

Automatic Bi-modal Question Title Generation for Stack Overflow with Prompt Learning

Shaoyu Yang · Xiang Chen · Ke Liu · Guang Yang · Chi Yu

Received: date / Accepted: date

Abstract When drafting question posts for Stack Overflow, developers may not accurately summarize the core problems in the question titles, which can cause these questions to not get timely help. Therefore, improving the quality of question titles has attracted the wide attention of researchers. An initial study aimed to automatically generate the titles by only analyzing the code snippets in the question body. However, this study ignored the helpful information in their corresponding problem descriptions. Therefore, we propose an approach SOTitle+ by considering bi-modal information (i.e., the code snippets and the problem descriptions) in the question body. Then we formalize the title generation for different programming languages as separate but related tasks and utilize multi-task learning to solve these tasks. Later we fine-tune the pre-trained language model CodeT5 to automatically generate the titles. Unfortunately, the inconsistent inputs and optimization objectives between the pre-training task and our investigated task may make fine-tuning hard to fully explore the knowledge of the pre-trained model. To solve this issue, SOTitle+ further prompt-tunes CodeT5 with hybrid prompts (i.e., mixture of hard and soft prompts). To verify the effec-

Shaoyu Yang
School of Information Science and Technology,
Nantong University, Nantong, China
E-mail: shaoyuyoung@gmail.com

Xiang Chen
School of Information Science and Technology
Nantong University, Nantong, China
E-mail: xchencs@ntu.edu.cn

Ke Liu
School of Information Science and Technology
Nantong University, Nantong, China
E-mail: aurora.ke.liu@outlook.com

Guang Yang
College of Computer Science and Technology,
Nanjing University of Aeronautics and Astronautics, Nanjing, China
E-mail: novelyg@outlook.com

Chi Yu
School of Information Science and Technology
Nantong University, Nantong, China
E-mail: yc_struggle@163.com

tiveness of SOTitle+, we construct a large-scale high-quality corpus from recent data dumps shared by Stack Overflow. Our corpus includes 179,119 high-quality question posts for six popular programming languages. Experimental results show that SOTitle+ can significantly outperform four state-of-the-art baselines in both automatic evaluation and human evaluation. In addition, our ablation studies also confirm the effectiveness of component settings (such as bi-modal information, prompt learning, hybrid prompts, and multi-task learning) of SOTitle+. Our work indicates that considering bi-modal information and prompt learning in Stack Overflow title generation is a promising exploration direction.

Keywords Question Title Generation · Bi-modal Information · Code Snippet · Problem Description · Prompt Learning · Multi-task Learning

1 Introduction

Recently, developers widely relied on Stack Overflow (SO) as a primary resource for coding and related information online. With millions utilizing the platform, it serves as a go-to source for finding well-crafted solutions to programming problems. Moreover, SO has evolved into a comprehensive knowledge repository, enabling developers to enhance their programming skills by accessing a wealth of high-quality questions and corresponding answers. The success of SO largely depends on the user’s willingness to answer others’ questions. Generally speaking, high-quality questions can increase opportunities for getting assistance. It is also beneficial for developers to find solutions to similar problems, and ultimately for the entire community as it promotes knowledge sharing behavior (Anderson et al 2012; Jin and Servant 2019; Cao et al 2021; Gao et al 2023).

In support of developers crafting high-quality questions, Stack Overflow has developed comprehensive quality assurance guidelines¹ for its community members. However, a significant number of questions posted on SO still fall short of meeting the platform’s quality standards. These low-quality questions lack clarity, precision, or completion, deterring potential experts from providing immediate answers. Consequently, this hampers the advancement and sharing of knowledge within the SO community. One reason for these low-quality questions is that users cannot draft informative question titles (Correa and Sureka 2013; Trienes and Balog 2019; Tóth et al 2019; Gao et al 2020) as users may not be familiar with the knowledge and terminology related to the problem or have weak writing skills. Therefore, to improve title quality and reduce the manual efforts for quality maintenance of the SO community, automatic question title generation for Stack Overflow is urgently needed. For the convenience of subsequent descriptions, we define **Stack Overflow Question Title Generation** as **SOQTG**. For the SOQTG task, we aim to automatically generate titles, which can serve as concise and informative summaries for capturing the essence of a question posted on the Stack Overflow platform. SOQTG is more challenging as compared to general text summarization (El-Kassas et al 2021) or source code summarization (Ahmad et al 2020; Iyer et al 2016; Li et al 2022). First, question titles need to be succinct yet informative, capturing the essence of the problem described in the post body. Second, question titles often involve specialized technical terms and domain-specific language prevalent in software development, making them distinct from general text summarization. Third, generating informative titles requires understanding the context provided in the post body and associated code snippets, highlighting the core issue to attract relevant answers. Finally, question ti-

¹ <https://stackoverflow.com/help/how-to-ask>

titles should align with community norms and conventions prevalent on the Stack Overflow platform to ensure increased visibility and relevance.

To our best knowledge, (Gao et al 2020) were the first to automatically generate question titles by analyzing the code snippets in the question body. They formalized the SOQTG task as a sequence-to-sequence learning problem and then proposed an approach Code2Que. Specifically, an attention mechanism is enhanced to perform better content selection, a copy mechanism is utilized to handle the out-of-vocabulary problem, and a coverage mechanism is used to avoid meaningless repetition. However, they ignored the valuable information in relevant natural language descriptions, which are called problem descriptions in our study. As shown in Fig. 1, two posts in Stack Overflow have the same code snippet. However, these two posts have different question titles due to different problem descriptions. Based on this motivation example, we find the problem description of the question post can provide valuable information for question title generation, which sometimes cannot be provided by the code snippet. Therefore, considering bi-modal information (i.e., code snippet and problem description) in the question body can help to generate high-quality question titles.

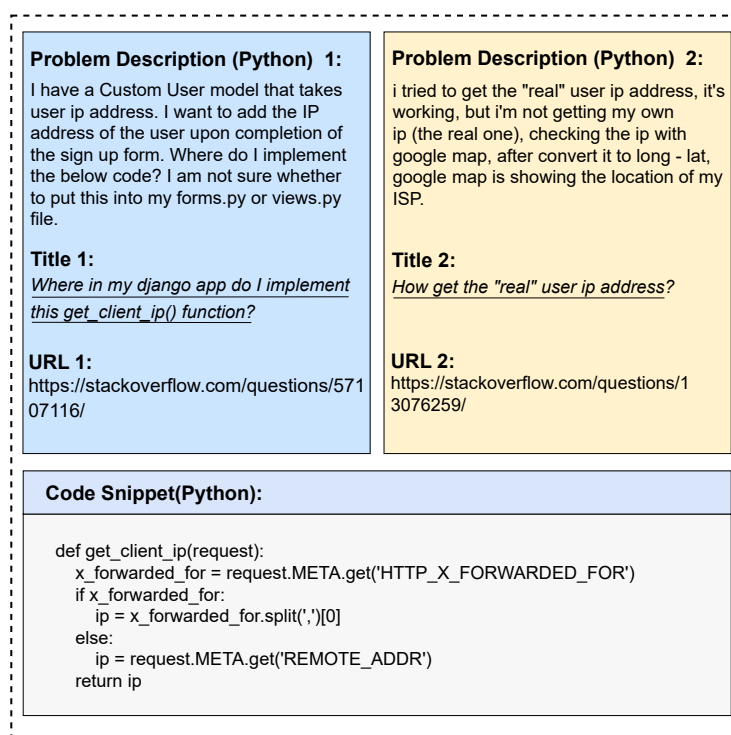


Fig. 1 Two posts from Stack Overflow, which have the same code snippet but different problem descriptions

For the SOQTG task, a popular approach is to treat this task as the downstream task and then fine-tune a Pre-trained Language Model (PLM) on the collected corpus of this task. However, there exist several differences between the pre-training task and the downstream task (Wang et al 2022). As shown in Fig. 2 (a), we use CodeT5 (Wang et al 2021), which is a PLM for code understanding and generation, as an example. CodeT5 is pre-trained using

the MSP (Masked Span Prediction) optimization objective. The input to MSP involves bi-modal information, which is a mixture of code snippets and corresponding natural language descriptions. Then the model is trained to predict randomly masked input tokens. However, when CodeT5 is fine-tuned for the downstream task SOQTG, it loses the special token `[MASK]` and fails to identify the bi-modal information as CodeT5 is trained to a sequence-to-sequence model, which is shown in Fig. 2 (b). Moreover, the optimization objective of the SOQTG task shown in Section 2.2.3 is also task-specific. The inconsistent inputs and optimization objectives between the pre-training task and the downstream task may make fine-tuning hard to elicit the knowledge of the pre-trained model.

Recently, prompt tuning (Liu et al 2023c) has been introduced to bridge the gap between pre-training and fine-tuning (Schick and Schütze 2021; Lester et al 2021; Wang et al 2022; Huang et al 2022; Zhu et al 2023). Specifically, prompt tuning aims to explicitly guide the PLM on how to adapt its learned representations for the specific downstream task. Taking the SOQTG task as an example, instead of directly mixing bi-modal information for fine-tuning (Liu et al 2022; Zhang et al 2022), we use prompt tuning to make the model identify bi-modal information by adding additional natural language descriptions. Fig. 2(c) demonstrates the concept of prompt-tuning on our investigated SOQTG task in a visual way. In this example, we use two natural language sentences (i.e., “The problem description is:” and “The code snippet is:”) to distinguish the problem description and the code snippet. Then, a prompt (such as “Generate the question title: `[MASK]`”) is added after the input sentence, and `[MASK]` is the ground-truth title to be predicted. The PLM predicts the title at the masked token position, as it does in the pre-training stage. Therefore, by adding the prompt, we can reformulate the SOQTG task into an MSP task, aligning the objective with the pre-training stage. This task of reformulation can help fully utilize the task-specific knowledge hidden in the pre-trained model.

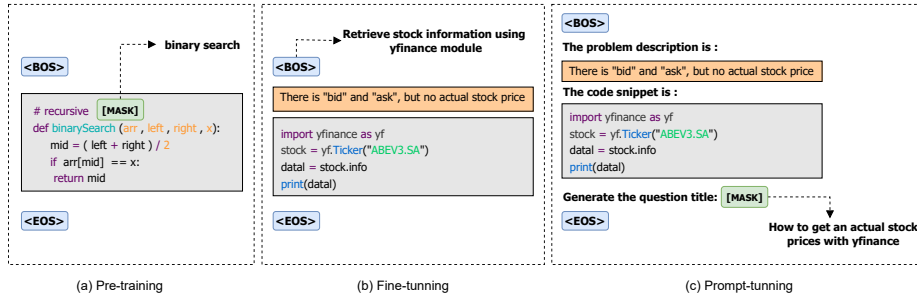


Fig. 2 Illustration on the process of pre-training, fine-tuning, and prompt tuning of the SOQTG task. We use `[MASK]` and `[BOS]` to denote two special tokens in CodeT5.

We propose the approach SOTitle+ based on the above research motivation analysis. Specifically, we first construct the corpus from Stack Overflow. To ensure the quality of the corpus, we use three heuristic selection rules to filter low-quality question posts. Since we formalize question title generation for different programming languages as separate but related tasks, we utilize multi-task learning to simultaneously solve these tasks, which can alleviate the insufficient training data issue for low-resource programming languages. When constructing the model, we use bi-modal information as the input. Moreover, we use hybrid prompts, which is a trade-off between keeping task-related tokens and introducing additional

training parameters, to perform prompt-tuning on CodeT5. In the model prediction phase, based on the trained pre-trained model based on CodeT5 and hybrid prompts, the pre-trained model is prompted to generate the question title with the input of the problem description and code snippet.

