

GENERATIVE AUTO-BIDDING IN LARGE-SCALE AUCTIONS VIA DIFFUSION COMPLETER-ALIGNER

Anonymous authors

Paper under double-blind review

ABSTRACT

[Bid optimization strategy in auto-bidding](#) is central to computational advertising, achieving notable commercial success by optimizing advertisers' bids within constraints. Recently, generative models have revolutionizing auto-bidding by directly learning a policy from large-scale datasets. Among them, the diffuser is superior in tackling sparse-reward challenges, along with its trajectory stitching and explainability capabilities, making it well-suited for industrial auto-bidding. However, its performance could be limited by generation uncertainty, particularly regarding generations' dynamic illegitimacy and preference misalignment, which can lead to suboptimal bids and further cause poor performance when competing with other advertisers in highly competitive auctions. To address it, we propose a **Causal auto-Bidding** method based on a **Diffusion completer-aligner** framework, termed **CBD**. Firstly, we conduct a theoretical analysis and propose a completer to augment the training process with an extra random variable t for enhancing the dynamic legitimacy between adjacent states. Then, we employ a trajectory-level return model as an aligner to refine the generated trajectories in inference, aligning more closely with advertisers' objectives. Experiments across diverse settings demonstrate that our approach not only achieves superior performance on large-scale auto-bidding benchmarks, such as a 29.9% improvement of conversion value in the challenging sparse-reward setting, but also delivers significant improvements on an online advertising platform, including a 2.0% increase in target cost.

1 INTRODUCTION

The rapid digital transformation of commerce has significantly expanded the reach of online ad platforms, making auto-bidding essential by automatically determining bid prices on behalf of advertisers in large-scale auctions (Ren et al., 2017). [Bid optimization strategy in auto-bidding is to maximize achieved impression value for an advertiser while adhering to some economic constraints](#), such as cost-per-conversion (CPC), throughout an entire bidding period (Guo et al., 2024). Previous auto-bidding strategies have evolved from simple rule-based methods (Zhu et al., 2017; Yang et al., 2019) to the advanced methods based on reinforcement learning (RL) (He et al., 2021; Cai et al., 2017; Sutton, 2018), continuously pushing the boundaries of performance. With the vast amount of bidding logs available from the online system, it is getting more promising to leverage generative models to revolutionize auto-bidding by directly learning a strategy from a large-scale offline dataset that can handle the complex advertising environment while adapting to advertisers' diverse preferences.

In existing literature, two major paradigms employ generative models for decision-making tasks based on offline datasets: Decision Transformer (DT) (Chen et al., 2021; Liu et al., 2023) and Diffuser (Ajay et al., 2023; Janner et al., 2022a). DT learns to autoregressively generate bid actions upon receiving the return-to-go and observed states and actions, which relies on the single-step reward design. However, this approach is not suitable for auto-bidding because conversions, *i.e.*, the results of bidding actions, are typically very sparse in online ad environments such as in-app purchases (IAP) (Guo et al., 2024; Su et al., 2024), making it difficult to design an accurate single-step reward (Zhu et al., 2023). Diffuser offers a promising alternative by generating a trajectory of states based on accumulated rewards over multiple steps as a condition, then actions are executed according to the generation (Ajay et al., 2023). This accumulation approach significantly alleviates the single-step sparse reward issue, making it more suitable for realistic auto-bidding. Moreover, the diffuser has advantages in trajectory stitching ability and explainability (Ajay et al., 2023). DT and Diffuser were

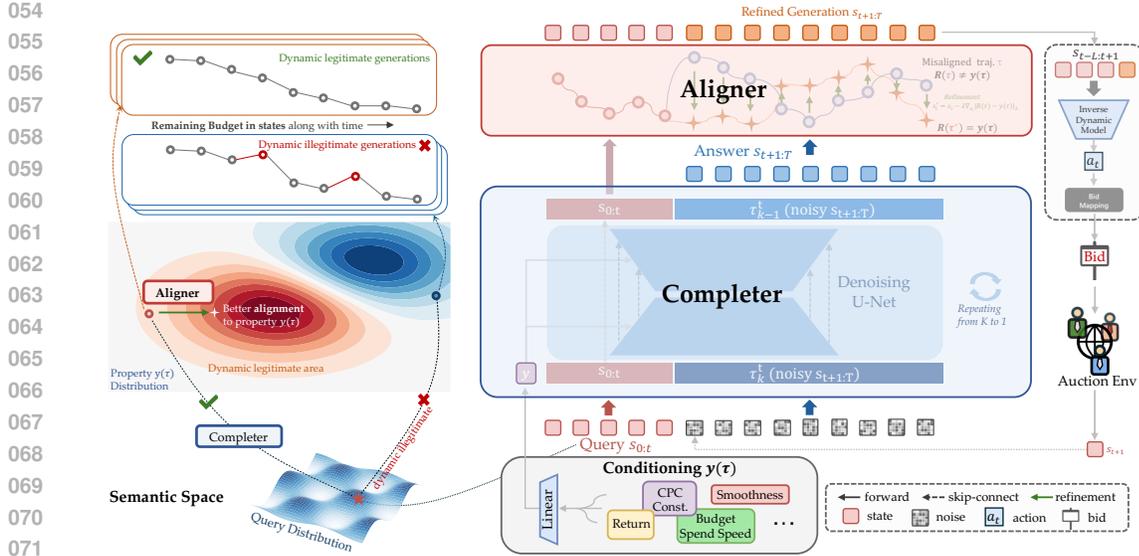


Figure 1: Overview of our CBD method: At each timestep t , a completer first receives a “Query” (observed states $s_{0:t}$ reflecting ad impression features and auction status) and outputs an “Answer” (generated future states $s_{t+1:T}$) within the dynamic legitimate area; The aligner then refines the generation to achieve closer alignment with the expected trajectory property; Finally, a bid is determined based on the refined trajectory and sent to the auction to compete with other advertisers.

initially developed to address offline RL tasks, aiming to stitch suboptimal sub-trajectories in offline datasets into an optimal trajectory. DT achieves this through dynamic programming, *i.e.*, generating actions and interacting with the environment to get the next state. However, this method has proven inefficient due to error accumulation in long-term decision-making compared to Diffuser, which, on the other hand, can generate an optimal trajectory by stitching via its powerful policy representation ability. Additionally, Diffuser is more explainable by revealing its planning process first and then executing strictly according to it, while DT operates as a black-box by directly outputting actions. Explainability is crucial for industrial auto-bidding services, as advertisers may discontinue the service if they observe abnormal and unexplainable behavior from the auto-bidding model.

However, we found that directly applying diffusers could result in limited performance in large-scale auctions with numerous advertisers and millions of impression opportunities (Su et al., 2024). This arises from the inherent uncertainty issue in generative models, which can occur in models ranging from simple generative models like VAEs (Kingma and Welling, 2014) to the most advanced LLMs based on auto-regressive transformers, where this phenomenon is specifically referred to as *hallucination* (Huang et al., 2025). Specifically, the generated next state may not be achievable from the current state, leading to a dynamic legitimacy issue (Ni et al., 2023). Additionally, as it is hard to train a perfect model covering all modes, the generations of diffusers can sometimes be misaligned with the advertiser’s preferences (Li et al., 2024). The situation worsens in large-scale auctions, where the tolerance for bids is minimal due to high competitiveness among advertisers. Suboptimal bids caused by these issues can lead to significant loss of impression opportunities or budget waste. Previous works have tried to address these issues in offline RL tasks, typically involving single-agent, non-competitive tasks like locomotion control (Fu et al., 2020). AdaptDiffuser (Liang et al., 2023) and MetaDiffuser (Ni et al., 2023) use a forward dynamic model (FDM) for enhancing legitimacy, but a poor base model could make FDM useless (Ajay et al., 2023). DiffuserLite (Dong et al., 2024a) revises the training process as generative interpolation by training multiple diffusers, which is inefficient and could limit overall planning. Adopting causal network architectures for training a diffuser, such as DiT (Peebles and Xie, 2023) with causal attention (Yang et al., 2021), is closely related to addressing this issue, but also raises concerns about accumulated errors.

To alleviate the generation uncertainty issue of diffusers in auto-bidding for large-scale competitive auctions, we conduct an in-depth analysis of the diffuser’s training and inference processes within the context of auto-bidding. Our theoretical findings indicate that directly applying diffusers can lead to unnecessary distribution mismatch errors, as defined in Eq. 21, due to a discrepancy between

the inputs used during the training and inference stages. Therefore, we begin by augmenting the denoising process during the training stage by introducing an additional random variable t . This reformulates model training as learning a completer, thereby mitigating the discrepancy that could lead to illegitimacy. Drawing an analogy to LLM, this process involves observing random t -length historical sequences as a “query” and then completing the remaining sequence as an “answer”, thus enhancing the dynamic legitimacy between randomly adjacent states. However, due to the model’s limited capabilities and the constantly evolving preferences of advertisers, the generated trajectory may still be suboptimal, leading to potential misalignment issues. Thus, during inference, we utilize a trajectory-level return model as an efficient aligner to adaptively refine the final generated trajectories, ensuring they align more closely with advertisers’ preferences. This approach also highlights the flexibility and explainability of using a planning-based method with diffusers. Combining them leads to our Causal auto-Bidding method based on a Diffusion completer-aligner, termed CBD. To summarize, our contributions are as follows:

- We highlight the necessity of adopting diffusers in the context of industrial auto-bidding and theoretically analyze the factors hindering their development.
- We propose a novel CBD method to improve the diffuser’s auto-bidding performance by alleviating the impact of generation uncertainty in both training and inference.
- Experiments show that our method excels in large-scale auto-bidding benchmarks, like achieving a notable 29.9% improvement in the challenging sparse-reward scenario similar to industrial reality. Given the sequential decision-making nature of our method, it could also provide insights and advantages for general offline RL tasks. More importantly, we successfully deploy our CBD method to an online advertising platform and achieve a significant improvement, *e.g.*, a 2.0% increase in the target cost.

2 PRELIMINARY AND RELATED WORK

2.1 PROBLEM STATEMENT OF AUTO-BIDDING

Consider a scenario where there are H impression opportunities arriving sequentially, each labeled by an index i . An advertiser wins an impression if its bid b_i surpasses those advertisers and incurs a cost c_i . The goal is to maximize the total value of the impressions won, represented by $\sum_i o_i v_i$, where v_i denotes the impression value and o_i is a binary variable indicating whether the advertiser wins impression i . Additionally, it is essential to consider the budget and various economic constraints, such as limiting the unit cost of specific advertising events like CPC and CPA (He et al., 2021). Therefore, under the budget constraint, the objective of auto-bidding is:

$$\text{maximize } \sum_i o_i v_i \quad \text{s.t.} \quad \frac{\sum_i o_i c_{ij}}{\sum_i o_i p_{ij}} \leq C_j, \forall j. \quad (1)$$

where p_{ij} is the performance indicator to j -th constraint provided by the advertiser and C_j is the upper bound. A previous study (He et al., 2021) has shown the optimal solution:

$$b_i^* = \lambda_0 v_i + \sum_{j=1}^J \lambda_j p_{ij} C_j, \quad (2)$$

where b_i^* is the optimal bid for impression i . The parameters λ_j represent the optimal bidding parameters. Please note that in industrial ad platforms, **auto-bidding methods focus on adjusting bidding parameters at a low frequency, such as every 30 minutes, rather than for every single impression that arises in milliseconds**. Therefore, there is ample time for an advanced auto-bidding method to respond to bidding parameter adjustment requirements, even for LLMs. Appendix A also provides a detailed literature review that highlights the evolution of auto-bidding methods.

2.2 AUTO-BIDDING AS DECISION-MAKING

We consider a sequential decision-making framework where an auto-bidding agent interacts with the advertising environment \mathcal{E} at discrete timesteps. At each timestep t , the agent receives a state $\mathbf{s}_t \in \mathcal{S}$ that describes the ad impression features and auction status, and then outputs an action $a_t \in \mathcal{A}$. The action is the adjustment to the bidding parameters λ_j , *i.e.*, $a_t := \{\lambda_t^{\lambda_0}, \dots, \lambda_t^{\lambda_J}\}$, which is then mapped to the final bid using Eq. (2). After transitioning to the next state, the env \mathcal{E} emits a reward r_t , typically the acquired conversion value. This process repeats until the end of a bidding period, such as one day, which could be seen as a fixed-length decision-making. A detailed description of these elements can be seen in Appendix C.

Related Work. As the complexity of online ad environments increased, RL became a major approach for auto-bidding. USCB (He et al., 2021) and SORL (Mou et al., 2022) employ RL to maximize the value under multiple constraints. Recently, DiffBid (Guo et al., 2024), a pioneering generative auto-bidding method, directly utilizes a decision diffuser to create a planning and control scheme, but not tested in large-scale auctions. GAS (Li et al., 2024) and GAVE (Gao et al., 2025a) propose data augmentation techniques to enhance the quality of offline datasets, but they still rely on the single-step reward and could introduce instabilities from the extra value estimation (Ajay et al., 2023). Therefore, this paper focuses on exploring the potential of diffusers in large-scale auctions for contributing a superior generative backbone.

