



西北工业大学
NORTHWESTERN POLYTECHNICAL UNIVERSITY

基于对比表示学习的因果 关系抽取技术研究

汇报人：刘磊

指导老师：梁韵基老师



目录

CONTENTS

01

选题背景及意义

02

相关工作

03

研究思路与方法

04

研究成果与应用

05

未来计划与论文总结



研究背景



□ 人们的社会活动以事件为驱动，事件是社会活动的载体。

在自然语言处理中，事件作为一种信息表示形式显得十分重要，受到越来越多的重视。事件中包含了大量的内部组成结构（如参与者、时间、地点等）和外部关联（如因果、共指、时序等语义关系）。

□ 因果关系在诸多领域发挥着越来越重要的作用

在**人机交互、航空安全、医学、经济、政治、社会发现和自然科学**等领域应用日益增长



航空安全领域本体因果关系主要概念的构成



人机对话中文本关系信息抽取



精神疾病成因分析



“将因果关系整合到人工智能中是一件大事！” ——Yoshua Bengio



突发生物安全事件成因分析



研究背景



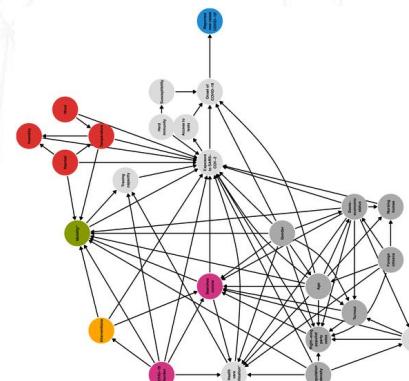
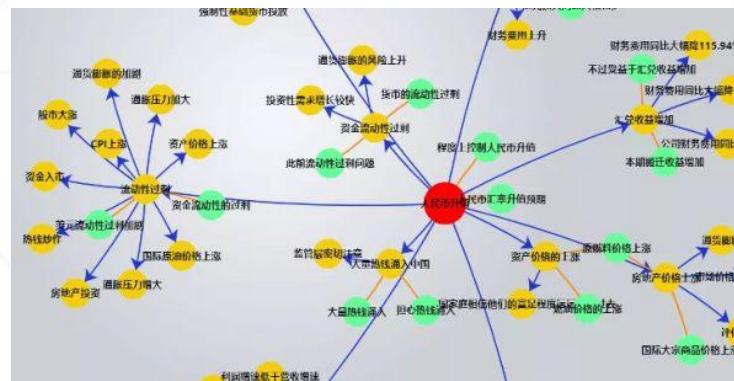
□ 什么是因果关系?

因果关系是一种由现象A导致现象B发生的**单向关系**。

□ 现有的因果研究有哪些?

- 因果分析 (Causal Analysis)
- 因果推断 (Causal Inference)
- 因果关系抽取 (Causality Extraction)

□ 现有的因果关系组织形式?

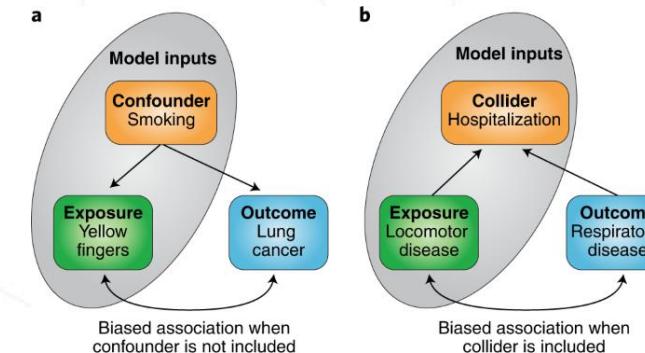


	flu	causes	cold	.
flu	0	0	0	0
causes	1	0	1	1
cold	0	0	0	0
.	0	0	0	0

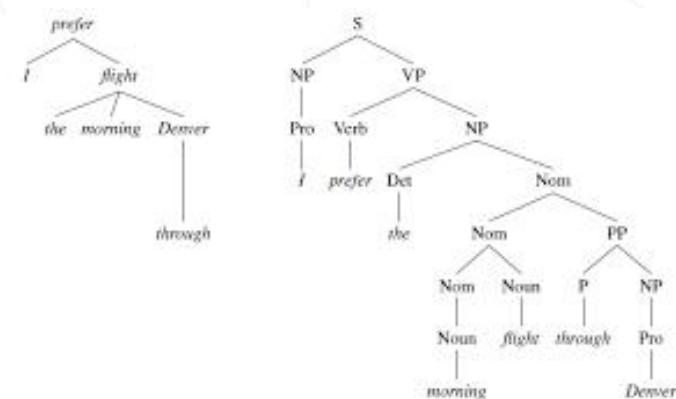
Nsubj Root Xcomp Punct

I got so **depressed** because I **lost my favorite toy.**

{[lost my favorite toy], C-E, [depressed]}



因果推断



因果分析知识图谱

Causal Graph

稀疏矩阵

语法依存树



研究背景



□ 因果关系的类型有哪些？

- 是否包含因果触发词：

- 显式因果关系：包含因果连接词
- 隐式因果关系：不包含因果连接词

- 因果数目：

- 1→1：1个原因导致1个结果
- 1→多：1个原因导致多个结果
- 多→1：多个原因导致多个结果
- 多→多：多个原因交叉导致的多个结果

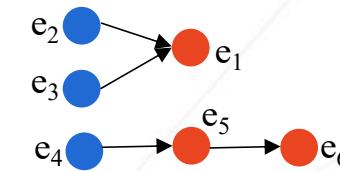
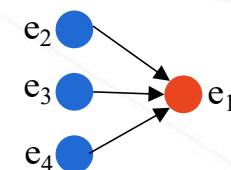
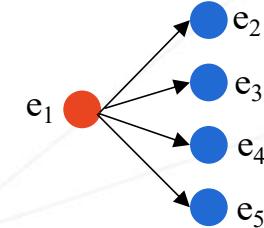
复杂因果

<e1>I don' t believe in myself</e1>
due to <e2>my anxiety</e2>.

