586 APPENDIX

587 A Preprocessing and tokenization details

588 A.1 Search space primitives

Below are the exact descriptions of the hyperparameter primitives used to define a given $\mathcal{X}^{(d)}$.

- Double: Specifies a continuous range of possible values in the closed interval $[x_{\min}, x_{\max}]$ for some real values $x_{\min} \le x_{\max}$.
- Integer: Specifies an integer range of possible values in $[x_{\min}, x_{\max}] \in \mathbb{Z}$ for some integers $x_{\min} \leq x_{\max}$.
- Discrete: Specifies a finite, ordered set of values from \mathbb{R} .
- Categorical: Specifies an unordered list of strings.

596 A.2 Data preprocessing and tokenization

- ⁵⁹⁷ We list out the full set of preprocessing steps (from Section 4.1) below:
- Omit parameter and metric names in all trials, remove redundant keywords ("parameter", "trial", etc.), order trial parameters according to those in metadata m, and replace keywords and enumerations in metadata m with short strings and integers respectively.
 These replacements are:
 - "name" \rightarrow "N"
- 603 "metric" → "O"

602

604

605

609

610

611

612

614

- "goal" \rightarrow "G", with values "MAXIMIZE" \rightarrow 1, "MINIMIZE" \rightarrow 2
 - "type" \rightarrow "P", with values "DOUBLE" \rightarrow 1, "INTEGER" \rightarrow 2, "DISCRETE"
- 606 $\rightarrow 1$, "CATEGORICAL" $\rightarrow 4$
- 607 "algorithm" \rightarrow "A"
- 608 "min_value" \rightarrow "m", "max_value" \rightarrow "M"
 - "scale_type" \rightarrow "S", with values "LINEAR" \rightarrow 1, "LOG" \rightarrow 2
 - "categories" ightarrow "C"
 - Insert short separator symbols, e.g. * between parameter/metrics in a trial, "I" between trials, and "&" between experiment description and parameter configurations in metadata.
- Convert all values in history h to single integers.
 - Represent discrete and categorical parameters with their index in the set of values.
- 615 Normalize float and integer parameter values in $x_t^{(d)}$ with their value range and the 616 function values y_t with their minimum and maximum seen values in the entire study. 617 Then quantize the normalized values to an integer, e.g., "0.12345" \rightarrow "123" with a 618 quantization level of Q = 1000. More formally, we apply the following transformation 619 $q(\cdot)$:

$$q(z) = \operatorname{int}[z_{\operatorname{norm}} * Q], \text{ where } z_{\operatorname{norm}} = (z - z_{\min})/(z_{\max} - z_{\min})$$
(7)

The shortened text string is then converted to a sequence of tokens via the SentencePiece tokenizer [44] with a vocabulary of 33000 words. Quantized numbers in h are always converted into single tokens. As long as Q is sufficiently large, there is no concern from the loss of precision over numerical quantizations, and thus the serialized study contains nearly the same amount of information as the original data. For comparison, the naive tokenization for the example of Table 1 with t = 100 trials will produce 8221 tokens which can overload GPU memory, while our proposed tokenization will only produce 584 tokens, a 14x reduction.

627 **B** Algorithm and baseline details

628 **B.1 Dataset algorithms**

Grid Search: DOUBLE parameters are first log-transformed if specified. They are then converted into DISCRETE parameters by discretizing their ranges into 100 equidistant points. Suggestions are outputted using lexicographic ordering from the cartesian product over all parameters' feasible points. The traversal order follows the alphabetical ordering of parameter names. That is, given two parameters "Foo" and "Bar", both in [0,1] range, the sequence of trials looks like: {"Foo": 0, "Bar":0}, {"Foo": 0, "Bar":0.01}, ..., {"Foo": 0, "Bar":1}, {"Foo": 0.01, "Bar":0}, {"Foo": 0.01, "Bar":0.01},

Shuffled Grid Search: Shuffled grid search is the same as Grid Search in how it handles DOUBLE
 parameters. Instead of traversing the grid in a deterministic order, it selects without replacement a
 random point from the grid at each iteration.

