

An Exploratory Study of Bayesian Prompt Optimization for Test-Driven Code Generation with Large Language Models

SHLOK TOMAR, Washington State University, USA

ARYAN DESHWAL, University of Minnesota, USA

ETHAN VILLALOVOZ, Washington State University, USA

MATTIA FAZZINI, University of Minnesota, USA

HAIPENG CAI, University at Buffalo, USA

JANARDHAN RAO DOPPA, Washington State University, USA

We consider the task of generating functionally correct code using large language models (LLMs). The correctness of generated code is influenced by the prompt used to query the given base LLM. We formulate the problem of finding the appropriate prompt as combinatorial search process and propose a Bayesian optimization (BO) approach referred to as *BO for Code GENeration (BODE-GEN)*. BODE-GEN performs an adaptive data-driven search over prompts guided by training data in the form of prompts tried and the functional accuracy of the generated code over a set of given test cases. The key insight is to perform BO in continuous embedding space by using an auxiliary LLM to bridge the gap between discrete prompt space and continuous embedding space. We leverage two synergistic ideas, namely, random projections and dimensionality scaled priors, to build effective Gaussian process based surrogate models over the high-dimensional embedding space. Our experiments on the HumanEval+ benchmark using multiple base LLMs show that BODE-GEN can improve performance in terms of code generation accuracy compared to fixed prompts and manual prompt engineering. Additionally, we demonstrate that BODE-GEN is sample-efficient, requiring relatively few iterations of BO to demonstrate improvements in code accuracy.

CCS Concepts: • **Software and its engineering** → **Source code generation**.

Additional Key Words and Phrases: AI for software development, large language models, program synthesis, Bayesian optimization

ACM Reference Format:

Shlok Tomar, Aryan Deshwal, Ethan Villalovoz, Mattia Fazzini, Haipeng Cai, and Janardhan Rao Doppa. 2025. An Exploratory Study of Bayesian Prompt Optimization for Test-Driven Code Generation with Large Language Models. 1, 1 (December 2025), 20 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

Large language models (LLMs) have emerged as transformative tools in various domains, including software development. Their ability to assist with code-related tasks has positioned them as indispensable coding assistants for software developers today[8, 9, 25]. However, as developers increasingly rely on code generated by LLMs, the functional correctness of this code has become an important factor in ensuring the overall quality of software products.

Authors' Contact Information: Shlok Tomar, Washington State University, USA, shlok.tomar@wsu.edu; Aryan Deshwal, University of Minnesota, USA, adeshwal@umn.edu; Ethan Villalovoz, Washington State University, USA, ethan.villalovoz@wsu.edu; Mattia Fazzini, University of Minnesota, USA, mfazzini@umn.edu; Haipeng Cai, University at Buffalo, USA, haipengc@buffalo.edu; Janardhan Rao Doppa, Washington State University, USA, jana.doppa@wsu.edu.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM XXXX-XXXX/2025/12-ART

<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

In modern software development, the software supply chain comprises of various components, including low-level systems software, application frameworks, third-party libraries, and build tools. These components are often inter-dependent and any flaw in one can propagate through the entire system, leading to significant and widespread challenges [11, 12]. When developers incorporate LLM-generated code into these components, the functional correctness of the code becomes vital. Incorrect code generated by LLMs can introduce subtle bugs that are difficult to detect and diagnose, potentially causing cascading side-effects for the entire software supply chain [4]. Therefore, ensuring the correctness of LLM-generated code before its incorporation into the software development process is of paramount importance.

Recent studies have highlighted that LLMs can produce incorrect code, for various reasons, including hallucinations and insufficient understanding of coding tasks [20, 25, 33]. To address this challenge, some potential solutions include better pre-training datasets and improved training/fine-tuning methods to create high-performing LLMs specifically for coding tasks [14]. There is also an inherent trade-off between resource cost and performance of LLMs. On one end of the spectrum, large LLMs require huge amount of high-quality training data and compute resources for both training and inference [13, 36]. On the other end, training and inference with smaller LLMs is feasible and computationally cheap, but these models may struggle to generalize well to handle the complexities of real-world code synthesis [34, 35]. In this context, leveraging large foundational LLMs such as ChatGPT through *prompting* has emerged as a promising avenue for developers seeking to generate high-quality code [7, 19, 23, 24]. However, recent studies have shown that treating LLMs as black-boxes and relying on standard/manual prompting can result in the generation of incorrect code [3, 16, 21, 30, 31].

Improving automated prompting approaches is complementary to alternative approaches of improving the code generation capabilities of LLMs. Prior work in this direction include knowledge augmentation [2, 16] and reasoning elicitation [1, 23]. Enhanced by self-consistency [37], the correctness of generated code repair can be further improved with CoT [1]. However, it is not clear if similar improvement can be achieved for synthesizing programs from scratch. Increasing the specificity of the code-generating prompts can provide additional help [24], but determining the specificity level is currently a manual process.

Test-driven software development [6] is an effective software engineering paradigm where developers write tests before code to write correct code. Inspired by the practical success of this paradigm, this paper asks the following question: *Given a base LLM and test cases for a coding task, can we develop an automated prompting approach to generate correct code by minimizing the number of tried prompts (sample-efficiency)?* There are two key challenges in answering this question. First, the search space of prompts is combinatorial and huge. Second, querying the base LLM with a candidate prompt and evaluating the accuracy of code on test cases is expensive. Therefore, we have a challenging search problem. This paper answers this question by proposing a prompt search approach based on the framework of Bayesian optimization (BO) [29]. The key idea behind BO is to learn a surrogate model from the past evaluations (prompt and code accuracy pairs) and use it to intelligently select a sequence of prompts to achieve our goal (generating code with 100 percent accuracy on the given test case).

Our proposed *BO for Code GENeration (BODE-GEN)* approach performs search in continuous embedding space by using an auxiliary LLM to bridge the gap between discrete prompt space and continuous embedding space. The popular Gaussian process (GP) [26] based surrogate models work well in low dimensions, but the embedding dimension for auxiliary LLM such as LLama2 is 4096 which poses significant challenges. To address this high-dimensional challenge, we leverage two synergistic ideas, namely, random projections and dimensionality scaled priors to build effective GP based surrogate models which are critical for the success of BO. BODE-GEN addresses the core

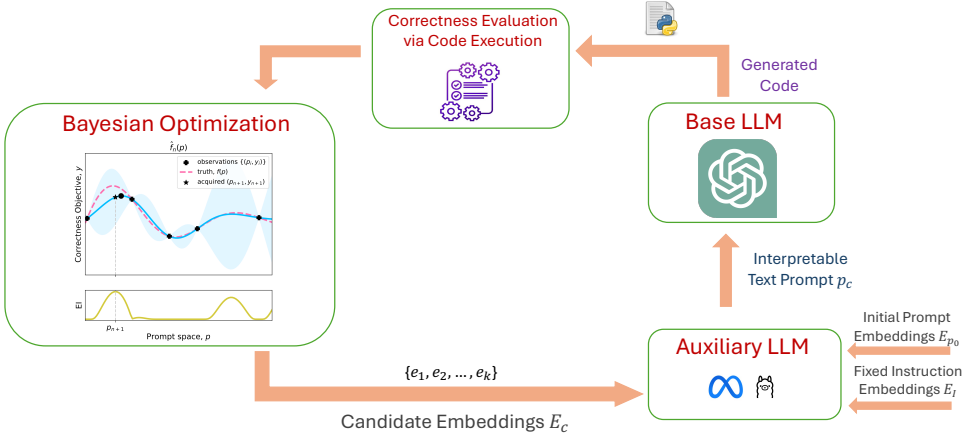


Fig. 1. High-level overview of our BODE-GEN approach. The method begins with a set of candidate embeddings $E_c = \{e_1, e_2, \dots, e_m\}$ proposed by the Bayesian Optimization algorithm. These embeddings combined with the initial prompt embeddings E_{p_0} and fixed instruction embeddings E_I are passed to the auxiliary LLM which generates an interpretable text prompt p_c . Subsequently, the base LLM is queried with this prompt p_c to generate the code which is evaluated for functional correctness through code execution on a set of developer provided test cases. The percentage of test cases passed by the code is used as the objective function value for the BO procedure. This overall procedure is repeated for a fixed number of iterations or until we find a prompt that generates code with maximizing code generation accuracy on the given test cases.

issue of incorrect code generation using LLMs at its source and provides a potential solution for developers. **Contributions.** The key contributions of this paper include:

- Development of a Bayesian optimization approach (BODE-GEN) to iteratively improve prompts for a given base LLM to solve code synthesis tasks.
- Demonstrating the effectiveness of BODE-GEN in generating code that meets functional requirements with higher correctness (fraction of passed test cases).

