

LiNR: Model Based Neural Retrieval on GPUs at LinkedIn

Fedor Borisyyuk*
Qingquan Song*
Mingzhou Zhou*
Ganesh

Parameswaran*
LinkedIn

Mountain View, CA, USA
fedorvb@gmail.com

Madhu Arun
Siva Popuri
Tugrul Bingol
LinkedIn

Mountain View, CA, USA
maarun@linkedin.com

Zhuotao Pei
Kuang-Hsuan Lee
Lu Zheng
LinkedIn

Mountain View, CA, USA
zpei@linkedin.com

Qizhan Shao
Ali Naqvi
Sen Zhou

Aman Gupta
LinkedIn

Mountain View, CA, USA
hshao@linkedin.com

Abstract

This paper introduces *LiNR*, LinkedIn’s large-scale, GPU-based retrieval system. *LiNR* supports a billion-sized index on GPU models. We discuss our experiences and challenges in creating scalable, differentiable search indexes using TensorFlow and PyTorch at production scale. In *LiNR*, both items and model weights are integrated into the model binary. Viewing index construction as a form of model training, we describe scaling our system for large indexes, incorporating full scans and efficient filtering. A key focus is on enabling attribute-based pre-filtering for exhaustive GPU searches, addressing the common challenge of post-filtering in KNN searches that often reduces system quality. We further provide multi-embedding retrieval algorithms and strategies for tackling cold start issues in retrieval. Our advancements in supporting larger indexes through quantization are also discussed. We believe *LiNR* represents one of the industry’s first Live-updated model-based retrieval indexes. Applied to out-of-network post recommendations on LinkedIn Feed, *LiNR* has contributed to a 3% relative increase in professional daily active users. We envisage *LiNR* as a step towards integrating retrieval and ranking into a single GPU model, simplifying complex infrastructures and enabling end-to-end optimization of the entire differentiable infrastructure through gradient descent.

CCS Concepts

• **Information systems** → **Similarity measures**; **Search engine indexing**; **Learning to rank**.

Keywords

information retrieval, recommender systems, candidate generation, nearest neighbor search, neural retrieval

ACM Reference Format:

Fedor Borisyyuk, Qingquan Song, Mingzhou Zhou, Ganesh Parameswaran, Madhu Arun, Siva Popuri, Tugrul Bingol, Zhuotao Pei, Kuang-Hsuan Lee, Lu Zheng, Qizhan Shao, Ali Naqvi, Sen Zhou, and Aman Gupta. 2024. *LiNR*:

*Authors contributed equally to this research.

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 the author(s) 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.
CIKM '24, October 21–25, 2024, Boise, ID, USA

© 2024 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 979-8-4007-0436-9/24/10
<https://doi.org/10.1145/3627673.3680091>

Model Based Neural Retrieval on GPUs at LinkedIn. In *Proceedings of the 33rd ACM International Conference on Information and Knowledge Management (CIKM '24), October 21–25, 2024, Boise, ID, USA*. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3627673.3680091>

1 Introduction

LinkedIn, the world’s largest professional network, serves over a billion members globally, offering services from job searches to content engagement. This paper explores *LiNR*, LinkedIn’s model-based GPU retrieval system, focusing on embedding-based retrieval (EBR). Traditional EBR uses unsupervised nearest neighbor search solutions [9, 14], indexing item vectors for fast retrieval. Our paper presents an innovative approach, combining exhaustive search with pre-filtering in a differentiable GPU model, using neural networks for distance learning and ranking. In *LiNR*, item vectors and model weights coexist within the same model binary, unlike traditional search indexing methods.

We believe the future of search and recommender systems lies in differentiable model-based serving, enabling joint optimization of retrieval and ranking. The K-nearest neighbor (KNN) search algorithm, an essential embedding-based retrieval method, uses learned query and item embeddings with a specific similarity metric to select the top-K closest items. Typically, KNN uses dot-product similarity, a form of matrix multiplication with normalized embeddings, which has been significantly sped up on modern GPUs (A100, H100, etc.) in frameworks like PyTorch and TensorFlow. Several challenges motivate us to propose model-based KNN algorithms implemented on GPUs:

- **Liquidity challenge:** Real-time search systems rely on specific attributes to filter relevant items. In job recommendation systems, for example, filters like company names, locations, and skills are essential. Items meeting these conditions must be prioritized to avoid exclusion due to low KNN scores from embeddings alone.
- **Low latency requirement:** Reducing retrieval latency and increasing throughput is a constant priority.
- **Huge memory cost:** As number item embeddings and clauses increase, finding ways to lower memory usage and boost computational speed without compromising retrieval quality presents a significant challenge.
- **Freshness:** Demonstrate that model-based approaches can enhance traditional nearest neighbor searches in quality and latency while supporting functionalities like live updates.

In this paper we discuss deployment of large-scale, neural model-based retrieval system, highlighting key challenges and solutions. A

major challenge was the absence of efficient pre-filtering in PyTorch and TensorFlow, addressed by our custom indexing and filtering methods detailed in §3.1, which also tackle latency issues. We also cover memory cost management through quantization techniques for larger indexes in §3.2. *LiNR* enhanced search quality, utilizing multi-embedding retrieval algorithms discussed in §3.3. Our work positions us among the pioneers in the industry in introducing a retrieval model-based serving infrastructure (§4.2), showcasing the capability of such model-based retrieval systems to be effectively live-updated at scale (§4.3). We perform our study of model-based index serving focuses on interest-based recommendations on LinkedIn’s Feed, also known as out-of-network (OON) recommendations. These recommendations leverage member profiles and previous interactions with the Feed, enabling LinkedIn members to access highly relevant content. We integrate OON content into various LinkedIn surfaces, like Feed and Notifications, based on predicted user engagement likelihood. The effectiveness of OON recommendations is gauged by member interactions with OON content. We use two-tower neural networks to create embeddings for members and Feed Posts, forming a candidate selection vertical for OON in the Feed through EBR with a differentiable model-based search index. *LiNR* significantly outperforms FAISS-based [8] retrieval system in OON recommendations. We support full-scan model-based index serving on GPUs with latencies as low as 4 ms, handling indexes from 15 million to a billion entries. This capability, along with modeling enhancements, significantly boosts quality as detailed in §5.2.

