

Comparative Analysis of Different Prompting Methods for Large Language Model

Hyoyeon Lee¹

¹ School of Computer Science, University of Bristol

1 Introduction

This document presents a comparative analysis of various prompting techniques based on their memory usage and execution time when applied to Gemini-Flash 1.5. The study evaluates three distinct methods: one-shot prompting, Chain-of-Thought (CoT) prompting, and few-shots prompting. Notably, the current implementation does not incorporate counterexamples from the CBMC (C Bounded Model Checker) with the new prompts. Instead, it uses 'simple restart' method which uses the same prompt and relies on the possibility of obtaining a different response from the LLM than previously returned. Additionally, the number of reattempt has been restricted to three to avoid the hallucination of the LLM and endless reiteration of the software. This comparison aims to provide insights into the efficiency and performance of each prompting strategy, thereby guiding the appropriate usage of prompting methods.

2 Experimental Setup

In this section, I define the overall procedure of the experiment and its environment settings. For the exact prompt, check [Section 3](#).

2.1 Implementation

I tested three different prompting methods to input code into Gemini Flash-1.5, receive the generated response, and verify its correctness using CBMC. The prompt has been refined from looping until it individually checks all proof harness functions generated to create a new function, 'combined_proof_harness()', to verify its correctness all at once. This approach reduced the running time as it removed the need of creating a separate JSON file of CBMC results per every proof harness functions. A 'simple restart' method has been applied when reattempting to query the LLM.

2.2 LLM

Gemini Flash 1.5 has been used for generating proof harness functions, where its temperature has been fixed to default temperature of 0.25 to ensure its determinism. Other hyperparameters have been set to default as well.

2.3 CBMC

CBMC 6.0.1 has been used to verify the correctness of generated proof harness functions. Identical CBMC command '`cbmc FILE_NAME.c -function combined_proof_harness -no-standard-checks -no-malloc-may-fail -verbosity 8 -unwind 3 -trace -json-ui`' has been used for every prompting method. As the updated version of CBMC was released, new flags were used to restore the default value of the previous version of CBMC.

2.4 Benchmarks

The benchmarks were collected from real-world datasets on Kaggle, primarily related to LeetCode problems. All programs were written in C, using only standard libraries to avoid errors.

Prompt Method	Project	Success	#Iteration	Runtime	Memory Usage	LoC	#Functions
One shot	prompt.c	True	2	28.7236	90093	42	5
Few shots	prompt.c	False	3	42.7236	318866	42	5
CoT	prompt.c	True	1	4.8829	81507	42	5

Table 1: Benchmark details and format

Details on the benchmark are provided in Table 1. The verification result, indicated by 'Success', is marked as 'True' for successful verification from the CBMC and 'False' otherwise. The total number of iterations for each project is listed under '#Iteration'. The iteration count is limited to three; projects that reached the maximum iteration number without 'True' success result are considered unsuccessful. The total runtime and memory usage are reported by 'Runtime' and 'Memory Usage' in seconds and bytes, respectively. The runtime and memory usage have not been averaged to better reflect the LLM's capability to generate correct and safe C codes in a given number of iteration.

3 Prompting Methods

The following are the three different sample prompt formats: one-shot prompting, few-shots prompting, and Chain-of-Thought prompting, to query Gemini Flash 1.5 for generating proof harness functions. New constraint regarding library declarations has been added due to frequent failures in proof harness function generation, 'assert.h' in particular.

