
DUO: Diffusion Models for Universal Offline Black-Box Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Offline black-box optimization aims to find high-performing designs from a fixed
2 dataset without online evaluations, with applications spanning protein design, ma-
3 terials discovery, and robot learning. Most existing methods are typically designed
4 for a single task with fixed dimensionality, leaving universal offline optimization—
5 learning one shared model across heterogeneous search spaces, mixed variable
6 types, and scarce-data transfer settings—largely unresolved. In this paper, we study
7 universal offline optimization from a generative perspective for the first time and
8 propose DUO (Diffusion model for Universal offline black-box Optimization),
9 which bridges universal representation learning with trajectory-level diffusion
10 modeling. DUO uses a Transformer-based variational autoencoder to embed both
11 continuous and discrete designs into a shared latent space, avoiding task-specific
12 architectures for incompatible native domains. Within this unified space, we syn-
13 thesize optimization-oriented trajectories and train a conditional diffusion model,
14 with task-level semantics injected through frozen text-metadata embeddings and
15 classifier-free guidance. A cross-entropy consistency term further aligns continuous
16 training with discrete evaluation objectives. Evaluated on the Design-Bench and
17 SOO-Bench benchmarks, DUO demonstrates strong performance across diverse
18 continuous and discrete tasks under multitask joint training. Our experiments
19 highlight its robust zero-shot and few-shot transfer capabilities, suggesting that
20 metadata-aware latent trajectory diffusion provides a highly effective framework
21 for universal offline black-box optimization.

22 1 Introduction

23 Black-box optimization (BBO) seeks high-performing designs when the objective is unknown
24 and can only be observed through input–output evaluations. Such problems arise in high-cost
25 domains including biological engineering, materials screening, robot learning, and chip parameter
26 tuning [33, 13]. Classical tools such as Bayesian optimization [8] remain effective when the optimizer
27 can query the objective repeatedly, but this assumption breaks down in wet-lab pipelines, large-scale
28 simulations, and other settings where each measurement is expensive. In these regimes one must
29 instead optimize from a *fixed* historical dataset with no online queries during training, the setting
30 known as offline BBO [33].

31 Most offline BBO methods, however, follow a single-task regime: each model is tied to one search
32 space, one dimensionality, and one task-specific dataset. This is restrictive because real tasks often
33 provide only a small number of labeled examples, or even metadata alone, while related tasks may
34 share useful structure [21]. Universal offline BBO addresses this gap by training one shared model
35 across heterogeneous objectives, dimensions, and variable types, typically by serializing designs and
36 conditioning on natural-language metadata so that one predictor can cover many tasks [30, 29]. The

37 result is cross-task knowledge transfer in zero-shot and few-shot regimes that single-task models
38 cannot reach.

39 Yet the optimization interface in these universal methods remains forward-oriented: the model predicts
40 rewards, and a separate inner-loop search or ranking step is responsible for actually improving
41 designs. UniSO [30] is representative of this pattern—string representations and metadata-guided
42 representation learning unify different search spaces, but optimization still consists of fitting a reward
43 predictor on offline tuples and then searching in model space. Such an interface inherits the well-
44 known fragility of offline reward models on out-of-distribution designs, and the search step itself is
45 hard to share or amortize across tasks whose coordinates and step sizes have no common meaning. A
46 generative *backward* approach offers a natural alternative: directly learn how designs move toward
47 higher scores, rather than scoring them and searching afterwards. GTG [43] demonstrates the strength
48 of this idea in the single-task setting by constructing pseudo optimization paths and modeling them
49 with diffusion.

50 Extending trajectory generation to universal offline BBO, however, cannot be achieved by simply
51 pooling trajectories from different tasks. GTG-style models are trained and sampled in one native
52 design space of fixed shape, whereas universal BBO must handle spaces that differ in variable type,
53 dimensionality, and objective semantics. Because these spaces share no common coordinate system,
54 neighborhoods, noise perturbations, and improvement directions are not comparable across tasks.
55 A universal backward optimizer must therefore (i) expose a common interface for heterogeneous
56 designs, (ii) build improvement trajectories on which “moving forward” has consistent meaning, and
57 (iii) condition generation on the task semantics that define what “improvement” is.

58 We instantiate these requirements in **DUO** (**D**iffusion model for **U**niversal offline black-box
59 **O**ptimization), which addresses these limitations with three matching components: a Transformer-
60 based VAE that maps mixed continuous and discrete designs into a shared latent space, synthetic
61 improvement trajectories constructed in that latent space, and a metadata-conditioned 1D U-Net
62 diffusion model that generates better designs backward along these paths. DUO thus replaces the
63 forward-regression-plus-search interface used by universal forward methods such as UniSO [30]
64 with a single generative backbone, and lifts the single-task trajectory paradigm of GTG [43] to
65 heterogeneous multitask settings.

66 Our contributions are threefold. (1) **A backward formulation of universal offline BBO:** we reframe
67 universal offline BBO as a backward generative problem over improvement trajectories, removing the
68 reward-predictor-plus-inner-search interface that current universal methods inherit from single-task
69 offline BBO. (2) **Metadata-aware latent trajectory diffusion:** we realize this formulation with a
70 Transformer-based VAE shared latent interface, a conditional 1D U-Net over synthetic improvement
71 trajectories, frozen text-metadata guidance [10], and a grouped cross-entropy term for discrete tasks.
72 (3) **Empirical gains over forward baselines:** on Design-Bench and SOO-Bench, DUO improves
73 over the representative universal forward optimizer UniSO and the single-task trajectory baseline
74 GTG under multitask training, while retaining zero-shot and few-shot transfer ability on held-out
75 tasks.

76 2 Related Work

77 **Offline black-box optimization.** Offline BBO seeks high-performing designs from a fixed dataset
78 without online queries [33]. Existing methods are commonly divided into *forward* approaches, which
79 fit a surrogate objective and then search against it, and *backward* approaches, which directly learn to
80 generate promising designs [35]. Representative forward methods include COMs [32], RoMA [40],
81 ICT [41], and RaM-ListNet [31]; representative backward methods include MINs [17], DDOM [15],
82 and RGD [5]. These methods have shown strong single-task results but are tied to one fixed search
83 space and cannot directly reuse experience from other tasks.

84 **Universal and pretrained optimizers.** Universal offline BBO asks one model to operate across
85 tasks with different dimensions, variable types, objectives, and metadata [30]. UniSO [30] is the clos-
86 est universal offline baseline: it serializes heterogeneous designs, attaches natural-language metadata,
87 and trains token-targeted or numeric-targeted regressors with semantic alignment. Related pretrained
88 optimizers and predictors, including ExPT [21], POM [18], OptFormer [6], and OmniPred [29], fur-
89 ther show that synthetic diversity, serialized histories, and language-model embeddings can improve

90 transfer. However, most universal methods are forward types: they learn a reward predictor or ranking
 91 interface and still require a separate search procedure. In particular, UniSO covers heterogeneous
 92 multitask learning with metadata but offers no trajectory- or generation-based optimization interface;
 93 DUO instead keeps the universal, metadata-aware training setting and replaces forward regression
 94 with a backward generative model of improvement trajectories.

95 **Trajectory-based generative optimization.** Trajectory-level methods model how designs move
 96 toward better regions rather than only modeling isolated high-score samples. BONET [16] sorts offline
 97 examples by score, stitches synthetic optimization paths, and trains an autoregressive Transformer.
 98 GTG [43] improves this idea with diffusion and neighborhood-based trajectory construction, reducing
 99 autoregressive error accumulation and providing a strong single-task trajectory baseline. Trajectory
 100 diffusion is also well established in offline RL [12, 2, 19]. The key limitation for universal BBO is
 101 geometric: these models are normally trained in one native design space of fixed shape. BONET
 102 and GTG motivate trajectory-centric generation but remain single-task and lack any mechanism
 103 for joint multitask training, heterogeneous design spaces, or task-metadata conditioning. DUO
 104 extends the trajectory diffusion paradigm by first mapping heterogeneous designs into a shared latent
 105 interface and then conditioning the denoising process on task metadata, thereby combining trajectory-
 106 based backward generation with the multitask, metadata-aware setting that current universal forward
 107 optimizers occupy.

108 **Latent and metadata-conditioned generative interfaces.** DUO combines two auxiliary ideas that
 109 make universal trajectory diffusion practical. First, latent generative modeling uses a compressed
 110 representation as the common optimization space, as in latent diffusion [27], latent Bayesian opti-
 111 mization [34], latent energy-based optimization [39], and transfer optimization across spaces [7].
 112 Second, language-model features provide task-level semantics for optimization, either as Bayesian-
 113 optimization features [23, 22], multitask prediction inputs [29], or fully agentic optimization inter-
 114 faces [20]. DUO uses these ideas in a narrower way: a Transformer-VAE supplies the shared latent
 115 space, while frozen metadata embeddings steer diffusion through classifier-free guidance [10]; the
 116 optimization backbone remains a trajectory denoiser rather than an LLM or a forward surrogate.

117 3 DUO

118 In this section, we introduce DUO, a diffusion-based framework for universal offline black-box
 119 optimization. The main difficulty is not merely how to condition a model on different tasks. A more
 120 basic question is: *what should a single diffusion model denoise when different tasks have different*
 121 *dimensions, variable types, and native geometries?* Directly defining one diffusion process over
 122 raw designs is ill-posed: a robot morphology vector, a material descriptor, a DNA sequence, and a
 123 trajectory-planning vector do not share a common coordinate system.

124 DUO addresses this issue in three steps, as shown in Figure 1. First, it maps heterogeneous designs
 125 into a shared latent interface with a Transformer-VAE. Second, it constructs synthetic improvement
 126 trajectories in this latent space, so that the model learns how designs move from lower-score to
 127 higher-score regions rather than only modeling isolated good samples. Third, it trains a metadata-
 128 conditioned diffusion model over these latent trajectories and decodes the generated latents back to
 129 task-native designs.

130 3.1 Universal offline optimization setting

131 We consider a set of training tasks \mathcal{T}_{tr} . Each task $\tau \in \mathcal{T}_{\text{tr}}$ has a search space \mathcal{X}_τ , an unknown
 132 objective function $f_\tau : \mathcal{X}_\tau \rightarrow \mathbb{R}$, and a fixed offline dataset

$$133 \mathcal{D}_\tau = \{(\mathbf{x}_i^\tau, y_i^\tau)\}_{i=1}^{N_\tau}, \quad y_i^\tau = f_\tau(\mathbf{x}_i^\tau). \quad (1)$$

134 No online query to f_τ is allowed during training. The search spaces \mathcal{X}_τ may be continuous, discrete,
 135 or mixed, and their dimensionalities can be different across tasks. Each task also comes with metadata
 136 m_τ , such as the task name, variable type, description, and optimization objective.

