

Contextual Label Projection for Cross-Lingual Structured Prediction

Anonymous ACL submission

Abstract

Label projection, which involves obtaining translated labels and texts jointly, is essential for leveraging machine translation to facilitate cross-lingual transfer in structured prediction tasks. Prior research exploring label projection often compromises translation accuracy in favor of simplified label identification or suffers from inaccuracies by relying solely on word alignment for constructing label phrases. In this paper, we introduce a novel label projection approach, CLAP, which translates text to the target language and performs *contextual translation* on the labels using the translated text as the context, ensuring better accuracy for the translated labels. We leverage instruction-tuned language models with multilingual capabilities as our contextual translator, imposing the constraint of the presence of translated labels in the translated text via instructions. We compare CLAP with other label projection techniques on zero-shot cross-lingual transfer across 39 languages on two representative structured prediction tasks — event argument extraction (EAE) and named entity recognition (NER). Experiments reveal that CLAP improves by 1.7 F1 points for EAE and by 1.4 F1 points for NER.

1 Introduction

Cross-lingual transfer for structured prediction tasks such as named entity recognition, relation extraction, and event extraction, has gained considerable attention recently (Huang et al., 2022; Cao et al., 2023; Tedeschi and Navigli, 2022; Cabot et al., 2023; Fincke et al., 2022; Jenkins et al., 2023; Ahmad et al., 2021b). It generalizes models trained in a source languages to other target languages, broadening the scope of these applications to more languages (Chen and Ritter, 2021; Subburathinam et al., 2019; Pouran Ben Veyseh et al., 2022).

One effective and simple way to improve cross-lingual transfer performance is translate-train,

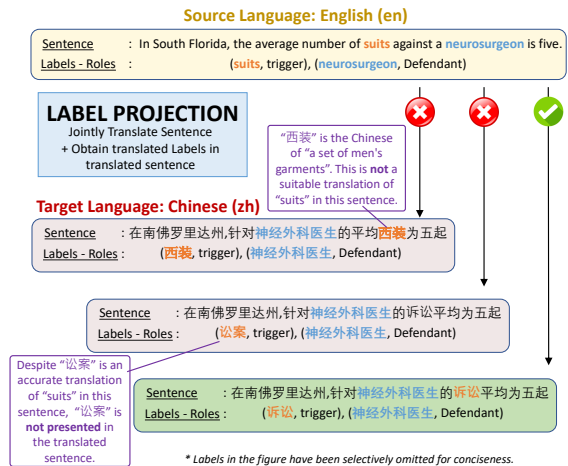


Figure 1: Illustration of the task of *label projection* from English to Chinese. Label projection converts sentences from a source to a target language while translating the associated labels jointly. Failures in this process occur when labels are either inaccurately translated or missing in the translated sentence in the target language.

which leverages machine translation to generate pseudo-training data in the target languages by translating source language training data (Xue et al., 2021; Ruder et al., 2021; Yu et al., 2023). However, applying this technique to structured prediction tasks necessitates a *label projection* step, which involves jointly translating input sentences and labels (Chen et al., 2023). Label projection requires not only *accurate translation* of the labels but also *maintaining the association* between the translated texts and labels. As illustrated in Figure 1, while “suits” can have multiple valid translations, only “诉讼” is presented in the translated sentence and is a proper translation at the same time.

Prior works have dealt with label projection through two primary frameworks. The first one, illustrated in Figure 2(a), performs machine translation on modified source sentences that incorporate label annotations using special markers (Chen

et al., 2023; Hennig et al., 2023). Translated labels can be extracted if special markers are retained in the translations. In this approach, the quality of the translation is *inherently compromised* due to the inclusion of special markers (Chen et al., 2023). The other framework uses word similarity to procure word alignments between the source and translated sentences. Label translations are further constructed by combining mapped tokens in the translated sentence (Stengel-Eskin et al., 2019; Akbik et al., 2015; Aminian et al., 2019), as shown in Figure 2(b). However, it is hard for this framework to ensure *accurate* label translation by merely using word alignments, as we will show in Section 4.4.

In this work, we introduce CLAP (Contextual Label Projection), which obtains projected label annotations by utilizing contextual machine translation for the labels. We first acquire the translation of the whole input sentence by any plug-and-play machine translation model. Then, inspired by the idea of contextual machine translation (Wong et al., 2020; Voita et al., 2018), we use the translated input text as context to perform label translation, as shown in Figure 2(c). Exploiting contextual machine translation strongly enhances the *accuracy* of the translated labels while preserving their *association* to the translated sentence. Furthermore, translating the input sentence in an unmodified manner better exploits machine translators, and in turn, assures high quality of the translated sentence.

To implement contextual machine translation, we utilize an instruction-tuned language model with multilingual capabilities, Llama-2 (Touvron et al., 2023). We encode the translated input sentence and the constraint for the presence of labels in the form of instruction prompts. Despite sacrificing some translation ability compared to supervised machine translation models (Zhu et al., 2023), instruction-tuned language models provide better understanding of contextual constraints.

We experiment on the tasks of event argument extraction (EAE) and named entity recognition (NER) using the ACE dataset (Doddington et al., 2004) and the WikiANN dataset (Pan et al., 2017), covering 39 different languages in total. Our experiments show that utilizing label-projected data from CLAP for translate-train yields an average improvement of 1.7 and 1.4 F1 scores over strong baselines for EAE and NER respectively. We also perform an intrinsic evaluation using human study in Chinese, Arabic, Hindi, and Spanish to assess the projected

labels’ quality which shows how CLAP provides more accurate label translations while preserving the label presence in the translated sentence. Further analyses also reveal how CLAP generalizes for different translation models and works effectively for the translate-test paradigm as well. These evaluations and robust analyses underscore the effectiveness of CLAP for label projection.

2 Background

2.1 Structure Prediction Tasks

Given an input sentence \mathbf{x} , structure prediction models aim to predict structure output $\mathbf{y} = [\mathbf{x}[i_1 : j_1], \mathbf{x}[i_2 : j_2], \dots, \mathbf{x}[i_n : j_n]]$ (where $\mathbf{x}[i_1 : j_1]$ is an input sentence span from token i_1 to j_1) corresponding to a set of roles $\mathbf{r} = [r_1, r_2, \dots, r_n]$ (where $r_i \in \mathcal{R}$, a pre-defined set of roles). This vastly differs from standard classification-based tasks wherein the output prediction y is a singular value from a fixed set of classes independent of the input sentence \mathbf{x} .

2.2 Zero-shot Cross-Lingual Transfer

Zero-shot cross-lingual transfer (Hu et al., 2020; Ahmad et al., 2019; Huang et al., 2021; Hsu et al., 2023b) aims to train a downstream model for the target language l_{tgt} using supervised data \mathcal{D}_{src} from a source language l_{src} without using any data in the target language (i.e. $\mathcal{D}_{tgt} = \phi$). The paradigm has effectively advanced language technologies for under-resourced languages.

2.3 Translate-Train

Translate-train (Hu et al., 2020; Ruder et al., 2021) is a popular and powerful zero-shot cross-lingual transfer technique that leverages machine translators \mathcal{T} to boost downstream model performance. Specifically, in translate-train, \mathcal{D}_{src} is translated into the target language as pseudo training data \mathcal{D}_{src}^{tgt} and the downstream model is trained using a combination of $\{\mathcal{D}_{src}, \mathcal{D}_{src}^{tgt}\}$.

Utilizing translate-train for structured prediction tasks requires *Label Projection*, which includes two sets of translations: (1) Sentence translation ($\mathbf{x}^{src} \xrightarrow{\mathcal{T}} \mathbf{x}^{tgt}$), where we use $\xrightarrow{\mathcal{T}}$ to denote the translation from l_{src} to l_{tgt} using \mathcal{T} ; and (2) Label translation ($\mathbf{y}^{src} \rightarrow \mathbf{y}^{tgt}$), such that the translated label \mathbf{y}^{tgt} is appropriately *associated with* \mathbf{x}^{tgt} . This demand makes translate-train for structure prediction tasks more complex than that for classification

