

# 论文阅读笔记

## Step5

MF1833063, 史鹏, spwannasing@gmail.com

2019 年 7 月 8 日

# 1 ERNIE: Enhanced Language Representation with Informative Entities

此文章是对bert的一次扩展，提出了知识图谱中的多信息实体（informative entity）可以作为外部知识改善语言表征。

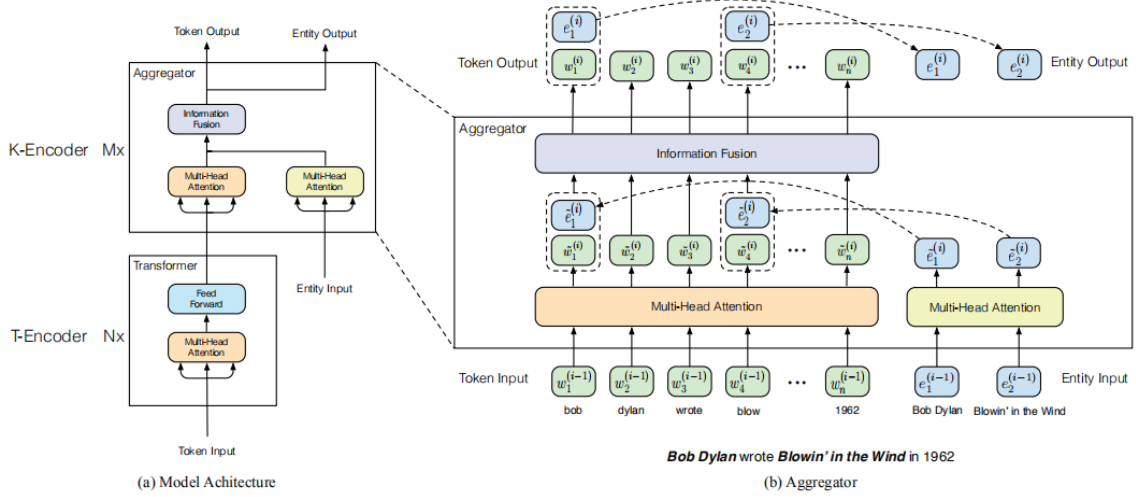


图 1: overview

**Mark Twain** wrote **The Million Pound Bank Note** in 1893.

Input for Common NLP tasks:

[CLS] [ ] mark twain [ ] wrote [ ] the million pound bank note [ ] in 1893 . [SEP]

Input for Entity Typing:

[CLS] [ENT] mark twain [ENT] wrote [ ] the million pound bank note [ ] in 1893 . [SEP]

Input for Relation Classification:

[CLS] [HD] mark twain [HD] wrote [TL] the million pound bank note [TL] in 1893 . [SEP]

图 2: finetune

Knowledgeable Encoder:

$$\{\tilde{w}_1^{(i)}, \dots, \tilde{w}_n^{(i)}\} = \text{MH-ATT}(\{w_1^{(i-1)}, \dots, w_n^{(i-1)}\}) \quad (1.1)$$

$$\{\tilde{e}_1^{(i)}, \dots, \tilde{e}_m^{(i)}\} = \text{MH-ATT}(\{e_1^{(i-1)}, \dots, e_m^{(i-1)}\}) \quad (1.2)$$

对于和entity对齐的token:

$$\begin{aligned} h_j &= \sigma(\tilde{W}_t^{(i)} \tilde{w}_j^{(i)} + \tilde{W}_e^{(i)} \tilde{e}_k^{(i)} + \tilde{b}^{(i)}) \\ w_j^{(i)} &= \sigma(W_t^{(i)} h_j + b_t^{(i)}) \\ e_k^{(i)} &= \sigma(W_e^{(i)} h_j + b_e^{(i)}) \end{aligned} \quad (1.3)$$

else:

$$\begin{aligned} \mathbf{h}_j &= \sigma \left( \tilde{\mathbf{W}}_t^{(i)} \tilde{\mathbf{w}}_j^{(i)} + \tilde{\mathbf{b}}^{(i)} \right) \\ \mathbf{w}_j^{(i)} &= \sigma \left( \mathbf{W}_t^{(i)} \mathbf{h}_j + \mathbf{b}_t^{(i)} \right) \end{aligned} \quad (1.4)$$

对于引入的信息的pre-training目标:

$$p(e_j | w_i) = \frac{\exp(\text{linear}(\mathbf{w}_i^o) \cdot \mathbf{e}_j)}{\sum_{k=1}^m \exp(\text{linear}(\mathbf{w}_i^o) \cdot \mathbf{e}_k)} \quad (1.5)$$

## 2 ERNIE: Enhanced Representation through Knowledge Integration

ERNIE 通过建模海量数据中的词、实体及实体关系，学习真实世界的语义知识。相较于 BERT 学习原始语言信号，ERNIE 直接对先验语义知识单元进行建模，增强了模型语义表示能力。这里我们举个例子：

Learnt by BERT：哈 [mask] 滨是 [mask] 龙江的省会，[mask] 际冰 [mask] 文化名城。

Learnt by ERNIE：[mask] [mask] [mask] 是黑龙江的省会，国际 [mask] [mask] 文化名城。

在 BERT 模型中，我们通过『哈』与『滨』的局部共现，即可判断出『尔』字，模型没有学习与『哈尔滨』相关的任何知识。而 ERNIE 通过学习词与实体的表达，使模型能够建模出『哈尔滨』与『黑龙江』的关系，学到『哈尔滨』是『黑龙江』的省会以及『哈尔滨』是个冰雪城市。

