bo Han, Young-Bum Kim, Alan Ritter, and Wei Xu. 2015. [Shared tasks of the 2015 workshop on noisy user-generated text: Twitter lexical normalization and named entity recognition](#). In *Proceedings of the Workshop on Noisy User-generated Text*, pages 126–135, Beijing, China. Association for Computational Linguistics.
- Anab Maulana Barik, Rahmad Mahendra, and Mirna Adriani. 2019. [Normalization of Indonesian-English code-mixed twitter data](#). In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 417–424, Hong Kong, China. Association for Computational Linguistics.
- Irshad Bhat, Riyaz A. Bhat, Manish Shrivastava, and Dipti Sharma. 2017. [Joining hands: Exploiting monolingual treebanks for parsing of code-mixing data](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 324–330, Valencia, Spain. Association for Computational Linguistics.
- Irshad Bhat, Riyaz A. Bhat, Manish Shrivastava, and Dipti Sharma. 2018. [Universal Dependency parsing for Hindi-English code-switching](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 987–998, New Orleans, Louisiana. Association for Computational Linguistics.
- Leo Breiman. 2001. Random forests. *Machine learning*, 45(1):5–32.
- Özlem Çetinoğlu. 2016. [A Turkish-German code-switching corpus](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 4215–4220, Portorož, Slovenia. European Language Resources Association (ELRA).
- Özlem Çetinoğlu and Çağrı Çöltekin. 2016. [Part of speech annotation of a Turkish-German code-switching corpus](#). In *Proceedings of the 10th Linguistic Annotation Workshop held in conjunction with ACL 2016 (LAW-X 2016)*, pages 120–130, Berlin, Germany. Association for Computational Linguistics.
- Amitava Das and Björn Gambäck. 2014. [Identifying languages at the word level in code-mixed Indian social media text](#). In *Proceedings of the 11th International Conference on Natural Language Processing*, pages 378–387, Goa, India. NLP Association of India.
- Leon Derczynski, Alan Ritter, Sam Clark, and Kalina Bontcheva. 2013. [Twitter part-of-speech tagging for all: Overcoming sparse and noisy data](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing RANLP 2013*, pages 198–206, Hissar, Bulgaria. INCOMA Ltd. Shoumen, BULGARIA.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Jacob Eisenstein. 2013. [What to do about bad language on the internet](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 359–369, Atlanta, Georgia. Association for Computational Linguistics.
- Rob van der Goot. 2019. [MoNoise: A multi-lingual and easy-to-use lexical normalization tool](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 201–206, Florence, Italy. Association for Computational Linguistics.
- Rob van der Goot and Gertjan van Noord. 2017. [Parser adaptation for social media by integrating normalization](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 491–497, Vancouver, Canada. Association for Computational Linguistics.

