

Detecting Ongoing Events Using Contextual Word and Sentence Embeddings

Mariano Maisonnave
Universidad Nacional del Sur
Argentina
`mariano.maisonnave@cs.uns.edu.ar`

Fernando Delbianco
Universidad Nacional del Sur
Argentina
`fernando.delbianco@uns.edu.ar`

Fernando Tohmé
Universidad Nacional del Sur
Argentina
`ftohme@criba.edu.ar`

Ana Maguitman
Universidad Nacional del Sur
Argentina
`agm@cs.uns.edu.ar`

Evangelos Milios
Dalhousie University
Canada
`eem@cs.dal.ca`

Abstract

This paper introduces the Ongoing Event Detection (OED) task, which is a specific Event Detection task where the goal is to detect ongoing event mentions only, as opposed to historical, future, hypothetical, or other forms of events that are neither fresh nor current. Any application that needs to extract structured information about ongoing events from unstructured texts can take advantage of an OED system. The main contribution of this paper are the following: (1) it introduces the OED task along with a dataset manually labeled for the task; (2) it presents the design and implementation of an RNN model for the task that uses BERT embeddings to define contextual word and contextual sentence embeddings as attributes, which to the best of our knowledge were never used before for detecting ongoing events in news; (3) it presents an extensive empirical evaluation that includes (i) the exploration of different architectures and hyperparameters, (ii) an ablation test to study the impact of each attribute, and (iii) a comparison with a replication of a state-of-the-art model. The results offer several insights into the importance of contextual embeddings and indicate that the proposed approach is effective in the OED task, outperforming the baseline models.

Keywords: *Ongoing Event Detection; Information Extraction; Contextual Embeddings; BERT; RNN; CNN.*

31 1 Introduction

32 The Information Extraction (IE) task consists in extracting structured infor-
33 mation from unstructured natural language texts (Grishman, 1997). Event
34 Extraction (EE) is a subtask of IE, in which the goal is to detect and re-
35 trieve real-world events from those texts. An EE system usually performs
36 two different steps to complete the extraction of the events (Ahn, 2006).
37 The first step is to identify the event trigger (or anchor), which is the word
38 that most clearly expresses the occurrence of an event, and to classify it into
39 one of the predefined event types. This step is called Event Detection (ED).
40 The second step is to extract the arguments of the events, which refer to
41 participants (e.g. people) or attributes (e.g. time and place) associated with
42 the event. In this work, we focus only on a variation of the event detection
43 task and not on the argument extraction task.

44 The ED task has been addressed in the literature, both as a standalone
45 task (Duan et al., 2017; Jagannatha and Yu, 2016; Liu et al., 2018a; Nguyen
46 et al., 2016b; Nguyen and Grishman, 2015,1,1; Wu et al., 2014) and as a
47 part of an EE system (Ahn, 2006; Chen et al., 2015; Huang et al., 2017;
48 Li et al., 2013; Liu et al., 2018b; Nguyen et al., 2016a; Sha et al., 2018;
49 Zhang et al., 2019). There is a great incentive to study ED systems not
50 only because of their direct use in several applications but also because any
51 improvement in an ED system will impact directly on the performance of
52 any EE system implemented with it. EE and ED systems are crucial for
53 any application or domain that needs structured information and relies on
54 a large corpus of unstructured data. Some examples of this are question-
55 answering systems (Sha et al., 2018) and text summarization systems (Lee
56 et al., 2003). These systems are also useful for generating reports of the
57 information available for a domain (Adedoyin-Olowe et al., 2016; Chen et al.,
58 2018; Liu et al., 2009). These reports can help an expert to make decisions
59 or create policies to address an issue.

60 In this paper, we define and address the Ongoing Event Detection (OED)
61 task as part of a broader project that aims at detecting ongoing real-world
62 events and other relevant variables from news and social media with the

ultimate goal of learning causal models (Maisonnavé et al., 2020). In the broader project, we use time series-related techniques from Econometrics to learn such causal models (Granger, 1969). To build those time series, we require (i) all the event mentions in the news and (ii) the time those events took place. These two requirements defined the needs of our event detection task. Because of (i), we require only the detection of event mentions (trigger detection), not the event arguments nor the participants. Because of (ii), we require the event to be ongoing when reported (discarding past, future and hypothetical events). If the event is ongoing at the moment it is being reported, the news article’s date can be used as the moment of occurrence of the event. Because we consider an event as ongoing if it was in progress at the time it was reported, an old news article could refer to an ongoing event (if it was happening when the news article was released). Using this definition and the publication date, we can keep track of the beginning and ending of events in time, and we can try to find causal links between them (which is part of the broader project of which we made a pilot study in (Maisonnavé et al., 2020)). These two requirements result in the definition of a proposed novel task, to which we refer to as the Ongoing Event Detection (OED) task.

We divided the broader project into two steps, (1) the OED task and (2) the time series construction and causal structure learning from those time series. The two steps involved in building the tool are depicted in Figure 1. The work reported in this paper focuses on the first step only, i.e., on the OED task.

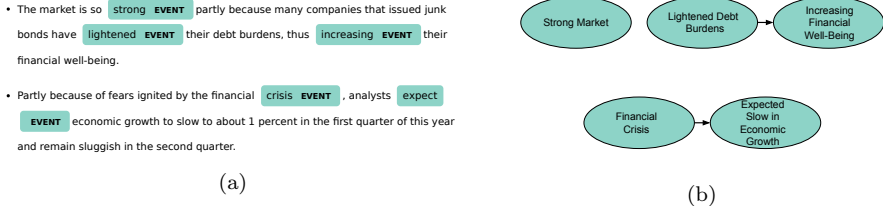


Figure 1: Two use cases for the OED tool, applied to two sentences (a). A causal graph manually extracted from the examples (b).

87 Motivated by the limitations of previous proposals to address our specific
88 requirements, we outlined a definition of event (described in Section 3) that,
89 on the one hand, is not limited by a fixed taxonomy of events and, on the
90 other hand, is centered on ongoing events only. Guided by this definition,
91 we manually labeled a dataset for training and testing purposes (described
92 in Section 4.1). Using this dataset, we developed an RNN model for event
93 prediction (described in Section 4.2). Since there are no previous studies
94 on this dataset, we also implemented two baselines for comparison. First,
95 we implemented a simple OED model based on a classical approach (SVM
96 model). Second, we replicated a baseline (CNN model) from the state-of-the-
97 art in ED to apply to our task (Nguyen and Grishman, 2015). Both baselines
98 are described in Section 4.3. The code for the proposed and baseline models,
99 as well as the dataset, are made available to the research community for
100 reproducibility and data reuse.¹ Results, discussion and conclusions are
101 presented at the end of the paper in Sections 5.2, 5.3, 6, respectively.

102 The contributions of this paper can be summarized as follows.

- 103 1. Firstly, we define the OED task for the first time, and we present a
104 manually labeled dataset for the task. An active learning tool was
105 developed and used for the construction of the dataset for the task.
106 According to the annotators, this tool (which was fine-tuned as new
107 labels became available) made the annotation process easier, there-
108 fore shortening the labeling time. The proposed task proved to be a
109 promising direction for implementing the first step of the framework
110 for causal structure learning depicted in Figure 1.
- 111 2. Secondly, we design and implement an RNN model for the OED task
112 that includes BERT embeddings (Devlin et al., 2018) as features. The
113 use of BERT embeddings for the task is extensively studied in this
114 work, and the results demonstrate the usefulness of context-sensitive
115 embeddings for this task. We also elaborate on a baseline by repli-

¹The code is available at https://cs.ums.edu.ar/~mmaisonnave/resources/ED_code/, and the dataset is available at https://cs.ums.edu.ar/~mmaisonnave/resources/ED_data/

cating a state-of-the-art model (Nguyen and Grishman, 2015) on our data, as well as a classical SVM model. The proposed RNN model outperforms the evaluated baselines.

3. Finally, we present an extensive empirical evaluation with different architectures and hyperparameters for both the baseline and our model. We also perform an ablation test to measure the impact of each feature on the performance. The full code of the experiments is made available to allow reproducibility.

2 Background and Related Work

The most widely used dataset for EE is the ACE 2005 Multilingual Training Corpus (Walker et al., 2006), which contains the complete set of English, Arabic, and Chinese training data for the 2005 Automatic Content Extraction (ACE) technology evaluation (Doddington et al., 2004). There are other EE and ED datasets available, such as TimeML (Pustejovsky et al., 2003) and SentiFM (Jacobs et al., 2018). However, to the best of our knowledge, all these datasets have a fixed taxonomy of valid event types.

Three main disadvantages arise when using existing EE datasets from the literature. First, as already mentioned, these datasets limit the number of events to a fixed set of possible events (those available in the taxonomy). This problem arises because datasets for the EE task require specific details about arguments for each new event type added to the taxonomy. However, for the ED task in general, and the OED task in particular, there is no need to consider arguments since the goal is to detect event triggers only. Second, when using an existing dataset developed for a particular new task or domain, the event taxonomy and the source of the news articles usually do not suit precisely the needs for the new task or domain. Third, existing datasets do not distinguish ongoing events (which are the focus of our work) from historical, future, hypothetical, or other forms of events. Because of these three disadvantages, we rely on creating our own dataset for our specific task.

146 The state-of-the-art approach for ED is to generate a neural-based clas-
147 sifier with a class for each possible event type, and an extra class for non-
148 events (Chen et al., 2015; Feng et al., 2018; Huang and Ji, 2020; Lai et al.,
149 2021; Li et al., 2020; Nguyen et al., 2016a; Nguyen and Grishman, 2015,1;
150 Tong et al., 2020b; Yao et al., 2021) (or one binary classification for each pos-
151 sible event type (Nguyen et al., 2016b)). Traditionally, these classifiers use a
152 wide variety of features representing lexical, syntactic or entity information,
153 which are typically the result of applying NLP tools (e.g. (Li et al., 2013)).
154 With the advent of neural network language models and the availability of
155 distributed continuous representations for words and sentences learned in
156 an unsupervised fashion from large corpora (such as word and sentence em-
157 beddings), the features used in the models changed radically. A lot of the
158 input features used in state-of-the-art approaches are unsupervised repre-
159 sentations automatically learned from big corpora (e.g., Word2vec (Mikolov
160 et al., 2013b) and Fasttext (Bojanowski et al., 2017)), and many of these rep-
161 resentations can keep improving during training through gradient descent.