2 Related Work

Industry focus has predominantly been on approximate neighbor search systems, with FAISS [8], ScaNN [4], SONG [24], RAFT [15] among notable examples. These support algorithms like HNSW [14], IVFPQ [9], CAGRA [16] on CPU and GPU platforms. Termed model-free, these methods use unsupervised algorithms for partitioning space using existing item embeddings, offering flexibility for any item set. In contrast, our approach employs deep neural networks for a model-based search index, fully operational on GPUs. We integrate item indexes with neural network weights within a PyTorch or TensorFlow model, training during index construction and using the model for retrieval.

Recently with more performance and memory available on GPUs several publications have appeared considering model based nearest neighbor search such as [18, 21–23]. Mixture of logits (MoL) [22] in its production deployed form implements weighted combination of cosine similarities with neural network gates used to infer per distance component weights. The MoL paper does not provide information on examples of implementation of logits components, and which embeddings have been used in production. We extend on top of MoL and introduce practical algorithms on how to learn components of MoL. Conversely, research by [18, 21, 23] has explored using transformers and generative techniques for search indexes. Unlike our system, which stores item embeddings directly, these studies create semantic structures through clustering and transformers to generate document IDs.

A lot of research has been focused on representation learning with works representing posts and users in social networks

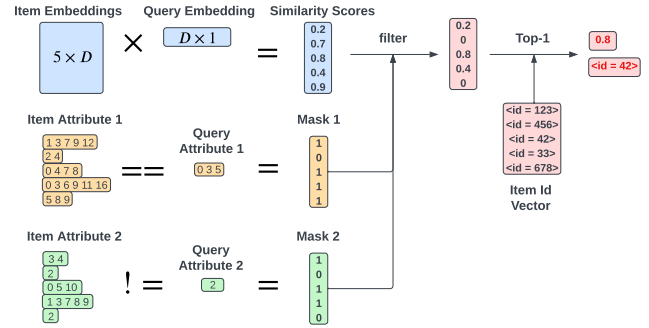


Figure 1: KNN with Similarity Masking. An example of five items with single query is used for illustration. Item similarities are computed and masked with two 0-1 vectors returned from the two clause checking. For each item, as long as one attribute is matched with the query attribute, the clause checking is passed (return one) in the masking matrix. The 2nd clause is a reverse matching clause. Top-1 selection is used in this example. D is dimension of item embedding.

[12, 17, 19]. As one of the components in MoL we have used approaches similar to [12, 17, 19], and additionally extended it with approaches for cold start infrequent users using clustering representations. Several previous works have explored the concept of model live-updates, which we expand on in this paper. These works include Monolith [13], PERSIA [11], and XDL [7]. In contrast, traditional search engines, as seen in Facebook Search EBR [6, 12] via Unicorn [3], and Lucene [2], have primarily focused on live-update functionality for unsupervised indexing techniques such as [8]. To the best of our knowledge, our paper represents one of the pioneering efforts in the realm of retrieval-based techniques for live-updating TensorFlow (TF) or PyTorch model-based retrieval indexes at a large-scale production level, with high QPS demands.

3 Modeling Technology

In this section we will describe how we modeled and developed exhaustive embedding-based search on GPU with attribute-based matching. We will provide details on how we scaled our model-based index to billion size on a single GPU with quantization. We extend Mixture of Logits (MoL) [22] by automatically training cluster embedding components and experimenting with different gating functions and variety of embedding components.

3.1 Exhaustive Search with Attribute-Based Matching (ABM)

Considering the post-filtering (filter after similarity-based retrieval) often suffers from the liquidity issue especially combining with ANN algorithms implemented on GPUs [25], we first focus on the KNN-based algorithm with attribute-based pre-filtering and introduce several basic approaches adopted to tackle the above challenges. Strategies to further improve the algorithm and tackle other online serving challenges including the live update problem will be introduced in §4.

3.1.1 KNN with Similarity Masking. Our first KNN algorithm with ABM is a two-step similarity masking approach. As shown in Figure 1, given a query, we first compute the similarity between the query embedding with all item embeddings stored in a matrix to

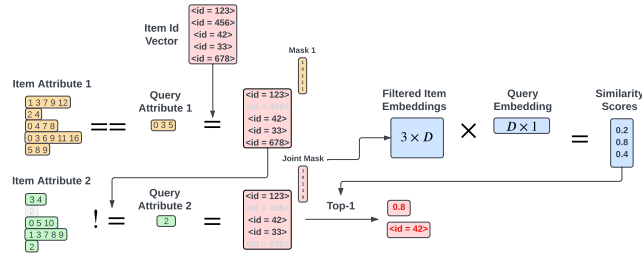


Figure 2: KNN with Explicit Pre-Filtering. Clauses are checked one by one and a joint 0-1 mask vector is returned to retrieve the feasible items for matrix multiplication and top-K selection ($K=1$ here).

capture their semantic relationships in a similarity vector. Then, we filter out irrelevant items by multiplying it with the 0-1 mask vectors given by each clause to map the similarity scores of the filtered items to zero before the top-K selection. Each query clause could contain multiple attributes. Feasible items should satisfy all clauses, requiring at least one of the attribute in each clauses is matched. Reverse clauses are also supported (such as the company name attributes in Figure 1). As each item could contain different number of attributes for each clause, to effectively utilize the GPU memory for saving and update the clauses, we store all clause attributes in a single matrix in practice and have an extra counting matrix to record the number of attributes for each item in each clause similar to the counting matrix in a CSR format but for each item separately without having the indexing vector. Each item clause is sorted before the concatenation for faster judgement (as we can stop checking early as long as one attribute is matched). We implement the algorithm in CUDA and registered the clause filtering kernel as TensorFlow and PyTorch operations to integrate and serve with other modules.

