

Extracting Cause-Effect Pairs from a Sentence with a Dependency-Aware Transformer Model

Anonymous EACL submission

Abstract

001 Extracting cause and effect phrases from a sen-
002 tence is an important NLP task, with numer-
003 ous applications in various domains, includ-
004 ing legal, medical, education, and scientific
005 research. There are many unsupervised and
006 supervised methods proposed for solving this
007 task. Among these, unsupervised methods uti-
008 lize various linguistic tools, including syntac-
009 tic patterns, dependency tree, dependency re-
010 lations, etc. among different sentential units
011 for extracting the cause and effect phrases. On
012 the other hand, the contemporary supervised
013 methods use various deep learning based mask
014 language models equipped with a token clas-
015 sification layer for extracting cause and effect
016 phrases. Linguistic tools, specifically, depen-
017 dency tree, which organizes a sentence into
018 different semantic units have been shown to
019 be very effective for extracting semantic pairs
020 from a sentence, but existing supervised meth-
021 ods do not have any provision for utilizing such
022 tools within their model framework. In this
023 work, we propose DEPBERT, which extends a
024 transformer-based model by incorporating de-
025 pendency tree of a sentence within the model
026 framework. Extensive experiments over three
027 datasets show that DEPBERT is better than
028 various state-of-the-art supervised causality ex-
029 traction methods.

030 1 Introduction

031 Automatic extraction of cause and effect phrases
032 from natural language text is an important task with
033 enormous applications in various fields. In medi-
034 cal field, causal sentences are used for providing
035 the cause and the effects associated with diseases,
036 treatments, and side effects. Say, the sentence,
037 “Vitamin D deficiency contributes to both the initial
038 insulin resistance and the subsequent onset of
039 diabetes”, reflects disease causality; effectively ex-
040 tracting many such cause (deficiency of Vitamin D)
041 and effect (insulin resistance, diabetes) pairs (Wald
042 et al., 2002; Azagi et al., 2020) from medical text

can advance medical research. In fact, in medical
research, causality analysis provides the founda-
tion for generating complex hypotheses on which
new research can be designed. Causal sentences
are also used in legal fields for determining liability
and responsibility. Consider a sentence like, “The
company’s failure to adhere to safety regulations
resulted in a workplace accident.” Here, the causal
link between the company’s actions (not adhering
to safety regulations) and the accident is evident.
Methodologies for automatic extraction of such
causal relationships from legal texts can be use-
ful for building an AI-based legal assistant. In the
service field, an AI-based chatbot can provide diag-
nostic services to customers if cause-effect phrases
from instruction manuals can be mined effectively
and accurately.

Due to the importance of causality extraction,
several works (Atkinson and Rivas, 2008; Lee and
Shin, 2017; Zhao et al., 2018; An et al., 2019b;
Sorgente et al., 2013b; Kabir et al., 2022; An et al.,
2019a) have been proposed, which either identify
causal relations from sentences or extract cause
and effect pairs to build causal networks. These
existing works either use a set of known syntactic
patterns to extract the cause and effect entities or
use traditional supervised learning models (such as,
SVM) to identify those entity pairs. Considering
the cause and effect pairs as name entities, exist-
ing methods focus on entity extraction, and they
performed well when the causes and effects are
name entities, or noun phrases, such as, the name
of diseases, medications, or genes. There is another
group of research focusing on deep learning based
methods for causality extraction (Li et al., 2019;
Dasgupta et al., 2018). They mainly transform the
causality extraction into a token-classification task.
Within this group, Chansai et al. (Chansai et al.,
2021) fine tunes different learning methods (De-
vlin et al., 2018; Bojanowski et al., 2016; Touvron
et al., 2023) to make those amenable for the token

084 classification task. However, the major bottleneck
085 for all these supervised models is the lack of mech-
086 anisms for incorporating linguistic tools, such as,
087 dependency tree, syntactic patterns, etc. in those
088 models.

089 Dependency tree is an important linguistic tool
090 encompassing grammatical structure, syntax, se-
091 mantics, POS tags, and tag-to-tag interactions. It
092 plays a pivotal role in the extraction of cause-and-
093 effect relationships from textual data. Recent re-
094 search studies (Kabir et al., 2021, 2022) highlight
095 the critical significance of these linguistic compo-
096 nents in facilitating precise semantic relation ex-
097 traction. Dependency parsers, such as SpaCy (Hon-
098 nibal and Montani, 2020) and Stanza (Qi et al.,
099 2020), provide a powerful framework for dissect-
100 ing the intricate connections between words and
101 phrases within a sentence. Through syntactic analy-
102 sis these parsers ease identification of causal verbs,
103 subjects, and their corresponding objects. Addi-
104 tionally, POS tags and tag-to-tag interactions offer
105 valuable contextual information that aids in disam-
106 biguating causal relationships, thereby enhancing
107 the accuracy and reliability of extracted informa-
108 tion. The integration of these dependency-based ap-
109 proaches into supervised causality extraction mod-
110 els would improve the extraction of cause and effect
111 phrases, as we show in this paper.

112 In this paper we propose a transformer based
113 supervised method, DEPBERT, which seamlessly
114 integrates the dependency structure of sentences
115 into the bidirectional encoding representation of the
116 transformer model. Our method effectively merges
117 word-word co-occurrence, sentence semantics, and
118 the syntactic dependency structure within the do-
119 main of the transformer’s self-attention mechanism.
120 Through this integration, DEPBERT consistently
121 outperforms existing baseline methods, providing
122 clear evidence that the incorporation of dependency
123 structures significantly enhances the foundational
124 building blocks of the transformer architecture for
125 enhanced language understanding.

