Cross-Lingual Event Detection via Optimized Adversarial Training

Luis F. Guzman-Nateras, Minh Van Nguyen and Thien Huu Nguyen

Department of Computer and Information Science University of Oregon

{lfguzman, minhnv, thien}@cs.uoregon.edu

Abstract

In this work, we focus on Cross-Lingual Event Detection where a model is trained on data from a source language but its performance is evaluated on data from a second, target, language. Most recent works in this area have harnessed the language-invariant qualities displayed by pre-trained Multi-lingual Language Models. Their performance, however, reveals there is room for improvement as the crosslingual setting entails particular challenges. We employ Adversarial Language Adaptation to train a Language Discriminator to discern between the source and target languages using unlabeled data. The discriminator is trained in an adversarial manner so that the encoder learns to produce refined, language-invariant representations that lead to improved performance. More importantly, we optimize the adversarial training process by only presenting the discriminator with the most informative samples. We base our intuition about what makes a sample informative on two disparate metrics: sample similarity and event presence. Thus, we propose leveraging Optimal Transport as a solution to naturally combine these two distinct information sources into the selection process. Extensive experiments on 8 different language pairs, using 4 languages from unrelated families, show the flexibility and effectiveness of our model that achieves state-of-the-art results.

1 Introduction

Event Detection (ED) is an important sub-task within the broader Information Extraction (IE) task. Event detection consists of being able to identify the words, commonly referred to as *triggers*, that denote the occurrence of events in a sentence, and classify them into a discrete set of event types. For example, in the sentence "*Jamie bought a car yesterday.*", *bought* is considered the trigger of a TRANSACTION:TRANSFER-OWNERSHIP¹

event type. It is a very well studied task in which there have been lots of previous research efforts that have recently been primarily deep learning-based (Nguyen and Grishman, 2015; Chen et al., 2015; Nguyen et al., 2016a,b; Sha et al., 2018; Wadden et al., 2019; Zhang et al., 2019; Nguyen and Nguyen, 2019; Zhang et al., 2020; Liu et al., 2020).

Nonetheless, ED remains quite a challenging task as the context in which a trigger occurs can change its corresponding type completely. Furthermore, the same event might also be expressed by entirely different words/phrases. Additionally, the vast majority of the aforementioned efforts are limited to a monolingual setting — performing ED on text belonging to a single language.

Alternatively, Cross-Lingual ED (CLED) proposes the scenario of creating models that effectively perform ED on data belonging to more than one language, which brings about additional challenges. For instance, trigger words present in one language might not exist in another one. An frequent example of this phenomenon are verb conjugations where some tenses only exist in some languages. Accurate verb handling is of particular importance for the ED task as event triggers are usually related to the verbs in a sentence. Some recent work (Majewska et al., 2021) has attempted to address this issue by injecting external verb knowledge into the training process. Another similar problematic issue for CLED are triggers with different meanings that are each distinct words in different languages. For instance, the word "juicio" in Spanish can either mean "judgement" or "trial" in English, depending on the context.

A compelling approach to creating a crosslingual model is to use *transfer learning* which carries the performance of a model trained on a *source* language over onto a second *target* language. The general idea is leveraging the existing high-quality annotated data available for a high-

¹Event type taken from the ACE05 dataset.

resource language to train a model in a way that allows it to learn the language-invariant characteristics of the task at hand, ED in this case, so that it also performs effectively on text from a second language. Prior works on transfer learning for CLED have relied on pre-trained Multilingual Language Models (MLMs), such as multilingual BERT (mBERT) (Devlin et al., 2019), to take advantage of their innate language-invariant qualities. Yet, their performance still shows room for improvement as they sometimes struggle to handle the difficult instances, unique to cross-lingual settings, mentioned earlier. We identify a significant shortcoming of previous CLED efforts in that they do not exploit the abundant supply of unlabeled data: even though MLMs are trained on immense amounts of it, unlabeled data is not used when fine-tuning for the ED task. It is our intuition that by integrating unlabeled target-language data into the training process, the model is exposed to more language context which should help deal with issues such as verb variation and multiple connotations.

As such, we propose making use of Adversarial Language Adaptation (ALA) (Joty et al., 2017; Chen et al., 2018) to train a CLED model. The key idea is to generate language-invariant representations that are not indicative of language but remain informative for the ED task. Unlabeled data from both the source and target languages is used to train a Language Discriminator (LD) network that learns to discern between the two. The adversarial part comes from the fact that the encoder and discriminator are trained with opposing objectives: as the LD becomes better at distinguishing between languages, the encoder learns to generate more language-invariant representations in an attempt to *fool* the LD. To the best of our knowledge, our work is the first one proposing the use of ALA for the CLED task.

Nonetheless, contrary to past uses of ALA where the same importance is given to all unlabeled samples, we recognize that such course of action is suboptimal as certain samples are bound to be more informative for the discriminator than others. For example, we would like to present the LD with the samples that allow it to learn the fine-grained distinctions between the source and target languages, instead of relying on syntactic differences. Moreover, in the context of ED, we suggest it would be beneficial for the LD to be trained with examples containing events, instead of non-event samples, as

the presence of an event can then be incorporated into the generated representations.

