

# ReCode: Robustness Evaluation of Code Generation Models

Shiqi Wang<sup>1,\*,\ddagger</sup> Zheng Li<sup>2,\*,\ddagger</sup> Haifeng Qian<sup>1</sup> Chenghao Yang<sup>3,\ddagger</sup> Zijian Wang<sup>1</sup>  
Mingyue Shang<sup>1</sup> Varun Kumar<sup>1</sup> Samson Tan<sup>4</sup> Baishakhi Ray<sup>1</sup> Parminder Bhatia<sup>1</sup>  
Ramesh Nallapati<sup>1</sup> Murali Krishna Ramanathan<sup>1</sup> Dan Roth<sup>1</sup> Bing Xiang<sup>1</sup>  
<sup>1</sup>AWS AI Labs <sup>2</sup>Cornell University <sup>3</sup>University of Chicago <sup>4</sup>AWS AI Research & Education  
{wshiqi,qianhf,zijwan,bxiang}@amazon.com z1634@cornell.edu

## Abstract

Code generation models have achieved impressive performance. However, they tend to be brittle as slight edits to a prompt could lead to very different generations; these robustness properties, critical for user experience when deployed in real-life applications, are not well understood. Most existing works on robustness in text or code tasks have focused on classification, while robustness in generation tasks is an uncharted area and to date there is no comprehensive benchmark for robustness in code generation. In this paper, we propose ReCode, a comprehensive robustness evaluation benchmark for code generation models. We customize over 30 transformations specifically for code on docstrings, function and variable names, code syntax, and code format. They are carefully designed to be natural in real-life coding practice, preserve the original semantic meaning, and thus provide multifaceted assessments of a model’s robustness performance. With human annotators, we verified that over 90% of the perturbed prompts do not alter the semantic meaning of the original prompt. In addition, we define robustness metrics for code generation models considering the worst-case behavior under each type of perturbation, taking advantage of the fact that executing the generated code can serve as objective evaluation. We demonstrate ReCode on SOTA models using HumanEval, MBPP, as well as function completion tasks derived from them. Interesting observations include: better robustness for CodeGen over InCoder and GPT-J; models are most sensitive to syntax perturbations; more challenging robustness evaluation on MBPP over HumanEval.

## 1 Introduction

Code generation has emerged as an important AI application. Multiple models (Nijkamp et al., 2022;

Fried et al., 2022; Wang and Komatsuzaki, 2021) have been proposed and achieved impressive performance on generating code using a natural-language description, on completing partial lines and functions, and even on solving complex coding-contest problems. They can offer real-life help to software engineers and enhance their productivity, and multiple commercial offerings exist today for AI-powered code generation (Chen et al., 2021).

However, one important aspect, robustness of the code generation models, is commonly overlooked. Anecdotally, people know that these models are sensitive to perturbations over prompts: sometimes just an extra space in a line or a slight change to a function name would lead to completely different generations, with potentially negative impacts to usability. In Fig. 1 and Fig. 2, we show two failure cases on InCoder-6B (Fried et al., 2022) and CodeGen-16B-mono (Nijkamp et al., 2022) where they perform correctly on regular prompts but fail on our perturbed ones after docstring paraphrasing and function camel case renaming in our ReCode benchmark. The perturbed prompts are natural and retain the original meaning, indicating weakness of these models if deployed in real-life applications.

There exists no comprehensive and quantitative robustness benchmark for code generation models. Li et al. (2022) includes a brief study on robustness but it has limited perturbation types and is in a setting with massive numbers of samples, unrealistic in practice. Other existing works on robustness in text or code tasks have focused on classification and are not directly applicable to code generation (Zhang et al., 2020; Jha and Reddy, 2022).

In this paper, we present **ReCode**, a **Robustness Evaluation** framework for **Code**, aiming to provide comprehensive assessment for robustness of code generation models. ReCode includes only transformations that (1) appear naturally in practice and (2) preserve the semantic meaning of the original inputs. We carefully collect and customize a com-

<sup>\ddagger</sup>Corresponding author.

<sup>\*</sup>Equal contribution.

<sup>\ddagger</sup>Work done while the authors were at Amazon.

```

def test_distinct(data):
    """
    Write a python function to determine whether all the
    numbers are different from each other are not.
    >>> test_distinct([1,5,7,9])
    True
    >>> test_distinct([2,4,5,5,7,9])
    False
    >>> test_distinct([1,2,3])
    True
    """
    return len(set(data)) == len(data)

```

Original docstring

Original completion

```

def test_distinct(data):
    """
    Write a Python function to see if all
    numbers differ from each other.
    >>> test_distinct([1,5,7,9])
    True
    >>> test_distinct([2,4,5,5,7,9])
    False
    >>> test_distinct([1,2,3])
    True
    """
    return len(set(data)) != len(data)

```

Perturbed docstring

New completion

Figure 1: InCoder-6B predicts correctly on nominal prompt (left) but fails on the prompt where docstrings are paraphrasing with BackTranslation (right). We underline the perturbed positions and wrong model completions.

```

def remove_lowercase(str1):
    """
    Write a function to remove lowercase
    substrings from a given string.
    >>> remove_lowercase("PYTHon")
    ('PYTH')
    >>> remove_lowercase("FInD")
    ('FID')
    >>> remove_lowercase("STRInG")
    ('STRG')
    """
    return "".join([i for i in str1 if i.isupper()])

```

Original Function name

Original completion

```

def removeLowercase(str1):
    """
    Write a function to remove lowercase
    substrings from a given string.
    >>> removeLowercase("PYTHon")
    ('PYTH')
    >>> removeLowercase("FInD")
    ('FID')
    >>> removeLowercase("STRInG")
    ('STRG')
    """
    str2 = str1.lower()
    return str2

```

Perturbed function name

New completion

Figure 2: CodeGen-16B-mono is correct on nominal prompt (left) but fails when function name is perturbed (right).

prehensive list of natural transformations on docstrings, function and variable names, code syntax, and code format, providing multifaceted assessments of a model’s robustness performance. We verify the quality of the perturbed data using both human evaluation and objective similarity scores. We take advantage of the fact that executing the generated code can serve as objective evaluation and define three robustness evaluation metrics that aggregate a model’s correctness across randomized transformations and transformation types. These metrics quantify a model’s accuracy on perturbed prompts, its relative accuracy drop from original prompts, as well as its general instability.

We summarize our contributions below:

- We present the first robustness evaluation benchmark ReCode for code generation tasks. Our evaluation framework is general and can be easily extended to any code generation datasets and models.<sup>1</sup>
- We collect and customize over 30 natural transformations from the aspects of docstrings, function and variable names, code syntax, and code format. Human evaluation shows that most of the perturbed prompts do not alter the semantic meaning and that their level of naturalness is close to the originals. Quantitative

similarity metrics confirm the same.

- We propose robustness evaluation metrics for code-generation tasks: Robust Pass<sub>s</sub>@k, Robust Drop<sub>s</sub>@k, and Robust Relative<sub>s</sub>@k.
- We demonstrate the ReCode benchmark on HumanEval and MBPP datasets and present extensive empirical robustness comparisons on state-of-the-art models including CodeGen, InCoder, and GPT-J across different sizes. We find that 1) diverse pretraining corpus and larger model size can help improve the model worst-case robustness, but models may learn to generalize in a non-robust way; 2) code generation models are most sensitive to syntax perturbations; 3) due to diversity, MBPP poses greater changes than HumanEval.

## 2 Related Work

**Robustness for NLP.** Recent research have identified the severe robustness problem in Large Language Models (LLMs) using adversarial examples. For example, LLMs can be easily fooled by synonym replacement (Jin et al., 2020; Zang et al., 2020). To better illustrate the severity of adversarial robustness problems for NLP models, existing works (Nie et al., 2020; Gardner et al., 2020; Kiela et al., 2021; Wang et al., 2021a) build robustness benchmarks, which encourage people to further build robust and trustworthy models. Zhang

<sup>1</sup>Code and datasets released at <https://github.com/amazon-science/recode>.

et al. (2020) presents a comprehensive overview of works in this field. Most existing works in this field focus on **classification tasks** rather than **generation tasks**. The main challenge for benchmarking robustness over generation tasks is that the evaluation of text generation is highly subjective and is usually hard to quantify. However, code generation provides a special opportunity because we can do objective and quantitative evaluation on generated codes, and code generation models use similar model architecture as NLP models.

**Robustness for code.** There are a series of previous work on different aspects of robustness problems for code. Specifically, [Bielik and Vechev \(2020\)](#) studies the adversarial robustness problem for type inference in programming languages. [Yang et al. \(2022\)](#) focuses on improving the naturalness of adversarial examples in code vulnerability prediction, clone detection and authorship attribution. [Zhou et al. \(2022\)](#) focuses on the adversarial robustness problems of source code comment generation and [\(Jha and Reddy, 2022\)](#) focuses on code translation, repair and summarization. These papers mainly focus on proposing attack and defense methods for different tasks in code domain, but there is no previous work on a comprehensive robustness benchmark for code generation domain.

**Code generation.** Code generation, also known as program synthesis, is a task of generating code based on natural language statements or code from context. Researchers have adapted transformer-based large language models to the code generation field. Various architectures have been explored: For example, CodeBERT ([Feng et al., 2020](#)), PLBART ([Ahmad et al., 2021](#)), CodeGPT ([Lu et al., 2021](#)) explore BERT, BART and GPT architectures for language models pretrained on code corpus. There are also works that propose to incorporate code structures for models to better understand the semantic information, including GraphCodeBERT ([Guo et al., 2021](#)) and CodeT5 ([Wang et al., 2021b](#)). Most recently, models with much larger size (i.e., billion-scale parameter numbers) are shown to significantly improve the performance on code generation benchmarks. Codex-12B ([Chen et al., 2021](#)) and CodeGen-16B ([Nijkamp et al., 2022](#)) are two representative very large pretrained code generation models and have established new state of the arts. However, few works have systematically explored robustness in code generation.

### 3 Methodology

In this section, we introduce the transformations to perturb prompts on both text (docstring) and code. We then propose new robust evaluation metrics.

#### 3.1 Problem Formulation

We consider the end-to-end model-based code generation task. The input prompt can include natural language statements that describe the functionality, signature of the function to generate, helper functions, and possibly a half-written function. The goal is left-to-right generation that creates or completes the function. This setting is agnostic to model architectures and is applicable to encoder-decoder or decoder-only models.

We perturb the input prompt with transformations. We focus on natural transformations that preserve the semantic meaning of the original prompt and that are likely to appear in practice, e.g., frequent typos in docstrings, tab to four spaces, function name style changes, and many more. We do not consider adversarial attacks that require model feedbacks in this paper because it is non-trivial to control the naturalness of adversarial attacks and they often require higher computational cost. Instead, we randomly generate perturbed prompts based on the restrictions for each type of perturbations and propose new metrics to evaluate model robustness based on these prompts. We leave adversarial attacks for future work.

#### 3.2 Natural Transformations on Docstrings

Docstring describes the target function to generate. Since docstrings can vary greatly when written by different users, robustness against changes in docstrings is critical for usability in applications.

For docstrings, we use the NL-Augmenter ([Dhole et al., 2021](#)) library which is designed for data augmentation and robustness evaluation on text.<sup>2</sup> We carefully select ten transformations, including character-level, word-level and sentence-level ones, that are likely to preserve semantic similarity. The selected perturbations include CharCaseChange, where random characters are replaced with their upper cases, SynonymSubstitution, where random words are substituted with their WordNet synonyms ([Miller, 1992](#)), BackTranslation, where sentences are translated to a different language (e.g., German by default) then back to English for paraphrasing the

<sup>2</sup><https://github.com/GEM-benchmark/NL-Augmenter>

Perturbations	MBPP Docstrings
Nominal	Write a function to find all words which are at least 4 characters long in a string by using regex.
BackTranslation	Write a function to find all words <b>in a string at least 4 characters long</b> using regex.
ButterFingers	Write a function to find all words which are at least 4 characters long in a string by using regex.
ChangeCharCase	Write a function to find All words which are at least 4 Characters Long in a string by using regex.
EnglishInflectionalVariation	Writes a functions to found all word which was at least 4 character long in a string by use regex.
SwapCharacters	Write a function to find all words which are at least 4 characters long in a string by using regex.
SynonymInsertion	Write a function to find <b>discover</b> all words which are at least 4 characters long in a string by using regex.
SynonymSubstitution	Write a function to find all words which <b>equal</b> at least 4 character long in a <b>chain</b> by using regex.
TenseTransformationPast	Write a function to find all words which <b>was</b> at least 4 characters long in a string by using regex.
TenseTransformationFuture	Write a function to find all words which <b>will be</b> at least 4 characters long in a string by using regex.
Whitespace	Write a function to find all words <b>w hich</b> are at least 4 characters long in a string by using regex.

Table 1: Illustrations for docstring perturbations on a MBPP sample.

whole sentence (Li and Specia, 2019; Sugiyama and Yoshinaga, 2019), and more. To perform perturbations, we extract docstring sentences from the input prompt and then put the perturbed version back to the prompt. See Appendix A for details.

We observe that directly applying NL-Augmenter to docstrings without constraints can potentially lead to low quality due to keywords in the programming languages. For example, "Create a list a[][]" could be perturbed by "Create a list [a][]" by character case swap, which is not natural. Therefore, to guarantee naturalness of perturbations, we use tree-sitter to parse the whole code snippet (the prompt & the canonical solution) to extract any existing function names, variable names ("a"), and type names ("list"). We then exclude them from being perturbed by the transformations. In Tab. 1, we list all ten transformations that are customized from NL-Augmenter and are included in our robustness benchmark along with sample illustrations.

