

# Appendices

## A Data Collection

In this section, we detail the data collection steps used for creating each of the tasks in design-bench. We answer (1) where is the data from, and (2) what pre-processing steps are used?

### A.1 TF Bind 8

The TF Bind 8 task is a derivative of the transcription factor binding activity survey performed by [5], where the binding activity scores of every possible length eight DNA sequence was measured with a variety of human transcription factors. We filter the dataset by selecting a particular transcription factor SIX6\_REF\_R1, and defining an optimization problem where the goal is to synthesize a length 8 DNA sequence with high binding activity with human transcription factor SIX6\_REF\_R1. This particular transcription factor for TF Bind 8 was recently used for optimization in [2, 3]. TF Bind 8 is a fully characterized dataset containing 65792 samples, representing every possible length 8 combination of nucleotides  $\mathbf{x}_{\text{TFBind8}} \in \{0, 1\}^{8 \times 4}$ . The training set given to offline MBO algorithms is restricted to the bottom 50%, which results in a visible training set of 32898 samples.

### A.2 GFP

The GFP task provided is a derivative of the GFP dataset [30]. The dataset we use in practice is that provided by [7] at the url <https://github.com/dhbrookes/CbAS/tree/master/data>. We process the dataset such that a single training example consists of a protein represented as a tensor  $\mathbf{x}_{\text{GFP}} \in \{0, 1\}^{237 \times 20}$ . This tensor is a sequence of 237 one-hot vectors corresponding to which amino acid is present in that location in the protein. We use the dataset format of [7] with no additional processing. The data was originally collected by performing laboratory experiments constructing proteins similar to the Aequorea victoria green fluorescent protein and measuring fluorescence. We employ the full dataset of 56086 proteins when learning approximate oracles for evaluating offline MBO methods, but restrict the training set given to offline MBO algorithms to 5000 samples drawn from between the 50th percentile and 60th percentile of proteins in the GFP dataset, sorted by fluorescence values. This subsampling procedure is consistent with prior work [7].

### A.3 UTR

The UTR task is derived from work by Sample et al. [29] who trained a CNN model to predict the expressive level of a particular gene from a corresponding 5'UTR sequence. Our use of the UTR task for model-based optimization follows Angermüller et al. [3], where the goal is to design a length 50 DNA sequence to maximize expression level. We follow the methodology set by Sample et al. [29] to sort all length 50 DNA sequences in the unprocessed UTR dataset by total reads, and then select the top 280,000 DNA sequences with the most total reads. The result is a dataset containing 280,000 samples of length 50 DNA sequences  $\mathbf{x}_{\text{UTR}} \in \{0, 1\}^{50 \times 4}$  and corresponding ribosome loads. When training offline MBO algorithms, we subsequently eliminate the top 50% of sequences ranked by their ribosome load, resulting in a visible dataset with only 140,000 samples.

### A.4 ChEMBL

The ChEMBL task is a derivative of a much larger dataset that is derived from ChEMBL [13], a large database of chemicals and their properties. The data, similar to GFP, was originally collected by performing physical experiments on a large number of molecules, and measuring their activity with a target assay. We have processed the ChEMBL database—available at <https://www.ebi.ac.uk/chembl/g/#browse/activities>—into collections of smaller datasets mapping particular molecules to measured values, determined by a target assay that accompanies each set. We choose the assay specified by `ASSAY_CHEMBL_ID = CHEMBL1964047` and select the standard type of GI50 as the measurement to maximize with offline model-based optimization. The resulting dataset has 40516 samples in total. We preprocess the dataset by converting each molecule into a SMILES string using RDKit, and then apply the DeepChem SmilesTokenizer to convert each SMILES string into

a sequence of integer tokens. The resulting processed samples have a maximum sequence length of 425 and the vocabulary has 591 elements,  $\mathbf{x}_{\text{ChEMBL}} \in \{0, 1\}^{425 \times 591}$ . The final step is to remove the top 47% of molecules sorted by their GI50 value. This step removes a mode of the dataset that we found to hurt the performance of learned oracle models trained on this processed dataset.

## A.5 Superconductor

The Superconductor task is inspired by recent work [10] that applies offline MBO to optimize the properties of superconducting materials for high critical temperature. The data we provide in our benchmark is real-world superconductivity data originally collected by [16], and subsequently made available to the public at <https://archive.ics.uci.edu/ml/datasets/Superconductivity+Data#>. The original dataset consists of superconductors featurized into vectors containing measured physical properties like the number of chemical elements present, or the mean atomic mass of such elements. One issue with the original dataset that was used in [10] is that the numerical representation of the superconducting materials did not lend itself to recovering a physically realizable material that could be synthesized in a lab after performing model-based optimization. In order to create an *invertible* input specification, we deviate from prior work and encode superconductors as vectors whose components represent the number of atoms of specific chemical elements present in the superconducting material—a serialization of the chemical formula of each superconductor. The result is a real-valued design space with 86 components  $\mathbf{x}_{\text{Superconductor}} \in \mathcal{R}^{86}$ . The full dataset used to learn approximate oracles for evaluating MBO methods has 21263 samples, but we restrict this number to 17010 (the 80th percentile) for the training set of offline MBO methods to increase difficulty.

