# Multi-task Offline Reinforcement Learning for Online Advertising in Recommender Systems

Langming Liu*
Taobao & Tmall Group of Alibaba
City University of Hong Kong
Hangzhou, China

Wanyu Wang*
Southern University of Science and
Technology, City University of Hong
Kong
Hong Kong, China

Chi Zhang
City University of Hong Kong
Hong Kong, China

Bo Li
City University of Hong Kong
Hong Kong, China

Hongzhi Yin
University of Queensland
Brisbane, Australia

Xuetao Wei†
Southern University of Science and
Technology
Shenzhen, China

Wenbo Su
Taobao & Tmall Group of Alibaba
Beijing, China

Bo Zheng
Taobao & Tmall Group of Alibaba
Beijing, China

Xiangyu Zhao†
City University of Hong Kong
Hong Kong, China

## Abstract

Online advertising in recommendation platforms has gained significant attention, with a predominant focus on channel recommendation and budget allocation strategies. However, current offline reinforcement learning (RL) methods face substantial challenges when applied to sparse advertising scenarios, primarily due to severe overestimation, distributional shifts, and overlooking budget constraints. To address these issues, we propose MTORL, a novel multi-task offline RL model that targets two key objectives. First, we establish a Markov Decision Process (MDP) framework specific to the nuances of advertising. Then, we develop a causal state encoder to capture dynamic user interests and temporal dependencies, facilitating offline RL through conditional sequence modeling. Causal attention mechanisms are introduced to enhance user sequence representations by identifying correlations among causal states. We employ multi-task learning to decode actions and rewards, simultaneously addressing channel recommendation and budget allocation. Notably, our framework includes an automated system for integrating these tasks into online advertising. Extensive experiments on offline and online environments demonstrate MTORL's superiority over state-of-the-art methods. The code is available online at https://github.com/Applied-Machine-Learning-Lab/MTORL.

## CCS Concepts

• **Information systems → Recommender systems**.

*Both authors contributed equally to this research.
†Corresponding Authors

## Keywords

Advertising, Offline Reinforcement Learning, Multi-task Learning

## 1 Introduction

In recent years, advertising has emerged as a prominent approach for promoting products and services across various domains, encompassing e-commerce, short-video platforms, social media, and insurance [8, 24, 41, 58, 93]. Unlike traditional advertising, which confines promotion to fixed channels, online advertising systems offer dynamic and personalized ad placement policies [19, 21] for online advertisers. At the core of online advertising lies channel recommendation and budget allocation. **Channel recommendation** aims to autonomously identify suitable advertising channels (e.g., search ads, display ads, short videos, social media) for each exposure based on user preferences [75, 78]. Conversely, judicious **budget allocation** at the channel and other levels becomes imperative to maximize revenue [5]. While prior works have strived to elaborate corresponding advertising strategies [11, 20, 57, 78], the proliferation of users and user features presents a formidable challenge for traditional approaches in accurately capturing user preferences, jeopardizing advertising performance.

Deep learning (DL) has revolutionized the advertising landscape [17, 46, 49, 54, 70, 81, 82, 85–87, 92], especially **channel recommendation**, leveraging deep networks to extract features from intra- and inter-channel across various ads and provide recommendation results through suitable channels [26, 52, 73]. However, most DL-based methods lack dynamic decision-making capabilities and ignore long-term rewards [3, 34]. Deep reinforcement learning (DRL) [3, 71] provides the solution to address the issues above, as

its nature is to solve the MDP problem, making dynamic decisions to maximize long-term revenue [55, 77, 80, 83, 84, 88–90]. However, directly deploying untrained DRL into the online-serving advertising system is destructive for revenue and user experience [80]. Therefore, pre-training of DRL on the offline dataset is necessary for painless deployment, where the paradigm of offline RL [40, 42, 43] perfectly adapts to this environment. Unfortunately, the issues of overestimation and distributional shifts frequently emerge in offline environments. Inspired by existing advancements in the offline RL community [4, 42], some works strive to address such issues in advertising scenarios [38, 39, 51] by constraining the policy to adhere to the data distribution. However, this approach proves unsuitable due to the heightened sparsity of rewards in advertising datasets compared to general RL scenarios [16, 65, 72]. Another line of offline RL methods [31, 32] treats states as features and actions as labels, lifting the training efficiency by applying a sequential modeling paradigm, e.g., recent prevailing Transformer-based models [15, 69], including decision Transformer (DT). Despite the success of this approach, a new challenge arises: Transformer-based models do not adhere to the Markov property and exhibit suboptimal performance in capturing the temporal dependencies inherent in advertising (i.e., short-term data) [76].

In addition to the above challenges, existing works fall short in addressing critical considerations of **budget allocation** in advertising. One prominent approach is multi-touch attribution (MTA) [33, 44, 62], a data-driven DL-based method that attributes contributions to individual touchpoints and channels. Subsequently, the attribution scores serve as the basis for further budget allocation. However, it is noteworthy that MTA methods lack dynamic decision-making capabilities since they merely consider static allocation instead of modeling user dynamic preferences [33, 44]. Subsequently, constrained RL [12, 13, 25, 47, 79] is introduced in advertising to address the limits of MTA methods in dynamic decision-making, wherein budget limitation is incorporated as constraints within the constrained MDP (CMDP) framework. Specifically, the Lagrangian multiplier is applied to relax the problem into an unconstrained or soft-constrained one. Nevertheless, existing methods neglect dynamically selecting target users for advertising to satisfy the budget constraint, thereby overlooking an essential feature of advertising: a few users with high conversion tendencies can yield more returns than many users with low conversion tendencies [27]. Moreover, users' conversion tendencies vary over time [45].

In response to the aforementioned challenges, we present a novel multi-task offline RL model, MTORL, specifically tailored for addressing advertising concerns encompassing channel recommendation and budget allocation. We first formulate the advertising problem within the framework of offline RL. Subsequently, we employ a sequence modeling strategy to tackle the offline RL problem, where we significantly designate actions and rewards as labels to rich supervised signals, addressing the influence of overestimation and distributional shifts. In this context, we propose a causal state encoder [10] to capture temporal dependencies. Furthermore, we introduce the causal attention module [60, 68] to enhance prior information collection in the sequence. The enhanced sequence representation is directed into two branches of multi-task learning: one devoted to the action decoder, directly influencing the channel recommendation policy, and the other to the reward decoder. We

also propose direct preference optimization (DPO) loss to mitigate issues caused by highly sparse rewards. The reward predictions play a pivotal role in channel- and user-level budget allocation. The automated integration of channel recommendation and budget allocation modules guides the advertising procedure.

The major contributions to our work are summarized as follows:

- We formulate the channel recommendation and budget allocation in advertising into a well-structured offline RL problem, where each fundamental concept within advertising is rigorously defined within the framework of CMDP.
- We propose an innovative multi-task offline RL model (MTORL) for channel recommendation and budget allocation. In addition, we develop an automated advertising procedure that integrates channel recommendation and budget allocation modules. This practical framework offers an efficient solution for advertising.
- Our work encompasses extensive experiments on two public benchmark datasets. The results unequivocally demonstrate the superior performance of MTORL compared to other state-of-the-art baselines. We also conduct online experiments to validate its effectiveness in the online environment.

## 2 Preliminary

This section introduces the necessary preliminaries and formally presents the problem of multi-task offline RL for advertising.

### 2.1 Multi-channel Advertising

Given a set of users $\mathcal{U} = \{u_1, u_2, \cdots, u_{|\mathcal{U}|}\}$ and a set of channels $C = \{c_1, c_2, \cdots, c_{|C|}\}$, denote $u_i$'s static feature (e.g., user profiles) as $f_i$ and historical journey as $\mathcal{J}_i = \{o_t^i\}_{t=1}^{T_i} = \{c_t^i, q_t^i, g_t^i, w_t^i\}_{t=1}^{T_i}$, where $o_t^i$ is the observation of ad exposure at time step $t$, including advertising channel $c_t^i$, touch point feature vector $q_t^i$, touch point gain $g_t^i$ (e.g., click, conversion), and cost $w_t^i$.

**Channel Recommendation.** Supposing $C$ is fixed, online service providers (e.g., advertising platforms) are interested in making an automated strategy for the advertiser that recommends a personalized channel $c_t^i \in C$ for user $u_i$ based on her historical journey $\mathcal{J}_i$. In addition, the strategy should adjust based on new observations $o_t^i$ after each ad exposure to meet the user's dynamic preferences. The goal is to select the most suitable channel to maximize each user's total gain (e.g., clicks, conversions), thereby boosting revenue.

**Budget Allocation.** In practice, online advertisers typically have a limited budget for advertising. Therefore, online service providers cannot mindlessly advertise but must constrain the channel recommendation and accurately filter out non-target customers to control cost risks, corresponding to channel-level and user-level budget allocation modules in Section 3.7.2.

