

QUERT: Continual Pre-training of Language Model for Query Understanding in Travel Domain Search

Jian Xie*

Shanghai Key Laboratory of Data Science, School of Computer Science, Fudan University
Shanghai, China
jianx0321@gmail.com

Yidan Liang

Alibaba Group
Hangzhou, China
liangyidan.lyd@alibaba-inc.com

Jingping Liu†

School of Information Science and Engineering, East China University of Science and Technology
Shanghai, China
jingpingliu@ecust.edu.cn

Yanghua Xiao†

Shanghai Key Laboratory of Data Science, School of Computer Science, Fudan University
Shanghai, China
shawyh@fudan.edu.cn

Baohua Wu

Alibaba Group
Hangzhou, China
zhengmao.wbh@alibaba-inc.com

Shenghua Ni

Alibaba Group
Hangzhou, China
shenghua.nish@alibaba-inc.com

ABSTRACT

In light of the success of the pre-trained language models (PLMs), continual pre-training of generic PLMs has been the paradigm of domain adaption. In this paper, we propose **QUERT**, A Continual Pre-trained Language Model for **QUERY** Understanding in Travel Domain Search. QUERT is jointly trained on four tailored pre-training tasks to the characteristics of query in travel domain search: Geography-aware Mask Prediction, Geohash Code Prediction, User Click Behavior Learning, and Phrase and Token Order Prediction. Performance improvement of downstream tasks and ablation experiment demonstrate the effectiveness of our proposed pre-training tasks. To be specific, the average performance of downstream tasks increases by 2.02% and 30.93% in supervised and unsupervised settings, respectively. To check on the improvement of QUERT to online business, we deploy QUERT and perform A/B testing on Fliggy APP. The feedback results show that QUERT increases the Unique Click-Through Rate and Page Click-Through Rate by 0.89% and 1.03% when applying QUERT as the encoder. Resources are available at <https://github.com/hsaest/QUERT>.

CCS CONCEPTS

• **Information systems** → **Query representation.**

KEYWORDS

Continual Pre-training, Query Understanding, Travel Domain Search

*Work done when interned at Alibaba Group.

†Jingping Liu and Yanghua Xiao are corresponding authors.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](https://permissions.acm.org).
KDD '23, August 6–10, 2023, Long Beach, CA, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0103-0/23/08...\$15.00
<https://doi.org/10.1145/3580305.3599891>

ACM Reference Format:

Jian Xie, Yidan Liang, Jingping Liu, Yanghua Xiao, Baohua Wu, and Shenghua Ni. 2023. QUERT: Continual Pre-training of Language Model for Query Understanding in Travel Domain Search. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '23)*, August 6–10, 2023, Long Beach, CA, USA. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3580305.3599891>

1 INTRODUCTION

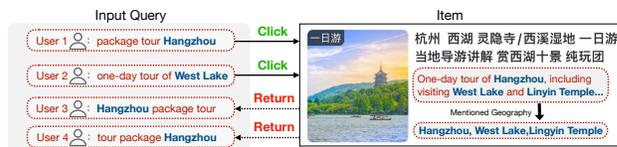


Figure 1: Examples of query and click item. “Click” in green means users click the item, and “Return” in red means the item is expected to return for the query.

Pre-trained language models (PLMs) [6, 12, 17, 30] have become the backbone models in the field of natural language processing (NLP) due to their superior representation ability. Thus, PLMs are widely used in various domains and provide significant performance gains for downstream tasks (e.g., text classification [7] and information extraction [15]).

However, commonly used pre-trained models often perform poorly when directly transferred to specific domains (e.g., travel domain search). This is caused by the mismatch between the corpora in the pre-training stage and the data in the target task. In order to address this problem, previous work proposes to use domain corpora to continually pre-train domain-specific models, like BioBERT [13] for bio-medicine, FinBERT [1] for finance, and COMUS [8] for math problem understanding, etc. With the domain adaption model, the performance of downstream tasks has been significantly improved, which validates the importance of continual pre-training on domain data.

In this paper, we focus on the continual pre-training for query understanding in travel domain search. Travel domain search is the basis of an online travel APP, and there are many studies devoted to this aspect (e.g., named entity recognition [3] and hotel search ranking [29]). Nevertheless, to the best of our knowledge, there is no work focused on PLM in travel domain search. In addition, due to the characteristics of the query in travel domain search, the previous pre-training tasks are not applicable. We analyze the query in travel domain search and summarize three characteristics: (1) **Geography Awareness**. Most user input queries consist of geographical location and intention (e.g., in “package tour Hangzhou”, “package tour” is the intention, and “Hangzhou” is a tourist city in China). But classical MLMs fail to grasp the importance of geography because the random masking strategy treats all tokens equally. Besides, the representation produced by MLM is based on contextual understanding, which means it can not reflect the real physical geography information. (2) . In the search logs, many different queries link to the same click item, which means they have similarity driven by user click behavior. Therefore, PLM in travel domain search is expected to model such a potential similarity. As shown in Figure 1, “package tour Hangzhou” and “one-day tour of West Lake” are two literally different queries, but they point to the same click item. However, the conventional PLMs do not have such an ability because of lacking related pre-training tasks. (3) **Robustness to Phrase and Token Order**. First, due to the user input habit, users might freely permute the phrase order in a query. For example, “package tour | Hangzhou” (“|” denotes the phrase separator) might be entered as “Hangzhou | package tour”. Second, the token orders in the phrase might be transposed because of the user misinput (e.g., “tour package | Hangzhou”). We define the above phenomena as Phrase Permutation and Token Permutation, respectively. In fact, the intentions in these two permutation cases are the same, and the returned results are expected to be the same. Therefore, the model in travel domain is expected to be robust to phrase and token order. However, due to lacking specific pre-training, conventional language models cannot understand the logical consistency in permutation. According to the query characteristics in travel domain search, we propose QUERT to effectively learn query representations through four customized pre-training tasks. Given a query, we introduce its click item as the additional information. Specifically, to solve problem (1), we design a masking strategy called **Geography-aware Mask Prediction (Geo-MP)** to force the pre-trained model to pay more attention to the geographical location phrases. In addition to semantic understanding, we introduce geohash in **Geohash Code Prediction (Geo-CP)** task to model real physical geographic information for the language model. As for problem (2), in order to build a connection between the different queries linking to the same click item, we propose **User Click Behavior Learning (UCBL)** to learn the potential similarity. To solve problem (3), we propose **Phrase and Token Order Prediction (PTOP)**. We shuffle the phrases and tokens to simulate the permutation. QUERT is expected to predict the original phrases and tokens order of the shuffled query. This task aims to enable QUERT to learn the logical consistency in permutation and be robust to phrase and token order.

Our contributions are summarized as follows: 1) To the best of our knowledge, we are the first to explore continual pre-training

for query understanding in the travel domain search. 2) We propose four tailored pre-training tasks: Geography-aware Mask Prediction, Geohash Code Prediction, User Click Behavior Learning, and Phrase and Token Order Prediction. 3) The experimental results on five downstream tasks related to travel domain search prove the effectiveness of our method. In particular, model performance improves by 2.02% and 30.93% under supervised and unsupervised settings, respectively. And the online A/B testing on Fliggy APP¹ demonstrates that QUERT improves the Unique Click-Through Rate and Page Click-Through Rate by 0.89% and 1.03% when applying QUERT as the feature encoder.