## A.6 Hopper Controller

The HopperController task is one that we provide ourselves. The goal of this task is to design a set of weights for a neural network policy, in order to achieve high expected return when evaluating that policy. The data collected for HopperController was taken by training a three layer neural network policy with 64 hidden units and 5126 total weights on the Hopper-v2 MuJoCo task using Proximal Policy Optimization [31]. Specifically, we use the default parameters for PPO provided in stable baselines [18]. The dataset we provide with this benchmark has 3200 unique weights. In order to collect this many, we run 32 experimental trials of PPO, where we train for one million steps, and save the weights of the policy every 10,000 environment steps. The policy weights are represented originally as a list of tensors. We first traverse this list and flatten each of the tensors, and we then concatenate each of these flattened tensors into a single training example  $\mathbf{x}_{\text{Hopper}} \in \mathcal{R}^{5126}$ . The result is an optimization problem over neural network weights. After collecting these weights, we perform no additional pre-processing steps. In order to collect scores we perform a single rollout for each  $x$  using the Hopper-v2 MuJoCo environment. The horizon length for training and evaluation is limited to 1000 simulation time steps, which is standard practice for this MuJoCo environment.

## A.7 Ant & D’Kitty Morphology

Both morphology tasks are collected by us, and share methodology. The goal of these tasks is to design the morphology of a quadrupedal robot—an ant or a D’Kitty—such that the agent is able to crawl quickly in a particular direction. In order to collect data for this environment, we create variants of the MuJoCo Ant and the ROBEL D’Kitty agents that have parametric morphologies. The goal is to determine a mapping from the morphology of the agent to the average return of a pre-trained morphology conditioned agent. We implement this by pre-training a morphology conditioned neural network policy using SAC [15]. For both the Ant and the D’Kitty, we train the agents for more than ten million environment steps, and a maximum episode length of 200, with all other settings as default. These agents are pre-trained on Gaussian distributions of morphologies. The Gaussian distributions are obtained by adding Gaussian noise with standard deviation 0.03 for Ant and 0.01 for D’Kitty the design-space range to the default morphologies. After obtaining trained morphology-conditioned policies, we create a dataset of morphologies for model-based optimization by sampling initialization points randomly, and then using CMA-ES to optimize for morphologies that attain high reward using the pretrained morphology-conditioned policy. To obtain initialization points, we add Gaussian random noise to the default morphology for the Ant with standard deviation 0.075 and D’Kitty with standard deviation 0.1, and then apply CMA-ES with standard deviation 0.02. We ran CMA-ES for 250 iterations and then restart, until a minimum of 25000 morphologies were collected, resulting in a

730 final dataset size of 25009 for both the Ant and D’Kitty. The design space for Ant Morphologies is  
731  $\mathbf{x}_{\text{Ant}} \in \mathcal{R}^{60}$ , whereas for D’Kitty morphologies is  $\mathbf{x}_{\text{D’Kitty}} \in \mathcal{R}^{56}$ . We subsample the dataset to its  
732 40th percentile when training offline MBO algorithms, resulting in 10004 samples.

## 733 B Oracle Functions

734 We detail oracle functions for evaluating ground truth scores for each of the tasks in design-bench. A  
735 common thread among these is that the oracle, if trained, is fit to a larger static dataset containing  
736 higher performing designs than observed by a downstream MBO algorithm.

### 737 B.1 TF Bind 8

738 TF Bind 8 is a fully characterized discrete offline MBO task, which means that all possible designs  
739 have been evaluated [5] and are contained in the full hidden TF Bind 8 dataset. The oracle for TF  
740 Bind 8 is therefore implemented as a lookup table that returns the score corresponding to a particular  
741 length 8 DNA sequence from the dataset. By restricting the size of the training set visible to an offline  
742 MBO algorithm, it is possible for the algorithm to propose a design that achieves a higher score than  
743 any other DNA sequence visible to the offline MBO algorithm during training.

### 744 B.2 GFP

745 GFP uses the oracle function derived from Rao et al. [27]. This oracle is a Transformer regression  
746 model with 4 attention blocks and a hidden size of 64. The Transformer is fit to the entire hidden  
747 GFP dataset, making it possible to sample a protein design that achieves a higher score than any other  
748 protein visible to an offline MBO algorithm. Our Transformer has a Spearman’s rank correlation  
749 coefficient of 0.8497 with a held-out validation set derived from the GFP dataset.

