# Pfeed: Generating near real-time personalized feeds using precomputed embedding similarities

Binyam Gebre
Bol
The Netherlands
bgebre@bol.com

Karoliina Ranta*
Booking.com
The Netherlands
karoliina.ranta@booking.com

Stef van den Elzen
Eindhoven University of Technology
The Netherlands
s.j.v.d.elzen@tue.nl

Ernst Kuiper
Bol
The Netherlands
ekuiper@bol.com

Thijs Baars*
Last Mile Solutions
The Netherlands
thijs.baars@lastmilesolutions.com

Tom Heskes
Radboud University Nijmegen
The Netherlands
tom.heskes@ru.nl

## ABSTRACT

In personalized recommender systems, embeddings are often used to encode customer actions and items, and retrieval is then performed in the embedding space using approximate nearest neighbor search. However, this approach can lead to two challenges: 1) user embeddings can restrict the diversity of interests captured and 2) the need to keep them up-to-date requires an expensive, real-time infrastructure. In this paper, we propose a method that overcomes these challenges in a practical, industrial setting. The method dynamically updates customer profiles and composes a feed every two minutes, employing precomputed embeddings and their respective similarities. We tested and deployed this method to personalize promotional items at Bol, one of the largest e-commerce platforms of the Netherlands and Belgium. The method enhanced customer engagement and experience, leading to a significant 4.9% uplift in conversions.

## KEYWORDS

deep learning, joint embeddings, dual encoders, contrastive learning, personalization, e-commerce

## 1 INTRODUCTION

Bol, like many other e-commerce platforms, faces the challenge of providing customers with an easy and efficient way to navigate their vast catalog and find products that match their customers' interests. The traditional approach of relying on customer controlled text-based search engines or browsing through categories is often limited and cumbersome, particularly during the customer's discovery phase. To overcome these limitations and enhance customers' overall discovery experience, Bol has launched personalized feeds called *Top deals for you*, *Top picks for you*, and *New for you*.

These personalized feed systems utilize a combination of the customer's historical and recent behavior to display the best recommendations on the customer's home page across both app and desktop platforms. In this paper, we present the methodology behind these feeds. We begin by presenting the challenges inherent to creating personalized feed systems. Subsequently, we delve into the prevailing industry approach (related work) that tackles these challenges, concluding with the presentation of our proposed solution and the evaluation outcomes.

---

*Both authors contributed to this work while working at Bol.

### 1.1 Four Challenges in Personalized Feed Systems

Personalized feed systems can be viewed as search engines, where customers are the search queries and items in the catalog are the search results. In this view, there are four challenges that need to be overcome to provide customers with a personalized set of items that align with their interests and preferences: customer, item, candidate retrieval and ranking challenges.

*1.1.1 Customer representation challenge.* Customers show complex behaviors while shopping on e-commerce sites before making a purchase, e.g., searching for items, viewing items, reading reviews, and making item comparisons. The challenge is distilling these interactions into a concise customer representation. In addition to their dynamic interactions, the representation may also need to incorporate static attributes of customers, such as customer ID, gender, and clothing-size.

*1.1.2 Item representation challenge.* Items have rich structured information such as item ID, title, description, specifications, and other metadata. Items also have historical customer interactions: views, clicks, customer ratings, reviews, etc. The item representation challenge is identifying the most relevant data for representing various items, a task complicated by two factors. The first is the diversity of item attributes. For instance, *author* and *title* are key attributes for books, whereas *size* and *gender* are more critical for fashion items. The second factor is the cold-start problem associated with new products; these have no historical interactions.

*1.1.3 Candidate retrieval challenge.* Candidate retrieval entails determining which items best match a given customer's preferences. Here, the challenges are of two varieties: 1) training customer and item representations in the same embedding space and 2) the inference challenge, which aims to efficiently retrieve the best matches from a corpus containing millions to billions of items.

*1.1.4 Ranking challenge.* The candidate retrieval stage is followed by a ranking stage, where the retrieved candidates are re-ranked using a more complex model and more complex features of both the retrieved candidates and queries. The goal of this stage is to select and rank the top K items per customer (for example, the top 100 items) using learning-to-rank algorithms [3, 6, 23, 32].

In this paper, we focus on addressing the first three challenges: the *customer representation*, *item representation*, and *candidate retrieval* challenges.

## 1.2 Our Contributions

The most dominant approach to tackling the aforementioned challenges relies on a user-item framework (see figure 1). Two neural networks, called dual encoders, are each trained to generate embeddings for user and item data [4, 18, 19, 24, 27]. The user embedding model receives input in the form of a sequence or bag of interactions on items, along with context and user data [24]. On the other hand, the item embedding model utilizes various item metadata types including item IDs [2, 4, 19] or output embeddings from pre-trained models [18, 19, 28].

However, despite its widespread use, the user encoding model in this framework has two significant drawbacks: the *single vector representation bottleneck* and the *high infrastructure and maintenance costs*.

*1.2.1 Single vector representation bottleneck.* Using a single vector to represent users introduces challenges due to the diversity and complexity of their interests, compromising both the capacity to accurately represent users and the interpretability of the representation by obscuring which interests are represented and which are not. While attempts to use multiple embeddings have been made to overcome these limitations, the exact number of vectors needed and the method for obtaining them remain topics of research [16, 18].