Regularized Evolution [48]: In summary, this algorithm at every iteration randomly selects a tournament subset from the current population, and mutates the argmax member of the tournament. When inserting a new trial, the oldest trial will be removed. We use a population size of 25 and tournament size of 5. The mutation operation uniformly at random selects one of the parameters $x^{(r)}$ from x, and mutates $x^{(r)}$ based on the following: for DOUBLE, INTEGER, the new value is uniformly sampled from $\left[x_{\min}^{(r)}, x_{\max}^{(r)}\right]$, while for DISCRETE, CATEGORICAL, the new value is uniformly sampled from the feasible list.

Hill Climbing: This is a naive implementation, where at every iteration t, the current x_{pivot} is mutated (using the same operation as Regularized Evolution) to x_{mutated} , and evaluated. If $f(x_{\text{mutated}}) > f(x_{\text{pivot}})$, then we reassign x_{pivot} to be the mutated x_{mutated} . An extension of this method can be "batched", as seen in [60], although we not include this for the sake of clarity and presentation.

Eagle Strategy [49]: Eagle strategy is a metaheuristics algorithm that is a slight variation of Particle Swarm Optimization [61].

⁶⁵³ The algorithm is originally formulated for continuous search spaces only. The reason is that it involves

a subroutine (move step) where we take a convex combination of a particle (called *firefly* in [49]) and

another particle that has a better objective value. Mathematically, given two particle vectors x and x'and the coefficient $c \in [0, 1]$, the move step generates cx + (1 - c)x'.

The algorithm is extended to support DISCRETE and CATEGORICAL parameters by applying a separate move operation for each non-continuous dimension *d*:

 $\mathrm{move}(x^{(d)}, x'^{(d)}, c, \alpha) = \begin{cases} x^{(d)} & \text{with probability } (1 - \alpha)c \\ x'^{(d)} & \text{with probability } (1 - \alpha)(1 - c) \\ \mathrm{random \ value} & \text{with probability } \alpha \end{cases}$

where α is a small perturbation coefficient that decreases in the dimension of the search space.

Vizier [2]: Vizier's default algorithm is available via Google Cloud as Vertex Vizier. We have contacted the authors of the algorithm and received the the following details on its implementation.

In summary, the algorithm uses a particular implementation of GP-UCB with trust regions. The GP regressor model consists of the following:

- $\alpha \sim \text{TruncatedLogNormal controls the amplitude of Matern5/2 kernel.}$
 - $\lambda_i \sim \text{TruncatedLogNormal}$ (i.i.d. for each dimension *i*) controls the length scale for the *i*-th dimension.
- $\sigma \sim \text{TruncatedLogNormal controls the Gaussian noise}.$
- $z \sim \text{Normal}(0, \sigma)$ is the observation noise.
- $f \sim GP(\lambda, \alpha)$ is the function.

663

664

• $y(x) \sim f(x) + z$ is the noisy function. 668

where the prior distribution parameters are chosen depending on the user's estimate of the observation 669 noise. 670

The algorithm then uses gradient descent with line search for step sizes to obtain the MAP estimate 671 of α, λ and σ . Furthermore, the algorithm uses a variation of Eagle Strategy (explained above) to 672 optimize the UCB acquisition function with coefficient of 1.8. In order to prevent overexploration 673 that may result from the large UCB coefficient, the algorithm optimizes acquisition functions inside 674 trust region. The trust region is the union of L_{∞} -norm balls around explored points. The radius of 675 the L_{∞} -norm ball grows in the number of explored points. The algorithm also starts at the center of 676 the search space (unless user specifies an alternative initial batch). 677

GP-UCB: It is the same as Vizier's GP-UCB, except for the model definition. We used the model 678 definition from the github repository of the authors of "Heteroscedastic and Evolutionary Bayesian 679 Optimisation solver" (HEBO) [62], the winner of 2020 Blackbox Optimization challenge [63]. It is 680 worth noting that HEBO uses multi-dimensional acquisition functions derived from the GP model. 681 The priors over hyperparameters are thus not tuned to optimize the performance of GP-UCB algorithm, 682 which explains its suboptimal performance. 683

Gaussian Process for uncertainty estimation 684 B.2

We use the same GP model as GP-UCB. 685

When comparing the function prediction performance with the OPTFORMER, we choose $[y_{\min}, y_{\max}]$ 686 to normalize function value token based on the range of observed value in the sampled sequence 687

 $(x_1, y_1, \dots, x_t, y_t)$, and therefore the real value of y_t always resides in the prediction support of the 688 OPTFORMER. 689