2 Problem Setup and Challenges

Let LLM_{base} denote a base large language model (LLM) that can be queried using textual prompts to solve coding tasks. Given an initial prompt p_0 (e.g., “write Python code for testing whether a given number is prime or not”) for a coding task T and a set of n developer-provided test cases to verify the correctness of the generated code, our goal is to find a prompt \hat{p} that when used with LLM_{base} will generate functionally correct code, i.e., produces correct outputs on all n test cases.

Suppose $Accuracy(LLM_{base}, p)$ represent the the functional accuracy of the code generated by the given base LLM using the prompt p on n test cases (e.g., 90 percent accuracy means the code passes 90 percent of the given n test cases). Our goal is to find a prompt \hat{p} which maximizes this accuracy and ideally achieves 100 percent accuracy. This problem can be mathematically formulated as follows.

$$\hat{p} = \arg \max_p Accuracy(LLM, T) \quad (1)$$

Key Challenges. There are two main challenges in solving this optimization problem.

- *Large combinatorial space of prompts.* Each prompt is a sequence of valid tokens. Given a token vocabulary and maximum size of the sequence, the search space of all candidate prompts is combinatorial and very large.
- *Expensive-to-evaluate objective function.* To evaluate the accuracy of a candidate prompt p , we need to query the base LLM to generate code and execute it on all n test cases. Each query to base LLM is expensive in terms of dollar cost or computational cost.

Therefore, our goal is to minimize the number of queries to the objective function (i.e., the number of tried prompts) to solve this optimization problem. Random search and trial-and-error approaches are not compatible with this goal because their exploration strategy doesn't incorporate machine learning and decision-theoretic reasoning to achieve the target goal.

3 Bayesian Optimization for Prompt Search

In this section, we describe an approach for sample-efficient prompt search based on the framework of Bayesian optimization (BO). First, we provide the necessary background on BO. Next, we describe the key challenges in using BO for prompt search and our proposed BODE-GEN approach based on continuous embeddings to address those challenges.

3.1 Background on Bayesian Optimization

BO [29] is a derivative-free method to *adaptively* and *efficiently* search a given input space X (e.g., search space of prompts) to optimize expensive-to-evaluate objective function $f(x \in X)$. BO is an adaptive procedure because it intelligently selects inputs from X by iterating between querying the objective function $y=f(x \in X)$ and making a decision about which input to query next x_{next} . BO is sample-efficient because it makes a data-driven decision to select the next input to query the objective function by taking into account all input-output pairs from previous query evaluations.

Each decision to select the next input from X to evaluate with f must trade-off two conflicting goals: 1) *Exploitation* suggests to use our current, but uncertain, approximation of the input-output relationship, based on the past query evaluations, to select the most promising input in terms of objective function value; and 2) *Exploration* suggests to select the input that we are most uncertain about to improve our approximation of the input-output relationship.

The key ingredients of BO for data-driven decision making are: 1) *surrogate model* that captures our beliefs, based on past objective function evaluations, about the input-output relationship; and 2) *acquisition function* that scores each input according to the utility of querying the objective function on it next. The acquisition function uses the surrogate model of the true input-output relationship $f(x \in X)$ to decide which input to evaluate next by trading-off exploration and exploitation.

The surrogate model $\hat{f}(x)$ is a probabilistic model of the input-output relationship $f(x)$ trained on all input-output pairs from past objective function evaluations. It reflects our current beliefs about $f(x)$ and serves two purposes in BO. First, to guide exploitation, it allows us to cheaply estimate the objective function value of all unevaluated inputs. Second, to guide exploration, variance quantifies the uncertainties in the predicted objective function value for the unevaluated inputs. This makes us aware of regions in input space we need to explore to improve our surrogate model and reduce the uncertainty in our beliefs about $f(x)$. Gaussian processes (GPs) [26] are the most commonly used surrogate models in BO owing to their flexibility as function approximators and principled uncertainty quantification.

The acquisition function scores the utility of evaluating the next input with the expensive objective function f . Here, "utility" is defined in terms of our ultimate goal of finding the optimal input with the minimum number of objective function evaluation queries. The acquisition function employs the prediction of the objective function value and the associated uncertainty from the

Algorithm 1 BODE-GEN Algorithm for Prompt Optimization

Require: Coding task T and n test cases; Initial prompt p_0 ; Base LLM LLM_{base} ; and Auxiliary LLM LLM_{aux}

Ensure: Optimized prompt \hat{p}

- 1: $E_{p_0} \leftarrow$ Embedding of initial prompt p_0 using LLM_{aux}
 - 2: $E_I \leftarrow$ Embedding of fixed instruction using LLM_{aux}
 - 3: Initialize surrogate model \mathcal{M} with random projections and dimensionality-scaled priors [15] on a set of randomly initialized points.
 - 4: **for** iteration $t = 1$ to T_{max} **do**
 - 5: $E_c \leftarrow \arg \max_E \text{ACQUISITIONFUNCTION}(\mathcal{M}, E)$
 - 6: $E_{comb} \leftarrow E_I \circ E_c \circ E_{p_0}$ ▷ Concatenated embedding
 - 7: $p_c \leftarrow LLM_{aux}(E_{comb})$ ▷ Generate discrete prompt
 - 8: $C \leftarrow LLM_{base}(p_c)$ ▷ Generate code using prompt p_c
 - 9: $Acc \leftarrow \text{EVALUATE}(C, n)$ ▷ Accuracy on n test cases
 - 10: Update surrogate model \mathcal{M} with (E_c, Acc)
 - 11: **end for**
 - 12: **return** best found prompt \hat{p} in terms of code accuracy
-

surrogate model to assign a utility score to each candidate input that balances exploitation and exploration, respectively. The decision of which input to evaluate next is made by maximizing the acquisition function. Importantly, the acquisition function is cheap to evaluate. Some popular acquisition functions include expected improvement (EI) and upper confidence bound (UCB) [29].

To summarize, BO is an iterative procedure that is executed until we reach our goal or maximum iterations are reached. In each iteration, we select the input that maximizes the acquisition function for objective function evaluation and then update the surrogate model based on new training example.

3.2 BO-based Prompt Search via Continuous Embeddings

Much of the BO success is on continuous spaces with small number of input dimensions. There are two intertwined surrogate modeling challenges in applying BO for prompt search. First, as opposed to continuous inputs, modeling of combinatorial objects (e.g., sequences) is quite challenging because of a lack of general notion of smoothness on such objects. This is especially exacerbated in the *small-data regime* where we have access to only a small number of supervised examples from the input space. Second, the search space of prompts is high-dimensional. We provide principled solutions to address these challenges as part of our proposed BODE-GEN approach and explain their details below.

BODE-GEN performs search for optimized prompts in a continuous embedding space as opposed to the discrete prompt space. The key insight is to leverage an auxiliary open-source LLM LLM_{aux} (e.g., LLaMA 2) to bridge the gap between continuous embedding space and discrete prompt space. Specifically, we perform BO (both surrogate modeling and acquisition function optimization) in the continuous embedding space. The continuous search space for our BO method is parameterized as a set of d -dimensional embeddings $E = \{e_1, e_2, \dots, e_m\}$, where each $e_i \in \mathbb{R}^d$ is a continuous vector in the d -dimensional embedding space of a local auxiliary LLM LLM_{aux} .

We perform the following sequence of steps in each iteration of BODE-GEN (see Algorithm 1 for pseudo-code and Figure 1 for illustration) given a surrogate model trained on the continuous embedding space.

