

SCLA: Automated Smart Contract Summarization via LLMs and Control Flow Prompt

Anonymous ACL submission

Abstract

Smart contract code summarization is crucial for efficient maintenance and vulnerability mitigation. While many studies use Large Language Models (LLMs) for summarization, their performance still falls short compared to fine-tuned models like CodeT5+ and CodeBERT. Some approaches combine LLMs with data flow analysis but fail to fully capture the hierarchy and control structures of the code, leading to information loss and degraded summarization quality. We propose SCLA, an LLMs-based method that enhances summarization by integrating a Control Flow Graph (CFG) and semantic facts from the code’s control flow into a semantically enriched prompt. SCLA uses a control flow extraction algorithm to derive control flows from semantic nodes in the Abstract Syntax Tree (AST) and constructs the corresponding CFG. Code semantic facts refer to both explicit and implicit information within the AST that is relevant to smart contracts. This method enables LLMs better to capture the structural and contextual dependencies of the code. We validate the effectiveness of SCLA through comprehensive experiments on a dataset of 40,000 real-world smart contracts. The experiment shows that SCLA significantly improves summarization quality, outperforming the SOTA baselines with improvements of 26.7%, 23.2%, 16.7%, and 14.7% in BLEU-4, METEOR, ROUGE-L, and BLEURT scores, respectively.

1 Introduction

Smart contracts (Liao et al., 2023) are self-executing programs on Ethereum, and the blockchain’s immutability complicates vulnerability maintenance (Zhang et al., 2022). Solidity, designed specifically for smart contract development, compiles code into bytecode and ABI for execution on the Ethereum Virtual Machine (EVM). Unlike general-purpose languages like Java and Python, Solidity emphasizes security with strict type safety

and single-threaded execution. Analyzing Solidity code requires examining syntax, semantics, and state management. Even minor vulnerabilities can result in financial losses (Kushwaha et al., 2022), making smart contract code summarization essential for improving efficiency and reducing security risks. Smart contract summarization has received less attention than Java and Python, with traditional methods relying on deep learning and fine-tuning. Yang et al. (Yang et al., 2021) proposed MMTrans, integrating deep learning with structure-based traversal (SBT) and Abstract Syntax Tree (AST) graphs for code summarization. Lei et al. (Lei et al., 2024) introduced FMCF, a Transformer-based method that fuses multi-scale features to preserve both semantic and syntactic information. Zhao et al. (Zhao et al., 2024) proposed SCCLLM, combining context learning with information retrieval to improve summarization.

However, these approaches pose security risks. Fine-tuning is limited by the quality and size of training data, which may lack emerging vulnerabilities. Additionally, fine-tuned models may suffer from knowledge forgetting (De Lange et al., 2022), reducing adaptability. In contrast, LLMs-based approaches leverage extensive pre-trained datasets to identify vulnerabilities and enhance security. However, prior methods mainly use isolated code snippets, failing to capture the full semantic context. Ahmed et al. (Ahmed et al., 2024) demonstrated that LLMs struggle with implicit semantics, leading to the loss of critical information. These methods also neglect control structures, hindering accurate summarization and security.

To address the limitations of existing methods and generate more secure smart contract code summarization, we propose SCLA (Smart Contract summarization with LLMs and Semantic Augmentation). SCLA integrates LLMs with control flow analysis to enhance scalability and summarization quality. Incorporating control flow-based semantic

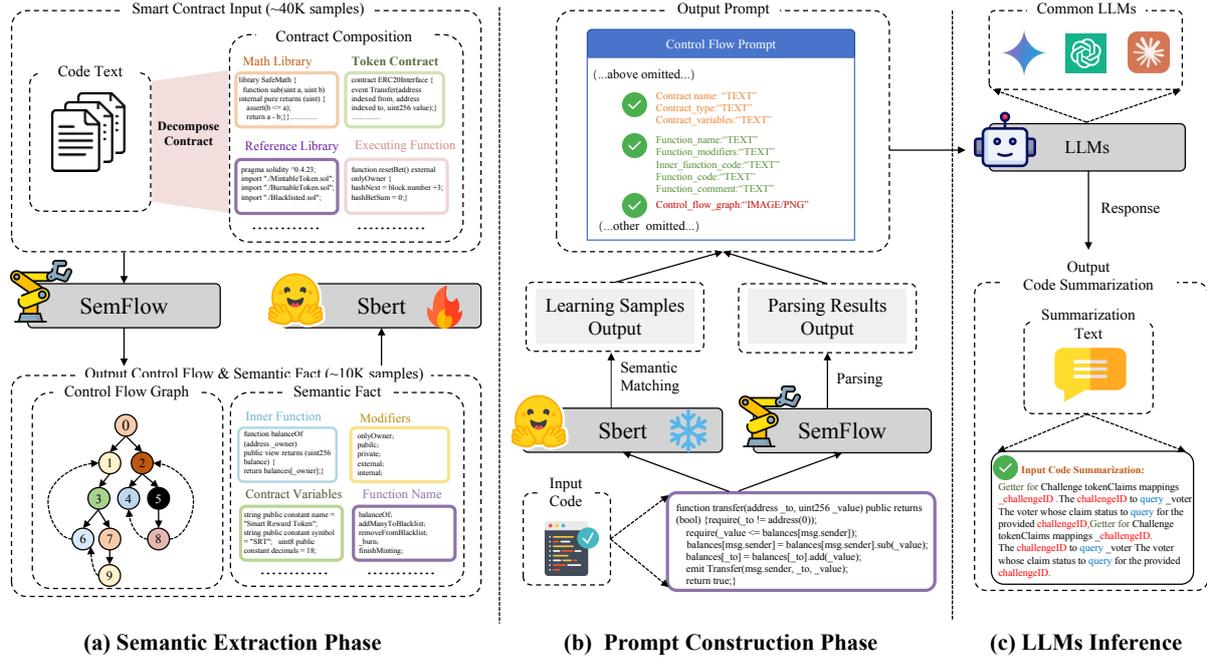


Figure 1: Overview of our proposed framework, SCLA, powered by Google’s Gemini-1.5-Pro, performs automated generation of smart contract code summarization. SCLA extracts control flow semantic facts from smart contract code and uses Gemini-1.5-Pro to generate code summarization from control flow semantic facts.

