## **TextGNN: Improving Text Encoder via Graph Neural Network in Sponsored Search**

Jason Yue Zhu\*<sup>†</sup> Stanford University Stanford, CA, USA jzhu121@stanford.edu

Hao Sun Microsoft Beijing, China hasun@microsoft.com

Tiangi Yang Microsoft Beijing, China tianqi.yang@microsoft.com

Yanling Cui<sup>†</sup> Microsoft Beijing, China yanling.cui@microsoft.com

> Xue Li Microsoft Sunnyvale, CA, USA xeli@microsoft.com

Liangjie Zhang Microsoft Beijing, China liazha@microsoft.com

Huasha Zhao Microsoft Sunnyvale, CA, USA huasha.zhao@microsoft.com

Yuming Liu Microsoft Beijing, China yumliu@microsoft.com

Markus Pelger Stanford University Stanford, CA, USA mpelger@stanford.edu

#### Ruofei Zhang

Microsoft Sunnyvale, CA, USA bzhang@microsoft.com

### ABSTRACT

Text encoders based on C-DSSM or transformers have demonstrated strong performance in many Natural Language Processing (NLP) tasks. Low latency variants of these models have also been developed in recent years in order to apply them in the field of sponsored search which has strict computational constraints. However these models are not the panacea to solve all the Natural Language Understanding (NLU) challenges as the pure semantic information in the data is not sufficient to fully identify the user intents. We propose the TextGNN model that naturally extends the strong twin tower structured encoders with the complementary graph information from user historical behaviors, which serves as a natural guide to help us better understand the intents and hence generate better language representations. The model inherits all the benefits of twin tower models such as C-DSSM and TwinBERT so that it can still be used in the low latency environment while achieving a significant performance gain than the strong encoder-only counterpart baseline models in both offline evaluations and online production system. In offline experiments, the model achieves a 0.14% overall increase in ROC-AUC with a 1% increased accuracy for long-tail low-frequency Ads, and in the online A/B testing, the model shows

\*This work was completed during the 1st author's internship at Microsoft <sup>†</sup>Authors contributed equally to this work

© 2021 IW3C2 (International World Wide Web Conference Committee), published under Creative Commons CC-BY 4.0 License.

ACM ISBN 978-1-4503-8312-7/21/04.

https://doi.org/10.1145/3442381.3449842

a 2.03% increase in Revenue Per Mille with a 2.32% decrease in Ad defect rate.

### **CCS CONCEPTS**

 Information systems → Recommender systems; Language models; Similarity measures; Learning to rank; Query representation.

#### **KEYWORDS**

Ad Relevance; Sponsored Search; Text Encoder; Graph Neural Network; Transformers; C-DSSM; BERT; Knowledge Distillation

#### ACM Reference Format:

Jason Yue Zhu, Yanling Cui, Yuming Liu, Hao Sun, Xue Li, Markus Pelger, Tianqi Yang, Liangjie Zhang, Ruofei Zhang, and Huasha Zhao. 2021. TextGNN: Improving Text Encoder via Graph Neural Network in Sponsored Search. In Proceedings of the Web Conference 2021 (WWW '21), April 19-23, 2021, Ljubljana, Slovenia. ACM, New York, NY, USA, 10 pages. https: //doi.org/10.1145/3442381.3449842

#### 1 **INTRODUCTION**

Sponsored search refers to the business model of search engine platforms where third-party sponsored information is shown to targeted users along with other organic search results. This allows the advertisers such as manufacturers or retailers to increase the exposure of their products to more targeted potential buyers, and at the same time gives users a quicker access to solutions for their needs. Hence it has become an indispensable part of our modern web experience. While many of the existing models are very powerful for various tasks in sponsored search, there still remain three main challenges for future developments in this field: 1) while

This paper is published under the Creative Commons Attribution 4.0 International (CC-BY 4.0) license. Authors reserve their rights to disseminate the work on their personal and corporate Web sites with the appropriate attribution. WWW '21, April 19-23, 2021, Ljubljana, Slovenia

the existing models have strong performances on matching common queries with popular products, they usually still find long-tail low-frequency queries/Ads to be more challenging. The worse embedding representations in rare items are potentially caused by under-training due to naturally scarce data on these low-frequency examples. 2) while many modern models improve in implicit feature engineering on the existing input data, finding new and easily accessible data with complement information is still a promising route to greatly improve the model performance but is rarely explored. 3) the search engine systems generally have very strict constraints on computational resources and latency requirements. Many recently developed large powerful models are simply infeasible to deploy onto the highly constrained online search engine systems.

Representation learning for queries, products, or users has been a key research field with many breakthroughs over the last years and has been adopted in many production sponsored search systems [8][18][4][1]. Convolutional Deep Structured Semantic Model (C-DSSM) [21] is among the first powerful solutions to encode text data into low-dimensional representation vectors which can be applied to downstream tasks and have efficient inference performance, but its NLU performance has been surpassed by many recently developed NLP models. The pre-trained language models emerged in recent years, such as transformers [22] and BERT [3], have demonstrated far superior performance in many NLU tasks and even reach human level performance on many tasks. These models are better at capturing contextual information in the sentences and generate better language representation embeddings, leading to much stronger performance in downstream tasks. However, due to the complexity, these models are unfortunately not feasible to run in low latency systems without modifications. Recently, the transformer model has been modified and trained with special techniques such as knowledge distillation [7], which allows us to use similar transformers structure but much smaller model called TwinBERT [17] to run with reasonable computational cost in the production systems while having little or no performance loss compared to the full size BERT models. This breakthrough significantly improves the user Information Retrieval experience when using search engines. However, while both C-DSSM and TwinBERT are specifically designed to be applied to the low latency systems with strong performance, they are not the panacea to fully solve all the problems in sponsored search. Their model ability is sometimes hindered by the limited information in the original input texts and hence still suffers in understanding many challenging low frequency inputs.

