
FSEO: A Few-Shot Evolutionary Optimization Framework for Expensive Multi-Objective Optimization and Constrained Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Meta-learning has been demonstrated to be useful to improve the sampling effi-
2 ciency of Bayesian optimization (BO) and surrogate-assisted evolutionary algo-
3 rithms (SAEAs) when solving expensive optimization problems (EOPs). However,
4 existing studies focuses on only single-objective optimization, leaving other ex-
5 pensive optimization scenarios unconsidered. We propose a generalized few-shot
6 evolutionary optimization (FSEO) framework and focus on its performance on two
7 common expensive optimization scenarios: multi-objective EOPs (EMOPs) and
8 constrained EOPs (ECOPs). We develop a novel meta-learning modeling approach
9 to train surrogates for our FSEO framework, an accuracy-based update strategy is
10 designed to adapt surrogates during the optimization process. The surrogates in
11 FSEO framework combines neural network with Gaussian Processes (GPs), their
12 network parameters and some parameters of GPs represent useful experience and
13 are meta-learned across related optimization tasks, the remaining GPs parameters
14 are task-specific parameters that represent unique features of the target task. We
15 demonstrate that our FSEO framework is able to improve sampling efficiency on
16 both EMOP and ECOP. Empirical conclusions are made to guide the application of
17 our FSEO framework.

18 1 Introduction

19 Expensive optimization problems (EOPs) aim to find as good as possible solutions within a budget
20 of limited solution evaluations. Conventional Bayesian optimization (BO) and surrogate-assisted
21 evolutionary algorithms (SAEAs) have been widely used to solve EOPs, but they train surrogate
22 models from the scratch. To further improve the sampling efficiency and optimization performance,
23 many efforts have been made to pre-train surrogates with the prior experience gain from related
24 optimization tasks, resulting in experience-based optimization algorithms [1, 21, 36, 35].

25 This work considers solving EOPs on the context of few-shot problems [5, 41], where plenty of
26 expensive related tasks are available and each of them can provide a small dataset for experience
27 learning. Therefore, many experience-based optimization approaches such as multi-tasking optimiza-
28 tion [43, 2, 47], transfer optimization [35, 17, 16] are not considered as they cannot learn experience
29 from small related tasks (A discussion is available in Appendix A). In comparison, meta-learning
30 [14] has been proved to be powerful in solving few-shot problems, leading to a new subcategory of
31 experience-based optimization, namely few-shot optimization (FSO) [46].

32 Existing studies on FSO are mainly few-shot Bayesian optimization (FSBO) where meta-learning
33 approaches are combined with BO to solve EOPs with only one objective. In this paper, we propose
34 a generalized few-shot evolutionary optimization (FSEO) framework to address EOPs from the

35 perspective of SAEAs and consider two expensive optimization scenarios which have been limited
36 studied: multi-objective EOPs (EMOPs) and constrained EOPs (ECOPs). Major contributions are
37 summarized as follows.

- 38 • A novel meta-learning method, namely Meta Deep Kernel Learning (MDKL), is developed
39 to gain prior experience from related expensive tasks. Our model architecture and parameter
40 designs make it possible to generate a regression-based surrogate on the prior experience
41 and then continually adapt the surrogate to approximate the target task.
- 42 • We propose a FSEO framework to solve EOPs from the perspective of SAEAs. FSEO
43 framework is applicable to regression-based SAEAs since FSEO embed our meta-learning
44 models in these SAEAs as their surrogates. In addition, an update strategy is designed to
45 adapt surrogates constantly during the optimization. Note that our FSEO framework is a
46 general framework but we focus on its performance on EMOPs and ECOPs in this paper.
- 47 • Experiments are conducted on EMOPs and ECOPs to show our FSEO framework is effective.
48 Our comprehensive ablation studies discover the influence of some factors on FSEO
49 performance and provide empirical guidance to the application of FSEO framework.

50 2 Related Work

51 Experience-based optimization can be divided into several subcategories according to the techniques
52 of learning prior experience from related tasks. A detailed classification and discussion on these
53 subcategories is available in Appendix A. This subsection focuses on related work on FSO.

54 FSO studies in the literature can be classified based on their model architectures. Most studies meta-
55 learn parameters for Gaussian Processes (GPs) [44], namely FSBO or Meta Bayesian Optimization
56 (MBO) [32, 42, 26, 38]. In addition, [23] meta-learns with transformer neural processes and [46, 6]
57 meta-learn parameters for the architecture of deep kernel learning (DKL) [45]. The MDKL model in
58 our FSEO belongs to the last category as its model architecture is relevant to DKL.

59 Our work is different from existing studies in three points: Firstly, many studies [46] use existing
60 meta-learning models [27] as their surrogates. No further adaptations are made to these surrogates
61 during optimization since they are not originally designed for optimization. In comparison, we try to
62 develop a meta-learning model, MDKL, for optimization purpose. MDKL has explicit task-specific
63 parameters, which allows continually model adaptations during the optimization. Secondly, existing
64 work investigated only global optimization, leaving other optimization scenarios such as EMOP and
65 ECOP still awaiting for investigation. As our MDKL is designed for optimization and is capable of
66 continually adaptation, we pay attention on EMOPs and ECOPs which require more effective models
67 than global optimization. Our work widens the scope of existing FSO research and it focuses on the
68 perspective of SAEAs instead of BO. Lastly, in-depth ablation studies are lacking in the literature,
69 making it unclear which factors affect the performance of FSO. Our extensive ablation studies fill
70 this gap and we conclude some empirical rules to improve the performance of FSO.

71 3 Background

72 This section gives preliminaries about meta-learning and DKL. The former is the method of experience
73 learning, the latter is the underlying structure of experience representation.

74 3.1 Meta-Learning in Few-Shot Problems

75 In the context of few-shot problems, we have plenty of related tasks, each task T contributes a couple
76 of small datasets $D = \{(S, Q)\}$, namely support dataset S and query dataset Q , respectively. After
77 learning from datasets of random related tasks, a support set S_* from new unseen task T_* is given
78 and one is asked to estimate the labels or values of a query set Q_* . The problem is called 1-shot or
79 5-shot when only 1 data point or 5 data points are provided in S_* . A comprehensive definition of
80 few-shot problems is available in [5, 41].

81 Meta-learning methods have been widely used to solve few-shot problems [41]. They learn domain-
82 specific features that are shared among related tasks as experience, such experience is used to
83 understand and interpret the data collected from new tasks encountered in the future.

84 3.2 Deep Kernel Learning (DKL)

85 DKL aims at constructing kernels that encapsulate the expressive power of deep architectures for GPs.
 86 To create expressive and scalable closed form covariance kernels, DKL combines the non-parametric
 87 flexibility of kernel methods and the structural properties of deep neural networks. In practice, a deep
 88 kernel $k(\mathbf{x}^i, \mathbf{x}^j | \gamma)$ transforms the inputs \mathbf{x} of a base kernel $k(\mathbf{x}^i, \mathbf{x}^j | \theta)$ through a non-linear mapping
 89 given by a deep architecture $\phi(\mathbf{x} | \mathbf{w}, \mathbf{b})$:

$$k(\mathbf{x}^i, \mathbf{x}^j | \gamma) = k(\phi(\mathbf{x}^i | \mathbf{w}, \mathbf{b}), \phi(\mathbf{x}^j | \mathbf{w}, \mathbf{b}) | \theta), \quad (1)$$

90 where θ and (\mathbf{w}, \mathbf{b}) are parameter vectors of the base kernel and the deep architecture, respectively.
 91 $\gamma = \{\theta, \mathbf{w}, \mathbf{b}\}$ is the set of all parameters in this deep kernel. Note that in DKL, all parameters γ of a
 92 deep kernel $k(\mathbf{x}^i, \mathbf{x}^j | \gamma)$ are learned jointly by using the log marginal likelihood function of GPs as a
 93 loss function. Such a jointly learning strategy has been shown to make a DKL algorithm outperform
 94 a combination of a deep neural network and a GP model, where a trained GP model is applied to the
 95 output layer of a trained deep neural network [45].

96 3.3 Meta-Learning on DKL

97 An important distinction between DKL algorithms and the applications of meta-learning to DKL is
 98 that DKL algorithms learn their deep kernels from single tasks instead of collections of related tasks.
 99 Such a difference alleviates two drawbacks of single task DKL [39]: First, the scalability of deep
 100 kernels is no longer an issue as datasets in meta-learning are small. Second, the risk of overfitting is
 101 decreased since diverse data points are sampled across tasks.

102 4 Few-Shot Evolutionary Optimization (FSEO) Framework

103 In this paper, T_* denotes the target optimization task, and plenty of small datasets D_i sampled from
 104 related tasks T_i are available for experience learning. A complete list of notations is available at the
 105 beginning of Appendix.

106 4.1 Overall Working Mechanism

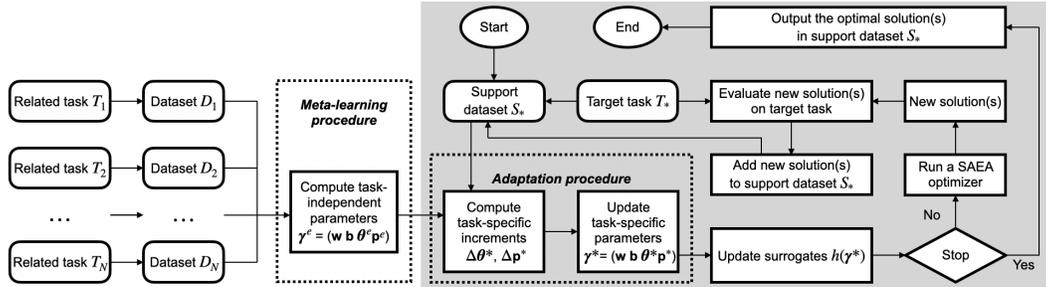


Figure 1: Diagram of our FSEO framework.

107 As illustrated in Fig. 1, all modules covering the optimization of target task T_* are included in a
 108 grey block. The modules beyond the grey block are associated with related tasks T_i and experience
 109 learning, which distinguishes our FSEO framework from conventional SAEAs and BO. The MDKL
 110 surrogate modeling method consists of two procedures: meta-learning procedure and adaptation
 111 procedure. The former learns prior experience from T_i , the latter uses experience to adapt surrogates
 112 to fit T_* . The framework of FSEO is depicted in Alg. 1, it consists of the following major steps.

- 113 1. **Experience learning:** Before expensive optimization starts, a meta-learning procedure is
 114 conducted to train task-independent parameters γ^e for MDKL surrogates (line 2). N_m
 115 datasets $\{D_{m1}, \dots, D_{mN_m}\}$ collected from N related tasks $\{T_1, \dots, T_N\}$ are used to train
 116 γ^e . γ^e is the experience that represents the domain-specific features of related tasks.
- 117 2. **Initialize surrogates with experience:** Optimization starts when a target optimization task
 118 T_* is given. An initial dataset S_* is sampled (line 3) to adapt task-specific parameters γ^* on
 119 the basis of experience γ^e . After that, MDKL surrogates are updated (line 4).

- 120 3. **Reproduction:** MDKL surrogates $h(\gamma^*)$ are combined with a SAEA optimizer Opt to
 121 search for optimal solution(s) \mathbf{x}^* on $h(\gamma^*)$ (line 7). This is implemented by replacing the
 122 original (regression-based) surrogates in a SAEA with $h(\gamma^*)$.
- 123 4. **Update archive and surrogates:** New optimal solution(s) \mathbf{x}^* is evaluated on target task T_*
 124 (line 8). The evaluated solutions will be added to dataset S_* (line 9) which serves as an
 125 archive. Then, surrogate adaptation is triggered, surrogates $h(\gamma^*)$ are updated (line 10).
- 126 5. **Stop criterion:** Once the evaluation budget has run out, the evolutionary optimization will
 127 be terminated and the optimal solutions in dataset S_* will be outputted. Otherwise, the
 128 algorithm goes back to step 3.

Algorithm 1 FSEO Framework.

- 1: **Input:** D_i : Datasets collected from related tasks $T_i, i=\{1, \dots, N\}$; N_m : Number of subsets D_m
 for meta-learning; $|D_m|$: Size of subsets D_m . $|D_m| \leq |D_i|$ due to $D_m \subseteq D_i$; Batch size B ;
 Surrogate learning rates α, β ; Target task T_* ; A SAEA optimizer Opt ; Fitness evaluation budget
 FE_{max} .
- 2: Experience $\gamma^e \leftarrow \text{Meta-learning}(D_i, N_m, |D_m|, B, \alpha)$. /* Alg. 2. */
- 3: $S_* \leftarrow \text{Sampling } 1d \text{ solutions from } T_*$.
- 4: $h(\gamma^*) \leftarrow \text{Adaptation}(\gamma^e, S_*, \beta)$. /* Initialize surrogate. */
- 5: Set evaluation counter $FE = |S_*|$.
- 6: **while** $FE < FE_{max}$ **do**
- 7: Candidate solution(s) $\mathbf{x}^* \leftarrow \text{Surrogate-assisted optimization}(Opt, h(\gamma^*))$.
- 8: $f(\mathbf{x}^*) \leftarrow \text{Evaluate } \mathbf{x}^* \text{ on } T_*$.
- 9: $S_* \leftarrow S_* \cup \{(\mathbf{x}^*, f(\mathbf{x}^*))\}$.
- 10: $h(\gamma^*) \leftarrow \text{Update}(\gamma^*, S_*, \beta)$. /* Alg. 4. */
- 11: Update FE .
- 12: **end while**
- 13: **Output:** Optimal solutions in S_* .
-

129 **4.2 Learning and Using Experience by MDKL**

130 In MDKL, the domain-specific features of related tasks are used as experience, which are represented
 131 by the task-independent parameters γ^e learned across related tasks. To make MDKL more capable of
 132 expressing complex domain-specific features, the base kernel $k(\mathbf{x}^i, \mathbf{x}^j | \boldsymbol{\theta})$ in GP is combined with a
 133 neural network $\phi(\mathbf{w}, \mathbf{b})$ to construct a deep kernel (see Eq.(1)). The modeling of a MDKL model
 134 consists of two procedures: meta-learning procedure and adaptation procedure. To make a clear
 135 illustration, we introduce frameworks of two procedures and then explain them in detail.

136 **Meta-learning procedure: Learning experience**

137 Our MDKL model uses the kernel in [18] as its base kernel:

$$k(\mathbf{x}^i, \mathbf{x}^j | \boldsymbol{\theta}, \mathbf{p}) = \exp\left(-\sum_{k=1}^d \theta_k |x_k^i - x_k^j|^{p_k}\right). \quad (2)$$

138 Therefore, the deep kernel will be:

$$k(\mathbf{x}^i, \mathbf{x}^j | \boldsymbol{\gamma}) = \exp\left(-\sum_{k=1}^d \theta_k |\phi(x_k^i) - \phi(x_k^j)|^{p_k}\right), \quad (3)$$

139 where $\boldsymbol{\gamma} = \{\mathbf{w}, \mathbf{b}, \boldsymbol{\theta}, \mathbf{p}\}$ is a set of deep kernel parameters. ϕ, \mathbf{w} and \mathbf{b} are neural network and its
 140 parameters (see Eq.(1)). Details about alternative base kernels are available in [44].

141 The aim of meta-learning procedure is to learn experience γ^e from related tasks $\{T_1, \dots, T_N\}$,
 142 including neural network parameters \mathbf{w}, \mathbf{b} , and task-independent base kernel parameters $\boldsymbol{\theta}^e, \mathbf{p}^e$. The
 143 pseudo-code of meta-learning procedure is given in Alg. 2. Ideally, the experience γ^e is learned from
 144 plenty of (N_m) small datasets D_m collected from different related tasks. However, in practice, the
 145 number of available related tasks N may be much smaller than N_m . Hence, the meta-learning is
 146 conducted gradually over U update iterations (line 3). During each update iteration, a small batch of

related tasks contribute
 B small datasets
 $\{D_{m1}, \dots, D_{mB}\}$ for
 meta-learning purpose
 (lines 5 and 7). Note that
 if $N < N_m$, a related task
 T_i can be used multiple
 times in the meta-learning
 procedure.
 For a given dataset D_{mi} ,
 we denote $\theta^i = \theta^e + \Delta\theta^i$
 and $\mathbf{p}^i = \mathbf{p}^e + \Delta\mathbf{p}^i$ as
 the task-specific kernel pa-
 rameters, where $\Delta\theta^i, \Delta\mathbf{p}^i$
 are the distance we need
 to move from the task-
 independent parameters to
 the task-specific parame-
 ters (line 9). The loss func-
 tion L of MDKL is the like-
 lihood function defined as
 follows [18]:

Algorithm 2 Meta-learning($D_i, N_m, |D_m|, B, \alpha$)

1: **Input:** D_i : Datasets collected from related tasks $T_i, i=\{1, \dots, N\}$; N_m :
 Number of subsets D_m for meta-learning; $|D_m|$: Size of subsets D_m .
 $|D_m| \leq |D_i|$ due to $D_m \subseteq D_i$; Batch size B ; Learning rate for priors α .
 2: Randomly initialize $\mathbf{w}, \mathbf{b}, \theta^e, \mathbf{p}^e$.
 3: Set the number of update iterations $U = N_m/B$.
 4: **for** $j = 1$ to U **do**
 5: $\{D'_1, \dots, D'_B\} \leftarrow$ Randomly select a batch of datasets from
 $\{D_1, \dots, D_N\}$.
 6: **for all** D'_i in the batch **do**
 7: $D_{mi} \leftarrow$ A subset of size $|D_m|$ from D'_i .
 8: Initialize task-specific increment $\Delta\theta^i, \Delta\mathbf{p}^i$.
 9: Compute task-specific parameters: $\theta^i = \theta^e + \Delta\theta^i, \mathbf{p}^i = \mathbf{p}^e + \Delta\mathbf{p}^i$.
 10: Obtain deep kernel $k(\mathbf{x}^i, \mathbf{x}^j | \gamma)$ based GP: $h(\gamma)$, where $\gamma =$
 $\{\mathbf{w}, \mathbf{b}, \theta^i, \mathbf{p}^i\}$ (Eq.(3)).
 11: Compute the loss function $L(D_{mi}, h(\gamma))$ (Eq.(4)).
 12: **end for**
 13: Update $\mathbf{w}, \mathbf{b}, \theta^e, \mathbf{p}^e$ via gradient descent: $\alpha \nabla L(D_{mi}, h(\gamma))$ (Eq.(5)).
 14: **end for**
 15: **Output:** Task-independent parameters: $\gamma^e = \{\mathbf{w}, \mathbf{b}, \theta^e, \mathbf{p}^e\}$.

$$\frac{1}{(2\pi)^{n/2}(\sigma^2)^{n/2}|\mathbf{R}|^{1/2}} \exp\left[-\frac{(\mathbf{y} - \mathbf{1}\mu)^T \mathbf{R}^{-1}(\mathbf{y} - \mathbf{1}\mu)}{2\sigma^2}\right], \quad (4)$$

148 where $|\mathbf{R}|$ is the determinant of the correlation matrix \mathbf{R} , each element in the matrix is computed
 149 through Eq.(3). \mathbf{y} is the fitness vector of D_{mi} . μ and σ^2 are the mean and the variance of the prior
 150 distribution, respectively. Experience $\gamma^e = \{\mathbf{w}, \mathbf{b}, \theta^e, \mathbf{p}^e\}$ is updated by gradient descent (line 13),
 151 take θ^e as an example:

$$\theta^e \leftarrow \theta^e - \frac{\alpha}{B} \sum_{i=1}^B \nabla_{\theta^e} L(D_{mi}, h(\gamma)). \quad (5)$$

152 After U iterations, γ^e has been trained sufficiently by N_m small datasets D_m and will be used in
 153 target task T_* later.

