

Personalized Retrieval over Millions of Items

ABSTRACT

Personalized retrieval (PPR) seeks to retrieve items relevant to a user event (e.g. a page visit or a query) that are adapted to the user's personal preferences. For example, two users who happen to perform the same event such as visiting the same product page or asking the same query should receive potentially distinct recommendations adapted to their individual tastes. Personalization is seldom attempted over catalogs of millions of items since the cost of existing personalization routines scale linearly in the number of candidate items. For example, performing *two-sided* personalized retrieval (with both event and item embeddings personalized to the user) incurs prohibitive storage and compute costs. Instead, it is common to use non-personalized retrieval to obtain a small shortlist of items over which personalized re-ranking can be done quickly. Despite being scalable, this strategy risks losing items uniquely relevant to a user that fail to get shortlisted during non-personalized retrieval. This paper bridges this gap by developing the XPERT algorithm that identifies a form of two-sided personalization that can be scalably implemented over millions of items and hundreds of millions of users. Key to overcoming the computational challenges of PR is a novel concept of morph operators that offers multi-intent retrieval, can be used with arbitrary encoder architectures, completely avoids the steep memory overheads of two-sided personalization and offers millisecond-time inference. On multiple public and proprietary datasets, XPERT offered upto 5% superior recall and AUC than state-of-the-art techniques and upto 4% gains in click-through rates in live flights within a popular ad-serving network. Code for XPERT is available at the [supplementary \[link\]](#).

ACM Reference Format:

. 2023. Personalized Retrieval over Millions of Items. In *Proceedings of ACM Conference (Conference'17)*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Personalized Retrieval (PPR): the goal in a typical online retrieval setting is to take a user *event* such as a page visit or a query and retrieve items most relevant to that event from a large catalog of candidate items. Performing retrieval on the basis of such trigger or seed events is a popular paradigm in ranking and recommendation systems and several tasks such as product-to-product recommendation [18] and search [32] can be cast as “per-event” retrieval tasks. When such retrieval is non-personalized, we call it non-personalized retrieval (NPPR). Personalized Retrieval (PPR) on the other hand, seeks to adapt the retrieved item set to the user's

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
Conference'17, July 2017, Washington, DC, USA

© 2023 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM...\$15.00
<https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

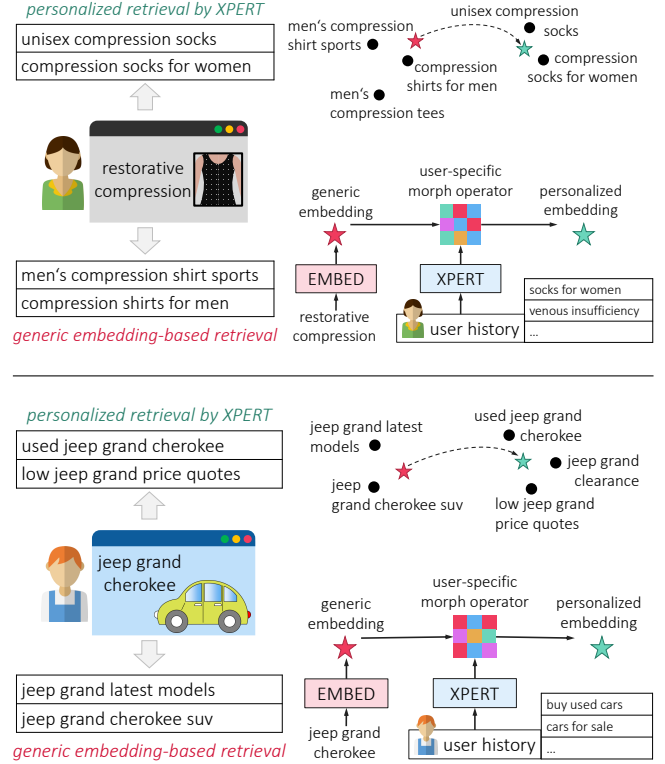


Figure 1: Actual personalized retrievals by XPERT (details in Table 12 in the [supplementary \[link\]](#)). XPERT personalizes the generic embedding of a user event using the *morph* operator for that user. The morphed event embedding is used to execute a MIPS query and offers superior personalized retrievals. (Left) A user visited a webpage titled “Restorative Compression” and selling compression garments. Prior browsing history indicated the user’s gender as well as their desire to purchase socks which XPERT correctly inferred. Suggestions by NPPR using the generic event embedding focussed on shirts for men instead. (Right) A user visited a webpage selling a vehicle but browsing history revealed the user’s desire to purchase discounted or low-cost products. XPERT successfully incorporated this into its suggestions but NPPR failed to do so.

long-term preferences such as can be inferred from their browsing history. This allows two users with distinct preferences to receive distinct retrievals even if they perform the same event e.g. visiting the same product page or asking the same query, thus enhancing user experience.

The Key to Efficient NPPR: Advances in dense retrieval such as Siamese networks [10, 14, 31] have enabled efficient NPPR. Consider a user u who performs event e and let a be a candidate item. Dense NPPR techniques learn a deep architecture ϕ to embed user events and candidate items into a shared vector space so that the

dot product $\langle \phi(a), \phi(e) \rangle$ is indicative of the relevance of item a to event e . Note that the embedding function ϕ does not depend on the user u i.e. it is non-personalized. At test time, the items most relevant to an event e_t can be retrieved simply by running a Maximum Inner Product Search (MIPS) query [15] with the test event embedding $\phi(e_t)$. This is usually orders of magnitude faster than explicitly computing $\langle \phi(a), \phi(e_t) \rangle$ for all catalog items a and can offer millisecond-time retrieval even when the catalog contains millions of items.

Challenges in PPR: It is challenging to accelerate PPR using the MIPS trick. One possibility is *two-sided* personalization that uses a personalized embedding function ϕ_u that is unique to every user i.e. use $\langle \phi_u(a), \phi_u(e) \rangle$ as the relevance score. Experiments show (see Table 5 and Figure 2) that two-sided personalization can indeed offer flexible and accurate retrieval. This could offer accelerated retrieval if a separate MIPS structure over the entire candidate set is established separately for every user. For a million items and a billion users, creating these structures alone would require million \times billion compute and storage. Another possibility is to use non-decomposable personalized scoring functions [17] of the form $f(a, e, u)$ indicating the relevance of a candidate item a for an event e performed by user u . Note that this score is personalized as it depends on the user u (e.g. their browsing history). Existing methods use nonlinear architectures e.g. MLP [9] or transformers [21] (see Section 2) that offer effective personalization but also preclude the MIPS trick and require applying f explicitly to each candidate item to find the most relevant ones. This does not require enormous storage but incurs infeasible inference times when candidate items are in the millions. The popular workaround is to shrink the candidate set by first performing NPPR and then applying personalized ranking only to the shortlisted items. However, items that fail to get shortlisted by NPPR step are irrevocably lost even if personalization could have later found them to be relevant to a particular user's tastes. Scalable PPR over millions of items poses opportunities and computational challenges.

