

Impact of Large Language Models on Generating Software Specifications

Danning Xie, Byoungwoo Yoo, Nan Jiang, Mijung Kim, Lin Tan, Xiangyu Zhang, Judy S. Lee

Abstract—Software specifications are essential for ensuring the reliability of software systems. Existing specification extraction approaches, however, suffer from limited generalizability and require manual efforts. The recent emergence of Large Language Models (LLMs), which have been successfully applied to numerous software engineering tasks, offers a promising avenue for automating this process. In this paper, we conduct the *first* empirical study to evaluate the capabilities of LLMs for generating software specifications from software comments or documentation. We evaluate LLMs’ performance with Few-Shot Learning (FSL), enabling LLMs to generalize from a small number of examples, as well as different prompt construction strategies, and compare the performance of LLMs with traditional approaches. Additionally, we conduct a comparative diagnosis of the failure cases from both LLMs and traditional methods, identifying their unique strengths and weaknesses. Lastly, we conduct extensive experiments on 15 state-of-the-art LLMs, evaluating their performance and cost-effectiveness for generating software specifications.

Our results show that with FSL, LLMs outperform traditional methods (by 5.6%), and more sophisticated prompt construction strategies can further enlarge this performance gap (to 5.1 – 10.0%). Yet, LLMs suffer from their *unique* challenges, such as ineffective prompts and the lack of domain knowledge, which together account for 53 – 60% of LLM unique failures. The strong performance of open-source models (e.g., StarCoder) makes closed-source models (e.g., GPT-3 Davinci) less desirable due to size and cost. Our study offers valuable insights for future research to improve specification generation.

Index Terms—software specifications, large language models, few-shot learning

I. INTRODUCTION

ACCURATE and comprehensive software specifications are essential for ensuring the correctness, dependability, and quality of software systems [1]–[4]. Common software specifications include pre- and post-conditions to a target function that describe the constraints of input parameters and the expected behaviors or output values, which are often required or crucial for generating effective test cases and test oracles, symbolic execution, and abnormal behavior identification [3], [5]–[7].

Corresponding author: Mijung Kim

Danning Xie, Nan Jiang, Lin Tan, and Xiangyu Zhang are with the Computer Science Department, Purdue University, USA (email: {xie342, jiang719, lintan}@purdue.edu and xyzhang@cs.purdue.edu).

Byoungwoo Yoo and Mijung Kim are with the Computer Science and Engineering Department, UNIST, South Korea (email: {captainnemo9292, mijungk}@unist.ac.kr).

Judy S. Lee is with IBM Chief Analytics Office, Armonk, NY, USA (email: leesj@us.ibm.com).

The first two authors contributed equally to this paper.

Numerous approaches have been proposed to advance automation in extracting specifications from software texts (e.g., documents or comments), including rule-based methods [1], [3], [8], ML-based methods [4], [9], search-based methods [10], etc. For example, Jdoctor [1] leverages patterns, lexical, and semantic matching to translate code comments into machine-readable specifications of pre-/post-conditions, which enables automated test generation that leads to fewer false alarms and the discovery of more defects. Several other attempts have been made to further improve these processes in various domains [9]–[12]. However, most of existing work is domain-specific, relying on heuristics [1], [8] or a large amount of manually annotated data [3], [10]. This reliance makes it challenging to generalize these approaches to other domains.

With the emergence of Large Language Models (LLMs) [13]–[15], pre-trained on a tremendous amount of documents and source code [16]–[19], they have been applied to various Software Engineering (SE) tasks such as code generation [17], [20]–[23] program repair [24], and reasoning [25]. These models have demonstrated competitive performance compared to traditional approaches [17]–[19], [24], [26]. Given that software specification extraction predominantly involves the analysis and extraction from software texts, such as comments or documents, and the translation of natural language into (semi-)formal specifications, two research questions naturally arise: **(1) Are LLMs effective in generating software specifications from documentation? (2) What are the inherent strengths and weaknesses of LLMs for software specification generation compared to traditional approaches?**

A. Our Study

To fill in the gap, we conduct the *first* empirical study to evaluate the capabilities of LLMs in generating software specifications, in comparison with traditional approaches. First, due to the scarcity of labeled data in software specification extraction, we leverage LLMs with *Few-Shot Learning* (FSL) [27], a technique that enables LLMs to generalize from a limited number of examples. Second, we explore the potential of combining LLMs and FSL by investigating different prompt construction strategies and assessing their effectiveness. Third, we conduct an in-depth comparative diagnosis of the failure cases from both LLM and traditional approaches. This allows us to pinpoint their unique strengths and weaknesses, providing valuable insights to guide future research and improvement of

arXiv:2306.03324v2 [cs.SE] 2 Oct 2023

LLM applications. Lastly, we conduct extensive experiments, involving 15 state-of-the-art LLMs, and evaluate their performance and cost-effectiveness to facilitate model selection in software specification generation. We then describe our methodology and key findings.

FSL with random examples outperforms traditional methods: To assess the performance of LLMs with FSL, we first collect available datasets from the previous specification extraction research, Jdoctor [1] and DocTer [3]. These datasets contain software documents and comments, as well as the corresponding ground-truth specifications. We start with a basic prompt construction method that *randomly* selects examples for FSL, which helps LLMs learn to generate specifications from software texts. We then compare the results with the state-of-the-art specification extraction techniques, Jdoctor [1] and DocTer [3]. Our Finding 1 reveals that *with 10 – 60 randomly selected examples, LLMs achieve results that are comparable to (2.2% lower) or better than (5.6% higher) the state-of-the-art specification extraction techniques.*

Advanced prompt strategy enlarges the performance gap: To further explore the potential of LLMs with FSL for specification generation, we evaluate and compare different prompt construction strategies in terms of their impact on the performance of LLMs. Such prompt construction strategies include the above-mentioned random selection and a *semantics-based* selection strategy. Finding 2 shows that *with a more sophisticated prompt construction method, the performance gap between LLMs and traditional approaches is enlarged (to 5.1 – 10.0%).*

LLMs and traditional techniques exhibit unique strengths and weaknesses: Despite LLMs’ compelling performance over traditional techniques, our analysis of failure cases reveals noticeable differences in their failure patterns – Finding 3: *traditional approaches are more likely to generate empty specifications, while the failing cases of LLM approaches are dominant with ill-formed or incomplete ones.*

This variation in failures prompts us to analyze the distinct capabilities of LLM and traditional methods. We hence conduct an in-depth comparative diagnosis of the failure cases from both ends and investigate their root causes. This comprehensive analysis provides insights into their respective strengths and weaknesses. Notably, we identify several *unique* challenges of LLM. For example, Findings 4&5: *the two prevalent root causes combined (ineffective prompts and missing domain knowledge) result in 53 – 60% of unique LLM failures. In contrast, traditional rule-based methods experience unique failures (73 – 93%) due to insufficient or incorrect rules, as they are extracted manually or semi-automatically from limited datasets.*

Open-sourced StarCoder is the most competitive model: Lastly, given the vast spectrum of LLMs in terms of their open-source availability, cost, and model sizes, it becomes imperative to understand their capabilities. We perform rigorous experiments on 15 popular and state-of-the-art LLMs, e.g., StarCoder, GPT-3, etc., varying in designs and sizes, and evaluate their performance and cost-effectiveness in generating software specifications. Remarkably, our Findings 6–8 show that *most LLMs achieve better or comparable performance*

compared to traditional techniques. StarCoder is the overall most competitive model for generating specifications, with high performance, open-sourced flexibility and long prompt support. Its strong performance makes GPT-3 Davinci less desirable due to size and cost.

Identifying areas for future enhancement: These findings enable us to identify challenges for further improvement in LLM applications, i.e., *hybrid approaches*, that integrate LLMs and traditional methods, and *improving prompts effectiveness.*

B. Contributions

This paper makes the following key contributions:

- 1) We conduct the *first* empirical study comparing the effectiveness of LLMs and traditional methods in generating software specifications from comments or documents and find that LLMs with FSL achieve results that are comparable to (2.2% lower) or better than (5.6% higher) the traditional methods with only 10 – 60 randomly selected examples (Section V-A).
- 2) We evaluate the impact of different prompt construction strategies on the FSL performance and find that the advanced strategy can further enlarge the performance gap between LLMs and traditional methods (to 5.1 – 10.0%) (Section V-B).
- 3) We deliver a comprehensive failure diagnosis and identify *unique* strengths and weaknesses for both traditional methods and LLMs, facilitating future research directions (Section VI).
- 4) We extensively experiment on 15 state-of-the-art LLMs, assessing their performance and cost-effectiveness. StarCoder emerges as the most competitive model for its high performance, open-source flexibility, and long prompt support (Section VII).
- 5) We discuss inspirations of future improvements for LLMs’ performance on generating software specifications: hybrid approaches and improving prompt effectiveness (Section VIII).
- 6) We release the data and code for replicating our results ¹.

II. BACKGROUND

1) *Large Language Models (LLMs)*: LLMs are deep learning models pre-trained on a large corpus that contains natural language texts and programming language source code. They acquire general knowledge of natural languages (e.g., grammar, syntax, and semantics) and programming languages (e.g., code syntax and semantics) from the pre-training tasks such as masked span prediction [13], [28] and next token prediction [14], [15], [29]–[31].

While LLMs demonstrate strong performance on pre-training tasks, they are less effective in downstream tasks such as text summarization and specification extraction as they are different from the pre-trained tasks and domain-specific. Techniques like fine-tuning [13], [24], [29], [32] and few-shot learning [14], [15] are commonly used to adapt and enhance LLMs for these tasks. Fine-tuning requires additional

¹<https://github.com/DNXie/LLM4Spec-Data>

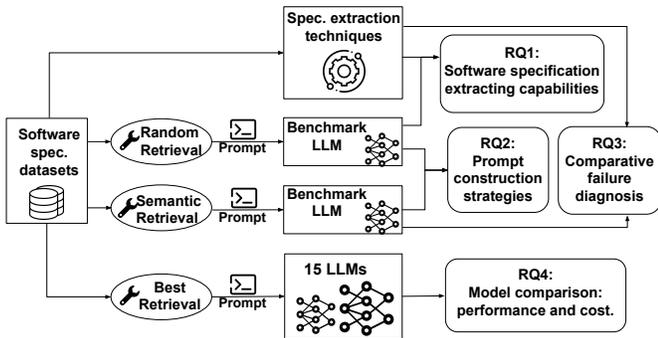


Fig. 1: Study overview

substantial amount of labeled data, while few-shot learning requires only a limited number of labeled instances.

2) *Few-Shot Learning (FSL)*: Few-shot learning (FSL) is crucial for enhancing pre-trained LLMs’ comprehension and performance in downstream tasks, especially when limited labeled data is available [15], [27], [33]–[35]. FSL provides the model with a small set of examples, demonstrating the task, allowing the model to generalize from these examples to complete the task. In this work, we focus on in-context few-shot learning [15], [36], where LLMs’ weights remain unchanged. FSL involves a small amount of labeled data into the LLMs’ input as prompts, guiding LLMs to effectively learn the specific task without modifying their weights or architecture.

