

PinRec: Outcome-Conditioned, Multi-Token Generative Retrieval for Industry-Scale Recommendation Systems

Anirudhan Badrinath*
abadrinath@pinterest.com
Pinterest
Palo Alto, USA

Jaewon Yang
jaewonyang@pinterest.com
Pinterest
Palo Alto, USA

Prabhakar Agarwal*
pagarwal@pinterest.com
Pinterest
Palo Alto, USA

Jiajing Xu
jiajing@pinterest.com
Pinterest
Palo Alto, USA

Laksh Bhasin*
lbhasin@pinterest.com
Pinterest
Palo Alto, USA

Charles Rosenberg
crosenberg@pinterest.com
Pinterest
Palo Alto, USA

Abstract

Generative retrieval methods utilize generative sequential modeling techniques, such as transformers, to generate candidate items for recommender systems. These methods have demonstrated promising results in academic benchmarks, surpassing traditional retrieval models like two-tower architectures. However, current generative retrieval methods lack the scalability required for industrial recommender systems, and they are insufficiently flexible to satisfy the multiple metric requirements of modern systems.

This paper introduces PinRec, a novel generative retrieval model developed for applications at Pinterest. PinRec utilizes outcome-conditioned generation, enabling modelers to specify how to balance various outcome metrics, such as the number of saves and clicks, to effectively align with business goals and user exploration. Additionally, PinRec incorporates multi-token generation to enhance output diversity while optimizing generation. Our experiments demonstrate that PinRec can successfully balance performance, diversity, and efficiency, delivering a significant positive impact to users using generative models. This paper marks a significant milestone in generative retrieval, as it presents, to our knowledge, the first rigorous study on implementing generative retrieval at the scale of Pinterest.

CCS Concepts

- Information systems → Retrieval models and ranking; Personalization.

Keywords

generative retrieval, recommender systems, transformer, item recommendation

1 Introduction

Pinterest is a visual discovery platform where over 550 million monthly active users search for, save, and shop for ideas, represented by Pins [17]. Given the massive volume of Pins – for instance, users save over 1.5 billion Pins per week [18] – it is essential to provide personalized recommendations in Pinterest feeds and search results. To achieve this, recommender systems often leverage a two-stage approach: retrieval and ranking. Retrieval models select thousands of potentially relevant candidate items, then ranking

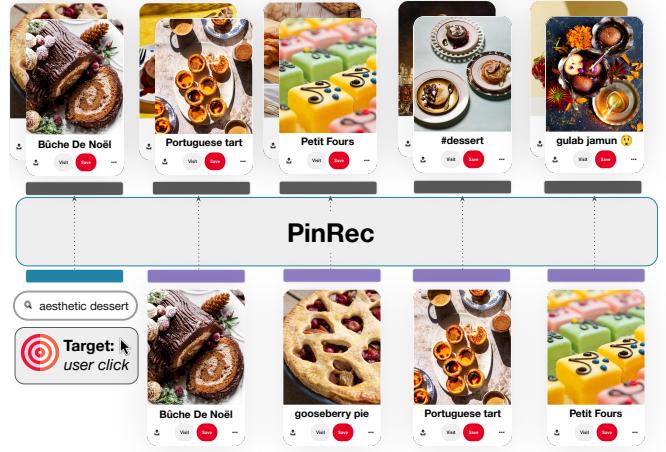


Figure 1: Illustration of PinRec, a generative item retrieval technique for heterogeneous user journeys on Pinterest. Sequences of user searches, engagements, and outcome-conditioning (bottom) are used to recommend Pins (top).

models prioritize the top items to display to the user. Since ranking models depend on the retrieval results, it is important for retrieval models to identify diverse, fresh, and relevant items.

Traditionally, retrieval methods employ “two-tower” models. The query tower computes query embeddings, while the item tower computes item embeddings. These models retrieve items with the highest similarity scores given a query embedding. As an alternative to two-tower models, recent work has proposed *generative retrieval*. Generative retrieval models leverage a sequence of tokens using generative models, such as transformers [25], and then use those tokens to retrieve items. Generative retrieval has shown promising results in academic settings given its powerful ability to understand user activity history and generate candidates that existing two-tower models may not be able to retrieve [10, 13, 30]. Nevertheless, to our knowledge, generative retrieval has not been attempted as successfully for large-scale industrial applications yet.

Using generative retrieval models for large scale recommender systems poses many challenges at the scale of Pinterest. Given their increased complexity and sequential generation, generating a set of recommended items requires higher computational cost

*Authors contributed equally to this research.

than computing an embedding with two-tower models. Further, the generative model should be capable of efficiently optimizing for different engagement outcomes because (a) business needs are constantly evolving to optimize combinations of different metrics and (b) as importantly, we wish for users to engage with content in the most optimal way for their growth on the platform. For instance, we would want to approach recommendations for a user with a propensity for low-effort engagement differently than for core users; through which recommendations could we enable the maximum progress in their user journey on Pinterest? To answer this question requires a malleable technique that can adapt to different circumstances, goals, and demands. Finally, the generated item sequences must possess sufficient diversity to facilitate users’ exploration of varied ideas; one-step generation of a single item representation may not suffice.

To our knowledge, there has been no large scale generative retrieval model that has addressed all of the aforementioned challenges. Most methods developed within academic settings [9, 10, 13] have not proven effective in meeting the latency and cost requirements necessary to handle web-scale traffic. While there are a few recent studies to apply generative retrieval to production, they focus on optimizing one objective [32] and do not demonstrate optimal optimization for different target outcomes.

In this paper, we present PinRec, a generative retrieval technique for real-world recommender systems. Unlike previous methods, PinRec addresses and resolves three crucial criteria for generation of optimal recommendations: multifaceted generative capabilities, computational efficiency, and item diversity. Our main contributions in the paper are summarized below.

- **Outcome-Conditioned Generative Retrieval:** Given evolving business needs and to allow for different types of user exploration, we propose multifaceted generative capabilities through *outcome-conditioned generation*. Across offline and online evaluation, we demonstrate that customizing item recommendations to the desired outcome allows for designers to provide recommendations that successfully result in those specified outcomes.
- **Windowed, Multi-Token Prediction:** We propose an efficient multi-token prediction objective tailored to the structure of the sequential recommendation task, where prior next-token objectives have often lacked representational power. Instead, by optimizing multiple per-step recommendations with respect to a diverse set of future targets, we achieve notable offline and online benefits in performance, diversity, and efficiency.
- **Efficient Serving:** We detail how we optimize autoregressive generation for optimal item recommendation, implement the serving infrastructure, and allow for horizontal scaling to serve hundreds of millions of users. To our knowledge, this study is the first to disclose the details of infrastructure optimization for industry-scale generative retrieval systems.

2 Related Work

2.1 Sequential Recommendation Systems

Sequential recommendation systems operate on a sequence of items with which a user has interacted [10]. Earlier work in academia has used recurrent neural networks (RNN) [8], but later work has often

adopted transformers [10, 24]. Recent studies have harnessed the extensive knowledge of pre-trained language models [9, 13].

In industrial applications, a common approach is to apply sequential encoders to learn user representations [16]. The user representation can then be used in traditional recommender systems (e.g., embedding-based retrieval or ranking models). Due to its seamless integration with existing recommendation models, this approach has been widely adopted industry-wide [7, 22, 26, 29].

Another approach in sequential recommendation systems is the generative approach, wherein the model learns to generate the items to recommend. Rajput et al. [20] used transformers to generate items by generating a sequence of semantic ID tokens that are obtained by Residual Quantization with Variational Auto Encoder (RQ-VAE) [12, 23]. Zhai et al. [32] augmented the sequence by adding heterogeneous items and also showed that the models can perform two tasks: ranking and retrieval. However, these studies have primarily concentrated on generating a single item rather than exploring the generation of sequences. In contrast, our research seeks to generate a sequence of multiple items to fully realize the potential of the generative approach. Our method is capable of recommending a diverse set of items that can drive various types of user actions. While similar approaches have been explored in academic settings [14], such strategies have not yet been widely adopted in industry.

2.2 Multi-Objective Generation

Given the multifaceted nature of user exploration and business needs on a platform like Pinterest, we explore generation techniques that consider multiple objectives. While reinforcement learning (RL) algorithms have demonstrated this capability in industry settings [27], research in language modeling indicates that RL often requires significant amounts of tuning for optimal generative capability [15]. Recent research has introduced more streamlined and scalable techniques grounded in supervised learning. These approaches employ goal-conditioned methods [5] or utilize prompting strategies [28]. Building on these, we develop an outcome-conditioned sequential recommendation system to serve billions of items and millions of users.

3 Problem Setup

In this section, we outline the problem setup for the task of a sequential recommendation system for a set of users \mathcal{U} , items \mathcal{I} , and actions \mathcal{A} . Our aim is to produce an ordered set of item recommendations $R(u, t) \subset \mathcal{I}$ of size N for a user $u \in \mathcal{U}$ at a given timestep t . Each item i has an item type $T(i)$, accounting for all kinds of items that the user can interact with (e.g., Pins). We define the set of all item types \mathcal{T} as $\mathcal{T} = \{T(i) \mid i \in \mathcal{I}\}$.

Conditioned on the user u , we denote their historical user activity sequence ending at timestep t_{\max} as $H(u, t_{\max})$, a chronologically-ordered heterogeneous sequence of items. For each item in the sequence $i_t \in H(u, t_{\max})$, the item i_t is interacted by user u at timestep $t \leq t_{\max}$ and its item type $T(i) \in \mathcal{T}$. We assume each interaction with an item i_t corresponds to an action $a(i_t) \in \mathcal{A}$.