### 2.2 Multi-task Offline Reinforcement Learning

*2.2.1 Markov Decision Process.* The problem can be formulated as a MDP, which consists of $(\mathcal{S}, \mathcal{A}, P, r)$. $\mathcal{S}$ is the state space, $\mathcal{A}$ is the action space, $P(s'|s, a)$ represents the transition dynamics, $r(s, a)$ is the reward function. In the user sequence, we denote $s_t, a_t, r_t$ as state, action, and reward at time step $t$. The policy $\pi$ can be deterministic or stochastic, which outputs a unique $a$ or $P(a|s)$
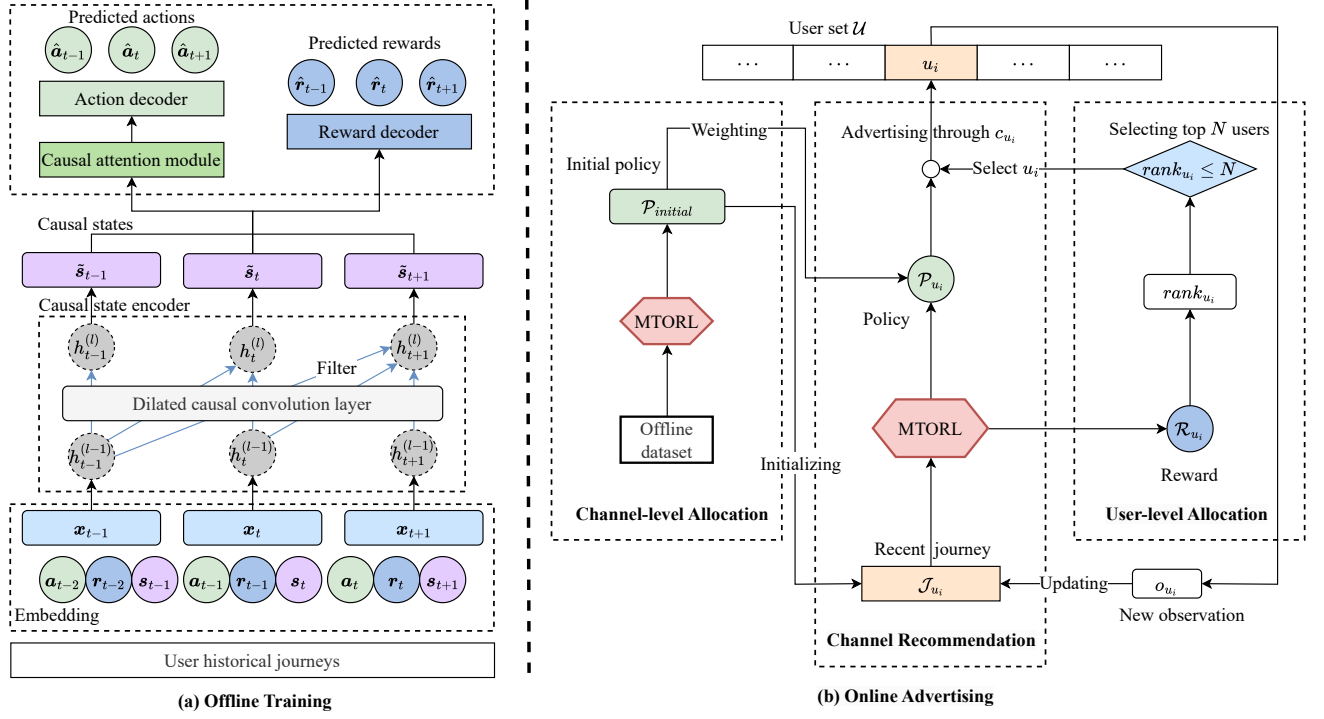
**Figure 1: The overview of MTORL. Part (a) shows the network architecture, which inputs the historical journeys and outputs predicted actions and rewards. Part (b) demonstrates the workflow of using MTORL in online advertising, where the channel recommendation module generates ad exposure policy. The channel- and user-level allocation modules control the ad exposures by reweighting policy and selecting target users, respectively.**

given state $s$, respectively. In our problem, $s$ represents the user state, $a$ means selecting ad channels for the user, and $r$ is the gain.

*2.2.2  **Optimization Objective.*** When considering the limit, the MDP problem converts to the CMDP problem [2, 47], which is to maximize the expected cumulative rewards while constraining cumulative costs. In the advertising domain, the user interests are heterogeneous, and the corresponding click-through rate (CTR) and conversion rate (CVR) are different, so we consider the user dimension additionally. We define a new optimization target that is formed as

$$\max_{\pi, \pi_u} J_r(\pi, \pi_u) := \mathbb{E}\bigg[\sum_{t=1}^{T} \sum_{u \in \mathcal{U}_t} r_{t,u}\bigg], J_w(\pi_u) = \mathbb{E}\bigg[\sum_{t=1}^{T} \sum_{u \in \mathcal{U}_t} w_{t,u}] \le W, \tag{1}$$

where $\pi$ is the policy, $W$ is the budget of the advertiser. $\pi_u = \{\mathcal{U}_1, \cdots, \mathcal{U}_t, \cdots\}$ is a filter that selects the target users $\mathcal{U}_t \subset \mathcal{U}$ for advertising at $t$, and $r_{t,u}, w_{t,u}$ are the reward and cost of ad exposure for user $u$. To solve the problem smoothly, we penalize the reward. Specifically, we transform the problem into an unconstrained one by Lagrangian multiplier [2, 67, 79]:

$$L(\pi, \pi_u, s) = J_r(\pi, \pi_u) - s \cdot J_w(\pi_u) = \mathbb{E}\bigg[\sum_{t=1}^{T} \sum_{u \in \mathcal{U}_t} r_{t,u}^*\bigg] + s \cdot W, \tag{2}$$

where $r_{t,u}^* = r_{t,u} - s \cdot w_{t,u}$ is defined as the penalized reward, controlled by the penalty strength $s$. Therefore, our target is simplified to maximize the cumulative penalized reward. We penalize the reward, and for simplicity, we still name it the *reward* in this paper.

*2.2.3  **Problem Statement.*** We formally describe our tasks as: We learn an optimal policy $\hat{\pi}(a|s)$ and reward prediction $\hat{r}(s)$ simultaneously, which means the model inputs a state $s_t$, and then outputs the action prediction $\hat{a}_t$ and reward prediction $\hat{r}_t$. The prediction $\hat{a}_t, \hat{r}_t$ will be leveraged for channel recommendation policy (i.e., $\pi$) and user filter (i.e., $\pi_u$), respectively, introduced in Section 3.7.

## 3  Methodology

As illustrated in Figure 1, MTORL consists of five major components: the causal state encoder modeling sequential dependencies, the causal attention module enhancing user preferences modeling capabilities, action and reward encoders to predict corresponding actions and rewards for the agent, respectively, and the multi-task policy optimization. In the following, we detail each component.

### 3.1  Embedding Module

Considering the conceptual discrepancy between advertising and offline RL, building a bridge to connect them is necessary. We propose a vital embedding module to map user historical exposure

features into compact embedding vectors—representing states, actions, and rewards in MDP [7, 30]. Employing sequence modeling, we generate user sequence embeddings, which lay the foundation of the following encoder-decoder framework.

### 3.1.1 *Markov Decision Process.*

Specifically, for a user $u_i$, the state of $u_i$ at time step $t$ consists of user dynamic feature $\{c_t^i, q_t^i, g_t^i, w_t^i\}$ and user static feature $f_i$. We generate states, actions, and rewards (i.e., $s_t$, $a_t$, and $r_t$) embeddings at each time step $t$ by

$$\textbf{State: } s_t = \text{Concat}(q_t, f),$$
$$\textbf{Action: } a_t = \text{OnehotEncoder}(c_t), \tag{3}$$
$$\textbf{Reward: } r_t = \text{MinMaxNorm}(g_t - sw_t).$$

OnehotEncoder($\cdot$) (e.g., dummy variable) and MinMaxNorm($\cdot$) (e.g., min-max normalization) are embedding functions generating action and reward, respectively. We default penalty strength $s$ in Equation (3) as a constant 0.5 for consistency [67, 79]. We will detail our implementations under different settings in Appendix A.1.

### 3.1.2 *Sequence Modeling.*

We consequently merge the embeddings to model users' sequential behaviors (i.e., generate a sequence embedding $X$) as follows

$$X = \{x_t\}_{t=1}^n = \{W_e \cdot \text{Concat}(a_{t-1}, r_{t-1}, s_t)\}_{t=1}^n, \tag{4}$$

where $W_e \in \mathbb{R}^{d \times F}$ is the embedding matrix, $F$ is the feature size, $n$ is the sequence length. To study how different lengths impact learning performance, we conduct comprehensive experiments in Section 4.4.1, representing an interesting trade-off between the quality and quantity of captured preferences: a moderate $n$ could introduce more reliable sequential dependencies and less noise.

### 3.2 Causal State Encoder

In advertising, the return (e.g., conversion) is influenced by current exposure and its connection to a series of preceding user exposure events. The challenge is integrating the causality between previous and current states into sequence embeddings to simulate users' dynamic interests. Towards this end, we propose the causal state encoder, for each user $u$, aggregating her previous behaviors and current state by the causal temporal convolutional network (TCN) [10]. The causal TCN leverages the dilated causal convolution to capture temporal dependencies in the sequence. Mathematically, we generate the hidden layer $H^{(l)} = (h_1^{(l)}, \cdots, h_n^{(l)})$ as

$$H^{(l)} = \text{DilatedCausalConv}(H^{(l-1)}) ; l = 1, 2, \cdots, L_1, \tag{5}$$

where $H^{(0)} = X$ is the input sequence embeddings. The causal convolution for generating the hidden state is formulated as $h_t^{(l)} = \sum_{i=0}^{k_c-1} f(i) \cdot h_{t-d_c \cdot i}^{(l-1)}$, where $d_c, k_c$ represents the dilation factor and filter size of causal convolution, and $f$ is the filter function. According to the formulation, the hidden state at time $t$ only depends on the past hidden state, i.e., the hidden state at time $k(k \leq t)$, of the previous layer instead of any future information, providing the causality. After that, we add residual connections [29] and other layers (e.g., Dropout and Normalization [63]) to enhance the sequence representation learning. Finally, we make a linear transformation for the causal TCN's output $H^{(L_1)}$ to get the causal state sequence:

$$\tilde{S} = \text{LeakyReLU}(W_s H^{(L_1)} + B_s), \tag{6}$$

where $W_s \in \mathbb{R}^{d \times d}, B_s \in \mathbb{R}^{d \times n}$ are the weight and bias, respectively. In particular, $\tilde{s}_t$ represents the causal state at time step $t$, LeakyReLU [56] is a variant of ReLU activation.

