
프롬프트 학습 기반 한국어 개체 중의성 해결

민진우, 나승훈
인지컴퓨팅 연구실
전북대학교



목차

- 개체 연결
- 관련 연구
- 프롬프트 학습 기반 한국어 개체 중의성 해결
- 실험결과



개체 연결

- 정의(두 가지의 과정으로 구분)
 - 멘션 탐지 : 문서에서 등장하는 지명, 인명, 기관명 등을 나타내는 개체 표현인 개체 멘션을 찾는 과정
 - 중의성 해결 : 개체 멘션이 가질 수 있는 후보 엔티티 중 단 하나의 엔티티로 연결하는 과정

이들은 개미들에게 손실을 떠넘기면서 큰 수익을 챙겨간다.

개미

개미는 개미과에 속하는 진사회성 곤충의 총칭으로, 말벌상과, 벌과 더불어 벌목에 속한다.

투자자

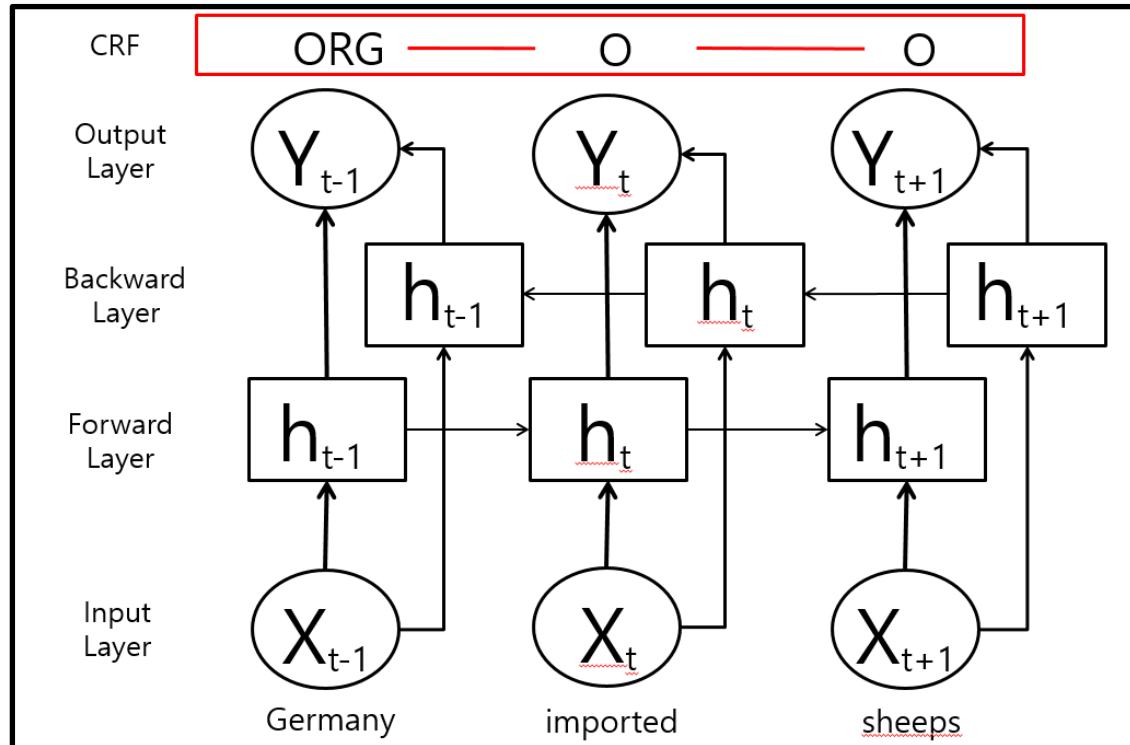
투자자는 주식이나 채권·파생상품·부동산·통화·상품 등에 투자하는 개인 또는 법인을 말한다.

개미(소설)

《개미》는 프랑스의 작가 베르나르 베르베르의 등단작이자 가장 유명한 작품이다.