To learn a sequential recommendation system, we propose a sequence model f_θ parameterized by θ . To optimize the model parameters θ , we leverage a set of features $\mathcal{F}(t)$ for each item type

$t \in \mathcal{T}$. These may either be learned end-to-end with the sequence model or provided as learned or raw signals. The sequence model is consequently optimized using these features across the set of user sequences $H(u, t_{\max})$ for all $u \in \mathcal{U}$.

4 PinRec

In this section, we discuss the architecture, methodology, and serving strategies for the proposed recommender system, PinRec. Leveraging multi-token, autoregressive generation with outcome conditioning, we present an expressive, diverse, and performant transformer-based generative retrieval technique. We further discuss how PinRec is served in real-time through an in-house implementation of GPT-2 powered by NVIDIA Triton.

4.1 Representation Learning for Items

A critical prerequisite for effective optimization of a sequential recommendation system is an effective input and output representation for the heterogeneous set of items. We denote the representation of some item $i \in \mathcal{I}$ of type $T_i = T(i)$ as $f_{T_i}(i)$, where there exist representation models $f_t : \mathcal{F}(t) \rightarrow \mathcal{R}$ for each item type $t \in \mathcal{T}$ mapping the input feature space to the chosen representation space \mathcal{R} . For our purposes, we consider two item types (\mathcal{T}): Pins and search queries, which represent two core entities at Pinterest.

While some work has often leveraged ID-based representation spaces [23], recent work by Liu et al. [14] and Yang et al. [30] demonstrates that semantic IDs can lead to poor performance due to representational collapse. To compound onto this, given a larger set of item types $|\mathcal{T}|$, it may be challenging to optimize representations across many different types in a compatible representational space. On the contrary, Agarwal et al. [1] observe that real-valued vector representations (i.e., $\mathcal{R} \subseteq \mathbb{R}^d$) generalize across item types.

For these reasons, we opt for a real-valued vector representation for all items. Based on a set of features $\mathcal{F}(t)$ for each item type $t \in \{\text{pin, search query}\}$, we construct MLPs f_t for each item type $t \in \mathcal{T}$ to generate a real-valued embedding. For simplicity, the architecture of each input embedder consists of 2-3 fully connected layers with ReLU activation and layer norms.

As input to each embedder, we featurize the item i given the chosen set of features $\mathcal{F}(T(i))$ and concatenate the features into a vector v_i . To featurize search queries, we use the trained OmniSearchSage query representation [1], and to featurize Pins, we use the trained PinSage representation [31]. Given that Pins are the core entity of Pinterest, we choose to implement a large, sharded set of hash-based embedding tables to store a learned ID embedding for each item. The output of the input embedder for item i is a real-valued, L_2 normalized embedding $\mathbf{i} = f_t(v_i)$.

4.2 Sequential Modeling Architecture

To construct a sequential recommendation system, we leverage a causal transformer decoder architecture. Provided a sequence of input embeddings, we repeatedly apply multi-head attention followed by feedforward layers. Critically, we apply a causal attention mask at each token, i.e., preventing tokens at time t from attending to any token at $t' > t$. The output of the transformer $f_\theta(H(u, t_{\max}))$ is a sequence of hidden states from the last layer, and we apply an output head to obtain representations of recommended items.

We base our transformer decoder architecture on the GPT-2 model. We use an in-house implementation that optimizes existing solutions (e.g., HuggingFace¹), while allowing for additional functionality. For example, we incorporate the surface for each user interaction and the raw timestamp, following Pancha et al. [16].

To generate the output representations for item recommendation, we construct an output head O as a small MLP. While the output head is typically conditioned on the transformer’s hidden state, its input can incorporate other forms of conditioning. To that end, we denote any “auxiliary conditioning” as $\Phi(H(u, t_{\max}))$, and we explore this further in Section 4.3.1. We normalize the outputs so they are in the same representational space as the item representations \mathbf{i} . To summarize, the output representation $\hat{\mathbf{i}}_{u,t}$ at timestep t for user u , applies the output head to the transformer’s last hidden state and any auxiliary conditioning (Equation 1).

$$\hat{\mathbf{i}}_{u,t} = O(f_\theta(H(u, t_{\max})), \Phi(H(u, t_{\max}))) \quad (1)$$

4.3 Optimizing Item Recommendations

In this section, we detail our approach to optimizing our item representations with a sequential objective. The standard choice to learn a distribution over personalized item recommendations is simply to employ a “next-token” prediction objective. Beyond directly taking into account the ground-truth user activity and engagement distribution, the next-token objective is effective, efficient, and ubiquitous. Following Pancha et al. [16], we employ a generic sampled softmax loss to optimize our output item representations against our next item targets, using in-batch and random negatives, and with temperature scaling. Since in-batch negatives are sampled from a biased distribution based on our sequences, we employ sample probability correction through estimating the conditional frequency of each of the items $Q(\mathbf{i}_c) = P(\mathbf{i}_c | \hat{\mathbf{i}}_{u,t})$ [2] with a count-min sketch [3].

Based on the vector representations of the output item representation $\hat{\mathbf{i}}_{u,t}$ and a candidate item \mathbf{i}_c , we define their similarity function $s(\hat{\mathbf{i}}_{u,t}, \mathbf{i}_c) = \lambda \cdot \hat{\mathbf{i}}_{u,t}^\top \mathbf{i}_c - Q(\mathbf{i}_c)$, applying a correction to account for sampling bias. The sampled softmax loss defined in Equation 2 maximizes the similarity of our output item representation $\hat{\mathbf{i}}_{u,t}$ and the next item representation $\mathbf{i}_{u,t+1}$, relative to in-batch and random negative examples N .

$$L_s(\hat{\mathbf{i}}_{u,t}, \mathbf{i}_{u,t+1}) = -\log \left(\frac{\exp(s(\hat{\mathbf{i}}_{u,t}, \mathbf{i}_{u,t+1}))}{\exp(s(\hat{\mathbf{i}}_{u,t}, \mathbf{i}_{u,t+1})) + \sum_{\mathbf{i}_n \in N} \exp(s(\hat{\mathbf{i}}_{u,t}, \mathbf{i}_n))} \right) \quad (2)$$

Our empirical next item loss function, over a user sequence $H(u, t_{\max})$ and sampled negatives N , maximizes next item similarity across the entire sequence, as shown in Equation 3 [16].

$$L_\theta^{\text{next}}(H(u, t_{\max})) = \frac{1}{|H(u, t_{\max})|} \sum_t L_s(\hat{\mathbf{i}}_{u,t}, \mathbf{i}_{u,t+1}) \quad (3)$$

To optimize the model parameters θ and consequently learn an optimal set of item recommendations, we minimize the empirical

¹<https://huggingface.co/openai-community/gpt2>

loss across user sequences and sampled negatives in our dataset \mathcal{D} .

$$\arg \min_{\theta} \mathbb{E}_{H(u, t_{\max}) \sim \mathcal{D}} [L_{\theta}^{\text{next}}(H(u, t_{\max}))] \quad (4)$$

4.3.1 Outcome-Conditioned Generation. Although we optimize with respect to item targets conditioned on the user’s history, the next item loss provides no additional contextual information for the output representation by itself. For instance, given a user who has a propensity for low-effort engagement, the standard next token objective tends to generate recommendations that reinforce this same behavior. Though prior approaches leverage outcome-based loss weighting to resolve this, it requires domain knowledge and nevertheless requires significant amounts of tuning to achieve the desired outcome for each user. Alternatively, RL approaches are fraught with training instability and/or inefficiency and a significant need for hyperparameter tuning [21]. On the contrary, we propose reframing this problem by simply conditioning the generation of item recommendations *directly* on the desired outcomes (i.e., outcome or goal-conditioning).

To that end, we additionally condition the output head O of the model on the desired outcome. To enable representing different outcomes, we employ a set of learnable embeddings for each possible action outcome (e.g., user repins, user clicks, etc.). By conditioning the output generation on a representation of the desired outcome (i.e., $\Phi(H(u, t_{\max}))$ represents $a_{\text{desired}} \in \mathcal{A}$), we can generate outcome-specific, personalized item representations $\hat{i}_{u,t}$.

At training time, we leverage user sequences to optimize the output representation with respect to the targets. Since we cannot dictate the desired outcome in offline data post-hoc, we optimize with respect to outcomes achieved in the data [5]. Hence, we can define the auxiliary conditioning during training time as $\Phi(H(u, t_{\max})) = a(i_{t+1})$, corresponding to the target item i_{t+1} .

During inference, we can specify the desirability of each outcome and condition our generation according to these preferences. When the system is deployed, this allows us to induce these desirable outcomes from users. Formally, we define a per-action budget distribution $\mathcal{B} : \mathcal{A} \rightarrow \mathbb{R}$ that maps each action to a non-negative budget such that $\sum_{a \in \mathcal{A}} \mathcal{B}(a) = 1$. Hence, as the budget for any action grows, it becomes more desirable. Based on the budget \mathcal{B} and the desired numbers of recommended items N , we can either deterministically or stochastically allocate a fixed number of items for each action during generation (e.g., based on the expectation, each action a can deterministically be allocated $\mathcal{B}(a) \cdot N$ items).

4.3.2 Windowed Multi-Token Generation. One of the principal assumptions in the next token objective is that there is strict ordering between the elements of the sequence. However, that reasoning is not necessarily applicable to social media platforms, where a user may engage with suggestions on their feed in an arbitrary order within a session. Moreover, a user may engage with certain items later on within a session (or perhaps at the very end) that they may have come across at the very beginning of the session. The next token objective simply ignores these crucial facets of recommendation systems for platforms such as Pinterest.

