

# Hallucinations in Code Change to Natural Language Generation: Prevalence and Evaluation of Detection Metrics

Anonymous ACL submission

## Abstract

Language models have shown strong capabilities across a wide range of tasks in software engineering, such as code generation, yet they suffer from hallucinations. While hallucinations have been studied independently in natural language and code generation, their occurrence in tasks involving code changes which have a structurally complex and context-dependent format of code remains largely unexplored. This paper presents the first comprehensive analysis of hallucinations in two critical tasks involving code change to natural language generation: commit message generation and code review comment generation. We quantify the prevalence of hallucinations in recent language models and explore a range of metric-based approaches to automatically detect them. Our findings reveal that approximately 50% of generated code reviews and 20% of generated commit messages contain hallucinations. Whilst commonly used metrics are weak detectors on their own, combining multiple metrics substantially improves performance. Notably, model confidence and feature attribution metrics effectively contribute to hallucination detection, showing promise for inference-time detection.<sup>1</sup>

## 1 Introduction

AI-based software engineering tools are becoming increasingly ubiquitous due to their potential to improve developer productivity (Jain et al., 2022; Fan et al., 2023; Hou et al., 2024). While such tools can accelerate software development, their reliance on underlying language models exposes the risk of hallucination—the phenomenon where models generate outputs that are inconsistent with their inputs or fabricate non-existent information (Ji et al., 2023; Huang et al., 2025). Such behavior may decrease developer productivity or even mislead junior developers (Ferino et al., 2025), allowing errors to propagate through to the software. Although prior

research has focused on the effects of hallucination during code generation (Liu et al., 2024; Tian et al., 2024; Agarwal et al., 2024), these effects remain largely unexplored in generation tasks involving code changes. Unlike complete code files, code changes present snippets of both the old and new versions simultaneously, which could potentially amplify hallucinations due to the model’s need to process and reason about multiple code states with partial context.

Indeed, code changes commonly used in the software engineering workflows (Tao et al., 2012; Grazia et al., 2023). Recent work also leveraged code changes as primary inputs of language models for automated software engineering tasks such as code reviews (Li et al., 2022; Lin et al., 2023). Given the increasing use of code changes in generation tasks, there is a need to understand the prevalence and effectiveness of the current detection metrics. The fragmented and context-dependent nature of code changes may increase hallucination risk and hinder detection.

In this paper, we present a comprehensive study of hallucinations in code change to natural language (CodeChange2NL) generation tasks. We focus on two key tasks: (1) automated commit message generation, which aids developers in documenting what and why code was changed, and (2) automated code review generation, which assists reviewers in identifying potential issues in code changes and suggesting improvements. To systematically analyze hallucination in CodeChange2NL, we first develop a hallucination annotation workflow specific to the CodeChange2NL context based on the outputs from task-specific models. We then empirically evaluate the effectiveness of various metric-based approaches for automatically detecting these hallucinations. In particular, we examine both reference-based metrics (which compare against human-written references) and reference-free metrics (without the references).

<sup>1</sup>All code and data will be released upon acceptance.

Our findings reveal the severity of the hallucination problem in CodeChange2NL tasks. We found that nearly 50% of model-generated code reviews and 20% of generated commit messages contain hallucinations. The three predominant categories of hallucinations are input inconsistency (where the generated NL is inconsistent with the code change), logic inconsistency (where the NL contains internally contradictory reasoning), and intention violation (where the generation fails for the specific task, e.g., it is not a review comment for code review but just a summary of the code change). Furthermore, we demonstrate that individual metrics for hallucination detection perform only marginally better than random chance (56.6% ROC-AUC for code review and 61.7% for commit messages). However, combining multiple metrics yields substantial improvements (69.1% and 75.3% respectively). Notably, reference-free metrics show promising results comparable to using all available metrics, suggesting the feasibility of detecting hallucinations without ground truth references.

This work makes three primary contributions: (1) the first systematic characterization of hallucinations in code change to natural language tasks, revealing the severity and patterns of the problem; (2) a comprehensive evaluation of automatic hallucination detection methods, demonstrating that combining multiple metrics significantly improves detection capability; and (3) identification of key reference-free metrics (model confidence and attribution scores) that effectively predict hallucinations, facilitating real-time detection in production environments without requiring reference text.

## 2 Related Work

### Hallucination in Natural Language Generation

Initially, Maynez et al. (2020) categorized hallucinations in summarization into two types: **intrinsic** hallucinations (where models misinterpret information present in the input, generating content that contradicts the source document) and **extrinsic** hallucinations (where models forge information absent from the input that cannot be verified using available information). Recently, Huang et al. (2025) identified three subcategories of intrinsic hallucinations in LLMs: instruction-inconsistent (outputs are not consistent with the instruction), logic inconsistency (output itself exhibits internal logical contradictions), and context inconsistency (outputs are not consistent with the provided input context).

Huang et al. (2025) further refined these factual hallucinations by distinguishing between factual contradiction (outputs that can be grounded but contradict real-world knowledge) and factual fabrication (outputs that are completely made up with no basis in reality or verifiable facts). Research on hallucination in code generation tasks also grounds hallucination types based on these categories (Liu et al., 2024). This taxonomy aligns closely with our CodeChange2NL tasks and serves as a foundation to determine the hallucination types in Section 3.2.

### Hallucination in Code to Natural Language Generation

Different from hallucination research in natural language to code generation, which primarily focuses on incorrect code generations e.g., dead-/unreachable code, syntactic incorrectness (Liu et al., 2024; Agarwal et al., 2024), hallucination in code to natural language generation focuses on natural language utterances that are incorrect with respect to the code/task at hand. Whilst many hallucinations in code generation can be verified by static analysis and execution (Tian et al., 2024), these solutions are not applicable for natural language outputs. Recent work examined hallucination in code-to-natural language tasks (Zhang, 2024; Maharaj et al., 2024; Kang et al., 2024). However, they primarily focus on compilable code implementations (e.g., the full body of a method). For example, Maharaj et al. (2024) studied entity-level hallucination in code summarization, where the input consists of a method-level function containing adequate contextual information. Yet, other code-to-natural language tasks involving snippets of code changes remain largely overlooked, despite their common use in real-world scenarios like commit message generation and code review (Lin et al., 2023). Moreover, due to the technical constraints of long-context modeling, snippets of code changes are often used as inputs for generation tasks instead of the complete code context (Lu et al., 2025; Berabi et al., 2024). The fragmented, context-dependent nature of code changes may increase hallucination risk and hinder detection, motivating our investigation into their prevalence and the effectiveness of existing metrics.

**Automatic Hallucination Detection** Automatic hallucination detection methods fall into two broad categories: reference-based and reference-free. *Reference-based* metrics use ground truth to gauge the quality of the generated outputs, using this quality as an estimation of hallucination. This in-

cludes lexical overlap such as BLEU (Papineni et al., 2002), which evaluates n-gram similarity between generated and reference texts. This is widely used in both Code2NL and NL2NL tasks (Liu et al., 2018a; Tufano et al., 2021; Li et al., 2022; Liu et al., 2025). More advanced metrics use Natural Language Inference (NLI): the model output is treated as a “hypothesis” to be validated against the reference. An entailment classifier labels output as entailment or contradiction, which maps to faithful or hallucinated content (Manakul et al., 2023; Elaraby et al., 2023; Hu et al., 2024; Valentin et al., 2024). *Reference-free* methods operate in many open-ended generation settings, where a reference is unavailable, by analyzing internal model behaviors and input-output relationships. One family of approaches estimates uncertainty inside models during generation (Guerreiro et al., 2023; Huang et al., 2024), with hallucinations typically exhibiting lower confidence in probability distributions and higher entropy. Another promising line is feature attribution techniques (Tang et al., 2022; Chen et al., 2025), which examine how inputs influence outputs, e.g., when a model hallucinates, its attention patterns or hidden states behave anomalously. While these metrics have been used to detect hallucinations in various NL2NL tasks, such as machine translation and question answering (Guerreiro et al., 2023; Dale et al., 2023), their capabilities in CodeChange2NL tasks remain unknown.

### 3 Study Design

#### 3.1 Research Questions

**RQ1: To what extent do task-specific language models hallucinate in code change to natural language tasks?** Prior work on hallucination in software engineering has focused on code generation, which can be verified deterministically. However, little attention has been paid to hallucinations in CodeChange2NL generation tasks, such as code review comment generation and commit message generation.

**RQ2: How effectively can existing hallucination detection methods perform on code change to natural language tasks?** While prior work in NLP have developed various methods (Dale et al., 2023; Huang et al., 2025; Ji et al., 2023) to detect hallucinations in natural language generation, their applicability to the bi-modal scenario of CodeChange2NL remains unknown. Effective detection in such contexts requires an understand-

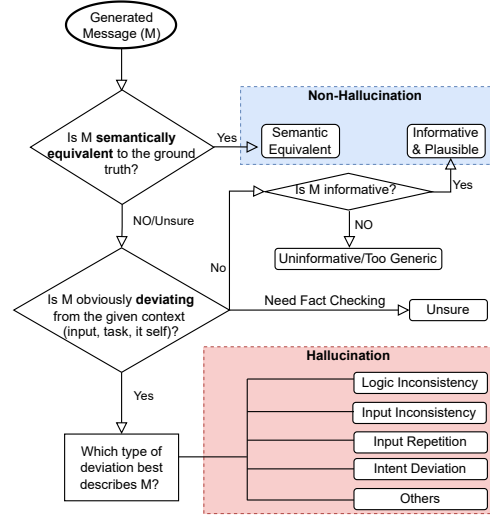


Figure 1: Hallucination Annotation Flowchart

ing of the semantics behind both code, natural language, and their interaction.

#### 3.2 Hallucination Annotation Workflow

Since no existing work addresses hallucinations in the CodeChange2NL context, we developed a decision-tree-based hallucination detection workflow by adapting taxonomies from both code generation (Liu et al., 2024) and natural language hallucination (Huang et al., 2025). Our workflow<sup>2</sup> (see Figure 1) evaluates a generated NL as follows:

**Semantic Equivalence.** We first determine whether the generated NL is semantically equivalent to the ground truth (i.e., conveying the same intent with similar framing and emphasis). If equivalent, the output is classified as non-hallucination.

**Contextual Faithfulness.** For semantically non-equivalent outputs, we assess whether the NL deviates from the context (source code, task specification, and generated text itself). Non-deviating outputs are classified as either *Informative & Plausible* (valid alternatives) or *Uninformative* (truisms).

**Hallucination Type Classification.** When context deviation exists, we categorize the hallucination into five types:<sup>3</sup> 1) *Input Inconsistency*, where the generation conflicts with the source code, e.g., pointing out a non-existent issue in code review or speculating intent that contradicts the code change in commit messages; 2) *Logic Inconsistency*, where the generation is internally illogical, independent of the input; 3) *Input Repetition*, where the generation directly copies from the input; 4) *Intent Deviation*,

<sup>2</sup>See Appendix A for definition and annotation guidelines.

<sup>3</sup>Examples are provided in Appendix A.3.

where the generation deviates from the task’s goal, e.g., not identifying issues in a code review or not explaining the code change in a commit message; and 5) *Others* for cases that are not covered by the above types. Cases requiring additional project specific fact-checking are labeled as *Unsure*.

### 3.3 Datasets and CodeChange2NL Generation

**Datasets.** We choose the widely used CodeReviewer (Li et al., 2022) dataset for code review comment generation and CommitBench (Schall et al., 2024) for commit message generation. The CodeReviewer corpus contains code diff and natural language review pairs, across 9 popular programming languages and over 1k GitHub projects. It includes 118k training, 10k validation, and 10K testing examples. CommitBench contains code diffs paired with natural language commit messages, spanning over 72k GitHub repositories and 6 programming languages. It includes 1.16 million training examples and 250k examples each for validation and testing. While related, the two tasks are different in nature—commit messages are primarily descriptive, whereas code reviews require deeper reasoning about functional correctness and potential impacts across the codebase.

