

# Translate & Fill: Improving Zero-Shot Multilingual Semantic Parsing with Synthetic Data

Massimo Nicosia, Zhongdi Qu, Yasemin Altun

Google Research

{massimon, dqu, altun}@google.com

## Abstract

While multilingual pretrained language models (LMs) fine-tuned on a single language have shown substantial cross-lingual task transfer capabilities, there is still a wide performance gap in semantic parsing tasks when target language supervision is available. In this paper, we propose a novel Translate-and-Fill (TaF) method to produce silver training data for a multilingual semantic parser. This method simplifies the popular Translate-Align-Project (TAP) pipeline and consists of a sequence-to-sequence filler model that constructs a full parse conditioned on an utterance and a view of the same parse. Our filler is trained on English data only but can accurately complete instances in other languages (i.e., translations of the English training utterances), in a zero-shot fashion. Experimental results on three multilingual semantic parsing datasets show that data augmentation with TaF reaches accuracies competitive with similar systems which rely on traditional alignment techniques.

## 1 Introduction

Semantic parsing is a core task in virtual assistants, popular applications that require accurate natural language understanding (NLU). User utterances are parsed into a structured representation made of intents and slots that is interpreted to initiate an action on the user device. For example, the sentence “set an 8 am alarm” could lead to the following interpretation – *Create\_alarm(time=“8 am”)* – and result in an alarm being created.

As in many NLP tasks, numerous English parsing datasets are available and well studied (Price, 1990; Banarescu et al., 2013; Williams et al., 2016; Fan et al., 2017; Gupta et al., 2018; Goo et al., 2018; Qin et al., 2019; Rongali et al., 2020). Supporting new domains and schemas requires a sizeable data collection effort and while English is receiving the most attention, it is also important to extend NLU to other languages in order to provide users con-

sistent experiences across languages. Multilingual pretrained language models (LMs) are an excellent starting point for enabling cross-lingual transfer in a parser but they are no substitute for using high quality, albeit costly to annotate, training data in the target languages. Without such data, we can translate the available annotated examples to other languages and slot annotations can be transferred (Yarowsky et al., 2001; Shah et al., 2010). Traditionally, annotation transfer requires (i) token alignment models (Brown et al., 1993), which may have been trained on text tokenized differently from the annotated training data, and (ii) label projection logic that can be complex, especially if it includes heuristics for fixing systematic alignment errors, or if nested structures need to be mapped.

In this work, we propose an alternative approach to the classical Translate-Align-Project (TAP) pipeline: we leverage multilingual pretrained representations and a sequence-to-sequence (seq2seq) model to directly generate the parse of translated examples in a zero-shot fashion. Our model is trained on English data only and it is able to reconstruct the full parse while having access to the English utterance and to a signature (or view) of the full parse. At inference, we substitute the English utterance with its translation and our model, pulling content from the latter, is able to construct a high quality silver parse. The main contributions of this paper can be summarized as follows:

- We propose a novel approach, Translate-and-Fill (TaF), for generating synthetic data to train multilingual semantic parsers that is robust to tokenization, is inherently generative and makes use of the intent and slot schema to potentially learn label-specific alignment rules. TaF replaces the alignment and projection modules of the TAP approach with a learned component that generates full parses of examples translated from English, removing the need of aligners.

- We analyze the zero-shot capabilities of TaF in terms of quality of the silver parses.
- We evaluate the impact of the synthetic data generated with our approach on three multilingual semantic parsing datasets, showing that data augmentation with TaF on multilingual pretrained seq2seq models sets new state-of-the-art (SOTA) results in multiple scenarios and in some cases, closes the gap with respect to full in-language supervision.

## 2 Related Work

Our work is closely related to two research areas: (i) multilingual representations and models, and (ii) annotation projection methods.

Cross-lingual transfer has been studied in several structured prediction tasks such as part-of-speech tagging (Yarowsky et al., 2001; Täckström et al., 2013; Plank and Agić, 2018; Kann et al., 2020), named entity recognition (Zirikly and Hagiwara, 2015; Tsai et al., 2016; Xie et al., 2018) and dependency parsing (Guo et al., 2015; Ahmad et al., 2019; Zhang et al., 2019a).

One way to achieve cross-lingual transfer is by adopting multilingual representations and models pretrained on a large amount of text in different languages. This way, similar languages with overlapping vocabularies at word or subword level can benefit from information sharing. These models can encode the input using words (Mikolov et al., 2013; Pennington et al., 2014), characters or subwords (Sennrich et al., 2016; Kudo and Richardson, 2018; Wu et al., 2016; Clark et al., 2021). With the latter, interesting zero-shot performance (i.e., training on a language and evaluating on a different target language) can be achieved, especially between similar languages (Lauscher et al., 2020).

Multilingual representations can be obtained from encoders pretrained on multilingual corpora with tasks such as masked language modeling (MLM), or trained on supervised tasks such as neural machine translation (NMT) (Eriguchi et al., 2018; Yu et al., 2018; Singla et al., 2018; Siddhant et al., 2020). After the success of fill-in-the-blank-style denoising objectives and BERT/mBERT (Devlin et al., 2019), other multilingual encoders achieved a similar level of popularity. These models include XLM (Lample and Conneau, 2019), XLM-R (Conneau et al., 2020) and a recent multilingual version of T5 (Raffel et al., 2020), named mT5 (Xue et al., 2021). T5 based models differ

from the others by their seq2seq architecture where both the encoder and the decoder are pretrained with the MLM task. In this work, we leverage the multilinguality and the generative capabilities of mT5 to produce interpretations and create synthetic internationalization (i18n) data for semantic parsing.

A second way to improve cross-lingual transfer is data augmentation. Typically, annotated data is available in at least one language, and more often than not, this is a high-resource language such as English. NMT is a strong data augmentation baseline, as shown in recent cross-lingual evaluation benchmarks (Hu et al., 2020; Ladhak et al., 2020). NMT can be used to translate training examples from a source to a target language (translate-train), creating training data in the target language. Otherwise, it can be used to translate the test data to the language of the trained model (translate-test).

While translating works quite well for classification tasks where the label is at instance level, for sequence tagging or parsing tasks the reality is more challenging since the labels are at token level and they have to be transferred from the tokens of the original text to the tokens of its translation.

Prior work relies on word aligners to establish a match between the tokens of source and translated text, and to transfer the labels (Ni et al., 2017; Jain et al., 2019; Daza and Frank, 2020; Fei et al., 2020). Alignment methods include unsupervised word alignment (Brown et al., 1993; Vogel et al., 1996; Och and Ney, 2000, 2003), the use of attention weights from NMT models (Schuster et al., 2019; Chen et al., 2020; Zenkel et al., 2020) or computing the similarity between word embeddings (Jalili Sabet et al., 2020; Dou and Neubig, 2021).