- (1) Select the candidate input E_c from the continuous embedding space by maximizing the expected improvement (EI) acquisition function.
- (2) The selected embedding E_c is added as a suffix to the continuous embedding E_{p_0} of the initial prompt p_0 . We also prepend the embeddings E_I for a simple instruction I : “Your task is to rephrase/reformulate the code prompt given below to achieve a higher score on code generation by a large language model. Please provide the rephrased prompt in one block.” given to the auxiliary LLM LLM_{aux} in order to rephrase the original prompt p_0 . The resulting combined embedding $E_{comb}=E_I \circ E_c \circ E_{p_0}$ where \circ stands for concatenation operation.
- (3) The combined continuous embedding input E_{comb} is passed to the auxiliary LLM LLM_{aux} to generate a human-interpretable discrete prompt p_c .
- (4) The discrete prompt p_c is passed as input to the base LLM LLM_{base} to generate code C . The code C is executed on all n test cases to measure the functional accuracy, namely, $ACCURACY(LLM_{base}, T)$.
- (5) If the functional accuracy of code C is 100 percent, we return code C as output. Otherwise, the surrogate model is updated using the new training example: input is the continuous embedding E_c and output is code accuracy.

Surrogate Modeling over High-Dimensions. Gaussian Process (GP) [26] based surrogate models are commonly used in real-world BO applications with small number of input dimensions (typically less than 50). However, the high-dimensional embedding space of auxiliary LLM poses a significant challenge for standard GP models. For example, the embedding dimension for LLaMA2 is 4096. GP models that are directly fitted on such a high-dimensional continuous space struggle to generalize, especially when the amount of available supervised data is limited, as is considered in this paper. We apply two synergistic techniques to tackle this challenge: random projections followed by dimensionality-scaled priors [15] for kernel hyper-parameters.

- **Random Projections:** We employ random projections [18] to reduce the dimensionality of our search space. The key intuition behind this approach is that, in high-dimensional spaces, most of the interesting structure in the data lies in a lower-dimensional manifold. Random projections can capture this structure while preserving important properties of the data, such as pairwise distances between points (as formalized by Johnson-Lindenstrauss lemma [17]). In the context of our prompt optimization setting, random projections allow us to work with a more manageable representation of the embedding space without excessively reducing the information content. Let $x \in \mathbb{R}^d$ be a point in our original high-dimensional embedding space, where d is large (e.g., 4096 for LLaMA 2). We aim to project this point onto a lower-dimensional space $\mathcal{Z} \in \mathbb{R}^k$, where $k \ll d$. The random projection is defined by a matrix $A \in \mathbb{R}^{k \times d}$, where each entry a_{ij} of matrix A is sampled independently from a standard normal distribution:

$$a_{ij} \sim \mathcal{N}(0, 1) \quad (2)$$

The projected point $z \in \mathcal{Z}$ is then obtained by:

$$z = Ax \quad (3)$$

After applying random projections, our GP surrogate model is defined on the low-dimensional space \mathcal{Z} .

- **Dimensionality scaled priors:** While random projections effectively reduce the dimensionality of our search space, the resulting projected space can still have hundreds of dimensions, posing challenges for standard GP models. GP surrogate models are entirely characterized by the choice of a kernel function (covariance function) $k(z, z')$ that measures the similarity

between two input points z and z' . The choice of kernel function is critical as it encodes our prior beliefs about the function we are trying to model. Many canonical kernels such as RBF (Radial Basis Function) Kernel and Matern Kernel depend on a lengthscale augmented squared Euclidean distance $d(z, z')$ between z and z' i.e.

$$d(z, z') = \sum_i^k \left(\frac{z_i - z'_i}{l_i^2} \right) \quad (4)$$

The lengthscale $\{l_i\}$ is a critical hyper-parameter that captures the smoothness of functions represented by the kernel. In small supervised data settings (as in our problem), careful prior specification for this parameter is critical to achieve good performance on high-dimensional inputs. In order to address this challenge, we consider the recently proposed idea of scaling the prior on the lengthscale hyper-parameters of the GP kernel with the square root of the input dimensionality of the search space [15]. Specifically, for a k -dimensional input space, dimensionality scaled prior for the length-scale parameter l_i for $i \in \{1, 2, \dots, k\}$ is described as:

$$l_i \sim \text{LogNormal}(\mu = \log(\sqrt{k}), \sigma^2 = 1) \quad (5)$$

Overall, after applying random projections, our GP surrogate model is defined on the reduced space \mathcal{Z} with dimensionality scaled priors on the lengthscale of the kernel $k(z, z')$.

4 Experiments and Results

To evaluate the effectiveness of BODE-GEN, we investigate the following research questions (RQs):

- **RQ1:** *How effective is BODE-GEN for code generation?*
- **RQ2:** *How does BODE-GEN compare with CoT and OPRO prompting methods for code generation?*
- **RQ3:** *What changes are introduced by BODE-GEN to the initial prompt and how do they affect the resulting code?*

In what follows, we first describe our experimental setup including benchmark dataset, details of LLMs, configuration of BODE-GEN and baseline methods, and evaluation methodology. Next, we discuss the results to answer the three RQs.

4.1 Experimental Setup

Dataset. We benchmark our proposed BODE-GEN approach on the HumanEval+ benchmark [22] which is a recent extension of the widely-adopted HumanEval benchmark [8] for coding tasks. HumanEval consists of 164 python programming tasks, each containing a function signature and an initial prompt which is written as a docstring. The correctness of each task is measured by evaluating it on a set of pre-specified test cases. HumanEval+ extends the number of test cases in HumanEval tasks by 80x making it a challenging program synthesis benchmark.

Large Language Models. We employ three different LLM_{base} for code generation: one closed-source ChatGPT 3.5 (Turbo) [7], and two open source LLMs CodeLlama (7B) [28] and DeepSeek-Coder-33b [14]. The instruction-tuned version of LLama2 (7B) [34] is used as the Auxiliary LLM LLM_{aux} for all the experiments.

Computing Machine. All our experiments were run on a machine featuring an Nvidia A40 GPU. The system is built on an x86_64 architecture with a 32-core AMD EPYC 7573X processor. The machine has 251 GB of memory.

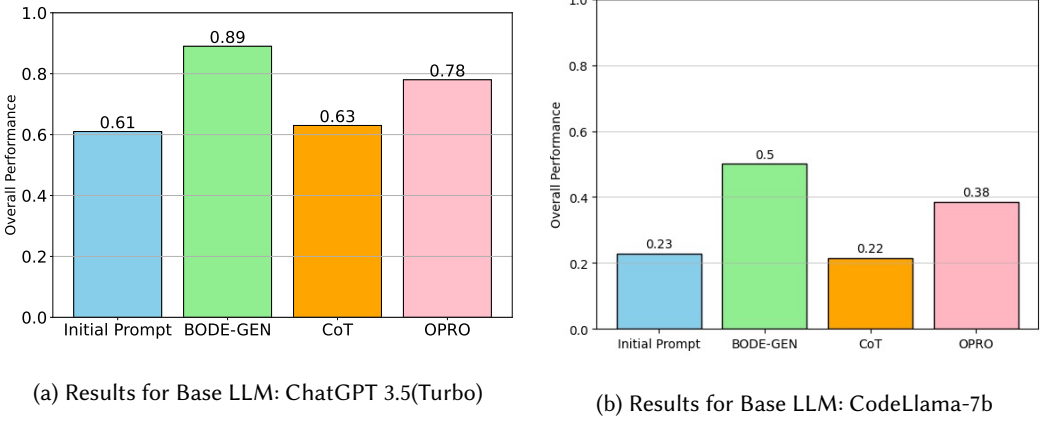


Fig. 2. Results comparing the overall performance of BODE-GEN and baselines with (a) ChatGPT 3.5(Turbo) and (b) CodeLlama-7b as the base LLM. Here, overall performance is computed as the percentage of tasks (out of total 164) from HumanEval+ solved to 100% correctness.

Configuration of BODE-GEN. The search space for BODE-GEN is set to four embeddings, i.e., $E = \{e_1, e_2, \dots, e_4\}$ where each $e_i \in \mathbb{R}^{4096}$. We randomly project each embedding to a 64 dimensional space resulting in a overall search space of $256 = (4 \times 64)$ dimensions. We ran the BO method for $T_{max}=50$ iterations after initializing the surrogate model with 20 randomly picked embedding points (i.e., BO performance curves are shown for 70 iterations). Expected Improvement was picked as the acquisition function due to its practical success without any hyper-parameters and optimized over a discrete set of 10K points. We ran BO for multiple seeds (five) and present results with mean and error bars.