### 3.3 Natural Transformations on Function Names

Perturbing function names also results in performance drops for code generation models. We summarize our perturbations in Tab. 2.

Some perturbations switch function names between naming conventions. For example, the perturbation called CamelCase transform function names between camel-case (e.g., "findCharLong") and snake-case ("find\_char\_long").

Other perturbations apply character-level or word-level natural text transformations on component words in a function name, including ChangeCharCase, InflectionalVariation, and SynonymSubstitution as discussed in Sect. 3.2.

### 3.4 Natural Transformations on Code Syntax

Code generation models are often used on function completion task where the prompt includes a partial

```
def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello", "l")
    "heo"
    >>> remove_Occ("abcda", "a")
    "bcd"
    >>> remove_Occ("PHP", "P")
    "H"
    """
    for i in range(len(s)):
        if s[i] == ch:
            s = s[0:i] + s[i + 1 :]
            break
```

(a) Baseline Partial Code

```
def remove_Occ(s, ch):
    # [same doc string]
    i = 0
    while i < len(s):
        if s[i] == ch:
            s = s[0:i] + s[i + 1 :]
            break
        i = i + 1
```

(b) For-While Switch

```
def remove_Occ(lines, ch):
    # [same doc string]
    for i in range(len(lines)):
        if lines[i] == ch:
            lines = lines[0:i] + lines[i + 1 :]
            break
```

(c) Variable Renaming with CodeBERT

Figure 3: An original prompt with partial code (a) and its perturbed versions (b, c).

Perturbations on Function Names	MBPP
Nominal	find_char_long
CamelCase	<b>findCharLong</b>
ButterFingers	finf_char_long
SwapCharacters	find_cahr_long
ChangeCharCase	finD_chaR_long
InflectionalVariation	<b>found_chars_long</b>
SynonymSubstitution	<b>discover_char_long</b>

Table 2: Illustrations for function name perturbations on a MBPP sample.

implementation of the target function and the goal is to complete it. In such scenarios, the partial



code in prompt is work in progress and can be subject to frequent editing, and ideally a model should be robust with respect to perturbations in the partial code. For this evaluation, we derive new customized datasets from HumanEval and MBPP by adding half<sup>3</sup> of the canonical solutions to the prompts (Fig. 3a). Then we perturb such partial code inside prompts. Details and examples for each perturbations can be found in Appendix A.

Transformations on partial code must be syntactically correct and must not alter semantic meaning. The next section will address code format, and let us first focus on code refactoring: these are syntactic changes that are semantically invariant.

We adopt three transformations from NatGen (Chakraborty et al., 2022): (1) Deadcode Insertion where dummy loops (0 iterations) or if conditions are randomly inserted; (2) Operand Swap where we randomly swap one operation (e.g.,  $a < b$  to  $b > a$ ); (3) For-While Switch where we randomly transform one for-loop structure in code to equivalent while-loop structure and vice versa.

Additionally, we implement three different schemes of variable renaming. We select the most frequent variable in the partial code and replace it using: (1) using CodeBERT (Feng et al., 2020) predictions with highest aggregated scores according to the context around all its appearance, a method inspired by (Jha and Reddy, 2022; Li et al., 2020), (2) using NatGen style renaming as "VAR\_0", and (3) random name generation with half alphabetic and half numeric characters. The first strategy tends to provide more natural variable names, yet names from the other two strategies are also plausible.

### 3.5 Natural Transformations on Code Format

A natural way to perturb partial code is by code format transformations as they preserve the original semantic meaning. We implement following code format transformations in ReCode.

**Newline Insertion:** We consider three methods of new line insertions: (1) empty lines at randomly selected positions, (2) an empty line inserted between docstring and partial code, and (3) an empty line inserted after partial code.

**Tab-Indent:** We randomly replace any space indent with tab or replace tab with 4 spaces for indent-sensitive languages like Python.

**Line Split:** We select the longest line of code and split it into two lines in the middle.

<sup>3</sup>add first  $\lfloor k/2 \rfloor$  lines given a  $k$ -line canonical solution.

**Docstrings to Comments:** We convert docstrings to comments (e.g., `""" docstring """` to `# docstring` for Python).

### 3.6 Evaluation Metrics

Many proposed transformations are randomized operations. Hence, we need to measure model robustness over multiple samples to reduce variance. Specifically, for each transformation and each prompt, we create  $s$  randomly perturbed prompts. The model under evaluation generates outputs for each of them. We measure the worst-case performance across each group of  $s$  perturbed prompts: the model is considered robust on a prompt if and only if it generates a correct solution for **all**  $s$  perturbed prompts, where correctness is measured by executing associated unit tests.

Based on such worst-case measurements, we propose three new metrics for robustness evaluation.

**Robust Pass<sub>s</sub>@k (RP<sub>s</sub>@k):** Pass@k is a widely used metric for measuring the performance of code generation tasks (Chen et al., 2021). We extend its definition to Robust Pass<sub>s</sub>@k (RP<sub>s</sub>@k) with  $s$  random perturbations. For an original prompt  $x$  and for each transformation, let the perturbed prompts be  $x_1, \dots, x_s$ . We sample  $n$  generations by the model for each prompt, and in total there are  $n \cdot s$  generations  $f_i(x_j)$ , where  $1 \leq i \leq n$  and  $1 \leq j \leq s$ . Instead of regular pass@k, we first consider the worst-case correctness across  $f_i(x_1), \dots, f_i(x_s)$  for  $1 \leq i \leq n$ : Let  $c_{i,s}(x) = 1$  if  $f_i(x_1), \dots, f_i(x_s)$  are all correct and  $c_{i,s}(x) = 0$  otherwise. Let  $rc_s(x) = \sum_{i=1}^n c_{i,s}(x)$ . Following definition of pass@k, we define the RP<sub>s</sub>@k metric as Eq. (1).

$$RP_s@k := \mathbb{E}_x \left[ 1 - \frac{\binom{n-rc_s(x)}{k}}{\binom{n}{k}} \right] \quad (1)$$

**Robust Drop<sub>s</sub>@k (RD<sub>s</sub>@k):** RP<sub>s</sub>@k directly measure worst-case robustness in absolute values. It provides a worst-case estimation for models under certain perturbation. But in some applications, users may care more about **relative performance change** to compare worst-case performance and average-case performance. We propose Robust Drop<sub>s</sub>@k defined in Eq. (2) as another important robustness metric to quantify relative changes.

$$RD_s@k := \frac{Pass@k - Robust\ Pass_s@k}{Pass@k} \quad (2)$$

**Robust Relative<sub>s</sub>@k (RR<sub>s</sub>@k):** Lastly, there are cases where models generate incorrect code on

original prompts yet predict correctly on perturbed ones. This can (arguably) be considered as non-robust behavior that we should include when reporting model robustness. Let’s first consider the case of greedy decoding with  $n = k = 1$ . Let  $RC_s^{[-]}$  denote the number of correct-to-incorrect changes under the worst-case measurement as discussed. Symmetrically, let  $RC_s^{[+]}$  denote the number of incorrect-to-correct changes under best-case measurement: if the prediction with the original prompt is incorrect yet is correct for any of the  $s$  perturbed prompts. We define the Robust Relative $_s@1$  metric as the fraction of changes in both directions out of the size of the dataset ( $N$ ):

$$RR_s@1 := \frac{RC_s^{[+]} + RC_s^{[-]}}{N} \quad (3)$$

This definition can be generalized to sampling. Let  $rc_s^{[-]}(x)$  and  $rc_s^{[+]}(x)$  be similarly defined as  $RC_s^{[-]}$  and  $RC_s^{[+]}$  except that they are the number of changes within  $n$  samples for a prompt  $x$  instead of counting across the dataset. We define

$$RR_s@k := \mathbb{E}_x \left[ 2 - \frac{\binom{n-rc_s^{[-]}(x)}{k}}{\binom{n}{k}} - \frac{\binom{n-rc_s^{[+]}(x)}{k}}{\binom{n}{k}} \right] \quad (4)$$

Eq. (4) falls back to Eq. (3) when  $n = k = 1$ .

**Discussion.**  $RP_s@k$ ,  $RD_s@k$  and  $RR_s@k$  focus on different robustness requirements in practice. High  $RP_s@k$  does not necessarily mean low  $RD_s@k$  or  $RR_s@k$ , because the model may learn to utilize spurious correlation in the datasets to demonstrate better Pass@ $k$  or RP@ $k$ , which is not robust. We advocate to report all of them to provide a comprehensive estimation of model robustness.

## 4 Evaluation

**Evaluation setup.** In this work, we use execution-based code generation benchmarks HumanEval (Chen et al., 2021) and MBPP (Austin et al., 2021) to demonstrate our ReCode robustness evaluation framework. We perform a comprehensive study of robustness evaluation on popular public models including CodeGen (Nijkamp et al., 2022), InCoder (Fried et al., 2022), and GPT-J (Wang and Komatsuzaki, 2021) to show the robustness comparisons across different model architectures and sizes. The perturbations and metrics implemented in ReCode are general and applicable to any code generation datasets and models.

### 4.1 Code Generation Robustness Evaluation

Tab. 3 and Tab. 4 show the general perturbation performances on all the models in terms of the four general perturbation categories including transformations on docstrings, function names, code syntax, and code format. The nominal baselines for docstrings and function name perturbations are the pass@ $k$  on nonperturbed datasets. For perturbations on code syntax and format, the nominal baseline is the pass@ $k$  on nonperturbed customized datasets with partial code (see Sect. 3.4). We use greedy sampling for all the models to eliminate randomness effect and enable fair comparisons as the default setting. We consider  $s = 5$ , i.e., we generate five different datasets with different random seeds for each type of perturbations and evaluate worst-case robustness performance according to the robustness evaluation metric defined in Sect. 3.6. To evaluate and compare model robustness in a unified fashion, we aggregate the worst performance across different perturbations under each category. Taking the docstring perturbation category as an example, we say the model is robust only when the model predicts correctly on all the  $s$  perturbed datasets for each transformation listed in Tab. 1. We present detailed numbers for each perturbation type in Appendix D, Tab. 11-18. In Appendix B, we showcase and analyze failure cases on CodeGen-16B-mono under three top perturbations, causing significant performance drops.

(1) **Diverse pretraining corpus helps with both generalization and worst-case robustness.** Comparing all code generation models with the same size 6B, CodeGen models have much better nominal performance, and have better robustness on  $RP_5@1$ , a very strict worst-case robustness metric. That is possibly because CodeGen models are pretrained over a more diverse corpus than InCoder and GPT-J and thus have more capacity to deal with unseen instances and perturbations. However, CodeGen models have worse performance on  $RD_5@1$  and  $RR_5@1$ , two robustness metrics relative to nominal performance, indicating that CodeGen models cannot generalize in a robust way (e.g., may learn to use spurious features in data).<sup>4</sup>

(2) **Larger model size brings improvement in**

<sup>4</sup>Although these models may have some subtle architecture-wise differences in details, we follow the benchmarking and evaluation strategies in previous works to focus more on pre-training parts and model sizes (e.g., BigBench (Srivastava et al., 2022), HELM (Liang et al., 2022)). We leave further ablation study for future work.

HumanEval	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Docstring	Nominal↑	0.232	0.140	0.262	0.195	<b>0.305</b>	0.195	0.104	0.152	0.122
	RP <sub>5</sub> @1↑	0.122	0.049	0.104	0.073	<b>0.128</b>	0.098	0.024	0.067	0.037
	RD <sub>5</sub> @1(%)↓	<b>47.37</b>	65.28	60.47	62.50	58.00	50.00	76.47	56.00	70.00
	RR <sub>5</sub> @1(%)↓	20.73	14.63	27.44	18.90	35.37	18.90	14.63	15.85	<b>10.98</b>
Function	Nominal↑	0.232	0.140	0.262	0.195	<b>0.305</b>	0.195	0.104	0.152	0.122
	RP <sub>5</sub> @1↑	0.140	0.061	0.146	0.116	<b>0.213</b>	0.116	0.055	0.098	0.073
	RD <sub>5</sub> @1(%)↓	39.47	56.52	44.19	40.63	<b>30.00</b>	40.63	47.06	36.00	40.00
	RR <sub>5</sub> @1(%)↓	14.02	10.37	18.90	12.20	19.51	9.146	8.537	9.756	<b>6.098</b>
Syntax	Nominal↑	0.402	0.293	0.518	0.366	<b>0.549</b>	0.390	0.189	0.323	0.250
	RP <sub>5</sub> @1↑	0.110	0.067	0.152	0.110	<b>0.159</b>	0.091	0.043	0.079	0.079
	RD <sub>5</sub> @1(%)↓	72.73	77.08	70.59	70.00	71.11	76.56	77.42	75.47	<b>68.29</b>
	RR <sub>5</sub> @1(%)↓	41.46	32.93	44.51	36.59	46.95	39.02	<b>21.34</b>	34.76	30.49
Format	Nominal↑	0.402	0.293	0.518	0.366	<b>0.549</b>	0.390	0.189	0.323	0.250
	RP <sub>5</sub> @1↑	0.268	0.207	0.274	0.195	<b>0.354</b>	0.232	0.091	0.171	0.104
	RD <sub>5</sub> @1(%)↓	33.33	<b>29.17</b>	47.06	46.67	35.56	40.63	51.61	47.17	58.54
	RR <sub>5</sub> @1(%)↓	23.17	16.46	32.93	23.78	25.00	22.56	<b>14.63</b>	23.78	21.95

Table 3: ReCode benchmark robustness evaluation on popular code generation models for HumanEval.