*1.2.2 High infrastructure and maintenance costs.* Generating and maintaining up-to-date user embeddings requires substantial investment in terms of infrastructure and maintenance (see, for example, the SOAP platform from Meta [31]). Each new user action necessitates executing the user encoder to generate fresh embeddings and recommendations. Furthermore, the user encoder must be large in order to effectively model a sequence of interactions, leading to expensive training and inference requirements.
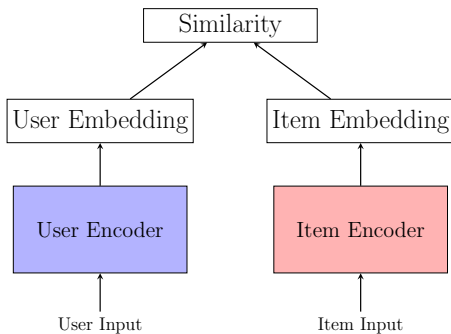


**Figure 1: User-to-item framework: Single vectors from the user encoder limit representation and interpretability. Keeping them fresh demands high-maintenance infrastructure.**

Our approach overcomes these drawbacks by modelling item-to-item relationships, as illustrated in figure 2. Here, the first item represents the query context (an item that has been bought or viewed), while the second item is the target (the item that is subsequently bought). We utilize dual encoders to effectively capture
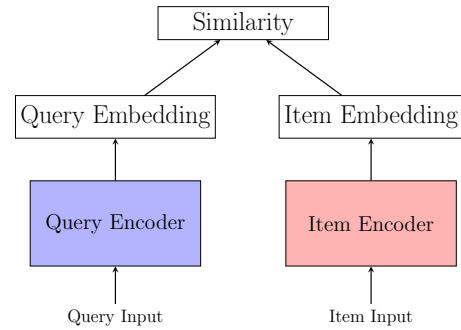


**Figure 2: Query-to-item framework: Query embeddings and their similarities are precomputed. Users are represented by a dynamic set of queries that can be updated as needed.**

relationships between viewed and bought items, as well as between items bought together. Specifically, our contributions include:

(1) We demonstrate how a transformer-based two-tower architecture, also known as dual encoders, can be utilized to generate multiple embeddings per item in one model run. Generating multiple embeddings is effective for capturing the various roles of items, and generating them with one model run provides inference efficiency.

(2) We show how we represent customers with multiple queries, where each query corresponds to a product that the customer has interacted with, either through a view or a buy. This approach of representing customers by a set of queries allows us to precompute query embeddings and their respective similarities, facilitating the generation of personalized feeds in near real-time (updates occurring every 2 minutes). This approach offers the benefits of efficiency, as queries are shared, and interpretability, as each recommendation is associated with a specific query.

(3) We showcase real-world applications of our approach in deployed systems at Bol, namely, *Top deals for you*, *Top picks for you*, and *New for you*. By indexing products that are on sale, new or popular and matching them with selected customer query representations, we generate the *Top deals for you*, *New for you*, and *Top picks for you* recommendations.

## 2 RELATED WORK

Pre-deep learning era, matrix factorization methods were used for personalized recommendations (see [9, 12, 13, 20]). Since the AlexNet paper [14], which showed the value of deep learning in image recognition, deep learning has also been applied in recommender systems [8, 30]. Among this rich literature (see survey [30]), the papers most related to our work come from industrial recommender systems such as those of eBay [24], Youtube [4, 27], Google Play [25], Pinterest[18, 19], and Alibaba [2, 16, 22]. We examine these papers on how they address the customer representation, item representation, and retrieval challenges.

## 2.1 Customer Representation Challenge

The YouTube paper [4] uses a Multilayer Perceptron (MLP) model to encode both user and video entities into the same space. The

user encoding model takes as inputs embedded video watches (50 recent watches), embedded search tokens (50 recent searches) and user attributes such as age and gender. A vocabulary of 1M videos and 1M search tokens is embedded with 256 floats.

The eBay paper [24] uses a recurrent (GRU) model to generate user embeddings. The inputs to the GRU model are item or query embeddings along with their respective event type embeddings. The event type embeddings are defined by four dimensions and serve to capture various actions on the items. The item embeddings are based on content-based features such as item titles, categories (e.g., mobile phones), and structured aspects (e.g., brand: Apple, network: Verizon). The user embedding has 64 dimensions.

The Pinterest paper [19] uses a transformer model to represent the user in 256 dimensions. The inputs to the model are: a sequence of Pins, represented by their PinSage embedding (256-dimensional) and metadata features: action type, surface, timestamp, and action duration [18].

To capture the diverse and multifaceted interests of users, prior work from Pinterest and JD.com used multiple embeddings per user [16, 18, 19]. While the notion of employing multiple embeddings to represent users is similar to our method, it also differs. In our solution, the embeddings that constitute customer representations are not unique to each individual customer but rather, are shared among users.

## 2.2 Item Representation Challenge

The YouTube paper [4] represents videos with embeddings of 256 dimensions based on Item IDs. The eBay study [24] employs a 3-layer MLP to create item embeddings with a 64-dimensional output. These embeddings are derived from inputs that include title, aspect, and category embeddings. Each of these embeddings is formulated as a Continuous-Bag-of-Words (CBOW) representation, corresponding to the tokens found in the title, aspect, and category. The Pinterest paper [19] uses an MLP model to represent items (more specifically, Pins) based only on PinSage embeddings of dimension 256.

Our work utilizes textual metadata (such as the title and category of a product) to embed item entities. In the YouTube paper, item IDs are used as input to the neural network model, leading to a larger model size due to the need to store an embedding table of significant size. In contrast, our approach generates embeddings directly from input metadata, eliminating the need for a separate table. This is similar to the eBay paper, which also utilizes metadata alone to represent items [24].

## 2.3 Candidate Retrieval Challenge

*2.3.1 Training challenge.* The most common training strategy for learning user and item embeddings is based on a two-tower user-item framework (see papers from eBay, YouTube and Pinterest [4, 19, 24, 27]). The user-item framework tackles the twin challenges of user representation and training using two neural networks in one go. The first network represents user activity of item views and searches whereas the second network represents target items. Variations exist in both the models employed for user and item representation, as well as in the input types fed into the model.

Additionally, variations arise in the negative sampling approach utilized during training.

Our training strategy also builds upon the two-tower model and negative sampling techniques. However, it emphasizes capturing item-to-item relationships, rather than the more common user-to-item relationships. During training for the retrieval stage, our work eliminates the necessity for user specific data and modeling the user, focusing solely on aggregated item-to-item relationships, specifically view-buy or buy-buy interactions.

*2.3.2 Inference challenge.* The approach to overcoming the inference challenge is essentially the same for all large-scale recommender systems. Embeddings of items are indexed and approximate nearest neighbor search is used to efficiently retrieve the most relevant items for given queries represented by user embeddings. Most systems differ in the tools used, e.g., the vector database. For example, eBay uses FAISS [24], an open source library from Facebook. Youtube and Pinterest use their own implementations [4, 7, 19]. Our work uses the FAISS library [11] for indexing and search operations. Since all potential query embeddings (item views and buys) are known in advance, we precompute their similarities and store the query results in a lookup table. Personalized recommendations are then generated by identifying relevant queries for a user and retrieving the corresponding recommendations.