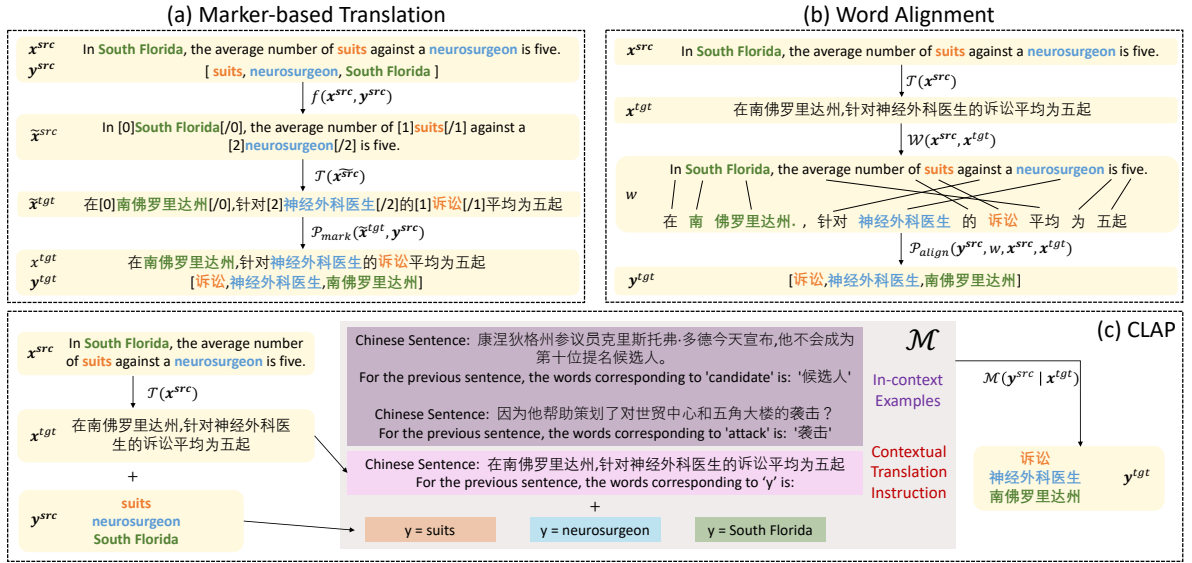


Figure 2: Illustration of the various techniques to conduct label projection: (a) **Marker-based Translation** use markers to transform the sentence and translate the transformed sentence with label markers jointly, (b) **Word Alignment** methods use external word alignment tools to locate the translated labels in the translated sentence, and (c) **CLAP** performs contextual translation on labels using \mathcal{M} (here we show instruction-tuned language model as \mathcal{M}) to locate translated label in the translated sentence.

tasks, as the latter only requires sentence translation (since y is a value independent of \mathbf{x}).¹

Translate-Test Besides translate-train, translate-test is another commonly used technique in zero-shot cross-lingual transfer. During testing time, it uses models solely trained on \mathcal{D}_{src} to make predictions on translated test sentences ($\mathbf{x}^{tgt} \xrightarrow{\mathcal{T}} \mathbf{x}^{src}$), and then uses label projection to map predictions on \mathbf{x}^{src} back to predictions on \mathbf{x}^{tgt} . Since it will cause additional error propagation issues during inference time, we mainly focus on translate-train in this paper. However, we discuss CLAP’s utilization and effectiveness on translate-test in Section 5.4.

2.4 Label Projection

We hereby technically define the problem of *label projection* (Akbik et al., 2015; Chen et al., 2023):

$$\begin{aligned} & \mathbf{x}^{src} \xrightarrow{\mathcal{T}} \mathbf{x}^{tgt} \\ & \& y_m^{src} \rightarrow y_m^{tgt} \quad \forall y_m^{src} \in \mathbf{y} \\ \text{s.t. } & y_m^{tgt} \in \mathbf{x}^{tgt} \quad \forall y_m^{tgt} \in \mathbf{y}^{tgt}. \end{aligned}$$

This problem requires optimizing two properties of **accuracy** and **faithfulness** on the translations. Accuracy ensures that $[\mathbf{x}^{tgt}, y_1^{tgt}, y_2^{tgt}, \dots, y_n^{tgt}]$ are

¹For certain structure prediction tasks like relation classification (determining the relationship between two entities in \mathbf{x}), even if the output y is scalar, translate-train necessitates label projection step due to the required projection of the two given entities into the translated sentence.

accurate translations of $[\mathbf{x}^{src}, y_1^{src}, y_2^{src}, \dots, y_n^{src}]$. On the other hand, faithfulness ensures that each y_m^{tgt} is associated with \mathbf{x}^{tgt} (the constraint of $y_m^{tgt} \in \mathbf{x}^{tgt}$). Standard translation models \mathcal{T} trained on supervised sentence translation pairs cannot simply impose the additional faithfulness constraint, such as the failure cases shown in Figure 1. This demonstrates the challenges of the label projection.

3 Methodology

In this section, we first formally define the previous attempts at label projection and later introduce CLAP, which provides a new perspective of using contextual machine translation for label projection.

3.1 Baseline Methods

As stated in Section 1, two primary frameworks, Marker-based translations and word-alignment-based methods, are primarily used in prior works.

Marker-based Translations solve the label projection by first marking labels to the input sentence \mathbf{x}^{src} , forming $\tilde{\mathbf{x}}^{src}$, and then use the translation model to obtain the potential translation of input sentence and labels jointly (Lewis et al., 2020; Hu et al., 2020; Chen et al., 2023). For example, in Figure 2(a), “South Florida” is delineated by markers [0] and [0]. Assuming the preservation of markers after translation of $\tilde{\mathbf{x}}^{src}$, a post-processing step, \mathcal{P}_{mark} , is performed to retain the translated labels

\mathbf{y}^{tgt} and translated sentence \mathbf{x}^{tgt} . Putting every step together, we have

$$\begin{aligned}\tilde{\mathbf{x}}^{src} &= f(\mathbf{x}^{src}, \mathbf{y}^{src}), & \tilde{\mathbf{x}}^{tgt} &= \mathcal{T}(\tilde{\mathbf{x}}^{src}) \\ \mathbf{x}^{tgt}, \mathbf{y}^{tgt} &= \mathcal{P}_{mark}(\tilde{\mathbf{x}}^{tgt}, \mathbf{y}^{src}),\end{aligned}$$

where f denotes the marker addition step and $\tilde{\mathbf{x}}^{tgt}$ is the translation of $\tilde{\mathbf{x}}^{src}$ using translator \mathcal{T} .

Despite their simplicity, these methods suffer from poor translation quality and reduced robustness to different translation models owing to their input sentence transformations and strong assumptions about the retention of markers in $\tilde{\mathbf{x}}^{tgt}$.

Word Alignment approaches (Akbik et al., 2015; Yarmohammadi et al., 2021) first translate the input sentence and acquire word alignments (Dyer et al., 2013; Dou and Neubig, 2021) between the translation pairs. Each translated label y_m^{tgt} is then procured by merging the aligned words of y_m^{src} in the translated sentence using the word mappings w . For example, in Figure 2(b), the translated label for ‘‘South Florida’’ is obtained by merging two aligned words, which is done by a heuristic post-processing algorithm \mathcal{P}_{align} . Formally, we have

$$\begin{aligned}\mathbf{x}^{tgt} &= \mathcal{T}(\mathbf{x}^{src}), & w &= \mathcal{W}(\mathbf{x}^{src}, \mathbf{x}^{tgt}) \\ y_m^{tgt} &= \mathcal{P}_{align}(y_m^{src}, w, \mathbf{x}^{src}, \mathbf{x}^{tgt}) & \forall y_m^{src} \in \mathbf{y}^{src}\end{aligned}$$

Although these approaches provide high-quality sentence translations, their translated labels can be error-prone as they use simple word alignment modules for capturing word-level translation relations without considering the entire label for translation (Akbik et al., 2015; Chen et al., 2023).

3.2 CLAP

We tackle the task of label projection through a new perspective — performing actual translation on labels instead of recovering them from translated text \mathbf{x}^{tgt} . This better ensures the accuracy of the translated labels \mathbf{y}^{tgt} . To accomplish this, we leverage the idea of *contextual machine translation* on the label translation with \mathbf{x}^{tgt} as context.

Contextual machine translation, which aims to perform phrase-level translations conditional on the context of the translated sentence, is tangentially explored for applications like anaphora resolution (Voita et al., 2018) and pronoun translations (Wong et al., 2020). The main goal of this task is to maintain the consistency of phrasal translations in the given context. In our work, we develop a novel

model CLAP to extend the idea of contextual translation to the application of label projection.

As illustrated in Figure 2(c), CLAP first utilizes machine translation model \mathcal{T} to translate input sentence \mathbf{x}^{src} to \mathbf{x}^{tgt} . Treating \mathbf{x}^{tgt} as the context, the contextual translation model \mathcal{M} translates the labels \mathbf{y}^{src} to \mathbf{y}^{tgt} . Contextual translation implicitly imposes the *faithfulness constraint* which requires $y_m^{tgt} \in \mathbf{x}^{tgt} \quad \forall y_m^{tgt} \in \mathbf{y}^{tgt}$, hence, slackly satisfying the requirement of label projection. These two steps can be formally described as

$$\begin{aligned}\mathbf{x}^{tgt} &= \mathcal{T}(\mathbf{x}^{src}) \\ y_m^{tgt} &= \mathcal{M}(y_m^{src} | \mathbf{x}^{tgt}) & \forall y_m^{src} \in \mathbf{y}^{src}\end{aligned}$$

where y_m^{tgt} is generated from $\mathcal{M}(y_m^{src} | \mathbf{x}^{tgt})$, drawing the significant difference from the previous works.

