

# An Off-Policy Reinforcement Learning Algorithm Customized for Multi-Task Fusion in Large-Scale Recommender Systems

Peng Liu  
Tencent Inc.  
Beijing, China  
liupengvswww@gmail.com

Cong Xu  
Tencent Inc.  
Beijing, China  
congxcu@tencent.com

Ming Zhao  
Tencent Inc.  
Beijing, China  
marcozhao@tencent.com

Jiawei Zhu  
Tencent Inc.  
Beijing, China  
erickjwzhu@tencent.com

Bin Wang  
Tencent Inc.  
Beijing, China  
hillmwang@tencent.com

Yi Ren  
Tencent Inc.  
Beijing, China  
henrybjren@tencent.com

## ABSTRACT

As the last critical stage of RSs, Multi-Task Fusion (MTF) is responsible for combining multiple scores outputted by Multi-Task Learning (MTL) into a final score to maximize user satisfaction, which determines the ultimate recommendation results. Recently, to optimize long-term user satisfaction within a recommendation session, Reinforcement Learning (RL) is used for MTF in the industry. However, the off-policy RL algorithms used for MTF so far have the following severe problems: 1) to avoid out-of-distribution (OOD) problem, their constraints are overly strict, which seriously damage their performance; 2) they are unaware of the exploration policy used for producing training data and never interact with real environment, so only suboptimal policy can be learned; 3) the traditional exploration policies are inefficient and hurt user experience. To solve the above problems, we propose a novel method named IntegratedRL-MTF customized for MTF in large-scale RSs. IntegratedRL-MTF integrates off-policy RL model with our online exploration policy to relax overstrict and complicated constraints, which significantly improves its performance. We also design an extremely efficient exploration policy, which eliminates low-value exploration space and focuses on exploring potential high-value state-action pairs. Moreover, we adopt progressive training mode to further enhance our model's performance with the help of our exploration policy. We conduct extensive offline and online experiments in the short video channel of Tencent News. The results demonstrate that our model outperforms other models remarkably. IntegratedRL-MTF has been fully deployed in our RS and other large-scale RSs in Tencent, which have achieved significant improvements.

## CCS CONCEPTS

• Information systems → Recommender systems.

## KEYWORDS

Recommender system; Reinforcement learning; Multi-task fusion; Long-term user satisfaction

## ACM Reference format:

Peng Liu, Cong Xu, Ming Zhao, Jiawei Zhu, Bin Wang and Yi Ren. 2024. Insert Your Title Here: Insert Subtitle Here. In *Proceedings of ACM conference (Conf'24)*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/1234567890>

## 1 INTRODUCTION

Recommender Systems (RSs) [1, 2] which provide personalized recommendation service based on user preference are widely used in various platforms such as short video platforms [3, 7, 14], video platforms [4, 5], E-commerce platforms [6, 8-11] and social networks [12, 13], serving billions of users every day. In brief, industrial RSs mainly include three stages: candidate generation, ranking and Multi-Task Fusion (MTF) [4, 15]. During candidate generation, thousands of candidates are selected from millions or even billions of items. Ranking typically uses a Multi-Task Learning model (MTL) [4, 8, 16-18] to estimate the scores of various user behaviors such as click, watching time, fast slide, like and sharing. Finally, a MTF model combines multiple scores outputted by MTL model into a final score to produce the final ranking of candidates [15], which decides the final recommendation results. However, there is little valuable research on MTF.

The goal of MTF is to maximize user satisfaction. User satisfaction is commonly evaluated by the weighted sum of a user's various feedbacks including watching time, valid click, like, sharing and other behaviors within a single recommendation or a recommendation session. A recommendation session is defined as the process from when a user starts accessing RS to leaving, which may include one or more consecutive requests, as shown in

\*corresponding author

†Author Footnote to be captured as Author Note

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

Conference '24, June, 2024, Texas USA

© 2024 Copyright held by the owner/author(s). 978-1-4503-0000-0/18/06...\$15.00

Figure 1. Early works such as Grid Search [36] and Bayesian Optimization [19] try to obtain the optimal fusion weights through parameter searching in candidate set composed of parameter combinations. The main drawbacks of these methods are not only they ignore user preference and generate the same fusion weights for all users but also they are inefficient. Currently, they are rarely used in large-scale RSs. Evolution Strategy (ES) [20-22] takes the features of user preference as input to produce personalized fusion weights for different users. However, the number of ES model parameters must be small due to its learning pattern, which limits its capability [23]. Furthermore, all the methods mentioned above only focus on the reward of current recommendation (instant reward or instant user satisfaction) but ignore long-term rewards.

In RSs such as Tencent News, TikTok and Kwai, current recommendation has an obvious impact on subsequent recommendations, especially within a recommendation session. Therefore, we need to consider both the reward of current recommendation and the rewards of subsequent recommendations within a session together. Recently, a few works [15, 24, 25] use off-policy Reinforcement Learning (RL) [26] to search for the optimal fusion weights to maximize long-term rewards. Compared with the aforementioned methods, RL considers the cumulative rewards within a session and recommends items that not only satisfy a user in current recommendation but also lead to positive long-term engagement. Moreover, RL has the advantages of much more powerful model performance and more sample-efficient than ES [23]. Currently, RL has been used for MTF in RSs in Tencent [15] and many other companies. However, the existing works on RL-MTF have the following serious problems [15, 26-31]: 1) to avoid out-of-distribution (OOD) problem, the constraints of the off-policy RL algorithms used for MTF in RSs are overly strict and complex, which significantly hurts their performance; 2) online exploration and offline model training are two independent processes, the off-policy RL algorithms are unaware of the exploration policy behind training data and no longer interact with real environment, therefore only suboptimal policy can be learned; 3) the existing exploration policies are inefficient and hurt user experience.