训练数据方面，除百科类、资讯类中文语料外，ERNIE 还引入了论坛对话类数据，利用 DLM（Dialogue Language Model）建模 Query-Response 对话结构，将对话 Pair 对作为输入，引入 Dialogue Embedding 标识对话的角色，利用 Dialogue Response Loss 学习对话的隐式关系，进一步提升模型的语义表示能力。

我们在自然语言推断，语义相似度，命名实体识别，情感分析，问答匹配 5 个公开的中文数据集上进行了效果验证，ERNIE 模型相较 BERT 取得了更好的效果。

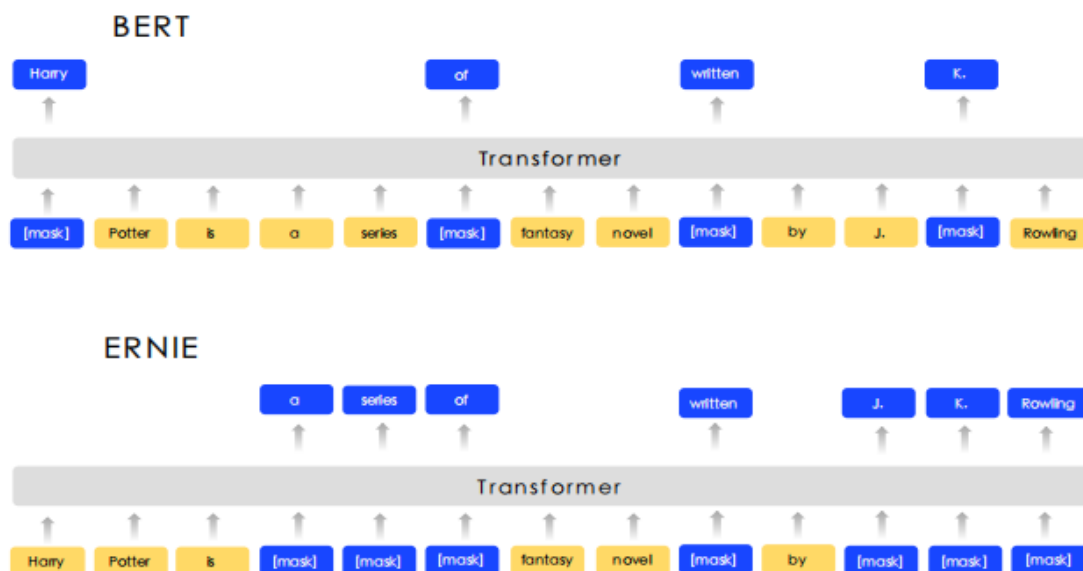


Figure 1: The different masking strategy between BERT and ERNIE

Sentence	Harry	Potter	is	a	series	of	fantasy	novels	written	by	British	author	J.	K.	Rowling
Basic-level Masking	[mask]	Potter	is	a	series	[mask]	fantasy	novels	[mask]	by	British	author	J.	[mask]	Rowling
Entity-level Masking	Harry	Potter	is	a	series	[mask]	fantasy	novels	[mask]	by	British	author	[mask]	[mask]	[mask]
Phrase-level Masking	Harry	Potter	is	[mask]	[mask]	[mask]	fantasy	novels	[mask]	by	British	author	[mask]	[mask]	[mask]

Figure 2: Different masking level of a sentence



Figure 3: Dialogue Language Model. Source sentence: [cls] How [mask] are you [sep] 8 . [sep] Where is your [mask] ? [sep]. Target sentence (words the predict): old, 8, hometown

### 3 Learning to Ask Unanswerable Questions for Machine Reading Comprehension

提出一种数据增强技术，根据与包含答案的相应段落配对的可回答问题自动生成相关的无法回答的问题。所提出的结构为“pair-to-sequence”。

<b>Title:</b> Victoria (Australia)
<b>Paragraph:</b> ...Public schools, also known as state or government schools, are funded and run directly by the <b>Victoria Department of Education</b> . Students do not pay tuition fees, but some extra costs are levied. Private fee-paying schools include parish schools ...
<b>Ans. Question:</b> What organization runs <i>the public schools</i> in Victoria?
<b>UnAns. Question:</b> What organization runs <i>the waste management</i> in Victoria?
<b>(Plausible) Answer:</b> <b>Victoria Department of Education</b>

Figure 1: An example taken from the SQuAD 2.0 dataset. The annotated (plausible) answer span in the paragraph is used as a pivot to align the pair of answerable and unanswerable questions.