After <e1> the discard and break up</e1>,
i was so <e2> confused </e2>,
<e3> heart broken </e3>,
<e4> horribly jealous </e4>,
<e5> angry </e5>.

I <e1> haven't been sleeping well </e1>
because i'm <e2> afraid of the dark </e2>,
<e3> sensitive to noise at night </e3>,
and <e4> my thoughts run wild </e4>.

I am so <e1> sad </e1>,
because my <e2> parents got divorced </e2>,
and my little <e3> brother got a fever </e3>.
I am so <e4> tired </e4> that <e5> have
no time to finish my homework </e5>.
So the <e6> teacher criticized me </e6>



目前因果关系抽取的研究主要集中在简单因果抽取，针对单对因果抽取涌现了各种各样的模型，对复杂因果研究较少。

图：因果关系的四种类型



研究挑战



□ 为什么对复杂因果研究较少?

- 缺少数据集：目前的数据集用于简单因果抽取，用于复杂因果抽取的**数据集规模小**
- 包含多个论元：论元之间的**交互关系复杂**，研究困难
- 论元长度复杂多样：复杂因果的论元涉及**单词、短语以及分句**等多种语法单元，难以用统一的模板衡量
- 关系跨度远：论元之间的**跨度增大**，增加了模型理解文本的难度，很难做出决策
- 无显式连接词：隐式因果占大多数，**缺乏因果关系连接词**，对抽取造成困难

I am so <e1>sad</e1>, because my <e2>parents got divorced</e2> last month, and my little <e3>brother got a fever</e3>, I must take him to the hospital. I'm so tired that have no time to finish my homework. So the <e4>teacher criticized me</e4> today. I want to <e5>commit suicide</e5>, because I can not breathe in such a <e6>high pressure</e6>, as well as the <e7>loneliness</e7>.

Cause-Effect((e2, e1), (e3, e1), (e4, e1), (e6, e5), (e7,e5))

I got so **depressed** **because** I **lost my favorite toy**.

{[lost my favorite toy], C-E, [depressed]}

[my parents got divorced] last month, but they didn't tell me. [I felt so sad] when I found their divorce certificate.

{[my parents got divorced], C-E, [I felt so sad]}



相关工作



因果关系抽取数据集调研

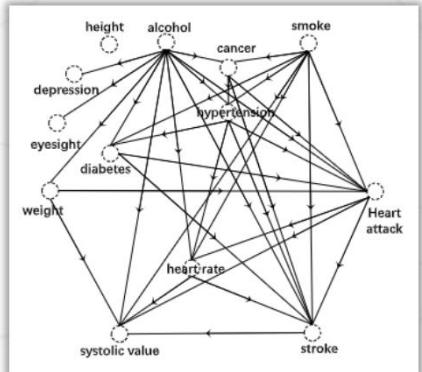
数据集	规模	简单因果	复杂因果	语言	备注
SemEval2010 Task-8	10717	1331	0	English	
BECause	5380	1634	587	English	
ADE	23516	7100	0	English	
SCIFI	5235	924	346	English	Task-8数据集重新标注
CREST	28878	28878	0	English	融合多个数据集
AItex	81325	6289	0	English	
Chinese Emergency Corpus	332	0	200	Chinese	
ECE Corpus(Gui et al 2016a)	2105	2046	59	Chinese	
航空安全	11392	11302	0	Chinese	955篇航空事故文本

现状：目前用于因果关系研究的数据集包含**航空安全、医疗、经济、新闻、情绪分析以及紧急事件等**领域，但是现有数据集主要针对简单因果关系，**复杂因果关系数据集规模小**。



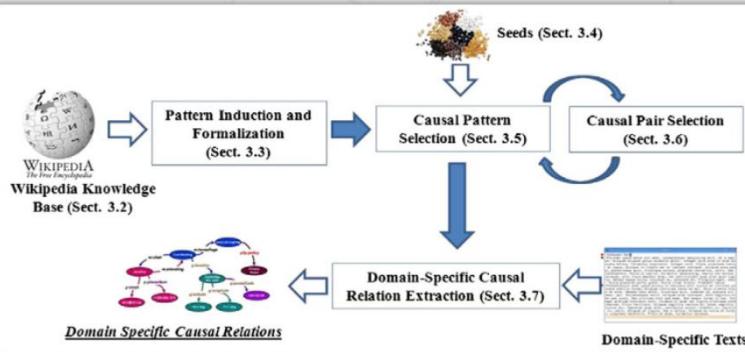
➤ 简单因果关系抽取

Causal Graph^[2]



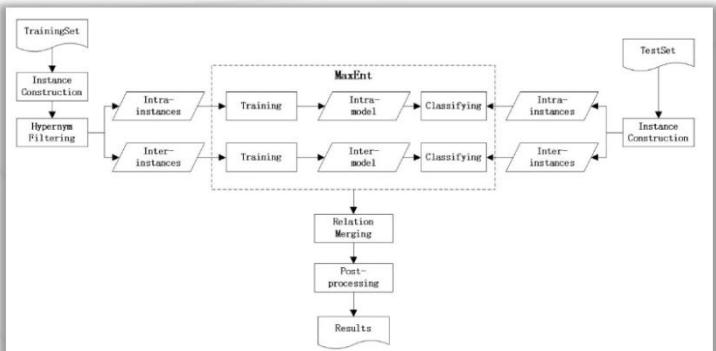
COMPUT BIOL MED 2019

Causal Pattern^[3]



