

Topology of Attention: Detecting Hallucinations in Code Generation Models

Anonymous ACL submission

Abstract

While the AI-code assistant tools become widespread, automatic assessment of the correctness of the generated code becomes a significant challenge. Code-generating LLMs are prone to hallucinations, which may lead to code that doesn't 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. 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 a great potential to transform natural language processing and artificial intelligence. As far as code generation is concerned, LLMs which are trained on large amounts of code, are capable to generate 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. These hallucination are of various types. Sometimes generated code has syntactic or logical errors, sometimes it is correct but do not solve a required problem. In some cases, the generated code might contain security issues or robustness issues, like a memory leak. While many definitions of hallucinations exist, in this paper we assume that code hallucination is a code which do not pass tests. For a wide adoption of code generating LLMs, there is a high need of automatic assessment of code quality. As for the current state of technologies, a

significant time is spent to 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 a illustration to model's "thought process". The research community actively studies approaches to mitigate hallucinations of LLMs by external knowledge bases (Peng et al., 2023) or reduce them to some degree (Elaraby et al., 2023). It is a highly desirable to evaluate to 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.

Our contributions are the following:

- We propose a new approach to detection of hallucinations in LLM generated code based on analysing a topology of attention maps;
- We carry out computational experiments with CodeLlama, StarCoder2, DeepSeek-Coder and Qwen2.5-Coder and two benchmarks –

081 HumanEval and MBPP, and show that the pro- 131
 082 posed method outperforms baselines; 132

- 083 • We empirically show that proposed classifier 133
 084 of hallucinations is transferable between code 134
 085 benchmarks. 135

086 2 Related work 136

087 Code generation via Large Language Models 137
 088 (LLMs) is the topic of active research. The popular 138
 089 projects are: CodeLlama (Roziere et al., 2023), 139
 090 StarCoder2 (Lozhkov et al., 2024), DeepSeek- 140
 091 Coder (Guo et al., 2024), Qwen2.5-Coder (Hui 141
 092 et al., 2024), to name a few. Code LLMs differ by 142
 093 data were used for training, by their training and 143
 094 fine-tuning protocols, including RLHF, tokenizers, 144
 095 variants of attention mechanism, etc. 145

096 Several works studied attention maps in 146
 097 transformer-based LLMs. (Clark, 2019) studied 147
 098 BERT’s attention patterns: attending to delimiter 148
 099 tokens, specific positional offsets, or broadly at- 149
 100 tending over the whole sentence, with heads in 150
 101 the same layer often exhibiting similar behaviors. 151
 102 (Clark, 2019) further showed that certain attention 152
 103 heads correspond well to linguistic notions of syn- 153
 104 tax and coreference. (Htut et al., 2019) found that 154
 105 for some universal dependency tree relation types, 155
 106 there exist heads that can recover the dependency 156
 107 type significantly better than baselines on parsed 157
 108 English text, suggesting that some self-attention 158
 109 heads act as a proxy for syntactic structure. (Michel 159
 110 et al., 2019) showed that for downstream tasks, a 160
 111 large proportion of attention heads can be removed 161
 112 at test time without significantly impacting perfor- 162
 113 mance, and that some layers can even be reduced 163
 114 to a single head. 164

115 The phenomenon of code hallucinations is stud- 165
 116 ied and categorized several papers. (Tian et al., 166
 117 2024) introduces a categorization of code halluci- 167
 118 nations into four main types: mapping, naming, 168
 119 resource, and logic hallucinations, with each cat- 169
 120 egory further divided into different subcategories. 170
 121 (Tian et al., 2024) proposed a CodeHalu dataset 171
 122 and studied frequencies of different types of hallu- 172
 123 cinations in popular code LLMs. (Liu et al., 2024) 173
 124 categorized hallucinations into: intent conflicting, 174
 125 inconsistency, repetition, knowledge conflicting, 175
 126 dead code. (Liu et al., 2024) released a HaluCode 176
 127 benchmark with labeled code hallucinations. (Jiang 177
 128 et al., 2024) proposed Collu-Bench, the benchark 178
 129 with localization of code hallucinations. (Jiang 179
 130 et al., 2024) found that code LLMs are less confi-

131 dent when hallucinating, as the hallucinated tokens 132
 133 have lower probability and hallucinated generation 134
 135 steps have higher entropy. 136