To compensate for the fact that GP's distribution is wider than the real support used by the Transformer, 690 we truncate the GP's prediction into $[y_{\min}, y_{\max}]$ for a fairer comparison. 691

B.3 Transfer learning baselines 692

We use the following methods as transfer-learning baselines for the HPO-B dataset from Section 6.3: 693

ABLR [12, 51]: BO with multi-task adaptive Bayesian linear regression. Our implementation of 694 ABLR is equivalent to a GP with 0 mean and a dot-product kernel with learned basis functions. 695 We use a neural net (NN) with (128, 128) hidden layers and tanh activation as the basis functions. 696 We then train ABLR by optimizing the negative log likelihood (NLL) over NN weights θ as well 697 covariance matrix SS^{\top} and bias parameters δ^2 that define the dot-product kernel k, i.e. 698

$$k(x, x') = \phi_{\vartheta}(x)^{\top} S S^{\top} \phi_{\vartheta}(x') + \delta^{2}, \qquad (8)$$

where matrix $S \in \mathbb{R}^{128 \times 256}$, basis function ϕ_{θ} is parameterized by NN weights ϑ and $\delta \in \mathbb{R}$. 699

FSBO [7]: Bare-bone few-shot BO. We did not include data-driven initialization due to lack of 700 reproducing details. Following [7], our implementation of FSBO is equivalent to BO using a GP 701 with 0 mean and a squared-exponential kernel on top of a NN with (128, 128) hidden layers and tanh 702 activation functions. We train the NN weights as well as the parameters in the squared-exponential 703 kernel. 704

HyperBO [52, 53]: BO with pre-trained GPs. Following [53], we pre-train a GP with Matérn32 705 kernel on top of a NN with one hidden layer of width $2 \times D$ and tanh activation functions. Here D is 706 the input dimension of the search space. 707

For training, we use the Adam optimizer with learning rate 0.001 and batch size 50 for all the 708 transfer-learning baselines. Notice that these transfer-learning methods require "pre-training" a GP 709 on the same search space. We sample 10000 random data points on each HPO-B surrogate functions 710

from each search space. We train a separate GP for each search space. 711

712 C Data details

713 C.1 Dataset details

751

752

753

RealWorldData dataset: The RealWorldData dataset contains a total of 750K studies collected from users, and thus each study has a variable number of trials. Since some user studies can potentially have an excessive number of trials (e.g. 10K+), for all studies we only consider the first 300 trials for experiments. Since the dataset also includes usernames (due its origin from a proprietary database), we made sure to anonymize every study first.

We split the dataset in temporal order to avoid information leak, use most studies for training, and select 16 studies generated by a different set of users for testing. To bootstrap these studies into actual functions to be evaluated, we fit a GP for each study and output the function value as the GP's median

⁷²² function (due to the use of output warping).

HPO-B dataset: For HPO-B dataset, a tuning task is identified with a (search space id, dataset id) pair, which refers to tuning the hyperparameters defined in a search space for some machine learning model trained on a dataset. we use the "v3-augmented" meta-training/validation/test splits that includes all the 16 test search spaces as well as less frequent search spaces in the meta-training split. There are uniquely 1730, 91, and 86 tasks for training, validation and testing respectively. For every tuning task, [5] fits an XGBoost model to the trial data of every tuning task as the objective function.

Similar to the BBOB dataset, we generate 10M, 500K studies for training and validation respectively, along with the same set of controlled algorithms. For each of the test tuning task, we run 5 optimizations each with a different initial set of observations provided in [5].

The HPO-B uses the Apache 2.0 open-source license.

BBOB dataset: The BBOB dataset contains a total of 10M studies for training, each containing exactly 300 trials. An additional 500K studies (using different randomization seeds) are used for validation. While the number of studies can be freely generated and effectively unlimited, we found that 10M studies were sufficient for the Transformer to train properly.