information and control flow graphs improves security and semantic accuracy in LLMs-generated summarization. The approach extracts this information via semantic analysis and generates control flow graphs. Using few-shot prompting and a fine-tuned Sentence-Transformer(Reimers and Gurevych, 2020), SCLA identifies semantically similar examples for task-specific prompts. **SemFlow**, a core component, extracts control flow graphs and semantic details, such as function call graphs and variables. To prevent overwhelming LLMs, non-control flow information is presented separately, while control flow graphs are provided in a tagged PNG format. Our experiments on 14,789 method-comment pairs from a GitHub repository of 40,000 smart contracts show that control flow graphs improve LLMs’ performance.

Our main contributions can be summarized as follows:

- We propose SCLA, the first framework that integrates LLMs with smart contract code summarization using control flow prompts. It extracts control flow graphs and associated semantic information from the AST, enhancing the LLM’s understanding of code structure.
- We conduct extensive experiments on a dataset with 14,795 method-comment pairs,

using BLEU-4, METEOR, ROUGE-L, and BLEURT as evaluation metrics. We perform a comparative analysis with state-of-the-art approaches, achieving a 37.53 BLEU-4 score, 52.54 METEOR score, 56.97 ROUGE-L score, and 63.4 BLEURT score.

- We thoroughly evaluated the generalizability of SCLA through extensive experiments on Java and Python datasets, offering valuable insights for future research on control flow-based prompts in other code domains.
- We have uploaded the related code and experimental data to Figshare, with plans to open-source them after the paper’s publication.

2 Related Work

Smart Contract Summarization

Deep learning models have made significant advances in smart contract code summarization. Yang et al. (Yang et al., 2021) proposed MMTrans, which extracts SBT sequences and AST-based graphs to capture global and local semantics using dual encoders and a joint decoder. Transformer models like CodeT5 (Wang et al., 2021) and CodeBERT (Feng et al., 2020) also enhance summarization quality but require extensive fine-tuning and

large datasets. LLMs, such as GPT-4o and Gemini-1.5-Pro, excel in few-shot or zero-shot summarization tasks, bypassing fine-tuning. Previous studies (Ahmed and Devanbu, 2023; Ahmed et al., 2024) highlight the benefits of few-shot learning. However, LLMs often produce suboptimal summarization, lacking conciseness and functional generalization. Ahmed et al. (Ahmed et al., 2024) proposed ASAP, incorporating data flow and GitHub context. Still, it fails to capture function call relationships and control flow, suggesting the need for improved semantic facts and control flow integration for better summarization.

3 METHODOLOGY

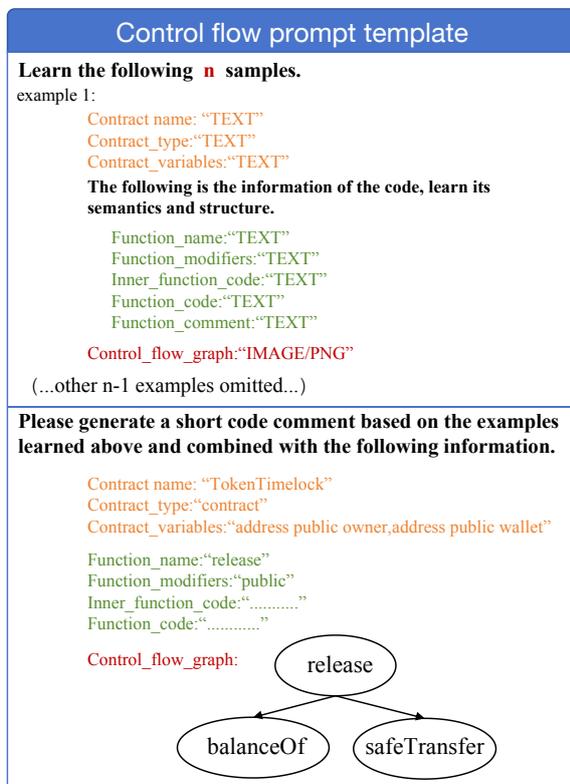


Figure 2: An Example of Control Flow Prompt.

3.1 Control Flow Prompt

In this section, we discuss the control flow prompt and the corresponding semantic facts utilized by SCLA, as illustrated in Figure 2. These semantic facts are carefully integrated into the prompts to enhance the LLMs’ ability to generate more accurate, relevant, and comprehensive summarizations.

Control Flow Graph & Inner Function. We define the set of inner functions as those invoked within the target function, with each element referred to as an inward function. The function call

graph captures the precise sequence of function calls, representing the control flow of the target code. This graph is used as control flow input for the LLMs, along with the set of inner functions, to provide valuable additional context about invoked functions. This approach mitigates misinterpretation based solely on function names, significantly enhancing semantic inference. Moreover, the function call graph helps the LLMs accurately determine the sequence and depth of function calls, thereby aiding in the understanding of complex functions and their interdependencies.

