# HALO: Hindsight-Augmented Learning for Online Auto-Bidding

**Pusen Dong, Chenglong Cao, Xinyu Zhou, Jirong You, Linhe Xu, Feifan Xu, Shuo Yuan**[*]

Shopee Company, China

{pusen.dong, chenglong.cao, xinyu.zhou, jirong.you, linhe.xu, feifan.xu, cody.tan}@shopee.com

## Abstract

Digital advertising platforms operate millisecond-level auctions through Real-Time Bidding (RTB) systems, where advertisers compete for ad impressions through algorithmic bids. This dynamic mechanism enables precise audience targeting but introduces profound operational complexity due to advertiser heterogeneity: budgets and ROI targets span orders of magnitude across advertisers, from individual merchants to multinational brands. This diversity creates a demanding adaptation landscape for Multi-Constraint Bidding (MCB). Traditional auto-bidding solutions fail in this environment due to two critical flaws: 1) severe sample inefficiency, where failed explorations under specific constraints yield no transferable knowledge for new budget-ROI combinations, and 2) limited generalization under constraint shifts, as they ignore physical relationships between constraints and bidding coefficients. To address this, we propose HALO: Hindsight-Augmented Learning for Online Auto-Bidding. HALO introduces a theoretically grounded hindsight mechanism that repurposes all explorations into training data for arbitrary constraint configuration via trajectory reorientation. Further, it employs B-spline functional representation, enabling continuous, derivative-aware bid mapping across constraint spaces. HALO ensures robust adaptation even when budget/ROI requirements differ drastically from training scenarios. Industrial dataset evaluations demonstrate the superiority of HALO in handling multi-scale constraints, reducing constraint violations while improving GMV.

## Introduction

In today's digital marketing landscape, advertisers strategically deploy multi-channel campaigns to reach consumers and drive conversions cost-effectively. These efforts primarily follow two distinct approaches based on campaign objectives (Vakratsas and Ambler 1999): Brand Advertising focuses on long-term brand building through immersive formats like promotional videos, prioritizing awareness and brand loyalty; while Performance Advertising emphasizes immediate results via e-commerce product ads and social media promotions, measured through quantifiable KPIs including click-through rate (CTR), conversion rate (CVR), gross merchandise volume (GMV), and return on investment (ROI). Performance advertising dominates the digital advertising landscape, accounting for 60% of industry expendi-

ture compared to brand advertising's 40% (Gartner 2024). This spending disparity establishes performance advertising as the primary focus of our research.

The quantifiable nature of performance advertising objectives necessitates infrastructure capable of real-time decision-making at scale. To address the diverse needs of advertisers, platforms like TikTok and Google Ads (TikTok 2025; Google 2025) utilize Real-Time Bidding (RTB) (Yuan, Wang, and Zhao 2013), the industry's standard infrastructure, which enables millisecond-level auctions for sequentially arriving ad impressions. Under this infrastructure, advertisers increasingly adopt multi-objective bidding strategies that maximize core outcomes (e.g., Gross Merchandise Volume (GMV)) while simultaneously satisfying multiple constraints, such as daily budgets and performance KPIs (e.g., ROI or Cost Per Action (CPA)). The bidding agent's main task is to dynamically compute a distinct bid for every opportunity, aligning real-time decisions with the advertiser's constraints and objectives. This real-time constraint satisfaction requirement formally characterizes the auto-bidding problem into two strategy classes: Budget-Constrained Bidding (BCB) (Zhou, Chakrabarty, and Lukose 2008) maximizes outcomes under only a total budget constraint. Multi-Constrained Bidding (MCB) (Yang et al. 2019) maximizes outcomes under both budget constraint and additional KPI requirements (e.g., ROI ≥ X). This paper focuses on solving both BCB and MCB challenges within a unified framework.

The design of a unified framework faces fundamental challenges due to the heterogeneity of advertisers. Digital advertising platforms serve heterogeneous advertisers, from multinational enterprises to small businesses, resulting in vastly divergent budget allocations and ROI targets. This operational diversity creates complex constraint landscapes: enterprise advertisers may deploy budgets orders of magnitude larger than small businesses while pursuing fundamentally different performance thresholds. Critically, advertisers frequently adjust constraints mid-flight, implementing sudden budget surges or tightening ROI thresholds during peak campaigns. Empirical evidence reveals that about a third of campaigns make such adjustments within a single auction day, triggering severe performance volatility.

Existing work (He et al. 2021; Guo et al. 2024; Wu et al. 2018) predominantly operates under an implicit homo-

---

[*]Corresponding author

geneity assumption that constraint magnitudes exhibit limited variation across advertisers. This foundational limitation manifests in two critical deficiencies: 1. Severe exploration inefficiency: existing methods require exhaustive trajectory exploration to find a viable bidding trajectory for each specific constraint configuration, while generating numerous failed attempts. Crucially, these failed attempts yield no transferable knowledge when constraints shift, forcing complete re-exploration for each new configuration and incurring substantial exploration waste. 2. Fundamental generalization failures: due to neglecting contextual relationships between constraints and bidding coefficients, existing methods fail to capture physical correlations among constraints and bidding coefficients. This oversight creates critical modeling gaps that cause catastrophic extrapolation errors under both distribution shifts and constraint adjustments. The compound effect of these deficiencies leads to concrete operational failures: significant extrapolation errors upon constraint shifts, resulting in unstable delivery pacing and suboptimal KPIs.

In this paper, we introduce an integrated framework that fundamentally transforms sample inefficiency and dynamic adaptation failures in constrained bidding. Drawing inspiration from human experiential learning, where even unsuccessful actions yield insights about alternative objectives, we establish a theoretically-grounded hindsight (Andrychowicz et al. 2017) sampling mechanism. This approach converts every exploration trajectory into productive training data, achieving near-perfect sample utilization efficiency. Crucially, we formally prove the efficacy of this mechanism and validate its optimization properties under constrained bidding conditions. Simultaneously, to resolve the extrapolation fragility, we abandon discrete point-wise mappings in favor of continuous functional mapping. Constraints-to-bid relationships are parameterized via adaptive B-spline (Prautzsch, Boehm, and Paluszny 2002). This ensures robustness against constraint shift by preserving contextual correlations between multi-scale constraints and bidding. Empirical validations on the industrial dataset conclusively demonstrate the superiority of our framework. In general, our contributions can be summarized into three aspects:

- We propose a theoretically-grounded hindsight sampling mechanism for the constrained bidding problem, which significantly improves sample utilization efficiency.

- We parameterize bid mappings via adaptive B-splines, achieving continuous adaptation to diverse constraint combinations while fundamentally enhancing extrapolation robustness against constraint shifts.

- Industrial validation on the real-world advertising dataset demonstrates the superiority of our solution.

## Problem Formulation

In a bidding period (normally one day), consider a sequentially arriving set $\mathcal{I}$ of $N$ impression opportunities. Advertisers participate auction by submitting a bid $b_i$ for each opportunity. The impression is won if $b_i$ exceeds the highest competing bid $c_i$. The cost of this impression is also $c_i$ determined by generalized second-price auction mechanism

(Edelman, Ostrovsky, and Schwarz 2007). The advertiser's objective during the bidding period is to maximize the total value of winning impressions $\max \sum_{i \in \mathcal{I}} v_i x_i$, where $v_i \in \mathbb{R}^+$ is the private impression value (the same opportunity has different values for different advertisers), $x_i$ is a binary variable indicating whether the campaign wins impression $i$. Beyond that, advertisers must satisfy specific operational constraints. When subject solely to a budget constraint, the problem is termed Budget-Constrained Bidding (BCB):

$$\max \quad \sum_{i \in \mathcal{I}} v_i x_i$$
$$s.t. \quad \sum_{i \in \mathcal{I}} c_i x_i \leq B$$

If additional KPI constraint (such as Return on Investment (ROI)) is imposed alongside the budget, the problem is termed Multi-Constrained Bidding (MCB):

$$\max \quad \sum_{i \in \mathcal{I}} v_i x_i$$
$$s.t. \quad \sum_{i \in \mathcal{I}} c_i x_i \leq B$$
$$\frac{\sum_{i \in \mathcal{I}} v_i x_i}{\sum_{i \in \mathcal{I}} c_i x_i} \geq r_{target},$$

where $r_{target}$ is the minimum allowable ROI.

Previous work (He et al. 2021) reformulated the constrained bidding problem as a linear programming problem, thus deriving the optimal bidding strategy for BCB and MCB scenarios.

$$BCB : b_i = \beta_0 v_i = \beta^{bcb} v_i$$
$$MCB : b_i = (\beta_0 - \beta_1 r_{target}) v_i = \beta^{mcb} v_i$$

where coefficient $\beta_0$ is relative to budget, and coefficient $\beta_1$ is relative to ROI. The fundamentally divergent optimal bidding formulations for BCB and MCB have prevented prior approaches from developing a unified optimization method for both paradigms. Our method decouples budget and ROI constraints to enable integrated processing of BCB and MCB within a single optimization architecture.

