# Event Causality Extraction
# from Natural Science Literature

Biswanath Barik, Erwin Marsi, Pinar Öztürk

Department of Computer and Information Science
Norwegian University of Science and Technology
{biswanath.barik,emarsi,pinar}@idi.ntnu.no

**Abstract.** We aim to develop a text mining framework capable of identifying and extracting *causal dependencies* among *changing variables* (or *events*) from scientific publications in the cross-disciplinary field of oceanographic climate science. The extracted information can be used to infer new knowledge or to find out unknown hypotheses through reasoning, which forms the basis of a knowledge discovery support system. Automatic extraction of causal knowledge from text content is a challenging task. Generally, the approaches of causal relation identification proposed in the literature target specific domain such as online news or biomedicine as the domain has significant influence on causality expressions found in the domain texts. Therefore, the existing models of causality extraction may not be directly portable to other/new domains. In this paper, we describe the nature of causation observed in climate science domain, review the state-of-the-art approaches in causal knowledge extraction from text and carefully select the methods and resources most likely to be applicable to the considered domain.

**Keywords:** causal relation, information extraction, relation extraction, knowledge discovery

## 1 Introduction

Climate change is a major concern in recent years and has various direct and indirect effects in day-to-day life. Global warming, intensified with various anthropogenic activities, have serious impacts on the precious climate system of this blue planet. As a result, the pattern of the climate has been changing rapidly. The other associated natural processes and systems are, in consequence with climate change, being affected significantly. Changes in the ecological system of the marine environment is one such system.

Various studies in the domain of marine science, climate science, environmental science and other related fields of Earth Science report significant *changes* in recent years in several *parameters* (i.e., quantitative variables) of the ocean environment. *Sea-surface temperature*, *bottom-water temperature*, *direction of ocean circulation*, *acidity*, *pH*, *alkalinity* and *$CO_2$ concentration level in water* are a few examples of variables observed to change significantly in past years.

Significant changes of such variables have impact on *phytoplankton growth rate* which indirectly affects marine *food web* - a complex feeding network of different species (i.e., who eats whom) living in the marine environment. Along with this, the efficiency of the *biological pump*, the ocean's ability to absorb and store $CO_2$ as the food web transfer parts of the biomass deeper into the ocean, is degraded.

The research publications in climate science, marine science and environmental science are the authentic sources of information describing various theories and models consisting of changing variables (or *events*) and their complex interactions. The elementary interactions among the events may be in the form of correlations, causal relations or the positive/negative feedback cycles consisting of sequence of events. Identifying and extracting valuable interactions from the scientific articles and combining them to explore various hidden connections among the events can help to better understand the functionality of various processes of the domain and their dependencies. However, the sheer volume of the articles limits scientists and policy makers to collect useful information by reading the articles in due time. Human cognition, on the other hand, may be another limit for recognizing and interpreting various cross-domain knowledge fragments. Therefore, an automated knowledge discovery support system is needed to quickly process the vast collection of research articles, extract useful knowledge fragments and produce new insights, hypothesis or discover unknown knowledge by combining the extracted knowledge units.

To distill essential factual knowledge from unstructured text content of research papers, the text mining techniques are successfully applied in the domain of bio-medicine. Significant advancements have been observed in identifying named entities [9], detecting events [28], coreference resolution [1] and causal relations extraction [26]. With the use of domain dependent Natural Language Processing(NLP) tools like Part-Of-Speech (POS) tagger, shallow and full sentence parsers for syntactic analysis of the text content, biomedical text mining is capable of providing a platform where researchers can query on the vast database of research articles of the domain. The researchers don't need to bother how many papers the system needs to process to find or infer the required answer, or how did the system do it?

In the domain of Natural Science, specifically in the cross-disciplinary domain of oceanographic climate science, our goal is to develop a (literature-based) knowledge discovery support system to facilitate a large community of scientists and researchers. The major challenges we face are the lack of resources like task specific annotated corpora, indexed literature databases covering the entire field, domain dependent NLP tools with good accuracies and knowledge resources (ontologies) as our target domain is almost unexplored. The resources and tools developed in biomedicine domain are not directly usable due to domain difference as shown in [24, 25]. Therefore, to meet this goal, a constant effort is being employed to develop resources and tools. In [25], authors describe an annotation scheme to annotate quantitative variables, their change events, correlations and causal relations among change events, and feedback loops from the abstracts and full-text journal papers collected from the nature publication. In [24], authors

automatically identify and extract variables and their direction of changes using a tree pattern matching technique and generalise these variables by progressive pruning of syntax tree using tree transformation operations.

