

# OMNISEARCHSAGE: Multi-Task Multi-Entity Embeddings for Pinterest Search

Prabhat Agarwal\*  
pagarwal@pinterest.com  
Pinterest  
USA

Minhazul Islam SK\*  
msk@pinterest.com  
Pinterest  
USA

Nikil Pancha  
npancha@pinterest.com  
Pinterest  
USA

Kurchi Subhra Hazra  
ksubhrahazra@pinterest.com  
Pinterest  
USA

Jiajing Xu  
jiajing@pinterest.com  
Pinterest  
USA

Chuck Rosenberg  
crosenberg@pinterest.com  
Pinterest  
USA

## ABSTRACT

In this paper, we present OMNISEARCHSAGE, a versatile and scalable system for understanding search queries, pins, and products for Pinterest search. We jointly learn a unified query embedding coupled with pin and product embeddings, leading to an improvement of > 8% relevance, > 7% engagement, and > 5% ads CTR in Pinterest’s production search system. The main contributors to these gains are improved content understanding, better multi-task learning, and real-time serving. We enrich our entity representations using diverse text derived from image captions from a generative LLM, historical engagement, and user-curated boards. Our multi-task learning setup produces a single search query embedding in the same space as pin and product embeddings and compatible with pre-existing pin and product embeddings. We show the value of each feature through ablation studies, and show the effectiveness of a unified model compared to standalone counterparts. Finally, we share how these embeddings have been deployed across the Pinterest search stack, from retrieval to ranking, scaling to serve 300k requests per second at low latency. Our implementation of this work is available at this link<sup>1</sup>.

## CCS CONCEPTS

• **Information systems** → **Retrieval models and ranking; Web searching and information discovery.**

## KEYWORDS

Multi-Task Learning, Multimodal Embeddings, Representation Learning, Search Recommendation Systems

## ACM Reference Format:

Prabhat Agarwal, Minhazul Islam SK, Nikil Pancha, Kurchi Subhra Hazra, Jiajing Xu, and Chuck Rosenberg. 2024. OMNISEARCHSAGE: Multi-Task Multi-Entity Embeddings for Pinterest Search. In *Companion Proceedings of the*

ACM Web Conference 2024 (WWW ’24 Companion), May 13–17, 2024, Singapore, Singapore. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3589335.3648309>

## 1 INTRODUCTION

Pinterest’s mission is to bring everyone the inspiration to create a life they love. Search is one of the key surfaces on Pinterest where users seek inspiration spanning a wide range of interests, such as decorating their homes, planning weddings, or keeping up with the latest trends in beauty and fashion. In order to enhance the search experience, modern search systems aim to incorporate various types of content such as web documents, news, shopping items, videos, and more. Similarly, Pinterest’s search feed encompasses a diverse range of content, including pins, shopping items, video pins, and related queries. To construct an inspiring feed for each of the more than 6 billion searches per month on Pinterest we must uncover relevant content from billions of pins and products. We must also find relevant queries to help users refine their queries and navigate their search journey.

As an additional challenge, Pinterest search is global and multi-lingual with searchers using more than 45 languages to find inspirational content.

Embeddings are useful building blocks in recommendation systems, especially search, where natural language understanding is key [11, 23, 24]. Embeddings can power retrieval use cases via approximate nearest neighbor (ANN) search [14, 22], enable detailed content and query understanding in ranking models without the overhead of processing raw data, and serve as a strong base to learn in low-data use-cases [31]. Despite their utility, embeddings come with their own challenges: if we learn a separate embedding for every use-case, there is an explosion of potentially expensive models that must be inferred on every request and used in downstream models. This also may lead to suboptimal recommendation quality – some use-cases may not have enough labels to learn an optimal representation. In practice, it could entail additional maintenance costs and technical debt for upgrading to new versions of embeddings in certain applications, as some data may have been collected over the course of months or years.

Through rigorous offline experimentation, we show the impact of our key decisions in building embeddings for web-scale search at Pinterest:

- Pin and product representations can be substantially enriched using diverse text derived from image captions from

\*Both authors contributed equally to this work.

<sup>1</sup><https://github.com/pinterest/atg-research/tree/main/omnisearchsage>



This work is licensed under a Creative Commons Attribution International 4.0 License.

a generative LLM, historical engagement, and user-curated boards.

- A single query embedding can be used to retrieve queries, products, and Pins with nearly the same effectiveness as task-specific embeddings.
- A single query embedding can learn compatibility with multiple pre-existing embeddings *and* learned entity embeddings, and perform well when compared across tasks.

OMNISEARCHSAGE has been deployed at Pinterest and is an integral component of the search stack. It powers embedding-based retrieval for standard and product pins, queries and ads. It is also one of the most important feature in multi-stage ranking models and various query classification models. These gains all arise despite the existence of other features enabling pin and product understanding, which highlights the importance optimizing embeddings end-to-end for search.

## 2 RELATED WORK

Our work to build multi-task multi-entity embeddings for search draws upon broad areas of work. Our representation of pins and products extends existing work on multi-modal learning and two tower models for search retrieval. These have been extensively applied in the context of search and recommendation systems as an efficient way to retrieve results not purely related to the search query based on text. In OMNISEARCHSAGE, we demonstrate that the embeddings generated by these models can also serve as features in ranking and relevance models. Additionally, we offer a brief examination of specific embeddings within the Pinterest ecosystem.

### 2.1 Model-based Search Retrieval

Historically, search systems have been powered by two stages: token-based matching, or candidate generation, and then scoring with a complex model. These have drawbacks, especially when users make complex queries or content is not primarily textual. This has led to the exploration of two tower models, which encode a query into a single embedding or a small set of embeddings, and then use those to retrieve relevant documents with approximate or exact nearest neighbor search [5, 11, 18, 20, 21, 24, 40].

Two natural topics in learning embeddings for search are document representation, and query representation. Depending on the learning objective, this query representation could be personalized, or it could be a pure text embedding model. Many architectures for query embeddings in industry have been proposed based on simple CNNs [12], bag of words models [11, 23], transformers [19], and more, but they share a basic structure involving query understanding and sometimes context understanding. Document representation is also a major challenge. The text associated directly with an item is popular as a key feature, but depending on the task, other sources have been found to provide great value, including queries where other users have engaged with a given item [5, 24, 25] and image content embeddings [19].