Hence, we propose refining the adversarial training process by only keeping the most informative examples while disregarding less useful ones. Our intuition as to what makes samples more informative for CLED is two-fold: First, we presume that presenting the LD with examples that are too different makes the discrimination task too simple. As mentioned previously, we would like the LD to learn a fine-grained distinction between the source and target languages which, in turn, improves the language-invariance of the encoder's representations. Thus, we suggest presenting the LD with examples that have similar contextual semantics, i.e., similar contextualized representations. Second, we consider that sentences containing events should provide an ED system with additional task-relevant information when compared against non-event samples. Accordingly, we argue that event-containing sentences should have a larger probability of being selected for ALA training.

With these intuitions in mind, we propose Optimal Transport (OT) (Villani, 2008) as a natural solution to simultaneously incorporate both the similarity between sample representations and the likelihood of the samples containing an event into a single framework. Therefore, we cast sample selection as an OT problem in which we attempt to find the best alignment between the samples from the source and target languages.

For our experiments, we focus on the widely used ACE05 (Walker et al., 2006) and ERE (Song et al., 2015) datasets which, in conjuction, contain event-annotations in 4 different languages: English, Spanish, Chinese, and Arabic. We work on 8 different language pairs by selecting different languages as the source and target. Our proposed model obtains new state-of-the-art results with considerable performance improvements (+ 2-3% in F1 scores) over competitive baselines and previously published results (M'hamdi et al., 2019). We believe these results demonstrate our model's efficacy and applicability at creating CLED systems.

The rest of this paper is organized as follows: section 2 provides an thorough description of our proposed model, section 3 presents and analyses the results from our experiments, section 4 provides a brief review of related work, and section 5 includes our conclusions.

2 Model

2.1 Problem Definition

Following prior works (M'hamdi et al., 2019; Majewska et al., 2021), we treat ED as a sequence labeling problem. Given a set \mathcal{D} of word sequences $w_i = \{w_{i1}, w_{i2}, ..., w_{in-1}, w_{in}\}$ and their corresponding label sequences $y_i = \{y_{i1}, y_{i2}, ..., y_{in-1}, y_{in}\}$, we use an encoder network E to obtain a contextualized vector representation of the words in the input sequence $\mathbf{h_i} = E(w_i) = \{h_{i1}, h_{i2}, ..., h_{in-1}, h_{in}\}$. Using such representations as input, a prediction network P computes a distribution over the set of possible labels and is trained in a supervised manner using the negative log-likelihood function \mathcal{L}_P :

$$\mathcal{L}_{P} = -\sum_{i=1}^{|D|} \sum_{j=1}^{n} \log P(y_{ij}|h_{ij})$$
 (1)

In the cross-lingual transfer-learning setting, the data used to train the model and the data on which the model is tested come from different languages known as the *source* and *target*, respectively. As such, we deal with two datasets \mathcal{D}_{src} and \mathcal{D}_{tgt} . We assume that we do not have access to the gold labels of the target language y_{tgt} , other than to evaluate our CLED model at testing time.

Our goal is to define a model able to generate language-invariant word representations that are refined enough so that cross-lingual issues, such as the ones described in section 1, are properly handled.

2.2 Baseline Model

Here, we briefly describe the BERT-CRF model proposed by M'hamdi et al. (2019) which was the previous state-of-the-art and serves as our main baseline. Using multilingual BERT (mBERT, (Devlin et al., 2019)) as its encoder, BERT-CRF generates robust, contextualized representations for words from different languages. For words that are split into multiple word-pieces, the average of the representation vectors for all comprising sub-pieces is used as the representation of the full word.

For classification purposes, instead of assigning the labels of each token independently, BERT-CRF uses a Conditional Random Field (CRF) (Lafferty et al., 2001) layer on top of the prediction network to better capture the interactions between the label

sequences. In summary, the contextualized representation vectors \boldsymbol{h}_i generated by the mBERT encoder from the words in the sequence are then fed to a CRF layer which finds the optimal label sequence.

2.3 Adversarial Language Adaptation

The pre-trained versions of MLMs like mBERT or XLM-RoBERTa (Conneau et al., 2019) generate contextualized representations with a certain degree of language-invariance. This can be confirmed by their successful application in cross-lingual settings (M'hamdi et al., 2019; Majewska et al., 2021). However, a lingering issue is the difficulty of learning the nuances of the target language such as verb variations that do not exist in the source language used to train them. Majewska et al. (2021), for instance, propose to address this issue by injecting external verb knowledge into the encoder via task-specific adapter modules (Pfeiffer et al., 2020).

It is our intuition, however, that these issues can be mitigated by achieving a more refined level of language-invariance in the word representations. As such, we propose using Adversarial Language Adaptation (ALA) (Joty et al., 2017), a technique used to create language-invariant models. The ALA framework consists in including a *Language Discriminator* (LD) whose purpose is to learn language-dependent features and be able to differentiate between the samples from either the source or the target languages.

A fundamental characteristic of the ALA approach is its lack of requirements for annotated data in the target language. As such, we can use data from both \mathcal{D}_{src} and \mathcal{D}_{tgt} . An auxiliary dataset $D_{aux} = \{(w_1, l_1), \ldots, (w_{2m}, l_{2m})\}$ is created where w_i is a text sequence from either \mathcal{D}_{src} or \mathcal{D}_{tgt} , and l_i is a language label. The cardinality of D_{aux} is $|D_{aux}| = 2m$, where m is equal to the batch size. Text samples $w_1 \ldots w_m \in \mathcal{D}_{src}$, and samples $w_{m+1} \ldots w_{2m} \in \mathcal{D}_{tgt}$. As described earlier, the encoder E receives the text sequences and produces a sequence of contextualized representations $E(w_i) = h_i = \{h_{i0}, h_{i1}, h_{i2}, \ldots, h_{in}\}$ where h_{i0} is the representation of the [CLS] token added at the beginning of every input sequence.

