
Automated Kernel Discovery Towards Understanding High-dimensional Bayesian Optimization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Gaussian Process (GP) kernels are central to Bayesian optimization (BO), yet
2 designing effective kernels for high-dimensional problems still relies on extensive
3 manual engineering. Existing automated approaches struggle in high dimensions
4 for two bottlenecks: their kernel search space is limited to additions and multi-
5 plications of base kernels, and LLM-based approaches require conditioning on
6 raw observations, which becomes infeasible due to context-length limits and the
7 difficulty of extracting meaningful patterns. We introduce **Kernel Discovery**, a
8 LLM-driven evolutionary framework for high-dimensional BO that searches a
9 broader kernel space beyond predefined composition rules and does not require
10 conditioning on observations. Motivated by the observation that directly prompt-
11 ing an LLM to generate kernel code yields syntactically varied but functionally
12 identical kernels, we adopt a two-stage approach: an LLM first proposes novel
13 mathematical forms, then a second LLM call converts each form into validated,
14 executable code. We also propose a leave-one-out continuous ranked probability
15 score (LOO-CRPS) as a selection criterion that penalizes overfitted kernels. On
16 five high-dimensional BO benchmarks, our method achieves an average rank of **1.2**
17 **out of 17**, outperforming competitive baselines. We further analyze the discovered
18 kernels to identify which kernels lead to improvements in high-dimensional BO.

19 1 Introduction

20 Optimizing high-dimensional black-box functions spans a wide range of machine learning problems,
21 including hyperparameter optimization [1, 2], neural architecture search [3, 4], biological sequence
22 design [5, 6], and control tasks [7, 8]. Bayesian optimization (BO) is the de facto paradigm for
23 black-box optimization, which iteratively selects promising candidates based on surrogate models
24 fitted to observations so far [9, 10]. However, it often suffers from the curse of dimensionality, as
25 large distances between observations make accurate surrogate modeling difficult and exacerbate
26 overexploration toward the boundary of the search space.

27 Prior works in high-dimensional BO have exploited structural assumptions such as additive decom-
28 positions [11, 12, 13, 14], low-dimensional subspaces [15, 16, 17], or sparsity [18, 19]. However,
29 these approaches perform well only on synthetic benchmarks with explicit underlying structures,
30 which rarely hold in practice. Another family of methods employs trust regions to effectively search
31 high-dimensional spaces [7, 20]. While locality achieves promising results across diverse tasks
32 [5, 21, 22], it requires a large number of evaluations and may struggle to escape local optima.

33 Recently, Hvarfner et al. [23] and Xu et al. [24] showed that even with basic kernels, such as RBF or
34 Matérn52 kernels, we can match or exceed specialized methods in high-dimensional BO by proper
35 scaling of the lengthscale prior with dimensionality. Moreover, Oh et al. [25] and Doumont et al.
36 [26] demonstrated that geometric input warpings (e.g., hypercylinder or hypersphere) with simple

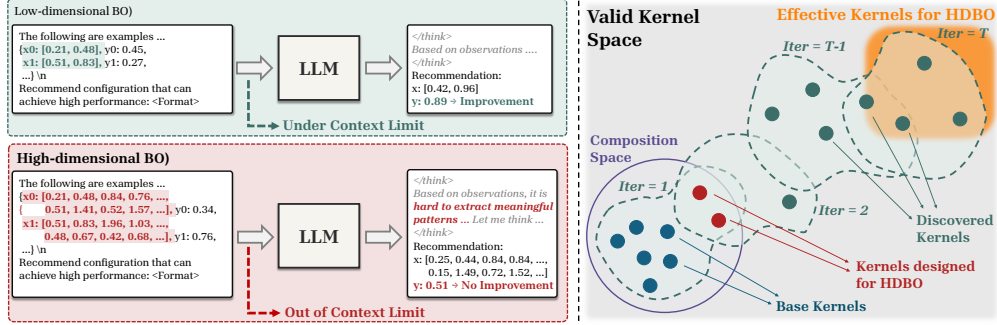


Figure 1: Motivation Figure. **(Left):** In high-dimensional BO, a conventional LLM-based BO does not work due to the out-of-context limits and the difficulty of extracting meaningful patterns. **(Right):** Our pipeline enables us to discover effective kernels for high-dimensional BO beyond compositions.

37 models can surpass sophisticated baselines on high-dimensional, real-world benchmarks. Although
 38 these breakthroughs are promising, identifying such effective priors or transformations still requires
 39 extensive manual engineering by domain experts. This suggests that automating kernel design with a
 40 system capable of mathematical reasoning could substantially reduce this barrier.

41 This naturally prompts us to leverage Large Language Models (LLMs, 27, 28, 29), which have
 42 internalized vast mathematical knowledge and can propose novel functional forms when incorporated
 43 into an evolutionary algorithm (EA) pipeline [30, 31, 32, 33, 34]. However, existing LLM-based
 44 BO methods do not directly scale to high-dimensional problems due to two bottlenecks. First, an
 45 **expressive bottleneck**: Existing kernel search methods are limited to additive and multiplicative
 46 compositions of base kernels, thereby limiting their ability to discover novel kernels suitable for
 47 high-dimensional BO. Second, an **interface bottleneck**: most LLM-based BO methods require
 48 conditioning on raw observations, which becomes infeasible in high dimensions due to context-length
 49 limits and the difficulty of extracting patterns from a long stream of numbers, as shown in Figure 1.

50 To address these limitations, we introduce **Kernel Discovery**, a novel framework for designing
 51 effective kernel structures for high-dimensional BO. First, we **initialize** the population with kernels
 52 tailored for high-dimensional BO. At each BO iteration, we **discover** new kernels through a two-stage
 53 approach: an LLM first proposes novel mathematical forms, and a second LLM converts the form
 54 into executable code. Through this decomposition, we encourage LLMs to explore the broader space
 55 of valid kernels by leveraging mathematical reasoning capabilities. We then **validate** the discovered
 56 kernels via agnostic execution and PSD checks, retaining only valid GP kernels. To **choose** the most
 57 promising kernel, we introduce LOO-CRPS, which is less prone to in-sample overfitting than the
 58 marginal log-likelihood metric. Finally, we evaluate the objective at the selected point and **update**
 59 the dataset and population, repeating until the evaluation budget is exhausted.

60 We conduct experiments on five high-dimensional BO benchmarks and achieve an average rank of
 61 **1.2 out of 17**, outperforming competitive baselines. We also conduct several ablation studies on
 62 the design choices of our framework. Finally, we analyze the discovered kernels and uncover that
 63 unexpected kernels (e.g., compositions of geometric warping or non-stationary components) might
 64 lead to improvements, providing insights into what makes a kernel effective in high dimensions.

65 2 Related Work

66 **High-dimensional Bayesian Optimization.** Several approaches have been proposed to push the
 67 limit of BO in high dimensions. One line of work exploits explicit structural assumptions, such as
 68 low-dimensional subspaces [15, 16, 17, 35], sparsity [18, 19], variable selection [7, 36, 37], and
 69 additive decompositions [11, 12, 13, 14]. However, such assumptions rarely hold in real-world
 70 problems. Another line of work employs trust regions, in which we constrain the search space to
 71 prevent the search among boundary points in high-dimensional space [5, 20, 22, 38, 39, 40]. While it
 72 achieves superior performance across various domains, particularly in scalable settings, it requires
 73 too large a number of evaluations and remains susceptible to being trapped in local optima.

74 More recently, a series of surprising findings have challenged longstanding intuitions about high-
 75 dimensional BO. Hvarfner et al. [23] and Xu et al. [24] showed that simple kernels with well-chosen

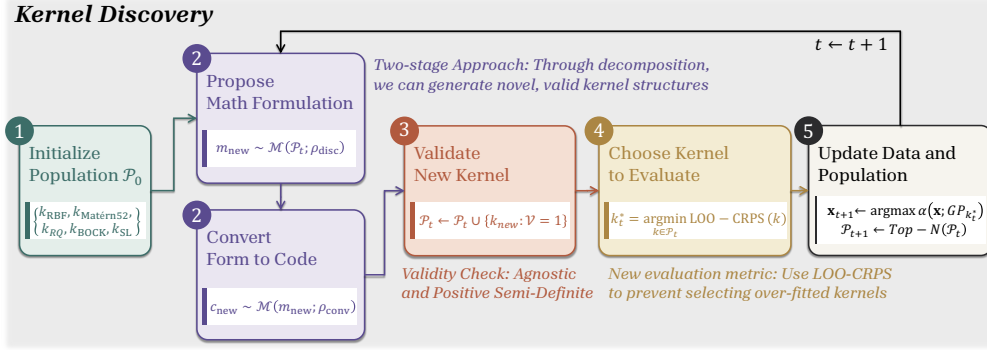


Figure 2: Overview of the method. Given an initial population, we instruct an LLM to generate novel kernels via a two-stage approach: mathematical formulation and code conversion. We then perform a sanity check to retain only valid kernels and select the kernel with the lowest LOO-CRPS value. Finally, we propose a next query using the chosen kernel and update the dataset and population.

lengthscale priors can match or exceed sophisticated methods. Papenmeier et al. [41] also claimed that scaling the initial lengthscale and random axis-aligned subspace perturbation sampling (RAASP, 42) can mitigate the vanishing gradient issue in acquisition function optimization in high dimensions. Doumont et al. [26] further demonstrated that the smoothness of surrogate models is a critical factor, and that spherical input mappings with a simple linear basis suffice for competitive performance. Together, these results reveal that the community’s prevailing assumptions about what drives high-dimensional BO may be incomplete. This underscores the need for a more in-depth investigation of kernel design in high-dimensional BO, which is precisely the gap our work aims to fill.

LLMs for Black-box Optimization. Recent advancements in LLMs have demonstrated strong capabilities for complex reasoning [27, 28, 29], prompting several works to integrate them into the BO pipeline. Aglietti et al. [32] and Ngo et al. [43] utilize LLM to discover novel acquisition functions, while Li et al. [44] and Liu et al. [45] replace the entire BO loop with an LLM-driven system. However, these methods are primarily evaluated in low-dimensional settings and rely on prepending prior observations as context, which becomes infeasible in high dimensions due to context-length constraints. Furthermore, extracting meaningful structure from dense observations can be challenging even for recent, closed-source LLMs.

Most closely related to our work is CAKE [46], which employs an LLM to select and combine base kernels to improve BO surrogates. While it shares our motivation for LLM-guided kernel design, it also relies on prior trials as context and restricts the search space to additive and multiplicative compositions of base kernels. In contrast, our framework does not condition on raw observations and actively explores a broader space of valid kernels beyond compositions, thereby enabling more expressive surrogate modeling for high-dimensional problems.

3 Preliminaries

Bayesian Optimization. In BO, we aim to find an input $\mathbf{x} \in \mathcal{X}$ that maximizes the unknown black-box function $f : \mathcal{X} \rightarrow \mathbb{R}$, which is typically expensive to evaluate and not differentiable. BO can be modularized into three orthogonal components: a surrogate model, an acquisition function, and an optimizer for the acquisition function [47]. At each iteration, we fit the surrogate model to the observations so far and select the next query that maximizes the chosen acquisition function using the optimizer. While there are several possible options for each component, the conventional approach is to use a Gaussian Process (GP, 48) as a surrogate model, Expected Improvement (EI, 49) as an acquisition function, and L-BFGS-B [50] as an optimizer.

Gaussian Processes and Kernels. A GP defines a distribution over functions and is specified by a mean function (constant in most cases) and a covariance matrix \mathbf{K} . The covariance matrix can be derived by using a kernel function $k(\mathbf{x}, \mathbf{x}') : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$. The most frequently used kernel for BO is the RBF kernel [51], which can be defined as follows:

$$k_{\text{RBF}}(\mathbf{x}, \mathbf{x}') := \sigma^2 \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}')^\top \mathbf{L}^{-1}(\mathbf{x} - \mathbf{x}')\right) \quad (1)$$

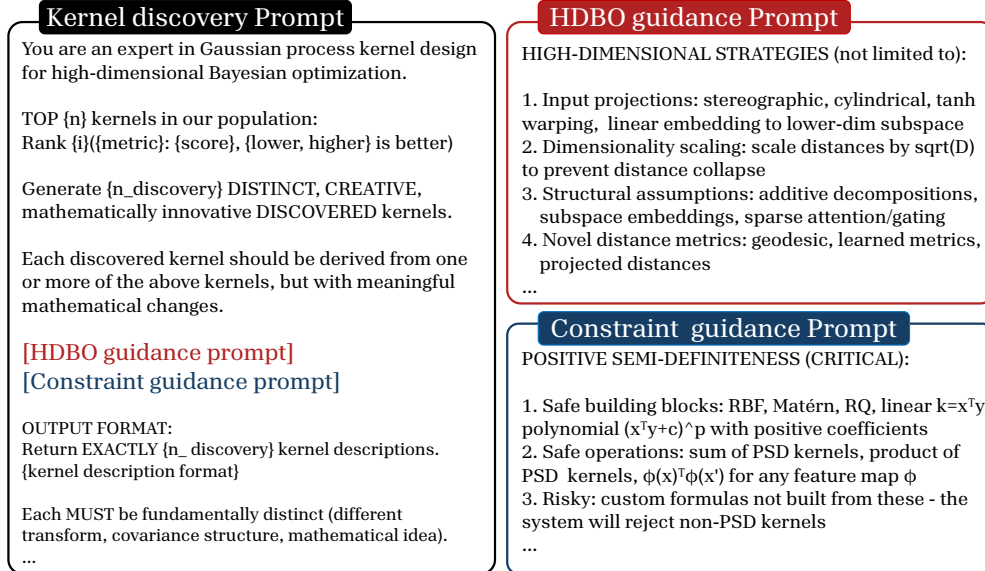


Figure 3: Prompts for Kernel Discovery. To discover novel kernels for high-dimensional BO, we start with a prompt that simulates the reasoning of a human expert. Then, we instruct the LLM to understand prior strategies for high-dimensional BO and the constraints it should satisfy.

111 where $\sigma^2 > 0$ is the output scale and $\mathbf{L} = \text{diag}(\ell_1^2, \dots, \ell_d^2)$ is a diagonal matrix of per-dimension
 112 lengthscale hyperparameters $\ell_i > 0$.

113 Kernels can be combined or constructed in several ways. First, the sum and product of two valid
 114 kernels are also valid kernels, enabling the construction of richer covariance structures. Second,
 115 for any feature map $\phi : \mathcal{X} \rightarrow \mathbb{R}^m$, the inner product $k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$ defines a positive
 116 semi-definite kernel by construction. By injecting these compositional and feature-map properties
 117 as inductive biases into an LLM-driven kernel construction process, we can discover novel kernel
 118 structures tailored for high-dimensional BO.

119 4 Methods

120 In this section, we introduce **Kernel Discovery**, a novel framework for discovering effective kernels
 121 for high-dimensional BO by leveraging LLMs as an evolutionary operator. We first initialize the
 122 population with well-known base kernels tailored for high dimensions. Then, we guide the LLM to
 123 generate novel kernel structures as mathematical forms and convert them into executable Python code.
 124 We conduct a sanity check to determine whether the generated kernel is valid. We also introduce a
 125 new selection criterion to choose the most promising kernel for acquisition function optimization.
 126 We repeat the process until convergence. We summarize the overview of our framework in Figure 2.

127 **Notations.** We represent each kernel $k = (m, c)$, where m is its mathematical expression and c is
 128 the corresponding code. Let $\mathcal{P}_t = \{k_1^{(t)}, \dots, k_N^{(t)}\}$ denote the population of size N at iteration t . We
 129 denote by $\mathcal{M}(\cdot; \rho)$ a call to the LLM conditioned on a role-specific prompt ρ . We also introduce an
 130 auxiliary verifier \mathcal{V} to check whether the discovered kernel is valid.

131 **Initialize Population.** Similar to several approaches that integrate LLMs into the evolutionary loop
 132 [31], we initialize the base population with existing kernels that are widely used for BO. As our focus
 133 is on high-dimensional BO, we need to carefully select base kernels that will encourage LLMs to
 134 extract useful knowledge from them and give insights into how to propose better kernels. To this end,
 135 we set the initial population as follows:

$$\mathcal{P}_0 = \{k_{\text{RBF}}, k_{\text{Matérn52}}, k_{\text{RQ}}, k_{\text{BOCK}}, k_{\text{SL}}\}, \quad (2)$$

136 For the RBF, Matérn52, and RQ, we use the lengthscale prior $\ell_i \sim \mathcal{LN}(\sqrt{2} + \log(\sqrt{D}), \sqrt{3})$,
 137 proposed by Hvarfner et al. [23]. BOCK [25] and SL [26] propose geometric mappings to avoid
 138 overexploration. We do not include Linear and Periodic as they perform poorly in high dimensions.

139 **Discover Kernels with LLMs.** Starting from the initial population \mathcal{P}_0 , we leverage LLMs to propose
 140 new kernel structures. Specifically, our discovery pipeline consists of two stages: mathematical
 141 formulation and code generation.

142 While one can directly prompt an LLM to generate novel kernels in the code generation space, this is
 143 undesirable, as LLMs often produce functionally redundant kernel codes (e.g., by renaming variables
 144 or substituting $x-y$ with $x.\text{sub}(y)$) rather than genuinely novel structures. To mitigate this, we first
 145 express the existing kernels in mathematical form $\{m_i^{(t)}\}_{i=1}^N$ and prompt the LLM to generate new
 146 mathematical forms that satisfy the constraints and are suitable for high-dimensional BO:

$$m_{\text{new}} \sim \mathcal{M}\left(\{m_i^{(t)}\}_{i=1}^N; \rho_{\text{disc}}\right) \quad (3)$$

147 where ρ_{disc} is a prompt for discovery stage as depicted in Figure 3. After discovering a new mathe-
 148 matical form, we introduce another LLM to convert the formulation into executable code:

$$c_{\text{new}} \sim \mathcal{M}(m_{\text{new}}; \rho_{\text{conv}}) \quad (4)$$

149 where ρ_{conv} is a prompt for conversion stage as depicted in Figure 9. Rather than conditioning on
 150 raw observations, we provide only the mathematical forms from previous iterations as context. This
 151 dramatically reduces the context length, allowing the LLM to focus on structural novelty rather than
 152 sifting through high-dimensional numerical data for patterns.

153 Following prior works, we find that composing kernels is the simplest way to discover new valid
 154 kernels. We therefore add a composition stage in which the LLM takes the top- k kernels as input and
 155 combines them via addition or multiplication. You can find the prompt for composition in Figure 10.

