# Ongoing Event Detection Task's Annotation Guidelines

## 1. Introduction

The Information Extraction (IE) task consists in extracting structured information from unstructured natural language texts. Event Extraction (EE) is a subtask of IE, in which the goal is to detect and retrieve real-world events from those texts. An EE system usually performs two different steps to complete the extraction of the events. The first step is to identify the event trigger, which is the word that most clearly expresses the occurrence of an event, and classify it into one of the predefined event types. This step is called EventDetection (ED). The second step is to extract the arguments of the events.

There is a great incentive to study ED systems not only because of their direct use in several applications but also because any improvement in an ED system will impact directly on the performance of any EE system implemented with it. EE and ED systems are crucial for any application or domain that needs structured information and relies on a large corpus of unstructured data. Some examples of this are question-answering systems and text summarization systems. These systems are also useful for generating reports of the information available for a domain. These reports can help an expert to make decisions or create policies to address an issue.

As part of this work, we define the Ongoing Event Detection (OED) task. The OED task could be briefly described as a specific ED task whose goal is to detect ongoing event mentions only, as opposed to historical, future, hypothetical, or other forms or events that are neither fresh nor current. States of affairs and changes of states reported in the news are also considered ongoing events. As part of our work, we also plan to build a dataset for the task, in which news article fragments are annotated with ongoing events. This guideline is aimed at the annotators of the news article fragments, providing useful information about the task and the labeling process. The remaining of this document is organized as follows. First, all the required definitions are provided, and the task itself is defined and exemplified to allow a better understanding by the annotators. Second, the data to be used is described, and the raw data collection is described. Lastly, the system to be used by annotators in the labeling process is described, and screenshots of the system illustrating the annotation process are presented.

# 2. The Ongoing Event Detection Task

In this section, we present all the definitions required for the OED task, and we present several key examples to facilitate the understanding of the task.

**Definition 1** (Ongoing Event). An ongoing event is any text fragment in a news article reporting a real-world event that meets any of the following conditions: (i) it as a fresh event, (ii) it happened a time ago and is still ongoing, or (iii) it is the current state of affairs for a given entity.

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

Ongoing Event Trigger and Ongoing Event Detection Task are defined analogously to Event Trigger and Event Detection Task[1], respectively, as follows:

**Definition 2** (Ongoing Event Trigger). An Ongoing Event Trigger is the word that most clearly states the occurrence of an ongoing event (**Definition 1**).

**Definition 3** (Ongoing Event Detection Task). The Ongoing Event Detection (OED) Task is the task of detecting the ongoing event trigger (**Definition 2**).

In the OED task, the context of a word is crucial to determine if it refers to an ongoing event or not. For example, take the word "*crisis*" in the following sentences:
1. The current *crisis* will accelerate digital technology.
2. There will not be a *crisis* in the foreseeable future.
3. The same trend could be observed during the Global Financial *Crisis* more than a decade ago.
4. Any financial *crisis* is catastrophic, and we must mitigate the risks of a future *crisis*.

Only the reference to a "*crisis*" in sentence 1 is considered an ongoing event trigger, while the other mentions of the same word are not. The annotators' task is to identify ongoing event triggers in each text fragment that the system will be presenting, assisted by the system suggestion. The labeling process is described in detail in Section 4.

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*

---

[1] https://catalog.ldc.upenn.edu/LDC2006T06

*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.

# 3. Data Collection

To build our dataset, we tokenized the full New York Times (NYT) archive (1987-2007) using the Spacy NLP library and divided the news into sentences. From the full set of sentences extracted from the corpus (64 million), we selected a subset for labeling. We chose four episodes of real-world crises: the **Mexican peso crisis of 1994**, the **Russian financial crisis of 1998**, the **Asian financial crisis of 1997**, and the **Argentine financial crisis of 2001**. We set up the search engine Lucene (with the default configuration) to search sentences related to these four episodes. We performed a search using keywords manually selected by experts. Examples of these keywords include "Mexico", "crisis", "debt", "capital flight", and "devaluation" (the full list is available in Annex A). From the obtained results, we randomly selected 2,000 sentences. Also, we randomly selected from these results a separate set of 200 sentences for testing purposes. The statistics of the raw dataset are presented in Tables 1 and 2.