The functions we use for data are from [47], and consist of separable functions (Sphere,
Ellipsoid Separable, Rastrigin Separable, Bueche Rastrigin,

{Linear Slope}), moderately conditioned, potentially multi-modal functions (Attractive 739 Sector, Step Elllipsoid, {Rosenbrock Rotated}), ill-conditioned functions 740 (Discus, Bent Cigar, Sharp Ridge, {Sum of Powers}), multi-modal functions 741 (Weierstrass, Schaffers F7, Schaffers F7 Illconditioned, {Greiwank 742 Rosenbrock}), and functions with weak global structures (Schwefel, Gallagher 21, 743 Gallagher 101, Katsuura, {Lunacek}). The functions noted with the extra "{}" are for 744 testing and excluded from the training data. We apply significant randomization over the functions 745 for both the training dataset and test-time evaluation. In order, we randomize the following: 746

- Function dimension *D*, which is uniformly selected from a range. For training data generation, this range is [1, 20].
- Orthonormal rotation matrix Γ , which is applied to the input first, i.e. producing a new function $f'(\boldsymbol{x}) = f(\Gamma \boldsymbol{x})$.
 - Shift vector \boldsymbol{x}_{shift} which is also applied to the input first, i.e. producing a new function $f'(x) = f(x \boldsymbol{x}_{shift})$, where \boldsymbol{x}_{shift} has all of its coordinate-wise entries sampled from [-4, 4], while the domain is [-5, 5].
- Discretization, in which the parameter space $\mathcal{X}^{(d)}$ is uniformly at random chosen to be either a DOUBLE, DISCRETE, CATEGORICAL parameter. The DOUBLE parameter "discretization" is actually a no-op, as it allows the original continuous space $\mathcal{X}^{(d)} \subset \mathbb{R}$. Otherwise, a number L of feasible points is uniformly selected from the range [2, 8], and used to divide the original [-5, 5] range into L equally-spaced points. If DISCRETE was chosen, then the ordering of the grid points is preserved, otherwise if CATEGORICAL was chosen, then all of the gridpoints become effectively unordered strings.
- Noise Type, in which one of 10 noise settings (including no noise) is uniformly chosen.
 Noise consists of either Gaussian (multiplier sampled from a random Gaussian of varying

scale is applied), Uniform (multiplier sampled from uniform distribution of varying scale is
 applied), or Cauchy (additive noise which only occurs at a probabilistic frequency, with a
 varying fixed strength is applied).

For evaluation, we randomly sample 100 configurations for each of the five test functions, resulting in 500 optimization trajectories in total.

For BBOB, as all parameters are named as "x_i" with $i \in [0, D)$ and always have value range in [-5, 5], significantly different from the other two datasets, we omit their parameter names and value in the metadata m and only keep parameter type information.

771 **D** Model and training details

The open-sourced T5 model codebase we use can be found at https://github.com/ google-research/t5x.

774 D.1 Conditional probability decomposition

From Section 4.2, the joint distribution of the optimization history h conditioned on metadata m can be written using the chain rule as

$$P(\bar{\boldsymbol{h}}|\bar{m}) = P\left(\bar{x}_{1}^{(1)}, \bar{x}_{1}^{(2)}, \dots, \bar{x}_{1}^{(D)}, \star, \bar{y}_{1}, "|", \dots, \bar{x}_{T}^{(1)}, \bar{x}_{T}^{(2)}, \dots, \bar{x}_{T}^{(D)}, \star, \bar{y}_{T}|\bar{m}\right)$$

$$= \prod_{t=1}^{T} \left(\prod_{d=1}^{D} P\left(\bar{x}_{t}^{(d)}|\bar{m}, \bar{\boldsymbol{h}}_{t-1}, \bar{\boldsymbol{x}}_{t}^{(1:d-1)}\right)\right) P\left(\star|\bar{m}, \bar{\boldsymbol{h}}_{t-1}, \bar{\boldsymbol{x}}_{t}\right) P\left(\bar{y}_{t}|\bar{m}, \bar{\boldsymbol{h}}_{t-1}, \bar{\boldsymbol{x}}_{t}\right) P\left("|"|\bar{m}, \bar{\boldsymbol{h}}_{t}\right)$$
(9)

We note that this correctly formalizes the prediction of objects we are most interested in, which are parameter values $P\left(\bar{x}_{t}^{(d)}|\bar{m}, \bar{h}_{t-1}, \bar{x}_{t}^{(1:d-1)}\right)$ and function values $P\left(\bar{y}_{t}|\bar{m}, \bar{h}_{t-1}, \bar{x}_{t}\right)$.