In this work, we propose an alternative and novel label projection method that leverages the signatures of available parses for internationalization, in the spirit of sketch or template decoding (Dong and Lapata, 2018; Zhang et al., 2019b; Wiseman et al., 2018)). Our method avoids alignment models altogether and leverages multilingual representations and instance labels to generate high quality silver data that can be finally used to train accurate multilingual semantic parsers. In addition and differently from NMT attention-based aligners, our method does not access the internals of neural translation models and therefore has a wider applicability.

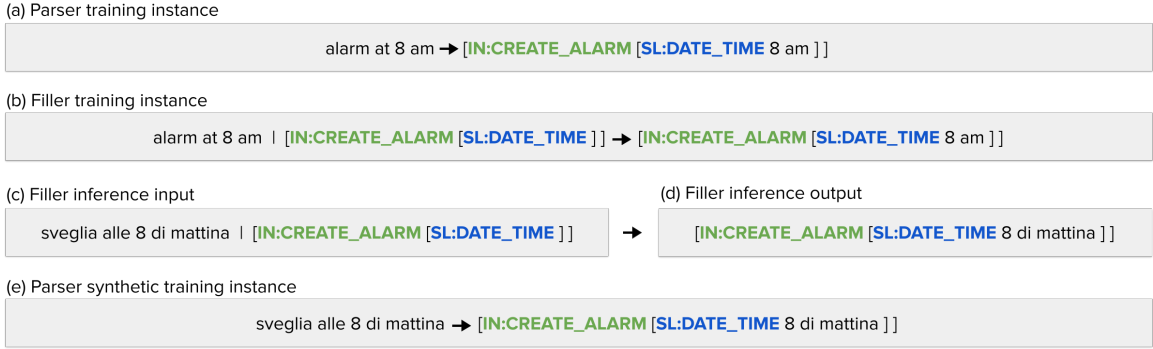


Figure 1: Example instances for training the semantic **parser** (a) and the **filler** (b). The filler is trained to produce a full parse from the concatenation of an English utterance and the corresponding parse signature (b). At inference, we replace the English utterance with its (Italian in this case) translation (c), and obtain a silver parse where the slots contain text from the translation (d). The latter is used to assemble a synthetic training instance (e) for a multilingual semantic parser.

### 3 Translate-and-Fill (TaF)

We address the problem of the i18n of semantic parsers when (i) English training data is available and (ii) high quality and cost-effective training data in other languages is needed. We translate English data to a target language using NMT. This leaves us with the problem of mapping original slot annotations to translations. Our solution is a novel method that we call Translate-and-Fill (TaF), which replaces the align and project modules of the popular Translate-Align-Project (TAP) pipeline while leveraging multilingual pretraining. In our approach, we use two seq2seq models trained differently: one is the usual **semantic parser** and the other is what we call the **filler**.

Figure 1 shows the example instances used to train the semantic parser and our filler, and then to run inference with the latter. The example parse has a `CREATE_ALARM` intent (`IN:`) and a single `DATE_TIME` slot (`SL:`). We transform a training instance for the semantic parser that maps an utterance to its parse (a) into a training instance for the filler. A filler training instance (b) maps the English utterance concatenated with its parse signature to a full parse (target from a). To obtain the parse signature, we simply remove all the slot values from the parse. The filler must then reproduce the input signature while filling the signature slots with words from the input utterance.

We leverage pretrained multilingual seq2seq models (in particular mT5) to train the filler model with only English filler instances. A trained filler can be used to obtain labeled semantic parsing data in other languages, thanks to the cross-lingual transfer capabilities of the pretrained seq2seq model, as

well as the slot-filling capabilities gathered from the English training filler instances. We construct an inference example for the filler from the same examples used to train it (b) by simply replacing the English utterance in the input part with its corresponding translation (c). The filler will now reproduce the input signature but fill the slots using words from the translation (d).

Finally, we create a synthetic i18n instance for training a parser for the target language. The synthetic instance maps the translated utterance to the parse produced by the filler at inference (e).

Similar to TAP, our basic assumption is that the parse structure of a translated sentence does not change. The proposed approach (i) can be applied to any language supported by NMT and by the pretrained seq2seq model; (ii) can handle nested interpretations naturally thanks to the seq2seq formulation; and (iii) since it has access to the interpretation, it can learn label specific projection strategies as opposed to the handcrafted TAP projection rules.

## 4 Experimental Setup

### 4.1 Datasets

We experiment with three multilingual task-oriented semantic parsing datasets.

**MTOP** (Li et al., 2021) is an almost parallel dataset covering 6 languages and 11 domains. Each utterance has associated intent and slots, but also comes with a decoupled compositional representation similar to the parses in Figure 1. Compositional instances will have nested intents. The seq2seq nature of our model lets

us handle such cases without any specialized component. For our experiments, we use the provided train/validation/test splits and focus on predicting the decoupled representations.

**Multilingual ATIS** (Upadhyay et al., 2018) is a dataset for the travel-planning domain that extends the popular ATIS dataset (Price, 1990) to two other languages: Hindi and Turkish. Differently from MTOP, there are no nested intents and therefore just flat span annotations.

**MultiATIS++** (Xu et al., 2020) adds six new languages to Multilingual ATIS bringing the number of non-English languages to 8. For both Multilingual ATIS and MultiATIS++, we create an MTOP-style interpretation by converting the BIO-tagged sentences into an intent/slot structure (as in Figure 1). For both datasets we use the train/validation/test split ratios reported in Xu et al. (2020).

## 4.2 Models

Our parser and filler are trained using mT5 (Xue et al., 2021), a multilingual version of the text-to-text T5 model (Raffel et al., 2020), pretrained on the mC4 corpus and 101 languages. We experiment with two mT5 versions, `large` and `xxl` in different settings. In the *gold data* setting, we train multilingual parsers with all the available training data. In the *zero-shot* setting, we train our models on English data only. In the *+TaF* setting, we train our models on English gold data and on the synthetic data produced by our filler for all the other languages. We do not do any hyper-parameter tuning and use a batch size of 512 and a constant 0.001 learning rate. We train all our models for 3k steps saving checkpoints every 200 steps. The parser produces structured interpretations and we run Unicode normalization on the tokens. The filler is an mT5-xxl model trained for 400 steps since its output does not significantly change after that. We then run inference on the « translation | signature » inputs and generate synthetic training data for the parser. Apart from discarding a negligible number of outputs that cannot be parsed into a tree, we do not apply any additional quality filter. According to our experience, this is an advantage w.r.t. alignment-based methods that require complex filtering to suppress systematic alignment errors and improve synthetic data quality.

Language	en	de	es	fr	hi	th
Match %	93.50	93.75	96.39	94.61	98.35	42.08

Table 1: % of MTOP training instances where our tokenization matches the original MTOP tokenization.