Additionally, due to the unpredictability of future impressions during real-time bidding, a fixed global optimal bid coefficient does not exist; bidding agents need to dynamically adjust bidding coefficients throughout the auction period to ensure cumulative constraints are satisfied. However, a single bidding period typically involves billions of impression opportunities, and it is not feasible to adjust for every impression. To address this, the common solution (Wu et al. 2018) is to divide the bidding period into $T$ discrete decision steps. At each step, the agent derives a bidding coefficient $\beta$, which is then applied to all impression opportunities within that interval.

## Method

We reformulate bidding strategies through dual-level innovations: at the data level, we propose a hindsight mechanism that enhances data efficiency; at the model level, we propose a parameterized B-spline to enable robust adaptation to multi-scale constraint configurations while reducing extrapolation errors. We now elaborate on these two innovations in detail. Due to space constraints, detailed proofs and formal definitions are provided in the appendix.

## Hindsight Mechanism

Existing works (He et al. 2021; Wang et al. 2022) typically solve the constrained bidding problem by independently optimizing the bidding coefficient for each ($B$, $r_{target}$) pair, employing trial-and-error (Kaelbling, Littman, and Moore 1996a) to determine the optimal coefficient. However, these works suffer from a critical limitation of sample inefficiency. Specifically, the exploration data collected to optimize one budget, say $B_k$, contribute minimally to the optimization process for a different budget $B_m$. This inefficiency arises because the target constraint $B_m$ differs from the budget scale that governed the data collection history.

To improve sample efficiency, we propose a hindsight mechanism. To demonstrate its effectiveness, we provide a series of theoretical guarantees. First, we define a naive strategy called Fixed Coefficient Strategy:

**Definition 1 (Fixed Coefficient Strategy (FCS))** *At any decision step $\tau$, a fixed coefficient strategy commits to the coefficient $\beta$ is constant for all remaining steps $[\tau, T]$.*

FCS consolidates sequential decision-making across steps $[\tau, T]$ into a single optimization step, significantly reducing the decision space from $O(T - \tau)$ to $O(1)$. Furthermore, we define:

**Definition 2 (Omniscient Optimal Value)** *The omniscient optimal value, denoted $V^{oracle}$, is the maximum achievable total value under perfect information with budget and ROI constraints.*

We prove that the maximum obtainable value can be achieved by FCS, exhibiting an ignorable error bound relative to the omniscient optimal value (solved by Mixed Integer Linear Programming (Floudas and Lin 2005)), as established in Lemma 1.

**Lemma 1 (FCS's Error Bound)** *The total potential obtainable value $V^{fixed}$ using a fixed coefficient strategy satisfying*

$$V^{fixed} > V^{oracle} - v_{max},$$

*where $v_{max}$ is the maximum individual impression value.*

Empirical evidence demonstrates that the value gap between the omniscient optimal value $V^{oracle}$ and the maximum individual impression value $v_{max}$ spans at least 2 orders of magnitude in practical advertising systems. Beyond that, under typical campaign conditions, the actual error remains much smaller than $v_{max}$, rendering the theoretical upper bound practically negligible. Consequently, we establish that FCS can achieve near-optimal performance, effectively approximating the omniscient optimum value while reducing optimization difficulty.

**Lemma 2 (Optimality Condition for FCS Under BCB)** *Under budget-constrained bidding with fixed coefficient Strategy, at any decision step $\tau$ with remaining budget $B_\tau$, if $\exists \beta^*$ such that*

$$\sum_{t=\tau}^{T} C_t(\beta^*) = B_\tau,$$

*then $\beta^*$ is the optimal coefficient on $[\tau, T]$.*

---

Algorithm 1: Hindsight Experience Collection

1: **Input:** total timesteps $T$, advertiser feature, number of sample coefficients $L$, uniform distribution $\mathcal{U}$
2: **Output:** Dataset $\mathcal{D} = \{(\mathbf{s}_i, \hat{\beta}_i, \text{cost}_i, \text{value}_i)\}$
3: **for** $l = 1$ to $L$ **do**
4:     **for** $\tau = 1$ to $T$ **do**
5:         $\hat{\beta}_\tau \sim \mathcal{U}$
6:         $\mathbf{s}_\tau \leftarrow$ compute_features
7:         $\hat{C}_\tau^T = 0$
8:         $\hat{V}_\tau^T = 0$
9:         **for** $t = \tau$ to $T$ **do**
10:             Bid with $\hat{\beta}_\tau$ at time-step $t$ and get $C_t, V_t$;
11:             $\hat{C}_\tau^T \leftarrow \hat{C}_\tau^T + C_t$
12:             $\hat{V}_\tau^T \leftarrow \hat{V}_\tau^T + V_t$
13:         **end for**
14:         $\mathcal{D}$.append $\left(\mathbf{s}_\tau, \hat{\beta}_\tau, \hat{C}_\tau^T, \hat{V}_\tau^T\right)$
15:     **end for**
16: **end for**
17: **return** $\mathcal{D}$

---

We establish that under BCB, the FCS achieves optimality if and only if the budget is fully exhausted over the interval $[\tau, T]$. Lemma 2 reveals a fundamental insight: multi-scale constrained bidding can be reformulated as a multi-objective trajectory alignment problem. Specifically, for any exploratory trajectory generated under a fixed coefficient $\beta$ (even when violating the target budget $B$), replacing $B$ with the realized cost $C_{realized}(\beta)$ transforms the trajectory into an optimal exploration for the posterior constraint configuration $B' = C_{realized}(\beta)$. Through this mechanism, every exploratory trial becomes an effective sample for some budget-constrained scenario. We formalize this paradigm as the Hindsight Mechanism, which establishes a bijection between historical explorations and valid constraint configurations. Through this mechanism, we can collect the hindsight experience through Algorithm 1.

**Definition 3 (Hindsight Experience)** *For any bidding trial over $[\tau, T]$ with fixed coefficient, we record*

$$\mathcal{H}_\tau = (s_\tau, \hat{\beta}_\tau, \hat{C}_\tau^T(\beta_\tau), \hat{V}_\tau^T(\beta_\tau)),$$

*as one hindsight experience tuple, where $s_\tau$ is the state feature, $\hat{C}_\tau^T(\hat{\beta}_\tau)$ is the realized cost over $[\tau, T]$ and $\hat{V}_\tau^T(\hat{\beta}_\tau))$ is the realized value.*

After elucidating how budget constraint learning benefits from the Hindsight Mechanism, we now systematically extend this framework to ROI constraint satisfaction. When introducing ROI target, the bid coefficient $\beta^*$ given by Lemma 2 is likely to cause the realized ROI to violate ROI constraint. Define the expected ROI:

$$R(\beta) = \mathbb{E}[\frac{\sum_{i \in \mathcal{I}} v_i x_i(\beta)}{\sum_{i \in \mathcal{I}} c_i x_i(\beta)}]$$

**Lemma 3 (ROI Monotonicity)** *The expected ROI function $R(\beta)$ is non-increasing in $\beta$.*

Lemma 3 indicates that the expected ROI typically increases as the bidding coefficient $\beta$ decreases. Therefore,

when the $\beta^*$ derived from Lemma 2 results in the realized ROI violating the ROI constraint, we need to search for an adjustment factor $\alpha < 1$ to shade $\beta^*$. This scaling reduces bid intensity, filtering out low-value traffic while ensuring the adjusted ROI falls within the target threshold. These two core mechanisms transform constrained optimization into a dual-control system: Lemma 2 governs budget exhaustion, while $\alpha$ enforces ROI compliance through strategic bid suppression. Then we get the following lemma:

**Lemma 4 (MCB Coefficient Adjustment)** *Under budget and ROI constraint with fixed coefficient Strategy, at any decision step $\tau$ with history cost $C_1^{\tau-1}$ and history value $V_1^{\tau-1}$, if $\exists \alpha \leq 1$ such that $\beta^{opt} = \alpha\beta^*$ applied on $[\tau, T]$ can make the total realized ROI equal to ROI constraint, then $\alpha$ satisfys*

$$\beta^* \int_\alpha^1 \left[ \gamma\left(\beta^* u\right) - r_{\text{target}} \eta\left(\beta^* u\right) \right] du = \Delta_{\text{ROI}}$$

*where $\gamma = \frac{\partial V}{\partial \beta}$ is derivative of value with respect to coefficient and $\eta = \frac{\partial C}{\partial \beta}$ is derivative of cost with respect to coefficient, $\Delta_{ROI} = V_1^{\tau-1} + V_\tau^T(\beta^*) - r_{target}(C_1^{\tau-1} + C_\tau^T(\beta^*))$ with $\beta^*$.*

Lemma 4 establishes the functional form of the adjustment factor $\alpha$ as dependent on seven key determinants: historical aggregate cost $C_1^{\tau-1}$, historical cumulative value $V_1^{\tau-1}$, future total cost with $\beta^*$ ($C_\tau^T(\beta^*)$), future total value with $\beta^*$ ($V_\tau^T(\beta^*)$), the derivative of value ($\gamma$), the derivative of cost ($\eta$), and the target ROI ($r_{target}$). This relationship is formally expressed as:

$$\alpha = f(C_1^{\tau-1}, V_1^{\tau-1}, C_\tau^T(\beta^*), V_\tau^T(\beta^*), \gamma, \eta, r_{target}).$$

The formulation exposes a fundamental limitation in traditional bidding methods: their failure to explicitly model the value derivative $\gamma$ and cost derivative $\eta$ critically undermines bidding performance. This omission prevents accurate sensitivity quantification of cost-value responses to coefficient variations, directly causing ROI constraint violations. Consequently, we fundamentally depart from conventional point-wise mapping approaches that naively fit constraint-to-coefficient mappings without physical interpretability. Instead, we establish continuous functional mappings that intrinsically encode derivative relationships.

