

Context-Aware Interaction Network for Question Matching



Zhe Hu¹, Zuohui Fu², Yu Yin³, and Gerard de Melo⁴

¹Baidu Inc, ²Rutgers University ³Northeastern University, ⁴HPI/University of Potsdam





Digital Engineering • Universität Potsdam

TASK

MOTIVATION

- Question matching aims to predict the semantic relationship given two questions
 - How do I know if my phone is tapped ?
 - How do I check if my phone is tapped ? dupli



• Cross-attention is widely adopted for text matching



• Computes a **word-by-word attention matrix** to obtain alignments between two sequences

How does a landline call a cell phone ? non-duplicate

METHOD

• We propose COIN: a COntext-aware Interaction Network





cross-attention

- Each value of the attention matrix is based on just two **individual tokens** from the sequences
- Mostly focus on word-level local matching and fail to fully account for the overall semantics
- We aim to contextualize **cross-attention** for better interaction
 - Accurate matching requires a deeper understanding of the two sentences along with pertinent linguistic patterns and constructions
 - We enables the model to consult contextual information while computing the attention matrix to measure the word relevance
 - Yields better contextualized alignments for

Context-Aware Cross-Attention Layer

- We want to integrate **contextual features** C
- We apply a **self-alignment layer** to aggregate pertinent contextual information

Original cross-attention

Context-aware cross-attention

Step 1: compute attention matrix:

 $\begin{aligned} \mathbf{E}_{ij} &= \operatorname{Att}(\mathbf{h}_{a_i}, \mathbf{h}_{b_j}) \\ &= FFN(\mathbf{h}_{a_i})^{\mathrm{T}}FFN(\mathbf{h}_{b_j}) \end{aligned}$

Step 2: compute alignments

$$\mathbf{a}_{i} = \operatorname{softmax}(\mathbf{E}_{i:}), \quad \mathbf{b}_{j} = \operatorname{softmax}(\mathbf{E}_{:j})$$
$$\mathbf{h}_{b_{j}}' = \sum_{k=1}^{m} \mathbf{b}_{kj} \mathbf{h}_{a_{k}}, \quad \mathbf{h}_{a_{i}}' = \sum_{k=1}^{n} \mathbf{a}_{ik} \mathbf{h}_{b_{k}}$$

 $\mathbf{C}_{a} = \operatorname{Layer}_{\operatorname{self}-\operatorname{align}}(\mathbf{H}_{a})$ $\mathbf{C}_{b} = \operatorname{Layer}_{\operatorname{self}-\operatorname{align}}(\mathbf{H}_{b})$

 $\begin{aligned} \mathbf{E}_{ij}^{\mathrm{c}} &= \mathrm{Att}_{\mathrm{context}}(\mathbf{h}_{a_i}, \mathbf{h}_{b_j}, \mathbf{c}_{a_i}, \mathbf{c}_{b_j}) \\ &= FFN(\mathbf{h}_{a_i} + \mathbf{c}_{a_i})^{\mathrm{T}}FFN(\mathbf{h}_{b_j} + \mathbf{c}_{b_j}) \end{aligned}$

Contextual features ($C_a \& C_b$) of two sequences are directly considered when computing the attention matrix



semantic reasoning

Model

EXPERIMENTS

• Datasets: Quora Question Pairs & LCQMC

• **Results** :

How

does

Model	Acc (%)	F1 (%)
Lattice-CNN	82.1	82.4
ESIM	82.0	84.0
BiMPM	83.3	84.9
GMN	84.6	86.0
COIN (Ours)	85.6	86.5
BERT	85.7	86.8
Sentence-BERT	85.4	86.6
COIN (5-run ensemble)	86.2	87.0

BiMPM 88.2 1.6M DIIN 89.0 4.4M CAFE 88.7 4.7M **OSOA-DFN** 89.0 10.0M RE2 89.2 2.8M **89.4** 6.5M COIN (ours) BERT 90.1 109.5M Sentence-BERT 90.6 109.5M COIN (ensemble) **90.7** 32.5M

Acc. (%)

Params

Table 1: Experimental results on LCQMC.

Table 2: Experimental results on Quora dataset.

- Better results than non-pretrained methods
- Comparable results with pre-trained Methods with fewer parameters

3rd alignment

• 5-run ensemble model outperforms BERT and Sentence-BERT

This mirrors human behavior that people tend to <u>first read each sentence and</u> <u>pay attention to the salient contents</u>, and then <u>compare and match two</u> <u>sentences</u>.

Gated Fusion Layer

• update sequence representations by blending alignments

 $\begin{array}{c} \boldsymbol{f}_{i} = \sigma(\mathbf{W}_{1}\mathbf{h}_{a_{i}} + \mathbf{W}_{2}\widetilde{\mathbf{h}}_{a_{i}} + \mathbf{b}_{g}) & \longleftarrow & gate \\ \begin{array}{c} \textit{updated} \\ \textit{representation} \end{array} & \widehat{\mathbf{h}}_{a_{i}} = \boldsymbol{f}_{i} \odot \mathbf{h}_{a_{i}} + (\mathbf{1} - \boldsymbol{f}_{i}) \odot \widetilde{\mathbf{h}}_{a_{i}} \end{array}$



- We first compare the original representations and the aligned ones from difference perspectives
- We use gate operation to enable the model to flexibly incorporate aligned features by **controlling gates**;
- Gate operation is similar to a **skip connection** in mitigating the additional model complexity coming from the deeper structure

• Visualization of Attention:



1st alignment

- model learns to refine the alignments from low-level to high-level;
- the structured phrase such as " what do you think of " is also connected.

Ablation:

Model	Quora	LCQMC
original	89.6	85.4
w/o context	89.1	84.8
simple fusion	88.8	85.2
w/o aggregat.	89.2	84.9
simple pool	89.4	85.2

- Ablation results confirms the effectiveness of each module;
- Removing context-aware alignments brings perform decrease on both datasets.