### 2.2 Multi-task, multi-modal, and multi-entity embeddings

The area of learning embeddings isn't exclusive to the realm of recommendation systems and has been studied extensively [4, 6,

29, 30]. Multi-task learning is a technique commonly utilized in ranking models to optimize for multiple objectives concurrently, aiming for enhanced performance or more efficient information sharing [33, 41]. A less frequently encountered approach involves the joint learning of embeddings for more than two entities. Though this methodology is sometimes implemented in graph learning scenarios, it can also be perceived as an extension of multi-task learning [39].

Multi-modal embeddings are of substantial interest in the industry since the majority of web content is multi-modal, typically including at both text and images [18, 19, 38]. One can take embeddings or raw data from each modality as inputs, and merge them at any stage of the model. The methodology typically involves utilizing embeddings or raw data from each mode as inputs, which are then merge at different stages in the model. Early-stage fusion can pose computational hurdles; therefore, in cases where performance is indifferent, utilizing embeddings instead of raw data is generally the preferred course of action [38].

### 2.3 Embeddings at Pinterest

PinSage [37] is a scalable GNN-based embedding representing pins. It is based on the GraphSage GCN algorithm [10], sampling neighborhoods with personalized PageRank to augment pin understanding, instead of simple heuristics like  $n$ -hop neighbors. It aggregates some basic visual [2] and text information into a single dense representation, and is a critical feature in many models.

To represent products, we have an embedding, ItemSage [1], which aggregates raw data about products, including metadata from product pages, and potentially many images of the product. ItemSage is trained for compatibility with PinSage, and the search query embedding preceding OMNISEARCHSAGE, meaning that the distance between ItemSage and these two embeddings can be used for retrieving or ranking content [27].

## 3 METHOD

### 3.1 Problem Formulation

In order to enhance the search experience, modern search systems aim to incorporate various types of content such as web documents, news, shopping items, videos, and more. Similarly, Pinterest's search feed encompasses a diverse range of content, including pins, shopping items, video pins, and related queries. Training separate query embedding models for each content type and its representation proves to be resource-intensive and inefficient. To address this issue, we introduce OMNISEARCHSAGE, which offers a unified query embedding model that jointly trains query embeddings for query-query, query-pin, and query-product retrieval and ranking.

Another requirement in production systems is compatibility with existing embeddings, which is essential for purposes such as cost-efficiency and simplified migration. Hence we also train the query embeddings to be compatible with the corresponding pre-existing embeddings for the entities. As a side effect, we also get compatibility with some embeddings due to the triangle inequality property inherent to cosine similarity.

### 3.2 Enriching Entity Representations

On Pinterest, each pin or product is associated with an image and title, along with an optional text (known as description) and link. Beyond these typical attributes, products may carry additional meta-data, such as brand information, color description, and more. Document expansion techniques has been empirically demonstrated to significantly enhance the performance of not just token-based, but also embedding-based search retrieval systems [8, 25, 26, 28, 34]. Hence, in OMNISEARCHSAGE, we enrich our entity representations using diverse text derived from image captions from a generative LLM, historical engagement, and user-curated boards as described below. In the dataset, 71% of pins and products feature a title or description, 91% include non-empty board titles, and 65% contain non-empty engaged queries. Synthetic GenAI captions are generated for all pins and products, ensuring full coverage. Section 4.3.2 discusses the importance of each of these enrichment.

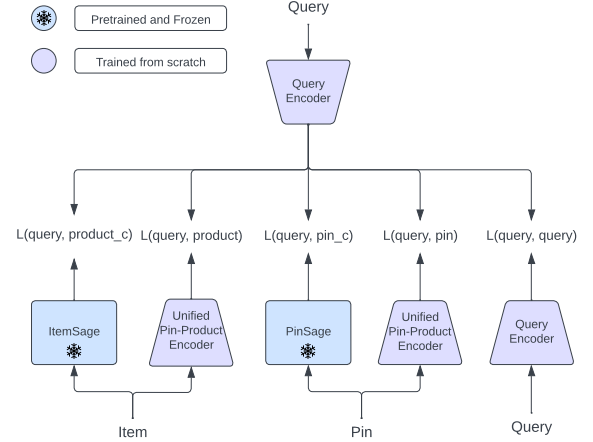
**3.2.1 Synthetic GenAI Captions.** On our platform, a substantial volume of pins (about 30%) lack associated titles or descriptions, or possess noisy and/or irrelevant title or description. We address this issue by employing an off-the-shelf image captioning model, BLIP [17], to generate synthetic descriptions for these images.

To assess the quality of these synthetically generated descriptions, we enlisted human evaluators to judge their relevance and quality. For a robust assessment, three distinct ratings were collected for each image within a sample of 10k images, curated uniformly across various broad pin categories. The results indicated that an overwhelming 87.84% of the generated descriptions were both relevant and of high quality, while a meager 1.16% were deemed irrelevant and of poor quality.

These synthetically generated descriptions serve as an added feature in our model, enriching the diversity of data associated with each entity. Despite not being directly visible to the users, their addition significantly contributes to a deeper understanding of the pins' content.

**3.2.2 Board Titles.** On Pinterest, users explore and save pins to their personal collections, referred to as boards. Each board carries an associated title, reflecting the topic or theme of the collection. Most often, these user-crafted boards are meticulously organized, each focusing on a distinct theme or purpose. A user might, for instance, create discrete boards for "Social Media Marketing" and "Graphic Design". Consequently, these board titles provide valuable, user-generated descriptors for the pins within the respective boards.

We exploit this user-curated information by accumulating the titles of all boards each pin has been saved to. We limit our selection to a maximum of 10 unique board titles for each pin/product, systematically eliminating any potentially noisy or redundant titles as described next. First, each title is assigned a score influenced by two factors: its frequency of occurrence and the prevalence of its comprising words. Following this, titles are then ranked based on a hierarchy of their score (ascending), word count (descending), and character length (descending). The resulting top 10 board titles are subsequently incorporated as a feature in our model. This process eliminates any potentially noisy or redundant titles from the feature.