Specifically, we model the bid coefficient function $\beta_\tau = f_\theta(B_\tau)$ and future value function $V_\tau^T = f_\phi(B_\tau)$, where $f_\theta, f_\phi : \mathbb{R}^+ \to \mathbb{R}^+$ are parameterized continuous functions learned from hindsight experiences $\mathcal{D} = \{\mathcal{H}_\tau\}$. It is worth noting that we transition from coefficient-centric (in Lemma 4) to budget-centric parameterization, because budget serves as the directly observable state variable that intrinsically defines the operational constraint boundary. This transformation preserves derivative relationships through the chain rule (Swokowski 1979) $\frac{\partial V}{\partial \beta} = \frac{\partial f_\phi}{\partial B} \cdot \frac{\partial B}{\partial \beta} = \frac{\partial f_\phi}{\partial B} \cdot \left(\frac{\partial f_\theta}{\partial B}\right)^{-1}$. This functional mapping intrinsically encodes gradient dynamics while ensuring operational feasibility, enabling derivative-informed computation of the adjustment factor $\alpha$ for precise constraint satisfaction. This fundamental idea delivers
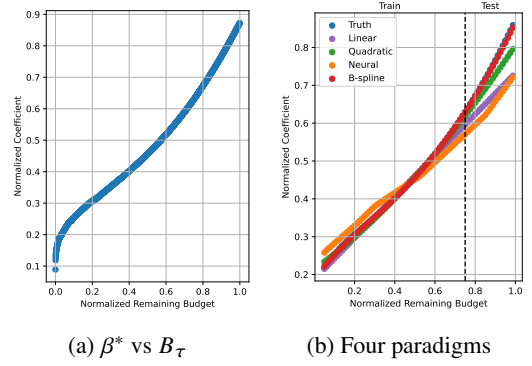


(a) $\beta^*$ vs $B_\tau$      (b) Four paradigms

Figure 1: (a) For a specific advertiser at decision step $\tau$, the relationship between remaining budget $B_\tau$ and $\beta^*$. (b) Partition the budget domain into training and testing to evaluate the generalization capabilities of four modeling paradigms under extrapolation conditions.

a breakthrough advance and makes our method outperform state-of-the-art alternatives

Combining the above lemma and definition, we can get the theorem below.

**Theorem 1 (Optimal Bidding Strategy)** *At any step $\tau$, with remaining budget $B_\tau$, target ROI $r_{target}$, historical cost $C_1^{\tau-1}$, historical value $V_1^{\tau-1}$, solve for $\beta_\tau = f_\theta(B_\tau)$, $V_\tau^T = f_\phi(B_\tau)$ and calculate $\Delta_{ROI} = V_1^{\tau-1} + V_\tau^T - r_{target}(C_1^{\tau-1} + B_\tau)$*

- *if $\Delta_{ROI} \geq 0$, $\beta_\tau$ is the optimal bidding coefficient.*
- *if $\Delta_{ROI} < 0$, $\alpha\beta_\tau$ is the optimal bidding coefficient.*

### B-spine Policy

As we discuss above, the continuous functions of value and cost are important to compute the adjustment factor $\alpha$ in Lemma 4. However, the properties of continuous value/cost functions exhibit significant variation across bidding environments. To ensure modeling tractability while preserving fidelity, functional assumptions become necessary because the assumptions critically govern both approximation accuracy and the precision of $\alpha$. To determine optimal functional representations, we categorize derivatives into three types with corresponding functional assumptions:

1. Constant Derivative ($\frac{\partial \beta}{\partial B} = $ constant): model the relationship function via linear polynomial $f_\theta(B_\tau) = \theta_0 + \theta_1 B_\tau$.
2. Linear Derivative ($\frac{\partial \beta}{\partial B} = a + b \cdot B$): model the relationship function via quadratic polynomial $f_\theta(B_\tau) = \theta_0 + \theta_1 B_\tau + \theta_2 B_\tau^2$.
3. Nonlinear Derivative ($\frac{\partial \beta}{\partial B} = g(B)$): model the relationship function via neural network, $f_\theta(B_\tau) = \text{MLP}(B_\tau; \theta)$

This taxonomy directly dictates computational methodology: constant derivatives permit closed-form $\alpha$ through algebraic manipulation, while linear and nonlinear derivatives need to compute $\alpha$ via iterative numerical techniques (Ypma 1995). The choice of assumptions inherently balances expressivity against computational complexity. Linear and quadratic polynomial functions impose a strong parametric prior that fails to capture the inherent heterogeneity in real-world online bidding environments. As empir-
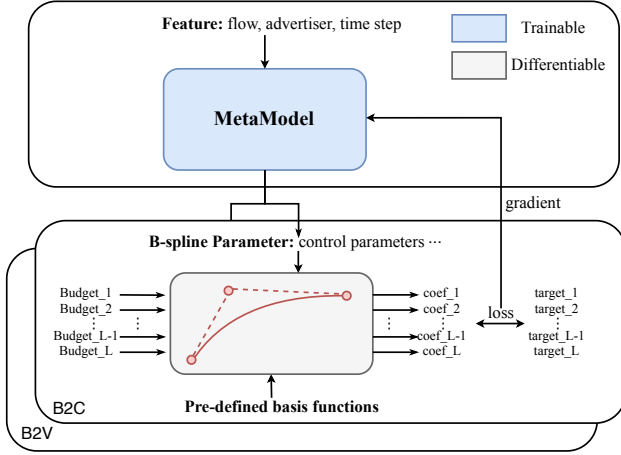
Figure 2: The training process of the Hindsight-Augmented B-Spline Policy, the blue Meta-Model is trainable, while the gray B-splines curve determined by the output of the Meta-Model. This B-spline module supports differentiable gradient backpropagation but remains non-trainable.

ically demonstrated in Figure 1 (a), the relationship between remaining budget and optimal coefficient exhibits non-stationary complexity that cannot be adequately modeled by rigid functional forms. While neural networks offer greater flexibility, their decision boundaries exhibit fragility and high sensitivity to input noise (Ghorbani, Abid, and Zou 2019).
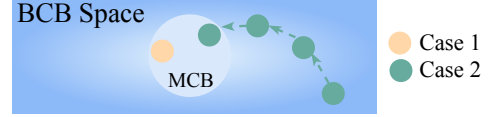
We therefore propose modeling the mapping relationship using B-splines (Liu et al. 2024), which provide adaptive local control through piecewise polynomial segments connected at knots. Crucially, this formulation ensures local modification independence, unlike neural networks. This locality property enables precise approximation of complex mappings while maintaining structural stability. The comparative analysis of four modeling paradigms (linear polynomial, quadratic polynomial, neural network, and B-spline) is visually summarized in Figure 1 (b), providing empirical validation of B-spline's superiority.

Based upon the analysis, we propose a Dual B-Spline Policy for the constrained bidding problem, where a shared MetaModel processes multi-dimensional features (traffic patterns, advertiser features, left time) to generate parameters ($\theta, \phi$ = MetaModel($s_\tau$)) for two B-spline curves: $\beta_\tau$ = B-Spline$_\theta(B_\tau)$, $V_\tau$ = B-Spline$_\phi(B_\tau)$ where $\theta$ and $\phi$ are learnable control points for their respective splines. We employ a meta-learning-inspired (Hospedales et al. 2021) paradigm to train the MetaModel, wherein predicted bidding coefficients $\beta_\tau$ and future values $V_\tau^T$ are simultaneously evaluated against ground-truth optimal labels (obtained by Algo. 1) to compute a composite loss function:

$$\mathcal{L} = \frac{1}{|\mathcal{S}|} \sum_{s_\tau \in \mathcal{S}} \frac{1}{L} \sum_{k=1}^{L} (\|\beta_{\tau,k} - \hat{\beta}_{\tau,k}\|^2 + \|V_{\tau,k} - \hat{V}_{\tau,k}\|^2).$$

Gradients are backpropagated through the dual spline modules to update the MetaModel end-to-end, enabling joint

**Algorithm 2: Optimal bidding Algorithm**



**Input:** current step $\tau$, state $s_\tau$, remaining budget $B_\tau$, target ROI $r_{\text{target}}$, historical cost $C_1^{\tau-1}$, historical value $V_1^{\tau-1}$, MetaModel $f_\Theta$, update scale $\lambda$, max iterations $K$
**Output:** optimal bidding coefficient $\beta_\tau^{\text{opt}}$