126 We claim the following two specific contribu-
127 tions:

- 128 • We introduce DEPBERT, a transformer model
129 that is sensitive to dependencies. It con-
130 currently learns from dependency relation
131 graphs, parts-of-speech tag sequences, and
132 token-token co-occurrences for token clas-
133 sification tasks, outperforming conventional
134 transformer-based language models in terms

of performance.

- We develop a dataset, referred to as CAUSAL-
GPT comprising 22,273 instances that include
cause terms, effect terms, and the sentences
containing both terms. The primary objec-
tive of this dataset is to alleviate the scarcity
of annotated datasets in the field of causality
extraction.

2 Related Works

The tasks of extracting causal relations can gen-
erally be classified into three main categories:
unsupervised (Khoo et al., 1998, 2001), super-
vised (Dasgupta et al., 2018; Li et al., 2019), and
hybrid approaches (Chang and Choi, 2006; Sor-
gente et al., 2013a; An et al., 2019a). Unsuper-
vised methods primarily rely on pattern-based ap-
proaches, employing causative verbs, causal links,
and relations between words or phrases to extract
cause-effect pairs. Supervised approaches, on the
other hand, require a labeled training dataset con-
taining pairs of cause and effect phrases, allowing
the training of supervised learning models for the
extraction of causal relationships between phrases.

Do et al. (Do et al., 2011) introduced a minimally
supervised approach based on a constrained condi-
tional model framework, incorporating discourse
connectives into their objective function. Dasgupta
et al. (Dasgupta et al., 2018) utilized word embed-
dings and selected linguistic features to construct
entity representations, serving as input for a bidirec-
tional Long-Short Term Memory (LSTM) model to
predict causal entity pairs. Nguyen et al. (Nguyen
and Grishman, 2015) harnessed pre-trained word
embeddings to train a convolutional Neural Net-
work (CNN) for classifying given causal pairs. In
contrast, Peng et al. (Peng et al., 2017) presented
a model-based approach that leverages deep learn-
ing architectures to classify relations between pairs
of drugs and mutations, as well as triplets involv-
ing drugs, genes, and mutations with N-ary rela-
tions across multiple sentences extracted from the
PubMed corpus.

Despite the various supervised methods avail-
able, existing causality extraction techniques often
fall short in incorporating dependency relations
within deep transformer architectures. While some
approaches do consider dependency relations (Ah-
mad et al., 2021; Song and King, 2022; Sachan
et al., 2020) to enhance language understanding
and introduce external sequential knowledge into

deep learning models (Wang et al., 2021), this work represents a novel fusion of learning from sequential knowledge, dependency relations, and co-occurring tokens.

3 Methodology

In this section, we first provide a formal discussion of token classification framework for solving the cause-effect pairs extraction. Then we provide the motivation and an overall framework of DEP-BERT, our proposed model. Finally, we discuss DEP-BERT’s architecture in details.

3.1 Problem formulation

Given, a sentence S and two phrases u (cause) and w (effect) in S , such that they exhibit a causality relation within the sentential context. Let S contains N number of tokens which are $s_1, s_2 \dots s_N$. Say, the cause phrase u consists of the token $s_i \dots s_K$, and effect phrase w consists of the tokens $s_j \dots s_L$, and there is no overlap between these two sequences of tokens. We also consider special tokens, such as, start token, end token and the padding tokens. Then we label all the tokens in S based on the following:

$$l_t = \begin{cases} 1, & \text{if } s_t \text{ is a **Special** token} \\ 2, & \text{if } s_t \text{ is a **Cause** token} \\ 3, & \text{if } s_t \text{ is an **Effect** token} \\ 4 & \text{otherwise} \end{cases}$$

we transform l_i into a one-hot encoding vector of size K , denoted as c_i , using the one-hot-encoding method (Harris and Harris, 2012; Brownlee, 2017) Suppose there is a model, θ , which predicts $p_1, p_2 \dots p_M$, where p_i represents the probability values predicted by the model for l_i . The token classification task objective of θ is to minimize the following multinomial cross-entropy loss function denoted by \mathcal{L}

$$\mathcal{L} = -\frac{1}{M} \sum_{i=1}^M \sum_{j=1}^K c_{i,j} \log p_{i,j}$$

3.2 DEP-BERT: Motivation and Design Justification

Syntactic patterns are important to extract semantic relation (Kabir et al., 2021, 2022) due to the fact that dependency relation plays an important role to identify semantic pairs. For instance, let there be a pattern u precipitates w , which can be applied to extract two semantic pairs u and w from

a sentence where u and w exhibit a cause effect semantic relation. However, the dependency relation and parts of speech tag for patterns can be crucial. For instance, here u needs to be a subject for the verb *causes* and u needs to be a *NOUN*, *causes* a *VERB*, and w be another *NOUN*. So for this particular pattern, dependency relation as well as *NOUN*, *VERB*, *NOUN* sequence can be important as well.

Traditional BERT (Devlin et al., 2018) lacks the capability to consider syntax and dependency relations when learning token embeddings, a shortcoming addressed by DEP-BERT. This innovative approach incorporates dependency relations and POS tag sequences into a transformer model designed to be acutely aware of these linguistic dependencies. While several dependency parsers are available for converting sentences into dependency trees, some previous research (Kabir et al., 2021)) has indicated that the performance of syntactic dependency pattern extraction does not depend much on the specific format of dependency tree. In this work, we have used Spacy dependency parser Spacy(Honnibal and Montani, 2020) for converting a sentence to a dependency tree, as the API of this library was convenient.