III. STUDY SETUP

Fig. 1 presents an overview of our study. We collect available datasets from previous specification extraction work, containing software documents and comments and the corresponding ground-truth specifications. We then answer four research questions (RQs):

RQ1: How do LLMs with FSL compare with traditional rule-based approaches to extracting software specifications?

We apply the benchmark LLM to the collected datasets and compare its performance with the traditional approaches. For the LLM, we use a basic prompt construction strategy, *random retrieval*, to construct prompts with random examples (Section III-B2a). These prompts coach the LLMs to generate specifications from texts by examples.

RQ2: How do prompt construction strategies affect the performance of the LLM approaches? We compare the performance of different prompt construction strategies, i.e., *Random Retrieval* and *Semantic Retrieval*. Semantic retrieval selects examples based on semantic similarity to the target context (discussed in Section III-B2b).

RQ3: What are the unique strengths and weaknesses of LLMs and traditional approaches? To provide better insights into the capabilities of different approaches and shed light on future research, we conduct a comparative failure diagnosis. In particular, we sample a set of cases that an LLM approach succeeds while the traditional approach fails and vice versa. We analyze the symptoms and root causes of these failures, and identify their unique strengths and limitations, especially for the LLM approach.

Function signature	<code>isEmptyOrNull(java.lang.String string)</code>
Javadoc comment	<code>@return true if the string is null or is the empty string</code>
Specification	<code>string==null string.isEmpty() --methodResultID=true</code>

Fig. 2: An example data point from Jdoctor-data.

Function signature	<code>tf.image.extract_glimpse(input, size, offsets, ...)</code>
Document description	<code>input: A 'Tensor' of type 'float32'. A 4-D float tensor of shape [batch_size, height, width, channels].</code>
Specifications	<code>dtype: float32 structure: tensor shape: [batch_size, height, width, channels] ndim: 4 range: Null enum: Null</code>

Fig. 3: An example data point from DocTer-data.

RQ4: How do different LLMs compare in terms of their performance and cost for generating software specifications?

To assess performance and cost-effectiveness, we conduct extensive experiments with 15 state-of-the-art LLMs of different sizes, designs, and so forth. We employ the best-performing prompt construction strategy (“Best Retrieval” in the figure) based on the results of RQ2.

A. Existing Specification Extraction and Data

We compare with two state-of-the-art rule-based specification extraction approaches, Jdoctor [1] and DocTer [3], on their respective datasets, consisting of annotated comments or documents with associated specifications. To avoid confusion, we use the terms Jdoctor-data and DocTer-data to refer to the datasets, while Jdoctor and DocTer for the two approaches.

Table I presents the number of data points in each dataset. Rows “#Annotation” lists the number of document-specification pairs annotated for Jdoctor-data, and the parameter-specification pairs from each library of DocTer-data.

1) *Jdoctor-data*: It contains pre-/post-conditions written as executable Java expressions translated from corresponding Javadoc comments of `@return`, `@param`, and `@throws` tags. Fig. 2 shows an example: a Javadoc comment of `@return` tag describing the post-condition of function `isEmptyOrNull`, and the specification, with corresponding parts highlighted in matching colors.

2) *Jdoctor*: It translates Javadoc comments into specifications leveraging pattern, lexical, and semantic matching. It first identifies (subject, predicate) pairs in sentences. Fig. 2 presents an example, where the phrase “the string is null” contains such a pair with the subject “the string” and predicate “is null”. It then translates these pairs into executable procedures by applying manually defined heuristics (e.g., the phrase “is null” to Java expression `==null`). Next, it applies lexical matching to match a subject or a predicate to the corresponding code element based on lexical similarity (e.g., matching “the string” to input parameter `string`). Finally, it uses semantic matching to process meanings such as negation.

3) *DocTer-data*: It contains DL-specific specifications from the API documents of DL libraries, TensorFlow, PyTorch,

<pre>Signature: <x_j-signature> Javadoc comment: <x_j-comment> Condition: <y_j> ... Signature: <x_k-signature> Javadoc comment: <x_k-comment> Condition: <y_k> Signature: <x_{target}-signature> Javadoc comment: <x_{target}-comment> Condition:</pre>	<pre>Signature: <x_j-signature> Description: <x_j-param.> - <x_j-descp.> Constraints: <y_j> ... Signature: <x_k-signature> Description: <x_k-param.> - <x_k-descp.> Constraints: <y_k> Signature: <x_{target}-signature> Description: <x_{target}-param.> - <x_{target}-descp.> Constraints:</pre>
---	--

(a) Jdoctor-data

(b) DocTer-data

Fig. 4: Prompt structures with target highlighted in orange.

TABLE I: Software specification datasets

Jdoctor-data [1]	Tag	Type	@param	@return	@throws	Total
	#Annotations		243	139	472	854
DocTer-data [3]	Library	TensorFlow	PyTorch	MXNet	Total	
	#Annotations	1008	484	1,384	2,876	

and MXNet, with an example data point in Fig. 3. DocTer categorizes specifications into four groups: *dtype* for data types, *structure* for data structures, *shape* for the parameter’s shape or the number of dimensions (ndim), and *valid value* for the set of valid values (enum) or valid range.

4) *DocTer*: It automatically constructs rules that map syntactic patterns of API descriptions to specifications using a portion of annotated data. These rules are then applied to API documents to extract parameter specifications. For instance, DocTer automatically constructs a rule that maps the syntactic pattern “of type D_TYPE” to the *dtype* specification, where “D_TYPE” is a placeholder for any data type keyword. By applying this rule to the description of parameter input “A Tensor of type float32” (shown in Fig. 3), the *dtype* specification “float32” is extracted.

B. Specification Generation with LLMs and FSL

To extract specifications using an LLM, we construct prompts with examples to coach the LLM, i.e., via few-shot learning.

1) *Few-Shot Learning (FSL)*: Consider a dataset $D = \{(x_i, y_i)\}_{i=1}^I$, where x represents the *context* and y denotes the expected *completion*. In our study, x corresponds to the document or comments while y refers to the target software specifications. For each data point (x_{target}, y_{target}) in D , we first select K indices, denoted by S . Then, $\{(x_k, y_k)\}_{k \in S}$ are K (i.e., $|S|$) examples of input context x_k and their corresponding ground-truth completions y_k .

There are multiple ways to sample the examples; for instance, one can split the dataset into D_{train} and D_{test} , and only select examples from D_{train} for the target in D_{test} . In our study, to utilize as many cases as possible and explore the upper limit of FSL, we use leave-one-out cross-validation [37]–[39], where we select the K examples from all other $|I| - 1$ data points in D except for the target itself. These examples are integrated into the prompt $P = \{(x_k, y_k)\}_{k \in S} + x_{target}$, where the target context is appended to the end of the prompt. Subsequently, the LLMs generate the completion y_{out} for the target context. This output is then compared against the ground-truth completion, y_{target} .

<pre>Signature: min(float[] array) Javadoc: @param array a nonempty array of float values Annotation: (array.length==0)==false Generated: array.length>0</pre>

Fig. 5: Equivalent conditions: a Jdoctor-data example with semantically correct LLM completion.

2) *Prompt Construction*: Fig. 4 presents example prompts for Jdoctor-data and DocTer-data. In the prompt for Jdoctor-data, each of the K examples includes a function signature $\langle x_k\text{-signature} \rangle$, a Javadoc comment $\langle x_k\text{-comment} \rangle$, and the corresponding annotated condition $\langle y_k \rangle$. Then, the target x_{target} is appended to the end of the prompt (highlighted in orange).

Similarly, the prompt for DocTer-data (Fig. 4b) includes a function signature $\langle x_k\text{-signature} \rangle$, a document description $\langle x_k\text{-param.} \rangle - \langle x_k\text{-descp.} \rangle$, where $\langle x_k\text{-descp.} \rangle$ is the corresponding document description for the parameter $\langle x_k\text{-param.} \rangle$, and the annotated constraints $\langle y_k \rangle$.

We treat the three tag types of Jdoctor-data (Table I) and the DocTer-data from three libraries as separate sub-datasets. For example, for a target of @param tag, we only select K examples from the other 242 data points of @param tag, excluding itself. Similarly, in DocTer-data, for a target, we select examples from the same library, excluding itself. We study two commonly used strategies to choose the K examples: *random retrieval* and *semantic retrieval*.

a) *Random Retrieval*: The *random retrieval* strategy is randomly selecting K examples from the dataset, excluding itself, i.e., $D \setminus \{(x_{target}, y_{target})\}$. For example, when experiment on a target with $K = 20$, we randomly select 20 instances from the dataset excluding the target itself and use the 20 instances as examples in the prompt.

b) *Semantic Retrieval (SR)*: The *SR* strategy aims to select examples semantically close to the target context, and are proven by previous studies [40], [41] to be more effective than random retrieval. For example selection, BERT-based models are more effective compared to traditional methods, such as TF-IDF [40], [41]. We choose RoBERTa [42] as the SR retrieval model, specifically, Roberta-large [43] fine-tuned for sentence similarity using a contrastive learning objective, due to its outstanding performance on the Semantic Textual Similarity (STS) dataset [44]. We do not study the impact of different retrieval models since previous research suggests that the differences in results are small [45].

3) *Post-Processing of LLM Output*: For Jdoctor-data, LLM-generated completions may be semantically correct, albeit not identical to the annotations provided in the dataset. Fig. 5 illustrates an example of a Jdoctor-data data point, including the function signature, Javadoc comment, annotated specification (in yellow), and the specification generated by LLM (in blue). The annotated and generated specifications convey the same requirement that the variable `array` must be non-empty. An automatic script that checks for a perfect match would label it as incorrect. Such equivalent but syntactically different specifications frequently occur in LLM results and pose a

TABLE II: Studied LLMs: size (#Param, * for estimates), token limit, price per 1,000 tokens, and open-source status.

Model	#Param	Token limit	Price (per 1K)	Open-source?
StarCoder [46]	15.5B	8,192	-	✓
GPT-3 [15]	Davinci (175B*) Curie (Unknown)	4,000 2,048	\$0.02 \$0.002	✗ ✗
GPT3.5 [47]	Turbo (Unknown)	4,096	0.002	✗
BLOOM [48]	176B	1,000	-	✓
CodeGen [17]	16B, 6B, 2B, 350M	2,048	-	✓
CodeGen2 [49]	16B, 7B, 3.7B, 1B	2,048	-	✓
InCoder [18]	6B, 1B	2,048	-	✓

challenge for automatic assessment, especially when domain-specific knowledge is needed (e.g., the length of an array is non-negative). Therefore, we manually inspect the generated completions for Jdoctor-data that deviate from the annotation and report both raw accuracy automatically calculated by the script and the final accuracy after manual correction. Here, two of the authors conducted the manual correction independently, and resolve any disagreements with a third author.

