# Generative Auto-Bidding with Value-Guided Explorations

Jingtong Gao
City University of Hong Kong
Hong Kong, China
jt.g@my.cityu.edu.hk

Yewen Li
Nanyang Technological University
Singapore, Singapore
yewen001@e.ntu.edu.sg

Shuai Mao
The Chinese University of Hong Kong
Hong Kong, China
smao@mae.cuhk.edu.hk

Peng Jiang
Kuaishou Technology
Beijing, China
jiangpeng07@kuaishou.com

Nan Jiang
Kuaishou Technology
Beijing, China
jiangnan07@kuaishou.com

Yejing Wang
City University of Hong Kong
Hong Kong, China
yejing.wang@my.cityu.edu.hk

Qingpeng Cai*
Kuaishou Technology
Beijing, China
caiqingpeng@kuaishou.com

Fei Pan
Kuaishou Technology
Beijing, China
panfei05@kuaishou.com

Peng Jiang
Kuaishou Technology
Beijing, China
jiangpeng@kuaishou.com

Kun Gai
Unaffiliated
Beijing, China
gai.kun@qq.com

Bo An
Nanyang Technological University
Singapore, Singapore
boan@ntu.edu.sg

Xiangyu Zhao*
City University of Hong Kong
Hong Kong, China
xianzhao@cityu.edu.hk

## Abstract

Auto-bidding, with its strong capability to optimize bidding decisions within dynamic and competitive online environments, has become a pivotal strategy for advertising platforms. Existing approaches typically employ rule-based strategies or Reinforcement Learning (RL) techniques. However, rule-based strategies lack the flexibility to adapt to time-varying market conditions, and RL-based methods struggle to capture essential historical dependencies and observations within Markov Decision Process (MDP) frameworks. Furthermore, these approaches often face challenges in ensuring strategy adaptability across diverse advertising objectives. Additionally, as offline training methods are increasingly adopted to facilitate the deployment and maintenance of stable online strategies, the issues of documented behavioral patterns and behavioral collapse resulting from training on fixed offline datasets become increasingly significant. To address these limitations, this paper introduces a novel offline **G**enerative **A**uto-bidding framework with **V**alue-Guided **E**xplorations (**GAVE**). GAVE accommodates various advertising objectives through a score-based Return-To-Go (RTG) module. Moreover, GAVE integrates an action exploration mechanism with an RTG-based evaluation method to explore novel actions while ensuring stability-preserving updates. A learnable value function is also designed to guide the direction of action exploration and

mitigate Out-of-Distribution (OOD) problems. Experimental results on two offline datasets and real-world deployments demonstrate that GAVE outperforms state-of-the-art baselines in both offline evaluations and online A/B tests. By applying the core methods of this framework, we proudly secured first place in the NeurIPS 2024 competition, 'AIGB Track: Learning Auto-Bidding Agents with Generative Models' [1]. The implementation code is publicly available to facilitate reproducibility and further research [2].

## CCS Concepts

• **Information systems → Computational advertising**.

## Keywords

Auto-bidding, Generative Model, Decision Transformer

## 1 Introduction

Bidding remains a fundamental component of modern online advertising platforms [11, 14, 58, 62], enabling businesses to engage target audiences and increase sales. As digital advertising evolves, traditional manual bid management methods have become increasingly inadequate and cost-ineffective, failing to address the demands of today's dynamic market [18, 34]. The complexity of modern advertising systems, characterized by fluctuating market conditions and diverse user behaviors [1], requires automated solutions that can

---

*Corresponding authors.

---

[1] https://tianchi.aliyun.com/competition/entrance/532236/rankingList
[2] https://github.com/Applied-Machine-Learning-Lab/GAVE

adapt to these variations while aligning with advertisers' diverse objectives [4, 48]. This need is further intensified by the massive volume of ad impressions requiring real-time processing, where human intervention becomes impractical and often suboptimal for achieving advertising goals [4, 16].

To meet these requirements, existing solutions have evolved into two primary categories: predefined rule-based strategies [7, 55] and Reinforcement Learning (RL)-based methods [5, 13, 54, 59]. While rule-based strategies are computationally lightweight and straightforward to deploy, their static nature makes them ill-suited for dynamic markets and unable to accommodate advertisers' diverse demands [5, 15, 24, 30]. RL-based approaches, though employing Markov Decision Processes (MDPs) [35, 54, 60, 61] to adapt to environmental changes and obtain better performance, face a critical structural constraint: the MDP framework's state-independence assumption inherently disregards temporal dependencies and observations within bidding sequences [6, 25, 27]. This limitation obstructs the identification of evolving behavioral patterns and market fluctuations, substantially undermining RL's real-world applicability in highly volatile real-time bidding environments.

Recently, Decision Transformer (DT) [6, 10, 12] has emerged as a powerful framework, effectively capturing temporal dependencies and historical context. Therefore, applying DT to offline bidding modeling offers a promising direction for improving strategies. Specifically, by adopting an offline training paradigm, DT circumvents the risks and implementation challenges of online training, ensuring broader applicability across diverse scenarios [9, 26, 31, 49]. The generative modeling foundation of DT further enables explicit capture of temporal dependencies and historical bidding context, allowing adaptive decision-making that aligns with the dynamic nature of real-world advertising environments.

However, several critical challenges emerge when implementing a DT-based auto-bidding approach in real-world scenarios. First, practical deployment requires accommodating complex advertising objectives where evaluation metrics extend beyond elementary indicators like total clicks or conversions. These objectives typically involve sophisticated functions with different preferences on interdependent parameters—such as Cost Per Action (CPA) thresholds and Cost Per Click (CPC) ceilings [16, 27]—requiring DT modeling with adaptive optimization objectives to align with diverse operational criteria. Second, directly training DT models in offline environments may restrict to documented behavioral patterns [21, 29] and suffer from behavioral collapse [9], which requires amplified action explorations with stable updates.

To tackle these challenges, we propose a unified framework that enhances DT for offline **G**enerative **A**uto-bidding through **V**alue-guided **E**xplorations (GAVE). First, to accommodate complex advertising objectives [16, 17], we design **a score-based Return-To-Go (RTG)** module with customizable score functions, enabling adaptive modeling of various objective requirements like CPA constraints through differentiable programming. Second, **an action exploration mechanism** is proposed alongside an RTG-based evaluation method to explore and evaluate actions outside the fixed dataset while ensuring stability-preserving updates between explored and original actions. However, it is rather challenging to learn a beneficial strategy through random explorations and avoid

Out-of-Distribution (OOD) risks in such a sensitive bidding environment with large action space [20, 32, 53]. Thus, we introduce **a learnable value function** [21, 46] to guide the action exploration process, directing exploration toward potentially optimal actions. This mechanism anchors explorations within plausible regions while enabling controlled extrapolation, thereby facilitating strategy improvement and further mitigating OOD issues.

Our contributions are summarized as follows:

- We introduce an innovative framework GAVE that leverages DT to optimize auto-bidding strategies, which is designed for seamless adaptability to various real-world scenarios.
- This paper presents three technical innovations: (1) A score-based RTG module with customizable functions for various advertising objectives through differentiable programming; (2) An action exploration mechanism with RTG-based evaluation to ensure stability-preserving updates; (3) A learnable value function to anchor exploration to plausible regions, thus mitigating OOD risks and enabling controlled extrapolation for strategy improvement.
- Experiments on two public datasets, along with results from online deployments, demonstrate the effectiveness of GAVE compared to various state-of-the-art offline bidding baselines. Additionally, by applying the core methods of this framework, we proudly secured first place in the NeurIPS 2024 competition, 'AIGB Track: Learning Auto-Bidding Agents with Generative Models' hosted by Alimama [1].

## 2 Preliminary

In this section, we first illustrate the auto-bidding problem, and then introduce the DT-based decision-making process for modeling.

### 2.1 Auto-Bidding Problem

Consider a sequence of $I$ impression opportunities arriving over a discrete time period $i = 1, ..., I$. Advertisers engage in real-time competition by submitting bids $\{b_i\}_{i=1}^{I}$ for these impressions.