Identifiers. Previous studies highlight that identifiers play a critical role in helping language models retrieve valuable information for code summarization (Ahmed and Devanbu, 2022). Identifiers, including modifiers, local variables, and function names, offer essential context about the code’s operations. By understanding an identifier’s role, the language model can better interpret the code. In our approach, we use a tree-sitter to traverse the AST of the function, collect identifiers and their roles, and incorporate them into the prompt.

Contract Name & Global Member Variables. Incorporating domain-specific information into prompts greatly enhances LLMs’ overall performance and effectiveness, particularly in specialized tasks such as smart contract analysis. For instance, smart contract names (Kong et al., 2024) often reflect their functional roles or token names (Chen et al., 2021), providing valuable contextual information for the LLMs. Additionally, global member variables, such as contract addresses and account balances, assist LLMs in more effectively understanding contract functions and their interrelations. This significantly reduces the need for LLMs to infer complex operations from variable names, leading to more precise descriptions and significantly improved summarization accuracy.

3.2 Semantic-based Retrieval

In this paper, we use the **Sentence-Transformer** (Sbert)(Reimers and Gurevych, 2020) model to semantically match the code samples in the repository that are most similar to the target code snippet, which is then used as few-shot learning examples in the prompt. First, we divide the samples in the code repository into three subsets: training set, test set, and validation set (as shown in Table 1), and fine-tune the SBERT model using the training set. We begin by vectorizing the given sentences S_1 and

S_2 , as described by the following formula:

$$\begin{aligned} \mathbf{v}_1 &= \text{Pooling}(BERT(S_1)) \\ \mathbf{v}_2 &= \text{Pooling}(BERT(S_2)) \end{aligned} \quad (1)$$

Sbert is trained using contrastive learning or triplet loss, optimizing sentence embeddings such that similar sentences are closer in vector space and dissimilar sentences are farther apart. Given a positive pair (S_1, S_2) and a negative pair (S_1, S_3) , the model optimizes the following loss function:

$$\mathcal{L} = \max(0, \text{cosine_similarity}(\mathbf{v}_1, \mathbf{v}_2) - \text{cosine_similarity}(\mathbf{v}_1, \mathbf{v}_3) + \Delta) \quad (2)$$

where Δ is a margin hyperparameter that controls the minimum desired similarity difference, and \mathbf{v}_1 and \mathbf{v}_2 are the vector representations of sentences S_1 and S_2 .

Finally, we compute the cosine similarity between the target code vector and the repository code vectors to identify the most semantically similar samples. The formula is as follows:

$$\text{cosine_similarity}(\mathbf{v}_1, \mathbf{v}_2) = \frac{\mathbf{v}_1 \cdot \mathbf{v}_2}{\|\mathbf{v}_1\| \|\mathbf{v}_2\|} \quad (3)$$

For each target sample, we rank the repository samples by cosine similarity in descending order. The top k matches, as specified by the parameter `number_top_matches`, are selected and stored in a result dictionary, which contains the matched code snippets and their similarity scores. If a file path is provided, the results are serialized and saved in JSON format for further analysis or review.

3.3 SCLA Framework

Figure 1 illustrates the overall framework of SCLA. We outline the three stages of the SCLA process for generating smart contract code summarization.

Semantic Extraction: We use local extraction methods to segment `.sol` files, avoiding parsing errors. Regular expressions extract code and comments, which are then passed to **SemFlow** for semantic extraction. The function call graph and semantic facts are stored in a repository, indexed by contract file path, and named UUIDs.

Prompt Construction: SCLA uses few-shot learning to enhance LLMs' code summarization performance. Sentence-Transformer (Reimers and Gurevych, 2020) retrieves the top k semantically similar code samples. The extracted semantic information, including function call graphs, function

arguments, function modifiers, and contract metadata, is integrated into the prompt.

LLMs Inference: The semantically enhanced prompt, including the function call graph, is input into the LLMs interface to improve understanding of the function call sequence, resulting in higher-quality code summarization.

Algorithm 1 Source Data to Function Call Tree

```

1: Input: Source code  $f$  to be parsed by Solparser; initialized empty dictionary  $T$ 
2: Output: Function call tree  $T$ 
3:  $AST \leftarrow \text{Solparser.parser}(f)$ 
4:  $T \leftarrow \{\}$ 
5: for each  $c$  in  $AST$  do
6:   for each  $g$  in  $c$  do
7:     for each  $x$  in  $g$ .calls do
8:        $n \leftarrow x$ .name
9:       if  $n \notin T[c][g]$  then
10:         $T[c][g][n] \leftarrow \{c : c, \text{count} : 1\}$ 
11:       else
12:         $T[c][g][n].\text{count} \leftarrow T[c][g][n].\text{count} + 1$ 
13:       end if
14:     end for
15:   end for
16: end for
17: for each  $c, g$  in  $T$  do
18:   CreateCallTree( $c, g, T[c][g], T$ )
19: end for
20: return  $T$ 
21: function CREATECALLTREE( $p, k, n, T$ )
22:   for each  $m$  in  $n$ .keys do
23:      $o \leftarrow n[m]$ 
24:     if  $m \notin T[p][k]$ .keys then
25:        $T[p][k][m] \leftarrow T[o.c][m]$ 
26:       CreateCallTree( $k, m, T[o.c][m], T$ )
27:     end if
28:   end for
29: end function

```

3.4 Control Flow Extraction

We use **SemFlow**, a component integrated with a control flow extraction algorithm, to extract function call graphs from the AST as control flow input in the prompt." The algorithm in 1 demonstrates the entire extraction process. It first uses an AST parsing tool to parse the input code into an AST. The AST is traversed in a depth-first manner to remove irrelevant nodes, such as imports. Function nodes with calls are marked in the "FunctionCall" field, allowing the construction of a reference tree (lines 5-20 of Algorithm 1). The depth of the reference tree ranges from 2 to 3 layers, depending on the presence of function calls. When a third-level call points to a second or third-level node, the reference tree is transformed into a complete call tree by grafting branch nodes (lines 21-29 of Algorithm 1). The call tree is then visualized using Graphviz and saved to the code sample repository.