## 4.3 Translation and Postprocessing

We translate the English utterances to different target languages and tokenize them with in-house translation and tokenization systems. The datasets used in our experiments come with tokenized gold data but no tokenized translations. In the MTOP paper an in-house tokenizer is used, while the other dataset papers do not contain details about tokenization. This is a common issue, as also reported in Kaliamoorthi et al. (2021). This implies a tokenization mismatch between our synthetic data and the synthetic data used in the original dataset papers. To quantify this, we compare our tokenization of MTOP utterances with gold tokenization. Table 1 shows that we can reasonably match the original tokenization for all languages except for Thai. In the synthetic data setting, this could potentially disadvantage our results due to the noise introduced by the dissimilar tokenization. In one experiment, we do not tokenize the translations to test the quality of the final synthetic data.

In Multilingual ATIS and MultiATIS++, Spanish and Turkish eval data is all lowercase and we lowercase our translations too. In addition, Turkish data does not contain special characters, so we replace the latter in the translations according to the following mapping:  $\check{g}\check{G}\check{i}\check{I}\check{o}\check{O}\check{u}\check{U}\check{s}\check{S}\check{c}\check{C} \Rightarrow gGiIoOuUsScC$ .

## 4.4 Evaluation

For MTOP, we report Exact Match (EM) accuracy as in Li et al. (2021). For Multilingual ATIS and MultiATIS++, we report EM accuracy, Intent Accuracy and Slot F1 (micro) computed with the `seqeval` toolkit (Nakayama, 2018). Since we predict structured interpretations, we reconvert our outputs to a sequence of BIO-tagged tokens before computing Slot F1. We first map slots which can be unambiguously identified in the input utterance by full or partial string matching. The remaining slots are aligned using the Needleman-Wunsch alignment algorithm (Needleman and Wunsch, 1970), a strategy shown to be robust to small generation errors (Paolini et al., 2021). In the *Avg* columns of the tables, we report the evaluation metrics averaged over the non-English languages.



For the *gold* data setting, we do model selection by computing the best average (across non-en languages) EM on the dev set. For the *zero-shot* and *+TaF* settings we compute metrics using the last checkpoint, assuming the unavailability of a development set. To keep the amount of compute required for running the experiments reasonable, all our numbers are averaged over three runs, and we report standard deviation.

#### 4.5 Translate-and-Align (TAP) Baseline

We also experiment with synthetic data produced via TAP, aligning tokens with an implementation of the IBM Model (Brown et al., 1993) and HMM (Vogel et al., 1996). To achieve the best alignment quality, we tokenize both the English input and the translations with our in-house tokenizer (also used to train the aligner), and discard examples for which our tokenization of the English utterance differs from the original. We apply heuristic filters to the synthetic data, discarding examples where a span is split into non-consecutive tokens in the target, and examples where the target has a set of slots different from the source.

Two significant sources of error in TAP data are prepositions and determiners. When those are introduced in a translation, they are often aligned to the adjacent nouns in the original English utterance and are therefore included in the nouns’ slots. Take an example from the MTOP dataset, “Play some Elvis for me”. Its French translation is “Jouez à Elvis pour moi”, where “à” is a preposition with no direct correspondence in the English utterance. As a result, our aligner maps it to “Elvis”, and the value for the slot `MUSIC_ARTIST_NAME` becomes “à Elvis”, instead of “Elvis”. To mitigate this problem, we run the translated utterances through an in-house parts-of-speech tagger and exclude prepositions and determiners from the slots when they appear at slot boundaries (except for the slot `DATE_TIME`, for which prepositions and determiners are generally kept as a part of the slot values in the MTOP dataset). The POS tagger also performs tokenization and we discard examples for which the POS tokenization differs from the aligner tokenization so that the data left have both high-quality alignments and POS tags.

We also observed that the aligner performs poorly around punctuations that are introduced in the target utterances to function as word connectors. Take an example from MTOP, “will there be

fog in the morning”. Its French translation is “y aura-t-il du brouillard le matin”, where “il” translates to “it” and serves the same function as “it” in English sentences about the weather such as “it is raining”. Our aligner maps both the second “-” and “il” to “fog”, and as a result the value for the slot `WEATHER_ATTRIBUTE` becomes “- il du brouillard”, instead of just “brouillard”. To obtain high-quality synthetic data without these issues, we have experimented with training using only the part of data where our tokenizer does simple white-space tokenization on the target utterances. These data points, which do not contain punctuations as individual tokens, are easier to align and ultimately leads to better synthetic data.

The fraction of examples discarded during the TaF filtering stage ranges between 0.01%-0.4% for both MTOP and MultiATIS++. For TAP, significantly more filtering was required: for MTOP, 33.1% of examples were filtered because the aligner tokenizes the source queries differently from the dataset tokenizer, 30.4% because target queries cannot be simply tokenized by white-space, 0.8% due to span splitting, and 3.1% because projected labels have a different set of slots from the original signature; for MultiATIS++, 10.0% were filtered because the aligner tokenization differs from the provided source tokens, 32.9% because our tokenizer and the aligner tokenize the translations differently, 0.8% because of span splitting, and 5.8% because projected labels have a different set of slots from the original signature.

## 5 Results and Discussion

**MTOP.** Table 2 contains the results on MTOP. *XLM-R* from Li et al. (2021) is a seq2seq model that uses XLM-R as encoder and it is extended with a pointer network. This and the *mt5-xxl* model have a comparable average EM accuracy when trained multilingually with all the available gold data, although *mt5-xxl* has more parameters. In the zero-shot setting, *mt5-large* lags behind *XLM-R* by 7.5 EM points, while *mt5-xxl* already improves over *XLM-R* by 16.3 EM points. When *+TaF* synthetic data is added, *mt5-large+TaF* reaches *XLM-R+TAP*, and *mt5-xxl+TaF* surpasses it by 2.5 points, indicating that TaF is effective for i18n over strong and weak base models. While we could not run Li et al. (2021) model on our data, we can see that *mt5-large+TaF* is able to close all the gap with *XLM-R+TAP*, despite starting from a