**Figure 1: Diagrammatic Representation of OMNISEARCHSAGE's Multi-Entity, Multi-Task Architecture.**

**3.2.3 Engaged Queries.** When multiple users interact with a specific pin or product for a certain query within a search feed, it signifies that pin's relevance to that query. We can use these queries to expand our understanding of the pin/product. For every pin, we generate a list of queries that have attracted user engagements, along with the counts and types of such engagements. This list of queries is then sorted using a function based on the count for each type of engagement. We use the top 20 queries from these sorted lists as a feature in our model.

Through experimentation with diverse time-windows of query logs for feature creation, we discovered that larger windows yield superior performance. Consequently, we have opted for a two-year window for feature calculation. However, the complexity of computing this from scratch every time presents a challenge. To mitigate this, we deploy an incremental approach. Every  $n$  days, we examine new query logs, create a list of queries for every pin, and then blend it with the previously existing top 20 queries, thereby updating the latest value of the feature.

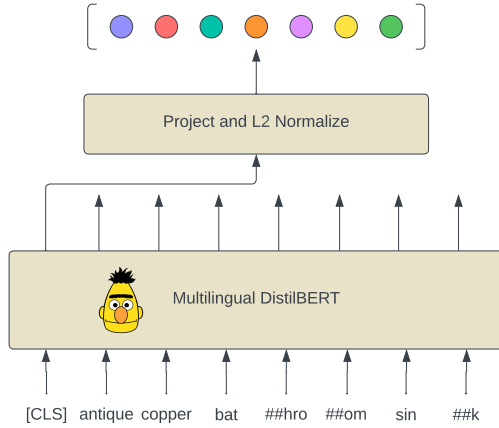
### 3.3 Entity Features

The features we incorporate include PinSage [37] and unified image embeddings [2] to capture the essence of each pin. Additionally, for product pins, we use ItemSage [1] given its capability in effectively representing product-related pins. Text-based features such as the title and description of each pin are also integral to our feature set. Furthermore, we augment the text associated with each pin with the inclusion of synthetic captions, board titles, and engagement queries as outlined earlier. By integrating all these features, we attain a comprehensive and multi-dimensional representation of each pin, hence facilitating enhanced learning of representations.

### 3.4 Encoders

In our work, we consider 3 entity types, namely, pin, product and query. Our model consists of an encoder for query, a unified learned encoder for both pin and product, and dedicated compatibility encoders for pin and product, respectively.

**3.4.1 Query Encoder.** The query encoder in our model (depicted in Figure 2) is based on a multilingual version of the DistilBERT



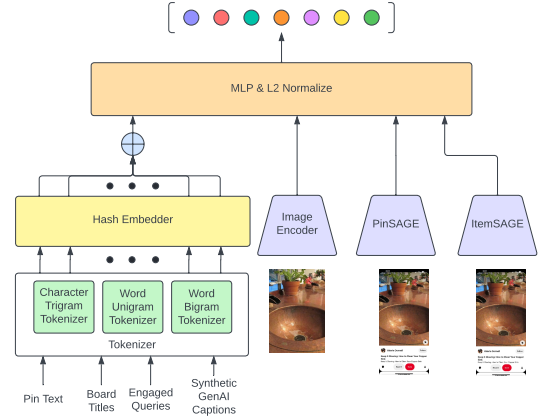
**Figure 2: Overview of the query encoder architecture.** The encoder takes the output from the last layer associated with the ‘CLS’ token, projects it onto a 256-dimensional vector space, and finally L2-normalizes the output to generate the final embedding.

(*distilbert-base-multilingual-cased*<sup>2</sup>) [32]. This choice facilitates efficient handling of queries across a variety of languages. The encoder utilizes the output from the last layer corresponding to the *CLS* token and thereafter projects it to a 256-dimensional vector space. Post projection, we apply a L2 normalization on the 256-dimensional vectors to obtain the final embedding. This normalization greatly simplifies the calculation of cosine-distance in downstream applications, allowing for a straightforward dot product operation.

**3.4.2 Unified Pin and Product Encoder.** In our model, we utilize a single unified encoder for both pins and products (depicted in Figure 3), and this encoder is jointly trained with the query embeddings. Designed to process both textual features and continuous features, it plays a crucial role in learning the respective embeddings of pins and products. In cases where certain features are defined for one entity but not the other, we substitute them with zero, ensuring a consistent data input.

As detailed in section 3.5, we utilize in-batch negatives to train our model. Prior research [9, 15, 16, 29] has empirically demonstrated that larger batches with a substantial number of negatives help in learning better representations. Therefore, to accommodate a larger batch size in the GPU memory, we employ a simple pin encoder model. The following encoder design has been determined through numerous ablation studies. These studies have allowed us to select the most effective configuration for each of the components, while still considering the importance of both training and serving efficiencies.

The encoder uses three distinct tokenizers to process the textual features associated with a pin [1, 13, 23]. These include (i) a word unigram tokenizer that uses a vocabulary encompassing the 200k most frequent word unigrams, (ii) a word bigram tokenizer that makes use of a vocabulary comprising the 1M most frequent word bigrams, and (iii) a character trigram tokenizer that utilizes a vocabulary of 64k character trigrams. The tokens are mapped to their respective IDs in the vocabulary  $\mathcal{V}$  which constitute all three



**Figure 3: Schematic of the unified encoder model for pins and products, illustrating the use of three different tokenizers, a hash embedding table, and an MLP layer for combining text embeddings with other continuous features.**

tokenizers. Any token that falls out of this combined vocabulary gets discarded. The use of these combined tokenizers effectively helps in capturing the semantics of various texts associated with a pin/product.

For token embedding learning, we use a 2-hash hash embedding table of size 100,000 [1, 35]. Each identified token’s ID  $i$  is hashed into two places within the embedding table using hash functions  $h_1(i)$  and  $h_2(i)$ . The ultimate embedding of a token with ID  $i$  is a weighted interpolation of the two locations:  $W_1 h_1(i) + W_2 h_2(i)$ , where  $W_1$  and  $W_2$  are learned weight vectors of size  $|\mathcal{V}|$  each.

The sum of all token embeddings and the embedding features are concatenated and fed into a 3-layer MLP, with layer sizes of 1024, 1024, 256. Following this, the output of the MLP layer undergoes L2-normalization just like the query embedding.

**3.4.3 Compatibility Encoders.** In our model, we employ two discrete compatibility encoders individually dedicated to pins and products. These encoders leverage the pre-existing pin and product embeddings, represented by PinSage for pins and ItemSage for products. This allows the model to adeptly learn query embeddings that align effectively with PinSage and ItemSage embeddings.