**Models.** To analyze hallucination behaviors, we conduct experiments to select language models that are highly capable in both tasks. This is determined by BLEU-4 results, which is the most commonly used metric (Li et al., 2022; Schall et al., 2024). We choose two recent LLM families (Qwen2.5 and Llama3.1)<sup>4</sup> with varied model sizes for both direct prompting (7-8B, 70-72B) and task-specific fine-tuning (7-8B). We also fine-tune CCT5 (Lin et al., 2023), which is a 220M T5-based model pre-trained on 1.5M code change to commit message pairs. We used the original training data in two datasets to fine-tune the models. We found that fine-tuned models performed the best for both tasks.<sup>5</sup> Table 1 (Overall columns) presents the experimental results. Thus, we select the three fine-tuned models to generate outputs for hallucination analysis in Sections 4 and 5.

### 3.4 Hallucination Detection Methodology

We use both reference-based and reference-free hallucination detection approaches: the former for model development where the ground truth

Model	CodeReview		CommitBench	
	Overall	Sample	Overall	Sample
Llama3.1-8B	5.28	5.25	15.06	15.29
Qwen2.5-7B	5.43	5.73	15.37	15.57
CCT5	<b>5.58</b>	<b>6.53</b>	<b>17.45</b>	<b>17.46</b>

Table 1: Performance (BLEU-4 in %) of fine-tuned models on CodeReview and CommitBench benchmarks.

is available, and the latter for real-world deployment where references are unavailable. Table 2 presents a summary of the metrics we used, including two types of reference-based (BLEU-4 and NLI), and three types of reference-free (similarity, uncertainty, and feature-attribution). Uncertainty and feature-attribution metrics are calculated with either LLaMA3.1-8B-Instruct (Grattafiori et al., 2024), Qwen2.5-7B-Instruct (Yang et al., 2025) or CCT5 (Lin et al., 2023). Due to space limitations, detailed descriptions and formulas are provided in Appendix D.1. In total, 26 unique methods were considered: 2 reference-based metrics + 3 similarity scores + 3 models  $\times$  7 feature attribution and uncertainty metrics.

## 4 To what extent do task-specific language models hallucinate in CodeChange2NL tasks?

To address RQ1, we manually categorize the messages generated by the three fine-tuned models into our CodeChange2NL hallucination annotation workflow introduced in Section 3.2 to identify the presence and types of hallucinations. Using the annotated samples, we further analyze the overall prevalence of hallucinations and their distributional patterns across models and two datasets.

### 4.1 Manual Annotation

We selected the top 3 fine-tuned models (lama3.1-8B, Qwen2.5-7B, and CCT5) to generate messages in the test set. To address RQ1, we manually labeled a subset of samples that were randomly selected from the test set of each task, constituting a statistically significant sample size with a confidence level of 90% and a margin of error of  $\pm 5\%$ . This results in 264 samples for CodeReviewer comments and 268 samples for CommitBench. In total, we annotated 1,596 samples, including  $264 \times 3$  model outputs for CodeReviewer comments and  $268 \times 3$  for CommitBench messages.

Two annotators (authors of the paper) with 5+ years of experience in computer science and software engineering annotated all samples. We con-

<sup>4</sup>These were the latest models at the time of experiment.

<sup>5</sup>See Appendix B for details on prompting and fine-tuning.

Metric	Type	Description
BLEU-4	Lexical-Overlap	The n-gram overlap between the generation $y$ and reference $\hat{y}$ .
Entailment	NLI	The probability that a NLI classifier predicts $\hat{y}$ entails $y$ . We used nli-deberta-v3 <sup>6</sup> as the classifier.
Similarity	Similarity	The embedding-based cosine similarity between the generation $y$ and source code $x$ . We used three embedding models: codebert-base <sup>7</sup> , codet5p-220m-bimodal <sup>8</sup> , and codet5p-770m <sup>9</sup> .
SeqLogProb	Uncertainty	The average negative log-probability of the generated tokens in $y$ as assigned by a language model $M$ .
SeqLogit	Uncertainty	The average raw logit score (pre-Softmax) of the generated tokens in $y$ from a model $M$ .
SeqEntropy	Uncertainty	The average entropy of the generated tokens in $y$ from a model $M$ .
Source Attribution	Feature Attribution	The average of the maximum attribution scores from source tokens to each generated token in $y$ (i.e., $\frac{1}{T} \sum_{t=1}^T \max_{i \in [1, N]} A_{i,t}$ , where $A_{i,t} = x_i \times \frac{\partial y_t}{\partial x_i}$ is the importance of $x_i$ to $y_t$ from a model $M$ ). A higher score represents source contributes more strongly to $y$ .
Target Attribution	Feature Attribution	The average of the maximum attribution scores from previously generated tokens ( $y_1, \dots, y_{t-1}$ ) to each current token $y_t$ . A higher score represents the reliance on previously generated tokens.
Changed Attribution	Feature Attribution	The average of the maximum attribution scores from source tokens that are changed (in +, - lines) to each generated token in $y$ . A high score represents changed tokens contributes strongly to $y$ .
Unchanged Attribution	Feature Attribution	The average of the maximum attribution scores from source tokens that are unchanged to each generated token in $y$ . A high score represents unchanged snippets in source contributes strongly to $y$ .

Table 2: Descriptions of hallucination detection metrics, including into *reference-based* (BLEU-4 and Entailment) and *reference-free* (all others). For uncertainty and feature attribution, the model  $M \in \{\text{LLaMA3.1-8B}, \text{Qwen2.5-7B}, \text{and CCT5}\}$ . We apply both self-attribution (generator attributes its own output) and cross-attribution (external model attributes generator’s output). See Appendix D.1 for a detailed description.

Category	Type	CodeReviewer			CommitBench		
		CCT5	Llama3.1	Qwen2.5	CCT5	Llama3.1	Qwen2.5
Non-Hallucination	Semantic_Equivalent Informative	1.5	1.1	1.5	11.2	12.3	16.4
		9.5	9.8	8.7	48.1	42.5	44.4
Uninformative	Uninformative	20.1	1.5	3.8	15.7	7.1	9.7
Unsure	Unsure	22.0	41.3	43.2	5.6	16.4	15.3
Hallucination	Input_Inconsistency	26.5	23.9	24.6	17.2	19.8	13.1
	Input_Repetition	4.2	0.0	0.0	0.0	0.7	0.7
	Intent_Deviation	0.8	17.4	15.9	0.4	0.4	0.0
	Logic_Inconsistency	14.0	4.5	1.9	1.9	0.7	0.4
	Others	1.5	0.4	0.4	0.0	0.0	0.0
Total Hallucination		47.0	46.2	42.8	19.5	21.6	14.2

Table 3: The distribution (percentage) of hallucination categories and types for annotated samples. The Category column is the high-level category in Figure 1. The “Total Hallucination” is the sum of the four hallucination types.

ducted two pilot rounds (150 samples each) to refine the taxonomy and guidelines. Cohen’s  $\kappa$  improved from 0.36/0.30 (CodeReviewer/CommitBench) in the first round to 0.56/0.38 in the second. Final disagreements were resolved through discussion, achieving near-perfect agreement ( $\kappa = 0.98 / 0.96$ ). The annotators then divided the remaining samples (half-half), cross-examining each other’s work to ensure consistent labeling.

## 4.2 Hallucination Prevalence and Patterns

Table 3 shows that hallucination rates vary significantly across tasks. For the code review task, all models exhibit high hallucination rates ranging from 42.8% to 47.0%. Surprisingly, although CCT5 achieves the highest BLEU score on the CodeReviewer dataset among the three models (Table 8), it also exhibits the highest hallucination rate at 47.0%. This highlights the risk of hallucinations

even in models with strong BLEU performance. On the other hand, the commit message generation task has a lower hallucination rate than code review (14.2% to 21.6%), where Qwen2.5 has the lowest rate at 14.2%. This may be because code review is more challenging than commit message generation, as it requires identifying problems and providing specific feedback beyond what is directly observable in the code changes. Such added complexity might lead to increased hallucination behavior.

The overall distribution of hallucination types varies between tasks. Notably, the Input Inconsistency emerges as the dominant hallucination type for both tasks. This suggests that models frequently generate messages that contradict or misrepresent the actual code changes. One frequent issue in code review is that the generated messages tend to fabricate non-existent code tokens. For example, CCT5 suggests “*I think this should be orderPath instead*

of orderPathKey”. However, orderPathKey does not appear in the code change:<sup>10</sup> `+~public static final String ORDER_PATH = "orderPath";` This suggests that the model does not fully understand the meaning of newly introduced code. In the commit message task, models also often misunderstand the code changes. For example, the generated message “nomad: fix peers.json recovery for protocol version 3” misrepresents the change, which actually adds support for Nomad versions **below 3**, as indicated by the code line `+ if s.config.RaftConfig.ProtocolVersion < 3 {`.<sup>11</sup>

Intent deviation and logic inconsistency appear as another two pronounced hallucination types in the code review task, but they are rare in the commit message generation, suggesting that commit message generation models generate messages that better align with the task and suffer less logic inconsistency. Interestingly, we observe many cases where the generated review comment reads more like a commit message—for example, “*This is a temporary fix.*”, which describes the code change rather than providing a review.

Different models exhibit different type of hallucinations. CCT5, which is the specialized fine-tuned model demonstrates higher logic inconsistencies (14.0% in CodeReviewer) but significantly lower intent deviation (0.8%) than general-purpose LLMs. On the other hand, larger models (Llama3.1, Qwen2.5) frequently have intent deviation ( $\geq 15.9\%$  average) but fewer logic inconsistencies ( $\leq 4.5\%$ ). This pattern likely reflects the difference between specialized and general-purpose pretraining. Despite fine-tuning, general models retain broad task knowledge from pretraining, which can lead them to apply reasoning patterns from unrelated tasks—resulting in higher intent deviation.

## 5 How well do existing metrics detect hallucinations in CodeChange2NL tasks?

RQ1 showed that models often exhibit hallucinations and misinterpretations of code changes. In RQ2, we examine how effective automated approaches are at detecting these hallucinations in code review and commit message generation. Using our manually annotated dataset, we evaluate both reference-based and reference-free metrics described in Section 3.4. Our goal is to assess how well existing metrics detect hallucinations in Code-

to-NL tasks, particularly for code changes. We evaluate both individual metrics and combinations of complementary ones to determine whether they can approximate human judgment.

We use ROC-AUC to evaluate the hallucination detection capability of each metric. The positive class is the hallucination samples that we annotated. The negative class is the non-Hallucination samples. A ROC-AUC score of 1 indicates perfect discrimination between hallucinated and non-hallucinated cases, while a score of 0.5 suggests no discriminatory power equivalent to random guessing. For individual metrics, we calculate the ROC-AUC to assess discrimination power.<sup>12</sup> To combine metrics, we use logistic regression and evaluate its performance using accuracy and ROC-AUC.

### 5.1 How do individual metrics perform in detecting hallucinations?

*Metric Effectiveness.* Based on the the generator-agnostic results, the current metrics achieve modest ROC-AUC scores ranging from 0.538–0.566 on CodeReviewer and 0.562–0.617 on CommitBench (see Figures 2 and 3). Based on the generator-specific results, hallucinations in CCT5 are more detectable on the CodeReviewer dataset (ROC-AUC 0.65-0.71), while hallucinations in Llama3.1 are most detectable on the CommitBench dataset (ROC-AUC 0.62-0.68). This suggests that the effectiveness on hallucination detection of the metrics may vary across generation models and datasets.

Table 4 shows the metrics with the highest ROC-AUC scores in each studied dataset. In addition, we observe that on CodeReviewer, uncertainty-based metrics (logit and entropy) perform best, while embedding similarity and reference-based metrics are best on CommitBench. Nonetheless, the ROC-AUC scores suggest the limited effectiveness of current metrics on hallucination detection, which are slightly better than random guessing, highlighting the challenges of automated hallucination detection in these tasks.

