TOPOLOGY OF ATTENTION DETECTS HALLUCINA-TIONS IN CODE LLMS

Anonymous authors

Paper under double-blind review

ABSTRACT

While the AI-code assistant tools become widespread, automatic assessment of the correctness of the generated code becomes a significant challenge. Code LLMs are prone to hallucinations, which may lead to code that does not solve a required problem, or even to code with severe security vulnerabilities. In this paper, we propose a new approach to assessment of code correctness. Our solution is based on topological data analysis (TDA) of attention maps of code LLMs. We carry out experiments with two benchmarks – HumanEval, MBPP and 5 code LLMs: StarCoder2-7B, CodeLlama-7B, DeepSeek-Coder-6.7B, Qwen2.5-Coder-7B, Magicoder-S-DS-6.7B. The experimental results show that the proposed method is better than several baselines. Moreover, the trained classifiers are transferable between coding benchmarks.

1 Introduction

Large Language Models (LLMs) are now widespread and have great potential to transform natural language processing and artificial intelligence. As far as code generation is concerned, LLMs that are trained on large amounts of code are capable of generating human-level code for a plethora of simple problems and are expected to revolutionize software engineering. At the same time, code-generating LLMs are prone to hallucinations of various types. For example, syntactic and runtime errors prevent proper program execution, while logical errors lead to incorrect solution of the problem. In some cases, the generated code might contain security issues or robustness issues, such as a memory leak. While many definitions of hallucinations exist, in this paper we assume that code hallucination is a code which is not functionally correct, that is, does not pass functional tests. For a wide adoption of code LLMs, there is a high need for automatic assessment of code quality. As for the current state of technology, significant time is spent on debugging and rewriting automatically generated code (Liang et al., 2024).

We hypothesize that code quality can be inferred before its execution from an internal state of LLM, in particular its attention maps. Previous studies have shown that attention maps of transformers are useful for artificial text detection (Kushnareva et al., 2021), acceptability judgments (Cherniavskii et al., 2022) and speech classification (Tulchinskii et al., 2022).

Attention maps of LLMs are shown to capture semantically meaningful information and might be an illustration of model's "thinking process". The research community actively studies approaches to mitigate hallucinations of LLMs by external knowledge bases (Peng et al., 2023) or reducing them to some degree (Elaraby et al., 2023). It is a highly desirable to evaluate a code quality before its execution and a running of tests since the code might contain security vulnerabilities.

The study of hallucinations in LLMs is intrinsically tied to generalization in NLP models. Both challenges stem from how models learn, represent, and apply knowledge. Improving generalization – through robust training, diverse data, and better uncertainty handling – reduces hallucinations by ensuring models produce contextually appropriate, factually grounded outputs. Conversely, analyzing hallucinations provides insights into generalization failures, guiding the development of more reliable NLP systems. This symbiotic relationship underscores the importance of addressing both issues holistically in AI research.

Out contributions are the following:

- 056
- 060
- 061 062 063
- 064 065 066

068

- 069 071
- 072 073 074 075 076
- 077 079
- 080 081 083 084
- 085 086 087 880
- 089

091

- 092 094 095 096
- 097 098 099
- 100 101 102

103 104 105

106

107

3

BACKGROUND

(Vaswani, 2017).

- We propose a new approach to detection of hallucinations in LLM-generated code based on analyzing a topology of attention maps;
- We carry out computational experiments with CodeLlama, StarCoder2, DeepSeek-Coder and Qwen2.5-Coder and two benchmarks - HumanEval and MBPP, and show that the proposed method outperforms baselines;
- We empirically show that proposed classifier of hallucinations is transferable between code benchmarks.

RELATED WORK

Code generation via LLMs is a topic of active research. The popular projects are: CodeLlama (Roziere et al., 2023), StarCoder2 (Lozhkov et al., 2024), DeepSeek-Coder (Guo et al., 2024), Owen2.5-Coder (Hui et al., 2024), to name a few. Code LLMs differ by the data used for training, by their training and fine-tuning protocols, including RLHF, tokenizers, variants of attention mechanism, etc.

Several works studied attention maps in transformer-based LLMs. Clark et al. (2019) studied BERT's attention patterns: attending to delimiter tokens, specific positional offsets, or broadly attending over the whole sentence, with heads in the same layer often exhibiting similar behaviors. Clark et al. (2019) further showed that certain attention heads correspond well to the linguistic notions of syntax and coreference. Htut et al. (2019) found that for some universal dependency tree relation types, there exist heads that can recover the dependency type significantly better than baselines on parsed English text, suggesting that some self-attention heads act as a proxy for syntactic structure. Michel et al. (2019) showed that for downstream tasks, a large proportion of attention heads can be removed at test time without significantly affecting performance and that some layers can even be reduced to a single head.

The phenomenon of code hallucinations is studied and categorized in several papers. Tian et al. (2024) introduces a categorization of code hallucinations into four main types: mapping, naming, resource, and logic hallucinations, with each category further divided into different subcategories. Tian et al. (2024) proposed the CodeHalu dataset and studied the frequencies of different types of hallucinations in popular code LLMs. Liu et al. (2024) categorized hallucinations as intent conflicting, inconsistency, repetition, knowledge conflicting, dead code. Liu et al. (2024) released a HaluCode benchmark with code hallucinations labeled. Jiang et al. (2024) proposed Collu-Bench, the benchmark with localized code hallucinations. Jiang et al. (2024) found that code LLMs are less confident when hallucinating, since hallucinated tokens have a lower probability and hallucinated generation steps have a higher entropy.

In the broader context of NLP, several works introduced methods for preventing and detecting hallucinations. Peng et al. (2023) proposed to mitigate hallucination with an LLM-AUGMENTER, a system that allows the LLM to generate responses grounded in external knowledge, for example, stored in task-specific databases. Zhang et al. (2024b) proposed Self-Eval, a self-evaluation component, to prompt an LLM to validate the factuality of its own generated responses solely based on its internal knowledge. Feng et al. (2024) proposed two novel approaches for hallucination detection that are based on model collaboration, i.e., LLMs that investigate other LLMs for knowledge gaps, cooperatively or competitively. Zhang et al. (2024a) proposed to improve the truthfulness of LLMs by editing their internal representation during inference in the "truthful" space. Yehuda et al. (2024) introduced InterrogateLLM, a method which prompts the model multiple times to reconstruct the input query using the generated answer. Subsequently, InterrogateLLM quantifies the inconsistency level between the original query and the reconstructed queries.

3.1 Transformer-based LLMs

All of the state-of-the art code LLMs are based on different variants of the transformer architecture

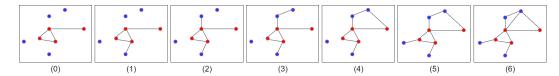
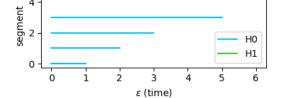


Figure 1: An example of MTD evaluation for a graph having two groups of vertices – red and blue. (0): initially, only edges connecting red vertices are present. (1)-(6): the rest of edges are added sequentially in an ascending order by their weights. While adding edges, connected components merge with each other. These moments are depicted by H_0 bars in Fig. 2. At moment (4) a cycle appears, at moment (6) this cycle disappears. These moments are depicted by the H_1 bar in Fig. 2.

A transformer architecture comprises L layers of multi-head self-attention blocks each of them having H heads. Each attention head takes $X \in \mathbb{R}^{n \times d}$ matrix as an input, and an output is $X^{out} = A(XW^V)$, where

$$A = \operatorname{softmax} \left(\frac{(XW^Q)(XW^K)^T}{\sqrt{d}} \right),$$



and $W^Q, W^K, W^V \in \mathbb{R}^{d \times d}$ are projection matrices, and $A \in [0,1]^{n \times n}$ is an **attention** map. In the self-attention block, the attention

Figure 2: Cross-Barcode for a filtration from Fig. 1.

map shows how each token in the input sequence "interacts" to every other token in the same sequence. A token might attend more to other tokens that are contextually related. We interpret each element $a_{i,j}$ of an attention map as an "interaction force" between tokens i and j.

3.2 Representing an attention map by a weighted graph

While attention map is typically presented as a matrix, we treat it as a weighted graph. For n tokens in a sequence, we consider a fully-connected weighted graph with n vertices, where weights of edges are related to the "interaction force" between tokens (vertices). The natural idea is to leave only the most interacting tokens, that is, attending to each other higher than some threshold. However, the optimal threshold is not known in advance. Moreover, topology of such graph changes discontinuously with the change of a threshold (or weights). Topological Data Analysis (TDA) (Chazal & Michel, 2017) introduces a principled way to access topology of such graphs for all thresholds simultaneously.

3.3 Manifold Topology Divergence

MTD (Manifold Topology Divergence) (Barannikov et al., 2021) is a tool of TDA which can be used to evaluate the "dissimilarity" between two sets of vertices in a weighted graph $\mathcal{G}=(V,E,W)$ or, in other words, to which degree one set of vertices is covered by another set.

