

---

# CodeSCM: Causal Analysis for Multi-Modal Code Generation

---

Mukur Gupta,\* Noopur Bhatt,\* and Suman Jana  
Columbia University  
{mukur.gupta, noopur.bhatt}@columbia.edu  
suman@cs.columbia.edu

## Abstract

In this paper, we propose CodeSCM, a Structural Causal Model (SCM) for analyzing multi-modal code generation using large language models (LLMs). By applying interventions to CodeSCM, we define the causal effects of different prompt modalities, such as natural language, code, and input-output examples, on the model. CodeSCM introduces latent mediator variables to separate the code and natural language semantics of a multi-modal prompt. Using the principles of Causal Mediation Analysis on these mediators, we define direct effects through targeted interventions, quantitatively representing the model’s spurious leanings. We find that, in addition to natural language instructions, input-output examples significantly influence model generation, and total causal effects evaluations from CodeSCM also reveal the memorization of code generation benchmarks.

## 1 Introduction

Modern Large Language Models (LLMs) have shown remarkable effectiveness in code reasoning tasks, particularly code generation [29, 37, 4]. This task involves generating code that meets specific multi-modal requirements, constrained by natural language instructions, code snippets, and input-output example pairs [8, 2]. Furthermore, some multi-modal prompt components contain information from both code and natural language modalities [6], such as function signatures and variable names, where code structure and natural language coexist. This enriched coding context, combining programming and natural language semantics, helps LLMs better understand both the semantics and syntactic requirements of the desired code.

Prior research has shown the effectiveness of prompt tuning in improving generation performance [5, 25, 44]. These works have shown that multi-modal prompts can be highly sensitive, where small adjustments might result in drastically different responses from the model [7, 51, 38]. However, the interactions between the multi-modal components of code generation prompts and how they directly or indirectly affect the generated code, are still not well understood. In this paper, we systematically explore these complex multi-modal effects using a causal approach. We propose a novel causal framework, CodeSCM, to measure the causal effects of different modalities in the prompt on the performance of code generation LLMs. CodeSCM defines a Structural Causal Model [33], shown in Figure 1, where each modal component of the prompt is treated as an independent variable that causally affects the code generated by the model. To account for the semantics of natural language and code, we introduce two latent mediator variables to capture the code semantics and natural language semantics of the input prompt, mimicking a human mental model to generate correct code snippets given a multi-modal input problem.

Specifically, we make three key contributions in this paper: (i) we introduce CodeSCM, a novel framework for causal inference in multi-modal code generation tasks, enhancing interpretability and

---

\*Denotes equal contribution

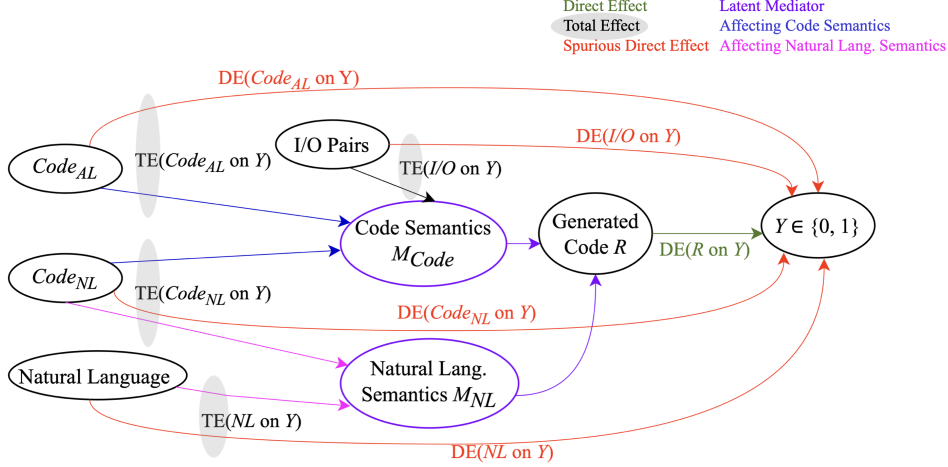


Figure 1: CodeSCM causal graph representing the total and direct effects of the modal variable nodes on the response variable  $Y$  representing the correctness of the generated code.

causal understanding of codeLLMs. While CodeSCM is designed for the code generation task in this paper, it can be extended to other modalities, tasks, and transformations. (ii) using CodeSCM, we define the Total Effects of the modalities on code generation, highlighting that I/O pairs and natural language code components, like function headers, are significant modal components alongside natural language instructions. Our Total Effect analysis also provides strong evidence of benchmark memorization in LLMs like GPT-4T, and (iii) through targeted interventions on CodeSCM, we measure the Direct Effects of each modality, representing spurious model correlations, and show that simple semantics-preserving transformations to input-output example pairs lead to a significant drop in accuracy.

## 2 Background

**Structural Causal Model.** A Structural Causal Model (SCM)  $\mathcal{M}$  [33] is defined by a 4-tuple  $\mathcal{M} = \langle \mathbf{U}, \mathbf{V}, \mathbf{F}, P(\mathbf{U}) \rangle$ , where  $\mathbf{U}$  represents a set of exogenous variables,  $P(\mathbf{U})$  is a joint probability distribution over the set  $\mathbf{U}$ ,  $\mathbf{V}$  is a set of endogenous variables determined by  $\mathbf{U} \cup \mathbf{V}$ , and  $\mathbf{F}$  is a set of functions from  $\mathbf{U} \cup \mathbf{V}$  to  $\mathbf{V}$ . The SCM  $\mathcal{M}$  can be represented by a causal graph  $\mathcal{G}$ , which employs nodes to represent both exogenous and endogenous variables. Causal effects on any response variable  $Y \in \mathbf{V}$  are quantitatively measured using interventions. An intervention on  $X \in \mathbf{V}$ , represented by  $do(x)$ , creates a sub-SCM  $\mathcal{M}^x = \langle \mathbf{U}, \mathbf{V}, \mathbf{F}_x, P(\mathbf{U}) \rangle$  where  $\mathbf{F}_x$  represents a subset of function mappings in  $\mathbf{F}$  which do not have  $X$  in their co-domain, and  $X$  is replaced by a constant  $x$ . Following the above notation, we can formally define causal effects:

**Definition 2.1. (Total Effect)** [31]: The causal effect of two distinct realizations of the variable  $X$  with  $do(x')$  and  $do(x'')$ :

$$TE(x', x'') = \mathbb{E}[Y | do(X = x')] - \mathbb{E}[Y | do(X = x'')] \quad (1)$$

Causal Mediation Analysis [36, 32, 35] involves understanding the effects of a mediator  $M \in \mathbf{V}$  in explaining changes in  $Y$ . All the effects from  $X$  to  $Y$  where all  $Z \in \mathbf{V}$ , representing the parents of  $Y$  excluding  $X$ , remain fixed are called Direct Effects. We measure the direct effect of modalities to define spurious learnings that are not mediated by the latent mediators. We define the Path Effect that allows us to systematically measure the direct effect of a modality averaged across the complete dataset.