The auction mechanism operates under the following rules: An advertiser wins impression $i$ if his bid $b_i$ exceeds $b_i^-$, the highest competing bid from other participants. The winning advertiser then incurs a cost $c_i$, which is determined by the auction mechanism. Following standard industry practice [2, 8], we adopt the generalized second-price auction mechanism where the winning cost equals the second-highest bid. The advertiser's objective is to maximize the total acquired value through won impressions during the period. This optimization problem can be formally expressed as:

$$\max \sum_{i=1}^{I} x_i v_i \tag{1}$$

where $v_i \in \mathbb{R}^+$ represents the advertiser's private valuation for impression $i$ such as conversion or click-through rate, and $x_i \in \{0, 1\}$ denotes the binary decision variable indicating auction outcome:

$$x_i = \begin{cases} 1 & \text{if } b_i > b_i^- \\ 0 & \text{otherwise} \end{cases} \tag{2}$$

Simultaneously, advertisers must satisfy multiple constraints to ensure efficient campaign management. The fundamental constraint is the total budget limitation:

$$\sum_{i=1}^{I} x_i c_i \leq B \tag{3}$$

where $B \in \mathbb{R}^+$ represents the advertiser's total budget. Other Key Performance Indicator (KPI) constraints, exemplified by Cost Per Acquisition (CPA), can be formulated as follows:

$$\frac{\sum_{i=1}^{I} x_i c_i}{\sum_{i=1}^{I} x_i v_i} \leq C \tag{4}$$

where $C \in \mathbb{R}^+$ denotes the maximum allowable CPA. This ratio quantifies the efficiency of advertising expenditure relative to value creation. Since most other KPI constraints can be modeled similarly, we consider only the CPA constraint for simplicity in this paper. However, unlike budget constraints, which are directly managed by the auction platform, these KPI constraints are generally not strict in practical scenarios. This is because calculating these constraints requires the advertiser's $v_i$ for all bidding impressions, making it possible to determine the true CPA only after the entire bidding process concludes. Nevertheless, we still hope to use them as soft constraints in modeling.

Therefore, the whole bidding process could be expressed as:

$$\max_{b_1,\cdots,b_I} \sum_i x_i v_i$$
$$\text{s.t. } \sum_i x_i c_i \leq B \tag{5}$$
$$\frac{\sum_i x_i c_i}{\sum_i x_i v_i} \leq C$$

Solving this optimization problem presents inherent challenges stemming from both the high cardinality of impressions and the fundamental uncertainty about future auction performance. Previous research [17] reformulates this problem as a Linear Programming problem to yield a simplified optimal bidding strategy:

$$b_i^* = \lambda_0^* v_i - \Sigma_j \lambda_j^* \left( \mathbb{q}_{ij} \left( 1 - \mathbb{1}_{CR_j} \right) - \mathbb{k}_j \mathbb{p}_{ij} \right) \tag{6}$$

where $b_i^*$ denotes the theoretically optimal bid for impression $i$, $\mathbb{q}_{ij}$ can be any performance indicator or constant and $\mathbb{1}_{CR_j}$ is the indicator function of whether constraint $j$ is cost-related. $\mathbb{p}_{ij}$ and $\mathbb{k}_j$ can be treated as expanded expressions of $v_i$ and $C$ in Equation (5) under multiple KPI consdtions. This reformulation transforms the auto-bidding problem into the identification of the optimal $\lambda_0^*$ and $\lambda_j^*$ that satisfy all constraints. By substituting Equation (6) into Equation (5) with $j = 1$, $\mathbb{1}_{CR_j} = 1$, $\mathbb{p}_{ij} = v_i$, $\mathbb{k}_j = C$ and $\mathbb{q}_{ij}$ being any performance indicator or constant, we can obtain:

$$b_i^* = (\lambda_0^* + \lambda_1^* C) v_i = \lambda^* v_i \tag{7}$$

where $\lambda^* = \lambda_0^* + \lambda_1^* C$ serves as the unified bidding parameter. Therefore, many recent studies have sought to address the bidding problem by iteratively identifying the optimal $\lambda^*$ within the bidding process [17, 27, 44]. Additionally, it is worth noting that when solving a bidding problem according to Equation (7), the first condition, i.e., s.t. $\sum_i x_i c_i \leq B$, is always satisfied. This is because the bidding

platform will automatically control $x_i$ when the advertiser's budget is insufficient to ensure that the advertiser does not owe money. However, the second condition is not always satisfied since there's a gap between our predicted $\lambda$ and the optimal $\lambda^*$. A simple solution to this problem is to add a penalty term about the CPA condition to the objective function in Equation (7) in the evaluation stage for model selection [44], which will be discussed later in Section 3.2.

## 2.2 DT-based Auto-bidding

To solve the auto-bidding problem, existing approaches employ either rule-based policies [3, 42] or RL methods [51, 56] for optimization. However, rule-based policies often fail to adapt to the highly dynamic nature of real-world bidding environments [17], and RL approaches [47] rely on state transitions defined by $s_{t+1} = f(s_t, a_t)$, which complicates the modeling of essential temporal dependencies and historical observations inherent in auction ecosystems [57].

Recent advancements in transformer architectures [23, 28] have led to the emergence of DT [6, 52, 63], positioning them as state-of-the-art for sequential decision-making. DTs excel in capturing long-range dependencies, making them ideal for bidding environments where auction outcomes display significant temporal correlations. Building on this framework, we approach auto-bidding as a sequence modeling task [27, 44] under DT settings. The bidding period is divided into discrete time steps, with each step configured under specific environmental settings:

- **state $s_t$**: The state vector $s_t$ encompasses a collection of features that characterizes the bidding conditions at timestep $t$. For advertising scenarios, these features could be the remaining time, unused budget, historical bidding statistics, etc.
- **action $a_t$**: The action $a_t$ denotes the bidding variables that could be iteratively adjusted through the whole bidding period. In this paper, according to Equation (7), the optimal action is $a = \lambda^*$. Therefore, we denote the real action at time step $t$ as:

$$a_t = \lambda_t \tag{8}$$

- **reward $rw_t$**: Suppose there are $N_t$ candidate impressions coming between $t$ and $t + 1$. The reward $rw_t$ could then be defined as:

$$rw_t = \sum_{n=0}^{N_t} x_{n_t} v_{n_t} \tag{9}$$

where $x_{n_t}$ and $v_{n_t}$ are the binary indicator and value of the $n^{th}$ impression at time step $t$.

- **Return-To-Go (RTG) $r_t$**: The RTG value indicates the total amount of rewards to be obtained in the future time steps:

$$r_t = \sum_{t'=t}^{T} rw_{t'} \tag{10}$$

where $T$ is the final time step.

These settings result in the following trajectory representation, which is well-suited for autoregressive training and inference:

$$\tau = (r_1, s_1, a_1, r_2, s_2, a_2, \ldots, r_T, s_T, a_T) \tag{11}$$

## 3 Framework

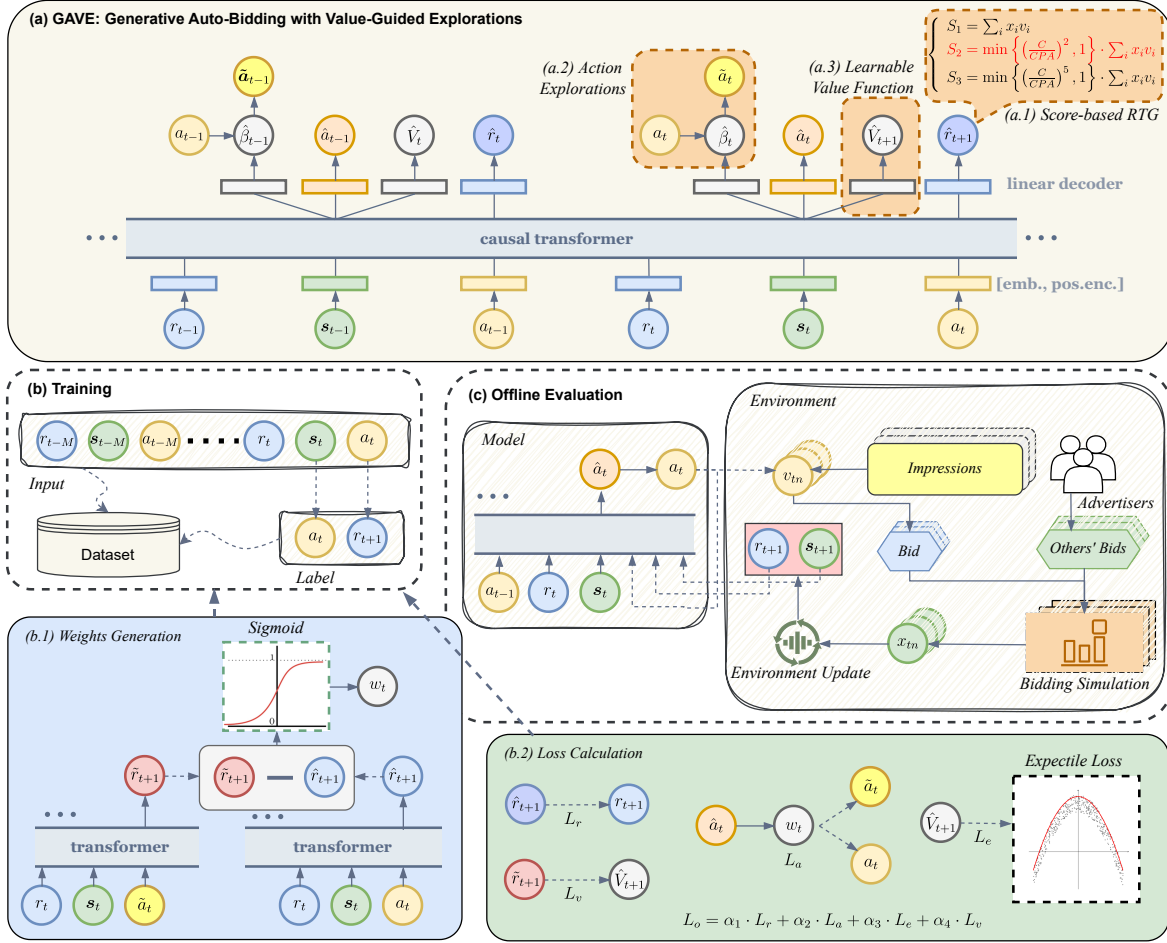Here, we will detail the GAVE's overview and key components.

Figure 1: Overall structure of GAVE.

## 3.1 GAVE Overview

As shown in Figure 1, GAVE adopts a DT architecture, where the pair of RTGs, states, and actions form the input sequence, i.e., $(r_t, s_t, a_t)$ for the timestamp $t$. Unlike conventional DTs, GAVE introduces several key innovations to achieve adaptive optimization, enhance stability, and facilitate strategy improvement. These include a **score-based RTG** (Figure 1 (a.1)) for aligning with diverse advertising objectives, an **action exploration module** (Figure 1 (a.2)) equipped with an RTG-based evaluation mechanism for discovering and evaluating new actions and stabilizing updates, and a **learnable value function** (Figure 1 (a.3)) to steer exploration for strategy improvement while mitigating Out-of-Distribution (OOD) risks. The training of GAVE follows an offline paradigm (Figure 1 (b)), using sequence samples as input to generate predicted labels. For evaluation, a simulated bidding environment (Figure 1 (c)) is employed, where the test model interacts with fixed-policy agents.

Specifically, GAVE employs an adaptive score-based RTG function that can align optimization objectives with varying advertising objectives. During action exploration at time step $t$, in addition to predicting the action $\hat{a}_t$, GAVE predicts a coefficient $\hat{\beta}_t$, a $\hat{V}_{t+1}$ for estimating the learnable value function $V_{t+1}$, and an RTG value $\hat{r}_{t+1}$. The process is formally expressed as follows:

$$\begin{cases} (\hat{\beta}_t, \hat{a}_t, \hat{V}_{t+1}) = GAVE(r_{t-M}, s_{t-M}, a_{t-M}, \ldots, r_t, s_t) \\ \hat{r}_{t+1} = GAVE(r_{t-M}, s_{t-M}, a_{t-M}, \ldots, r_t, s_t, a_t) \\ \tilde{a}_t = \hat{\beta}_t a_t \end{cases} \tag{12}$$

where $M$ is a hyper-parameter, indicating a sequence with $M + 1$ input time steps.

By evaluating the explored action $\tilde{a}_t$ and the action label $a_t$ with an RTG-based evaluation method, GAVE applies a balanced update strategy to reconcile $\tilde{a}_t$ and $a_t$. This ensures a stability-preserving update process. Additionally, the learnable value function $V_{t+1}$ is introduced to direct the model toward potentially optimal strategies while further reducing OOD risks. These innovations collectively enable GAVE to achieve improved performance and robustness.

## 3.2 Score-based RTG

As illustrated in Section 2.1, directly optimizing the cumulative value of won impressions may cause the CPA constraint to significantly exceed its limit. To address this issue, objective functions incorporating penalty terms can be established as evaluation metrics, allowing the degree of emphasis on CPA restrictions to be tailored to specific advertising objectives. This approach facilitates

the evaluation and selection of the optimal model. For example, previous work [44] proposes using a score $S$ to assess the model's actual performance during the testing phase, thereby enabling the selection of higher-performing models. This score integrates a penalty term for CPA constraints to evaluate the overall performance of the bidding model throughout the entire bidding period with $\gamma = 2$:

$$
\begin{cases}
CPA = \frac{\sum_i x_i c_i}{\sum_i x_i v_i} \\
\mathbb{P}(CPA; C) = \min\left\{ \left(\frac{C}{CPA}\right)^\gamma, 1 \right\} \\
S = \mathbb{P}(CPA; C) \cdot \sum_i x_i v_i
\end{cases}
\tag{13}
$$

In this paper, we integrate the constraints directly into the training stage, moving beyond reliance on pre-trained model selection to achieve improved evaluation scores. To align with evaluation metrics for various advertising objectives, we propose employing a constrained score function instead of the unconstrained $\sum_{i=1}^{I} x_i v_i$ for RTG modeling in GAVE, as illustrated in Figure 1 (a.1). For example, based on the evaluation metric defined in Equation (13), the following score-based RTG function can be utilized to synchronize training with evaluation:

$$
\begin{cases}
CPA_t = \frac{\sum_i^{I_t} x_i c_i}{\sum_i^{I_t} x_i v_i} \\
\mathbb{P}(CPA_t; C) = \min\left\{ \left(\frac{C}{CPA_t}\right)^\gamma, 1 \right\} \\
S_t = \mathbb{P}(CPA_t; C) \cdot \sum_i^{I_t} x_i v_i \\
r_t = S_T - S_{t-1}
\end{cases}
\tag{14}
$$

Here, $I_t$ denotes the number of impressions from time step 0 to time step $t$, $S_t$ represents the generalized score function at time step $t$, and $T$ signifies the final time step in a bidding period. By generalizing the score calculation to each time step, the RTG $r_t$ is derived to represent the future score yet to be obtained, guiding the optimization direction of GAVE.

Furthermore, in practical applications, different advertising objectives may exhibit varying degrees of dependence on CPA constraints, resulting in different evaluation metrics. Nonetheless, training and evaluation can remain aligned by generalizing $S$ to each time step (i.e., $S_t$) in a similar way, leading to:

$$
r_t = S_T - S_{t-1}
\tag{15}
$$

This score-based RTG function enhances the flexibility of GAVE, ensuring its applicability across diverse advertising objectives.

## 3.3 Action Explorations

The primary objective of this section is to explore novel actions during training to discover strategies potentially absent from the offline dataset, thereby enabling better model optimization. However, in an offline setting where environment interaction is impossible, learning merely from the fixed dataset may lead to documented behavioral patterns. Conversely, exploring actions beyond the dataset can introduce inherent distribution shifts, potentially leading to behavioral collapse [6, 21]. Moreover, compared to real action labels, the impact of explored actions on model performance can be either beneficial or detrimental, presenting significant challenges in developing a stability-preserving update procedure.

To address these challenges, GAVE introduces a novel action exploration mechanism in conjunction with an RTG-based evaluation method as illustrated in Figure 1 (a.2). This enables GAVE to adaptively adjust both the exploration and update directions of actions by identifying their significance, thereby achieving stability-preserving updates. Specifically, at time step $t$, GAVE predicts a coefficient $\hat{\beta}_t$ with the same dimensionality as $a_t$ to generate a new action $\tilde{a}_t$. This process is formally expressed as:

$$
\begin{cases}
\hat{\beta}_t = \sigma(FC_\beta(DT(r_{t-M}, \mathbf{s}_{t-M}, a_{t-M}, \ldots, r_t, \mathbf{s}_t))) \\
\tilde{a}_t = \hat{\beta}_t a_t
\end{cases}
\tag{16}
$$

where $DT()$ represents the DT backbone, $FC_\beta()$ denotes a fully-connected layer, and $\sigma$ is scaling function. To mitigate OOD issues, the scaling function is defined as:

$$
\sigma(x) = Sigmoid(x) + 0.5
\tag{17}
$$

This formulation constrains $\hat{\beta}_t$ to the interval $(0.5, 1.5)$, ensuring the explored action $\tilde{a}_t$ remains in proximity to the action label $a_t$.