*Metric Complementarity.* Different metrics may capture distinct aspects of hallucinations, potentially flagging different instances. To assess this, we selected the three highest-performing metrics based on ROC-AUC and examined their top 25% ranked samples (see the analysis details in Appendix D.2). Figure 4 shows small overlap in the

<sup>10</sup>The full code context is provided in Appendix C.1.

<sup>11</sup>The code patch is provided in Appendix C.2.

<sup>12</sup>The point-biserial correlation confirms a similar trend between metric scores and hallucination labels. Detailed results are provided in Appendix D.3.

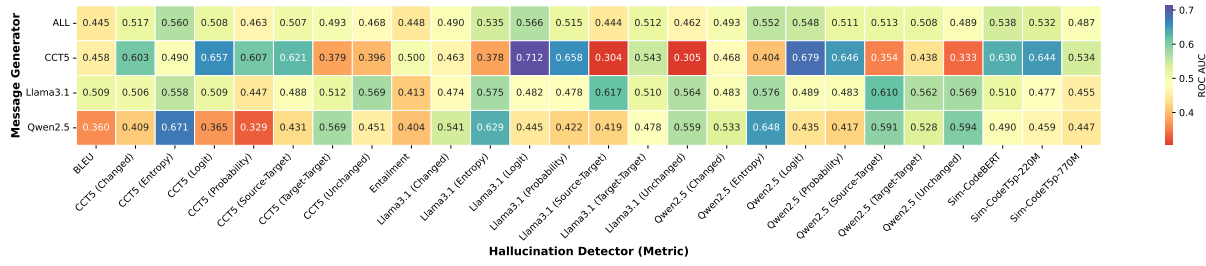


Figure 2: ROC-AUC Scores of Metrics for Hallucination Detection Across Generators on CodeReviewer. The ALL row represents the generator-agnostic result, using all outputs from CCT5, Llama3.1, and Qwen2.5. The remaining rows show performance in the generator-specific result, based on outputs from each model individually.

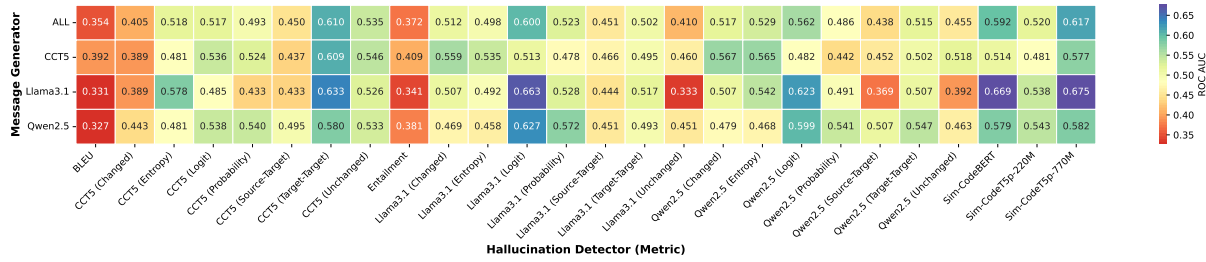


Figure 3: ROC-AUC Scores of Metrics for Hallucination Detection Across Generators on CommitBench.

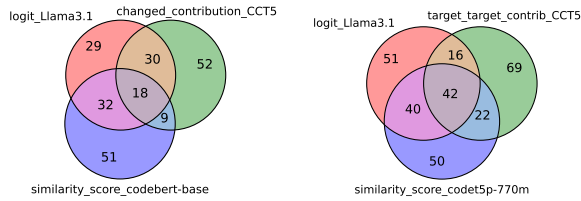


Figure 4: Top 3 individual metrics complement to each other on CodeReviewer (left) and CommitBench (right)

Type	CodeReviewer		CommitBench	
	Acc	AUC	Acc	AUC
Top Performing Individual Metrics				
logit_Llama3.1	-	0.57	-	0.60
Sim-CodeT5p-770M	-	0.48	-	0.62
Sim-Codebase	-	0.54	-	0.59
changed_contrib_CCT5	-	0.52	-	0.41
target_target_contrib_CCT5	-	0.49	-	0.61
Multiple Metrics on Logistic Regression				
Reference-based	81.6	0.59	76.0	0.68
Reference-free	81.6	0.66	78.9	0.75
ALL	<b>82.7</b>	<b>0.69</b>	<b>77.8</b>	<b>0.75</b>

Table 4: Logic regression results (Acc (%) and AUC) on hallucination prediction using multiple metrics.

top 25% samples ranked by these three metrics, indicating these metrics flag different instances as hallucinated. This highlights the potential complementarity between metrics.

## 5.2 Can combining multiple metrics enhance the accuracy of hallucination detection?

The results in section 5.1 highlight the potential complementarity between metrics. Thus, we explore whether combining them can improve performance. Prior work (Snyder et al., 2024) also shows that combining multiple signals improves hallucination detection in question-answering tasks. To analyze the discrimination power of combined metrics for hallucination detection, we use a logistic regression model fitted to our annotated samples. For each generation task, we combine all samples from the three models, resulting in 440 samples for CodeReviewer and 717 samples for CommitBench.

To understand the capability of different types of

metrics, we build three logistic regression models using: 1) all metrics, 2) reference-based metrics only, and 3) reference-free metrics only. Since some metrics may capture similar signals or redundant, leading to multicollinearity and overfitting, we use the Akaike Information Criterion (AIC) (Akaike, 1974) to identify metrics that meaningfully contribute to the prediction. Then, we use the selected metrics as features to fit the logistic regression model and analyze the coefficients to identify which metrics are most important for hallucination detection.

Table 4 shows the logistic regression results. Combining multiple metrics substantially improves ROC-AUC scores for hallucination detection on both datasets, compared to individual metrics alone.

Type	Metric	Coeff	Sign
Uncertainty	logit_Llama3.1	6.00*	+
Uncertainty	entropy_Qwen2.5	3.33*	+
Attribution	source_target_Qwen2.5	2.83*	+
Attribution	source_target_Llama3.1	2.78*	-
N-gram	BLEU	1.94*	-

Table 5: Top-5 important features on predicting hallucinations (vs. non-hallucinations) in CodeReviewer. \* indicates the coef is significant ( $p < 0.05$ ).

Type	Metric	Coeff	Sign
Uncertainty	logit_Llama3.1	6.86*	+
Uncertainty	logit_Qwen2.5	5.93*	-
Attribution	changed_CCT5	4.71*	-
N-gram	BLEU	3.49*	-
Similarity	similarity_score_codebert	2.41*	+

Table 6: Top-5 importance features on predicting hallucinations (vs. non-hallucinations) in CommitBench.

For CodeReviewer, the ROC-AUC increased from the best individual score of 0.57 (logit\_Llama3.1) to 0.69 when using all metrics. For CommitBench, it improved from 0.62 (similarity\_score\_codet5p-770m) to 0.75. Surprisingly, using reference-free metrics alone achieved ROC-AUC scores close to that of using all metrics. In contrast, reference-based metrics achieved lower performance, possibly because they are fewer in number or inherently less predictive. This highlights a potential benefit of hallucination detections in these CodeChange2NL tasks without ground-truth.

Tables 5 and 6 present the most important features along with their coefficients. The SeqLogit calculated with Llama3.1 (Logit\_Llama3.1) emerges as the most important feature for both tasks. Uncertainty metrics from Llama3.1 and Qwen2.5 consistently appear among the top features, demonstrating strong predictive power. For CommitBench dataset, the  $-$  coefficient in Qwen2.5 aligns with prior findings, i.e., when not hallucinating, a model is more confident (Dale et al., 2023). On the other hand, the  $+$  coefficient in Llama3.1 could be due to its overconfidence as the distribution of logit in Llama3.1 is skewed towards high scores in CommitBench (See Appendix D.4). Feature attribution metrics rank next in predictive strength, indicating that hallucinations can be detected by analyzing how models utilize source code during generation.

Figure 5 presents an example generated for code review.<sup>13</sup> The generated review suggests passing a parameter that is already being passed in both

<sup>13</sup>See an example of commit message in Appendix Figure 8.

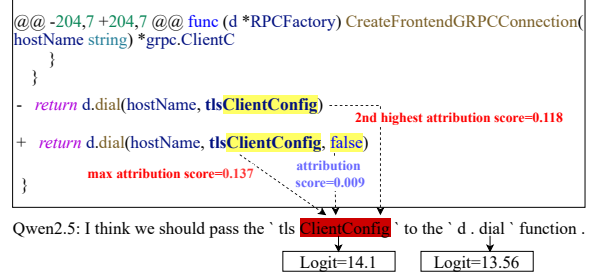


Figure 5: An example of feature attribution on a hallucinated code review comment generated by Qwen2.5. Attribution model: Llama3.1.

old and new code, while ignoring the actual code change. This hallucinated generation has high logit and high attribution from source code. Particularly, the generated tokens appearing in the input context have high confidence based on elevated logit values. For example, based on uncertainty calculated with Llama3.1, particularly API method names like `tlsClientConfig` and `dial` have logit values of 14.1 and 13.6. However, based on the attribution scores, critical changes (i.e., the addition of the “false” parameter) that should be the primary focus of the review has minimal contribution to the generation. Instead, these common tokens like `tlsClientConfig` have large attribution scores, meaning that they contributing significantly to the generation.

For non-hallucinations, we observed that the correct input in the code changes contributes significantly to the relevant generation compared to other code snippets (e.g., in the generated comment, “Why is this needed?” the “this” token was mainly contributed by the changed line of code “+ from databricks import koalas as ks”). This indicates that the balance between the contribution from changed code and unchanged code is one important cause of hallucination in code review tasks.

## 6 Conclusion

Hallucinations are prevalent in CodeChange2NL tasks, occurring in 50% of code reviews and 20% of commit messages. We identify three common types—input/logic inconsistency, and intention violation. Our findings show that individual metrics are insufficient for effective detection, while a multi-metric approach significantly improves performance, particularly combining model confidence and feature attribution.

## 7 Limitations

While our study advances the understanding of hallucination severity and automatic detection capabilities in CodeChange2NL tasks, several limitations remain.

**Dataset Size.** Despite using statistically representative samples from the test set, our annotated dataset is relatively small due to the significant effort required for manual annotation. To mitigate this limitation, we analyzed both model-specific and aggregated samples across models to increase effective sample sizes.

**Hallucination Granularity.** We primarily focused on instance-level (whole sequence) hallucination analysis to establish a foundational understanding of the phenomenon. Our feature attribution analysis showed promise for token-level hallucination detection, revealing cases where generation heavily relied on unchanged code snippets while ignoring critical changes. Future work should explore finer-grained token-level hallucination analysis with appropriate annotations and develop techniques for more precisely identifying hallucinations at different levels of granularity.

**Model Recency and Coverage.** Due to cost constraints, we excluded commercial models (e.g., GPT-4o, Claude 3.7) from our analysis and focused on the latest open-source language models available at the time of our experiments. However, the landscape is evolving rapidly, with newer models such as LLaMA 4 and Qwen2.5-Coder emerging since our evaluation. As a result, our findings may not fully generalize to these newer or commercial models, or to different model families such as Gemini, which could exhibit different hallucination patterns in Code2NL tasks. Also, our study focuses on the hallucination in task-specific fine-tuned models since they perform better than zero-shot prompting. The hallucination prevalence in zero-shot prompting may be different. Our work lays the foundation for future research in this space, highlighting the need for ongoing evaluation as models continue to evolve and diversify.

## References

Vibhor Agarwal, Yulong Pei, Salwa Alamir, and Xiaomo Liu. 2024. Codemirage: Hallucinations in code generated by large language models. *arXiv preprint arXiv:2408.08333*.