137 The goal is to train one shared generator that can output high-scoring candidates for both seen and
 138 unseen tasks. At deployment time, a new task may provide only metadata $m_{\tau_{\text{new}}}$, corresponding
 to the zero-shot case, or metadata plus a small dataset $\mathcal{D}_{\tau_{\text{new}}}$, corresponding to the few-shot case.

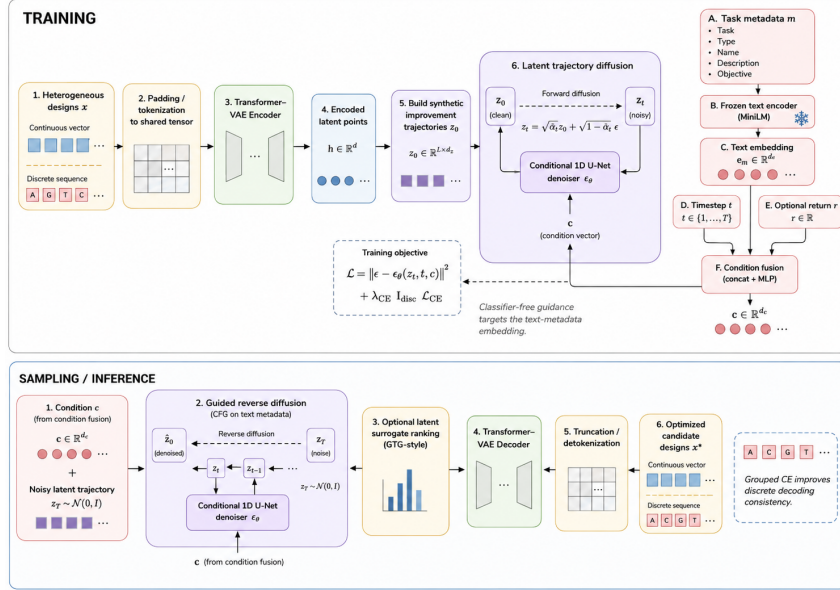


Figure 1: Overview of DUO. **Training:** Heterogeneous continuous/discrete designs are tokenized into masked tensors, encoded by a Transformer-VAE, and stitched into score-improving latent trajectories. A metadata-conditioned 1D U-Net denoises these trajectories, with optional return conditioning and grouped CE for discrete variables. **Sampling:** Guided reverse diffusion generates latent trajectories, which are optionally reranked by an offline surrogate and decoded back to task-native designs.

139 In the few-shot case, DUO can optionally finetune the shared parameters, while keeping the same
 140 architecture and latent interface.

141 3.2 A shared latent interface via VAE

142 A natural but problematic solution is to pad all raw designs to the same length and train diffusion
 143 directly in that padded space. This makes the tensors compatible, but it does not make the geometry
 144 meaningful. For example, Euclidean distance between two padded DNA one-hot tensors is not
 145 comparable to Euclidean distance between two robot morphology vectors. A universal denoiser
 146 trained in such a raw space would therefore mix incompatible notions of locality and smoothness.

147 DUO instead learns a shared latent interface. For each task, we first convert a native design x_i^T into a
 148 unified tensor

$$u_i^T = \Pi_\tau(x_i^T) \in \mathbb{R}^{L \times C}, \quad (2)$$

149 where Π_τ denotes task-specific preprocessing. Continuous vectors are normalized and padded or
 150 truncated to length L with masks. Discrete designs are represented by grouped one-hot tensors, where
 151 each group corresponds to a categorical decision such as a sequence position. The mask records
 152 which entries are valid and which entries are introduced only for padding.

153 A Transformer-VAE then maps u_i^T into a latent vector:

$$q_\phi(z | u) = \mathcal{N}(\mu_\phi(u), \text{diag}(\sigma_\phi^2(u))), \quad z = \mu_\phi(u) + \sigma_\phi(u) \odot \xi, \quad \xi \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (3)$$

154 The decoder reconstructs the unified tensor from the latent code:

$$\hat{u} = \text{Dec}_\psi(z). \quad (4)$$

155 The VAE is trained over all tasks with a masked reconstruction objective and a KL regularizer:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{\tau, i} [\ell_{\text{rec}}(\text{Dec}_\psi(z_i^T), u_i^T) + \beta D_{\text{KL}}(q_\phi(z | u_i^T) \| \mathcal{N}(\mathbf{0}, \mathbf{I}))]. \quad (5)$$

156 Here ℓ_{rec} is evaluated only on valid, non-padded entries. After training, each offline design is
 157 represented by a latent point $z_i^T = \text{Enc}_\phi(u_i^T)$. This gives DUO a common object to model: every
 158 task now provides latent points of the same dimension.

159 **3.3 Constructing latent improvement trajectories**

160 DUO does not train diffusion on independent latent points. Instead, it trains on short trajectories that
 161 imitate the behavior of an optimizer moving toward better regions. This trajectory view is useful
 162 because offline BBO cares about improvement, not just reconstruction of the offline data distribution.

163 For each task τ , we first encode all offline designs into latent points $\{(z_i^\tau, y_i^\tau)\}_{i=1}^{N_\tau}$. We then build
 164 synthetic trajectories of fixed length H by chaining local score-improving transitions. Let $\mathcal{N}_\kappa(i)$ be
 165 the set of κ nearest neighbors of z_i^τ in the latent space. A candidate successor of i must be both close
 166 in latent space and better in objective value:

$$\mathcal{S}_\tau(i) = \{j \in \mathcal{N}_\kappa(i) : \tilde{y}_j^\tau > \tilde{y}_i^\tau + \epsilon_{\text{traj}}\}, \quad (6)$$

167 where \tilde{y} denotes the normalized score and ϵ_{traj} filters out transitions with negligible improvement.
 168 Starting from an offline point, we repeatedly sample a successor from $\mathcal{S}_\tau(i)$ and append it to the path.
 169 This produces a latent trajectory

$$\mathbf{Z}_0^\tau = [z_{i_1}^\tau, z_{i_2}^\tau, \dots, z_{i_H}^\tau] \in \mathbb{R}^{H \times d_z}, \quad \tilde{y}_{i_1}^\tau \leq \tilde{y}_{i_2}^\tau \leq \dots \leq \tilde{y}_{i_H}^\tau. \quad (7)$$

170 The important point is that every trajectory is constructed within one task, but after VAE encoding all
 171 such trajectories have the same shape. Thus, the diffusion model can be trained jointly across tasks
 172 without requiring a native-space diffusion process for each task.

173 **3.4 Metadata-conditioned latent diffusion**

174 Given the latent trajectories, DUO trains a conditional diffusion model to generate new improvement
 175 paths. For a clean latent trajectory \mathbf{Z}_0 , the forward diffusion process adds Gaussian noise:

$$\mathbf{Z}_t = \sqrt{\bar{\alpha}_t} \mathbf{Z}_0 + \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim \mathcal{N}(\mathbf{0}, \mathbf{I}). \quad (8)$$

176 A conditional 1D U-Net $\boldsymbol{\epsilon}_\theta$ predicts the noise from (\mathbf{Z}_t, t) and task-level condition \mathbf{c}_τ .

177 The condition \mathbf{c}_τ is built from natural-language metadata. For each task, we write a short text prompt
 178 containing fields such as Task, Type, Name, Description, and Objective. This prompt is encoded
 179 by the frozen `sentence-transformers/all-MiniLM-L6-v2` model, producing a 384-dimensional
 180 metadata vector:

$$\mathbf{e}_\tau^{\text{text}} = E_{\text{text}}(m_\tau). \quad (9)$$

181 We project this vector to the U-Net width and fuse it with the diffusion timestep embedding and
 182 optional return information:

$$\mathbf{c}_\tau = g_\eta(W_m \mathbf{e}_\tau^{\text{text}}, \mathbf{e}_t, \mathbf{e}_r). \quad (10)$$

183 Here \mathbf{e}_r is used only when return conditioning is enabled. The diffusion loss is

$$\mathcal{L}_{\text{diff}}^{(\tau)} = \mathbb{E}_{\mathbf{Z}_0^\tau, t, \boldsymbol{\epsilon}} \left[\|\boldsymbol{\epsilon} - \boldsymbol{\epsilon}_\theta(\mathbf{Z}_t^\tau, t, \mathbf{c}_\tau)\|_2^2 \right]. \quad (11)$$

184 To make metadata useful at sampling time, DUO uses classifier-free guidance on the text-metadata
 185 channel. During training, the metadata embedding is randomly dropped with probability p_{drop} .
 186 During sampling, the conditional and text-dropped predictions are combined as

$$\tilde{\boldsymbol{\epsilon}}_\theta = (1 + s)\boldsymbol{\epsilon}_\theta(\mathbf{Z}_t, t, \mathbf{c}_\tau) - s\boldsymbol{\epsilon}_\theta(\mathbf{Z}_t, t, \mathbf{c}_\emptyset), \quad s \geq 0. \quad (12)$$

187 In this default form, \mathbf{c}_\emptyset removes the text-metadata vector while keeping the diffusion timestep and
 188 other non-text conditioning signals unchanged. This makes the guidance weight s directly control
 189 how strongly the generated trajectory follows the task semantics.

190 **3.5 Discrete consistency**

191 Latent diffusion is trained with an MSE noise-prediction objective, but this objective alone is not
 192 always aligned with discrete evaluation. For example, in a DNA task, each sequence position must
 193 decode to exactly one nucleotide. Two relaxed vectors can be close under MSE while still producing
 194 different categorical decisions after argmax decoding. This mismatch is especially harmful when the
 195 benchmark oracle evaluates the final discrete sequence rather than the relaxed tensor.

Algorithm 1 DUO training and sampling

Require: Offline datasets $\{\mathcal{D}_\tau\}_{\tau \in \mathcal{T}_{\text{tr}}}$, metadata $\{m_\tau\}_{\tau \in \mathcal{T}_{\text{tr}}}$, trajectory length H , diffusion steps T , CE weight λ_{CE} , guidance scale s .

Ensure: Generated candidates for each target task.