To minimize distribution shift and obtain stability-preserving updates during training, rather than directly utilizing $\tilde{a}_t$ for generating new samples, we employ it as an additional label to balance action updates in conjunction with the original label $a_t$. This approach necessitates estimating the relative significance of $\tilde{a}_t$ and $a_t$ to determine an optimal update direction for the predicted action $\hat{a}_t$. Following reinforcement learning conventions [21, 39, 41], we define the action-value of $a_t$ as $r_{t+1}$ (the RTG at time step $t + 1$), as it represents the cumulative future returns after executing action $a_t$. This enables the design of $w_t$ as illustrated in Figure 1 (b.1) to balance the update direction:

$$
\begin{cases}
\tilde{r}_{t+1} = GAVE(r_{t-M}, \mathbf{s}_{t-M}, a_{t-M}, \ldots, r_t, \mathbf{s}_t, \tilde{a}_t))) \\
w_t = Sigmoid(\alpha_r \cdot (\tilde{r}_{t+1} - \hat{r}_{t+1}))
\end{cases}
\tag{18}
$$

where $\tilde{r}_{t+1}$ and $\hat{r}_{t+1}$ represent the estimated RTG for $\tilde{a}_t$ and $a_t$ respectively. The corresponding loss function for action explorations is defined as:

$$
\begin{cases}
L_r = \frac{1}{M+1} \sum_{t-M}^{t} (\hat{r}_{t+1} - r_{t+1})^2 \\
L_a = \frac{1}{M+1} \sum_{t-M}^{t} ((1 - w_t') \cdot (\hat{a}_t - a_t)^2 + w_t' \cdot (\hat{a}_t - \tilde{a}_t')^2)
\end{cases}
\tag{19}
$$

where $w'$ and $\tilde{a}_t'$ denote $w$ and $\tilde{a}_t$ with frozen gradients. Through $L_r$, GAVE ensures accurate RTG prediction, enabling reliable estimation of the RTG for both $\tilde{a}_t$ and $a_t$. Through $L_a$, GAVE maintains a balanced and stability-preserving updating process between $\tilde{a}_t$ and $a_t$, directing updates toward $\tilde{a}_t$ when it proves superior ($w_t > 0.5$), and toward $a_t$ otherwise to mitigate OOD issues and potential negative impacts from exploration.

## 3.4 Learnable Value Function

While the action exploration mechanism ensures explorations outside the dataset and a stability-preserving update process, randomly generated $\tilde{a}_t$ cannot guarantee improved model performance. To address this limitation, we propose a learnable value function that

---

**Algorithm 1** Optimization algorithm of GAVE

---

**Input**: A training dataset $\mathcal{D} = \left\{ \left( z_j, y_j \right) \right\}_{j=1}^{|\mathcal{D}|}$ with $|\mathcal{D}|$ samples sampled from a bidding environment. $z_j$ is a sequence with $M + 1$ time steps and $y_j$ are the label set.
$z_j = \{ r_{t-M}, s_{t-M}, a_{t-M}, \ldots, r_t, s_t, a_t \}$; $y_j = \{ a_t, r_{t+1} \}$
**Output**: A well-trained model $f$ with parameters $\Phi$

1: Randomly initialize parameters $\Phi$ of the model $f$
2: **for** Step 1,..., Max Step **do**
3:     Sample a training batch $B$ from $\mathcal{D}$
4:     Obtain $\hat{\beta}_t, \hat{a}_t, \hat{V}_{t+1}, \hat{r}_{t+1}$ and $\tilde{a}_t$ with $f(B)$ via Equation (12)
5:     Obtain $\tilde{r}_{t+1}$ and $w_t$ via Equation (18)
6:     Calculate loss $L_r, L_a, L_e$ and $L_v$ based on Equation (19), (21) and (22)
7:     Calculate the overall loss $L_o$ based on Equation (23)
8:     Update $\Phi$ via minimizing the loss $L_o$
9: **end for**
10: return $f$

---

facilitates the discovery of superior actions for strategy improvement, as illustrated in Figure 1 (a.3). Specifically, drawing inspiration from reinforcement learning conventions [21, 39, 41], we propose a sequence-value function $V_{t+1}$ analogous to the optimal state-value function in RL, which represents the upper bound of $r_{t+1}$ as follows:

$$V_{t+1} = \arg\max_{a_t \in \mathbb{A}} r_{t+1} \tag{20}$$

where $\mathbb{A}$ denotes the available action space. Due to the extensive action space and limited real actions within offline datasets, the direct statistical computation of $V_{t+1}$ is infeasible. However, we can learn this value through an expectile regression process with $r_{t+1}$:

$$
\begin{aligned}
L_e &= \frac{1}{M+1} \sum_{t-M}^{t} \left( L_2^\tau (r_{t+1} - \hat{V}_{t+1}) \right) \\
&= \frac{1}{M+1} \sum_{t-M}^{t} \left( |\tau - \mathbb{1}\left( (r_{t+1} - \hat{V}_{t+1}) < 0 \right)| (r_{t+1} - \hat{V}_{t+1})^2 \right)
\end{aligned}
\tag{21}
$$

where $L_2^\tau (y - m(x))$ represents the loss function for predicting the expectile $\tau \in (0, 1)$ of labels $y$ with model $m(x)$ [21]. $\hat{V}_{t+1}$ represents the predicted value of $V_{t+1}$. Following Equation (20), we set $\tau = 0.99$ to learn the upper bound of $r_{t+1}$, effectively estimating $V_{t+1}$.

By estimating $V_{t+1}$ with $\hat{V}_{t+1}$ and utilizing it to guide the updating directions of $\tilde{r}_{t+1}$, GAVE implicitly steers the update direction of the explored $\tilde{a}_t$ toward potentially optimal actions. This process is illustrated in Figure 1 (b.2) and can be formalized as:

$$L_v = \frac{1}{M+1} \sum_{t-M}^{t} (\tilde{r}_{t+1} - \hat{V}'_{t+1})^2 \tag{22}$$

where $\hat{V}'_{t+1}$ represents the gradient-frozen version of $\hat{V}_{t+1}$. Through the application of $L_v$, GAVE implicitly guides the update direction of $\tilde{a}_t$ toward the optimal action by anchoring their RTGs near $\hat{V}_{t+1}$. This approach mitigates OOD risks and enables controlled extrapolation for strategy improvement.

**Table 1: Data statistics.**

| Params | AuctionNet | AuctionNet-Sparse |
|---|---|---|
| Trajectories | 479,376 | 479,376 |
| Delivery Periods | 9,987 | 9,987 |
| Time steps in a trajectory | 48 | 48 |
| State dimension | 16 | 16 |
| Action dimension | 1 | 1 |
| Return-To-Go Dimension | 1 | 1 |
| Action range | [0, 493] | [0, 589] |
| Impression's value range | [0, 1] | [0, 1] |
| CPA range | [6, 12] | [60, 130] |
| Total conversion range | [0, 1512] | [0, 57] |

## 3.5 Optimization Algorithm

Through the mechanisms described above, GAVE implements an offline generative auto-bidding framework incorporating value-guided explorations to enhance strategy learning. The comprehensive loss function is formulated as a weighted combination of the components defined in Equations (19), (21), and (22):

$$L_o = \alpha_1 \cdot L_r + \alpha_2 \cdot L_a + \alpha_3 \cdot L_e + \alpha_4 \cdot L_v \tag{23}$$

where $\{\alpha_1, \alpha_2, \alpha_3, \alpha_4\}$ are hyperparameters controlling the relative contribution of each loss component. The complete optimization procedure of GAVE is detailed in Algorithm 1, with the training process visualized in Figure 1 (b).

