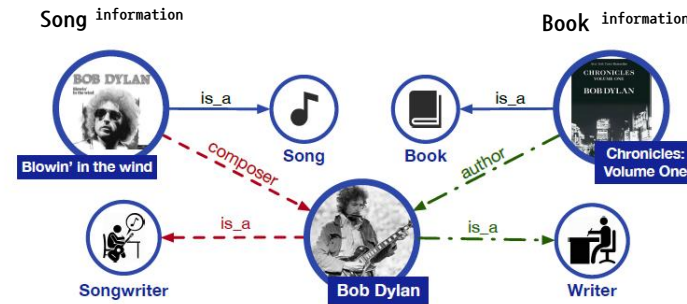


**ERNIE,  
Knowledge Enhanced Contextual Word Representations,  
and LUKE**

**Sunkyung Park**

- **Limitations of Existing Research**

- pre-trained LM neglects knowledge information for language understanding
- **(not included)** Entity Typing Assigning types to mentions in docs
- **(not included)** Relation Classification Predicts a relation between two entities in a sentence



*Bob Dylan* wrote *Blowin' in the Wind* in 1962, and wrote *Chronicles: Volume One* in 2004.

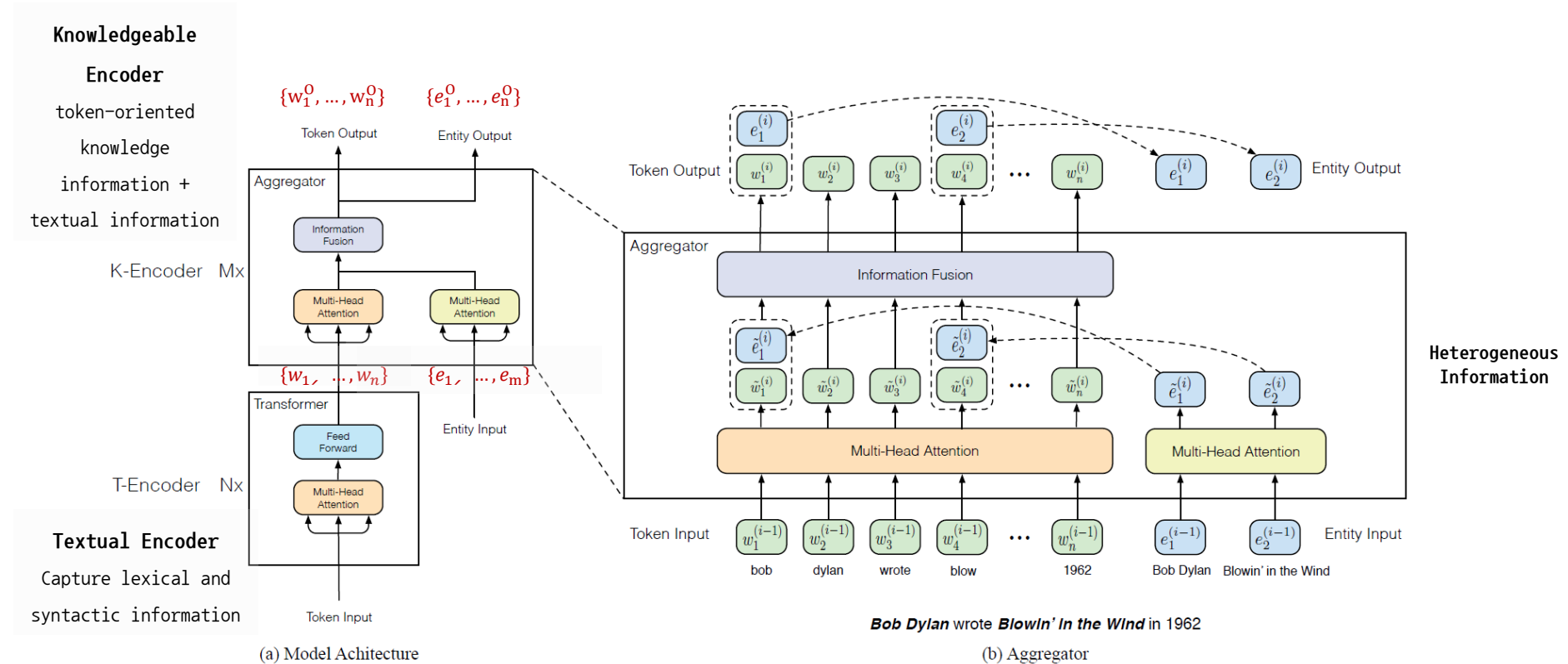
Existing pre-trained LM : UNK wrote UNK in UNK. → syntactically ambiguous

Rich Knowledge Information: Language Understanding → entity typing & relation classification

- **Challenges** knowledge into language representation models

- **Structured Knowledge Encoding**
  - How to extract and encode related informative facts in KGs
- **Heterogeneous Information Fusion**
  - How to **design** a special **pre-training objective** to fuse lexical, syntactic and knowledge information

- **Proposed Model** ERNIE
  - Recognize **mentions** in text → Align mentions to **entities** in KGs
  - **Informative entity embeddings**
    - Encode **graph structure** of KGs w/ **knowledge embedding**
  - Integrates **entity representations** in the knowledge module into layers of the **semantic module**
- **New objectives**
  - Randomly **mask** some of the **named entity alignments** in the input
  - Ask the model to select appropriate entities from KGs
  - Require models to aggregate **context** and **knowledge facts**
  - Predict **tokens** and **entities**



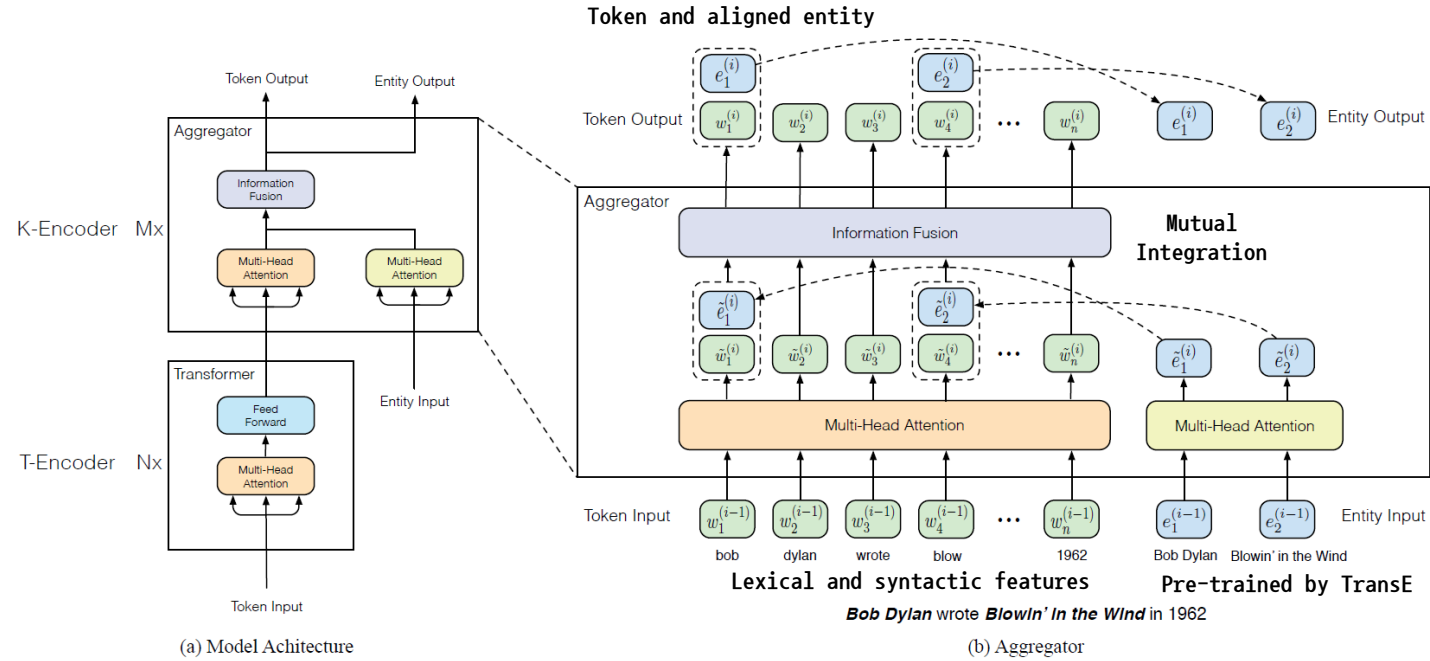
### 1. T-Encoder does

- Sum token · segment · positional embedding
- Compute  $\{w_1, \dots, w_n\}$  lexical and syntactic features

### 3. Output Embeddings ( $\{w_1^0, \dots, w_n^0\}, \{e_1^0, \dots, e_m^0\}$ ) will be used for specific tasks

### 2. Both $\{w_1, \dots, w_n\}$ and $\{e_1, \dots, e_m\}$ pre-trained by TransE fed into K-Encoder

- Knowledgeable Encoder > Information fusion layer



W/ corresponding entities

$$\begin{aligned}
 h_j &= \sigma(\tilde{W}_t^{(i)} \tilde{w}_j^{(i)} + \tilde{W}_e^{(i)} \tilde{e}_k^{(i)} + \tilde{b}^{(i)}), \\
 w_j^{(i)} &= \sigma(W_t^{(i)} h_j + b_t^{(i)}), \\
 e_k^{(i)} &= \sigma(W_e^{(i)} h_j + b_e^{(i)}).
 \end{aligned}$$

W/O corresponding entities

$$\begin{aligned}
 h_j &= \sigma(\tilde{W}_t^{(i)} \tilde{w}_j^{(i)} + \tilde{b}^{(i)}), \\
 w_j^{(i)} &= \sigma(W_t^{(i)} h_j + b_t^{(i)}).
 \end{aligned}$$

- Pre-training for Injecting Knowledge

- New pre-training task

- Randomly masks some token-entity alignments
    - Requires the system to predict all corresponding entities based on aligned tokens

dEA denoising entity auto-encoder

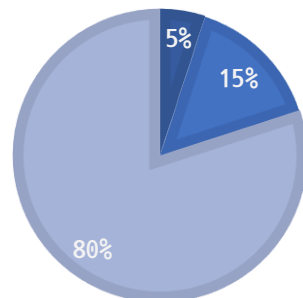
$$p(e_j|w_i) = \frac{\exp(\text{linear}(w_i^o) \cdot e_j)}{\sum_{k=1}^m \exp(\text{linear}(w_i^o) \cdot e_k)},$$

