

Abstract

Supervisors: Chunchuan Lyu Authors: Yanzhuo Jin, Zhaohui Yang, Jieran He, Zichen Qiu, Xu Chen

Program synthesis generates computer programs based on given problem specifications. Large Language Models have been proven to be effective. However, models are still struggling to handle complicated tasks involving multiple steps. Recently, a multi-turn code generation benchmark has been proposed. They provide pre-defined individual steps in natural language. Instead of starting from a pre-decomposed steps, we aim to automatically decompose the task into individual steps. We reverse-engineered the multi-turn benchmark to include a short summary instruction. We empirically examined the ability of current models to automatically decompose the summary specification and complete the task given the decomposed steps.

Introduction

Our contributions :

- We construct a dataset based on MTPB by adding a summary instruction that summarized the multi-step instructions in MTPB
- We use prompt engineering to automatically decompose our summary prompt and enhance our result with the chain-of-thought prompting.

Findings:

- LLMs can automatically decompose a short instruction into individual steps
- LLMs can complete the code better with both generated and gold individual steps.
- LLMs may prefer copying example answers, but this can be alleviated via corrective prompting

Related Work

A. Large Language Models

Recently, autoregressive Transformer models[2] with billions of parameters pretrained on terabytes of data have been used to model the conditional distribution of natural language. By further training on a wide range of tasks, they can complete new tasks without training[3]. We will use GPT3.5-turbo in our study.

B. Multi-Turn Program Synthesis

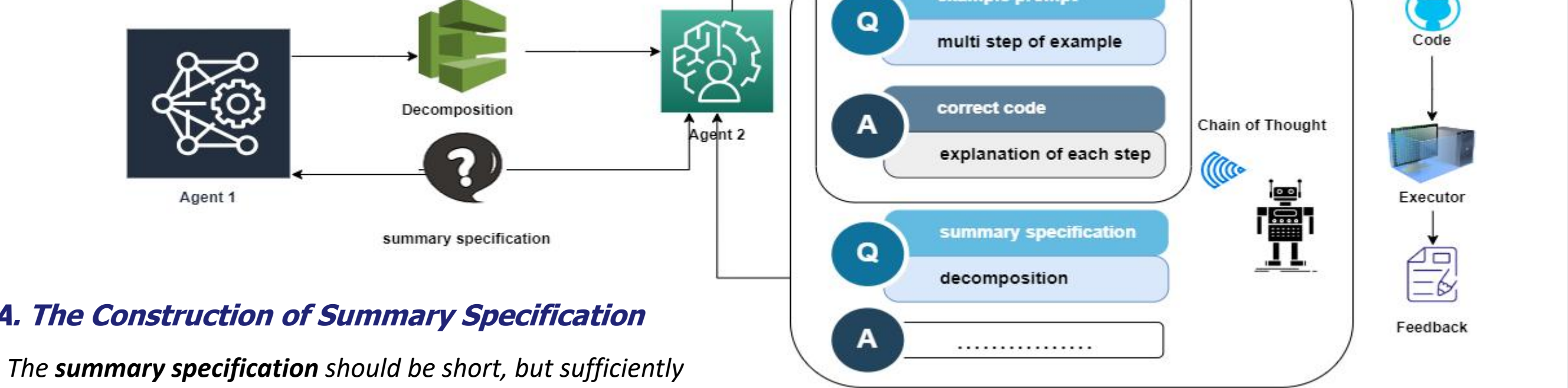
Program synthesis is the generation of computer programs based on problem specifications. A program specification could involve multi-turn steps. To investigate the multi-step paradigm for program synthesis, researchers introduced Multi-Turn Programming Benchmark (MTPB) which consists of 115 diverse problem sets that are factorized into multi-turn prompts[1]. While MTPB offers a manual breakdown of each step for a given problem, rarely users will specify coding steps in such detail. We aim to provide a more realistic setting where LLMs start with a short but sufficient problem specification.