Given the strong performance of the baseline models in NLU tasks, it would be extremely difficult to further improve them solely based on the structural changes of the model without introducing new complement information. The newly developed NLP models achieve relatively small improvements with exponentially growth in model complexity, and hence reach the margin of diminishing returns making it harder to satisfy all the latency constraints. A real improved model in this field should then be able to take in additional information beyond the tradition semantic text inputs, demonstrate stronger performance over the harder low-frequency inputs, and at the same time should not significantly increase the inference time.

A natural and easily accessible data source that provides information beyond semantic text in the search engine system is users' implicit feedbacks recorded in logs in the form of clicks through the links shown to them. A click signals a connection between a query and an Ad and hence a large behavior graph based on clicks can be easily built. In the recent years, various Graph Neural Network (GNN) structures [27] have been proposed to deal with the abundant graph-typed data and demonstrated strong performance and breakthroughs in social networks, recommendations, or natural science tasks. Motivated by the recent developments in GNN community, we are aiming to identify ways to include complementary and abundant graph-type data into the text model in a natural way. Most existing GNN models focus only on the aggregation of pre-existing neighbor features that are fixed throughout training. Instead of training the language model and the graph model separately, we want the two models to work in conjunction with each other to generate better query/Ad representations that can help understanding users' needs in a deeper way.

The main contributions of this work are three-folds:

- (1) We propose TextGNN<sup>1</sup>, a general end-to-end framework for NLU that combines the strong language modeling text encoders with graph information processed by Graph Neural Networks to achieve stronger performance than each of its individual components.
- (2) We find a systematical way to leverage graph information that greatly improves the robustness and performance by 1% on hard examples. These samples are very challenging when only using semantic information.
- (3) We trained TextGNN with knowledge distillation to get a compact model. The model has been adopted in the production system that has strict computational and latency constraints while achieving a 2.03% increase in Revenue Per Mille with a 2.32% decrease in Ad defect rate in the online A/B testing.

The rest of this paper is organized as follows. Section 2 is a brief introduction of sponsored search and Ad relevance task. Section 3 reviews related literature. Section 4 discusses the details of the model, including the architecture, the construction of graph-type data, and the training methodology. Section 5 reports the experimental results of TextGNN in comparison to the baseline model under both offline and online settings with a few illustrative case study examples. Section 6 concludes the paper and briefly discusses the future directions of this work.

#### 2 SPONSORED SEARCH AND AD RELEVANCE

The TextGNN model is developed to improve the existing Ad Relevance model at a major Sponsored Search platform. In a typical sponsored search ecosystem, there are often three parties: user, advertiser and search engine platform. When the user types a query into the search engine, the goal of the platform is to understand the underlying intent of the user behind the semantic meanings of the query, and then try to best match it with a short list of Ads submitted by the advertisers alongside other organic search results.

<sup>&</sup>lt;sup>1</sup>The BERT version implementation of the model may be found at: https://github.com/microsoft/TextGNN

TextGNN: Improving Text Encoder via Graph Neural Network in Sponsored Search

In the back-end when a query is received by the platform, the system will first conduct a quick but crude recall step using highly efficient Information Retrieval algorithms (such as TF-IDF [11] or BM25 [19]) to retrieve an initial list of matched candidates. The relatively long list is then passed to the downstream components for a finer filtering and final ranking using much more sophisticated but slightly less efficient models to serve the users. In both of the later steps, Deep Learning based Ad Relevance models play a key role in delivering high quality contents to the user and match advertisers' products with the potential customers. For the Ad Relevance task, our model usually relies only on the query from a user and keywords provided by the advertiser. A **query** refers to a short text that a user typed into the search engine when he/she is looking for relevant information or product, and the model needs to identify the user's intent based on the short query. A keyword is a short text submitted by an advertiser that is chosen to express their intent about potential customers. The keyword is in general not visible from end users, but it is crucial for the search engine platform to match user intents.

When an Ad is displayed to a user, we call this an **impression**. The platform does not receive anything from an impression but earns revenue only when the displayed Ad is **click**ed by the user. Because of this mechanism, the search engine platform has an incentive to display the Ads that best match user intents, which directly affects the revenue. Lastly, given the scale of the traffic of the search engine, Ad Relevance models are such an indispensable component of the system and any improvement of the performance of the model can lead to huge impact on the business side of the search engine.

#### **3 RELATED WORK**

**Text Encoders** including C-DSSM and Pre-trained Transformerbased Language Models (such as BERT) have achieved impressive state-of-the-art performance in many NLP tasks for their effective language or contextual word representations, hence have become one of the most important and most active research areas.

