# Within-Document Event Coreference with BERT-Based Contextualized Representations

## Shafiuddin Rehan Ahmed, James H. Martin,

University of Colorado, Boulder shah7567, james.martin@colorado.edu

#### Abstract

Event coreference continues to be a challenging problem in information extraction. With the absence of any external knowledge bases for events, coreference becomes a clustering task that relies on effective representations of the context in which event mentions appear. Recent advances in contextualized language representations have proven successful in many tasks, however, their use in event linking been limited. Here we present a three part approach that (1) uses representations derived from a pretrained BERT model to (2) train a neural classifier to (3) drive a simple clustering algorithm to create coreference chains. We achieve state of the art results with this model on two standard datasets for within-document event coreference task and establish a new standard on a third newer dataset.

## Introduction

Event linking, or event coreference resolution, is the task of recognizing mentions of the same event either within a document or across different documents. Event linking is a critical component in information extraction pipelines since the facts about an event tend to be spread over many mentions, with each mention contributing partial information. Thus, a complete picture of an event can only be produced by accumulating information across many mentions.

Event linking is a challenging task due to the lexical diversity of event triggers and that, unlike entities, events typically lack explicit proper names. Moreover, unlike entity linking where named entities can often be linked to external resources such as Wikipedia or domain specific ontologies, event coreference resolution is typically based entirely on information gleaned from the documents themselves (Shen, Wang, and Han 2015; Raiman and Raiman 2018).

Further complicating the study of event coreference are the widely varying criteria that have been used to create annotated datasets. Over the years, relevant information about at events has included the event type (with respect to some ontology), predicate, argument fillers, and realis status (actual event, hypothetical, future, etc.). Early annotation efforts required matching types, predicates, arguments and realis status. This precise approach improves annotator agreement and classifier performance but results in a highly con-

Copyright © 2021, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

strained set of coreference judgments that conflicts with the purpose of gathering information about events that may be dispersed among mentions. More recent efforts using the Event Hopper approach for both within and across document coreference (Song et al. 2018) rely on annotators intuitions that two events are the same given looser guidelines for matching arguments and realis status.

Existing methods rely on feature engineering to capture relevant aspects of the context of event mentions, which are then used to train systems to make pair-wise coreference decisions. Recent neural methods have combined static word embeddings with explicit features in scoring event mentions. With the advent of deep contextualized language representations (Vaswani et al. 2017), there is an opportunity to leverage their strengths in the task of event coreference resolution. Specifically, these approaches (1) leverage large amounts of training data to address the issue of lexical diversity among event triggers, and (2) provide a natural semantic composition model to capture local context.

Here, we use BERT (Devlin et al. 2019) to generate joint contextualized representations of event pairs using the sentences in which they appear. Then, we train classifiers to score mention pairs. The resulting scores are then used to drive a clustering algorithm that produces the required set of coreferring mentions. We show that this approach paves the way for a fully neural approach that surpasses previous SoTA results on two standard event datasets and establishes a new standard for a newer dataset producing using the Event Hopper approach.

## **Related Work**

A great deal of previous work is based on modeling the probability of coreference between pairs of mentions. Such models are based on a supervised classifier trained over features extracted from the coreferent and non-coreferent pairs. The probabilities are then fed into a clustering method to find overall clusters of coreferent event mentions. Many methods (Bejan and Harabagiu 2010; Lee et al. 2012; Liu et al. 2014; Araki and Mitamura 2015; Cybulska and Vossen 2015) have relied on engineering features that model the relation between coreferent events. Liu et al. (2014) model it with a rich set 105 features that include lexical, syntactic, discourse and semantic features.

More recent work approach event coreference by

embedding the context surrounding the event mention (Krause et al. 2016; Kenyon-Dean, Cheung, and Precup 2018). The target is to train the embeddings to learn the implicit relations between coreferent events. They typically use static word2vec word embeddings (Mikolov et al. 2013) at the input layer. This way they eliminate the use of any extrinsic features.

Krause et al. (2016) use CNNs to generate contextualized representations of event mentions. While their approach is neural, they use features specific to the ACE 2005 corpus thereby making their approach domain specific. Kenyon-Dean, Cheung, and Precup (2018) also use word embeddings to embed the context (sentence and document) and the event mention in a way that maximizes the cosine similarity between coreferent events and minimizes it for non-coreferent ones.

#### **Datasets**

ACE 2005 (Walker et al. 2006), albeit an old dataset, is still a standard for evaluating event coreference. Among its annotations, the corpus provides gold standard event mentions, links between mentions, and event type information from the ACE ontology. We use the data split<sup>1</sup> of Krause et al. (2016) to compare our results against theirs and Liu et al. (2014).