MBPP	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Docstring	Nominal↑	0.317	0.191	0.361	0.221	<b>0.407</b>	0.241	0.128	0.199	0.133
	RP <sub>5</sub> @1↑	0.137	0.050	0.147	0.042	<b>0.163</b>	0.045	0.011	0.031	0.013
	RD <sub>5</sub> @1(%)↓	<b>56.96</b>	73.66	59.38	80.93	59.85	81.28	91.20	84.54	90.00
	RR <sub>5</sub> @1(%)↓	36.86	34.39	41.89	36.76	46.72	44.66	<b>25.57</b>	35.32	30.08
Function	Nominal↑	0.317	0.191	0.361	0.221	<b>0.407</b>	0.241	0.128	0.199	0.133
	RP <sub>5</sub> @1↑	0.221	0.101	0.252	0.110	<b>0.279</b>	0.139	0.047	0.087	0.043
	RD <sub>5</sub> @1(%)↓	30.42	47.31	<b>30.40</b>	50.23	31.31	42.55	63.20	56.19	67.69
	RR <sub>5</sub> @1(%)↓	19.51	20.43	24.13	22.79	24.95	23.51	<b>16.22</b>	20.02	17.56
Syntax	Nominal↑	0.450	0.285	0.535	0.331	<b>0.571</b>	0.379	0.219	0.292	0.176
	RP <sub>5</sub> @1↑	0.027	0.008	0.027	0.008	<b>0.038</b>	0.017	0.008	0.006	0.004
	RD <sub>5</sub> @1(%)↓	<b>94.06</b>	97.12	95.01	97.52	93.34	95.39	96.24	97.89	97.66
	RR <sub>5</sub> @1(%)↓	59.03	45.07	64.17	47.74	67.04	54.21	35.42	45.79	<b>30.60</b>
Format	Nominal↑	0.450	0.285	0.535	0.331	<b>0.571</b>	0.379	0.219	0.292	0.176
	RP <sub>5</sub> @1↑	0.333	0.146	0.289	0.166	<b>0.403</b>	0.214	0.091	0.130	0.080
	RD <sub>5</sub> @1(%)↓	<b>26.03</b>	48.92	46.07	49.69	29.32	43.63	58.22	55.28	54.39
	RR <sub>5</sub> @1(%)↓	19.82	25.15	31.11	27.00	25.26	26.59	19.61	28.54	<b>18.28</b>

Table 4: ReCode benchmark robustness evaluation on popular code generation models for MBPP.

### worst-case robustness, but may risk overfitting.

In general, we observe higher RP<sub>5</sub>@1 for larger models within the same model family (e.g., improved from 0.174 to 0.217 for CodeGen-mono 2B to 16B on average across all perturbations), indicating larger model helps improve worst-case robustness. Similarly, we observe that larger models usually have larger RR<sub>5</sub>@1 (e.g., increased from 27.90% to 35.91% for CodeGen-mono 2B to 16B on average), indicating that larger models may risk overfitting as the relative performance drops under perturbations are significant.

(3) **Code generation models are most sensitive to syntax perturbation.** Among all perturbation types and across MBPP and HumanEval, we observe that syntax perturbations often result in the most performance drops. That reveals a significant limitation of syntax understanding ability of the state-of-the-art code generation models.

### (4) Datasets having more variances in code

### style poses more challenges on model robustness.

In Tab. 5, we can see that models show better robustness on HumanEval over MBPP on average. MBPP has more variances in code style (e.g., indent with 1 space), closer to natural code distribution hence more challenging for model robustness.

Category	Metric	HumanEval	MBPP
Docstring	RP <sub>5</sub> @1↑	<b>0.078</b>	0.071
	RD <sub>5</sub> @1(%)↓	<b>60.67</b>	75.31
	RR <sub>5</sub> @1(%)↓	<b>19.72</b>	36.92
Function	RP <sub>5</sub> @1↑	0.113	<b>0.142</b>
	RD <sub>5</sub> @1(%)↓	<b>41.61</b>	46.59
	RR <sub>5</sub> @1(%)↓	<b>12.06</b>	21.01
Syntax	RP <sub>5</sub> @1↑	<b>0.100</b>	0.025
	RD <sub>5</sub> @1(%)↓	<b>72.58</b>	93.40
	RR <sub>5</sub> @1(%)↓	<b>33.88</b>	47.86
Format	RP <sub>5</sub> @1↑	<b>0.211</b>	0.206
	RD <sub>5</sub> @1(%)↓	<b>43.30</b>	45.73
	RR <sub>5</sub> @1(%)↓	<b>22.70</b>	24.60

Table 5: Average robustness numbers across all models. MBPP is more challenging for robustness evaluation.

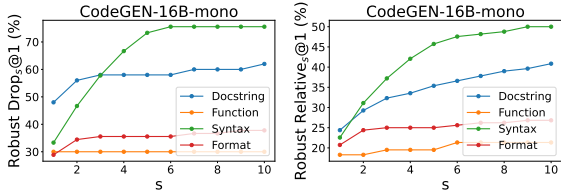


Figure 4: Robust Drop<sub>s</sub>@1 and Robust Relative<sub>s</sub>@1 under different  $s$ . Larger  $s$  indicates stronger perturbations evaluated and larger performance drops.

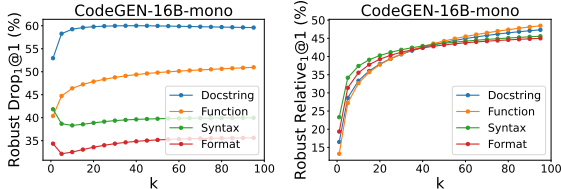


Figure 5: Robust Drop<sub>1</sub>@ $k$  and Robust Relative<sub>1</sub>@ $k$  under different  $k$  using sampling  $n = 100$ . Robust Drop remains stable while Robust Relative increases with  $k$ .

## 4.2 Ablation Study

**Robustness with  $s$  perturbed datasets.** As described in Sect. 3.6, our robustness metrics consider worst-case performance across  $s$  perturbed datasets for each perturbation. Larger  $s$  leads to stronger perturbations evaluated, larger performance drops, and more extensive coverage to practical failures. The performance drops will start converging when large enough  $s$  evaluated. We can clearly see such trends in Fig. 4 where we evaluate CodeGen-16B-mono RD<sub>s</sub>@1 and RR<sub>s</sub>@1 under greedy sampling with  $s = 1, \dots, 10$ . Perturbation categories like docstring and syntax that involve larger searching space and more randomness tend to benefit more with larger  $s$  (see Appendix A for details). As a trade-off, evaluation cost linearly increase with  $s$ . Thus, we recommend  $s = 5$  as a good balance between cost and evaluation strength. We summarize the ablation study in terms of larger sampling  $n$  in Appendix D.3 which can also benefit our proposed robustness estimation with additional sampling cost.

**Stable RD@ $k$  and increasing RR@ $k$  under different  $k$ .** Pass@ $k$  allows the model to have  $k$  trials and model performance is often reported with different  $k$ . With the sampling setting of  $n = 100$ , we plot the RD<sub>1</sub>@ $k$  and RR<sub>1</sub>@ $k$  in Fig. 5. Interestingly, we observe that RD@ $k$  stays stable across different  $k$  while RR@ $k$  increases with  $k$ . This is because larger  $k$  leads to higher nominal pass@ $k$  and RP@ $k$  but their relative ratio stays similar leading to stable RD. On the other hand, larger  $k$  involves more samples potentially changing results on perturbed datasets causing larger RR. Similar

	HumanEval	MBPP
Naturalness (Nominal) $\uparrow$	0.92	0.92
Naturalness (Perturbed) $\uparrow$	0.75	0.80
Semantics Similarity $\uparrow$	0.92	0.92

Table 6: Human evaluation for practical naturalness and semantic similarity by 5 annotators. Either 0, 0.5, or 1 is assigned to each data point indicating quality level.

	HumanEval		MBPP	
	Syntax	Format	Syntax	Format
CodeBLEU (syntax) $\uparrow$	0.95	0.98	0.93	0.96
CodeBLEU (dataflow) $\uparrow$	0.94	1.00	0.92	1.00

Table 7: Average CodeBLEU syntax and format scores between non-perturbed codes and perturbed ones with our syntax and format transformations.

trends on CodeGen-2B and 6B in Appendix D.2 further confirm the observations.

## 4.3 Perturbation Sample Quality

**Human evaluation.** To verify the naturalness of the perturbations in ReCode, we randomly sample and shuffle 100 and 50 perturbed and non-perturbed MBPP and HumanEval data points and create a shuffle mix of 300 samples. Each sample is shown to 5 human annotators who are familiar with Python and who are asked to rate naturalness out of 0: not natural, 0.5: possible to appear in practice but rare, and 1: natural. The scores for naturalness drop 14% on average for our perturbed data where drops mainly come from typos by Butterfingers, CharCaseChanges, SwapCharacter, etc.

In addition, we randomly sample 100 and 50 pairs perturbed and non-perturbed MBPP and HumanEval data points. Each pair is shown to 5 human annotators who are asked to rate semantics out of 0: totally changed, 0.5: slightly changed, and 1: exactly preserved. We summarize the main results in Tab. 6, and we present statistic details and setup in Appendix C.1. Notably, the majority vote (at least three out of five) is 1 for 90% of data points. We further provide automatic evaluation below to support the quality of our perturbed datasets, but human evaluation is in general more reliable.

**Docstring/function names similarity.** We measure the sentence cosine similarity between perturbed and non-perturbed docstrings and function names. We obtain the embeddings by sentence transformers using model all-mpnet-base-v2<sup>5</sup> (Song et al., 2020). Note that we split each function name into words to get

<sup>5</sup>Model embedding quality in <https://www.sbert.net>



sentence embeddings. On average, we have 0.93 and 0.81 for docstring and function name perturbations, showing that they well preserve the semantics. Scores for some function name perturbations are sensitive to typos due to the lack of sentence context (e.g., 0.21 for `interperse` and `intErpErse`). Appendix C.2 summarizes detailed numbers for each perturbation.

**Code syntax/format similarity.** In Tab. 7, we also measure the code similarity using CodeBLEU scores (Lu et al., 2021) for perturbed and non-perturbed data involving code syntax/format transformations. Here we consider the CodeBLEU score with syntax and dataflow separately as the evaluation metrics. On average, we have score 0.96 and 0.97 for CodeBLEU syntax and dataflow, showing good quality of perturbed datasets. Note that a few perturbations should expect low CodeBLEU scores: `doc2comments` transforms docstrings into comments causing changes of syntax; `Deadcode` insertion and `for-while` switch involve new if-conditions, loops, and new variables causing changes of code syntax and dataflow. Please refer to Appendix C.3 for details.

## 5 Conclusion

In this paper, we propose ReCode, a comprehensive robustness evaluation benchmark for code generation models. We collect and customize over 30 natural transformations under categories of docstrings, function names, code syntax, and code format perturbations. These transformations are carefully selected and designed to be natural in practice and preserve the semantic meaning after perturbations. We further propose general worst-case robustness metrics to give a unified overview of the model robustness performance. We empirically demonstrate our ReCode benchmark on popular models including CodeGen, InCoder, and GPT-J using HumanEval and MBPP datasets and function completion tasks derived from them. With human evaluation, over 90% of our perturbed data are confirmed to preserve the original semantic meaning; sentence similarity and CodeBLEU scores additionally support the quality of perturbations in ReCode.

## Limitations

ReCode benchmark has several limitations: (1) It contains perturbed datasets based on HumanEval and MBPP which focuses on Python function completion use cases. Therefore, we only perform

evaluation on Python language and not be able to capture robustness in a wide variety of code completion use cases. However, our transformations are generalizable and could be easily extended to other languages and also other code-related datasets (Athiwaratkun et al., 2023). We encourage researchers to apply and extend ReCode benchmark to additional languages and other code-related tasks; (2) ReCode benchmark is designed for robustness evaluation and cannot mitigate the lack of robustness. Given that our benchmark can be used to generate comprehensive collection of perturbed data, we believe that it can be used for training data augmentation to enhance model robustness. We will consider corresponding robust training strategy design and evaluation in the future work.

## Ethics Statement

Our ReCode robustness benchmark aims to provide a comprehensive robustness evaluation framework for any code-generation models, which we believe is critical towards building robust and user-friendly language models for code. With the new robustness evaluation metrics, users can rely on ReCode and assess model predictions with more confidence. The model trainers, on the other hand, will be aware of the potential vulnerabilities that might cause mispredictions in practice and mitigate them before deployments. Therefore, we believe our ReCode benchmark is beneficial in terms of broader impact.