**Definition 2.2. (Path Effect)** [3, 48]: The causal effect of variable  $X$  along a path  $\alpha$  can be represented in an edge subgraph  $\mathcal{G}_\alpha$

$$PE_\alpha(x', x'') = \mathbb{E}[Y | Z_\alpha(do(x'')), Z_{\bar{\alpha}}(do(x'))] - \mathbb{E}[Y(x')] \quad (2)$$

where,  $Z_\alpha$  is set of all the mediators  $\in \mathcal{G}_\alpha$  and  $Z_{\bar{\alpha}}$  is complementary mediator set. Hence, all the variables on the path  $\alpha$  take values with  $do(x'')$  and other mediators which do not lie on  $\alpha$  take values with  $do(x')$ . We use the definition of Path Effect to quantify the Direct Effect, as it is a special case of the Path Effect.

<pre>def noprofit_noloss(actual_cost, sale_amount):     if(sale_amount == actual_cost):         return True     else:         return False</pre>	<pre>def func_noprofit_noloss(actual_cost, sale_amount):</pre>	do( $X_{C_{NL}} = 1$ )
<p>Write a function to check whether the given amount has no profit and no loss</p>	<pre>def noprofit_noloss(actual_cost, sale_amount):     if False:         x=[_ for i in range(42)]</pre>	do( $X_{C_{AL}} = 1$ )
<pre>[ "assert noprofit_noloss(1500,1200)==False",   "assert noprofit_noloss(100,100)==True", "assert   noprofit_noloss(2000,5000)==False" ]</pre>	<p><b>DOCSTRING:</b> Write a function to check whether the given amount has no profit and no loss</p>	do( $X_{NL} = 1$ )
<pre>assert noprofit_noloss(1500,1200) &gt;= False assert noprofit_noloss(1500,1200) &lt;= False</pre>		do( $X_{IO} = 1$ )

Figure 2: Modalities in an example from mMBPP+ dataset. Original prompt breakdown (right) and semantics preserving transformations with respective interventions (left): **red** for natural language, **blue** for  $Code_{AL}$  and  $Code_{NL}$ , **orange** for transformed  $Code_{AL}$ , and **green** for Input/Output examples.

**Multi-Modal Code Generation.** Natural language instructions alone are often insufficient to meet strict context-based syntax requirements in code generation, such as variable or function names that dependent on surrounding code. Thus, natural language prompts are augmented with code modality, guiding the generation into appropriate syntactical space [8, 2]. Additionally, some prompt components appear as a single entity but contribute to multiple modalities in terms of model understanding, as observed by Casalnuovo et al. [6]. For instance, the function header name in Figure 2 is primarily a code component, but its natural language name also conveys information about the desired output. As highlighted later in Section 3, we call this a natural language channel of Code. Similarly, input-output (I/O) example pairs also carry information about code correctness and logic beyond the syntactical structure. We believe that future codeLLMs might rely heavily on these components to ensure the correctness of intermediary variables, akin to code debugging process of longer code fragments. Therefore, we consider I/O pairs and natural language channels of code as separate modalities, in addition to natural language instructions and code.

### 3 Problem setup

**CodeSCM.** Each prompt  $\mathcal{P}$  in dataset  $\mathcal{D}$  is decomposed into its multi-modal components, which are represented as variables in the structural causal model, as shown in Figure 1. We use the extended Backus–Naur form (Equation 7) to represent the multi-modal prompt. In CodeSCM, as shown in Equation 7, we consider four modalities: Natural Language ( $NL$ ), algorithmic channel of Code ( $Code_{AL}$ ), natural language channel of Code ( $Code_{NL}$ ), and input-output example pairs ( $I/O$ ). We define the multi-modal structural causal model (CodeSCM) to model the causal relationship between prompt modalities and the model-generated code. Since different code snippets and similar natural language texts can convey the same semantics for a human mental model, we introduce two latent mediators:  $M_{Code}$  for code semantics and  $M_{NL}$  for natural language semantics. Following the Causal Mediation Analysis, we assume each modality’s effect on the output is mediated through these variables. As shown in Figure 1,  $Code_{AL}$  and  $NL$  directly affect  $M_{Code}$  and  $M_{NL}$  respectively;  $I/O$  affects  $M_{Code}$ , and,  $Code_{NL}$  directly affects both mediators. The output code generated,  $R$ , is tested for correctness against the test cases, where code correctness is the response variable  $Y \in 0, 1$ , with  $\mathbb{E}(Y)$  representing accuracy over dataset  $\mathcal{D}$ .

**Modal Causal Effects and Interventions.** Using CodeSCM, we define the causal effects of each modal variable in  $\mathcal{P}$  on the generated code. We measure the Total Effect (Definition 2.1) of each modality on the response variable  $Y$ , reflecting the model’s sensitivity. Additionally, we examine the Direct Effect (Definition 2.2) of modalities on  $Y$  that bypass  $M_{Code}$  and  $M_{NL}$ , capturing spurious correlations learned during training. We also define additional variables and interventions to quantify these effects. Direct Effect (DE) and Total Effect (TE) for  $Code_{AL}$  are presented here and the detailed derivations for other modalities are in Appendix A.

**Causal effects of  $Code_{AL}$ .** We consider the  $Code_{AL}$  variable as an output of a structural equation  $F_C \in \mathbf{F}$  on  $C_{AL}$ ,  $C_{DC}$  and  $X_{AL}$ , shown in equation 3.  $C_{AL}$  is the actual prompt component  $\mathcal{P}_{Code_{AL}} \in \mathcal{D}$ . To measure the DE of code on  $Y$   $DE(Code_{AL} \text{ on } Y)$ , quantifying the spurious correlations, we vary  $Code_{AL}$  variable while keeping mediator  $M_C$  constant i.e  $M_C(Code_{AL}) = M_C(Code'_{AL})$ . We do this by inserting Dead Code (DC), a semantics-preserving transformation, represented by  $C_{DC}$ . We use the categorical variable  $X_{AL}$  to represent the interaction between the

Table 1: Statistics and prompt modalities of HumanEval+, MBPP+, mMBPP+, and CoderEval datasets.