137 In the broader context of NLP, several works in- 138
 139 troduced methods to hallucination preventing and 139
 140 detection. (Peng et al., 2023) proposed to miti- 140
 141 gate hallucination by an LLM-AUGMENTER, a 141
 142 system which makes the LLM generate responses 142
 143 grounded in external knowledge, e.g., stored in 143
 144 task-specific databases. (Zhang et al., 2024b) pro- 144
 145 posed Self-Eval, a self-evaluation component, to 145
 146 prompt an LLM to validate the factuality of its 146
 147 own generated responses solely based on its in- 147
 148 ternal knowledge. (Feng et al., 2024) proposed 148
 149 two novel approaches for hallucination detection 149
 150 that are based on model collaboration, i.e., LLMs 150
 151 probing other LLMs for knowledge gaps, either co- 151
 152 operatively or competitively. (Zhang et al., 2024a) 152
 153 proposed to improve truthfulness of LLMs by edit- 153
 154 ing their internal representation during inference 154
 155 in the “truthful” space. (Yehuda et al., 2024) in- 155
 156 troduced InterrogateLLM, a method which prompts 156
 157 the model multiple times to reconstruct the input 157
 158 query using the generated answer. Subsequently, 158
 159 InterrogateLLM quantifies the inconsistency level 159
 160 between the original query and the reconstructed 160
 161 queries. 161

158 3 Background 158

159 3.1 Transformer-based LLMs 159

160 All of the state-of-the art LLMs for code genera- 160
 161 tion networks are based on different variants of the 161
 162 transformer architecture (Vaswani, 2017). A trans- 162
 163 former architecture comprises L layers of multi- 163
 164 head self-attention blocks each of them having H 164
 165 heads. Each attention head takes $X \in \mathbb{R}^{n \times d}$ ma- 165
 166 trix as an input, and an output of attention head in 166
 167 X^{out} : 167

$$168 \quad X^{out} = A(XW^v), \quad 168$$

$$169 \quad A = \text{softmax} \left(\frac{(XW^Q)(XW^K)^T}{\sqrt{d}} \right), \quad 169$$

170 where $W^Q, W^K, W^V \in \mathbb{R}^{d \times d}$ are projection 170
 171 matrices and $A \in [0, 1]^{n \times n}$ is an **attention map**. 171
 172 In self-attention block, the attention map shows 172
 173 how each token in the input sequence “interacts” to 173
 174 every other token in the same sequence. A token 174
 175 might attend more to other tokens that are contex- 175
 176 tually related. We interpret each element $a_{i,j}$ of 176

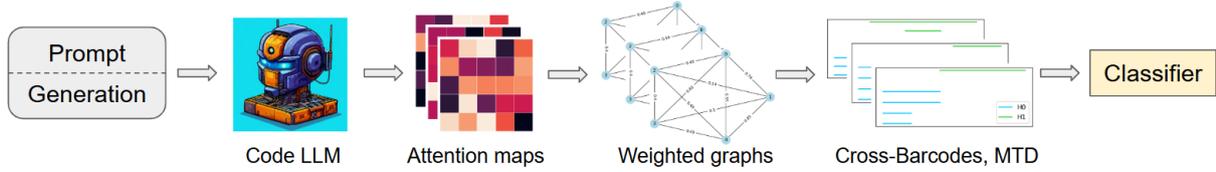


Figure 1: A pipeline of the proposed method for hallucination detection.

an attention map as an “interaction force” between tokens i and j .

3.2 Representing 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 and 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 $\mathcal{G}_0 \subset \dots \subset \mathcal{G}_i \subset \mathcal{G}_{i+1} \subset \dots \subset \mathcal{G}$ 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 2. During this process, graphs’ topology naturally changes: connected components are merged, cycles appear and disappear, etc. This process is rigorously described by the persistence barcodes theory (Barannikov, 1994; Chazal and Michel, 2017). Each topological feature like connected component or cycle has “birth time” and “death time”, by a corresponding edge weight. The multi-set of these birth-death pairs (intervals) al-

together is called a Cross-Barcode $_k$, see Figure 3. 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).

4 Methods

In the context of code generation, we naturally have two sets of tokens – a prompt and a generation. Intuitively, hallucination happens when code LLM *doesn’t pay much attention* to the prompt. 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 calculated. To predict code hallucinations, we use the following set of features:

- $\text{MTD}_0(P, G)/|P|, \text{MTD}_0(G, P)/|G|$ 245
- $\text{MTD}_1(P, G)/|P|, \text{MTD}_1(G, P)/|G|$ 246
- $\sum_{i \in P} a_{i,i}/|P|, \sum_{i \in G} a_{i,i}/|G|$ 247

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 and Guestrin, 2016) for a classification. The high-level pipeline of the proposed method is shown in Figure 1.