### 750 B.3 UTR

751 UTR uses a Transformer as the oracle function, which differs from the CNN that was originally  
752 used by [3]. Our reasoning for making this change is that the Transformer is a newer and possibly  
753 higher capacity model that may be less prone to mistakes than the shallower CNN model proposed by  
754 Sample et al. [29]. This Transformer has 4 attention blocks and a hidden size of 64. The Transformer  
755 is fit to the entire hidden UTR dataset, making it possible to sample a DNA sequence that achieves a  
756 higher score than any other sequence visible to an offline MBO algorithm. The resulting Transformer  
757 has a spearman’s rank correlation coefficient of 0.6424 with a held-out validation set.

### 758 B.4 ChEMBL

759 We tested several models as candidate oracle functions for ChEMBL [13], including Gaussian Process,  
760 Random Forest, CNN, and Transformer regression models. We ultimately found the CNN to result in  
761 the highest validation performance, achieving a spearman’s rank correlation coefficient of 0.3208  
762 with a held-out validation set. These models were trained on the entire hidden ChEMBL dataset  
763 encoded into SMILES and tokenized. While this rank correlation is low in comparison to the previous  
764 tasks, such illustrates the challenge learning of predicting molecule properties directly from SMILES,  
765 and is a welcome avenue for future work in the design an training of such models. Our CNN consists  
766 of four residual blocks each with two convolution layers with kernel size 3 and 64 filters.

### 767 B.5 Superconductor

768 The Superconductor oracle function is also a random forest regression model. The model we use  
769 is the model described by [16]. We borrow the hyperparameters described by them, and we use  
770 the RandomForestRegressor provided in scikit-learn. Similar to the setup for the previous set of  
771 tasks, this oracle is trained on the entire hidden dataset of superconductors. The random forest has a  
772 spearman’s rank correlation coefficient with a held-out validation set of 0.9155.

## 773 B.6 HopperController

774 Unlike the previously described tasks, HopperController and the remaining tasks implement an  
775 exact oracle function. For HopperController the oracle takes the form of a single rollout using the  
776 Hopper-v2 MuJoCo environment. The designs for HopperController are neural network weights, and  
777 during evaluation, a policy with those weights is instantiated—in this case that policy is a three layer  
778 neural network with 11 input units, two layers with 64 hidden units, and a final layer with 3 output  
779 units. The intermediate activations between layers are hyperbolic tangents. After building a policy,  
780 the Hopper-v2 environment is reset and the reward for 1000 time-steps is summed. That summed  
781 reward constitutes the score returned by the HopperController-v0 oracle. The limit of performance is  
782 the maximum return that an agent can achieve in Hopper-v2 over 1000 steps.

## 783 B.7 Ant & D’Kitty Morphology

784 The final two tasks in design-bench use an exact oracle function, using the MuJoCo simulator. For  
785 both morphology tasks, the simulator performs a rollout and returns the sum of rewards at every  
786 timestep in that rollout. Each task is accompanied by a pre-trained morphology-conditioned policy. To  
787 perform evaluation, a morphology is passed to the Ant or D’Kitty MuJoCo environments respectively,  
788 and a dynamic-morphology agent is initialized inside these environments. These simulations can be  
789 time consuming to run, and so we limit the rollout length to 100 steps. The morphology conditioned  
790 policies were trained using the reinforcement learning algorithm SAC for 10 million steps for each  
791 task, and are ReLU networks with two hidden layers of size 64.

## 792 C Experimental Details

793 In this section we present additional details for the experiments, including the score normalization  
794 process and 50th percentile performance.

### 795 C.1 Score Normalization

796 In order to report performance on the same order of magnitude for each offline model-based opti-  
797 mization task in Design-Bench, we normalize the performance reported in Table 2 by calculating  
798 the minimum objective value  $y_{\min}$  and the the maximum objective value  $y_{\max}$  in the full unobserved  
799 dataset associated with each offline model-based optimization problem. *Crucially*, note that this is  
800 not the same as normalizing with respect to the best and worst samples in the training dataset used  
801 by the offline MBO algorithm, but rather a bigger dataset of designs and objective values. We then  
802 report performance by calculating what fraction of the distance between  $y_{\min}$  and  $y_{\max}$  is attained by  
803 a particular offline MBO baseline.

$$y_{\text{normalized}}(y) = \frac{y - y_{\min}}{y_{\max} - y_{\min}} \quad (3)$$

804 The final performance  $y_{\text{normalized}}$  is the normalized performance of an offline MBO method that  
805 achieved an unprocessed objective value of  $y$ . The result is larger than one when the offline MBO  
806 method finds a solution more performance than all solutions in the full unobserved dataset associated  
807 with the corresponding task. The result is less than zero when the offline MBO method finds a  
808 solution attaining less performance than all samples in the full unobserved dataset.