156 **Validate Discovered Kernels.** Given a new candidate $k_{\text{new}} = (m_{\text{new}}, c_{\text{new}})$, we filter it with two
 157 empirical checks: that the code runs correctly across varying input shapes (\mathcal{V}_{agn}) and that the Cholesky
 158 decomposition succeeds on random test inputs (\mathcal{V}_{psd}). Formally:

$$\mathcal{V}(k_{\text{new}}) = \mathcal{V}_{\text{agn}}(k_{\text{new}}) \wedge \mathcal{V}_{\text{psd}}(k_{\text{new}}) \quad (5)$$

159 where \mathcal{V}_{agn} tests the forward function on inputs of varying batch size and dimensionality, and \mathcal{V}_{psd}
 160 checks for a successful Cholesky decomposition. Only candidates with $\mathcal{V}(k_{\text{new}}) = 1$ are added to the
 161 current population pool \mathcal{P}_t .

$$\mathcal{P}_t \leftarrow \mathcal{P}_t \cup \{k^{\text{new}} : \mathcal{V}(k^{\text{new}}) = 1\} \quad (6)$$

162 **Choose the Promising Kernel.** After generating new kernels, we need to select which kernel $k \in \mathcal{P}_t$
 163 should be used for the acquisition function maximization. While the marginal log-likelihood (MLL) is
 164 a straightforward approach, it may prefer overly complex kernels that overfit the current observations
 165 and sacrifice predictive calibration. To this end, we introduce the Leave-one-out Continuous Ranked
 166 Probability Score (LOO-CRPS), a selection criterion that scores held-out predictive distributions and
 167 thereby penalizes overconfident fits more directly than MLL, defined as follows:

$$\text{CRPS}(\mathcal{N}(\mu, \sigma^2), y) = \sigma \left[\frac{y - \mu}{\sigma} \left(2\Phi\left(\frac{y - \mu}{\sigma}\right) - 1 \right) + 2\phi\left(\frac{y - \mu}{\sigma}\right) - \frac{1}{\sqrt{\pi}} \right] \quad (7)$$

168 where ϕ and Φ are the standard normal PDF and CDF, respectively. While CRPS is more robust than
 169 MLL, it can still overfit when evaluated on training data. To alleviate this, we aggregate CRPS over
 170 leave-one-out predictives across all n observations:

$$\text{LOO-CRPS}(\mathbf{K}_k, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \text{CRPS}(\mathcal{N}(\mu_{-i}, \sigma_{-i}^2), y_i), \quad \mu_{-i} = y_i - \frac{[\mathbf{K}_k^{-1}\mathbf{y}]_i}{[\mathbf{K}_k^{-1}]_{ii}}, \quad \sigma_{-i}^2 = \frac{1}{[\mathbf{K}_k^{-1}]_{ii}} \quad (8)$$

171 where \mathbf{K}_k is a kernel matrix derived from the kernel k . By leveraging the standard LOO identity for
 172 GP, the predictive mean μ_{-i} and variance σ_{-i}^2 can be computed analytically from a single matrix
 173 inversion. At each round, we select

$$k_t^* = \arg \min_{k \in \mathcal{P}_t} \text{LOO-CRPS}(\mathbf{K}_k, \mathbf{y}) \quad (9)$$

174 Finally, we use the selected kernel k_t^* and fit its hyperparameters (e.g., lengthscale, scale, variance),
 175 maximize the acquisition function $\alpha(\mathbf{x}; \text{GP}_{k_t^*})$ to obtain the next query point \mathbf{x}_{t+1} , evaluate the
 176 corresponding objective $y_{t+1} = f(\mathbf{x}_{t+1})$, and augment the dataset $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup \{(\mathbf{x}_{t+1}, y_{t+1})\}$.

Table 1: Performance comparison of various methods across standard benchmarks. D denotes the dimensionality of the task. **Blue** denotes the best entry in the column, and **Violet** denotes the second best. Experiments are conducted with 4 random seeds.

Method	Rover (\uparrow) ($D = 100$)	Mopta08 (\downarrow) ($D = 124$)	Lasso-DNA (\downarrow) ($D = 180$)	SVM (\downarrow) ($D = 388$)	Humanoid (\uparrow) ($D = 6392$)	Average Rank
Base Kernels						
RBF	3.552 \pm 0.304	217.75 \pm 2.33	0.291 \pm 0.001	0.061 \pm 0.004	503.00 \pm 67.82	6.6 / 17
Matérn52	3.453 \pm 0.366	215.77 \pm 0.57	0.292 \pm 0.003	0.063 \pm 0.003	509.79 \pm 102.32	6.8 / 17
Linear	3.950 \pm 0.441	274.67 \pm 8.87	0.312 \pm 0.004	0.226 \pm 0.003	435.49 \pm 42.00	13.0 / 17
Periodic	3.664 \pm 0.231	309.10 \pm 5.01	0.332 \pm 0.005	0.226 \pm 0.002	413.83 \pm 71.30	15.0 / 17
BOCK	3.725 \pm 0.593	225.36 \pm 1.82	0.291 \pm 0.003	0.068 \pm 0.007	669.52 \pm 50.15	6.1 / 17
SL	4.096 \pm 0.473	246.78 \pm 3.92	0.297 \pm 0.001	0.112 \pm 0.007	637.72 \pm 37.86	8.2 / 17
RQ	3.858 \pm 0.515	217.53 \pm 4.25	0.294 \pm 0.002	0.069 \pm 0.017	600.80 \pm 152.49	6.6 / 17
Search-based						
Greedy Search	3.606 \pm 0.270	225.36 \pm 1.82	0.291 \pm 0.003	0.066 \pm 0.004	687.12 \pm 119.52	6.1 / 17
Compositional Search	3.869 \pm 0.531	218.76 \pm 1.17	0.288 \pm 0.003	0.067 \pm 0.007	700.58 \pm 81.55	4.4 / 17
CAKE	3.412 \pm 0.586	231.40 \pm 15.67	0.300 \pm 0.011	0.131 \pm 0.036	667.49 \pm 24.99	9.8 / 17
LLM-based						
ARM	3.981 \pm 0.378	238.17 \pm 3.65	0.319 \pm 0.003	0.228 \pm 0.003	457.14 \pm 44.96	12.0 / 17
ATRBO	3.054 \pm 0.816	267.65 \pm 4.63	0.328 \pm 0.010	0.221 \pm 0.016	428.80 \pm 295.37	15.2 / 17
TREvol	3.398 \pm 0.620	236.41 \pm 3.96	0.301 \pm 0.007	0.079 \pm 0.007	469.39 \pm 144.43	11.6 / 17
TROpt	4.114 \pm 0.409	234.23 \pm 1.86	0.300 \pm 0.006	0.110 \pm 0.005	520.38 \pm 65.42	7.8 / 17
TRPareto	3.872 \pm 0.467	253.45 \pm 4.73	0.296 \pm 0.001	0.185 \pm 0.008	509.90 \pm 114.06	10.0 / 17
LMABO	3.953 \pm 0.369	240.70 \pm 5.24	0.304 \pm 0.006	0.226 \pm 0.001	390.23 \pm 44.78	12.6 / 17
Kernel Discovery (Ours)	4.353 \pm 0.319	216.81 \pm 0.87	0.286 \pm 0.001	0.056 \pm 0.003	762.78 \pm 72.39	1.2 / 17

177 **Update the Population.** To maintain a fixed population size N across rounds, we then truncate \mathcal{P}_{t+1}
178 by retaining the N members with the lowest LOO-CRPS:

$$\mathcal{P}_{t+1} \leftarrow \text{Top-N}(\mathcal{P}_t; \text{LOO-CRPS}). \quad (10)$$

179 This elitist update ensures that the population progressively favors kernels with more calibrated
180 predictives while preventing unbounded growth. We additionally discard kernels that fail to improve
181 performance over consecutive rounds. We provide full details of the algorithm in Algorithm 1.

182 5 Experiments

183 We evaluate our framework on five high-dimensional BO benchmarks and compare against 16
184 baselines, including base kernels, search-based methods, and LLM-based methods. Our code is here.

185 **Benchmarks.** Following prior works on high-dimensional BO [23, 26], we evaluate our framework
186 on five standard benchmarks with dimensionality ranging from $D = 100$ to $D = 6392$: (i) **Rover**
187 ($D = 100, 52$), (ii) **Mopta08** ($D = 124, 18$), (iii) **Lasso-DNA** ($D = 180, 53$), (iv) **SVM** ($D = 388,$
188 18), (v) **Humanoid** ($D = 6392, 7$). We exclude **Ant** since the function can be easily exploited by
189 random initialization. Please refer to Appendix A for the detailed description of each task.

190 **Baselines.** We compare against 16 baselines across three categories: (1) BO methods that rely on a
191 single fixed kernel, (2) search-based kernel methods, and (3) LLM-based BO methods. Please refer
192 to Appendix B for the detailed description of each baseline.

- 193 • **Base kernels:** BO methods that rely on a single kernel: RBF, Matérn52, Linear, Periodic, RQ,
194 BOCK, and SL.
- 195 • **Search-based Methods:** Methods that search over kernel candidates, including Greedy Search
196 [54], Compositional Search [55], and CAKE [46]. For a fair comparison, all search-based
197 methods share the same initial population \mathcal{P}_0 defined in Equation (2).
- 198 • **LLM-based Methods:** LLM-based BO methods, including end-to-end approaches (ARM,
199 ATRBO, TREvol, TROpt, and TRPareto) from the LLAMEA-BO [44], and LMABO [43],
200 which adaptively selects the acquisition function.

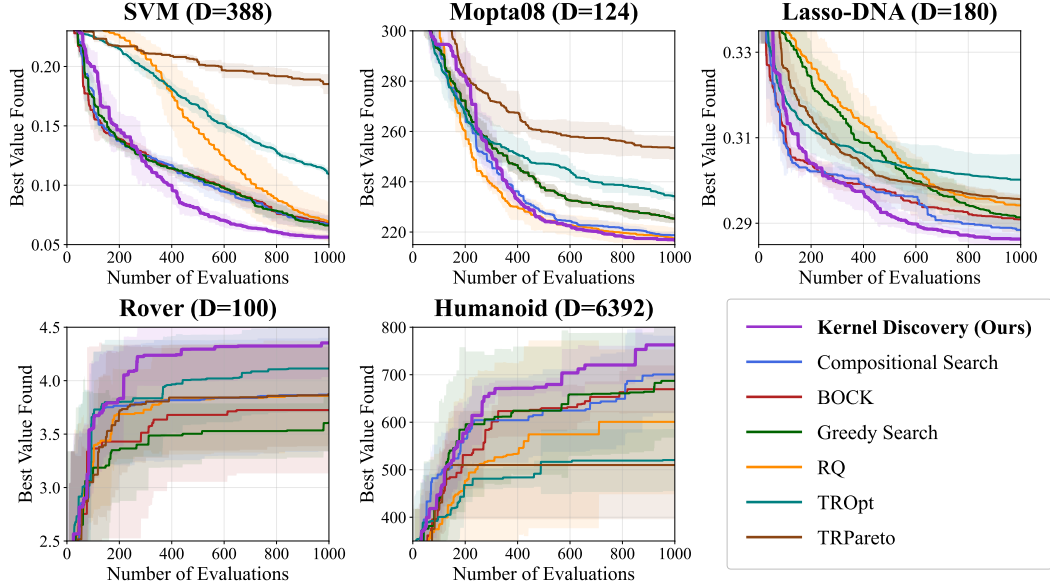


Figure 4: Learning curve of various methods across standard benchmarks. We visualize the top-2 baselines by average rank per category for clarity. Please refer to Figure 11 for the full results.

201 **Implementation.** For all methods (except for LLM-based approaches that change other parts of the
 202 BO pipeline), we use LogEI [56] as an acquisition function and L-BFGS-B as an optimizer. We use
 203 the Sobol sequence [57] to generate initial samples with size $|\mathcal{D}_0| = 20$, use a batch size of $q = 20$,
 204 and set the evaluation budget $T = 1000$. For our method, we set the population size to $N = 10$,
 205 generate two new kernel candidates (one from the discovery stage and one from the composition
 206 stage) per BO iteration. For a fair comparison, we use GPT-4o as the LLM backbone for both our
 207 method and all LLM-based baselines. Please refer to Appendix C for more details on implementation.

208 **Main Results.** We summarize the results of all methods in Table 1. As shown in the table, Kernel
 209 Discovery achieves an average rank of **1.2** across benchmarks, the best among all 17 methods.
 210 The runner-up is Compositional Search (rank 4.4), confirming that actively searching over kernel
 211 candidates is beneficial. However, due to the restricted search space, it is insufficient to discover
 212 effective kernels in high dimensions. Notably, CAKE (rank 9.8) performs substantially worse than
 213 competitive base kernels, indicating that injecting observations as context for the LLM is ineffective
 214 and degrades performance in high dimensions. Similarly, other LLM-based baselines underperform,
 215 indicating that leveraging LLMs to discover kernels within the BO pipeline is more effective than
 216 leveraging the LLM itself as a black-box optimizer.

217 We visualize the learning curve of kernel discovery against several baselines in Figure 4. Our approach
 218 consistently improves throughout training by discovering novel kernel functions at each BO iteration.
 219 Compositional Search and BOCK sometimes perform well in early rounds but stagnate in suboptimal
 220 regions. Our method avoids this by continually evolving the kernel population. We also conduct a
 221 runtime analysis of our method and the baselines in Appendix F.

222 6 Ablation Studies

223 We discuss the key design choices in our pipeline, justifying the need for the two-stage approach and
 224 analyzing the importance of each component.

225 **Importance of Two-Stage Approach.** We observed that when we directly instruct an LLM to
 226 generate code for novel kernel structures, it sometimes changes the variable names and replaces
 227 operators with equivalent operators, as shown in Figure 5a. We further demonstrate this phenomenon
 228 by measuring the functional diversity between generated kernels from direct code generation and
 229 those from our two-stage approach. As depicted in Figure 5b, our approach consistently produces
 230 more genuinely diverse kernel candidates. We also found that this functional diversity translates
 231 to consistently better performance (See Table 8), confirming that the two-stage approach discovers
 232 genuinely novel and effective kernel structures. Further details on the definition of the diversity metric,
 233 performance comparison, and functionally redundant code examples are provided in Appendix G.

Code Snippet 1

```

x1s = x1 / self.lengthscale
x2s = x2 / self.lengthscale

dist = torch.cdist(x1s, x2s, p=2).clamp(
    min=1e-15)
sqrt5_d = math.sqrt(5.0) * dist

exp_comp = torch.exp(-sqrt5_d)
covar = (1.0 + sqrt5_d + (5.0/3.0) *
    (dist**2)) * exp_comp

```

Code Snippet 2

```

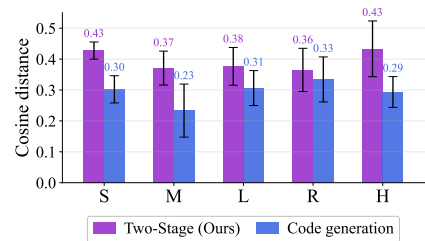
x1, x2 = x1 / self.lengthscale, x2 /
self.lengthscale

dist = torch.cdist(x1, x2, p=2).clamp(
    min=1e-15)
scaled_d = math.sqrt(5.0) * dist

covar = (1.0 + scaled_dist +
    scaled_dist.pow(2) / 3.0) *
    torch.exp(-scaled_d)

```

— Variable Renaming
— Syntax Variation
— Algebraic Rewriting



(a) Example code snippets from direct code generation. Both snippets result in the same function.

(b) Cosine distance among kernels from direct code generation vs two-stage approach.

Figure 5: Ablation studies on the two-stage approach for kernel discovery.

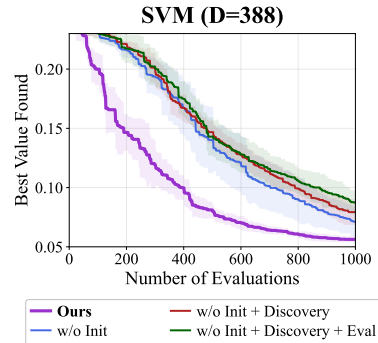


Figure 6: Ablation on each component of Kernel Discovery.

Table 2: Analysis on prompt template. Experiments are conducted with four random seeds in the SVM benchmark.

	Best Value Found	Failure Prob. (%)
Original Prompt (Ours)	0.056 ± 0.002	39.5
w/o HDBO Guidance	0.068 ± 0.007	35.7
w/o Constraint Guidance	0.064 ± 0.003	44.9

Table 3: Robustness on different LLMs. Experiments are conducted with four random seeds in the SVM benchmark.

	Best Value Found	Failure Prob. (%)
GPT-4o (Ours)	0.056 ± 0.002	39.5
GPT-4o-mini	0.060 ± 0.003	43.1
GPT-5-mini	0.056 ± 0.003	30.4

234 **Ablation on Each Component.** We conduct a systematic ablation study by removing each component.
235 For the initial population, we replace BOCK and SL with Linear and Periodic kernels, matching the
236 base kernel set used in CAKE. For the evaluation metric, we replace LOO-CRPS with the standard
237 marginal log-likelihood (MLL). As shown in Figure 6, performance consistently drops as each
238 component is removed. In particular, the initial population significantly affects sample efficiency,
239 validating that kernels specifically designed for high-dimensional BO encode inductive biases that are
240 valuable for effective surrogate modeling. We extend this ablation with more variants in Appendix H.

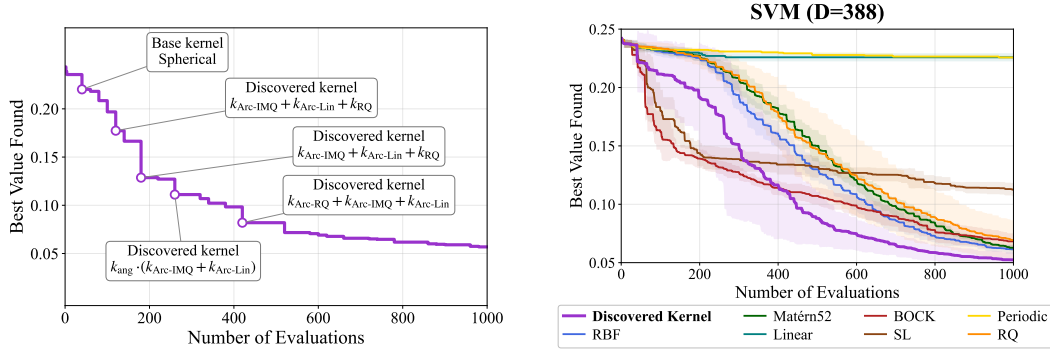
241 **Analysis on Prompt Template.** As shown in Figure 3, our prompt includes two key guidelines for
242 the LLM: HDBO and Constraint guidelines. We ablate each by removing it in turn. As shown in
243 Table 2, removing the HDBO guideline degrades performance, whereas removing the Constraint
244 guideline leads to a high failure rate.