Motivated by the absence of strict ordering in sequential user engagement, we extend the next-token prediction objective. We assume that there exists a timestep window of size Δ within which user activity is roughly time-equivalent (i.e., their order within that

window is not significant). To address the latter issue of delayed user engagement of an item relative to the current timestep, we employ an efficient multi-token objective which enables prediction beyond the next token at subsequent timesteps $t' \geq t + 1$.

Often, multi-token prediction for language models involves predicting a handful (e.g., 5-10) of successive tokens at timesteps $t + 1$ to $t + k$ after the current token at timestep t [6]. However, this leads to a similar issue of penalizing any predictions that do not match the target at exactly some timestep $t' \in [t + 1, t + k]$, even if a matching target may exist within the window $[t + 1, t + \Delta]$. To incorporate the time-equivalent window size Δ , we propose a windowed multi-token objective. That is, for any future timestep window of size Δ starting at some timestep $t' > t$, all items that are engaged within that window are *targets* for our prediction task.

We choose to condition the output generation at timestep t on both the target timestep $t' > t$ and its action outcome (though it can work without outcome-conditioning). That is, we extend $\Phi(H(u, t_{\max}))$ to represent $\{t', a(i_{t'})\}$ at training time. Similarly to outcome-conditioning, we can decide t' (and the corresponding desired outcomes based on \mathcal{B}) at inference time as desired.

Our mathematical formulation for the windowed multi-token objective is shown in Equation 5, where $p_{t'}$ denotes the chosen distribution over target timestep t' given t (e.g., a uniform distribution from $[t + 1, t_{\text{lim}}]$ for some t_{lim}). Note that a non-zero likelihood of sampling $t' = t + 1$ additionally enables this multi-token objective to perform next-token prediction and autoregressive generation.

$$L_{\theta}^{\text{mt}}(H(u, t_{\max})) = \frac{1}{|H(u, t_{\max})|} \sum_t \mathbb{E}_{t' \sim p_{t'}(\cdot | t)} \left[\min\{L_s(\hat{i}_{u,t}, i_{u,i}) \mid \forall i \in [t', t' + \Delta]\} \right] \quad (5)$$

Intuitively, the optimization problem associated with the windowed multi-token objective can be posed as a question: can we sequentially predict any future items from the set of items that the user will engage with between timestep t' and $t' + \Delta$? We find that this optimization is more intuitively aligned with the properties of user interaction on Pinterest compared to next-token prediction, while allowing for significantly reduced serving latency due to multiple items being predicted at a single token.

4.4 Serving

PinRec is implemented with an in-house, optimized implementation of GPT-2, and served with NVidia Triton Inference Server’s ensemble setting and Python backend, to perform autoregressive generative retrieval. This multi-stage serving pipeline includes separate stages for fetching signals, featurizing Pins with learned ID embeddings, autoregressive inference for item generation, and Faiss retrieval (Figure 2).

4.4.1 Signal Fetching. Our Signal Service uses a Lambda architecture, which combines both batch and real-time processing. The batch signal is updated daily through a Spark pipeline that processes the past year’s positive user engagements and search queries, and stores the sequences and features for each user ID. The real-time component is backed by a RocksDB key-value store [19] that allows us to efficiently fetch user engagements since the end time of the batch signal. The batch and real-time signals are de-duplicated

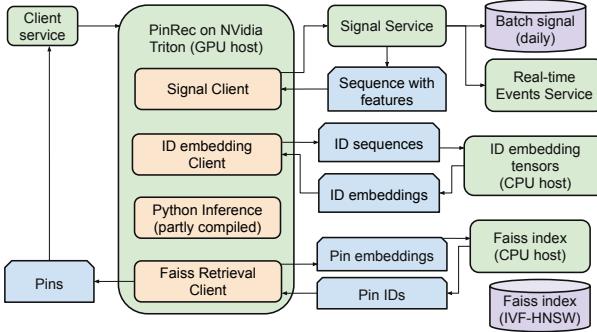


Figure 2: Serving flow for PinRec system (green boxes represent services, blue rectangles represent transmitted data, purple cylinders represent indexed data, and beige boxes are steps within the PinRec Nvidia Triton ensemble).

and returned to PinRec. On the search surface, model inference also includes the search query issued by the user at request time as the final element in the sequence. To limit transport costs, all embeddings associated with sequence elements are quantized to INT8 precision. These are later de-quantized to FP16 for inference.

4.4.2 Learned ID Embeddings. As discussed in Section 4.1, our pin embedder includes an ID embedding table. Effectively, it is represented as a large in-memory table that maps each pin ID’s hashes to an embedding, followed by some MLP layers. As this ID embedder takes up significant space in memory, especially when replicated, it is separated out from the main PinRec model and runs in its own CPU memory-optimized service. It is served with torchscript to limit latencies to single-digit milliseconds.

4.4.3 Autoregressive Item Generation. To generate an ordered set of item recommendations, we employ autoregressive generation on the transformer with a KV Cache. Depending upon whether the model is outcome-conditioned or multi-token, we are able to sample or generate multiple representations respectively. After autoregressive generation, we use Faiss to retrieve a set of recommended items of size N corresponding to the generated representations. Several optimizations are required for efficient inference:

CUDA Graphs. To limit overhead from kernel launches and include other in-built PyTorch optimizations, we use CUDA Graphs with `torch.compile` to maximally compile the main components of the model.

KV Cache. During the prefill stage of the transformer, the K and V tensors for the user’s known sequence are cached in a per-batch-request KV Cache. On subsequent decode steps, these are reused, and only new sequence elements are processed to produce K and V for the additional sequence dimension.

Compression. Given that we generate multiple output representations with autoregressive generation, we observe there are often embeddings that are similar to one another. During Faiss retrieval, these are likely to collapse into a similar set of output items. To optimize this, we “compress” each generated embedding into a

prior uncompressed embedding, in the order it is generated, if it is more similar than a given cosine similarity threshold. Compression significantly helps reduce the median number of embeddings generated by each model, with an average ratio of 6:1 for an unconditioned PinRec model and a smaller ratio of 3:1 for the multi-token, outcome-conditioned variant.

Since we wish to fetch N total items, we allocate these proportionally to each compressed embedding based on the total number of embeddings that were compressed into it. In the case of outcome-conditioning, we must perform this allocation independently for each action, in adherence with the chosen budget \mathcal{B} . The set of compressed embeddings alongside the allocations are passed to the Faiss retrieval stage.

4.4.4 Faiss Retrieval. Our nearest-neighbor retrieval involves a Faiss IVF-HNSW index running on CPU hosts, with inner product as the similarity metric. As with the main model, this is also served with NVidia Triton with batching enabled. We run kNN retrieval to fetch the allocated number of Pins for each embedding, with an overfetch factor for later de-duplication. Sample results are visualized for a manually specified user history in Figure 3.

5 Offline Evaluation

In this section, we dissect the performance of the proposed sequential recommendation system, PinRec, on generative retrieval across a range of offline evaluations. Our experimental setting includes interactions on three major surfaces for the Pinterest platform, namely Homefeed, Related Pins, and Search. We demonstrate that outcome-conditioned generation allows for improved controllability and that multi-token generation leads to significantly improved diversity, performance, and efficiency.

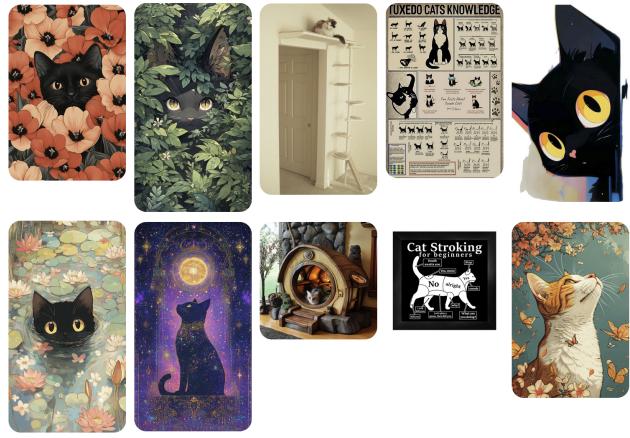
5.1 Offline Evaluation Metrics

We construct the evaluation set using interactions for a set of users disjoint from the training set and evaluate whether the proposed technique can produce a representation to retrieve future item embeddings among a set of 1M random items (“index”). In such tasks, recall @ k is a standard performance metric, defined as the proportion of target items correctly retrieved by the model within the top- k ranked items across an index. However, standard recall assumes the existence of a single target per sample. As mentioned in Section 4.3.2, such an approach fails to capture the assumption of a lack of strict ordering between user activity within a small range of interactions Δ . Further, when the goal of the retrieval model is candidate generation for a downstream ranking model, the task naturally has more than one unique plausible target, making standard recall less relevant.

We address these shortcomings by evaluating model performance on an unordered generalization of recall @ k . Given the autoregressively predicted item representations across all users $\hat{\mathbf{I}}$, target item embeddings across all users \mathbf{I} , and random item embeddings \mathbf{N} . Differently from standard recall, for each user $u \in \mathcal{U}$, unordered recall compares a set of predicted $\hat{\mathbf{I}}_u$ with a set of targets \mathbf{I}_u of fixed size Δ , calculating the proportion of targets \mathbf{I}_u where there exists some $\hat{\mathbf{I}}_u$ that is within the top- k ranked items. We



(a) Recent user history specified to produce results in Subfigure b.



(b) Results for the user sequence specified in Subfigure a.

Figure 3: Item recommendations from PinRec-OC for a user history with cat-related search queries and pin interactions. History in reverse chronological order. Additional examples in Section A.2.