- H. Akaike. 1974. [A new look at the statistical model identification](#). *IEEE Transactions on Automatic Control*, 19(6):716–723.
- Berkay Berabi, Alexey Gronskey, Veselin Raychev, Gishor Sivanrupan, Victor Chibotaru, and Martin Vechev. 2024. Deepcode ai fix: Fixing security vulnerabilities with large language models. *arXiv preprint arXiv:2402.13291*.
- Yuyan Chen, Zehao Li, Shuangjie You, Zhengyu Chen, Jingwen Chang, Yi Zhang, Weinan Dai, Qingpei Guo, and Yanghua Xiao. 2025. Attributive reasoning for hallucination diagnosis of large language models. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 39, pages 23660–23668.
- David Dale, Elena Voita, Loic Barrault, and Marta R. Costa-jussà. 2023. [Detecting and mitigating hallucinations in machine translation: Model internal workings alone do well, sentence similarity Even better](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 36–50, Toronto, Canada. Association for Computational Linguistics.
- 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*.
- Angela Fan, Beliz Gokkaya, Mark Harman, Mitya Lyubarskiy, Shubho Sengupta, Shin Yoo, and Jie M. Zhang. 2023. [Large Language Models for Software Engineering: Survey and Open Problems](#). In *2023 IEEE/ACM International Conference on Software Engineering: Future of Software Engineering (ICSE-FoSE)*, pages 31–53, Los Alamitos, CA, USA. IEEE Computer Society.
- Samuel Ferino, Rashina Hoda, John Grundy, and Christoph Treude. 2025. Junior software developers’ perspectives on adopting llms for software engineering: a systematic literature review. *arXiv preprint arXiv:2503.07556*.
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Luca Di Grazia, Paul Bredl, and Michael Pradel. 2023. [Diffsearch: A scalable and precise search engine for code changes](#). *IEEE Trans. Softw. Eng.*, 49(4):2366–2380.
- Nuno M. Guerreiro, Elena Voita, and André Martins. 2023. [Looking for a needle in a haystack: A comprehensive study of hallucinations in neural machine translation](#). In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics*, pages 1059–1075, Dubrovnik, Croatia. Association for Computational Linguistics.

694	Xinyi Hou, Yanjie Zhao, Yue Liu, Zhou Yang, Kailong	Hong Yi Lin, Patanamon Thongtanunam, Christoph	751
695	Wang, Li Li, Xiapu Luo, David Lo, John Grundy,	Treude, and Wachiraphan Charoenwet. 2024. <a href="#">Im-</a>	752
696	and Haoyu Wang. 2024. <a href="#">Large language models for</a>	<a href="#">proving automated code reviews: Learning from ex-</a>	753
697	<a href="#">software engineering: A systematic literature review.</a>	<a href="#">perience.</a> In <i>Proceedings of the 21st International</i>	754
698	<i>ACM Trans. Softw. Eng. Methodol.</i> , 33(8).	<i>Conference on Mining Software Repositories</i> , MSR	755
699	Xiangkun Hu, Dongyu Ru, Lin Qiu, Qipeng Guo,	'24, page 278–283, New York, NY, USA. Association	756
700	Tianhang Zhang, Yang Xu, Yun Luo, Pengfei Liu,	for Computing Machinery.	757
701	Yue Zhang, and Zheng Zhang. 2024. <a href="#">Refchecker:</a>	Chunhua Liu, Hong Yi Lin, and Patanamon Thongta-	758
702	<a href="#">Reference-based fine-grained hallucination checker</a>	nunam. 2025. Too noisy to learn: Enhancing data	759
703	<a href="#">and benchmark for large language models.</a>	quality for code review comment generation. In <i>Pro-</i>	760
704	Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,	<i>ceedings of the 21st International Conference on Min-</i>	761
705	Zhangyin Feng, Haotian Wang, Qianglong Chen,	<i>ing Software Repositories.</i>	762
706	Weihua Peng, Xiaocheng Feng, Bing Qin, and Ting	Fang Liu, Yang Liu, Lin Shi, Houkun Huang, Ruifeng	763
707	Liu. 2025. <a href="#">A survey on hallucination in large lan-</a>	Wang, Zhen Yang, Li Zhang, Zhongqi Li, and Yuchi	764
708	<a href="#">guage models: Principles, taxonomy, challenges, and</a>	Ma. 2024. Exploring and evaluating hallucinations	765
709	<a href="#">open questions.</a> <i>ACM Trans. Inf. Syst.</i> , 43(2).	in llm-powered code generation. <i>arXiv preprint</i>	766
710	Yuheng Huang, Jiayang Song, Zhijie Wang, Shengming	<i>arXiv:2404.00971.</i>	767
711	Zhao, Huaming Chen, Felix Juefei-Xu, and Lei Ma.	Zhongxin Liu, Xin Xia, Ahmed E Hassan, David Lo,	768
712	2024. Look before you leap: An exploratory study of	Zhenchang Xing, and Xinyu Wang. 2018a. Neural-	769
713	uncertainty measurement for large language models.	machine-translation-based commit message genera-	770
714	In <i>International Conference on Software Engineering</i>	tion: how far are we? In <i>Proceedings of the 33rd</i>	771
715	(ICSE).	<i>ACM/IEEE International Conference on Automated</i>	772
716	Naman Jain, Skanda Vaidyanath, Arun Iyer, Nagarajan	<i>Software Engineering</i> , pages 373–384.	773
717	Natarajan, Suresh Parthasarathy, Sriram Rajamani,	Zhongxin Liu, Xin Xia, Ahmed E. Hassan, David Lo,	774
718	and Rahul Sharma. 2022. <a href="#">Jigsaw: large language</a>	Zhenchang Xing, and Xinyu Wang. 2018b. <a href="#">Neural-</a>	775
719	<a href="#">models meet program synthesis.</a> In <i>Proceedings of</i>	<a href="#">machine-translation-based commit message genera-</a>	776
720	<i>the 44th International Conference on Software Engi-</i>	<a href="#">tion: how far are we?</a> In <i>Proceedings of the 33rd</i>	777
721	<i>neering</i> , ICSE '22, page 1219–1231, New York, NY,	<i>ACM/IEEE International Conference on Automated</i>	778
722	USA. Association for Computing Machinery.	<i>Software Engineering</i> , ASE '18, page 373–384, New	779
723	Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan	York, NY, USA. Association for Computing Machin-	780
724	Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea	ery.	781
725	Madotto, and Pascale Fung. 2023. Survey of hal-	Junyi Lu, Lili Jiang, Xiaojia Li, Jianbing Fang, Fengjun	782
726	lucination in natural language generation. <i>ACM com-</i>	Zhang, Li Yang, and Chun Zuo. 2025. <a href="#">Towards prac-</a>	783
727	<i>puting surveys</i> , 55(12):1–38.	<a href="#">tical defect-focused automated code review.</a> In <i>Forty-</i>	784
728	Sungmin Kang, Louis Milliken, and Shin Yoo. 2024.	<i>second International Conference on Machine Learn-</i>	785
729	<a href="#">Identifying inaccurate descriptions in llm-generated</a>	<i>ing.</i>	786
730	<a href="#">code comments via test execution.</a> <i>Preprint,</i>	Kishan Maharaj, Vitobha Munigala, Srikanth G Tamil-	787
731	<i>arXiv:2406.14836.</i>	selvam, Prince Kumar, Sayandeep Sen, Palani	788
732	Jiawei Li, David Faragó, Christian Petrov, and Iftekhar	Kodeswaran, Abhijit Mishra, and Pushpak Bhat-	789
733	Ahmed. 2024. Only diff is not enough: Generating	tacharyya. 2024. Etf: An entity tracing framework	790
734	commit messages leveraging reasoning and action of	for hallucination detection in code summaries. <i>arXiv</i>	791
735	large language model. <i>Proceedings of the ACM on</i>	<i>preprint arXiv:2410.14748.</i>	792
736	<i>Software Engineering</i> , 1(FSE):745–766.	Potsawee Manakul, Adian Liusie, and Mark Gales. 2023.	793
737	Zhiyu Li, Shuai Lu, Daya Guo, Nan Duan, Shailesh	<a href="#">SelfCheckGPT: Zero-resource black-box hallucina-</a>	794
738	Jannu, Grant Jenks, Deep Majumder, Jared Green,	<a href="#">tion detection for generative large language models.</a>	795
739	Alexey Svyatkovskiy, Shengyu Fu, and Neel Sun-	In <i>Proceedings of the 2023 Conference on Empiri-</i>	796
740	daresan. 2022. Automating code review activities	<i>cal Methods in Natural Language Processing</i> , pages	797
741	by large-scale pre-training. In <i>Proceedings of ES-</i>	9004–9017, Singapore. Association for Computa-	798
742	<i>EC/FSE</i> , page 1035–1047.	tional Linguistics.	799
743	Bo Lin, Shangwen Wang, Zhongxin Liu, Yepang Liu,	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and	800
744	Xin Xia, and Xiaoguang Mao. 2023. <a href="#">Cct5: A code-</a>	Ryan McDonald. 2020. <a href="#">On faithfulness and factu-</a>	801
745	<a href="#">change-oriented pre-trained model.</a> In <i>Proceedings</i>	<a href="#">ality in abstractive summarization.</a> In <i>Proceedings</i>	802
746	<i>of the 31st ACM Joint European Software Engineer-</i>	<i>of the 58th Annual Meeting of the Association for</i>	803
747	<i>ing Conference and Symposium on the Foundations</i>	<i>Computational Linguistics</i> , pages 1906–1919, On-	804
748	<i>of Software Engineering</i> , ESEC/FSE 2023, page	line. Association for Computational Linguistics.	805
749	1509–1521, New York, NY, USA. Association for	Sabrina J. Mielke, Arthur Szlam, Emily Dinan, and Y-	806
750	Computing Machinery.	Lan Boureau. 2022. <a href="#">Reducing conversational agents'</a>	807