2 RELATED WORK

The related work in this paper can be divided into three groups: open domain pre-trained language models, domain adaption pre-trained language models and query search.

Open Domain PLMs. Based on the Transformer [26] architecture, BERT [5] has proven the effectiveness of large-scale corpora pre-training in natural language processing. After that, other Transformer-based models [2, 12, 14, 30] followed, either structural improvements or changes in pre-training tasks. Instead of only considering representations based on token-level masking strategy, Zhang et al. [31] adopt entity-level masking and phrase-level masking as the masking strategy. Likewise, SpanBERT [11] no longer masks individual tokens but continuous spans to enhance the model’s text comprehension.

Domain Adaption PLMs. Gururangan et al. [9] point out that domain adaption PLMs provide large gains in domain tasks. Based on the large corpora of financial texts, FinBERT [1] gains outstanding performance in financial sentiment classification. For Tweets, Nguyen et al. [22] propose BERTweet to improve the performance on several Tweet NLP Tasks. Besides, COVID-Twitter-BERT [21] achieves state-of-the-art (SOTA) performance on five classification tasks. And the largest performance gap is in the COVID-19 classification task. Excellent performance reflects the importance of domain adaption pre-training.

Query Search. For query search, PROP [19] and B-PROP [20] introduce representative words prediction task in the pre-training phase to model query and document. Furthermore, it turns out that incorporating geographic information into a language model can make the model geo-sensitive. For example, Baidu Maps apply ERNIE-GeoL [10], a language model pre-trained on whole word mask and geocoding prediction task, to improve the performance on geographical tasks. In the Point of Interest (POI) search, Liu et al. [16] propose Geo-BERT to learn graph embeddings that simulate the POIs distribution. However, due to the pre-defined single objective, the above work only considers a single characteristic of the query and thus cannot be generalized to other tasks in the search domain. Therefore, we propose custom tasks based on the common characteristics of the travel domain search to improve the generalization of the model.

3 QUERT

Given a query $q = [w_{q,1}, \dots, w_{q,m}]$ and its click item title $c = [w_{c,1}, \dots, w_{c,n}]$, where m and n are the numbers of tokens in q and

¹Fliggy is an online travel agency in China. The official website is www.fliggy.com

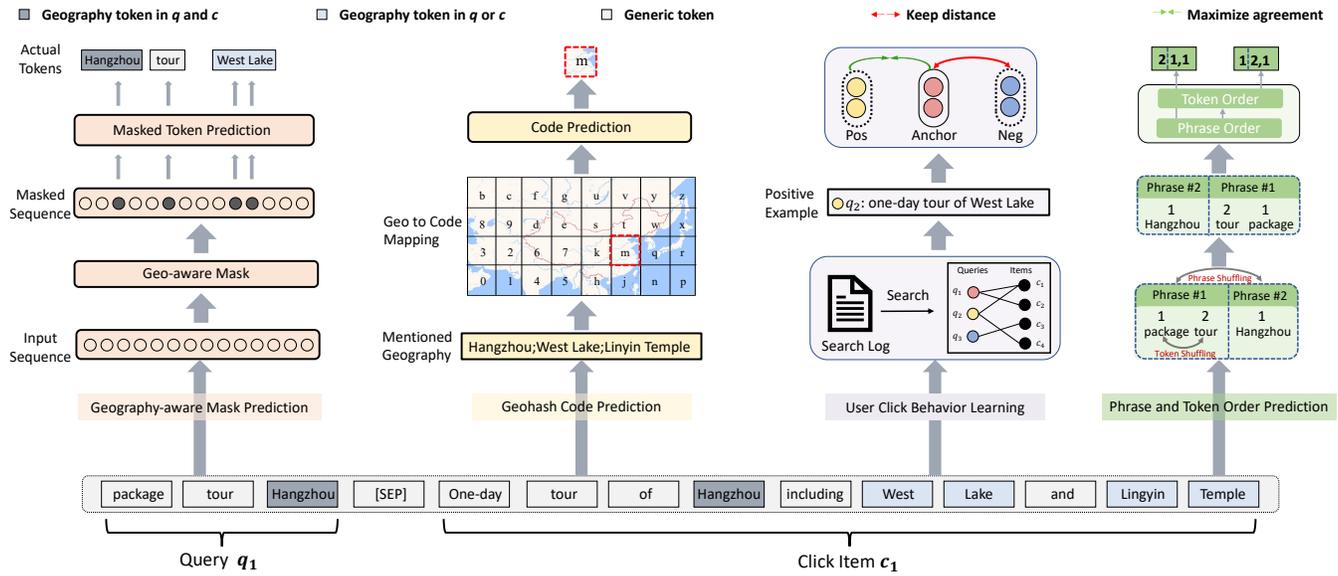


Figure 2: Our framework of QUERT. The Geohash Codes in the figure are only for presentation and do not represent the final code in our implementation.

c , QUERT obtains the contextual representation for each token. As shown in Figure 2, “package tour Hangzhou” is the query, and “One-day tour of Hangzhou including visiting West Lake and Lingyin Temple” is the related click item title. For convenience, we define the phrase as a group of tokens standing for a grammatical unit. For example, in the mentioned query, “package” is a token, while “package tour” is a phrase. According to the characteristics of query in travel domain search, we design four pre-training tasks: **Geography-aware Mask Prediction** (Section 3.1), **Geohash Code Prediction** (Section 3.2), **User Click Behavior Learning** (Section 3.3), and **Phrase and Token Order Prediction** (Section 3.4).

3.1 Geography-aware Mask Prediction

Geography-aware Mask Prediction (Geo-MP) aims to enable QUERT to be aware of geography information. The reason is that we observe that most of the queries contain geography information in the travel domain search. To verify this point, we randomly sample 1,000 queries from Fliggy APP and identify the geography information through inner geography dictionary mapping. The statistical result shows queries containing geography information account for 65%. Therefore, the language models are expected to be good at the representation of geography-related queries. However, the existing pre-trained language models (e.g., BERT) do not have this capability because of lacking specific pre-training tasks. Thus, we propose the Geography-aware Mask Prediction to improve the situation.

Given query q and its corresponding click item c , we use a special token “[SEP]” to combine them. The reasons for incorporating click items are summarized as two points. First, in our statistics, most of the query lengths are concentrated between 1 and 10. Short text length is not conducive to pre-training but increases the risk of weakening the model’s representation ability (detailed experimental proofs are described in Section 4.2.1). Second, the click items

are derived from real user behavior, which verifies the query and item are highly matched. Besides, the length of the item title is appropriate (concentrating between 25 and 35). Thus, the click item is the optimal solution to enrich query information.