The annotators' task is to identify ongoing event triggers in those 2,200 sentences. Each of the 2,000 training sentences will be presented, one by one, to the annotators with suggestions of ongoing event triggers. The annotators will have to correct the suggestions and agree on the correct set of ongoing event triggers for that sentence. Note that each sentence could have zero, one, or more ongoing event triggers. The final 200 sentences, the test set, will be presented to the annotators without suggestion, and they will have to agree on the correct set of ongoing event triggers for that sentence. The labeling process is described in more detail in Section 4.

| Full Dataset (training/validation + test) | |
|---|---:|
| Sentence Count | 2200 |
| Word Vocab. Size | 8647 |
| 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.

| Metric | Training | | 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 |

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

# 4. The Labeling Process

We developed a simple active learning tool to assist the annotators in the process of labeling the 2,200 sentences with ongoing event trigger information. The system will present the 2,000 training sentences with suggested labels, and the annotators will have to agree on the correct set of labels and correct the suggestions accordingly. The 200 test sentences will have no suggestions from the active learning tool, and the annotators will have to agree on the correct set of labels from the sentence alone. Each token (word) of each of the 2,200 sentences will have to be individually assigned to one of the two possible labels: ongoing event trigger or non ongoing event trigger.

Because the whole set of 2,000 training sentences has to be labeled entirely, we applied an active learning strategy specifically developed to facilitate the annotation task. Instead of choosing the complicated instances for labeling, we first chose the instances that were the easy ones according to the system's confidence. The intuition behind this preference is that, as more complicated instances start to appear, the system already has a pool of labeled instances to make smarter suggestions. So the idea is that the system goes incrementally from easy to hard instances, with the idea that in the beginning, when the support provided by the system is limited, the instances are easy. Afterward, as the instances start to get complicated, the system could deliver better suggestions. In the remainder of this section, we will present the active learning tool and describe its usage.

The tool has only one view from which all the functionality is available. An overview of the tool's view is presented in Figure 1. The tool is divided into five horizontal stacked panels: (1) sentence details, (2) labeled and unlabeled lists, (3) the sentence, (4) buttons and (5) log information.

**Panel 1**
Panel 1 displays information about the current number of sentences and the details on the selected sentence. On the column of the left, it shows how many labeled and unlabeled instances are in each list displayed in panel 2. The system displays the file identification, sentence number, and date of the currently selected sentence on the middle and right column.