3.1.2 KNN with Explicit Pre-Filtering. The second iteration of our KNN with ABM uses a new approach, incorporating explicit filtering before embedding multiplication, as illustrated in Figure 2. Initially, we slice the matrix to filter out irrelevant items, removing them early from subsequent computations. This method speeds up the process by reducing the computational burden during matrix multiplication and top-K selection, especially beneficial when the query filters result in a significantly smaller item set. We found that with custom CUDA implementation [20] to merge the kernels, the speed could be generally faster than the first version introduced above. However, without customizing the masked matrix multiplication and kernel merging, simply adopting the matrix slicing in TensorFlow and PyTorch will introduce extra matrix copy and creation overhead, causing it to be slower than the first version when the pass rate is high.

3.2 Quantized KNN

Addressing memory constraints, we adopt a quantized KNN strategy using the Sign One Permutation One Random Projection (Sign-OPORP) method to compress embeddings to 1-bit and approximate dot-products via bitwise matching. This technique balances prediction accuracy and search speed, akin to typical ANN methods, but as an exhaustive search, it seamlessly integrates with attribute-based pre-filtering, circumventing liquidity issues. OPORP is a variant of

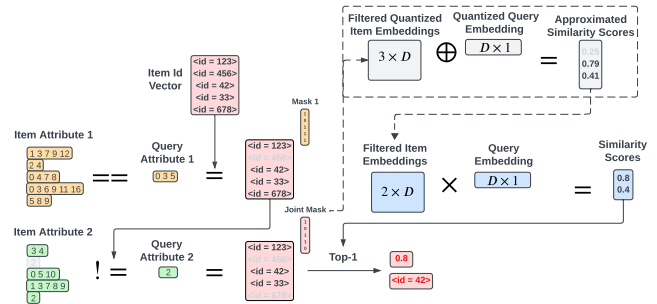


Figure 3: KNN with Quantized Filtering helps to reduce the number of retrieved items before the full precision similarity computation. A bit-wise matching is used to measure the approximated similarity between 1-bit quantized embedding obtained via Sign-OPORP method. We use bit-wise XOR operation and perform an integer bit-wise NOT conversion for query or item embedding in advance to measure the number of matched bits in the packed integer vector. The quantized KNN module can be used without full precision matrix multiplication when K is large in top-K selection.

count-sketch method. It leverages single random projection with fixed-length binning scheme to efficiently project embedding to a low-dimensional embedding. Sign-OPORP takes the sign of the projected embedding to generate 1-bit embedding that could accurately approximate the cosine similarity of the original floating-point embedding [10] via bit-wise matching, i.e., counting the number of matched bits of two quantized 1-bit embedding.

As the bit-wise matching operation is often much faster than regular matrix multiplication, we can replace the original embedding with quantized embedding and adopt the bit-wise matching operations in the above-mentioned KNN algorithm with pre-filtering, which can help greatly reduce the memory consumption. Compressing 1 billion fp16 embedding of dimension 64 to 1-bit embedding of the same dimension can reduce the memory by 16 times and help serve 1 billion items in single V100 GPU. We could adjust the size of the quantized embedding to balance the trade-off between the memory/speed and accuracy. Besides, if memory is not the concern, we could leverage the approximated similarity as an extra pre-filtering step reduce the computation of the full-precision matrix multiplication (see Figure 3), offering a unique perspective on exhaustive KNN with ABM. Note that, we call it exhaustive KNN to discriminate it from the regular ANN method with clustering such as HNSW and IVFPQ since our approach still computes all item similarity based on the quantized embeddings, which is easier to be combined with the pre-filtered ABM and live updates (see §4.3).

3.3 Similarity Modeling

3.3.1 Hadamard MLP. Dot Product or cosine similarity has been common in retrieval and it's computationally efficient. On the other hand, the multilayer perceptron (MLP)-based learned similarity functions has been reported inferior compared to properly tuned dot product. To balance the computation cost/latency and retrieval metrics, we attempted to boost the MLP-based learned similarity function through hadamard product. The architecture is shown on the left of Figure 4. A MLP block is applied to member and item embedding respectively, whose output performs hadamard product

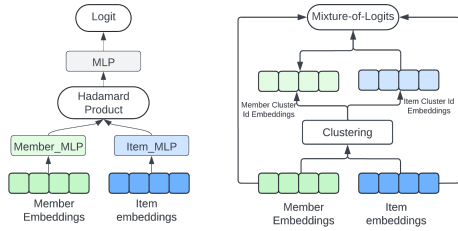


Figure 4: Illustration of Hadamard MLP (left) and learning cluster id embedding for Mixture-of-Logits(right)

and then passes to another MLP block to output the final logit. With proper hyper-parameter tuning, it can reliably outperform dot product.

3.3.2 Mixture-of-Logits with Clustering. Mixture-of-logits [22] defines a model for computing high rank similarity based on adaptive gating of elementary logits across multiple embedding components $\phi_{MoL}(x, u) = \sum_k \pi_{k,\theta}(x, u)\delta_{k,\theta}(x, u)$, where $\pi_{k,\theta}(x, u)$ represents a learnt gating function, which gives per component weight using soft-max gate given input of user and item features. The parameters θ are learnt through Adam optimization of gradients of a sampled soft-max loss.