C. Studied Large Language Models

Table II summarizes 15 LLMs from seven series that we study in this paper, including the best generic LLMs, code-specific LLMs, and open-sourced LLMs. We select GPT-3.5 over ChatGPT, as the changes made for ChatGPT focus on conversations [50].

a) *StarCoder* [46]: StarCoder is the state-of-the-art open-sourced LLM with 15.5B parameters. It uses an autoregressive transformer [51] architecture that supporting multi-query attention mechanism [52], and is pre-trained on The Stack [53] dataset using next token prediction and code infilling tasks.

b) *GPT-3* [15] and *GPT-3.5* [47]: They are two series of autoregressive language models developed by OpenAI. We study the two best GPT-3 models: GPT-3-Davinci (175B), the largest and most capable one, and GPT-3-Curie with unknown size. GPT-3.5 is subsequently fine-tuned with “Reinforcement Learning from Human Feedback (RLHF)” [54]–[57] technique to improve the quality of text completion and chat. We access them via OpenAI’s API through model `text-davinci-003`, `text-curie-001`, and `gpt-3.5-turbo`.

c) *BLOOM* [48]: It is an open-sourced model similar to GPT-3 in structure. BLOOM is pre-trained on a 1.61 TB corpus containing both natural language and source code data [48]. We select BLOOM-176B, which has 176 billion parameters. We study this model by accessing its Hugging Face API [58], which has a 1,000 token limit.

d) *CodeGen* [17] and *CodeGen2* [49]: They are two series of autoregressive LLMs for program synthesis, pre-trained on THEPILE [59] and BigQuery [60]. We study CodeGen-350M/2B/6B/16B-Multi models, which are expected to generalize well to both Jdoctor-data and DocTer-data, and CodeGen-350M/2B/6B/16B-Mono pre-trained additionally on BigPython [17], providing an advantage for DocTer-data’s Python-related tasks. We also examine CodeGen2-1/3.7/7/16B,

which use the same architecture as CodeGen models, but are pre-trained on The Stack [53] dataset and additionally support code infilling task.

e) *InCoder* [18]: InCoder is a series of LLMs based on the XGLM [61] architecture, employing code infilling as its pre-training task. The pre-training data contains 57 GB of natural language text and 159 GB of source code collected from platforms like GitHub and StackOverflow. We study the InCoder-1B and InCoder-6B models.

We run StarCoder, CodeGen, CodeGen2, and InCoder locally using 4 Nvidia RTX A6000 GPUs with 48GB memory.

D. Benchmark LLM

We choose to use StarCoder as the benchmark LLM for the study of RQ1, RQ2, and RQ3 (Fig. 1). StarCoder is a state-of-the-art open-source code LLM, ensuring reproducibility of our results. Additionally, it supports more input tokens than all other models (Table II), allowing a wide range of experimental settings, such as different numbers of examples in the prompt.

After identifying the best prompt construction strategies (“Best Retrieval” in Fig. 1) using StarCoder, we apply it to all other LLMs listed in Table II for the study of RQ4.

IV. EXPERIMENTAL SETTINGS

A. Model Settings

As per our experimental design, we truncate examples from the beginning of the prompts to fit the token limit of each model (as shown in Table II). If the median number of tokens in the prompts exceeds the limit, we skip that experiment. For instance, we skip the experiment for Jdoctor-data when $K = 60$ for BLOOM (with a token limit of 1,000) since the median number of tokens in the prompts is 3,737. To provide comprehensive results, we run the benchmark model (StarCoder) for RQ1 and RQ2 in all settings.

We follow the recommended practice [62] of using the separator `\n\n###\n\n` between examples in the prompts. For the generative models, we use the same separator as the stopword. We set the models’ `temperature` to 0 for minimal randomness and to reflect the models’ most confident answers. The `max_tokens` parameter is set to 100 for Jdoctor-data and 200 for DocTer-data, larger than the number of tokens in the longest completion in each dataset.

B. Accuracy and F1 Metrics

To evaluate the correctness of the generated specifications, we employ the accuracy metric for the generated specifications for Jdoctor-data. As LLMs under consideration inherently generate a specification for all given Javadoc comments, standard metrics such as precision and recall do not provide additional information in this setting. We define $Accuracy = \frac{C}{|D|}$, where C is the number of accurately generated specifications, and $|D|$ is the total number of specifications annotated in D .

For DocTer-data, we use precision, recall, and F1 metrics to evaluate the generated results for each specification category (e.g., *dtype*). For category t , let C_t be the number of correctly generated specifications, N_t be the total number of annotated

TABLE III: Comparison of StarCoder with random prompt construction and Jdoctor: Accuracy (%).

Approach	K	@param		@return		@throws		Overall	
		Raw	Processed	Raw	Processed	Raw	Processed	Raw	Processed
Jdoctor	-	97.0	97.0	69.0	69.0	79.0	79.0	83.0	83.0
StarCoder	10	82.7	91.4	28.8	39.6	74.2	83.7	69.2	78.7
StarCoder	20	84.4	93.0	35.5	51.8	76.5	84.1	72.1	81.4
StarCoder	40	88.9	93.8	46.0	59.0	80.7	90.7	77.4	86.4
StarCoder	60	90.5	94.7	54.0	67.6	82.0	91.7	79.9	88.6

TABLE IV: Comparison of StarCoder with random prompt construction and DocTer: Precision/Recall/F1 (%).

Approach	K	TensorFlow	PyTorch	MXNet	Overall
DocTer	898	90.0/74.8/81.7	78.4/77.4/77.9	87.9/82.4/85.1	85.4/78.2/81.6
StarCoder	10	66.4/71.2/68.7	69.0/64.9/66.9	68.6/66.1/67.3	67.9/67.7/67.7
StarCoder	20	73.5/69.7/71.5	72.2/69.3/70.7	70.6/66.1/68.3	71.9/67.9/69.8
StarCoder	40	73.4/79.7/76.4	76.2/76.6/76.4	79.6/74.5/77.0	76.9/76.7/76.7
StarCoder	60	78.2/78.8/78.5	78.6/78.2/78.4	81.1/79.8/80.5	79.7/79.2/79.4

specifications in the dataset for category t (i.e., D_t), and G_t denote the number of generated specifications for category t . We define precision as $P_t = \frac{C_t}{G_t}$, recall as $R_t = \frac{C_t}{N_t}$, and F1 score as $F_t = 2 \cdot \frac{P_t \cdot R_t}{P_t + R_t}$. We report the overall precision, recall, and F1 across all four categories (*dtype*, *structure*, *shape*, and *valid value*) for each library (e.g., TensorFlow), as no significant performance differences were observed among different categories.

As discussed in Section III-B2, we treat data from Jdoctor-data of different tag types and data from DocTer-data of different libraries as separate datasets. We report the accuracy and F1 metrics for them separately, as well as the overall accuracy/F1.

V. EVALUATION RESULTS

A. RQ1: Specification Extracting Capabilities

We evaluate LLMs on both Jdoctor-data and DocTer-data using random retrieval strategy for prompt construction (Section III-B2a) with StarCoder as the subject model (Section III-D). Results for both StarCoder and baseline methods (Jdoctor and DocTer) are presented in Table III and Table IV.

As described in Section III-B3, we manually post-process the specification generated for Jdoctor-data and present the raw accuracy (automatically calculated by the script) in col. ‘‘Raw’’ and the final accuracy (after manual correction) in col. ‘‘Processed’’ in Table III. For Jdoctor (row *Jdoctor*), these two values are the same.

The columns K in Table III and IV represent the number of examples in the prompts (Section III-B1). As Jdoctor relies on manually defined heuristics, we leave the corresponding cell blank (‘‘-’’) in Table III. For DocTer, we list the average number of annotated examples used for rule construction for each library ($2,696 / 3 = 898$) as K .

a) Results for Jdoctor-data: As shown in Table III, StarCoder achieves an overall accuracy of 86.4% using only 40 randomly chosen examples from each comment type (i.e., 120 in total) and outperforms Jdoctor, which has an accuracy of 83.0%. StarCoder outperforms Jdoctor in terms of accuracy for @throws comments, using just 10 randomly selected

TABLE V: Comparison of StarCoder using SR prompt construction and Jdoctor: Accuracy (%).

Approach	K	@param		@return		@throws		Overall	
		Raw	Processed	Raw	Processed	Raw	Processed	Raw	Processed
Jdoctor	-	97.0	97.0	69.0	69.0	79.0	79.0	83.0	83.0
StarCoder + SR	10	91.8	95.4	63.3	74.1	83.9	90.0	82.8	88.9
StarCoder + SR	20	94.2	97.9	65.5	78.4	86.2	91.1	85.1	91.0
StarCoder + SR	40	95.1	97.9	63.3	74.8	88.8	93.4	86.4	91.7
StarCoder + SR	60	95.9	98.4	64.7	77.7	90.3	94.7	87.7	93.0

TABLE VI: Comparison of StarCoder using SR prompt construction and DocTer: Precision/Recall/F1 (%).

Approach	K	TensorFlow	PyTorch	MXNet	Overall
DocTer	898	90.0/74.8/81.7	78.4/77.4/77.9	87.9/82.4/85.1	85.4/78.2/81.6
StarCoder + SR	10	81.9/81.7/81.8	79.7/77.1/78.4	86.3/86.1/86.2	83.6/83.0/83.3
StarCoder + SR	20	83.2/84.4/83.8	81.3/80.7/81.0	86.8/88.1/87.5	84.6/85.6/85.1
StarCoder + SR	40	84.8/86.1/85.4	82.6/83.1/82.8	87.6/89.0/88.3	85.8/87.0/86.4
StarCoder + SR	60	85.0/86.7/85.8	83.4/84.6/84.0	87.4/89.2/88.3	85.9/87.5/86.7

examples, and achieves comparable performance to Jdoctor for @param and @return comments.

b) Results for DocTer-data: As shown in Table IV, with only 60 random examples from each library (180 in total), StarCoder achieves an overall F1 score of 79.4%, only 2.2% lower than DocTer which has an F1 score of 81.6% and requires 2,696 annotated examples. Furthermore, StarCoder surpasses DocTer for PyTorch in terms of F1 score using only 60 examples, and achieves comparable performance for TensorFlow and MXNet with DocTer.

Finding 1: StarCoder, with a small number (10 – 60) of randomly selected examples, achieves comparable results (2.2% lower) with the state-of-the-art specification extraction technique DocTer, and outperforms Jdoctor by 5.6%.

B. RQ2: Prompt Construction Strategies

Table V and Table VI reveal that StarCoder, when employing the SR strategy, outperforms both Jdoctor and DocTer, even with only 10 examples selected in the prompt from each type/category.