Inspired by BERT, we apply masked language model (MLM) to predict the masked tokens, where the difference is we assign a higher probability to geography tokens. To be specific, we use the NER tool² as a detector to identify geography phrases that appear in q and c . We design three mask strategies. First, for geography phrases in both q and c , we set a mask probability of 50% for them. We consider that when common geography phrases are masked at the same time, the randomness of the prediction is excessive due to the lack of context. Ideally, with this probability, for geography phrases in q and c , when one of them is masked, the other is visible. This enables the model to infer the masked phrase from the visible one. Second, we assign a mask probability of 30% to the geography phrases that only appear in q or c . Third, for the rest of the tokens, as with BERT, we mask 15% of them. The masked language model loss is given by:

$$\mathcal{L}_{Geo-MP} = - \sum_{i \in M_w} \log p(w_i | w_{\setminus M_w}), \quad (1)$$

where M_w is the positions of masked tokens.

3.2 Geohash Code Prediction

In addition to Geo-MP, we introduce Geohash Code Prediction (Geo-CP) to enhance the geography sensitivity of QUERT. The reason is that in travel domain search, semantic understanding is not enough for downstream tasks (e.g., Query Rewriting and Query-POI Retrieval). For example, in the Query-POI retrieval task, when the query is “Hangzhou Tour”, “West Lake” is one of the potential

²In this paper, we use AliNLP which is developed by Alibaba Damo Academy.

recall because West Lake is a famous POI located in Hangzhou. However, MLMs based on mask strategy can only understand geographical location from the semantic level and cannot capture the hierarchy or distance relationship between geographical locations. So MLM may recall POIs mentioning “Hangzhou” while ignoring “West Lake”. Thus, QUERT is expected to have the ability to model real geographic locations (e.g., physical location distance and geographical hierarchy). Therefore, we put forward Geo-CP.

Given click query and click items, the objective of Geo-CP is to predict the geohash code produced by the geohash algorithm. The geohash algorithm divides geographic blocks into several grids and encodes them according to latitude and longitude. The code is represented as a string, with precision controlled by the number of bits in the string. Each bit represents a different granularity of geographic information. And the adjacent grids share the same prefix. We assume that the length of geohash code is N bits. In order to encode the geography text, first, we locate the latitude and longitude of every geography entity in items. Geography data uploaded by the service provider is considered to be of high confidence. This is why we only consider the geography in item. Second, we encode latitude and longitude into geohash code. Finally, in order to get a unique encoding for each input, we process them separately according to the number of parsed geography units. 1) For item including no geography unit, we use N bits special token “*” to stand for it. 2) For item including only one geography unit, we adopt its geohash code as the final geohash code. 3) For item including several different geography units, we adopt their longest prefix as the final geohash code. And the parts that are short of N bits are filled with special tokens “*”. In terms of model architecture, we use N independent multi-layer perceptrons (MLPs) to predict bits at different positions. In other words, each MLP has its own granularity prediction capability. The Geo-CP loss is defined as:

$$\mathcal{L}_{Geo-CP} = -\frac{1}{N} \sum_{i=1}^L y_i \log p_i, \quad (2)$$

where L is the number of potential characters in one bit.

3.3 User Click Behavior Learning

In the search logs, we observe that literally different queries may point to the same click item. For instance, “package tour Hangzhou” and “one-day tour of West Lake” are both related to the item “One-day tour of Hangzhou including visiting West Lake and Lingyin Temple”. These two queries are not literally similar, but they have an implicit query similarity driven by user click behavior. Conventional MLMs are unable to model this implicit similarity because of lacking specific pre-training. Therefore, we propose User Click Behavior Learning (UCBL) based on contrastive learning.

Formally, given a query q_i and its click item c_i , according to the click rate, we select the top K from the queries linking to c_i and combine them into a group $G = \{q_{i1}, q_{i2}, \dots, q_{iK}\}$. Then, in order to guarantee the diversity, we randomly choose one from G as the positive example q_i^{pos} of q_i . As for other queries q_j in batch, we regard them and the corresponding positive examples q_j^{pos} as the negative examples for q_i . Feeding q_i and q_i^{pos} into the encoder (i.e., QUERT), we adopt the “[CLS]” embedding to represent the input. The embeddings are represented as \mathcal{R}_i and \mathcal{R}_i^{pos} respectively. The

optimizing objective can be expressed as following:

$$\mathcal{L}_{UCBL} = -\log \frac{\exp(\text{sim}(\mathcal{R}_i, \mathcal{R}_i^{pos})/\tau)}{\sum_{j=1}^N \exp(\text{sim}(\mathcal{R}_j, \mathcal{R}_j^{pos})/\tau)}, \quad (3)$$

The $\text{sim}(\cdot)$ function which we use is cosine similarity, and τ is used for smoothing the distribution.

3.4 Phrase and Token Order Prediction

Phrase and Token Order Prediction (PTOP) is designed to enable QUERT to learn the logical consistency in permutation, thus being robust to phrase and token order. According to our observation, the permutation query in the travel domain has the following two types. 1) Phrase Permutation. Queries are presented in different forms because of users’ personalized input habits. For example, the query “package tour Hangzhou” would be entered as “Hangzhou package tour”. This would lead to differences in the returned results due to the change of query. In essence, for the same intention queries, the returned results are expected to be the same. 2) Tokens in Phrase Permutation. User misinput causes the token permutation in the query. Taking the same case “package tour Hangzhou” for example, wrong results will be returned if the user enters “tour package Hangzhou”. In our statistics, the inappropriate returned results caused by token permutation account for 5.3% in 5000 randomly selected bad cases. However, conventional PLMs are not good at modeling such logical consistency in permutation. To this end, we propose Phrase and Token Order Prediction in this paper.

Given the original query “package tour Hangzhou”, QUERT is expected to predict every token’s corresponding phrase order and token order after permutation. We pre-define the phrase order for every token as (1, 1, 2). “package” and “tour” are both in the first phrase, so their phrase orders are both “1”. And the token order for every token is (1, 2, 1). Note that every token order is limited to the phrase to which the token belongs. In other words, the maximum token order is no greater than the length of its corresponding phrase. Specifically, in the phrase “package tour”, “package” is the first token, and “tour” is the last token, so their token order is “1” and “2” respectively. And “Hangzhou” belongs to the second phrase, so its token order is back to “1”.

To simulate the permutation, first, we randomly shuffle phrases in a query. So the “package tour Hangzhou” would be “Hangzhou package tour”, and the ground truth phrase order is defined as $y = (2, 1, 1)$. Second, we randomly select phrases with a specific probability and shuffle the tokens in them. Under this setting, we assume the selected phrase is “package tour”. After shuffling, the phrase would be “tour package”. And the final shuffled query is “Hangzhou tour package”. Therefore, in this case, the ground truth of the token order is $y = (1, 2, 1)$. Lastly, the output order is computed by the tokens layer and phrase layer, respectively. Since token order is predicted based on the phrase, we design the tokens layer following phrase layer:

$$\begin{aligned} [r_1, \dots, r_i, \dots, r_m] &= \text{QUERT}([w_{q,1}, w_{q,2}, \dots, w_{q,m}]), \\ p_i^\alpha &= \text{Softmax}(\text{MLP}_{\text{phrase}}(r_i)), \\ p_i^\beta &= \text{Softmax}(\text{MLP}_{\text{token}}(\text{MLP}_{\text{phrase}}(r_i))). \end{aligned} \quad (4)$$

The training formulation is given by:

$$\mathcal{L}_{TPOP} = -\frac{1}{m} \sum_{i=1}^m \left(\sum_{c=1}^Q y_{i,c}^\alpha \log(p_{i,c}^\alpha) + \sum_{d=1}^R y_{i,d}^\beta \log(p_{i,d}^\beta) \right) \quad (5)$$

where Q is the pre-defined maximum number of phrases, and R is the pre-defined maximum token number in phrase. Note that we only predict the order of the tokens and phrases in query.

