# Enhancing Embedding via Two-level Features for Machine Reading Comprehension

Shuyi Wang, Hui Song, Bo Xu, Hongkuan Zhang

School of Computer Science and Technology, Donghua University, Shanghai, 201620, China
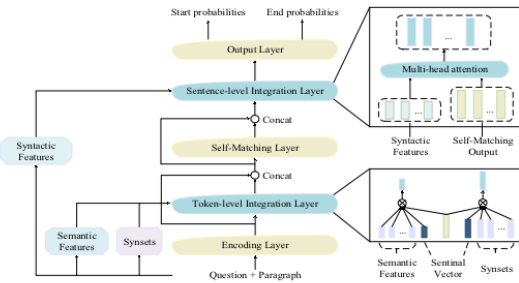
## Contribution

- Through the algorithm analysis and case study of MRC task, we confirm the issues of existing models and propose a two-level approach for integrating token features and grammatical structure of token-pairs into sentence encoding.
- For token-level, we select the valid features from the candidates for different tokens according to the actual context, sentinel vector is added. And for sentence level, We transform the dependency tree of sentence into m-hop matrixes, and then integrate the information through multi-head attention mechanism.
- We combine different features that need to be introduced to conduct experiments, find that the impact of features varies on different datasets. We choose DocQA, pre-trained model Google BERT_base, KT-NET_base as baseline methods. Our model achieves +0.63 $EM$ /+0.68 $F_1$ improvements in SQuAD1.1 dataset and +6.43 $EM$ /+6.79 $F_1$ improvements in ReCoRD dataset than the best result of baseline method, better than the current state-of-the-art models.

## Modeling

### Method

- **Task:** Span extraction of MRC task, which predicts the answer by locating the start and the end position in a context.
  - Given a question $Q = \{q_1, q_2, ..., q_n\}$ with $n$ tokens and a context $C = \{c_1, c_2, ..., c_m\}$ with $m$ tokens.
  - The extracted answer is a span in $C$.
- **Construction:** TLE-BERT model is based on BERT, and enhance the BERT embedding in two levels: token-level and sentence-level.
  - In token-level, for each token $t \in Q \cup C$, we add external knowledge $syn$ from the synonym dictionary and the semantic features $s$ of the token. The features are introduced by sentinel mechanism.
  - In sentence-level, the dependency relations $D$ is fused by multi-head attention mechanism with the token embedding.



### Encoding Layer based on BERT Model

- **Input:** $X = (x_1, x_2, ..., x_n)$, $n$ is the number of tokens.
  - Input contains the question sequence $Q$, the context sequence $C$ and three tags [CLS], [SEP], [SEP].
    - [CLS] is the classification label used in the classification task
    - [SEP] is used to distinguish the question and the context sequence
    - $X = \{[CLS], Q, [SEP], C, [SEP]\}$
- **Encoding:** for each token $x_t$, the initial text embedding $x_t^{token}$, position embedding $x_t^{pos}$ and segment embedding $x_t^{seg}$ are added together to obtain the initial input $x_t^0$.
  - $x_t^0 = x_t^{token} + x_t^{pos} + x_t^{seg}$
- **Output:** input embedding $x_t^0$ passes through the N-layer *encoder* of BERT to get context-aware representation $x_t^N$, and the output embedding sequence $x$ of this layer.
  - $x = \{x_1^N, x_2^N, ..., x_n^N\}, x_t^k = encoder(x_t^{k-1})$