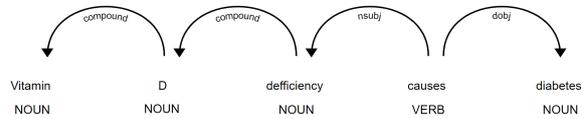


Figure 1: Dependency tree for the sentence, “Vitamin D deficiency causes diabetes”

As shown in Figure 1, a dependency tree is a dependency-aware representation of a sentence. If $\mathcal{G} = (V, E)$ is a dependency tree of S , then vertex-set V is associated with the tokens from S , and each of the vertices is labeled with POS tag of the corresponding token. For example, for the dependency tree in Figure 1, the sequence of POS tags are *NOUN*, *NOUN*, *NOUN*, *VERB*, and *NOUN*. Edgeset E represents the connecting pairs of tokens; an edge between two nodes, i and j , denote dependency relation between the token s_i and s_j . DEP-BERT’s main motivation is to utilize the dependency information in the transformer model. In any transformer architecture, a token receives attention from all other tokens. Likewise, in DEP-BERT, a token receives attention from other tokens in a traditional ways; besides, in a distinct channel, a token also receives attentions from other tokens connected through dependency tree. Since the to-

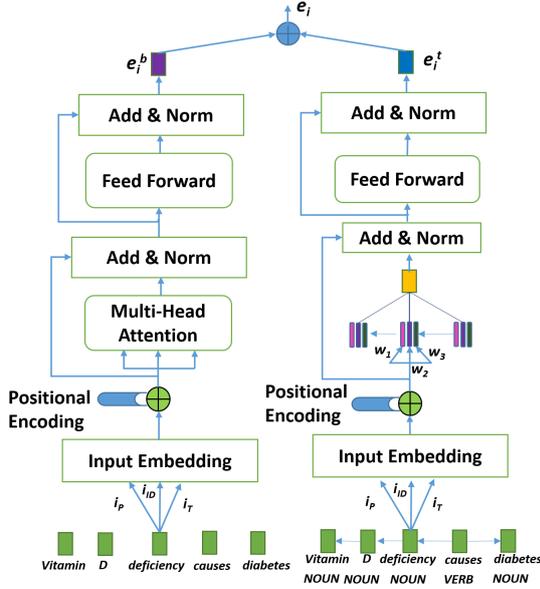


Figure 2: DEPBERT Model Architecture for Token Classification

270 kens in second channel holds POS information,
 271 these POS tags are utilized in DEPBERT, allowing
 272 it to incorporate semantic information of the tokens
 273 into the model. Last, we combine the two repre-
 274 sentations of token embedding coming from two
 275 channels: traditional token-based, and dependency
 276 graph based, which is then passed to a layer for to-
 277 ken classification. In this way, DEPBERT brings a
 278 flavor of graph attention network (Veličković et al.,
 279 2017) in its attention propagation mechanism to
 280 learn a better embedding of the tokens of a sen-
 281 tence.

282 3.3 Model Architecture

283 Figure 2 presents the token classification model
 284 of DEPBERT. The diagram showcases two main
 285 components: the token embedding encoder of a
 286 traditional BERT (Devlin et al., 2018) on the left
 287 side and the encoder (Raganato and Tiedemann,
 288 2018) architecture with a modified self-attention
 289 mechanism incorporating a dependency graph on
 290 the right side. To facilitate our discussion, we uti-
 291 lize the running example of the text “Vitamin D
 292 deficiency causes diabetes”.

293 At the bottom of the right part of the architec-
 294 ture, we show the dependency association graph
 295 of this sentence from Spacy. As we can see the to-
 296 kens for this sentence, the POS tags are: {NOUN,
 297 NOUN, NOUN, VERB, and NOUN}, from left to
 298 right. In the dependency graph there is an edge
 299 for each of the following token pairs – (D, Vita-

300 min), (deficiency, D), (causes, deficiency), (causes,
 301 diabetes).

302 Our final model is a two tower model where both
 303 left and right sides learn embedding of tokens in
 304 parallel. For both towers, the standard tokenizer
 305 is modified to incorporate POS tags. To provide
 306 a more detailed illustration, DEPBERT generates
 307 three types of unique IDs for each token: i_T for the
 308 parts of speech tag, i_P for the positional embed-
 309 ding, and i_{ID} as the input ID. These IDs are then
 310 passed to an embedding layer, and the resulting
 311 embeddings are summed to create a single vector
 312 per token, denoted as v_i .

313 For left side of the model, all representation vec-
 314 tors are passed through a multi-head attention layer.
 315 The multi-head attention layer learns attention for
 316 all the tokens, and the output of this layer is calcu-
 317 lated based on the attention scores of all the tokens.
 318 This output is then normalized using a layer nor-
 319 malization layers and a feed forward layer to form
 320 e_i^b which is the output embedding for token s_i from
 321 the left part of the model. Meanwhile, the gener-
 322 ated tree for S is fit to the right side of the diagram.
 323 Like the BERT encoder architecture, embedding
 324 of input id, and token type id, and positional id are
 325 gathered. These three embeddings are then added
 326 to construct a single vector v_i . The embedding of
 327 other tokens are also calculated in the similar fash-
 328 ion. Let the tokens connected with s_i through edges
 329 is the set \mathcal{V} , and s_k be any token in \mathcal{V} . Additionally,
 330 let there be three trainable matrices \mathbf{W}_1 , \mathbf{W}_2 , and
 331 \mathbf{W}_3 . $v_i \mathbf{W}_1$, $v_i \mathbf{W}_2$, and $v_i \mathbf{W}_3$ are then query, key,
 332 and value vectors respectively for the token s_i . The
 333 affinity score of two connected tokens s_i , and s_k is
 334 then calculated by the following equation.

$$335 a_{ik} = (v_i \mathbf{W}_1) * (v_k \mathbf{W}_2)^T$$

336 The attention value, α_{ik} (a scalar) is the softmax-
 337 score of these affinity values for all the tokens in \mathcal{V}
 338 which actually represents how important the con-
 339 nected token is with respect to current token.