Baselines. We compare BODE-GEN with two baseline methods for code generation: Zero-shot Chain of Thought (CoT) [38] and Optimization by PROMpting (OPRO) [39]. CoT is a popular prompting approach that is shown to illicit reasoning from LLMs for complex tasks. Unlike the original formulation which requires providing multiple intermediate sequence of steps towards solving the problem, we consider the Zero-shot version where we append the key prompt "Let's think step by step" as a suffix to the initial prompt. OPRO is a recent iterative prompt optimization technique that leverages LLMs to iteratively suggest better candidates conditioned on previously found prompts.

Evaluation Metric. For each coding task T , base language model LLM_{base} is queried with a candidate prompt to generate python code. We employ the pass @1 metric defined below [8] to evaluate the correctness of the generated code:

$$\text{pass @1} = \mathbb{E}_{\text{tasks}} \left[1 - \frac{\binom{n-c}{1}}{\binom{n}{1}} \right] \quad (6)$$

where we generate $n = 3$ code samples per prompt to reduce the variance of the metric and c is the number of code samples that pass all unit tests for that task.

4.2 RQ1: Effectiveness of BODE-GEN

We measure the effectiveness of BODE-GEN in two ways: the code generation accuracy aggregated over all coding tasks in the HumanEval+ benchmark (higher the better) and the number of BO iterations to achieve high code generation accuracy (smaller the better). BODE-GEN demonstrates

improved performance in finding prompts that improve the generated code's correctness accuracy compared to the baseline methods. The results, as depicted in Figure 2a, show that BODE-GEN achieves an average code generation accuracy of **0.89 on ChatGPT-3.5Turbo**[7], which is higher than that of initial prompts (0.61), CoT (0.63), and OPRO (0.78). Figure 2b, shows that BODE-GEN archives an average code generation accuracy of **0.5 on CodeLlama-7b**[28], which is higher than that of initial prompts (0.23), CoT (0.22), and OPRO (0.38). Figure 2c, shows that BODE-GEN archives an average code generation accuracy of **0.94 on DeepSeeker-DeepSeek-Coder-33b**[14], which is notably higher than that of initial prompts (0.66), and OPRO (0.78). This increase in performance (code accuracy) suggests the utility of BODE-GEN in finding prompts that generate functionally correct code on the HumanEval+ benchmark.

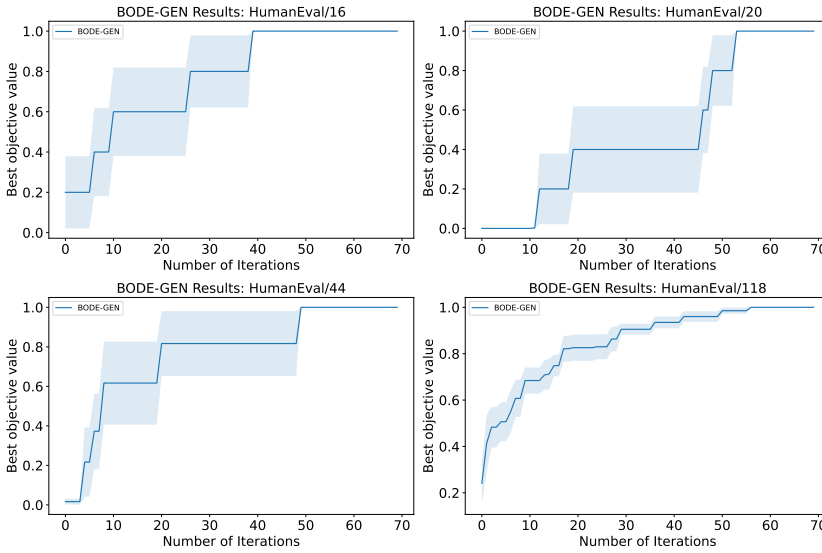


Fig. 3. Results showing BODE-GEN's performance on ChatGPT 3.5(Turbo) as base-LLM as a function of number of iterations (number of base LLM calls with different prompts) on a subset of representative tasks from HumanEval+ benchmark. Note that the objective value for BO is the percentage of test cases passed by the generated code for a given coding task. As shown in the figure, prompts suggested by BODE-GEN are often able to reach high 100% code generation correctness. Each BO iteration corresponds to roughly one query to the base LLM (precisely it is three queries per iteration since we generate three samples for each prompt to compute pass @1).

We also show the progress of BODE-GEN in terms of the generated code's accuracy as a function of number of BO iterations (i.e., number of tried prompts) in Figure 3, 4, 5 on some representative coding tasks noting that our findings are similar on other coding tasks. In all cases, BODE-GEN demonstrates improvement in the objective value (percentage of test cases passed) within the first 20-30 BO iterations, often achieving near-optimal performance by the 50th iteration. This convergence behavior suggests that BODE-GEN is able to explore the prompt space to find optimized prompts for code generation. As evident from the figures, BODE-GEN iteratively finds better prompts that are able to reach 100% correctness accuracy.

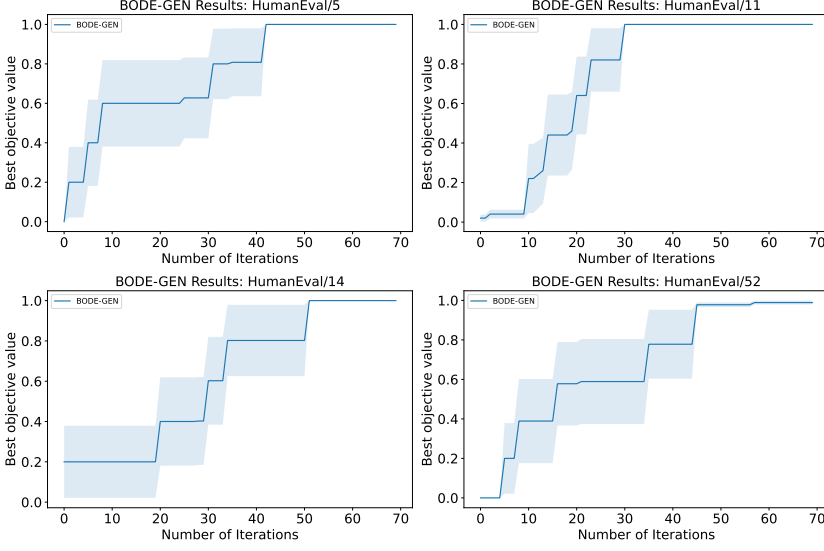


Fig. 4. Results showing BODE-GEN's performance on CodeLlama-7b as base-LLM as a function of number of iterations (number of base LLM calls with different prompts) on a subset of representative tasks from HumanEval+ benchmark.

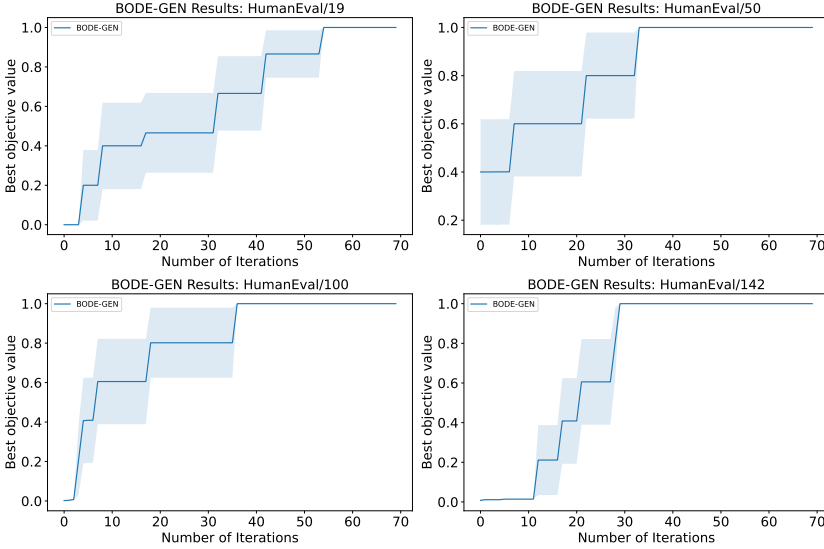


Fig. 5. Results showing BODE-GEN's performance on DeepSeek-Coder as base-LLM as a function of number of iterations (number of base LLM calls with different prompts) on a subset of representative tasks from HumanEval+ benchmark.