Dataset	Size	NL	$Code_{AL}$	$Code_{NL}$	I/O Pairs
HumanEval+	164	✓	✓	✓	✓
MBPP+	399	✓	×	✓	✓
mMBPP+	399	✓	✓	✓	✓
CoderEval	460	✓	✓	✓	×
CoderEval-SCP	35	✓	✓	✓	×
CoderEval-SCJ	55	✓	✓	✓	×

actual code and the dead code, which also allows us to calculate causal effects over all prompts in  $\mathcal{D}$ .

$$Code_{AL} \leftarrow \mathbb{1}_{\{X_{AL}=1\}}(C_{AL} + C_{DC}) + \mathbb{1}_{\{X_{AL}=0\}}(C_{AL}) + \mathbb{1}_{\{X_{AL}=-1\}}(NULL) \quad (3)$$

where  $\mathbb{1}(\cdot)$  is an indicator function;  $(C_{AL} + C_{DC})$  represents the concatenation of a snippet of dead code with the actual code. We measure the TE of  $Code_{AL}$ ,  $TE(Code_{AL} \text{ on } Y)$ , by computing the expected change in the  $Y$  by setting the  $Code_{AL}$  component as NULL in the prompt using the variable  $X_{AL}$ .

$$\begin{aligned} TE(Code_{AL} \text{ on } Y) &= TE(do(X_{AL} = 0), do(X_{AL} = -1)) \\ &= \mathbb{E}[Y|do(X_{AL} = 0)] - \mathbb{E}[Y|do(X_{AL} = -1)] \quad (\text{Definition 2.1}) \\ &\stackrel{(i)}{=} \mathbb{P}[Y = 1|do(X_{AL} = 0)] - \mathbb{P}[Y = 1|do(X_{AL} = -1)] \\ &\stackrel{(ii)}{=} Acc(\mathcal{D}) - Acc(\mathcal{D}; \mathcal{P}_{Code_{AL}} = NULL) \end{aligned}$$

where, equality (i) follows because  $Y$  follows Bernoulli distribution;  $Acc(\mathcal{D})$  is the accuracy of the model over the dataset  $\mathcal{D}$ . The DE of  $Code_{AL}$  on  $Y$ ,  $DE(Code_{AL} \text{ on } Y)$  is measured by the expected change in  $Y$  with varying  $Code_{AL}$ , while keeping  $M_{Code}$  unchanged with dead code insertion. We calculate  $DE(Code_{AL} \text{ on } Y)$  using the Path Effect of  $X_{AL}$  on  $Y$ , along a path from  $X_{AL}$  to  $Y$  which goes through  $Code_{AL}$  but skips  $M_{Code}$ . Using Definition 2.2:

$$\begin{aligned} DE(Code_{AL} \text{ on } Y) &= \mathbb{E} \left[ Y \left( X_{AL} = 1, Code_{AL}(X_{AL} = 1), M_{Code}(Code_{AL}(X_{AL} = 0)) \right) \right] \\ &\quad - \mathbb{E} [Y(Code_{AL}(X_{AL} = 0))] \\ &\stackrel{(i)}{=} \mathbb{E}[Y|do(X_{AL} = 1)] - \mathbb{E}[Y|do(X_{AL} = 0)] \\ &\stackrel{(ii)}{=} Acc(\mathcal{D}) - Acc(\mathcal{D}; \mathcal{P}_{Code_{AL}} = C_{AL} + C_{DC}) \end{aligned}$$

where, equality (i) follows from  $M_{Code}(Code_{AL}(X_{AL} = 0)) = M_{Code}(Code_{AL}(X_{AL} = 1))$  where the dead code insertion in Equation 3 keeps the code semantics  $M_{Code}$  unchanged; equality (ii) is similar to equality (i) and (ii) used in TE.

**Causal Effects of Other Modalities.** Similarly, the  $NL$ ,  $I/O$ , and  $Code_{NL}$  modal variables are considered as outputs of structural Equations 4, 5 and 6 respectively. For  $Code_{NL}$ , the direct effect requires bypassing two mediators,  $M_{NL}$  and  $M_{Code}$ . Therefore, we define a transformation that preserves semantics for both. As seen in Figure 2, our transformation adds a prefix DN (Dead Name) to the function header, preserving semantics in both the natural language and code domains. For  $I/O$  transformations, each assertion equality is replaced by two inequalities ( $\geq$  and  $\leq$ ). While we demonstrate one specific transformation for each modality in our work to compute the respective Direct Effects, we note that CodeSCM can be extended to any other transformations, provided that i) the mediator variables remain unchanged, and ii) the transformations are independent of the input prompt. In addition to DE experiments in Section 4, we show DE computation with one additional transformation in Appendix D. We use simple prefix/suffix transformations to ensure independence between variables like  $S$  and  $DS$  or  $C$  and  $DC$ , to avoid correlation introduced by common transformations such as back-translation for  $NL$ .

Table 2: Total Effect (TE) and Direct Effect (DE) of modalities on code generation. Pass@1 accuracy on Full prompt for each model and dataset is reported, followed by accuracy drop, indicating TE or DE. "\*" denotes an increase in accuracy with the respective intervention. Bold highlights top TE and DE for each dataset and model. Accuracy results are averaged across three runs.

Model	Modality	HumanEval+		mMBPP+		CoderEval-SCP		CoderEval-SCJ		Mean	
		TE	DE	TE	DE	TE	DE	TE	DE	TE	DE
	Full	81.71		72.68		48.57		43.64		61.65	
GPT-4T	NL	<b>42.08</b>	1.22	19.05	4.26	<b>20.00</b>	<b>2.86*</b>	3.64	1.82	21.19	3.64
	<i>Code<sub>AL</sub></i>	1.83	1.22	1.25	4.01	8.57	0.00	<b>43.64</b>	<b>18.18*</b>	13.82	5.86
	<i>Code<sub>NL</sub></i>	18.91	1.83	<b>42.86</b>	2.76	0.00	<b>2.86*</b>	1.82	0.00	15.90	1.52
	I/O Pairs	5.49	<b>2.44</b>	12.28	<b>6.26</b>	N/A	N/A	N/A	N/A	8.89	4.35
	Full	53.05		52.63		37.14		47.27		47.52	
WizCoder	NL	<b>30.49</b>	5.49	<b>13.53</b>	0.50	5.71	<b>8.57*</b>	10.91	<b>3.64</b>	15.16	3.70
	<i>Code<sub>AL</sub></i>	4.27	9.76	2.00	<b>2.50</b>	2.86	2.86*	<b>45.45</b>	0.00	13.65	3.78
	<i>Code<sub>NL</sub></i>	6.10	2.44	4.01	0.50	<b>8.57*</b>	<b>8.57*</b>	3.64	0.00	5.58	3.34
	I/O Pairs	12.20	<b>12.20</b>	5.26	0.75	N/A	N/A	N/A	N/A	8.73	6.48
	Full	55.49		58.64		31.43		0		36.39	
LLaMa-3	NL	<b>33.54</b>	3.66	<b>16.54</b>	0.00	<b>11.43</b>	<b>5.71*</b>	0.00	<b>3.64*</b>	15.38	3.54
	<i>Code<sub>AL</sub></i>	0.61	3.66	1.76	1.51	0.00	2.86*	0.00	0.00	0.59	2.01
	<i>Code<sub>NL</sub></i>	10.98	3.05	6.02	2.01	8.57	0.00	0.00	0.00	6.39	0.98
	I/O Pairs	6.10	<b>4.27</b>	6.27	<b>2.76</b>	N/A	N/A	N/A	N/A	6.19	3.52