Contributions: This paper presents the XPERT method that makes PPR possible on catalogs of millions of items with millisecond-scale inference time. The key contributions of XPERT include:

- (1) Introducing the concept of personalized *morph operators* that make PPR with a form of two-sided personalization possible with storage and inference costs similar to NPPR.
- (2) Integrating a scalable channel mechanism into XPERT that offers improved diversity in the retrieved items
- (3) Developing the XPERT technique that could train on a PPR task with 1+ million catalog items and 1+ billion user-item interactions within 2 hours on a single P40 GPU and offer 2.1 millisecond inference time on a single CPU.
- (4) 5% superior recall and AUC than state-of-the-art personalized retrieval methods and 4% gains in click-through rates on live A/B tests on a popular ad-serving network.

Figure 1 shows two real-life examples where XPERT retrieved personalized items relevant to the user event. Live A/B tests indicate that XPERT remarkably improves the overall accuracy of a recommendation pipeline. However, this paper focuses on improving retrieval through personalization and, as is usual in retrieval literature, other layers in a recommendation pipeline e.g. downstream

re-ranking over shortlisted items, item auctions, display-layout selection etc. are beyond its scope.

2 RELATED WORKS

As noted in Section 1, several personalized ranking and recommendation methods can only be applied to small item shortlists. Since they do not have to deal with catalogs of millions of items, these methods implement the scoring function $f(a, e, u)$ using elaborate architectures such as MLP [9], transformers [21], personalized attention [29], graph convolutions on user-item interaction graphs [11] and meta-path variants thereof [8]. However, as discussed in Section 1, these methods rely on items shortlisted via non-personalized NPPR routines and risk losing items, especially rare ones. This section instead focuses on methods most related to PPR tasks.

Single User Representation (SUR) Methods: SUR methods do not use a specific user event to trigger or seed the retrieval process (i.e. do not perform “per-event” retrieval) but rather embed the user itself as a vector and use a scoring function of the form $f(a, u) = \langle \phi(a), \xi(u) \rangle$ with a non-personalized item embedding function ϕ . Diverse models have been used to implement ξ as a function of user profile and browsing history including MLPs (feed-forward networks) [4], sequential models to encode browsing order such as GRU [1, 27] and LSTM [20], hierarchical attention networks [30] and transformers [25]. Although these methods offer scalable retrieval using the MIPS trick, they need to update the user embedding rapidly in response to browsing activities in order to stay relevant to the user's current intents which is a non-trivial overhead for real-time systems. They can also suffer from lack of diversity by relying on a single embedding. In experiments, XPERT could offer upto 15% better recall than SUR methods such as [4].

Channel-based Methods: A common technique [13, 22] to improve retrieval diversity, especially in SUR methods is to learn multiple user representations. User events are partitioned into *channels* each of which offers a distinct user representation. A variety of channeling techniques exist in literature including parametric ones such as [13] that learn multi-head architectures to effect channels as well as non-parametric ones such as [22] that simply cluster the event embeddings of a user with each cluster becoming a separate channel. In experiments, an NPPR variant that used channels could offer as much as 8.7% higher recall than SUR methods such as [4]. XPERT itself incorporates a scalable channeling architecture. It is notable however that channels by themselves do not offer sufficient personalization. For instance, XPERT offered as much as 8% higher recall compared to the NPPR+channels variant.

PPR Methods: The DPSR method [32] uses a decomposed scoring function of the form $f(a, e, u) = \langle \phi(a), \psi(e, u) \rangle$ that uses a non-personalized item embedding ϕ but a personalized event embedding ψ . The event embedding is personalized by concatenating a non-personalized embedding of the event e with an average of embeddings of historical events from the users history and applying several MLP layers. However, this approach to event embedding personalization does not seem very effective as noted by [32] themselves. In contrast, XPERT shows that by making use of personalized linear operators, one can effectively use a scoring function of the form $f(a, e, u) = \langle \psi(a, u), \psi(e, u) \rangle$ but in a manner that requires computational expense similar to those required by scoring

functions of the form $f(a, e, u) = \langle \phi(a), \psi(e, u) \rangle$. This implicitly allows personalization of not just the event embedding but candidate item embeddings as well and can offer upto 9% higher recall but millisecond-time inference at the same time. The use of linear operators to affect personalization has been studied in the past, such as the COT method [12] that implements personalization by using a linear (tensor) operator. However, the method works in the matrix completion setting and faces challenges in incorporating new items and new users inherent to all matrix completion-style methods.

3 PROBLEM SETTING AND NOTATION

Let \mathcal{U} denote the set of users and \mathcal{A} denote the catalog of items available for recommendation. \mathcal{A} can have millions of items and both \mathcal{U} and \mathcal{A} may dynamically evolve over time. For any user $u \in \mathcal{U}$, let \mathcal{H}_u denote the (ordered) list of their historical events. We will use the term *event* to denote user activity such as visiting a webpage, clicking on an ad, submitting a query, etc. Note that an event can be naturally identified with an item as well, say by considering the product for an ad click event or the webpage title for a page visit event. Events and items will be represented using their textual descriptions such as webpage title or query text. A shared encoder \mathcal{E} will be used to obtain generic (non-personalized) embeddings of events and items as D -dimensional unit norm vectors i.e. for any event e or item $a \in \mathcal{A}$, we have $\mathcal{E}(e), \mathcal{E}(a) \in S^{D-1}$ where S^{D-1} is the surface of the D -dimensional unit sphere. As is standard in dense retrieval applications, we assume access to a maximum inner product search (MIPS) structure such as [15] that offers log-time retrieval. A MIPS structure over a set of unit vectors $W \subset S^{D-1}$ can take a query vector say $\mathbf{v} \in S^{D-1}$ and $k \in \mathbb{N}$ and return the k vectors from W with the largest dot product $\langle \mathbf{w}, \mathbf{v} \rangle$ with \mathbf{v} (equivalently the nearest neighbors of \mathbf{v} since the vectors are unit norm). This $O(kD \log |W|)$ time and milliseconds in practice even when $|W|$ is in the millions. MIPS structures can also be updated to ingest new items. Given this, PPR requires the following two subtasks to be solved:

- (1) Identify a small subset $\hat{\mathcal{H}}_u \subset \mathcal{H}_u$ of *seed events* from the user history to seed the retrieval process.
- (2) For each seed event $e \in \hat{\mathcal{H}}_u$ retrieve a set of items from \mathcal{A} relevant to e and personalized to u .

Seed events are expected to be historical events that reflect the unique content/product preferences of the user and can, for example, include the most latest event performed by the user. Also note that the items (implicitly) associated with historical events may or may not be a part of the candidate set \mathcal{A} (e.g. a historical product item may no longer be on sale) but these seed events are key to initiating the retrieval process as personalized embeddings of seed events will be used by XPERT to retrieve candidate items most relevant to the user's preferences.

4 XPERT: EXTREME PERSONALIZED PER-EVENT RETRIEVAL

Due to lack of space, additional details and code for XPERT is provided with the supplementary at https://www.dropbox.com/sh/p4ifzbedogvwl1/AACbzSIF4ykXI_iCsp1_Cil9a?dl=0.