Compared to word alignment approaches using simple word-similarity aligners \mathcal{W} , we use models with translation capabilities \mathcal{M} , to improve the accuracy of translated labels. Furthermore, the independence of \mathcal{T} and \mathcal{M} for translating \mathbf{x}^{src} and \mathbf{y}^{src} respectively assures that CLAP has better translation quality for \mathbf{x}^{tgt} and is more robust than the marker-based baselines. We empirically back these intuitions in § 4.4.

3.3 Implementing CLAP

Putting our idea into practice, we configure \mathcal{T} to be a modular component that can be replaced by any third-party translation model. For \mathcal{M} , we use an instruction-tuned language model with multilingual capabilities. Instruction-tuned language models can accept conditional information in their natural language prompt. Specifically, we encode the translated target sentence \mathbf{x}^{tgt} as well as the faithfulness constraint $y_m^{tgt} \in \mathbf{x}^{tgt}$ implicitly in the form of natural language instructions (highlighted as ‘‘Contextual Translation Instruction’’ in Figure 2(c)). Following Brown et al. (2020), we also provide n randomly chosen in-context examples (highlighted as ‘‘In-context examples’’ in Figure 2(c)) to improve the instruction-understanding capability of the model.² Lastly, we use simple string-matching algorithms to get the exact span index of y_m^{tgt} in \mathbf{x}^{tgt} . Although this may not be the optimal solution when duplicated strings exist in \mathbf{x}^{tgt} , it works well in practice as stated in prior word-alignment methods (Dou and Neubig, 2021).

²The in-context examples are generated using Google translation and initial prediction from instruction-tuned LMs. The label predictions are further verified by back-translation.

	ACE	WikiANN
# Train Instances	4,202	20,000
# Dev Instances	450	10,000
# Avg. Test Instances	194	6,469
# Test Languages	2	39

Table 1: High-level data statistics for ACE and WikiANN datasets for EAE and NER tasks respectively. # = ‘number of’ and Avg. = average.

4 Experiments and Results

This section describes our experimental setup comprising the datasets, baselines, and implementation details. Later, we present both intrinsic and extrinsic evaluations of CLAP.

4.1 Task and Dataset

We choose two structure prediction tasks, event argument extraction (EAE) (Sundheim, 1992; Hsu et al., 2023a) and named entity recognition (NER) (Tjong Kim Sang, 2002; Tjong Kim Sang and De Meulder, 2003) for evaluating our label projection method. EAE requires the extraction of text segments serving as arguments corresponding to an event and mapping them to their corresponding argument roles. NER aims to identify and categorize named entities from the input sentence. We use the multilingual ACE dataset (Doddington et al., 2004) and the WikiANN (Pan et al., 2017; Rahimi et al., 2019) for benchmarking EAE and NER, respectively. We consider the zero-shot cross-lingual transfer using English (*en*) as the source language for both tasks. For ACE, we follow the pre-processing by Huang et al. (2022) to retain 33 event types and 22 argument roles. For WikiANN, we utilize pre-processing by Hu et al. (2020). We provide the high-level statistics for these datasets in Table 1. More details can be found in § A.

4.2 Baselines

We select two label projection models as baselines, each representing the two baseline frameworks we covered in Section 3.1, respectively: (1) **EasyProject** (Chen et al., 2023), a recent marker-based translation technique, utilizes numbered square braces (e.g. [0] and [/0]) to mark the labels in the input sentence. (2) **Awesome-Align** (Dou and Neubig, 2021), a neural bilingual word alignment model, uses multilingual language models to find word similarities to derive word alignments, which are later used for label projection.

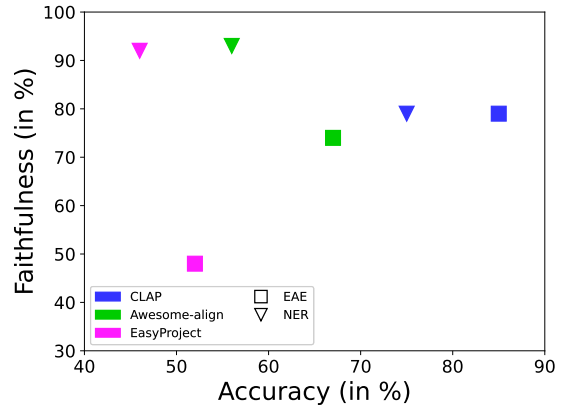


Figure 3: Reporting faithfulness and accuracy (in %) for the different label projection models on EAE and NER datasets. The closer the model is to the top-right, the better it is.

	mBART		mT5		mT5+Copy		Avg
	ar	zh	ar	zh	ar	zh	
Zero-shot*	36.3	47.3	36.7	51.0	40.3	51.9	43.9
Awesome-align	45.2	49.4	46.8	53.7	48.6	54.5	49.7
EasyProject	37.9	52.3	34.5	54.6	38.5	56.3	45.7
CLAP (ours)	46.0	53.4	44.3	56.5	49.3	58.6	51.4

Table 2: Extrinsic evaluation of the different label projection techniques regarding downstream model performance using translate-train for EAE. Avg = Average. * indicates the reproduced results of our base zero-shot cross-lingual EAE model, X-Gear (Huang et al., 2022).

4.3 Implementation Details

For the translation model \mathcal{T} , we experiment with the Google Machine Translation (GMT).³ For CLAP, we use the text-completion version of Llama-2 (Touvron et al., 2023) with 13B parameters as \mathcal{M} . We use $n = 2$ in-context examples for CLAP prompts. For Awesome-align, we use the unsupervised version of their model utilizing multilingual BERT (Devlin et al., 2019) as it provides better results (Chen et al., 2023).

4.4 Intrinsic Evaluation

We first evaluate CLAP by directly evaluating the label projection quality, mainly focusing on evaluating the accuracy and faithfulness of the translated labels, with the definition stated in Section 2.4.

Accuracy is measured in terms of translation quality by comparing the source labels with the translated labels ($\mathbf{y}^{src} \leftrightarrow \mathbf{y}^{tgt}$). We hire native human speakers to rank the translated labels by the

³<https://cloud.google.com/translate>

Lang	af	ar	bg	bn	de	el	es	et	eu	fa	fi	fr	he	hi	hu	id	it	ja	jv	ka
Zero-shot	77.4	48.1	82.8	77.0	78.8	80.6	74.5	78.7	61.4	69.2	79.3	79.4	57.3	70.6	80.8	53.1	79.4	19.1	58.5	72.3
Awesome-align	77.9	46.0	81.0	81.2	78.8	71.7	65.3	78.0	66.8	46.4	77.4	78.2	55.3	73.9	77.4	52.8	79.3	20.3	56.3	70.4
EasyProject	76.1	34.4	81.0	78.6	78.8	69.3	70.5	73.9	54.8	49.1	77.8	78.8	61.1	73.0	75.6	51.0	79.0	41.3	62.4	66.4
CLAP	74.4	48.7	81.0	78.1	78.4	75.9	74.7	77.4	68.8	59.0	75.9	79.4	58.4	73.1	72.4	56.1	80.1	45.3	64.8	70.5
	kk	ko	ml	mr	ms	my	nl	pt	ru	sw	ta	te	th	tl	tr	ur	vi	yo	zh	Avg
Zero-shot	51.9	57.5	66.4	65.3	53.4	65.8	83.0	80.0	74.2	68.4	60.3	62.1	0.4	74.5	65.6	62.2	75.0	34.1	24.6	64.2
Awesome-align	47.7	57.7	63.4	62.4	70.7	54.1	83.0	75.8	64.8	70.1	62.4	55.4	2.4	80.9	62.8	53.7	66.4	61.5	45.4	63.5
EasyProject	31.7	48.2	56.5	59.8	71.7	60.3	81.9	79.6	66.3	71.5	53.2	54.2	11.4	78.2	66.8	63.8	65.6	68.8	42.0	63.2
CLAP	42.8	60.1	60.3	61.4	73.5	61.5	82.2	78.2	68.3	70.6	59.6	53.1	13.2	74.6	62.9	32.9	75.8	59.6	49.7	64.9

Table 3: Extrinsic evaluation of the different label projection techniques in terms of downstream model performance using translate-train for NER. Avg = Average.

different models based on their translation quality. We conduct this evaluation on 50 data samples for four languages - Chinese, Arabic, Hindi, and Spanish, respectively. The final accuracy score for each model is the average percentage when the given methods provided the best quality translation for the labels among the other competitors.

Faithfulness measures the fulfillment of the label projection constraint. It is measured as a percentage of projected data points when all the translated labels are present in the translated input sentence ($y_m^{tgt} \in \mathbf{x}^{tgt}$, $\forall y_m^{tgt} \in \mathbf{y}^{tgt}$). The statistics use the complete test set on ACE and WikiANN.