- 1: Convert each native design \mathbf{x}_i^τ into a unified masked tensor $\mathbf{u}_i^\tau = \Pi_\tau(\mathbf{x}_i^\tau)$.
 - 2: Train the Transformer-VAE with \mathcal{L}_{VAE} and encode all offline designs into latent points \mathbf{z}_i^τ .
 - 3: For each task, construct synthetic improvement trajectories $\mathbf{Z}_0^\tau \in \mathbb{R}^{H \times d_z}$ by chaining local score-improving latent points.
 - 4: Encode metadata m_τ with the frozen text encoder and build conditions \mathbf{c}_τ .
 - 5: Train the conditional 1D U-Net denoiser with \mathcal{L}_{DUO} .
 - 6: **for** each target task τ **do**
 - 7: Sample latent trajectories from Gaussian noise using guided reverse diffusion.
 - 8: Decode generated latent points into native designs.
 - 9: Optionally rank the decoded candidates with an offline-trained surrogate.
 - 10: Return the top candidates.
 - 11: **end for**
-

196 DUO therefore adds a grouped cross-entropy consistency term for discrete tasks. From the predicted
197 noise, we recover the predicted clean latent trajectory

$$\hat{\mathbf{Z}}_0 = \frac{\mathbf{Z}_t - \sqrt{1 - \bar{\alpha}_t} \boldsymbol{\epsilon}_\theta(\mathbf{Z}_t, t, \mathbf{c}_\tau)}{\sqrt{\bar{\alpha}_t}}. \quad (13)$$

198 We then decode $\hat{\mathbf{Z}}_0$ into grouped logits $\hat{\mathbf{U}} = \text{Dec}_\psi(\hat{\mathbf{Z}}_0)$ and compare them with the original grouped
199 one-hot tensor \mathbf{U}_0 . For a discrete task with groups \mathcal{G}_τ , the auxiliary loss is

$$\mathcal{L}_{\text{CE}}^{(\tau)} = -\frac{1}{H|\mathcal{G}_\tau|} \sum_{h=1}^H \sum_{g \in \mathcal{G}_\tau} \sum_{c=1}^{C_g} U_{0,h,g,c} \log \frac{\exp(\hat{U}_{h,g,c})}{\sum_{c'=1}^{C_g} \exp(\hat{U}_{h,g,c'})}. \quad (14)$$

200 The full training objective is

$$\mathcal{L}_{\text{DUO}} = \mathbb{E}_{\tau \sim \mathcal{T}_{\text{tr}}} \left[\mathcal{L}_{\text{diff}}^{(\tau)} + \lambda_{\text{CE}} \mathbb{I}[\tau \in \mathcal{M}_{\text{disc}}] \mathcal{L}_{\text{CE}}^{(\tau)} \right], \quad (15)$$

201 where $\mathcal{M}_{\text{disc}}$ denotes the set of discrete tasks. The CE term is inactive for purely continuous tasks.

202 3.6 Sampling and candidate selection

203 At inference time, DUO starts from Gaussian noise $\mathbf{Z}_T \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ and runs the guided reverse
204 diffusion process using Eq. (12). This yields a generated latent trajectory

$$\hat{\mathbf{Z}}_0 = [\hat{\mathbf{z}}_1, \hat{\mathbf{z}}_2, \dots, \hat{\mathbf{z}}_H]. \quad (16)$$

205 Each latent point is decoded by the VAE decoder and mapped back to the native task space:

$$\hat{\mathbf{x}}_h = \Pi_\tau^{-1}(\text{Dec}_\psi(\hat{\mathbf{z}}_h)), \quad h = 1, \dots, H. \quad (17)$$

206 For continuous tasks, Π_τ^{-1} removes padding and reverses normalization. For discrete tasks, it applies
207 group-wise argmax or the benchmark-specified discretization rule.

208 Following the GTG evaluation pipeline, DUO may optionally train a lightweight task-specific
209 surrogate on the offline dataset and use it only as a final filter over generated candidates. This step
210 is not part of the core generative model: it is used to rank a finite candidate pool under matched
211 evaluation budgets. Without this filter, DUO simply returns decoded samples from the metadata-
212 conditioned latent diffusion model.

213 3.7 Algorithm summary

214 Algorithm 1 summarizes the full DUO pipeline. Line 1 standardizes heterogeneous raw designs into a
215 common masked tensor format. Line 2 learns the shared Transformer-VAE interface and stores each
216 offline design as a fixed-dimensional latent point. Line 3 turns these points into short score-improving

217 trajectories, which provide the training data for diffusion. Line 4 converts task metadata into condition
 218 vectors through the frozen text encoder and condition-fusion module. Line 5 trains the conditional
 219 1D U-Net with the DUO objective in Eq. (15). Lines 6–10 describe inference: for each target task,
 220 DUO samples latent trajectories by guided reverse diffusion, decodes the generated latents into native
 221 designs, optionally reranks them with an offline surrogate, and returns the top candidates.

222 4 Experiments

223 Our experiments investigate that whether one offline generator can be trained once on heterogeneous
 224 tasks and still produce high-scoring designs without any online feedback. We evaluate this question
 225 on Design-Bench [33] and SOO-Bench [25], following the universal multitask offline protocol [30]:
 226 nine shared training tasks, held-out tasks for transfer, and no oracle queries during learning.

227 4.1 Main results

228 Table 1 reports maximum rewards on the nine multitask training tasks (absolute simulator units; see
 229 caption). DUO achieves the best mean rank among UniSO, single-task GTG, and itself, with the
 230 clearest margins on continuous robotics and several GTOPIX missions; UniSO¹ remains competitive
 231 on some DNA and trajectory rows, so the comparison is best read task-by-task.

Table 1: Multitask training corpus at **raw** scale (mean±std): UniSO [30], single-task GTG [43], and our method—multitask DUO with full text-metadata embeddings. $\mathcal{D}(\text{best})$ is the best logged reward. Mean rank averages per-task ranks among UniSO, GTG, and DUO (lower is better). Ablations and single-task latent DUO are in Table 5.

Task	$\mathcal{D}(\text{best})$	UniSO	GTG	DUO
Ant	165.326	455.658 ± 39.188	452.330 ± 61.502	565.873 ± 13.651
D’Kitty	199.363	222.007 ± 33.677	264.857 ± 17.571	312.350 ± 6.353
Superconductor	74.000	82.642 ± 3.467	89.802 ± 6.763	112.880 ± 8.088
TF Bind 8	0.439	0.857 ± 0.069	0.934 ± 0.030	0.933 ± 0.055
TF Bind 10	0.005	0.944 ± 0.794	0.632 ± 0.106	0.955 ± 0.516
GTOPIX 2	-196.663	-90.106 ± 8.223	-72.163 ± 8.495	-44.477 ± 4.761
GTOPIX 3	-152.364	-44.427 ± 13.456	-53.866 ± 9.877	-29.310 ± 2.237
GTOPIX 4	-216.709	-74.779 ± 11.032	-72.749 ± 6.632	-51.019 ± 6.536
GTOPIX 6	-111.666	-48.244 ± 5.766	-55.834 ± 9.655	-33.449 ± 12.583
Mean rank	/	2.556 ± 0.497	2.333 ± 0.667	1.111 ± 0.314

232 Table 2 evaluates the same multitask DUO model (full text-metadata conditioning, our method) on
 233 the five Design-Bench tasks following [33, 30]; preprocessing is as in the table caption. We list list
 234 Design-Bench baselines together with offline BBO methods since 2024 onward and recompute **Avg.**
 235 **Rank** within this subset. The row $\mathcal{D}(\text{best})$ is the best objective observed in each task’s offline training
 236 log without a learned generative policy. Competing methods report mean±std when their releases
 237 include multiple runs.

238 Taken together, the two views support the same conclusion: DUO improves over the universal baseline
 239 UniSO overall (better mean rank on the training corpus and the best average rank in Table 2) and over
 240 representative single-task trajectory diffusion (GTG) and classic Design-Bench optimizers, while
 241 staying competitive with the strongest recent entries (e.g., LTR on DNA). Task-specific noise remains,
 242 especially on DNA binding, where Appendix H studies the cross-entropy auxiliary. Section 4.3
 243 and Appendix I add an illustrative reverse-diffusion diagnostic (Superconductor, seed 0) with full
 244 mean/top-8/max panels. The full Design-Bench leaderboard is Table 6.

245 4.2 Zero-shot and few-shot transfer on held-out tasks

246 We evaluate whether metadata and shared latents transfer to held-out control-style tasks: LunarLander,
 247 Rover, and RobotPush [4, 36, 24]. Table 3 reports maximum episode reward under zero-shot and

¹Unless stated otherwise, we write **UniSO** for the universal offline optimizer of Tan et al. [30]; Tables use the results of *Improved UniSO-T*, which is the best-performing model reported in that framework.

Table 2: Design-Bench (rescaled): Design-Bench baselines [33] plus learned optimizers from 2024 onward. Publication venues for all methods appear in Table 6. **Avg. Rank** reports the mean per-task rank divided by the number of competing methods in this table (16 learned optimizers; $\mathcal{D}(\text{best})$ excluded). Evaluation follows Appendix D.

Method	Ant	D’Kitty	Superconductor	TF-Bind-8	TF-Bind-10	Avg. Rank
$\mathcal{D}(\text{best})$	0.565	0.884	0.400	0.439	0.467	/
BO- q EI	0.812 \pm 0.000	0.896 \pm 0.000	0.382 \pm 0.013	0.802 \pm 0.081	0.628 \pm 0.036	12.0 / 16
CMA-ES	1.712 \pm 0.754	0.725 \pm 0.002	0.463 \pm 0.042	0.944 \pm 0.017	0.641 \pm 0.036	8.2 / 16
REINFORCE	0.248 \pm 0.039	0.541 \pm 0.196	0.478 \pm 0.017	0.935 \pm 0.049	0.673 \pm 0.074	10.0 / 16
Grad. Ascent	0.273 \pm 0.023	0.853 \pm 0.018	0.510 \pm 0.028	0.969 \pm 0.021	0.646 \pm 0.037	8.2 / 16
Grad. Ascent Mean	0.306 \pm 0.053	0.875 \pm 0.024	0.508 \pm 0.019	0.985 \pm 0.008	0.633 \pm 0.030	7.6 / 16
Grad. Ascent Min	0.282 \pm 0.033	0.884 \pm 0.018	0.514 \pm 0.020	0.979 \pm 0.014	0.632 \pm 0.027	7.8 / 16
PGS	0.715 \pm 0.046	0.954 \pm 0.022	0.444 \pm 0.020	0.889 \pm 0.061	0.634 \pm 0.040	9.6 / 16
FGM	0.923 \pm 0.023	0.944 \pm 0.014	0.481 \pm 0.024	0.811 \pm 0.079	0.611 \pm 0.008	9.8 / 16
MATCH-OPT	0.933 \pm 0.016	0.952 \pm 0.008	0.504 \pm 0.021	0.824 \pm 0.067	0.655 \pm 0.050	6.8 / 16
GTG	0.855 \pm 0.044	0.942 \pm 0.017	0.480 \pm 0.055	0.910 \pm 0.040	0.619 \pm 0.029	9.4 / 16
GABO	0.038 \pm 0.012	0.719 \pm 0.001	0.374 \pm 0.020	0.926 \pm 0.038	0.619 \pm 0.043	14.2 / 16
LTR	0.949 \pm 0.025	0.962 \pm 0.015	0.517 \pm 0.029	0.981 \pm 0.012	0.670 \pm 0.035	3.0 / 16
ROOT	0.958 \pm 0.012	0.971 \pm 0.005	0.451 \pm 0.032	0.977 \pm 0.015	0.653 \pm 0.030	5.4 / 16
DynAMO-Adam	0.113 \pm 0.085	0.789 \pm 0.059	0.413 \pm 0.106	0.719 \pm 0.142	0.556 \pm 0.090	14.8 / 16
UniSO	0.850 \pm 0.062	0.915 \pm 0.015	0.489 \pm 0.062	0.947 \pm 0.036	0.673 \pm 0.136	6.4 / 16
DUO	0.975 \pm 0.014	0.977 \pm 0.005	0.610 \pm 0.044	0.933 \pm 0.055	0.706 \pm 0.129	2.8 / 16