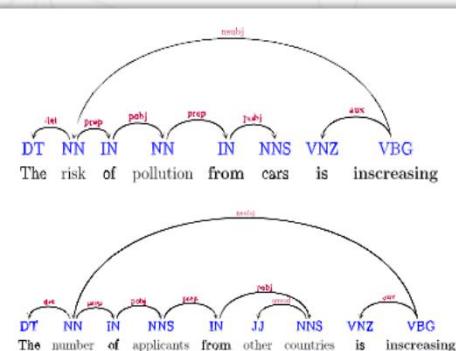
DATA KNOWL ENG 2013

Maximum Entropy^[5]



Database 2020

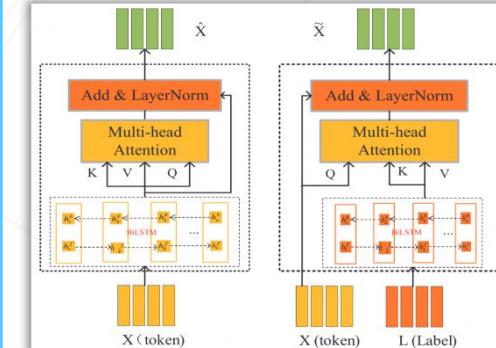
RHNB^[4]



Neurocomputing 2016

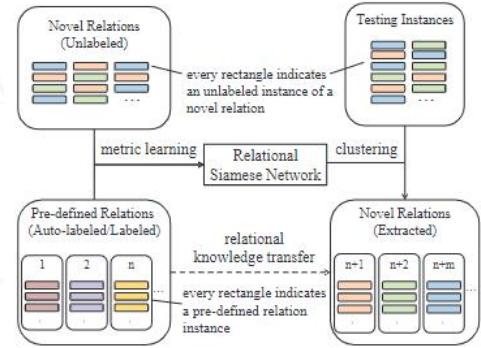
深度学习

CMAN^[6]



IJCAI 2020

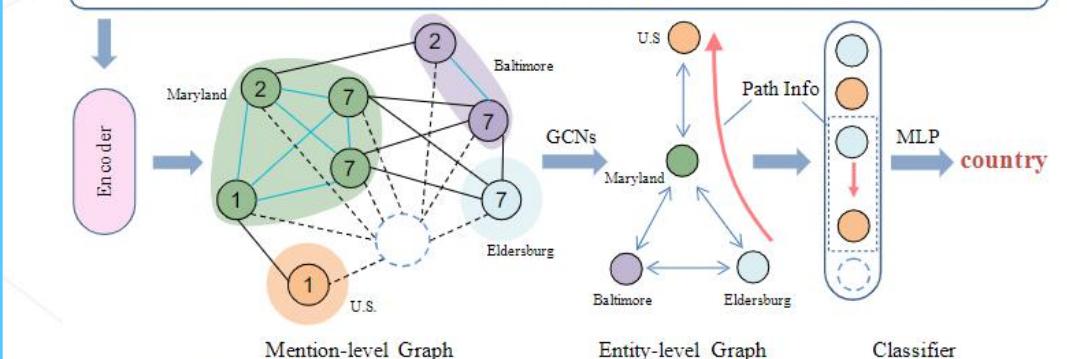
RSN^[8]



EMNLP 2019

GAIN^[7]

[1] Elias Brown ... was a U.S. Representative from Maryland. [2] Born near Baltimore, Maryland, ...
[7] He died near Baltimore, Maryland, and is interred in a private cemetery near Eldersburg, Maryland.

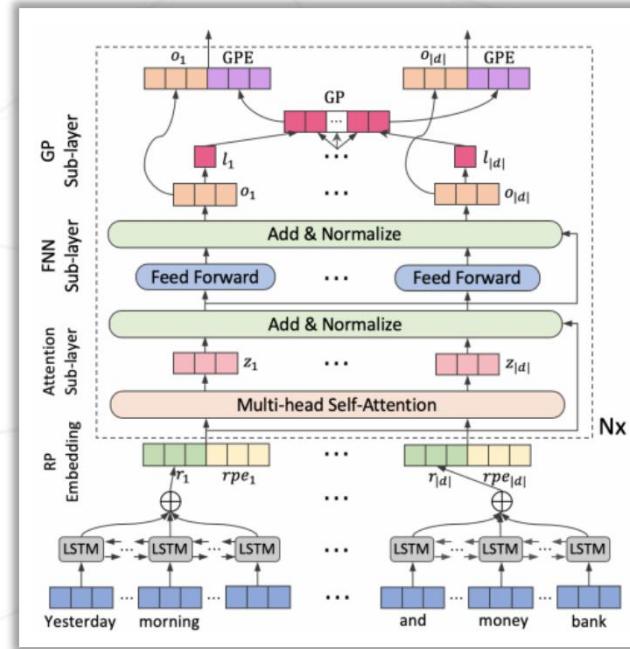


1 2 7 1 Mention Node 2 Document Node 3 Intra-Entity Edge — Inter-Entity Edge -- Document Edge