## 3 METHODOLOGY

Our method for creating personalized feed recommendations, which we call Pfeed, involves two phases. In the first phase, we train and produce multi-vector item embeddings (see figures 3a and 3b). In the second phase, these embeddings are applied to generate personalized product recommendations (see figures 3c and 3d). The goal of the first phase is to capture item-to-item relationships through embeddings. We use "query-to-item" and "query-to-target" interchangeably to refer to the same concept of item-to-item relationships.

## 3.1 Representing an Item with Three Embeddings

In Pfeed, an item can play one of three roles: 1) view query, 2) buy query, and 3) target item. View queries are items clicked during a session leading to the purchase of specific items, thus creating view-buy relationships. Buy queries, on the other hand, are items frequently purchased in conjunction with or shortly before other items, establishing buy-buy relationships. The items that come after view or buy queries are the target items. Our goal is to capture the three roles of an item - view query, buy query, and target - using three distinct embeddings, all generated by a single encoder.

## 3.2 Model Architecture - Generating Three Item Embeddings with One Model Run

We use a transformer encoder [21] to generate three embeddings for a given item, each corresponding to the view, buy, or target role. To achieve this, we first tokenize the item metadata into a sequence of tokens using the sentencepiece library [15]. We then prepend three special tokens: `[Q_V]`, `[Q_B]` and `[TGT]` as shown in figure 4. These special tokens play a similar role as the `[CLS]` special token

(a) Step 1: Contrastive pre-training

(b) Step 2: Generating embeddings

(c) Step 3: Indexing and precomputing similarities

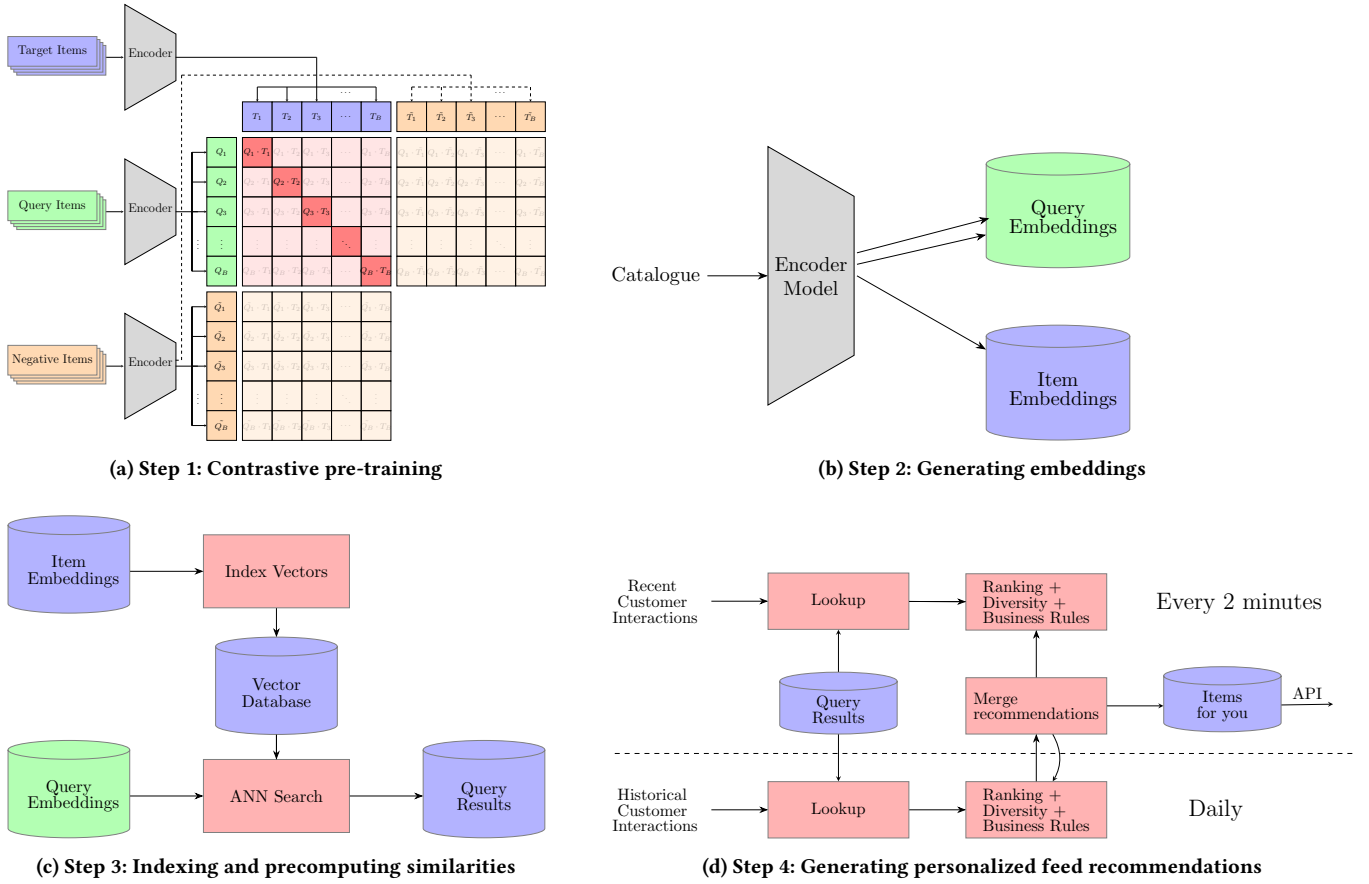(d) Step 4: Generating personalized feed recommendations

Figure 3: The major steps involved in generating near real-time personalized recommendations

in BERT [5]. The first three embeddings from the transformer's final layer, corresponding to the special tokens [Q_V], [Q_B], and [TGT], respectively represent the item's view query, buy query, and target embeddings. Because all these three embeddings are generated in one model run, we call the model a Single Input Multi Output (SIMO) embedding model. The SIMO model achieves threefold efficiency compared to a SISO (Single Input and Single Output) embedding model, which requires executing the model three times with distinct prompts for each of the three item roles (view query, buy query, and target roles).