Bidirectional LSTM-CRF Models for Sequence Tagging (Zhiheng Huang et al, '16)

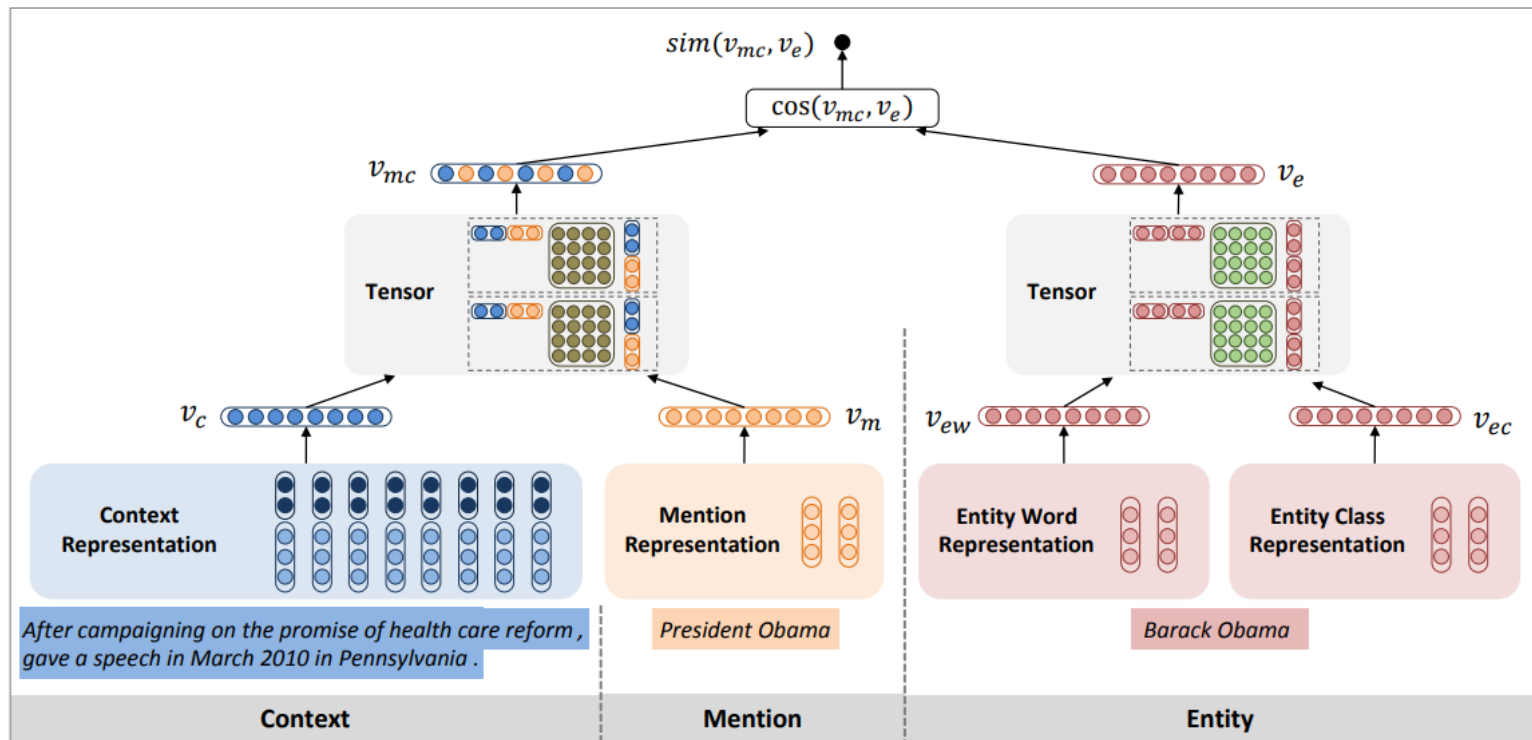


- Bi-LSTM-CRF 기반의 개체명 인식

- 주어진 sequence에 대해서 label을 부여하는 sequence labeling(순차 태깅) 방식
- 개체명 인식은 멘션 추출 과정의 또 다른 이름