779 D.2 Training

During training, the encoder (denoted as \mathbf{E}_{θ}) input sequence length is selected to be the maximum length of the tokenized metadata \bar{m} from a dataset, ranging from 256 to 1024. The decoder (denoted as \mathbf{D}_{θ}) input sequence is fixed at 1024, which means it can model up to 1024//(D+3) trials where *D* is the number of parameters. We use Adam optimizer with a rsqrt learning rate schedule and a mini-batch size of 256, and train each model up to 1M steps, with early stopping according to the validation loss. Each model is trained with a 4x4 TPU-v3 slice.

Thus the prediction for $\bar{h}^{(n)}$ is:

$$P_{\theta}\left(\bar{h}^{(n)}\Big|m, \bar{\boldsymbol{h}}^{(1:n-1)}\right) = \operatorname{SoftMax}\left[\mathbf{D}_{\theta}(\mathbf{E}_{\theta}(\bar{m}), \bar{\boldsymbol{h}}^{(1:n-1)})\right]$$
(10)

787 D.3 Data augmentation

788 We adopt the following three data augmentations to reduce overfitting to the offline datasets:

- 1. In order for the model to be invariant to parameter ordering, we apply random parameter permutations over metadata \bar{m} and every suggestion \bar{x}_t .
- 791 2. In order for the model to be robust to a different normalization range given a new function, 792 we apply random scaling and shifting to the normalized function value $y_{\text{norm}} = (y - y_{\text{min}})/(y_{\text{max}} - y_{\text{min}})$ before quantization:

$$y'_{\text{norm}} = y_{\text{norm}} * s + c, \ s \sim \text{Uniform}[0.3, 1], \ c \sim \text{Uniform}[0, 1 - s]$$
 (11)

and thus $y'_{norm} \in [c, c+s] \subseteq [0, 1]$ after transformation.

⁷⁹⁵ 3. Randomly drop textual and parameter value range information in metadata.

Table 5: Example of studies in RealWorldData (left), BBOB (middle) and HPO-B (right).

"name": "gan1d 500 iters -"2022-05-18" "parameter": { "name": "learning_rate", "min_value": 1e-06,
"max_value": 0.01, "type": "DOUBLE", "scale_type": "LOG", "parameter": "name": "modifier", "min_value": 0.1,
"max_value": 1000000.0, "type": "DOUBLE", "scale_type": "LOG", "parameter": { "name": "weight_init_std", "min_value": 0.01, "max_value": 2.0, "type": "DOUBLE", "parameter": { "name": "optimizer", "type": "CATEGORICAL", "categories": "sgd", "categories": "adam", "categories": "rmsprop", "goal": "MINIMIZE", "max_num_trials": 500, "metric": "", "observation_noise": "HIGH", "trial"• { "parameter": { "learning_rate": 0.0001, "modifier": 316.2277660168381, "optimizer": "sgd", "weight_init_std": 1.005, "metric": { "": -0.946908021738347, "trial": { "parameter": { "learning_rate": 0.000504, "modifier": 12.346786652749216, "optimizer": "rmsprop", "weight_init_std": 1.2192566347109868, "metric": { "": -1.5144472008077585, }

"name": "SCHAFFERS_F7", "algorithm": "gp", "parameter": { "parameter": { "name": "x1", "min_value": -5.0, "max_value": 5.0, "type": DOUBLE, "scale_type": UNIT_LINEAR_SCALE, 1. "parameter": { "name": "x2", "min_value": -5.0,
"max_value": 5.0, "type": DOUBLE, "scale_type": UNIT_LINEAR_SCALE, "parameter": { "name": "x3", "type": DISCRETE, "values": [-5.0, 5.0], "parameter": { "name": "x4", "type": CATEGORICAL, "categories": ["5.0", "-1.666666666667", "-5.0", "1.666666666667"], "parameter": { "name": "x5", "min_value": -5.0, "max_value": 5.0,
"type": DOUBLE, "scale_type": UNIT_LINEAR_SCALE, "metric": "", "goal": MAXIMIZE, "observation_noise": HIGH "trial": { "parameter": { "x0": "0.0", "x1": 0.0, "x2": 0.0, "x3": 5.0, "x4": "-5.0", "x5": 0.0, "metric": { "": -334.4782223514127, "trial": { "parameter": { "x0": "5.0", "x1": -1.9867479768748013, "x2": -1.7665621302793095, "x3": -5.0, "x4": "1.6666666666666667", "x5": -1.7634306558106605, "metric": { "": -323.84900527589326,

