DemoCraft: Using In-Context Learning to Improve Code Generation in Large Language Models

Kapu Nirmal Joshua Department of Electrical Engineering Indian Institute of Technology Kanpur nirmalj21@iitk.ac.in

Abstract—Generating executable code from natural language instructions using Large Language Models (LLMs) poses challenges such as semantic ambiguity and understanding taskspecific contexts. To address these issues, we propose a system called *DemoCraft*, which enhances code generation by leveraging in-context learning and demonstration selection, combined with latent concept learning. Latent concept learning introduces additional concept tokens, which are trainable embeddings that capture task-specific knowledge. We then test our system on two major datasets: *MBPP* and *Humaneval*. Our experimental results demonstrate that the proposed system achieves an approximate 2x increase in the *pass@k* metric compared to baseline models. Furthermore, we introduce two novel evaluation metrics: *correctness@k* and *similarity@k*. Our empirical studies indicate that our system attains nearly a 3x improvement in these metrics as well.

Index Terms—in-context learning, code generation, latent concept learning, demonstration selection, large language models

I. INTRODUCTION

The problem of generating code from natural language using Large Language Models (LLMs) involves creating systems capable of translating human language instructions into executable code accurately. This requires the LLM to understand the semantics of the natural language input, grasp the intent behind the instructions, and convert it into syntactically correct and functional code in a specified programming language. Key challenges include handling ambiguous or imprecise language, ensuring the generated code is both correct and efficient, and covering a wide range of programming scenarios and languages.



Fig. 1. Large Language Models struggling at Code Generation

Code generation remains a significant challenge for large language models, as evidenced by Google's *AlphaCode* [1], developed specifically for competitive programming tasks. When evaluated on the *CodeContests* benchmark, AlphaCode achieves a maximum Codeforces rating of only **1238**, placing Mihit Sreejith

Department of Computer Science and Engineering Indian Institute of Technology Guwahati s.mihit@iitg.ac.in

it in approximately the top 28th percentile. Furthermore, a comprehensive survey on code generation using large language models [2] reports a maximum *pass@1* rate of around 30%. These studies have been conducted under zero-shot conditions, highlighting the necessity for few-shot learning approaches. Few-shot learning allows models to leverage relevant demonstrations associated with the prompt prior to generating the output, potentially improving performance.

II. PROBLEM BEHIND SELECTING DEMONSTRATIONS

In-context learning operates by pre-pending a series of demonstrations—examples of prompts and corresponding answers—before the final prompt that the model needs to solve. This setup effectively guides the model, allowing it to leverage patterns from prior examples to generate improved responses. By selecting demonstrations that closely match the problem at hand, we can significantly enhance the model's performance on complex tasks like code generation.



Fig. 2. Few Shot Learning Pipeline

However, selecting relevant demonstrations is a challenging task in itself. Semantic similarity-based selection, a commonly used approach, attempts to identify demonstrations that share high textual similarity with the prompt. While this method may capture surface-level relationships, it often fails to consider the deeper task requirements.

For instance, in competitive programming contexts like Codeforces, problem statements frequently involve recurring character names like "Alice" and "Bob," often engaging in a hypothetical game. A semantic similarity-based approach might assume that any problem mentioning "Alice and Bob playing a game" is contextually relevant to another problem with similar phrasing. However, while these problems may seem alike in language, they can differ significantly in their underlying algorithms. One "Alice and Bob" problem may require a dynamic programming approach, while another could involve graph theory or combinatorial analysis. As a result, semantically similar demonstrations might mislead the model, offering examples that match the language but fail to provide the right procedural insights.

This is where our system, *DemoCraft*, becomes instrumental. *DemoCraft* utilizes a latent concept-based selection algorithm to analyze and select demonstrations that are aligned not only in linguistic features but also in conceptual depth. By focusing on the intrinsic structure of computational problems, *DemoCraft* identifies demonstrations that share the same reasoning paradigms or algorithmic strategies necessary to solve the target prompt. For instance, when presented with a complex binary search or dynamic programming problem, *DemoCraft* is capable of prioritizing demonstrations that involve these specific techniques over those with mere superficial similarity, thereby ensuring that the model is provided with the most contextually relevant guidance.



Fig. 3. Demonstration Selection with Latent Concept Learning

III. DEMOCRAFT: SYSTEM DETAILS

In this section, we provide a detailed technical description of our system architecture, which consists of three primary components: the Latent Concept Learning module, the Task Concept Probability Calculation module, and the Demonstration Selector.

A. Latent Concept Learning

In this stage, we introduce additional tokens [6], referred to as *concept tokens*, to enable the model to learn task-specific features for a given task. These concept tokens function as specialized units within the language model, representing knowledge specific to the task. Incorporating these tokens allows the model to predict the structure and requirements of the task more effectively.