Adaptation procedure: Using experience

154 The meta-learning of experience γ^e enables MDKL to handle a family of related tasks in general. To
 155 approximate a specific task T_* well, surrogate $h(\gamma^e)$ needs to adapt task-specific increments $\Delta\theta^*$
 156 and $\Delta\mathbf{p}^*$ in the way described in Alg. 3. A diagram of the deep kernel implemented in our MDKL
 158 model is illustrated in Fig. 2.

159

Algorithm 3 Adaptation(γ^*, S_*, β)

1: **Input:** Current surrogate parameters γ^* ; A dataset S_* sam-
 pled from target task T_* (Archive); Learning rate for adap-
 tation β .
 2: **if** $\gamma^* == \gamma^e$ **then**
 3: Initialize task-specific increments $\Delta\theta^*, \Delta\mathbf{p}^*$.
 4: Compute task-specific parameters: $\theta^* = \theta^e + \Delta\theta^*$,
 $\mathbf{p}^* = \mathbf{p}^e + \Delta\mathbf{p}^*$.
 5: Obtain deep kernel $k(\mathbf{x}^i, \mathbf{x}^j | \gamma^*)$ based GP: $h(\gamma^*)$, where
 $\gamma^* = \{\mathbf{w}, \mathbf{b}, \theta^*, \mathbf{p}^*\}$ (Eq.(3)).
 6: **end if**
 7: Compute the loss function $L(S_*, h(\gamma^*))$ (Eq.(4)).
 8: Update $\Delta\theta^*, \Delta\mathbf{p}^*$ using gradient descent: $\beta \nabla$
 $L(S_*, h(\gamma^*))$.
 9: **Output:** Adapted MDKL $h(\gamma^*)$.

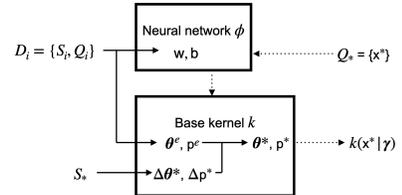


Figure 2: Diagram of our deep kernel implementation. The solid lines depict the training process, the dotted lines depict the inference process. Q_* denotes query samples to be evaluated on our surrogates.

160 **Surrogate prediction.** Due to the nature of a GP, when predicting the fitness of a solution \mathbf{x}^* , a
 161 MDKL surrogate produces a predictive Gaussian distribution $\mathcal{N}(\hat{y}(\mathbf{x}^*), \hat{s}^2(\mathbf{x}^*))$, the predicted mean
 162 $\hat{y}(\mathbf{x}^*)$ and covariance $\hat{s}^2(\mathbf{x}^*)$ are specified as [18]:

$$\hat{y}(\mathbf{x}^*) = \mu + \mathbf{r}'\mathbf{R}^{-1}(\mathbf{y} - \mathbf{1}\mu), \quad (6)$$

163

$$\hat{s}^2(\mathbf{x}^*) = \sigma^2(1 - \mathbf{r}'\mathbf{R}^{-1}\mathbf{r}), \quad (7)$$

164 where \mathbf{r} is a correlation vector consisting of covariances between \mathbf{x}^* and S_* , other variables are
 165 explained in Eq.(4).

166 4.3 Surrogate Update Strategy

167 In this subsection, we describe the update strategy in our FSEO framework. To properly integrate
 168 experience and data from T_* , our update strategy is designed to determine whether a MDKL surrogate
 169 should be adapted in the current iteration or not, ensuring an optimal update frequency of surrogates.

170 As illustrated in Alg. 4, the surrogate update starts when a new optimal solution(s) has been evaluated on
 expensive functions and an updated archive S_* is available. For a given surrogate $h(\gamma^*)$, its mean squared
 error (MSE) on S_* is selected as the update criterion: If the MSE after an adaptation e_1 (line 4) is larger than
 the MSE without an adaptation e_0 (line 2), then the surrogate will roll back to the status before the adaptation.
 This indicates the surrogate update has been refused and $h(\gamma^*)$ will not be adapted in the current iteration.
 Otherwise, the adapted surrogate will be chosen (line 6). Note that no matter whether surrogate adaptations
 are accepted or refused, the resulting surrogates will be treated as updated surrogates, which are employed
 to assist the SAEA optimizer in the next iteration.

Algorithm 4 Update(γ^*, S_*, β)

- 1: **Input:**
 Current surrogate parameters γ^* ;
 Updated archive S_* ;
 Learning rate for further adaptations β .
 - 2: $e_0 \leftarrow \text{MSE}(h(\gamma^*), S_*)$.
 - 3: $h(\gamma') \leftarrow \text{Adaptation}(\gamma^*, S_*, \beta)$.
 /*Temporary surrogate, Alg. 3.*/
 - 4: $e_1 \leftarrow \text{MSE}(h(\gamma'), S_*)$.
 - 5: **if** $e_0 > e_1$ **then**
 - 6: update $\gamma^* = \gamma'$, obtain new $h(\gamma^*)$.
 - 7: **end if**
 - 8: **Output:** Surrogate $h(\gamma^*)$.
-

171 5 Computational Studies

172 Our computational studies can be divided into three parts: (1). Appendix D evaluates the effectiveness
 173 of learning experience through a synthetic problem and a real-world engine modeling problem. (2).
 174 Sections 5.1 to 5.2 use EMOPs as examples to investigate the performance of our FSEO framework
 175 in depth. Empirical evidence is provided to guide the use of our FSEO framework. (3). Section
 176 5.3 investigates the performance of our FSEO framework on a real-world ECOP. The source code
 177 is available online¹ For all meta-learning methods used in our experiments, their basic setups are
 178 listed in Table 1. The neural network structure is suggested by [10, 27], and the learning rates are the
 179 default values that have been widely used in many meta-learning methods [13, 27].

180 5.1 Performance on EMOPs

181 In the following subsections, we aim to demonstrate the effectiveness of our FSEO framework. The
 182 experiment in this subsection is designed to answer the question below: With the experience learned
 183 from related tasks, can our FSEO framework helps a SAEA to save $9d$ solutions without a loss of
 184 optimization performance?

185 The computational study is conducted on the DTLZ test problems [8]. All the DTLZ problems have
 186 $d = 10$ decision variables and 3 objectives, as the setups that have been widely used in [25, 33].
 187 The details of generating DTLZ variants (related tasks) are provided in Appendix C. We test our
 188 FSEO framework using an instantiation on MOEA/D-EGO, resulting MOEA/D-FS. Details of the
 189 comparison algorithms are given in Appendix E.1.

¹A link will be disclosed here once the paper is accepted.

Table 1: Parameter setups for meta-learning methods.

| Module | Parameter | Value |
|----------------|--|--------------|
| Meta-learning | Number of meta-learning datasets N_m | 20000 |
| | Number of update iterations U | 2000 |
| | Batch size B | 10 |
| Neural network | Number of hidden layers | 2 |
| | Number of units in each hidden layer | 40 |
| | Learning rates α, β | 0.001, 0.001 |
| | Activation function | ReLU |

Table 2: Parameter setups for DTLZ optimization.

| Parameter | MOEA/D-FS | Comparisons |
|---|---------------------------|---------------|
| Number of related tasks N | 20000 (N_m in Table 1) | - |
| Size of datasets from related tasks $ D_i $ | 20 ($2d$) | - |
| Size of datasets for meta-learning $ D_m $ | $ D_i $ | - |
| Evaluations for initialization | 10 ($1d$) | 100 ($10d$) |
| Evaluations for further optimization | 50 | 50 |
| Total evaluations | 60 | 150 |

190 5.1.1 Experimental setups

191 The parameter setups for this multi-objective optimization experiment are listed in Table 2. During
 192 the optimization process, an initial dataset S_* is sampled using Latin-Hypercube Sampling (LHS)
 193 method [24], then extra evaluations are conducted until the evaluation budget has run out. Note that
 194 we aim to use related tasks to save $9d$ evaluations without a loss of SAEA optimization performance.
 195 Hence, the total evaluation budgets for MOEA/D-FS and comparison algorithms are different.

196 Since the test problems have 3 objectives, we employ inverted generational distance plus (IGD+) [15]
 197 as our performance indicator, where smaller IGD+ values indicate better optimization results. 5000
 198 reference points are generated for computing IGD+ values, as suggested in [25]. More results in IGD
 199 [4] and HV [55] metrics are reported in Appendix E.3.

200 5.1.2 Results and analysis

201 The statistical test results are reported in Fig. 3 and Appendix E.2 (Table 5). It can be seen from Fig. 3
 202 that, although 90 fewer evaluations are used in surrogate initialization, MOEA/D-FS can still achieve
 203 competitive or even smaller IGD+ values than MOEA/D-EGO on all DTLZ problems except for
 204 DTLZ7. In addition, the IGD+ values obtained by MOEA/D-FS drop rapidly, especially during the
 205 first few evaluations, implying the experience learned from DTLZ variants are effective. Therefore,
 206 in most situations, our FSEO framework is able to assist MOEA/D-EGO in reaching competitive
 207 or even better optimization results, with the number of evaluations used for surrogate initialization
 reduced from $10d$ to only $1d$.

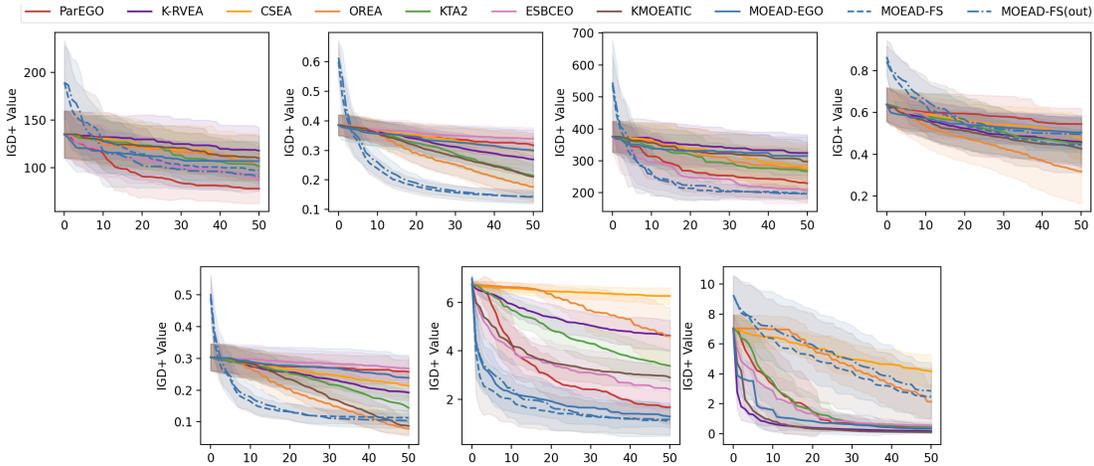


Figure 3: IGD+ curves averaged over 30 runs on the DTLZ problems. Solid lines are mean values, while shadows are error regions. **Upper:** DTLZ1, DTLZ2, DTLZ3, DTLZ4. **Lower:** DTLZ5, DTLZ6, DTLZ7. MOEA/D-FSs and comparison algorithms initialize their surrogates with 10, 100 samples, respectively. X-axis denotes the extra 50 evaluations allowed in the further optimization. Note that ‘FS(out)’ indicates the target task is excluded from the range of related tasks during the meta-learning procedure) (see Appendix F).

208

209 MOEA/D-FS is less effective on DTLZ7 than on other DTLZ problems, which might be attributed to
 210 the discontinuity of the Pareto front on DTLZ7. Note that MOEA/D-FS learns experience from small
 211 datasets such as D_m and S_* . The solutions in these small datasets are sampled at random, hence, the

212 probability of having optimal solutions being sampled is small. However, it is difficult to learn the
213 discontinuity of the Pareto front from the sampled non-optimal solutions. As a result, the knowledge
214 of ‘there are four discrete optimal regions’ cannot be learned from such small datasets ($|D_m| = 20$)
215 collected from related tasks. The performance analysis between MOEA/D-FS and other comparison
216 algorithms are available in Appendix E.2.

217 5.1.3 More comparison experiments

218 We also compared the performance of our FSEO framework when only 10 evaluations are used for
219 surrogate initialization for comparison algorithms. The results are reported in Table 8 in Appendix
220 E.4. In addition, the performance of our FSEO framework in the context of extremely expensive
221 optimization has been investigated in Appendix H (Table 11 and Fig. 7).

222 The question raised at the beginning of this subsection can be answered by the results discussed so
223 far. Due to the integration of the experience learned from related tasks (DTLZ variants), although the
224 evaluation cost of surrogates initialization has been reduced from $10d$ to $1d$, our FSEO framework is
225 still capable of assisting regression-based SAEAs to achieve competitive or even better optimization
226 results in most situations.

227 5.2 Ablation Studies

228 We conduct two ablation studies to investigate the influence of task similarity and that of the dataset
229 size used in meta-learning, results and analysis are reported in Appendixes F and G, respectively.

230 5.3 Performance on Real-World ECOPs

231 The experiments on EMOPs have investigated the performance of our FSEO framework in depth. In
232 this subsection, we evaluate our FSEO framework on a real-world gasoline motor engine calibration
233 problem, which is an ECOP.

234 The calibration problem has 6 adjustable engine parameters, namely the throttle angle, waste gate
235 orifice, ignition timing, valve timings, state of injection, and air-fuel-ratio. The calibration aims at
236 minimizing the BSFC while satisfying 4 constraints in terms of temperature, pressure, CA50, and
237 load simultaneously [53].

238 5.3.1 Comparison algorithms

239 Since the comparison algorithms in the DTLZ optimization experiments are not designed for handling
240 constrained optimization, our comparison is conducted with 3 state-of-the-art constrained optimization
241 algorithms used in industry: A variant of EGO designed to handle constrained optimization problems
242 (cons_EGO) [53], a GA customized for this calibration problem (adaptiveGA) [53], and a bilevel
243 constrained SAEA (SAB-DE) [50]. The settings of the comparison algorithms are the same as
244 suggested in the literature. In this experiment, we apply our FSEO framework to cons_EGO and
245 investigate its optimization performance. The GP surrogates in cons_EGO are replaced by our MDKL
246 surrogates to conduct the comparison, and the resulting algorithm is denoted as cons_FS.

247 5.3.2 Experimental setups

248 The setup of related tasks (N, D_i) is the same as described in Appendix D. In the meta-learning
249 procedure, both the support set and the query set contain 6 data points, thus $|D_m| = 12$. The total
250 evaluation budget for all algorithms is set to 60. For adaptiveGA, all evaluations are used in the
251 optimization process as it is not a SAEA. For cons_EGO and SAB-DE, 40 samples are used to
252 initialize the surrogates and 20 extra evaluations are used in the optimization process. For cons_FS,
253 only 6 samples are used to initialize MDKL surrogates, and the remaining evaluations are used for
254 further optimization.