3.5 Loss Function

In order to enable QUERT to integrate the above capabilities, we use joint training to combine four tasks. We consider each task to be equally important, so no additional weight is assigned. And the overall training loss function is defined as:

$$\mathcal{L} = \mathcal{L}_{Geo-MP} + \mathcal{L}_{Geo-CP} + \mathcal{L}_{UCBL} + \mathcal{L}_{TPOP}, \quad (6)$$

The pre-training target is to minimize the \mathcal{L} .

4 EXPERIMENTS

In this section, we demonstrate the effectiveness of our model through extensive offline experiments. Then, we deploy the model on Fliggy APP to test its performance in real online scenarios.

4.1 Experimental Setup

4.1.1 Pre-training corpora. We collect data from the past three years of Fliggy’s real online business scenarios. The data consists of user search query q , click items c , and their confidence scores (i.e., unique visitor click hits UV_C and unique visitor payment hits UV_P). Since the original data from online business contains numerous noise, we rank the data according to the weighted score $S = UV_C + 10 \times UV_P$. To ensure the quality of the data, we select the top 5 million highest weighted score pairs $P = (q, c)$ and discard others.

4.1.2 Downstream Tasks. The downstream tasks are described as follows.

Query Rewriting (QR). QR aims to reduce the expression gap between user queries and items. In other words, QR rewrites an unparsed query into a parsed query which is stored in the database. To build a high confidence dataset, first, we collect unparsed user input queries and the corresponding click items from search logs. Then we select the parsed query linked to the same click item as the candidate samples. Finally, to ensure the quality of the dataset, we invite people to annotate the candidate samples and filter those that are marked as untrustworthy.

For evaluation, we select 181,405 standard queries as candidates. We adopt accuracy (Acc) and hit rate (Hits@K) as evaluation metrics. Hits@K means the success rate of finding the ground truth in top K candidates. In line with the actual business, we set K to 20 in our experiments. In order to retrieve, we feed the query into an encoder (e.g., BERT, RoBERTa, and QUERT) and gain its embedding. We calculate the cosine similarity between the tested query and all candidates. The top K highest cosine similarity sample would be used for evaluation. In the supervised setting, we use contrastive learning to fine-tune the similarity. Specifically, the ground truth is regarded as a positive example, and negative examples are randomly sampled in the training batch.

Query-POI Retrieval (QPR). Given a query, QPR aims to provide several POI recommendations to improve users’ search experience.

Table 1: Statistics of downstream tasks dataset. “#” indicates the number of samples.

TASK	#			Metric
	Train	Dev	Test	
Query Rewriting	20,000	2,500	2,500	Acc, Hits@20
Query-POI Retrieval	26,854	3,836	2,947	Acc, Hits@20
Query Intention Classification	54,003	6,712	6,709	P, R, F1
Query Destination Matching	67,241	7,982	8,443	P, R, F1
Query Error Detection	48,333	6,890	6,794	P, R, F1

In order to construct the dataset, for every query, we select the top 20 click rate POIs as ground truths. All data is manually verified again.

For evaluation, we select 201,184 POIs as candidates for retrieval. Accuracy (Acc) and hit rate (Hits@K) are adopted as evaluation metrics. And K is set to 20. The implementation detail of training is the same as QR.

Query Intention Classification (QIC). For precise recommendation of relevant items, QIC aims to predict the user’s intention based on the input query. Based on business practices, we pre-define 20 categories of intentions. For dataset building, first, we extract the item category with the highest click rate corresponding to the query as the intention pseudo-label. For these data, we perform additional human checks and construct the dataset.

We regard this task as a sentence multi-class classification task. In terms of evaluation metrics, we adopt the widely used precision (P), recall (R) and F1.

Query Destination Matching (QDM). Given a query and candidate city, QDM aims to determine whether the city is the intended destination of the query. This task adds constraints to the search recall module, avoiding irrelevant recall results, thereby improving the accuracy and relevance of the entire search engine.

Similar to the intention classification task, we select the location of the click item with the top K highest click rate as the destination label. We invite people to do a further evaluation of the data and build the dataset. Given a pair consisting of a query and a city, the model performs binary classification of the pair. We use precision, recall and F1 to evaluate the quality.

Query Error Detection (QED). QED identifies specific token errors within queries, such as typos or incorrect order, rather than classifying the entire query as wrong. Improved error detection through QED enables more accurate search results. We collect wrong queries from real data, including redundant tokens, transpositions of token order, and typos, to build the final dataset. There are four kinds of labels in total. Labels “0” to “3” signify no error, a typo, token transposition, and redundant token respectively. We regard QED as a sequence-labeling task, so the model is expected to predict the specific error for each token. Precision, recall and F1 of token-level are used as evaluation metrics.

The detailed data scale and evaluation metrics of five downstream tasks are presented in Table 1.

4.1.3 Baselines. We select several pre-trained models widely used in Chinese NLP tasks as the baselines for comparison with QUERT.

- BERT[5] is a strong PLM based on Transformer architecture. We use the Chinese version released by Google.

Table 2: Supervised performance comparison of different setting models on downstream tasks. All results are based on the average of five replicates. BERT_q is BERT continually pre-trained on only query corpora, and BERT_{q+c} is continually pre-trained on query and item. The subscript of QUERT indicates the backbone model. The best results are marked in bold and the second best are underlined.

Models	QR		QPR		QIC			QDM			QED			Average
	Acc	Hits@20	Acc	Hits@20	P	R	F1	P	R	F1	P	R	F1	
BERT [5]	47.94	92.53	59.71	79.97	41.70	35.08	35.56	90.18	88.36	89.20	89.55	87.39	88.35	71.19
RoBERTa [17]	47.93	92.72	59.43	80.37	42.07	33.55	34.56	90.21	88.13	89.07	89.96	86.94	88.27	71.02
ERNIE [24]	47.52	92.81	60.54	82.23	41.84	35.54	36.06	89.22	88.78	88.98	89.57	87.22	88.27	71.43
StructBERT [27]	47.51	92.35	60.00	81.34	41.90	35.44	35.87	90.00	88.90	89.42	89.79	85.75	87.63	71.22
Mengzi [32]	47.87	92.05	59.60	79.57	40.98	34.46	34.60	88.89	87.08	87.92	89.95	<u>87.48</u>	88.57	70.69
BERT _q	47.80	92.96	59.76	79.68	41.09	35.00	35.54	89.55	86.97	88.10	89.09	86.19	87.42	70.70
BERT _{q+c}	48.47	<u>93.04</u>	59.94	80.68	<u>45.67</u>	35.90	36.93	89.24	<u>90.11</u>	89.63	90.17	87.19	88.53	71.96
QUERT _{BERT}	<u>48.66</u>	93.08	<u>61.26</u>	82.49	<u>45.44</u>	39.20	39.58	91.78	<u>90.11</u>	90.88	<u>91.25</u>	88.30	89.64	73.21
QUERT _{ERNIE}	<u>48.67</u>	92.99	<u>60.98</u>	<u>82.25</u>	<u>43.24</u>	<u>38.41</u>	<u>38.30</u>	<u>90.77</u>	90.96	<u>90.87</u>	91.86	87.00	<u>89.20</u>	<u>72.73</u>

Table 3: Unsupervised performance comparison of different setting models on downstream tasks.