2.3 DIFFUSER FOR AUTO-BIDDING

The diffuser is developed for decision-making based on a diffusion model (Janner et al., 2022b; Ajay et al., 2023). For simplicity, we introduce basic elements of the decision diffuser (DD) employed by DiffBid, which aims to maximize likelihood estimation of a trajectory $\mathbf{x}_0(\tau)$ given a condition $\mathbf{y}(\tau)$:

$$\max_{\theta} \mathbb{E}_{\tau \sim \mathcal{D}} [\log p_{\theta}(\mathbf{x}_0(\tau) | \mathbf{y}(\tau))], \quad (3)$$

where τ is the trajectory index and $\mathbf{y}(\tau)$ can be the property of the trajectory $\mathbf{x}_0(\tau)$, such as the accumulated reward of the trajectory. For brevity, we index the states from 0 to T , i.e., $\mathbf{x}_0 = \{s_0, s_1, \dots, s_T\}$. Specifically, DD is built with a predefined forward noising process $q(\mathbf{x}_{k+1}(\tau) | \mathbf{x}_k(\tau)) := \mathcal{N}(\mathbf{x}_{k+1}(\tau) | \sqrt{\alpha_k} \mathbf{x}_k(\tau), (1 - \alpha_k) \mathbf{I})$ and a trainable reverse denoising process $p_{\theta}(\mathbf{x}_{k-1}(\tau) | \mathbf{x}_k(\tau), \mathbf{y}(\tau)) := \mathcal{N}(\mathbf{x}_{k-1}(\tau) | \mu_{\theta}(\mathbf{x}_k(\tau), k, \mathbf{y}(\tau)), \sigma_k^2 \mathbf{I})$, where \mathcal{N} denotes a Gaussian distribution, and $\mathbf{x}_k(\tau) := \{s_0^k, s_1^k, \dots, s_T^k\}$ denotes the noisy trajectory at step k . We could get such a DD by **training** on a tractable variational evidence lower bound of $\log p_{\theta}(\mathbf{x}_0(\tau) | \mathbf{y}(\tau))$ with a surrogate loss (Ho et al., 2020):

$$\mathcal{L}_{\text{denoise}}(\theta) := \mathbb{E}_{k \sim U\{1, \dots, K\}, \mathbf{x}_0 \sim \mathcal{D}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \|\epsilon - \epsilon_{\theta}(\mathbf{x}_k(\tau), k, \mathbf{y}(\tau))\|^2, \quad (4)$$

where ϵ is the noise added to the data sample \mathbf{x}_0 to produce \mathbf{x}_k , and we use a deep neural network $\epsilon_{\theta}(\mathbf{x}_k, k, \mathbf{y}(\tau))$ to estimate the noise ϵ .

At a timestep t during **inference**, the diffuser combines the observed states and padding noise as inputs $\tilde{\mathbf{x}}_K = \{s_0, \dots, s_t, s_{t+1}^K, \dots, s_T^K\}$. It then generates a trajectory $\tilde{\mathbf{x}}_0 = \{s_0, \dots, s_t, \tilde{s}_{t+1}, \dots, \tilde{s}_T\}$, which consists of a sequence of multi-step states with expected property $\mathbf{y}(\tau)$. Subsequently, an action \mathbf{a}_t is typically obtained by inputting the observed states $s_{t-L:t}$ over a horizon L and the generated next state \tilde{s}_{t+1} into an inverse dynamics model f_{ϕ} , i.e., $\mathbf{a}_t = f_{\phi}(s_{t-L:t}, \tilde{s}_{t+1})$. **An illustration of this process can be seen in Fig. 7.** The inverse dynamic model could be independently trained by

$$\mathcal{L}_a(\phi) = \mathbb{E}_{\{s_{t-L:t+1}, \mathbf{a}_t\} \sim \mathcal{D}} \|\mathbf{a}_t - f_{\phi}(s_{t-L:t}, s_{t+1})\|^2. \quad (5)$$

However, the training objective described in Eq. (4) does not provide explicit guidance on maintaining dynamic legitimacy between adjacent states. This may lead to the generation of subsequent states $\tilde{s}_{>t}$ that are not achievable from the previously observed states $s_{\leq t}$, as demonstrated in Fig. 1 and further supported by the experiments shown in Fig. 2. Such issue can cause the inverse dynamic model to produce suboptimal actions due to these anomalous transitions.

Related Work. AdaptDiffuser (AD) (Liang et al., 2023) and MetaDiffuser (Ni et al., 2023) use a forward dynamic model for guiding the generation to pursue consistent trajectories, but poor base model generation can hinder effectiveness and necessitate training revision (Ajay et al., 2023). DiffuserLite (DL) (Dong et al., 2024a) reformulates the training process as generative interpolation by training multiple diffusers to generate segment trajectories, but it is inefficient and could limit overall planning ability. Adopting causal network architectures for training diffusers, such as diffusers based on RNN (Chen et al., 2024) or DiT (Peebles and Xie, 2023) with causal attention (Yang et al., 2021), is closely related to addressing this issue by enforcing the generation of next states based on historical states. **A more detailed discussion to the auto-regressive modeling and masked trajectory modeling can be seen in Appendix M.** Therefore, in this work, we firstly aim to propose an effective solution specifically in the context of auto-bidding for addressing the dynamic legitimacy issue without introducing extra modules.

3 CAUSAL AUTO-BIDDING VIA DIFFUSION COMPLETER-ALIGNER

We propose a novel causal auto-bidding method, termed CBD, to enhance the diffuser’s bidding performance in large-scale auctions by *i*) reformulating training a diffuser as learning a completer and *ii*) aligning the generated states to advertisers’ objectives in the inference process by refining the generation via an aligner. A detailed illustration of our CBD method could be seen in Fig. 1.

3.1 COMPLETER: LEARNING TO COMPLETE IN TRAINING

We now provide an insightful analysis of the compounding error that arises from directly applying diffusers for auto-bidding, which leads to the necessity of reformulating the bidding process as a completion process. We denote $\mathbf{o}_t := \{\mathbf{s}_{\leq t}, \mathbf{y}(\tau)\}$ as all the observations and \mathbf{s}' as the generated future states. Then, we can represent a simple optimal bidding strategy in inference as $\mathbf{a}_t = f_\phi(\mathbf{o}_t, \mathbf{s}')$, $\mathbf{s}' \sim p_{\theta^*}(\mathbf{s}'|\mathbf{o}_t, \mathbf{z}) = \mu_{\theta^*}(\mathbf{z}|\mathbf{o}_t) + \sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t)\epsilon$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, where \mathbf{z} is the input padding noise that will be denoised to the next states \mathbf{s}' conditioned on \mathbf{o}_t . We have a diffuser trained via Eq. 4, *i.e.*, denoising \mathbf{z} to \mathbf{s}' conditioned on a diffused observation $\tilde{\mathbf{o}}_t = a_t \mathbf{o}_t + b_t \boldsymbol{\eta}$, $\boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For brevity, we assume $a_t = b_t = 1$. Thus, the trained diffuser is actually performing $p_\theta(\mathbf{s}'|\tilde{\mathbf{o}}_t, \mathbf{z}) = \mu_\theta(\mathbf{z}|\mathbf{o}_t + \boldsymbol{\eta}) + \sigma_\theta(\mathbf{z}|\mathbf{o}_t + \boldsymbol{\eta})\epsilon = [\mu_\theta(\mathbf{z}|\mathbf{o}_t) + d_\mu(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta})] + [\sigma_\theta(\mathbf{z}|\mathbf{o}_t) + d_\sigma(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta})]\epsilon$, where d denotes the residual part. Now, we can conclude that if the diffuser is well-trained in transitioning to the distribution on \mathbf{s}' of optimal value, we have the *distribution matching error*, denoted as \mathcal{E}_t , of applying p_θ under the diffuser inference scheme, *i.e.*,

$$\begin{aligned} \mathcal{E}_t &= \mathbb{D}_{\text{KL}}[p_{\theta^*}(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})||p_\theta(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})] \\ &= \frac{1}{2} \left[\frac{\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t)^2 + d_\mu(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta})^2}{(\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t) - d_\sigma(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta}))^2} - 1 + \ln \frac{(\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t) - d_\sigma(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta}))^2}{\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t)^2} \right]. \end{aligned} \quad (6)$$

Proof can be seen in Appendix B. Intuitively, this is because of the unnecessary denoising of observations during inference, leading to $\mathcal{E}_t > 0$. As there are multiple steps of a bidding period, the final compounding error could be extremely large. This could lead to suboptimal performance of the diffuser, as it risks deviating from the optimal distribution of next states in p_{θ^*} , potentially entering dynamic illegitimate areas, as illustrated in Fig. 1.

To address this issue in auto-bidding, a practical approach to mitigate the compounding error is to eliminate the discrepancy between training and inference phases. This can be achieved by leveraging observed information to enhance generation, *i.e.*, learning to directly minimize the divergence $\mathbb{D}_{\text{KL}}[p_{\theta^*}(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})||p_\theta(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})]$. Specifically, the generation procedure of training and inference phase should be unified as follows: at a timestep t , we start generation by sampling $\mathbf{x}_K(\tau) = \{\mathbf{s}_0^K, \mathbf{s}_1^K, \dots, \mathbf{s}_T^K\}$ from $\mathcal{N}(\mathbf{0}, \mathbf{I})$ and assign the observed states $\mathbf{s}_{0:t}$ to the corresponding places to get the input of the diffuser, *i.e.*,

$$\tilde{\mathbf{x}}_K(\tau, t) = \underbrace{\{\mathbf{s}_0, \dots, \mathbf{s}_t\}}_{\text{Query}} \underbrace{\{\mathbf{s}_{t+1}^K, \dots, \mathbf{s}_T^K\}}_{\text{padding with noise}}. \quad (7)$$

Then, according to denoising process $p_\theta(\mathbf{x}_{k-1}(\tau)|\mathbf{x}_k(\tau), \mathbf{y}(\tau))$, we could sample the predicted next step's trajectory by

$$\tilde{\mathbf{x}}_{k-1}(\tau, t) = \mu_\theta(\tilde{\mathbf{x}}_k(\tau, t), k, \mathbf{y}(\tau)) + \sigma_k \mathbf{z}, \quad (8)$$

where $\mathbf{z} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The σ_k could be typically set as $\sigma_k := 1 - \alpha_k$ and the μ_θ could be parameterized by $\mu_\theta(\tilde{\mathbf{x}}_k(\tau, t), k, \mathbf{y}(\tau)) = \frac{1}{\sqrt{\bar{\alpha}_k}}(\tilde{\mathbf{x}}_k(\tau, t) - \frac{1 - \alpha_k}{\sqrt{1 - \bar{\alpha}_k}} \hat{\epsilon}_k)$, where $\bar{\alpha}_k := \prod_{i=1}^k \alpha_i$ and the $\hat{\epsilon}_k$ could be modeled with a conditional model $\epsilon_\theta(\tilde{\mathbf{x}}_k(\tau, t), k, \mathbf{y}(\tau))$ by randomly dropping the condition $\mathbf{y}(\tau)$ in a classifier-free guidance scheme (Ho and Salimans, 2022). By recursively applying the above procedure and assigning the first t -length of the generation with the query, we could get the final generated trajectory

$$\tilde{\mathbf{x}}(\tau, t) = \underbrace{\{\mathbf{s}_0, \dots, \mathbf{s}_t\}}_{\text{Query}} \underbrace{\{\tilde{\mathbf{s}}_{t+1}, \dots, \tilde{\mathbf{s}}_T\}}_{\text{Answer}}. \quad (9)$$

Additionally, as the dynamic legitimacy issue may occur in every position of the generated trajectory, we augment the previous diffuser training procedure by setting t as an extra random variable to construct the observations query and reformulate training as learning to complete. Specifically, with the augmented random variable t , the completer p_θ should capture the causality between the observed part $\{\mathbf{s}_0, \dots, \mathbf{s}_t\}$ and the remaining part $\{\mathbf{s}_{t+1}, \dots, \mathbf{s}_T\}$ with an ‘‘instruction’’ $\mathbf{y}(\tau)$, resembling a query-and-answer style like LLM. Therefore, the training objective is changed into

$$\max_{\theta} \mathbb{E}_{\tau \sim \mathcal{D}} [\log p_\theta(\mathbf{s}_{t+1:T}|\mathbf{s}_{0:t}, \mathbf{y}(\tau))]. \quad (10)$$

To achieve this objective, we can make a small modification by introducing an augmented new completion training loss \mathcal{L}_c , while still utilizing the powerful generation ability of diffusers, *i.e.*,

$$\mathcal{L}_c(\theta) := \mathbb{E}_{k \sim U\{1, \dots, K\}, \mathbf{x}_0 \sim \mathcal{D}, t \sim U\{0, \dots, T-1\}, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} \|\epsilon - \epsilon_\theta(\tilde{\mathbf{x}}_k(\tau, t), k, \mathbf{y}(\tau))\|_{t+1:T}^2 \quad (11)$$

where $\|\cdot\|_{t+1:T}^2$ means we only compute the loss on the position index $t+1$ to T in a trajectory. With such an additional variable $t \sim \text{Uniform}\{0, \dots, T-1\}$ and replacing $\mathbf{x}_k(\tau)$ to $\tilde{\mathbf{x}}_k(\tau, t)$ of Eq. (7), the model is learned to be aware of the dynamic legitimacy of adjacent states in every positions of the trajectory and thus enhancing the dynamic legitimacy of the whole trajectory.

3.2 ALIGNER: REFINING GENERATIONS IN INFERENCE

After training as a completer, the learned bidding policy, which includes the diffuser p_θ and the inverse dynamic model f_ϕ , is capable of producing actions with improved dynamic legitimacy. However, its alignment with the expected advertiser preferences, which could be represented by $\mathbf{y}(\tau)$, may still be suboptimal. This misalignment can be attributed to factors such as the limited capacity of the model, as we cannot get a perfect model, and the dynamic nature of advertiser preferences, which may include previously unseen preferences in datasets. Additionally, it remains challenging to explicitly control the sampling process during denoising to precisely obtain the expected optimal action a_t^* (Zheng et al., 2023a;b). Consequently, these concerns may lead to misalignment issues, failing to meet the advertiser’s expected properties, as illustrated in Fig. 2.

To enhance the quality of generated trajectories and improve subsequent auto-bidding performance, we can leverage a trajectory-level return model to introduce an aligner during the inference stage for refining the generated trajectories. The return model could be independently trained to regress on the conditions $\mathbf{y}(\tau)$ in the datasets expressed as

$$\mathcal{L}_r(\varphi) := \mathbb{E}_{\{\mathbf{x}_0, \mathbf{y}\} \sim \mathcal{D}} \|R_\varphi(\mathbf{x}_0(\tau)) - \mathbf{y}(\tau)\|^2. \quad (12)$$

This regression objective is not restricted by the issue of sparse rewards. Although the immediate reward r_t for most steps may be zero, leading to a sparse reward challenge, the accumulated rewards $\sum_{t=0:T} r_t$ are not zero and can be easy to learn.