248 few-shot protocols (mean \pm std), alongside offline references $\mathcal{D}_{\text{all}}(\text{best})$ and $\mathcal{D}_{\text{fs}}(\text{best})$ and UniSO with
 249 metadata for direct comparison [30].

Table 3: Held-out control tasks: maximum episode reward (mean \pm std). UniSO uses the same metadata convention as [30]. $\mathcal{D}_{\text{all}}(\text{best})$ and $\mathcal{D}_{\text{fs}}(\text{best})$ are the best logged rewards on the merged offline dataset and on the few-shot evaluation pool, respectively.

Task	$\mathcal{D}_{\text{all}}(\text{best})$	$\mathcal{D}_{\text{fs}}(\text{best})$	UniSO zero-shot	UniSO few-shot	DUO zero-shot	DUO few-shot
RobotPush	9.764	-4.989	3.171 \pm 0.984	7.067 \pm 0.169	2.715 \pm 1.602	15.109 \pm 4.495
Rover	1.216	-32.127	-8.888 \pm 2.119	-8.239 \pm 1.270	33.489 \pm 0.016	39.600 \pm 3.183
LunarLander	280.180	-236.793	31.186 \pm 27.971	248.573 \pm 45.386	-216.627 \pm 4.602	81.193 \pm 44.657
Mean rank	/	/	3.333 \pm 0.577	2.000 \pm 1.000	3.333 \pm 1.155	1.333 \pm 0.577

250 Few-shot DUO leads RobotPush and Rover here, ahead of UniSO few-shot and the offline-pool
 251 references—consistent with trajectory diffusion exploiting a small extra log when metadata stay
 252 informative. LunarLander remains hard for all methods, suggesting a train–test dynamics gap beyond
 253 missing text. Zero-shot orderings vary by task, so we emphasize comparisons to offline pools and a
 254 matched universal baseline rather than a single headline.

255 4.3 Qualitative visualization along reverse diffusion

256 We visualize oracle **max** along reverse diffusion: at each sampling step we oracle-evaluate the current
 257 partial trajectory and take the maximum objective over trajectory positions (axes in Appendix I).
 258 Figure 2 overlays four regimes on one protocol (**Superconductor**, evaluation seed 0).

259 **ST** (single-task DUO without text-metadata embeddings) and **ST+text** (single-task with frozen
 260 text-metadata embeddings) curves remain close throughout, which matches what we expect under
 261 *single-task* training: without multitask data there is little cross-task manifold structure or transferable
 262 metadata knowledge to exploit, so adding diffusion conditioning on text alone barely moves the
 263 objective. Our full model **MT+text** (our default setting) clearly separates above **ST+text**, which sup-
 264 ports the claim that *multitask* training is doing real work. **MT+text** also tracks well above **MT+label**
 265 (multitask with task-ID conditioning only), indicating that rich textual metadata is substantially more
 266 useful than task-id conditioning alone. The oracle-**max** trace often peaks in the middle of reverse
 267 diffusion and weakens toward the final step because we deliberately maximize over *intermediate* chain
 268 states where stochasticity is higher—spuriously high spot maxima are easier—whereas the terminal
 269 decode is what we actually deploy; a mid-chain peak therefore does *not* imply that the intermediate

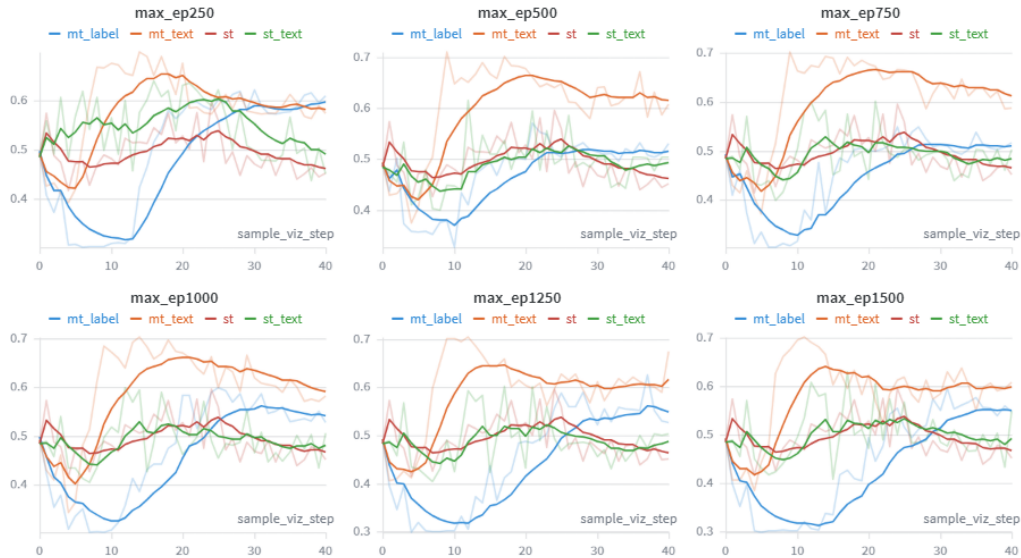


Figure 2: **Inference-time diagnostic.** Oracle **max** vs. reverse-diffusion step on **Superconductor** (evaluation seed 0). Six checkpoints in a 3×2 grid (early to late training, roughly 250–1500 epochs; learning rate 2×10^{-4}).

270 latent is the better model output. Vertically across checkpoints, the earliest row corresponds to ~ 250
 271 epochs and is still under-trained; rows in the ~ 500 –1000 epoch band look strongest, while ~ 1500
 272 epochs shows a mild decline consistent with overfitting to the finite offline log.

273 4.4 Additional experiments

274 We include three auxiliary analyses to clarify which parts of DUO matter beyond the headline reward
 275 tables. First, Appendix G sweeps the classifier-free guidance weight on text metadata; moderate
 276 guidance generally balances continuous robotics and DNA objectives, while overly large weights
 277 can hurt some GTOPIX tasks, consistent with guidance pushing samples away from the offline
 278 support. Second, Appendix H shows that the cross-entropy auxiliary mainly improves the discrete
 279 TF Bind tasks, with smaller changes on continuous and GTOPIX rows, supporting its role as a
 280 targeted correction for discrete decoding rather than a global regularizer. Finally, Appendix I expands
 281 Figure 2 with mean and top-8 oracle traces, which helps distinguish broad trajectory improvement
 282 from isolated high-scoring spikes; Appendix E lists the trajectory-construction hyperparameters used
 283 throughout.

284 5 Conclusion

285 In this paper, we present DUO, coupling metadata-aware multitask conditioning with latent trajectory
 286 diffusion behind a shared Transformer–VAE bottleneck. Empirically it compares favorably to UniSO
 287 and strong single-task optimizers on the benchmarks we report, suggesting that universal offline
 288 optimization benefits from combining semantic task conditioning with latent trajectory diffusion
 289 rather than relying only on forward reward prediction followed by search. One limitation is that our
 290 current evaluation still focuses on a moderate collection of Design-Bench and SOO-Bench tasks rather
 291 than Vizier-scale suites [6], and the approach depends on informative task metadata and diffusion
 292 sampling. Natural next steps include broader prompt studies, comparisons with diffusion–LLM
 293 hybrids [42], stronger discrete encoders, multi-objective extensions [38], improved samplers [44],
 294 and distillation for deployment.

295 References

296 [1] Michael Ahn, Henry Zhu, Kristian Hartikainen, Hugo Ponte, Abhishek Gupta, Sergey Levine,
 297 and Vikash Kumar. ROBEL: Robotics Benchmarks for Learning with Low-Cost Robots, 2019.