162 According to (Boroş, 2018), ED approaches fall into one of three cat-
163 egories: pattern-based (Hobbs et al., 1992; Krupka et al., 1991; Riloff,
164 1996a,9; Yangarber et al., 2000), feature-based (Bronstein et al., 2015; Chieu
165 et al., 2003; Freitag, 1998; Hong et al., 2011; Huang and Riloff, 2011; Ji and
166 Grishman, 2008; Li et al., 2013; Liao and Grishman, 2010; Patwardhan and
167 Riloff, 2009; Surdeanu et al., 2006), and neural-based (Chen et al., 2015;
168 Feng et al., 2018; Nguyen et al., 2016a,1; Nguyen and Grishman, 2015). The
169 first two rely on sophisticated handcrafted rules, patterns, and features, as
170 well as NLP tools. The first two categories present several disadvantages.
171 First, creating the handcrafted rules is very time-consuming and will likely
172 result in high precision and low recall (since it is unlikely to cover all the
173 possible patterns manually). Second, selecting a good set of features is not
174 a trivial task, and building such features relies on sophisticated NLP tools
175 that have a certain level of error associated. The third category relies on
176 neural networks for both feature representation and token prediction (which
177 eliminates the need for handcrafted rules, NLP tools, and – in many cases –
178 feature selection). In this paper, we focus on the last category, which solves

a lot of the problems with the first two categories mentioned above, while 179
achieving state-of-the-art results for both ED and EE. 180

EE systems typically follow one of two possible architectures: pipelined 181
(Chen et al., 2015; Hong et al., 2011; Ji and Grishman, 2008; Liao and Grish- 182
man, 2010) or joint (Liu et al., 2018b; Sha et al., 2018; Zhang et al., 2019). 183
In the pipelined architecture, the ED task is the first step, in which the trig- 184
ger word is detected and classified. Afterward, the system performs the rest 185
of the EE task by extracting the arguments for those triggers. In the joint 186
architecture, the argument extraction is part of the event detection (trigger 187
detection) phase and vice versa. A description of this two architectures is 188
provided in Figure 2. 189

The notion of context information is not new for the ED task. In (Feng 190
et al., 2018), the authors use an RNN model to capture a simple notion of 191
context for each word to improve the performance on the task while using 192
classical embeddings. In (Fan et al., 2020), the authors show how BERT 193
embeddings improve the performance on the adverse drug event (ADE) de- 194
tection task in the medical domain. In (Tong et al., 2020a) BERT-based 195
text features are integrated with image features to address the ED task in 196
news. Even though these approaches are either applied to domains different 197
from news or incorporate non-textual information, they provide evidence 198
supporting the usefulness of exploiting contextual word representations for 199
ED. 200

On the other hand, the notion of ongoing events as defined in this paper 201
– to the best of our knowledge – is introduced here for the first time. In 202
(Huang et al., 2016), the authors use document classification to categorize 203
news articles containing events as past, ongoing, future planned, future alert, 204
and future possible. Although they point to the importance of distinguish- 205
ing ongoing from past and future events, their approach is different from 206
ours. They work at the document level, classifying full news articles, while 207
we work at the level of event triggers. Furthermore, they classify events into 208
different categories (past, ongoing, future), while we only determine if the 209
event is ongoing or not. Additionally, the dataset they introduce, EventSta- 210
tus corpus, is defined for a different context than ours. Their dataset consists 211

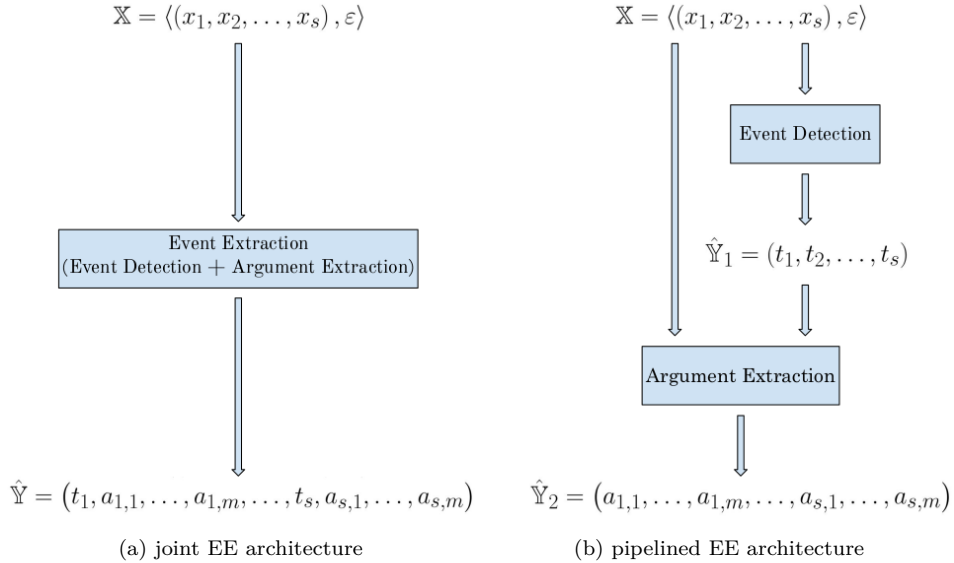


Figure 2: Depiction of the joint (a) and pipeline (b) architectures used for the EE task as defined in ACE (Walker et al., 2006), notation adapted from (Li et al., 2013). The input, in both cases, is a sentence separated in tokens (x_1, x_2, \dots, x_s) and a set of m candidate arguments $\varepsilon = \{e_k\}_{k=1}^m$ (which can be a gold-standard or predicted with a Named Entity Tagger). In the pipelined architecture, first, the input is used to predict the event trigger type (event detection), the output of this stage is represented with $\hat{\mathbb{Y}}_1 = (t_1, t_2, \dots, t_s)$. Here each $t_i \in \mathcal{L} \cup \{\perp\}$, where \mathcal{L} is the set of all possible event types and \perp is used for identifying non-events. Secondly, a argument extraction stage is performed using \mathbb{X} and $\hat{\mathbb{Y}}_1$ as input and producing $\hat{\mathbb{Y}}_2 = (a_{1,1}, \dots, a_{1,m}, \dots, a_{s,1}, \dots, a_{s,m})$ as output. Here $a_{i,j} \in \mathcal{R} \cup \perp$, where \mathcal{R} is the set of possible argument roles and \perp denotes that the current candidate argument is not an argument for the current trigger. Therefore $a_{i,j}$ represents the role of the j -th argument candidate for the i -th token. The joint architecture performs both steps at the same time (Event Detection + Argument Extraction), resulting in the $\hat{\mathbb{Y}}$ output which consists exactly of the $\hat{\mathbb{Y}}_1$ and $\hat{\mathbb{Y}}_2$ information combined. In the pipelined architecture the tasks are addressed independently and the error propagates downstream. In the joint architecture a technique is proposed for both tasks at the same time and usually takes advantage of the information of one task to complete the other.

212 of civil unrest events, such as protests, demonstrations and strikes.

213 3 The Ongoing Event Detection Task

214 The OED task is a specific ED task whose goal is to detect ongoing event
 215 mentions only, as opposed to historical, future, hypothetical, or other forms
 216 of events that are neither fresh nor current. Current states of affairs re-

ported in the news are also considered ongoing events. Since we need to 217
establish the date on which each real-world event took place, we use the 218
best approximation available, the publication date. For example, if a crisis 219
took place in 1994 and we can detect that an article from 1994 mentions 220
that crisis as current, we can use the publication date as the date on which 221
the event happened. On the other hand, if an article published years later 222
describes that crisis as a past one, we can ignore that mention (because it is 223
not ongoing). Using each real-world event and the time at which it occurred, 224
we can build a time series dataset to perform causal discovery (which is the 225
broader goal of our work (Maisonnavé et al., 2020)). 226

In this section, we present all the definitions required for the OED task, 227
and we present several key examples to facilitate the understanding of the 228
task. All these definitions, the examples for the task, and the annotators’ 229
guidelines are available on the dataset website.² 230

Definition 3.1 (Ongoing Event). An ongoing event is any text fragment in 231
a news article reporting a real-world event that meets any of the following 232
conditions: (i) it is a fresh event, (ii) it happened a time ago and is still 233
ongoing, or (iii) it is the current state of affairs for a given entity. 234

An example of (i) is when an earthquake just took place, and a news 235
article fragment covers that event. An example of (ii) is when a riot started 236
in a city some days ago and is again reported in the news while it is still 237
happening (because there is some new information or is a recapitulation of 238
what happened). Lastly, an example of (iii) is when a news article reports 239
a crisis or a recession that is taking place in some country or region. It is 240
important to notice that the same fragment could contain more than one 241
ongoing event. For example, some news article fragments could be reporting 242
an event caused by another. For instance, an ongoing crisis (an example of 243
(iii)) could cause a riot (an example of (ii)). 244

Ongoing Event Trigger and Ongoing Event Detection Task are defined 245
analogously to Event Trigger and Event Detection Task (Walker et al., 2006), 246
respectively, as follows: 247

²https://cs.unc.edu/~mmaisonnave/resources/ED_data/.

248 **Definition 3.2** (Ongoing Event Trigger). An Ongoing Event Trigger is the
249 word (verb, noun, or any other part of speech) that most clearly states the
250 occurrence of an ongoing event (Definition 3.1).

251 **Definition 3.3** (Ongoing Event Detection Task). The Ongoing Event De-
252 tection (OED) Task is the task of detecting the ongoing event trigger (Def-
253 inition 3.2).

254 In the OED task, the context of a word is crucial to determine if it
255 refers to an ongoing event or not. For example, take the word “crisis” in the
256 following sentences:

- 257 1. The current *crisis* will accelerate digital technology.
- 258 2. There will not be a *crisis* in the foreseeable future.
- 259 3. The same trend could be observed during the Global Financial *Crisis*
260 more than a decade ago.
- 261 4. Any financial *crisis* is catastrophic, and we must mitigate the risks of
262 a future *crisis*.

263 Only the reference to a “crisis” in sentence 1 is considered an event trig-
264 ger, while the other mentions of the same word are not. This is guided by
265 the need to distinguish ongoing events from those that are not, which is a
266 requirement of the ultimate goal of our broader project of detecting causal
267 relations between events reported in the news. Most existing approaches do
268 not adopt this context-sensitive definition of event. Furthermore, our defi-
269 nition of ongoing event also accounts for current states of affairs, which are
270 not taken into consideration by existing proposals. To address this context-
271 sensitive definition of event, we include BERT-based features. The context-
272 sensitive nature of the BERT model allows finding different representations
273 of the same word (for example, the word “crisis”) according to the context.
274 We believe these different representations will allow the model to correctly
275 distinguish ongoing events from non-events, as well as from hypothetical,
276 future or past events.