## References

- Wasi Ahmad, Saikat Chakraborty, Baishakhi Ray, and Kai-Wei Chang. 2021. [Unified pre-training for program understanding and generation](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2655–2668, Online. Association for Computational Linguistics.
- Ben Athiwaratkun, Sanjay Krishna Gouda, Zijian Wang, Xiaopeng Li, Yuchen Tian, Ming Tan, Wasi Uddin Ahmad, Shiqi Wang, Qing Sun, Mingyue Shang, Sujan Kumar Gonugondla, Hantian Ding, Varun Kumar, Nathan Fulton, Arash Farahani, Siddhartha Jain, Robert Giaquinto, Haifeng Qian, Murali Krishna Ramanathan, Ramesh Nallapati, Baishakhi Ray, Parminder Bhatia, Sudipta Sengupta, Dan Roth, and Bing Xiang. 2023. [Multi-lingual evaluation of code generation models](#). In *The 11th International Conference on Learning Representations (ICLR)*.
- Jacob Austin, Augustus Odena, Maxwell I. Nye, Maarten Bosma, Henryk Michalewski, David Dohan,

- Ellen Jiang, Carrie J. Cai, Michael Terry, Quoc V. Le, and Charles Sutton. 2021. [Program synthesis with large language models](#). *ArXiv preprint*, abs/2108.07732.
- Pavol Bielik and Martin T. Vechev. 2020. [Adversarial robustness for code](#). In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 896–907. PMLR.
- Saikat Chakraborty, Toufique Ahmed, Yangruibo Ding, Premkumar Devanbu, and Baishakhi Ray. 2022. [Naten: Generative pre-training by "naturalizing" source code](#). *ArXiv preprint*, abs/2206.07585.
- 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*, abs/2107.03374.
- Kaustubh D Dhole, Varun Gangal, Sebastian Gehrmann, Aadesh Gupta, Zhenhao Li, Saad Mahamood, Abinaya Mahendiran, Simon Mille, Ashish Srivastava, Samson Tan, et al. 2021. [NL-augmenter: A framework for task-sensitive natural language augmentation](#). *ArXiv preprint*, abs/2112.02721.
- Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. 2020. [CodeBERT: A pre-trained model for programming and natural languages](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1536–1547, Online. Association for Computational Linguistics.
- Daniel Fried, Armen Aghajanyan, Jessy Lin, Sida Wang, Eric Wallace, Freda Shi, Ruiqi Zhong, Wen-tau Yih, Luke Zettlemoyer, and Mike Lewis. 2022. [InCoder: A generative model for code infilling and synthesis](#). *ArXiv preprint*, abs/2204.05999.
- Matt Gardner, Yoav Artzi, Victoria Basmov, Jonathan Berant, Ben Bogin, Sihao Chen, Pradeep Dasigi, Dheeru Dua, Yanai Elazar, Ananth Gottumukkala, Nitish Gupta, Hannaneh Hajishirzi, Gabriel Ilharco, Daniel Khashabi, Kevin Lin, Jiangming Liu, Nelson F. Liu, Phoebe Mulcaire, Qiang Ning, Sameer Singh, Noah A. Smith, Sanjay Subramanian, Reut Tsarfaty, Eric Wallace, Ally Zhang, and Ben Zhou. 2020. [Evaluating models' local decision boundaries via contrast sets](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 1307–1323, Online. Association for Computational Linguistics.
- Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang, Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy, Shengyu Fu, Michele Tufano, Shao Kun Deng, Colin B. Clement, Dawn Drain, Neel Sundaresan, Jian Yin, Daxin Jiang, and Ming Zhou. 2021. [Graphcodebert: Pre-training code representations with data flow](#). In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net.
- Kilem L Gwet. 2014. *Handbook of inter-rater reliability: The definitive guide to measuring the extent of agreement among raters*. Advanced Analytics, LLC.
- Akshita Jha and Chandan K Reddy. 2022. [Codeattack: Code-based adversarial attacks for pre-trained programming language models](#). *ArXiv preprint*, abs/2206.00052.
- Di Jin, Zhijing Jin, Joey Tianyi Zhou, and Peter Szolovits. 2020. [Is BERT really robust? A strong baseline for natural language attack on text classification and entailment](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*, pages 8018–8025. AAAI Press.
- Douwe Kiela, Max Bartolo, Yixin Nie, Divyansh Kaushik, Atticus Geiger, Zhengxuan Wu, Bertie Vidgen, Grusha Prasad, Amanpreet Singh, Pratik Ringshia, Zhiyi Ma, Tristan Thrush, Sebastian Riedel, Zeerak Waseem, Pontus Stenetorp, Robin Jia, Mohit Bansal, Christopher Potts, and Adina Williams. 2021. [Dynabench: Rethinking benchmarking in NLP](#). In *North American Association for Computational Linguistics (NAACL)*, pages 4110–4124.
- Linyang Li, Ruotian Ma, Qipeng Guo, Xiangyang Xue, and Xipeng Qiu. 2020. [BERT-ATTACK: Adversarial attack against BERT using BERT](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6193–6202, Online. Association for Computational Linguistics.
- Yujia Li, David Choi, Junyoung Chung, Nate Kushman, Julian Schrittwieser, Rémi Leblond, Tom Eccles, James Keeling, Felix Gimeno, Agustin Dal Lago, Thomas Hubert, Peter Choy, Cyprien de Masson d'Autume, Igor Babuschkin, Xinyun Chen, Po-Sen Huang, Johannes Welbl, Sven Gowal, Alexey Cherepanov, James Molloy, Daniel Mankowitz, Esme Sutherland Robson, Pushmeet Kohli, Nando de Freitas, Koray Kavukcuoglu, and Oriol Vinyals. 2022. [Competition-level code generation with alphacode](#). *ArXiv preprint*, abs/2203.07814.
- Zhenhao Li and Lucia Specia. 2019. [Improving neural machine translation robustness via data augmentation: Beyond back-translation](#). In *Proceedings of the 5th Workshop on Noisy User-generated Text (W-NUT 2019)*, pages 328–336, Hong Kong, China. Association for Computational Linguistics.
- Percy Liang, Rishi Bommasani, Tony Lee, Dimitris Tsipras, Dilara Soylu, Michihiro Yasunaga, Yian

- Zhang, Deepak Narayanan, Yuhuai Wu, Ananya Kumar, et al. 2022. [Holistic evaluation of language models](#). *ArXiv preprint*, abs/2211.09110.
- Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, et al. 2021. Codexglue: A machine learning benchmark dataset for code understanding and generation. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- Mitchell P. Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. [Building a large annotated corpus of English: The Penn Treebank](#). *Computational Linguistics*, 19(2):313–330.
- Simon Mille, Kaustubh Dhole, Saad Mahamood, Laura Perez-Beltrachini, Varun Gangal, Mihir Kale, Emiel van Miltenburg, and Sebastian Gehrmann. 2021. Automatic construction of evaluation suites for natural language generation datasets. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- George A. Miller. 1992. [WordNet: A lexical database for English](#). In *Speech and Natural Language: Proceedings of a Workshop Held at Harriman, New York, February 23-26, 1992*.
- Yixin Nie, Adina Williams, Emily Dinan, Mohit Bansal, J. Weston, and Douwe Kiela. 2020. Adversarial nli: A new benchmark for natural language understanding. In *Association for Computational Linguistics (ACL)*.
- Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. 2022. [A conversational paradigm for program synthesis](#). *ArXiv preprint*, abs/2203.13474.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2020. [Mpnnet: Masked and permuted pre-training for language understanding](#). In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*.
- Aarohi Srivastava, Abhinav Rastogi, Abhishek Rao, Abu Awal Md Shoeb, Abubakar Abid, Adam Fisch, Adam R Brown, Adam Santoro, Aditya Gupta, Adrià Garriga-Alonso, et al. 2022. [Beyond the imitation game: Quantifying and extrapolating the capabilities of language models](#). *ArXiv preprint*, abs/2206.04615.
- Amane Sugiyama and Naoki Yoshinaga. 2019. [Data augmentation using back-translation for context-aware neural machine translation](#). In *Proceedings of the Fourth Workshop on Discourse in Machine Translation (DiscoMT 2019)*, pages 35–44, Hong Kong, China. Association for Computational Linguistics.
- Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>.
- Boxin Wang, Chejian Xu, Shuohang Wang, Zhe Gan, Yu Cheng, Jianfeng Gao, Ahmed Hassan Awadallah, and Bo Li. 2021a. Adversarial glue: A multi-task benchmark for robustness evaluation of language models. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Yue Wang, Weishi Wang, Shafiq Joty, and Steven C.H. Hoi. 2021b. [CodeT5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 8696–8708, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Zhou Yang, Jieke Shi, Junda He, and David Lo. 2022. [Natural attack for pre-trained models of code](#). *ArXiv preprint*, abs/2201.08698.
- Yuan Zang, Fanchao Qi, Chenghao Yang, Zhiyuan Liu, Meng Zhang, Qun Liu, and Maosong Sun. 2020. [Word-level textual adversarial attacking as combinatorial optimization](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 6066–6080, Online. Association for Computational Linguistics.
- Wei Emma Zhang, Quan Z Sheng, Ahoud Alhazmi, and Chenliang Li. 2020. Adversarial attacks on deep-learning models in natural language processing: A survey. *TIST*, 11(3):1–41.
- Yu Zhou, Xiaoqing Zhang, Juanjuan Shen, Tingting Han, Taolue Chen, and Harald Gall. 2022. Adversarial robustness of deep code comment generation. *ACM Transactions on Software Engineering and Methodology (TOSEM)*, 31(4):1–30.

## A Transformation Details and Qualitative Examples

In this section, we give detailed descriptions and settings for each type of perturbations that are included in our ReCode benchmark with qualitative examples for illustrations.

### A.1 Natural Transformations on Docstrings

For natural transformations on docstrings, we aim to perturb the docstrings to their variances that preserve the semantics and also appear natural in practice. Specifically, we will first extract and perturb the docstrings with the following natural transformations in each prompt, and then attach their perturbed versions to the prompt. To preserve semantics for the code generation task prompts, we extract a blacklist of program keywords using tree-sitter as discussed in Sect. 3.2 that are excluded from perturbations. We extend most transformations from NL-Augmenter (Dhole et al., 2021), a standard library designed for data augmentation and robustness evaluation on text. We list some qualitative examples in Tab. 1.

**BackTranslation.** BackTranslation paraphrases the docstrings by translating them to another language (in this case, German) and then back to English. It is a common method for data augmentation in generating sentence variances with the same semantics (Li and Specia, 2019; Sugiyama and Yoshinaga, 2019). Overall, it can reliably generate high quality perturbed docstrings. We use the default implementation in NL-Augmenter (Dhole et al., 2021). BackTranslation contains no randomness in transformations.

**ButterFingers.** ButterFingers transformation randomly selects characters of the docstrings and perturbs each of them to a random subset of similar characters, it is from (Dhole et al., 2021) and is also used in (Mille et al., 2021). Since this transformation tends to introduce character-level typos, we set randomness for perturbing each character to be low as 0.05 for naturalness consideration.

**ChangeCharCase.** ChangeCharCase transformation randomly changes the selected characters to upper case in the docstrings. We use the default probability 0.35 where majority annotators vote 0.5 for naturalness in the setting of Sect. 4.3.

**EnglishInflectionalVariation.** This transformation randomly selects words in the docstring and change them to a random inflection variance. This can be from plural to singular (or vice versa) for nouns and tense changes for verbs. To maintain naturalness, the perturbation is constrained to be the same Part of Speech (POS) tag in the Penn Treebank (Marcus et al., 1993).

**SwapCharacters.** This transformation randomly selects pairs of adjacent characters in the docstring and swap them. This represents a common type of typos by humans. To ensure naturalness, we set the probability as 0.05 for making the swap.

**SynonymInsertion.** This transformation randomly select words in the docstrings and inserts their synonyms in WordNet (Miller, 1992). Punctuations and stopwords are excluded. We set the probability to be 0.35 considering low success rate after keywords filtering.

**SynonymSubstitution.** This transformation randomly selects words in the docstring and replaces each one with a synonym from WordNet (Miller, 1992). Similar to SynonymInsertion, we set the probability as 0.35 to balance naturalness and perturbation success rates.

**TenseTransformationPast.** This is a deterministic transformation that converts sentences in the docstring to past tense.

**TenseTransformationFuture.** This is a deterministic transformation that converts sentences in the docstring to future tense.

**Whitespace.** This transformation inserts or deletes a single white space at randomly selected locations in the docstring. This represents a common type of typos by humans. Following NL-Augmenter, we use probability 0.1 for adding whitespaces and 0.05 for removing whitespaces.

### A.2 Natural Transformations on Function Names

These transformations modify the name of the target function to generate. Any references to the function name in the prompt, e.g., in docstring, are also modified to maintain consistency. Qualitative examples can be found in Tab. 2.

**CamelCase.** A function name is often composed of multiple words. If the original function name concatenates the words in camel-case style, this



transformation changes it to snake-case, and vice versa. This transformation is deterministic.

**ButterFingers.** ButterFingers transformation randomly selects characters of the docstrings and perturbs each of them to a random subset of similar characters, it is from (Dhole et al., 2021) and is also used in (Mille et al., 2021). Since this transformation tends to introduce character-level typos, we set randomness for perturbing each character to be low as 0.05 for naturalness consideration.

**SwapCharacters.** This transformation randomly selects pairs of adjacent characters in the function name and swap each pair. This represents a common type of typos by humans. To control naturalness, we set the probability to be 0.05, same setting as the docstring perturbations.

**ChangeCharCase.** ChangeCharCase transformation randomly changes the selected characters to upper case in the docstrings. We use the default probability 0.35 where majority annotators vote 0.5 for naturalness in the setting of Sect. 4.3.

**InflectionalVariation.** This transformation randomly selects words in the function name and applies a random inflection on them. This can be from plural to singular (or vice versa) for nouns and tense change for verbs. To control naturalness, the perturbation is constrained to be the same Part of Speech (POS) tag in the Penn Treebank (Marcus et al., 1993).

**SynonymSubstitution.** This transformation randomly selects words in the docstring and replaces each one with a synonym from WordNet (Miller, 1992). Similar to SynonymInsertion, we set the probability as 0.35 to balance naturalness and perturbation success rates.