### 3.5 Multi-Task Sampled Softmax Loss

Taking inspiration from Itemsage [1], the problem of learning query and entity embeddings is treated as an extreme classification problem, with the aim of predicting entities relevant to a given query [7]. We employ the sampled softmax loss with logQ correction [36] to train our model.

We use multitasking to jointly train entity embeddings and train the query embeddings to be compatible with existing entity embeddings.

Formally, we define a task  $T \in \mathcal{T}$  as a tuple of a dataset of query-entity pairs ( $\mathcal{D} = \{(x, y)_i\}$ ) and an entity encoder  $\mathcal{E}$ .

$$T \triangleq \{\mathcal{D}, \mathcal{E}\}.$$

For a batch of data,  $\mathcal{B} = \{(x, y)_i\} \subset \mathcal{D}$ , for task  $T \in \mathcal{T}$ , the aim is to learn query embedding  $q_{x_i}$  and entity embedding  $p_{y_i} = \mathcal{E}(y_i)$  such that the cosine similarity of the embeddings  $q_{x_i} \cdot p_{y_i}$  is maximized.

<sup>2</sup><https://huggingface.co/distilbert-base-multilingual-cased>

This is achieved by minimizing the softmax loss:

$$L_T = -\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \log \frac{\exp(q_{x_i} \cdot p_{y_i})}{\sum_{y \in C} \exp(q_{x_i} \cdot p_y)}, \quad (1)$$

where  $C$  is the catalog of all entities of the same type as  $y_i$ . To ensure problem tractability, the normalization term in the denominator is approximated using a sample of the catalog  $C$ . We use (i) positives in the batch,  $BN = \{y_i | (x_i, y_i) \in \mathcal{B}\}$ , and (ii) a random sample of the catalog,  $C'$ . To rectify any bias that might have been introduced through sampling, we utilize the logQ correction technique. This method operates by deducting the sampling probability of the negative, represented as  $\log Q(y|x_i)$ , from the existing logits. This is crucial to ensure that popular entities aren't disproportionately penalized.

$$L_T = L_T^{Sbn} + L_T^{Srn} \quad (2)$$

$$L_T^{Sbn} = -\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \log \frac{\exp(q_{x_i} \cdot p_{y_i} - \log Q(y_i|x_i))}{\sum_{z \in BN} \exp(q_{x_i} \cdot p_z - \log Q(z|x_i))} \quad (3)$$

$$L_T^{Srn} = -\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \log \frac{\exp(q_{x_i} \cdot p_{y_i} - \log Q(y_i|x_i))}{\sum_{y \in C'} \exp(q_{x_i} \cdot p_y - \log Q(y|x_i))} \quad (4)$$

$$= -\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \log \frac{\exp(q_{x_i} \cdot p_{y_i} - \log Q(y_i|x_i))}{\sum_{y \in C'} \exp(q_{x_i} \cdot p_y - \log Q_n(y))}, \quad (5)$$

since  $y$  is sampled independently

The total loss is defined as the sum of all individual task losses,

$$L = \sum_{T \in \mathcal{T}} L_T. \quad (6)$$

We mix together different tasks together in one batch and control the influence of each task on the model through this composition. To increase training efficiency, we share the pairs in the batch across all tasks with the same dataset.

### 3.6 Model Serving

OMNISEARCHSAGE query embeddings are integral to numerous applications in the search stack, which necessitates us to maintain a strict latency budget. For real-time inference with minimized latency, our query encoder is served on GPUs by our in-house C++-based machine learning model server, the Scorpion Model Server (SMS). Factoring in that query distribution complies with Zipf's law, we have instituted a cache-based system to curb costs and shorten response times. The query embedding server first verifies if a query is cached before resorting to the query inference server should it be absent from the cache. After testing various Cache Time-To-Live (TTL) periods, a TTL of 30 days was established as optimal. The system is equipped for handling 300k requests per second, maintaining a median (p50) latency of just 3ms, and 90 percentile (p90) latency of 20ms. The implementation of this cache-based system efficiently reduces the load on the inference server to approximately 500 QPS, leading to substantial cost and latency reductions.

The pin and product embeddings are derived offline on a daily basis through batch inference on GPUs and are subsequently published to our signal store for consumption.

Pair	Source	Actions	Size
Query-Pin	Query Logs	repin, longclick	1.5B
Query-Product	Query Logs	repin, longclick	136M
Query-Product	Offsite logs	add-to-cart, checkout	2.5M
Query-Query	Query Logs	click	195M

**Table 1: Summary of the different training datasets.**

## 4 EXPERIMENTS

### 4.1 Dataset

Our dataset is primarily constructed by extracting unique query-entity pairs from one year of search query logs. We consider various forms of engagement on the platform when extracting these pairs, including 'saves' (when a user saves a pin to a board) and 'long clicks' (instances where users browse the linked page for more than 10 seconds before returning to Pinterest). For products, we enrich our dataset by incorporating offsite actions as well. Thus, we also include anonymized pairs tied to significant actions like 'add to cart' and 'checkout'. A common challenge in recommendation systems is the popularity bias, where certain pins are overrepresented due to their high appeal. To counteract this bias, we impose a limit on the number of times the same pin can be paired. This limit is capped at 50 pairs for pins and is extended to 200 pairs for products (since products have lower volume and engagement). By adopting this strategy, we ensure our dataset is robust and truly representative of the user's activity on the platform.

Our model training is further extended to encompass query-query pairs. On Pinterest, users are presented with similar query suggestions, and engagements with these recommendations are recorded in the search logs. We leverage these records, extracting such pairs from an entire year's logs, thus enriching our training dataset.

A detailed breakdown of the positive labels in the dataset is provided in Table 1.

### 4.2 Offline Evaluation Metrics

Our evaluation of the model encompasses both user engagement data and human-labeled relevance data.

Relevance gets measured using human-labeled pairs of queries and pins, sampled from production traffic from four distinct countries: US, UK, France, and Germany. This strategy serves to assess the model's performance in handling multiple languages and cultural contexts.

Evaluation of user engagement considers a selected 7-day period. We ensure no data leakage—possible due to the inclusion of engagement features such as engaged queries—by maintaining a 15-day separation between the end of the training dataset and the beginning of the evaluation phase. We sample 80k pairs from the defined evaluation duration to represent repins and long clicks for both pins and products. Another 80k pairs, corresponding to clicks for queries, are also included for comprehensive performance evaluation.