C. Chain of Thought prompting

Chain of Thought (CoT) prompting provides demonstrations as exemplars in prompting. CoT teaches LLMs to solve problems by leveraging LLMs' ability to in-context learn new patterns. We use CoT exemplars to teach LLMs to decompose the summary specification and provide well-formatted code. [4]

Leveraging Large Language Models for Specification Decomposition in Code Generation

Dataset and Methods



A. The Construction of Summary Specification

The **summary specification** should be short, but sufficiently descriptive of the tasks defined by multiple steps. Five annotators manually read all the individual steps in all **MTPB** and one lead annotator checked the validity of summary specification with the help of LLMs. The average length of our summary specification is X and average length of all individual steps is Y. We find that some problems in MTPB are simple concatenation of almost irrelevant processing steps that cannot be summarized. We split the data into two groups: summarized and un-summarizable. There are X summarized cases and Y un-summarizable cases.

B. Auto Decomposition with Chain of Thought

Baseline: The GPT3.5 model is given only the summary specification to complete the task, and it should directly produce the code.

Automatic decomposition (Auto-D): The GPT3.5 model is given only the summary specification to complete the task, but it is asked to provide a decomposition of summary specification into each steps then provide the code.

Gold decomposition (Gold-D): Given both the summary specification and the individual steps in original MTBP.

Automatic decomposition with Chain of Thought (Auto-D + CoT): In addition to 2, provided with a example of decomposition process and answer.

Gold decomposition with Chain of Thought (Gold-D + CoT): In addition to 3, provided with a example of decomposition and answer. (Figure5)

C. Corrective Prompting to Alleviate Copying from CoT Examples

Experiments and Results

We find that when using CoT, GPT3.5 may prefer copying from examples. This cause the model to deviate from solving the given problem. We explicitly add “You must not copy the code about the question example” into the prompts to avoid this issue.

