
End-to-end Learnable Clustering for Intent Learning in Group Recommendation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Intent learning, which aims to learn users' intents for user understanding and item
2 recommendation, has become a hot research spot in recent years. However, the
3 existing methods suffer from complex and cumbersome alternating optimization,
4 limiting the performance and scalability. To this end, we propose a novel intent
5 learning method termed ELCRec, by unifying behavior representation learning
6 into an End-to-end Learnable Clustering framework, for effective and efficient
7 Recommendation. Concretely, we encode users' behavior sequences and initialize
8 the cluster centers (latent intents) as learnable neurons. Then, we design a novel
9 learnable clustering module to separate different cluster centers, thus decoupling
10 users' complex intents. Meanwhile, it guides the network to learn intents from
11 behaviors by forcing behavior embeddings close to cluster centers. This allows
12 simultaneous optimization of recommendation and clustering via mini-batch data.
13 Moreover, we propose intent-assisted contrastive learning by using cluster centers
14 as self-supervision signals, further enhancing mutual promotion. Both experimental
15 results and theoretical analyses demonstrate the superiority of ELCRec from six
16 perspectives. Compared to the runner-up, ELCRec improves NDCG@5 by 8.9%
17 and reduces computational costs by 22.5% on Beauty dataset. Furthermore, due to
18 the scalability and universal applicability, we deploy this method on the industrial
19 recommendation system with 130 million page views and achieve promising results.
20 The codes are available at [Anonymous GitHub](#).

21 1 Introduction

22 Sequential Recommendation (SR), which aims to recommend relevant items to users by learning
23 patterns from users' historical behavior sequences, is a vital and challenging task in machine learning
24 domain. In recent years, benefiting the strong representation learning ability of deep neural networks
25 (DNNs), DNN-based sequential recommendation methods[95, 32, 85, 111, 43, 98, 45, 58] have
26 achieved promising recommendation performance and attracted researchers' high level of attention.

27 More recently, intent learning has become a hot topic in both research and industrial field of recom-
28 mendation. It aims to model users' intents by learning from users' historical behaviors. For example,
29 a user interacted the shoes, bag, and racket in history. Thus, the user's potential intent can be inferred
30 as playing badminton. Then, the system may recommend the intent-relevant items to the user. Follow
31 this principle, various intent learning methods [37, 11, 38, 15, 42, 46, 5] have been proposed to
32 achieve better user understanding and item recommendation.

33 The optimization paradigm of the recent representative intent learning methods can be summarized
34 as a generalized Expectation Maximization (EM) framework. To be specific, at the E-step, clustering
35 algorithms are adopted to learn the latent intents from users' behavior embeddings. And, at the

36 M-step, the self-supervised learning methods are utilized to embed behaviors. The optimizations of
 37 these two steps are conducted alternately, achieving promising performance.

38 However, we highlight two issues in this complex and tedious alternating optimization. (1) At
 39 the E-step, we need to apply the clustering algorithm on the whole data, limiting the model’s
 40 scalability, especially in large-scale industrial scenarios, e.g., apps with billion users. (2) In the EM
 41 framework, the optimization of behavior learning and the clustering algorithm are separated, leading
 42 to sub-optimal performance and increasing the implementation difficulty.

43 To this end, we propose a novel intent learning model named ELCRec via integrating represen-
 44 tation learning into an End-to-end Learnable Clustering framework, for effective and efficient
 45 Recommendation. Specifically, the user’s behavioral process is first embedded into the latent space.
 46 Cluster centers, recognized as the users’ latent intents, are initialized as learnable neural network
 47 parameters. Then, a simple yet effective learnable clustering module is proposed to decouple users’
 48 complex intents into different simple intent units by separating the cluster centers. Meanwhile, it
 49 makes the behavior embeddings close to cluster centers to guide the models to learn more accurate
 50 intents from users’ behaviors. This improves the model’s scalability and alleviates the issue (1) by
 51 optimizing the cluster distribution on mini-batch data. Furthermore, to further enhance the mutual
 52 promotion of representation learning and clustering, we present intent-assisted contrastive learning to
 53 integrate the cluster centers as self-supervision signals for representation learning. These settings
 54 unify behavior learning and clustering optimization in an end-to-end optimizing framework, improv-
 55 ing recommendation performance and simplifying deployment. Therefore, the issue (2) has been also
 56 solved. The contributions of this paper are summarized as follows.

- 57 • We innovatively promote the existing optimization framework of intent learning by unifying
 58 behavior representation learning and clustering optimization.
- 59 • A new intent learning model termed ELCRec is proposed with a simple yet effective learnable
 60 cluster module and intent-assisted contrastive learning.
- 61 • Comprehensive experiments and theoretical analyses show advantages of ELCRec from six
 62 aspects, including superiority, effectiveness, efficiency, sensitivity, convergence, and visualization.
- 63 • We successfully deployed it on industrial recommendation system with 130 million page views
 64 and achieve promising results, providing various practical insights.

65 2 Related Work

66 We provide a brief overview of the related work for this paper. It can be divided into three parts,
 67 including sequential recommendation, intent learning, and clustering algorithms. At first, Sequential
 68 Recommendation (SR) focuses on recommending relevant items to users based on their historical
 69 behavior sequences. In addition, intent learning has emerged as a promising and practical technique in
 70 recommendation systems. It aims to capture users’ latent intents to achieve better user understanding
 71 and item recommendation. Lastly, clustering algorithms play a crucial role in recommendation
 72 systems since they can identify patterns and similarities in the users or items. Due to the limitation of
 73 the pages, we introduce the detailed related methods in the Appendix 7.9.

74 3 Methodology

75 We present our proposed framework, ELCRec, in this section. Firstly, we provide the necessary
 76 notations and task definition. Secondly, we analyze and identify the limitations of existing intent
 77 learning. Finally, we propose our solutions to address these challenges.

78 3.1 Basic Notation

79 In a recommendation system, \mathcal{U} denotes the user set, and \mathcal{V} denotes the item set. For each user $u \in \mathcal{U}$,
 80 the historical behaviors are described by a sequence of interacted items $S^u = [s_1^u, s_2^u, \dots, s_t^u, \dots, s_{|S^u|}^u]$.
 81 S^u is sorted by time. $|S^u|$ denotes the interacted items number of user u . s_t^u denotes the item which
 82 is interacted with user u at t step. In practice, during sequence encoding, the historical behavior
 83 sequences are limited with a maximum length T [29, 32, 15]. The sequences truncated and remain
 84 the most recent T interacted items if the length is greater than T . Besides, the shorter sequences are

filled with “padding” items on the left until the length is T . Due to the limitation of the pages, we list the basic notations in Table 5 of the Appendix 7.1.

3.2 Task Definition

Given the user set \mathcal{U} and the item set \mathcal{V} , the recommendation system aims to precisely model the user interactions and recommend items to users. Take user u for an example, the sequence encoder firstly encodes the user’s historical behaviors S^u to the latent embedding \mathbf{E}^u . Then, based on the historical behavior embedding, the target of the recommendation task is to predict the next item that is most likely interacted with by user u at $|S^u| + 1$ step.

3.3 Problem Analyses

Among the techniques in recommendation, intent learning has become an effective technique to understand users. We summarize the optimization procedure of the intent learning as the Expectation Maximization (EM) framework. It contains two steps including E-step and M-step. These two steps are conducted alternately, mutually promoting each other. However, we find two issues of the existing optimization framework as follows.

- (1) In the process of E-step, it needs to perform a clustering algorithm on the full data, easily leading to out-of-memory or long-running time problems. It restricts the scalability of the model on large-scale industrial data.
- (2) The alternative optimization approach within the EM framework separates the learning process for behaviors and intents, leading to sub-optimal performance and increased implementation complexity. Also, it limits the training and inference on the real-time data. That is, when users’ behaviors and intents change over time, there is a long lag in the training and inference process

Therefore, we aim to develop a new optimization framework for intent learning to solve issue (1) and issue (2). For the issue (1), a new learnable online clustering method is the key solution. For the issue (2), we aim to break the alternative optimization in the EM framework.

3.4 Proposed Method

To this end, we present a new intent learning method termed **ELCRec** by unifying sequence representation learning into an **End-to-end Learnable Clustering** framework, for **Recommendation**. It contains three parts, including behavior encoding, end-to-end learnable cluster module (ELCM), and intent-assisted contrastive learning (ICL).

3.4.1 Behavior Encoding

In this process, we aim to encoder the users’ behavior sequences. Concretely, given the user set \mathcal{U} , the item set \mathcal{V} , and the users’ historical behavior sequence set $\{S^u\}_{u=1}^{|\mathcal{U}|}$, the behavior encoder \mathcal{F} embeds the behavior sequences of each user u into the latent space as follows.

$$\mathbf{E}^u = \mathcal{F}(S^u), \quad (1)$$

where $\mathbf{E}^u \in \mathbb{R}^{|S^u| \times d'}$ denotes the behavior sequence embedding of user u , d' is the dimension number of latent features, and $|S^u|$ denotes the length of behavior sequence of user u . Note that the behavior sequence lengths of different users are different. Therefore, all user behavior sequences are pre-processed to the sequences with the same length T by padding or truncating. The encoder \mathcal{F} is designed as a Transformer-based [91] architecture. Subsequently, to summarize the behaviors over different time of each user, the behavior sequence embedding is aggregated by the concatenate pooling function \mathcal{P} as follows.

$$\mathbf{h}_u = \mathcal{P}(\mathbf{E}^u) = \text{concat}(\mathbf{e}_1^u || \dots || \mathbf{e}_i^u || \dots || \mathbf{e}_T^u), \quad (2)$$

where $\mathbf{e}_i^u \in \mathbb{R}^{1 \times d'}$ denotes the embedding of user behavior at i -th step and $\mathbf{h}_u \in \mathbb{R}^{1 \times Td'}$ denotes the aggregated behavior embedding of user u . We re-denote Td' as d for convenience. By encoding and aggregation, we obtain the behavior embeddings of all users $\mathbf{H} \in \mathbb{R}^{|\mathcal{U}| \times d}$.

3.4.2 End-to-end Learnable Cluster Module

After behavior encoding, we guide the model to learn the users' latent intents from the behavior embeddings. To this end, an end-to-end learnable cluster module (ELCM) is proposed to break the alternative optimization in the previous mentioned EM framework. This module can group the users' behaviors embeddings into various clusters, which represent the users' latent intents or interests. Concretely, at first, the cluster centers $\mathbf{C} \in \mathbb{R}^{k \times d}$ are initialized as the learnable neural parameters, i.e., the tensors with gradients. Then, we design a simple yet effective clustering loss to train the networks and cluster centers as formulated as follows.

$$\mathcal{L}_{\text{cluster}} = \underbrace{\frac{-1}{(k-1)k} \sum_{i=1}^k \sum_{j=1, j \neq i}^k \|\hat{\mathbf{c}}_i - \hat{\mathbf{c}}_j\|_2^2}_{\text{Intent Decoupling}} + \underbrace{\frac{1}{bk} \sum_{i=1}^b \sum_{j=1}^k \|\hat{\mathbf{h}}_i - \hat{\mathbf{c}}_j\|_2^2}_{\text{Intent-behavior Alignment}}, \quad (3)$$

where $\hat{\mathbf{h}}_i = \mathbf{h}_i / \|\mathbf{h}_i\|_2$, $\hat{\mathbf{c}}_i = \mathbf{c}_i / \|\mathbf{c}_i\|_2$. In Eq. (3), k denotes the number of clusters (intents), and b denotes the batch size. $\mathbf{h}_i \in \mathbb{R}^{1 \times d}$ denotes the i -th user's behavior embedding and $\mathbf{c}_j \in \mathbb{R}^{1 \times d}$ denotes the j -th cluster center. For better network convergence, we constrain the behavior embeddings and cluster center embeddings to distribute on a unit sphere. Concretely, we apply the l_2 normalization to both the user behavior embeddings \mathbf{H} and the cluster centers \mathbf{C} during calculating $\mathcal{L}_{\text{cluster}}$.

In the proposed clustering loss, the first term is designed to disentangle the complex users' intents into simple intent units. Technically, it pushes away different cluster centers, therefore reducing the overlap between different clusters (intents). The time complexity and space complexity of this term are $\mathcal{O}(k^2d)$ and $\mathcal{O}(kd)$, respectively. The number of users' intents is vastly less than the number of users, i.e., $k \ll |\mathcal{U}|$. Therefore, the first term will not bring significant time or space costs.

In addition, the second term of the proposed clustering loss aims to align the users' latent intents with the behaviors by pulling the behavior embeddings to the cluster centers. This design makes the in-class cluster distribution more compact and guides the network to condense similar behaviors into one intention. Also, on another aspect, it forces the model to learn users' intents from behavior embeddings. Note that the behavior embedding \mathbf{h}_i is pulled to all center centers $\mathbf{c}_j, j = 1, \dots, k$ rather than the nearest cluster center. The main reason is that the practical clustering algorithm is imperfect, and pulling to the nearest center easily leads to the confirmation bias problem [67]. To this end, the proposed clustering loss $\mathcal{L}_{\text{cluster}}$ aims to optimize the clustering distribution in an adversarial manner by pulling embeddings together to cluster centers while pushing different cluster centers away. Besides, it enables the optimization of this term via mini-batch samples, avoiding performance clustering algorithms on the whole data. Time complexity and space complexity of the second term are $\mathcal{O}(bkd)$ and $\mathcal{O}(bk + bd + kd)$, respectively. Since the batch size is essentially less than the number of users, namely, $b \ll |\mathcal{U}|$, the second term of clustering loss $\mathcal{L}_{\text{cluster}}$ alleviates the considerable time or space costs. Besides, theoretically, based on the Rademacher complexity, we investigate the generalization bounds of $\mathcal{L}_{\text{cluster}}$ in the Appendix 7.3.

In the existing EM optimization framework, the clustering algorithm needs to be applied on the entire users' behavior embeddings $\mathbf{H} \in \mathbb{R}^{|\mathcal{U}| \times d}$. Take the classical k -Means clustering as an example, at each E-step, it leads to $\mathcal{O}(t|\mathcal{U}|kd)$ time complexity and $\mathcal{O}(|\mathcal{U}|k + |\mathcal{U}|d + kd)$ space complexity, where t denote the iteration steps of k -Means clustering algorithm. We find that, at each step, the time and space complexity is linear to the number of users, thus leading to out-of-memory or running time problems (issue (1)), especially on large-scale industrial data with millions or billions of users.

Fortunately, our proposed end-to-end learnable cluster module can solve this issue (1). By summarising previous analyses, we draw that the overall time and space complexity of calculating the clustering loss $\mathcal{L}_{\text{cluster}}$ are $\mathcal{O}(bkd + k^2d + bd)$ and $\mathcal{O}(bk + bd + kd)$, respectively. They are both linear to the batch size b at each step, enabling the model's scalability. Besides, the proposed module is plug-and-play and easily deployed in real-time large-scale industrial systems. We provide detailed evidence and practical insights in Section 5. The proposed ELCM can not only improve the recommendation performance (See Section 4.2 & 4.3) but also promote efficiency (See Section 4.4).