3.1 One-shot Prompting

```
1 # Preamble
2 You are given a C program. We need to create a proof harness function.
3 # Code generation example - You are given three examples.
4
5 Q: Write a method "void proof_harness_withdraw()" that tests method withdraw below for all possible
6 inputs.
7
8 // Define the Account structure
9 struct Account {{
10     unsigned short bal;
11 }};
12
13 // Function to withdraw an amount from an account
14 void withdraw(struct Account *account, unsigned short amount) {{
15     unsigned short de = account->bal;
16     account->bal = de - amount;
17 }}
18
19 A:
20 struct Account {{
21     unsigned short bal;
22 }};
23
24 void withdraw(struct Account *account, unsigned short amount) {{
25     unsigned short de = account->bal;
26     account->bal = de - amount;
27 }}
28
29 void proof_harness_withdraw() {{
30     struct Account *account = (struct Account *)malloc(sizeof(struct Account));
31     __CPROVER_assume(account != NULL); // Ensure account is not NULL
32
33     unsigned short amount;
34
35     __CPROVER_assume(account->bal >= 0);
36     __CPROVER_assume(amount > 0);
37     __CPROVER_assume(account->bal >= amount);
38
39     unsigned short initial_balance = account->bal;
40
41     withdraw(account, amount);
42     assert(account->bal == initial_balance - amount);
43
44     free(account);
45 }}
46
47 # Instruction
48 Give me a proof harness code of the below C code.
49
50 # Query
51 Q: Write method "void proof_harness_{{function_name}}()" that tests every methods including 'main()'
52 below for all possible inputs. Then write a method "void combined_proof_harness()" that calls every
53 other proof
54 harness methods generated.
55
56 {query_code}.
57
58 # Constraints
```

```

59 Here are some constraints that you should respect:
60 - Give me only the translated code, don't add explanations or anything else.
61 - Use only safe C.
62 - Do not use custom generics. # fuzzer limitation
63 - Do not remove any code from the original code you have received.
64 - Ensure that all libraries needed, including <assert.h>, are declared at the beginning of the code.

```

Listing 1: One-shot Prompting Format

3.2 Few-shots prompting

```

1
2 # Preamble
3 You are given a C program. We need to create a proof harness function.
4
5 # Code generation example
6 ## Example 1
7 Q: Write a method "void proof_harness_withdraw()" that tests method withdraw below for all
8 possible inputs.
9 // Define the Account structure
10 struct Account {{
11     unsigned short bal;
12 }};
13
14 // Function to withdraw an amount from an account
15 void withdraw(struct Account *account, unsigned short amount) {{
16     unsigned short de = account->bal;
17     account->bal = de - amount;
18 }}
19
20 A:
21 struct Account {{
22     unsigned short bal;
23 }};
24
25 void withdraw(struct Account *account, unsigned short amount) {{
26     unsigned short de = account->bal;
27     account->bal = de - amount;
28 }}
29
30 void proof_harness_withdraw() {{
31     struct Account *account = (struct Account *)malloc(sizeof(struct Account));
32     __CPROVER_assume(account != NULL); // Ensure account is not NULL
33
34     unsigned short amount;
35
36     __CPROVER_assume(account->bal >= 0);
37     __CPROVER_assume(amount > 0);
38     __CPROVER_assume(account->bal >= amount);
39
40     unsigned short initial_balance = account->bal;
41
42     withdraw(account, amount);
43     assert(account->bal == initial_balance - amount);
44
45     free(account);
46 }}
47 ## end Example 1
48
49 ## Example 2
50 Q: Write method "void proof_harness_{{function_name}}()" that tests every methods below for all
51 possible inputs.
52 #include <stdio.h>
53 #include <stdlib.h>
54
55 struct node{{
56     struct node *leftNode;

```

```

56     int data;
57     struct node *rightNode;
58 };
59
60 struct node *newNode(int data){{
61     struct node *node = (struct node *)malloc(sizeof(struct node));
62     node->leftNode = NULL;
63     node->data = data; node->rightNode = NULL;
64     return node;
65 }}
66
67 int main(void){{
68     return 0;
69 }}
70
71 A:
72 #include <stdio.h>
73 #include <stdlib.h>
74 #include <assert.h>
75
76 struct node{{
77     struct node *leftNode;
78     int data;
79     struct node *rightNode;
80 }};
81
82 struct node *newNode(int data){{
83     struct node *node = (struct node *)malloc(sizeof(struct node));
84     node->leftNode = NULL;
85     node->data = data; node->rightNode = NULL;
86     return node;
87 }}
88
89 int main(void){{
90     return 0;
91 }}
92
93 void proof_harness_newNode() {{
94     int data;
95     __CPROVER_assume(data >= 0);
96     __CPROVER_assume(data <= 2147483647);
97     struct node *node = newNode(data);
98     assert(node != NULL);
99     assert(node->leftNode == NULL);
100    assert(node->data == data);
101    assert(node->rightNode == NULL);
102    free(node);
103 }}
104
105 ## end Example 2
106
107 ## Example 3
108 Q: Write method "void proof_harness_{{function_name}}()" that tests every methods below for all
    possible inputs.
109
110     #include <stdio.h>
111
112     struct node{{
113         struct node *leftNode;
114         int data;
115         struct node *rightNode;
116     }};
117
118     void inOrderTraversal(struct node *node){{
119         if (node == NULL) return;
120         inOrderTraversal(node->leftNode);
121         printf("\t%d\t", node->data);
122         inOrderTraversal(node->rightNode);