### Token-level Integration Layer

- This layer integrates the semantic features $s$ and external knowledge $syn$ with the context-aware embedding representation $x_t^N$.
- Semantic features $s$ are identified from input sequence $X$.
  - include three embedding sequences, POS tag $s_{POS} = (s_t^{POS})$, entity type tag $s_{ET} = (s_t^{ET})$ and noun phrase tag $s_{NP} = \{s_t^{NP}\}$
  - enhance semantics on the basis of BERT embedding
- External knowledge $syn$ are synonyms retrieved from **WordNet**.
  - enrich the BERT embedding representation with external knowledge
- **Input:** semantic features$(s_t^{POS}, s_t^{ET}, s_t^{NP})$, some synonyms $\{syn\}_t$ and BERT embedding $x_t^N$. (for each token $x_t$ in $X$)
- **Integration:**
  - calculate two attention weights of $t$-token between $t$-semantic feature and $t$-Synonym, $\alpha_t^s$ and $\alpha_t^H$
    - $\alpha_t^s = s_t^T \cdot x_t^N, \alpha_t^H = syn_t^T \cdot x_t^N$
  - employ sentinel attention mechanism to organically integrate these features and knowledge
    - $\beta_t^s = softmax \cdot \alpha_t^s, \beta_t^H = softmax \cdot \alpha_t^H$
  - calculation of the final semantic feature $\hat{s}_t$ and knowledge $\hat{syn}_t$
    - $\hat{s}_t = \sum_h \beta_{t,h}^s + \beta_t^s s_t, \hat{syn}_t = \sum_h \beta_{t,h}^H syn_t^h + \beta_t^H$
- **Output:** concatenate the $x_t^N, \hat{s}_t$ and $\hat{syn}_t$ as the output of $t$-token in this layer
  - $u_t^k = [x_t^N, \hat{s}_t, \hat{syn}_t]$
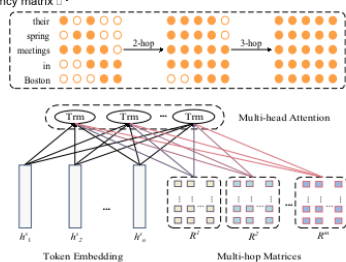
### Self-Matching Layer

- **Input:** token-level enriched representation $\{u_t\}_{t=1}^n$ as input.
- **Method:** use self-attention mechanism to interact context-aware embedding $\{x_t^N\}$, semantic features $\{s_t\}$ and external knowledge $\{syn_t\}$. This operation mainly discovers the direct and indirect interactions between tokens.
- **Output:** $v = \{v_t\}_{t=1}^n$

### Sentence-level Layer

- This layer introduces the dependencies of sentence-level into token embedding to get the structure message of the tokens.
- **Input:** Self-Matching layer output $v$ and multi-hop dependency matrices
- **Method:**
  - convert the pair-wise dependency of tokens into a 1-hop dependency matrix $D^1$
    - $D_{ij} = \begin{cases} 0, & c_i \nrightarrow c_j \rightarrow \text{ no relation} \\ 1, & c_i \cdot c_j \rightarrow \text{ exist relation} \end{cases}$
  - m-hop dependency matrix based on $D^1$
    - $D^m = (D^1)^m$
    - The example of multi-hop matrices is on the right.
  - get three input-based tensor $Query, Key, Value$ through $v$
    - $\{Query, Key, Value\} = \{v \cdot v_{Query}^z, v \cdot v_{Key}^z, v \cdot v_{Value}^z\}$
  - calculate the $z$-th head self-attention matrix $A^z$
    - $A^z = softmax(\frac{Query \cdot Key^T}{\sqrt{d}})$
  - calculate the $z$-th head self-attention based on $m$-hop matrix
    - $b^z = A^z \cdot Value$
  - concatenate all heads to get the embedding representation
    - $b_t = [b_t^1, ..., b_t^M]$
- **Output:** product the output $b = \{b_t\}$ through $v$
  - $o = b \cdot v$



### Output Layer

- **Input:** sentence layer output $o$
- **Output:** start probability $p_t^1$ and end probability $p_t^2$ for each token
  - $p_t^1 = \frac{e^{S_t^1}}{\sum_j e^{S_j^1}}, p_t^2 = \frac{e^{S_t^2}}{\sum_j e^{S_j^2}}$

### Objective Function

- Objective function is the log-likelihood
  - $J = -\frac{1}{N} \sum_{i=1}^N (\log p_{y_i^1}^1 + \log p_{y_i^2}^2)$

## Datasets

- **Source:**
  - **SQuAD1.1**: StandFord Question Answering Datasets
  - **ReCoRD**: Reading Comprehension with Commonsense Reasoning Datasets
- **Content:**
  - SQuAD1.1
    - 10,000+ questions based on 500 articles in Wikipedia
    - questions are in the form of interrogative sentences
    - the 11 types of samples are summarized in following table

| Question type | Number | Question | Number |
|---|---|---|---|
| What | 56905 | Where | 4090 |
| Who | 9900 | Be/Do/etc. | 1671 |
| Which | 6620 | Why | 1353 |
| When | 6258 | Whom | 394 |
| How many | 5735 | Whose | 350 |
| How | 4893 | | |

  - ReCoRD
    - 110,730 samples with a large portion of queries requiring commonsense reasoning
    - questions are a sentence with @placeholder instead of the missing text span that needs to be predicted

## Data Preprocessing

- **Method:**
  - Features of the sentences we used are mainly obtained through the spaCy tool.
  - Modify the tokenized method of spaCy to make it consistent with BERT.
- **Analysis:**
  - The tokenization mechanism of the pre-trained language model BERT is different from the spaCy's.