As another example, consider the following news extract: “devaluation
is not a realistic option to the current account deficit since it would only
contribute to weakening the credibility of economic policies as it did during
the last crisis.” The only word that is labeled as ongoing event trigger in
this example is “deficit” because it is the only ongoing event referred to
in the news. The word “devaluation” is not an ongoing event trigger as
a devaluation may not take place. Similarly, the word “weakening” is not
an ongoing event trigger as it is a hypothetical event. Finally, the word
“crisis” is not considered an ongoing event trigger as the news refers to a
crisis from the past. Note that the words “devaluation”, “weakening” and
“crisis” could be labeled as ongoing event triggers in other news extracts,
where the context of use of these words is different, but not in the given
example.

As stated in Definition 3.3 and illustrated with the provided examples,
the OED task can be addressed using a binary classification model, which
predicts for each token in a text fragment if it is an ongoing event trigger or
not. We adopt this approach and propose a binary classification model for
the task along with several baselines for comparison. The proposed model
is introduced in Section 4.2, and a high-level depiction of the proposal is
presented in Figure 5.

We provide additional details on the OED dataset building process and
the active learning tool in the next section (Section 4.1). A document with
the annotation guidelines is available at the dataset website (https://cs.uns.edu.ar/~mmaisonnave/resources/ED_data/). In that document a
brief introduction to the topic is given, and the OED task is presented (as
presented in this section), and the active learning tool is fully described to
the annotators.

304 4 Experimental Setup

305 4.1 Dataset

306 To build our dataset we tokenized the full New York Times (NYT) archive
307 (1987-2007) (Sandhaus, 2008) using the Spacy NLP library and divided the
308 news into sentences. From the full set of sentences extracted from the corpus
309 (~64 million), we selected a subset for labeling. We chose three episodes
310 of real-world crises: the Mexican peso crisis of 1994, the Russian financial
311 crisis of 1998, and the Asian financial crisis of 1997. We selected financial
312 episodes because, in the broader project (Maisonnave et al., 2020), we work
313 with economic experts who have an interest in applying these tools to the
314 economic domain. We set up the search engine Lucene³ (with the default
315 configuration) to search sentences related to these three episodes. We per-
316 formed a search using keywords manually selected by experts. Examples
317 of these keywords include “Mexico”, “crisis”, “debt”, “capital flight”, and
318 “devaluation”. From the obtained results, we randomly selected 2,000 sen-
319 tences. Also, we randomly selected from these results a separate set of 200
320 sentences for testing purposes.

321 We chose to use sentences instead of full texts to favor diversity in our
322 dataset. Because the annotation task is labor-intensive, by labeling sen-
323 tences instead of full texts we covered different events from different news
324 articles rather than repeated mentions of the same event discussed in vari-
325 ous parts of the same article. We handle each sentence as an independent
326 instance (independently of the article it came from). Therefore, the Lucene
327 search engine could return zero, one, or more sentences from each given ar-
328 ticle. However, given that there are approximately 64 million sentences in
329 2 million full-length articles, and we only selected 2,200 sentences in total,
330 it is likely that many of the sentences are drawn from different articles. We
331 believe the decision of using sentences instead of full texts does not simplify
332 the task. Instead, it changed the task to a different one, with different ad-
333 vantages and disadvantages. In general terms, and for the text classification

³<https://lucene.apache.org/>

task in particular (Cekik and Uysal, 2020), dealing with short-text presents
several problems like incomplete, inconsistent and sparse data, together with
high dimensionality. On the other hand, there are several advantages of us-
ing short texts. For example, since usually short-texts require less effort
to write, we can expect short-text datasets to have a higher number of in-
stances. For instance, while the New York Times corpus has roughly 2
million full-length articles, on Twitter, 500 million tweets are generated in
a single day (İnanç Arın et al., 2018). Another advantage of the short-text
context is that simpler models can be used without compromising perfor-
mance. Usually, this is not the case for full-length articles when the model
needs to keep track of long dependencies between words in the text. In this
context, we have short texts (sentences), which means we could use simpler
models. However, we have the problem that context was missing for some
instances (because important information was not part of the sentence being
labeled). In that sense, this task is similar to short-text classification.

Since in the proposed OED task we are not limited to a fixed set of events,
any ongoing real-world event or situation reported in the news articles was
considered an event. Consequently, each word in our dataset was labeled as
event trigger or *non-event trigger*, transforming the OED task into a binary
classification task. As mentioned earlier, it is important to distinguish those
events and situations that are in progress (or are reported as fresh events) at
the moment the news is delivered from past events that are simply brought
back, future events, hypothetical events, or events that will not take place.
In our dataset we only labeled as *event trigger* the first type of event. Based
on this criterion, some words that describe real-world events (like crisis,
recession, earthquake, flood, etc.) are labeled as *non-event triggers* if they
do not refer to ongoing events at the time the analyzed news is released. This
is illustrated by the example presented in Section 3 with different mentions
of the word “crisis”.

We developed a simple active learning tool to assist the labeling process
of the training and validation data (but not the test data). In Figure 3, we
present a screenshot of the active learning tool used by the annotators. This
tool used an early prototype of the RNN model used for event prediction

367 (described in Section 4.2) to suggest labels. Each sentence was presented
368 to a group of four annotators, along with the corresponding suggestions
369 generated by the model. Then the four annotators had to agree on keeping
370 the suggested labels, removing some of them, or adding new ones. The whole
371 process took a total of fifteen sessions of approximately two hours each. Since
372 the labeling process was carried out following a consensus-based approach,
373 no inner-agreement could be calculated. Therefore, such a metric is not
374 reported. Because we initially do not have any labeled instances for training,
375 we had a *cold start*. As new instances were labeled, the model was re-trained.
376 The re-training process took place each time 50 news instances were labeled,
377 and each re-training process was performed using the whole set of instances
378 labeled at that moment (including the newly labeled 50 instances). It is
379 important to notice that to avoid biasing the users' decisions when tagging
380 events in the test set, the four users were presented with sentences with no
381 suggestions provided by the tool and had to reach a consensus on which
382 words had to be marked as event triggers. Because the annotation process
383 is so labor-intensive, the size of the test set is relatively small. However, for
384 the results to be reliable, statistical tests were used on the results reported
385 in Section 5 to guarantee that the hypotheses tested in this work are valid.

386 The statistics of the dataset are presented in Tables 1 and 2. Additional
387 details on how the dataset was built can be found on the dataset website
388 (https://cs.uns.edu.ar/~mmaisonnave/resources/ED_data/).

389 4.2 RNN Model

390 We designed and trained a Recurrent Neural Network (RNN) for the OED
391 task. We conducted several experiments combining different hyperparame-
392 ters, features, and architectures. In this section, we will review the different
393 configurations tested and the intuition behind each selection.

394 **Features.** To be able to detect ongoing events in natural language text,
395 we hypothesize that syntactic, semantic, and grammatical information is
396 needed. To represent the semantics of each token we use Word2vec (Mikolov
397 et al., 2013a) (W). Note that we use Word2vec instead of other word em-



Figure 3: Overview of the active learning tool running for the first time. There are 2,000 sentences to be labeled, and 0 labeled. The system sorts the unlabeled list and presents it to the annotators. Because the system has not been trained yet (as there is no label information yet), the system’s suggestions have a performance no better than a random generator.

beddings (e.g., FastText (Bojanowski et al., 2017)) for the sake of compar- 398
 ison, as Word2vec is used in our baseline (Nguyen and Grishman, 2015). 399
 Word2vec embeddings represent each word by means of a 300-dimension 400
 vector that is fine-tuned as the model is trained for OED. We use the pre- 401

Full Dataset (training/validation + test)	
Sentence Count	2,200
Word Vocab. Size	8,647
Entity (E) Vocab. Size	34
Part-Of-Speech Simplified (P) Tag Vocab. Size	16
Dependency Parser (D) Tag Vocab. Size	47
Part-Of-Speech Detailed (T) Tag Vocab. Size	47

Table 1: Statistics about the OED dataset vocabulary.

402 trained Word2vec⁴ word embeddings to initialize the weights of a Keras
403 embedding layer. Fine-tuning is performed to update those weights as we
404 train the model for the OED task. We also tested context-sensitive tensors
405 provided by the Spacy Library (*Sp*). A Spacy context-sensitive tensor is a
406 96-dimension vector, which is the internal state of the neural model used for
407 NLP by the Spacy library⁵. This feature was not fine-tuned during training.

408 To encode the syntactic information, we first used the Spacy library to
409 identify the dependency tree (*D*). For each token, we extracted the de-
410 pendency relation with the head token in the dependency tree. Using that
411 information, we trained a 10-dimension supervised Keras embedding layer.
412 This layer was initialized with random weights and the representation was
413 incrementally fine-tuned along with the training of the rest of the network.
414 The grammatical information was encoded using two different embeddings

⁴<https://code.google.com/archive/p/word2vec/>

⁵<https://spacy.io/>

Metric	Training/Validation		Test	
	Total	Avg. per Sent	Total	Avg. per Sent
Token Count	76,629	38.31	7,382	36.91
Word Count	67,032	33.52	6,442	32.21
Entity Count	11,502	5.75	950	4.75
Event Count	5,119	2.56	416	2.08

Table 2: Total number of tokens, words, entities, and events found in the dataset.

representing information retrieved through Spacy’s Part-Of-Speech (POS) 415
tagger. We used both the simplified POS tag version (P) and the detailed 416
one (T) of the Spacy NLP library. As we did in the case of the syntactic 417
information, we used a first embedding layer to transform these two cate- 418
gorical variables into two vectors of 10 dimensions each and fine-tuned the 419
representation during training. 420

We used the Spacy Named Entity Recognizer tagger to include entity 421
information for the OED task in the form of a separate feature (E). We 422
represented each word as a two-part tag. The first part is the IOB notation 423
and the second is the type of the entity. As we did for other features (T , 424
 P , D), we used an embedding layer to transform a one-hot encoding of this 425
categorical input into a 10-dimension vector. This layer is initialized with 426
random weights and is fine-tuned during training. 427