- 298 [2] Anurag Ajay, Yilun Du, Abhi Gupta, Joshua Tenenbaum, Tommi Jaakkola, and Pulkit Agrawal.
299 Is Conditional Generative Modeling all you need for Decision-Making?, 2023.
- 300 [3] Luis A. Barrera, Anastasia Vedenko, Jesse V. Kurland, Julia M. Rogers, Stephen S. Gisselbrecht,
301 Elizabeth J. Rossin, Jaie Woodard, Luca Mariani, Kian Hong Kock, Sachi Inukai, Trevor Siggers,
302 Leila Shokri, Raluca Gordân, Nidhi Sahni, Chris Cotsapas, Tong Hao, Song Yi, Manolis Kellis,
303 Mark J. Daly, Marc Vidal, David E. Hill, and Martha L. Bulyk. Survey of variation in human
304 transcription factors reveals prevalent DNA binding changes. *Science*, 351(6280):1450–1454,
305 2016. doi: 10.1126/science.aad2257.
- 306 [4] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang,
307 and Wojciech Zaremba. OpenAI Gym. <https://arxiv.org/abs/1606.01540v1>, 2016.
- 308 [5] Can Sam Chen, Christopher Beckham, Zixuan Liu, Xue Liu, and Christopher Pal. Robust
309 Guided Diffusion for Offline Black-Box Optimization, 2024.
- 310 [6] Yutian Chen, Xingyou Song, Chansoo Lee, Zi Wang, Qiuyi Zhang, David Dohan, Kazuya
311 Kawakami, Greg Kochanski, Arnaud Doucet, Marc’ aurelio Ranzato, Sagi Perel, and Nando
312 de Freitas. Towards Learning Universal Hyperparameter Optimizers with Transformers, 2022.
- 313 [7] Zhou Fan, Xinran Han, and Zi Wang. Transfer Learning for Bayesian Optimization on Hetero-
314 geneous Search Spaces, 2023.
- 315 [8] Roman Garnett. *Bayesian Optimization*. Cambridge University Press, 2023. ISBN 978-1-108-
316 62355-1.
- 317 [9] Kam Hamidieh. A Data-Driven Statistical Model for Predicting the Critical Temperature of a
318 Superconductor, 2018.
- 319 [10] Jonathan Ho and Tim Salimans. Classifier-Free Diffusion Guidance, 2022.
- 320 [11] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising Diffusion Probabilistic Models, 2020.
- 321 [12] Michael Janner, Yilun Du, Joshua B. Tenenbaum, and Sergey Levine. Planning with Diffusion
322 for Flexible Behavior Synthesis, 2022.
- 323 [13] Minsu Kim, Jiayao Gu, Ye Yuan, Taeyoung Yun, Zixuan Liu, Yoshua Bengio, and Can Chen.
324 Offline Model-Based Optimization: Comprehensive Review, 2026.
- 325 [14] Diederik P. Kingma and Max Welling. Auto-Encoding Variational Bayes, 2022.
- 326 [15] Siddarth Krishnamoorthy, Satvik Mehul Mashkaria, and Aditya Grover. Diffusion Models for
327 Black-Box Optimization, 2023.
- 328 [16] Siddarth Krishnamoorthy, Satvik Mehul Mashkaria, and Aditya Grover. Generative Pretraining
329 for Black-Box Optimization, 2023.
- 330 [17] Aviral Kumar and Sergey Levine. Model Inversion Networks for Model-Based Optimization,
331 2019.
- 332 [18] Xiaobin Li, Kai Wu, Yujian Betterest Li, Xiaoyu Zhang, Handing Wang, and Jing Liu. Pretrained
333 Optimization Model for Zero-Shot Black Box Optimization. <https://arxiv.org/abs/2405.03728v2>,
334 2024.
- 335 [19] Zhixuan Liang, Yao Mu, Mingyu Ding, Fei Ni, Masayoshi Tomizuka, and Ping Luo. AdaptDif-
336 fuser: Diffusion Models as Adaptive Self-evolving Planners, 2023.
- 337 [20] Natalie Maus, Yimeng Zeng, Haydn Thomas Jones, Yining Huang, Gaurav Ng Goel, Alden
338 Rose, Kyurae Kim, Hyun-Su Lee, Marcelo Der Torossian Torres, Fangping Wan, Cesar de la
339 Fuente-Nunez, Mark Yatskar, Osbert Bastani, and Jacob R. Gardner. Purely Agentic Black-Box
340 Optimizational for Biological Design, 2026.
- 341 [21] Tung Nguyen, Sudhanshu Agrawal, and Aditya Grover. ExPT: Synthetic Pretraining for
342 Few-Shot Experimental Design, 2023.

- 343 [22] Tung Nguyen, Qiuyi Zhang, Bangding Yang, Chansoo Lee, Jorg Bornschein, Yingjie Miao, Sagi
344 Perel, Yutian Chen, and Xingyou Song. Predicting from Strings: Language Model Embeddings
345 for Bayesian Optimization, 2024.
- 346 [23] Tung Nguyen, Qiuyi Zhang, Bangding Yang, Chansoo Lee, Jorg Bornschein, Yingjie Miao,
347 Sagi Perel, Yutian Chen, and Xingyou Song. Language Model Embeddings Can Be Sufficient
348 for Bayesian Optimization, 2025.
- 349 [24] Ian Parberry. *Introduction to Game Physics with Box2D*. CRC Press, 2013. ISBN 978-1-4665-
350 6577-7.
- 351 [25] Hong Qian, Yiyi Zhu, Xiang Shu, Shuo Liu, Yaolin Wen, Xin An, Huakang Lu, Aimin Zhou,
352 Ke Tang, and Yang Yu. SOO-Bench: Benchmarks for Evaluating the Stability of Offline Black-
353 Box Optimization. In *The Thirteenth International Conference on Learning Representations*,
354 2024.
- 355 [26] Nils Reimers and Iryna Gurevych. Sentence-BERT: Sentence Embeddings using Siamese BERT-
356 Networks. In Kentaro Inui, Jing Jiang, Vincent Ng, and Xiaojun Wan, editors, *Proceedings*
357 *of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th*
358 *International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages
359 3982–3992, Hong Kong, China, 2019. Association for Computational Linguistics. doi: 10.
360 18653/v1/D19-1410.
- 361 [27] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer.
362 High-Resolution Image Synthesis with Latent Diffusion Models, 2022.
- 363 [28] Martin Schlueter, Mehdi Neshat, Mohamed Wahib, Masaharu Munetomo, and Markus Wagner.
364 GTOPIX space mission benchmarks. *SoftwareX*, 14:100666, 2021. doi: 10.1016/j.softx.2021.
365 100666.
- 366 [29] Xingyou Song, Oscar Li, Chansoo Lee, Bangding Yang, Daiyi Peng, Sagi Perel, and Yutian
367 Chen. OmniPred: Language Models as Universal Regressors, 2025.
- 368 [30] Rong-Xi Tan, Ming Chen, Ke Xue, Yao Wang, Yaoyuan Wang, Sheng Fu, and Chao Qian.
369 Towards Universal Offline Black-Box Optimization via Learning Language Model Embeddings,
370 2025.
- 371 [31] Rong-Xi Tan, Ke Xue, Shen-Huan Lyu, Haopu Shang, Yao Wang, Yaoyuan Wang, Sheng Fu,
372 and Chao Qian. Offline Model-Based Optimization by Learning to Rank, 2025.
- 373 [32] Brandon Trabucco, Aviral Kumar, Xinyang Geng, and Sergey Levine. Conservative Objective
374 Models for Effective Offline Model-Based Optimization, 2021.
- 375 [33] Brandon Trabucco, Xinyang Geng, Aviral Kumar, and Sergey Levine. Design-Bench: Bench-
376 marks for Data-Driven Offline Model-Based Optimization, 2022.
- 377 [34] Austin Tripp, Erik Daxberger, and José Miguel Hernández-Lobato. Sample-Efficient Optimiza-
378 tion in the Latent Space of Deep Generative Models via Weighted Retraining, 2020.
- 379 [35] Masatoshi Uehara, Yulai Zhao, Ehsan Hajiramezani, Gabriele Scalia, Gökçen Eraslan, Avan-
380 tika Lal, Sergey Levine, and Tommaso Biancalani. Bridging Model-Based Optimization and
381 Generative Modeling via Conservative Fine-Tuning of Diffusion Models, 2024.
- 382 [36] Shukuan Wang, Ke Xue, Lei Song, Xiaobin Huang, and Chao Qian. Monte Carlo Tree Search
383 based Space Transfer for Black-box Optimization, 2024.
- 384 [37] Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. MiniLM:
385 Deep Self-Attention Distillation for Task-Agnostic Compression of Pre-Trained Transformers.
386 <https://arxiv.org/abs/2002.10957v2>, 2020.
- 387 [38] Ke Xue, Rong-Xi Tan, Xiaobin Huang, and Chao Qian. Offline Multi-Objective Optimization,
388 2024.

- 389 [39] Peiyu Yu, Dinghuai Zhang, Hengzhi He, Xiaojian Ma, Ruiyao Miao, Yifan Lu, Yasi Zhang,
390 Deqian Kong, Ruiqi Gao, Jianwen Xie, Guang Cheng, and Ying Nian Wu. Latent Energy-Based
391 Odyssey: Black-Box Optimization via Expanded Exploration in the Energy-Based Latent Space,
392 2024.
- 393 [40] Sihyun Yu, Sungsoo Ahn, Le Song, and Jinwoo Shin. RoMA: Robust Model Adaptation for
394 Offline Model-based Optimization, 2021.
- 395 [41] Ye Yuan, Can Chen, Zixuan Liu, Willie Neiswanger, and Xue Liu. Importance-aware Co-
396 teaching for Offline Model-based Optimization, 2023.
- 397 [42] Ye Yuan, Can, Chen, Zipeng Sun, Dinghuai Zhang, Christopher Pal, and Xue Liu. Diffusion
398 Large Language Models for Black-Box Optimization, 2026.
- 399 [43] Taeyoung Yun, Sujin Yun, Jaewoo Lee, and Jinkyoo Park. Guided Trajectory Generation with
400 Diffusion Models for Offline Model-based Optimization, 2024.
- 401 [44] Taeyoung Yun, Kiyoungh Om, Jaewoo Lee, Sujin Yun, and Jinkyoo Park. Posterior Inference
402 with Diffusion Models for High-dimensional Black-box Optimization, 2025.

403 A Extended preliminaries

404 This appendix collects formal definitions and standard objectives that parallel the treatment in our
405 thesis-style exposition, so that the main paper can stay focused on modeling choices. Readers
406 comfortable with latent diffusion and offline RL-style notation may skim and proceed to Appendix B.

407 A.1 Single-task offline black-box optimization

408 Let \mathcal{X} denote a search space and $f : \mathcal{X} \rightarrow \mathbb{R}$ an unknown objective accessible only through past
409 evaluations. Given a static dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$ with $y_i = f(x_i)$, offline black-box optimization
410 seeks

$$x^* \in \arg \max_{x \in \mathcal{X}} f(x) \quad (18)$$

411 using only information contained in \mathcal{D} and inductive biases encoded in the algorithm [33]. Crucially,
412 no additional queries $f(x)$ are permitted during learning: this models expensive simulators, wet-lab
413 measurements, or any setting where online interaction is prohibitively costly.

414 A.2 Universal offline black-box optimization

415 Universal offline BBO [30] considers a finite suite of tasks \mathcal{T} . Each task $\tau \in \mathcal{T}$ is equipped with
416 its own space \mathcal{X}_τ , dimension d_τ , objective f_τ , dataset \mathcal{D}_τ , and metadata m_τ (names, descriptions,
417 goals, variable types, ...). The goal is to learn a *single* shared model that, when conditioned on
418 m_τ , proposes high-scoring designs for every training task and—when metadata are available at
419 deployment—generalizes to held-out tasks in zero- or few-shot regimes.