In this research direction, our target is to develop a *causality extraction model* in oceanographic climate science domain where the causal relation among change events (as described in Section 2) can be automatically extracted from the scientific publications. The causal relations extracted from a collection of research papers of the domain can then be used for causal reasoning with the help of domain knowledge to discover new facts or unknown hypotheses of the domain. Such a reasoning system can provide a better information search (semantic search) capability to the scientists and researchers to efficiently access vast database of publications. In this paper, we explore the existing methods and algorithms of event causality identification and carefully select the methods and resources most likely to be applicable to oceanographic climate science domain with a proposed work plan.

The paper is organized in the following way. Section 2 describes the causality in oceanographic climate science domain. Section 3 shows the nature of causation observed in the considered domain. Section 4 describes the existing approaches for handling causality in text content. Section 5 discusses the suitability of the existing algorithms in this domain and the proposed work plan. Section 6 concludes the paper.

## 2 Problem Description

Causality extraction from text content is a fundamental task towards the desire of developing literature-based knowledge discovery support system. In climate science domain, an *event* is defines as: "*a change is an event in which the value of a quantitative variable is changing*" [25]. Causal relation, in general, is a semantic relation between two events where the occurrence of one event (called *cause event*) causes the occurrence of the other event (referred as the *effect event*). Figure 1 shows a typical example of causal relation between two events $E_1$ and $E_2$, where $E_1$ is "Reduced calcification of marine plankton" and $E_2$ is "increased atmospheric $CO_2$."

The causality between these two events is expressed explicitly by the causality marker (or cue phrase) "in response to". All the examples (Fig. 1 - Fig. 4) of causal relation in climate science domain are taken from the pilot annotation described in [25].

We have mentioned earlier that climate science is a new text mining domain. The necessary resources and domain specific NLP tools are not available in order to immediately and effortlessly build the pipeline for analyzing larger context and to just focus only on causal relation extraction module. Keeping in mind this limitation, we are interested to develop our causal relation extraction model incrementally. In the first step, we focus on identifying intra-sentence explicit causal relations. Some issues related to identifying causal relations within a single sentence where causation is explicitly expressed are described in the following
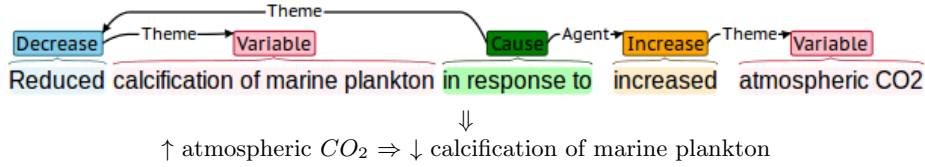
$$\Downarrow$$

$$\uparrow \text{atmospheric } CO_2 \Rightarrow \downarrow \text{calcification of marine plankton}$$

**Fig. 1.** Causality between two events in Natural Science literature

section. We have a plan to deal with inter-sentence causal relations by identifying discourse causal markers and resolving coreference issues, but this is not a focus in this paper.

## 3 Nature of Causation in Climate Science

We have studied various diversities and issues related to the notion of causation and how it may be expressed in natural language specifically in the research articles of climate science collected from Nature publications. These diversities need to be discussed in detail in order to develop a causal relation extraction model in this domain. The issues are described below:

1. **Multi-Event Participation in Causation:** Causal relation is, in general, a binary relation. It is a relation between two events: the cause event and the effect event as described in figure 1. However, it is observed that in causal implication, events can participate in the antecedent and/or in the consequence part as shown in figure 2 below.
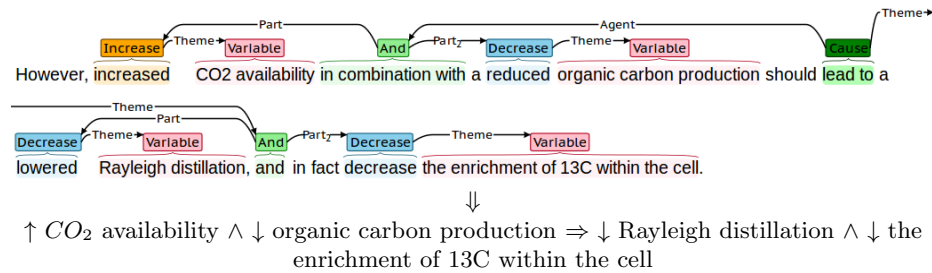


$$\Downarrow$$

$$\uparrow CO_2 \text{ availability} \wedge \downarrow \text{organic carbon production} \Rightarrow \downarrow \text{Rayleigh distillation} \wedge \downarrow \text{the enrichment of 13C within the cell}$$

**Fig. 2.** Many-to-Many events causality

In this example, two events "increased $CO_2$ availability" and "reduced organic carbon production" are connected through the conjunction "in combination with" to form a larger (or composite) event and serve as a cause event to the causal relation signaled by "lead to". Similarly, the events "lowered Rayleigh distillation" and "decrease the enrichment of 13C within the cell" jointly forms the effect event. Therefore, causal relation in this example is a Many-to-Many(M:M) relation.

2. **Event Participating in Multiple-Causation:** Causal relations are sometimes expressed in a cascaded style where a single event can participate in more than one causal expressions in the same sentence. Figure 3 shows an example of such causation.
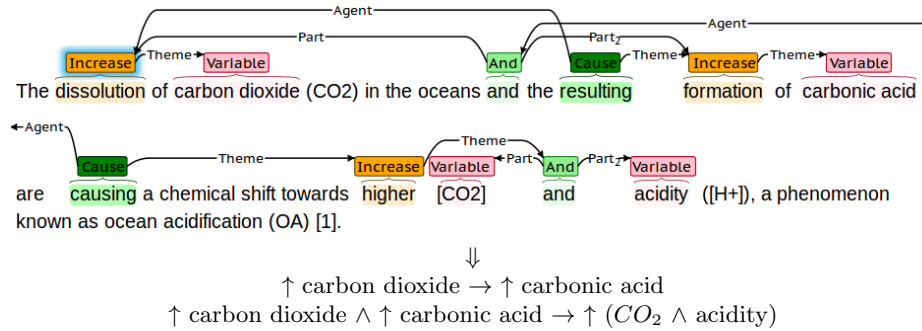


$$\Downarrow$$
$$\uparrow \text{carbon dioxide} \rightarrow \uparrow \text{carbonic acid}$$
$$\uparrow \text{carbon dioxide} \wedge \uparrow \text{carbonic acid} \rightarrow \uparrow (CO_2 \wedge \text{acidity})$$

**Fig. 3.** Many-to-many event causality

In this example, the first causal relation expresses an 1:1 relation between events "↑ carbon dioxide" (cause event) and "↑ carbonic acid" (effect event). The second causal relation in the same sentence represented as a M:M event relation where the events of the first relation collectively constitute the cause event and the events "$\uparrow CO_2$" and "↑ acidity" are the effect events.

3. **Double Role of Causal Markers:** It is observed in other domains that the causal markers are often ambiguous i.e., they do not always express causality. Under certain context and semantic orientation, the markers express causality between events. This phenomenon is also true in climate science domain. However, we have noticed that the causal markers also *trigger* the change events along with its causation indication role in certain contexts.
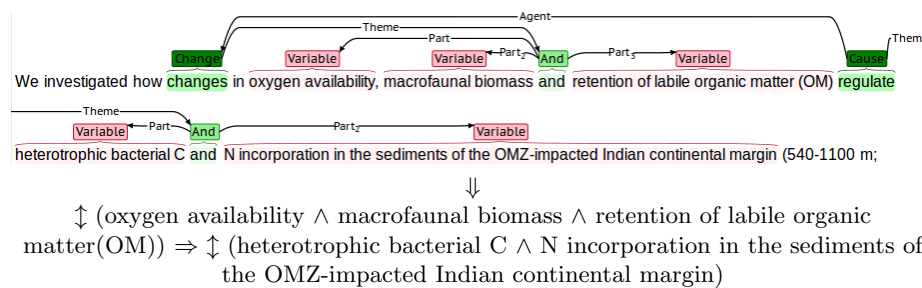


$$\Updownarrow$$
$$\updownarrow (\text{oxygen availability} \wedge \text{macrofaunal biomass} \wedge \text{retention of labile organic}$$
$$\text{matter(OM)}) \Rightarrow \updownarrow (\text{heterotrophic bacterial C} \wedge \text{N incorporation in the sediments of}$$
$$\text{the OMZ-impacted Indian continental margin})$$