MTOP	en	es	fr	de	hi	th	Avg(5 langs)
<i>Multilingual models (trained on all data from all languages)</i>							
XLM-R	83.6	79.8	78	74	74	73.4	75.8
mT5-large	83.8 $\pm 0.2$	76.9 $\pm 0.1$	75.2 $\pm 0.2$	72.8 $\pm 0.3$	73.2 $\pm 0.4$	73.3 $\pm 0.2$	74.3 $\pm 0.2$
mT5-xxl	86.0 $\pm 0.4$	79.3 $\pm 0.6$	77.5 $\pm 0.5$	75.5 $\pm 0.9$	75.7 $\pm 0.3$	75.1 $\pm 0.3$	76.6 $\pm 0.5$
<i>Zero-shot models (trained on English only)</i>							
XLM-R	N/A	50.3	43.9	42.3	30.9	26.7	38.8
mT5-large	83.2 $\pm 0.2$	40.0 $\pm 0.7$	41.1 $\pm 1.8$	36.2 $\pm 1.5$	16.5 $\pm 3.3$	23.0 $\pm 2.1$	31.3 $\pm 1.8$
mT5-xxl	86.7 $\pm 0.1$	62.4 $\pm 2.1$	63.7 $\pm 1.3$	57.1 $\pm 1.2$	43.3 $\pm 0.2$	49.2 $\pm 0.8$	55.1 $\pm 1.0$
<i>Augmented data models</i>							
XLM-R + TAP	N/A	<b>71.9</b>	70.3	62.4	<b>63</b>	60	65.5
mT5-large + TaF	83.5 $\pm 0.6$	69.6 $\pm 0.7$	71.1 $\pm 0.6$	70.5 $\pm 0.4$	58.1 $\pm 1.1$	57.5 $\pm 0.5$	65.4 $\pm 0.6$
mT5-xxl + TaF	85.9 $\pm 0.1$	71.5 $\pm 0.2$	74.0 $\pm 1.1$	<b>72.4</b> $\pm 0.2$	61.9 $\pm 0.4$	60.2 $\pm 0.3$	68.0 $\pm 0.1$
mT5-xxl + TaF, untokenized	85.9 $\pm 0.2$	71.5 $\pm 0.1$	<b>74.6</b> $\pm 0.2$	71.9 $\pm 0.1$	61.5 $\pm 0.4$	<b>62.2</b> $\pm 0.4$	<b>68.3</b> $\pm 0.1$
mT5-xxl + TAP	86.2 $\pm 0.1$	69.3 $\pm 0.4$	71.5 $\pm 0.3$	62.1 $\pm 0.3$	57.8 $\pm 0.3$	58.2 $\pm 0.9$	63.8 $\pm 0.4$

Table 2: Exact Match (EM) accuracies on the MTOP dataset. XLM-R results are from Li et al. (2021). In bold, we mark best performances in the data augmentation scenario.

MTOP	es	fr	de	hi	th	Avg
mT5-xxl (zero-shot)	62.4	63.7	57.1	43.3	49.2	55.1
mT5-xxl + TAP	54.2	55.8	57.4	55.3	39.8	52.5
+ POS-based postprocessing	68.5	67.2	62.2	59.6	46.0	60.7
+ white-space tokenization	69.3	71.5	62.1	57.8	58.2	63.8

Table 3: Exact Match (EM) on the MTOP dataset with different TAP configurations.

MultiAtis++	en	es	de	zh	ja	pt	fr	hi	tr	Avg(8 langs)
<i>Multilingual Intent Accuracy</i>										
mBERT	97.20	96.77	96.86	95.54	96.44	96.48	97.24	92.70	92.2	95.44
mT5-xxl	97.84 $\pm 0.13$	97.57 $\pm 0.17$	97.16 $\pm 0.17$	97.13 $\pm 0.26$	97.50 $\pm 0.17$	97.72 $\pm 0.26$	97.98 $\pm 0.22$	95.97 $\pm 0.51$	94.87 $\pm 0.40$	96.99 $\pm 0.27$
<i>Multilingual Slot F1</i>										
mBERT	95.90	87.95	95.00	93.67	92.04	91.96	90.39	86.73	86.04	91.02
mT5-xxl	96.29 $\pm 0.04$	89.31 $\pm 0.39$	95.48 $\pm 0.16$	94.59 $\pm 0.21$	93.54 $\pm 0.03$	93.00 $\pm 0.27$	90.12 $\pm 0.11$	89.83 $\pm 0.25$	87.88 $\pm 0.20$	91.72 $\pm 0.20$
<i>Zero-Shot and Augmented Intent Accuracy</i>										
mBERT	N/A	96.35	95.27	86.27	79.42	94.96	95.92	80.96	69.59	87.34
mBERT + fastalign	N/A	97.02	96.77	96.10	88.82	96.55	96.89	93.12	93.77	94.88
mBERT + softalign	N/A	97.20	96.66	95.99	88.33	96.78	97.49	92.81	93.71	94.87
mT5-xxl	97.87 $\pm 0.11$	96.90 $\pm 0.34$	93.06 $\pm 1.62$	92.53 $\pm 0.55$	89.18 $\pm 0.64$	96.75 $\pm 0.22$	96.83 $\pm 0.42$	92.46 $\pm 0.32$	86.67 $\pm 1.07$	93.05 $\pm 0.47$
mT5-xxl + TaF	97.65 $\pm 0.11$	97.65 $\pm 0.22$	96.79 $\pm 0.13$	96.75 $\pm 0.11$	<b>95.41</b> $\pm 0.19$	<b>97.61</b> $\pm 0.17$	<b>97.61</b> $\pm 0.17$	<b>96.53</b> $\pm 0.11$	<b>95.06</b> $\pm 0.21$	<b>96.68</b> $\pm 0.12$
mT5-xxl + TAP	97.76 $\pm 0.11$	<b>97.69</b> $\pm 0.06$	<b>97.76</b> $\pm 0.11$	<b>97.72</b> $\pm 0.26$	94.66 $\pm 0.53$	96.79 $\pm 0.06$	97.13 $\pm 0.13$	95.71 $\pm 0.17$	93.85 $\pm 0.37$	96.41 $\pm 0.01$
<i>Zero-Shot and Augmented Slot F1</i>										
mBERT	N/A	74.98	82.61	62.27	35.75	74.05	75.71	31.21	23.75	57.54
mBERT + fastalign	N/A	79.18	87.21	81.82	79.53	78.26	70.18	69.42	23.61	71.15
mBERT + softalign	N/A	76.42	<b>89.00</b>	83.25	79.10	76.30	79.64	78.56	61.70	78.00
mBERT + TMP	N/A	83.98	87.54	85.05	82.60	81.73	79.80	77.24	44.80	77.84
mT5-xxl	96.19 $\pm 0.19$	84.60 $\pm 1.20$	77.03 $\pm 0.59$	81.00 $\pm 1.31$	59.29 $\pm 3.76$	81.62 $\pm 1.06$	81.72 $\pm 1.20$	66.28 $\pm 5.12$	50.50 $\pm 3.37$	72.76 $\pm 1.25$
mT5-xxl + TaF	95.35 $\pm 0.17$	<b>88.26</b> $\pm 0.05$	86.78 $\pm 0.10$	<b>87.49</b> $\pm 0.41$	<b>88.66</b> $\pm 0.43$	<b>87.30</b> $\pm 0.37$	<b>86.19</b> $\pm 0.25$	<b>88.06</b> $\pm 0.08$	<b>84.47</b> $\pm 0.27$	<b>87.15</b> $\pm 0.14$
mT5-xxl + TAP	95.77 $\pm 0.18$	85.40 $\pm 0.13$	84.25 $\pm 0.19$	81.65 $\pm 0.21$	82.05 $\pm 0.24$	82.85 $\pm 0.70$	84.48 $\pm 0.57$	86.11 $\pm 0.21$	82.05 $\pm 1.05$	83.61 $\pm 0.27$

Table 4: Intent Accuracy and Slot F1 of our mT5 models on MultiAtis++. Multilingual BERT (mBERT) results are from Xu et al. (2020). In bold, the best models in the data augmentation scenario.