Budget is: {repin: 0.4, outbound_click: 0.1, grid_click: 0.4, long_grid_click: 0.05, other: 0.05}. Generation runs for 6 steps (1 prefill, 5 decode).

formalize the notion of unordered recall in Equation 6.

$$\text{ur}_k(\hat{\mathbf{I}}, \mathbf{I}, \mathbf{N}) := \frac{1}{|\hat{\mathbf{I}}|} \times \sum_{i=1}^{|\hat{\mathbf{I}}|} \mathbf{1} \left\{ \min_{j, z \in [\Delta]} \left| \{ \mathbf{n} \in \mathbf{N} \mid \hat{\mathbf{I}}_{ij}^\top \mathbf{n} \geq \hat{\mathbf{I}}_{iz}^\top \mathbf{I}_{iz} \} \right| \leq k \right\} \quad (6)$$

Note that since we perform autoregressive generation without teacher forcing, the transformer is not conditioned on any of the targets at any point.

Beyond unordered retrieval performance, we examine the diversity of the recommendations. After retrieving the top- k predicted items for each user, we compute the proportion of unique items retrieved across the evaluation set. This allows us to quantify the degree to which learned representations are biased towards popular or frequently occurring Pins in the dataset.

5.2 Offline Results

In this section, we compare PinRec with an existing sequence-based user modeling representations, Pinnerformer [16], SASRec [11], and TIGER [20], in terms of offline evaluation metrics. Further, we quantify the impacts of our design on the performance and diversity of recommendations.

5.2.1 Comparison with Prior Method. To validate that PinRec outperforms existing sequential techniques in terms of item recommendation in offline evaluation, we compare its performance to PinnerFormer. We evaluate three different variants of PinRec.

- PinRec-UC: We do not condition on anything (i.e., $\Phi = \emptyset$) and do not employ multi-token prediction.
- PinRec-OC: We condition on the target action (“outcome conditioned”) and do not employ multi-token prediction.
- PinRec-{MT, OC}: We condition on the target action and employ multi-token prediction.

To ensure that the comparison between PinRec-{MT, OC} and other methods is fair, we adjust the number of autoregressive generation steps for multi-token generation to ensure that the total number of generated embeddings is held constant. As in training (where we cannot modify the evaluation dataset post-hoc), we condition on the target item’s action for outcome-conditioned variants.

Our comparison in Table 1 shows that PinRec outperforms a similarly sized Pinnerformer across the board. Importantly, outcome-conditioning demonstrates notable improvements over PinRec-UC (unconditioned) in unordered recall @ 10 (e.g., +2-3% on Homefeed and Related Pins). Above all, PinRec-{MT, OC} shows improvements over every method, with over **+10%** lift on Homefeed, over **+20%** on Related Pins, and **+30%** on search over PinRec-UC.

Table 1: Comparison of baselines and PinRec variants (unordered recall @ 10) across major surfaces at Pinterest.

Model	Homefeed	Related Pins	Search
TIGER [20]	0.208	0.230	0.090
SASRec [11]	0.382	0.426	0.142
PinnerFormer [16]	0.461	0.412	0.257
PinRec-UC	0.608	0.521	0.350
PinRec-OC	0.625	0.537	0.352
PinRec-{MT, OC}	0.676	0.631	0.450

5.2.2 Outcome-Conditioned Generation. We seek to understand whether we can achieve “controllability” between the desired action and true action; for instance, does conditioning on a “repin” generate an output representation that truly aligns with content that will be repinned? If we are able to control the outcome through our item recommendations, then we will observe that for the specific action that we condition on, there are significant gains in offline recall. Conversely, we expect that the recall for items corresponding to

other actions will decrease since our action-conditioned generation should ideally induce only that chosen action. As a baseline, we compare to the unconditioned variant of PinRec, which does not take into account the desired action.

For simplicity, we allocate the entire budget entirely to one action and perform offline evaluation conditioning across each action. We show the result of this analysis in Figure 4. Whenever we condition on an action outcome, we show notable gains over an unconditioned baseline for that action. Specifically, matching outcome conditioning yields 1.9% recall lift for true repins, 1.4% on true grid clicks, 4.8% on long grid clicks, and 6.2% on outbound clicks.

On the other hand, we observe decreases in recall compared to an unconditioned baseline if the true action does not correspond to the desired action, indicating that our conditioning achieves the desired result of producing outcome-specific representations. Especially for outbound clicks, we observe that conditioning on any other actions leads to massive drops in unordered recall.

	Repin	Grid Click	Long Grid Click	Outbound Click
Desired Action	1.9%	-8.6%	-7.2%	-35.2%
Grid Click	-4.3%	1.4%	-0.7%	-30.8%
Long Grid Click	-16.4%	-5.7%	4.8%	-32.6%
Outbound Click	-71.6%	-64.2%	-62.7%	6.2%
Actual Action	Repin	Grid Click	Long Grid Click	Outbound Click

Figure 4: Percentage lift in unordered recall for PinRec-OC over PinRec-UC when conditioning on the desired action, stratified by the actual action taken by the user.

5.2.3 Utility of Multi-Token Generation. Beyond the improvement in performance, we evaluate the benefits of multi-token generation in terms of efficiency, diversity, and performance of generations relative to the number of chosen generations per timestep. We examine the tradeoff between generating more representations per timestep versus longer autoregressive generation, while keeping the number of output embeddings constant. Since the strongest offline performance aside from the multi-token variant is by PinRec-OC, we consider it as a baseline. For simplicity, we average the unordered recall @ 10 across the three surfaces.

We show the tradeoff between performance, diversity, and efficiency with the number of embeddings we generate in each generation step using multi-token generation in Figure 5. While next-token generation (i.e., 1 embedding per step) is the most inefficient and yields the poorest performance (where we perform similarly to PinRec-OC), we observe significant benefits across the board in increasing the number of embeddings per step. At 16 generations per step, we demonstrate $\sim 10x$ reduction in latency over PinRec-OC, with improved unordered recall (+16.0%) and diversity (+21.3%).

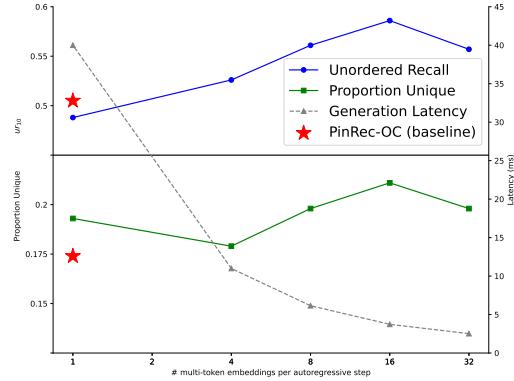


Figure 5: Tradeoff between unordered recall @ 10 (top), proportion unique (bottom), and generation latency (right) with the number of multi-token generations per step.

6 Online A/B Experiments

We perform online A/B experiments to evaluate PinRec as a candidate generator, measuring its impact on user engagement and confirming its online controllability. The items recommended by PinRec are passed through (unmodified) ranking and blending stages on each surface, alongside other candidate generators, to produce the final results shown to users. Since we employ a heterogeneous sequence of Pins and search queries, we focus on the Homefeed and Search surfaces. On these surfaces, we provide personalized item recommendations through PinRec, leading to significant lifts in action metrics, session metrics, and intent fulfillment.

6.1 Homefeed

6.1.1 Experiment Configuration. In the Homefeed online experiment, we run five groups exposed to global traffic for PinRec-OC, with varying action budgets \mathcal{B} . The budgets are fixed for outbound click (0.1), long grid clicks (0.05), and other unspecified actions (0.05). The remaining budget is split between the repin and grid click engagements, in steps of 0.2 ranging from 0 to 0.8. We similarly configure five groups for PinRec-{OC, MT}, with the same budgets. We include control groups without PinRec.

6.1.2 Metrics. Pinterest uses “session fulfillment” as an informative and holistic whole-funnel metric to capture (a) the upper funnel of user intent expression (e.g., when a user clicks on a pin in the grid), (b) the middle funnel of curating Pins (e.g., repin), and (c) the lower funnel of realizing intent (e.g., clicking on a product pin to visit its website). Any of the aforementioned actions results in a session being marked as “fulfilled”, whereas the absence of these results in an “unfulfilled” session.

6.1.3 Results. In Figure 6, we visualize a Pareto front of the repin rate (relative to impressions) and grid click rate for the PinRec candidate generator by itself, with absolute numbers hidden due to non-public information. By reducing the grid click budget and correspondingly increasing the repin budget, we see the expected result as the grid click rate drops and repin rate increases (i.e., we demonstrate online “controllability”). Similarly, we observe a significant correlation between the budgets and the overall A/B

experiment metric lifts for Homefeed grid clicks and repins in Figure 7. Further, we note that across both variants in Figure 7, for $\beta(\text{repin}) \geq 0.6$, there are statistically significant overall online lifts in Homefeed repin. We achieve statistically significant lifts in Homefeed grid click for all budget groups on each variant.

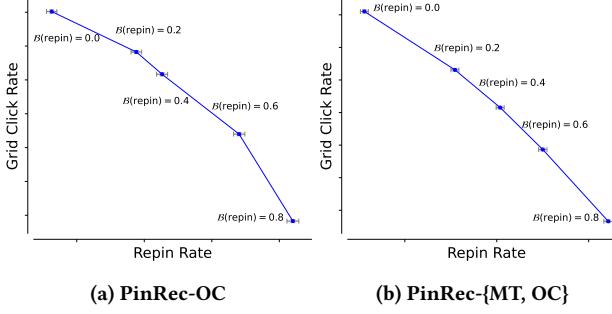


Figure 6: Visualization of Pareto fronts for online rate of grid click versus repin for items recommended only by PinRec variants, with error bars denoted in grey.