Models	QR		QPR		Average
	Acc	Hits@20	Acc	Hits@20	
BERT	17.89	40.59	21.95	33.70	28.53
ERNIE	14.49	32.27	17.48	25.45	22.42
QUERT _{BERT}	41.75	85.79	<u>44.15</u>	<u>63.76</u>	59.46
QUERT _{ERNIE}	<u>39.79</u>	<u>85.03</u>	44.32	64.54	<u>58.42</u>

- **RoBERTa**[17] is a PLM whose architecture is the same as BERT. We use the Chinese version RoBERTa-wwm [4].
- **ERNIE**[24] is also a Transformer-based PLM. In addition to the token-level mask, it introduces the entity-level and phrase-level masking strategies.
- **StructBERT**[27] is the variant of BERT. It adds a new word structural objective in the training stage to force the model to reconstruct the right order of sequence, which is similar to our PTOp task.
- **Mengzi**[32] is also the variant of BERT, which is dedicated to being lightweight but powerful. Mengzi gains outstanding performance on multiple Chinese NLP tasks.

4.1.4 Implementation Details. Our implementation is based on Transformers framework [28] and Pytorch [23]. To verify the effectiveness of our proposed pre-training task on different model architectures, we apply BERT [5] and ERNIE [31] as backbone models to implement continual pre-training. We adopt AdamW [18] as the optimizer and set the initial learning rate to 5e-5 for pre-training and 1e-5 for fine-tuning. With steps increasing, we decrease the learning rate linearly. The size of geohash code is 6 bits. And we adopt Base32 as the numeral system in geohash coding. The temperature τ of contrastive learning is set to 0.1. And the shuffling probability in TPOP is 0.15. The training takes about 72 hours for 386460 steps on 8 Tesla V100 GPUs with 16 batch size per device.

4.2 Offline Results

4.2.1 Supervised results. We compare QUERT with baseline models on five downstream tasks in the supervised setting (the model is fine-tuned on train set). Table 2 shows the results.

First, **QUERT achieves SOTA results in all downstream tasks**, which demonstrates the effectiveness of our proposed pre-training tasks. Specifically, compared with BERT, **QUERT_{BERT}** improves the average performance by 2.02%. Furthermore, we find that the advantage of QUERT is more pronounced on difficult tasks. In detail, QUERT has a huge performance advantage on QIC, which has an average improvement is nearly 4%.

Second, **the corpora composed of single query information bring negative effects**. We test directly continual pre-training on raw query corpora with the masking strategy used in BERT [5], which is presented as BERT_q in Table 2. Experimental results show that BERT_q brings the risk of negative effects. Specifically, BERT_q achieves an average score of 70.90% which is even lower than the original BERT. We analyze that the model may not be able to learn knowledge representation from the short query but weaken the text understanding ability. The results demonstrate that the regular pre-training task cannot be directly generalized to the pre-trained model focusing on the query. This verifies that our proposed pre-training tasks are more applicable to travel domain.

Third, **the integration of click items information improves the representation ability of the model**. We concatenate the title text information of items into the query to construct the new corpora. Compared to BERT_q, the performance of BERT_{q+c} improves nearly 1.3%. This verifies our conjecture that in the pre-training stage, the model is informed of the items' information with high confidence, which is helpful for the model to acquire more knowledge of the query.

4.2.2 Unsupervised results. We compare QUERT with two backbone models in the unsupervised setting (without fine-tuning). Table 3 reports the unsupervised results. We select QR and QPR as unsupervised test tasks because they can directly obtain results by calculating the embedding similarity. We obtain the prediction by calculating the cosine similarity between the embeddings that are

Table 4: Ablation results of QUERT_{BERT}. We repeat five times and report the average scores.

Models	QR		QPR		QIC		QDM			QED		Average		
	Acc	Hits@20	Acc	Hits@20	P	R	F1	P	R	F1	P		R	F1
QUERT _{BERT}	48.66	93.08	61.26	82.49	45.44	39.20	39.58	91.78	90.11	90.88	91.25	88.30	89.64	73.21
- w/o Geo-MP	48.05	93.05	60.81	81.94	45.89	36.93	38.44	90.82	90.58	90.68	90.91	87.16	88.84	72.62
- w/o Geo-CP	47.61	92.31	61.32	82.35	43.02	37.00	37.35	89.65	91.08	90.32	90.57	87.37	88.80	72.21
- w/o UCBL	48.08	92.87	60.70	82.44	44.00	36.50	37.01	90.04	88.72	89.33	90.66	86.57	88.39	71.95
- w/o PTOp	48.61	92.95	61.42	82.22	44.95	37.23	38.23	90.92	90.77	90.84	88.73	85.60	86.99	72.25

gained from the last layer of hidden states. “[CLS]” token is used for representing the whole query.

From the table, we observe that QUERT has significant performance advantages in unsupervised setting. Both QUERT_{BERT} and QUERT_{ERNIE} significantly outperform the backbone model. For QR, our proposed QUERT_{BERT} outperforms the baseline BERT by 45.20% on Hits@20. As for QPR, QUERT_{BERT} and QUERT_{ERNIE} significantly outperform the two baselines, with a maximum performance gap of 39.09% on Hits@20. In terms of average scores, QUERT_{BERT} achieves a 30.93% performance improvement over BERT. We analyze that the reason for the large performance gap is that QUERT is more potent in query understanding for travel domain search. In unsupervised setting, the tailored pre-training tasks empower QUERT to give better query representation.

4.2.3 Ablation Studies. To check whether the customized pre-training tasks can effectively improve the performance of downstream tasks, we perform a series of ablation experiments. Specifically, we remove tasks one at a time to evaluate the impact of their presence or absence on performance. Table 4 reports the results.

First, the removal of any tasks leads to performance loss. The elimination of every component, i.e., Geo-MP, Geo-CP, UCBL, and PTOp, causes the 0.59%, 1.00%, 1.26%, and 0.96% drop in the score.

Second, we find that the removal of Geo-MP (i.e., degenerate to the original masking strategy) and Geo-CP results in performance degradation on all tasks. This reveals that these two geography awareness pre-training tasks make an effective contribution to improving the ability of pre-trained model to perceive geography.

Third, UCBL plays the most important role in sentence-level tasks (i.e., QIC and QDM). On the one hand, the construction of similarities in user behavior enables QUERT to understand queries better. On the other hand, in the study of negative examples, the differentiation of sentence-level embedding representation is enlarged, which is beneficial to model recognition in the prediction.