There are several options for implementing the aligner scheme. One approach is a best-of-N sampling method that leverages the diverse generation ability of diffusers. This method generates multiple future trajectories in parallel and employs the trajectory-level return model to select the best trajectory, as illustrated below:

$$\tilde{\mathbf{s}}'_{t+1:T} \leftarrow \arg \min_{\tilde{\mathbf{s}}^i_{t+1:T}} \|R_\varphi(\mathbf{s}_{0:t}, \tilde{\mathbf{s}}^i_{t+1:T}) - \mathbf{y}(\tau)\|^2, i = 1 \sim N. \quad (13)$$

Alternatively, the aligner can be implemented using a more efficient revision method that performs a gradient update on $\tilde{\mathbf{s}}_{t+1:T}$ after the generation, which could be expressed as

$$\tilde{\mathbf{s}}'_{t+1:T} \leftarrow \tilde{\mathbf{s}}_{t+1:T} - \lambda \nabla_{\tilde{\mathbf{s}}_{t+1:T}} \|R_\varphi(\tilde{\mathbf{x}}_0(\tau, t)) - \mathbf{y}(\tau)\|^2. \quad (14)$$

When the generated trajectory has a higher estimated property value $R_\varphi(\tilde{\mathbf{x}}_0(\tau, t))$ than the expected property $\mathbf{y}(\tau)$ provided by the advertiser, the refinement operation would adjust the trajectory to a new one with a lower property value, and vice versa. Note that this approach differs from the classifier-guidance denoising technique used in diffusers, which iteratively interacts with every denoising step by $\mathbf{x}_{k-1} \sim \mathcal{N}(\boldsymbol{\mu}_k + \alpha \nabla_{\theta} \mathcal{J}(\boldsymbol{\mu}_k), \sigma_k^2 \mathbf{I})$. This iterative interaction creates challenges in extending it to various denoising sampling methods, such as DDIM (Song et al., 2021) and ODE (Zheng et al., 2023a) methods. In contrast, our approach refines only the final output to maintain computational efficiency and enable flexible future development, regardless of the training and inference techniques employed by the diffuser. More implementations of the aligner scheme could include multi-round refinement or exploring non-gradient methods, such as directly adding a residual to the output like LoRA (Hu et al., 2022) for optimal alignment. However, in this work, we found the gradient-based method in Eq. 14 to be effective and robust enough. We will report the aligner’s performance based on this gradient-based refinement method.

Finally, we feed the updated $\tilde{\mathbf{s}}'_{t+1}$ into the inverse dynamic model to obtain the final action. For a better understanding of our CBD method, we provide the detailed training and inference procedure in Algorithm 1 and 2 in Appendix D.

4 EXPERIMENT

4.1 SETUP

Datasets. To evaluate the performance of bidding strategy decision-making in large-scale ad auctions, we employ AuctionNet and its sparse variant, a publicly available benchmark from Alibaba designed for ad auctions, based on a real-world online advertising platform (Su et al., 2024). Each dataset comprises 21 advertising delivery periods, with each period containing approximately 500,000 impression opportunities, divided into 48 intervals. Details of the AuctionNet are in Appendix C.

Table 1: Comparison with baselines under the metric of Value in different budget settings in the MCB task. The best results are bolded and the base policy model’s results are underlined to demonstrate the performance improvement of our CBD method in the last column (improvements are statistically significant, *i.e.*, two-sided t-test with $p < 0.05$).

Dataset	Budget	RL				Transformer			Diffuser					improve
		USCB	BCQ	CQL	IQL	DT	CDT	DT-S	DiffBid	DiT	DiT-causal	CBD-Completer	CBD	
AuctionNet	50%	133.1	184.2	180.1	185.0	186.2	190.3	196.8	<u>161.9</u>	156.4	160.0	191.9	198.6 ±1.8	+22.3%
	75%	201.9	260.1	256.6	259.9	234.1	290.8	298.0	<u>233.7</u>	228.3	229.1	280.3	298.3 ±2.6	+27.8%
	100%	267.4	270.6	336.6	322.2	364.0	366.5	373.1	<u>316.5</u>	307.0	311.8	370.4	374.0 ±3.0	+18.3%
	125%	335.0	333.0	400.2	398.6	393.0	400.2	418.6	<u>384.0</u>	375.1	381.8	420.5	427.2 ±4.0	+11.1%
	150%	415.9	388.5	465.3	457.3	445.2	429.5	437.6	<u>452.5</u>	435.1	440.6	473.7	480.5 ±4.6	+6.19%
AuctionNet-sparse	50%	15.42	18.10	16.79	20.21	18.20	18.01	19.62	<u>14.54</u>	14.12	14.33	20.08	20.56 ±0.21	+41.4%
	75%	21.63	27.00	21.80	26.90	27.50	27.54	28.70	<u>23.25</u>	22.58	22.60	29.90	29.97 ±0.26	+28.9%
	100%	28.48	30.54	30.15	36.06	31.30	34.31	36.82	<u>31.33</u>	31.02	31.27	40.40	40.71 ±0.31	+29.9%
	125%	34.88	35.19	37.67	43.92	43.51	42.52	43.91	<u>40.54</u>	39.40	39.14	45.54	46.04 ±0.45	+13.5%
	150%	43.08	36.97	44.43	49.27	50.04	50.93	49.65	<u>47.18</u>	45.81	46.07	51.10	51.35 ±0.47	+8.91%

Evaluation Metrics. We test the baselines and our method on the maximize conversions bidding (MCB) task¹, and we adopt three metrics to evaluate the performance:

- **Value:** the total received impression value, *i.e.*, number of conversions $\sum_i o_i v_i$;
- **ER** (Exceeding rate of the KPI constraints): In a specific cost-cap setting, advertisers are concerned with how closely the strategy’s actual KPI performance matches the target KPI limit, such as CPA. To assess this, we introduce the ER metric for a delivery period, defined by $ER = \frac{1}{J} \sum_j C_j^{real} / C_j = \frac{1}{J} \sum_j (\sum_i c_{ij} o_i / \sum_i p_{ij} o_i) / C_j$.
- **Score:** by introducing a penalty term $penalty_j = \min\{(\frac{C_j}{\sum_i c_{ij} o_i / \sum_i p_{ij} o_i})^2, 1\}$, we can get a score as $score = (\sum_i o_i v_i) \times \min\{penalty_j\}_{j=1 \sim J}$.

Baselines. We conduct a comprehensive comparison of our method against a variety of baselines, including offline RL, DTs, and diffusers, to assess the potential of CBD as a new generative backbone. For RL, we evaluate against **USCB** (He et al., 2021), **BCQ** (Fujimoto et al., 2019), **CQL** (Kumar et al., 2020), and **IQL** (Kostrikov et al., 2022). For DTs, our comparisons include **DT** (Chen et al., 2021), **CDT** (Liu et al., 2023), and **DT-S**, which is a DT-based approach that incorporates the reward function by considering both the winning value and KPI constraints. **DT-S** represents the most advanced DT backbone utilized in prior works such as GAS (Li et al., 2024) and GAVE (Gao et al., 2025a). For diffusers, we compare with **DiffBid** (Guo et al., 2024), **DiT** (Peebles and Xie, 2023), which enhances DiffBid by replacing its U-Net (Ronneberger et al., 2015) backbone with a diffusion transformer, and **DiT-causal**, which further refines DiT by integrating a causal-attention mechanism (Waswani et al., 2017; Yang et al., 2021). The diffuser backbone of CBD is also based on Diffbid. For each method, the reported results are averaged over 5 random seeds.

We move the **implementation details** and detailed experimental setup to Appendix C.

4.2 PERFORMANCE COMPARISON WITH BASELINES

In this experiment, we conducted a performance comparison of various baseline methods under different settings, including varying datasets and budget constraints in the MCB task. The results are indicated by the value as a comprehensive assessment of the performance, as presented in Table 1. The CBD method consistently outperforms other models across all budget levels, with the most notable improvements in lower budget and sparse reward scenarios. Even without employing a trajectory-level return model for updating the predicted next state (referred to as CBD-Completer), its performance still significantly exceeds the baseline Diffbid, highlighting the benefits of employing a completer to address the dynamic legitimacy issue. Additionally, the DiT and its variant DiT-causal do not enhance the diffuser’s performance, underscoring the effectiveness of our training objective, *i.e.*, learning to perform auto-bidding as learning to complete. This insight could be inspiring for future research. To better understand the generation uncertainty issue, we provide a visualization in Fig. 2 and in Appendix in Appendix K by randomly sampling trajectories from the test set and comparing the generated trajectories between DiffBid and our methods. To enable a more comprehensive comparison with diffusion-based methods across a wider range of metrics, specifically considering the KPI constraints common in the cost-cap bidding mode, we present the

¹<https://support.google.com/google-ads/answer/7381968?hl=en>

Table 2: Comparison under more metrics of the diffuser-based methods and CBD variants, where arrow directions indicate performance improvement direction.

Dataset	Metrics	Diffbid	DiT	DiT-causal	CBD-Transformer	CBD-CVAE	CBD-A	CBD-Distillation	CBD-Completer	CBD
AuctionNet	Value \uparrow	316	307	311	313	314	354	365	370	374
	ER \downarrow	1.37	1.42	1.40	1.35	0.90	1.13	1.13	1.12	1.10
	Score \uparrow	214	196	205	210	284	294	285	292	298
AuctionNet-sparse	Value \uparrow	31.3	31.0	31.2	31.2	31.2	32.0	39.8	40.4	40.7
	ER \downarrow	1.23	1.32	1.30	1.21	1.25	1.24	0.95	0.90	0.86
	Score \uparrow	23.1	22.0	22.5	21.9	29.2	29.5	34.0	35.5	37.0

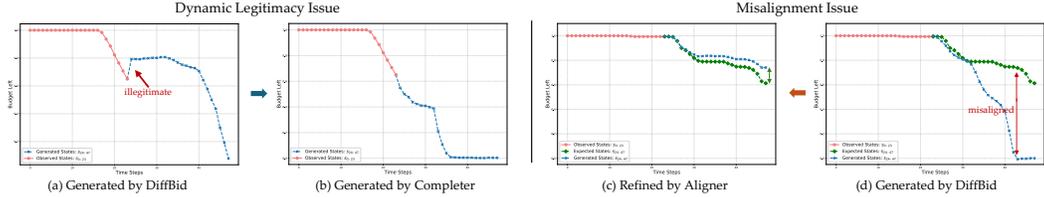


Figure 2: visualization of the generation uncertainty issue using an example case. In all subfigures, the red and green lines depict the ground-truth state information of the remaining budget, and the blue lines represent the generation results. Observing from $t = 0$ to 23 , DiffBid generates states where the remaining budget at $t = 24$ is greater than at $t = 23$, indicating a **dynamic legitimacy issue**. Additionally, DiffBid could have significant deviations of the generated trajectories from the expected trajectory, highlighting a **misalignment issue**.

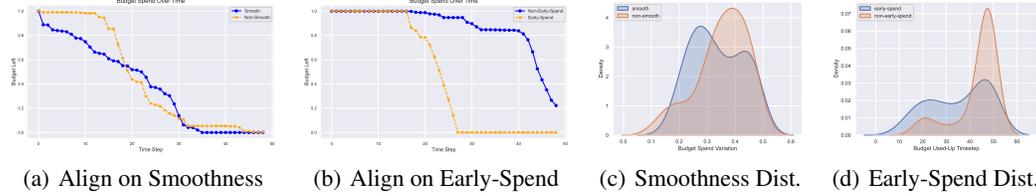


Figure 3: Results for property alignment. (a-b) show the generated trajectories with different properties, including smoothness and early spending; (c-d) illustrate the property distribution after aligning the generated trajectories to the desired properties at each bidding timestep.

results in Table 2. The CBD method consistently exhibits superior performance in ER, suggesting a more cost-efficient bidding strategy. Overall, CBD excels particularly in sparse scenarios, strongly endorsing the diffuser-based CBD as an exceptional generative backbone for industrial auto-bidding, where conversions are typically sparse.

4.3 PROPERTY ALIGNMENT PERFORMANCE

In addition to the aligner’s improvements shown in Table 1, we also examine other preferences depicted in Figures 3(a) and 3(b). We can generate diverse trajectories and utilize the aligner to refine these trajectories with desired characteristics, such as achieving a smoother trajectory alignment. Specifically, Figure 3(c) illustrates the distribution of budget spending variation, which represents smoothness; the lower the variation, the smoother the strategy. Figure 3(d) displays the distribution of budget used-up timesteps to demonstrate the property of early spending. The smooth budget spending is more concentrated and exhibits low variation, while the early-spend (depicted in blue) and non-early-spend (depicted in orange) differ in their timestep distributions. These results highlight the aligner’s effectiveness in refining trajectories. This flexibility is a significant advantage of the planning-based auto-bidding strategy, as it enhances the reliability of the bidding process and aids in diagnosing and correcting errors or biases.

4.4 ABLATION STUDY

Alternative Backbones compared to Diffusion Model. We conducted a comparison experiment with other simpler generative models, including a Transformer-based CBD method and a Conditional VAE (CVAE)-based CBD method, noted as CBD-Transformer and CBD-CVAE, respectively. Specifically, the CBD-Transformer utilizes a transformer where a masked sequence of states serves as the input, and the return of the trajectory acts as a condition interacting with the transformer block via AdaLN-zero. The CBD-CVAE employs a Temporal-UNet as the encoder and decoder, using the returns as the

condition. The results in Table 2 clearly demonstrate that employing diffusion models can achieve a better performance, demonstrating the necessity of adopting diffusion models for CBD.

Diffusion Steps and Distillation. The diffusion model typically requires repeating the recursive denoising step in Eq. (8) over multiple iterations, often up to 1000 steps. Therefore, we examined how the number of denoising steps affects performance to determine if a larger number of steps could improve outcomes. As shown in Fig. 4, we trained the diffusion model with varying denoising steps, including 10, 50, 100, 500, and 1000. Our findings indicate that 100 denoising steps are sufficient for achieving good performance, resulting in a computationally efficient approach. Though industrial auto-bidding services have sufficient time for model’s response as stated in Section 2.1, we also introduce a distillation method as a post-training trick, named SiD (Zhou et al., 2024), to distill our trained diffusion model into a one-step generator using the same U-Net while maintaining similar performance as the results of CBD-Distillation shown in Table 2. As these acceleration techniques can be directly applied to our method, the inference latency can be significantly reduced. Therefore, we highlight the importance of identifying the root causes of failures and developing advanced methods for training the original diffusion models for auto-bidding, which is the urgent need in this field.

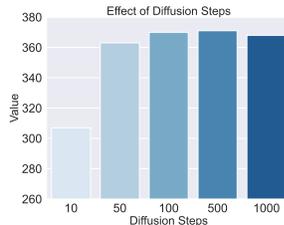


Figure 4: Effect of the diffusion steps.