Contextual embeddings, such as ELMo (Peters et al., 2018) and BERT 428
(Devlin et al., 2018), were proposed to solve the problem of mixed seman- 429
tics when the same word has a context-sensitive meaning. These embeddings 430
achieve ground-breaking performance on a wide range of NLP tasks. We hy- 431
pothesize that having vector sensitivity to the context is crucial for the OED 432
task, and it will bring a boost in performance. To illustrate this hypothesis, 433
consider the following two sentences: “The firm had to fire employees” and 434
“A fire burns a home near Brainerd airport.” They both use the word “fire”, 435
written in the same way, but with totally different meanings. Word2vec will 436
only have one vector for this word, which can limit the performance of an 437
ED model for identifying these cases. While Word2vec has the same rep- 438
resentation for the two mentions of the word “fire” in these two sentences, 439
the BERT model will produce different vectors for the two different uses 440
of the word “fire” (because it uses the context to build the representation). 441
Motivated by the intuition that the context is crucial for detecting ongoing 442
events, we incorporated pre-trained contextual word embeddings, in partic- 443
ular BERT embeddings (B), as an input to the model. We expect BERT 444
embeddings to have a favorable impact on the performance for instances 445
where the same word could be or not be an ongoing event depending on 446
context (like the example of the word “crisis” in Section 3). 447

448 Another problematic situation occurs with a word that preserves its se-
449 mantics across different sentences but that can refer to an ongoing event
450 trigger or not, depending on its context of use. This was illustrated in Sec-
451 tion 1 with the word “crisis”, which can be used to refer to a recent or
452 ongoing event (event trigger) or to a historical, future, hypothetical or other
453 forms of events that are neither fresh nor current (non-event triggers). Fol-
454 lowing this intuition, we incorporate for each token a 768-dimensional vector
455 feature representing a contextual embedding for the whole sentence. Note
456 that each token in the same sentence will have the same contextual sentence
457 embedding as input. The intuition behind using a common contextual sen-
458 tence embedding for each token in a sentence is that sentences usually give
459 information either about past, future and hypothetical situations (which
460 were irrelevant to us), or about current events (which are the ones that we
461 want to detect). It was uncommon for a sentence to mix these two types of
462 information. Guided by this intuition, we incorporate a contextual sentence
463 embedding (S) built by adding up the BERT embeddings for each token in
464 the sentence. The B and S inputs are not modeled with a Keras layer and
465 therefore they are fine-tuned during training for the OED task. A summary
466 of the eight above-described features used for the RNN models is presented
467 in Table 3. Since we explore different subsets of features, not all the fea-
468 tures are present in all the RNN models. The subset of features used in each
469 model is concatenated in one vector. For example, for the model with all
470 the features, the resulting concatenation is a vector of 1962 dimensions.

471 We implicitly use transfer learning in the various inputs of the model.
472 Firstly, by using the pre-trained Word2vec, Spacy and BERT embeddings,
473 we are incorporating semantic information extracted from large corpora of
474 unlabeled text. Secondly, by using the POS tagger and the Dependency
475 tagger from the Spacy library, we are including grammar and syntactic in-
476 formation coming from other IE fields, where there are large labeled corpora.

477 **Architecture.** We chose an RNN architecture based on Long-Short
478 Term Memory (LSTM) cells to exploit the dependencies between previous
479 and subsequent tokens for the classification of the current word. The inputs
480 for each word, which are eight embeddings, are all concatenated to form a

Abbr.	Size	Description
<i>W</i>	300	Pre-trained word embeddings.
<i>P</i>	10	Part-Of-Speech tag embeddings, simplified version.
<i>T</i>	10	Part-Of-Speech tag embeddings, detailed version.
<i>D</i>	10	Dependency parser tag embeddings.
<i>E</i>	10	Entity tag embeddings.
<i>Sp</i>	96	Spacy contextual word embeddings.
<i>B</i>	768	Pre-trained contextual word embeddings.
<i>S</i>	768	Contextual sentence embeddings.

Table 3: **Features** used by the RNN models. The subset of features used for each model is concatenated into one embedding. For example, if a model uses all the features, it will have a 1962-dimension input vector consisting of all the embeddings concatenated.

vector of 1962 dimensions. We added a Dropout layer after the concatenated embeddings. Following the dropout layer, we added a Bidirectional LSTM (Bi-LSTM) layer with 15 hidden units. A final dense layer with only one hidden unit was added to the Bi-LSTM Layer. The output of this final layer is the prediction. The RNN architecture is outlined in Figure 4. Although this architecture is the one used during the experiments on the data held-out for testing, we tested other architectures (different numbers of Bi-LSTM layers and hidden units). The results of these additional experiments with various architectures can be found in the supplementary materials.

Hyperparameters. For the embedding layers we use the default configuration, only changing the random initialization for the Word2vec embeddings. We use $p = 0.1$ for the Dropout layer and we set up every Bi-LSTM using the default configuration, adding only L1L2 regularization with values of 0.001 for both L1 and L2. For the final experiments on the data held-out for testing, we used a single Bi-LSTM layer architecture with 15 hidden units, which was the architecture with the best performance during the preliminary studies reported in the supplementary materials. Finally, we set up the dense layer using the sigmoid activation function.

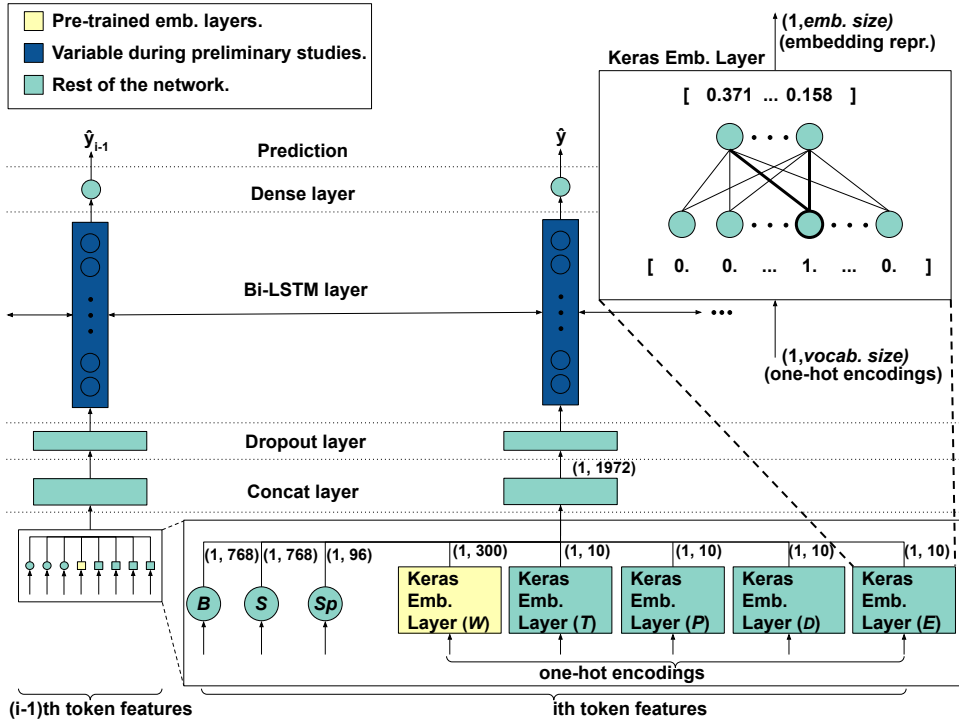


Figure 4: **Outline of the proposed RNN architecture.** From bottom to top. A representation of the eight inputs, for the i th and the $(i - 1)$ th token. For each token, three of the inputs are vector representations (B , S , Sp); the other five are one-hot encodings that enter five Keras embedding layers. Each of the dense layers is depicted in the top-right corner diagram. One of these embedding layers starts with pre-trained weights (W), the others with random weights. After the embedding layer, the eight vectors are concatenated together in a single 1972-dimensional vector. This vector enters a dropout layer. The output of the dropout layer enters a Bi-LSTM layer with fifteen hidden units. During the preliminary studies, we tested different numbers of Bi-LSTM layers and hidden units. For the held-out data, we use the configuration depicted in this Figure. The output of the Bi-LSTM layer is the input to a last dense layer, which makes the prediction.

499 **4.3 Baseline Models**

500 Since no previous model was validated using our dataset, existing metrics
 501 and results from the ED field are not directly comparable to ours. For this
 502 reason, we replicated an ED model from the state-of-the-art to use as a
 503 baseline, as well as a classical SVM model. Since the focus of our work
 504 is ED (and more specifically OED), we evaluated our model against state-

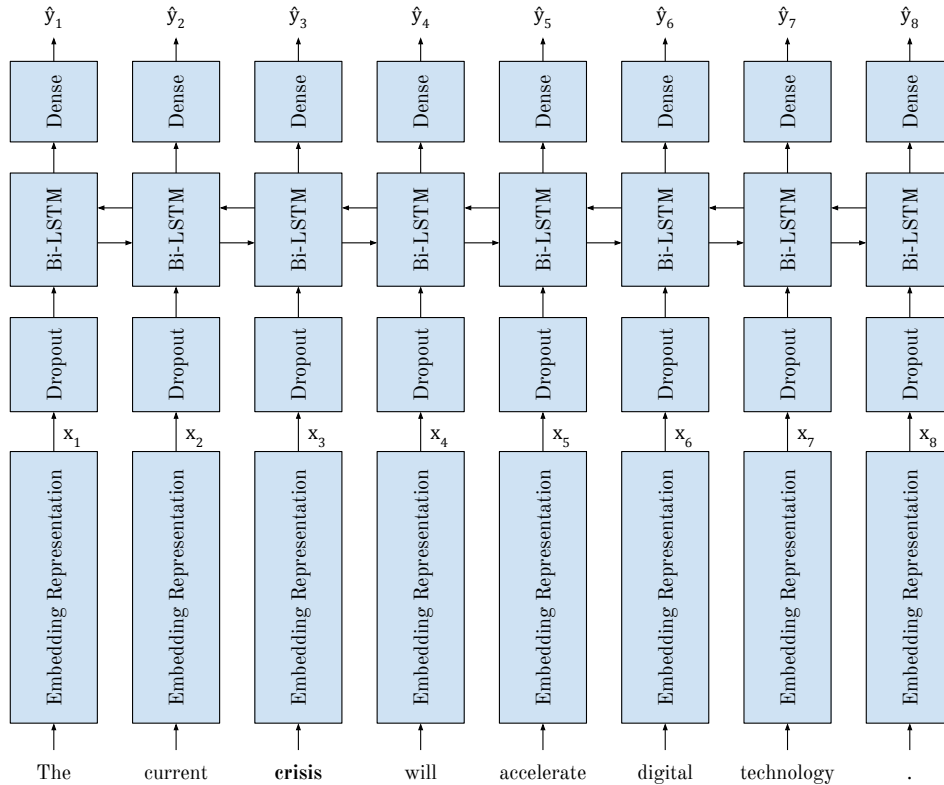


Figure 5: High-level description of the proposed methodology to solve the OED task. At the bottom of the Figure, the complete sentence is fed as input to an embedding representation layer. This first layer transforms each token into a 1962-dimension vector representation (x_i). Each of these 1962-dimension vectors is built concatenating eight embedding representations from different sources (Word2vec, BERT, and embeddings built using NLP tools (NER, dependency parser)). See Table 3 for a description of the eight embeddings used. A dropout layer is applied to x_i for regularization, followed by a layer of Bi-LSTM cells. Finally, a single-hidden-unit dense layer produces the prediction (\hat{y}_i). In this example sentence, we expect to have the word “crisis” labeled as an event trigger ($y_3 = 1$) and all the other words predicted as non-event-trigger ($y_i = 0$ for $i \neq 3$).

