LREC-COLING 2024



Typos Correction Training Against Misspellings from Text-to-Text Transformers

Guicai Xie¹, Ke Zhang¹, Lei Duan², Wei Zhang¹, Zeqian Huang¹

¹Machine Learning Platform Department, Tencent Inc., China ²School of Computer Science, Sichuan University, Chengdu, China

20-25 May, 2024

Contents

- 1. Background
- 2. Motivation
- 3. Our method: ToCoTR
- 4. Results & Analysis
- 5. Conclusions

Background

MLPD

- Misspelling Queries

- **Rand**-(**Insert**, **Delete**, **Sub**): Randomly inserts, deletes, or substitutes a random character. e.g., typo → {typos, typ, type}
- **SwapNeighbor**: Randomly swaps a character with one of its neighbor characters, e.g., typo → tyop.
- **SwapAdjacent**: Randomly swaps a character with one of its adjacent letter on the QWERTY keyboard, e.g., typo → typp.

- Ranking Performance in "Drop"

Typoed queries resulting from the users' mistyping words or phonetic typing errors exist widely in search behaviors.





MRR@10 and R@1000 results on MSMARCO. A significant drop in effectiveness across different types of simulated typos on queries.

Background

MLPD

- Misspelling Queries

- **Rand**-(**Insert**, **Delete**, **Sub**): Randomly inserts, deletes, or substitutes a random character. e.g., typo → {typos, typ, type}
- SwapNeighbor: Randomly swaps a character with one of its neighbor characters, e.g., typo → tyop.
- **SwapAdjacent**: Randomly swaps a character with one of its adjacent letter on the QWERTY keyboard, e.g., typo → typp.

- Ranking Performance in "Drop"

Typoed queries resulting from the users' mistyping words or phonetic typing errors exist widely in search behaviors.





MRR@10 and R@1000 results on MSMARCO. A significant drop in effectiveness across different types of simulated typos on queries.

Motivation

MLPD

Original Query

- Existing Methods



- Contrastive learning to push a typoed query close to its original variation: BERT + Aug + CL
- Self-teaching with supervised label in KL-divergence: CharacterBERT + ST
- Local ranking alignment: RoDR

Observation 1: Simply aligning the latent embeddings or ranking differences between the original and misspelled queries is inadequate for sophisticated training retrievers such as coCondenser and SimLM.

Observation 2: The spell-checker and dense retriever are optimized as separate models. If the spell-checker is subpar, this will result in a decreased ranking performance.

Motivation

MLPD

- Existing Methods



How to effectively incorporate the spelling correction objective into the dual-encoder dense retriever?

- Base model selection?
- How to improve the effectiveness of the spell correction?
- Training strategy?

Observation 1: Simply aligning the latent embeddings or ranking differences between the original and misspelled queries is inadequate for sophisticated training retrievers such as coCondenser and SimLM.

Observation 2: The spell-checker and dense retriever are optimized as separate models. If the spell-checker is subpar, this will result in a decreased ranking performance.

[2] Ni J, Abrego G H, Constant N, et al. Sentence-T5: Scalable Sentence Encoders from Pre-trained Text-to-Text Models. ACL (Findings) 2022.

- Base model selection

У1

- Spelling correction is formulated as a monolingual translation task and treated with an encoder-decoder based model^{[1].}
- The encoder-decoder sentence embedding model have proved to be a promising architecture^[2].

y₁



y₂

T5 Encoder-Decoder (Spell Correction)

• The T5-style model does not place a special symbol (e.g., [CLS] in BERT) at the beginning of the text sequence.

• To obtain the decoder output, the input text is fed into the encoder, and the standard "start" symbol (first token) is fed as the first decoder input.

La coder La

y₂

T5 Encoder-Decoder First (Sentence Embedding)

^[1] Rothe S, Mallinson J, Malmi E, et al. A Simple Recipe for Multilingual Grammatical Error Correction. ACL 2021.