Fig. 6 showcases the effectiveness of random and SR strategies across prompt sizes. SR strategy consistently outperforms both the random strategy and traditional specification techniques across different prompt sizes, highlighting the importance of an appropriate prompt construction strategy for improved FSL performance.

Finding 2: The semantic retrieval strategy further improves StarCoder’s performance, leading to a 5.9 – 10.0% improvement over Jdoctor and a 1.7 – 5.1% increase over DocTer.

VI. RQ3: COMPARATIVE FAILURE DIAGNOSIS

In light of the outstanding performance of LLM compared with traditional techniques, it is crucial to delve deeper into their strengths and limitations. We manually examine failing cases of both LLM-based and baseline approaches (i.e., Jdoctor and DocTer), and study their failure symptoms and root

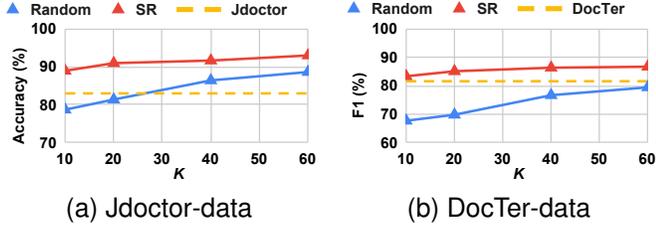


Fig. 6: Comparison of FSL performance using Random and Semantic Retrieval (SR) for prompt sizes (K) from 10 to 60.

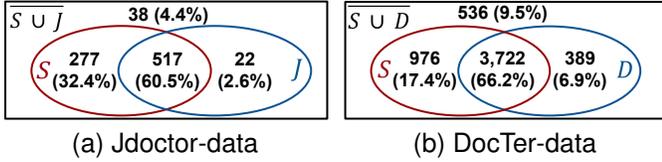


Fig. 7: Venn diagrams of specification generation for the two datasets. S : StarCoder; J : Jdoctor; D : DocTer.

causes in a comparative manner, aiming to provide insights and directions for future techniques.

Fig. 7 presents the comparative performance of the StarCoder-based LLM method and baseline methods as Venn diagrams. Fig. 7a and 7b respectively illustrate the number of unique and shared successes and failures of specifications for StarCoder (S) against Jdoctor (J) and DocTer (D). The numbers in the intersections ($S \cap J$ and $S \cap D$) denote cases where both methods are correct, while numbers in sections $\overline{S \cup J}$ and $\overline{S \cup D}$ indicate cases they both fail. The presented results are derived from the experiment using StarCoder with SR and $K = 60$.

Fig. 7a shows that both the LLM and Jdoctor perform well on the majority of cases (60.5%), indicating that the LLM quickly learns most specification extraction rules from a small number of examples in the prompts. The LLM has more (29.8%) unique correct cases than Jdoctor, indicating the generalizability of LLMs from extensive pretraining. There are a few cases that both methods fail, possibly due to inherent difficulties such as incomplete software text.

To better understand the pros and cons of different methods, we investigate the symptoms and the underlying causes of the failing cases. We randomly sample 15 cases from each section of Fig. 7 where at least one of the methods fails, i.e., $S \cap \overline{J}$, $\overline{S} \cap J$, and $\overline{S \cup J}$ in Fig. 7a. The sampling results in 90 failing cases of StarCoder and 90 of the baseline methods. Two authors categorize these cases independently, with a third author resolving disagreements.

A. Failure Symptom Analysis

We conduct further analysis on the distributions of failure symptoms in various sections of the Venn diagrams (Fig. 7), and present our results in Fig. 8. The failure symptoms are classified into four categories, “ill-formed”, “incorrect”, “incomplete”, and “empty”. The bars illustrate the distributions

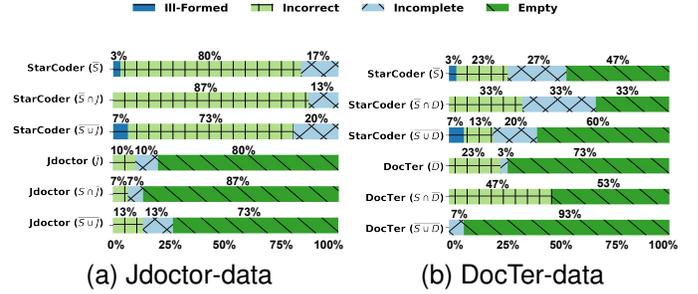


Fig. 8: Distributions of symptoms in failing cases across approaches and datasets. The y-axis is “approach (section)”.

of failure symptoms for these methods (StarCoder, Jdoctor, and DocTer) from different sections of Fig. 7. For example, bar “StarCoder (\overline{S})” in Fig. 8a denotes the distributions of StarCoder’s failure symptoms on Jdoctor-data cases where it fails. Bar “StarCoder ($\overline{S} \cap J$)” indicates the distributions of StarCoder’s failure symptoms within the $\overline{S} \cap J$ section (StarCoder fails, and Jdoctor succeeds). When both fail, as in $\overline{S \cup J}$, “StarCoder ($\overline{S \cup J}$)” and “Jdoctor ($\overline{S \cup J}$)” shows the failure symptoms categories for StarCoder and Jdoctor outputs, respectively.

Category “*ill-formed*” refers to invalidly formed generated specifications. For Jdoctor-data, this denotes syntactical errors or grammatical mistakes. For DocTer-data, it refers to improper specifications forms, such as generating a boolean instead of an expected numerical range. “*Incorrect*” indicates specification errors, while “*Incomplete*” denotes the specification is a strict subset of the ground truth. For DocTer-data, if the ground-truth specifications for the *dtype* category include both `int32` and `int64`, generating only `int32` is incomplete, whereas generating `bool` is incorrect. “*Empty*” denotes a missing specification. For DocTer-data, each specification contains four categories (Section III-A3). If the specification for one of the categories is empty while the ground truth is not, it is considered an empty type of failure.

a) *LLMs are more likely to generate ill-formed and incomplete specifications:* A small fraction (3%) of StarCoder’s failing cases are ill-formed, while traditional methods do not exhibit this issue. Traditional methods are rule-based and therefore guarantee the validity of the generated specifications. Whereas LLMs are generative and do not ensure the correctness of their output. Fig. 8 also shows that LLMs are more likely to generate incomplete specifications. For instance, when given the description “tuple/list of 1 int”, StarCoder only generates `tuple` for the *structure* specification and misses the specification `list`. Such errors are unlikely to occur with rule-based methods, as they match the entire sequence of types and extract them directly from the document as specifications, while LLMs use sampling to decode outputs from a distribution, which may occasionally miss tokens.

b) *Traditional techniques are much more likely to generate empty specifications than LLMs:* For both datasets, the most common failure symptom among baseline methods is “Empty” (80% and 73%), where the methods fail to generate any specifications. For StarCoder, none of the failing cases for

```

...
Signature: tf.sparse.fill_empty_rows(sp_input,...)
Description: sp_input - A 'SparseTensor' with shape '[N, M]'.
Constraints: shape: [N, M]
ndim: Null
Prompt

+ Signature: tf.math.lbeta(x,...)
+ Description: x - A rank 'n + 1' 'Tensor', 'n >= 0' with type 'float', or 'double'.
+ Constraints: shape: Null
ndim: >=1
Additional Example

+ Signature: tf.sparse.softmax(sp_input,...)
Description: sp_input - 'N'-D 'SparseTensor', where 'N >= 2'.
- Constraints: shape: Null
ndim: Null
Completion

+ ndim: >=2
New completion
    
```

Fig. 9: Example of “Ineffective Prompt”. Yellow, blue, and green denote the original prompt (simplified), generated completion, and an added example that enables LLM to generate the correct specification.

Jdoctor-data are empty, while only 47% of the failing cases for DocTer-data are empty. The reason is that the Jdoctor-data examples provided to the LLMs are never empty and hence the LLMs always generate some results for queries for Jdoctor-data. In contrast, some of the specification categories (e.g., dtype) of DocTer-data examples (provided to the LLMs) may be empty. The LLMs learn that empty is a possible result for DocTer-data. Since the baselines are rule-based, when rules are inapplicable, they tend to produce empty results. In contrast, LLMs produce results by predicting missing tokens. They may produce empty results when empty is a legitimate token. The results are in line with the precision and recall in Table VI, that StarCoder has a much higher recall than DocTer (e.g., 84.0% versus 72.8%) with a comparable or slightly lower precision (e.g., 84.3% versus 86.8%). Although with a lower empty rate, LLM tends to generate incomplete and ill-formed specifications as discussed in (a).

Finding 3: Compared to LLMs, traditional specification extraction approaches are much more likely to generate empty specifications (e.g., 73 – 80% versus 0 – 47%), while LLMs are more likely to generate ill-formed or incomplete specifications.

B. Root Cause Analysis

In this section, we categorize the root causes of failures and study their distributions. At the end, we perform a comparative study based on the sections in the Venn diagrams (Fig. 7).

1) *LLM Failure Root Causes:* Since LLM results are difficult to interpret, it is in general difficult to determine the root causes of failing cases by LLMs. We hence determine the root cause by finding a fix for it. The nature of the fix indicates the root cause. In some cases, the failure may be fixed in multiple ways. We consider the one requiring the least effort as the root cause and categorize them into five categories. Fig. 11 presents the distributions of the root causes of StarCoder in different sections of the Venn diagrams (Fig. 7). We now explain the five categories.

a) *Ineffective Prompts:* It means that the failure is due to the ineffectiveness of the examples in prompt, even with SR. Although SR significantly improves FSL’s performance (Section V-B), it occasionally falls short in selecting the appropriate examples. If we can fix a failing case by manually

```

+ Relevant functions:
+ searchForDanger(int range,float threat)
+ ...
Additional domain knowledge

... (K examples)
Signature: isInDanger(int range,float threat)
Javadoc comment: @return True if a threat was found.
prompt

- Condition: this.getFile().isInDanger(range,threat) -> methodResultID==true
+ Condition: this.searchForDanger(range,threat) -> methodResultID==true
    
```

Fig. 10: Example of “Missing Domain Knowledge”. Yellow, blue, and green denote the original prompt (simplified), generated completion, added domain knowledge, and the new completion.

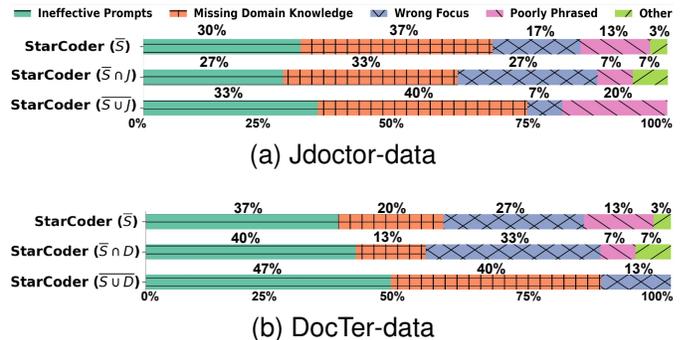


Fig. 11: Distributions of root causes of StarCoder’s failures.

selecting more relevant example(s) to the prompt, or simply altering the order of the examples in the original prompt, we consider the failure is due to ineffective prompts.