**Fig. 4.** Causality marker holds implicit trigger to event

In figure 4., the word *regulate* serves as a causal marker. Also, it serves as a *change* indicator (denoted as "↕") to the variables "heterotrophic bacterial C" and "N incorporation in the sediments of the OMZ-impacted Indian continental margin".

## 4 Causal Relation Extraction - A Brief Review

The general notion of causality is very broad. It has been studied in various fields of research like philosophy, statistics, cognitive science, linguistics, physics, economics, biology, medicine and so on. In computational linguistics, considerable amount of work has been done on automatic extraction of causal knowledge from text in general [5, 8, 22, 32] and in specialized domains like biomedical science [20, 26], or online news domain [11, 31]. However, causality knowledge extraction is a non-trivial problem till date. Many questions remain unsolved about the nature of causation. Also, causation is subjective - human judgement about causation is even conflicting in many instances (shows low inter-annotators agreement) and subject to the realization of the context [14]. In the following sub-sections we broadly categorize existing approaches of causality identification from text content found in literature.

### 4.1 Causality Extraction using Handcrafted Patterns

The initial attempts of causal relation extraction rely on knowledge-based inference techniques [18,22,35]. These works used linguistic patterns of causation along with manually crafted resources to detect causal relation hidden the context. In this research direction, Kaplan *et* al. [19] proposed a linguistic pattern-based approach for causal knowledge extraction where the resources like grammar, lexicon and domain-knowledge are hand-crafted for the target domain. Garcia [10] develop an approach where the causative verb patterns are extracted from *French* texts using handcrafted rules. In this experiment, the author found 25 causal relations and classified them with a precision of 85% using a semantic model based on "Force Dynamics" of Leonard Talmy [37]. Explicit causal relations are also identified from MEDLINE text database by Khoo *et* al. [20] using predefined linguistic patterns and achieved a precision about 68%. In this work, partially parsed verb linguistic patterns indicating causality relationships are matched on text to extract cause-effect information.

The causal relation extraction models based on linguistic patterns perform pretty well in restricted domains. However, rule creation is expensive and time consuming and it suffers from domain portability issue.

### 4.2 Semi-automatic Causal Pattern Learning

The other research direction explores semi-automatic learning of causal patterns from corpus with minimal (or no) domain knowledge. In this direction, Khoo *et* al. [21] developed an automatic system for extracting cause-effect relation

from newspaper texts using simple pattern-matching and without using linguistic clues and domain knowledge. In [12], Girju and Moldovan describe syntactic and semantic classification of cause-effect lexico-syntactic patterns found in English texts. They developed an approach to automatically identify lexico-syntactic patterns consisting of a pair of noun phrases connected by causative verbs ($<NP_1$ verb $NP_2>$) that express the causal relations. Finally, a semi-automatic validation method is proposed to evaluate the extracted causal patterns. Marcu and Echihabi [23] classifies a sentence pair as 'causal' or '¬causal' by training a Naïve Bayes classifier on inter-sentence lexical pair probability. Girju [11] developed a decision tree based classifier on causality-annotated corpora, where the cue phrases are automatically extracted from WordNet [27] and also from the corpus, and achieved a precision of 73.91%. In [5], Blanco *et* al. first manually identify the syntactic patterns that may encode marked and explicit causation and found that the four most common relators encoding causation are *because, since, as* and *after*. Then they used decision tree based learning algorithm (an implementation of Bagging with C4.5 decision trees) to decide whether or not a pattern instance encodes a causation. However, this method is not able to detect the causes and the effects. Ittoo and Bouma [17] present a semi-supervised method for automatic extraction of high quality causal relations from domain-specific, sparse corpora. In this work, they initially acquire a set of explicit and implicit lexico-syntactic patterns from Wikipedia. Using some *seed* cause-effect patterns, the extracted patterns are then classified as causal or non-causal by measuring their reliability through computing point-wise mutual information between extracted patterns and seed patterns and ranking the extracted patterns accordingly. Finally, the extracted causal patterns are used to identify domain-specific causal relations.