- Rob van der Goot, Barbara Plank, and Malvina Nissim. 2017. [To normalize, or not to normalize: The impact of normalization on part-of-speech tagging](#). In *Proceedings of the 3rd Workshop on Noisy User-generated Text*, pages 31–39, Copenhagen, Denmark. Association for Computational Linguistics.
- Rob van der Goot, Alan Ramponi, Tommaso Caselli, Cafagna Michele, and Lorenzo De Mattei. 2020. [Norm it! lexical normalization for Italian and its downstream effects for dependency parsing](#). In *Proceedings of the Twelfth International Conference on Language Resources and Evaluation (LREC 2020)*. European Language Resources Association (ELRA).
- Rob van der Goot, Ahmet Üstün, Alan Ramponi, Ibrahim Sharaf, and Barbara Plank. 2021. [Massive Choice, Ample tasks \(MaChAmp\): A toolkit for multi-task learning in NLP](#). In *2005.14672v3*.
- John J Gumperz. 1982. *Discourse strategies*, volume 1. Cambridge University Press.
- Bo Han. 2014. [Improving the utility of social media with Natural Language Processing](#). Ph.D. thesis, The University of Melbourne.
- Bo Han and Timothy Baldwin. 2011. [Lexical normalisation of short text messages: Makn sens a #twitter](#). In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 368–378, Portland, Oregon, USA. Association for Computational Linguistics.
- Ning Jin. 2015. [NCSU-SAS-Ning: Candidate generation and feature engineering for supervised lexical normalization](#). In *Proceedings of the Workshop on Noisy User-generated Text*, pages 87–92, Beijing, China. Association for Computational Linguistics.
- Armand Joulin, Edouard Grave, Piotr Bojanowski, and Tomas Mikolov. 2017. [Bag of tricks for efficient text classification](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 427–431, Valencia, Spain. Association for Computational Linguistics.
- Guillaume Lample, Alexis Conneau, Marc’Aurelio Ranzato, Ludovic Denoyer, and Hervé Jégou. 2018. [Word translation without parallel data](#). In *International Conference on Learning Representations*.
- Nikola Ljubešić, Katja Zupan, Darja Fišer, and Tomaz Erjavec. 2016. [Normalising Slovene data: historical texts vs. user-generated content](#). *Bochumer Linguistische Arbeitsberichte*, page 146.
- Ismini Lourentzou, Kabir Manghnani, and Chengxiang Zhai. 2019. [Adapting sequence to sequence models for text normalization in social media](#). In *Proceedings of the International AAAI Conference on Web and Social Media*, volume 13, pages 335–345.
- Ryan McDonald, Joakim Nivre, Yvonne Quirnbach-Brundage, Yoav Goldberg, Dipanjan Das, Kuzman Ganchev, Keith Hall, Slav Petrov, Hao Zhang, Oscar Täckström, Claudia Bedini, Núria Bertomeu Castelló, and Jungmee Lee. 2013. [Universal dependency annotation for multilingual parsing](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 92–97, Sofia, Bulgaria. Association for Computational Linguistics.
- Giovanni Molina, Fahad AlGhamdi, Mahmoud Ghoneim, Abdelati Hawwari, Nicolas Rey-Villamizar, Mona Diab, and Thamar Solorio. 2016. [Overview for the second shared task on language identification in code-switched data](#). In *Proceedings of the Second Workshop on Computational Approaches to Code Switching*, pages 40–49, Austin, Texas. Association for Computational Linguistics.
- Thomas Mueller, Helmut Schmid, and Hinrich Schütze. 2013. [Efficient higher-order CRFs for morphological tagging](#). In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 322–332, Seattle, Washington, USA. Association for Computational Linguistics.
- Benjamin Muller, Benoit Sagot, and Djamel Seddah. 2019. [Enhancing BERT for lexical normalization](#). In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 297–306, Hong Kong, China. Association for Computational Linguistics.
- Carol Myers-Scotton. 1995. *Social motivations for codeswitching: Evidence from Africa*. Oxford University Press.
- Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Yoav Goldberg, Jan Hajič, Christopher D. Manning, Ryan McDonald, Slav Petrov, Sampo Pyysalo, Natalia Silveira, Reut Tsarfaty, and Daniel Zeman. 2016. [Universal dependencies v1: A multilingual treebank collection](#). In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*.
- Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajič, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. [Universal Dependencies v2: An evergrowing multilingual treebank collection](#). In *Proceedings of The 12th Language Resources and Evaluation Conference*, pages 4034–4043, Marseille, France. European Language Resources Association.
- Franz Josef Och and Hermann Ney. 2003. [A systematic comparison of various statistical alignment models](#). *Computational Linguistics*, 29(1):19–51.
- Brendan O’Connor, Michel Krieger, and David Ahn. 2010. [Tweetmotif: Exploratory search and topic summarization for twitter](#). In *Fourth International AAAI Conference on Weblogs and Social Media*.