During inference, as illustrated in Figure 1 (c), GAVE processes each input sequence to predict $\hat{a}_t = \lambda_t$, which serves as the bid parameter at time step $t$. The bidding price for the $n$-th impression at time step $t$ is then computed according to Equation (7) as $b_{tn} = \lambda_t v_{tn}$, enabling real-time bidding simulation.

## 4 Offline Experiments

In this section, we conduct experiments on two public datasets to investigate the following questions:

- **RQ1:** How does GAVE perform compared to state-of-the-art auto-bidding baselines?
- **RQ2:** Can GAVE adapt to diverse advertising objectives?
- **RQ3:** How effective is the proposed learnable value function in facilitating action exploration?
- **RQ4:** How do the proposed components in GAVE contribute to the final bidding performance?

In the following subsections, we begin by outlining the evaluation settings. Then, we address the relevant questions by concisely analyzing our experimental findings.

## 4.1 Experimental Setup

*4.1.1 Dataset.* Previous auto-bidding research has predominantly relied on proprietary bidding logs for evaluation, with problem formulations often specific to particular scenarios. This heterogeneity in evaluation methodologies has hindered fair and systematic comparisons across different approaches. Recently, Alimama introduced AuctionNet [3] [44], the industry's first standardized large-scale simulated bidding benchmark, enabling comprehensive model evaluation under consistent conditions.

In this study, we utilize two datasets from the AuctionNet framework: (i) **AuctionNet**: The primary dataset containing comprehensive bidding trajectories, and (ii) **AuctionNet-Sparse**: A sparse variant of AuctionNet featuring reduced conversion rates. Both

---

[3] https://github.com/alimama-tech/AuctionNet

**Table 2: Performance comparison. The boldface denotes the highest score. The underline indicates the best result of baselines. "*" indicates the statistically significant improvements (i.e., two-sided t-test with $p < 0.05$) over the best baseline.**

| Dataset | Budget | DiffBid | USCB | CQL | IQL | BCQ | DT | CDT | GAS | GAVE | Improve |
|---------|--------|---------|------|-----|-----|-----|-----|-----|-----|------|---------|
| AuctionNet | 50% | 54 | 86 | 113 | 164 | 190 | 191 | 174 | 193 | **201*** | 4.15% |
| | 75% | 100 | 135 | 139 | 232 | 259 | 265 | 242 | 287 | **296*** | 3.14% |
| | 100% | 152 | 157 | 171 | 281 | 321 | 329 | 326 | 359 | **376*** | 4.74% |
| | 125% | 193 | 220 | 201 | 355 | 379 | 396 | 378 | 409 | **421*** | 2.93% |
| | 150% | 234 | 281 | 238 | 401 | 429 | 450 | 433 | 461 | **467*** | 1.30% |
| AuctionNet-Sparse | 50% | 9.9 | 11.5 | 12.8 | 16.5 | 17.7 | 14.8 | 11.2 | 18.4 | **19.6*** | 6.52% |
| | 75% | 15.4 | 14.9 | 16.7 | 22.1 | 24.6 | 22.9 | 18.0 | 27.5 | **28.3*** | 2.91% |
| | 100% | 19.5 | 17.5 | 22.2 | 30.0 | 31.1 | 29.6 | 31.2 | 36.1 | **37.2*** | 3.05% |
| | 125% | 25.3 | 26.7 | 28.6 | 37.1 | 34.2 | 34.3 | 31.7 | 40.0 | **42.7*** | 6.75% |
| | 150% | 30.8 | 31.3 | 35.8 | 43.1 | 37.9 | 44.5 | 39.1 | 46.5 | **47.4*** | 1.94% |

datasets comprise approximately 500k bidding trajectories collected across 10k distinct delivery periods, each consisting of 48 time steps and interactions derived from millions of impression opportunities. Detailed statistics are presented in Table 1.

*4.1.2 Evaluation Protocol.* Our evaluation methodology follows the AuctionNet benchmark [44] and employs a simulated environment [27] that emulates real-world advertising systems, as illustrated in Figure 1 (c). The evaluation spans a 24-hour delivery period, discretized into 48 uniform time steps, during which the predicted action is utilized for bidding ($\hat{a}_t = a_t$). Within this simulated environment, 48 bidding agents with distinct strategies compete for incoming impression opportunities, with performance measured using Equation (13) with $\gamma = 2$.

To ensure comprehensive evaluation, we employ a round-robin testing strategy: the test model sequentially replaces each of the 48 agents, competing against the remaining agents in each round. The final performance is computed as the average score across all evaluations, providing a robust measure of the model's effectiveness.

*4.1.3 Baselines.* To assess the effectiveness of GAVE, we perform a thorough comparison with multiple baseline approaches:

- **DiffBid** [16]: applies the diffusion frameworks to simulate bidding trajectories and model bidding sequences.
- **USCB** [17]: dynamically adjust bidding parameters for optimal bidding performance in an online RL bidding environment.
- **CQL** [22]: learns a conservative value function to mitigate the overestimation problems in offline RL.
- **IQL** [21]: applies an expectile regression method to enable policy improvement without evaluating the out-of-scope actions.
- **BCQ** [9]: applies a restriction on the action space for a typical offline RL learning process.
- **DT** [6]: employs a transformer architecture for sequential decision-making modeling and utilizes a behavior cloning method to learn the average strategy from the dataset.
- **CDT** [33]: tries to train a constraint satisfaction policy in the offline settings for a balance of safety and task performance.
- **GAS** [27]: tries to model a DT-based offline bidding framework with post-training search by applying Monte Carlo Tree Search (MCTS) in modeling.

*4.1.4 Implementation Details.* Following previous studies [27, 44], we conduct evaluations using varying budget ratios from the original dataset. Performance is measured using the scoring metric:

$$S = \mathbb{P}(CPA; C) \cdot \sum_i x_i v_i \tag{24}$$

as defined in Equation (13) with $\gamma = 2$.

All experiments are conducted on NVIDIA H100 GPUs, utilizing a fixed batch size 128 for a maximum of 400k training steps. The GAVE implementation employs a causal transformer architecture with 8 layers and 16 attention heads. Model parameters are optimized using the AdamW optimizer [36] with a learning rate of $1e^{-5}$. Additional hyper-parameters are determined through a comprehensive grid search to maximize performance. To ensure statistical significance, we conduct 10 independent runs using the optimal parameter configuration and report the average performance metrics.

## 4.2 Overall Performance (RQ1)

We present a comprehensive comparison between GAVE and various baseline approaches across different budget settings, with results summarized in Table 2. Our experimental analysis reveals several key findings:

- GAVE demonstrates superior performance across all budget and dataset configurations, consistently outperforming existing methods. This superiority can be attributed to our novel action exploration method with the value function's guidance, which enables the discovery of novel, potentially optimal actions beyond the offline dataset while maintaining robust training through a stability-preserving update process balancing the exploration benefits and risks.
- Among all baselines, DT-based methods (GAS, DT, and CDT) exhibit superior performance, highlighting the effectiveness of DT structure in capturing temporal dependencies and facilitating sequential decision-making in bidding scenarios. Notably, GAS achieves better results compared to DT and CDT, validating the efficacy of its MCTS implementation in strategy optimization.
- DiffBid does not perform well on the datasets, probably due to the long sequence and highly dynamic environment posing extra challenges for DiffBid in accurately predicting trajectories and learning from its reverse process.

## 4.3 Alignment Analysis (RQ2)

As discussed in Section 3.2, advertising objectives may necessitate different evaluation metrics. To address this, GAVE employs an adaptive score-based RTG modeling that accommodates various optimization objectives, thereby aligning training objectives with

**Table 3: Alignment analysis on AuctionNet-Sparse with 100% budget. "Train" denotes applying the score to RTG modeling, while "Eval" denotes its utilization as the evaluation metric.**

| Train \ Eval | $S_1$ | $S_2$ | $S_3$ |
|---|---|---|---|
| $S_1$ | 41.4 | 33.0 | 23.6 |
| $S_2$ | 39.9 | 37.2 | 33.3 |
| $S_3$ | 39.1 | 36.8 | 33.5 |

evaluation metrics as illustrated in Equation (15). In this section, we explore the performance of GAVE under different RTG and evaluation metric configurations to answer RQ2. Specifically, we consider three evaluation metrics defined as follows:

$$\begin{cases} S_1 = \sum_i x_i v_i, \\ S_2 = \min\left\{\left(\frac{C}{CPA}\right)^2, 1\right\} \cdot \sum_i x_i v_i, \\ S_3 = \min\left\{\left(\frac{C}{CPA}\right)^5, 1\right\} \cdot \sum_i x_i v_i, \end{cases} \quad (25)$$

where $S_1$ only considers the total impression values obtained and represents the business scenarios with elastic restrictions on CPA conditions. $S_2$ is the optimization goal and evaluation score of this paper. It adds a punishment on CPA constraints. $S_3$ further enhances the penalty coefficient of CPA to represent the business scenarios with strict restrictions on CPA conditions. These metrics are utilized either for RTG modeling during training or as evaluation criteria during testing. The results are presented in Table 3.