## 4 Experiments

### 4.1 Settings

We analyze causal effects on codeLLMs across three benchmarks: HumanEval+, mMBPP+, and CoderEval. To address the absence of the *Code<sub>AL</sub>* modality in MBPP+, we create mMBPP+ by adding a function header to the original prompts. HumanEval and mMBPP are evaluated using evalplus [23], which includes additional challenging test cases. CoderEval provides coding tasks ranging from self-contained functions to more complex ones requiring a full project setup [49]. We focus on self-contained subsets, CoderEval-SCP (Python) and CoderEval-SCJ (Java), with detailed statistics in Table 1. Using CodeSCM, we evaluate causal effects on three CodeLLMs: OpenAI GPT-4 Turbo [30], WizardCoder-15B [27], and Llama-3-8B [1]. Further implementation details are included in the Appendix C.

### 4.2 Total Effects of Modalities

**Natural Language.** The Natural Language (NL) component, often a docstring in code completion tasks and containing the core logic of the generated code, shows the highest Total Effect across all models on HumanEval+, mMBPP+, and CoderEval-SCP. As shown in the error analysis, in Table 2, removing NL (row corresponding to NL for each model and dataset) increases semantic errors, where the model generates syntactically correct but fails test cases. The Total Effect of NL is highest for HumanEval+, followed by mMBPP+ and CoderEval-SCP, likely due to the greater detail in HumanEval+ docstrings compared to the shorter ones in CoderEval-SCP. However, the  $do(X_{NL} = -1)$  intervention still maintains a non-zero accuracy. Given that generating correct code output without NL semantics is not possible, we hypothesize that the model either infers the correct NL semantics from the *Code<sub>NL</sub>*, or relies on its memory, representing memorization.

***Code<sub>NL</sub> TE.*** The latter hypothesis is confirmed with the TE computation of *Code<sub>NL</sub>*, which emerges as an important prompt component in the HumanEval+ and mMBPP+ datasets. For GPT-4T on mMBPP+, the TE of *Code<sub>NL</sub>* is 42.86%, surpassing the NL modality. Natural language chat models, like GPT-4T and LLaMa-3, consistently show higher TE for *Code<sub>NL</sub>*, with GPT-4T reaching

18.91% on HumanEval+. This suggests that natural language models may prioritize NL semantics ( $M_{NL}$ ) more than code-focused models.

**Code<sub>AL</sub> TE.** In CoderEval-SCJ,  $Code_{AL}$  has a high TE across all models, with GPT-4T and WizardCoder performance dropping nearly to zero. We observe limited code generation capabilities in the Java programming language exhibited by codeLLMs, particularly evident with zero performance from LLaMa-3. We observe that the low performance in Java is because the models hallucinate code entry points when  $Code_{AL}$  is absent. For instance, in all 55 examples, LLaMa-3 places the required function name in a hallucinated class, as illustrated in Figure 3. On Python subsets,  $Code_{AL}$ , which contains minimal syntax information like function headers and variable names (Figure 2), has the lowest TE across all models. However,  $Code_{AL}$  in all three datasets under consideration is limited to the function header and input variable names along with function syntax (Figure 2); the TE of  $Code_{AL}$  where it may contain essential generation logic is yet to be explored.

**I/O Pairs.** The TE of I/O pairs surpasses that of  $Code_{NL}$  with WizardCoder and holds equal significance with LLaMa-3 and GPT-4. This underscores the syntactic information encoded within I/O pairs, potentially aiding the model in reasoning over correct code structures. We speculate that this might be analogous to human programmers employing unit testing for iterative code writing, the TE of I/O pairs suggests a similar process within codeLLMs. Future versions of codeLLMs may even utilize intermediate I/O values for handling complex code, akin to the debugging process in software engineering.

**Memorization of Code Benchmarks.** Given that codeLLMs are trained on open-source datasets, we explore the potential for benchmark memorization. The non-zero  $pass@1$  accuracy, even without NL instructions, indicates strong memorization. Furthermore, even after standardizing function header names (Figure 2), GPT-4T still generated original function names in 11.5% of HumanEval+ and 7.2% of mMBPP+ cases (Figure 3). LLaMa-3 showed similar behavior, with 10.3% of HumanEval+ examples despite standardization or prefix transformation of function names. The notably high memorization figures for GPT-4T also raise concerns regarding its performance on the EvalPlus leaderboard[23]. While prior studies, such as [20], have addressed memorization concerns, our Causal Analysis also reveals significant dataset memorization, emphasizing the need for new code evaluation benchmarks.

### 4.3 Direct Effects of Modalities

We define direct effects (DE) by noting the drop in  $pass@1$  accuracy of the model under the semantics-preserving transformations of modalities where the latent mediators remain unchanged (Section 3). These effects also represent the spurious correlations, as any non-spurious learning process must be mediated through  $M_{NL}$  and  $M_{Code}$ . From Table 2, I/O pairs exhibit the strongest direct effect (DE) on HumanEval+ and mMBPP+ across models, except for mMBPP+ on WizardCoder. As seen in Figure 3, replacing a single assert equality in each I/O example with two inequalities makes it harder for the model to reason correctly over the code logic.

The DE of I/O pairs is then followed by the DE of  $Code_{AL}$ , where WizardCoder shows a very high DE of 9.76% on HumanEval+. For CoderEval-SCJ, GPT-4T’s accuracy increased by 18.18% under the  $do(X_{AL} = 1)$  intervention. As shown in Figure 3, a Java code snippet in the form of dead code reduces the class name hallucinated by the model. With this finding, we speculate that dead code might help control hallucinations, but we leave the detailed analysis to future work. In general, we observe that the DEs of  $NL$  and  $Code_{NL}$  are comparatively lower, implying models are more robust to natural language than code semantics, likely due to instruction tuning stages.

## 5 Related Work

**Automatic Code Generation.** Some of the earlier works on code generation with natural language and I/O pairs include [12, 16]. Recent works have either adopted the transformer architecture [14, 43] or leveraged the GPT [5] skeleton with massive pretraining for Code Generation [37, 21, 27, 29].

**Prompt-Tuning.** Various approaches to prompt-tuning [46] have been explored for various domains and modalities [28, 47], such as Chain-of-Thought reasoning [44], discrete prompt optimization

[45, 39], and few-shot learning [5]. In the context of code generation, prompt engineering has been leveraged for human-in-loop debugging [11], correctness evaluation of generated code [24], multistep planning, and generation [50]. Our work explores the effects of modalities in prompts on code generation, which can be further used to guide the prompt-tuning process for better performance.