808	<a href="#">overconfidence through linguistic calibration</a> . <i>Transactions of the Association for Computational Linguistics</i> , 10:857–872.	
809		
810		
811	Kishore Papineni, Salim Roukos, Todd Ward, and Wei-	
812	Jing Zhu. 2002. Bleu: a method for automatic evalu-	
813	ation of machine translation. In <i>Proceedings of ACL</i> ,	
814	pages 311–318.	
815	Gabriele Sarti, Nils Feldhus, Ludwig Sickert, and Os-	
816	kar van der Wal. 2023. <a href="#">Inseq: An interpretability</a>	
817	<a href="#">toolkit for sequence generation models</a> . In <i>Proceed-</i>	
818	<i>ings of the 61st Annual Meeting of the Association</i>	
819	<i>for Computational Linguistics (Volume 3: System</i>	
820	<i>Demonstrations)</i> , pages 421–435, Toronto, Canada.	
821	Association for Computational Linguistics.	
822	Maximilian Schall, Tamara Czinczoll, and Gerard De	
823	Melo. 2024. <a href="#">Commitbench: A benchmark for com-</a>	
824	<a href="#">mit message generation</a> . In <i>2024 IEEE Interna-</i>	
825	<i>tional Conference on Software Analysis, Evolution</i>	
826	<i>and Reengineering (SANER)</i> , pages 728–739, Pots-	
827	dam, Germany. IEEE.	
828	Avanti Shrikumar, Peyton Greenside, and Anshul Kun-	
829	daje. 2017. <a href="#">Learning important features through</a>	
830	<a href="#">propagating activation differences</a> . In <i>Proceedings of</i>	
831	<i>the 34th International Conference on Machine Learn-</i>	
832	<i>ing</i> , volume 70 of <i>Proceedings of Machine Learning</i>	
833	<i>Research</i> , pages 3145–3153. PMLR.	
834	Ben Snyder, Marius Moisescu, and Muhammad Bilal	
835	Zafar. 2024. <a href="#">On early detection of hallucinations in</a>	
836	<a href="#">factual question answering</a> . In <i>Proceedings of the</i>	
837	<i>30th ACM SIGKDD Conference on Knowledge Dis-</i>	
838	<i>covery and Data Mining</i> , KDD ’24, page 2721–2732,	
839	New York, NY, USA. Association for Computing	
840	Machinery.	
841	Joël Tang, Marina Fomicheva, and Lucia Specia. 2022.	
842	Reducing hallucinations in neural machine trans-	
843	lation with feature attribution. <i>arXiv preprint</i>	
844	<i>arXiv:2211.09878</i> .	
845	Yida Tao, Yingnong Dang, Tao Xie, Dongmei Zhang,	
846	and Sunghun Kim. 2012. <a href="#">How do software engineers</a>	
847	<a href="#">understand code changes? an exploratory study in</a>	
848	<a href="#">industry</a> . In <i>Proceedings of the ACM SIGSOFT 20th</i>	
849	<i>International Symposium on the Foundations of Soft-</i>	
850	<i>ware Engineering</i> , FSE ’12, New York, NY, USA.	
851	Association for Computing Machinery.	
852	Yingchen Tian, Yuxia Zhang, Klaas-Jan Stol, Lin Jiang,	
853	and Hui Liu. 2022. What makes a good commit	
854	message? In <i>Proceedings of the 44th International</i>	
855	<i>Conference on Software Engineering</i> , pages 2389–	
856	2401.	
857	Yuchen Tian, Weixiang Yan, Qian Yang, Xuandong	
858	Zhao, Qian Chen, Wen Wang, Ziyang Luo, Lei Ma,	
859	and Dawn Song. 2024. Codehalu: Investigating code	
860	hallucinations in llms via execution-based verifica-	
861	tion. <i>arXiv preprint arXiv:2405.00253</i> .	
	Rosalia Tufano, Luca Pascarella, Michele Tufano,	862
	Denys Poshyvanyk, and Gabriele Bavota. 2021. To-	863
	wards automating code review activities. In <i>Proceed-</i>	864
	<i>ings of ICSE</i> , pages 163–174.	865
	Simon Valentin, Jinmiao Fu, Gianluca Detommaso,	866
	Shaoyuan Xu, Giovanni Zappella, and Bryan Wang.	867
	2024. Cost-effective hallucination detection for llms.	868
	In <i>KDD 2024 GenAI Evaluation Workshop</i> .	869
	Lanxin Yang, Jinwei Xu, Yifan Zhang, He Zhang, and	870
	Alberto Bacchelli. 2023. <a href="#">Evacrc: Evaluating code</a>	871
	<a href="#">review comments</a> . In <i>Proceedings of the 31st ACM</i>	872
	<i>Joint European Software Engineering Conference</i>	873
	<i>and Symposium on the Foundations of Software Engi-</i>	874
	<i>neering</i> , ESEC/FSE 2023, page 275–287, New York,	875
	NY, USA. Association for Computing Machinery.	876
	Qwen An Yang, Baosong Yang, Beichen Zhang,	877
	Binyuan Hui, Bo Zheng, Bowen Yu, Chengyuan Li,	878
	Dayiheng Liu, Fei Huang, Haoran Wei, Huan Lin,	879
	Jian Yang, Jianhong Tu, Jianwei Zhang, Jianxin Yang,	880
	Jiaxi Yang, Jingren Zhou, Junyang Lin, Kai Dang,	881
	Keming Lu, Keqin Bao, Kexin Yang, Le Yu, Mei Li,	882
	Mingfeng Xue, Pei Zhang, Qin Zhu, Rui Men, Runji	883
	Lin, Tianhao Li, Tianyi Tang, Tingyu Xia, Xingzhang	884
	Ren, Xuancheng Ren, Yang Fan, Yang Su, Yichang	885
	Zhang, Yu Wan, Yuqiong Liu, Zeyu Cui, Zhenru	886
	Zhang, and Zihan Qiu. 2025. <a href="#">Qwen2.5 technical</a>	887
	<a href="#">report</a> . <i>Preprint</i> , arXiv:2412.15115.	888
	Yichi Zhang. 2024. <a href="#">Detecting code comment incon-</a>	889
	<a href="#">sistencies using llm and program analysis</a> . In <i>Com-</i>	890
	<i>panion Proceedings of the 32nd ACM International</i>	891
	<i>Conference on the Foundations of Software Engineer-</i>	892
	<i>ing</i> , FSE 2024, page 683–685, New York, NY, USA.	893
	Association for Computing Machinery.	894
	Kaitlyn Zhou, Dan Jurafsky, and Tatsunori Hashimoto.	895
	2023. <a href="#">Navigating the grey area: How expressions</a>	896
	<a href="#">of uncertainty and overconfidence affect language</a>	897
	<a href="#">models</a> . In <i>Proceedings of the 2023 Conference on</i>	898
	<i>Empirical Methods in Natural Language Processing</i> ,	899
	pages 5506–5524, Singapore. Association for Com-	900
	putational Linguistics.	901
	<b>A Hallucination Annotation</b>	902
	We used the annotation workflow described in Sec-	903
	tion 3.2 to guide the process of identifying and	904
	labeling hallucinations. Detailed definitions for	905
	each node (both non-hallucination and hallucina-	906
	tion classes) are provided in Table 7.	907
	To help annotators understand the essential ele-	908
	ments of commit messages and code review com-	909
	ments, task definitions were also provided in A.1.	910
	Through initial pilot rounds and discussions among	911
	annotators, we distilled a set of rules to guide the	912
	annotation process, which is provided in A.2.	913

Type: Definition

**Semantic Equivalent (SE):** The generated message is semantically equivalent to the ground truth.

- In code review, a semantically equivalent comment should share the same intentions regarding both the issues identified and the solutions proposed as in the ground truth.
- In commit message, we should consider both the “What” and “Why” together to decide the semantic equivalence. Semantic equivalent commit messages should convey the same intents with similar framing and emphasis.

**Not\_SE\_Informative:** M is different from ground truth but it is informative for the task at hand.

- In code review, M is considered as informative if it points out a concern and/or provide suggestions for improvement.
- For commit messages, M captures some aspects of the code change but may overlook certain points compared to the ground truth. For instance, ‘Add ‘scheme’ to sys path in ok\_test/scheme.py’ indicates where the change occurs but lacks the ‘why.’ In contrast, the ground-truth message ‘Add ‘scheme’ to path to handle zip archive case’ provides (why) context on the purpose of the modification. Note (simple way): M must contain “What”, but can be incomplete or slightly different from ground truth; “Why” can be missing.

**Not\_SE\_Uninformative:** M is different from the ground truth and it doesn’t provide useful information for the task at hand.

- In code review, M is considered uninformative if it merely seeks information to understand the code design or implementation choices, presents a general question without rationale, serves as self-justification for the code change, or acts as a compliment to the code. Note (simple way): if the What (issue) is missing, then it’s not informative.
- In commit messages, vague and general wording fails to clearly communicate the specifics of the change, such as the ‘what’ (the nature of the modification) and the ‘why’ (the reason for the modification). For example, the message ‘Minor refactoring in VRaptor’ lacks detail about what parts were refactored and the intended impact of those changes, making it difficult for reviewers to understand the significance or context of the update. Note (simple way): “What” is essential, it’s uninformative if it lacks specifics of “What”.

**Unsure\_or\_Looks\_Applicable:** M appears relevant to the context but needs further fact-checking, as its factual accuracy cannot be directly verified from the given context

- In code review, this can involve M using context such as historical background, rationale beyond the given input, or the need for fact-checking the provided solution.
- In a commit message, the rationale for explaining the issue or objectives in M might need fact-checking.

**Input Inconsistency :** M conflicts with the provided input.

- In code review, this means M points out a non-existent issue or provides a solution that is already exists in the code change or violates with programming commonsense.
- In commit message, this means that M contains information that’s not included in the code change, or misinterpret code change.

**Logic Inconsistency:** M itself doesn’t make logical sense.

**Context Repetition:** M is completely or largely copied from the input.

**Intent Deviation:** M deviates with the goal of the task at hand: not providing a review in code review task or not providing a commit message that covers what is being changed and why it’s being changed.

**Others:** This is used to capture any other types that’s not covered in the above categories

Table 7: The definitions for each of the type in our annotation. M denotes the model generated message.

## A.1 Essential Elements in Code Reviews and Commit Messages

**Code Review Comments** The primary purpose of code review comments is to offer constructive feedback from reviewers to code authors, aiming to improve code quality and maintain coding standards. A review comment often covers three elements:

- What (Evaluation): A review comment should point out what is the concern or issue in the code (Yang et al., 2023).
- How (Suggestion): An ideal review comment provides suggestions for correction or prevention since code review is expected to help fix defects, improve quality, and address developers’ quality concerns (Yang et al., 2023).

- Why: Explain the reasoning behind the concern and/or the suggested improvement (Lin et al., 2024).

**Commit Messages** The primary purpose of commit messages is to provide developers (both current and future) with a summary of code changes, enabling them to understand how the code of a project has changed and why. Two elements have been shown to be essential for a commit message (Liu et al., 2018b; Tian et al., 2022).

- What (Changes): A summary of what changes were made in the code. This often includes:
  - A summary of code object change that shows the object of change, characteristics of changes, or contrast before and after. For example, “this commit removes the following deprecated prop-

947	<i>erties: * 'server.connection-timeout' *</i>	it won't overwrite, so it's Input Inconsistency.	995
948	<i>'server.use-forward-headers' [...]</i> ". An-	Base on the two labels, we choose Input In-	996
949	other example, " <i>rename HeldCertificate.Builder.issuedBy() to signedBy()</i> ".	consistency for this message.	997
950			
951	– An illustration of function. For example,	4. How to distinguish it's a review or a justifi-	998
952	<i>Rename preferred-mapper property so its</i>	cation? A review should contain the basic	999
953	<i>clear it only applies to JSON)</i>	components of issue/concern, with optional	1000
954		suggestion and explanation, while a justifica-	1001
955	– Description of implementation princi-	tion is a message aligned with the code change	1002
956	ples. For example, " <i>SslContextBuilder</i>	(no concern or suggestion, no new informa-	1003
957	<i>was using InetAddress.getByName(null)</i>	tion inside). For example, this message "This	1004
958	<i>[...] On Android, null returns IPv6 loop-</i>	is a bit of a hack, but I think it's the best we	1005
959	<i>back, which has the name 'ip6-localhost'</i>	can do for now" should be labeled as Intent	1006
	<i>"</i>	Deviation since there is no any issue or con-	1007
960		cern.	1008
961	• Why: A justification of the motivation behind	5. Cases where the model suggests changing	1009
962	the code change. This often includes describ-	back to the older version without explanation,	1010
963	ing objectives or issues, illustrating require-	we don't know whether the suggestion is bet-	1011
	ments, or implying necessity.	ter or not. If know exactly what to fact check,	1012
964		we label it Unsure (needs fact checking); oth-	1013
965	<b>A.2 Summarized rules for annotation</b>	erwise, if it's not violating the context, then	1014
	<b>Rules for Annotating Generated Code Reviews</b>	we choose NO context deviation and then de-	1015
966		cide whether it's Informative or Uninforma-	1016
967	1. Unsure → Knowledge_Overreach: a note	tive. The following message should be labeled	1017
968	of Knowledge_Overreach should be left for	as Context Deviation → No and Informative,	1018
969	cases that contain code snippets or software	because it's sensible given the code context:	1019
970	evolution (maintains, process related), we are	"I think this is a bit of a misnomer. I think	1020
971	not sure whether the generated content is true	it should be "Gets or sets JSON serialization	1021
972	or not. E.g., " <i>I think it would be better to use</i>	<i>'getId' here.</i> "	1022
973	<i>'getId' here.</i> "		
974	2. For a composite review that contains multiple	6. In cases where the review is ambiguous, it	1023
975	sentences, there might be some sentences not	might refer back to multiple places in the code	1024
976	functioning as review. As long as there is at	patch, we label it as No-context deviation if	1025
977	least one review exist, we consider it as review	it's possible to apply in at least one kinds of	1026
	(not intent deviation).	scenario. Leave a comment of "Can be inter-	1027
978		preted as another wrong way". In the example	1028
979	3. A review might have multiple sentences and	of: " <i>Layout/EmptyLinesAroundBlockBody:</i>	1029
980	each sentence has different labels, we decide	<i>Extra empty line detected at block body end.</i> ",	1030
981	the final label based on most severe one (label	where the 'block body end' can be mapped to	1031
982	hallucination types if it exists).	different places, one with an extra empty line	1032
983	For example, given this message " <i>I think</i>	and one without.	1033
984	<i>this is a bug. The 'm_indirectKernelMem'</i>		
985	<i>is a 'std::vector&lt;usm::memory&gt;'.</i> The	7. A review can apply to multiple places in	1034
986	<i>'usm_mem' is a single element of that vec-</i>	the code patch, we prioritize mapping it to	1035
987	<i>tor. So this line is going to overwrite the</i>	the code change part (-/+ lines) unless the	1036
988	<i>'m_indirectKernelMem' with a single ele-</i>	review explicitly mentions other unchanged	1037
989	<i>ment.</i> ". We have two labels: (a) we can-	code snippets. For example, in this message	1038
990	not tell that the <i>m_indirectKernelMem' is a</i>	<i>"I think this is a bit of overkill. We can just</i>	1039
991	<i>'std::vector&lt;usm::memory&gt;' or not, which is</i>	<i>use 'Fatal' and 'Warning' directly.</i> ", the 'Fa-	1040
992	<i>'Unsure' requires fact checking; and (b) we</i>	<i>tal' and 'Warning' exist in both code changed</i>	1041
993	<i>know that "So this line is going to overwrite</i>	parts and unchanged parts, but we prioritize	1042
994	<i>the 'm_indirectKernelMem' with a single el-</i>	the changed part.	1043
	<i>ement.</i> " is wrong based on the code context,		