Given  $\{w_1, \dots, w_n\}$  token sequence and  $\{e_1, \dots, e_m\}$  corresponding entity sequence,

This is an aligned entity distribution for token  $w_i$

- Operations for dEA

- Replace w/ Entity with Another Random Entity
  - Mask token-entity alignments



- Tokens in Text

- MLM
  - NSP

- Fine-tuning for Specific Tasks

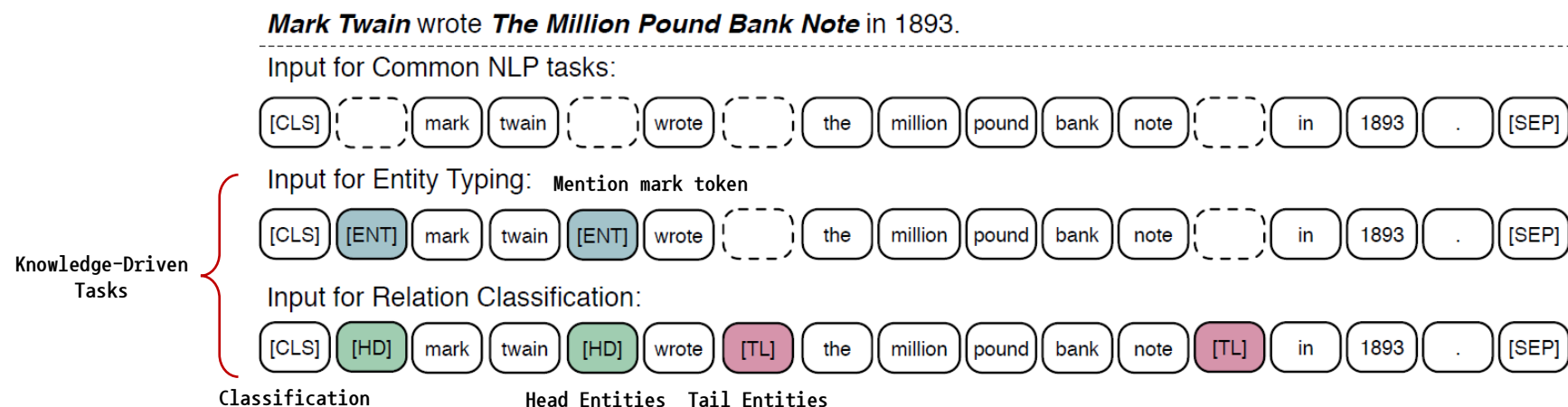


Figure 3: Modifying the input sequence for the specific tasks. To align tokens among different types of input, we use dotted rectangles as placeholder. The colorful rectangles present the specific mark tokens.

- Task 1. Relation Classification

- Modifies the input token sequence by adding two marktokens to highlight entity mention

- Task 2. Entity Typing

- Modified input sequence guides ERNIE to combine both context information and entity-mention information

## Pre-training Dataset

- Before pre-training ERNIE, **Adopts** the knowledge embeddings by TransE
- Multi-Tasks: NSP, MLM, and dEA
- Use Eng Wikipedia and Align text to Wiki-data Convert corpus into the formatted data

## (Downstream)Entity Typing

- (Task) Given an entity mention and its context, **Label** the entity mention w/ corresponding semantic types
- (Evaluate) Fine-tunes ERNIE on 2 datasets

- FIGER
  - Train set : Labeled
  - Test set : Annotated by Human
- Open Entity
  - Completely manually-annotated

Dataset	Train	Develop	Test	Type
FIGER	2,000,000	10,000	563	113
Open Entity	2,000	2,000	2,000	6

Table 1: The statistics of the entity typing datasets FIGER and Open Entity.

ACC: Prediction = Human annotations

Model	Acc.	Macro	Micro
NFGEC (Attentive)	54.53	74.76	71.58
NFGEC (LSTM)	55.60	75.15	71.73
BERT	52.04	75.16	71.63
ERNIE	<b>57.19</b>	<b>76.51</b>	<b>73.39</b>

SOTA on FIGER  
2 Bi-LSTMs

(Compare) ERNIE vs. BERT

(Acc) **External Knowledge** regularizes ERNIE to fitting the noisy labels

Precision Recall

Model	P	R	F1
NFGEC (LSTM)	68.80	53.30	60.10
UFET	77.40	60.60	68.00
BERT	76.37	70.96	73.56
ERNIE	<b>78.42</b>	<b>72.90</b>	<b>75.56</b>

1 Bi-LSTM

(Compare) BERT & ERNIE vs. Others

Pre-training models **make full use of** manually-annotated data for entity-typing

(Compare) ERNIE vs. BERT

**informative entities** help predict the labels

Table 3: Results of various models on Open Entity (%). precisely



**(Downstream)Relation Classification**

- (Task) Determine the correct relation between two entities<sup>Head and Tail entities</sup> in a sentence
- (Evaluate) Fine-tunes ERNIE on 2 datasets

- FewRel
  - W/O any null instances
- TACRED
  - W/ 80% null instances

Dataset	Train	Develop	Test	Relation
FewRel	8,000	16,000	16,000	80
TACRED	68,124	22,631	15,509	42

Table 4: The statistics of the relation classification datasets FewRel and TACRED.

Graph Conv

Model	FewRel			TACRED		
	P	R	F1	P	R	F1
CNN	69.51	69.64	69.35	70.30	54.20	61.20
PA-LSTM	-	-	-	65.70	64.50	65.10
C-GCN	-	-	-	69.90	63.30	66.40
BERT	85.05	85.11	84.89	67.23	64.81	66.00
ERNIE	88.49	88.44	<b>88.32</b>	69.97	66.08	<b>67.97</b>

Table 5: Results of various models on FewRel and TACRED (%).

**FewRel**

(Feature) Training data w/ not enough instances

(Compare) CNN vs. BERT & ERNIE

(F1) Pre-training model outperforms

(Compare) BERT vs. ERNIE

(F1) Fusing external knowledge effective

**TACRED**

(Compare) C-GCN vs. BERT

(F1) C-GCN utilizes entity mask strategy

(Compare) BERT vs. ERNIE

(R, F1) effectiveness of the knowledgeable module

- **Limitations of Existing Research**
  - Do not contain any **explicit ground** to real word entities → difficult to recover **factual knowledge**
  - Integrate external knowledge into task-specific models
- **(Good Stuff) KB** knowledge bases
  - Rich source
  - Complementary information to that in raw text
  - Encode **factual knowledge** impossible in infrequent mentions and long-range dependencies

- **Proposed Model** <sup>KAR</sup>
  - Method to insert **multiple KBs** into a large pretrained model w/ KAR <sup>Knowledge Attention and Recontextualization</sup>
  - Idea in **KAR**
    - Model explicitly **entity spans** in the input text
    - Use **entity linker** to retrieve relevant entity embeddings from a KB
  - Recontextualizes **entity-span representations** w/ **word-to-entity attention** → **word representation & all entity spans in the context**
  - **KAR is inserted between 2 Layers** in the middle of a pre-trained model.
  - Learns **the entity linkers** w/ self-supervision on unlabeled data
- **Benefits**
  - Leaves Original model **unchanged** → Swap out BERT for KnowBert in any downstream tasks.
  - KAR is **Lightweight**
  - **Easy to incorporate** additional KBs by inserting them at other locations

## KnowBert

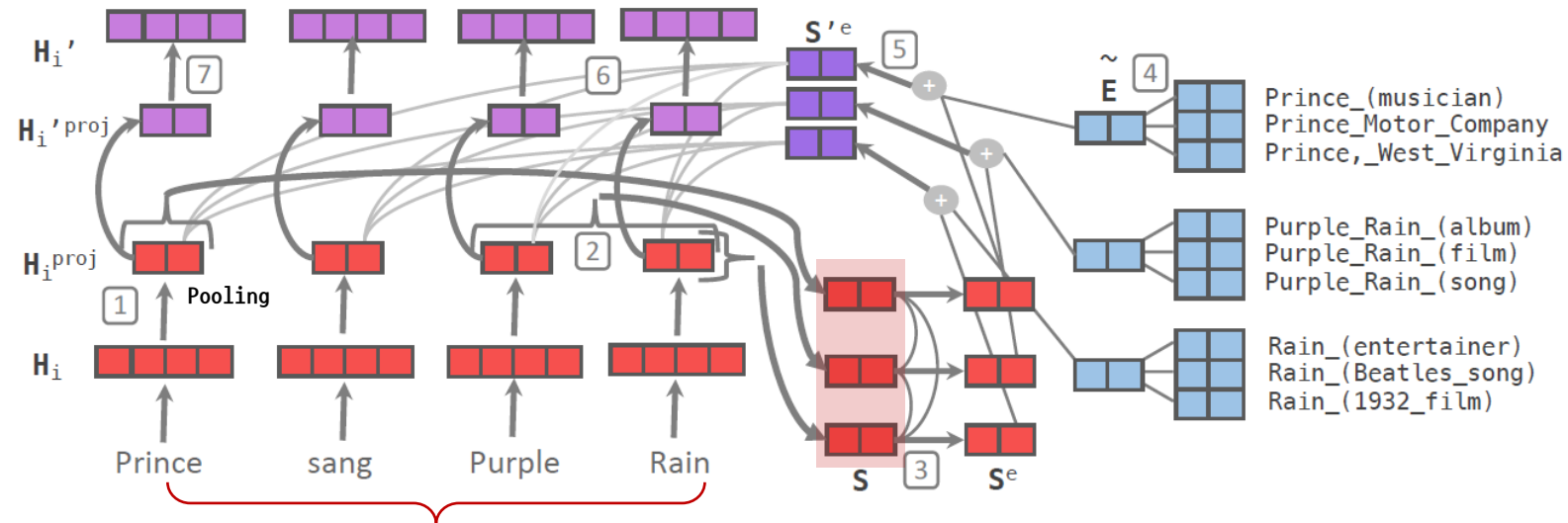
- (Definition) incorporates **knowledge bases** into pretrained **BERT** using the **KAR** Knowledge Attention and Recontextualization component
- **Knowledge Bases**
  - w/ **fixed** collection of  $K$  entity nodes  $\rightarrow$  computes entity embeddings  $e_k \in \mathbb{R}^E$
  - **entity embeddings** includes
    - KBs w/ a (**subj, rel, obj**) graph structure
    - KBs w/ only entity metadata **w/o** graph
    - KBs w/ a **graph** and **entity metadata**
- **Assumption**
  - The entities are not typed
  - KB accompanied by an **entity candidate selector**
- **Entity Candidate Selector**
  - **(input)** some text
  - **(output)** list of  $C$  potential **entity links**  $\rightarrow$  passed to a context dependent **entity linker** to **discriminate** the candidates
 

$$C = \{ \langle (\overset{\text{start / end indices of mention span}}{\text{start}_m, \text{end}_m}), (\overset{\text{candidate entities in KG}}{e_{m,1}, \dots, e_{m,M_m}}) \rangle \mid m \in 1 \dots C, e_k \in 1 \dots K \}.$$
- **Over-generate** potential candidates and add a **NULL entity** to each candidate list  $\rightarrow$  allows linker to **discriminate** between actual links and false positive candidates