The primary metric we used for evaluation is named 'Recall@10'. This metric denotes the likelihood of the occurrence of the engaged entity within the top 10 entities when these entities are sorted in descending order based on their similarity to the query.

Metric	SearchSage	OMNISEARCHSAGE	Gain
Pin			
Save	0.39	0.65	+67%
Long-Click	0.45	0.73	+62%
Relevance (US)	0.25	0.45	+80%
Relevance (UK)	0.29	0.51	+76%
Relevance (FR)	0.23	0.43	+87%
Relevance (DE)	0.28	0.46	+64%
Product			
Save	0.57	0.73	+28%
Long-Click	0.58	0.73	+26%
Query			
Click	0.54	0.78	+44%

**Table 2: Comparative analysis of OMNISEARCHSAGE and the baseline SearchSage across various tasks - Pin, Product, and Query.**

Consider a dataset  $D = (q_i, e_i)_{i=1}^n$ , where each  $(q_i, e_i)$  denotes a query-engaged entity pair, and also consider a random corpus  $C$  with  $m$  entities. The Recall@10 metric can then be defined as the average over all queries of the indicator function 1, where 1 equals 1 if the engaged entity  $e_i$  is amongst the top 10 entities in  $C$  when ranked by their dot product with the query  $q_i$ .

$$\text{Recall@10} = \frac{1}{|D|} \sum_{i=1}^{|D|} 1\left(\sum_{y \in C} x_i \cdot y > x_i \cdot y_i > 10\right)$$

For every pin, query, and product, we employ a uniformly distributed random sample of  $m = 1.5M$  entities from our corpus.

### 4.3 Offline Results

In this section, we provide a comprehensive comparison between our proposed model, OMNISEARCHSAGE, and the existing baselines, which helps showcase its performance enhancements. Subsequently, we undertake an in-depth exploration of key influential aspects such as the significance of text enrichments, the pros and cons of adopting multitasking approaches, and the operational efficacy of compatibility encoders in the context of our model.

**4.3.1 Comparison with Baselines.** In this study, the existing version of SearchSage [27] serves as our comparison baseline. It operates using fixed PinSage and ItemSage embeddings for pins and products, respectively. For OMNISEARCHSAGE, we utilize the query encoder to derive query embeddings and the unified pin and product encoder to generate pin and product embeddings.

In Table 2, comparisons are drawn between OMNISEARCHSAGE and SearchSage, with both models being trained and evaluated on the same dataset. It is important to highlight that the baseline model, SearchSage, does not involve query-query pairs for training purposes.

On the pin dataset, OMNISEARCHSAGE shows a significant gain, between 60% and 90%, over SearchSage across all metrics. Recall is relatively consistent across different countries, reflecting the multilingual robustness of OMNISEARCHSAGE.

Analysis of the product dataset reveals that OMNISEARCHSAGE outperforms the baseline model by about 27% in predicting product

	save	long-click	relevance
No captions	0.51	0.60	0.36
With captions	0.66	0.76	0.36
Improvement	+30.43%	+25.58%	0%

**Table 3: Comparative assessment displaying the influence of Synthetic GenAI Captions on pins lacking titles and descriptions.**

engagement. This increment is less prominent as compared to the pins dataset, mainly because ItemSage, upon which this comparison is based, has already undergone training on search tasks. Nevertheless, the observed improvement shows the positive impact of incorporating new features as well as the benefit of multi-tasking.

Interestingly, SearchSage is able to predict related query clicks substantially better than random despite not being trained on this task. However, when we directly optimize for this objective in OMNISEARCHSAGE, we see a substantial +44% improvement. We show this improvement can be attributed to both training on related queries, and multi-task learning in Section 4.3.3.

**4.3.2 Importance of content enrichment.** In this section, we delve into an analysis of the importance of various text enhancements described in Section 3.2. To maintain brevity, the evaluation focuses solely on the metrics related to the query-pin task.

Our first direction of investigation centers around the impact of integrating synthetic captions for pins that lack both a title and description. For this purpose, we extracted pairs from the evaluation dataset in which the engaged pin was missing a title or a description. This resulted in a narrowed evaluation dataset of 24k pairs. The model’s performance, initially based on solely continuous features and native text, was then compared to a model additionally enriched with captions.

Table 3 presents the results of this comparison. When synthetic captions were added, both ‘save’ and ‘long-click’ metrics saw substantial improvements — approximately +30% and +26% respectively. However, the relevance metric remained unchanged.

This suggests that adding synthetic captions can significantly enhance the model’s performance for certain metrics when representing pins that lack a title and description.

Table 4 illustrates the impact of adding different text enrichments on the model’s performance. Each percentage increase is relative to the previous row, displaying the additional improvement from each additional feature.

Our baseline model utilizes only continuous features for training and its performance values are reflected in the first row. Upon adding ‘Title’, ‘Description’, and ‘Synthetic GenAI Captions’ to the baseline model, we notice a robust improvement across all metrics.

	save	long-click	relevance
Continuous Features Only	0.43	0.53	0.30
Adding Title, Description and Synthetic GenAI Captions	0.52 (+21%)	0.63 (+19%)	0.39 (+30%)
Adding Board Titles	0.61 (+17%)	0.68 (+8%)	0.44 (+13%)
Adding Engaged Queries	0.65 (+7%)	0.73 (+7%)	0.46 (+5%)

**Table 4: Impact of adding different text enrichments on the model’s performance. Each percentage increase is relative to the previous row, displaying the additional improvement from each additional feature.**



Dataset		Pin Only	Product only	Query Only	OmniSearchSage
pin	save	0.68	-	-	0.65
	long-click	0.75	-	-	0.73
	avg relevance	0.45	-	-	0.46
product	save	-	0.73	-	0.73
	long-click	-	0.73	-	0.73
query	click	-	-	0.73	0.78

**Table 5: Comparative analysis illustrating the contrasts between our unified multi-task model and models trained individually for each task - pin, product, and query.**

There is a 20% improvement in the engagement datasets, while the relevance metric improves by a notable 30%, demonstrating the substantial impact of these text features.