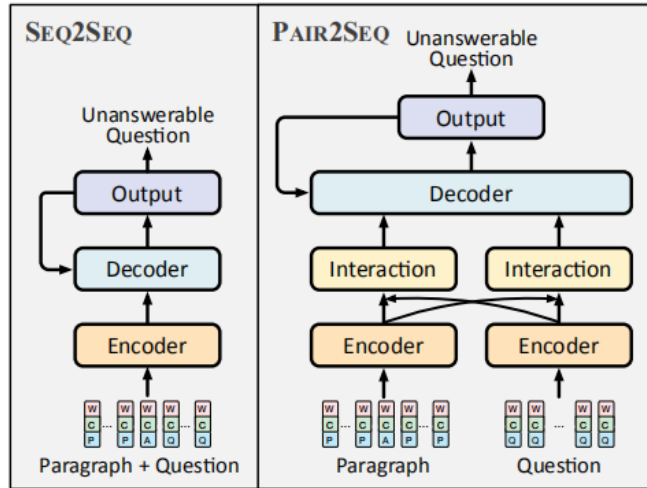


Figure 2: Diagram of the proposed pair-to-sequence model and sequence-to-sequence model. The input embeddings is the sum of the word embeddings, the character embeddings and the token type embeddings. The input questions are all answerable.

## 4 Course Concept Expansion in MOOCs with External Knowledge and Interactive Game

随着大规模在线开放课程(MOOC)的日益普及,为MOOC用户自动提供课外知识成为可能。语义漂移和知识缺乏在复杂的MOOC环境下,现有的方法不能有效地扩展课程概念。本文首先在通过外部知识库搜索新概念的过程中建立一个新的边界,然后利用异构特征来验证高质量的结果。具体算法没有特别仔细地看。

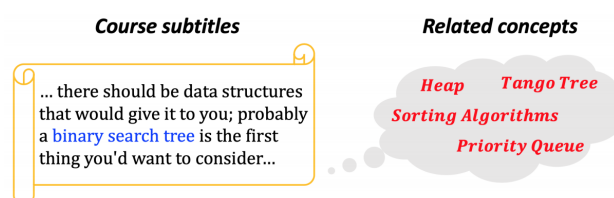


Figure 1: An example of “out-of-teaching” concepts in the course “*Data Structure and Algorithm*”.

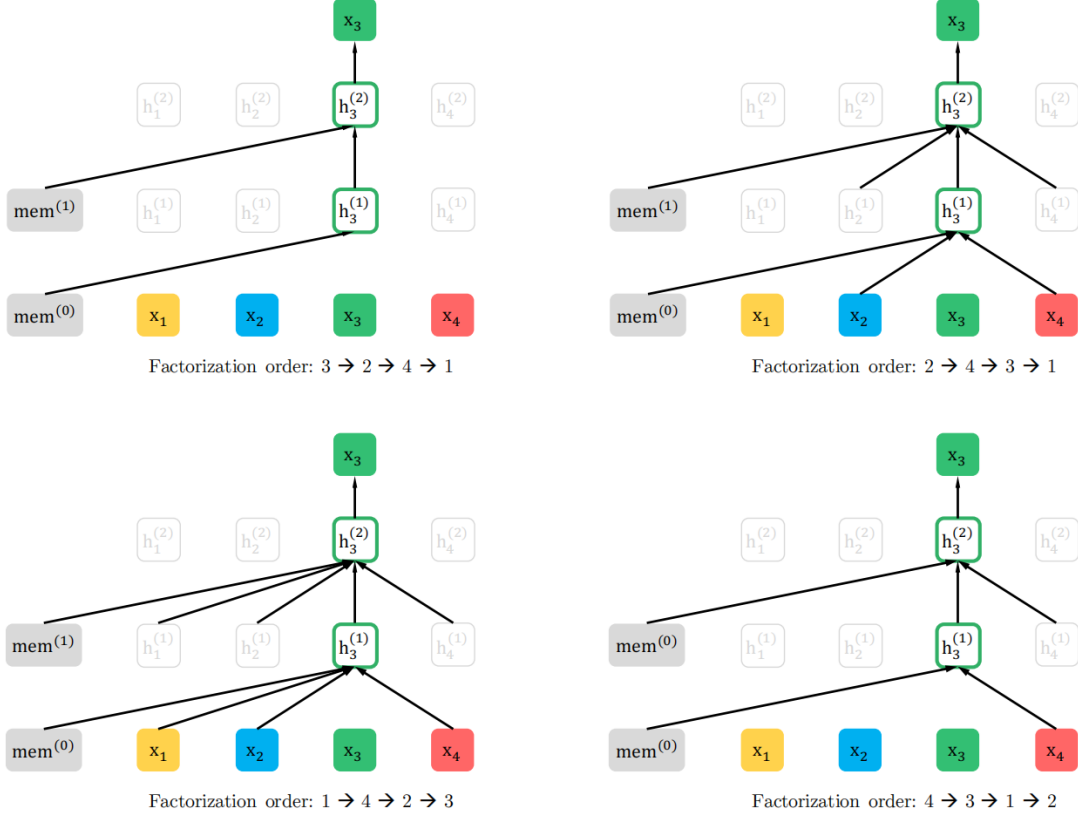
## 5 XLNet: Generalized Autoregressive Pretraining for Language Understanding

本文提出了一个新的预训练模型，在20个任务上刷新了Bert的记录。其相对bert的主要改进是结合了AR(Transformer-XL)和AE的优点，抛弃了bert的Masked language model，使用了Permutation language model。对于每一种排列的序列，并不是改变原始的序列（即position embedding不变），只是改变mask（bert里的mask是直接mask一个word，这里的是在Transformer的计算过程中“mask”住不需要的部分）。e.g 在图二中，对于序列3241，1能“看见”324，所以都不用mask掉，而对于2，只能看见3，所以要mask掉14。