KAR Knowledge Attention and Recontextualization component

- Mention-span representations

$$\mathbf{H}_i^{\text{proj}} = \mathbf{H}_i \mathbf{W}_1^{\text{proj}} + \mathbf{b}_1^{\text{proj}}.$$



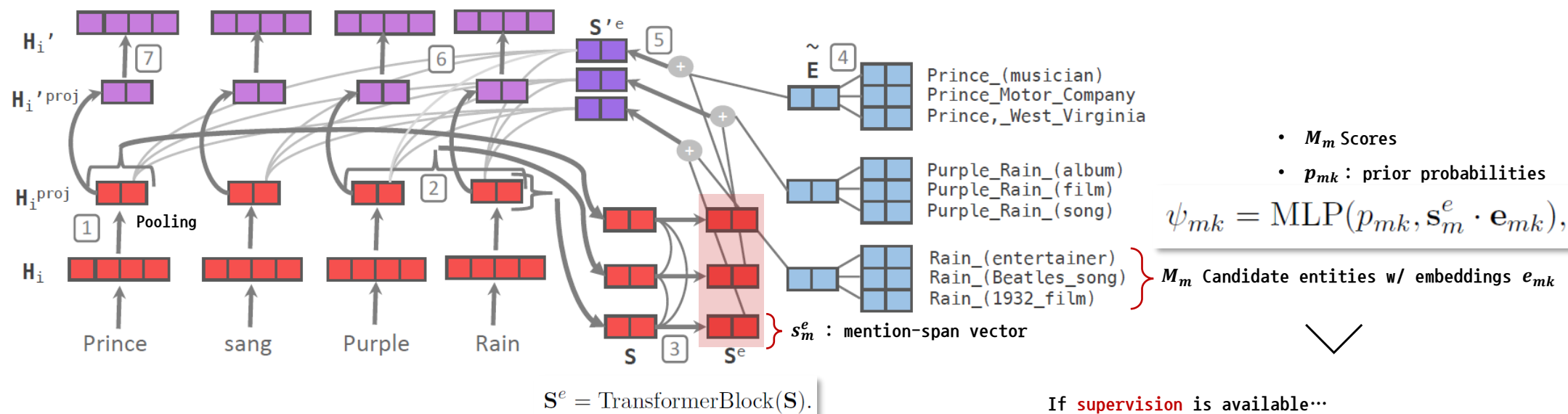
$C$  mention-span representations

- (candidate mention-span)  $s_m \in \mathbb{R}^E$
- $S \in \mathbb{R}^{C \times E}$

**KAR** Knowledge Attention and Recontextualization component

- Entity-linker

- **(Responsibility)** Performs entity disambiguation for each potential mention from candidates



## Mention-span self attention

KnowBert resolves which of several overlapping candidate mentions should be linked

If **supervision** is available...

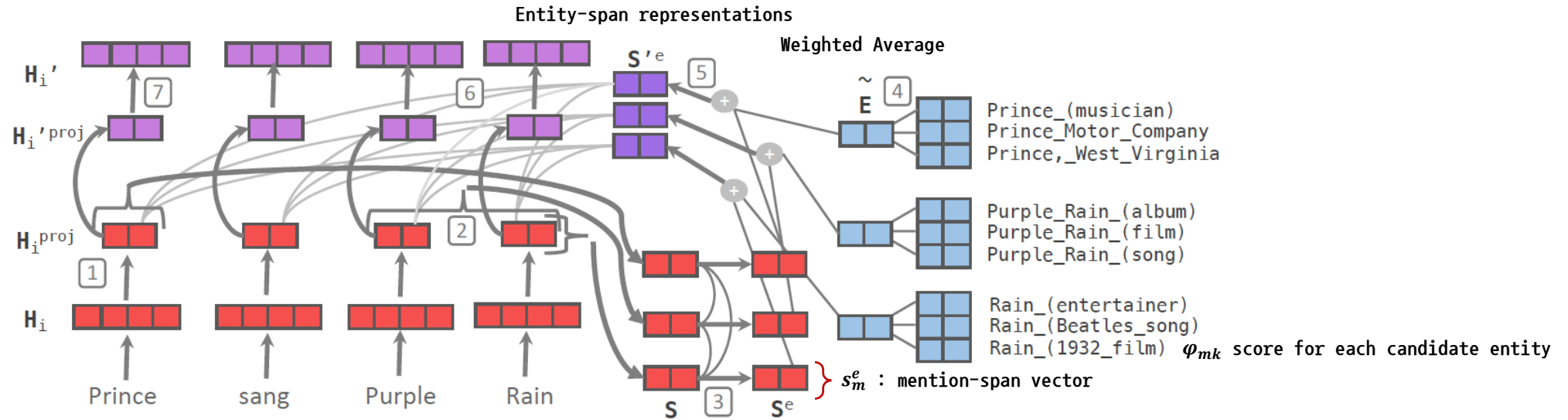
$$\mathcal{L}_{\text{EL}} = - \sum_m \log \left( \frac{\exp(\psi_{mg})}{\sum_k \exp(\psi_{mk})} \right)$$

- $e_{mg}$  gold entity
- Loss depending on the KB

KAR Knowledge Attention and Recontextualization component

- Knowledge enhanced entity-span representations

- Injects the KB entity info into the mention-span representations to form entity-span representations



$$\tilde{\psi}_{mk} = \begin{cases} \frac{\exp(\psi_{mk})}{\sum_{\psi_{mk} \geq \delta} \exp(\psi_{mk})}, & \psi_{mk} \geq \delta \\ 0, & \psi_{mk} < \delta. \end{cases}$$

>

$$\tilde{\mathbf{e}}_m = \sum_k \tilde{\psi}_{mk} \mathbf{e}_{mk}.$$

>

$$\mathbf{s}_m'^e = \mathbf{s}_m^e + \tilde{\mathbf{e}}_m,$$

- Disregard all candidates below a threshold

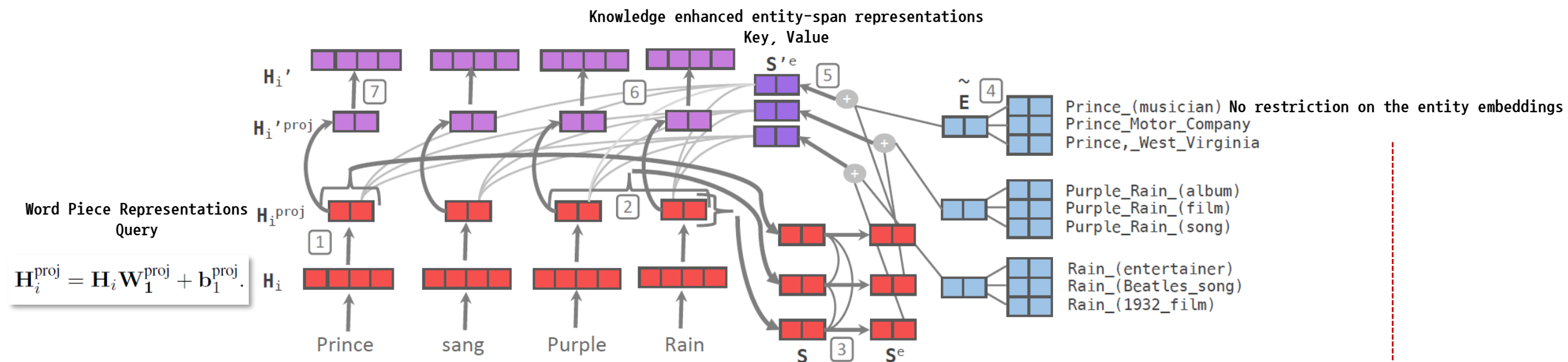
- Weighted entity embedding
- $\tilde{\mathbf{e}}_m <_{\text{below}} \delta \rightarrow \text{NULL embedding}$

- Entity-span representations updated
- $\mathbf{S}'^e \in \mathbb{R}^{C \times E}$

# KAR Knowledge Attention and Recontextualization component

## Recontextualization

- Uses entity-span representations to recontextualize the word piece representations



$$H_i^{\text{'proj}} = \text{MLP}(\text{MultiHeadAttn}(H_i^{\text{proj}}, S'^e, S'^e)).$$

- Word-to-entity span attention
- Propagates entity information over long contexts

$$H'_i = H_i^{\text{'proj}} W_2^{\text{proj}} + b_2^{\text{proj}} + H_i$$

- Projected back to the BERT dimension
- Alignment of BERT and entity vectors

- Initialize  $W_2^{\text{proj}}$  as the inverse of  $W_1^{\text{proj}}$



## Experiment

- **Setup**
  - **BERT** <sub>BASE</sub>
    - (+) biasing dataset to shorter sequences of 128 words pieces
  - **KnowBert-Wiki**
    - Insert the KB <sub>Knowledge Base</sub> between layers 10 and 11 of 12 layers BERT
  - **KnowBert-WordNet**
    - Synset metadata, Lemma metadata, and the relational graph
    - Insert the KB <sub>Knowledge Base</sub> between layers 10 and 11 of 12 layers BERT
  - **KnowBert-W** <sub>Wiki</sub> + **W** <sub>WordNet</sub>
    - Insert the Wikipedia KB <sub>Knowledge Base</sub> between layers 10 and 11, with WordNet between 11 and 12.

## Experiment

### • Downstream Tasks

- **(evaluate)** Validates that the knowledge addition improves performance on tasks

	System	LM	P	R	F <sub>1</sub>
	Zhang et al. (2018)	—	69.9	63.3	66.4
	Alt et al. (2019)	GPT	70.1	65.0	67.4
	Shi and Lin (2019)	BERT <sub>BASE</sub>	73.3	63.1	67.8
ERNIE	Zhang et al. (2019)	BERT <sub>BASE</sub>	70.0	66.1	68.0
	Soares et al. (2019)	BERT <sub>LARGE</sub>	—	—	70.1
Matching in the blanks	Soares et al. (2019)	BERT <sub>LARGE</sub> <sup>†</sup>	—	—	<b>71.5</b>
	KnowBert-W+W	BERT <sub>BASE</sub>	71.6	71.4	<b>71.5</b>

Table 4: Single model test set results on the TACRED relationship extraction dataset. <sup>†</sup> with MTB pretraining.