Finally, the removal of PTOp leads to the F1 score on QED (86.99%) being lower than the original BERT (88.35%). We guess that other tasks may introduce additional logical bias, which results in a performance penalty for the model on QED. However, the introduction of PTOp allows QUERT to refactor its understanding of the logical consistency of query, resulting in significant performance improvements. In addition, we find that StructBERT, which aims at reconstructing order in MLM, did not achieve superior performance in QED. This verifies that our proposed PTOp is more suitable for the real downstream task of travel domain search.

In conclusion, the results verify that our designed pre-training task in a tailored way endows the pre-trained model with powerful query representation capabilities in travel domain search.

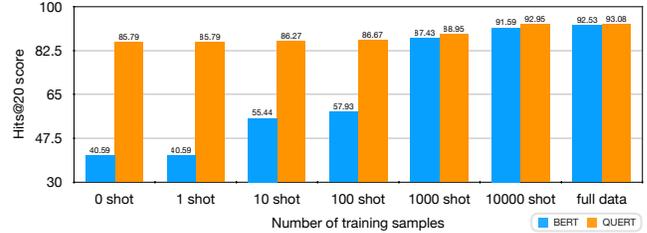


Figure 3: The comparison of few-shot learning on QR between BERT and QUERT_{BERT}. We use Hits@20 as the evaluation metric.

4.2.4 Few-shot Learning. In order to test the performance of QUERT in the data scarcity scenario, we conduct the few-shot learning experiment. We choose QR as the experimental object because it is the most dependent on annotated data in real business, which can be verified in Table 3 (only 17.89% Acc score in zero-shot). Figure 3 reports the results. We conclude that: 1) QUERT can achieve better results than BERT under the same training data size. 2) When the amount of labeled data is no more than 100, the performance gap between QUERT and BERT is noticeable, ranging from 28.74% to 45.20%. 3) BERT does not surpass the zero-shot performance of QUERT until the training sample number reaches 1000. The experimental results show that, compared with BERT, QUERT is competent for downstream tasks under sparse data scenarios, even under the extreme conditions of zero-shot and one-shot.

4.3 Task Study

In this section, we continually pre-train BERT with a single task to verify the true effect of each task.

4.3.1 TASK 1: Geo-MP. In order to verify the effectiveness of Geo-MP, we evaluate the sensitivity to the geography of BERT and BERT+Geo-MP. We randomly collect 500 queries from Fliggy APP. First, for geography token in query, we use a special token “[mask]” to hide it. Then we let the model predict these masked tokens. We adopt the hit rate score (denoted as Hits@K) and mean rank (i.e., MR and MRR) as the metric. As shown in Table 5, the best score of BERT is only 2.2% on Hits@5. In contrast, BERT+Geo-MP gains 13.2% on Hits@5, lower MR and higher MR.

In order to further analyze the advantages of Geo-MP in location prediction, we present the cases in Table 6. In case 1 and 2, BERT+Geo-MP predict the right masked token, but BERT output the nonstandard city name or phrase. Besides, in case 3, although BERT+Geo-MP gives a wrong geography phrase, it still outputs a

Table 5: Geography sensitivity comparison for BERT and BERT+Geo-MP. Hits@K means the success rate of finding the ground truth in K candidates. MR is the mean rank. MRR is the mean reciprocal rank.

	Hits@1	Hits@3	Hits@5	MR	MRR
BERT	1.2%	1.8%	2.2%	1.91	0.71
BERT+Geo-MP	9.8%	12.8%	13.2%	1.41	0.85

Table 6: Geography mask cases. “*” is the placeholder of “MASK” token. “X” indicates the prediction result is not a phrase in Chinese.

No	Text	BERT+Geo-MP	BERT	Answer
1	[**]天涯海角景区	三亚	北州	三亚
	[**]Tianya Haijiao	Sanya	X	Sanya
2	北京[**]长城	八达岭	的长山	八达岭
	The Great Wall of [**] Section, Beijing	Ba da ling	X	Ba da ling
3	杭州[**]门票	灵隐寺	天物园	岳王庙
	The ticket of [**], Hangzhou	Ling yin Temple	X	Yue Fei Temple
4	扬州[**]二十四桥	瘦西湖	市桥第	瘦西湖
	Yangzhou [**] Twenty-four Bridges	Slender West Street	X	Slender West Lake

reasonable prediction. In fact, “Ling yin Temple” is exactly a POI located in Hangzhou, and it makes sense in the given context. As another bad case, case 4 shows that the output of BERT+Geo-MP is closer to the answer.

The results verify that classical MLMs are indifferent to geography information. By contrast, Geo-MP enables BERT to be more sensitive to geography.

4.3.2 TASK 2: Geo-CP. To verify the physical geography representation ability of Geo-CP, we select 500 hot POIs in 10 popular cities and visualize their embeddings through t-SNE[25]. We use the embedding of “[CLS]” token to represent the POI.

As shown in Figure 4, the embeddings produced by BERT are distributed out of order in the space. This proves that BERT does not reflect real location information of the geography query. However, it can be observed that the embedding space of BERT+Geo-CP is more orderly, and the POIs in the same city are in the same cluster. In addition, we also notice that the space presented by Geo-CP does have a real physical geographical location relationship. For example, in the real world, Shanghai, Hangzhou and Nanjing are close to each other, and in the figure, the relationship among them is indeed the same. We analyze that through Geo-CP, the language model QUERT is endowed with geographical location representation ability. This demonstrates the effectiveness of our proposed pre-training task.

4.3.3 TASK 3: UCBL. To verify the effectiveness of UCBL, we compare the similarity between the click behavior related queries in the BERT and BERT+UCBL settings. First, we select 500 queries q and their click behavior related queries q^{pos} from search logs. With BERT and BERT+UCBL as encoders, we feed q and q^{pos} into the encoders and gain their embeddings. For convenience, we define the embeddings of q and q^{pos} produced by BERT as R_B^q and R_B^{pos} , and these gained from BERT+UCBL as R_{B+U}^q and R_{B+U}^{pos} . Then we calculate cosine similarity of (R_B^q, R_B^{pos}) and $(R_{B+U}^q, R_{B+U}^{pos})$. The results are 0.7758 and 0.8278, respectively, which proves that BERT+UCBL can perceive the potential similarities of user behavior in query.

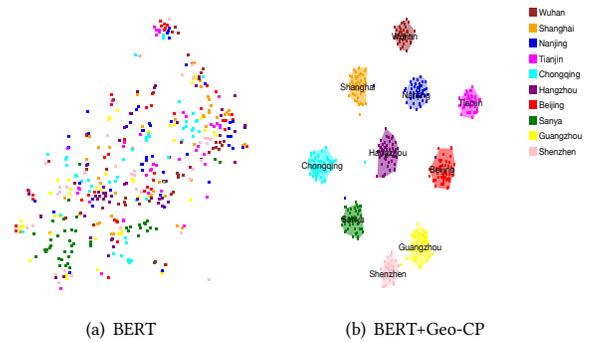


Figure 4: The t-SNE space comparison between BERT and BERT+Geo-CP.

Table 7: The performance comparison of the order prediction subtask in QED.