From Table 3, we observe that GAVE consistently achieves the highest performance when the training RTG corresponds to the same function used as the evaluation metric. This finding underscores the importance of aligning training objectives with specific evaluation metrics through our score-based RTG approach.

### 4.4 Parameter Analysis (RQ3)

To address RQ3, we conduct a parameter analysis of the weight $w_t$, as depicted in Figure 1 (b.1), to elucidate the distinction between $\tilde{a}_t$ and $a_t$ during the training process. Specifically, Figure 2 visualizes the average overall loss $L_o$ and the weight $w_t$ across training steps, allowing us to monitor the disparity between $\tilde{r}_{t+1}$ and $\hat{r}_{t+1}$. A larger $w_t$ signifies a greater influence of $\tilde{r}_{t+1}$ over $\hat{r}_{t+1}$, thereby demonstrating that $\tilde{a}_t$ is superior to $a_t$. This result underscores the effectiveness of the value function in guiding action exploration.

From Figure 2, it is evident that as training progresses, the parameter $w_t$ increases from approximately 0.5 to a stable position above 0.5. The stable position is influenced by both the dataset distribution and the model's hyper-parameters. This trend confirms the efficacy of the learnable value function in directing action exploration. With the guidance of the value function, the model consistently explores actions $\tilde{a}$ with higher RTG values $\tilde{r}_{t+1}$ near the estimated optimal value $\hat{V}_{t+1}$. This approach facilitates learning potential optimal strategies while mitigating OOD issues.

### 4.5 Ablation Study (RQ4)

To further elucidate the contribution of each module within GAVE for answering RQ4, we conduct an ablation study by evaluating the following modified versions of GAVE:

- **GAVE-V**: excludes the learnable value function described in Section 3.4. In this configuration, the loss functions $L_v$ and $L_e$ are
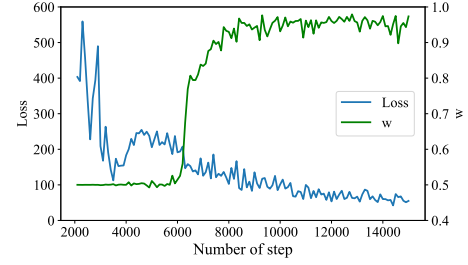


**Figure 2: Parameter analysis of $w$ on AuctionNet.**



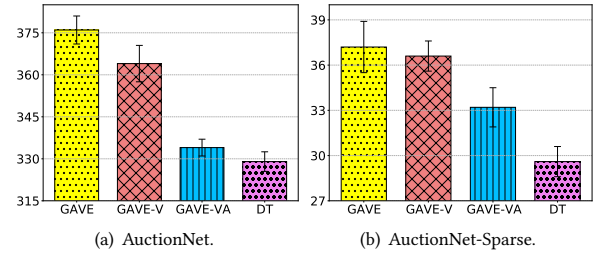(a) AuctionNet.  (b) AuctionNet-Sparse.

**Figure 3: Ablation study with 100% budget.**

replaced with the following update rule to ensure that the explored actions generally surpass the original labels by enhancing their RTG values $\tilde{r}_{t+1}$:

$$L_w = 1 - \text{Sigmoid}(\alpha_r \cdot (\tilde{r}_{t+1} - \hat{r}'_{t+1})) \quad (26)$$

where $\hat{r}'_{t+1}$ is a gradient-frozen version of $\hat{r}_{t+1}$. However, without the value function, the update direction of $\tilde{r}_{t+1}$ becomes unbounded, leading to OOD issues and suboptimal performance.

- **GAVE-VA**: omits both the value function from Section 3.4 and the action exploration mechanism detailed in Section 3.3.
- **DT**: removes all designed modules related to GAVE, including those described in Sections 3.4, 3.3, and Section 3.2. Consequently, this configuration aligns with the pure DT framework [6] with $S = \sum_i x_i v_i$ for RTG modeling.

Figure 3 presents the evaluation results. The findings indicate that: (i) Aligning the optimization objectives with the evaluation metrics using the score-based RTG modeling allows GAVE-VA to outperform DT, demonstrating the importance of objective alignment in training; (ii) Incorporating the action exploration mechanism and RTG-based evaluation in GAVE-V enables the model to discover potential strategies beyond the offline dataset and evaluate their significance for a stability-preserving update process, thereby achieving better performance than GAVE-VA; and (iii) Fully integrating the value function to guide action explorations within GAVE leverages potential optimal strategies, further alleviating OOD issues and enhancing overall performance.

## 5 Online Application

We evaluate GAVE's effectiveness through A/B tests in two industrial live bidding scenarios: Nobid [4] (maximizing conversions within daily budget) and Costcap (maximizing conversions with CPA/ROI limits). The experimental settings are detailed below.

- **State**: 20-step sequence with features including budget, CPA limit, predictions, traffic/cost speeds, time-phased budget, remaining time, and window-averaged bid coefficient.

---

[4]https://support.google.com/google-ads/answer/7381968?hl=en

**Table 4: Online A/B test.**

|  | Cost | Conversion | Target cost | CPA valid ratio |
|---|---|---|---|---|
| Nobid | +0.8% | +8.0% | +3.2% | / |
| Costcap | +2.0% | +3.6% | +2.2% | +1.9% |

- **Action**: To stabilize bidding results, the bid coefficient $\lambda$ is determined based on a windowed-average of the preceding two hours containing $E$ time steps, $\lambda_t = a_t + \frac{1}{|E|} \sum_{t'=t-E}^{t-1} \lambda_{t'}$, where $a_t$ is GAVE's output action at time step $t$.
- **Return-To-Go (RTG)**: Given the sparsity of real conversions, we utilize the expected total conversions, $\sum_i pcvr_i$ during training, where $pcvr_i$ is the predicted conversion rate for winning traffic $i$. During inference, the RTG for the entire sequence is set to the campaign's total expected conversions from the previous day.

We compare GAVE with the offline reinforcement learning algorithm IQL [21], currently in production. Evaluation metrics include cost, conversions, target cost, and CPA valid ratio, with bidding strategies focused on maximizing conversions within budget and CPA constraints. To account for varying campaign targets, target cost serves as a value-weighted conversion measure. For Costcap campaigns, the conversion value equals the CPA limit, while Nobid campaigns use the average real CPA from total traffic. A Costcap campaign is CPA valid if its CPA remains below the limit, assessed solely for Costcap campaigns. Our five-day online A/B testing allocates 25% of each campaign's budget and traffic to the baseline bidding model and GAVE, with results summarized in Table 4.

For both Nobid and Costcap, GAVE improves cost and conversions. In Nobid campaigns, GAVE achieves a 0.8% increase in cost, 8.0% in conversions, and 3.2% in target cost. For Costcap campaigns, advertising revenue and advertiser value rise alongside a significant improvement in CPA validity, with +2.0% cost, +3.6% conversions, +2.2% target cost, and +1.9% valid CPA ratio.

## 6 Related Works

This section provides a brief review of relevant research topics, i.e., offline reinforcement learning, and auto-bidding.