To solve the above problems, we propose a novel method named IntegratedRL-MTF, which is tailored for MTF by leveraging the characteristics of RSs. Firstly, IntegratedRL-MTF integrates off-policy RL model with our online exploration policy. When training model offline, the distribution of the training data generated by our exploration policy can be directly obtained. Therefore, the overstrict constraints used to avoid OOD problem can be relaxed, which significantly improves our RL model’s performance. Secondly, we design a simple but extremely efficient exploration policy. Our exploration policy not only accelerates model iteration speed but also prevents excessive negative impact on user experience, which are of great value for commercial companies. Finally, we propose progressive training mode to further improve our RL model’s performance with the help of our efficient exploration policy, which enables target policy swiftly to converge toward the optimal policy through multiple iterations of online exploration and offline model training. We conduct offline

evaluation to compare our model with other models on the same dataset using the new metric defined by us which is simpler and more accurate for RL-MTF evaluation. Furthermore, we conduct online experiments in the short video channel of Tencent News which serves hundreds of millions of users and the result demonstrates that our RL model outperforms other models remarkably. IntegratedRL-MTF has been fully deployed in our RS for almost one year. Moreover, it has also been adopted in other large-scale RSs in Tencent and achieved remarkable improvements.

The major contributions of our work include:

- We investigate the existing works on RL-MTF and point out the main problems of them, including their overstrict constraints significantly affecting their performance, online exploration and offline training are two independent processes and merely suboptimal policy can be learned, the traditional exploration policies are not efficient and have a negative impact on user experience.
- We propose a novel RL algorithm customized for MTF in large-scale RSs. Our solution integrates off-policy RL algorithm with our exploration policy to relax excessively strict constraints, which significantly improves our RL-MTF model’s performance. Furthermore, our solution adopts progressive training mode to learn the optimal policy. By iteratively exploring the environment multiple times with the help of our efficient exploration policy, the learned policy will be refined repeatedly and swiftly converge toward the actual optimal policy of the environment.
- We conduct offline experiments on the datasets of Tencent News to compare the performance of different models using the new metric designed by us. Moreover, we also conduct online A/B testing in the short video channel of Tencent News and the result demonstrates that our model significantly better than other models, improving +4.64% user valid consumption and +1.74% user duration time compared to baseline.

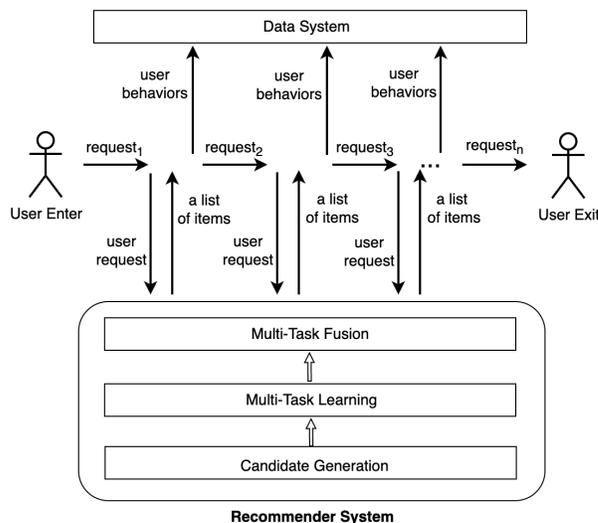


Figure 1: The interactions between user and RS within a recommendation session.

## 2 PROBLEM DEFINITION

This section gives the problem definition of RL-MTF in the short video channel of Tencent News which is similar to TikTok. As mentioned above, current recommendation has an evident influence on subsequent recommendations within a recommendation session. At each time step  $t$ , after RS receives a user request, firstly, thousands of candidates are selected from millions of items; secondly, MTL model predicts scores of multiple user behaviors for each candidate; thirdly, MTF model generates fusion weights to combine multiple scores outputted by MTL model into a final score using Formula 1; finally, a list of items are sent to the user and the user's feedbacks are reported to data system of platform.

$$final\_score = \prod_{i=1}^m (pred\_score_i + bias_i)^{pow_i} \quad (1)$$

We model the above fusion problem as a Markov Decision Process (MDP) within a recommendation session. In this MDP, RS acts as an agent that interacts with a user (environment) and makes sequential recommendations, aiming to maximize the cumulative reward within a session. The MDP framework has the following key components [26]:

- **State Space (S):** is a set of state  $s$  which includes user profile feature (e.g., age, gender, top k interests, flush num, etc.) and user history behavior sequence (e.g., watching, valid click, like, etc.).

- **Action Space (A):** is a set of action  $a$  generated by RL model. In our problem, action  $a$  is a fusion weight vector  $(\alpha_1, \dots, \alpha_k)$ , of which each element corresponds to different pow or bias term in Formula 1.

- **Reward (R):** After RS takes an action  $a_t$  at state  $s_t$  and sends a list of items to a user, the user's various behaviors to those items will be reported to RS and the instant reward  $r(s_t, a_t)$  will be calculated based on these behaviors.

- **Transition Probability (P):** the transition probability  $p(s_{t+1}|s_t, a_t)$  represents the likelihood of transitioning from state  $s_t$  to state  $s_{t+1}$  when action  $a_t$  is taken. In our problem, a state includes user profile feature and user history behavior sequence, so the next state  $s_{t+1}$  depends on user feedback and is deterministic.

- **Discount Factor ( $\gamma$ ):** determines how much weight the agent assigns to future rewards compared to instant reward,  $\gamma \in [0, 1]$ .

With the above definitions, the objective of applying RL for MTF in RS can be defined as follows: Given the interaction history between RS and a user in MDP form within a session, how to learn the optimal policy to maximize the cumulative reward.

## 3 PROPOSED SOLUTION

### 3.1 Reward Function

Within a recommendation session, RS takes an action  $a_t$  at state  $s_t$  to compute final score for each candidate and sends a list of items to a user, then the user's various feedbacks are reported to RS, as

shown in Figure 1. To evaluate instant reward, we define the instant reward function as shown in Formula 2.

$$r(s_t, a_t) = \sum_{i=1}^k w_i v_i \quad (2)$$

where  $w_i$  is the weight of behavior  $v_i$ . In our recommendation scenario, user behaviors  $v_1, \dots, v_k$  contains watching time, valid consumption (watching a video longer than 10 seconds) and interaction behaviors such as liking, sharing, collecting, etc. By analyzing the correlations between different user behaviors and user duration time, we set different weights for these behaviors.

### 3.2 Online Exploration

Before training RL model, a large amount of exploration data needs to be collected first, which has a critical impact on model performance. However, there are two challenges for traditional exploration policies [15, 32]:

- **Low Efficiency:** In practice, it usually requires a long time to collect sufficient exploration data in a large-scale RS using traditional exploration policies. For example, it often takes five or more days to collect exploration data once using action-noise exploration policy on our platform. This affects the speed of model iteration and means a loss of income.

