# Tell Me How to Ask Again: Question Data Augmentation with Controllable Rewriting in Continuous Space
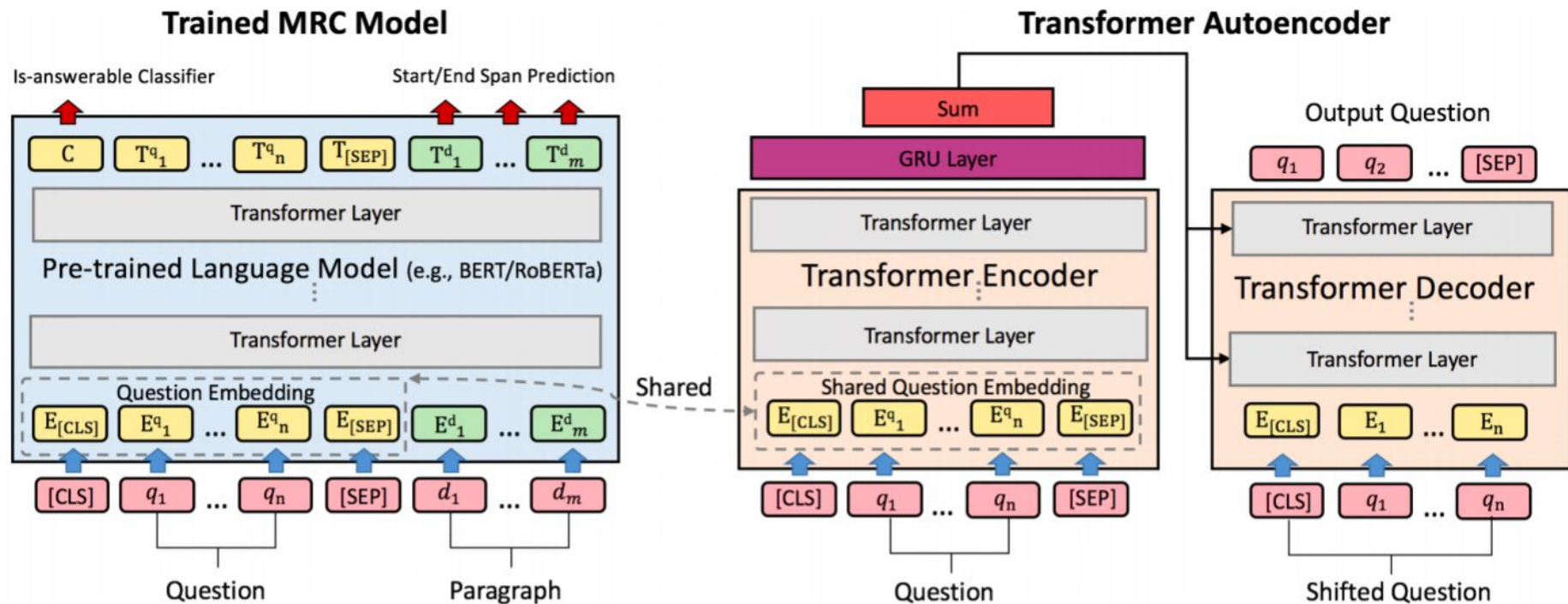
Dayiheng Liu, Yeyun Gong, Jie Fu, Yu Yan,
Jiusheng Chen, Jiancheng Lv, Nan Duan, Ming Zhou

EMNLP 2020

# Task

- Question answering:
  - Given a question
  - Given a document
  - Find the span in the document in which answer is provided
  - Binay classification if the question is answerable or not


- They look for a method to augment the questions:
  - To have the same answer span
  - To be answerable/unanswerable from the document
  - To be similar to the original question

# Model

# Pre-trained Language Model based MRC Model

- A BERT-based model to classify the input question is answerable or not and to find the answer span

$$\mathbf{E}^q, \mathbf{E}^d = \text{BertEmbedding}(q, d),$$

$$P_a(\text{is-answerable}) = \text{Sigmoid}(\mathbf{CW}_c^T + \mathbf{b}_c),$$

$$P_s(i =< \text{start} >) = \text{Sigmoid}(\mathbf{T}_i^d \mathbf{W}_s^T + b_s),$$
$$P_e(i =< \text{end} >) = \text{Sigmoid}(\mathbf{T}_i^d \mathbf{W}_e^T + b_e),$$

$$\mathcal{L}_{\text{mrc}} = \lambda \mathcal{L}_a(t) + \mathcal{L}_s(s) + \mathcal{L}_e(e), \qquad (5)$$
$$= -\lambda \log P_a(t) - \log P_s(s) - \log P_e(e),$$