### 3.3 Training with Contrastive Learning

*3.3.1 Training data.* We train the SIMO embedding model with query-target pairs consisting of the two types of relationships. The first set consists of item pairs of view-buy relationship (i.e., {$q$, view, $t$}). The second set consists of items pairs of buy-buy relationship (i.e., {$q$, buy, $t$}). We combine the sets to form one set $\{(q_i, r_i, t_i)\}_{i=1}^{N}$, where $(q_i, r_i, t_i)$ corresponds to a positive example, indicating that item $q_i$ and interaction (or relation) $r_i$ led to the purchase of item $t_i$. In addition to query-target item pairs, we also sample random items to reduce bias [26].
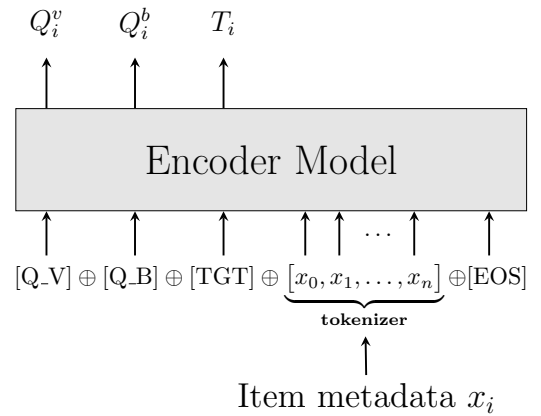


Figure 4: The SIMO (Single Input Multi Ouput) embedding model generates three embeddings per item in one model run using three special tokens: **[Q_V]**, **[Q_B]**, and **[TGT]**.

*3.3.2 Dual encoders.* The objective of our training is to get a model that produces similar embeddings for matching query-target $(q_i, r_i, t_i)$ inputs and dissimilar embeddings for non-matching inputs such as $(\tilde{q}_i, r_i, t_i)$ or $(q_i, r_i, \tilde{t}_i)$. To achieve this objective, we

employ dual encoders. We feed the query input $q_i$ and the target input $t_i$ into two instances of the transformer encoder $E$. The encoder $E$ maps $q_i$ and $t_i$ independently and outputs three embeddings of $Q_i^v$, $Q_i^b$, and $T_i$. From target encoder, we take $T_i$ embedding and do a dot product with the $Q_i$ embedding of the query encoder, which is $Q_i^v$ or $Q_i^b$, depending on the relation $r_i$ (see figure 5). When the training samples also include randomly sampled items, called random negatives, we use the same encoder $E$ to generate embeddings: $\tilde{Q}_i^v$, $\tilde{Q}_i^b$, and $\tilde{T}_i$ . These embeddings are mixed with the embeddings of in-batch negatives during training [26].
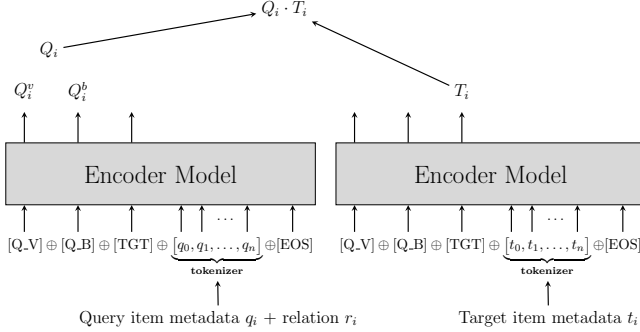


**Figure 5: Inputs to dual SIMO encoders: the query encoder takes in the metadata of the query item and generates three embeddings and the target encoder takes in the metadata of the target item and generates three embeddings. During training, the loss is determined by the target embedding derived from the target item $t_i$ encoder and pairing it with a query embedding from the query item $q_i$ encoder, selected by the relation $r_i$ indicator.**

*3.3.3 Training objectives.* The training objective consists of two contrastive loss terms. The first loss term employs a query-target softmax formulation (see equation 1). In this formulation, we sample negative targets for a given query-target pair. The second loss term employs a target-query softmax (see equation 2), where negative queries are sampled for the same query-target pair. We use four types of negative sampling strategies: 1) in-batch negatives, 2) uniformly sampled negatives, and 3) mixed negatives [26] which is a combination of in-batch negatives and uniformly sampled negatives, and 4) self-negatives.

$$\mathcal{L}_1 = -\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \log \underbrace{\frac{e^{\beta \mathbf{Q}_i \cdot \mathbf{T}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{\beta \mathbf{Q}_i \cdot \mathbf{T}_j} + \sum_{j=1}^{|\mathcal{N}|} e^{\beta \mathbf{Q}_i \cdot \tilde{\mathbf{T}}_j}}}_{\text{query} \rightarrow \text{target softmax}} \quad (1)$$

$$\mathcal{L}_2 = -\frac{1}{|\mathcal{B}|} \sum_{i=1}^{|\mathcal{B}|} \log \underbrace{\frac{e^{\beta \mathbf{Q}_i \cdot \mathbf{T}_i}}{\sum_{j=1}^{|\mathcal{B}|} e^{\beta \mathbf{Q}_j \cdot \mathbf{T}_i} + \sum_{j=1}^{|\mathcal{N}|} e^{\beta \tilde{\mathbf{Q}}_j \cdot \mathbf{T}_i}}}_{\text{target} \rightarrow \text{query softmax}} \quad (2)$$

In equations 1 and 2, $\mathcal{B}$ represents a batch of embedding pairs for positive samples: $\{(\mathbf{Q}_1, \mathbf{T}_1), (\mathbf{Q}_2, \mathbf{T}_2), \ldots, (\mathbf{Q}_{|\mathcal{B}|}, \mathbf{T}_{|\mathcal{B}|})\}$. $\mathcal{N}$ represents a set of embeddings from negative items that are uniformly

sampled from the catalog and appear as $\tilde{\mathbf{T}}_j$ or $\tilde{\mathbf{Q}}_j$, depending on the direction of the softmax computation (query-to-target or target-to-query). Each embedding is L2 normalized (i.e., $\|\mathbf{Q}\|_2 = 1$ and $\|\mathbf{T}\|_2 = 1$). The scale parameter $\beta$ is a parameter that is trained with the model parameters. Initially, we tried a few manually fixed values (e.g., 10, 100) and found it to affect performance significantly.