255 5.3.3 Optimization results and analysis

256 The left side and right side of Fig. 4 plot the normalized BSFC results and the number of feasible
257 solutions found over the number of used evaluations, respectively. Solid lines are mean lines, while
258 shadows are error regions. From the left side of Fig. 4, it can be observed that the minimal BSFC

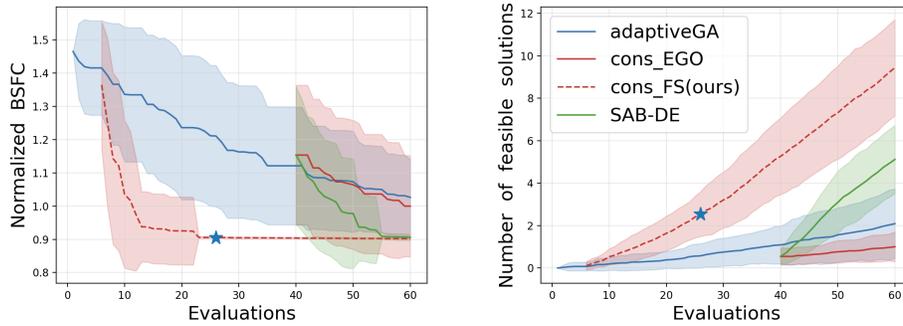


Figure 4: Results of 30 runs on the real-world engine calibration problem, all BSFC values are normalized. Solid lines are mean values, while shadows are error regions. Left figure shows how BSFC varies with the number of evaluations. The star markers highlight the results achieved when 20 evaluations are used in the optimization process. Right figure illustrates how the number of feasible solutions varies with the number of evaluations.

259 obtained by cons_FS decreases drastically in the first few evaluations, implying that the experience
 260 learned from related tasks is effective. In comparison, the minimal BSFC obtained by adaptiveGA
 261 and cons_EGO drops in a relatively slow rate, even though cons_EGO has used 34 more samples
 262 to initialize its surrogates. The star marker denotes the point at which cons_FS has evaluated 20
 263 samples after surrogate initialization. It is worth noting that when 20 samples have been evaluated
 264 in the optimization, cons_FS achieves a smaller BSFC value than cons_EGO. After the star marker,
 265 the decrease of BSFC becomes slow as cons_FS has reached the optimal region. Therefore, further
 266 improvement in the normalized BSFC value is not significant and thus hard to be observed. The
 267 advantages of our FSEO framework can also be observed in constraint handling. In the right side of
 268 Fig. 4, cons_FS finds more feasible solutions than the 3 comparison algorithms. These results indicate
 269 that our FSEO framework improves the performance of cons_EGO on both objective function and
 270 constraint functions. Meanwhile, only 1d evaluations are used to initialize surrogates.

271 5.3.4 Discussion on runtime

272 It should be noted that real engine performance evaluations on engine facilities are very costly in
 273 terms of both time and financial budget [49]. Since a single real engine performance evaluation can
 274 cost several hours [22, 49], the time cost of the meta-learning procedure is negligible as it takes only
 275 a few minutes. Savings from reduced real engine performance evaluations on engine facilities and the
 276 reduced development cycle due to our FSEO framework could amount to millions of dollars [49]. our
 277 FSEO framework is an effective and efficient method to solve this real-world calibration problem.

278 6 Conclusion and further work

279 **Conclusion.** In this paper, we present a FSEO framework to address EMOPs and ECOPs from the
 280 perspective of SAEAs. A novel meta-learning approach MDKL is proposed to learn prior experience
 281 from related expensive tasks. Our MDKL model is designed for optimization and has explicit
 282 task-specific parameters, which allows continually update of task-specific parameters during the
 283 optimization process. Our empirical experiments show that the FSEO framework is able to improve
 284 the sampling efficiency and thus save expensive evaluations for existing regression-based SAEAs.
 285 Ablation studies reveal the influence between optimization performance and solutions similarity as
 286 well as the size of datasets for meta-learning.

287 **Limitation and further work.** The limitations of this work can be summarized as the following
 288 two points: First, we do not have a mathematical definition of related tasks. Second, the proposed
 289 framework is currently for regression-based SAEAs only.

References

- [1] Tianyi Bai, Yang Li, Yu Shen, Xinyi Zhang, Wentao Zhang, and Bin Cui. Transfer learning for bayesian optimization: A survey. *arXiv preprint arXiv:2302.05927*, 2023.
- [2] Kavitesh Kumar Bali, Yew-Soon Ong, Abhishek Gupta, and Puay Siew Tan. Multifactorial evolutionary algorithm with online transfer parameter estimation: MFEA-II. *IEEE Transactions on Evolutionary Computation*, 24(1):69–83, 2019.
- [3] Hongli Bian, Jie Tian, Jialiang Yu, and Han Yu. Bayesian co-evolutionary optimization based entropy search for high-dimensional many-objective optimization. *Knowledge-Based Systems*, 274:110630, 2023.
- [4] Peter AN Bosman and Dirk Thierens. The balance between proximity and diversity in multiobjective evolutionary algorithms. *IEEE Transactions on Evolutionary Computation*, 7(2):174–188, 2003.
- [5] Wei-Yu Chen, Yen-Cheng Liu, Zsolt Kira, Yu-Chiang Frank Wang, and Jia-Bin Huang. A closer look at few-shot classification. In *Proceedings of the 7th International Conference on Learning Representations (ICLR'19)*, 2019.
- [6] Wenlin Chen, Austin Tripp, and José Miguel Hernández-Lobato. Meta-learning adaptive deep kernel gaussian processes for molecular property prediction. In *Proceedings of the 11th International Conference on Learning Representations (ICLR'23)*, 2023.
- [7] Tinkle Chugh, Yaochu Jin, Kaisa Miettinen, Jussi Hakanen, and Karthik Sindhya. A surrogate-assisted reference vector guided evolutionary algorithm for computationally expensive many-objective optimization. *IEEE Transactions on Evolutionary Computation*, 22(1):129–142, 2016.
- [8] Kalyanmoy Deb, Lothar Thiele, Marco Laumanns, and Eckart Zitzler. Scalable test problems for evolutionary multiobjective optimization. In *Evolutionary Multiobjective Optimization*, pages 105–145. Springer, London, U.K., 2005.
- [9] Jinliang Ding, Cuie Yang, Yaochu Jin, and Tianyou Chai. Generalized multitasking for evolutionary optimization of expensive problems. *IEEE Transactions on Evolutionary Computation*, 23(1):44–58, 2017.
- [10] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. In *Proceedings of the 34th International Conference on Machine Learning (ICML'17)*, pages 1126–1135, 2017.
- [11] Zhendong Guo, Haitao Liu, Yew-Soon Ong, Xinghua Qu, Yuzhe Zhang, and Jianmin Zheng. Generative multiform Bayesian optimization. *IEEE Transactions on Cybernetics*, 53(7):4347–4360, 2022.
- [12] Abhishek Gupta, Yew-Soon Ong, and Liang Feng. Insights on transfer optimization: Because experience is the best teacher. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 2(1):51–64, 2017.
- [13] James Harrison, Apoorva Sharma, and Marco Pavone. Meta-learning priors for efficient online Bayesian regression. In *Proceedings of the 13th Workshop on the Algorithmic Foundations of Robotics (WAFR'18)*, pages 318–337, 2018.
- [14] Timothy M. Hospedales, Antreas Antoniou, Paul Micaelli, and Amos J Storkey. Meta-learning in neural networks: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2021.
- [15] Hisao Ishibuchi, Hiroyuki Masuda, Yuki Tanigaki, and Yusuke Nojima. Modified distance calculation in generational distance and inverted generational distance. In *Proceedings of the 8th International Conference on Evolutionary Multi-criterion Optimization (EMO'15)*, pages 110–125, 2015.
- [16] Min Jiang, Zhenzhong Wang, Shihui Guo, Xing Gao, and Kay Chen Tan. Individual-based transfer learning for dynamic multiobjective optimization. *IEEE Transactions on Cybernetics*, 51(10):4968–4981, 2020.

- 339 [17] Min Jiang, Zhenzhong Wang, Liming Qiu, Shihui Guo, Xing Gao, and Kay Chen Tan. A fast
340 dynamic evolutionary multiobjective algorithm via manifold transfer learning. *IEEE Transactions*
341 *on Cybernetics*, 51(7):3417–3428, 2020.
- 342 [18] Donald R. Jones, Matthias Schonlau, and William J. Welch. Efficient global optimization of
343 expensive black-box functions. *Journal of Global Optimization*, 13(4):455–492, 1998.
- 344 [19] Joshua Knowles. ParEGO: A hybrid algorithm with on-line landscape approximation for
345 expensive multiobjective optimization problems. *IEEE Transactions on Evolutionary Computation*,
346 10(1):50–66, 2006.
- 347 [20] Rung-Tzuo Liaw and Chuan-Kang Ting. Evolutionary manytasking optimization based on
348 symbiosis in biocoenosis. In *Proceedings of the 33rd AAAI Conference on Artificial Intelligence*
349 *(AAAI’19)*, pages 4295–4303, 2019.
- 350 [21] Shengcai Liu, Ke Tang, and Xin Yao. Experience-based optimization: A coevolutionary
351 approach. *arXiv preprint arXiv:1703.09865*, 2017.
- 352 [22] He Ma. *Control Oriented Engine Modeling and Engine Multi-objective Optimal Feedback*
353 *Control*. PhD thesis, University of Birmingham, 2013.
- 354 [23] Alexandre Maraval, Matthieu Zimmer, Antoine Grosnit, and Haitham Bou Ammar. End-to-end
355 meta-bayesian optimisation with transformer neural processes. *arXiv preprint arXiv:2305.15930*,
356 2023.
- 357 [24] Michael D. McKay, Richard J. Beckman, and William J. Conover. A comparison of three
358 methods for selecting values of input variables in the analysis of output from a computer code.
359 *Technometrics*, 42(1):55–61, 2000.
- 360 [25] Linqiang Pan, Cheng He, Ye Tian, Handing Wang, Xingyi Zhang, and Yaochu Jin. A
361 classification-based surrogate-assisted evolutionary algorithm for expensive many-objective opti-
362 mization. *IEEE Transactions on Evolutionary Computation*, 23(1):74–88, 2018.
- 363 [26] Jiarong Pan, Stefan Falkner, Felix Berkenkamp, and Joaquin Vanschoren. MALIBO: Meta-
364 learning for likelihood-free bayesian optimization. *arXiv preprint arXiv:2307.03565*, 2023.
- 365 [27] Massimiliano Patacchiola, Jack Turner, Elliot J Crowley, Michael O’Boyle, and Amos Storkey.
366 Bayesian meta-learning for the few-shot setting via deep kernels. In *Advance in Neural Information*
367 *Processing Systems 33 (NeurIPS’20)*, 2020.
- 368 [28] Shufen Qin, Chaoli Sun, Farooq Akhtar, and Gang Xie. Expensive many-objective evolutionary
369 optimization guided by two individual infill criteria. *Memetic Computing*, pages 1–15, 2023.
- 370 [29] Gan Ruan, Leandro L. Minku, Stefan Menzel, Bernhard Sendhoff, and Xin Yao. When and how
371 to transfer knowledge in dynamic multi-objective optimization. In *Proceedings of the 2019 IEEE*
372 *Symposium Series on Computational Intelligence (SSCI’19)*, pages 2034–2041, 2019.
- 373 [30] Gan Ruan, Leandro L. Minku, Stefan Menzel, Bernhard Sendhoff, and Xin Yao. Computa-
374 tional study on effectiveness of knowledge transfer in dynamic multi-objective optimization. In
375 *Proceedings of the 22nd IEEE Congress on Evolutionary Computation (CEC’20)*, pages 1–8,
376 2020.
- 377 [31] Jerome Sacks, William J. Welch, Toby J. Mitchell, and Henry P. Wynn. Design and analysis of
378 computer experiments. *Statistical Science*, 4(4):409–423, 1989.
- 379 [32] Gresa Shala, Thomas Elsken, Frank Hutter, and Josif Grabocka. Transfer NAS with meta-
380 learned bayesian surrogates. In *Proceedings of the 11th International Conference on Learning*
381 *Representations (ICLR’23)*, 2023.
- 382 [33] Zhenshou Song, Handing Wang, Cheng He, and Yaochu Jin. A Kriging-assisted two-archive
383 evolutionary algorithm for expensive many-objective optimization. *IEEE Transactions on Evolu-*
384 *tionary Computation*, 25(6):1013–1027, 2021.
- 385 [34] Michael L. Stein. *Interpolation of Spatial Data: Some Theory for Kriging*. Springer Science &
386 Business Media, New York, NY, 1999.

- 387 [35] Kay Chen Tan, Liang Feng, and Min Jiang. Evolutionary transfer optimization—a new frontier
388 in evolutionary computation research. *IEEE Computational Intelligence Magazine*, 16(1):22–33,
389 2021.
- 390 [36] Ke Tang, Shengcai Liu, Peng Yang, and Xin Yao. Few-shots parallel algorithm portfolio
391 construction via co-evolution. *IEEE Transactions on Evolutionary Computation*, 25(3):595–607,
392 2021.
- 393 [37] Ye Tian, Ran Cheng, Xingyi Zhang, and Yaochu Jin. PlatEMO: A MATLAB platform for
394 evolutionary multi-objective optimization [educational forum]. *IEEE Computational Intelligence
395 Magazine*, 12(4):73–87, 2017.
- 396 [38] Petru Tighineanu, Lukas Grossberger, Paul Baireuther, Kathrin Skubch, Stefan Falkner, Julia
397 Vinogradska, and Felix Berkenkamp. Scalable meta-learning with gaussian processes. *arXiv
398 preprint arXiv:2312.00742*, 2023.
- 399 [39] Prudencio Tossou, Basile Dura, Francois Laviolette, Mario Marchand, and Alexandre Lacoste.
400 Adaptive deep kernel learning. *arXiv preprint arXiv:1905.12131*, 2019.
- 401 [40] Michael Volpp, Lukas P. Fröhlich, Kirsten Fischer, Andreas Doerr, Stefan Falkner, Frank
402 Hutter, and Christian Daniel. Meta-learning acquisition functions for transfer learning in Bayesian
403 optimization. In *Proceedings of the 8th International Conference on Learning Representations
404 (ICLR’20)*, 2020.
- 405 [41] Yaqing Wang, Quanming Yao, James T. Kwok, and Lionel M. Ni. Generalizing from a few
406 examples: A survey on few-shot learning. *ACM Computing Surveys*, 53(3):1–34, 2020.
- 407 [42] Zi Wang, George E. Dahl, Kevin Swersky, Chansoo Lee, Zachary Nado, Justin Gilmer, Jasper
408 Snoek, and Zoubin Ghahramani. Pre-trained gaussian processes for bayesian optimization. *arXiv
409 preprint arXiv:2109.08215*, 2021.
- 410 [43] Tingyang Wei, Shibin Wang, Jinghui Zhong, Dong Liu, and Jun Zhang. A review on evolutionary
411 multi-task optimization: Trends and challenges. *IEEE Transactions on Evolutionary Computation*,
412 26(5):941–960, 2021.
- 413 [44] Christopher KI Williams and Carl Edward Rasmussen. *Gaussian Processes for Machine
414 Learning*. MIT press, Cambridge, MA, 2006.
- 415 [45] Andrew Gordon Wilson, Zhiting Hu, Ruslan Salakhutdinov, and Eric P. Xing. Deep kernel
416 learning. In *Proceedings of the 19th International Conference on Artificial Intelligence and
417 Statistics (AISTATS’16)*, pages 370–378, 2016.
- 418 [46] Martin Wistuba and Josif Grabocka. Few-shot Bayesian optimization with deep kernel surro-
419 gates. In *Proceedings of the 9th International Conference on Learning Representations (ICLR’21)*,
420 2021.
- 421 [47] Xiaoming Xue, Kai Zhang, Kay Chen Tan, Liang Feng, Jian Wang, Guodong Chen, Xinggang
422 Zhao, Liming Zhang, and Jun Yao. Affine transformation-enhanced multifactorial optimization
423 for heterogeneous problems. *IEEE Transactions on Cybernetics*, pages 1–15, 2020.
- 424 [48] Xunzhao Yu, Xin Yao, Yan Wang, Ling Zhu, and Dimitar Filev. Domination-based ordinal
425 regression for expensive multi-objective optimization. In *Proceedings of the 2019 IEEE Symposium
426 Series on Computational Intelligence (SSCI’19)*, pages 2058–2065, 2019.
- 427 [49] Xunzhao Yu, Ling Zhu, Yan Wang, Dimitar Filev, and Xin Yao. Internal combustion engine
428 calibration using optimization algorithms. *Applied Energy*, 305:117894, 2022.
- 429 [50] Xunzhao Yu, Yan Wang, Ling Zhu, Dimitar Filev, and Xin Yao. Engine calibration with
430 surrogate-assisted bilevel evolutionary algorithm. *IEEE Transactions on Cybernetics (Early
431 Access)*, 2023.
- 432 [51] Qingfu Zhang, Wudong Liu, Edward Tsang, and Botond Virginas. Expensive multiobjective
433 optimization by MOEA/D with gaussian process model. *IEEE Transactions on Evolutionary
434 Computation*, 14(3):456–474, 2010.

- 435 [52] Liangjie Zhang, Yuling Xie, Jianjun Chen, Liang Feng, Chao Chen, and Kai Liu. A study on
436 multiform multi-objective evolutionary optimization. *Memetic Computing*, 13(3):307–318, 2021.
- 437 [53] Ling Zhu, Yan Wang, Anuj Pal, and Guoming Zhu. Engine calibration using global optimization
438 methods with customization. Technical Report 2020-01-0270, SAE Technical Paper, 2020.
- 439 [54] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui
440 Xiong, and Qing He. A comprehensive survey on transfer learning. *Proceedings of the IEEE*,
441 109(1):43–76, 2020.
- 442 [55] Eckart Zitzler and Lothar Thiele. Multiobjective optimization using evolutionary algorithms - a
443 comparative case study. In *Proceedings of the 5th International Conference on Parallel Problem
444 Solving from Nature (PPSN V)*, pages 292–301, 1998.