C-DSSM is developed specifically for extracting semantic information into a low-dimension representation vector by combining convolutional layers that extract local contextual semantic information in the string with max-pooling layers that helps identifying globally important features. It is still a workhorse model used extensively in the stacks of many production search engine systems.

The large and expensive BERT model has recently become very popular. The model is usually learned in two steps. First the model is trained on extremely large corpus with unsupervised tasks such as masked language model (MLM) and next sentence prediction (NSP) to learn the general language, and then in a second step fine-tuned on the task-specific labelled data to be used in downstream tasks. Despite the strong performance of the BERT models on language representations, they are in general too expensive to be deployed in the real-time search engine systems where there are strict constraints on computation costs and latency.

**Distilled TwinBERT** is one successful model that adapts the Transformer family models to the sponsored search applications and achieves comparable performance at reasonable inference time





Figure 1: Architecture of the twin tower TwinBERT model

cost compared with heavy stacked transformer layers. The Twin-BERT model as demonstrated in Figure 1 benefits from two important techniques: 1) given two input texts, a query and a keyword, a vanilla transformer encoder would concatenate them into one input sequence, while TwinBERT has a twin tower structure to decouple the two-sentence input. Such twin tower structure is first proposed in the DSSM model [9] for web document ranking. Given that the keywords are already known to the platform, the encoded outputs of the keyword-side tower could then be pre-generated offline and fetched efficiently during inference time. Without concatenating the keyword strings, the input to the query-side tower can also be set with a low maximum length, and hence greatly reduce the inference time complexity compared to a large BERT model. 2) knowledge distillation technique is used to transfer the knowledge learnt by a teacher model to a much smaller student model. Our teacher model can be seen as a stronger version of the BM25 signal in the previous weak supervision method [2]. While the teacher model has strong performance, it is usually too costly and infeasible to be directly used in a production system. Knowledge distillation enables us to train a smaller model that is much faster when inference with only little or no significant loss in performance [14][20]. When a TwinBERT model with only 3 layers of encoders is used, with all the optimizations it is possible to be deployed in the real-world production systems that satisfies the strict limit from computational resources and latency requirement.

However, as a pure language model, TwinBERT can only rely on the semantic meanings of the query-keyword pairs to infer the relationships, and in many cases when we encounter uncommon words it is still very challenging to correctly infer relevance for our main applications based on the limited input information.

Graph Neural Network has also become a hot research area in recent years due to its efficacy in dealing with complex graph data. Graph Convolutional Networks (GCN) [13], GraphSage [5], and Graph Attention Networks (GAT) [23] are among the most popular GNN models that can effectively propagate neighbor information in a graph through connected edges and hence are able to generate convincing and highly interpretable results on many graph specific tasks such as node/edge/graph property predictions. Recently there are also attempts to bring GNN to the sponsored search area such as click-through rate (CTR, ratio of the number of clicks to the number of impressions) prediction [15][26], but so far these attempts have only focused on using GNN to generalize the interactions among the existing fixed features. There is no strong convincing story why these features naturally form a graph and the GNN itself has no impact on the generation of the features. Alternatively people have also proposed to utilize the graph information implicitly through label-propagation to unlabeled examples[12], but explicitly using the neighbor features in the model structure will be more efficient in aggregating complementary information as demonstrated in the experiments.

To the best of our knowledge, we are the first to extend various text encoders with a graph in a natural way, and co-train both text encoders and GNN parameters at the same time to achieve stronger performance in our downstream tasks.

#### 4 TEXTGNN

In this section we will discuss the architecture of the proposed TextGNN model in Section 4.1. Then we describe the graph we used to naturally augment the semantic information of the input query-keyword sentence pairs in Section 4.2. Lastly in Section 4.3 we briefly recap knowledge distillation and its application in our model.

#### 4.1 Model Architecture

The architecture of the TextGNN model is discussed in detail in this subsection and also illustrated in Figure 2. The proposed model is a natural extension of the high-performance C-DSSM/TwinBERT baseline model with additional information from graph structured data. In sponsored search scenario, we have tens of millions candidate Ads. It is infeasible to use a complex text encoder to compute the similarity between a search query and each Ad one-by-one. Twin tower structure is a good choice for us where we could compute Ads representation vectors in advance and when a query comes, we then compute the representation vector of the query online. Notice that we only need to run the complex text encoder once for each incoming search query, compared with vanilla BERT which requires this for each unique pair. For transformer encoders, the computation cost in self-attention is also quadratic to the length of the input string. Hence, splitting the query and keyword strings for separate calculation is also much less costly than calculating the concatenated string. With these benefits in mind, our model also follows the twin tower structure of the baseline models with small encoder structure layers so that all the benefits of the twin

tower structured model are inherited and hence can be deployed in the production system. Taking the query-side tower as an example, given a query and its three neighbors (defined later in the graph construction section) will all go through any general Text Encoder blocks to each generate a vector representation for the short sentence. The information from the four representation vectors is then aggregated by a GNN Aggregator to generate a single output vector. This output vector is then connected with the direct output of the text encoder of the query sentence through either concatenation or addition, similar to the idea of a Residual Connection Network [6]. The combined output vector is considered as the final output of the query-side tower and can then be interacted with keyword-side output (generated from the very similar structured keyword-side tower) in the crossing layer to get the final output similar to a C-DSSM/TwinBERT model.