Let a set of vertices $V = P \sqcup G$, be split into disjoint sets P, G. We consider a nested sequence of graphs $\mathcal{G}_0 \subset \ldots \subset \mathcal{G}_i \subset \mathcal{G}_{i+1} \subset \ldots \subset \mathcal{G}$ in the following way. \mathcal{G}_0 has all the vertices P, G and all the edges connecting vertices from P. The sequence \mathcal{G}_i is obtained by adding the rest of edges one by one in an ascending order by their weights, see Figure 1. During this process, graphs' topology naturally changes: connected components are merged, cycles appear and disappear, etc. This process is rigorously described by the theory of *persistence barcodes* (Chazal & Michel, 2017). Each topological feature like connected component or cycle has a "birth time" and a "death time", by a corresponding edge weight. The multi-set of these birth-death pairs (intervals) altogether is called a Cross-Barcode $_k$, see Figure 2. Here k is an index of a *persistence homology*, each of them reflects a kind of topological feature: 0 - connected components, 1 - cycles, 2 - voids, etc. MTD $_k$ is an integral characteristic of a Cross-Barcode $_k$ and it is defined as a sum of birth-death intervals' lengths. The higher MTD $_k$ is, the bigger is a "dissimilarity" between sets of tokens. Note, that according to a definition, MTD $_k$ is not symmetric. Also, MTD $_k$, as a kind of persistence barcode, enjoys stability w.r.t. small perturbations of weights (Cohen-Steiner et al., 2005).

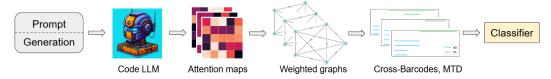


Figure 3: A pipeline of the proposed method for hallucination detection: (1) a prompt concatenated with a generated code is fed into a Code LLM. (2) Attention maps from the Code LLM are obtained. (3) Attention maps are transformed into fully-connected weighted graphs. (4) Cross-Barcodes and MTD features for weighted graphs are calculated. (5) On the top of the generated features a binary classifier of hallucinations is fitted.

Table 1: Characteristics of generated data: Pass@1, number of correct (#Pos.) and incorrect (#Neg.) solutions for each of the selected code LLMs.

Model	HumanEval			MBPP		
	Pass@1	#Pos.	#Neg.	Pass@1	#Pos.	#Neg.
StarCoder2-7B	28.9	1186	2914	42.8	1071	1429
CodeLlama-7B	25.9	1064	3036	35.2	879	1621
DeepSeek-Coder-6.7B	40.3	1653	2447	52.6	1315	1185
Qwen2.5-Coder-7B	47.8	1961	2139	52.1	1302	1198
Magicoder-S-DS-6.7B	65.5	2689	1411	61.3	1533	967

4 METHODS

In the context of code generation, we naturally have two sets of tokens – a prompt and a generation. A common cause of hallucination is when the model's attention drifts away from the prompt¹. Our topology-based features quantify this mismatch, yet the same signal helps to identify other error patterns as shown below (Section 5.6). As was pointed in Section 3.2, attention matrices can be analyzed as weighted graphs. Specifically, for n tokens in a sequence, we consider a fully-connected weighted graph with n vertices, where weights of edges are obtained by a symmetrization of an attention map: $w_{i,j} = 1 - \max(a_{i,j}, a_{j,i})$, for $i \neq j$. Then, Cross-Barcode and MTD for a weighted "attention graph" can be calculated². To predict code hallucinations, we use the following set of features:

- $MTD_0(P, G)/|G|$, $MTD_0(G, P)/|P|$
- $MTD_1(P, G)/|G|, MTD_1(G, P)/|P|$
- $\sum_{i \in P} a_{i,i}/|P|$, $\sum_{i \in G} a_{i,i}/|G|$

Here all the features are normalized by a size of corresponding vertices set for better transferability. Additionally, sums of diagonal values of attention matrices which are not directly present in edge weights are included. These features are calculated for every layer and head of a code LLM. At the top of the proposed topological features, we applied XGBoost (Chen & Guestrin, 2016) for a classification³. The high-level pipeline of the proposed method is shown in Figure 3.

5 EXPERIMENTS

5.1 GENERATION OF DATASETS

To assess the efficacy of the proposed method for hallucination prediction, we carry out a set of computational experiments. In the main experiments, we use the following popular code LLMs:

¹The definition of a hallucination is discussed in Appendix G.

²Appendix D contains examples of Cross-Barcodes and corresponding attention maps.

³Appendix E contains an ablation study of a classifier model.

StarCoder2-7B (Lozhkov et al., 2024), CodeLlama-7B (Roziere et al., 2023), DeepSeek-Coder-6.7B (Guo et al., 2024), Qwen2.5-Coder-7B (Hui et al., 2024), Magicoder-S-DS-6.7B (Wei et al., 2024). We adapted two public benchmarks for evaluation of code generation: HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) 4. In order to account for various possible code generations, for each of the coding problems several solutions were generated by each of the selected code LLMs: we obtained 25 generations per task for HumanEval and 5 generations per task for MBPP unless otherwise specified. To address the quality of the proposed approach in different LLM prompting regimes, we used 0-shot prompt for the HumanEval dataset and 1-shot prompt for the MBPP dataset. To enable diversity of generated solutions, a sampling with non-zero temperature of was done. Thus, we obtain 4100 samples for HumanEval and 2500 samples for MBPP for each code LLM (see Appendix A for further details). Table 1 presents a summary of generated code solutions. The correctness of code is evaluated via functional tests provided together with the coding benchmarks. Functional tests check that the function called with the certain arguments has the corresponding output (examples are shown in Figures 6, 7 in Appendix A). Incorrect code is considered a "hallucination"; prediction of code's correctness is a binary classification problem. The pass@1 metric is slightly lower that reported in original papers, mostly because we have used sampling with non-zero temperature instead of greedy search. Before moving further, note that there is a strong negative dependency between prompt and generation lengths and code quality, see Figure 8, 9. The longer the prompt (i.e. task description) and generation (i.e. task solution) are, the lower is the probability of code's correctness. This dependency is more pronounced for HumanEval than MBPP, because MBPP employed more complicated 1-shot prompts. These attributes are natural baselines for hallucination's prediction.

5.2 Analyzing method's classification quality

Using the generated data, we estimated the classification quality of the proposed approach. We applied 5-fold stratified group cross-validation where different solutions of the same coding problem belonged to the same group. In this way, training and testing were performed always at non-overlapping coding problems (prompts). The reported results are the mean and standard deviation estimated over the 5 folds.

As baselines for comparison, we used XGBoost classifier trained on simple features: tokenized prompt length, tokenized generation length, and mean log-probability of generated tokens Chen et al. (2021). Also, we trained a linear classification head on top of a frozen CodeT5-base Wang et al. (2021) encoder. Furthermore, we have adapted the Self-Eval Zhang et al. (2024b) and Interrogate-LLM Yehuda et al. (2024) to detect hallucinations in LLM-generated code and utilized them as baselines as well. Finally, we utilize a combination of all features, i.e. tokenized prompt length, tokenized generation length, mean log-probability and the proposed attention features, to train a classifier (we refer to it 'All Feat.' for brevity). Training details are provided in Appendix B. Table 2 presents the main results. Table 12 in Appendix presents the additional results for the 32-34B models.

In the majority of cases, the proposed classifier based on features of attention maps performed significantly better than the baselines and demonstrated stable results for all models and datasets as measured by ROC-AUC score. Further analysis revealed that some features made a high contribution to the classification quality, see Figure 4. Usage of additional features can both decrease and increase the overall performance. Thus, we suppose the proposed attention features are strong enough to capture the most important information for code hallucination detection.

Moreover, the classifier detects diverse failure modes: SyntaxError, ZeroDivisionError, NameError, etc., across six categories (Section 5.6, Table 5), confirming that the approach generalizes beyond the prompt-insufficient attention failure mode.

Furthermore, we evaluate the proposed approach on Java (high resource), Go (medium resource), Rust (low resource) and Lua (niche) programming languages of the HumanEval subdivision of the MultiPL-E dataset Cassano et al. (2023). Tables 14, 16, 17 demonstrate that the performance of the proposed approach is comparable to the main experiments. Also, we increase the linguistic complexity of the prompt via two-shot prompts on MBPP dataset, see Table 13 for the results. Finally, for additional comparison with zero-shot classification via larger LLMs, see Appendix J.

⁴Licenses of pretrained models and benchmarks permit use for research purposes.

Table 2: Quality of code hallucination detection for HumanEval and MBPP datasets. **Bold** and <u>underline</u> denote the first and second best results.