245 **Robustness to Different LLMs.** By default, we use GPT-4o for kernel formulation and code
246 generation. To assess robustness across different LLMs, we run experiments with different LLMs:
247 GPT-4o-mini and GPT-5-mini. As shown in Table 3, we find no significant performance gap
248 across models, though GPT-4o-mini has a slightly higher kernel rejection rate due to weaker code
249 reasoning, while GPT-5-mini generates fewer invalid kernels. We also conduct ablation studies with
250 open-source LLMs and different LLM backbones for other LLM-based methods in Appendix I.

251 7 Analysis of Discovered Kernels

252 In this section, we characterize which kernels our pipeline discovers, how they evolve over training
253 iterations, and whether they transfer across benchmarks.

254 **Evolution of Discovered Kernels.** A key question is whether our pipeline discovers genuinely novel
255 structures or merely recovers known base kernels. To investigate this, we conduct a systematic study
256 of the kernels selected by our pipeline over training iterations in Figure 7a. In the initial rounds, the
257 pipeline favors simple base kernels with geometric input warpings (e.g., Spherical). As optimization
258 progresses, more advanced structures emerge: the pipeline introduces geometric transformations such
259 as \arctan warping and composes them with existing kernels (e.g., $RQ \times \text{Linear}$ with \arctan ,



(a) Summary of selected kernels for the top-5 performance improvements.

(b) Evaluation of discovered kernel from SVM benchmark directly without kernel search.

Figure 7: Post-analysis on kernel discovery. Both experiments are conducted in SVM benchmark.

260 further combined with BOCK or IMQ). This suggests that compositions of multiple geometrically-
 261 warped kernels are a promising direction for high-dimensional BO, going beyond the single-warping
 262 strategies explored by prior works. We provide detailed explanations of the discovered kernels, along
 263 with the same evolution analysis across all benchmarks, in Appendix J.

264 **Transferability of Discovered Kernels.** While our method performs kernel discovery for a single
 265 black-box function, we can extract kernels discovered by one benchmark and determine whether they
 266 can be transferred to other benchmarks. To verify this, we choose one kernel discovered from the
 267 SVM benchmark as follows:

$$k_{\text{discover}}(\mathbf{x}, \mathbf{x}') = k_{\text{Matérn52}}(\mathbf{x}, \mathbf{x}') \cdot [k_{\text{tanh-Poly}}(\mathbf{x}, \mathbf{x}') + k_{\text{RQ}}(\mathbf{x}, \mathbf{x}')] \quad (11)$$

268 Note that this kernel cannot be obtained via naive compositions since k_{poly} does not exist in the base
 269 population and the $\text{tanh}(\cdot)$ operation is not allowed. As shown in Figure 7b, the discovered kernel
 270 outperforms base kernels on the SVM benchmark, and achieves an average rank of 2.8 across 8 fixed-
 271 kernel methods (full results in Table 12). To understand why the discovered kernel achieves superior
 272 performance on the SVM benchmark, we analyze the boundary hit ratio and the observation traveling
 273 salesman distance (OTSD) and find that the discovered kernel does not exhibit boundary-seeking
 274 behavior (See Figure 24). It is noteworthy that $k_{\text{tanh-Poly}}$ is a non-stationary kernel that is not typically
 275 preferred for BO, yet it performs well in high dimensions. This suggests that the role of non-stationary
 276 kernels in high-dimensional BO warrants further investigation. We provide detailed explanations of
 277 the discovered kernel and also present other frequently discovered kernels in Appendix K.

278 8 Conclusion

279 We presented Kernel Discovery, an LLM-driven evolutionary framework for designing effective
 280 GP kernels in high-dimensional BO. We adopt a two-stage pipeline in which one LLM proposes
 281 novel mathematical kernel forms and the other LLM converts them into executable code. To select
 282 among discovered kernels without overfitting to current observations, we introduce LOO-CRPS, an
 283 overfit-resistant criterion. On five high-dimensional BO benchmarks, our method achieves an average
 284 rank of 1.2 out of 17 methods. Beyond performance, our analysis of the discovered kernels provides
 285 new insights into what drives effective surrogate modeling in high-dimensional BO.

286 **Limitations and Future Work.** Our current framework runs the discovery process independently
 287 for each benchmark, optimizing kernels for each benchmark. While we demonstrate that kernels
 288 discovered on one benchmark can transfer to others, a natural extension is to develop a general
 289 discovery strategy that operates across multiple benchmarks simultaneously, leveraging shared
 290 structure to amortize the search cost. Additionally, some kernels generated by LLMs do not pass the
 291 verifiers for positive semi-definiteness or dimension-agnosticity. We believe that the failure rate could
 292 be further reduced by detailed prompt engineering or by imposing constraints during generation,
 293 although we view these as engineering directions orthogonal to our main contribution, which is to
 294 demonstrate that LLM-driven kernel discovery is a promising strategy for high-dimensional BO.

295 **References**

- 296 [1] Ryan Turner, David Eriksson, Michael McCourt, Juha Kiili, Eero Laaksonen, Zhen Xu, and
297 Isabelle Guyon. Bayesian optimization is superior to random search for machine learning
298 hyperparameter tuning: Analysis of the black-box optimization challenge 2020. In *NeurIPS*
299 *2020 Competition and Demonstration Track*, 2021.
- 300 [2] Zi Wang, George E Dahl, Kevin Swersky, Chansoo Lee, Zelda Mariet, Zachary Nado, Justin
301 Gilmer, Jasper Snoek, and Zoubin Ghahramani. Pre-training helps bayesian optimization too.
302 In *ICML Workshop on Adaptive Experimental Design and Active Learning in the Real World*,
303 2022.
- 304 [3] Kevin Swersky, David Duvenaud, Jasper Snoek, Frank Hutter, and Michael A Osborne. Raiders
305 of the lost architecture: Kernels for bayesian optimization in conditional parameter spaces.
306 *arXiv preprint arXiv:1409.4011*, 2014.
- 307 [4] Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabas Poczos, and Eric P Xing.
308 Neural architecture search with bayesian optimisation and optimal transport. In *Advances in*
309 *Neural Information Processing Systems (NeurIPS)*, 2018.
- 310 [5] Natalie Maus, Haydn Jones, Juston Moore, Matt J Kusner, John Bradshaw, and Jacob Gard-
311 ner. Local latent space bayesian optimization over structured inputs. In *Advances in Neural*
312 *Information Processing Systems (NeurIPS)*, 2022.
- 313 [6] Seunghun Lee, Jinyoung Park, Jaewon Chu, Minseo Yoon, and Hyunwoo J Kim. Latent
314 bayesian optimization via autoregressive normalizing flows. In *International Conference on*
315 *Learning Representations (ICLR)*, 2025.
- 316 [7] Linnan Wang, Rodrigo Fonseca, and Yuandong Tian. Learning search space partition for
317 black-box optimization using monte carlo tree search. In *Advances in Neural Information*
318 *Processing Systems (NeurIPS)*, 2020.
- 319 [8] Felix Berkenkamp, Andreas Krause, and Angela P Schoellig. Bayesian optimization with safety
320 constraints: safe and automatic parameter tuning in robotics. *Machine Learning*, 2023.
- 321 [9] Harold J Kushner. A new method of locating the maximum point of an arbitrary multipeak
322 curve in the presence of noise. *Journal of Basic Engineering*, 1964.
- 323 [10] Roman Garnett. *Bayesian optimization*. Cambridge University Press, 2023.
- 324 [11] David K Duvenaud, Hannes Nickisch, and Carl Rasmussen. Additive gaussian processes. In
325 *Advances in Neural Information Processing Systems (NeurIPS)*, 2011.
- 326 [12] Kirthevasan Kandasamy, Jeff Schneider, and Barnabás Póczos. High dimensional bayesian
327 optimisation and bandits via additive models. In *International Conference on Machine Learning*
328 *(ICML)*, 2015.
- 329 [13] Jacob Gardner, Chuan Guo, Kilian Weinberger, Roman Garnett, and Roger Grosse. Discovering
330 and exploiting additive structure for bayesian optimization. In *International Conference on*
331 *Artificial Intelligence and Statistics (AISTATS)*, 2017.
- 332 [14] Xiaoyu Lu, Alexis Boukouvalas, and James Hensman. Additive gaussian processes revisited.
333 In *International Conference on Machine Learning (ICML)*, 2022.
- 334 [15] Ziyu Wang, Frank Hutter, Masrour Zoghi, David Matheson, and Nando De Freitas. Bayesian
335 optimization in a billion dimensions via random embeddings. *Journal of Artificial Intelligence*
336 *Research*, 2016.
- 337 [16] Amin Nayebi, Alexander Munteanu, and Matthias Poloczek. A framework for bayesian
338 optimization in embedded subspaces. In *International Conference on Machine Learning*
339 *(ICML)*, 2019.
- 340 [17] Ben Letham, Roberto Calandra, Akshara Rai, and Eytan Bakshy. Re-examining linear em-
341 beddings for high-dimensional bayesian optimization. In *Advances in Neural Information*
342 *Processing Systems (NeurIPS)*, 2020.

- 343 [18] David Eriksson and Martin Jankowiak. High-dimensional bayesian optimization with sparse
344 axis-aligned subspaces. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 2021.
- 345 [19] Leonard Papenmeier, Luigi Nardi, and Matthias Poloczek. Increasing the scope as you learn:
346 Adaptive bayesian optimization in nested subspaces. In *Advances in Neural Information
347 Processing Systems (NeurIPS)*, 2022.
- 348 [20] David Eriksson, Michael Pearce, Jacob Gardner, Ryan D Turner, and Matthias Poloczek.
349 Scalable global optimization via local bayesian optimization. In *Advances in Neural Information
350 Processing Systems (NeurIPS)*, 2019.
- 351 [21] David Eriksson and Matthias Poloczek. Scalable constrained bayesian optimization. In *Internat-
352 ional Conference on Artificial Intelligence and Statistics (AISTATS)*, 2021.
- 353 [22] Natalie Maus, Zhiyuan Jerry Lin, Maximilian Balandat, and Eytan Bakshy. Joint composite
354 latent space bayesian optimization. In *International Conference on Machine Learning (ICML)*,
355 2024.
- 356 [23] Carl Hvarfner, Erik Orm Hellsten, and Luigi Nardi. Vanilla bayesian optimization performs
357 great in high dimensions. In *International Conference on Machine Learning (ICML)*, 2024.
- 358 [24] Zhitong Xu, Haitao Wang, Jeff M Phillips, and Shandian Zhe. Standard gaussian process is all
359 you need for high-dimensional bayesian optimization. In *International Conference on Learning
360 Representations (ICLR)*, 2025.
- 361 [25] ChangYong Oh, Efstratios Gavves, and Max Welling. Bock: Bayesian optimization with
362 cylindrical kernels. In *International Conference on Machine Learning (ICML)*, 2018.
- 363 [26] Colin Doumont, Donney Fan, Natalie Maus, Jacob R Gardner, Henry Moss, and Geoff Pleiss.
364 We still don't understand high-dimensional bayesian optimization. In *International Conference
365 on Artificial Intelligence and Statistics (AISTATS)*, 2026.
- 366 [27] Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian,
367 Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. The llama
368 3 herd of models. *arXiv preprint arXiv:2407.21783*, 2024.
- 369 [28] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu,
370 Chang Gao, Chengen Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv preprint
371 arXiv:2505.09388*, 2025.
- 372 [29] Gheorghe Comanici, Eric Bieber, Mike Schaekermann, Ice Pasupat, Noveen Sachdeva, Inderjit
373 Dhillon, Marcel Blistein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5: Pushing the
374 frontier with advanced reasoning, multimodality, long context, and next generation agentic
375 capabilities. *arXiv preprint arXiv:2507.06261*, 2025.
- 376 [30] Bernardino Romera-Paredes, Mohammadamin Barekatin, Alexander Novikov, Matej Balog,
377 M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,
378 Omar Fawzi, et al. Mathematical discoveries from program search with large language models.
379 *Nature*, 2024.
- 380 [31] Alexander Novikov, Ngân Vũ, Marvin Eisenberger, Emilien Dupont, Po-Sen Huang, Adam Zsolt
381 Wagner, Sergey Shirobokov, Borislav Kozlovskii, Francisco JR Ruiz, Abbas Mehrabian,
382 et al. Alphaevolve: A coding agent for scientific and algorithmic discovery. *arXiv preprint
383 arXiv:2506.13131*, 2025.
- 384 [32] Virginia Aglietti, Ira Ktena, Jessica Schrouff, Eleni Sgouritsa, Francisco Ruiz, Alan Malek,
385 Alexis Bellot, and Silvia Chiappa. Funbo: Discovering acquisition functions for bayesian
386 optimization with funsearch. In *International Conference on Machine Learning (ICML)*, 2025.
- 387 [33] Zhikai Zhao, Chuanbo Hua, Federico Berto, Kanghoon Lee, Zihan Ma, Jiachen Li, and Jinkyoo
388 Park. Trajevo: Trajectory prediction heuristics design via llm-driven evolution. In *Proceedings
389 of the AAAI Conference on Artificial Intelligence (AAAI)*, 2026.

- 390 [34] Sein Kim, Sangwu Park, Hongseok Kang, Wonjoong Kim, Jimin Seo, Yeonjun In, Kanghoon
391 Yoon, and Chanyoung Park. Self-evolverec: Self-evolving recommender systems with llm-based
392 directional feedback. *arXiv preprint arXiv:2602.12612*, 2026.
- 393 [35] Johannes Kirschner, Mojmir Mutny, Nicole Hiller, Rasmus Ischebeck, and Andreas Krause.
394 Adaptive and safe bayesian optimization in high dimensions via one-dimensional subspaces. In
395 *International Conference on Machine Learning (ICML)*, 2019.
- 396 [36] Lei Song, Ke Xue, Xiaobin Huang, and Chao Qian. Monte carlo tree search based variable
397 selection for high dimensional bayesian optimization. In *Advances in Neural Information*
398 *Processing Systems (NeurIPS)*, 2022.
- 399 [37] Erik Orm Hellsten, Carl Hvarfner, Leonard Papenmeier, and Luigi Nardi. High-dimensional
400 bayesian optimization with group testing. *arXiv preprint arXiv:2310.03515*, 2023.
- 401 [38] Samuel Daulton, David Eriksson, Maximilian Balandat, and Eytan Bakshy. Multi-objective
402 bayesian optimization over high-dimensional search spaces. In *Conference on Uncertainty in*
403 *Artificial Intelligence (UAI)*, 2022.
- 404 [39] Natalie Maus, Kaiwen Wu, David Eriksson, and Jacob Gardner. Discovering many diverse
405 solutions with bayesian optimization. In *International Conference on Artificial Intelligence and*
406 *Statistics (AISTATS)*, 2023.
- 407 [40] Paolo Ascia, Elena Raponi, Thomas Bäck, and Fabian Duddeck. Feasibility-driven trust region
408 bayesian optimization. In *International Conference on Automated Machine Learning (AutoML)*,
409 2025.
- 410 [41] Leonard Papenmeier, Matthias Poloczek, and Luigi Nardi. Understanding high-dimensional
411 bayesian optimization. In *International Conference on Machine Learning (ICML)*, 2025.
- 412 [42] Bahador Rashidi, Kerrick Johnstonbaugh, and Chao Gao. Cylindrical thompson sampling for
413 high-dimensional bayesian optimization. In *International Conference on Artificial Intelligence*
414 *and Statistics (AISTATS)*, 2024.
- 415 [43] Giang Ngo, Dat Phan Trong, Dang Nguyen, Sunil Gupta, and Svetha Venkatesh. Adaptive
416 acquisition selection for bayesian optimization with large language models. In *International*
417 *Conference on Learning Representations (ICLR)*, 2026.
- 418 [44] Wenhui Li, Niki van Stein, Thomas Bäck, and Elena Raponi. Llamea-bo: A large language
419 model evolutionary algorithm for automatically generating bayesian optimization algorithms.
420 *arXiv preprint arXiv:2505.21034*, 2025.
- 421 [45] Tennison Liu, Nicolás Astorga, Nabeel Seedat, and Mihaela van der Schaar. Large language
422 models to enhance bayesian optimization. In *International Conference on Learning Representa-*
423 *tions (ICLR)*, 2024.
- 424 [46] Richard Cornelius Suwandi, Feng Yin, Juntao Wang, Renjie Li, Tsung-Hui Chang, and Sergios
425 Theodoridis. Adaptive kernel design for bayesian optimization is a piece of cake with llms. In
426 *Advances in Neural Information Processing Systems (NeurIPS)*, 2025.
- 427 [47] Austin Tripp. We are underselling the modularity of Bayesian optimization. [https://www.
428 austintripp.ca/blog/2026-02-09-bo-modularity/](https://www.austintripp.ca/blog/2026-02-09-bo-modularity/), 2026. Accessed: 2026-04-14.
- 429 [48] Carl Edward Rasmussen. Gaussian processes in machine learning. In *Summer School on*
430 *Machine Learning*. Springer, 2003.
- 431 [49] Jonas Mockus. The application of bayesian methods for seeking the extremum. *Towards Global*
432 *Optimization*, 1998.
- 433 [50] Richard H Byrd, Peihuang Lu, Jorge Nocedal, and Ciyou Zhu. A limited memory algorithm for
434 bound constrained optimization. *SIAM Journal on Scientific Computing*, 1995.
- 435 [51] Christopher Williams and Carl Rasmussen. Gaussian processes for regression. In *Advances in*
436 *Neural Information Processing Systems (NeurIPS)*, 1995.

- 437 [52] Zi Wang, Clement Gehring, Pushmeet Kohli, and Stefanie Jegelka. Batched large-scale bayesian
438 optimization in high-dimensional spaces. In *International Conference on Artificial Intelligence*
439 *and Statistics (AISTATS)*, 2018.
- 440 [53] Kenan Šehić, Alexandre Gramfort, Joseph Salmon, and Luigi Nardi. Lassobench: A high-
441 dimensional hyperparameter optimization benchmark suite for lasso. In *International Confer-*
442 *ence on Automated Machine Learning (AutoML)*, 2022.
- 443 [54] Ibai Roman, Roberto Santana, Alexander Mendiburu, and Jose A Lozano. An experimental
444 study in adaptive kernel selection for bayesian optimization. *IEEE Access*, 2019.
- 445 [55] David Duvenaud, James Lloyd, Roger Grosse, Joshua Tenenbaum, and Ghahramani Zoubin.
446 Structure discovery in nonparametric regression through compositional kernel search. In
447 *International Conference on Machine Learning (ICML)*, 2013.
- 448 [56] Sebastian Ament, Samuel Daulton, David Eriksson, Maximilian Balandat, and Eytan Bakshy.
449 Unexpected improvements to expected improvement for bayesian optimization. In *Advances in*
450 *Neural Information Processing Systems (NeurIPS)*, 2023.
- 451 [57] Ilya M Sobol', Danil Asotsky, Alexander Kreinin, and Sergei Kucherenko. Construction and
452 comparison of high-dimensional sobol'generators. *Wilmott*, 2011.