**Causal Inference in Code/NLP.** Recent research has applied causal inference to the NLP domain [42, 15, 40] to better understand model behavior, which is now formalized as causal NLP [19, 13]. In the context of code, prior approaches have applied causal framework for various classification tasks such as vulnerability detection [34, 17] and code performance prediction [10]. To the best of our knowledge, we are the first to apply causal inference to the multi-modal code generation task.

## 6 Conclusion

We propose CodeSCM, a Structural Causal Model for analyzing multi-modal code generation using LLMs. Our analysis revealed that input-output examples and natural language code components significantly influence model generation, and our evaluations demonstrate the memorization of code generation benchmarks. Additionally, our interventions show that semantics-preserving changes can impact accuracy and can also lead to fewer hallucinations in some cases. Our work highlights the relative importance and direct effect of each modality in the prompt, which should also guide the prompt engineering process for code LLMs.

**Limitations and Future Work.** We can calculate causal effects in CodeSCM with the assumption of no confounders. We believe that in the future, our causal formulation of code generation could be extended to account for confounders using the backdoor criterion.

## References

- [1] AI@Meta. Llama 3 model card. 2024. URL [https://github.com/meta-llama/llama3/blob/main/MODEL\\_CARD.md](https://github.com/meta-llama/llama3/blob/main/MODEL_CARD.md).
- [2] Jacob Austin, Augustus Odena, Maxwell Nye, Maarten Bosma, Henryk Michalewski, David Dohan, Ellen Jiang, Carrie Cai, Michael Terry, Quoc Le, and Charles Sutton. Program synthesis with large language models, 2021.
- [3] Chen Avin, Ilya Shpitser, and Judea Pearl. Identifiability of path-specific effects. 2005.
- [4] Jinze Bai, Shuai Bai, Yunfei Chu, Zeyu Cui, Kai Dang, Xiaodong Deng, Yang Fan, Wenbin Ge, Yu Han, Fei Huang, Binyuan Hui, Luo Ji, Mei Li, Junyang Lin, Runji Lin, Dayiheng Liu, Gao Liu, Chengqiang Lu, Keming Lu, Jianxin Ma, Rui Men, Xingzhang Ren, Xuancheng Ren, Chuanqi Tan, Sinan Tan, Jianhong Tu, Peng Wang, Shijie Wang, Wei Wang, Shengguang Wu, Bofeng Xu, Jin Xu, An Yang, Hao Yang, Jian Yang, Shusheng Yang, Yang Yao, Bowen Yu, Hongyi Yuan, Zheng Yuan, Jianwei Zhang, Xingxuan Zhang, Yichang Zhang, Zhenru Zhang, Chang Zhou, Jingren Zhou, Xiaohuan Zhou, and Tianhang Zhu. Qwen technical report. *arXiv preprint arXiv:2309.16609*, 2023.
- [5] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- [6] Casey Casalnuovo, Earl T. Barr, Santanu Kumar Dash, Prem Devanbu, and Emily Morgan. A theory of dual channel constraints. In *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering: New Ideas and Emerging Results, ICSE-NIER '20*, page 25–28, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 9781450371261. doi: 10.1145/3377816.3381720. URL <https://doi.org/10.1145/3377816.3381720>.
- [7] Patrick Chao, Alexander Robey, Edgar Dobriban, Hamed Hassani, George J. Pappas, and Eric Wong. Jailbreaking black box large language models in twenty queries, 2023.

- [8] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code. 2021.
- [9] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Ponde de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Josh Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating large language models trained on code, 2021.
- [10] Jürgen Cito, Isil Dillig, Vijayaraghavan Murali, and Satish Chandra. Counterfactual explanations for models of code, 2021. URL <https://arxiv.org/abs/2111.05711>.
- [11] Paul Denny, Viraj Kumar, and Nasser Giacaman. Conversing with copilot: Exploring prompt engineering for solving cs1 problems using natural language, 2022.
- [12] Aditya Desai, Sumit Gulwani, Vineet Hingorani, Nidhi Jain, Amey Karkare, Mark Marron, Sailesh R, and Subhajit Roy. Program synthesis using natural language, 2015.
- [13] Amir Feder, Katherine A. Keith, Emaad Manzoor, Reid Pryzant, Dhanya Sridhar, Zach Wood-Doughty, Jacob Eisenstein, Justin Grimmer, Roi Reichart, Margaret E. Roberts, Brandon M. Stewart, Victor Veitch, and Diyi Yang. Causal inference in natural language processing: Estimation, prediction, interpretation and beyond, 2022.
- [14] Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xiaocheng Feng, Ming Gong, Linjun Shou, Bing Qin, Ting Liu, Daxin Jiang, and Ming Zhou. Codebert: A pre-trained model for programming and natural languages, 2020.
- [15] Matthew Finlayson, Aaron Mueller, Sebastian Gehrmann, Stuart Shieber, Tal Linzen, and Yonatan Belinkov. Causal analysis of syntactic agreement mechanisms in neural language models. *arXiv preprint arXiv:2106.06087*, 2021.
- [16] Sumit Gulwani, Oleksandr Polozov, Rishabh Singh, et al. Program synthesis. *Foundations and Trends® in Programming Languages*, 4(1-2):1–119, 2017.
- [17] Jingzhu He, Yuhang Lin, Xiaohui Gu, Chin-Chia Michael Yeh, and Zhongfang Zhuang. Perfsig: Extracting performance bug signatures via multi-modality causal analysis. In *2022 IEEE/ACM 44th International Conference on Software Engineering (ICSE)*, pages 1669–1680, 2022. doi: 10.1145/3510003.3510110.
- [18] Dan Hendrycks, Steven Basart, Saurav Kadavath, Mantas Mazeika, Akul Arora, Ethan Guo, Collin Burns, Samir Puranik, Horace He, Dawn Song, and Jacob Steinhardt. Measuring coding challenge competence with apps. *NeurIPS*, 2021.
- [19] Zhijing Jin, Amir Feder, and Kun Zhang. CausalNLP tutorial: An introduction to causality for natural language processing. In Samhaa R. El-Beltagy and Xipeng Qiu, editors, *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing: Tutorial Abstracts*, pages 17–22, Abu Dubai, UAE, December 2022. Association for Computational Linguistics. doi: 10.18653/v1/2022.emnlp-tutorials.4. URL <https://aclanthology.org/2022.emnlp-tutorials.4>.