- **Negative Impact on User Experience:** Excessive exploration actions generated by traditional exploration policies including unusual actions have a significantly negative impact on user experience and even lead to user churn, which is unacceptable.

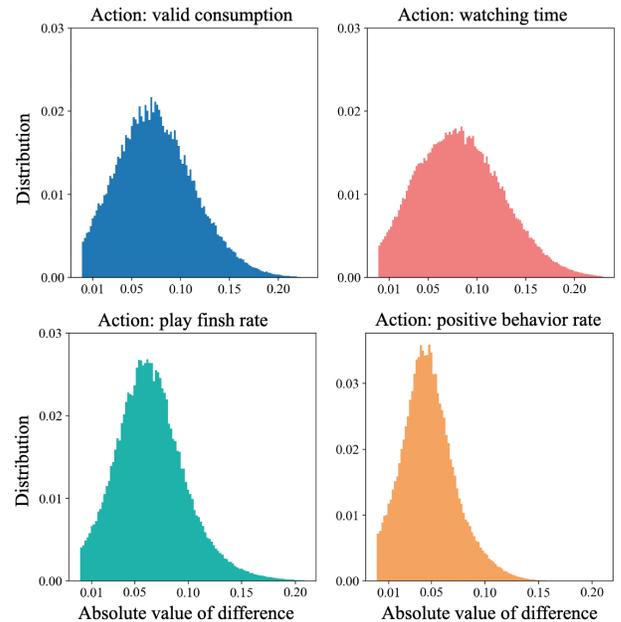
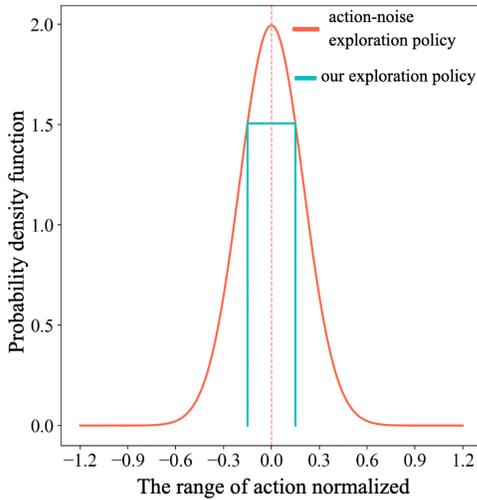


Figure 2: The distribution of absolute difference between the actions generated by new learned RL policy and baseline RL policy for the same state.

To address the above issues, we first perform analysis on the distribution of absolute difference between the actions generated by new learned RL policy and baseline RL policy for the same state on the dataset of our recommendation scenario. For simplicity, we normalize the entire value space of each dimension of the action to the  $[-1, 1]$  range and select the most important 4-dimensional actions including valid consumption, watching time, play finish rate and positive behavior rate for illustration, as shown in Figure 2. We notice that the action generated by new learned RL policy will typically not deviate far from the action generated by baseline RL policy for the same state, which is also consistent with our intuition.

$$\mu_{ep}(s) = \mu_{bp}(s) + \varepsilon, \quad \varepsilon \sim \mathbf{U}(\text{lower}_b, \text{upper}_b) \quad (3)$$



**Figure 3: The action distributions of our exploration policy and action-noise exploration policy.**

Therefore, inspired by this finding, we propose a simple but extremely efficient exploration policy, which defines personalized exploration upper and lower bounds for each user based on baseline policy, as shown in Formula 3. Exploration action is generated by the action outputted by baseline policy plus a random perturbation generated by uniform distribution defined by  $\text{lower}_b$  and  $\text{upper}_b$ . We conduct statistical analysis and select the values of the  $\text{lower}_b$  and  $\text{upper}_b$  elaborately. The basic idea of our exploration policy is to eliminate low-value exploration space and only focus on exploring potential high-value state-action pairs, as shown in Figure 3. By this way, our exploration policy exhibits extremely high efficiency compared to traditional exploration policies. For simplicity, in this paper, we take action-noise exploration policy as an example, which is commonly used to produce exploration data, as shown by the coral curve in Figure 3. In our recommendation scenario, under the same requirement for exploration density, the efficiency of our exploration policy is about  $2^{10}$  times higher than that of action-noise exploration policy, which is analyzed in Section 4. In addition, compared to action-noise exploration policy, our exploration policy can reduce the interference of data dis-

tribution on RL-MTF model training. Moreover, the progressive training mode detailed in Section 3 further expands the exploration space of our exploration policy and the upper and lower bounds of personalized exploration space can be set smaller.

### 3.3 IntegratedRL-MTF: A RL Algorithm Customized for MTF in Large-scale RSs

To solve the problems mentioned previously, we propose a novel method named IntegratedRL-MTF. Next, we will introduce actor network, critic network and progressive training mode of IntegratedRL-MTF.

**3.3.1 Actor Network.** The Actor network aims to output the optimal action for a specific state. Following the common setting, we build two actor networks during the learning process, including one current actor network  $\mu(s)$  and one target actor network  $\mu'(s)$ .  $\mu(s)$  integrates actor network with our exploration policy to relax too strict constraints and introduce an additional penalty term based on the consistency of multiple critics to alleviate extrapolation error, as shown in Formula 4 and 5.

$$\begin{aligned} \theta_{k+1} \leftarrow \arg \min_{\theta_k} E_{s_t \sim \mathcal{D}, a \sim \mu(s_t | \theta_k)} [ & \\ -1/m * \sum_{i=1}^m Q_i(s_t, \mu(s_t | \theta_k)) & \\ + \eta * d(\mu(s_t | \theta_k), \mu_{bp}(s_t), \text{lower}_b, \text{upper}_b) & \\ + \lambda * \sqrt{1/m * \sum_{i=1}^m (Q_i - \text{mean}(Q))^2} & ] \end{aligned} \quad (4)$$

$$d(\mu(s_t | \theta_k), \mu_{bp}(s_t), \text{lower}_b, \text{upper}_b) = \begin{cases} 0, & \mu(s_t | \theta_k) \in (\mu_{bp}(s_t) + \text{lower}_b, \mu_{bp}(s_t) + \text{upper}_b) \\ e^{-\frac{\mu(s_t | \theta_k) - (\mu_{bp}(s_t) + \text{upper}_b)}{\beta * (\text{upper}_b - \text{lower}_b)}}, & \mu(s_t | \theta_k) > \mu_{bp}(s_t) + \text{upper}_b \\ e^{-\frac{(\mu_{bp}(s_t) + \text{lower}_b) - \mu(s_t | \theta_k)}{\beta * (\text{upper}_b - \text{lower}_b)}}, & \mu(s_t | \theta_k) < \mu_{bp}(s_t) + \text{lower}_b \end{cases} \quad (5)$$