### A.3 Natural Transformations on Code Syntax

These transformations modify the code content in the prompt. We derived function completion tasks with half the code from the canonical solutions such that the following code transformations and robustness evaluation can be performed. To guarantee fair comparisons to the nominal baseline, we make sure that we have the same block of code before and after code perturbations. In the following part we show qualitative examples on the same MBPP sample baseline ( Fig. 6).

```
def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello", "l")
    "heo"
    >>> remove_Occ("abcda", "a")
    "bcd"
    >>> remove_Occ("PHP", "P")
    "H"
    """
    for i in range(len(s)):
        if s[i] == ch:
            s = s[0:i] + s[i + 1 :]
            break
```

Figure 6: An baseline example of the prompt with partial code derived from original MBPP prompt for robustness evaluation on code.

**DeadCodeInserter.** This transformation inserts a block of useless code at a random location. The added block can be a loop of zero iteration or an if condition that is always false. The code content inside the dummy loop or if condition is randomly selected from the adjacent code statements with limited tree-sitter node sizes.

```
def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello", "l")
    "heo"
    >>> remove_Occ("abcda", "a")
    "bcd"
    >>> remove_Occ("PHP", "P")
    "H"
    """
    if False:
        break
    for i in range(len(s)):
        if s[i] == ch:
            s = s[0:i] + s[i + 1 :]
            break
```

Figure 7: An example of the DeadCodeInserter perturbation.

**For-While Switch.** This transformation randomly selects a for-loop or while-loop in the prompt and transforms it to its equivalent counterpart.

**OperandSwap.** This transformation randomly selects a binary logical operation, swaps the two operands, and modifies the operator if necessary to maintain semantic equivalence.

**VarRenamerCB.** This transformation selects the most frequently referenced variable name in the partial code and replaces it throughout the prompt with a new name obtained by CodeBERT (Feng

```

def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    i = 0
    while i in list(range(len(s))):
        if s[i] == ch:
            s = s[0:i] + s[i + 1 :]
            break

```

Figure 8: An example of the For-While Switch perturbation.

```

def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    for i in range(len(s)):
        if ch == s[i]:
            s = s[0:i] + s[i + 1 :]
            break

```

Figure 9: An example of the OperandSwap perturbation.

et al., 2020). Specifically, we replace all occurrence of the variable name with a mask token, and then run CodeBERT inference to obtain candidate names at each location, where each candidate name comes with a probability score. We pick the candidate name with the highest aggregated score across locations. This transformation is inspired by (Jha and Reddy, 2022; Li et al., 2020).

```

def remove_Occ(lines, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    for i in range(len(lines)):
        if lines[i] == ch:
            lines = lines[0:i] + lines[i + 1 :]
            break

```

Figure 10: An example of the VarRenamerCB perturbation.

**VarRenamerNaive.** This transformation selects the most frequently referenced variable name in the partial code and replaces it with "VAR\_0". This is the original implementation in the NatGen package. This transformation is deterministic.

```

def remove_Occ(VAR_0, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    for i in range(len(VAR_0)):
        if VAR_0[i] == ch:
            VAR_0 = VAR_0[0:i] + VAR_0[i + 1 :]
            break

```

Figure 11: An example of the VarRenamerNaive perturbation.

**VarRenamerRN.** This transformation selects the most frequently referenced variable name in the partial code and replaces it with a random string with half alphabetic and half numeric characters.

```

def remove_Occ(z5, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    for i in range(len(z5)):
        if z5[i] == ch:
            z5 = z5[0:i] + z5[i + 1 :]
            break

```

Figure 12: An example of the VarRenamerRN perturbation.

#### A.4 Natural Transformations on Code Format

**Tab-Indent.** This transformation replaces any space indents with tabs or replaces tabs with 4 spaces for indent-sensitive languages like Python. This transformation is deterministic.

**Line Split.** This transformation splits the longest line in the partial code into two lines. This transformation is deterministic.

**Doc2Comments.** This transformation changes the style of the documentation in the prompt. For Python, it converts docstring (e.g., """ docstring """) to commented lines (e.g., # docstring) and vice versa. For Java, it converts comments in the

```

def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    for i in range(len(s)):
        if s[i] == ch:
            s = s[0:i] + s[i + 1 :]
            break

```

Figure 13: An example of the Tab-Indent perturbation.

```

def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    (new line)
    for i in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
    (new line)
    break

```

Figure 16: An example of the NewlineRandom perturbation.

```

def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    for i \
in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break

```

Figure 14: An example of the Line Split perturbation.

```

def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    for i in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break
    (new line)

```

Figure 17: An example of the NewlineAfterCode perturbation.

format of `/* docstring */` to `// docstring` and vice versa. This transformation is deterministic.

```

def remove_Occ(s, ch):
    # Write a python function to remove
    # first and last occurrence of a
    # given character from the string.
    # >>> remove_Occ("hello","l")
    # "heo"
    # >>> remove_Occ("abcda","a")
    # "bcd"
    # >>> remove_Occ("PHP","P")
    # "H"
    for i in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break

```

Figure 15: An example of the Doc2Comments perturbation.

```

def remove_Occ(s, ch):
    """
    Write a python function to remove
    first and last occurrence of a
    given character from the string.
    >>> remove_Occ("hello","l")
    "heo"
    >>> remove_Occ("abcda","a")
    "bcd"
    >>> remove_Occ("PHP","P")
    "H"
    """
    (new line)
    for i in range(len(s)):
        if (s[i] == ch):
            s = s[0 : i] + s[i + 1:]
            break

```

Figure 18: An example of the NewlineAfterDoc perturbation.

**NewlineRandom.** This transformation inserts empty lines at randomly selected positions.

**NewlineAfterCode.** This transformation inserts an empty line at the end of the prompt. This transformation is deterministic.

**NewlineAfterDoc.** This transformation inserts an empty line between the docstring and the partial code. This transformation is deterministic.

## B Failure Case Study under Perturbations

In this section, we showcase and analyze some failure cases on CodeGen-16B-mono and perturbed HumanEval datasets under three top perturbations that will cause significant performance drops.

DeadCode insertion is one of the most effective

perturbations. It can commonly mislead the model predictions with the inserted dead code, especially when the completions are required right after the inserted dead code. Fig. 19 shows an failure example where CodeGen-mono-16B only predicts a newline after inserted meaningless for loop, which might be misled by the inserted return statement.

```
def change_base(x: int, base: int):
    """Change numerical base of input number
    x to base. return string representation
    after the conversion.
    base numbers are less than 10.
    >>> change_base(8, 3)
    '22'
    >>> change_base(8, 2)
    '1000'
    >>> change_base(7, 2)
    '111'
    """
    ret = ""
    while x > 0:
        ret = str(x % base) + ret
        x = x // base
    return ret
```

} Original completion

(a) Correct completion without perturbation.

```
def change_base(x: int, base: int):
    """Change numerical base of input number
    x to base. return string representation
    after the conversion.
    base numbers are less than 10.
    >>> change_base(8, 3)
    '22'
    >>> change_base(8, 2)
    '1000'
    >>> change_base(7, 2)
    '111'
    """
    ret = ""
    while x > 0:
        for _i_3 in range(0):
            return ret
    [empty line]
```

} Inserted deadcode  
} New completion

(b) Wrong completion perturbed by deadcode insertion.

Figure 19: HumanEval showcase 1 illustrating failure case under deadcode insertion.

Fig. 20 shows a failure example of CodeGen-16B-mono on a prompt where an empty newline is inserted right before completion. Such simple perturbation causes wrong predictions for the following if-else conditions. It is especially effective when the required completion code is complicated.

ButterFingers perturbation on docstring causes large performance drops as well. Fig. 21 shows another failure example on CodeGen-16B-mono. The typos introduced in the perturbation might cause the model to misunderstand the targeted docstrings, leading to wrong model completions.

```
def sum_squares(lst):
    """
    This function will take a list of integers.
    For all entries in the list, the function
    shall square the integer entry if its index
    is a multiple of 3 and will cube the integer
    entry if its index is a multiple of 4 and
    not a multiple of 3. The function will not
    change the entries in the list whose indexes
    are not a multiple of 3 or 4. The function
    shall then return the sum of all entries.

    Examples:
    For lst = [1,2,3] the output should be 6
    For lst = [] the output should be 0
    For lst = [-1,-5,2,-1,-5] the output should
    be -126
    """
    result = []
    for i in range(len(lst)):
        if i % 3 == 0:
            result.append(lst[i]**2)
        elif i % 4 == 0 and i % 3 != 0:
            result.append(lst[i]**3)
        else:
            result.append(lst[i])
    return sum(result)
```

} Original completion

(a) Correct completion without perturbation.

```
def sum_squares(lst):
    """
    This function will take a list of integers.
    For all entries in the list, the function
    shall square the integer entry if its index
    is a multiple of 3 and will cube the integer
    entry if its index is a multiple of 4 and
    not a multiple of 3. The function will not
    change the entries in the list whose indexes
    are not a multiple of 3 or 4. The function
    shall then return the sum of all entries.

    Examples:
    For lst = [1,2,3] the output should be 6
    For lst = [] the output should be 0
    For lst = [-1,-5,2,-1,-5] the output should
    be -126
    """
    result = []
    for i in range(len(lst)):
        if i % 3 == 0:
            result.append(lst[i]**2)
        (new line)
        elif i % 4 == 0 and i % 3 != 0:
            result.append(lst[i]**3)
    return sum(result)
```

} New completion

(b) Wrong completion perturbed by NewLineAfterCode insertion.

Figure 20: HumanEval showcase 2 illustrating failure case under NewLineAfterCode insertion.

## C Perturbation Sample Quality

### C.1 Details for Human Evaluation

The annotators are all recruited from software engineers online who have good experience in Python via strict coding interview. To guarantee the reliability of the human evaluation results, we first conducted annotation trials with our annotators. We gave them clear definitions for each level of naturalness and semantic similarity.



```
def truncate_number(number: float) -> float:
    """ Given a positive floating point number,
    it can be decomposed into and integer part
    (largest integer smaller than given number)
    and decimals (leftover part always smaller
    than 1).
    Return the decimal part of the number.
    >>> truncate_number(3.5)
    0.5
    """
    return number - int(number)
```

(a) Correct completion without perturbation.

```
def truncate_number(number: float) -> float:
    """ Givfn a positive floating point number,
    it can be decimposev into and integer part
    (largest integer smeller than given number)
    and decimals (leftover part always wnaller
    thai 1).
    Retuxn the decimal pert of the number.
    >>> truncate_number(3.5)
    0.5
    """
    return int(number) + (number - int(number))
```

(b) Wrong completion perturbed by ButterFingers.

Figure 21: HumanEval showcase 3 illustrating failure case under ButterFingers perturbations on docstrings.

We measure the inter-annotator agreement rate Fleiss Kappa in Tab. 8. The overall average Fleiss Kappa for the annotations is 0.52, 0.36 for semantic and naturalness measurements on perturbed samples. The confidence interval (95%) with bootstrap sampling (10K samples) is [0.515, 0.528] and [0.358, 0.364], indicating that our annotation reaches “moderate agreement” and thus our annotations are reliable (Gwet, 2014). The scores from annotators are not perfectly consistent especially for naturalness since people have different preferences for code.

Fleiss Kappa	HumanEval	MBPP
Naturalness (Nominal) ↑	0.362	0.301
Naturalness (Perturbed) ↑	0.435	0.326
Semantics Similarity ↑	0.658	0.461

Table 8: Fleiss Kappa of human evaluation.

## C.2 Sentence Transformers for Docstring/Function Names Similarity

In this subsection, we give experimental details for measuring the sentence similarity of perturbed and unperturbed data points using sentence transformers.

To measure the similarity scores for the docstring perturbations, we first extract docstrings from each pair of perturbed and unperturbed data points, and we use sentence transformer all-mpnet-base-v2 (Song et al., 2020) to predict

an embedding vector for each docstring. Then cosine similarity is calculated and reported for each pair of perturbed and unperturbed datapoints.

Same process cannot be directly applied to function name perturbations since function names are concatenations of words instead of common sentences, barely seen by the sentence transformer training data. In order get more accurate sentence embeddings for function names, we first split each name into words (e.g., has\_close\_elements to has close elements) and then calculate the corresponding cosine similarities.

In Table 9, we present the detailed results for each type of perturbations for sentence similarity. On average, we can have 0.93 and 0.92 similarity scores for docstring perturbations and 0.80 and 0.81 for function name perturbations on the HumanEval and MBPP datasets. The overall high similarity numbers provide support that our perturbations have good quality in naturalness and semantic preservation from the unperturbed inputs.

Some function name perturbations including ButterFinger, SynonymSubstitution, and Char-CaseChange have relatively low sentence similarity. This is mainly because the function names only include keywords without complete sentence context and thus minor changes to each words could potentially cause large change in measured cosine similarity. For instance, character case changes on function name intersperse to intErspErse which lacks of context only has 0.21 similarity. On the other hand, the function names with more context has much higher scores, e.g., 1.0 similarity score for has\_close\_elements and has\_Close\_Elements.

## C.3 CodeBLEU Scores for Code Similarity

Here we present the experimental details for the CodeBLEU syntax and dataflow scores to quantitatively measure the quality of our code syntax and format transformations.