Our method: ToCoTR

MLPD

- Prompt-based typos correction training



$$\mathcal{L}_{toco} = -\log \sum_{k=1}^{|C|} (\mathbf{P}(Y_k \mid X_k^{typo}; \theta))$$

- S1: Given a set of context $C = \{X_1, X_2, ..., X_{|C|}\}$ and randomly sampling samples to conduct the prompt-based typos generation at a predetermined rate value (in our experiment, the rate equals 80%);
- S2: For each selected source text $X_k = (x_1, x_2, ..., x_t)$, choose $\alpha |X_k|$ token positions at random to simulate typos;
- S3: If the *t*-th token is chosen then use a randomly selected typos generator to inject the typos, including RandInsert, RandDelete, RandSub, SwapNeighbor, SwapAdjacent;



S4: Add special symbols to the selected typo token according to the augmentation template "<h>[X] </h>" to highlight the errors;

There'is \longrightarrow <h> There'is </h> abswer. \longrightarrow <h> abswer. </h>

S5: Training with \mathcal{L}_{toco}

MLPD

- Dual-encoder architecture

Relevance Score Distribution



- Given a query q_i, passage retrieval aims to return a sorted list of the n most relevant passages L = [p₁, p₂, ..., p_n] from a large set D = {p_i}^m_{i=1} according to the relevance score of the retrieval model.
- For dense retriever training, we assume a set of binary positive correlation judgments as supervised signals, denoted by

 $R = \langle q_i, p_i^+, \{p_{i,1}^-, p_{i,2}^-, \dots, p_{i,s}^-\} \rangle$ where p_i^+ denotes the relevant passages and p_i^- denotes the irrelevant passages for query q_i .

• To optimize the dense retriever, the negative log-likelihood (NLL) loss is applied:

$$\mathcal{L}_{nll}\left(\left\langle q_i, p_i^+, \left\{p_{i,1}^-, p_{i,2}^-, ..., p_{i,s}^-\right\}\right\rangle\right) \\ = -\log\frac{e^{sim(q_i, p_i^+)}}{e^{sim(q_i, p_i^+)} + \sum_{j=1}^s e^{sim(q_i, p_{i,j}^-)}}$$

Our method: ToCoTR



 \mathcal{L}_{toco}

 \mathcal{L}_{nll}

Initialize





- Datasets

Dataset	Train	Dev	Test	Avg. q length
MSMARCO	400,782	6,980	6,837	6.10
TREC 2019§	-	-	43	5.60
ANTIQUE	2,426	-	200	10.82

• MSMARCO

Source: queries sampled from Bing search logs and annotated with binary relevant passages **Typo Query**: randomly generated 10 sets by repeating the typos simulation

• TREC 2019

Source: queries sampled from Bing search logs and annotated with four-level relevant annotations (Same corpus with MSMARCO) **Typo Query**: randomly generated 10 sets by repeating the typos simulation

• ANTIQUE (zero-shot validate)

Source: non-factual questions and answers from a community answering service, where questions and answers are manual fourlevel relevance annotations

Typo Query: sampling three misspelling variations from manually validated typoed questions by released researchers, including SwapNeighbor, SwapAdjacent, RandSub