	Train	Valid	Test	Total
# documents	500	49	60	599
# event instances	3186	425	479	4090
# event mentions	4158	574	617	5349

Table 1: ACE 2005 Corpus Statistics

ECB+ Corpus (Cybulska and Vossen 2014) is an extension of Event Coref Bank corpus (Bejan and Harabagiu 2008) with within and cross document event coreference annotations. This dataset includes ments from broader range of topics than ACE 2005. Cybulska and Vossen (2015) published a subset of the original ECB+ corpus that has been manually checked for correctness. We use their data split, and use thttps://www.overleaf.com/project/5f5144cbc622fa000136c12ehe processed input files made available To Barhom et al. (2019).our knowledge, Kenyon-Dean, Cheung, and Precup (2018) are the only ones that report within document event coreference results on this subset.

	Train	Valid	Test	Total
# documents	574	196	206	976
# event instances	1464	409	805	2741
# event mentions	3808	1245	1780	6833

Table 2: ECB+ Corpus Statistics

**LDC2019E77 English** (LDC 2020a,b) is the evaluation dataset for SM-KBP TAC challenge <sup>2</sup>. It is multilingual (En-

glish and Russian) and revolves around specific scenarios involving global conflicts. It has within document and cross-document coreference annotations. It follows the Event Hopper (Song et al. 2015) guidelines to annotate coreferent event mentions using a more fine-grained ontology. For evaluation of our work we use the English documents consisting 69 articles with 846 event mentions and 511 event instances.

Since, LDC2019E77 follow Event Hopper guidelines, the definition of corefering events is largely different from ACE 2005 and ECB+. The earlier datasets defined event coreference in very strict manner, where the ontological type of the events need to be the same, the realis status should match and they should not have conflicting or non-corefering arguments. Event Hopper defines event coreference in a more intuitive manner, which may not follow some or all restrictions imposed in the earlier datasets. This way, many more corefering events are annotated when compared to the earlier datasets.

## **Approach**

Following Krause et al. (2016), our approach has three components. Firstly, generation of the contextualized representation of pair of event mentions. With the generated representations, we try two techniques as coreference scorers. In the first method, we use the cosine similarity between the mention pairs as a metric to make the coreferencing decision. The second method is to train a logistic regressor using the joint representation of the pairs that scores for coreference. With the pair-wise coreference predictions that comes from these scorers, we finally cluster the event mentions by finding the transitive closure.

#### **Mention Pair Generation**

We generate the pairs of event mentions the same way as explained in Krause et al. (2016). Each mention is paired with all preceding mentions in the same document. For example if a mention  $v_2$  appears after the mention  $v_1$  in a document, the example pair is generated as  $(v_1, v_2)$ . The binary class labels follow from the gold standard coreference annotations. We also try two more pair generation methods where we sample only the pairs having the same (1) type and (2) lemma.

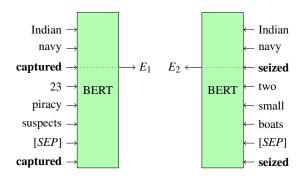


Figure 1: Generation of the contextualized representation of a pair of coreferent event mentions using BERT.

<sup>&</sup>lt;sup>1</sup>https://git.io/vwEEP

<sup>&</sup>lt;sup>2</sup>https://tac.nist.gov/2019/SM-KBP/

## **Contextualized Representation of Mentions**

We use a frozen pretrained BERT-Large model to generate contextualized representations of event mentions using the sentences in which they appear. Events are signaled in the input by adding the event trigger to end of the sentence separated by a SEP tag, a technique adapted from the semantic role labeling work of Shi and Lin (2019). The modified sentence is passed through the BERT tokenizer and transformer to generate a contextualized embedding for each token of the sentence. We then pick the output embedding of the event trigger or take the mean of all the individual word-piece tokens the mention may have been split into. In Figure 1, we illustrate the process of generation of mention representations for two coreferent event mentions. The embeddings  $E_1$ and  $E_2$ , extracted at the output layer of BERT, are the ones associated with the mentions captured and seized, respectively.

$\mathbf{Dev} \rightarrow$	ACE	ECB	LDC		
Model ↓	3 *	( cos <sub>+</sub> /cos_/ c	$\cos_\delta)$		
RobertA Large	0.96 0.93 0.03	0.99 0.99 0.00	0.98 0.97 0.01		
XLNET Large	0.96 0.96 0.00	0.96 0.96 0.00	0.99 0.99 0.00		
BERT Base	0.64 0.48 0.16	0.66 0.53 0.13	0.59 0.49 0.10		
BERT Large	0.54 0.30 <b>0.24</b>	0.51 0.31 <b>0.20</b>	0.46 0.34 <b>0.12</b>		

Table 3: The average cosine similarity values of the representations of coreferent  $(\cos_+)$  and non-coreferent pairs  $(\cos_-)$ , and their difference  $(\cos_\delta)$  generated by various language models.