We aim to find the optimal value of the variable θ_d for each task d in the set of tasks T. The variable θ_d , referred to as the *latent concept variable*, is intended to capture the essential characteristics of each task to maximize the model's predictive accuracy. Mathematically, the optimal θ_d maximizes the probability of the correct output given the input, achieved through the Bayes optimal classifier defined as

$$\theta_d = \arg\max_{\theta_d} P^d_M(Y \mid \theta_d, X) \tag{1}$$

where $P_M^d(Y \mid \theta_d, X)$ is the probability that the model M assigns to the output Y given the input X and task-specific variable θ_d .

To train the model to make better predictions, we aim to find θ_d that minimizes the cross-entropy loss. This involves minimizing the negative expected log probability:

$$\hat{\theta}_d = \arg\min_{\theta_d} -\mathbb{E}_{X,Y,d} \left[\log P_M^d(Y \mid \theta_d, X) \right]$$
(2)

We align $\hat{\theta}_d$ with the token embedding space by introducing new tokens—our concept tokens—into the model's vocabulary. These tokens represent the task concept θ_d , allowing the model to utilize them within its regular vocabulary. Following methods proposed by Lester et al. [3], we add c new concept tokens, denoted as $\hat{\theta}_d$, to represent each task's concept. The embeddings of these new tokens, $E_{\text{new}}(\hat{\theta}_d)$, are fine-tuned specifically for the task while keeping the rest of the language model's parameters frozen. This approach enables the model to focus on learning the nuances of θ_d without altering its general language capabilities. The parameter c, representing the number of concept tokens, is treated as a hyperparameter adjustable based on task requirements.

During training, the c concept tokens associated with $\hat{\theta}_d$ are prepended to the input X (or output Y) to condition the model on the specific task, providing task-specific context that enhances predictive performance.



Fig. 4. Latent Concept Learning Module

This process is illustrated in Figure 4, which provides a flowchart for the latent concept learning method. The flow depicts how, starting from a dataset D, the input X_i is fed into the model along with the updated concept tokens $\hat{\theta}$. The model generates the output Y'_i , and the cross-entropy loss $\log P_M(Y_i \mid \theta, X_i)$ is computed to update θ . This iterative training process enables the model to understand and adapt to the task-specific requirements embedded in θ , leading to more relevant demonstration selections in *DemoCraft*.

B. Task Concept Probability Calculation

In the Task Concept Probability Calculation stage, our objective is to quantify how well each demonstration aligns with the target task. This involves calculating the relevance of each input-output pair (X_i, Y_i) within the context of the task's specific requirements.

Leveraging the previously trained concept tokens θ , we evaluate the suitability of input-output pairs from our dataset \mathcal{D} . For each pair (X_i, Y_i) , we compute the probability $P_M(\theta \mid Y_i, X_i)$, which measures the degree to which the demonstration aligns with the task-specific concept encapsulated by θ . This probability serves as an evaluative metric, where higher values indicate stronger alignment with the task.

Formally, the task concept probability is calculated using Bayes' theorem:

$$P_M(\theta \mid Y_i, X_i) = \frac{P_M(Y_i, X_i \mid \theta) P_M(\theta)}{P_M(Y_i, X_i)},$$
(3)

where:

- $P_M(\theta \mid Y_i, X_i)$ is the posterior probability of the concept tokens given the demonstration pair.
- $P_M(Y_i, X_i \mid \theta)$ is the likelihood of the demonstration pair given the concept tokens.
- $P_M(\theta)$ is the prior probability of the concept tokens.
- $P_M(Y_i, X_i)$ is the marginal probability of the demonstration pair.

In this stage, the large language model M operates in an evaluative capacity; it computes the task concept probabilities based on its learned representations without undergoing further fine-tuning. By assigning task concept probabilities to each demonstration, we gain insights into their relative relevance, which is crucial for selecting the most appropriate demonstrations in subsequent stages.



Fig. 5. Task Concept Probability Calculation Module

This process is illustrated in Figure 5, which outlines how input-output pairs, along with the trained concept tokens θ , are processed through the model to compute the task concept probabilities $P_M(\theta \mid Y_i, X_i)$ for each pair (X_i, Y_i) .

C. Demonstration Selection

In the *Demonstration Selection* stage, our objective is to identify the most relevant demonstrations for a given task prompt. Having computed the task concept probability $P_M(\theta \mid Y_i, X_i)$ for each demonstration pair (X_i, Y_i) in our dataset \mathcal{D} , we proceed to select the top k demonstrations that align most closely with the task-specific concept θ .