The measurement is straightforward. The unperturbed baseline is each data point from our customized partial code datasets derived from HumanEval and MBPP. The perturbed one is the same data point transformed by each type of our perturbations. The CodeBLEU syntax and dataflow scores are then directly measured using the CodeXGLUE (Lu et al., 2021) implementation.<sup>6</sup>

<sup>6</sup><https://github.com/microsoft/CodeXGLUE>

Categories	Perturbations	HumanEval	MBPP
Docstring	BackTranslation	0.91	0.95
	ButterFingers	0.87	0.89
	ChangeCharCase	1.00	1.00
	EnglishInflectionalVariation	0.96	0.93
	SwapCharacters	0.90	0.87
	SynonymInsertion	0.91	0.88
	SynonymSubstitution	0.88	0.84
	TenseTransformationPast	0.98	1.00
	TenseTransformationFuture	0.97	0.97
	Whitespace	0.90	0.86
Function	CamelCase	1.00	1.00
	ButterFingers	0.57	0.57
	SwapCharacters	0.75	0.75
	ChangeCharCase	0.86	0.96
	InflectionalVariation	0.94	0.93
	SynonymSubstitution	0.68	0.64

Table 9: Cosine similarity for each type of perturbations where perturbed and unperturbed docstrings/function names are embedded by the SOTA sentence transformer.

In Table 10, we present the detailed CodeBLEU results for each type of perturbations. The average numbers are summarized in Table 7. Overall, 77% and 89% of our transformations have over 0.9 CodeBLEU syntax and dataflow scores, showing good quality in preserving semantics from the unperturbed code.

However, CodeBLEU syntax and dataflow are not perfect in quantitatively measuring naturalness and semantic preservation for the perturbations and thus some perturbations have expected relatively low scores: Doc2Comments transforms docstrings into comments causing changes of syntax; Deadcode insertion and for-while switch involve new if-conditions, loops, and new variables causing changes of code syntax and dataflow.

## D Additional Results

### D.1 Fine-grained Robustness Evaluation

We present the robustness evaluation for each type of perturbations from Table 11 to 18. The evaluation setting is the same as Table 3 and 4 where we evaluate various sizes of CodeGen (Nijkamp et al., 2022), InCoder (Fried et al., 2022), and GPT-J (Wang and Komatsuzaki, 2021) with greedy sampling. For each type of perturbations, we randomly generate  $s = 5$  different perturbed datasets derived from HumanEval and MBPP. For perturbations without randomness, only one single version of perturbed dataset is evaluated. The list of indeterminate perturbations can be found in Appendix A.

### D.2 Additional Results for Different k

As discussed in Sect. 4.2, we observe that Robust Drop stays stable across different k while Robust Relative increases linearly with k. We present additional results on CodeGen-2B-mono, CodeGen-6B-mono along with CodeGen-16B-mono in Fig. 22. We evaluate each model with large  $n$  ( $n = 100$ ) using top-p sampling strategy with probability 0.95 and temperature 0.2.

### D.3 Additional Results for Large Sampling n

Larger sampling  $n$  is commonly used for preventing model generation variances and providing accurate estimations. The evaluation cost increases linearly to  $n$ . Here we show that larger  $n$  can also benefit our proposed three robustness metrics but not causing significant differences. In specific, we measure Robust Pass<sub>1</sub>@1, Robust Drop<sub>1</sub>@1, and Robust Relative<sub>1</sub>@1 on CodeGen-16B-mono and HumanEval dataset. The model is run with  $n = 100$  using top-p sampling strategy with probability 0.95 and temperature 0.2. We present detailed results in Tab. 19.

Categories	Perturbations	HumanEval		MBPP	
		CodeBLEU (syntax)	CodeBLEU (dataflow)	CodeBLEU (syntax)	CodeBLEU (dataflow)
Syntax	DeadCodeInserter	0.85	0.79	0.72	0.67
	For-While Switch	0.92	0.90	0.84	0.86
	OperandSwap	0.91	1.00	0.90	1.00
	VarRenamerCB	1.00	0.99	0.93	0.99
	VarRenamerNaive	1.00	0.99	0.93	0.99
	VarRenamerRN	1.00	0.99	0.93	0.99
Format	Tab-Indent	1.00	1.00	1.00	1.00
	Line Split	1.00	1.00	1.00	1.00
	Doc2Comments	0.84	1.00	0.76	1.00
	NewlineRandom	1.00	1.00	1.00	1.00
	NewlineAfterCode	1.00	1.00	1.00	1.00
	NewlineAfterDoc	1.00	1.00	1.00	1.00

Table 10: CodeBLEU syntax and format similarity scores between unperturbed codes and perturbed ones with each type of our syntax and format transformations.

HumanEval	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Nominal	RP <sub>5</sub> @1↑	0.232	0.140	0.262	0.195	<b>0.305</b>	0.195	0.104	0.152	0.122
BackTranslation	RP <sub>5</sub> @1↑	0.213	0.116	0.238	0.159	<b>0.244</b>	0.152	0.098	0.134	0.098
	RD <sub>5</sub> @1(%)↓	7.89	17.39	9.30	18.75	20.00	21.88	<b>5.88</b>	12.00	20.00
	RR <sub>5</sub> @1(%)↓	4.27	6.10	8.54	6.10	10.98	5.49	<b>3.05</b>	3.05	3.66
ButterFingers	RP <sub>5</sub> @1↑	0.165	0.098	0.171	0.122	<b>0.189</b>	0.116	0.067	0.098	0.067
	RD <sub>5</sub> @1(%)↓	<b>28.95</b>	30.43	34.88	37.50	38.00	40.62	35.29	36.00	45.00
	RR <sub>5</sub> @1(%)↓	10.37	7.32	15.85	10.37	20.12	12.20	7.32	9.15	<b>6.71</b>
ChangeCharCase	RP <sub>5</sub> @1↑	0.152	0.079	0.152	0.104	<b>0.177</b>	0.122	0.037	0.098	0.049
	RD <sub>5</sub> @1(%)↓	<b>34.21</b>	43.48	41.86	46.88	42.00	37.50	64.71	36.00	60.00
	RR <sub>5</sub> @1(%)↓	12.80	10.98	15.85	9.76	17.68	9.15	10.37	<b>7.32</b>	7.93
EnglishInflectional Variation	RP <sub>5</sub> @1↑	0.207	0.134	0.226	0.171	<b>0.268</b>	0.177	0.091	0.146	0.104
	RD <sub>5</sub> @1(%)↓	10.53	4.35	13.95	12.50	12.00	9.38	11.76	<b>4.00</b>	15.00
	RR <sub>5</sub> @1(%)↓	3.66	3.05	8.54	6.10	7.93	4.27	<b>1.22</b>	1.83	3.05
SwapCharacters Perturbation	RP <sub>5</sub> @1↑	0.159	0.098	0.183	0.128	<b>0.207</b>	0.134	0.085	0.104	0.067
	RD <sub>5</sub> @1(%)↓	31.58	30.43	30.23	34.38	32.00	31.25	<b>17.65</b>	32.00	45.00
	RR <sub>5</sub> @1(%)↓	12.20	7.32	12.80	8.54	17.07	10.37	<b>4.88</b>	10.37	6.10
Synonym Insertion	RP <sub>5</sub> @1↑	0.183	0.104	0.159	0.128	<b>0.226</b>	0.128	0.067	0.104	0.079
	RD <sub>5</sub> @1(%)↓	<b>21.05</b>	26.09	39.53	34.38	26.00	34.38	35.29	32.00	35.00
	RR <sub>5</sub> @1(%)↓	7.32	<b>4.88</b>	14.63	8.54	15.85	9.15	6.10	9.15	5.49
Synonym Substitution	RP <sub>5</sub> @1↑	0.146	0.091	0.159	0.104	<b>0.201</b>	0.140	0.073	0.079	0.061
	RD <sub>5</sub> @1(%)↓	36.84	34.78	39.53	46.88	34.00	<b>28.12</b>	29.41	48.00	50.00
	RR <sub>5</sub> @1(%)↓	10.37	6.71	17.07	10.98	15.24	7.93	<b>4.88</b>	9.76	6.71
TenseTransformation Past	RP <sub>5</sub> @1↑	0.250	0.146	0.238	0.189	<b>0.305</b>	0.171	0.110	0.134	0.110
	RD <sub>5</sub> @1(%)↓	<b>-7.89</b>	-4.35	9.30	3.13	0.00	12.50	-5.88	12.00	10.00
	RR <sub>5</sub> @1(%)↓	3.05	1.83	6.10	5.49	7.32	2.44	1.83	1.83	<b>1.22</b>
TenseTransformation Future	RP <sub>5</sub> @1↑	0.238	0.122	0.250	0.183	<b>0.311</b>	0.171	0.085	0.146	0.110
	RD <sub>5</sub> @1(%)↓	<b>-2.63</b>	13.04	4.65	6.25	-2.00	12.50	17.65	4.00	10.00
	RR <sub>5</sub> @1(%)↓	4.27	4.27	4.88	4.88	6.71	3.66	1.83	1.83	<b>1.22</b>
Whitespace Perturbation	RP <sub>5</sub> @1↑	0.146	0.085	0.146	0.122	<b>0.177</b>	0.122	0.073	0.091	0.049
	RD <sub>5</sub> @1(%)↓	36.84	39.13	44.19	37.50	42.00	37.50	<b>29.41</b>	40.00	60.00
	RR <sub>5</sub> @1(%)↓	14.02	9.76	15.85	9.76	22.56	10.37	<b>6.10</b>	10.98	7.32

Table 11: Robustness evaluation for each type of docstring perturbations on HumanEval.

MBPP	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Nominal	RP <sub>5</sub> @1↑	0.317	0.191	0.361	0.221	<b>0.407</b>	0.241	0.128	0.199	0.133
BackTranslation	RP <sub>5</sub> @1↑	0.304	0.186	0.360	0.222	<b>0.387</b>	0.230	0.119	0.177	0.128
	RD <sub>5</sub> @1(%)↓	4.21	2.69	0.28	<b>-0.47</b>	4.80	4.68	7.20	11.34	3.85
	RR <sub>5</sub> @1(%)↓	6.26	6.06	7.91	6.06	6.26	7.08	<b>4.00</b>	5.95	5.44
ButterFingers Perturbation	RP <sub>5</sub> @1↑	0.210	0.092	0.240	0.100	<b>0.280</b>	0.126	0.044	0.082	0.057
	RD <sub>5</sub> @1(%)↓	33.66	51.61	33.52	54.88	<b>31.06</b>	47.66	65.60	58.76	56.92
	RR <sub>5</sub> @1(%)↓	20.43	21.66	23.72	22.07	25.26	23.31	<b>14.48</b>	20.12	16.32
ChangeCharCase	RP <sub>5</sub> @1↑	0.187	0.087	0.220	0.105	<b>0.266</b>	0.124	0.053	0.074	0.055
	RD <sub>5</sub> @1(%)↓	41.10	54.30	39.20	52.56	<b>34.60</b>	48.51	58.40	62.89	58.46
	RR <sub>5</sub> @1(%)↓	22.07	20.74	27.21	20.84	26.28	25.87	<b>13.24</b>	21.46	17.45
EnglishInflectional Variation	RP <sub>5</sub> @1↑	0.306	0.161	0.334	0.198	<b>0.399</b>	0.214	0.103	0.179	0.113
	RD <sub>5</sub> @1(%)↓	3.56	15.59	7.67	10.23	<b>1.77</b>	11.49	20.00	10.31	15.38
	RR <sub>5</sub> @1(%)↓	8.93	10.78	11.40	9.65	10.68	12.53	7.08	9.45	<b>6.78</b>
SwapCharacters Perturbation	RP <sub>5</sub> @1↑	0.232	0.115	0.266	0.123	<b>0.304</b>	0.149	0.059	0.108	0.063
	RD <sub>5</sub> @1(%)↓	26.86	39.78	26.42	44.19	<b>25.25</b>	38.30	54.40	45.88	53.08
	RR <sub>5</sub> @1(%)↓	16.53	15.81	19.61	18.99	20.84	20.43	<b>12.42</b>	14.99	15.30
Synonym Insertion	RP <sub>5</sub> @1↑	0.238	0.111	0.263	0.103	<b>0.290</b>	0.121	0.052	0.101	0.055
	RD <sub>5</sub> @1(%)↓	<b>24.92</b>	41.94	27.27	53.49	28.79	49.79	59.20	49.48	58.46
	RR <sub>5</sub> @1(%)↓	16.63	18.99	21.25	20.53	24.44	24.85	<b>13.14</b>	17.56	14.99
Synonym Substitution	RP <sub>5</sub> @1↑	0.193	0.099	0.213	0.079	<b>0.233</b>	0.092	0.027	0.064	0.031
	RD <sub>5</sub> @1(%)↓	<b>39.16</b>	48.39	41.19	64.19	42.68	61.70	79.20	68.04	76.92
	RR <sub>5</sub> @1(%)↓	22.79	18.17	27.31	23.61	30.18	26.90	<b>16.22</b>	22.38	17.45
TenseTransformation Past	RP <sub>5</sub> @1↑	0.318	0.190	0.362	0.214	<b>0.402</b>	0.238	0.120	0.197	0.141
	RD <sub>5</sub> @1(%)↓	-0.32	0.54	-0.28	3.26	1.01	1.28	6.40	1.03	<b>-5.38</b>
	RR <sub>5</sub> @1(%)↓	2.16	2.36	4.00	3.18	3.29	2.16	1.64	2.26	<b>1.54</b>
TenseTransformation Future	RP <sub>5</sub> @1↑	0.314	0.197	0.369	0.218	<b>0.400</b>	0.242	0.122	0.185	0.125
	RD <sub>5</sub> @1(%)↓	0.97	<b>-3.23</b>	-1.99	1.40	1.52	-0.43	4.80	7.22	6.15
	RR <sub>5</sub> @1(%)↓	3.18	3.29	5.65	4.21	4.52	4.00	2.46	3.49	<b>2.26</b>
Whitespace Perturbation	RP <sub>5</sub> @1↑	0.214	0.107	0.252	0.106	<b>0.287</b>	0.134	0.057	0.094	0.054
	RD <sub>5</sub> @1(%)↓	32.69	44.09	30.40	52.09	<b>29.29</b>	44.26	55.20	52.58	59.23
	RR <sub>5</sub> @1(%)↓	20.64	17.66	21.56	20.64	25.46	23.31	<b>12.42</b>	17.86	16.02

Table 12: Robustness evaluation for each type of docstring perturbations on MBPP.

