

Table S1. Detailed collection list of peer-reviewed papers regarding offline BBO.

Paper title	Venue	Year
Conditioning by adaptive sampling for robust design.	ICML	2019
Autofocused oracles for model-based design.	NeurIPS	2020
Model inversion networks for model-based optimization.	NeurIPS	2020
Offline model-based optimization via normalized maximum likelihood estimation.	ICLR	2021
Conservative objective models for effective offline model-based optimization.	ICML	2021
RoMA: Robust model adaptation for offline model-based optimization.	NeurIPS	2021
Data-driven offline optimization for architecting hardware accelerators.	ICLR	2022
Design-Bench: Benchmarks for data-driven offline model-based optimization.	ICML	2022
Data-driven offline decision-making via invariant representation learning.	NeurIPS	2022
Bidirectional learning for offline infinite-width model-based optimization.	NeurIPS	2022
Generative pretraining for black-box optimization.	ICML	2023
Diffusion models for black-box optimization.	ICML	2023
Bidirectional learning for offline model-based biological sequence design.	ICML	2023
Parallel-mentoring for offline model-based optimization.	NeurIPS	2023
Importance-aware co-teaching for offline model-based optimization.	NeurIPS	2023
ExPT: Synthetic pretraining for few-shot experimental design.	NeurIPS	2023
Bootstrapped training of score-conditioned generator for offline design of biological sequences.	NeurIPS	2023
Degradation-resistant offline optimization via accumulative risk control.	ECAI	2023
Offline model-based optimization via policy-guided gradient search.	AAAI	2024
Functional graphical models: Structure enables offline data-driven optimization.	AISTATS	2024
Learning surrogates for offline black-box optimization via gradient matching.	ICML	2024
Boosting offline optimizers with surrogate sensitivity.	ICML	2024
Offline multi-objective optimization.	ICML	2024
Incorporating surrogate gradient norm to improve offline optimization techniques.	NeurIPS	2024
Generative adversarial model-based optimization via source critic regularization.	NeurIPS	2024
Guided trajectory generation with diffusion models for offline model-based optimization.	NeurIPS	2024
Exploring validation metrics for offline model-based optimisation with diffusion models.	TMLR	2024
Robust guided diffusion for offline black-box optimization.	TMLR	2024
Offline model-based optimization by learning to rank.	ICLR	2025
SOO-Bench: Benchmarks for evaluating the stability of offline black-box optimization.	ICLR	2025
ParetoFlow: Guided flows in multi-objective optimization.	ICLR	2025

Table S2. Comparison of different model-inner BBO optimizers (BO- $q$ EI, CMA-ES, and EA) in improved UniSO-T. As CMA-ES cannot operate in categorical space, we conduct experimental comparison on continuous tasks from SOO-Bench, where the best and runner-up results on each task are **Blue** and **Violet**.  $\mathcal{D}(\text{best})$  denotes the best score in the offline dataset. Here all optimizers are allowed for an evaluation budget of 1000 (i.e., the optimizer can access the model’s outputs for upmost 1000 times) for fair comparison.

Task	$\mathcal{D}(\text{best})$	BO- $q$ EI	CMA-ES	EA
GTOPX 2	-195.586	<b>-80.220 <math>\pm</math> 11.852</b>	-138.716 $\pm$ 45.523	<b>-99.778 <math>\pm</math> 18.512</b>
GTOPX 3	-151.190	<b>-48.493 <math>\pm</math> 3.745</b>	-81.524 $\pm$ 19.054	<b>-69.278 <math>\pm</math> 20.009</b>
GTOPX 4	-215.716	<b>-80.232 <math>\pm</math> 13.582</b>	-177.559 $\pm$ 29.664	<b>-132.543 <math>\pm</math> 38.303</b>
GTOPX 6	-112.599	<b>-73.306 <math>\pm</math> 12.892</b>	-101.525 $\pm$ 16.986	<b>-62.024 <math>\pm</math> 23.828</b>
Avg. Rank	/	<b>1.250 <math>\pm</math> 0.433</b>	3.000 $\pm$ 0.000	<b>1.750 <math>\pm</math> 0.433</b>

Table S3. Ablation studies on different metadata components (name, description, and objective) on various tasks. All experiments are conducted based on improved UniSO-T. The first part includes tasks appear in training dataset, i.e., unconstrained tasks from Design-Bench and SOO-Bench. The remaining two parts contain unseen tasks during training, where we compare them under both zero-shot and few-shot settings. The best and runner-up results on each task are **Blue** and **Violet**.  $\mathcal{D}(\text{best})$  denotes the best score in the offline dataset.