$$340 \alpha_{ik} = \frac{e^{a_{ik}}}{\sum_{j \in [1, |\mathcal{V}|]} e^{a_{ij}}}$$

341 The attention scores are used for attention based
 342 weighted sum for the output vector, o_i from atten-
 343 tion layer, as shown in the next equation (\mathbf{W}_4 is
 344 another trainable matrix and $v_j \mathbf{W}_3$ is a value vec-
 345 tor for the corresponding token.). Unlike BERT,
 346 o_i is calculated for only those tokens which are
 347 connected to s_i .

$$o_i = \sum_{k \in [1, |\mathcal{V}|]} \alpha_{ik} (v_j \mathbf{W}_3) \mathbf{W}_4$$

o_i is then passed through a normalization layer, and the output from **Add and Norm** layer, \bar{o}_i is calculated using the following equation where γ , β are trainable scalars, ϵ is a very small scalar constant, μ_i , and σ^2 are the mean, and variance of the vector o_i .

$$\bar{o}_i = v_i + \gamma \odot \frac{o_i - \mu_i}{\sigma^2 + \epsilon} + \beta$$

The output \bar{o}_i is furthermore passed through a feed forward layer with GELU activation (Hendrycks and Gimpel, 2016) using another trainable matrix \mathbf{W}_5 , and bias b .

$$FFN(\bar{o}_i) = GELU(\bar{o}_i \mathbf{W}_5 + b)$$

Additionally, e_i^t is calculated passing the feed forward layer’s output to another **Add and Norm** layer. Meanwhile, e_i^b and e_i^c are passed through another gate to form e_i which is the token embedding of s_i using DEPBERT. This embedding is passed to another neural network for the token classification task.

$$e_i^s = \sigma(e_i^b \mathbf{W}_6 + c)$$

$$e_i = e_i^s \odot e_i^b + (1 - e_i^s) \odot e_i^c$$

4 Experiment and Result

We perform comprehensive experiment to show the effectiveness of DEPBERT for token classification task. Below we first discuss the dataset, and competing methods, followed by experimental results.

4.1 Dataset

Semeval: This is a popular benchmark dataset, built by combining the SemEval 2007 Task 4 dataset (Girju et al., 2007) and the SemEval 2010 Task 8 datasets (Hendrickx et al., 2010). A row for SemEval datasets contains a term pair, a sentence containing this pair, and the semantic relation. The SemEval 2007 Task 4 possesses 7 semantic relations whereas the SemEval 2010 Task 8 describes 9 relations. However, cause-effect relation is common among these two tasks. The datasets include predefined train and test partitions. For building validation partition, we borrow from the train partition. Train, test and validation partitions are then merged to concatenate into a single dataset. Note

that, the merged dataset contains 14 relations of which we consider the instances of cause-effect relations only. The preprocessed dataset, as a result, comprises a total of 1,427 instances, all exclusively related to cause-effect relations. In terms of partition distribution, the training, test, and validation segments account for 60%, 30%, and 10%, respectively. **SCITE**: This is another dataset of the paper SCITE (Li et al., 2021). The dataset contains 1079 sentences exhibiting cause-effect relation. We split the dataset into training, test and validation subsets maintaining the same ratio, 6:3:1 like Semeval. Additionally, each row of this dataset also maintains the same format like the Semeval.

CAUSALGPT: The previously described datasets have a very limited number of sentences as those are annotated by human. However, motivated from the recent generative models such as Bard and ChatGPT, we create a new dataset containing adequate number of sentences. To create this dataset, we harness the power of Large Language Models (LLMs) to generate cause terms, effect terms, and sentences that preserve the cause-effect relationship. It’s important to note that the generated sentences may exhibit duplication and may sometimes contain special Unicode strings that require preprocessing. To ensure a wide variety of sentence structures, we develop a program capable of generating both active and passive sentence constructions. Additionally, our focus during sentence generation is on the medical domain. In total, our dataset comprises 22,273 sentences, making it a valuable resource for in-depth research into cause-effect phrase extraction. Much like SCITE, we partition this dataset randomly into training, test, and validation subsets while maintaining the same proportional distribution. In Table 1 we show some instances of CAUSALGPT dataset; The cause term, effect term, and the sentences are in Column one, two, and three respectively.

4.2 Competing Methods

For comparison, we consider a collection of neural architectures, including BERT, Sentence-BERT, LLaMA, Dargupta, and SCITE methods. We have also developed several models, such as BERT plus dependency and BERT plus POS tags. We discuss all of the competing methods below.

Table 1: Example Sentences from the CAUSALGPT

Cause Term	Effect Term	Sentence
Diabetes	Blindness	Diabetes can lead to blindness if left uncontrolled.
Vitamin D deficiency	Osteoporosis	A deficiency in vitamin D can result in osteoporosis.
Smoking	Lung Cancer	Smoking is a major risk factor for developing lung cancer.
High cholesterol	Heart Disease	Elevated cholesterol levels are associated with an increased risk of heart disease.
Obesity	Type 2 Diabetes	Obesity is a significant risk factor for developing type 2 diabetes.

4.2.1 Dasgupta

Dasgupta’s (Dasgupta et al., 2018) method is one of the first deep learning methods to extract cause-effect pairs from sentences. They design the method for token classification task. Each token is labelled either as a cause word, a effect word, causal connects or None. They learn word embedding by both word2vec (Mikolov et al., 2013) and linguistic feature vector (Dasgupta et al., 2018). Each of the embedding is fit to a bidirectional Long Short Term Memory (Hochreiter and Schmidhuber, 1997) architecture for token classification.