Mixture-of-logits requires the availability of multiple features to leverage the gates, because the gates will collapse to a value of 1 if there is only one feature for user and item pair. We augment the feature with learnt cluster id embedding that obviate the necessity for having multiple features. In §5.1 we show that learnt cluster id embedding leveraged through Mixture-of-Logits can significantly improve on top of dot product in production settings.

Across LinkedIn we observed variety of member behaviour with some members coming frequently and some coming from time to time. For the infrequent members we aimed to improve retrieval system performance. To achieve this we learn cluster id embeddings, which represent interests of cohorts of members and topics of posts. We describe the process on the right of Figure 4.

For training *LiNR*, we obtain two-tower embeddings for posts and members as part of the training data, along with available engagement labels. We initialize cluster ID embeddings using K-means on millions of post embeddings. During training for both members and posts, we find the closest cluster ID embedding based on cosine similarity to their two-tower embedding. These cluster IDs for members and posts are integrated into Mixture-of-Logits, along with the original two-tower embedding and other embeddings we developed for our use cases. We experimented with using K-means-initialized cluster ID embeddings as is and fine-tuning them through back propagation. We report the experiment results in §5.1.

4 System Architecture

4.1 Out-of-Network Recommendations

Out of Network Recommendations is one of the many sources (first pass rankers) of LinkedIn Homepage Feed. When a member visits LinkedIn feed, a request is triggered from the front end and sent to feed service. Feed service passes this request to many first pass rankers including feed-OON mid tier (a.k.a. interest discovery).

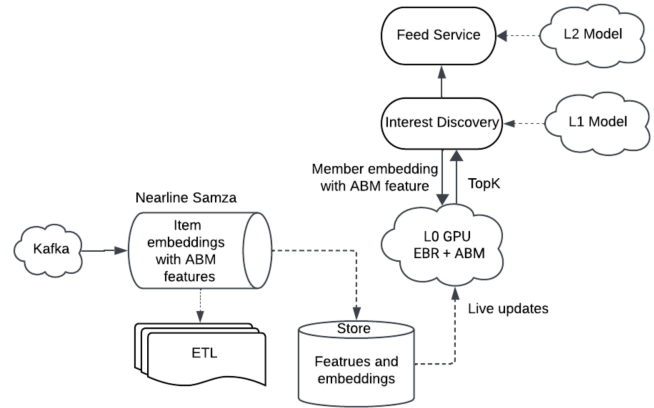


Figure 5: Feed OON Architecture.

This service is responsible for retrieving the top-K most eligible items for the member to send back to feed. Today, the underlying index used for retrieval is a lucene based index. The runtime of the query of OON is depicted at Figure 5. For every member query, a embedding based search is performed across all eligible item embeddings, followed by a layer 1 (L1) ranking model, which decides top-K items. These are then sent to feed service and ranked by more sophisticated layer 2 (L2) ranking model for members consumption. *LiNR* aims to provide an online service to run model based retrieval algorithms that can outperform our baselines: (1) dot-product based EBR, and (2) FAISS-IVFPQ, which is supported at LinkedIn for lucene systems.

As shown in the figure, interest-discovery will call model-cloud-L0 to fetch candidate items for the member. Model-cloud-L0 hosts the RAR model that does (1) item attribute-based filtering (2) embedding based retrieval with ranking using model. The model consists of the item embeddings, features needed for filtering and the trained model weights. Item embeddings are generated on a nearline fashion as and when a document is created at LinkedIn so as to keep the index up to date. The filters required for filtering are also ingested nearline.

4.2 ML Infra Architecture

We enhance Model Cloud, our hosted solution for serving model inferences, to support retrieval as ranking as shown in Figure 5.

4.2.1 Retriever. This component performs attribute-based filtering and embedding-based retrieval of the top-k documents for a query. At startup, retriever initializes with the retrieval model and bootstrapped data. Its framework-agnostic design allows easy extension to any framework, such as Torch or TensorFlow. AI engineers can experiment with new methods by developing and deploying corresponding models to this system.

4.2.2 Ingestor. Model-based retrieval requires the entire document corpus to reside in GPU memory for low latency. To provide fresh results, this corpus must be updated near real-time (nearline). Several following components work together to achieve this functionality.

Index Store: Attributes and embeddings come from offline sources and nearline data streams. We use Apache Beam to join and transform feature data for the entire document corpus. Offline, the full

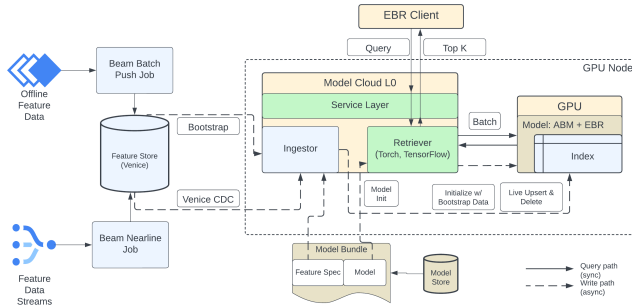


Figure 6: Model Cloud LiNR Architecture

corpus is batch-pushed to a Venice Store. Nearline updates are also written to this Venice Store.

Updater: Updater subscribes to the Index Store’s Change-Data-Capture (CDC) Stream. As the feature data gets batch pushed and live updated, the Updater gets notified to further process them and write to the model.

Bootstrapper: At startup, the Ingestor bootstraps from the Venice CDC client by replaying all data from the beginning. The entire data corpus is transformed into the required format and copied to the GPU. To minimize bootstrapping time, we regularly compact the bootstrap data and store a snapshot on disk for a fast warm start.

4.2.3 Service. To meet our performance needs, we avoid the latency and unpredictability of managed, garbage-collected languages. We also minimize network hops, data copies, and transformations. Our Model Cloud L0 service is written in a native language with minimal data transformations. User queries from the L0 client land directly on our service, ensuring we meet latency requirements.