Motivation: The key desirables of personalized retrieval include:

- (1) *Computational Efficiency:* given a new user event, PPR-based retrieval must happen within milliseconds.
- (2) *Storage Overheads:* the model should not require more than $O(|\mathcal{A}|)$ storage. Note that $O(|\mathcal{A}|)$ storage is required for a MIPS structure over the set \mathcal{A} of candidate items even in NPPR.
- (3) *Long-short Term Adaptation:* the model should adapt to long-term user preferences e.g. those based on demographics but not neglect current (short-term) interests.
- (4) *Ease of Maintenance:* a trained model should not require frequent re-training.

PPR naturally addresses points 3 and 4 as choosing the last few events of a user as seed events to trigger the retrieval process is expected to cover the short-term interests of the user as well as offers multi-intent retrieval. Thus, personalization need only capture the long-term preferences of a user. However, these are by definition not expected to change rapidly and thus the model need not be re-trained frequently. To address points 1 and 2, we notice that a decomposed function of the form $f(a, e, u) = \langle \phi(a, u), \phi(e, u) \rangle$ that performs two-sided personalization using personalized embeddings of both the event e and candidate item a does allow MIPS-based retrieval. Table 5 and Figure 2 show that two-sided personalization can adapt to user tastes very flexibly. However, this would also require maintaining a separate, personalized MIPS structure for each user $u \in \mathcal{U}$ containing embeddings of each candidate item personalized to that user. Since each such MIPS structure would need $\Omega(|\mathcal{A}|)$ space, this would require $\Omega(|\mathcal{U}| \cdot |\mathcal{A}|)$ space in total which is infeasible with millions of users. This is why existing techniques such as DPSR [32] settle for a simplified scoring function $f(a, e, u) = \langle \phi(a), \psi(e, u) \rangle$ that uses personalized event embeddings but non-personalized item embeddings that requires a single MIPS structure and $O(|\mathcal{A}|)$ storage. XPERT's key novelty is a technique that enables a form of two-sided personalization i.e. using scoring functions of the form $f(a, e, u) = \langle \phi(a, u), \phi(e, u) \rangle$, but at $O(|\mathcal{A}|)$ storage cost and millisecond retrieval time.

Morph Operators: XPERT considers personalized embedding functions of the form $\phi(a, u) = P_u \cdot \mathcal{E}(a)$ and $\phi(e, u) = Q_u \cdot \mathcal{E}(e)$. Here \mathcal{E} is a text embedding model shared by all users that embeds events and items to D -dimensional unit vectors (XPRT uses a 6-layered DistilBERT base [24, 28] model) and user-specific *morph operators* $P_u, Q_u \in SO(D)$ for $u \in \mathcal{U}$. These morph operators are expected to encode user preferences for a certain type of product, e.g. low-priced ones, and personalize generic embeddings along those directions. For now, let P_u, Q_u be orthonormal matrices since they ensure that for any unit-norm vector \mathbf{v} , $M_u \mathbf{v}$ is also unit norm. We will relax this requirement soon. Thus, $f(a, e, u) = \langle \phi(a, u), \phi(e, u) \rangle = \mathcal{E}(a)^T P_u^T Q_u \mathcal{E}(e)$. However, since the set of orthonormal matrices $SO(D)$ is closed under transposition and product it must be that $P_u^T Q_u = L_u$ for some $L_u \in SO(D)$. Thus, $f(a, e, u) = \langle \mathcal{E}(a), \psi(e, u) \rangle$ where $\psi(e, u) = L_u \cdot \mathcal{E}(e)$ and L_u is the user-specific morph operator. Note that this re-parameterized scoring function uses non-personalized item embeddings and requires a single MIPS structure but actually embodies a two-sided personalization model. Since it is expensive to perform optimization over the rotation group $SO(D)$ directly [6], XPERT relaxes the orthonormality requirement by reparameterizing ψ as $\psi(e, u) = \Re((R_u + I_D) \cdot \mathcal{E}(e))$ with

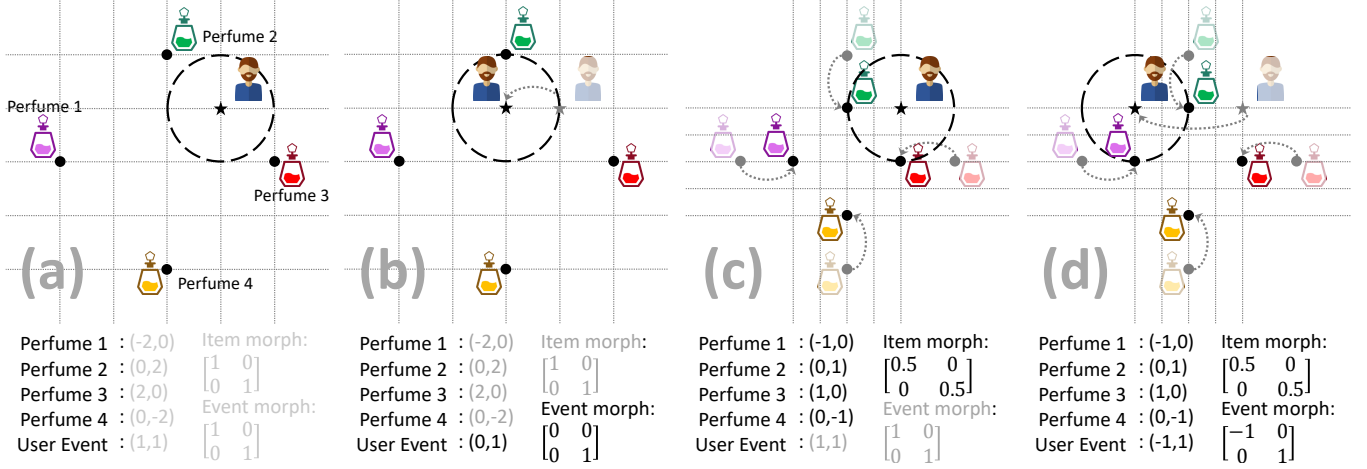


Figure 2: Illustrating the power of two-sided personalization, never mind its exorbitant cost. In this toy example, it can be shown that no matter what subset of items is desired by the user, there exist item and event morph operators P_u, Q_u that cause exactly that subset to be retrieved. In each of Figures 2(a,b,c,d), the five pointed star denotes the user-event embedding and the solid black circles denote embeddings for the 4 candidate items (perfumes). The grid lines illustrate the action of the item morph (event morph action is not shown to avoid clutter). Only items within a unit radius of the user event are retrieved (denoted by the dashed circle around the star). Coordinates of all embeddings are listed below each figure. Gray coordinates are generic and those in black are morphed. The corresponding (linear) morph operators are listed alongside. Identity morphs (equivalent to no morph) are typeset in gray. Figure 2(a) shows the generic embeddings for the user-event and the items – no perfume is retrieved if generic embeddings are used. Figure 2(b) shows that morphing just the event embedding can allow retrieval of any single perfume by changing the event morph operator but may struggle if the user is interested in say two or three perfumes. Figure 2(c) shows that morphing just the item embeddings can allow retrieval of Perfumes 2 and 3 together. Figure 2(d) shows that if morphing both the event and item embeddings are allowed (two-sided personalization), it becomes convenient to retrieve other pairs of perfumes. The complete example in Figure 5 in the [supplementary link](#) shows that two-sided personalization allows retrieval of any subset (pair or triplet or all) of perfumes. By implementing personalization using orthonormal morph operators, XPERT is able to offer the benefits of two-sided personalization without incurring the prohibitive storage cost. We note that the morphed vectors in this illustrative example are non-unit norm to prevent clutter given the 2-dimensional canvas available for the illustration.