### 4.3   Causality Prediction by Supervised Learning

In the supervised learning set-up, the domain corpus is needed to be annotated with events (or entities) and their causal relationships. Causality annotated corpora is then used to train supervised model for classifying a pair of events as causal or non-causal pair. In [11], the author manually annotates *Loss Angles Times corpus* based on explicit causal verbs (e.g., "to cause"). Using this annotated corpora and WordNet [27], the cue phrases are extracted, automatically. A decision tree classifier is then trained which detects the causality relation in news events with 73.9% precision and 88.7% recall. Beamer *et* al. [2] develop a support vector machine (SVM) based classifier trained on SemEval 2007 Task 4 corpus [13] and report an accuracy of 77.5% in identifying cause-effect noun pairs. SVM classifiers are also trained on annotated *Wall Street Journal (WSJ)* texts by [4,34]. In [29], the verb-pair rules are used to train Naïve Bayes (NB) and SVM classifier to identify causality from multiple Elementary Discourse Units (EDUs) and reported precision of 88% with NB and 89% with SVM.

## 4.4 Statistical Approaches

Existing corpus-based approaches to causality extraction use distributional characteristics of events, like co-occurrence features, object-sharing features, temporal features, distance features and so on. Machine learning based approaches are more robust than rule-based approaches and require less linguistic information and domain knowledge. Torisawa [38] developed a model for extracting commonsense inference rules from coordinate verb phrases based on co-occurrence and object sharing features. This work is further extended in [39] where the occurrence frequency of a single verb is emphasized and reported a relative improvement of 60 precision. Other approaches use predicate semantics [15, 16] and shared arguments [6, 7].

Since causation can be expressed in many different ways in natural language, the automatic recognition of causal relations is challenging. In recent works, researchers try to overcome this challenge by considering specific constructions like causation between verbs [2, 4, 8, 32], between verb-noun pairs [8, 33] and between two discourse segments [30, 36].

## 5 Proposed Work Plan

The approach of causal relation extraction using handcrafted rules/patterns is not suitable in our domain as our domain is cross-disciplinary in nature and creating causal patterns requires sufficient expertise of the domains. Also, in the rule-based causality identification approach, the rules created manually work well when the causality is obvious i.e., there has no ambiguity in identifying the causal indicators and the participating events. However, in Section 3, we have seen that causation in the considered domain is often complex in nature. Therefore, hand-made pattern-based causality detection is not feasible in this domain.

The supervised learning approach of causal relation is also not applicable as it requires large amount of causality annotated corpora of the domain which is a rather costly process. However, as we showed in Section 3 that an event can participate in more than one causal relation in some contexts, supervised classification approaches can be good candidates where the lexical, contextual, syntactic and semantic features can be well exploited for classifying a pair of events. We are under process of developing a small amount of causality annotated corpora for developing a baseline causal relation extraction model in our domain using supervised learning method like SVM. In the next step, we will experiment with various semi-supervised algorithms to improve the baseline performance using the annotated corpus and with large collection of unannotated texts.

The unsupervised way of measuring causal association between an event pair based on mutual information between them (PMI) and its variations like causal potential [3], Cause-Effect-Association (CEA) [8] reported good accuracy. We should explore the opportunities of such unsupervised approaches and evaluate their performance.

Our hypothesis pertinent to extraction of cascaded relations (issue 2 in Section 3) is that joint extraction of events and causal relations may better suit to cope with the complexity introduced by the inter-dependency between events and the causal expressions.

## 6    Conclusion

In this paper, we discuss about the necessity of the development causal relation extraction model in the cross-disciplinary field of climate science, marine science and environmental science. We describe the causality expressions found in oceanographic climate science domain and the issues need to be handled to develop a causal relation extraction model. We describe a brief survey of existing approaches of causal relation extraction from text data. Finally, we discuss about the suitability of existing causal relation extraction models in the considered domain and present the work plan.