### 6.1 Offline Reinforcement Learning and Decision Transformers

Reinforcement Learning (RL) trains decision-making agents through environment interactions, evolving from foundational works [37, 38] to advanced methods like policy gradient [45], deep Q-learning [39], and deterministic policy optimization [43]. While effective, their reliance on frequent online interactions poses risks and costs in real-world applications [19]. Offline RL addresses this by learning policies from static datasets, with methods like BCQ [9], CQL [22], and IQL [21] offering robust solutions for stable continuous control, mitigating value overestimation, and reducing distributional shifts. However, their dependence on Markov Decision Processes (MDPs) limits access to prior observations and modeling long-range dependencies, critical in sequential tasks with strong temporal patterns. Decision Transformers (DT) [6] overcome these limitations by reframing RL as sequential modeling, leveraging transformer architectures to capture historical patterns and long-term dependencies, achieving state-of-the-art offline RL performance. Extensions like CDT [33] further enable zero-shot constraint adaptation, balancing

safety and performance without online fine-tuning. Building on DT, we propose an auto-bidding framework that aligns with DT's strengths by modeling bid adjustments as trajectory-based sequence generation, effectively capturing intricate temporal correlations in high-stakes, data-sensitive environments.

### 6.2 Auto-Bidding at Online Advertising Platforms

Auto-bidding plays a vital role in managing large-scale ad auctions by automatically optimizing bids per impression to meet advertisers' goals. Early methods like PID [7] and OnlineLP [55] focus on rule-based approaches, using feedback loops and stochastic programming to address budget pacing and bid optimization, though they depend on simplified assumptions. To tackle the complexity of modern advertising ecosystems, RL-based solutions such as RLB [5], USCB [17], MAAB [50], and SORL [40] enable adaptive decision-making with capabilities for handling high-dimensional states and multi-agent coordination. However, due to the risks of real-time bidding, offline RL methods like BCQ [9], CQL [22], and IQL [21] have gained prominence for leveraging historical data. These approaches, while effective under MDP frameworks, struggle with modeling long-range temporal dependencies crucial for sequential bid optimization. Generative sequential modeling methods such as DiffBid [16] and GAS [27] address these challenges using diffusion models and transformers with MCTS for improved bid trajectory generation. In this work, we propose a score-based RTG, action exploration mechanisms, and a learnable value function framework to align optimization objectives with evaluation metrics, enhance action exploration, and learn optimal strategies for improving bidding performance with Decision Transformer.

## 7 Conclusion

In this work, we propose GAVE to enhance DTs for offline generative auto-bidding through value-guided explorations. To accommodate complex advertising objectives, we design a customizable score-based RTG mechanism, enabling adaptive modeling of diverse optimization objectives to align with different evaluation metrics. Moreover, we integrate an action exploration mechanism with an RTG-based evaluation method to explore actions outside the offline dataset while ensuring a stability-preserving update process. To further guide the exploration and mitigate OOD risks, we employ a learnable value function to anchor RTG updates to distributionally plausible regions while allowing controlled extrapolation for strategy improvement. Extensive experiments, online deployments and NeurIPS competition [1] demonstrate the effectiveness of our GAVE framework in enhancing the adaptability and performance of auto-bidding strategies, providing a versatile solution for optimizing digital advertising campaigns in dynamic environments.

# References

[1] Gagan Aggarwal, Ashwinkumar Badanidiyuru, Santiago R Balseiro, Kshipra Bhawalkar, Yuan Deng, Zhe Feng, Gagan Goel, Christopher Liaw, Haihao Lu, Mohammad Mahdian, et al. 2024. Auto-bidding and auctions in online advertising: A survey. *ACM SIGecom Exchanges* (2024), 159–183.

[2] Gagan Aggarwal, Shan Muthukrishnan, Dávid Pál, and Martin Pál. 2009. General auction mechanism for search advertising. In *Proc. of WWW*. 241–250.

[3] Raju Balakrishnan and Rushi P Bhatt. 2014. Real-time bid optimization for group-buying ads. *ACM Transactions on Intelligent Systems and Technology (TIST)* (2014), 1–21.

[4] Nikolay Borissov, Dirk Neumann, and Christof Weinhardt. 2010. Automated bidding in computational markets: an application in market-based allocation of computing services. *Autonomous Agents and Multi-Agent Systems* (2010), 115–142.

[5] Han Cai, Kan Ren, Weinan Zhang, Kleanthis Malialis, Jun Wang, Yong Yu, and Defeng Guo. 2017. Real-time bidding by reinforcement learning in display advertising. In *Proc. of WSDM*. 661–670.

[6] 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. *Proc. of NeurIPS* (2021), 15084–15097.

[7] Ye Chen, Pavel Berkhin, Bo Anderson, and Nikhil R Devanur. 2011. Real-time bidding algorithms for performance-based display ad allocation. In *Proc. of KDD*. 1307–1315.

[8] Benjamin Edelman, Michael Ostrovsky, and Michael Schwarz. 2007. Internet advertising and the generalized second-price auction: Selling billions of dollars worth of keywords. *American economic review* (2007), 242–259.

[9] Scott Fujimoto, David Meger, and Doina Precup. 2019. Off-policy deep reinforcement learning without exploration. In *Proc. of ICML*. 2052–2062.

[10] Jingtong Gao, Bo Chen, Xiangyu Zhao, Weiwen Liu, Xiangyang Li, Yichao Wang, Wanyu Wang, Huifeng Guo, and Ruiming Tang. [n. d.]. LLM4Rerank: LLM-based Auto-Reranking Framework for Recommendations. In *Proc. of WWW*.

[11] Jingtong Gao, Bo Chen, Menghui Zhu, Xiangyu Zhao, Yuhao Wang, Yichao Wang, Huifeng Guo, and Ruiming Tang. 2024. HierRec: Scenario-Aware Hierarchical Modeling for Multi-scenario Recommendations. In *Proc. of CIKM*. 653–662.

[12] Jingtong Gao, Zhaocheng Du, Xiaopeng Li, Yichao Wang, Xiangyang Li, Huifeng Guo, Ruiming Tang, and Xiangyu Zhao. 2025. SampleLLM: Optimizing Tabular Data Synthesis in Recommendations. *arXiv preprint arXiv:2501.16125* (2025).

[13] Jingtong Gao, Xiangyu Zhao, Bo Chen, Fan Yan, Huifeng Guo, and Ruiming Tang. 2023. AutoTransfer: Instance transfer for cross-domain recommendations. In *Proc. of SIGIR*. 1478–1487.

[14] Jingtong Gao, Xiangyu Zhao, Muyang Li, Minghao Zhao, Runze Wu, Ruocheng Guo, Yiding Liu, and Dawei Yin. 2024. SMLP4Rec: an Efficient all-MLP architecture for sequential recommendations. *ACM Transactions on Information Systems* (2024), 1–23.

[15] Chaojie Guo, Russell G Thompson, Greg Foliente, and Xiaoshuai Peng. 2021. Reinforcement learning enabled dynamic bidding strategy for instant delivery trading. *Computers & Industrial Engineering* (2021), 107596.

[16] Jiayan Guo, Yusen Huo, Zhilin Zhang, Tianyu Wang, Chuan Yu, Jian Xu, Bo Zheng, and Yan Zhang. 2024. AIGB: Generative Auto-bidding via Conditional Diffusion Modeling. In *Proc. of KDD*. 5038–5049.

[17] Yue He, Xiujun Chen, Di Wu, Junwei Pan, Qing Tan, Chuan Yu, Jian Xu, and Xiaoqiang Zhu. 2021. A unified solution to constrained bidding in online display advertising. In *Proc. of KDD*. 2993–3001.

[18] Arti Jha, Harshit Jain, Parikshit Sharma, Yashvardhan Sharma, and Kamlesh Tiwari. 2024. Optimizing Real-Time Bidding Strategies: An Experimental Analysis of Reinforcement Learning and Machine Learning Techniques. *Procedia Computer Science* (2024), 2017–2026.

[19] 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).

[20] Ilya Kostrikov, Rob Fergus, Jonathan Tompson, and Ofir Nachum. 2021. Offline reinforcement learning with fisher divergence critic regularization. In *Proc. of ICML*. 5774–5783.

[21] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. [n. d.]. Offline Reinforcement Learning with Implicit Q-Learning. In *Proc. of ICLR*.

[22] Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. 2020. Conservative q-learning for offline reinforcement learning. *Proc. of NeurIPS* (2020), 1179–1191.

[23] Pranjal Kumar. 2024. Large language models (LLMs): survey, technical frameworks, and future challenges. *Artificial Intelligence Review* (2024), 260.

[24] Chengxi Li, Yejing Wang, Qidong Liu, Xiangyu Zhao, Wanyu Wang, Yiqi Wang, Lixin Zou, Wenqi Fan, and Qing Li. 2023. STRec: Sparse transformer for sequential recommendations. In *Proc. of RecSys*. 101–111.