# Transformer-based Autoencoder

- Encode the question
  - Use the same embedding as the pre-trained model
- Compute a vector representation of the input question
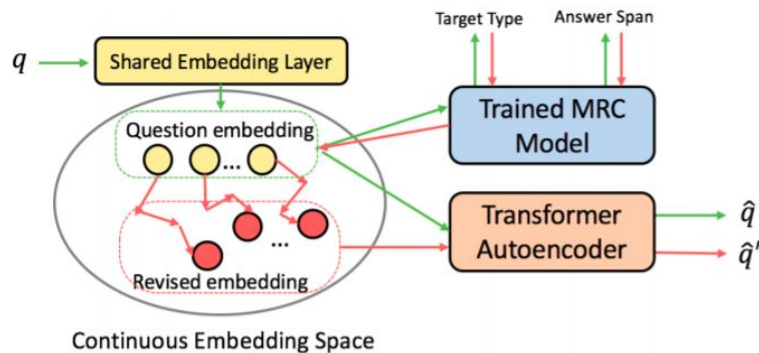- Decode the question vector using a Decoder

$$\mathbf{H}_{enc} = \text{TransformerEncoder}(q),$$
$$\mathbf{z} = \text{Sum}(\text{GRU}(\mathbf{H}_{enc})),$$
$$\hat{q} = \text{TransformerDecoder}(\mathbf{z}).$$

# Rewriting Question with Gradient-based Optimization

- Three objectives for rewriting:
  - Be unanswerable or have the same span
  - Should not be trapped by local optimum
  - Should be similar to Q



- Unanswerable question: $\mathbf{E}^{q'} = \mathbf{E}^q - \eta(\nabla_{\mathbf{E}^q}\mathcal{L}_a(t'))$
- Same Span: $\mathbf{E}^{q'} = \mathbf{E}^q - \eta(\nabla_{\mathbf{E}^q}(\lambda\mathcal{L}_a(t) + \mathcal{L}_s(s) + \mathcal{L}_e(e)))$
- Update step-size for avoiding local optimum
- Use unigram overlap rate for choosing similar questions: $\mathcal{J}(q, \hat{q}') = \dfrac{\text{count}(w_q \cap w_{\hat{q}})}{\text{count}(w_q \cup w_{\hat{q}})},$

# Rewriting Question with Gradient-based Optimization

**Algorithm 1** Question Rewriting with Gradient-based Optimization.

---

**Input:** Data tuple $(q, d, s, e, t)$; Original question embedding $\mathbf{E}^q$; pre-trained MRC model and Transformer autoencoder; A set of step size $S_\eta = \{\eta_i\}$; Step size decay coefficient $\beta_s$; the target answerable or unanswerable label $t'$; Threshold $\beta_t, \beta_a, \beta_b$;

**Output:** a set of new answerable and unanswerable question data tuples $\mathcal{D}' = \{(\hat{q}', d, s, e, t'), .., (\hat{q}', d, s, e, t)\}$;

1: $\mathcal{D}' = \{\}$;
2: **for** each $\eta \in S_\eta$ **do**
3:     **for** max-steps **do**
4:         revise $\mathbf{E}^{q'}$ by Eq. (10) or Eq. (9)
5:         $\hat{q}' = \textbf{TransformerAutoencoder}\left(\mathbf{E}^{q'}\right)$
6:         **if** $P_a(t') > \beta_t$ and $\mathcal{J}(q, \hat{q}') \in [\beta_a, \beta_b]$ **then**
7:             add $(\hat{q}', d, s, e, t')$ to $\mathcal{D}'$;
8:         **end if**
9:         $\eta = \beta_s \eta$;
10:     **end for**
11: **end for**
12: **return** $\mathcal{D}'$;

---

# Results

| Methods | EM | F1 |
|---|---|---|
| BERT$_{large}$ (Devlin et al., 2018) (original) | 78.7 | 81.9 |
| + EDA (Wei and Zou, 2019) | 78.3 | 81.6 |
| + Back-Translation (Yu et al., 2018) | 77.9 | 81.2 |
| + Text-VAE (Liu et al., 2019a) | 75.3 | 78.6 |
| + AE with Noise | 76.7 | 79.8 |
| + 3M synth (Alberti et al., 2019) | 80.1 | 82.8 |
| + UNANSQ (Zhu et al., 2019) | 80.0 | 83.0 |
| + CRQDA (ours) | **80.6** | **83.3** |

# Results

| Methods | EM | F1 |
|---|---|---|
| BERT$_{base}$ | 73.7 | 76.3 |
| + CRQDA | **75.8** (+2.1) | **78.7** (+2.4) |
| BERT$_{large}$ | 78.7 | 81.9 |
| + CRQDA | **80.6** (+1.9) | **83.3** (+1.4) |
| RoBERTa$_{base}$ | 78.6 | 81.6 |
| + CRQDA | **80.2** (+1.6) | **83.1** (+1.5) |
| RoBERTa$_{large}$ | 86.0 | 88.9 |
| + CRQDA | **86.4** (+0.4) | **89.5** (+0.6) |

# Thanks