3.4.3 Intent-assisted Contrastive Learning

Next, we aim to enhance further the mutual promotion of behavior learning and clustering. To this end, Intent-assisted contrastive learning (ICL) is proposed by adopting cluster centers as self-supervision

signals for behavior learning. Firstly, we conduct contrastive learning among the behavior sequences. The new views of the behavior sequences are constructed via sequential augmentations, including mask, crop, and reorder. The two views of behavior sequence of user u are denoted as $(S^u)^{v1}$ and $(S^u)^{v2}$. According to Section 3.4.1, the behaviors are encoded to the behavior embeddings $\mathbf{h}_u^{v1}, \mathbf{h}_u^{v2} \in \mathbb{R}^{1 \times d}$. Then, the sequence contrastive loss of user u is formulated as follows.

$$\mathcal{L}_{\text{seq_cl}}^u = - \left(\log \frac{e^{\text{sim}(\mathbf{h}_u^{v1}, \mathbf{h}_u^{v2})}}{\sum_{\text{neg}} e^{\text{sim}(\mathbf{h}_u^{v1}, \mathbf{h}_{\text{neg}})}} + \log \frac{e^{\text{sim}(\mathbf{h}_u^{v1}, \mathbf{h}_u^{v2})}}{\sum_{\text{neg}} e^{\text{sim}(\mathbf{h}_u^{v2}, \mathbf{h}_{\text{neg}})}} \right), \quad (4)$$

where “sim” denotes the dot-product similarity, “neg” denotes the negative samples. Here, the same sequence with different augmentations is recognized as the positive sample pairs, and the other sample pairs are recognized as the negative sample pairs. By minimizing $\mathcal{L}_{\text{seq_cl}} = \sum_u \mathcal{L}_{\text{seq_cl}}^u$, the similar behaviors are pulled together, and the others are pushed away from each other, therefore enhancing the representation capability of users’ behaviors. The learned cluster centers $\mathbf{C} \in \mathbb{R}^{k \times d}$ are adopted as the self-supervision signals. Index of the assigned cluster of \mathbf{h}_u^{v1} is queried as follows.

$$idx = \arg \min_i (\|\mathbf{c}_i - \mathbf{h}_u^{v1}\|_2^2), \quad (5)$$

where $\mathbf{c}_i \in \mathbb{R}^{1 \times d}$ denotes the i -th cluster (intent) center embedding. Then, the intent information is fused to the user behavior during the sequence contrastive learning. Here, we consider two optional fusion strategies, including the concatenate fusion $\mathbf{h}_u^{v1} = \text{concat}(\mathbf{h}_u^{v1} \parallel \mathbf{c}_{idx})$ and the shift fusion $\mathbf{h}_u^{v1} = \mathbf{h}_u^{v1} + \mathbf{c}_{idx}$. A similar operation is applied to the second view of the behavior embedding \mathbf{h}_u^{v2} . After fusing the intent information to user behaviors, the networks are trained by minimizing $\mathcal{L}_{\text{seq_cl}}$.

In addition, to further collaborate intent learning and sequential representation learning, we conduct contrastive learning between the user’s behaviors and the learnable intent centers. The intent contrastive loss is formulated as follows.

$$\mathcal{L}_{\text{intent_cl}}^u = - \left(\log \frac{\min_i e^{\text{sim}(\mathbf{h}_u^{v1}, \mathbf{c}_i)}}{\sum_{\text{neg}} e^{\text{sim}(\mathbf{h}_u^{v1}, \mathbf{c}_{\text{neg}})}} + \log \frac{\min_i e^{\text{sim}(\mathbf{h}_u^{v2}, \mathbf{c}_i)}}{\sum_{\text{neg}} e^{\text{sim}(\mathbf{h}_u^{v2}, \mathbf{c}_{\text{neg}})}} \right), \quad (6)$$

where $\mathbf{h}_u^{v1}, \mathbf{h}_u^{v2}$ are two-view behavior embedding of the user u . Besides, “neg” denotes the negative behavior-intent pairs among all pairs. Here, we regard the behavior embedding and the corresponding nearest intent center as the positive pair and others as negative pairs. By minimizing the intent contrastive loss $\mathcal{L}_{\text{intent_cl}} = \sum_u \mathcal{L}_{\text{intent_cl}}^u$, behaviors with the same intents are pulled together, but behaviors with different intents are pushed away. The objective of ICL is formulated as follows.

$$\mathcal{L}_{\text{icl}} = \mathcal{L}_{\text{seq_cl}} + \mathcal{L}_{\text{intent_cl}}. \quad (7)$$

The effectiveness of ICL is verified in Section 4.3. With the proposed ELCM and ICL, we develop a new end-to-end optimization framework for intent learning, improving performance and convenience. By these designs, the issue (2) is also solved.

3.4.4 Overall Objective

The neural networks and learnable clusters are trained with multiple tasks, including intent learning, intent-assisted contrastive learning, and next-item prediction. The intent learning task aims to capture the users’ underlying intents. Besides, intent-assisted contrastive learning aims to collaborate with intent learning and behavior learning. In addition, the next-item prediction task is a widely used task for recommendation systems. The overall objective of ELCRec is formulated as follows.

$$\mathcal{L}_{\text{overall}} = \mathcal{L}_{\text{next_item}} + 0.1 \times \mathcal{L}_{\text{icl}} + \alpha \times \mathcal{L}_{\text{cluster}}, \quad (8)$$

where $\mathcal{L}_{\text{next_item}}$, \mathcal{L}_{icl} , and $\mathcal{L}_{\text{cluster}}$ denotes the next item prediction loss, intent-assisted contrastive learning loss, and clustering loss, respectively. α is a trade-off hyper-parameter. We present the overall algorithm process of the proposed ELCRec method in Algorithm 1 in Appendix.

4 Experiment

This section aims to comprehensively evaluate ELCRec by answering research questions (RQs).

- 215 (i) Superiority: does it outperform the state-of-the-art sequential recommendation methods?
 - 216 (ii) Effectiveness: are the ELCM and ICL modules effective?
 - 217 (iii) Efficiency: how about the time and memory efficiency of the proposed ELCRec?
 - 218 (iv) Sensitivity: what is the performance of the proposed method with different hyper-parameters?
 - 219 (v) Convergence: have the loss function and recommendation performance converged?
 - 220 (vi) Visualization: Can the visualized learned embeddings reflect the promising results?
- 221 We answer RQ(i), (ii), (iii) in Section 4.2, 4.3, 4.4, respectively. Due to the limited pages, RQ(iv), (v),
 222 (vi) are answered in the Appendix 7.5, 7.6, and 7.7 respectively.

223 4.1 Experimental Setup

224 4.1.1 Experimental Environment

225 Experimental results on the public benchmarks are obtained from the desktop computer with one
 226 NVIDIA GeForce RTX 4090 GPU, six 13th Gen Intel(R) Core(TM) i9-13900F CPUs, and the
 227 PyTorch platform. During training, we monitored the training process via the Weights & Biases.

228 4.1.2 Public Benchmark

229 We performed our experiments on four public benchmarks: Sports, Beauty, Toys, and Yelp¹. The
 230 Sports, Beauty, and Toys datasets are subcategories of the Amazon Review Dataset [62]. The Sports
 231 dataset contains reviews for sporting goods, the Beauty dataset contains reviews for beauty products,
 232 and the Toys dataset contains toy reviews. On the other hand, the Yelp dataset focuses on business
 233 recommendations and is provided by Yelp company. Table 6 summarizes the datasets' details. We
 234 only kept datasets where all users and items have at least five interactions. Besides, we adopted the
 235 dataset split settings used in the previous method [15].

236 4.1.3 Evaluation Metric

237 To evaluate ELCRec, we adopt two groups of metrics, including Hit Ratio@ k (HR@ k) and Normal-
 238 ized Discounted Cumulative Gain@ k (NDCG@ k), where $k \in \{5, 20\}$.

239 4.1.4 Compared Baseline

240 We compare our method with nine baselines including BPR-MF [79], GRU4Rec [29], Caser [87],
 241 SASRec [32], DSSRec [60], BERT4Rec [85], S3-Rec [111], CL4SRec [98], and ICLRec [15].
 242 Detailed introductions to these methods are in the Appendix 7.9.2.

243 4.1.5 Implementation Detail

244 For the baselines, we adopt their original code with the original settings to reproduce the results on
 245 four benchmarks. Due to page limitation, the detailed implementation of the baselines are listed in
 246 Appendix 7.10. The proposed method, ELCRec, was implemented using the PyTorch deep learning
 247 platform. In the Transformer encoder, we employed self-attention blocks with two attention heads.
 248 The latent dimension, denoted as d , was set to 64, and the maximum sequence length, denoted as T ,
 249 was set to 50. We utilized the Adam optimizer with a learning rate of 1e-3. The decay rate for the
 250 first moment estimate was set to 0.9, and the decay rate for the second moment estimate was set to
 251 0.999. The cluster number, denoted as k , was set to 256 for the Yelp and Beauty datasets and 512
 252 for the Sports and Toys datasets. The trade-off hyper-parameter, denoted as α , was set to 1 for the
 253 Sports and Toys datasets, 0.1 for the Yelp dataset, and 10 for the Beauty dataset. During training, we
 254 monitored the training process via the Weights & Biases.

255 4.2 Superiority

256 In this section, we aim to answer the research question (i) and demonstrate the superiority of
 257 ELCRec. To be specific, we compare ELCRec with nine state-of-the-art recommendation baselines

¹<https://www.yelp.com/dataset>

Table 1: Recommendation performance on benchmarks. **Bold values** and underlined values denote the best and runner-up results. * indicates that, in the t -test, the best method significantly outperforms the runner-up with $p < 0.05$. "-" indicates models do not converge.

Dataset	Metric	BPR-MF [79]	GRU4Rec [29]	Caser [87]	SASRec [32]	BERT4Rec [85]	DSSRec [60]	S3-Rec [111]	CL4SRec [98]	DCRec [100]	MAERec [102]	IOCRec [42]	ICLRec [15]	ELCRec Ours	Impr.	p-value
Sports	HR@5	0.0141	0.0162	0.0154	0.0206	0.0217	0.0214	0.0121	0.0217	0.0172	0.0225	0.0246	<u>0.0263</u>	0.0286	8.75% \uparrow	2.34e-6*
	HR@20	0.0323	0.0421	0.0399	0.0497	0.0604	0.0495	0.0344	0.0540	0.0357	0.0488	<u>0.0641</u>	0.0630	0.0648	1.09% \uparrow	2.29e-4*
	NDCG@5	0.0091	0.0103	0.0114	0.0135	0.0143	0.0142	0.0084	0.0137	0.0118	0.0152	0.0162	<u>0.0173</u>	0.0185	6.94% \uparrow	3.54e-5*
	NDCG@20	0.0142	0.0186	0.178	0.0216	0.0251	0.0220	0.0146	0.0227	0.0170	0.0225	<u>0.0280</u>	0.0276	0.0286	2.14% \uparrow	7.87e-3*
Beauty	HR@5	0.0212	0.0111	0.0251	0.0374	0.0360	0.0410	0.0189	0.0423	0.0368	0.0414	0.0408	<u>0.0495</u>	0.0529	6.87% \uparrow	3.18e-6*
	HR@20	0.0589	0.0478	0.0643	0.0901	0.0984	0.0914	0.0487	0.0994	0.0674	0.0854	0.0916	<u>0.1072</u>	0.1079	0.65% \uparrow	3.30e-3*
	NDCG@5	0.0130	0.0058	0.0145	0.0241	0.0216	0.0261	0.0115	0.0281	0.0269	0.0283	0.0245	<u>0.0326</u>	0.0355	8.90% \uparrow	4.48e-6*
	NDCG@20	0.0236	0.0104	0.0298	0.0387	0.0391	0.0403	0.0198	0.0441	0.0357	0.0407	0.0444	<u>0.0491</u>	0.0509	3.67% \uparrow	9.08e-6*
Toys	HR@5	0.0120	0.0097	0.0166	0.0463	0.0274	0.0502	0.0143	0.0526	0.0399	0.0477	0.0311	0.0586	<u>0.0585</u>	0.17% \downarrow	1.22e-1
	HR@20	0.0312	0.0301	0.0420	0.0941	0.0688	0.0975	0.0235	0.1038	0.0679	0.0904	0.0781	<u>0.1130</u>	0.1138	0.71% \uparrow	4.20e-3*
	NDCG@5	0.0082	0.0059	0.0107	0.0306	0.0174	0.0337	0.0123	0.0362	0.0296	0.0336	0.0197	<u>0.0397</u>	0.0403	1.51% \uparrow	2.87e-4*
	NDCG@20	0.0136	0.0116	0.0179	0.0441	0.0291	0.0471	0.0162	0.0506	0.0374	0.0458	0.0330	<u>0.0550</u>	0.0560	1.82% \uparrow	3.72e-5*
Yelp	HR@5	0.0127	0.0152	0.0142	0.0160	0.0196	0.0171	0.0101	0.0229		0.0166	0.0222	<u>0.0233</u>	0.0236	1.29% \uparrow	7.81e-3*
	HR@20	0.0346	0.0371	0.0406	0.0443	0.0564	0.0464	0.0314	0.0630		0.0460	0.0640	<u>0.0645</u>	0.0653	1.24% \uparrow	3.73e-4*
	NDCG@5	0.0082	0.0091	0.0080	0.0101	0.0121	0.0112	0.0068	0.0144		0.0105	0.0137	<u>0.0146</u>	0.0150	2.74% \uparrow	1.23e-2*
	NDCG@20	0.0143	0.0145	0.0156	0.0179	0.0223	0.0193	0.0127	0.0256		0.0186	<u>0.0263</u>	0.0261	0.0266	1.14% \uparrow	6.82e-3*

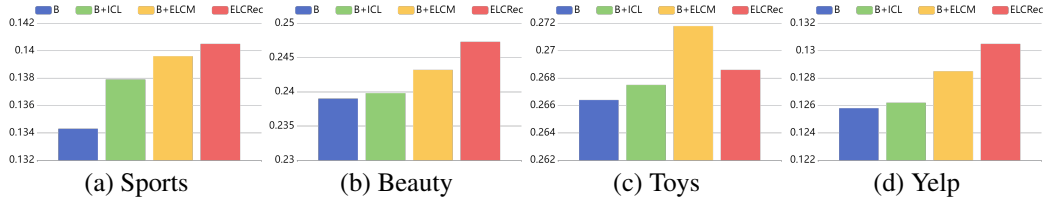


Figure 1: Ablation studies of the proposed end-to-end learnable cluster module (ELCM) and the intent-assisted contrastive learning (ICL). The results are the sum of four metrics, including HR@5, HR@20, NDCG@5, and NDCG@20.