4 EXPERIMENT

In the empirical study, we conducted comparison, ablation, and generalization experiments. First, we

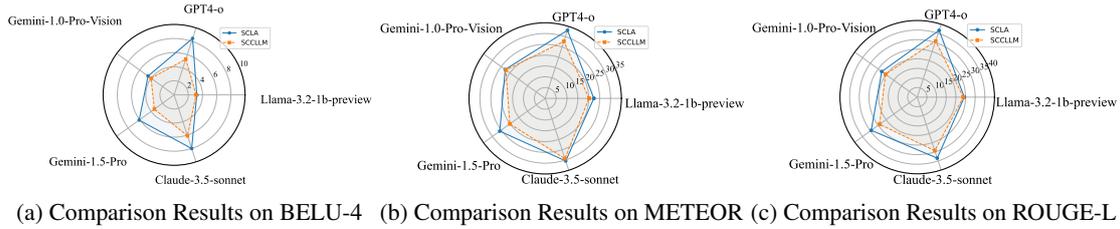


Figure 3: The Comparison of BLEU, METEOR, and ROUGE-L Scores on Our Test Set Under Five Different LLMs, Using the SCCLLM and the Proposed SCLA for Zero-Shot Summarization Tasks.

Type	Train	Validation	Test
Number	11032	2758	1000
Avg. tokens in codes	42.44	42.08	41.95
Avg. tokens in comments	26.34	26.16	26.66

Table 1: Statistics of Experimental Dataset.

used **SemFlow** to process the raw data and generate semantic facts, data flow graphs, and the semantic sample library. The code snippets were then input into SCLA for summarization and evaluation. In the comparison experiment, we varied the number of few-shot learning samples and compared the evaluation scores with baseline methods. Ablation experiments assessed the contribution of different semantic components, while generalization experiments extended SCLA to Java and Python code summarization tasks. The results and expert evaluations validate the effectiveness of SCLA in generating smart contract code summarizations.

4.1 Experiment Settings

All our experiments are performed on a computer equipped with an Nvidia GeForce RTX 4070Ti GPU (12GB graphic memory), Gen Intel (R) Core (TM) i9-13900K, running Ubuntu 22.04 LTS.

4.2 DataSet

The raw data for our study, provided by Liu et al. (Liu et al., 2021), consists of 40,000 Solidity smart contracts from Etherscan.io¹, authored by professional developers and deployed on Ethereum. Building on Yang et al.’s work (Yang et al., 2021), we developed an extraction method using AST location data to segment the code into distinct domains and extract functions with their documentation via regular expressions. After filtering low-quality summarizations, we removed poor entries, resulting in 14,790 method-comment pairs. The dataset is divided into 11,032 training, 2,758 validation, and 1,000 test sets, with average token counts

¹<https://etherscan.io/>

reported. Notably, unlike baselines, our approach, SCLA, does not require a validation set.

4.3 Baseline

We compare our proposed SCLA with six state-of-the-art methods, including general code summarization models such as **CodeT5** (Wang et al., 2021), **CodeT5+** (Wang et al., 2023), and **CodeBERT** (?), deep learning-based smart contract code summarization methods **MMTran** (Yang et al., 2021) and **FMCF** (Lei et al., 2024), and smart contract-specific code summarization methods based on the latest LLMs, such as **SCCLLM** (Zhao et al., 2024).

4.4 Performance Metrics

To evaluate SCLA performance against baselines, we adopted various automatic performance metrics, including **BLEU-4** (Papineni et al., 2002), **METEOR** (Banerjee and Lavie, 2005), **ROUGE-L** (Lin, 2004), and **BLEURT** (Sellam et al., 2020). These metrics effectively assess the similarity between the automatically generated smart contract summarization and the real human-generated summarization. BLEURT, in particular, calculates similarity based on sentence semantics by using a pre-trained BERT model, providing a more accurate reflection of semantic meaning.

4.5 Main Results

We conducted a comprehensive evaluation of the LLMs-driven SCLA framework in two different experimental setups. The framework demonstrated significant performance improvements in smart contract code summarization tasks in both zero-shot and few-shot settings. These findings provide valuable insights and contributions to the research community. The specific results are as follows:

Zero-shot Results. To evaluate the impact of function call graphs and internal functions on LLMs-generated code summarization, we conducted experiments using GPT-4o, Gemini-1.5-Pro, and Claude-3.5-Sonnet under zero-shot conditions.

Model	# of sample	BLEU-4			METEOR			ROUGE-L			p-value
		Zero-Shot	+CFG +IF	Gain(%)	Zero-Shot	+CFG +IF	Gain(%)	Zero-Shot	+CFG +IF	Gain(%)	
Llama-3.2-1b-preview	11032	3.03	5.43	+79.21%	19.58	23.97	+22.42	18.88	23.49	+24.42	<0.01
GPT-4o	11032	5.34	7.45	+39.51%	22.32	26.62	+19.27	25.32	32.62	+28.83	<0.01
Gemini-1.0-Pro-Vision	11032	3.01	5.32	+76.74%	16.89	20.73	+22.73	18.46	20.31	+10.02	<0.01
Gemini-1.5-Pro	11032	3.21	5.87	+82.87%	19.89	25.61	+28.76	23.95	27.42	+14.49	<0.01
Claude-3.5-sonnet	11032	3.31	5.32	+60.73%	23.42	28.62	+22.20	25.82	30.12	+16.65	<0.01

Table 2: Performance of different LLMs on smart contract code summarization, measured using BLEU-4, METEOR, ROUGE-L. p-values are calculated applying a one-sided pairwise Wilcoxon signed-rank test and B-H corrected.