of-the-art in ED, instead of considering approaches that perform ED+EE. 505
Based on these considerations and because of its simplicity and comparable 506
performance with the state-of-the-art in ED, we selected and replicated the 507
Convolutional Neural Network (CNN) model proposed in (Nguyen and Gr- 508
ishman, 2015) as a baseline for comparison (in addition to the SVM baseline 509
model). Despite being proposed several years ago, this model is still compet- 510

511 itive with the available state-of-the-art methods. For example, it presents
512 an F1-score of only 4.1 percentage points lower than more recent models,
513 such as (Nguyen and Grishman, 2018). The selected baseline reports an
514 F1-score of 69% for the ED task using gold-standard entity annotations,
515 while the model from (Nguyen and Grishman, 2018) reports an F1-score of
516 73.1% on the same dataset (ACE 2005 Corpus) and using the same gold-
517 standard entity annotations. The authors do not report in (Nguyen and
518 Grishman, 2018) the result of the more recent neural model for predicted
519 entities. However, they do report in (Nguyen and Grishman, 2015) the F1-
520 score for predicted entities for the model we use as our baseline. For this
521 setting, the model achieves an F1-score of 67.6%.

522 As mentioned earlier, in order to compare our work with a classical
523 approach, we also implemented a SVM-based ED model using the Scikit-
524 learn library. The SVM approach used Word2vec embeddings as features
525 (W) and it was tested with four different kernels and default parameters.

526 The CNN model from (Nguyen and Grishman, 2015) was replicated as
527 faithfully as possible. However, although we tried to replicate the model ex-
528 actly as it is in the original paper, to adapt it to our dataset, we had to make
529 some minor changes. Furthermore, since the code was not available, some
530 minor implementation details may not be the same as in the original model.
531 We had to make some decisions about some aspects and configuration of
532 the models which were not explicit in the original paper. For example, we
533 had to decide which NLP tools to use for entity extraction. In our work, we
534 use the Spacy library for all the NLP related tasks. In the remainder of this
535 section, we describe in detail the CNN model used. We also describe some
536 minor changes and decisions that we had to make, explaining the rationale
537 behind them.

538 **CNN features.** As in the original paper, we use three input features for
539 the baseline model. First, we use the Word2vec embeddings as in our RNN
540 model (W). Second, we use entities embeddings (E). For building these
541 features, we use the Spacy Named Entity Recognition (NER) tagger and a
542 Keras embedding layer for building the embeddings, using the former to get
543 categorical variables (tags) and the latter for transforming these variables

into vectors. We chose Spacy as our NER tagger, as in the original paper 544
there is no mention of a specific tool, and because of the state-of-the-art 545
performance of the Spacy library in several NLP tools. The last feature 546
employed for the baseline model is a position embedding (Po), which rep- 547
resents the relative position of each word with respect to the current token 548
under classification. The word embedding layer starts with a representation 549
learned from the pre-trained Word2vec vectors, while the other two start 550
with random weights, as in (Nguyen and Grishman, 2015). The three layers 551
are updated while training for the ED task. 552

CNN architecture. The architecture used in (Nguyen and Grishman, 553
2015) is a one layer CNN, followed by a max-pooling layer and lastly a 554
dropout layer followed by a dense layer for the prediction. The inputs of this 555
network are three lookup tables, with the three types of embeddings. The 556
representation in these lookup tables improves along with the training for 557
the ED task. We replicated the same architecture and behavior, by replacing 558
the lookup tables with embedding layers that meet the same purpose: storing 559
and providing a vectorized representation, and improving the representation 560
while training for the ED task. The remaining of the network follows the 561
same architecture as in (Nguyen and Grishman, 2015). An outline of the 562
architecture of the CNN model is presented in Figure 6. 563

CNN hyperparameters. While many of the hyperparameters remain 564
the same, since our data is different from the one used in (Nguyen and 565
Grishman, 2015), the window size used had to be changed to better suit the 566
data. Since in our dataset, each data item is a sentence and not a whole 567
document like in ACE 2005 (the dataset used in (Nguyen and Grishman, 568
2015)), the window sizes had to be adjusted. We tried with smaller window 569
sizes, 1, 3, 5, 11, 21, and also for comparison sake we tested the model 570
with window size 31. We use the same number of filters (150), and the 571
same size for the filters (2, 3, 4, 5), as in the original paper. We used 572
the sigmoid activation function for the final dense layer to use the same 573
performance metrics as in our RNN model. We use, as in the original paper, 574
a batch size of 50, a probability for the dropout layer of 0.5, and we set the 575
hyperparameter for the l2 norms to 3. We used the binary Cross-Entropy 576

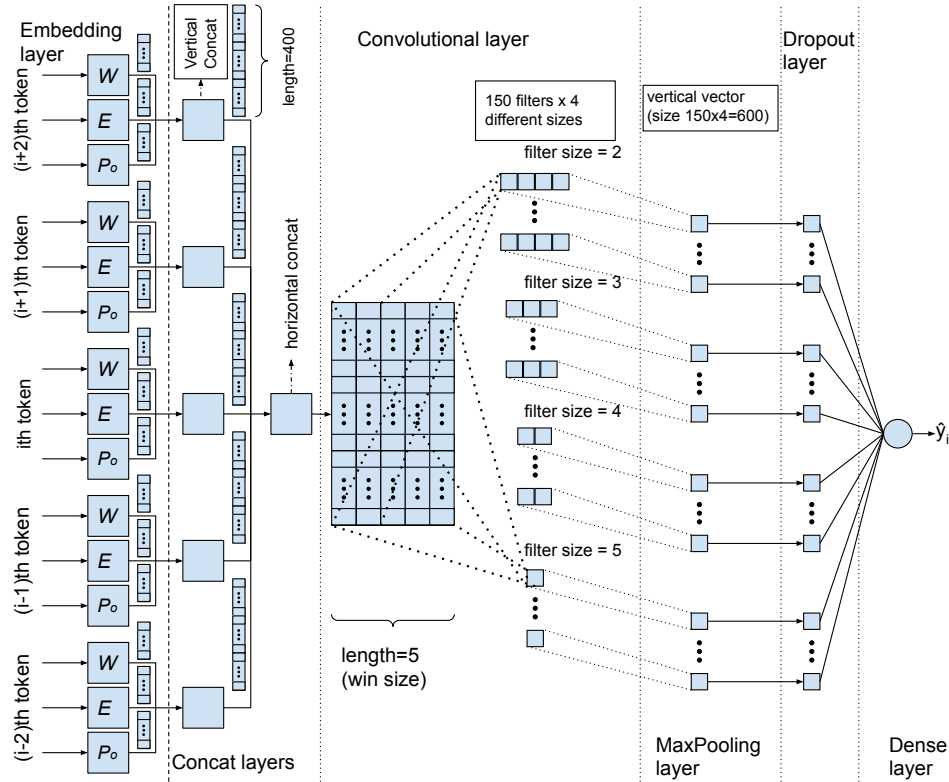


Figure 6: **Outline of the baseline model (CNN)** with window size 5. Descriptions are given from left to right. (a) Inputs. A representation of the three different inputs (Word embedding W , Entity embedding E and Position Embedding P_o) for the five tokens in the window. Each input is a one-hot encoding. (b) Embedding layer. Input enters a Keras embedding layer. (c) Concatenation layers. We concatenated the resulting embeddings in a single 400-dimension vector that represents each token. Another concatenation layer stacks these five vectors into one matrix. (d) A convolutional layer applies 600 filters of four different sizes (2, 3, 4, and 5) with 150 filters for each size. (e) MaxPooling layer. A max-pooling layer is applied afterward. (f) Dropout and Dense layers. Finally, the last two layers are a dropout and a dense layer. The dense layer is the last one of the network and is the one that makes the prediction.

5 Results and Discussion

578

5.1 Experiments Configuration and Reported Metrics

579

In this section, we present the results and discussions for different variations of the proposed (RNNs) and baseline models (CNNs and SVM). For each variation and each model (excluding the SVM models), we randomly split the training/validation part of the dataset into different training and validation subsets. For each model variation, we trained the model until it did not improve the F1-score on the validation set for 400 epochs (early stopping with patience 400). Although we set a hard limit on the number of epochs (2,000), the models were never trained for that many epochs. All the models stopped training due to the early stopping criterion (when they do not improve for 400 epochs, they stop automatically). We used consecutive seeds from 1 to 10 to guarantee replicability. The model with the best performance on the validation set was selected and used to make the predictions on the test set. In the case of the SVM models, since early stopping was not necessary, we use the training and validation sets together (training+validation) for training.

580
581
582
583
584
585
586
587
588
589
590
591
592
593
594

It is important to notice that we used the same held-out data to test all the models and this data was never used during the training stage. In this section, we present and discuss the performance of each model on the validation splits and test data. We report results for a total of fifteen different model variations. For selecting these model variations, we run several preliminary studies. The preliminary studies and their discussion are presented in the supplementary material. For each variation of the neural-network models, we report the average metrics for the ten trials, which include the average sensitivity, specificity, and harmonic mean between these two metrics, namely F1-score. We do not report the average accuracy for the model, because, in the presence of highly unbalanced data (93.41% are non-events), the accuracy is a misleading metric. To thoroughly analyze the F1-score metric on the test set we also report the confidence intervals (CI) at 95% level of confidence and the p-value for a t-test between the model considered and the best model of each table. For example, in the first row of Table 4, we

595
596
597
598
599
600
601
602
603
604
605
606
607
608
609

610 report the p-value of a single-tail t-test between the F1-score values achieved
611 by the first model (model RNN₁) and the best model of the table (model
612 RNN₇). We used the F1-score metric on the test data for the t-test. Since
613 the SVM model is not stochastic, only one run on the training+validation
614 and test data was executed for each of the four kernels. Consequently the
615 results reported are from a single run (not averages), and therefore no t-test
616 could be performed.