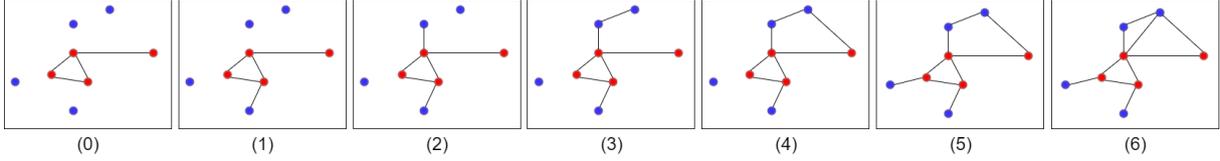


Figure 2: An example of MTD evaluation for a graph having two groups of vertices – red and blue. (0): initially, only edges connecting red vertices are present. (1)-(6): the rest of edges are added sequentially in an ascending order by their weights. While adding edges, connected components merge with each other. These moments are depicted by H_0 bars in Fig. 3. At moment (4) a cycle appears, at moment (6) this cycle disappears. These moments are depicted by the H_1 bar in Fig. 3.

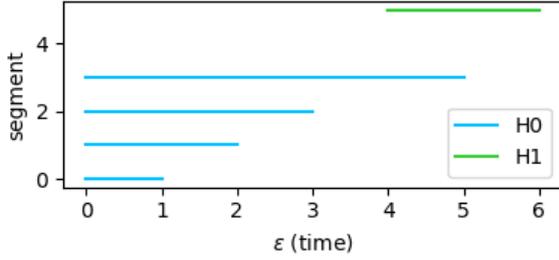


Figure 3: Cross-Barcode for a filtration from Fig. 2.

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

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

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. We use the following popular code LLMs: 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), Magocoder-S-DS-6.7B (Wei et al., 2024). We adapted two public benchmarks for evaluation of code genera-

tion: HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021)¹. 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. 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 tests provided together with the coding benchmarks. 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 5, 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

¹Licenses of pretrained models and benchmarks permit use for research purposes.

Method	ROC-AUC	F1-Score
StarCoder2-7B		
Prompt Len.	54.5 ± 6.6	24.6 ± 11.2
Gen. Len.	57.7 ± 5.6	13.5 ± 4.8
Mean Log. Prob.	70.9 ± 1.3	32.4 ± 4.8
CodeT5-base ft.	70.1 ± 7.1	33.3 ± 10.1
Attn. Feat. (ours)	82.9 ± 2.7	54.2 ± 6.9
CodeLlama-7B		
Prompt Len.	61.6 ± 4.4	25.7 ± 15.0
Gen. Len.	60.1 ± 5.3	10.6 ± 7.0
Mean Log. Prob.	64.1 ± 2.0	25.4 ± 6.2
CodeT5-base ft.	74.5 ± 6.3	43.6 ± 13.2
Attn. Feat. (ours)	85.6 ± 3.9	56.4 ± 7.2
DeepSeek-Coder-6.7B		
Prompt Len.	56.2 ± 4.6	44.4 ± 4.3
Gen. Len.	57.9 ± 2.4	34.4 ± 4.9
Mean Log. Prob.	69.8 ± 2.5	51.1 ± 3.4
CodeT5-base ft.	69.1 ± 4.2	52.6 ± 6.5
Attn. Feat. (ours)	85.6 ± 2.8	68.9 ± 5.5
Qwen2.5-Coder-7B		
Prompt Len.	54.3 ± 8.7	51.0 ± 5.7
Gen. Len.	57.6 ± 3.6	48.9 ± 5.1
Mean Log. Prob.	63.1 ± 2.4	55.6 ± 5.5
CodeT5-base ft.	65.9 ± 3.7	58.2 ± 4.5
Attn. Feat. (ours)	81.7 ± 2.8	70.2 ± 4.2
Magocoder-S-DS-6.7B		
Prompt Len.	57.3 ± 5.4	70.4 ± 7.0
Gen. Len.	52.5 ± 2.1	76.3 ± 2.6
Mean Log. Prob.	71.0 ± 5.3	78.4 ± 2.9
CodeT5-base ft.	64.7 ± 2.7	77.5 ± 2.7
Attn. Feat. (ours)	82.3 ± 4.9	80.7 ± 3.6

Table 2: Code hallucination detection for HumanEval dataset.

307 belonged to the same group. In this way, training
308 and testing were performed always at non-
309 overlapping coding problems (prompts). The re-
310 ported results are the mean and standard deviation
311 estimated over the 5 folds. As baselines for compar-
312 ison, we used XGBoost classifier trained on simple
313 features: tokenized prompt length, tokenized gener-
314 ation length, and mean log-probability of generated
315 tokens (Chen et al., 2021). Also, we trained a linear
316 classification head on top of a frozen CodeT5-base
317 (Wang et al., 2021) encoder. Training details are
318 provided in Appendix B. Tables 2, 3 present re-
319 sults. In the majority of cases, the proposed classi-
320 fier based on features of attention maps performed
321 significantly better than the baselines and demon-
322 strated stable results for all models and datasets as
323 measured by ROC-AUC score. Further analysis re-
324 vealed that some features made a high contribution