- Comparison with the retrievers

Queries	Methods	PLMs	TREC 2019		MSMARCO		ANTIQUE			
duonice			nDCG@10	MRR	MAP	MRR@10	R@1000	nDCG@10	MAP	R@1000
w/o typos	 a) BM25 (pyserini) b) DPR c) CharacterBERT d) Sentence-T5 e) coCondenser[†] f) SimLM[†] g) BERT+Aug h) BERT+Aug+CL i) CharacterBERT+ST j) RoDR 	- BERT $_{base}$ CharacterBERT T5 $_{base}$ Condenser BERT $_{base}$ BERT $_{base}$ BERT $_{base}$ CharacterBERT BERT $_{base}$	50.6 59.7 61.6 64.3 71.5 71.4 61.8 61.1 <u>63.9</u> 62.1	70.4 72.5 78.5 82.4 86.8 87.9 80.7 77.1 80.7 78.8	30.1 35.2 33.0 36.4 45.3 46.9 35.3 34.5 33.6 34.6	18.7 32.6 32.1 32.1 <u>38.3</u> 41.1 [‡] 32.7 <u>32.8</u> 32.6 32.8	85.7 95.2 94.8 95.9 98.4 98.7 [‡] <u>95.1</u> 94.8 94.6 95.1	23.7 26.5 25.4 <u>30.0</u> 32.4 [‡] - 26.6 <u>26.7</u> 26.6 26.1	15.9 19.0 17.9 <u>22.1</u> 24.4 [‡] - <u>19.6</u> 19.5 19.0 19.0	46.1 57.9 55.2 61.8 66.9 [‡] - <u>59.7</u> 58.5 53.8 59.7
	k) ToCoTR (two-stage)	T5 _{base}	69.4	87.8	41.0 [‡]	34.4 ‡	96.6 [‡]	28.6	20.1	60.5
w/ typos	 l) BM25 (pyserini) m) DPR m) CharacterBERT o) Sentence-T5 p) coCondenser[†] q) SimLM[†] 	- BERT _{base} CharacterBERT T5 _{base} Condenser BERT _{base}	25.9 28.4 35.9 40.5 46.1 45.9	35.3 41.3 53.5 56.4 60.2 60.2	15.0 15.9 18.9 21.8 <u>27.6</u> 27.7	9.5 13.5 16.0 18.5 <u>22.1</u> 23.6 [‡]	61.1 68.2 72.2 80.6 84.5 [‡] 83.8	16.9 16.4 17.2 <u>20.8</u> 24.5 [‡]	11.3 11.4 12.4 <u>15.0</u> 18.1 [‡]	38.0 42.9 42.7 52.3 56.2 [‡]
	 r) BERT+Aug s) BERT+Aug+CL t) CharacterBERT+ST u) RoDR v) ToCoTR (two-stage) 	$\begin{array}{c} BERT_{base} \\ BERT_{base} \\ CharacterBERT \\ BERT_{base} \\ T5_{base} \end{array}$	42.9 43.9 <u>52.0</u> 43.9 63.4 [‡]	59.4 59.8 70.2 59.0 83.5	23.8 24.1 <u>27.0</u> 23.9 36.6 [‡]	21.8 22.9 <u>26.4</u> 23.2 31.3 [‡]	84.2 85.6 <u>89.2</u> 86.2 94.6 [‡]	20.6 21.0 <u>22.2</u> 21.0 26.1 [‡]	14.8 15.1 <u>15.8</u> 15.0 18.6 [‡]	50.2 49.7 48.0 51.0 57.1 [‡]

 \ddagger indicate significant differences with the second-best score (underlined) at *p*-value < 0.05



- Comparison with the retrievers involving spell-checkers

Methods	w/o ty	ypos	w/ typos		
methode	MRR@10	R@1000	MRR@10	R@1000	
CharacterBERT	32.1	94.8	16.0	72.2	
pyspellchecker $ ightarrow$ CharacterBERT	27.3	88.5	23.0	81.9	
MS-Spellchecker \rightarrow CharacterBERT	32.0	94.6	29.9	91.3	
$\textbf{GG-Spellchecker} \rightarrow \textbf{CharacterBERT}$	32.2	94.8	29.4	90.4	
Sentence-T5	32.1	95.9	18.5	80.6	
pyspellchecker $ ightarrow$ Sentence-T5	28.4	91.9	24.7	87.0	
MS-Spellchecker \rightarrow Sentence-T5	32.0	95.9	30.5	93.6	
GG-Spellchecker \rightarrow Sentence-T5	32.0	95.9	29.9	93.1	
ToCoTR (two-stage)	34.4 ‡	96.6 [‡]	31.3	94.6 ‡	

 \ddagger indicate significant differences with the second-best score (underlined) at *p*-value < 0.05

- pyspellchecker: a rule-based spell-checking toolkit that relies on dictionary-based rule sets;
- MS-Spellchecker: Microsoft Bing Spell Check API, it utilizes machine learning and statistical machine translation to provide corrections;
- GG-Spellchecker: Google Search API, it has been shown in previous research to be possibly the most useful spell corrections^[1];

^[1] Hagen M, Potthast M, Gohsen M, et al. A large-scale query spelling correction corpus. ACM SIGIR 2017.

- Ablation study

Methods	w/o ty	ypos	w/ typos		
motrioue	MRR@10	R@1000	MRR@10	R@1000	
ToCoTR	34.4	96.6	31.3	94.6	
w/o Hard Negative	33.5↓	95.7↓	29.8↓	93.2↓	
w/o Self-Teaching	34.1	97.0	27.3↓	91.9↓	
w/o Prompt	34.1	96.7	30.6↓	94.3	
w/o ToCo	33.9	96.6	28.7↓	92.4↓	

Statistically significant drops at *p*-value < 0.05 are marked with \downarrow .