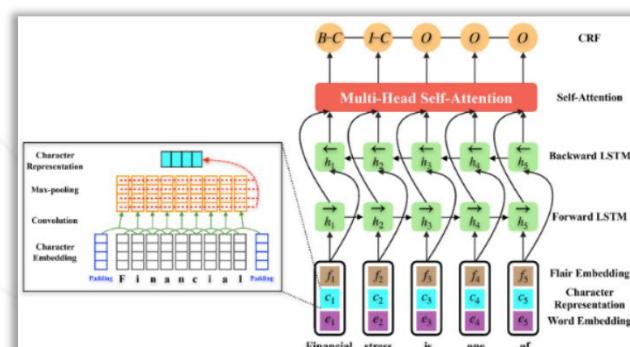
EMNLP 2020



➤ 复杂因果关系抽取

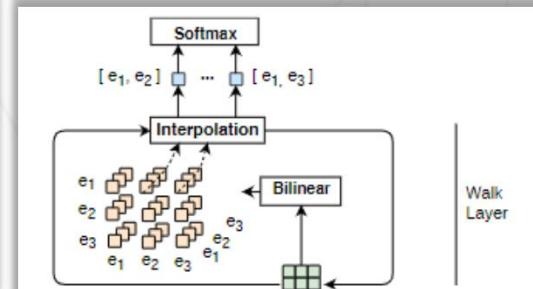
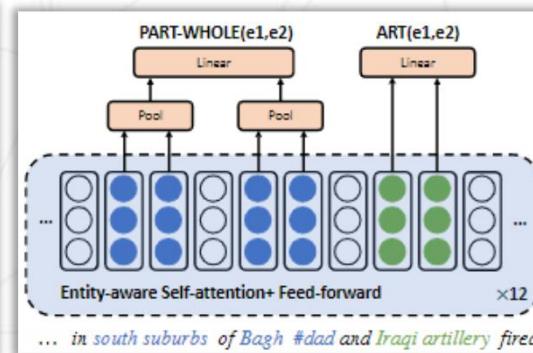


RTHN^[9] IJCAI 2019



SCITE^[11] Neurocomputing 2021

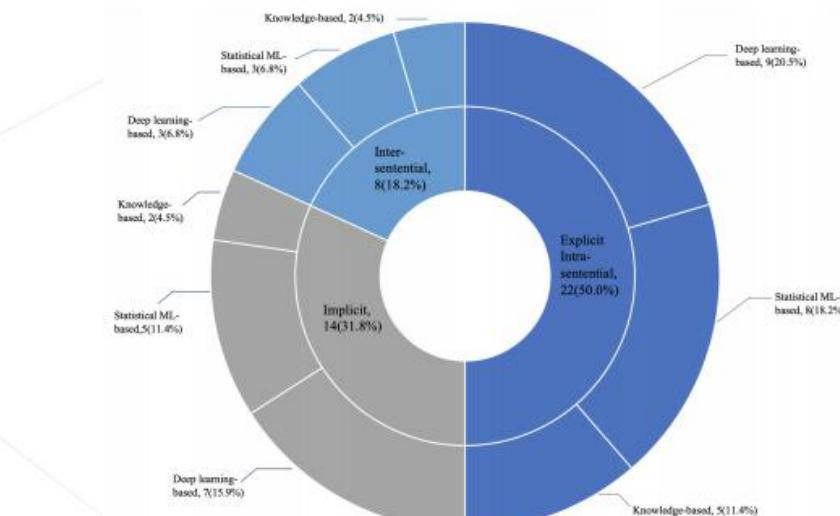
Transformer^[9] ACL 2019



Walk-based^[9]

□ 目前，研究复杂因果的工作较少，主要分为关系推理、分割文本以及Tag2Triplet的方法

- 关系推理对因果关系借鉴较少
- 分割文本不利于句子之间、跨句因果抽取
- Tag2Triplet需要大量的人工标注工作





解决思路



■ RedditCE数据集构建

采用网络爬虫技术，基于社交媒体Reddit抓取关于12种精神疾病的5w+条文本数据，对其进行筛选、清洗、标注，最终形成35202条复杂因果关系数据集。



I am so **<e1>sad</e1>**, because my **<e2>parents got divorced</e2>** last month, and my little **<e3>brother got a fever</e3>**, I must take him to the hospital. I'm so tired that have no time to finish my homework. So the **<e4>teacher criticized me</e4>** today. I want to **<e5>commit suicide</e5>**, because I can not breathe in such a **<e6>high pressure</e6>**, as well as the **<e7>loneliness</e7>**.

label: cause-effect((e2,e1),(e3,e1),(e4,e1),(e6,e5),(e7,e5))

type: 3

显/隐: 2/3

标注样例

	Single		Multiple		
	train	test	train	test	no cause
RedditCE	21553	5446	18942	4716	15455