Choosing the right Language Model: There has been a surge in the kinds of transformer based language models since the introduction of BERT. We experimented with the cased versions BERT-Large, BERT-Base, XLNet (Yang et al. 2019) and RoBERTa (Liu et al. 2019). The task is to choose the language model that encapsulates the similarities between coreferent and the dissimilarities between non-coreferent mentions the best. To achieve this task, we adopt a technique of comparing the average cosine similarities between the representations of coreferent pairs (cos<sub>+</sub>) and that of non-coreferent pairs (cos\_). We say that the representations from a model is of high quality if the differece between  $cos_+$  and  $cos_-$  ( $cos_\delta$ ) is high. We hypothesize the model of greater quality, i.e., one with higher  $\cos \delta$ , is more suitable for event coreference. As shown in Table 3, BERT-Large achieves the highest  $\cos_{\delta}$ , and therefore is used for subsequent tasks. This finding also leads us to believe that BERT's use of next sentence prediction during training makes it appropriate to the task of recognizing alternative expressions of the same event.

## **Coreference Using Cosine Similarity**

Our first approach of using the BERT representations to do event coreference is by using the cosine similarity. We say a pair of mentions is coreferent if the cosine similarity between their corresponding vectors is greater than a certain threshold. The task is to find a threshold that gives the best performance on the development set. We can quickly see why the choice of the language model is vital for deciding the threshold. As mentioned in the previous section, we want the cosine similarity of coreferent pairs to be much higher than that of non-coreferent ones. That way we might be able to reduce the false positives and negatives with the threshold.

$\mathbf{Dev} \rightarrow$	ACE	ECB	LDC					
Train $\downarrow$	3 *	3 * ( cos <sub>+</sub> / cos <sub>-</sub> / cos <sub>6</sub>						
BERT Large	0.54 0.30 0.24	0.51 0.31 0.20	0.46 0.34 0.12					
ACE	0.47 0.05 0.42	0.40 0.12 0.28	0.28 0.11 0.17					
ECB	0.57 0.19 0.38	0.49 0.07 <b>0.42</b>	0.41 0.15 <b>0.26</b>					
ACE + ECB	0.50 0.05 <b>0.45</b>	0.45 0.08 0.37	0.30 0.10 0.20					

Table 4: The average cosine similarity values that come out from the Cosine model trained on various datasets.

Improving Cosine Similarity: To make the metric more reliable we want the  $\cos_{\delta}$  of the representations to be higher. We train a linear regressor that transforms the representations of the mentions so that the  $\cos_{+}$  is close to 1 and  $\cos_{-}$  is close to -1. Since we are using BERT-Large, the hidden units of the hidden layer is 1024. We initialize the weights of the hidden layer as an identity matrix. We use Mean Squared Loss as the loss function for this regressor.

To train the regressor, we experiment with the training sets of ACE 2005 and ECB+ corpi. We report the change in the cosine similarity values of the various development sets in Table 4. The row corresponding to BERT-Large is the same as the row BERT-Large in Table 3. ACE + ECB represents the model trained using the merged training sets of ACE 2005 and ECB+. We observe that the model is consistent in distancing the cos<sub>+</sub> and cos<sub>-</sub> across all datasets. As far as quality of dataset is concerned, the model trained with ECB+ corpus generalizes well for LDC2019E77. Also, ACE + ECB, i.e., the model trained with merged datasets is beneficial for ACE.

## **Coreference Using Logistic Regressor**

In the second method, we train a logistic regressor on the joint representation of the pair of mentions to make the coreferencing decision. The joint representation of each pair is generated by concatenating both the individual mention representations and also the point-wise vector product of them. So in Figure 1, the joint representation is produced as  $[E_1, E_2, E_1 \circ E_2]$ . The joint representation is passed through a logistic regressor that makes the coreference decision.

The logistic regressor is a 2-layered neural network having 512 hidden units at the first layer with square activation. Activation at the output layer is LogSoftMax and the loss function is a negative log likelihood loss. We use AdamW Optimizer (Loshchilov and Hutter 2017) with  $lr = 5e^{-6}$ . We train the model for 50 epochs. The strategy of picking the best model will be explained in the later section. We implemented both the models using pytorch and huggingface<sup>3</sup> libraries.

## Clustering

The core algorithm is to find the transitive closure of coreferent mentions using the pairwise coreference decision. To do that, we find the connected components from the adjacency matrix of the mentions created by these pairwise decisions. We use the following techniques to create the adjacency matrix:

**Baselines** We use the following rules as our baselines:

**Singletons**. Assign an adjacency of 0 to each pair of event mentions. This way each mention becomes a singleton cluster.

**Type**. Adjacency between a mention pair is 1 if they are of the same type (if the type information is available).

**Lemma**. Adjacency between a mention pair is 1 if the head lemma of each mention is either exactly same or partially matching (one is included in other). We use spaCy.io's dependency parser and lemmatizer.