- Analysis on typos correction training

Joint or Two-stage?

Methods	w/o ty	/pos	w/ typos		
methode	MRR@10	R@1000	MRR@10	R@1000	
Two-stage Joint	33.5 ‡ 32.8	95.7 95.7	29.8 ‡ 26.8	93.2 ‡ 90.8	

Table 5: The results of two different ways for incorporating typos correction training on MSMARCO. The hyper-parameter β of balance \mathcal{L}_{nll} and \mathcal{L}_{toco} is set as 0.1 in joint training. Statistically significant differences at *p*-value < 0.05 are marked with \ddagger .

- Without typos correction training (namely, w/o ToCo), the performance of ToCoTR on both original and typoed queries is greatly affected.
- The adaption of hard negative mining can demonstrate why advanced retrievers with sophisticated training can outperform most typos-aware retrievers.



The impact of the number of typos (α).

The more typos injected, the worse performance becomes.

The proportion of typos per input text α is set as 0.2.

MLPD

- Latency Analysis

	API call	Query encoding	Index search
MS-Spellchecker	271.3ms	-	-
GG-Spellchecker*	437.2ms	-	-
DPR	-	0.2ms	226.1ms
CharacterBERT	-	0.4ms	233.6ms
Sentence-T5	-	1.2ms	223.5ms
ToCoTR	-	1.2ms	235.7ms

Table 6: Latency analysis for two spell-checker APIs and four retrieval systems. We retrieve the top 1000 results for MSMARCO dev queries with a single thread and then average over all the queries. The latency for two spell-checker API calls depends on server load and is difficult to precisely measure. For example, the latency of GG-spellchecker to call the API to return the full web page (not only the corrected queries) is marked \star .

Typoed Queries	ToCoTR		coCor	ndenser	SimLM		
.,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,,	MRR@10	Recall@1000	MRR@10	Recall@1000	MRR@10	Recall@1000	
w/o Typos	34.4	96.6	38.3	98.4	41.1	98.7	
w/ Typos (Avg.)	31.2 /-9.20%	94.5 /-2.17%	22.4/-41.5%	84.9/-13.7%	23.9/-41.9%	84.2/-14.7%	
RandInsert	32.4 /-5.79%	95.5 /-1.20%	22.9/-40.1%	85.5/-13.1%	23.8/-42.0%	84.2/-14.7%	
RandDelete	30.6 /-11.1%	94.5/-2.22%	22.7/-40.6%	86.2/-12.4%	24.7/-40.0%	86.0/-12.9%	
RandSub	30.1 /-12.3%	93.6/-3.17%	22.0/-42.6%	84.6/-14.0%	24.2/-41.1%	84.0/-14.9%	
SwapNeighbor	32.1 /-6.69%	95.1/-1.57%	21.6/-43.7%	83.5/-15.1%	22.9/-44.4%	82.5/-16.4%	
SwapAdjacent	30.9 /-10.1%	94.0 /-2.68%	22.8/-40.6%	84.5/-14.1%	23.9/-42.0%	84.2/-14.7%	

- Drop rate in different simulated typos

- ToCoTR can reach a smaller drop rate from original query set to various misspelled query variations.
- Different retrievers exhibit varying degrees of performance degradation with different simulated typos: ToCoTR with RandSub, coCondenser and SimLM with SwapNeighbor.

- We explore different incorporating strategies to conduct typos correction training. This establishes a feasible and effective approach that explicitly incorporates typos correction training into the training pipeline of dense retrieval.
- We propose a simple yet effective prompt-based augmentation technique to enhance the typos correction training. It adaptively realizes the alignment of typoed words to correct words and reduces the difference between typed query embedding and its corresponding correct query embedding.
- We conduct a comprehensive comparison study to show the retrieval effectiveness of ToCoTR on queries with typos across three benchmark datasets. It outperforms those with typos-aware training retrievers and even outperforms the combining solutions involving some advanced spell checkers for dense retrievals.