445 **A Discussion on Experience-based Optimization**

446 In the past decade, experience-based optimization has attracted much attention as it uses the experience
447 gained from other optimization problems to improve the optimization efficiency of target problems,
448 which mimics human capabilities of cognitive and knowledge generalization [12]. The optimization
449 problems that provide experience or knowledge are regarded as source tasks, while the target
450 optimization problems are regarded as target tasks. To obtain useful experience, the tasks that are
451 related to target tasks are chosen as source tasks since they usually share domain-specific features
452 with target tasks. Diverse experience-based optimization methods have been proposed to use the
453 experience gained from related tasks to tackle target tasks. They can be divided into two categories
454 based on the direction of experience transformation.

455 In the first category, experience is transformed mutually. Every considered optimization problem is a
456 target task and also is one of the source tasks of other optimization problems. In other words, the
457 roles of source task and target task are compatible. One representative tributary is EMTO that aims to
458 solve multiple optimization tasks concurrently [9, 43, 20, 2, 47]. In EMTO, experience is learned,
459 updated, and spontaneously shared among target tasks through multi-task learning techniques. A
460 variant of EMTO is multiforms optimization [12, 52, 11]. In multiforms optimization, multi-task
461 learning methods are employed to learn experience from distinct formulations of a single target task.

462 In the second category, experience is transformed unidirectionally. The roles of source task and
463 target task are not compatible, an optimization problem cannot serve as a source task and a target
464 task simultaneously. One popular tributary is transfer optimization which employs transfer learning
465 techniques to transform experience from source tasks to target tasks [35, 17, 16, 40]. In transfer
466 learning, experience can be transformed from a single source task, multiple source tasks, or even
467 source tasks from a different domain [54]. However, these transfer learning techniques pay more
468 attention to experience transformation instead of experience learning. Despite diverse and complex
469 situations of experience transformation have been studied [29, 30], the difficult of learning experience
470 from small (expensive) source tasks has not been well studied. Actually, a common scenario in
471 transfer learning is that the source task(s) is/are large enough such that useful experience can be
472 obtained easily through solving source task(s) [54]. In contrast to transfer optimization, recently, some
473 experience-based optimization algorithms attempted to use meta-learning methods to learn experience
474 from small source tasks, which are known as few-shot optimization (FSO)[46]. Since meta-learning
475 only works for related tasks in the same domain, the situations of experience transformation are less
476 complex than that of transfer learning. As a result, meta-learning pays more attention to experience
477 learning instead of experience transformation. Domain-specific features are extracted as experience
478 and no related task needs to be solved.

479 Our work belongs to the FSO in the second category discussed above since our experience is
480 transformed unidirectionally. More importantly, our experience is learned across many related
481 expensive tasks, rather than gained through solving more or less source tasks. Therefore, our work is
482 different from transfer optimization.

483 **B Discussion on Framework Compatibility and Limitation**

484 Our FSEO framework is applicable to regression-based SAEAs as our MDKL surrogates can be
485 embedded in these SAEAs directly. Classification-based SAEAs are not compatible with our FSEO
486 framework. The classification surrogates in these SAEAs are employed to learn the relation between
487 pairs of solutions, or the relation between solutions and a set of reference solutions. The class labels
488 used for surrogate training can be fluctuating during the optimization and thus hard to be learned
489 over related tasks. Similarly, in ordinal-regression-based SAEAs, the ordinal relation values to be
490 learned are not as stable as the fitness of expensive functions. So ordinal-regression-based SAEAs are
491 also not compatible with our FSEO framework. In this paper, we focus on FSO for regression-based
492 SAEAs, while other SAEA categories are left to be discussed in future work.

493 **C Generation of DTLZ variants**

494 Our DTLZ optimization experiments generate m -objective DTLZ variants in the following ways:

495 **DTLZ1:**

$$f_1 = (a_1 + g)0.5 \prod_{i=1}^{m-1} x_i, \quad (8)$$

496

$$f_{j=2:m-1} = (a_j + g)(0.5 \prod_{i=1}^{m-j} x_i)(1 - x_{m-j+1}), \quad (9)$$

497

$$f_m = (a_m + g)0.5(1 - x_1), \quad (10)$$

498

$$g = 100 \left(k + \sum_{i=1}^k ((z_i - 0.5)^2 - \cos(20\pi(z_i - 0.5))) \right), \quad (11)$$

499 where \mathbf{z} is a vector consisting of the last $k = d - m + 1$ variables in \mathbf{x} . In other words,
500 $\mathbf{z} = \{z_1, \dots, z_k\} = \{x_m, \dots, x_d\}$. The variants of DTLZ1 introduce only one variable
501 $\mathbf{a} \in [0.1, 5.0]^m$ in Eq.(8), Eq.(9), and Eq.(10), where $\mathbf{a} = \mathbf{1}$ in the original DTLZ1. For out-of-range
502 test, $\mathbf{a} \in [1.5, 5.0]^m$.

503

504 **DTLZ2:**

$$f_1 = (a_1 + g) \prod_{i=1}^{m-1} \cos\left(\frac{x_i \pi}{b_1}\right), \quad (12)$$

505

$$f_{j=2:m-1} = (a_j + g) \left(\prod_{i=1}^{m-j} \cos\left(\frac{x_i \pi}{b_j}\right) \right) \sin\left(\frac{x_{m-j+1} \pi}{b_j}\right), \quad (13)$$

506

$$f_m = (a_m + g) \sin\left(\frac{x_1 \pi}{b_m}\right), \quad (14)$$

507

$$g = \sum_{i=1}^k (z_i - 0.5)^2. \quad (15)$$

508 The variants of DTLZ2 introduce two variables $\mathbf{a} \in [0.1, 5.0]^m$ and $\mathbf{b} \in [0.5, 2.0]^m$ in Eq.(12),
509 Eq.(13), and Eq.(14), where $\mathbf{a} = \mathbf{1}$ and $\mathbf{b} = \mathbf{2}$ in the original DTLZ2. For out-of-range test,
510 $\mathbf{a} \in [1.5, 5.0]^m$, $\mathbf{b} \in [0.5, 1.5]^m$.

511

512 **DTLZ3:** The variants of DTLZ3 are generated using the same way as described in DTLZ2, except
513 the equation g from Eq.(15) is replaced by the one from Eq.(11).

514

515 **DTLZ4:** The variants of DTLZ4 are generated using the same way as described in DTLZ2, except
516 all x_i are replaced by x_i^{100} .

517

518 **DTLZ5:** The variants of DTLZ5 are generated using the same way as described in DTLZ2, except
519 all x_2, \dots, x_{m-1} are replaced by $\frac{1+2gx_i}{2(1+g)}$.

520

521 **DTLZ6:**

$$g = \sum_{i=1}^k z_i^{0.1}. \quad (16)$$

522 The variants of DTLZ6 are generated using the same way as described in DTLZ5, except the equation
523 g from Eq.(15) is replaced by the one from Eq.(16).

524

525 **DTLZ7:**

$$f_{j=1:m-1} = x_j + a_j, \quad (17)$$

526

$$f_m = (1 + g) \left(m - \sum_{i=1}^{m-1} \left(\frac{f_i}{1 + g} (1 + \sin(3\pi f_i)) \right) \right), \quad (18)$$

527

$$g = a_m + 9 \sum_{i=1}^k \frac{z_i}{k}. \quad (19)$$

528 The variants of DTLZ7 introduce one variable $\mathbf{a} \in [0.1, 5.0]^m$ in Eq.(17) and Eq.(19), where
 529 $a_{j=1:m-1} = 0$ and $a_m = 1$ in the original DTLZ7. For out-of-range test, $\mathbf{a} \in [1.5, 5.0]^m$.

530 D Effectiveness of Learning Experience

531 Evaluating the effectiveness of learning experiences aims to demonstrate that our MDKL model can
 532 learn experience from related tasks and outperforms other meta-learning models. For this reason, the
 533 experiment is designed to answer the following questions:

- 534 • Given a small dataset S_* from target task T_* , can MDKL learn experience from related tasks
 535 and then generate a model that has the smallest MSE?
- 536 • If yes, which components of MDKL contribute to the effectiveness of learning experience?
 537 Meta-learning or/and deep kernel learning? If not, why not?

538 To answer the two questions above, we consider two experiments to evaluate the effectiveness
 539 of learning experience: amplitude prediction for unknown periodic sinusoid functions, and fuel
 540 consumption prediction for a gasoline motor engine. The former is a few-shot regression problem
 541 that motivates many meta-learning studies [10, 13, 39, 27], while the latter is a real-world regression
 542 problem [53].

543 D.1 Effectiveness of Learning Experience: Sinusoid Function Regression

544 In the sinusoid regression experiment, we learn experience from a series of 1-dimensional sinusoid
 545 functions:

$$y = A \sin(wx + b) + \epsilon, \quad (20)$$

546 where the amplitude A and phase w of sine waves are varied between functions. The target is to
 547 approximate an unknown sinusoid function with a small support dataset S_* and the learned experience.
 548 Clearly, by integrating experience with S_* , we estimate parameters (A, w, b) for an unknown sinusoid
 549 function. As a result, the output y of the given sinusoid function can be predicted once a query data x
 550 is inputted.

551 D.1.1 Generation of Sinusoid Function Variants

552 As suggested in [10, 13], we set amplitude $A \in [0.1, 5.0]$, frequency $w \in [0.999, 1.0]$, phase $b \in [0,$
 553 $\pi]$, and Gaussian noise $\epsilon \sim (0, 0.1)$. Therefore, a sinusoid function can be generated by sampling
 554 three parameters (A, w, b) from their ranges uniformly. In total, $N_m = N = 20000$ related sinusoid
 555 functions are generated at random.

556 D.1.2 Experimental Setups

557 All data points x are sampled from the range $\in [-5.0, 5.0]$. In the meta-learning procedure, both
 558 support set and query set contain 5 data points. Hence, a dataset D_i is sampled from each (related)
 559 sinusoid function T_i , and $|D_i| = |D_m| = 10$. Six experiments are conducted where $|S_*| =$
 560 $\{2, 3, 5, 10, 20, 30\}$ data points are sampled from the target function. Considering Gaussian noise ϵ
 561 could be relatively large when amplitude A is close to 0.1, normalized mean squared error (NMSE) is
 562 chosen as a performance indicator. NMSE is measured using a dataset that contains 100 data points
 563 sampled uniformly from the x range.

Table 3: Mean NMSE and standard deviation (in parentheses) of 30 runs on the amplitude regression of sinusoid function. GP [34] is a widely used surrogate in SAEAs, MAML [10], ALPaCA [13], and DKT [27] are meta-learning methods. GP_Adam is a GP model fitted by Adam optimizer. DKL is a deep kernel learning algorithm that adds a neural network to GP_Adam. MDKL_NN applies meta-learning to DKL, but no task-independent base kernel parameters are shared between related tasks. Support data points are used to train non-meta surrogates or adapt meta-learning surrogates. ‘+’, ‘≈’, and ‘-’ denote MDKL is statistically significantly superior to, almost equivalent to, and inferior to the compared modelling methods in the Wilcoxon rank sum test (significance level is 0.05), respectively. The last row counts the total win/tie/loss results. It shows that MDKL and DKT have lower NMSE than other models. The effectiveness of meta-learning on both the neural network and the base kernel has been demonstrated on this example.

| Support data points $ S_* $ | GP [34] | GP_Adam | DKL | MDKL_NN | MDKL (ours) | DKT [27] | MAML [10] | ALPaCA [13] |
|-----------------------------|-------------------|-------------------|-------------------|-------------------|------------------|-------------------|-------------------|-------------------|
| 2 | 1.63e-1(9.18e-2)≈ | 1.93e-1(9.72e-2)+ | 1.63e-1(9.05e-2)≈ | 1.57e-1(9.26e-2)≈ | 1.56e-1(9.49e-2) | 1.56e-1(9.49e-2)≈ | 2.09e-1(3.63e-1)≈ | 1.07e+0(2.57e+0)≈ |
| 3 | 1.27e-1(6.04e-2)≈ | 1.62e-1(6.53e-2)+ | 1.21e-1(5.96e-2)≈ | 1.16e-1(5.95e-2)≈ | 1.10e-1(6.20e-2) | 1.10e-1(6.20e-2)≈ | 2.09e-1(3.60e-1)≈ | 4.36e-1(8.57e-1)≈ |
| 5 | 6.76e-2(4.62e-2)≈ | 1.09e-1(5.61e-2)+ | 7.52e-2(4.40e-2)+ | 6.38e-2(3.91e-2)≈ | 4.79e-2(3.73e-2) | 4.79e-2(3.70e-2)≈ | 2.08e-1(3.59e-1)+ | 4.31e-1(8.04e-1)≈ |
| 10 | 1.70e-2(1.87e-2)≈ | 6.13e-2(4.58e-2)+ | 2.87e-2(1.89e-2)+ | 1.89e-2(1.61e-2)+ | 1.07e-2(1.16e-2) | 1.09e-2(1.17e-2)≈ | 2.08e-1(3.58e-1)+ | 6.59e-1(2.14e+0)+ |
| 20 | 5.42e-3(7.64e-3)+ | 3.92e-2(4.29e-2)+ | 9.64e-3(1.02e-2)+ | 5.24e-3(6.57e-3)+ | 2.57e-3(4.53e-3) | 2.63e-3(4.61e-3)≈ | 2.08e-1(3.58e-1)+ | 1.13e-1(3.39e-1)+ |
| 30 | 3.97e-3(7.40e-3)+ | 3.32e-2(4.18e-2)+ | 4.81e-3(6.68e-3)+ | 3.20e-3(3.85e-3)+ | 1.68e-3(3.61e-3) | 1.60e-3(3.39e-3)≈ | 2.08e-1(3.58e-1)+ | 7.59e-2(2.01e-1)+ |
| + / ≈ / - | 2/4/0 | 6/0/0 | 4/2/0 | 3/3/0 | -/- | 0/6/0 | 4/2/0 | 3/3/0 |

564 D.1.3 Comparison methods

565 In this experiment, three families of modeling methods are compared with our MDKL model:

- 566 • **Meta-learning methods** that were proposed for regression tasks: MAML [10], ALPaCA
567 [13], and DKT [27]. The configurations of MAML, ALPaCA, and DKT are the same as
568 suggested in the original literature.
- 569 • **Non-meta-learning method** that is widely used for regression tasks: the GP model. We
570 choose a GP as a baseline since it is effective and more relevant to MDKL than other
571 non-meta-learning modeling methods. We set the range of base kernel parameters in the GP
572 model as $\theta \in [10^{-5}, 10]$ and $p \in [1, 2]$.
- 573 • **MDKL related methods** that are designed to investigate which components of MDKL
574 contribute to the modeling performance: GP_Adam, DKL, and MDKL_NN. GP_Adam is a
575 GP model fitted by Adam optimizer. The combination of GP_Adam and a neural network
576 results in a kind of DKL algorithm. MDKL_NN is a meta-learning version of DKL, but it
577 learns only neural network parameters through meta-learning and has no task-independent
578 base kernel parameters.

579 D.1.4 Results and Analysis

580 Table 3 reports the statistical test results of the NMSE values achieved by comparison algorithms
581 in sinusoid function regression experiments. Each row lists the results obtained when the same
582 number of fitness evaluations $|S_*|$ are used to train models. The results of Wilcoxon rank sum test
583 between MDKL and other compared algorithms are listed in the last row. It can be observed that both
584 MDKL and DKT have achieved the smallest NMSE values on all tests in the comparison with other
585 meta-learning and non-meta-learning modeling methods.

586 Contributions of MDKL components are analyzed through statistical tests between MDKL related
587 methods. The statistical test results between DKL and GP_Adam are 5/1/0, showing that DKL is
588 preferable to GP_Adam when only a few data points are available for modeling. Hence, using a
589 neural network to build a deep kernel for GP is able to enhance the performance of modeling. When
590 meta-learning technique is applied to DKL, the statistical test results between MDKL_NN and DKL
591 are 3/3/0. The meta-learning of neural network parameters is necessary since it contributes to the
592 performance of MDKL. Further statistical test between MDKL and MDKL_NN gives results of 3/3/0,
593 indicating that the meta-learning of base kernel parameters is effective on this regression problem.

594 D.2 Effectiveness of Learning Experience: Engine Performance Regression

595 In this subsection, we focus on a Brake Special Fuel Consumption (BSFC) regression task for a
596 gasoline motor engine [53], where BSFC is evaluated on a gasoline engine simulation (denoted by
597 T_*).