MultiAtis++	en	es	de	zh	ja	pt	fr	hi	tr	Avg(8 langs)
<i>Intent Accuracy</i>										
mBERT, Zero-shot	96.53	82.31	86.9	85.89	81.08	84.43	92.72	75.59	71.61	82.57
mBERT, TAP	97.12	95.00	96.15	93.92	91.68	96.38	96.15	94.55	79.67	92.94
mBERT, TaF	97.16	93.32	96.60	95.29	93.58	95.19	95.67	95.11	93.48	94.78
mBERT, Gold	96.75	94.4	96.53	93.17	94.29	95.97	97.31	92.95	90.63	94.41
mBERT, Gold (es lowercased)	96.75	93.73	95.41	90.48	91.38	95.97	96.53	92.61	90.63	93.34
<i>Slot F1</i>										
mBERT, Zero-shot	95.65	43.83	31.25	67.2	50.8	48.71	45.32	40.36	29.74	44.65
mBERT, TAP	95.79	77.48	76.05	78.69	70.25	79.38	77.89	79.36	60.24	74.92
mBERT, TaF	95.78	81.18	81.80	84.11	86.97	82.14	79.21	86.13	84.99	83.31
mBERT, Gold	95.91	72.41	90.61	90.76	89.91	87.03	87.29	86.49	85.65	86.27
mBERT, Gold (es lowercased)	96.11	80.63	91.22	89.96	88.53	88.25	87.9	86.98	85.05	87.32

Table 5: Intent Accuracy and Slot F1 of our multilingual BERT (mBERT) model on MultiAtis++.

Multilingual ATIS	hi	tr
<i>Multilingual models (trained on all data from all languages)</i>		
XLM-R	62.3 / 85.9 / 87.8	65.7 / 92.7 / 86.5
mT5-xxl	73.01 $\pm$ 0.30 / 95.04 $\pm$ 0.06 / 88.93 $\pm$ 0.09	70.68 $\pm$ 0.63 / 94.13 $\pm$ 0.37 / 87.69 $\pm$ 0.26
<i>Zero-shot models (trained on English only)</i>		
XLM-R	40.3 / 80.2 / 76.2	15.7 / 78 / 51.8
mT5-xxl	40.87 $\pm$ 8.91 / 91.41 $\pm$ 0.28 / 68.69 $\pm$ 7.47	14.78 $\pm$ 2.18 / 84.99 $\pm$ 0.53 / 51.29 $\pm$ 3.31
<i>Augmented data models</i>		
XLM-R + translate align	53.2 / 85.3 / 84.2	49.7 / 91.3 / 80.2
mT5-xxl + TaF	<b>65.29</b> $\pm$ 0.22 / <b>96.23</b> $\pm$ 0.17 / <b>84.85</b> $\pm$ 0.09	<b>67.41</b> $\pm$ 0.92 / <b>95.15</b> $\pm$ 0.21 / <b>85.30</b> $\pm$ 0.18
mT5-xxl + TAP	63.94 $\pm$ 0.30 / 96.04 $\pm$ 0.50 / 84.00 $\pm$ 0.39	58.41 $\pm$ 0.91 / 95.10 $\pm$ 0.14 / 82.40 $\pm$ 0.46

Table 6: Results of our mT5 models on Multilingual ATIS. Metrics are Exact Match (EM) accuracy / Intent Accuracy / Slot F1 respectively. XLM-R results are from Li et al. (2021).

much lower zero-shot accuracy. *mT5-xxl+TaF* is only 8.6 points behind *mT5-xxl* trained on all the available gold multilingual data and covers 60% of the gap between zero-shot and full multilingual supervision. *mT5-xxl+TaF* shows a remarkable improvement on German w.r.t. *XLM-R+TAP* and *mT5-xxl+TAP*, probably due to alignment errors caused by the heavy compounding nature of German, as Li et al. (2021) report in their paper too. In the *mT5-xxl+TaF, untokenized* experiment, we do not tokenize the translations for the filler. The results do not significantly change, suggesting that our approach is robust to tokenization and therefore tokenizers and aligners are not necessary.

Table 3 contains the results on MTOP when training *mt5-xxl* with English gold data and synthetic data generated by TAP for all the other languages. Out-of-the-box TAP is well behind zero-shot. With POS-based postprocessing, we see a significant improvement in all languages. Except for Thai, all languages are well above zero-shot performance. This shows that human-engineering is essential for TAP to perform well. Note that the preposition and determiner exclusion rule

is not being applied for the DATE\_TIME slots, according to the labeling trend we have observed in the MTOP dataset. On the other hand, the filler is able to learn this trend by itself and no heuristics are needed. The experiment where we keep only the whitespace-segmented synthetic data reaches the best performance with a significant bump in Thai, but it is still  $\sim$ 4 EM points below that of the filler on average. This shows that high-quality alignments are paramount for TAP to work well. The filler completely eliminates the need of an aligner and achieves better results. Note that for the other tables, we only included the results from the best TAP configuration.

**MultiATIS++.** In Table 4, we compare our approach with the mBERT models from Xu et al. (2020) that use synthetic data obtained by projecting labels with *fastalign* alignments (Dyer et al., 2013), attention weights (*softalign*) and *TMP* linguistic features (Jain et al., 2019). mT5-xxl has remarkable zero-shot Intent Accuracy and Slot F1, even without synthetic data. With the latter, the average Slot F1 is  $\sim$ 10 points higher than the

best mBERT baselines and the result variance decreases significantly w.r.t. zero-shot. Data produced with fastalign degrades performance for French and Turkish, while TaF synthetic data always leads to better results and contributes to set SOTA performance in data augmentation settings for mT5. We achieve a 52.83% Relative Reduction in Error (RRIE) in Slot F1 w.r.t. zero-shot, and compared to using all gold data, we close the full gap in Intent Accuracy and reduce the difference in Slot F1 to only  $\sim 4.5$  points. TaF consistently outperforms TAP in all languages, with more pronounced differences in Slot F1.