We rank all demonstration pairs based on their computed task concept probabilities and select the top k pairs with the highest values of $P_M(\theta \mid Y_i, X_i)$. This selection process ensures that we retain demonstrations that are most



Fig. 6. Demonstration Selection Module

contextually relevant to the task at hand. By focusing on the highest probability values, we choose examples that the model has identified as highly aligned with the desired taskspecific features. This maximizes the likelihood that these demonstrations will enhance the model's understanding and performance when generating responses for the target prompt.

This process is illustrated in Figure 6, which shows how we systematically select the top k demonstrations with the highest alignment scores, ultimately constructing a refined set of examples tailored to optimize the model's responses for the given prompt.

D. Final System Diagram

DemoCraft extends the foundational concepts discussed—namely, latent concept learning and task concept probability calculation—to operate across multiple datasets. This enables the model to learn a comprehensive set of concept tokens, each corresponding to distinct task types denoted by $\theta_1, \theta_2, \ldots, \theta_k$. Once trained, these concept tokens allow the system to retrieve relevant demonstrations from a diverse range of sources.

When a new prompt Q is provided, *DemoCraft* evaluates it by calculating probabilities over both the learned concept tokens and potential demonstration pairs (X_j, Y_j) from the dataset \mathcal{D} . This involves a two-step process:

- 1) For each concept token θ_i , compute the probability $P_M(\theta_i \mid X_j, Y_j)$ for all demonstration pairs $(X_j, Y_j) \in \mathcal{D}$.
- 2) Maximize this probability over both θ_i and (X_j, Y_j) to select the top k demonstrations:

$$\{(X_{i^*}, Y_{i^*})\} = \arg \max_{\theta_i, (X_j, Y_j)} P_M(\theta_i \mid X_j, Y_j), \quad (4)$$

where $\{(X_{i^*}, Y_{i^*})\}$ denotes the set of top k demonstrations that best align with the task-specific requirements of Q. This approach leverages both the learned task-specific knowledge encapsulated in the concept tokens and the diversity of the dataset, ensuring a refined and targeted selection process.

The overall system flowchart is provided in Figure 7, illustrating how the trained concept tokens, task probability calculator, and demonstration selector operate in unison to choose the most relevant examples for each new prompt.

IV. EXPERIMENTS

In this section, we highlight our experimental metrics and the conditions under which we conducted the experiments.



Fig. 7. DemoCraft System Flowchart

A. Evaluation Metrics

We evaluate our model using three primary metrics:

 pass@k: This metric measures the probability that at least one of the top k generated code samples passes all the test cases for a given problem. Suppose for each problem we generate n code samples, out of which c samples are correct (i.e., they pass all the unit tests). The pass@k is calculated as:

$$\operatorname{pass}@k = \mathbb{E}_D\left[1 - \frac{\binom{n-c}{k}}{\binom{n}{k}}\right],\tag{5}$$

where \mathbb{E}_D denotes the expectation over the dataset D, and $\binom{n}{k}$ is the binomial coefficient representing the number of ways to choose k samples out of n.

2) correctness@k: This metric is defined as the average precision of the model over the entire dataset when k outputs are generated per prompt. For each prompt, if the model generates k outputs and c of them are correct, the correctness for that prompt is calculated as:

correctness@
$$k = \mathbb{E}_D\left[\frac{c}{k}\right],$$
 (6)

where \mathbb{E}_D denotes the expectation over the dataset D.

3) similarity@k: This metric measures the average similarity between the working codes generated by the model and the golden solution provided in the dataset. For each prompt, let S be the set of all generated codes that pass all the test cases (i.e., working codes), and let y be the golden solution from the dataset. The similarity@k is defined as:

similarity @
$$k = \mathbb{E}_D\left[\frac{1}{|S|}\sum_{y_i \in S} \sin(y_i, y)\right],$$
 (7)

where $sim(y_i, y)$ is a similarity function between the generated code y_i and the golden solution y, and |S| is the number of working codes for that prompt. The outer expectation \mathbb{E}_D is taken over all prompts in the dataset D. The similarity function used over here is the *edit distance* metric, provided in the standard *NLTK* library.

B. Datasets and Models

We conducted our experiments using the following datasets and model:

- MBPP: The Mostly Basic Python Problems (MBPP) dataset [4] consists of 427 programming problems designed for code generation tasks. Each problem includes a natural language description, the corresponding code solution, and three unit tests. The programming language used is Python.
- 2) HumanEval: The HumanEval dataset [5] comprises 164 programming tasks focused on code completion. Each task provides a function signature and a docstring describing the desired functionality. The solutions are written in C++, and each problem includes approximately seven unit tests, making it a stricter benchmark than MBPP.

