

# CATE: CAusality Tree Extractor from Natural Language Requirements

Noah Jadallah<sup>1</sup>, Jannik Fischbach<sup>2</sup>, Julian Frattini<sup>3</sup>, and Andreas Vogelsang<sup>4</sup>

<sup>1</sup>Technical University of Munich, Germany, noah.jadallah@tum.de

<sup>2</sup>Qualicen GmbH, Germany, jannik.fischbach@qualicen.de

<sup>3</sup>Blekinge Institute of Technology, Sweden, julian.frattini@bth.se

<sup>4</sup>University of Cologne, Germany, vogelsang@cs.uni-koeln.de

**Abstract**—Causal relations (If A, then B) are prevalent in requirements artifacts. Automatically extracting causal relations from requirements holds great potential for various RE activities (e.g., automatic derivation of suitable test cases). However, we lack an approach capable of extracting causal relations from natural language with reasonable performance. In this paper, we present our tool *CATE* (CAusality Tree Extractor), which is able to parse the composition of a causal relation as a tree structure. *CATE* does not only provide an overview of causes and effects in a sentence, but also reveals their semantic coherence by translating the causal relation into a binary tree. We encourage fellow researchers and practitioners to use *CATE* at <https://causalitytreeextractor.com/>

## I. INTRODUCTION

**Motivation:** Black-box behavior of a system is often described by causal relations: *If A and B, then the system shall <functionality>*. Recent studies proved that causal relations occur in traditional requirements documents [1] as well as agile requirement artifacts [2] such as user stories and acceptance criteria. Automatically extracting causal relations from requirements can help to increase the automation of RE activities. For example, we see great potential for the automatic test case derivation from requirements and the automatic detection of dependencies between requirements [3].

**Problem:** We lack an approach capable of extracting causal relations from natural language with reasonable performance. Existing approaches [4] do not consider the combinatorics (e.g., disjunctions) between causes and effects. They also do not allow splitting causes and effects into more granular text fragments (e.g., variable and condition), making the extracted relations unsuitable for the mentioned use cases.

**Contributions:** We present our tool *CATE* (acronym for CAusality Tree Extractor), which is able to parse the composition of a causal relation as a tree structure. *CATE* does not only provide an overview of the causes and effects in a sentence, but also reveals their semantic coherence through the binary tree structure. *CATE* is based on Recursive Neural Networks (RNN) [5], which we trained on the *Causality Treebank* - our self-annotated gold standard corpus of fully labeled binary parse trees representing the composition of 1,571 causal requirements. In this paper, we only present *CATE* as a tool. A detailed description of the *Causality Treebank* and an explanation of why and how we use RNN

for causality extraction is available in our AIRE research paper [6].<sup>1</sup>

## II. CATE: CAUSALITY TREE EXTRACTOR

**Principal Idea** A RNN embraces the idea of natural language as a recursive structure. For example, the syntax of a sentence is recursively structured, with noun phrases containing relative phrases, which in turn contain further noun phrases, and so on [7]. We understand a causal relation also as a recursive structure, since it consists of causes and effects, which in case of conjunctions and disjunctions consist of further causes and effects, and so on [3]. A RNN can be trained to rebuild this recursive structure. Specifically, it recovers the composition of a sentence by identifying related words and merging them into pre-defined segments (e.g., causes and conditions). This results in a binary tree structure (see Fig. 1).

**Left vs. Right-branching** A RNN builds up the binary tree of segments in a bottom-up fashion and supports two different branching methods: left-branching and right-branching. Let us consider the three tokens “set to true”, which constitute a condition segment (see Fig. 1). A RNN can only join adjacent token pairs in each recursion, resulting in two merge options for these tokens. Either “set” and “to” are merged first (left-branching) and then connected with “true” or “to” is first connected with “true” (right-branching) and then joined with “set”. Our first experiments revealed that the RNN seems to learn left-branching better than right-branching [6].

**Word Embeddings** To process causal sentences with a RNN, their words must be first converted into a word embedding. *CATE* is trained on two different types of word embeddings: pre-trained word vectors and randomly initialized word vectors that are trained jointly with the other parameters of the RNN. In the case of pre-trained word vectors, we utilize embeddings established in practice and academia: FastText [8] (1M word vectors), GloVe [9] (2.2M word vectors) and 768-dimensional BERT embeddings [10].

**Training vs. Inference** Emphasis is often put on the training of a model and less on its inference. Since we intend to use *CATE* in practice, we considered two ways of optimizing its inference: ① temperature scaling and ② beam search. We

<sup>1</sup>Please note that *CATE* is currently only based on the vanilla RNN. We are working on also integrating Recursive Neural Tensor Networks described in our research paper into *CATE*.

compare both methods to the naive approach of creating the binary tree, which can be understood as a simple greedy algorithm. Let us assume a sentence consisting of the tokens:  $(a, b, c, d, e)$ . The naive approach to create the binary tree for this sentence consists of the following steps:

- 1) Compute the softmax scores for all adjacent nodes  $\{(a, b), (b, c), (c, d), (d, e)\}$
- 2) Concatenate the two nodes that achieved the highest posterior probability and update the list of adjacent nodes.
- 3) Repeat until only one parent node is left.