4.4.1 Results

The accuracy and faithfulness of the models are plotted together in Figure 3. An ideal model should optimize both these metrics and thus, the closer the models are to the top-right, the better they are deemed. Overall, this figure shows how CLAP performs the best intrinsically as it is the closest to the top-right for both the tasks. For EAE, CLAP is better than all models in both the metrics, while for NER, CLAP compromises faithfulness slightly for stronger accuracy. Awesome-align and EasyProject are both great at attaining higher projection rates but produce more inaccurate label translations. Overall, intrinsic evaluation reveals how CLAP provides the best balance of accuracy and faithfulness.

4.5 Extrinsic Evaluation

Extrinsic evaluation implicitly evaluates the label projection techniques’ ability to generate good-quality data for downstream tasks. The projected data is utilized to train downstream models using the translate-train paradigm together with the original training data in English. For translate-train, we only retain the projected datapoints that satisfy the

faithfulness constraint as part of the target pseudo-training data \mathcal{D}_{src}^{tgt} .

EAE For EAE downstream model, we use the state-of-the-art model for zero-shot cross-lingual EAE: X-Gear (Huang et al., 2022). We explore three versions of the X-Gear model: mBART without copy (mBART), mT5 without copy (mT5), and mT5 with copy mechanism (mT5+Copy). We present the results in terms of argument classification F1 scores⁴ in Table 2. For reference, we also include the zero-shot baseline (training only on \mathcal{D}_{src}). Evidently, CLAP performs the best providing an average gain of 1.7 F1 points over the next best baseline of Awesome-align and a net gain of 7.5 F1 points over the zero-shot baseline. This result is in sync with our intrinsic evaluation wherein CLAP performed the best for EAE.

NER For NER, we utilize XLM-RoBERTa_{large} (Conneau et al., 2020) as our downstream model and use the XTREME (Hu et al., 2020) setup for implementation. The main results for entity classification F1 scores are presented in Table 3 along with the zero-shot baseline. Overall, CLAP performs the best with an absolute improvement of 0.7 F1 points over the zero-shot baseline and 1.4-1.7 F1 points over the previous works. The strong downstream model performance using CLAP combined with our learnings from intrinsic evaluation underscores the importance of prioritizing accuracy over faithfulness for NER.

5 Analysis

5.1 Qualitative Analysis

Diving deeper, we qualitatively study typical error cases for the translated labels in four languages

⁴Averaged over five model runs

Source Sentence	Source Label	Target Lang	Technique	Translated Label	Explanation
Born in Castelvetro, Trapani and raised in Catania, he moved to Madrid to keep up his busy career.	Castelvetro	hi	Awesome-align	कैस्टेलवेद्रानो ट्रापानी	Extra word
			EasyProject	Castelvetro	No translation
			CLAP	कैस्टेलवेद्रानो	Perfect
Unilaterally leading a coalition featuring tyrannies, effect such change remains a bad idea, Iraq's elections notwithstanding.	Iraq	zh	Awesome-align	伊拉	Incomplete
			EasyProject	尽管伊拉克	Extra word
			CLAP	伊拉克	Perfect

Table 4: Qualitative examples highlighting the error-cases of the baseline models along with explanations for Hindi (hi) and Chinese (zh). We also show how CLAP performs better and fixes the errors.

by different label projection techniques. In 200 examples of our study, we found that 18% of the time, EasyProject predicts nothing due to markers dropped in the translated sentence, and for 19%, EasyProject simply copies the English label failing to translate it to the target language. For Awesome-align, the majority of errors are due to additional words or incomplete label translations, similar to the observation presented in (Chen et al., 2023). This could be because it is hard for the word-alignment module to decide alignments between sub-words, leading to over-alignment or under-alignment. We show two selected examples of our study from Hindi (hi) and Chinese (zh) in Table 4, where we show how Awesome-align predicts extra words or incomplete words owing to misalignments, and EasyProject fails to translate the word for Hindi while producing extra tokens for Chinese. In both cases, we show how CLAP makes accurate predictions and is more robust in maintaining accurate label translations.

5.2 Generalization to other translation models

To verify the generalizability of our approach to other translation models, we perform an extrinsic evaluation of the label projection techniques on the EAE task using the mBART-50 many-to-many (MMT) (Kong et al., 2021) translation model. We show the results for this evaluation in Table 5. We see that CLAP performs the best with an average improvement of 2 F1 points over the next best baseline of Awesome-align and 6.5 F1 points over the zero-shot baseline. This result shows our CLAP is a generalizable label projection technique and agnostic to the underlying translation model.

5.3 Ablation Study for CLAP

To study the impact of using instruction-tuned models for *contextual translation*, we conduct an ablation study comparing CLAP with the follow-

	mBART		mT5		mT5+Copy		Avg
	ar	zh	ar	zh	ar	zh	
Zero-shot	36.3	47.3	36.7	51.0	40.3	51.9	43.9
Awesome-align	45.7	48.6	43.1	52.1	47.1	53.8	48.4
EasyProject	37.3	53.6	35.3	54.0	36.5	55.6	45.4
CLAP (ours)	45.5	52.0	44.8	54.7	48.2	56.9	50.4

Table 5: Extrinsic evaluation of the different label projection techniques using translate-train for EAE using the mBART-50 many-to-many translation model.

	mBART		mT5		mT5+Copy		Avg
	ar	zh	ar	zh	ar	zh	
Zero-shot	36.3	47.3	36.7	51.0	40.3	51.9	43.9
Independent	44.8	49.5	41.3	50.6	44.8	54.3	47.6
Constrained	44.5	51.2	42.3	53.5	45.6	55.6	48.8
CLAP (ours)	45.5	52.0	44.8	54.7	48.2	56.9	50.4
Supervised	60.7	66.4	61.4	68.6	63.2	69.7	65.0

Table 6: Ablation study comparing different contextual translation techniques for label projection. Performance is measured by downstream EAE performance.

ing strong baselines: (1) **Independent** translation uses the translation model \mathcal{T} to independently (without any context of the input sentence) translate the source text labels to the target language (i.e. $\mathbf{y}^{tgt} = \mathcal{T}(\mathbf{y}^{src})$), (2) **Constrained** translation which uses a decoding constraint to carry out the faithfulness requirements. More specifically, during translation, it limits the generation vocabulary to the tokens in the translated sentence x^{tgt} . We follow De Cao et al. (2022); Lu et al. (2022) for implementing these constraints.

We extrinsically evaluate the model performances of the techniques on the task of EAE using the MMT translation model⁵ and show the results

⁵Since decoding-time constraints for the Constrained model can't be applied to GMT

	EAE		NER			Avg
	ar	zh	it	es	id	
Zero-shot	36.3	47.3	79.4	74.5	53.1	58.1
Awesome-align	32.8	30.1	77.5	69.6	51.4	52.3
EasyProject	17.0	11.5	65.9	62.6	51.8	41.8
CLAP (ours)	34.3	39.5	73.4	75.0	57.4	55.9

Table 7: Extrinsic evaluation of the different label projection techniques in terms of downstream model performance using translate-test using GMT for EAE and NER. Avg = Average

in Table 6. We notice how simple independent translations can provide strong gains over the zero-shot model, but contextual translation can provide higher gains. The improvement of 1.6 F1 points of CLAP over the Constrained model highlights the significance of using an instruction-tuned model for contextual translation.

5.4 Using CLAP for Translate-Test

Another popular technique for cross-lingual transfer is translate-test (Hu et al., 2020; Ruder et al., 2021) which was discussed in Section 2.3. As part of this analysis, we study the applicability of CLAP for translate-test using extrinsic evaluation on Arabic (ar) and Chinese (zh) for EAE and Italian (it), Spanish (es), and Indonesian (id) for NER. We show the results in Table 7. Overall, we see how CLAP outperforms both the other methods significantly achieving the best scores for 4 out of the 5 languages. EasyProject performs the worst as it uses the translation model twice causing higher error propagation. We also note how translate-test doesn’t yield improvements over the zero-shot baseline, especially for EAE as it requires using label projection twice (once for trigger and once for arguments), thus leading to error propagation.