```

```

123 }}
124
125 void preOrderTraversal(struct node *node){{
126     if (node == NULL) return;
127     printf("\t%d\t", node->data);
128     preOrderTraversal(node->leftNode);
129     preOrderTraversal(node->rightNode);
130 }}
131
132 void postOrderTraversal(struct node *node){{
133     if (node == NULL) return;
134     postOrderTraversal(node->leftNode);
135     postOrderTraversal(node->rightNode);
136     printf("\t%d\t", node->data);
137 }}
138
139 int main(void){{
140     return 0;
141 }}
142
143 A:
144
145 #include <stdio.h>
146 #include <assert.h>
147
148 struct node{{
149     struct node *leftNode;
150     int data;
151     struct node *rightNode;
152 }};
153
154 void inOrderTraversal(struct node *node){{
155     if (node == NULL) return;
156     inOrderTraversal(node->leftNode);
157     printf("\t%d\t", node->data);
158     inOrderTraversal(node->rightNode);
159 }}
160
161 void preOrderTraversal(struct node *node){{
162     if (node == NULL) return;
163     printf("\t%d\t", node->data);
164     preOrderTraversal(node->leftNode);
165     preOrderTraversal(node->rightNode);
166 }}
167
168 void postOrderTraversal(struct node *node){{
169     if (node == NULL) return;
170     postOrderTraversal(node->leftNode);
171     postOrderTraversal(node->rightNode);
172     printf("\t%d\t", node->data);
173 }}
174
175 int main(void){{
176     return 0;
177 }}
178
179 void proof_harness_inOrderTraversal() {{
180     struct node *node = (struct node *)malloc(sizeof(struct node));
181     __CPROVER_assume(node != NULL);
182     node->leftNode = NULL;
183     node->rightNode = NULL;
184     __CPROVER_assume(node->data >= 0 && node->data <= 100);
185     inOrderTraversal(node);
186     free(node);
187 }}
188
189 void proof_harness_preOrderTraversal() {{
190     struct node *node = (struct node *)malloc(sizeof(struct node));

```

```

191     __CPROVER_assume(node != NULL);
192     node->leftNode = NULL;
193     node->rightNode = NULL;
194     __CPROVER_assume(node->data >= 0 && node->data <= 100);
195     preOrderTraversal(node);
196     free(node);
197 }
198
199 void proof_harness_postOrderTraversal() {{
200     struct node *node = (struct node *)malloc(sizeof(struct node));
201     __CPROVER_assume(node != NULL);
202     node->leftNode = NULL;
203     node->rightNode = NULL;
204     __CPROVER_assume(node->data >= 0 && node->data <= 100);
205     postOrderTraversal(node);
206     free(node);
207 }}
208
209 ## end Example 3
210
211 # Instruction
212     Give me a proof harness code of the below C code.
213
214 # Query
215     Q: Write method "void proof_harness_{{function_name}}()" that tests every methods including 'main()'
216     below for all possible inputs. Then write a method "void combined_proof_harness()" that calls every
217     other proof
218     harness methods generated.
219
220     {query_code}
221
222 # Constraints
223     Here are some constraints that you should respect:
224     - Give me only the translated code, don't add explanations or anything else.
225     - Use only safe C.
226     - Do not use custom generics. # fuzzer limitation
227     - Do not remove any code from the original code you have received.
228     - Ensure that all libraries needed, including <assert.h>, are declared at the beginning of the code.