We provide more ablation studies regarding the **Aligner’s Robustness** in Appendix E, **Necessity of Inverse Dynamic Model** in Appendix F, and **Extension to Offline RL Tasks** in Appendix G.

5 ONLINE A/B TEST

To verify the effectiveness of CBD, we have deployed it to an online ad platform. Under the MCB setting, the advertisers set the budget with or without the CPA/ROI constraint, and the bidding strategy aims to get as many conversions as possible under constraints. We compared CBD with the in-production DT-based method of a dedicated reward design. Over 7 days A/B test, there have been an average of 14700 campaigns per day. With the same budget allocated for each method’s experiment, each method achieved a daily average of approximately 80,000,000 impressions, which is substantial enough to demonstrate the effectiveness of the methods.

As shown in Table 3, the average achieved conversions (Target Cost) for advertisers are improved by an average of +2.0% while spending a similar budget (Cost), and the CPA valid ratio is also improved by +1.85%, which are significant improvements for such a large-scale auto-bidding production. Though our method incurs an additional 6ms compared to DT (6ms), it is still acceptable and cost-effective without needing additional computational resources, which could support a maximum of 26ms at the same frequency (around 20s) of calling the service to update the bidding parameters in the current computational resource. Given the substantial commercial value achieved, the additional inference latency is justifiable.

Table 3: A/B test results.

Metric	-Compare
Inference Latency	+6 ms
Cost	+0.0%
Target Cost	+2.0%
CPA Valid Ratio	+1.85%

6 LIMITATIONS AND CONCLUSION

This paper focuses exclusively on the auto-bidding strategy, which represents only one aspect of the real-time bidding system. Therefore, it is a limitation that we do not deeply consider other components, such as auction mechanism design and cooperation between advertisers, and their influence on the auto-bidding strategy. Given the difficulties in replicating these components within a public benchmark, we believe the effectiveness of CBD in handling real-world complexity could be empirically supported by the online results.

In conclusion, this paper explores the potential of diffusers for industrial auto-bidding. After analyzing failures in applying diffusers, we identified generation uncertainty as a key root cause and introduced a causal auto-bidding method based on a novel diffusion completer-aligner framework. Our method addresses the dynamic legitimacy issue through the completer and refines the generations to address misalignment using an aligner. Experiments support our analysis and demonstrate the superiority of CBD as a generative backbone for auto-bidding.

REFERENCES

- 486
487
488 Kan Ren, Weinan Zhang, Ke Chang, Yifei Rong, Yong Yu, and Jun Wang. Bidding machine: Learning
489 to bid for directly optimizing profits in display advertising. *IEEE Transactions on Knowledge and*
490 *Data Engineering*, 30(4):645–659, 2017.
- 491 Jiayan Guo, Yusen Huo, Zhilin Zhang, Tianyu Wang, Chuan Yu, Jian Xu, Bo Zheng, and Yan Zhang.
492 Generative auto-bidding via conditional diffusion modeling. In *Proceedings of the 30th ACM*
493 *SIGKDD Conference on Knowledge Discovery and Data Mining, KDD 2024, Barcelona, Spain,*
494 *August 25-29, 2024*, pages 5038–5049. ACM, 2024.
- 495
496 Han Zhu, Junqi Jin, Chang Tan, Fei Pan, Yifan Zeng, Han Li, and Kun Gai. Optimized cost per click
497 in taobao display advertising. In *Proceedings of the 23rd ACM SIGKDD international conference*
498 *on knowledge discovery and data mining*, pages 2191–2200, 2017.
- 499
500 Xun Yang, Yasong Li, Hao Wang, Di Wu, Qing Tan, Jian Xu, and Kun Gai. Bid optimization by
501 multivariable control in display advertising. In *Proceedings of the 25th ACM SIGKDD international*
502 *conference on knowledge discovery & data mining*, pages 1966–1974, 2019.
- 503
504 Yue He, Xiujun Chen, Di Wu, Junwei Pan, Qing Tan, Chuan Yu, Jian Xu, and Xiaoqiang Zhu. A
505 unified solution to constrained bidding in online display advertising. In *Proceedings of the 27th*
506 *ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, pages 2993–3001, 2021.
- 507
508 Han Cai, Kan Ren, Weinan Zhang, Kleantlis Malialis, Jun Wang, Yong Yu, and Defeng Guo. Real-
509 time bidding by reinforcement learning in display advertising. In *Proceedings of the tenth ACM*
510 *international conference on web search and data mining*, pages 661–670, 2017.
- 511
512 Richard S Sutton. Reinforcement learning: An introduction. *A Bradford Book*, 2018.
- 513
514 Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel,
515 Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence
516 modeling. In *Advances in Neural Information Processing Systems 34: Annual Conference on*
517 *Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages
518 15084–15097, 2021.
- 519
520 Zuxin Liu, Zijian Guo, Yihang Yao, Zhepeng Cen, Wenhao Yu, Tingnan Zhang, and Ding Zhao.
521 Constrained decision transformer for offline safe reinforcement learning. In *International Confer-*
522 *ence on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of
523 *Proceedings of Machine Learning Research*, pages 21611–21630. PMLR, 2023.
- 524
525 Anurag Ajay, Yilun Du, Abhi Gupta, Joshua B. Tenenbaum, Tommi S. Jaakkola, and Pulkit Agrawal.
526 Is conditional generative modeling all you need for decision making? In *The Eleventh International*
527 *Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*, 2023.
- 528
529 Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey Levine. Planning with diffusion for
530 flexible behavior synthesis. In *International Conference on Machine Learning, ICML 2022, 17-23*
531 *July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*,
532 pages 9902–9915. PMLR, 2022a.
- 533
534 Kefan Su, Yusen Huo, Zhilin Zhang, Shuai Dou, Chuan Yu, Jian Xu, Zongqing Lu, and Bo Zheng.
535 Auctionnet: A novel benchmark for decision-making in large-scale games. In *The Thirty-eight*
536 *Conference on Neural Information Processing Systems Datasets and Benchmarks Track*, 2024.
- 537
538 Tianchen Zhu, Yue Qiu, Haoyi Zhou, and Jianxin Li. Towards long-delayed sparsity: Learning a
539 better transformer through reward redistribution. In *Proceedings of the Thirty-Second International*
Joint Conference on Artificial Intelligence, IJCAI 2023, 19th-25th August 2023, Macao, SAR,
China, pages 4693–4701. ijcai.org, 2023.
- 538
539 Diederik P. Kingma and Max Welling. Auto-encoding variational bayes. In *2nd International*
Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014,
Conference Track Proceedings, 2014.

- 540 Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong
541 Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting Liu. A survey on hallucination in large
542 language models: Principles, taxonomy, challenges, and open questions. *ACM Trans. Inf. Syst.*, 43
543 (2):42:1–42:55, 2025.
- 544 Fei Ni, Jianye Hao, Yao Mu, Yifu Yuan, Yan Zheng, Bin Wang, and Zhixuan Liang. Metadiffuser:
545 Diffusion model as conditional planner for offline meta-rl. In *International Conference on Machine
546 Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of
547 Machine Learning Research*, pages 26087–26105. PMLR, 2023.
- 548 Yewen Li, Shuai Mao, Jingtong Gao, Nan Jiang, Yunjian Xu, Qingpeng Cai, Fei Pan, Peng Jiang, and
549 Bo An. GAS: generative auto-bidding with post-training search. *CoRR*, abs/2412.17018, 2024.
- 550 Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4RL: datasets for deep
551 data-driven reinforcement learning. *CoRR*, abs/2004.07219, 2020.
- 552 Zhixuan Liang, Yao Mu, Mingyu Ding, Fei Ni, Masayoshi Tomizuka, and Ping Luo. Adaptdiffuser:
553 Diffusion models as adaptive self-evolving planners. In *International Conference on Machine
554 Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of
555 Machine Learning Research*, pages 20725–20745. PMLR, 2023.
- 556 Zibin Dong, Jianye Hao, Yifu Yuan, Fei Ni, Yitian Wang, Pengyi Li, and Yan Zheng. Diffuserlite:
557 Towards real-time diffusion planning. In *Advances in Neural Information Processing Systems 38:
558 Annual Conference on Neural Information Processing Systems 2024, NeurIPS 2024, Vancouver,
559 BC, Canada, December 10 - 15, 2024*, 2024a.
- 560 William Peebles and Saining Xie. Scalable diffusion models with transformers. In *IEEE/CVF
561 International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023*,
562 pages 4172–4182. IEEE, 2023.
- 563 Xu Yang, Hanwang Zhang, Guojun Qi, and Jianfei Cai. Causal attention for vision-language tasks. In
564 *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2021, virtual, June 19-25,
565 2021*, pages 9847–9857. Computer Vision Foundation / IEEE, 2021.
- 566 Zhiyu Mou, Yusen Huo, Rongquan Bai, Mingzhou Xie, Chuan Yu, Jian Xu, and Bo Zheng. Sustain-
567 able online reinforcement learning for auto-bidding. In *Advances in Neural Information Processing
568 Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022,
569 New Orleans, LA, USA, November 28 - December 9, 2022*, 2022.
- 570 Jingtong Gao, Yewen Li, Shuai Mao, Peng Jiang, Nan Jiang, Yejing Wang, Qingpeng Cai, Fei Pan,
571 Kun Gai, Bo An, et al. Generative auto-bidding with value-guided explorations. *arXiv preprint
572 arXiv:2504.14587*, 2025a.
- 573 Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey Levine. Planning with diffusion for
574 flexible behavior synthesis. In *International Conference on Machine Learning, ICML 2022, 17-23
575 July 2022, Baltimore, Maryland, USA*, volume 162 of *Proceedings of Machine Learning Research*,
576 pages 9902–9915. PMLR, 2022b.
- 577 Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In Hugo
578 Larochelle, Marc’Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin,
579 editors, *Advances in Neural Information Processing Systems 33: Annual Conference on Neural
580 Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 581 Boyuan Chen, Diego Marti Monso, Yilun Du, Max Simchowitz, Russ Tedrake, and Vincent Sitzmann.
582 Diffusion forcing: Next-token prediction meets full-sequence diffusion. In *Advances in Neural
583 Information Processing Systems 38: Annual Conference on Neural Information Processing Systems
584 2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024.
- 585 Jonathan Ho and Tim Salimans. Classifier-free diffusion guidance. *CoRR*, abs/2207.12598, 2022.
- 586 Kaiwen Zheng, Cheng Lu, Jianfei Chen, and Jun Zhu. Improved techniques for maximum likelihood
587 estimation for diffusion odes. In *International Conference on Machine Learning, ICML 2023,
588 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning
589 Research*, pages 42363–42389. PMLR, 2023a.

- 594 Kaiwen Zheng, Cheng Lu, Jianfei Chen, and Jun Zhu. Dpm-solver-v3: Improved diffusion ODE
595 solver with empirical model statistics. In *Advances in Neural Information Processing Systems 36: Annual Conference on Neural Information Processing Systems 2023, NeurIPS 2023, New Orleans, LA, USA, December 10 - 16, 2023*, 2023b.
- 598
599 Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*, 2021.
- 602
603 Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang,
604 and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*. OpenReview.net, 2022.
- 606
607 Scott Fujimoto, David Meger, and Doina Precup. Off-policy deep reinforcement learning without
608 exploration. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 2052–2062. PMLR, 2019.
- 611
612 Aviral Kumar, Aurick Zhou, George Tucker, and Sergey Levine. Conservative q-learning for offline
613 reinforcement learning. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020.
- 615
616 Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-
617 learning. In *The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, April 25-29, 2022*, 2022.
- 619
620 Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical
621 image segmentation. In *Medical Image Computing and Computer-Assisted Intervention - MICCAI 2015 - 18th International Conference Munich, Germany, October 5 - 9, 2015, Proceedings, Part III*, volume 9351 of *Lecture Notes in Computer Science*, pages 234–241. Springer, 2015.
- 623
624 A Waswani, N Shazeer, N Parmar, J Uszkoreit, L Jones, A Gomez, L Kaiser, and I Polosukhin.
625 Attention is all you need. In *NIPS*, 2017.
- 626
627 Mingyuan Zhou, Huangjie Zheng, Zhendong Wang, Mingzhang Yin, and Hai Huang. Score identity
628 distillation: Exponentially fast distillation of pretrained diffusion models for one-step generation.
629 In *Forty-first International Conference on Machine Learning, ICML 2024, Vienna, Austria, July 21-27, 2024*, 2024.
- 630
631 Jun Wang and Shuai Yuan. Real-time bidding: A new frontier of computational advertising research.
632 In *Proceedings of the Eighth ACM International Conference on Web Search and Data Mining, WSDM 2015, Shanghai, China, February 2-6, 2015*, pages 415–416. ACM, 2015.
- 633
634 David S Evans. The online advertising industry: Economics, evolution, and privacy. *Journal of economic perspectives*, 23(3):37–60, 2009.
- 635
636 International Advertising Bureau. Iab internet advertising revenue report 2024. <https://www.iab.com/>, 2024.
- 638
639 Kushal S. Dave and Vasudeva Varma. Computational advertising: Techniques for targeting relevant
640 ads. *Found. Trends Inf. Retr.*, 8(4-5):263–418, 2014.
- 641
642 Jun Wang, Weinan Zhang, and Shuai Yuan. Display advertising with real-time bidding (RTB) and
643 behavioural targeting. *Found. Trends Inf. Retr.*, 11(4-5):297–435, 2017.
- 644
645 S. Muthukrishnan. Ad exchanges: Research issues. In *Internet and Network Economics, 5th International Workshop, WINE 2009, Rome, Italy, December 14-18, 2009. Proceedings*, volume 5929 of *Lecture Notes in Computer Science*, pages 1–12. Springer, 2009.
- 646
647