### 809 C.2 50th Percentile Experiment Results

810 In this section, we present the 50th percentile performance of the runs presented in main paper  
811 in Table 2. Similar to the 100th percentile performance reported in the main text, performance is  
812 calculated by evaluating solutions to each task found by an optimization method, subtracting the  
813 minimum objective value present in the corresponding task dataset, and dividing by the range of  
814 objective values present in the corresponding task dataset. The result is a performance of greater than  
815 one if optimization converges to a solution with a higher objective value than the best observed design  
816 in the corresponding task dataset.

	GFP	TF Bind 8	UTR	ChEMBL
Autofocus.-CbAS	0.848 $\pm$ 0.007	0.419 $\pm$ 0.007	0.517 $\pm$ 0.010	0.216 $\pm$ 0.000
CbAS	0.852 $\pm$ 0.004	0.428 $\pm$ 0.010	0.515 $\pm$ 0.014	0.243 $\pm$ 0.006
BO-qEI	0.246 $\pm$ 0.341	0.439 $\pm$ 0.000	0.574 $\pm$ 0.000	0.298 $\pm$ 0.034
CMA-ES	0.047 $\pm$ 0.000	0.537 $\pm$ 0.014	0.385 $\pm$ 0.019	0.316 $\pm$ 0.031
Gradient Ascent	0.838 $\pm$ 0.004	0.609 $\pm$ 0.019	0.489 $\pm$ 0.009	0.300 $\pm$ 0.031
Grad. Min.	0.837 $\pm$ 0.001	0.645 $\pm$ 0.030	0.514 $\pm$ 0.007	0.306 $\pm$ 0.023
Grad. Mean	0.838 $\pm$ 0.002	0.616 $\pm$ 0.023	0.508 $\pm$ 0.005	0.326 $\pm$ 0.016
MINs	0.820 $\pm$ 0.018	0.421 $\pm$ 0.015	0.560 $\pm$ 0.007	0.347 $\pm$ 0.001
REINFORCE	0.844 $\pm$ 0.003	0.462 $\pm$ 0.021	0.530 $\pm$ 0.008	0.246 $\pm$ 0.004
	Superconductor	Ant Morphology	D’Kitty Morphology	Hopper Controller
Autofocus.-CbAS	0.131 $\pm$ 0.010	0.362 $\pm$ 0.015	0.736 $\pm$ 0.025	0.019 $\pm$ 0.008
CbAS	0.111 $\pm$ 0.017	0.382 $\pm$ 0.016	0.753 $\pm$ 0.008	0.015 $\pm$ 0.002
BO-qEI	0.300 $\pm$ 0.015	0.566 $\pm$ 0.000	0.883 $\pm$ 0.000	0.343 $\pm$ 0.010
CMA-ES	0.379 $\pm$ 0.003	-0.050 $\pm$ 0.004	0.684 $\pm$ 0.016	-0.033 $\pm$ 0.005
Gradient Ascent	0.476 $\pm$ 0.022	0.130 $\pm$ 0.018	0.509 $\pm$ 0.200	0.092 $\pm$ 0.084
Grad. Min	0.471 $\pm$ 0.016	0.182 $\pm$ 0.008	0.746 $\pm$ 0.034	0.222 $\pm$ 0.065
Grad. Mean	0.469 $\pm$ 0.022	0.184 $\pm$ 0.010	0.748 $\pm$ 0.024	0.243 $\pm$ 0.064
MINs	0.336 $\pm$ 0.016	0.619 $\pm$ 0.041	0.887 $\pm$ 0.004	0.352 $\pm$ 0.058
REINFORCE	0.463 $\pm$ 0.016	0.134 $\pm$ 0.033	0.356 $\pm$ 0.131	-0.064 $\pm$ 0.003

Table 3: **50th percentile** evaluations for baselines on every task. Results are averaged over 8 trials, and the  $\pm$  indicates the standard deviation of the reported performance. This table corresponds to the normalized performance, using the normalization methodology described in Appendix C.1

### C.3 Unnormalized Experimental Results

In this section, we present the raw 100th percentile performance of the runs presented in main paper in Table 2. These values, presented in Table 4, represent the mean raw objective values and the standard deviation of the objective values attained by various offline MBO methods.