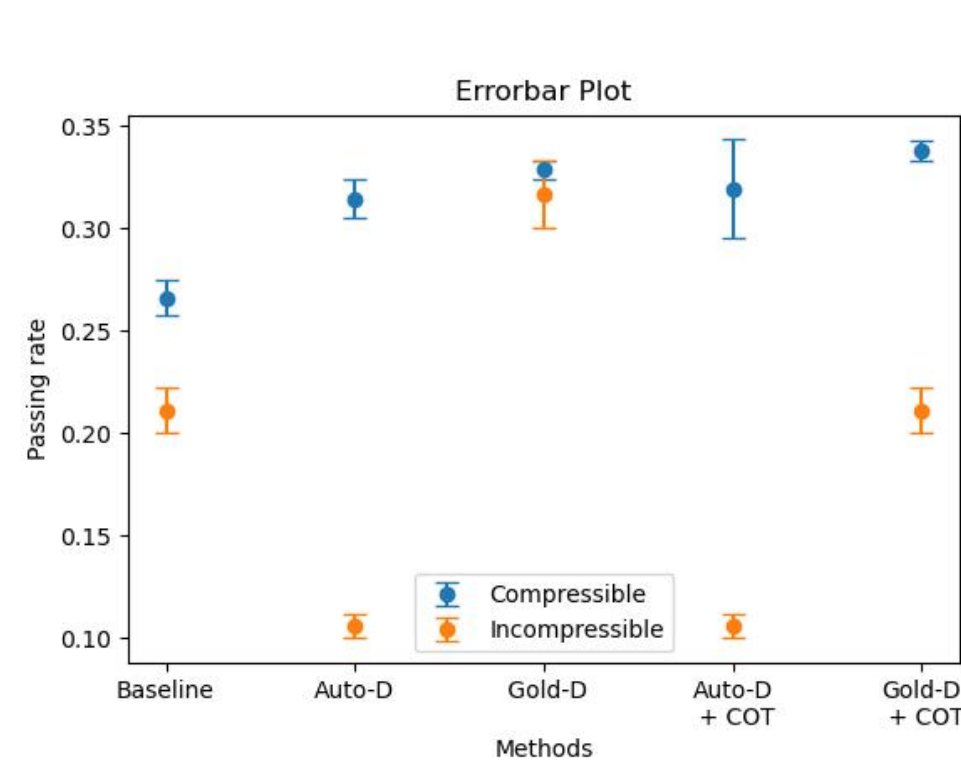


Figure8: The error bar of the experiments in two epoches

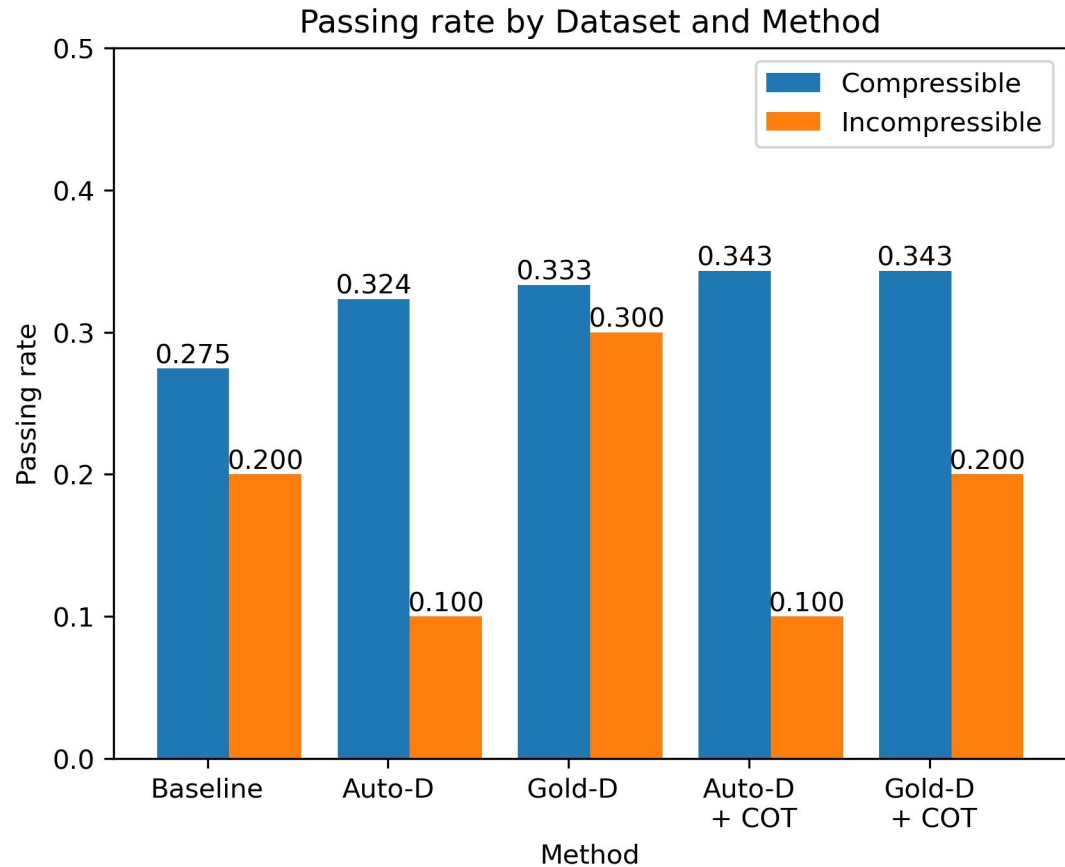


Figure9: The accuracy Pass@2 of five methods based on the EMTPB.

Conclusion and Discussion

In this study, we show that LLMS can automatically decompose user intention into individual steps, and its’ performance of code completion can benefit from relying on decomposed steps. However, our study is limited in terms of the models we have tested and the size and variety of coding task. It is possible that better models such as GPT4 can decompose the instruction better or it might be more capable or solving the problem without explicit decomposition.