- 648 Kareem Amin, Michael J. Kearns, Peter B. Key, and Anton Schwaighofer. Budget optimization for
649 sponsored search: Censored learning in mdps. In *Proceedings of the Twenty-Eighth Conference on*
650 *Uncertainty in Artificial Intelligence, Catalina Island, CA, USA, August 14-18, 2012*, pages 54–63.
651 AUAI Press, 2012.
- 652
653 Weitong Ou, Bo Chen, Xinyi Dai, Weinan Zhang, Weiwen Liu, Ruiming Tang, and Yong Yu. A
654 survey on bid optimization in real-time bidding display advertising. *ACM Trans. Knowl. Discov.*
655 *Data*, 18(3):58:1–58:31, 2024.
- 656
657 Stuart Bennett. Development of the pid controller. *IEEE Control Systems Magazine*, 13(6):58–62,
658 1993.
- 659
660 Ye Chen, Pavel Berkhin, Bo Anderson, and Nikhil R. Devanur. Real-time bidding algorithms for
661 performance-based display ad allocation. In *Proceedings of the 17th ACM SIGKDD International*
662 *Conference on Knowledge Discovery and Data Mining, San Diego, CA, USA, August 21-24, 2011*,
663 pages 1307–1315. ACM, 2011.
- 664
665 Nicolas Grislain, Nicolas Perrin, and Antoine Thabault. Recurrent neural networks for stochastic
666 control in real-time bidding. In *Proceedings of the 25th ACM SIGKDD International Conference*
667 *on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019*,
668 pages 2801–2809. ACM, 2019.
- 669
670 Roger B. Myerson. Optimal auction design. *Math. Oper. Res.*, 6(1):58–73, 1981.
- 671
672 Sahin Cem Geyik, Sergey Faleev, Jianqiang Shen, Sean O’Donnell, and Santanu Kolay. Joint
673 optimization of multiple performance metrics in online video advertising. In *Proceedings of the*
674 *22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pages
675 471–480, 2016.
- 676
677 Brendan Kitts, Michael Krishnan, Ishadutta Yadav, Yongbo Zeng, Garrett Badeau, Andrew Potter,
678 Sergey Tolkachov, Ethan Thornburg, and Satyanarayana Reddy Janga. Ad serving with multiple
679 kpis. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery*
680 *and Data Mining*, pages 1853–1861, 2017.
- 681
682 Kan Ren, Weinan Zhang, Yifei Rong, Haifeng Zhang, Yong Yu, and Jun Wang. User response
683 learning for directly optimizing campaign performance in display advertising. In *Proceedings of*
684 *the 25th ACM International Conference on Information and Knowledge Management, CIKM 2016,*
685 *Indianapolis, IN, USA, October 24-28, 2016*, pages 679–688. ACM, 2016.
- 686
687 Ning li, Sai Kumar Arava, Chen Dong, Zhenyu Yan, and Abhishek Pani. Deep neural net with
688 attention for multi-channel multi-touch attribution. *CoRR*, abs/1809.02230, 2018.
- 689
690 Kan Ren, Yuchen Fang, Weinan Zhang, Shuhao Liu, Jiajun Li, Ya Zhang, Yong Yu, and Jun
691 Wang. Learning multi-touch conversion attribution with dual-attention mechanisms for online
692 advertising. In *Proceedings of the 27th ACM International Conference on Information and*
693 *Knowledge Management, CIKM 2018, Torino, Italy, October 22-26, 2018*, pages 1433–1442.
694 ACM, 2018.
- 695
696 Daisuke Moriwaki, Yuta Hayakawa, Isshu Munemasa, Yuta Saito, and Akira Matsui. Unbiased
697 lift-based bidding system. *CoRR*, abs/2007.04002, 2020.
- 698
699 Jian Xu, Xuhui Shao, Jianjie Ma, Kuang-chie Lee, Hang Qi, and Quan Lu. Lift-based bidding in ad
700 selection. In *Proceedings of the aaai conference on artificial intelligence*, volume 30, 2016.
- 701
702 Chao Wen, Miao Xu, Zhilin Zhang, Zhenzhe Zheng, Yuhui Wang, Xiangyu Liu, Yu Rong, Dong
703 Xie, Xiaoyang Tan, Chuan Yu, Jian Xu, Fan Wu, Guihai Chen, Xiaoqiang Zhu, and Bo Zheng. A
704 cooperative-competitive multi-agent framework for auto-bidding in online advertising. In *WSDM*
705 *’22: The Fifteenth ACM International Conference on Web Search and Data Mining, Virtual Event /*
706 *Tempe, AZ, USA, February 21 - 25, 2022*, pages 1129–1139. ACM, 2022.
- 707
708 Carl Doersch. Tutorial on variational autoencoders. *arXiv preprint arXiv:1606.05908*, 2016.

- 702 Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair,
703 Aaron Courville, and Yoshua Bengio. Generative adversarial networks. *Communications of the*
704 *ACM*, 63(11):139–144, 2020.
- 705
706 Ling Pan, Nikolay Malkin, Dinghuai Zhang, and Yoshua Bengio. Better training of gflownets with
707 local credit and incomplete trajectories. In *International Conference on Machine Learning, ICML*
708 *2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning*
709 *Research*, pages 26878–26890. PMLR, 2023.
- 710 Ziru Liu, Shuchang Liu, Zijian Zhang, Qingpeng Cai, Xiangyu Zhao, Kesen Zhao, Lantao Hu, Peng
711 Jiang, and Kun Gai. Sequential recommendation for optimizing both immediate feedback and
712 long-term retention. In *Proceedings of the 47th International ACM SIGIR Conference on Research*
713 *and Development in Information Retrieval*, pages 1872–1882, 2024.
- 714 Chongming Gao, Kexin Huang, Ziang Fei, Jiaju Chen, Jiawei Chen, Jianshan Sun, Shuchang
715 Liu, Qingpeng Cai, and Peng Jiang. Future-conditioned recommendations with multi-objective
716 controllable decision transformer. *arXiv preprint arXiv:2501.07212*, 2025b.
- 717
718 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Confer-*
719 *ence on Learning Representations*.
- 720
721 Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-
722 based offline reinforcement learning. In *Advances in Neural Information Processing Systems 33:*
723 *Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December*
724 *6-12, 2020, virtual*, 2020.
- 725
726 Michael Janner, Qiyang Li, and Sergey Levine. Offline reinforcement learning as one big sequence
727 modeling problem. In *Advances in Neural Information Processing Systems 34: Annual Conference*
728 *on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*,
pages 1273–1286, 2021.
- 729
730 Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi,
731 Quoc V. Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language
732 models. In *Advances in Neural Information Processing Systems 35: Annual Conference on Neural*
733 *Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 -*
December 9, 2022, 2022.
- 734
735 Zibin Dong, Yifu Yuan, Jianye Hao, Fei Ni, Yi Ma, Pengyi Li, and Yan Zheng. Cleandiffuser: An
736 easy-to-use modularized library for diffusion models in decision making. In *Advances in Neural*
737 *Information Processing Systems 38: Annual Conference on Neural Information Processing Systems*
738 *2024, NeurIPS 2024, Vancouver, BC, Canada, December 10 - 15, 2024*, 2024b.
- 739
740 Philipp Wu, Arjun Majumdar, Kevin Stone, Yixin Lin, Igor Mordatch, Pieter Abbeel, and Aravind
741 Rajeswaran. Masked trajectory models for prediction, representation, and control. In *International*
Conference on Machine Learning, pages 37607–37623. PMLR, 2023.
- 742
743 Xinzhi Zhang, Yifan Gao, Yaqing Hou, Xuechuan Liu, Zengyang Wang, Yi Cai, Bo An, and
744 Mengchen Zhao. Efficient decision sequence modeling via feature-level masking. 2025.
- 745
746
747
748
749
750
751
752
753
754
755

APPENDIX

A DETAILED BACKGROUND

A.1 AUTO-BIDDING IN COMPUTATIONAL ADVERTISING

The rapid digital transformation of commerce has significantly expanded the reach of online advertising platforms, making them essential for engaging audiences and driving sales (Wang and Yuan, 2015; Evans, 2009). As reported by the International Advertising Bureau, U.S. digital ad revenue from computational advertising revenue has reached \$114.2 billion (Bureau, 2024), significantly contributing to the success and sustainability of information systems and technologies. The primary goal of computational advertising is to effectively match ads with their relevant contexts on the web, such as the content where the ad is displayed and users’ search queries or browsing behavior (Dave and Varma, 2014; Wang et al., 2017). With the continuous influx of ad impression opportunities to the platform, real-time bidding (RTB) based display advertising, which emerged in 2009 (Muthukrishnan, 2009), has become a major paradigm of computational advertising, which enables online advertising platforms to sell individual ad impressions via hosting a real-time auction and facilitates advertisers to bid for the impression opportunity based on estimated value. Specifically, in real-time auto-bidding (Ren et al., 2017), when an impression opportunity is generated from a user’s visits, a bid request for the ad display impression opportunity, including its features like user, context, and auction status, is sent to numerous advertisers via an ad exchange. Then, advertisers employ specific auto-bidding algorithms to assess the bid request’s potential value and determine a bid price via Eq. 2 in real-time on demand-side platforms (DSPs). **Please note that the bidding parameters in Eq. 2 are adjusted by auto-bidding methods at a low frequency, such as every 30 minutes, rather than for every single impression that arises in milliseconds.** The ad exchanger then selects the highest bidder to display the ad, charging them the market prices, typically the second-highest bid price in a second-price auction setting (Amin et al., 2012). A detailed illustration can be seen in Figure 5.

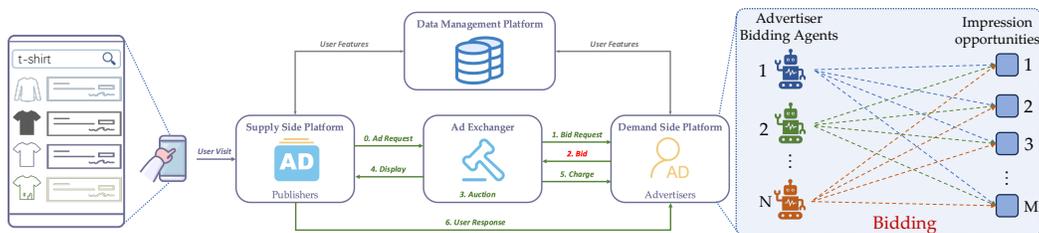


Figure 5: Real-time bidding system.

To maximize the total value of impressions for an advertiser while adhering to certain economic constraints, research into advanced auto-bidding strategies has progressed from rule-based approaches to prediction-based methods, and more recently, to reinforcement learning-based techniques (Ou et al., 2024). For rule-based methods, the proportional-integral-differential (PID) is the most widely used controller method (Bennett, 1993). The control signal consists of proportional, integral, and differential terms of the error, accounting for current, past, and future trends, resulting in comprehensive control of the variable. A typical example of applying PID in online advertising is controlling the campaign’s Cost-Per-Action (CPA) to a reference level (Chen et al., 2011). The controller could also be used to control the budget pacing to spend the budget smoothly (Grislain et al., 2019; Myerson, 1981). In addition to the PID controller, custom designs exist to adjust bids based on feedback data by classifying states as comfortable or risky and updating controlled variables accordingly (Geyik et al., 2016), while Multi-KPI (Kitts et al., 2017) controlled multiple KPIs simultaneously based on observed and target errors. For prediction-based techniques integrating auto-bidding strategies with upstream prediction tasks, such as CTR prediction, CVR prediction, and market price prediction, bid optimization is no longer an independent task but is correlated with these prediction tasks and jointly optimized. It is because a previous work, Bidding Machine (Ren et al., 2017), argues that interconnected tasks and sequential optimization lead to suboptimal solutions, and joint optimization allows the learning process to focus on valuable and competitive impressions. Some methods have tried this view, such as an EM-like (Expectation-Maximization) algorithm proposed to iteratively

optimize the CTR learning and bid optimization (Ren et al., 2016) and its extension incorporates market price estimation into the framework (Ren et al., 2017). Besides joint optimization, some methods tried introducing the lift/attribution prediction into the bid optimization, which is to predict the difference value before and after the ad is displayed (li et al., 2018; Ren et al., 2018; Moriwaki et al., 2020; Xu et al., 2016). For the reinforcement learning based methods, they repeatedly interact with the environment, allowing bidding agents to gain value from winning ad impression opportunities and update their subsequent strategies. The entire bidding process for advertisers is essentially a sequential decision-making process, which has recently demonstrated superior effectiveness.

A.2 SEQUENTIAL DECISION-MAKING FOR AUTO-BIDDING

As the complexity of online bidding environments increased, reinforcement learning (RL) algorithms like USCB (He et al., 2021), SORL (Mou et al., 2022), and MAAB (Wen et al., 2022) became essential for bidding. USCB (He et al., 2021) maximizes the value under multiple constraints and discovers a recursive structure to accelerate convergence. MAAB (Wen et al., 2022) considers the multi-agent modeling to maximize value and social welfare, proposing to mix cooperation and competition among multiple advertisers and improving the platform’s revenue. SORL (Mou et al., 2022) proposes a sustainable online RL framework that trains the auto-bidding policy by directly interacting with the RTB system, using safe and efficient online exploration instead of a simulated environment. However, online RL poses risks to the RTB system and cannot fully leverage the extensive historical bidding logs, resulting in a preference for bidding methods based on offline RL. Offline RL methods, which derive policies from existing datasets without requiring online interaction, have demonstrated considerable success. Noteworthy approaches include BCQ (Fujimoto et al., 2019), which constrains the action space to encourage on-policy behavior; CQL (Kumar et al., 2020), which regularizes Q-values for conservative estimates; and IQL (Kostrikov et al., 2022), which facilitates multi-step dynamic programming updates without querying Q-values of out-of-sample actions during training.

Despite their effectiveness, these methods are limited by the Markov Decision Process (MDP) assumption, whereas generative models (Doersch, 2016; Goodfellow et al., 2020; Pan et al., 2023; Ho et al., 2020; Chen et al., 2021) offer greater potential for bidding. Some methods try to apply decision transformers (DTs) (Chen et al., 2021; Liu et al., 2024; Gao et al., 2025b) for auto-bidding. GAS (Li et al., 2024) and GAVE (Gao et al., 2025a) propose data augmentation techniques via introducing an extra value prediction module to enhance the quality of the offline datasets. However, they still rely on dense sparse reward signals, and the introduced value function could suffer from instability issues, making them limited in the industrial auto-bidding tasks. (Ajay et al., 2023). As this paper focuses on proposing a superior generative backbone, we compare our CBD method to DT and its variants directly. In contrast, DiffBid (Guo et al., 2024) employs a decision diffuser to generate a long trajectory based on conditions like return, and then uses an inverse dynamic model to output the final action. However, in large-scale auctions, DiffBid faces challenges related to generation uncertainty, which is the primary focus of our work.

B PROOF OF THE DISTRIBUTION MISMATCH ERROR

To facilitate understanding the motivation behind the diffusion completer, we provide a detailed derivation of the distribution matching error in this section.