Modeling Mention, Context and Entity with Neural Networks for Entity Disambiguation (Sun, IJCAI '15)

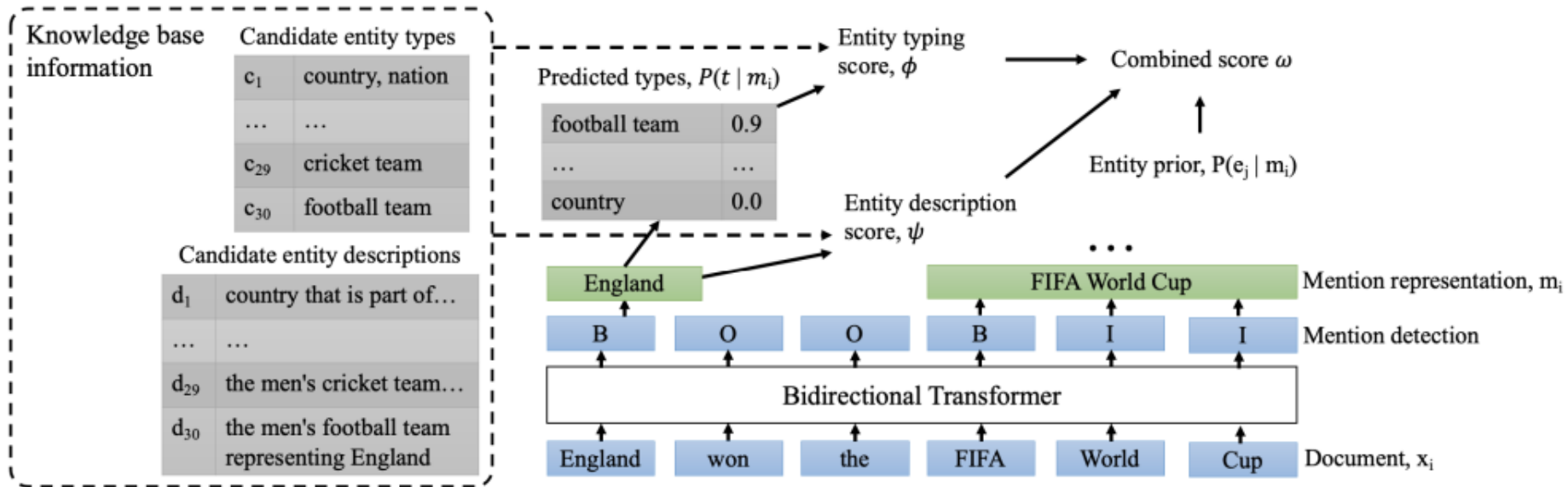


중의성 해결 모델

- 신경망을 통해 mention, context, 후보 엔티티에 대한 vector를 구성하고 결합된 mention-context vector와 후보 entity vector 사이의 유사도를 구하고 유사도가 가장 높은 entity를 선택



ReFinED: An Efficient Zero-shot-capable Approach to End-to-End Entity Linking (Tom Ayoola et al, '21)