### 3.3 Causal Attention Module

In advertising, user historical behaviors showcase different impacts of user current preferences, e.g., meaningless accidental exposures or repeated behavior patterns [14, 59]. To better reflect user preferences, we introduce the causal attention mechanism, which is the foremost part of the GPT [60]. Different from the standard self-attention mechanism [68], causal attention uses a causal mask on self-attention to guarantee auto-regressive modeling and cut off information from future time steps. Formally, the causal attention module takes the causal state sequence $\tilde{S} \in \mathbb{R}^{d \times n}$ as input and outputs causal attention scores $A$ as follows

$$A = \text{CausalAttn}(\tilde{S}^{\text{T}} W_Q, \tilde{S}^{\text{T}} W_K, \tilde{S}^{\text{T}} W_V), \tag{7}$$

where $W_Q, W_K, W_V \in \mathbb{R}^{d \times d}$ are projection matrices. To alleviate overfitting issues and accelerate the training process [50], we introduce the residual connection [29] and layer normalization [9] as follows:

$$M^{(l)} = \text{LayerNorm}(M^{(l-1)} + \text{Dropout}(S^{(l-1)})),$$
$$S^{(l-1)} = \text{LeakyReLU}(W_M^{(l-1)} M^{(l-1)} + B_M^{(l-1)}), \tag{8}$$
$$M^{(0)} = A^{\text{T}} ; l = 1, 2, \cdots, L_2,$$

where $W_M^{(l)} \in \mathbb{R}^{d \times d}, B_M^{(l)} \in \mathbb{R}^{d \times n}$ are the weight and bias of layer $l$. We can learn behavior-aware sequential dependencies based on the above modules (i.e., causal state encoder and causal attention module). To further accommodate advertising, we devise two decoders (i.e., action decoder and reward decoder) based on the encoder-decoder architecture in the following.

### 3.4 Action Decoder

The channel recommendation problem entails discerning the user's optimal course of action (i.e., channel). To address this, we employ an action decoder designed to extract and interpret acquired prior information from preceding modules. This process facilitates the generation of a refined action policy, thereby enhancing the precision and effectiveness of the decision-making mechanism. With the aforementioned casual attention module, we can calculate the probability $\hat{y}_{tj}$ that a user chooses an action $j$ at time step $t$ by

$$Z = (W_Z M^{(L_2)} + B_Z)^{\text{T}},$$
$$\hat{y}_{tj} = \frac{\exp(z_{tj})}{\sum_{j=1}^m \exp(z_{tj})}, \tag{9}$$

where $W_Z \in \mathbb{R}^{m \times d}$ linearly transforms representations from hidden space into action space, $B_Z \in \mathbb{R}^{m \times n}$ is bias, and $z_{tj}$ is the action score at time step $t$ of action $j$. For the implementation, such a design choice allows us to suit different settings to make action decisions. For example, if we apply a deterministic policy, the predicted action $\hat{a}_t$ at time step $t$ should be the action with the highest $\hat{y}_{tj}$. Therefore, we can formulate the deterministic policy by

$$\hat{\pi}(a|x_{t-n+1:t}) = \underset{a^j \in \{a^1, \cdots, a^m\}}{\arg\max} \hat{y}_{tj}. \tag{10}$$

Moreover, we can seamlessly adopt the stochastic policy as follows

$$\hat{\pi}(\boldsymbol{a}|\boldsymbol{x}_{t-n+1:t}) = \boldsymbol{a}^j \quad \text{with probability } \hat{y}_{tj}. \tag{11}$$

Therefore, during the training process, we can predict actions (i.e., $\hat{\boldsymbol{a}}_t$) at each time step and force the agent to simulate users' sequential behaviors for subsequent channel recommendation inferences.

## 3.5 Reward Decoder

Predicting rewards is more challenging due to the offline RL settings. In practical scenarios, we have to learn rewards based on actions to optimize the agent, but we cannot obtain the ground truth (i.e., $\boldsymbol{a}_t$). To tackle the issues, previous offline RL studies [22, 31] used predicted actions (i.e., $\hat{\boldsymbol{a}}_t$) to further predict rewards. However, such an approach is unsuitable for recent advantages (i.e., leveraging the causal attention module to encode states). Theoretically, the attention mechanism is *not* Lipschitz continuous [35], meaning that potential noise in the inputted causal state may significantly change the attention's output of action predictions as the layer deepens.

To prevent cumulative errors, we directly leverage the causal states $\tilde{s}_t$ instead of $\hat{\boldsymbol{a}}_t$ to predict reward, retaining the flexibility of the auxiliary model meanwhile exploiting the learned prior information. The logic behind our design choice is that the proposed causal attention module can be treated as an "information filter", emphasizing highly related user exposures by assigning larger attention weights. Therefore, leveraging causal states (i.e., $\{\tilde{s}_t\}_{t=1}^{n}$) to predict rewards can inherently prevent the overestimated rewards caused by the exaggeration of the attention mechanism for some noisy patterns. Specifically, we formulate the reward decoder to predict rewards (i.e., $\hat{r}$) as follows:

$$\begin{aligned} \boldsymbol{N}^{(l)} &= \text{LeakyReLU}(\boldsymbol{W}_N^{(l-1)}\boldsymbol{N}^{(l-1)} + \boldsymbol{b}_N^{(l-1)}), \\ \boldsymbol{N}^{(0)} &= \tilde{\boldsymbol{S}}; \ l = 1, 2, \cdots, L_3, \\ \hat{\boldsymbol{r}} &= f(\boldsymbol{W}_r \boldsymbol{N}^{(L_3)} + \boldsymbol{B}_r), \end{aligned} \tag{12}$$

where $\boldsymbol{W}_N^{(l)}, \boldsymbol{B}_N^{(l)}$ are the weight and bias of layer $l$. $\boldsymbol{W}_r, \boldsymbol{b}_r$ are weight and bias at the output layer, $f(\cdot)$ is an activation function. $f(\cdot)$ is an activation function (e.g., ReLU, Sigmoid) to obtain the final reward prediction. We detail $f(\cdot)$ under different settings (e.g., binary, continuous, or multi-class rewards) in Appendix A.2.

The insight here is that using the predicted action $\hat{\boldsymbol{a}}_t$ to further predict rewards is a two-hop approach, which may bring cumulative errors. So, we predict reward based on causal states, which are obtained in a shallower layer, implicitly leveraging prior information of action prediction while mitigating the risk of error cumulation.

## 3.6 Multi-task Policy Optimization

Then, the critical challenge is to improve predictions of action and reward decoders jointly. In addition, it is necessary to learn the policy with high cumulative rewards. In order to solve these problems together, we propose a multi-task optimization objective that comprises three components: a primary focus on policy learning and two auxiliary parts—reward learning and DPO.

**Policy Learning Optimization.** Specifically, we formulate the policy learning task as cross-entropy (CE) loss by minimizing the

information entropy between predicted actions and ground truth:

$$\mathcal{L}_{Policy} = -\mathbb{E}_{\boldsymbol{x}=x_n,x_{n-1}\cdots,x_1\in\mathcal{D}}\left[\frac{1}{n}\sum_{t=1}^{n}\sum_{j=1}^{m}y_{tj}\log(\hat{y}_{tj})\right], \tag{13}$$

where $\boldsymbol{x}$ are continuous exposures from the user journeys $\mathcal{D}$.

**Reward Learning Optimization.** For binary rewards (e.g., $r_t = \{0, 1\}$), we define the auxiliary task of predicting reward as

$$\mathcal{L}_{Reward} = -\mathbb{E}_{\boldsymbol{x}\in\mathcal{D}}\left[\frac{1}{n}\sum_{t=1}^{n}\boldsymbol{r}_t\log(\hat{\boldsymbol{r}}_t)+(1-\boldsymbol{r}_t)(1-\log(\hat{\boldsymbol{r}}_t))\right]. \tag{14}$$

Note that for other settings (e.g., multi-class or continuous rewards), we also formulate variants of Eq. (14) in Appendix A.3.

**Direct Preference Optimization Loss.** To directly optimize model preferences, i.e., maximize the cumulative reward, we design a DPO loss [61] for the proposed model, formulated as

$$\mathcal{L}_{DPO} = -\mathbb{E}_{\boldsymbol{x}^w,\boldsymbol{x}^l\in\mathcal{D}}\left[\frac{1}{n}\sum_{t=1}^{n}\log\sigma(\beta\log\left(\pi(a_t^w|\boldsymbol{x}_t^w)-\beta\log(\pi(a_t^l|\boldsymbol{x}_t^l))\right)\right]. \tag{15}$$

where $\{x^w, x^l\}$ is a pair of user sequences sampled from the dataset $\mathcal{D}$ with total rewards $\sum_{t=1}^{n}\boldsymbol{r}_t^w > \sum_{t=1}^{n}\boldsymbol{r}_t^l$, and we use the default value $\beta = 0.1$. Therefore, minimizing DPO loss makes the model prefer to apply a policy that brings higher long-term returns.

**Final Optimization Objective.** According to the above three objectives, we can formulate our final optimization objective by

$$\mathcal{L} = \mathcal{L}_{Policy} + \mu\mathcal{L}_{Reward} + \lambda\mathcal{L}_{DPO}, \tag{16}$$

where $\mu, \lambda$ are the tuning parameters, adjusting contributions of each loss function for the multi-task policy optimization. To verify the impact of different parameter values, we carefully tune the parameters (i.e., $\mu$ and $\lambda$) in Section 4.4, which shows small values (e.g., 0.08 and 1.4) can affect the final performance, demonstrating the effectiveness of our multi-task policy optimization.

## 3.7 Online Advertising

Based on the above components, we can pre-train an MTORL model on the offline dataset, which can well simulate user behaviors. This section will technically detail how to deploy MTORL for online advertising. The pseudo-codes are specified in Algorithm 1.

### 3.7.1 *Channel Recommendation*.
We need to recommend an appropriate channel for each user to maximize business revenue. Specifically, for user $u_i$, we input her recent journey memory $\mathcal{J}_{u_i}$ into MTORL, as shown in Figure 1(b), and infer the policy $\mathcal{P}_{u_i}$, which is determined by the predicted action (Eq. (10)). The action prediction is specified in Figure 1(a), and notice that only the last prediction is used for inference, different from the training process. Then, we recommend channel $c_{u_i}$ for $u_i$ based on policy $\mathcal{P}_{u_i}$.

### 3.7.2 *Budget Allocation*.
It is crucial to adjust the weights of different channels from a global view and dynamically filter target users from the entire user set simultaneously to maximize the gains under the budget constraint. To address these two challenges, we propose a budget allocation strategy from both channel- and user-level, as shown in Figure 1(b).