[79, 29, 87, 32, 60, 85, 111, 98, 15]. Experimental results are the mean values of three runs. As shown in Table 1, the **bold values** and underlined values denote the best and runner-up results, respectively. From these results, we have four conclusions as follows. (a) The non-sequential model BPR-MF [79] has not achieved promising performance since the shallow method lacks the representation learning capability of users' historical behaviors. (b) The conventional sequential methods [29, 87, 32] improve the recommendation via different DNNs such as CNN [35], RNN [105], and Transformer [91]. But they perform worse since limiting self-supervision. (c) The recent methods [85, 111, 98] enhance the self-supervised capability of models via the self-supervised learning techniques. However, they neglect the underlying users' intent, thus leading to sub-optimal performance. (d) More recently, the intent learning methods [37, 11, 38, 15, 42, 46, 5] have been proposed to mine users' underlying intent to assist recommendation. Motivated by their success, we propose a new intent learning method termed ELCRec. Befitting from the strong intent learning capability of ELCRec, it surpasses all other intent learning methods.

To further verify the superiority of ELCRec, we conduct the t -test between the best and runner-up methods. As shown in Table 1, the most p -value is less than 0.05 except HR@5 on the Toys dataset. It indicates that ELCRec significantly outperforms runner-up methods. Overall, the extensive experiments demonstrate the superiority of ELCRec. In addition, we also conduct comparison experiments on recommendation datasets of other domains, including movie recommendation and news recommendation, as shown in the Appendix 7.4.1 and 7.4.2. These experimental results demonstrate a broader applicability of our proposed ELCRec.

4.3 Effectiveness

This section is dedicated to answering the research question (ii) and evaluating the effectiveness of the End-to-end Learnable Cluster Module (ELCM) and Intent-assisted Contrastive Learning (ICL). To achieve this, we conducted meticulous ablation studies on four benchmarks. Figure 1 illustrates the experimental results. In each sub-figure, "B", "B+ICL", "B+ELCM", and "ELCRec" correspond to the backbone, backbone with ICL, backbone with ELCM, and backbone with both ICL and ELCM,

respectively. Through the ablation studies, we draw three key conclusions. (a) “B+ICL” outperforms the backbone “B” on all four benchmarks. It indicates that the proposed ICL effectively improves behavior learning. (b) “B+ELCM” surpasses the backbone “B” significantly on all benchmarks. This phenomenon demonstrates that our proposed end-to-end learnable cluster module helps the model better capture the users’ underlying intents, thus improving recommendation performance. (c) ELCRec achieves the best performance on three out of four datasets. It shows the effectiveness of the combination of these two modules. On the Toys dataset, ELCRec can outperform the “B” and “B+ICL” but perform worse than “B+ELCM”. This phenomenon indicates it is worth researching the better collaboration of these two modules in the future. To summarize, these extensive ablation studies verify the effectiveness of the proposed intent-assisted contrastive learning and end-to-end learnable cluster module in ELCRec.

4.4 Efficiency

We test the efficiency of ELCRec on four benchmarks and answer the research question (iii). Concretely, the efficiency contains two perspectives, including running time costs (in second) and GPU memory costs (in MB). Note that we use the same epoch number of our method and the baseline when we test the running time. Besides, we calculate the average GPU memory cost during the training process. We have two observations as follows. (a) ELCRec can speed up ICLRec on three out of four datasets (See Table 2). Overall, on four datasets, the running time is decreased by 7.18% on average. The reason is that our proposed end-to-end optimization of intent learning breaks the alternative optimization of the EM framework, saving computation costs. (b) The results demonstrate that the GPU memory costs of our ELCRec are lower than that of ICLRec on four datasets (See Table 2). On average, the GPU memory costs are decreased by 9.58%. It is because we enable the model to conduct intent learning via the mini-batch users’ behaviors. Therefore, in summary, we demonstrate the efficiency of ELCRec from both time and memory aspects. Please note that, due to the relatively small size of the open benchmarks, the efficiency improvements are not particularly significant. However, on large-scale data, our method can achieve more substantial improvements.

Table 2: Running time and memory costs. **Bold values** denote better results.

Cost	Dataset	Sports	Beauty	Toys	Yelp	Average
Time	ICLRec	5282	3770	4374	4412	4460
	ELCRec	5360	2922	4124	4151	4139
	Improvement	1.48% ↑	22.49% ↓	5.72% ↓	5.92% ↓	7.18% ↓
Memory	ICLRec	1944	1798	2887	3671	2575
	ELCRec	1781	1594	2555	3383	2328
	Improvement	8.38% ↓	11.35% ↓	11.50% ↓	7.85% ↓	9.58% ↓

5 Application

Our proposed ELCRec is versatility and plug-and-play. Benefiting its advantages, we aim to apply it to real-time large-scale industrial recommendation systems with millions of users. First, we introduce the background and settings of the application. Then, we conduct extensive A/B testing and analyze the experimental results. Besides, due to the page limitation, we provide deployment details and practical insights in Appendix 7.11 and 7.8, respectively.

5.1 Application Background

The applied scenario is the livestreaming recommendation on the front page of the Alipay app. The user view (UV) and page view (PV) of this application are about 50 million and 130 million, respectively. Note that most users are new to this application, therefore leading to the sparsity of users’ behaviors. To solve this cold-start problem in the recommendation system, we adopt our proposed method to group users and recommend items based on the groups. Concretely, due to the sparsity of users’ behaviors, we first replace the users’ behavior with the users’ activities features in this application and model them via the multi-gate mixture-of-expert (MMOE) model [59]. Then we aim

Table 3: A/B testing on real-time large-scale industrial recommendation. **Bold values** denotes the significant improvements with $p < 0.05$. The symbol “-” denotes business secret.

Method	Livestreaming Metrics		Merchandise Metrics	
	PVCTR	VV	PVCTR	UVCTR
Baseline	-	-	-	-
Impro.	2.45% ↑	2.28% ↑	2.41% ↑	1.62% ↑

to group the users into various groups. For the existing intent learning methods, they are easily lead to the long-running time or the out-of-memory problems. To solve this problem we adopt the end-to-end learnable cluster module to group the users into various groups effectively and efficiently. Through this module, the high-activity users and new users are grouped into different clusters, alleviating the cold-start issue and assisting in better recommendations. Besides, during the learning process of the cluster embeddings, the low-activity users can transfer to high-activity users, improving the overall users’ activities in the application. Eventually, the networks are trained with multiple tasks. In the next section, we conduct experiments to demonstrate the effectiveness of our proposed method on real-time large-scale industrial data.

5.2 A/B Testing on Real-time Large-scale Data

We conduct A/B testing on the real-time large-scale industrial recommendation system. The experimental results are listed in Table 3. We evaluate the models with two metric systems, including livestreaming metrics and merchandise metrics. livestreaming metrics contain Page View Click Through Rate (PVCTR) and Video View (VV). Merchandise metrics contain PVCTR and User View Click Through Rate (UVCTR). The results indicate that our method can improve the recommendation performance of the baseline by about 2%. Besides, the improvements are significant with $p < 0.05$ in three out of four metrics.

In addition, to further explore why our method can work well in real-time large-scale recommendation systems, we further analyze the recommendation performance on different user groups. The results are shown in Table 4. Based on the users’ activity, we classify them into five groups, including Pure New users (PN), New users (N), Low-Activity users (LA), Medium-Activity users (MA), and High-Activity users (HA). Compared with the general recommendation algorithms that are unfriendly to new users, the experimental results show that our module not only improves the recommendation performance of high-activity users but also improves the recommendation performance of new users. Therefore, it can alleviate the cold-start problem and construct a more friendly user ecology.

Table 4: Results on different user groups. **Bold values** denotes improvements with $p < 0.05$.

Metric	PN	N	LA	MA	HA
PVCTR	6.96% ↑	1.67% ↑	1.98% ↑	0.35% ↑	19.02% ↑
VV	6.81% ↑	1.50% ↑	1.50% ↑	0.04% ↑	16.90% ↑

6 Conclusion

In this paper, we explore intent learning in recommendation systems. To be specific, we summarize and analyze two drawbacks of the existing EM optimization framework of intent learning. The complex and cumbersome alternating optimization limits the scalability and performance of existing methods. To this end, we propose a novel intent learning method termed ELCRec with an end-to-end learnable cluster module and intent-assisted contrastive learning. Extensive experiments on four benchmarks demonstrate ELCRec’s six abilities. In addition, benefiting from the versatility of ELCRec, we successfully apply it to the real-time large-scale industrial scenario and also achieve promising performance. Due to the limited pages, We discuss the limitations and future work of this paper in Appendix 7.12, such as pre-defined cluster number, limited recommendation domains, and uncontrollable update rate of cluster centers.

References

- [1] Abadi, M., Barham, P., Chen, J., Chen, Z., Davis, A., Dean, J., Devin, M., Ghemawat, S., Irving, G., Isard, M., et al. {TensorFlow}: a system for {Large-Scale} machine learning. In *12th USENIX symposium on operating systems design and implementation (OSDI 16)*, pp. 265–283, 2016.
- [2] Aflalo, A., Bagon, S., Kashti, T., and Eldar, Y. Deepcut: Unsupervised segmentation using graph neural networks clustering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 32–41, 2023.
- [3] Aggarwal, C. C. and Zhai, C. A survey of text clustering algorithms. *Mining text data*, pp. 77–128, 2012.
- [4] Asano, Y., Rupprecht, C., and Vedaldi, A. Self-labelling via simultaneous clustering and representation learning. In *International Conference on Learning Representations*, 2019.
- [5] Bai, Y., Zhou, Y., Dou, Z., and Wen, J.-R. Intent-oriented dynamic interest modeling for personalized web search. *ACM Transactions on Information Systems*, 2024.
- [6] Bartlett, P. L. and Mendelson, S. Rademacher and gaussian complexities: Risk bounds and structural results. *Journal of Machine Learning Research*, 3(Nov):463–482, 2002.
- [7] Brost, B., Mehrotra, R., and Jehan, T. The music streaming sessions dataset. In *The World Wide Web Conference*, pp. 2594–2600, 2019.
- [8] Caron, M., Bojanowski, P., Joulin, A., and Douze, M. Deep clustering for unsupervised learning of visual features. In *Proc. of ECCV*, 2018.
- [9] Caron, M., Misra, I., Mairal, J., Goyal, P., Bojanowski, P., and Joulin, A. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in neural information processing systems*, 33:9912–9924, 2020.
- [10] Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., and Joulin, A. Emerging properties in self-supervised vision transformers. In *Proceedings of the IEEE/CVF international conference on computer vision*, pp. 9650–9660, 2021.
- [11] Cen, Y., Zhang, J., Zou, X., Zhou, C., Yang, H., and Tang, J. Controllable multi-interest framework for recommendation. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 2942–2951, 2020.
- [12] Chang, B., Karatzoglou, A., Wang, Y., Xu, C., Chi, E. H., and Chen, M. Latent user intent modeling for sequential recommenders. In *Companion Proceedings of the ACM Web Conference 2023*, pp. 427–431, 2023.
- [13] Chang, J., Wang, L., Meng, G., Xiang, S., and Pan, C. Deep adaptive image clustering. In *Proceedings of the IEEE international conference on computer vision*, pp. 5879–5887, 2017.
- [14] Chang, J., Gao, C., Zheng, Y., Hui, Y., Niu, Y., Song, Y., Jin, D., and Li, Y. Sequential recommendation with graph neural networks. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*, pp. 378–387, 2021.
- [15] Chen, Y., Liu, Z., Li, J., McAuley, J., and Xiong, C. Intent contrastive learning for sequential recommendation. In *Proceedings of the ACM Web Conference 2022*, pp. 2172–2182, 2022.
- [16] Comaniciu, D. and Meer, P. Mean shift: A robust approach toward feature space analysis. *IEEE Transactions on pattern analysis and machine intelligence*, 24(5):603–619, 2002.
- [17] Dang, Y., Yang, E., Guo, G., Jiang, L., Wang, X., Xu, X., Sun, Q., and Liu, H. Ticoserec: Augmenting data to uniform sequences by time intervals for effective recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [18] del Barrio, E., Inouzhe, H., and Loubes, J.-M. Attraction-repulsion clustering with applications to fairness. *arXiv preprint arXiv:1904.05254*, 2019.

- [19] Dong, X., Song, X., Liu, T., and Guan, W. Prompt-based multi-interest learning method for sequential recommendation. *arXiv preprint arXiv:2401.04312*, 2024.
- [20] Ester, M., Kriegel, H.-P., Sander, J., Xu, X., et al. A density-based algorithm for discovering clusters in large spatial databases with noise. In *kdd*, volume 96, pp. 226–231, 1996.
- [21] Fan, L., Pu, J., Zhang, R., and Wu, X.-M. Neighborhood-based hard negative mining for sequential recommendation. *arXiv preprint arXiv:2306.10047*, 2023.
- [22] Fan, Z., Liu, Z., Wang, Y., Wang, A., Nazari, Z., Zheng, L., Peng, H., and Yu, P. S. Sequential recommendation via stochastic self-attention. In *Proceedings of the ACM Web Conference 2022*, pp. 2036–2047, 2022.
- [23] Guo, X., Gao, L., Liu, X., and Yin, J. Improved deep embedded clustering with local structure preservation. In *Proc. of IJCAI*, 2017.
- [24] Harper, F. M. and Konstan, J. A. The movielens datasets: History and context. *Acm transactions on interactive intelligent systems (tiis)*, 5(4):1–19, 2015.
- [25] Hartigan, J. A. and Wong, M. A. Algorithm as 136: A k-means clustering algorithm. *Journal of the royal statistical society. series c (applied statistics)*, 1979.
- [26] He, K., Chen, X., Xie, S., Li, Y., Dollár, P., and Girshick, R. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 16000–16009, 2022.
- [27] He, R. and McAuley, J. Fusing similarity models with markov chains for sparse sequential recommendation. In *2016 IEEE 16th international conference on data mining (ICDM)*, pp. 191–200. IEEE, 2016.
- [28] He, R. and McAuley, J. Ups and downs: Modeling the visual evolution of fashion trends with one-class collaborative filtering. In *proceedings of the 25th international conference on world wide web*, pp. 507–517, 2016.
- [29] Hidasi, B., Karatzoglou, A., Baltrunas, L., and Tikk, D. Session-based recommendations with recurrent neural networks. *arXiv preprint arXiv:1511.06939*, 2015.
- [30] Jacksi, K., Ibrahim, R. K., Zeebaree, S. R., Zebari, R. R., and Sadeeq, M. A. Clustering documents based on semantic similarity using hac and k-mean algorithms. In *2020 International Conference on Advanced Science and Engineering (ICOASE)*, pp. 205–210. IEEE, 2020.
- [31] Jing, M., Zhu, Y., Zang, T., and Wang, K. Contrastive self-supervised learning in recommender systems: A survey. *arXiv preprint arXiv:2303.09902*, 2023.
- [32] Kang, W.-C. and McAuley, J. Self-attentive sequential recommendation. In *2018 IEEE international conference on data mining (ICDM)*, pp. 197–206. IEEE, 2018.
- [33] Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [34] Kodinariya, T. M., Makwana, P. R., et al. Review on determining number of cluster in k-means clustering. *International Journal*, 1(6):90–95, 2013.
- [35] Krizhevsky, A., Sutskever, I., and Hinton, G. E. Imagenet classification with deep convolutional neural networks. *Advances in neural information processing systems*, 25, 2012.
- [36] Lee, S.-H. and Kim, C.-S. Deep repulsive clustering of ordered data based on order-identity decomposition. In *International Conference on Learning Representations*, 2020.
- [37] Li, C., Liu, Z., Wu, M., Xu, Y., Zhao, H., Huang, P., Kang, G., Chen, Q., Li, W., and Lee, D. L. Multi-interest network with dynamic routing for recommendation at tmall. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pp. 2615–2623, 2019.