Method	ROC-AUC	F1-Score
StarCoder2-7B		
Prompt Len.	51.2 ± 2.3	40.0 ± 3.8
Gen. Len.	57.7 ± 0.9	45.4 ± 3.5
Mean Log. Prob.	62.0 ± 2.0	47.5 ± 3.4
CodeT5-base ft.	58.5 ± 3.5	43.3 ± 9.0
Attn. Feat. (ours)	81.9 ± 2.4	68.4 ± 5.3
CodeLlama-7B		
Prompt Len.	59.1 ± 4.2	35.4 ± 2.9
Gen. Len.	60.8 ± 2.5	24.2 ± 5.3
Mean Log. Prob.	61.0 ± 3.7	27.2 ± 1.5
CodeT5-base ft.	61.7 ± 3.0	19.1 ± 7.1
Attn. Feat. (ours)	83.4 ± 3.3	64.0 ± 4.4
DeepSeek-Coder-6.7B		
Prompt Len.	52.5 ± 2.5	56.4 ± 3.6
Gen. Len.	54.6 ± 1.9	59.4 ± 1.3
Mean Log. Prob.	61.0 ± 1.9	62.3 ± 1.6
CodeT5-base ft.	55.7 ± 3.0	64.8 ± 2.7
Attn. Feat. (ours)	82.6 ± 1.9	76.5 ± 2.7
Qwen2.5-Coder-7B		
Prompt Len.	51.8 ± 3.6	56.2 ± 4.4
Gen. Len.	55.6 ± 2.1	59.7 ± 4.8
Mean Log. Prob.	61.5 ± 1.3	60.4 ± 1.8
CodeT5-base ft.	56.0 ± 1.3	65.2 ± 2.0
Attn. Feat. (ours)	82.2 ± 2.2	75.4 ± 1.7
Magocoder-S-DS-6.7B		
Prompt Len.	52.5 ± 2.5	56.4 ± 3.6
Gen. Len.	58.7 ± 1.1	60.7 ± 2.1
Mean Log. Prob.	60.6 ± 3.7	72.1 ± 1.4
CodeT5-base ft.	61.0 ± 3.7	74.8 ± 1.4
Attn. Feat. (ours)	77.8 ± 2.5	73.4 ± 3.4

Table 3: Code hallucination detection for MBPP dataset.
to the classification quality, see Figure 4.

5.3 Analyzing method’s ranking quality

325
326
327 Next, we assess the usefulness of the proposed code
328 hallucination classifier for ranking of code gener-
329 ations. For each problem, all generations were
330 ranked via probability of correctness predicted by
331 the classifier and one with the highest probability
332 was selected. A baseline was random picking of a
333 code generation. The usage of a classifier is always
334 significantly better by a pass@1 score, see Table 4.

5.4 Method’s transferability between benchmarks

335
336
337 We study further the transferability of the classi-
338 fiers, based on topological features. In this setting,
339 hallucination classifiers for a fixed code LLM are
340 trained on data for one benchmark (HumanEval,
341 MBPP) and evaluated on another, then repeated

Model	Random	Clf. Prob.
HumanEval		
StarCoder2-7B	28.6 ± 5.5	43.3 ± 9.0
CodeLlama-7B	26.0 ± 5.1	39.7 ± 7.2
DeepSeek-Coder-6.7B	39.1 ± 4.9	56.7 ± 7.4
Qwen2.5-Coder-7B	51.8 ± 8.0	64.0 ± 7.3
Magicoder-S-DS-6.7B	72.5 ± 10.0	74.3 ± 6.1
MBPP		
StarCoder2-7B	43.0 ± 3.6	49.6 ± 4.6
CodeLlama-7B	35.2 ± 3.3	43.6 ± 3.4
DeepSeek-Coder-6.7B	53.0 ± 2.5	61.4 ± 2.3
Qwen2.5-Coder-7B	52.6 ± 3.6	62.0 ± 2.4
Magicoder-S-DS-6.7B	61.4 ± 3.4	63.8 ± 2.0

Table 4: pass@1 for random choice vs argmax of classifier probability

342 vice versa. Tables 5, 6 shows results: the proposed
343 classifiers are transferable, albeit the performance
344 is lower when training and testing is done on the
345 same benchmark.

346 5.5 Ablation study

347 The proposed approach is based on the two types
348 of attention features: the diagonal elements of at-
349 tention maps corresponding to the prompt and gener-
350 ation and topological features computed for the
351 corresponding ‘‘attention graph’’ (see Section 4 for
352 details). In this Section, we provide an ablation
353 study to estimate the contribution of each type of
354 attention features. For this purpose, we trained
355 the XGBoost classifier using only MTD features
356 (i.e. without the diagonal elements of attention
357 maps) or using only diagonal attention values (i.e.
358 without MTD features) and compared its perfor-
359 mance with the initial setup where both types of
360 attention features were used. As demonstrated
361 with Tables 7, 8, the DeepSeek-Coder-6.7B and
362 Qwen2.5-Coder-7B achieved the best performance
363 when both types of attention features were used
364 for both HumanEval and MBPP datasets. In con-
365 trast, the best performance of StarCoder2-7B and
366 Magicoder-S-DS-6.7B was achieved with different
367 sets of attention features dependent on dataset and
368 metric choices. In order to account for various in-
369 formation available via attention maps, we propose
370 to use both types of features as the most universal
371 choice. Nevertheless, we note that for some code
372 LLM one certain type of attention features may
373 result in better performance than combination of
374 both types.