4.3 RQ2: Comparison to Baselines CoT and OPRO

In comparison to CoT and OPRO, BODE-GEN outperforms both methods in our setup as shown in Figure 2a, and 2b. CoT shows only small improvement over the initial prompt and performs worse than BODE-GEN. OPRO achieves better results compared to CoT since it is an iterative approach (similar to BODE-GEN) that finds better prompts iteratively. We define a notion of task difficulty as

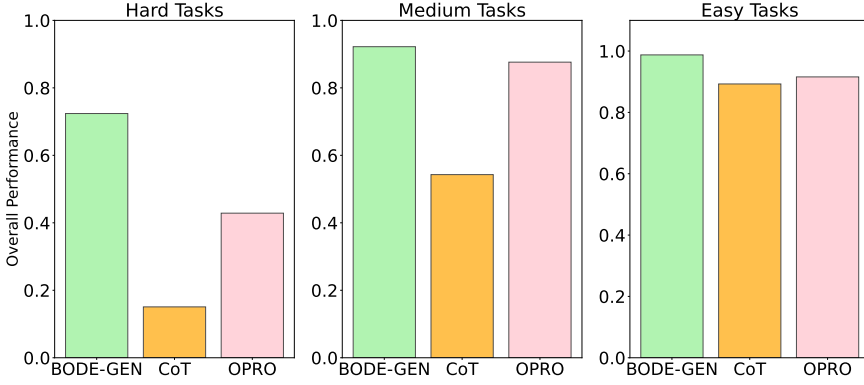


Fig. 6. Results comparing the overall performance of BODE-GEN with zero-shot CoT and OPRO with ChatGPT 3.5(Turbo) as base-LLM on a grouping of HumanEval+ tasks based on a notion of difficulty measured as the correctness of the code generated by initial prompts given for each task. For example, the easy/hard class refers to all tasks for which the code generated via initial prompt achieves correctness (above 67%/below 30%) respectively. The medium class contains all tasks with their initial prompts' correctness between 30-67%.

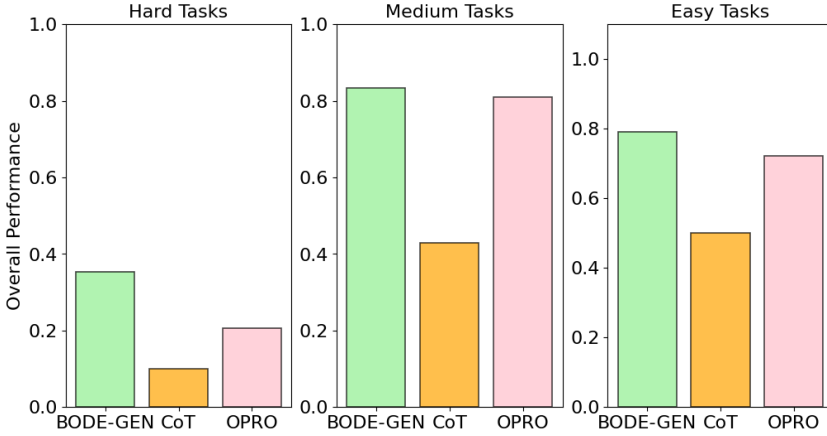


Fig. 7. Results of BODE-GEN and baselines with CodeLlama-7b as base-LLM on a grouping of HumanEval+ tasks based on a notion of difficulty measured as the correctness of the code generated by initial prompts given for each task. For example, the easy/hard class refers to all tasks for which the code generated via initial prompt achieves correctness (above 67%/below 30%) respectively. The medium class contains all tasks with their initial prompts' correctness between 30-67%.

the correctness of the generated code achieved by initial fixed prompt given for each task and create three groups (easy, medium, and hard) of increasing difficult to evaluate all three methods. The results are shown in Figure 6, 7 which shows that the gap between BODE-GEN and the performance of baselines increases as we increase the difficulty of the task. This means that BODE-GEN is more effective than baselines for hard coding tasks.

4.4 RQ3: Qualitative Analysis

To understand how our BO-based prompt optimization approach improves LLM-driven program synthesis in terms of the correctness of the resulting code, we conducted qualitative analysis on 15 randomly selected coding task cases among those in which the correctness improvement was the most substantial (i.e., the most challenging cases for the LLMs with the *original prompts* used in the HumanEval+ dataset [21]). By examining these cases, we aim to (1) identify the key changes in the prompt that our optimization approach makes, and (2) based on these changes and how the generated code differs between the initial prompt and optimized prompt, distill common patterns of and main insights into what makes a prompt effective leading LLMs to produce correct code.

4.4.1 Key Changes Induced in Prompt Optimizations. For each of these chosen cases, we carefully compare the two versions of the code-generation prompt, aiming to identify key differences between them in all possible aspects that may affect the base LLM’s ability to generate correct code.

Change 1: Use of Examples. While providing (e.g., input/output) examples generally help LLMs generate correct code, how the examples are used in the prompt matter. The original prompts typically include examples within the docstring, which might be less visible. In optimized prompts from BO, the examples are clearly separated from the main instructions, making them more prominent and easier to reference. For example, in the case of the `Multiply` function, shown in Figure 8, the original prompt provides examples within the docstring, versus our optimized prompt listing the examples in plain text after describing the task. By clearly separating instructions from examples rather than combining them in one information block, these changes improve LLMs’ code correctness through *separation of concerns*.

Original prompt	BO-optimized prompt
<pre>def multiply(a, b): """Complete the function that takes two integers and returns the product of their unit digits. Assume the input is always valid. Examples: multiply(148, 412) should return 16. multiply(19, 28) should return 72. multiply(2020, 1851) should return 0. multiply(14,-15) should return 20.</pre>	<pre>Write a function named `multiply` that takes two integer arguments `a` and `b`. The function should return the product of the unit digits of `a` and `b`. For example, `multiply(148, 412)` should return `16`, `multiply(19, 28)` should return `72`, `multiply(2020, 1851)` should return `0`, and `multiply(14, -15)` should return `20`. Ensure that the input values are always valid and handle edge cases appropriately.</pre>

Fig. 8. Original prompt versus BO-optimized prompt for generating the `multiply` function (HumanEval+ Case 97).

Change 2: Instruction Clarity and Detail. In several studied cases, the original prompt often includes instructions within the code’s docstring, which might be concise but less explicit. In the optimized version, the prompt provides a detailed, narrative description of the task, often in plain English, making the instructions clearer and more explicit. For example, in the case of `max_fill` function generation, as shown in Figure 9, the original prompt (left) describes the task as “You are given a rectangular grid of wells..”, while the improved prompt provides more details and clarity with “Please write a function called `max_fill` that takes in a rectangular grid of wells and a capacity as input...”.

Change 3: Structure and Organization. In most of the studied cases, the original prompts have a less organized structure, while the optimization-resulted prompts use a more structured, step-by-step format, often with clear separations between different parts of the instructions. As seen in Figure 10, for generating the `get_closest_vowel` function, our optimized prompt improves

Original prompt	BO-optimized prompt
<pre>import math def max_fill(grid, capacity): """ You are given a rectangular grid of wells. Each row represents a single well, and each 1 in a row represents a single unit of water. Each well has a corresponding bucket that can be used to extract water from it, and all buckets have the same capacity. Your task is to use the buckets to empty the wells. Output the number of times you need to lower the buckets. Example 1: Input: grid : [[0,0,1,0], [0,1,0,0], [1,1,1,1]] bucket_capacity : 1 Output: 6 Example 2: Input: grid : [[0,0,1,1], [0,0,0,0], [1,1,1,1], [0,1,1,1]] bucket_capacity: 2 Output: 5 Example 3: Input: grid : [[0,0,0], [0,0,0]] bucket_capacity : 5 Output: 0 Constraints: * all wells have the same length * 1 <= grid.length <= 10^2 * 1 <= grid[:,1].length <= 10^2 * grid[i][j] -> 0 1 * 1 <= capacity <= 10 """</pre>	<p>Please write a function called <code>'max_fill'</code> that takes in a rectangular grid of wells and a capacity as input, and outputs the number of times you need to lower the buckets to empty the wells. The grid is represented as a 2D array of booleans, where each row represents a single well and each 1 in a row represents a single unit of water. Each well has a corresponding bucket that can be used to extract water from it, and all buckets have the same capacity. Your task is to use the buckets to empty the wells. The function should solve the problem by iterating through the rows of the grid, keeping track of the number of times each bucket needs to be lowered to empty each well, and returning the total number of bucket lowerings needed. The function should take into account the constraints provided in the problem statement, including the length of the grid, the length of each row, and the capacity of each bucket.</p>

Fig. 9. Original prompt versus BO-optimized prompt for generating the `max_fill` function (HumanEval+ Case 115).
how the entire prompt is structured: both the assumptions and examples are clearly organized in addition to the task description.