- Deana Pennell and Yang Liu. 2011. [A character-level machine translation approach for normalization of SMS abbreviations](#). In *Proceedings of 5th International Joint Conference on Natural Language Processing*, pages 974–982, Chiang Mai, Thailand. Asian Federation of Natural Language Processing.
- Barbara Plank, Anders Søgaard, and Yoav Goldberg. 2016. [Multilingual part-of-speech tagging with bidirectional long short-term memory models and auxiliary loss](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 412–418, Berlin, Germany. Association for Computational Linguistics.
- Royal Sequiera, Monojit Choudhury, Parth Gupta, Paolo Rosso, Shubham Kumar, Somnath Banerjee, Sudip Kumar Naskar, Sivaji Bandyopadhyay, Gokul Chittaranjan, Amitava Das, et al. 2015. [Overview of FIRE-2015 shared task on mixed script information retrieval](#). In *FIRE Workshops*, volume 1587, pages 19–25.
- Zareen Sharf and Saif Ur Rahman. 2017. [Lexical normalization of roman Urdu text](#). *International Journal of Computer Science and Network Security*, 17(12):213–221.
- Annav Sharma, Sakshi Gupta, Raveesh Motlani, Piyush Bansal, Manish Shrivastava, Radhika Mamidi, and Dipti M. Sharma. 2016. [Shallow parsing pipeline - Hindi-English code-mixed social media text](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 1340–1345, San Diego, California. Association for Computational Linguistics.
- Victor Soto and Julia Hirschberg. 2018. [Joint part-of-speech and language ID tagging for code-switched data](#). In *Proceedings of the Third Workshop on Computational Approaches to Linguistic Code-Switching*, pages 1–10, Melbourne, Australia. Association for Computational Linguistics.
- Umut Sulubacak, Memduh Gokirmak, Francis Tyers, Çağrı Çöltekin, Joakim Nivre, and Gülşen Eryiğit. 2016. [Universal Dependencies for Turkish](#). In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*, pages 3444–3454, Osaka, Japan.
- Almeida Jacqueline Toribio and Barbara E Bullock. 2012. *The Cambridge handbook of linguistic code-switching*. Cambridge University Press.
- Congle Zhang, Tyler Baldwin, Howard Ho, Benny Kimelfeld, and Yunyao Li. 2013. [Adaptive parser-centric text normalization](#). In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1159–1168, Sofia, Bulgaria. Association for Computational Linguistics.
- Yuan Zhang, Jason Riesa, Daniel Gillick, Anton Bakalov, Jason Baldrige, and David Weiss. 2018. [A fast, compact, accurate model for language identification of codemixed text](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 328–337, Brussels, Belgium. Association for Computational Linguistics.

Model	ID	EN	TR	DE
LAI	66.92	71.33	70.21	66.92
MFR	87.58	88.10	76.26	69.55
Monoling1 (tr/id)	92.71	95.82	78.24	70.17
Monoling2 (de/en)	91.78	95.75	76.81	77.50
Frag	84.71	82.03	70.10	63.92
Multiling.	92.91	95.80	78.27	78.73
Lang-aware	92.77	95.78	78.10	79.32

Table 8: Normalization accuracies per language (with gold language labels)

Model	Id-En		Tr-De	
	recall	precision	recall	precision
LAI	0.00	0.00	0.00	0.00
MFR	57.50	98.20	21.42	84.39
Monoling1 (tr/id)	82.15	97.96	25.55	88.50
Monoling2 (de/en)	80.35	98.03	29.49	87.32
Frag	82.05	97.88	31.15	90.21
Multiling.	82.51	97.87	32.98	90.85
Lang-aware	82.21	97.98	32.95	90.34

Table 9: precision and recall for both datasets, we follow the definitions of (van der Goot, 2019)

## Appendix

### A Breakdown of performance per language

Table 8 show the accuracy of all the proposed models per language. The LAI scores show that most of the normalization replacements are necessary for ID and DE. Interestingly, performance of the last two models is highest on respectively EN and DE, which is probably due to the original model being developed mostly with a focus European languages.