Table 4: Mean MSE and standard deviation (in parentheses) of 30 runs on the regression of engine fuel consumption. Support data points are used to train non-meta surrogates or adapt meta-learning surrogates. All results are normalized since the actual engine data is unable to be disclosed. The symbols ‘+’, ‘ \approx ’, ‘-’ denote the win/tie/loss result of Wilcoxon rank sum test (significance level is 0.05) between MDKL and comparison modeling methods, respectively. The last row counts the total win/tie/loss results.

| Support data points $ S_i $ | GP [34] | GP_Adam | DKL | MDKL_NN | MDKL (ours) | DKT [27] | MAML [10] | ALPaCA [13] |
|-----------------------------|-------------------|-------------------|-------------------|----------------------------|------------------|----------------------------|----------------------------|----------------------------|
| 2 | 2.23e+1(3.20e+0)+ | 2.37e+1(6.30e+0)+ | 2.30e+1(5.87e+0)+ | 1.73e+1(6.33e+0) \approx | 1.72e+1(6.34e+0) | 1.81e+1(5.68e+0) \approx | 1.87e+1(6.37e+0) \approx | 1.91e+1(1.02e+1) \approx |
| 3 | 2.14e+1(3.74e+0)+ | 2.41e+1(1.38e+1)+ | 2.20e+1(3.74e+0)+ | 1.45e+1(7.13e+0) \approx | 1.45e+0(7.01e+0) | 1.55e+1(6.66e+0) \approx | 1.80e+1(4.69e+0) \approx | 2.13e+1(1.97e+1) \approx |
| 5 | 2.13e+1(3.27e+0)+ | 2.46e+1(1.00e+1)+ | 2.07e+1(3.95e+0)+ | 1.12e+1(6.65e+0) \approx | 1.10e+1(6.58e+0) | 1.21e+1(6.49e+0) \approx | 1.84e+1(6.05e+0)+ | 1.99e+1(2.29e+1)+ |
| 10 | 1.84e+1(1.89e+0)+ | 2.06e+1(1.19e+1)+ | 2.10e+1(5.79e+0)+ | 7.19e+0(4.82e+0) \approx | 7.08e+0(4.77e+0) | 7.99e+0(4.87e+0) \approx | 1.70e+1(5.54e+0)+ | 1.38e+1(8.12e+0)+ |
| 20 | 1.56e+1(2.00e+0)+ | 2.38e+1(1.05e+1)+ | 1.76e+1(2.42e+0)+ | 5.03e+0(1.82e+0) \approx | 4.86e+0(1.71e+0) | 5.74e+0(1.91e+0)+ | 1.50e+1(2.59e+0)+ | 1.01e+1(5.52e+0)+ |
| 40 | 1.28e+1(2.03e+0)+ | 1.48e+1(7.35e+0)+ | 1.67e+1(3.73e+0)+ | 4.13e+0(7.90e-1) \approx | 4.00e+0(8.59e-1) | 4.92e+0(1.09e+0)+ | 1.45e+1(1.85e+0)+ | 8.01e+0(3.35e+0)+ |
| + / \approx / - | 6/0/0 | 6/0/0 | 6/0/0 | 0/6/0 | -/- | 2/4/0 | 4/2/0 | 4/2/0 |

598 D.2.1 Experimental setups

599 The related tasks T_i used in our experiment are $N = 100$ gasoline engine models. These engine
600 models have different behaviors when compared with T_* , but they share the basic features of gasoline
601 engines. All related tasks and the target task have the same six decision variables. Each related task
602 T_i provides only 60 solutions, forming a dataset D_i . The size of datasets used for meta-learning
603 $|D_m|$ is set to 40. Six tests are conducted where $|S_*| = \{2, 3, 5, 10, 20, 40\}$ data points are sampled
604 from the real engine simulation T_* . MSE is chosen as an indicator of modeling accuracy, which is
605 measured using a dataset consisting of 12500 data points that are sampled uniformly from the engine
606 decision space. The comparison algorithms are the same as described in Appendix D.1.

607 D.2.2 Results and analysis

608 The statistical test results of the MSE values achieved by comparison algorithms in BSFC regression
609 experiments are summarized in Table 4. Each row lists the results obtained when the same number
610 of fitness evaluations $|S_*|$ are used to train models. The results of Wilcoxon rank sum test between
611 MDKL and other compared algorithms are listed in the last row. It can be observed that MDKL and
612 MDKL_NN outperform other comparison modeling methods since they have achieved the smallest
613 MSE on all tests.

614 Additional Wilcoxon rank sum tests have been conducted between MDKL related algorithms to
615 answer our second question (results are not reported in Table 4). The statistical test results between
616 DKL and GP_Adam are 1/5/0, indicating that the neural network in DKL makes some contributions
617 to the performance of MDKL. The statistical test results between MDKL_NN and DKL are 6/0/0,
618 demonstrating that the meta-learning of neural network parameters constructs a useful deep kernel
619 and contributes to the improvement of modeling accuracy. However, there is no significant difference
620 between the performance of MDKL and that of MDKL_NN, the meta-learning on base kernel
621 parameters does not play a critical role on this engine problem. In comparison, the meta-learning on
622 base kernel parameters is effective in sinusoid function regression experiments (see Appendix D.1).
623 In addition, the statistical test results between MDKL_NN and MAML are 4/2/0. Considering that
624 MAML is a neural network regressor learned through meta-learning, we can conclude that GP is an
625 essential component of our MDKL. In summary, all components in MDKL are necessary, they all
626 contribute to the effectiveness of learning experience.

627 The comparison experiments on sinusoid functions and the gasoline motor engine have demonstrated
628 the effectiveness of our MDKL modeling method in the learning of experience. Given a small
629 dataset of the target task, the model learned through MDKL method has the smallest MSE among
630 all comparison models. Additionally, the investigation between MDKL and its variants shows that
631 all components in MDKL have made their contributions to the effectiveness of learning experience.
632 However, similar to other meta-learning studies [10, 13], we have not defined the similarity between
633 tasks. In other words, the boundary between related tasks and unrelated tasks has not been defined.
634 This should be a topic of further study on meta-learning. Moreover, the relationship between task
635 similarity and modeling performance has not been investigated. Instead, we study the relationship
636 between task similarity and SAEA optimization performance in Section F, since our main focus is
637 the surrogate-assisted evolutionary optimization.

Table 5: Mean IGD+ values and standard deviation (in parentheses) obtained from 30 runs on the DTLZ problems. MOEA/D-FS and the comparison algorithms initialize their surrogates with 10, 100 samples, respectively. Extra 50 evaluations are allowed in the further optimization.

| Problem | MOEA/D-EGO | MOEA/D-FS (ours) | ParEGO | K-RVEA | KTA2 | CSEA | OREA | ESBCEO | KMOEA-TIC |
|-----------|-------------------|------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| DTLZ1 | 1.07e+2(2.05e+1)+ | 9.70e+1(1.87e+1) | 7.82e+1(1.54e+1)- | 1.18e+2(2.45e+1)+ | 1.01e+2(2.38e+1)≈ | 1.10e+2(2.50e+1)+ | 1.02e+2(1.97e+1)≈ | 8.81e+1(1.18e+1)≈ | 1.10e+2(2.29e+1)+ |
| DTLZ2 | 2.99e-1(7.01e-2)+ | 1.43e-1(2.29e-2) | 3.17e-1(4.12e-2)+ | 2.69e-1(5.97e-2)+ | 2.14e-1(3.84e-2)+ | 2.98e-1(5.25e-2)+ | 1.76e-1(4.69e-2)+ | 3.39e-1(3.78e-2)+ | 2.10e-1(7.10e-2)+ |
| DTLZ3 | 3.15e+2(6.04e+1)+ | 1.97e+2(1.64e+1) | 2.30e+2(5.99e+1)≈ | 3.24e+2(5.90e+1)+ | 2.67e+2(6.70e+1)+ | 2.82e+2(6.97e+1)+ | 2.72e+2(6.88e+1)+ | 2.09e+2(4.23e+1)≈ | 2.98e+2(6.14e+1)+ |
| DTLZ4 | 5.04e-1(8.25e-2)≈ | 4.44e-1(1.35e-1) | 5.44e-1(7.58e-2)+ | 4.57e-1(1.14e-1)≈ | 4.51e-1(9.54e-2)≈ | 4.75e-1(1.09e-1)≈ | 3.18e-1(1.54e-1)- | 4.99e-1(7.37e-2)≈ | 4.26e-1(9.19e-2)≈ |
| DTLZ5 | 2.39e-1(7.17e-2)+ | 1.13e-1(2.24e-2) | 2.58e-1(3.68e-2)+ | 1.92e-1(5.97e-2)+ | 1.44e-1(4.60e-2)+ | 2.14e-1(4.05e-2)+ | 7.84e-2(2.42e-2)- | 2.68e-1(3.62e-2)+ | 8.73e-2(2.77e-2)- |
| DTLZ6 | 1.29e+0(4.74e-1)≈ | 1.11e+0(5.71e-1) | 1.67e+0(6.77e-1)+ | 4.62e+0(6.42e-1)+ | 3.37e+0(6.71e-1)+ | 6.26e+0(3.40e-1)+ | 4.60e+0(1.19e+0)+ | 2.41e+0(7.97e-1)+ | 2.90e+0(5.34e-1)+ |
| DTLZ7 | 3.31e-1(3.11e-1)- | 2.47e+0(1.89e+0) | 3.66e-1(1.31e-1)- | 1.74e-1(3.57e-2)- | 4.34e-1(2.20e-1)- | 4.17e+0(1.13e+0)+ | 2.14e+0(1.15e+0)≈ | 5.47e-1(2.46e-1)- | 9.44e-2(1.23e-2)- |
| + / ≈ / - | 4/2/1 | -/- | 4/1/2 | 5/1/1 | 4/2/1 | 6/1/0 | 3/2/2 | 3/3/1 | 4/1/2 |

638 E Additional Details on Expensive Multi-Objective Optimization

639 E.1 Comparison algorithms

640 As explained in Section B, our FSEO framework is compatible with regression-based SAEAs. Hence,
641 we select MOEA/D-EGO [51] as an example and replace its GP surrogates by our MDKL surrogates.
642 The resulting algorithm is denoted as MOEA/D-FS. Note that it is not necessary to specially select a
643 newly proposed regression-based SAEA as our example, our main objective is to save evaluations
644 with experience and observe if there is any damage to the optimization performance caused by the
645 saving of evaluations. Therefore, it does not make any difference which regression-based SAEA
646 or BO we choose as our example. Additionally, to demonstrate the improvement of optimization
647 performance caused by using experience on DTLZ problems is significant, several state-of-the-art
648 SAEAs and MOBO are also compared as baselines, including ParEGO [19], K-RVEA [7], CSEA [25],
649 OREA [48], KTA2 [33], ESBCEO [3], and KMOEA-TIC [28]. Among these algorithms, ParEGO,
650 K-RVEA, KTA2, KMOEA-TIC use regression-based surrogates, CSEA uses a classification-based
651 surrogate, OREA employs an ordinal-regression-based surrogate, and ESBCEO is a recently proposed
652 MOBO.

653 We implemented the FSEO framework, MOEA/D-EGO, ParEGO, and OREA, while the code of
654 K-RVEA, CSEA, KTA2, and ESBCEO [3] is available on PlatEMO [37], an open source MATLAB
655 platform for evolutionary multi-objective optimization. The code of KMOEA-TIC [28] is obtained
656 from its authors. To make a fair comparison, all comparison algorithms share the same initial dataset
657 S_* in an independent run. We also set $\theta \in [10^{-5}, 100]^d$ and $\mathbf{p} = \mathbf{2}$ for all GP surrogates as suggested
658 in [33], these GP surrogates are implemented through DACE [31]. Other configurations are the same
659 as suggested in their original literature.

660 E.2 Result Table and Analysis of Expensive Multi-Objective Optimization

661 The experience learned from related tasks makes MOEA/D-EGO more competitive when compared
662 to other SAEAs. The use of MDKL surrogates results in significantly smaller IGD+ values on DTLZ1,
663 DTLZ2, DTLZ3, and DTLZ5 than before. As a result, MOEA/D-FS achieves the smallest IGD+
664 values on DTLZ2 and DTLZ3, and its optimization results on DTLZ1 and DTLZ5 are much closer
665 to the best optimization results (e.g. results obtained by ParEGO and OREA) than MOEA/D-EGO.
666 Although MOEA/D-FS does not achieve the smallest IGD+ values on all DTLZ problems, it should
667 be noted that MOEA/D-FS is still the best algorithm among comparison SAEAs due to its overall
668 performance. Table 5 shows that no comparison SAEA outperforms MOEA/D-FS on three DTLZ
669 problems, but MOEA/D-FS outperforms all comparison SAEAs on at least three DTLZ problems.
670 Furthermore, the IGD+ values of MOEA/D-FS are achieved with an evaluation budget of 60, while
671 the IGD+ values of other SAEAs are reached with a cost of 150 evaluations (see Table 2).

672 E.3 Result Tables and Figures in IGD and HV Metrics

673 The performance of our method and the comparison algorithms are also evaluated on inverted
674 generational distance (IGD) [4] and Hypervolume (HV) [55] metrics.

675 Results in IGD values are reported in Table 6 and Fig. 5. A smaller IGD value indicates a better
676 optimization result.

677 Results in HV values are reported in Table 7 and Fig. 6. A larger HV value indicates a better
678 optimization result.

Table 6: Mean IGD values and standard deviation (in parentheses) obtained from 30 runs on 7 DTLZ problems. MOEA/D-FS and comparison algorithms initialize their surrogates with 10, 100 samples, respectively. Extra 50 evaluations are allowed in the further optimization. ‘+’, ‘ \approx ’, and ‘-’ denote MOEA/D-FS is statistically significantly superior to, almost equivalent to, and inferior to the compared algorithms in the Wilcoxon rank sum test (significance level is 0.05), respectively. The last row counts the total win/tie/loss results.

| Problems | MOEA-EGO | MOEA-FS | ParEGO | K-RVEA | KTA2 | CSEA | OREA | ESBCEO | KMOEAIC |
|-------------------|----------------------------|------------------|----------------------------|-------------------|----------------------------|-------------------|----------------------------|----------------------------|----------------------------|
| DTLZ1 | 7.07e+2(2.05e+1)+ | 9.70e+1(1.87e+1) | 7.82e+1(1.54e+1)- | 1.18e+2(2.41e+1)+ | 1.01e+2(2.34e+1) \approx | 1.10e+2(2.46e+1)+ | 1.02e+2(1.97e+1) \approx | 8.81e+1(1.18e+1) \approx | 1.10e+2(2.29e+1)+ |
| DTLZ2 | 3.30e+1(7.23e-2)+ | 1.72e-1(2.41e-2) | 3.59e-1(2.82e-2)+ | 3.08e-1(4.93e-2)+ | 2.45e-1(3.57e-2)+ | 3.36e-1(3.96e-2)+ | 2.14e-1(4.10e-2)+ | 3.64e-1(3.29e-2)+ | 2.86e-1(6.31e-2)+ |
| DTLZ3 | 3.15e+2(6.04e+1)+ | 1.97e+2(1.64e+1) | 3.24e+2(5.99e+1) \approx | 3.24e+2(5.80e+1)+ | 2.67e+2(6.58e+1)+ | 2.82e+2(6.85e+1)+ | 2.72e+2(6.88e+1)+ | 2.09e+2(4.23e+1) \approx | 2.98e+2(6.14e+1)+ |
| DTLZ4 | 7.51e-1(1.50e-1) \approx | 7.96e-1(2.25e-1) | 7.65e-1(1.14e-1) \approx | 5.94e-1(1.28e-1)- | 6.30e-1(1.51e-1)- | 7.00e-1(1.48e-1)- | 5.64e-1(2.01e-1)- | 6.70e-1(8.05e-2)- | 5.23e-1(8.60e-2)- |
| DTLZ5 | 2.47e-1(7.21e-2)+ | 1.17e-1(2.08e-2) | 2.83e-1(3.13e-2)+ | 2.13e-1(5.55e-2)+ | 1.61e-1(4.60e-2)+ | 2.33e-1(3.65e-2)+ | 8.64e-2(2.48e-2)- | 2.83e-1(3.00e-2)+ | 1.18e-1(3.17e-2) \approx |
| DTLZ6 | 1.36e+0(4.10e-1) \approx | 1.18e+0(5.35e-1) | 1.78e+0(6.29e-1)+ | 4.63e+0(6.26e-1)+ | 3.37e+0(6.50e-1)+ | 6.26e+0(3.28e-1)+ | 4.61e+0(1.18e+0)+ | 2.45e+0(7.92e-1)+ | 2.92e+0(5.35e-1)+ |
| DTLZ7 | 4.22e-1(3.16e-1)- | 2.56e+0(1.86e+0) | 5.34e-1(1.25e-1)- | 2.55e-1(4.36e-2)- | 5.54e-1(2.38e-1)- | 4.20e+0(1.11e+0)+ | 2.21e+0(1.11e+0) \approx | 6.21e-1(2.43e-1)- | 1.85e-1(1.81e-2)- |
| + / \approx / - | 4/2/1 | -/- | 3/2/2 | 5/0/2 | 4/1/2 | 6/0/1 | 3/2/2 | 3/2/2 | 4/1/2 |

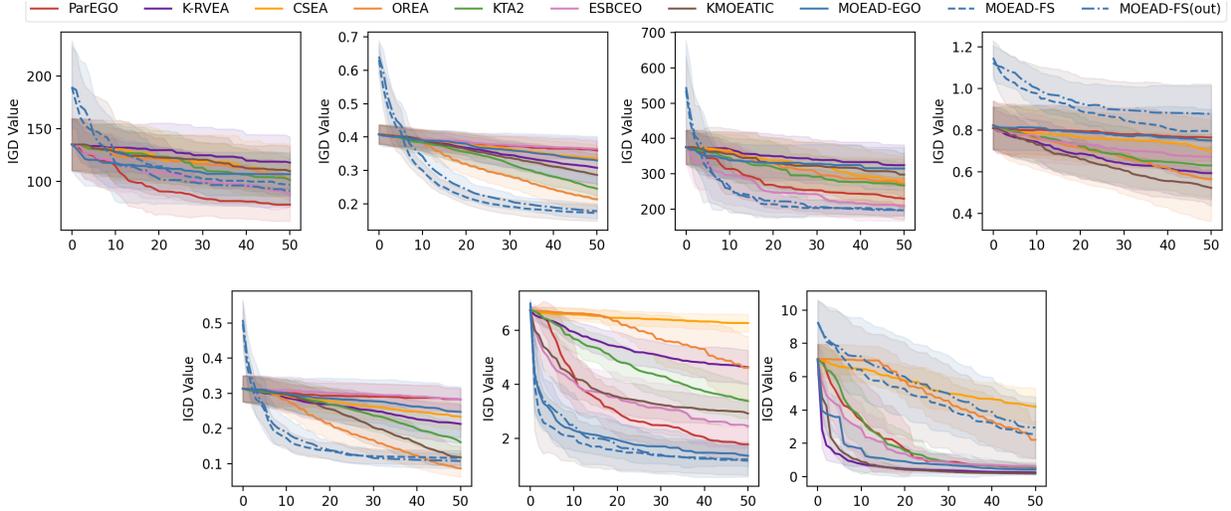


Figure 5: IGD curves averaged over 30 runs on 7 DTLZ problems. Solid lines are mean values, while shadows are error regions. **Upper:** DTLZ1, DTLZ2, DTLZ3, DTLZ4. **Lower:** DTLZ5, DTLZ6, DTLZ7. MOEA/D-FSs and comparison algorithms initialize their surrogates with 10, 100 samples, respectively. Extra 50 evaluations are allowed in the further optimization. Note that ‘FS(out)’ indicates the target task is excluded from the range of related tasks during the meta-learning procedure). X-axis denotes the number of evaluations used after the surrogate initialization.