To verify the effectiveness of SOTitle+, we construct a large-scale high-quality corpus from recent shared data dumps provided by Stack Overflow. After using the selection rules, our corpus contains 179,119 high-quality question posts for six popular programming languages (i.e., Python, Java, C#, JavaScript, HTML, and PHP). In our empirical study, we compare SOTitle+ with state-of-the-art baselines via automatic evaluation (i.e., Rouge (LIN 2004), BLEU (Papineni et al 2002), METEOR (Banerjee and Lavie 2005), CIDEr (Vedantam et al 2015)) and human study. Specifically, we first select Code2Que proposed by (Gao et al 2020) as the first state-of-the-art question title generation baseline. Then we select two recent question title generation approaches proposed by Zhang et al. (i.e., CCBERT (Zhang et al 2022) and M3NSCT5 (Zhang et al 2023)) as the other two baselines. Finally, we consider our preliminary work SOTitle (Liu et al 2022) as the final baseline. Final comparison results of both experimental study and human study show that SOTitle+ can generate higher-quality titles than these four baselines. Moreover, a set of ablation studies also confirms the effectiveness of using bi-modal information, multi-task learning, prompt tuning, and hybrid prompts in SOTitle+.

We extend our preliminary study, which appears as a research paper in SANER 2022 (Liu et al 2022), as follows: (1) **Approach Extension.** Compared with SOTitle (Liu et al 2022), We use prompt tuning to replace fine tuning. Specifically, we separate the bi-modal information in the question body with prompt templates, which can help to effectively identify bi-modal information. Then we use the pre-trained model CodeT5 (Wang et al 2021) for prompt tuning. Finally, we conduct the ablation study to confirm the effectiveness of using the hybrid prompts in SOTitle+. (2) **Updated Corpus.** We evaluate the effectiveness of SOTitle+ by considering recent shared data dumps from Stack Overflow. Moreover, we also consider two more popular programming languages (i.e., HTML and PHP). Finally, We evaluate the performance of the trained model by considering the temporal relationships of posts when we split the constructed corpus. This can guarantee that the test set contains the newest posts and relieve the data leakage problem. (3) **More Baselines.** We further consider three baselines, including our preliminary study (Liu et al 2022) and two other state-of-the-art question title generation approaches CCBERT (Zhang et al 2022) and M3NSCT5 (Zhang et al 2023). (4) **More Performance Measures.** We further consider three text overlap-based measures to evaluate the model performance: BLEU (Papineni et al 2002), METEOR (Banerjee and Lavie 2005), and CIDEr (Vedantam et al 2015). (5) **More Discussions.** We first evaluate the effectiveness of SOTitle+ for two low-resource programming languages (i.e., Ruby and Go). We second compare our approach SOTitle+ with the Large Language Model (i.e., ChatGPT) on the SOQTG task in the zero-shot scenario. Finally, we analyze the limitations of SOTitle+ and discuss potential improvement directions for the SOQTG task.

To our best knowledge, the main contributions of our preliminary study and the extensions in this study are summarized as follows:

- **Direction.** We provide two feasible research directions for improving Stack Overflow question title generation. The first direction is to consider multi-modal information (such as problem description, and code snippet) in the question posts. The second direction is to consider prompt learning. Based on the promising performance of our study, we call for more follow-up studies on these two directions for further exploration and advancements in this task.

- **Approach.** We propose a novel approach SOTitle+ based on the bi-modal information . We apply prompt learning to SOTitle+ and design the novel hybrid prompts. SOTitle+ adopts multi-task learning to this task for different programming languages and a pre-trained model CodeT5 is utilized.
- **Corpus.** We construct a high-quality corpus, which contains 179,119 high-quality question posts for six popular programming languages.
- **Evaluation.** We conduct both automatic evaluation and human evaluation to show the competitiveness of SOTitle+. Moreover, we conduct a set of ablation studies to show the component setting rationality of SOTitle+.
- **Tool.** Based on SOTitle+, we developed a browser tool, which can assist developers in generating question titles for Stack Overflow. We share the video demonstration at YouTube².

Open Science. To support the open science community, we share our scripts, detailed experimental results, and the developed plugin tool in the GitHub repository³. Moreover, we also shared our dataset and trained model on Zenodo⁴.

The rest of this paper is organized as follows. Section 2 shows the framework of SOTitle+ and the details of each phase. Section 3 shows the experimental setup, including research questions and their design motivation, experimental subject, performance measures, baselines, implementation details, and running platform. Section 4 performs a comparison with baselines in the automatic evaluation and human evaluation and conducts a set of ablation studies. Section 5 discusses the effectiveness of using multi-task learning, question title generation in the low-resource scenario, the comparison with ChatGPT, the limitations of SOTitle+, and the main threats to the validity of our empirical study. Section 6 analyzes related work for the SOQTG task, and prompt learning and its applications to software engineering tasks. Section 7 concludes our work and shows potential future directions for SOQTG.

2 Our Proposed Approach SOTitle+

The overall framework of SOTitle+ is shown in Fig. 3. Specifically, SOTitle+ consists of three phases: corpus construction phase, model construction phase, and model prediction phase. In the rest of this section, we show the details for each phase.

2.1 Corpus Construction Phase

In this phase, we aim to construct the corpus by mining shared data dumps from Stack Overflow⁵. Due to a large number of posts in Stack Overflow, we refer to popularity statistics of tags⁶ and then determine the top-six popular programming languages (i.e., Python, Java, JavaScript, C#, PHP, and HTML). Moreover, we select two other programming languages (i.e., Go and Ruby), which can be used to evaluate the effectiveness of SOTitle+ in the low-resource scenario (discussed in Section 5.2).

² https://www.youtube.com/watch?v=_KgUISAT74M

³ <https://github.com/shaoyuyoung/SOTitlePlus>

⁴ <https://zenodo.org/records/10656359>

⁵ <https://archive.org/download/stackexchange>, downloaded in March 2023

⁶ <https://stackoverflow.com/tags?tab=popular>, accessed in March 2023

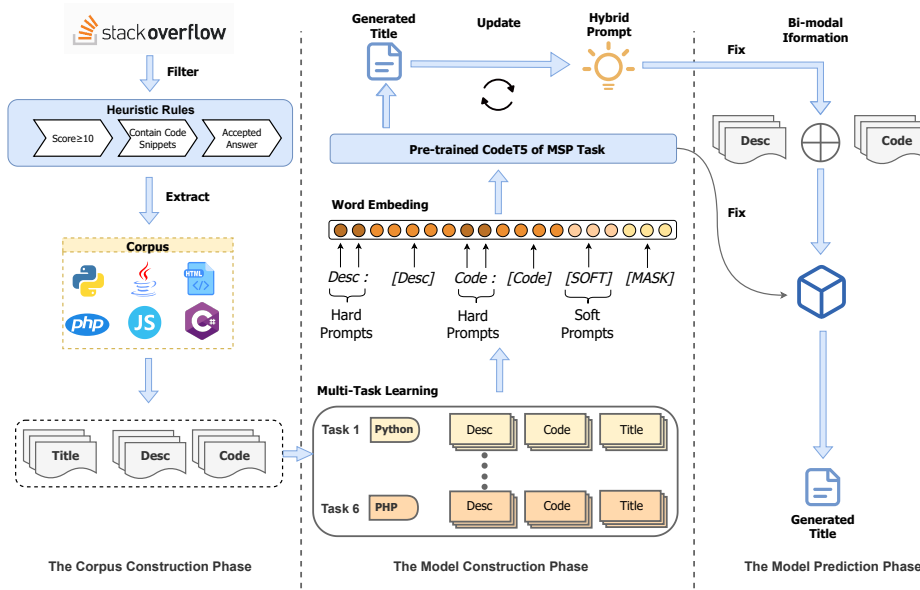


Fig. 3 Framework of our proposed approach SOTitle+

Since there exist many low-quality posts in our initial gathered corpus, we consider three heuristic selection rules by following previous studies (Islam et al 2019; Yazdaniya et al 2021; Yin et al 2018). For the convenience of introducing these selection rules, we use an example post⁷ shown in Fig. 4 to illustrate the main components of a Stack Overflow post. Specifically, this question post is related to the Python programming language, which contains a short question title, the problem description with the corresponding code snippet, the number of votes, and multiple tags. After introducing these components, we show the details of these three selection rules and their design motivations as follows.

- **Rule 1:** The votes of the question post should not be smaller than 10. This rule is motivated by the previous study on Stack Overflow post analysis on deep learning bug characteristic analysis (Islam et al 2019). They reduced the number of posts by considering the number of votes to focus on the high-quality posts. Through our preliminary empirical analysis, we found setting this value to 10 can help to ensure a certain level of post quality for the top six different popular programming languages.
- **Rule 2:** The question post should contain code snippets. The reason for setting this rule is to verify the effectiveness of our approach by considering bi-modal information.
- **Rule 3:** The question post should have an accepted answer. This rule is motivated by previous studies on mining Stack Overflow (Yazdaniya et al 2021; Yin et al 2018). We used the existence of an accepted answer as a criterion to potentially indicate question post quality, assuming that such questions might have clearer question titles.

After applying these three selection rules, we select 1.62% of posts (i.e., 179,119 posts from 11,023,365 posts). Then we extract the related programming language, the question title, the problem description, and the code snippet as the quadruplet (Lang, Title, Description, Code) for each question post. Notice as question posts were stored as a unified HTML

⁷ <https://stackoverflow.com/questions/51560850>

How to unit test a POST method in python? ← Question Title

Asked 4 years, 11 months ago Modified 4 years, 10 months ago Viewed 23k times

Problem Description

10
Votes

I have a method that sends a POST containing a JSON to an Elasticsearch instance. I am trying to write a unit test that verify the contents of the sent JSON, but I am not sure how to go about that. Should I create a local server in python and have it verify the contents of the POST or something else? I currently have this:

```
class TestAnalytics(BaseTest):
    def test_post(self):
        info = {"test1": "value1", "test2": "value2"}
        resp = requests.post(config.tool_repo_urls['es_url'], data=json.dumps(info), headers={'Content-Type': 'application/json'})
        assert_equal(resp.status_code, 200) # verify valid response code
```

python unit-testing post ← Code Snippet

Share Edit Follow Tags

asked Jul 27, 2018 at 15:11
Citut
807 ●2 ●10 ●25

Fig. 4 A question post related to Python programming language

format in the data dump, we use the Python Library `lxml`⁸ to extract the text wrapped by the “<code>” tag as the code snippet and the remaining text as the problem description.

2.2 Model Construction Phase

In this phase, we aim to train the model for question title generation. Specifically, we first formalize question title generation for different programming languages as separate but related tasks and utilize multi-task learning (Liu et al 2019b; Yang et al 2023b) to solve these tasks. Then we consider the bi-modal information (i.e., problem description and code snippet), which is distinguished by hard prompts (e.g., “The problem description” and “The code snippet”). Meanwhile, we create soft prompts by replacing unimportant tokens with [SOFT] tokens. Based on the hybrid prompts (i.e., a mixture of the hard prompts and the soft prompts), our model is prompt-tuned based on a pre-trained Transformer model CodeT5 (Wang et al 2021).

2.2.1 Hybrid Prompt Template Construction

Prompt learning (Liu et al 2023c) uses the prompt template to modify the original input to generate the new input with a prompt. However, designing a high-quality prompt template for the specific downstream task is still a challenging task. According to the types of prompt tokens, prompt templates can be classified into three types: hard prompt, soft prompt, and hybrid prompt.