4.3 Model Live Update

Live Update Ingestor subscribes to Venice CDC [5] from the bootstrapped offset, classifying changes into upserts and deletes, then transforming and copying them to the GPU.

The system’s effectiveness depends on the quality of the document index, which must remain fresh. This can be done by either regularly rebuilding the index or updating it in near-real-time via a data change stream. We chose the latter for two reasons: it keeps the corpus current, reflecting changes within seconds, and it’s more efficient, avoiding the cost of rebuilding and replacing the entire index. To implement this, we modified the PyTorch model to expose Upsert and Delete APIs, ensuring safe and efficient concurrent index updates during inference. Techniques used include pre-allocating larger tensors, using a high-water mark to track the working set, and making thread-safe in-memory tensor manipulations with minimal data access serialization. These methods ensure that modifications have minimal to no impact on the inference path, as detailed in the Model Inference Benchmarking section.

4.4 Inference on Native Stack

We built a native serving system as we made performance and efficiency our top priorities. To serve the PyTorch model in this system, we had to convert it to a compatible format. There are a few alternatives for this purpose such as TorchScript and torch.export.

PyTorch supports two execution modes: eager mode and graph mode. Optimal performance is achieved by executing everything in graph mode as the operators are first synthesized into a graph, which are compiled and executed as a whole. We picked TorchScript for our initial implementation to execute the model in graph mode. However, by doing so we traded off the performance with ease of development. TorchScript is a subset of Python and comes up with some constraints. It requires static typing and does not support things like exceptions and data-dependent control flows. We found executing this conversion quite challenging and concluded that it should be a part of the model development rather than an afterthought. We also decided to pursue other options which are deemed to be more recent technologies such as torch.export.

5 Experiments

In this section we provide results on modeling ablation studies, online A/B experiments with OON application and infrastructure model-based retrieval inference benchmarking for production indexes.

5.1 Model offline evaluation

We evaluated LiNR model on our internal dataset. The dataset consists of millions of examples where a member interacted with an item. The member and item are represented by embeddings learnt from a two-tower model. The two-tower model contains variety of features including member interaction history modeled by [19] and member profile features, which usually contains member job title, job location, company, skills and professional summary of the member. Posts usually contain text, image, video or external link information. Therefore, such content features from member and posts can help us to identify topics of interests of posts or professional topics of members.

We used Hit Rate @ 400 over evaluation dataset to report the metrics. We use cosine similarity with exhaustive search as baseline to evaluate against LiNR. We report results of Hadamard MLP for single embedding feature and extended Mixture-of-Logits with clustering for both single and multi-embedding features in Table 1.

Hadamard MLP is favored for production due to its simplicity for deployment and low latency. However, we found Hadamard MLP is very sensitive to weight initialization and general initialization methods such as GlorotNormal or HeNormal can’t stabilize the performance. Empirically we observed that the initial few steps determine the overall training trend, Thus we reinitialize the model if the first 100 steps go south.

In addition to the two-tower model and cluster ID features, we enhanced Mixture-of-Logits by introducing multiple embedding features developed at LinkedIn. We incorporated Graph Neural Network (GNN) embeddings for members and posts mapped to the same space using a heterogeneous GNN [1]. We found that adding more embeddings improved the Hit Rate @ 400 (see Table 1).

Our extended Mixture-of-Logits with clustering perform well for both single and multi-embedding features. One surprise finding is that fixed clusters (non-trainable) outperform trainable clusters in all cases we explored, one possible explanation is the convergence pace of the clustering and other trainable parameters are different, we’ll further investigate it in our future work. Another interesting

observation is that it was important to carefully tune the number of clusters: having either too high or too low a value can cause performance to degrade.

Method	Gain in Hit Rate @ 400
Cosine similarity	-
Single Embedding Feature	
Hadamard(Member & Item MLP [50]+[10, 1])	10.21%
MoL with 70 trained non-trained clusters	1.33% 10.11%
MoL with 100 trained non-trained clusters	11.97% 15.16%
MoL with 150 trained non-trained clusters	4.26% 11.17%
Multiple Embedding Features	
MoL without clustering	12.80%
MoL with 140 trained non-trained clusters	20.75% 22.61%
MoL with 200 trained non-trained clusters	16.49% 22.34%
MoL with 300 trained non-trained clusters	19.04% 23.67%

Table 1: Hit Rate @ 400 for single embedding of two-tower model alone in Hadamard or combined with cluster id in MoL, and multiple embeddings combining two-tower, GNN, and cluster ID in MoL.

5.2 A/B test of LiNR

For our baseline a dot-product EBR is done across all eligible items given member embedding query. We enabled cache on a cloud based storage for online lookup of computed results. This top K is retrieved by mid tier service when an online feed request is received. We leveraged this retrieval framework to test RAR based algorithms to understand relevance impact. The baseline for these experiments are full scan dot-product.

Metric Name	Metric Lift
Total professional interactions	+7%
Daily Unique Gold Professional Interactors	+3%
Feed Update Views With 30+ Secs Dwell	+2%
Feed Update Viewers With 30+ Secs Dwell	+5%
Skipped Update Rate	-20%

Table 2: LiNR A/B test relative metric improvements.

Table 2 shows the A/B test results from ramp of LiNR. *Total professional interactions* are the total amount of high quality interactions in the form of reshares, reposts, comments, message responses, reacts, votes, saves, and long dwells. *Daily Unique Gold Professional Interactors* is the daily moving average of the number of members or companies generating high quality interactions. *Feed Update Views With 30+ Secs Dwell* counts the total number of feed updates viewed with 30+ secs dwell time. *Feed Update Viewers With 30 Plus Secs Dwell* counts the total number of unique members that viewed a feed update for 30+ secs. *Skipped Update Rate* is the ratio of updates that are skipped (viewed for less than 2 seconds) compared to all viewed updates.