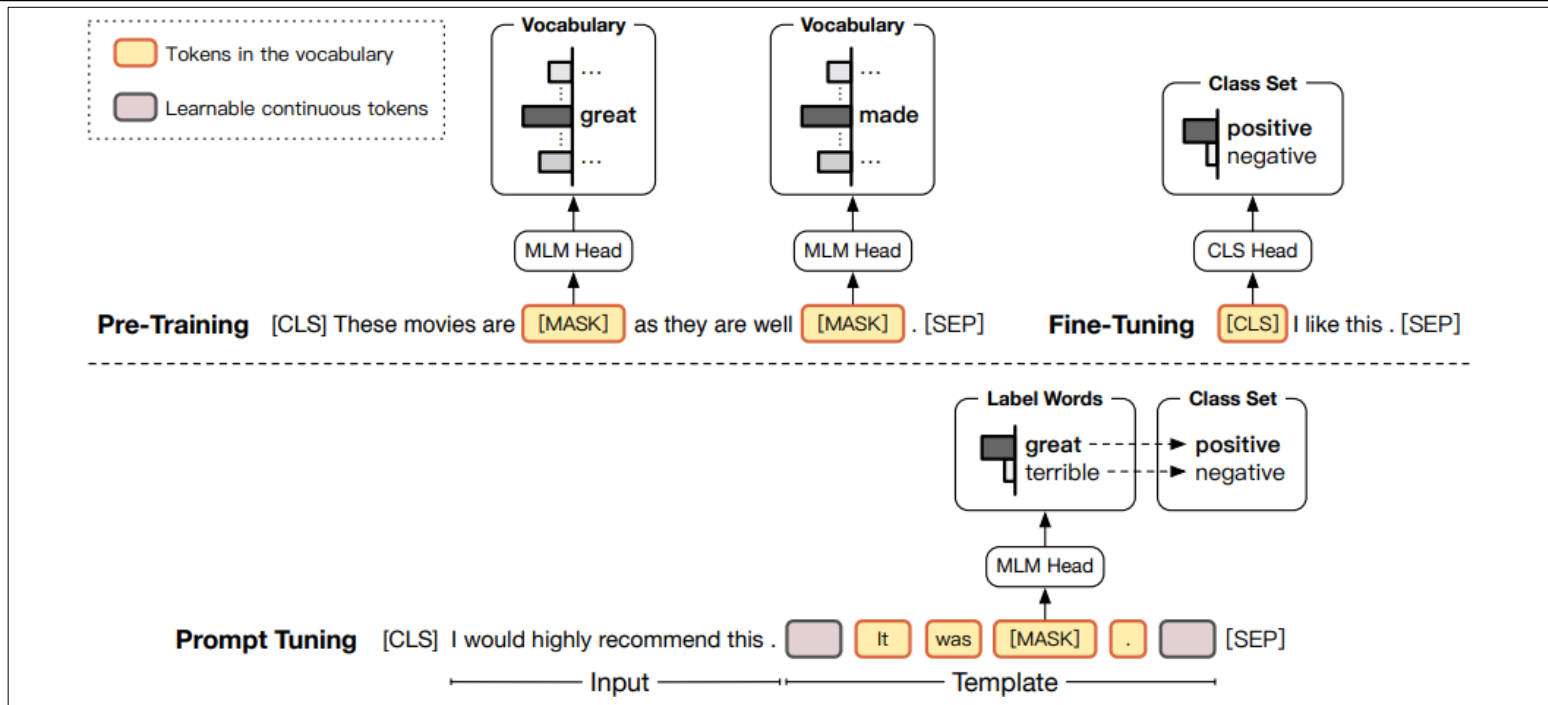


- 멘션 추출 및 중의성 해결 통합 모델

- 문서 내의 모든 멘션에 대한 멘션 탐지(추출)와 엔티티 타이핑과 중의성 해결을 동시에 해결하는 end-to-end 방식의 개체 연결 모델
- 엔티티 타입 정보와 description 정보를 이용하여 중의성 해결



PTR: Prompt Tuning with Rules for Text Classification(Xu Han et al '22)