System	Accuracy
ELMo <sup>†</sup>	57.7
BERT <sub>BASE</sub> <sup>†</sup>	65.4
BERT <sub>LARGE</sub> <sup>†</sup>	65.5
BERT <sub>LARGE</sub> <sup>††</sup>	69.5
KnowBert-W+W	<b>70.9</b>

Table 6: Test set results for the WiC dataset (v1.0).

<sup>†</sup>Pilehvar and Camacho-Collados (2019)

<sup>††</sup>Wang et al. (2019a)

### • Relation extraction

#### • Settings

- (dataset) TACRED and SemEval 2010 Task 8
- (given) a sentence w/ marked a Sub and Obj → (pred) which realtation

#### • Model

- [E1][ /E1] <sub>Loc of Sub</sub>, [E2] [ /E2] <sub>Loc of Obj</sub>
- Concatenate the contextual representations of them to predict relation

#### • Results

- Better than ERNIE, BERT<sub>LARGE</sub>
- Match the MTB <sub>relation specific</sub> pretraining

### • Words in Context

#### • Definition

- 2 sentences containing a word with the same lemma → determine if they are the same sense or not → Test contextual word representations

#### • Settings

- [CLS] fine tuning Sentence 1 [SEP] Sentence 2

#### • Results

- Achieves SOTA in Accuracy

## Experiment

- Downstream Tasks

System	P	R	F <sub>1</sub>
UFET	68.8	53.3	60.1
BERT <sub>BASE</sub>	76.4	71.0	73.6
ERNIE	78.4	72.9	75.6
KnowBert-W+W	78.6	73.7	<b>76.1</b>

Table 7: Test set results for entity typing using the nine general types from (Choi et al., 2018).

- Entity-Typing

- Settings

- (dataset) Choi et al.(2018)

- Model

- Evaluation protocol w/ 9 general entity types
    - [E][/E] mark the loc of a target span

- Results

- Shows improvements F1 over ERNIE and BERT<sub>BASE</sub>

- Limitations of Existing Research

- Requires entity linking to represent entities in a text, and cannot **represent** entities **not in the KB**
- **(Alternative)** CWRs provides word representations trained w/ **unsupervised pretraining task**
- **CWR** Contextualized word representation
  - Many studies solved entity-related tasks w/ **contextualized representations** of entities **computed on CWR**
  - CWR <sup>Transformer based</sup> **not suited** to representing **entities**
    - (Why) CWRs do not output **the span-level representations** of entities
    - (Why) CWRs are **difficult** to perform **reasoning** between entities <sup>multiple tokens</sup>
      - (In contrary) Possible to capture the complex relationships between words

*“The Lord of the [MASK]”*

- (e.g) **Word-based** pretraining task of CWRs
  - Predicting a masked word given other words in the entity
  - **Not suitable** for learning the representations of entities <sup>Book, Movie</sup>

- **Proposed Model** <sup>LUKE</sup>

- New pretrained **contextualized representations of words and entities**
- **Difference between LUKE and CWR**
  - **Treats** not only words but also **entities** as **independent tokens**
  - Computes **intermediate** and **output representations** for **all tokens** w/ transformer
- **Task**
  - Randomly **masking entities** by [MASK] entities
  - Trains model by predicting **the originals of the masked entities.**
- **Contribution: Model**
  - Base pre-trained model : **RoBERTa**
  - (+)conduct pretraining of the model by optimizing **MLM** and **proposed task**
  - (downstream task) the resulting model **computes representations of arbitrary entities** in the text using [MASK] entities
- **Contribution: entity-aware self-attention**
  - **Extends** the transformer using entity-aware self-attention mechanism.
  - **Determine easily** the types of tokens.

- Input Representation of a token

- Token Embedding

- (Word token embedding)  $A \in \mathbb{R}^{V_w \times D}$
- (Entity token embedding)  $B \in \mathbb{R}^{V_e \times H}$  and  $U \in \mathbb{R}^{H \times D}$

- Position Embedding  $i$ -th pos in seq

- (Word)  $C_i \in \mathbb{R}^D$
  - (Entity)  $D_i \in \mathbb{R}^D$ 
    - Multiple words
- avg the embeddings of the corresponding positions

- Entity Type Embedding

- Single vector  $e \in \mathbb{R}^D$

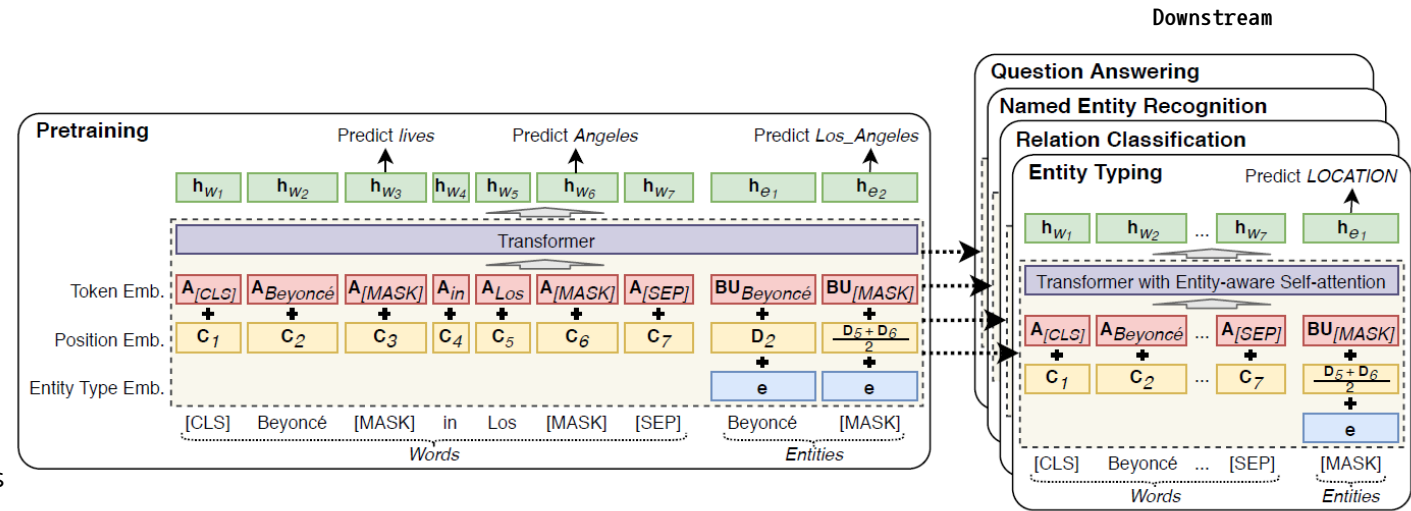


Figure 1: Architecture of LUKE using the input sentence “Beyoncé lives in Los Angeles.” LUKE outputs contextualized representation for each word and entity in the text. The model is trained to predict randomly masked words (e.g., *lives* and *Angeles* in the figure) and entities (e.g., *Los\_Angeles* in the figure). Downstream tasks are solved using its output representations with linear classifiers.

- Entity-aware Self-attention

- Relates token each other based on the attention score between each pair of tokens

A seq of input or output vector

- $x_1, \dots, x_k$  where  $x_i \in \mathbb{R}^D$
- $y_1, \dots, y_k$  where  $y_i \in \mathbb{R}^L$
- $k = m + n$  ( word or entity )

$$y_i = \sum_{j=1}^k \alpha_{ij} \mathbf{V} \mathbf{x}_j$$

$$e_{ij} = \frac{\mathbf{K} \mathbf{x}_j^\top \mathbf{Q} \mathbf{x}_i}{\sqrt{L}}$$

$$\alpha_{ij} = \text{softmax}(e_{ij})$$

How?

$$e_{ij} = \begin{cases} \mathbf{K} \mathbf{x}_j^\top \mathbf{Q} \mathbf{x}_i, & \text{if both } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ are words} \\ \mathbf{K} \mathbf{x}_j^\top \mathbf{Q}_{w2e} \mathbf{x}_i, & \text{if } \mathbf{x}_i \text{ is word and } \mathbf{x}_j \text{ is entity} \\ \mathbf{K} \mathbf{x}_j^\top \mathbf{Q}_{e2w} \mathbf{x}_i, & \text{if } \mathbf{x}_i \text{ is entity and } \mathbf{x}_j \text{ is word} \\ \mathbf{K} \mathbf{x}_j^\top \mathbf{Q}_{e2e} \mathbf{x}_i, & \text{if both } \mathbf{x}_i \text{ and } \mathbf{x}_j \text{ are entities} \end{cases}$$

- Assumption

- Beneficial to use the information of target token types when computing  $e_{ij}$

- Entity-aware query mechanism

- Uses a different query matrix for each pair of  $x_i$  and  $x_j$

- Different Query Matrices

- $Q_{w2e}, Q_{e2w}, Q_{e2e} \in \mathbb{R}^{L \times D}$

- Pretraining Task

- Conventional MLM

- An extension of the MLM to learn entity representations

- Settings

- (entity annotation) Hyperlinks in Wikipedia
      - (entity annotated corpus) Retrieved from Wikipedia

- Operations

- Randomly mask by replacing entities w/ [MASK]
    - The original entity to a masked one is predicted by the softmax function **over all entities in voca.**

$$\hat{y} = \text{softmax}(\mathbf{B}\mathbf{T}\mathbf{m} + \mathbf{b}_o)$$

$$\mathbf{m} = \text{layer\_norm}(\text{gelu}(\mathbf{W}_h \mathbf{h}_e + \mathbf{b}_h))$$

B : all entities

T : representation corresponding to the masked entity

- Final Loss

- MLM loss + Cross entropy loss on MLM extension



## • Experiments

Name	Prec.	Rec.	F1
BERT (Zhang et al., 2019)	67.2	64.8	66.0
C-GCN (Zhang et al., 2018b)	69.9	63.3	66.4
ERNIE (Zhang et al., 2019)	70.0	66.1	68.0
SpanBERT (Joshi et al., 2020)	70.8	70.9	70.8
MTB (Baldini Soares et al., 2019)	-	-	71.5
KnowBERT (Peters et al., 2019)	<b>71.6</b>	71.4	71.5
KEPLER (Wang et al., 2019b)	70.4	73.0	71.7
K-Adapter (Wang et al., 2020)	68.9	<b>75.4</b>	72.0
RoBERTa (Wang et al., 2020)	70.2	72.4	71.3
LUKE	70.4	75.1	<b>72.7</b>

Primary baseline

Table 2: Results of relation classification on the TA-CRED dataset.

