

Co reference resolution is crucial for matching pronouns to named entities. Introduction of coref149 [10] and senticoref [11], which is comparable to English-based corpora allowed advancement in more complex neural based methods, which require a lot of training data.

## Methods

### 0.1 Corpus

Part of this paper was focused on construction of new corpus containing annotated Slovene and English literary works. Our new corpus consists of 22 literary works, 11 of them are in English and 11 in Slovene. For each work we analyzed basic information such as number of words, language, title, author and release year.

In addition to that we denoted important characters and rank them by importance. If applicable we determined protagonists, antagonists and relations between characters, which were classified as positive, negative or neutral.

We tried to select literary works of differing lengths and from different time periods, as that can effect the used language. Statistics about length can be seen in Table 1. Books range from year 1597 to 1997, with average year being 1874. This is a result of primarily picking fairy tales.

Fairy tales often contain plot twists and deception. For example in Little Red Cap, the Wolf pretends to be Little Red Cap. This provides a challenging task of determining positive and negative relations between characters.

**Table 1.** Corpus statistics

language	number of works	shortest (num. words)	longest (num. words)	avg. number of words
slovenian	11	681	73799	26987
english	11	1068	349736	63592
both	22	681	349736	45289

### 0.2 Entity co-occurrence knowledge graph

To extract usefull information from text we used Classla [9] and Stanza [12] pipelines. Stanza NLP pipeline allows us to perform tokenization, part-of-speech (POS) tagging, lemmatization, dependency parsing as well as named entity recognition (NER). We used it for english language, while we used Classla, a fork of Stanza, which is adapted for processing Slovenian, Croatian, Serbian and Bulgarian language, for processing Slovene text.

Majority of named entities are captured via Classla and Stanza pipelines, but there still exists a problem of detecting coreferences that refer to detected named entities. Coreference resolution for Slovene language was made possible by introduction of coref149 [10] and senticoref [11] dataset, which is comparable to English-based corpora. Large datasets allow models to learn more general patterns. Better generalization of models trained on larger datasets was shown in Klemen and Žitnik [13], where they used BERT [3] contextual embedding as an input to the neural based coreference scorer. This model

was used in our paper for coreference resolution of Slovene text.

Extracted named entities and their coreferences are often ambiguous. When it comes to people, first and last names are often only used the first time the entity occurs and are in the future referred by only the first or last name. Organizations are also often referred to by only a part of the full name or only their initials [14].

To combat in-text ambiguities we rely on string matching using Levenshtein distances between strings. We adapted the deduplication method from FuzzyWuzzy library[15] to be better suited for named entity deduplication. We improved the selection of the canonical example in which all the similar entities are merged, by accounting for number of occurrences in the text, references that only use first and last name and organizations, which are referred by only their initials. As a result we get proper entity names as well as number of occurrences of entity in text.

Deduplication using Levenshtein distance matching is slow when it comes to large collections of items. We mitigate this issue by reducing search space with simstring [16] method, which is approximate dictionary matching using different similarity measures (Dice, Jaccard and cosine). We use 2-grams as a feature and cosine distance as similarity measure. Entities that have similarity higher than threshold are then used with the FuzzyWuzzy deduplication that is explained earlier. Deduplication process significantly reduces the number of nodes referring to the same entity and thus improves performance of the algorithm.

After deduplication process we use dependency parsed text as to generate subject-verb-object (SVO) triplets. This is done using pattern-based framework for language-neutral predicate-argument extraction patterns presented in White et al. [17]. Evaluation done by Zhang et al. [18] showed state of the art performance compared to other Open Information Extraction (Open IE) tools.

SVO triplets can be viewed as an event connection two entities. After correcting entity names using deduplication from previous section and resolving co references we get named entity pairs, which can be used to construct a entity co-occurrence knowledge graph (ECKG).

ECKG is than used to perform entity relation analysis. From co-occurrences we can gather interactions between different entities, which can be seen as character, when it comes to literary text.