The naive approach runs in linear run-time. However, it holds some major drawbacks which can result in poor parsing results, especially for complex sentences. Firstly, the softmax score can not be interpreted reliably as a confidence score since the predicted probabilities tend to be too high even if the input data does not make any sense. Consequently, the probability associated with the predicted class label does not reflect its ground truth correctness likelihood [11]. As a solution, Guo et al. [11] propose to use temperature scaling in order to calibrate the predicted probabilities. Specifically, the neural network is complemented with a parameter  $T$  in the softmax function, with the objective of achieving confidence scores that closer represent the “certainty/uncertainty” of the model. Thus, for a classification into 27 classes, we compute the calibrated softmax score as follows:

$$\text{Softmax}(x_i) = \frac{\exp(x_i)}{\sum_{j=1}^{27} \exp(\frac{x_j}{T})} \quad (1)$$

The second drawback is the inability of the naive approach to correct parsing mistakes. In many cases, the model only detects at a later stage that the underlying sub-tree is incorrect, as the softmax probabilities are decreasing. One solution to determine the optimal tree would be to simply parse all possible trees and compare their calibrated softmax scores. However, this operation is computationally intensive, which is why we identify the best tree by using beam search. Specifically, we sort the parent nodes according to their softmax scores at each stage and cache a predefined number of the most promising nodes (*beam width*).

### III. DEMO PLAN

**User Interface** During the workshop we will use our online demo of *CATE*. The UI consists of two parts. The left side features a text input field, where a causal sentence can be entered. In addition, the configuration of *CATE* can be adjusted. It is possible to set the beam width and to select whether temperature scaling should be used for prediction. Additionally, the user can specify whether the binary tree should be built using left or right branching and which word embeddings should be utilized. After clicking on the `Predict` button, the parsing result is displayed on the right side. The user can interact with the binary tree by expanding or collapsing certain segments.

**Evaluation on unseen data** In the workshop, we will demonstrate that *CATE* is suitable for practical use. For this purpose, we evaluate *CATE* on unseen real word data to simulate its application in the intended context. We use a

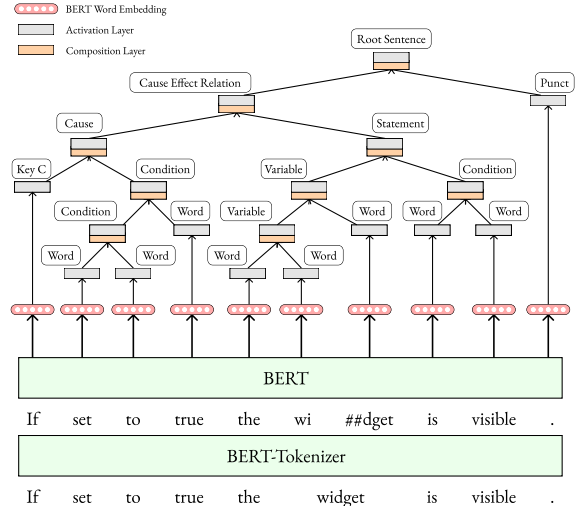


Fig. 1: RNN combined with BERT embeddings.

publicly available data set of acceptance criteria provided by SAP<sup>2</sup>, which describe the functionality of the German “Corona-Warn-App”. The data set consists of 32 user stories, which contain a total of 61 acceptance criteria. We found that 32 of these acceptance criteria exhibit causality [12], making them well suited for the evaluation of our approach. During the workshop, we will present both true and false predictions and discuss the performance, challenges, and possibilities of *CATE* with the other participants.