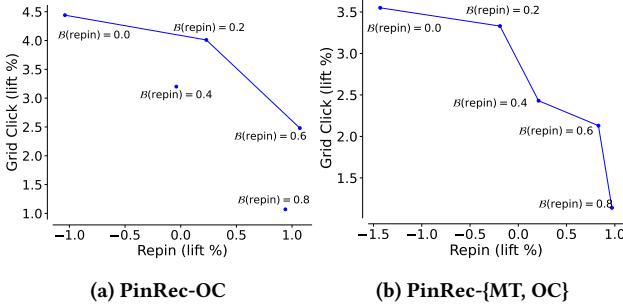


Figure 7: Visualization of Pareto fronts for online lift % in Homefeed grid clicks versus repin for PinRec variants.

To examine the holistic impact of our candidate generators, we examine the lift in fulfilled sessions for each PinRec variant, selecting the group which maximizes this metric. To demonstrate notable site-wide, session-level, and global metrics improvements, we present the CUPED-adjusted [4] percentage lift in metrics across Homefeed and the entire site (Table 2) for the groups that maximize fulfilled sessions. Relative to the control group, the unconditioned variant of PinRec (PinRec-UC) demonstrates improvements in unfulfilled sessions and in Homefeed and site-wide grid clicks. In addition to those lifts, the outcome-conditioned variant (PinRec-OC) shows significant lift in site-wide fulfilled sessions. The multi-token variant, PinRec-{MT, OC}, significantly increased successful site sessions, time spent on site, and site-wide grid clicks, which we believe represent the most comprehensive site-wide improvements.

Table 2: CUPED-adjusted metrics improvements from online A/B experiments for PinRec variants as candidate generators in Homefeed.

Type	Metric	UC	OC	{MT, OC}
Overall Metrics	Fulfilled Sessions	+0.02%	+0.21%	+0.28%
	Time Spent	-0.02%	+0.16%	+0.55%
	Unfulfilled Sessions	-0.28%	-0.89%	-0.78%
Action Metrics	Site-wide Grid Clicks	+0.58%	+1.76%	+1.73%
	Homefeed Grid Clicks	+1.87%	+4.01%	+3.33%

6.2 Search

On the Search surface, we evaluate the unconditioned variant of PinRec. As shown in Table 3, PinRec-UC improves search metrics significantly relative to a control group. We observe session-wide search “fulfillment” increase, i.e., when a search session results in positive interaction such as repins or long outbound clicks.

Table 3: CUPED-adjusted metrics improvements from online A/B experiments for PinRec as candidate generator in Search retrieval.

Metric Type	Metric	PinRec-UC
Overall Metrics	Search Fulfillment Rate	+0.48%
Action Metrics	Search Repins	+0.84%
	Search Grid Clicks	+1.46%

7 Conclusion

We introduced PinRec, a generative retrieval system developed for Pinterest’s Homefeed and Search feed. To the best of our knowledge, this paper represents the first comprehensive study on deploying generative retrieval in a web-scale application. The incorporation of two key innovations – outcome-conditioned generation and multi-token inference – enabled PinRec to enhance several business metrics without substantial cost increases.

Several promising future directions can be pursued. First, applying generative approaches to ranking, as explored by Zhai et al. [32], could be beneficial. Second, optimizing input sequences to the model, akin to prompt engineering in large language models, is another avenue worth exploring. Third, developing language model alignment techniques similar to Ouyang et al. [15] could enhance the use of engagement signals from users. Overall, this paper highlights a promising avenue for utilizing generative retrieval in large-scale recommender systems.

Acknowledgments

We thank collaborators such as Deepak Agarwal, Divyansh Agarwal, Edoardo Botta, Bee-Chung Chen, Bowen Deng, Pong Eksombatchai, Kurchi Subhra Hazra, Se Won Jang, Saurabh Vishwas Joshi, Krishna Kamath, Sujay Khandagale, James Li, Cristian Lopez, Matt Madrigal, Lei Pan, and Ruijia Wang for their critical contributions to our work and for facilitating productive discussions.

References

- [1] Prabhat Agarwal, Minhzul Islam SK, Nikil Pancha, Kurchi Subhra Hazra, Jiajing Xu, and Chuck Rosenberg. 2024. OmniSearchSage: Multi-Task Multi-Entity Embeddings for Pinterest Search. In *Companion Proceedings of the ACM Web Conference 2024*. Association for Computing Machinery, New York, NY, USA, 121–130. doi:10.1145/3589335.3648309
- [2] Yoshua Bengio and Jean-Sébastien Senecal. 2008. Adaptive Importance Sampling to Accelerate Training of a Neural Probabilistic Language Model. *IEEE Transactions on Neural Networks* 19, 4 (2008), 713–722. doi:10.1109/TNN.2007.912312
- [3] Graham Cormode and S. Muthukrishnan. 2005. An Improved Data Stream Summary: The Count-Min Sketch and Its Applications. *Journal of Algorithms* 55, 1 (April 2005), 58–75. doi:10.1016/j.jalgor.2003.12.001
- [4] Alex Deng, Ya Xu, Ron Kohavi, and Toby Walker. 2013. Improving the Sensitivity of Online Controlled Experiments by Utilizing Pre-Experiment Data. In *Proceedings of the Sixth ACM International Conference on Web Search and Data Mining*. Association for Computing Machinery, New York, NY, USA, 123–132. doi:10.1145/243396.2433413
- [5] Dibya Ghosh, Abhishek Gupta, Ashwin Reddy, Justin Fu, Coline Manon Devin, Benjamin Eysenbach, and Sergey Levine. 2021. Learning to Reach Goals via Iterated Supervised Learning. In *Proceedings of the 9th International Conference on Learning Representations*. OpenReview.net. <https://openreview.net/forum?id=rALA0Xo6yNj>
- [6] Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. 2024. Better & Faster Large Language Models via Multi-Token Prediction. In *Proceedings of the 41st International Conference on Machine Learning*. JMLR.org, Article 629, 29 pages. doi:10.1109/ICDM.2018.00035
- [7] Casper Hansen, Christian Hansen, Lucas Maystre, Rishabh Mehrotra, Brian Brost, Federico Tomasi, and Mounia Lalmas. 2020. Contextual and Sequential User Embeddings for Large-Scale Music Recommendation. In *Proceedings of the 14th ACM Conference on Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 53–62. doi:10.1145/3383313.3412248
- [8] Balázs Hidasi, Alexandros Karatzoglou, Linas Baltrunas, and Domonkos Tikk. 2016. Session-based Recommendations with Recurrent Neural Networks. In *4th International Conference on Learning Representations, ICLR, Conference Track Proceedings*. San Juan, Puerto Rico. arXiv:1511.06939
- [9] Yupeng Hou, Shanlei Mu, Wayne Xin Zhao, Yaliang Li, Bolin Ding, and Ji-Rong Wen. 2022. Towards Universal Sequence Representation Learning for Recommender Systems. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, 585–593. doi:10.1145/3534678.3539381
- [10] Wang-Cheng Kang and Julian J. McAuley. 2018. Self-Attentive Sequential Recommendation. In *2018 IEEE International Conference on Data Mining (ICDM)*. IEEE Computer Society, Los Alamitos, CA, USA, 197–206. doi:10.1109/ICDM.2018.00035
- [11] Wang-Cheng Kang and Julian McAuley. 2018. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*. IEEE, 197–206.
- [12] Doyup Lee, Chiheon Kim, Saehoon Kim, Minsu Cho, and Wook-Shin Han. 2022. Autoregressive Image Generation using Residual Quantization. In *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 11513–11522. doi:10.1109/CVPR52688.2022.01123
- [13] Jiacheng Li, Ming Wang, Jin Li, Jimmiao 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*. Association for Computing Machinery, New York, NY, USA, 1258–1267. doi:10.1145/3580305.3599519
- [14] Zihan Liu, Yupeng Hou, and Julian McAuley. 2024. Multi-Behavior Generative Recommendation. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*. Association for Computing Machinery, New York, NY, USA, 1575–1585. doi:10.1145/3627673.3679730
- [15] Long Ouyang, Jeff Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askelian, Peter Welinder, Paul Christiano, Jan Leike, and Ryan Lowe. 2022. Training Language Models to Follow Instructions with Human Feedback. In *Proceedings of the 36th International Conference on Neural Information Processing Systems*. Curran Associates Inc., Red Hook, NY, USA, Article 2011, 15 pages.
- [16] 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*. Association for Computing Machinery, New York, NY, USA, 3702–3712. doi:10.1145/3534678.3539156
- [17] Pinterest. 2024. *Pinterest Announces Fourth Quarter and Full Year 2024 Results*. Pinterest. <https://investor.pinterestinc.com/news-and-events/press-releases/press-releases-details/2025/Pinterest-Announces-Fourth-Quarter-and-Full-Year-2024-Results-Delivers-First-Billion-Dollar-Revenue-Quarter/default.aspx>
- [18] Pinterest. 2024. *Pinterest Internal Data, Global, June 2023. Weekly average over last twelve months*. Pinterest. <https://newsroom.pinterest.com/company/>
- [19] Rajath Prasad. 2024. *Building Pinterest's New Wide-Column Database Using RocksDB*. Pinterest. <https://medium.com/pinterest-engineering/building-pinterests-new-wide-column-database-using-rocksdb-f5277ee4e3d2> *Pinterest Engineering on Medium*.
- [20] Shashank Rajput, Nikhil Mehta, Anima Singh, Raghuandan Keshavan, Trung Vu, Lukasz Heidt, Lichan Hong, Yi Tay, Vinh Q. Tran, Jonah Samost, Maciej Kula, Ed H. Chi, and Maheswaran Sathiamoorthy. 2023. Recommender Systems with Generative Retrieval. In *Proceedings of the 37th International Conference on Neural Information Processing Systems*. Curran Associates Inc., Red Hook, NY, USA, Article 452, 17 pages. http://papers.nips.cc/paper_files/paper/2023/hash/20dcab0f14046a5cb6b2b61da9f13229-Abstract-Conference.html
- [21] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. 2017. Proximal Policy Optimization Algorithms. arXiv:1707.06347
- [22] Zihua Si, Lin Guan, Zhongxiang Sun, Xiaoxue Zang, Jing Lu, Yiqun Hui, Xingchao Cao, Zeyi Yang, Yichen Zheng, Dewei Leng, Kai Zheng, Chenbin Zhang, Yanan Niu, Yang Song, and Kun Gai. 2024. TWIN V2: Scaling Ultra-Long User Behavior Sequence Modeling for Enhanced CTR Prediction at Kuaishou. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management*. Association for Computing Machinery, New York, NY, USA, 4890–4897. doi:10.1145/3627673.3680030
- [23] Anima Singh, Trung Vu, Nikhil Mehta, Raghuandan Keshavan, Maheswaran Sathiamoorthy, Yilin Zheng, Lichan Hong, Lukasz Heldt, Li Wei, Devansh Tandon, Ed Chi, and Xinyang Yi. 2024. Better Generalization with Semantic IDs: A Case Study in Ranking for Recommendations. In *Proceedings of the 18th ACM Conference on Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 1039–1044. doi:10.1145/3640457.3688190
- [24] 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 *Proceedings of the 28th ACM International Conference on Information and Knowledge Management*. Association for Computing Machinery, New York, NY, USA, 1441–1450. doi:10.1145/3357384.3357895
- [25] 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 *Advances in Neural Information Processing Systems*, Vol. 30. Curran Associates, Inc. https://proceedings.neurips.cc/paper_files/paper/2017/file/3f5ee243547dee1fb0d53c1c4a845aa-Paper.pdf
- [26] Xue Xia, Pong Eksombatchai, Nikil Pancha, Dhruvil Deven Badani, Po-Wei Wang, Neng Gu, Saurabh Vishwas Joshi, Nazanin Farahpour, Zhiyuan Zhang, and Andrew Zhai. 2023. TransAct: Transformer-based Realtime User Action Model for Recommendation at Pinterest. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. Association for Computing Machinery, New York, NY, USA, 5249–5259. doi:10.1145/3580305.3599918
- [27] Xin Xin, Alexandros Karatzoglou, Ioannis Arapakis, and Joemon M. Jose. 2020. Self-Supervised Reinforcement Learning for Recommender Systems. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval*. Association for Computing Machinery, New York, NY, USA, 931–940. doi:10.1145/3397271.3401147
- [28] Xin Xin, Tiago Pimentel, Alexandros Karatzoglou, Pengjie Ren, Konstantina Christakopoulou, and Zhaochun Ren. 2022. Rethinking Reinforcement Learning for Recommendation: A Prompt Perspective. In *Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval*. Association for Computing Machinery, New York, NY, USA, 1347–1357. doi:10.1145/3477495.3531714
- [29] Jiajing Xu, Andrew Zhai, and Charles Rosenberg. 2022. Rethinking Personalized Ranking at Pinterest: An End-to-End Approach. In *Proceedings of the 16th ACM Conference on Recommender Systems*. Association for Computing Machinery, New York, NY, USA, 502–505. doi:10.1145/3523227.3547394
- [30] Liu Yang, Fabian Paisher, Kaveh Hassani, Jiacheng Li, Shuai Shao, Zhang Gabriel Li, Yun He, Xue Feng, Nima Noorshams, Sem Park, Bo Long, Robert D. Nowak, Xiaoli Gao, and Hamid Eghbalzadeh. 2024. Unifying Generative and Dense Retrieval for Sequential Recommendation. arXiv:2411.18814
- [31] Rex Ying, Ruining He, Kaifeng Chen, Pong Eksombatchai, William L. Hamilton, and Jure Leskovec. 2018. Graph Convolutional Neural Networks for Web-Scale Recommender Systems. In *Proceedings of the 24th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, New York, NY, USA, 974–983. doi:10.1145/3219819.3219890
- [32] Jiaqi Zhai, Lucy Liao, Xing Liu, Yueming Wang, Rui Li, Xuan Cao, Leon Gao, Zhaojie Gong, Fangda Gu, Jiayuan He, Yinghai Lu, and Yu Shi. 2024. Actions Speak Louder than Words: Trillion-Parameter Sequential Transducers for Generative Recommendations. In *Proceedings of the 41st International Conference on Machine Learning*, Vol. 235. Proceedings of Machine Learning Research, 58484–58509. <https://proceedings.mlr.press/v235/zhai24a.html>