**Lemma & Type**. Adjacency created by taking the logical and of Lemma and Type.

**Our methods using BERT** We vary our approaches by changing the training set and also the pair generation process:

**Cosine**. We use the cosine similarity generated by models from our first approach. We use the threshold from Table 5 for the corresponding models. So, the adjacency between a mention pair is 1 if their cosine similarity generated by the model is greater than the model's threshold.

**Regressor**. The vanilla regressor method that uses the representations of a frozen BERT-Large as is. This method does not make any assumptions about the dataset, i.e., it does not use any annotated features like ontological type, modality, etc. In other words, all the pairs generated by the process as explained in earlier section are used during training and prediction.

**Regressor-Type**. Variation of the vanilla method where only pairs of the same type are used.

**Regressor - Lemma**. Variation of the vanilla method where only pairs having either same or partially matching head lemmas are used .

**Regressor - Cosine.** This is a combination of our two methods in which the output representation of the Cosine model is the input representation of the Regressor. We freeze the Cosine model while training the Regressor. We do not report

the results for the Type and Lemma sampling approaches because the Cosine model isn't trained using those sampling methods.

### **Evaluation Metrics**

Since most event coreference datasets lack linking information to external knowledge resources, event coreference systems are typically treated as clustering systems, with clustering inspired metrics used for evaluation. Over the years a number of metrics have been proposed, here we report results on the major metrics of CONNL 2012 (Pradhan et al. 2014):

**B**<sup>3</sup> (Bagga and Baldwin 1998): It is a mention based scoring metric that measures the average purity of the mentions, and the mentions' cluster with respect to the gold standard.

MUC (Vilain et al. 1995): It is a link based scoring metric that measures the number of missing and extra links in the system's output. Clusters are represented as transitive closure of the mentions during the calculation.

**CEAF-E** (**Luo 2005**): It is an entity/cluster based scorer that measures the alignment of the clusters of system output with the gold standard clusters.

BLANC (Recasens and Hovy 2011; Luo et al. 2014): It is a mention based scorer that calculates the rand index of the coreferent and non-coreferent pairs of mentions from the system output with respect to the coreferent and non-coreferent pairs of gold standard.

Issues with BLANC and CEAF-E: The evaluation of testing set is an estimate on the average performance of within document coreference on a number of documents. BLANC and CEAF-E measure the overall performance across all document. In BLANC, to get the non-coreferent pairs, each mention is paired against all the other mentions in all the documents of the testing set. Similarly, CEAF-E attempts to align the clusters of a document with all other clusters in all the documents. By definition of within document coreference, mentions or clusters of a document are never to be compared with those of other documents. These extraneous comparisons make these metrics unreliable for within document coreference.

Cosine Model	Dev Threshold
BERT Large	0.67 / 0.67
ACE	0.56
ECB	0.55
ACE + ECB	0.61 / 0.51

Table 5: The tuned thresholds on the development set of the model trained on corresponding training set.

<sup>&</sup>lt;sup>3</sup>https://github.com/huggingface/transformers

## Results

**Model Selection:** For the Cosine model, we need to determine an optimal threshold for coreference. We try values ranging from 0 to 1.0 and pick the threshold that yields the best average of B<sup>3</sup> and MUC F1 scores on the development set. While B<sup>3</sup> and MUC F1 scores are not reliable by itself, the average strongly correlates to better clustering performance. There are two threshold values for BERT-Large and ACE + ECB, one each determined using the ACE and ECB+development sets respectively. The threshold for the development set is then used for the corresponding test set to get the final results. We report the thresholds tuned for the models trained on the datasets as mentioned earlier in Table 5. Similarly, while training the Regressor we save the parameters of the iteration that gives the best average B<sup>3</sup> and MUC F1 scores on the development set.

ACE 2005	$\mathbf{B}^3$	MUC	Average
Method	F	'I score	in %
Cosine (BERT-Large)	82.36	45.25	63.81
Cosine (ACE)	84.76	49.68	67.22
Cosine (ACE + ECB)	85.84	50.50	68.17
Regressor	86.43	53.79	70.11
Regressor - Cosine (ACE)	88.39	55.51	71.95
Regressor - Cosine (ACE + ECB)	88.62	57.46	73.04

Table 6: Clustering results for ACE 2005 test set for different variations of our methods.

ECB+	$\mathbf{B}^3$	MUC	Average
Method	F	'I score	in %
Cosine (BERT-Large)	88.54	37.90	63.22
Cosine (ECB)	89.32	51.22	70.27
Cosine (ACE + ECB)	89.89	55.44	72.66
Regressor	90.53	53.22	71.88
Regressor - Cosine (ECB)	90.28	58.90	74.59
Regressor - Cosine (ACE + ECB)	90.54	59.88	75.21

Table 7: Clustering results for ECB+ test set for different variations of our methods.