Model	ROC-AUC	F1-Score
StarCoder2-7B		
Prompt Len.	48.6	0.0
Gen. Len.	56.0	14.6
Mean Log. Prob.	63.7	36.2
CodeT5-base ft.	53.7	0.0
Attn. Features	67.5	0.14
CodeLlama-7B		
Prompt Len.	51.7	0.0
Gen. Len.	61.5	4.2
Mean Log. Prob.	57.7	15.2
CodeT5-base ft.	54.9	0.0
Attn. Features	69.5	0.2
DeepSeek-Coder-6.7B		
Prompt Len.	48.0	15.6
Gen. Len.	55.3	41.4
Mean Log. Prob.	62.5	56.3
CodeT5-base ft.	53.4	0.0
Attn. Features	67.7	70.4
Qwen2.5-Coder-7B		
Prompt Len.	49.9	34.1
Gen. Len.	51.6	46.1
Mean Log. Prob.	60.3	60.4
CodeT5-base ft.	49.1	52.3
Attn. Features	70.6	63.3
Magicoder-S-DS-6.7B		
Prompt Len.	48.1	56.3
Gen. Len.	54.8	75.2
Mean Log. Prob.	63.7	74.9
CodeT5-base ft.	49.3	76.0
Attn. Features	73.5	78.4

Table 5: Transferability of code hallucination detectors. Each classifier was trained on HumanEval (HE) dataset and tested on MBPP dataset.

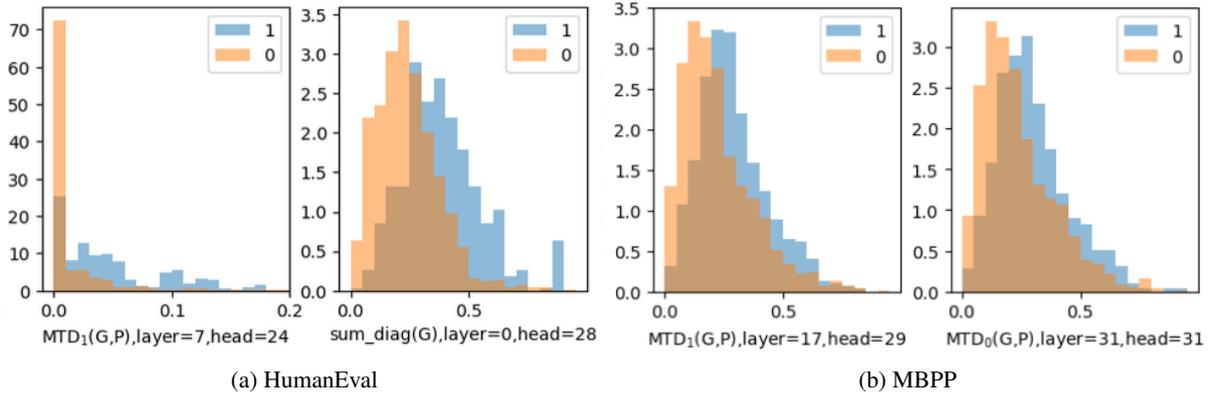


Figure 4: Distribution of classes (0-code is not correct, hallucination, 1-code is correct) vs. features from attention maps. Some of the most discriminative features are presented. Features are normalized with MinMaxScaler. CodeLlama.

Model	ROC-AUC	F1-Score
StarCoder2-7B		
Prompt Len.	52.1	45.0
Gen. Len.	52.4	38.3
Mean Log. Prob.	71.8	45.4
CodeT5-base ft.	59.1	0.0
Attn. Features	67.2	25.5
CodeLlama-7B		
Prompt Len.	53.4	42.9
Gen. Len.	50.0	41.0
Mean Log. Prob.	65.0	34.9
CodeT5-base ft.	62.4	0.0
Attn. Features	80.3	34.1
DeepSeek-Coder-6.7B		
Prompt Len.	52.2	58.2
Gen. Len.	54.0	51.3
Mean Log. Prob.	69.1	58.3
CodeT5-base ft.	55.9	57.4
Attn. Features	72.4	20.4
Qwen2.5-Coder-7B		
Prompt Len.	51.1	64.1
Gen. Len.	54.5	54.3
Mean Log. Prob.	64.7	60.8
CodeT5-base ft.	51.6	65.6
Attn. Features	64.2	54.3
Magocoder-S-DS-6.7B		
Prompt Len.	54.5	79.5
Gen. Len.	56.1	74.6
Mean Log. Prob.	69.8	76.5
CodeT5-base ft.	45.9	79.2
Attn. Features	56.9	37.6