4.1.1 Text Encoder Block. The Text Encoder block is very similar to a single tower in the C-DSSM/TwinBERT model. For example, for a transformer type text encoder, a sentence is first tokenized using the BERT WordPiece tokenizer. Trainable token embedding vectors is combined with BERT style positional embedding through addition before it go through three BERT encoder layers. The only difference with a BERT-style model is that the segment embeddings in the BERT are no longer needed as all inputs will be from the same sentence. With this structure so similar to a BERT-type one, we can conveniently load the weights from the first three layers of the pre-trained large BERT model to get a good starting point that leads to much better performance, faster model convergence, and requires significantly less training data compared to a random initialization. After the text encoder layers, we get a sequence of vectors corresponding to each token in the sentence. The vectors are then combined using a weighted-average pooling layer similar to the TwinBERT model which has demonstrated better performance in generating a single vector representation for a sentence.

The four Text Encoder blocks within a single tower are set to share the same parameters. However, the model is flexible enough to allow the two towers to have all different Text Encoder blocks, but as the TwinBERT paper shows that shared encoder blocks generally lead to slightly better performance we use that approach.

4.1.2 GNN Aggregator. In one tower of our TextGNN, the four text encoder blocks generate four vector representations, one for the center node (query/keyword) and the other three for its three one-hop neighbors. To aggregate the information from four vectors into one, we adopt a GNN aggregation layer, where we take the query/keyword as the central node and perform one-hop aggregation using the three neighbor nodes. The aggregators such as GCN, GraphSAGE, and GAT. In our experiments we found that GAT, which assigns learnable weights to the neighbors to generated a weighted average, demonstrates the strongest performance and is used in our experiments.

4.1.3 Skip Layer. The output vector of the query/keyword encoder is connected to the output of GNN Aggregator as the final output of the query-/keyword-side tower. This layer can be thought as a skip layer [6] so that the additional GNN outputs serve as a complementary information to the text semantic representation



Figure 2: TextGNN Architecture: twin tower structure for decoupled generation of query/keyword embeddings

vector. In this sense the encoder-only-models can also be considered as a special case of the TextGNN model when the GNN output is completely skipped. The two vectors are combined using either concatenation or addition. In case they have different dimensions an additional dense layer is applied after the GNN Aggregator to up/downscale the GNN output dimension to match the Text Encoder output.

4.1.4 *Crossing Layer.* Given the final outputs of the query-/keywordside tower, the two vectors are first combined through concatenation, and then compute the similarity score using the Residual network proposed in the TwinBERT model. Formally, the residual function is defined as:

$$\mathbf{y} = \mathcal{F}(\mathbf{x}, W, b) + \mathbf{x},\tag{1}$$

where **x** is the concatenation of the query-side vector **q** and keywordside vector **k** and  $\mathcal{F}$  is the mapping function from **x** to the residual with parameters *W* and *b*. A logistic regression layer is then applied to the output vector **y** to predict the binary relevance label.

#### 4.2 Graph Construction

On top of the powerful structure of the model, it is also crucial to get access to high quality graph-type data. Such data should satisfy the following properties:

- (1) **Relevant:** since the graph neural networks propagate information along the edges, we are looking for neighbors that are highly relevant to the intent of the center node (query/keyword).
- (2) Complementary: we expect the GNN to excel the most in situations where the language modeling part struggles to infer the intention only from the semantic meanings of the sentence, but the additional neighbors might be extremely valuable to provide complementary information that help the model to better understand the inputs. This situation happens most frequently on rare and low frequency items

where the language models usually struggles on these longtail inputs.

(3) Accessible: in sponsored search system, there are large amount of user input queries and candidate keywords. We try to find their neighbors in a graph. As a large graph is preferred, the neighbors need to be found with little effort and constructing the graph data should be feasible without heavy manual work, strong assumptions, or complicated structures.

Given the requirements, we find that the user behavior graph generated from historical user clicking logs is a great candidate for our purpose. It is based on the insight that when a user inputs a query a and then clicks the Ad b, then b has to sufficiently fit the user's intent from a to trigger the click. In the next two subsections, we discuss such behavior graph and its extension to address the sparse coverage issue of the behavior graph.



Figure 3: Click Graph Construction: use ANN proxy neighbor if no native neighbor available

4.2.1 User Click Graph. The eligible neighbors of a query are the keyword of Ads that have been shown to be relevant to the query and received explicit positive feedback by a click. One general assumption to sort all the candidates is that the empirically observable CTR is highly correlated to the relevance between the query and

Table 1: Example of neighbors of a query from the Click Graph

Query	Clicked Neigh Keyword	Neigh # Impress	Neigh # Click
usps com	united state postal service jobs	59	18
careers login	usps com employment	344	92
	postal service hiring	1721	384

the keyword. Based on this assumption, as illustrated in Figure 3(a), we take all clicked Ads that have been shown to users at least 50 times in the past year (to partially address the issue of noisy estimates of CTR on Ads with small number of impressions) and take the top three as the neighbors.

Table 1 shows an illustrative example, where the search query is "usps com careers login". Its top three neighbors, which are the keywords of the corresponding Ads, are listed with their historical total number of impressions and clicks. Although the first keyword "united state postal service jobs" is only shown 59 times which is significantly fewer than the third keyword "postal service hiring" with 1,721 impressions, it has a much higher CTR of 30.5% compared to 22.3%, indicating that users who searched for this query are more likely to find the first keyword useful, which is a strong indication of higher relevance.