We denote $\mathbf{o}_t := \{s_{\leq t}, \mathbf{y}(\tau)\}$ as all the observations and s' as the generated future states. As there are multiple denoising steps of a diffusion model, for simplicity to get an insightful analysis, we assume the number of the denoising step as 1. Then, we can represent a simple optimal bidding strategy in inference as

$$a_t = f_\phi(\mathbf{o}_t, s'), s' \sim p_{\theta^*}(s'|\mathbf{o}_t, z) = \mu_{\theta^*}(z|\mathbf{o}_t) + \sigma_{\theta^*}(z|\mathbf{o}_t)\epsilon, \epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I}), \quad (15)$$

where z is the input padding noise that will be denoised to the next states s' conditioned on \mathbf{o}_t . The optimal distribution $p_{\theta^*}(s'|\mathbf{o}_t, z)$ signifies that the generated s' will lead to optimal auto-bidding performance.

We have a diffuser trained via Eq. 4, *i.e.*, denoising z to s' conditioned on a diffused observation $\tilde{\mathbf{o}}_t = a_t \mathbf{o}_t + b_t \boldsymbol{\eta}, \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For brevity, we assume $a_t = b_t = 1$. Thus, the trained diffuser is

864 actually performing

$$865 \begin{aligned} p_\theta(\mathbf{s}'|\tilde{\mathbf{o}}_t, \mathbf{z}) &= \mu_\theta(\mathbf{z}|\mathbf{o}_t + \boldsymbol{\eta}) + \sigma_\theta(\mathbf{z}|\mathbf{o}_t + \boldsymbol{\eta})\boldsymbol{\epsilon} \\ 866 &= [\mu_\theta(\mathbf{z}|\mathbf{o}_t) + d_\mu(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta})] + [\sigma_\theta(\mathbf{z}|\mathbf{o}_t) + d_\sigma(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta})]\boldsymbol{\epsilon}, \end{aligned} \quad (16)$$

868 where d denotes the residual part. Intuitively, p_θ learns to denoise the diffused $\tilde{\mathbf{o}}_t$ into clean \mathbf{o}_t first, then project \mathbf{o}_t to \mathbf{s}' . This approach could be effective in image generation tasks, where we typically start generation from a completely uninformative prior distribution, such as a Gaussian distribution, because we cannot observe any part of the generated images before generation. However, the first step may be unnecessary in the context of auto-bidding, as we can obtain a portion of the trajectory, *i.e.*, $\mathbf{s}_{0:t}$, during inference. Therefore, denoising $\mathbf{s}_{0:t}$ could be redundant and might even hinder the utilization of historical information.

875 Now, we can conclude that if the diffuser is well-trained in transitioning to the distribution on \mathbf{s}' of optimal value, *i.e.*,

$$876 \mathbb{D}_{\text{KL}}[p_{\theta^*}(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})||p_\theta(\mathbf{s}'|\tilde{\mathbf{o}}_t, \mathbf{z})] = 0, \quad (17)$$

879 we can have a solution that

$$880 \mu_{\theta^*}(\mathbf{z}|\mathbf{o}_t) = \mu_\theta(\mathbf{z}|\tilde{\mathbf{o}}_t) = \mu_\theta(\mathbf{z}|\mathbf{o}_t) + d_\mu(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta}) \quad (18)$$

$$881 \sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t) = \sigma_\theta(\mathbf{z}|\tilde{\mathbf{o}}_t) = \sigma_\theta(\mathbf{z}|\mathbf{o}_t) + d_\sigma(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta}). \quad (19)$$

883 With this solution, we apply the trained diffuser to the inference stage by receiving the observation \mathbf{o}_t as inputs, *i.e.*,

$$884 \begin{aligned} p_\theta(\mathbf{s}'|\mathbf{o}_t, \mathbf{z}) &= \mu_\theta(\mathbf{z}|\mathbf{o}_t) + \sigma_\theta(\mathbf{z}|\mathbf{o}_t)\boldsymbol{\epsilon} \\ 885 &= [\mu_{\theta^*}(\mathbf{z}|\mathbf{o}_t) - d_\mu(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta})] + [\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t) - d_\sigma(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta})]\boldsymbol{\epsilon}. \end{aligned} \quad (20)$$

888 Finally, we can get the *distribution matching error*, denoted as \mathcal{E}_t , of applying p_θ under the diffuser inference scheme, *i.e.*,

$$891 \begin{aligned} \mathcal{E}_t &= \mathbb{D}_{\text{KL}}[p_{\theta^*}(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})||p_\theta(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})] \\ 892 &= \frac{1}{2} \left[\frac{\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t)^2 + d_\mu(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta})^2}{(\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t) - d_\sigma(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta}))^2} - 1 + \ln \frac{(\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t) - d_\sigma(\mathbf{z}|\mathbf{o}_t, \boldsymbol{\eta}))^2}{\sigma_{\theta^*}(\mathbf{z}|\mathbf{o}_t)^2} \right] > 0. \end{aligned} \quad (21)$$

895 The \mathcal{E}_t can only be 0 when $d_\mu = d_\sigma = 0, \forall \boldsymbol{\eta} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, which intuitively means the trained diffuser is ideal in denoising any $\tilde{\mathbf{o}}_t$ with any noise $\boldsymbol{\eta}$ to \mathbf{o}_t . However, this is difficult and even impossible with the existing training methods for diffusers. Actually, this could be unnecessary, as we can model the training as a completion task to directly minimize the divergence $\mathbb{D}_{\text{KL}}[p_{\theta^*}(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})||p_\theta(\mathbf{s}'|\mathbf{o}_t, \mathbf{z})]$.

900 C DETAILS OF EXPERIMENTS

901 **Datasets.** AuctionNet effectively simulates the complexity and integrity of real-world ad auctions through the interaction of several key modules: the ad opportunity generation module, the auto-bidding module, and the auction mechanism module. The AuctionNet benchmark comprises two datasets: AuctionNet and AuctionNet-sparse. AuctionNet-sparse is a sparser version of AuctionNet, featuring fewer conversions. Each dataset includes 21 advertising delivery periods, with each period containing approximately 500,000 impression opportunities, divided into 48 intervals. The models are trained on these two datasets, while 5,000 randomly selected trajectories from each dataset are used exclusively for evaluating the visualization of the generation. Detailed parameters are provided in Table 4. Additionally, the state at each time step includes the following information:

- 911 • time_left: The remaining time steps left in the current advertising period.
- 912 • budget_left: The remaining budget that the advertiser has available to spend in the current advertising period.
- 913 • historical_bid_mean: The average values of bids made by the advertiser over past time steps.
- 914 • last_three_bid_mean: The average values of bids over the last three time steps.
- 915 • historical_LeastWinningCost_mean: The average of the least cost required to win an impression over previous time steps.

- `historical_pValues_mean`: The average of historical p-values over past time steps.
- `historical_conversion_mean`: The average number of conversions (e.g., sales, clicks, etc.) the advertiser achieved in previous time steps.
- `historical_xi_mean`: The average winning status of advertisers in impression opportunities, where 1 represents winning and 0 represents not winning.
- `last_three_LeastWinningCost_mean`: The average of the least winning costs over the last three time steps.
- `last_three_pValues_mean`: The average of conversion probability of advertising exposure to users over the last three time steps.
- `last_three_conversion_mean`: The average number of conversions over the last three time steps.
- `last_three_xi_mean`: The average winning status of advertisers over the last three time steps.
- `current_pValues_mean`: The mean of p-values at the current time step.
- `current_pv_num`: The number of impressions served at the current time step
- `last_three_pv_num_total`: The total number of impressions served over the last three time steps.
- `historical_pv_num_total`: The total number of impressions served over past time steps.

Table 4: The parameters of AuctionNet and AuctionNet-sparse.

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
Action range	[0, 493]	[0, 8178]
Impression’s value range	[0, 1]	[0, 1]
CPA range	[6, 12]	[60, 130]
Total conversion range	[0, 1512]	[0, 57]

Baseline Details. For RL methods, **USCB** (He et al., 2021), an online RL approach that adjusts parameters dynamically to the optimal bids; **BCQ** (Fujimoto et al., 2019), a typical offline RL method that updates policies with only access to a fixed dataset; **CQL** (Kumar et al., 2020), an offline RL method for learning conservative value function by regularizing Q-values; **IQL** (Kostrikov et al., 2022), an offline RL method that enables multi-step dynamic programming updates without querying out-of-sample actions;

Evaluation Details. To assess the auto-bidding performance, we perform experiments in a simulated environment that mirrors a real-world advertising system, as provided by Alibaba (Su et al., 2024). During the evaluation, an episode—also known as an advertising delivery period—consists of a day divided into 48 intervals, each lasting 30 minutes. Each episode includes approximately 500,000 impression opportunities that occur sequentially. An advertising delivery period features 48 advertisers from various categories, each with distinct budgets and CPAs, competing for all impression opportunities during the period. In each evaluation, our well-trained model acts as a specific advertiser, participating in bidding with a given budget and CPA. To thoroughly evaluate the model’s performance across different advertisers, we employ various advertiser configurations and advertising periods, conducting multiple evaluations in the simulated environment and averaging the results to obtain the evaluation score. [The reported results of each method is averaged over trained models with 5 random seeds.](#)

Implementation Details. For implementing the baselines, we use the default hyperparameters suggested in their respective papers and further fine-tune them to the best of our ability. For our CBD method, the diffuser backbone is based on Diffbid (Guo et al., 2024). The number of diffusion steps is set to 100, the horizon length of historical states is 3, and the batch size is 512 with a total of 500 training epochs. For the backbone of the diffuser, we adopt a temporal U-Net (Ronneberger et al.,

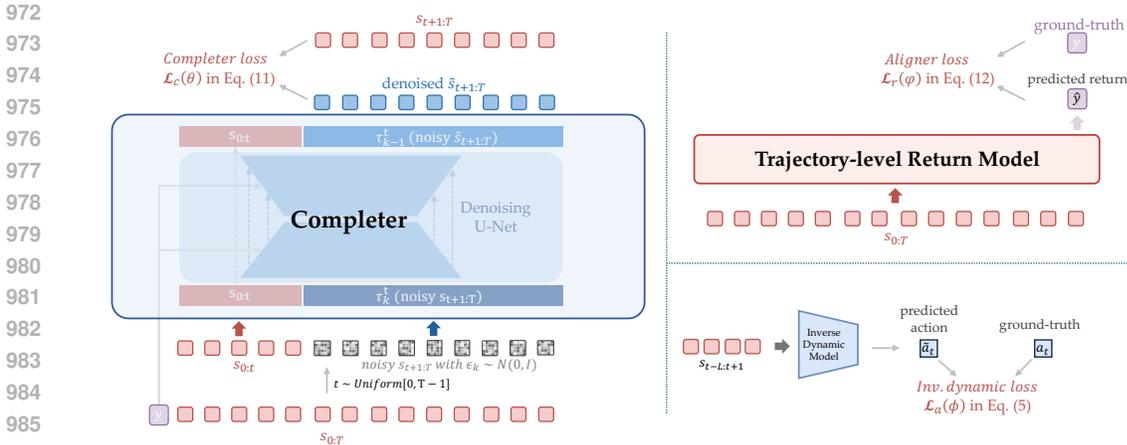


Figure 6: Demonstration of the training losses for each modules of our CBD method.

2015) with hidden sizes of [128, 256, 512]. The inverse dynamic model is implemented as an MLP with hidden sizes of [1024, 1024, 512]. The gradient guidance scale λ is set as 0.1 following the previous typical setting of diffusers (Janner et al., 2022b) for all tasks. We train the model using the Adam optimizer (Loshchilov and Hutter) with a learning rate of $1e-4$ on an H100 GPU. In this work, we adopt the setting from the previous decision diffuser (Ajay et al., 2023), where $y(\tau)$ is a scalar. In the main results shown in Tables 1 and 2, $y(\tau)$ is defined as $R(\tau) = \sum_{i=0}^T r_i$, which is scaled to $R(\tau) \in [0, 1]$ during training. During inference, $y(\tau) = 1$ is used to sample a high-return trajectory. In the general max conversion bidding setting, this approach is reasonable, as achieving more conversions implies utilizing more budget while maintaining a lower CPA.

Details of Extension to Offline RL Tasks. The offline reinforcement learning (RL) tasks under evaluation include three widely recognized locomotion challenges: HalfCheetah, Hopper, and Walker2d. These tasks involve maneuvering three distinct Mujoco robots to achieve optimal speed while ensuring energy efficiency and maintaining stability. The D4RL (Fu et al., 2020) benchmark offers offline datasets at three quality tiers: "medium," which includes demonstrations of moderate proficiency; "medium-replay," which encompasses all replay buffer recordings observed during training up to the point where the policy attains medium-level performance; and "medium-expert," which equally blends medium and expert-level performances. We also compared some additional non-diffuser-based methods on this benchmark, with details as follows: MOREL (Kidambi et al., 2020) addresses common pitfalls of model-based reinforcement learning, such as model exploitation, by learning a pessimistic Markov Decision Process (P-MDP) and training a near-optimal policy within this P-MDP. On the other hand, TT (Janner et al., 2021) employs a Transformer architecture to model distributions over trajectories and utilizes beam search as a planning algorithm.

D ALGORITHM

To provide a comprehensive understanding of our CBD method, we present a detailed algorithm outlining both the training and inference procedures in this section, as shown in Algorithm 1 and 2.

To offer a higher-level overview of the proposed method, we present a demonstration of the inference procedure in Figure 1 of the manuscript and a demonstration of the training procedure in Figure 6.