### 3.4 Inference

Pfeed has three inference steps: precomputing embeddings, precomputing similarities and generating personalized feeds.

*3.4.1 Precomputing embeddings.* After successful training using the approach described above, we use the resulting trained encoder to generate embeddings for all items in the catalog (see figure 3b). For each item, we generate three embeddings. The first two embeddings are query embeddings for when the item is viewed (indicated as embedding $Q_i^v$ in figure 4) or bought (indicated as embedding $Q_i^b$ in figure 4). The third embedding is for when the item is used as a target item (indicated as embedding $T_i$ in figure 4).

*3.4.2 Precomputing similarities.* The target embeddings of all items in the catalog (or selected part of it) are indexed with a vector indexing library (in our case, we use FAISS) and we search against the index using the view query and buy query embeddings of all items in the catalog. If the catalog has $N$ items, then we get $2 \times N$ queries (view and buy for every item in the catalog). For each of the $2 \times N$ queries, we get the $M$ most similar items, resulting in a table with $2 \times N \times M$ entries (see figure 3c). Only entries with a score greater than a prefixed threshold are stored in a lookup table. We fix this threshold from known item-to-item scores (validation data split). Similarity scores above the first percentile (approximately 15% of the original set) are stored in the lookup database.

*3.4.3 Generating personalized feeds.* The process for generating a ranked list of items per customer includes: 1) selecting queries for each customer (up to 100), 2) retrieving up to 10 potential next items-to-buy for each query, and 3) combining these items and applying ranking, diversity, and business criteria (see figure 3d). This process is executed daily for all customers and every two minutes for those active in the last two minutes. Recommendations resulting from recent queries are prioritized over those from historical ones.

### 3.5 Case Study: Personalized Item Feeds at Bol

We applied Pfeed to generate multiple personalized feeds at Bol, one of the largest e-commerce platforms of the Netherlands and Belgium. The feeds can be seen on the app or website and have titles such as *Top deals for you*, *Top picks for you*, and *New for you*. These feeds differ on at least one of two factors: the specific items targeted for personalization and/or the particular queries selected to represent customer interests.

*3.5.1 Top deals for you.* This feed personalizes items with promotional offers or discounted prices. Pfeed takes the most recent 100 unique customer item views/buys (per category) as query keys. And for each key, it retrieves up to 10 potential discounted items for the customer to buy. This is achieved by accessing precomputed query results and merging them, ensuring near real-time response in the

process. This is done daily for all customers and every 2 minutes for recently active customers (see figure 3d).

### 3.5.2 New for you.
This feed personalizes newly released items. New items, often marked by limited interactions, present a challenge to recommender systems reliant on item IDs or interaction data. However, Pfeed circumvents this cold-start issue because it generates item embeddings using textual catalog metadata [17]. The *New for you* feed works similarly to the *Top deals for you* feed, with the distinction being the type of items selected for personalization. In *New for you*, items are designated as new if they fall within the most recent 10% of items based on their release date, relative to their specific categories. This approach guarantees that each category features its own set of new items, accommodating the varying time scales across different categories.

### 3.5.3 X for you.
In general, Pfeed generates *X for you* by limiting the search index or the search output to consist of only items of *X*. In addition to *Top deals for you* and *New for you*, Pfeed has been used to generate other feeds, namely *Top picks for you* and *Select deals for you*. Items for *Top picks for you* come from those that have a certain level of popularity and match the customers' most recent queries from their most frequently interacted with categories. Items for *Select deals for you* come from items that are curated to reward customer loyalty and apply only to customers who are Select members.

## 4 EXPERIMENTS

To evaluate Pfeed, we run both offline and online experiments. The offline experiments are used to evaluate the quality of the embeddings and to illustrate the effects of different design choices on performance. To understand the impact of the embeddings on the personalized feed system, we report results from an online A/B testing experiment. The experiments are specifically designed to answer the following questions.

**Q1:** How does the model that produces three embeddings in one run (SIMO model) compare in terms of performance to the model that generates each embedding in three separate runs (SISO model)?

**Q2:** How effective is the SIMO model for cold-start product recommendation? And popular items?

**Q3:** How sensitive is the SIMO model to the training strategy, particularly concerning negative sampling and model sizes.

**Q4:** How effective are these query-target relationships in generating personalized feeds (online A/B testing)?

## 4.1 Dataset

We create view-buy and buy-buy datasets, comprising of approximately two million positive training/testing samples from around a million unique items (see table 1). These datasets are constructed from customer item views and item buys.

**Table 1: Bol dataset statistics**

| Dataset | # of positive pairs | # of distinct items |
| --- | --- | --- |
| view-buy | 0.99M | 1.08M |
| buy-buy | 0.96M | 0.27M |
| Negative | - | 2.00M |
| Combined | 1.95M | 3.28M |

### 4.1.1 view-buy dataset.
The view-buy dataset consists of item pairs with view-buy relationships. The pairs are constructed from converting customer sessions. Items that are purchased become target items and the items that were viewed in the same session become the view queries. Of all the view-buy pairs aggregated from sessions from the last four months, we choose the top one million pairs that meet a minimum occurrence threshold and have a high association strength as measured by a cosine metric [10]).

### 4.1.2 buy-buy dataset.
The buy-buy dataset consists of item pairs with buy-buy relationships. The pairs are constructed from customer purchases. Items that are purchased later in time become target items and the items that were purchased earlier in time become the buy queries. From all the possible buy-buy pairs constructed from the customer purchases, we select the top one million pairs that meet a minimum occurrence threshold and have a high association strength as measured by a cosine metric.

### 4.1.3 Negative dataset.
In addition to view-buy or buy-buy datasets, we also use a negative dataset that consists of uniformly sampled random items (about two millions). The purpose of this dataset is to reduce selection bias [26].