5.3 Model Inference Benchmarking

We conducted offline experiments to benchmark the effectiveness of different framework implementations of two KNN variants with ABM on datasets with different pass-rate scenarios. V1 represents KNN with similarity masking, and V2 represents KNN with explicit pre-filtering. Both are implemented in TF and PyTorch with CUDA kernel for attribute matching registered as a custom operator. We selected the Job recommendation index for benchmarking due to its

variety of filters, providing multi-dimensional performance insights for LiNR. The high-pass-rate and low-pass-rate datasets are derived from job search tasks, containing around 15.5 million jobs with 25 thousand queries. The high-pass-rate dataset includes two clauses: geo-location matching and company name reverse matching (mismatched items are returned), with an average of 1.7 million items passing the clauses. The low-pass-rate dataset includes an additional job title exact matching clause, with a maximum pass rate of 1.2 million for single title matching and most queries having only thousands of passed items. Each item and query has one attribute per clause, converted to 64-bit integers before GPU comparison. The embedding dimension is 128, stored as fp16 values. Performance is measured by average latency, p95 latency in milliseconds, and recall label@2000.

Method	Batch	Avg. Latency (ms/batch)	P95 Latency (ms/batch)	Recall@2k
TF-V1	1	6.3	6.9	0.688
TF-V2	1	6.9	14.4	0.688
PyTorch-V1	1	4.8	4.9	0.688
PyTorch-V2	1	14.6	47.8	0.688
TF-V1	16	34.8	36.6	0.688
PyTorch-V1	16	22.8	23.1	0.688

Table 3: Comparison of implementations on high-pass-rate dataset.

5.3.1 High-Pass-Rate ABM Dataset Benchmarking. From Table 3, we see that on the high-pass-rate dataset, both TF and PyTorch implementations of V1 (exhaustive search) are faster than V2 (explicit pre-filtering). This is likely due to the native slicing and copying operations in TF and PyTorch, which are especially slow for large matrices, as in V2 with high-pass-rate filters. Benchmarking individual operations revealed that the top-K selection in the latest TF version is slower than in PyTorch, while large-matrix slicing is slower in PyTorch than in TF, leading to performance differences between frameworks. V1’s implementation in the high-pass-rate dataset benefits more from the PyTorch implementation with increased batch sizes. Testing the V3 quantized KNN version showed further latency improvements with a trade-off in recall. In this experiment, we used 512-bit quantized embeddings and explored filtering different percentages of items based on quantized embedding similarity before full-precision similarity calculation. The trade-off between latency and recall, correlated with the filter size hyperparameter, is shown in Figure 7. By retaining 1% of items with an additional approximate ranking stage, we achieved around 10% further latency improvement with nearly parity performance.

Method	Batch	Avg. Latency (ms/batch)	P95 Latency (ms/batch)
TF-V2	1	3.4	4.5
PyTorch-V2	1	1.9	2.1
TF-V2	16	14.2	14.8
PyTorch-V2	16	21.4	21.9

Table 4: Comparison of implementations on low-pass-rate dataset.

5.3.2 Low-Pass-Rate ABM Dataset Benchmarking. As V2 version has the superiority on the low-pass-rate dataset compared to the other two versions (V3 may introduce redundant quantize matrix computation and filtering in the low-pass-rate case), we compare

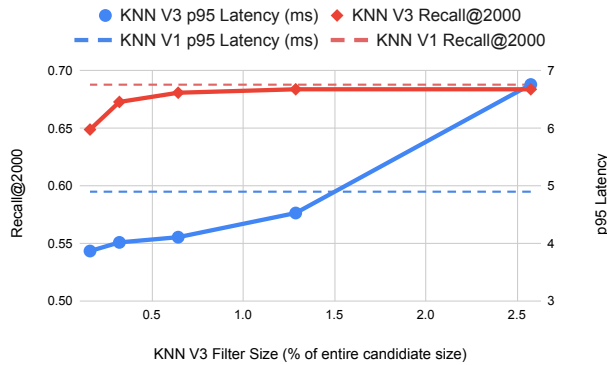


Figure 7: Trade-off of latency and recall correlated with the filter size of V3 quantized KNN with ABM on high-pass-rate dataset.

the its performance implemented with TF and PyTorch in Table 5. One single query, PyTorch still shows its advantage, but TF performs better on larger batch size. We attribute this to the fact the TF has better parallel schema for our case to conduct the retrieval in parallel. Though the queries are fetched in batch, since each query has different filters leading to different sets and number of retrieved items, we split the query batch and conduct the retrieval in parallel for each query independently. Considering that V2 is an exhaustive KNN search without liquidity issue, no recall drop and results are reported here.

Update Per Sec	Batch	QPS	Avg. Latency (ms/batch)	P95 Latency (ms/batch)
0	1	218	4.57	4.79
300	1	215	4.64	4.93
600	1	217	4.58	4.80
0	5	93	10.66	11.10
300	5	93	10.70	11.15
600	5	93	10.70	11.16

Table 5: Inference latency with concurrent model update.

We also conduct a stretch testing on single A100 GPU to measure the capacity of the exhaustive search method on handling large amount of items. For plain KNN with ABM (V1 & V2) on the high-pass-rate dataset, we are able to handle upto 240 million embeddings with 128 dim and fp16 precision for top-2k selection with single query. For the quantized KNN on an internal notification use case, which we select top-50million members from 1 billion members (64 dimensional embedding saved as fp16) to send relevant notifications, the 1-bit quantized KNN method with 64 bits quantized embedding size can reduce the original 120GB embedding memory to 7.5GB. When processing single query on an A100 GPU, it achieves maximum 21GB high-bandwidth memory with 97.6ms p95 latency.