- [38] Li, H., Wang, X., Zhang, Z., Ma, J., Cui, P., and Zhu, W. Intention-aware sequential recommendation with structured intent transition. *IEEE Transactions on Knowledge and Data Engineering*, 34(11):5403–5414, 2021.
- [39] Li, J., Zhou, P., Xiong, C., and Hoi, S. Prototypical contrastive learning of unsupervised representations. In *International Conference on Learning Representations*, 2020.
- [40] Li, M., Zhao, X., Lyu, C., Zhao, M., Wu, R., and Guo, R. Mlp4rec: A pure mlp architecture for sequential recommendations. *arXiv preprint arXiv:2204.11510*, 2022.
- [41] Li, M., Zhang, Z., Zhao, X., Wang, W., Zhao, M., Wu, R., and Guo, R. Automlp: Automated mlp for sequential recommendations. In *Proceedings of the ACM Web Conference 2023*, pp. 1190–1198, 2023.
- [42] Li, X., Sun, A., Zhao, M., Yu, J., Zhu, K., Jin, D., Yu, M., and Yu, R. Multi-intention oriented contrastive learning for sequential recommendation. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*, pp. 411–419, 2023.
- [43] Li, Y., Chen, T., Zhang, P.-F., and Yin, H. Lightweight self-attentive sequential recommendation. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 967–977, 2021.
- [44] Li, Y., Hu, P., Liu, Z., Peng, D., Zhou, J. T., and Peng, X. Contrastive clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 35, pp. 8547–8555, 2021.
- [45] Li, Y., Hao, Y., Zhao, P., Liu, G., Liu, Y., Sheng, V. S., and Zhou, X. Edge-enhanced global disentangled graph neural network for sequential recommendation. *ACM Transactions on Knowledge Discovery from Data*, 17(6):1–22, 2023.
- [46] Li, Z., Xie, Y., Zhang, W. E., Wang, P., Zou, L., Li, F., Luo, X., and Li, C. Disentangle interest trend and diversity for sequential recommendation. *Information Processing & Management*, 61(3):103619, 2024.
- [47] Liu, B., Bai, B., Xie, W., Guo, Y., and Chen, H. Task-optimized user clustering based on mobile app usage for cold-start recommendations. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 3347–3356, 2022.
- [48] Liu, Q., Wen, Y., Han, J., Xu, C., Xu, H., and Liang, X. Open-world semantic segmentation via contrasting and clustering vision-language embedding. In *European Conference on Computer Vision*, pp. 275–292. Springer, 2022.
- [49] Liu, Y., Tu, W., Zhou, S., Liu, X., Song, L., Yang, X., and Zhu, E. Deep graph clustering via dual correlation reduction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 7603–7611, 2022.
- [50] Liu, Y., Zhou, S., Liu, X., Tu, W., and Yang, X. Improved dual correlation reduction network. *arXiv preprint arXiv:2202.12533*, 2022.
- [51] Liu, Y., Liang, K., Xia, J., Yang, X., Zhou, S., Liu, M., Liu, X., and Li, S. Z. Reinforcement graph clustering with unknown cluster number. In *Proceedings of the 31st ACM International Conference on Multimedia*, pp. 3528–3537, 2023.
- [52] Liu, Y., Liang, K., Xia, J., Zhou, S., Yang, X., Liu, X., and Li, S. Z. Dink-net: Neural clustering on large graphs. *arXiv preprint arXiv:2305.18405*, 2023.
- [53] Liu, Y., Yang, X., Zhou, S., Liu, X., Wang, S., Liang, K., Tu, W., and Li, L. Simple contrastive graph clustering. *IEEE Transactions on Neural Networks and Learning Systems*, 2023.
- [54] Liu, Y., Yang, X., Zhou, S., Liu, X., Wang, Z., Liang, K., Tu, W., Li, L., Duan, J., and Chen, C. Hard sample aware network for contrastive deep graph clustering. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, pp. 8914–8922, 2023.
- [55] Liu, Z., Li, X., Fan, Z., Guo, S., Achan, K., and Philip, S. Y. Basket recommendation with multi-intent translation graph neural network. In *2020 IEEE International Conference on Big Data (Big Data)*, pp. 728–737. IEEE, 2020.

- [56] Liu, Z., Chen, Y., Li, J., Yu, P. S., McAuley, J., and Xiong, C. Contrastive self-supervised sequential recommendation with robust augmentation. *arXiv preprint arXiv:2108.06479*, 2021.
- [57] Liu, Z., Fan, Z., Wang, Y., and Yu, P. S. Augmenting sequential recommendation with pseudo-prior items via reversely pre-training transformer. In *Proceedings of the 44th international ACM SIGIR conference on Research and development in information retrieval*, pp. 1608–1612, 2021.
- [58] Ma, H., Xie, R., Meng, L., Chen, X., Zhang, X., Lin, L., and Kang, Z. Plug-in diffusion model for sequential recommendation. *arXiv preprint arXiv:2401.02913*, 2024.
- [59] Ma, J., Zhao, Z., Yi, X., Chen, J., Hong, L., and Chi, E. H. Modeling task relationships in multi-task learning with multi-gate mixture-of-experts. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining*, pp. 1930–1939, 2018.
- [60] Ma, J., Zhou, C., Yang, H., Cui, P., Wang, X., and Zhu, W. Disentangled self-supervision in sequential recommenders. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 483–491, 2020.
- [61] Ma, J., Zhou, C., Yang, H., Cui, P., Wang, X., and Zhu, W. Disentangled self-supervision in sequential recommenders. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 483–491, 2020.
- [62] McAuley, J., Targett, C., Shi, Q., and Van Den Hengel, A. Image-based recommendations on styles and substitutes. In *Proceedings of the 38th international ACM SIGIR conference on research and development in information retrieval*, pp. 43–52, 2015.
- [63] Mei, K. and Zhang, Y. Lightlm: A lightweight deep and narrow language model for generative recommendation. *arXiv preprint arXiv:2310.17488*, 2023.
- [64] Min, E., Guo, X., Liu, Q., Zhang, G., Cui, J., and Long, J. A survey of clustering with deep learning: From the perspective of network architecture. *IEEE Access*, 2018.
- [65] Mohri, M., Rostamizadeh, A., and Talwalkar, A. *Foundations of machine learning*. MIT press, 2018.
- [66] Nema, P., Karatzoglou, A., and Radlinski, F. Disentangling preference representations for recommendation critiquing with β -vae. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 1356–1365, 2021.
- [67] Nickerson, R. S. Confirmation bias: A ubiquitous phenomenon in many guises. *Review of general psychology*, 2(2):175–220, 1998.
- [68] Pan, S., Hu, R., Long, G., Jiang, J., Yao, L., and Zhang, C. Adversarially regularized graph autoencoder for graph embedding. *arXiv preprint arXiv:1802.04407*, 2018.
- [69] Pan, Z., Cai, F., Ling, Y., and de Rijke, M. An intent-guided collaborative machine for session-based recommendation. In *Proceedings of the 43rd international ACM SIGIR conference on research and development in information retrieval*, pp. 1833–1836, 2020.
- [70] Petrov, A. V. and Macdonald, C. gsasrec: Reducing overconfidence in sequential recommendation trained with negative sampling. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pp. 116–128, 2023.
- [71] Qian, Q. Stable cluster discrimination for deep clustering. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 16645–16654, 2023.
- [72] Qian, Q., Xu, Y., Hu, J., Li, H., and Jin, R. Unsupervised visual representation learning by online constrained k-means. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 16640–16649, 2022.
- [73] Qin, X., Yuan, H., Zhao, P., Liu, G., Zhuang, F., and Sheng, V. S. Intent contrastive learning with cross subsequences for sequential recommendation. In *Proceedings of the ACM international conference on web search and data mining*, 2024.

- [74] Qiu, R., Huang, Z., Yin, H., and Wang, Z. Contrastive learning for representation degeneration problem in sequential recommendation. In *Proceedings of the fifteenth ACM international conference on web search and data mining*, pp. 813–823, 2022.
- [75] Ren, X., Xia, L., Yang, Y., Wei, W., Wang, T., Cai, X., and Huang, C. Sslrec: A self-supervised learning library for recommendation. *arXiv preprint arXiv:2308.05697*, 2023.
- [76] Ren, X., Xia, L., Zhao, J., Yin, D., and Huang, C. Disentangled contrastive collaborative filtering. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pp. 1137–1146, 2023.
- [77] Rendle, S. Factorization machines. In *2010 IEEE International conference on data mining*, pp. 995–1000. IEEE, 2010.
- [78] Rendle, S., Freudenthaler, C., and Schmidt-Thieme, L. Factorizing personalized markov chains for next-basket recommendation. In *Proceedings of the 19th international conference on World wide web*, pp. 811–820, 2010.
- [79] Rendle, S., Freudenthaler, C., Gantner, Z., and Schmidt-Thieme, L. Bpr: Bayesian personalized ranking from implicit feedback. *arXiv preprint arXiv:1205.2618*, 2012.
- [80] Reynolds, D. A. Gaussian mixture models. *Encyclopedia of biometrics*, 2009.
- [81] Rodriguez, A. and Laio, A. Clustering by fast search and find of density peaks. *science*, 344(6191):1492–1496, 2014.
- [82] Ronen, M., Finder, S. E., and Freifeld, O. Deepdpm: Deep clustering with an unknown number of clusters. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 9861–9870, 2022.
- [83] Sabour, S., Frosst, N., and Hinton, G. E. Dynamic routing between capsules. *Advances in neural information processing systems*, 30, 2017.
- [84] Saeed, M. Y., Awais, M., Talib, R., and Younas, M. Unstructured text documents summarization with multi-stage clustering. *IEEE Access*, 8:212838–212854, 2020.
- [85] Sun, F., Liu, J., Wu, J., Pei, C., Lin, X., Ou, W., and Jiang, P. Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer. In *Proceedings of the 28th ACM international conference on information and knowledge management*, pp. 1441–1450, 2019.
- [86] Syakur, M., Khotimah, B., Rochman, E., and Satoto, B. D. Integration k-means clustering method and elbow method for identification of the best customer profile cluster. In *IOP conference series: materials science and engineering*, volume 336, pp. 012017. IOP Publishing, 2018.
- [87] Tang, J. and Wang, K. Personalized top-n sequential recommendation via convolutional sequence embedding. In *Proceedings of the eleventh ACM international conference on web search and data mining*, pp. 565–573, 2018.
- [88] Tanjim, M. M., Su, C., Benjamin, E., Hu, D., Hong, L., and McAuley, J. Attentive sequential models of latent intent for next item recommendation. In *Proceedings of The Web Conference 2020*, pp. 2528–2534, 2020.
- [89] Tran, N.-T. and Lauw, H. W. Learning multi-faceted prototypical user interests. In *The Twelfth International Conference on Learning Representations*, 2023.
- [90] Van der Maaten, L. and Hinton, G. Visualizing data using t-sne. *Journal of machine learning research*, 9(11), 2008.
- [91] Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [92] Von Luxburg, U. A tutorial on spectral clustering. *Statistics and computing*, 2007.

- [93] Wang, C., Pan, S., Hu, R., Long, G., Jiang, J., and Zhang, C. Attributed graph clustering: A deep attentional embedding approach. *arXiv preprint arXiv:1906.06532*, 2019.
- [94] Wang, S., Hu, L., Wang, Y., Sheng, Q. Z., Orgun, M., and Cao, L. Modeling multi-purpose sessions for next-item recommendations via mixture-channel purpose routing networks. In *International Joint Conference on Artificial Intelligence*. International Joint Conferences on Artificial Intelligence, 2019.
- [95] Wu, C.-Y., Ahmed, A., Beutel, A., Smola, A. J., and Jing, H. Recurrent recommender networks. In *Proceedings of the tenth ACM international conference on web search and data mining*, pp. 495–503, 2017.
- [96] Wu, F., Qiao, Y., Chen, J.-H., Wu, C., Qi, T., Lian, J., Liu, D., Xie, X., Gao, J., Wu, W., et al. Mind: A large-scale dataset for news recommendation. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pp. 3597–3606, 2020.
- [97] Xie, J., Girshick, R., and Farhadi, A. Unsupervised deep embedding for clustering analysis. In *Proc. of ICML*, 2016.
- [98] Xie, X., Sun, F., Liu, Z., Wu, S., Gao, J., Zhang, J., Ding, B., and Cui, B. Contrastive learning for sequential recommendation. In *2022 IEEE 38th international conference on data engineering (ICDE)*, pp. 1259–1273. IEEE, 2022.
- [99] Yang, J., Parikh, D., and Batra, D. Joint unsupervised learning of deep representations and image clusters. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 5147–5156, 2016.
- [100] Yang, Y., Huang, C., Xia, L., Huang, C., Luo, D., and Lin, K. Debiased contrastive learning for sequential recommendation. In *Proceedings of the ACM Web Conference 2023*, pp. 1063–1073, 2023.
- [101] Ye, Y., Xia, L., and Huang, C. Graph masked autoencoder for sequential recommendation. *arXiv preprint arXiv:2305.04619*, 2023.
- [102] Ye, Y., Xia, L., and Huang, C. Graph masked autoencoder for sequential recommendation. *arXiv preprint arXiv:2305.04619*, 2023.
- [103] Yu, J., Yin, H., Xia, X., Chen, T., Li, J., and Huang, Z. Self-supervised learning for recommender systems: A survey. *IEEE Transactions on Knowledge and Data Engineering*, 2023.
- [104] Yue, L., Jun, X., Sihang, Z., Siwei, W., Xifeng, G., Xihong, Y., Ke, L., Wenxuan, T., Wang, L. X., et al. A survey of deep graph clustering: Taxonomy, challenge, and application. *arXiv preprint arXiv:2211.12875*, 2022.
- [105] Zaremba, W., Sutskever, I., and Vinyals, O. Recurrent neural network regularization. *arXiv preprint arXiv:1409.2329*, 2014.
- [106] Zhai, S., Liu, B., Yang, D., and Xiao, Y. Group buying recommendation model based on multi-task learning. In *2023 IEEE 39th International Conference on Data Engineering (ICDE)*, pp. 978–991. IEEE, 2023.
- [107] Zhang, D. J., Hu, M., Liu, X., Wu, Y., and Li, Y. Netease cloud music data. *Manufacturing & Service Operations Management*, 24(1):275–284, 2022.
- [108] Zhang, X., Xu, S., Lin, W., and Wang, S. Constrained social community recommendation. In *Proceedings of the 29th ACM SIGKDD conference on knowledge discovery and data mining*, pp. 5586–5596, 2023.
- [109] Zhang, Y., Liu, Y., Xu, Y., Xiong, H., Lei, C., He, W., Cui, L., and Miao, C. Enhancing sequential recommendation with graph contrastive learning. *arXiv preprint arXiv:2205.14837*, 2022.