## Rules for Annotating Commit Messages

1. A message is considered as semantically equivalent to the ground truth message if the information you can get are equal after reading both. Specifically, both “what” changed in the code and and “why” it is changed should be aligned.
2. For semantic equivalence, we don’t not over-infer the meanings, if the message doesn’t explicit mention about it then it’s not. E.g., “Added support for CircleMarker” we don’t infer the CircleMarker is a type/instance of Marker unless the code explicitly defined it.
3. For cases where we are not sure and cannot understand the message based on the given context, our prior knowledge and external web search, label it as Unsure, leave a note of “Difficult to comprehend the message”.
4. The <I> symbol comes from training data, where they mask out information referring to a different platform such as issue IDS, URLs, and version numbers. For example, the message “Bump to <I> (#<I>)” is not hallucinating, but it’s Uninformative based on the code change as it doesn’t tell specifics of what bump to <I>. This message “removed unused imports from rfc<I>” is considered informative based on the code context.

### A.3 red

Examples with Annotated Hallucination Types We use the code review task as a representative example to illustrate the distinctions between different types of hallucinations.

**Input Inconsistency** We annotate the model output code review “*I think we should use Trim() here as well.*” as an Input Inconsistency, since the corresponding code change (see below) already uses Trim(). While the comment may seem reasonable on its own, it conflicts with the actual code context and addresses a non-existent issue.

```
@@ -144,6 +144,12 @@ namespace OpenTelemetry.  
    Instrumentation.AspNetCore.Implementation  
    {  
        activity.SetTag(SemanticConventions.  
            AttributeHttpUserAgent, userAgent);  
    }  
+  
+    var xForwardedFor = request.  
    Headers["X-Forwarded-For"].FirstOrDefault();  
+    if (!string.IsNullOrEmpty(  
        xForwardedFor))
```

```
+    {  
+        activity.SetTag(  
            SemanticConventions.AttributeHttpClientIP,  
            xForwardedFor.Split(',').First().Trim());  
+    }  
}
```

**Logic Inconsistency** We annotate the following model output as Logic Inconsistency: “*I think it would be better to test that `resp.json['data']` is a `dict` and not a `dict`.*” The logic inside the review comment itself is problematic, regardless of the code input.

**Input Repetition** For the code change below, the model output ‘*input[form] polyfill, cannot load conditionally*’ is annotated as *Input Repetition*, as the output merely copies snippets from the code change.

```
@@ -15,6 +15,7 @@ return array(  
    'vendor/bootstrap.min.js',  
    'vendor/bootstrap-accessibility.min.js',  
    'vendor/validator.min.js',  
+    'vendor/form-attr-polyfill.js', // input  
    [form] polyfill, cannot load conditionally,  
    since we need all versions of IE  
    'autocomplete.js',  
    'common.js',  
    'lightbox.js',
```

### Intent Deviation

For the code change below, the model output “*This is a bug fix. The `period.fraction` was being used instead of `period.nanosec` in the `DeadlineQosPolicy` constructor.*” is annotated as *Intent Deviation*, as it describes the code change rather than providing a review comment. It fails to identify issues or offer suggestions, thus violating the task objective.

```
@@ -42,7 +42,7 @@ bool DeadlineQosPolicy::  
    addToCDRMessage(CDRMessage_t* msg)  
    bool valid = CDRMessage::addUInt16(msg, this  
        ->Pid);  
    valid &= CDRMessage::addUInt16(msg, this->  
        length); //this->length);  
    valid &= CDRMessage::addInt32(msg, period.  
        seconds);  
-    valid &= CDRMessage::addUInt32(msg, period.  
    fraction);  
+    valid &= CDRMessage::addUInt32(msg, period.  
    nanosec);  
    return valid;  
}
```

## B Prompting and Fine-tuning Models

**Zero-shot prompting** We use vLLM<sup>14</sup> for zero-shot prompting. The model temperature was set

<sup>14</sup><https://docs.vllm.ai/en/latest/>

to 0 to make the output deterministic. We used the following prompts for code review and commit message generation.

Below is a code diff submitted during a code review process.  
Please write a commit message within 50 words.  
[code\_diff]: {code\_diff}  
Respond only with valid JSON. Do not write an introduction or summary.

Below is a code diff submitted during a code review process. Please write a code review comment within 50 words to identify the concerns and suggest improvements.  
[code\_diff]: {code\_diff}  
Respond only with valid JSON. Do not write an introduction or summary.

**Fine-tuning models** We fine-tuned the three models on task-specific training data, including two general language models (Llama3.1-8B-Instruct<sup>15</sup> and Qwen2.5-7B-Instruct<sup>16</sup>) and one specialized small language model pre-trained on code and commit message generation (Lin et al., 2023). The experiment was conducted on 1 NVIDIA H100 GPU.

For CCT5 (Lin et al., 2023), we reused the code and original scripts from their replication package<sup>17</sup> to fine-tune the model on our dataset. The hyperparameters are: train\_batch\_size= 32, learning\_rate = 3e-4, max\_source\_length = 512, max\_target\_length = 128 and warmup\_steps = 500, gradient\_accumulation\_steps = 4, maximum\_train\_steps = 150000, optimizer=AdamW.

For LLaMA3.1-8B-Instruct and Qwen2.5-7B-Instruct, we perform instruction fine-tuning to further update the models parameters for the tasks at hand. We use full fine-tuning rather than parameter-efficient methods such as LoRA, as our preliminary experiments found that full fine-tuning performed better. The following instruction templates are used during training:

Below is an instruction that describes a task, paired with an input that provides further context. Write an Output that appropriately completes the request.

### Instruction: Review the code diff and provide a constructive comment highlighting any issues and suggesting improvements.

### Input:

Code diff: {code\_diff}

### Output:

{code\_review}

Below is an instruction that describes a task, paired with an input that provides further context. Write an Output that appropriately completes the request.

### Instruction: You are a programmer who makes the below code changes. Please write a commit message for the below code diff

### Input:

Code diff: {code\_diff}

### Output:

{commit\_message}

Regarding the hyperparameters used to fine-tune the two LLMs (Llama3.1-8B-Instruct and Qwen2.5-7B), we set the learning\_rate = 5e-5, max\_sequence\_length = 1024, batch\_size = 4. We set the max\_steps of fine-tuning to be 30000 and choose the best performing model on the validation set. The optimiser is Adamw.

**Results** We evaluated seven models in total, including four zero-shot and 2 fine-tuned models,<sup>18</sup> on their capability of generating task-specific messages using the traditional BLEU-4 metric (Papineni et al., 2002). Table 8 presents the experimental results on code review comment generation and commit message generation across prompting and fine-tuning approaches.

The experimental results reveal several key patterns. First, zero-shot prompting approaches consistently underperform fine-tuned models, with BLEU scores ranging from 3.88-4.70% for code review and 8.62-9.72% for commit messages. In contrast, fine-tuned models achieve substantially higher performance, with the specialized CCT5 model reaching 5.58% on code review and 17.45% on commit messages. This highlights the neces-

<sup>15</sup><https://huggingface.co/meta-llama/Llama-3.1-8B-Instruct>

<sup>16</sup><https://huggingface.co/Qwen/Qwen2.5-7B-Instruct>

<sup>17</sup><https://github.com/Ringbo/CCT5>

<sup>18</sup>We consider the fine-tuned LLMs as different models from the ones before fine-tuning, as their weights have been updated for the tasks.

Setting	Model	CodeReview		CommitBench	
		Overall	Sample	Overall	Sample
Zero-shot prompt	Llama3.1-8B-Instruct	4.22	3.28	9.21	8.89
	Qwen2.5-7B-Instruct	4.70	4.00	8.99	8.62
	Llama3.1-70B-Instruct	3.88	4.09	9.72	9.88
	Qwen2.5-72B-Instruct	4.29	4.31	8.62	8.06
Fine-tuned	Llama3.1-8B-Instruct	5.28	5.25	15.06	15.29
	Qwen2.5-7B-Instruct	5.43	5.73	15.37	15.57
	CCT5	<b>5.58</b>	<b>6.53</b>	<b>17.45</b>	<b>17.46</b>

Table 8: Performance (BLEU-4 measured in %) comparison of different models on CodeReview and CommitBench benchmarks under zero-shot and fine-tune settings.

sity of fine-tuning for generating higher-quality CodeChange2NL messages.

## C Examples of Code Changes

### C.1 Example from CodeReview

Second, code review proves to be a more challenging task compared to commit message generation, with BLEU scores approximately 2-3 times lower across all model configurations. This is sensible given that code review comments require models to critically analyze and provide constructive feedback on code changes, representing a higher cognitive demand than the descriptive nature of commit messages.

The performance on our manually sampled subset closely mirrors the overall dataset performance, with sample BLEU scores showing similar trends (e.g., CCT5 achieving 6.53% vs 5.58% overall for code review), validating the representativeness of our evaluation approach.

```
@@ -65,6 +65,7 @@ public class
    SmartStorePlugin extends ForcePlugin {
    public static final String LIKE_KEY = "
        likeKey";
    public static final String MATCH_KEY = "
        matchKey";
    public static final String SMART_SQL = "
        smartSql";
+ public static final String ORDER_PATH = "
    orderPath";
    public static final String ORDER = "order";
    public static final String PAGE_SIZE = "
        pageSize";
    public static final String QUERY_TYPE = "
        queryType";
```

### Human Code Review: Previously you could only order by the field you were searching by. But for full-text search, you can search across all indexed fields, and it didn't make sense not to have a sorting. For backward compatibility, the javascript code uses indexPath as the orderPath when no orderPath is provided.