"name": "5859_145853", "algorithm": "GP UCB", "parameter": { "name": "minsplit", "max_value": 60.0, "type": "DOUBLE", "scale_type": "LINEAR", "parameter": {, "name": "minsplit.na", "max_value": 1.0, "type": "DOUBLE", "parameter": { "name": "minbucket", "min_value": 1.0, "max_value": 60.0, "type": "DOUBLE", "scale_type": "LINEAR", 'parameter": { "name": "cp", "min_value": 0.000100788830221, "max_value": 1.000092678873241, "type": "DOUBLE", "scale_type": "LOG", "parameter": { "name": "maxdepth", "max_value": 29.0, "type": "DOUBLE", "scale_type": "LINEAR", "parameter": { "name": "maxdepth.na", "max_value": 1.0,
"type": "DOUBLE", "observation_noise": AUTOMATIC, "metric": "objective_value", "goal": "MAXIMIZE" "trial": { "parameter": { "minsplit": 4.0, "minsplit.na": 0.0, "minbucket": 18.0, "cp": 0.7342895964927976, "maxdepth": 3.0, "maxdepth.na": 0.0, "metric": { "objective_value": 0.500024080276, "trial": { "parameter": { "minsplit": 8.0, "minsplit.na": 0.0, "minbucket": 32.0, "cp": 0.30972302652187583, "maxdepth": 4.0, "maxdepth.na": 0.0, "metric": { "objective_value": 0.50002408028, . . .

ł

796 D.4 Inference

At inference time, we choose the decoder input sequence length according to the maximum number of trials to run. E.g. to optimize a function with 18 parameters (highest possible dimension D over our test functions) over 105 trials, we set the input sequence length to be at least (18 + 3) * 105 = 2205.

We compute the (y_{\min}, y_{\max}) range for function value normalization in the tokenization process with the current minimal and maximum observations. We set c = 0.2, s = 0.6 so that all normalized observations fall in the range of $y'_{norm} \in [0.2, 0.8]$, and the model's y value predicted distribution support, [0, 1], is sufficiently large.

We also use a softmax temperature hyperparameter when predicting function values. We choose the

temperature to maximize the log-likelihood of the validation split of each dataset seperately. On

- RealWorldData, the function prediction temperature is set as 1.1 and on HPO-B it is 1.5. The policy
- ⁸⁰⁷ prediction temperature is always set to be 1.

808 E Additional experimental results

We provide additional experimental results in this section. The most important result to note is when we compare different acquisition function choices in Appendix E.1 and observe that the **augmented policy with EI outperforms all baselines** on both RealWorldData and HPO-B benchmarks across all metrics.

E.1 Ablation on acquisition functions: EI augmented policy outperforms all baselines

We provide additional ablations on acquisition function choices on both the RealWorldData and HPO-B datasets, and show that the OPTFORMER augmented with Expected Improvement (EI) produces the best performance across the board, over all other OPTFORMER variants and baselines.

In Fig. 6, we compare the Thompson Sampling (TS) used in the main body with Probability of Improvement (PI), Expected Improvement (EI) and Upper Confidence Bound (UCB) with a confidence level of 0.9. We also include the best performing standalone baseline, Vizier, and transfer learning baseline, HyperBO, for reference. We observe that the prior policy is improved by all the acquisition functions. Particularly, OPTFORMER (EI) is the best among all acquisition functions and clearly outperforms all the baseline methods (HyperBO and Vizier) on both datasets across all trial steps.



Figure 6: Ablation on the choice of acquisition functions. The plot shows the best normalized function values averaged over HPO-B test functions. Ablation curves are shown with \bigcirc markers.

To further bolster this hypothesis, we also compare using other metrics such as performance profiles 823 [64], which are widely used performance evaluation tools for comparing optimization methods. In our 824 case, the y-axis is the fraction of tasks that each method succeeds in, for a given trial index (x-axis). 825 The specific criteria of success depends on the problem itself, and we present natural performance 826 profiles based on the following criterion: outperforming 90% of the best function value obtained 827 by all methods at the 50th iteration. As this metric depends on the set of methods being compared, 828 we include all baselines from the main body. As we can see, Fig. 7 demonstrates that augmented 829 OPTFORMER policies, especially OPTFORMER (EI), produce superior performance compared to 830 other baselines. 831



Figure 7: Ablation on the choice of acquisition functions. The plot shows the performance profile metric with success threshold: median best function value at 50th iteration.