4.2.2 SCITE

Another method SCITE (Li et al., 2021) uses a multi head self attention mechanism (Vaswani et al., 2017), and a conditional random field (Fields, 2001) along with a bi-LSTM architecture. Additionally, flair embedding (Akbik et al., 2018) is learned in a large context and transferred the string embedding for the task of causality extraction.

4.2.3 BERT

Bidirectional Encoder Representations from Transformers (BERT) is proposed by researchers from Google (Devlin et al., 2018), which is not trained on any specific downstream task but instead on a more generic task called Masked Language Modeling. The idea is to leverage huge amounts of unlabeled data to pre-train a model, which can be fine-tuned to solve different kinds of NLP tasks by adding a task specific layer which maps the contextualized token embeddings into the desired output function. In this work we use the pre-trained model “bert- base-uncased” which has a vocabulary of 30K tokens and 768 dimension for each token. We use the BERT embedding for token classification.

4.2.4 BERT plus dependency

We design this baseline model to incorporate the dependency structure, but unlike DEPBERT, it does not consider POS tags. In both towers of DEPBERT, we utilize the tokenizer from the original BERT model. However, this baseline model differs

from BERT in that it is partially pretrained. While we do not have pretrained embeddings specifically for the dependency relation, we utilize the pre-trained model “bert- base-uncased” for the left tower instead of training the model on an extensive corpus like Wikipedia or Google corpus.

4.2.5 BERT plus POS Tags

We develop this baseline as well to investigate the relative importance of POS tags compared to the dependency relation. In contrast to the previous baseline, this model involves modifying the existing BERT tokenizer to include POS tokens. However, the dependency relation is not taken into account, resulting in a single tower model. Similar to the previous baseline, this model is semi-pretrained for the input ID and positional embedding.

4.2.6 Sentence-BERT

Sentence-BERT (Reimers and Gurevych, 2019) is another pretrained model designed to capture mainly semantic textual similarity. The model uses siamese and triplet network structures to derive semantically meaningful sentence embedding. The model can also be used to build token representation for token classification. For comparison we use the pretrained weights from the model “all-MiniLM-L6-v2” available in Huggingface, which produces 384 dimensional vectors for each token.

4.2.7 LLaMA

LLaMA (Touvron et al., 2023) is another transformer based model developed by researchers in Meta. The model contains 7B to 65B parameters and it is trained on trillions of tokens. The model outperforms GPT 3 and other state of the art methods. The pre-trained model is available online. We use the pretrained model to represent word tokens. The representation of each tokens (4096 dimensional) is learned by LLaMA for token classification.

Table 2: Performance of DEPBERT compared to baseline methods in CAUSALGPT Dataset

Method	Prec	Rec	F ₁ Score (% imp.)	Acc (Exact) (% imp.)
Bi-LSTM (Dasgupta)	0.911	0.828	0.847 (-)	0.778 (-)
Bi-LSTM-CRF (SCITE)	0.899	0.834	0.849 (0.23)	0.781 (0.4)
BERT	0.942	0.958	0.938 (10.7)	0.811 (4.2)
BERT plus dependency	0.956	0.967	0.958 (13.1)	0.833 (7.1)
BERT plus POS tags	0.921	0.911	0.897 (5.9)	0.831 (6.8)
Sentence-BERT	0.946	0.964	0.948 (11.9)	0.822 (5.7)
LLaMA	0.953	0.956	0.954 (12.6)	0.828 (6.4)
DEPBERT (Gated)	0.967	0.969	0.963 (13.7)	0.858 (10.3)

Table 3: Performance of DEPBERT compared to baseline methods in SemEval Dataset

Method	Prec	Rec	F ₁ Score (% imp.)	Acc (Exact) (% imp.)
Bi-LSTM (Dasgupta)	0.917	0.825	0.844 (-)	0.768 (-)
Bi-LSTM-CRF (SCITE)	0.896	0.851	0.86 (2)	0.771 (0.4)
BERT	0.936	0.941	0.932 (10.4)	0.809 (5.3)
BERT plus dependency	0.951	0.959	0.954 (13)	0.841 (9.5)
BERT plus POS tags	0.937	0.947	0.941 (11.5)	0.831 (8.2)
Sentence-BERT	0.926	0.936	0.933(10.5)	0.818 (6.5)
LLaMA	0.938	0.921	0.937 (11)	0.819 (6.6)
DEPBERT (Gated)	0.942	0.962	0.957 (13.4)	0.842 (9.63)

Table 4: Performance of DEPBERT compared to baseline methods in SCITE Dataset

Method	Prec	Rec	F ₁ Score (% imp.)	Acc (Exact) (% imp.)
Bi-LSTM (Dasgupta)	0.811	0.825	0.817 (-)	0.747 (-)
Bi-LSTM-CRF (SCITE)	0.832	0.849	0.831 (1.7)	0.751 (0.53)
BERT	0.884	0.9	0.893 (9.3)	0.768 (2.8)
BERT plus dependency	0.911	0.923	0.916 (12.1)	0.811 (8.5)
BERT plus POS tags	0.908	0.917	0.906 (10.9)	0.796 (6.5)
Sentence-BERT	0.887	0.889	0.886 (8.4)	0.773 (3.5)
LLaMA	0.9	0.913	0.909 (11.3)	0.79 (5.7)
DEPBERT (Gated)	0.932	0.943	0.939 (14.9)	0.834 (11.6)

4.3 Experimental Setup

Our DEPBERT model contains 227 millions parameters, and all of them are trainable. The left tower of the DEPBERT architecture is initialized with pretrained BERT parameters sourced from *bert-uncased* model. It is important to note that, in our model architecture, no additional external hyperparameters are introduced. We consistently emphasize the default setup for all variations of our methods. Specifically, for LLaMA, BERT, BERT plus Dependency, BERT plus POS tags, SentenceBERT, the number of trainable parameters stands at 524 million, 109 million, 110 million, and 110 million, respectively. In contrast, both Dasgupta’s method and SCITE feature a relatively smaller number of hyperparameters, around 400K for each. In all of our models, we make use of the Adam optimizer with a batch size of 128 and a default learning rate of 0.001. Additionally, we implement early stopping with 1000 epochs and a tolerance rate of 10, with the majority of the models concluding training within the first 100 epochs. It’s worth mentioning that all the results presented in this research are derived from the initial stable runs.