During training  $\mu(s)$ , the upper and lower bounds of the exploration data distribution for each user can be directly obtained, as mentioned in Section 3.2. Therefore, we can utilize this characteristic to simplify the overly strict constraints and fully exploit the capacity of  $\mu(s)$ . If the action generated by  $\mu(s)$  at state  $s_t$  is within the user's upper and lower bounds, the value of the second term in Formula 4 is zero, namely no penalty is imposed to avoid affecting the capacity of the model. Otherwise, a penalty will be applied based on the deviation that exceeds either the upper or lower bounds of the user. By this way, the performance of the current actor network is significantly improved compared to the existing methods, which is demonstrated by the experiments in Section 4. Furthermore, we also introduce a penalty mechanism which is defined as the standard deviation of the estimated values outputted by multiple independent critics [33] to mitigate extrapolation er-

ror, which is the third term in Formula 4. Due to extremely high efficiency of our exploration policy, the exploration actions collected within the user's upper and lower bounds have a significantly higher average density compared to traditional action-noise exploration policy, which is highly valuable for model optimization. Moreover, compared with Gaussian perturbation, Random perturbation within personalized upper and lower bounds mitigates the interference of data distribution on model training. If the action outputted by  $\mu(s)$  is within the exploration space of the user, the value of the third term in Formula 4 will be small and can even be ignored. Otherwise, there will be a corresponding penalty used to mitigate extrapolation error.

The target actor network  $\mu'(s)$  is an auxiliary network responsible for generating the next optimal action based on next state to alleviate the overestimation problem caused by bootstrapping. Its parameters are periodically soft updated using the current actor network.

**3.3.2 Critic Network.** The Critic network  $Q(s, a)$  is responsible for estimating the cumulative reward of a state-action pair  $(s, a)$  within a recommendation session.  $Q(s, a)$  also integrates critic network with our exploration policy to avoid extrapolation error. In our solution, multiple independent critic networks are created, which are initialized randomly and trained independently. The goal of each critic network is to minimize TD-error, as shown in Formula 6. If the next action generated by  $\mu'(s)$  at next state  $s_{t+1}$  is within the user's upper and lower bounds, the value of the second term in Formula 6 is zero. Otherwise, a penalty will be applied based on the deviation that surpasses the prescribed upper or lower bounds of the user. In practice, we often set the number of critic networks to 24, which is sufficient to achieve good result in our recommendation scenario. To achieve better performance, we define a target network for each critic, of which the parameters are soft updated periodically using the corresponding critic network.

$$\phi_{k+1} \leftarrow \underset{\phi_k}{\operatorname{arg\,min}} E_{s_t \sim \mathcal{D}, a_t \sim \mathcal{D}, a_{t+1} \sim \mu'(s_{t+1})} [ (Q_i(s_t, a_t | \phi_k) - (r_t + \gamma * Q_i(s_{t+1}, \mu'(s_{t+1}))))^2 + \delta * d(\mu'(s_{t+1}), \mu_{bp}(s_{t+1}), lower_b, upper_b) ] \quad (6)$$

for  $i = 1, 2, \dots, m$

**3.3.3 Progressive Training Mode.** One severe drawback of off-policy RL is that when the model is trained offline, it relies solely on the data collected before without further interaction with real environment. The absence of real-time interaction during offline training can lead to the discrepancy between the learned policy and the actual environment, which has a significant negative impact on the performance of off-policy RL algorithms [15, 26-31].

To alleviate this problem in large-scale RSs, our solution adopts progressive training mode to learn the optimal policy through multiple iterations of online exploration and offline model training with the help of our efficient exploration policy, which enables target policy to converge rapidly towards the optimal policy. Due to high efficiency of our exploration policy, we divide the previous single data exploration and offline model training into five

rounds of online data exploration and offline model training. The latest learned policy is used as the baseline policy of the next online exploration. By iteratively and efficiently exploring the environment, the learned policy will be improved repeatedly, which further enhances the performance of our RL model.

### 3.4 Recommender System with RL-MTF

We implement IntegratedRL-MTF in the short video channel of Tencent News, as shown in Figure 4. Our RL-MTF framework is composed of two components: offline model training and online model serving. The offline model training component is in charge of preprocessing exploration data and training RL-MTF model. The online model serving component is mainly responsible for generating the personalized optimal action when receiving a user request to calculate the final score for each candidate. In addition, the online model serving component also takes charge of online exploration to collect training data.

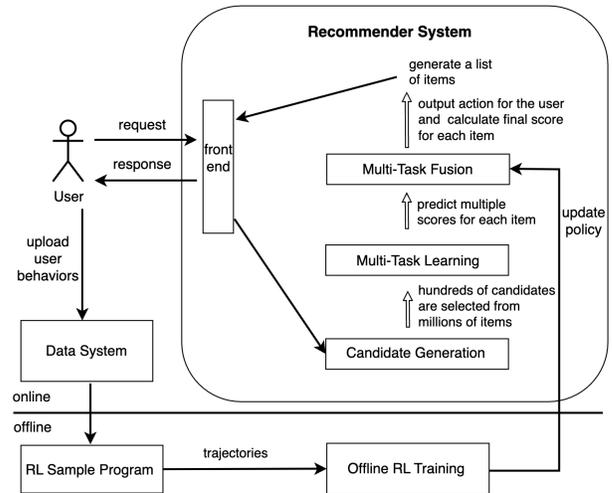


Figure 4: RL-MTF framework in our recommender system.