In Tables 6 and 7, we report only the B<sup>3</sup> and MUC F1 scores and their average. A higher average correlates to better performance on other metrics, namely CEAF-E and BLANC. The first takeaway from these results is that neural regressor is a better approach than using just the cosine similarity for coreference. Indeed, the Regressor is able to model

many more intrinsic relations of event coreference that a single metric like cosine similarity cannot.

Secondly, looking at the improvements by using the Cosine model, we can confirm the hypothesis that improving the quality of the input representations improves coreference. Remember, the quality of the representation is determined by the cosine similarities of coreferent and noncoreferent mentions as shown in Table 4. For the Cosine model, it is intuitive that if the difference between the cosine similarity of coreferent and non-coreferent mentions is large, it is more reliable in terms of avoiding false positives. Seeing a similar improvement in the Regressor model confirms that neural scorers for coreference also benefit from better quality of mention representations.

LDC2019E77	$B^3$	MUC	CE	CoNLL	BLANC		
Method	F1 score in %						
Singletons	75.31	0	65.60	46.97	49.93		
Type	57.37	49.83	45.25	50.81	63.86		
Lemma	76.26	45.56	71.45	64.42	63.54		
Lemma-Type	75.05	34.21	68.67	59.31	59.58		
Regressor - Cosine (ACE + ECB)	77.56	49.22	70.28	65.69	65.14		
Cosine (ACE + ECB)	76.93	52.34	71.01	66.76	63.94		
Cosine (ECB)	77.03	60.15	72.7	69.99	66.61		

Table 8: Clustering results for the LDC2019E77 dataset for only English documents.

In Table 8 we report the baselines and results using our methods on the LDC2019E77 dataset. From the baselines, we observe that Lemma is pretty strong and addition of Type reduces the baseline's performance. This shows one of the many differences the annotations guidelines of this dataset when compared to ACE 2005. Remember, we want to use this dataset as an unknown test bench to compare our method against. From Tables 6 and 7 we confirmed that the best variation of our method is when we use ACE + ECB version of Cosine model. So the obvious choice is to choose a method for an unknown set that is built by combining the two datasets. To that end, we use the Regressor - Cosine (ACE + ECB) models trained on ECB+ as a coreference scorer for this dataset. We use the Regressor trained on ECB+ corpus because of its broad coverage on varied topics. We also, use Cosine (ACE + ECB) model with the threshold determined with ECB+ development set (0.51).

Unfortunately, the Regressor - Cosine (ACE + ECB) models fail short in achieving the best results. We achieve the best performance with the Cosine (ECB) model instead. Although, this result goes against our finding with the Regressor model, the finding from Table 4 still holds, i.e., the representations with better quality with respect to the dataset achieves better results. However, for an unknown dataset, there is no way of determining the quality before testing.

While the results of the neural modes fall short, it is in

ACE 2005		$B^3$			MUC		CEAF-E	CONNL	BLANC
<b>Model</b> ↓	P	R	F	P	R	F	F	F	F
Liu et al. (2014)	89.90	88.86	89.38	53.42	48.75	50.98	86.47	75.61	70.43
Krause et al. (2016)	90.52	86.12	88.26	61.54	45.16	52.09	-	-	73.31
Lemma & Type	75.60	90.15	82.23	42.41	68.84	52.48	75.96	70.22	70.44
Regressor - Cosine (ACE + ECB)	88.94	88.30	88.62	59.23	55.79	57.46	85.78	77.29	78.84
Regressor-Type	88.21	90.52	89.35	62.04	61.59	61.81	85.74	78.97	79.67
Perfect Disambiguator	100.0	90.15	94.82	100.0	68.84	81.54	93.8	90.05	80.42

Table 9: Comparison of clustering metrics of our best performing approach with previous bests for ACE 2005 corpus.

part expected because of the vast differences in the annotation guidelines of LDC2019E77 with the older datasets. There are many more events annotated as coreferent using the Event Hopper approach, that might not be marked as coreferent following the strict annotation guidelines of earlier dataset. The hope is to find a common ground between these datasets, and Cosine (ECB) seems to be able to achieve that the best amongst the various models.

Finally we compare our results with previous work for ACE 2005 and ECB+ corpuses. For ACE 2005, we compare against Krause et al. (2016) and Liu et al. (2014). Note that, we use the same data split of Krause et al. (2016). Since the data split of Liu et al. (2014) is different, we cannot fairly compare our results with them. As shown in Table 9, our approach Regressor - Cosine (ACE + ECB) surpass the previous best without using any extrinsic features on all metrics except B³ and CEAF-E. While we cannot directly compare our results with Liu et al. (2014), it wouldn't be wrong to say that our approach performs better on average. By using the Type information, we gain slightly more in terms of overall performance.