6 Related Works

Zero-shot Cross-lingual Structure Extraction

Since the emergence of strong multilingual models (Devlin et al., 2019; Conneau et al., 2020), various works have focused on zero-shot cross-lingual learning (Hu et al., 2020; Ruder et al., 2021) for various structure extraction tasks like named entity recognition (Li et al., 2021; Yang et al., 2022), relation extraction (Ni and Florian, 2019; Subburathinam et al., 2019), slot filling (Krishnan et al., 2021), and semantic parsing (Nicosia et al., 2021; Sherborne and Lapata, 2022). Recent works have

focussed on building datasets (Pouran Ben Veyseh et al., 2022; Parekh et al., 2023) as well as developing novel modeling designs exploring the usage of parse trees (Subburathinam et al., 2019; Ahmad et al., 2021a; Hsu et al., 2023c), data projection (Yarmohammadi et al., 2021), pooling strategies (Agarwal et al., 2023) and generative models (Hsu et al., 2022; Huang et al., 2022) to improve cross-lingual transfer. We utilize the state-of-the-art model X-Gear (Huang et al., 2022) and XLM-R (Conneau et al., 2020) as the downstream models for EAE and NER respectively, and improve them further using CLAP-guided translate-train.

Label Projection Techniques

Several works have attempted to solve label projection for various structure extraction tasks such as semantic role labeling (Aminian et al., 2017; Fei et al., 2020), slot filling (Xu et al., 2020), semantic parsing (Moradshahi et al., 2020; Awasthi et al., 2023), NER (Ni et al., 2017; Stengel-Eskin et al., 2019), and question-answering (Lee et al., 2018; Lewis et al., 2020; Bornea et al., 2021). The earliest works (Yarowsky et al., 2001; Akbik et al., 2015) utilized statistical word-alignment techniques like GIZA++ (Och and Ney, 2003) or fast-align (Dyer et al., 2013) for locating the labels in the translated sentence. Recent works (Chen et al., 2023) have also explored the usage of neural word aligners like QA-align (Nagata et al., 2020) and Awesome-align (Dou and Neubig, 2021). Another set of works has explored the paradigm of mark-then-translate using special markers like quote characters (") (Lewis et al., 2020), XML tags (<a>) (Hu et al., 2020), and square braces ([0]) (Chen et al., 2023) to locate the translated labels. Overall, both these techniques can be error-prone and have poorer translation quality (Akbik et al., 2015), as shown in § 4.4 and 5.1.

7 Conclusion and Future Work

In our work, we propose a novel approach CLAP for label projection, which utilizes contextual machine translation using instruction-tuned language models. Experiments on two structure prediction tasks of EAE and NER demonstrate the effectiveness of CLAP compared to other label projection techniques. Furthermore, intrinsic evaluation provides insights to justify our model improvements. Overall, we lay the foundation for exploring the utilization of contextual translation and future works can use it for various other applications as well.

566 Limitations

567 In our work, we show the effectiveness of our
568 model CLAP on two representative structure pre-
569 diction tasks of EAE and NER. Its effectiveness for
570 other structure prediction tasks remains unknown
571 and can be extended in future works. For CLAP,
572 we utilized the 13B version of the Llama-2 model
573 as the base instruction-tuned language model as a
574 proof-of-concept for the effectiveness of CLAP. Fu-
575 ture works can explore the usage of other stronger
576 LLMs to enhance the model performance. Lastly,
577 we would like to point out that our model doesn't
578 improve over the zero-shot model for several lan-
579 guages, mainly owing to the limited language un-
580 derstanding and poor translation quality. However,
581 the focus of our work has been to show the effec-
582 tiveness of our model with other used label pro-
583 jection techniques. With growing model sizes and
584 enhanced coverage of languages, we posit that our
585 model will eventually be able to provide significant
586 improvements for all languages.

587 Ethical Concerns

588 We use an instruction-tuned language model
589 (specifically LLama-2) as the base model for
590 CLAP. Since these instruction-tuned models are
591 not trained equitably in all languages, the model
592 generation quality may vary drastically for each
593 language. Furthermore, since these models are not
594 trained on filtered safe content data, the model may
595 potentially generate harmful content.

596 References

597 Shantanu Agarwal, Steven Fincke, Chris Jenkins, Scott
598 Miller, and Elizabeth Boschee. 2023. [Impact of sub-
599 word pooling strategy on cross-lingual event detec-
600 tion](#). *CoRR*, abs/2302.11365.

601 Wasi Ahmad, Haoran Li, Kai-Wei Chang, and Yashar
602 Mehdad. 2021a. [Syntax-augmented multilingual
603 BERT for cross-lingual transfer](#). In *Proceedings
604 of the 59th Annual Meeting of the Association for
605 Computational Linguistics and the 11th International
606 Joint Conference on Natural Language Processing
607 (Volume 1: Long Papers)*, pages 4538–4554, Online.
608 Association for Computational Linguistics.

609 Wasi Uddin Ahmad, Nanyun Peng, and Kai-Wei Chang.
610 2021b. [GATE: graph attention transformer encoder
611 for cross-lingual relation and event extraction](#). In
612 *Thirty-Fifth AAAI Conference on Artificial Intelli-
613 gence, AAAI 2021, Thirty-Third Conference on In-
614 novative Applications of Artificial Intelligence, IAAI*

2021, *The Eleventh Symposium on Educational Ad-
vances in Artificial Intelligence, EAAI 2021, Vir-
tual Event, February 2-9, 2021*, pages 12462–12470.
AAAI Press. 615
616
617
618

Wasi Uddin Ahmad, Zhisong Zhang, Xueze Ma, Kai-
Wei Chang, and Nanyun Peng. 2019. [Cross-lingual
dependency parsing with unlabeled auxiliary lan-
guages](#). In *Proceedings of the 23rd Conference on
Computational Natural Language Learning (CoNLL)*,
pages 372–382, Hong Kong, China. Association for
Computational Linguistics. 619
620
621
622
623
624
625

Alan Akbik, Laura Chiticariu, Marina Danilevsky, Yun-
yao Li, Shivakumar Vaithyanathan, and Huaiyu Zhu.
2015. [Generating high quality proposition Banks for
multilingual semantic role labeling](#). 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 397–407, Beijing,
China. Association for Computational Linguistics. 626
627
628
629
630
631
632
633
634

Maryam Aminian, Mohammad Sadegh Rasooli, and
Mona Diab. 2017. [Transferring semantic roles using
translation and syntactic information](#). In *Proceedings
of the Eighth International Joint Conference on Nat-
ural Language Processing (Volume 2: Short Papers)*,
pages 13–19, Taipei, Taiwan. Asian Federation of
Natural Language Processing. 635
636
637
638
639
640
641

Maryam Aminian, Mohammad Sadegh Rasooli, and
Mona Diab. 2019. [Cross-lingual transfer of semantic
roles: From raw text to semantic roles](#). In *Proceed-
ings of the 13th International Conference on Com-
putational Semantics - Long Papers*, pages 200–210,
Gothenburg, Sweden. Association for Computational
Linguistics. 642
643
644
645
646
647
648

Abhijeet Awasthi, Nitish Gupta, Bidisha Samanta,
Shachi Dave, Sunita Sarawagi, and Partha Talukdar.
2023. [Bootstrapping multilingual semantic parsers
using large language models](#). In *Proceedings of the
17th Conference of the European Chapter of the As-
sociation for Computational Linguistics*, pages 2455–
2467, Dubrovnik, Croatia. Association for Computa-
tional Linguistics. 649
650
651
652
653
654
655
656

Mihaela A. Bornea, Lin Pan, Sara Rosenthal, Radu
Florian, and Avirup Sil. 2021. [Multilingual transfer
learning for QA using translation as data augmenta-
tion](#). In *Thirty-Fifth AAAI Conference on Artificial
Intelligence, AAAI 2021, Thirty-Third Conference
on Innovative Applications of Artificial Intelligence,
IAAI 2021, The Eleventh Symposium on Educational
Advances in Artificial Intelligence, EAAI 2021, Vir-
tual Event, February 2-9, 2021*, pages 12583–12591.
AAAI Press. 657
658
659
660
661
662
663
664
665
666

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie
Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
Neelakantan, Pranav Shyam, Girish Sastry, Amanda
Askell, Sandhini Agarwal, Ariel Herbert-Voss,
Gretchen Krueger, Tom Henighan, Rewon Child,
Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, 667
668
669
670
671
672