HumanEval	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Nominal	RP <sub>5</sub> @1↑	0.232	0.140	0.262	0.195	<b>0.305</b>	0.195	0.104	0.152	0.122
CamelCase	RP <sub>5</sub> @1↑	0.238	0.140	0.256	0.201	<b>0.293</b>	0.165	0.098	0.152	0.116
	RD <sub>5</sub> @1(%)↓	-2.63	0.00	2.33	<b>-3.13</b>	4.00	15.62	5.88	0.00	5.00
	RR <sub>5</sub> @1(%)↓	1.83	1.22	3.05	3.05	3.66	3.05	3.05	1.22	<b>0.61</b>
ButterFinger	RP <sub>5</sub> @1↑	0.195	0.104	0.232	0.177	<b>0.274</b>	0.159	0.098	0.140	0.091
	RD <sub>5</sub> @1(%)↓	15.79	26.09	11.63	9.38	10.00	18.75	<b>5.88</b>	8.00	25.00
	RR <sub>5</sub> @1(%)↓	4.88	4.88	9.76	4.88	9.15	3.66	3.05	<b>2.44</b>	3.05
SwapChar	RP <sub>5</sub> @1↑	0.226	0.116	0.226	0.177	<b>0.299</b>	0.183	0.073	0.146	0.116
	RD <sub>5</sub> @1(%)↓	2.63	17.39	13.95	9.38	<b>2.00</b>	6.25	29.41	4.00	5.00
	RR <sub>5</sub> @1(%)↓	3.05	3.05	4.88	4.27	4.88	2.44	3.05	2.44	<b>0.61</b>
ChangeCharCase	RP <sub>5</sub> @1↑	0.207	0.122	0.213	0.140	<b>0.256</b>	0.146	0.098	0.152	0.091
	RD <sub>5</sub> @1(%)↓	10.53	13.04	18.60	28.12	16.00	25.00	5.88	<b>0.00</b>	25.00
	RR <sub>5</sub> @1(%)↓	7.32	5.49	10.37	7.93	10.98	4.88	<b>4.27</b>	7.32	5.49
Inflectional Variation	RP <sub>5</sub> @1↑	0.232	0.134	0.262	0.195	<b>0.305</b>	0.201	0.110	0.128	0.110
	RD <sub>5</sub> @1(%)↓	0.00	4.35	0.00	0.00	0.00	-3.13	<b>-5.88</b>	16.00	10.00
	RR <sub>5</sub> @1(%)↓	3.66	3.05	4.27	3.66	2.44	<b>0.61</b>	1.83	2.44	1.22
Synonym Substitution	RP <sub>5</sub> @1↑	0.195	0.098	0.232	0.159	<b>0.305</b>	0.159	0.085	0.128	0.098
	RD <sub>5</sub> @1(%)↓	15.79	30.43	11.63	18.75	<b>0.00</b>	18.75	17.65	16.00	20.00
	RR <sub>5</sub> @1(%)↓	7.32	6.71	7.93	6.10	7.32	3.66	3.05	4.88	<b>2.44</b>

Table 13: Robustness evaluation for each type of function name perturbations on HumanEval.



MBPP	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Nominal	RP <sub>5</sub> @1↑	0.317	0.191	0.361	0.221	<b>0.407</b>	0.241	0.128	0.199	0.133
CamelCase	RP <sub>5</sub> @1↑	0.316	0.196	0.367	0.219	<b>0.408</b>	0.245	0.116	0.194	0.134
	RD <sub>5</sub> @1(%)↓	0.32	<b>-2.69</b>	-1.42	0.93	-0.25	-1.70	9.60	2.58	-0.77
	RR <sub>5</sub> @1(%)↓	5.44	5.44	7.29	5.34	7.08	4.52	5.75	5.03	<b>3.18</b>
ButterFinger	RP <sub>5</sub> @1↑	0.312	0.185	0.370	0.203	<b>0.412</b>	0.231	0.110	0.175	0.117
	RD <sub>5</sub> @1(%)↓	1.62	3.23	<b>-2.27</b>	7.91	-1.26	4.26	14.40	12.37	12.31
	RR <sub>5</sub> @1(%)↓	7.19	8.62	9.65	10.99	8.73	9.86	<b>6.67</b>	8.11	6.98
SwapChar	RP <sub>5</sub> @1↑	0.309	0.189	0.342	0.202	<b>0.399</b>	0.237	0.116	0.171	0.113
	RD <sub>5</sub> @1(%)↓	2.59	<b>1.08</b>	5.40	8.37	1.77	1.70	9.60	13.92	15.38
	RR <sub>5</sub> @1(%)↓	4.41	4.52	7.29	6.88	6.06	4.52	<b>3.18</b>	5.24	4.21
ChangeCharCase	RP <sub>5</sub> @1↑	0.295	0.179	0.346	0.192	<b>0.400</b>	0.244	0.093	0.171	0.111
	RD <sub>5</sub> @1(%)↓	7.12	6.45	4.26	13.02	1.52	<b>-1.28</b>	27.20	13.92	16.92
	RR <sub>5</sub> @1(%)↓	9.55	10.88	11.91	12.22	12.73	9.75	<b>8.32</b>	10.57	9.45
Inflectional Variation	RP <sub>5</sub> @1↑	0.318	0.187	0.343	0.202	<b>0.402</b>	0.243	0.128	0.188	0.125
	RD <sub>5</sub> @1(%)↓	-0.32	2.15	5.11	8.37	1.01	<b>-0.85</b>	0.00	5.67	6.15
	RR <sub>5</sub> @1(%)↓	3.08	4.31	6.88	5.75	5.95	4.31	<b>2.46</b>	2.98	3.49
Synonym Substitution	RP <sub>5</sub> @1↑	0.316	0.186	0.346	0.197	<b>0.384</b>	0.243	0.105	0.164	0.117
	RD <sub>5</sub> @1(%)↓	0.32	2.69	4.26	10.70	5.56	<b>-0.85</b>	18.40	17.53	12.31
	RR <sub>5</sub> @1(%)↓	6.88	7.49	10.88	10.47	9.96	9.86	7.70	8.52	<b>6.88</b>

Table 14: Robustness evaluation for each type of function name perturbations on MBPP.

HumanEval	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Nominal	RP <sub>5</sub> @1↑	0.402	0.293	0.518	0.366	<b>0.549</b>	0.390	0.189	0.323	0.250
DeadCodeInsertor	RP <sub>5</sub> @1↑	0.116	0.079	0.152	0.110	<b>0.159</b>	0.091	0.055	0.079	0.079
	RD <sub>5</sub> @1(%)↓	71.21	72.92	70.59	70.00	71.11	76.56	70.97	75.47	<b>68.29</b>
	RR <sub>5</sub> @1(%)↓	37.80	30.49	41.46	32.93	45.12	37.20	<b>17.07</b>	30.49	27.44
ForWhile TransformerFirst	RP <sub>5</sub> @1↑	0.384	0.226	0.500	0.305	<b>0.537</b>	0.384	0.159	0.280	0.213
	RD <sub>5</sub> @1(%)↓	4.55	22.92	3.53	16.67	2.22	<b>1.56</b>	16.13	13.21	14.63
	RR <sub>5</sub> @1(%)↓	<b>5.49</b>	6.71	9.15	8.54	6.10	<b>5.49</b>	<b>5.49</b>	6.71	9.76
OperandSwap	RP <sub>5</sub> @1↑	0.402	0.274	0.500	0.348	<b>0.512</b>	0.354	0.171	0.311	0.220
	RD <sub>5</sub> @1(%)↓	<b>0.00</b>	6.25	3.53	5.00	6.67	9.38	9.68	3.77	12.20
	RR <sub>5</sub> @1(%)↓	6.71	<b>4.27</b>	6.71	6.10	5.49	6.71	6.10	7.93	7.32
VarRenamerCB	RP <sub>5</sub> @1↑	0.415	0.268	0.476	0.329	<b>0.518</b>	0.354	0.146	0.287	0.238
	RD <sub>5</sub> @1(%)↓	<b>-3.03</b>	8.33	8.24	10.00	5.56	9.38	22.58	11.32	4.88
	RR <sub>5</sub> @1(%)↓	<b>4.88</b>	6.10	6.71	8.54	5.49	7.32	7.93	8.54	4.88
VarRenamerNaive	RP <sub>5</sub> @1↑	0.396	0.244	0.482	0.348	<b>0.494</b>	0.341	0.177	0.280	0.220
	RD <sub>5</sub> @1(%)↓	<b>1.52</b>	16.67	7.06	5.00	10.00	12.50	6.45	13.21	12.20
	RR <sub>5</sub> @1(%)↓	<b>4.27</b>	9.76	7.32	9.15	6.71	8.54	9.76	10.37	5.49
VarRenamerRN	RP <sub>5</sub> @1↑	0.366	0.207	0.421	0.280	<b>0.470</b>	0.280	0.085	0.152	0.177
	RD <sub>5</sub> @1(%)↓	<b>9.09</b>	29.17	18.82	23.33	14.44	28.12	54.84	52.83	29.27
	RR <sub>5</sub> @1(%)↓	12.20	14.63	14.02	12.80	<b>11.59</b>	17.07	16.46	24.39	12.20

Table 15: Robustness evaluation for each type of code syntax perturbations on HumanEval.

MBPP	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Nominal	RP <sub>5</sub> @1↑	0.450	0.285	0.535	0.331	<b>0.571</b>	0.379	0.219	0.292	0.176
DeadCodeInserter	RP <sub>5</sub> @1↑	0.043	0.020	0.044	0.024	<b>0.055</b>	0.025	0.015	0.015	0.009
	RD <sub>5</sub> @1(%)↓	90.41	93.17	91.75	92.86	<b>90.29</b>	93.50	92.96	94.72	94.74
	RR <sub>5</sub> @1(%)↓	52.05	37.99	57.39	39.12	60.57	44.87	29.26	37.78	<b>24.95</b>
ForWhile TransformerFirst	RP <sub>5</sub> @1↑	0.432	0.259	0.497	0.303	<b>0.532</b>	0.346	0.182	0.245	0.149
	RD <sub>5</sub> @1(%)↓	<b>3.88</b>	9.35	7.10	8.39	6.83	8.67	16.90	15.85	15.20
	RR <sub>5</sub> @1(%)↓	13.66	12.94	11.60	13.45	11.70	13.35	12.73	16.53	<b>9.24</b>
OperandSwap	RP <sub>5</sub> @1↑	0.450	0.275	0.506	0.321	<b>0.544</b>	0.379	0.225	0.276	0.211
	RD <sub>5</sub> @1(%)↓	0.00	3.60	5.37	2.80	4.68	0.00	-2.82	5.28	<b>-20.47</b>
	RR <sub>5</sub> @1(%)↓	13.24	11.81	<b>10.57</b>	13.45	11.81	12.32	12.32	15.50	11.91
VarRenamerCB	RP <sub>5</sub> @1↑	0.428	0.263	0.475	0.307	<b>0.511</b>	0.359	0.194	0.247	0.207
	RD <sub>5</sub> @1(%)↓	4.79	7.91	11.13	7.14	10.43	5.15	11.27	15.14	<b>-18.13</b>
	RR <sub>5</sub> @1(%)↓	15.30	13.96	15.20	15.50	14.17	14.07	13.14	16.12	<b>12.83</b>
VarRenamerNaive	RP <sub>5</sub> @1↑	0.417	0.240	0.461	0.286	<b>0.513</b>	0.338	0.171	0.226	0.172
	RD <sub>5</sub> @1(%)↓	7.31	15.83	13.82	13.35	10.07	10.84	21.60	22.54	<b>1.75</b>
	RR <sub>5</sub> @1(%)↓	16.63	15.61	17.04	17.97	14.78	16.43	14.17	18.07	<b>13.24</b>
VarRenamerRN	RP <sub>5</sub> @1↑	0.355	0.191	0.405	0.205	<b>0.426</b>	0.259	0.114	0.168	0.114
	RD <sub>5</sub> @1(%)↓	<b>21.00</b>	33.09	24.38	37.89	25.36	31.71	47.89	42.25	35.09
	RR <sub>5</sub> @1(%)↓	22.90	22.90	23.82	26.59	24.95	26.28	23.82	25.87	<b>19.40</b>

Table 16: Robustness evaluation for each type of code syntax perturbations on MBPP.