617 Additionally, we conducted Tukey’s tests between the performances (in
618 terms of F1-score in the validation set) for all pairs of RNN models (Fig-
619 ure 7a). We also conducted Tukey’s tests between all pairs of CNN models
620 and also on all CNN models against the best RNN model (in terms of aver-
621 age F1-score in the validation test), model RNN₇ (Figure 7b). Since Tukey’s
622 tests are used here to reduce the probability of joint errors in the analysis
623 of differences among different means in an ANOVA analysis, they allowed
624 us to obtain more robust results.

625 5.2 Experiments on the Test Set

626 **For the proposed RNN model**, we run ten trials for seven different sets of
627 features (models RNN₁ to RNN₇) as depicted in Table 4. Although we tried
628 different numbers of Bi-LSTM layers and hidden units, we only evaluated
629 on the test set the architecture that achieved the best performance during
630 the preliminary studies, i.e., the architecture with a single Bi-LSTM layer
631 with 15 hidden units. The results are reported in Table 4.

632 We conducted these experiments for two main reasons. First, to test the
633 hypothesis that the contextual word embeddings (B) are having a signifi-
634 cant positive impact on the overall performance. We tested this hypothesis
635 during the preliminary studies, and we found further evidence during the
636 experiments on the test dataset. The significant drop in performance (of
637 11.7 percentage points) from including or excluding the contextual word
638 embeddings (model RNN₁ vs. model RNN₂) is indicating the importance
639 of these features for the OED task. Second, we carried out the rest of the
640 experiments (models RNN₃ to RNN₇) as an ablation study to assess the

impact of the proposed features to the overall performance. We measure 641
the contribution of each feature during the preliminary studies. During the 642
experiments on the test dataset, we found further evidence that several in- 643
puts are not contributing to the performance in the presence of the other 644
features. Model RNN₇, which excludes all the attributes except the contex- 645
tual word and sentence embeddings, is the model with the best performance. 646
This result shows that the other features are not relevant in the presence 647
of these two other features, and they were adding unnecessary complexity 648
to the model. The statistical analysis between model RNN₇ and models 649
RNN₁, RNN₂, RNN₃, RNN₅, and RNN₆ shows that all the differences are 650
significant at the level of 95% (p-value<0.05). On the other hand, there is 651
less statistical evidence to support that the models RNN₇ and RNN₄ have 652
different performances (significant at the level of 90%, p-value<0.1). 653

To further analyze the differences in performance between the models, 654
we conducted an ANOVA analysis. The results of the ANOVA analysis 655
are presented in Figures 7a and 7b. The main conclusions (analyzed in the 656
caption of the figures) are that all models that contain contextual BERT 657
embeddings are statistically superior to model RNN₂ (the only model with- 658
out BERT embeddings) and that model RNN₇ (the best RNN model) is 659
statistically superior to all the CNN models. 660

The baseline CNN model was run for ten trials with different window 661
sizes. Here, we analyze windows of size 1 and 11 (other window sizes are 662
examined in the supplementary material). For both window sizes, we run 663
experiments with and without the entity embedding, but always including 664
the word embeddings. For window size 1, we excluded the position embed- 665
dings. The results of these experiments are presented in Table 5. The best 666
CNN baseline (model CNN₁) shows an F1-score of 0.575 on the test data, 667
which is significantly lower than the performance of all the proposed RNN 668
models, except for the one that excludes the contextual word embeddings 669
(model RNN₂). 670

To compare the RNN models and the best CNN baselines, we evaluated 671
the p-values of single-tail t-tests between the average F1-score for the worst 672
and the best RNN models (models RNN₂ and RNN₇, respectively) against 673

Model	Features	validation			test				
		\overline{sens}	\overline{spec}	$\overline{F1}$	\overline{sens}	\overline{spec}	$\overline{F1}$	F1 CI	p-value
RNN ₁	all	0.696	0.945	0.712	0.667	0.931	0.676	± 0.026	0.024
RNN ₂	all- $\{B\}$	0.545	0.957	0.591	0.522	0.948	0.559	± 0.046	0.000
RNN ₃	all- $\{T,P,D\}$	0.677	0.949	0.697	0.654	0.935	0.666	± 0.031	0.012
RNN ₄	all- $\{T,P,D,W\}$	0.703	0.948	0.718	0.686	0.935	0.690	± 0.018	0.076
RNN ₅	all- $\{T,P,D,Sp\}$	0.689	0.945	0.706	0.672	0.931	0.683	± 0.013	0.007
RNN ₆	all- $\{T,P,D,E\}$	0.686	0.949	0.706	0.662	0.936	0.675	± 0.023	0.012
RNN ₇	all- $\{T,P,D,Sp,W,E\}$	0.726	0.943	0.734	0.706	0.928	0.704	± 0.012	—

Table 4: Average performance of the **RNN** model in the validation and test set for seven different sets of features using an architecture with a single Bi-LSTM layer and 15 hidden units. We describe the features used by each model by indicating which features are removed from the full set of features (all), where the full set of features is $\{T,P,D,Sp,W,E,B,S\}$. We adopted this notation to simplify our experiments’ interpretation as an ablation study (comparing the model with all the features against the models with some features removed). The training set is used for training and the validation set for early stopping, while the test set is held out for testing. The resulting average performance during training is omitted from the reported analysis.

674 the two best CNN models (models CNN₁ and CNN₂). (1) Model RNN₇ is
675 statistically significantly better than all the others (RNN₂, CNN₁, CNN₂)
676 having a p-value of 0.0 for all the t-tests. The t-test between models CNN₁
677 and CNN₂ has a p-value of 0.384. (2) This high p-value indicates that the
678 difference in performance of these two models is not statistically significant,
679 which points to the fact the entity embeddings (E) have a negligent im-
680 pact on performance. (3) The difference in performance of the worst RNN
681 (model RNN₂) with respect to the best CNNs (models CNN₁ and CNN₂) is
682 not statistically significant either (p-value of 0.264 and 0.344, respectively).

Model	Win Size	Features	validation			test				
			\overline{sens}	\overline{spec}	$\overline{F1}$	\overline{sens}	\overline{spec}	$\overline{F1}$	F1 CI	p-value
CNN ₁	1	$\{W,E\}$	0.477	0.971	0.570	0.499	0.968	0.575	± 0.031	—
CNN ₂	1	$\{W\}$	0.472	0.972	0.565	0.493	0.969	0.569	± 0.031	0.384
CNN ₃	11	$\{W,E,Po\}$	0.507	0.963	0.596	0.326	0.960	0.394	± 0.028	0.000
CNN ₄	11	$\{W,Po\}$	0.499	0.965	0.589	0.345	0.942	0.406	± 0.035	0.000

Table 5: Average performance of variations of the **CNN** model on the validation and test sets. The training set is used for training and the validation for early stopping, while the test set is held out for testing. The average performance during training is omitted from the reported analysis.

Conclusions (1), (2) and (3) are confirmed not only by t-tests but also by Tukey’s tests.

The baseline SVM model was run on the training and test data with linear, polynomial, RBF and sigmoid kernels. The results of these experiments are presented in Table 6. The best SVM baseline (model SVM₄) shows an F1-score of 0.564 on the test data, which is significantly lower than the performance of all the RNN models except for the one that excludes the contextual word embeddings (model RNN₂). On the other hand, its performance is only slightly inferior to that achieved by the two best CNN baselines (models CNN₁ and CNN₂).

5.3 Discussion

In this section, we review and discuss the behavior of five representative models. The models selected for this discussion are model RNN₂ (RNN with all the features except for the contextual word embeddings), model RNN₇ (best RNN model in terms of average F1-score in the validation set), models CNN₁ (best CNN model for the same metric in the same set), and model SVM₄ (best SVM model for the same metric in the same set).

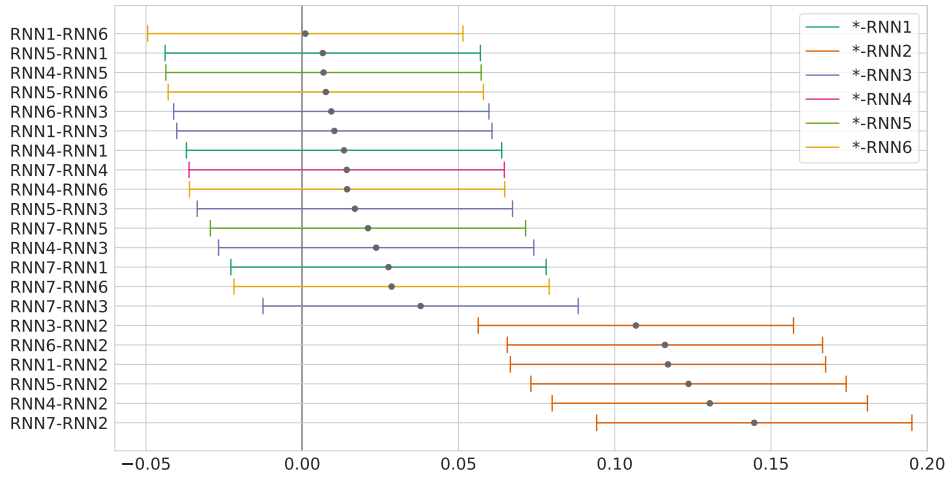
We hypothesize that the contextual embeddings are a crucial factor for boosting the performance of the proposed models. Therefore, in model RNN₂, we exclude only the contextual word embeddings to assess their impact on the performance. Since this model lacks the contextual word embeddings, we expected a poor performance in comparison to model RNN₇ (which has the contextual word and sentence embeddings). The poor performance of model RNN₂, which is lower than the best baseline models, provides evidence to support our hypothesis. Moreover, the superior performance achieved by model RNN₇ (containing only the contextual embeddings B and S as features) indicates that the inclusion of the other features (W , Sp , T , P , D , E) are not useful for the OED task once the contextual embeddings are included. It is also interesting to observe that model SVM₄, which has a *sigmoid* kernel, obtained an F1-score of 0.564 on the test set, outperforming some of the neural-based models (models RNN₂, CNN₃ and