453 **Appendix**

454 **A Task Details**

455 We evaluate our method on five standard high-dimensional BO benchmarks covering a range of
 456 problem types and dimensionalities.

457 **Rover** ($D = 100$). The Rover benchmark [52] is a trajectory optimization task in which a rover must
 458 navigate from a fixed start position to a fixed goal on a 2D terrain. The trajectory is parameterized by
 459 50 intermediate 2D waypoints, yielding a $D = 100$ -dimensional continuous search space in $[0, 1]^D$.
 460 The objective rewards proximity to the goal while penalizing path irregularity; higher values are
 461 better. The reward landscape contains many local optima corresponding to different trajectories.

462 **Mopta08** ($D = 124$). The Mopta08 benchmark is a structural car body engineering design problem
 463 [18]. The task is to minimize the mass of a car structural component while satisfying 68 physical safety
 464 and manufacturing constraints. The 124-dimensional search space encodes continuous geometric
 465 design parameters. Constraint violations are aggregated into a penalty term, making the effective
 466 objective landscape highly non-smooth. Lower values are better.

467 **Lasso-DNA** ($D = 180$). Lasso-DNA is from the LassoBench suite [53], a benchmark collection
 468 for high-dimensional hyperparameter optimization of regularized regression models on genomic
 469 data. The input is a $D = 180$ -dimensional vector of group-wise LASSO regularization coefficients;
 470 the objective is the mean squared prediction error on a validation set from a DNA microarray gene
 471 expression study (lower is better). The structured sparsity of the optimal regularizer makes this
 472 benchmark representative of practical bioinformatics hyperparameter search.

473 **SVM** ($D = 388$). The SVM benchmark [18] involves hyperparameter optimization of a Support
 474 Vector Machine classifier on a high-dimensional tabular dataset. The $D = 388$ -dimensional search
 475 space encodes regularization and kernel hyperparameters of the SVM. The objective is the test
 476 misclassification error (lower is better). The high sensitivity of SVM performance to specific
 477 hyperparameter regions makes this a particularly challenging task.

478 **Humanoid** ($D = 6392$). The Humanoid benchmark [7] is a continuous locomotion control task
 479 from the MuJoCo physics simulator. A simulated humanoid robot is controlled by a linear policy that
 480 maps the 376-dimensional state observation vector to 17 joint torques, yielding a weight matrix of
 481 size $17 \times 376 = D = 6392$ parameters. The objective is the cumulative reward accumulated over a
 482 fixed-length simulation episode (higher is better).

483 **Ant** ($D = 888$, **excluded from main evaluation**). We first considered the Ant locomotion bench-
 484 mark from MuJoCo, in which a quadruped robot is controlled by a linear policy mapping its state
 485 observations to joint torques. However, we excluded Ant from the main evaluation due to its abnormal
 486 reward structure. Specifically, we found out that the trivial zero-initialization baseline, which sets
 487 all policy parameters to $\mathbf{x} = \mathbf{0}$, achieves a cumulative reward of $+1000.40 \pm 4.40$. By contrast, all
 488 kernel-based GP-BO methods obtain dramatically lower or negative rewards, as shown in Table 4.
 489 This results makes the benchmark unsuitable for comparing kernel quality in BO.

Table 4: Final best values across base kernels compared to the zero-init-parameter policy baseline. Experiments are conducted with 4 random seeds.

Arm	Best Value Found
$\mathbf{x} = \mathbf{0}$ (zero-init-parameter policy)	1000.40 ± 4.40
RBF	$- 778.67 \pm 176.89$
Matérn52	$- 810.20 \pm 112.79$
RQ	-1011.04 ± 261.02
SL	$- 660.03 \pm 86.11$
BOCK	151.70 ± 358.00

490 **B Baseline Details**

491 We compare against 16 baselines across three categories: base kernels, search-based methods, and
 492 LLM-based methods.

493 **Base Kernels**

494 The following methods use the standard GP-BO setup with a single fixed kernel. For RBF,
 495 Matérn52, and RQ, we apply the data-dependent lengthscale prior of Hvarfner et al. [23]. Let
 496 $\mathbf{L} = \text{diag}(\ell_1, \dots, \ell_D)$ denote the ARD lengthscale matrix and $r_{\mathbf{L}} = \sqrt{(\mathbf{x} - \mathbf{x}')^\top \mathbf{L}^{-1} (\mathbf{x} - \mathbf{x}')}$.

- 497 • **RBF.** The squared exponential kernel with ARD:

$$k_{\text{RBF}}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{1}{2} r_{\mathbf{L}}^2\right). \quad (12)$$

- 498 • **Matérn52.** The Matérn kernel with smoothness $\nu = 5/2$ and ARD:

$$k_{\text{Matérn52}}(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{5} r_{\mathbf{L}} + \frac{5}{3} r_{\mathbf{L}}^2\right) \exp\left(-\sqrt{5} r_{\mathbf{L}}\right). \quad (13)$$

499 Models twice mean-square differentiable functions; often better suited than RBF to practical
 500 optimization landscapes.

- 501 • **Linear.** The dot-product kernel with an unconstrained scale parameter $v \in \mathbb{R}$:

$$k_{\text{Lin}}(\mathbf{x}, \mathbf{x}') = v \mathbf{x}^\top \mathbf{x}'. \quad (14)$$

- 502 • **Periodic.** The ARD periodic kernel with per-dimension lengthscales $\ell_i > 0$ and periods $p_i > 0$:

$$k_{\text{Per}}(\mathbf{x}, \mathbf{x}') = \exp\left(-2 \sum_{i=1}^D \frac{\sin^2(\pi(x_i - x'_i)/p_i)}{\ell_i^2}\right). \quad (15)$$

- 503 • **RQ (Rational Quadratic).** A scale mixture of RBF kernels with Gamma-distributed inverse
 504 lengthscales:

$$k_{\text{RQ}}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{r_{\mathbf{L}}^2}{2\alpha}\right)^{-\alpha}, \quad \alpha > 0. \quad (16)$$

505 The parameter α controls the relative contribution of short-range vs. long-range variation.

- 506 • **BOCK.** Bayesian Optimization with Cylindrical Kernels (BOCK) [25] decomposes each normal-
 507 ized input $\tilde{\mathbf{x}} = (\mathbf{x} - \mathbf{c})/R$ into a radial component $r = \|\tilde{\mathbf{x}}\|$ and an angular component $\hat{\mathbf{x}} = \tilde{\mathbf{x}}/r$.
 508 The radial coordinate is transformed by the Kumaraswamy CDF $\kappa(r; \alpha, \beta) = 1 - (1 - r^\alpha)^\beta$
 509 before applying a Matérn52 kernel, and the angular component is modeled by a polynomial
 510 kernel in the inner product:

$$k_{\text{BOCK}}(\mathbf{x}, \mathbf{x}') = k_{\text{Matérn52}}(\kappa(r), \kappa(r')) \cdot \sum_{p=0}^{P-1} w_p (\hat{\mathbf{x}}^\top \hat{\mathbf{x}}')^p, \quad w_p > 0. \quad (17)$$

511 This product structure separates the smooth radial variation from the angular geometry, mitigating
 512 boundary-seeking behavior.

- 513 • **SL (Spherical Linear).** The kernel of Doumont et al. [26] first applies ARD lengthscale scaling
 514 followed by a global scale g , then maps the result to the unit sphere via the stereographic
 515 projection $\psi : \mathbb{R}^D \rightarrow S^D$:

$$\psi(\mathbf{z}) = \frac{1}{1 + \|\mathbf{z}\|^2} \begin{pmatrix} 2\mathbf{z} \\ \|\mathbf{z}\|^2 - 1 \end{pmatrix} \in S^D. \quad (18)$$

516 Letting $\tilde{\mathbf{x}} = (\mathbf{x} - \mathbf{c}) / (\ell \cdot g)$, the kernel is a softmax-weighted mixture of a linear term on the
 517 sphere and a constant bias:

$$k_{\text{SL}}(\mathbf{x}, \mathbf{x}') = \lambda_1 \psi(\tilde{\mathbf{x}})^\top \psi(\tilde{\mathbf{x}}') + \lambda_0, \quad \lambda_0 + \lambda_1 = 1, \quad \lambda_0, \lambda_1 \geq 0. \quad (19)$$

518 The stereographic projection ensures that inputs at any distance from the origin are mapped to a
 519 bounded region on the sphere, avoiding the unbounded extrapolation that afflicts linear kernels
 520 in high dimensions.

521 **Search-based Methods**

522 All search-based baselines share the same initial population \mathcal{P}_0 as our method for a fair comparison.

- 523 • **Greedy Search.** At each BO iteration, it evaluates every kernel in \mathcal{P}_0 and selects the one with
524 the lowest Bayesian Information Criterion (BIC). No new kernels are created.
- 525 • **Compositional Search.** At each BO iteration, use greedy search to discover kernel structures
526 that best explain the observed data via pairwise additive and multiplicative compositions of
527 kernels in the current population. Following the paper, we use BIC as a selection metric and
528 directly use hyperparameters fitted in a single kernel to composite kernels without refitting. We
529 use the max depth of the composition as 3.
- 530 • **CAKE.** Context-Aware Kernel Evolution [46] uses an LLM to select and compose base kernels
531 conditioned on prior observations. Its search is restricted to additive and multiplicative composi-
532 tions of a fixed base set, and passing raw observations to the LLM induces context-length issues
533 in high dimensions. If the dimensions are too high, we truncate the observations to prevent
534 context limit violation. We use the max depth of the composition as 3. While CAKE uses
535 GPT-4o-mini for the LLM API, we evaluate the baseline with GPT-4o for a fair comparison.

536 **LLM-based Methods**

537 For all LLM-based baselines, we initialize the dataset with Sobol sequences of size of 20 as our
538 method for a fair comparison.

- 539 • **LLaMEA-BO (ARM, ATRBO, TREvol, TROpt, TRPareto).** Five end-to-end LLM-driven
540 BO algorithms produced by the LLaMEA-BO pipeline [44], which uses an LLM in an evo-
541 lutionary loop to generate and refine BO algorithm code. Each variant is a distinct algorithm
542 discovered by the pipeline. These methods may not use a GP surrogate and directly propose
543 candidates from prior observations.
- 544 • **LMABO.** An LLM-based BO method that adaptively selects the acquisition function at each
545 iteration [43]. For a fair comparison, we use GP with an RBF as the base kernel, following
546 Hvarfner et al. [23], and use L-BFGS-B as the optimizer.
- 547 • **LLAMBO (excluded from main evaluation).** LLAMBO [45] uses LLMs to propose and
548 evaluate candidates conditioned on historical evaluations, making its prompt size scale with
549 both history length and input dimensionality. On Rover ($D = 100$), the full history of 980
550 evaluations already requires approximately 600K input tokens, exceeding the 128k context
551 window of GPT-4o. Even after truncating the history to 20 observations, our GPT-4o probe costs
552 about \$1.05 per BO iteration, corresponding to about \$1,024 for a single full run. We therefore
553 exclude LLAMBO due to context-length and budget limitations.

554 **Additional Baselines: TuRBO and SAASBO.** We compare against two additional baselines,
555 TuRBO [20] and SAASBO [18], which are not included in the main table as they employ different
556 acquisition function maximization strategies and inference procedures (TuRBO uses a trust-region
557 approach and SAASBO uses sparse axis-aligned priors with MCMC-based inference). However, we
558 compare the results of those baselines as they are widely used in high-dimensional BO. As shown in
559 Table 5, Kernel Discovery outperforms both methods across all benchmarks.

Table 5: Comparison with additional baselines across standard benchmarks. SAASBO results use $T = 500$ due to its higher computational cost. Experiments are conducted with 4 random seeds.

Method	Rover (\uparrow) ($D = 100$)	Mopta08 (\downarrow) ($D = 124$)	Lasso-DNA (\downarrow) ($D = 180$)	SVM (\downarrow) ($D = 388$)	Humanoid (\uparrow) ($D = 6392$)
Kernel Discovery (Ours)	4.353 \pm 0.319	216.81 \pm 0.87	0.286 \pm 0.001	0.056 \pm 0.003	762.78 \pm 72.39
TuRBO	3.584 \pm 0.175	242.39 \pm 3.60	0.289 \pm 0.001	0.142 \pm 0.015	449.26 \pm 85.56
SAASBO ($T = 500$)	3.981 \pm 0.362	223.30 \pm 3.44	0.290 \pm 0.001	0.148 \pm 0.012	449.08 \pm 53.58

560 C Implementation Details

561 **Shared BO Setup.** For all methods except end-to-end LLM-based approaches (ARM, ATRBO,
562 TREvol, TROpt, TRPareto), we use a shared GP-based BO framework built on GPyTorch and
563 BoTorch. The surrogate is an exact Gaussian Process with homoscedastic Gaussian observation noise.
564 Kernel hyperparameters are optimized by maximizing the marginal log-likelihood (MLL) using
565 L-BFGS-B with 3 random restarts. Inputs are normalized to $[0, 1]^D$ and outputs are standardized to
566 zero mean and unit variance before fitting.

567 We use LogEI as the acquisition function and optimize it using L-BFGS-B with 4 random restarts
568 initialized from 512 Sobol quasi-random candidates. RAASP sampling is implemented by tuning
569 the option `sample_around_best=True`. The initial dataset \mathcal{D}_0 is generated by a Sobol sequence of
570 size $|\mathcal{D}_0| = 20$. The total evaluation budget is $T = 1000$ function queries.

571 All experiments use a batch BO setting with batch size $q = 20$. At each BO round, the acquisition
572 function is maximized to select a batch of q candidate points, which are evaluated simultaneously.
573 This yields $T/q = 50$ BO rounds per run. All baselines also use the same $q = 20$ batch size to ensure
574 a fair comparison. Please refer to Table 6 for the comparison of the performance of base kernels
575 under different batch sizes. All experiments are done with a single RTX NVIDIA L40S GPU.

576 **Full Prompt Templates.** We present the full prompt templates for each stage, mathematical form
577 discovery, code conversion, and composition, in Figure 8 to 10, respectively. In addition to generating
578 solely mathematical expressions or code blocks, we instruct LLM to justify its answer for exploring
579 its reasoning capabilities and the robustness of the pipeline.

```
1 You are an expert in Gaussian process kernel design for high-dimensional
2 Bayesian optimization.
3
4 TOP {len(parent_formulas_with_loss)} kernels in our population:
5 {context}
6
7 Generate {n_discoveries} DISTINCT, CREATIVE, mathematically innovative
8 DISCOVERED kernels.
9 Each discovery should be derived from one or more of the above kernels, but
10 with meaningful mathematical changes.
11 {highdim}{psd}{simplicity}
12 OUTPUT FORMAT:
13 Return EXACTLY {n_discoveries} kernel descriptions. Enclose EACH in a ““
14 formula block:
15
16 ““formula
17 KERNEL: [name]
18
19 PARAMETERS:
20 - [name] ([shape], [constraint]): [description]
21
22 INPUT TRANSFORM:
23 [step-by-step math]
24
25 COVARIANCE FUNCTION:
26 [k(x1, x2) = ...]
27
28 PSD GUARANTEE:
29 [why this kernel is PSD]
30 ““
31
32 Each MUST be fundamentally distinct (different transform, different
33 covariance structure, different mathematical idea).””
```

Figure 8: Mathematical form discovery prompt template

```

1 You are a GPyTorch kernel code generator. Convert this kernel formula into a
  GPyTorch EvolvedKernel class.
2
3 {formula}
4
5 REFERENCE TEMPLATE A - Distance-based kernel {RBF code}
6
7 REFERENCE TEMPLATE B - Feature-map kernel {SL code}
8
9 CRITICAL RULES:
10 1. __init__ SIGNATURE: '__init__' MUST accept 'ard_num_dims: int' as the ONLY
    required parameter. The system calls 'EvolvedKernel(ard_num_dims=D)'. All
    other args MUST have defaults. NEVER add required args like 'q: int', 'center
    : Tensor', 'M: int'.
11 2. FORWARD SIGNATURE: 'def forward(self, x1, x2, diag=False, **params)' -
    always accept **params.
12 3. BATCH-SAFE SHAPES (CRITICAL): The kernel is tested with these exact shapes
    - your code MUST handle ALL of them:
13   - 2D: x1=(5, D) vs x2=(1, D) -> output must be (5, 1), NOT (5, 5)
14   - 2D: x1=(3, D) vs x2=(7, D) -> output must be (3, 7)
15   - 3D batched (BoTorch acqf optimization): x1=(1, 4, D) vs x2=(1, 3, D) ->
    output must be (1, 4, 3)
16   Therefore: ALL dim indexing must be relative (-1 for D, -2 for N). NEVER
    use '.size(0)', '.expand(x.size(0), ...)', or '.shape[0]'. For broadcasting
    scalars to match batch+N dims, use 'torch.ones_like(x[..., :1])'.
17 4. OUTPUT SHAPE: return (... , N1, N2) for ANY N1 != N2. If your output is (N1
    , N1) instead of (N1, N2), you have a bug.
18 5. DISTANCE: 'torch.cdist(x1s, x2s, p=2)' - NEVER 'x1.unsqueeze(-2) - x2.
    unsqueeze(-3)' (OOM).
19 6. INNER PRODUCT: 'x1 @ x2.transpose(-1, -2)'.
20 7. DEVICE: new tensors in forward use 'device=x1.device'. Use '
    register_buffer' in __init__.
21 8. NO IN-PLACE: 'x = x + y' not 'x += y'. 'x = x.clamp(...)' not 'x.clamp_
    (...)'
22 9. NO DETACH/NUMPY: never '.detach()', '.numpy()', '.item()', 'torch.no_grad
    ()' in forward.
23 10. NUMERICAL: '.clamp(min=1e-15)' before sqrt/log, '.clamp(max=20.0)' for
    exp.
24 11. LENGTHSCALE: if 'has_lengthscale=True', do NOT 'register_parameter("
    raw_lengthscale", ...)'
25 12. PARAM ORDER: 'register_parameter(name, ...)' BEFORE 'register_constraint(
    name, ...)'
26 13. DIAG: 'if diag: return covar.diagonal(dim1=-2, dim2=-1)'.
27 14. IMPORTS: include 'import math' if using math.sqrt/pi/log.
28 15. NO SELF-REASSIGNMENT in forward(): never 'self.x = ...' - use local
    variables.
29 16. ONLY USE EXISTING gpytorch classes.
30 17. NO SUB-KERNELS: do NOT instantiate gpytorch kernel objects (e.g. '
    MaternKernel()', 'RBFKernel()') as sub-components. They produce
    LazyEvaluatedKernelTensor with unpredictable shapes. Instead, implement the
    formula directly (e.g. for Matern-5/2: '(1 + sqrt(5*d + 5/3*d**2)) * exp(-sqrt(5
    *d))').
31 18. NO DATA-DEPENDENT BOUNDS in forward(): never compute bounds/center from
    the input data (e.g. 'x1.max(dim=-2)' or 'x1.min(dim=-2)'). These change with
    N and break shape invariance. Store fixed bounds via 'register_buffer' in
    __init__, or use constants (e.g. 0 and 1).
32 19. SCALAR PAIRWISE DISTANCE: when computing distance between scalar values (
    e.g. warped radii of shape (... , N, 1)), use 'torch.cdist(r1, r2, p=2)' which
    gives (... , N1, N2). NEVER use 'r1 - r2' or 'torch.abs(r1 - r2)' - this
    gives (... , N1, 1) not (... , N1, N2) and breaks cross-pair computation.
33
34 Return ONLY the Python code in a "'python block.'"