The model enhancement continues with adding board titles to the feature set, leading to a further increase of 8 – 15% in different metrics. This affirms the relevance of board titles in improving predictive accuracy.

Finally, we incorporated engaged queries feature into the model, resulting in a consistent, albeit smaller growth across all three metrics. Although the incremental relative gain appears smaller, it still constitutes a significant improvement when compared to the baseline model. In summary, each text enrichment feature contributes significantly to improving model performance as seen by the increment in metrics compared to their immediate preceding state.

**4.3.3 Effect of multi-tasking.** In Table 5, we present a comparative analysis between models trained independently for each task (pin, product, and query) and our consolidated multitask model. For this comparison, both the independent and multitask models were trained under equivalent conditions - with matching batch sizes, computational power, and iterations. The datasets used for both training and evaluation were also identical, with the sole difference that the individual models were trained on their respective subset of pairs from the dataset. This systematic approach ensures the fair and accurate assessment of the performance of the multitask model in relation to the independent task models.

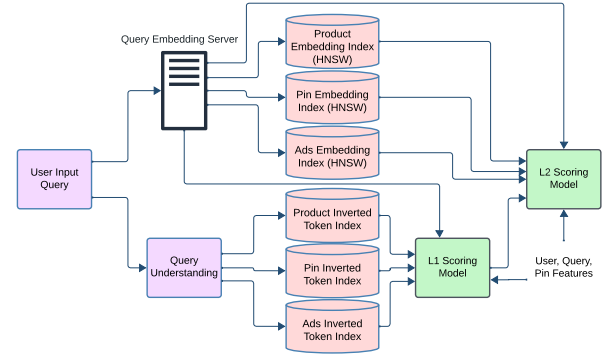
On the pin task, we see slight degradation in quality from multi-task learning, but, on product and query tasks, results are neutral to positive. This aligns with general notions about multi-task learning: low-data tasks are unlikely to see regressions from multi-task learning, while the pin task using 1.5B pairs sees a very slight drop in performance. Despite this drop, the simplification benefits of multi-task learning outweigh the metric loss.

**4.3.4 Effect of compatibility encoders.** We examine the influence of incorporating compatibility encoders on the effectiveness of the learned pin/product embeddings. We train a model that comprises only the query and unified pin and product encoder. Subsequently, this model is compared with another model that fully incorporates all the encoders. Interestingly, there is almost no noticeable degradation in the metrics of the learned encoder, thereby essentially achieving seamless compatibility of the query embedding with pre-existing embeddings at no substantial cost.

Furthermore, as demonstrated in Table 6, the performance of the compatibility encoders in the OMNISEARCHSAGE model is either on par with or surpasses that of the SearchSage model, which is trained utilising only compatibility encoders.

Dataset		SearchSage	OMNISEARCHSAGE
pin	save	0.39	0.39
	long-click	0.45	0.43
	avg relevance	0.26	0.26
product	save	0.57	0.57
	long-click	0.58	0.57

**Table 6: Comparison of co-trained compatibility encoders with independently trained compatibility encoders.**



**Figure 4: A simplified depiction of the search retrieval and ranking stack at Pinterest highlighting the integration points for OMNISEARCHSAGE embeddings.**

## 5 APPLICATIONS IN PINTEREST SEARCH

OMNISEARCHSAGE embeddings find wide applications throughout the Pinterest search stack, primarily in retrieval and ranking tasks. Figure 4 presents a simplified depiction of the search retrieval and ranking stack at Pinterest and highlights the integration points for OMNISEARCHSAGE embeddings.

These embeddings are employed to power the retrieval of pins and products using HNSW [22]. They are also instrumental in the L1 scoring model, where they enhance the efficiency of token-based retrieval sources. Moreover, OMNISEARCHSAGE embeddings serve as one of the most critical features in the L2 scoring and relevance models.

In this section, we delineate the results derived from the A/B tests we conducted. In these tests, production SearchSage embeddings were replaced with OMNISEARCHSAGE embeddings, resulting in boosted performance in both organic and promoted content (Ads) in search. Additionally, we provide results from a human relevance assessment conducted on actual production-sampled traffic. This evaluation further confirms the improved performance derived from the utilization of OMNISEARCHSAGE embeddings. Finally, we demonstrate how employing query embeddings also enhances performance in other tasks, such as classification, particularly in situations where data availability is limited. This highlights the ability of the OMNISEARCHSAGE model to generalize to tasks different from its original training objectives.

### 5.1 Human Relevance Evaluation

To understand advantages of OMNISEARCHSAGE, we enlisted human evaluators to assess the relevance of candidates retrieved via two



**Figure 5: Comparative display of pins retrieved in response to the query 'antique copper bathroom sink' from the token-based system and the OMNISEARCHSAGE-based system. Pins deemed relevant are outlined in green, while those considered irrelevant are encircled in red.**

methods: OMNISEARCHSAGE embeddings-based pin retrieval and token-based pin retrieval.

For this evaluation, we selected a set of 300 queries, deliberately stratified across both head and tail queries. The top 8 candidate pins were then retrieved from each system using these queries, and human evaluators determined the relevance of the pins to the corresponding query. Every query-pin pair received three judgments, with an inter-annotator agreement rate of 0.89. Evaluation results revealed a noticeable improvement with OMNISEARCHSAGE, showing a 10% increase in relevance compared to the token-based system.

Figure 5 offers a distinct comparison of retrieved pins for the query 'antique copper bathroom sink' between the candidates retrieved by the token-based system and the OMNISEARCHSAGE-based system. The token-based retrieval system often fetches pins related to only part of the query and fails to fetch consistently relevant results. In striking contrast, nearly all pins retrieved by the OMNISEARCHSAGE-based system are highly relevant to the specified query, underlining the efficacy of the OMNISEARCHSAGE model in understanding the query and aligning similar pins and queries in the same space together.

## 5.2 Organic Search

In this section, we outline the results of the A/B testing conducted to substitute the existing production SearchSage query and entity embeddings with OMNISEARCHSAGE embeddings for organic content within Pinterest search. Within the context of search experiments at Pinterest, our attention is largely concentrated on two key metrics: the search fulfillment rate and relevance. The search fulfillment rate is defined as the proportion of searches that result in a user engagement action of significance. Relevance is calculated as the weighted average relevance of the top eight pins for each query, assessed across different query segments. This is measured through human evaluation.