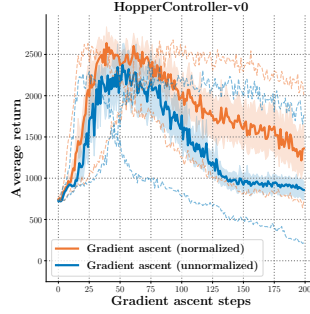
	GFP	TF Bind 8	UTR	ChEMBL
Autofocus CbAS	3.739 $\pm$ 0.001	0.910 $\pm$ 0.044	7.794 $\pm$ 0.070	42467.285 $\pm$ 0.000
CbAS	3.738 $\pm$ 0.001	0.927 $\pm$ 0.051	7.800 $\pm$ 0.023	46681.988 $\pm$ 4987.456
BO-qEI	2.005 $\pm$ 1.000	0.798 $\pm$ 0.083	7.903 $\pm$ 0.000	30069.684 $\pm$ 3187.300
CMA-ES	1.438 $\pm$ 0.005	0.953 $\pm$ 0.022	7.993 $\pm$ 0.052	31607.031 $\pm$ 1578.222
Gradient Ascent	3.737 $\pm$ 0.002	0.977 $\pm$ 0.025	7.669 $\pm$ 0.113	32514.541 $\pm$ 2612.903
Grad. Min	3.738 $\pm$ 0.000	0.984 $\pm$ 0.012	7.758 $\pm$ 0.081	32617.006 $\pm$ 370.390
Grad. Mean	3.738 $\pm$ 0.000	0.986 $\pm$ 0.012	7.764 $\pm$ 0.058	33715.059 $\pm$ 1136.034
MINs	3.741 $\pm$ 0.002	0.905 $\pm$ 0.052	7.787 $\pm$ 0.053	42732.578 $\pm$ 5126.862
REINFORCE	3.739 $\pm$ 0.001	0.948 $\pm$ 0.028	7.747 $\pm$ 0.059	41448.012 $\pm$ 3220.380
	Superconductor	Ant Morphology	D’Kitty Morphology	Hopper Controller
Autofocus CbAS	77.910 $\pm$ 8.361	474.888 $\pm$ 44.424	226.156 $\pm$ 7.043	261.820 $\pm$ 6.366
CbAS	93.078 $\pm$ 12.695	469.499 $\pm$ 30.570	209.412 $\pm$ 9.593	267.623 $\pm$ 15.833
BO-qEI	74.322 $\pm$ 6.347	413.084 $\pm$ 0.000	213.816 $\pm$ 0.000	788.990 $\pm$ 149.878
CMA-ES	86.072 $\pm$ 4.508	799.394 $\pm$ 715.702	4.290 $\pm$ 1.505	857.178 $\pm$ 273.980
Grad.	95.789 $\pm$ 4.436	-100.265 $\pm$ 22.118	187.206 $\pm$ 27.274	1406.413 $\pm$ 613.631
Grad. Min	93.590 $\pm$ 1.719	80.853 $\pm$ 62.308	205.639 $\pm$ 13.427	1860.192 $\pm$ 750.026
Grad. Mean	92.265 $\pm$ 3.206	48.064 $\pm$ 78.555	209.355 $\pm$ 13.928	2108.004 $\pm$ 578.350
MINs	86.702 $\pm$ 4.171	505.515 $\pm$ 34.934	273.479 $\pm$ 14.184	628.436 $\pm$ 210.953
REINFORCE	88.996 $\pm$ 2.389	-127.440 $\pm$ 30.831	-194.540 $\pm$ 238.857	62.183 $\pm$ 84.937

Table 4: **Unnormalized 100th percentile** unnormalized evaluations for baselines on every task. Results are averaged over 8 trials, and the  $\pm$  indicates the standard deviation of the reported performance. This table corresponds to the unnormalized performance.

## 821 C.4 Computation Resources

822 The amount of computation resources required to produce the experiments in this paper is relatively  
 823 modest. We ran our experiments on a single server with 2 Intel Xeon E5-2698 v4 CPUs and 8 Nvidia  
 824 Tesla V100 GPUs. All our experiments can be completed within 96 hours on this single machine.

## 825 D Normalization Of Inputs and Outputs For Gradient Ascent Baseline



826 An important component for the the good performance of gradient ascent baseline is the normalization  
 827 of design space. We found that the identical gradient-ascent baseline performed a factor 1.4x worse  
 828 on HopperController, when optimizing in the space of unnormalized designs and objective values, as  
 829 seen in Figure D. This indicates that normalization is key in obtaining good performance with a naïve  
 830 gradient ascent baseline.