- [20] Yuhang Lai, Chengxi Li, Yiming Wang, Tianyi Zhang, Ruiqi Zhong, Luke Zettlemoyer, Scott Wen tau Yih, Daniel Fried, Sida Wang, and Tao Yu. Ds-1000: A natural and reliable benchmark for data science code generation. *ArXiv*, abs/2211.11501, 2022.
- [21] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, Rudra Murthy, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan Dey, Zhihan Zhang, Nour Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. Starcoder: may the source be with you!, 2023.
- [22] 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. Competition-level code generation with alphacode. *arXiv preprint arXiv:2203.07814*, 2022.
- [23] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatGPT really correct? rigorous evaluation of large language models for code generation. In *Thirty-seventh Conference on Neural Information Processing Systems*, 2023. URL <https://openreview.net/forum?id=1qvx610Cu7>.
- [24] Jiawei Liu, Chunqiu Steven Xia, Yuyao Wang, and Lingming Zhang. Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation, 2023.
- [25] Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing, 2021.
- [26] Shuai Lu, Daya Guo, Shuo Ren, Junjie Huang, Alexey Svyatkovskiy, Ambrosio Blanco, Colin Clement, Dawn Drain, Daxin Jiang, Duyu Tang, Ge Li, Lidong Zhou, Linjun Shou, Long Zhou, Michele Tufano, Ming Gong, Ming Zhou, Nan Duan, Neel Sundaresan, Shao Kun Deng, Shengyu Fu, and Shujie Liu. Codexglue: A machine learning benchmark dataset for code understanding and generation, 2021.
- [27] Ziyang Luo, Can Xu, Pu Zhao, Qingfeng Sun, Xiubo Geng, Wenxiang Hu, Chongyang Tao, Jing Ma, Qingwei Lin, and Daxin Jiang. Wizardcoder: Empowering code large language models with evol-instruct, 2023.
- [28] Ankan Mullick, Mukur Gupta, and Pawan Goyal. Intent detection and entity extraction from biomedical literature. *arXiv preprint arXiv:2404.03598*, 2024.
- [29] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan Wang, Yingbo Zhou, Silvio Savarese, and Caiming Xiong. Codegen: An open large language model for code with multi-turn program synthesis, 2023.
- [30] OpenAI. Gpt-4 turbo. <https://help.openai.com/en/articles/8555510-gpt-4-turbo>, 2024.
- [31] Judea Pearl. *Causality: Models, Reasoning, and Inference*. Cambridge University Press, 2000.
- [32] Judea Pearl. Direct and indirect effects. In *Probabilistic and causal inference: the works of Judea Pearl*, pages 373–392. 2022.

- [33] Judea Pearl et al. Models, reasoning and inference. *Cambridge, UK: CambridgeUniversityPress*, 19(2):3, 2000.
- [34] Md Mahbubur Rahman, Ira Ceka, Chengzhi Mao, Saikat Chakraborty, Baishakhi Ray, and Wei Le. Towards causal deep learning for vulnerability detection. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering*, pages 1–11, 2024.
- [35] James M Robins. Semantics of causal dag models and the identification of direct and indirect effects. *Highly structured stochastic systems*, pages 70–82, 2003.
- [36] James M Robins and Sander Greenland. Identifiability and exchangeability for direct and indirect effects. *Epidemiology*, 3(2):143–155, 1992.
- [37] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Ellen Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, Ivan Evtimov, Joanna Bitton, Manish Bhatt, Cristian Canton Ferrer, Aaron Grattafiori, Wenhan Xiong, Alexandre Défossez, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. Code llama: Open foundation models for code, 2023.
- [38] Melanie Sclar, Yejin Choi, Yulia Tsvetkov, and Alane Suhr. Quantifying language models’ sensitivity to spurious features in prompt design or: How i learned to start worrying about prompt formatting, 2023.
- [39] Taylor Shin, Yasaman Razeghi, Robert L. Logan IV au2, Eric Wallace, and Sameer Singh. Autoprompt: Eliciting knowledge from language models with automatically generated prompts, 2020.
- [40] Alessandro Stolfo, Zhijing Jin, Kumar Shridhar, Bernhard Schölkopf, and Mrinmaya Sachan. A causal framework to quantify the robustness of mathematical reasoning with language models. *arXiv preprint arXiv:2210.12023*, 2022.
- [41] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten, Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. Llama 2: Open foundation and fine-tuned chat models, 2023.
- [42] Jesse Vig, Sebastian Gehrmann, Yonatan Belinkov, Sharon Qian, Daniel Nevo, Yaron Singer, and Stuart Shieber. Investigating gender bias in language models using causal mediation analysis. *Advances in neural information processing systems*, 33:12388–12401, 2020.
- [43] Yue Wang, Weishi Wang, Shafiq Joty, and Steven C. H. Hoi. Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation, 2021.
- [44] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
- [45] Yuxin Wen, Neel Jain, John Kirchenbauer, Micah Goldblum, Jonas Geiping, and Tom Goldstein. Hard prompts made easy: Gradient-based discrete optimization for prompt tuning and discovery, 2023.
- [46] Jules White, Quchen Fu, Sam Hays, Michael Sandborn, Carlos Olea, Henry Gilbert, Ashraf Elnashar, Jesse Spencer-Smith, and Douglas C. Schmidt. A prompt pattern catalog to enhance prompt engineering with chatgpt, 2023.

- [47] Qi Wu, Yuyao Zhang, and Marawan Elbatel. Self-prompting large vision models for few-shot medical image segmentation. In *MICCAI workshop on domain adaptation and representation transfer*, pages 156–167. Springer, 2023.
- [48] Yongkai Wu, Lu Zhang, Xintao Wu, and Hanghang Tong. Pc-fairness: A unified framework for measuring causality-based fairness, 2019.
- [49] Hao Yu, Bo Shen, Dezhi Ran, Jiabin Zhang, Qi Zhang, Yuchi Ma, Guangtai Liang, Ying Li, Qianxiang Wang, and Tao Xie. Codereval: A benchmark of pragmatic code generation with generative pre-trained models. In *Proceedings of the IEEE/ACM 46th International Conference on Software Engineering, ICSE '24*, New York, NY, USA, 2024. Association for Computing Machinery. ISBN 9798400702174. doi: 10.1145/3597503.3623316. URL <https://doi.org/10.1145/3597503.3623316>.
- [50] Wenqing Zheng, S P Sharan, Ajay Kumar Jaiswal, Kevin Wang, Yihan Xi, Deji Xu, and Zhangyang Wang. Outline, then details: Syntactically guided coarse-to-fine code generation, 2023.
- [51] Kaijie Zhu, Jindong Wang, Jiaheng Zhou, Zichen Wang, Hao Chen, Yidong Wang, Linyi Yang, Wei Ye, Yue Zhang, Neil Zhenqiang Gong, and Xing Xie. Promptbench: Towards evaluating the robustness of large language models on adversarial prompts, 2023.

## A Causal Effects

In this section, we present the causal effects of three modalities: Natural Language (NL), I/O Pairs, and Code with NL component ( $Code_{NL}$ ). For each modality, we provide the corresponding structural equation, followed by the total effect (TE) and direct effect (DE).