**Panel 2**

When labeling the training set, in panel 2, the 2,000 sentences will be available for selecting them. In the beginning, the 2,000 sentences will be on the unlabeled list. Similarly, when labeling the test set, the 200 sentences will initially be available in the unlabeled list. As the sentences are labeled, they will move from the unlabeled to the labeled list. The annotators can select any sentence from the unlabeled list for labeling; when selected, the sentence will appear on panel 3 with the system suggested labels. Similarly, any sentences from the labeled list can be selected for the annotators to modify the previously assigned labels. Sentences selected from the labeled list will appear on panel 3 with the previously assigned labels. When a sentence is displayed in panel 3 by using its functionality, the annotators can change any label from the sentence. One key component of the unlabeled list is the order in which the sentences appear. The unlabeled list is sorted by the underlying active learning tool following the previously mentioned approach (from easy to hard to label sentences). So, although the annotators can potentially choose any sentence, the recommendation is that they pick from the first elements of the list. Whenever the underlying active learning tool changes (after each re-training process), the score for each unlabeled instance changes. Therefore, after each re-training process, the unlabeled list is re-sorted.
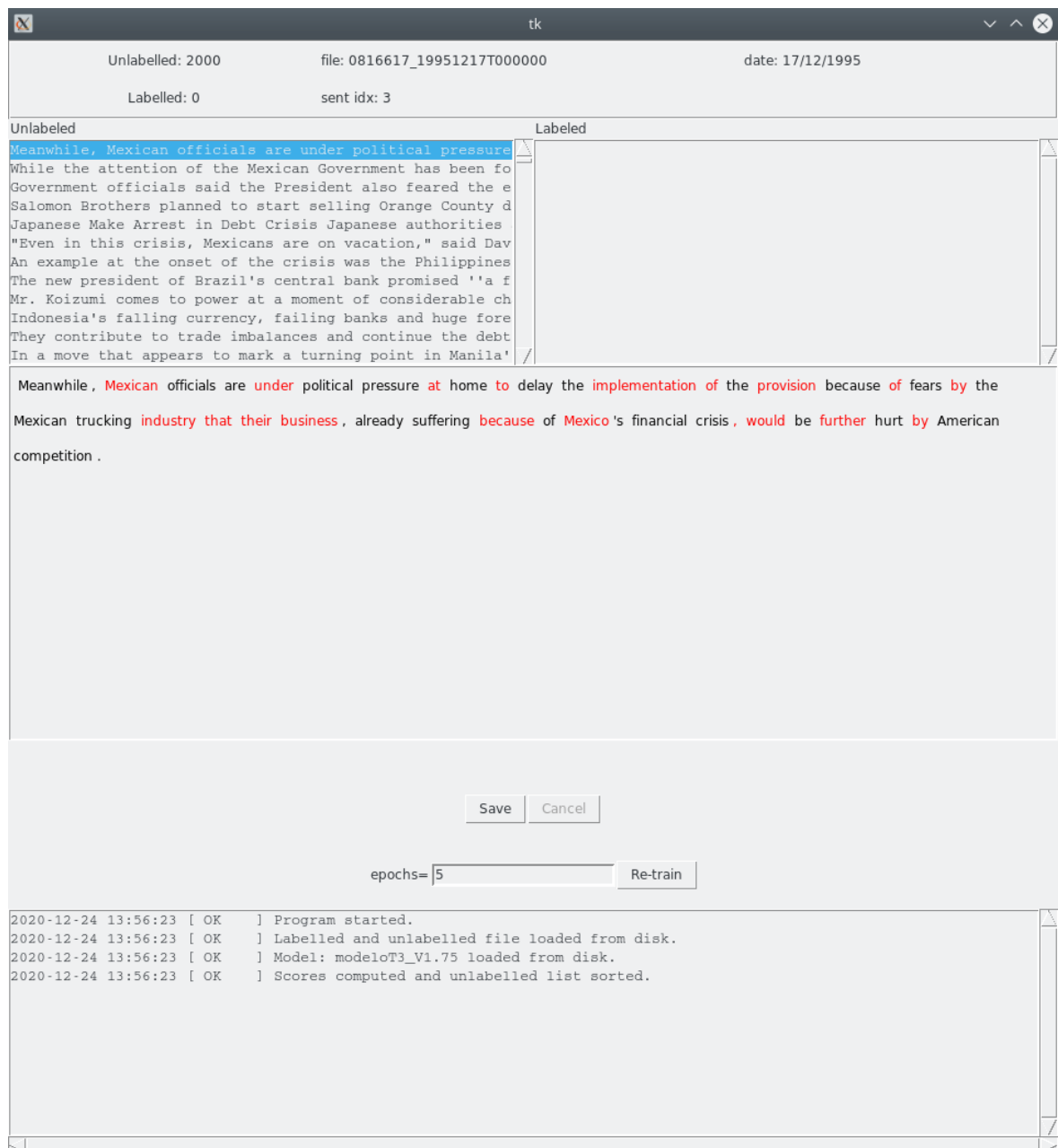
Figure 1

**Panel 3**

Panel 3 is used to display the selected sentence's text, which can be selected from the labeled or unlabeled list. Each word's text will appear in red or black, depending on whether it is an ongoing event trigger or a non ongoing event trigger, respectively. If it is a sentence from the unlabeled list, the red words will correspond to the suggested ongoing event triggers suggested by the tool. If it is a sentence from the labeled list, the red words will correspond to the ongoing event triggers detected by the annotators in a previous step. When the system starts or there is a re-training process, or a sentence is labeled, panel 3 will automatically show the next sentence from the top of the unlabeled list. When any sentence is on display on panel 3, the annotators can click on any word to change from non ongoing event trigger (black) to ongoing event trigger (red) and vice versa.

**Panel 4**

Panel 4 consists of three buttons and a text field for saving changes, undo modifications and starting the re-training process. By pressing the "save" button, the system will save the new set of labels for the current sentence on display on panel 3. If it is a sentence from the unlabeled list, it will be moved to the labeled list. If it is a sentence from the labeled list, the new set of labels will override the existing ones. The "cancel" button will undo all the changes performed on panel 3 (all the clicks produced). If it is a sentence from the unlabeled list, it will restore the suggestions from the system. If it is a sentence from the labeled list, it will restore the previously assigned labels selected by the annotators. By clicking on the "Re-train" button, the system will take the number written on the text field, and it will re-train the system for that number of epochs using the labeled instances to fine-tune the suggestions from that point onwards.