Another way to determine relation between two entities is their location in the text. We construct another ECKG, which determines relation based on co-occurrence in the same

sentence.

We use two groups of approaches of determining entity importance:

1. Importance based on number of occurrences of entity in text
2. Importance based on node importance in the ECKG.

Node importance is determined by different centrality measures performed on ECKG. We evaluated performance of different centralities, which can be seen in the results section 0.4.5.

In addition to character detection and importance evaluation we want to determine protagonists and antagonists. Determining positive character can be done by analyzing the difference between its positively and negatively weighted edges, where positive weight represents positive relation between two characters and vice versa. Large difference between negative and positive edges determines positive characters, while small difference can determine negative characters. For determining positive and negative relations between entities we perform sentiment analysis on SVO triplets and sentences containing named entities.

### 0.3 Relationship extraction

Our next task was to extract semantic relationships from novels in order to obtain information that would allow us to grow some knowledge. Extracted relationships usually occur between two or more entities of a certain type and fall into a number of semantic categories. But to gain knowledge that makes sense, the removal of repeated relations (disambiguation) and generally refers to the extraction of many relationships is required.

On extracted entities, we performed relationship extraction using two different approaches, rule-based and model-based. Usually, good precision values are achieved with Rule-based approaches. Using partial syntactic parsing can simplify the linguistic structures containing instances of semantic relations.

Relation statements can be represented as:

$$r^i = (x^i, s_1^i, s_2^i), \quad (1)$$

$x$  is the tokenized sentence,  $s_1$  and  $s_2$  are the spans of the two entities within that sentence. Two statements can consist of two different sentences, but they can both contain the same entity pair. If both contain the same entity pair, they should have the same  $s_1$ - $s_2$  relation.

#### 0.3.1 Rule-based

In this approach, different patterns were used to extract relationships between entities [19]. We used spaCy framework [20], an open-source software library for advanced natural language processing, that features NER, POS tagging, dependency parsing, word vectors and more. After that, we

constructed a knowledge graph, where nodes represent entities and links represent relationships.

Defining the right set of patterns is the most difficult part. Because the approach is manual, it often results in large sets of noisy patterns. There must be some sort of semantic categorization of relationships between the entities. Taking the verb, which has the time component, may not be a good way to have build a knowledge graph. There could be synonyms of verbs that have different relationship types with similar, identical semantic values.

#### 0.3.2 Model-based

To avoid noisy patterns, NLP model can be used for event extraction. In this step, we used an OpenNRE [21], open-source unified framework for relation extraction models. This framework provides different pre-trained models, dataset on which the models were trained and tools for performing additional learning of the models on custom data. GPU parallel computing is also supported. We tried CNN and BERT encoder models.

Convolutional neural network models are commonly used in information extraction tasks. But CNN just considers the correlation between consecutive words and ignores the correlation between discontinuous words. CNNs help us with parallelization, local dependencies and distance between positions. OpenNRE CNN model is pre-trained on WIKI80 dataset.

Bidirectional Encoder Representations from Transformers is a transformer-based machine learning technique. Transformers were introduced to deal with the long-range dependency challenge, which cannot be extracted using CNNs. They provide a self-attention mechanism, that allows models to look at other words in the input sequence to get a better understanding of a certain word in the sequence.

#### 0.4 Event causality identification

For event causality identification we implemented the method introduced by Jian Liu et. al. [22]. They propose a mixture model using a knowledge aware reasoner and a mention masking reasoner, which can leverage both external knowledge to enrich the representation of events as well as learn event-agnostic, context specific patterns which grants the model a decent ability to generalize on previously unseen data.

##### 0.4.1 Knowledge aware reasoner

Given a pair of events  $e_1$  and  $e_2$  the knowledge aware reasoner first retrieves the related knowledge in ConceptNet [23] and then encodes the knowledge into contexts for reasoning. We only consider 18 semantic relations that are potentially useful for event causality identification: CapableOf, IsA, HasProperty, Causes, MannerOf, Causes-Desire, UsedFor, HasSubevent, HasPrerequisite, NotDesires, PartOf, HasA, Entails, ReceivesAction, UsedFor, CreatedBy, MadeOf, and Desires. We also limit the total knowledge retrieved.