**Channel-level Allocation.** In practice, the importance of channels to gains (e.g., click) is different [33, 62]. Learning such prior

---

**Algorithm 1** Online advertising procedure

1: **Input:** A trained MTORL model, initial policy $\mathcal{P}_{initial}$, budget $W$, channel cost $\{w_j\}_{j=1}^m$, factor $\eta$, total rounds $K$
2: **Exploration:** Initialize journey memory $\mathcal{J}$ by $\mathcal{P}_{initial}$
3: $\mathcal{P}, \mathcal{R}^s \leftarrow MTORL(\mathcal{J})$    // Initialize policy and reward
4: **Exploitation:**
5: **while** $k \le K$ and $W > 0$ **do**
6:     Rank users $\mathcal{U}$ using $\mathcal{R}$ and select top-$N$ user
7:     **for** $i = 1, \cdots, N$ **do**
8:       Use policy $\mathcal{P}_{u_i}$ to recommend a channel $c_{u_i}$
9:       Advertising for $u_i$ from channel $c_{u_i}$ and observe $o_{u_i}$
10:      Update $u_i$'s memory $\mathcal{J}_{u_i}$ by new observation $o_{u_i}$
11:      $\mathcal{P}_{u_i}, \mathcal{R}_{u_i} \leftarrow MTORL(\mathcal{J}_{u_i})$    // predict action and reward
12:      $\mathcal{P}_{u_i} \leftarrow \eta \mathcal{P}_{u_i} + (1 - \eta)\mathcal{P}_{initial}$    // weight averaging
13:      $W \leftarrow W - w_{u_i}, k \leftarrow k + 1$    // reduce the budget
14:     **end for**
15: **end while**

---

information from offline datasets and exploiting it in online advertising is challenging. To tackle this problem, we propose to learn both explicit and implicit policies. Specifically, we straightly leverage the CTR [37] to represent channels' importance so as to learn the explicit policy. Accordingly, we calculate the explicit policy by $\mathcal{P} = [p_1, p_2, \cdots, p_m]$, where $p_j = CTR_j / \sum_{i=1}^m CTR_i$ and $CTR_j$ is the CTR of channel $j$. Despite being effective, ubiquitous sparse datasets may lead to bias issues (e.g., popularity bias [1]), making the explicit policy exaggerate some channels' importance. To overcome this challenge, we learn the implicit policy by filtering out unreliable conversions. Technically, we leverage the predicted rewards (i.e., Eq. (12)) serving as the threshold, which is commonly used in transfer learning scenarios (e.g., cross-domain recommendations [66]) to adapt unbalanced data distributions. Therefore, we formulate the implicit budget ratio as $\hat{\mathcal{P}} = [\hat{p}_1, \hat{p}_2, \cdots, \hat{p}_m]$, where $p_j = \hat{CTR}_j / \sum_{i=1}^m \hat{CTR}_i$, $\hat{CTR}_j = \sum \mathbb{I}(\hat{r}_t > \tau | a^j) / N_j$, $N_j$ is the channel $j$'s ad exposures and $\mathbb{I}(\cdot)$ is the indicator function filtering reliable conversions. Afterward, we merge the explicit and implicit ratios to obtain the initial policy:

$$\mathcal{P}_{initial} = (1 - \alpha)\mathcal{P} + \alpha\hat{\mathcal{P}}, \tag{17}$$

where $\alpha$ is a tuning parameter, and $\mathcal{P}_{initial}$ is leveraged in the initial exploration and adjusting the recommendation policy.

**User-level Allocation.** We select the target users by sorting the user's reliability to reduce meaningless costs. Like the generation of policies, we use MTORL to predict rewards $\mathcal{R}_{u_i}$ for each user $u_i$ according to her journey memory $\mathcal{J}_{u_i}$. Then, we apply $\mathcal{R}$ as the sorting criterion of reliability. Since the learned reward model can directly reflect users' conversion tendencies, which are highly related to CTR and GMV. Moreover, in a scenario with massive users and new streaming users, it is difficult to model the action space for selecting users to conduct RL (the potential space is too large and complicated). Therefore, we use the predicted value of $\mathcal{R}$ as a ranking indicator for users, which can more accurately and efficiently filter high-conversion-tendency users.

*3.7.3* **Advertising Procedure.** Then, we design an automated advertising procedure integrating the above components. At the

**Table 1: Statistics of the datasets.**

| Datasets | # Users | # Interactions |
|---|---|---|
| **KuaiRand-Pure** | $27,285$ | $1,436,609$ |
| **Criteo** | $6,142,256$ | $16,468,027$ |

**Exploration** phase, we conduct advertising to all users $\mathcal{U}$ applying the stochastic policy $\mathcal{P}_{initial}$ to initialize journey memory $\mathcal{J} = \{\mathcal{J}_u, u \in \mathcal{U}\}$ for alleviating the cold-start problem [48]. Next, we feed $\mathcal{J}$ into MTORL to generate policy memory $\mathcal{P}$ and reward memory $\mathcal{R}$. Then, we get into the **Exploitation** phase. At each round, we use $\mathcal{R}^s$ to rank top-$N$ as this round's target users, and utilize $\mathcal{P}$ to recommend the advertising channels for them. We advertise to target users based on recommendations, get user feedback, and update their journey memories (by adding new observations and dropping early memories). Subsequently, we update their policy and reward memories using MTORL, prepared for the next round. The procedure will be terminal while the maximum budget or rounds is reached.

In the real-time serving environments, we could store the streaming data (i.e., user journey memory $\mathcal{J}$) and leverage incremental training methods to update the model hourly or daily (e.g., minibatch training or fine-tuning). This is efficient and stable for real-time serving, not requiring entire re-training.

## 4 Experiments

In this section, we aim to answer the following research questions:
- **RQ1**: How does MTORL perform compared with other methods regarding behavioral policy learning?
- **RQ2**: How do different components contribute to MTORL?
- **RQ3**: How does MTORL influenced by tuning parameter?
- **RQ4**: How effectiveness of MTORL in online environment?

### 4.1 Experimental Settings

*4.1.1* **Datasets and Metrics.** To evaluate the effectiveness of MTORL, we conduct experiments on two benchmark datasets: (1) **KuaiRand-Pure**[1]: it contains users' unbiased journey data with random exposure. (2) **Criteo**[2]: it contains Criteo live traffic data. Each sample represents one impression exposed to a user. For each user, we form her historical exposures chronologically as journey $\mathcal{J}_i$, specified in Section 2.1. We continue to process datasets by extracting the user exposure sequences with the minimal length of 10. In particular, to simulate multi-channel advertising, we regard the video types and campaign categories in three datasets as channels, respectively, following previous work [44, 62, 74]. The statistical information is provided in Table 1. More information on datasets is introduced in Appendix B.1. We consider widely used metrics, *F1-score*, *Precision*, and *Recall*.

*4.1.2* **Baselines.** We compare MTORL with state-of-the-art methods and briefly introduce them.
***DL for Advertising***: **Wide&Deep** [17]: consists of a wide part and a deep part, which linearly processes cross-product features and

---

[1]https://kuairand.com/
[2]http://ailab.criteo.com/criteo-attribution-modeling-bidding-dataset/

**Table 2: Overall accuracy performance comparison. The best result is bold and the second-best result is underlined in each row. All improvements are statistically significant (i.e., two-sided t-test with $p < 0.05$) over baselines.**

| Datasets | Metrics | W&D | DIN | STAR | MIREC | R-BCQ | AG | BC | CQL | IQL | DT | ODT | **MTORL** |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | F1-score | 0.5210 | 0.5649 | 0.5731 | 0.5723 | 0.6148 | 0.6060 | 0.5862 | 0.5970 | 0.6203 | 0.6296 | 0.6319 | **0.6793** |
| KuaiRand-Pure | Precision | 0.5464 | 0.5692 | 0.5744 | 0.5735 | 0.6154 | 0.6081 | 0.5863 | 0.5980 | 0.6251 | 0.6327 | 0.6390 | **0.6811** |
| | Recall | 0.5269 | 0.5669 | 0.5738 | 0.5727 | 0.6150 | 0.6065 | 0.5862 | 0.5972 | 0.6203 | 0.6308 | 0.6335 | **0.6798** |
| | F1-score | 0.2712 | 0.2974 | 0.2956 | 0.3034 | 0.3440 | 0.2982 | 0.2850 | 0.3395 | 0.3834 | 0.4303 | 0.4168 | **0.5002** |
| Criteo | Precision | 0.2966 | 0.3118 | 0.3092 | 0.3170 | 0.3854 | 0.3239 | 0.3311 | 0.3906 | 0.4230 | 0.4922 | 0.4738 | **0.5424** |
| | Recall | 0.2663 | 0.2808 | 0.2878 | 0.2991 | 0.3326 | 0.2934 | 0.2659 | 0.3173 | 0.3655 | 0.4101 | 0.3918 | **0.4893** |

captures non-linear relations using deep networks. (2) **DIN** [92]: incorporates a local activation unit, which automatically learns user representations of interests from interaction sequences tailored to specific ads. (3) **STAR** [64]: trains a unified model across all channels by simultaneously leveraging their data to capture channel-specific and channel-shared features. (4) **MIREC** [73]: introduces an effective online allocation algorithm that optimizes the exposure distribution across various channels by leveraging a global view. (5) **R-BCQ** [79]: designs $\lambda$-generalization method to merge the constraints and combine the advantages of BCQ [23] and REM [4]. (6) **AG** [13]: uses the game-theoretic value-based method and selects the single best policy.

***Offline Q-learning***: (1) **CQL** [43]: estimates a conservative Q-function to alleviate reward overestimation. (2) **IQL** [40]: approximates the policy improvement step implicitly.

***Sequence Models***: (1) **BC** (behavior cloning): utilizes an imitation learning approach. (2) **DT** [15]: leverages Transformer to address RL problems by sequence modeling. (3) **ODT** [91] (online DT): is a DT variant that combines auto-regressive and entropy regularization.

*4.1.3* **Implementation Details**. Following previous studies [68], we use Gaussian distribution to initialize parameters. We optimize MTORL utilizing Adam [36] with learning rate 0.001 and batch size 512. From the suggestions of previous studies [10, 68, 79], we set the hyper-parameters as follows: the layer number of each module is $L_1 = L_2 = 2$, $L_3 = 3$, embedding size and hidden size are 512. We obtain the optimal values of important hyper-parameters (e.g., $n = 20$) in Section 4.4. Appendix C provides more information about hyper-parameter tuning. We truncate long-length sequences (i.e., $n_i > n$) and pad short-length sequences (i.e., $n_i < n$) to guarantee equal length, more details of data processing are in Appendix B.2. The maximum number of training epochs is 800. We implement our proposed model MTORL[3] in Python 3.10.11, Pytorch 2.0.1+cu118.