In Table 5 we report results on MultiAtis++ with an mBERT (Devlin et al., 2019) based parser, in order to understand how effective our synthetic data method is with models with less parameters and lower complexity than mT5. We use the cased mBERT encoder to obtain word representations: a linear layer on top of the [CLS] token performs intent classification, while a linear layer on top of the first wordpiece of each sentence token is used for slot tagging. Similarly to what we have observed in the mT5-xxl experiments, our synthetic data closes the full gap (between zero-shot and gold) in average Intent Accuracy, while the gap in average Slot F1 is reduced to less than 3 points. This shows the effectiveness of TaF for models with different power and capacity. It is worth noting that in the multilingual gold setting Spanish was underperforming TaF. The reason seemed to be that Spanish training data is mixcased while test data is lowercased. If we lowercase Spanish training data, we see a significant improvement in Spanish Slot F1. Note that the zero-shot and gold performance of our mBERT model is below that of the implementation in Xu et al. (2020). We suspect this is due mBERT model differences: our model has about 110M parameters while Xu et al. (2020) report more than 166M parameters. Despite the lower zero-shot performance, our mBERT model with TaF is more than 5 points better in average Slot F1 across 8 languages than the best data augmentation method from Xu et al. (2020).

**Multilingual ATIS.** Table 6 confirms the MultiAtis++ results also for this dataset. mT5-xxl is more effective than XLM-R at zero-shot intent classification but not at Slot F1. With the help of our synthetic data, *mT5-xxl+TaF* reaches SOTA

Language	Filler Errors	Zero-shot Parser Errors
de	459 (2.93%)	3607 (23.02%)
es	97 (0.62%)	2977 (19.00%)
fr	98 (0.63%)	3001 (19.15%)
hi	241 (1.53%)	4811 (30.71%)
th	1369 (8.74%)	5556 (35.46%)
Total	2264 (2.89%)	19952 (25.47%)

Table 7: Number and % of instances with errors that are matched by our heuristic filters.

performance on the task, recovering 51.6% and 69.8% of the full supervision gap in Slot F1 w.r.t zero-shot, on the Hindi and Turkish evaluation sets respectively. We observe TaF outperforms TAP on EM, particularly on highly agglutinative Turkish.

**General Remarks.** The relative deltas in performance across datasets on same languages may be explained by the heterogeneous domains and by the annotation structure. In addition, the starting pretrained models have different quality across languages as shown in “Zero-shot models” in our tables and as also noted in Conneau et al. (2020) (e.g., XLM-R performs particularly well on low-resource languages). Pretraining quality typically transfers to fine-tuned models.

## 6 Analysis of the Filler Output

In this section, we analyze the output of our filler trained on the English MTOP data and run on translated MTOP. We use two simple heuristic filters to understand how good the filler is at reproducing the signature provided in the input and how much it suffers from hallucination. Therefore we count (i) how many times the input and output signatures differ (ignoring slot orders); (ii) for how many utterances the output slots contain word spans which cannot be found in the input utterance.

Table 7 contains the number of examples (and the %s for each language) triggering our filters. The last row summarizes the numbers for the total of about 75k utterances ( $\sim 15$ k English MTOP training instances translated to 5 languages). In addition to the filler statistics, we compute the same numbers for a model that does not have access to the parse signature, i.e., a zero-shot parser trained on English. As we can see, the outputs of the filler contain errors in only 2.89% of cases. Of these, 0.5% parses are malformed, 3.7% have mistakes in the signatures and 96% have hallucination errors. We can conclude that the filler is able to reproduce input signatures and the only issues are due to wrong tokens put in the slots. On the contrary,



Lang	Utterance	Representation
<i>(1) Hallucination of pronouns</i>		
en	What reminders do <b>we</b> have this weekend ?	[IN: ... [SL:PERSON_REMINDED <b>we</b> ] [SL:DATE this weekend ] ]
es	Qué recordatorios hacemos este fin de semana ?	[IN: ... [SL:PERSON_REMINDED <b>nosotros</b> ] [SL:DATE este fin de semana ] ]
<i>(2) Confusion around prepositions and determiners</i>		
en	cancel reminder to call dentist	[IN: ... [SL:TODO [IN:CREATE_CALL [SL:CONTACT dentist ] ] ] ]
es	cancelar recordatorio para llamar <b>al</b> dentista	[IN: ... [SL:TODO [IN:CREATE_CALL [SL:CONTACT <b>el</b> dentista ] ] ] ]
<i>(3) Slot words reordering</i>		
en	What year did T. Woods <b>turn pro</b> ?	[IN: ... [SL:CONTACT T. Woods ] [SL:EVENT <b>turn pro</b> ] ]
es	En qué año <b>se convirtió</b> T. Woods <b>en profesional</b> ?	[IN: ... [SL:CONTACT T. Woods ] [SL:EVENT <b>se convirtió en profesional</b> ] ]
<i>(4) Hallucination of unaligned/missing words</i>		
en	Play some rap <b>music</b>	[IN: ... [SL:MUSIC_GENRE rap ] [SL:MUSIC_TYPE <b>music</b> ] ]
es	toca algo de rap	[IN: ... [SL:MUSIC_GENRE rap ] [SL:MUSIC_TYPE <b>music</b> ] ]
en	Will it be hot <b>out</b> today ?	[IN: ... [SL:WEATHER_ATTR hot ] [SL:LOC <b>out</b> ] [SL:TIME today ] ]
es	Va a hacer calor hoy ?	[IN: ... [SL:WEATHER_ATTR calor ] [SL:LOC <b>Alicante</b> ] [SL:TIME hoy ] ]
<i>(5) Compound or word splitting</i>		
en	Delete the <b>homework</b> reminder	[IN: ... [SL:TODO <b>homework</b> ] ]
de	Löschen Sie die <b>Hausaufgabenerinnerung</b>	[IN: ... [SL:TODO <b>Hausaufgaben</b> ] ]

Table 8: Examples where our filler generates spans that cannot be found in the input translation (in Spanish or German). *en* rows contain the original English utterance and parse. Intents are omitted and some slots are shortened for readability.

the 25% outputs with mistakes from the zero-shot parser are dominated by signature mistakes, which are 76% of the total. Hallucination errors amount to 28%.

Table 8 contains interesting examples matched by our heuristic filters. Hallucinations may happen when some words are dropped in the translation. In (1), the pronoun is dropped and the model generates the relevant first person plural pronoun in Spanish. In (4), the word “music” is not contained in the translation but still relevant, while “Alicante” is a quite random choice for the location slot. Other frequent issues are related to the choice of prepositions and determiners as in (2), where the latter is often preferred by the model. Example (3) is an interesting case of word reordering that highlights a well known issue in the i18n of span labeling annotations, namely span splitting. The generative filler is able to reorder the phrase back. Finally, we highlight example (5). In German, a compound rich language, the noun “homework” forms a compound with the noun “reminder”. The filler is able to split the compound noun, thanks to its subword output vocabulary, and put the relevant part in the TODO slot. How useful this is ultimately depends on the annotation guidelines defined for the i18n languages (e.g., allowing and supporting subword annotations). The relatively low number of errors and their nature explain why we are able to use all the synthetic data produced by our method to train the final parsers. We experimented by filtering out synthetic examples with the aforementioned heuris-

tics but we did not register any improvement on the final performance.