## 4 EXPERIMENTS

### 4.1 Dataset

The following datasets are collected from the short video channel of Tencent News, which serves hundreds of millions of users. We use four groups of users to collect exploration data. The users of each group are selected randomly and the number of them in each group is the same (about 2 million) for a fair comparison. The online exploration policy of each dataset is defined as follows:

- **Dataset 1:** It is generated by action-noise exploration policy, which adds Gaussian noise to the action outputted by baseline policy, as shown by the coral curve in Figure 3. In our experiment, the Gaussian distribution has a mean of 0.0 and a standard deviation of 0.2, which were used previously in our RS.

- **Dataset 2:** It is produced by our exploration policy, as shown in Formula 3. As mentioned before, we set the value of upper bound for a user as the action outputted by baseline policy plus

0.15 and set the value of lower bound for a user as the baseline action minus 0.15, as shown by the cyan curve in Figure 3.

- **Dataset 3:** It is also generated by our exploration policy. The difference is this dataset is collected through five rounds of online exploration. After completing model training with the exploration data collected in current round, the newly learned RL policy will serve as the baseline policy to start next round of online exploration. The duration for data collection of each round is one day.

- **Test Dataset:** The dataset generated by baseline policy, which serves as the test data for different models.

All the exploration policies are used to explore the environment for five days. And each dataset includes about 6.8 million sessions. For offline experiments, we train different RL models with different datasets and evaluate these models on the same test dataset using the novel metric defined by us. For online experiments, we deploy different models in the short video channel of Tencent News for one week to conduct A/B tests.

## 4.2 Implementation Details

In RL-MTF, user state  $s$  includes the user's profile features (e.g., age, gender, top  $k$  interests, flush num, etc.) and user history behavior sequences, of which the maximum length is limited to 100. The action generated by MTF model is a 10-dimensional vector representing the fusion weights in Formula 1. All networks of RL-MTF models are Multi-Layer Perceptron (MLP) and are optimized based on Adam optimizer. The reward discount factor  $\gamma$  is set to 0.9. The values of  $\eta$ ,  $\beta$  and  $\lambda$  are set to 1.2, 0.3 and 0.2 separately. The number of critic networks of our RL-MTF model is set to 24. The soft update rate and the delay update step for target networks are  $w = 0.08$  and  $L = 15$ . In addition, the mini-batch size and training epochs are set to 256 and 300, 000 respectively.

## 4.3 Evaluation Setting

**4.3.1 Offline policy evaluation.** Offline evaluation has the advantages of low cost and no negative impact on users. The Group Area Under the Curve (GAUC) [34] is an evaluation metric commonly used in RSs to measure the quality of the ranking produced by ranking algorithms [4, 8, 16-18]. Since the goal of a MTF model is to provide a ranking of candidates based on final score to maximize reward, similar evaluation metrics can be used as well. Building on the GAUC metric, we design a new metric for offline RL-MTF model evaluation. Firstly, we define the label of each user-item sample according to the user's watching time. If the user watches an item more than 10 seconds (a valid consumption), the label of the user-item sample is 1; otherwise, its label is 0. Secondly, we define the reward of each user-item sample according to Formula 2. And the weight of each user-item sample is its reward plus a constant value which is typically set to 1. Thirdly, we use RL policy to generate the action for each user-item sample and calculate its final score. Then we normalize all the final scores in test dataset to the  $[0, 1]$  range, which are used as predict scores. Finally, we calculate weighted GAUC metric to evaluate the performance of different RL-MTF models. By considering the re-

ward of each user-item sample, the weighted GAUC metric provides a more comprehensive measure of a RL policy's ability to final rank and discriminate between different rewards. Using this metric, we can easily measure the quality of the final ranking generated by each MTF model on the same dataset. Compared to the existing evaluation metrics [15, 26], our evaluation metric is simpler, more accurate and works well in practice, which is more appropriate for RL-MTF evaluation in RSs.

**4.3.2 Online A/B testing.** We adopt user valid consumption and user duration time to evaluate each MTF model, which are the two most important online metrics in our recommendation scenario.

- **User valid consumption** is the average of all users' total valid consumptions during a day. A valid consumption is defined as a user watching a video more than 10 seconds.

- **User duration time** is the average of all users' total watching time within a day.

## 4.4 Compared Methods

We compare IntegratedRL-MTF with ES and other advanced off-policy RL algorithms. In addition, we also design two variants of IntegratedRL-MTF to illustrate the effects of integrating off-policy RL model with our exploration policy and progressive training mode on model performance respectively.

- **ES** [20-23] takes user profile features as its model input to generate personalized fusion weights. Since ES is simple and works well, it is used as the benchmark in this paper.

- **BatchRL-MTF** [15] is proposed for MTF in RSs and generates actions based on BCQ [28], which has been deployed in multiple RSs in Tencent for several years and achieved excellent improvements. The other compared RL-MTF methods are implemented based on the framework proposed by BatchRL-MTF.

- **DDPG** (Deep Deterministic Policy Gradient) [27, 35] is a classical off-policy actor-critic algorithm that can learn policies in high-dimensional, continuous action spaces.

- **CQL+SAC** (Conservative Q-Learning with Soft Actor-Critic) [30, 39] learns a conservative, lower-bound Q function by regularizing the Q value of OOD action-state pairs to reduce extrapolation error.

- **IQL** (Implicit Q-learning) [31] never needs to evaluate actions outside of the dataset, but still enables the learned policy to improve substantially over the best behavior in the data through generalization.

- **IntegratedRL-MTF without PTM** (Progressive Training Mode is abbreviated as PTM for simplicity) integrates off-policy RL model with our online exploration policy to relax overly strict constraints, which significantly improves its model performance. To validate its effectiveness, we train the model of IntegratedRL-MTF without PTM using dataset 2.

- **IntegratedRL-MTF** integrates off-policy RL model with online exploration policy and learns the optimal policy through

multiple iterations of online exploration and offline training with the help of our efficient exploration policy using dataset 3.

## 4.5 Offline Evaluation

In this section, extensive offline experiments are conducted to demonstrate the remarkable performance of IntegratedRL-MTF compared to other MTF algorithms. In addition, we also analyze the efficiency of our online exploration policy.

*4.5.1 Effectiveness of IntegratedRL-MTF.* To compare the performance of the above algorithms, we train their models separately and evaluate their performance on the same test dataset using the weighted GAUC metric. ES updates its model parameters through a process of alternating mutation and selection and had been fully deployed in our RS for a long time, which is the benchmark of offline evaluation. DDPG, CQL+SAC, BatchRL-MTF and IQL are trained by dataset 1. IntegratedRL-MTF without PTM is trained by dataset 2. IntegratedRL-MTF is trained by dataset 3 using progressive training mode. All the models are tested on the same dataset and the results are shown in Table 1.