A Appendix

A.1 Additional Analysis

In this section, we perform additional experiments to quantitatively and qualitatively show that the design decisions made for PinRec are justified and that it performs across the state-of-the-art architectures. Additionally, we provide more online analyses of the performance across different user cohorts and tail items.

A.1.1 HSTU Architecture. Though PinRec relies primarily on the well-known transformer architecture, we note that our technique is completely architecture-agnostic. Specifically, outcome conditioning and multi-token prediction can be applied regardless of the specific architecture (e.g., transformer blocks) or algorithms used (e.g., attention), allowing them to be generic for any future use case of generative retrieval despite architectural or algorithmic improvements. To demonstrate this, we reproduce significant gains over baselines using the HSTU architecture [32], as opposed to the multi-head attention-based causal transformer. Note that we do not leverage the ranking task or any sequence sampling explored in [32] as it is not relevant for our specific use case. For simplicity, we reproduce results on the outcome conditioned variant of PinRec only, as we believe similar results can be extrapolated to other variants.

The results are shown in Table 4, which clearly shows that HSTU is roughly 3% worse on HomeFeed (effectively negating the improvement with respect to the unconditioned transformer-based PinRec) and competitive on the related pins surface. However, there is a substantial drop in performance on the search surface, which is a critical component of the user experience at Pinterest. We believe that this may be related to specific architectural elements of HSTU.

Table 4: Comparison of HSTU and transformer architectures for PinRec, where PinRec-OC leverages the transformer as presented in the main text.

Model	Homefeed	Related Pins	Search
PinRec-OC	0.625	0.537	0.352
HSTU (OC)	0.596	0.539	0.179

Though we observe a sharp drop in the search surface, HSTU performs roughly on the same order of magnitude as the traditional transformer architecture, demonstrating that our ideas perform across different architectures. This validates that our ideas are robust to architectures across the board.

A.1.2 Ablating Negative Sampling. Following from Pancha et al. [16], we perform ablations on negative sampling, as it is a critical component of the chosen sampled softmax loss function. Though we ultimately end up with similar conclusions as Pancha et al. [16], we show these findings for maximum clarity regarding our design choices.

Random Negatives. Though in-batch negatives are a critical component of our negative corpus, we construct a random negative corpus consisting of a random set of possible items. Given that our retrieval problem is ultimately across a random selection of pins,

we wish for PinRec to show retrieval capability across this corpus. We perform several ablations with random negatives, with in-batch negatives constant at a non-zero value. For these analyses, we focus on Homefeed (the most critical surface) and present results aggregated across all surfaces as well.

We show results in Table 5, where there are moderate gains in unordered recall if there is a non-zero number of in-batch negatives present. Specifically, on the most critical surface of Homefeed, we see notable improvement over the baseline in terms of unordered recall.

Table 5: Percentage lift in unordered recall using 16K random negatives, relative to using only 8K random negatives.

Surface	% Lift in Recall @ 10
Homefeed	+2.5%
All	+1.3%

To summarize, while random negatives are certainly an important piece of the puzzle, they cannot substitute in-batch negatives. When used in conjunction, scaling the number of negatives yields notable improvement, especially on Homefeed.

In-Batch Negatives. As shown in the previous subsection, disabling in-batch negatives entirely leads to significantly poorer metrics, showing their importance. Can we improve performance offline by scaling the batch size and consequently, the number of in-batch negatives? To examine this, we fix the number of random negatives and ablate the number of in-batch negatives through the batch size.

Table 6: Percentage lift in unordered recall @ 10 through scaling the number of in-batch negatives (relative to baseline of 1x) with random negatives fixed.

# In-Batch Negs.	% Lift (HF)	% Lift (all)
1.5x	+4.2%	+3.3%
2x	+8.4%	+5.4%
12x	+18.3%	+12.1%

We show results in Table 6, which demonstrates a clear trend that shows that a greater number of in-batch negatives and batch size leads to better performance, especially for the Homefeed surface. These are similar in directional improvements to random negatives, but we observe greater improvements with in-batch negatives, which are likely more relevant as negative examples.

A.1.3 Ablating Parameter Scaling. One of the recent explorations within recommendation systems has been scaling to improve performance, in an attempt to replicate the success of these efforts for language modeling. To explore this hypothesis, we scale the number of parameters within PinRec across two fronts: the “vocabulary” and the model itself. While model scaling is clear and involves scaling up the embedding size, the number of transformer layers, and so on, scaling the vocabulary for an embedding-based generative retrieval technique involves improving the input representations.

Specifically, as mentioned in the main text, we implement “ID” embeddings for pins, that effectively serve as additional representation for the “pin vocabulary”.

Model Parameters. As presented in the main text, PinRec (all variants) have roughly 100 million transformer parameters, corresponding to the base form or size of GPT-2. To ablate the effect of the size of the transformer, we attempt to follow the different scales employed by GPT-2 (base, medium, large, XL) as they have shown success in terms of the factors by which the embedding size, number of layers, and number of heads have scaled. Note that due to computational constraints, we reduce the scaling in terms of the number of layers, which increases GPU memory footprint significantly. Note that we present lift in offline unordered recall @ 10 for the outcome conditioned variant of PinRec below, but we observe similar trends for the unconditioned model as well.