## 7 Conclusions and Future Work

In this paper, we proposed a novel Translate-and-Fill synthetic data generation approach which requires less engineering effort than TAP. TaF leverages NMT, multilingual pretrained seq2seq models and task labels, at the same time removing the need of aligners and tokenizers. Our filler model, trained on English data only, works remarkably well on other languages and enables improvements on multiple semantic parsing datasets in synthetic data scenarios. As future work, we plan to explore applications of the filler to (i) other i18n synthetic data generation tasks that require span alignment and to (ii) in-language data augmentation, e.g., using paraphrases to improve parsing accuracy of intent and slots with little annotated data.

## Acknowledgments

We would like to thank Melvin Johnson and the anonymous reviewers for their constructive feedback, useful comments and suggestions.

## Ethical Considerations

This work does not use sensitive data, or language models for uncontrolled generation. The output of the parser is not user facing, therefore engineers can easily intervene to eliminate potentially harmful hallucinations of the seq2seq model by simple

filtering. One concern could be the energy consumption of the experiments and to mitigate that, we did not perform hyper-parameter tuning and limited the experiment reruns.

## References

- Wasi Uddin Ahmad, Zhisong Zhang, Xuezhe Ma, Kai-Wei Chang, and Nanyun Peng. 2019. [Cross-lingual dependency parsing with unlabeled auxiliary languages](#). In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*, pages 372–382, Hong Kong, China. Association for Computational Linguistics.
- Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. [Abstract Meaning Representation for sembanking](#). In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 178–186, Sofia, Bulgaria. Association for Computational Linguistics.
- Peter F. Brown, Stephen A. Della Pietra, Vincent J. Della Pietra, and Robert L. Mercer. 1993. [The mathematics of statistical machine translation: Parameter estimation](#). *Computational Linguistics*, 19(2):263–311.
- Yun Chen, Yang Liu, Guanhua Chen, Xin Jiang, and Qun Liu. 2020. [Accurate word alignment induction from neural machine translation](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 566–576, Online. Association for Computational Linguistics.
- Jonathan H. Clark, Dan Garrette, Iulia Turc, and John Wieting. 2021. [CANINE: pre-training an efficient tokenization-free encoder for language representation](#). *CoRR*, abs/2103.06874.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. [Unsupervised cross-lingual representation learning at scale](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451, Online. Association for Computational Linguistics.
- Angel Daza and Anette Frank. 2020. [X-SRL: A parallel cross-lingual semantic role labeling dataset](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 3904–3914, Online. Association for Computational Linguistics.
- 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. Association for Computational Linguistics.
- Li Dong and Mirella Lapata. 2018. [Coarse-to-fine decoding for neural semantic parsing](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 731–742, Melbourne, Australia. Association for Computational Linguistics.
- Zi-Yi Dou and Graham Neubig. 2021. [Word alignment by fine-tuning embeddings on parallel corpora](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2112–2128, Online. Association for Computational Linguistics.
- Chris Dyer, Victor Chahuneau, and Noah A. Smith. 2013. [A simple, fast, and effective reparameterization of IBM model 2](#). In *Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 644–648, Atlanta, Georgia. Association for Computational Linguistics.
- Akiko Eriguchi, Melvin Johnson, Orhan Firat, Hideto Kazawa, and Wolfgang Macherey. 2018. [Zero-shot cross-lingual classification using multilingual neural machine translation](#). *CoRR*, abs/1809.04686.
- Xing Fan, Emilio Monti, Lambert Mathias, and Markus Dreyer. 2017. [Transfer learning for neural semantic parsing](#). In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pages 48–56, Vancouver, Canada. Association for Computational Linguistics.
- Hao Fei, Meishan Zhang, and Donghong Ji. 2020. [Cross-lingual semantic role labeling with high-quality translated training corpus](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7014–7026, Online. Association for Computational Linguistics.
- Chih-Wen Goo, Guang Gao, Yun-Kai Hsu, Chih-Li Huo, Tsung-Chieh Chen, Keng-Wei Hsu, and Yun-Nung Chen. 2018. [Slot-gated modeling for joint slot filling and intent prediction](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 753–757, New Orleans, Louisiana. Association for Computational Linguistics.
- Jiang Guo, Wanxiang Che, David Yarowsky, Haifeng Wang, and Ting Liu. 2015. [Cross-lingual dependency parsing based on distributed representations](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages

- 1234–1244, Beijing, China. Association for Computational Linguistics.
- Sonal Gupta, Rushin Shah, Mrinal Mohit, Anuj Kumar, and Mike Lewis. 2018. [Semantic parsing for task oriented dialog using hierarchical representations](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 2787–2792, Brussels, Belgium. Association for Computational Linguistics.
- Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. [XTREME: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation](#). In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4411–4421. PMLR.
- Alankar Jain, Bhargavi Paranjape, and Zachary C. Lipton. 2019. [Entity projection via machine translation for cross-lingual NER](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1083–1092. Association for Computational Linguistics.
- Masoud Jalili Sabet, Philipp Dufter, François Yvon, and Hinrich Schütze. 2020. [SimAlign: High quality word alignments without parallel training data using static and contextualized embeddings](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1627–1643, Online. Association for Computational Linguistics.
- Prabhu Kaliamoorthi, Aditya Siddhant, Edward Li, and Melvin Johnson. 2021. [Distilling large language models into tiny and effective students using pqrrn](#). *CoRR*, abs/2101.08890.
- Katharina Kann, Ophélie Lacroix, and Anders Søgaard. 2020. [Weakly supervised pos taggers perform poorly on truly low-resource languages](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8066–8073.
- Taku Kudo and John Richardson. 2018. [SentencePiece: A simple and language independent subword tokenizer and detokenizer for neural text processing](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 66–71, Brussels, Belgium. Association for Computational Linguistics.
- Faisal Ladhak, Esin Durmus, Claire Cardie, and Kathleen McKeown. 2020. [WikiLingua: A new benchmark dataset for cross-lingual abstractive summarization](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4034–4048, Online. Association for Computational Linguistics.
- Guillaume Lample and Alexis Conneau. 2019. [Cross-lingual language model pretraining](#). *CoRR*, abs/1901.07291.
- Anne Lauscher, Vinit Ravishankar, Ivan Vulić, and Goran Glavaš. 2020. [From zero to hero: On the limitations of zero-shot language transfer with multilingual Transformers](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 4483–4499, Online. Association for Computational Linguistics.
- Haoran Li, Abhinav Arora, Shuohui Chen, Ankit Gupta, Sonal Gupta, and Yashar Mehdad. 2021. [MTOP: A comprehensive multilingual task-oriented semantic parsing benchmark](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2950–2962, Online. Association for Computational Linguistics.
- Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119.
- Hiroki Nakayama. 2018. [seqeval: A python framework for sequence labeling evaluation](#). Software available from <https://github.com/chakki-works/seqeval>.
- S. B. Needleman and C. D. Wunsch. 1970. A general method applicable to the search for similarities in the amino acid sequence of two proteins. *Journal of molecular biology*, 48 3:443–53.
- Jian Ni, Georgiana Dinu, and Radu Florian. 2017. [Weakly supervised cross-lingual named entity recognition via effective annotation and representation projection](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1470–1480. Association for Computational Linguistics.
- Franz Josef Och and Hermann Ney. 2000. [Improved statistical alignment models](#). In *Proceedings of the 38th Annual Meeting of the Association for Computational Linguistics*, pages 440–447, Hong Kong. Association for Computational Linguistics.
- Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models. *Computational Linguistics*, 29(1):19–51.
- Giovanni Paolini, Ben Athiwaratkun, Jason Krone, Jie Ma, Alessandro Achille, Rishita Anubhai, Cicero Nogueira dos Santos, Bing Xiang, and Stefano Soatto. 2021. Structured prediction as translation between augmented natural languages. In *9th International Conference on Learning Representations, ICLR 2021*.
- Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. [GloVe: Global vectors for word representation](#). In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language*