**Table 1: The weighted GAUC of the compared methods on the same test dataset.**

Compared Methods	Weighted GAUC
ES	0.7836
DDPG	0.7881
CQL+SAC	0.7892
BatchRL-MTF	0.7894
IQL	0.7906
<b>IntegratedRL-MTF without PTM</b>	<b>0.7941</b>
<b>IntegratedRL-MTF</b>	<b>0.7953</b>

In offline evaluation, the weighted GAUC of IntegratedRL-MTF without PTM is significantly higher than that of ES model and other existing off-policy RL models. As previously mentioned, to avoid OOD problem, the constraints of the existing off-policy RL algorithms are overly strict and intricate, which significantly damage their performance. This is because traditional off-policy RL algorithms are limited to training with a fixed dataset but are unaware of the exploration policy behind the dataset. Therefore, those off-policy RL algorithms can only avoid OOD problem through strong constraints. In RSs, we design a novel online exploration policy and during offline model training, the upper and lower bounds of the exploration data distribution for each user can be directly obtained. We utilize this characteristic to simplify the overly strict constraints and integrate off-policy model algorithm with our efficient exploration policy. By this way, IntegratedRL-MTF without PTM significantly improves its model performance and that is also proved in offline evaluation. Furthermore, IntegratedRL-MTF is superior to IntegratedRL-MTF without PTM. By iteratively exploring the environment multiple times with the

help of our efficient exploration policy, the learned policy will be enhanced repeatedly.

*4.5.2 Efficiency of Our Exploration policy.* The action of our online exploration policy is generated by the action outputted by baseline policy plus a random perturbation, which must be within the user's personalized upper and lower bounds, as shown by the cyan curve in Figure 3. The idea of this approach is to exclude low-value exploration space and merely focus on exploring potential high-value state-action pairs. We conduct statistical analysis on the distribution of the actions generated by different RL policies for the same state and carefully select the values of  $upper_b$  and  $lower_b$ . For simplicity, we assume the exploration range of each dimension of our exploration policy is the same value in this paper. Typically, we set the values of  $upper_b$  and  $lower_b$  as 0.15 and -0.15 respectively, which is sufficient in our recommendation scenario. The action-noise exploration policy previously used in our recommendation scenario has a mean of 0.0 and a standard deviation of 0.2, as shown by the coral curve in Figure 3. The action generated by RL-MTF model is a 10-dimensional vector. Therefore, in our recommendation scenario, under the same requirement for exploration density, the efficiency of our exploration policy is about  $2^{10}$  times higher than that of the action-noise exploration policy. In practice, because IntegratedRL-MTF adopts progressive training mode, the upper and lower bounds of our exploration policy can be set smaller.

## 4.6 Online Evaluation

In this section, we deploy all the compared models in the short video channel of Tencent News for one week to conduct A/B tests. We take ES as the benchmark for comparison and show the improvements of user valid consumption and user duration time for other models. The online experiment results are shown in Table 2 and all the improvements have statistical significance with p-value less than 0.05.

**Table 2: The online results of the compared methods in the short video channel of Tencent News.**

Compared Methods	Valid Consumption	Duration Time
ES	*	*
DDPG	+1.39%	+0.81%
CQL+SAC	+1.62%	+0.95%
BatchRL-MTF	+1.79%	+0.98%
IQL	+2.09%	+1.15%
<b>IntegratedRL-MTF</b>	<b>+4.64%</b>	<b>+1.74%</b>

DDPG increases +1.39% user valid consumption and +0.81% user duration time by considering long-term reward and much stronger model performance compared to ES. CQL+SAC is superior to DDPG and BatchRL-MTF slightly outperforms CQL+SAC. IQL outperforms DDPG, CQL+SAC and BatchRL-MTF, increasing +2.09% user valid consumption and +1.15% user

duration time compared to ES. IntegratedRL-MTF significantly outperforms the other algorithms, increasing +4.64% user valid consumption and +1.74% user duration time compared to benchmark. In addition, the online evaluation results are consistent with the offline evaluation results given by weighted GAUC metric defined by us in Section 4.3.1.

## 5 RELATED WORK

MTF [4, 15] is responsible for merging multiple scores generated by MTL [4, 8, 16-18] into a final score to produce the final ranking of candidates, which is critical for final recommendation results. However, there is little valuable research on MTF so far.

Initially, Grid Search [36] is used to search the optimal fusion weights via parameter searching in the candidate set which is composed of lots of parameter combinations. Then Bayesian Optimization [19, 37] is proposed to accelerate the speed of parameter search. The primary drawback of these two methods is that they ignore user preference and merely produce the same fusion weights for all users. In addition, these methods are also not efficient. Therefore, they are seldom used in industrial RSs now. ES [20-22] takes the features of user preference as model input to output personalized fusion weights for different users, which is a simple and common MTF method. ES updates its model parameters through a process of alternating mutation and selection. The main drawback of ES is that due to its parameter update pattern, the number of its model parameters must be small, which significantly limits its capability [23]. Furthermore, all the methods mentioned above only focus on instant reward but ignore long-term reward.

In personalized RSs such as Tencent News, TikTok, Kwai and Little Red Book [3, 7, 14], current recommendation has an obvious influence on subsequent recommendations, especially within a recommendation session. Therefore, in recent years, some works try to search the optimal fusion weights via RL to maximize long-term reward. For example, [38] uses off-policy RL algorithm to search the optimal weights between the predicted click-through rate and bid price of an advertiser. [15] proposes BatchRL-MTF to find out the optimal fusion weights, which has been fully deployed in multiple RSs in Tencent for several years and achieved excellent improvements. And some teams in the industry have also attempted alternative off-policy RL algorithms such as DDPG [27], CQL+SAC [30, 39], IQL [31], etc. However, the existing works on RL-MTF so far have the following serious problems. Firstly, to avoid OOD problem, the constraints of the existing off-policy RL algorithms are too strict and complicated, which significantly hurt their performance. For example, in BCQ, the action generated by target policy for a given state must closely resemble the distribution of action for this state contained in the exploration data. How to relax overstrict constraints but avoid OOD problem in off-policy RL is a hot research area. Secondly, online exploration and offline model training are two independent processes. During offline training, the existing off-policy RL algorithms are unaware of the exploration policy used for generating training data and no longer interact with environment, so only suboptimal

policy can be learned. Thirdly, the efficiency of traditional online exploration policies [15, 32] are low efficiency and significantly hurt user experience.