```

Listing 2: Prompt Format

The prompt has been set to include three examples to prevent excessive length and minimise the risk of hallucinations from the LLM. As it will be demonstrated in the results, for few-shot prompting, which is the longest prompt format, longer codes significantly increase the likelihood of hallucinations and failure rates.

3.3 Chain-of-Thought (CoT)

```

1 # Preamble
2     You are given a C program. We need to create a proof harness function.
3 # Code generation example
4     Q: Write a method "void proof_harness_withdraw()" that tests method withdraw below for all possible
5     inputs.
6
7     // Define the Account structure
8     struct Account {{
9         unsigned short bal;
10    }};
11
12    // Function to withdraw an amount from an account
13    void withdraw(struct Account *account, unsigned short amount) {{
14        unsigned short de = account->bal;
15        account->bal = de - amount;
16    }}
17
18    A: You first identify which functions there are. In this case, we have a single function called '
19    withdraw'. Since

```

```

18  withdraw is a function that needs to be tested, we create a proof_harness_withdraw function for this
19  purpose and
20  implement to test the correctness of the function 'withdraw'.
21
22  struct Account {{
23      unsigned short bal;
24      }};
25
26  void withdraw(struct Account *account, unsigned short amount) {{
27      unsigned short de = account->bal;
28      account->bal = de - amount;
29  }}
30
31  void proof_harness_withdraw() {{
32      struct Account *account = (struct Account *)malloc(sizeof(struct Account));
33      __CPROVER_assume(account != NULL); // Ensure account is not NULL
34
35      unsigned short amount;
36
37      __CPROVER_assume(account->bal >= 0);
38      __CPROVER_assume(amount > 0);
39      __CPROVER_assume(account->bal >= amount);
40
41      unsigned short initial_balance = account->bal;
42
43      withdraw(account, amount);
44      assert(account->bal == initial_balance - amount);
45
46      free(account);
47  }}
48
49
50  # Instruction
51  Give me a proof harness code of the below C code.
52
53  # Query
54  Q: Write method "void proof_harness_{{function_name}}()" that tests every methods including 'main()'
55  below for all possible inputs. Then write a method "void combined_proof_harness()" that calls every
56  other proof
57  harness methods generated.
58
59  {query_code}
60
61  # Constraints
62  Here are some constraints that you should respect:
63  - Give me only the translated code, don't add explanations or anything else.
64  - Use only safe C.
65  - Do not use custom generics. # fuzzer limitation
66  - Do not remove any code from the original code you have received.
67  - Ensure that <assert.h> is declared when using assert.
68  """

```

Listing 3: Generated Proof Harness Function

4 Results

I run an experiment for each of the three prompts for a total of 366 proof harness functions generation and correctness verification experiment. The input C programs were validated to ensure them using safe C codes by identifying those that returned conversion errors and the pre-processing them accordingly. The results have been rounded to four decimal places.