虽然排列语言模型能满足目前的目标，但是对于普通的transformer结构来说是存在一定的问题的，为什么这么说呢，看个例子，假设我们要求这样的一个对数似然， $p_\theta(X_{z_t}|x_{z_{<t}})$ 如果采用标准的softmax的话，那么

$$p_\theta(X_{z_t}|x_{z_{<t}}) = \frac{\exp(e(x)^T h_\theta(x_{z_t}))}{\sum_{x'} \exp(e(x')^T h_\theta(x_{z_{<t}}))}$$

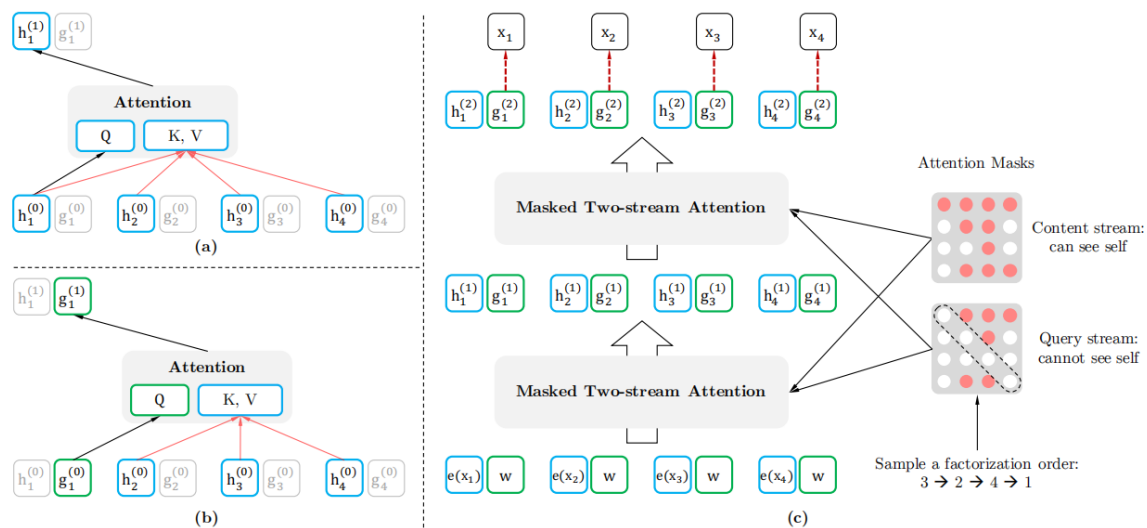
其中 $h_\theta(x_{z_{<t}})$ 表示的是添加了mask后的transformer的输出值，可以发现 $h_\theta(x_{z_{<t}})$ 并不依赖于其要预测的内容的位置信息，因为无论预测目标的位置在哪里，因式分解后得到的所有情况都是一样的，并且transformer的权重对于不同的情况是一样的，因此无论目标位置怎么变都能得到相同的分布结果，如下图所示，假如我们的序列index表示为[1,2,3]，对于目标2与3来说，其因式分解后的结果是一样的，那么经过transformer之后得到的结果肯定也是一样的。





因此就提出了基于目标感知表征的双流自注意力。

$$\begin{aligned} g_{z_t}^{(m)} &\leftarrow \text{Attention} \left( Q = g_{z_t}^{(m-1)}, KV = \mathbf{h}_{z_{<t}}^{(m-1)}; \theta \right) \\ h_{z_t}^{(m)} &\leftarrow \text{Attention} \left( Q = h_{z_t}^{(m-1)}, KV = \mathbf{h}_{z_{\leq t}}^{(m-1)}; \theta \right) \end{aligned}$$



reference:

<https://zhuanlan.zhihu.com/p/70257427>

<https://blog.csdn.net/u012526436/article/details/93196139>

## 6 Star-Transformer

本文提出了一种对于Transformer改进的结构“star-Transformer”。将标准Transformer中全连接的部分调整为星型的拓扑结构，在实际任务中速度提升了4.5倍。

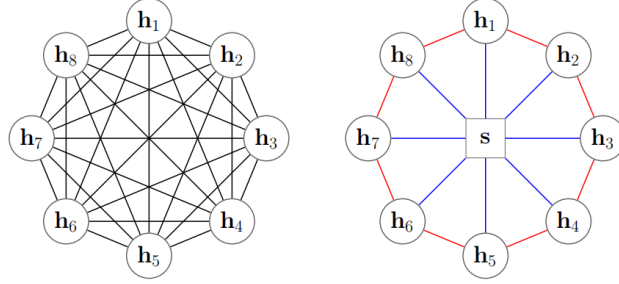


Figure 1: Left: Connections of one layer in Transformer, circle nodes indicate the hidden states of input tokens. Right: Connections of one layer in Star-Transformer, the square node is the virtual relay node. Red edges and blue edges are ring and radical connections, respectively.