### CCT5 Code Review: I think this should be orderPath instead of orderPathKey.

## C.2 Example from CommitBench

```
diff --git a/nomad/server.go b/nomad/server.
go
index <HASH>..<HASH> 100644
--- a/nomad/server.go
+++ b/nomad/server.go
@@ -1169,7 +1169,12 @@ func (s *Server)
    setupRaft() error {
    }
    } else if _, err := os.Stat(peersFile);
        err == nil {
        s.logger.Info("found peers.json file,
            recovering Raft configuration...")
-    configuration, err := raft.ReadPeersJSON
(peersFile)
+    var configuration raft.Configuration
+    if s.config.RaftConfig.ProtocolVersion <
3 {
+    configuration, err = raft.
ReadPeersJSON(peersFile)
    } else {
+    configuration, err = raft.
ReadConfigJSON(peersFile)
+    }
    if err != nil {
        return fmt.Errorf("recovery failed to
            parse peers.json: %v", err)
    }

### Human Commit Message: Add support in
nomad for supporting raft 3 protocol peers.json
### CCT5 Commit Message: nomad: fix
peers.json recovery for protocol version 3
```

## D Hallucination Detection

### D.1 Hallucination Detection Methodology Details

We adopt existing hallucination measurement metrics, including reference-based and reference-free hallucination detection approaches to address different practical needs. Reference-based metrics serve as valuable benchmarks during model training and evaluation when gold standards are available, while reference-free methods enable hallucination detection in real-world deployment scenarios where reference texts are typically unavailable.

#### D.1.1 Reference-based Metrics

In reference-based metrics, hallucination is estimated by the quality of a generation  $y$ , which is evaluated by comparing against the reference  $\hat{y}$  using certain metrics. The hypothesis is that the lower the quality is, the more likely  $y$  it is to be a hallucination. We use two metrics that are widely used for quality estimation: Lexical overlap with BLEU, and Natural Language Inference.

**Lexical overlap** metrics such as BLEU evaluate the n-gram overlap between the  $y$  and  $\hat{y}$ . This type of metric has been widely used in prior work to evaluate the quality of generated commit messages (Liu et al., 2018a; Li et al., 2024) and review comments (Tufano et al., 2021; Li et al., 2022). Recently, it has also been adapted to study the correlation with hallucinations in natural language generation tasks, such as machine translation (Guerreiro et al., 2023; Dale et al., 2023).

**Natural Language Inference (NLI).** NLI is a standard NLP task that evaluates the logic relationship between a pair of premise and hypothesis sentences, determining whether it is entailment, contradiction, or neutral, which has been widely used to evaluate the factual consistency (Hu et al., 2024; Valentin et al., 2024) and hallucination detection (Manakul et al., 2023; Elaraby et al., 2023). We use NLI to measure the probability of the reference  $y$  entails the the generated NL  $\hat{y}$ . The intuition is that if the  $y$  can be directly inferred from the reference  $\hat{y}$ , then it is high quality and less likely to hallucinate. We used the best performing model nli-deberta-v3<sup>19</sup> based on the performance on Sentence Transformer<sup>20</sup> to obtain the entailment logit.

#### D.1.2 Reference-free Metrics

In reference-free measurements, reference is not accessed, only information from the source input or from the model behaviors while generating a sequence is used. We use three types of measurements: similarity-based, uncertainty-based, and feature-attribution based.

**Similarity between the generation and the source** We estimate semantic similarity between source and generation using cosine similarity  $\cos(E_y, E_x)$  between embeddings of generated NL  $y$  and source code  $x$ . The intuition is that irrelevant generations are less similar and more likely to hallucinate. To obtain the embeddings, we use three models pre-trained on both code and natural language corpora: codebert-base<sup>21</sup>, codet5p-220m-bimodal<sup>22</sup>, and codet5p-770m<sup>23</sup>.

**Sequence-level confidence scores** A sequence-level confidence score has been used in machine

<sup>19</sup><https://huggingface.co/cross-encoder/nli-deberta-v3-base>

<sup>20</sup><https://sbnet.net/>

<sup>21</sup><https://huggingface.co/microsoft/codebert-base>

<sup>22</sup><https://huggingface.co/Salesforce/codet5p-220m-bimodal>

<sup>23</sup><https://huggingface.co/Salesforce/codet5p-770m>

translation for hallucination detection (Guerreiro et al., 2023; Huang et al., 2024), where it is calculated via aggregating token-level uncertainty into sentence level by taking the average across the sequence. Token-level confidence can be measured in various ways. The intuition is when a model hallucinates, it tends to be less confident. Several metrics have been proposed to estimate the token-level uncertainty, including probability, logit and entropy (Guerreiro et al., 2023; Huang et al., 2024; Valentin et al., 2024).

We also use entropy to measure uncertainty: a more uniform token distribution (higher entropy) indicates lower model certainty. This can be formulated as follows:

$$\text{SeqEntropy} = \frac{1}{L} \sum_{i=1}^L H_i, \quad (1)$$

where  $H_i$  is the entropy of the token distribution.

**Feature attribution** In a transformer-based model  $M$ , generating a token  $y_t$  involves both the input  $x$  and previously generated target tokens ( $y_1$  to  $y_{t-1}$ ). Prior work has shown that the interaction between  $y_t$  and these sources reveals hallucination patterns (Tang et al., 2022; Chen et al., 2025; Snyder et al., 2024), which can be detected through feature attribution in NL hallucinations. We conduct both feature attribution for both the input source  $x$  and the previously generated target tokens.

We employ a widely used feature attribution method Input X Gradient (Shrikumar et al., 2017), which calculates the gradient of the output with respect to the input and considers the impact of input magnitudes on generation. The attribution score from  $x_i$  to  $y_t$  can be formulated as:

$$A_{i,t} = x_i \times \frac{\partial y_t}{\partial x_i} \quad (2)$$

where  $A_{i,t}$  is the attribution score, and  $\frac{\partial y_t}{\partial x_i}$  denotes the gradient of  $y_t$  in an attribution model  $M$  with respect to the input  $x_i$ . A higher  $A_{i,t}$  indicates that  $x_i$  is more important for generating  $y_t$ .

**Source Attribution Score.** To investigate hallucinations on sequence level, we apply an aggregation function on  $A$  to convert a sequence of token-level attribution scores into a single attribution value. We first compute the maximum attribution value across all input tokens for each output token  $y_t$ , then take the average of these maximum values. The attribution score of the source to the

generated sequence.

$$\text{SourceAttr} = \frac{1}{T} \sum_{t=1}^T \max_{i \in [1, N]} A_{i,t}, \quad (3)$$

where  $T$  is the length of the generated sequence, SourceAttr represents final sequence-level overall source contribution score. The intuition is that when the maximum input contribution is small, the generated  $y$  is likely to be a hallucination as the model didn't generate based on the input.

Given our input is a code change consisting of both old and new code, human developers primarily focus on the changed parts when generating commit messages and code review comments. Based on this observation and the assumption that models should similarly emphasize code changes, we designed variations of the aggregation methods that separate attribution scores for changed and unchanged code. Our hypothesis is that lower attribution scores on the changed parts indicate a higher likelihood of hallucination.

$$\text{ChangedAttr} = \frac{1}{T} \sum_{t=1}^T \max_{i \in C} A_{i,t}, \quad (4)$$

$$\text{UnchangedAttr} = \frac{1}{T} \sum_{t=1}^T \max_{i \in [1, N] \setminus C} A_{i,t}, \quad (5)$$

where  $C \subset [1, N]$  represents the indices of tokens in the changed code (all - and + lines), and  $[1, N] \setminus C$  represents the indices of unchanged code tokens.

**Target Attribution.** We also calculate the attribution score from previously generated tokens:

$$\text{TargetAttr} = \frac{1}{T} \sum_{t=1}^T \max_{j \in 1, \dots, t-1} \hat{A}_{j,t}, \quad (6)$$

where  $\hat{A}_{j,t}$  is the attribution score from  $y_1$  to  $y_j$  ( $j$  ranges from 1 to  $t-1$ ). The final TargetAttr score denotes the overall maximum attribution score from previously generated tokens to the current token.

To obtain attribution scores for generated sequences, we use constrained attribution (Sarti et al., 2023) through the Inseq library.<sup>24</sup> Constrained attribution works by providing an attribution model  $M$  with both the input code  $x$  and the generated output  $y$ , then analyzing how the model associates each input token with each output token step by step. Rather than generating text freely, the model is

<sup>24</sup><https://inseq.org/en/latest/>

constrained to follow the specified target sequence, allowing us to measure which parts of the input most strongly influence each token in the output. This reveals the model’s implicit justification for each output token based on the input.

As the attribution model  $M$ , we use the same three models fine-tuned in our RQ1 experiments for each task: LLaMA3.1-8B-Instruct, Qwen2.5-7B-Instruct, and CCT5. For each generation, we apply both self-attribution (where the generator attributes its own output, e.g., CCT5 attributes its own generation) and cross-attribution (where a different model attributes the output, e.g., CCT5 attributes LLaMA3.1-8B’s generation). This dual perspective helps us understand whether a model is aware of its own hallucinations and whether external models can detect hallucinations based on attribution signals. While attributing each output token, we also extract uncertainty scores based on logit, probability, and entropy.

## D.2 Complementarity Among Individual Detection Metrics

To examine how different types of metrics complement each other, we select the top three individual metrics (one from each category) based on ROC-AUC.

For CodeReviewer, we choose `logit_Llama3.1`, `similarity_score_codebert-base`, and `changed_contribution_CCT5`. For CommitBench, we select `similarity_score_codet5p-770m`, `target_target_contrib_CCT5`, and `logit_Llama3.1`.

From each metric, we extract the top 25% samples ranked by their metric score, indicating that they are highly correlated with hallucination labels. We then analyze the overlaps and unions of these sets.

Figure 4 shows the Venn diagrams of the selected metrics. On CODEREVIEWER, the three metrics capture almost disjoint sets. On COMMITBENCH, only three samples are shared across all three metrics, suggesting strong complementarity.

## D.3 Correlation between Detection Metrics and Hallucination

In addition to ROC-AUC, we also analyzed the correlation between each individual metric and the hallucination labels we annotated (hallucination = 1, non-hallucination = 0). To evaluate the correlation, we use the point-biserial correlation coefficient ( $r_{pb}$ ), which measures the strength and direction of the relationship between a continuous

variable (i.e., metric scores) and a dichotomous variable (i.e., the binary hallucination label).

The results are presented in Figures 6 and 7. Overall, the correlation is weak ( $|r_{pb}| \in [0, 0.2)$ ) across all samples for individual metrics. However, when examining generator-specific results, the correlation between certain generator-metric pairs increases ( $|r_{pb}| \in [0.2, 0.3)$ ).

These findings further motivate our exploration of how combining multiple metrics can improve hallucination detection.

## D.4 Signs of Coefficients in LR model

In Section 5.2 (Table 6), we observed that the two uncertainty-based metrics—`logit_Llama3.1` and `logit_Qwen2.5`—both contribute significantly to hallucination prediction, but with opposite coefficient signs: positive for `logit_Llama3.1` and negative for `logit_Qwen2.5`. The signs of the coefficients indicate that higher logits from LLaMA3.1 are associated with hallucinations, whereas higher logits from Qwen2.5 are associated with non-hallucinations. We hypothesize that Qwen’s confidence is more reliable, while LLaMA3.1 tends to be overconfident. To further explore this, we plot the joint distribution of the two logits in Figure 9. When Qwen2.5 is more confident than LLaMA3.1 (above the diagonal), hallucinations are less frequent; conversely, when LLaMA3.1 is more confident (below the diagonal), hallucinations occur more often. This pattern supports our hypothesis.

This observation aligns with prior work (Zhou et al., 2023; Mielke et al., 2022), which shows that models can be overconfident when generating outputs due to differences in training data and strategies. In our study, both models were fine-tuned on the same data, so we suspect this difference is partly due to pre-training.

## D.5 LR Model Predictions by Hallucination Type

To understand which hallucination types are correctly detected, we examine samples predicted as hallucinations by our best logistic regression models on CodeReviewer and CommitBench (Section 5).