In our work, the LD is a a simple Multi-Layer Perceptron(MLP) network that takes h_{i0} as input and produces a single sigmoid output. It's trained with the usual *binary cross-entropy* loss function objective: $LD_{loss} = \arg\min_{LD} \mathcal{L}(LD(h_{i0}), l_i)$.

As the LD learns to distinguish between the source and target languages, we concurrently train the encoder to "fool" the discriminator. In other words, the encoder must learn to generate representations that are language-invariant enough that the LD is unable to classify them while still remaining predictive for event-trigger classification. We optimize the following loss:

$$\arg\min_{E,C} \sum_{j=1}^{n} (\mathcal{L}(C(h_{ij}), y_{ij})) - \lambda \mathcal{L}(LD(h_{i0}, l_i))$$

Where C refers to the CRF-based classifier network and λ is a hyperparameter.

Equation 2 is implemented by using a Gradient-Reversal Layer (GRL) (Ganin and Lempitsky, 2015) which acts as the identity during the forward pass, but reverses the direction of the gradients during the backward pass. The first term in Equation 2 can, of course, only be applied for annotated data from the source language.

The GRL is applied to the input vectors, h_{i0} , of the LD. This way, the LD is being trained to differentiate between the two languages while the encoder is trained in the opposite direction, i.e. to generate sequence representations that are harder to discriminate.

2.4 Adversarial Training Optimization

ALA has already been shown to be effective at generating language-invariant models (Joty et al., 2017; Chen et al., 2018). However, in regular ALA training, all samples in a batch, from both the source and target domains, are treated equally. That is, all samples are used as examples for the discriminator to learn how to better discern between the two domains. We propose that ALA effectiveness can be further improved by carefully selecting the samples with which to train the discriminator. We argue that some samples might be more informative than others and that, by only using such informative samples during training, better adaptation results can be achieved.

We base our notion as to *what* makes a sample more informative on two factors. First, we argue that presenting the LD with examples from the source and target language that are too dissimilar makes its task easier which, in turn, leads to the LD not learning the fine-grained distinctions between the languages. Instead, we propose using samples whose vector representations h_{i0} are

close to each other in the embedding space. The intuition for this being that, as representations capture the contextual semantics of the samples, closer representations correspond to more similar examples. Second, we suggest that presenting the LD with samples containing events should make the encoder incorporate task-specific information into its representations.

2.4.1 Optimal Transport

One challenge of using the two mentioned criteria for the ALA sample selection process is that they come with two different measures which are hard to combine. To address this, we propose using Optimal Transport (OT) (Villani, 2008) as a natural way to combine these two metrics into a single framework for sample selection. Optimal transport is, in broad terms, the problem of finding out the cheapest transformation between two discrete probability distributions. It requires a cost function to determine the cost of transforming a data point in one distribution into a data point in the second distribution. When the cost function is based on a valid distance function, the minimum cost is known as the Wasserstein distance. Formally, it solves the following optimization problem:

$$\pi^*(s,t) = \min_{\pi \in \prod(s,t)} \sum_{s \in \mathcal{S}} \sum_{t \in \mathcal{T}} \pi(s,t) \ C(s,t) \ ds \ dt$$

$$\mathbf{s.t.} \ s \sim p(s) \ \text{and} \ t \sim q(t)$$
(3)

where \mathcal{S} and \mathcal{T} are the two domains to be transformed; p(s) and q(t) are the probability distributions of \mathcal{S} and \mathcal{T} , respectively; C is a cost function for mapping \mathcal{S} to \mathcal{T} , $C(s,t): \mathcal{S} \times \mathcal{T} \longrightarrow \mathbb{R}_+$; and finally, $\pi^*(s,t)$ is the optimal joint distribution over the set of all joint distributions $\prod(s,t)$. The problem described by Equation 3 is, of course, intractable. Therefore, we use instead the Sinkhorn algorithm (Cuturi, 2013) which is an entropy-based relaxation of the discrete OT problem.

2.4.2 Problem Formulation

We formulate the OT problem as follows: the domains \mathcal{S} and \mathcal{T} are defined as the representation vectors of the text samples in either the source h_{i0}^s or the target h_{j0}^t languages. We use the L2 distance between these representations as the cost function:

$$C(h_{i0}^s, h_{j0}^t) = ||h_{i0}^s - h_{j0}^t||_2^2$$
 (4)

To define the marginal probability distributions p(s) and q(t) for the $\mathcal S$ and $\mathcal T$ domains, we propose including an Event-Presence (EP) prediction module and use its normalized likelihood scores as the probability distributions for $\mathcal S$ and $\mathcal T$. Thus, the auxiliary dataset D_{aux} is augmented to include an event-presence label e_i for each sample. Of course, this can only be done for samples in the source language as the labels for the target-language data are unavailable:

$$D_{aux} = \{(w_1, l_1, e_1), \dots, (w_m, l_m, e_m), (w_{m+1}, l_{m+1}), \dots, (w_{2m}, l_{2m})\}$$

The EP module is then trained to optimize the following loss:

$$EP_{loss} = \arg\min_{EP} \mathcal{L}(EP(h_{i0}), e_i)$$
 (5)

where $i \le m$, i.e., only using samples from the source language.