Approach	#of training sample	#of test sample	BLEU-4	METEOR	ROUGE-L	BLEURT	p-values
CodeT5+	11032	1000	28.95	45.62	49.77	57.79	/
CodeT5	11032	1000	27.24	43.31	49.03	52.61	/
CodeBERT	11032	1000	26.31	39.57	44.52	52.74	/
MMTran	11032	1000	22.12	38.92	40.12	54.73	/
FMCF	11032	1000	29.98	36.67	51.21	51.73	/
SCLLM (One-Shot)	/	1000	19.45	20.12	19.12	36.56	<0.01
SCLLM (Three-Shot)	/	1000	29.73	35.33	49.44	50.91	<0.01
SCLLM (Five-Shot)	/	1000	31.73	48.12	60.44	58.74	<0.01
SCLA (Zero-Shot)	/	1000	6.09	25.80	29.45	46.63	<0.01
SCLA (One-Shot)	/	1000	25.46	42.78	47.55	57.07	<0.01
SCLA (Three-Shot)	/	1000	35.15	51.80	55.89	63.11	<0.01
SCLA (Five-Shot)	/	1000	37.53	52.54	56.97	63.44	<0.01

Table 3: The impact of different few-shot learning quantities on SCLA performance with Gemini-1.5-Pro. p-values are calculated applying a one-sided pairwise Wilcoxon signed-rank test and B-H corrected.

Type	Zero-Shot	One-Shot	Three-Shot	Five-Shot
Avg. tokens in prompt	561.4	1154.8	2242.5	3330.3

Table 4: Number of Tokens Consumed with Different Numbers of Learning Sample for SCLA.

The experiment had two phases: first, the target code was embedded into the prompt and evaluated with standard metrics; second, the prompt was enhanced with internal functions and control flow graphs, followed by re-evaluation. Table 2 shows that incorporating internal functions and call graphs improved summarization. GPT-4o improved by 39.51%, 19.27%, and 28.83%; Gemini-1.5-Pro by 82.87%, 28.76%, and 14.49%; and Claude-3.5-Sonnet by 60.73%, 22.20%, and 16.65%. However, Gemini-1.5-Pro underperformed compared to GPT-4o and Claude-3.5-Sonnet. **These results validate our hypothesis that control flow graphs enhance smart contract summarization.** To further validate the control flow prompt’s effectiveness, we compared SCLA with SCLLM using five multimodal models on the test set. Results in Figure 3 show SCLA outperforming SCLLM in BLEU, METEOR, and ROUGE-L scores. **This demonstrates that SCLA with the control flow prompt outperforms SCLLM, confirming the effectiveness of control flow prompts.**

Few-shot Results. To evaluate the performance of SCLA against SOTA baseline models, we conducted a validation experiment. Since SCLA em-

plloys few-shot learning, we tested its performance under Zero-Shot, One-Shot, Three-Shot, and Five-Shot conditions to investigate the number of learning samples required for optimal performance. The results (see Table 3) indicate that SCLA initially lags behind the baseline models in Zero-Shot and One-Shot settings. However, starting from Three-Shot, SCLA outperforms the baseline models across all four evaluation metrics: BLEU-4, METEOR, ROUGE-L, and BLEURT. Compared to FMCF, SCLA improved by 17.24%, 41.26%, 9.14%, and 22.00%, and compared to CodeT5+, the improvements were 21.42%, 13.55%, 12.30%, and 9.21%. **Compared to all baseline models, SCLA showed average improvements of 26.7%, 23.2%, 16.7%, and 14.7% in these metrics.** Performance continued to improve under Five-Shot, although the gains were modest. We also analyzed token consumption to determine the optimal number of few-shot samples (see Table 4). The token consumption for Five-Shot was 48.51% higher than for Three-Shot, but the average improvement in generated code summarization metrics was only 2.20%. Therefore, **Three-Shot provides the best balance between performance and efficiency.**

4.6 Ablation Study

We conducted We conducted ablation experiments to assess the contributions of different semantic facts in SCLA. The goal was to quantify the impact

Language	Model	#of Test Sample	BLEU-4			BLEURT			p-values
			SCCLLM	SCLA	Gain (%)	SCCLLM	SCLA	Gain (%)	
Java	GPT-4o	1000	28.59	38.43	+34.42	50.34	68.89	+36.85	<0.01
	Gemini-1.5-Pro	1000	23.22	31.43	+35.36	56.33	63.67	+13.03	<0.01
	Claude-3.5-sonnet	1000	31.05	39.13	+26.02	58.89	70.90	+20.40	<0.01
Python	GPT-4o	1000	22.78	29.56	+29.76	55.90	64.23	+14.90	<0.01
	Gemini-1.5-Pro	1000	20.15	26.06	+29.33	51.78	61.03	+17.86	<0.01
	Claude-3.5-sonnet	1000	25.45	33.77	+32.69	58.21	73.56	+26.37	<0.01
Overall	/	/	25.21	33.06	+31.14	55.24	67.05	+21.38	<0.01

Table 5: The performance of SCLA and SCCLLM on the Java and Python tasks, driven by three different LLMs, was evaluated using BLEU-4 and BLEURT as metrics. To assess the statistical significance of the results, p-values were calculated using a one-sided pairwise Wilcoxon signed-rank test, with Benjamini-Hochberg (B-H) correction applied for multiple comparisons.