添加了一个虚拟的节点作为relay node。由此可以分为两种连接，Radical Connections和Ring Connections。环连接可以有效地减少局部和非局部成分的非偏置学习负担，提高模型的泛化能力。将标准的Transformer复杂度  $O(n^2)$  降低为  $O(n)$ 。

根式连接集中在非局部构图上，环连接集中在局部构图上。因此，star-Transformer适用于尺寸适中的数据集，不依赖于繁重的预培训。

令  $\mathbf{s}^t \in \mathbb{R}^{1 \times d}$  and  $\mathbf{H}^t \in \mathbb{R}^{n \times d}$  分别为relay node和所有的satellite node。

初始化  $\mathbf{H}^0 = \mathbf{E}$  和  $s^0 = \text{average}(\mathbf{E})$

$$\begin{aligned}
 \mathbf{C}_i^t &= [\mathbf{h}_{i-1}^{t-1}; \mathbf{h}_i^{t-1}; \mathbf{h}_{i+1}^{t-1}; \mathbf{e}^i; \mathbf{s}^{t-1}] \\
 \mathbf{h}_i^t &= \text{MultiAttn}(\mathbf{h}_i^{t-1}, \mathbf{C}_i^t) \\
 \mathbf{h}_i^t &= \text{LayerNorm}(\text{ReLU}(\mathbf{h}_i^t)), i \in [1, n] \\
 \mathbf{s}^t &= \text{MultiAttn}(\mathbf{s}^{t-1}, [\mathbf{s}^{t-1}; \mathbf{H}^t]) \\
 \mathbf{s}^t &= \text{LayerNorm}(\text{ReLU}(\mathbf{s}^t))
 \end{aligned}$$

## 7 Retrieve, Read, Rerank: Towards End-to-End Multi-Document Reading Comprehension

这篇文章是基于Multi-Document阅读理解的，主要改进是将原始的retriever-reader-reranker的管道式结构，改为一个统一的end2end的模型。因为管道式方法，相互割裂，上游一些优秀的representation无法有效的传递到下游，因为下游会对其重新encoding。

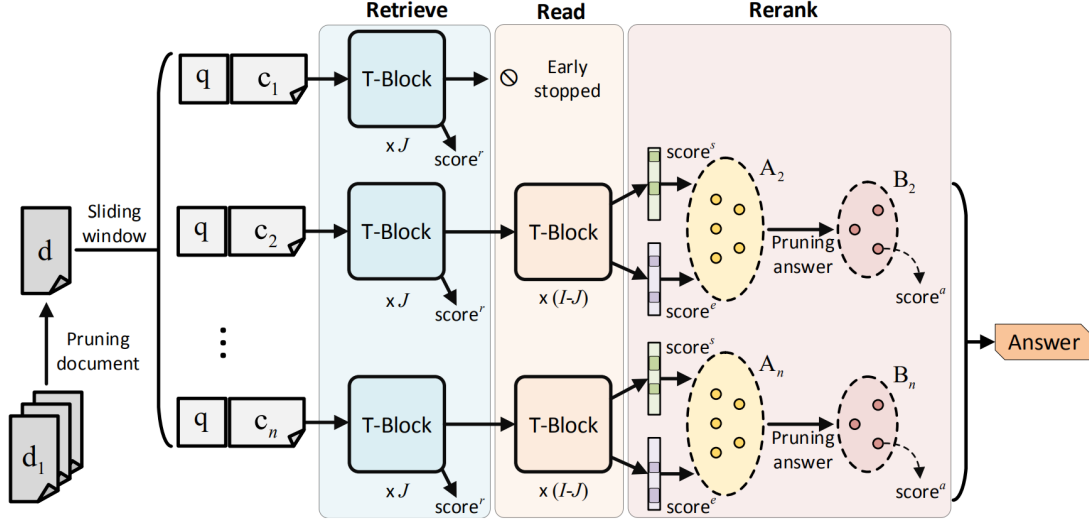


Figure 1: RE<sup>3</sup>QA architecture. The input documents are pruned and splitted into multiple segments of text, which are then fed into the model<sup>2</sup>. Few top-ranked segments are retrieved and the rest are early stopped. Multiple candidate answers are proposed for each segment, which are later pruned and reranked. RE<sup>3</sup>QA has three outputs per candidate answer: the retrieving, reading, and reranking scores. The network is trained end-to-end with a multi-task objective. “T-Block” refers to pre-trained Transformer block (Devlin et al., 2018).

1. Document Pruning:利用TF-IDF余弦距离。

2. Segment Encoding

$$\mathbf{h}^i = \text{Transformer Block}(\mathbf{h}^{i-1}), \forall i \in [1, I] \quad (7.1)$$

3. Early-Stopped Retriever

$$\begin{aligned} \mu &= \text{softmax}(\mathbf{w}_\mu \mathbf{h}^J) \\ \text{score}^r &= \mathbf{w}_r \tanh\left(\mathbf{w}_r \sum_{i=1}^{L_x} \mu_i \mathbf{h}_i^J\right) \end{aligned} \quad (7.2)$$

$$\mathcal{L}_I = - \sum_{i=1}^2 \mathbf{y}_i^r \log(\text{softmax}(\text{score}^r)_i) \quad (7.3)$$

4. Distantly-Supervised Reader

$$\text{score}^s = \mathbf{w}_s \mathbf{h}^I, \text{score}^e = \mathbf{w}_e \mathbf{h}^I \quad (7.4)$$

$$\begin{aligned} \mathcal{L}_{II} &= - \sum_{i=1}^{L_x} \mathbf{y}_i^s \log(\text{softmax}(\text{score}^s)_i) \\ &\quad - \sum_{j=1}^{L_x} \mathbf{y}_j^e \log(\text{softmax}(\text{score}^e)_j) \end{aligned} \quad (7.5)$$