Name	F1
LSTM-CRF (Lample et al., 2016)	91.0
ELMo (Peters et al., 2018)	92.2
BERT (Devlin et al., 2019)	92.8
Akbik et al. (2018)	93.1
Baevski et al. (2019)	93.5
RoBERTa	92.4
LUKE	<b>94.3</b>

Table 3: Results of named entity recognition on the CoNLL-2003 dataset.

## • Entity-Typing

### • Settings

- (dataset) Open Entity
- **9** general entity types
- micro-precision, recall, and F1

### • Model

- Treat the task as **multi-label classification**
- Train w/ **Binary-cross entropy** averaged over **all** entity types

### • Results

- Achieves **SOTA** in F1 points

## • Relation Classification

### • Settings

- (dataset) TACRED
- micro-precision, recall, and F1

### • Model

- [HEAD], [TAIL]
- Use a concatenated representations of the head and tail entities
- Traub w/ **cross-entropy**

### • Results

- Achieves **SOTA** in F1 points

## ERNIE vs. KAR vs. LUKE

- **ERNIE** 4 Jun 2019

- **(Purpose)** incorporate knowledge information into language representation models
- **(Proposed)** For **better fusion** of heterogeneous info from both text and KGs
  - knowledgeable aggregator
  - Pre-training task dEA
- (Experimental Result) **Better** abilities than BERT
- (Future research) Injecting knowledge into **feature-based** pre-training models, introducing diverse **structured knowledge** into **different** language representation model

- **KAR** 31 Oct 2019

- **(Purpose)** presents an efficient and general method to insert **prior knowledge** into a Neural Net.
- **(Proposed)** incorporates KB into BERT using **Knowledge Attention** and **Recontextualization** component
- (Experimental Result) **Improvements** for relationship extraction, entity typing and word sense disambiguation datasets
- (Future research) **incorporates a diverse set** of domain specific **KBs**

- **LUKE** 2 Oct 2020

- **(Limitation)** unsupervised pretraining task, **CWR**
- **(Purpose)** **newly pretrained contextualized** representations of words and entities
- **(Proposed)** **entity-aware self-attention mechanism**
- (Experimental Result) **proves** its effectiveness on tasks
- (Future research) **apply** LUKE to **domain-specific** <sup>biomedical, legal</sup> tasks

NLP Seminar

---

# KEPLER & GreaseLM

---

2022. 06. 03

KISTI - UST **IKJE CHOI**



# CONTENTS

01

KEPLER

02

GreaseLM

## KEPLER

- Title : KEPLER: A Unified Model for Knowledge Embedding and Pre-trained Language Representation
- Google Scholar

### KEPLER: A **unified** model for knowledge embedding and pre-trained language representation

[X Wang](#), [T Gao](#), [Z Zhu](#), [Z Zhang](#), [Z Liu](#), [J Li](#)... - Transactions of the ..., 2021 - [direct.mit.edu](#)

... In this paper, we propose **KEPLER**, a simple but effective **unified** model for knowledge embedding and pre-trained language representation. We train **KEPLER** with both the KE and ...

☆ Save  Cite Cited by 160 Related articles All 12 versions

# KEPLER

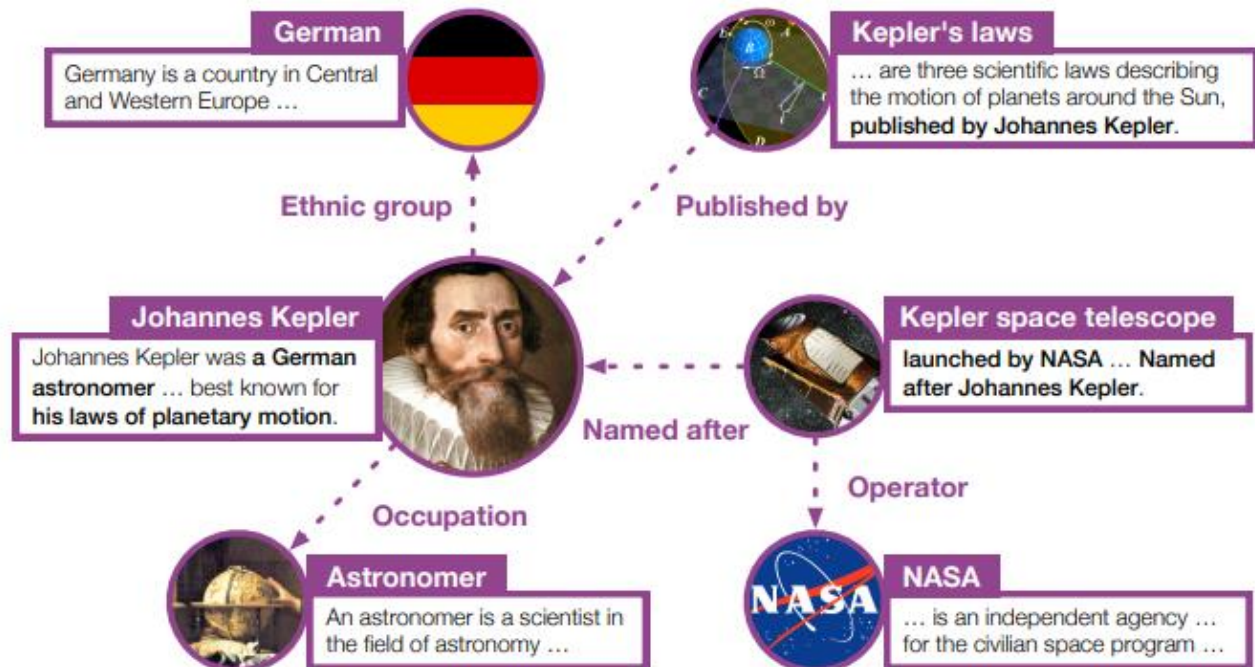
## ➤ KEPLER

- ✓ Pre-trained language representation model :  
not good to capture factual knowledge from text
- ✓ Knowledge graph :  
not good no abundant textual information.
- ✓ Propose a unified model KEPLER which integrates factual knowledge into pre-trained language representation model by using Knowledge graph

## KEPLER's Entity description

### ➤ KEPLER

- ✓ Entity description is used to fill the gap between KE and PLM.
- ✓ Using PLM encoder, we learn the model by encoding text and entity description into the unified semantic space and by jointly optimizing KE objective and MLM objective.



## KEPLER's strengths

### ➤ As a PLM

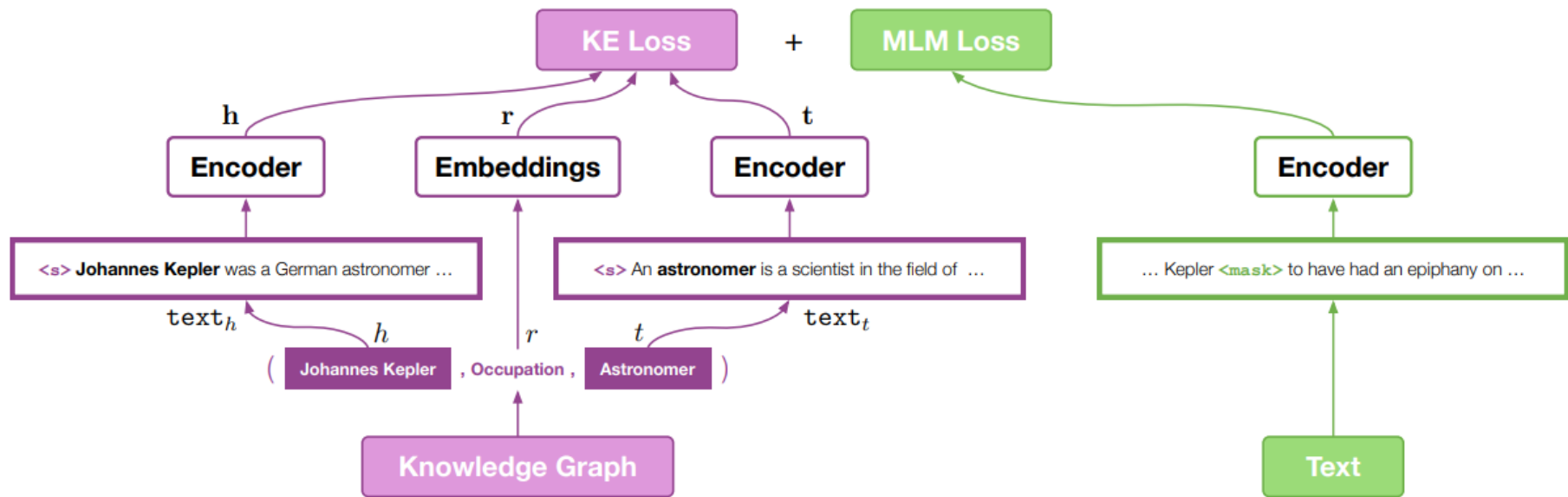
- ✓ KEPLER can integrate the actual knowledge into the language presentation using KE objective.
- ✓ It uses MLM objective to embed a strong language understanding capability of PLM.
- ✓ KE objective requires a model to encode an entity from an entity corresponding description, thus improving the possibility of extracting knowledge from text.
- ✓ KEPLER does not modify the model structure but only adds training objectives, so it can be applied directly to NLP tasks like other PLMs



## KEPLER's strengths

- As a KE model
  - ✓ KEPLER will be able to make good use of abundant information from the entity description using MLM objective.
  - ✓ KEPLER can also perform knowledge embedding in inductive setting.

## KEPLER framework



- KE Loss and MLM Loss are learned in different mini batches.

# KEPLER

## ➤ Knowledge Embedding in KEPLER

- ✓ Knowledge Embedding (KE) is an encoding of the entity and relation of a knowledge graph (KG) as distributed representation.
- ✓ In KEPLER, Using corresponding text as the entity and encode it to vector. There are three kinds of corresponding text that can be used.