```

Figure 9: Code conversion prompt template.

```

1 You are combining GP kernels for high-dimensional Bayesian optimization.
2
3 Population of {len(population_formulas_with_loss)} kernels:
4 {context}
5
6 Generate {n_compositions} DISTINCT composed kernels. For EACH composition:
7 - Pick exactly 2 kernels from the population that COMPLEMENT each other
  mathematically
8 - Combine them: sum their covariance functions, multiply them, share a
  transform with different distance metrics, or any other valid composition
9 - Explain WHY these two kernels complement each other
10 {psd}{simplicity}
11 OUTPUT FORMAT:
12 Return EXACTLY {n_compositions} kernel descriptions. Enclose EACH in a ““
  formula block:
13
14 ““formula
15 KERNEL: [name]
16 COMPOSED FROM: [Kernel X] + [Kernel Y] (or *, or other operation)
17 WHY: [1 sentence on why these complement each other]
18
19 PARAMETERS:
20 - [param] ([shape], [constraint]): [description]
21
22 INPUT TRANSFORM:
23 [step-by-step math]
24
25 COVARIANCE FUNCTION:
26 k(x1, x2) = ...
27
28 PSD GUARANTEE:
29 [why this kernel is PSD]
30 ““””

```

Figure 10: Composition prompt template.

580 **Kernel Validation Failure Rates.** The failure probability in Tables 2 and 3 denotes the empirical
581 fraction of LLM-generated kernels rejected by our structural validator. A generated kernel k_{new} is
582 accepted only when

$$\mathcal{V}(k_{\text{new}}) = \mathcal{V}_{\text{agn}}(k_{\text{new}}) \wedge \mathcal{V}_{\text{psd}}(k_{\text{new}}) = 1.$$

583 Here, \mathcal{V}_{agn} checks whether the generated kernel implementation is dimension- and shape-agnostic, as
584 required by the BO pipeline. Concretely, we instantiate the kernel with different input dimensionalities
585 and evaluate its forward function on several synthetic input configurations, including self-covariance
586 inputs (N, D) , cross-covariance inputs (N_1, D) and (N_2, D) with $N_1 \neq N_2$, and batched inputs used
587 during acquisition optimization. The kernel is rejected if any call raises a runtime error, assumes a
588 hard-coded dimensionality, fails to broadcast correctly, or returns a covariance tensor with an invalid
589 shape. The PSD check \mathcal{V}_{psd} then evaluates the kernel on random inputs and tests whether the resulting
590 Gram matrix is numerically positive semi-definite via a Cholesky decomposition with a small jitter.
591 We report the fraction of generated kernels for which at least one of these checks fails, aggregated
592 over BO iterations and random seeds.

593 **Effect of Batch Size (q).** The main experiments use $q = 20$, reducing the number of BO rounds to
594 $T/q = 50$. To assess whether this setting disproportionately favors our method, we compare base
595 kernels under $q = 20$ vs. $q = 1$ on the SVM benchmark. Table 6 shows that the relative ordering of
596 kernels is largely preserved across both settings, and that the optimal kernel (RBF in the $q = 1$ setting)
597 is the same as in the $q = 20$ setting. This confirms that the batch setting does not systematically alter
which kernels are competitive, and that our comparison is not confounded by the choice of q .

Table 6: Comparison of final best values across base kernels under different oracle query batch sizes.

	RBF	Matérn52	RQ	BOCK	SL	Linear	Periodic
$q = 20$	0.061 ± 0.003	0.062 ± 0.002	0.069 ± 0.014	0.068 ± 0.006	0.111 ± 0.006	0.226 ± 0.003	0.226 ± 0.001
$q = 1$	0.078 ± 0.003	0.117 ± 0.021	0.078 ± 0.007	0.070 ± 0.009	0.064 ± 0.002	0.226 ± 0.002	0.223 ± 0.002

598

599 **D Algorithm Details**

600 **Discard Logic.** The difference between the standard evolutionary pipeline and our task is that
 601 one kernel that suits a certain round may not suit the other rounds, as the data keeps changing. To
 602 mitigate this, we discard the kernel if it is selected but does not improve the best value found so far.
 603 For kernels in the initial population, which are the basic building blocks of our discovery pipeline,
 604 we discard them if they are selected and no improvement is observed after $P = 3$ iterations. If all
 605 kernels in the population are discarded, we reset the population to the initial population.

606 **Humanoid Benchmark.** In the humanoid benchmark, the dimension is too high, and we observed
 607 that even with the LOO-CRPS evaluation metric, it sometimes leads to selecting too complex kernels
 608 that may not lead to an improvement. To remedy this, we regularize the LOO-CRPS score by the
 609 BIC metric as follows:

$$\text{LOO-CRPS-BIC}(\mathbf{K}_k, \mathbf{y}) = \frac{1}{n} \sum_{i=1}^n \text{CRPS}(\mathcal{N}(\mu_{-i}, \sigma_{-i}^2), y_i) + \frac{|\boldsymbol{\theta}_k| \log n}{n}, \quad (20)$$

610 where $|\boldsymbol{\theta}_k|$ is the number of hyperparameters in the kernel k . We found that this regularization works
 well, especially in extremely high dimensions, $D \gg n$.

Algorithm 1 Kernel Discovery for High-Dimensional BO

Require: Objective f , initial dataset \mathcal{D}_0 , population size N , BO budget T , initial-population patience
 P , LLM \mathcal{M} with prompts $\rho_{\text{disc}}, \rho_{\text{conv}}, \rho_{\text{comp}}$, verifier \mathcal{V}

- 1: Initialize $\mathcal{P}_0 \leftarrow \{k_{\text{RBF}}, k_{\text{Mat52}}, k_{\text{RQ}}, k_{\text{BOCK}}, k_{\text{SL}}\}$, $\mathcal{P} \leftarrow \mathcal{P}_0$
- 2: $\mathcal{D} \leftarrow \mathcal{D}_0$, $y^* \leftarrow \max_{(\mathbf{x}, y) \in \mathcal{D}_0} y$, $\text{fail}[k] \leftarrow 0$ for all $k \in \mathcal{P}_0$
- 3: **for** $t = 0, 1, \dots, T - 1$ **do**
- 4: **// Discovery stage**
- 5: $m^{\text{new}} \sim \mathcal{M}(\{m_i\}_{k_i \in \mathcal{P}}; \rho_{\text{disc}})$, $c^{\text{new}} \leftarrow \mathcal{M}(m^{\text{new}}; \rho_{\text{conv}})$, $k^{\text{new}} \leftarrow (m^{\text{new}}, c^{\text{new}})$
- 6: **if** $\mathcal{V}(k^{\text{new}}) = 1$ **and** GP fitting time ≤ 60 s **then**
- 7: $\mathcal{P} \leftarrow \mathcal{P} \cup \{k^{\text{new}}\}$
- 8: **end if**
- 9: **// Composition stage**
- 10: $m^{\text{comp}} \sim \mathcal{M}(\{m_{(i)}\}_{i=1}^N; \rho_{\text{comp}})$, $c^{\text{comp}} \leftarrow \mathcal{M}(m^{\text{comp}}; \rho_{\text{conv}})$, $k^{\text{comp}} \leftarrow (m^{\text{comp}}, c^{\text{comp}})$
- 11: **if** $\mathcal{V}(k^{\text{comp}}) = 1$ **and** GP fitting time ≤ 60 s **then**
- 12: $\mathcal{P} \leftarrow \mathcal{P} \cup \{k^{\text{comp}}\}$
- 13: **end if**
- 14: **// Kernel selection and BO step**
- 15: $k^* \leftarrow \arg \min_{k \in \mathcal{P}} \text{LOO-CRPS}(\mathbf{K}_k, \mathbf{y})$
- 16: Fit GP with k^* on \mathcal{D} ; $\mathbf{x}_{t+1} \leftarrow \arg \max_{\mathbf{x}} \alpha(\mathbf{x}; \text{GP})$
- 17: $y_{t+1} \leftarrow f(\mathbf{x}_{t+1})$; $\mathcal{D} \leftarrow \mathcal{D} \cup \{(\mathbf{x}_{t+1}, y_{t+1})\}$
- 18: **// Discard logic**
- 19: **if** $y_{t+1} \leq y^*$ **then**
- 20: $\text{fail}[k^*] \leftarrow \text{fail}[k^*] + 1$
- 21: **if** $k^* \notin \mathcal{P}_0$ **or** $\text{fail}[k^*] \geq P$ **then**
- 22: $\mathcal{P} \leftarrow \mathcal{P} \setminus \{k^*\}$
- 23: **end if**
- 24: **end if**
- 25: **if** $\mathcal{P} = \emptyset$ **then**
- 26: $\mathcal{P} \leftarrow \mathcal{P}_0$; $\text{fail}[k] \leftarrow 0$ for all $k \in \mathcal{P}_0$
- 27: **end if**
- 28: **// Population update**
- 29: $\mathcal{P} \leftarrow \text{Top-N}(\mathcal{P}; \text{LOO-CRPS}(\cdot, \mathbf{y}))$, $y^* \leftarrow \max(y^*, y_{t+1})$
- 30: **end for**
- 31: **return** $\arg \max_{(\mathbf{x}, y) \in \mathcal{D}} y$

611

612 **E Extended Experiment Results**

613 **Full Learning Curves.** Figure 11 extends Figure 4 to show the full learning curves for all 16
 614 baselines across all five benchmarks.

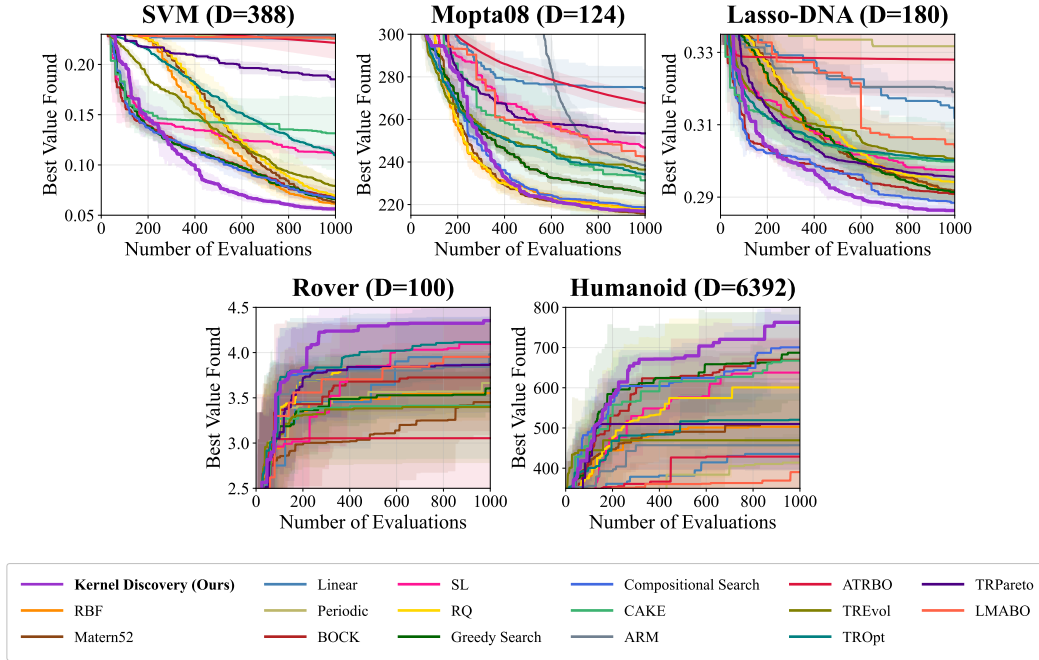


Figure 11: Full learning curves for all methods across five benchmarks.

615 **F Time Complexity Analysis**

616 We analyze the per-iteration wall-clock time of Kernel Discovery and compare it against all baselines
 617 in Table 7. To systematically identify the bottleneck of our running time, we decompose each iteration
 618 into the following phases:

- 619 • **GP Fitting:** Fitting all kernels in the current population in parallel using a process pool. The kernel
 620 selection score (e.g., LOO-CRPS) is computed within this phase, so the cost of choosing the best
 621 candidate is already included.
- 622 • **LLM Inference:** Two API calls per iteration - one for generating a novel mathematical kernel
 623 form and one for converting that form into executable code.
- 624 • **Kernel Validation:** Lightweight structural checks on each newly generated kernel before it is
 625 admitted to the population.
- 626 • **Search:** An additional kernel search procedure that applies only to search-based methods. Com-
 627 positional Search enumerates pairwise additive and multiplicative compositions, while CAKE
 628 employs an LLM-guided composition with a BAKER-based selection procedure. Our method and
 629 Greedy Search select the best candidate directly from the evaluated pool, so this cost is negligible.
- 630 • **Acquisition Function Optimization:** Maximizing the acquisition function via L-BFGS-B to obtain
 631 the next query point.
- 632 • **Oracle Evaluation:** Querying the black-box benchmark function at the selected point, shared
 633 identically across all methods.

634 For LLM-based end-to-end methods, the pipeline does not follow the standard BO decomposition,
 635 so we report only the total per-iteration time and oracle evaluation cost. As shown in Table 7, our
 636 method incurs longer time than fixed base kernel methods but achieves comparable total running time
 637 to search-based methods. Among LLM-based baselines, TREvol and TRPareto require significantly
 638 more time per iteration. We also observe that LLM inference is the primary bottleneck in our pipeline,
 639 partially due to the two-stage approach that issues separate API calls for formula generation and
 640 code conversion. While this decomposition is critical for producing functionally diverse kernels
 (Section 6), developing more efficient inference strategies is a promising direction for future work.

Table 7: Per-iteration wall-clock time breakdown (mean \pm std, seconds) on the SVM benchmark with 4 random seeds. For LLM-based methods, individual phases are not separately tracked; we report only the total time and oracle evaluation time.

Method	GP Fit	LLM	Validation	Search	Acqf Opt	Oracle	Total
Base Kernels							
RBF	0.68 \pm 0.02	N/A	N/A	0.00 \pm 0.00	2.70 \pm 0.97	0.28 \pm 0.04	5.18 \pm 1.00
Matérn52	0.69 \pm 0.04	N/A	N/A	0.00 \pm 0.00	2.76 \pm 0.99	0.28 \pm 0.04	5.26 \pm 1.00
Linear	0.95 \pm 0.58	N/A	N/A	0.00 \pm 0.00	3.14 \pm 0.88	0.29 \pm 0.04	6.49 \pm 2.38
Periodic	0.85 \pm 0.20	N/A	N/A	0.00 \pm 0.00	1.35 \pm 1.52	0.25 \pm 0.04	4.28 \pm 1.40
BOCK	1.14 \pm 0.44	N/A	N/A	0.00 \pm 0.00	4.03 \pm 1.62	0.31 \pm 0.06	7.02 \pm 1.87
SL	1.41 \pm 0.85	N/A	N/A	0.00 \pm 0.00	3.57 \pm 1.10	0.30 \pm 0.04	6.83 \pm 1.69
Search-based							
Greedy Search	5.38 \pm 1.80	N/A	N/A	0.00 \pm 0.00	5.70 \pm 0.56	0.35 \pm 0.05	11.51 \pm 2.04
Compositional Search	4.90 \pm 1.55	N/A	N/A	13.80 \pm 6.78	5.88 \pm 1.33	0.34 \pm 0.05	20.70 \pm 8.87
CAKE	7.15 \pm 3.61	9.09 \pm 3.23	N/A	17.44 \pm 3.73	19.34 \pm 4.02	0.35 \pm 0.06	37.14 \pm 4.47
LLM-based							
ARM	N/A	N/A	N/A	N/A	N/A	0.30 \pm 0.04	4.54 \pm 0.02
ATRBO	N/A	N/A	N/A	N/A	N/A	0.28 \pm 0.05	14.51 \pm 0.06
TREvol	N/A	N/A	N/A	N/A	N/A	0.31 \pm 0.02	209.88 \pm 3.40
TROpt	N/A	N/A	N/A	N/A	N/A	0.29 \pm 0.05	3.67 \pm 0.10
TRPareto	N/A	N/A	N/A	N/A	N/A	0.30 \pm 0.02	66.29 \pm 0.26
LMABO	3.56 \pm 0.50	1.92 \pm 0.45	N/A	0.00 \pm 0.00	1.92 \pm 0.28	0.27 \pm 0.04	5.76 \pm 0.61
Kernel Discovery (Ours)	4.02 \pm 1.22	14.24 \pm 2.05	1.90 \pm 0.78	0.00 \pm 0.00	4.56 \pm 0.11	0.30 \pm 0.01	28.05 \pm 1.81

641

642 G Direct Code Generation Analysis

643 **Details on Cosine Distance Metric.** To quantify the functional diversity of generated kernels, we
 644 compare their induced Gram matrices on a shared reference input set. For each benchmark and
 645 each method, we collect all valid LLM-generated kernels produced during the BO run. We then
 646 sample a fixed reference batch $\mathbf{X}_{\text{ref}} \in [0, 1]^{N \times D}$ with $N = 80$, where D is the dimensionality of the
 647 corresponding benchmark. The same reference batch is used for all kernels within a benchmark.