Hard Prompt. A prompt template mainly includes two types of slots: input slot and answer slot. For the SOQTG task, there are two input slots. The input slot [X] will be filled with the problem description and the input slot [Y] will be filled with the code snippet. The answer slot [Z] is eventually filled with the text generated by the PLM. Except for these

⁸ <https://lxml.de/>

input slots and the answer slot, the prompt tokens in the hard prompt (Liu et al 2023c) are words that humans can understand and are often designed by domain experts. Notice that there is no verbalizer for the generation task. An example of the hard prompt can be found in Fig. 2 (c). Therefore, we define the hard prompt for the SOQTG task as follows.

$$f_{\text{hard}} = \text{The problem description: } [X] \text{The code snippet: } [Y] \text{Generate the question title:} [Z] \quad (1)$$

Soft Prompt. The prompt tokens in the soft prompt (Liu et al 2023c) are words that humans cannot understand but continuous embeddings that can be tuned during the downstream task training. Compared to the hard prompt, the soft prompt can overcome the difficulty of designing appropriate hard prompts for a specific task. There are two types of soft prompt: vanilla soft prompt and prefix soft prompt (Wang et al 2022). For the SOQTG task, we consider the vanilla soft prompt and define the soft prompt as follows.

$$f_{\text{soft}} = [\text{SOFT}] [X] [\text{SOFT}] [Y] [\text{SOFT}] [Z] \quad (2)$$

Notice, these [SOFT] tokens are initialized by using the embeddings of natural words in the hard prompt template in our study. For example, the third [SOFT] token can be initialized by “generate the question title:”⁹. This setting can ensure the reproduction of our experimental results and help to find the optimal embeddings when compared with the random initialization method.

Hybrid Prompts. Hybrid Prompt is a combination of the hard prompt and the soft prompt. Specifically, tokens in the hard prompt are usually important tokens related to the investigated task, which cannot be tuned during the training. On the contrary, tokens in the soft prompt are unimportant and can be tuned during the downstream task training. Recently, the effectiveness of using the hybrid prompt has been confirmed in a recent study for vulnerability characteristic prediction (Li et al 2023). For the SOQTG task, we treat the task-related tokens (i.e., “The problem description:” and “The code snippet:”) as the important tokens, since these tokens can instruct the model to clearly separate the problem description and the code snippet from the bi-modal information. For the remaining tokens (i.e., “Generate the question title:”), we consider the vanilla soft prompt and replace them with [SOFT] tokens since there are many alternatives to this sentence, such as “Create a question heading:”, “Draft a post title:” or “Summarize a question post”. Based on the above analysis, we define the hybrid prompt for the SOQTG task as follows.

$$f_{\text{hybrid}} = \text{The problem description: } [X] \text{The code snippet: } [Y] [\text{SOFT}] [Z] \quad (3)$$

2.2.2 Multi-task Learning

Multi-task learning (Liu et al 2019b) is a technique that trains a single model to perform multiple related tasks simultaneously, which can help to improve scalability and efficiency. Multi-task learning has been shown to improve the generalization capabilities of natural language pre-training models by reusing the majority of model weights across multiple tasks (Zhang and Yang 2021). Since we treat question title generation for different programming languages as separate but related tasks, we prefix the input X of each programming language with a task-specific prefix: [LANG] (e.g., the [LANG] “JS:” denotes JavaScript)

⁹ An example used to clarify this process can be found in <https://github.com/shaoyuyoung/SOTitlePlus/blob/main/embeddings.md>

to make the model distinguish different tasks. Then the format of the input can be defined as follows:

$$X = [\text{LANG}] \oplus X \quad (4)$$

2.2.3 Prompt Tuning on Pre-trained Model CodeT5

The encoder of CodeT5 is made up of a stack of blocks, each of which mainly consists of two components: a self-attention layer which is followed by a feed-forward network.

Self-attention (Cheng et al 2016) is calculated based on a set of queries (Q), keys (K), and values (V). Then dot product is calculated between the queries and keys, Finally, each is divided by $\sqrt{d_k}$, and the softmax function is applied to get the weight of the corresponding value.

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad (5)$$

The encoder includes a fully connected feed-forward network (FFN) at each layer. The structure of the decoder is similar to the encoder but with a modification to its self-attention mechanism, which can avoid positions attending to subsequent positions. To facilitate deep architecture, the transformer incorporates residual connections and layer normalization between layers.

To capture subtle semantic differences in the input sequence, multi-head self-attention is employed across different representation subspaces. Multi-head self-attention is defined as follows:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h)W^O \quad (6)$$

where i -th head head_i can be calculated as follows:

$$\text{head}_i = \text{Attention}\left(QW_i^Q, KW_i^K, VW_i^V\right) \quad (7)$$

where W_i^Q, W_i^K, W_i^V are linear projection matrices.

Recently, (Niu et al 2023) showed that using CodeT5 (Wang et al 2021) could achieve state-of-the-art performance in many downstream tasks (such as source code summarization and code translation) for software engineering. CodeT5 was proposed by (Wang et al 2021), which was trained on a large-scale corpus consisting of bi-modal information (NL-PL). Specifically, Wang et al. randomly masked spans with arbitrary lengths and then predicted these masked spans combined with some sentinel tokens at the decoder, They referred to this task as **Masked Span Prediction (MSP)**, as illustrated in Fig. 2(a). The loss function of MSP is defined as follows:

$$\mathcal{L}_{MSP}(\theta) = \sum_{t=1}^k -\log P_{\theta}(x_t^{\text{mask}} | \mathbf{x}^{\setminus \text{mask}}, \mathbf{x}_{<t}^{\text{mask}}) \quad (8)$$

where θ denotes the model parameters, $\mathbf{x}^{\setminus \text{mask}}$ denotes the masked input, \mathbf{x}^{mask} denotes the masked sequence to predict from the decoder with k denoting the number of tokens in \mathbf{x}^{mask} , and $\mathbf{x}_{<t}^{\text{mask}}$ is the span sequence generated so far.

Since we utilize multi-task learning for processing these six different programming languages, we define the final loss function as follows:

$$\mathcal{L}_{Loss}(\theta) = (\sum_{i=1}^6 \mathcal{L}_{MSP}^{Task_i}(\theta))/6 \quad (9)$$

where each $Task_i$ represents different tasks (i.e., different programming languages).

2.3 Model Prediction Phase

In this phase, for the problem posts that require title generation, we extract the problem description and the code snippet from the post body. Subsequently, based on the tuned hybrid prompt, we can construct the input and feed it into the PLM CodeT5 to obtain the generated question title.

3 Experiment Setup

3.1 Research Questions

In our empirical study, we want to answer the following five research questions (RQs):

RQ1: Can our proposed approach SOTitle+ generate higher-quality question titles than four state-of-the-art baselines via automatic evaluation?

Motivation. To show the competitiveness of SOTitle+ when compared to the baselines, we consider four performance measures to compare the similarity between the generated question title and the ground-truth question title in this RQ.

RQ2: Whether using the bi-modal information in the question body can improve the performance of SOTitle+?

Motivation. In this RQ, we want to conduct the ablation study to investigate which input modal contributes the most to the performance of SOTitle+.

RQ3: Whether using the prompt-tuning paradigm can improve the performance of SOTitle+ than using the fine-tuning paradigm?

Motivation. In this RQ, we want to conduct the ablation study to investigate whether using the prompt-tuning paradigm can achieve better performance of SOTitle+ than using the fine-tuning paradigm.

RQ4: Whether using hybrid prompts can improve the performance of SOTitle+?

Motivation. Prompt template design is a challenging and open problem in both natural language processing and software engineering tasks. In this RQ, we want to conduct the ablation study to show the competitiveness of using hybrid prompts in our proposed approach SOTitle+.

RQ5: Can our proposed approach SOTitle+ outperform state-of-the-art baselines via human study?

Motivation. There exist some disadvantages to automatic evaluation measures. For example, the text overlap-based measures used in our study may not effectively reflect the similarity between the generated question title and the ground-truth question title from a semantic perspective (Hu et al 2020). Therefore, we want to conduct a human study to evaluate the competitiveness of SOTitle+ in this RQ.

3.2 Experimental Subject

After the corpus construction phase (introduced in Section 2.1), we finally select 43,056 posts for Python, 34,297 posts for Java, 37,900 posts for JavaScript, 34,846 posts for C#, 14,288 posts for HTML, and 14,732 posts for PHP. In the process of our corpus split, we consider the temporal relationship of posts. Specifically, splitting posts in chronological order is more applicable to the real-world application scenario and this setting can guarantee that the test set contains the latest posts. Therefore, this corpus split setting can relieve the data leakage problem caused by the homogeneous questions between the training set and the test set. Finally, we split the gathered corpus as the training set, the validation set, and the testing set according to 80%: 10%: 10% respectively. In Table 1, we show the detailed statistical information of our corpus split results. Finally, we also show the length statistics of code snippets, problem descriptions, and question titles in Table 2. Notice, some posts do not contain any problem description¹⁰, so the minimum value of the description length is 0.

Table 1 Statistical information of our corpus split results

Language	Training Set	Validation Set	Testing Set	Total
Python	34,444	4,306	4,306	43,056
Java	27,437	3,430	3,430	34,297
JavaScript	30,320	3,790	3,790	37,900
C#	27,876	3,485	3,485	34,846
HTML	11,430	1,429	1,429	14,288
PHP	11,785	1,473	1,474	14,732
Total	143,292	17,913	17,914	179,119

Table 2 Length statistics of code snippets, problem descriptions, and question titles in our corpus

Length Metrics	Values	Python	Java	C#	JavaScript	HTML	PHP
Code Length	Maximum	17,629	27,348	18,575	17,568	17,568	18,192
	Average	190.82	271.49	165.07	164.95	173.23	174.89
	Mode	29	17	27	40	20	22
	Median	92	96	86	86	86	79
	Minimum	2	2	2	2	2	2
	<256	81.90%	77.13%	83.99%	84.37%	83.97%	84.12%
Description Length	Maximum	1,776	2,158	1,851	1,579	1,360	1,139
	Average	87.99	94.81	101.14	86.38	83.12	86.37
	Mode	41	50	49	40	33	46
	Median	70	73	77	68	66	66
	Minimum	0	0	0	0	0	0
	<256	97.19%	95.93%	94.85%	97.87%	97.47%	96.57%
Title Length	Maximum	56	60	50	45	45	49
	Average	8.56	8.48	8.61	8.42	8.55	8.29
	Mode	7	8	8	7	7	7
	Median	8	8	8	8	8	8
	Minimum	2	2	2	2	2	2
	<16	96.63%	96.56%	95.70%	97.06%	96.65%	97.35%

¹⁰ <https://stackoverflow.com/questions/1478248>

3.3 Performance Measures

To show the competitiveness of SOTitle+, we consider four performance measures: Rouge-L (LIN 2004), BLEU (Papineni et al 2002), METEOR (Banerjee and Lavie 2005), and CIDEr (Vedantam et al 2015). Using these measures, we can measure the lexical similarity between the generated question title and the ground-truth question title. Moreover, these measures have been widely used in previous text generation tasks (Lin et al 2023; Liu et al 2019a; Li et al 2022; Raffel et al 2020; Gros et al 2020; Yang et al 2023a, 2022b; Li et al 2021).

- **ROUGE-L.** ROUGE (LIN 2004) is a recall-based measure, which is used to calculate the lexical overlap between the generated sentence and the ground-truth sentence. In our study, we use ROUGE-L, which is based on the LCS (Longest Common Subsequence).
- **BLEU.** BLEU (Papineni et al 2002) is employed to assess the similarity degree between the generated sentence and the ground-truth sentence. Specifically, BLEU- n calculates the number of n -grams that appear in the generated sentence.
- **METEOR.** METEOR (Banerjee and Lavie 2005) has been widely used to measure the quality of machine translation. It incorporates many aspects (e.g., vocabulary, grammar, and semantics) to evaluate the quality of the generated sentences.
- **CIDEr.** CIDEr (Vedantam et al 2015) is a measure that evaluates the quality of the generated sentences by analyzing the frequency of n -grams in the ground-truth sentences. To accomplish this goal, CIDEr employs TF-IDF weighting for each n -gram. Moreover, CIDEr calculates the $CIDEr_n$ score for each n -gram by averaging the cosine similarity between the generated sentence and the ground-truth sentence.