- [110] Zhang, Y., Wang, X., Chen, H., and Zhu, W. Adaptive disentangled transformer for sequential recommendation. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, pp. 3434–3445, 2023.
- [111] Zhou, K., Wang, H., Zhao, W. X., Zhu, Y., Wang, S., Zhang, F., Wang, Z., and Wen, J.-R. S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization. In *Proceedings of the 29th ACM international conference on information & knowledge management*, pp. 1893–1902, 2020.
- [112] Zhou, K., Yu, H., Zhao, W. X., and Wen, J.-R. Filter-enhanced mlp is all you need for sequential recommendation. In *Proceedings of the ACM web conference 2022*, pp. 2388–2399, 2022.
- [113] Zhou, P., Gao, J., Xie, Y., Ye, Q., Hua, Y., Kim, J., Wang, S., and Kim, S. Equivariant contrastive learning for sequential recommendation. In *Proceedings of the 17th ACM Conference on Recommender Systems*, pp. 129–140, 2023.
- [114] Zou, D., Zhao, S., Wei, W., Mao, X.-l., Li, R., Chen, D., Fang, R., and Fu, Y. Towards hierarchical intent disentanglement for bundle recommendation. *IEEE Transactions on Knowledge and Data Engineering*, 2023.

7 Appendix

7.1 Notation and Dataset

We list the basic notations in Table 5. And Table 6 summarizes the datasets’ details.

Table 5: Basic notations.

Notation	Meaning
\mathcal{U}	User set
\mathcal{V}	Item set
$\{S^u\}_{u=1}^{ \mathcal{U} }$	Users’ behavior sequence set
$(S^u)^{v_k}$	Users’ behavior sequence set in view k
d'	Dimension number of latent features
d	Dimension number of aggregated latent features
b	Batch size
k	Cluster number
T	Maximum sequence length
$\mathcal{L}_{\text{cluster}}$	Clustering loss
$\mathcal{L}_{\text{seq_cl}}$	Behavior sequence contrastive loss
$\mathcal{L}_{\text{intent_cl}}$	Intent contrastive loss
\mathcal{L}_{icl}	intent-assisted contrastive learning loss
$\mathcal{L}_{\text{next_item}}$	Next item prediction loss
$\mathcal{L}_{\text{overall}}$	Overall loss of the proposed ELCRec
\mathcal{F}	Behavior Encoder
\mathcal{P}	Concatenate pooling function
$\mathbf{E}^u \in \mathbb{R}^{ \mathcal{S}^u \times d'}$	Behavior sequence embedding of user u
$\mathbf{H} \in \mathbb{R}^{ \mathcal{U} \times d}$	Behavior embeddings of all users
$\hat{\mathbf{H}} \in \mathbb{R}^{ \mathcal{U} \times d}$	Normalized Behavior embeddings of all users
$\mathbf{H}^{v_k} \in \mathbb{R}^{ \mathcal{U} \times d}$	Behavior embeddings of all users in view v_k
$\mathbf{C} \in \mathbb{R}^{k \times d}$	Learnable cluster center embeddings
$\hat{\mathbf{C}} \in \mathbb{R}^{k \times d}$	Normalized Learnable cluster center embeddings

Table 6: Statistical information of four public datasets.

Dataset	#User	#Item	#Action	Avg. Len.	Sparsity
Sports	35,598	18,357	0.3M	8.3	99.95%
Beauty	22,363	12,101	0.2M	8.9	99.95%
Toys	19,412	11,924	0.17M	8.6	99.93%
Yelp	30,431	20,033	0.3M	8.3	99.95%

659 7.2 Algorithm Table

660 We summarize the overall process of the ELCRec method in Algorithm 1.

Algorithm 1 End-to-end Learnable Clustering Framework for Recommendation (ELCRec)

Input: user set \mathcal{U} ; item set \mathcal{V} ; historical behavior sequences $\{\mathcal{S}^u\}_{u=1}^{|\mathcal{U}|}$; cluster number k ; epoch number E ; learning rate; trade-off parameter α .

Output: Trained ELCRec.

- 1: Initialize model parameters in encoders.
 - 2: **for** epoch = 1, 2, ..., E **do**
 - 3: **for** $u = 1, 2, \dots, |\mathcal{U}|$ **do**
 - 4: Obtain u -th user’s behavior sequence embedding $\mathbf{E}^u \in \mathbb{R}^{|\mathcal{S}^u| \times d'}$ via encoding \mathcal{S}^u in Eq. (1).
 - 5: Obtain u -th user’s aggregated behavior embedding $\mathbf{h}_u \in \mathbb{R}^{1 \times d}$ via aggregating \mathbf{E}^u in Eq. (2)
 - 6: **end for**
 - 7: Obtain behavior embeddings of all users $\mathbf{H} \in \mathbb{R}^{|\mathcal{U}| \times d}$.
 - 8: Initialize cluster centers $\mathbf{C} \in \mathbb{R}^{k \times d}$ as learnable.
 - 9: Calculate clustering loss to conduct intent learning.
 - 10: Generate two views of behaviors via data augmentations.
 - 11: Encode the two views of the behavior sequences.
 - 12: Calculate $\mathcal{L}_{\text{seq_cl}}$ to conduct behavior contrastive learning.
 - 13: Query cluster index of the behavior embeddings via Eq. (5).
 - 14: Fuse the intent information to behavior embeddings.
 - 15: Calculate $\mathcal{L}_{\text{intent_cl}}$ to conduct intent contrastive learning.
 - 16: Calculate $\mathcal{L}_{\text{next_item}}$ to conduct next item prediction task.
 - 17: Minimize $\mathcal{L}_{\text{overall}}$ to train the model in Eq. (8).
 - 18: **end for**
 - 19: **Return** Well-trained ELCRec model.
-

662 7.3 Theoretical Analyses

663 In this subsection, we investigate the generalization bounds of the proposed clustering loss. Our
664 analysis is based on the Rademacher complexity and investigates how it improves the generalization
665 bound of the algorithm.

666 Without loss of generality, we have the following notation. Let $\mathbf{x} \in \mathcal{X}$ be the input, where \mathbf{x} are
667 generated from a underlying distribution $\mathbf{x} \sim \mathcal{P}$. Given n training samples $\mathcal{S} \triangleq \{\mathbf{x}_i\}_{i \in [n]}$ generated
668 from distribution \mathcal{P} , we denote its empirical distribution by \mathcal{P}^n . For every hyperparameter $\omega \in \Omega$,
669 we define \mathcal{F}_ω as a distribution-dependent hypothesis space corresponding to the ω , where Ω is a finite
670 set of hyperparameters. \mathcal{F}_ω is defined as $\{f_\omega | f_\omega = \mathcal{A}_\omega(\mathcal{S}), \mathcal{S} \in \mathcal{S}\}$, where \mathcal{A}_ω is an algorithm that
671 outputs the hypothesis f_ω given a dataset \mathcal{S} .

672 In the subsequent analysis, we denote $\mathcal{L}_{\text{cluster}}(\mathcal{S}, f_\omega) = \ell(f_\omega(\mathbf{x}, \mathbf{c}))$ as the proposed cluster loss
673 $\mathcal{L}_{\text{cluster}}$ with the embedding \mathbf{c} . Let u, v are the upper and lower bounds of the cluster loss re-
674 spectively. In other words, $u \geq \ell(f_\omega(\mathbf{x}, \mathbf{c})) \geq v$. In this paper, $u = 4$ and $v = -4$.

675 $\mathcal{R}_n^\ell(\mathcal{F}_\omega)$ is the rademacher complexity of the set $\{\mathbf{x} \mapsto \ell(f_\omega(\mathbf{x}, \mathbf{c}) : f_\omega \in \mathcal{F}_\omega\}$. Besides, we have
 676 $\mathbb{E}_{\mathbf{x} \sim \mathcal{P}^n} [\ell(f_\omega(\mathbf{x}, \mathbf{c}))] = \frac{1}{n} \sum_{i=1}^n \ell(f_\omega(\mathbf{x}_i, \mathbf{c}))$.

677 With the notation above, we have the following theorem.

678 **Theorem 7.1.** *For any $\delta > 0$ and $\omega \in \Omega$, for all $f_\omega \in \mathcal{F}_\omega$, with the probability at least $1 - \delta$, we*
 679 *have:*

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{P}} [\ell(f_\omega(\mathbf{x}, \mathbf{c}))] - \mathbb{E}_{\mathbf{x} \sim \mathcal{P}^n} [\ell(f_\omega(\mathbf{x}, \mathbf{c}))] \\ & \leq 2\sqrt{\frac{2\ln\Pi_{\mathcal{F}_\omega}(n)}{n}} + (u - v)\sqrt{\frac{\ln(1/\delta)}{2n}}. \end{aligned} \quad (9)$$

680 where $\ln\Pi_{\mathcal{F}_\omega}(n)$ denotes the growth function.

681 *Remark 7.2.* For each fixed \mathcal{F}_ω , the generalization bound in Theorem 1 goes to zero since
 682 $\ln\Pi_{\mathcal{F}_\omega}(n)/n \rightarrow 0$ and $\ln(1/\delta)/n \rightarrow 0$ when $n \rightarrow \infty$. In conclusion, the generation gap is ap-
 683 proximately $\mathcal{O}(1/\sqrt{n})$. Therefore, the generalization bound is promised.

684 To prove the above theorem, we need the following lemma.

685 **Lemma 7.3.** [6] *Let \mathcal{F} be a class of real-valued function that map from \mathcal{X} to $[v, u]$. Let \mathcal{D} be a*
 686 *probability distribution on $\mathcal{X} \times [v, u]$, and suppose that sample set $\mathbf{X} = \{x_1, x_2, \dots, x_n\}$ are chosen*
 687 *independently according to the distribution \mathcal{D} . For all $f \in \mathcal{F}$, with probability at least $1 - \delta$, we*
 688 *have:*

$$\Phi(S) \leq 2\mathcal{R}_n(\mathcal{F}) + (u - v)\sqrt{\frac{\ln(1/\delta)}{2n}}, \quad (10)$$

689 where $\Phi(S) = \sup_{f \in \mathcal{F}} (\mathbb{E}_{\mathbf{x} \sim \mathcal{P}} [f] - \mathbb{E}_{\mathbf{x} \sim \mathcal{P}^n} [f])$, $\mathcal{R}_n(\cdot)$ is the correspondent rademacher complex-
 690 ity.

691 **Lemma 7.4.** [65] *Let \mathcal{F} be the hypothesis space. The Rademacher complexity $\mathcal{R}_n(\mathcal{F})$ and the*
 692 *growth function $\Pi_{\mathcal{F}}(n)$ have:*

$$\mathcal{R}_n(\mathcal{F}) \leq \sqrt{\frac{2\ln\Pi_{\mathcal{F}}(n)}{n}}. \quad (11)$$

693 *Proof.* With the above lemma, we have the following derivation

$$\begin{aligned} \text{Let } \Phi(S) &= \sup_{f_\omega \in \mathcal{F}_\omega} (\mathbb{E}_{\mathbf{x} \sim P} [\mathcal{L}(f_\omega(\mathbf{x}, \mathbf{c}))] - \mathbb{E}_{\mathbf{x} \sim P^n} [\mathcal{L}(f_\omega(\mathbf{x}, \mathbf{c}))]) \\ &= \sup \left(\mathbb{E}_{\mathbf{x} \sim P} [\mathcal{L}(f_\omega(\mathbf{x}, \mathbf{c}))] - \frac{1}{n} \sum_{i=1}^n [\mathcal{L}(f_\omega(x_i, \mathbf{c}))] \right). \end{aligned} \quad (12)$$

694 We first provide an upper bound on $\Phi(S)$ by using McDiarmid's inequality. To apply McDiarmid's
 695 inequality, we compute an upper bound on $|\Phi(S) - \Phi(S')|$ where S and S' be two training datasets
 696 differing by exactly one point of an arbitrary index i_0 ; i.e., $\mathbf{x}_i = \mathbf{x}'_i$ for all $i \neq i_0$ and $\mathbf{x}_{i_0} \neq \mathbf{x}'_{i_0}$.

$$\begin{aligned} \text{Then, } |\Phi(S) - \Phi(S')| &= \left| \sup(\mathbb{E}_{\mathbf{x} \sim P} [\mathcal{L}(f_\omega(\mathbf{x}, \mathbf{c}))] - \frac{1}{n} \sum_{i=1}^n [\mathcal{L}(f_\omega(\mathbf{x}_i, \mathbf{c}))]) - \right. \\ & \quad \left. \sup(\mathbb{E}_{\mathbf{x} \sim P} [\mathcal{L}(f_\omega(\mathbf{x}, \mathbf{c}))] + \frac{1}{n} \sum_{i=1}^n [\mathcal{L}(f_\omega(\mathbf{x}'_i, \mathbf{c}))]) \right| \\ &\leq \frac{1}{n} \sup_{f_\omega \in \mathcal{F}} (|\mathcal{L}(f_\omega(\mathbf{x}_{i_0}, \mathbf{c})) - \mathcal{L}(f_\omega(\mathbf{x}'_{i_0}, \mathbf{c}))|) \\ &\leq \frac{u - v}{n}. \end{aligned} \quad (13)$$

697

□

698 In this way, $\Phi(S') - \Phi(S) \leq \frac{u-v}{n}$. We could obtain the similar bound $\Phi(S) - \Phi(S') \leq \frac{u-v}{n}$.
699 Therefore, for any $\delta > 0$, with Lemma A.3, at least the probability $1 - \delta$:

$$\Phi(S) \leq 2\mathcal{R}_n(\mathcal{F}_\omega) + (u - v)\sqrt{\frac{\ln(1/\delta)}{2n}}. \quad (14)$$

700 Furthermore, with Lemma A.4, we have:

$$\Phi(S) \leq 2\sqrt{\frac{2\ln\Pi_{\mathcal{F}}(n)}{n}} + (u - v)\sqrt{\frac{\ln(1/\delta)}{2n}}. \quad (15)$$

701 Based on above proof, we obtain that for any $\delta > 0$ and all $f_\omega \in \mathcal{F}_\omega$, with probability at least $1 - \delta$:

$$\begin{aligned} & \mathbb{E}_{\mathbf{x} \sim \mathcal{P}}[\ell(f_\omega(\mathbf{x}, \mathbf{c}))] - \mathbb{E}_{\mathbf{x} \sim \mathcal{P}^n}[\ell(f_\omega(\mathbf{x}, \mathbf{c}))] \\ & \leq 2\sqrt{\frac{2\ln\Pi_{\mathcal{F}}(n)}{n}} + (u - v)\sqrt{\frac{\ln(1/\delta)}{2n}}. \end{aligned} \quad (16)$$

702 7.4 Applicability on Diverse Domains

703 To further demonstrate the applicability of ELCRec on different recommendation domains, we
704 conduct additional experiments on movie recommendation and news recommendation.