	P	R	F1
BERT	84.77	77.88	81.18
StructBERT	81.49	84.33	82.89
BERT+PTOP	83.48	84.39	84.20

4.3.4 TASK 4: PTOF. To verify the effectiveness of PTOF, we analyze the performance of different models on the order prediction subtask in QED. To be specific, this subtask aims to judge whether the token order is transposed. The result is reported in Table 7.

Compared to BERT, StructBERT gains better performance because of its pre-training task which aims to reconstruct the token order. However, BERT+PTOP has a performance advantage over StructBERT. According to our analysis, PTOF directly focuses on order judgment, which is more matching with the downstream task of travel domain search. The performance advantage of PTOF demonstrates the effectiveness of our proposed pre-training task.

4.4 Case Study

According to the different characteristics of queries targeted by different tasks, we conduct the case study in this section. As shown in Figure 5, we compare the representation ability of BERT and QUERT by evaluating the cosine similarity of embedding.

Geography Awareness. We evaluate the sensitivity of QUERT to geography information (i.e., POI and City). In Figure 5 (a), we calculate the cosine similarity between different POI in the same city. For example, “Happy Valley” and “Disney” are both names of parks, and BERT gives them a high degree of similarity. However, QUERT recognizes that they are two completely different POIs and distinguishes them. Similarly, we evaluate the same POI in different cities. Although the two queries both contain the same POI “West Lake”, QUERT senses that the exact essential part is the city and gives a low similarity score. These results indicate that QUERT has an awareness of geography information, thus differentiating queries for different locations.

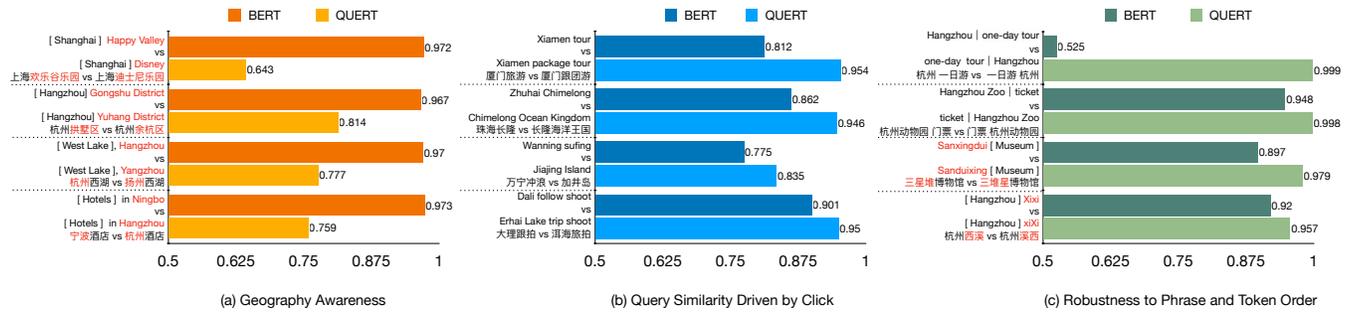


Figure 5: The comparison of cosine similarity. We compare the cosine similarity of query embeddings obtained from BERT and QUERT_{BERT}, respectively. “[CLS]” token is used to represent the whole query.

Query Similarity Driven by Click. We also evaluate the understanding of click behavior similarity in Figure 5 (b). For the examples “Xiamen tour” and “Xiamen package tour”, BERT assign low similarity for them because of literal difference. However, QUERT understands the potential user behavior similarity for the same intention “Searching information about the tour of Xiamen” and assigns high similarity. Such a phenomenon confirms that QUERT is able to understand the similarity driven by user click behavior.

Robustness to Phrase and Token Order. In order to test whether QUERT is robust to phrase and token order, we choose four permutation queries as experimental objectives in Figure 5 (c). For the example of “Hangzhou one-day tour” and “one-day tour Hangzhou”, QUERT has a higher tolerance for phrase permutation. As for token permutation, QUERT identifies “Sanxingdui” as the correct form of “Sanduixing” and assigns high similarity. These cases show that QUERT is truly robust to phrase and token order.

4.5 Online Application

We perform online A/B testing on Fliggy APP to validate QUERT’s capabilities in the real business scenario. To be specific, given an unparsed query (e.g., misinput or emerging query), we gain its embedding \mathcal{R} by feeding it to an encoder. Then we calculate the cosine similarity between \mathcal{R} and other embeddings of parsed queries in the database. We select the top 20 queries with the highest similarity scores as similar queries and get the search items results corresponding to these similar queries. And these results are considered as the potential recommendation for the unparsed query.

In our online A/B testing, we compare two buckets, each containing 10% randomly-selected users. We use BERT as the encoder for one bucket, and for the other, we adopt QUERT. To measure the actual change in online business, we use two well-known metrics: **Unique Click-Through Rate (U-CTR)** and **Page Click-Through Rate (P-CTR)**. After running 7 days, the feedback results show that U-CTR and P-CTR increase by 0.89% and 1.03%, respectively. This suggests that QUERT’s unparsed query representations retrieve more relevant similar queries, signifying its aptness for embedding-based retrieval in travel domain search.

In Figure 6, we provide two cases. The query “Xia Shangri-La” is unparsed because it misses a token “men”. And the real intention of this query is “Searching for the hotel named Shangri-La located in Xiamen”. As shown in Figure 6 (a), with BERT, the system mistakenly recalls items related to “Yunnan” because the token

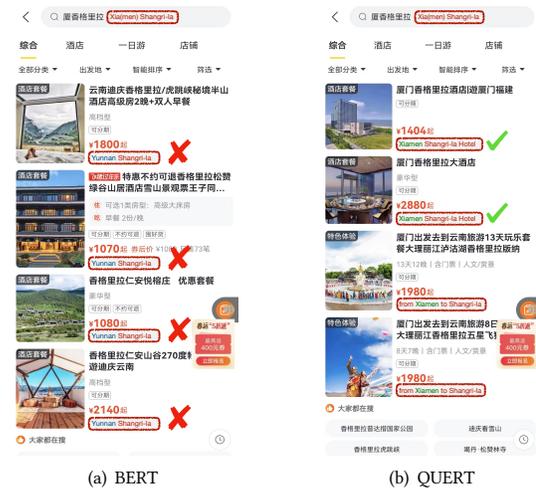


Figure 6: Comparison of online application between BERT and QUERT.

“Shangri-La” also refers to a city located in Yunnan. Due to missing “men”, BERT fails to capture the potential constraint “Xiamen”. However, with QUERT, the system recalls the right items mentioning “Shangri-La Hotel in Xiamen” and high-related items “Trip from Xiamen to Shangri-La”, which proves the effectiveness of QUERT.

5 CONCLUSION AND FUTURE WORK

In this paper, we focus on the continual pre-training for query understanding in travel domain search. We analyze the causes of query representation difficulties and propose a solution: QUERT, a continual pre-trained language model. To be specific, we propose four tailored pre-training tasks: Geography-aware Mask Prediction, Geohash Code Prediction, User Click Behavior Learning, and Phrase and Token Order Prediction. We evaluate offline performance on five downstream tasks in the travel domain. Experimental results show that compared to BERT, the performance of QUERT on downstream tasks improves by 2.02% and 30.93% in supervised and unsupervised settings, respectively. Furthermore, the online A/B testing on Fliggy APP demonstrates that U-CTR and P-CTR increase by 0.89% and 1.03% when applying the QUERT as the feature encoder.