## 5. Answer Reranker

$$\begin{aligned} \eta &= \text{softmax}(\mathbf{w}_\eta \mathbf{h}_{\alpha_i; \beta_i}^I) \\ \text{score}_i^a &= \mathbf{w}_a \tanh\left(\mathbf{W}_a \sum_{j=\alpha_i}^{\beta_i} \eta_{j-\alpha_i+1} \mathbf{h}_j^I\right) \end{aligned} \quad (7.6)$$

$$\mathcal{L}_{III} = - \sum_{i=1}^{M^*} \mathbf{y}_i^{\text{hard}} \log(\text{softmax}(\text{score})_i) + \sum_{i=1}^{M^*} \left\| \mathbf{y}_i^{\text{soft}} - \frac{\text{score}_i^a}{\sum_{j=1}^{M^*} \text{score}_j^a} \right\|^2 \quad (7.7)$$

## 6. Training and Inference

$$\mathcal{J} = \mathcal{L}_I + \mathcal{L}_{II} + \mathcal{L}_{III} \quad (7.8)$$

## 8 Explicit Utilization of General Knowledge in Machine Reading Comprehension

提出了一个叫做Knowledge Aided Reader (KAR)的模型，本质是在Attention上做文章。利用了外部数据：利用WordNet来作为data enrichment的方法。显式的利用了人类的通用知识，即：不同词之间的多义词、次词、全音词、合音词、属性等。

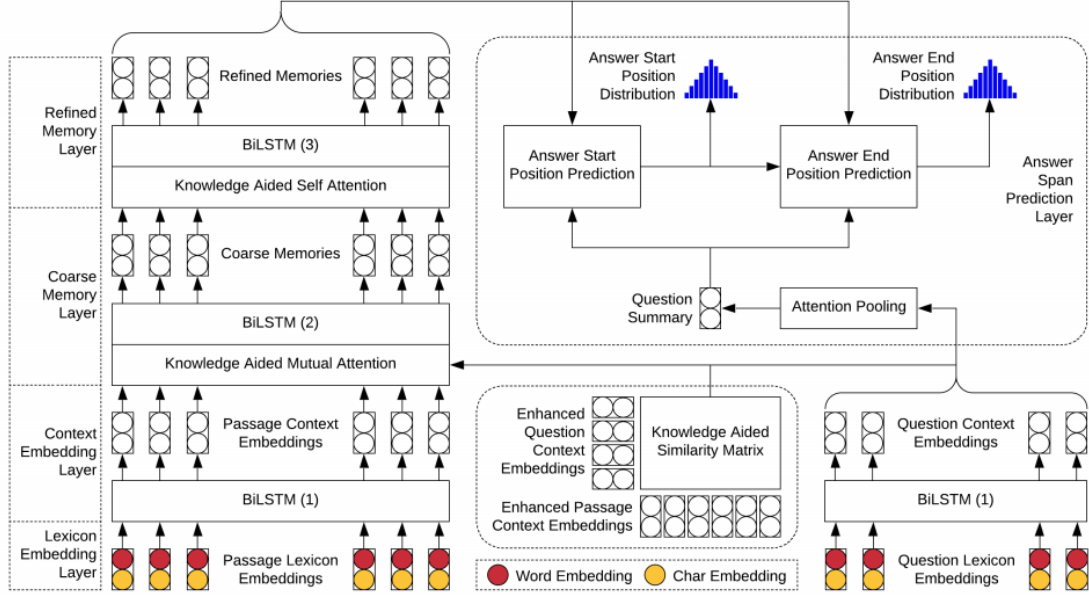


Figure 1: An end-to-end MRC model: Knowledge Aided Reader (KAR)

创新点有二：

### 1. Knowledge Aided Mutual Attention:

传统的计算passage和questions的相似度：

$$f(c_{p_i}, c_{q_j}) = v_f^\top [c_{p_i}; c_{q_j}; c_{p_i} \odot c_{q_j}] \in \mathbb{R} \quad (8.1)$$

本文提出的：

对于每一个单词 $w$ ，首先去构建它的enhanced context embedding:  $c_w^*$ 。对于一个单词，根据Wordnet得到一个集合 $E_w$ ，从而得到 $Z \in \mathbb{R}^{d \times |E_w|}$ 。

$$t_i = v_c^\top \tanh(W_c z_i + U_c c_w) \in \mathbb{R} \quad (8.2)$$

$$c_w^+ = Z \text{softmax}(\{t_1, \dots, t_{|E_w|}\}) \in \mathbb{R}^d$$

然后将 $c_w$ 和 $c_w^+$ 连接起来送入有着ReLU激活的dense layer，输出为d-dimension，得到 $c_w^* \in \mathbb{R}^d$ 。

$$f^*(c_{p_i}, c_{q_j}) = v_f^\top [c_{p_i}^*; c_{q_j}^*; c_{p_i}^* \odot c_{q_j}^*] \in \mathbb{R} \quad (8.3)$$

$$A_{i,j} = f^*(c_{p_i}, c_{q_j}) \quad (8.4)$$

$$R_Q = C_Q \text{softmax}_r^\top(A) \in \mathbb{R}^{d \times n} \quad (8.5)$$

$$R_P = C_P \text{softmax}_c(A) \text{softmax}_r^\top(A) \in \mathbb{R}^{d \times n}$$