4.2.2 User Click Graph with Semantic ANN. For rare and low frequency queries/keywords, we observe by construction substantially less feedback from clicks logs. Furthermore, to avoid the noise of selecting neighbors with high CTR, we have criteria to exclude neighbors that are shown less than 50 times in the past year and this unfortunately eliminates a number of neighbors and makes the situation even worse for long-tail inputs. To address this issue, we propose a neighbor completion technique based on Approximate Nearest Neighbor (ANN) [10] using Neighborhood Graph Search (NGS) [24]. As illustrated in Figure 3(b), first we infer vector representations by a powerful C-DSSM (which is used extensively in a major sponsored search system) for all nodes in user click graph. Next, for a query that we could not identify any eligible clicked keywords, we infer its vector representation by the same C-DSSM. Then, we leverage the ANN search tool to find another query that is supposed to be semantically close enough to the original query and has the click neighbors and use its clicked keywords as approximate neighbors for the original query. This has the same spirit as the common technique of query rewriting in search engine systems but does so in a more implicit way. For keywords without any clicked queries, we find neighbors for them in a similar way.

In Table 2 we show another example that we are not able to find any eligible neighbors for the query "video games computers free", but its ANN query "no internet games" has user behavior feedback and the three approximate neighbors are obviously relevant to the original query.

For both types of graphs, we only take at most the top three neighbors. The number of neighbors can be set as a hyper-parameter of the model framework. We choose three for following reasons:

Query	ANN Query	Clicked Neigh Keyword	Neigh # Impress	Neigh # Click
video		free games	58	1
games	no	online games	260	4
computers	internet	online		
free	games	computer games	67	1

- More than one neighbor to provides additional complementary information while also adds robustness.
- (2) Each additional neighbor means an extra run of the text encoder. Even though the encoder blocks can be run in parallel a large number of neighbors can still be computationally challenging for the system.
- (3) We do not want to include more neighbors that are less relevant and introduce additional noisy information to "pollute" the encoded representation.

Therefore, choosing three neighbors balances all the requirements and concerns.

#### 4.3 Knowledge Distillation

In order to have a high performance but compact model that satisfies the computation and latency constraints, the teacher-student training framework via knowledge distillation is used. We use an expensive but high-performance RoBERTa model [16] as the teacher model to label a very large query-keyword pair dataset, the label scores are between 0 and 1. Our model is relatively data-hungry and without this teacher model to automatically label the huge dataset, our existing human-labelled data is not sufficient to train a strong model that gets close to teacher model level performance. Since the model target, the RoBERTa score, is a continuous value, it provides more fine-grained information than the traditional binary labels. For example, a score of 0.99 indicates a stronger relevance than a score of 0.51, although both will be categorized as relevant pairs. We use mean squared error to measure the difference between the model output and the RoBERTa teacher scores.

With such a strong teacher model, we train the student Twin-BERT/TextGNN model with small encoder blocks (only 3 transformer layers). Hence the student models are much more feasible in inference time but are able to achieve close to teacher model performance with only very minor performance loss. We could even further finetune the student model on a smaller human-labelled dataset with binary labels and achieve a performance surpassing the much larger teacher model. Hence, the performance of our model is not capped/limited by the teacher model.

#### **5 EXPERIMENTS**

In this section we present experiment results of TextGNN on various tasks. We also show the comparison with the strong baseline models to show the superiority of the proposed new model and the efficacy of introducing graph information. In Section 5.1 we discuss some key statistics of the complementary graph data, and some related details of our training methods. Section 5.2 compares the TextGNN: Improving Text Encoder via Graph Neural Network in Sponsored Search

Table 3: Coverage Summary of Two Graph Construction Methods: almost full coverage after adopting ANN Neighbors

	Click Only		ANN	
	Q	Κ	Q	K
1 Neighbor	4%	7%	5%	7%
2 Neighbors	3%	4%	3%	4%
3 Neighbors	30%	76%	92%	88%
Coverage	37%	87%	100%	99%

performance with the baseline encoder-only models. Section 5.3 shows a more detailed sub-group analysis. Section 5.4 presents case studies of typical examples with false positive and false negative examples for TwinBERT which are correctly classified by the new TextGNN model and provide intuitive insights why the additional graph information can be valuable. Lastly in Section 5.5 we present an initial effort to apply our model to online production system and show the significant improvement over the baseline in online A/B testings.

#### 5.1 Data and Training Details

For our knowledge distillation training, 397 million query-keyword pairs are scored by the teacher RoBERTa model. The student models are initialized using the parameters of the first three transformer layers of the 12-layer uncased BERT-base checkpoint [25]. The models are evaluated on a small evaluation dataset consisting of 243 thousand human labelled samples. The query and keyword pairs were given labels with five different levels of relevance: excellent, perfect, good, fair, and bad. In the evaluation stage the first four levels excellent, perfect, good, and fair are mapped as positive samples (label 1) where the bad category is kept as negative category (label 0). The model ROC-AUC is our main metric for evaluation.