648 For each generated kernel k_i , we compute its Gram matrix $K_i = k_i(\mathbf{X}_{\text{ref}}, \mathbf{X}_{\text{ref}}) \in \mathbb{R}^{N \times N}$, and
 649 flatten it into a vector $\mathbf{v}_i = \text{vec}(K_i) \in \mathbb{R}^{N^2}$. We then compute the mean pairwise cosine similarity
 650 among all generated kernels in the pool \mathcal{P} :

$$\bar{s} = \frac{2}{|\mathcal{P}|(|\mathcal{P}| - 1)} \sum_{i < j} \frac{\mathbf{v}_i^\top \mathbf{v}_j}{\|\mathbf{v}_i\|_2 \|\mathbf{v}_j\|_2}.$$

651 Finally, we report the cosine distance as $\text{CosineDistance} = 1 - \bar{s}$. A higher value indicates that the
 652 generated kernels induce more diverse covariance structures on the same reference inputs.

653 **Performance Comparison.** We further compare direct code generation with our two-stage gener-
 654 ation under the same experimental setup. As shown in Table 8, the two-stage approach achieves
 655 better final performance on all benchmarks. This suggests that generating kernels through an explicit
 656 mathematical formulation step leads to more functionally meaningful and consistently effective kernel
 candidates than directly generating code.

Table 8: Performance comparison between direct code generation and our two-stage generation. Results are averaged over 4 random seeds. Bold denotes the best entry in each column.

Method	Rover (\uparrow) ($D = 100$)	Mopta08 (\downarrow) ($D = 124$)	Lasso-DNA (\downarrow) ($D = 180$)	SVM388 (\downarrow) ($D = 388$)	Humanoid (\uparrow) ($D = 6392$)
Two-Stage Generation (Ours)	4.353 \pm 0.319	216.81 \pm 0.87	0.286 \pm 0.001	0.056 \pm 0.003	762.78 \pm 72.39
Direct Code Generation	3.872 \pm 0.321	217.32 \pm 0.75	0.288 \pm 0.001	0.057 \pm 0.004	651.85 \pm 85.93

657

658 **Additional Examples of Functional Redundancy in Direct Code Generation.** As shown in
 659 Figure 5a, direct code generation often produces syntactically different implementations that cor-
 660 respond to the same kernel function. This is a common pitfall of searching directly in code space:
 661 changes such as variable renaming, syntax-level variation, or algebraic rewriting can be counted as
 662 new candidates. In Figures 12 to 14, we provide examples of this behavior. Across these pairs, the
 two forward() implementations look different at the code level but share the same kernel function.

<pre>def forward(self, x1, x2, diag=False, ** params): x1s = x1_scaled / self.lengthscale x2s = x2_scaled / self.lengthscale dist_sq = torch.cdist(x1s, x2s, p=2). pow(2) covar1 = torch.exp(-0.5 * dist_sq). clamp(min=1e-15) covar2 = (1.0 + dist_sq / (2.0 * alpha)) .pow(-alpha).clamp(min=1e-15) covar = covar1 + covar2 if diag: return covar.diagonal(dim1=-2, dim2=-1) return covar</pre>	<pre>def forward(self, x1, x2, diag=False, ** params): x1s = x1 / self.lengthscale x2s = x2 / self.lengthscale dist_sq = torch.cdist(x1s, x2s, p=2). pow(2) covar_exp = torch.exp(-0.5 * dist_sq). clamp(min=1e-15) covar_matern = (1 + dist_sq / (2 * self. alpha)).pow(-self.alpha).clamp(min=1e-15) covar_combined = covar_exp + covar_matern if diag: return covar_combined.diagonal(dim1=-2, dim2=-1) return covar_combined</pre>
---	--

Figure 12: Functional redundancy example from direct code generation. Both snippets compute $k(x, x') = \exp(-d^2/2) + (1 + d^2/(2\alpha))^{-\alpha}$. Green denotes variable renaming, and red denotes syntax-level variation.

663

```

def forward(self, x1, x2, diag=False, **
params):
    x1s = x1 / self.lengthscale
    x2s = x2 / self.lengthscale
    sq_dist = torch.cdist(x1s, x2s).pow(2)

    covar1 = torch.exp(-0.5 * sq_dist)
    sqrt5_d = math.sqrt(5) * torch.cdist(
x1s, x2s).clamp(min=1e-15)

    covar2 = (1.0 + sqrt5_d + (5.0 / 3.0) *
sq_dist) * torch.exp(-sqrt5_d)

    covar = self.alpha * covar1 + (1 -
self.alpha) * covar2

    if diag:
        return covar.diagonal(dim1=-1,
dim2=-2)
    return covar

```

```

def forward(self, x1, x2, diag=False, **
params):
    x1s = x1 / self.lengthscale
    x2s = x2 / self.lengthscale
    dist2 = torch.cdist(x1s, x2s).pow(2)

    exp_neg_half_dist2 = torch.exp(-0.5 *
dist2)
    scaled_dist = torch.cdist(x1s, x2s).
clamp(min=1e-15)
    covar1 = exp_neg_half_dist2
    covar2 = (1.0 + math.sqrt(5) *
scaled_dist + (5.0 / 3.0) * dist2) * torch
.exp(-math.sqrt(5) * scaled_dist)

    covar = self.alpha * covar1 + (1 -
self.alpha) * covar2

    if diag:
        return covar.diagonal(dim1=-1,
dim2=-2)
    return covar

```

Figure 13: Functional redundancy example from direct code generation. Both snippets compute $k(x, x') = \alpha \exp(-d^2/2) + (1 - \alpha)(1 + \sqrt{5}d + \frac{5}{3}d^2) \exp(-\sqrt{5}d)$. Green denotes variable renaming, and red denotes syntax-level variation.

```

def forward(self, x1, x2, diag=False, **
params):
    x1s = x1 / self.lengthscale
    x2s = x2 / self.lengthscale

    dist2 = torch.cdist(x1s, x2s, p=2).pow
(2)

    radial_covar = torch.exp(-dist2 / self
.alpha).clamp(max=1e15)
    negative_power = (1 + dist2).pow(-self
.beta).clamp(min=1e-15)

    covar = radial_covar * negative_power

    if diag:
        return covar.diagonal(dim1=-2,
dim2=-1)
    return covar

```

```

def forward(self, x1, x2, diag=False, **
params):
    device = x1.device
    x1s = x1 / self.lengthscale.to(device)
    x2s = x2 / self.lengthscale.to(device)

    dist2 = torch.cdist(x1s, x2s, p=2).pow
(2)

    radial_covar = torch.exp(-dist2 / self
.alpha.to(device))
    reciprocal_power = 1.0 / (1.0 + dist2).
pow(self.beta.to(device))

    combined_covar = radial_covar *
reciprocal_power

    if diag:
        return combined_covar.diagonal(
dim1=-2, dim2=-1)
    return combined_covar

```

Figure 14: Functional redundancy example from direct code generation. Both snippets compute $k(x, x') = \exp(-d^2/\alpha)(1 + d^2)^{-\beta}$. Green denotes variable renaming, red denotes syntax-level variation, and blue denotes algebraic rewriting.

664 **H Extended Ablation Studies**

665 **Component Ablation on Other Benchmarks** In Figure 6, we conduct an ablation study on each
 666 component of our method on the SVM benchmark. In this section, we summarize the ablation studies
 667 on other benchmarks, extending several combinations of our methods. As shown in Table 9, removing
 668 any single component consistently degrades performance. Notably, omitting the Discovery stage or
 669 both the Discovery and the Initialization stage results in the most severe performance drops, second
 670 only to removing all components. This suggests that these components are crucial design choices for
 671 the efficacy of the kernel discovery pipeline in high-dimensional BO.

Table 9: Ablation study of Kernel Discovery components across standard benchmarks. D denotes the dimensionality of the task. **Bold** denotes the best entry in the column. Experiments are conducted with 4 random seeds.

Method	Rover (\uparrow) ($D = 100$)	Mopta08 (\downarrow) ($D = 124$)	Lasso-DNA (\downarrow) ($D = 180$)	SVM388 (\downarrow) ($D = 388$)	Humanoid (\uparrow) ($D = 6392$)	Avg Rank (\downarrow)
Kernel Discovery (Ours)	4.353 \pm 0.319	216.81 \pm 0.87	0.286 \pm 0.001	0.056 \pm 0.003	762.78 \pm 72.39	1.0 / 8
w.o. Init	3.837 \pm 0.390	217.60 \pm 2.34	0.287 \pm 0.001	0.071 \pm 0.011	608.40 \pm 63.92	4.9 / 8
w.o. Discovery	3.810 \pm 0.448	220.42 \pm 0.95	0.288 \pm 0.001	0.064 \pm 0.003	720.41 \pm 41.77	5.2 / 8
w.o. Eval	3.928 \pm 0.402	219.42 \pm 2.37	0.289 \pm 0.001	0.057 \pm 0.003	712.91 \pm 99.62	4.0 / 8
w.o. Init + Discovery	4.153 \pm 0.471	217.47 \pm 0.65	0.293 \pm 0.003	0.079 \pm 0.009	557.07 \pm 117.39	5.4 / 8
w.o. Init + Eval	4.132 \pm 0.493	218.70 \pm 1.55	0.291 \pm 0.004	0.065 \pm 0.010	610.89 \pm 73.68	5.0 / 8
w.o. Discovery + Eval	3.908 \pm 0.415	220.20 \pm 0.59	0.287 \pm 0.001	0.063 \pm 0.002	705.71 \pm 10.51	4.3 / 8
w.o. Init + Discovery + Eval	3.905 \pm 0.321	218.56 \pm 1.83	0.290 \pm 0.001	0.087 \pm 0.009	569.78 \pm 88.46	6.2 / 8

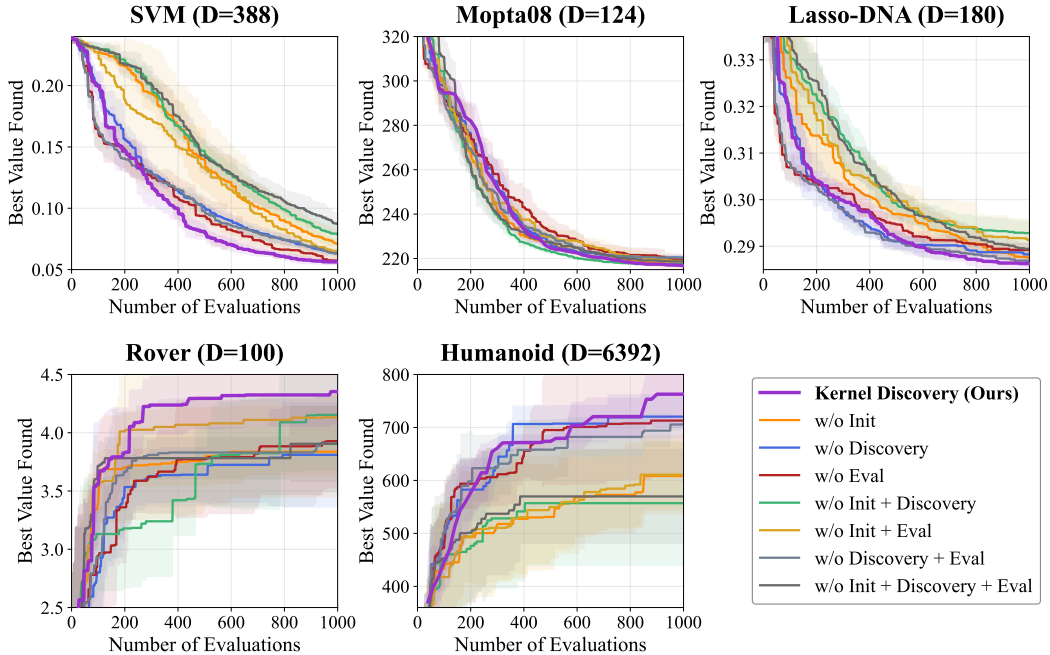


Figure 15: Ablation studies on each component of our method across other benchmarks.

672 **More Ablations on Base Population.** For a fair comparison, we also initialize the base populations
 673 of the search-based methods, Compositional Search and CAKE, in the same way as in our method.
 674 We also conduct ablation studies on the base population for those baselines. As shown in Figure 16,
 675 the performance of those methods significantly degrades when we remove BOCK and SL kernels
 676 from the base population. While our method also exhibits low sample efficiency without those kernels,
 677 it consistently improves the performance through discovering novel kernel structures.

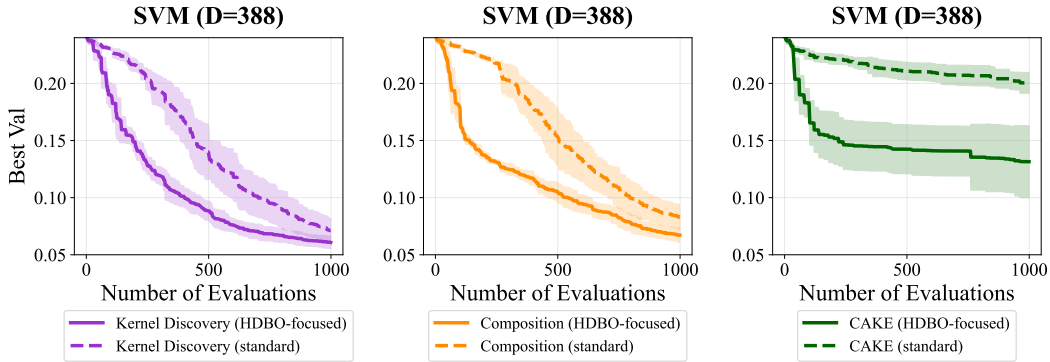


Figure 16: Ablation studies on the base population for Compositional Search and CAKE baselines.

678 **More Ablations on Evaluation Metric.** To follow the original implementation, we use the marginal
 679 log-likelihood (MLL) to select a kernel for Compositional Search and BIC-Acquisition Kernel
 680 Ranking (BAKER) to select a kernel for CAKE. To analyze the effect of the proposed evaluation
 681 metric. We compare performance across different evaluation metrics for both our kernel discovery
 682 pipeline and search-based baselines. As shown in Figure 17, performance degrades when we replace
 683 the evaluation metric with MLL, which tends to favor overly complex kernels that are likely to overfit
 684 to the current dataset. We also observe that LOO-CRPS improves the performance of search-based
 685 baselines, indicating that it is a powerful metric for selecting kernels in high-dimensional BO.

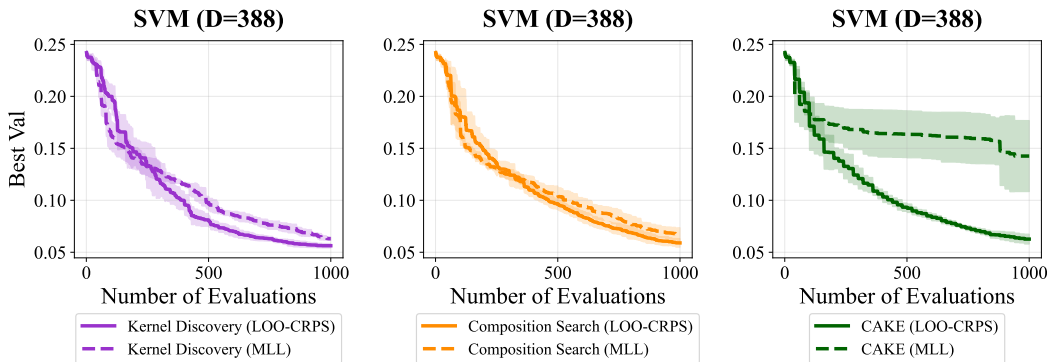


Figure 17: Ablation studies on the evaluation metric for Compositional Search and CAKE baselines.

686 **Representative Kernel Validation Failures.** For \mathcal{V}_{agn} , we test whether a generated implementation
 687 satisfies the covariance shape requirements. GP fitting only requires the self-covariance matrix
 688 $K(\mathbf{X}, \mathbf{X})$, but posterior prediction and acquisition evaluation also require the generally rectangular
 689 cross-covariance matrix $K(\mathbf{X}_*, \mathbf{X})$. Accordingly, a common failure mode is that the generated code
 690 works for self-covariance calls but breaks under cross-covariance inputs. Figure 18 presents two
 691 representative examples.

```
def forward(self, x1, x2, diag=False, **
params):
    # x1: (N1, D), x2: (N2, D)
    # ... input normalization, sphere
    projection,
    # and tanh warping ...
    geodesic_kernel = ... # (N1, N2)
    rq_component_TW = ... # (N1, N2)
    base_kernel = geodesic_kernel *
    rq_component_TW # (N1, N2)

    I = torch.eye(x1.size(-2), device=x1.
device)[None] # (1, N1, N1)
    regularized_kernel = base_kernel +
self.regularization_strength * I
    # ERROR: shape mismatch when N1 != N2

    if diag:
        return regularized_kernel.diagonal
(dim1=-2, dim2=-1)
    return regularized_kernel
```

```
def forward(self, x1, x2, diag=False, **
params):
    # x1: (N1, D), x2: (N2, D)
    x1 = (x1 - center) / lengthscale
    x2 = (x2 - center) / lengthscale

    dist2 = torch.cdist(x1, x2).square()
    # dist2: (N1, N2)

    r1 = dist2.diagonal(dim1=-2, dim2=-1)
    # r1: (min(N1, N2),)

    angular_kernel = ...
    # angular_kernel: (N1, N2)

    radial_kernel = self.
radial_base_kernel(r1.sqrt(), **params)
    # radial_kernel : (min, min)

    return angular_kernel * radial_kernel
    # ERROR: shape mismatch when N1 != N2
```

(a) Identity-matrix regularization. The added `torch.eye(N1)` creates a square $N_1 \times N_1$ term, which is only compatible with the covariance tensor shape $N_1 \times N_2$ when $N_1 = N_2$.

(b) Implicit square term. The diagonal extraction and sub-kernel call produce an $M \times M$ block with $M = \min(N_1, N_2)$, which is incompatible with the $N_1 \times N_2$ angular term.

Figure 18: Representative cross-covariance failures caught by \mathcal{V}_{agn} . Both implementations pass $K(\mathbf{X}, \mathbf{X})$ but fail on rectangular $K(\mathbf{X}_*, \mathbf{X})$ calls. Problematic lines are highlighted in red.