673	Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners . <i>CoRR</i> , abs/2005.14165.	4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.	731 732
679	Pere-Lluís Huguet Cabot, Simone Tedeschi, Axel-Cyrille Ngonga Ngomo, and Roberto Navigli. 2023. Red^{fm}: a filtered and multilingual relation extraction dataset . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4326–4343, Toronto, Canada. Association for Computational Linguistics.	George Doddington, Alexis Mitchell, Mark Przybocki, Lance Ramshaw, Stephanie Strassel, and Ralph Weischedel. 2004. The automatic content extraction (ACE) program – tasks, data, and evaluation . In <i>Proceedings of the Fourth International Conference on Language Resources and Evaluation (LREC’04)</i> , Lisbon, Portugal. European Language Resources Association (ELRA).	733 734 735 736 737 738 739 740
686	Pengfei Cao, Zhuoran Jin, Yubo Chen, Kang Liu, and Jun Zhao. 2023. Zero-shot cross-lingual event argument extraction with language-oriented prefix-tuning . In <i>Thirty-Seventh AAAI Conference on Artificial Intelligence, AAAI 2023, Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence, IAAI 2023, Thirteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2023, Washington, DC, USA, February 7-14, 2023</i> , pages 12589–12597. AAAI Press.	Zi-Yi Dou and Graham Neubig. 2021. Word alignment by fine-tuning embeddings on parallel corpora . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2112–2128, Online. Association for Computational Linguistics.	741 742 743 744 745 746
696	Yang Chen, Chao Jiang, Alan Ritter, and Wei Xu. 2023. Frustratingly easy label projection for cross-lingual transfer . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5775–5796, Toronto, Canada. Association for Computational Linguistics.	Chris Dyer, Victor Chahuneau, and Noah A. Smith. 2013. A simple, fast, and effective reparameterization of IBM model 2 . In <i>Proceedings of the 2013 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 644–648, Atlanta, Georgia. Association for Computational Linguistics.	747 748 749 750 751 752 753
702	Yang Chen and Alan Ritter. 2021. Model selection for cross-lingual transfer . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021</i> , pages 5675–5687. Association for Computational Linguistics.	Hao Fei, Meishan Zhang, and Donghong Ji. 2020. Cross-lingual semantic role labeling with high-quality translated training corpus . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7014–7026, Online. Association for Computational Linguistics.	754 755 756 757 758 759
709	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 <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 8440–8451, Online. Association for Computational Linguistics.	Steven Fincke, Shantanu Agarwal, Scott Miller, and Elizabeth Boschee. 2022. Language model priming for cross-lingual event extraction . In <i>Thirty-Sixth AAAI Conference on Artificial Intelligence, AAAI 2022, Thirty-Fourth Conference on Innovative Applications of Artificial Intelligence, IAAI 2022, The Twelfth Symposium on Educational Advances in Artificial Intelligence, EAAI 2022 Virtual Event, February 22 - March 1, 2022</i> , pages 10627–10635. AAAI Press.	760 761 762 763 764 765 766 767 768 769
718	Nicola De Cao, Ledell Wu, Kashyap Papat, Mikel Artetxe, Naman Goyal, Mikhail Plekhanov, Luke Zettlemoyer, Nicola Cancedda, Sebastian Riedel, and Fabio Petroni. 2022. Multilingual autoregressive entity linking . <i>Transactions of the Association for Computational Linguistics</i> , 10:274–290.	Leonhard Hennig, Philippe Thomas, and Sebastian Möller. 2023. MultiTACRED: A multilingual version of the TAC relation extraction dataset . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3785–3801, Toronto, Canada. Association for Computational Linguistics.	770 771 772 773 774 775 776
724	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding . In <i>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)</i> , pages	I-Hung Hsu, Kuan-Hao Huang, Elizabeth Boschee, Scott Miller, Prem Natarajan, Kai-Wei Chang, and Nanyun Peng. 2022. DEGREE: A data-efficient generation-based event extraction model . In <i>Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies</i> , pages 1890–1908, Seattle, United States. Association for Computational Linguistics.	777 778 779 780 781 782 783 784 785
730		I-Hung Hsu, Kuan-Hao Huang, Shuning Zhang, Wenxin Cheng, Prem Natarajan, Kai-Wei Chang, and Nanyun	786 787

788	Peng. 2023a. TAGPRIME: A unified framework for relational structure extraction . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 12917–12932, Toronto, Canada. Association for Computational Linguistics.		
789			
790			
791			
792			
793			
794	I-Hung Hsu, Avik Ray, Shubham Garg, Nanyun Peng, and Jing Huang. 2023b. Code-switched text synthesis in unseen language pairs . In <i>Findings of the Association for Computational Linguistics: ACL 2023</i> , pages 5137–5151, Toronto, Canada. Association for Computational Linguistics.		
795			
796			
797			
798			
799			
800	I-Hung Hsu, Zhiyu Xie, Kuan-Hao Huang, Prem Natarajan, and Nanyun Peng. 2023c. AMPERE: AMR-aware prefix for generation-based event argument extraction model . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 10976–10993, Toronto, Canada. Association for Computational Linguistics.		
801			
802			
803			
804			
805			
806			
807			
808	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 generalization . <i>CoRR</i> , abs/2003.11080.		
809			
810			
811			
812			
813	Kuan-Hao Huang, Wasi Ahmad, Nanyun Peng, and Kai-Wei Chang. 2021. Improving zero-shot cross-lingual transfer learning via robust training . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 1684–1697, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.		
814			
815			
816			
817			
818			
819			
820	Kuan-Hao Huang, I-Hung Hsu, Prem Natarajan, Kai-Wei Chang, and Nanyun Peng. 2022. Multilingual generative language models for zero-shot cross-lingual event argument extraction . In <i>Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 4633–4646, Dublin, Ireland. Association for Computational Linguistics.		
821			
822			
823			
824			
825			
826			
827			
828	Chris Jenkins, Shantanu Agarwal, Joel Barry, Steven Fincke, and Elizabeth Boschee. 2023. Massively multi-lingual event understanding: Extraction, visualization, and search . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)</i> , pages 247–256, Toronto, Canada. Association for Computational Linguistics.		
829			
830			
831			
832			
833			
834			
835			
836	Xiang Kong, Adithya Renduchintala, James Cross, Yuqing Tang, Jiatao Gu, and Xian Li. 2021. Multilingual neural machine translation with deep encoder and multiple shallow decoders . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, EACL 2021, Online, April 19 - 23, 2021</i> , pages 1613–1624. Association for Computational Linguistics.		
837			
838			
839			
840			
841			
842			
843			
844			
	Jitin Krishnan, Antonios Anastasopoulos, Hemant Purohit, and Huzefa Rangwala. 2021. Multilingual code-switching for zero-shot cross-lingual intent prediction and slot filling . In <i>Proceedings of the 1st Workshop on Multilingual Representation Learning</i> , pages 211–223, Punta Cana, Dominican Republic. Association for Computational Linguistics.		845 846 847 848 849 850 851
	Kyungjae Lee, Kyoungho Yoon, Sunghyun Park, and Seung-won Hwang. 2018. Semi-supervised training data generation for multilingual question answering . In <i>Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC 2018)</i> , Miyazaki, Japan. European Language Resources Association (ELRA).		852 853 854 855 856 857 858
	Patrick Lewis, Barlas Oguz, Ruty Rinott, Sebastian Riedel, and Holger Schwenk. 2020. MLQA: Evaluating cross-lingual extractive question answering . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 7315–7330, Online. Association for Computational Linguistics.		859 860 861 862 863 864 865
	Bing Li, Yujie He, and Wenjin Xu. 2021. Cross-lingual named entity recognition using parallel corpus: A new approach using xlm-roberta alignment . <i>CoRR</i> , abs/2101.11112.		866 867 868 869
	Keming Lu, I-Hung Hsu, Wenxuan Zhou, Mingyu Derek Ma, and Muhao Chen. 2022. Summarization as indirect supervision for relation extraction . In <i>Findings of the Association for Computational Linguistics: EMNLP 2022</i> , pages 6575–6594, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.		870 871 872 873 874 875 876
	Mehrad Moradshahi, Giovanni Campagna, Sina Semnani, Silei Xu, and Monica Lam. 2020. Localizing open-ontology QA semantic parsers in a day using machine translation . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 5970–5983, Online. Association for Computational Linguistics.		877 878 879 880 881 882 883
	Masaaki Nagata, Katsuki Chousa, and Masaaki Nishino. 2020. A supervised word alignment method based on cross-language span prediction using multilingual BERT . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 555–565, Online. Association for Computational Linguistics.		884 885 886 887 888 889 890
	Jian Ni, Georgiana Dinu, and Radu Florian. 2017. Weakly supervised cross-lingual named entity recognition via effective annotation and representation projection . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1470–1480, Vancouver, Canada. Association for Computational Linguistics.		891 892 893 894 895 896 897
	Jian Ni and Radu Florian. 2019. Neural cross-lingual relation extraction based on bilingual word embedding mapping . In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing</i>		898 899 900 901