### A.1 Natural Language (NL)

**NL Structural Equation.** The  $NL$  variable is defined by the following structural equation:

$$NL \leftarrow \mathbb{1}_{\{X_{NL}=1\}}(S + DS) + \mathbb{1}_{\{X_{NL}=0\}}(S) + \mathbb{1}_{\{X_{NL}=-1\}}(NULL) \quad (4)$$

where  $S$  is the actual natural language prompt component  $\mathcal{P}_S \in \mathcal{D}$ ,  $DS$  is a Dead String that does not alter the semantics of the natural language, and  $X_{NL}$  is used to control whether to allow the original  $\mathcal{P}_S$ , concatenate a dead string, or remove the natural language modality. Similar to dead code insertion in  $Code_{AL}$ , dead string insertion is a semantics-preserving transformation such that  $M_{NL}(S) = M_{NL}(S + DS)$ .

**Total Effect of NL.**

$$\begin{aligned} TE(NL) &= TE(do(X_{NL} = 0), do(X_{NL} = -1)) \\ &= Acc(\mathcal{D}) - Acc(\mathcal{D}; \mathcal{P}_{NL} = NULL) \end{aligned}$$

**Direct Effect of NL.**

$$\begin{aligned} DE(NL) &= \mathbb{E} \left[ Y \left( X_{NL} = 1, NL(X_{NL} = 1), M_{NL}(NL(X_{NL} = 0)) \right) \right] \\ &\quad - \mathbb{E} [Y(NL(X_{NL} = 0))] \\ &= Acc(\mathcal{D}) - Acc(\mathcal{D}; \mathcal{P}_{NL} = S + DS) \end{aligned}$$

Here,  $DS$  represents the Dead String. We use the prefix ‘DOCSTRING: ’ concatenated to each natural language instruction to preserve semantics. Other transformations such as back-translation are possible but introduce correlations between variables, so we prefer simpler prefix or suffix transformations that keep  $S$  and  $DS$  independent.

### A.2 I/O Pairs

**I/O Structural Equation.** The I/O modality is defined by the following structural equation:

$$I/O \leftarrow \mathbb{1}_{\{X_{IO}=0\}}(I^r = I^r) + \mathbb{1}_{\{X_{IO}=1\}}((I^l \leq I^r) + (I^r \geq I^r)) + \mathbb{1}_{\{X_{IO}=-1\}}(NULL) \quad (5)$$

where  $I^l$  and  $I^r$  represent the left-hand side (LHS) and right-hand side (RHS) of the assertion equality statement in the original prompt, respectively. For semantics-preserving transformations, we replace each assertion equality with two inequalities,  $\leq$  and  $\geq$ .

**Total Effect of I/O.**

$$\begin{aligned} TE(IO) &= TE(do(X_{IO} = 0), do(X_{IO} = -1)) \\ &= Acc(\mathcal{D}) - Acc(\mathcal{D}; \mathcal{P}_{IO} = NULL) \end{aligned}$$

**Direct Effect of I/O.**

$$\begin{aligned} DE(IO) &= \mathbb{E} \left[ Y \left( X_{IO} = 1, M_{Code}(X_{IO} = 1), M_{IO}(M_{Code}(X_{IO} = 0)) \right) \right] \\ &\quad - \mathbb{E} [Y(M_{IO}(X_{IO} = 0))] \\ &= Acc(\mathcal{D}) - Acc(\mathcal{D}; \mathcal{P}_{IO} = (I^l \leq I^r) + (I^r \geq I^r)) \end{aligned}$$

### A.3 Code with NL Component ( $Code_{NL}$ )

**Code<sub>NL</sub> Structural Equation.** The  $Code_{NL}$  modality is defined by the following structural equation:

$$Code_{NL} \leftarrow \mathbb{1}_{\{X_{NL}=1\}}(C_{NL} + DN) + \mathbb{1}_{\{X_{NL}=0\}}(C_{NL}) + \mathbb{1}_{\{X_{NL}=-1\}}(NULL) \quad (6)$$

where  $C_{NL}$  is the code prompt component  $\mathcal{P}_{Code_{NL}} \in \mathcal{D}$ , and  $DN$  is a Dead Name added to the function header. This transformation preserves semantics for both the natural language and code domains. For instance,  $M_{NL}(C_{NL}) = M_{NL}(C_{NL} + DN)$ .

**Total Effect of  $Code_{NL}$ .**

$$\begin{aligned} TE(Code_{NL}) &= TE(do(X_{CN} = 0), do(X_{CN} = -1)) \\ &= Acc(\mathcal{D}) - Acc(\mathcal{D}; \mathcal{P}_{Code_{NL}} = NULL) \end{aligned}$$

**Direct Effect of  $Code_{NL}$ .**

$$\begin{aligned} DE(Code_{NL}) &= \mathbb{E} \left[ Y(X_{CN} = 1, Code_{NL}(X_{CN} = 1), \right. \\ &\quad \left. M_{NL}(Code_{NL}(X_{CN} = 0)), M_{Code}(Code_{NL}(X_{CN} = 0)) \right) \\ &\quad - \mathbb{E} \left[ Y(Code_{NL}(X_{CN} = 0)) \right] \\ &= Acc(\mathcal{D}) - Acc(\mathcal{D}; \mathcal{P}_{Code_{NL}} = C_{NL} + DN) \end{aligned}$$

Here,  $DN$  represents Dead Name, and we use the prefix ‘func\_’ in Python and ‘Method’ in Java to maintain semantic preservation. Other transformations, like capitalization, are possible but avoided to keep  $C_{NL}$  and  $DN$  independent.

## B Multi-Modal Prompt

The multi-modal prompt  $\mathcal{P}$  can be expressed as an equation comprising one or more prompt components  $P^j$  of modality  $M_i$ , where different prompt components are concatenated using one of the defined separators:

$$\mathcal{P} = P_{M_1}^1 \left[ \text{sep } P_{M_i}^j \right], \text{ sep} = ' ' | \backslash n | \backslash t | : | , | < \text{sep\_token} > | ' | ; \quad (7)$$

In this equation, different prompt components are concatenated using one of the defined separators.

Furthermore, the figure below (Fig. 2) shows an example of a prompt broken down into its constituent modalities and the semantics-preserving transformations applied to each.

## C Implementation details

All datasets used are evaluation-only subsets, with no training involved in our experiments. For inference on all LLMs, we use a temperature of 0.01, a top\_p value of 0.95, and a batch size of 8. The open-source model experiments were conducted on a single A100 GPU with 40 GB VRAM and GDDR5 memory. We selected ‘\n’ as our modal-separator, though different choices of modal-separators (equation 7) may result in varying performance [38]. During experiments with self-contained CoderEval functions in Python and Java, we ensured transformations were equivalent across both languages. For example, dead code added for the  $Code_{AL}$  modality was semantically equivalent in Python and Java, depending on the prompts.