Table 7: Mean HV values and standard deviation (in parentheses) obtained from 30 runs on 7 DTLZ problems. MOEA/D-FS and comparison algorithms initialize their surrogates with 10, 100 samples, respectively. Extra 50 evaluations are allowed in the further optimization. ‘+’, ‘ \approx ’, and ‘-’ denote MOEA/D-FS is statistically significantly superior to, almost equivalent to, and inferior to the compared algorithms in the Wilcoxon rank sum test (significance level is 0.05), respectively. The last row counts the total win/tie/loss results.

| Problems | MOEA-EGO | MOEA-FS | ParEGO | K-RVEA | KTA2 | CSEA | OREA | ESBCEO | KMOEAIC |
|-------------------|----------------------------|------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|----------------------------|
| DTLZ1 | 0.00e+0(0.00e+0) \approx | 0.00e+0(0.00e+0) | 0.00e+0(0.00e+0) \approx |
| DTLZ2 | 2.02e-1(1.28e-1)+ | 4.69e-1(3.70e-2) | 1.21e-1(4.31e-2)+ | 1.93e-1(8.93e-2)+ | 3.19e-1(6.49e-2)+ | 1.59e-1(5.39e-2)+ | 3.80e-1(7.64e-2)+ | 1.39e-1(4.55e-2)+ | 2.91e-1(1.29e-1)+ |
| DTLZ3 | 0.00e+0(0.00e+0) \approx | 0.00e+0(0.00e+0) | 0.00e+0(0.00e+0) \approx |
| DTLZ4 | 6.25e-2(5.53e-2)+ | 1.43e-1(7.17e-2) | 2.03e-2(2.71e-2)+ | 3.81e-2(4.24e-2)+ | 6.49e-2(7.42e-2)+ | 4.30e-2(5.29e-2)+ | 2.11e-1(1.37e-1) \approx | 2.27e-2(2.65e-2)+ | 6.50e-2(7.07e-2)+ |
| DTLZ5 | 4.50e-2(4.17e-2)+ | 1.63e-1(1.60e-2) | 1.29e-2(1.30e-2)+ | 4.82e-2(2.78e-2)+ | 7.98e-2(3.80e-2)+ | 3.08e-2(1.61e-2)+ | 1.49e-1(2.88e-2) \approx | 1.64e-2(1.42e-2)+ | 1.58e-1(3.69e-2) \approx |
| DTLZ6 | 1.24e-3(3.77e-3) \approx | 1.59e-2(3.46e-2) | 2.02e-5(1.09e-4) \approx | 0.00e+0(0.00e+0) \approx |
| DTLZ7 | 3.83e-1(8.50e-2)- | 1.12e-1(1.27e-1) | 2.21e-1(9.60e-2)- | 3.79e-1(2.61e-2)- | 3.70e-1(3.88e-2)- | 3.14e-1(5.69e-4)+ | 2.97e-2(4.45e-2) \approx | 1.71e-1(8.33e-2)- | 4.67e-1(1.27e-2)- |
| + / \approx / - | 3/3/1 | -/- | 3/3/1 | 3/3/1 | 3/3/1 | 4/3/0 | 1/6/0 | 3/3/1 | 2/4/1 |

679 E.4 Performance on Expensive Multi-Objective Optimization Under the Same Evaluation Budget

680
681 The statistical test results reported in the last row of Table 5 show that ParEGO [19] and OREA [48]
682 are the best two comparison algorithms when compared with our MOEA/D-FS. In this subsection,
683 we evaluate the performance of MOEA/D-FS when no extra evaluation is saved. For this purpose, we
684 compare the optimization performance of these three SAEAs under the same evaluation budget: 10
685 evaluations (1d) for surrogate initialization and 50 evaluations for further optimization. The statistical
686 test results are reported in Table 8. It can be seen that our MOEA/D-FS generally outperforms the

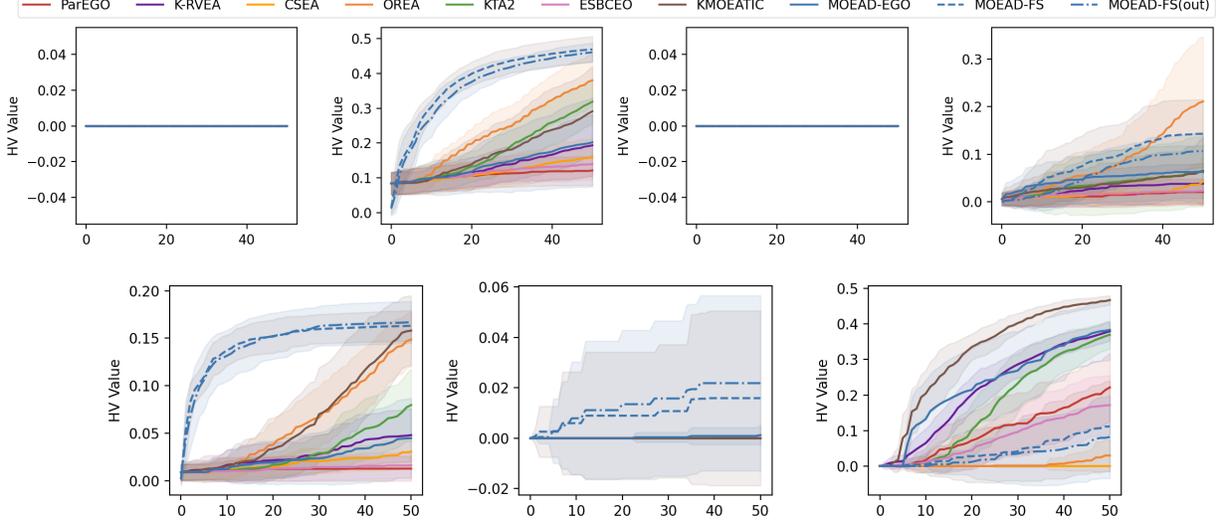


Figure 6: HV curves averaged over 30 runs on 7 DTLZ problems. Solid lines are mean values, while shadows are error regions. **Upper:** DTLZ1, DTLZ2, DTLZ3, DTLZ4. **Lower:** DTLZ5, DTLZ6, DTLZ7. MOEA/D-FSs and comparison algorithms initialize their surrogates with 10, 100 samples, respectively. Extra 50 evaluations are allowed in the further optimization. Note that ‘FS(out)’ indicates the target task is excluded from the range of related tasks during the meta-learning procedure). X-axis denotes the number of evaluations used after the surrogate initialization.

Table 8: Mean IGD+ values and standard deviation (in parentheses) obtained from 30 runs on DTLZ problems. MOEA/D-FS is compared with ParEGO and OREA under the same evaluation budget: 10 evaluations for surrogate initialization and 50 evaluations for the optimization process. ‘+’, ‘≈’, and ‘-’ denote MOEA/D-FS is statistically significantly superior to, almost equivalent to, and inferior to the compared two algorithms in the Wilcoxon rank sum test (significance level is 0.05), respectively. The last row counts the total win/tie/loss results.

| Problem | MOEA/D-FS | ParEGO | OREA |
|-----------|------------------|-------------------|-------------------|
| DTLZ1 | 9.70e+1(1.87e+1) | 6.70e+1(4.75e+0)- | 1.10e+2(3.65e+1)≈ |
| DTLZ2 | 1.43e-1(2.29e-2) | 5.51e-1(5.37e-2)+ | 4.28e-1(6.68e-2)+ |
| DTLZ3 | 1.97e+2(1.64e+1) | 1.84e+2(8.86e+0)≈ | 2.72e+2(6.59e+1)+ |
| DTLZ4 | 4.44e-1(1.35e-1) | 6.29e-1(7.99e-2)+ | 6.45e-1(1.24e-1)+ |
| DTLZ5 | 1.13e-1(2.24e-2) | 4.32e-1(8.88e-2)+ | 3.02e-1(7.63e-2)+ |
| DTLZ6 | 1.11e+0(5.71e-1) | 1.03e+0(4.78e-1)≈ | 5.71e+0(6.73e-1)+ |
| DTLZ7 | 2.47e+0(1.89e+0) | 4.38e-1(1.39e-1)- | 7.12e+0(1.77e+0)+ |
| + / ≈ / - | - / - / - | 3 / 2 / 2 | 6 / 1 / 0 |

687 compared SAEAs when only $1d$ evaluations are used to initialize their surrogates. The effectiveness
688 of our FSEO framework has been demonstrated. Note that OREA is an evolutionary algorithm
689 assisted by ordinal-regression-based surrogates. Currently, our FSEO framework is applicable to the
690 SAEAs working with regression-based surrogates. The meta-learning of ordinal-regression models
691 can be a topic of further research.

692 F Influence of Task Similarity

693 In real-world applications, it is optimistic to assume that some related tasks are very similar to the
694 target task. A more common situation is that all related tasks have limited similarity to the target task.
695 To investigate the relationship between task similarity and FSEO optimization performance, we also
696 test the performance in an ‘out-of-range’ situation, where the original DTLZ is excluded from the
697 range of DTLZ variants during the MDKL meta-learning procedure. As a result, only the DTLZ
698 variants that are quite different from the original DTLZ problem can be used to learn experience. The
699 ‘out-of-range’ situation eliminates the probability that MDKL surrogates benefit greatly from the

Table 9: Mean IGD+ values and standard deviation (in parentheses) obtained from 30 runs on DTLZ problems. Both MOEA/D-FSs initialize their surrogates with 10 samples, extra 50 evaluations are allowed in the further optimization. The last two rows count the statistical test results between MOEA/D-FSs and other compared algorithms.

| MOEA/D-FSs | In-range | Out-of-range |
|------------------|--------------------|------------------|
| DTLZ1 | 9.70e+1(1.87e+1)≈ | 9.11e+1(1.53e+1) |
| DTLZ2 | 1.43e-1(2.29e-2)≈ | 1.41e-1(1.75e-2) |
| DTLZ3 | 1.97e+2 (1.64e+1)≈ | 1.98e+1(1.51e+1) |
| DTLZ4 | 4.44e-1(1.35e-1)≈ | 4.96e-1(8.63e-2) |
| DTLZ5 | 1.13e-1(2.24e-2)≈ | 1.03e-1(2.39e-2) |
| DTLZ6 | 1.11e+0(5.71e-1)≈ | 1.17e+0(6.88e-1) |
| DTLZ7 | 2.47e+0(1.89e+0)≈ | 2.86e+0(1.87e+0) |
| + / ≈ / - | 0/7/0 | -/-/- |
| vs MOEA/D-EGO | 4/2/1 | 4/2/1 |
| vs 6 Comparisons | 26/9/7 | 27/7/8 |

700 DTLZ variants that are very similar to the original DTLZ problem. Detailed definitions of the related
701 tasks used in the ‘out-of-range’ situation are given in Appendix C. Apart from the related tasks used,
702 the remaining experimental setups are the same as the setups described in Section 5.1. For the sake of
703 convenience, we denote the situation we tested in Section 5.1 as ‘in-range’ below.

704 The statistical test results reported in Table 9 show that the ‘out-of-range’ situation achieves competi-
705 tive IGD+ values to the ‘in-range’ situation on all 7 test instances. This suggests that related tasks
706 that are very similar to the target task have a limited impact on the optimization performance of our
707 FSEO framework. Useful experience can be learned from related tasks that are not very similar to
708 the target task. Crucially, when comparing the performance of the ‘out-of-range’ situation and that
709 of MOEA/D-EGO, we can still observe competitive or improved optimization results on 6 DTLZ
710 problems (see Table 9, the row titled by ‘vs MOEA/D-EGO’, or Fig. 3). Moreover, it can be seen
711 from the last row of Table 9 that the ‘out-of-range’ situation achieves better/competitive/worse IGD+
712 values than all compared SAEAs on 27/7/8 test instances. In comparison, the corresponding statistical
713 test results for the ‘in-range’ situation are 26/9/7. The difference between these statistical test results
714 is not significant.

715 A study on the ‘out-of-range’ situation in the context of extremely expensive multi-objective opti-
716 mization is presented in Appendix H.2. Consistent results can be observed from Table 12 and Fig.
717 7.

718 Consequently, related tasks that are very similar to the target task are not essential to the optimization
719 performance of our FSEO framework. In the ‘out-of-range’ situation, our MOEA/D-FS can still
720 achieve competitive or better optimization results than MOEA/D-EGO while using only $1d$ samples
721 for surrogate initialization.

722 G Influence of the Size of Datasets Used in Meta-Learning

723 We also investigated the performance of our FSEO framework when different sizes of datasets $|D_m|$
724 are used in the meta-learning procedure. The experimental setups are the same as the setups of
725 MOEA/D-FS in Section 5.1 except for $|D_m|$.

726 It is evident from Table 10 that when each DTLZ variant provides $|D_m| = 60$ samples for the
727 meta-learning of MDKL surrogates, the performance of both MOEA/D-FSs are improved on 2 or
728 3 DTLZ problems. Particularly, a significant improvement can be observed from the optimization
729 results of DTLZ7. As we discussed in Section 5.1, the poor performance of our experience-based
730 optimization on DTLZ7 is caused by the small size of D_m . Optimal solutions have few chances to
731 be included in a small D_m , which makes D_m fails to provide the experience about the discontinuity
732 of optimal regions. In comparison, the experience of ‘optimal regions’ can be learned from large
733 datasets D_m and thus the optimization results are improved significantly.

734 In conclusion, for our FSEO framework, a large D_m for the meta-learning procedure indicates
735 more useful experience can be learned from related tasks, which further improves the performance
736 of experience-based optimization. Therefore, when applying our FSEO framework to real-world

Table 10: Mean IGD+ values and standard deviation (in parentheses) obtained from 30 runs on DTLZ problems. 10 samples are used for initialization and extra 50 evaluations are allowed in the further optimization. $|D_m|$ is the size of the dataset collected from each related task.

| Problem | In-range | | Out-of-range | |
|-----------|-------------------|------------------|-------------------|------------------|
| | $ D_m =20$ | $ D_m =60$ | $ D_m =20$ | $ D_m =60$ |
| DTLZ1 | 9.70e+1(1.87e+1)≈ | 9.77e+1(1.73e+1) | 9.11e+1(1.53e+1)≈ | 9.93e+1(1.87e+1) |
| DTLZ2 | 1.43e-1(2.29e-2)+ | 1.24e-1(2.11e-2) | 1.41e-1(1.75e-2)+ | 1.29e-1(2.36e-2) |
| DTLZ3 | 1.97e+2(1.64e+1)≈ | 1.98e+2(2.21e+1) | 1.98e+1(1.51e+1)≈ | 1.93e+2(1.19e+1) |
| DTLZ4 | 4.44e-1(1.35e-1)≈ | 5.17e-1(5.68e-2) | 4.96e-1(8.63e-2)≈ | 5.17e-1(5.38e-2) |
| DTLZ5 | 1.13e-1(2.24e-2)+ | 9.96e-2(2.18e-2) | 1.03e-1(2.39e-2)≈ | 1.05e-1(2.73e-2) |
| DTLZ6 | 1.11e+0(5.71e-1)≈ | 1.04e+0(6.06e-1) | 1.17e+0(6.88e-1)≈ | 1.22e+0(6.41e-1) |
| DTLZ7 | 2.47e+0(1.89e+0)+ | 7.49e-1(2.61e-1) | 2.86e+0(1.87e+0)+ | 6.96e-1(2.41e-1) |
| + / ≈ / - | 3/4/0 | -/- | 2/5/0 | -/- |

Table 11: Mean IGD+ values and standard deviation (in parentheses) obtained from 30 runs on the DTLZ problems. MOEA/D-FS and the comparison algorithms initialize their surrogates with 10, 60 samples, respectively. Extra 30 evaluations are allowed in the further optimization. ‘+’, ‘≈’, and ‘-’ denote MOEA/D-FS is statistically significantly superior to, equivalent to, and inferior to the compared algorithms in the Wilcoxon rank sum test (significance level is 0.05), respectively. The last row is the total win/tie/loss results. Performance improvement can be observed from the comparisons between MOEA/D-FS and MOEA/D-EGO, while 50 evaluations are saved from surrogate initialization.