## 6 CONCLUSION

In this paper, we first point out the problems of the existing RL-MTF methods and propose a novel RL-MTF solution by utilizing the characteristics of RSs. Our solution integrates off-policy RL algorithm with our online exploration policy, which significantly improves its performance. Furthermore, we use progressive training mode to learn the optimal policy with the help of our efficient exploration policy, which further improves the performance of our RL model. In addition, we analyze the efficiency of our exploration policy, which not only accelerates model iteration speed but also prevents excessive negative impact on user experience, which are highly valuable for industrial RSs. We also design a new metric building on GAUC for offline RL-MTF model evaluation. Offline experiments show that our model significantly exceeds other models. Finally, we also conduct online A/B testing in the short video channel of Tencent News and the result demonstrates our model performs significantly better than other models, resulting in a +4.64% improvement in user valid consumption and a +1.74% improvement in user duration time compared to the baseline. Our RL-MTF model has been fully deployed in the short video channel of Tencent News for about one year. Moreover, our solution has also been used in other RSs in Tencent and achieved remarkable improvements. Applying RL to optimize long-term reward in large-scale online RSs requires careful consideration of details and is a challenging task. We have conducted extensive work on RL-MTF and will release other relevant works in the future. We welcome communication with professionals from both industry and academia.

## REFERENCES

- [1] Yang Li, Kangbo Liu, Ranjan Satapathy, Suhang Wang, and Erik Cambria. 2024. Recent Developments in Recommender Systems: A Survey [Review Article]. *IEEE Comput. Intell. Mag.* 19, 2 (May 2024), 78–95. <https://doi.org/10.1109/MCI.2024.3363984>.
- [2] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2020. Deep Learning Based Recommender System: A Survey and New Perspectives. *ACM Comput. Surv.* 52, 1 (January 2020), 1–38. <https://doi.org/10.1145/3285029>.
- [3] Jianxin Chang, Chenbin Zhang, Yiqun Hui, Dewei Leng, Yanan Niu, Yang Song, and Kun Gai. 2023. PEPNet: Parameter and Embedding Personalized Network for Infusing with Personalized Prior Information. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, August 06, 2023. ACM, Long Beach CA USA, 3795–3804. <https://doi.org/10.1145/3580305.3599884>.
- [4] Paul Covington, Jay Adams, and Emre Sargin. 2016. Deep Neural Networks for YouTube Recommendations. In *Proceedings of the 10th ACM Conference on Recommender Systems*, September 07, 2016. ACM, Boston Massachusetts USA, 191–198. <https://doi.org/10.1145/2959100.2959190>.
- [5] Carlos A. Gomez-Urbe and Neil Hunt. 2016. The Netflix Recommender System: Algorithms, Business Value, and Innovation. *ACM Trans. Manage. Inf. Syst.* 6, 4 (January 2016), 1–19. <https://doi.org/10.1145/2843948>.
- [6] Greg Linden, Brent Smith, and Jeremy York. 2003. Amazon.com recommendations: Item-to-item collaborative filtering. *IEEE Internet computing* 7, 1 (2003), 76–80.
- [7] Zhuoran Liu, Leqi Zou, Xuan Zou, Caihua Wang, Biao Zhang, Da Tang, Bolin Zhu, Yijie Zhu, Peng Wu, Ke Wang, and Youlong Cheng. 2022. Monolith: Real Time Recommendation System With Collisionless Embedding Table. Retrieved April 11, 2024 from <http://arxiv.org/abs/2209.07663>.
- [8] Guorui Zhou, Xiaoqiang Zhu, Chenru Song, Ying Fan, Han Zhu, Xiao Ma, Yanghui Yan, Junqi Jin, Han Li, and Kun Gai. 2018. Deep Interest Network for Click-Through Rate Prediction. In *Proceedings of the 24th ACM SIGKDD In-*