##### 0.4.2 Mention masking reasoner

With the mention masking reasoner we aim to explore event agnostic, context-specific patterns for reasoning. We replace  $e_1$  and  $e_2$  with a '[MASK]' token to exclude event information. Then a BERT encoder is used to obtain embedded representations of events  $F_{MASK}^{(e_1, e_2)}$ .

By using the mention masking reasoner, we force our model to predict whether  $e_1$  and  $e_2$  form a casual relation based on context specific clues as the masked representation does not contain any event-specific learning.

$$L = -\delta_{A,B} * \log(p(l=1|A, B)) + (1 - \delta_{A,B}) * \log(1 - p(l=1|A, B)), \quad (2)$$

where  $\delta_{A,B}$  is the Kronecker delta which takes the value 1 when both A and B express a causal relation and 0 otherwise.  $p(l=1|A, B) = \frac{1}{1 + \exp(F_{MASK}^A \cdot F_{MASK}^B)}$  defines the distributional similarity score.

##### 0.4.3 The attentive sentinel

With the attentive sentinel we aim to learn a trade-off between the knowledge aware reasoner and the mention masking reasoner, by learning an attentive gate as their combination of weights:

$$g_{e_1, e_2} = \sigma(W(F_{KG}^{(e_1, e_2)} \oplus F_{MASK}^{(e_1, e_2)}) + b), \quad (3)$$

where  $W$  and  $b$  are model parameters,  $\oplus$  is the concatenation operator. We then adopt a weighted summation to integrate  $F_{KG}^{(e_1, e_2)}$  and  $F_{MASK}^{(e_1, e_2)}$  as the final feature:

$$F_{e_1, e_2} = g_{e_1, e_2} * F_{KG}^{(e_1, e_2)} + (1 - g_{e_1, e_2}) * F_{MASK}^{(e_1, e_2)}. \quad (4)$$

The attentive sentinel balances the knowledge aware reasoner and the mention masking reasoner to make the final prediction.

##### 0.4.4 Model prediction and training

To make the final prediction, we perform a binary classification by taking  $F_{e_1, e_2}$  as input:

$$o_{e_1, e_2} = \sigma(W_o F_{e_1, e_2} + b_o), \quad (5)$$

where  $o_{e_1, e_2}$  is the probability of there being a causal relationship between the two events;  $w_o$  and  $b_o$  are model parameters. For training we use cross-entropy as the loss function:

$$J(\Theta) = - \sum_s \sum_{e_i, e_j \in E_s, e_i \neq e_j} y_{e_i, e_j} \log(o_{e_i, e_j}) + (1 - y_{e_i, e_j}) \log(1 - o_{e_i, e_j}), \quad (6)$$

where  $\Theta$  denotes the parameter set of our model;  $s$  ranges over each sentence;  $e_i$  and  $e_j$  range over each event in  $s$ . We used the Adam [24] algorithm to optimize model parameters.

##### 0.4.5 Datasets used for training and evaluation

For training and evaluating the model we used the Causal-TimeBank dataset [25], which contains annotated causal relations of events within a single sentence.

Because of the lack of annotated datasets on causal event relations extracted from literary works, we used the model trained on the Causal-TimeBank dataset. As shown in Liu et. al. [22] this method generalizes very well and is suitable for transfer learning and cross-task adaptation.

We then manually evaluated the performance of the model on the Litbank event dataset [26], which contains 100 annotated works of English-language fiction.