In Table 10, interestingly, our Lemma matching approach outperforms Kenyon-Dean, Cheung, and Precup (2018) which also uses a different lemma approach. This discrepancy might be due to improvements in the lemmatization model in the sPacy library. Our best model from Table 7 falls short against both Kenyon-Dean, Cheung, and Precup (2018) and the Lemma baseline approach. However, once we sample pairs by matching the lemma (Regressor-Lemma), we achieve marginal gains over the baseline on almost all the F1 metrics, thus setting a new state of the art for this dataset.

Error Analysis: We also do some error analysis on the ACE 2005 test set with the best model from Table 9, i.e., Regressor - Type. In general, even with the number of metrics, clustering results are hard to interpret. So, to get a feel of how well the model is performing, we change the problem to see how well the model is able to disambiguate mentions with the same lemma and also, how well the model does in finding alternate expressions of the same event. Go-

ing by the results in Table 11, it is clear that the model does reasonably well in the disambiguation part with a ratio of 0.82, but falls way short in finding positive pairs with different lemmas with a ratio of 0.32. Now, a model can always achieve 100% true negatives by never clustering mentions, or achieve 100% true positives by clustering all mentions in a single cluster. Finding the sweet spot that doesn't cluster mentions for the most part, but when it does cluster it does so correctly, is still an eluding problem. Although, our model achieves SotA results, we think there is a lot of scope for improvement in finding this sweet spot.

## Conclusion

We presented two fully neural approaches to do within document event coreference using BERT that advances previous SoTA results for ACE2005 and ECB+ corpora. We also showed the importance of the quality of mention representations and its overall effect on clustering performance. By combining the two approaches, we were able to establish a single model that performs well across datasets with similar definition of event coreference. We also showed the effectiveness of the neural models for transfer learning event coreference across vastly different datasets (LDC2019E77). Looking at the SoTA results, it is clear that generation or sampling of the pairs is crucial for the performance of the models presented. We believe use of better negative sampling approaches (Cai and Wang 2018; Sun et al. 2019) will boost the performance and hence is a possible area to work on in the future.

While we have showed the effectiveness of BERT for within document coreference, trying it for cross-document coreference would be the next obvious step. This task will give us an opportunity to explore how well these BERT representations can be composed when dealing with clusters of many mentions. Finally, with the presence of multilingual datasets like LDC2019E77, the use of multi-lingual BERT (Devlin et al. 2019) is an interesting research direction for cross-lingual event coreference.

ECB+ Corpus		$B^3$			MUC		CEAF-E	CONNL	BLANC
Model ↓	P	R	F	P	R	F	F	F	F
Lemma	95.71	88.89	92.71	76.32	57.29	65.45	88.78	82.31	77.62
Kenyon-Dean et al. (2018)	94	90	92	69	57	63	88	81	75
Regressor - Cosine (ACE + ECB)	90.82	90.26	90.54	63.04	57.02	59.88	86.22	78.88	75.93
Regressor-Lemma	96.12	90.00	92.97	78.49	58.09	66.76	85.74	83.03	78.06

Table 10: Comparison of clustering metrics of our best performing approach with previous best for ECB+ corpus.

ACE 2005	+ve pairs with same lemma	+ve pairs with different lemma	-ve pairs with same lemma	-ve pairs with different lemma
Actual	219	281	353	6249
Predicted	124	90	288	6082
Ratio	0.57	0.32	0.82	0.97
Perfect Disambi Ratio	1.0	0	1.0	1.0

Table 11: An error analysis on the test set of ACE 2005 to check how well the best model (Regressor - Type) is disambiguating mentions of the same lemma and also finding alternate expressions of the same event.

### References

Araki, J.; and Mitamura, T. 2015. Joint Event Trigger Identification and Event Coreference Resolution with Structured Perceptron. In *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing*, 2074–2080. Lisbon, Portugal: Association for Computational Linguistics. doi:10.18653/v1/D15-1247. URL https://www.aclweb.org/anthology/D15-1247.

Bagga, A.; and Baldwin, B. 1998. Algorithms for Scoring Coreference Chains. In *In The First International Conference on Language Resources and Evaluation Workshop on Linguistics Coreference*, 563–566.

Barhom, S.; Shwartz, V.; Eirew, A.; Bugert, M.; Reimers, N.; and Dagan, I. 2019. Revisiting Joint Modeling of Cross-document Entity and Event Coreference Resolution. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 4179–4189. Florence, Italy: Association for Computational Linguistics. doi:10.18653/v1/P19-1409. URL https://www.aclweb.org/anthology/P19-1409.