420 A.3 Variational autoencoders (VAE)

421 A VAE [14] consists of an encoder and decoder, inducing an approximate posterior $q_\phi(z|x)$ and
422 reconstructing from latent z . A common training objective combines reconstruction and a KL penalty
423 to a prior $p(z)$:

$$\mathcal{L}_{\text{VAE}} = \mathbb{E}_{z \sim q_\phi(z|x)} \left[\|x - \text{Dec}(z)\|_2^2 \right] + D_{\text{KL}}(q_\phi(z|x) \parallel p(z)). \quad (19)$$

424 In DUO, the VAE maps heterogeneous native tensors into a fixed-format latent tensor so that a single
425 diffusion backbone can be shared across tasks.

426 A.4 Denoising diffusion (DDPM)

427 DDPMs [11] define a forward noising process that gradually corrupts data x_0 into Gaussian noise
428 over T steps. A common Gaussian parameterization is

$$q(x_t | x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta_t} x_{t-1}, \beta_t I), \quad (20)$$

429 with variance schedule $\{\beta_t\}_{t=1}^T$. Let $\bar{\alpha}_t = \prod_{i=1}^t (1 - \beta_i)$; then one may sample

$$x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, \quad \epsilon \sim \mathcal{N}(0, I), \quad (21)$$

430 which is the same reparameterization used in the latent trajectory model of Section 3. The reverse
431 process is learned by a network ϵ_θ that predicts noise from (x_t, t) (optionally conditioned on task
432 signals). A widely used simplified noise-prediction objective is

$$\mathcal{L}_{\text{diff}} = \mathbb{E}_{t, x_0, \epsilon} \left[\|\epsilon - \epsilon_\theta(x_t, t)\|_2^2 \right], \quad (22)$$

433 where x_t is sampled from the forward schedule given x_0 .

434 A.5 Classifier-free guidance (CFG)

435 CFG [10] trains a score model with randomly dropped conditions and at sampling time mixes
436 conditional and unconditional predictions:

$$\hat{c} = (1 + s) \epsilon_\theta(x_t, t, c) - s \epsilon_\theta(x_t, t, \emptyset), \quad s \geq 0. \quad (23)$$

437 In DUO’s default setup, the guided component c in Eq. (23) should be read as the *textual metadata*
 438 *embedding*; training randomly drops this vector while keeping timestep, task index, and optional
 439 return features in place, and sampling uses a scalar weight s to interpolate between conditional and
 440 text-dropped predictions. The public codebase can also factor guidance across additional condition
 441 splits for controlled comparisons (Appendix G), but we do not treat that engineering variant as a core
 442 contribution.

443 A.6 Discrete cross-entropy grouping

444 For a length- L sequence with alphabet size C , let $\mathbf{u}_i \in \mathbb{R}^C$ be logits at position i and \mathbf{y}_i the one-hot
 445 label. The auxiliary loss is

$$\mathcal{L}_{\text{CE}} = -\frac{1}{L} \sum_{i=1}^L \sum_{c=1}^C y_{i,c} \log \frac{\exp(u_{i,c})}{\sum_{c'} \exp(u_{i,c'})}. \quad (24)$$

446 A.7 Trajectory-based generative offline optimization

447 **Synthetic improvement trajectories.** Instead of modeling only the marginal distribution of high-
 448 scoring designs, trajectory-based methods construct ordered sequences $(x^{(0)}, \dots, x^{(H)})$ that mimic
 449 iterative improvement inside the offline log [16, 43]. The generative model is trained to reproduce
 450 the *dynamics* of moving from lower- to higher-reward regions, which often yields better optimization
 451 inductive bias than static density modeling when data are limited.

452 **BONET and GTG.** BONET [16] sorts offline points by score and samples trajectories toward high
 453 values, using an autoregressive Transformer to model the sequence. GTG [43] refines this idea by
 454 chaining neighbors in design space so that successive states resemble local optimization moves, and
 455 replaces autoregressive decoding with diffusion to mitigate compounding errors. Both demonstrate
 456 that “learning how improvement unfolds” is a strong recipe for offline BBO.

457 **Limitations for universal settings.** Existing trajectory pipelines are typically bound to a single
 458 \mathcal{X}_τ and do not natively share parameters across heterogeneous spaces, nor do they consume free-
 459 form metadata beyond coarse task IDs. DUO builds on the GTG-style trajectory viewpoint but
 460 lifts trajectories into a shared latent tensor space with Transformer-VAE encoding and metadata-
 461 conditioned diffusion, directly targeting the universal formulation in Appendix A.2.

462 B DUO training and sampling pseudocode

463 Algorithm 2 mirrors the implementation order in our training scripts: optional GTG-style [43]
 464 surrogates and the VAE are fit before trajectory construction, diffusion training consumes the latent
 465 trajectories together with metadata embeddings, and the sampling loop may close with the same
 466 surrogate ranking step as in the public GTG pipeline when we seek matched evaluation budgets.
 467 Pseudocode omits engineering details (mixed-precision flags, EMA weights, gradient clipping) that
 468 do not change the scientific narrative but are documented alongside the released code.

469 C Tasks, datasets, and metadata

470 **Benchmarks and splits.** We follow the unconstrained task selection commonly used in universal
 471 offline BBO work [30]: **Design-Bench** [33] provides Ant Morphology, D’Kitty Morphology, Super-
 472 conductor, TF Bind 8, and TF Bind 10; **SOO-Bench** [25] provides GTOPIX instances 2, 3, 4, and
 473 6. These **nine** tasks form the default multitask *training* corpus for DUO in the main experiments
 474 (Section ??). Together they cover continuous robotics morphologies [4, 1], a surrogate materials
 475 objective [9], discrete DNA binding design [3], and continuous trajectory planning [28].

476 **Offline logs.** Every task ships with a static design–score dataset; we never query the true simulator
 477 or oracle during DUO training, matching the offline BBO definition in Design-Bench. When a
 478 benchmark exposes extremely large logs, we follow the same fixed-size subsampling practice as in
 479 the community implementations we compare against (details mirror our thesis experiments).

Algorithm 2 DUO: diffusion models for universal offline black-box optimization (outline)

Require: Task datasets $\{\mathcal{D}_k\}$, metadata $\{m_k\}$, diffusion steps T , trajectory length H , weight λ_{CE} , text embedder E

Ensure: Per-task candidates $\{x_k^*\}$

- 1: Train task surrogates $f_{\text{surr},k}$
 - 2: Fit VAE (Enc, Dec)
 - 3: Encode points $\mathbf{z} \leftarrow \text{Enc}(\mathbf{x})$; build synthetic trajectories \mathbf{Z}
 - 4: Embed metadata $\mathbf{c}_k \leftarrow E(m_k)$; train conditional U-Net ϵ_θ with loss (15)
 - 5: **for** each task k **do**
 - 6: Sample latent trajectories with classifier-free guidance on text metadata to obtain $\{\hat{\mathbf{z}}_k\}$
 - 7: Rank with $f_{\text{surr},k}$, decode top- N : $x_k^* \leftarrow \text{Dec}(\hat{\mathbf{z}}_k)$
 - 8: **end for**
 - 9: **return** $\{x_k^*\}$
-

480 **Task metadata for conditioning.** As in UniSO [30], we concatenate short textual fields—typically
481 task name, variable type, natural-language description, and stated optimization objective—into one
482 metadata string per task, encode it with the frozen `sentence-transformers/all-MiniLM-L6-v2`
483 encoder (MiniLM [37]; Sentence-BERT-style training [26]; 384-dimensional outputs), and supply
484 the projected embedding to DUO’s diffusion backbone (together with task-index embeddings where
485 used). Field-level templates and longer Chinese explanations appear in the companion thesis.

486 **Held-out transfer tasks.** For zero- and few-shot transfer we additionally evaluate LunarLander,
487 Rover, and RobotPush [4, 36, 24], using the same offline logs and metadata conventions as in [30].

488 D Experimental setup details

489 This section expands the condensed setup in Section ???. Every experiment respects the offline
490 protocol: the true objective is never queried during training; all learning uses static logs \mathcal{D}_τ . When
491 we mirror the GTG [43] decoding stack, lightweight surrogates may rank candidates before oracle
492 calls, matching that codebase rather than defining a new protocol element.

493 **Design-Bench evaluation.** Following the experimental setting of Tan et al. [30] on Design-
494 Bench [33], for each task we independently sample 128 candidate designs from the learned policy,
495 evaluate each with the oracle, and report the maximum objective. Mean and standard deviation
496 aggregate evaluation seeds as in Section ???. **Avg. Rank** is computed from per-task ranks among the
497 listed competing methods, normalized as in each table caption. Numbers taken from prior work retain
498 the statistics reported in the original publications.

499 **Random seeds.** All DUO runs that we train and evaluate ourselves use **eight** random seeds
500 $\{0, 1, \dots, 7\}$; reported means and standard deviations pool these runs. Tables that include external
501 baselines may quote $\text{mean} \pm \text{std}$ only when the original release provides multiple seeds, but our own
502 rows always follow the eight-seed protocol above.

503 **Multitask training corpus.** Unless we explicitly ablate single-task variants, DUO is trained *once*
504 on the **nine** tasks listed in Appendix C, using textual metadata and task-index conditioning. Any table
505 with fewer than nine task columns still uses this same checkpoint and simply omits tasks that fall
506 outside that evaluation protocol (e.g., standard Design-Bench scaling reports five columns). Held-out
507 control tasks in Section 4.2 are *not* part of this nine-task fit.

508 **Text embeddings.** Runs that condition on textual task metadata use the frozen Hugging Face
509 checkpoint `sentence-transformers/all-MiniLM-L6-v2` (a 384-dimensional MiniLM sentence
510 encoder [37] in the Sentence-BERT lineage [26]); see Section 3 for fusion into diffusion.

511 **E Trajectory construction hyperparameters**

512 Synthetic trajectories are controlled by four quantities: the number of sampled paths per refresh n_{traj} ,
 513 the neighborhood size k used when chaining latent states, the edge-quality threshold ε , and the fixed
 514 trajectory length $H=64$. Here k limits how many nearest neighbors are considered when forming a
 515 synthetic edge, and ε drops edges whose score improvement is too small so that trajectories do not
 516 wander through flat or misleading regions.

517 **Design-Bench tasks.** For Ant, D’Kitty, Superconductor, TF Bind 8, and TF Bind 10 we **reuse**
 518 **the trajectory hyperparameter recommendations from GTG** [43] (i.e., the same k and ε choices
 519 advocated there for single-task trajectory diffusion), so that our multitask DUO remains aligned with
 520 the best-studied single-task trajectory baseline on these domains.