HumanEval	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Nominal	RP <sub>5</sub> @1↑	0.402	0.293	0.518	0.366	<b>0.549</b>	0.390	0.189	0.323	0.250
Tab-Indent	RP <sub>5</sub> @1↑	0.415	0.305	0.518	0.354	<b>0.561</b>	0.396	0.146	0.299	0.244
	RD <sub>5</sub> @1(%)↓	-3.03	<b>-4.17</b>	0.00	3.33	-2.22	-1.56	22.58	7.55	2.44
	RR <sub>5</sub> @1(%)↓	<b>3.66</b>	4.88	8.54	4.88	3.66	4.27	7.93	9.76	7.93
Line Split	RP <sub>5</sub> @1↑	0.384	0.274	0.500	0.378	<b>0.524</b>	0.390	0.171	0.305	0.244
	RD <sub>5</sub> @1(%)↓	4.55	6.25	3.53	<b>-3.33</b>	4.44	0.00	9.68	5.66	2.44
	RR <sub>5</sub> @1(%)↓	3.05	4.27	4.27	4.88	3.66	<b>2.44</b>	3.05	6.71	4.27
Doc2Comments	RP <sub>5</sub> @1↑	0.335	0.287	0.433	0.293	<b>0.457</b>	0.335	0.146	0.293	0.195
	RD <sub>5</sub> @1(%)↓	16.67	<b>2.08</b>	16.47	20.00	16.67	14.06	22.58	9.43	21.95
	RR <sub>5</sub> @1(%)↓	11.59	5.49	14.63	8.54	12.80	7.93	<b>4.27</b>	5.49	10.37
NewlineRandom	RP <sub>5</sub> @1↑	0.360	0.220	0.390	0.250	<b>0.457</b>	0.299	0.152	0.232	0.171
	RD <sub>5</sub> @1(%)↓	<b>10.61</b>	25.00	24.71	31.67	16.67	23.44	19.35	28.30	31.71
	RR <sub>5</sub> @1(%)↓	12.20	10.98	17.68	15.85	11.59	13.41	<b>7.32</b>	15.85	12.20
NewlineAfterCode	RP <sub>5</sub> @1↑	0.409	0.262	0.494	0.311	<b>0.537</b>	0.335	0.165	0.287	0.183
	RD <sub>5</sub> @1(%)↓	<b>-1.52</b>	10.42	4.71	15.00	2.22	14.06	12.90	11.32	26.83
	RR <sub>5</sub> @1(%)↓	6.71	7.93	8.54	9.15	<b>3.66</b>	7.93	4.88	8.54	9.15
NewlineAfterDoc	RP <sub>5</sub> @1↑	0.396	0.274	0.518	0.348	<b>0.549</b>	0.384	0.183	0.311	0.244
	RD <sub>5</sub> @1(%)↓	1.52	6.25	0.00	5.00	<b>0.00</b>	1.56	3.23	3.77	2.44
	RR <sub>5</sub> @1(%)↓	4.27	4.27	6.10	4.27	1.22	1.83	<b>0.61</b>	3.66	4.27

Table 17: Robustness evaluation for each type of code format perturbations on HumanEval.

MBPP	Metric	CodeGen 2B mono	CodeGen 2B multi	CodeGen 6B mono	CodeGen 6B multi	CodeGen 16B mono	CodeGen 16B multi	InCoder 1B	InCoder 6B	GPT-J 6B
Nominal	RP <sub>5</sub> @1↑	0.450	0.285	0.535	0.331	<b>0.571</b>	0.379	0.219	0.292	0.176
Tab-Indent	RP <sub>5</sub> @1↑	0.452	0.302	0.530	0.339	<b>0.566</b>	0.385	0.208	0.325	0.176
	RD <sub>5</sub> @1(%)↓	-0.46	-5.76	0.96	-2.48	0.90	-1.63	4.69	<b>-11.62</b>	0.00
	RR <sub>5</sub> @1(%)↓	6.37	6.98	<b>5.85</b>	8.01	6.88	6.78	9.24	12.22	7.60
Line Split	RP <sub>5</sub> @1↑	0.445	0.275	0.524	0.326	<b>0.556</b>	0.378	0.187	0.283	0.163
	RD <sub>5</sub> @1(%)↓	1.14	3.60	2.11	1.24	2.52	<b>0.27</b>	14.55	2.82	7.02
	RR <sub>5</sub> @1(%)↓	4.41	6.37	4.41	6.16	5.54	6.26	6.06	6.78	<b>3.90</b>
Doc2Comments	RP <sub>5</sub> @1↑	0.435	0.269	0.476	0.299	<b>0.529</b>	0.342	0.169	0.264	0.172
	RD <sub>5</sub> @1(%)↓	3.20	5.76	10.94	9.63	7.37	9.76	22.54	9.51	<b>1.75</b>
	RR <sub>5</sub> @1(%)↓	<b>6.16</b>	8.62	8.32	9.14	8.93	11.29	7.19	8.32	6.47
NewlineRandom	RP <sub>5</sub> @1↑	0.375	0.181	0.335	0.198	<b>0.470</b>	0.262	0.123	0.159	0.104
	RD <sub>5</sub> @1(%)↓	<b>16.67</b>	36.69	37.43	40.06	17.63	30.89	43.66	45.42	40.94
	RR <sub>5</sub> @1(%)↓	12.94	16.32	23.72	19.40	15.81	17.56	12.73	16.63	<b>10.37</b>
NewlineAfterCode	RP <sub>5</sub> @1↑	0.406	0.238	0.379	0.240	<b>0.525</b>	0.291	0.165	0.207	0.150
	RD <sub>5</sub> @1(%)↓	9.82	16.55	29.17	27.33	<b>8.09</b>	23.31	24.41	28.87	14.62
	RR <sub>5</sub> @1(%)↓	9.14	10.27	19.51	13.55	10.99	14.17	9.03	12.73	<b>7.91</b>
NewlineAfterDoc	RP <sub>5</sub> @1↑	0.449	0.274	0.518	0.305	<b>0.570</b>	0.378	0.180	0.242	0.153
	RD <sub>5</sub> @1(%)↓	0.23	3.96	3.07	7.76	<b>0.18</b>	0.27	17.84	16.90	12.87
	RR <sub>5</sub> @1(%)↓	<b>2.36</b>	4.41	4.11	7.08	4.62	4.41	4.72	6.57	4.31

Table 18: Robustness evaluation for each type of code format perturbations on MBPP.

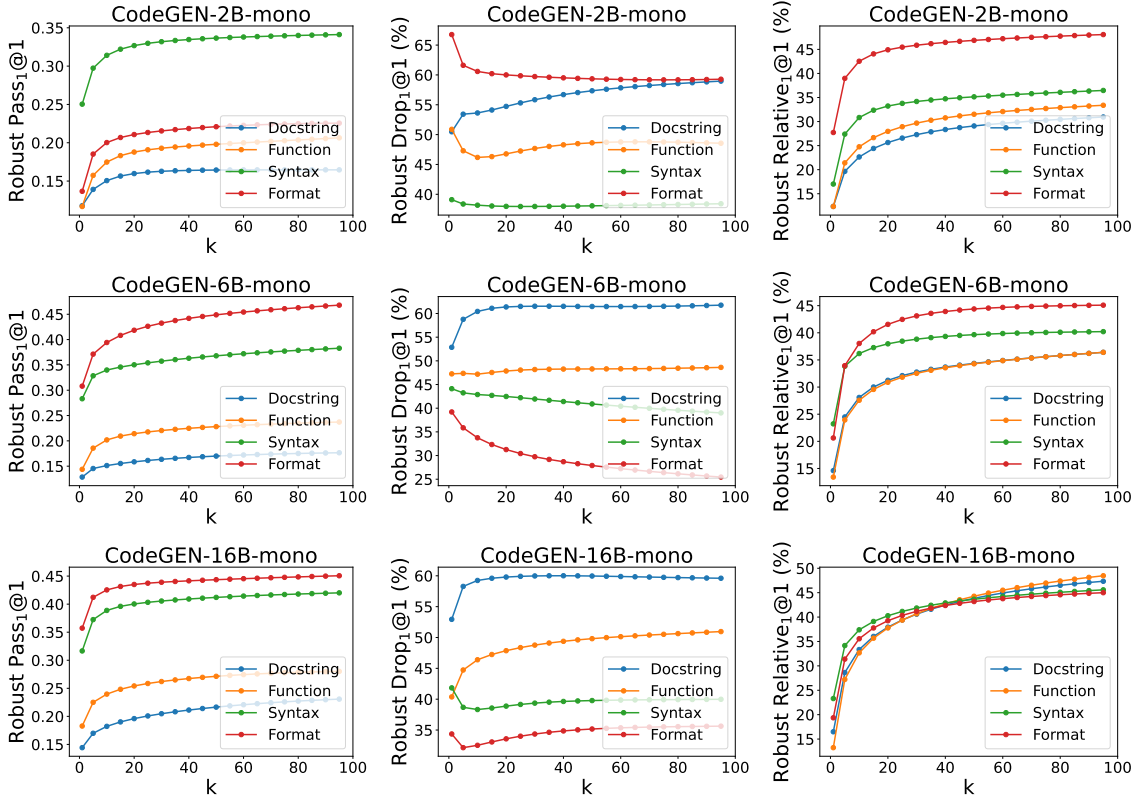


Figure 22: Robust Drop<sub>1</sub>@1 and Robust Relative<sub>1</sub>@1 on CodeGen-16B-mono under different  $k$  using sampling  $n = 100$ . Robust Drop remains stable while Robust Pass and Robust Relative increases with  $k$ .

Category	Metric	$n = 1$	$n = 10$	$n = 100$
Docstring	Nominal $\uparrow$	0.287	0.308	0.306
	RP <sub>1</sub> @1 $\uparrow$	0.128	0.140	0.143
	RD <sub>1</sub> @1(%) $\downarrow$	55.32	54.46	53.34
	RR <sub>1</sub> @1(%) $\downarrow$	15.85	16.77	16.55
Function	Nominal $\uparrow$	0.287	0.308	0.306
	RP <sub>1</sub> @1 $\uparrow$	0.183	0.180	0.183
	RD <sub>1</sub> @1(%) $\downarrow$	36.17	41.39	40.37
	RR <sub>1</sub> @1(%) $\downarrow$	10.37	12.99	13.24
Syntax	Nominal $\uparrow$	0.561	0.542	0.544
	RP <sub>1</sub> @1 $\uparrow$	0.220	0.234	0.244
	RD <sub>1</sub> @1(%) $\downarrow$	60.87	56.81	55.19
	RR <sub>1</sub> @1(%) $\downarrow$	34.15	31.04	30.39
Format	Nominal $\uparrow$	0.561	0.542	0.544
	RP <sub>1</sub> @1 $\uparrow$	0.341	0.352	0.357
	RD <sub>1</sub> @1(%) $\downarrow$	39.13	34.98	34.36
	RR <sub>1</sub> @1(%) $\downarrow$	21.95	19.70	19.36

Table 19: Generation variances for three robustness metrics with different sampling  $n$  on CodeGen-16B-mono and HumanEval.



## ACL 2023 Responsible NLP Checklist

---

### A For every submission:

- A1. Did you describe the limitations of your work?  
*Left blank.*
- A2. Did you discuss any potential risks of your work?  
*Left blank.*
- A3. Do the abstract and introduction summarize the paper’s main claims?  
*Left blank.*
- A4. Have you used AI writing assistants when working on this paper?  
*Left blank.*

### B Did you use or create scientific artifacts?

*Left blank.*

- B1. Did you cite the creators of artifacts you used?  
*Left blank.*
- B2. Did you discuss the license or terms for use and / or distribution of any artifacts?  
*Left blank.*
- B3. Did you discuss if your use of existing artifact(s) was consistent with their intended use, provided that it was specified? For the artifacts you create, do you specify intended use and whether that is compatible with the original access conditions (in particular, derivatives of data accessed for research purposes should not be used outside of research contexts)?  
*Left blank.*
- B4. Did you discuss the steps taken to check whether the data that was collected / used contains any information that names or uniquely identifies individual people or offensive content, and the steps taken to protect / anonymize it?  
*Left blank.*
- B5. Did you provide documentation of the artifacts, e.g., coverage of domains, languages, and linguistic phenomena, demographic groups represented, etc.?  
*Left blank.*
- B6. Did you report relevant statistics like the number of examples, details of train / test / dev splits, etc. for the data that you used / created? Even for commonly-used benchmark datasets, include the number of examples in train / validation / test splits, as these provide necessary context for a reader to understand experimental results. For example, small differences in accuracy on large test sets may be significant, while on small test sets they may not be.  
*Left blank.*

### C Did you run computational experiments?

*Left blank.*

- C1. Did you report the number of parameters in the models used, the total computational budget (e.g., GPU hours), and computing infrastructure used?  
*Left blank.*

*The Responsible NLP Checklist used at ACL 2023 is adopted from NAACL 2022, with the addition of a question on AI writing assistance.*

- C2. Did you discuss the experimental setup, including hyperparameter search and best-found hyperparameter values?

*Left blank.*

- C3. Did you report descriptive statistics about your results (e.g., error bars around results, summary statistics from sets of experiments), and is it transparent whether you are reporting the max, mean, etc. or just a single run?

*Left blank.*

- C4. If you used existing packages (e.g., for preprocessing, for normalization, or for evaluation), did you report the implementation, model, and parameter settings used (e.g., NLTK, Spacy, ROUGE, etc.)?

*Left blank.*

**D**  **Did you use human annotators (e.g., crowdworkers) or research with human participants?**

*Left blank.*

- D1. Did you report the full text of instructions given to participants, including e.g., screenshots, disclaimers of any risks to participants or annotators, etc.?

*Left blank.*

- D2. Did you report information about how you recruited (e.g., crowdsourcing platform, students) and paid participants, and discuss if such payment is adequate given the participants' demographic (e.g., country of residence)?

*Left blank.*

- D3. Did you discuss whether and how consent was obtained from people whose data you're using/curating? For example, if you collected data via crowdsourcing, did your instructions to crowdworkers explain how the data would be used?

*Left blank.*

- D4. Was the data collection protocol approved (or determined exempt) by an ethics review board?

*Left blank.*

- D5. Did you report the basic demographic and geographic characteristics of the annotator population that is the source of the data?

*Left blank.*