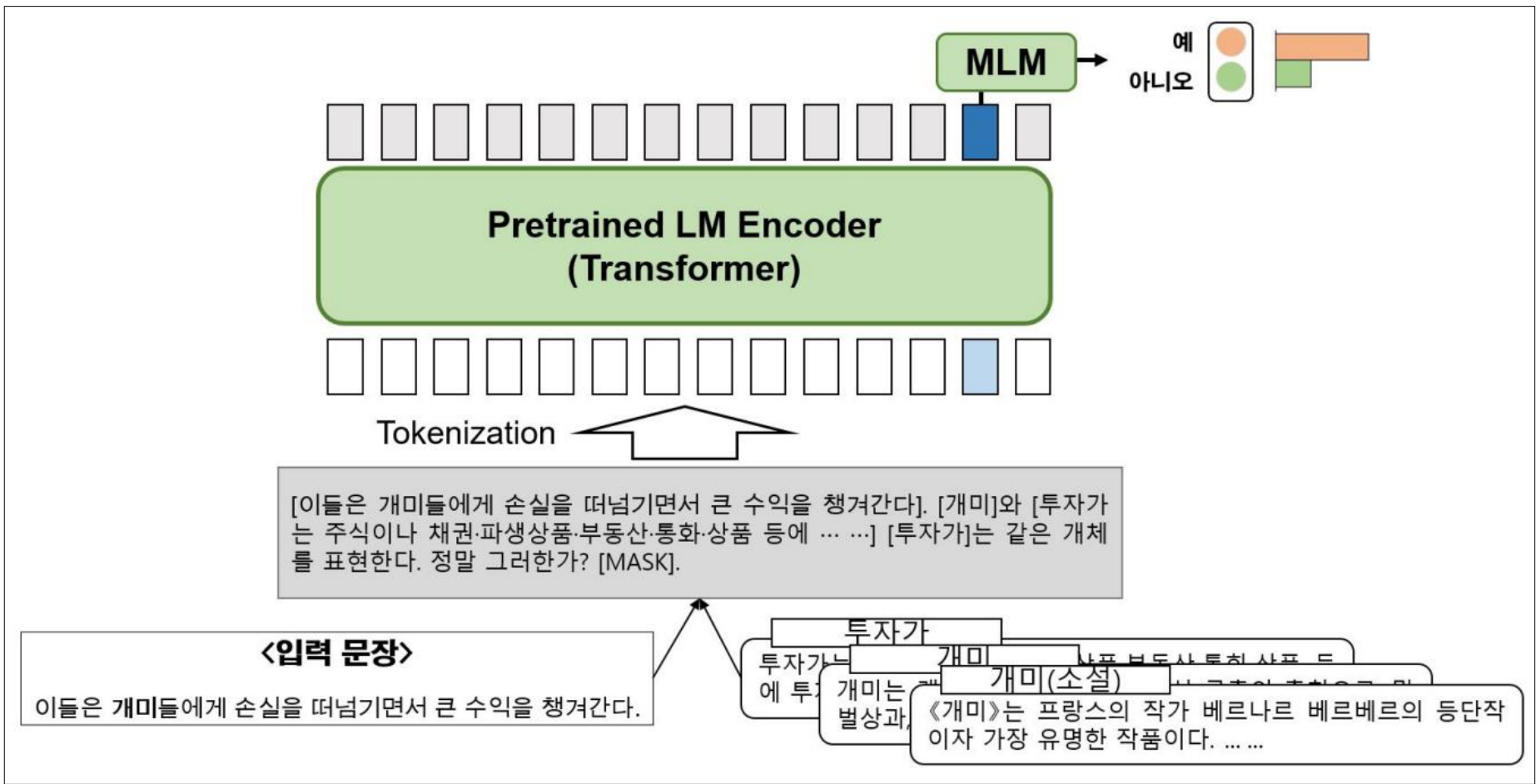


- PLM에 친화적인 튜닝 방식

- a)는 PLM(Pre-trained Language Model) 의 [CLS] 에 추가적인 선형 층을 쌓아 학습하는 방식
- b) 주어진 문장에 템플릿 함수를 통하여 프롬프트(최소 하나의 [MASK] 를 포함하는)를 PLM 에 넣고 MLM Head 로 학습하는 방식



프롬프트 학습 기반 한국어 개체 중의성 해결



- 템플릿 구성
 - 템플릿 정의

[Context]. [Mention] 과 [Entity Description] [Entity Title]은 같은 개체를 표현한다. 정말 그러한가? [MASK].