1→1	1→N	N→1	N→N
11544	864	1071	6288

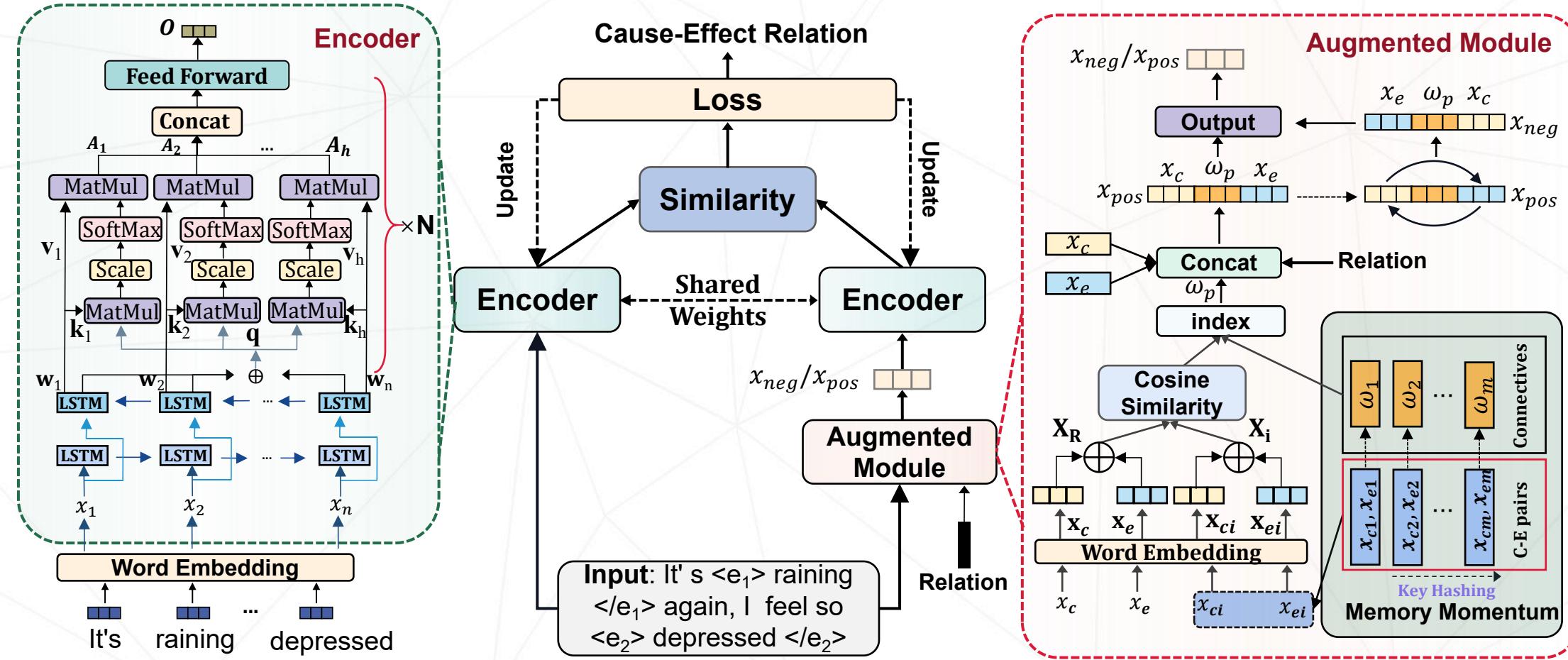


解决思路



■ 基于对比表示学习的复杂因果关系抽取模型

使用对比表示学习框架，融合BiLSTM-Transformer的深度神经网络学习多对原因和结果之间的复杂交互关系；使用Memory Momentum构造正负样例

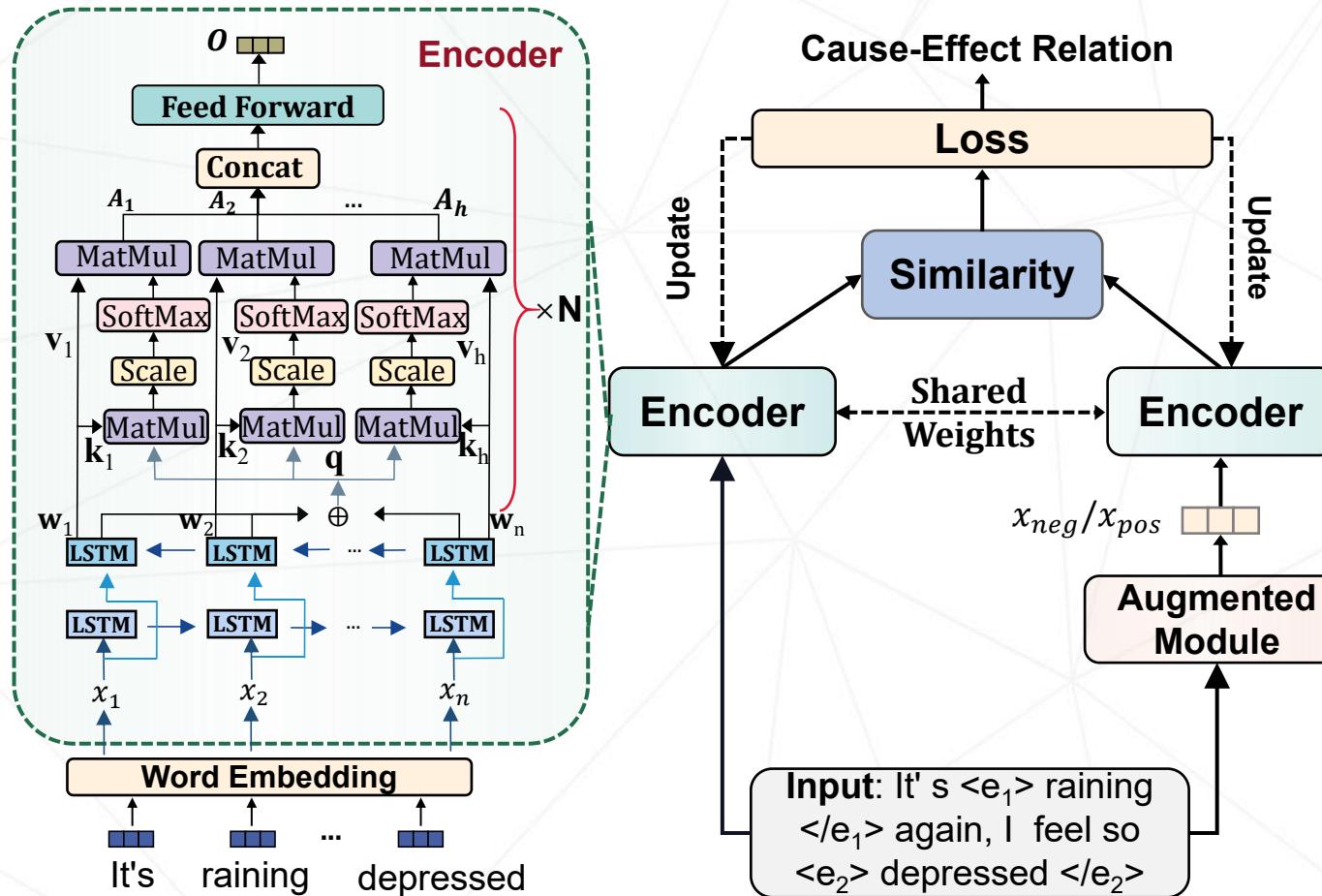




解决思路



■ 基于对比表示学习的复杂因果关系抽取模型



- BiLSTM获取分句级别的文本信息获取分句中单词的上下文信息，获得文本的句子级表示
- Transformer捕获BiLSTM所有时间步长中重要的信息。多头自注意力可以更加关注多对因果之间的信息。
- 增强模块进行数据增强的同时生成正负样例，从而满足对比学习框架的学习需求

$$\vec{h}_t = LSTM(\vec{h}_{t-1}, x_t, c_{t-1})$$

$$\underline{h}_t = LSTM(\underline{h}_{t-1}, x_t, c_{t-1})$$

$$w_i = [\vec{h}_t, \underline{h}_t]$$

$$z_i = Concat(A_1, \dots, A_h)W_O$$