692 For \mathcal{V}_{psd} , a generated implementation may satisfy the required output shape but still fail to define
 693 a valid covariance function. A valid GP kernel must produce a positive semi-definite Gram matrix
 694 for any finite input set. A common failure mode is that the generated code adds terms that are not
 695 generally PSD-preserving, such as distance-increasing components or oscillatory transformations of
 696 pairwise distances. Such terms can yield indefinite Gram matrices, as illustrated in Figure 19.

```
def forward(self, x1, x2, diag=False, **
params):
    x1s = x1 / self.lengthscale
    x2s = x2 / self.lengthscale

    dist_sq = torch.cdist(x1s, x2s).pow(2)

    rbf_term = torch.exp(-dist_sq)
    dist_term = torch.sqrt(dist_sq)

    covar = self.alpha * rbf_term + self.
beta * dist_term + self.gamma

    if diag:
        return covar.diagonal(dim1=-2,
dim2=-1)
    return covar
```

```
def forward(self, x1, x2, diag=False, **
params):
    x1s = x1 / self.lengthscale
    x2s = x2 / self.lengthscale

    dist_sq = torch.cdist(x1s, x2s).pow(2)

    base_term = torch.exp(-0.5 * dist_sq)
    osc_term = torch.cos(math.pi*dist_sq)

    covar = self.alpha * base_term + \
(1 - self.alpha) * osc_term

    if diag:
        return covar.diagonal(dim1=-2,
dim2=-1)
    return covar
```

(a) Distance-increasing term. This corresponds to $k(x, x') = \alpha e^{-d^2} + \beta d + \gamma$. The added distance term $d = \|x - x'\|$ is not generally PSD-preserving.

(b) Oscillatory distance transform. This corresponds to $k(x, x') = \alpha e^{-d^2/2} + (1 - \alpha) \cos(\pi d^2)$. The oscillatory term $\cos(\pi d^2)$ is not generally PSD-preserving.

Figure 19: Representative PSD failures caught by \mathcal{V}_{psd} . Both include terms that do not generally preserve positive semi-definiteness and can yield indefinite Gram matrices. Problematic terms are highlighted in red.

697 **I Robustness to Different LLMs**

698 **Different LLMs for Other Baselines.** We additionally evaluate LLM-based baselines under
 699 different LLM backbones to examine whether the performance gap is sensitive to the choice of LLM.
 700 Specifically, we compare CAKE and LMABO using GPT-4o and GPT-4o-mini, and report the results
 701 together with our method in Table 10. Across the evaluated settings, our method consistently achieves
 702 stronger BO performance, suggesting that the performance gain comes from the kernel discovery
 703 framework rather than merely from the choice of LLM backbone.

Table 10: Performance of LLM-based methods under different LLM backbones. Results are averaged over 4 random seeds. Bold denotes the best entry in each column.

Method	Rover (\uparrow) ($D = 100$)	Mopta08 (\downarrow) ($D = 124$)	Lasso-DNA (\downarrow) ($D = 180$)	SVM388 (\downarrow) ($D = 388$)	Humanoid (\uparrow) ($D = 6392$)
Kernel Discovery (GPT-4o, Ours)	4.353 \pm 0.319	216.81 \pm 0.87	0.286 \pm 0.001	0.056 \pm 0.003	762.78 \pm 72.39
Kernel Discovery (GPT-4o-mini, Ours)	3.582 \pm 0.230	219.16 \pm 1.54	0.289 \pm 0.002	0.061 \pm 0.004	589.28 \pm 107.87
CAKE (GPT-4o)	3.412 \pm 0.586	231.40 \pm 15.67	0.300 \pm 0.011	0.131 \pm 0.036	667.49 \pm 24.99
CAKE (GPT-4o-mini)	3.863 \pm 0.497	229.43 \pm 3.88	0.290 \pm 0.001	0.103 \pm 0.058	648.13 \pm 29.11
LMABO (GPT-4o)	3.953 \pm 0.369	240.70 \pm 5.24	0.304 \pm 0.006	0.226 \pm 0.001	390.23 \pm 44.78
LMABO (GPT-4o-mini)	4.078 \pm 0.299	244.21 \pm 3.97	0.302 \pm 0.005	0.218 \pm 0.016	410.09 \pm 42.34

704 **Open-source LLMs for Kernel Discovery.** We further evaluate whether our kernel discovery
 705 pipeline can operate with a smaller open-source LLM. To this end, we replace the kernel-proposing
 706 LLM with Qwen3-8B while keeping the rest of the pipeline unchanged. As shown in Table 11,
 707 GPT-4o achieves the best performance, but the Qwen3-8B variant still obtains competitive results
 708 across benchmarks. Compared with the main results in Table 1, the open-source variant remains
 709 comparable to strong baseline methods, indicating that our framework does not rely exclusively on a
 710 proprietary LLM. These results suggest that stronger LLMs improve the quality of discovered kernels,
 711 while the proposed pipeline remains effective even with a relatively small open-source model.

Table 11: Performance of Kernel Discovery with proprietary and open-source LLMs. Results are averaged over 4 random seeds. Bold denotes the best entry in each column.

Method	Rover (\uparrow) ($D = 100$)	Mopta08 (\downarrow) ($D = 124$)	Lasso-DNA (\downarrow) ($D = 180$)	SVM388 (\downarrow) ($D = 388$)	Humanoid (\uparrow) ($D = 6392$)
Kernel Discovery (GPT-4o, Ours)	4.353 \pm 0.319	216.81 \pm 0.87	0.286 \pm 0.001	0.056 \pm 0.003	762.78 \pm 72.39
Kernel Discovery (Qwen3-8B, Ours)	3.838 \pm 0.293	221.67 \pm 1.53	0.288 \pm 0.002	0.058 \pm 0.003	622.15 \pm 107.11

712 J Evolution of Discovered Kernels

713 **Detailed Explanation of Discovered Kernels in Figure 7a.** We provide an explanation of the dis-
 714 covered kernels highlighted in the figure. We first define the kernel components used in the discovered
 715 kernels, and then describe each discovered kernel by its decomposition and its interpretations.

716 *Kernel components.* We first define the kernel components that appear in the discovered kernels.
 717 For the RQ component evaluated in the original input space, let $\ell_B \in \mathbb{R}_{>0}^d$ be an ARD lengthscale,
 718 $r_{\text{RQ}} = \|\mathbf{x} \odot \ell_B - \mathbf{x}' \odot \ell_B\|_2$, and $\alpha > 0$. Then

$$k_{\text{RQ}}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{r_{\text{RQ}}^2}{2\alpha}\right)^{-\alpha}.$$

719 This component captures multi-scale smooth variation in a separately ARD-scaled input space.

720 Several other components are evaluated after an arctangent-based input transformation. Given an
 721 input $\mathbf{x} \in \mathbb{R}^d$ and an ARD lengthscale $\ell \in \mathbb{R}_{>0}^d$, define

$$\mathbf{z}_0(\mathbf{x}) = \mathbf{x} \odot \ell, \quad \mathbf{z}_i(\mathbf{x}) = \frac{1}{\sqrt{i}} \arctan(sw \mathbf{z}_{i-1}(\mathbf{x})), \quad i = 1, \dots, D,$$

722 where $s, w > 0$ are scalar transformation parameters and $D \in \mathbb{N}$ is the number of transformation
 723 layers. We denote the final transformed feature by $t(\mathbf{x}) = \mathbf{z}_D(\mathbf{x})$. Using this transformed feature, we
 724 define

$$k_{\text{Arc-IMQ}}(\mathbf{x}, \mathbf{x}') = (1 + \|t(\mathbf{x}) - t(\mathbf{x}')\|_2^2)^{-1}, \quad k_{\text{Arc-Lin}}(\mathbf{x}, \mathbf{x}') = t(\mathbf{x})^\top t(\mathbf{x}'),$$

725 and

$$k_{\text{Arc-RQ}}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{\|t(\mathbf{x}) - t(\mathbf{x}')\|_2^2}{2\alpha}\right)^{-\alpha}.$$

726 Here, the Arc-IMQ component captures heavy-tailed distance-based similarity in the arctangent-
 727 transformed feature space, the Arc-Linear component captures global alignment in the same trans-
 728 formed space, and the Arc-RQ component captures multi-scale smooth variation after the arctangent
 729 transformation.

730 One discovered kernel additionally uses an angular similarity component. Let

$$\tilde{\mathbf{x}} = \frac{\mathbf{x} - \mathbf{c}}{\mathbf{r}}, \quad \rho(\mathbf{x}) = \|\tilde{\mathbf{x}}\|_2, \quad \mathbf{a}(\mathbf{x}) = \frac{\tilde{\mathbf{x}}}{\rho(\mathbf{x})}.$$

731 Here $\mathbf{a}(\mathbf{x})$ is the unit direction vector of the normalized input. The polynomial angular kernel is

$$k_{\text{ang}}(\mathbf{x}, \mathbf{x}') = \sum_{p=0}^3 w_p (\mathbf{a}(\mathbf{x})^\top \mathbf{a}(\mathbf{x}'))^p, \quad w_p > 0.$$

732 This component captures directional similarity relative to the normalized center.

733 *Kernel 1:*

$$k_1(\mathbf{x}, \mathbf{x}') = k_{\text{Arc-IMQ}}(\mathbf{x}, \mathbf{x}') + k_{\text{Arc-Lin}}(\mathbf{x}, \mathbf{x}') + k_{\text{RQ}}(\mathbf{x}, \mathbf{x}').$$

734 This kernel combines heavy-tailed distance-based similarity and linear alignment in the arctangent-
 735 transformed feature space with multi-scale smooth variation in a separate ARD-scaled input space.

736 *Kernel 2:*

$$k_2(\mathbf{x}, \mathbf{x}') = k_{\text{ang}}(\mathbf{x}, \mathbf{x}') \cdot [k_{\text{Arc-IMQ}}(\mathbf{x}, \mathbf{x}') + k_{\text{Arc-Lin}}(\mathbf{x}, \mathbf{x}')].$$

737 This kernel is a product of a polynomial angular kernel and an arc-transformed additive kernel.
 738 The angular component captures directional similarity relative to a normalized center, while the arc
 739 component captures heavy-tailed distance-based similarity and linear alignment in the arctangent-
 740 transformed feature space.

741 *Kernel 3:*

$$k_3(\mathbf{x}, \mathbf{x}') = k_{\text{Arc-RQ}}(\mathbf{x}, \mathbf{x}') + k_{\text{Arc-IMQ}}(\mathbf{x}, \mathbf{x}') + k_{\text{Arc-Lin}}(\mathbf{x}, \mathbf{x}').$$

742 This kernel evaluates all three components in the same arctangent-transformed feature space. It
 743 combines multi-scale smooth variation, heavy-tailed distance-based similarity, and global transformed-
 744 space alignment under a shared nonlinear input warping.

745 We also provide the corresponding forward implementations for the discovered kernels.

```

1 def forward(self, x1, x2, diag=False, **params):
2     if x1.dim() == 1: x1 = x1.unsqueeze(-1)
3     if x2.dim() == 1: x2 = x2.unsqueeze(-1)
4
5     # Kernel A: Arc-transformed IMQ + Arc-transformed Linear
6     x1_scaled_a = x1 / self.lengthscale_a
7     x2_scaled_a = x2 / self.lengthscale_a
8     t1, t2 = x1_scaled_a, x2_scaled_a
9     for i in range(1, self.depth.item() + 1):
10        t1 = torch.atan(self.global_scale * self.arc_weight * t1) / math.sqrt
(i)
11        t2 = torch.atan(self.global_scale * self.arc_weight * t2) / math.sqrt
(i)
12    linear_term = t1 @ t2.transpose(-1, -2)
13    r_squared_a = torch.cdist(t1, t2, p=2).pow(2)
14    k_a = (1 + r_squared_a).pow(-1) + linear_term
15
16    # Kernel B: Rational Quadratic
17    x1_scaled_b = x1 / self.lengthscale_b
18    x2_scaled_b = x2 / self.lengthscale_b
19    r_squared_b = torch.cdist(x1_scaled_b, x2_scaled_b, p=2).pow(2)
20    k_b = (1 + r_squared_b / (2 * self.alpha)).pow(-self.alpha)
21
22    covar = k_a + k_b
23    if diag:
24        return covar.diagonal(dim1=-2, dim2=-1)
25    return covar

```

Figure 20: Forward code for the discovered kernel k_1 .

```

1 def forward(self, x1, x2, diag=False, **params):
2     if x1.dim() == 1: x1 = x1.unsqueeze(-1)
3     if x2.dim() == 1: x2 = x2.unsqueeze(-1)
4     arc_weight = self.raw_arc_weight_constraint.transform(self.raw_arc_weight
)
5     global_scale = self.raw_global_scale_constraint.transform(self.
raw_global_scale)
6     angular_weights = self.raw_angular_weights_constraint.transform(self.
raw_angular_weights)
7
8     x1_normalized = (x1 - self.center) / self.radius
9     x2_normalized = (x2 - self.center) / self.radius
10    r1 = x1_normalized.norm(dim=-1, keepdim=True)
11    r2 = x2_normalized.norm(dim=-1, keepdim=True)
12    a1 = x1_normalized / r1.clamp(min=1e-15)
13    a2 = x2_normalized / r2.clamp(min=1e-15)
14
15    t1 = x1_normalized / self.lengthscale
16    t2 = x2_normalized / self.lengthscale
17    for i in range(1, self.depth + 1):
18        t1 = torch.atan(global_scale * arc_weight * t1) / math.sqrt(i)
19        t2 = torch.atan(global_scale * arc_weight * t2) / math.sqrt(i)
20
21    angular_cov = self._angular_kernel(a1, a2, angular_weights)
22    t1_dist = torch.cdist(t1, t2, p=2)
23    arc_cov = (1 + t1_dist.pow(2)).pow(-1) + (t1 @ t2.transpose(-1, -2))
24
25    covar = angular_cov * arc_cov
26    if diag:
27        return covar.diagonal(dim1=-2, dim2=-1)
28    return covar

```

Figure 21: Forward code for the discovered kernel k_2 .

```

1 def forward(self, x1, x2, diag=False, **params):
2     if x1.dim() == 1: x1 = x1.unsqueeze(-1)
3     if x2.dim() == 1: x2 = x2.unsqueeze(-1)
4
5     # Kernel A: Arc-transformed IMQ + Arc-transformed Linear
6     x1_scaled_a = x1 / self.lengthscale_a
7     x2_scaled_a = x2 / self.lengthscale_a
8     t1, t2 = x1_scaled_a, x2_scaled_a
9     for i in range(1, self.depth.item() + 1):
10        t1 = torch.atan(self.global_scale * self.arc_weight * t1) / math.sqrt
(i)
11        t2 = torch.atan(self.global_scale * self.arc_weight * t2) / math.sqrt
(i)
12    linear_term = t1 @ t2.transpose(-1, -2)
13    r_squared_a = torch.cdist(t1, t2, p=2).pow(2)
14    k_a = (1 + r_squared_a).pow(-1) + linear_term
15
16    # Kernel B: Rational Quadratic
17    x1_scaled_b = x1 / self.lengthscale_b
18    x2_scaled_b = x2 / self.lengthscale_b
19    r_squared_b = torch.cdist(x1_scaled_b, x2_scaled_b, p=2).pow(2)
20    k_b = (1 + r_squared_b / (2 * self.alpha)).pow(-self.alpha)
21
22    covar = k_a + k_b
23    if diag:
24        return covar.diagonal(dim1=-2, dim2=-1)
25    return covar

```

Figure 22: Forward code for the discovered kernel k_3 .

746 **Evolving Population on Other Benchmarks** Figure 23 extends the post-analysis of Figure 7a to
747 four additional benchmarks: Rover, Mopta08, Lasso-DNA, and Humanoid. For each benchmark,
748 we visualize the best-so-far value trajectory and annotate the kernels selected at the top-5 largest
749 improvements. The components of the annotated kernels are described using the same kernel
750 vocabulary introduced above.

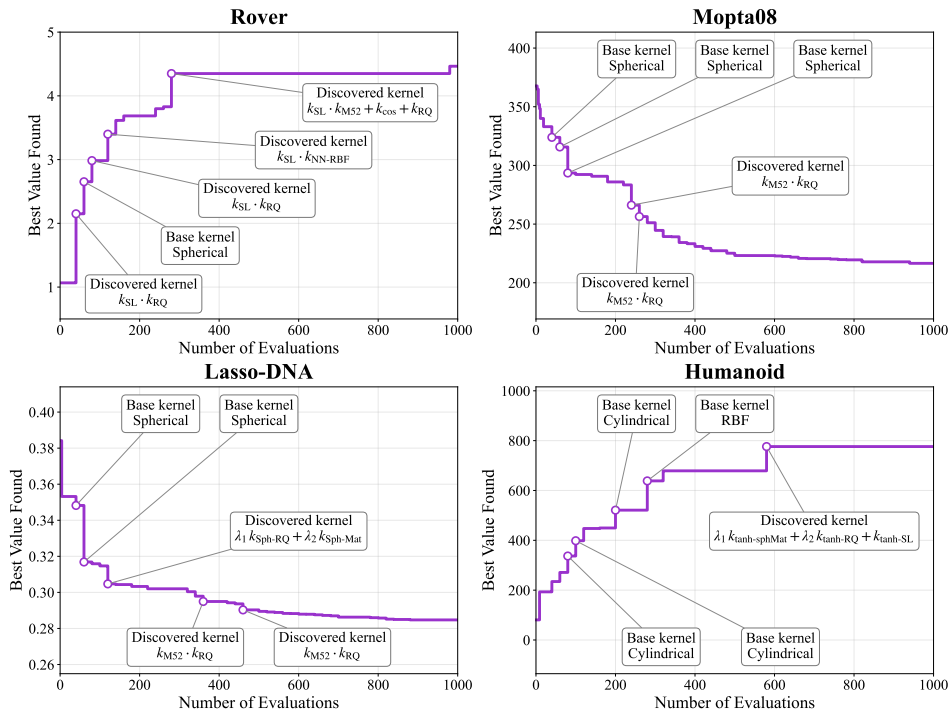


Figure 23: Post-analysis of kernel discovery on four additional benchmarks. For each benchmark, we annotate the kernels selected at the top-5 largest improvements in the best-so-far value.

751 For the additional discovered kernels annotated in Figure 23, we provide the mathematical definitions
 752 of the component kernels used in the decompositions.