Notice the higher the value of these performance measures, the better the performance of the corresponding approach. In our study, we use the implementations provided by nlgeval¹¹ to calculate the value of these four performance measures.

3.4 Baselines

For the Stack Overflow question title generation, we select four state-of-the-art approaches as our baselines to show the competitiveness of SOTitle+. We briefly introduce the characteristics of these baselines as follows.

- **Code2Que.** Code2Que is the first study for Stack Overflow question title generation, which was proposed by (Gao et al 2020). Code2Que automatically generates question titles for Stack Overflow posts through an LSTM-based deep learning approach. The model’s encoder is a multi-layer, bidirectional LSTM network that processes the input code tokens sequentially. While the decoder is a single-layer LSTM that recursively generates predicted tokens.
- **SOTitle.** To show the effectiveness of SOTitle+, we also consider our previously proposed approach SOTitle (Liu et al 2022) as the baseline. SOTitle utilizes the pre-trained model T5 (Raffel et al 2020), which follows a Transformer encoder-decoder architecture and employs multi-task learning to unify text-based language problems into a text-to-text paradigm.

¹¹ <https://github.com/Maluuba/nlg-eval>

- **CCBERT.** CCBERT (Zhang et al 2022) is another approach, which also considers bi-modal information. CCBERT considers the CodeBERT model and introduces customized copy attention layers on the encoder and decoder.
- **M3NSCT5.** M3NSCT5 (Zhang et al 2023) is a recently proposed approach that can improve the quality and diversity of the question titles generated from the code snippets. Zhang et al. fine-tuned the codeT5 and utilized the nucleus sampling and maximal marginal ranking strategy to make the model return a set of high-quality and diverse title candidates.

To alleviate the internal threats, we utilize the scripts shared by these four baselines and follow the hyperparameter settings used in their original studies. Since the baselines Code2Que and M3NSCT5 only consider code snippets as their inputs, we still use bi-modal information as their inputs to guarantee a fair comparison.

3.5 Implementation Details

In our empirical study, we use Pytorch library¹² and transformers library¹³ to construct our model. For prompt tuning, we use OpenPrompt library¹⁴. Specifically, by following the tutorial for OpenPrompt (Ding et al 2022), we use the *manual_template*, *SoftTemplate*, *MixedTemplate* APIs to construct the hard prompt, the soft prompt, and the hybrid prompt templates, respectively. Based on the constructed prompt templates, we use *PromptForGeneration* API to perform prompt tuning. More details about our implementation can be found in our shared GitHub repository.

For the model hyperparameter setting, the word embedding dimension and hidden size are set to 768, and the number of attention heads and layers is set to 12. All parameters were optimized with AdamW (Loshchilov and Hutter 2018), and the initial learning rate is set to 0.00005. During training, the batch size is set to 16. The maximum length of the encoder and decoder is set to 512 and 64 respectively. To implement heuristic search, we adopt beam search and set the beam size to 10. Finally, we also use the early stop method (Prechelt 1998) to further alleviate the overfitting problem, and the weights with the highest performance on the validation set are taken as the final weights of the neural network.

3.6 Running Platform and Model Training Time

We run all the experiments on a computer with an Intel(R) Core(TM) i5-13600K and a GeForce RTX4090 GPU with 24 GB memory. The running operation system platform is Windows Operation System. Based on our platform, we need about 24 hours to train the model by using SOTitle+.

¹² <https://pytorch.org/>

¹³ <https://github.com/huggingface/transformers>

¹⁴ <https://github.com/thunlp/OpenPrompt>

4 Result Analysis

4.1 RQ1: Comparison with Baselines via Automatic Evaluation

The comparison results between SOTitle+ and baselines can be found in Table 3. In this table, we emphasize the best value in bold for each column. Based on the comparison results, we first find that our proposed approach SOTitle+ can outperform four baselines in terms of the performance measure Rouge-L for different programming languages. In particular, for our considered baselines, we find fine-tuning-based baselines (i.e., CCBERT (Zhang et al 2022), SOTitle (Liu et al 2022) and M3NSCT5 (Zhang et al 2023)) can achieve better performance than Code2Que (Gao et al 2020), which trains an LSTM model from scratch. The potential reason is that these three baselines are based on PLMs, which have been trained on the large-scale corpus. Moreover, CCBERT and SOTitle achieve considerable performance, which is better than M3NSCT5. We attribute the good performance of CCBERT and SOTitle to the reason that these two baselines were specifically designed to process bi-modal information. Finally, our approach can outperform the fine-tuning-based baselines. Specifically, compared to the baseline CCBERT, SOTitle+ can improve the performance of ROUGE-L by 18.04%, 14.65%, 9.41%, 13.43%, 20.19% and 7.98% for Python, Java, C#, JavaScript, PHP, and HTML respectively. Moreover, we also perform the Wilcoxon signed-rank test (Wilcoxon 1992) at the confidence level of 95% to check whether the performance differences between SOTitle+ and baselines are significant. All the p -values are smaller than 0.05, which means SOTitle+ can significantly outperform the baselines.

Except for Rouge-L performance measure, we also compare SOTitle+ with baselines in terms of METEOR (Banerjee and Lavie 2005), BLEU (Papineni et al 2002) and CIDEr (Vedantam et al 2015) performance measures. Final comparison results and the Wilcoxon signed-rank test results also show the competitive performance of SOTitle+ (i.e., the performance improvement is statistically significant).

Fig. 5 shows the question titles generated by SOTitle+ and baselines for a question post related to Python programming language¹⁵. In this example, we can find that the question titles generated by Code2Que and M3NSCT5 fail to capture the developer’s intent, and they mislocate key information to the “HTTP request”. The question titles generated by CCBERT and SOTitle lose the keywords “default” and “Asyncio”, respectively. However, the question title generated by SOTitle+ can accurately and fluently express the key information in this post.

Summary for RQ1: SOTitle+ can achieve better performance than four state-of-the-art baselines for different programming languages via automatic evaluation.

4.2 RQ2: Ablation Study on Using the Bi-modal Information

In this RQ, we aim to investigate the contribution of different modal information to the performance of SOTitle+. Here we use different subscripts to distinguish different control approaches and introduce the meaning of these subscripts as follows.

- **code.** The corresponding control approach only uses the code snippet along with its hard prompt as the input modality. We use *w/o desc* to denote this control approach.

¹⁵ <https://stackoverflow.com/questions/66287470>

Table 3 Comparison results between our proposed approach SOTitle+ and state-of-the-art baselines

Language	Approach	ROUGE-L (%)	METEOR (%)	BLEU-1 / 2 / 3 / 4 (%)	CIDEr
Python	Code2Que	21.75	14.87	23.80 / 15.54 / 10.92 / 8.10	0.75
	M3NSCT5	25.39	17.84	27.72 / 19.26 / 14.35 / 11.20	0.98
	CCBERT	26.50	15.52	28.07 / 19.09 / 13.94 / 10.63	1.03
	SOTitle	26.49	17.87	26.75 / 19.31 / 14.89 / 11.98	1.11
	SOTitle+	31.27	22.00	33.28 / 25.03 / 19.98 / 16.59	1.43
Java	Code2Que	19.67	13.77	23.80 / 15.54 / 10.92 / 7.14	0.66
	M3NSCT5	22.83	16.46	25.43 / 17.37 / 12.70 / 9.89	0.86
	CCBERT	24.17	14.33	26.11 / 17.68 / 12.90 / 9.93	0.91
	SOTitle	23.70	16.14	24.10 / 17.34 / 13.32 / 10.68	0.99
	SOTitle+	27.71	19.46	30.04 / 21.99 / 17.18 / 14.01	1.21
C#	Code2Que	20.00	14.29	21.19 / 13.71 / 9.43 / 6.94	0.64
	M3NSCT5	23.45	16.83	24.75 / 17.09 / 12.70 / 9.98	0.85
	CCBERT	25.28	15.32	27.84 / 19.87 / 14.86 / 11.40	0.94
	SOTitle	24.22	16.39	25.20 / 18.31 / 14.14 / 11.43	0.97
	SOTitle+	27.66	19.37	30.69 / 22.94 / 18.27 / 15.24	1.20
JavaScript	Code2Que	22.40	15.45	24.18 / 15.75 / 11.22 / 8.51	0.77
	M3NSCT5	25.97	18.11	27.95 / 19.58 / 14.92 / 11.97	1.02
	CCBERT	27.70	16.23	29.50 / 20.50 / 15.34 / 12.05	1.12
	SOTitle	27.19	18.13	27.56 / 20.18 / 15.81 / 12.93	1.15
	SOTitle+	31.42	21.90	33.69 / 25.29 / 20.27 / 16.91	1.41
PHP	Code2Que	22.50	15.63	24.17 / 15.60 / 11.22 / 8.51	0.82
	M3NSCT5	26.32	18.68	28.53 / 19.82 / 14.95 / 11.83	1.10
	CCBERT	26.60	16.26	28.39 / 19.94 / 15.21 / 12.06	1.09
	SOTitle	27.69	19.25	28.63 / 21.07 / 16.45 / 13.35	1.23
	SOTitle+	31.97	23.20	34.53 / 25.80 / 20.82 / 17.59	1.51
HTML	Code2Que	22.80	16.67	23.39 / 14.87 / 10.11 / 7.16	0.79
	M3NSCT5	25.06	18.41	25.96 / 16.81 / 11.65 / 8.39	0.90
	CCBERT	26.91	16.98	29.38 / 19.59 / 13.46 / 9.54	0.98
	SOTitle	25.84	18.54	28.36 / 18.83 / 13.37 / 9.76	0.98
	SOTitle+	29.00	21.03	30.74 / 20.82 / 15.11 / 11.51	1.14

- **desc.** The corresponding control approach only uses the problem description along with its hard prompt as the input modality. We use *w/o code* to denote this control approach.

We show the performance of SOTitle+ with different modal information in Table 4. When only focusing on code snippets (*w/o desc*), the performance of SOTitle+ causes a significant decrease. In particular, in terms of Rouge-L, the average performance of SOTitle+ drops by 51.85%. On the other hand, when only focusing on problem descriptions (*w/o code*), there is a slight decrease in the performance of SOTitle+. These findings show that problem descriptions can provide more assistance in improving the performance of SOTitle+ than code snippets. Moreover, using both problem descriptions and code snippets simultaneously can help to achieve the best performance due to the complementarity between these two different modal information.

Post body:	Question Title:
<p>Problem Description</p> <p>I'm using Python Asyncio to do a lot of HTTP requests. I'm wondering what is the default concurrency level of Asyncio or how many HTTP requests would be happening in parallel at any given time?</p> <p>The code I'm using to do the HTTP requests you can find below:</p> <p>Code Snippet</p> <pre> async def call_url(self, session, url): response = await session.request(method='GET', url=url) return response async def main(self, url_list): async with aiohttp.ClientSession() as session: res = await asyncio.gather(*[self.call_url(session, url) for url in url_list]) return res </pre>	<p>Ground Truth: What is the default concurrency level with asyncio in Python?</p> <p>Code2Que: Using Python Asyncio to do a lot of HTTP requests</p> <p>M3NSCT5: How many HTTP requests would be happening?</p> <p>CCBERT: What is the concurrency of Asyncioioio ?</p> <p>SOTitle: What is default concurrency level in Python?</p> <p>SOTitle+: What is the default concurrency level of Asyncio in Python?</p>

Fig. 5 The question titles generated by SOTitle+ and baselines for a question post related to the Python programming language

Summary for RQ2: Problem description contributes more to the performance of SOTitle+ when compared to code snippet. Moreover, Considering both of these two different modal information can achieve the best performance for SOTitle+.