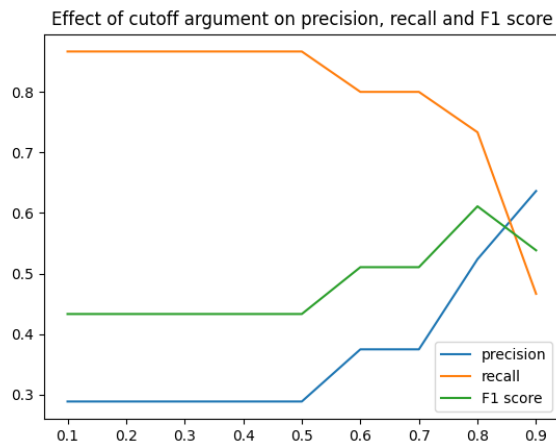
## Results

We evaluated three approaches presented in 0.2. For better readability we will use following names through this chapter:

- Importance based on number of occurrences in the text will be called method 1
- Importance based on graph importance measures for graph constructed by denoting co-occurrences in the same sentence as relation, will be called method 2
- Importance based on graph importance measures for graph constructed by denoting co-occurrences in the SVO triplet as relation, will be called method 3

### 0.5 Character recognition

First we evaluate character recognition, which is evaluated using recall, precision and F1 value. The most important value in this case is F1, because precision and recall are inversely correlated. They can be manipulated by cutoff argument, which determines entity importance score at which it is not considered a character any more. This effect can be seen on figure 1



**Figure 1.** Effect of cutoff argument on precision, recall and F1 score (evaluated on one literary work with method 1)

We compared all three methods on 5 different books. Cutoff argument was set to 0.9 and methods 2 and 3 were using a combination of eigenvector, degree and closeness centrality for importance calculation. Average result can be seen in Table ???. To better evaluate methods best general cutoff should be determined for each method separately. Importance of deciding on how many characters to keep can be seen in Table ??, where there is big difference between mean average precision when ranking importance and F1 score. This shows that too many or too little entities were kept.

	method 1	method 2	method 3
precision	0.277	0.327	0.440
recall	0.918	0.786	0.490
F1 score	0.392	0.412	0.353
mAP	0.850	0.778	0.382

**Table 2.** Average score of evaluation of methods on 5 Slovene books

### 0.6 Character ranking

Our corpus contains characters, which are ranked by their importance. All three methods can be in addition to character recognition used to determine ranking of importance for recognized characters. To evaluate ranking we will be using mean average precision (mAP) metric. Results of mAP average over 5 different books can be seen in Table ??.

### 0.7 Protagonist/antagonist detection

#### Discussion

Use the Discussion section to objectively evaluate your work, do not just put praise on everything you did, be critical and exposes flaws and weaknesses of your solution. You can also explain what you would do differently if you would be able to start again and what upgrades could be done on the project in the future.

#### Acknowledgments

Here you can thank other persons (advisors, colleagues ...) that contributed to the successful completion of your project.

#### References

- [1] Paramita Mirza and Sara Tonelli. An analysis of causality between events and its relation to temporal information. In *COLING*, 2014.
- [2] Tommaso Caselli and P. Vossen. The event storyline corpus: A new benchmark for causal and temporal relation extraction. In *NEWS@ACL*, 2017.
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019.
- [4] Pedram Hosseini, David A. Broniatowski, and Mona T. Diab. Predicting directionality in causal relations in text. *ArXiv*, abs/2103.13606, 2021.
- [5] Lei Gao, Prafulla Kumar Choubey, and Ruihong Huang. Modeling document-level causal structures for event causal relation identification. In *NAACL*, 2019.
- [6] Jannik Strötgen and Michael Gertz. A baseline temporal tagger for all languages. In *EMNLP*, 2015.
- [7] Sam Wei, Igor Korostil, Joel Nothman, and Ben Hachey. English event detection with translated language features.