**Panel 5**

Panel 5 outputs the log of the system showing informative messages. This Panel is for debugging purposes, and the annotators will not require to use it during the normal course of labeling.

**Overview of the process of labeling**

The process of labeling will be carried out in several 2-hour sessions. At each session, the annotators will have to choose one sentence at a time for labeling from the unlabeled list. The recommendation is that they choose from the top of the unlabeled list. The annotators should carefully read the selected sentence and agree on the correct set of ongoing event triggers. The suggestions should help the annotators to save time. As the system learns, the number of clicks required to correct the suggested set of labels should go down. Whenever all the annotators agree on the set of ongoing event triggers, they will have to save the current instance with the button "save", this will move the sentence from the unlabeled to the labeled list. At any point, the annotators could go back to an already labeled sentence (from the labeled list) and correct it if they realize they made a mistake. At any point, the annotators could click on the "cancel" button to undo all the changes made on panel 3. The cancel button will bring back the original set of labels (that could be suggestions or labels previously selected by the annotators).

When 50 new instances are labeled, the system should be re-trained during five epochs. Once the system is re-trained, the unlabeled list will be re-sorted and the labeling process will continue with the remaining unlabeled sentences. When the session is over, the annotators should exit the program with the X button at the top left. A dialog will pop-up asking if they want to save and exit, exit without saving, or cancel the exit order. They will have to click on save and exit to store all the changes made during the session. During each session, the system developer will be present to start the system and assist in any problem that may arise. As part of this guideline, we designed a use case of the labeling and re-training process to guide the annotators in these two essential processes. This use case is presented in Annex B.

# Annex A: The Keywords used in Lucene

| Asian Financial Crisis 1997 | Mexican Debt Disclosure Act 1994 | Russian Financial Crisis 1998 | Argentina Financial Crisis 2001 |
|---|---|---|---|
| thai baht | Mexico | Russia | argentine peso |
| indonesia | Crisis | Russian flu | convertibility |
| thailand | Debt | Ruble crisis | fixed exchange rate |
| south korea | capital flight | fixed exchange rate | capital flight |
| Hong Kong | devaluation | fiscal deficit | devaluation |
| Laos | financial crisis | Financial package | IMF |
| Malaysia | emerging markets | International Monetary Fund | default |
| Philippines | cost of borrowing | IMF | fiscal deficit |
| foreign debt | country risk premium | domestic debt | trade deficit |
| Association of Southeast Asian Nations | tequila effect | long term capital management | |
| ASEAN | default | foreign debt | |
| debt-to-GDP | recession | | |
| IMF | Stock exchange | | |
| developing countries | Financial contagion | | |
| capital flows | bailout | | |
| Credit bubbles | balance of payments | | |
| leverage | Mexican bonds | | |
| credit crunch | | | |

Table A1

# Annex B: A Use Case of the System

In this Annex, we present several figures showing a system's use case in which we label one sentence, and then we re-train the system after labeling ten sentences.



Figure B1: Overview of the system starting for the first time. There are 2,000 sentences for labeling, 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 random.

Figure B2: An example of the second step of the use case. In this case, we have chosen the third sentence from the unlabeled list for labeling. The text of these sentences appears on panel 3 with the suggested labels from the system. Due to the lack of training, the suggestions are no good, and several clicks are required to correct the instance. Because we chose the third sentence without labeling any other, we still have an empty labeled list.

Figure B3: In this figure, by clicking on the words in panel 3, we corrected the set of labels initially proposed by the system. However, until we do not press on the "save" button, the label for the instances is not updated; therefore, the labeled list is still empty.

Figure B4: In this figure, the previously corrected instance was saved with the button "save". The annotated instance was moved from the unlabeled list to the labeled list. Afterward, the system automatically loads the next sentence for labeling, which by default is the first of the unlabeled list. This new sentence is presented with the corresponding suggested labels.

Figure B5: This figure illustrates the system's state after five instances have been labeled and the "Re-train" button has been pressed. All further interactions are disabled, and the system starts the re-training process.

Figure B6: This figure illustrates the system's state after the re-training process finishes. When the system ends the re-training process, it is informed in panel 5, and the buttons are enabled again. The unlabeled list is re-sorted according to the newly computed scores of the system. A checkpoint of the model is saved. Lastly, a new sentence is presented to the annotators, which again, by default, is the first sentence of the unlabeled list.