## 2. Knowledge Aided Self Attention:

总体上思想类似:

$$\begin{aligned} t_i &= v_g^\top \tanh(W_g z_i + U_g g_{p_i}) \in \mathbb{R} \\ g_{p_i}^+ &= Z \operatorname{softmax}\left(\left\{t_1, \dots, t_{|E_{p_i}|}\right\}\right) \in \mathbb{R}^d \end{aligned} \quad (8.6)$$

## 9 Multi-Style Generative Reading Comprehension

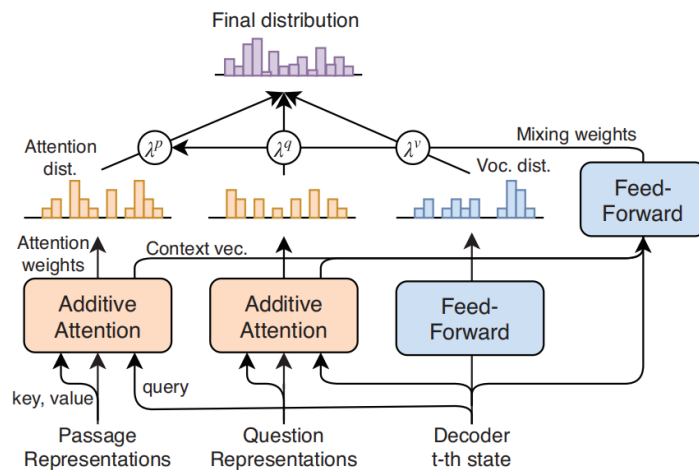
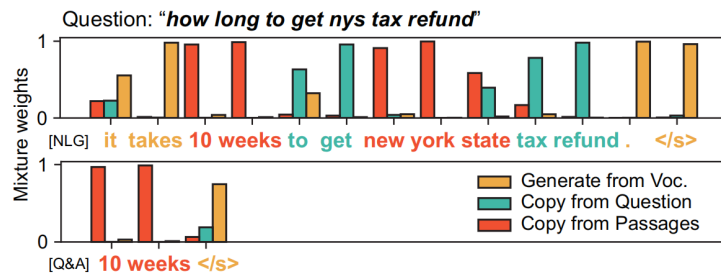
目前，本文是MARCO数据集NLG任务的第一名，文章突破了抽取式阅读理解模型的束缚，采用摘要的方式生成指定风格的答案，这两个突破算是GMRC的一个里程碑。本文针对的是MARCO数据集，简单介绍一下该数据集的特点。用一句话来概括就是——开放域、多文档、生成式的大规模阅读理解数据集，问题和答案都来自于真实数据，所有问题的答案都是人类生成的，有一定的答案还需要额外的人工评估，也就是well-formed的答案。

作者发现之前在MARCO上刷榜的模型大都是抽取式的模型，但是先抽取再生成的pipeline框架是有点弱了。并且现有模型不能根据给定的风格生成相应的答案。论文发现，通常Q&A任务的答案比较简短，NLG任务的答案比较详细。本文想从这两个点入手，提出了两个很新奇的点子：

1. 把生成式阅读理解当作摘要问题来做，生成问题、段落的摘要作为最终的答案。
2. 给定模型风格，训练模型使得其有能力生成相应风格的答案。

提出了一个端到端的模型——masque【风格可控的多源头摘要式模型】，模型根据给定的风格生成答案序列，生成词的来源可以是问题、段落和词表。一共分为四个部分：

1. The question-passages reader: 建模问题与段落之间的关系；
2. The passage ranker: 找寻与问题相关的段落；
3. The answer possibility classifier: 识别出可回答的问题；
4. The answer sentence decoder: 根据给定的风格生成词序列。