- Processing (EMNLP)*, pages 1532–1543, Doha, Qatar. Association for Computational Linguistics.
- Barbara Plank and Željko Agić. 2018. [Distant supervision from disparate sources for low-resource part-of-speech tagging](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 614–620. Association for Computational Linguistics.
- P. J. Price. 1990. [Evaluation of spoken language systems: the ATIS domain](#). In *Speech and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990*.
- Libo Qin, Wanxiang Che, Yangming Li, Haoyang Wen, and Ting Liu. 2019. [A stack-propagation framework with token-level intent detection for spoken language understanding](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2078–2087, Hong Kong, China. Association for Computational Linguistics.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Subendhu Rongali, Luca Soldaini, Emilio Monti, and Wael Hamza. 2020. [Don’t parse, generate! a sequence to sequence architecture for task-oriented semantic parsing](#). In *Proceedings of The Web Conference 2020, WWW ’20*, page 2962–2968, New York, NY, USA. Association for Computing Machinery.
- Sebastian Schuster, Sonal Gupta, Rushin Shah, and Mike Lewis. 2019. [Cross-lingual transfer learning for multilingual task oriented dialog](#). 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 3795–3805. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Rushin Shah, Bo Lin, Anatole Gershman, Robert Frederking, and Microsoft Bing Translatortm. 2010. [Synergy: A named entity recognition system for resource-scarce languages such as swahili using online machine translation](#). In *Proceedings of International Conference on Language Resource and Evaluation Workshop on African Language Technology*.
- Aditya Siddhant, Melvin Johnson, Henry Tsai, Naveen Ari, Jason Riesa, Ankur Bapna, Orhan Firat, and Karthik Raman. 2020. [Evaluating the cross-lingual effectiveness of massively multilingual neural machine translation](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(05):8854–8861.
- Karan Singla, Dogan Can, and Shrikanth Narayanan. 2018. [A multi-task approach to learning multilingual representations](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 214–220. Association for Computational Linguistics.
- Oscar Täckström, Dipanjan Das, Slav Petrov, Ryan McDonald, and Joakim Nivre. 2013. [Token and type constraints for cross-lingual part-of-speech tagging](#). *Transactions of the Association for Computational Linguistics*, 1:1–12.
- Chen-Tse Tsai, Stephen Mayhew, and Dan Roth. 2016. [Cross-lingual named entity recognition via wikification](#). In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, pages 219–228. Association for Computational Linguistics.
- Shyam Upadhyay, Manaal Faruqui, Gokhan Tur, Dilek Hakkani-Tur, and Larry Heck. 2018. (almost) zero-shot cross-lingual spoken language understanding. In *Proceedings of the IEEE ICASSP*.
- Stephan Vogel, Hermann Ney, and Christoph Tillmann. 1996. [HMM-based word alignment in statistical translation](#). In *COLING 1996 Volume 2: The 16th International Conference on Computational Linguistics*.
- Jason Williams, Antoine Raux, and Matthew Henderson. 2016. [The dialog state tracking challenge series: A review](#). *Dialogue & Discourse*, 7:4–33.
- Sam Wiseman, Stuart Shieber, and Alexander Rush. 2018. [Learning neural templates for text generation](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3174–3187, Brussels, Belgium. Association for Computational Linguistics.
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. [Google’s neural machine translation system: Bridging the gap between human and machine translation](#). *CoRR*, abs/1609.08144.
- Jiateng Xie, Zhilin Yang, Graham Neubig, Noah A. Smith, and Jaime Carbonell. 2018. [Neural cross-lingual named entity recognition with minimal resources](#). In *Proceedings of the 2018 Conference on*



- Empirical Methods in Natural Language Processing*, pages 369–379. Association for Computational Linguistics.
- Weijia Xu, Batool Haider, and Saab Mansour. 2020. [End-to-end slot alignment and recognition for cross-lingual NLU](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5052–5063, Online. Association for Computational Linguistics.
- Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and Colin Raffel. 2021. [mt5: A massively multilingual pre-trained text-to-text transformer](#).
- David Yarowsky, Grace Ngai, and Richard Wicentowski. 2001. [Inducing multilingual text analysis tools via robust projection across aligned corpora](#). In *Proceedings of the First International Conference on Human Language Technology Research*.
- Katherine Yu, Haoran Li, and Barlas Oguz. 2018. [Multilingual seq2seq training with similarity loss for cross-lingual document classification](#). In *Proceedings of The Third Workshop on Representation Learning for NLP*, pages 175–179. Association for Computational Linguistics.
- Thomas Zenkel, Joern Wuebker, and John DeNero. 2020. [End-to-end neural word alignment outperforms GIZA++](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1605–1617, Online. Association for Computational Linguistics.
- Meishan Zhang, Yue Zhang, and Guohong Fu. 2019a. [Cross-lingual dependency parsing using code-mixed TreeBank](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 997–1006, Hong Kong, China. Association for Computational Linguistics.
- Xiang Zhang, Shizhu He, Kang Liu, and Jun Zhao. 2019b. [AdaNSP: Uncertainty-driven adaptive decoding in neural semantic parsing](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4265–4270, Florence, Italy. Association for Computational Linguistics.
- Ayah Zirikly and Masato Hagiwara. 2015. [Cross-lingual transfer of named entity recognizers without parallel corpora](#). In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 390–396. Association for Computational Linguistics.