Task	$\mathcal{D}(\text{best})$	UniSO-T	w/o name	w/o desc.	w/o obj.	w/o metadata
Ant	165.326	<b>374.665 <math>\pm</math> 56.057</b>	345.369 $\pm$ 40.675	362.842 $\pm$ 33.008	<b>363.812 <math>\pm</math> 59.380</b>	358.379 $\pm$ 64.211
D’Kitty	199.363	225.752 $\pm$ 8.521	<b>235.715 <math>\pm</math> 11.617</b>	226.479 $\pm$ 15.311	<b>245.135 <math>\pm</math> 15.571</b>	227.169 $\pm$ 12.278
Superconductor	74.000	92.200 $\pm$ 15.209	85.207 $\pm$ 6.261	<b>99.113 <math>\pm</math> 12.742</b>	<b>97.610 <math>\pm</math> 10.297</b>	90.871 $\pm$ 10.611
TF Bind 8	0.439	0.903 $\pm$ 0.041	<b>0.954 <math>\pm</math> 0.025</b>	0.935 $\pm$ 0.041	0.937 $\pm$ 0.012	<b>0.950 <math>\pm</math> 0.025</b>
TF Bind 10	0.005	<b>0.823 <math>\pm</math> 0.542</b>	<b>0.696 <math>\pm</math> 0.126</b>	0.664 $\pm$ 0.141	0.596 $\pm$ 0.148	0.651 $\pm$ 0.121
GTOPX 2	-195.586	<b>-72.848 <math>\pm</math> 9.576</b>	-97.806 $\pm$ 40.646	<b>-63.670 <math>\pm</math> 20.381</b>	-80.789 $\pm$ 8.908	-79.864 $\pm$ 13.338
GTOPX 3	-151.190	<b>-45.602 <math>\pm</math> 8.433</b>	-50.788 $\pm$ 8.706	<b>-45.981 <math>\pm</math> 4.211</b>	-50.660 $\pm$ 10.713	-48.178 $\pm$ 12.638
GTOPX 4	-215.716	-84.271 $\pm$ 8.307	-84.962 $\pm$ 11.300	-92.163 $\pm$ 9.529	<b>-75.233 <math>\pm</math> 5.734</b>	<b>-79.887 <math>\pm</math> 14.729</b>
GTOPX 6	-112.599	-47.794 $\pm$ 11.943	<b>-42.181 <math>\pm</math> 11.671</b>	<b>-45.591 <math>\pm</math> 12.310</b>	-48.050 $\pm$ 13.901	-45.764 $\pm$ 7.685
RobotPush (zero-shot)	0.102	<b>3.171 <math>\pm</math> 0.984</b>	2.747 $\pm$ 1.455	<b>3.416 <math>\pm</math> 1.455</b>	2.634 $\pm$ 0.953	2.517 $\pm$ 1.640
Rover (zero-shot)	-16.148	<b>-8.888 <math>\pm</math> 2.119</b>	-11.009 $\pm$ 0.598	-10.854 $\pm$ 0.859	<b>-9.099 <math>\pm</math> 2.202</b>	-9.089 $\pm$ 3.070
LunarLander (zero-shot)	7.038	<b>31.186 <math>\pm</math> 27.971</b>	30.105 $\pm$ 57.577	30.892 $\pm$ 54.657	<b>52.108 <math>\pm</math> 47.941</b>	6.251 $\pm$ 53.042
RobotPush (few-shot)	0.102	<b>7.067 <math>\pm</math> 0.169</b>	7.026 $\pm$ 0.219	<b>7.129 <math>\pm</math> 0.486</b>	6.310 $\pm$ 1.677	6.155 $\pm$ 1.495
Rover (few-shot)	-16.148	<b>-8.239 <math>\pm</math> 1.270</b>	-8.850 $\pm$ 0.703	<b>-8.084 <math>\pm</math> 0.569</b>	-8.342 $\pm$ 1.573	-10.511 $\pm$ 2.070
LunarLander (few-shot)	7.038	<b>248.573 <math>\pm</math> 45.386</b>	226.726 $\pm$ 65.244	<b>244.252 <math>\pm</math> 38.329</b>	233.169 $\pm$ 51.037	233.919 $\pm$ 60.467
Avg. Rank	/	<b>2.333 <math>\pm</math> 1.350</b>	3.600 $\pm$ 1.451	<b>2.467 <math>\pm</math> 1.310</b>	3.067 $\pm$ 1.340	3.533 $\pm$ 1.087

---

*Table S4.* References of response to Reviewer 4Ajl.

No.	Paper title	Venue	Year
1	Position: Leverage foundational models for black-box optimization.	ICML	2024
2	A systematic survey on large language models for algorithm design.	arXiv	2024
3	LLM4AD: A platform for algorithm design with large language model.	arXiv	2024
4	Mathematical discoveries from program search with large language models.	Nature	2024
5	Evolution of heuristics: Towards efficient automatic algorithm design using large language model.	ICML	2024
6	Can GPT-4 perform neural architecture search?	arXiv	2023
7	OmniPred: Language models as universal regressors.	TMLR	2024
8	Understanding LLM embeddings for regression.	TMLR	2025
9	Predicting from strings: Language model embeddings for Bayesian optimization.	arXiv	2024
10	Towards learning universal hyperparameter optimizers with Transformers.	NeurIPS	2022
11	LICO: Large language models for in-context molecular optimization.	ICLR	2025
12	Understanding the behaviour of contrastive loss.	CVPR	2021
13	Advancing Bayesian optimization via learning correlated latent space.	NeurIPS	2023
14	pymoo: Multi-objective optimization in python.	IEEE Access	2020
15	Offline multi-objective optimization.	ICML	2024