The probability distributions p(s) and p(t) are the computed as follows:

$$p(s) = Softmax(EP(h_{i0}^s) \mid l_i == s)$$
 (6)

$$p(t) = Softmax(EP(h_{i0}^t) \mid l_i == t)$$
 (7)

2.4.3 Sample Selection

We use the OT solution matrix π^* , where an entry $\pi^*(s,t)$ represents the optimal cost of transforming data point $s \in \mathcal{S}$ into $t \in \mathcal{T}$, to compute an the overall similarity score v_i of a sample $h_{i0} \in \mathcal{S}$ to the samples in the target domain \mathcal{T} by using the average distance:

$$v_i = \frac{\sum_{j}^{m} \pi^*(h_{i0}^s, h_{j0}^t)}{m}$$
 (8)

Correspondingly, we compute an overall similarity score v_j of each sample $h_{j0} \in \mathcal{T}$ to the samples in the source domain \mathcal{S} :

$$v_j = \frac{\sum_{i=1}^{m} \pi^*(h_{i0}^s, h_{j0}^t)}{m}$$
 (9)

Lastly, we select a fraction, hyperparameter γ , of samples with the best similarity scores from both the source and target languages, and only use these selected samples during ALA training.

2.5 OACLED Model

We train our Optimized Adversarial Cross-Lingual Event Detection (OACLED) model end-to-end with the following loss objective:

$$L_{full} = CRF_{loss} + \alpha LD_{loss} + \beta EP_{loss}$$
 (10) where α and β are trade-off hyperparameters.

3 Experiments

3.1 Datasets

We evaluate our model on the ACE05 (Walker et al., 2006) dataset which includes annotated event-trigger data in 3 languages: English, Chinese and Arabic. To include an additional language in our experiments, we also evaluate on the ERE dataset which has annotated data in English and Spanish. Note that the ACE05 and ERE datasets do not share the same label set: ACE05 involves 33 distinct event types while ERE involves 38 event types. We follow the same data pre-processing and splits as in previous work (M'hamdi et al., 2019) to ensure a fair comparison. Table 1 presents the data statistics.

Dataset	Language	Split	Sentences	Events
	English	Train	19,240	4,419
		Dev	902	468
		Test	676	424
	Chinese	Train	6,841	2,926
ACE05		Dev	526	217
		Test	547	190
	Arabic	Train	2,555	1,793
		Dev	301	230
		Test	262	247
ERE		Train	14,219	6,419
	English	Dev	1,162	552
		Test	1,129	559
	Spanish	Train	7,067	3,272
		Dev	556	210
		Test	546	269

Table 1: Dataset statistics.

3.2 Hyper-parameters

We fine-tune the hyper-parameters for our OA-CLED model using the development data. We apply the following values based on the fine-tuning process:

- · AdamW as the optimizer.
- 5 warm up epochs.
- A learning rate of $1e^{-5}$ for the transformer parameters and of $1e^{-4}$ for the rest of the parameters.
- A batch size of 16.
- 300 for the dimensionality of the layers in feed-forwards networks.
- A $\gamma=0.5$ for the percentage of samples used in adversarial training.

- A $\lambda = 0.001$ as the scaling factor of the GRL layer.
- An $\alpha=1$ and $\beta=0.001$ as the trade-off parameters of the LD loss and ED loss, respectively.
- A dropout of 10% for added regularization during training.

3.3 Main Results

In our experiments, we work with 8 distinct language pairs by selecting each of the available languages as either the source or target language: English-Chinese, Chinese-English, English-Arabic, Arabic-English, Chinese-Arabic, Arabic-Chinese, English-Spanish, and Spanish-English. The Chinese-Spanish, Spanish-Chinese, Arabic-Spanish, and Spanish-Arabic language combinations are unavailable due the previously mentioned incompatibility between the event type sets in ACE05 and ERE.

We compare our OACLED model against 3 relevant baselines. First, the previous state-of-the-art CLED model BERT-CRF (M'hamdi et al., 2019) as described in section 2.2. Second, the mBERT-2TA model (Majewska et al., 2021) which aims at improving cross-lingual performance by incorporating language-independent verb knowledge via task-specific adapters. And third, XLM-R-CRF which is equivalent in all regards to BERT-CRF except that it uses XLM-RoBERTa (Conneau et al., 2019) as the encoder.

Table 2 and Table 3 show the results of our experiments on the ACE05 and ERE datasets, respectively. In all our experiments, we use the base transformer versions *bert-base-cased* and *xlm-roberta-base* as the encoders, parameters are tuned on the development data of the source language, and all entries are the average of five runs.

From Tables 2 and 3, it should be noted that there is a substantial performance increase by performing the trivial change of replacing mBERT with XLM-RoBERTa as the encoder. Furthermore, our OACLED model clearly and consistently outperforms the baselines for all language pairings, with the exception of the *Chinese-Arabic* pair. We attribute this to the impaired performance of XLM-RoBERTa as the encoder for that specific pair as can be confirmed by the poor performance of the XLM-R-CRF baseline on the same configuration. Most importantly, OACLED's improvement over

		Target		
Source	Model	English	Chinese	Arabic
	BERT-2TA	X	46.9*	29.3*
English	BERT-CRF	X	68.5*	30.9*
English	XLM-R-CRF	X	70.49 ± 0.85	43.54±2.77
	OACLED	X	74.64 ±0.73	44.86 ±3.1
	BERT-CRF	37.52±1.73	X	35.05 ±2.85
Chinese	XLM-R-CRF	41.72±1.4	X	32.76 ± 2.31
	OACLED	45.77 ±1.45	X	34.48 ± 2.43
Arabic	BERT-CRF	40.1 ± 3.26	58.78±2.33	X
	XLM-R-CRF	45.22±1.82	61.76±1.57	X
	OACLED	47.98 ±2.07	63.13 ±1.7	X

Table 2: Results on the ACE05 dataset with standard deviation across random seeds. Entries marked * are taken directly from the original papers.