753 *Matérn-5/2 component.* Let

$$r_{\text{Matérn52}} = \|\mathbf{x} \odot \boldsymbol{\ell}_{\text{Matérn52}} - \mathbf{x}' \odot \boldsymbol{\ell}_{\text{Matérn52}}\|_2.$$

754 Then

$$k_{\text{Matérn52}}(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{5}r_{\text{Matérn52}} + \frac{5}{3}r_{\text{Matérn52}}^2\right) \exp(-\sqrt{5}r_{\text{Matérn52}}).$$

755 *Spherical lifted linear component.* Given a scaled input $\mathbf{u}(\mathbf{x})$, define the stereographic projection

$$\psi(\mathbf{u}(\mathbf{x})) = \frac{[2\mathbf{u}(\mathbf{x}), \|\mathbf{u}(\mathbf{x})\|_2^2 - 1]}{\|\mathbf{u}(\mathbf{x})\|_2^2 + 1}.$$

756 Then

$$k_{\text{SL}}(\mathbf{x}, \mathbf{x}') = \lambda_0 + \lambda_1 \psi(\mathbf{u}(\mathbf{x}))^\top \psi(\mathbf{u}(\mathbf{x}')), \quad \lambda_0 + \lambda_1 = 1, \lambda_0, \lambda_1 \geq 0.$$

757 *Neural-network RBF component.* Let

$$\mathbf{h}(\mathbf{x}) = \text{NN}(\mathbf{x} \odot \boldsymbol{\ell}_{\text{NN}}).$$

758 Then

$$k_{\text{NN-RBF}}(\mathbf{x}, \mathbf{x}') = \exp\left(-\frac{\|\mathbf{h}(\mathbf{x}) - \mathbf{h}(\mathbf{x}')\|_2^2}{2\sigma_{\text{NN}}^2}\right).$$

759 *Cosine warping component.* Define

$$\mathbf{v}(\mathbf{x}) = \tanh(\mathbf{x}/r_{\text{cos}}), \quad \rho_{\text{cos}}(\mathbf{x}) = \|\mathbf{v}(\mathbf{x})\|_2.$$

760 Then

$$k_{\text{cos}}(\mathbf{x}, \mathbf{x}') = \cos(2\pi\alpha_{\text{cos}} |\rho_{\text{cos}}(\mathbf{x}) - \rho_{\text{cos}}(\mathbf{x}')|).$$

761 *Spherical RQ and spherical Matérn-like components.* Using the stereographic projection $\psi(\mathbf{u}(\mathbf{x}))$,
 762 define

$$d_{\text{sph}} = \|\psi(\mathbf{u}(\mathbf{x})) - \psi(\mathbf{u}(\mathbf{x}'))\|_2.$$

763 The spherical RQ component is

$$k_{\text{sphRQ}}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{d_{\text{sph}}^2}{2\alpha}\right)^{-\alpha},$$

764 and the spherical Matérn-like component is

$$k_{\text{sphMat}}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{\sqrt{3}}{\nu}d_{\text{sph}} + \frac{3}{\nu^2}d_{\text{sph}}^2\right) \exp\left(-\frac{\sqrt{3}}{\nu}d_{\text{sph}}\right).$$

765 where $\nu > 0$ is a learnable parameter that controls the effective length scale and smoothness of the
 766 Matérn-like component.

767 *Tanh-transformed components.* Define the tanh embedding

$$\mathbf{h}_{\text{tanh}}(\mathbf{x}) = \tanh(\alpha(\mathbf{x} \odot \boldsymbol{\ell})).$$

768 For the spherical branch, define

$$\mathbf{u}_{\text{tanh}}(\mathbf{x}) = (\mathbf{h}_{\text{tanh}}(\mathbf{x}) - \mathbf{c}) \odot \boldsymbol{\ell},$$

769 and let $\psi_{\text{tanh}}(\mathbf{x}) = \psi(\mathbf{u}_{\text{tanh}}(\mathbf{x}))$ be the stereographic projection of $\mathbf{u}_{\text{tanh}}(\mathbf{x})$. Define

$$r_{\text{tanh-sph}} = \|\psi_{\text{tanh}}(\mathbf{x}) - \psi_{\text{tanh}}(\mathbf{x}')\|_2.$$

770 Then

$$k_{\text{tanh-sphMat}}(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{2\nu}r_{\text{tanh-sph}} + \frac{2\nu}{3}r_{\text{tanh-sph}}^2\right) \exp\left(-\sqrt{2\nu}r_{\text{tanh-sph}}\right),$$

771

$$k_{\text{tanh-RQ}}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{\|\mathbf{h}_{\text{tanh}}(\mathbf{x}) - \mathbf{h}_{\text{tanh}}(\mathbf{x}')\|_2^2}{2\gamma}\right)^{-\gamma},$$

772 and

$$k_{\text{tanh-SL}}(\mathbf{x}, \mathbf{x}') = \psi_{\text{tanh}}(\mathbf{x})^\top \psi_{\text{tanh}}(\mathbf{x}').$$

773 **K Transferability of Discovered Kernels**

774 **Experiment Results on Other Benchmarks.** In Figure 7b, we visualize the results of the kernel
 775 discovered by our pipeline on the SVM benchmark. In this section, we report the kernel’s results
 776 across the other benchmarks. As shown in Table 12, it consistently achieves competitive and often
 777 outperforms other base kernels. We also visualize the code snippet of the kernel in Figure 25.

Table 12: Transferability of the discovered kernel across other benchmarks. We compare the kernel discovered on the SVM benchmark, when transferred to the other four standard benchmarks, against base kernels. D denotes the dimensionality of the task. **Blue** denotes the best entry in the column, and **Violet** denotes the second best. Experiments are conducted with 4 random seeds.

Method	Rover (\uparrow) ($D = 100$)	Mopta08 (\downarrow) ($D = 124$)	Lasso-DNA (\downarrow) ($D = 180$)	SVM (\downarrow) ($D = 388$)	Humanoid (\uparrow) ($D = 6392$)	Average Rank
Base Kernels						
RBF	3.552 \pm 0.304	217.75 \pm 2.33	0.291 \pm 0.001	0.061 \pm 0.004	503.00 \pm 67.82	3.8 / 8
Matérn52	3.453 \pm 0.366	215.77 \pm 0.57	0.292 \pm 0.003	0.063 \pm 0.003	509.79 \pm 102.32	4.0 / 8
Linear	3.950 \pm 0.441	274.67 \pm 8.87	0.312 \pm 0.004	0.226 \pm 0.003	435.49 \pm 42.00	6.4 / 8
Periodic	3.664 \pm 0.231	309.10 \pm 5.01	0.332 \pm 0.005	0.226 \pm 0.002	413.83 \pm 71.30	7.4 / 8
BOCK	3.725 \pm 0.593	225.36 \pm 1.82	0.291 \pm 0.003	0.068 \pm 0.007	669.52 \pm 50.15	3.4 / 8
SL	4.096 \pm 0.473	246.78 \pm 3.92	0.297 \pm 0.001	0.112 \pm 0.007	637.72 \pm 37.86	4.4 / 8
RQ	3.858 \pm 0.515	217.53 \pm 4.25	0.294 \pm 0.002	0.069 \pm 0.017	600.80 \pm 152.49	3.8 / 8
Discovered Kernel (Ours)	4.238 \pm 0.640	218.30 \pm 1.00	0.296 \pm 0.004	0.052 \pm 0.006	629.12 \pm 122.84	2.8 / 8

778 **Analysis on Boundary-seeking Behavior.** To understand why the discovered kernel k_{discover} in
 779 Equation (11) achieves superior performance on the SVM benchmark, we analyze whether its non-
 780 stationary component $k_{\text{tanh-Poly}}$ induces boundary-seeking behavior. We measure two complementary
 781 statistics over BO iterations: (1) the boundary hit ratio, i.e., the fraction of queried points lying on
 782 the boundary of the search domain, and (2) the observation traveling salesman distance (OTSD),
 783 which measures the total path length traversed across sequential queries. As shown in Figure 24, the
 784 discovered kernel maintains a boundary hit ratio comparable to RBF throughout optimization, and
 785 substantially lower than the Linear kernel, on both the SVM and Mopta08 benchmarks. The OTSD of
 786 the discovered kernel similarly remains low and closely tracks that of RBF, indicating that consecutive
 787 queries tend to be spatially concentrated rather than widely scattered. Together, these results suggest
 788 that the non-stationary component is regularized during GP hyperparameter optimization, and that the
 789 combination of several kernel components encourages a locally concentrated query behavior, which
 790 is a key factor enabling the discovered kernel to consistently identify higher-quality candidates.

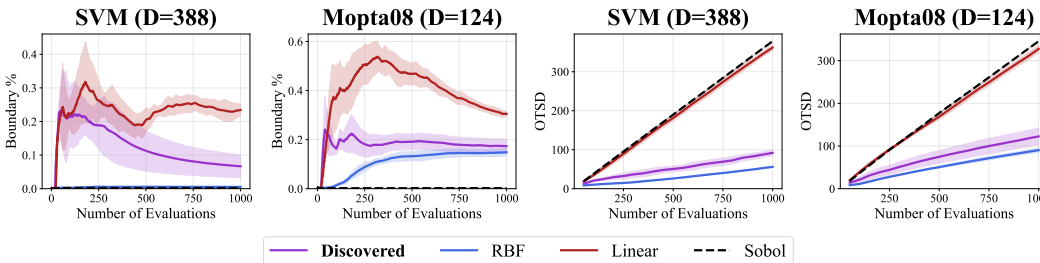


Figure 24: Analysis on boundary-seeking behavior of several kernels.

791 **Detailed Explanation of the Discovered Kernel in Equation (11).** We provide a detailed explana-
 792 tion of the discovered kernel in Equation (11). This kernel combines a Matérn-5/2 component with
 793 an additive composition of a tanh-polynomial component and an RQ component:

$$k_{\text{discover}}(\mathbf{x}, \mathbf{x}') = k_{\text{Matérn52}}(\mathbf{x}, \mathbf{x}') \cdot [k_{\text{tanh-Poly}}(\mathbf{x}, \mathbf{x}') + k_{\text{RQ}}(\mathbf{x}, \mathbf{x}')].$$

794 *Matérn-5/2 component.* Let $\ell \in \mathbb{R}_{>0}^d$ be an ARD lengthscale and define

$$r_{\text{Matérn52}}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} \odot \ell - \mathbf{x}' \odot \ell\|_2.$$

795 The Matérn-5/2 component is

$$k_{\text{Matérn52}}(\mathbf{x}, \mathbf{x}') = \left(1 + \sqrt{5}r_{\text{Matérn52}} + \frac{5}{3}r_{\text{Matérn52}}^2 \right) \exp\left(-\sqrt{5}r_{\text{Matérn52}}\right).$$

796 This component captures moderately smooth local correlation in the ARD-scaled input space.

797 *Tanh-polynomial component.* Using the same ARD-scaled input, define the tanh-transformed feature

$$\mathbf{h}_{\text{tanh}}(\mathbf{x}) = \tanh(\sigma(\mathbf{x} \odot \boldsymbol{\ell})),$$

798 where $\sigma > 0$ controls the strength of the tanh projection. The tanh-polynomial component is

$$k_{\text{tanh-Poly}}(\mathbf{x}, \mathbf{x}') = \sum_{n=0}^2 w_n (\mathbf{h}_{\text{tanh}}(\mathbf{x})^\top \mathbf{h}_{\text{tanh}}(\mathbf{x}'))^n, \quad w_n > 0.$$

799 This component captures polynomial similarity after mapping inputs into a bounded nonlinear feature
800 space.

801 *Rational Quadratic component.* Let

$$r_{\text{RQ}}(\mathbf{x}, \mathbf{x}') = \|\mathbf{x} \odot \boldsymbol{\ell} - \mathbf{x}' \odot \boldsymbol{\ell}\|_2.$$

802 The RQ component is

$$k_{\text{RQ}}(\mathbf{x}, \mathbf{x}') = \left(1 + \frac{r_{\text{RQ}}^2}{2\alpha} \right)^{-\alpha}, \quad \alpha > 0.$$

803 This component captures multi-scale smooth variation in the ARD-scaled input space.

804 Overall, the product structure modulates the additive tanh-polynomial and RQ similarities by the
805 Matérn-5/2 local smoothness component.

806 **Illustrative Example of Discovered Kernel.** As a representative example, we present the full
807 implementation of the kernel discovered from the SVM benchmark, which corresponds to k_{discover} in
808 Equation (11).

```

1 def forward(self, x1, x2, diag=False, **params):
2     if x1.dim() == 1: x1 = x1.unsqueeze(-1)
3     if x2.dim() == 1: x2 = x2.unsqueeze(-1)
4     x1s = x1 / self.lengthscale
5     x2s = x2 / self.lengthscale
6
7     # Matern-5/2 Component
8     r = torch.cdist(x1s, x2s, p=2).clamp(min=1e-15)
9     sqrt5_r = math.sqrt(5.0) * r
10    k_matern = (1 + sqrt5_r + 5.0/3.0 * r.pow(2)) \
11              * torch.exp(-sqrt5_r)
12
13    # Tangential Polynomial Component
14    t1 = torch.tanh(self.projection_sigma * x1s)
15    t2 = torch.tanh(self.projection_sigma * x2s)
16    t_ip = t1 @ t2.transpose(-1, -2)
17    k_tanh_poly = sum(
18        self.tangent_weights[n] * t_ip.pow(n)
19        for n in range(3))
20
21    # Rational Quadratic Component
22    r2 = torch.cdist(x1s, x2s, p=2).pow(2)
23    k_rq = (1 + r2 / (2 * self.alpha)).pow(-self.alpha)
24
25    covar = k_matern * (k_tanh_poly + k_rq)
26    if diag:
27        return covar.diagonal(dim1=-2, dim2=-1)
28    return covar

```

Figure 25: Code snippet of the discovered kernel from SVM benchmark.

809 We additionally present two discovered kernels that show strong transferability. Using the same
 810 fixed-kernel evaluation protocol as in Table 12, both achieve an overall average rank of 2.8 across
 811 five benchmarks, outperforming all standard base kernels in aggregate.

```

1 def forward(self, x1, x2, diag=False, **params):
2     if x1.dim() == 1: x1 = x1.unsqueeze(-1)
3     if x2.dim() == 1: x2 = x2.unsqueeze(-1)
4
5     # Rational Quadratic Transform
6     x_centered_RQ1 = x1 / self.lengthscale_RQ
7     x_centered_RQ2 = x2 / self.lengthscale_RQ
8     x_scaled_RQ1 = x_centered_RQ1 / math.sqrt(self.D)
9     x_scaled_RQ2 = x_centered_RQ2 / math.sqrt(self.D)
10
11    # Bezier Path Projection
12    t = self.harmonic_alpha
13    path_bezier = sum((1-t)**(self.control_point_count-i) * t**i * self.
bezier_control_points[... , i, :])
14        for i in range(self.control_point_count))
15    x_path_BZ1 = x1 / self.lengthscale_BZ
16    x_path_BZ2 = x2 / self.lengthscale_BZ
17    x_proj_BZ1 = self._inv_stereographic((x_path_BZ1 / x_path_BZ1.norm(dim
=-1, keepdim=True) + path_bezier) / 2)
18    x_proj_BZ2 = self._inv_stereographic((x_path_BZ2 / x_path_BZ2.norm(dim
=-1, keepdim=True) + path_bezier) / 2)
19    x_scaled_BZ1 = x_proj_BZ1 / math.sqrt(self.D)
20    x_scaled_BZ2 = x_proj_BZ2 / math.sqrt(self.D)
21
22    # Rational Quadratic Kernel
23    r_RQ2 = torch.cdist(x_scaled_RQ1, x_scaled_RQ2, p=2).pow(2)
24    k_RQ = (1 + r_RQ2 / (2 * self.alpha_RQ)).pow(-self.alpha_RQ)
25
26    # Bezier-Projected Polynomial Kernel
27    k_BZ = (x_proj_BZ1 @ x_proj_BZ2.transpose(-1, -2) + 1).pow(3)
28
29    # Hybrid Kernel Covariance
30    k = k_RQ + self.harmonic_alpha * k_BZ
31
32    if diag:
33        return k.diagonal(dim1=-2, dim2=-1)
34    return k

```

Figure 26: Code snippet of an additional discovered kernel with strong transferability. This kernel combines an RQ component with a Bézier-projected polynomial component.

```

1 def forward(self, x1, x2, diag=False, **params):
2     if x1.dim() == 1: x1 = x1.unsqueeze(-1)
3     if x2.dim() == 1: x2 = x2.unsqueeze(-1)
4
5     # Rational Quadratic component
6     x1_rq = x1 / self.lengthscale
7     x2_rq = x2 / self.lengthscale
8     r_rq2 = torch.cdist(x1_rq, x2_rq, p=2).pow(2)
9     k_rq = (1 + r_rq2 / (2 * self.alpha)).pow(-self.alpha)
10
11    # Spherical component
12    x1_sphere = self._project_to_sphere(x1 / self.lengthscale_sphere)
13    x2_sphere = self._project_to_sphere(x2 / self.lengthscale_sphere)
14    coeffs = self.poly_coeffs_sphere
15    phi1 = torch.cat([x1_sphere * coeffs[1].sqrt(), coeffs[0].sqrt() * torch.
ones_like(x1_sphere[... , :1])], dim=-1)
16    phi2 = torch.cat([x2_sphere * coeffs[1].sqrt(), coeffs[0].sqrt() * torch.
ones_like(x2_sphere[... , :1])], dim=-1)
17    k_sphere = phi1 @ phi2.transpose(-1, -2)
18
19    # Matern-5/2 component
20    x1_m = x1 / self.lengthscale_matern52
21    x2_m = x2 / self.lengthscale_matern52
22    r_m = torch.cdist(x1_m, x2_m, p=2).clamp(min=1e-15)
23    k_m = (1 + math.sqrt(5) * r_m + 5 / 3 * r_m.pow(2)) * torch.exp(-math.
sqrt(5) * r_m)
24
25    # Cosmic Warping component
26    x1_cosmic = torch.tanh(x1 / self.radius_cosmic)
27    x2_cosmic = torch.tanh(x2 / self.radius_cosmic)
28    cosmic_phase1 = x1_cosmic.norm(dim=-1, keepdim=True)
29    cosmic_phase2 = x2_cosmic.norm(dim=-1, keepdim=True)
30    k_cosmic = torch.cos(2 * math.pi * self.alpha_cosmic * torch.cdist(
cosmic_phase1, cosmic_phase2, p=2))
31
32    # Combined Kernel
33    covar = k_rq + (k_sphere * k_m) + k_cosmic
34
35    if diag:
36        return covar.diagonal(dim1=-2, dim2=-1)
37    return covar

```

Figure 27: Code snippet of another additional discovered kernel with strong transferability. This kernel combines RQ, spherical, Matérn-5/2, and warped periodic-like components.