[25] Xinhang Li, Zhaopeng Qiu, Xiangyu Zhao, Zihao Wang, Yong Zhang, Chunxiao Xing, and Xian Wu. 2022. Gromov-wasserstein guided representation learning for cross-domain recommendation. In *Proc. of CIKM*. 1199–1208.

[26] Xiaopeng Li, Fan Yan, Xiangyu Zhao, Yichao Wang, Bo Chen, Huifeng Guo, and Ruiming Tang. 2023. Hamur: Hyper adapter for multi-domain recommendation. In *Proc. of CIKM*. 1268–1277.

[27] Yewen Li, Shuai Mao, Jingtong Gao, Nan Jiang, Yunjian Xu, Qingpeng Cai, Fei Pan, Peng Jiang, and Bo An. 2024. GAS: Generative Auto-bidding with Post-training Search. *arXiv preprint arXiv:2412.17018* (2024).

[28] Tianyang Lin, Yuxin Wang, Xiangyang Liu, and Xipeng Qiu. 2022. A survey of transformers. *AI open* (2022), 111–132.

[29] Weilin Lin, Xiangyu Zhao, Yejing Wang, Yuanshao Zhu, and Wanyu Wang. 2023. Autodenoise: Automatic data instance denoising for recommendations. In *Proc. of WWW*. 1003–1011.

[30] Haochen Liu, Da Tang, Ji Yang, Xiangyu Zhao, Hui Liu, Jiliang Tang, and Youlong Cheng. 2022. Rating distribution calibration for selection bias mitigation in recommendations. In *Proc. of WWW*. 2048–2057.

[31] Haochen Liu, Xiangyu Zhao, Chong Wang, Xiaobing Liu, and Jiliang Tang. 2020. Automated embedding size search in deep recommender systems. In *Proc. of SIGIR*. 2307–2316.

[32] Shuchang Liu, Qingpeng Cai, Bowen Sun, Yuhao Wang, Ji Jiang, Dong Zheng, Peng Jiang, Kun Gai, Xiangyu Zhao, and Yongfeng Zhang. 2023. Exploration and regularization of the latent action space in recommendation. In *Proc. of WWW*. 833–844.

[33] Zuxin Liu, Zijian Guo, Yihang Yao, Zhepeng Cen, Wenhao Yu, Tingnan Zhang, and Ding Zhao. 2023. Constrained decision transformer for offline safe reinforcement learning. In *Proc. of ICML*. 21611–21630.

[34] Ziru Liu, Shuchang Liu, Zijian Zhang, Qingpeng Cai, Xiangyu Zhao, Kesen Zhao, Lantao Hu, Peng Jiang, and Kun Gai. 2024. Sequential recommendation for optimizing both immediate feedback and long-term retention. In *Proc. of SIGIR*. 1872–1882.

[35] 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 *Proc. of WWW*. 1273–1282.

[36] I Loshchilov. 2017. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101* (2017).

[37] Donald Michie and Roger A Chambers. 1968. BOXES: An experiment in adaptive control. *Machine intelligence* (1968), 137–152.

[38] Marvin Minsky. 1961. Steps toward artificial intelligence. *Proceedings of the IRE* (1961), 8–30.

[39] Volodymyr Mnih. 2013. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602* (2013).

[40] Zhiyu Mou, Yusen Huo, Rongquan Bai, Mingzhou Xie, Chuan Yu, Jian Xu, and Bo Zheng. 2022. Sustainable online reinforcement learning for auto-bidding. *Proc. of NeurIPS* (2022), 2651–2663.

[41] Ofir Nachum, Mohammad Norouzi, Kelvin Xu, and Dale Schuurmans. 2017. Bridging the gap between value and policy based reinforcement learning. *Proc. of NeurIPS* (2017).

[42] Amin Sayedi. 2018. Real-time bidding in online display advertising. *Marketing Science* (2018), 553–568.

[43] David Silver, Guy Lever, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. 2014. Deterministic policy gradient algorithms. In *Proc. of ICML*. 387–395.

[44] Kefan Su, Yusen Huo, Zhilin Zhang, Shuai Dou, Chuan Yu, Jian Xu, Zongqing Lu, and Bo Zheng. 2024. AuctionNet: A Novel Benchmark for Decision-Making in Large-Scale Games. In *Proc. of NeurIPS*.

[45] Richard S Sutton, David McAllester, Satinder Singh, and Yishay Mansour. 1999. Policy gradient methods for reinforcement learning with function approximation. *Proc. of NeurIPS* (1999).

[46] Csaba Szepesvári and Michael L Littman. 1999. A unified analysis of value-function-based reinforcement-learning algorithms. *Neural computation* (1999), 2017–2060.

[47] X Wang, S Wang, X Liang, D Zhao, J Huang, X Xu, B Dai, and Q Miao. 2022. Deep Reinforcement Learning: A Survey. *IEEE Transactions on Neural Networks and Learning Systems* (2022).

[48] Yuhao Wang, Ha Tsz Lam, Yi Wong, Ziru Liu, Xiangyu Zhao, Yichao Wang, Bo Chen, Huifeng Guo, and Ruiming Tang. 2023. Multi-task deep recommender systems: A survey. *arXiv preprint arXiv:2302.03525* (2023).

[49] 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 *Proc. of SIGIR*. 1498–1507.

[50] Chao Wen, Miao Xu, Zhilin Zhang, Zhenzhe Zheng, Yuhui Wang, Xiangyu Liu, Yu Rong, Dong Xie, Xiaoyang Tan, Chuan Yu, et al. 2022. A cooperative-competitive multi-agent framework for auto-bidding in online advertising. In *Proc. of WSDM*. 1129–1139.

[51] Di Wu, Xiujun Chen, Xun Yang, Hao Wang, Qing Tan, Xiaoxun Zhang, Jian Xu, and Kun Gai. 2018. Budget constrained bidding by model-free reinforcement learning in display advertising. In *Proc. of CIKM*. 1443–1451.

[52] Yueh-Hua Wu, Xiaolong Wang, and Masashi Hamaya. 2024. Elastic decision transformer. *Proc. of NeurIPS* (2024).

[53] Haoran Xu, Li Jiang, Jianxiong Li, Zhuoran Yang, Zhaoran Wang, Victor Wai Kin Chan, and Xianyuan Zhan. [n. d.]. Offline RL with No OOD Actions: In-Sample Learning via Implicit Value Regularization. In *Proc. of ICLR*.

[54] Yujian Ye, Dawei Qiu, Mingyang Sun, Dimitrios Papadaskalopoulos, and Goran Strbac. 2019. Deep reinforcement learning for strategic bidding in electricity markets. *IEEE Transactions on Smart Grid* (2019), 1343–1355.

[55] Hao Yu, Michael Neely, and Xiaohan Wei. 2017. Online convex optimization with stochastic constraints. *Proc. of NeurIPS* (2017).

[56] Congde Yuan, Mengzhuo Guo, Chaoneng Xiang, Shuangyang Wang, Guoqing Song, and Qingpeng Zhang. 2022. An actor-critic reinforcement learning model for optimal bidding in online display advertising. In *Proc. of CIKM*. 3604–3613.

[57] Shuai Yuan, Jun Wang, and Xiaoxue Zhao. 2013. Real-time bidding for online advertising: measurement and analysis. In *Proceedings of the seventh international workshop on data mining for online advertising*. 1–8.

[58] Jun Zhao, Guang Qiu, Ziyu Guan, Wei Zhao, and Xiaofei He. 2018. Deep reinforcement learning for sponsored search real-time bidding. In *Proc. of KDD*. 1021–1030.

[59] 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 *Proc. of AAAI*. 750–758.

[60] 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), 1–15.

[61] Xiangyu Zhao, Liang Zhang, Zhuoye Ding, Long Xia, Jiliang Tang, and Dawei Yin. 2018. Recommendations with negative feedback via pairwise deep reinforcement learning. In *Proc. of KDD*. 1040–1048.

[62] Xiangyu Zhao, Xudong Zheng, Xiwang Yang, Xiaobing Liu, and Jiliang Tang. 2020. Jointly learning to recommend and advertise. In *Proc. of KDD*. 3319–3327.

[63] Qinqing Zheng, Amy Zhang, and Aditya Grover. 2022. Online decision transformer. In *Proc. of ICML*. 27042–27059.