## 4.2 Overall Performance Comparison (RQ1)

To demonstrate the effectiveness of MTORL for learning the behavioral policy of the dataset, we compare it with state-of-the-art baselines in prediction accuracy and reward.

*4.2.1* **Prediction Accuracy**. We report the main experimental prediction accuracy results in Table 2, averaging results of 10 time runs. Accordingly, we list some interesting observations:

---
[3]https://github.com/Applied-Machine-Learning-Lab/MTORL

**Table 3: Overall performance comparison of average reward.**

| Datasets | DIN | MIREC | R-BCQ | AG | IQL | DT | **MTORL** |
|---|---|---|---|---|---|---|---|
| KuaiRand-Pure | 0.1744 | 0.1955 | 0.1996 | 0.2086 | 0.1835 | 0.2306 | **0.2591** |
| Criteo | 0.2865 | 0.2917 | 0.2996 | 0.2968 | 0.2789 | 0.3083 | **0.3203** |

- CQL, IQL outperform W&D, DIN, STAR, MIREC, and BC on KuaiRand-Pure since they restrict policy close to user behavior, alleviating distributional shift. R-BCQ and AG perform par to CQL and IQL since they are essentially value-based RL methods.
- On KuaiRand-Pure, Compared with CQL and IQL, the Transformer-based methods (DT, ODT) have further improvements. Because of the strong capability of Transformers to learn prior information, they can discover the optimal paths from datasets.
- MTORL performs significantly better than the baselines in all cases, demonstrating its superiority. We attribute such improvements to the fact that it can capture temporal dependency using the causal state encoder. In addition, the causal attention module can detect node correlation in user sequence, which is conducive to user sequence representation learning and action decoding.
- The comparison results on Criteo dataset are similar to KuaiRand-Pure, and the main difference is that the value-based baselines (R-BCQ, AG, CQL, IQL) perform significantly worse on Criteo. The reason is that the reward signals (e.g., clicks and conversion) on this dataset are sparser, where value-based methods cannot learn user behavior validly. In contrast, MTORL performs steadily, demonstrating its robustness in dealing with sparser datasets.
- MTORL consistently achieves the best performance on both public benchmark datasets against the all baselines, validating our framework's potential benefits in real-world applications.

*4.2.2* **Reward**. We conduct experiments on MTORL and baselines, ensuring the same budget for all models, and report the average reward in Table 3. Specifically, R-BCQ and AG outperform DIN, MIREC, and IQL since they design corresponding CMDP methods to integrate cost as penalization in the Lagrangian problem, improving reward under the limited budget. DT outperforms R-BCQ and AG since it leverages conditional sequence modeling and designs return-to-go prompts to detect the journeys with high rewards. Our model, MTORL, obtains the highest average reward among all models since we devise an accurate, personalized advertising strategy.
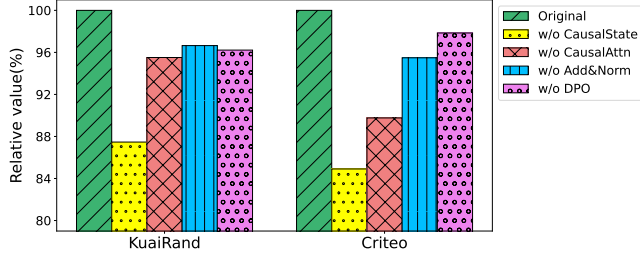
Figure 2: Impact of different components

## 4.3 Ablation Study (RQ2)

To verify the contribution of each module in the main framework, we conduct the ablation study with four variants of MTORL, including (1) *w/o CausalState*: without the causal state encoder, (2) *w/o CausalAttn*: without the causal attention module, (3) *w/o Add&Norm*: without the residual connection and layer normalization, and (4) *w/o DPO*: without the DPO loss. We record relative values of Recall to embody the performance comparison of variants and the original model, demonstrated in Figure 2. All modules of MTORL are indispensable, and the causal state encoder is the most critical part of MTORL, which contributes the most to performance, demonstrating its capability of capturing temporal dependency and local features. In addition, causal attention and DPO loss improve the performance of MTORL by capturing global features and optimizing model preference. Leveraging the residual connection and layer normalization can improve performance by enhancing sequence representation learning.

We also evaluate two variants in the budget allocation: (1) w/o Channel-level Allocation (CA), meaning random initialization policy, and (2) w/o User-level Allocation (UA), on KuaiRand. The results of average rewards are 0.2487 and 0.2353 for variants w/o CA and w/o UA, respectively, inferior to the original design of the allocation module (i.e., 0.2591). The experimental results demonstrate the effectiveness of each component. UA contributes the most by filtering users to reduce meaningless costs.

## 4.4 Parameter Analysis (RQ3)

In this section, we conduct experiments on tuning hyper-parameters $n, \mu, \lambda$ of proposed MTORL, where $n$ is the time step length, $\mu, \lambda$ are the tuning parameters of auxiliary losses in Equation (16).

*4.4.1 Time step length $n$.* Selecting a suitable $n$ is crucial for our sequence modeling, so we first tune $n$ in $\{5, 10, \cdots, 30\}$ to determine the best $n$ of KuaiRand-Pure and Criteo datasets. We demonstrate the experimental results of Recall in Figure 3(a) and Figure 3(b). We can observe that the best results are achieved when $n = 20$ and the trade-off explained in Section 3.1.2: a larger $n$ value leads to more information but less reliability (more noise), harming performance.

*4.4.2 Tuning parameter $\lambda$.* After selecting $n$, we consider learning the best $\lambda$ by tuning it in $\{0, 0.2, \cdots, 2.4\}$. Figure 3(c) and Figure 3(d) demonstrate the changing performance (i.e., Recall) when tuning $\lambda$. We determine $\lambda = 1.4$ as the final choice, as it performs well in two datasets. An appropriate $\lambda$ can reinforce main task learning by increasing user preference learning.
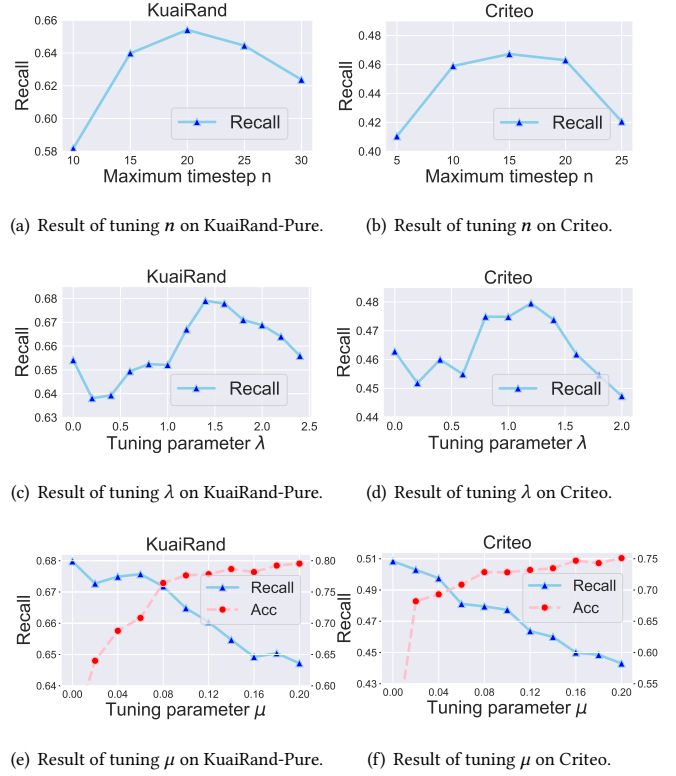


(a) Result of tuning $n$ on KuaiRand-Pure.

(b) Result of tuning $n$ on Criteo.

(c) Result of tuning $\lambda$ on KuaiRand-Pure.

(d) Result of tuning $\lambda$ on Criteo.

(e) Result of tuning $\mu$ on KuaiRand-Pure.

(f) Result of tuning $\mu$ on Criteo.

Figure 3: Results of tuning parameters.

*4.4.3 Tuning parameter $\mu$.* Then, we consider the auxiliary task: tuning $\mu$ to find an optimal point that balances reward and action prediction performance. We tune $\mu$ in $\{0, 0.02, 0.04, \cdots, 0.2\}$ and record the Recall (of action prediction) and Accuracy (of reward prediction) of MTORL. The results are shown in Figure 3(e), 3(f). The overall trend of Recall and Accuracy is down and up, demonstrating the trade-off between main and auxiliary tasks. We finally select $\mu = 0.08$ and $\mu = 0.04$ for KuaiRand-Pure and Criteo, respectively.

## 4.5 Online Experiments (RQ4)

We validate the effectiveness of MTORL in the online advertising system of Taobao, a major e-commerce platform serving millions of recommendations daily. We conduct a multi-stage training and inference process to integrate the proposed model with the deployed advertising system (i.e., the production version). Specifically, we first pre-train MTORL and leverage the pre-trained user encoder (i.e., causal state encoder and causal attention module) to produce the user embeddings in the offline environment. The produced user embeddings are integrated into the main flow of CTR prediction, and we co-train the integrated framework for a few days, obtaining +0.2pt AUC in the offline test. Then, we conduct the online A/B test with 1% online traffic for the baseline (i.e., DIN-based framework) and experimental buckets (i.e., the integrated version) for 48 hours.

We observe a gain of 0.08 CTR[4] and 0.23% RPM (Revenue Per Mille) for the experimental bucket.

## 5 Related works

### 5.1 Channel Recommendation

Channel recommendation is a recently emerging and swiftly advancing field driven by industrial challenges [52], capturing intra-channel and inter-channel features within diverse ads and delivering recommendation results from user-satisfied channels. Previous works make complicated advertising mechanisms [20, 78] based on domain-specific prior information. However, the surge of users reduces the performance of such heavy strategies. DRL formularizes this problem as an MDP to maximize long-term revenue [80, 83, 84]. Its variant, offline RL, is introduced for pre-training the ad systems on the offline datasets to reduce revenue loss during the initial launch phase [38, 39, 51], where the offline RL advancements are leveraged to address critical issues such as distributional shift. However, these approaches are less effective due to the sporadic reward signals in advertising compared to typical RL environments [16]. A line of works [31, 32] treats actions as supplementary signals, boosting training efficiency with sequential models like Transformers [15, 69]. Due to the highly sparse rewards in advertising datasets, most existing models encounter difficulty learning user behavior. Our proposed MTORL uses causal states to collect substantial prior information for sequence modeling, capturing user preferences.