5.3.3 Impact of Live Model Update on Inference. We run a benchmark to measure the impact of live model update on the inference latency on a single A100 GPU using our native serving system and bench marking tool. The bench marking tool uses a client for the native serving service and issues requests serially. We use plain KNN model with ABM (V1) and repeat the runs with various concurrent

update rates and request batch sizes. We observe no measurable impact on the latency with increased update rate.

6 Deployment lessons

Freshness: We initially deployed *LiNR* with offline inference and found it missed some fresh candidates. A/B tests revealed that live updates are crucial for serving newly created LinkedIn posts. Enabling live updates resulted in a +6% gain in our production systems, highlighting their importance for improved performance.

Pre-filtering: Our system employs EBR with pre-filtering, significantly enhancing retrieval quality. Many existing EBR infrastructures use KNN search with post-filtering, where results are first retrieved by KNN distance and then filtered. This post-filtering approach reduces system recall and quality by wasting candidate slots on items that don’t meet attribute constraints. By enabling pre-filtering on GPU retrieval, we greatly improved the quality of results compared to our production FAISS and lucene-based systems.

Custom filtering kernel: One lesson we learned early is that native TF or PyTorch do not effectively support filtering operations because deep learning frameworks weren’t initially designed for model-based retrieval indexes. Native boolean masking and indexing cause a 100X latency increase, making them impractical for production. Therefore, we implemented a custom CUDA solution for pre-filtering on the GPU, which scans items in memory to find those that meet constraints. One approach is to create a CUDA filtering kernel and fuse it with the matrix multiplication kernel to perform masked-matrix multiplication for KNN with pre-filtering. However, this solution is hard to generalize to other similarity measures or operations, as each new architecture would require re-implementation and fine-tuning, slowing down model development and deployment. In practice, we could make a trade-off of fully fused kernels and separately implemented kernels. For products needing regular KNN support, fusing the entire kernel with top-k selection and quantization improves serving speed. For general use cases, we create individual custom operations, like pre-filtering and quantization, to allow flexible development and deployment of advanced selection strategy with native neural network operations supported by TF and PyTorch. It is noteworthy that all the results reported in the paper are based on the second solution (even for the three plain KNN version) without extra custom kernel fusion for the purpose of self-consistency and generalizability.

7 Conclusion

In this paper, we introduced *LiNR*, a state-of-the-art model-based embedding retrieval solution for LinkedIn’s production system. Deploying *LiNR* to our online systems resulted in significant improvements in Out-Of-Network post recommendations on the LinkedIn Feed. We believe we are among the first in the industry to support live-updated, differentiable model-based indexing for recommendation and search applications. Looking forward, *LiNR* paves the way for unifying retrieval and ranking into a single GPU model, simplifying complex infrastructure and allowing end-to-end optimization of the entire differentiable system with gradient descent.

8 Acknowledgements

The authors would like to thank Jerry Shen, Yuchin Juan, Xiaobing Xue, Souvik Ghosh, Amol Ghoting, Vivek Hariharan, Ping Li, Luke Simon and others who collaborated with us.