$$\mathbf{e}_i = ReLU(z_i W_1 + b_1)W_2 + b_2$$

$$\mathbf{o}_i = Normalize(\mathbf{e}_i + \mathbf{q}_i)$$

$$\mathbf{q}_i = ReLU(w_i W_Q)$$

$$\mathbf{k}_i = ReLU(w_i W_K)$$

$$\mathbf{v}_i = ReLU(w_i W_V)$$



解决思路



➤ 正负样例构造模块

输入: 文本序列 $X = \{\dots, \langle e_i \rangle x_c, \langle /e_i \rangle, \dots, \langle e_j \rangle x_e, \langle /e_j \rangle\}$

输出: 正负样例

Step 1: 初始化 Memory Momentum, 以<原因, 结果, 连接词>形式存储;

Step 2: 获得输入中的因果对 $R = (x_c, x_e)$, 并连接 $X_R = \text{Concat}(x_c, x_e)$

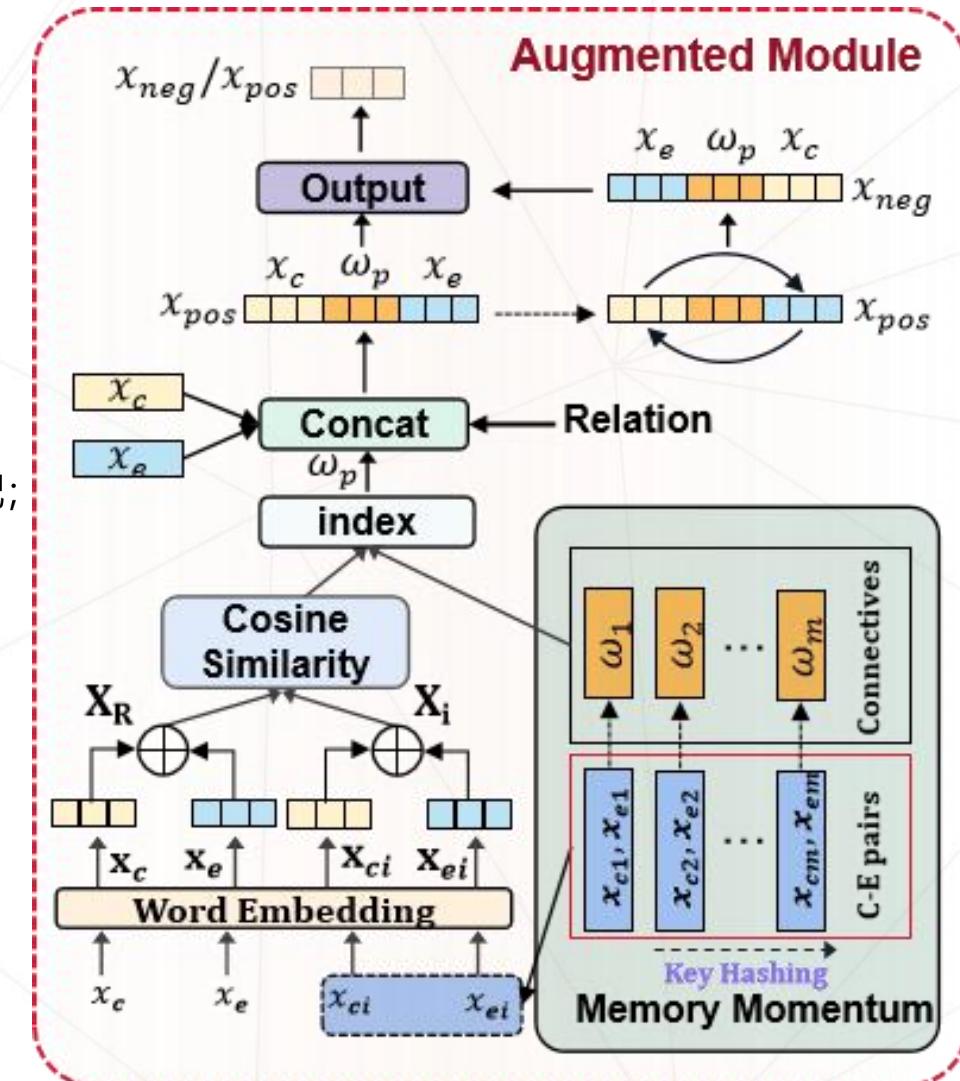
Step 3: 与 Memory Momentum 中的论元对 $X_i = \text{Concat}(x_{ci}, x_{ei})$ 进行相似度匹配;