521 **SOO-Bench (GTOPX) tasks.** For GTOPX 2/3/4/6 we do not rely on an external recipe; instead
 522 we selected (k, ε) via **lightweight single-task screening** over a small candidate grid on each mission,
 523 picked the setting with the best validation proxy on that task alone, and then froze those values for all
 524 multitask training and reporting (matching the procedure documented in our thesis).

525 Table 4 lists the resulting $(n_{\text{traj}}, k, \varepsilon)$ for every training task.

Table 4: Per-task trajectory construction hyperparameters (n_{traj} , neighborhood size k , threshold ε); latent trajectory length $H=64$. Design-Bench rows follow GTG-recommended (k, ε) ; GTOPX rows use values from single-task screening (Section text).

Task	n_{traj}	k	ε
Ant	1000	20	0.05
D’Kitty	1000	20	0.01
Superconductor	1000	20	0.05
TF Bind 8	1000	50	0.05
TF Bind 10	1000	50	0.05
GTOPX 2	1000	20	0.05
GTOPX 3	1000	10	0.01
GTOPX 4	1000	20	0.05
GTOPX 6	1000	20	0.01

526 **F Training-corpus comparison and ablations**

527 Section 4.1 reports aggregate multitask performance together with the Design-Bench view; this
 528 appendix gives the full training-corpus grid used for that comparison. Table 5 places the universal
 529 regressor **UniSO** and single-task trajectory diffusion **GTG** next to four **DUO** training regimes:
 530 **ST** and **ST+text** (single-task, without vs. with frozen text-metadata embeddings), and **MT+label**
 531 vs. **MT+text** (multitask with task-ID conditioning only vs. multitask with full text metadata—our
 532 default). The mean-rank row summarizes cross-task ordering among these six learned columns
 533 together with the two multitask baselines (see table caption).

534 The **ST** column diffuses entirely in the shared VAE latent while GTG diffuses in native coordinates;
 535 scores are often comparable and trade wins row-wise, which supports latent trajectory diffusion as
 536 a workable interface rather than a bookkeeping trick before multitask sharing. Under single-task
 537 training, **ST+text** moves only modestly relative to **ST**, consistent with Section 4.3: without multitask
 538 context, metadata has little structure to exploit and text-only conditioning rarely reshapes the objective.
 539 By contrast, **MT+text** improves markedly over **ST+text** on average and attains the best mean rank
 540 among the four DUO columns, matching the diagnostic where **MT+text** separates upward while
 541 **ST** and **ST+text** remain entangled—evidence that shared training, not text alone, drives the gain.
 542 **MT+text** also dominates **MT+label**, paralleling Section 4.3: task IDs are a weak substitute for full
 543 textual descriptions when steering a universal denoiser. Finally, relative to Section 4.1, **UniSO** and
 544 **GTG** remain strong anchors: **UniSO** stays competitive on several DNA and trajectory rows, whereas
 545 **MT+text** tends to show the clearest margins on continuous robotics and selected GTOPX tasks, so
 546 the grid is best read task-by-task rather than as a single global ordering [30].

Table 5: Multitask training corpus at **raw** simulator scale (mean \pm std). Baselines: UniSO [30] and single-task GTG [43]. The four DUO columns are latent trajectory diffusion ablations: (i) single-task without metadata embeddings; (ii) single-task with frozen text-metadata embeddings; (iii) multitask with task-ID conditioning only; (iv) our method (multitask with full text-metadata embeddings). $\mathcal{D}(\text{best})$ is the best logged reward. Mean rank averages per-task ranks across the six learned columns (lower is better). A compact three-way comparison (UniSO, GTG, our method) is in Table 1.

Task	$\mathcal{D}(\text{best})$	UniSO	GTG	DUO (single-task)	DUO (single-task + metadata emb.)	DUO (multitask + task label)	DUO (multitask + metadata emb.)
Ant	165.326	455.658 \pm 39.188	452.330 \pm 61.502	474.207 \pm 11.743	490.691 \pm 17.299	565.699 \pm 8.628	565.873 \pm 13.651
D’Kitty	199.363	222.007 \pm 33.677	264.857 \pm 17.571	257.494 \pm 17.376	252.830 \pm 10.977	300.777 \pm 13.149	312.350 \pm 6.353
Superconductor	74.000	82.642 \pm 3.467	89.802 \pm 6.763	78.774 \pm 4.163	76.990 \pm 3.889	100.741 \pm 12.287	112.880 \pm 8.088
TF Bind 8	0.439	0.857 \pm 0.069	0.934 \pm 0.030	0.953 \pm 0.035	0.952 \pm 0.038	0.724 \pm 0.073	0.933 \pm 0.055
TF Bind 10	0.005	0.944 \pm 0.794	0.632 \pm 0.106	0.641 \pm 0.144	0.787 \pm 0.076	0.632 \pm 0.171	0.955 \pm 0.516
GTOPX 2	-196.663	-90.106 \pm 8.223	-72.163 \pm 8.495	-66.664 \pm 15.951	-75.231 \pm 20.343	-110.869 \pm 13.664	-44.477 \pm 4.761
GTOPX 3	-152.364	-44.427 \pm 13.456	-53.866 \pm 9.877	-47.522 \pm 6.089	-51.900 \pm 13.122	-51.175 \pm 14.311	-29.310 \pm 2.237
GTOPX 4	-216.709	-74.779 \pm 11.032	-72.749 \pm 6.632	-57.550 \pm 12.458	-67.212 \pm 11.006	-63.167 \pm 13.669	-51.019 \pm 6.536
GTOPX 6	-111.666	-48.244 \pm 5.766	-55.834 \pm 9.655	-59.971 \pm 9.638	-58.092 \pm 2.617	-61.664 \pm 10.439	-33.449 \pm 12.583
Mean rank	/	4.111 \pm 1.691	4.222 \pm 1.481	3.333 \pm 1.414	4.000 \pm 1.225	4.000 \pm 1.803	1.333 \pm 1.000

547 **Design-Bench leaderboard (complete table).** Table 6 lists all competing methods with **Source**
548 venues; Section 4.1 uses a shorter method list in Table 2.

Table 6: **Full** Design-Bench leaderboard under standard per-task rescaling (mean \pm std when reported). **Source** lists the publication venue (or Design-Bench baselines [33]). **Avg. Rank:** mean per-task rank divided by the number of competing methods (26). Evaluation: Appendix D. Short leaderboard: Table 2.

Method	Source	Ant	D’Kitty	Superconductor	TF-Bind-8	TF-Bind-10	Avg. Rank
$\mathcal{D}(\text{best})$	/	0.565	0.884	0.400	0.439	0.467	/
BO- q EI	Design-Bench baselines [33]	0.812 \pm 0.000	0.896 \pm 0.000	0.382 \pm 0.013	0.802 \pm 0.081	0.628 \pm 0.036	20.4 / 26
CMA-ES		1.712 \pm 0.754	0.725 \pm 0.002	0.463 \pm 0.042	0.944 \pm 0.017	0.641 \pm 0.036	12.8 / 26
REINFORCE		0.248 \pm 0.039	0.541 \pm 0.196	0.478 \pm 0.017	0.935 \pm 0.049	0.673 \pm 0.074	15.7 / 26
Grad. Ascent		0.273 \pm 0.023	0.853 \pm 0.018	0.510 \pm 0.028	0.969 \pm 0.021	0.646 \pm 0.037	12.8 / 26
Grad. Ascent Mean		0.306 \pm 0.053	0.875 \pm 0.024	0.508 \pm 0.019	0.985 \pm 0.008	0.633 \pm 0.030	12.4 / 26
Grad. Ascent Min		0.282 \pm 0.033	0.884 \pm 0.018	0.514 \pm 0.020	0.979 \pm 0.014	0.632 \pm 0.027	12.6 / 26
CbAS	ICML’19	0.846 \pm 0.032	0.896 \pm 0.009	0.421 \pm 0.049	0.921 \pm 0.046	0.630 \pm 0.039	17.7 / 26
MINs	ICML’19	0.906 \pm 0.024	0.939 \pm 0.007	0.464 \pm 0.023	0.910 \pm 0.051	0.633 \pm 0.034	14.6 / 26
DDOM	ICML’23	0.908 \pm 0.024	0.930 \pm 0.005	0.452 \pm 0.028	0.913 \pm 0.047	0.616 \pm 0.018	16.4 / 26
BONET	ICML’23	0.921 \pm 0.031	0.949 \pm 0.016	0.390 \pm 0.022	0.798 \pm 0.123	0.575 \pm 0.039	17.3 / 26
GTG	NeurIPS’24	0.855 \pm 0.044	0.942 \pm 0.017	0.480 \pm 0.055	0.910 \pm 0.040	0.619 \pm 0.029	15.6 / 26
COMs	ICML’21	0.916 \pm 0.026	0.949 \pm 0.016	0.460 \pm 0.040	0.953 \pm 0.038	0.644 \pm 0.052	10.5 / 26
RoMA	ICML’21	0.430 \pm 0.048	0.767 \pm 0.031	0.494 \pm 0.025	0.665 \pm 0.000	0.553 \pm 0.000	20.9 / 26
IOM	NeurIPS’22	0.889 \pm 0.034	0.928 \pm 0.008	0.491 \pm 0.034	0.925 \pm 0.054	0.628 \pm 0.036	14.5 / 26
BDI	NeurIPS’22	0.963 \pm 0.000	0.941 \pm 0.000	0.508 \pm 0.013	0.973 \pm 0.000	0.658 \pm 0.000	6.5 / 26
ICT	NeurIPS’23	0.915 \pm 0.024	0.947 \pm 0.009	0.494 \pm 0.026	0.897 \pm 0.050	0.659 \pm 0.024	10.6 / 26
Tri-Mentoring	NeurIPS’23	0.891 \pm 0.011	0.947 \pm 0.005	0.503 \pm 0.013	0.956 \pm 0.000	0.662 \pm 0.012	8.3 / 26
PGS	AAAI’24	0.715 \pm 0.046	0.954 \pm 0.022	0.444 \pm 0.020	0.889 \pm 0.061	0.634 \pm 0.040	15.4 / 26
FGM	AISTATS’24	0.923 \pm 0.023	0.944 \pm 0.014	0.481 \pm 0.024	0.811 \pm 0.079	0.611 \pm 0.008	15.0 / 26
MATCH-OPT	ICML’24	0.933 \pm 0.016	0.952 \pm 0.008	0.504 \pm 0.021	0.824 \pm 0.067	0.655 \pm 0.050	9.4 / 26
GABO	NeurIPS’24 Oral	0.038 \pm 0.012	0.719 \pm 0.001	0.374 \pm 0.020	0.926 \pm 0.038	0.619 \pm 0.043	22.1 / 26
LTR	ICLR’25	0.949 \pm 0.025	0.962 \pm 0.015	0.517 \pm 0.029	0.981 \pm 0.012	0.670 \pm 0.035	3.2 / 26
ROOT	NeurIPS’25 Spotlight	0.958 \pm 0.012	0.971 \pm 0.005	0.451 \pm 0.032	0.977 \pm 0.015	0.653 \pm 0.030	7.8 / 26
DynAMO-Adam	ICML’25	0.113 \pm 0.085	0.789 \pm 0.059	0.413 \pm 0.106	0.719 \pm 0.142	0.556 \pm 0.090	24.0 / 26
UniSO	ICML’25	0.850 \pm 0.062	0.915 \pm 0.015	0.489 \pm 0.062	0.947 \pm 0.036	0.673 \pm 0.136	11.1 / 26
DUO	/	0.975 \pm 0.014	0.977 \pm 0.005	0.610 \pm 0.044	0.933 \pm 0.055	0.706 \pm 0.129	3.4 / 26