| Problems | MOEA-D-EGO | MOEA-D-FS (ours) | ParEGO | K-RVEA | RTA2 | CSEA | OREA | ESBCEO | KMOEATIC |
|-----------|-------------------|------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| DTLZ1 | 1.07e+2(2.73e+1)≈ | 1.03e+2(2.34e+1) | 8.70e+1(2.53e+1)- | 1.22e+2(3.20e+1)+ | 1.15e+2(3.03e+1)≈ | 1.08e+2(2.64e+1)≈ | 1.11e+2(2.25e+1)+ | 1.00e+2(2.07e+1)≈ | 1.20e+2(2.71e+1)+ |
| DTLZ2 | 3.49e-1(5.82e-2)+ | 1.57e-1(2.29e-2) | 3.51e-1(5.01e-2)+ | 3.72e-1(4.32e-2)+ | 3.57e-1(4.60e-2)+ | 3.55e-1(5.14e-2)+ | 3.14e-1(3.76e-2)+ | 3.83e-1(3.83e-2)+ | 3.79e-1(4.46e-2)+ |
| DTLZ3 | 3.07e+2(5.32e+1)+ | 2.03e+2(2.42e+1) | 2.16e+2(4.89e+1)≈ | 3.53e+2(7.76e+1)+ | 3.23e+2(8.67e+1)+ | 3.35e+2(6.83e+1)+ | 3.39e+2(7.72e+1)+ | 2.41e+2(5.51e+1)+ | 3.27e+2(8.10e+1)+ |
| DTLZ4 | 5.45e-1(1.09e-1)≈ | 4.91e-1(1.24e-1) | 6.36e-1(8.67e-2)+ | 5.53e-1(9.79e-2)≈ | 5.47e-1(1.02e-1)≈ | 5.84e-1(9.59e-2)+ | 5.14e-1(1.21e-1)≈ | 5.47e-1(7.55e-2)≈ | 4.53e-1(1.03e-1)≈ |
| DTLZ5 | 2.79e-1(5.69e-2)+ | 1.18e-1(2.25e-2) | 2.78e-1(5.59e-2)+ | 2.82e-1(5.42e-2)+ | 2.60e-1(5.50e-2)+ | 2.77e-1(4.34e-2)+ | 1.99e-1(4.53e-2)+ | 2.94e-1(4.92e-2)+ | 2.69e-1(6.11e-2)+ |
| DTLZ6 | 2.04e+0(7.33e-1)+ | 1.29e+0(6.44e-1) | 2.47e+0(7.39e-1)+ | 5.23e+0(6.17e-1)+ | 4.58e+0(6.36e-1)+ | 6.44e+0(3.53e-1)+ | 5.79e+0(6.70e-1)+ | 3.04e+0(9.46e-1)+ | 3.55e+0(6.90e-1)+ |
| DTLZ7 | 1.90e+0(9.19e-1)- | 4.16e+0(2.54e+0) | 1.39e+0(1.49e+0)- | 3.13e-1(6.07e-2)- | 2.05e+0(2.16e+0)- | 5.47e+0(1.31e+0)+ | 5.51e+0(1.32e+0)+ | 9.57e-1(5.40e-1)- | 2.68e-1(1.47e-1)- |
| + / ≈ / - | 4/2/1 | -/- | 4/1/2 | 5/1/1 | 4/2/1 | 6/1/0 | 6/1/0 | 4/2/1 | 5/1/1 |

737 optimization problems, it is preferable to collect more data from related tasks for experience learning.
738

739 H Experiments on Extremely Expensive Multi-Objective Optimization

740 In this section, we investigate the performance of our FSEO framework in the context of extremely
741 expensive optimization, where the allowed fitness evaluations on target problems are fewer than that
742 in the experiment carried out in Sections 5.1 of the main file and Appendix F.

743 H.1 Performance between Comparison Algorithms

744 We conduct the experiment described in Section 5.1 of the main file, but with a smaller evaluation
745 budget than the budget listed in Table 2. The size of the initial dataset S_* is set to 10, 60 for our
746 MOEA/D-FS and comparison algorithms, respectively. 30 extra evaluations for further optimization
747 are allowed. The total evaluation budget is 40, 90 for our MOEA/D-FS and comparison algorithms,
748 respectively.

749 The aim of this subsection is to answer the question below:

- 750 • Is our FSEO framework more suitable for the optimization problems in which evaluations are
751 extremely expensive? In other words, will the advantage of our FSEO framework become
752 more prominent if the optimization problems allow a smaller evaluation budget?

753 The comparison results reported in Fig. 7 and Table 11 show that MOEA/D-FS has achieved
754 competitive or smaller IGD+ values than MOEA/D-EGO on all DTLZ problems except for DTLZ7.
755 Meanwhile, 5*d* evaluations have been saved.

756 Consistent with the results discussed in Section 5.1 of the main file, MOEA/D-FS fails to achieve a
757 competitive result compared to MOEA/D-EGO on DTLZ7 since experience is learned from small
758 datasets collected from related tasks. Although we set a different evaluation budget for all SAEAs,
759 the size of datasets for meta-learning $|D_m|$ has not been modified. However, it can be observed from
760 the statistical test results (see the last row of Tables 5 and 11) that our MOEA/D-FS outperforms
761 the comparison algorithms on 26, 29 test instances when the total evaluation budget of comparison

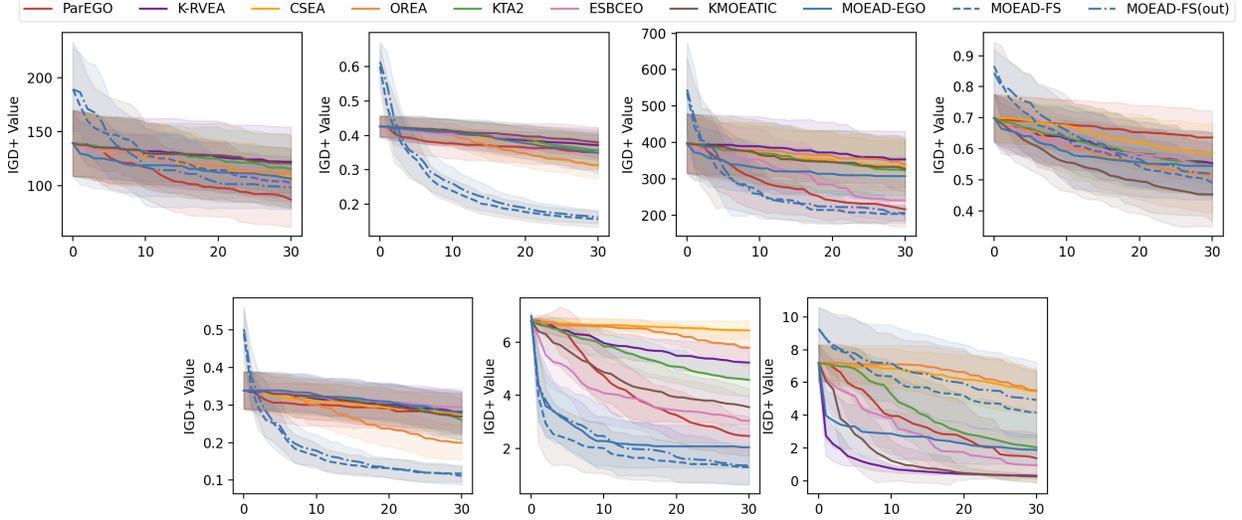


Figure 7: IGD+ curves averaged over 30 runs on 7 DTLZ problems. Solid lines are mean values, while shadows are error regions. **Upper:** DTLZ1, DTLZ2, DTLZ3, DTLZ4. **Lower:** DTLZ5, DTLZ6, DTLZ7. MOEA/D-FSs and comparison algorithms initialize their surrogates with 10, 60 samples, respectively. Extra 30 evaluations are allowed in the further optimization. Note that ‘FS(out)’ indicates the target task is excluded from the range of related tasks during the meta-learning procedure. X-axis denotes the number of evaluations used after the surrogate initialization. In comparison to MOEA/D-EGO, both MOEA/D-FSs achieve smaller or competitive IGD+ values on all DTLZ test problems except for DTLZ7, while 50 evaluations are saved with the assistance from related tasks. Moreover, MOEA/D-FSs achieve the smallest IGD+ values on DTLZ2, DTLZ3, DTLZ4, DTLZ5 and DTLZ6.

762 algorithms is set to 150, 90, respectively. This answers the question we raised before: The advantage
 763 of our FSEO framework is more prominent in the extremely expensive problems where a smaller
 764 evaluation budget is allowed. The comparison between the results obtained from Tables 5 and 11 has
 765 demonstrated that our FSEO framework is preferable when solving optimization problems within a
 766 very limited evaluation budget.

767 H.2 Out-Of-Range Analysis on Extremely Expensive Optimization

768 In Section F of the main file, we carried out an experiment to study the influence of task similarity
 769 on the performance of experience-based expensive multi-objective optimization. The optimization
 770 results obtained from the ‘in-range’ and the ‘out-of-range’ situations are compared. In this subsection,
 771 we conduct an experiment to investigate the difference between the ‘in-range’ and the ‘out-of-range’
 772 situations for extremely expensive multi-objective optimization. The experimental setups are the
 773 same as the setups described in Section F of the main file, except the allowed fitness evaluation budget
 774 is the same as described in Appendix H.1.

775 Table 12 gives the statistical test results, it can be seen that the ‘out-of-range’ situation achieves
 776 competitive IGD+ values to the ‘in-range’ situation on all 7 test instances. In comparison to MOEA/D-
 777 EGO, the experience gained in the ‘out-of-range’ situation leads to competitive or smaller IGD+
 778 values on 6 DTLZ problems. Furthermore, similar results can be observed in the last row of Table 12,
 779 the ‘out-of-range’ situation achieves better/competitive/worse IGD+ values than all compared SAEAs
 780 on 28/9/5 test instances. In comparison, the ‘in-range’ situation achieves better/competitive/worse
 781 IGD+ values than all compared SAEAs on 29/8/5 test instances. There is only a minor difference
 782 between the optimization results obtained in two situations. These observations are consistent with
 783 the conclusions we made in Section F of the main file.

Table 12: Mean IGD+ values and standard deviation (in parentheses) obtained from 30 runs on DTLZ problems. ‘Out-of-range’ indicates the target task is excluded from the range of related tasks during the meta-learning procedure. Both MOEA/D-FSs initialize their surrogates with 10 samples, extra 30 evaluations are allowed in the further optimization. ‘+’, ‘≈’, and ‘-’ denote the result of the ‘out-of-range’ situation is statistically significantly superior to, almost equivalent to, and inferior to that of the ‘in-range’ situation in the Wilcoxon rank sum test (significance level is 0.05), respectively. The last two rows count the statistical test results between MOEA/D-FSs and other compared algorithms.

| MOEA/D-FSs | In-range | Out-of-range |
|------------------|-------------------|------------------|
| DTLZ1 | 1.03e+2(2.34e+1)≈ | 9.84e+1(2.04e+1) |
| DTLZ2 | 1.57e-1(2.29e-2)≈ | 1.62e-1(1.90e-2) |
| DTLZ3 | 2.03e+2(2.42e+1)≈ | 2.06e+2(2.13e+1) |
| DTLZ4 | 4.91e-1(1.24e-1)≈ | 5.20e-1(6.92e-2) |
| DTLZ5 | 1.18e-1(2.25e-2)+ | 1.11e-1(2.41e-2) |
| DTLZ6 | 1.29e+0(6.44e-1)≈ | 1.36e+0(7.36e-1) |
| DTLZ7 | 4.16e+0(2.54e+0)≈ | 4.94e+0(2.31e+0) |
| + / ≈ / - | 0/7/0 | -/-/- |
| vs MOEA/D-EGO | 4/2/1 | 4/2/1 |
| vs 6 Comparisons | 29/8/5 | 28/9/5 |

784 H.3 Result Tables and Figures in IGD and HV Metrics

785 Results in IGD values are reported in Table 13 and Fig. 8. A smaller IGD value indicates a better optimization result.

Table 13: Mean IGD values and standard deviation (in parentheses) obtained from 30 runs on 7 DTLZ problems. MOEA/D-FS and comparison algorithms initialize their surrogates with 10, 60 samples, respectively. Extra 30 evaluations are allowed in the further optimization. ‘+’, ‘≈’, and ‘-’ denote MOEA/D-FS is statistically significantly superior to, almost equivalent to, and inferior to the compared algorithms in the Wilcoxon rank sum test (significance level is 0.05), respectively. The last row counts the total win/tie/loss results.

| Problems | MOEA-D-EGO | MOEA-D-FS | ParEGO | K-RVEA | KTA2 | CSEA | OREA | ESBCEO | KMOEAATIC |
|-----------|-------------------|------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| DTLZ1 | 1.07e+2(2.73e+1)≈ | 1.03e+2(2.34e+1) | 8.70e+1(2.53e+1)- | 1.22e+2(3.20e+1)+ | 1.15e+2(3.03e+1)≈ | 1.08e+2(2.64e+1)≈ | 1.11e+2(2.25e+1)+ | 1.00e+2(2.07e+1)≈ | 1.20e+2(2.71e+1)+ |
| DTLZ2 | 3.69e-1(5.72e-2)+ | 1.91e-1(2.19e-2) | 3.95e-1(3.55e-2)+ | 3.96e-1(3.55e-2)+ | 3.80e-1(4.24e-2)+ | 3.84e-1(4.05e-2)+ | 3.38e-1(3.44e-2)+ | 4.05e-1(3.07e-2)+ | 4.07e-1(3.85e-2)+ |
| DTLZ3 | 3.07e+2(5.32e+1)+ | 2.03e+2(2.42e+1) | 2.16e+2(4.89e+1)≈ | 3.53e+2(7.76e+1)+ | 3.23e+2(8.67e+1)+ | 3.35e+2(6.83e+1)+ | 3.39e+2(7.72e+1)+ | 2.41e+2(5.51e+1)+ | 3.27e+2(8.10e+1)+ |
| DTLZ4 | 8.36e-1(1.51e-1)≈ | 8.47e-1(1.87e-1) | 9.14e-1(1.22e-1)≈ | 7.28e-1(1.16e-1)- | 7.95e-1(1.49e-1)≈ | 8.41e-1(1.48e-1)≈ | 7.89e-1(1.67e-1)≈ | 7.68e-1(1.21e-1)- | 5.97e-1(1.25e-1)- |
| DTLZ5 | 2.88e-1(5.64e-2)+ | 1.22e-1(2.10e-2) | 3.10e-1(4.36e-2)+ | 2.99e-1(5.02e-2)+ | 2.73e-1(5.06e-2)+ | 2.97e-1(3.77e-2)+ | 2.12e-1(4.27e-2)+ | 3.10e-1(4.29e-2)+ | 2.93e-1(5.32e-2)+ |
| DTLZ6 | 2.08e+0(7.16e-1)+ | 1.36e+0(6.03e-1) | 2.54e+0(7.09e-1)+ | 5.24e+0(6.15e-1)+ | 4.58e+0(6.36e-1)+ | 6.45e+0(3.51e-1)+ | 5.79e+0(6.67e-1)+ | 3.10e+0(8.82e-1)+ | 3.57e+0(6.85e-1)+ |
| DTLZ7 | 2.02e+0(8.97e-1)- | 4.22e+0(2.52e+0) | 1.53e+0(1.42e+0)- | 4.03e-1(7.19e-2)- | 2.12e+0(2.13e+0)- | 5.49e+0(1.31e+0)+ | 5.53e+0(1.32e+0)+ | 1.02e+0(5.29e-1)- | 3.59e-1(1.49e-1)- |
| + / ≈ / - | 4/2/1 | -/-/- | 3/2/2 | 5/0/2 | 4/2/1 | 5/2/0 | 6/1/0 | 4/1/2 | 5/0/2 |

786

787 Results in HV values are reported in Table 14 and Fig. 9. A larger HV value indicates a better optimization result.

Table 14: Mean HV values and standard deviation (in parentheses) obtained from 30 runs on 7 DTLZ problems. MOEA/D-FS and comparison algorithms initialize their surrogates with 10, 60 samples, respectively. Extra 30 evaluations are allowed in the further optimization. ‘+’, ‘≈’, and ‘-’ denote MOEA/D-FS is statistically significantly superior to, almost equivalent to, and inferior to the compared algorithms in the Wilcoxon rank sum test (significance level is 0.05), respectively. The last row counts the total win/tie/loss results.