Original prompt	BO-optimized prompt
<pre>def get_closest_vowel(word): """You are given a word. Your task is to find the closest vowel that stands between two consonants from the right side of the word (case sensitive). Vowels in the beginning and ending doesn't count. Return empty string if you didn't find any vowel met the above condition. You may assume that the given string contains English letter only. Example: get_closest_vowel("yogurt") ==> "u" get_closest_vowel("FULL") ==> "U" get_closest_vowel("quick") ==> "" get_closest_vowel("ab") ==> "" ""</pre>	<p>Given a word, your task is to find the closest vowel that is positioned between two consonants from the right side of the word, taking into account case sensitivity. Vowels at the beginning or end of the word do not count. Return an empty string if no vowel meets the specified condition.</p> <p>You may assume that the input string consists solely of English letters.</p> <p>Example: <code>get_closest_vowel("yogurt") ==> "u"</code> Example: <code>get_closest_vowel("FULL") ==> "U"</code> Example: <code>get_closest_vowel("quick") ==> ""</code> Example: <code>get_closest_vowel("ab") ==> ""</code></p>

Fig. 10. Original prompt versus BO-optimized prompt for generating the `get_closest_vowel` function (HumanEval+ Case 118).

Change 4: Language and Readability. Another main prompt change induced by our optimization lies in the language use in the prompt that affects its readability. In particular, the original prompts tend to use a technical and formal style typical of in-code documentation, while the optimized

prompts use plain English and a more narrative style, improving readability and comprehension. Take the `max_fill` function (Figure 9) as an example again, the improved version of the prompt uses a narrative and explanatory style, versus the more technical style of the original prompt. Intuitively, LLMs are trained on more natural language corpus than technical documents, which may justify why these changes that improve readability help produce LLMs produce more correct code.

Original prompt	BO-optimized prompt																														
<pre>def numerical_letter_grade(grades): """It is the last week of the semester and the teacher has to give the grades to students. The teacher has been making her own algorithm for grading. The only problem is, she has lost the code she used for grading. She has given you a list of GPAs for some students and you have to write a function that can output a list of letter grades using the following table: GPA Letter grade 4.0 A+ > 3.7 A > 3.3 A- > 3.0 B+ > 2.7 B > 2.3 B- > 2.0 C+ > 1.7 C > 1.3 C- > 1.0 D+ > 0.7 D > 0.0 D- 0.0 E Example: grade_equation([4.0, 3, 1.7, 2, 3.5]) ==> ['A+', 'B', 'C-', 'C', 'A-'] """</pre>	<p>Given a list of GPAs (Grade Point Average) for a group of students, generate a function that outputs a list of letter grades corresponding to the GPA values. The function should use a table of letter grades and their corresponding GPA values to determine the appropriate letter grade for each GPA. The table should be represented as follows:</p> <table><tr><th>GPA</th><th>Letter Grade</th></tr><tr><td>---</td><td>---</td></tr><tr><td>4.0</td><td>A+</td></tr><tr><td>> 3.7</td><td>A</td></tr><tr><td>> 3.3</td><td>A-</td></tr><tr><td>> 3.0</td><td>B+</td></tr><tr><td>> 2.7</td><td>B</td></tr><tr><td>> 2.3</td><td>B-</td></tr><tr><td>> 2.0</td><td>C+</td></tr><tr><td>> 1.7</td><td>C</td></tr><tr><td>> 1.3</td><td>C-</td></tr><tr><td>> 1.0</td><td>D+</td></tr><tr><td>> 0.7</td><td>D</td></tr><tr><td>> 0.0</td><td>D-</td></tr><tr><td>0.0</td><td>E</td></tr></table> <p>Example: <code>numerical_letter_grade([4.0, 3, 1.7, 2, 3.5]) ==> ['A+', 'B', 'C-', 'C', 'A-']</code></p>	GPA	Letter Grade	---	---	4.0	A+	> 3.7	A	> 3.3	A-	> 3.0	B+	> 2.7	B	> 2.3	B-	> 2.0	C+	> 1.7	C	> 1.3	C-	> 1.0	D+	> 0.7	D	> 0.0	D-	0.0	E
GPA	Letter Grade																														
---	---																														
4.0	A+																														
> 3.7	A																														
> 3.3	A-																														
> 3.0	B+																														
> 2.7	B																														
> 2.3	B-																														
> 2.0	C+																														
> 1.7	C																														
> 1.3	C-																														
> 1.0	D+																														
> 0.7	D																														
> 0.0	D-																														
0.0	E																														

Fig. 11. Original prompt versus BO-optimized prompt for generating the `numerical_letter_grade` function (HumanEval+ Case 81).

Change 5: Descriptive Guidance. The original prompts often provide guidance on what the function should do in a general manner. In contrast, the optimized prompt has more specific guidance on how to approach the task, including iterative processes or specific logic to follow. As shown in Figure 9, if we put aside the detailed examples, the original task description itself is overall general. In the optimized version, much more specific guidance is included on iterating through the grid and tracking bucket usage, which helped the LLMs correct the errors in the code generated with the original prompt, even without using those examples.

To summarize, the common differences between the two versions of each prompt in our studied cases highlight the importance of clarity, detail, explicitness, and structure in writing prompts for code generation. By ensuring that *instructions are clear, detailed, and well-organized*, and by *providing explicit guidance and examples*, the quality and accuracy of the generated code can be significantly improved.

4.4.2 Patterns of Correct-Code-Generating Prompts. Based on these prompt changes made by our optimization approach as summarized above, together with comparing the code generated by the

two prompt versions, we further identified the following major patterns of correct-code-generating prompts.

Pattern 1: Clear and Detailed Instructions. We observed that providing detailed and explicit instructions helps LLMs understand the requirements better, leading to more accurate code generation. In the case of `max_fill` function (Figure 9), for example, the original prompt states “Your task is to use the buckets to empty the wells”, which is much less elaborate and explicit than the optimized version: “The function should solve the problem by iterating through the rows of the grid, keeping track of the number of times each bucket needs to be lowered to empty each well, and returning the total number of bucket lowerings needed.”

The detailed instructions in the optimized prompt clarify the method to solve the problem, resulting in better code generation.

Pattern 2: Step-by-Step Breakdown. From multiple cases, it appears clear that breaking down the task into smaller, clear steps helps the model generate code that follows the intended logic more closely. One example is found in the prompt for synthesizing the `numerical_letter_grade` function, a shown in Figure 11. The original version simply provides a context and a list of GPAs, offering no intermediate reasoning. The optimized version of the prompt clearly describes the function’s purpose, parameters, expected behavior, and gives a table for GPA to letter grade mapping. *The step-by-step breakdown ensures that the model accurately follows the logic needed to map GPAs to letter grades.*

Pattern 3: Explicit Constraints and Assumptions. Inspection of the optimized prompts reveals that specifying constraints and assumptions ensures LLMs adhere to the necessary conditions and edge cases. As illustrated in Figure 10, the original prompt does provide assumptions and examples but not explicitly (i.e., within a docstring), while the optimized version clearly lists the constraints and examples in plain text (i.e., explicitly). This shows that *explicitly mentioning constraints and assumptions helps the model generate code that respects those conditions.*

Pattern 4: Separation of Examples and Instructions. Almost all of our studied cases confirm that including (e.g., input/output) examples in code-generating prompts help LLMs produce correct code, but it is essential to separate examples from instructions. The separation ensures clarity and helps the model focus on understanding both the requirements and the examples. As shown in Figure 10, the optimized prompt, which lists examples separately after the task description, improves over the original prompt which combines examples and the task description (within the docstring). Apparently, *clearly separating examples from instructions enhances LLMs’ ability to parse and understand both parts effectively.*

Pattern 5: Plain English Descriptions. As evidenced in many cases, using plain English to describe the task makes it easier for LLMs to parse the requirements and generate correct code. As seen in Figure 9, the original prompt uses a technical and formal style, while the improved prompt provides a narrative and explanatory style. *Plain English task descriptions make it easier for the model to understand the task and generate correct logic.*

In summary, to write better prompts for code generation with LLMs, it is crucial to *provide clear, detailed instructions, break down tasks into smaller steps, specify constraints and assumptions explicitly, use multiple examples, describe tasks in plain English, and separate examples from instructions.* These patterns help the model understand the requirements better and generate more accurate and correct code.