831 For continuous design-space tasks, we normalize both the designs, and the scores to have unit  
 832 Gaussian statistics. For discrete design-space tasks, we normalize only the scores to have unit  
 833 Gaussian statistics. This is a necessary part of the optimization workflow because scores vary by  
 834 several orders of magnitude in the dataset, for example, 0.864 for GFP similar to other methods and  
 835 as high as 1.586 for HopperController. The specific normalization equation for continuous-valued  
 836 designs is given below.

$$\tilde{\mathbf{x}}_{i,j} = \frac{\mathbf{x}_{i,j} - \mu(\mathbf{x},j)}{\sigma(\mathbf{x},j)} : \mathbf{x} \in \mathbb{R}^{N \times D} \quad (4)$$

837 We also normalize the objective values in a similar fashion to have unit Gaussian statistics. The result  
 838 in a new set of designs  $\tilde{\mathbf{x}}$  and objective values  $\tilde{y}$  that is optimized over

$$\tilde{y}_{i,j} = \frac{y_{i,j} - \mu(y,j)}{\sigma(y,j)} : y \in \mathbb{R}^{N \times 1} \quad (5)$$

839 The gradient ascent procedure is performed in the space of these normalized designs. Suppose  $T$   
 840 steps of gradient ascent have been taken, and a final normalized solution  $\tilde{\mathbf{x}}_T^*$  is found. This solution is  
 841 de-normalized using the following transformation.

$$(\mathbf{x}_T^*)_{ij} = (\tilde{\mathbf{x}}_T^*)_{ij} \cdot \sigma(\mathbf{x},j) + \mu(\mathbf{x},j) \quad (6)$$

842 This normalization strategy is heavily inspired by data whitening, which is known to reduce the  
 843 variance of machine learning algorithms that learn discriminative mappings on that data. The learned  
 844 model of the objective function is one such discriminative model, and normalization likely improves  
 845 the consistency of Gradient Ascent across independent experimental trials.

## 846 E Hyperparameter Selection Workflow

847 Hyperparameter tuning under a restricted computational budget is emerging as an import research  
 848 domain in optimization [34, 9, 19]. Care must be taken when tuning each of the prescribed algorithms  
 849 so that only offline information about the task is used for hyperparameter selection. Formally, this  
 850 means that the hyperparameters,  $\mathcal{H}$ , are conditionally independent of the particular value of the  
 851 performance metric  $\mathcal{M}$ , given the offline task dataset  $\mathcal{D}$ . Examples of hyperparameter selection  
 852 strategies that violate this requirement might, for example, perform a grid search over  $\mathcal{H}$  and take the  
 853 set that maximizes the performance metric, but this is not offline. An example of a tuning strategy



that is fully offline is tuning the parameters of a learned model such that it is a good fit for the task dataset  $\mathcal{D}$ . One can choose  $\mathcal{H}$  that minimizes a validation loss, such as negative log likelihood. A detailed record of hyperparameters can be found in the experiment scripts located alongside our reference implementations: <https://github.com/brandontrabucco/design-baselines>.

We now present specific guidelines for hyperparameter selection (i.e. *workflow*) for some of the methods evaluated in our benchmark. These principles are general principles that can be used to tune the hyperparameters of these methods on a new task in a completely offline fashion. While we only present workflow details for methods we benchmark in Section 7, we expect that these general strategies will allow users to devise analogous schemes for tuning hyperparameters of new offline MBO methods with shared components.

## E.1 Strategy For Autofocused CbAS

The main tunable components of Autofocused methods [10] are the learned objective function, and the generative model fit to the data distribution. When training the learned objective function, tracking a validation performance metric like rank correlation is helpful to ensure that the resulting learned model is able to generalize beyond its training dataset. This tracking is especially important for Autofocused methods because re-fitting the learned objective model during importance can lead to divergence if the importance weights generated by Autofocusing are very large or very small in magnitude. The algorithm is tuned well if, for example, the validation rank correlation stays above a positive threshold, such as a threshold of 0.9.

The second component of Autofocused methods is the fit of the generative model used for sampling designs. The algorithm has the best chance of success if the generative model can generalize beyond the dataset in which it was trained. This can be monitored by holding out a validation set and tracking a metric such as negative log likelihood on this held-out set. In the case when the generative model is not an exact likelihood-based generative model—for example, a VAE—other validation metrics can be used that measure the fit of the generative model on a validation set. The generative model is especially impacted by the importance sampling procedure used by Estimation of Distribution Algorithms (EDAs), and tracking the effective sample size of the importance weights can help diagnose when the generative model is failing to generalize to a validation set.