4.4 Results

We conducted comprehensive experiments that covered all baseline methods, including DEPBERT, across the three previously described datasets. In these experiments, accuracy was determined without allowing for partial matches. Each sentence contains a causal entity and an effect entity, and each of them may consist of one or multiple tokens. To register a correct prediction, a model needs to accurately predict all the causal and effect tokens. The results for all three datasets are conveniently presented in Tables 2, 3, and 4. The last column in each table highlights the exact accuracy score. Additionally, our evaluation encompassed a comprehensive range of metrics, such as precision, recall, and F_1 score, to ensure a thorough assessment.

Table 2 presents the performance of all the methods on our CAUSALGPT dataset. It is observed that DEPBERT achieves a 10.3% higher exact matching accuracy compared to Dasgupta’s method, indicating that approximately 86% of the extracted pairs precisely match the actual pairs. Furthermore, DEPBERT exhibits the best precision, recall, and F_1 score. Another noteworthy finding is that the BERT plus dependency method, which we designed, outperforms other baselines. This clearly

demonstrates the significance of the dependency relation over POS tags. However, the combination of POS tags and the dependency relation yields even more meaningful results than solely incorporating POS tags into the model.

Similarly, Table 3 presents the results on the SemEval dataset. The performance of all the methods on this dataset is slightly lower compared to the previous dataset. This could be attributed to the nature of cause-effect sentences. Moreover, the limited number of sentences in this dataset is insufficient to effectively train deep learning models. Nonetheless, DEPBERT and BERT plus dependency still outperform other methods, following a similar pattern as observed previously. In terms of F_1 score, DEPBERT achieves a 13.4% improvement compared to Dasgupta’s method.

Furthermore, Table 4 displays the performance of all the baseline methods on the SCITE dataset. Due to the insufficient number of sentences in this dataset as well, the performance of all the methods falls short of expectations. However, DEPBERT outperforms all other baseline methods even for this dataset. The accuracy and F_1 score are 0.834 and 0.939, respectively, marking an improvement of 11.6% and 14.9% compared to Dasgupta’s method.

Clearly, DEPBERT’s dependency and parts-of-speech aware attention mechanism contribute to its superiority over other methods. Moreover, the combination of POS tags and the dependency relation proves to be more effective in extracting cause-effect pairs from sentences.

5 Conclusion

In this study, we have effectively unveiled a pioneering method for extracting causal relationships, drawing inspiration from the sentence’s underlying dependency structure. Our model, named DEPBERT, stands out by fusing the transformer architecture with dependency graph networks, harnessing the power of dependency relations and parts-of-speech markers. This amalgamation yields a marked improvement in the precision of causal relationship extraction across a multitude of domains. Looking ahead, the expansive utility of such models across various domains presents a promising path for advancing information extraction methodologies.

616 **6 Limitations**

617 While the DEPBERT model demonstrates superior
618 performance compared to baseline methods, it’s
619 important to note that its performance is intricately
620 tied to the characteristics of the dataset. While large
621 language models (LLMs) like ChatGPT can extract
622 cause-effect pairs, even in zero-shot corpora, it’s
623 crucial to clarify that our paper does not intend to
624 diminish the significance of LLMs. Instead, our
625 primary emphasis lies in enhancing the founda-
626 tional elements of transformer architectures to in-
627 corporate sentence dependency structures. Another
628 limitation of our research pertains to the newly cre-
629 ated dataset, CausalGPT. Unlike DEPBERT, which
630 can extract multiple semantic pairs from sentences
631 simultaneously, the CausalGPT dataset is intention-
632 ally constructed so that each sentence contains only
633 one semantic pair. Consequently, if DEPBERT is
634 trained with this dataset, it cannot extract multiple
635 semantic pairs from sentences. That’s why, in the
636 context of this specific study, we did not conduct an
637 evaluation of its performance in extracting multiple
638 semantic pairs. Furthermore, we specifically em-
639 phasize semantic pairs within the English language.
640 While semantic pairs and dependency relationships
641 can be of great importance in all languages, it’s
642 worth noting that our research did not include ex-
643 periments in languages other than English.

644 **7 Ethical Impacts**

645 This research plays a pivotal role in the extraction
646 of cause-effect relationships, negating the reliance
647 on syntactic dependency patterns. The cause-and-
648 effect connection serves as a foundational and in-
649 dispensable element within linguistics and logic,
650 acting as the linchpin for comprehending the intri-
651 cate web of associations between events and their
652 consequences. The extraction of causality holds
653 paramount importance across a diverse spectrum
654 of fields, encompassing law, medicine, and event
655 analysis, for it provides the means to unearth the
656 concealed mechanisms and repercussions that un-
657 derlie a wide array of phenomena. Ultimately, this
658 research empowers us to make well-informed deci-
659 sions, pinpoint root causes, and enhance outcomes.

660 Within the legal domain, the capacity to extract
661 cause-effect relationships without being bound by
662 syntactic dependencies is nothing short of indis-
663 pensable. The legal system hinges on precise com-
664 prehension and documentation of causality, as it
665 is fundamental to establishing liability, attributing

666 fault, and ensuring accountability. This research
667 equips legal professionals with the tools to unravel
668 the causal connections embedded in intricate legal
669 cases, thereby simplifying the process of identify-
670 ing the factors leading to specific events or circum-
671 stances. Consequently, it bolsters the pursuit of
672 justice, whether in civil or criminal proceedings.