### REFERENCES

- [1] J. Fischbach, J. Frattini, A. Spaans, M. Kummeth, A. Vogelsang, D. Mendez, and M. Unterkalmsteiner, “Automatic detection of causality in requirement artifacts: the cira approach,” in *REFSQ’21*.
- [2] J. Fischbach, A. Vogelsang, D. Spies, A. Wehrle, M. Junker, and D. Freudenstein, “Specmate: Automated creation of test cases from acceptance criteria,” in *ICST’20*.
- [3] J. Fischbach, B. Hauptmann, L. Konwitschny, D. Spies, and A. Vogelsang, “Towards causality extraction from requirements,” in *RE’20*.
- [4] J. Yang, S. C. Han, and J. Poon, “A survey on extraction of causal relations from natural language text,” *CoRR*, vol. abs/2101.06426, 2021.
- [5] R. Socher, A. Perelygin, J. Wu, J. Chuang, C. D. Manning, A. Ng, and C. Potts, “Recursive deep models for semantic compositionality over a sentiment treebank,” in *EMNLP’13*.
- [6] J. Fischbach, T. Springer, J. Frattini, H. Femmer, and A. Vogelsang, “Fine-grained causality extraction from natural language requirements using recursive neural tensor networks,” in *AIRE’21*.
- [7] R. Socher, C. C.-Y. Lin, A. Y. Ng, and C. D. Manning, “Parsing natural scenes and natural language with recursive neural networks,” ser. *ICML’11*.
- [8] T. Mikolov, E. Grave, P. Bojanowski, C. Puhersch, and A. Joulin, “Advances in pre-training distributed word representations,” in *LREC’18*.
- [9] J. Pennington, R. Socher, and C. D. Manning, “Glove: Global vectors for word representation,” in *EMNLP’14*.
- [10] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” *CoRR*, vol. abs/1810.04805, 2018.
- [11] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, “On calibration of modern neural networks,” *CoRR*, vol. abs/1706.04599, 2017.
- [12] J. Fischbach, J. Frattini, and A. Vogelsang, “Cira: A tool for the automatic detection of causal relationships in requirements artifacts,” in *NLP4RE’21*.

<sup>2</sup>The data set can be found at [https://github.com/corona-warn-app/cwa-documentation/blob/master/scoping\\_document.md](https://github.com/corona-warn-app/cwa-documentation/blob/master/scoping_document.md).

#### IV. SCREENSHOTS OF CATE

**CATE: CAUsality Tree Extractor from Natural Language Requirements**

Your Sentence  
If Event A happens, then Event B shall be triggered.

Beam Width: 5

Temperature Scaling:  (BERT/Left Branching only)

Dataset:

Word Embeddings:

Enter the requirement you want **CATE** to parse. Please note that the requirement must follow a causal pattern.

Configure **CATE** according to your needs:

- Choose the beam width.
- Shall temperature scaling be used for the prediction?
- Shall the tree be built by left or right branching?
- Which word embeddings should be used?

Fig. 2: Configuration options of *CATE*.

**CATE: CAUsality Tree Extractor from Natural Language Requirements**

Your Sentence  
If the system detects an error, a warning window shall be shown.

Beam Width: 8

Temperature Scaling:  (BERT/Left Branching only)

Dataset:

Word Embeddings:

Got result in 11.27s

GoJit 2.1 evaluation  
(c) 1999-2021 Northwoods Software  
Not for distribution or production use  
gaj.net

```

    graph TD
      ROOT_SENTENCE[ROOT_SENTENCE] --- CAUSE_EFFECT_RELATION[CAUSE_EFFECT_RELATION]
      ROOT_SENTENCE --- PUNCT1[PUNCT]
      CAUSE_EFFECT_RELATION --- SEPARATEDCAUSE[SEPARATEDCAUSE]
      CAUSE_EFFECT_RELATION --- STATEMENT1[STATEMENT]
      SEPARATEDCAUSE --- CAUSE[CAUSE]
      SEPARATEDCAUSE --- PUNCT2[PUNCT]
      CAUSE --- KEY_C[KEY_C]
      CAUSE --- STATEMENT2[STATEMENT]
      STATEMENT2 --- VARIABLE1[VARIABLE]
      STATEMENT2 --- CONDITION1[CONDITION]
      VARIABLE1 --- WORD1[WORD]
      VARIABLE1 --- WORD2[WORD]
      CONDITION1 --- CONDITION2[CONDITION]
      CONDITION1 --- WORD3[WORD]
      CONDITION2 --- WORD4[WORD]
      CONDITION2 --- WORD5[WORD]
      KEY_C --- IF[If]
      STATEMENT2 --- THE[the]
      STATEMENT2 --- SYSTEM[system]
      STATEMENT2 --- DETECT[detect]
      STATEMENT2 --- HASH[##s]
      STATEMENT2 --- AN[an]
      STATEMENT2 --- ERROR[error]
      STATEMENT2 --- COMMA[.]
      PUNCT2 --- COMMA2[.]
      STATEMENT1 --- VARIABLE2[VARIABLE]
      STATEMENT1 --- WORD6[WORD]
      STATEMENT1 --- CONDITION3[CONDITION]
      STATEMENT1 --- WORD7[WORD]
      VARIABLE2 --- WARNING[warning]
      VARIABLE2 --- WINDOW>window
      WORD6 --- SHALL>shall
      WORD7 --- BE>be
      WORD8[WORD] --- SHOWN>shown
      PUNCT1 --- END[.]
  
```

Fig. 3: Binary parse of the requirement “If the system detects an error, a warning window shall be shown.”