Approach	#of Training Sample	#of Test Sample	Java			Python		
			BLEU-4	METEOR	ROUGE-L	BLEU-4	METEOR	ROUGE-L
CodeBERT	8000	1000	19.91	25.11	34.34	20.56	33.37	33.19
CodeT5	8000	1000	22.45	28.98	41.98	28.82	37.98	39.52
CodeT5+	8000	1000	28.82	39.79	49.31	34.67	46.98	47.34
SCLA	/	1000	34.34	50.66	60.71	37.34	52.61	57.49

Table 6: The performance of our proposed method and the baseline model was evaluated on Java and Python datasets.

Approach	Prompt Component	BLEU-4	METEOR	ROUGE-L	BLEURT
SCLA	ALL	6.09	25.80	29.45	46.63
	-CFG	5.21	27.77	27.94	45.90
	-IF	4.42	25.43	26.23	44.56
	-Id&MGV	5.62	25.47	29.01	46.32
	-ALL	0.85	20.28	18.53	43.53

Table 7: Ablation study. Effect of Semantic Augmentation on Gemini-1.5-Pro Generated Summarization. CFG is Control Flow Graph, IF is Inner Function, Id&MGV is Identifiers&Global Member Variables.

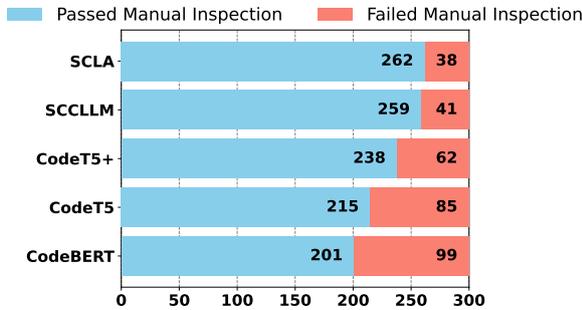


Figure 4: Human Evaluation Results of 300 Code Summarizations Generated by SCLA and the Baseline.

of each component on Gemini-1.5-Pro’s code summarization performance. Five experiments were performed by selectively removing semantic elements used in semantically enhanced prompts, all under Zero-Shot learning. The results, shown in Table 7, highlight the importance of inner function ordering, function call graphs, identifiers, and global member variables. Removing inner functions resulted in performance drops of 27.42%, 1.43%, 10.93%, and 4.44% across key metrics. Excluding the function call graph caused declines of 14.45

(BLEU-4), 5.13 (ROUGE-L), and 1.57 (BLEURT). The removal of identifiers and global member variables led to decreases of 7.72%, 1.28%, 1.49%, and 0.65%. These findings confirm that inline functions and call graphs are crucial for improving Gemini-1.5-Pro’s code summarization. Additionally, **experiments underscore the significance of global member variables in maintaining semantic consistency, while function call graphs provide a structural understanding that directly influences summarization accuracy.** Removing these elements disrupts both the coherence and comprehensiveness of the generated summarization. Ultimately, the ablation study validates the necessity of each semantic component for optimizing code summarization performance in scenarios.

4.7 Human Evaluation of Summarization Generated by SCLA and the Baseline

To assess the summarization generated by SCLA, we randomly selected 300 samples from the smart contract code summarization generated by SCLA and baseline models for manual evaluation. This evaluation focused on similarity, conciseness, and completeness, categorizing the summarization as usable or unusable. To reduce subjectivity and bias, six volunteer evaluators, all Chinese graduate students with experience in smart contract development, were recruited and briefed on the research and evaluation standards. The results, shown in Figure 4, reveal that SCLA generated the fewest unusable summarization (38), outperforming all

Example				
	<pre> #Input Function Code function transferDataOwnership (address _addr) onlyOwner public { data.transferOwnership(_addr); } #inner function code function transferOwnership(address _newOwner) public onlyOwner { _transferOwnership(_newOwner); } </pre>			
Approach	Coment	BLEU-4	METEOR	ROUGE-L
Ground Truth	Transfer ownership of data contract to <code>_addr</code> . <code>_addr</code> address.	NA	NA	NA
CodeBERT	Transfer ownership of an address to another. <code>_addr</code> address The address to transfer to.	51.20	28.00	67.00
CodeT5	Allows the owner to transfer control of the contract to an address. <code>_addr</code> The address to transfer ownership to.	42.48	26.64	47.62
CodeT5+	Allows the owner to transfer control of data contract to <code>_addr</code> . <code>_addr</code> The address.	60.68	41.67	60.00
SCLA	Transfers ownership of the data contract to <code>_addr</code> .	70.77	72.70	75.00

Table 8: An example illustrating the effectiveness of SCLA.

451 baseline models. **These findings demonstrate**
452 **that SCLA is more likely to generate satisfactory**
453 **smart contract code summarization, reducing**
454 **the chances of low-quality outputs.**

455 4.8 Case Study

456 Upon reviewing the results, we found that the
457 SCLA prompt includes crucial information for
458 effective summarization. Table 8 highlights the
459 differences between real-world smart contract ab-
460 stracts and Summarization generated by Code-
461 BERT, CodeT5, CodeT5+, and SCLA. CodeBERT
462 identifies key terms like "transfer," "ownership,"
463 and "address," but lacks clarity, with ambiguous
464 pronoun references and repetition of the transfer
465 concept. CodeT5 captures "onlyOwner" but over-
466 looks broader global semantics, rendering the sec-
467 ond sentence redundant. CodeT5+ addresses this
468 limitation with more precise terminology, such
469 as identifying the object as a "data contract." **In**
470 **contrast, SCLA’s Summarization aligns more**
471 **closely with real-world Summarization, being**
472 **both more concise and semantically accurate,**
473 **omitting redundancy for a much clearer, more**
474 **refined, and contextually precise structure.**