(1) Entity description

(2) Entity and relation description

(3) Entity embeddings conditioned on relations

## KEPLER

### ➤ Knowledge Embedding in KEPLER

- ✓ Three kinds of corresponding text

#### (1) Entity description

- Encodes the description text related to the head entity and uses it as an embedding of the head entity.
- Tail entity is also used by obtaining embeddings in the same way.
- $\langle s \rangle$  means that special token  $\langle s \rangle$  at the beginning and end of the text.

$$\mathbf{h} = E_{\langle s \rangle}(\text{text}_h),$$

$$\mathbf{t} = E_{\langle s \rangle}(\text{text}_t),$$

$$\mathbf{r} = \mathbf{T}_r,$$

## KEPLER

### ➤ Knowledge Embedding in KEPLER

- ✓ Three kinds of corresponding text

#### (2) Entity and relation description

- The embeddings of the head entity and tail entity are obtained by using the same method as the previous entity description.
- The embeddings of the relation type are also obtained by encoding the relation corresponding text description.
- $\langle s \rangle$  means that special token  $\langle s \rangle$  at the beginning and end of the text.

$$\mathbf{h} = E_{\langle s \rangle}(\text{text}_h),$$

$$\mathbf{t} = E_{\langle s \rangle}(\text{text}_t),$$

$$\hat{\mathbf{r}} = E_{\langle s \rangle}(\text{text}_r),$$

## KEPLER

### ➤ Knowledge Embedding in KEPLER

- ✓ Three kinds of corresponding text

(3) Entity embeddings conditioned on relations

- It is a method derived from the intention that different relations will have different embeddings.
- Concatenate head entity and the relation entity, and uses it as an embedding of head entity.

$$\mathbf{h}_r = E_{\langle s \rangle}(\text{text}_{h,r}),$$

$$\mathbf{t} = E_{\langle s \rangle}(\text{text}_t),$$

$$\mathbf{r} = \mathbf{T}_r,$$

## KEPLER

➤ Loss function in KEPLER

$$\mathcal{L}_{\text{KE}} = -\log \sigma(\gamma - d_r(\mathbf{h}, \mathbf{t})) \\ - \sum_{i=1}^n \frac{1}{n} \log \sigma(d_r(\mathbf{h}'_i, \mathbf{t}'_i) - \gamma),$$

- Gamma means margin, and  $d_r(\mathbf{h}, \mathbf{t})$  means score function.
- The  $d_r(\mathbf{h}, \mathbf{t})$  used here is derived from the TransE model.
- The case of ' $\mathbf{h}_i$ ' and ' $\mathbf{t}_i$ ' means negative sample.
- Using negative sampling, related entities are learned close to each other and unrelated entities are learned far.

## KEPLER

➤ Loss function in KEPLER

$$d_r(\mathbf{h}, \mathbf{t}) = \|\mathbf{h} + \mathbf{r} - \mathbf{t}\|_p,$$

- The  $d_r(h,t)$  used here is derived from the TransE model.
- In the case of TransE, the entity is used as a vector, and the equation is designed that the sum of the head entity and the relation type is similar to the tail entity, resulting in the following score function. ( $h + r = t$ )



# KEPLER

## ➤ Conclusion

- ✓ In this paper, proposing a simple but effective unified model KEPLER for knowledge embedding and pre-trained language presentation.
- ✓ KEPLER uses KE and MLM objective to learn to align factual knowledge and language repression in the same semantic space, showing good performance in many NLP and KE applications.



## GreaseLM

- Title : GreaseLM: Graph REASoning Enhanced Language Models for Question Answering
- Google Scholar

### GreaseLM: Graph REASoning Enhanced Language Models for Question Answering

[X Zhang](#), [A Bosselut](#), [M Yasunaga](#), [H Ren](#)... - arXiv preprint arXiv ..., 2022 - arxiv.org

... In this work, we propose **GREASELM**, a new model that fuses encoded representations ...  
-USMLE) domains demonstrate that **GREASELM** can more reliably answer questions that ...

☆ Save    Cite   Cited by 1   Related articles   All 2 versions   

## GreaseLM

### ➤ Main Task : Solving QA task

- ✓ Existing Method : Knowledge Graph + Language Models fuse in a shallow.

It means that they encode separately and fuse them at the output for a prediction.

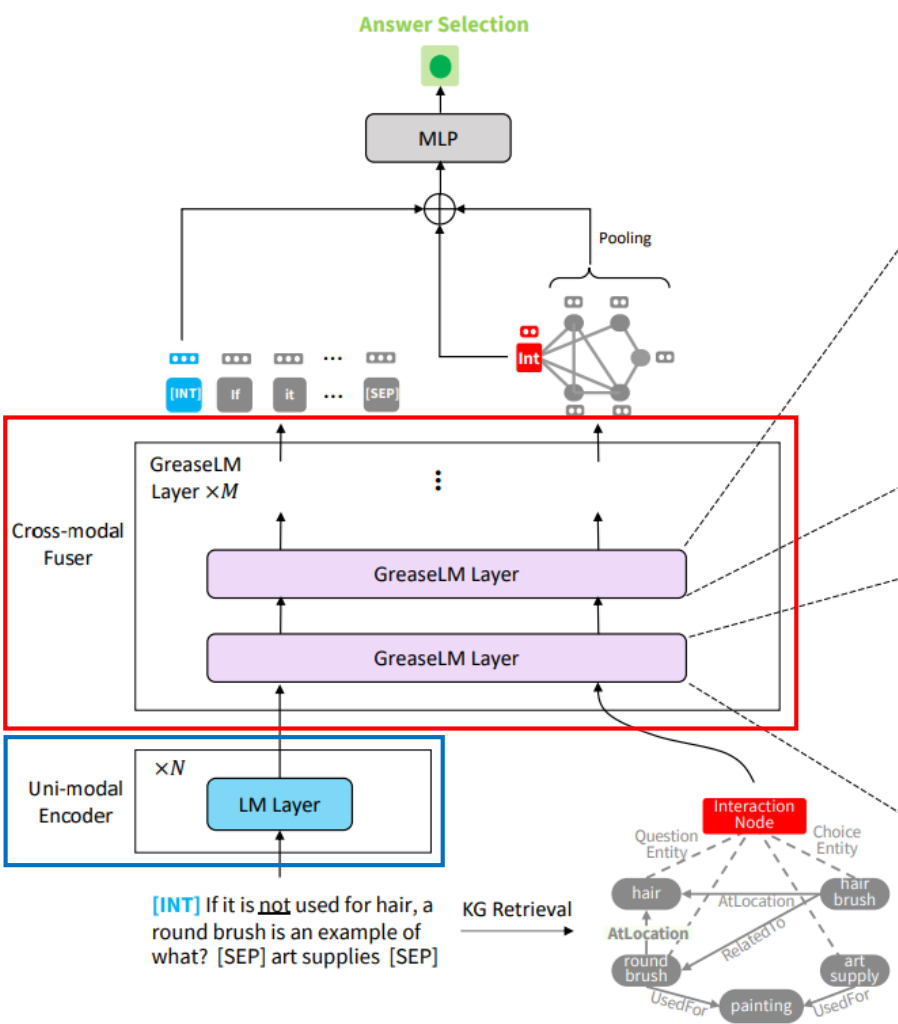
-> so, they do not exchange useful information each other.

- ✓ Proposed Model : GreaseLM

Fuses pre-trained language models (PLMs) and Knowledge graphs(KGs) then encodes and operates.

# GreaseLM Architecture

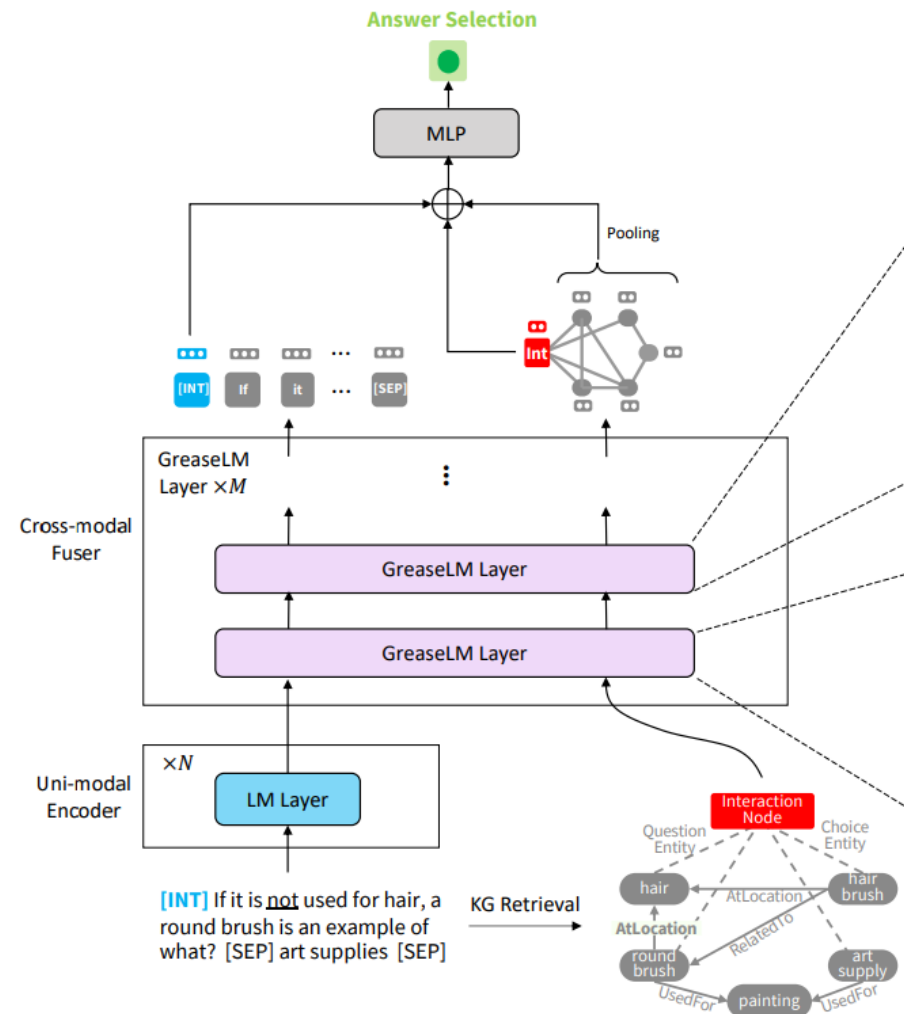
- GreaseLM consists of two stacked components
  - ✓ Unimodal LM layers (N) : initial representation of input tokens
  - ✓ Cross-modal GreaseLM layers (M) : LM layer + GNN layer



## GreaseLM Architecture

### ➤ Input Representation

- ✓ Concatenate cotext paragraph  $c$ , question  $q$ , candidate answer  $a$  with SEP tokens
- ✓ Use input sequence to retrieve subgraph of the KG
- ✓ Special interaction token  $W$  : prepend to token sequence
- ✓ Special interaction node  $E$  : connect to all linked nodes in  $V$  in  $G$