Data and Code: <https://github.com/nakupenda-c/Auto-Prompt-Decomposition-with-Chain-of-Thought-in-Large-Language-Models-for-Program-Synthesis>



code report

Reference

- [1] Nijkamp E, Pang B, Hayashi H, et al. Codegen: An open largelanguage model for code with multi-turn program synthesis[J]. arXiv preprint arXiv:2203.13474, 2022.
- [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, L ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In Advances in Neural Information Processing Systems, volume 30, 2017.
- [3] Brown T, Mann B, Ryder N, et al. Language models are few-shot learners[J]. Advances in neural information processing systems, 2020, 33: 1877-1901.
- [4] Wei J, Wang X, Schuurmans D, et al. Chain-of-thought prompting elicits reasoning in large language models[J]. Advances in Neural Information Processing Systems, 2022, 35:24824-24837.

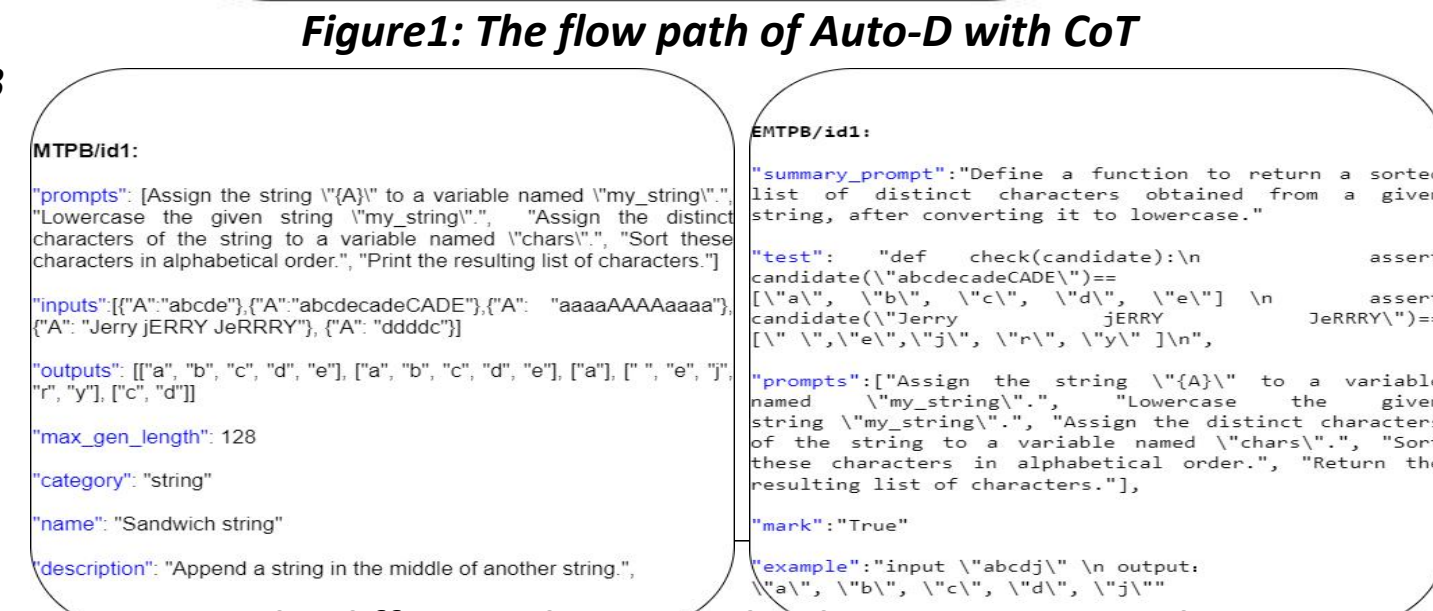


Figure2: The difference between the datasets MTPB and EMTPB

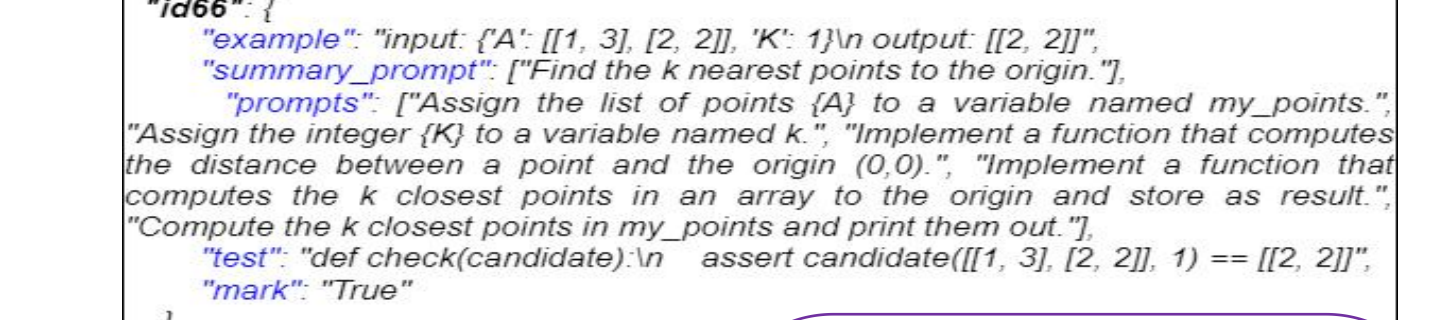


Figure3: The datasets of EMTPB: id66 to verify the method feasibility