We construct the behavior click graph based on the historical search engine click logs from July 2019 to June 2020. Here in Table 3 we present some statistics on the neighbor coverage comparing the two ways of graph constructions. Here are some key observations:

- Without the added ANN neighbors, almost 2/3 of the queries miss neighbors from the user click graph. The situation is significantly better for keywords as the majority of the Ads have been shown and clicked by users.
- (2) With the ANN search, we essentially increase the neighbor coverage to almost 100%.
- (3) Among all nodes, the majority of them have at least three eligible neighbors. For the examples with less than 3 neighbors, dummy padding are added.

#### 5.2 Model Performance Results

In the experiment we train the baseline TwinBERT model and the new TextGNN model with the same common hyper-parameters for a fair comparison. The same training dataset files were used by both models, but the additional neighbor information is not read by the baseline TwinBERT model as it does not have the mechanism to process the additional information. Tabel 4 presents the ROC-AUC values of the baseline model and TextGNN based on two different types of graphs. We see that the addition of GNN has significantly improves the performance of the baseline model and the performance increase of this magnitude will lead to a huge difference in revenue for large scale systems.

 Table 4: ROC-AUC Comparison: TextGNN with ANN Neighbor Graph significantly outperform baseline TwinBERT

Model	AUC
TwinBERT	0.8459
TextGNN	0.8461
TextGNN with ANN Neighbor	0.8471

#### 5.3 Sub-group Analysis

In addition to showing the stronger overall performance of the TextGNN models over the baseline, we also conduct a more detailed sub-group analysis on inference results to confirm that the TextGNN models indeed improve on the tail examples just as expected.

We split the validation data into three bins by the Ads frequency in the dataset (as a proxy for their population frequency of impressions). 43% of the samples are Ads that have been shown only once (among 243k samples) which are the rare examples, and 12% of the samples have been shown twice. Even though the tail Ads individually are rarely recalled and shown to users, they consist of the majority portion of the total traffic and the improvements on these long-tail examples can lead to significant benefits.

We see the results in Figure 4 that the TextGNN model based on vanilla click graph shows an extremely large improvement in the most rare Ads, but the performance downgrades in common ones. Our hypothesis is that in the more common examples the semantic information is already good, and the limited additional information from a sparse graph is not enough to offset the potential underfitting from a more complex model. Once we adopt ANN to generate a more complete graph, we see the TextGNN model demonstrates stronger performance than baseline across the board.

Lastly, we note that the non-ANN version is still much stronger than the ANN version in the bin of the most rare Ads, potentially because the ANN proxy neighbors are on average having lower quality than the native neighbors, and hence introduce noise to the model. This analysis also reveals a future direction to further improve the model where we can potentially use the sample frequency as a simple indicator to switch between various candidate models based on their strength within different sub-groups.

#### 5.4 Case Studies

We expect the introduction of graph data to improve the model performance especially on tail inputs that are often seen as "hard" samples for the baseline models. In table 5, we present some "hard" cases to demonstrate the value that graph data could bring.

5.4.1 False-positive Examples of TwinBERT. The first example shows that the user searched for the Greek methology "achilles heel", which was incorrectly determined by TwinBERT as relevant to plantar fasciitis shoes. From the semantic meaning, heel is very

False Positive Examples			
Query	Query Neighbors	Keyword	Keyword Neighbors
achilles heel	what is an achilles heel		shoes plantar fasciitis heel pain
	what is achilles heel	plantar fasciitis shoes	work shoes plantar fasciitis
	causes heel spurs		tennis shoes good plantar fasciitis
animal repellent products	animal repeller		best cleaning remove & product home
	keep squirrel out attic	animal odor	air fresheners home
	animal repellent		best air fresheners
False Negative Examples			
Query	Query Neighbors	Keyword	Keyword Neighbors
sharding	mongodb cluster		sql server download windows 10
	database sharding	sql server	sql server hosting
	N/A		sequel server database
use imovie	imovies		adobe premiere pro mac
	imovie 11 tutorials	adobe premiere	adobe premier mac
	imovie video editor		use imovie

	- 1 • 1 1		1
highle by Case stud	v Evamples, peighbor	s provide criicial com	nlementary information
Table J. Case stud	y Linampies, neighbor	s provide er delar com	prementary mormation



Figure 4: Performance on Different Subgroups of Data by Ads Frequency: TextGNN with vanilla click neighbor achieves extremely large gain in low frequency Ads, while the ANN version outperforms the baseline across the board

close to shoes and the achilles ankle is highly related to the pain of tendon. However, the neighbors strongly indicate that people who search for this query are actually looking for the story from Greek mythology and not the foot injury.

The second example shows that TwinBERT determines that "animal repellent products" is highly relevant to animal cleaning product. From the semantic meaning it is true that repellent is close in meaning to the word "remove" but the two products are used for completely different purposes. When averaging over the neighbors it is very clear that this is a negative example.

5.4.2 False-negative Examples of TwinBERT. The query "sharding" is a very specific concept in database systems on how large data are split and stored. Without the domain knowledge it is very hard to understand such an uncommon word. Furthermore, the word is

tokenized to: [CLS], sha, ##rdi, ##ng, [SEP] by the BERT WordPiece tokenizer, making it essentially an impossible task for TwinBERT to identify the relevance. However, from the historical user behaviors we clearly see both sides taking the very important common words "database", hence allowing the TextGNN model to leverage on the user behavior to identify domain specific connections and find the hidden relevance.