### B Precision and recall

Table 9 show the precision and recall of all models on both datasets. LAI has 0.0 on all metrics, because it never finds a correct normalization.

### C Confusions of POS taggers

We conducted an analysis of POS tagging confusions for the setting described in Section 5.3. In Table 10 and Table 11 the error frequencies of respectively MaChAmp and Bilty are shown. The tables report the frequency of the top-10 most frequent errors of the baseline (LAI), and the difference in counts observed using a variety of normalization strategies. In Figure 3 and Figure 3 the full confusion matrices for respectively MaChAmp

	LAI	Multiling.	Lang-aware	Gold
SYM-PUNCT	529	+2	+2	+34
NOUN-PROPN	310	+10	+10	+36
PROPN-NOUN	307	-24	-11	-39
NOUN-ADJ	244	-40	-38	-40
PROPN-PUNCT	220	-5	-8	-16
VERB-NOUN	174	-21	-27	-78
PROPN-ADJ	122	-13	-12	-20
ADV-ADJ	108	-23	-23	-23
ADJ-NOUN	104	-22	-23	-30
ADJ-PROPN	103	-2	+0	+5

Table 10: 10 most common POS tagging errors for LAI baseline, counted for all normalization strategies for the MaChAmp tagger. Counts are all relative compared to the baseline (LAI). The tag on the left is gold, right is predicted.

and Bilty are shown. For both of these analyses, we do not report the other baselines, the fragment based model and the MarMot tagger, because performance of these was inferior and this would make comparisons more complex.

	LAI	Multiling.	Lang-aware	Gold
PROPN-VERB	507	+67	-417	+65
PROPN-NOUN	310	-95	+138	-205
VERB-NOUN	239	+32	-44	-134
NOUN-ADJ	200	-59	-44	-65
SYM-PUNCT	194	+19	+83	+152
NOUN-PROPN	162	-19	+17	+14
INTJ-NOUN	152	+12	-10	-30
SYM-ADJ	127	-90	-126	-96
ADJ-NOUN	124	+2	-16	-20
NOUN-VERB	122	-12	-27	-39

Table 11: 10 most common POS tagging errors for LAI baseline, counted for all normalization strategies for the Bilty tagger. Counts are all relative compared to the baseline (LAI). The tag on the left is gold, right is predicted.