705 7.4.1 Movie Recommendation

706 For the movie recommendation, we conducted experiments on the MovieLens 1M dataset (ML-1M)
707 [24]. This dataset contains 1M ratings from about 6K users on about 4K movies, as shown in Table 7.
708 In this experiment, we compared our proposed ELCRec with the most related baseline ICLRec. The
709 experimental results are presented in the Table 8.

Table 7: Statistical information of ML-1M dataset.

Dataset	#User	#Movie	#Rating	Rating per User	Rating per Movie
ML-1M	6,040	3,706	1,000,209	166	270

Table 8: Recommendation performance on ML-1M dataset. **Bold values** denote the best results. * indicates the p -value <0.05 .

Method	HR@5	HR@20	NDCG@5	NDCG@20
ICLRec	0.0293	0.0777	0.0186	0.0320
ELCRec	0.0333	0.0836	0.0208	0.0347
Impro.	13.65% \uparrow	7.59% \uparrow	11.83% \uparrow	8.44% \uparrow
p -value	4.03e-6*	6.68e-9*	6.36e-6*	1.66e-6*

710 From these experimental results, we draw two conclusions as follows.

- 711 (a) ELCRec achieves better recommendation performance, as evidenced by higher values for all
712 four metrics: HR@5, HR@20, NDCG@5, and NDCG@20. For example, with the HR@5
713 metric, ELCRec outperforms ICLRec by 13.65%.
- 714 (b) We calculated the p -value between our method and the runner-up. The results indicate that all
715 the p -values are less than 0.05, suggesting that our ELCRec significantly outperforms ICLRec.
- 716 (c) We demonstrate the applicability and superiority of the proposed ELCRec in the movie recom-
717 mendation domain.

7.4.2 News Recommendation

In addition, for news recommendation, we aim to conduct experiments on the MIND-small dataset [96]. MIND contains about 160k English news articles and more than 15 million impression logs generated by 1 million users. Every news article contains rich textual content including title, abstract, body, category and entities. Each impression log contains the click events, non-clicked events and historical news click behaviors of this user before this impression. To protect user privacy, each user was de-linked from the production system when securely hashed into an anonymized ID. MIND-small is a small version of the MIND dataset by randomly sampling 50,000 users and their behavior logs from the MIND dataset. We list the experimental results in Table 9.

Table 9: Recommendation performance on MIND-small dataset. **Bold values** denote the best results. * indicates the p -value <0.05 .

Method	HR@5	HR@20	NDCG@5	NDCG@20
ICLRec	0.0890	0.2128	0.0578	0.0926
ELCRec	0.0944	0.2332	0.0603	0.0994
Impr.	6.07% \uparrow	9.59% \uparrow	4.33% \uparrow	7.34% \uparrow
p -value	7.09e-17*	9.57e-09*	6.11e-7*	1.09e-7*

From these experimental results, we have three conclusions as follows.

- ELCRec supasses the runner-up for all four metrics, including HR@5, HR@20, NDCG@5, and NDCG@20. Significantly, ELCRec improve the runner-up by 9.59% with HR@20.
- We conduct t -test for ELCRec and the runner-up method and find all the p -values are less than 0.05. It indicates that our method significantly outperform the runner-up method.
- We demonstrate the applicability and superiority of the proposed ELCRec in the news recommendation domain.

Overall, we further demonstrate the applicability of ELCRec on diverse domains from the news and movie aspects.

7.5 Sensitivity

This section aims to answer the research question (iv). To verify the sensitivity of the proposed ELCRec to hyper-parameters, we test the performance on four datasets with different hyper-parameters. The experimental results are demonstrated in Figure 2. The x-axis denotes the values of hyper-parameters, and the y-axis denotes the values of the HR@5 metric. We obtain two conclusions as follows.

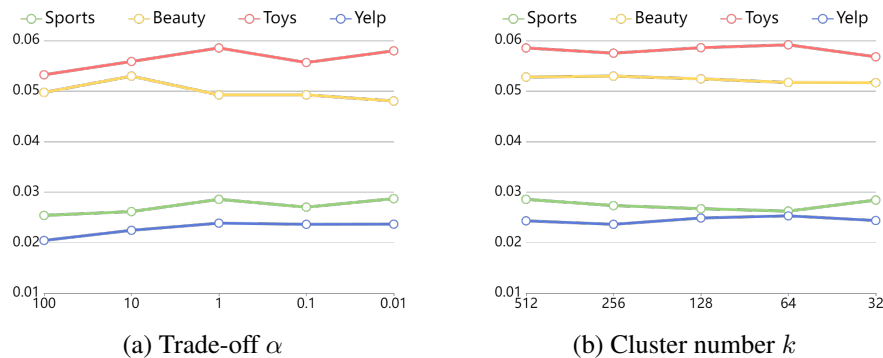


Figure 2: Sensitivity analyses of ELCRec. The results are evaluated by the HR@5 metric.

- For the trade hyper-parameter α , we test the performance with $\alpha \in \{0.01, 0.1, 1, 10, 100\}$. We find that our proposed ELCRec is not very sensitive to trade-off α . And ELCRec can achieve promising performance when $\alpha \in [0.1, 10]$.

(b) For the cluster number k , we test the recommendation performance with $\alpha \in \{32, 64, 128, 256, 512\}$. The results show that ELCRec is also not very sensitive to cluster number k and can perform well when $k \in [256, 512]$.

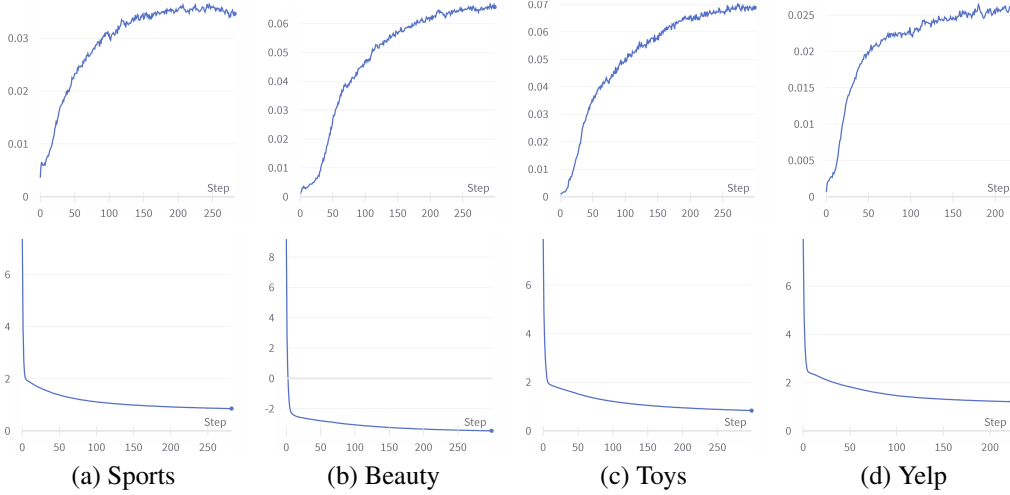


Figure 3: Convergence analyses. The first and second row denotes HR@5 on the evaluation set and training loss, respectively.

7.6 Convergence

To answer the research question (v), we monitor the recommendation performance and training loss as shown in Figure 3. We find that the losses gradually decrease and eventually converge. Besides, during the training process, the recommendation performance gradually increases and eventually reaches a promising value.

7.7 Visualization

We conduct visualization experiments on four public datasets to further demonstrate ELCRec’s capability to capture users’ underlying intents. Concretely, the learned behavior embeddings are visualized via t -SNE during training. As shown in Figure 6, the first row to the fourth row denotes the results on Sports, Beauty, Toys, and Yelp, respectively. From these experimental results, we have three observations as follows.

7.8 Practical Insights

In this section, we provide practical experiences and insights for the deployment of our proposed method. They contain three parts, including case study, solutions to rapid shift problem, and solutions to balance problem.

7.8.1 Case Study

To explore how our proposed method works well, we conduct case studies on large-scale industrial data. They contain two parts: case studies on user group distribution and case studies on the learned clusters.

Firstly, for the user group distribution, the results are demonstrated in Figure 4. We visualize the cluster distribution of different user groups. “top” denotes the cluster IDs that have the highest proportion in the user group. “bottom” denotes the cluster IDs that have the lowest proportion in the user group. From these analyses, we have two findings as follows.

- (a) As the user activity increases, the distribution becomes sharper. Namely, the users who have higher activities tend to distribute to one or two clusters. For example, about 60% of the high-activity users are attributed to cluster 10.

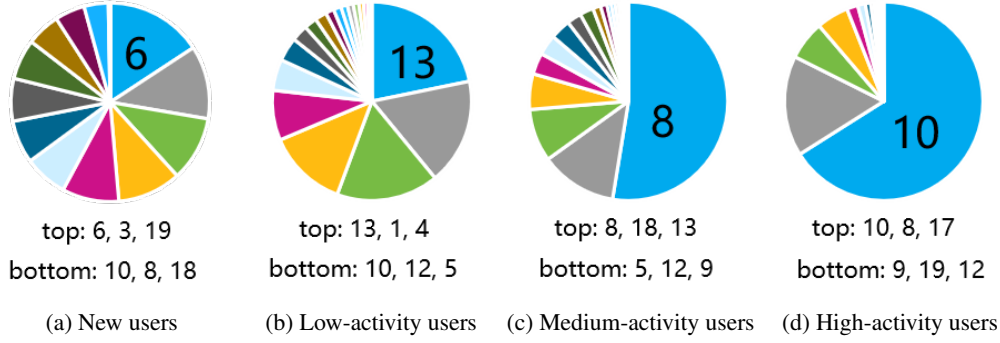


Figure 4: Case studies on different user groups. The distributions of different user groups are visualized. “top” denotes the cluster IDs, which have the highest proportion in the user group. “bottom” denotes the cluster IDs, which have the lowest proportion in the user group.

(b) The “top” cluster IDs of the high-activity user group, such as cluster 10 and cluster 8, are exactly the “bottom” cluster IDs of the low-activity user group. Similarly, the “bottom” cluster IDs of the high-activity user group, such as cluster 9, are exactly the “top” cluster IDs of the low-activity user group. This indicates that the learned cluster centers can well separate different user groups.

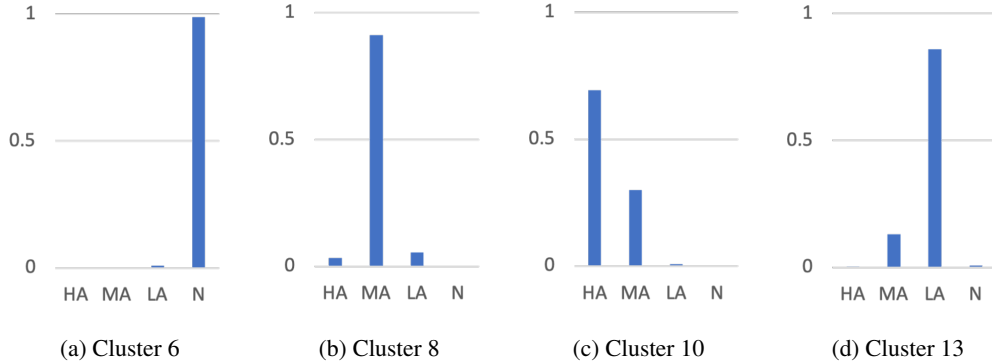


Figure 5: Case studies on the learned cluster. We visualize the distribution of the learned clusters. “HA”, “MA”, “LA”, and “N” denotes the high-activity, medium-activity, low-activity, and new user groups, respectively.

Secondly, we also conduct extensive case studies on the learned clusters. To be specific, we analyze the user distribution of each cluster, as shown in Figure 5. From the results, we can observe that, in cluster 6, most users are new. Besides, in the cluster 8, the most users are with medium activity. In addition, in cluster 10, most users are with high activity and medium activity. Moreover, in cluster 13, most users are with low activity and medium activity. Previous observations show that the learned centers can separate the users into different groups based on their activities.

In summary, these case studies further verify the effectiveness of ELCRec. Also, they provide insights for future work.

7.8.2 Solutions to Rapid Shift Problem

On real-time large-scale industrial data, the users’ behaviors and intents will shift rapidly. Therefore, we argue that the existing EM optimization can not capture the latest users’ intents, thus easily misunderstanding users and harming recommendation performance. Fortunately, our proposed ELCRec method can alleviate this problem. Concretely, the end-to-end learnable cluster module can guide the network to learn users’ intents dynamically. It can update the learned clusters (intents) at each batch, satisfying the requirement of rapid update. However, our method is hard to control the

794 update rate of the users' intents. That is one of drawbacks of ELCRec, we will discuss it and the
795 potential solution in 7.12.

796 7.8.3 Solutions to Balance Problem

797 Balancing the different loss functions in our model is indeed an important challenge. Our overall loss
798 function consists of next-item prediction loss, intent-assisted contrastive loss, and cluster loss. It is
799 formulated as follows: $\mathcal{L}_{\text{overall}} = \mathcal{L}_{\text{next_item}} + 0.1 \times \mathcal{L}_{\text{icl}} + \alpha \times \mathcal{L}_{\text{cluster}}$. We set the weight of \mathcal{L}_{icl} as
800 0.1 to maintain it in the same order of magnitude as the first term. This reduces the number of hyper-
801 parameters and simplifies the selection process. The weight of $\mathcal{L}_{\text{cluster}}$ is set as a hyper-parameter α .
802 We test different values of $\alpha \in \{0.01, 0.1, 1, 10, 100\}$ and find that our ELCRec method is not very
803 sensitive to the trade-off α . Promising performance is achieved when $\alpha \in [0.1, 10]$. The sensitivity
804 analysis experiments are presented in Figure 2 (b). In our proposed model, we set α to 1 for the
805 Sports and Toys datasets, 0.1 for the Yelp dataset, and 10 for the Beauty dataset. The selection of α is
806 mainly based on the model performance. We provide several practical strategies to balance multiple
807 losses in multi-task learning.

- 808 • **Weighted Balancing.** Assign weights to each loss function to control their contribution to the
809 overall loss. By adjusting the weights, a balance can be achieved between different loss functions.
810 This can be determined through prior knowledge, empirical rules, or methods like cross-validation.
- 811 • **Dynamic Weight Adjustment.** Adjust the weights of the loss functions in real time based on the
812 model's training progress or the characteristics of the data. For example, dynamically adjust the
813 weights based on the model's performance on a validation set, giving relatively smaller weights to
814 underperforming loss functions.
- 815 • **Multi-objective Optimization.** Treat different loss functions as multiple optimization objectives
816 and use multi-objective optimization algorithms to balance these objectives. This allows for the
817 simultaneous optimization of multiple loss functions and seeks balance between them.
- 818 • **Gradient-based Adaptive Adjustment.** Adaptively adjust the weights of loss functions based on
819 their gradients. If a loss function has a larger gradient, it may have a greater impact on the model's
820 training, and its weight can be increased accordingly.
- 821 • **Ensemble Methods.** Train multiple models based on different loss functions and use ensemble
822 learning techniques to combine their prediction results. By combining the predictions of different
823 models, a balance between different loss functions can be achieved.