Model	Huma	anEval	MBPP	
	ROC-AUC	F1-Score	ROC-AUC	F1-Score
	Sta	rCoder2-7B		
Prompt Len.	54.5 ± 6.6	24.6 ± 11.2	51.2 ± 2.3	40.0 ± 3.8
Gen. Len.	57.7 ± 5.6	13.5 ± 4.8	57.7 ± 0.9	45.4 ± 3.5
Mean Log. Prob.	70.9 ± 1.3	32.4 ± 4.8	62.0 ± 2.0	47.5 ± 3.4
CodeT5-base ft.	70.1 ± 7.1	33.3 ± 10.1	58.5 ± 3.5	43.3 ± 9.0
Self-Eval	50.0 ± 2.9	0.0 ± 0.0	58.9 ± 2.6	0.0 ± 0.0
Interrogate-LLM	63.9 ± 2.1	48.1 ± 4.2	58.6 ± 2.6	60.9 ± 1.3
Attn. Feat. (ours)	82.9 ± 2.7	54.2 ± 6.9	81.9 ± 2.4	68.4 ± 5.3
All Feat.	$\textbf{83.7} \pm \textbf{4.1}$	$\textbf{57.5} \pm \textbf{7.9}$	81.8 ± 2.1	68.7 ± 5.1
		deLlama-7B		
Prompt Len.	61.6 ± 4.4	25.7 ± 15.0	59.1 ± 4.2	35.4 ± 2.9
Gen. Len.	60.1 ± 5.3	10.6 ± 7.0	60.8 ± 2.5	24.2 ± 5.3
Mean Log. Prob.	64.1 ± 2.0	25.4 ± 6.2	61.0 ± 3.7	27.2 ± 1.5
CodeT5-base ft.	74.5 ± 6.3	43.6 ± 13.2	61.7 ± 3.0	19.1 ± 7.1
Self-Eval	49.7 ± 1.7	39.6 ± 4.4	50.0 ± 0.0	0.0 ± 0.0
Interrogate-LLM	67.9 ± 2.5	44.5 ± 4.2	57.6 ± 2.4	52.5 ± 1.6
Attn. Feat. (ours)	85.6 ± 3.9	56.4 ± 7.2	83.4 ± 3.3	64.0 ± 4.4
All Feat.	85.7 ± 3.8	57.9 ± 9.7	83.8 ± 2.0	$\overline{65.2 \pm 4.3}$
		eek-Coder-6.71		F 0.4.1.0.0
Prompt Len.	56.2 ± 4.6	44.4 ± 4.3	52.5 ± 2.5	56.4 ± 3.6
Gen. Len.	57.9 ± 2.4	34.4 ± 4.9	54.6 ± 1.9	59.4 ± 1.3
Mean Log. Prob.	69.8 ± 2.5	51.1 ± 3.4	61.0 ± 1.9	62.3 ± 1.6
CodeT5-base ft.	69.1 ± 4.2	52.6 ± 6.5	55.7 ± 3.0	64.8 ± 2.7
Self-Eval	56.7 ± 2.5	0.0 ± 0.0	51.1 ± 0.6	69.2 ± 1.6
Interrogate-LLM	71.8 ± 2.4	62.9 ± 2.7	60.1 ± 5.1	69.0 ± 1.8
Attn. Feat. (ours)	85.6 ± 2.8	$\frac{68.9 \pm 5.5}{60.3 \pm 5.3}$	$\frac{82.6 \pm 1.9}{22.1 \pm 1.2}$	$\frac{76.5 \pm 2.7}{77.0 \pm 1.3}$
All Feat.	85.4 ± 2.7	$\frac{69.3 \pm 5.3}{\text{n}2.5\text{-Coder-7B}}$	83.1 ± 1.2	77.0 ± 1.3
Drompt I on	0.00000000000000000000000000000000000	$\frac{12.3 - \text{Coder-} 7B}{51.0 \pm 5.7}$	51.8 ± 3.6	56.2 ± 4.4
Prompt Len. Gen. Len.	54.3 ± 6.7 57.6 ± 3.6	48.9 ± 5.1	51.6 ± 3.0 55.6 ± 2.1	50.2 ± 4.4 59.7 ± 4.8
Mean Log. Prob.	63.1 ± 2.4	46.9 ± 5.1 55.6 ± 5.5	61.5 ± 1.3	60.4 ± 1.8
CodeT5-base ft.	65.1 ± 2.4 65.9 ± 3.7	58.0 ± 3.5 58.2 ± 4.5	56.0 ± 1.3	65.2 ± 2.0
Self-Eval	73.8 ± 1.6	75.6 ± 3.8	66.6 ± 2.3	76.0 ± 1.3
Interrogate-LLM	60.1 ± 3.7	65.7 ± 5.1	55.0 ± 1.8	$\frac{70.0 \pm 1.5}{68.3 \pm 2.1}$
Attn. Feat. (ours)	81.7 ± 2.8	70.2 ± 4.2	82.2 ± 2.2	75.4 ± 1.7
All Feat.	$rac{81.7 \pm 2.8}{81.8 \pm 2.3}$	$\frac{10.2 \pm 4.2}{69.6 \pm 3.0}$	$\frac{62.2\pm2.2}{82.3\pm2.0}$	76.4 ± 1.7 76.7 ± 1.7
7 Hi i Cut.		oder-S-DS-6.71		10.1 ± 1.1
Prompt Len.	57.3 ± 5.4	60007000000000000000000000000000000000	54.0 ± 2.6	69.2 ± 3.0
Gen. Len.	52.5 ± 2.1	76.3 ± 2.6	52.3 ± 2.1	71.4 ± 2.4
Mean Log. Prob.	71.0 ± 5.3	78.4 ± 2.9	60.5 ± 3.7	71.4 ± 2.1
CodeT5-base ft.	64.7 ± 2.7	77.5 ± 2.7	61.0 ± 3.7	74.8 ± 1.4
Self-Eval	44.3 ± 3.2	79.1 ± 3.1	49.2 ± 1.4	75.5 ± 2.4
Interrogate-LLM	65.3 ± 2.1	79.6 ± 3.1	59.0 ± 7.1	76.1 ± 2.1
Attn. Feat. (ours)	82.3 ± 4.9	80.7 ± 3.6	81.5 ± 1.6	80.6 ± 2.0
All Feat.	81.8 ± 3.3	80.7 ± 3.0	81.2 ± 2.0	80.6 ± 2.2

5.3 Analyzing method's ranking quality

Next, we assess the usefulness of the proposed code hallucination classifier for ranking of code generations. For each problem, all generations were ranked via predicted probability of correctness and one with the highest probability was selected. A baseline was random picking of a code generation. The usage of a classifier is always significantly better by a pass@1 score, see Table 3.

Table 3: pass@1 scores for variants of ranking of code generations.

Model	Huma	nEval	MBPP	
	Random	Clf. Prob.	Random	Clf. Prob.
StarCoder2-7B	28.6 ± 5.5	43.3 ± 9.0	43.0 ± 3.6	$\textbf{49.6} \pm \textbf{4.6}$
CodeLlama-7B	26.0 ± 5.1	39.7 ± 7.2	35.2 ± 3.3	43.6 ± 3.4
DeepSeek-Coder-6.7B	39.1 ± 4.9	56.7 ± 7.4	53.0 ± 2.5	61.4 ± 2.3
Qwen2.5-Coder-7B	51.8 ± 8.0	64.0 ± 7.3	52.6 ± 3.6	62.0 ± 2.4
Magicoder-S-DS-6.7B	72.5 ± 10.0	$\textbf{74.3} \pm \textbf{6.1}$	61.4 ± 3.4	64.8 ± 2.1

Table 4: Transferability of code hallucination detectors. Each classifier was trained on HumanEval (MBPP) dataset and tested on MBPP (HumanEval) dataset.

Model	HumanEval	\rightarrow MBPP	$MBPP \to H$	umanEval
	ROC-AUC	F1-Score	ROC-AUC	F1-Score
St	arCoder2-7B			
Prompt Len.	48.6	0.0	52.1	45.0
Gen. Len.	56.0	14.6	52.4	$\overline{38.3}$
Mean Log. Prob.	63.7	36.2	71.8	45.4
CodeT5-base ft.	53.7	0.0	59.1	0.0
Attn. Features	67.5	0.14	67.2	25.5
Attn. Feat. (ours) + Mean. Log. Prob	$\overline{71.0}$	24.0	$\overline{61.7}$	15.7
	odeLlama-7B			
Prompt Len.	51.7	0.0	53.4	42.9
Gen. Len.	61.5	4.2	50.0	41.0
Mean Log. Prob.	57.7	15.2	65.0	34.9
CodeT5-base ft.	54.9	0.0	62.4	0.0
Attn. Features	<u>69.5</u>	0.2	80.3	34.1
Attn. Feat. (ours) + Mean. Log. Prob.	72.3	<u>15.1</u>	<u>79.0</u>	34.9
Deep	Seek-Coder-6.	.7B		
Prompt Len.	48.0	15.6	52.2	58.2
Gen. Len.	55.3	41.4	54.0	51.3
Mean Log. Prob.	62.5	56.3	69.1	58.3
CodeT5-base ft.	53.4	0.0	55.9	57.4
Attn. Features	67.7	70.4	72.4	20.4
Attn. Feat. (ours) + Mean. Log. Prob.	68.0	71.0	72.5	22.6
Qwe	en2.5-Coder-7			
Prompt Len.	49.9	34.1	51.1	64.1
Gen. Len.	51.6	46.1	54.5	54.3
Mean Log. Prob.	60.3	60.4	64.7	60.8
CodeT5-base ft.	49.1	52.3	51.6	65.6
Attn. Features	70.6	63.3	64.2	54.3
Attn. Feat. (ours) + Mean. Log. Prob.	<u>70.0</u>	55.7	65.4	56.1
	coder-S-DS-6.			
Prompt Len.	48.1	56.3	54.5	79.5
Gen. Len.	54.8	75.2	56.1	74.6
Mean Log. Prob.	63.7	74.9	69.8	76.5
CodeT5-base ft.	49.3	76.0	45.9	79.2
Attn. Features	73.5	78.4	56.9	37.6
Attn. Feat. (ours) + Mean. Log. Prob.	71.8	79.1	63.7	45.8

5.4 METHOD'S TRANSFERABILITY BETWEEN BENCHMARKS

We study further the transferability of the classifiers, based on topological features. In this setting, hallucination classifiers for a fixed code LLM were trained on data for one benchmark (HumanEval, MBPP) and evaluated on another, then repeated vise versa. Table 4 shows results: the proposed classifiers are transferable, albeit the performance is lower when training and testing is done on the

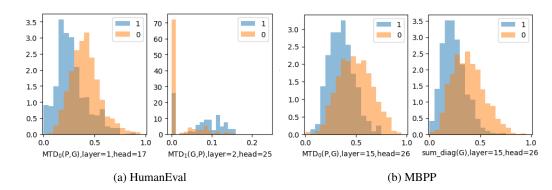


Figure 4: Distribution of classes (0-code is not correct, hallucination, 1-code is correct) vs. features from attention maps. Some of the most discriminative features are presented. Features are normalized with MinMaxScaler. Features come from CodeLlama-7B.