Table 4 Comparison results between our proposed approach SOTitle+ and SOTitle+ with different modal information

Language	Approach	ROUGE-L (%)	METEOR (%)	BLEU-1 / 2 / 3 / 4 (%)	CIDEr
Python	<i>w/o desc</i>	17.73	12.88	19.08 / 12.58 / 9.36 / 7.52	0.60
	<i>w/o code</i>	27.17	18.64	29.08 / 20.46 / 15.35 / 11.99	1.06
	SOTitle+	31.27	22.00	33.28 / 25.03 / 19.98 / 16.59	1.43
Java	<i>w/o desc</i>	14.73	10.83	16.13 / 10.42 / 7.74 / 6.26	0.48
	<i>w/o code</i>	24.25	17.03	26.29 / 18.27 / 13.58 / 10.62	0.96
	SOTitle+	27.71	19.46	30.04 / 21.99 / 17.18 / 14.01	1.21
C#	<i>w/o desc</i>	12.06	8.44	13.35 / 7.69 / 5.25 / 4.10	0.38
	<i>w/o code</i>	26.25	18.23	28.94 / 20.89 / 16.21 / 13.23	1.07
	SOTitle+	27.66	19.37	30.69 / 22.94 / 18.27 / 15.24	1.20
JavaScript	<i>w/o desc</i>	14.33	10.29	15.48 / 9.77 / 7.10 / 5.66	0.42
	<i>w/o code</i>	28.42	19.53	30.41 / 21.78 / 16.82 / 13.64	1.17
	SOTitle+	31.42	21.90	33.69 / 25.29 / 20.27 / 16.91	1.41
PHP	<i>w/o desc</i>	15.67	12.13	17.08 / 10.74 / 8.02 / 6.55	0.50
	<i>w/o code</i>	28.51	20.21	30.93 / 21.88 / 16.87 / 13.71	1.21
	SOTitle+	31.97	23.20	34.53 / 25.80 / 20.82 / 17.59	1.51
HTML	<i>w/o desc</i>	11.84	8.95	13.12 / 6.67 / 3.89 / 2.59	0.26
	<i>w/o code</i>	27.64	20.15	29.66 / 19.51 / 13.69 / 10.00	1.02
	SOTitle+	29.00	21.03	30.74 / 20.82 / 15.11 / 11.51	1.14

4.3 RQ3: Ablation Study on the Prompt-tuning Paradigm

In this RQ, we want to investigate how much can **prompt tuning** contribute to the performance of SOTitle+. We design a control approach SOTitle+_{ft}, which is fine-tuned on CodeT5. For SOTitle+_{ft}, we remove all the prompts and use a special identifier (`<code>`) to distinguish problem descriptions and code snippets. The final input X is defined as follows:

$$X = [LANG] \oplus X_{desc} \oplus \langle code \rangle \oplus X_{code} \quad (10)$$

where X_{desc} denotes the problem description and x_{code} denotes the code snippet.

We show the comparison results between SOTitle+_{ft} and SOTitle+ in Table 5. In this table, We find that using the prompt-tuning paradigm can outperform using the fine-tuning paradigm for all programming languages. For example, for the Python programming language, using prompt tuning can improve the performance by 14.13%, 17.95%, 32.00%, and 24.47% in terms of ROUGE-L, METEOR, BLEU-4, and CIDEr respectively. The performance improvement is particularly obvious in the low-resource scenario (i.e., programming languages with fewer question posts). For example, for the PHP programming language, prompt tuning shows the highest performance improvement, with increase rate at 13.89%, 19.04%, 47.94%, and 27.97% in terms of Rouge-L, METEOR, BLEU-4, and CIDEr, respectively. Finally, we also show examples on our GitHub repository¹⁶ that can demonstrate the prompt-tuning contributions to bridge the gap between the fine-tuning paradigm and the prompt-tuning paradigm.

Summary for RQ3: Compared with the fine-tuning paradigm, the prompt-tuning paradigm can help SOTitle+ achieve better performance, especially for low-resource programming languages.

4.4 RQ4: Ablation Study on Hybrid Prompts

The prompt template design is a challenging and open problem. As (Liu et al 2023d) reported even a single word difference in the hard prompt can make a huge performance difference. In this RQ, we want to investigate how much can our designed hybrid prompts contribute to SOTitle+. To answer this RQ, we design two control approaches that use only hard prompts or soft prompts, respectively. We show the details of them as follows:

- **Only Hard Prompts.** This control approach only uses hard prompts as prompt templates. Specifically, we use Equation (1) as our hard prompt templates and use *w/ hard* to denote this control approach.
- **Only Soft Prompts.** This control approach only uses soft prompts as prompt templates. Specifically, we use Equation (2) as our soft prompt templates. We use *w/ soft* to denote this control approach.

We show the performance of SOTitle+ with different types of prompt templates in Table 6. When only using hard prompts (*w/ hard*), the performance of SOTitle+ has a slight

¹⁶ <https://github.com/shaoyoung/SOTitlePlus/blob/main/Appendix.md>

Table 5 Comparison results between SOTitle+ with the prompt-tuning paradigm and SOTitle+ with the fine-tuning paradigm

Language	Approach	ROUGE-L (%)	METEOR (%)	BLEU-1 / 2 / 3 / 4 (%)	CIDEr
Python	SOTitle+ _{ft}	28.53	20.04	28.88 / 21.08 / 16.52 / 13.66	1.24
	SOTitle+	31.27	22.00	33.28 / 25.03 / 19.98 / 16.59	1.43
Java	SOTitle+ _{ft}	26.36	18.57	27.52 / 19.57 / 14.86 / 11.94	1.09
	SOTitle+	27.71	19.46	30.04 / 21.99 / 17.18 / 14.01	1.21
C#	SOTitle+ _{ft}	26.37	18.37	28.09 / 20.05 / 15.30 / 12.38	1.05
	SOTitle+	27.66	19.37	30.69 / 22.94 / 18.27 / 15.24	1.20
JavaScript	SOTitle+ _{ft}	29.01	19.87	29.70 / 21.62 / 16.95 / 14.00	1.23
	SOTitle+	31.42	21.90	33.69 / 25.29 / 20.27 / 16.91	1.41
PHP	SOTitle+ _{ft}	29.38	21.08	30.37 / 21.79 / 17.06 / 14.10	1.29
	SOTitle+	31.97	23.20	34.53 / 25.80 / 20.82 / 17.59	1.51
HTML	SOTitle+ _{ft}	27.55	19.68	29.91 / 19.67 / 13.72 / 10.16	1.02
	SOTitle+	29.00	21.03	30.74 / 20.82 / 15.11 / 11.51	1.14

performance decrease. The potential reason is that considering soft prompts in the hybrid prompt may extend the ability of prompt templates. It is similar to only using soft prompts (*w/ soft*) for SOTitle+, which causes a slight performance decrease. The potential reason is that without original hard prompts (i.e., “The problem description:” and “The code snippet:”), which are the task-related important tokens, our proposed approach may lose the ability to effectively separate the problem description and the code snippet from the bi-modal information.

Summary for RQ4: Compared to only using hard prompts or only using soft prompts, Our results show that using hybrid prompts can achieve the best performance for SOTitle+.

4.5 RQ5: Comparison with baselines via human study

Since CCBERT (Zhang et al 2022) and SOTitle (Liu et al 2022) can achieve the best performance in our considered baselines (discussed in Section 4.1), we conduct a human study to assess the quality of the question titles generated by SOTitle+, CCBERT, and SOTitle respectively. To evaluate the quality of the generated question titles, we followed a human study methodology used by the previous studies (Wei et al 2020; Liu et al 2019a) for similar tasks. In our human study, we consider three quality criteria (i.e., similarity, naturalness, and informativeness) and the score values range from 1 to 4, where a higher score indicates better quality of the generated question titles. We show the meaning of these three criteria as follows:

- **Similarity.** This criterion assesses the similarity between the generated question title and the ground-truth title.
- **Naturalness.** This criterion evaluates the grammaticality and fluency of the generated question title.

Table 6 Comparison results between our proposed approach SOTitle+ and the control approaches with different types of prompt templates

Language	Approach	ROUGE-L (%)	METEOR (%)	BLEU-1 / 2 / 3 / 4 (%)	CIDEr
Python	<i>w/ hard</i>	30.51	21.19	32.87 / 24.23 / 19.03 / 15.99	1.34
	<i>w/ soft</i>	29.12	20.49	30.99 / 22.55 / 17.52 / 14.19	1.25
	SOTitle+	31.27	22.00	33.28 / 25.03 / 19.98 / 16.59	1.43
Java	<i>w/ hard</i>	27.04	18.89	29.39 / 21.19 / 16.21 / 12.99	1.14
	<i>w/ soft</i>	26.88	18.78	28.76 / 20.44 / 15.57 / 12.42	1.12
	SOTitle+	27.71	19.46	30.04 / 21.99 / 17.18 / 14.01	1.21
C#	<i>w/ hard</i>	27.25	19.17	30.17 / 22.29 / 17.56 / 14.45	1.15
	<i>w/ soft</i>	26.80	18.76	29.43 / 21.50 / 16.87 / 13.87	1.11
	SOTitle+	27.66	19.37	30.69 / 22.94 / 18.27 / 15.24	1.20
JavaScript	<i>w/ hard</i>	30.76	21.43	33.18 / 24.56 / 19.42 / 16.06	1.34
	<i>w/ soft</i>	29.53	20.53	31.14 / 22.72 / 17.80 / 14.68	1.27
	SOTitle+	31.42	21.90	33.69 / 25.29 / 20.27 / 16.91	1.41
PHP	<i>w/ hard</i>	31.14	22.76	33.88 / 25.02 / 19.79 / 16.39	1.43
	<i>w/ soft</i>	30.17	22.20	32.95 / 23.93 / 18.75 / 15.39	1.37
	SOTitle+	31.97	23.20	34.53 / 25.80 / 20.82 / 17.59	1.51
HTML	<i>w/ hard</i>	28.05	20.49	29.30 / 19.70 / 14.11 / 10.56	1.08
	<i>w/ soft</i>	28.30	20.82	30.45 / 20.15 / 14.35 / 10.77	1.08
	SOTitle+	29.00	21.03	30.74 / 20.82 / 15.11 / 11.51	1.14

- **Informativeness.** This criterion quantifies the amount of content conveyed by the generated question title. Notice this criterion is irrespective of its naturalness.

We recruit six developers, including one Ph.D. student, and five Master students who are familiar with the usage of Stack Overflow and have at least three years of software development experience. They are not co-authors of this paper, which can avoid the bias of the human study results.

We randomly selected 100 posts for each programming language in the testing set, resulting in a total of 600 posts. For each post, we collected three question titles generated by SOTitle+, SOTitle, and CCBERT, respectively. We split posts for six programming languages into three groups (i.e., Python and Java, C#, and PHP, HTML, and JavaScript). Each group of posts is assigned to two different developers. Later, each developer manually evaluated 200 posts by considering the aforementioned three criteria. Notice the participants do not know which question titles are generated by SOTitle+. Since the evaluation process is labor-intensive, we limit each developer to only evaluating 100 posts within a half day. Finally, the evaluation scores provided by the two developers for each post are averaged. In our human study, to improve the score quality, the recruited developers are allowed to utilize the Internet to explore unfamiliar concepts related to the posts.