- [Context], [Mention]는 입력 문장과 문장에서 나타난 멘션, [Entity Description] 과 [Entity Title]은 각각 후보 엔티티에 대한 지식 베이스 상의 개체 설명문과 타이틀을 나타냄.
 - [MASK] 위치의 토큰은 출력 레이블에 매핑되는 토큰으로 이는 개체 멘션에 대한 후보 엔티티의 연관도를 표현
- Verbalizer는 프롬프트 기반 분류 학습에서 중요한 요소로 클래스의 레이블을 전체 레이블 단어장에 투영
 - 본 논문에서는 긍정/부정의 2개의 클래스이며 각각의 클래스에 대해 "예", "아니오"의 단일 단어만을 단어장으로 취함
 - 즉, [MASK] 위치의 단어는 최종 레이블에 연결되며 여기서 해당 후보 엔티티가 정답이면 Positive 레이블 단어로 "예"를 설정하고 부정적인 레이블 단어로 "아니오"를 설정



• 학습 및 디코딩

- 학습

- 출력층에서 [MASK]에 해당하는 토큰 위치의 벡터에 대해 Masked Language Model(MLM) 방식으로 각 레이블에 대한 확률을 얻어냄

$$p([\text{MASK}] = w \in V | T(x)) = \frac{\exp(v_w \cdot h_{[\text{MASK}]})}{\sum_{v_i \in V} \exp(v_i \cdot h_{[\text{MASK}]})}$$

- $h_{[\text{MASK}]}$ 는 언어 모델의 출력에서 [MASK]토큰 위치에 해당하는 은닉 표상이고 v_i 는 단어장의 i 번째 레이블에 대한 임베딩
- 둘의 내적을 통해 softmax 함수를 통해 각 레이블의 확률을 얻음. 교차 엔트로피 손실 함수를 통해 정답 레이블의 확률이 최대가 되도록 학습

- 디코딩

- 주어진 멘션에 대한 모든 후보 엔티티에 대해 긍정 레이블에 대한 점수를 모은 후 긍정 레이블의 점수가 최대가 되는 후보 엔티티를 선택하여 중의성 해결.
- 본 연구에서는 멘션에 후보 엔티티는 구축된 멘션-후보 엔티티 사전을 통해 얻어진 멘션의 후보 엔티티들을 빈도수 상위 30개로 제한하여 사용



실험 세팅

- 데이터 셋
 - 모두의 말뭉치 개체 연결 데이터 셋
 - 전체 문서 집합 중 1678개의 문서를 추출 후 학습셋 1378 문서, 개발셋 100 문서, 평가셋 200문서로 구성.
 - 사전 구축
 - 멘션-후보 엔티티 사전 구축을 위해 말뭉치 전체 문서 내 각 멘션에 대해 링크되어 있는 엔티티 문서의 빈도수를 카운트한 후 전체 후보 엔티티를 빈도수의 내림차순으로 정렬하여 구축
 - 3가지 세팅의 실험을 구성
 - 1) Bi-Encoder 구조 + fine-tuning
 - 2) Cross-Encoder 구조 + fine-tuning
 - 3) Cross-Encoder 구조 + 프롬프트 tuning
 - Bi-Encoder 구조 : 멘션-컨텍스트 인코더와 별도로 개체 정보를 인코딩하기 위한 별도의 인코더를 두어 인코딩. 내적하여 주어진 멘션에 대한 해당 후보 엔티티의 점수를 얻어냄
 - Cross-Encoder + fine-tuning : 멘션-컨텍스트와 후보 엔티티 description을 [SEP] 특수토큰을 두고 연결하여 인코딩



실험 결과

- 실험 결과

	ACC
Bi-Encoder + fine-tuning	87.60%
Cross-Encoder + fine-tuning	92.90%
Cross-Encoder + 프롬프트 tuning	93.70%

- 결과 분석

- 멘션-컨텍스트와 후보 엔티티 정보를 함께 인코딩하는 Cross-Encoder 기반 방법이 별도로 인코딩하는 Bi-Encoder 기반 방법 대비 5.3% 높은 성능을 보임
- 프롬프트 튜닝 방식이 fine-tuning 방식 대비 0.8%의 추가적인 성능 향상을 보임 프롬프트 방식이 태스크에 대한 이해도를 높여 좋은 성능을 보임을 확인

- 향후 연구

- 개체 타입, prior 정보 및 지식 그래프 정보 등을 활용하여 성능을 높이는 실험을 진행할 예정
- 프롬프트 방식을 End-to-End 개체 연결로의 확장 진행



Q&A

감사합니다.