Step 4: 获得最合适连接词构造正样例 $X_{pos} = \langle x_c, \omega_p, x_e \rangle$

Step 5: 调换正样例的因果位置得到负样例 $X_{neg} = \langle x_e, \omega_p, x_c \rangle$

$$S_i = \frac{\mathbf{X}_R \cdot \mathbf{X}_i}{\|\mathbf{X}_R\| \cdot \|\mathbf{X}_i\|}$$

$$p = \text{index}(\max_{i \in [0, m]} (S_i))$$





实验结果

➤ 与现有模型的对比结果

选取涵盖CNN, BiLSTM, Attention在内的5个Baseline，并对深度模型进行替换，进行性能评估，结果如下：

Model	RedditCE				SCIFI				CREST	
	1 → 1		N → N		1 → 1		N → N		1 → 1	
	F1 (%)	Acc (%)								
TextCNN (Li and Mao, 2019)	71.26	77.14	-	-	87.57	96.87	-	-	71.64	81.92
RCNN (Lai et al., 2015)	70.80	78.43	75.89	71.12	89.57	96.87	73.08	98.08	59.13	73.96
PCNN (Lai et al., 2015)	78.36	71.48	-	-	79.75	80.48	-	-	47.67	36.00
BiLSTM (Bekoulis et al., 2018)	79.35	84.21	-	-	88.52	97.52	-	-	67.99	79.46
BiLSTM-ATT (Zhao et al., 2021) (Li et al., 2021)	84.29	93.71	79.58	80.19	88.33	97.09	73.87	94.85	58.84	74.36
CL-Transformer	67.10	50.37	68.66	58.99	84.57	81.63	75.60	61.29	68.52	57.39
CL-CNN	82.39	80.37	75.81	77.54	92.68	90.82	93.10	91.40	76.19	76.64
CL-RNN	89.65	89.44	87.60	88.26	94.83	93.87	95.72	94.62	80.35	80.79
N2NCause	93.26	93.33	90.09	90.72	96.66	95.92	96.55	95.70	83.79	83.74

¹红色表示该项指标的最好结果；蓝色表示次好结果

²CL-Transformer/CNN/RNN表示模型采用的深度学习模型分别是Transformer/CNN/RNN



实验结果



➤ 不同损失函数对实验结果的影响

对比损失: $\mathcal{L}_{con} = Y \cdot \frac{1}{2}d^2 + (1 - Y) \cdot \frac{1}{2} \max(0, (t - d))^2$

MSE: $\mathcal{L}_{mse} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$ ➔ 有监督对比损失: $\mathcal{L}_{scl} = (1 - \lambda)L_{con} + \lambda L_{mse}$

Model	Loss	RedditCE				SCIFI			
		1 → 1		N → N		1 → 1		N → N	
		F1 (%)	Acc (%)						
CL-RNN	MSE	90.31	90.19	86.86	87.68	92.85	91.84	92.72	91.40
	MAE	82.21	81.48	83.11	85.21	94.02	92.86	94.55	93.55
	SCL	89.65	89.44	87.60	88.26	94.83	93.87	95.72	94.62
CL-Transformer	MSE	63.09	70.72	62.47	69.71	90.09	88.78	90.91	89.24
	MAE	68.73	75.74	59.00	73.62	92.59	91.84	91.42	90.32
	SCL	67.10	50.37	68.66	58.99	84.57	81.63	75.60	61.29
CL-CNN	MSE	74.32	75.56	71.36	74.64	91.23	89.79	94.44	93.54
	MAE	73.53	72.40	73.65	75.22	83.50	82.65	76.60	76.34
	SCL	82.39	80.37	75.81	77.54	92.68	90.82	93.10	91.40
N2NCause	MSE	91.59	91.67	89.51	90.14	87.72	85.71	91.43	90.32
	MAE	70.25	75.37	74.33	72.17	91.74	90.82	92.45	91.40
	SCL	93.26	93.33	90.09	90.72	96.66	95.92	96.55	95.70

有监督对比损失
性能提升明显



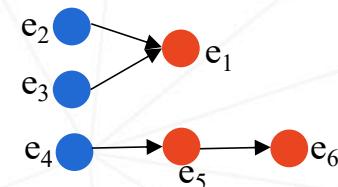
案例分析



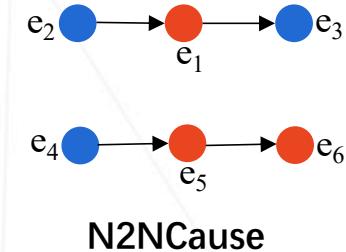
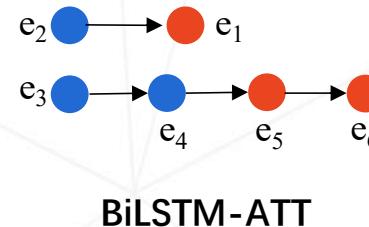
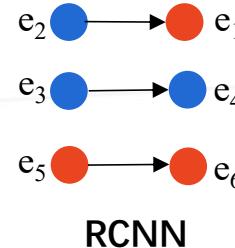
➤ 不同模型抽取复杂因果关系结果示例

为了展示模型对复杂因果关系抽取性能的提升，选取RCNN、BiLSTM-ATT和N2NCause三个模型，**多→多的因果关系类型**作为测试案例进行分析，结果如下：