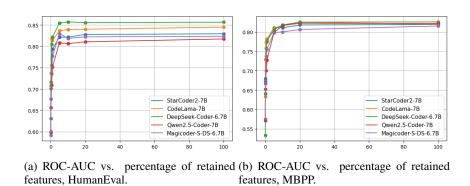


Figure 5: Analysis of feature importance of the proposed method, ROC-AUC.

same benchmark. The proposed attention features achieve better transferability in 70% of cases as measured by ROC-AUC for both HE \rightarrow MBPP and MBPP \rightarrow HE transfer. Next, it is possible to further improve the transferability by combining the proposed attention features with mean log. probability of a generation. The combined classifier outperforms the Mean. Log. Prob. baseline in 80% of cases. This indicates that these two approaches cover different aspects, and the proposed attention features are indeed useful.

5.5 ANALYZING FEATURE IMPORTANCE

In its base setup, the proposed approach requires computation of attention features from attention maps for all layers and heads. However, we observed that the trained XGBoost classifier experienced a natural sparsity with only about 25% of meaningful features as measured by classifiers' feature importance. To explore further the feature importance and feature selection, we followed the two-stage pipeline. First, for a given sparsity level, we selected the most important features as measured by feature importance of the classifier trained on all attention features simultaneously. Second, we trained a new XGBoost classifier on the selected set of attention features. As indicated by Figures 5, 10, the proposed feature selection procedure could retain only 5% of all attention features without significant loss of classification quality highlighting that only a limited number of all attention heads is relevant for the hallucination detection.

We carry out additional experiments benchmarks different programming languages (Python, Go, Rust, Java) and find that topological features of some heads exhibit a high predictive performance consistently across all programming languages, see Appendix K.

Table 5: Performance of multi-classification of error types.

Model	Accuracy	F1-Score
StarCoder2-7B	0.7 ± 0.02	0.68 ± 0.02
CodeLlama-7B	0.66 ± 0.02	0.62 ± 0.02
DeepSeek-Coder-6.7B	0.7 ± 0.03	0.68 ± 0.04
Qwen2.5-Coder-7B	0.64 ± 0.02	0.6 ± 0.02
Magicoder-S-DS-6.7B	0.73 ± 0.02	0.71 ± 0.02

5.6 A FINE-GRAINED CLASSIFICATION OF ERROR TYPES

We carried out additional experiments to study the ability of the proposed approach to detect specific types of hallucinations. A Python exception can be considered as a hallucination type. Here are the common exceptions from HumanEval and MBPP benchmarks: AssertionError, AttributeError, IndentationError, IndexError, ModuleNotFoundError, NameError, RecursionError, SyntaxError, TypeError, UnboundLocalError, ValueError, ZeroDivisionError, timed out

Therefore, we did multi-classification instead of binary classification in XGBoost. Table 5 shows results. Here is a breakdown of detection accuracy of particular types of errors for CodeLlama-7B: AssertionError: 82.0%, IndexError: 97.8%, NameError: 75.7%, RecursionError: 100%, Syntax-Error: 93.3%, TypeError: 80.6%, ValueError: 87.5%, ModuleNotFoundError, ZeroDivisionError, UnboundLocalError, IndentationError, AttributeError, timed out: 80%. Some error types were grouped together because of a very low frequency.

5.7 ABLATION STUDY

The proposed approach is based on the two types of attention features: the diagonal elements of attention maps corresponding to a prompt and a generation, and topological features (MTD) computed for the corresponding "attention graph" (see Section 4 for details). In this Section, we provide an ablation study to estimate the contribution of each type of attention features. For this purpose, we trained the XGBoost classifier using 1) only MTD features 2) only diagonal attention values 3) both types of features (our initial setup). As demonstrated in Table 9 in Appendix, the DeepSeek-Coder-6.7B and Qwen2.5-Coder-7B achieved the best performance when both types of attention features were used for both HumanEval and MBPP datasets. In contrast, the best performance of StarCoder2-7B and Magicoder-S-DS-6.7B was achieved with different sets of attention features dependent on dataset and metric choices. In order to account for various information available via attention maps, we propose to use both types of features as the most universal choice. Nevertheless, we note that for some code LLM one certain type of attention features may result in better performance than combination of both types⁵.

6 Conclusions

In this paper, we have proposed a new hallucination detection approach for code-generating LLMs. Our approach is based on the introspection of a LLM: we get attention maps for a prompt and generation and study their topology after transforming to weighed graphs. The proposed topological features of these graphs are empirically shown to be relevant to detection of code hallucinations. A classifier built on top of these features outperformed several baselines. These classifiers are transferable across coding benchmarks. The natural extension of our research is detection of specific places of code with bugs, we leave it for a further research. We believe that our work may lead to a wider application of code generating LLMs by making them more reliable. In a wider context, our work contributes to study of interpretation and generalization in NLP models since hallucinations and generalization ability are intrinsically tied.

⁵Additionl ablation study in Appendix I shows that for small models (QwenCoder-1.5b, QwenCoder-3b), and transfer learning setting, the differences in contribution of diagonal and MTD features is more pronounced.

7 REPRODUCIBILITY STATEMENT

The source code to reproduce the presented results is provided in the supplementary material. Hyperparameters are either disclosed in the main text and Appendices A, B, or were equal to default values in the code.

REFERENCES

- Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, et al. Program synthesis with large language models. arXiv preprint arXiv:2108.07732, 2021.
- Serguei Barannikov, Ilya Trofimov, Grigorii Sotnikov, Ekaterina Trimbach, Alexander Korotin, Alexander Filippov, and Evgeny Burnaev. Manifold topology divergence: a framework for comparing data manifolds. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 7294–7305, 2021. URL https://proceedings.neurips.cc/paper/2021/hash/3bc31a430954d8326605fc690ed22f4d-Abstract.html.
- Federico Cassano, John Gouwar, Daniel Nguyen, Sydney Nguyen, Luna Phipps-Costin, Donald Pinckney, Ming-Ho Yee, Yangtian Zi, Carolyn Jane Anderson, Molly Q Feldman, Arjun Guha, Michael Greenberg, and Abhinav Jangda. Multipl-e: A scalable and polyglot approach to benchmarking neural code generation. *IEEE Transactions on Software Engineering*, 49(7):3675–3691, 2023. doi: 10.1109/TSE.2023.3267446.
- Frédéric Chazal and Bertrand Michel. An introduction to topological data analysis: fundamental and practical aspects for data scientists. *CoRR*, abs/1710.04019, 2017.
- Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde De Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, et al. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*, 2021.
- Tianqi Chen and Carlos Guestrin. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pp. 785–794, 2016.
- Daniil Cherniavskii, Eduard Tulchinskii, Vladislav Mikhailov, Irina Proskurina, Laida Kushnareva, Ekaterina Artemova, Serguei Barannikov, Irina Piontkovskaya, Dmitri Piontkovski, and Evgeny Burnaev. Acceptability judgements via examining the topology of attention maps. *arXiv preprint arXiv:2205.09630*, 2022.
- Kevin Clark, Urvashi Khandelwal, Omer Levy, and Christopher D Manning. What does BERT look at? an analysis of BERT's attention. *Proceedings of the 2019 ACL Workshop BlackboxNLP*, 2019.
- David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. Stability of persistence diagrams. In *Proceedings of the twenty-first annual symposium on Computational geometry*, pp. 263–271, 2005.
- Mohamed Elaraby, Mengyin Lu, Jacob Dunn, Xueying Zhang, Yu Wang, Shizhu Liu, Pingchuan Tian, Yuping Wang, and Yuxuan Wang. Halo: Estimation and reduction of hallucinations in open-source weak large language models. *arXiv preprint arXiv:2308.11764*, 2023.
- Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan Ding, Vidhisha Balachandran, and Yulia Tsvetkov. Don't hallucinate, abstain: Identifying llm knowledge gaps via multi-llm collaboration. *arXiv* preprint arXiv:2402.00367, 2024.
- Daya Guo, Qihao Zhu, Dejian Yang, Zhenda Xie, Kai Dong, Wentao Zhang, Guanting Chen, Xiao Bi, Yu Wu, YK Li, et al. Deepseek-coder: When the large language model meets programming—the rise of code intelligence. *arXiv* preprint arXiv:2401.14196, 2024.