4.1 Success Rate with regards to Line of Code and number of functions

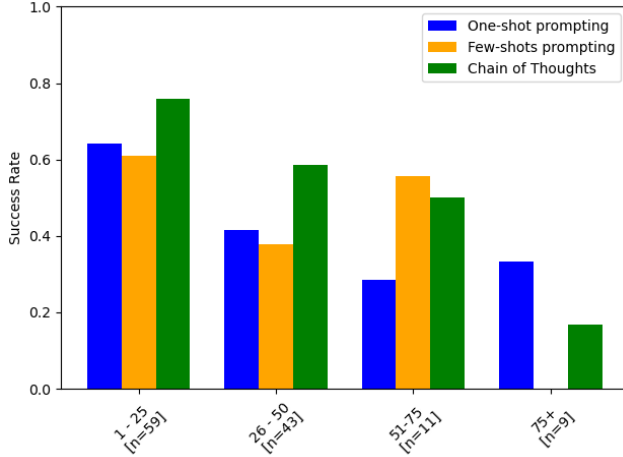


Figure 1: Success Rate of Prompting Methods by Line of Code

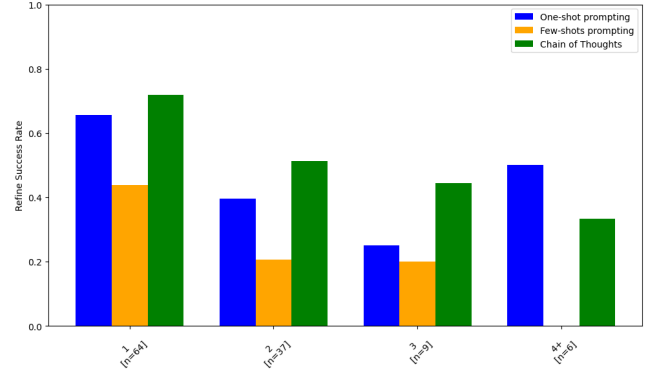


Figure 2: Success Rate of Prompting Methods by Number of Functions

Prompting methods \ LoC	1 ~ 25	26 - 50	51-75	75+
One-shot prompting	35.9375	58.5366	71.4286	66.6667
Few-shots prompting	38.8889	62.0690	44.4444	100
Chain of Thoughts	24.0741	41.4634	50	83.3333

Table 2: Failure Rates of Prompting Methods by Line of Code (percentage)

Figure 1 illustrates the success rates of different prompting methods in relation to the length of the input code, measured in lines of code. The success rates are presented for four ranges of code length: 1-25 lines, 26-50 lines, 51-75 lines, and 75+ lines. The results indicate a general trend where longer prompts and more extensive input code correlate with lower success rates. Specifically:

- **ONE-SHOT PROMPTING** demonstrates a significant decrease in success rate as the line of code increase. It achieves its highest success rate of 64.06% with the shortest code line group and drops to its lowest rate of 28.57% within the 51-75 lines range. Regarding the number of functions, one-shot prompting shows consistent performance across different numbers of functions, with the highest success rate observed for single-function code (64.71%) and the lowest for code with three functions (33.33%).
- **FEW-SHOTS PROMPTING** exhibits its highest success rate of 61.11% with 1-25 lines and the lowest with the longest line of codes. Notably, there exists an anomaly in the 51-75 lines group, where the success rate increases by approximately 18.25% from the previous group. This indicates that the few-shots prompting performs better performance with moderately long code. Apart from the line of code, few-shots prompting shows lower performance as the number of functions increases. Its success rate significantly drops compared to other prompting methods, with the lowest performance observed for code with four or more functions (16.67%).
- **CHAIN-OF-THOUGHT PROMPTING** maintains the highest success rates among the tested prompting methods for the first three groups. Similarly to other prompting methods, the performance gradually declines for longer code, but shows a highest average success rate among all three. In terms of the number of functions, Chain-of-Thought consistently shows high performance across different function groups. Although there is a decline in success rate as the number of functions increases, Chain-of-Thought remains the most effective method compared to the others.

4.2 Refine Success Rate

One-shot prompting	Few-shots prompting	CoT prompting
5.8333 %	1.8692 %	3.7383 %

Table 3: Refine success rate

Adopting the 'simple restart' strategy, the success rate of refined proof harness function verification was much lower than expected for all three methods. One-shot prompting showed the best ability to return valid proof harness functions, followed by Chain-of-Thoughts prompting with a 3.74% lower success rate. The few-shots prompting method demonstrated the lowest probability of fixing bugs in its proof harness function generation.