Table 7: Percentage lift in unordered recall through expanding the transformer at various GPT-2 scales relative to base (100 million parameters).

Model	% Lift in Recall @ 10
Medium	-0.8%
Large	-0.9%
XL	+2.2%

Though we observe some gains in Table 7 from scaling the number of transformer parameters to around 1 billion (XL-sized), these become intractable to serve in real time, requiring offline batch inference or significantly higher cost and more optimization. For the remaining sizes, we observe insignificant lifts (note that less than 1% in lift magnitude can be often explained through noise). Hence, we opt to not use larger sized transformers.

Vocabulary Parameters. As presented in the main text, the ID embeddings consist of roughly 10 billion parameters. These are a core component of the model since they provide additional representational power for the most important entity type at Pinterest: Pins. To ablate the importance of these ID embeddings, we examine the performance in offline recall across all surfaces with and without ID embeddings. For this analysis, we focus on Homefeed (the most critical surface) and present results aggregated across all surfaces as well in Table 8.

Table 8: Percentage lift through adding ID embeddings (10 billion parameters) over no ID embedding-based vocabulary.

Surface	% Lift in Recall @ 10
Homefeed	+14.0%
All	+14.4%

Evidently, the improvement in recall is substantial and indicates that the scaling of vocabulary to learn an embedding for each unique ID (roughly, as it is hash-based) rather than only content-based features is beneficial.

A.1.4 Tail User and Item Analysis. One of the potential drawbacks of outcome conditioning would be for “tail” or “fresh” users, i.e., those which have lesser activity or less propensity for high-effort engagement. Critically, our method intends to optimize for outcomes across all cohorts of users, including such users. In the same vein, we wish to optimize for “tail” items, which may receive a lesser number of impressions, in an effort to distribute content to improve corpus coverage.

User Cohort	Grid Clicks		Outbound CTs	
	UC	OC	UC	OC
New	+2.60%	+2.62%	+3.98%	+5.78%
Resurrected	+1.56%	+1.86%	+1.98%	+1.72%

(a) Online lift (%) for grid clicks and outbound clickthroughers for new and resurrected users.

Metric	UC	OC
Corpus Coverage (unique pins)	+6.2%	+17.3%
1-Imp Pin Rec. (tail items)	+5.8%	+7.7%

(b) Relative lift (vs. 2-tower model) for tail item and corpus coverage metrics.

Table 9: Online experimental results for different user cohorts and tail item metrics, across PinRec-UC and PinRec-OC.

In terms of new users, PinRec-UC/OC lift grid clicks and outbound clickthroughers (outbound CTs) in Table 9a. Critically, the lift for outcome conditioning-based modeling is either neutral or positive with respect to both control and the unconditioned variant of PinRec, demonstrating that outcome conditioning does not suffer on such cohorts of users. In fact, across new users, OC performs better on outbound CTs.

To examine performance for tail items, we compare to a production 2-tower model in an online setting in Table 9. On corpus coverage, adding PinRec-UC/OC surfaces 6.2%/17.3% more unique pins relative to the 2-tower. For long-tail items with 1-impression, PinRec-UC/OC achieve 5.8%/7.7% absolute lifts respectively, relative to the aforementioned 2-tower model. Evidently, across the board, PinRec surfaces more of the corpus and surfaces more tail items compared to control, with outcome conditioning showing even more significant lifts over control in both regards.

We believe the evidence suggests that our generative retrieval solution is equal to or better across inactive/new users and tail items relative to existing control, and perhaps more importantly, it shows that outcome conditioning can help in boosting performance in these criteria, which is one of its goals.

A.1.5 Serving Efficiency. One of the critical aspects of our proposed generative retrieval technique is an effort to maintain serving efficiency. With our custom infrastructure built from off-the-shelf serving technologies such as NVidia Triton Inference Server and optimization techniques like torch.compile, we achieve a system that can be served in real-time. In this section, we delve deeper into the memory and time requirements, comparing our method with prior methods deployed in production.

On the memory front, all variants of PinRec (unconditioned, outcome-conditioned with or without multi-token generation) use about 1.6GiB of GPU memory including KVCache and CUDA graphs, without any significant memory usage increase for multi-token prediction or outcome conditioning. Excluding the KVCache and CUDA graphs, we use roughly the same amount of memory (~250 MiB) relative to other techniques (i.e., two-tower models).

In terms of latency, at 80 QPS, the online and offline latency for PinRec-OC/OC, MT are p50 40ms/p90 65ms (6 steps). On the other hand, across only 3 steps for PinRec-UC, we observe a latency of p50 23ms and p90 46ms (note that this produces significantly less embeddings and less compressed embeddings than either of the above techniques, though it is lower latency). Though the latency is roughly 3-4x higher (p50) than a traditional two-tower solution, we note that parallelizing it with other RPCs results in <1% added end-to-end increase in latency.

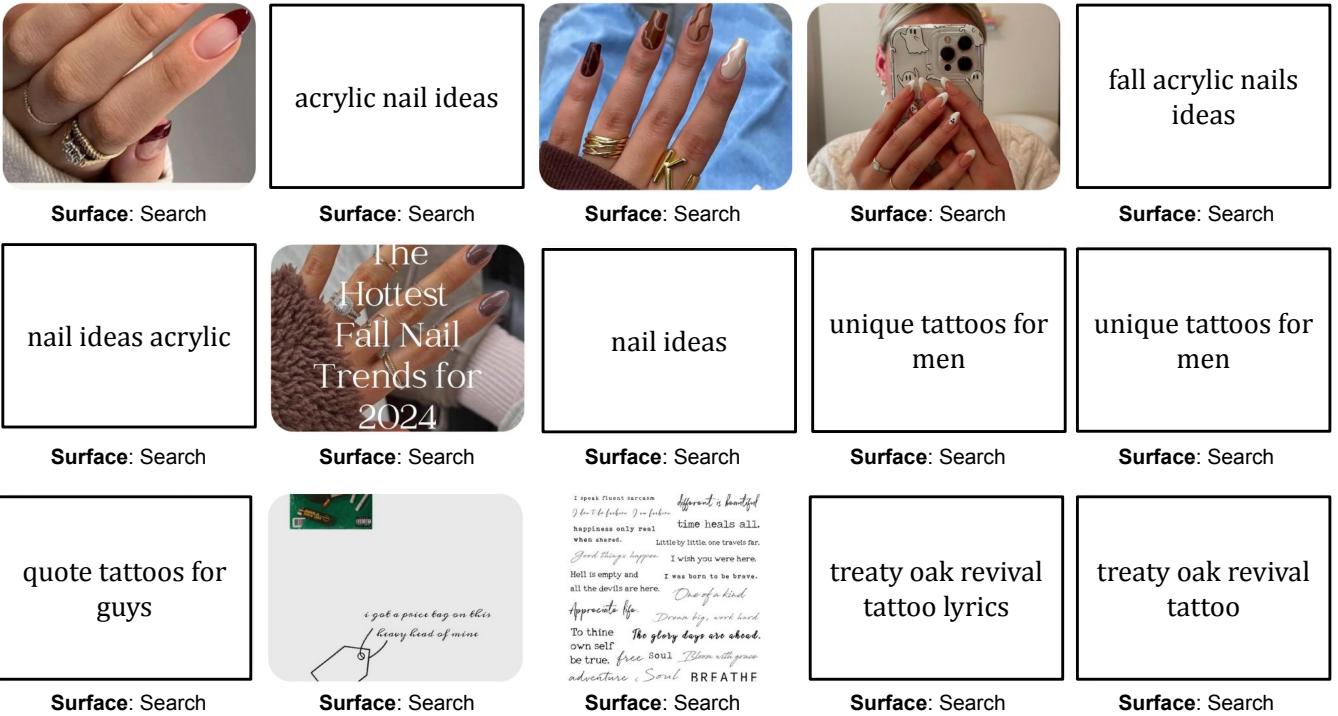
A.2 Visualization of Multi-Token Impact

In this section, we present some additional qualitative examples through visualizations of the retrieved examples for some sample