Model	Kernel	training + validation			test		
		<i>sens</i>	<i>spec</i>	<i>F1</i>	<i>sens</i>	<i>spec</i>	<i>F1</i>
SVM ₁	linear	0.247	0.995	0.395	0.231	0.994	0.375
SVM ₂	polynomial	0.494	0.993	0.660	0.380	0.980	0.548
SVM ₃	RBF	0.459	0.993	0.628	0.346	0.984	0.512
SVM ₄	sigmoid	0.372	0.957	0.536	0.401	0.949	0.564

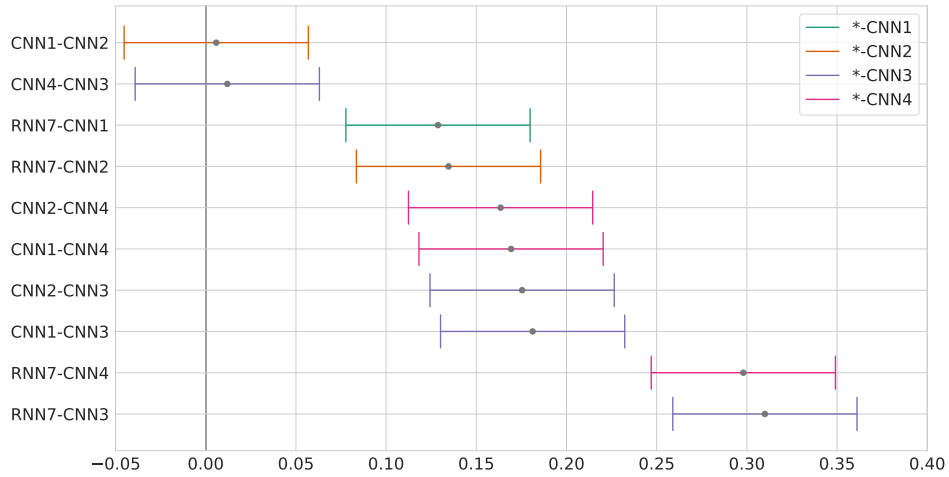
Table 6: Performance of variations of the **SVM** model on the training+validation and test sets. Since SVM has a deterministic output and there is no need for early stopping, there was no need to use the training and validation tests separately. Therefore, unlike the RNN and CNN models, the SVM models were trained using the training and validation sets together (training+validation set).

714 CNN₄).

715 We draw two major conclusions from the discussed results. First, in
716 an OED task trained on a small labeled datasets, the use of pre-trained
717 features is crucial. This follows from the fact that BERT-based contextual
718 embeddings (B and S) and Word2vec embeddings (W) had a major impact
719 on the performance of models RNN₇ and SVM₄, respectively. This provides
720 further evidence of the well-known fact that the use of transfer learning
721 in small data settings contributes significantly to improving performance.
722 Second, in the presence of a small labeled dataset, the choice of features, and
723 in particular the use of transfer learning, is more important than having a
724 large or complex model. This follows from the fact that a classical approach
725 (model SVM₄) that relies on pre-trained features outperformed or had a
726 performance comparable to that of some of the more complex neural-network
727 models.



(a)



(b)

Figure 7: The figure shows the difference in performance between each pair of models along with the confidence interval (CI) obtained using Tukey’s test at the level of significance of 95%. For each pair of models, if the CI contains the value zero, then there is not enough statistical evidence to conclude that the models’ performances are different. On the other hand, if the CI does not contain the value zero, then we can conclude that the difference in performance is significant. Unlike the pairwise test in Table 4 and 5, Tukey’s tests take into consideration the number of different means under an ANOVA analysis, i.e., the CIs are larger than those that could be obtained with the t-tests. Therefore, some pairs of models that we found statistically different using the t-test are probably found not statistically different due to larger confidence intervals. However, we still found evidence to support the main claim about the models (that the BERT embeddings provide the greatest improvement in performance for the task). This claim is supported by the fact that all the RNN models are significantly better than model RNN₂ (the only model without BERT embeddings) (Figure 7a) and that model RNN₇ (with BERT) has a significant leap in performance compared with all CNN models (Figure 7b).

728 6 Conclusions

729 The main contribution of this work is the definition of the OED task and an
730 extensive evaluation of an OED framework that combines several features.
731 The best of the proposed models, based on an RNN architecture and BERT-
732 based contextual word and sentence embeddings, shows an improvement of
733 13.3% in F1-score with respect to the best baseline model on the test data.
734 Considering that some more recent approaches for ED are slightly above our
735 baseline (4.1% for (Nguyen and Grishman, 2018)), we have enough evidence
736 to believe that our model achieves a performance competitive to the state-
737 of-the-art, even outperforming some of the most advanced models. It is also
738 worth mentioning that the code and data for replicating our model are fully
739 available.

740 Two main conclusions follow from our extensive analysis. First, contex-
741 tual embeddings are better suited for the OED task than the other analyzed
742 embeddings and features. In particular, some features proved to have a
743 negligible impact on performance, in the presence of other features. For ex-
744 ample, grammar information (captured using POS and Dependency Parser
745 taggers), entity tags and non-contextual word embeddings showed to have
746 no positive impact on performance as long as contextual embeddings are
747 accounted for in the model. The absence of a positive impact is not because
748 these features do not carry useful information, but probably because this in-
749 formation is already provided by contextual embeddings. Also, the analysis
750 of a classical baseline model showed that, in this context of small data, it is
751 sometimes more important to have useful features than to have a complex
752 or large model.

753 The second conclusion we derive from these results is that although the
754 proposed RNN models (which are best suited for text processing) show more
755 flexibility and effectiveness for the OED task than the analyzed CNN models
756 (which are best suited for computer vision tasks), the major difference in
757 performance between the evaluated models relies on the features used, in
758 particular the BERT-based contextual embeddings.

759 Another important contribution is the construction of a public dataset

for the OED task. The labeling of the dataset was assisted by an active
learning tool. The resulting dataset is particularly useful for the OED task
as it focuses on ongoing events only. It also differs from other ED and EE
datasets in being independent of a fixed event type taxonomy.

As part of our future work, we plan to evaluate the impact of contextual
embeddings on other ED settings (using other datasets and other models).

Acknowledgment

This work was enabled by support provided by CONICET, Universidad
Nacional del Sur (PGI-UNS 24/N051 and PGI-UNS 24/E145), ANPCyT
(PICT 2019-01640, PICT 2019-02302, and PICT 2019-03944), a LARA
project (Google Research Award for Latin America 2019-2021), Compute
Canada (<https://www.computecanada.ca>), a New Frontiers in Research
Fund Exploration Grant, an ELAP scholarship by the Department of For-
eign Affairs, Trade and Development Canada and ACENET ([https://
www.ace-net.ca/](https://www.ace-net.ca/)).

References

- Adedoyin-Olowe, M., Gaber, M. M., Dancausa, C. M., Stahl, F., and Gomes,
J. B. (2016). A rule dynamics approach to event detection in twitter
with its application to sports and politics. *Expert Systems with Appli-
cations*, 55:351 – 360.
- Ahn, D. (2006). The stages of event extraction. In *Proceedings of the Work-
shop on Annotating and Reasoning about Time and Events*, pages 1–8,
Sydney, Australia. Association for Computational Linguistics.
- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2017). Enriching
word vectors with subword information. *Transactions of the Associa-
tion for Computational Linguistics*, 5:135–146.
- Boroş, E. (2018). *Neural Methods for Event Extraction*. PhD thesis, Uni-
versité Paris-Saclay.
- Bronstein, O., Dagan, I., Li, Q., Ji, H., and Frank, A. (2015). Seed-based

- 789 event trigger labeling: How far can event descriptions get us? In
790 *Proceedings of the 53rd Annual Meeting of the Association for Com-*
791 *putational Linguistics and the 7th International Joint Conference on*
792 *Natural Language Processing (Volume 2: Short Papers)*, pages 372–
793 376.
- 794 Cekik, R. and Uysal, A. K. (2020). A novel filter feature selection method
795 using rough set for short text data. *Expert Systems with Applications*,
796 160:113691.
- 797 Chen, X., Zhou, X., Sellis, T., and Li, X. (2018). Social event detection with
798 retweeting behavior correlation. *Expert Systems with Applications*,
799 114:516–523.
- 800 Chen, Y., Xu, L., Liu, K., Zeng, D., and Zhao, J. (2015). Event extraction
801 via dynamic multi-pooling convolutional neural networks. In *Proceed-*
802 *ings of the 53rd Annual Meeting of the Association for Computational*
803 *Linguistics and the 7th International Joint Conference on Natural Lan-*
804 *guage Processing (Volume 1: Long Papers)*, pages 167–176, Beijing,
805 China. Association for Computational Linguistics.
- 806 Chieu, H. L., Ng, H. T., and Lee, Y. K. (2003). Closing the gap: Learning-
807 based information extraction rivaling knowledge-engineering methods.
808 In *Proceedings of the 41st Annual Meeting on Association for Compu-*
809 *tational Linguistics-Volume 1*, pages 216–223. Association for Compu-
810 tational Linguistics.
- 811 Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2018). Bert: Pre-
812 training of deep bidirectional transformers for language understanding.
813 *arXiv preprint arXiv:1810.04805*.
- 814 Doddington, G. R., Mitchell, A., Przybocki, M. A., Ramshaw, L. A.,
815 Strassel, S. M., and Weischedel, R. M. (2004). The automatic con-
816 tent extraction (ace) program-tasks, data, and evaluation. In *Lrec*,
817 volume 2, page 1. Lisbon.
- 818 Duan, S., He, R., and Zhao, W. (2017). Exploiting document level infor-
819 mation to improve event detection via recurrent neural networks. In
820 *Proceedings of the Eighth International Joint Conference on Natural*
821 *Language Processing (Volume 1: Long Papers)*, pages 352–361.