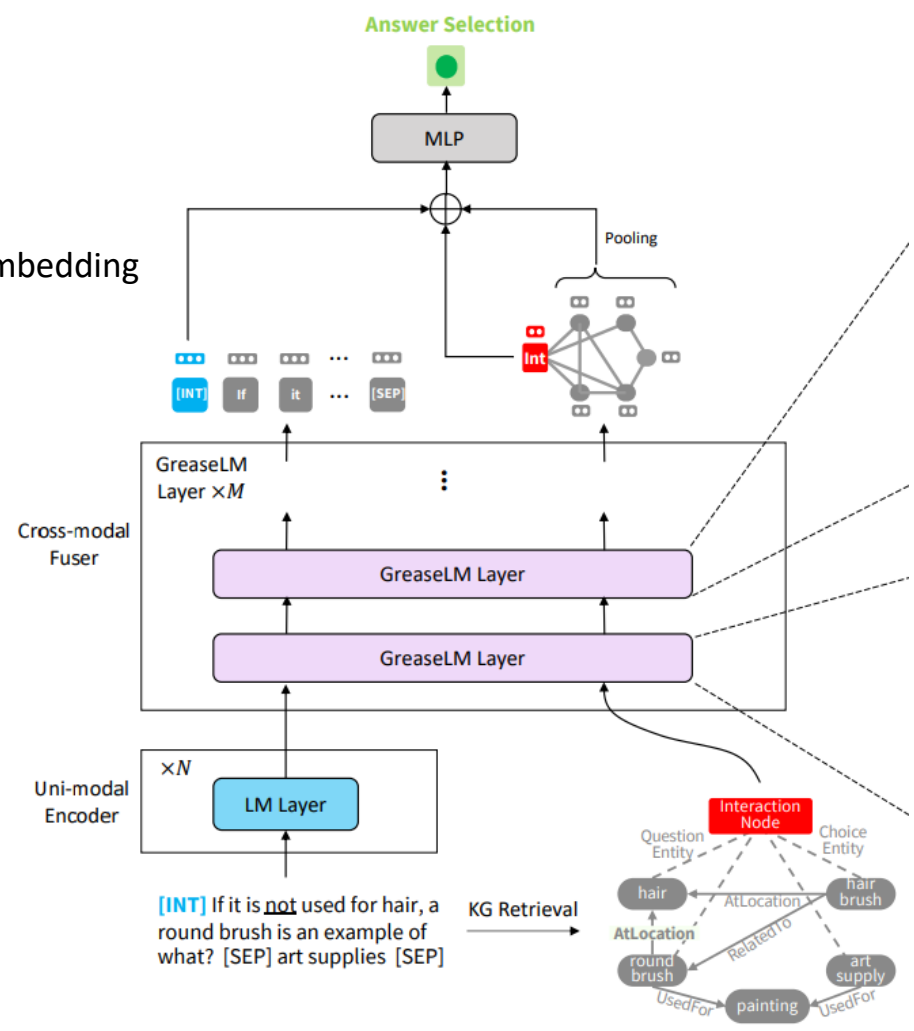
# GreaseLM Architecture

➤ Language Pre-encoding

$$\{h_{int}^{(0)}, h_1^{(0)}, \dots, h_T^{(0)}\} = \{w_{int}, w_1, \dots, w_T\} + \text{segment} + \text{positional embedding}$$

$$\{h_{int}^{(\ell)}, h_1^{(\ell)}, \dots, h_T^{(\ell)}\} = \text{LM-Layer}(\{h_{int}^{(\ell-1)}, h_1^{(\ell-1)}, \dots, h_T^{(\ell-1)}\})$$

for  $\ell = 1, \dots, N$

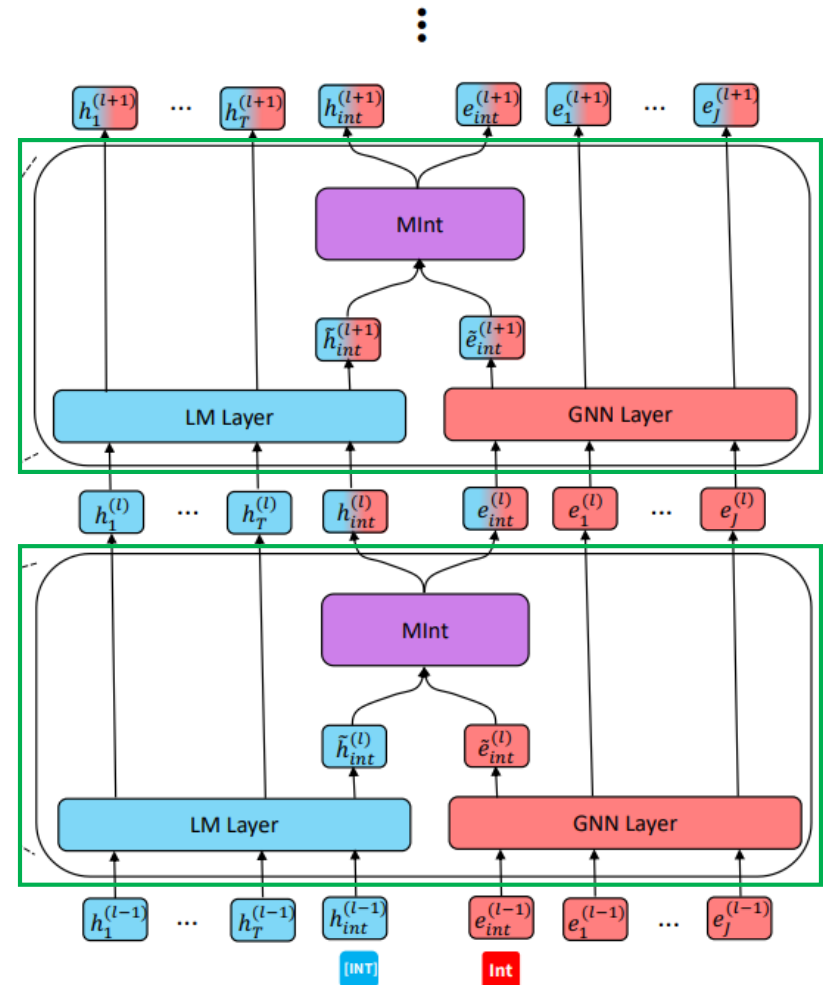


## ➤ What GreaseLM Layer does

- ✓ Separately encode information from both modalities (LM and GNN layer)
- ✓ Fuse both representations using the special interaction token and node

## ➤ Three Components of a GreaseLM layer

- ✓ A transformer LM encoder block : continues to encode language context
- ✓ A GNN layer : reason over KG
- ✓ A modality interaction layer : two modalities fuse info



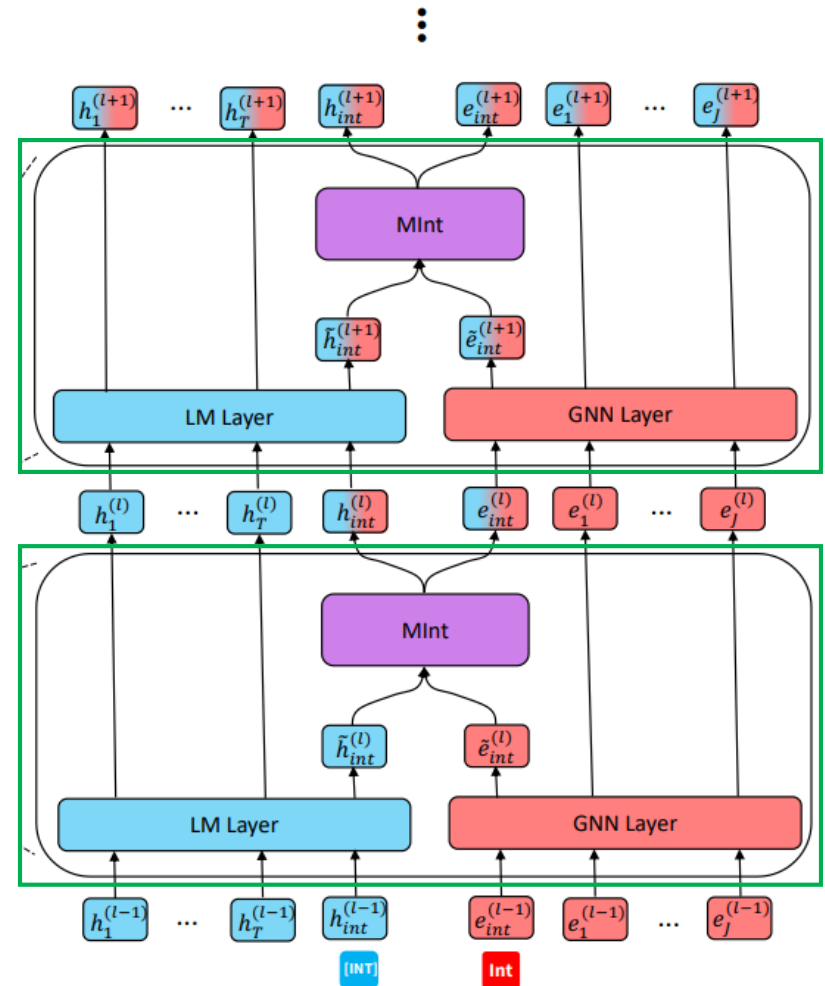


## ➤ Three Components of a GreaseLM layer

- ✓ A transformer LM encoder block :  
continues to encode language context

$$\{\tilde{h}_{int}^{(N+\ell)}, \tilde{h}_1^{(N+\ell)}, \dots, \tilde{h}_T^{(N+\ell)}\} = \text{LM-Layer}(\{h_{int}^{(N+\ell-1)}, h_1^{(N+\ell-1)}, \dots, h_T^{(N+\ell-1)}\})$$

for  $\ell = 1, \dots, M$



## ➤ Three Components of a GreaseLM layer

✓ A GNN layer : reason over KG

- (1) Compute initial node embedding
- (2) Initial embedding of the interaction node is initialized randomly
- (3) In each layer of GNN

$$\{\tilde{e}_{int}^{(\ell)}, \tilde{e}_1^{(\ell)}, \dots, \tilde{e}_J^{(\ell)}\} = \text{GNN}(\{e_{int}^{(\ell-1)}, e_1^{(\ell-1)}, \dots, e_J^{(\ell-1)}\})$$

for  $\ell = 1, \dots, M$

$$\tilde{e}_j^{(\ell)} = f_n\left(\sum_{e_s \in \mathcal{N}_{e_j} \cup \{e_j\}} \alpha_{sj} m_{sj}\right) + e_j^{(\ell-1)}$$