		Target	
Source	Model	English	Spanish
English	BERT-CRF	X	43.28 ± 2.01
	XLM-R-CRF	X	46.79 ± 1.34
	OACLED	X	47.69 ±1.63
Spanish	BERT-CRF	39.8±2.27	X
	XLM-R-CRF	45.61±1.76	X
	OACLED	47.5 ±1.89	X

Table 3: Results on ERE dataset with standard deviation across random seeds.

the XLM-R-CRF baseline is present in every configuration, which validates the effectiveness of our optimized approach to ALA training.

3.4 Ablation Study

We identify 2 main components in our approach: using ALA to create refined language-invariant representations, and optimizing the adversarial training process by selecting a subset of samples chosen with OT to incorporate our measures of informativeness into the sample selection process. Of course, removing ALA training entirely restores the model to the baseline. However, adversarial training optimization via OT has various aspects to it. In order to understand the contribution of these aspects, we explore four different models: OACLED-OT presents the effects of removing sample selection entirely and using all available samples to train the LD; OACLED-L2 uses a constant distance between the unlabeled samples instead the standard L2 distance used in the Sinkhorn algorithm; OACLED-EP completely removes the EP module and a uniform distribution is used as the probability distributions for both languages; finally, OACLED-ED-Loss keeps the EP module, but removes its EP_{loss} term from Equation 10. The performance results of these models is presented in Table 4. In this and the following sections (3.5,

3.6.2), we present the results of experiments using English as the sole source language as it is the source language most ubiquitously used. We, however, found consistency in the displayed effects for different source/target language configurations.

Model version	Target Language		
English	Chinese	Arabic	Spanish
OACLED-OT	70.94	40.55	44.96
OACLED-L2	71.35	41.79	44.39
OACLED-EP	73.08	42.81	46.99
OACLED-EP-Loss	72.93	43.4	46.35
OACLED	74.64	44.86	47.69

Table 4: Ablation experiment results

As expected, removing the sample selection through OT leads to the worst performance drop. This highlights the importance of selecting informative examples for the LD. Furthermore, removing the cost function also hurts performance greatly, which shows that a proper distance function is needed for the OT algorithm to work effectively. While the effects of removing the EP module and its corresponding loss term are not of the same magnitude, they are still significant. These results support our claim for the need and utility of all the components in our approach, showing that their inclusion is crucial in achieving state-of-the-art performance.

3.5 Language Model Finetuning

The key contribution of our approach is to exploit unlabeled data in the target language, which is usually abundant, by introducing it into the training process to improve our model's language-invariant qualities.

To confirm the utility of our approach, Table 5 contrasts our model's performance against a baseline whose encoder has been finetuned with the same unlabeled data using the standard masked language model objective.

Model Version	Target Language		
English	Chinese	Arabic	Spanish
Finetuned XLM-R	71.06	43.71	47.82
OACLED	74.64	44.86	47.69

Table 5: OACLED performance versus a baseline using an encoder finetuned with unlabeled data.

It can be observed that our model outperforms the finetuned baseline in two out of the three target languages. Additionally, the difference in performance in those two instances is considerably larger (3.58% and 1.15%), than the setting in which the baseline performs better (0.13%).

3.6 Analysis

3.6.1 Learned Representation Distances

First, we look at the distance between the sentence-level representations h_{i0} generated by the encoder for different source/target language pairs. Figure 1 shows a plot of such distances using cosine distance as the distance function.

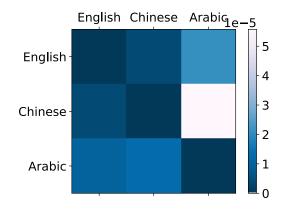


Figure 1: Distance between sentence representations for different language pairs.

When computing the correlation with the performance results in Table 2, we obtain a score R=-0.6616, meaning there is moderate negative correlation between the distance of the representations and model performance, i.e. closer representations lead to better performance.

Similarly, Table 6 shows a comparison of the distances between the representations generated by OACLED and those obtained by the XLM-R-CRF baseline.

	Cosine Distance		
Source/Target	Baseline	OACLED	
English/Chinese	3.64e-3	3.93e-6	
English/Arabic	7.71e-2	2.08e-5	
English/Spanish	5.4e-3	5.3e-6	
Chinese/English	3.62e-3	3.87e-6	
Arabic/English	4.16e-2	1.02e-5	
Spanish/English	6.87e-3	1.49e-5	

Table 6: Comparison of representation-vector distances for language pairs between our model and the baseline.

We observe that OACLED representations are closer, by several orders of magnitude, than those obtained by the baseline. This supports our claim that our model's encoder generates more refined

language-invariant representations than those obtained by the default version of XLM-RoBERTa.

3.6.2 Access to Labeled Target Data

Previously, we discussed how a key feature of our approach is that it does not require annotated data in the target language and, instead, leverages the use of unlabeled data which is readily available. Nonetheless, we also explore the performance of our model in the event that there exists a small amount of annotated target data available. Figure 2 shows the results of our experiments when using different amounts of labeled target data during training.