E ROBUSTNESS OF THE ALIGNER

For the principle of choosing a proper step size, as the offline data collected from online systems can be assumed to have a similar distribution to the online data, it is straightforward to use a portion of the offline data to evaluate the validity of different step sizes and employ the trajectory-level return model to assess alignment improvements. In this work, we choose the gradient step size of the aligner as 0.1, which generally brings improvement on various tasks. Furthermore, we have

Algorithm 1 Training Procedure of CBD

- 1: Initialize model parameters of completer ϵ_θ , aligner R_φ , and inverse dynamic model f_ϕ
- 2: **for** each iteration **do**
- 3: Sample $k \sim [1, K]$, $x_0 \sim \mathcal{D}$, $t \sim U\{0, \dots, T - 1\}$, $\epsilon \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$
- 4: Compute the losses:
- 5: $\mathcal{L}_c(\theta) = \mathbb{E} \|\epsilon - \epsilon_\theta(\tilde{\mathbf{x}}_k(\tau, t), k, \mathbf{y}(\tau))\|_{t+1:T}^2$
- 6: $\mathcal{L}_r(\varphi) = \|R_\varphi(\mathbf{x}_0(\tau)) - \mathbf{y}(\tau)\|^2$
- 7: $\mathcal{L}_a(\phi) = \|\mathbf{a}_t - f_\phi(\mathbf{s}_{t-L:t}, \mathbf{s}_{t+1})\|^2$
- 8: Update θ , φ , ϕ using gradient descent on $\mathcal{L}_c(\theta)$, $\mathcal{L}_r(\varphi)$, and $\mathcal{L}_a(\phi)$, respectively
- 9: **end for**

Algorithm 2 Inference Procedure of CBD

- 1: Initialize: initial state $\{\mathbf{s}_0\}$, completer ϵ_θ , aligner R_φ , inverse dynamic model f_ϕ
- 2: **for** bidding timestep $t = 0$ to $T - 1$ **do**
- 3: Generate future states $\{\tilde{\mathbf{s}}_{t+1}, \dots, \tilde{\mathbf{s}}_T\}$ using completer ϵ_θ by Eq. 8
- 4: Do refinement by aligner: $\tilde{\mathbf{s}}'_{t+1} \leftarrow \tilde{\mathbf{s}}_{t+1:T} - \lambda \nabla_{\tilde{\mathbf{s}}_{t+1:T}} \|R_\varphi(\tilde{\mathbf{x}}_0(\tau, t))\|^2$
- 5: Get action: $\mathbf{a}_t = f_\varphi(\mathbf{s}_{t-L:t}, \tilde{\mathbf{s}}'_{t+1})$
- 6: Execute bids and get next state \mathbf{s}_{t+1} from the ad system
- 7: Update query: $\{\mathbf{s}_0, \dots, \mathbf{s}_t, \mathbf{s}_{t+1}\} \leftarrow \{\mathbf{s}_0, \dots, \mathbf{s}_t\} \oplus \mathbf{s}_{t+1}$
- 8: **end for**

included experimental results demonstrating the robustness of the gradient step size in Table 5. The validity is measured by the generated remaining budget’s legitimacy. The GS-Align is an alternative alignment technique without gradient refinement, which is based on greedy selection of the best-aligned generated trajectory among parallelly sampled 10 generations, *i.e.*, best-of-N.

Table 5: Robustness and efficiency of the aligner’s gradient step size.

AuctionNet-Sparse	DiffBid	$\lambda = 0.001$	$\lambda = 0.01$	$\lambda = 0.05$	$\lambda = 0.1$	$\lambda = 1.0$	$\lambda = 5.0$	GS-Align
Value \uparrow	31.3	40.4	40.4	40.6	40.7	40.6	39.4	40.6
ER \downarrow	1.23	0.91	0.91	0.87	0.86	0.91	0.87	0.88
Score \uparrow	23.1	35.5	35.5	36.2	37.0	35.7	34.2	35.8
Validity \uparrow	77.4%	98.7%	98.7%	98.7%	98.7%	98.7%	96.0%	98.7%

F NECESSITY OF INVERSE DYNAMIC MODEL.

Current diffuser-based methods, such as Diffbid and our CBD method, generate states first and then derive an action, which is not straightforward since our major interest is in the action itself. This raises the question of whether we can directly generate the actions together with states, thereby avoiding the need for an additional inverse dynamic model. However, as shown in Table 2, our method still outperforms this variant (**CBD-A**), empirically demonstrating the effectiveness and stability of training the model to learn planning first, resembling a chain-of-thought process (Wei et al., 2022) before determining the final action.

G EXTENSION TO OFFLINE RL TASKS

Our CBD method could also bearing insights and benefit for general offline RL tasks like the locomotion tasks. We compare to the existing non-diffuser and diffuser-based methods together in Table 6, and a detailed introduction to these baselines can be seen in Appendix C. In some offline RL tasks, previous methods always set the first position of inputs as \mathbf{s}_t and let the diffuser generate the remaining sequence (Dong et al., 2024b) and we also follow this setting in these tasks. However, we emphasize that this is not necessary for auto-bidding since the trajectory usually has a fixed length within a given period, *e.g.*, 48 steps in a day with a bidding frequency of every 30 minutes. Additionally, historical information plays a crucial role in bid assessment (Guo et al., 2024). As shown in Table 6, the CBD method achieves state-of-the-art performance, demonstrating its effectiveness and generalizations in offline RL tasks. Notably, our method performs slightly better in the *Medium-Expert* setting than in others, suggesting that transition quality also impacts performance, *i.e.*, the *Medium-Expert* data trains the model to generate both dynamic legitimate and expert-level transitions, whereas other settings only teach the model about dynamic legitimacy. This situation differs from auto-bidding tasks, where the offline training log is derived from real advertisers’ behaviors, which inherently have a basic quality guarantee. Poor strategies will be immediately discontinued to prevent significant economic losses. However, this also presents unique challenges in developing new methods that can outperform existing competitive strategies in large-scale auctions.

Table 6: Performance comparison across different datasets and environments of the D4RL benchmark. Results correspond to the mean and standard error over 150 episode seeds, and the best scores of diffuser-based methods are emphasized in **bold**.

Dataset	Environment	non-Diffuser-based						Diffuser-based				
		BC	CQL	IQL	DT	TT	MOReL	Diffuser	DD	AD	DL	CBD
Medium-Expert	HalfCheetah	55.2	91.6	86.7	86.8	95	53.3	79.8	90.6	90.4	88.5	93.5±0.83
	Hopper	52.5	105.4	91.5	107.6	110.0	108.7	107.2	111.8	109.3	111.6	113.8±0.13
	Walker2d	107.5	108.8	109.6	108.1	101.9	95.6	108.4	108.8	108.2	107.1	109.1±0.33
Medium	HalfCheetah	42.6	44.0	47.4	42.6	46.9	42.1	44.2	49.1	44.3	48.9	49.4±0.20
	Hopper	52.9	58.5	66.3	67.6	61.1	95.4	58.5	79.3	96.6	100.9	97.7±1.38
	Walker2d	75.3	72.5	78.3	74.0	79.0	77.8	79.7	82.5	84.4	88.8	83.0±0.42
Medium-Replay	HalfCheetah	36.6	45.5	44.2	36.6	41.9	40.2	42.2	39.3	38.3	41.6	42.8±0.31
	Hopper	18.1	95	94.7	82.7	91.5	93.6	96.8	100.0	92.2	96.6	96.5±0.18
	Walker2d	26.0	77.2	73.9	66.6	82.6	49.8	61.2	75.0	84.7	90.2	75.2±0.26

H ETHICS STATEMENT

To the best of our knowledge, there are no ethical issues with this paper.

I REPRODUCIBILITY STATEMENT

We provide the code attached in the supplementary materials to ensure reproducibility, with detailed implementation details in Appendix C.

J STATEMENT OF THE USE OF LARGE LANGUAGE MODELS

We use the LLMs to polish the writing by checking for grammatical errors.

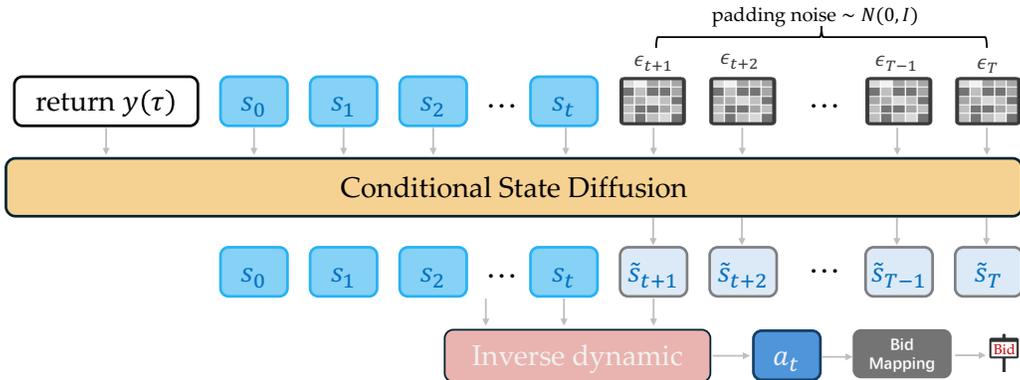


Figure 7: Overview of a typical diffuser-based decision-making procedure for auto-bidding: At each timestep for updating the bidding parameter, as introduced in Section 2.1, the diffuser first receives a trajectory consisting of a sequence of historical states $s_{0:t}$ and a sequence of padding noise $\epsilon_{t+1:T} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$, along with an expected return $y(\tau)$ as the condition. The conditional state diffusion model then denoises the padding noise to generations $\tilde{s}_{t+1:T}$. An inverse dynamic model projects the generated \tilde{s}_{t+1} , together with $s_{t-L:t}$, to the corresponding action a_t , which is the bidding parameter. During the extended period before the next decision-making timestep of bidding parameter update, which could last for a significant duration, such as half an hour, the bidding parameter a_t is continuously projected to the final bid for each impression opportunity using the bid mapping function in Eq. (2).

K MORE VISUALIZATION OF GENERATED TRAJECTORIES

We show more generated trajectories for visualization of the generation uncertainty issue in Fig. 8 and Fig. 9. The blue lines represent the actual trajectories found in the dataset, while the yellow lines depict the generated trajectories, which incorporate observations from the trajectory spanning from $t = 0$ to 23. A red line highlights the specific timestep at which decision-making occurs. As the results demonstrate, our CBD method effectively alleviates the issue of generation uncertainty by producing dynamic legitimate and well-aligned generations.

L BACKGROUND OF DIFFUSER-BASED DECISION-MAKING FOR AUTO-BIDDING

To facilitate understanding of the diffuser-based decision-making process for auto-bidding, we include an illustrative figure in Fig. 7, which is based on a commonly used decision diffuser (Ajay et al., 2023). The process consists of two main components. First, the decision diffuser takes in historical states and a return condition to generate future states. The action, specifically the bidding parameter λ as defined in Eq. (2), is then derived using an inverse dynamic model. The latency discussed in Section 4.4 refers to this component. Second, once the action a_t is determined at timestep t , it is continuously applied to the bidding of each impression opportunity until the next timestep $t + 1$ when the bidding parameters are needed to be updated. This application phase can extend over a longer duration, potentially lasting up to half an hour.

M DISCUSSION ON AUTO-REGRESSIVE AND MASKED TRAJECTORY MODELING METHODS

To distinguish our proposed method from prior auto-regressive and masked trajectory modeling approaches, we have included a detailed discussion in this section.

Auto-regressive methods, such as the Decision Transformer (DT), generate a single action using a causal transformer $p_\theta(a_t | \mathbf{R}t_{g_{\leq t}}, \mathbf{s}_{\leq t}, \mathbf{a}_{< t})$, where $\mathbf{R}t_{g_{\leq t}}$ represents the *return-to-go*. They are not directly modeling Eq. (10), *i.e.*, $p_\theta(\mathbf{s}_{t+1:T} | \mathbf{s}_{\leq t}, \mathbf{R})$. In contrast, our CBD method, which falls under the category of diffusers, generates all future states in **one step** rather than using an auto-regressive approach, which is the key distinction from DT. Additionally, CBD can effectively leverage historical observations $\mathbf{s}_{\leq t}$ to generate future states $\mathbf{s}_{t+1:T}$, due to our augmented training process that learns to complete sequences of random lengths. A significant challenge in applying dynamic programming methods like DT to auto-bidding is its inability to handle sparse-reward settings. In DT, actions \mathbf{a}_t are learned through the step-wise reward signals $r_t(\mathbf{s}_t, \mathbf{a}_t)$. However, r_t could always be zero in the sparse training set. A zero reward could indicate either a poor action or a good action that did not receive a reward due to conversion randomness, making it difficult to accurately assess the value of single steps. Our method overcomes this limitation by conditioning on a trajectory-level accumulative rewards, alleviating the need for accurate single-step rewards. This advantage is supported by previous works (Ajay et al., 2023; Janner et al., 2022a) and our experiments in Table 1, where our method significantly outperforms DT in sparse settings. Although one could adapt DT to function as a “one-step” generator by auto-regressively generating states until the end of the trajectory, essentially implementing CBD with a decision transformer, we conducted experiments labeled “**CBD-Transformer**” in Table 2 and the experiment result of CBD-Transformer showed poor performance, underscoring the necessity of using a diffuser to achieve Eq. (10).

Masked trajectory modeling (MTM) methods (Wu et al., 2023; Zhang et al., 2025) are designed to mask random positions within a trajectory and train a model to reconstruct the original trajectory using a transformer. The primary goal is to enhance the representation learning of a backbone model, which can be applied to a variety of downstream tasks such as dynamic modeling, imitation learning, and action generation. Applying this MTM approach to diffusers based on transformers can be seen as a variant of DiT-causal. Recall that DiT-causal enforces the generation of next states strictly based on historical states for each timestep during training, accomplished using the masking mechanism

Table 7: Performance comparison between Completer and Completer (MTM) under Dense and Sparse version of AuctionNet benchmarks.

Method	Dense			Sparse		
	Value \uparrow	ER \downarrow	Score \uparrow	Value \uparrow	ER \downarrow	Score \uparrow
Completer	370	1.12	292	40.4	0.90	35.5
Completer (MTM)	323	1.25	280	32.5	1.24	29.5

of causal attention. However, we found that DiT-causal performs poorly in experiments compared to our diffusion completer of CBD, as shown in Table 1 and Table 2. Intuitively, this could be due to the fact that a generated \tilde{s}_{t+1} in CBD not only receives information from real historical states $s_{\leq t}$ but also from the noisy future generated $\tilde{s}_{>t+1}^k$ during the denoising steps $k = \{T, \dots, 1\}$ in training and inference, forming a consistent trajectory based on a temporal U-Net backbone, which is a widely used and effective backbone in diffusers (Dong et al., 2024b). In contrast, DiT-causal strictly abandons the information from noisy generated $\tilde{s}_{>t+1}^k$ due to the masking mechanism, limiting its dynamic legitimacy. Therefore, instead of implementing MTM in transformers, we implemented MTM based on the same U-Net backbone as our CBD method. The masking ratio of MTM is set to 50% to keep a fair comparison with our Completer, which can observe an average of half the length of the trajectory. However, as shown in Table 7, results indicate that MTM could still limit the auto-bidding performance of diffusers. This could be due to the core motivation of this paper: the inconsistency between training and inference can affect performance. When applying MTM to a diffuser, it is trained to reconstruct a masked state \tilde{s}_t by partially utilizing information from the **incomplete** non-noisy historical states $masked(s_{<t})$. However, this is inconsistent with inference, where complete non-noisy historical states are available.