673 In the sphere of medicine, the extraction of
674 causality fulfills a pivotal role in diagnostics and
675 treatment. It empowers medical practitioners to dis-
676 cern the underlying causes of diseases and ailments,
677 thereby facilitating more accurate and timely diag-
678 noses. This, in turn, not only elevates the quality
679 of patient care but also expedites the development
680 of more effective treatment strategies. Moreover,
681 comprehending the cause-and-effect relationships
682 between various medical variables proves instru-
683 mental in public health endeavors, including epi-
684 demiological studies and the management of dis-
685 ease outbreaks.

686 In the context of event extraction, this research
687 emerges as an indispensable tool across various
688 applications, spanning disaster response, business
689 analytics, and social science research. When grap-
690 pling with extensive datasets, the identification of
691 causality allows organizations and researchers to
692 fathom the core drivers of specific events. In the
693 realm of disaster response, it aids in comprehending
694 the triggers and consequences of natural dis-
695 asters, thus enhancing preparedness and response
696 strategies. In the domain of business analytics, it
697 facilitates the identification of factors influencing
698 financial performance and market trends. For so-
699 cial science research, it provides a foundational
700 framework for unraveling the complex dynamics
701 that govern society, shedding light on the causes
702 and effects of a multitude of social phenomena.

703 In summary, the extraction of cause-effect re-
704 lationships, liberated from syntactic dependency
705 patterns, emerges as a cornerstone in the domains
706 of law, medicine, and event extraction, primarily
707 due to its role in enhancing precision, accuracy, and
708 the depth of understanding of causality within these
709 fields. This, in turn, results in more well-informed
710 decision-making, improved outcomes, and an in-
711 creased capacity to effectively address complex
712 challenges.

713 **References**

714 Wasi Ahmad, Haoran Li, Kai-Wei Chang, and Yashar
715 Mehdad. 2021. [Syntax-augmented multilingual](#)

716	BERT for cross-lingual transfer . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 4538–4554, Online. Association for Computational Linguistics.	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding.	771 772 773 774
722	Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence labeling . In <i>Proceedings of the 27th International Conference on Computational Linguistics</i> , pages 1638–1649, Santa Fe, New Mexico, USA. Association for Computational Linguistics.	Quang Do, Yee Seng Chan, and Dan Roth. 2011. Minimally supervised event causality identification. In <i>Proceedings of the 2011 Conference on Empirical Methods in Natural Language Processing</i> , pages 294–303.	775 776 777 778 779
728	Ning An, Yongbo Xiao, Jing Yuan, Yang Jiaoyun, and Gil Alterovitz. 2019a. Extracting causal relations from the literature with word vector mapping . volume 115, page 103524.	Conditional Random Fields. 2001. Probabilistic models for segmenting and labeling sequence data. In <i>ICML 2001</i> .	780 781 782
732	Ning An, Yongbo Xiao, Jing Yuan, Jiaoyun Yang, and Gil Alterovitz. 2019b. Extracting causal relations from the literature with word vector mapping . <i>Computers in biology and medicine</i> , 115:103524.	Roxana Girju, Preslav Nakov, Vivi Nastase, Stan Szpakowicz, Peter Turney, and Deniz Yuret. 2007. SemEval-2007 task 04: Classification of semantic relations between nominals . In <i>Proceedings of the Fourth International Workshop on Semantic Evaluations (SemEval-2007)</i> , pages 13–18, Prague, Czech Republic. Association for Computational Linguistics.	783 784 785 786 787 788 789
736	John Atkinson and Alejandro Rivas. 2008. Discovering novel causal patterns from biomedical natural-language texts using bayesian nets. <i>IEEE Transactions on Information Technology in Biomedicine</i> , 12(6):714–722.	David Harris and Sarah L Harris. 2012. <i>Digital design and computer architecture</i> . Morgan Kaufmann.	790 791
741	Tal Azagi, Hein Sprong, Dieuwertje Hoornstra, and Joppe Hovius. 2020. Evaluation of disease causality of rare ixodes ricinus-borne infections in europe . volume 9, page 150.	Iris Hendrickx, Su Nam Kim, Zornitsa Kozareva, Preslav Nakov, Diarmuid Ó Séaghdha, Sebastian Padó, Marco Pennacchiotti, Lorenza Romano, and Stan Szpakowicz. 2010. SemEval-2010 task 8: Multiway classification of semantic relations between pairs of nominals . In <i>Proceedings of the 5th International Workshop on Semantic Evaluation</i> , pages 33–38, Uppsala, Sweden. Association for Computational Linguistics.	792 793 794 795 796 797 798 799 800
745	Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2016. Enriching word vectors with subword information. volume 5.	Dan Hendrycks and Kevin Gimpel. 2016. Bridging nonlinearities and stochastic regularizers with gaussian error linear units . <i>CoRR</i> , abs/1606.08415.	801 802 803
748	Jason Brownlee. 2017. Why one-hot encode data in machine learning. <i>Machine Learning Mastery</i> , pages 1–46.	Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. <i>Neural Computation</i> , 9(8):1735–1780.	804 805 806
751	Du-Seong Chang and Key-Sun Choi. 2006. Incremental cue phrase learning and bootstrapping method for causality extraction using cue phrase and word pair probabilities . volume 42, pages 662–678.	Matthew Honnibal and Ines Montani. 2020. spaCy 2.2.3: Industrial-strength natural language processing. In https://spacy.io/ .	807 808 809
755	Terapat Chansai, Ruksit Rojpaisarnkit, Teerakarn Bori-boonsub, Suppawong Tuarob, Myat Su Yin, Peter Haddawy, Saeed-Ul Hassan, and Mihai Pomarlan. 2021. Automatic cause-effect relation extraction from dental textbooks using bert. In <i>Towards Open and Trustworthy Digital Societies: 23rd International Conference on Asia-Pacific Digital Libraries, ICADL 2021, Virtual Event, December 1–3, 2021, Proceedings 23</i> , pages 127–138. Springer.	Md Kabir, AlJohara Almulhim, Xiao Luo, and Mohammad Hasan. 2022. Informative causality extraction from medical literature via dependency-tree based patterns . volume 6, pages 295–316.	810 811 812 813
764	Tirthankar Dasgupta, Rupsa Saha, Lipika Dey, and Abir Naskar. 2018. Automatic extraction of causal relations from text using linguistically informed deep neural networks . In <i>Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue</i> , pages 306–316, Melbourne, Australia. Association for Computational Linguistics.	Md. Ahsanul Kabir, Typer Phillips, Xiao Luo, and Mohammad Al Hasan. 2021. Asper: Attention-based approach to extract syntactic patterns denoting semantic relations in sentential context .	814 815 816 817
766		C. Khoo, Jaklin Kornfilt, R. Oddy, and Sung-Hyon Myaeng. 1998. Automatic extraction of cause-effect information from newspaper text without knowledge-based inferencing. volume 13, pages 177–186.	818 819 820 821
767		Christopher Khoo, Sung-Hyon Myaeng, and Robert Oddy. 2001. Using cause-effect relations in text to improve information retrieval precision . volume 37, pages 119–145.	822 823 824 825