- In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 293–298, Vancouver, Canada, July 2017. Association for Computational Linguistics.
- [8] Simon Krek, Kaja Dobrovoljc, Tomaž Erjavec, Sara Može, Nina Ledinek, Nanika Holz, Katja Zupan, Polona Gantar, Taja Kuzman, Jaka Čibej, Špela Arhar Holdt, Teja Kavčič, Iza Škrjanec, Dafne Marko, Lucija Jezeršek, and Anja Zajc. Training corpus ssj500k 2.2, 2019. Slovenian language resource repository CLARIN.SI.
  - [9] Nikola Ljubešić and Kaja Dobrovoljc. What does neural bring? analysing improvements in morphosyntactic annotation and lemmatisation of Slovenian, Croatian and Serbian. In *Proceedings of the 7th Workshop on Balto-Slavic Natural Language Processing*, pages 29–34, Florence, Italy, August 2019. Association for Computational Linguistics.
  - [10] Slavko Žitnik. Slovene coreference resolution corpus coref149, 2018. Slovenian language resource repository CLARIN.SI.
  - [11] Slavko Žitnik. Slovene corpus for aspect-based sentiment analysis - SentiCoref 1.0, 2019. Slovenian language resource repository CLARIN.SI.
  - [12] Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. Stanza: A Python natural language processing toolkit for many human languages. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, 2020.
  - [13] Matej Klemen and Slavko Žitnik. Neural coreference resolution for slovene language. *Computer Science and Information Systems*, 2021.
  - [14] Mónica Marrero, Julián Urbano, Sonia Sánchez-Cuadrado, Jorge Morato, and Juan Miguel Gómez-Berbís. Named entity recognition: Fallacies, challenges and opportunities. *Computer Standards Interfaces*, 35(5):482–489, 2013.
  - [15] Adam Cohen. Fuzzywuzzy library. Accessible: <https://github.com/seatgeek/fuzzywuzzy>. [Accessed: 28. 3. 2022].
  - [16] Naoaki Okazaki and Jun’ichi Tsujii. Simple and efficient algorithm for approximate dictionary matching. In *Proceedings of the 23rd International Conference on Computational Linguistics (Coling 2010)*, pages 851–859, Beijing, China, August 2010.
  - [17] Aaron Steven White, Drew Reisinger, Keisuke Sakaguchi, Tim Vieira, Sheng Zhang, Rachel Rudinger, Kyle Rawlins, and Benjamin Van Durme. Universal Decompositional Semantics on Universal Dependencies. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1713–1723, Austin, Texas, November 2016. Association for Computational Linguistics.
  - [18] Sheng Zhang, Rachel Rudinger, and Ben Van Durme. An Evaluation of PredPatt and Open IE via Stage 1 Semantic Role Labeling. In *Proceedings of the 12th International Conference on Computational Semantics (IWCS)*, Montpellier, France, September 2017.
  - [19] Kang Liu, Yubo Chen, Jian Liu, Xinyu Zuo, and Jun Zhao. Extracting events and their relations from texts: A survey on recent research progress and challenges. *AI Open*, 1:22–39, 2020.
  - [20] Industrial-strength natural language processing. <https://spacy.io>. Accessed: 2022-05-01.
  - [21] Xu Han, Tianyu Gao, Yuan Yao, Deming Ye, Zhiyuan Liu, and Maosong Sun. OpenNRE: An open and extensible toolkit for neural relation extraction. In *Proceedings of EMNLP-IJCNLP: System Demonstrations*, pages 169–174, 2019.
  - [22] Jian Liu, Yubo Chen, and Jun Zhao. Knowledge enhanced event causality identification with mention masking generalizations. In *Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence*, pages 3608–3614, 2021.
  - [23] Robyn Speer, Joshua Chin, and Catherine Havasi. Conceptnet 5.5: An open multilingual graph of general knowledge. In *Thirty-first AAAI conference on artificial intelligence*, 2017.
  - [24] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
  - [25] Paramita Mirza, Rachele Sprugnoli, Sara Tonelli, and Manuela Speranza. Annotating causality in the tempeval-3 corpus. In *Proceedings of the EACL 2014 Workshop on Computational Approaches to Causality in Language (CAtoCL)*, pages 10–19, 2014.
  - [26] Matthew Sims, Jong Ho Park, and David Bamman. Literary event detection. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3623–3634, 2019.