N EXPERIMENTS ON CLASSIFIER-GUIDANCE GENERATION

In the very beginning, the diffuser was implemented as a classifier-guidance scheme. Here, a trajectory-level return model $R(\tau)$ acts as the classifier providing gradient guidance during the denoising steps, expressed as $x_{k-1} \sim \mathcal{N}(\mu_k + \alpha \nabla_{\theta} \mathcal{J}(\mu_k), \sigma_k^2 I)$ (Janner et al., 2022a). However, this method has a severe drawback: The return model $R(\tau)$ is independently trained using real **non-noisy** trajectories and thus has not seen noisy trajectories during the denoising steps. These can be seen as out-of-distribution (OOD) cases, leading to gradient guidance failure. Unlike categorical distribution in image classifiers, the trajectory return distribution can be complex, and offline RL tasks often suffer from OOD errors (Kostrikov et al., 2022). In contrast, we apply the return model only in the final generation output, which is in the non-noisy trajectory distribution, exploiting the generative model’s ability in generating in-distribution data. To support this claim, we provide the best achieved experimental results of adopting an classifier-guidance (CG) generation scheme as below, noted as “CBD(CG)”.

Table 8: Performance comparison between CBD and CBD (CG) under Dense and Sparse version of AuctionNet benchmarks.

Method	Dense			Sparse		
	Value \uparrow	ER \downarrow	Score \uparrow	Value \uparrow	ER \downarrow	Score \uparrow
CBD	374	1.10	298	40.7	0.86	37.0
CBD (CG)	350	1.15	290	32.2	1.23	29.7

O ERROR BAR

For a clearer demonstration of each method’s robustness, Table 9 displays error bars calculated from 5 random seeds, corresponding to the main results in Table 1.

1242
 1243
 1244
 1245
 1246
 1247
 1248
 1249
 1250
 1251
 1252
 1253
 1254
 1255
 1256
 1257
 1258
 1259
 1260
 1261
 1262
 1263
 1264
 1265
 1266
 1267
 1268
 1269
 1270
 1271
 1272
 1273
 1274
 1275
 1276
 1277
 1278
 1279
 1280
 1281
 1282
 1283
 1284
 1285
 1286
 1287
 1288
 1289
 1290
 1291
 1292
 1293
 1294
 1295

Table 9: Performance comparison with error bar across various methods under different budget conditions for AuctionNet and AuctionNet-sparse datasets.

Dataset	Budget	RL				Transformer			Diffuser					Improve
		USCB	BCQ	CQL	IQL	DT	CDT	DT-S	DiffBid	DiT	DiT-causal	CBD-Completer	CBD	
Dense	50%	133.1±3.9	184.2±5.3	180.1±5.6	185.0±5.2	186.2±1.9	190.3±1.8	196.8±2.1	<u>161.9±3.1</u>	156.4±3.3	160.0±3.0	191.9±2.0	198.6±1.8	+22.3%
	75%	201.9±5.8	260.1±7.5	256.6±8.0	259.9±7.3	234.1±2.2	290.8±2.7	298.0±3.1	<u>233.7±4.9</u>	228.3±4.3	229.1±4.7	280.3±3.4	298.3±2.6	+27.8%
	100%	267.4±7.6	270.6±8.4	336.6±9.8	322.2±9.1	364.0±3.5	366.5±3.8	373.1±3.6	<u>316.5±6.5</u>	307.0±5.9	311.8±6.1	370.4±3.2	374.0±3.0	+18.3%
	125%	335.0±9.7	333.0±10.3	400.2±11.6	398.6±11.2	393.0±4.1	400.2±3.7	418.6±4.3	<u>384.0±7.9</u>	375.1±7.2	381.8±7.5	420.5±4.1	427.2±4.0	+11.1%
	150%	415.9±12.1	388.5±12.0	465.3±13.5	457.3±13.9	445.2±4.3	429.5±4.5	437.6±4.6	<u>452.5±8.8</u>	435.1±9.0	440.6±8.5	473.7±5.2	480.5±4.6	+6.19%
Sparse	50%	15.42±0.45	18.10±0.52	16.79±0.52	20.21±0.58	18.20±0.19	18.01±0.17	19.62±0.21	<u>14.54±0.30</u>	14.12±0.27	14.33±0.29	20.08±0.29	20.56±0.21	+41.4%
	75%	21.63±0.63	27.00±0.78	21.80±0.68	26.90±0.77	27.50±0.26	27.54±0.29	28.70±0.30	<u>23.25±0.48</u>	22.58±0.43	22.60±0.47	29.90±0.28	29.97±0.26	+28.9%
	100%	28.48±0.83	30.54±0.94	30.15±0.88	36.06±1.05	31.30±0.33	34.31±0.32	36.82±0.35	<u>31.33±0.64</u>	31.02±0.60	31.27±0.65	40.40±0.38	40.71±0.31	+29.9%
	125%	34.88±1.01	35.19±1.09	37.67±1.10	43.92±1.27	43.51±0.42	42.52±0.44	43.91±0.45	<u>40.54±0.79</u>	39.40±0.81	39.14±0.76	45.54±0.48	46.04±0.45	+13.5%
	150%	43.08±1.25	36.97±1.14	44.43±1.29	49.27±1.43	50.04±0.52	50.93±0.49	49.65±0.51	<u>47.18±0.92</u>	45.81±0.94	46.07±0.90	51.10±0.49	51.35±0.47	+8.91%

1296
1297
1298
1299
1300
1301
1302
1303
1304
1305
1306
1307
1308
1309
1310
1311
1312
1313
1314
1315
1316
1317
1318
1319
1320
1321
1322
1323
1324
1325
1326
1327
1328
1329
1330
1331
1332
1333
1334
1335
1336
1337
1338
1339
1340
1341
1342
1343
1344
1345
1346
1347
1348
1349

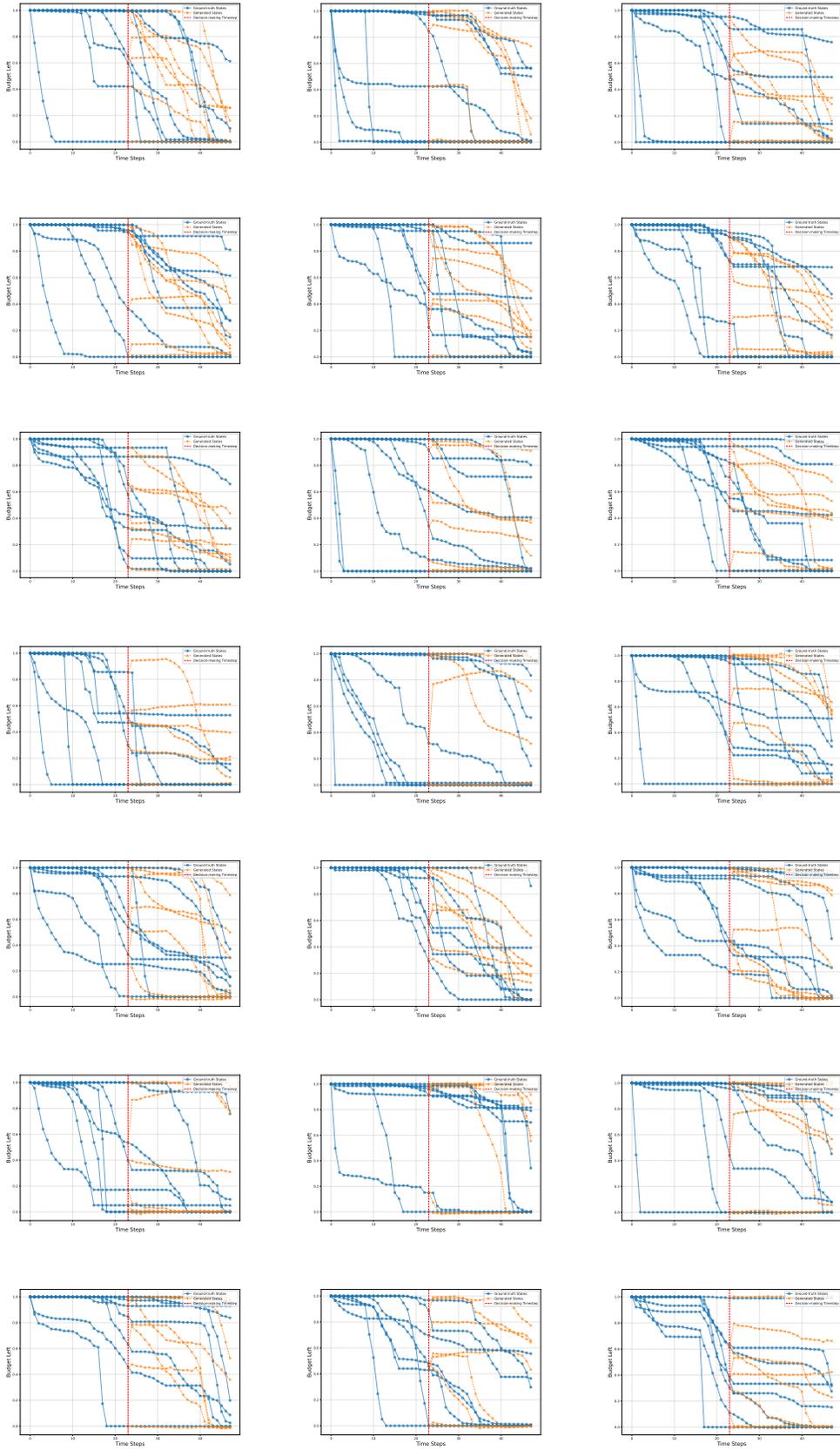


Figure 8: Generated trajectories of DiffBid.

1350
1351
1352
1353
1354
1355
1356
1357
1358
1359
1360
1361
1362
1363
1364
1365
1366
1367
1368
1369
1370
1371
1372
1373
1374
1375
1376
1377
1378
1379
1380
1381
1382
1383
1384
1385
1386
1387
1388
1389
1390
1391
1392
1393
1394
1395
1396
1397
1398
1399
1400
1401
1402
1403

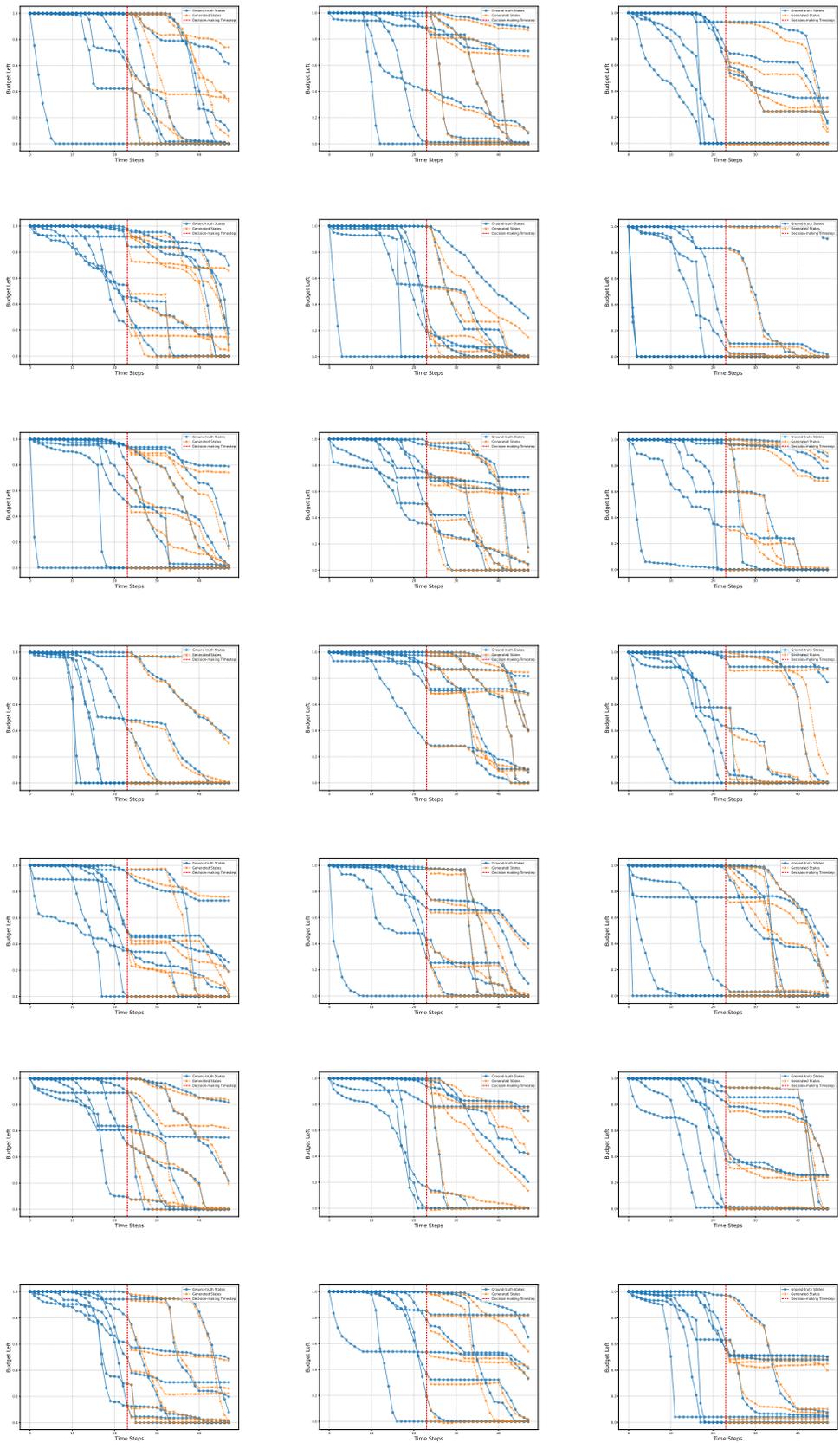


Figure 9: Generated trajectories of CBD.