- Phu Mon Htut, Jason Phang, Shikha Bordia, and Samuel R Bowman. Do attention heads in bert track syntactic dependencies? *arXiv preprint arXiv:1911.12246*, 2019.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.
 - Nan Jiang, Qi Li, Lin Tan, and Tianyi Zhang. Collu-bench: A benchmark for predicting language model hallucinations in code. *arXiv preprint arXiv:2410.09997*, 2024.
 - Laida Kushnareva, Daniil Cherniavskii, Vladislav Mikhailov, Ekaterina Artemova, Serguei Barannikov, Alexander Bernstein, Irina Piontkovskaya, Dmitri Piontkovski, and Evgeny Burnaev. Artificial text detection via examining the topology of attention maps. *arXiv preprint* arXiv:2109.04825, 2021.
 - Jenny T Liang, Chenyang Yang, and Brad A Myers. A large-scale survey on the usability of ai programming assistants: Successes and challenges. In *Proceedings of the 46th IEEE/ACM International Conference on Software Engineering*, pp. 1–13, 2024.
 - Fang Liu, Yang Liu, Lin Shi, Houkun Huang, Ruifeng Wang, Zhen Yang, Li Zhang, Zhongqi Li, and Yuchi Ma. Exploring and evaluating hallucinations in llm-powered code generation. *arXiv* preprint arXiv:2404.00971, 2024.
 - Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173*, 2024.
 - Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32, 2019.
 - Baolin Peng, Michel Galley, Pengcheng He, Hao Cheng, Yujia Xie, Yu Hu, Qiuyuan Huang, Lars Liden, Zhou Yu, Weizhu Chen, et al. Check your facts and try again: Improving large language models with external knowledge and automated feedback. *arXiv preprint arXiv:2302.12813*, 2023.
 - Baptiste Roziere, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Romain Sauvestre, Tal Remez, et al. Code llama: Open foundation models for code. *arXiv preprint arXiv:2308.12950*, 2023.
 - Yuchen Tian, Weixiang Yan, Qian Yang, Qian Chen, Wen Wang, Ziyang Luo, and Lei Ma. Codehalu: Code hallucinations in llms driven by execution-based verification. *arXiv* preprint *arXiv*:2405.00253, 2024.
 - Eduard Tulchinskii, Kristian Kuznetsov, Laida Kushnareva, Daniil Cherniavskii, Serguei Barannikov, Irina Piontkovskaya, Sergey Nikolenko, and Evgeny Burnaev. Topological data analysis for speech processing. *arXiv preprint arXiv:2211.17223*, 2022.
 - A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017.
 - Yue Wang, Weishi Wang, Shafiq Joty, and Steven C. H. Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation, 2021. URL https://arxiv.org/abs/2109.00859.
 - Yuxiang Wei, Zhe Wang, Jiawei Liu, Yifeng Ding, and Lingming Zhang. Magicoder: Empowering code generation with OSS-instruct. In *Proceedings of the 41st International Conference on Machine Learning*, volume 235 of *Proceedings of Machine Learning Research*, pp. 52632–52657. PMLR, 21–27 Jul 2024. URL https://proceedings.mlr.press/v235/wei24h.html.
 - Yakir Yehuda, Itzik Malkiel, Oren Barkan, Jonathan Weill, Royi Ronen, and Noam Koenigstein. Interrogatellm: Zero-resource hallucination detection in llm-generated answers. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 9333–9347, 2024.

Shaolei Zhang, Tian Yu, and Yang Feng. Truthx: Alleviating hallucinations by editing large language models in truthful space. *arXiv preprint arXiv:2402.17811*, 2024a.

Xiaoying Zhang, Baolin Peng, Ye Tian, Jingyan Zhou, Lifeng Jin, Linfeng Song, Haitao Mi, and Helen Meng. Self-alignment for factuality: Mitigating hallucinations in llms via self-evaluation. *arXiv preprint arXiv:2402.09267*, 2024b.

Figure 6: Example of prompt (problem description) and model generation for the HumanEval dataset.

```
You are an expert Python programmer, and here is your task: Write a function to find the
similar elements from the given two tuple lists. Your code should pass these tests:
assert similar_elements((3, 4, 5, 6), (5, 7, 4, 10)) == (4, 5)
 assert similar_elements((1, 2, 3, 4), (5, 4, 3, 7)) == (3, 4)
 assert similar_elements((11, 12, 14, 13),(17, 15, 14, 13)) == (13, 14)
 def similar_elements(test_tup1, test_tup2):
  res = tuple(set(test_tup1) & set(test_tup2))
  return (res)
                                                                           one-shot
[DONE]
 You are an expert Python programmer, and here is your task: Write a python function to
 remove first and last occurrence of a given character from the string. Your code should
 pass these tests:
 assert remove_Occ("hello","l") == "heo"
 assert remove_0cc("abcda","a") == "bcd"
 assert remove_Occ("PHP","P") == "H"
                                                                problem description
 [BEGIN]
 def remove_0cc(s,c):
   return s.replace(c,'',s.count(c)-1)
                                                                          generation
 [DONE]
             _____
```

Figure 7: Example of prompt (one-shot example and problem description) and model generation for the MBPP dataset.

A DETAILS ON GENERATION PROCEDURE

We generated solutions for the coding problems with a temperature of 0.8. For the HumanEval dataset, the maximum length of model output (i.e. input prompt + generation) was limited to 512 tokens. For the MBPP dataset, the maximum number of new tokens to generate was set to 256. Figures 6, 7 provide examples of prompts and generations for HumanEval and MBPP datasets. We followed the guidelines⁶ to post process the model output and extract the valid problem solution. To compute attention features according to the proposed method in Section 4, we used the attention submatrix corresponding to input prompt and valid problem solution. For computational experiments we used NVIDIA TITAN RTX.

⁶https://github.com/bigcode-project/bigcode-evaluation-harness

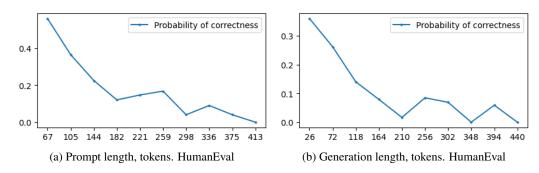


Figure 8: The individual conditional expectations for prompt and generation lengths, CodeLlama-7B.

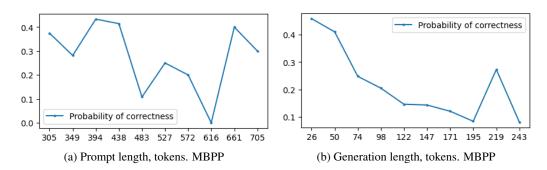


Figure 9: The individual conditional expectations for prompt and generation lengths, CodeLlama-7B.

B DETAILS ON TRAINING PROCEDURE

For the code hallucination detectors, based on the XGBoost classifier training, we utilized the XGBClassifier with an approximation tree method "hist" from the XGBoost library⁷. For the code hallucination detector based on the embeddings from CodeT5-base, we used the pretrained frozen CodeT5-base encoder with a trainable classification head consisting of 2 linear layers with hidden dimensionality 768. The classification head was trained for 100 epochs with batch size 32 and learning rate 3e-5.

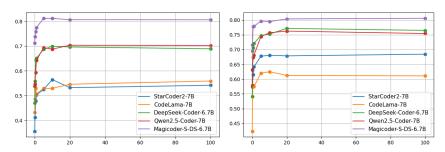
Self-Eval (Zhang et al., 2024b) is a way to evaluate responses of an LLM using its internal knowledge. Self-Eval extracts a list of atomic claims from the responses and then prompts an LLM itself to validate the factuality of claims. Self-Eval is not directly applicable to Code LLMs since there are no "facts" in code. However, we applied the core idea of Self-Eval by prompting Code LLMs to evaluate the functional correctness of a generated code. Also, we have adapted Interrogate-LLM (Yehuda et al., 2024) to hallucination detection in LLM-generated code. As an embedding model we used CodeT5+ 110M Embedding model, K=5 and a fixed temperature.

C METRICS USED FOR EVALUATION

In order to account for possible class imbalance, we use ROC-AUC and F1-Score to evaluate the code hallucination detectors. We briefly introduce them below.

ROC curve demonstrates the quality of a binary classifier for all possible classification thresholds. X-axis corresponds to False Positive Rate (FPR) and Y-axis corresponds to True Positive Rate (TPR) that can be defined as follows: $FPR = \frac{FP}{FP+TN}$, $TPR = \frac{TP}{FP+FN}$, where TP – true positive samples, FP – false positive samples, TN – true negative samples, FN – false negative samples.

⁷https://xgboost.readthedocs.io/en/latest/index.html



(a) F1 vs. percentage of retained features, (b) F1 vs. percentage of retained features, HumanEval. MBPP.

Figure 10: Analysis of feature importance of the proposed method, F1-score

ROC-AUC is defined as the area under the ROC curve. ROC-AUC of a random model is equal to 0.5, ROC-AUC of a perfect model is 1.