## E.2 Strategy For CbAS

The main tunable components of CbAS methods [7] are the learned objective function, and the generative model fit to the data distribution. While the learned objective function is not affected by the importance sampling weights generated by CbAS, the same tuning strategy described in section E.1 that focuses on generalization to a validation set is effective. Generative model tuning can also follow an identical strategy to that described in section E.1, which focuses on the ability for the generative model to represent samples outside of its training set. In the case of a  $\beta$ -VAE, which is used with CbAS in this work, the main parameter for controlling this generalization ability is the  $\beta$  parameter. We found that  $\beta$  is task specific, and must be found in order for the CbAS optimizer using  $\beta$ -VAE to generate samples that are in the same distribution as its validation set. This value can be tuned in practice using a validation metric like that in section E.1.

## E.3 Strategy For MINs

The main tunable components of MINs [22] are the learned objective function, and the generative model fit to the data distribution. The learned objective function is typically trained using a maximum likelihood objective, and the validation log-likelihood (or regression error) can be directly tracked. The learned objective function should train until a minimum validation loss is reached, which ensured that the model will generalize as well as possible beyond its training set. Since only the static task dataset is used for this—it may be split into train/validation sets—this tuning strategy is fully offline.

The generative model for MINs is an inverse mapping  $\mathbf{x} = f^{-1}(y, \mathbf{z})$ , conditioned on the objective value  $y$ . Training conditional generative models is considerably less stable than unconditional generative models, so in addition to monitoring the fit of a validation set recommended in section E.1, it is also necessary to track the extent of the dependence of the generative model’s predictions on the objective value  $y$ . This can be evaluated in practice by comparing the distribution of  $x$  from the

conditional generative model  $p(\mathbf{x}|y)$  to an unconditional generative model  $p(\mathbf{x})$  with an identical initialization, or by comparing if  $p(\mathbf{x}|y)$  is independent of  $y$  by querying the inverse model for different values of  $y$  and visualizing the similarity in the predictions of  $\mathbf{x}$ . One metric for more formally studying the extent of the dependence of  $\mathbf{x}$  on  $\mathbf{z}$  is the mutual information  $I(\mathbf{x}; \mathbf{z})$ . The conditional generative model has an appropriate fit if for some positive threshold  $c$  we have that  $I(\mathbf{x}; \mathbf{z}) > c$ .

#### 911 **E.4 Strategy For Gradient Ascent**

912 The main tunable components of Gradient Ascent MBO methods are the learned objective function,  
913 and the parameters for gradient ascent. The learned objective function is typically trained using a  
914 maximum likelihood objective under a Gaussian distribution, and the methodology for obtaining a  
915 high-performing learned objective function is identical to that in section E.3. The second aspect of  
916 gradient ascent MBO algorithms are the parameters of the gradient-based optimizer for the designs—  
917 such as its learning rate, and the number of gradient steps it performs. The learning rate should be  
918 small enough that the gradient steps taken increase the prediction of the learned objective function—if  
919 the learning rate is too large, gradient steps may not follow the path of steepest ascent of the objective  
920 function. The number of gradient steps is more difficult to tune. The strategy we used is a fixed  
921 number of steps, and an offline model-selection criterion to select this parameter is future work.

#### 922 **E.5 Strategy For REINFORCE**

923 The main tunable components of REINFORCE-based MBO methods are the learned objective  
924 function, and the parameters for the policy gradient estimator. The learned objective function is  
925 typically trained using a maximum likelihood objective, and the methodology for obtaining a high-  
926 performing learned objective function is identical to that in section E.3. The remaining parameters  
927 to tune are specific to REINFORCE. The distribution of the policy should be carefully selected to  
928 be able to model the distribution of designs. For continuous MBO tasks, a Gaussian distribution is  
929 appropriate, and for discrete MBO tasks, a categorical distribution is appropriate. In addition, the  
930 learning rate, and optimizer should be selected so that policy updates improve the model-predicted  
931 score. If the stability of the gradient estimator is suffering, due to high-variance updates, a density  
932 ratio threshold like in PPO [31] can be applied, and a baseline can be subtracted.

#### 933 **E.6 Strategy For Bayesian Optimization**

934 The main tunable components of Bayesian Optimization MBO methods [4] are the learned objective  
935 function, and the parameters for the bayesian optimization loop. The learned objective function  
936 is typically trained using a maximum likelihood objective, and the methodology for obtaining a  
937 high-performing learned objective function is identical to that in section E.3. For a detailed review of  
938 the strengths and weaknesses of various Bayesian Optimization strategies and their hyperparameters,  
939 we refer the reader to the BoTorch documentation, available at the BoTorch website [https://  
940 botorch.org/docs/overview](https://botorch.org/docs/overview). In this work we employ a Gaussian Process as the model, and the  
941 quasi-Monte Carlo Expected Improvement acquisition function, which has the advantage of scaling  
942 up to our high-dimensional optimization problems.