The impact on these two metrics, from replacing SearchSage with OMNISEARCHSAGE, is presented in Table 7. The table provides data drawn from experiments for three distinct use-cases: (i) retrieval of pins and products, (ii) L1 scoring model, and (iii) L2 scoring model and relevance model.

	Search Fulfillment Rate	Relevance
Pin and Product Retrieval	+4.1%	+0.5%
L1 Scoring	+0.5%	+0.0%
L2 Scoring and Relevance Model	+2.8%	+3.0%

**Table 7: Online A/B experiment results of OMNISEARCHSAGE in Organic Search.**

	gCTR
Product Ads Retrieval	+5.27%
Ads Search Engagement Model	+2.96%
Ads Search Relevance Model	+1.55%

**Table 8: Online A/B experiment results of OMNISEARCHSAGE for Ads in Search.**

## 5.3 Ads in Search

The OMNISEARCHSAGE embeddings have also successfully replaced the SearchSage embeddings in various applications within Ads on Search surface. We present the results of three use cases: search engagement model, search relevance model, and product ads retrieval.

Uniformly, we noted substantial improvements in engagement and relevance within Ads across all use cases. These increments, specifically in the long clickthrough rate (gCTR), are outlined in Table 8. Furthermore, OMNISEARCHSAGE led to a noteworthy 4.95% increase in Ads relevance within the Search Ads relevance model. These gains highlight the positive impact of transitioning to OMNISEARCHSAGE embeddings for Ads on Search.

## 5.4 Classification

One of the primary advantages of developing robust query representation such as OMNISEARCHSAGE is its utility in powering downstream applications, particularly when there is a lack of labels for learning large models. One example of this at Pinterest is interest classification, where we classify queries into a hierarchical taxonomy. Using OMNISEARCHSAGE query embeddings for query representation, we were able to increase performance when compared to the baseline FastText [3] model. Precision increased by 30% on average across levels, with the larger gains coming from more granular levels.

## 6 CONCLUSION

In this work, we presented OMNISEARCHSAGE, an end-to-end optimized set of query, pin, and product embeddings for Pinterest search, which have shown value across many applications.

In contrast to other work focused on learning embeddings for search, we demonstrate the value of unified query, pin, and product embeddings as both candidate generators and features in Pinterest search. We show a great improvement over previous solutions at Pinterest can be attributed to rich document text representations, which improved offline evaluation metrics by > 50%. We also describe practical decisions enabling serving and adoption, including compatibility encoders, multi-task learning, and long-TTL caching.

Lastly, we summarize results from online A/B experiments across organic and ads applications, which have directly led to cumulative gains of +7.4% fulfillment rate on searches, and +3.5% relevance.