## References

1. Batista-Navarro, R., Ananiadou, S.: Adapting the cluster ranking supervised model to resolve coreferences in the drug literature. In: Proceedings of the Fourth International Symposium on Languages in Biology and Medicine (LBM 2011) (2011)
2. Beamer, B., Bhat, S., Chee, B., Fister, A., Rozovskaya, A., Girju, R.: Uiuc: A knowledge-rich approach to identifying semantic relations between nominals. In: Proceedings of the 4th International Workshop on Semantic Evaluations. pp. 386–389. Association for Computational Linguistics (2007)
3. Beamer, B., Girju, R.: Using a bigram event model to predict causal potential. In: Proceedings of the 10th International Conference on Computational Linguistics and Intelligent Text Processing. pp. 430–441. CICLing '09, Springer-Verlag, Berlin, Heidelberg (2009)
4. Bethard, S., Corvey, W.J., Klingenstein, S., Martin, J.H.: Building a corpus of temporal-causal structure. In: LREC (2008)
5. Blanco, E., Castell, N., Moldovan, D.I.: Causal relation extraction. In: LREC (2008)
6. Chambers, N., Jurafsky, D.: Unsupervised learning of narrative schemas and their participants. In: Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2-Volume 2. pp. 602–610. Association for Computational Linguistics (2009)
7. Chambers, N., Jurafsky, D.: Unsupervised learning of narrative event chains. In: ACL. vol. 94305, pp. 789–797. Citeseer (2008)
8. Do, Q.X., Chan, Y.S., Roth, D.: Minimally supervised event causality identification. In: Proceedings of the Conference on Empirical Methods in Natural Language Processing. pp. 294–303. Association for Computational Linguistics (2011)
9. Fukuda, K., Tamura, A., Tsunoda, T., Takagi, T.: Toward information extraction: identifying protein names from biological papers. Pacific Symposium on Biocomputing. Pacific Symposium on Biocomputing pp. 707–718 (1998)

10. Garcia, D.: Coatis, an nlp system to locate expressions of actions connected by causality links. In: Proceedings of the 10th European Workshop on Knowledge Acquisition, Modeling and Management. pp. 347–352. EKAW '97, Springer-Verlag, London, UK, UK (1997)

11. Girju, R.: Automatic detection of causal relations for question answering. In: Proceedings of the ACL 2003 workshop on Multilingual summarization and question answering-Volume 12. pp. 76–83. Association for Computational Linguistics (2003)

12. Girju, R., Moldovan, D.I., et al.: Text mining for causal relations. In: Proceedings of Florida Artificial Intelligence Research Society (FLAIRS). pp. 360–364 (2002)

13. Girju, R., Nakov, P., Nastase, V., Szpakowicz, S., Turney, P., Yuret, D.: Semeval-2007 task 04: Classification of semantic relations between nominals. In: Proceedings of the 4th International Workshop on Semantic Evaluations. pp. 13–18. Association for Computational Linguistics (2007)

14. Grivaz, C.: Automatic extraction of causal knowledge from natural language texts. Ph.D. thesis, University of Geneva (2012)

15. Hashimoto, C., Torisawa, K., De Saeger, S., Oh, J.H., Kazama, J.: Excitatory or inhibitory: A new semantic orientation extracts contradiction and causality from the web. In: Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning. pp. 619–630. EMNLP-CoNLL '12, Association for Computational Linguistics, Stroudsburg, PA, USA (2012)

16. Hashimoto, C., Torisawa, K., Kloetzer, J., Sano, M., Varga, I., Oh, J.H., Kidawara, Y.: Toward future scenario generation: Extracting event causality exploiting semantic relation, context, and association features. In: ACL (1). pp. 987–997 (2014)

17. Ittoo, A., Bouma, G.: Extracting explicit and implicit causal relations from sparse, domain-specific texts. In: International Conference on Application of Natural Language to Information Systems. pp. 52–63. Springer (2011)

18. Joskowicz, L., Ksiezyk, T., Grishman, R.: Deep domain models for discourse analysis. In: Proc Annu AI Syst Gov Conf. pp. 195–200 (1989)

19. Kaplan, R.M., Berry-Rogghe, G.: Knowledge-based acquisition of causal relationships in text. Knowledge Acquisition 3(3), 317–337 (1991)