According to bars “StarCoder (\bar{S})” from Fig. 11a (Jdoctor-data) and Fig. 11b (DocTer-data), 30% and 37% StarCoder’s failures on the two datasets are due to this reason. We find that the order of examples plays a crucial role, as 21% of the failure cases in this category are resolved by rearranging the order of the examples.

Fig. 9 presents a portion of the prompt for the target parameter `sp_input`. StarCoder fails to generate the specification `ndim:>=2`, which is not explicitly stated in the description and requires StarCoder to comprehend the implicit relationship between `N` and its value range, i.e., “`N ≥ 2`”. Adding an example with such implicit constraints enables StarCoder to generate the correct specification.

b) *Missing Domain Knowledge:* This refers to LLM failures due to insufficient domain knowledge. For instance, in Fig. 10, StarCoder generates a specification using a non-existent function, `isInDanger`, instead of `searchForDanger`. This issue arises as the LLM lacks relevant context, such as the methods in the class, while Jdoctor employs a search-based approach examining all methods in the relevant classes. This result uncovers LLMs’ limitation compared to traditional search-based methods: a deficiency in domain knowledge. To validate our hypothesis, we manually incorporate relevant domain knowledge into the prompt, alongside the provided examples. StarCoder then successfully generates the accurate specification, utilizing the correct function (in green). Fig. 11 shows that 37% and 20% of StarCoder’s failures are due to missing domain knowledge.

c) *Wrong Focus:* It denotes instances where LLMs fail to focus on crucial keywords or are misguided or diverted

values: 1-D or higher numeric `Tensor`.
 values: 1-D or higher `numeric` `Tensor`.

Fig. 12: Example of “Wrong Focus”. Yellow and blue highlight the original document and the modified version. A minor adjustment (quoting keyword) enables the LLM to generate the correct specification.

by other content. For example, Fig. 12 shows the LLM fails to generate the correct specification from the original document (in yellow), `numeric`, which specifies the data type of the input tensor. By employing a slightly revised description that merely quotes the keyword, StarCoder successfully generates the specification. Fig. 11 reveals that 17% and 27% of StarCoder’s failures are due to the wrong focus. To identify such failures, we apply three input mutation strategies: *minor modifications* by simply adding quotation marks to the keywords; *rewriting the sentence* while preserving the same syntactic structure (e.g., changing “A or B” to “B or A”); and *deleting redundant content* to help the LLM concentrate on the essential parts. All three strategies involve simple and semantics-preserving mutations and do not have impacts on the rule-based methods like DocTer as they rely on syntactic structure. We find that 42% of such cases can be resolved by merely quoting the keyword(s).

d) *Poorly Phrased*: It refers to instances where the original documents or comments are poorly written, ambiguous, or hard to understand even for humans. Rewriting the sentence to clarify its meaning enables LLMs to generate correct answers. According to Fig. 11, 13% of StarCoder’s failures are due to this reason.

e) *Others*: We group the non-prevalent categories as “others”, including “contradictory document” and “Unclear”, which together takes only 3% of the StarCoder’s failures. The former refers to buggy or self-contradictory comments or documents, causing discrepancies between dataset annotations and LLM-generated specifications. We identified one buggy document and reported it to the developers. “Unclear” indicates failures with unclear root causes, which we fail to fix despite various attempts.

Finding 4: The two dominant root causes combined (ineffective prompts and missing domain knowledge) result in 57 – 67% of StarCoder failures on the two datasets.

2) *Jdoctor and DocTer Failure Root Causes*: We manually investigated the sampled failing cases for Jdoctor and DocTer and identified three root causes: *missing rule*, *incomplete semantic comprehension*, and *incorrect rule*. The distributions are in Fig. 13.

a) *Missing Rule*: It refers to the absence of relevant rules or patterns, usually resulting in “empty” specifications (Section VI-A). A notable 93% and 63% failing cases of Jdoctor and DocTer respectively fall into this category, exposing a limitation of rule-based methods: heavily dependent on manually defined or limited rules.

b) *Incomplete Semantic Comprehension*: It describes instances where rule-based methods successfully match a part

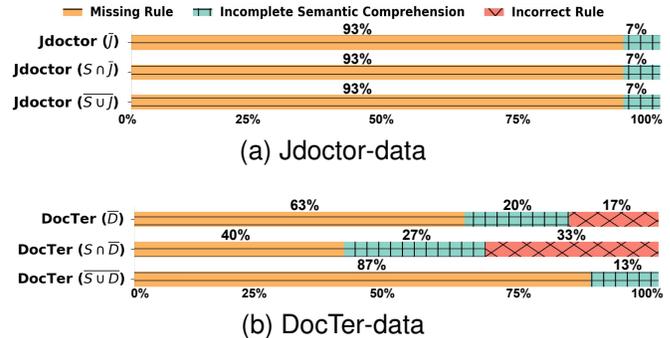


Fig. 13: Distributions of root causes for baselines’ failures.

of the sentence but fail to grasp the semantics of the entire sentence. This lack of comprehensive understanding leads to incorrect results. For instance, DocTer extracts a *structure* specification `vector` from the description “Initializer for the bias vector” but does not consider the full context or relationships among elements, ultimately impacting the correctness of the extracted specifications. In Fig. 13, 7% and 20% of the failing cases in the two datasets are due to this reason.

c) *Incorrect Rules*: This category denotes the cases where the applied rules are incorrect. DocTer automatically constructs the rules (from syntactic patterns to specifications) based on their co-occurrence in the annotated dataset, which can potentially introduce incorrect rules, leading to incorrect extractions. In Fig. 13, 17% of DocTer’s failing cases are due to *incorrect rules*, while Jdoctor does not have any of such failures since their rules are all manually defined.

3) *Comparative Root Cause Analysis*: We compare the root cause distributions of StarCoder (Fig. 11), with those of the baseline methods (Fig. 13).

Bars “StarCoder ($\bar{S} \cap J$)” and “StarCoder ($\bar{S} \cap D$)” show that in cases where the baseline methods succeed, StarCoder’s dominating failure causes are ineffective prompts and wrong focus. Notably, the wrong focus is particularly prevalent here, in contrast to sections $S \cup J$ and $S \cup D$ where both approaches fail.

From bars “Jdoctor ($S \cap \bar{J}$)” and “DocTer ($S \cap \bar{D}$)” where StarCoder succeeds and the baseline methods fail, we observe that the unique baseline failures are primarily due to missing rules and incomplete semantic comprehension. In other words, when prompts and software texts are of high quality, *LLMs demonstrate outstanding generalizability*, unbounded by rule sets. They make predictions based on entire descriptions rather than just parts of them. Notably, 17% of DocTer’s unique failures are due to incorrect rules, all of which can be addressed by StarCoder as suggested in the figure. We suspect that any automatic rule inference techniques may suffer from such problems if human corrections are not in place.

Finding 5: Compared to traditional methods, StarCoder suffers from ineffective prompts and wrong focuses, which are inherent to LLMs and cause 54 – 73% *unique* StarCoder failures. In contrast, LLMs show outstanding *generalizability* when rules, extracted manually or semi-

TABLE VII: Comparison of different LLMs with SR on Jdoctor-data: Accuracy (%) and Cost.

Approach/ Model (+SR)	K	@param		@return		@throws		Overall		Cost (\$)	
		Raw	Processed	Raw	Processed	Raw	Processed	Raw	Processed		
Jdoctor	-	97.0	97.0	69.0	69.0	79.0	79.0	83.0	83.0	-	
StarCoder	60	95.9	98.4	64.7	77.7	90.3	94.7	87.7	93.0	0	
GPT-3	davinci	10	89.1	98.4	64.6	78.5	85.4	93.2	84.0	92.9	12.6
		20	93.4	98.9	67.1	83.5	87.1	93.0	86.5	93.5	25.2
		40	94.0	98.4	68.4	83.5	90.0	94.7	88.6	94.4	50.4
	curie	60	95.1	98.4	68.4	83.5	91.3	96.6	89.7	95.6	75.6
		10	60.1	61.2	35.4	44.3	51.2	53.2	51.8	54.3	1.3
		20	76.0	78.1	31.6	44.3	64.6	65.5	63.8	66.4	2.6
GPT-3.5	turbo	10	86.9	93.5	60.8	70.9	81.6	91.0	80.6	89.3	1.3
		20	82.0	89.1	60.8	72.2	83.5	90.3	80.4	87.9	2.6
		40	83.6	89.6	59.5	74.7	84.5	88.8	81.3	87.4	5.2
		60	86.3	89.1	59.5	72.2	80.1	84.7	79.4	84.4	7.8
BLOOM	10	91.3	95.6	60.8	73.4	80.8	85.4	81.3	86.8	0	
CodeGen (Multi)	16B	10	91.8	97.3	62.0	69.6	81.1	84.7	81.8	86.4	0
		20	94.5	97.8	64.6	73.4	83.7	87.1	84.4	88.4	0
		10	91.3	98.4	60.8	69.6	76.9	83.7	78.9	86.0	0
	2B	20	94.0	98.4	60.8	72.2	81.6	87.1	82.5	88.4	0
		10	92.3	98.4	59.5	68.4	75.0	78.6	77.9	82.8	0
		20	93.4	98.9	64.6	72.2	82.5	85.2	83.4	87.4	0
350M	10	80.9	84.7	59.5	65.8	60.7	62.1	66.0	68.7	0	
	20	90.7	96.2	59.5	67.1	71.6	72.8	75.4	78.5	0	
CodeGen2	16B	10	93.4	97.8	59.5	68.4	81.6	85.0	82.2	86.5	0
		20	93.4	98.9	60.8	70.9	84.2	88.1	84.0	89.0	0
		10	89.6	95.1	64.6	70.9	77.9	80.8	79.5	83.5	0
	3.7B	20	91.8	96.7	64.6	72.2	84.5	87.6	84.2	88.3	0
		10	90.7	93.4	60.8	67.1	57.5	60.2	66.9	70.0	0
		20	92.3	96.7	59.5	69.6	73.3	75.2	76.8	80.4	0
1B	10	86.9	93.4	59.5	63.3	68.7	70.2	72.6	75.7	0	
	20	90.7	96.7	58.2	62.0	76.7	78.9	78.3	81.7	0	
InCoder	6B	10	32.8	40.4	45.6	54.4	54.6	57.8	47.6	52.7	0
		20	37.2	38.3	50.6	60.8	68.7	72.1	58.0	61.6	0
		10	32.8	37.7	48.1	53.2	58.5	61.7	50.3	54.2	0
1B	20	34.4	40.4	44.3	53.2	71.8	74.8	58.4	62.9	0	

automatically from limited datasets, are often insufficient or incorrect for traditional rule-based approaches, which causes 73 – 93% *unique* baseline failures.