#### 943 **E.7 Strategy For Covariance Matrix Adaptation (CMA-ES)**

944 The main tunable components of Covariance Matrix Adaptation MBO methods are the learned  
945 objective function, and the parameters for the evolution strategy. The learned objective function is  
946 typically trained using a maximum likelihood objective, and the methodology for obtaining a high-  
947 performing learned objective function is identical to that in Subsection E.3. For a detailed review of the  
948 strengths and weaknesses of various Bayesian Optimization strategies and their hyperparameters, we  
949 refer the reader to an open-source implementation of CMA-ES and its corresponding documentation  
950 <https://github.com/CMA-ES/pycma>. In this work we employ the default settings for CMA-ES  
951 reported in this open source implementation, with  $\sigma = 0.5$ .



## F Fidelity of Expert Model Oracle Functions

In this section we present additional studies on the quality of expert model oracle functions we use in our benchmark, as high fidelity oracles are used for evaluation in this benchmarks on tasks where real evaluation using the actual ground truth objective requires running real experiments. We first analyze the train-test discrepancy of the oracle, where we use the oracle’s predicted  $y$  values instead of the dataset values for the MBO algorithm. We measure the correlation of the performance ranks of all baseline algorithms in our benchmark, and we also measure of the change of rank of different algorithms. High correlation and low rank shift indicate that the use of oracle predicted values does not affect the relative performance order of different offline MBO algorithms, suggesting the benchmark is robust. We present the results in Table 5. For the majority of MBO tasks where an exact oracle is not available, the impact of the train-test discrepancy is minimal. For UTR, the task most sensitive to the train-test discrepancy, an option in the Design-Bench code-base allows for the training set to be relabelled with the predictions of the oracle model and the train-test discrepancy to be removed.

	GFP	UTR	Superconductor
Rank Correlation	0.895405	0.116667	0.666667
Max Rank Shift	2.000000	7.000000	4.000000
Avg Rank Shift	0.888889	2.666667	1.333333

Table 5: Spearman’s rank correlation coefficient of MBO algorithms when (1) using raw measurements in the training set and a learned oracle during evaluation versus (2) using predictions from a learned oracle during training and during evaluation.

In addition to the train-test discrepancy, we also measure the oracle agreement between different expert oracle models on the same dataset. To quantitatively analyze the agreement on different oracles, we use them to evaluate all offline MBO algorithms in this dataset, and measure the changes ranks of offline MBO algorithms between oracles. We present the results in Figure 4. The results suggest that high quality expert oracle models do agree with each other on the solutions found by offline MBO algorithms, indicating that the expert oracles in our benchmark is robust.

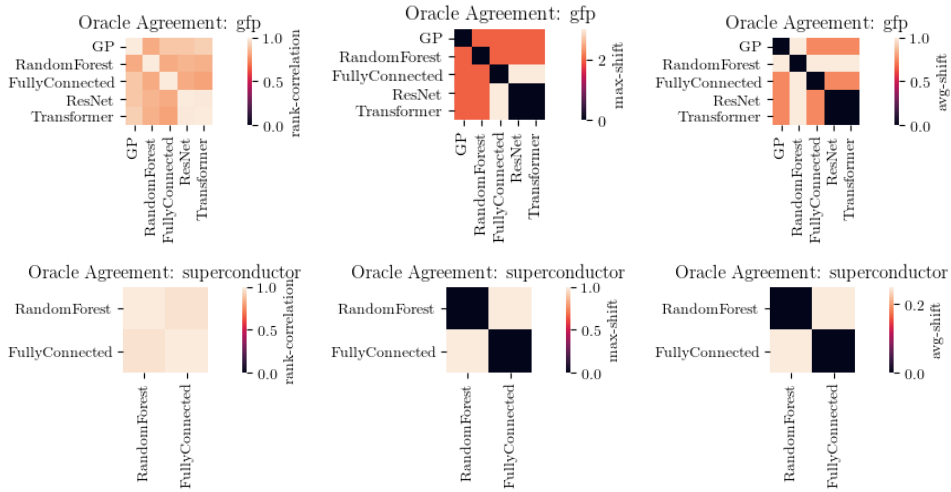


Figure 4: **Agreements of the ranks of offline MBO methods using different oracles on the same dataset.** We evaluate all offline MBO algorithms benchmarked in this paper on the same dataset using different oracles, and compare their performance ranking under different oracles. The rank correlation measures Spearman’s rank correlation coefficient between the order of the algorithms, and the mean and max rank shift shows the average and the max rank change of the same algorithm between different oracles. We can see that the rank correlations are high and rank shifts are low, suggesting that different oracles tend to agree on the relative order of the performance of offline MBO algorithms. This result indicates that our expert oracle models are robust.