Bejan, C.; and Harabagiu, S. 2008. A Linguistic Resource for Discovering Event Structures and Resolving Event Coreference. In *LREC* 2008. URL http://www.lrec-conf.org/proceedings/lrec2008/pdf/734\_paper.pdf

Bejan, C.; and Harabagiu, S. 2010. Unsupervised Event Coreference Resolution with Rich Linguistic Features. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, 1412–1422. Uppsala,

Sweden: Association for Computational Linguistics. URL https://www.aclweb.org/anthology/P10-1143.

Cai, L.; and Wang, W. Y. 2018. KBGAN: Adversarial Learning for Knowledge Graph Embeddings. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, 1470–1480. New Orleans, Louisiana: Association for Computational Linguistics. doi:10.18653/v1/N18-1133. URL https://www.aclweb.org/anthology/N18-1133.

Cybulska, A.; and Vossen, P. 2014. Using a sledgehammer to crack a nut? Lexical diversity and event coreference resolution. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)*, 4545–4552. Reykjavik, Iceland: European Languages Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2014/pdf/840\_Paper.pdf.

Cybulska, A.; and Vossen, P. 2015. Translating Granularity of Event Slots into Features for Event Coreference Resolution. In *Proceedings of the The 3rd Workshop on EVENTS: Definition, Detection, Coreference, and Representation*, 1–10. Denver, Colorado: Association for Computational Linguistics. doi:10.3115/v1/W15-0801. URL https://www.aclweb.org/anthology/W15-0801.

Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Lan-*

- guage Technologies, Volume 1 (Long and Short Papers), 4171–4186. Minneapolis, Minnesota: Association for Computational Linguistics. doi:10.18653/v1/N19-1423. URL https://www.aclweb.org/anthology/N19-1423.
- Kenyon-Dean, K.; Cheung, J. C. K.; and Precup, D. 2018. Resolving Event Coreference with Supervised Representation Learning and Clustering-Oriented Regularization. In *Proceedings of the Seventh Joint Conference on Lexical and Computational Semantics*, 1–10. New Orleans, Louisiana: Association for Computational Linguistics. doi:10.18653/v1/S18-2001. URL https://www.aclweb.org/anthology/S18-2001.
- Krause, S.; Xu, F.; Uszkoreit, H.; and Weissenborn, D. 2016. Event Linking with Sentential Features from Convolutional Neural Networks. In *Proceedings of The 20th SIGNLL Conference on Computational Natural Language Learning*, 239–249. Berlin, Germany: Association for Computational Linguistics. doi:10.18653/v1/K16-1024. URL https://www.aclweb.org/anthology/K16-1024.
- LDC. 2020a. AIDA Phase 1 Evaluation Source Data. *Linguistic Data Consortium* LDC2019E42.
- LDC. 2020b. AIDA Phase 1 Evaluation Topic Annotations Unsequestered V2.0. *Linguistic Data Consortium* LDC2019E77.
- Lee, H.; Recasens, M.; Chang, A.; Surdeanu, M.; and Jurafsky, D. 2012. Joint Entity and Event Coreference Resolution across Documents. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, 489–500. Jeju Island, Korea: Association for Computational Linguistics. URL https://www.aclweb.org/anthology/D12-1045.
- Liu, Y.; Ott, M.; Goyal, N.; Du, J.; Joshi, M.; Chen, D.; Levy, O.; Lewis, M.; Zettlemoyer, L.; and Stoyanov, V. 2019. RoBERTa: A Robustly Optimized BERT Pretraining Approach. *CoRR* abs/1907.11692. URL http://arxiv.org/abs/1907.11692.
- Liu, Z.; Araki, J.; Hovy, E.; and Mitamura, T. 2014. Supervised Within-Document Event Coreference using Information Propagation. In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC-2014)*, 4539–4544. Reykjavik, Iceland: European Languages Resources Association (ELRA). URL http://www.lrec-conf.org/proceedings/lrec2014/pdf/646\_Paper.pdf
- Loshchilov, I.; and Hutter, F. 2017. Fixing Weight Decay Regularization in Adam. *CoRR* abs/1711.05101. URL http://arxiv.org/abs/1711.05101.
- Luo, X. 2005. On Coreference Resolution Performance Metrics. In *Proceedings of the Conference on Human Language Technology and Empirical Methods in Natural Language Processing*, HLT '05, 25–32. USA: Association for Computational Linguistics. doi:10.3115/1220575.1220579. URL https://doi.org/10.3115/1220575.1220579.
- Luo, X.; Pradhan, S.; Recasens, M.; and Hovy, E. 2014. An Extension of BLANC to System Mentions. In *Proceedings of the 52nd Annual Meeting of the Associa-*