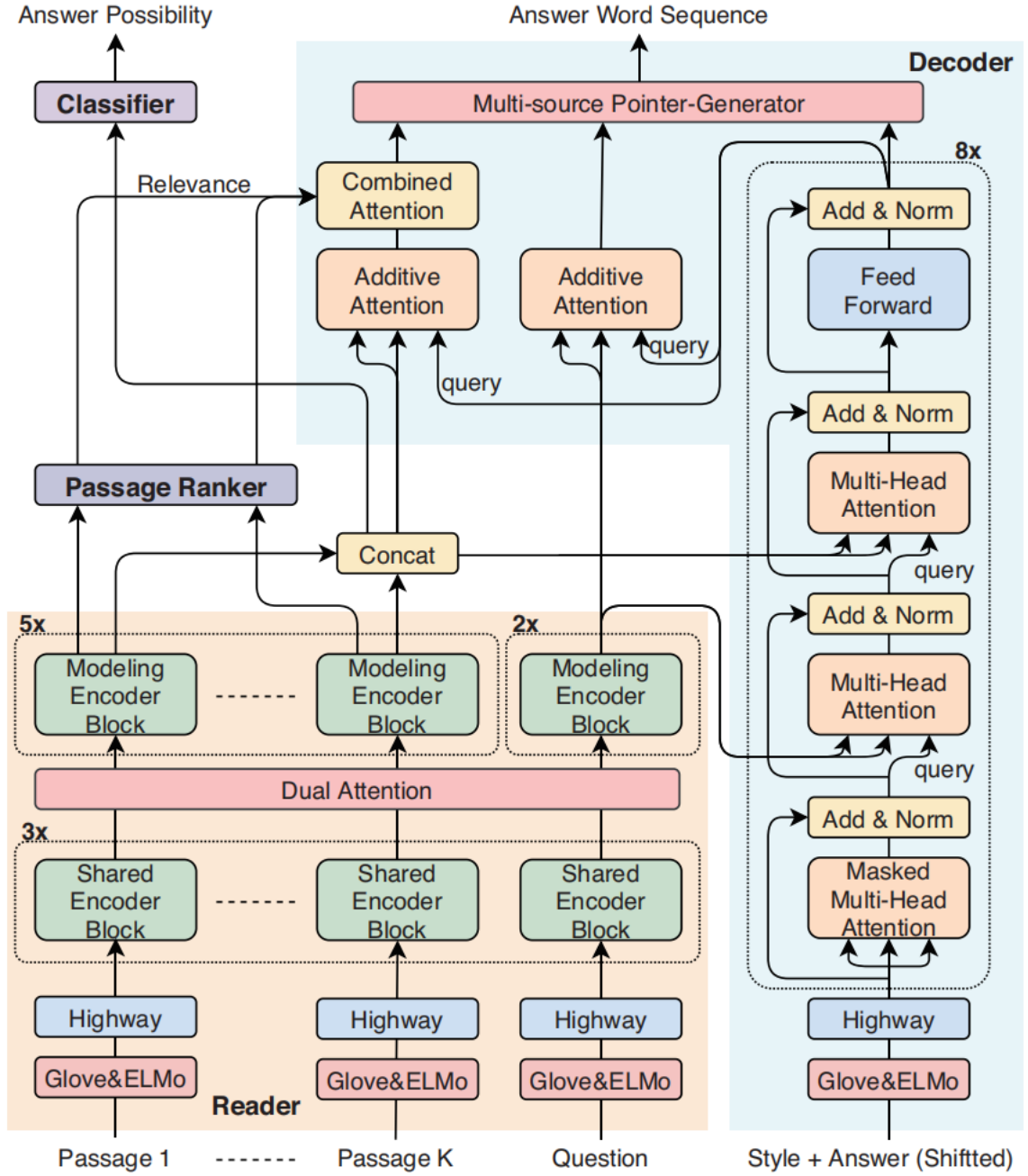
$$P^v(y_t) = \text{softmax}\left(W^{2^\top}(W^1 s_t + b^1)\right) \quad (9.1)$$

$$P^q(y_t) = \sum_{j: x_j^q = y_t} \alpha_{tj}^q \quad (9.2)$$

$$P^p(y_t) = \sum_{l: x_l^{pk(l)} = y_t} \alpha_{tl}^p \quad (9.3)$$

$$P(y_t) = \lambda^v P^v(y_t) + \lambda^q P^q(y_t) + \lambda^p P^p(y_t) \quad (9.4)$$

$$\lambda^v, \lambda^q, \lambda^p = \text{softmax}(W^m[s_t; c_t^q; c_t^p] + b^m)$$





## 10 How Much Reading Does Reading Comprehension Require? A Critical Investigation of Popular Benchmarks

EMNLP2018最佳短论文。

最近在阅读理解问题上有很多研究，它们一般都包含 (question, passage, answer) 元组。大概而言，阅读理解模型必须结合来自问题和文章的信息以预测对应的回答。然而，尽管这一主题非常受关注，且有数百篇论文都希望更好地解决该问题，但许多流行基准的测试难度问题仍未得到解决。在本论文中，我们为 bAbI、SQuAD、CBT、CNN 和 Whodid-What 数据集建立了合理的基线模型，并发现仅带有问题或文章的模型通常有更好的表现。在 20 个 bAbI 任务的 14 个中，仅带有文章的模型实现了高达 50 % 的准确度，它有时能与全模型的性能相匹配。有趣的是，虽然 CBT 提供了 20-sentence 的故事，但只有最后一句能进行相对准确的预测。

bAbI Tasks 1-10										
Dataset	1	2	3	4	5	6	7	8	9	10
True dataset	<b>100%</b>	<b>100%</b>	39%	<b>100%</b>	<b>99%</b>	<b>100%</b>	<b>94%</b>	<b>97%</b>	<b>99%</b>	<b>98%</b>
Question only	18%	17%	22%	22%	34%	50%	48%	34%	64%	44%
Passage only	53%	86%	<b>60%</b>	59%	31%	48%	85%	79%	63%	47%
$\Delta(\min)$	-47	-14	+21	-41	-65	-52	-9	-18	-35	-51

bAbI Tasks 11-20										
	11	12	13	14	15	16	17	18	19	20
True dataset	<b>94%</b>	<b>100%</b>	<b>94%</b>	<b>96%</b>	<b>100%</b>	<b>48%</b>	<b>57%</b>	<b>93%</b>	<b>30%</b>	<b>100%</b>
Question only	17%	15%	18%	18%	34%	26%	48%	91%	10%	70%
Passage only	71%	74%	<b>94%</b>	50%	64%	<b>47%</b>	48%	53%	21%	<b>100%</b>
$\Delta(\min)$	-23	-26	0	-46	-36	-1	-9	-2	-9	0

Table 1: Accuracy on bAbI tasks using our implementation of the Key-Value Memory Networks