$\theta, \phi \leftarrow f_\Theta(s_\tau)$
$\beta_\tau \leftarrow f_\theta(B_\tau)$
$V_\tau^T \leftarrow f_\phi(B_\tau)$
$C_{\text{target}} \leftarrow B_\tau$
$\Delta_{\text{ROI}} \leftarrow V_1^{\tau-1} + V_\tau^T - r_{\text{target}}(C_1^{\tau-1} + C_{\text{target}})$
**if** $\Delta_{\text{ROI}} \geq 0$ **then**
    $\beta_\tau^{\text{opt}} \leftarrow \beta_\tau$                       ▷ Case 1
**else**
    **for** $k = 1$ to $K$ **do**
        $g \leftarrow \frac{\partial(\Delta_{\text{ROI}})}{\partial C_{\text{target}}}$
        $C_{target} \leftarrow C_{target} + \lambda \cdot g$
        $V_\tau^T \leftarrow f_\phi(C_{target})$
        $\Delta_{\text{ROI}} \leftarrow V_1^{\tau-1} + V_\tau^T - r_{\text{target}}(C_1^{\tau-1} + C_{\text{target}})$
        **if** $\Delta_{\text{ROI}} \geq 0$ **then**
            **break**
        **end if**
    **end for**
    $\beta_\tau^{\text{opt}} \leftarrow f_\theta(C_{target})$          ▷ Case 2
**end if**
**return** $\beta_\tau^{\text{opt}}$

functional optimization of both budget-to-coefficient (B2C) and budget-to-value (B2V) mappings. As illustrated in Figure 2, the training process performs parallelized evaluation of mapping at $L$ anchor budgets $\{B_k\}_{k=1}^{L}$ under identical contextual features. This bath comparison against the corresponding coefficient target and value targets holistically quantifies the B-spline's functional approximation quality across the budget domain. This design breaks through the limitations of traditional point-wise optimization, accelerating convergence through Multi-Anchor Joint Gradients while enhancing model generalization (Finn, Abbeel, and Levine 2017).

In practical deployment (Algo. 2), at each decision timestep $\tau$, we input multi-dimensional features into the MetaModel. This MetaModel generates control parameters for B-spline curves. Cost-to-Coefficient B-spline subsequently maps the advertiser's remaining budget $C_{target} = B_\tau$ to a bidding coefficient $\beta_\tau$ under BCB assumptions. Then we leverage Cost-to-Value B-spline to map the target cost $C_{target}$ to the expected future value $V_\tau^T$. If $\Delta_{\text{ROI}} = V_1^{\tau-1} + V_\tau^T - r_{\text{target}}(C_1^{\tau-1} + C_{\text{target}}) \geq 0$, $\beta_\tau$ inherently fall in the feasible solution space of MCB and is deployed without modification. Conversely, if $\Delta_{\text{ROI}} < 0$, we initiate a gradient-based correction protocol (Zhang et al. 2023) to dynamically adjust the target cost $C_{target}$ until ROI feasibility is

achieved:

$$\psi(C_{target}^k) = C_{target}^k + \lambda \cdot \frac{\partial}{\partial C_{target}^k} \left[ \Delta_{\text{ROI}} \right]^-$$

where $\lambda$ is the update scale. Then we get $\beta_\tau^{opt} = f_\theta(C_{target})$.

This framework intrinsically supports both BCB and MCB bidding modes. Furthermore, the robust extrapolation capability of B-spline parameterization enables real-time adaptation to multi-scale constraint configurations and constraint modifications mid-period. Such adaptability guarantees continuous constraint satisfaction under non-stationary advertising environments while maintaining solution optimality.

# Experiments

## Experiments Setup

**Dataset.** We evaluate our approach using AuctionNet (Su et al. 2024), a standardized large-scale advertising bidding dataset that abstracts real-world auction dynamics while leveraging authentic bidding data to provide a unified benchmark for evaluating diverse bidding strategies. This dataset includes 21 days of bidding records, with each day containing approximately 500,000 impression opportunities temporally partitioned into 48 distinct bidding time-steps. Critically, each daily auction period involves participation from 48 unique advertisers with different scales of constraints. All evaluation methods emulate 48 distinct advertisers, conducting bidding under heterogeneous constraint magnitudes. And we partition the dataset into a 14-day training set (678 bidding trajectories ) and a 7-day test set (339 trajectories).

**Baselines.** We benchmark our approach against three categories of existing methods:

1. **Rule-based Approaches**:
   - **PID** (Johnson and Moradi 2005): Adjusts bids dynamically using real-time error feedback.
   - **MPC** (Morari and Lee 1999): Solves constrained optimization over a finite time horizon.
   - **LP** (Dantzig 2002): Formulates bidding as a constrained linear programming problem.

2. **Reinforcement Learning (RL) Methods** (Kaelbling, Littman, and Moore 1996b):
   - **IQL** (Kostrikov, Nair, and Levine 2021): Offline RL leveraging value function approximation.
   - **USCB** (He et al. 2021): A method for learning to dynamically adjust bid coefficients using the DDPG (Lillicrap et al. 2015) algorithm.

3. **Imitation (Hussein et al. 2017) & Generative (Bond-Taylor et al. 2021) Approaches**:
   - **BC** (Torabi, Warnell, and Stone 2018): Mimics offline bidding decisions through supervised learning.
   - **DT** (Chen et al. 2021): Generates actions using return-conditioned sequence modeling.
   - **GAVE** (Gao et al. 2025): DT-based model, accommodates various advertising objectives through a score-based Return-To-Go (RTG) module.
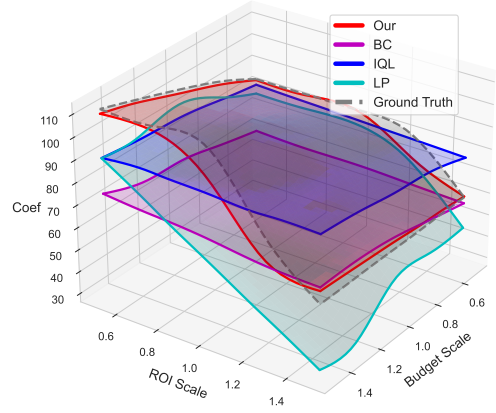


Figure 3: The variation of average bidding coefficients with different budget and ROI constraints. Each surface is interpolated from a grid of 25 discrete points (5 budget scales × 5 ROI scales).

**Evaluation Metrics.** For each advertiser, the core objective within every auction period is to maximize the aggregate value (conversions here) of winning impressions while adhering to predefined constraints. We categorize the evaluation paradigm into BCB and MCB. Under MCB settings, we formalize two evaluation metrics: **Conversions (Conv)**: Average conversions exclusively during periods where ROI constraints are satisfied. Conversions from constraint-violating periods are strictly excluded. **Compliance Rate (C.R.):** Proportion of evaluation periods where ROI constraints are fully satisfied. For BCB scenarios (absent ROI constraints), two metrics are defined: **Conversions (Conv):** Average conversions. **Cost/Budget Ratio (C/B):** Measures actual expenditure relative to allocated budget.

**Implementation Details.** All experiments are conducted on NVIDIA V100 GPU. The sample number of $\beta$ in Algo. 1 is 10. We use Adam optimizer to optimize MetaModel parameters with a learning rate of $1e^{-2}$. The state feature specifications and additional model hyperparameters are provided in the appendix. To ensure statistical significance, we conducted five independent experiments and report aggregated performance metrics.

## Evaluation Results

**Main Result.** Table 1 presents a comprehensive performance comparison of our method against baselines across varying bidding modes and budget scales. The results demonstrate that our approach achieves dominant performance across all bidding modes and budget scales. Under MCB settings, our method consistently maintains constraint compliance rates exceeding 90% while maximizing value delivery. In BCB scenarios, our method exhibits precise budget-scale adaptability and achieves near-perfect budget utilization with cost/budget ratios around 99% regardless of budget scales. The dual-mode efficacy establishes unprecedented robustness in real-world advertising operations. Traditional rule-based methods (PID, MPC, LP) underperformed due to relying exclusively on real-time error feedback. Reinforcement learning (IQL, USCB) and generative

Table 1: **Evaluation results.** For value-related and budget-related metrics, the top-performing method is highlighted in boldface, while ROI constraint compliance rates exceeding 90% are denoted with underlining.