- ternational Conference on Knowledge Discovery & Data Mining, July 19, 2018. ACM, London United Kingdom, 1059–1068. <https://doi.org/10.1145/3219819.3219823>.
- [9] Yulong Gu, Zhuoye Ding, Shuaiqiang Wang, and Dawei Yin. 2020. Hierarchical User Profiling for E-commerce Recommender Systems. In WSDM. 223–231.
- [10] Yulong Gu, Zhuoye Ding, Shuaiqiang Wang, and Dawei Yin. 2020. Hierarchical User Profiling for E-commerce Recommender Systems. In Proceedings of the 13th International Conference on Web Search and Data Mining, January 20, 2020. ACM, Houston TX USA, 223–231. <https://doi.org/10.1145/3336191.3371827>.
- [11] Lixin Zou, Long Xia, Zhuoye Ding, Jiaying Song, Weidong Liu, and Dawei Yin. 2019. Reinforcement learning to optimize long-term user engagement in recommender systems. In SIGKDD. 2810–2818.
- [12] Yulong Gu, Jiaying Song, Weidong Liu, and Lixin Zou. 2016. HLGPS: a home location global positioning system in location-based social networks. In ICDM. IEEE, 901–906.
- [13] XinranHe, JunfengPan, OuJin, TianbingXu, BoLiu, TaoXu, YanxinShi, Antoine Atallah, Ralf Herbrich, Stuart Bowers, et al. 2014. Practical lessons from predicting clicks on ads at facebook. In ADKDD. 1–9.
- [14] Yanhua Huang, Weikun Wang, Lei Zhang, and Ruiwen Xu. 2021. Sliding Spectrum Decomposition for Diversified Recommendation. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, August 14, 2021. ACM, Virtual Event Singapore, 3041–3049. <https://doi.org/10.1145/3447548.3467108>.
- [15] Qihua Zhang, Junning Liu, Yuzhuo Dai, Yiyang Qi, Yifan Yuan, Kunlun Zheng, Fan Huang, and Xianfeng Tan. 2022. Multi-Task Fusion via Reinforcement Learning for Long-Term User Satisfaction in Recommender Systems. In Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, August 14, 2022. ACM, Washington DC USA, 4510–4520. <https://doi.org/10.1145/3534678.3539040>.
- [16] Jiaqi Ma, Zhe Zhao, Xinyang Yi, Jilin Chen, Lichan Hong, and Ed H. Chi. 2018. 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*, July 19, 2018. ACM, London United Kingdom, 1930–1939. <https://doi.org/10.1145/3219819.3220007>.
- [17] Zhe Zhao, Lichan Hong, Li Wei, Jilin Chen, Aniruddh Nath, Shawn Andrews, Aditee Kumthekar, Maheswaran Sathiamoorthy, Xinyang Yi, and Ed Chi. 2019. Recommending what video to watch next: a multitask ranking system. In Proceedings of the 13th ACM Conference on Recommender Systems, September 10, 2019. ACM, Copenhagen Denmark, 43–51. <https://doi.org/10.1145/3298689.3346997>.
- [18] Hongyan Tang, Junning Liu, Ming Zhao, and Xudong Gong. 2020. Progressive Layered Extraction (PLE): A Novel Multi-Task Learning (MTL) Model for Personalized Recommendations. In Fourteenth ACM Conference on Recommender Systems, September 22, 2020. ACM, Virtual Event Brazil, 269–278. <https://doi.org/10.1145/3383313.3412236>.
- [19] Jonas Moćkus. 1975. On Bayesian methods for seeking the extremum. In Optimization techniques IFIP technical conference. Springer, 400–404.
- [20] Nikolaus Hansen. 2016. The CMA Evolution Strategy: A Tutorial. (2016). <https://doi.org/10.48550/ARXIV.1604.00772>.
- [21] Hans-Georg Beyer and Hans-Paul Schwefel. 2002. Evolution strategies a comprehensive introduction. *Natural computing* 1, 1 (2002), 3–52.
- [22] Daan Wierstra, Tom Schaul, Jan Peters, and Juergen Schmidhuber. 2008. Natural Evolution Strategies. In 2008 IEEE Congress on Evolutionary Computation (IEEE World Congress on Computational Intelligence), June 2008. IEEE, Hong Kong, China, 3381–3387. <https://doi.org/10.1109/CEC.2008.4631255>.
- [23] Amjad Yousef Majid, Serge Saaybi, Vincent Francois-Lavet, R. Venkatesh Prasad, and Chris Verhoeven. 2024. Deep Reinforcement Learning Versus Evolution Strategies: A Comparative Survey. *IEEE Trans. Neural Netw. Learning Syst.* (2024), 1–19. <https://doi.org/10.1109/TNNLS.2023.3264540>.
- [24] ChanghuaPei, XinruYang, QingCui, XiaoLin, FeiSun, PengJiang, WenwuOu, and Yongfeng Zhang. 2019. Value-aware recommendation based on reinforced profit maximization in e-commerce systems. *arXiv preprint arXiv:1902.00851* (2019).
- [25] Jianhua Han, Yong Yu, Feng Liu, Ruiming Tang, and Yuzhou Zhang. 2019. Optimizing Ranking Algorithm in Recommender System via Deep Reinforcement Learning. In AIAM. IEEE, 22–26.
- [26] Rafael Figueiredo Prudencio, Marcos R. O. A. Maximo, and Esther Luna Colombini. 2024. A Survey on Offline Reinforcement Learning: Taxonomy, Review, and Open Problems. *IEEE Trans. Neural Netw. Learning Syst.* (2024), 1–0. <https://doi.org/10.1109/TNNLS.2023.3250269>.
- [27] Timothy P. Lillicrap, Jonathan J. Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In Yoshua Bengio and Yann LeCun, editors, 4th International Conference on Learning Representations, ICLR 2016, San Juan, Puerto Rico, May 2–4, 2016, Conference Track Proceedings, 2016.
- [28] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-policy deep reinforcement learning without exploration. In ICML. PMLR, 2052–2062.
- [29] Aviral Kumar, Justin Fu, George Tucker, and Sergey Levine. 2019. Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction. (2019). <https://doi.org/10.48550/ARXIV.1906.00949>.
- [30] Aviral Kumar, Aurick Zhou, George Tucker, Sergey Levine. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems*, 33:1179–1191, 2020.
- [31] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. 2021. Offline Reinforcement Learning with Implicit Q-Learning. (2021). <https://doi.org/10.48550/ARXIV.2110.06169>.
- [32] Pawel Ladosz, Lilian Weng, Minwoo Kim, and Hyondong Oh. 2022. Exploration in Deep Reinforcement Learning: A Survey. (2022). <https://doi.org/10.48550/ARXIV.2205.00824>.
- [33] Kamyar Ghasemipour, Shixiang Shane Gu, and Ofir Nachum. Why so pessimistic? estimating uncertainties for offline rl through ensembles, and why their independence matters. *Advances in Neural Information Processing Systems*, 35:18267–18281, 2022.
- [34] Han Zhu, Junqi Jin, Chang Tan, Fei Pan, Yifan Zeng, Han Li, and Kun Gai. 2017. Optimized cost per click in taobao display advertising. In SIGKDD. 2191–2200.
- [35] Silver David, Lever Guy, Heess Nicolas, Degris Thomas, Wierstra Daan, and Riedmiller Martin. 2014. Deterministic policy gradient algorithms. In International Conference on Machine Learning. PMLR, 387–395.
- [36] Petro Liashchynskyi and Pavlo Liashchynskyi. 2019. Grid Search, Random Search, Genetic Algorithm: A Big Comparison for NAS. (2019). <https://doi.org/10.48550/ARXIV.1912.06059>.
- [37] Bruno G Galuzzi, Ilaria Giordani, Antonio Candelieri, Riccardo Perego, and Francesco Archetti. 2020. Hyperparameter optimization for recommender systems through Bayesian optimization. *CMS* 17, 4 (2020), 495–515.
- [38] Jianhua Han, Yong Yu, Feng Liu, Ruiming Tang, and Yuzhou Zhang. 2019. Optimizing Ranking Algorithm in Recommender System via Deep Reinforcement.
- [39] Tuomas Haarnoja, Aurick Zhou, Pieter Abbeel, and Sergey Levine. 2018. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In International conference on machine learning. PMLR, 1861–1870.