- tion for Computational Linguistics (Volume 2: Short Papers), 24–29. Baltimore, Maryland: Association for Computational Linguistics. doi:10.3115/v1/P14-2005. URL https://www.aclweb.org/anthology/P14-2005.
- Mikolov, T.; Sutskever, I.; Chen, K.; Corrado, G. S.; and Dean, J. 2013. Distributed Representations of Words and Phrases and their Compositionality. In Burges, C. J. C.; Bottou, L.; Welling, M.; Ghahramani, Z.; and Weinberger, K. Q., eds., *Advances in Neural Information Processing Systems* 26, 3111–3119. Curran Associates, Inc. URL http://papers.nips.cc/paper/5021-distributed-representations-of-words-and-papers.nips.cc/paper/5021-distributed-representations-of-words-and-papers.nips.cc/paper/5021-distributed-representations-of-words-and-papers.nips.cc/paper/5021-distributed-representations-of-words-and-papers.nips.cc/paper/5021-distributed-representations-of-words-and-papers.nips.cc/paper/5021-distributed-representations-of-words-and-papers-paper/5021-distributed-representations-of-words-and-papers-papers-paper/5021-distributed-representations-of-words-and-papers-
- Pradhan, S.; Luo, X.; Recasens, M.; Hovy, E.; Ng, V.; and Strube, M. 2014. Scoring Coreference Partitions of Predicted Mentions: A Reference Implementation. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, 30–35. Baltimore, Maryland: Association for Computational Linguistics. URL http://www.aclweb.org/anthology/P14-2006.
- Raiman, J.; and Raiman, O. 2018. DeepType: Multilingual Entity Linking by Neural Type System Evolution. *CoRR* abs/1802.01021. URL http://arxiv.org/abs/1802.01021.
- Recasens, M.; and Hovy, E. 2011. BLANC: Implementing the Rand index for coreference evaluation. *Natural Language Engineering* 17: 485 510. doi:10.1017/S135132491000029X.
- Shen, W.; Wang, J.; and Han, J. 2015. Entity Linking with a Knowledge Base: Issues, Techniques, and Solutions. *Knowledge and Data Engineering, IEEE Transactions on* 27: 443–460. doi:10.1109/TKDE.2014.2327028.
- Shi, P.; and Lin, J. 2019. Simple BERT Models for Relation Extraction and Semantic Role Labeling. *CoRR* abs/1904.05255. URL http://arxiv.org/abs/1904.05255.
- Song, Z.; Bies, A.; Mott, J.; Li, X.; Strassel, S.; and Caruso, C. 2018. Cross-Document, Cross-Language Event Coreference Annotation Using Event Hoppers. In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*. Miyazaki, Japan: European Languages Resources Association (ELRA). URL https://www.aclweb.org/anthology/L18-1558.
- Song, Z.; Bies, A.; Strassel, S.; Riese, T.; Mott, J.; Ellis, J.; Wright, J.; Kulick, S.; Ryant, N.; and Ma, X. 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*, 89–98. Denver, Colorado: Association for Computational Linguistics. doi:10.3115/v1/W15-0812. URL https://www.aclweb.org/anthology/W15-0812.
- Sun, Z.; Deng, Z.; Nie, J.; and Tang, J. 2019. RotatE: Knowledge Graph Embedding by Relational Rotation in Complex Space. *CoRR* abs/1902.10197. URL http://arxiv.org/abs/1902.10197.
- Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A. N.; Kaiser, L. u.; and Polosukhin, I. 2017. Attention is All you Need. In Guyon, I.; Luxburg, U. V.; Bengio, S.; Wallach, H.; Fergus, R.; Vishwanathan, S.; and

Garnett, R., eds., *Advances in Neural Information Processing Systems 30*, 5998–6008. Curran Associates, Inc. URL http://papers.nips.cc/paper/7181-attention-is-all-you-need.pdf.

Vilain, M.; Burger, J.; Aberdeen, J.; Connolly, D.; and Hirschman, L. 1995. A Model-Theoretic Coreference Scoring Scheme. In *Proceedings of the 6th Conference on Message Understanding*, MUC6 '95, 45–52. USA: Association for Computational Linguistics. ISBN 1558604022. doi:10.3115/1072399.1072405. URL https://doi.org/10.3115/1072399.1072405.

Walker, C.; Strassel, S.; Medero, J.; and Maeda, K. 2006. ACE 2005 Multilingual Training Corpus. *Linguistic Data Consortium* LDC2006T06. URL https://catalog.ldc.upenn.edu/LDC2006T06.

Yang, Z.; Dai, Z.; Yang, Y.; Carbonell, J.; Salakhutdinov, R. R.; and Le, Q. V. 2019. XLNet: Generalized Autoregressive Pretraining for Language Understanding. In Wallach, H.; Larochelle, H.; Beygelzimer, A.; d'Alché-Buc, F.; Fox, E.; and Garnett, R., eds., Advances in Neural Information Processing Systems 32, 5753–5763. Curran Associates, Inc. URL

http://papers.nips.cc/paper/8812-xlnet-generalized-autoregressive-pretraining-for-language-understanding.pdf.