The evaluation results are shown in Fig. 6. Regarding the similarity criterion, SOTitle+ outperforms SOTitle and CCBERT in terms of generating question titles with higher similarity scores, with an average value of 2.95. This indicates that the question titles generated by SOTitle+ are of good quality and can be used with minimal modifications. As for the naturalness criterion, all of these approaches achieve similar performance as expected. In terms of the informativeness criterion, SOTitle+ demonstrates the ability to generate more comprehensive question titles compared to SOTitle and CCBERT. To measure the consistency of

the scoring results among different developers, we use Cohen’s Kappa (Cohen 1960). The final Kappa value was 0.83, which indicates a substantial agreement among the developers. In addition to these scores, we also provided examples, of which SOTitle+ outperforms or does not outperform these two baselines, in our GitHub repository.

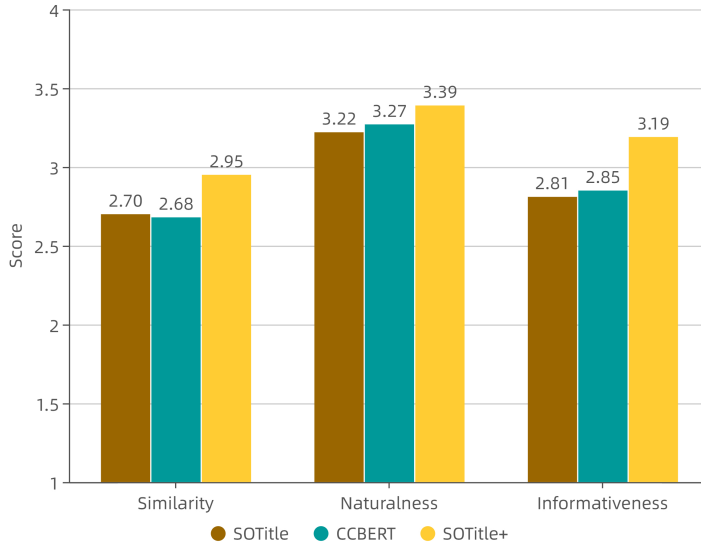


Fig. 6 The average score value of our human study by evaluating similarity, naturalness, and informativeness of the generated question titles

Summary for RQ5: Our human study shows that SOTitle+ outperforms the baselines SOTitle and CCBERT by evaluating similarity, naturalness, and informativeness of the generated question titles.

5 Discussions

5.1 Effectiveness of Using Multi-task Learning

In this subsection, to investigate the advantage of using multi-task learning in SOTitle+, we design *w/o multitask* as the control approach. In this control approach, we train separate models for each programming language. The comparison results can be found in Table 7. In this table, we find that by employing multi-task learning, SOTitle+ can outperform the models trained individually for each programming language. For example, SOTitle+ achieves a performance improvement of 17.06% at most in terms of BLEU₄ when compared to *w/o multitask* for six programming languages. Notice that using multi-task learning can achieve higher performance improvement for low-resource programming languages (such as HTML and PHP). This advantage can be attributed to the limited training data available for these low-resource languages, which can hinder the control approach *w/o multitask* from

effectively learning helpful knowledge. However, by utilizing multi-task learning, SOTitle+ can effectively acquire related knowledge from programming languages with sufficient training data.

Table 7 The ablation results on using multi-task learning for SOTitle+

Language	Approach	ROUGE-L (%)	METEOR (%)	BLEU-1 / 2 / 3 / 4 (%)	CIDEr
Python	<i>w/o multitask</i>	29.95	20.76	31.99 / 23.63 / 18.51 / 15.16	1.32
	SOTitle+	31.27	22.00	33.28 / 25.03 / 19.98 / 16.59	1.43
Java	<i>w/o multitask</i>	25.96	18.41	28.03 / 20.28 / 15.55 / 12.45	1.09
	SOTitle+	27.71	19.46	30.04 / 21.99 / 17.18 / 14.01	1.21
C#	<i>w/o multitask</i>	25.76	18.13	28.04 / 20.48 / 16.17 / 13.41	1.07
	SOTitle+	27.66	19.37	30.69 / 22.94 / 18.27 / 15.24	1.20
JavaScript	<i>w/o multitask</i>	29.25	20.73	31.63 / 23.19 / 18.25 / 15.00	1.27
	SOTitle+	31.42	21.90	33.69 / 25.29 / 20.27 / 16.91	1.41
PHP	<i>w/o multitask</i>	28.98	21.33	31.95 / 23.35 / 18.35 / 15.03	1.30
	SOTitle+	31.97	23.20	34.53 / 25.80 / 20.82 / 17.59	1.51
HTML	<i>w/o multitask</i>	25.53	19.11	26.71 / 17.05 / 11.71 / 8.40	0.90
	SOTitle+	29.00	21.03	30.74 / 20.82 / 15.11 / 11.51	1.14

5.2 Question Title Generation in Low-resource Scenario

In Stack Overflow, there are some programming languages whose training data is limited and we call this scenario the low-resource scenario. In this subsection, we want to investigate the performance of SOTitle+ in this low-resource scenario. Specifically, we select other two low-resource programming languages (i.e., Ruby and Go). We use the same rules (introduced in Section 2) for corresponding corpus construction and corpus split. The statistical information of these two corpora is shown in Table 8.

Table 8 Statistical information of two corpora for low-resource programming languages (i.e., Ruby and Go)

Language	Training Set	Validation Set	Testing Set	Total
Ruby	4,876	610	610	6,096
Go	2,017	252	253	2,522
Total	6,893	862	863	8,618

In this experiment, we only consider two baselines (i.e., CCBERT (Zhang et al 2022) and SOTitle (Liu et al 2022)) due to their higher performance (discussed in Section 4.1). For SOTitle+, we directly use the model trained for answering RQ1 since we want to investigate the generalization of SOTitle+ on new programming languages. For the two baselines, we train SOTitle and CCBERT on corpus related to Ruby and Go. We show the comparison results in Table 9. For the Ruby programming language, SOTitle+ can achieve 29.98% and 16.69% in terms of Rouge-L and BLEU-4, respectively. Similarly, for the Go programming

language, SOTitle+ can achieve 31.22% and 15.02% in terms of Rouge-L and BLEU-4, respectively. Taking the Go programming language as an example, SOTitle+ demonstrates the performance improvements of 13.98%, 16.97%, and 31.87% in terms of the ROUGE-L, METEOR, and BLEU-4 measures, respectively when compared to SOTitle, which can achieve the best performance among these two baselines. These results verify the effectiveness of SOTitle+ for posts related to low-resource programming languages.

Table 9 Comparison results between our proposed approach SOTitle+ and two baselines for two low-resource programming languages (i.e., Ruby and Go)

Language	Approach	ROUGE-L (%)	METEOR (%)	BLEU-1 / 2 / 3 / 4 (%)	CIDEr
Ruby	SOTitle	26.71	18.28	27.39 / 19.95 / 15.56 / 12.72	1.12
	CCBERT	26.85	14.30	25.79 / 17.65 / 13.14 / 10.10	0.93
	SOTitle+	29.98	21.94	32.04 / 24.44 / 19.83 / 16.69	1.37
Go	SOTitle	27.39	18.44	29.34 / 20.74 / 15.00 / 11.39	1.11
	CCBERT	27.47	15.50	28.50 / 19.99 / 15.18 / 11.97	0.99
	SOTitle+	31.22	21.57	32.80 / 23.79 / 18.45 / 15.02	1.46

5.3 Can SOTitle+ outperform ChatGPT?

Recently, LLMs (Brown et al 2020) (e.g., ChatGPT) have attracted great attention and achieved promising performance for many software engineering tasks, such as program repair (Xia and Zhang 2023) and code generation (Liu et al 2023a). In this subsection, we want to investigate the effectiveness of using ChatGPT in the SOQTG task and whether SOTitle+ can outperform ChatGPT on the posts related to Go and Ruby. Specifically, for ChatGPT, we use OpenAI official API¹⁷ and the hard prompt introduced in Section 2.2.1 since ChatGPT does not support hybrid prompt. For SOTitle+, we directly use the model trained for answering RQ1. Notice both SOTitle+ and ChatGPT are not trained on the corpus related to Go and Ruby. Therefore, they are compared in the zero-shot scenario.

Table 10 Comparison results between our proposed approach SOTitle+ and ChatGPT

Language	Approach	ROUGE-L (%)	METEOR (%)	BLEU-1 / 2 / 3 / 4 (%)	CIDEr
Ruby	ChatGPT	20.59	19.37	14.93 / 8.72 / 5.71 / 3.96	0.58
	SOTitle+	29.98	21.94	32.04 / 24.44 / 19.83 / 16.69	1.37
Go	ChatGPT	19.95	19.68	15.97 / 8.97 / 5.18 / 3.23	0.46
	SOTitle+	31.22	21.57	32.80 / 23.79 / 18.45 / 15.02	1.46

Table 10 shows the performance comparison results of ChatGPT and SOTitle+. In this table, we can find SOTitle+ can achieve better performance than ChatGPT in terms of all the

¹⁷ <https://platform.openai.com/docs/api-reference>. The version of ChatGPT we use is GPT-3.5-turbo.

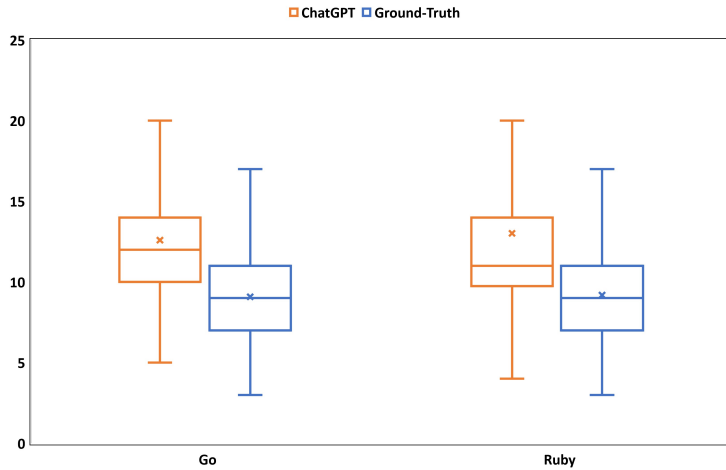


Fig. 7 Length distribution of question titles generated by ChatGPT and ground-truth question titles

performance measures. After our manual analysis, we found that ChatGPT tends to generate semantically rich but lengthy question titles. We use Figure 7 to show the length of question titles generated by ChatGPT and ground-truth question titles. Therefore, a potential solution is to design customized prompts, which can crop the rich but lengthy generated titles, in the future.

5.4 User Study for Improving Low-quality Question Titles

By following the study of Code2Que (Gao et al 2021), we conduct a user study to investigate whether SOTitle+ can generate better question titles for low-quality questions. Specifically, we sampled 20 low-quality questions for each programming language. Notice the votes for these questions are less than 10 and these questions are not contained in our corpus. For each question, we gather the title written by users (i.e., ground-truth title) and the title generated by SOTitle+. Then we recruited the same six developers, who were also recruited in our human study, to determine which title was better in terms of three criteria (i.e., Clearness, Fitness, and Willingness). The detailed information of these six developers can be found in Section 4.5. Specifically, Clearness evaluates the clarity of a question title. Fitness measures whether a question title is reasonable in logic with the provided code snippet and the problem description. Willingness measures whether a user is willing to respond to this question.

We show the comparison results (i.e., Win, Loss, and Tie) for two programming languages in Table 11 and other programming languages in our GitHub repository. In this table, we find that titles generated by SOTitle+ can improve the title quality in terms of all three criteria in most cases. Moreover, we find that the titles generated by SOTitle+ can effectively increase the willingness of users to answer questions. Finally, we find that not all the titles of the low-quality questions can be improved. One of the reasons is that the post itself did not provide enough information. In summary, our preliminary user study can show the potential of SOTitle+ in title quality improvement.