826	Dong-gi Lee and Hyunjung Shin. 2017. Disease causality extraction based on lexical semantics and document-clause frequency from biomedical literature. <i>BMC medical informatics and decision making</i> , 17(1):53.	Antonio Sorgente, Giuseppe Vettigli, and Francesco Mele. 2013b. Automatic extraction of cause-effect relations in natural language text. <i>DART@ AI* IA</i> , 2013:37–48.	879
827			880
828			881
829			882
830			
831	Zhaoning Li, Qi Li, Xiaotian Zou, and Jiangtao Ren. 2019. Causality extraction based on self-attentive bilstm-crf with transferred embeddings.	Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. <i>arXiv preprint arXiv:2302.13971</i> .	883
832			884
833			885
834	Zhaoning Li, Qi Li, Xiaotian Zou, and Jiangtao Ren. 2021. Causality extraction based on self-attentive bilstm-crf with transferred embeddings. <i>Neurocomputing</i> , 423:207–219.		886
835			887
836			888
837			
838	Tomas Mikolov, G.s Corrado, Kai Chen, and Jeffrey Dean. 2013. Efficient estimation of word representations in vector space. pages 1–12.	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. <i>Advances in neural information processing systems</i> , 30.	889
839			890
840			891
841	Thien Huu Nguyen and Ralph Grishman. 2015. Relation extraction: Perspective from convolutional neural networks. In <i>Proceedings of the 1st Workshop on Vector Space Modeling for Natural Language Processing</i> , pages 39–48.	Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua Bengio. 2017. Graph attention networks. <i>arXiv preprint arXiv:1710.10903</i> .	892
842			893
843			
844			894
845			895
846	Nanyun Peng, Hoifung Poon, Chris Quirk, Kristina Toutanova, and Wen-tau Yih. 2017. Cross-sentence n-ary relation extraction with graph lstms. <i>Transactions of the Association for Computational Linguistics</i> , 5:101–115.	David Wald, Malcolm Law, and Joan Morris. 2002. H omocysteine and cardiovascular disease: Evidence on causality from a meta-analysis. <i>BMJ (Clinical research ed.)</i> , 325:1202.	896
847			897
848			
849			898
850			899
851	Peng Qi, Yuhao Zhang, Yuhui Zhang, Jason Bolton, and Christopher D. Manning. 2020. Stanza: A Python natural language processing toolkit for many human languages. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> .	Yaxuan Wang, Hanqing Lu, Yunwen Xu, Rahul Goutam, Yiwei Song, and Bing Yin. 2021. Q ueen: Neural query rewriting in e-commerce. In <i>The Web Conference 2021</i> .	900
852			901
853			
854			902
855			903
856			904
857	Alessandro Raganato and Jörg Tiedemann. 2018. An analysis of encoder representations in transformer-based machine translation. In <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP</i> . The Association for Computational Linguistics.	Sendong Zhao, Meng Jiang, Ming Liu, Bing Qin, and Ting Liu. 2018. Causaltriad: toward pseudo causal relation discovery and hypotheses generation from medical text data. In <i>Proceedings of the 2018 ACM International Conference on Bioinformatics, Computational Biology, and Health Informatics</i> , pages 184–193.	905
858			906
859			907
860			908
861			909
862			910
863	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert: Sentence embeddings using siamese bert-networks. <i>arXiv preprint arXiv:1908.10084</i> .		911
864			912
865			
866	Devendra Singh Sachan, Yuhao Zhang, Peng Qi, and William Hamilton. 2020. Do syntax trees help pre-trained transformers extract information? <i>arXiv preprint arXiv:2008.09084</i> .		
867			
868			
869			
870	Zixing Song and Irwin King. 2022. Hierarchical heterogeneous graph attention network for syntax-aware summarization. In <i>Proceedings of the AAAI Conference on Artificial Intelligence</i> , volume 36, pages 11340–11348.		
871			
872			
873			
874			
875	Antonio Sorgente, G. Vettigli, and Francesco Mele. 2013a. Automatic extraction of cause-effect relations in natural language text. volume 1109, pages 37–48.		
876			
877			
878			