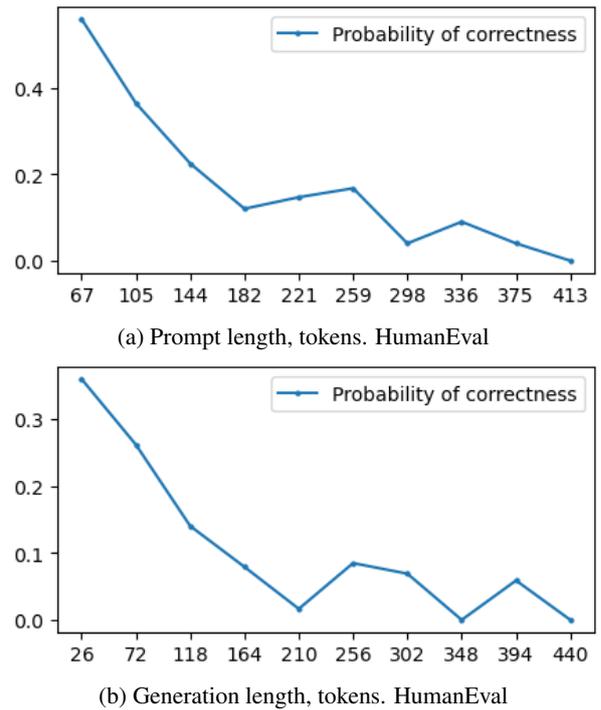


Figure 5: The individual conditional expectations for prompt and generation lengths, CodeLlama.

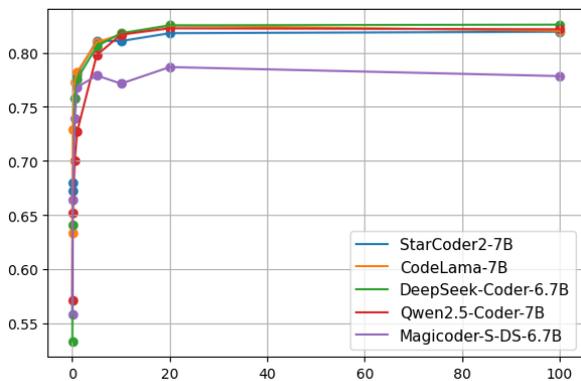
Table 6: Transferability of code hallucination detectors. Each classifier was trained on MBPP dataset and tested on HumanEval (HE) dataset.

Method	ROC-AUC	F1-Score
StarCoder2-7B		
Attn. Feat. (ours)	82.9 ± 2.7	54.2 ± 6.9
- w/o Diag. Feat.	82.2 ± 4.5	56.1 ± 9.7
- w/o MTD Feat.	83.8 ± 2.7	52.5 ± 8.4
CodeLlama-7B		
Attn. Feat. (ours)	85.6 ± 3.9	56.4 ± 7.2
- w/o Diag. Feat.	83.5 ± 4.8	50.0 ± 6.5
- w/o MTD Feat.	85.5 ± 4.4	58.3 ± 10.1
DeepSeek-Coder-6.7B		
Attn. Feat. (ours)	85.6 ± 2.8	68.9 ± 5.5
- w/o Diag. Feat.	85.1 ± 2.2	67.0 ± 5.9
- w/o MTD Feat.	84.4 ± 2.2	67.1 ± 3.8
Qwen2.5-Coder-7B		
Attn. Feat. (ours)	81.7 ± 2.8	70.2 ± 4.2
- w/o Diag. Feat.	80.6 ± 2.3	68.9 ± 3.9
- w/o MTD Feat.	78.9 ± 1.9	66.4 ± 1.4
Magocoder-S-DS-6.7B		
Attn. Feat. (ours)	82.3 ± 4.9	80.7 ± 3.6
- w/o Diag. Feat.	79.8 ± 2.7	81.1 ± 1.8
- w/o MTD Feat.	82.1 ± 3.4	81.6 ± 2.6

Table 7: HumanEval features ablation

5.6 Analyzing method’s pruning ability

In its base setup, the proposed approach requires computation of attention features from attention maps for all layers and heads. However, we explored 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 pruning ability of our approach, 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



(a) ROC-AUC vs. percentage of retained features, MBPP.

Figure 6: Pruning ability of the proposed method.