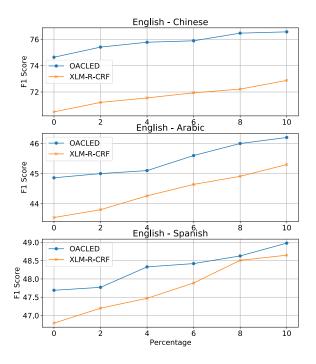


Figure 2: Model performance when training on small quantities of labeled target data. The X axis presents the percentage (0 - 10%) of data used out of the entire training set of the target language.

It can be observed that OACLED consistently outperforms the baseline even when there is some availability of annotated data. Additionally, performance steadily increases as more and more data is used. This conforms to expectations, and confirms that having labeled data in the target language available for training is ultimately beneficial to the model's performance.

3.6.3 Case Study

Next, we look into our model's predictions and analyse instances where it outperforms the baseline to exemplify the advantages of dealing with optimized language-invariant representations. We identify two important patterns.

First, our model seems to better classify events in the target language that involve trigger words that have distinct connotations that depend on context. Specially those that are two distinct words in the source language. For example, the Spanish word "juicio" can have two distinct meanings that are different words in English: "trial" and "judgement". Our model correctly classifies it as a JUSTICE:TRIAL-HEARING trigger in the sentence "Dos llamados a juicio fueron hechos por un jurado federal investigador". Meanwhile, the baseline fails to even recognize it as a trigger. Another example is the word "detenido", an adjective that can mean both "detained", in a criminal context, and "stopped", as in halted. Our model correctly classifies it in the sentence "Padilla no debería permanecer detenido durante meses alejado de otros reos" as a JUSTICE:ARREST-JAIL trigger while the baseline fails to detect the event. We manually identified 23 of these polysemous triggers in the Spanish² test set and found that 19 (82.6%) were correctly classified by our OACLED model versus 14 (60.8%) by the baseline (27.8%) improvement).

Additionally, we found our model correctly classifies verb conjugation variants that do not exist in the source language. For instance, our model correctly recognizes the words "venderlos", "vender", "vendes", and "vendedor" (variants of the verb "to buy") as TRANSACTION:TRANSFER-OWNERSHIP triggers whereas the baseline incorrectly classifies them as being of the TRANSACTION:TRANSFER-MONEY type. As previously mentioned, Majewska et al. (2021) propose injecting external verb-knowledge into the training to help with verb interpretation for event extraction. Our empirical results, however, outperform their reports which appears to imply that, at least for CLED, holistically learning the language-invariant features shared between the target and source languages works better than injecting language-specific verb knowledge.

We believe these findings illustrate how, by introducing additional context in the form of unlabeled data, the model is able to learn fine-grained word representations that better capture the semantics of the words in the target language, and successfully deal with difficult cross-lingual issues.

²We use Spanish for the analysis as it is the mother tongue of the first author.

4 Related Work

Research efforts on monolingual ED are extensive and varied. Hand-crafted, feature-based, languagespecific methods were the basis of early ED approaches (Ahn, 2006; Ji and Grishman, 2008; Patwardhan and Riloff, 2009; Liao and Grishman, 2010a,b; Hong et al., 2011; McClosky et al., 2011; Li et al., 2013; Miwa et al., 2014; Yang and Mitchell, 2016). More recent efforts have primarily made use of deep learning techniques such as convolutional neural networks (Nguyen and Grishman, 2015; Chen et al., 2015; Nguyen et al., 2016b), recurrent neural networks (Nguyen et al., 2016a; Sha et al., 2018; Lai et al., 2020), graph convolutional networks (Nguyen and Grishman, 2018; Yan et al., 2019; Nguyen et al., 2021a), adversarial networks (Hong et al., 2018; Zhang et al., 2019b), and pre-trained language models (Wadden et al., 2019; Zhang et al., 2019a; Yang et al., 2019; Zhang et al., 2020; Liu et al., 2020; Pouran Ben Veyseh et al., 2021b,a).

Works on cross-lingual ED are not as prevalent and generally make use of cross-lingual resources employed to address the differences between languages such as bilingual dictionaries or parallel corpora (Muis et al., 2018; Liu et al., 2019) and, more recently, pre-trained multilingual language models (M'hamdi et al., 2019; Hambardzumyan et al., 2020; Majewska et al., 2021). Unlike these previous efforts, our method leverages unlabeled data to further refine the language-invariant qualities of the language models.

Adversarial Language Adaptation, inspired by models in domain adaptation research (Ganin and Lempitsky, 2015; Naik and Rose, 2020; Ngo Trung et al., 2021), has been successfuly applied at generating language-invariant models (Joty et al., 2017; Chen et al., 2018; Nguyen et al., 2021b). Our method improves upon these approaches optimizing the adversarial training process by selecting the most informative examples from the unlabeled data.

Additional examples of downstream applications of cross-lingual learning are document classification (Holger and Xian, 2018), named entity recognition (Xie et al., 2018) and part-of-speech tagging (Cohen et al., 2011). For a thorough review on cross-lingual learning, we refer the reader to Pikuliak et al. (2021).

5 Conclusion

We present OACLED, a new model for crosslingual event detection that learns fine-grained language-invariant representations by optimizing the standard ALA training through optimaltransport-based sample selection. Our model achieves new state-of-the-art performance in our experiments on 8 different language pairs which demonstrate its robustness and effectiveness at generating refined language-invariant representations that allow for better event detection results. Our analysis of its intermediate outputs and predictions confirm that OACLED's representations are indeed closer to each other and that this proximity translates into better handling of difficult cross-lingual instances. We also note that, while this work focuses on the event detection task, our proposed optimization of the adversarial training process is task independent and can be generalized to other related IE tasks when leveraging ALA is deemed beneficial.