### 5.2 Budget Allocation

Budget allocation in advertising has previously been explored within online convex optimization [28]. Previous works circumvent computationally intensive projection operations by imposing penalties on constraint breaches, leveraging duality principles [11, 57]. The more efficient alternative, MTA [33, 44, 62], assigns attribution to touchpoints and channels based on historical interactions. The attribution scores are used to guide budget distribution. However, MTA approaches are limited by their static allocations, as they do not account for evolving user preferences or enable dynamic decision-making [33, 44]. Therefore, constrained RL [6] is introduced to address the issues, making dynamic policies to maximize cumulative rewards with constraints. Some works [13, 27, 72, 79] leverage the Lagrangian multiplier to integrate the constraints into the main objective and relax the constraints to reduce the computational costs. However, such a method may lead to an unstable learning process, and the policy may not consistently guarantee the constraint [18, 53]. Compared with them, MTORL integrates the knapsack problem into the main framework without jeopardizing policy learning and utilizes a memory buffer, striking a good balance between efficiency and effectiveness.

## 6 Conclusion

This paper studied the problem of multi-task offline reinforcement learning for online advertising, including two scenario-specific tasks, channel recommendation and budget allocation, from a joint learning perspective. We devise causal states to encode the temporal dependencies in user sequences and capture both local and global features in the sequences, cooperating with the causal attention module. Then, the action and reward decoder extracts the useful prior information from the encoder to obtain the corresponding predictions of each time step for further processing. To better address the online advertising problem, we present a complete and automated advertising procedure within the proposed MTORL framework. In addition, we offer channel- and user-level budget allocation for macro and micro control, benefiting accurate and personalized advertising. Eventually, we conduct extensive offline and online experiments to demonstrate the effectiveness of our model in terms of online advertising.

## References

[1] Himan Abdollahpouri, Masoud Mansoury, Robin Burke, and Bamshad Mobasher. 2019. The unfairness of popularity bias in recommendation. *arXiv preprint arXiv:1907.13286* (2019).

[2] Joshua Achiam, David Held, Aviv Tamar, and Pieter Abbeel. 2017. Constrained policy optimization. In *International conference on machine learning*. PMLR, 22–31.

[3] M Mehdi Afsar, Trafford Crump, and Behrouz Far. 2022. Reinforcement learning based recommender systems: A survey. *Comput. Surveys* 55, 7 (2022), 1–38.

[4] Rishabh Agarwal, Dale Schuurmans, and Mohammad Norouzi. 2020. An optimistic perspective on offline reinforcement learning. In *International Conference on Machine Learning*. PMLR, 104–114.

[5] Noga Alon, Iftah Gamzu, and Moshe Tennenholtz. 2012. Optimizing budget allocation among channels and influencers. In *Proceedings of the 21st international conference on World Wide Web*. 381–388.

[6] Eitan Altman. 1999. *Constrained Markov decision processes*. Vol. 7. CRC press.

[7] Kai Arulkumaran, Marc Peter Deisenroth, Miles Brundage, and Anil Anthony Bharath. 2017. Deep reinforcement learning: A brief survey. *IEEE Signal Processing Magazine* 34, 6 (2017), 26–38.

[8] Vashist Avadhanula, Riccardo Colini Baldeschi, Stefano Leonardi, Karthik Abinav Sankararaman, and Okke Schrijvers. 2021. Stochastic bandits for multi-platform budget optimization in online advertising. In *Proceedings of the Web Conference 2021*. 2805–2817.

[9] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. 2016. Layer normalization. *arXiv preprint arXiv:1607.06450* (2016).

[10] Shaojie Bai, J Zico Kolter, and Vladlen Koltun. 2018. An empirical evaluation of generic convolutional and recurrent networks for sequence modeling. *arXiv preprint arXiv:1803.01271* (2018).

[11] Santiago Balseiro, Haihao Lu, and Vahab Mirrokni. 2020. Dual mirror descent for online allocation problems. In *International Conference on Machine Learning*. PMLR, 613–628.

[12] Qingpeng Cai, Zhenghai Xue, Chi Zhang, Wanqi Xue, Shuchang Liu, Ruohan Zhan, Xueliang Wang, Tianyou Zuo, Wentao Xie, Dong Zheng, et al. 2023. Two-Stage Constrained Actor-Critic for Short Video Recommendation. In *Proceedings of the ACM Web Conference 2023*. 865–875.

[13] Tianchi Cai, Jiyan Jiang, Wenpeng Zhang, Shiji Zhou, Xierui Song, Li Yu, Lihong Gu, Xiaodong Zeng, Jinjie Gu, and Guannan Zhang. 2023. Marketing Budget Allocation with Offline Constrained Deep Reinforcement Learning. In *Proceedings of the Sixteenth ACM International Conference on Web Search and Data Mining*. 186–194.

[14] Huiyuan Chen, Yusan Lin, Menghai Pan, Lan Wang, Chin-Chia Michael Yeh, Xiaoting Li, Yan Zheng, Fei Wang, and Hao Yang. 2022. Denoising self-attentive sequential recommendation. In *Proceedings of the 16th ACM Conference on Recommender Systems*. 92–101.

[15] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter Abbeel, Aravind Srinivas, and Igor Mordatch. 2021. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems* 34 (2021), 15084–15097.

---

[4] Gain at 0.001-level is regarded as significant for the CTR prediction task [17]

[16] Xiaocong Chen, Lina Yao, Julian McAuley, Guanglin Zhou, and Xianzhi Wang. 2021. A survey of deep reinforcement learning in recommender systems: A systematic review and future directions. *arXiv preprint arXiv:2109.03540* (2021).

[17] Heng-Tze Cheng, Levent Koc, Jeremiah Harmsen, Tal Shaked, Tushar Chandra, Hrishi Aradhye, Glen Anderson, Greg Corrado, Wei Chai, Mustafa Ispir, et al. 2016. Wide & deep learning for recommender systems. In *Proceedings of the 1st workshop on deep learning for recommender systems*. 7–10.

[18] Yinlam Chow, Ofir Nachum, Aleksandra Faust, Edgar Duenez-Guzman, and Mohammad Ghavamzadeh. 2019. Lyapunov-based safe policy optimization for continuous control. *arXiv preprint arXiv:1901.10031* (2019).

[19] Peter J Danaher, Janghyuk Lee, and Laoucine Kerbache. 2010. Optimal internet media selection. *Marketing Science* 29, 2 (2010), 336–347.

[20] Evert De Haan, Thorsten Wiesel, and Koen Pauwels. 2016. The effectiveness of different forms of online advertising for purchase conversion in a multiple-channel attribution framework. *International journal of research in marketing* 33, 3 (2016), 491–507.

[21] Ming Fan, Subodha Kumar, and Andrew B Whinston. 2007. Selling or advertising: Strategies for providing digital media online. *Journal of Management Information Systems* 24, 3 (2007), 143–166.

[22] Chelsea Finn, Sergey Levine, and Pieter Abbeel. 2016. Guided cost learning: Deep inverse optimal control via policy optimization. In *International conference on machine learning*. PMLR, 49–58.

[23] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-policy deep reinforcement learning without exploration. In *International conference on machine learning*. PMLR, 2052–2062.

[24] Jingtong Gao, Yewen Li, Shuai Mao, Peng Jiang, Nan Jiang, Yejing Wang, Qingpeng Cai, Fei Pan, Kun Gai, Bo An, et al. 2025. Generative Auto-Bidding with Value-Guided Explorations. *arXiv preprint arXiv:2504.14587* (2025).

[25] Yingqiang Ge, Shuchang Liu, Ruoyuan Gao, Yikun Xian, Yunqi Li, Xiangyu Zhao, Changhua Pei, Fei Sun, Junfeng Ge, Wenwu Ou, et al. 2021. Towards long-term fairness in recommendation. In *Proceedings of the 14th ACM international conference on web search and data mining*. 445–453.

[26] Qi Hao, Tianze Luo, and Guangda Huzhang. 2021. Re-ranking with constraints on diversified exposures for homepage recommender system. *arXiv preprint arXiv:2112.07621* (2021).

[27] Xiaotian Hao, Zhaoqing Peng, Yi Ma, Guan Wang, Junqi Jin, Jianye Hao, Shan Chen, Rongquan Bai, Mingzhou Xie, Miao Xu, et al. 2020. Dynamic knapsack optimization towards efficient multi-channel sequential advertising. In *International Conference on Machine Learning*. PMLR, 4060–4070.

[28] Elad Hazan et al. 2016. Introduction to online convex optimization. *Foundations and Trends® in Optimization* 2, 3-4 (2016), 157–325.

[29] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.

[30] Peter Henderson, Riashat Islam, Philip Bachman, Joelle Pineau, Doina Precup, and David Meger. 2018. Deep reinforcement learning that matters. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 32.

[31] Jonathan Ho and Stefano Ermon. 2016. Generative adversarial imitation learning. *Advances in neural information processing systems* 29 (2016).

[32] Ahmed Hussein, Mohamed Medhat Gaber, Eyad Elyan, and Chrisina Jayne. 2017. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)* 50, 2 (2017), 1–35.

[33] Wendi Ji and Xiaoling Wang. 2017. Additional multi-touch attribution for online advertising. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 31.

[34] Leslie Pack Kaelbling, Michael L Littman, and Andrew W Moore. 1996. Reinforcement learning: A survey. *Journal of artificial intelligence research* 4 (1996), 237–285.

[35] Hyunjik Kim, George Papamakarios, and Andriy Mnih. 2021. The lipschitz constant of self-attention. In *International Conference on Machine Learning*. PMLR, 5562–5571.

[36] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).

[37] Simon Kingsnorth. 2022. *Digital marketing strategy: an integrated approach to online marketing*. Kogan Page Publishers.