I am so <e₁> sad </e₁>, because my <e₂> parents got divorced </e₂>, and my little <e₃> brother got a fever </e₃>. I am so <e₄> tired </e₄> that <e₅> have no time to finish my homework </e₅>. So the <e₆> teacher criticized me </e₆>



Ground Truth



- ✓ RCNN仅识别出2对，错误的认为e₃和e₄也具有因果关系，忽略了e₄和e₅， e₃和e₁的关系
- ✓ BiLSTM-ATT错误的识别出e₃和e₄之间的因果关系，但是忽略了e₁和e₃的关系
- ✓ N2NCause忽略了e₁与e₃之间的方向性



➤ 模型对方向是否敏感

为了展示模型对方向性是否敏感，选取两个案例，对CL-RNN、CL-Transformer和N2NCause三个模型进行方向敏感性分析，结果如下：

Case 1

< e_1 > the vascular dilatation </ e_1 > was caused by < e_2 > the sympathetic </ e_2 > and < e_3 > the course </ e_3 > of vaso-motor nerves.

$e_2 \rightarrow e_1$
 $e_3 \rightarrow e_1$

generally it appears that most of < e_1 > the damage </ e_1 > was caused by < e_2 > the winds </ e_2 > and < e_3 > the rough seas </ e_3 >

$e_2 \rightarrow e_1$
 $e_3 \rightarrow e_1$

Positive	CL-RNN	CL-Transformer	N2NCause
----------	--------	----------------	----------

$e_2 \rightarrow e_1$	✗	✓	✓
$e_3 \rightarrow e_1$	✗	✓	✓

Positive	CL-RNN	CL-Transformer	N2NCause
----------	--------	----------------	----------

$e_2 \rightarrow e_1$	✗	✓	✓
$e_3 \rightarrow e_1$	✓	✓	✓

Negative	CL-RNN	CL-Transformer	N2NCause
----------	--------	----------------	----------

$e_1 \rightarrow e_2$	✓	✓	✗
$e_1 \rightarrow e_3$	✓	✓	✗

Negative	CL-RNN	CL-Transformer	N2NCause
----------	--------	----------------	----------

$e_1 \rightarrow e_2$	✗	✓	✗
$e_1 \rightarrow e_3$	✗	✓	✗

¹ ✓ means that the triplet is recognized as a cause-effect pair; ✗ means that there is no causal relationship.



下一步完善目标包括现有**实验结果扩充、数据集扩建以及方法提升**三方面

2022.08

实验扩充

- 参数选取消融实验
- Memory Momentum改进实验 (**连接词的丰富**)
- 样本复杂度对实验结果的研究 (**自步学习**)
- 目标期刊：TKDE

2022.11

数据集扩建

- 数据集规模扩大 (目标数据量**10w+**)
- 数据集完善度 (包括**论元标注和因果关系**的完备程度)
- 提升复杂因果占比(**70%**)

2022.12

方法提升

- 对比表示学习样本丰富度提升 (主动对比学习)
- 在精神疾病成因分析的具体应用案例
- 解耦表示学习和因果关系抽取的结合



参考文献



- [1]. WANG Hong, ZHU Han, LIN Haizhou. Research on causality extraction of civil aviation accident. *Computer Engineering and Applications*, 2020, 56 (11) : 265-270.
- [2]. An N, Xiao Y, Yuan J, et al. Extracting causal relations from the literature with word vector mapping[J]. *Computers in biology and medicine*, 2019, 115: 103524.
- [3]. Ittoo A, Bouma G. Minimally-supervised learning of domain-specific causal relations using an open-domain corpus as knowledge base[J]. *Data & Knowledge Engineering*, 2013, 88: 142-163.
- [4]. Zhao S, Liu T, Zhao S, et al. Event causality extraction based on connectives analysis[J]. *Neurocomputing*, 2016, 173: 1943-1950.
- [5]. Gu J, Qian L, Zhou G. Chemical-induced disease relation extraction with various linguistic features[J]. *Database*, 2016, 2016.
- [6]. Zhao S, Hu M, Cai Z, et al. Modeling dense cross-modal interactions for joint entity-relation extraction[C]//Proceedings of the Twenty-Ninth International Conference on International Joint Conferences on Artificial Intelligence. 2021: 4032-4038.
- [7]. Zeng S, Xu R, Chang B, et al. Double graph based reasoning for document-level relation extraction[J]. arXiv preprint arXiv:2009.13752, 2020.
- [8]. Wu R, Yao Y, Han X, et al. Open relation extraction: Relational knowledge transfer from supervised data to unsupervised data[C]//Proceedings of the 2019 *Conference on Empirical Methods in Natural Language Processing* and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP). 2019: 219-228.
- [9]. Wang H, Tan M, Yu M, et al. Extracting multiple-relations in one-pass with pre-trained transformers[J]. arXiv preprint arXiv:1902.01030, 2019.
- [10]. Dasgupta T, Saha R, Dey L, et al. Automatic extraction of causal relations from text using linguistically informed deep neural networks[C]//*Proceedings of the 19th Annual SIGdial Meeting on Discourse and Dialogue*. 2018: 306-316.
- [11]. Li Z, Li Q, Zou X, et al. Causality extraction based on self-attentive BiLSTM-CRF with transferred embeddings[J]. *Neurocomputing*, 2021, 423: 207-219.
- [12]. Xia R, Zhang M, Ding Z. RTHN: A RNN-transformer hierarchical network for emotion cause extraction[J]. arXiv preprint arXiv:1906.01236, 2019.
- [13]. Christopoulou F, Miwa M, Ananiadou S. A walk-based model on entity graphs for relation extraction[J]. arXiv preprint arXiv:1902.07023, 2019.

THANKS

请各位老师多批评指正！