Acknowledgement

This research has been supported by the Army Research Office (ARO) grant W911NF-21-1-0112 and the NSF grant CNS-1747798 to the IU-CRC Center for Big Learning. This research is also based upon work supported by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via IARPA Contract No. 2019-19051600006 under the Better Extraction from Text Towards Enhanced Retrieval (BETTER) Program. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of ARO, ODNI, IARPA, the Department of Defense, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes notwithstanding any copyright annotation therein. This document does not contain technology or technical data controlled under either the U.S. International Traffic in Arms Regulations or the U.S. Export Administration Regulations. We thank the anonymous reviewers and Tracy King for their helpful feedback.

References

- David Ahn. 2006. The stages of event extraction. In *Proceedings of the Workshop on Annotating and Reasoning about Time and Events*.
- Xilun Chen, Yu Sun, Ben Athiwaratkun, Claire Cardie, and Kilian Weinberger. 2018. Adversarial Deep Averaging Networks for Cross-Lingual Sentiment Classification. In *Transactions of the Association for Computational Linguistics*.
- Yubo Chen, Liheng Xu, Kang Liu, Daojian Zeng, and Jun Zhao. 2015. Event extraction via dynamic multipooling convolutional neural networks. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Shay B. Cohen, Dipanjan Das, and Noah Smith. 2011. Unsupervised structure prediction with nonparallel multilingual guidance. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. In *CoRR*.
- Marco Cuturi. 2013. Sinkhorn distances: Lightspeed computation of optimal transport. In *Proceedings of the 26th International Conference on Neural Information Processing Systems Volume 2.*
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT*.
- Yaroslav Ganin and Victor Lempitsky. 2015. Unsupervised domain adaptation by backpropagation. In *Proceedings of the 32nd International Conference on Machine Learning*, pages 1180–1189.
- Karen Hambardzumyan, Hrant Khachatrian, and Jonathan May. 2020. The role of alignment of multilingual contextualized embeddings in zero-shot crosslingual transfer for event extraction. In *Collaborative Technologies and Data Science in Artificial Intelligence Applications*.
- Schwenk Holger and Li Xian. 2018. A corpus for multiligual document classification in eight languages. In Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC).
- Yu Hong, Jianfeng Zhang, Bin Ma, Jianmin Yao, Guodong Zhou, and Qiaoming Zhu. 2011. Using cross-entity inference to improve event extraction. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.

- Yu Hong, Wenxuan Zhou, Jingli Zhang, Guodong Zhou, and Qiaoming Zhu. 2018. Self-regulation: Employing a generative adversarial network to improve event detection. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 515–526, Melbourne, Australia. Association for Computational Linguistics.
- Heng Ji and Ralph Grishman. 2008. Refining event extraction through cross-document inference. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Shafiq Joty, Preslav Nakov, Lluís Màrquez, and Israa Jaradat. 2017. Cross-language learning with adversarial neural networks. In *Proceedings of the 21st Conference on Computational Natural Language Learning (CoNLL)*, pages 226–237.
- John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning*.
- Viet Dac Lai, Tuan Ngo Nguyen, and Thien Huu Nguyen. 2020. Event detection: Gate diversity and syntactic importance scores for graph convolution neural networks. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 5405–5411, Online. Association for Computational Linguistics.
- Qi Li, Heng Ji, and Liang Huang. 2013. Joint event extraction via structured prediction with global features. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Shasha Liao and Ralph Grishman. 2010a. Filtered ranking for bootstrapping in event extraction. In *Proceedings of the International Conference on Computational Linguistics (COLING)*.
- Shasha Liao and Ralph Grishman. 2010b. Using document level cross-event inference to improve event extraction. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Jian Liu, Yubo Chen, Kang Liu, Wei Bi, and Xiaojiang Liu. 2020. Event extraction as machine reading comprehension. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP).
- Jian Liu, Yubo Chen, Kang Liu, and Jun Zhao. 2019. Neural cross-lingual event detection with minimal parallel resources. 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 738–748, Hong Kong, China. Association for Computational Linguistics.
- Olga Majewska, Ivan Vulić, Goran Glavaš, Edoardo Maria Ponti, and Anna Korhonen. 2021. Verb knowledge injection for multilingual

- event processing. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing.*
- David McClosky, Mihai Surdeanu, and Christopher Manning. 2011. Event extraction as dependency parsing. In *BioNLP Shared Task Workshop*.
- Meryem M'hamdi, Marjorie Freedman, and Jonathan May. 2019. Contextualized cross-lingual event trigger extraction with minimal resources. In *Proceedings of the 23rd Conference on Computational Natural Language Learning (CoNLL)*.
- Makoto Miwa, Paul Thompson, Ioannis Korkontzelos, and Sophia Ananiadou. 2014. Comparable study of event extraction in newswire and biomedical domains. In *Proceedings of the International Conference on Computational Linguistics (COLING)*.
- Aldrian Obaja Muis, Naoki Otani, Nidhi Vyas, Ruochen Xu, Yiming Yang, Teruko Mitamura, and Eduard Hovy. 2018. Low-resource cross-lingual event type detection via distant supervision with minimal effort. In *Proceedings of the 27th International Conference on Computational Linguistics*.
- Aakanksha Naik and Carolyn Rose. 2020. Towards open domain event trigger identification using adversarial domain adaptation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*.
- Nghia Ngo Trung, Duy Phung, and Thien Huu Nguyen. 2021. Unsupervised domain adaptation for event detection using domain-specific adapters. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 4015–4025, Online. Association for Computational Linguistics.
- Minh Van Nguyen, Viet Lai, and Thien Huu Nguyen. 2021a. Cross-task instance representation interactions and label dependencies for joint information extraction with graph convolutional networks. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 27–38, Online. Association for Computational Linguistics.
- Minh Van Nguyen, Tuan Ngo Nguyen, Bonan Min, and Thien Huu Nguyen. 2021b. Crosslingual transfer learning for relation and event extraction via word category and class alignments. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Thien Huu Nguyen, Kyunghyun Cho, and Ralph Grishman. 2016a. Joint event extraction via recurrent neural networks. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*.