549 G Additional analysis: metadata guidance strength

550 The main text (Section 4) reports primary benchmarks, a qualitative reverse-diffusion diagnostic,
551 and held-out transfer; Table 5, additional discussion, and the full Design-Bench leaderboard appear
552 in Appendix F. The discrete cross-entropy ablation is in Appendix H; full diffusion-curve panels
553 are in Appendix I. Here we sweep the scalar classifier-free guidance weight w on the text-metadata
554 embedding, holding the remainder of the conditioning stack to its default settings. Other condition
555 splits can be ablated in code; those variants are not emphasized in the main narrative.

556 Table 7 aggregates mean \pm std over all evaluation runs in the ablation. Moderate weights tend to
557 balance continuous robotics tasks and DNA objectives; very large weights occasionally hurt GTOPX

558 instances, suggesting that overly aggressive semantic guidance can push sampling outside the offline
 559 data support—analogueous to excessive guidance in offline RL.

Table 7: Ablation of classifier-free guidance strength on the text-metadata embedding (w): mean \pm std of maximum episode reward over evaluation runs. Baselines: $\mathcal{D}(\text{best})$, UniSO, GTG (single-task).

Task	$\mathcal{D}(\text{best})$	UniSO	GTG ST	$w=0$	$w=1$	$w=2$	$w=4$	$w=8$	$w=16$	$w=32$	$w=64$
Ant	165.326	455.658 \pm 39.188	452.330 \pm 61.502	568.346 \pm 10.406	574.388 \pm 8.549	566.263 \pm 8.668	561.212 \pm 5.810	566.659 \pm 9.082	564.747 \pm 3.107	567.857 \pm 8.753	565.456 \pm 7.107
D’Kitty	199.363	222.007 \pm 33.677	264.857 \pm 17.571	311.384 \pm 3.566	310.822 \pm 1.918	308.203 \pm 8.994	310.155 \pm 11.200	307.139 \pm 12.211	311.659 \pm 9.848	305.326 \pm 13.748	304.468 \pm 14.558
Superconductor	74.000	82.642 \pm 3.467	89.802 \pm 6.763	90.914 \pm 3.655	90.833 \pm 3.615	92.202 \pm 3.700	98.833 \pm 6.457	104.827 \pm 7.582	108.596 \pm 12.243	104.092 \pm 8.325	106.923 \pm 7.899
TF Bind 8	0.439	0.857 \pm 0.069	0.934 \pm 0.030	0.718 \pm 0.095	0.718 \pm 0.095	0.780 \pm 0.071	0.780 \pm 0.061	0.800 \pm 0.091	0.717 \pm 0.113	0.821 \pm 0.034	0.816 \pm 0.116
TF Bind 10	0.005	0.944 \pm 0.794	0.632 \pm 0.106	0.627 \pm 0.227	0.627 \pm 0.227	0.650 \pm 0.249	0.652 \pm 0.300	0.730 \pm 0.356	0.597 \pm 0.108	0.649 \pm 0.082	0.663 \pm 0.058
GTOPX 2	-196.663	-90.106 \pm 8.223	-72.163 \pm 8.495	-94.874 \pm 24.696	-99.556 \pm 25.617	-126.143 \pm 27.620	-82.143 \pm 37.817	-48.779 \pm 9.781	-65.063 \pm 37.999	-69.415 \pm 47.726	-76.584 \pm 38.446
GTOPX 3	-152.364	-44.427 \pm 13.456	-53.866 \pm 9.877	-85.963 \pm 32.525	-86.938 \pm 33.593	-65.658 \pm 15.598	-48.934 \pm 26.540	-29.687 \pm 1.663	-34.348 \pm 9.069	-29.687 \pm 1.663	-47.244 \pm 24.893
GTOPX 4	-216.709	-74.779 \pm 11.032	-72.749 \pm 6.632	-136.050 \pm 84.354	-135.972 \pm 84.455	-80.577 \pm 12.241	-81.501 \pm 26.936	-49.997 \pm 7.345	-49.997 \pm 7.345	-49.997 \pm 7.345	-57.299 \pm 17.860
GTOPX 6	-111.666	-48.244 \pm 5.766	-55.834 \pm 9.655	-67.402 \pm 4.817	-69.440 \pm 4.137	-65.996 \pm 9.502	-48.609 \pm 17.321	-30.053 \pm 1.303	-30.053 \pm 1.303	-34.417 \pm 10.093	-37.325 \pm 15.452
Mean rank	/	6.000 \pm 3.317	6.556 \pm 2.833	7.111 \pm 3.008	7.444 \pm 3.196	6.667 \pm 1.732	6.000 \pm 1.500	3.889 \pm 1.746	4.167 \pm 3.775	3.611 \pm 1.799	4.556 \pm 1.740

560 H Discrete cross-entropy auxiliary loss

561 Discrete tasks are often trained with continuous relaxations in latent space; mean-squared trajec-
 562 tory losses alone may not penalize one-hot decoding errors sharply enough. Table 8 compares
 563 multitask DUO with text metadata when $\lambda_{\text{CE}}=0$ versus $\lambda_{\text{CE}}=0.005$ in the auxiliary term (24). All
 564 entries use the same classifier-free guidance strength $w=8$ on the text-metadata embedding; other
 565 hyperparameters are held fixed across the two columns.

Table 8: Auxiliary cross-entropy ablation on multitask DUO with text metadata (maximum episode reward, mean \pm std; $\lambda_{\text{CE}}=0$ vs. 0.005). **Bold** marks TF Bind rows: CE sharpens discrete-task performance, while co-trained continuous tasks show only small shifts relative to the reported dispersion (no consistent degradation).

Task	$\lambda_{\text{CE}}=0$	$\lambda_{\text{CE}}=0.005$
Ant	559.863 \pm 8.221	565.873 \pm 13.651
D’Kitty	313.803 \pm 4.345	312.350 \pm 6.353
Superconductor	108.143 \pm 8.186	112.880 \pm 8.088
TF Bind 8	0.855 \pm 0.052	0.933 \pm 0.055
TF Bind 10	0.642 \pm 0.174	0.955 \pm 0.516
GTOPX 2	-44.477 \pm 4.761	-44.477 \pm 4.761
GTOPX 3	-29.310 \pm 2.237	-29.310 \pm 2.237
GTOPX 4	-54.560 \pm 14.187	-51.019 \pm 6.536
GTOPX 6	-28.580 \pm 2.489	-33.449 \pm 12.583
Mean rank	1.667 \pm 0.433	1.333 \pm 0.433

566 Turning on the cross-entropy term yields clear improvements on the discrete TF Bind 8/10 rows in
 567 Table 8, while Ant, D’Kitty, Superconductor, and GTOPX entries move only modestly and remain
 568 on the same scale once standard deviations are considered—consistent with CE acting as a targeted
 569 correction for discrete decoding rather than a global regularizer that reshapes the entire multitask
 570 solution. The mean rank row shifts slightly toward $\lambda_{\text{CE}}=0.005$, matching the view that the auxiliary
 571 loss primarily aligns latent training with discrete evaluation.

572 I Additional analysis: reverse-diffusion diagnostic curves

573 This appendix complements Section 4.3 with the full 3×6 panel of reverse-diffusion diagnostics
 574 on **Superconductor** (evaluation seed 0), using the same four regimes as Figure 2: **ST**, **ST+text**,
 575 **MT+label**, and **MT+text** (our method: multitask with full text-metadata embeddings). Figure 2
 576 shows only the **oracle max** column; Figure 3 adds **mean** and **top-8** alongside **max** for rows R1–R6
 577 (checkpoints in serialized export order). All runs use learning rate 2×10^{-4} . Section 4.3 interprets
 578 regime differences, checkpoint aging, and the oracle-**max** shape; here we focus on how **mean** and
 579 **top-8** should be read.

580 The horizontal axis counts reverse-diffusion steps from noisy latents to a denoised trajectory of
 581 horizon $H=64$. **Mean** averages the oracle objective over all generated trajectory positions *except*
 582 any context tokens from the conditioning stack, summarizing *global* trajectory quality rather than a

583 single extremum. **Top-8** averages the oracle values of the eight highest-scoring positions (top 8/ H
584 by oracle at the current step), which stresses whether strong objectives appear in multiple coordinates
585 instead of one outlier slot. **Max** is the best oracle anywhere on the partial chain at the current step
586 (the column duplicated in the main figure); we do not reinterpret its mid-chain peaks here. These
587 traces call the *true* simulator oracle for visualization only and omit the deployment-time surrogate
588 filter, so they are not directly comparable to proxy-ranked entries in the main tables.

589 When **top-8** rises in step with or ahead of **mean**, high oracle mass is spreading across the trajectory
590 rather than concentrating in a lone spike; flat **mean** with volatile **top-8** instead suggests occasional
591 lucky coordinates without broad improvement. Comparing the two columns across R1–R6 therefore
592 highlights whether later checkpoints improve *distributed* quality or mainly sharpen a few positions;
593 we defer narrative about training duration and overfitting to Section 4.3.



Figure 3: **Reverse-diffusion diagnostics (illustrative).** Superconductor, seed 0: oracle mean, top-8 mean, and max along reverse-diffusion steps (see text). Rows R1–R6 are exported training checkpoints from early to late in dashboard order; learning rate 2×10^{-4} . Curves compare ST, ST+text, MT+label, and MT+text as defined in Figure 2.