Table 11 User Study (SOTitle+ vs. Human)

Language	Evaluation Criterion	Win (%)	Lose (%)	Tie (%)
Python	Clearness	50.83	30.83	18.33
	Fitness	52.50	25.83	21.67
	Willingness	56.67	26.67	16.67
Java	Clearness	44.17	39.17	16.67
	Fitness	47.50	49.00	12.50
	Willingness	39.17	35.83	25.00

5.5 Limitations of SOTitle+ by Analyzing Failed Cases

Although SOTitle+ achieves promising performance via both automatic evaluation and human study, it should be noted that for some cases, SOTitle+ may generate question titles that show low similarity to the ground-truth question titles. To gain deeper insights into this issue, we randomly selected 100 posts and performed a manual analysis. This analysis aims to identify the specific types of posts that pose challenges for SOTitle+ and help to point out potential future studies.

The first limitation arises from the possibility of missing certain question words (e.g., *how to*, *what*, *why*) in the generated question titles. For example, Given this question post¹⁸, the ground-truth title is “*How to use OpenCV ConnectedComponents to get the images*” while the title generated by SOTitle+ is “*OpenCV ConnectedComponents in Python*”. It appears that the model fails to capture the intention of this question post. One potential solution is to automatically identify the type of posts and provide this information for title generation.

The second challenge type is that the generated title and the ground truth convey the same meaning but are expressed using different ways. For example, given this question post¹⁹, the ground-truth title is “*Remove seaborn lineplot legend title*”, the title generated by SOTitle+ is “*Remove title from seaborn lineplot legend*”. The sentence-level ROUGE-L score between these two titles is 26.56%. However, these two titles are semantically equivalent. This type of challenge can lead to lower scores in automatic evaluation metrics. One possible solution is to develop semantic-based evaluation metrics that can more accurately assess the semantic similarity between the generated title and the ground-truth title.

The third challenge type is that the problem description may contain noise and then mislead the model to generate the wrong question title. Given this question post²⁰, the ground-truth title is “*How to run UVICORN in Heroku?*”, while the title generated by SOTitle+ is “*How to run a fastapi script in Heroku?*”. After our manual analysis, we find that the keyword “*fastapi*” appears four times in the problem description and the keyword “*UVICORN*” appears only one time. However, in this post, the developer wants to ask how to use a command (i.e., “*UVICORN*”) but not a script (i.e., “*fastapi*”). One possible solution is to design an effective approach to identify the developer’s intent and use this information to guide title generation.

¹⁸ <https://stackoverflow.com/questions/51523765>

¹⁹ <https://stackoverflow.com/questions/51579215>

²⁰ <https://stackoverflow.com/questions/59391560>

5.6 Threats to validity

In this section, we mainly analyze the potential threats to the validity of our empirical study.

5.6.1 Internal Threats

The first threat to internal validity is related to potential faults in the implementation of our proposed approach SOTitle+ and baselines. To alleviate this threat, we use mature frameworks (such as transformers, and OpenPrompt) as much as possible to implement SOTitle+. We also use software testing and code inspection to guarantee the code quality of SOTitle+. For baselines, we use their shared scripts. We replicated these scripts on the same dataset shared by these studies to confirm that they can achieve comparable results reported in the original studies. After replication, we train the models by these baselines on our gathered corpus. The second threat to internal validity is related to prompt template construction for the SOQTG task. However, designing an optimal prompt template is still a challenging task due to the large-scale search space. To alleviate this threat, we considered three different types of prompt templates. Our empirical results show that using the hybrid prompt can achieve the best performance for SOTitle+ and outperform state-of-the-art baselines. Thus, the performance reported in our study can serve as a lower bound of SOTitle+, which can be even further improved by designing more high-quality prompt templates in the future. The third threat to internal validity is related to the data leakage problem, which is caused by the usage of pre-trained language models. Our study used CodeT5 (Wang et al 2021) as the pre-trained language model. Specifically, Wang et al. employed CodeSearchNet (Husain et al 2019) to pre-train CodeT5, which consists of six program languages with both unimodal data and bimodal data. Moreover, they additionally considered C/CSharp from BigQuery²¹. It is not hard to find that the dataset in their pre-training task mainly comes from GitHub, while our dataset comes from Stack Overflow. Therefore, using CodeT5 as the pre-trained model can avoid the problem of dataset leakage. However, if the pre-trained model (such as PLBART (Ahmad et al 2021)) used by SOTitle+ considered Stack Overflow data, this threat exists. A practical solution is to remove the data considered in the pre-trained model from the testing set.

5.6.2 External Threats

The threat to external validity concerns the quality and generalization ability of our constructed corpus. To alleviate this threat, we first consider question posts for six popular programming languages based on their popularity. Then we consider three heuristic selection rules, which aim to identify and remove low-quality posts. However, even if these rules are used, there may still exist some low-quality posts in the selected posts and more effective rules need to be designed. Finally, to alleviate the data leakage problem, we consider the temporal relationships of posts when we split the corpus, which can guarantee that the test set contains the latest posts. To facilitate the experiment’s reproducibility, We provide our corpus process scripts on our GitHub repository.

5.6.3 Construct Threats

The threat to construct validity comes from human studies and performance evaluation measures. For our human study, it may introduce evaluation bias. To alleviate this threat and en-

²¹ <https://console.cloud.google.com/marketplace/details/github/github-repos>

sure a fair comparison, all the developers did not know which question title is generated by SOTitle+ and they are not co-authors of this paper. Moreover, before their engagement in the study, we ensured that all the developers received a tutorial to familiarize them with the tasks and the evaluation procedures. This training aimed to minimize potential biases and inconsistencies in their assessments. For automatic evaluation, we consider four commonly used evaluation measures, which have been widely used in text-generation tasks and question title-generation studies (Liu et al 2022; Hu et al 2020; Yang et al 2022a; El-Kassas et al 2021). To avoid implementation errors, we use the mature framework nlg-eval to calculate these performance measure values.

6 Related Work

In this section, we first summarize the related work for Stack Overflow question title generation. Meanwhile, since we use the prompt learning paradigm in SOTitle+, we also analyze the related work for using this paradigm in software engineering tasks. Finally, we emphasize the novelty of our study.

6.1 Stack Overflow Question Title Generation

Ensuring the quality of question posts is essential for maintaining the usefulness of Q&A websites (such as Stack Overflow) (Yang et al 2014; Ponzanelli et al 2014; Duijn et al 2015; Arora et al 2015). Recently, researchers have focused on question title generation for Stack Overflow, as high-quality titles play a significant role in attracting potential experts to provide immediate answers. (Gao et al 2020) were the first to study this problem and proposed Code2Que, which utilizes the LSTM network with the copy and coverage mechanisms to generate question titles. Similar to our previous study (Liu et al 2022), (Zhang et al 2022) found that considering both problem description and code snippet could significantly improve the quality of generated titles. Notice this study was conducted simultaneously with our previous study. But compared to our previous study, their study adopts a different architecture, which considers the CodeBERT model and introduces customized copy attention layers on the encoder and decoder. In addition, they simply merged the problem description and code snippet without distinguishing them with special identifiers explicitly in bi-modal construction. Later, (Zhang et al 2023) further proposed the approach M3NSCT5, which utilizes the nucleus sampling and maximal marginal ranking strategy. In M3NSCT5, they fine-tuned CodeT5 and showed the promising performance of their proposed approach. Recently, (Liu et al 2023b) further proposed an automatic question title reformulation approach QETRA by mining modification logs from Stack Overflow. Specifically, QETRA can further polish the draft titles generated by previous question title generation approaches (Gao et al 2020; Liu et al 2022; Zhang et al 2022).

6.2 Prompt Learning and its Application to Software Engineering Tasks

Prompt-tuning is a rapidly evolving technology, which directly adapts the pre-trained language model to *cloze*-style prediction or sequence-to-sequence generation through some additional prompts. As a new paradigm, prompt learning has achieved promising performances on various software engineering tasks. For example, (Wang et al 2022) conducted

an empirical study on three code intelligence tasks. They compared prompt-tuning with fine-tuning and reported that prompt-tuning has better performance on both classification tasks and generation tasks. (Zhu et al 2023) proposed a context-aware prompt-tuning approach and applied this approach to two method naming tasks (i.e., method name recommendation and method name consistency checking). Recently, (Li et al 2023) considered various types of prompt templates and combined prompt ensembling and transfer learning to improve the performance of prediction of vulnerability characteristics.

6.3 Novelty of Our Study

As the example shown in Fig. 1, two posts have the same code snippet, but if the problem description is different, it will result in different titles. Compared to the previous study (Gao et al 2020), our study aims to generate question titles by considering both the code snippet and the problem description. Our ablation study confirms the effectiveness of considering problem description and finds the problem description contributes more to the performance of SOTitle+. Moreover, as discussed in Fig. 2, the inconsistent inputs and optimization objectives between the pre-training task and our investigated task may make fine-tuning hard to fully explore the knowledge of the pre-trained model. Therefore, SOTitle+ further prompt-tune CodeT5 with hybrid prompts (i.e., a mixture of hard and soft prompts) when compared to our preliminary study SOTitle (Liu et al 2022). The effectiveness of using prompt tuning and hybrid prompts is confirmed in our ablation studies. Finally, to alleviate the insufficient training data issue for low-resource programming languages and fully utilize the knowledge in the related tasks, we formalize question title generation for different programming languages as separate but related tasks. Then we utilize multi-task learning to solve these tasks. The effectiveness of using multi-task learning is also confirmed in our ablation study and the promising performance is achieved when applying SOTitle+ to other two low-resource programming languages (i.e., Ruby and Go).

7 Conclusion

In this study, we propose the prompt learning-based approach SOTitle+ for automatic question title generation by leveraging bi-modal information (i.e., the code snippets and the problem description) in the question posts and hybrid prompts. To verify the effectiveness of our proposed approach, we gathered 179,119 high-quality problem posts for six popular programming languages from Stack Overflow as our experimental subject. Experimental results show the competitiveness of our proposed approach when compared with the four state-of-the-art baselines in terms of four performance measures (i.e., Rouge, METEOR, BLEU, and CIDEr). Moreover, we also conduct a human study to verify the effectiveness of SOTitle+ by considering the similarity, naturalness, and informativeness of the generated question titles.

In the future, we first want to further improve the performance of SOTitle+ by designing more useful prompts via prompt engineer. We second want to research how to fully exploit ChatGPT on the SOQTG task. The challenging task is how to crop the rich but lengthy titles that are generated by ChatGPT. Finally, We want to investigate whether vector representation of Stack Overflow posts (Xu et al 2021) can be used for SOTitle+ to generate question titles with higher quality.

Acknowledgements The authors would like to thank the editors and the anonymous reviewers for their insightful comments and suggestions, which can substantially improve the quality of this work. Shaoyu Yang and Xiang Chen have contributed equally to this work and they are co-first authors. Xiang Chen is the corresponding author. This work is supported in part by the National Natural Science Foundation of China (Grant no. 61202006 and 61702041) and the Innovation Training Program for College Students (Grant no. 2023214 and 2023356).

Conflict of Interest

The authors declare that they have no conflict of interest.

Data Availability Statements

The datasets generated during and analyzed during the current study are available in the Github repository (<https://github.com/shaoyuyoung/SOTitlePlus>).