$R_u \in \mathbb{R}^{D \times D}$ as the morph operator (that need not be orthonormal) and $\mathfrak{N} : \mathbf{v} \mapsto \mathbf{v} / \|\mathbf{v}\|_2$ being normalization. The skip connection with the identity matrix I_D provides regularization and also allows training to commence with non-personalized embeddings by initializing with $R_u = 0$.

Prediction Pipeline: XPERT uses a 3-segment pipeline (see Fig. 3): **Segment S1:** For a user $u \in \mathcal{U}$, embeddings of events in the user history $\{\mathcal{E}(e) : e \in \mathcal{H}_u\}$ are aggregated using a 2-layer transformer with 8 attention heads [26] without positional encoding to obtain an *intermediate user embedding* $\hat{\mathbf{z}}_u \in \mathbb{R}^D$.

Segment S2: The intermediate user embedding $\hat{\mathbf{z}}_u$ is passed through a single feed-forward layer with ReLU non-linearity and the D^2 -dimensional output is reshaped into a $D \times D$ matrix that serves as the morphing operator R_u for user u . XPERT uses $D = 64$ for speed. Mild accuracy boosts were observed by using a multi-head architecture that passed $\hat{\mathbf{z}}_u$ through multiple feed-forward layers and averaged the resulting matrices to obtain R_u .

Segment S3: Using a seed selection strategy discussed below, a set $\hat{\mathcal{H}}_u \subset \mathcal{H}_u$ of s seed events is selected from user history and the personalized embedding $\psi(e, u) = \mathfrak{N}((R_u + I_D) \cdot \mathcal{E}(e))$ of each seed event $e \in \hat{\mathcal{H}}_u$ is used to query a MIPS structure over (generic embeddings of) the candidate items $\{\mathcal{E}(a) : a \in \mathcal{A}\}$ to obtain a set $\hat{\mathcal{A}}(e, u)$ of items. These are combined $\bigcup_{e \in \hat{\mathcal{H}}_u} \hat{\mathcal{A}}(e, u)$ to make the

final prediction. Note that all events in \mathcal{H}_u are used to learn R_u and not just the seed events $\hat{\mathcal{H}}_u$.

XPERT Training: The model was trained to predict the next event of a user $u \in \mathcal{U}$ based on their user history \mathcal{H}_u . To compute loss over a user $u \in \mathcal{U}$, the item involved in the next event $a_u^* \in \mathcal{A}$ was taken as a ground truth positive, a set $\hat{\mathcal{H}}_u$ of s seed events was chosen from the user history (seed selection discussed in Appendix B in the [supplementary link](#)) and a set $\mathcal{N}_u \subset \mathcal{A}$ of n hard-negative items was identified (negative mining discussed below). s and n were tuned as hyperparameters (see Appendix D). Next, the seed event whose personalized embedding most closely resembled the clicked item/ad was chosen as the *most promising* seed event $e^* = \arg \max_{e \in \hat{\mathcal{H}}_u} \langle \mathcal{E}(a_u^*), \psi(e, u) \rangle$. The hard negative item most similar to the most promising seed event $b^* = \arg \max_{b \in \mathcal{N}_u} \langle \mathcal{E}(b), \psi(e^*, u) \rangle$ was also chosen. Then, using margin hyperparameters λ_+, λ_- , loss on a user is computed as

$$\ell(u) = \max \{ \lambda_+ - \langle \psi(e^*, u), \mathcal{E}(a_u^*) \rangle, 0 \} + \max \{ \langle \psi(e^*, u), \mathcal{E}(b^*) \rangle - \lambda_-, 0 \}.$$

The loss is averaged over all training users i.e. $\mathcal{L} = \frac{1}{|\mathcal{U}|} \sum_{u \in \mathcal{U}} \ell(u)$ and is used to train the transformer layers in segment S1 and the feed-forward layer in S2 using the Adam optimizer.

Hard Negative Mining: Items that seem deceptively likely to get clicked/viewed as the next event but were not a part of the ground

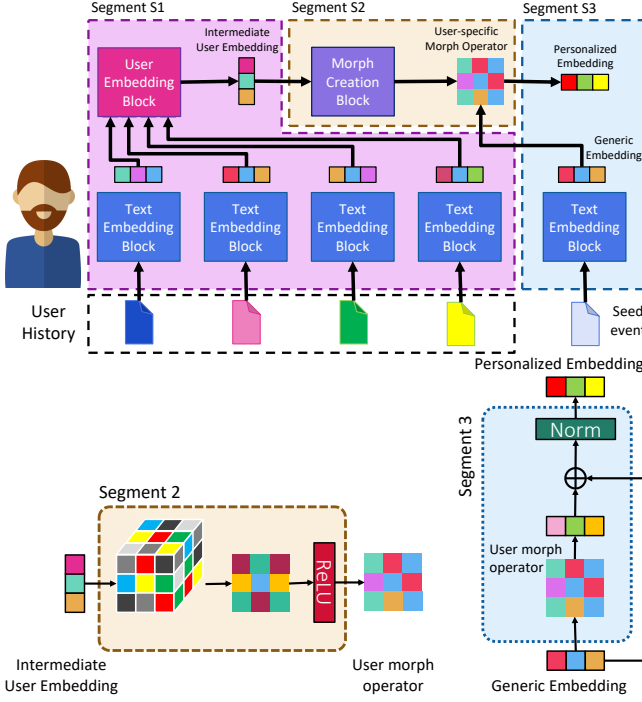


Figure 3: (Left) XPERT’s pipeline consists of 3 segments. Segment S1 aggregates user history to obtain an intermediate user embedding that is used by S2 to extract a user-specific morph operator. S3 takes a seed event from user history and offers a personalized embedding that is used to perform retrieval by making a MIPS call. These three light-weight segments offer personalized retrieval within milliseconds. **(Center)** A schematic of segment S2. **(Right)** A schematic of segment S3.

truth are termed hard negatives. Negative sampling is necessary [3, 5, 7, 31] with large output spaces since evaluating \mathcal{L} with respect to *all* negative items in \mathcal{A} would take $\Omega(|\mathcal{U}| \cdot |\mathcal{A}|)$ time. XPERT considered two types of negative mining: in-batch negative mining and global negative mining (see details in Appendix A in the [supplementary \[link\]](#)).