$$r_{sj} = f_r(\tilde{r}_{sj}, u_s, u_j)$$

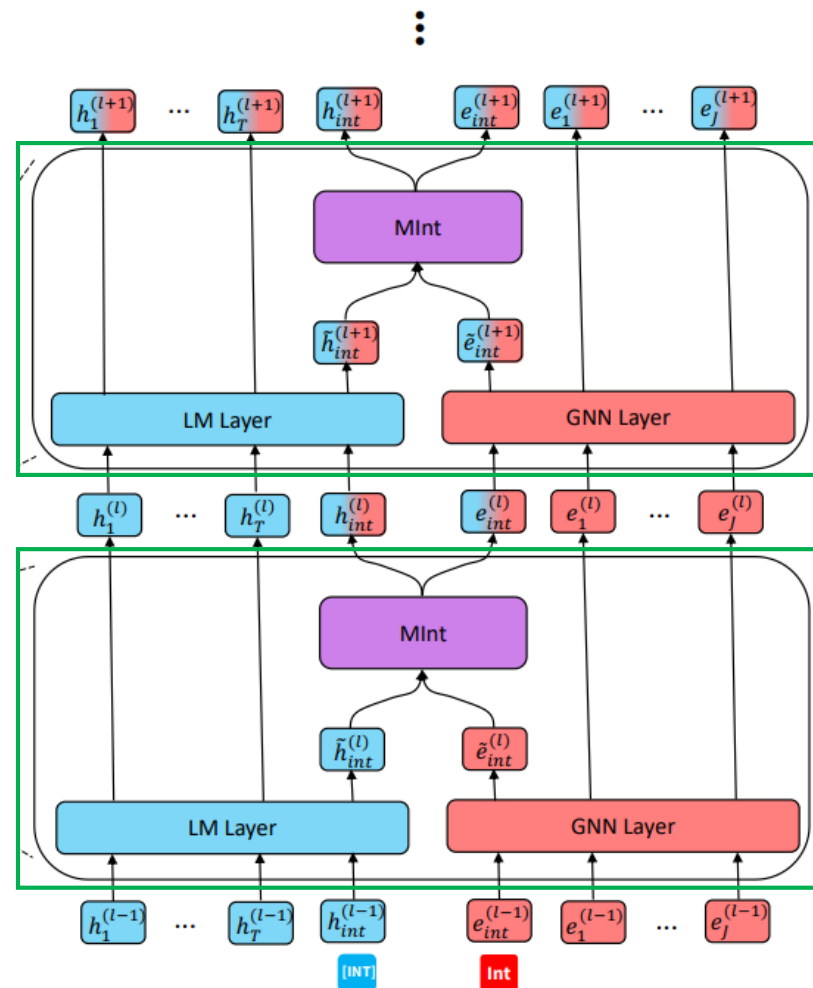
$$m_{sj} = f_m(e_s^{(\ell-1)}, u_s, r_{sj})$$

$$q_s = f_q(e_s^{(\ell-1)}, u_s)$$

$$k_j = f_k(e_j^{(\ell-1)}, u_j, r_{sj})$$

$$\gamma_{sj} = \frac{q_s^\top k_j}{\sqrt{D}}$$

$$\alpha_{sj} = \frac{\exp(\gamma_{sj})}{\sum_{e_s \in \mathcal{N}_{e_j} \cup \{e_j\}} \exp(\gamma_{sj})}$$



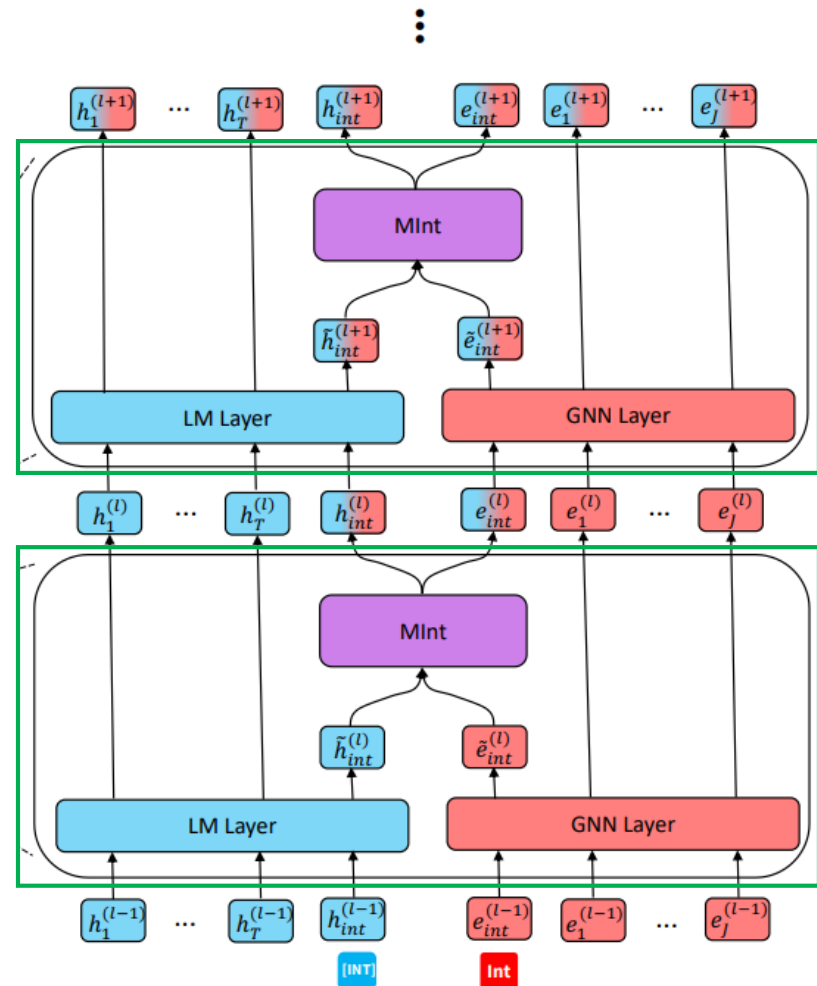
## ➤ Three Components of a GreaseLM layer

### ✓ A modality interaction layer :

two modalities fuse info

- (1) Concatenate pre-fused embeddings of the interaction token and interaction node
- (2) Pass joint representation through a mixing operation (MInt)
- (3) Split output post-fused embeddings into interaction token and interaction node

$$[h_{int}^{(\ell)}; e_{int}^{(\ell)}] = \text{MInt}([\tilde{h}_{int}^{(\ell)}; \tilde{e}_{int}^{(\ell)}])$$



## GreaseLM Architecture

### ➤ Learning & Inference

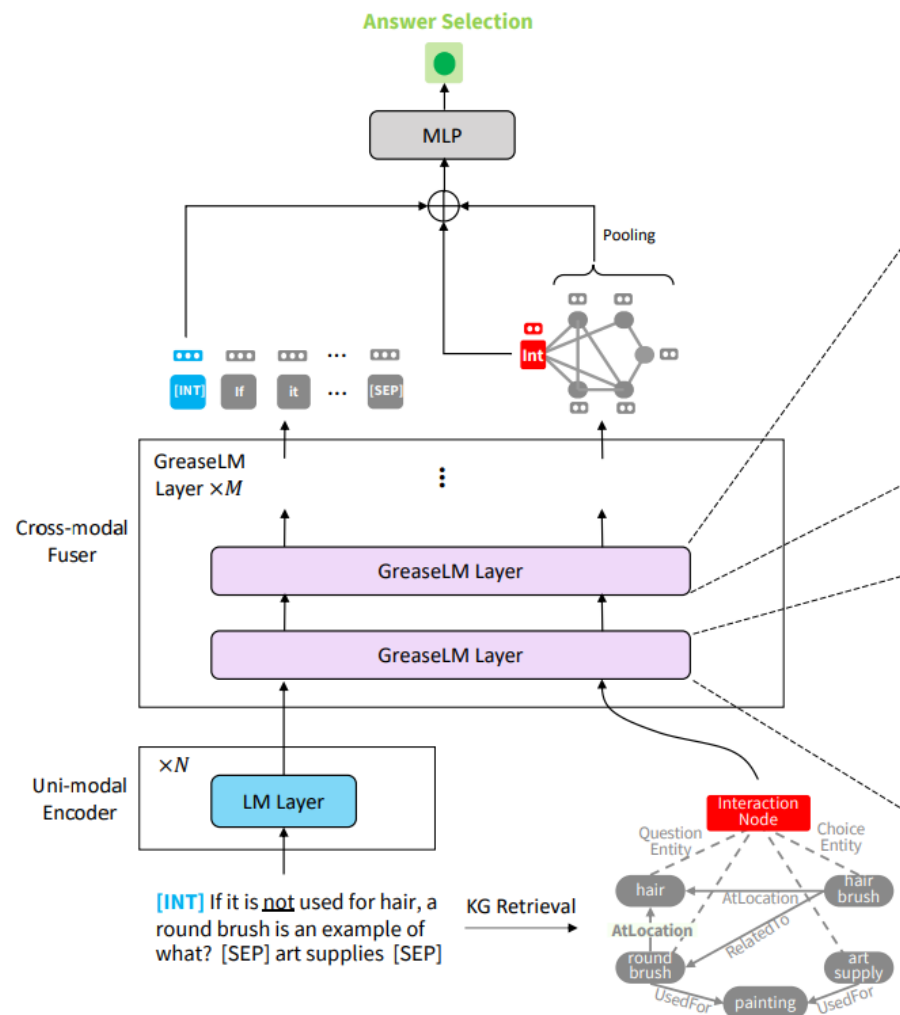
- ✓ Compute the probability of a being correct answer

$$\exp(\text{MLP}(h_{int}^{(N+M)}, e_{int}^{(M)}, g))$$

$g$  : attention-based pooling of

$$\{e_j^{(M)} \mid e_j \in \{e_1, \dots, e_J\}\} \text{ using } h_{int}^{(N+M)}$$

- ✓ Loss function : Cross Entropy Loss
- ✓ At inference :  
predict the most plausible answer as  
 $\arg \max_{a \in \mathcal{A}} p(a \mid q, c).$



## GreaseLM

### ➤ Setup

- ✓ CommonsenseQA, OpenBookQA, MedQA-USMLE
- ✓ Knowledge graph : ConceptNet and self-constructed KG for MedQA-USMLE
- ✓ LM : CommonsenseQA(RoBERTa-Large), OpenBookQA(AristoRoBERTa), MedQA-USMLE(SapBERT)



**THANK  
YOU**