The second false-negative one is an example of two video editing softwares on the Mac platform. Without the domain knowledge is it impossible to conclude from the semantic meaning that adobe premier mac is a video editing software. However, since the query string is identified as a neighbor of the keyword, our graph model can use this information to find the correct connection.

#### 5.5 Online A/B Test

A slightly simplified version of our TextGNN model has already been successfully deployed in a major sponsored search platform and demonstrated significant performance gains. We have evaluated the performance of the models on the sponsored product advertising system where user search queries are matched with products with rich information provided by advertisers. In this initial effort we choose C-DSSM as the text encoder for its much faster inference time in the application of large-scale Ads corpus and use graph aggregators only on the product side of the tower. Note again that the product side representations can be generated offline in advance and hence at online service stage the latency is identical to a traditional C-DSSM model. We use the TextGNN model outputs as features to be feed into a downstream online product advertising system and evaluated the efficacy of this simple model in both offline and online settings.

For evaluation, we randomly sampled examples from online logs and labeled the data manually by human experts and observe on average 1.3% (we only show normalized relative numbers due to business confidentiality) PR-AUC lift across different validation sets when comparing the simplified TextGNN model with the baseline C-DSSM model. TextGNN: Improving Text Encoder via Graph Neural Network in Sponsored Search

The online A/B testing results of the TextGNN model are summarized in Table 6 as we applied the model to both recall and relevance stage of the Ads serving in the system, where we observe significant gains in several normalized key online metrics numbers that are crucial for our sponsored search system. The two most important metrics are:

- (1) **Revenue Per Mille (RPM):** the revenue gained for every thousand search requests, which is one of the most important online metrics for sponsored search.
- (2) Ad Defect Rate: the ratio of irrelevant Ad impressions with respect to total number of Ad impressions. In online A/B test, this ratio is approximated by sampling Ad impressions and submitting them for human-evaluated labels. This is highly correlated to user satisfaction and hence is considered as a very crucial metric.

As shown in the table, the TextGNN model yields very impressive results as it can greatly boost the RPM and reduce the Ad Defect Rate, which is a strong sign that model could help to improve revenue and user experience simultaneously. It's worthy pointing out that current production model already contains many advanced sub-models and features so the magnitude of the improvement in the online KPI here is considered as a significant gain for our system at the large scale.

# Table 6: Online A/B Testing: significant improvements in production product advertising systems

Tasks	Relative RPM	Relative Ad Defect Rate
TextGNN Relevance	+2.03%	-2.32%
TextGNN Selection	+1.21%	-0.34%

### 6 CONCLUSION

We present a powerful NLP model TextGNN that combines two strong model structures, text encoders and GNN, into a single endto-end framework and shows strong performance in the task of Ad relevance. The model retains the strong natural language understanding ability from the existing powerful text encoders, while complements text encoders with additional information from graphtype data to achieve stronger performance than what could be achieved from only pure semantic information. We demonstrate with experiments that the TextGNN model show overall much stronger performance than a great baseline model based only on text encoders, and that the new model demonstrates the big gains in the most difficult task of low-frequency Ads. In our next step, the ensemble model idea could be explored to automatically mix different representation model outputs based on Ads frequency to achieve even better performance.

#### REFERENCES

[1] Xiao Bai, Erik Ordentlich, Yuanyuan Zhang, Andy Feng, Adwait Ratnaparkhi, Reena Somvanshi, and Aldi Tjahjadi. 2018. Scalable Query N-Gram Embedding for Improving Matching and Relevance in Sponsored Search. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (London, United Kingdom) (KDD '18). Association for Computing Machinery, New York, NY, USA, 52–61. https://doi.org/10.1145/3219819.3219897