VII. RQ4: MODEL COMPARISON

Table VII and Table VIII compare the performance and cost of 15 LLMs with baselines. We only list the best results for StarCoder with SR ($K = 60$), and the full results can be found in Table V and VI. Model response time for generating specifications, subject to various factors such as environment, is omitted. Generally, the response time is reasonably quick for practical usage, ranging between 0.6 to 26.6 seconds. Due to the token limitation discussed in Section IV-A, some experiments are skipped. More experiments on DocTer-data are skipped since prompts for DocTer-data are much longer than those for Jdoctor-data, making them inapplicable for certain settings.

Overall, generic LLMs with as few as 10 domain examples achieve better or comparable performance as custom-built state-of-the-art specification extraction techniques such as DocTer. Specifically, nine out of the 15 models outperform Jdoctor’s 83.0% accuracy. For DocTer-data, four out of the 15 models outperform DocTer’s 81.6% F1 score and all remaining models’ F1 scores are only short of a few percentage points from DocTer’s.

Finding 6: Most LLMs achieve better or comparable performance as custom-built traditional specification extraction techniques.

a) StarCoder: Among all 15 models, StarCoder achieves the best performance on DocTer-data, with an overall F1

TABLE VIII: Comparison of different LLMs with SR on DocTer-data: Precision/Recall/F1 (%), and Cost.

Approach / Model (+SR)	K	TensorFlow	PyTorch	MXNet	Overall	Cost (\$)	
DocTer	898	90.0/74.8/81.7	78.4/77.4/77.9	87.9/82.4/85.1	85.4/78.2/81.6	-	
StarCoder	60	85.0/86.7/85.8	83.4/84.6/84.0	87.4/89.2/88.3	85.9/87.5/ 86.7	0	
GPT-3	davinci	10	81.8/83.9/82.9	81.4/83.7/82.5	85.2/87.9/86.5	83.4/85.8/ 84.6	26.2
		20	84.2/86.5/85.4	83.4/84.9/84.1	86.8/89.2/88.0	85.3/87.5/ 86.4	52.4
		10	68.6/73.9/71.2	69.5/75.6/72.4	81.4/83.4/82.4	74.9/78.8/76.8	2.6
GPT-3.5	turbo	10	79.4/81.7/80.5	73.6/79.8/76.6	81.2/85.6/83.3	79.3/83.3/81.2	2.6
		20	81.9/83.7/82.8	77.3/81.2/79.2	84.3/87.3/85.8	82.3/85.0/ 83.6	5.2
BLOOM	10	74.3/79.0/76.5	73.3/76.6/74.9	78.6/82.2/80.3	76.2/80.1/78.1	0	
CodeGen (Multi)	16B	10	79.4/77.7/78.5	77.5/76.2/76.8	83.2/82.7/82.9	80.9/79.9/80.3	0
	6B	10	73.9/76.0/74.9	76.4/76.3/76.4	81.5/81.9/81.7	78.0/78.9/78.4	0
	2B	10	69.4/73.4/71.4	70.1/76.2/73.0	76.0/80.9/78.4	72.7/77.5/75.0	0
	350M	10	65.1/74.2/69.4	65.8/76.6/70.8	76.9/82.5/79.6	70.9/78.6/74.5	0
CodeGen2	16B	10	79.7/81.4/80.6	77.6/79.3/78.5	85.0/86.3/85.7	81.9/83.4/82.7	0
	7B	10	77.7/77.3/77.5	76.6/77.2/76.9	82.2/84.0/83.1	79.7/80.5/80.1	0
	3.7B	10	75.0/78.6/76.8	73.9/77.3/75.6	83.5/85.3/84.4	78.9/81.8/80.3	0
InCoder	6B	10	73.0/73.2/73.1	72.0/73.0/72.5	81.8/78.4/80.1	77.1/75.7/76.4	0
	1B	10	71.4/75.2/73.2	74.5/76.0/75.2	85.2/82.5/83.8	78.6/78.8/78.6	0

score of 86.7% (Table VIII). In addition, it achieves one of the highest overall accuracies of 93.0% on Jdoctor-data (Table VII). Compared to StarCoder, which is free, open-source, and with much fewer parameters (Table II), paid models GPT-3.5 and GPT-3 add no F1 gains on DocTer-data. Specifically, GPT-3 Davinci’s total cost of \$78.6 on DocTer-data causes a 0.3% F1 score degradation. GPT3’s tiny accuracy improvement of 2.6% comes at an additional \$163.8 total cost on Jdoctor-data. As Table II shows, StarCoder supports a larger number of input tokens (8,192) than other models, which enables us to add more examples to prompts for better performance. Moreover, compared with closed-source LLM APIs, open-sourced StarCoder offers enhanced flexibility and customizability for research and development.

Finding 7: StarCoder, an open-source model, is the most competitive model for extracting specifications, with among the highest performance of 86.7% and 93.0%, \$0 cost, and long prompt support, facilitating its adaptability and customization.

b) Paid Models (GPT-3 Davinci, GPT-3 Curie, and GPT-3.5): GPT-3 Davinci (175B parameters, \$0.02 per 1,000 tokens) achieves comparable or slightly worse performance than StarCoder given the same number of examples. GPT-3 Davinci supports fewer prompt tokens (Table II), preventing it from achieving higher performance with more examples on DocTer-data. In addition, it is much more costly and larger in size than StarCoder. The cheaper GPT-3 model, Curie, has a relatively poor performance. Despite being more recent, GPT-3.5 Turbo has lower performance on both Jdoctor (6.3% accuracy decrease) and DocTer (2.8% F1 decrease) compared to GPT-3 Davinci, possibly due to its optimization for conversations.

c) Other Open-Source Models (BLOOM, CodeGen, CodeGen2, and InCoder): An open-source GPT-3 alternative, BLOOM underperforms, with 8.3 – 8.8% lower performance than GPT-3 Davinci. Non-GPT-3 models, CodeGen and CodeGen2 achieve similar performance, which is comparable with BLOOM despite their much smaller sizes. We also conducted the experiment for DocTer-data with CodeGen’s Mono series, specifically trained for Python (Section III-C), and

they underperform its Multi series by 2.7%, potentially due to less generalizability. InCoder significantly underperforms CodeGen and CodeGen2 on Jdoctor-data, likely due to its pre-training data containing more Python (50 GB) than Java code (6 GB).

Finding 8: StarCoder’s strong performance makes GPT-3 Davinci less desirable given its size and cost. CodeGen and CodeGen2 are reasonable open-source alternatives.

VIII. CHALLENGES AND FUTURE DIRECTIONS

Our detailed analysis of the failure cases reveals several challenges and provides insights for future research directions, which are grouped into two interconnected topics:

a) Hybrid Approaches: Our results (e.g., Section VI-B3) suggest that the distinct strengths and weaknesses of LLMs and traditional methods can be complementary, providing an opportunity to explore hybrid approaches for specification generation. Such approaches can leverage the excellent generalizability of LLMs, supplemented by the context-rich and domain-specific insights of traditional methods, thus mitigating major challenges we identified in LLM applications (e.g., missing domain knowledge). Some studies have already made promising progress, such as the integration of LLMs with software testing and program analysis [21], [22], [63], demonstrating the potential of converging LLMs and traditional methods in future research for enhanced performance in SE tasks.

Specifically, to enhance LLMs for a hybrid approach for generating more accurate specifications, addressing the lack of domain knowledge and ineffective prompts is crucial, e.g., by providing sufficient context, and leveraging domain-specific example selection and ordering methods for prompt construction.

b) Improving Prompt Effectiveness: Our study identify another key direction of improving the effectiveness of prompts for LLMs to improve specification extraction. Recent research on designing expressive, customizable, and domain-specific prompts [64]–[66] have emerged as effective techniques to guide LLMs more accurately.

In addition, our study highlights the challenge of utilizing LLMs within the constraint token limit. It poses a barrier to providing sufficient context and examples for achieving better performance. Exploring more effective prompts for complex tasks, e.g., chain-of-thought [67], and iterative interaction [21]–[23], [68], offers exciting potentials for future research.

IX. THREATS TO VALIDITY

a) Reproducibility: Some of the evaluated models are not open-sourced, posing a threat to reproducibility. Nonetheless, our conclusion that LLMs with FSL outperform traditional approaches remains valid. The open-sourced StarCoder, served as our benchmark model (Section III-D), achieves one of the best performance and surpasses the baseline methods, Jdoctor and DocTer (Section VII). We expect that future models will continue to improve upon these results.

b) Manual Evaluation: To address the equivalence specification issue in Jdoctor-data (Section III-B3), we manually evaluated results for RQ1, RQ2, and RQ4 (Sections V-A, V-B, VII), and analyzed 180 samples in RQ3 (Section VI). To minimize biases, two authors *independently* conducted evaluations, resolving disagreements with a third author.

X. RELATED WORK

1) Software Specification Datasets and Extraction Methods: Traditional techniques for extracting software specifications from text, such as rule-based [1]–[4], [8], [12] or ML-based methods [1], [4], [11], require manual effort and domain knowledge, and lack generalizability across domains. Our work is the first to study LLMs’ capability on this task, leveraging FSL that offers improved generalizability and requires little annotated data. We evaluate LLMs on two datasets, Jdoctor [1] and DocTer [3]. Other than Jdoctor, techniques like @tcomment [8], Toradocu [2], and C2S [10] also extract specifications from Javadoc. They are excluded from this study as C2S is unavailable and the others are outdated or less effective [1], [10]. Advance [11] and DRONE [12] are excluded due to the absence of ground-truth specifications.

2) LLMs and FSL: LLMs such as GPT-3 have demonstrated their strong capabilities in numerous fields. We evaluate 15 state-of-the-art LLMs, varying in their design, sizes, etc., and discussed them in Section II. While other LLMs exist, they are not explored in this study as they are unsuitable for our task [13], less effective [19], [30], [31], [69], or unable to fit in our devices [70]. We exclude GPT-4 [71] due to its limited availability and high cost, which is 30x that of GPT-3.5 [72]. Moreover, open-source alternatives like StarCoder already exhibit impressive performance, reinforcing our primary conclusion that FSL-equipped LLMs outperform traditional methods.

One of the most effective way of utilizing LLMs, aside from zero-shot, and one-shot learning, is FSL [27], which adapts pre-trained LLMs to specific tasks with limited labeled data, yielding performance improvements across various tasks [15], [33]–[35], [73], [74]. Our work leverages the powerful generalizability of LLMs and evaluates their capabilities on generating software specifications.