F1-score is a harmonic mean of Precision and Recall:

$$F_1 = \frac{2}{\frac{1}{Precision} + \frac{1}{Recall}}.$$

To study the ranking ability of the hallucinations detector, we used pass@1 metric. pass@1 is a proportion of coding problems from a benchmark for which a Code LLM generated the correct solution passing all the tests, with a restriction that only one solution (among 25 generations for HumanEval and 5 for MBPP) is executed.

D EXAMPLES OF CROSS-BARCODES

Figure 11 shows examples of $Cross-Barcode_0$ for a fixed attention head. $Cross-Barcode_1(P, G)$ is empty for these attention maps. Correct generations (a), (b) tend to have more and more H_0 bars than not-correct ones (c), (d).

Figure 12 shows examples of Cross-Barcode₀, Cross-Barcode₁ for a fixed attention head. Correct generations (a), (b) tend to have more and longer H₁ bars than not-correct ones (c), (d).

Attention maps for the correspondind head are shown in Figures 13, 14.

E EXPERIMENTS ON THE CLASSIFIER MODEL SELECTION

We carried out additional experiments with the feed-forward network (MLP), logistic regression, and support vector classifier (SVC) instead of XGBoost as a hallucinations classifier (the rightmost block in Figure 3). We used MLP with two hidden layers of size 256 and ReLU activations. This configuration was selected after a moderate optimization of an architecture. For logistic regression and SVC, we tuned the value of regularization strength. Table 6 presents results. MLP tends to have lower ROC-AUC but sometimes it has higher F1-score. While F1-score is a threshold-dependent metric, ROC-AUC integrates over all thresholds. In an unbalanced setting, ROC-AUC is often more stable thus a classifier may have a higher ROC-AUC even when its F1-score is lower. XGBoost offers (a) strong average performance across all the code LLMs, (b) negligible training cost (\approx 30s per fold), and (c) no hyper-parameter tuning in our setting. XGBoost guarantees a low computational overhead while providing a single, robust baseline for subsequent work.

F A NOTE ON APPLICABILITY OF GNNs TO ATTENTION MAPS

To train a GNN-based approach on the graphs with edge weights obtained from attention matrices, these attention matrices need to be stored. We can estimate the approximate memory footprint to store attention matrices of size (seq_len_k)² for a model with n_layers and n_heads for a dataset of

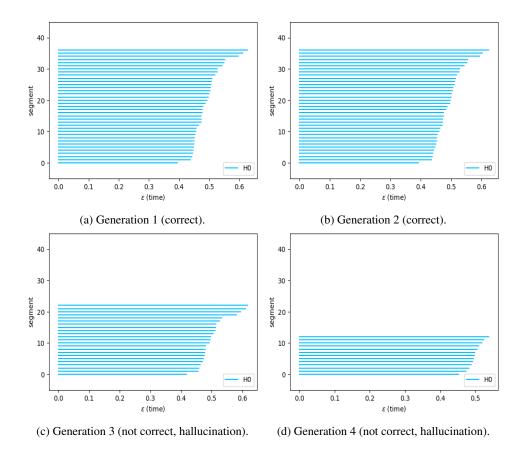


Figure 11: Examples of $Cross-Barcode_0(P, G)$, CodeLLama-7B, HumanEval dataset, problem 14, layer 4, head 18. $Cross-Barcode_1(P, G)$ is empty for these attention maps.

size N using the formula: $n_layers \times n_heads \times s \times \sum_{i=1}^k (seq_len_k)^2$ where s is the size of the float type. We assume s=4 bytes. If one uses only attention matrices from the last layer of the model, we obtain the memory footprint approximately 20.7 - 31.1 Gb for the Human Eval dataset and 36.6 - 53.7 Gb for MBPP (depending on the model). However, to store the attention matrices for all layers and all heads, the memory footprint is about 578.2 - 996.4 Gb for Human Eval and 1023.5 - 1718.7 Gb for MBPP. We highlight that even for datasets of moderate size (i.e. 4100 generations for HumanEval and 2500 generations for MBPP), the memory footprint becomes prohibitively high. Thus, it is not always feasible to store such features. In contrast, in our approach, we do not need to store the attention matrices since we compute all the features immediately during generation. Hence, the size of our training dataset is negligible. Moreover, our approach demonstrates high performance without hyperparameter tuning. Therefore, we suppose the proposed approach has better scalability and is more practical.

G ON DEFINITION OF A CODE HALLUCINATION

In our approach, the topological features obtained from attention maps account for dissimilar structures in the prompt and generation subsets. Our intuition is that a correct solution should correspond to the prompt's structure as their high-level semantic meanings correspond. While other reasons behind the hallucinations are possible, our approach estimates the code correctness based only on the model's internal information flow not requiring additional resources. Nevertheless, the proposed approach can be further integrated with other tools of code hallucination detection to achieve better performance. Also, the most popular benchmarks like HumanEval and MBPP check only functional correctness, that is, whether a code solves the corresponding problem as it is stated in a prompt. The verification of a code is done by running functional tests, as explained in Section 5.1.

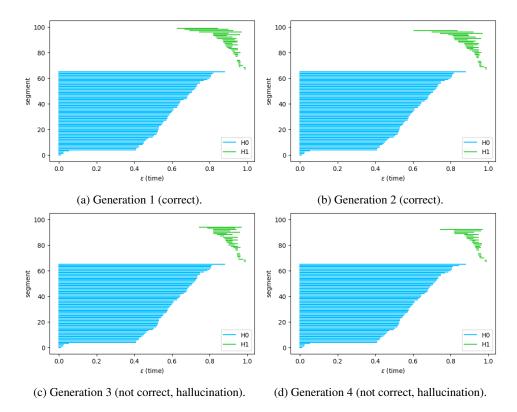


Figure 12: Examples of $Cross-Barcode_0(G, P)$, $Cross-Barcode_1(G, P)$. CodeLLama-7B, HumanEval dataset, problem 14, layer 15, head 5.

H LIMITATIONS

Although we have achieved good experimental results, we realize that our research have several limitations. Our method targets hallucinations that manifest in the model's attention geometry; it does not disentangle all root causes (e.g. spurious pre-training correlations, decoding drift, RLHF bias). Extending the analysis to those factors is left for future work. We mostly explored code LLMs having no more than 7B parameters. Information in larger models are more distributed in attention heads and results might differ. Also, processing more attention heads is computationally costly. Next, the proposed classifiers of hallucinations are based on the attention maps of the same code LLMs as for code generations. We leave more general setting to a further research. Finally, our approach can predict whether a code is correct as a whole but can not point to a specific place with a bug.

I ADDITIONAL ABLATION STUDY

Our main argument in the Ablation study is that different models benefit from different types of attention features. In this section, we provide additional evaluation to further support our claim. We observe that the smaller QwenCoder-1.5b and QwenCoder-3b models experience substantial performance decrease when removing the topological features, see Table 7. To further support our observations, we provide some examples of transferability from MBPP to HE datasets where the contribution of diagonal and MTD features is different for different models, see Table 8. Therefore, we propose to account for both types of features in our approach or at least be aware of potential benefits of using each type.

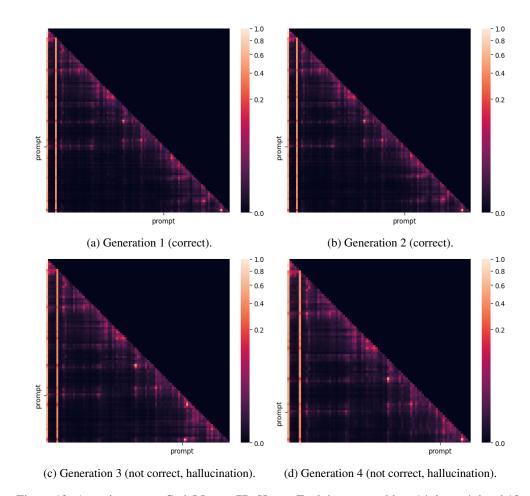


Figure 13: Attention maps. CodeLLama-7B, HumanEval dataset, problem 14, layer 4, head 18.