4.5 Additional Results on DeepSeek-Coder

In this section, we present additional results comparing BODE-GEN with the baseline approaches on DeepSeek-Coder base LLMs. As shown in Figure 12, with DeepSeek-Coder as base LLM, BODE-GEN

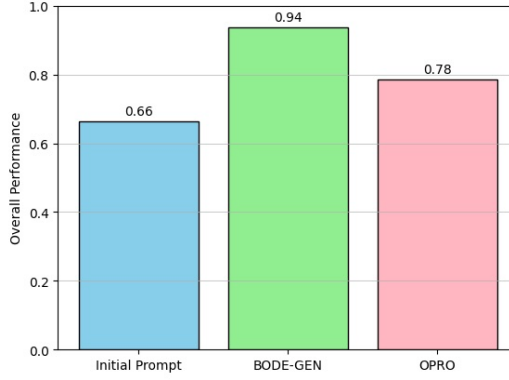


Fig. 12. Results comparing the overall performance of BODE-GEN with other baselines on the HumanEval+ benchmark with DeepSeek-Coder as the base LLM (analogous to Figure 2).

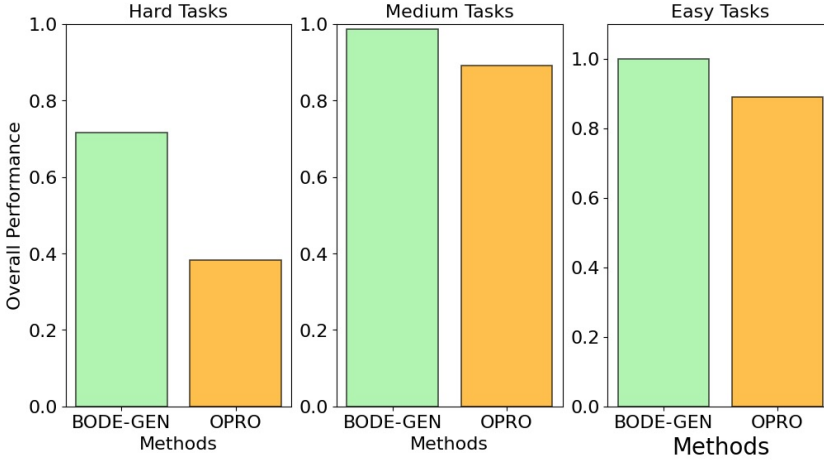


Fig. 13. Results comparing the overall performance of BODE-GEN with OPRO with DeepSeek-Coder as base-LLM on a grouping of HumanEval+ tasks based on a notion of difficulty measured as the correctness of the code generated by initial prompts given for each task (analogous to Figure 2).

outperforms OPRO in finding prompts that generate code with higher correctness and the gap between the baselines and BODE-GEN increases on difficult tasks as shown in Figure 13. Since CoT was consistently performing worse than BODE-GEN and OPRO, we do not include it in this comparison. Similar to Figure 3 in the main paper, we also show the progress of BODE-GEN as a function of number of BO iterations with DeepSeek-Coder (Figure 5).

5 Threats to Validity

The main internal validity threat lies in possible implementation errors in our tool and experimental scripts. To minimize this threat, we have performed careful code reviews of these implementations against manageable testing scenarios. Another major issue LLMs are commonly subject to are hallucination [33], which in our study may cause inconsistent and unreliable prompt improvements. To deal with such issues, we set the temperature of the auxiliary LLM to zero in our experiments

and compute correctness accuracy (pass @1) by generating multiple samples from the base LLM. We also ran each experiment multiple times and took results that are consistent among the runs. Nevertheless, this mitigation may not have entirely ruled out the possibility that developers may not always get the same correctness improvement when using our technique for prompt optimization, or they may need slightly more iterations to get the same improvement as shown in our evaluation.

The main external validity threat concerns the datasets we used, as well as LLMs and baseline approaches chosen, in our experiments. We used a benchmark (HumanEval+) that is popularly used in LLM evaluations. Yet the prompts in it may not represent real-world code-synthesis prompts developers use. Meanwhile, we chose capable LLMs that are affordable to use with respect to the scale of our experiments. With the rapid evolution of LLMs, more advanced models may not guaranteed to be improved as much as what we present. Similarly, we selected CoT as one of our baselines, which is a state-of-the-art prompting strategy on LLMs. Even more advanced prompting could have achieved better performance (although how to instantiate them for code generation remains an open problem as of now).

6 Related Work

We provided an overview of the broader literature on LLMs for code generation in the introduction section. Below we discuss closely related work to our specific problem setting.

LLM-driven Code Generation. To address the challenge that LLMs often generate incorrect code [3, 21], prior work has explored various strategies. One category includes knowledge augmentation [2, 16] and reasoning elicitation [1, 23]. Chain-of-thoughts (CoT) prompting [38] has been shown to be useful to improve LLM-driven code generation [23]. However, our experiments on HumanEval benchmark show that the improvement with CoT is small. Jigsaw [16] performs post-processing of the generated code to check and calibrate correctness, by augmenting the LLMs with knowledge about the syntax and semantics of programs based on program analysis and synthesis techniques. This approach is not automatic and maybe incomplete for some coding tasks. Ahmed et al., [1] leverage the self-consistency technique, which has previously been shown promising for improving the reasoning capabilities of LLMs via CoT, for generating bug fixes by leveraging associated commit logs as explanations. However, it is not clear how to apply this method for synthesizing programs from scratch. Recent work by Murr et al., [24] found that the specificity of prompts has a major impact on the quality of LLM-generated code: more specific prompts tend to produce accurate code meeting the functional requirements although compromising the generation diversity. However, determining the specificity level is currently a manual process.

A recent prompting technique, AceCoder [19], is proposed particularly for code generation with LLMs. It starts with asking the LLM to analyze the given requirements and output an intermediate preliminary output (e.g., test cases), which is then utilized to retrieve similar programs that meet the requirements as exemplars in the code-generation prompt. However, such code examples may not always be available for a given arbitrary requirement. Following a different strategy, AlphaCodium [27] is proposed to leverage test-based iterative refinement of the code generation process based on LLMs, where the tests are from what is already available and generated by the AI models. The effectiveness of AlphaCodium relies on the quality and availability of such tests, which itself is an unresolved challenge. Tao et al. [32] employs grammar-guided evolutionary search to improve LLM-based program synthesis. It uses the LLM-generated code as an initial input for the evolutionary search process after a grammar-mapping phase, which allows for program development and fixing errors. Subsequently, it employs different similarity metrics related to the LLM-generated code to steer the multi-objective evolutionary search process. However, the

evaluation is only limited to grammar validation of the generated code as opposed to functional correctness.

Optimization by PRoMpting (OPRO) is the state-of-the-art technique for optimizing tasks specified in natural language. OPRO uses LLMs to generate new solutions based on previously evaluated solutions and their scores, iteratively improving the objective function. Our experiments on HumanEval benchmark demonstrates that OPRO is better than CoT but it is less effective than our proposed BODE-GEN approach in terms of code generation accuracy.

In contrast to prior work, our proposed BODE-GEN approach for optimizing prompts is *fully automatic* (no manual intervention) and *sample-efficient* (minimized the number of prompts tried) by leveraging the test cases for coding task as per the test-driven software development paradigm [5].

Bayesian Optimization. BO has shown a lot of success in optimizing continuous spaces with small number of dimensions (typically less than 50). There is relatively less work on BO over discrete spaces which is significantly challenging than BO over continuous spaces [10] and is limited to a small number of dimensions. We solve the problem of BO over discrete prompt space by reducing it to BO over high-dimensional continuous embedding space using an auxiliary LLM. We study algorithmic solutions to build effective surrogate models in the small training data setting to handle challenges of the high-dimensional continuous embedding space.

7 Summary and Future Work

This paper developed and studied BODE-GEN, a Bayesian optimization approach for prompt search in large language models (LLMs) for code generation tasks. We formulate this search as optimization over high-dimensional continuous embedding space of an auxiliary LLM. By leveraging random projections and dimensionality-scaled priors, our method addresses the challenges of high dimensionality in the embedding space. Experiments on the HumanEval+ benchmark demonstrate BODE-GEN’s effectiveness to significantly improve code generation accuracy across a wide variety of coding tasks when compared to strong baselines. As LLMs continue to play a crucial role in software development, BODE-GEN offers a framework for improving their reliability and effectiveness, ultimately improving developer’s productivity. Future work includes reducing the dependency of BODE-GEN on the requirement of test cases for coding tasks.