902				
903				
904				
905				
906	Massimo Nicosia, Zhongdi Qu, and Yasemin Altun.			
907	2021. Translate & Fill: Improving zero-shot multilingual semantic parsing with synthetic data . In <i>Findings of the Association for Computational Linguistics: EMNLP 2021</i> , pages 3272–3284, Punta Cana, Dominican Republic. Association for Computational Linguistics.			962
908				963
909				964
910				965
911				966
912				967
				968
913	Franz Josef Och and Hermann Ney. 2003. A systematic comparison of various statistical alignment models . <i>Computational Linguistics</i> , 29(1):19–51.			969
914				970
915				
916	Xiaoman Pan, Boliang Zhang, Jonathan May, Joel Nothman, Kevin Knight, and Heng Ji. 2017. Cross-lingual name tagging and linking for 282 languages . In <i>Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 1946–1958, Vancouver, Canada. Association for Computational Linguistics.			971
917				972
918				973
919				974
920				975
921				976
922				977
				978
				979
923	Tanmay Parekh, I-Hung Hsu, Kuan-Hao Huang, Kai-Wei Chang, and Nanyun Peng. 2023. GENEVA: Benchmarking generalizability for event argument extraction with hundreds of event types and argument roles . In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</i> , pages 3664–3686, Toronto, Canada. Association for Computational Linguistics.			980
924				981
925				982
926				983
927				984
928				
929				985
930				986
				987
931	Amir Pouran Ben Veyseh, Javid Ebrahimi, Franck Dernoncourt, and Thien Nguyen. 2022. MEE: A novel multilingual event extraction dataset . In <i>Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing</i> , pages 9603–9613, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.			988
932				989
933				990
934				991
935				
936				992
937				993
				994
				995
				996
938	Afshin Rahimi, Yuan Li, and Trevor Cohn. 2019. Massively multilingual transfer for NER . In <i>Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics</i> , pages 151–164, Florence, Italy. Association for Computational Linguistics.			997
939				998
940				999
941				1000
942				1001
				1002
943	Sebastian Ruder, Noah Constant, Jan Botha, Aditya Siddhant, Orhan Firat, Jinlan Fu, Pengfei Liu, Junjie Hu, Dan Garrette, Graham Neubig, and Melvin Johnson. 2021. XTREME-R: Towards more challenging and nuanced multilingual evaluation . In <i>Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing</i> , pages 10215–10245, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.			1003
944				1004
945				1005
946				1006
947				1007
948				1008
949				
950				1009
951				1010
				1011
952	Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks . <i>arXiv preprint arXiv:1704.04368</i> .			1012
953				1013
954				1014
955				1015
956	Tom Sherborne and Mirella Lapata. 2022. Zero-shot cross-lingual semantic parsing . In <i>Proceedings of the</i>			958
957				959
				960
				961
				962
				963
				964
				965
				966
				967
				968
				969
				970
				971
				972
				973
				974
				975
				976
				977
				978
				979
				980
				981
				982
				983
				984
				985
				986
				987
				988
				989
				990
				991
				992
				993
				994
				995
				996
				997
				998
				999
				1000
				1001
				1002
				1003
				1004
				1005
				1006
				1007
				1008
				1009
				1010
				1011
				1012
				1013
				1014
				1015

- 1016 KayYen Wong, Sameen Maruf, and Gholamreza Haf-
1017 fari. 2020. [Contextual neural machine translation](#)
1018 [improves translation of cataphoric pronouns](#). In *Pro-*
1019 *ceedings of the 58th Annual Meeting of the Asso-*
1020 *ciation for Computational Linguistics*, pages 5971–
1021 5978, Online. Association for Computational Lin-
1022 guistics.
- 1023 Weijia Xu, Batoool Haider, and Saab Mansour. 2020.
1024 [End-to-end slot alignment and recognition for cross-](#)
1025 [lingual NLU](#). In *Proceedings of the 2020 Conference*
1026 *on Empirical Methods in Natural Language Process-*
1027 *ing (EMNLP)*, pages 5052–5063, Online. Association
1028 for Computational Linguistics.
- 1029 Linting Xue, Noah Constant, Adam Roberts, Mihir Kale,
1030 Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and
1031 Colin Raffel. 2021. [mT5: A massively multilingual](#)
1032 [pre-trained text-to-text transformer](#). In *Proceedings*
1033 *of the 2021 Conference of the North American Chap-*
1034 *ter of the Association for Computational Linguistics:*
1035 *Human Language Technologies*, pages 483–498, On-
1036 line. Association for Computational Linguistics.
- 1037 Jian Yang, Shaohan Huang, Shuming Ma, Yuwei Yin,
1038 Li Dong, Dongdong Zhang, Hongcheng Guo, Zhou-
1039 jun Li, and Furu Wei. 2022. [CROP: Zero-shot cross-](#)
1040 [lingual named entity recognition with multilingual](#)
1041 [labeled sequence translation](#). In *Findings of the Asso-*
1042 *ciation for Computational Linguistics: EMNLP 2022*,
1043 pages 486–496, Abu Dhabi, United Arab Emirates.
1044 Association for Computational Linguistics.
- 1045 Mahsa Yarmohammadi, Shijie Wu, Marc Marone,
1046 Haoran Xu, Seth Ebner, Guanghui Qin, Yunmo
1047 Chen, Jialiang Guo, Craig Harman, Kenton Murray,
1048 Aaron Steven White, Mark Dredze, and Benjamin
1049 Van Durme. 2021. [Everything is all it takes: A multi-](#)
1050 [pronged strategy for zero-shot cross-lingual informa-](#)
1051 [tion extraction](#). In *Proceedings of the 2021 Confer-*
1052 *ence on Empirical Methods in Natural Language Pro-*
1053 *cessing*, pages 1950–1967, Online and Punta Cana,
1054 Dominican Republic. Association for Computational
1055 Linguistics.
- 1056 David Yarowsky, Grace Ngai, and Richard Wicentowski.
1057 2001. [Inducing multilingual text analysis tools via](#)
1058 [robust projection across aligned corpora](#). In *Proce-*
1059 *edings of the First International Conference on Human*
1060 *Language Technology Research*.
- 1061 Pengfei Yu, Jonathan May, and Heng Ji. 2023. [Bridg-](#)
1062 [ing the gap between native text and translated text](#)
1063 [through adversarial learning: A case study on cross-](#)
1064 [lingual event extraction](#). In *Findings of the Associ-*
1065 *ation for Computational Linguistics: EACL 2023*,
1066 pages 754–769, Dubrovnik, Croatia. Association for
1067 Computational Linguistics.
- 1068 Wenhao Zhu, Hongyi Liu, Qingxiu Dong, Jingjing Xu,
1069 Lingpeng Kong, Jiajun Chen, Lei Li, and Shujian
1070 Huang. 2023. [Multilingual machine translation with](#)
1071 [large language models: Empirical results and analy-](#)
1072 [sis](#). *CoRR*, abs/2304.04675.

A Data Statistics

We present the extensive data statistics for the ACE and WikiANN datasets used for downstream model evaluation on EAE and NER respectively. For ACE, Table 8 provides details about the number of events and arguments for each language. For WikiANN, we present the statistics in Table 9

Language	Train	Dev	Test	
	English	English	Arabic	Chinese
# Events	4,202	450	198	190
# Arguments	4,859	605	287	336

Table 8: Data Statistics in terms of events and arguments of the ACE dataset for the downstream task of EAE. # indicates ‘number of’.

B Complete Results for Intrinsic Evaluation

B.1 Accuracy Evaluation

Accuracy evaluation is done by 5 native bilingual speakers for Chinese, Arabic, Hindi, and Spanish by ranking the translation quality of the translated labels. The native speakers were undergraduate and graduate students who were well-versed in their respective native languages. We present the interface of the google sheets along with the instructions shown to the annotators for Chinese in Figure 4. Similarly, annotation was performed for the other languages as well. We present the complete results as an A/B comparison of the different techniques in terms of their win rates (i.e. percentage when A is better than B) in Table 10. We note how CLAP is more accurate than previous baselines of Awesome-align and EasyProject while at par with the Independent baseline.

B.2 Faithfulness Evaluation

We present the complete results for the faithfulness evaluation per language in Tables 11 and 12 for EAE and NER tasks respectively. For EAE, CLAP has the best faithfulness followed by Awesome-align. For NER, Awesome-align and EasyProject have the highest faithfulness.