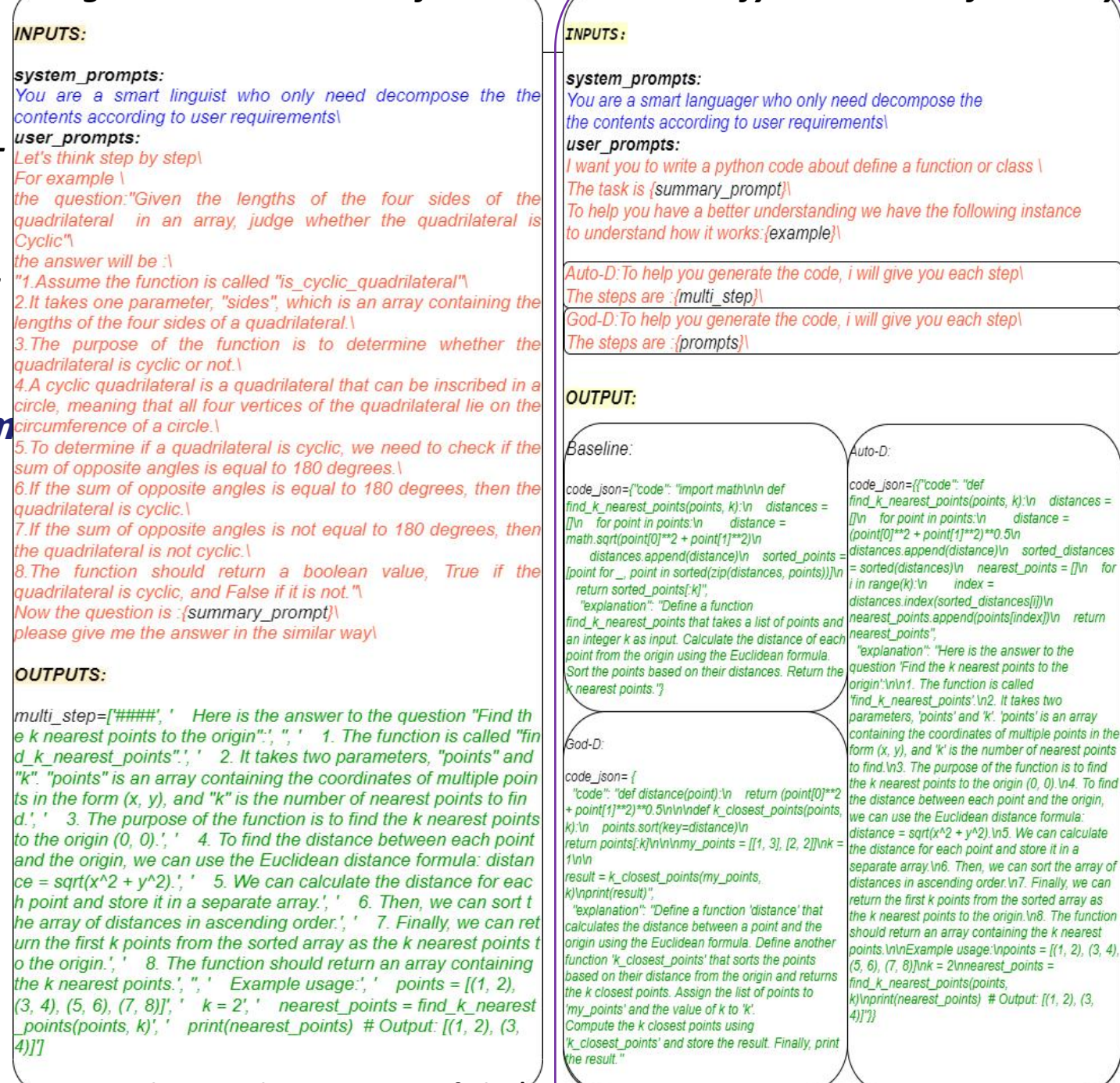


Figure4: The Auto decomposition of id66

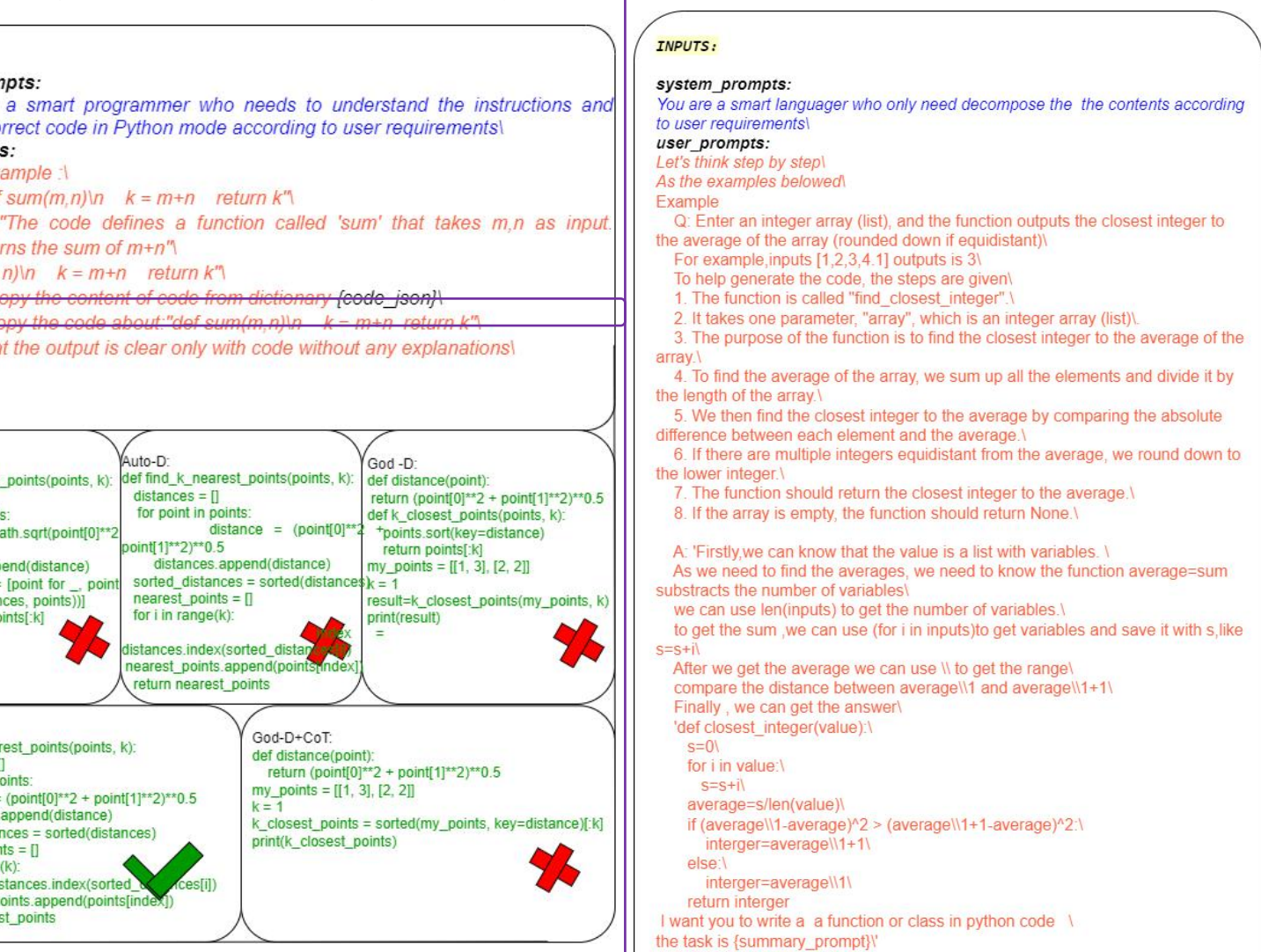


Figure6: The five code of id66

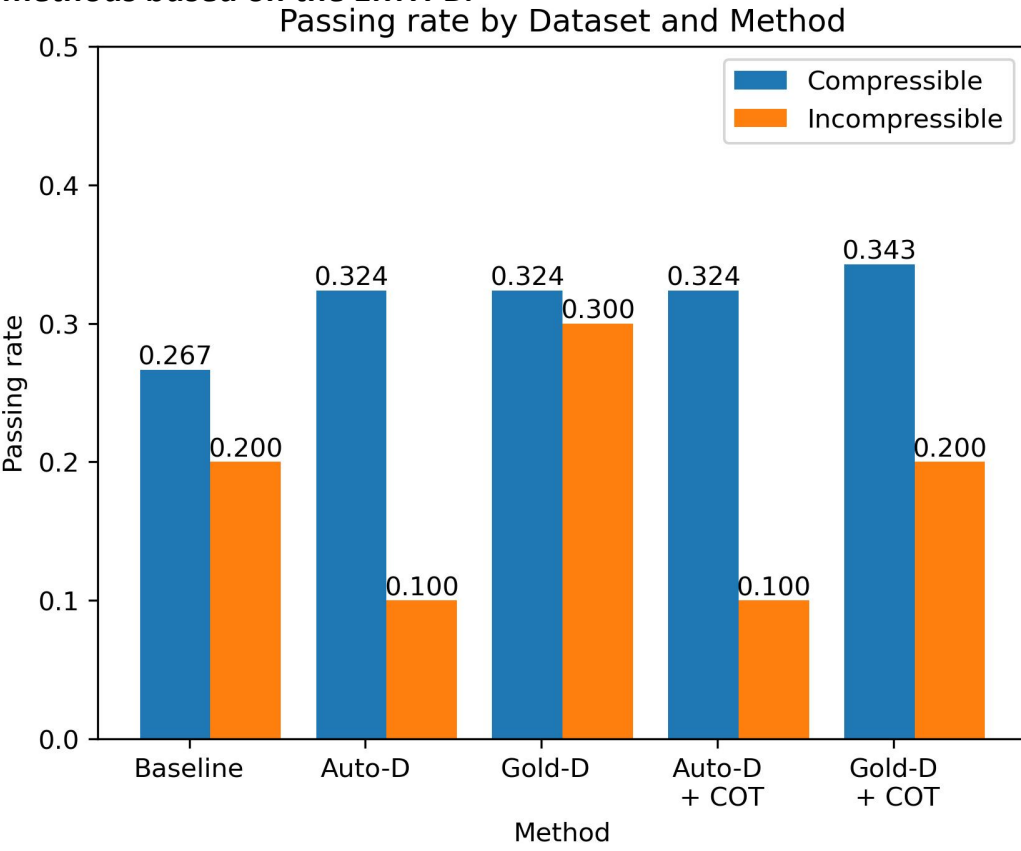


Figure9: The accuracy Pass@1 of five methods without “You must not copy the code about the question example” based on the EMTPB.

Figure5: The five code string outputs of id66