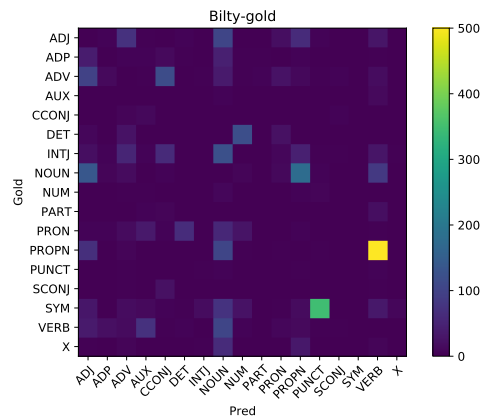
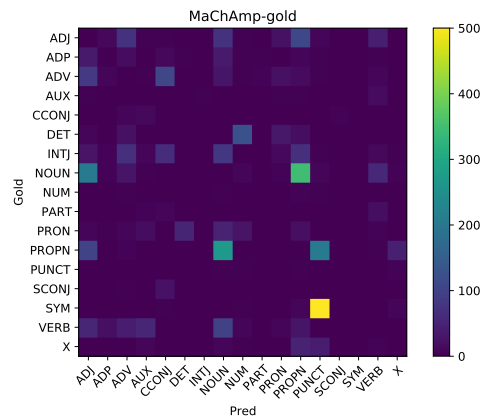
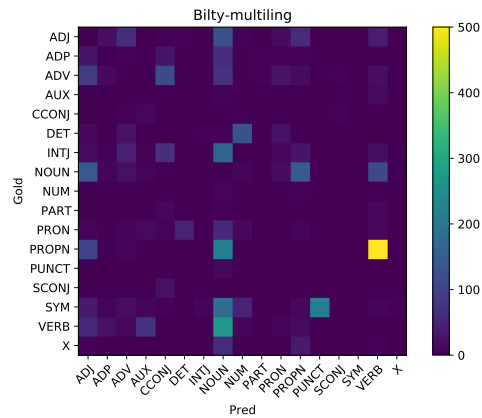
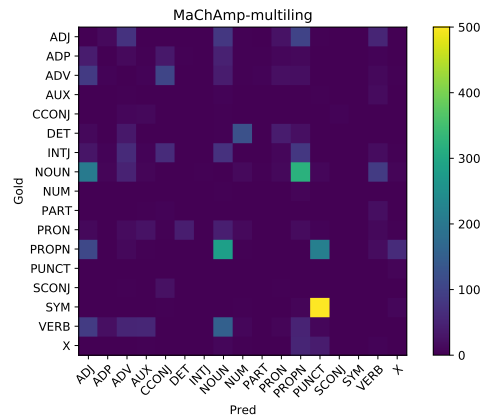
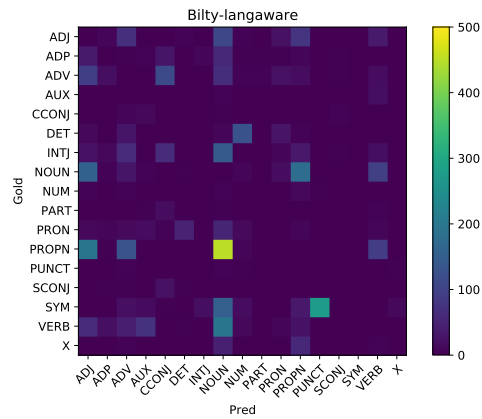
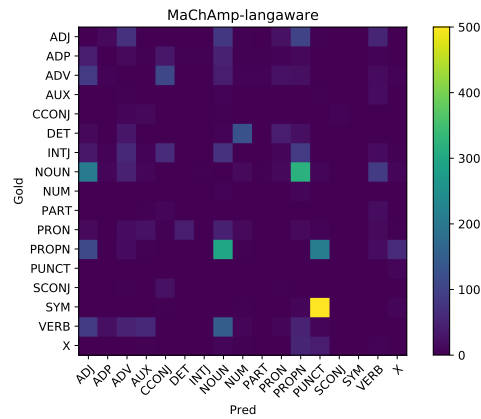
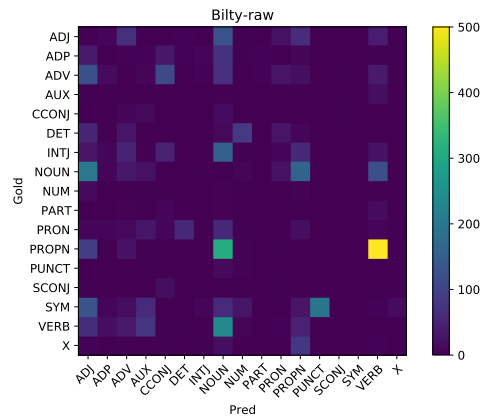
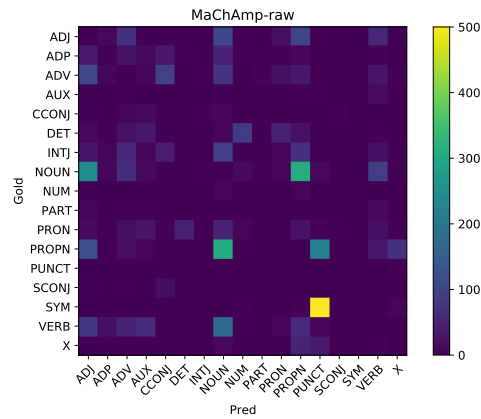


Figure 3: Confusion matrices for MaChAmp using a variety of normalization strategies

Figure 3: Confusion matrices for Bilty using a variety of normalization strategies