## 4.2 Offline Evaluation

We use the recall metric to compare different design choices. Our dataset is split into training, validation and test sets in the proportions of 80%, 10%, and 10%. To the target items $t_i$ in the test samples $(q_i, r_i, t_i)$, we add a distractor set $\tilde{C}$ of one million items, randomly sampled from the item catalog (a similar approach is used in ItemSage from Pinterest [1]). We consider a design choice to be better when its recall@K is higher, i.e., the proportion of $(q_i, r_i, t_i)$ samples for which the retrieved item $t_i$ is ranked within the top K among $\tilde{C} \cup t_i$.

## 4.3 Model Architecture Details

We use a transformer encoder model with four layers and eight attention heads. The model is identified as SIMO-128, where 128 represents the size of the hidden dimension. Depending on the input sequence we feed to the model, we have either a SIMO or a SISO embedding model.

## 4.4 Model Training Details

We use Pytorch and Pytorch lightning for training the transformer model. The model is optimized with Lamb optimizer [29] with a learning rate of 0.01 on four V100 GPUs using Distributed Data Parallel (DDP) strategy. Each GPU runs three instances of the model, each handling a batch size of 1024. These instances handle input sequences from query, target, and negative item sequences after

tokenization using the sentencepiece library [15] using a vocabulary size of 20k. Prior to loss computation, all forward passes from each GPU are gathered, resulting in a total batch size of $1024 \times 4 (= 4096)$. The loss is computed by incorporating both in-batch and uniformly sampled negative samples, amounting to a total of 8192 minus 1 negatives per positive sample [26]. To stabilize training, gradients are clipped to 0.5. The context length of the input sequence is fixed to a maximum of 64 tokens, sufficient for encoding item titles and essential metadata such as categories but excluding descriptions.

## 4.5 Retrieval Performance and Efficiency (Q1)

The query-target retrieval system, based on the embeddings generated by a transformer model that generates three embeddings with a single run (SIMO embedding model), performs comparably to the model that generates the embeddings separately (SISO embedding model). The SIMO embedding model generates embeddings three times faster than the SISO embedding model (see table 2).

**Table 2: Recall@K on view-buy and buy-buy datasets. The SIMO-128 model performs comparably to the SISO-128 model while being 3 times more efficient during inference.**

| Model | Recall@10 (%) | | |
|---|---|---|---|
| | view-buy dataset | buy-buy dataset | efficiency |
| SIMO-128 | 41.86 | 36.41 | 3x |
| SISO-128 | 41.57 | 36.12 | x |

## 4.6 Retrieval Performance on Cold-start and Popular Items (Q2)

The query-target retrieval system, based on the SIMO-128 model, shows varying performance depending on the nature of the dataset and the level of popularity of the items. On the buy-buy dataset, recall scores are lower for head items. On the view-buy dataset, recall scores are slightly higher for head items (see table 3). This recall score difference between the two datasets is attributed to the differing distributions of query-to-target relationship categories. On the buy-buy dataset, approximately 75% of the relationships are either one-to-many, many-to-one, or many-to-many (complex relationships). In contrast, on the view-buy dataset, such relationships constitute less than 21% (see table 4). A detailed analysis of recall scores segmented by relationship category reveals a consistent trend across both datasets: scores on item pairs with complex relationships are lower (see table 5). The reasons for this are twofold: First, single vectors face difficulties in capturing complex relationships. Second, during training, the model is inaccurately penalized for failing to replicate the exact query-target pairs provided, rather than being evaluated on its ability to identify any valid query-target pairs.

**Table 3: Impact of item popularity on Recall@K. Performance on popular items is lower than on tail items on the buy-buy dataset. This is due to a higher proportion of complex relations on the buy-buy dataset, indicated in table 4.**

| Popularity | Recall@10 (%) | |
|---|---|---|
| | view-buy dataset | buy-buy dataset |
| Cold-start | 38.52 | 59.76 |
| Tail | 41.66 | 55.88 |
| Head | 42.32 | 25.54 |
| All | 41.86 | 36.41 |

**Table 4: Relationship categories and their distributions. The buy-buy dataset has a higher distribution of complex relations ($1 \times n$, $m \times 1$ and $m \times n$ relations).**

| Relationship category | Distribution (%) | |
|---|---|---|
| | view-buy dataset | buy-buy dataset |
| $1 \times 1$ | 80.5 | 24.7 |
| $1 \times n$ | 6.9 | 16.2 |
| $m \times 1$ | 11.5 | 20.5 |
| $m \times n$ | 1.1 | 38.6 |
| All | 100.0 | 100.0 |

**Table 5: Relationship categories and Recall@K. Performance is higher on test data with simple $1 \times 1$ relations than with complex relations. The buy-buy dataset has a higher proportion of complex relations ($\sim 75\%$), see table 4.**

| Relationship category | Recall@10 (%) | |
|---|---|---|
| | view-buy dataset | buy-buy dataset |
| $1 \times 1$ | 42.08 | 58.01 |
| $1 \times n$ | 40.22 | 41.98 |
| $m \times 1$ | 41.71 | 35.55 |
| $m \times n$ | 37.63 | 20.72 |
| All | 41.86 | 36.41 |

## 4.7 Sensitivity of the Retrieval Performance (Q3)

We conduct a sensitivity analysis of our method by varying the hidden dimensions of the SIMO model and altering particular aspects of the training strategy, particularly the negative sampling strategy.

*4.7.1 Hidden dimension.* We vary the hidden dimension of the model between 64, 128, 256, 384, and 512 while keeping the rest of the transformer model and training strategy the same. Performance increases as the dimension increases until 384. At dimension 512, the model's performance drops (see table 6).

**Table 6: Impact of hidden dimension vector size on Recall@K**

| Vector size | Parameter # | Recall@10 (%) | |
|---|---|---|---|
| | | view-buy | buy-buy dataset |
| 64 | 1.5M | 37.87 | 32.09 |
| 128 | 3.6M | 41.86 | 36.41 |
| 256 | 9.1M | 44.31 | 40.73 |
| 384 | 16.6M | 44.71 | 41.61 |
| 512 | 26.0M | 41.23 | 38.93 |

*4.7.2 Negative sampling strategy.* We use four types of negative sampling strategies: in-batch negative sampling, uniform negative sampling, mixed negative sampling, and self-negative sampling. The best performance is achieved with mixed negative sampling, where both in-batch and uniform sampled negatives are used [26]. In-batch negative sampling is second best (see table 7).