Seed Selection and Channels: XPERT explored two seed selection strategies: recency- and channel-based. The latter offered moderate boosts in retrieval accuracy but required an online algorithm for channel maintenance (see Appendix B). Channel-based architectures are popular [13, 22] and increase diversity but do not offer significant personalization by themselves.

5 EXPERIMENTAL RESULTS

Datasets: Experiments were conducted on multiple public (AmazonReviews) and proprietary (U2A) PPR datasets (see Table 1) on a 24-core Intel Xeon 2.6 GHz machine with a single Nvidia P40 GPU.

- **AmazonReviews** - The AmazonReviews-(1M/10M) datasets were created out of the Amazon Review Data dump [19] and are publicly available in the [supplementary \[link\]](#). The PPR task in these datasets is predicting the product a user will review next, given the past products they have reviewed. The $n - 2$ history items (where n is the user history length) of a user were used for training

Table 1: Dataset Statistics. A ‡ indicates data redacted for proprietary datasets. U2A datasets are proprietary, the AmazonReviews datasets are available at the [supplementary \[link\]](#).

Dataset	$ \mathcal{U} $	Total # of interactions	$ \mathcal{H}_u $	$ \mathcal{A} $
U2A-4M	3.96M	1.093B	‡	1.02M
U2A-300M	316M	23B	‡	19.4M
AmazonReviews-1M	920K	9.67M	10.42	286K
AmazonReviews-10M	9.71M	156M	16.06	3.7M

while the last two items were used for evaluation. Each history item was represented by the 768-dimensional embeddings of the product title from a pretrained 6-layered DistilBERT base model [24, 28]. Two versions of this dataset were created, namely AmazonReviews-1M and AmazonReviews-10M. For the smaller dataset, reviews from the following 5 categories were considered: CDs and Vinyl, Games, Electronics, Beauty, Grocery and Gourmet Food whereas all categories were considered for the AmazonReviews-10M dataset. It is notable that that user history events and candidate set items came from a shared domain in these tasks.

- **U2A** The U2A datasets were created from ad-click logs mined for two consecutive months from a popular ad-serving network. For each click, the user’s browsing history for the previous month was used to predict the final ad clicked by the user. We note that browsing history was used and not just ad-click history since ad-clicks are scarce and browsing history offered a much more informative signal to model the user’s preferences. Clicks with no accompanying historical browsing activity were filtered out. The first month’s data was used for training and second month’s data was used for testing. The final data was randomly subsampled to form U2A-300M and U2A-4M datasets. Each ad and webpage was represented by 64-dimensional embeddings of the ad/page title from a pretrained 6-layered DistilBERT base model. Note that user history events and candidate set items did not come from a shared domain in these datasets since the candidate set items were ads whereas user history events included their entire browsing activity.

Implementation and Hyperparameter Details: This is discussed in Appendix D in the [supplementary \[link\]](#)

Baselines: Several SUR/PPR baselines were used. Appendix E gives implementation details for these baselines.

- **NPPR:** This is a popular technique in the retrieval layers of commercial ad placement engines for item-to-item or product-to-product recommendation tasks [16, 18, 22, 23] and serves as a non-personalized competitor to XPERT as it uses a generic event embedding to perform retrieval.

- **SUR:** These use a single user embedding and do not use a seed event by the user to trigger the retrieval process. Various architectures have been used to implement SUR methods e.g. transformers (BERT4Rec [25]), feedforward networks (YouTube-DNN [4]) and GRUs ([27]) – see Section 2. We compare XPERT with the SUR-BERT and SUR-DNN baselines that closely resemble BERT4Rec and YouTube-DNN.

- **DPSR** [32]: This method was adapted to PPR tasks by seeking retrieval with respect to all events in a seed event set. DPSR was

offered advantages such as channel-based seed event selection and transformer-based user history aggregation ([32] use a simple average instead).

- **PinnerSage** [22]: This is a non-SUR method that clusters a user's history into multiple channels followed by NPPR-based retrieval for every channel.
- **XPert (morph only)**: This variant used morph operators but only the last s events were chosen as seeds.
- **NPPR+channel**: This variant offered channel-based seed event selection to NPPR and closely resembles the PinnerSage method [22] but for differences in clustering algorithm.

Note that NPPR offers retrieval but which is not personalized while SUR variants offer personalized retrieval but that may suffer from diversity and coverage issues since they do not trigger the retrieval process using a user event. In contrast, XPert, PinnerSage and DPSR offer personalized retrieval that is diverse and offers good coverage of user interests since retrieval is triggered by a user event. Certain baselines were implemented by closely following their corresponding publications due to lack of publicly available code.

Text Encoder Training: XPert used a 6-layered DistilBERT base architecture [24, 28] pretrained in a Siamese fashion [5] as its text encoder \mathcal{E} . For U2A, a query-to-ad recommendation dataset mined from search engine click-logs was used to pretrain \mathcal{E} on the (non-personalized) task of predicting the ad that would be clicked in response to a search engine query. For AmazonReviews, \mathcal{E} was pretrained on the LF-AmazonTitles-1.3M dataset [2] on a product-to-product recommendation task. \mathcal{E} was not fine-tuned while training XPert although doing so could yield minor improvements. The AmazonReview datasets are available at the [supplementary \[link\]](#).

Evaluation Metrics: Performance was evaluated using standard retrieval performance measures: Recall@ k , nDCG@ k , MRR@ k ($k \in \{10, 50, 100\}$) and AUC (see Appendix C for metric definitions).

Results: Table 2 presents results on all datasets. XPert was found to beat all baselines across all metrics on both public and proprietary datasets. XPert's Recall@100 was at least 10% higher than NPPR which is a widely adopted retrieval method in commercial settings [16, 18, 22, 23]. In addition to this, XPert outperforms the personalized retrieval baselines by up to 5% in terms of Recall@100. XPert could also be up to 4%, 7.5% and 5% better in terms MRR@100, AUC@100 and nDCG@100 respectively compared to DPSR on U2A-4M dataset. Similar trends were observed on the two AmazonReviews datasets as well.

Analysis of Results: Fig 4 shows the contributions of each item decile to the overall recall of various baselines. On U2A-4M, NPPR offered similar performance on both rare and popular items whereas SUR-BERT, SUR-DNN and DPSR all performed worse than NPPR on rare items. PinnerSage was better than NPPR across all deciles. However, XPert offers the best of both worlds and performs better than NPPR on rare items and better than SUR-BERT/DPSR/PinnerSage for popular items. On AmazonReviews-1M, XPert again leads on all deciles whereas baselines underperform either on popular items or rare items or both. NPPR+channel and XPert (morph only) were found to offer almost 4% and 7% gains over NPPR on U2A-4M revealing the benefits of morph operators and careful seed selection. It is notable that gains for XPert (morph only) are attributable to