| Mode | Budget | Metrics | BC | DT | IQL | PID | MPC | LP | GAVE | USCB | Our |
|---|---|---|---|---|---|---|---|---|---|---|---|
| MCB | 50% | Conv | 15.31 | 16.54 | 13.44 | 8.1 | 10.71 | 11.62 | 15.08 | 14.92 | **18.58** |
| | | C.R. | 0.938 | 0.854 | 0.792 | 0.729 | 0.708 | 0.833 | 0.833 | 0.938 | 0.917 |
| | 75% | Conv | 19.85 | 22.96 | 18.25 | 10.33 | 15.29 | 17.19 | 19.21 | 19.69 | **25.46** |
| | | C.R. | 0.938 | 0.854 | 0.771 | 0.667 | 0.688 | 0.812 | 0.812 | 0.979 | 0.938 |
| | 100% | Conv | 23.15 | 28.5 | 22.31 | 13.1 | 18.67 | 22.1 | 23.58 | 23.5 | **31.35** |
| | | C.R. | 0.938 | 0.854 | 0.75 | 0.729 | 0.729 | 0.792 | 0.833 | 0.938 | 0.938 |
| | 125% | Conv | 25.94 | 33.81 | 25.77 | 17.46 | 24.23 | 25.56 | 28.48 | 27.6 | **36.06** |
| | | C.R. | 0.938 | 0.875 | 0.729 | 0.646 | 0.708 | 0.833 | 0.875 | 0.979 | 0.917 |
| | 150% | Conv | 27.6 | 37.73 | 29.98 | 17.75 | 29.6 | 31.65 | 29.58 | 29.5 | **42.25** |
| | | C.R. | 0.938 | 0.875 | 0.729 | 0.792 | 0.875 | 0.917 | 0.833 | 0.938 | 0.938 |
| BCB | 50% | Conv | 16.31 | 19.94 | 18.71 | 11.33 | 16.54 | 15.35 | 16.14 | 14.25 | **24.15** |
| | | C/B | 0.594 | 0.842 | 0.827 | 0.739 | 0.899 | 0.591 | 0.611 | 0.624 | **0.992** |
| | 75% | Conv | 21.6 | 27.62 | 26.81 | 16.44 | 21.04 | 23.31 | 19.27 | 19.23 | **35.21** |
| | | C/B | 0.508 | 0.759 | 0.792 | 0.715 | 0.847 | 0.612 | 0.575 | 0.547 | **0.993** |
| | 100% | Conv | 23.42 | 30.67 | 34.21 | 16.9 | 27.06 | 27.56 | 24.35 | 23.31 | **43.83** |
| | | C/B | 0.448 | 0.692 | 0.762 | 0.572 | 0.820 | 0.621 | 0.543 | 0.510 | **0.992** |
| | 125% | Conv | 26.88 | 36.81 | 38.69 | 23.75 | 32.06 | 33.23 | 29.17 | 28.9 | **51.27** |
| | | C/B | 0.402 | 0.640 | 0.737 | 0.632 | 0.772 | 0.606 | 0.485 | 0.470 | **0.990** |
| | 150% | Conv | 30.17 | 39.04 | 42.94 | 26.19 | 37.02 | 35.35 | 30.81 | 30.38 | **58.96** |
| | | C/B | 0.348 | 0.591 | 0.715 | 0.610 | 0.761 | 0.559 | 0.426 | 0.416 | **0.989** |

approaches (DT, GAVE[1]) outperformed imitation methods (BC), as both categories incorporate guidance mechanisms during training or inference. However, imitation methods cannot intrinsically discern suboptimal actions, resulting in indiscriminate behavioral replication without value discrimination.

**Constraint Awareness.** Our method demonstrates superior constraint-awareness in dynamically adapting coefficients to varying budget and ROI configurations, as visualized in Fig. 3. While data-driven approaches (BC, IQL) exhibit limited adjustment responsiveness due to the limitation of the optimization algorithm, model-based methods (LP) fail to achieve competitive scores despite detecting constraint shift due to the lack of environmental dynamics modeling capabilities for fine adaptation.

**Coefficient Adjustment.** As shown in Fig. 4, our method achieves precise $\alpha$ prediction ($\beta^{opt}/\beta^*$), enabling proactive bid suppression during early auction stages to ensure terminal constraint satisfaction. In contrast, baseline approaches exhibit delayed response, resulting in violating the ROI constraint. This occurs because late-period corrective adjustments cannot offset the irrecoverable ROI deficit accumulated during initial phases. This mechanistic divergence manifests quantitatively in $\text{ROI}_{realized}/\text{ROI}_{target}$ ra-

---

[1]GAVE requires 400k+ trajectories in original work, due to training time constraint, our dataset size is insufficient to meet its training requirements.
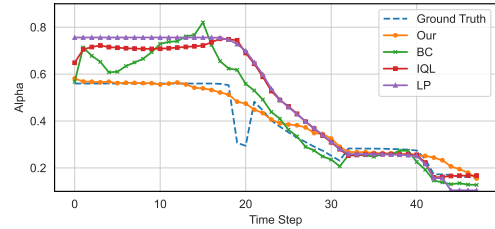


Figure 4: During a bidding period, the variation in $\alpha$ (defined in Lemma 4).

tios: Our: 1.0204; BC: 0.9265; IQL: 0.7657; LP: 0.7632.

## Conclusions

This paper introduces a theoretically grounded framework for multi-scale constrained advertising bidding. We first derive the optimality conditions under budget constraints, which inspires the development of the Hindsight Mechanism. Subsequently, we establish that the adjustment factor $\alpha$ for ROI constraints depends critically on the continuous relationship of cost and value. This insight motivates us to use the parameterized B-spline to model continuous value-/cost relationships. The resulting framework accurately predicts both base bidding coefficients and adjustment factors to satisfy budget and ROI constraints simultaneously. Comprehensive experiments demonstrate significant improvements over state-of-the-art methods on advertising datasets.

# References

Andrychowicz, M.; Wolski, F.; Ray, A.; Schneider, J.; Fong, R.; Welinder, P.; McGrew, B.; Tobin, J.; Pieter Abbeel, O.; and Zaremba, W. 2017. Hindsight experience replay. *Advances in neural information processing systems*, 30.

Bond-Taylor, S.; Leach, A.; Long, Y.; and Willcocks, C. G. 2021. Deep generative modelling: A comparative review of vaes, gans, normalizing flows, energy-based and autoregressive models. *IEEE transactions on pattern analysis and machine intelligence*, 44(11): 7327–7347.

Chen, L.; Lu, K.; Rajeswaran, A.; Lee, K.; Grover, A.; Laskin, M.; Abbeel, P.; Srinivas, A.; and Mordatch, I. 2021. Decision transformer: Reinforcement learning via sequence modeling. *Advances in neural information processing systems*, 34: 15084–15097.

Dantzig, G. B. 2002. Linear programming. *Operations research*, 50(1): 42–47.

Edelman, B.; Ostrovsky, M.; and Schwarz, M. 2007. Internet advertising and the generalized second-price auction: Selling billions of dollars worth of keywords. *American economic review*, 97(1): 242–259.

Finn, C.; Abbeel, P.; and Levine, S. 2017. Model-agnostic meta-learning for fast adaptation of deep networks. In *International conference on machine learning*, 1126–1135. PMLR.

Floudas, C. A.; and Lin, X. 2005. Mixed integer linear programming in process scheduling: Modeling, algorithms, and applications. *Annals of Operations Research*, 139(1): 131–162.

Gao, J.; Li, Y.; Mao, S.; Jiang, P.; Jiang, N.; Wang, Y.; Cai, Q.; Pan, F.; Jiang, P.; Gai, K.; et al. 2025. Generative Auto-Bidding with Value-Guided Explorations. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 244–254.

Gartner. 2024. Performance Marketing and Brand Building. https://www.rainforgrowth.com/insights-updates/performance-marketing-and-brand-building-whats-the-right-ratio-for-growth/. Accessed: 2025-07-30.

Ghorbani, A.; Abid, A.; and Zou, J. 2019. Interpretation of neural networks is fragile. In *Proceedings of the AAAI conference on artificial intelligence*, volume 33, 3681–3688.

Google. 2025. Google Ad. https://ads.google.com/. Accessed: 2025-07-30.

Guo, J.; Huo, Y.; Zhang, Z.; Wang, T.; Yu, C.; Xu, J.; Zheng, B.; and Zhang, Y. 2024. Generative auto-bidding via conditional diffusion modeling. In *Proceedings of the 30th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 5038–5049.

He, Y.; Chen, X.; Wu, D.; Pan, J.; Tan, Q.; Yu, C.; Xu, J.; and Zhu, X. 2021. A unified solution to constrained bidding in online display advertising. In *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*, 2993–3001.

Hospedales, T.; Antoniou, A.; Micaelli, P.; and Storkey, A. 2021. Meta-learning in neural networks: A survey. *IEEE transactions on pattern analysis and machine intelligence*, 44(9): 5149–5169.

Hussein, A.; Gaber, M. M.; Elyan, E.; and Jayne, C. 2017. Imitation learning: A survey of learning methods. *ACM Computing Surveys (CSUR)*, 50(2): 1–35.

Johnson, M. A.; and Moradi, M. H. 2005. *PID control*. Springer.

Kaelbling, L. P.; Littman, M. L.; and Moore, A. W. 1996a. Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4: 237–285.

Kaelbling, L. P.; Littman, M. L.; and Moore, A. W. 1996b. Reinforcement learning: A survey. *Journal of artificial intelligence research*, 4: 237–285.

Kostrikov, I.; Nair, A.; and Levine, S. 2021. Offline reinforcement learning with implicit q-learning. *arXiv preprint arXiv:2110.06169*.

Lillicrap, T. P.; Hunt, J. J.; Pritzel, A.; Heess, N.; Erez, T.; Tassa, Y.; Silver, D.; and Wierstra, D. 2015. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*.

Liu, Z.; Wang, Y.; Vaidya, S.; Ruehle, F.; Halverson, J.; Soljačić, M.; Hou, T. Y.; and Tegmark, M. 2024. Kan: Kolmogorov-arnold networks. *arXiv preprint arXiv:2404.19756*.

Morari, M.; and Lee, J. H. 1999. Model predictive control: past, present and future. *Computers & chemical engineering*, 23(4-5): 667–682.