References

- [1] Toufique Ahmed and Premkumar Devanbu. 2023. Better patching using LLM prompting, via Self-Consistency. In *2023 38th IEEE/ACM International Conference on Automated Software Engineering (ASE)*. IEEE, 1742–1746.
- [2] Toufique Ahmed, Kunal Suresh Pai, Premkumar Devanbu, and Earl Barr. 2024. Automatic semantic augmentation of language model prompts (for code summarization). In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–13.
- [3] Jacob Austin, Augustus Odena, Max Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. 2021. Program Synthesis with Large Language Models. *arXiv preprint arXiv:2108.07732* (2021).
- [4] Agathe Balayn, Mireia Yurrita, Fanny Rancourt, Fabio Casati, and Ujwal Gadiraju. 2024. An Empirical Exploration of Trust Dynamics in LLM Supply Chains. *arXiv preprint arXiv:2405.16310* (2024).
- [5] Kent Beck. 2002. *Test-Driven Development by Example*. Addison Wesley. ISBN 978-0-321-14653-3.
- [6] Kent Beck. 2003. *Test-driven development: by example*. Addison-Wesley Professional.
- [7] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems*. 1877–1901.
- [8] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Neel Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Michael Pavlov, Alethea Power, Lukasz Kaiser,

- Mohammad Bavarian, Clemens Winter, Phil Tillet, Felipe Petroski Such, David Cummings, Matthias Plappert, Filippos Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William H Guss, Alex Nichol, Alexander Paino, Nikolai Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, Andrew Saunders, Brandon Houghton, Jacob Pfau, Diego de Las Casas, Leon Bottou, Charles Choi, Adam Coates, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2021. Evaluating Large Language Models Trained on Code. *arXiv preprint arXiv:2107.03374* (2021).
- [9] Tristan Coignion, Clément Quinton, and Romain Rouvoy. 2024. A Performance Study of LLM-Generated Code on LeetCode. In *Proceedings of the 28th International Conference on Evaluation and Assessment in Software Engineering*. 79–89.
- [10] Aryan Deshwal, Syrine Belakaria, and Janardhan Rao Doppa. 2021. Mercer features for efficient combinatorial Bayesian optimization. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 7210–7218.
- [11] Robert J Ellison, John B Goodenough, Charles B Weinstock, and Carol Woody. 2010. Evaluating and mitigating software supply chain security risks. *Software Engineering Institute, Tech. Rep. CMU/SEI-2010-TN-016* (2010).
- [12] William Enck and Laurie Williams. 2022. Top five challenges in software supply chain security: Observations from 30 industry and government organizations. *IEEE Security & Privacy* 20, 2 (2022), 96–100.
- [13] Tianyu Gao, Adam Fisch, and Danqi Chen. 2021. Making Pre-trained Language Models Better Few-shot Learners. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics (ACL)*. 3816–3830.
- [14] Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. 2024. DeepSeek-Coder: When the Large Language Model Meets Programming—The Rise of Code Intelligence. *arXiv preprint arXiv:2401.14196* (2024).
- [15] Carl Hvarfner, Erik Orm Hellsten, and Luigi Nardi. 2024. *Vanilla Bayesian Optimization Performs Great in High Dimensions*. <https://openreview.net/forum?id=OfT8MgIqHT>
- [16] Naman Jain, Skanda Vaidyanath, Arun Iyer, Nagarajan Natarajan, Suresh Parthasarathy, Sriram Rajamani, and Rahul Sharma. 2022. Jigsaw: Large language models meet program synthesis. In *Proceedings of the 44th International Conference on Software Engineering*. 1219–1231.
- [17] Kasper Green Larsen and Jelani Nelson. 2017. Optimality of the Johnson-Lindenstrauss lemma. In *2017 IEEE 58th Annual Symposium on Foundations of Computer Science (FOCS)*. IEEE, 633–638.
- [18] Ben Letham, Roberto Calandra, Akshara Rai, and Eytan Bakshy. 2020. Re-examining linear embeddings for high-dimensional Bayesian optimization. *Advances in neural information processing systems* 33 (2020), 1546–1558.
- [19] Jia Li, Yunfei Zhao, Yongmin Li, Ge Li, and Zhi Jin. 2024. AceCoder: An Effective Prompting Technique Specialized in Code Generation. *ACM Transactions on Software Engineering and Methodology* (2024).
- [20] Fang Liu, Yang Liu, Lin Shi, Houkun Huang, Ruifeng Wang, Zhen Yang, and Li Zhang. 2024. Exploring and evaluating hallucinations in llm-powered code generation. *arXiv preprint arXiv:2404.00971* (2024).
- [21] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2024. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems* 36 (2024).
- [22] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. 2024. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation. *Advances in Neural Information Processing Systems* 36 (2024).
- [23] Yingwei Ma, Yue Yu, Shanshan Li, Yu Jiang, Yong Guo, Yuanliang Zhang, Yutao Xie, and Xiangke Liao. 2023. Bridging Code Semantic and LLMs: Semantic Chain-of-Thought Prompting for Code Generation. *arXiv preprint arXiv:2310.10698* (2023).
- [24] Lincoln Murr, Morgan Grainger, and David Gao. 2023. Testing LLMs on Code Generation with Varying Levels of Prompt Specificity. *arXiv preprint arXiv:2311.07599* (2023).
- [25] Daye Nam, Andrew Macvean, Vincent Hellendoorn, Bogdan Vasilescu, and Brad Myers. 2024. Using an llm to help with code understanding. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*. 1–13.
- [26] Carl Edward Rasmussen and Christopher K. I. Williams. 2006. *Gaussian processes for machine learning*. MIT Press.
- [27] Tal Ridnik, Dedy Kredo, and Itamar Friedman. 2024. Code generation with alphacodium: From prompt engineering to flow engineering. *arXiv preprint arXiv:2401.08500* (2024).
- [28] Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, et al. 2023. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950* (2023).
- [29] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P Adams, and Nando De Freitas. 2016. Taking the Human Out of the Loop: A Review of Bayesian Optimization. *Proc. IEEE* 104, 1 (2016), 148–175.
- [30] Dominik Sobania, Justyna Petke, Martin Briesch, and Franz Rothlauf. 2024. A Comparison of Large Language Models and Genetic Programming for Program Synthesis. *IEEE Transactions on Evolutionary Computation* (2024).
- [31] Claudio Spiess, David Gros, Kunal Suresh Pai, Michael Pradel, Md Rafiqul Islam Rabin, Susmit Jha, Prem Devanbu, and Toufique Ahmed. 2024. Quality and Trust in LLM-generated Code. *arXiv preprint arXiv:2402.02047* (2024).

- [32] Ning Tao, Anthony Ventresque, Vivek Nallur, and Takfarinas Saber. 2024. Enhancing Program Synthesis with Large Language Models Using Many-Objective Grammar-Guided Genetic Programming. *Algorithms* 17, 7 (2024), 287.
- [33] Yuchen Tian, Weixiang Yan, Qian Yang, Qian Chen, Wen Wang, Ziyang Luo, and Lei Ma. 2024. CodeHalu: Code Hallucinations in LLMs Driven by Execution-based Verification. *arXiv preprint arXiv:2405.00253* (2024).
- [34] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. 2023. Llama: Open and efficient foundation language models. *arXiv preprint arXiv:2302.13971* (2023).
- [35] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. 2023. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288* (2023).
- [36] Pengcheng Wang, Richard Shin, Xiaodong Liu, Yao Jin, Prafulla Sharma, Nitish Keskar, Grace Fung, Mayur Naik, and Sameer Yu. 2021. Codex: A Large-Scale Neural Network Model for Code Generation. *arXiv preprint arXiv:2106.01482* (2021).
- [37] Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171* (2022).
- [38] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems* 35 (2022), 24824–24837.
- [39] Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V Le, Denny Zhou, and Xinyun Chen. 2024. Large Language Models as Optimizers. In *The Twelfth International Conference on Learning Representations*. <https://openreview.net/forum?id=Bb4VGOWELI>