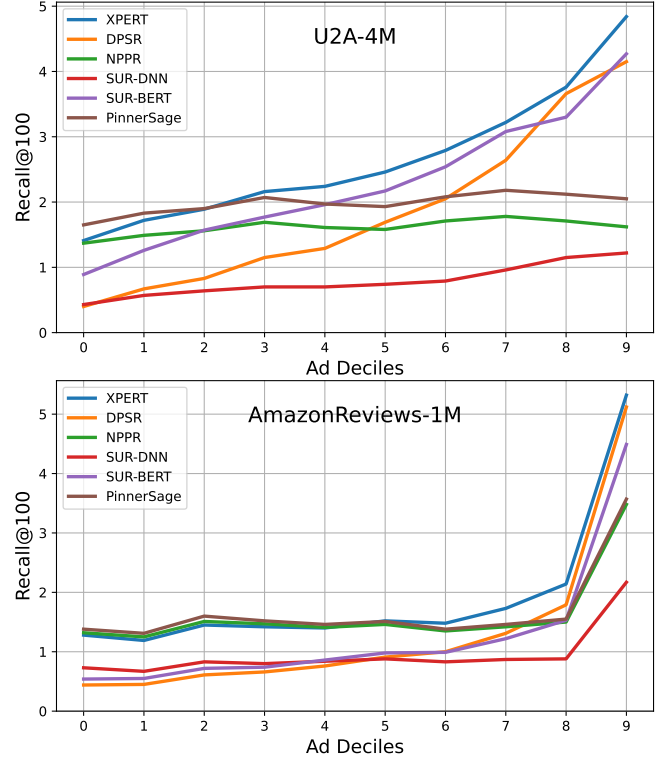


Figure 4: Decile-wise breakdown of Recall@100. Equi-voluminous item deciles were created on the basis of item popularity with decile 0 containing the most rarely clicked items and decile 9 containing the most popularly clicked items. XPert leads on all deciles, often by a margin. In contrast, other methods such as DPSR do well on popular deciles but not so on rare deciles or the other way round.

the personalized embeddings whereas those for NPPR+channel are attributable to channelling.

Live A/B testing: XPert was deployed on a popular ad-serving network and A/B tests were performed against an ensemble of leading (proprietary) information retrieval, NPPR and SUR techniques. Downstream layers (e.g. layout, auctions) were shared by all methods to ascertain the differential contribution of XPert's personalized retrieval. Performance was measured in terms of Click-Through-Rate (CTR) defined as the fraction of impressed ads that were clicked by users. XPert was found to increase CTR by 3.95% while showing the same number of ads indicating the alignment XPert offered to users' personal tastes. This also revealed that NPPR was severely limiting the amount of personalization possible in downstream ranking layers.

Scalability and Inference Time: XPert could train on the largest U2A-300M dataset with 316M users and 23 billion user-item interactions within 48 hours on a single P100 GPU. XPert supports scalable inference essential for live applications. Segment S1 intermediate user embeddings \hat{z}_u were updated nightly and did not affect live inference latency. XPert needed 256 bytes of additional memory per user to store the 64 dimensional vectors \hat{z}_u . Upon any user activity, updating the seed events and applying the user-specific morph operator to obtain personalized seed event embedding took

Table 2: Performance of XPERT vs baseline and competing algorithms across datasets and metrics. XPERT was found to offer up to 5% superior recall and 5% superior AUC compared to state-of-the-art in personalized retrieval techniques.

Method	Recall@10	Recall@50	Recall@100	nDCG @10	nDCG @50	nDCG @100	AUC@100	MRR@100
U2A-4M dataset								
NPPR	6.450	15.502	16.155	3.726	5.762	5.872	13.451	3.364
NPPR+channel	8.504	17.879	19.715	4.617	6.713	7.019	15.988	3.938
SUR-DNN	4.257	9.149	11.036	2.346	3.422	3.730	8.320	2.023
SUR-BERT	12.635	20.042	22.864	8.487	10.116	10.575	18.686	7.592
DPSR	8.371	17.040	18.352	4.687	6.617	6.836	15.147	4.016
PinnerSage	8.633	17.859	19.855	4.535	6.597	6.929	16.027	3.790
XPERT	14.852	25.778	27.189	9.318	11.790	12.026	23.390	8.212
AmazonReviews-1M dataset								
NPPR	3.467	7.620	9.034	2.115	3.224	3.509	11.997	2.468
SUR-DNN	2.616	4.591	5.413	1.702	2.242	2.406	7.327	1.973
SUR-BERT	4.285	6.227	7.213	3.684	4.208	4.404	10.333	4.888
DPSR	3.879	6.539	7.401	3.000	3.725	3.899	10.590	3.932
PinnerSage	3.638	7.707	9.342	2.268	3.355	3.683	12.254	2.713
XPERT	5.505	9.344	10.70	4.422	5.453	5.724	14.849	5.799
AmazonReviews-10M dataset								
NPPR	2.383	4.793	5.449	1.454	2.106	2.239	7.484	1.662
SUR-DNN	1.368	2.039	2.270	0.930	1.115	1.161	3.182	1.054
SUR-BERT	3.277	4.329	4.710	2.665	2.956	3.032	7.017	3.316
DPSR	2.795	4.240	4.426	2.114	2.516	2.554	6.688	2.568
PinnerSage	2.499	5.131	6.027	1.512	2.218	2.399	8.063	1.742
XPERT	3.943	6.106	6.550	3.136	3.724	3.815	9.740	3.992
U2A-300M dataset								
NPPR	7.140	15.901	17.819	3.413	5.891	6.190	15.225	3.478
XPERT	9.682	20.010	22.291	6.013	9.348	10.001	20.119	8.012

0.1 ms and the MIPS call took upto 2 ms. All latency results are reported on single core CPU.

Benefits of Channels: Appendix F presents a discussion. Tab 10 shows examples of channels created by Algo 1 that evidently capture distinct product categories. Tab 11 shows that using channels allows XPERT to use superior seed events to increase retrieval diversity compared to NPPR. Note that Algo 1, Tab 11 and Tab 10 are all available in the [supplementary \[link\]](#).

Illustrative Examples: As noted earlier, XPERT draws its gains from two primary sources (1) personalization by the user-specific morph operators and (2) diversity offered by channels. Table 12 shows two examples where ads retrieved by XPERT were better aligned to user tastes than NPPR. In the first example, XPERT could retrieve the ad actually clicked by the user. In the first example in Table 12, user history indicated searches such as "compression socks for women". Note that Tables 12 and 12 are both available in the [supplementary \[link\]](#). However, the seed event namely "Restorative Compression - Aqua Confetti - Primes Compression" was actually related to compression socks for men although the title does not reveal this. Thus, NPPR recommended "compression shirts for men".

In contrast, XPERT accurately captured both the gender and the intent of the user and recommended ads such as "unisex compression socks" that indicate a far superior alignment to user preferences. It is notable that the candidate set of items did not contain any web page such as "compression socks for women". The second example in Table 12 shows a user interested in used or low-price cars. However, the seed event page selected for retrieval, "2019 Jeep Grand Cherokee Limited RWD - \$27,995 - CarGurus" did not contain this information and hence NPPR recommended ads for new cars such as "new Jeep Grand Cherokee". However XPERT successfully latched on to the user intent and made recommendations such as "Used Jeep Grand Cherokee". Table 11 shows how channeling can increase retrieval diversity by selecting diverse seed events.