[38] Haruka Kiyohara, Kosuke Kawakami, and Yuta Saito. 2021. Accelerating offline reinforcement learning application in real-time bidding and recommendation: Potential use of simulation. *arXiv preprint arXiv:2109.08331* (2021).

[39] Dmytro Korenkevych, Frank Cheng, Artsiom Balakir, Alex Nikulkov, Lingnan Gao, Zhihao Cen, Zuobing Xu, and Zheqing Zhu. 2024. Offline reinforcement learning for optimizing production bidding policies. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 5251–5259.

[40] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. 2021. Offline reinforcement learning with implicit q-learning. *arXiv preprint arXiv:2110.06169* (2021).

[41] Yueh-Ning Ku, Mikhail Kuznetsov, Shaunak Mishra, and Paloma de Juan. 2023. Staging e-commerce products for online advertising using retrieval assisted image generation. *arXiv preprint arXiv:2307.15326* (2023).

[42] Aviral Kumar, Justin Fu, Matthew Soh, George Tucker, and Sergey Levine. 2019. Stabilizing off-policy q-learning via bootstrapping error reduction. *Advances in Neural Information Processing Systems* 32 (2019).

[43] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. 2020. Conservative q-learning for offline reinforcement learning. *Advances in Neural Information Processing Systems* 33 (2020), 1179–1191.

[44] Sachin Kumar, Garima Gupta, Ranjitha Prasad, Arnab Chatterjee, Lovekesh Vig, and Gautam Shroff. 2020. Camta: Causal attention model for multi-touch attribution. In *2020 International Conference on Data Mining Workshops (ICDMW)*. IEEE, 79–86.

[45] Kuang-chih Lee, Burkay Orten, Ali Dasdan, and Wentong Li. 2012. Estimating conversion rate in display advertising from past erformance data. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining*. 768–776.

[46] Xiaopeng Li, Fan Yan, Xiangyu Zhao, Yichao Wang, Bo Chen, Huifeng Guo, and Ruiming Tang. 2023. Hamur: Hyper adapter for multi-domain recommendation. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*. 1268–1277.

[47] Guogang Liao, Ze Wang, Xiaoxu Wu, Xiaowen Shi, Chuheng Zhang, Yongkang Wang, Xingxing Wang, and Dong Wang. 2022. Cross dqn: Cross deep q network for ads allocation in feed. In *Proceedings of the ACM Web Conference 2022*. 401–409.

[48] Blerina Lika, Kostas Kolomvatsos, and Stathes Hadjiefthymiades. 2014. Facing the cold start problem in recommender systems. *Expert systems with applications* 41, 4 (2014), 2065–2073.

[49] Weilin Lin, Xiangyu Zhao, Yejing Wang, Yuanshao Zhu, and Wanyu Wang. 2023. Autodenoise: Automatic data instance denoising for recommendations. In *Proceedings of the ACM Web Conference 2023*. 1003–1011.

[50] Langming Liu, Liu Cai, Chi Zhang, Xiangyu Zhao, Jingtong Gao, Wanyu Wang, Yifu Lv, Wenqi Fan, Yiqi Wang, Ming He, et al. 2023. Linrec: Linear attention mechanism for long-term sequential recommender systems. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*. 289–299.

[51] Tao Liu, Qi Xu, Wei Shi, Zhigang Hua, and Shuang Yang. 2025. Session-Level Dynamic Ad Load Optimization using Offline Robust Reinforcement Learning. *arXiv preprint arXiv:2501.05591* (2025).

[52] Weiwen Liu, Yunjia Xi, Jiarui Qin, Fei Sun, Bo Chen, Weinan Zhang, Rui Zhang, and Ruiming Tang. 2022. Neural re-ranking in multi-stage recommender systems: A review. *arXiv preprint arXiv:2202.06602* (2022).

[53] Yongshuai Liu, Avishai Halev, and Xin Liu. 2021. Policy learning with constraints in model-free reinforcement learning: A survey. In *The 30th International Joint Conference on Artificial Intelligence (IJCAI)*.

[54] Ziru Liu, Shuchang Liu, Bin Yang, Zhenghai Xue, Qingpeng Cai, Xiangyu Zhao, Zijian Zhang, Lantao Hu, Han Li, and Peng Jiang. 2024. Modeling User Retention through Generative Flow Networks. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*. 5497–5508.

[55] Ziru Liu, Jiejie Tian, Qingpeng Cai, Xiangyu Zhao, Jingtong Gao, Shuchang Liu, Dayou Chen, Tonghao He, Dong Zheng, Peng Jiang, et al. 2023. Multi-task recommendations with reinforcement learning. In *Proceedings of the ACM web conference 2023*. 1273–1282.

[56] Andrew L Maas, Awni Y Hannun, Andrew Y Ng, et al. 2013. Rectifier nonlinearities improve neural network acoustic models. In *Proc. icml*, Vol. 30. Atlanta, GA, 3.

[57] Joseph Naor and David Wajc. 2018. Near-optimum online ad allocation for targeted advertising. *ACM Transactions on Economics and Computation (TEAC)* 6, 3-4 (2018), 1–20.

[58] Junwei Pan, Yizhi Mao, Alfonso Lobos Ruiz, Yu Sun, and Aaron Flores. 2019. Predicting different types of conversions with multi-task learning in online advertising. In *Proceedings of the 25th acm sigkdd international conference on knowledge discovery & data mining*. 2689–2697.

[59] Yuqi Qin, Pengfei Wang, and Chenliang Li. 2021. The world is binary: Contrastive learning for denoising next basket recommendation. In *Proceedings of the 44th international ACM SIGIR conference on research and development in information retrieval*. 859–868.

[60] Alec Radford, Karthik Narasimhan, Tim Salimans, Ilya Sutskever, et al. 2018. Improving language understanding by generative pre-training. (2018).

[61] Rafael Rafailov, Archit Sharma, Eric Mitchell, Christopher D Manning, Stefano Ermon, and Chelsea Finn. 2024. Direct preference optimization: Your language model is secretly a reward model. *Advances in Neural Information Processing Systems* 36 (2024).

[62] Kan Ren, Yuchen Fang, Weinan Zhang, Shuhao Liu, Jiajun Li, Ya Zhang, Yong Yu, and Jun Wang. 2018. Learning multi-touch conversion attribution with dual-attention mechanisms for online advertising. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*. 1433–1442.

[63] Tim Salimans and Durk P Kingma. 2016. Weight normalization: A simple reparameterization to accelerate training of deep neural networks. *Advances in neural information processing systems* 29 (2016).

[64] Xiang-Rong Sheng, Liqin Zhao, Guorui Zhou, Xinyao Ding, Binding Dai, Qiang Luo, Siran Yang, Jingshan Lv, Chi Zhang, Hongbo Deng, et al. 2021. One model to serve all: Star topology adaptive recommender for multi-domain ctr prediction. In

*Proceedings of the 30th ACM International Conference on Information & Knowledge Management.* 4104–4113.

[65] Jing-Cheng Shi, Yang Yu, Qing Da, Shi-Yong Chen, and An-Xiang Zeng. 2019. Virtual-taobao: Virtualizing real-world online retail environment for reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 33. 4902–4909.

[66] Jie Tang, Sen Wu, Jimeng Sun, and Hang Su. 2012. Cross-domain collaboration recommendation. In *Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining.* 1285–1293.

[67] Chen Tessler, Daniel J Mankowitz, and Shie Mannor. 2018. Reward constrained policy optimization. *arXiv preprint arXiv:1805.11074* (2018).

[68] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

[69] Siyu Wang, Xiaocong Chen, Dietmar Jannach, and Lina Yao. 2023. Causal decision transformer for recommender systems via offline reinforcement learning. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval.* 1599–1608.

[70] Yuhao Wang, Xiangyu Zhao, Bo Chen, Qidong Liu, Huifeng Guo, Huanshuo Liu, Yichao Wang, Rui Zhang, and Ruiming Tang. 2023. PLATE: A prompt-enhanced paradigm for multi-scenario recommendations. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval.* 1498–1507.

[71] Marco A Wiering and Martijn Van Otterlo. 2012. Reinforcement learning. *Adaptation, learning, and optimization* 12, 3 (2012), 729.

[72] Shuai Xiao, Le Guo, Zaifan Jiang, Lei Lv, Yuanbo Chen, Jun Zhu, and Shuang Yang. 2019. Model-based constrained MDP for budget allocation in sequential incentive marketing. In *Proceedings of the 28th ACM International Conference on Information and Knowledge Management.* 971–980.

[73] Yue Xu, Qijie Shen, Jianwen Yin, Zengde Deng, Dimin Wang, Hao Chen, Lixiang Lai, Tao Zhuang, and Junfeng Ge. 2023. Multi-channel Integrated Recommendation with Exposure Constraints. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining.* 5338–5349.

[74] Di Yao, Chang Gong, Lei Zhang, Sheng Chen, and Jingping Bi. 2022. CausalMTA: Eliminating the User Confounding Bias for Causal Multi-touch Attribution. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining.* 4342–4352.

[75] Daniel Zantedeschi, Eleanor McDonnell Feit, and Eric T Bradlow. 2017. Measuring multichannel advertising response. *Management Science* 63, 8 (2017), 2706–2728.

[76] Ailing Zeng, Muxi Chen, Lei Zhang, and Qiang Xu. 2023. Are transformers effective for time series forecasting?. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 37. 11121–11128.

[77] Weinan Zhang, Xiangyu Zhao, Li Zhao, Dawei Yin, Grace Hui Yang, and Alex Beutel. 2020. Deep reinforcement learning for information retrieval: Fundamentals and advances. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval.* 2468–2471.

[78] Xubing Zhang. 2009. Retailers' multichannel and price advertising strategies. *Marketing Science* 28, 6 (2009), 1080–1094.

[79] Yang Zhang, Bo Tang, Qingyu Yang, Dou An, Hongyin Tang, Chenyang Xi, Xueying Li, and Feiyu Xiong. 2021. BCORLE ($\lambda$): An Offline Reinforcement Learning and Evaluation Framework for Coupons Allocation in E-commerce Market. *Advances in Neural Information Processing Systems* 34 (2021), 20410–20422.

[80] Jun Zhao, Guang Qiu, Ziyu Guan, Wei Zhao, and Xiaofei He. 2018. Deep reinforcement learning for sponsored search real-time bidding. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining.* 1021–1030.