J EXPERIMENTS WITH LARGER MODELS

We carried out additional experiments, where code hallucinations are detected by a recent reasoning DeepSeek R1 model having 671 billion parameters with Chain of Thought inference on MBPP dataset. We used the following prompt:

You are provided with two coding tasks with solutions. The first one is just an example and it does not need an assessment. Tell whether the second task is correctly solved by the code provided in the second [BEGIN] [DONE] block. The answer must be Yes or No

For code generated with DeepSeek-6.7B model, we prompted DeepSeek R1 using the prompt above, extracted the final answers (i.e. "Yes" or "No") from the generated responses and trained XGboost classifier using these features. The quality of hallucination detection via such zero-shot prompting is in Table 10, row "DeepSeek R1, (671B model)". Quality of the proposed approach is in Table 10, row "Attn. Feat (ours, from 6.7B model)". Note that in this case we use *only* attention features from DeepSeek-6.7B model obtained during code generation. Finally, we combine both type of features to train a classifier and report its performance in Table 10, row "Attn. Feat. (ours, from 6.7B model) + DeepSeek R1". The larger DeepSeek R1 model demonstrates a higher performance than the classifier trained on attention features. However, by adding an output of DeepSeek R1 to our attention-based features and training XGBoost classifier we can achieve the best ROC-AUC score. Study of DeepSeek R1's Chain of Thoughts shows that this LLM is doing verification of code by interpreting Python code step by step for unit tests. This can explain the high accuracy of Deep Seek R1. At the same time, our method opens opportunities for deeper understanding of inner working and information flow inside transformer models. Our attention features were based on a

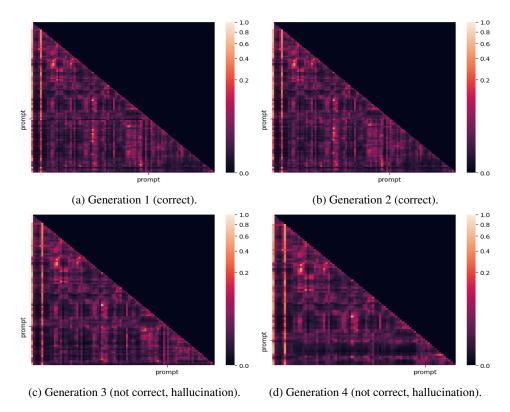


Figure 14: Attention maps. CodeLLama-7B, HumanEval dataset, problem 14, layer 15, head 5.

small 6.7B model in this experiment, however, our features were able to improve the performance of DeepSeek R1.

Additionally, we perform similar experiment with QwenCoder2.5-32B. In this case, we use the same QwenCoder2.5-32B to generate code, extract attention features and evaluate its performance with zero-shot prompting. Table 11 provides the experimental results. The proposed approach is able to achieve better detection quality than zero-shot prompting which supports the applicability of the proposed approach to larger models.

K ADDITIONAL RESULTS ON CONTRIBUTION OF INDIVIDUAL HEADS

We carry out additional experiments with 7B-models and MultiPL-E benchmark⁸ which is a translation of HumanEval to several popular programming languages; we used Go, Java, Rust, Lua among them. We find that features of some heads has quite high correlation with target value (presence of a hallucination) and can be used as individual predictors. We report in Table 14 ROC AUC scores of the top-performing features.

⁸https://github.com/nuprl/MultiPL-E

Table 6: Ablation study for the choice of a classification model for code hallucination detection.

Model	Huma	anEval	ME	3PP
	ROC-AUC	F1-Score	ROC-AUC	F1-Score
	Star	Coder2-7B		
XGBoost	82.9 ± 2.7	54.2 ± 6.9	81.9 ± 2.4	68.4 ± 5.3
MLP	80.9 ± 4.7	60.7 ± 7.6	81.1 ± 1.7	68.7 ± 3.0
Logistic Regression	84.3 ± 3.9	64.2 ± 8.0	82.9 ± 2.3	69.7 ± 3.4
SVC	84.9 ± 4.0	65.6 ± 7.4	82.3 ± 2.6	69.1 ± 4.4
		eLlama-7B		
XGBoost	85.6 ± 3.9	56.4 ± 7.2	83.4 ± 3.3	64.0 ± 4.4
MLP	81.8 ± 7.2	58.1 ± 11.8	81.6 ± 2.2	64.4 ± 5.6
Logistic Regression	81.8 ± 7.1	58.2 ± 7.9	82.6 ± 2.1	63.6 ± 4.7
SVC	83.2 ± 5.4	59.8 ± 9.2	82.4 ± 1.2	60.9 ± 4.6
		ek-Coder-6.7B		
XGBoost	85.6 ± 2.8	68.9 ± 5.5	82.6 ± 1.9	$\textbf{76.5} \pm \textbf{2.7}$
MLP	84.3 ± 2.4	69.6 ± 3.9	81.7 ± 1.1	74.8 ± 1.9
Logistic Regression	85.8 ± 3.2	71.1 ± 4.7	81.8 ± 2.4	75.7 ± 2.2
SVC	87.0 ± 2.4	72.9 ± 3.7	81.4 ± 1.6	75.2 ± 1.7
	Qwen2	2.5-Coder-7B		
XGBoost	81.7 ± 2.8	70.2 ± 4.2	82.2 ± 2.2	75.4 ± 1.7
MLP	81.3 ± 1.6	70.8 ± 2.5	79.5 ± 2.0	72.9 ± 3.0
Logistic Regression	80.0 ± 1.6	71.3 ± 3.8	81.0 ± 1.9	75.9 ± 2.0
SVC	79.6 ± 1.7	68.9 ± 2.4	81.0 ± 1.7	$\textbf{76.7} \pm \textbf{3.3}$
		der-S-DS-6.7B		
XGBoost	82.3 ± 4.9	80.7 ± 3.6	77.8 ± 2.5	73.4 ± 3.4
MLP	76.5 ± 3.3	78.9 ± 1.8	78.6 ± 3.6	73.3 ± 3.5
Logistic Regression	81.7 ± 1.6	81.6 ± 2.4	81.3 ± 2.8	82.1 ± 1.7
SVC	80.5 ± 2.4	82.0 ± 3.1	82.5 ± 2.6	81.3 ± 1.9

Table 7: HumanEval and MBPP features ablation for smaller models.

Model	HumanEval		MBPP	
	ROC-AUC	F1-Score	ROC-AUC	F1-Score
	Qwer	Coder2.5-1.5b		
Attn. Feat. (ours)	$\textbf{85.9} \pm \textbf{4.1}$	58.1 ± 9.6	82.8 ± 1.4	$\textbf{59.2} \pm \textbf{2.7}$
 w/o Diag. Feat. 	85.0 ± 3.8	57.6 ± 6.8	82.1 ± 1.8	58.6 ± 4.2
- w/o MTD Feat.	82.9 ± 3.9	53.6 ± 10.8	79.5 ± 1.5	57.0 ± 4.5
QwenCoder2.5-3b				
Attn. Feat. (ours)	80.4 ± 6.6	65.5 ± 8.5	78.0 ± 3.0	73.3 ± 3.6
- w/o Diag. Feat.	79.0 ± 6.1	67.2 ± 6.4	79.6 ± 3.0	74.8 ± 3.9
- w/o MTD Feat.	74.9 ± 4.4	59.0 ± 5.7	75.5 ± 2.2	71.0 ± 3.1

Table 8: MBPP transferability ablation.

Method	DeepSeek-Coder-6.7b	QwenCoder2.5-3b	StarCoder2-7b
Attn. Feat. (ours)	$\bf 72.4$	65.6	67.2
- w/o Diag. Feat.	71.3	67.4	64.2
- w/o MTD Feat.	63.5	61.3	66.8

Table 9: HumanEval and MBPP features ablation

Model	Hum	anEval	MBPP	
	ROC-AUC	F1-Score	ROC-AUC	F1-Score
	St	arCoder2-7B		
Attn. Feat. (ours)	82.9 ± 2.7	54.2 ± 6.9	81.9 ± 2.4	$\textbf{68.4} \pm \textbf{5.3}$
 w/o Diag. Feat. 	82.2 ± 4.5	56.1 ± 9.7	80.5 ± 2.8	66.3 ± 5.3
- w/o MTD Feat.	83.8 ± 2.7	52.5 ± 8.4	81.1 ± 2.6	67.7 ± 5.0
	Co	deLlama-7B		
Attn. Feat. (ours)	$\textbf{85.6} \pm \textbf{3.9}$	56.4 ± 7.2	83.4 ± 2.2	$\overline{64.0 \pm 4.4}$
 w/o Diag. Feat. 	83.5 ± 4.8	50.0 ± 6.5	81.5 ± 2.6	60.2 ± 4.2
- w/o MTD Feat.	85.5 ± 4.4	58.3 ± 10.1	83.5 ± 1.8	63.9 ± 4.3
	DeepS	Seek-Coder-6.7E	3	
Attn. Feat. (ours)	85.6 ± 2.8	68.9 ± 5.5	82.6 ± 1.9	$\textbf{76.5} \pm \textbf{2.7}$
 w/o Diag. Feat. 	85.1 ± 2.2	67.0 ± 5.9	81.3 ± 2.6	74.9 ± 3.2
- w/o MTD Feat.	84.4 ± 2.2	67.1 ± 3.8	82.2 ± 1.7	75.9 ± 1.7
	Qwe	n2.5-Coder-7B		
Attn. Feat. (ours)	81.7 ± 2.8	$\textbf{70.2} \pm \textbf{4.2}$	82.2 ± 2.2	$\textbf{75.4} \pm \textbf{1.7}$
 w/o Diag. Feat. 	80.6 ± 2.3	68.9 ± 3.9	80.8 ± 2.1	$\textbf{75.4} \pm \textbf{0.4}$
- w/o MTD Feat.	78.9 ± 1.9	66.4 ± 1.4	76.9 ± 2.2	71.6 ± 1.7
Magicoder-S-DS-6.7B				
Attn. Feat. (ours)	82.3 ± 4.9	80.7 ± 3.6	81.5 ± 1.6	80.6 ± 2.0
 w/o Diag. Feat. 	79.8 ± 2.7	81.1 ± 1.8	81.2 ± 1.9	80.1 ± 1.4
- w/o MTD Feat.	82.1 ± 3.4	81.6 ± 2.6	80.3 ± 2.5	79.7 ± 2.2

Table 10: MBPP features for larger models. See Section J for details.