Due to resource constraints, we evaluated the performance of our system using the *SantaCoder* model. SantaCoder is a transformer-based language model with 1.1 billion parameters, pretrained on a large corpus of code in multiple programming languages, including Python and C++. It is designed to generate syntactically correct and functionally meaningful code snippets. We conducted our experiments using Google Colab's T4 GPU, which provided sufficient computational resources for our evaluations without compromising performance.

C. Baselines

We compare our system against the following baseline methods:

1) Semantic Selection: In this baseline, we select demonstrations from the dataset purely based on their semantic similarity to the given prompt x. Let the dataset be $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$, where x_i are the prompts and y_i are the corresponding outputs. For each x_i in the dataset, we compute the similarity score $sim(x, x_i)$ between the given prompt x and each dataset prompt x_i . We then select the top k demonstrations with the highest similarity scores:

$$\{(x_{i^*}, y_{i^*})\} = \arg\max_{\substack{(x_i, y_i) \in \mathcal{D} \\ i=1,\dots, n}} \sin(x, x_i), \qquad (8)$$

where $\{(x_{i^*}, y_{i^*})\}$ denotes the set of top k demonstrations selected. The *sim(.)* function used here is the standard edit distance, implemented using the *NLTK* library.

2) **Random Selection**: In this baseline, we randomly select k demonstrations from the dataset \mathcal{D} without considering their relevance to the given prompt x. This method serves as a control to evaluate the impact of demonstration selection strategies on the model's performance.

V. RESULTS

In this section, we present the results of our experiments on both the *MBPP* and *HumanEval* datasets. Table 1 shows the results for the *MBPP* dataset, while Table 2 presents the results for the *HumanEval* dataset.

The results show that demonstrations chosen by *DemoCraft* consistently outperform other selection methods. This superiority arises from *DemoCraft*'s encoding of task-specific





Fig. 10.

TABLE I Evaluation Results on MBPP

Parameter	Semantic	DemoCraft	Random
correctness@5	2%	7.2%	1.5%
correctness@20	0.5%	6.0%	0.3%
correctness@100	0.3%	5.0%	0.2%
similarity@5	0.77%	3.0%	0.5%
similarity@20	0.771%	3.5%	0.4%
similarity@100	2.7%	7.0%	1.8%
pass@1	0.6%	4.0%	0.2%
pass@10	6.07%	11.5%	5.0%
pass@100	20%	27.0%	15.0%

TABLE II Evaluation Results on HumanEval

Parameter	Semantic	DemoCraft	Random
correctness@5	0.1%	1.2%	0.2%
correctness@20	0.04%	1.1%	0.03%
correctness@100	0.008%	1.0%	0.005%
similarity@5	0.91%	3.5%	0.8%
similarity@20	0.92%	4.0%	0.7%
similarity@100	3%	7.5%	2%
pass@1	0.3%	2.0%	0.4%
pass@10	4.56%	8.0%	3%
pass@100	13.2%	18.5%	10%







knowledge through specialized token embeddings tailored to each task.

VI. CONCLUSION

In this paper, we presented *DemoCraft*, a demonstration selection framework that enhances code generation models by leveraging task-specific knowledge through latent concept learning. *DemoCraft* introduces specialized token embeddings tailored to each task, enabling the model to internalize underlying concepts effectively. Our evaluations on the *MBPP* and *HumanEval* datasets, utilizing the metrics *pass@k*, *correctness@k*, and *similarity@k*, demonstrate that *DemoCraft* consistently outperforms baseline methods, including semantic similarity-based and random selection approaches. These results highlight the efficacy of targeted demonstration selection in improving code generation accuracy and functionality. Future work will explore the integration of *DemoCraft* with larger language models and its application to diverse domains, including software engineering and competitive programming.

Acknowledgements

We acknowledge Dr. Amar Prakash Azad and Dr. Brij Kumar Chavda from IBM Research Bangalore for their invaluable support and mentorship, which were instrumental to the success of this project.

REFERENCES

- [1] Yujia Li, David Choi, Junyoung Chung, Nate Kushman et al. Competition-Level Code Generation with AlphaCode. arXiv preprint.
- [2] Daoguang Zan, Bei Chen, Fengji Zhang, Dianjie Lu et al. Large Language Models Meet NL2Code: A Survey. arXiv preprint. 2022
- [3] B. Lester, R. Al-Rfou, and N. Constant. The power of scale for parameter-efficient prompt tuning. In Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pages 3045–3059, 2021
- [4] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma et al. Program Synthesis with Large Language Models. arXiv preprint. 2021
- [5] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto et al. Evaluating Large Language Models Trained on Code. arXiv preprint. 2021.
- [6] Xinyi Wang, Wanrong Zhu, Michael Saxon, Mark Steyvers, William Yang Wang. Large Language Models Are Latent Variable Models: Explaining and Finding Good Demonstrations for In-Context Learning. arXiv preprint. 2024.