[81] Kesen Zhao, Lixin Zou, Xiangyu Zhao, Maolin Wang, and Dawei Yin. 2023. User retention-oriented recommendation with decision transformer. In *Proceedings of the ACM Web Conference 2023.* 1141–1149.

[82] Xiangyu Zhao. 2022. Adaptive and automated deep recommender systems. *ACM SIGWEB Newsletter* 2022, Spring (2022), 1–4.

[83] Xiangyu Zhao, Changsheng Gu, Haoshenglun Zhang, Xiwang Yang, Xiaobing Liu, Jiliang Tang, and Hui Liu. 2021. Dear: Deep reinforcement learning for online advertising impression in recommender systems. In *Proceedings of the AAAI conference on artificial intelligence*, Vol. 35. 750–758.

[84] Xiangyu Zhao, Long Xia, Jiliang Tang, and Dawei Yin. 2019. "Deep reinforcement learning for search, recommendation, and online advertising: a survey" by Xiangyu Zhao, Long Xia, Jiliang Tang, and Dawei Yin with Martin Vesely as coordinator. *ACM sigweb newsletter* 2019, Spring (2019), 1–15.

[85] Xiangyu Zhao, Long Xia, Liang Zhang, Zhuoye Ding, Dawei Yin, and Jiliang Tang. 2018. Deep reinforcement learning for page-wise recommendations. In *Proceedings of the 12th ACM conference on recommender systems.* 95–103.

[86] Xiangyu Zhao, Long Xia, Lixin Zou, Hui Liu, Dawei Yin, and Jiliang Tang. 2020. Whole-chain recommendations. In *Proceedings of the 29th ACM international*

[87] Xiangyu Zhao, Long Xia, Lixin Zou, Hui Liu, Dawei Yin, and Jiliang Tang. 2021. Usersim: User simulation via supervised generativeadversarial network. In *Proceedings of the Web Conference 2021.* 3582–3589.

[88] Xiangyu Zhao, Liang Zhang, Zhuoye Ding, Long Xia, Jiliang Tang, and Dawei Yin. 2018. Recommendations with negative feedback via pairwise deep reinforcement learning. In *Proceedings of the 24th ACM SIGKDD international conference on knowledge discovery & data mining.* 1040–1048.

[89] Xiangyu Zhao, Liang Zhang, Long Xia, Zhuoye Ding, Dawei Yin, and Jiliang Tang. 2017. Deep reinforcement learning for list-wise recommendations. *arXiv preprint arXiv:1801.00209* (2017).

[90] Xiangyu Zhao, Xudong Zheng, Xiwang Yang, Xiaobing Liu, and Jiliang Tang. 2020. Jointly learning to recommend and advertise. In *Proceedings of the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining.* 3319–3327.

[91] Qinqing Zheng, Amy Zhang, and Aditya Grover. 2022. Online decision transformer. In *international conference on machine learning.* PMLR, 27042–27059.

[92] 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 international conference on knowledge discovery & data mining.* 1059–1068.

[93] Hao Zhou, Shaoming Li, Guibin Jiang, Jiaqi Zheng, and Dong Wang. 2023. Direct heterogeneous causal learning for resource allocation problems in marketing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 37. 5446–5454.

## A Reward Construction

### A.1 Embedding Module

Since the practice scenarios are complex (e.g., multi-class or continuous gains for reward embedding learning), we detail our implementations of reward embedding under different settings. The multi-class gains mean that we have various categories of gain. Usually, the multi-class gains are hard to merge. For example, for video exposure, gains consist of clicks, likes, follows, comments, and so on, which are independent and all significant. In that case, we can easily leverage the one-hot encoder. The multi-class gains can sometimes be merged. For example, if gains are clicks and conversions, we can design a simple rule to merge them, such as a weighted sum. We leverage different embedding functions in response to different situations. The embedding functions of reward can be formulated as follows

$$r_t = \begin{cases} g_{t+1}, & \text{if one-class and binary,} \\ \text{Normalize}(g_{t+1}), & \text{if one-class and continuous,} \\ \text{Normalize}(\text{Fusion}(g_{t+1})), & \text{if multi-class and compatible,} \\ \text{OnehotEncoder}(g_{t+1}), & \text{if multi-class and incompatible,} \end{cases}$$
(18)

where $\text{Normalize}(\cdot)$ represents min-max normalization, $\text{Fusion}(\cdot)$ means we add categories of gain by weight. For example, we weigh clicks as 1 and conversions as 10 and add them together for ad exposure. $\text{onehotEncoder}(\cdot)$ is generated directly by dummy variables when the number of categories is small. If the number of categories is large, we will conduct word2vec or other embedding to reduce the reward dimension.

### A.2 Reward Decoder

If the reward is continuous, we utilize a mapping function as the activation function $f(\cdot)$ to map the reward into $(0, 1)$; if the reward is binary, for example, Click/Conversion or not, we apply $\sigma(\cdot)$ function as $f(\cdot)$ to map values into $(0, 1)$ to get reward probability; if the reward is multi-class, we use the Softmax function as $f(\cdot)$.

**Table 4: Hyper-parameters and Searching Ranges**

| Hyper-parameters | Search range |
|---|---|
| batch size | [32, 64, 128, 256, 512, 1024, 2048] |
| learning rate | [1e-2, 5e-3, 1e-3, 5e-4, 1e-4] |
| dropout rate | [0.1, 0.2, 0.5] |
| weight decay rate | [0, 1e-4, 1e-5, 1e-6] |
| embedding size | [64, 128, 256, 512] |
| hidden size | [128, 256, 512, 768] |
| number of layers | [2, 3, 4] |
| number of heads | [1, 2, 4, 8] |
| time step length | [5, 10, 15, 20, 25, 30, 35, 40] |

## A.3 Multi-task Policy Optimization

We define the auxiliary task to predict reward similarly by a mean squared error (MSE) loss if the reward is continuous:

$$\mathcal{L}_{Reward} = \mathbb{E}_{\boldsymbol{x} \in \mathcal{D}} \Big[ \frac{1}{n} \sum_{t=1}^{n} (\boldsymbol{r}_t - \hat{\boldsymbol{r}}_t)^2 \Big], \tag{19}$$

or a binary cross-entropy (BCE) loss if the reward is binary (e.g., $r_t = \{0, 1\}$):

$$\mathcal{L}_{Reward} = \mathbb{E}_{\boldsymbol{x} \in \mathcal{D}} \Big[ -\frac{1}{n} \sum_{t=1}^{n} \boldsymbol{r}_t \log(\hat{\boldsymbol{r}}_t) + (1 - \boldsymbol{r}_t)(1 - \log(\hat{\boldsymbol{r}}_t)) \Big], \tag{20}$$

or a CE loss if the reward is multi-class:

$$\mathcal{L}_{Reward} = \mathbb{E}_{\boldsymbol{x} \in \mathcal{D}} \Big[ -\frac{1}{n} \sum_{t=1}^{n} \sum_{j} \boldsymbol{r}_{tj} \log(\hat{\boldsymbol{r}}_{tj}) \Big]. \tag{21}$$

## B Datasets

### B.1 Raw Data

We evaluate our proposed model on two benchmark datasets, both are large-scale enough and possess affluent user features that are frequently used in online advertising tasks, such as CTR prediction, bidding, and MTA. Here, we give more detailed information about each dataset.

- **KuaiRand-Pure**: It is an open-sourced, unbiased dataset that records abundant random video exposure and user feedback in the Kuaishou app. The KuaiRand-Pure dataset consists of three log files: a file of user features, a file of video features, and a file of video statistics. The entire time range of log files is 30 days, and each log file contains specific details of video exposures, such as user ID, video ID, time, click, like, follow, playtime, and so on, which are suitable for online advertising tasks. Significantly, the timestamp can guide us in reordering each user's interactions chronologically for sequential modeling. In addition, the file of user features provides sufficient information to complement the

interaction information. The descriptions of video features, such as video type, music ID, and tag, are conducive to modeling and simulating the advertising channels.
- **Criteo**: It is a public commercial dataset named Criteo Attribution Modeling for Bidding Dataset. Namely, it additionally contains attribution data compared to standard online advertising datasets for P value (CTR or CVR prediction). The dataset consists of natural flow data in 30 days. It also contains sufficient information on ad exposures, such as conversion, click, conversion timestamp, click position, cost, and cpo, which is conducive to online advertising modeling. The dataset contains about 700 ad campaigns, often used in channel or touchpoint attribution. It also provides nine contextual features related to ad exposures, which can be leveraged in modeling clicks or conversion.

### B.2 Data Processing

After choosing the datasets, we need to convert them into a form that can be used for model training and evaluation. The first thing to do is classify the ad/video exposures by user ID. Furthermore, according to the timestamp value, we form the exposures into user journeys in chronological order, benefiting sequential modeling. We filter out very short user journeys without generality to reduce noise in model training. We merge the corresponding user and item features into ad exposures to enhance information richness. The merge method is simply a concatenation operation. Moreover, some automated methods (e.g., sklearn.feature_selection) and regulation (e.g., filter out features with too many Na) of feature selection are used to refine the data and reduce the noise. The data are then processed to the same length $n$, aligning with the models' inputs. Specifically, we truncate the long and pad short user journeys to guarantee all journeys are the same length. The last thing is to split the processed same-length user journeys into training, validation, and test sets, where we set $0.8 : 0.1 : 0.1$ as their ratio.

## C Hyper-parameter

The selection of the hyper-parameters is critical to the deep learning model performance. To guarantee fairness and explore the optimal performance of each model, we conduct thorough and meticulous hyper-parameter tuning for both proposed MTORL and baselines. Besides the special hyper-parameter of MTORL that has been analyzed in Section 4.4, here we target tuning general hyper-parameters for both the proposed method and baseline methods. Following the previous work, the general hyper-parameters we consider for tuning are batch size, learning rate, dropout rate, weight decay rate, embedding size, hidden size, number of layers, number of heads, and time step length. Moreover, we set the searching ranges for all hyperparameters, illustrated in Table 4. Notice that the optimal number of layers is searched for all deep neural network architectures (e.g., MLP, TCN, ResNet, Transformer).