Method	ROC-AUC
DeepSeek-Coder-6.7B	
Attn. Feat (ours, from 6.7B model)	82.6 ± 1.9
DeepSeek R1 (671B model)	92.3 ± 1.1
Attn. Feat. (from 6.7B model) + DeepSeek R1 (671B)	95.1 ± 0.7

Table 11: HumanEval features for larger models. See Section J for details.

Method	ROC-AUC	F1-Score
Qwen	Coder2.5-32B	
Attn. Feat (ours)	85.0 ± 2.7	77.0 ± 5.5
Zero-shot prompt	67.4 ± 3.5	71.8 ± 3.6

Table 12: Code hallucination detection for HumanEval dataset for larger models. For each task, 10 candidate solutions were generated. **Bold** and <u>underline</u> denote the first and second best results.

Method	ROC-AUC	F1-Score			
CodeLlama-34B					
Prompt Len.	57.3 ± 6.3	37.7 ± 13.6			
Gen. Len.	62.8 ± 1.8	38.2 ± 4.0			
Mean Log. Prob.	74.1 ± 5.0	51.8 ± 7.6			
CodeT5-base ft.	54.7 ± 5.5	0.0 ± 0.0			
Attn. Feat. (ours)	83.2 ± 3.8	62.5 ± 7.1			
All. Feat.	$\overline{84.8 \pm 2.8}$	$\overline{62.8 \pm 6.7}$			
Qwen	2.5-Coder-32E	}			
Prompt Len.	53.2 ± 9.2	53.4 ± 12.8			
Gen. Len.	58.0 ± 3.1	62.5 ± 3.8			
Mean Log. Prob.	65.8 ± 4.8	63.6 ± 2.8			
CodeT5-base ft.	53.3 ± 4.5	59.5 ± 15.2			
Attn. Feat. (ours)	85.0 ± 2.7	77.0 ± 5.5			
All. Feat.	84.9 ± 3.1	$\overline{\textbf{77.7} \pm \textbf{5.7}}$			
DeepS	eek-Coder-331	3			
Prompt Len.	53.3 ± 6.5	44.0 ± 5.3			
Gen. Len.	56.6 ± 4.2	48.2 ± 3.4			
Mean Log. Prob.	66.4 ± 3.4	57.1 ± 2.9			
CodeT5-base ft.	57.6 ± 3.9	1.1 ± 2.2			
Attn. Feat. (ours)	88.0 ± 2.9	77.9 ± 2.2			
All. Feat.	$\overline{88.9 \pm 2.8}$	$\overline{78.6 \pm 3.1}$			

Table 13: Code hallucination detection with attention features for MBPP dataset with increasing complexity of prompt: two-shot prompts.

Method	ROC-AUC					
Qwen2.5-Coder-7B						
Prompt Len.	53.5 ± 2.9	60.2 ± 1.5				
Gen. Len.	56.9 ± 3.3	63.2 ± 2.0				
Mean Log. Prob.	58.1 ± 2.6	61.1 ± 3.3				
Attn. Feat. (ours)	78.3 ± 3.1	74.2 ± 1.8				
All. Feat.	76.7 ± 2.7	72.2 ± 2.1				
Self-Eval	68.4 ± 1.0	77.4 ± 2.4				
Interrogate-LLM	56.6 ± 1.5	69.1 ± 2.1				
DeepSe	DeepSeek-Coder-6.7B					
Prompt Len.	56.0 ± 2.0	58.0 ± 2.4				
Gen. Len.	54.4 ± 1.3	59.6 ± 1.6				
Mean Log. Prob.	64.9 ± 1.7	64.2 ± 2.4				
Attn. Feat. (ours)	81.4 ± 2.2	75.0 ± 5.3				
All. Feat.	80.4 ± 0.7	74.5 ± 2.6				
Self-Eval	50.0 ± 0.0	69.3 ± 5.1				
Interrogate-LLM	61.1 ± 3.1	69.4 ± 5.2				

Feature

avg. prompt's self-attention, layer 14, head 0

avg. prompt's self-attention, layer 15, head 5

avg. prompt's self-attention, layer 23, head 20

 $-MTD_1(P, G)/|P|$, layer 15, head 27

avg. prompt's self-attention, layer 7, head 22

 $MTD_0(P, G)/|P|$, layer 11, head 23

 $MTD_0(P, G)/|P|$, layer 30, head 11

avg. prompt's self-attention, layer 12, head 17

avg. prompt's self-attention, layer 12, head 23

avg. generation's self-attention, layer 24, head 2

 $MTD_0(P, G)/|P|$, layer 11, head 5

avg. prompt's self-attention, layer 9, head 26

 $MTD_0(P, G)/|P|$, layer 30, head 11

avg. prompt's self-attention, layer 13, head 13

avg. prompt's self-attention, layer 12, head 17

5 candidate solutions for MBPP.

11881189

Table 14: ROC AUC scores of top-performing features across several benchmarks.

StarCoder2-7B

CodeLlama-7B

DeepSeek-Coder-6.7B

Qwen2.5-Coder-7B

Magicoder-S-DS-6.7B

Python

70.8

71.1

69.2

67.6

74.6

69.1

67.5

65.0

64.1

65.8

66.2

65.4

68.8

63.1

63.0

HumalEval

Rust

74.6

75.1

74.2

73.7

76.5

71.9

67.0

70.0

66.8

58.8

59.5

66.9

60.0

69.0

68.0

Go

76.8

78.4

72.2

71.1

75.7

71.0

65.7

73.6

71.0

60.2

68.8

75.8

70.2

64.9

63.2

MBPP

Python

55.2

58.1

50.7

66.1

63.1

69.1

58.1

58.2

58.2

60.8

59.4

61.0

58.2

58.8

58.9

Lua

70.3

72.4

67.7

69.1

69.8

67.6

66.7

72.7

71.1

63.8

63.2

60.8

Java

68.9

72.9

64.7

71.0

69.7

65.4

66.3

69.8

64.4

60.2

61.6

67.0

65.1

67.3

67.1

•	1	v	9
1	1	9	0
1	1	9	1
1	1	9	2
1	1	9	3
1	1	9	4
1	1	9	5
1	1	9	6
1	1	9	7
1	1	9	8
1	1	9	9
		0	
1	2	0	1
1	2	0	2
1	2	0	3
1	2	0	4
1	2	0	5
1	2	0	6
1	2	0	7
1	2	0	8
1		0	
1		1	
1	2	1	1

Table 15: Characteristics of generated data, Pass@1. For each problem, we generated 25 candidate solutions for HumanEval, Python, 10 candidate solutions for HumanEval, Java, Go, Rust, Lua, and

1217
1218
1219
1220
1221

1212

1213

1214

1215 1216

Model	HumalEval MB			MBPP		
	Python	Java	Go	Rust	Lua	Python
StarCoder2-7B	28.9	24.5	17.5	20.9	19.1	42.8
CodeLlama-7B	25.9	25.8	17.6	20.8	0.0	35.2
DeepSeek-Coder-6.7B	40.3	33.5	23.6	28.7	16.6	52.6
Qwen2.5-Coder-7B	47.8	22.7	11.9	22.6	23.5	52.1
Magicoder-S-DS-6.7B	65.5	48.9	40.3	44.2	34.6	61.3

1222 1223 1224 1225

1226 1227 1228

Table 16: Performance of the proposed method on different programming languages, ROC-AUC.

1	229
1	230
1	231
1	232

Model	Java	Go	Rust	Lua
StarCoder2-7B	82.7 ± 4.9	86.5 ± 3.6	82.4 ± 5.2	82.0 ± 3.5
CodeLlama-7B	76.8 ± 6.1	81.9 ± 6.2	77.5 ± 10.8	_
DeepSeek-Coder-6.7B	84.9 ± 4.2	84.7 ± 2.9	82.8 ± 7.1	86.7 ± 4.3
Qwen2.5-Coder-7B	82.6 ± 4.1	91.8 ± 0.9	87.3 ± 2.6	90.2 ± 3.0
Magicoder-S-DS-6.7B	77.8 ± 4.5	80.8 ± 2.0	75.1 ± 5.9	79.0 ± 3.2

1232 1233 1234

1235

Table 17: Performance of the proposed method on different programming languages, F1-Score.

1236
1237
1238
1239
1240

```
Model
                              Java
                                               Go
                                                              Rust
                                                                              Lua
StarCoder2-7B
                           50.4 \pm 7.6
                                          39.7 \pm 11.9
                                                          36.5 \pm 9.2
                                                                          39.0 \pm 8.7
CodeLlama-7B
                           36.4 \pm 20.6
                                          26.5 \pm 10.2
                                                          34.8 \pm 17.3
DeepSeek-Coder-6.7B
                           64.1 \pm 4.5
                                           53.2 \pm 8.7
                                                          56.8 \pm 6.6
                                                                         41.6 \pm 15.1
                                                          49.9 \pm 6.9
Qwen2.5-Coder-7B
                           45.8 \pm 9.6
                                          37.6 \pm 11.5
                                                                          64.6 \pm 7.9
                                                          63.1 \pm 8.7
Magicoder-S-DS-6.7B
                           68.0 \pm 4.9
                                           63.3 \pm 2.0
                                                                          58.4 \pm 6.7
```