20. Khoo, C.S.G., Chan, S., Niu, Y.: Extracting causal knowledge from a medical database using graphical patterns. In: Proceedings of the 38th Annual Meeting on Association for Computational Linguistics. pp. 336–343. ACL '00, Association for Computational Linguistics, Stroudsburg, PA, USA (2000)

21. Khoo C., K.J., et al.: Automatic extraction of cause-effect information from newspaper text without knowledge-based inferencing. In: Literary & Linguistic Computing. pp. 177–186 (1998)

22. Kontos, J., Sidiropoulou, M.: On the acquisition of causal knowledge from scientific texts with attribute grammars. Internat. Journ. of Appl. Exp. Sys. 4(1), 31–48 (1991)

23. Marcu, D., Echihabi, A.: An unsupervised approach to recognizing discourse relations. In: Proceedings of the 40th Annual Meeting on Association for Computational Linguistics. pp. 368–375. Association for Computational Linguistics (2002)

24. Marsi, E., Ozturk, P.: Extraction and generalisation of variables from scientific publications. In: Proc. of EMNLP. pp. 505–511 (2015)

25. Marsi, E., Oztürk, P., Aamot, E., Sizov, G., Ardelan, M.V.: Towards text mining in climate science: Extraction of quantitative variables and their relations. In: Proceedings of the Fourth Workshop on Building and Evaluating Resources for Health and Biomedical Text Processing, Reykjavik, Iceland (2014)

26. Mihăilă, C., Ananiadou, S.: Semi-supervised learning of causal relations in biomedical scientific discourse. Biomedical engineering online 13(2), 1 (2014)

27. Miller, G.A.: Wordnet: a lexical database for english. Communications of the ACM 38(11), 39–41 (1995)

28. Miwa, M., Thompson, P., Ananiadou, S.: Boosting automatic event extraction from the literature using domain adaptation and coreference resolution. Bioinformatics 28(13), 1759–1765 (2012)

29. Pechsiri, C., Kawtrakul, A.: Mining causality from texts for question answering system. IEICE - Trans. Inf. Syst. E90-D(10), 1523–1533 (2007)

30. Pitler, E., Nenkova, A.: Using syntax to disambiguate explicit discourse connectives in text. In: Proceedings of the ACL-IJCNLP 2009 Conference Short Papers. pp. 13–16. Association for Computational Linguistics (2009)

31. Radinsky, K., Davidovich, S., Markovitch, S.: Learning causality for news events prediction. In: Proceedings of the 21st International Conference on World Wide Web. pp. 909–918. WWW '12, ACM, New York, NY, USA (2012)

32. Riaz, M., Girju, R.: Toward a better understanding of causality between verbal events: Extraction and analysis of the causal power of verb-verb associations. In: Proceedings of the SIGDIAL 2013 Conference. p. 21–30. Association for Computational Linguistics, Association for Computational Linguistics, Metz, France (August 2013)

33. Riaz, M., Girju, R.: Recognizing causality in verb-noun pairs via noun and verb semantics. EACL 2014 p. 48 (2014)

34. Rink, B., Bejan, C.A., Harabagiu, S.M.: Learning textual graph patterns to detect causal event relations. In: Guesgen, H.W., Murray, R.C. (eds.) FLAIRS Conference. AAAI Press (2010)

35. Selfridge, M.: Toward a natural language-based causal model acquisition system. Applied Artificial Intelligence an International Journal 3(2-3), 191–212 (1989)

36. Sporleder, C., Lascarides, A.: Using automatically labelled examples to classify rhetorical relations: An assessment. Nat. Lang. Eng. 14(3), 369–416 (Jul 2008)

37. Talmy, L.: Semantic causative types. the grammar of causative constructions. Syntax and Semantics 6 (1976)

38. Torisawa, K.: An unsupervised learning method for associative relationships between verb phrases. In: Proceedings of the 19th international conference on Computational linguistics-Volume 1. pp. 1–7. Association for Computational Linguistics (2002)

39. Torisawa, K.: Acquiring inference rules with temporal constraints by using japanese coordinated sentences and noun-verb co-occurrences. In: Proceedings of the Main Conference on Human Language Technology Conference of the North American Chapter of the Association of Computational Linguistics. pp. 57–64. HLT-NAACL '06, Association for Computational Linguistics, Stroudsburg, PA, USA (2006)