**Table 7: Impact of negative sampling strategy on Recall@K**

| Negative Sampling | Recall@10 (%) | |
|---|---|---|
| | view-buy dataset | buy-buy dataset |
| Mixed | 41.86 | 36.41 |
| In-batch | 40.87 | 35.88 |
| Uniform | 39.15 | 31.73 |
| Mixed + self-negatives | 40.45 | 31.24 |

Self-negatives refer to instances where the target embeddings of query items serve as their own negatives (or the query embeddings of target items serve as their own negatives). Self-negatives are advantageous for handling asymmetrical buy-buy relationships or instances of non-repeat purchases. When we add self-negatives from query-item pairs having buy-buy relations to the mixed negatives, we observe a decline in the overall recall score. This suggests that such relationships are less prevalent in the dataset.

## 4.8 Online A/B testing (Q4)

We ran an online A/B testing experiment where we compared a treatment group receiving personalized *Top deals for you* item lists (generated by Pfeed) against a control group that received a non-personalized *Top deals* list, curated by promotion specialists. This experiment was conducted over a two-week period with an even 50-50 split between the two groups. The results showed a statistically significant increase in performance for the treatment group: there was a 4.9% increase in conversion rates and a 27% increase in the number of items added to wish lists (see table 8). Following these results, Pfeed has been deployed and can be found on both the mobile app and the website of Bol.

**Table 8: Online A/B test**

| Model | Wish list additions | Conversion |
|---|---|---|
| Non-personalized deals | 0.00 | 0.00 |
| Top deals for you | +27% | +4.9% |

## 5 CONCLUSIONS

In this paper, we introduced Pfeed, a method for generating personalized product feeds on e-commerce platforms. The method has been deployed at Bol with services called: *Top deals for you*, *Top picks for you*, and *New for you* and achieved a significant conversion uplift. Pfeed uses a query-to-item framework as opposed to user-item, the framework most dominant for personalized recommender systems. We highlighted three benefits of the query-to-item framework. 1) Simplification of real-time deployment, as query results can be precomputed and user interests can dynamically be updated in real-time, all without requiring model inference or the unlearning of past preferences. 2) Enhanced interpretability, as each recommendation in the feed can be traced to specific queries. 3) Increased computational efficiency due to the reuse of queries among users. Additionally, we demonstrated the use of multiple special tokens as input in the transformer model, enabling a single model run to generate multiple embeddings for each item.

## 6 FUTURE WORK

Pfeed's embedding approach can be enhanced in two ways: 1) better handling of query-to-item training samples having complex relations and 2) explicit modeling of memorization and generalization features.

**Modeling complex query-to-item relations**: Pfeed's current method of representing users with a set of individual queries provides flexibility but falls short in modeling sequential user behavior. This isn't inherently an issue, as the ranking phase can incorporate sequential information. However, it requires the embedding-based retrieval phase to be expressive enough to handle an increased set of relevant items, including those that might otherwise be excluded by sequential modeling. For example, if a user buys diapers, there are numerous potential next purchases such as items related to baby toys or clothes. Pfeed's embedding strategy struggles to model such complex relations (one-to-many, many-to-one and many-to-many relations). In practice, Pfeed settles with the most probable next purchase and thus provides less variety per query. Future enhancements could involve multi-vector query representations, allowing for a wider range of item choices.

**Modeling memorization and generalization features**: Pfeed's embedding strategy leverages item content, like titles, which is good for generalization, but it does not explicitly incorporate memorization features such as item IDs or popularity. This limitation could impact the system's performance, particularly with popular items. Future work could focus on designing an architecture that can adaptively use memorization features when available, while still relying on generalization features in their absence. This improvement would enable the system to more accurately predict next-item choices, covering both popular and long tail items.

## 7 ACKNOWLEDGMENTS

# 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* (Washington DC, USA) *(KDD '22)*. Association for Computing Machinery, New York, NY, USA, 2703–2711. https://doi.org/10.1145/3534678.3539170

[2] Yukuo Cen, Jianwei Zhang, Xu Zou, Chang Zhou, Hongxia Yang, and Jie Tang. 2020. Controllable Multi-Interest Framework for Recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (Virtual Event, CA, USA) *(KDD '20)*. Association for Computing Machinery, New York, NY, USA, 2942–2951. https://doi.org/10.1145/3394486.3403344

[3] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, Rohan Anil, Zakaria Haque, Lichan Hong, Vihan Jain, Xiaobing Liu, and Hemal Shah. 2016. Wide & Deep Learning for Recommender Systems. In *Proceedings of the 1st Workshop on Deep Learning for Recommender Systems* (Boston, MA, USA) *(DLRS 2016)*. Association for Computing Machinery, New York, NY, USA, 7–10. https://doi.org/10.1145/2988450.2988454

[4] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *Proceedings of the 10th ACM conference on recommender systems*. 191–198.

[5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, Jill Burstein, Christy Doran, and Thamar Solorio (Eds.). Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. https://doi.org/10.18653/v1/N19-1423

[6] Huifeng Guo, Ruiming Tang, Yunming Ye, Zhenguo Li, and Xiuqiang He. 2017. DeepFM: A Factorization-Machine Based Neural Network for CTR Prediction. In *Proceedings of the 26th International Joint Conference on Artificial Intelligence* (Melbourne, Australia) *(IJCAI'17)*. AAAI Press, 1725–1731.

[7] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. 2020. Accelerating Large-Scale Inference with Anisotropic Vector Quantization. https://arxiv.org/abs/1908.10396

[8] Xiangnan He, Lizi Liao, Hanwang Zhang, Liqiang Nie, Xia Hu, and Tat-Seng Chua. 2017. Neural collaborative filtering. In *Proceedings of the 26th international conference on world wide web*. 173–182.

[9] Yifan Hu, Yehuda Koren, and Chris Volinsky. 2008. Collaborative filtering for implicit feedback datasets. In *Data Mining, 2008. ICDM'08. Eighth IEEE International Conference on*. Ieee, 263–272.