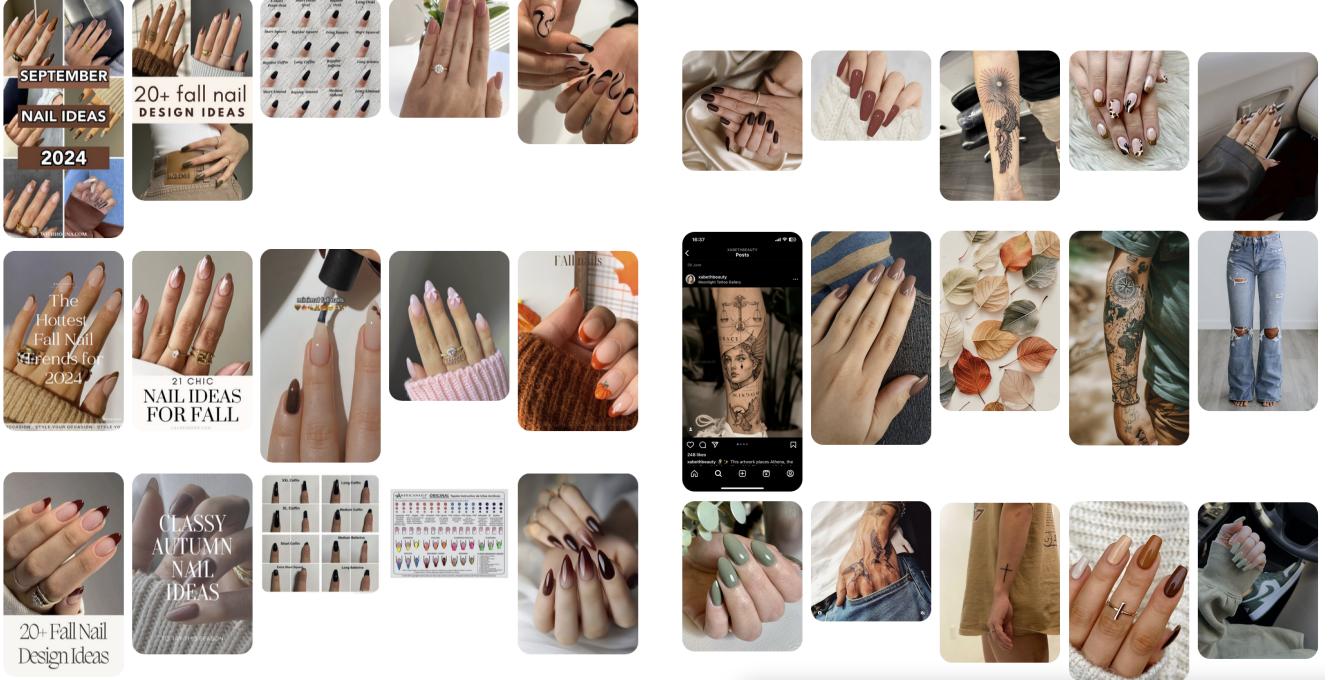
user histories across the windowed multi-token and non-multi-token variants. Through these visualizations, we are able to demonstrate a reasonable and consistent qualitative difference in the types of pins retrieved, justifying the improved session-wide and sitewide online metrics for PinRec-{MT, OC}.

Given a user history, we display two examples of retrieved pins in Figure 8 and Figure 9 for PinRec-{MT, OC} and PinRec-{OC}. Specifically, note that PinRec-{MT, OC} captures less-recent interactions for users compared to PinRec-OC. This is expected as the multi-token model predicts over a longer future time window and is less biased toward recency. PinRec-{MT, OC} can therefore produce more diverse results capturing more medium-term interests. On the other hand, since the PinRec-{OC} is naturally biased towards the short term, it fails to capture potential longer term trends that may induce engagement.

Note that these trends reflect the same takeaways from offline and online evaluation. Offline, we observe an improvement in diversity as the number of multi-token predictions is increased, which is consistently above the non-multi-token variant of PinRec. Online, we see that the “feed” presented by PinRec-{MT, OC} is of higher quality due to its diversity and visual appeal/aesthetic. These are critical qualities for a platform like Pinterest.



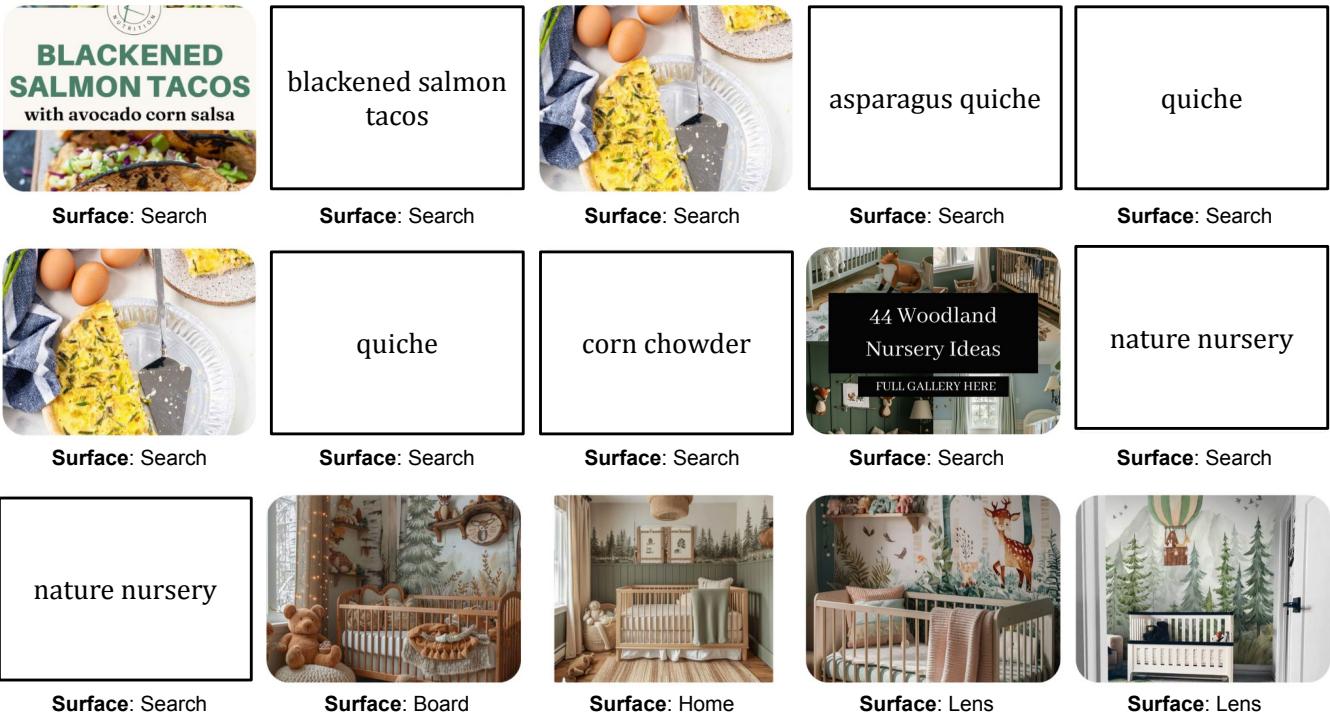
(a) Recent user history manually specified to produce results in subfigures. In reverse chronological order. Note a mix of nail and tattoo-related queries and engagements.



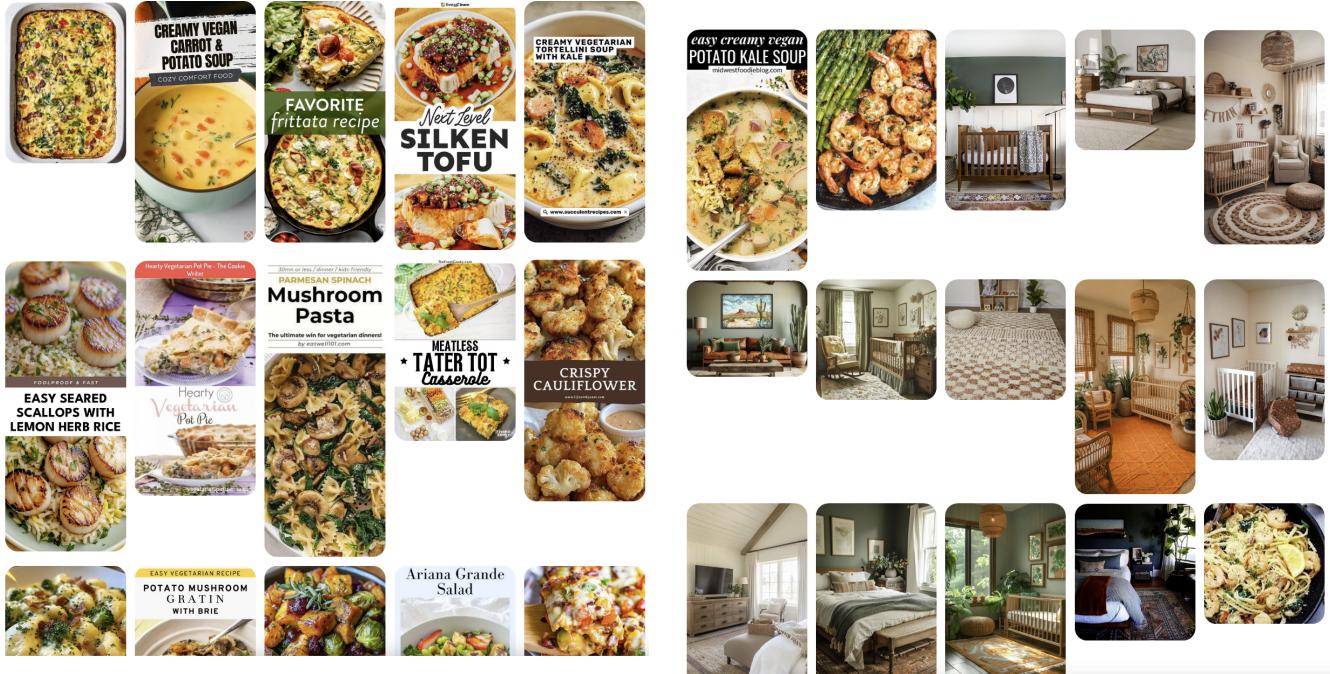
(b) Results for PinRec-OC candidate generator for user history in subfigure above. Notably only includes Pins related to nails, the user's most recent interaction.

(c) Results for PinRec-{MT, OC} candidate generator for user history in subfigure above. Includes a mix of Pins related to both tattoos and nails.

Figure 8: Comparison of results for a manually-specified user history.



(a) Recent user history manually specified to produce results in subfigures. In reverse chronological order. Note a mix of food and nursery-related queries and engagements.



(b) Results for PinRec-OC candidate generator for given user history.

Notably, the model is attending more to recent food-related interactions than earlier nursery-related interactions.

(c) Results for PinRec-{MT, OC} candidate generator for given user history. Includes a mix of Pins related to nurseries and food.

Figure 9: Comparison of results for another manually-specified user history.