## REFERENCES

- [1] Paul Baltescu, Haoyu Chen, Nikil Pancha, Andrew Zhai, Jure Leskovec, and Charles Rosenberg. 2022. ItemSage: Learning Product Embeddings for Shopping Recommendations at Pinterest. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, 2703–2711.
- [2] Josh Beal, Hao-Yu Wu, Dong Huk Park, Andrew Zhai, and Dmitry Kislyuk. 2022. Billion-Scale Pretraining with Vision Transformers for Multi-task Visual Representations. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 564–573.
- [3] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics* 5 (2017), 135–146.
- [4] Daniel Cer, Yinfei Yang, Sheng-yi Kong, Nan Hua, Nicole Limtiaco, Rhomni St. John, Noah Constant, Mario Guajardo-Cespedes, Steve Yuan, Chris Tar, Brian Strope, and Ray Kurzweil. 2018. Universal Sentence Encoder for English. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Brussels, Belgium, 169–174.
- [5] Nurendra Choudhary, Nikhil Rao, Karthik Subbian, and Chandan K. Reddy. 2022. Graph-Based Multilingual Language Model: Leveraging Product Relations for Search Relevance. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, 2789–2799.
- [6] Alexis Conneau, Douwe Kiela, Holger Schwenk, Loïc Barrault, and Antoine Bordes. 2017. Supervised Learning of Universal Sentence Representations from Natural Language Inference Data. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, Copenhagen, Denmark, 670–680.
- [7] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep Neural Networks for YouTube Recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 191–198.
- [8] Miles Efron, Peter Organisciak, and Katrina Fenlon. 2012. Improving Retrieval of Short Texts through Document Expansion. In *Proceedings of the 35th International ACM SIGIR Conference on Research and Development in Information Retrieval*. Association for Computing Machinery, New York, NY, USA, 911–920.
- [9] John Giorgi, Osvald Nitski, Bo Wang, and Gary Bader. 2021. DeCLUTR: Deep Contrastive Learning for Unsupervised Textual Representations. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Association for Computational Linguistics, Online, 879–895.
- [10] William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *Proceedings of the 31st International Conference on Neural Information Processing Systems*. Curran Associates Inc., Red Hook, NY, USA, 1025–1035.
- [11] Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padmanabhan, Giuseppe Ottaviano, and Linjun Yang. 2020. Embedding-Based Retrieval in Facebook Search. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, New York, NY, USA, 2553–2561.
- [12] 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. In *Proceedings of the 22nd ACM International Conference on Information & Knowledge Management*. Association for Computing Machinery, New York, NY, USA, 2333–2338.
- [13] 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. In *Proceedings of the 22nd ACM International Conference on Information & Knowledge Management*. Association for Computing Machinery, New York, NY, USA, 2333–2338.
- [14] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-Scale Similarity Search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.
- [15] Prannay Khosla, Piotr Teterwak, Chen Wang, Aaron Sarna, Yonglong Tian, Phillip Isola, Aaron Maschinot, Ce Liu, and Dilip Krishnan. 2020. Supervised Contrastive Learning. In *Advances in Neural Information Processing Systems*, H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin (Eds.), Vol. 33. Curran Associates, Inc., 18661–18673.
- [16] Kenton Lee, Ming-Wei Chang, and Kristina Toutanova. 2019. Latent Retrieval for Weakly Supervised Open Domain Question Answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, Florence, Italy, 6086–6096.
- [17] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. 2023. BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation. In *Proceedings of the 40th International Conference on Machine Learning*. JMLR.org, New York, NY, USA, Article 814, 13 pages.
- [18] Sen Li, Fuyu Lv, Taiwei Jin, Guli Lin, Keping Yang, Xiaoyi Zeng, Xiao-Ming Wu, and Qianli Ma. 2021. Embedding-based Product Retrieval in Taobao Search. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, New York, NY, USA, 3181–3189.
- [19] Yiqun Liu, Kaushik Rangadurai, Yunzhong He, Siddarth Malreddy, Xunlong Gui, Xiaoyi Liu, and Fedor Borisjuk. 2021. Que2Search: Fast and Accurate Query and Document Understanding for Search at Facebook. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, New York, NY, USA, 3376–3384.
- [20] Wenhao Lu, Jian Jiao, and Ruofei Zhang. 2020. TwinBERT: Distilling Knowledge to Twin-Structured Compressed BERT Models for Large-Scale Retrieval. In *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*. Association for Computing Machinery, New York, NY, USA, 2645–2652.
- [21] Alessandro Magnani, Feng Liu, Suthee Chaidaroon, Sachin Yadav, Praveen Reddy Suram, Ajit Puthenpuussery, Sijie Chen, Min Xie, Anirudh Kashi, Tony Lee, and Ciya Liao. 2022. Semantic Retrieval at Walmart. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, 3495–3503.
- [22] Yu A. Malkov and Dmitry A. Yashunin. 2018. Efficient and Robust Approximate Nearest Neighbor Search Using Hierarchical Navigable Small World Graphs. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 42, 4 (2018), 824–836.
- [23] Priyanka Nigam, Yiwei Song, Vijai Mohan, Vihan Lakshman, Weitian (Allen) Ding, Ankit Shingavi, Choon Hui Teo, Hao Gu, and Bing Yin. 2019. Semantic Product Search. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, New York, NY, USA, 2876–2885.
- [24] Xichuan Niu, Bofang Li, Chenliang Li, Rong Xiao, Haochuan Sun, Hongbo Deng, and Zhenzhong Chen. 2020. A Dual Heterogeneous Graph Attention Network to Improve Long-Tail Performance for Shop Search in E-Commerce. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, New York, NY, USA, 3405–3415.
- [25] Rodrigo Nogueira, Jimmy Lin, and AI Epistemic. 2019. From doc2query to docTTTTTquery. *Online preprint* 6 (2019), 2.
- [26] Rodrigo Frassetto Nogueira, Wei Yang, Jimmy Lin, and Kyunghyun Cho. 2019. Document Expansion by Query Prediction. *CoRR abs/1904.08375* (2019).
- [27] Nikil Pancha, Andrew Zhai, Chuck Rosenberg, and Jure Leskovec. 2021. SearchSage: Learning Search Query Representations at Pinterest. <https://medium.com/pinterest-engineering/searchsage-learning-search-query-representations-at-pinterest-654f2bb887fc>
- [28] Jeremy Pickens, Matthew Cooper, and Gene Golovchinsky. 2010. Reverted Indexing for Feedback and Expansion. In *Proceedings of the 19th ACM International Conference on Information and Knowledge Management*. Association for Computing Machinery, New York, NY, USA, 1049–1058.
- [29] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. Learning Transferable Visual Models From Natural Language Supervision. In *Proceedings of the 38th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 139)*. PMLR, 8748–8763.
- [30] Nils Reimers and Iryna Gurevych. 2019. Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan (Eds.). Association for Computational Linguistics, Hong Kong, China, 3982–3992.
- [31] Sebastian Ruder, Matthew E. Peters, Swabha Swayamdipta, and Thomas Wolf. 2019. Transfer Learning in Natural Language Processing. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Tutorials*. Association for Computational Linguistics, Minneapolis, Minnesota, 15–18.
- [32] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. DistilBERT, a Distilled Version of BERT: Smaller, Faster, Cheaper, and Lighter. *ArXiv abs/1910.01108* (2019).
- [33] Hongyan Tang, Junning Liu, Ming Zhao, and Xudong Gong. 2020. Progressive Layered Extraction (PLE): A Novel Multi-Task Learning (MTL) Model for Personalized Recommendations. In *Proceedings of the 14th ACM Conference on Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 269–278.
- [34] Tao Tao, Xuanhui Wang, Qiaozhu Mei, and ChengXiang Zhai. 2006. Language Model Information Retrieval with Document Expansion. In *Proceedings of the Human Language Technology Conference of the NAACL, Main Conference*. Association for Computational Linguistics, New York City, USA, 407–414.
- [35] Dan Tito Svenstrup, Jonas Hansen, and Ole Winther. 2017. Hash Embeddings for Efficient Word Representations. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. Curran Associates, Inc.

- [36] Xinyang Yi, Ji Yang, Lichan Hong, Derek Zhiyuan Cheng, Lukasz Heldt, Aditee Kumthekar, Zhe Zhao, Li Wei, and Ed Chi. 2019. Sampling-Bias-Corrected Neural Modeling for Large Corpus Item Recommendations. In *Proceedings of the 13th ACM Conference on Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 269–277.
- [37] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure Leskovec. 2018. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, New York, NY, USA, 974–983.
- [38] Licheng Yu, Jun Chen, Animesh Sinha, Mengjiao Wang, Yu Chen, Tamara L. Berg, and Ning Zhang. 2022. CommerceMM: Large-Scale Commerce MultiModal Representation Learning with Omni Retrieval. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, 4433–4442.
- [39] Chuxu Zhang, Dongjin Song, Chao Huang, Ananthram Swami, and Nitesh V. Chawla. 2019. Heterogeneous Graph Neural Network. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, New York, NY, USA, 793–803.
- [40] Jianjin Zhang, Zheng Liu, Weihao Han, Shitao Xiao, Ruicheng Zheng, Yingxia Shao, Hao Sun, Hanqing Zhu, Premkumar Srinivasan, Weiwei Deng, Qi Zhang, and Xing Xie. 2022. Uni-Retriever: Towards Learning the Unified Embedding Based Retriever in Bing Sponsored Search. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, 4493–4501.
- [41] Zhe Zhao, Lichan Hong, Li Wei, Jilin Chen, Aniruddh Nath, Shawn Andrews, Aditee Kumthekar, Maheswaran Sathiamoorthy, Xinyang Yi, and Ed Chi. 2019. Recommending What Video to Watch Next: A Multitask Ranking System. In *Proceedings of the 13th ACM Conference on Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 43–51.