[10] Yanxiang Huang, Bin Cui, Wenyu Zhang, Jie Jiang, and Ying Xu. 2015. TencentRec: Real-Time Stream Recommendation in Practice. In *Proceedings of the 2015 ACM SIGMOD International Conference on Management of Data* (Melbourne, Victoria, Australia) *(SIGMOD '15)*. Association for Computing Machinery, New York, NY, USA, 227–238. https://doi.org/10.1145/2723372.2742785

[11] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. *IEEE Transactions on Big Data* 7, 3 (2019), 535–547.

[12] Yehuda Koren, Robert Bell, and Chris Volinsky. 2009. Matrix factorization techniques for recommender systems. *Computer* 42, 8 (2009).

[13] Yehuda Koren, Steffen Rendle, and Robert Bell. 2022. Advances in collaborative filtering. *Recommender systems handbook* (2022), 91–142.

[14] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. ImageNet Classification with Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*, F. Pereira, C.J. Burges, L. Bottou, and K.Q. Weinberger (Eds.), Vol. 25. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2012/file/c399862d3b9d6b76c8436e924a68c45b-Paper.pdf

[15] Taku Kudo and John Richardson. 2018. SentencePiece: A simple and language independent subword tokenizer and detokenizer for Neural Text Processing. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*. Association for Computational Linguistics, Brussels, Belgium, 66–71. https://doi.org/10.18653/v1/D18-2012

[16] Chao Li, Zhiyuan Liu, Mengmeng Wu, Yuchi Xu, Huan Zhao, Pipei Huang, Guoliang Kang, Qiwei Chen, Wei Li, and Dik Lun Lee. 2019. Multi-Interest Network with Dynamic Routing for Recommendation at Tmall. 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, 2615–2623. https://doi.org/10.1145/3357384.3357814

[17] Jiacheng Li, Ming Wang, Jin Li, Jinmiao Fu, Xin Shen, Jingbo Shang, and Julian McAuley. 2023. Text Is All You Need: Learning Language Representations for Sequential Recommendation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (, Long Beach, CA, USA,) *(KDD '23)*. Association for Computing Machinery, New York, NY, USA, 1258–1267. https://doi.org/10.1145/3580305.3599519

[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 and 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] Nikil Pancha, Andrew Zhai, Jure Leskovec, and Charles Rosenberg. 2022. PinnerFormer: Sequence Modeling for User Representation at Pinterest. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining* (Washington DC, USA) *(KDD '22)*. Association for Computing Machinery, New York, NY, USA, 3702–3712. https://doi.org/10.1145/3534678.3539156

[20] Xiaoyuan Su and Taghi M Khoshgoftaar. 2009. A survey of collaborative filtering techniques. *Advances in artificial intelligence* 2009 (2009).

[21] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. 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. https://proceedings.neurips.cc/paper/2017/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf

[22] Jizhe Wang, Pipei Huang, Huan Zhao, Zhibo Zhang, Binqiang Zhao, and Dik Lun Lee. 2018. Billion-Scale Commodity Embedding for E-Commerce Recommendation in Alibaba. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (London, United Kingdom) *(KDD '18)*. Association for Computing Machinery, New York, NY, USA, 839–848. https://doi.org/10.1145/3219819.3219869

[23] Ruoxi Wang, Rakesh Shivanna, Derek Cheng, Sagar Jain, Dong Lin, Lichan Hong, and Ed Chi. 2021. DCN V2: Improved Deep & Cross Network and Practical Lessons for Web-Scale Learning to Rank Systems. In *Proceedings of the Web Conference 2021* (Ljubljana, Slovenia) *(WWW '21)*. Association for Computing Machinery, New York, NY, USA, 1785–1797. https://doi.org/10.1145/3442381.3450078

[24] Tian Wang, Yuri M Brovman, and Sriganesh Madhvanath. 2021. Personalized embedding-based e-commerce recommendations at ebay. *arXiv preprint arXiv:2102.06156* (2021).

[25] Ji Yang, Xinyang Yi, Derek Zhiyuan Cheng, Lichan Hong, Yang Li, Simon Xiaoming Wang, Taibai Xu, and Ed H. Chi. 2020. Mixed Negative Sampling for Learning Two-Tower Neural Networks in Recommendations. In *Companion Proceedings of the Web Conference 2020* (Taipei, Taiwan) *(WWW '20)*. Association for Computing Machinery, New York, NY, USA, 441–447. https://doi.org/10.1145/3366424.3386195

[26] Ji Yang, Xinyang Yi, Derek Zhiyuan Cheng, Lichan Hong, Yang Li, Simon Xiaoming Wang, Taibai Xu, and Ed H Chi. 2020. Mixed negative sampling for learning two-tower neural networks in recommendations. In *Companion Proceedings of the Web Conference 2020*. 441–447.

[27] 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 *(RecSys '19)*. Association for Computing Machinery, New York, NY, USA, 269–277. https://doi.org/10.1145/3298689.3346996

[28] 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 and Data Mining* (London, United Kingdom) *(KDD '18)*. Association for Computing Machinery, New York, NY, USA, 974–983. https://doi.org/10.1145/3219819.3219890

[29] Yang You, Jing Li, Sashank Reddi, Jonathan Hseu, Sanjiv Kumar, Srinadh Bhojanapalli, Xiaodan Song, James Demmel, Kurt Keutzer, and Cho-Jui Hsieh. 2019. Large batch optimization for deep learning: Training bert in 76 minutes. *arXiv preprint arXiv:1904.00962* (2019).

[30] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep learning based recommender system: A survey and new perspectives. *ACM Computing Surveys (CSUR)* 52, 1 (2019), 1–38.

[31] Wei Zhang, Dai Li, Chen Liang, Fang Zhou, Zhongke Zhang, Xuewei Wang, Ru Li, Yi Zhou, Yaning Huang, Dong Liang, et al. 2023. Scaling User Modeling: Large-scale Online User Representations for Ads Personalization in Meta. *arXiv preprint arXiv:2311.09544* (2023).

[32] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep Interest Network for Click-Through Rate Prediction. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining* (London, United Kingdom) *(KDD '18)*. Association for Computing Machinery, New York, NY, USA, 1059–1068. https://doi.org/10.1145/3219819.3219823