References

- [1] Fedor Borisyyuk, Shihai He, Yunbo Ouyang, Morteza Ramezani, Peng Du, Xiaochen Hou, Chengming Jiang, Nitin Pasumarthy, Priya Bannur, Birjodh Tiwana, Ping Liu, Siddharth Dangi, Daqi Sun, Zhoutao Pei, Xiao Shi, Sirou Zhu, Qianqi Shen, Kuang-Hsuan Lee, David Stein, Baolei Li, Haichao Wei, Amol Ghoting, and Souvik Ghosh. 2024. LiGNN: Graph Neural Networks at LinkedIn. In *KDD*.
- [2] Haonan Chen, Carlos Lassance, and Jimmy Lin. 2023. End-to-End Retrieval with Learned Dense and Sparse Representations Using Lucene. arXiv:2311.18503 [cs.IR]
- [3] Michael Curtiss, Iain Becker, Tudor Bosman, Sergey Doroshenko, Lucian Grijincu, Tom Jackson, Sandhya Kunnatur, Soren Lassen, Philip Pronin, Sriram Sankar, Guanghao Shen, Gintaras Woss, Chao Yang, and Ning Zhang. 2013. Unicorn: A System for Searching the Social Graph. In *VLDB*.
- [4] Ruiqi Guo, Philip Sun, Erik Lindgren, Quan Geng, David Simcha, Felix Chern, and Sanjiv Kumar. 2020. Accelerating Large-Scale Inference with Anisotropic Vector Quantization. *ICML*.
- [5] Félix GV. 2008. Open Sourcing Venice – LinkedIn’s Derived Data Platform. <https://www.linkedin.com/blog/engineering/open-source/open-sourcing-venice-linkedin-s-derived-data-platform>.
- [6] Jui-Ting Huang, Ashish Sharma, Shuying Sun, Li Xia, David Zhang, Philip Pronin, Janani Padmanabhan, Giuseppe Ottaviano, and Linjun Yang. 2020. Embedding-Based Retrieval in Facebook Search. In *KDD*.
- [7] Biye Jiang, Chao Deng, Huimin Yi, Zelin Hu, Guorui Zhou, Yang Zheng, Sui Huang, Xinyang Guo, Dongyue Wang, Yue Song, Liqin Zhao, Zhi Wang, Peng Sun, Yu Zhang, Di Zhang, Jinhui Li, Jian Xu, Xiaoqiang Zhu, and Kun Gai. 2019. XDL: An Industrial Deep Learning Framework for High-Dimensional Sparse Data. *KDD*.
- [8] Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2021. Billion-Scale Similarity Search with GPUs. *IEEE Transactions on Big Data* (2021).
- [9] Herve Jégou, Matthijs Douze, and Cordelia Schmid. 2011. Product Quantization for Nearest Neighbor Search. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2011).
- [10] Ping Li and Xiaoyun Li. 2023. OPORP: One permutation+ one random projection. *arXiv preprint arXiv:2302.03505* (2023).
- [11] Xiangru Lian, Binhang Yuan, Xuefeng Zhu, Yulong Wang, Yongjun He, Honghuan Wu, Lei Sun, Haodong Lyu, Chengjun Liu, Xing Dong, Yiqiao Liao, Mingnan Luo, Congfei Zhang, Jingru Xie, Haonan Li, Lei Chen, Renjie Huang, Jianying Lin, Chengchun Shu, Xuezhong Qiu, Zhishan Liu, Dongying Kong, Lei Yuan, Hai Yu, Sen Yang, Ce Zhang, and Ji Liu. 2022. Persia: An Open, Hybrid System Scaling Deep Learning-Based Recommenders up to 100 Trillion Parameters. In *KDD*.
- [12] Yiqun Liu, Kaushik Rangadurai, Yunzhong He, Siddarth Malreddy, Xunlong Gui, Xiaoyi Liu, and Fedor Borisyyuk. 2021. Que2Search: Fast and Accurate Query and Document Understanding for Search at Facebook. *KDD*.
- [13] Zhuoran Liu, Leqi Zou, Xuan Zou, Caihua Wang, Biao Zhang, Da Tang, Bolin Zhu, Yijie Zhu, Peng Wu, Ke Wang, and Youlong Cheng. 2022. Monolith: Real Time Recommendation System With Collisionless Embedding Table. arXiv:2209.07663 [cs.IR]
- [14] Yu A. Malkov and D. A. Yashunin. 2020. Efficient and Robust Approximate Nearest Neighbor Search Using Hierarchical Navigable Small World Graphs. *IEEE Trans. Pattern Anal. Mach. Intell.* (2020).
- [15] Corey Nolet. 2023. *Reusable Computational Patterns for Machine Learning and Information Retrieval with RAPIDS RAFT*. <https://developer.nvidia.com/blog/reusable-computational-patterns-for-machine-learning-and-data-analytics-with-rapids-raft/>
- [16] Hiroyuki Ootomo, Akira Naruse, Corey Nolet, Ray Wang, Tamas Feher, and Yong Wang. 2023. CAGRA: Highly Parallel Graph Construction and Approximate Nearest Neighbor Search for GPUs. arXiv:2308.15136 [cs.DS]
- [17] Aditya Pal, Chantat Eksombatchai, Yitong Zhou, Bo Zhao, Charles Rosenberg, and Jure Leskovec. 2020. PinnerSage: Multi-Modal User Embedding Framework for Recommendations at Pinterest. *KDD*.
- [18] Shashank Rajput, Nikhil Mehta, Anima Singh, Raghunandan H. Keshavan, Trung Vu, Lukasz Heldt, Lichan Hong, Yi Tay, Vinh Q. Tran, Jonah Samost, Maciej Kula, Ed H. Chi, and Maheswaran Sathiamoorthy. 2023. Recommender Systems with Generative Retrieval. arXiv:2305.05065 [cs.IR]
- [19] Kaushik Rangadurai, Yiqun Liu, Siddarth Malreddy, Xiaoyi Liu, Piyush Maheshwari, Vishwanath Sangale, and Fedor Borisyyuk. 2022. NxtPost: User To Post Recommendations In Facebook Groups. In *KDD*.
- [20] Jianqiang Shen, Yuchin Juan, Shaobo Zhang, Ping Liu, Wen Pu, Sriram Vasudevan, Qingquan Song, Fedor Borisyyuk, Kay Qianqi Shen, Haichao Wei, Yunxiang Ren, Yeou S. Chiou, Sicong Kuang, Yuan Yin, Ben Zheng, Muchen Wu, Shaghayegh Gharghabi, Xiaoqing Wang, Huichao Xue, Qi Guo, Daniel Hewlett, Luke Simon, Liangjie Hong, and Wenjing Zhang. 2024. Learning to Retrieve for Job Matching. arXiv:2402.13435 [cs.IR] <https://arxiv.org/abs/2402.13435>
- [21] Yi Tay, Vinh Q. Tran, Mostafa Dehghani, Jianmo Ni, Dara Bahri, Harsh Mehta, Zhen Qin, Kai Hui, Zhe Zhao, Jai Gupta, Tal Schuster, William W. Cohen, and Donald Metzler. 2022. Transformer Memory as a Differentiable Search Index. *NEURIPS*.
- [22] Jiaqi Zhai, Zhaojie Gong, Yueming Wang, Xiao Sun, Zheng Yan, Fu Li, and Xing Liu. 2023. Revisiting Neural Retrieval on Accelerators. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining (KDD '23)*. Association for Computing Machinery, New York, NY, USA, 5520–5531. <https://doi.org/10.1145/3580305.3599897>
- [23] Hailin Zhang, Yujing Wang, Qi Chen, Ruiheng Chang, Ting Zhang, Ziming Miao, Yingyan Hou, Yang Ding, Xupeng Miao, Haonan Wang, Bochen Pang, Yuefeng Zhan, Hao Sun, Weiwei Deng, Qi Zhang, Fan Yang, Xing Xie, Mao Yang, and Bin Cui. 2023. Model-enhanced Vector Index.
- [24] Weijie Zhao, Shulong Tan, and Ping Li. 2020. SONG: Approximate Nearest Neighbor Search on GPU. In *ICDE*.
- [25] Weijie Zhao, Shulong Tan, and Ping Li. 2022. Constrained Approximate Similarity Search on Proximity Graph. *arXiv preprint arXiv:2210.14958* (2022).