| Problems | MOEA-D-EGO | MOEA-D-FS | ParEGO | K-RVEA | KTA2 | CSEA | OREA | ESBCEO | KMOEAATIC |
|-----------|-------------------|------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|-------------------|
| DTLZ1 | 0.00e+0(0.00e+0)≈ | 0.00e+0(0.00e+0) | 0.00e+0(0.00e+0)≈ |
| DTLZ2 | 1.63e-1(8.93e-2)+ | 4.37e-1(3.48e-2) | 9.85e-2(3.44e-2)+ | 1.05e-1(4.43e-2)+ | 1.25e-1(4.84e-2)+ | 1.17e-1(5.59e-2)+ | 1.73e-1(4.75e-2)+ | 1.25e-1(5.19e-2)+ | 9.43e-2(4.62e-2)+ |
| DTLZ3 | 0.00e+0(0.00e+0)≈ | 0.00e+0(0.00e+0) | 0.00e+0(0.00e+0)≈ |
| DTLZ4 | 6.44e-2(6.93e-2)+ | 1.00e-1(6.58e-2) | 8.65e-3(1.74e-2)+ | 2.28e-2(4.11e-2)+ | 2.18e-2(3.52e-2)+ | 1.01e-2(2.38e-2)+ | 5.58e-2(6.13e-2)+ | 1.55e-2(2.64e-2)+ | 4.77e-2(5.93e-2)+ |
| DTLZ5 | 2.62e-2(2.46e-2)+ | 1.60e-1(1.54e-2) | 7.89e-3(1.16e-2)+ | 1.51e-2(1.58e-2)+ | 2.60e-2(1.91e-2)+ | 1.08e-2(1.14e-2)+ | 4.57e-2(2.76e-2)+ | 1.43e-2(1.32e-2)+ | 2.04e-2(2.38e-2)+ |
| DTLZ6 | 3.82e-4(2.06e-3)≈ | 1.07e-2(2.64e-2) | 0.00e+0(0.00e+0)≈ |
| DTLZ7 | 6.98e-2(1.00e-1)≈ | 4.14e-2(8.25e-2) | 8.22e-2(8.32e-2)- | 2.65e-1(3.94e-2)- | 1.31e-1(1.20e-1)- | 0.00e+0(0.00e+0)≈ | 0.00e+0(0.00e+0)≈ | 8.06e-2(8.74e-2)- | 3.58e-1(4.00e-2)- |
| + / ≈ / - | 2/5/0 | -/-/- | 3/3/1 | 3/3/1 | 3/3/1 | 3/4/0 | 3/4/0 | 3/3/1 | 3/3/1 |

788

789 I Summary of Experiments

790 Our computational studies have demonstrated the following: First, we provide empirical evidence to
791 show the effectiveness of learning experience: The meta-learning of neural network parameters and

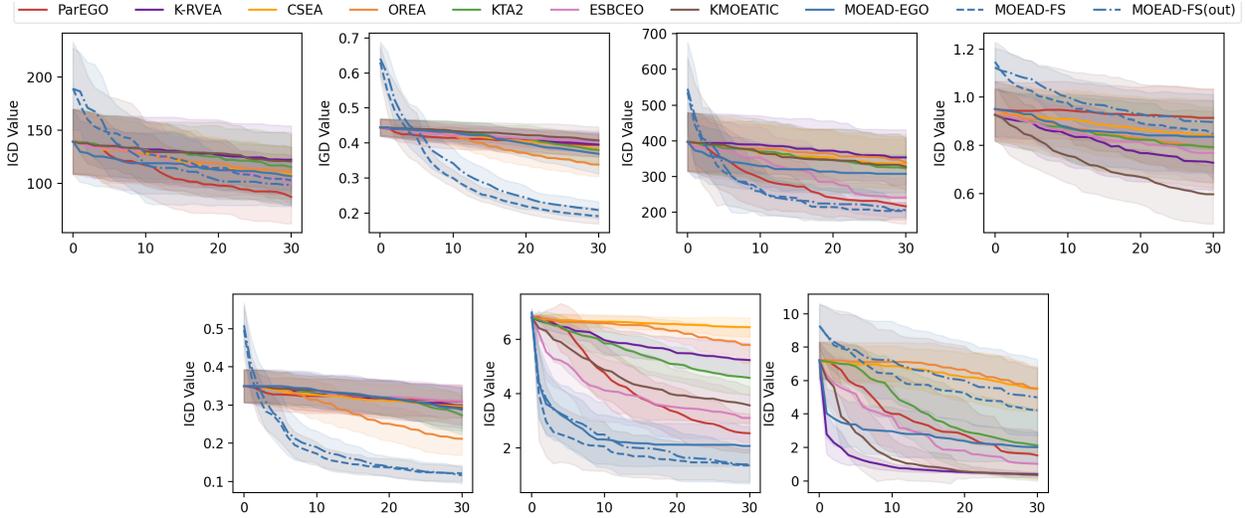


Figure 8: IGD curves averaged over 30 runs on 7 DTLZ problems. Solid lines are mean values, while shadows are error regions. **Upper:** DTLZ1, DTLZ2, DTLZ3, DTLZ4. **Lower:** DTLZ5, DTLZ6, DTLZ7. MOEA/D-FSs and comparison algorithms initialize their surrogates with 10, 60 samples, respectively. Extra 30 evaluations are allowed in the further optimization. Note that ‘FS(out)’ indicates the target task is excluded from the range of related tasks during the meta-learning procedure. X-axis denotes the number of evaluations used after the surrogate initialization.

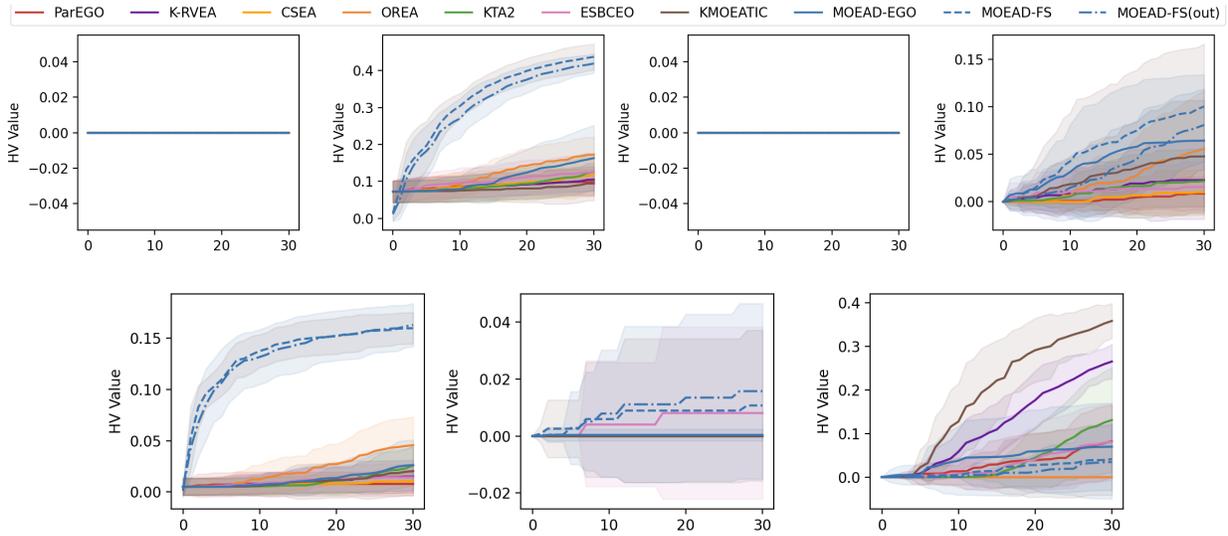


Figure 9: HV curves averaged over 30 runs on 7 DTLZ problems. Solid lines are mean values, while shadows are error regions. **Upper:** DTLZ1, DTLZ2, DTLZ3, DTLZ4. **Lower:** DTLZ5, DTLZ6, DTLZ7. MOEA/D-FSs and comparison algorithms initialize their surrogates with 10, 60 samples, respectively. Extra 30 evaluations are allowed in the further optimization. Note that ‘FS(out)’ indicates the target task is excluded from the range of related tasks during the meta-learning procedure. X-axis denotes the number of evaluations used after the surrogate initialization.

792 base kernel parameters are essential to the modeling accuracy of a MDKL model. As a result, our
 793 MDKL model outperforms the compared meta-learning modeling and non-meta-learning modeling
 794 methods on both the engine fuel consumption regression task and the sinusoid function regression
 795 task.

796 Second, we demonstrate the main contribution of this work: In most situations, the proposed FSEO
 797 framework can assist regression-based SAEAs to reach competitive or even better optimization

798 results, while the cost of surrogate initialization is only $1d$ samples. Due to the effectiveness of
799 saving evaluations, our FSEO framework is preferable to other SAEAs when solving problems within
800 a very limited evaluation budget. Moreover, some empirical guidelines are concluded to help the
801 application of our FSEO framework. For the influence of task similarity, we find that related tasks that
802 are very similar to the target task are not necessary to the application of our approach. The influence
803 of these similar tasks on the optimization performance is limited. Our FSEO framework can achieve
804 competitive results without datasets from very similar related tasks. Besides, for the related tasks
805 used for meta-learning, we have demonstrated that more useful experience can be learned if more
806 data points are sampled from related tasks.

807 Third, the effectiveness of our FSEO framework is validated on a real-world engine calibration
808 problem. Competitive or better results are achieved on the objective and constraint functions, while
809 $1d$ samples are used to initialize surrogates. Therefore, our FSEO framework can also be applied to
810 optimization scenarios such as single-objective optimization and constrained optimization.

811 **NeurIPS Paper Checklist**

812 **1. Claims**

813 Question: Do the main claims made in the abstract and introduction accurately reflect the
814 paper's contributions and scope?

815 Answer: [\[Yes\]](#)

816 Justification: Claims we made accurately reflect the paper's contributions and scope.

817 Guidelines:

- 818 • The answer NA means that the abstract and introduction do not include the claims
819 made in the paper.
- 820 • The abstract and/or introduction should clearly state the claims made, including the
821 contributions made in the paper and important assumptions and limitations. A No or
822 NA answer to this question will not be perceived well by the reviewers.
- 823 • The claims made should match theoretical and experimental results, and reflect how
824 much the results can be expected to generalize to other settings.
- 825 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
826 are not attained by the paper.

827 **2. Limitations**

828 Question: Does the paper discuss the limitations of the work performed by the authors?

829 Answer: [\[Yes\]](#)

830 Justification: See Appendix B.

831 Guidelines:

- 832 • The answer NA means that the paper has no limitation while the answer No means that
833 the paper has limitations, but those are not discussed in the paper.
- 834 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 835 • The paper should point out any strong assumptions and how robust the results are to
836 violations of these assumptions (e.g., independence assumptions, noiseless settings,
837 model well-specification, asymptotic approximations only holding locally). The authors
838 should reflect on how these assumptions might be violated in practice and what the
839 implications would be.
- 840 • The authors should reflect on the scope of the claims made, e.g., if the approach was
841 only tested on a few datasets or with a few runs. In general, empirical results often
842 depend on implicit assumptions, which should be articulated.
- 843 • The authors should reflect on the factors that influence the performance of the approach.
844 For example, a facial recognition algorithm may perform poorly when image resolution
845 is low or images are taken in low lighting. Or a speech-to-text system might not be
846 used reliably to provide closed captions for online lectures because it fails to handle
847 technical jargon.
- 848 • The authors should discuss the computational efficiency of the proposed algorithms
849 and how they scale with dataset size.
- 850 • If applicable, the authors should discuss possible limitations of their approach to
851 address problems of privacy and fairness.
- 852 • While the authors might fear that complete honesty about limitations might be used by
853 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
854 limitations that aren't acknowledged in the paper. The authors should use their best
855 judgment and recognize that individual actions in favor of transparency play an impor-
856 tant role in developing norms that preserve the integrity of the community. Reviewers
857 will be specifically instructed to not penalize honesty concerning limitations.

858 **3. Theory Assumptions and Proofs**

859 Question: For each theoretical result, does the paper provide the full set of assumptions and
860 a complete (and correct) proof?

861 Answer: [\[NA\]](#)

862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915

Justification: Not applicable.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental Result Reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: Experimental setups are described in detail.

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967

Answer: [No]

Justification: Will release our code after acceptance, or we can provide the code if any reviewers are interested in it during the review process. Anyway, the details about the code have already described in the paper.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [Yes]

Justification: We have described all the details about of experiments.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [Yes]

Justification: We have conducted statistical tests in our experiments, error bars are plotted in figures.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)

- 968 • The assumptions made should be given (e.g., Normally distributed errors).
- 969 • It should be clear whether the error bar is the standard deviation or the standard error
- 970 of the mean.
- 971 • It is OK to report 1-sigma error bars, but one should state it. The authors should
- 972 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
- 973 of Normality of errors is not verified.
- 974 • For asymmetric distributions, the authors should be careful not to show in tables or
- 975 figures symmetric error bars that would yield results that are out of range (e.g. negative
- 976 error rates).
- 977 • If error bars are reported in tables or plots, The authors should explain in the text how
- 978 they were calculated and reference the corresponding figures or tables in the text.

979 8. Experiments Compute Resources

980 Question: For each experiment, does the paper provide sufficient information on the com-
981 puter resources (type of compute workers, memory, time of execution) needed to reproduce
982 the experiments?

983 Answer: [No]

984 Justification: We did not provide information about compute workers and memory since our
985 experiments do not have specific requirements on memory or other computation resource.

986 Guidelines:

- 987 • The answer NA means that the paper does not include experiments.
- 988 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
- 989 or cloud provider, including relevant memory and storage.
- 990 • The paper should provide the amount of compute required for each of the individual
- 991 experimental runs as well as estimate the total compute.
- 992 • The paper should disclose whether the full research project required more compute
- 993 than the experiments reported in the paper (e.g., preliminary or failed experiments that
- 994 didn't make it into the paper).

995 9. Code Of Ethics

996 Question: Does the research conducted in the paper conform, in every respect, with the
997 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

998 Answer: [NA]

999 Justification: Not applicable.

1000 Guidelines:

- 1001 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1002 • If the authors answer No, they should explain the special circumstances that require a
- 1003 deviation from the Code of Ethics.
- 1004 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
- 1005 eration due to laws or regulations in their jurisdiction).

1006 10. Broader Impacts

1007 Question: Does the paper discuss both potential positive societal impacts and negative
1008 societal impacts of the work performed?

1009 Answer: [No]

1010 Justification: We do not think optimization algorithm can cause any negative social impacts.

1011 Guidelines:

- 1012 • The answer NA means that there is no societal impact of the work performed.
- 1013 • If the authors answer NA or No, they should explain why their work has no societal
- 1014 impact or why the paper does not address societal impact.
- 1015 • Examples of negative societal impacts include potential malicious or unintended uses
- 1016 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 1017 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 1018 groups), privacy considerations, and security considerations.

- 1019
- 1020
- 1021
- 1022
- 1023
- 1024
- 1025
- 1026
- 1027
- 1028
- 1029
- 1030
- 1031
- 1032
- 1033
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
 - The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
 - If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

1034 11. Safeguards

1035 Question: Does the paper describe safeguards that have been put in place for responsible
1036 release of data or models that have a high risk for misuse (e.g., pretrained language models,
1037 image generators, or scraped datasets)?

1038 Answer: [No]

1039 Justification: Code will be released after acceptance, it would be open access, no safeguards
1040 are required.

1041 Guidelines:

- 1042
- 1043
- 1044
- 1045
- 1046
- 1047
- 1048
- 1049
- 1050
- 1051
- The answer NA means that the paper poses no such risks.
 - Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
 - Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
 - We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

1052 12. Licenses for existing assets

1053 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
1054 the paper, properly credited and are the license and terms of use explicitly mentioned and
1055 properly respected?

1056 Answer: [Yes]

1057 Justification: We cited the algorithm platform and the data we used in our paper.

1058 Guidelines:

- 1059
- 1060
- 1061
- 1062
- 1063
- 1064
- 1065
- 1066
- 1067
- 1068
- 1069
- 1070
- 1071
- The answer NA means that the paper does not use existing assets.
 - The authors should cite the original paper that produced the code package or dataset.
 - The authors should state which version of the asset is used and, if possible, include a URL.
 - The name of the license (e.g., CC-BY 4.0) should be included for each asset.
 - For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
 - If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, paperswithcode.com/datasets has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
 - For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.

1072 • If this information is not available online, the authors are encouraged to reach out to
1073 the asset’s creators.

1074 **13. New Assets**

1075 Question: Are new assets introduced in the paper well documented and is the documentation
1076 provided alongside the assets?

1077 Answer: [NA]

1078 Justification: We did not introduce any new assets.

1079 Guidelines:

- 1080 • The answer NA means that the paper does not release new assets.
- 1081 • Researchers should communicate the details of the dataset/code/model as part of their
1082 submissions via structured templates. This includes details about training, license,
1083 limitations, etc.
- 1084 • The paper should discuss whether and how consent was obtained from people whose
1085 asset is used.
- 1086 • At submission time, remember to anonymize your assets (if applicable). You can either
1087 create an anonymized URL or include an anonymized zip file.

1088 **14. Crowdsourcing and Research with Human Subjects**

1089 Question: For crowdsourcing experiments and research with human subjects, does the paper
1090 include the full text of instructions given to participants and screenshots, if applicable, as
1091 well as details about compensation (if any)?

1092 Answer: [NA]

1093 Justification: We do not have any experiments or research with human subjects.

1094 Guidelines:

- 1095 • The answer NA means that the paper does not involve crowdsourcing nor research with
1096 human subjects.
- 1097 • Including this information in the supplemental material is fine, but if the main contribu-
1098 tion of the paper involves human subjects, then as much detail as possible should be
1099 included in the main paper.
- 1100 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1101 or other labor should be paid at least the minimum wage in the country of the data
1102 collector.

1103 **15. Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**
1104 **Subjects**

1105 Question: Does the paper describe potential risks incurred by study participants, whether
1106 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1107 approvals (or an equivalent approval/review based on the requirements of your country or
1108 institution) were obtained?

1109 Answer: [NA]

1110 Justification: We do not have any experiments or research with human subjects.

1111 Guidelines:

- 1112 • The answer NA means that the paper does not involve crowdsourcing nor research with
1113 human subjects.
- 1114 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1115 may be required for any human subjects research. If you obtained IRB approval, you
1116 should clearly state this in the paper.
- 1117 • We recognize that the procedures for this may vary significantly between institutions
1118 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1119 guidelines for their institution.
- 1120 • For initial submissions, do not include any information that would break anonymity (if
1121 applicable), such as the institution conducting the review.