References

- Ahmad W, Chakraborty S, Ray B, Chang KW (2020) A transformer-based approach for source code summarization. In: Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pp 4998–5007
- Ahmad W, Chakraborty S, Ray B, Chang KW (2021) Unified pre-training for program understanding and generation. In: Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp 2655–2668
- Anderson A, Huttenlocher D, Kleinberg J, Leskovec J (2012) Discovering value from community activity on focused question answering sites: a case study of stack overflow. In: Proceedings of the 18th ACM SIGKDD international conference on Knowledge discovery and data mining, pp 850–858
- Arora P, Ganguly D, Jones GJ (2015) The good, the bad and their kins: Identifying questions with negative scores in stackoverflow. In: 2015 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM), IEEE, pp 1232–1239
- Banerjee S, Lavie A (2005) Meteor: An automatic metric for mt evaluation with improved correlation with human judgments. In: Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization, pp 65–72
- Brown TB, Mann B, Ryder N, Subbiah M, Kaplan J, Dhariwal P, Neelakantan A, Shyam P, Sastry G, Askell A, et al (2020) Language models are few-shot learners. In: Proceedings of the 34th International Conference on Neural Information Processing Systems, pp 1877–1901
- Cao K, Chen C, Baltes S, Treude C, Chen X (2021) Automated query reformulation for efficient search based on query logs from stack overflow. In: 2021 IEEE/ACM 43rd International Conference on Software Engineering (ICSE), IEEE, pp 1273–1285
- Cheng J, Dong L, Lapata M (2016) Long short-term memory-networks for machine reading. In: 2016 Conference on Empirical Methods in Natural Language Processing, Association for Computational Linguistics, pp 551–561
- Cohen J (1960) A coefficient of agreement for nominal scales. *Educational and psychological measurement* 20(1):37–46
- Correa D, Sureka A (2013) Fit or unfit: analysis and prediction of ‘closed questions’ on stack overflow. In: Proceedings of the first ACM conference on Online social networks, pp 201–212
- Ding N, Hu S, Zhao W, Chen Y, Liu Z, Zheng H, Sun M (2022) Openprompt: An open-source framework for prompt-learning. In: Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics: System Demonstrations, pp 105–113
- Duijn M, Kucera A, Bacchelli A (2015) Quality questions need quality code: Classifying code fragments on stack overflow. In: 2015 IEEE/ACM 12th Working Conference on Mining Software Repositories, IEEE, pp 410–413
- El-Kassas WS, Salama CR, Rafea AA, Mohamed HK (2021) Automatic text summarization: A comprehensive survey. *Expert systems with applications* 165:113679

- Gao Z, Xia X, Grundy J, Lo D, Li YF (2020) Generating question titles for stack overflow from mined code snippets. *ACM Transactions on Software Engineering and Methodology (TOSEM)* 29(4):1–37
- Gao Z, Xia X, Lo D, Grundy J, Li YF (2021) Code2que: A tool for improving question titles from mined code snippets in stack overflow. In: *Proceedings of the 29th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp 1525–1529
- Gao Z, Xia X, Lo D, Grundy J, Zhang X, Xing Z (2023) I know what you are searching for: Code snippet recommendation from stack overflow posts. *ACM Transactions on Software Engineering and Methodology* 32(3):1–42
- Gros D, Sezhiyan H, Devanpu B, Yu Z (2020) Code to comment “translation”: Data, metrics, baselining & evaluation. In: *2020 35th IEEE/ACM International Conference on Automated Software Engineering (ASE)*, IEEE, pp 746–757
- Hu X, Li G, Xia X, Lo D, Jin Z (2020) Deep code comment generation with hybrid lexical and syntactical information. *Empirical Software Engineering* 25(3):2179–2217
- Huang Q, Yuan Z, Xing Z, Xu X, Zhu L, Lu Q (2022) Prompt-tuned code language model as a neural knowledge base for type inference in statically-typed partial code. In: *37th IEEE/ACM International Conference on Automated Software Engineering*, pp 1–13
- Husain H, Wu HH, Gazit T, Allamanis M, Brockschmidt M (2019) Codesearchnet challenge: Evaluating the state of semantic code search. *arXiv preprint arXiv:190909436*
- Islam MJ, Nguyen G, Pan R, Rajan H (2019) A comprehensive study on deep learning bug characteristics. In: *Proceedings of the 2019 27th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, pp 510–520
- Iyer S, Konstas I, Cheung A, Zettlemoyer L (2016) Summarizing source code using a neural attention model. In: *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp 2073–2083
- Jin X, Servant F (2019) What edits are done on the highly answered questions in stack overflow? an empirical study. In: *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*, IEEE, pp 225–229
- Lester B, Al-Rfou R, Constant N (2021) The power of scale for parameter-efficient prompt tuning. In: *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pp 3045–3059
- Li X, Ren X, Xue Y, Xing Z, Sun J (2023) Prediction of vulnerability characteristics based on vulnerability description and prompt learning. In: *2023 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*, IEEE, pp 604–615
- Li Z, Wu Y, Peng B, Chen X, Sun Z, Liu Y, Yu D (2021) Secnn: A semantic cnn parser for code comment generation. *Journal of Systems and Software* 181:111036
- Li Z, Wu Y, Peng B, Chen X, Sun Z, Liu Y, Paul D (2022) Setransformer: A transformer-based code semantic parser for code comment generation. *IEEE Transactions on Reliability* 72(1):258–273
- LIN C (2004) Rouge: A package for automatic evaluation of summaries. In: *Proc. Workshop on Text Summarization Branches Out, Post-Conference Workshop of ACL 2004*
- Lin H, Chen X, Chen X, Cui Z, Miao Y, Zhou S, Wang J, Su Z (2023) Gen-fl: Quality prediction-based filter for automated issue title generation. *Journal of Systems and Software* 195:111513
- Liu C, Bao X, Zhang H, Zhang N, Hu H, Zhang X, Yan M (2023a) Improving chatgpt prompt for code generation. *arXiv preprint arXiv:230508360*
- Liu K, Yang G, Chen X, Yu C (2022) Sotitle: A transformer-based post title generation approach for stack overflow. In: *2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER)*, IEEE, pp 577–588
- Liu K, Chen X, Chen C, Xie X, Cui Z (2023b) Automated question title reformulation by mining modification logs from stack overflow. *IEEE Transactions on Software Engineering* 49(9):4390–4410
- Liu P, Yuan W, Fu J, Jiang Z, Hayashi H, Neubig G (2023c) Pre-train, prompt, and predict: A systematic survey of prompting methods in natural language processing. *ACM Computing Surveys* 55(9):1–35
- Liu Q, Liu Z, Zhu H, Fan H, Du B, Qian Y (2019a) Generating commit messages from diffs using pointer-generator network. In: *2019 IEEE/ACM 16th International Conference on Mining Software Repositories (MSR)*, IEEE, pp 299–309
- Liu X, He P, Chen W, Gao J (2019b) Multi-task deep neural networks for natural language understanding. In: *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp 4487–4496
- Liu X, Zheng Y, Du Z, Ding M, Qian Y, Yang Z, Tang J (2023d) Gpt understands, too. *AI Open*
- Loshchilov I, Hutter F (2018) Decoupled weight decay regularization. In: *International Conference on Learning Representations*
- Niu C, Li C, Ng V, Chen D, Ge J, Luo B (2023) An empirical comparison of pre-trained models of source code. *arXiv preprint arXiv:230204026*

- Papineni K, Roukos S, Ward T, Zhu WJ (2002) Bleu: a method for automatic evaluation of machine translation. In: Proceedings of the 40th annual meeting of the Association for Computational Linguistics, pp 311–318
- Ponzanelli L, Mocci A, Bacchelli A, Lanza M, Fullerton D (2014) Improving low quality stack overflow post detection. In: 2014 IEEE international conference on software maintenance and evolution, IEEE, pp 541–544
- Prechelt L (1998) Early stopping-but when? In: Neural Networks: Tricks of the trade, Springer, pp 55–69
- Raffel C, Shazeer N, Roberts A, Lee K, Narang S, Matena M, Zhou Y, Li W, Liu PJ (2020) Exploring the limits of transfer learning with a unified text-to-text transformer. *Journal of Machine Learning Research* 21:1–67
- Schick T, Schütze H (2021) Exploiting cloze-questions for few-shot text classification and natural language inference. In: Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pp 255–269
- Tóth L, Nagy B, Janthó D, Vidács L, Gyimóthy T (2019) Towards an accurate prediction of the question quality on stack overflow using a deep-learning-based nlp approach. In: Proceedings of the 14th International Conference on Software Technologies, pp 631–639
- Trienes J, Balog K (2019) Identifying unclear questions in community question answering websites. In: 41st European Conference on Information Retrieval, ECIR 2019, Springer, pp 276–289
- Vedantam R, Lawrence Zitnick C, Parikh D (2015) Cider: Consensus-based image description evaluation. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp 4566–4575
- Wang C, Yang Y, Gao C, Peng Y, Zhang H, Lyu MR (2022) No more fine-tuning? an experimental evaluation of prompt tuning in code intelligence. In: Proceedings of the 30th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering, pp 382–394
- Wang Y, Wang W, Joty S, Hoi SC (2021) Codet5: Identifier-aware unified pre-trained encoder-decoder models for code understanding and generation. In: Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, pp 8696–8708
- Wei B, Li Y, Li G, Xia X, Jin Z (2020) Retrieve and refine: exemplar-based neural comment generation. In: Proceedings of the 35th IEEE/ACM International Conference on Automated Software Engineering, pp 349–360
- Wilcoxon F (1992) Individual comparisons by ranking methods. Springer
- Xia CS, Zhang L (2023) Keep the conversation going: Fixing 162 out of 337 bugs for \$0.42 each using chatgpt. arXiv preprint arXiv:230400385
- Xu B, Hoang T, Sharma A, Yang C, Xia X, Lo D (2021) Post2vec: Learning distributed representations of stack overflow posts. *IEEE Transactions on Software Engineering* 48(9):3423–3441
- Yang G, Chen X, Zhou Y, Yu C (2022a) Dualsc: Automatic generation and summarization of shellcode via transformer and dual learning. In: 2022 IEEE International Conference on Software Analysis, Evolution and Reengineering (SANER), IEEE, pp 361–372
- Yang G, Liu K, Chen X, Zhou Y, Yu C, Lin H (2022b) Ccgir: Information retrieval-based code comment generation method for smart contracts. *Knowledge-Based Systems* 237:107858
- Yang G, Zhou Y, Chen X, Zhang X, Han T, Chen T (2023a) Exploitgen: Template-augmented exploit code generation based on codebert. *Journal of Systems and Software* 197:111577
- Yang G, Zhou Y, Chen X, Zhang X, Xu Y, Han T, Chen T (2023b) A syntax-guided multi-task learning approach for turducken-style code generation. *Empirical Software Engineering* 28(6):141
- Yang J, Hauff C, Bozzon A, Houben GJ (2014) Asking the right question in collaborative q&a systems. In: Proceedings of the 25th ACM conference on Hypertext and social media, pp 179–189
- Yazdani M, Lo D, Sami A (2021) Characterization and prediction of questions without accepted answers on stack overflow. In: 2021 IEEE/ACM 29th International Conference on Program Comprehension (ICPC), IEEE, pp 59–70
- Yin P, Deng B, Chen E, Vasilescu B, Neubig G (2018) Learning to mine aligned code and natural language pairs from stack overflow. In: Proceedings of the 15th international conference on mining software repositories, pp 476–486
- Zhang F, Yu X, Keung J, Li F, Xie Z, Yang Z, Ma C, Zhang Z (2022) Improving stack overflow question title generation with copying enhanced codebert model and bi-modal information. *Information and Software Technology* 148:106922
- Zhang F, Liu J, Wan Y, Yu X, Liu X, Keung J (2023) Diverse title generation for stack overflow posts with multiple-sampling-enhanced transformer. *Journal of Systems and Software* 200:111672
- Zhang Y, Yang Q (2021) A survey on multi-task learning. *IEEE Transactions on Knowledge and Data Engineering* 34(12):5586–5609
- Zhu J, Li L, Yang L, Ma X, Zuo C (2023) Automating method naming with context-aware prompt-tuning. arXiv preprint arXiv:230305771