3) Applications of LLMs to SE tasks: LLMs have been applied to many SE tasks, aside from the ones discussed in Section VIII, such as code completion [16]–[18], [26], test case generation [21], [22], [75], [76], etc. Despite LLMs outstanding performance in certain tasks, their limitations have also been exposed in areas such as code summarization [77], code suggestions [78], and software Q&A [79]. This further motivates our comprehensive evaluation of LLMs’ potential in generating software specifications, aimed at uncovering their inherent strengths and weaknesses.

XI. CONCLUSION

We present the first empirical study that assesses the effectiveness of 15 LLMs for software specifications generation. Our findings reveal that most LLMs achieve better or comparable performance compared to traditional methods. One

of the best-performing model, StarCoder, as an open-source model, outperforms traditional approaches by 5.1 – 10.0% with semantically similar examples. Its strong performance makes closed-source models (e.g., GPT-3 Davinci) less desirable due to size and cost. Additionally, we conduct a comprehensive failure diagnosis and identify strengths and weaknesses of both traditional methods and LLMs. The two dominant limitations of LLMs (ineffective prompts and missing domain knowledge) result in 57 – 67% of their failures. Our study offers insights for future research to improve LLMs' performance on specification generation including hybrid approaches of combining traditional methods and LLMs, and improving prompts effectiveness.

XII. ACKNOWLEDGEMENT

The authors thank the anonymous reviewers for their invaluable feedback and Guannan Wei for his input on the early draft. The research is partially supported by NSF 1901242 and 2006688 awards. It is also partially supported by the National Research Foundation of South Korea (grants No. 2022R1A2C1091419 and 2022H1D3A2A01092979), the Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Korean government (MSIT) (No. 2022-1113), and the 2021 Research Fund (1.210111.01) of UNIST (Ulsan National Institute of Science & Technology).

REFERENCES

- [1] A. Blasi, A. Goffi, K. Kuznetsov, A. Gorla, M. D. Ernst, M. Pezzè, and S. D. Castellanos, “Translating code comments to procedure specifications,” *Proceedings of the 27th ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2018.
- [2] A. Goffi, A. Gorla, M. D. Ernst, and M. Pezzè, “Automatic generation of oracles for exceptional behaviors,” in *Proceedings of the 25th international symposium on software testing and analysis*, 2016, pp. 213–224.
- [3] D. Xie, Y. Li, M. Kim, H. V. Pham, L. Tan, X. Zhang, and M. W. Godfrey, “Docter: documentation-guided fuzzing for testing deep learning api functions,” *Proceedings of the 31st ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2021.
- [4] L. Tan, D. Yuan, G. Krishna, and Y. Zhou, “/*comment: bugs or bad comments?*/,” in *Symposium on Operating Systems Principles*, 2007.
- [5] E. Wong, L. Zhang, S. Wang, T. Liu, and L. Tan, “Dase: Document-assisted symbolic execution for improving automated software testing,” *2015 IEEE/ACM 37th IEEE International Conference on Software Engineering*, vol. 1, pp. 620–631, 2015.
- [6] C. Boyapati, S. Khurshid, and D. Marinov, “Korat: Automated testing based on java predicates,” *SIGSOFT Software Engineering Notes*, vol. 27, no. 4, pp. 123–133, 2002.
- [7] C. Cadar, D. Dunbar, and D. Engler, “KLEE: Unassisted and automatic generation of high-coverage tests for complex systems programs,” in *Proceedings of the 8th USENIX conference on Operating systems design and implementation*, 2008, pp. 209–224.
- [8] S. H. Tan, D. Marinov, L. Tan, and G. T. Leavens, “@tcomment: Testing javadoc comments to detect comment-code inconsistencies,” *2012 IEEE Fifth International Conference on Software Testing, Verification and Validation*, pp. 260–269, 2012.
- [9] R. Pandita, X. Xiao, H. Zhong, T. Xie, S. Oney, and A. Paradkar, “Inferring method specifications from natural language api descriptions,” in *2012 34th international conference on software engineering (ICSE)*. IEEE, 2012, pp. 815–825.
- [10] J. Zhai, Y. Shi, M. Pan, G. Zhou, Y. Liu, C. Fang, S. Ma, L. Tan, and X. Zhang, “C2s: translating natural language comments to formal program specifications,” *Proceedings of the 28th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering*, 2020.
- [11] T. Lv, R. Li, Y. Yang, K. Chen, X. Liao, X. Wang, P. Hu, and L. Xing, “Rtfm! automatic assumption discovery and verification derivation from library document for api misuse detection,” *Proceedings of the 2020 ACM SIGSAC Conference on Computer and Communications Security*, 2020.
- [12] Y. Zhou, C. Wang, X. Yan, T. Chen, S. Panichella, and H. C. Gall, “Automatic detection and repair recommendation of directive defects in java api documentation,” *IEEE Transactions on Software Engineering*, vol. 46, pp. 1004–1023, 2020.
- [13] J. Devlin, M. Chang, K. Lee, and K. Toutanova, “BERT: pre-training of deep bidirectional transformers for language understanding,” *CoRR*, vol. abs/1810.04805, 2018. [Online]. Available: <http://arxiv.org/abs/1810.04805>
- [14] A. Radford, J. Wu, R. Child, D. Luan, D. Amodei, and I. Sutskever, “Language models are unsupervised multitask learners,” 2019.
- [15] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” *CoRR*, vol. abs/2005.14165, 2020. [Online]. Available: <https://arxiv.org/abs/2005.14165>
- [16] M. Chen, J. Tworek, H. Jun, Q. Yuan, H. P. de Oliveira Pinto, J. Kaplan, H. Edwards, Y. Burda, N. Joseph, G. Brockman, A. Ray, R. Puri, G. Krueger, M. Petrov, H. Khlaaf, G. Sastry, P. Mishkin, B. Chan, S. Gray, N. Ryder, M. Pavlov, A. Power, L. Kaiser, M. Bavarian, C. Winter, P. Tillet, F. P. Such, D. Cummings, M. Plappert, F. Chantzis, E. Barnes, A. Herbert-Voss, W. H. Guss, A. Nichol, A. Paino, N. Tezak, J. Tang, I. Babuschkin, S. Balaji, S. Jain, W. Saunders, C. Hesse, A. N. Carr, J. Leike, J. Achiam, V. Misra, E. Morikawa, A. Radford, M. Knight, M. Brundage, M. Murati, K. Mayer, P. Welinder, B. McGrew, D. Amodei, S. McCandlish, I. Sutskever, and W. Zaremba, “Evaluating large language models trained on code,” *CoRR*, vol. abs/2107.03374, 2021. [Online]. Available: <https://arxiv.org/abs/2107.03374>
- [17] E. Nijkamp, B. Pang, H. Hayashi, L. Tu, H. Wang, Y. Zhou, S. Savarese, and C. Xiong, “A conversational paradigm for program synthesis,” *arXiv preprint*, 2022.
- [18] D. Fried, A. Aghajanyan, J. Lin, S. Wang, E. Wallace, F. Shi, R. Zhong, W.-t. Yih, L. Zettlemoyer, and M. Lewis, “InCoder: A generative model for code infilling and synthesis,” 2022. [Online]. Available: <https://arxiv.org/abs/2204.05999>
- [19] W. Yue, W. Weishi, J. Shafiq, and C. H. Steven, “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, EMNLP 2021*, 2021.
- [20] J. Liu, C. S. Xia, Y. Wang, and L. Zhang, “Is your code generated by chatgpt really correct? rigorous evaluation of large language models for code generation,” *arXiv preprint arXiv:2305.01210*, 2023.
- [21] Z. Yuan, Y. Lou, M. Liu, S. Ding, K. Wang, Y. Chen, and X. Peng, “No more manual tests? evaluating and improving chatgpt for unit test generation,” *arXiv preprint arXiv:2305.04207*, 2023.
- [22] B. Chen, F. Zhang, A. Nguyen, D. Zan, Z. Lin, J.-G. Lou, and W. Chen, “Codet: Code generation with generated tests,” *arXiv preprint arXiv:2207.10397*, 2022.
- [23] M. Schäfer, S. Nadi, A. Eghbali, and F. Tip, “Adaptive test generation using a large language model,” *arXiv preprint arXiv:2302.06527*, 2023.
- [24] N. Jiang, T. Lutellier, and L. Tan, “Impact of code language models on automated program repair,” in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, 2023.
- [25] K. Pei, D. Bieber, K. Shi, C. Sutton, and P. Yin, “Can large language models reason about program invariants?” 2023.
- [26] W. Ahmad, S. Chakraborty, B. Ray, and K.-W. Chang, “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*. Online: Association for Computational Linguistics, Jun. 2021, pp. 2655–2668. [Online]. Available: <https://aclanthology.org/2021.naacl-main.211>
- [27] T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell, S. Agarwal, A. Herbert-Voss, G. Krueger, T. J. Henighan, R. Child, A. Ramesh, D. M. Ziegler, J. Wu, C. Winter, C. Hesse, M. Chen, E. Sigler, M. Litwin, S. Gray, B. Chess, J. Clark, C. Berner, S. McCandlish, A. Radford, I. Sutskever, and D. Amodei, “Language models are few-shot learners,” *ArXiv*, vol. abs/2005.14165, 2020.
- [28] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, “Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension,” in *Annual Meeting of the Association for Computational Linguistics*, 2019.
- [29] A. Radford, K. Narasimhan, T. Salimans, and I. Sutskever, “Improving language understanding by generative pre-training,” 2018.
- [30] B. Wang and A. Komatsuzaki, “GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model,” <https://github.com/kingoflolz/mesh-transformer-jax>, May 2021.
- [31] S. Black, G. Leo, P. Wang, C. Leahy, and S. Biderman, “GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow,” Mar. 2021. [Online]. Available: <https://doi.org/10.5281/zenodo.5297715>
- [32] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, “Exploring the limits of transfer learning with a unified text-to-text transformer,” *CoRR*, vol. abs/1910.10683, 2019. [Online]. Available: <http://arxiv.org/abs/1910.10683>
- [33] P. Riley, N. Constant, M. Guo, G. Kumar, D. Uthus, and Z. Parekh, “Textsettr: Few-shot text style extraction and tunable targeted restyling,” 01 2021, pp. 3786–3800.
- [34] O. Ram, Y. Kirstain, J. Berant, A. Globerson, and O. Levy, “Few-shot question answering by pretraining span selection,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*. Online: Association for Computational Linguistics, Aug. 2021, pp. 3066–3079. [Online]. Available: <https://aclanthology.org/2021.acl-long.239>
- [35] R. Wang, T. Yu, H. Zhao, S. Kim, S. Mitra, R. Zhang, and R. Henao, “Few-shot class-incremental learning for named entity recognition,” in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 571–582. [Online]. Available: <https://aclanthology.org/2022.acl-long.43>
- [36] S. M. Xie, A. Raghunathan, P. Liang, and T. Ma, “An explanation of in-context learning as implicit bayesian inference,” in *International Conference on Learning Representations*, 2022. [Online]. Available: <https://openreview.net/forum?id=RdJVFCHjUMI>