- Thien Huu Nguyen, Lisheng Fu, Kyunghyun Cho, and Ralph Grishman. 2016b. A two-stage approach for extending event detection to new types via neural networks. In *Proceedings of the 1st ACL Workshop on Representation Learning for NLP (RepL4NLP)*.
- Thien Huu Nguyen and Ralph Grishman. 2015. Event detection and domain adaptation with convolutional neural networks. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Thien Huu Nguyen and Ralph Grishman. 2018. Graph convolutional networks with argument-aware pooling for event detection. In *Association for the Advancement of Artificial Inteligence (AAAI)*.
- Trung Minh Nguyen and Thien Huu Nguyen. 2019. One for all: Neural joint modeling of entities and events. In Association for the Advancement of Artificial Inteligence (AAAI).
- Siddharth Patwardhan and Ellen Riloff. 2009. A unified model of phrasal and sentential evidence for information extraction. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Jonas Pfeiffer, Andreas Rücklé, Clifton Poth, Aishwarya Kamath, Ivan Vulić, Sebastian Ruder, Kyunghyun Cho, and Iryna Gurevych. 2020. AdapterHub: A framework for adapting transformers. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 46–54, Online. Association for Computational Linguistics.
- Matúš Pikuliak, Marián Šimko, and Mária Bieliková. 2021. Cross-lingual learning for text processing: A survey. In *Expert Systems with Applications*.
- Amir Pouran Ben Veyseh, Viet Lai, Franck Dernoncourt, and Thien Huu Nguyen. 2021a. Unleash GPT-2 power for event detection. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6271–6282, Online. Association for Computational Linguistics.
- Amir Pouran Ben Veyseh, Minh Van Nguyen, Nghia Ngo Trung, Bonan Min, and Thien Huu Nguyen. 2021b. Modeling document-level context for event detection via important context selection. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*.
- Lei Sha, Feng Qian, Baobao Chang, and Zhifang Sui. 2018. Jointly extracting event triggers and arguments by dependency-bridge rnn and tensor-based argument interaction. In *Proceedings of the Association for the Advancement of Artificial Intelligence (AAAI)*.
- Zhiyi Song, Ann Bies, Stephanie Strassel, Tom Riese, Justin Mott, Joe Ellis, Jonathan Wright, Seth Kulick,

- Neville Ryant, and Xiaoyi Ma. 2015. From light to rich ERE: Annotation of entities, relations, and events. In *Proceedings of the The 3rd Workshop on EVENTS: Definition, Detection, Coreference, and Representation*, pages 89–98, Denver, Colorado. Association for Computational Linguistics.
- C. Villani. 2008. *Optimal Transport: Old and New*. Grundlehren der mathematischen Wissenschaften. Springer Berlin Heidelberg.
- David Wadden, Ulme Wennberg, Yi Luan, and Hannaneh Hajishirzi. 2019. Entity, relation, and event extraction with contextualized span representations. 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).
- Christopher Walker, Stephanie Strassel, Julie Medero, and Kazuaki Maeda. 2006. Ace 2005 multilingual training corpus. In *Technical report, Linguistic Data Consortium*.
- Jiateng Xie, Zhilin Yang, Graham Neubig, Noah A. Smith, and Jaime G. Carbonell. 2018. Neural crosslingual named entity recognition with minimal resources. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Haoran Yan, Xiaolong Jin, Xiangbin Meng, Jiafeng Guo, and Xueqi Cheng. 2019. Event detection with multi-order graph convolution and aggregated attention. 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 5766–5770, Hong Kong, China. Association for Computational Linguistics.
- Bishan Yang and Tom M. Mitchell. 2016. Joint extraction of events and entities within a document context. In *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*.
- Sen Yang, Dawei Feng, Linbo Qiao, Zhigang Kan, and Dongsheng Li. 2019. Exploring pre-trained language models for event extraction and generation. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics (ACL)*.
- Junchi Zhang, Yanxia Qin, Yue Zhang, Mengchi Liu, and Donghong Ji. 2019a. Extracting entities and events as a single task using a transition-based neural model. In *IJCAI*.
- Tongtao Zhang, Heng Ji, and Avirup Sil. 2019b. Joint Entity and Event Extraction with Generative Adversarial Imitation Learning. *Data Intelligence*, 1(2):99–120.

Yunyan Zhang, Guangluan Xu, Yang Wang, Daoyu Lin, Feng Li, Chenglong Wu, Jingyuan Zhang, and Tinglei Huang. 2020. A question answering-based framework for one-step event argument extraction. In *IEEE Access*, vol 8, 65420-65431.