Ablation Experiments: Table 3 explores the impact of various design choices by XPERT,

- **seg1-mean:** replacing the transformer layers in S1 with a simple tf-idf weighted mean caused a drop of 1.4%, 1.6%, 2%, 1.4% in Recall@100, nDCG@100, AUC@100 and MRR@100 respectively, indicating the benefit of accurate user history aggregation.
- **seg3-no-residual:** removing the skip connection in S3 led to a drop in recall@100, AUC@100 and MRR@100. Residual connections

Table 3: Ablation experiments for XPert on the U2A-4M dataset

Ablation Name	Recall@10	Recall@50	Recall@100	nDCG @10	nDCG @50	nDCG @100	AUC@100	MRR@100
XPert	14.852	25.778	27.189	9.318	11.790	12.026	23.390	8.212
seg1-mean	12.783	23.088	24.644	7.643	9.969	10.229	20.888	6.628
seg3-no-residual	14.773	24.879	26.002	9.335	11.631	11.819	22.621	8.217
NPPR-BoE	5.235	12.184	12.712	3.135	4.702	4.791	10.632	2.854
XPert-BoE	12.339	19.600	21.024	7.612	9.270	9.506	18.120	6.570
neg-inbatch-hard	13.546	24.494	25.942	8.149	10.624	10.865	22.111	7.089
userwise-batching	13.065	23.784	25.340	7.746	10.167	10.428	21.484	6.700
loss-triplet	14.424	24.785	26.081	9.008	11.356	11.573	22.507	7.910

Table 4: Ablation experiments for DPSR on the U2A-4M dataset

Method	Recall@10	Recall@50	Recall@100	nDCG @10	nDCG @50	nDCG @100	AUC@100	MRR@100
DPSR	8.371	17.040	18.352	4.687	6.617	6.836	15.147	4.016
mean	7.806	15.409	16.645	4.379	6.073	6.279	13.763	3.735
no channels	7.385	14.734	15.624	4.186	5.828	5.977	13.074	3.593
DPSR*	10.899	21.522	23.201	6.144	8.533	8.813	19.282	5.267

stabilize training and prevent overpersonalization that could lead to a morphed embedding completely unrelated to the seed event.

- **NPPR-BoE & XPert-BoE:** In this ablation, the DistilBERT text embedding model \mathcal{E} was replaced with a simpler bag-of-embeddings model from [5]. It is notable that although overall accuracy does decrease, XPert-BoE continues to offer similar gains over NPPR-BoE indicating that XPert offers its benefits irrespective of the base embedding model used in the pipeline.

- **Neg-inbatch-hard & userwise-batching:** Using in-batch-hard negatives were used in place of global-hard negatives and user-wise batching instead of item-wise batching caused noticeable drop was observed in metrics suggesting the importance of item-wise batching and negative sampling.

- **loss-triplet:** A triplet loss with margin 0.3 instead of a contrastive loss led to a drop in all metrics.

Table 4 reports additional ablations performed on the DPSR method.

- **mean:** Replacing the transformer aggregation (offered to DPSR as an advantage) with mean aggregation as suggested in [32] reduced accuracy consistent with the seg1-mean ablation for XPert.

- **no channels:** Removing the channel advantage offered to DPSR led to a drop of 2% in recall@100 indicating the utility of channeling irrespective of the (N)PPR method

- **DPSR*:** This variant was offered advantages over the base DPSR method such as additional MLP layers, residual connection, hard-negative mining in addition to the transformer based aggregation as well as channel-based seed event selection. This meant that the most prominent distinction between DPSR and XPert was the personalization mechanism with XPert using morphing operators and DPSR using its MLP architecture. As such, this experiment allowed us to explore the limits of DPSR-style architectures. DPSR* achieves noticeable gains over DPSR but was still worse than XPert by 4%, 3% in Recall@100 and MRR@100 respectively. This corroborates an observation of [32] that leveraging user history resulted in only marginally better performance for DPSR than NPPR.

Two-sided vs One-sided personalization: Given its exorbitant storage costs, two-sided personalization was explicitly implemented on a randomly subsampled subset of the U2A-4M dataset with 10k

Table 5: Training loss value (see Section 4 for a definition of \mathcal{L}) for two-sided vs one-sided DPSR

Method	Loss Value
Two-sided DPSR	0.183
One-sided DPSR	0.435

Table 6: Recall@k values for NPPR and XPert for large values of $k = 200, 1000$.

Method	Recall@200	Recall@1000
NPPR	20.5	31.2
XPert	31.9	42.5

candidate ads and 150k users. Table 5 shows that a two-sided version of the DPSR method, where both event and item embeddings were personalized using a DPSR model, obtained significantly smaller loss value confirming the power of two-sided personalization.

Re-ranking after PPR: as noted in Section 1, XPert focuses on improved retrieval via personalization and other layers in a typical recommendation pipeline such as re-ranking, display-layout are beyond its scope. As indicated by the 3.95% CTR boost offered by XPert in live A/B tests, the use of XPert does not preclude the use of downstream re-rankers to perform further personalization, diversification etc but XPert instead boosts the accuracy of the entire pipeline. To gain insights into this aspect, further results are presented. Table 6 shows that even when comparing recall at large values of k such as recall@200 or recall@1000, XPert continues to outperform NPPR by a large margin i.e. merely increasing the shortlist size cannot overcome the drawbacks of non-personalized retrieval. Table 8 indicates that XPert outperforms NPPR even if NPPR results are re-ranked using XPert itself indicating the irreversible loss of items risked by non-personalized retrieval. Table 7 offers ranking metrics such as precision@k values for XPert and various competitors where XPert continues to offer gains.

Table 7: Precision@K metric for various methods on the U2A-4M and AmazonReviews-1M datasets. Despite being focused on retrieval, XPERT offers superior precision numbers than competing methods.

	U2A-4M dataset		AmazonReviews-1M	
Method	P@1	P@3	P@1	P@3
XPERT	5.10	2.84	3.99	2.33
NPPR	1.70	1.09	0.84	0.95
SUR-DNN	0.48	0.38	0.82	0.85
DPSR	2.10	1.45	2.57	1.55
PinnerSage	1.65	1.21	1.06	1.01

Table 8: Re-ranking after NPPR vs XPERT on the U2A-4M dataset. 1000 ads were retrieved using NPPR and re-ranked using XPERT by computing the cosine similarity with personalized embeddings of a user's 10 recent-most events. The top 100 ads were chosen. Even at 10x cost, this retrieve-then-rerank pipeline is outperformed by XPERT all by itself both in terms of recall and precision@1.

Method	Recall@100	P@1
NPPR -> XPERT	21.9	4.10
XPERT	27.1	5.10

6 CONCLUSION AND FUTURE WORK

This paper presented the concept of morphing operators that offer scalable yet effective personalization. Notably, this enables personalization at the retrieval stage itself without requiring a prior shortlisting to have taken place. The notion of morphing operators opens up several possibilities and whereas the XPERT algorithm only explores linear morphing operators, more powerful non-linear operators must be explored given the substantial improvements in retrieval yielded by XPERT's frugal architecture. XPERT can also be augmented to perform collaborative learning (e.g. via graph convolutional networks) by exploiting relational information in the form of user-user or item-item graphs. Alternate training strategies that improve gains on rarer items would be interesting and ensure that more items enjoy the benefits of PPR.