Method	ROC-AUC	F1-Score
StarCoder2-7B		
Attn. Feat. (ours)	81.9 ± 2.4	68.4 ± 5.3
- w/o Diag. Feat.	80.5 ± 2.8	66.3 ± 5.3
- w/o MTD Feat.	81.1 ± 2.6	67.7 ± 5.0
CodeLlama-7B		
Attn. Feat. (ours)	83.4 ± 2.2	64.0 ± 4.4
- w/o Diag. Feat.	81.5 ± 2.6	60.2 ± 4.2
- w/o MTD Feat.	83.5 ± 1.8	63.9 ± 4.3
DeepSeek-Coder-6.7B		
Attn. Feat. (ours)	82.6 ± 1.9	76.5 ± 2.7
- w/o Diag. Feat.	81.3 ± 2.6	74.9 ± 3.2
- w/o MTD Feat.	82.2 ± 1.7	75.9 ± 1.7
Qwen2.5-Coder-7B		
Attn. Feat. (ours)	82.2 ± 2.2	75.4 ± 1.7
- w/o Diag. Feat.	80.8 ± 2.1	75.4 ± 0.4
- w/o MTD Feat.	76.9 ± 2.2	71.6 ± 1.7
Magocoder-S-DS-6.7B		
Attn. Feat. (ours)	77.8 ± 2.5	73.4 ± 3.4
- w/o Diag. Feat.	77.8 ± 3.2	72.4 ± 2.9
- w/o MTD Feat.	78.4 ± 2.6	73.2 ± 2.1

Table 8: MBPP features ablation

trained a new XGBoost classifier on the selected set of attention features. As indicated by Figure 6, 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 hallucination detection.

6 Conclusions

In this paper, we have proposed a new approach to hallucination detection is 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.

7 Limitations

Although we have achieved good experimental results, we realize that our research have several limitations. First of all, we explored only 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't point to a specific place with a bug.

References

- Jacob Austin, Augustus Odena, Maxwell 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*.
- S. Barannikov. 1994. Framed Morse complexes and its invariants. *Adv. Soviet Math.*, 22:93–115.
- Serguei Barannikov, Ilya Trofimov, Grigorii Sotnikov, Ekaterina Trimbach, Alexander Korotin, Alexander Filippov, and Evgeny Burnaev. 2021. **Manifold topology divergence: a framework for comparing data manifolds**. In *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 7294–7305.
- Frédéric Chazal and Bertrand Michel. 2017. An introduction to topological data analysis: fundamental and practical aspects for data scientists. *CoRR*, abs/1710.04019.
- 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. 2021. Evaluating large language models trained on code. *arXiv preprint arXiv:2107.03374*.
- Tianqi Chen and Carlos Guestrin. 2016. Xgboost: A scalable tree boosting system. In *Proceedings of the 22nd acm sigkdd international conference on knowledge discovery and data mining*, pages 785–794.
- Daniil Cherniavskii, Eduard Tulchinskii, Vladislav Mikhailov, Irina Proskurina, Laida Kushnareva, Ekaterina Artemova, Serguei Barannikov, Irina Piontkovskaya, Dmitri Piontkovski, and Evgeny Burnaev. 2022. Acceptability judgements via examining the topology of attention maps. *arXiv preprint arXiv:2205.09630*.

- Kevin Clark. 2019. What does bert look at? an analysis of bert's attention. *arXiv preprint arXiv:1906.04341*.
- David Cohen-Steiner, Herbert Edelsbrunner, and John Harer. 2005. Stability of persistence diagrams. In *Proceedings of the twenty-first annual symposium on Computational geometry*, pages 263–271.
- Mohamed Elaraby, Mengyin Lu, Jacob Dunn, Xueying Zhang, Yu Wang, Shizhu Liu, Pingchuan Tian, Yuping Wang, and Yuxuan Wang. 2023. Halo: Estimation and reduction of hallucinations in open-source weak large language models. *arXiv preprint arXiv:2308.11764*.
- Shangbin Feng, Weijia Shi, Yike Wang, Wenxuan Ding, Vidhisha Balachandran, and Yulia Tsvetkov. 2024. Don't hallucinate, abstain: Identifying llm knowledge gaps via multi-llm collaboration. *arXiv preprint arXiv:2402.00367*.
- 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*.
- Phu Mon Htut, Jason Phang, Shikha Bordia, and Samuel R Bowman. 2019. Do attention heads in bert track syntactic dependencies? *arXiv preprint arXiv:1911.12246*.
- Binyuan Hui, Jian Yang, Zeyu Cui, Jiayi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. 2024. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*.
- Nan Jiang, Qi Li, Lin Tan, and Tianyi Zhang. 2024. Collu-bench: A benchmark for predicting language model hallucinations in code. *arXiv preprint arXiv:2410.09997*.
- Laida Kushnareva, Daniil Cherniavskii, Vladislav Mikhailov, Ekaterina Artemova, Serguei Barannikov, Alexander Bernstein, Irina Piontkovskaya, Dmitri Piontkovski, and Evgeny Burnaev. 2021. Artificial text detection via examining the topology of attention maps. *arXiv preprint arXiv:2109.04825*.
- Jenny T Liang, Chenyang Yang, and Brad A Myers. 2024. 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*, pages 1–13.
- Fang Liu, Yang Liu, Lin Shi, Houkun Huang, Ruifeng Wang, Zhen Yang, Li Zhang, Zhongqi Li, and Yuchi Ma. 2024. Exploring and evaluating hallucinations in llm-powered code generation. *arXiv preprint arXiv:2404.00971*.
- Anton Lozhkov, Raymond Li, Loubna Ben Allal, Federico Cassano, Joel Lamy-Poirier, Nouamane Tazi, Ao Tang, Dmytro Pykhtar, Jiawei Liu, Yuxiang Wei, et al. 2024. Starcoder 2 and the stack v2: The next generation. *arXiv preprint arXiv:2402.19173*.