Prautzsch, H.; Boehm, W.; and Paluszny, M. 2002. *Bézier and B-spline techniques*. Springer Science & Business Media.

Su, K.; Huo, Y.; Zhang, Z.; Dou, S.; Yu, C.; Xu, J.; Lu, Z.; and Zheng, B. 2024. Auctionnet: A novel benchmark for decision-making in large-scale games. *Advances in Neural Information Processing Systems*, 37: 94428–94452.

Swokowski, E. W. 1979. *Calculus with analytic geometry*. Taylor & Francis.

TikTok. 2025. TikTok Ad. https://tiktokfor.business/. Accessed: 2025-07-30.

Torabi, F.; Warnell, G.; and Stone, P. 2018. Behavioral cloning from observation. *arXiv preprint arXiv:1805.01954*.

Vakratsas, D.; and Ambler, T. 1999. How advertising works: what do we really know? *Journal of marketing*, 63(1): 26–43.

Wang, H.; Du, C.; Fang, P.; Yuan, S.; He, X.; Wang, L.; and Zheng, B. 2022. ROI-constrained bidding via curriculum-guided Bayesian reinforcement learning. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 4021–4031.

Wu, D.; Chen, X.; Yang, X.; Wang, H.; Tan, Q.; Zhang, X.; Xu, J.; and Gai, K. 2018. Budget constrained bidding by model-free reinforcement learning in display advertising. In *Proceedings of the 27th ACM International Conference on Information and Knowledge Management*, 1443–1451.

Yang, X.; Li, Y.; Wang, H.; Wu, D.; Tan, Q.; Xu, J.; and Gai, K. 2019. Bid optimization by multivariable control in display advertising. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, 1966–1974.

Ypma, T. J. 1995. Historical development of the Newton–Raphson method. *SIAM review*, 37(4): 531–551.

Yuan, S.; Wang, J.; and Zhao, X. 2013. Real-time bidding for online advertising: measurement and analysis. In *Proceedings of the seventh international workshop on data mining for online advertising*, 1–8.

Zhang, L.; Zhang, Q.; Shen, L.; Yuan, B.; Wang, X.; and Tao, D. 2023. Evaluating model-free reinforcement learning toward safety-critical tasks. In *Proceedings of the AAAI conference on artificial intelligence*, volume 37, 15313–15321.

Zhou, Y.; Chakrabarty, D.; and Lukose, R. 2008. Budget constrained bidding in keyword auctions and online knapsack problems. In *Proceedings of the 17th international conference on world wide web*, 1243–1244.

## The Algorithm of obtaining $V^{oracle}$

Consider the following ad examples:

Table 2: Ad Selection Example ($B = 4$, $r_{\text{target}} = 1.0$)

| Opportunity | $v_i$ | $c_i$ | Efficiency ($v_i/c_i$) |
|---|---|---|---|
| A | 4 | 3 | 1.33 |
| B | 3 | 2 | 1.50 |
| C | 2 | 1 | 2.00 |

The fixed coefficient greedy strategy selects ads in descending efficiency order:

1. Select C: cost=1, value=2, remaining budget=3

2. Select B: cost=2, value=3, total value=5, remaining budget=1

3. Skip A (cost=3 > remaining budget=1) **Final value: 5**

The globally optimal solution combines different ads: Select A+C: cost=3+1=4, value=4+2=6 **Optimal value: 6**

The greedy approach fails because: 1. **Budget fragmentation**: Selecting B first leaves insufficient budget for higher-value A 2. **Combinatorial effects**: A+C synergy (value=6) exceeds B+C (value=5) 3. **Efficiency deception**: A's lower efficiency (1.33) hides its value potential in combination

So the fixed coefficient strategy can not always obtain the global optimal value, we give the global optimal solution (Mixed Integer Linear Programming) for the constrained bidding problem in Algo. 3 However, this multi-variable algorithm is extremely inefficient.

## Proof of Lemma 1

**Lemma 5** *For any fixed coefficient bidding strategy with optimal coefficient $\beta^*$, the achieved total value $V^{fixed}$ satisfies:*

$$V^{fixed} > V^{oracle} - v_{\max} \tag{1}$$

*where $V^{oracle}$ is the global optimal value satisfying both budget and ROI constraints, and $v_{\max} = \max_{i \in I} v_i$ is the maximum value of any single impression opportunity.*

---

Algorithm 3: MILP Solver for $V^{\text{oracle}}$

---

**Require:** Value vector $\mathbf{v}$, cost vector $\mathbf{c}$, budget $B$, target ROI $r_{\text{target}}$
**Ensure:** Optimal value $V^*$, binary selection vector $\mathbf{x}^* \in \{0,1\}^n$

1: **function** MILP($\mathbf{v}, \mathbf{c}, B, r_{\text{target}}$)
2:     **Define binary variables:** $x_i \in \{0,1\}$   $\forall i = 1, \ldots, n$
3:     **Objective:** max $\sum_{i=1}^{n} v_i x_i$
4:     **Subject to:**
5:         $\sum_{i=1}^{n} c_i x_i \leq B$       ▷ Budget constraint
6:         $\sum_{i=1}^{n} (v_i - r_{\text{target}} c_i) x_i \geq 0$   ▷ ROI constraint
7:     **Solve MILP** using branch-and-cut
8:     **return** $V^* = $ objval, $\mathbf{x}^* = [x_i]$
9: **end function**

---

**Proof** *Define the fractional relaxation problem:*

$$\max \sum v_i y_i$$
$$s.t. \sum c_i y_i \leq B$$
$$\sum (v_i - r_{target} c_i) y_i \geq 0$$
$$0 \leq y_i \leq 1$$

*has optimal value $V^{frac}$. Since any integer solution is feasible for this relaxation ($x_i \in \{0,1\} \subset [0,1]$), we have:*

$$V^{oracle} \leq V^{frac} \tag{2}$$

*This follows from the fundamental principle of linear programming relaxation: expanding the feasible region cannot decrease the optimal value.*

*The fractional optimum is achieved by:*

$$V^{frac} = \sum_{i=1}^{k} v_i + r \cdot v_{k+1} \tag{3}$$

*where:*

- *Ads are sorted by efficiency $e_i = v_i/c_i$ descending*
- *$k$ is the largest integer satisfying:*

$$\sum_{i=1}^{k} c_i \leq B \quad and \quad \frac{\sum_{i=1}^{k} v_i}{\sum_{i=1}^{k} c_i} \geq r_{target} \tag{4}$$

- *$r \in (0,1)$ is the solution to:*

$$\sum_{i=1}^{k} c_i + r \cdot c_{k+1} = B \tag{5}$$

$$\frac{\sum_{i=1}^{k} v_i + r \cdot v_{k+1}}{\sum_{i=1}^{k} c_i + r \cdot c_{k+1}} \geq r_{target} \tag{6}$$

*This structure is optimal because:*

*1. Sorting by $e_i$ maximizes value per cost*

2. *Full selection of first k ads uses the most efficient impressions*

3. *Fractional selection of ad $k + 1$ exactly exhausts budget while maintaining ROI*

*The fixed coefficient strategy:*

1. *Sort ads by $e_i$ descending*

2. *Select maximal set $S = \{1, \ldots, m\}$ satisfying:*

$$\sum_{i \in S} c_i \leq B \quad and \quad \frac{\sum_{i \in S} v_i}{\sum_{i \in S} c_i} \geq r_{target} \quad (7)$$

*yields value $V^{fixed} = \sum_{i=1}^{m} v_i$. By construction, $m \leq k$ because: the fixed coefficient strategy selects a subset of the fractionally optimal ads: $S \subseteq \{1, \ldots, k+1\}$*

*Suppose ad $j > k+1$ is selected. Since $e_j \leq e_{k+1} \leq \cdots \leq e_1$, adding $j$ would:*

- *Violate ROI if $e_j < r_{target}$ (cumulative ROI decreases)*
- *Or be suboptimal (more efficient ads available)*

*Thus the strategy only selects from $\{1, \ldots, k+1\}$.*

*The value difference is:*

$$V^{frac} - V^{fixed} = \left( \sum_{i=1}^{k} v_i + r \cdot v_{k+1} \right) - \sum_{i=1}^{m} v_i \quad (8)$$

*Case analysis:*

- *If $m = k$ (strategy takes all full ads):*

$$V^{frac} - V^{fixed} = r \cdot v_{k+1}$$

- *If $m < k$ (ROI constraint prevents full selection):*

$$V^{frac} - V^{fixed} = \sum_{i=m+1}^{k} v_i + r \cdot v_{k+1} \quad (9)$$

$$< \sum_{i=m+1}^{k+1} v_i \quad (since \ r < 1) \quad (10)$$

*In both cases:*

$$V^{frac} - V^{fixed} < v_{k+1} + \sum_{i=m+1}^{k} v_i \quad (11)$$

$$\leq (k - m + 1) \cdot \max_{i \geq m+1} v_i \quad (12)$$

$$\leq n \cdot v_{max} \quad (worst\text{-}case) \quad (13)$$

*But tighter bound exists: $V^{frac} - V^{fixed} < v_{max}$*

*The strategy selects maximal $m$ satisfying constraints. If $m < k$, then:*

$$\frac{\sum_{i=1}^{m} v_i}{\sum_{i=1}^{m} c_i} \geq r_{target} > \frac{\sum_{i=1}^{m} v_i + v_j}{\sum_{i=1}^{m} c_i + c_j} \quad \forall j > m$$