As a language model, QUERT only relies on text information. In future work, we plan to introduce more information (e.g., images) and explore the pre-trained multimodal in travel domain search.

6 ACKNOWLEDGEMENTS

This work was supported by Alibaba Group through Alibaba Innovative Research Program, Shanghai Municipal Science and Technology Major Project (No.2021SHZDZX0103), Science and Technology Commission of Shanghai Municipality Grant (No. 22511105902), and Shanghai Sailing Program (No. 23YF1409400).

REFERENCES

- [1] Dogu Araci. 2019. Finbert: Financial sentiment analysis with pre-trained language models. *arXiv preprint arXiv:1908.10063* (2019).
- [2] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language Models are Few-Shot Learners. *arXiv preprint arXiv:2005.14165* (2020).
- [3] Brooke Cowan, Sven Zethelius, Brittany Luk, Teodora Baras, Prachi Ukarde, and Daodao Zhang. 2015. Named entity recognition in travel-related search queries. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 29, 3935–3941.
- [4] Yiming Cui, Wanxiang Che, Ting Liu, Bing Qin, and Ziqing Yang. 2021. Pre-training with whole word masking for chinese bert. *IEEE/ACM Transactions on Audio, Speech, and Language Processing* 29 (2021), 3504–3514.
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [6] Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. *Advances in Neural Information Processing Systems* 32 (2019).
- [7] Siddhant Garg and Goutham Ramakrishnan. 2020. BAE: BERT-based Adversarial Examples for Text Classification. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 6174–6181.
- [8] Zheng Gong, Kun Zhou, Wayne Xin Zhao, Jing Sha, Shijin Wang, and Ji-Rong Wen. 2022. Continual Pre-training of Language Models for Math Problem Understanding with Syntax-Aware Memory Network. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. 5923–5933.
- [9] Suchin Gururangan, Ana Marasović, Swabha Swayamdipta, Kyle Lo, Iz Beltagy, Doug Downey, and Noah A Smith. 2020. Don't Stop Pretraining: Adapt Language Models to Domains and Tasks. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. 8342–8360.
- [10] Jizhou Huang, Haifeng Wang, Yibo Sun, Yunsheng Shi, Zhengjie Huang, An Zhuo, and Shikun Feng. 2022. ERNIE-GeoL: A Geography-and-Language Pre-trained Model and its Applications in Baidu Maps. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 3029–3039.
- [11] Mandar Joshi, Danqi Chen, Yinhan Liu, Daniel S Weld, Luke Zettlemoyer, and Omer Levy. 2020. Spanbert: Improving pre-training by representing and predicting spans. *Transactions of the Association for Computational Linguistics* 8 (2020), 64–77.
- [12] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2019. ALBERT: A Lite BERT for Self-supervised Learning of Language Representations. In *International Conference on Learning Representations*.
- [13] Jinhyuk Lee, Wonjin Yoon, Sungdong Kim, Donghyeon Kim, Sunkyu Kim, Chan Ho So, and Jaewoo Kang. 2020. BioBERT: a pre-trained biomedical language representation model for biomedical text mining. *Bioinformatics* 36, 4 (2020), 1234–1240.
- [14] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising Sequence-to-Sequence Pre-training for Natural Language Generation, Translation, and Comprehension. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Online, 7871–7880. <https://doi.org/10.18653/v1/2020.acl-main.703>
- [15] Chen Liang, Yue Yu, Haoming Jiang, Siawpeng Er, Ruijia Wang, Tuo Zhao, and Chao Zhang. 2020. Bond: Bert-assisted open-domain named entity recognition with distant supervision. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 1054–1064.
- [16] Xiao Liu, Juan Hu, Qi Shen, and Huan Chen. 2021. Geo-BERT Pre-training Model for Query Rewriting in POI Search. In *Findings of the Association for Computational Linguistics: EMNLP 2021*. 2209–2214.
- [17] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692* (2019).
- [18] Ilya Loshchilov and Frank Hutter. 2018. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*.
- [19] Xinyu Ma, Jiafeng Guo, Ruqing Zhang, Yixing Fan, Xiang Ji, and Xueqi Cheng. 2021. PROP: pre-training with representative words prediction for ad-hoc retrieval. In *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*. 283–291.
- [20] Xinyu Ma, Jiafeng Guo, Ruqing Zhang, Yixing Fan, Yingyan Li, and Xueqi Cheng. 2021. B-PROP: bootstrapped pre-training with representative words prediction for ad-hoc retrieval. In *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 1513–1522.
- [21] Martin Müller, Marcel Salathé, and Per E Kummervold. 2020. Covid-twitter-bert: A natural language processing model to analyse covid-19 content on twitter. *arXiv preprint arXiv:2005.07503* (2020).
- [22] Dat Quoc Nguyen, Thanh Vu, and Anh Tuan Nguyen. 2020. BERTweet: A pre-trained language model for English Tweets. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Online, 9–14. <https://doi.org/10.18653/v1/2020.emnlp-demos.2>
- [23] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. <https://proceedings.neurips.cc/paper/2019/file/bdbca288fee7f92f2bfa9f7012727740-Paper.pdf>
- [24] Yu Sun, Shuohuan Wang, Yukun Li, Shikun Feng, Xuyi Chen, Han Zhang, Xin Tian, Danxiang Zhu, Hao Tian, and Hua Wu. 2019. ERNIE: Enhanced Representation through Knowledge Integration. *arXiv e-prints* (2019), arXiv–1904.
- [25] Laurens Van der Maaten and Geoffrey Hinton. 2008. Visualizing data using t-SNE. *Journal of machine learning research* 9, 11 (2008).
- [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).
- [27] Wei Wang, Bin Bi, Ming Yan, Chen Wu, Jiangnan Xia, Zuyi Bao, Liwei Peng, and Luo Si. 2019. StructBERT: Incorporating Language Structures into Pre-training for Deep Language Understanding. In *International Conference on Learning Representations*.
- [28] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Online, 38–45. <https://doi.org/10.18653/v1/2020.emnlp-demos.6>
- [29] Jia Xu, Fei Xiong, Zulong Chen, Mingyuan Tao, Liangyue Li, and Quan Lu. 2022. G2NET: A General Geography-Aware Representation Network for Hotel Search Ranking. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 4237–4247.
- [30] Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019. Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems* 32 (2019).
- [31] Zhengyan Zhang, Xu Han, Zhiyuan Liu, Xin Jiang, Maosong Sun, and Qun Liu. 2019. ERNIE: Enhanced Language Representation with Informative Entities. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. 1441–1451.
- [32] Zhuosheng Zhang, Hanqing Zhang, Keming Chen, Yuhang Guo, Jingyun Hua, Yulong Wang, and Ming Zhou. 2021. Mengzi: Towards lightweight yet ingenious pre-trained models for chinese. *arXiv preprint arXiv:2110.06696* (2021).