832 E.2 Imitating HPO policies



Figure 8: Policy distribution $p(x_{40}^{(d)}|m, h_{39}, x_{40}^{(1:d-1)})$ for d = 1, 2 on a 2D GRIEWANK ROSEN-BROCK function.



Figure 9: Best normalized function value with std, averaged over 5 test functions each with 100 runs.



Figure 10: Best normalized function value of LINEAR SLOPE with std, averaged over 100 runs.



Figure 11: Best normalized function value of ROSENBROCK ROTATED with std, averaged over 100 runs.



Figure 12: Best normalized function value of SUM OF POWERS with std, averaged over 100 runs.



Figure 13: Best normalized function value of GRIEWANK ROSENBROCK with std, averaged over 100 runs.



Figure 14: Best normalized function value for LUNACEK with std, averaged over 100 runs.

833 E.3 Learning priors for hyperparameter response functions

We apply the same goodness-of-fit analysis on function prediction from Section 6.2 to the test split of HPO-B. The results are shown in Fig. 15.



Figure 15: Fitness of predicted CDF(y) on HPO-B test set.

The ECE metric is defined for a classification model. To obtain a similar measurement for a continuous

regression model, we convert the continuous regression problem into a multi-class classification problem by discretizing the range $[y_{\min}, y_{\max}]$ for each study into 100 equal intervals. Then, we

⁸³⁹ follow the definition of ECE in [27] and estimate the metric using 10 confidence bins.

840 E.4 Augmenting a prior policy with function prediction

Transfer learning results on HPO-B Fig. 4 showed the best normalized function values observed
so far at each trial. Though HyperBO uses a smaller NN for feature extraction, HyperBO has a
flexible mean function, which captures important information that benefits BO in beginning trials.
While we implemented a bare-bone FSBO, its performance is still better than ABLR in part thanks to
FSBO's use of a squared exponential kernel instead of a dot-product one. Comparing to a dot-product
kernel with a finite feature space, a squared exponential kernel introduces infinite features.

In Fig. 16 and Fig. 17, we show the performance profiles of all compared methods over 2 different metrics: outperforming 90% of the best function value obtained by all methods at the 50th iteration, and outperforming the median of the best function values obtained by each method at the 50th iteration.

Despite the relatively better performance of HyperBO, FSBO and ABLR especially during earlier 851 trials, shown by Fig. 4, these methods do not achieve a high percentage success rate on the 86 HPO-B 852 test functions as reflected by Fig. 17. As pointed out by Wang et al. [53], ABLR, FSBO can be viewed 853 as special cases of HyperBO with specific settings of kernel and mean functions. These methods 854 have guarantees only if each function (corresponding to each task) is an i.i.d. sample from the same 855 GP. However, for some search spaces in HPO-B, there exist surrogate functions that return constant 856 values. The constant surrogate function is unlikely to be an i.i.d. sample from the same GP as other 857 surrogates in the same search space. This means, ABLR, FSBO and HyperBO can be sensitive to 858 how the data is generated and outliers in the training data. 859

Summarizing the results in Fig. 4, Fig. 16 and Fig. 17, HyperBO was able to achieve very good overall performance on a subset of all search spaces, which led to a better averaged best normalized function values. It is likely that these search spaces have surrogate functions that meet the i.i.d function sample assumption from Wang et al. [53]. However, if we only look at the fraction of tasks each method surpasses a success metric, HyperBO may not be a method with superior performance that is comparable to the OPTFORMER. This reveals another benefit of the OPTFORMER: robustness to function outliers.



Figure 16: Performance profile on RealWorldData and HPO-B test functions with success threshold: 90% best function value at 50th iteration.



Figure 17: Performance profile on RealWorldData and HPO-B test functions with success threshold: median best function value at 50th iteration.

HPO-B plotting We further compare the augmented policies from Section 6.3 to the provided baselines for HPO-B in [5], using the same plotting format from [5] for fair comparison.



Figure 18: (Lower is better) Aggregated comparisons of normalized regret and mean ranks across all search spaces on the continuous search spaces of HPO-B-v3.