597 from CodeT5-base, we used the pretrained frozen
598 CodeT5-base encoder with trainable classification
599 head consisting of 2 linear layers with hidden di-
600 mensionality 768. The classification head was
601 trained for 100 epochs with batch size 32 and learn-
602 ing rate $3e - 5$.

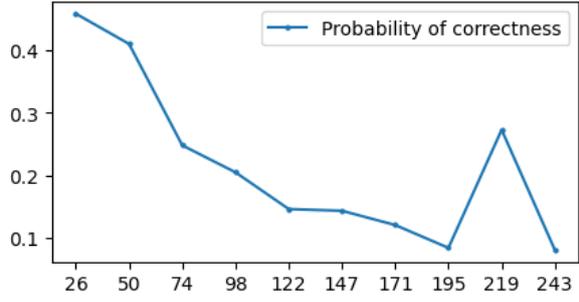
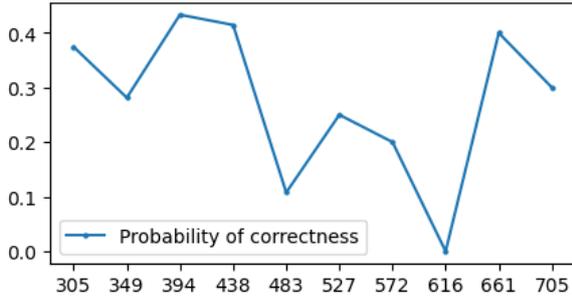
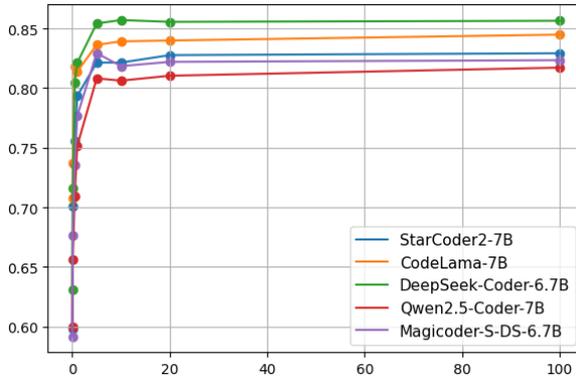
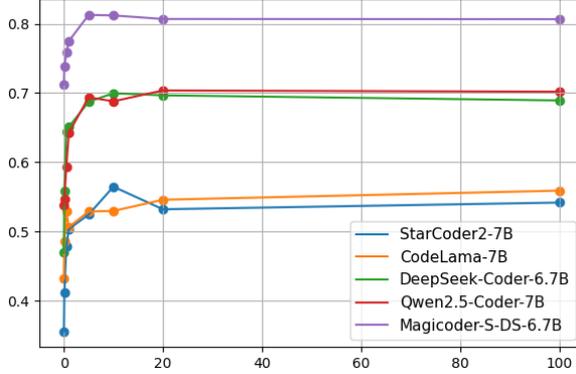


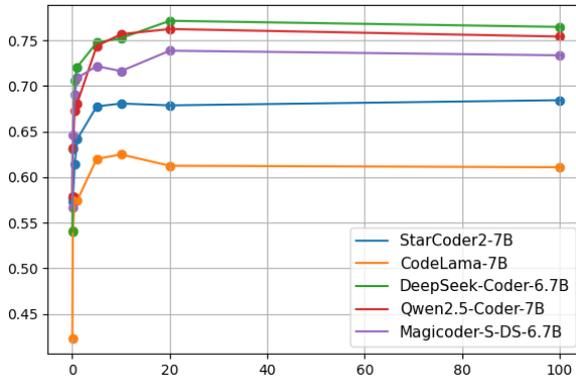
Figure 9: The individual conditional expectations for prompt and generation lengths, CodeLlama.



(a) ROC-AUC vs. percentage of retained features, HumanEval.



(b) F1, HumanEval



(c) F1, MBPP

Figure 10: Pruning ability of the proposed method.