*which implies $v_j < r_{target} c_j \leq v_{max}$. Thus:*

$$V^{frac} - V^{fixed} \leq \max_{j>m} v_j < v_{max}$$

*Chaining inequalities:*

$$V^{oracle} \leq V^{frac} \quad (14)$$

$$V^{frac} < V^{fixed} + v_{max} \quad (15)$$

*Thus:*

$$V^{oracle} < V^{fixed} + v_{max} \quad (16)$$

*Rearranging gives the guarantee:*

$$\boxed{V^{fixed} > V^{oracle} - v_{max}} \quad (17)$$

*The bound is tight when:*

- *$r \to 1^-$ and $v_{k+1} = v_{max}$*
- *Strategy selects $m = k$ (all full ads)*
- *Fractional solution gains $v_{max}$ while integer cannot*

*In advertising systems, $v_{max} \ll V^{oracle}$ (typically $v_{max}/V^{oracle} < 0.01$).*

# Proof of Lemma 2

**Lemma 6 (Optimality Condition for FCS Under BCB)**
*Under budget-constrained bidding with fixed coefficient Strategy, at any decision step $\tau$ with remaining budget $B_\tau$, if $\exists \beta^*$ such that*

$$\sum_{t=\tau}^{T} C_t(\beta^*) = B_\tau,$$

*then $\beta^*$ is the optimal coefficient on $[\tau, T]$.*

**Proof** *Define the efficiency metric $e_i = v_i/c_i$.*
*Consider the bid condition for winning opportunity $i$:*

$$\text{bid}_i \geq c_i \iff \beta \geq \frac{c_i}{v_i} = \frac{1}{e_i} \quad (18)$$

*Sort opportunities by decreasing efficiency: $e_1 \geq e_2 \geq \cdots \geq e_n$. A greedy algorithm selects opportunities in this order until the budget is exhausted. Let $k$ be the last selected opportunity satisfying:*

$$\sum_{i=1}^{k} c_i \leq B < \sum_{i=1}^{k+1} c_i \quad (19)$$

*The efficiency threshold is $e^* = e_k$. The optimal coefficient is:*

$$\beta^* = \frac{1}{e^*} = \frac{c_k}{v_k} \quad (20)$$

*With $\beta^*$, the advertiser wins all opportunities where:*

$$e_i \geq e^* \iff \frac{1}{e_i} \leq \frac{1}{e^*} \iff \beta^* \geq \frac{c_i}{v_i} \quad (21)$$

*The total cost is exactly $B$ by selection criterion.*
*Now consider any alternative coefficient $\beta' \neq \beta^*$:*
*Case 1: $\beta' > \beta^*$*
*The winning set under $\beta'$ expands to include inefficient opportunities:*

$$S' = \{i \mid e_i \geq 1/\beta'\} \quad (22)$$

$$S^* = \{i \mid e_i \geq e_{\min} = 1/\beta^*\} \quad (23)$$

$$S' \setminus S^* = \{j \mid 1/\beta' \leq e_j < e_{\min}\} \neq \emptyset \quad (24)$$

*The set $S' \setminus S^*$ is non-empty because $\beta' > \beta^* \implies 1/\beta' < 1/\beta^* = e_{\min}$.*

*Total cost exceeds budget:*

$$\sum_{i \in S'} c_i = \sum_{i \in S^*} c_i + \sum_{j \in S' \setminus S^*} c_j \tag{25}$$

$$= B + \underbrace{\sum_{j \in S' \setminus S^*} c_j}_{\epsilon > 0} > B \tag{26}$$

*This violates the budget constraint.*
*For any $j \in S' \setminus S^*$:*

$$v_j = e_j c_j < e_{\min} c_j = \frac{c_j}{\beta^*} \tag{27}$$

*The resources spent on $j$ ($c_j$) could have been allocated to more efficient opportunities in $S^*$. The marginal opportunity cost is:*

$$\Delta v_j = \max_{k \notin S'} \left( e_k - e_j \right) c_j \geq (e_{\min} - e_j) c_j > 0 \tag{28}$$

*since $e_k \geq e_{\min} > e_j$ for $k \in S^* \setminus S'$ (if any).*
*In real bidding systems, budget overrun causes:*

- *Auction disqualification for later opportunities*
- *Actual realized value $V_{actual} \leq \sum_{i \in S'} v_i - \sum_{j \in rejected} v_j$*

*Let $R \subseteq S'$ be the set of auctions rejected due to budget overrun. Then:*

$$V_{actual} = \sum_{i \in S' \setminus R} v_i \tag{29}$$

$$\leq \sum_{i \in S^*} v_i \tag{30}$$

*with strict inequality when $R \neq \emptyset$.*
*The inclusion of inefficient opportunities reduces overall value density:*

$$\frac{\sum_{i \in S'} v_i}{\sum_{i \in S'} c_i} < \frac{\sum_{i \in S^*} v_i}{B} \tag{31}$$

$$\text{since} \quad \sum_{j \in S' \setminus S^*} v_j < \frac{\sum_{j \in S' \setminus S^*} c_j}{B} \sum_{i \in S^*} v_i \tag{32}$$

*This follows from $e_j < e_{\min} \leq avg(e_i)_{i \in S^*}$.*
*Suppose $S'$ were optimal. But:*

$$\sum_{i \in S^*} v_i > \sum_{i \in S'} v_i - \epsilon \cdot v_{\max} \tag{33}$$

$$\geq V_{actual} \quad \text{(for small $\epsilon$)} \tag{34}$$

*contradicting optimality. Thus $S'$ cannot be optimal.*
*Case 2: $\beta' < \beta^*$*
*This raises the winning threshold to $1/\beta' > e_{\min}$. The advertiser loses some efficient opportunities in $\{1, \ldots, k\}$:*

$$\sum_{i \in S'} c_i < B \quad \text{(budget underutilization)} \tag{35}$$

$$\sum_{i \in S'} v_i < \sum_{i=1}^{k} v_i \quad \text{(since $S' \subset \{1, \ldots, k\}$)} \tag{36}$$

*where $S'$ is the winning set under $\beta'$.*
*Therefore $\beta^*$ uniquely satisfies both budget exhaustion and value maximization.*

## Proof of Lemma 3

**Lemma 7 (ROI Monotonicity)** *The expected ROI function $R(\beta)$ is non-increasing in $\beta$.*

**Proof** *Consider a fixed sample path (i.e., fixed set of projects with parameters). Sort projects by efficiency ratio $e_i = v_i/c_i$ in descending order: $e_{(1)} \geq e_{(2)} \geq \cdots \geq e_{(n)}$, with corresponding minimal coefficients $\beta_{(1)} \leq \beta_{(2)} \leq \cdots \leq \beta_{(n)}$. The decision rule is $x_i(\beta) = \mathbb{1}_{\{\beta \geq \beta_i\}}$. Divide the $\beta$-axis into intervals:*

- *For $\beta < \beta_{(1)}$, no projects are selected.*
- *For $\beta \in [\beta_{(k)}, \beta_{(k+1)})$, the first $k$ projects are selected.*
- *For $\beta \geq \beta_{(n)}$, all $n$ projects are selected.*

*In each interval $[\beta_{(k)}, \beta_{(k+1)})$, the ratio $R(\beta)$ is constant:*

$$R_k = \frac{\sum_{i=1}^{k} v_{(i)}}{\sum_{i=1}^{k} c_{(i)}} \tag{37}$$

*We prove that $\{R_k\}$ is non-increasing: $R_k \geq R_{k+1}$ for all $k$.*
*Base case ($k = 1$):*

$$R_1 = \frac{v_{(1)}}{c_{(1)}} = e_{(1)}, \quad R_2 = \frac{v_{(1)} + v_{(2)}}{c_{(1)} + c_{(2)}} \tag{38}$$

*Since $e_{(1)} \geq e_{(2)}$:*

$$v_{(1)} c_{(2)} \geq v_{(2)} c_{(1)} \implies v_{(1)}(c_{(1)} + c_{(2)}) \geq (v_{(1)} + v_{(2)})c_{(1)}$$
$$\implies R_1 \geq R_2 \tag{39}$$

*General case: For $k \geq 1$,*

$$R_{k+1} = \frac{V_k + v_{(k+1)}}{C_k + c_{(k+1)}}, \quad \text{where} \quad V_k = \sum_{i=1}^{k} v_{(i)}, \quad C_k = \sum_{i=1}^{k} c_{(i)} \tag{40}$$

*Since $e_{(k+1)} \leq \min_{i=1}^{k} e_{(i)} \leq R_k$:*

$$\frac{v_{(k+1)}}{c_{(k+1)}} \leq \frac{V_k}{C_k} \implies v_{(k+1)} C_k \leq V_k c_{(k+1)} \tag{41}$$