Figure 10 shows the type distributions. They largely mirror the overall dataset distribution, with INPUT INCONSISTENCY most frequent in both datasets, followed by INTENT DEVIATION in CODEREVIEWER, and LOGIC INCONSISTENCY thereafter.

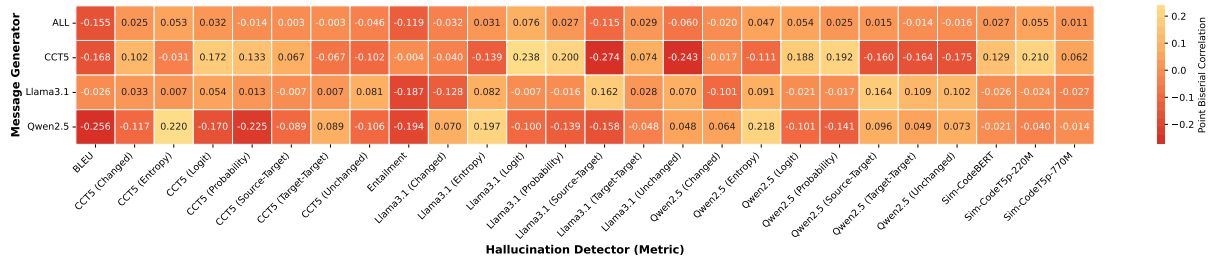


Figure 6: Point-biserial correlation between metrics and hallucinations on CodeReviewer.

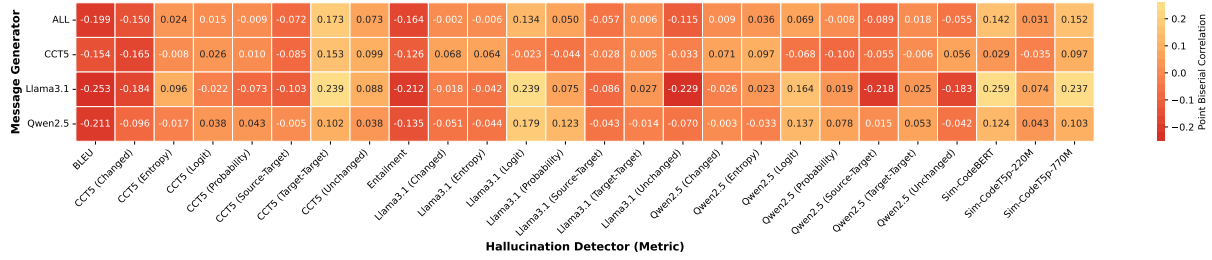


Figure 7: Point-biserial correlation between metrics and hallucinations on CommitBench.

## D.6 LR model prediction per programming language

While our hallucination detection approach is language-agnostic, model performance may still be influenced by programming language distributions in pre-training and fine-tuning data. To examine this, we analyze the distribution of programming languages among samples predicted as hallucinations by the logistic regression model and compare it to the distribution of samples labeled as hallucination in the full test set.

The results are shown in Figure 11 for CODEREVIEWER and Figure 12 for COMMITBENCH. In CODEREVIEWER, the language distribution of model predictions closely matches that of the test set, suggesting consistent detection across languages. In COMMITBENCH, the distributions also largely align, with one notable exception: JavaScript (js) is the most dominant in the test set but is not predicted (recalled) in the model’s predicted hallucinations.

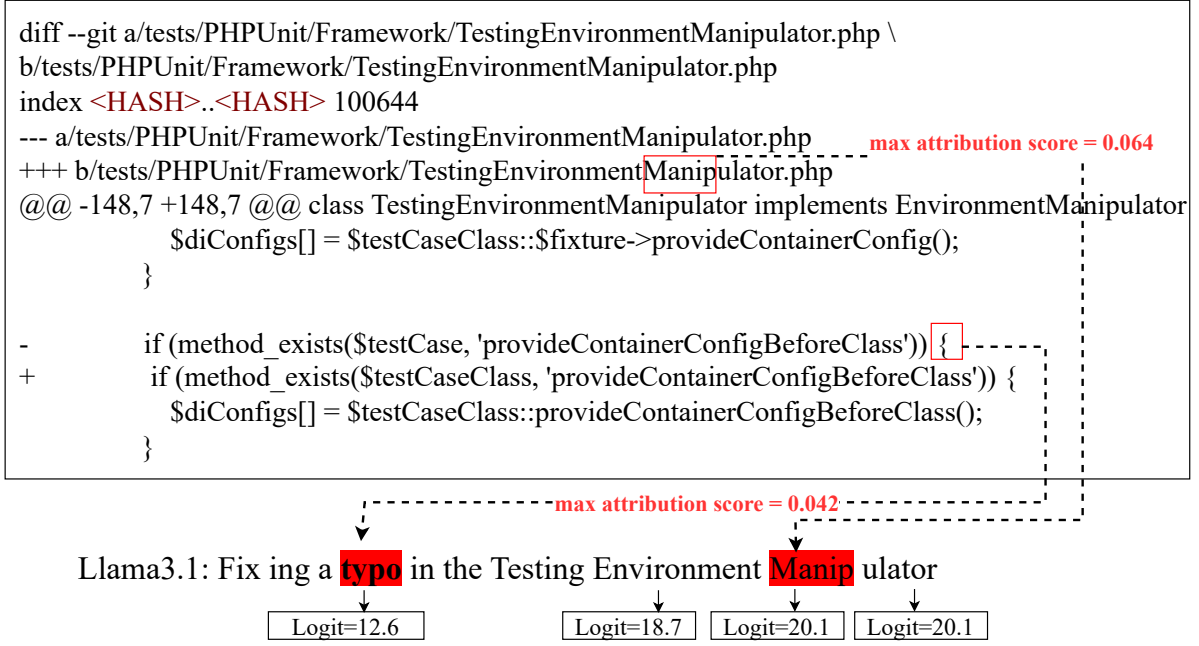


Figure 8: An example of feature attribution on a hallucinated commit message comment generated by Llama3.1. Attribution model: Llama3.1.

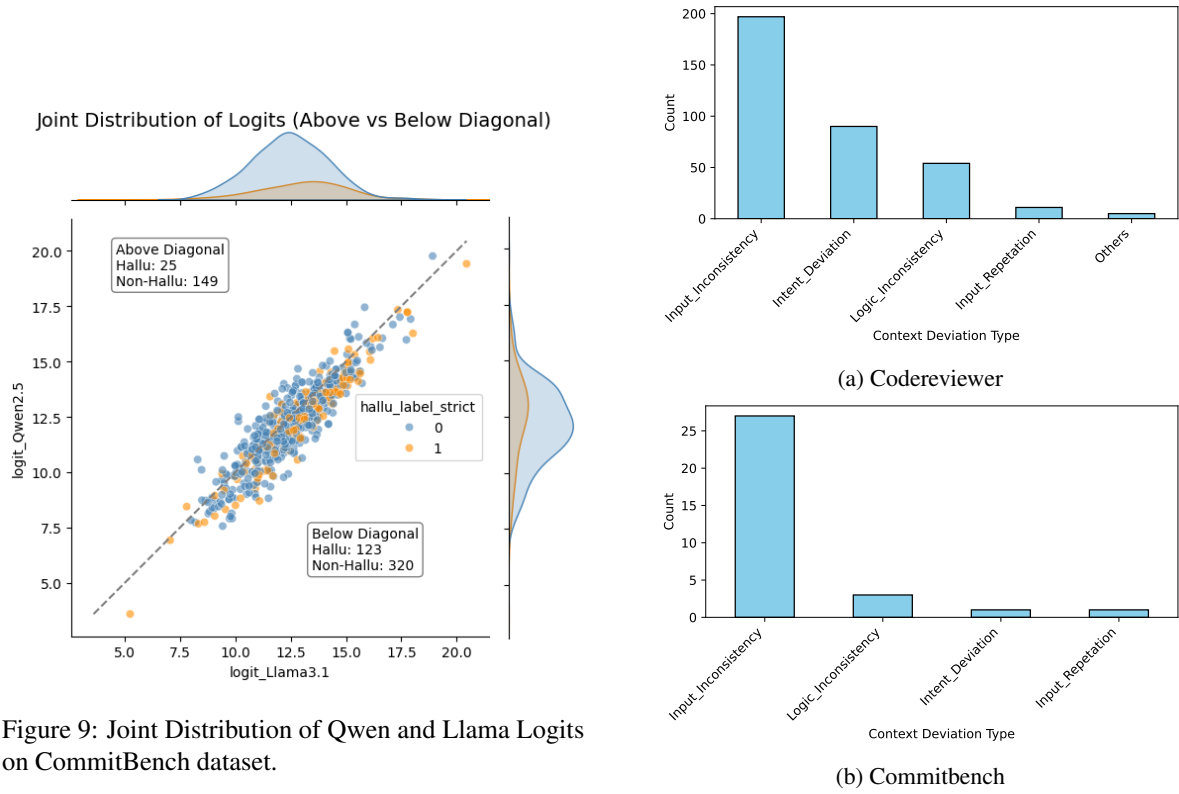
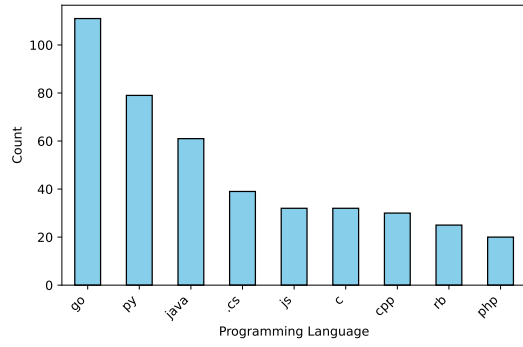
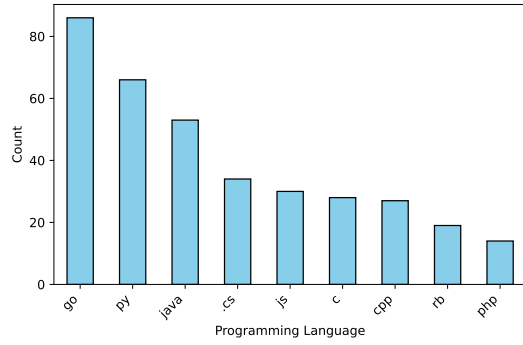


Figure 9: Joint Distribution of Qwen and Llama Logits on CommitBench dataset.

Figure 10: hallucination type distribution on LR models corrected predicted as hallucination.

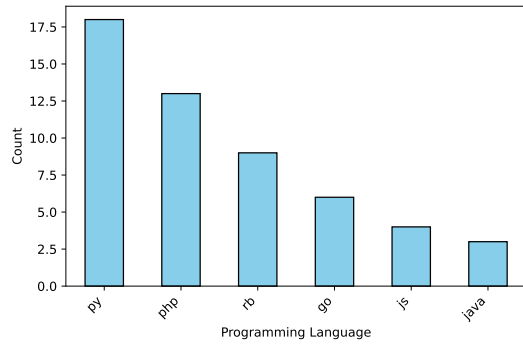


(a) Samples in model corrected predicted as hallucination

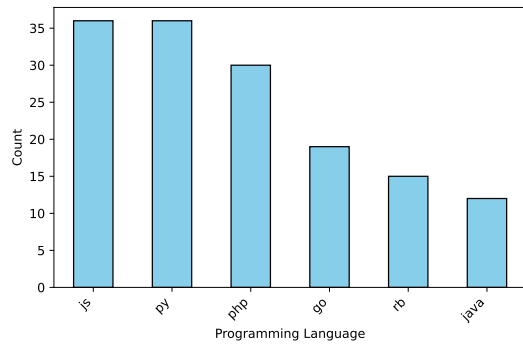


(b) Samples in test set

Figure 11: CodeReviewer: programming language distribution on the model corrected predicted as hallucinations and our test set.



(a) Samples in model corrected predicted as hallucination



(b) Samples in Test set

Figure 12: CommitBench: programming language distribution on the model corrected predicted as hallucinations and our test set.