475 4.9 Generalization Study

476 To evaluate the generalization and representa-
477 tiveness of SCLA, we selected 10,000 samples
478 each from Java and Python in the CodeSearch-
479 Net dataset (Husain et al., 2019). From these,
480 1,000 samples per language were randomly cho-
481 sen as test sets. We compared SCLA’s per-
482 formance against state-of-the-art fine-tuned mod-
483 els—CodeT5, CodeT5+, CodeBERT, and SC-

484 CLLM—using BLEU-4, METEOR, and ROUGE-
485 L. On the Java dataset, SCLA outperformed
486 CodeT5+ with improvements of 19%, 12%, and
487 23% in BLEU-4, METEOR, and ROUGE-L, re-
488 spectively. On the Python dataset, SCLA showed
489 gains of 7%, 18%, and 21%, respectively. We fur-
490 ther compared SCLA to SCCLLM on both datasets
491 using BLEU-4 and BLEURT. With LLMs GPT-
492 4o, Gemini-1.5-Pro, and Claude-3.5-Sonnet, SCLA
493 consistently outperformed SCCLLM, achieving an
494 average improvement of 31.14 in BLEU-4 and
495 21.38 in BLEURT. **These results demonstrate**
496 **that SCLA’s control flow-based prompts gener-**
497 **alize effectively to Java and Python, enhancing**
498 **the LLMs’ ability to capture code structure and**
499 **improve summarization quality.**

500 5 Conclusion

501 We propose that control flow graphs enhance
502 LLMs’ understanding of smart contract code se-
503 mantics, and experiments confirm their positive im-
504 pact on code comprehension. Ablation studies as-
505 sess the contribution of each prompt component to
506 summarization quality. SCLA is a framework that
507 combines LLMs with control flow prompts, out-
508 performing six baseline models. The experiments
509 show that, compared to other baseline models,
510 SCLA significantly improves BLEU-4, METEOR,
511 ROUGE-L, and BLEURT scores with improve-
512 ments of 30.34%, 23.15%, 16.74%, and 14.86%,
513 respectively. We also extended SCLA to Java
514 and Python code, further improving summarization
515 and providing new insights for advancing LLM-
516 generated code summarization.

517 Limitations

518 Our framework enhances Gemini-1.5-Pro’s under-
519 standing using function call graphs. However,
520 Gemini-1.5-Pro struggles with deep call stacks
521 or circular calls. Figure. 5 shows that circular
522 chains, like transferFrom → removeTokenFrom →
523 ownerOf → isApprovedOrOwner → transferFrom,
524 confuse the model, leading to misinterpretations
525 and incorrect summarization. In contrast, Gemini-
526 1.5-Pro handles typical tree structures even with a
527 depth of 5. Further research is needed to explore
528 the impact of loop calls and depth on-call interpre-
529 tation.

530 Another key challenge in using LLMs for smart
531 contract code summarization is the potential expo-
532 sure of test data during pre-training. Since general-
533 purpose LLMs like GPT-4o and Gemini-1.5-Pro
534 are not publicly accessible, direct verification of
535 this exposure is difficult. Additionally, LLMs’
536 memorization capability can produce artificially
537 high scores if prior summarizations are retained.
538 We also analyzed the effect of few-shot learning
539 on SCLA’s performance in Section 3. Our results
540 show that SCLA outperforms the baseline with a
541 Three-Shot setup, while performance gains plateau
542 at five shots, with a 1.5x increase in computational
543 cost.

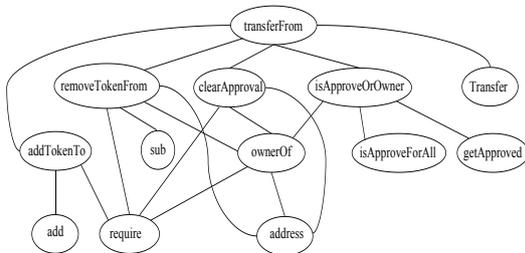


Figure 5: An Example of a Function Call Graph in Which Gemini-1.5-Pro Has Difficulty Understanding the Call Information.

544 References

545 Toufique Ahmed and Premkumar Devanbu. 2022. Mul-
546 tilingual training for software engineering. In *Pro-
547 ceedings of the 44th International Conference on
548 Software Engineering*, pages 1443–1455.

549 Toufique Ahmed and Premkumar Devanbu. 2023.
550 Few-shot training llms for project-specific code-
551 summarization. In *Proceedings of the 37th
552 IEEE/ACM International Conference on Automated
553 Software Engineering*, pages 1–5.