- Fan, B., Fan, W., Smith, C., and Garner, H. (2020). Adverse drug event
detection and extraction from open data: A deep learning approach.
Information Processing & Management, 57(1):102131.
- Feng, X., Qin, B., and Liu, T. (2018). A language-independent neural
network for event detection. *Science China Information Sciences*,
61(9):092106.
- Freitag, D. (1998). Information extraction from html: Application of a
general machine learning approach. In *AAAI/IAAI*, pages 517–523.
- Granger, C. W. (1969). Investigating causal relations by econometric models
and cross-spectral methods. *Econometrica: journal of the Econometric
Society*, pages 424–438.
- Grishman, R. (1997). Information extraction: Techniques and challenges.
In Pazienza, M. T., editor, *Information Extraction A Multidisci-
plinary Approach to an Emerging Information Technology*, pages 10–
27, Berlin, Heidelberg. Springer Berlin Heidelberg.
- Hobbs, J. R., Appelt, D., Tyson, M., Bear, J., and Israel, D. (1992). Sri in-
ternational: Description of the fastus system used for muc-4. Technical
report, SRI INTERNATIONAL MENLO PARK CA.
- Hong, Y., Zhang, J., Ma, B., Yao, J., Zhou, G., and Zhu, Q. (2011). Using
cross-entity inference to improve event extraction. In *Proceedings of
the 49th Annual Meeting of the Association for Computational Lin-
guistics: Human Language Technologies-Volume 1*, pages 1127–1136.
Association for Computational Linguistics.
- Huang, L. and Ji, H. (2020). Semi-supervised new event type induction and
event detection. In *Proceedings of the 2020 Conference on Empirical
Methods in Natural Language Processing (EMNLP)*, pages 718–724,
Online. Association for Computational Linguistics.
- Huang, L., Ji, H., Cho, K., and Voss, C. R. (2017). Zero-shot transfer
learning for event extraction. *arXiv preprint arXiv:1707.01066*.
- Huang, R., Cases, I., Jurafsky, D., Condoravdi, C., and Riloff, E. (2016).
Distinguishing past, on-going, and future events: The eventstatus cor-
pus. In *Proceedings of the 2016 Conference on Empirical Methods in
Natural Language Processing*, pages 44–54.

- 855 Huang, R. and Riloff, E. (2011). Peeling back the layers: detecting event role
856 fillers in secondary contexts. In *Proceedings of the 49th Annual Meeting
857 of the Association for Computational Linguistics: Human Language
858 Technologies-Volume 1*, pages 1137–1147. Association for Computa-
859 tional Linguistics.
- 860 Jacobs, G., Lefever, E., and Hoste, V. (2018). Economic event detection in
861 company-specific news text. In *Proceedings of the First Workshop on
862 Economics and Natural Language Processing*, pages 1–10.
- 863 Jagannatha, A. N. and Yu, H. (2016). Bidirectional rnn for medical event
864 detection in electronic health records. In *Proceedings of the conference.
865 Association for Computational Linguistics. North American Chapter.
866 Meeting*, volume 2016, page 473. NIH Public Access.
- 867 Ji, H. and Grishman, R. (2008). Refining event extraction through cross-
868 document inference. In *Proceedings of ACL-08: Hlt*, pages 254–262.
- 869 Krupka, G., Jacobs, P. S., Rau, L., and Iwanska, L. (1991). Ge: Description
870 of the nltoolset system as used for muc-3. In *THIRD MESSAGE UN-
871 DERSTANDING CONFERENCE (MUC-3): Proceedings of a Con-
872 ference Held in San Diego, California, May 21-23, 1991*.
- 873 Lai, V. D., Nguyen, M. V., Nguyen, T. H., and Dernoncourt, F. (2021).
874 Graph learning regularization and transfer learning for few-shot event
875 detection. In *Proceedings of the 44th International ACM SIGIR Con-
876 ference on Research and Development in Information Retrieval, SIGIR
877 ’21*, page 2172–2176, New York, NY, USA. Association for Computing
878 Machinery.
- 879 Lee, C.-S., Chen, Y.-J., and Jian, Z.-W. (2003). Ontology-based fuzzy event
880 extraction agent for chinese e-news summarization. *Expert Systems
881 with Applications*, 25(3):431 – 447.
- 882 Li, L., Jin, L., Zhang, Z., Liu, Q., Sun, X., and Wang, H. (2020). Graph
883 convolution over multiple latent context-aware graph structures for
884 event detection. *IEEE Access*, 8:171435–171446.
- 885 Li, Q., Ji, H., and Huang, L. (2013). Joint event extraction via structured
886 prediction with global features. In *Proceedings of the 51st Annual
887 Meeting of the Association for Computational Linguistics (Volume 1:*

- Long Papers*), pages 73–82. 888
- Liao, S. and Grishman, R. (2010). Using document level cross-event inference to improve event extraction. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 789–797. Association for Computational Linguistics. 889–892
- Liu, D.-R., Shih, M.-J., Liau, C.-J., and Lai, C.-H. (2009). Mining the change of event trends for decision support in environmental scanning. *Expert Systems with Applications*, 36(2, Part 1):972–984. 893–895
- Liu, J., Chen, Y., Liu, K., and Zhao, J. (2018a). Event detection via gated multilingual attention mechanism. In *Thirty-Second AAAI Conference on Artificial Intelligence*. 896–898
- Liu, X., Luo, Z., and Huang, H. (2018b). Jointly multiple events extraction via attention-based graph information aggregation. *arXiv preprint arXiv:1809.09078*. 899–901
- Maisonnave, M., Delbianco, F., Tohmé, F., Maguitman, A. G., and Milios, E. E. (2020). Assessing causality structures learned from digital text media. In *Proceedings of the ACM Symposium on Document Engineering 2020*, pages 1–4. 902–905
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. *arXiv preprint arXiv:1301.3781*. 906–908
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119. 909–912
- Nguyen, T. H., Cho, K., and Grishman, R. (2016a). Joint event extraction via recurrent neural networks. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 300–309. 913–916
- Nguyen, T. H., Fu, L., Cho, K., and Grishman, R. (2016b). A two-stage approach for extending event detection to new types via neural networks. In *Proceedings of the 1st Workshop on Representation Learning for NLP*, pages 158–165. 917–920

- 921 Nguyen, T. H. and Grishman, R. (2015). Event detection and domain adap-
922 tation with convolutional neural networks. In *Proceedings of the 53rd*
923 *Annual Meeting of the Association for Computational Linguistics and*
924 *the 7th International Joint Conference on Natural Language Process-*
925 *ing (Volume 2: Short Papers)*, pages 365–371.
- 926 Nguyen, T. H. and Grishman, R. (2016). Modeling skip-grams for event
927 detection with convolutional neural networks. In *Proceedings of the*
928 *2016 conference on empirical methods in natural language processing*,
929 pages 886–891.
- 930 Nguyen, T. H. and Grishman, R. (2018). Graph convolutional networks with
931 argument-aware pooling for event detection. In *Thirty-second AAAI*
932 *conference on artificial intelligence*.
- 933 Patwardhan, S. and Riloff, E. (2009). A unified model of phrasal and sen-
934 tential evidence for information extraction. In *Proceedings of the 2009*
935 *Conference on Empirical Methods in Natural Language Processing:*
936 *Volume 1-Volume 1*, pages 151–160. Association for Computational
937 Linguistics.
- 938 Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K.,
939 and Zettlemoyer, L. (2018). Deep contextualized word representations.
940 *arXiv preprint arXiv:1802.05365*.
- 941 Pustejovsky, J., Castano, J. M., Ingria, R., Sauri, R., Gaizauskas, R. J.,
942 Setzer, A., Katz, G., and Radev, D. R. (2003). Timeml: Robust
943 specification of event and temporal expressions in text. *New directions*
944 *in question answering*, 3:28–34.
- 945 Riloff, E. (1996a). Automatically generating extraction patterns from un-
946 tagged text. In *Proceedings of the national conference on artificial*
947 *intelligence*, pages 1044–1049.
- 948 Riloff, E. (1996b). An empirical study of automated dictionary construction
949 for information extraction in three domains. *Artificial intelligence*,
950 85(1-2):101–134.
- 951 Sandhaus, E. (2008). The new york times annotated corpus. *Linguistic Data*
952 *Consortium, Philadelphia*, 6(12):e26752.
- 953 Sha, L., Qian, F., Chang, B., and Sui, Z. (2018). Jointly extracting event

- triggers and arguments by dependency-bridge rnn and tensor-based
argument interaction. In *Thirty-Second AAAI Conference on Artificial
Intelligence*. 954
955
956
- Surdeanu, M., Turmo, J., and Ageno, A. (2006). A hybrid approach for the
acquisition of information extraction patterns. In *Proceedings of the
Workshop on Adaptive Text Extraction and Mining (ATEM 2006)*. 958
959
- Tong, M., Wang, S., Cao, Y., Xu, B., Li, J., Hou, L., and Chua, T.-S.
(2020a). Image enhanced event detection in news articles. In *AAAI*,
pages 9040–9047. 960
961
962
- Tong, M., Xu, B., Wang, S., Cao, Y., Hou, L., Li, J., and Xie, J. (2020b).
Improving event detection via open-domain trigger knowledge. In *Pro-
ceedings of the 58th Annual Meeting of the Association for Computa-
tional Linguistics*, pages 5887–5897, Online. Association for Compu-
tational Linguistics. 963
964
965
966
967
- Walker, C., Strassel, S., Medero, J., and Maeda, K. (2006). Ace 2005 mul-
tilingual training corpus. *Linguistic Data Consortium, Philadelphia*,
57. 968
969
970
- Wu, S., Bondugula, S., Luisier, F., Zhuang, X., and Natarajan, P. (2014).
Zero-shot event detection using multi-modal fusion of weakly super-
vised concepts. In *Proceedings of the IEEE Conference on Computer
Vision and Pattern Recognition*, pages 2665–2672. 971
972
973
974
- Yangarber, R., Grishman, R., Tapanainen, P., and Huttunen, S. (2000).
Automatic acquisition of domain knowledge for information extraction.
In *Proceedings of the 18th conference on Computational linguistics-
Volume 2*, pages 940–946. Association for Computational Linguistics. 975
976
977
978
- Yao, S., Shuang, K., Li, R., and Su, S. (2021). Fgcan: Filter-based gated
contextual attention network for event detection. *Knowledge-Based
Systems*, 228:107294. 979
980
981
- Zhang, T., Ji, H., and Sil, A. (2019). Joint entity and event extraction with
generative adversarial imitation learning. *Data Intelligence*, 1(2):99–
120. 982
983
984
- İnanç Arın, Erpam, M. K., and Saygın, Y. (2018). I-twec: Interactive clus-
tering tool for twitter. *Expert Systems with Applications*, 96:1–13. 985
986