C Additional Implementation Details

C.1 X-Gear

X-Gear is used as the downstream model for EAE for extrinsic evaluation of the label projection

Split	Language	# Sentences	# Entities
Train	English (en)	20,000	27,931
Dev	English (en)	10,000	14,146
	Afrikaans (af)	1,000	1,487
	Arabic (ar)	10,000	11,259
	Bulgarian (bg)	10,000	14,060
	Bengali (bn)	1,000	1,089
	German (de)	10,000	13,868
	Greek (el)	10,000	12,163
	Spanish (es)	10,000	12,260
	Estonian (et)	10,000	13,892
	Basque (eu)	10,000	13,459
	Farsi (fa)	10,000	10,742
	Finnish (fi)	10,000	14,554
	French (fr)	10,000	13,369
	Hebrew (he)	10,000	13,698
	Hindi (hi)	1,000	1,228
	Hungarian (hu)	10,000	14,163
	Indonesian (id)	10,000	11,447
	Italian (it)	10,000	13,749
	Japanese (ja)	10,000	13,446
	Javanese (jv)	100	117
Test	Georgian (ka)	10,000	13,057
	Kazakh (kk)	1,000	1,115
	Korean (ko)	10,000	14,423
	Malayalam (ml)	1,000	1,204
	Marathi (mr)	1,000	1,264
	Malay (ms)	1,000	1,115
	Burmese (my)	100	119
	Dutch (nl)	10,000	13,725
	Portuguese (pt)	10,000	12,823
	Russian (ru)	10,000	12,177
	Swahili (sw)	1,000	1,194
	Tamil (ta)	1,000	1,241
	Telugu (te)	1,000	1,171
	Thai (th)	10,000	16,970
	Tagalog (tl)	1,000	1,034
	Turkish (tr)	10,000	13,587
	Urdu (ur)	1,000	1,020
	Vietnamese (vi)	10,000	11,305
	Yoruba (yo)	100	111
	Chinese (zh)	10,000	12,049

Table 9: Data Statistics in terms of sentences and entities of the WikiANN dataset for the downstream task of NER. # indicates ‘number of’.

techniques. The original X-Gear work (Huang et al., 2022) explored two base multilingual models: mBART-50-large (mBART) (Kong et al., 2021) and the mT5-base (mT5) (Xue et al., 2021). They also explored the usage of copy mechanism (See et al., 2017) to prompt the models to predict strings from the input sentence. In our work, we utilized mBART without copy (mBART), mT5 without copy (mT5), and mT5 with copy mechanism (mT5+Copy) as the downstream models. We present details about the hyperparameter settings for these models in Table 13. We run experiments for CLAP on a NVIDIA GeForce RTX 2080 Ti machine with support for 8 GPUs.

System 1	v/s System 2	Arabic			Chinese			Hindi			Spanish		
		S1	Tie	S2	S1	Tie	S2	S1	Tie	S2	S1	Tie	S2
CLAP	Awesome-align	36%	58%	6%	45%	50%	5%	20%	74%	6%	12%	84%	4%
CLAP	EasyProject	52%	32%	16%	56%	39%	5%	42%	48%	10%	30%	66%	4%
CLAP	Independent	18%	60%	22%	12%	71%	17%	18%	64%	18%	24%	68%	8%
Independent	Awesome-align	44%	42%	14%	39%	57%	4%	28%	60%	12%	20%	64%	16%
Independent	EasyProject	50%	44%	6%	50%	46%	4%	52%	36%	12%	32%	52%	16%
Awesome-align	EasyProject	42%	26%	32%	34%	50%	16%	42%	42%	16%	26%	64%	10%

Table 10: A/B comparison of the various label projection techniques for accuracy evaluation for the Google Translation model. Accuracy is measured as the label translation quality by native human speakers. Here, **S1** = System 1 is better, **S2** = System 2 is better, and **Tie** = similar quality. The better systems are highlighted in **bold**.

Guidelines:
Looking at the **English word** in context of the **English sentence**, evaluate the word translations by System 1, 2 and 3 by giving them rankings - i.e. 1 / 2 / 3 / 4 (1 = best and 4 = worst)

SPECIAL NOTES
1. If two systems deserve the same rank, mark them with the same rank (e.g. 1 / 1 / 3 / 4 OR 1 / 2 / 2 / 2)
2. If a system translation has "-", that means the system was not able to translate the phrase at all. This is the worst kind of translation and should be ranked the worst.
3. If the word is not translated and in English itself, it would be considered a poorer translation than phonetic translation of the word in the target language. But the English translation should be considered better than random gibberish in the target language

English Sentence	English word	Translations				Rankings			
		System 1	System 2	System 3	System 4	System 1	System 2	System 3	System 4
happily watching tom and jerry on his mini television , his transformation from the pain - racked boy who left baghdad .	baghdad	巴格达	巴格达	巴格达	巴格达				
reporter : the kramers must wait and travel to another town for abby . on the next flight , passengers wear masks and their temperatures are taken for signs of sars	kramers	kramers	克莱默斯	克莱默斯	克莱默斯				
Allegations have come to light that several OSU players received illegal benefits including cash , access to cars , etc .	players	玩家	球员	球员	球员				
The first one was on Saturday and triggered intense gun battles , which according to some U.S. accounts , left at least 2,000 Iraq fighters dead .	gun	星期六的	枪	枪	枪战				
Now that armored columns of U.S. - led troops have reached the outskirts of Baghdad , eyewitnesses report fighting and shelling around Saddam Hussein International Airport .	Saddam Hussein International Airport	萨达姆·侯赛因国际机场	萨达姆·侯赛因国际机场	萨达姆·侯赛因国际机场	萨达姆·侯赛因国际机场				
we have eyewitnesses to his orders of execution of hundreds of people in 1991 during the shiite muslim uprising	people	-	人们	人	数百人				
I'm reminded of when I lived in another state and the local cop charged the town drunk in his driveway after following him home from the pub .	drunk	-	醉	醉	喝醉				

Figure 4: Annotation Interface for conducting the intrinsic evaluation for Accuracy. The shown examples are for Chinese, while the study was done for Hindi, Spanish, and Arabic as well.

Techniques	ar	zh	Avg.
Independent	33	38	35
Awesome-align	66	83	74
EasyProject	31	66	48
CLAP	74	85	79

on a NVIDIA GeForce RTX 2080 Ti machine with support for 8 GPUs.

1137
1138

Table 11: Faithfulness evaluation of the various label projection techniques for EAE as a percentage of the times the translated labels were present in the translated input sentence. Numbers are in percentage (%). Higher faithfulness is better and the best techniques are highlighted in **bold**.

C.2 XLM-R

XLM-R (Conneau et al., 2020) is used as the downstream model for NER for extrinsic evaluation of the label projection techniques. We mainly follow the XTREME (Hu et al., 2020) framework for setting up the task and model. We present details about the hyperparameter settings for this model in Table 14. We run experiments for CLAP on a NVIDIA GeForce RTX 2080 Ti machine with support for 8 GPUs.

C.3 CLAP

We report the hyperparameter settings for our model in Table 15. We run experiments for CLAP

Techniques	af	ar	bg	bn	de	el	es
Independent	78	66	67	74	79	57	70
Awesome-align	99	95	98	92	99	98	99
EasyProject	100	98	83	98	97	89	99
CLAP	94	75	63	93	79	46	84
	et	eu	fa	fi	fr	he	hi
Independent	70	64	61	71	71	71	65
Awesome-align	98	97	96	99	98	95	93
EasyProject	97	94	99	98	99	94	36
CLAP	92	91	72	92	74	80	90
	hu	id	it	ja	ju	ka	kk
Independent	68	77	74	68	66	64	56
Awesome-align	98	99	99	58	98	95	94
EasyProject	97	99	98	95	94	99	77
CLAP	93	84	78	67	53	70	85
	ko	ml	mr	ms	my	nl	pt
Independent	63	57	73	80	53	76	76
Awesome-align	96	88	92	99	90	99	97
EasyProject	93	87	73	98	62	100	99
CLAP	64	88	95	82	55	85	89
	ru	sw	ta	te	th	tl	tr
Independent	59	79	72	76	66	81	76
Awesome-align	97	96	91	91	51	99	98
EasyProject	99	97	91	87	99	99	98
CLAP	66	94	96	90	57	58	94
	vi	ur	yo	zh	Avg.		
Independent	74	74	45	66	69		
Awesome-align	83	97	92	92	93		
EasyProject	98	94	77	92	92		
CLAP	89	91	88	60	79		

Table 12: Faithfulness evaluation of the various label projection techniques for NER as a percentage of the times the translated labels were present in the translated input sentence. Numbers are in percentage (%). Higher faithfulness is better and the best techniques are highlighted in **bold**.

	mBART	mT5	mT5+Copy
Base Model	multilingual BART-Large	multilingual T5-Large	multilingual T5-Large
Usage of copy	No	No	Yes
Training Batch Size	16	16	16
Eval Batch Size	32	32	32
Learning Rate	2×10^{-5}	1×10^{-4}	2×10^{-5}
Weight Decay	1×10^{-5}	1×10^{-5}	1×10^{-5}
# Warmup Epochs	5	5	5
Gradient Clipping	5	5	5
Max Training Epochs	60	60	60
# Accumulation Steps	1	1	1
Beam Size	4	4	4
Max Sequence Length	350	350	350
Max Output Length	100	100	100

Table 13: Hyperparameter details for the EAE downstream X-Gear model.

Base Model	XLM - Roberta - Large
# Training Epochs	5
Training Batch Size	32
Evaluation Batch Size	32
Learning Rate	2×10^{-5}
Weight Decay	0
Max Sequence Length	128
# Accumulation Steps	1
# Saving Steps	1000

Table 14: Hyperparameter details for the NER downstream XLM-R model.

Base Model	LLAMA-2-13B
Temperature	0.6
Top-p	0.9
Maximum Generation Length	64-132
# In-context examples	2

Table 15: Hyperparameter details for the CLAP model.