Furthermore, we exclude APPS [18] and CodeContest [22], as they lack multi-modal prompts, making them unnecessary for multi-modal causal analysis. Similarly, while the CONCODE segment of the CodexCGLUE [26] benchmark includes multi-modal prompts, it measures code quality via natural language similarity metrics like BLEU, which is unsuitable for code generation tasks. Lastly,

DS-1000 [20] was excluded due to the need for manual screening of all examples to separate modal components for CodeSCM.

In terms of models, WizardCoder, a ‘code-aware’ model, is instruction-tuned using Evol-Instruct with 34B Python tokens, building on StarCoder [21]. LLaMa-3 [1] is pre-trained on 15T tokens, including four times more code data than its predecessor LLaMa-2 [41].

For our results, we use changes in mean  $pass@1$  accuracy ( $Pr(Y = 1)$ ) to measure the direct and total effects after interventions on CodeSCM [9]. For DE, simple string transformations (dead code, dead string, and dead name) are applied to  $Code_{AL}$ ,  $NL$ , and  $Code_{NL}$  modalities, and for I/O pairs, assert equalities are replaced with equivalent inequalities.

## D DE Additional Transformation

We demonstrate one specific transformation for each modality in the paper and compute the respective causal effects. CodeSCM can be directly extended to other transformations as well for DE computation. For example, in Table 3, along with original transformations from Table 2 (DE-1), we illustrate DE computation with an additional set of transformations (DE-2) for the mMBPP+ dataset using WizardCoder codeLLM. The following transformations are used for DE-2 - (dead string prefix, unused variable, dead name prefix, and negating the not assert statement):

- $DS = \text{Code Logic} \setminus \text{n}$  (in Equation 4)
- $C_{DC} = \text{tvar} = 42$  (in Equation 3)
- $DN = \text{header}_-$  (in Equation 6)
- $\text{assert } I^l == I^r$  is changed to  $\text{assert not } I^l = I^r$  (in Equation 5)

Table 3: Direct effects of WizardCoder on mMBPP+ dataset under an additional transformation. DE-1 values are the same as Table 2

Modality	DE-1	DE-2
Full	52.63	
NL	0.50	1.23
$Code_{AL}$	2.50	3.03
$Code_{NL}$	0.50	1.73
I/O Pairs	0.75	3.23

## E Error analysis

In this section, we categorize and analyze the errors encountered across different models and datasets, such as syntax errors, semantic errors, and runtime errors. Table 4 presents a breakdown of errors encountered by running the HumanEval+ and mMBPP+ datasets on GPT-4T, WizardCoder-15B, and Llama-3-8B. Negative numbers in this table mean that when an intervention is made upon a prompt, there is a decrease in error count from the full prompt. Positive numbers mean that upon intervention there is an increase in error count. For example, in the case of GPT-4T with the  $Code_{NL}$  modality on the mMBPP+ dataset, the intervention leads to an increase in error count for runtime errors by 32.21%. Similarly, when the NL modality is removed for HumanEval+ on WizardCoder-15B, there is a decrease in runtime errors by 11.65%. We classify errors into the following four categories:

**Syntax Errors.** These errors occur when the code does not conform to the syntactical rules of the programming language. They are typically detected during the parsing stage. An example of a syntax error might be a missing colon, unmatched parentheses, or incorrect indentation.

**Semantic Errors.** Semantic errors arise when the code is syntactically correct but fails to produce the intended output due to logical mistakes. This can include errors in the logic of the code, incorrect use of variables, or wrong implementation of algorithms. We broadly encounter two types of Semantic Errors, (i) **test case errors**: when the test cases in the respective dataset fail; (ii) **assertion errors**: when an input-output example assertion in the prompt fails.

Figure 3: Left figure shows a CoderEval-SCJ prompt where dead code insertion corrects the original prompt’s error of creating a hallucinated Java class (red box). The top right figure illustrates an mMBPP+ prompt where I/O pair transformations lead to a semantic error in lines 15-16. The bottom right figure shows GPT-4T’s memorization of a HumanEval+ prompt.

Table 4: Percentage of errors out of total passed cases for GPT-4T, WizardCoder-15B, and Llama-3-8B on HumanEval+ and mMBPP+. Negative percentages indicate a decrease in error count, while positive values indicate an increase in error count upon intervention. Syn represents Syntax errors, Sem represents Semantic Errors, and RunT represents runtime errors.

MODEL	MODALITY	HUMANEVAL+				MMBPP+			
		SYN	SEM	RUNT	OTHER	SYN	SEM	RUNT	OTHER
GPT-4T	NL	-8.22	26.13	-18.67	0.76	-5.58	14.32	-9.61	0.87
	Code <sub>AL</sub>	-0.75	1.49	-1.80	1.10	-1.42	2.31	-0.89	0.00
	Code <sub>NL</sub>	-3.12	-7.04	-8.47	18.63	-11.27	-30.64	32.21	9.70
	I/O PAIRS	-1.62	0.61	1.01	0.00	-3.85	-4.75	7.65	0.94
WIZARDCODER	NL	-3.37	15.01	-11.65	0.00	-1.29	7.72	-6.01	-0.42
	Code <sub>AL</sub>	-0.30	-0.65	-0.95	-0.00	-0.92	6.93	-5.58	-0.42
	Code <sub>NL</sub>	-0.59	-4.58	-0.21	4.96	-1.04	-2.08	2.74	0.39
	I/O PAIRS	-1.01	2.67	-1.65	0.0	-1.05	4.06	-3.38	0.37
LLAMA-3	NL	-4.11	20.96	-16.60	-0.25	-5.14	14.13	-7.98	-1.01
	Code <sub>AL</sub>	0.42	0.90	-1.30	-0.03	-1.32	4.38	-2.13	0.93
	Code <sub>NL</sub>	-1.57	-8.67	-1.21	11.44	-0.90	-10.04	7.68	3.26
	I/O PAIRS	4.43	1.75	-6.13	-0.05	-0.41	-9.85	10.80	-0.54

**Runtime Errors.** These errors occur during the execution of the code. They result from operations like division by zero, accessing out-of-bound indices, or other exceptional conditions that the code does not handle.

**Other Errors.** This category includes various errors that do not fit into the above classifications. It covers, (i) **resource errors:** these happen when the code exceeds the available resources, such as memory errors when the program tries to allocate more memory than what is available; (ii) **dependency errors:** these arise when the code fails to import necessary modules or packages. It could be due to missing dependencies or incorrect module names; (iii) **environment errors:** these are caused by issues in the execution environment, such as problems with file access, permissions, or environment-specific configurations; (iv) **timeout errors:** these occur when the execution of the code takes longer than the allowed time limit, indicating inefficiencies or infinite loops in the implementation.