*Then:*

$$v_{(k+1)} C_k + V_k C_k \leq V_k c_{(k+1)} + V_k C_k$$
$$\implies (V_k + v_{(k+1)}) C_k \leq V_k (C_k + c_{(k+1)}) \tag{42}$$

*Thus:*

$$\frac{V_k + v_{(k+1)}}{C_k + c_{(k+1)}} \leq \frac{V_k}{C_k} \implies R_{k+1} \leq R_k \tag{43}$$

*Therefore, $R(\beta)$ is non-increasing on each sample path. By the monotonicity of expectation, $\mathbb{E}[R(\beta)]$ is non-increasing in $\beta$.*

## Proof of Lemma 4

**Lemma 8 (MCB Coefficient Adjustment)** *Under budget and ROI constraint with fixed coefficient strategy, at any decision step $\tau$ with history cost $C_1^{\tau-1}$ and history value $V_1^{\tau-1}$, if $\exists \alpha \leq 1$ such that $\beta_{opt} = \alpha \beta^*$ applied on $[\tau, T]$ can make the total realized ROI equal to ROI constraint, then*

$$\beta^* \int_{\alpha}^{1} \left[ \gamma \left( \beta^* u \right) - r_{\text{target}} \eta \left( \beta^* u \right) \right] du = \Delta_{\text{ROI}}$$

where $\gamma = \frac{\partial V}{\partial \beta}$ is derivative of value with respect to coefficient, $\eta = \frac{\partial C}{\partial \beta}$ is derivative of cost with respect to coefficient, $\Delta_{ROI} = V_1^{\tau-1} + V_\tau^T(\beta^*) - r_{target}(C_1^{\tau-1} + C_\tau^T(\beta^*))$ with $\beta^*$ the reference coefficient.

**Proof** Let $V(\beta) = V_\tau^T(\beta)$ and $C(\beta) = C_\tau^T(\beta)$ denote future value and cost functions respectively. Define the net value function:

$$f(\beta) = V(\beta) - r_{\text{target}}C(\beta) \tag{44}$$

The ROI constraint requires:

$$V_1^{\tau-1} + V(\alpha\beta^*) = r_{\text{target}}(C_1^{\tau-1} + C(\alpha\beta^*)) \tag{45}$$

Which simplifies to:

$$f(\alpha\beta^*) = k \quad where \quad k = r_{\text{target}}C_1^{\tau-1} - V_1^{\tau-1} \tag{46}$$

By the Fundamental Theorem of Calculus:

$$f(\alpha\beta^*) - f(\beta^*) = \int_{\beta^*}^{\alpha\beta^*} \frac{df}{d\beta}d\beta \tag{47}$$

Where:

$$\frac{df}{d\beta} = \frac{\partial V}{\partial \beta} - r_{\text{target}}\frac{\partial C}{\partial \beta} = \gamma - r_{\text{target}}\eta \tag{48}$$

Substituting the constraint:

$$k - f(\beta^*) = \int_{\beta^*}^{\alpha\beta^*} \left[\gamma(\beta) - r_{\text{target}}\eta(\beta)\right] d\beta \tag{49}$$

Substitute $\beta = \beta^* u$, $d\beta = \beta^* du$:

$$\int_{\beta^*}^{\alpha\beta^*} \left[\gamma(\beta) - r_{\text{target}}\eta(\beta)\right] d\beta \tag{50}$$

$$= \int_1^\alpha \left[\gamma(\beta^* u) - r_{\text{target}}\eta(\beta^* u)\right] \beta^* du \tag{51}$$

$$= \beta^* \int_1^\alpha \left[\gamma(\beta^* u) - r_{\text{target}}\eta(\beta^* u)\right] du \tag{52}$$

Note that:

$$k - f(\beta^*) = (r_{\text{target}}C_1^{\tau-1} - V_1^{\tau-1}) - (V^* - r_{\text{target}}C^*) \tag{53}$$

$$= -[(V_1^{\tau-1} + V^*) - r_{\text{target}}(C_1^{\tau-1} + C^*)] \tag{54}$$

$$= -\Delta_{\text{ROI}} \tag{55}$$

Therefore:

$$\beta^* \int_1^\alpha \left[\gamma(\beta^* u) - r_{\text{target}}\eta(\beta^* u)\right] du = -\Delta_{\text{ROI}} \tag{56}$$

Reverse limits using $\int_a^b = -\int_b^a$:

$$\boxed{\beta^* \int_\alpha^1 \left[\gamma(\beta^* u) - r_{\text{target}}\eta(\beta^* u)\right] du = \Delta_{\text{ROI}}} \tag{57}$$

## Proof of Theorem 1

**Theorem 2 (Optimal Bidding Strategy)** *At any step $\tau$, with remaining budget $B_\tau$, target ROI $r_{target}$, historical cost $C_1^{\tau-1}$, historical value $V_1^{\tau-1}$, solve for $\beta_\tau = f_\theta(B_\tau)$, $V_\tau^T = f_\phi(B_\tau)$ and calculate $\Delta_{ROI} = V_1^{\tau-1} + V_\tau^T - r_{target}(C_1^{\tau-1} + B_\tau)$*

- *if $\Delta_{ROI} \geq 0$, $\beta_\tau$ is the optimal bidding coefficient.*
- *if $\Delta_{ROI} < 0$, $\alpha\beta_\tau$ is the optimal bidding coefficient.*

**Proof Case 1:** $\Delta_{\text{ROI}} \geq 0$
*By Lemma 1, $\beta^*$ satisfies:*

$$C_\tau^T(\beta^*) = B_\tau \tag{58}$$

*and by $\Delta_{ROI} \geq 0$:*

$$V_h + V_\tau^T(\beta^*) \geq r_{target}(C_h + C_\tau^T(\beta^*)) \tag{59}$$

*where $V_h = V_1^{\tau-1}$, $C_h = C_1^{\tau-1}$.*
 *Optimality follows from:*

1. **Budget feasibility**: $C_\tau^T(\beta^*) = B_\tau$ *(Lemma 1)*
2. **ROI satisfaction**: $\Delta_{ROI} \geq 0$
3. **Value maximization**: $(V_\tau^T)'(\beta) > 0$

**Case 2:** $\Delta_{\text{ROI}} < 0$
*By Lemma 4, $\exists \alpha \in (0,1)$ such that:*

$$\beta^* \int_\alpha^1 \left[\gamma(\beta^* u) - r_{\text{target}}\eta(\beta^* u)\right] du = -\Delta_{\text{ROI}} \tag{60}$$

*where $\gamma = (V_\tau^T)'$, $\eta = (C_\tau^T)'$. This implies:*

$$V_h + V_\tau^T(\alpha\beta^*) = r_{target}(C_h + C_\tau^T(\alpha\beta^*)) \tag{61}$$

 *Since $\alpha < 1$ and $(C_\tau^T)' > 0$:*

$$C_\tau^T(\alpha\beta^*) < C_\tau^T(\beta^*) = B_\tau \tag{62}$$

*Thus $C_h + C_\tau^T(\alpha\beta^*) \leq C_h + B_\tau$ (total budget).*
 *We prove $\alpha\beta^*$ maximizes value under dual constraints.*
 *From the adjustment equation:*

$$\frac{d}{d\alpha} \left[V_\tau^T(\alpha\beta^*) - r_{target}C_\tau^T(\alpha\beta^*)\right] = \beta^* \left(\gamma(\alpha\beta^*) - r_{target}\eta(\alpha\beta^*)\right)$$
$$< 0 \quad \text{(by Lemma 3)}$$

*Thus $V_\tau^T(\beta) - r_{target}C_\tau^T(\beta)$ is decrease in $\beta$. Therefore, for fixed ROI constraint, $\beta = \alpha\beta^*$ is unique.*
 *Consider the value:*

$$\frac{dV_\tau^T}{d\beta} \geq 0$$

*Thus $\alpha\beta^*$ maximizes value along the ROI constraint boundary.*
 *$\alpha\beta^*$ is the unique solution that: 1. Satisfies both constraints (budget and ROI) 2. Maximizes value on the ROI constraint boundary*
 *Thus it is globally optimal.*

# State Feature

Table 3: State Feature Vector

| Category | Feature Description |
|---|---|
| Temporal Context | Day of the week; Time left in campaign (steps) |
| Advertiser Identity | Advertiser ID; Advertiser category |
| Current Auction Metrics | Current pValue mean; Current pValue percentiles (10th,25th,40th,55th,70th,85th); Current pValue count |
| Historical pValues | Historical mean; Last 1-step: mean,10th,50th,90th percentile; Last 3-step: mean,10th,50th,90th percentile; All-time: 10th,50th,90th percentile |
| Winning Cost Metrics | Historical mean; Last 1-step: mean,10th,50th,90th percentile; Last 3-step: mean,10th,50th,90th percentile; All-time: 10th,50th,90th percentile |
| Volume Statistics | Historical total pValue count; Last 1-step pValue count; Last 3-step pValue count |

# Hyperparameters

Table 4: Model Hyperparameters

| Hyperparameter | Description | Value |
|---|---|---|
| $k$ | B-spline degree | 3 |
| num | Number of grid points | 16 |
| hidden_size | Hidden layer dimension | 128 |
| input_dim | input dimension | 39 |