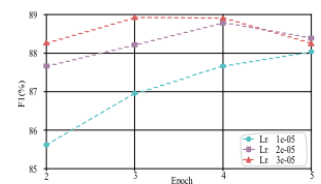| | |
|---|---|
| **Original:** | For whos glory did Father O'Hara believed that the Notre Dame football team played ? |
| **BERT:** | For who s glory did Father O'Hara believed that the Notre Dame football team played ? |
| **spaCy:** | For whos glory did Father O '. Hara believed that the Notre Dame football team played ? |

## Model Evaluation

- **Standard:** MRC task usually needs $F_1$ value and Extract Match ($EM$) for comprehensive evaluation.
  - $F_1$: calculate the degree of overlap between the predicted and the correct answer
    - $precision$: the ration of the number of overlapped characters to the number of correct answer characters
    - $recall$: he ration of the number of overlapped characters to the predicted answer characters
    - $F_1 = 2 * \frac{precision \cdot recall}{precision + recall}$
  - $EM$: precise matching
    - the value of $EM$ is 1 means that the predicted answer is exactly the same as the correct answer, otherwise is 0
    - in the SQuAD1.1 dataset, each question may have multiple candidate answers
    - the predicted answer only needs to exactly match one of the candidate answers

## Result Analysis

### Environment And Setup

- **Environment:**
  - python3.5.2
  - PaddlePaddle-gpu1.8.2
  - NLTK3.5
  - spaCy2.2.4
  - CUDA10.0.130
  - experiments are performed on 2 2080ti GPUs
- **Setup:**
  - parameters of the BERT_base in PaddlePaddle version
  - the trained model has 12 layers with 768 hidden layer dimension
  - maximum sequence length is set to 384
  - token-level feature's dimension is 20
  - learning rate is 3e-05
  - epoch 3
  - tokenized mechanism is WordPiece



### Results

- **General result:** compared TLE-BERT with three models DocQA, Google BERT_base and KT-NET (based on BERT_base) in our environment.
  - On SQuAD1.1, TLE-BERT gets a gain of +0.63 $EM$ /+0.68 $F_1$ over the BERT_base.
  - On ReCoRD, TLE-BERT achieves +6.43 $EM$ /+6.79 $F_1$ improvements.

| Model | SQuAD1.1 | | ReCoRD | |
|---|---|---|---|---|
| | EM | F1 | EM | F1 |
| DocQA(ELMo) | - | - | 44.13 | 45.39 |
| Google BERT_base(ours) | 81.25 | 88.41 | 55.99 | 57.99 |
| KT-NET_base(ours) | 81.56 | 88.75 | 60.79 | 62.91 |
| TLE-BERT(ours) | **81.88** | **89.09** | **62.42** | **64.78** |

- **Different feature result:** Introduce different features to the model to compare results.
  - NAN in following tabel is the KT-NET_base model containing synset knowledge, without any token or sentence level features.

| Model | SQuAD1.1 | | ReCoRD | |
|---|---|---|---|---|
| | EM | F1 | EM | F1 |
| NAN | 81.56 | 88.75 | 60.79 | 62.91 |
| +POS | 81.91 | **89.06** | 62.06 | 64.55 |
| +ET | 81.64 | 88.91 | 62.33 | **64.67** |
| +NP | 81.69 | 88.90 | 61.70 | 64.21 |
| +DEP | 81.83 | **89.05** | 62.23 | 64.69 |
| ALL | **81.88** | 89.09 | **62.42** | 64.78 |

- **Analysis:**
  - Token-level
    - POS achieve better performance in SQuAD1.1 than in ReCoRD, while ET is the opposite.
    - POS provides grammatical information that the SQuAD1.1 dataset lacks, and ReCoRD needs entity information for reasoning.
  - Sentence-level
    - Dependency relation improves the performance of the model on both datasets, which also confirms our previous conjecture that the pre-trained language model does not obtain the sentence dependency as well.

## Conclusion

- This paper proposes a new method of enhancing language representation through token-level and sentence-level features.
  - We can select appropriate features as external information according to the characteristics of the dataset or the application background to improve the model performance.
- We employ an attention mechanism to fuse multiple token-level features and make the model learn to pick a valid feature set by sentinel vectors.
- The sentence-level features enhance the language representation with syntactic information of the sentence.
- Experiments have shown that our method is better than the current SOTA BERT_base models.