- [2] Mostafa Dehghani, Hamed Zamani, Aliaksei Severyn, Jaap Kamps, and W. Bruce Croft. 2017. Neural Ranking Models with Weak Supervision. In Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval (Shinjuku, Tokyo, Japan) (SIGIR '17). Association for Computing Machinery, New York, NY, USA, 65–74. https://doi.org/10.1145/ 3077136.3080832
- [3] 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
- [4] Mihajlo Grbovic and Haibin Cheng. 2018. Real-Time Personalization Using Embeddings for Search Ranking at Airbnb. In Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (London, United Kingdom) (KDD '18). Association for Computing Machinery, New York, NY, USA, 311–320. https://doi.org/10.1145/3219819.3219885
- [5] Will Hamilton, Zhitao Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In Advances in Neural Information Processing Systems 30, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.). Curran Associates, Inc., 1024–1034. http://papers.nips.cc/ paper/6703-inductive-representation-learning-on-large-graphs.pdf
- [6] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Deep Residual Learning for Image Recognition. http://arxiv.org/abs/1512.03385 cite arxiv:1512.03385Comment: Tech report.
- [7] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the Knowledge in a Neural Network. arXiv:1503.02531 [stat.ML]
- [8] Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padmanabhan, Giuseppe Ottaviano, and Linjun Yang. 2020. Embeddingbased Retrieval in Facebook Search. Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (Aug 2020). https://doi.org/10.1145/3394486.3403305
- [9] Po-Sen Huang, Xiaodong He, Jianfeng Gao, Li Deng, Alex Acero, and Larry Heck. 2013. Learning Deep Structured Semantic Models for Web Search using Clickthrough Data. ACM International Conference on Information and Knowledge Management (CIKM). https://www.microsoft.com/enus/research/publication/learning-deep-structured-semantic-models-for-websearch-using-clickthrough-data/
- [10] Piotr Indyk and Rajeev Motwani. 1998. Approximate Nearest Neighbors: Towards Removing the Curse of Dimensionality. 604-613.
- [11] Karen Spärck Jones. 1972. A statistical interpretation of term specificity and its application in retrieval. *Journal of Documentation* 28 (1972), 11–21.
- [12] Soo-Min Kim, Patrick Pantel, Lei Duan, and Scott Gaffney. 2009. Improving Web Page Classification by Label-Propagation over Click Graphs. In Proceedings of the 18th ACM Conference on Information and Knowledge Management (Hong Kong, China) (CIKM '09). Association for Computing Machinery, New York, NY, USA, 1077–1086. https://doi.org/10.1145/1645953.1646090
- [13] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In Proceedings of the 5th International Conference on Learning Representations (Palais des Congrès Neptune, Toulon, France) (ICLR '17). https://openreview.net/forum?id=SJU4ayYgl
- [14] Xue Li, Zhipeng Luo, Hao Sun, Jianjin Zhang, Weihao Han, Xianqi Chu, Liangjie Zhang, and Qi Zhang. 2019. Learning Fast Matching Models from Weak Annotations. In *The World Wide Web Conference*. Association for Computing Machinery, 2985–2991.
- [15] Zekun Li, Zeyu Cui, Shu Wu, Xiaoyu Zhang, and Liang Wang. 2019. Fi-GNN: Modeling Feature Interactions via Graph Neural Networks for CTR Prediction. In Proceedings of the 28th ACM International Conference on Information and Knowledge Management (Beijing, China) (CIKM '19). Association for Computing Machinery, New York, NY, USA, 539–548. https://doi.org/10.1145/3357384.3357951
- [16] 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:1907.11692 [cs.CL]
- [17] Wenhao Lu, Jian Jiao, and Ruofei Zhang. 2020. TwinBERT: Distilling Knowledge to Twin-Structured BERT Models for Efficient Retrieval. arXiv:2002.06275 [cs.IR]
- [18] Aditya Pal, Chantat Eksombatchai, Yitong Zhou, Bo Zhao, Charles Rosenberg, and Jure Leskovec. 2020. PinnerSage: Multi-Modal User Embedding Framework for Recommendations at Pinterest. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining* (Virtual Event, CA, USA) (KDD '20). Association for Computing Machinery, New York, NY, USA, 2311–2320. https://doi.org/10.1145/3394486.3403280
- [19] Stephen Robertson, S. Walker, S. Jones, M. M. Hancock-Beaulieu, and M. Gatford. 1995. Okapi at TREC-3. In Overview of the Third Text REtrieval Conference (TREC-3) (overview of the third text retrieval conference (trec-3) ed.). Gaithersburg, MD: NIST, 109–126. https://www.microsoft.com/en-us/research/publication/okapiat-trec-3/
- [20] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. CoRR

abs/1910.01108 (2019). arXiv:1910.01108 http://arxiv.org/abs/1910.01108

- [21] Yelong Shen, Xiaodong He, Jianfeng Gao, Li Deng, and Gregoire Mesnil. 2014. Learning Semantic Representations Using Convolutional Neural Networks for Web Search. WWW 2014. https://www.microsoft.com/enus/research/publication/learning-semantic-representations-usingconvolutional-neural-networks-for-web-search/
- [22] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, undefinedukasz Kaiser, and Illia Polosukhin. 2017. Attention is All You Need. In Proceedings of the 31st International Conference on Neural Information Processing Systems (Long Beach, California, USA) (NIPS'17). Curran Associates Inc., Red Hook, NY, USA, 6000–6010.
- [23] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. *ICLR* (2018).
- [24] Jingdong Wang and Shipeng Li. 2012. Query-Driven Iterated Neighborhood Graph Search for Large Scale Indexing. In Proceedings of the 20th ACM International Conference on Multimedia (Nara, Japan) (MM '12). Association for Computing Machinery, New York, NY, USA, 179–188.
- [25] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, R'emi Louf, Morgan Funtowicz, and Jamie Brew. 2019. HuggingFace's Transformers: State-of-the-art Natural Language Processing. ArXiv abs/1910.03771 (2019).
- [26] Xiao Yang, Tao Deng, Weihan Tan, Xutian Tao, Junwei Zhang, Shouke Qin, and Zongyao Ding. 2019. Learning Compositional, Visual and Relational Representations for CTR Prediction in Sponsored Search. In Proceedings of the 28th ACM International Conference on Information and Knowledge Management (Beijing, China) (CIKM '19). Association for Computing Machinery, New York, NY, USA, 2851–2859. https://doi.org/10.1145/3357384.3357833
- [27] Jie Zhou, Ganqu Cui, Zhengyan Zhang, Cheng Yang, Zhiyuan Liu, Lifeng Wang, Changcheng Li, and Maosong Sun. 2019. Graph Neural Networks: A Review of Methods and Applications. arXiv:1812.08434 [cs.LG]