824 In the future, we will continue to improve our model based the above strategies.

- 825 (a) At the beginning of training, the behavior embeddings are disorganized and can not reveal the
826 underlying intents.
- 827 (b) During the training process, the latent distribution is optimized, and similar behaviors are
828 grouped into latent intents.
- 829 (c) After optimization, the users' underlying intents appear, and we highlight them with circles in
830 Figure 6. These intents can assist recommendation systems in better modeling users' behavior
831 and recommending items. In summary, the above experiments and observations verify the
832 effectiveness of our proposed ELCRec.

833 7.9 Detailed Related Work

834 7.9.1 Sequential Recommendation

835 Sequential Recommendation (SR) poses a significant challenge as it strives to accurately capture
836 users' evolving interests and recommend relevant items by learning from their historical behavior
837 sequences. In the early stages, classical techniques such as Markov Chains and matrix factorization
838 have assisted models [27, 77, 78] in learning patterns from past transactions. Deep learning has
839 garnered significant attention in recent years and has demonstrated promising advancements across
840 various domains, including vision and language. Inspired by the remarkable success of Deep
841 Neural Networks (DNNs), researchers have developed a range of deep Sequential Recommendation
842 (SR) methods. For instance, Caser [87] leverages Convolutional Neural Networks (CNNs) [35] to
843 embed item sequences into an "image" representation over time, enabling the learning of sequential

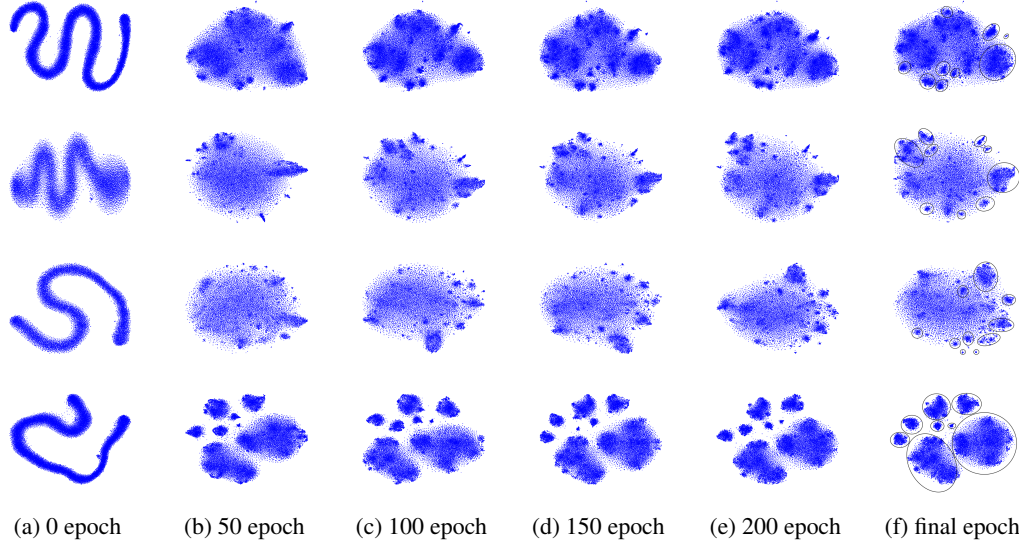


Figure 6: t -SNE visualization on four public datasets. The first row to the fourth row denotes the results on Sports, Beauty, Toys, and Yelp.

patterns through convolutional filters. Similarly, GRU4Rec [29] utilizes Recurrent Neural Networks (RNNs) [105], specifically the Gated Recurrent Unit (GRU), to model entire user sessions. The Transformer architecture [91] has also gained significant popularity and has been extended to the recommendation domain. For example, SASRec [32] employs a unidirectional Transformer to model users' behavior sequences, while BERT4Rec [85] utilizes a bidirectional Transformer to encode behavior sequences from both directions. To enhance the time and memory efficiency of Transformer-based SR models, LSAN [43] introduces aggressive compression techniques for the original embedding matrix. Addressing the cold-start issue in SR models, ASReP [57] proposes a pre-training and fine-tuning framework. Furthermore, researchers have explored the layer-wise disentanglement of architectures [110] and introduced novel modules like the Wasserstein self-attention module in STOSA [22] to model item-item position-wise relationships. In addition to Transformers, graph neural networks [101, 109, 45, 14] and multilayer perceptrons [41, 40, 112] have also found applications in recommendation systems. More recently, Self-Supervised Learning (SSL) [103, 75], particularly contrastive learning [31], has gained popularity due to its ability to learn patterns from large-scale unlabeled data. As a result, SSL-based SR models have been increasingly introduced. For instance, in CoSeRec [56], Liu et al. propose two informative augmentation operators that leverage item correlations to generate high-quality views. They then utilize contrastive learning to bring positive sample pairs closer while pushing negative pairs apart. Subsequently, TiCoSeRec [17] is designed by considering the time intervals in the behavior sequences. Another contrastive SR method, ECL-SR [113], ensures that the learned embeddings are sensitive to invasive augmentations while remaining insensitive to mild augmentations. Additionally, DCR [100] addresses the issue of popularity bias through a debiased contrastive learning framework. Moreover, DuoRec [74] is proposed to solve the representation degeneration problem in contrastive recommendation methods. Techniques such as hard negative mining [21, 70] have also proven beneficial for recommendation systems. Besides, motivated by the success of Mask Autoencoder (MAE) [26], MAERec [102] is proposed with the graph masked autoencoder.

7.9.2 Intent Learning for Recommendation

The preferences of users towards items are implicitly reflected in their intents. Recent studies [37, 11, 38, 15, 42, 46, 5] have highlighted the significance of users' intents in the user understanding and enhancing the performance of recommendation systems. For instance, MCPN [94] introduces a mixture-channel method to model subsets of items with multiple purposes. Inspired by capsule networks [83], MIND [37] utilizes dynamic routing to learn users' multiple interests. Furthermore, ComiRec [11] employs a multi-interest module to capture diverse interests from user behavior se-

quences, while the aggregation module combines items from different interests to generate overall recommendations. Besides, MITGNN [55] treats intents as translated tail entities and learns embeddings using graph neural networks. In addition, Pan et al. [69] propose an intent-guided neighbor detector to identify relevant neighbors, followed by a gated fusion layer that adaptively combines the current session with the neighbor sessions. Moreover, Ma et al. [60] aims to disentangle the intentions underlying users’ behaviors and construct sample pairs within the same intention. Meanwhile, the ASLI method [88] incorporates a temporal convolutional network layer to extract latent users’ intents. More recently, a general latent learning framework called ICLRec [15] is introduced, which utilizes contrastive learning and k -Means clustering to group the users’ behaviors to intents. Chang et al. [12] formulate users’ intents as latent variables and infer them based on user behavior signals using the Variational Auto Encoder (VAE) [33]. To mitigate noise caused by data augmentations in contrastive SR models, IOCRec [42] proposes building high-quality views at the intent level. Besides, ICSRec [73] is proposed to solve this issue by conducting contrastive learning on cross sub-sequences. DIMPS [5] aims to build dynamic and intent-oriented document representations for intent learning. PoMRec [19] insert the specific prompts into user interactions to make them adaptive to different learning objectives. Furthermore, Teddy [46] is proposed by utilizing the intent trend and diversity.

Firstly, we want to clearly claim the target of this paper and the demand of the industrial scenario as follows. 1) For the open benchmarks, we aim to develop an intent learning method to decoupling user’s intents for better recommendation based the appropriate intents of the user. 2) For the industrial data, we aim to design a user grouping method to cluster the users into different groups to solve the cold-start problem via mapping the new users into the user group, which contains more useful information. Therefore, the designed method needs to have the following abilities. 1) It can explicitly decouple users’ behaviours into different intents (grouping users into different clusters). 2) It can be easily adopted to the large-scale real-time industrial data, saving the memory and time costs. Secondly, we surveyed massive recent state-of-the-art methods to solve the above challenges in the related work part of this paper. We highlight the drawbacks of the related methods [42] [3] and claim why they will fail in our scenario. In the IOCRec method [42], they define the prototype intention of users as a $k \times d$ matrix. And the these prototype intention are directly used to calculate the relevance weights and the intentions. However, there are no designs for the initialization and optimization of the prototype intention, e.g., guiding the prototype intention to represent the users’ behaviours, and different intentions are separated. Therefore, it lacks explainability and persuasiveness, especially in the scenario where there is a need to conduct different recommendation strategies for different groups, i.e., user grouping recommendations. Also, we do not find theoretical or experimental evidence to support that the learned intents are separated well and reveal the representative behaviours of users in the original paper [42]. For the DCCF method [76], 1) it is based on the graph neural networks, limiting the model scalability and efficiency on large-scale data due to the large costs of graph constructing, graph storage, and neighbour sampling. And the sequential methods are more efficient since our data is naturally the sequences of the user behaviors. 2) Besides, in the DCCF method, the intents are randomly initialized via xavier normalization. Then, they are used to aggregate information. In the loss function part, we notice that there is only a penalty item to limit the complexity of the parameters of intent embeddings. Thus, there are no operations or loss functions to explicitly optimize the users’ intents, such as separating different intents, learning intents from behaviours, etc. We claim this intent decoupling is relatively weak and may not really learn well and separate the different intents of users. Also, in Figure 4 of the original paper [76], we find that the cluster pattern is not revealed well in the sampled data. We speculate the cluster pattern will also not be revealed well on the whole samples of the datasets. Thirdly, we explain why we chose ICLRec [15] as our baseline. 1) ICLRec is a sequential recommendation method, which is more suitable for our data. Compared to the GNN-based methods, it can save more time and memory costs. 2) ICLRec adopt the clustering algorithm to explicitly separate the users’ intents, which can also be adapted for user grouping. It explicitly optimizes the intents based on the users’ behaviour embeddings. We believe this technique can better separate the users’ intents well and also better obtain the users’ intents from their behaviors. In Figure 7 of the original paper [15], we find that ICLRec can reveal the cluster pattern well on the sampled data. Fourthly, we claim our motivation. Although ICLRec can achieve promising performance and effectively decouple users’ intents, the EM optimization framework limits the scalability and performance. 1) At the E-step, we need to apply the clustering algorithm on the whole data, limiting the model’s scalability, especially in large-scale industrial scenarios, e.g., apps with billion users. 2) In the EM framework, the optimization of behaviour learning and the clustering algorithm are separated, leading to sub-optimal performance

and increasing the implementation difficulty. We admit that our analyses of the problems start from ICLRec methods. But, actually, there are many intent learning methods [73, 61, 63, 66, 89] that adopt the clustering algorithms and the EM framework. They will meet the above problems and may fail when scaling to real-time large-scale data. Therefore, we claim our mentioned challenges are general recommendation systems, especially for intent decoupling methods. And we believe our proposed end-to-end learnable clustering module can bring performance improvement and saving time and space costs for these methods.

7.9.3 Clustering Algorithm

Clustering is a fundamental and challenging task that aims to group samples into distinct clusters without supervision. By leveraging the power of unlabeled data, clustering algorithms have found applications in various domains, including computer vision [13], natural language processing [3], graph learning [53], and recommendation systems [15, 73]. In the early stages, several traditional clustering methods [25, 92, 80, 20, 81] were proposed. For instance, the classical k -Means clustering [25] iteratively updates cluster centers and assignments to group samples. Spectral clustering [92] constructs a similarity graph and utilizes eigenvalues and eigenvectors to perform clustering. Additionally, probability-based Gaussian Mixture Models (GMM) [80] assume that the data distribution is a mixture of Gaussian distributions and estimate parameters through maximum likelihood. Moreover, the repulsive clustering methods [36, 18, 2] cluster data via the repulsive terms. In contrast, density-based methods [20, 81, 16] overcome the need for specifying the number of clusters as a hyperparameter. In recent years, the impressive performance of deep learning has sparked a growing interest in deep clustering [44, 82, 64, 4, 72, 39]. For instance, Xie et al. propose DEC [97], a deep learning-based approach for clustering. They initialize cluster centers using k -Means clustering and optimize the clustering distribution using a Kullback-Leibler divergence clustering loss [97]. IDEC [23] improves upon DEC by incorporating the reconstruction of original information from latent embeddings. JULE [99] and DeepCluster [8] both adopt an iterative approach, updating the deep network based on learned data embeddings and clustering assignments. SwAV [9], an online method, focuses on clustering data and maintaining consistency between cluster assignments from different views of the same image. DINO [10] introduces a momentum encoder to address representation collapse. Additionally, SeCu [71] proposes a stable cluster discrimination task and a hardness-aware clustering criterion. While deep clustering has been extensively applied to image data, it is also utilized in graph clustering [49, 50, 93, 104, 68, 53, 54, 52] and text clustering [3, 48, 30, 84]. However, the application of clustering-based recommendation [15, 73] is relatively unexplored. Leveraging the unsupervised learning capabilities of clustering could benefit intent learning in recommendation systems.

7.10 Implementation Details of Baselines

For the baseline methods, we adopt the public source code with the default parameter settings and reproduce their results on the used four benchmarks. The source codes of these methods are available at Table 10. Besides, for the used benchmarks, following [15], we only kept datasets where all users and items have at least five interactions. Besides, we adopted the dataset split settings used in [15]. The Sports, Beauty, and Toys datasets [62, 28] are obtained from: <http://jmcauley.ucsd.edu/data/amazon/index.html>. The yelp dataset is obtained from <https://www.yelp.com/dataset>.

For the results which have already existed in the original papers, we reuse them in our paper. For the results that do not exist in the original papers, we adopt the official codes of the baselines to reproduce the experimental results. Concretely, for the hyperparameters, we adopt and try several sets of the default hyperparameters on different datasets released by the original authors. We report the best result obtained from the best hyper-parameters. By the way, we also observe these results have already converged well. Besides, we conducted three runs on different random seeds for all experimental results and reported the average performance.

7.11 Deployment Details

We aim to apply our proposed method to the real-time large-scale industrial recommendation systems. Concretely, the ELCRec algorithm is applied to livestreaming recommendation in the front page of the Alipay app. The user view (UV) and page view (PV) of this application are about 50 million

Table 10: Implementation URLs of baseline methods.