REFERENCES

- [1] Mingxiao An, Fangzhao Wu, Chuhan Wu, Kun Zhang, Zheng Liu, and Xing Xie. 2019. Neural news recommendation with long-and short-term user representations. In *ACL*.
- [2] K. Bhatia, K. Dahiya, H. Jain, A. Mittal, Y. Prabhu, and M. Varma. 2016. The Extreme Classification Repository: Multi-label Datasets & Code. <http://manikvarma.org/downloads/XC/XMLRepository.html>
- [3] T. Chen, S. Kornblith, M. Norouzi, and G. Hinton. 2020. A simple framework for contrastive learning of visual representations. In *ICML*.
- [4] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep neural networks for youtube recommendations. In *RecSys*.
- [5] K. Dahiya, A. Agarwal, D. Saini, K. Gururaj, J. Jiao, A. Singh, S. Agarwal, P. Kar, and M. Varma. 2021. SiameseXML: Siamese Networks meet Extreme Classifiers with 100M Labels. In *ICML*.
- [6] Alan Edelman, Tomas A. Arias, and Steven T. Smith. 1998. The Geometry of Algorithms with Orthogonality Constraints. *SIAM J. Matrix Anal. Appl.* 20, 2 (1998), 303–353.

- [7] F. Faghri, D.-J. Fleet, J.-R. Kiros, and S. Fidler. 2018. VSE++: Improving Visual-Semantic Embeddings with Hard Negatives. In *BMVC*.
- [8] Shaohua Fan, Junxiong Zhu, Xiaotian Han, Chuan Shi, Linmei Hu, Biyu Ma, and Yongliang Li. 2019. Metapath-guided heterogeneous graph neural network for intent recommendation. In *KDD*.
- [9] Jiafeng Guo, Yixing Fan, Qingyao Ai, and W. Bruce Croft. 2016. A Deep Relevance Matching Model for Ad-hoc Retrieval. In *CIKM*.
- [10] P. S. Huang, X. He, J. Gao, L. Deng, A. Acero, and L. Heck. 2013. Learning Deep Structured Semantic Models for Web Search using Clickthrough Data. In *CIKM*.
- [11] Houyi Li, Zhihong Chen, Chenliang Li, Rong Xiao, Hongbo Deng, Peng Zhang, Yongchao Liu, and Haihong Tang. 2021. Path-based Deep Network for Candidate Item Matching in Recommenders. In *SIGIR*.
- [12] Qiang Liu, Shu Wu, and Liang Wang. 2015. COT: Contextual Operating Tensor for Context-Aware Recommender Systems. In *AAAI*.
- [13] Zheng Liu, Jianxun Lian, Junhan Yang, Defu Lian, and Xing Xie. 2020. Octopus: Comprehensive and Elastic User Representation for the Generation of Recommendation Candidates. In *SIGIR*.
- [14] W. Lu, J. Jiao, and R. Zhang. 2020. TwinBERT: Distilling Knowledge to Twin-Structured Compressed BERT Models for Large-Scale Retrieval. In *CIKM*.
- [15] A. Y. Malkov and D. A. Yashunin. 2020. Efficient and robust approximate nearest neighbor search using Hierarchical Navigable Small World graphs. *TPAMI* (2020).
- [16] Tharun Kumar Reddy Medini, Qixuan Huang, Yiqiu Wang, Vijai Mohan, and Anshumali Shrivastava. 2019. Extreme classification in log memory using count-min sketch: A case study of amazon search with 50m products. In *NeurIPS*.
- [17] Bhaskar Mitra and Nick Craswell. 2019. An Introduction to Neural Information Retrieval. *Foundations and Trends in Information Retrieval* 13, 1 (2019), 1–126.
- [18] A. Mittal, K. Dahiya, S. Agrawal, D. Saini, S. Agarwal, P. Kar, and M. Varma. 2021. DECAF: Deep Extreme Classification with Label Features. In *WSDM*.
- [19] Jianmo Ni, Jiacheng Li, and Julian McAuley. 2019. Justifying recommendations using distantly-labeled reviews and fine-grained aspects. In *EMNLP-IJCNLP*.
- [20] Yabo Ni, Dan Ou, Shichen Liu, Xiang Li, Wenwu Ou, Anxiang Zeng, and Luo Si. 2018. Perceive your users in depth: Learning universal user representations from multiple e-commerce tasks. In *KDD*.
- [21] Rodrigo Nogueira, Wei Yang, Kyunghyun Cho, and Jimmy Lin. 2019. Multi-stage document ranking with BERT. *arXiv:1910.14424*.
- [22] 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 *KDD*.
- [23] Deepak Saini, Arnav Kumar Jain, Kushal Dave, Jian Jiao, Amit Singh, Ruofei Zhang, and Manik Varma. 2021. GalaXC: Graph neural networks with labelwise attention for extreme classification. In *WWW*.
- [24] V. Sanh, L. Debut, J. Chaumond, and T. Wolf. 2019. DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv:1910.01108*.
- [25] Fei Sun, Jun Liu, Jian Wu, Changhua Pei, Xiao Lin, Wenwu Ou, and Peng Jiang. 2019. BERT4Rec: Sequential recommendation with bidirectional encoder representations from transformer. In *CIKM*.
- [26] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is All you Need. In *NeurIPS*.
- [27] Tian Wang, Yuri M Brovman, and Sriganesh Madhvanath. 2021. Personalized Embedding-based e-Commerce Recommendations at eBay. *arXiv:2102.06156*.
- [28] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz, J. Davison, S. Shleifer, P. von Platen, C. Ma, Y. Jernite, J. Plu, C. Xu, T. Le Scao, S. Gugger, M. Drame, Q. Lhoest, and A. Rush. 2020. Transformers: State-of-the-Art Natural Language Processing. In *EMNLP: System Demonstrations*.
- [29] Chuhan Wu, Fangzhao Wu, Mingxiao An, Jianqiang Huang, Yongfeng Huang, and Xing Xie. 2019. Npa: Neural news recommendation with personalized attention. In *KDD*.
- [30] Chuhan Wu, Fangzhao Wu, Junxin Liu, Shaojian He, Yongfeng Huang, and Xing Xie. 2019. Neural demographic prediction using search query. In *WSDM*.
- [31] L. Xiong, C. Xiong, Y. Li, K.-F. Tang, J. Liu, P. Bennett, J. Ahmed, and A. Overwijk. 2021. Approximate nearest neighbor negative contrastive learning for dense text retrieval. In *ICLR*.
- [32] Han Zhang, Songlin Wang, Kang Zhang, Zhiling Tang, Yunjiang Jiang, Yun Xiao, Weipeng Yan, and Wen-Yun Yang. 2020. Towards personalized and semantic retrieval: An end-to-end solution for e-commerce search via embedding learning. In *SIGIR*.