4.3 Average Runtime and Memory Usage

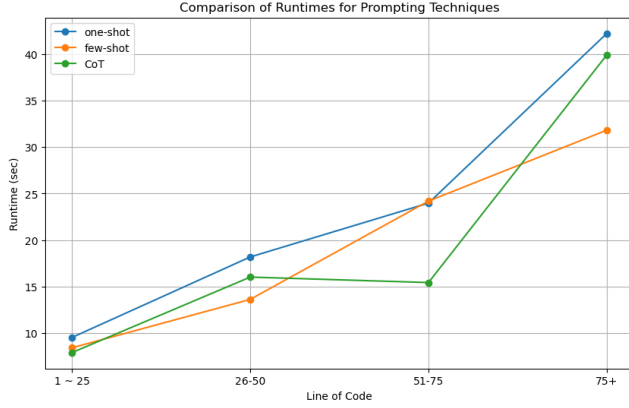


Figure 3: Average Runtime (sec)

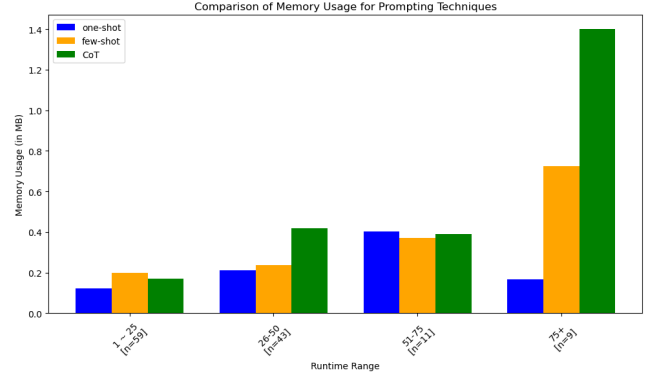


Figure 4: Average Memory Usage (MB)

Before the results are delved, it must be noted that the Chain-of-Thought prompting strategy had two extremely large outliers, which caused the average runtime to increase dramatically. The runtime and memory usage figure has been extracted from the very start until the program finished its final iteration.

The overall average runtime increased as the length of the input code got longer, as well as the memory usage for all three prompting strategies. One-shot prompting shows a gradual increase in runtime as the line of code increases, with a significant jump for the longest code group (75+ lines). Few-shots prompting maintains a relatively stable runtime, with it showing almost a linear relationship between the line of code and the runtime. Chain-of-Thought (CoT) prompting has the shortest runtime for most code lengths but also exhibits a sharp increase for the longest code group, potentially influenced by two extreme outliers. In terms of memory usage, one-shot prompting displays the lowest memory usage across all code lengths, few-shots prompting has moderate memory usage increasing steadily with longer code, and Chain-of-Thought (CoT) prompting shows the highest memory usage, especially for the longest code group, indicating a significant increase in resource consumption.

5 Conclusion

The comparative analysis of different prompting methods for proof harness function verification using Gemini Flash 1.5 reveals significant insights into their performance. One-shot prompting emerged as the most effective method in terms of both success rate and resource efficiency, although its performance notably decreased with longer lines of code and multiple functions. Chain-of-Thought prompting, despite having two large outliers that affected the average runtime, maintained a high success rate across varying code lengths and function numbers, showing consistent performance. Few-shots prompting, while relatively stable in runtime, demonstrated the lowest success rate and the highest memory usage, particularly as the code length increased. Given these results, when LLMs are not provided with feedback from previous failures, they are more likely to succeed if the input is of moderate length and contains fewer functions. The 'simple restart' strategy yielded lower-than-expected success rates for all three methods, with One-shot prompting leading, followed by Chain-of-Thought prompting, and Few-shots prompting trailing behind. Overall, these findings suggest that One-shot prompting is generally preferable for tasks requiring efficient and reliable proof harness function verification, while Chain-of-Thought prompting offers a balanced alternative, especially when runtime and memory usage are less of a concern.

6 Further Discussions

6.1 Using Different Feedback Strategies

Since the current feedback strategy does not show commendable performance in generating valid and safe proof harness functions, trying different feedback strategies such as providing counterexamples could be considered. Incorporating counterexamples into the

feedback prompt may help in guiding the model to correct its previous mistakes more effectively, potentially increasing the overall success rate and robustness of the generated proof harness functions. By leveraging more detailed and constructive feedback, it may be possible to enhance the model's learning process and improve the quality of its outputs.