Method	Url
BPR-MF [79]	https://github.com/xiangwang1223/neural_graph_collaborative_filtering
GRU4Rec [29]	https://github.com/slientGe/Sequential_Recommendation_Tensorflow
Caser [87]	https://github.com/graytowne/caser_pytorch
SASRec [32]	https://github.com/kang205/SASRec
BERT4Rec [85]	https://github.com/FeiSun/BERT4Rec
DSSRec [60]	https://github.com/abinashsinha330/DSSRec
S3-Rec [111]	https://github.com/RUCAIBox/CIKM2020-S3Rec
CL4SRec [98]	https://github.com/HKUDS/SSLRec
ICLRec [15]	https://github.com/salesforce/ICLRec
DCRec [100]	https://github.com/HKUDS/DCRec
MAERec [102]	https://github.com/HKUDS/MAERec
IOCRec [42]	https://github.com/LFM-bot/IOCRec

and 130 million, respectively. Since most of the users are new to this application, it easily leads to the sparsity of users’ behaviors, namely, the cold-start problem in recommendation systems. Our proposed ELCRec can alleviate this problem by grouping users and then making recommendations. This method can map a new user to a user group, which contains more intent behaviour information from similar users, such as other similar new users and similar users with low/middle activities. In this manner, we can guide the model to learn the behaviour of new users and provide more precise recommendations for them even with the sparse behaviours.

At first, we introduce the online baseline of this project. Since the sparsity of the users’ behaviors, we replaced the users’ behaviors with the users’ activities. Then, the online baseline multi-gate mixture-of-expert (MMOE) [59] models the users’ activities. In this model, the experts are designed to extract the features of users, and the multi-gates are designed to select specific experts. The inputs of the multi-gates are the activities of the users. This design aims to train an activity-awarded model to group different users and then conduct recommendations.

However, we found the performance of this model is limited, and the output of the gates is smooth, indicating that this model may fail to group users. Meanwhile, on the open benchmarks, extensive experiments demonstrate the proposed end-to-end learnable clustering module is effective and scalable. Thus, to solve the above problem, ELCRec is adopted in this project. It is designed to assist the gate to group users. For example, the high-activity users and new users are grouped into different clusters, and then the users in different groups will be recommended differently. Therefore, it alleviates the cold-start issue and further improves the recommendation performance. Besides, during the learning process of the cluster embeddings, the low-activity users can transfer to high-activity users, improving the overall users’ activities in the application. It is worth mentioning that the networks are trained with multi-task targets, e.g., CTR prediction, CVR prediction, etc. Following the previous online baseline, the method is implemented with the TensorFlow deep learning platform [1].

7.12 Limitations & Future Work

In this paper, we propose a novel intent learning method named ELCRec based on the end-to-end learnable clustering framework. It can better mine users’ underlying intents via unifying representation learning and clustering optimization. Besides, the end-to-end learnable clustering module optimizes the clustering distribution via mini-batch data, thus improving the scalability and convenience of deployment. Moreover, we demonstrate the superiority, effectiveness, efficiency, sensitivity, convergence, and visualization of ELCRec on four benchmarks. ELCRec is also successfully applied in the real-time large-scale industrial recommendation system. Although achieving promising results, we admit the proposed ELCRec algorithm has several limitations and drawbacks. We summarize them as follows.

- **Pre-defined Cluster Number.** The cluster number in ELCRec is a pre-defined hyper-parameter. In the real-time large scale data, it is hard to determine the cluster number, especially under the unsupervised conditions. In this paper, for the open benchmarks, we search the cluster number in

- 1025 {32, 64, 128, 256, 512}. Besides, for the industrial application, the cluster number is set to 20
 1026 based on the number of user groups. However, either the search method or the expert knowledge
 1027 can not determine the cluster number well at once. The cluster number may change dynamically
 1028 during model training, and the proposed method may fail to achieve promising performance.
- 1029 • Limited Recommendation Domains. In this paper, we adopt four recommendation benchmarks,
 1030 including Sports, Beauty, Toys, and Yelp, for the main experimental results. But, these four
 1031 datasets are all buying recommendation datasets. Besides, we adopt ML-1M [24] and MIND-
 1032 small [96] for the movie and news recommendation for the additional experiments. However, the
 1033 recommendation domains are still limited. In the future, we can further demonstrate the boarder
 1034 applicability of ELCRec in other domains.
 - 1035 • Uncontrollable Update Rate of Cluster Centers. In the real-time recommendation system, the users'
 1036 behaviors and intents usually change rapidly. Although our proposed ELCRec can dynamically
 1037 learn the users' intents, it is hard to control the update rate of the underlying clusters (intents).

1038 To solve these issues, we summarize several future works and the potential technical solutions as
 1039 follows.

- 1040 • Density-based Clustering. As mentioned above, the cluster number is a pre-defined value in this
 1041 paper, limiting the recommendation performance and flexibility of the method. To solve this
 1042 issue in the future, firstly, we can determine the cluster number based on some cluster number
 1043 estimation methods. They can help to determine the cluster number by performing multiple
 1044 clustering runs and selecting the best cluster number based on the unsupervised criterion. The
 1045 mainstream cluster number estimation methods [34] include the thumb rule, ELBOW [86], t -SNE
 1046 [90], etc. The thumb rule simply assigns the cluster number k with $\sqrt{n/2}$, where n is the number
 1047 of samples. This manual setting is empirical and can not be applicable to all datasets. Besides, the
 1048 ELBOW is a visual method. Concretely, they start the cluster number $k = 2$ and keep increasing
 1049 k in each step by 1, calculating the WSS (within-cluster sum of squares) during training. They
 1050 choose the value of k when the WSS drops dramatically, and after that, it reaches a plateau.
 1051 However, it will bring large computational costs since the deep neural network needs to be trained
 1052 with repeated times. Another visual method termed t -SNE visualizes the high-dimension data
 1053 into 2D sample points and helps researchers determine the cluster number. The effectiveness of
 1054 t -SNE heavily relies on the experience of researchers. Therefore, secondly, we can determine the
 1055 cluster number based on the data density [81, 82]. Concretely, the areas with high data density
 1056 are identified as the cluster centers, while the areas with low data density are identified as the
 1057 decision boundaries between cluster centers. Besides reinforcement learning is also a potential
 1058 solution [51]. Through these designs, the cluster number will be changeable during the training
 1059 process. It will be determined based on the embeddings itself, better revealing the users' behavior
 1060 and may achieve better recommendation performance.
- 1061 • More Recommendation Domains. As mentioned above, the applied recommendation domains
 1062 of our method are limited. We aim to test ELCRec on more recommendation domains, such as
 1063 music recommendation [107, 7], group recommendation [108, 47], group buying [106], bundle
 1064 recommendation [114], etc.
- 1065 • Controllable Intent Learning. As mentioned above, in the real-time recommendation system, the
 1066 intents of the users may change rapidly. Our method makes it hard to control the intent update
 1067 rate during training and inference. To this end, in the future, we can propose a controllable
 1068 cluster center learning method, such as the momentum update, to control the change rate of the
 1069 users' intents. Concretely, $\mathbf{C}_t = m \cdot \mathbf{C}_t + (1 - m) \cdot \mathbf{C}_{t-1}$. Here, \mathbf{C}_t denote the cluster center
 1070 embeddings at t and m denotes the momentum. Then, the cluster centers (intents of users) will
 1071 be changed rapidly when m is large, and the cluster centers (intents of users) will be changed
 1072 slowly when m is small. This strategy will control the change rate of the users' intent embeddings,
 1073 therefore alleviating the above problem.

1074 NeurIPS Paper Checklist

1075 The checklist is designed to encourage best practices for responsible machine learning research,
 1076 addressing issues of reproducibility, transparency, research ethics, and societal impact. Do not remove
 1077 the checklist: **The papers not including the checklist will be desk rejected.** The checklist should

1078 follow the references and precede the (optional) supplemental material. The checklist does NOT
1079 count towards the page limit.

1080 Please read the checklist guidelines carefully for information on how to answer these questions. For
1081 each question in the checklist:

- 1082 • You should answer [Yes], [No], or [NA].
- 1083 • [NA] means either that the question is Not Applicable for that particular paper or the
1084 relevant information is Not Available.
- 1085 • Please provide a short (1–2 sentence) justification right after your answer (even for NA).

1086 **The checklist answers are an integral part of your paper submission.** They are visible to the
1087 reviewers, area chairs, senior area chairs, and ethics reviewers. You will be asked to also include it
1088 (after eventual revisions) with the final version of your paper, and its final version will be published
1089 with the paper.

1090 The reviewers of your paper will be asked to use the checklist as one of the factors in their evaluation.
1091 While "[Yes]" is generally preferable to "[No]", it is perfectly acceptable to answer "[No]" provided a
1092 proper justification is given (e.g., "error bars are not reported because it would be too computationally
1093 expensive" or "we were unable to find the license for the dataset we used"). In general, answering
1094 "[No]" or "[NA]" is not grounds for rejection. While the questions are phrased in a binary way, we
1095 acknowledge that the true answer is often more nuanced, so please just use your best judgment and
1096 write a justification to elaborate. All supporting evidence can appear either in the main paper or the
1097 supplemental material, provided in appendix. If you answer [Yes] to a question, in the justification
1098 please point to the section(s) where related material for the question can be found.

1099 **IMPORTANT, please:**

- 1100 • **Delete this instruction block, but keep the section heading “NeurIPS paper checklist”,**
- 1101 • **Keep the checklist subsection headings, questions/answers and guidelines below.**
- 1102 • **Do not modify the questions and only use the provided macros for your answers.**

1103 1. Claims

1104 Question: Do the main claims made in the abstract and introduction accurately reflect the
1105 paper’s contributions and scope?

1106 Answer: [Yes]

1107 Justification: See the abstract and introduction part. We propose a novel intent learning
1108 method termed ELCRec, by unifying behavior representation learning into an end-to-end
1109 learnable clustering framework, for effective and efficient Recommendation. We clearly
1110 introduce the existing methods and their drawbacks. To solve the problem, we design the
1111 corresponding novel modules. And experimental results and theoretical analyses demonstrate
1112 ELCRec from six aspects.

1113 Guidelines:

- 1114 • The answer NA means that the abstract and introduction do not include the claims
1115 made in the paper.
- 1116 • The abstract and/or introduction should clearly state the claims made, including the
1117 contributions made in the paper and important assumptions and limitations. A No or
1118 NA answer to this question will not be perceived well by the reviewers.
- 1119 • The claims made should match theoretical and experimental results, and reflect how
1120 much the results can be expected to generalize to other settings.
- 1121 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
1122 are not attained by the paper.

1123 2. Limitations

1124 Question: Does the paper discuss the limitations of the work performed by the authors?

1125 Answer: [Yes]

Justification: See section 7.12: Limitations & Future work. We summarize the drawbacks of our proposed method, such as, pre-defined cluster number, limited recommendation domains, and uncontrollable update rate of cluster centers. And then we provide the potential solutions.

Guidelines:

- The answer NA means that the paper has no limitation while the answer No means that the paper has limitations, but those are not discussed in the paper.
- The authors are encouraged to create a separate "Limitations" section in their paper.
- The paper should point out any strong assumptions and how robust the results are to violations of these assumptions (e.g., independence assumptions, noiseless settings, model well-specification, asymptotic approximations only holding locally). The authors should reflect on how these assumptions might be violated in practice and what the implications would be.
- The authors should reflect on the scope of the claims made, e.g., if the approach was only tested on a few datasets or with a few runs. In general, empirical results often depend on implicit assumptions, which should be articulated.
- The authors should reflect on the factors that influence the performance of the approach. For example, a facial recognition algorithm may perform poorly when image resolution is low or images are taken in low lighting. Or a speech-to-text system might not be used reliably to provide closed captions for online lectures because it fails to handle technical jargon.
- The authors should discuss the computational efficiency of the proposed algorithms and how they scale with dataset size.
- If applicable, the authors should discuss possible limitations of their approach to address problems of privacy and fairness.
- While the authors might fear that complete honesty about limitations might be used by reviewers as grounds for rejection, a worse outcome might be that reviewers discover limitations that aren't acknowledged in the paper. The authors should use their best judgment and recognize that individual actions in favor of transparency play an important role in developing norms that preserve the integrity of the community. Reviewers will be specifically instructed to not penalize honesty concerning limitations.

3. Theory Assumptions and Proofs

Question: For each theoretical result, does the paper provide the full set of assumptions and a complete (and correct) proof?

Answer: [\[Yes\]](#)

Justification: See section 7.3: Theoretical analyses. This section provide the theoretical analyses and the complete and correct proof.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [\[Yes\]](#)

Justification: See section 7.10, 7.11, we provide the details about the experiments and deployments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [Yes]

Justification: The used benchmarks are opened. And the codes are released at Anonymous GitHub.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.

- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [\[Yes\]](#)

Justification: See section 7.10 and 7.11. All details are provided.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [\[Yes\]](#)

Justification: We calculate the p-value to demonstrate the significant improvement of the experiments. All experiments are obtained with three runs with different random seeds.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g. negative error rates).
- If error bars are reported in tables or plots, The authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments Compute Resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [\[Yes\]](#)

Justification: See section 4.1.1.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code Of Ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines>?

Answer: [Yes]

Justification: We check the NeurIPS Code of Ethics and our paper conform in every aspect with them.

Guidelines:

- The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer No, they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

10. Broader Impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes]

Justification: We demonstrate the practical application of our proposed method in real-world scenarios that directly impact people's lives.

Guidelines:

- The answer NA means that there is no societal impact of the work performed.
- If the authors answer NA or No, they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pretrained language models, image generators, or scraped datasets)?

Answer: [Yes]

Justification: We release the model's weights trained on open benchmarks and protect the model's weights trained on the sensitive data of users.

Guidelines:

- The answer NA means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

Answer: [Yes]

Justification: We have mentioned and cited their papers.

Guidelines:

- The answer NA means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New Assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [Yes]

Justification: We release the codes and models at Anonymous GitHub.

Guidelines:

- The answer NA means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.

- 1387 • At submission time, remember to anonymize your assets (if applicable). You can either
1388 create an anonymized URL or include an anonymized zip file.

1389 **14. Crowdsourcing and Research with Human Subjects**

1390 Question: For crowdsourcing experiments and research with human subjects, does the paper
1391 include the full text of instructions given to participants and screenshots, if applicable, as
1392 well as details about compensation (if any)?

1393 Answer: [NA]

1394 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1395 Guidelines:

- 1396 • The answer NA means that the paper does not involve crowdsourcing nor research with
1397 human subjects.
- 1398 • Including this information in the supplemental material is fine, but if the main contribu-
1399 tion of the paper involves human subjects, then as much detail as possible should be
1400 included in the main paper.
- 1401 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1402 or other labor should be paid at least the minimum wage in the country of the data
1403 collector.

1404 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
1405 **Subjects**

1406 Question: Does the paper describe potential risks incurred by study participants, whether
1407 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1408 approvals (or an equivalent approval/review based on the requirements of your country or
1409 institution) were obtained?

1410 Answer: [NA]

1411 Justification: The paper does not involve crowdsourcing nor research with human subjects.

1412 Guidelines:

- 1413 • The answer NA means that the paper does not involve crowdsourcing nor research with
1414 human subjects.
- 1415 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1416 may be required for any human subjects research. If you obtained IRB approval, you
1417 should clearly state this in the paper.
- 1418 • We recognize that the procedures for this may vary significantly between institutions
1419 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1420 guidelines for their institution.
- 1421 • For initial submissions, do not include any information that would break anonymity (if
1422 applicable), such as the institution conducting the review.