Toufique Ahmed, Kunal Suresh Pai, Premkumar De-
vanbu, and Earl Barr. 2024. Automatic semantic
augmentation of language model prompts (for code
summarization). In *Proceedings of the IEEE/ACM
46th International Conference on Software Engineer-
ing*, pages 1–13. 554 555 556 557 558 559

Satanjeev Banerjee and Alon Lavie. 2005. Meteor: An
automatic metric for mt evaluation with improved cor-
relation with human judgments. In *Proceedings of
the acl workshop on intrinsic and extrinsic evaluation
measures for machine translation and/or summariza-
tion*, pages 65–72. 560 561 562 563 564 565

Xiangping Chen, Peiyong Liao, Yixin Zhang, Yuan
Huang, and Zibin Zheng. 2021. Understanding code
reuse in smart contracts. In *Proceedings of the 2021
IEEE International Conference on Software Analysis,
Evolution and Reengineering*, pages 470–479. IEEE. 566 567 568 569 570

Matthias De Lange, Rahaf Aljundi, Marc Masana, Sarah
Parisot, Xu Jia, Aleš Leonardis, Gregory Slabaugh,
and Tinne Tuytelaars. 2022. A continual learning sur-
vey: Defying forgetting in classification tasks. *IEEE
Transactions on Pattern Analysis and Machine Intel-
ligence*, pages 3366–3385. 571 572 573 574 575 576

Zhangyin Feng, Daya Guo, Duyu Tang, Nan Duan, Xi-
aocheng Feng, Ming Gong, Linjun Shou, Bing Qin,
Ting Liu, Daxin Jiang, and Ming Zhou. 2020. Code-
BERT: A pre-trained model for programming and
natural languages. In *Findings of the Association
for Computational Linguistics: EMNLP 2020*, pages
1536–1547. 577 578 579 580 581 582 583

Hamel Husain, Ho-Hsiang Wu, Tiferet Gazit, Miltiadis
Allamanis, and Marc Brockschmidt. 2019. Code-
searchnet challenge: Evaluating the state of semantic
code search. *arXiv preprint arXiv:1909.09436*. 584 585 586 587

Dechao Kong, Xiaoqi Li, and Wenkai Li. 2024. Char-
acterizing the solana nft ecosystem. In *Companion
Proceedings of the ACM on Web Conference*, pages
766–769. 588 589 590 591

Satpal Singh Kushwaha, Sandeep Joshi, Dilbag Singh,
Manjit Kaur, and Heung-No Lee. 2022. System-
atic review of security vulnerabilities in ethereum
blockchain smart contract. *IEEE Access*, pages 6605–
6621. 592 593 594 595 596

Gang Lei, Donghua Zhang, Jianmao Xiao, Guodong
Fan, Yuanlong Cao, and Zhiyong Feng. 2024. Fmcf:
A fusing multiple code features approach based
on transformer for solidity smart contracts source
code summarization. *Applied Soft Computing*, page
112238. 597 598 599 600 601 602

Ze Qin Liao, Sicheng Hao, Yuhong Nan, and Zibin
Zheng. 2023. Smartstate: Detecting state-reverting
vulnerabilities in smart contracts via fine-grained
state-dependency analysis. In *Proceedings of the
32nd ACM SIGSOFT International Symposium on
Software Testing and Analysis*, pages 980–991. 603 604 605 606 607 608

- 609 Chin-Yew Lin. 2004. Rouge: A package for automatic
610 evaluation of summaries. In *Text summarization*
611 *branches out*, pages 74–81.
- 612 Zhenguang Liu, Peng Qian, Xiang Wang, Lei Zhu, Qin-
613 ming He, and Shouling Ji. 2021. Smart contract
614 vulnerability detection: From pure neural network to
615 interpretable graph feature and expert pattern fusion.
616 In *Proceedings of the Thirtieth International Joint*
617 *Conference on Artificial Intelligence*, pages 2751–
618 2759.
- 619 Kishore Papineni, Salim Roukos, Todd Ward, and Wei-
620 Jing Zhu. 2002. Bleu: a method for automatic evalu-
621 ation of machine translation. In *Proceedings of the*
622 *40th annual meeting of the Association for Computa-*
623 *tional Linguistics*, pages 311–318.
- 624 Nils Reimers and Iryna Gurevych. 2020. Making
625 monolingual sentence embeddings multilingual using
626 knowledge distillation. In *Proceedings of the Con-*
627 *ference on Empirical Methods in Natural Language*
628 *Processing*, pages 4512–4525.
- 629 Thibault Sellam, Dipanjan Das, and Ankur Parikh. 2020.
630 BLEURT: Learning robust metrics for text genera-
631 tion. In *Proceedings of the 58th Annual Meeting of*
632 *the Association for Computational Linguistics*, pages
633 7881–7892.
- 634 Yue Wang, Hung Le, Akhilesh Gotmare, Nghi Bui, Jun-
635 nan Li, and Steven Hoi. 2023. CodeT5+: Open code
636 large language models for code understanding and
637 generation. In *Proceedings of the Conference on*
638 *Empirical Methods in Natural Language Processing*,
639 pages 1069–1088.
- 640 Yue Wang, Weishi Wang, Shafiq Joty, and Steven C.H.
641 Hoi. 2021. Codet5: Identifier-aware unified pre-
642 trained encoder-decoder models for code understand-
643 ing and generation. In *Proceedings of the Conference*
644 *on Empirical Methods in Natural Language Process-*
645 *ing*, pages 8696–8708.
- 646 Zhen Yang, Jacky Keung, Xiao Yu, Xiaodong Gu,
647 Zhengyuan Wei, Xiaoxue Ma, and Miao Zhang. 2021.
648 A multi-modal transformer-based code summariza-
649 tion approach for smart contracts. In *Proceedings*
650 *of the IEEE/ACM 29th International Conference on*
651 *Program Comprehension*, pages 1–12.
- 652 Shenhui Zhang, Wenkai Li, Xiaoqi Li, and Boyi Liu.
653 2022. Authros: Secure data sharing among robot
654 operating systems based on ethereum. In *Proceed-*
655 *ings of the IEEE 22nd International Conference on*
656 *Software Quality, Reliability and Security*, pages 147–
657 156. IEEE.
- 658 Junjie Zhao, Xiang Chen, Guang Yang, and Yiheng
659 Shen. 2024. Automatic smart contract comment
660 generation via large language models and in-context
661 learning. 168.