- [37] L. Breiman and P. Spector, "Submodel selection and evaluation in regression. the x-random case," *International statistical review/revue internationale de Statistique*, pp. 291–319, 1992.
- [38] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning*, ser. Springer Series in Statistics. New York, NY, USA: Springer New York Inc., 2001.
- [39] M. Magnusson, M. Andersen, J. Jonasson, and A. Vehtari, "Bayesian leave-one-out cross-validation for large data," in *Proceedings of the 36th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. PMLR, 09–15 Jun 2019, pp. 4244–4253. [Online]. Available: <https://proceedings.mlr.press/v97/magnusson19a.html>
- [40] O. Rubin, J. Herzig, and J. Berant, "Learning to retrieve prompts for in-context learning," *ArXiv*, vol. abs/2112.08633, 2021.
- [41] J. Liu, D. Shen, Y. Zhang, B. Dolan, L. Carin, and W. Chen, "What makes good in-context examples for gpt-3?" *arXiv preprint arXiv:2101.06804*, 2021.
- [42] Y. Liu, M. Ott, N. Goyal, J. Du, M. Joshi, D. Chen, O. Levy, M. Lewis, L. Zettlemoyer, and V. Stoyanov, "Roberta: A robustly optimized BERT pretraining approach," *CoRR*, vol. abs/1907.11692, 2019. [Online]. Available: <http://arxiv.org/abs/1907.11692>
- [43] "all-roberta-large-v1," <https://huggingface.co/sentence-transformers/all-roberta-large-v1>, Accessed: 2023.
- [44] N. Reimers and I. Gurevych, "Sentence-bert: Sentence embeddings using siamese bert-networks," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. [Online]. Available: <http://arxiv.org/abs/1908.10084>
- [45] R. Shin, C. H. Lin, S. Thomson, C. C. Chen, S. Roy, E. A. Platanios, A. Pauls, D. Klein, J. Eisner, and B. V. Durme, "Constrained language models yield few-shot semantic parsers," *ArXiv*, vol. abs/2104.08768, 2021.
- [46] R. Li, L. B. Allal, Y. Zi, N. Muennighoff, D. Kocetkov, C. Mou, M. Marone, C. Akiki, J. Li, J. Chim *et al.*, "Starcoder: may the source be with you!" *arXiv preprint arXiv:2305.06161*, 2023.
- [47] OpenAI, "Gpt-3.5," 2022. [Online]. Available: <https://platform.openai.com/docs/models/gpt-3-5>
- [48] T. L. Scao *et al.*, "Bloom: A 176b-parameter open-access multilingual language model," *ArXiv*, vol. abs/2211.05100, 2022.
- [49] E. Nijkamp, H. Hayashi, C. Xiong, S. Savarese, and Y. Zhou, "Codegen2: Lessons for training llms on programming and natural languages," *arXiv preprint arXiv:2305.02309*, 2023.
- [50] OpenAI, "Chatgpt," 2022, accessed: 2023-03-10. [Online]. Available: <https://openai.com/chat-gpt/>
- [51] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *CoRR*, vol. abs/1706.03762, 2017. [Online]. Available: <http://arxiv.org/abs/1706.03762>
- [52] L. B. Allal, R. Li, D. Kocetkov, C. Mou, C. Akiki, C. M. Ferrandis, N. Muennighoff, M. Mishra, A. Gu, M. Dey *et al.*, "Santacoder: don't reach for the stars!" *arXiv preprint arXiv:2301.03988*, 2023.
- [53] D. Kocetkov, R. Li, L. Ben Allal, J. Li, C. Mou, C. Muñoz Ferrandis, Y. Jernite, M. Mitchell, S. Hughes, T. Wolf, D. Bahdanau, L. von Werra, and H. de Vries, "The stack: 3 tb of permissively licensed source code," *Preprint*, 2022.
- [54] D. M. Ziegler, N. Stiennon, J. Wu, T. B. Brown, A. Radford, D. Amodei, P. F. Christiano, and G. Irving, "Fine-tuning language models from human preferences," *CoRR*, vol. abs/1909.08593, 2019. [Online]. Available: <http://arxiv.org/abs/1909.08593>
- [55] L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin, C. Zhang, S. Agarwal, K. Slama, A. Ray, J. Schulman, J. Hilton, F. Kelton, L. Miller, M. Simens, A. Askell, P. Welinder, P. Christiano, J. Leike, and R. Lowe, "Training language models to follow instructions with human feedback," 2022. [Online]. Available: <https://arxiv.org/abs/2203.02155>
- [56] R. Nakano, J. Hilton, S. Balaji, J. Wu, L. Ouyang, C. Kim, C. Hesse, S. Jain, V. Kosaraju, W. Saunders, X. Jiang, K. Cobbe, T. Eloundou, G. Krueger, K. Button, M. Knight, B. Chess, and J. Schulman, "Webgpt: Browser-assisted question-answering with human feedback," *CoRR*, vol. abs/2112.09332, 2021. [Online]. Available: <https://arxiv.org/abs/2112.09332>
- [57] N. Lambert, L. Castriato, L. von Werra, and A. Havrilla, "Illustrating reinforcement learning from human feedback (rlhf)," *Hugging Face Blog*, 2022, <https://huggingface.co/blog/rlhf>.
- [58] "Hugging face api for bloom," <https://huggingface.co/bigscience/bloom>, Accessed: 2023.
- [59] L. Gao, S. Biderman, S. Black, L. Golding, T. Hoppe, C. Foster, J. Phang, H. He, A. Thite, N. Nabeshima, S. Presser, and C. Leahy, "The pile: An 800gb dataset of diverse text for language modeling," *CoRR*, vol. abs/2101.00027, 2021. [Online]. Available: <https://arxiv.org/abs/2101.00027>
- [60] "Bigquery dataset," <https://console.cloud.google.com/marketplace/details/github/github-repos?pli=1>, Accessed: 2023.
- [61] X. V. Lin, T. Mihaylov, M. Artetxe, T. Wang, S. Chen, D. Simig, M. Ott, N. Goyal, S. Bhosale, J. Du, R. Pasunuru, S. Shleifer, P. S. Koura, V. Chaudhary, B. O'Horo, J. Wang, L. Zettlemoyer, Z. Kozareva, M. T. Diab, V. Stoyanov, and X. Li, "Few-shot learning with multilingual language models," *CoRR*, vol. abs/2112.10668, 2021. [Online]. Available: <https://arxiv.org/abs/2112.10668>
- [62] "Preparing your dataset," <https://platform.openai.com/docs/guides/fine-tuning/preparing-your-dataset>, Accessed: 2023.
- [63] C. Lemieux, J. P. Inala, S. K. Lahiri, and S. Sen, "Codamosa: Escaping coverage plateaus in test generation with pre-trained large language models," in *International conference on software engineering (ICSE)*, 2023.
- [64] S. Abukhalaf, M. Hamdaqa, and F. Khomh, "On codex prompt engineering for ocl generation: An empirical study," *arXiv preprint arXiv:2303.16244*, 2023.
- [65] L. Beurer-Kellner, M. Fischer, and M. Vechev, "Prompting is programming: A query language for large language models," *Proceedings of the ACM on Programming Languages*, vol. 7, no. PLDI, pp. 1946–1969, 2023.
- [66] C. Wang, Y. Yang, C. Gao, Y. Peng, H. Zhang, and M. R. Lyu, "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*, 2022, pp. 382–394.
- [67] X. Wang, J. Wei, D. Schuurmans, Q. Le, E. Chi, S. Narang, A. Chowdhery, and D. Zhou, "Self-consistency improves chain of thought reasoning in language models," *arXiv preprint arXiv:2203.11171*, 2022.
- [68] C. S. Xia and L. Zhang, "Keep the conversation going: Fixing 162 out of 337 bugs for \$0.42 each using chatgpt," *arXiv preprint arXiv:2304.00385*, 2023.
- [69] H. Touvron, T. Lavril, G. Izacard, X. Martinet, M.-A. Lachaux, T. Lacroix, B. Rozière, N. Goyal, E. Hambro, F. Azhar *et al.*, "Llama: Open and efficient foundation language models," *arXiv preprint arXiv:2302.13971*, 2023.
- [70] S. Black, S. R. Biderman, E. Hallahan, Q. G. Anthony, L. Gao, L. Golding, H. He, C. Leahy, K. McDonell, J. Phang, M. M. Pieler, U. S. Prashanth, S. Purohit, L. Reynolds, J. Tow, B. Wang, and S. Weinbach, "Gpt-neox-20b: An open-source autoregressive language model," *ArXiv*, vol. abs/2204.06745, 2022.
- [71] OpenAI, "Gpt-4," 2022. [Online]. Available: <https://platform.openai.com/docs/models/gpt-4>
- [72] —, "Pricing," 2023. [Online]. Available: <https://openai.com/pricing>
- [73] J. Chen, Q. Liu, H. Lin, X. Han, and L. Sun, "Few-shot named entity recognition with self-describing networks," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 5711–5722. [Online]. Available: <https://aclanthology.org/2022.acl-long.392>
- [74] K. Krishna, D. Nathani, X. Garcia, B. Samanta, and P. Talukdar, "Few-shot controllable style transfer for low-resource multilingual settings," in *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*. Dublin, Ireland: Association for Computational Linguistics, May 2022, pp. 7439–7468. [Online]. Available: <https://aclanthology.org/2022.acl-long.514>
- [75] Y. Deng, C. S. Xia, H. Peng, C. Yang, and L. Zhang, "Large language models are zero-shot fuzzers: Fuzzing deep-learning libraries via large language models," in *Proceedings of the 32nd ACM SIGSOFT International Symposium on Software Testing and Analysis*, 2023, pp. 423–435.
- [76] S. Kang, J. Yoon, and S. Yoo, "Large language models are few-shot testers: Exploring llm-based general bug reproduction," in *2023 IEEE/ACM 45th International Conference on Software Engineering (ICSE)*, 2023.
- [77] W. Sun, C. Fang, Y. You, Y. Miao, Y. Liu, Y. Li, G. Deng, S. Huang, Y. Chen, Q. Zhang *et al.*, "Automatic code summarization via chatgpt: How far are we?" *arXiv preprint arXiv:2305.12865*, 2023.
- [78] A. M. Dakhel, V. Majdinasab, A. Nikanjam, F. Khomh, M. C. Desmarais, and Z. M. J. Jiang, "Github copilot ai pair programmer: Asset or liability?" *Journal of Systems and Software*, vol. 203, p. 111734, 2023.

- [79] B. Xu, T.-D. Nguyen, T. Le-Cong, T. Hoang, J. Liu, K. Kim, C. Gong, C. Niu, C. Wang, B. Le *et al.*, “Are we ready to embrace generative ai for software q&a?” *arXiv preprint arXiv:2307.09765*, 2023.